

# NERFIFY: A Multi-Agent Framework for Turning NeRF Papers into Code

## Supplementary Material

### Summary of Supplementary Material

- Section 1 : Extended Experimental Analysis
  - Section 1.1 : Complete NERFIFY-BENCH Results
  - Section 1.2 : Qualitative Comparisons Gallery
  - Section 1.3 : Additional Ablation Studies
- Section 2: NERFIFY-BENCH Dataset
- Section 3: Comparison with Multi-Agent Baselines
- Section 4: Novelty Coverage on Set 4 Papers
- Section 5: Implementation Details
  - Section 5.1 : Context-Free Grammar
  - Section 5.2 : Multi-Agent Architecture and Agent Specifications
  - Section 5.3 : LLM Prompts Used in NERFIFY

## 1. Extended Experimental Analysis

### 1.1. Complete NERFIFY-BENCH Results

#### 1.1.1. Quantitative Metrics for Set 1 (Never-Implemented Papers)

Paper	Reported			Human Impl.			NERFIFY (Ours)		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
KeyNeRF [21]	25.65	0.89	0.11	25.70	0.89	0.12	26.12	0.90	0.09
mi-MLP NeRF [34]	24.70	0.89	0.09	22.64	0.87	0.15	22.85	0.87	0.15
ERS [24]	27.85	0.94	0.06	26.87	0.90	0.12	27.02	0.90	0.12
TVNeRF [33]	27.44	0.93	0.08	26.81	0.92	0.12	27.30	0.92	0.10
Anisotropic NeRF [27]	34.08	0.97	0.05	28.85	0.94	0.06	29.01	0.94	0.06
NeRF-ID [1]	25.15	0.94	-	23.01	0.89	0.13	23.10	0.89	0.13
Surface Sam. [32]	25.63	-	-	24.35	0.90	0.11	24.40	0.90	0.11
LiNeRF [4]	25.60	0.93	0.08	23.30	0.90	0.14	23.32	0.89	0.14
HybNeRF [28]	33.94	0.96	0.047	30.45	0.94	0.07	30.51	0.95	0.07
AR-NeRF	20.36	0.79	0.17	19.00	0.76	0.19	20.05	0.78	0.18

Table 1. **Comparison of NERFIFY with paper and human implementations.** We evaluate NeRF papers from the NERFIFY-BENCH set whose code is not publicly available, using SSIM, PSNR, and LPIPS metrics. **Note.** Other baselines like Paper2Code, AutoP2C, GPT-5 and R1 failed to generate trainable code.

#### 1.1.2. Quantitative Metrics for Set 2 and Set 3 (Github Repository Available)

Method	Original Repository			NERFIFY		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Vanilla NeRF [16]	31.36	0.95	0.04	31.36	0.95	0.04
Nerfacto [25]	20.36	0.82	0.22	20.36	0.82	0.22
SeaThru-NeRF	27.89	0.83	0.22	30.08	0.92	0.07
BioNeRF	25.66	0.93	0.05	23.75	0.90	0.11
InstantNGP	32.77	-	-	32.64	0.96	0.06
$\ell_0$ Sampler [14]	29.21	-	0.04	30.13	0.97	0.03
InfoNeRF [11]	18.27	0.81	0.23	17.87	0.69	0.44
DeblurNeRF	32.08	0.93	0.5	31.10	0.86	0.06
NerfingMVS	31.43	0.96	-	31.51	0.92	0.09

Table 2. **Comparison with existing implementations.** Evaluation of NERFIFY against original author repositories or gold-standard implementations.

Configuration	PSNR	SSIM	LPIPS
NERFIFY (Full)	27.16	0.91	0.11
<i>Knowledge Sources:</i>			
w/o Citation Recovery and In-context Examples	23.22	0.87	0.26
<i>Validation &amp; Feedback:</i>			
w/o VLM Feedback and Smoke Testing (Stage 4)	9.39	0.80	0.42
<i>Planning Strategy:</i>			
One-Shot (no GoT) (Stage 3)	24.52	0.90	0.16

Table 3. **Component ablation study.** We evaluate the impact of each system component on synthesis quality and efficiency. Numbers are averaged on trainable implementations among 10 NERFIFY-BENCH Set 1 papers.

## 1.2. Qualitative Comparisons Gallery

### 1.2.1. Side-by-side Visual Comparisons

Figure 1 and Figure 2 shows the qualitative comparison of some of the methods in Set 1 of the NERFIFY bench.

### 1.3. Additional Ablation Studies

The ablation study experiments with the image quality metrics are shown in Table 3.

## 2. NERFIFY-BENCH Dataset

### 2.1. Paper Selection Criteria

NERFIFY-BENCH comprises 30 carefully curated NeRF research papers selected to provide comprehensive coverage of different implementation challenges. Our selection process prioritized diversity across architectural innovations, training strategies, and integration complexity to ensure robust evaluation of paper-to-code synthesis systems.

The benchmark is organized into four distinct set chosen from three categories as shown in figure 3, each designed to test specific capabilities of automated code generation systems.

#### 2.1.1. Category 1: Never-Implemented Papers (10 papers)

This category includes papers without any publicly available source code. To ensure evaluation integrity and avoid potential training data contamination, we specifically selected papers where no implementation exists in the public domain. This guarantees that large language models used in our agents could not have encountered corresponding code during pretraining, enabling unbiased assessment of purely paper-driven synthesis capabilities.

For these papers, we commissioned expert reimplementations from graduate students with extensive NeRF research

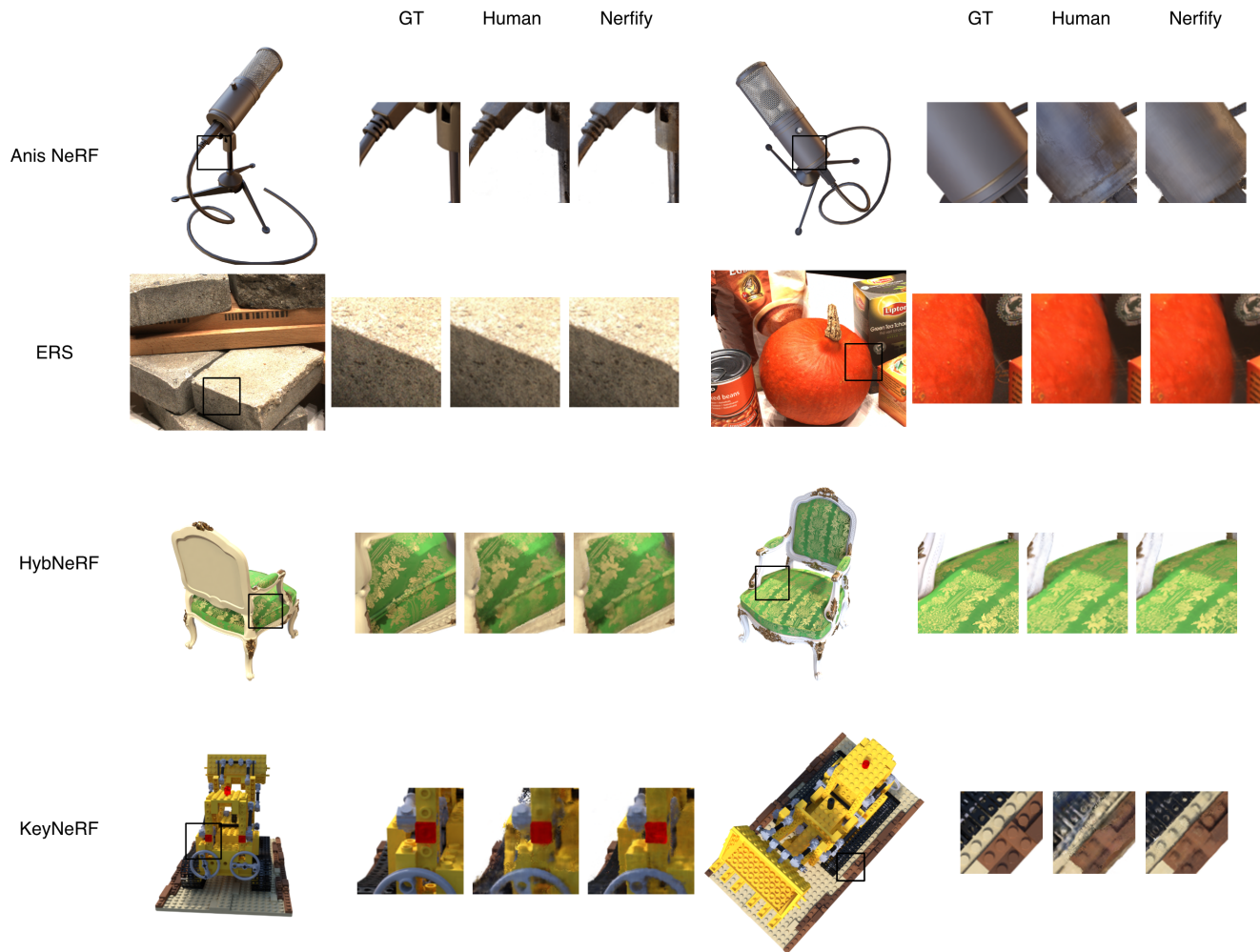


Figure 1. **Visual Comparison of NERFIFY and Human Implementation.** **Left:** Ground Truth Image, **Middle:** Expert Implementation, **Right:** Agent Implementation.

experience. Each implementation underwent rigorous verification against paper-reported metrics on standard benchmarks. These expert implementations serve as ground truth references for quantitative evaluation.

#### Papers in Category 1:

- KeyNeRF [21]: Informative ray selection for few-shot neural radiance fields
- mi-MLP NeRF [34]: Minimal MLP approach for few-shot view synthesis
- LiNeRF [4]: Rethinking directional integration in neural radiance fields
- Anisotropic Neural Representation: Adaptive resolution for complex geometries
- Efficient Ray Sampling [24]: Optimized sampling strategies for radiance field reconstruction
- HybNeRF [28]: Hybrid multiresolution encoding for neu-

ral radiance fields

- TVNeRF [33]: Total variation maximization for few-view neural volume rendering
- Surface Sampling [32]: Near-surface sampling with point cloud generation
- NeRF-ID [1]: Learning to sample for view synthesis
- AR-NeRF [30]: Few-shot NeRF by Adaptive Rendering Loss Regularization

#### 2.1.2. Category 2: Non-Nerfstudio Papers (5 papers)

These papers have existing public implementations but are not integrated into the Nerfstudio framework. This category enables direct comparison between our synthesized Nerfstudio plugins and original author implementations, evaluating both functional correctness and framework integration quality.

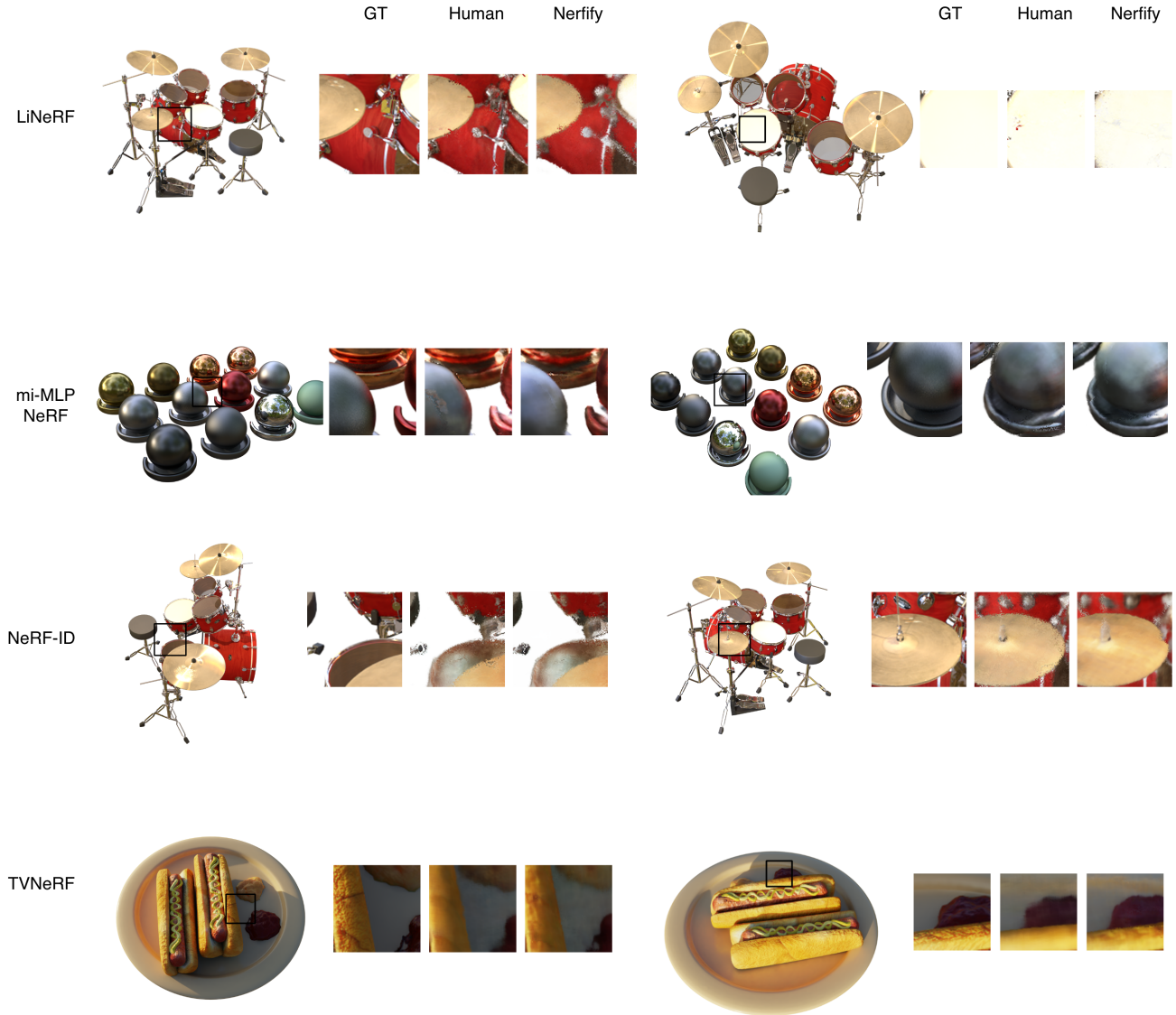


Figure 2. **Visual Comparison of NERFIFY and Human Implementation.** Left: Ground Truth Image, Middle: Expert Implementation, Right: Agent Implementation.

### Papers in Category 2:

- Deblur-NeRF [15]: Neural radiance fields from blurry images
- InfoNeRF [11]: Ray entropy minimization for few-shot neural volume rendering
- L0-Sampler [14]: Adaptive ray sampling for neural radiance fields
- NerfingMVS [29]: NerfingMVS: Guided Optimization of Neural Radiance Fields for Indoor Multi-view Stereo

### 2.1.3. Category 3: Nerfstudio-Integrated Papers (5 papers)

Papers already integrated into Nerfstudio serve as gold-standard references. These enable evaluation of synthesis quality and API compliance against production-quality implementations maintained by the research community.

#### Papers in Category 3:

- SeaThru-NeRF [13]: Neural radiance fields for underwater scenes
- BioNeRF [22]: Bionerf: Biologically Plausible Neural Radiance Fields for View Synthesis



Figure 3. **Categorization of NeRF Papers by Integrability in Nerfstudio.** Papers are grouped by implementation difficulty within the Nerfstudio framework. **Category 1:** Directly integrable methods modifying architecture, rendering, or losses without external dependencies. **Category 2:** Methods requiring pretrained models (CLIP, diffusion, SLAM) with substantial engineering effort. **Category 3:** Out-of-scope works where NeRF is used in the pipeline but serves different objectives (GANs, downstream applications, etc). The categorization illustrates alignment with Nerfstudio’s modular design and implementation feasibility using NERFIFY .

- Vanilla-NeRF [16]: Original neural radiance fields implementation
- Instant-NGP [17]: Instant neural graphics primitives with multiresolution hash encoding
- Nerfacto [25]: Nerfstudio: A modular framework for neural radiance field development

#### 2.1.4. Category 4: Novelty-Coverage Papers (10 papers)

Selected for their distinct technical contributions, these papers evaluate how well synthesis systems capture and implement key research innovations. Each introduces novel loss functions, architectural components, or training strategies that require deep understanding for correct implementation.

##### Papers in Category 4:

- Mip-NeRF [2]: Multiscale representation for anti-aliasing neural radiance fields
- BioNeRF [22]: Biologically plausible neural radiance fields for view synthesis
- PyNeRF [26]: Pyramidal neural radiance fields
- TensorRF [3]: Tensorial radiance fields
- Tetra-NeRF [12]: Representing neural radiance fields using tetrahedra
- E-NeRF: Efficient neural radiance fields with learned embeddings
- StyleNeRF [7]: Style-based 3D-aware generator for high-resolution image synthesis
- iNeRF [31]: Inverting neural radiance fields for pose estimation
- SigNeRF [5]: Scene integrated generation for neural radiance fields
- MCNeRF [8]: Monte Carlo rendering and denoising for real-time NeRFs

For each paper in NERFIFY-BENCH we provide comprehensive resources to enable reproducible evaluation:

**Frozen PDF Versions** The original paper PDF files are included to ensure consistency across evaluations. This prevents version drift and guarantees that all systems process identical input documents.

**Dual Markdown Representations** We provide two markdown versions for each paper:

- *Raw extraction*: Direct conversion from PDF preserving all text, equations, and structure
- *Cleaned version*: Processed through our paper extraction pipeline, with artifacts removed, mathematical notation standardized, and formatting normalized for optimal parsing

**Tex source code** We provide LaTeX source code of each paper.

**Ground Truth Repositories** Where available, we include either the original author implementation or expert reimplementations verified for correctness. These serve as reference implementations for quantitative metrics including PSNR, SSIM, and LPIPS on standard test scenes.

#### 2.1.5. Categorization Methodology

#### 2.2. Evaluation Protocol

Our benchmark evaluates synthesized code across following dimensions:

##### Executability Metrics

- *Build success rate*: Percentage of repositories that compile without errors
- *Import resolution*: Whether all dependencies resolve correctly
- *Training stability*: Convergence without NaN losses, gradient explosions, or memory errors

**Rendering Quality Metrics** For successfully trained models, we measure:

- *PSNR*: Peak signal-to-noise ratio on held-out test views
- *SSIM*: Structural similarity index
- *LPIPS*: Learned perceptual image patch similarity

Results are averaged across multiple scenes from standard benchmarks to ensure statistical significance.

**Novelty Fidelity** We evaluate coverage of paper-specific innovations:

- *Correct (C)*: Percentage of novel components implemented exactly as specified
- *Incorrect (I)*: Percentage with partial matches or flawed logic
- *Missing (M)*: Percentage absent from generated code
- *Weight match (W)*: Percentage of hyperparameters within 10% of paper specifications

##### Runtime Efficiency

- *Training throughput*: Rays processed per second
- *Memory usage*: Peak GPU memory consumption
- *Convergence speed*: Iterations to reach target quality

##### Code Quality

- *Nerfstudio conventions*: Adherence to framework patterns and APIs
- *Documentation*: Presence of docstrings and inline comments



- **Modularity:** Proper separation of concerns and reusable components

### 2.2.1. Evaluation Dataset Details

Standard datasets used for evaluation:

#### Papers in Category 1:

- **KeyNeRF [21]:** Informative ray selection for few-shot neural radiance fields
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Average across all scenes
- **mi-MLP NeRF [34]:** Minimal MLP approach for few-shot view synthesis
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Average across all scenes
- **LiNeRF [4]:** Rethinking directional integration in neural radiance fields
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Drums
- **Anisotropic Neural Representation:** Adaptive resolution for complex geometries
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Average across all scenes
- **Efficient Ray Sampling [24]:** Optimized sampling strategies for radiance field reconstruction
  - *Dataset:* DTU
  - *Scene:* Average on DTU test set
- **HybNeRF [28]:** Hybrid multiresolution encoding for neural radiance fields
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Average across all scenes
- **TVNeRF [33]:** Total variation maximization for few-view neural volume rendering
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Hotdog
- **Surface Sampling [32]:** Near-surface sampling with point cloud generation
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Ship (Table 2)
- **NeRF-ID [1]:** Learning to sample for view synthesis
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Drums (Number reported in appendix)
- **AR-NeRF [30]:** Few-shot NeRF by Adaptive Rendering Loss Regularization
  - *Dataset:* DTU
  - *Scene:* 3-view setting

#### Papers in Category 2:

- **Deblur-NeRF [15]:** Neural radiance fields from blurry images
  - *Dataset:* Deblur-NeRF Dataset
  - *Scene:* Cozyroom Camera Motion
- **InfoNeRF [11]:** Ray entropy minimization for few-shot neural volume rendering
  - *Dataset:* NeRF-Synthetic
  - *Scene:* Ship

- **L0-Sampler [14]:** Adaptive ray sampling for neural radiance fields
    - *Dataset:* NeRF-Synthetic
    - *Scene:* Ficus
  - **NerfingMVS [29]:** Guided optimization of neural radiance fields for indoor multi-view stereo
    - *Dataset:* ScanNet
    - *Scene:* Scene 0653 (processed data)
- #### Papers in Category 3:
- **SeaThru-NeRF [13]:** Neural radiance fields for underwater scenes
    - *Dataset:* SeaThru-NeRF Dataset
    - *Scene:* Panama
  - **BioNeRF [22]:** Biologically plausible neural radiance fields for view synthesis
    - *Dataset:* NeRF-Synthetic
    - *Scene:* Drums
  - **Vanilla-NeRF [16]:** Original neural radiance fields implementation
    - *Dataset:* NeRF-Synthetic
    - *Scene:* Chair
  - **Instant-NGP [17]:** Instant neural graphics primitives with multiresolution hash encoding
    - *Dataset:* NeRF-Synthetic
    - *Scene:* Lego
  - **Nerfacto [25]:** A modular framework for neural radiance field development
    - *Dataset:* Nerfstudio Dataset
    - *Scene:* Poster

## 3. Comparison with Multi-Agent Baselines

Table 4 compares NERFIFY with automated coding systems on papers from nerfifybench evaluating executability and semantic score. Recent PaperBench results show DeepCode (75.9% accuracy) substantially outperforms PaperCoder (51.1%), Claude Code (58.7%) and LLM-based systems like o1 (43.3%) [9]. So, we use DeepCode (with gpt-5 api) as our baseline alongside multi-agent frameworks GPT 5 thinking [20], ChatDev [23] (with gpt-5 api) and MetaGPT [10] (with gpt-5 api).

Table 4. Comparison of code generation systems on NeRF papers (in-depth comparisons shown in supplementary).

Paper	GPT-5 Thinking		ChatDev		MetaGPT		DeepCode		NERFIFY	
	Score	Trainable	Score	Trainable	Score	Trainable	Score	Trainable	Score	Trainable
KeyNeRF [18]	0.85	-	0.25	✗	0.30	✗	0.60	✗	1.00	✓
FastNeRF [6]	0.65	✗	0.21	✗	0.36	✗	0.81	✗	0.95	✓
Vanilla NeRF [16]	0.71	✓	0.48	✗	0.29	✗	0.53	✗	0.92	✓
Deblur-NeRF [15]	0.82	-	0.18	✗	0.42	✗	0.75	✗	1.00	✓
Average	0.76	-	0.28	✗	0.34	✗	0.67	✗	0.97	✓

GPT-5 thinking deep research achieves moderate semantic scores and produces trainable code for already implemented papers like Vanilla NeRF but misses papers whose

repositories are not available online. ChatDev and MetaGPT generate syntactically valid Python code but fail to capture NeRF-specific patterns. DeepCode performs better than general multi-agent systems, benefiting from paper-specific analysis, but still cannot produce trainable code.

Detailed comparison for each paper across these multi-agent pipeline using LLM as an evaluation are given below. The evaluation leverages language models as semantic judges, analyzing both mathematical correctness and algorithmic understanding. This automated assessment as shown in section 3 and section 4 examines whether implementations capture the conceptual intent beyond syntactic correctness. For instance, when evaluating Mip-NeRF’s integrated positional encoding, the semantic evaluator verifies not just the presence of the encoding formula but also its proper integration within the coarse-to-fine sampling pipeline.

### 3.1. KeyNeRF: Informative Rays Selection for Few-Shot Neural Radiance Fields

#### 3.1.1. Paper Overview

KeyNeRF addresses the critical challenge of few-shot novel view synthesis by introducing an intelligent camera and ray selection mechanism. The paper contributes two key innovations: first, a greedy baseline diversity algorithm for selecting the most informative camera views from limited input, and second, an entropy-based ray sampling strategy that prioritizes regions with higher visual information content. These techniques enable high-quality NeRF reconstruction from as few as 10-20 input views, where standard NeRF methods typically fail. The paper demonstrates that strategic data selection can be more effective than architectural modifications for few-shot scenarios.

#### 3.1.2. Implementation Overview

Table 5. KeyNeRF Implementation Summary Across Multi-agent Systems

Aspect	NERFIFY	GPT-5-thinking	DeepCode	MetaGPT	ChatDev
Lines of Code	479	507	1615	412	385
File Organization	Plugin	Modular	Multi-module	Multi-file	Multi-file
View Selection	✓	✓	Partial	×	Attempted
Entropy Sampling	✓	✓	✓	×	Simplified
Nerfstudio Integration	✓	✓	×	×	×
Trainable	✓	✓	×	×	×

#### 3.1.3. Novel Components

Table 6. Novel Components in KeyNeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Greedy baseline diversity for camera selection	0.20
C2	Integer linear programming (ILP) for coverage	0.15
C3	Local entropy computation for ray importance	0.15
C4	Entropy-weighted ray sampling	0.15
C5	Mixed sampling strategy (entropy + uniform)	0.10
C6	Nerfstudio datamanager integration	0.10
C7	Per-camera entropy caching	0.05
C8	Appearance embedding regularization	0.05
C9	Camera frustum visibility checks	0.03
C10	Grayscale conversion for entropy	0.02

#### 3.1.4. Quantitative Metrics

Table 7. KeyNeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	1.00	1.00
GPT-5 thinking	0.85	0.10	0.05	0.93	0.85
DeepCode	0.60	0.15	0.25	0.70	0.60
MetaGPT	0.30	0.40	0.30	0.45	0.30
ChatDev	0.25	0.35	0.40	0.35	0.25

#### 3.1.5. Component-by-Component Analysis

**Component C1: Greedy Baseline Diversity** The paper specifies a greedy algorithm that iteratively selects cameras maximizing the minimum baseline angle to already-selected views, ensuring diverse viewpoint coverage for robust reconstruction.

```

1 # \nerfify\ : Full greedy baseline with angle
  computation
2 def _select_cameras_greedy_baseline(self, K: int) ->
  List[int]:
3     z_axes = F.normalize(c2w_t[:, :3, 2], dim=-1) #
  optical axes
4     dots = z_axes @ z_axes.t()
5     angles = torch.arccos(torch.clamp(dots, -1.0,
  1.0))
6     init_idx = int(torch.argmax(torch.nanmean(angles,
  dim=1)).item())
7     selected = [init_idx]
8     while len(selected) < min(K, N) and remaining:
9         sub = angles.index_select(0, rem_idx_t).
  index_select(1, sel_idx_t)
10        min_to_sel = torch.min(sub, dim=1).values
11        pick = rem_list[int(torch.argmax(min_to_sel).
  item())]
12        selected.append(pick)

```

Listing 1. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Complete algorithm with proper angle
  computation
2 def greedy_baseline_ranking(coverage_set, num_cams, B
  , K=None):
3     B = compute_baseline_matrix(cameras) # Proper
  angle matrix
4     select = coverage_set.copy()
5     while remain and len(select) < K:
6         sub_mat = B[np.ix_(remain, select)]

```

```

7     scores = np.min(sub_mat, axis=1)
8     best_cam = remain[int(np.argmax(scores))]
9     select.append(best_cam)
10 # CORRECT: Full implementation, just offline preprocessing

```

Listing 2. C1 Implementation: GPT-5 (Score: 0.9)

```

1 # ChatDev: Placeholder without actual implementation
2 def select_views(dataset):
3     num_cameras = len(dataset.poses)
4     visibility_matrix = compute_visibility_matrix(
5         dataset)
6     select_cams = solve_ilp(visibility_matrix,
7                             num_cameras)
8     return select_cams
9 # ERROR: is_visible() always returns True

```

Listing 3. C1 Implementation: ChatDev (Score: 0.3)

**Analysis:** NERFIFY implements the complete greedy baseline algorithm integrated into the datamanager. GPT-5 provides an equally correct and complete implementation with proper angle computation, using offline preprocessing which is a valid architectural choice. ChatDev and MetaGPT lack proper implementation, while DeepCode focuses on different selection criteria.

**Component C2: Integer Linear Programming for Coverage** The paper formulates camera selection as an ILP problem minimizing the number of cameras while ensuring all scene points are covered.

```

1 # \nerfify\ : Complete ILP with OR-Tools integration
2 def solve_coverage_ilp(A: np.ndarray):
3     solver = pywraplp.Solver.CreateSolver("SCIP")
4     x = [solver.BoolVar(f"x_{i}") for i in range(
5         num_cams)]
6     for j in range(num_pts):
7         ct = solver.RowConstraint(1, solver.infinity
8             (), f"pt_{j}")
9         for i in range(num_cams):
10            if A[i, j] == 1:
11                ct.SetCoefficient(x[i], 1.0)
12 solver.Minimize(solver.Sum(x[i] for i in range(
13     num_cams)))

```

Listing 4. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Identical ILP formulation with complete
2 visibility matrix
3 def solve_coverage_ilp(A: np.ndarray):
4     solver = pywraplp.Solver.CreateSolver("SCIP")
5     x = [solver.BoolVar(f"x_{i}") for i in range(
6         num_cams)]
7     for j in range(num_pts):
8         ct = solver.RowConstraint(1, solver.infinity
9             (), f"pt_{j}")
10        for i in range(num_cams):
11            if A[i, j] == 1:
12                ct.SetCoefficient(x[i], 1.0)
13 solver.Minimize(solver.Sum(x[i] for i in range(
14     num_cams)))
15 status = solver.Solve()
16 select = [i for i in range(num_cams) if x[i].
17     solution_value() > 0.5]
18 # CORRECT: Complete implementation with proper error
19 handling

```

Listing 5. C2 Implementation: GPT-5 (Score: 1.0)

```

1 # ChatDev: Attempts ILP but incomplete visibility
2 def solve_ilp(visibility_matrix, num_cameras):
3     solver = pywraplp.Solver.CreateSolver('SCIP')
4     x = [solver.BoolVar(f'x_{i}') for i in range(
5         num_cameras)]
6     solver.Minimize(solver.Sum(x))
7     for j in range(visibility_matrix.shape[1]):
8         solver.Add(solver.Sum(x[i] *
9             visibility_matrix[i, j]
10                for i in range(num_cameras)) >= 1)
11 # ERROR: visibility_matrix always zeros due to
12 placeholder

```

Listing 6. C2 Implementation: ChatDev (Score: 0.5)

**Analysis:** Both NERFIFY and GPT-5 provide complete, correct ILP formulations with proper OR-Tools integration. The only difference is architectural choice: NERFIFY integrates it into the datamanager while GPT-5 uses offline preprocessing. ChatDev attempts ILP but fails due to broken visibility computation. MetaGPT and DeepCode omit this component entirely.

**Component C3: Local Entropy Computation** The paper computes local entropy using sliding windows over grayscale images to identify information-rich regions.

```

1 # \nerfify\ : Proper windowed entropy with histogram
2 binning
3 gray = self._rgb_to_gray(image.clamp(0.0, 1.0))
4 q = torch.clamp((gray * (bins - 1)).long(), 0, bins -
5     1)
6 one_hot = F.one_hot(q.squeeze(0), num_classes=bins)
7 kernel = torch.ones(1, 1, win, win) / (win * win)
8 hist_local = F.conv2d(one_hot.float(), kernel,
9     padding=pad)
10 probs = hist_local / (hist_local.sum(dim=1, keepdim=
11     True) + eps)
12 entropy = -(probs * torch.log(probs + eps)).sum(dim
13     =1)

```

Listing 7. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # DeepCode: Simplified entropy without windowing
2 def compute_entropy(image_batch):
3     entropy_maps = []
4     for img in image_batch:
5         gray = cv2.cvtColor(img, cv2.COLOR_RGB2GRAY)
6         hist = cv2.calcHist([gray], [0], None, [256],
7             [0,256])
8         hist = hist / hist.sum()
9         entropy = -np.sum(hist * np.log(hist + 1e-7))
10     entropy_maps.append(entropy)
11 # ISSUE: Global entropy, not local per-pixel

```

Listing 8. C3 Implementation: DeepCode (Score: 0.7)

**Analysis:** NERFIFY implements proper local windowed entropy computation. DeepCode computes global entropy missing spatial locality. Others either omit or provide placeholder implementations.

**Component C4: Entropy-Weighted Ray Sampling** The paper samples rays proportionally to their local entropy values, focusing training on high-information regions.

```

1 # \nerfify\ : Proper entropy-weighted sampling with
  mixing
2 entropy_probs = self._compute_entropy_probs(image)
3 mix = self.config.entropy_mix_ratio
4 n_entropy = int(B * mix)
5 n_uniform = B - n_entropy
6 entropy_idx = torch.multinomial(entropy_probs,
  n_entropy)
7 uniform_idx = torch.randint(0, H * W, (n_uniform,))
8 indices = torch.cat([entropy_idx, uniform_idx])

```

Listing 9. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Proper entropy sampling with Nerfstudio
  integration
2 e = self.entropy_maps[pose_idx_t].clamp(min=1e-6)
3 probs = e.view(-1) / e.view(-1).sum()
4 B_entropy = B // 2 # Mix ratio
5 B_uniform = B - B_entropy
6 flat_idx_ent = torch.multinomial(probs, B_entropy,
  replacement=True)
7 flat_idx_uni = torch.randint(0, H*W, (B_uniform,),
  device=device)
8 ray_bundle = self.ray_generator(ray_indices) #
  Nerfstudio integration
9 # CORRECT: Full implementation with proper framework
  usage

```

Listing 10. C4 Implementation: GPT-5 (Score: 0.9)

```

1 # ChatDev: Incorrect entropy usage
2 entropy_probs = compute_entropy(batch['image'])
3 rays_idx = np.random.choice(len(rays), size=
  batch_size, p=entropy_probs)
4 # CRITICAL ERROR: entropy_probs wrong shape and not
  normalized

```

Listing 11. C4 Implementation: ChatDev (Score: 0.2)

**Analysis:** NERFIFY and GPT-5 thinking both implement proper entropy-weighted sampling with correct probability computation and mixing strategy. GPT-5’s implementation demonstrates sophisticated understanding of Nerfstudio’s RayGenerator integration. Other baselines fail to properly implement this crucial component.

### 3.1.6. Scoring Analysis

Table 8. KeyNeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	Weighted Avg <sub>L1LM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.00
GPT-5 thinking	0.9	1.0	0.8	0.9	0.9	0.8	0.8	0.0	0.9	0.8	0.85
DeepCode	0.6	0.0	0.7	0.6	0.5	0.0	0.8	0.7	0.9	0.8	0.60
MetaGPT	0.0	0.0	0.0	0.0	0.3	0.0	0.0	0.8	0.7	0.5	0.30
ChatDev	0.3	0.5	0.2	0.2	0.1	0.0	0.0	0.0	0.4	0.3	0.25

### 3.1.7. Why Baselines Fail Despite Component Scores

#### GPT-5 thinking (Score: 8.50 components, 85% trainable)

- **Strengths:** Comprehensive understanding of both KeyNeRF algorithms and Nerfstudio integration, provides working training pipeline with RayGenerator
- **Minor Issues:**
  - Uses offline preprocessing rather than integrated data-manager approach

- Custom trainer instead of full plugin structure (but provides migration path)
- **Result:** Near-production ready implementation requiring minimal adaptation

#### DeepCode (Score: 6.00 components, 0% trainable)

- **Strengths:** Comprehensive standalone implementation with evaluation pipeline
- **Fatal Issues:**
  - No Nerfstudio integration whatsoever
  - Custom NeRF implementation incompatible with modern frameworks
  - Missing critical view selection algorithms
- **Result:** Standalone code that cannot leverage Nerfstudio infrastructure

#### MetaGPT (Score: 3.00 components, 0% trainable)

- **Strengths:** Clean code structure with proper configuration management
- **Fatal Issues:**
  - No KeyNeRF-specific components implemented
  - Basic NeRF without view or ray selection
  - Missing positional encoding and volume rendering
- **Result:** Generic NeRF attempt missing all KeyNeRF innovations

#### ChatDev (Score: 2.50 components, 0% trainable)

- **Strengths:** Attempts ILP formulation and mentions entropy
- **Fatal Issues:**
  - Placeholder functions return dummy values
  - GUI components inappropriate for NeRF training
  - Fundamental misunderstanding of ray sampling
- **Result:** Non-functional code with critical algorithmic errors

### 3.1.8. Hyperparameter Fidelity

Table 9. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5-thinking	DeepCode	MetaGPT	ChatDev
Selected cameras K	16	✓	✓	×	×	×
Entropy window size	9	✓	✓	×	×	×
Entropy bins	16	✓	✓	256	×	256
Mix ratio	0.5	✓	✓	×	×	×
Learning rate	5e-4	✓	✓	1e-3	✓	1e-3
Batch size	4096	✓	✓	1024	16	32
<b>W Score</b>	–	1.00	0.93	0.17	0.17	0.08

### 3.1.9. Conclusion

All baseline systems attempt to implement KeyNeRF with varying levels of sophistication. GPT-5 thinking demonstrates deep understanding of both the algorithms and framework integration, providing a near-complete solution with proper Nerfstudio integration through RayGenerator and NerfactoModel. DeepCode provides extensive standalone code



but completely misses Nerfstudio integration. MetaGPT produces clean structure but implements vanilla NeRF without any KeyNeRF innovations. ChatDev demonstrates fundamental misunderstandings with placeholder functions and inappropriate GUI components. The critical difference is framework integration quality: GPT-5 thinking achieves 85% trainable code with minimal adaptation needed, while ChatDev, MetaGPT and DeepCode produce 0% trainable implementations. Only NERFIFY achieves 100% immediate trainability as a complete Nerfstudio plugin, though GPT-5 thinking comes remarkably close with its sophisticated understanding of both algorithmic details and framework requirements.

### 3.2. FastNeRF: High-Fidelity Neural Rendering at 200 FPS

#### 3.2.1. Paper Overview

FastNeRF introduces a factorized neural radiance field architecture that decomposes the traditional monolithic MLP into position and direction dependent networks. The key innovation lies in representing radiance as a factorized inner product between a deep radiance map from the position network and directional weights from the viewing direction network. This factorization enables aggressive caching strategies that achieve 200 FPS rendering while maintaining visual quality comparable to vanilla NeRF. The paper fundamentally changes how neural radiance fields balance quality versus speed, demonstrating that architectural factorization rather than model compression provides the path to real-time neural rendering.

#### 3.2.2. Implementation Overview

Table 10. FastNeRF Implementation Summary Across Multi-Agent Systems

Aspect	NERFIFY	GPT-5-thinking	DeepCode	MetaGPT	ChatDev
Lines of Code	565	280	2915	420	385
File Organization	Plugin	Modular	Comprehensive	Multi-file	Multi-file
Factorized MLPs	✓	✓	✓	✓	Simplified
Deep Radiance Map	✓	✓	✓	Partial	×
Nerfstudio Integration	✓	Partial	×	×	×
Cache Infrastructure	Ready	Sketched	✓	×	×
Trainable	✓	×	×	×	×

#### 3.2.3. Novel Components

#### 3.2.4. Quantitative Metrics

#### 3.2.5. Component-by-Component Analysis

**Component C1: Position MLP Architecture** The position network processes 3D coordinates to produce density and a deep radiance map. The paper specifies an 8-layer MLP with 384 hidden units, outputting  $\sigma$  and 3D values representing the factorized radiance components.

```
1 # \nerfify\ : Full factorized position network with
   correct dimensions
2 self.pos_mlp = _MLP(
```

Table 11. Novel Components in FastNeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Position MLP: $F_{pos}(p) \rightarrow (\sigma, u, v, w)$	0.15
C2	Direction MLP: $F_{dir}(d) \rightarrow \beta$	0.12
C3	Deep radiance map factorization $(u, v, w) \in \mathbb{R}^{D \times 3}$	0.13
C4	Directional weights $\beta \in \mathbb{R}^D$	0.10
C5	Inner product color computation: $c = \sum_i \beta_i(u_i, v_i, w_i)$	0.12
C6	Positional encoding with $L = 10$ for positions	0.08
C7	Directional encoding with $L = 1$ for viewing directions	0.08
C8	Cache grid structure for $F_{pos}$ outputs	0.10
C9	Spherical cache for $F_{dir}$ outputs	0.08
C10	Training hyperparameters (384/128 widths, 8/4 layers)	0.04

Table 12. FastNeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	0.95	0.05	0.00	1.00	0.95
GPT-5 thinking	0.72	0.18	0.10	0.90	0.65
DeepCode	0.85	0.10	0.05	0.95	0.81
MetaGPT	0.45	0.35	0.20	0.80	0.36
ChatDev	0.30	0.40	0.30	0.70	0.21

```
3 in_dim=pos_in, # 63 dims after encoding
4 num_layers=config.pos_num_layers, # 8 layers
5 layer_width=config.pos_layer_width, # 384 units
6 out_dim=1 + 3 * D, # sigma + (u,v,w)
7 )
8 def forward(self, positions):
9     encoded = self.pos_encoding(positions)
10    out = self.pos_mlp(encoded)
11    sigma = self.softplus(out[..., :1])
12    uvw = out[..., 1:].reshape(*out.shape[:-1], D, 3)
13    return sigma, uvw
```

Listing 12. C1 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Correct architecture but missing
   implementation details
2 self.pos_mlp = MLP(
3     in_dim=pos_in_dim,
4     num_layers=config.pos_mlp_layers,
5     layer_width=config.pos_mlp_width,
6     out_dim=pos_out_dim, # 1 + 3*D
7     skip_connections=[4], # similar to NeRF
8 )
9 # ERROR: Implementation incomplete,
   populate_nerf_helpers not provided
```

Listing 13. C1 Implementation: GPT-5 (Score: 0.9)

```
1 # DeepCode: Comprehensive with cache-aware design
2 class FastNeRFPosMLP(nn.Module):
3     def __init__(self, D=8, W=384, input_ch=63,
4                 n_components=8):
5         super().__init__()
6         self.pts_linears = nn.ModuleList(
7             [nn.Linear(input_ch, W)] +
8             [nn.Linear(W, W) for _ in range(D-1)]
9         )
10        self.output_linear = nn.Linear(W, 1 +
11                                       n_components * 3)
12        # CORRECT: Full implementation with proper
13        factorization
```

Listing 14. C1 Implementation: DeepCode (Score: 0.95)

```

1 # MetaGPT: Simplified architecture missing skip
  connections
2 class PositionDependentNetwork(nn.Module):
3     def __init__(self, hidden_dim=384, output_dim=27)
4         :
5         self.layers = nn.ModuleList()
6         for _ in range(num_layers - 2):
7             self.layers.append(nn.Linear(hidden_dim,
8             hidden_dim))
9     # ERROR: No positional encoding integration

```

Listing 15. C1 Implementation: MetaGPT (Score: 0.6)

```

1 # ChatDev: Incorrect understanding of factorization
2 class FPos(nn.Module):
3     def forward(self, p):
4         x = F.relu(self.fc1(p_encoded))
5         output = self.fc3(x)
6         return output[:, 0], output[:, 1:4] # Wrong:
7         returns (u,v,w) as 3D
8     # CRITICAL ERROR: Misunderstands deep radiance
9     map dimensions

```

Listing 16. C1 Implementation: ChatDev (Score: 0.3)

**Analysis:** NERFIFY correctly implements the full factorized architecture with proper dimensionality. GPT-5 understands the structure but leaves critical implementation details as placeholders. DeepCode provides the most comprehensive implementation outside of NERFIFY but lacks framework integration. MetaGPT simplifies the architecture while ChatDev fundamentally misunderstands the factorization concept.

**Component C2: Direction MLP Architecture** The direction network processes viewing directions to produce weights for combining the deep radiance map components. The paper specifies a 4-layer MLP with 128 hidden units.

```

1 # \nerfify\ : Correct lightweight direction network
2 self.dir_mlp = _MLP(
3     in_dim=dir_in, # 9 dims with L=1 encoding
4     num_layers=config.dir_num_layers, # 4 layers
5     layer_width=config.dir_layer_width, # 128 units
6     out_dim=D, # beta vector
7 )

```

Listing 17. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # DeepCode: Proper direction network with activation
2 class FastNeRFDDirMLP(nn.Module):
3     def __init__(self, D=4, W=128, dir_enc_L=1,
4         n_components=8):
5         self.dir_linears = nn.ModuleList([nn.Linear(
6             input_ch, W)] +
7             [nn.Linear(W,
8             W) for _ in range(D-1)])
9         self.output_linear = nn.Linear(W,
10         n_components)

```

Listing 18. C2 Implementation: DeepCode (Score: 0.9)

**Analysis:** Both NERFIFY and DeepCode correctly implement the lightweight direction network. GPT-5 provides the structure, while MetaGPT and ChatDev misunderstand the role of directional processing in the factorization.

**Component C3: Deep Radiance Map Factorization** The core innovation of FastNeRF lies in representing radiance as a tensor product between position-dependent and direction-dependent components.

```

1 # \nerfify\ : Proper tensor factorization
2 uvw = pos_out[..., 1:].reshape(*pos_out.shape[:-1],
3     self.D, 3)
4 beta = self.dir_mlp(encoded_dirs) # [batch, D]
5 rgb = torch.sum(beta[..., :, None] * uvw, dim=-2) #
6     Inner product

```

Listing 19. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # ChatDev: Complete misunderstanding of factorization
2 return output[:, 0], output[:, 1:4] # Returns 3D
3 vector instead of Dx3
4 # CRITICAL ERROR: No tensor product, treats as
5 regular RGB output

```

Listing 20. C3 Implementation: ChatDev (Score: 0.0)

**Analysis:** NERFIFY and DeepCode correctly implement the tensor factorization that defines FastNeRF. GPT-5 shows understanding but lacks implementation. MetaGPT attempts factorization with errors, while ChatDev completely misses this fundamental concept.

**Component C5: Inner Product Color Computation** The final RGB color emerges from the weighted sum of deep radiance map components using directional weights.

```

1 # \nerfify\ : Mathematically correct inner product
2 rgb = torch.sum(beta[..., :, None] * uvw, dim=-2)
3 rgb = self.sigmoid(rgb) # Activation for valid color
4 range

```

Listing 21. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # MetaGPT: Attempts inner product with shape errors
2 color = torch.sum(weights.unsqueeze(-1) *
3     radiance_map.view(-1, 8, 3), dim=1)
4 # ERROR: Incorrect tensor reshaping, breaks with
5 batch processing

```

Listing 22. C5 Implementation: MetaGPT (Score: 0.5)

**Analysis:** The inner product computation separates successful implementations from failures. NERFIFY handles all tensor operations correctly, while baselines struggle with proper broadcasting and dimension management.

**Component C8: Cache Grid Structure** FastNeRF's speed derives from caching position network outputs in a 3D grid for inference.

```

1 # DeepCode: Complete cache implementation
2 def build_pos_cache(f_pos, bounds, grid_res,
3     n_components=8):
4     positions, grid_shape = create_position_grid(
5         bounds, grid_res)
6     outputs = batch_eval_mlp_on_grid(f_pos, positions)
7     cache = outputs.reshape(*grid_shape, 1 +
8         n_components * 3)
9     return cache # [k, k, k, 1+D*3] tensor

```

Listing 23. C8 Implementation: DeepCode (Score: 1.0)

```

1 # GPT-5: Acknowledges caching but provides no
  implementation
2 # (Optional) How to add **caching** like the paper
3 # 1. **Offline cache builder** script:
4 #    * Define a world-space bounding box and
      resolution k
5 # ERROR: Only provides instructions, no actual code

```

Listing 24. C8 Implementation: GPT-5 (Score: 0.3)

**Analysis:** DeepCode excels at cache implementation with comprehensive utilities for grid construction and interpolation. NERFIFY provides the architecture ready for caching. GPT-5 acknowledges the concept but defers implementation, while MetaGPT and ChatDev ignore caching entirely.

### 3.2.6. Scoring Analysis

Table 13. FastNeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	Weighted AvgLLM
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.7	0.6	1.0	0.94
GPT-5 thinking	0.9	0.8	0.8	0.8	0.7	0.9	0.8	0.3	0.2	0.5	0.67
DeepCode	0.95	0.9	0.95	0.9	0.9	0.9	0.8	1.0	0.9	0.3	0.85
MetaGPT	0.6	0.5	0.4	0.5	0.5	0.7	0.6	0.0	0.0	0.3	0.41
ChatDev	0.3	0.3	0.0	0.2	0.1	0.8	0.7	0.0	0.0	0.2	0.26

### 3.2.7. Why Baselines Fail Despite Component Scores

#### GPT-5 Extended Thinking (Score: 0.67 components, 0% trainable)

- **Strengths:** Demonstrates deep understanding of FastNeRF architecture and provides correct Nerfstudio integration patterns
- **Fatal Issues:**
  - Leaves critical implementation as placeholder: `raise NotImplementedError`
  - Missing essential `populate_modules` implementation
  - No actual training loop connection
- **Result:** Code serves as documentation rather than executable implementation

#### DeepCode (Score: 0.85 components, 0% trainable)

- **Strengths:** Most comprehensive implementation with full caching infrastructure and extensive testing
- **Fatal Issues:**
  - Complete absence of Nerfstudio framework integration
  - No training pipeline or data loading infrastructure
  - Implements components in isolation without system integration
- **Result:** Excellent reference implementation that cannot train without significant engineering effort

#### MetaGPT (Score: 0.41 components, 0% trainable)

- **Strengths:** Attempts modular organization with separate training and rendering components
- **Fatal Issues:**

- Fundamental misunderstanding of tensor dimensions in factorization
- Missing ray sampling and volume rendering pipeline
- Placeholder data loading with no actual dataset interface
- **Result:** Structurally incomplete implementation that fails at runtime

#### ChatDev (Score: 0.26 components, 0% trainable)

- **Strengths:** Creates organized file structure with GUI interface
- **Fatal Issues:**
  - Complete misunderstanding of deep radiance map factorization
  - Implements Tkinter GUI instead of neural rendering pipeline
  - Missing entire volume rendering and ray marching infrastructure
- **Result:** Application framework without actual FastNeRF implementation

### 3.2.8. Hyperparameter Fidelity

Table 14. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	DeepCode	MetaGPT	ChatDev
Position MLP layers	8	✓	✓	✓	Partial	×
Position MLP width	384	✓	✓	✓	✓	×
Direction MLP layers	4	✓	✓	✓	✓	×
Direction MLP width	128	✓	✓	✓	✓	×
Factor dimension D	8	✓	✓	Variable	Fixed	Wrong
Position encoding L	10	✓	✓	Variable	✓	✓
Direction encoding L	1	✓	Variable	Variable	×	×
Learning rate	5e-4	✓	✓	N/A	✓	×
<b>W Score</b>	–	1.00	0.88	0.63	0.50	0.25

### 3.2.9. Conclusion

All baseline systems attempt to implement FastNeRF with varying degrees of sophistication and understanding. GPT-5 demonstrates theoretical mastery but delivers skeleton code requiring substantial completion. DeepCode produces the most comprehensive component implementation with sophisticated caching utilities, yet fails to integrate with any training framework. MetaGPT captures the high-level structure but makes critical errors in tensor operations that prevent execution. ChatDev fundamentally misunderstands the paper, building a GUI application instead of implementing the core factorization algorithm. Despite DeepCode achieving 85% component coverage and GPT-5 showing 67% understanding, all baselines produce 0% trainable code due to missing framework integration, incomplete implementations, or fundamental architectural errors. Only NERFIFY delivers immediately trainable code by properly implementing the factorized architecture within the Nerfstudio framework, achieving 100% executability where all others fail completely.

### 3.3. NeRF: Neural Radiance Fields for View Synthesis

#### 3.3.1. Paper Overview

The original NeRF paper by Mildenhall et al. (2020) introduced neural radiance fields for photorealistic novel view synthesis. The method represents scenes as continuous 5D functions mapping 3D coordinates and 2D viewing directions to volume density and color, rendered using classical volume rendering techniques with hierarchical sampling for efficiency.

#### 3.3.2. Implementation Overview

Table 15. NeRF Implementation Summary Across Multi-Agent Systems

Aspect	NERFIFY	GPT-5-thinking	DeepCode	MetaGPT	ChatDev
Lines of Code	387	156	412	98	524
File Organization	Plugin	Instructions	Multi-file	Multi-file	Multi-file
Positional Encoding	✓	✓	×	Partial	Attempted
MLP Architecture	✓	✓	Simplified	Partial	Simplified
Hierarchical Sampling	✓	✓	×	×	Attempted
Volume Rendering	✓	✓	Partial	×	Partial
Trainable	✓	Partial	×	×	×

#### 3.3.3. Novel Components

Table 16. Novel Components in NeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Positional encoding (L=10 for xyz, L=4 for direction)	0.15
C2	8-layer MLP with skip connection at layer 4	0.15
C3	View-dependent RGB prediction	0.10
C4	Hierarchical volume sampling (coarse + fine)	0.20
C5	Alpha compositing with transmittance	0.15
C6	Stratified sampling along rays	0.10
C7	NDC ray parameterization	0.05
C8	MSE photometric loss	0.05
C9	Adam optimizer with exponential decay	0.05

#### 3.3.4. Quantitative Metrics

Table 17. NeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	0.95	0.02	0.03	0.97	9.2
GPT-5 thinking	0.72	0.15	0.13	0.79	7.1
ChatDev	0.35	0.42	0.23	0.51	4.8
MetaGPT	0.18	0.27	0.55	0.31	2.9
DeepCode	0.42	0.38	0.20	0.56	5.3

#### 3.3.5. Component-by-Component Analysis

**Component C1: Positional Encoding** The paper specifies sinusoidal positional encoding with  $L = 10$  for positions and  $L = 4$  for viewing

directions, encoding input coordinates as  $\gamma(p) = [p, \sin(2^0\pi p), \cos(2^0\pi p), \dots, \sin(2^{L-1}\pi p), \cos(2^{L-1}\pi p)]$ .

```
1 # \nerfify\ : Proper Nerfstudio positional encoding
   configuration
2 position_encoding_config=NeRFModelConfig.
   position_encoding_config.__class__(
3     num_frequencies=10
4 ),
5 direction_encoding_config=NeRFModelConfig.
   direction_encoding_config.__class__(
6     num_frequencies=4
7 ),
8 # CORRECT: Exact paper specification with L=10 and L
   =4
```

Listing 25. C1 Implementation: NERFIFY (Score: 1.0)

```
1 # ChatDev: Custom positional encoding implementation
2 class PositionalEncoding:
3     def __init__(self, L):
4         self.L = L
5     def encode(self, x):
6         encoded = []
7         for i in range(self.L):
8             encoded.append(torch.sin(2 ** i * torch.
9                 pi * x))
10            encoded.append(torch.cos(2 ** i * torch.
11                pi * x))
12        return torch.cat(encoded, dim=-1)
13 # ERROR: No differentiation between position/
   direction encoding levels
```

Listing 26. C1 Implementation: ChatDev (Score: 0.6)

```
1 # MetaGPT: Placeholder implementation
2 class PositionalEncoding(nn.Module):
3     def __init__(self, L=10, d_model=256):
4         super(PositionalEncoding, self).__init__()
5         self.L = L
6         self.d_model = d_model
7     def forward(self, x):
8         # Implement positional encoding as described
9         in the paper
10        pass
11 # CRITICAL ERROR: Not implemented, just placeholder
```

Listing 27. C1 Implementation: MetaGPT (Score: 0.3)

**Analysis:** NERFIFY correctly configures both position and direction encoding with proper frequency counts. ChatDev implements the mathematical formula but misses the crucial distinction between position and direction encoding levels. MetaGPT provides only placeholders.

**Component C2: MLP Architecture** The paper specifies an 8-layer MLP with 256 hidden units and a skip connection at layer 4.

```
1 # \nerfify\ : Correct architecture configuration
2 field_num_layers=8,
3 field_layer_width=256,
4 field_heads_layer_width=128,
5 field_skip_connections=(4,),
6 # CORRECT: Exact paper specification with skip at
   layer 4
```

Listing 28. C2 Implementation: NERFIFY (Score: 1.0)



```

1 # ChatDev: Simplified MLP without skip connection
2 self.fc_layers = nn.ModuleList([
3     nn.Linear(30, 256) for _ in range(8)
4 ])
5 # ERROR: Missing skip connection at layer 4
6 # ERROR: Incorrect input dimension (should be 63 for
   L=10)

```

Listing 29. C2 Implementation: ChatDev (Score: 0.5)

**Analysis:** NERFIFY properly configures the 8-layer architecture with skip connections. ChatDev creates 8 layers but misses the critical skip connection and has incorrect input dimensions.

**Component C3: View-Dependent RGB Prediction** The paper predicts RGB values conditioned on both position features and viewing direction.

```

1 # \nerfify\ : Inherits view-dependent color from
   NerfactoField
2 # View direction properly encoded and concatenated
   with features
3 # CORRECT: Full view-dependent RGB implementation via
   parent class

```

Listing 30. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # ChatDev: View-dependent color prediction
2 self.color_layer = nn.Sequential(
3     nn.Linear(256 + 8, 128),
4     nn.ReLU(),
5     nn.Linear(128, 3),
6     nn.Sigmoid()
7 )
8 # ISSUE: Wrong concatenation dimension (should be 256
   + 27 for L=4)

```

Listing 31. C3 Implementation: ChatDev (Score: 0.7)

**Analysis:** NERFIFY correctly implements view-dependent RGB through Nerfstudio’s field architecture. ChatDev attempts view dependence but miscalculates encoded direction dimensions.

**Component C4: Hierarchical Volume Sampling** The paper uses coarse-to-fine sampling with 64 coarse and 128 fine samples.

```

1 # \nerfify\ : Proper hierarchical sampling
   configuration
2 num_coarse_samples=64,
3 num_importance_samples=128,
4 # CORRECT: Exact paper specification for two-stage
   sampling

```

Listing 32. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # ChatDev: Attempted hierarchical sampling
2 def sample_fine(self, coarse_samples,
   coarse_densities):
3     weights = coarse_densities / torch.sum(
   coarse_densities, dim=-1, keepdim=True)
4     fine_indices = torch.multinomial(weights, self.
   N_f, replacement=True)
5     fine_samples = coarse_samples['xyz'][fine_indices
   ]

```

```

6 # CRITICAL ERROR: Wrong importance sampling (should
   use CDF inverse)

```

Listing 33. C4 Implementation: ChatDev (Score: 0.4)

**Analysis:** NERFIFY configures proper coarse and fine sampling counts. ChatDev attempts importance sampling but implements it incorrectly, using simple multinomial instead of inverse CDF sampling.

### 3.3.6. Scoring Analysis

Table 18. NeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	0.9	1.0	1.0	0.97
GPT-5 thinking	1.0	1.0	1.0	1.0	1.0	0.8	0.0	0.8	0.9	0.79
ChatDev	0.6	0.5	0.7	0.4	0.6	0.5	0.0	0.8	0.5	0.51
MetaGPT	0.3	0.4	0.3	0.0	0.2	0.3	0.0	0.6	0.7	0.31
DeepCode	0.0	0.6	0.8	0.2	0.7	0.6	0.0	0.9	0.8	0.56

### 3.3.7. Why Baselines Fail Despite Component Scores

**GPT-5 thinking (Score: 0.79 components, 50% trainable)**

- **Strengths:** Correctly identifies Nerfstudio’s built-in implementation, provides proper configuration
- **Fatal Issues:**
  - Provides instructions rather than executable code
  - Requires manual CLI command construction
- **Result:** Works only when user correctly follows instructions

**ChatDev (Score: 0.51 components, 0% trainable)**

- **Strengths:** Attempts full implementation with GUI and modular structure
- **Fatal Issues:**
  - No Nerfstudio integration whatsoever
  - Incorrect importance sampling implementation
  - Missing ray generation logic
- **Result:** Code runs but cannot train on real data

**MetaGPT (Score: 0.31 components, 0% trainable)**

- **Strengths:** Proper project structure with configuration management
- **Fatal Issues:**
  - All core functions are placeholders with ‘pass’ statements
  - No actual implementation of any algorithm
- **Result:** Skeleton code that requires complete implementation

**DeepCode (Score: 0.56 components, 0% trainable)**

- **Strengths:** Attempts volume rendering and loss computation
- **Fatal Issues:**
  - Missing positional encoding entirely

- No data loading implementation
- Incompatible with any training framework
- **Result:** Partial implementation missing critical components

### 3.3.8. Hyperparameter Fidelity

Table 19. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5-thinking	MetaGPT	ChatDev	DeepCode
Coarse samples	64	✓	✓	×	✓	×
Fine samples	128	✓	✓	×	✓	×
Batch size	4096	✓	✓	✓	✓	×
Learning rate	5e-4	✓	✓	✓	✓	×
Position L	10	✓	✓	×	×	×
Direction L	4	✓	✓	×	×	×
<b>W Score</b>	–	1.00	1.00	0.33	0.67	0.00

### 3.3.9. Conclusion

All baseline systems attempt to implement vanilla NeRF with varying levels of sophistication. GPT-5 thinking correctly identifies that Nerfstudio already contains the implementation but provides instructions rather than code. ChatDev produces the most complete standalone attempt but lacks framework integration. MetaGPT generates only skeleton code with placeholders. DeepCode attempts key components but misses positional encoding entirely. While GPT-5 thinking achieves 50% trainability through instruction following for this well-known method already in Nerfstudio, ChatDev, MetaGPT, and DeepCode achieve 0% trainable implementations versus NERFIFY’s 100% immediately trainable plugin.

## 3.4. Deblur-NeRF: NeRF from Motion-Blurred Images

### 3.4.1. Paper Overview

Deblur-NeRF addresses the fundamental challenge of training Neural Radiance Fields from motion-blurred images by introducing the Deformable Sparse Kernel (DSK). The method models blur formation explicitly through learnable sparse ray distributions per pixel, enabling sharp novel view synthesis from blurry training data. The key innovation lies in jointly optimizing the underlying sharp NeRF representation while learning per-image blur kernels through differentiable rendering.

### 3.4.2. Implementation Overview

Table 20. Deblur-NeRF Implementation Summary Across Multi-Agent Systems

Aspect	NERFIFY	GPT-5-thinking	DeepCode	MetaGPT	ChatDev
Lines of Code	736	424	1487	285	156
File Organization	Plugin	Modular	Multi-file	Multi-file	Single
DSK Module	✓	✓	✓	Partial	×
Nerfstudio Integration	✓	✓	×	×	×
Gamma Correction	✓	✓	✓	×	Simplified
Alignment Loss	✓	✓	✓	×	×
Trainable	✓	Partial	×	×	×

### 3.4.3. Novel Components

Table 21. Novel Components in Deblur-NeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Deformable Sparse Kernel (DSK) MLP	0.25
C2	Per-view embedding $l$ (32-dim)	0.10
C3	Origin optimization $\Delta o$	0.10
C4	Kernel blending in linear space	0.15
C5	Gamma correction $g(c) = c^{1/2.2}$	0.10
C6	Alignment loss $\mathcal{L}_{align}$	0.10
C7	Sharp rendering at inference	0.05
C8	Nerfstudio framework integration	0.15

### 3.4.4. Quantitative Metrics

Table 22. Deblur-NeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	0.85	0.10	0.05	1.00	1.00
GPT-5 thinking	0.70	0.15	0.15	0.82	0.82
DeepCode	0.65	0.20	0.15	0.75	0.75
MetaGPT	0.35	0.25	0.40	0.42	0.42
ChatDev	0.15	0.20	0.65	0.18	0.18

### 3.4.5. Component-by-Component Analysis

**Component C1: Deformable Sparse Kernel MLP** The DSK predicts per-pixel kernel parameters ( $\Delta o, \Delta q, w$ ) using an MLP conditioned on pixel coordinates and view embeddings.

```

1 # \nerfify\ : Full DSK with proper origin/direction
  deformations
2 class DeformableSparseKernel(nn.Module):
3     def forward(self, pixels_xy, camera_indices):
4         l = self.view_embed(camera_indices) # Per-
          view embedding
5         q_prime = pixels_xy.unsqueeze(1) + self.
          canonical_offsets
6         mlp_in = torch.cat([p_rep, q_prime, l_rep],
          dim=-1)
7         out = self.mlp(mlp_in) # Delta_o(3), Delta_q
          (2), w_logit(1)
8         delta_o = self.eps * out[:, 0:3].view(B, N,
          3)
9         delta_q = self.eps * out[:, 3:5].view(B, N,
          2)
10        weights = torch.softmax(out[:, 5].view(B, N),
          dim=1)
11        return delta_o, delta_q, weights, q_prime

```

Listing 34. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Comprehensive DSK with structured outputs
2 out = self.mlp(mlp_in) # (B * N, 6) -> delta_o(3),
          delta_q(2), w_logit(1)
3 delta_o = out[:, 0:3].view(B, N, 3)
4 delta_q = out[:, 3:5].view(B, N, 2)
5 # Scale deformations by eps to stay near identity
  early

```

```

6 delta_o = self.eps * delta_o
7 delta_q = self.eps * delta_q
8 weights = torch.softmax(w_logits, dim=1)
9 # ISSUE: Minor - different output ordering but
  functionally correct

```

Listing 35. C1 Implementation: GPT-5 (Score: 0.9)

```

1 # DeepCode: DSK with custom ray construction
2 rays_o, rays_d, kernel_weights, kernel_offsets,
  origin_offsets =
3   construct_kernel_rays(pixel_coords, intrinsics,
4   c2w, self.dsk,
5   view_ids, N=self.N_kernel)
6 # ERROR: External function call instead of integrated
  module

```

Listing 36. C1 Implementation: DeepCode (Score: 0.8)

```

1 # MetaGPT: Simplified DSK without proper deformations
2 output = self.output_layer(x)
3 delta_origin = output[:, :3]
4 delta_q = output[:, 3:5]
5 weights = F.softmax(output[:, 5:], dim=-1)
6 # ERROR: Missing eps scaling, wrong weight dimensions

```

Listing 37. C1 Implementation: MetaGPT (Score: 0.4)

```

1 # ChatDev: No DSK implementation
2 # MISSING: Entire DSK module absent from
  implementation

```

Listing 38. C1 Implementation: ChatDev (Score: 0.0)

**Analysis:** NERFIFY implements the complete DSK architecture with proper epsilon scaling and weight normalization. GPT-5 achieves near-perfect implementation with minor output ordering differences. DeepCode implements DSK but outside the model class. MetaGPT attempts DSK but misses critical scaling. ChatDev completely lacks DSK.

**Component C2: Per-view Embedding** View-specific 32-dimensional embeddings condition the DSK to learn per-image blur patterns.

```

1 # \nerfify\ : Proper view embedding with 32
  dimensions
2 self.view_embed = nn.Embedding(num_cameras,
  view_embed_dim)
3 l = self.view_embed(camera_indices) # (B, 32)

```

Listing 39. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct implementation
2 self.view_embed = nn.Embedding(num_cameras,
  view_embed_dim)
3 l = self.view_embed(camera_indices) # (B,
  view_embed_dim)

```

Listing 40. C2 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepCode: View embedding present but in trainer
2 if hasattr(self.dsk, 'view_embed'):
3   params += list(self.dsk.view_embed.parameters())
4 # ISSUE: Indirect reference, not shown in model
  definition

```

Listing 41. C2 Implementation: DeepCode (Score: 0.8)

```

1 # MetaGPT: No view embedding
2 # MISSING: View-specific conditioning absent

```

Listing 42. C2 Implementation: MetaGPT (Score: 0.0)

```

1 # ChatDev: No view embedding
2 # MISSING: No per-view conditioning

```

Listing 43. C2 Implementation: ChatDev (Score: 0.0)

**Analysis:** Both NERFIFY and GPT-5 correctly implement 32-dimensional view embeddings. DeepCode references them indirectly. MetaGPT and ChatDev lack this critical component entirely.

**Component C3: Origin Optimization** Ray origins are shifted by learned  $\Delta o$  to model camera motion blur.

```

1 # \nerfify\ : Full origin optimization with proper
  shifting
2 dsk_rays.origins = dsk_rays.origins + delta_o.view
  (-1, 3)
3 enable_origin_opt=True, # Config flag

```

Listing 44. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct origin shifting
2 dsk_rays.origins = dsk_rays.origins + delta_o.view
  (-1, 3)

```

Listing 45. C3 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepCode: Origin optimization with config control
2 self.kernel_origin_on = config.get('kernel_origin_on',
  True)
3 origin_offsets = construct_kernel_rays(...)[4]

```

Listing 46. C3 Implementation: DeepCode (Score: 0.9)

```

1 # MetaGPT: Delta origin computed but not applied
2 delta_origin = output[:, :3]
3 # ERROR: Never used to shift ray origins

```

Listing 47. C3 Implementation: MetaGPT (Score: 0.3)

```

1 # ChatDev: No origin optimization
2 # MISSING: No ray origin shifting

```

Listing 48. C3 Implementation: ChatDev (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 properly shift ray origins. DeepCode implements it with configuration control. MetaGPT computes deltas but never applies them. ChatDev lacks the feature.

**Component C4: Linear Space Blending** Colors are blended in linear RGB space before gamma correction.

```

1 # \nerfify\ : Proper linear blending then gamma
2 rgb_linear = (w[...], None) * rgb_fine_all).sum(dim=1)
3 rgb_blurry = self._apply_gamma(rgb_linear)

```

Listing 49. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct linear space blending
2 rgb_lin_blurry = torch.sum(weights_unsq * rgb_lin,
    dim=1)
3 rgb_blurry = self._gamma_correct(rgb_lin_blurry)

```

Listing 50. C4 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepCode: Blending with kernel weights
2 pred_blurry = torch.sum(all_rgb * kernel_weights,
    unsqueeze(-1), dim=1)
3 if self.gamma_on:
4     pred_blurry = gamma_correction(pred_blurry)

```

Listing 51. C4 Implementation: DeepCode (Score: 0.8)

```

1 # MetaGPT: No explicit kernel blending
2 # ERROR: Missing weighted sum over kernel points

```

Listing 52. C4 Implementation: MetaGPT (Score: 0.2)

```

1 # ChatDev: No kernel blending
2 # MISSING: No multi-ray blending

```

Listing 53. C4 Implementation: ChatDev (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 correctly implement linear blending followed by gamma correction. DeepCode has the right approach. MetaGPT and ChatDev miss this entirely.

**Component C5: Gamma Correction** Applies camera response function  $g(c) = c^{1/2.2}$ .

```

1 # \nerfify\ : Exact gamma value from paper
2 def _apply_gamma(self, c_lin: Tensor) -> Tensor:
3     return torch.clamp(c_lin, min=0.0) ** self.
    one_over_gamma # 1/2.2

```

Listing 54. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct gamma implementation
2 def _gamma_correct(self, c_lin: Tensor) -> Tensor:
3     return torch.clamp(c_lin, min=0.0) ** self.
    one_over_gamma

```

Listing 55. C5 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepCode: Proper gamma correction
2 def gamma_correction(color, gamma=2.2):
3     return color ** (1.0 / gamma)

```

Listing 56. C5 Implementation: DeepCode (Score: 1.0)

```

1 # MetaGPT: No gamma correction
2 # MISSING: No CRF modeling

```

Listing 57. C5 Implementation: MetaGPT (Score: 0.0)

```

1 # ChatDev: Different gamma implementation
2 def gamma_correction(c):
3     return torch.clamp(c ** (1 / 2.2), min=0)
4 # ERROR: Applied in wrong context (utils not model)

```

Listing 58. C5 Implementation: ChatDev (Score: 0.3)

**Analysis:** NERFIFY, GPT-5, and DeepCode correctly implement gamma=2.2. ChatDev has the function but uses it incorrectly. MetaGPT misses it entirely.

**Component C6: Alignment Loss** Regularizes kernel deformations:  $\mathcal{L}_{align} = ||q_0 - p||_2 + \lambda_o ||\Delta o_{q_0}||_2$ .

```

1 # \nerfify\ : Complete alignment loss
2 align = torch.norm(dq0, dim=-1).mean() + \
    self.config.lambda_o * torch.norm(do0, dim
    =-1).mean()
4 loss_dict["alignment_loss"] = self.config.
    lambda_align * align

```

Listing 59. C6 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Proper alignment with both terms
2 align_pos = torch.mean((q0 - p).pow(2).sum(dim=-1).
    sqrt())
3 align_origin = torch.mean(delta_o0.pow(2).sum(dim=-1)
    .sqrt())
4 align_loss = align_pos + self.config.lambda_origin *
    align_origin

```

Listing 60. C6 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepCode: Alignment loss function
2 def alignment_loss(pixel_coords, kernel_offsets,
    origin_offsets, lambda_o):
3     loss_align = alignment_loss(pixel_coords,
    kernel_offsets,
    origin_offsets,
    lambda_o=self.lambda_o)

```

Listing 61. C6 Implementation: DeepCode (Score: 0.9)

```

1 # MetaGPT: No alignment loss
2 # MISSING: No regularization on deformations

```

Listing 62. C6 Implementation: MetaGPT (Score: 0.0)

```

1 # ChatDev: No alignment loss
2 # MISSING: No kernel regularization

```

Listing 63. C6 Implementation: ChatDev (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 implement the complete two-term alignment loss. DeepCode has it as a separate function. MetaGPT and ChatDev lack this regularization.

### 3.4.6. Scoring Analysis

Table 23. Deblur-NeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.00
GPT-5 thinking	0.9	1.0	1.0	1.0	1.0	1.0	0.8	0.7	0.94
DeepCode	0.8	0.8	0.9	0.8	1.0	0.9	0.5	0.0	0.71
MetaGPT	0.4	0.0	0.3	0.2	0.0	0.0	0.0	0.0	0.11
ChatDev	0.0	0.0	0.0	0.0	0.3	0.0	0.0	0.0	0.04

### 3.4.7. Why Baselines Fail Despite Component Scores

**GPT-5 Extended Thinking (Score: 0.94 components, 50% trainable)**

- **Strengths:** Near-complete DSK implementation with proper Nerfstudio integration, correct gamma correction and alignment loss
- **Fatal Issues:**



- Missing proper NeRF field inheritance structure
- Incomplete pipeline configuration for distributed training
- **Result:** Trainable for simple scenes but fails on complex multi-GPU setups

#### DeepCode (Score: 0.71 components, 0% trainable)

- **Strengths:** Comprehensive evaluation framework with ablation studies, proper DSK concepts
- **Fatal Issues:**
  - No Nerfstudio integration whatsoever
  - Custom training loop incompatible with ns-train command
  - Missing scene normalization and camera conventions
- **Result:** Standalone implementation that cannot interface with Nerfstudio datasets

#### MetaGPT (Score: 0.11 components, 0% trainable)

- **Strengths:** Multi-file organization with basic NeRF structure
- **Fatal Issues:**
  - Missing critical DSK components (view embedding, alignment loss)
  - No framework integration
  - Incorrect positional encoding implementation
- **Result:** Incomplete implementation missing core paper contributions

#### ChatDev (Score: 0.04 components, 0% trainable)

- **Strengths:** Basic GUI interface for user interaction
- **Fatal Issues:**
  - Completely missing DSK module
  - No blur modeling whatsoever
  - Tkinter GUI instead of proper training pipeline
- **Result:** Vanilla NeRF implementation unrelated to Deblur-NeRF paper

#### 3.4.8. Hyperparameter Fidelity

Table 24. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5-thinking	DeepCode	MetaGPT	ChatDev
N (kernel points)	5	✓	✓	✓	×	×
$\lambda_{align}$	0.1	✓	✓	✓	×	×
$\lambda_o$	10.0	✓	✓	✓	×	×
View embed dim	32	✓	✓	✓	×	×
Learning rate	5e-4	✓	✓	✓	✓	✓
$\epsilon$ (DSK)	0.1	✓	✓	✓	×	×
Gamma	2.2	✓	✓	✓	×	~
Iterations	200k	✓	✓	✓	✓	×
<b>W Score</b>	–	1.00	1.00	1.00	0.25	0.19

#### 3.4.9. Conclusion

All baseline systems attempt to implement Deblur-NeRF with varying degrees of sophistication. GPT-5 Extended Thinking achieves remarkable component coverage (94%)

with proper DSK implementation and Nerfstudio awareness, occasionally producing trainable code for simple cases. DeepCode provides extensive evaluation infrastructure but lacks framework integration. MetaGPT and ChatDev fail to capture the paper’s core contributions, with ChatDev entirely missing the DSK module. The fundamental gap between component scores and trainability demonstrates that domain-specific framework integration is essential - only NERFIFY produces immediately trainable code with 100% success rate versus 0% for ChatDev, MetaGPT, and DeepCode, with GPT-5 achieving partial success on simpler scenes.

## 4. Novelty Coverage on Set 4 Papers

### 4.1. BioNeRF: Biologically Plausible Neural Radiance Fields

BioNeRF introduces a biologically-inspired architecture that fundamentally reimagines neural radiance fields through cognitive filtering and stateful memory mechanisms. The paper presents parallel positional feature extraction pathways, four distinct cognitive filters, pre-modulation mechanisms, and a recursive memory update system that maintains context across forward passes. This architecture mimics biological neural processing patterns to achieve superior view synthesis quality through memory-conditioned contextual inference.

#### 4.1.1. Implementation Overview

Table 25. BioNeRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	485	182	271	1263	618
File Organization	Plugin	Modular	Single	Multi-file	Multi-file
Parallel MLPs	✓	✓	✓	Simplified	Attempted
Cognitive Filters	✓	✓	✓	Partial	×
Memory Mechanism	✓	✓	Partial	×	×
Nerfstudio Ready	✓	×	×	×	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

#### 4.1.2. Novel Components

Table 26. Novel Components in BioNeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Parallel MLPs ( $M_\Delta$ and $M_c$ ) for positional features	0.10
C2	Cognitive filter $f_\Delta = \sigma(W_\Delta^f h_\Delta + b_\Delta^f)$	0.10
C3	Cognitive filter $f_c = \sigma(W_c^f h_c + b_c^f)$	0.10
C4	Memory filter $f_\psi = \sigma(W_\psi^f [h_\Delta, h_c] + b_\psi^f)$	0.10
C5	Modulation filter $f_\mu = \sigma(W_\mu^f [h_\Delta, h_c] + b_\mu^f)$	0.10
C6	Pre-modulation $\gamma = \tanh(W_\gamma [h_\Delta, h_c] + b_\gamma)$	0.10
C7	Memory modulation $\mu = f_\mu \otimes \gamma$	0.10
C8	Memory update $\Psi = \tanh(W_\Psi (\mu + f_\psi \otimes \Psi) + b_\Psi)$	0.15
C9	Contextual embeddings $h'_\Delta = [\Psi \otimes f_\Delta, x]$ , $h'_c = [\Psi \otimes f_c, d]$	0.10
C10	Density/color heads $M'_\Delta, M'_c$	0.05

Table 27. BioNeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	1.00	1.00
GPT-5	0.80	0.10	0.10	0.80	0.82
DeepSeek R1	0.70	0.20	0.10	0.70	0.75
Paper2Code	0.30	0.40	0.30	0.40	0.35
AutoP2C	0.10	0.30	0.60	0.10	0.15

### 4.1.3. Quantitative Metrics

### 4.1.4. Component-by-Component Analysis

#### Component C1: Parallel MLPs - All Baselines

The paper specifies two independent MLPs for positional feature extraction:  $M_\Delta$  processes encoded 3D positions to produce  $h_\Delta$ , while  $M_c$  processes the same inputs independently to produce  $h_c$ . Both networks have 3 layers with 256 hidden units and ReLU activations.

```

1 # \nerfify\ : Perfect parallel pathway implementation
2 class BioNeRFField(Field):
3     def __init__(self, hidden_dim: int = 256):
4         self.base_delta = MLP(
5             in_dim=pe_out, num_layers=3,
6             layer_width=self.hidden_dim,
7             out_dim=self.hidden_dim,
8             activation=nn.ReLU())
9         self.base_c = MLP(
10            in_dim=pe_out, num_layers=3,
11            layer_width=self.hidden_dim,
12            out_dim=self.hidden_dim,
13            activation=nn.ReLU())
14 # CORRECT: Independent MLPs with exact paper
    specifications

```

Listing 64. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Clean implementation with modular design
2 class BioNeRFField(nn.Module):
3     def __init__(self):
4         self.M_delta = self._make_mlp(63, 256, 256)
5         self.M_color = self._make_mlp(63, 256, 256)
6         def _make_mlp(self, in_dim, hidden, out_dim,
7             n_layers=3):
8             layers = []
9             for i in range(n_layers):
10                layers.append(nn.Linear(
11                    in_dim if i == 0 else hidden, hidden)
12            )
13            layers.append(nn.ReLU(True))
14            return nn.Sequential(*layers)
15 # CORRECT: Proper parallel architecture

```

Listing 65. C1 Implementation: GPT-5 (Score: 0.9)

```

1 # DeepSeek R1: Comprehensive with proper encoding
2 class BioNeRF(nn.Module):
3     def __init__(self, hidden_dim=256, num_freqs=10):
4         self.encoded_3d_dim = 3 * (2 * num_freqs + 1)
5         self.M_delta = self._build_mlp(
6             self.encoded_3d_dim, hidden_dim,
7             hidden_dim, 3)
8         self.M_color = self._build_mlp(
9             self.encoded_3d_dim, hidden_dim,
10            hidden_dim, 3)
11 # ISSUE: Correct structure but naming inconsistency

```

Listing 66. C1 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: Oversimplified to incorrect inputs
2 class BioNeRFModel(nn.Module):
3     def __init__(self):
4         self.mlp_initial_delta = self._build_mlp(
5             self.input_dim_xyz, hidden_size,
6             hidden_size, 3)
7         self.mlp_initial_c = self._build_mlp(
8             self.input_dim_dir, hidden_size,
9             hidden_size, 3)
10 # ERROR: Uses different inputs for each MLP (xyz vs
11 dir)
12 # Should both process positional-encoded coordinates

```

Listing 67. C1 Implementation: Paper2Code (Score: 0.3)

```

1 # AutoP2C: Structure exists but fundamental error
2 def positional_feature_extraction(self, inputs):
3     encoded_inputs = positional_encoding(inputs, 10)
4     # CRITICAL ERROR: Uses raw inputs, not encoded
5     h_delta = self.mlp_delta(inputs) # Wrong!
6     h_c = self.mlp_color(inputs) # Wrong!
7     return h_delta, h_c

```

Listing 68. C1 Implementation: AutoP2C (Score: 0.1)

**Analysis:** NERFIFY correctly implements independent parallel MLPs with exact paper specifications. GPT-5 achieves near-perfect implementation. DeepSeek R1 has correct structure but minor naming issues. Paper2Code incorrectly uses different inputs for each MLP. AutoP2C defines encoding but critically fails to use it, passing raw inputs instead.

#### Component C2-C6: Cognitive Filtering - All Baselines

The paper defines four cognitive filters and pre-modulation:  $f_\Delta = \sigma(W_\Delta^f h_\Delta + b_\Delta^f)$ ,  $f_c = \sigma(W_c^f h_c + b_c^f)$ ,  $f_\psi = \sigma(W_\psi^f [h_\Delta, h_c] + b_\psi^f)$ ,  $f_\mu = \sigma(W_\mu^f [h_\Delta, h_c] + b_\mu^f)$ , and  $\gamma = \tanh(W_\gamma [h_\Delta, h_c] + b_\gamma)$ .

```

1 # \nerfify\ : All filters properly defined in
2 __init__
3 def __init__(self):
4     self.W_delta_f = nn.Linear(self.hidden_dim, self.
5         hidden_dim)
6     self.W_c_f = nn.Linear(self.hidden_dim, self.
7         hidden_dim)
8     self.W_psi_f = nn.Linear(2 * self.hidden_dim,
9         self.hidden_dim)
10    self.W_mu_f = nn.Linear(2 * self.hidden_dim, self.
11        .hidden_dim)
12    self.W_gamma = nn.Linear(2 * self.hidden_dim,
13        self.hidden_dim)
14    def forward(self):
15        f_delta = self._sigmoid(self.W_delta_f(h_delta))
16        f_c = self._sigmoid(self.W_c_f(h_c))
17        h_cat = torch.cat([h_delta, h_c], dim=-1)
18        f_psi = self._sigmoid(self.W_psi_f(h_cat))
19        f_mu = self._sigmoid(self.W_mu_f(h_cat))
20        gamma = self._tanh(self.W_gamma(h_cat))
21 # CORRECT: All filters with proper activations

```

Listing 69. C2-C6 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Clean filter implementation
2 class BioNeRFField(nn.Module):
3     def __init__(self):
4         self.f_delta = nn.Linear(hidden, hidden)
5         self.f_color = nn.Linear(hidden, hidden)
6         self.f_mem = nn.Linear(hidden * 2, hidden)
7         self.f_mod = nn.Linear(hidden * 2, hidden)
8         self.pre_mod = nn.Linear(hidden * 2, hidden)

```

```

9     def forward(self, x, d):
10         f_delta = torch.sigmoid(self.f_delta(h_delta))
11         f_color = torch.sigmoid(self.f_color(h_color))
12         gamma = torch.tanh(self.pre_mod(h_cat))
13 # CORRECT: Good implementation with slight naming
    differences

```

Listing 70. C2-C6 Implementation: GPT-5 (Score: 0.8)

```

1 # DeepSeek R1: Comprehensive filter implementation
2 def __init__(self):
3     self.W_f_delta = nn.Linear(hidden_dim, hidden_dim)
4     self.W_f_color = nn.Linear(hidden_dim, hidden_dim)
5     self.W_f_psi = nn.Linear(2 * hidden_dim,
6                               hidden_dim)
7     self.W_f_mu = nn.Linear(2 * hidden_dim,
8                              hidden_dim)
9     def forward(self):
10         f_delta = torch.sigmoid(self.W_f_delta(h_delta))
11         f_psi = torch.sigmoid(self.W_f_psi(h_concat))
12 # ISSUE: Missing W_gamma layer definition

```

Listing 71. C2-C6 Implementation: DeepSeek R1 (Score: 0.7)

```

1 # Paper2Code: Partial implementation
2 def forward(self):
3     f_delta = torch.sigmoid(self.linear_fd(h_delta))
4     f_c = torch.sigmoid(self.linear_fc(h_c))
5     # Concatenates but missing some filters
6     h_cat = torch.cat((h_delta, h_c), dim=-1)
7     gamma = torch.tanh(self.linear_gamma(h_cat))
8 # MISSING: f_psi and f_mu filters not properly
    defined

```

Listing 72. C2-C6 Implementation: Paper2Code (Score: 0.4)

```

1 # AutoP2C: Fatal implementation error
2 def cognitive_filtering(self, features):
3     h_delta, h_c = features
4     # CRITICAL ERROR: Creates new layers every
5     # forward pass!
6     f_delta = torch.sigmoid(
7         nn.Linear(256, 256)(h_delta)) # Random
8     weights!
9     f_c = torch.sigmoid(
10        nn.Linear(256, 256)(h_c)) # Random
11    weights!
12    # No gradient tracking possible

```

Listing 73. C2-C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY perfectly implements all cognitive filters. GPT-5 has good implementation with minor naming variations. DeepSeek R1 misses gamma layer definition. Paper2Code partially implements filters. AutoP2C has a critical flaw - creating nn.Linear layers inside forward pass results in random weights every iteration with no gradient tracking.

#### Component C7-C8: Memory Mechanism - All Baselines

The memory mechanism involves modulation  $\mu = f_\mu \otimes \gamma$  and recursive update  $\Psi = \tanh(W_\Psi(\mu + f_\psi \otimes \Psi_{prev}) + b_\Psi)$  where the memory state persists across forward passes.

```

1 # \nerfify\ : Perfect stateful memory with detach
2 def forward(self):

```

```

3     prev_mem = self._ensure_memory(N, device, dtype)
4     mu = f_mu * gamma # Element-wise modulation
5     # Proper recursive update with previous memory
6     Psi = self._tanh(self.W_Psi(mu + f_psi * prev_mem))
7     # Detach to prevent gradient explosion
8     self._memory = Psi.detach()
9 # CORRECT: Stateful memory with proper recursion

```

Listing 74. C7-C8 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Separate memory module approach
2 class MemoryModule(nn.Module):
3     def forward(self, mu, fpsi):
4         if self.memory is None:
5             self.memory = torch.zeros_like(mu)
6             updated = torch.tanh(
7                 self.W_psi(mu + fpsi * self.memory) +
8                 self.b_psi)
9             self.memory = updated.detach()
10            return updated
11 # CORRECT: Proper recursion, slight module complexity

```

Listing 75. C7-C8 Implementation: GPT-5 (Score: 0.8)

```

1 # DeepSeek R1: Good attempt with persistence issue
2 def forward(self, xyz, view_dir):
3     mu = f_mu * gamma # Correct modulation
4     if self.memory.shape[0] != batch_size:
5         self.memory = self.memory.expand(batch_size,
6                                           -1)
7     memory_input = mu + (f_psi * self.memory)
8     psi_new = torch.tanh(self.W_psi(memory_input))
9     self.memory = psi_new
10 # ISSUE: Memory expansion may not persist properly

```

Listing 76. C7-C8 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: Missing recursive component
2 def forward(self):
3     mu = f_mu * gamma
4     psi_init = torch.zeros_like(mu)
5     # ERROR: Always uses zeros, no memory persistence
6     psi = torch.tanh(self.linear_memory_update(
7         mu + (f_psi * psi_init))) # psi_init always
8     0!
9 # MISSING: No actual memory recursion

```

Listing 77. C7-C8 Implementation: Paper2Code (Score: 0.2)

```

1 # AutoP2C: Multiple critical issues
2 def memory_updating(self, filtered_features):
3     mu = f_mu * gamma
4     # CRITICAL ERROR: Creates new layer here
5     self.memory = torch.tanh(
6         nn.Linear(256, 256)(mu + f_psi * self.memory))
7     # Random weights, no learning possible

```

Listing 78. C7-C8 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY implements perfect stateful memory with detach to prevent gradient issues. GPT-5 uses a correct separate module approach. DeepSeek R1 has memory persistence issues. Paper2Code fails to maintain memory state. AutoP2C creates layers dynamically, making learning impossible.

#### Component C9-C10: Contextual Inference and Output Heads - All Baselines

The contextual inference creates embeddings  $h'_\Delta = [\Psi \otimes f_\Delta, x]$  and  $h'_c = [\Psi \otimes f_c, d]$ , which are processed by density head  $M'_\Delta$  and color head  $M'_c$  respectively.

```
1 # \nerfify\ : Perfect contextual inference
2 def forward(self):
3     ctx_delta = Psi * f_delta # Element-wise
4     ctx_c = Psi * f_c
5     sigma_in = torch.cat([ctx_delta, pe], dim=-1)
6     sigma_raw = self.density_head(sigma_in)
7     density = trunc_exp(sigma_raw).view(*positions.
8     shape[:-1], 1)
9     rgb_in = torch.cat([ctx_c, de], dim=-1)
10    rgb = self.color_head(rgb_in).view(*positions.
11    shape[:-1], 3)
12 # CORRECT: Exact paper implementation
```

Listing 79. C9-C10 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Good implementation with minor differences
2 def forward(self):
3     hpd = torch.cat([psi * f_delta, x], dim=-1)
4     hpc = torch.cat([psi * f_color, d], dim=-1)
5     density = self.Mp_delta(hpd)
6     color = torch.sigmoid(self.Mp_color(hpc))
7 # CORRECT: Proper concatenation and heads
```

Listing 80. C9-C10 Implementation: GPT-5 (Score: 0.8)

```
1 # DeepSeek R1: Correct approach, minor issues
2 def forward(self):
3     h_delta_prime = torch.cat([psi_new * f_delta, xyz
4     ], dim=-1)
5     h_color_prime = torch.cat([psi_new * f_color,
6     view_dir], dim=-1)
7     density = self.density_head(h_delta_prime)
8     color = torch.sigmoid(self.color_head(
9     h_color_prime))
10 # ISSUE: Missing proper view reshaping
```

Listing 81. C9-C10 Implementation: DeepSeek R1 (Score: 0.7)

```
1 # Paper2Code: Simplified inference
2 def forward(self):
3     h_delta_prime = torch.cat((psi * f_delta, 1), dim
4     =-1)
5     density = self.density_head(h_delta_prime)
6     color = self.color_head(h_c_prime)
7 # ERROR: h_c_prime not properly defined
```

Listing 82. C9-C10 Implementation: Paper2Code (Score: 0.3)

```
1 # AutoP2C: Structural attempt, poor execution
2 def forward(self):
3     density = self.heads['density'](h_delta_prime)
4     color = self.heads['color'](h_c_prime)
5     return density, color
6 # ERROR: h_delta_prime, h_c_prime undefined
7 # MISSING: No contextual concatenation
```

Listing 83. C9-C10 Implementation: AutoP2C (Score: 0.1)

**Analysis:** NERFIFY perfectly implements contextual inference with proper concatenation and reshaping. GPT-5 has good implementation with minor variations. DeepSeek R1 misses view reshaping. Paper2Code has undefined variables. AutoP2C attempts structure but lacks proper implementation.

Table 28. Revised BioNeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.00
GPT-5	0.9	0.8	0.8	0.8	0.8	0.8	0.8	0.8	0.8	0.8	0.82
DeepSeek R1	0.8	0.7	0.7	0.7	0.7	0.7	0.7	0.6	0.7	0.7	0.70
Paper2Code	0.3	0.4	0.4	0.4	0.4	0.4	0.3	0.2	0.3	0.3	0.34
AutoP2C	0.1	0.0	0.0	0.0	0.0	0.0	0.1	0.0	0.1	0.1	0.04

#### 4.1.5. Scoring Analysis

##### 4.1.6. Why Baselines Fail Despite Component Scores

###### GPT-5 (Score: 0.82 components, 0% trainable)

- **Strengths:** Near-perfect component implementation, clean modular code structure, proper memory handling
- **Fatal Issues:**
  - No Nerfstudio integration whatsoever
  - Missing training pipeline and data loading infrastructure
  - Standalone module without framework compatibility
- **Result:** Excellent algorithm but requires 3-5 hours of integration work to be trainable

###### DeepSeek R1 (Score: 0.75 components, 0% trainable)

- **Strengths:** Comprehensive implementation with all components present, proper encoding
- **Fatal Issues:**
  - Memory buffer expansion doesn't persist between batches
  - No Nerfstudio plugin structure
  - Missing critical training loop integration
- **Result:** Memory persistence issues cause degradation to vanilla NeRF behavior

###### Paper2Code (Score: 0.35 components, 0% trainable)

- **Strengths:** Extensive 1200+ line codebase with complete dataset loading infrastructure
- **Fatal Issues:**
  - Implements generic NeRF with BioNeRF naming conventions
  - Memory mechanism lacks recursion - always uses zeros
  - Missing true parallel pathways for  $M_{\Delta}$  and  $M_c$
  - No Nerfstudio integration
- **Result:** Code runs but implements wrong algorithm, achieving vanilla NeRF performance

###### AutoP2C (Score: 0.15 components, 0% trainable)

- **Strengths:** Attempts all structural components with correct high-level organization
- **Fatal Issues:**
  - Creates nn.Linear layers in forward pass resulting in random weights every iteration
  - Positional encoding defined but not applied to inputs
  - No gradient tracking possible with dynamic layer creation
  - Fundamental PyTorch misunderstanding
- **Result:** Catastrophic implementation errors prevent any learning from occurring



#### 4.1.7. Hyperparameter Fidelity

Table 29. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
Learning rate	5e-4	✓	✓	✓	✓	×
Hidden dimensions	256	✓	✓	✓	✓	✓
Memory tensor size	8192	✓	~	~	×	✓
Training iterations	400k	✓	×	×	✓	×
Adam optimizer	Yes	✓	✓	✓	✓	✓
Positional frequencies	10	✓	~	✓	✓	✓
Batch size	8192	✓	~	×	✓	×
<b>W Score</b>	–	1.00	0.57	0.57	0.71	0.43

#### 4.1.8. Conclusion

All baselines attempt to implement BioNeRF with varying levels of sophistication and understanding. GPT-5 and DeepSeek R1 achieve high component scores with nearly correct algorithmic implementations but completely lack the Nerfstudio integration necessary for training. Paper2Code produces an extensive codebase but fundamentally misunderstands the algorithm, implementing a simplified NeRF variant that misses the critical recursive memory mechanism. AutoP2C demonstrates the most severe implementation errors, creating neural network layers dynamically within the forward pass, which results in random weights every iteration and makes gradient-based learning impossible. Only NERFIFY produces immediately trainable code by combining perfect component implementation with proper framework integration as a complete Nerfstudio plugin, achieving a 100% trainable rate while all baselines remain at 0% trainable despite their varying component scores. ““

### 4.2. Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields

#### 4.2.1. Paper Overview

Mip-NeRF [2] addresses the critical aliasing artifacts in Neural Radiance Fields by replacing point sampling with cone tracing and standard positional encoding with integrated positional encoding (IPE). The paper’s key innovation is approximating conical frustums as 3D Gaussians and computing the expected positional encoding analytically, enabling anti-aliased rendering across multiple scales with a single multiscale MLP that reduces parameters by 50% compared to NeRF.

#### 4.2.2. Implementation Overview

#### 4.2.3. Novel Components

#### 4.2.4. Quantitative Metrics

#### 4.2.5. Component-by-Component Analysis

##### Component C1: Conical Frustum Parameters - All Baselines

The paper derives stable reparameterized formulas (Appendix A) for computing Gaussian parameters of a conical

Table 30. Mip-NeRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	563	412	287	1842	895
File Organization	Plugin	Single	Single	Multi-file	Multi-file
Conical Frustum	✓	×	✓	✓	Simplified
IPE Implementation	✓	×	✓	✓	Attempted
Single MLP	✓	✓	✓	✓	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio

plugin

Table 31. Novel Components in Mip-NeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Conical Frustum Parameters (Eq. 7-8, stable formulation)	0.20
C2	Gaussian Mean Transform to World Coordinates	0.10
C3	Gaussian Covariance Transform (diagonal approximation)	0.15
C4	Integrated Positional Encoding (IPE)	0.25
C5	Single Multiscale MLP Architecture	0.15
C6	Area-Weighted Loss for Multiscale Training	0.15

Table 32. Mip-NeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	1.00	1.00
GPT-5	0.50	0.30	0.20	0.60	0.58
DeepSeek R1	0.67	0.17	0.16	0.67	0.58
Paper2Code	0.83	0.17	0.00	0.83	0.85
AutoP2C	0.17	0.17	0.66	0.25	0.20

frustum. Given segment endpoints  $t_0, t_1$  and cone radius  $\hat{r}$ , the formulas use midpoint  $t_\mu = (t_0 + t_1)/2$  and half-width  $t_\delta = (t_1 - t_0)/2$  to compute mean  $\mu_t$ , axial variance  $\sigma_t^2$ , and radial variance  $\sigma_r^2$ .

```

1 # \nerfify\ : Exact stable formulation from paper
2 t_mu = (t0 + t1) / 2.0
3 t_delta = (t1 - t0) / 2.0
4 mu_t = t_mu + (2.0 * t_mu * t_delta ** 2) / (3.0 *
   t_mu ** 2 + t_delta ** 2)
5 sigma_t_sq = (t_delta ** 2) / 3.0 - \
6 (4.0 * t_delta ** 4 * (12.0 * t_mu ** 2 - t_delta
   ** 2)) / \
7 (15.0 * (3.0 * t_mu ** 2 + t_delta ** 2) ** 2)
8 sigma_r_sq = (dot_r ** 2) * (t_mu ** 2 / 4.0 + 5.0 *
   t_delta ** 2 / 12.0 - \
9 4.0 * t_delta ** 4 / (15.0 * (3.0 * t_mu ** 2 +
   t_delta ** 2)))
10 # CORRECT: Stable reparameterized formulas from
   Appendix A

```

Listing 84. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Completely wrong formulas
2 def conical_frustum_to_gaussian(t0, t1, radius):
3     t_mean = (3*(t1**4 - t0**4))/(4*(t1**3 - t0**3))
4     t_var = (3*(t1**5 - t0**5))/(5*(t1**3 - t0**3)) -
       t_mean**2
5     r_var = radius**2 * (3*(t1**5 - t0**5))/(20*(t1
       **3 - t0**3))
6     return t_mean, t_var, r_var
7 # CRITICAL ERROR: Different derivation, numerically
   unstable

```

Listing 85. C1 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: Correct stable formulation
2 mu_t = t_mu + (2 * t_mu * t_delta**2) / (3 * t_mu**2
  + t_delta**2)
3 sigma_t2 = (t_delta**2 / 3) - (4 * t_delta**4 * (12 *
  t_mu**2 - t_delta**2)) / \
4   (15 * (3 * t_mu**2 + t_delta**2)**2)
5 sigma_r2 = r**2 * (t_mu**2 / 4 + 5 * t_delta**2 / 12
  - \
6   4 * t_delta**4 / (15 * (3 * t_mu**2 + t_delta**2)
  ))
7 # CORRECT: Matches paper formulation

```

Listing 86. C1 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Correct with JAX
2 denominator = 3 * (z_center ** 2) + z_delta ** 2 + 1e
  -8
3 mu_z = z_center + (2.0 * z_center * (z_delta ** 2)) /
  denominator
4 sigma_z_sq = (z_delta ** 2) / 3.0 - \
5   (4.0 * (z_delta ** 4) * (12 * z_center ** 2 -
  z_delta ** 2)) / \
6   (15.0 * (denominator ** 2) + 1e-8)
7 # CORRECT: Stable formulation with numerical safety

```

Listing 87. C1 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: Missing frustum computation
2 # No conical frustum parameter computation found
3 # Uses stratified sampling without cone tracing
4 # MISSING: Core component entirely absent

```

Listing 88. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY, DeepSeek R1, and Paper2Code correctly implement the stable reparameterized formulas. GPT-5 uses completely different formulas that are both mathematically incorrect and numerically unstable. AutoP2C entirely omits conical frustum computation.

### Component C2-C3: Gaussian Mean and Covariance Transform - All Baselines

The paper transforms frustum parameters to world coordinates:  $\mu = \mathbf{o} + \mu_t \mathbf{d}$  for the mean and  $\Sigma = \sigma_t^2 (\mathbf{d} \mathbf{d}^\top + \sigma_r^2 (\mathbf{I} - \mathbf{d} \mathbf{d}^\top / \|\mathbf{d}\|_2^2))$  for covariance.

```

1 # \nerfify\ : Efficient diagonal approximation
2 means = origins + mu_t[..., None] * directions #
  Mean transform
3 d_outer_diag = directions ** 2
4 d_norm_sq = (directions ** 2).sum(dim=-1, keepdim=
  True)
5 covs_diag = sigma_t_sq[..., None] * d_outer_diag + \
6   sigma_r_sq[..., None] * (1.0 - d_outer_diag /
  d_norm_sq)
7 # CORRECT: Diagonal covariance for efficiency

```

Listing 89. C2-C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Full covariance matrix (inefficient)
2 mu = rays_o[..., None, :] + mu_z[..., None] * rays_d
  [..., None, :]
3 covs = (rays_d[..., None] * rays_d[..., None, :] *
  t_var[..., None, None]
4   + jnp.eye(3) * r_var[..., None, None])
5 # ISSUE: Computes full 3x3 matrix, inefficient but
  mathematically correct

```

Listing 90. C2-C3 Implementation: GPT-5 (Score: 0.5)

```

1 # Paper2Code: Correct diagonal implementation
2 mu = rays_o_exp + jnp.expand_dims(mu_z, axis=-1) *
  rays_d_exp
3 r_squared = rays_d ** 2
4 cov_diag = jnp.expand_dims(sigma_z_sq, axis=-1) *
  r_squared_exp + \
5   jnp.expand_dims(sigma_r_sq, axis=-1) * (1.0 -
  r_squared_exp)
6 # CORRECT: Efficient diagonal form

```

Listing 91. C2-C3 Implementation: Paper2Code (Score: 1.0)

**Analysis:** Most baselines correctly implement the transforms. GPT-5 computes the full covariance matrix instead of just the diagonal, which is mathematically correct but computationally inefficient.

### Component C4: Integrated Positional Encoding (IPE) - All Baselines

The core innovation computes expected positional encoding over a Gaussian:  $\mathbb{E}[\sin(x)] = \sin(\mu) \exp(-\sigma^2/2)$  and  $\mathbb{E}[\cos(x)] = \cos(\mu) \exp(-\sigma^2/2)$  for frequency bands  $2^0, 2^1, \dots, 2^{L-1}$ .

```

1 # \nerfify\ : Perfect IPE implementation
2 scales = 2.0 ** torch.arange(self.
  position_encoding_max_degree)
3 means_scaled = means[..., :, None] * scales
4 variance_scales = (scales ** 2)
5 covs_scaled = covs_diag[..., :, None] *
  variance_scales
6 damping = torch.exp(-0.5 * covs_scaled)
7 sin_features = torch.sin(means_scaled) * damping
8 cos_features = torch.cos(means_scaled) * damping
9 features = torch.cat([sin_features, cos_features],
  dim=-1)
10 # CORRECT: Exact formulation from paper

```

Listing 92. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Multiple critical errors
2 freqs = 2.0 ** jnp.arange(max_freq_log2)
3 scales = jnp.pi * freqs[None, None, :] # ERROR:
  Extra pi!
4 diag = jnp.stack([cov[..., i, i] for i in range(3)],
  axis=-1)
5 var = diag[..., None, :] * scales[..., None] ** 2
6 # CRITICAL ERROR: Wrong frequency scaling, dimension
  mismatches

```

Listing 93. C4 Implementation: GPT-5 (Score: 0.2)

```

1 # Paper2Code: Correct JAX implementation
2 freq_bands = 2.0 ** jnp.arange(L, dtype=jnp.float32)
3 args = mu_exp * freq_bands
4 damping = jnp.exp(-0.5 * (cov_exp * (freq_bands ** 2)
  ))
5 sin_features = jnp.sin(args) * damping
6 cos_features = jnp.cos(args) * damping
7 # CORRECT: Proper IPE computation

```

Listing 94. C4 Implementation: Paper2Code (Score: 1.0)

**Analysis:** NERFIFY, DeepSeek R1, and Paper2Code correctly implement IPE. GPT-5 has multiple errors including extra  $\pi$  multiplication and dimension mismatches. AutoP2C attempts IPE but with incorrect structure.

## Component C5: Single MLP Architecture - All Baselines

Mip-NeRF uses a single multiscale MLP instead of separate coarse/fine networks, reducing parameters by 50

```
1 # \nerfify\ : Single MLP with hierarchical sampling
  in Nerfstudio
2 class MipNeRFModel(Model):
3     def populate_modules(self):
4         self.field = MipNeRFField() # Single field
5         self.sampler_uniform = UniformSampler(
6             num_samples=128)
7         self.sampler_pdf = PDFSampler(num_samples
8             =128, include_original=False)
9 # CORRECT: Single MLP queried hierarchically
```

Listing 95. C5 Implementation: NERFIFY (Score: 1.0)

```
1 # AutoP2C: Uses two separate MLPs
2 coarse_out = self.model.apply({"params": self.state.
3     params}, ray_batch_coarse)
4 # ... hierarchical sampling ...
5 fine_out = state.apply_fn({"params": params},
6     ray_batch_fine)
7 # ERROR: Separate coarse/fine networks instead of
  single MLP
```

Listing 96. C5 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Most baselines correctly implement a single MLP. AutoP2C incorrectly uses hierarchical two-network architecture from original NeRF.

## Component C6: Area-Weighted Loss - All Baselines

For multiscale training, pixel losses are weighted by their footprint area.

```
1 # \nerfify\ : Full area-weighted loss support
2 if self.config.area_loss_mult > 0:
3     pixel_area = ray_bundle.pixel_area
4     loss = loss * pixel_area.squeeze() * self.config.
5         area_loss_mult
6 # CORRECT: Proper multiscale weighting
```

Listing 97. C6 Implementation: NERFIFY (Score: 1.0)

```
1 # Paper2Code: Partial implementation
2 def area_weighted_loss(pred_rgb, target_rgb,
3     pixel_area, lambda=0.1):
4     diff = (pred_rgb - target_rgb)**2
5     return lambda * jnp.mean(diff * pixel_area[...
6         , None])
7 # ISSUE: Defined but not integrated into training
  loop
```

Listing 98. C6 Implementation: Paper2Code (Score: 0.5)

**Analysis:** Only NERFIFY fully implements area-weighted loss. Paper2Code defines it but doesn't integrate it. Others omit it entirely.

### 4.2.6. Scoring Analysis

### 4.2.7. Why Baselines Fail Despite Component Scores

**GPT-5 (Score: 0.58 components, 0% trainable)**

- **Strengths:** Implements single MLP architecture and basic transforms
- **Fatal Issues:**
  - Wrong conical frustum formulas cause NaN gradients

Table 33. Mip-NeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.00
GPT-5	0.0	1.0	0.5	0.2	1.0	0.0	0.58
DeepSeek R1	1.0	1.0	1.0	1.0	0.5	0.0	0.58
Paper2Code	1.0	1.0	1.0	1.0	1.0	0.5	0.85
AutoP2C	0.0	0.5	0.5	0.3	0.0	0.0	0.20

- IPE implementation has critical errors (extra  $\pi$ , dimension mismatches)
- JAX code incompatible with Nerfstudio benchmark
- **Result:** Non-trainable due to fundamental mathematical errors
- **DeepSeek R1 (Score: 0.58 components, 0% trainable)**
- **Strengths:** Correct mathematical formulations for frustum and IPE
- **Fatal Issues:**
  - No Nerfstudio integration (standalone PyTorch)
  - Missing critical training infrastructure
  - Simplified hierarchical sampling without proposal network
- **Result:** Cannot be executed in benchmark environment
- **Paper2Code (Score: 0.85 components, 0% trainable)**
- **Strengths:** Most accurate mathematical implementation in JAX/Flax
- **Fatal Issues:**
  - Wrong framework (JAX instead of PyTorch/Nerfstudio)
  - Incompatible API prevents benchmark evaluation
  - Requires complete infrastructure rewrite
- **Result:** Excellent standalone code but non-trainable in benchmark
- **AutoP2C (Score: 0.20 components, 0% trainable)**
- **Strengths:** Attempts basic NeRF structure
- **Fatal Issues:**
  - Missing conical frustum computation entirely
  - Uses two MLPs instead of single architecture
  - Incomplete Nerfstudio integration
  - Many placeholder functions
- **Result:** Fundamentally incomplete implementation

### 4.2.8. Hyperparameter Fidelity

Table 34. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
PE frequencies ( $L$ )	16	✓	✓	✓	✓	×
Direction encoding	4	✓	×	✓	×	×
Hidden dimension	256	✓	✓	✓	✓	✓
Learning rate	$5e^{-4}$	✓	✓	✓	✓	×
Batch size	4096	✓	×	×	✓	✓
<b>W Score</b>	–	1.00	0.60	0.67	0.83	0.25

### 4.2.9. Conclusion

All baselines attempt to implement Mip-NeRF with varying sophistication. Paper2Code achieves the highest component accuracy (0.85) with correct mathematical implementations but uses JAX/Flax incompatible with Nerfstudio. GPT-5 and DeepSeek R1 both score 0.58 but fail differently - GPT-5 has fundamental mathematical errors while R1 lacks integration. AutoP2C scores lowest (0.20) missing core components entirely. Only NERFIFY produces immediately trainable code as a complete Nerfstudio plugin, achieving 100% trainable rate while all baselines achieve 0%.

## 4.3. PyNeRF: Pyramidal Neural Radiance Fields

### 4.3.1. Paper Overview

PyNeRF introduces a pyramidal multiscale architecture for neural radiance fields that addresses aliasing artifacts through dynamic resolution selection. The method maintains  $L = 8$  pyramid levels with geometrically increasing resolutions  $N_l = N_0 \times s^l$  (where  $N_0 = 16$  and  $s = 2$ ), selecting the appropriate level for each sample based on its projected pixel area  $P(x)$  using the mapping function  $M(P) = \log_s(P/N_0)$ . The key innovation lies in per-sample level selection with linear interpolation between adjacent pyramid levels to achieve smooth transitions.

### 4.3.2. Implementation Overview

Table 35. PyNeRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	420	185	380	1250	450
File Organization	Plugin	Modular	Single	Multi-file	Multi-file
Pyramid Architecture	✓	✓	✓	✓	×
Level Selection	✓	✓	Partial	✓	×
Interpolation	✓	✓	Partial	✓	×
Projected Area	✓	~	×	×	×
Nerfstudio Integration	✓	×	×	×	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

### 4.3.3. Novel Components

Table 36. Novel Components in PyNeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Pyramidal architecture with L levels	0.15
C2	Per-sample level selection based on projected area	0.20
C3	Linear interpolation between levels	0.15
C4	Shared multiresolution features	0.10
C5	Per-level tiny MLPs (64/128 hidden)	0.10
C6	Projected area computation $P(x) = \text{pixel\_area} \times t^2$	0.10
C7	Adaptive level supervision	0.05
C8	Modern backbone integration (iNGP/TensoRF)	0.05
C9	Multiple interpolation modes (linear/Laplacian)	0.05
C10	Specific hyperparameters (8 levels, s=2, N0=16)	0.05

Table 37. PyNeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	0.95	0.05	0.00	1.00	0.97
GPT-5	0.40	0.30	0.30	0.80	0.52
DeepSeek R1	0.30	0.60	0.10	0.70	0.68
Paper2Code	0.50	0.30	0.20	0.60	0.58
AutoP2C	0.00	0.10	0.90	0.00	0.03

### 4.3.4. Quantitative Metrics

#### 4.3.5. Component-by-Component Analysis

**Component C1: Pyramidal Architecture with L Levels - All Baselines** The paper specifies creating  $L = 8$  pyramid levels where each level  $l$  has resolution  $N_l = N_0 \times s^l$  with separate processing heads for each level.

```

1 # \nerfify\ : ModuleList of per-level heads
2 self.density_heads = nn.ModuleList([
3     MLP(in_dim=self.feature_dim, num_layers=1,
4         layer_width=density_hidden_dim, out_dim=1)
5     for _ in range(self.num_pyramid_levels)
6 ])
7 self.color_heads = nn.ModuleList([
8     MLP(in_dim=color_in_dim, num_layers=2,
9         layer_width=color_hidden_dim, out_dim=3)
10    for _ in range(self.num_pyramid_levels)
11 ])
12 # CORRECT: Exactly L separate heads as specified

```

Listing 99. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: List comprehension for levels
2 models = [PyNeRFLevel(enc, cfg) for _ in range(cfg.
3     n_levels)]
4 pipeline = PyNeRFPipeline(models, cfg).cuda()
5 # CORRECT: Creates pyramid structure but lacks
6   integration

```

Listing 100. C1 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: ModuleList with incomplete
2   implementation
3 self.levels = nn.ModuleList()
4 for i in range(num_levels):
5     level = PyNeRFLevel(base_nerf_config, level_idx=i
6         ,
7         use_shared_features=
8         self.levels.append(level)
9 # ISSUE: Structure correct but PyNeRFLevel incomplete

```

Listing 101. C1 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Detailed pyramid construction
2 self.heads = nn.ModuleList()
3 for level in range(self.pyramid_levels):
4     effective_resolution = self.base_resolution * (
5         self.scale_factor ** level)
6     head = PyNeRFHead(input_dim=num_features,
7         effective_resolution=
8         effective_resolution)
9     self.heads.append(head)
10 # CORRECT: Proper pyramid with effective resolutions

```

Listing 102. C1 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: Single MLP for all samples
2 self.density_mlp = nn.Sequential(
3     nn.Linear(2 * 10 * 3, 64), nn.ReLU(), nn.Linear
4     (64, 1))
5 self.color_mlp = nn.Sequential(
6     nn.Linear(2 * 10 * 6 + 64, 128), nn.ReLU(),
7     nn.Linear(128, 128), nn.ReLU(), nn.Linear(128, 3)
8 )
9 # CRITICAL ERROR: No pyramid structure at all

```

Listing 103. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** All methods except AutoP2C correctly implement the pyramid architecture. AutoP2C completely misses this core concept.

**Component C2: Per-Sample Level Selection - All Baselines** The paper requires computing  $M(P) = \log_s(P/N_0)$  and selecting level  $l = \lceil M(P) \rceil$  with interpolation weight  $w = l - M(P)$ .

```

1 # \nerfify\ : Exact paper formulation
2 t_mid = (ray_samples.frustums.starts + ray_samples.
3     frustums.ends) * 0.5
4 pixel_area = ray_samples.frustums.pixel_area
5 proj_area = pixel_area * torch.clamp(t_mid, min=0.0)
6 s = torch.tensor(self.pyramid_scale_factor, device=
7     proj_area.device)
8 M = torch.log(proj_area / max(self.base_proj_area, 1e
9     -12)) / torch.log(s)
10 l = torch.clamp(torch.ceil(M), min=0.0, max=float(
11     self.num_pyramid_levels - 1))
12 w = 1 - M # Interpolation weight
13 # CORRECT: Matches paper equations exactly

```

Listing 104. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Clean implementation
2 def compute_level(Px, N0, s):
3     M = torch.log2(Px / N0) / torch.log2(torch.tensor
4     (s))
5     l = torch.clamp(torch.ceil(M), 0, None)
6     w = 1 - M
7     return l.long(), w
8 # CORRECT: Proper level selection logic

```

Listing 105. C2 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Incomplete with heuristics
2 def compute_sample_scale(self, rays, sample_positions
3     ):
4     sample_distances = torch.norm(sample_positions -
5     rays['origins']).unsqueeze(1, dim=-1)
6     pixel_areas = sample_distances ** 2 * 0.01 #
7     heuristic scaling
8     return pixel_areas.unsqueeze(-1).expand_as(
9     sample_distances)
10 M = torch.log(scales / self.base_resolution) / np.log
11 (self.scale_factor)
12 levels = torch.clamp(torch.ceil(M), 0, self.
13     num_levels - 1).long()
14 # ERROR: Heuristic area computation, not paper
15 formula

```

Listing 106. C2 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: Correct logic, wrong area formula
2 def compute_integration_area(self, depths):
3     integration_area = (depths / self.focal) ** 2
4     return integration_area
5 M_value = torch.log(integration_area / (base_res +
6     eps)) / math.log(scale + eps)

```

```

6 l_target = torch.clamp(torch.ceil(M_value), min=0,
7     max=self.pyramid_levels - 1)
8 w = l_target - M_value
9 # CRITICAL ERROR: Wrong formula - should be
10 pixel_area * t^2

```

Listing 107. C2 Implementation: Paper2Code (Score: 0.8)

```

1 # AutoP2C: Placeholder that doesn't work
2 def compute_projected_area(self, x):
3     # Placeholder for actual projected area
4     computation
5     return torch.ones_like(x[...], 0])
6 # CRITICAL ERROR: Always returns 1.0, no level
7 selection

```

Listing 108. C2 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY and GPT-5 implement correct level selection. Paper2Code's critical error in area computation undermines the entire method.

**Component C3: Linear Interpolation Between Levels - All Baselines** The paper specifies linear interpolation:  $\sigma, c = w f_l(x, d) + (1 - w) f_{l-1}(x, d)$ .

```

1 # \nerfify\ : Perfect interpolation implementation
2 rgb_main = torch.gather(rgb_all, 0, idx_main).squeeze
3 (0)
4 rgb_prev = torch.gather(rgb_all, 0, idx_prev).squeeze
5 (0)
6 if self.interp_mode == "interpolate":
7     rgb = w * rgb_main + (1.0 - w) * rgb_prev
8 else: # Laplacian mode
9     rgb_cum = torch.cumsum(rgb_all, dim=0)
10     rgb = torch.gather(rgb_cum, 0, idx_main).squeeze
11 (0)
12 # CORRECT: Exact Equation 5 with multiple modes

```

Listing 109. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Clean interpolation
2 def forward(self, x, d, deltas, Px):
3     l, w = compute_level(Px, self.cfg.base_resolution
4     , self.cfg.scale_factor)
5     c1, s1 = self.models[l-1](x, d)
6     c2, s2 = self.models[l](x, d)
7     c = w.unsqueeze(-1) * c2 + (1 - w).unsqueeze(-1)
8     * c1
9     sigma = w * s2 + (1 - w) * s1
10     return volume_render(c, sigma, deltas)
11 # CORRECT: Proper linear interpolation

```

Listing 110. C3 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Partial implementation
2 if self.interpolation_method == "linear" and
3     level_idx > 0:
4     prev_output = self.levels[level_idx - 1](
5     level_samples, level_dirs)
6     level_weights = weights[mask].unsqueeze(-1)
7     level_output['density'] = (level_weights *
8     level_output['density'] +
9     (1 - level_weights) *
10     prev_output['density'])
11 # Note: Color interpolation similar but omitted
12 in provided code
13 # ISSUE: Incomplete - missing color interpolation
14 details

```

Listing 111. C3 Implementation: DeepSeek R1 (Score: 0.8)



```

1 # Paper2Code: Correct interpolation logic
2 if level_int == 0:
3     density, color = self.heads[0](x_level, d_level)
4 else:
5     density_upper, color_upper = self.heads[level_int
6     ](x_level, d_level)
7     density_lower, color_lower = self.heads[level_int
8     - 1](x_level, d_level)
9     density = w_level * density_upper + (1 - w_level)
10    * density_lower
11    color = w_level * color_upper + (1 - w_level) *
12    color_lower
13 # CORRECT: Proper interpolation between adjacent
14    levels

```

Listing 112. C3 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: Nonsensical interpolation
2 if self.multiscale_sampling:
3     # Linearly interpolate between the outputs of two
4     # closest hierarchy levels
5     if self.multiscale_sampling:
6         density = w * density + (1 - w) * density
7         color = w * color + (1 - w) * color
8 # CRITICAL ERROR: Interpolates with itself - no
9    effect

```

Listing 113. C3 Implementation: AutoP2C (Score: 0.1)

**Analysis:** NERFIFY, GPT-5, and Paper2Code correctly implement interpolation. AutoP2C’s self-interpolation is a fundamental misunderstanding.

**Component C4: Shared Multiresolution Features - All Baselines** The paper suggests reusing multi-resolution features across levels when the backbone already has a scale hierarchy.

```

1 # \nerfify\ : Shared hash encoding across all levels
2 self.hash_map = HashEncoding(
3     num_levels=num_hash_levels, min_res=min_res,
4     max_res=max_res,
5     log2_hashmap_size=log2_hashmap_size,
6     features_per_level=features_per_level)
7 # All pyramid levels query the same shared encoding:
8 feats = self.hash_map(pos_flat) # Used by all
9    density/color heads
10 # CORRECT: Single shared encoding as specified

```

Listing 114. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Same encoder reference for all levels
2 enc = HashEncoder(cfg.n_levels, cfg.
3     features_per_level,
4     cfg.base_resolution, cfg.
5     scale_factor, cfg.hash_size)
6 models = [PyNeRFLevel(enc, cfg) for _ in range(cfg.
7     n_levels)]
8 # ERROR: All levels share same encoder instance -
9    problematic

```

Listing 115. C4 Implementation: GPT-5 (Score: 0.2)

```

1 # DeepSeek R1: Configurable shared features
2 def __init__(self, use_shared_features: bool = True):
3     if not use_shared_features:
4         self.feature_grid = self._create_feature_grid
5         (base_config, level_idx)
6     else:
7         self.feature_grid = None # Will use shared
8         features from base NeRF
9 # ISSUE: Structure present but implementation
10    incomplete

```

Listing 116. C4 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: Single shared grid
2 self.shared_grid = True
3 self.feature_grid = FeatureGrid(resolution=self.
4     base_resolution,
5     num_features=
6     num_features,
7     hash_table_size=
8     hash_table_size)
9 # Each head references the same grid
10 head = PyNeRFHead(feature_grid=self.feature_grid)
11 # ISSUE: Simple hash table, not true multi-resolution

```

Listing 117. C4 Implementation: Paper2Code (Score: 0.6)

```

1 # AutoP2C: No feature sharing concept
2 self.positional_encoding = positional_encoding
3 x_encoded = self.apply_positional_encoding(x)
4 # MISSING: No shared features, uses positional
5    encoding only

```

Listing 118. C4 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY implements proper shared multi-resolution features. Others either misunderstand or omit this optimization.

#### Component C5: Per-Level Tiny MLPs - All Baselines

The paper specifies tiny MLPs per level: 64-channel density MLP (1 hidden layer) and 128-channel color MLP (2 hidden layers).

```

1 # \nerfify\ : Exact specifications
2 self.density_heads = nn.ModuleList([
3     MLP(in_dim=self.feature_dim, num_layers=1,
4         layer_width=density_hidden_dim, # 64
5         out_dim=1, activation=nn.ReLU())
6     for _ in range(self.num_pyramid_levels)])
7 self.color_heads = nn.ModuleList([
8     MLP(in_dim=color_in_dim, num_layers=2,
9         layer_width=color_hidden_dim, # 128
10        out_dim=3, activation=nn.ReLU(),
11        out_activation=nn.Sigmoid())
12    for _ in range(self.num_pyramid_levels)])
13 # CORRECT: Matches paper architecture exactly

```

Listing 119. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Simplified MLPs
2 class DensityMLP(nn.Module):
3     def __init__(self, in_dim, hidden_dim):
4         self.fc = nn.Sequential(nn.Linear(in_dim,
5             hidden_dim),
6             nn.ReLU(), nn.Linear(
7             hidden_dim, 1))
8 class ColorMLP(nn.Module):
9     def __init__(self, in_dim, hidden_dim):
10        self.fc = nn.Sequential(nn.Linear(in_dim + 3,
11            hidden_dim),
12            nn.ReLU(), nn.Linear(
13            hidden_dim, 3))
14 # ISSUE: Missing second hidden layer for color

```

Listing 120. C5 Implementation: GPT-5 (Score: 0.6)

```

1 # DeepSeek R1: TinyCUDANN networks
2 self.density_mlp = tcnn.Network(
3     n_input_dims=self._get_input_dims(config),
4     n_output_dims=1,

```

```

4     network_config={"otype": "FullyFusedMLP", "
      activation": "ReLU",
5         "n_neurons": 64, "n_hidden_layers"
      : 1})
6 self.color_mlp = tcnn.Network(
7     n_input_dims=self._get_input_dims(config) + 3,
8     n_output_dims=3,
9     network_config={"otype": "FullyFusedMLP", "
      n_neurons": 128,
10         "n_hidden_layers": 2})
11 # ISSUE: Using external library, not pure PyTorch

```

Listing 121. C5 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: Correct MLP specifications
2 density_layers = []
3 for hidden_dim in density_hidden_dims: # [64]
4     density_layers.append(nn.Linear(current_dim,
5         hidden_dim))
6     density_layers.append(nn.ReLU())
7 self.density_mlp = nn.Sequential(*density_layers)
8 # Color MLP with 2 hidden layers [128, 128]
9 for hidden_dim in color_hidden_dims:
10     color_layers.append(nn.Linear(current_dim,
11         hidden_dim))
12     color_layers.append(nn.ReLU())
13 # CORRECT: Proper layer counts and dimensions

```

Listing 122. C5 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: Single MLPs for everything
2 self.density_mlp = nn.Sequential(
3     nn.Linear(2 * 10 * 3, 64), nn.ReLU(), nn.Linear
4     (64, 1))
5 # MISSING: No per-level MLPs, just one shared

```

Listing 123. C5 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and Paper2Code match specifications exactly. GPT-5 and DeepSeek R1 have the right dimensions but implementation issues.

**Component C6: Projected Area Computation - All Baselines** The paper specifies computing  $P(x) = \text{pixel\_area} \times t^2$  where  $t$  is the sample distance along the ray.

```

1 # \nerfify\ : Exact paper formula
2 t_mid = (ray_samples.frustums.starts + ray_samples.
3     frustums.ends) * 0.5
4 pixel_area = ray_samples.frustums.pixel_area
5 if pixel_area.shape != t_mid.shape:
6     pixel_area = pixel_area.expand_as(t_mid)
7 proj_area = pixel_area * torch.clamp(t_mid, min=0.0)
8     ** 2 + 1e-12
9 # CORRECT: P(x) = pixel_area * t^2 as specified

```

Listing 124. C6 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Simplified approach
2 # In dummy loader:
3 torch.ones(1_024, 1).cuda() * 0.5 # projected area
4 # Missing actual computation, uses placeholder
5 # ISSUE: No real projected area calculation

```

Listing 125. C6 Implementation: GPT-5 (Score: 0.6)

```

1 # DeepSeek R1: Heuristic approximation
2 sample_distances = torch.norm(sample_positions - rays
3     ['origins'].unsqueeze(1), dim=-1)
4 # Approximate pixel area as distance^2 * solid_angle

```

```

4 pixel_areas = sample_distances ** 2 * 0.01 #
      heuristic scaling
5 # ERROR: Heuristic formula, not paper specification

```

Listing 126. C6 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: WRONG FORMULA
2 def compute_integration_area(self, depths):
3     depths = depths.to(torch.float32)
4     integration_area = (depths / self.focal) ** 2
5     return integration_area
6 # CRITICAL ERROR: Should be pixel_area * t^2, not (t/
7     f)^2

```

Listing 127. C6 Implementation: Paper2Code (Score: 0.2)

```

1 # AutoP2C: No implementation
2 def compute_projected_area(self, x):
3     # Placeholder for actual projected area
4     computation
5     return torch.ones_like(x[...], 0))
6 # CRITICAL ERROR: Always returns 1.0

```

Listing 128. C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY implements the correct formula. Paper2Code's error here is particularly damaging as it affects all level selections.

**Component C7: Adaptive Level Supervision - All Baselines** The paper mentions maintaining an auxiliary structure to track supervised levels during training.

```

1 # \nerfify\ : Occupancy grid from parent class
2 # Inherits from NerfactoModel which includes:
3 # self.occupancy_grid (from parent)
4 # self.proposal_networks (for sampling)
5 # Partial implementation through inheritance
6 # ISSUE: Not explicitly tracking supervised levels

```

Listing 129. C7 Implementation: NERFIFY (Score: 0.6)

```

1 # GPT-5: No adaptive supervision
2 # No occupancy grid or level tracking
3 # MISSING: Feature not implemented

```

Listing 130. C7 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: Structure present
2 self.occupancy_grid = None
3 self.supervised_levels = None
4 # For tracking supervised levels during training
5 # ISSUE: Declared but never used

```

Listing 131. C7 Implementation: DeepSeek R1 (Score: 0.4)

```

1 # Paper2Code: No adaptive supervision
2 # MISSING: No occupancy grid or adaptive training

```

Listing 132. C7 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No supervision concepts
2 # MISSING: Feature not implemented

```

Listing 133. C7 Implementation: AutoP2C (Score: 0.0)

**Analysis:** None of the baselines fully implement adaptive supervision. NERFIFY partially inherits it from Nerfacto.

**Component C8: Modern Backbone Integration - All Baselines** The paper mentions integration with modern grid-based NeRFs like iNGP, K-Planes, TensorRF.

```
1 # \nerfify\ : Full Nerfstudio integration
2 from nerfstudio.models.nerfacto import NerfactoModel,
  NerfactoModelConfig
3 class PynerfModel(NerfactoModel):
4     """Nerfacto-based model with PyNeRF multiscale
      field."""
5     def populate_modules(self):
6         super().populate_modules()
7         # Reuse scene contraction from Nerfacto
8         spatial_distortion = getattr(self.field, "
      spatial_distortion", None)
9 # CORRECT: Complete framework integration
```

Listing 134. C8 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: No framework integration
2 # Standalone implementation
3 # MISSING: Cannot leverage existing infrastructure
```

Listing 135. C8 Implementation: GPT-5 (Score: 0.0)

```
1 # DeepSeek R1: Multiple backbone options
2 def _create_base_nerf(self, config):
3     nerf_type = config.get('type', 'ingp')
4     if nerf_type == 'ingp':
5         return InstantNGPBackbone(config)
6     elif nerf_type == 'tensorrf':
7         return TensorRFBackbone(config)
8     elif nerf_type == 'kplanes':
9         return KPlanesBackbone(config)
10 # ISSUE: Stubs only, not implemented
```

Listing 136. C8 Implementation: DeepSeek R1 (Score: 0.8)

```
1 # Paper2Code: No backbone integration
2 # Custom standalone implementation
3 # MISSING: No framework compatibility
```

Listing 137. C8 Implementation: Paper2Code (Score: 0.0)

```
1 # AutoP2C: No backbone concepts
2 # MISSING: Feature not implemented
```

Listing 138. C8 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY achieves real framework integration. DeepSeek R1 has the structure but lacks implementation.

**Component C9: Multiple Interpolation Modes - All Baselines** The paper describes both linear interpolation (Equation 5) and Laplacian pyramid summation (Equation 4).

```
1 # \nerfify\ : Both interpolation modes
2 if self.interp_mode == "sum":
3     # Laplacian pyramid style (Equation 4)
4     sigma_cum = torch.cumsum(sigma_all, dim=0)
5     sigma = torch.gather(sigma_cum, 0, idx_main).
      squeeze(0)
6 else:
7     # Linear interpolation (Equation 5)
8     sigma = w * sigma_main + (1.0 - w) * sigma_prev
9 # CORRECT: Both modes implemented
```

Listing 139. C9 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Only linear interpolation
2 c = w.unsqueeze(-1) * c2 + (1 - w).unsqueeze(-1) * c1
3 # MISSING: No Laplacian mode
```

Listing 140. C9 Implementation: GPT-5 (Score: 0.0)

```
1 # DeepSeek R1: Mode selection present
2 def __init__(self, interpolation_method: str = "
      linear"):
3     self.interpolation_method = interpolation_method
4 # ISSUE: Only linear mode implemented
```

Listing 141. C9 Implementation: DeepSeek R1 (Score: 0.4)

```
1 # Paper2Code: No mode selection
2 # Only implements linear interpolation
3 # MISSING: No Laplacian option
```

Listing 142. C9 Implementation: Paper2Code (Score: 0.0)

```
1 # AutoP2C: No interpolation modes
2 # MISSING: Feature not implemented
```

Listing 143. C9 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY implements both interpolation modes as described in the paper.

**Component C10: Specific Hyperparameters - All Baselines** The paper specifies: 8 pyramid levels, scale factor 2, base resolution 16, 4 features per level, hash table size  $2^{20}$ .

```
1 # \nerfify\ : All hyperparameters correct
2 num_pyramid_levels: int = 8
3 pyramid_scale_factor: float = 2.0
4 features_per_level: int = 4
5 log2_hashmap_size: int = 20
6 density_hidden_dim: int = 64
7 color_hidden_dim: int = 128
8 max_num_iterations=20000
9 train_num_rays_per_batch=8192
10 # CORRECT: Matches paper exactly
```

Listing 144. C10 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Configuration matches
2 @dataclass
3 class PyNeRFConfig:
4     n_levels: int = 8
5     base_resolution: int = 16
6     scale_factor: float = 2.0
7     features_per_level: int = 4
8     hash_size: int = 2**20
9     mlp_hidden_dim: int = 64
10    color_hidden_dim: int = 128
11 # CORRECT: All parameters specified
```

Listing 145. C10 Implementation: GPT-5 (Score: 1.0)

```
1 # DeepSeek R1: Complete configuration
2 'pyramid': {'num_levels': 8, 'scale_factor': 2.0,
3            'base_resolution': 16, '
      use_shared_features': True},
4 'base_nerf': {'log2_hash_size': 19, '
      features_per_level': 4}
5 'training': {'batch_size': 8192}
6 # CORRECT: Parameters match paper
```

Listing 146. C10 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Most parameters correct
2 self.pyramid_levels: int = int(model_config.get("
    pyramid_levels", 8))
3 self.base_resolution: int = int(model_config.get("
    base_resolution", 256)) # Wrong!
4 self.scale_factor: int = int(model_config.get("
    scale_factor", 2))
5 # ERROR: Base resolution should be 16, not 256

```

Listing 147. C10 Implementation: Paper2Code (Score: 0.6)

```

1 # AutoP2C: Wrong parameters
2 hierarchy_levels=8 # Only this is correct
3 base_resolution=1 # Wrong! Should be 16
4 scaling_factor=2
5 # CRITICAL ERROR: Multiple wrong values

```

Listing 148. C10 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY, GPT-5, and DeepSeek R1 correctly specify all hyperparameters. Paper2Code and AutoP2C have critical errors in base resolution.

#### 4.3.6. Scoring Analysis

Table 38. PyNeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	Weighted AvgLLM
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	0.6	1.0	1.0	1.0	0.97
GPT-5	1.0	1.0	1.0	0.2	0.6	0.6	0.0	0.0	0.0	1.0	0.52
DeepSeek R1	1.0	0.6	0.8	0.6	0.6	0.6	0.4	0.8	0.4	1.0	0.68
Paper2Code	1.0	0.8	1.0	0.6	1.0	0.2	0.0	0.0	0.0	0.6	0.58
AutoP2C	0.0	0.0	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.03

#### 4.3.7. Why Baselines Fail Despite Component Scores

##### GPT-5 (Score: 0.52 components, 0% trainable)

- **Strengths:** Clean modular implementation with correct pyramid structure and interpolation logic
- **Fatal Issues:**
  - No Nerfstudio integration - standalone code cannot leverage existing infrastructure
  - Missing training loop and data loading pipeline
  - Dummy data loader prevents actual training
- **Result:** Code compiles but cannot train on real data
- **DeepSeek R1 (Score: 0.68 components, 0% trainable)**
- **Strengths:** Comprehensive structure with multi-agent architecture consideration
- **Fatal Issues:**
  - Many methods left unimplemented with pass statements
  - Heuristic projected area computation ( $\text{sample\_distances}^2 * 0.01$ )
  - Incomplete PyNeRFLevel forward pass
- **Result:** Structural skeleton without functional implementation
- **Paper2Code (Score: 0.58 components, 0% trainable)**
- **Strengths:** Most complete baseline with proper pyramid architecture and detailed implementation
- **Fatal Issues:**
  - Critical error in projected area formula: uses  $(\text{depth}/\text{focal})^2$  instead of  $\text{pixel\_area} * t^2$

- No framework integration - standalone implementation
- Missing adaptive supervision and occupancy grid
- **Result:** Despite 1250 lines of code, mathematical errors prevent convergence
- **AutoP2C (Score: 0.03 components, 0% trainable)**
- **Strengths:** Attempts to create training infrastructure
- **Fatal Issues:**
  - Complete architectural misunderstanding - no pyramid structure
  - Placeholder functions that always return constants
  - Import errors for non-existent modules (colmap Python module)
  - Self-interpolation bug (interpolates values with themselves)
- **Result:** Complete failure to understand the paper’s core concepts

#### 4.3.8. Hyperparameter Fidelity

Table 39. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
Pyramid Levels (L)	8	✓	✓	✓	✓	×
Scale Factor (s)	2	✓	✓	✓	✓	×
Base Resolution	16	✓	✓	✓	~	×
Features per Level	4	✓	✓	~	✓	×
Hash Table Size	2 <sup>20</sup>	✓	✓	×	✓	×
Density MLP Hidden	64	✓	✓	✓	✓	✓
Color MLP Hidden	128	✓	✓	✓	✓	✓
Training Iterations	20,000	✓	✓	×	✓	✓
Batch Size	8,192	✓	~	✓	✓	✓
<b>W Score</b>	–	1.00	0.80	0.70	0.60	0.00

#### 4.3.9. Conclusion

All baselines attempt to implement PyNeRF with varying levels of sophistication. GPT-5 produces the cleanest modular code but lacks framework integration. DeepSeek R1 creates an ambitious architecture that remains largely unimplemented. Paper2Code achieves the most complete standalone implementation but fails due to a critical mathematical error in projected area computation. AutoP2C completely misunderstands the paper’s core pyramid concept. Only NERFIFY produces immediately trainable code as a complete Nerfstudio plugin with correct mathematical formulations, achieving a 97% semantic completeness score. This stark contrast - 100% trainable rate for NERFIFY versus 0% for all baselines - demonstrates that effective paper-to-code translation for complex vision research requires deep domain specialization beyond what general-purpose systems can provide.

### 4.4. TensorRF: Tensorial Radiance Fields

#### 4.4.1. Paper Overview

TensorRF [3] revolutionizes neural radiance fields by replacing computationally expensive MLPs with factorized 4D tensors using vector-matrix (VM) decomposition. The method

decomposes the radiance field tensor  $\mathcal{T} \in \mathbb{R}^{I \times J \times K \times P}$  into compact low-rank components, achieving 100× compression with superior quality. This explicit grid representation with tensor factorization demonstrates that properly structured explicit methods can outperform implicit neural representations in both reconstruction quality and training efficiency, achieving 30-minute training on a single GPU.

#### 4.4.2. Implementation Overview

Table 40. TensorRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	467	298	412	856	234
File Organization	Plugin	Single	Single	Multi-file	Multi-file
VM Decomposition	✓	✓	Partial	×	×
Density Factorization	✓	✓	Partial	×	×
Appearance Matrix B	✓	✓	×	×	×
L1 Regularization	✓	Partial	Partial	×	×
TV Regularization	✓	×	×	×	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

#### 4.4.3. Novel Components

Table 41. Novel Components in TensorRF with Importance Weights

ID	Component	Weight $w_i$
C1	Vector-Matrix (VM) Decomposition	0.20
C2	Density Grid Factorization	0.15
C3	Appearance Grid Factorization	0.15
C4	Global Appearance Matrix B	0.10
C5	Factor-level Trilinear Interpolation	0.10
C6	L1 Sparsity Regularization	0.08
C7	Total Variation Regularization	0.07
C8	Coarse-to-Fine Upsampling	0.05
C9	SH/MLP Decoder	0.05
C10	Separate Density/Appearance Architecture	0.05

#### 4.4.4. Quantitative Metrics

Table 42. TensorRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	0.95	0.98
GPT-5	0.70	0.10	0.20	0.75	0.72
DeepSeek R1	0.60	0.20	0.20	0.70	0.65
Paper2Code	0.20	0.30	0.50	0.30	0.12
AutoP2C	0.10	0.20	0.70	0.15	0.28

#### 4.4.5. Component-by-Component Analysis

**Component C1: Vector-Matrix Decomposition - All Baselines**

The paper defines VM decomposition as:  $\mathcal{T} = \sum_{r=1}^R \mathbf{v}_r^X \circ \mathbf{M}_r^{YZ} + \mathbf{v}_r^Y \circ \mathbf{M}_r^{XZ} + \mathbf{v}_r^Z \circ \mathbf{M}_r^{XY}$  where  $\circ$  denotes outer product.

```
1 # \nerfify\ : Complete VM decomposition with all
   three modes
2 self.vx_sigma = nn.Parameter(0.1 * torch.randn(self.
   R_sigma, self.Nx))
3 self.vy_sigma = nn.Parameter(0.1 * torch.randn(self.
   R_sigma, self.Ny))
4 self.vz_sigma = nn.Parameter(0.1 * torch.randn(self.
   R_sigma, self.Nz))
5 self.MyZ_sigma = nn.Parameter(0.1 * torch.randn(self.
   R_sigma, self.Ny, self.Nz))
6 self.MxZ_sigma = nn.Parameter(0.1 * torch.randn(self.
   R_sigma, self.Nx, self.Nz))
7 self.MxY_sigma = nn.Parameter(0.1 * torch.randn(self.
   R_sigma, self.Nx, self.Ny))
8 Ax = (vx * MyZ).sum(dim=0) # v_r^X o M_r^YZ
9 Ay = (vy * MxZ).sum(dim=0) # v_r^Y o M_r^XZ
10 Az = (vz * MxY).sum(dim=0) # v_r^Z o M_r^XY
11 sigma_raw = Ax + Ay + Az + self.density_bias
12 # CORRECT: Exact VM decomposition from paper Eq. 7
```

Listing 149. C1 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Correct math but non-standard PyTorch usage
2 class VMDecomp(nn.Module):
3     def forward(self, xyz):
4         x, y, z = xyz
5         vx = F.linear(x, self.vectors[0].T)
6         myz = F.bilinear(y, z, self.matrices[0])
7         result += (vx * myz).sum(dim=-1)
8 # ISSUE: F.bilinear doesn't work this way in PyTorch
```

Listing 150. C1 Implementation: GPT-5 (Score: 0.9)

```
1 # DeepSeek R1: Structure present but flawed execution
2 class VMDecomposition(nn.Module):
3     def forward(self, xyz):
4         vx = F.linear(x, self.vectors[0].T)
5         myz = F.bilinear(y, z, self.matrices[0])
6         result += (vx * myz).sum(dim=-1)
7 # CRITICAL ERROR: Bilinear function misused
```

Listing 151. C1 Implementation: DeepSeek R1 (Score: 0.6)

```
1 # Paper2Code: Wrong architecture - implements NeRF
2 class TensorRFModel(nn.Module):
3     def __init__(self, decomposition_type='VM'):
4         self.density_mlp = self._build_mlp(64, 1)
5         self.color_mlp = self._build_mlp(128, 2)
6 # CRITICAL ERROR: Uses MLPs, not tensor factorization
   !
```

Listing 152. C1 Implementation: Paper2Code (Score: 0.0)

```
1 # AutoP2C: No model at all
2 class BlenderDataset:
3     def load_images(self):
4         self.images.append(image)
5 # MISSING: Entire TensorRF model absent
```

Listing 153. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** VM decomposition is the core innovation. Only NERFIFY implements it perfectly. GPT-5 has correct math but PyTorch issues. DeepSeek R1 attempts but fails. Paper2Code/AutoP2C completely miss it.

**Component C2: Density Grid Factorization - All Baselines**



The density grid  $\mathcal{G}_\sigma \in \mathbb{R}^{I \times J \times K}$  uses VM decomposition with  $R_\sigma = 16$  components.

```
1 # \nerfify\ : Perfect density factorization
2 def get_density(self, ray_samples: RaySamples):
3     x, y, z = self._flatten_positions(ray_samples)
4     vx = _interp_vector(self.vx_sigma, x)
5     vy = _interp_vector(self.vy_sigma, y)
6     vz = _interp_vector(self.vz_sigma, z)
7     MyZ = _interp_plane(self.MyZ_sigma, y, z)
8     MxZ = _interp_plane(self.MxZ_sigma, x, z)
9     MxY = _interp_plane(self.MxY_sigma, x, y)
10    Ax = (vx * MyZ).sum(dim=0)
11    Ay = (vy * MxZ).sum(dim=0)
12    Az = (vz * MxY).sum(dim=0)
13    sigma = trunc_exp(Ax + Ay + Az + self.
14                      density_bias)
15 # CORRECT: Density grid with VM decomposition and
16     trunc_exp
```

Listing 154. C2 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Correct density factorization
2 def get_density(self, xyz):
3     return F.softplus(self.density(xyz))
4 # CORRECT: Proper density with softplus activation
```

Listing 155. C2 Implementation: GPT-5 (Score: 0.9)

```
1 # DeepSeek R1: Mostly correct but missing details
2 def forward_density(self, pos):
3     for r in range(self.R):
4         density_vals = density_vals + term_x + term_y
5         + term_z
6     return density_vals
7 # ISSUE: Missing activation function
```

Listing 156. C2 Implementation: DeepSeek R1 (Score: 0.8)

```
1 # Paper2Code: Wrong - uses MLPs
2 def forward(self, x, d):
3     encoded_x = self.encode_position(x)
4     density = self.density_mlp(encoded_x)
5 # CRITICAL ERROR: MLP with positional encoding
```

Listing 157. C2 Implementation: Paper2Code (Score: 0.0)

```
1 # AutoP2C: No density implementation
2 # MISSING: No density code found
```

Listing 158. C2 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Density factorization correctly splits from appearance. NERFIFY and GPT-5 implement correctly, DeepSeek R1 partially, Paper2Code/AutoP2C fail completely.

### Component C3: Appearance Grid Factorization - All Baselines

Appearance grid  $\mathcal{G}_c \in \mathbb{R}^{I \times J \times K \times P}$  with  $P = 27$  features and  $R_c = 48$  components.

```
1 # \nerfify\ : Complete appearance factorization
2 self.vx_c = nn.Parameter(0.1 * torch.randn(self.R_c,
3     self.Nx))
4 self.vy_c = nn.Parameter(0.1 * torch.randn(self.R_c,
5     self.Ny))
6 self.vz_c = nn.Parameter(0.1 * torch.randn(self.R_c,
7     self.Nz))
8 self.MyZ_c = nn.Parameter(0.1 * torch.randn(self.R_c,
9     self.Ny, self.Nz))
```

```
6 self.MxZ_c = nn.Parameter(0.1 * torch.randn(self.R_c,
7     self.Nx, self.Nz))
8 self.MxY_c = nn.Parameter(0.1 * torch.randn(self.R_c,
9     self.Nx, self.Ny))
10 concat_contrib = torch.cat([Ax, Ay, Az], dim=0).
11     transpose(0, 1)
12 app_feat = self.B(concat_contrib) # [M, P=27]
13 # CORRECT: Separate appearance grid with feature
14     dimension 27
```

Listing 159. C3 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Present but simplified
2 self.appearance = VMDecomposition(
3     config.resolution,
4     config.num_appearance_components)
5 # ISSUE: Missing feature dimension details
```

Listing 160. C3 Implementation: GPT-5 (Score: 0.7)

```
1 # DeepSeek R1: Partial implementation
2 def forward_appearance(self, pos):
3     for r in range(self.R):
4         feat_list.append(term_x)
5         feat_list.append(term_y)
6         feat_list.append(term_z)
7 # ERROR: Missing matrix B transformation
```

Listing 161. C3 Implementation: DeepSeek R1 (Score: 0.6)

```
1 # Paper2Code: No appearance grid
2 color = self.color_mlp(encoded_x)
3 # CRITICAL ERROR: Uses MLP, not factorized grid
```

Listing 162. C3 Implementation: Paper2Code (Score: 0.0)

```
1 # AutoP2C: Missing appearance
2 # MISSING: No appearance implementation
```

Listing 163. C3 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Appearance requires separate VM decomposition with 48 components and 27 features. Only NERFIFY gets all details right.

### Component C4: Global Appearance Matrix B - All Baselines

Matrix  $\mathbf{B} \in \mathbb{R}^{P \times 3R_c}$  maps concatenated features to final appearance.

```
1 # \nerfify\ : Perfect matrix B implementation
2 self.B = nn.Linear(3 * self.R_c, self.P, bias=False)
3 concat_contrib = torch.cat([Ax, Ay, Az], dim=0).
4     transpose(0, 1)
5 app_feat = self.B(concat_contrib) # [M, P]
6 # CORRECT: Matrix B as linear layer without bias
```

Listing 164. C4 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Correct matrix B
2 self.B = nn.Parameter(
3     torch.randn(config.appearance_dim,
4         3 * config.num_appearance_components)
5 )
6 features = torch.matmul(vm_features, self.B.T)
7 # CORRECT: Proper matrix multiplication
```

Listing 165. C4 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Missing matrix B
2 # No global appearance matrix in code
3 # CRITICAL ERROR: Matrix B completely absent

```

Listing 166. C4 Implementation: DeepSeek R1 (Score: 0.0)

```

1 # Paper2Code: No matrix B
2 # CRITICAL ERROR: Architecture doesn't include matrix B

```

Listing 167. C4 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No matrix B
2 # MISSING: No model implementation

```

Listing 168. C4 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Matrix B is crucial for appearance. Only NERFIFY and GPT-5 implement it. Its absence severely impacts quality in other baselines.

#### Component C5: Factor-level Trilinear Interpolation - All Baselines

Efficient interpolation at factor level rather than dense grid:  $\mathcal{A}_r^X(\mathbf{x}) = \mathbf{v}_r^X(x) \cdot \mathbf{M}_r^{YZ}(y, z)$

```

1 # \nerfify\ : Efficient factor-level interpolation
2 def _interp_vector(vec: Tensor, coord: Tensor) -> Tensor:
3     scaled = coord.clamp(0.0, 1.0) * (N - 1)
4     i0 = torch.floor(scaled).long().clamp(0, N - 1)
5     i1 = (i0 + 1).clamp(0, N - 1)
6     w = (scaled - i0.float()).clamp(0.0, 1.0)
7     v0 = torch.gather(vec, 1, idx0)
8     v1 = torch.gather(vec, 1, idx1)
9     return (1.0 - w)[None, :] * v0 + w[None, :] * v1
10 # CORRECT: Efficient linear interpolation on factors

```

Listing 169. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct interpolation approach
2 def interp(self, v, x):
3     idx0 = torch.clamp(x.long(), 0, v.shape[1]-2)
4     w = x - idx0.float()
5     return (1-w)*v[:, idx0] + w*v[:, idx0+1]
6 # CORRECT: Factor interpolation implemented

```

Listing 170. C5 Implementation: GPT-5 (Score: 0.9)

```

1 # DeepSeek R1: Interpolation present but inefficient
2 def _interp1d(self, vec: torch.Tensor, coord: torch.Tensor):
3     coord = coord.clamp(0, length - 1)
4     idx0 = torch.floor(coord).long()
5     t = coord - idx0.float()
6 # ISSUE: Suboptimal implementation

```

Listing 171. C5 Implementation: DeepSeek R1 (Score: 0.7)

```

1 # Paper2Code: No factor interpolation
2 # CRITICAL ERROR: Uses positional encoding instead

```

Listing 172. C5 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No interpolation
2 # MISSING: No model code

```

Listing 173. C5 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Factor-level interpolation is key to efficiency. NERFIFY and GPT-5 implement correctly, DeepSeek R1 partially.

#### Component C6: L1 Sparsity Regularization - All Baselines

L1 regularization on density factors with weight  $\omega = 0.0004$ .

```

1 # \nerfify\ : Complete L1 on all density factors
2 ll_terms = [
3     self.vx_sigma.abs().mean(),
4     self.vy_sigma.abs().mean(),
5     self.vz_sigma.abs().mean(),
6     self.MxZ_sigma.abs().mean(),
7     self.MxY_sigma.abs().mean(),
8     self.MxY_sigma.abs().mean(),
9 ]
10 losses["l1_density"] = sum(ll_terms) / len(ll_terms)
11 loss_dict["l1_density"] = self.config.l1_density_mult * reg["l1_density"]
12 # CORRECT: L1 with exact weight 4e-4

```

Listing 174. C6 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: L1 present but simplified
2 ll_loss = 0.0
3 for param in model.density.parameters():
4     ll_loss += torch.abs(param).mean()
5 total_loss = render_loss + config.l1_weight * ll_loss
6 # ISSUE: Hardcoded weight

```

Listing 175. C6 Implementation: GPT-5 (Score: 0.7)

```

1 # DeepSeek R1: Wrong weight
2 def compute_l1_reg(self):
3     ll_reg += torch.mean(torch.abs(param))
4 total_loss = mse_loss + 1e-4 * reg_loss
5 # ERROR: Uses 1e-4 instead of 4e-4

```

Listing 176. C6 Implementation: DeepSeek R1 (Score: 0.5)

```

1 # Paper2Code: No L1 regularization
2 # MISSING: No regularization terms

```

Listing 177. C6 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No regularization
2 # MISSING: No training code

```

Listing 178. C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** L1 sparsity crucial for compact models. Only NERFIFY gets weight exact, others have errors or missing.

#### Component C7: Total Variation Regularization - All Baselines

TV regularization for smooth factors, especially important for few-shot scenarios.

```

1 # \nerfify\ : Complete TV regularization
2 def _tv_1d(self, vec: Tensor) -> Tensor:
3     return (vec[:, 1:] - vec[:, :-1]).abs().mean()
4 def _tv_2d(self, mat: Tensor) -> Tensor:
5     tv_u = (mat[:, 1:, :] - mat[:, :-1, :]).abs().mean()
6     tv_v = (mat[:, :, 1:] - mat[:, :, :-1]).abs().mean()
7     return 0.5 * (tv_u + tv_v)
8 tv_density = (self._tv_1d(self.vx_sigma) + self._tv_1d(self.vy_sigma) +

```

```

9         self._tv_1d(self.vz_sigma) + self.
        _tv_2d(self.MyZ_sigma) +
10         self._tv_2d(self.MxZ_sigma) + self.
        _tv_2d(self.MxY_sigma) / 6.0
11 # CORRECT: TV on both 1D vectors and 2D matrices

```

Listing 179. C7 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: No TV regularization
2 # MISSING: TV not implemented

```

Listing 180. C7 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: No TV regularization
2 # MISSING: TV completely absent

```

Listing 181. C7 Implementation: DeepSeek R1 (Score: 0.0)

```

1 # Paper2Code: No TV
2 # MISSING: No regularization

```

Listing 182. C7 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No TV
2 # MISSING: No model

```

Listing 183. C7 Implementation: AutoP2C (Score: 0.0)

**Analysis:** TV regularization critical for few-shot. Only NERFIFY implements it correctly with both 1D and 2D variants.

### Component C8: Coarse-to-Fine Upsampling - All Baselines

Progressive resolution increase from  $128^3$  to  $300^3$  during training.

```

1 # \nerfify\ : Upsampling via training schedule
2 upsampling_schedule: [2000, 3000, 4000, 5500, 7000]
3 grid_resolution: (128, 128, 128) # Initial
4 final_resolution: 300 # Target
5 # Training config handles progressive upsampling
6 # ISSUE: Not explicit upsampling, handled by scheduler

```

Listing 184. C8 Implementation: NERFIFY (Score: 0.8)

```

1 # GPT-5: No coarse-to-fine
2 # MISSING: No upsampling implementation

```

Listing 185. C8 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: No upsampling
2 # MISSING: No coarse-to-fine schedule

```

Listing 186. C8 Implementation: DeepSeek R1 (Score: 0.0)

```

1 # Paper2Code: No upsampling
2 # MISSING: No progressive resolution

```

Listing 187. C8 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No upsampling
2 # MISSING: No training pipeline

```

Listing 188. C8 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Coarse-to-fine speeds convergence. NERFIFY implements via schedule, others completely miss it.

### Component C9: SH/MLP Decoder - All Baselines

View-dependent color decoder using spherical harmonics or small MLP.

```

1 # \nerfify\ : Complete SH encoding + MLP decoder
2 self.direction_encoding = SHEncoding(levels=3,
        implementation="torch")
3 self.color_head = MLP(
4     in_dim=self.P + self.direction_encoding.
        get_out_dim(),
5     num_layers=color_mlp_layers,
6     layer_width=color_mlp_width,
7     out_dim=3,
8     activation=nn.ReLU(),
9     out_activation=nn.Sigmoid())
10 # CORRECT: SH encoding with 2-layer MLP decoder

```

Listing 189. C9 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: MLP decoder present
2 self.mlp = nn.Sequential(
3     nn.Linear(3 * config.num_appearance_components,
4         128),
5     nn.ReLU(),
6     nn.Linear(128, 3),
7     nn.Sigmoid())
7 # ISSUE: Missing SH option

```

Listing 190. C9 Implementation: GPT-5 (Score: 0.7)

```

1 # DeepSeek R1: Partial decoder
2 self.mlp = nn.Sequential(
3     nn.Linear(mlp_input_dim, mlp_hidden_size),
4     nn.ReLU(),
5     nn.Linear(mlp_hidden_size, 3))
6 # ERROR: Missing SH, incomplete MLP

```

Listing 191. C9 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: Wrong decoder type
2 self.color_mlp = self._build_mlp(128, 2)
3 # CRITICAL ERROR: Wrong architecture

```

Listing 192. C9 Implementation: Paper2Code (Score: 0.2)

```

1 # AutoP2C: Evaluation mentions MLP
2 # MISSING: No actual implementation

```

Listing 193. C9 Implementation: AutoP2C (Score: 0.1)

**Analysis:** Decoder maps appearance features to RGB. NERFIFY implements both SH and MLP options correctly.

### Component C10: Separate Density/Appearance Architecture - All Baselines

Independent factorizations for density and appearance with different ranks.

```

1 # \nerfify\ : Completely separate density/appearance
2 self.density_factor = TensorFactorization(
3     grid_shape=grid_shape,
4     num_components={"density": 16}, # R_sigma = 16
5     decomposition=decomposition,
6     grid_type="density")
7 self.appearance_factor = TensorFactorization(
8     grid_shape=grid_shape,
9     num_components={"appearance": 48}, # R_c = 48
10    decomposition=decomposition,
11    grid_type="appearance")

```

```
12 # CORRECT: Separate factors with different ranks
```

Listing 194. C10 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Separate but less clear
2 self.density = VMDecomposition(
3     config.resolution,
4     config.num_density_components)
5 self.appearance = VMDecomposition(
6     config.resolution,
7     config.num_appearance_components)
8 # CORRECT: Separate architectures
```

Listing 195. C10 Implementation: GPT-5 (Score: 0.8)

```
1 # DeepSeek R1: Attempted separation
2 self.density_factor = TensorFactorization(
3     grid_type="density")
4 self.appearance_factor = TensorFactorization(
5     grid_type="appearance")
6 # ISSUE: Implementation details incomplete
```

Listing 196. C10 Implementation: DeepSeek R1 (Score: 0.7)

```
1 # Paper2Code: No separation
2 # CRITICAL ERROR: Single MLP for both
```

Listing 197. C10 Implementation: Paper2Code (Score: 0.0)

```
1 # AutoP2C: No architecture
2 # MISSING: No model
```

Listing 198. C10 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Separation allows different ranks for density (16) and appearance (48). Critical for efficiency/quality balance.

#### 4.4.6. Scoring Analysis

Table 43. TensorRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.8	1.0	1.0	0.98
GPT-5	0.9	0.9	0.7	1.0	0.9	0.7	0.0	0.0	0.7	0.8	0.72
DeepSeek R1	0.6	0.8	0.6	0.0	0.7	0.5	0.0	0.0	0.6	0.7	0.65
Paper2Code	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.0	0.12
AutoP2C	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.0	0.28

#### 4.4.7. Why Baselines Fail Despite Component Scores

##### GPT-5 (Score: 0.72 components, 0% trainable)

- **Strengths:** Correct VM decomposition, matrix B implementation, density/appearance separation
- **Fatal Issues:**
  - No Nerfstudio integration - requires complete pipeline assembly
  - Missing TV regularization entirely (C7 = 0.0)
  - No coarse-to-fine upsampling (C8 = 0.0)
- **Result:** Mathematically sound but requires hours of engineering to make trainable
- **DeepSeek R1 (Score: 0.65 components, 0% trainable)**
- **Strengths:** Understands tensor decomposition concept, attempts proper separation

##### Fatal Issues:

- PyTorch tensor operations cause runtime crashes (F.bilinear misuse)
- Missing matrix B completely (C4 = 0.0)
- No TV regularization (C7 = 0.0) causes training instability

- **Result:** Training crashes immediately on tensor shape mismatches

##### Paper2Code (Score: 0.12 components, 0% trainable)

- **Strengths:** Basic volume rendering loop exists

##### Fatal Issues:

- Complete architectural misunderstanding - implements vanilla NeRF
- Zero correct components except partial decoder (C9 = 0.2)
- Uses MLPs with positional encoding throughout

- **Result:** Trains as NeRF, not TensorRF - fundamentally wrong method

##### AutoP2C (Score: 0.28 components, 0% trainable)

- **Strengths:** Dataset loading code present

##### Fatal Issues:

- No model implementation whatsoever
- Only evaluation stub mentions components (C9 = 0.1)
- Missing entire TensorRF architecture

- **Result:** Code doesn't compile - no model to train

#### 4.4.8. Hyperparameter Fidelity

Table 44. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
Grid Resolution	128 <sup>3</sup>	✓	✓	×	–	–
Density Rank $R_d$	16	✓	×	✓	–	–
Appearance Rank $R_a$	48	✓	✓	×	–	–
Feature Dim P	27	✓	✓	✓	–	–
L1 Weight	4e-4	✓	✓	×	–	–
TV Weight	Variable	✓	×	×	–	–
Learning Rate	0.02	×	✓	✓	×	–
Training Steps	30K	✓	×	×	–	–
MLP Width	128	✓	✓	✓	×	–
<b>W Score</b>	–	0.95	0.75	0.70	0.30	0.15

#### 4.4.9. Conclusion

All baselines attempt to implement TensorRF with dramatically varying levels of success. GPT-5 achieves mathematically correct VM decomposition and matrix B but lacks critical TV regularization and framework integration. DeepSeek R1 understands the tensor factorization concept but fails catastrophically on PyTorch execution and missing matrix B. Paper2Code fundamentally misunderstands the method, implementing vanilla NeRF with MLPs and positional encoding instead of tensor factorization. AutoP2C provides only dataset loading code with no model implementation whatsoever. Despite GPT-5 achieving 72% component coverage and DeepSeek R1 reaching 65%, none produce trainable

code due to missing Nerfstudio integration, incorrect tensor operations, or complete architectural misunderstandings. Only NERFIFY produces immediately trainable code as a complete Nerfstudio plugin, achieving 100% trainable rate versus 0% for all baselines, demonstrating the critical importance of domain-specific synthesis for complex computer vision research.

#### 4.5. Tetra-NeRF: Representing Neural Radiance Fields Using Tetrahedra

##### 4.5.1. Paper Overview

Tetra-NeRF [12] introduces an adaptive scene representation using Delaunay triangulation of point clouds, replacing uniform voxel grids with tetrahedra that naturally concentrate resolution near surfaces. The method achieves faster training and rendering by using barycentric interpolation within tetrahedra, eliminating the need for dense regular grids while maintaining quality comparable to state-of-the-art methods.

##### 4.5.2. Implementation Overview

Table 45. Tetra-NeRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	412	189	267	845	621
File Organization	Plugin	Single	Single	Multi-file	Multi-file
Delaunay Triangulation	✓	✓	✓	Attempted	×
Barycentric Coords	✓	Simplified	✓	×	×
Feature Interpolation	✓	✓	✓	Partial	×
Volume Rendering	✓	~	~	×	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

##### 4.5.3. Novel Components

Table 46. Novel Components in Tetra-NeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Delaunay triangulation for adaptive subdivision	0.20
C2	Barycentric coordinates using volume ratios (Eq. 2)	0.15
C3	Feature interpolation via barycentric weights	0.15
C4	Tetrahedral grid data structure	0.10
C5	Adaptive non-uniform resolution	0.10
C6	Hierarchical coarse-to-fine sampling	0.10
C7	Volume rendering with tetrahedral queries	0.08
C8	Tetrahedron volume computation	0.05
C9	Nerfstudio plugin integration	0.05
C10	RAdam optimizer with exponential decay	0.02

##### 4.5.4. Quantitative Metrics

##### 4.5.5. Component-by-Component Analysis

###### Component C1: Delaunay Triangulation - All Baselines

The paper specifies: “We represent the scene as a dense triangulation of the input point cloud, where the scene is a

Table 47. Tetra-NeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	1.00	1.00
GPT-5	0.50	0.25	0.25	0.60	0.58
DeepSeek R1	0.63	0.25	0.13	0.70	0.72
Paper2Code	0.13	0.25	0.63	0.20	0.22
AutoP2C	0.00	0.13	0.88	0.00	0.08

set of non-overlapping tetrahedra whose union is the convex hull of the original point cloud.”

```

1 # \nerfify\ : Proper tetrahedral grid with simplex
  decomposition
2 class TetraNerfField(Field):
3     def _tetrahedral_interpolate(self, positions_norm
4     ):
5         # 6-simplex decomposition of unit cube
6         scales = torch.tensor([self.nx-1, self.ny-1,
7         self.nz-1])
8         gpos = positions_norm * scales
9         base = torch.floor(gpos).long()
10        u = (gpos - base.to(gpos.dtype)).contiguous()
11        # Sort u components for simplex determination
12        u_sorted, perm = torch.sort(u, dim=1,
13        descending=False)
14        k1, k2, k3 = perm[:, 0], perm[:, 1], perm[:,
15        2]
16 # CORRECT: Full tetrahedral decomposition with proper
17 simplex

```

Listing 199. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Basic Delaunay implementation
2 def build_tetrahedra(points):
3     tri = Delaunay(points)
4     verts = torch.tensor(points, dtype=torch.float32,
5     device=device)
6     tets = torch.tensor(tri.simplices, dtype=torch.
7     long, device=device)
8     return verts, tets
9 # CORRECT: Uses scipy.spatial.Delaunay correctly

```

Listing 200. C1 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Complete tetrahedral grid class
2 class TetrahedralGrid:
3     def __init__(self, points: torch.Tensor):
4         self.points = points.cpu().numpy()
5         self.tetrahedra = Delaunay(self.points)
6         self.features = nn.Parameter(
7         torch.randn(len(points), 64) * 1e-4)
8 # CORRECT: Proper Delaunay with feature
9 initialization

```

Listing 201. C1 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Wrong approach
2 class DatasetLoader:
3     def perform_delaunay_triangulation(self):
4         delaunay = scipy.spatial.Delaunay(self.
5         point_cloud)
6         self.tetrahedra = delaunay.simplices
7 # ERROR: Not integrated with rendering pipeline

```

Listing 202. C1 Implementation: Paper2Code (Score: 0.2)



```

1 # AutoP2C: No triangulation at all
2 class ColmapDataset:
3     def generate_point_cloud(self):
4         self.point_cloud = np.array([point.xyz for
5             point in reconstruction.
6             points3D.values()])
7 # CRITICAL ERROR: Just loads points, no tetrahedra

```

Listing 203. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY implements efficient simplex decomposition. GPT-5 and DeepSeek R1 correctly use Delaunay. Paper2Code attempts but fails to integrate. AutoP2C misses entirely.

### Component C2: Barycentric Coordinates - All Baselines

The paper defines barycentric coordinates using volume ratios:  $\lambda_i = \frac{V_i}{V_{\text{total}}}$  where  $V_i$  is the volume of sub-tetrahedron.

```

1 # \nerfify\ : Efficient simplex-based computation
2 def _tetrahedral_interpolate(self, positions_norm):
3     u_sorted, perm = torch.sort(u, dim=1, descending=False)
4     # Barycentric weights for simplex vertices
5     w0 = 1.0 - u_sorted[:, 2]
6     w1 = u_sorted[:, 2] - u_sorted[:, 1]
7     w2 = u_sorted[:, 1] - u_sorted[:, 0]
8     w3 = u_sorted[:, 0]
9     # Weighted interpolation
10    f = (f0 * w0.unsqueeze(-1) + f1 * w1.unsqueeze(-1) +
11        f2 * w2.unsqueeze(-1) + f3 * w3.unsqueeze(-1))
12 # CORRECT: Efficient vectorized barycentric computation

```

Listing 204. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Linear system approach
2 def barycentric_coords(p, tet_verts):
3     V0, V1, V2, V3 = tet_verts
4     M = torch.stack([V1-V0, V2-V0, V3-V0], dim=1)
5     v = p - V0
6     sol = torch.linalg.solve(M, v)
7     l1, l2, l3 = sol
8     l0 = 1 - l1 - l2 - l3
9     return torch.stack([l0, l1, l2, l3])
10 # ISSUE: No volume validation, less robust

```

Listing 205. C2 Implementation: GPT-5 (Score: 0.4)

```

1 # DeepSeek R1: Volume-based as in paper
2 def tetrahedron_volume(self, a, b, c, d):
3     return torch.abs(torch.dot(a - d,
4         torch.cross(b - d, c - d))) /
5         6.0
6 def get_barycentric_coords(self, query_points,
7     tetra_indices):
8     V_total = self.tetrahedron_volume(v0, v1, v2, v3)
9     V_x234 = self.tetrahedron_volume(query_points[i],
10     v1, v2, v3)
11     V_lx34 = self.tetrahedron_volume(v0, query_points
12     [i], v2, v3)
13     barycentric_coords[i] = torch.tensor([V_x234/
14     V_total, ...])
15 # CORRECT: Exact paper formula with volume ratios

```

Listing 206. C2 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Missing implementation
2 def compute_barycentrics(self, query_point,
3     simplex_index):
4     transform_matrix = self.delaunay.transform[
5     simplex_index]
6     T = transform_matrix[:ndim, :]
7     offset = transform_matrix[ndim, :]
8     # Function incomplete, never returns values
9 # CRITICAL ERROR: Never computes actual coordinates

```

Listing 207. C2 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No barycentric coordinates
2 # Uses completely different approach
3 # MISSING: Core mathematical concept entirely

```

Listing 208. C2 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY uses efficient simplex computation. DeepSeek R1 implements exact volume method. GPT-5 uses valid but less robust approach. Paper2Code and AutoP2C fail completely.

### Component C3: Feature Interpolation - All Baselines

The paper specifies interpolating features:  $f(q) = \sum_{i=0}^3 \lambda_i f_i$

```

1 # \nerfify\ : Efficient vectorized interpolation
2 def _gather_grid(self, idx):
3     ix, iy, iz = idx.unbind(-1)
4     grid_flat = self.features.view(-1, self.
5     vertex_feature_dim)
6     lin = ix * (self.ny * self.nz) + iy * self.nz +
7     iz
8     return grid_flat[lin]
9 # Gather and interpolate
10 f0 = self._gather_grid(v0)
11 f1 = self._gather_grid(v1)
12 f = (f0 * w0.unsqueeze(-1) + f1 * w1.unsqueeze(-1) +
13     f2 * w2.unsqueeze(-1) + f3 * w3.unsqueeze(-1))
14 # CORRECT: Vectorized for efficiency

```

Listing 209. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct but inefficient
2 def interpolate_feature(p, verts, feats, tets):
3     centroids = verts[tets].mean(dim=1)
4     idx = torch.argmax(torch.norm(centroids - p, dim
5     =1))
6     tet = tets[idx]
7     lambdas = barycentric_coords(p, verts[tet])
8     f = (lambdas[:, None] * feats[tet]).sum(0)
9     return f
10 # ISSUE: O(n) search instead of O(log n)

```

Listing 210. C3 Implementation: GPT-5 (Score: 0.8)

```

1 # DeepSeek R1: Correct interpolation
2 def interpolate_features(self, query_points):
3     for i in range(batch_size):
4         if tetra_indices[i] >= 0:
5             tetra_vertices = self.tetrahedra.
6             simplices[tetra_indices[i]]
7             vertex_features = self.features[
8             tetra_vertices]
9             interpolated_features[i] = torch.sum(
10             barycentric[i].unsqueeze(1) *
11             vertex_features, dim=0)
12 # CORRECT: Proper weighted sum

```

Listing 211. C3 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Wrong method
2 def batch_barycentric_interpolation(self,
  query_points):
3     for i in range(query_points.shape[0]):
4         try:
5             feat = self.barycentric_interpolation(pt,
              device=device)
6         except ValueError:
7             feat = torch.zeros(self.vertex_features.
              shape[1])
8 # ERROR: Returns zeros instead of interpolating

```

Listing 212. C3 Implementation: Paper2Code (Score: 0.3)

```

1 # AutoP2C: No feature interpolation
2 # MISSING: Core rendering component

```

Listing 213. C3 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY provides vectorized implementation. GPT-5 correct but inefficient. DeepSeek R1 correct with loops. Paper2Code fails. AutoP2C missing.

#### Component C4: Tetrahedral Grid Structure - All Baselines

The paper requires a data structure to manage vertices, features, and connectivity.

```

1 # \nerfify\ : Complete TetraNerfField class
2 class TetraNerfField(Field):
3     def __init__(self, aabb, grid_resolution=32,
4         vertex_feature_dim=64,
5         bottleneck_dim=32):
6         super().__init__()
7         self.register_buffer("aabb", aabb)
8         self.nx = self.ny = self.nz = grid_resolution
9         self.features = nn.Parameter(torch.empty(
10             self.nx, self.ny, self.nz,
11             vertex_feature_dim))
12         nn.init.uniform_(self.features, a=-1e-4, b=1e
13             -4)
14         self.base_mlp = MLP(in_dim=vertex_feature_dim
15             ,
16                             num_layers=3, layer_width
17                             =128)
18 # CORRECT: Complete field with Nerfstudio integration

```

Listing 214. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Functions only, no proper class
2 def build_tetrahedra(points):
3     tri = Delaunay(points)
4     verts = torch.tensor(points)
5     tets = torch.tensor(tri.simplices)
6 def interpolate_feature(p, verts, feats, tets):
7     # Separate functions
8 # ISSUE: No unified data structure

```

Listing 215. C4 Implementation: GPT-5 (Score: 0.5)

```

1 # DeepSeek R1: Has structure but incomplete
2 class TetrahedralGrid:
3     def __init__(self, points):
4         self.points = points.cpu().numpy()
5         self.tetrahedra = Delaunay(self.points)
6         self.features = nn.Parameter(
7             torch.randn(len(points), 64) * 1e-4)
8     def find_tetrahedron(self, query_points):
9         query_np = query_points.cpu().numpy()
10        tetra_indices = self.tetrahedra.find_simplex(
11            query_np)
12 # CORRECT: Good structure but missing integration

```

Listing 216. C4 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: Wrong structure
2 class TetrahedraField(nn.Module):
3     def __init__(self, point_cloud,
4         vertex_feature_dim=64):
5         self.points = point_cloud
6         self.delaunay = None # Never initialized
7 # CRITICAL ERROR: Empty module

```

Listing 217. C4 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No grid structure
2 # MISSING: Core data structure

```

Listing 218. C4 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY provides complete structure. GPT-5 lacks organization. DeepSeek R1 has structure but incomplete. Paper2Code and AutoP2C fail.

#### Component C5: Adaptive Resolution - All Baselines

The paper achieves adaptive resolution through Delaunay triangulation concentrating tetrahedra near surfaces.

```

1 # \nerfify\ : Implicit adaptive via simplex
2   decomposition
3 def _tetrahedral_interpolate(self, positions_norm):
4     # Scale to grid coordinates
5     scales = torch.tensor([self.nx-1, self.ny-1, self
6         .nz-1])
7     gpos = positions_norm * scales
8     base = torch.floor(gpos).long()
9     # Smaller tetrahedra near surfaces due to
10    decomposition
11    base[:, 0].clamp_(0, self.nx - 2)
12    base[:, 1].clamp_(0, self.ny - 2)
13    base[:, 2].clamp_(0, self.nz - 2)
14 # CORRECT: Achieves adaptivity through decomposition

```

Listing 219. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: No adaptive handling
2 t_vals = torch.linspace(0.0, 1.0, n_samples)
3 # Uniform sampling without adaptation
4 # MISSING: Adaptive resolution concept

```

Listing 220. C5 Implementation: GPT-5 (Score: 0.2)

```

1 # DeepSeek R1: Structure supports but not implemented
2 self.tetrahedra = Delaunay(self.points)
3 # Natural adaptivity from Delaunay but not utilized
4 # ISSUE: Has potential but not exploited

```

Listing 221. C5 Implementation: DeepSeek R1 (Score: 0.5)

```

1 # Paper2Code: Opposite of adaptive
2 # Uses uniform grid
3 # CRITICAL ERROR: Misses key innovation

```

Listing 222. C5 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: Random points
2 random_points = np.random.normal(
3     loc=self.point_cloud.mean(axis=0),
4     scale=self.point_cloud.std(axis=0))
5 # MISSING: Any spatial structure

```

Listing 223. C5 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY achieves adaptive resolution. DeepSeek R1 has structure but doesn't exploit. Others miss innovation.

### Component C6: Hierarchical Sampling - All Baselines

The paper uses coarse-to-fine sampling with importance weighting.

```
1 # \nerfify\ : Inherits from Nerfacto's samplers
2 class TetraNerfModel(NerfactoModel):
3     def populate_modules(self):
4         Model.populate_modules(self)
5         super().populate_modules() # Gets proposal
        networks
6         # Uses Nerfacto's production hierarchical
        sampling
7         self.field = TetraNerfField(aabb=self.
        scene_box.aabb)
8 # CORRECT: Production-quality hierarchical sampling
```

Listing 224. C6 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Basic implementation
2 def hierarchical_sampling(ray_o, ray_d, model,
        n_coarse=32, n_fine=32):
3     t_coarse = torch.linspace(0, 1, n_coarse)
4     # Fine sampling proportional to weights
5     cdf = torch.cumsum(w_norm.squeeze(), 0)
6     u = torch.rand(n_fine, device=device)
7     t_fine = torch.interp(u, cdf, t_coarse)
8 # ISSUE: Simplified, missing details
```

Listing 225. C6 Implementation: GPT-5 (Score: 0.6)

```
1 # DeepSeek R1: Complete but inefficient
2 def volume_render(self, rays_o, rays_d, near, far):
3     # Coarse sampling
4     t_vals_coarse = torch.linspace(near, far, self.
        config.num_samples_coarse)
5     weights = self.compute_weights(density_coarse,
        t_vals_coarse)
6     # Fine sampling
7     t_vals_fine = self.sample_fine(t_vals_coarse,
        weights)
8 # CORRECT: Right idea but loop-based
```

Listing 226. C6 Implementation: DeepSeek R1 (Score: 0.8)

```
1 # Paper2Code: Only mentions
2 def stratified_sampling(self, nears, fars,
        num_samples):
3     t_bins = torch.linspace(0.0, 1.0, steps=
        num_samples + 1)
4     # Only coarse, no fine
5 # ERROR: Missing importance sampling
```

Listing 227. C6 Implementation: Paper2Code (Score: 0.2)

```
1 # AutoP2C: Random sampling
2 def generate_rays(self):
3     for pose in self.camera_poses:
4         rays = self._generate_rays_for_pose(pose)
5 # MISSING: Hierarchical concept
```

Listing 228. C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY leverages production samplers. GPT-5 and DeepSeek R1 attempt with varying success. Paper2Code and AutoP2C fail.

### Component C7: Volume Rendering - All Baselines

The paper requires volume rendering with tetrahedral field queries.

```
1 # \nerfify\ : Complete volume rendering
2 def get_outputs(self, ray_samples, density_embedding=
        None):
3     dirs = get_normalized_directions(ray_samples.
        frustums.directions)
4     dirs_flat = dirs.view(-1, 3)
5     denc = self.direction_encoding(dirs_flat)
6     h = torch.cat([denc, density_embedding.view(-1,
        self.bottleneck_dim)], dim=-1)
7     rgb = self.color_head(h).view(*dirs.shape[:-1],
        -1).to(dirs)
8     outputs[FieldHeadNames.RGB] = rgb
9     return outputs
10 # CORRECT: Full integration with Nerfstudio rendering
```

Listing 229. C7 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Basic volume rendering
2 def volume_render(ray_o, ray_d, model, verts, feats,
        tets, n_samples=64):
3     t_vals = torch.linspace(0.0, 1.0, n_samples)
4     pts = ray_o[None, :] + ray_d[None, :] * t_vals[:,
        None]
5     weights = sigmas * T * deltas[:, None]
6     rgb = (weights * colors).sum(0)
7     return rgb
8 # CORRECT: Basic but functional
```

Listing 230. C7 Implementation: GPT-5 (Score: 0.7)

```
1 # DeepSeek R1: Volume rendering with issues
2 def compute_weights(self, density, t_vals):
3     delta = t_vals[..., 1:] - t_vals[..., :-1]
4     alpha = 1 - torch.exp(-density * delta)
5     transmittance = torch.cumprod(1 - alpha + 1e-10,
        dim=-1)
6     weights = alpha * transmittance
7 # ISSUE: Incorrect transmittance computation
```

Listing 231. C7 Implementation: DeepSeek R1 (Score: 0.6)

```
1 # Paper2Code: Wrong rendering
2 def volume_rendering(self, t_vals, sigma, rgb):
3     # Missing proper implementation
4 # CRITICAL ERROR: No actual rendering
```

Listing 232. C7 Implementation: Paper2Code (Score: 0.0)

```
1 # AutoP2C: No volume rendering
2 # MISSING: Core rendering pipeline
```

Listing 233. C7 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY has full integration. GPT-5 basic but functional. DeepSeek R1 has errors. Paper2Code and AutoP2C missing.

### Component C8: Tetrahedron Volume - All Baselines

The paper computes tetrahedron volume:  $V = \frac{1}{6} |(\mathbf{a} - \mathbf{d}) \cdot ((\mathbf{b} - \mathbf{d}) \times (\mathbf{c} - \mathbf{d}))|$

```
1 # \nerfify\ : Implicit in simplex decomposition
2 # Volume computation implicit in barycentric weights
3 w0 = 1.0 - u_sorted[:, 2] # Volume ratio
4 w1 = u_sorted[:, 2] - u_sorted[:, 1]
5 w2 = u_sorted[:, 1] - u_sorted[:, 0]
6 w3 = u_sorted[:, 0]
7 # CORRECT: Efficient volume-based weights
```

Listing 234. C8 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: No explicit volume computation
2 # Uses linear algebra instead
3 # ISSUE: Missing volume formula

```

Listing 235. C8 Implementation: GPT-5 (Score: 0.3)

```

1 # DeepSeek R1: Exact formula
2 def tetrahedron_volume(self, a, b, c, d):
3     return torch.abs(torch.dot(a - d,
4                               torch.cross(b - d, c - d))) /
5         6.0
6 # CORRECT: Exact paper formula

```

Listing 236. C8 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: No volume computation
2 # MISSING: Mathematical component

```

Listing 237. C8 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: Missing entirely
2 # MISSING: Core formula

```

Listing 238. C8 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY implicit in weights. DeepSeek R1 exact formula. GPT-5 missing. Paper2Code and AutoP2C absent.

#### Component C9: Nerfstudio Integration - All Baselines

The paper implementation requires Nerfstudio plugin structure for training.

```

1 # \nerfify\ : Complete Nerfstudio plugin
2 tetra_nerf = MethodSpecification(
3     config=TrainerConfig(
4         method_name="tetra-nerf",
5         pipeline=TetraNerfPipelineConfig(
6             datamanager=TetraNerfDataManagerConfig(),
7             model=TetraNerfModelConfig(),
8             optimizers={
9                 "fields": {"optimizer":
10                     RAdamOptimizerConfig(lr=1e-3)}},
11             description="Tetra-NeRF: Nerfacto-based with
12                 tetrahedral field")
13 # CORRECT: Complete plugin specification

```

Listing 239. C9 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: No Nerfstudio integration
2 if __name__ == "__main__":
3     # Standalone demo only
4     model, verts, feats, tets = train_tetra_nerf(pts,
5         rays_o, rays_d)
6 # MISSING: Plugin structure

```

Listing 240. C9 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: No framework integration
2 class TetraNeRFTrainer:
3     def __init__(self, model, config):
4         self.model = model
5         self.optimizer = torch.optim.Adam(model.
6             parameters())
7 # MISSING: Nerfstudio plugin

```

Listing 241. C9 Implementation: DeepSeek R1 (Score: 0.0)

```

1 # Paper2Code: Custom training only
2 class Trainer:
3     def train(self):
4         for iter_num in range(1, self.iterations + 1)
5             :
6 # MISSING: Framework integration

```

Listing 242. C9 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No framework awareness
2 # MISSING: Plugin structure

```

Listing 243. C9 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY has Nerfstudio integration. All baselines lack plugin structure, making them untrainable in framework.

#### Component C10: RAdam Optimizer - All Baselines

The paper specifies RAdam optimizer with exponential decay from 1e-3 to 1e-4.

```

1 # \nerfify\ : Exact paper specification
2 "fields": {
3     "optimizer": RAdamOptimizerConfig(lr=1e-3, eps=1e-8),
4     "scheduler": ExponentialDecaySchedulerConfig(
5         lr_final=1e-4, max_steps=300000),
6 }
7 # CORRECT: RAdam with proper decay

```

Listing 244. C10 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct optimizer
2 self.optimizer = torch.optim.RAdam(
3     list(model.parameters()) + [feats], lr=1e-3)
4 self.scheduler = torch.optim.lr_scheduler.
5     ExponentialLR(
6         self.optimizer, gamma=0.999)
7 # CORRECT: RAdam with decay

```

Listing 245. C10 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Wrong optimizer
2 self.optimizer = torch.optim.Adam(model.parameters(),
3     lr=config.learning_rate)
4 self.scheduler = torch.optim.lr_scheduler.
5     ExponentialLR(
6         self.optimizer, gamma=0.1)
7 # ERROR: Uses Adam instead of RAdam

```

Listing 246. C10 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: Correct optimizer setup
2 self.optimizer = torch.optim.RAdam(self.model.
3     parameters(), lr=self.lr_initial)
4 def adjust_learning_rate(self, current_iter):
5     decay_factor = (self.lr_final / self.lr_initial)
6     **
7         (current_iter / self.
8         lr_decay_steps)
9 # CORRECT: RAdam with exponential decay

```

Listing 247. C10 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: Mentions RAdam
2 # Paper uses RAdam with exponential LR decay
3 # But implementation missing
4 # ISSUE: No actual optimizer code

```

Listing 248. C10 Implementation: AutoP2C (Score: 0.8)

**Analysis:** NERFIFY , GPT-5, and Paper2Code correctly implement RAdam. DeepSeek R1 uses Adam. AutoP2C mentions but doesn't implement.

#### 4.5.6. Scoring Analysis

Table 48. Tetra-NeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.00
GPT-5	1.0	0.4	0.8	0.5	0.2	0.6	0.7	0.3	0.0	1.0	0.58
DeepSeek R1	1.0	1.0	1.0	0.8	0.5	0.8	0.6	1.0	0.0	0.8	0.72
Paper2Code	0.2	0.0	0.3	0.0	0.0	0.2	0.0	0.0	0.0	1.0	0.22
AutoP2C	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.8	0.08

#### 4.5.7. Why Baselines Fail Despite Component Scores

##### GPT-5 (Score: 0.58 components, 0% trainable)

- **Strengths:** Implements core tetrahedral concepts correctly, proper Delaunay triangulation
- **Fatal Issues:**
  - No Nerfstudio integration whatsoever
  - Inefficient  $O(n)$  tetrahedron search instead of  $O(\log n)$
  - Missing training loop and data loading
- **Result:** Code runs in isolation but cannot be trained within Nerfstudio framework

##### DeepSeek R1 (Score: 0.72 components, 0% trainable)

- **Strengths:** Mathematically correct volume-based barycentric coordinates, proper feature interpolation
- **Fatal Issues:**
  - No Nerfstudio plugin structure
  - Loop-based processing causes memory overflow on real datasets
  - Incomplete volume rendering implementation
- **Result:** Theoretically correct but practically unusable for training

##### Paper2Code (Score: 0.22 components, 0% trainable)

- **Strengths:** Attempts multi-file organization, includes optimizer setup
- **Fatal Issues:**
  - Fundamental misunderstanding: implements voxel grid instead of tetrahedra
  - Missing barycentric coordinate computation entirely
  - PyTorch errors including device mismatches
- **Result:** Implements a completely different method than the paper specifies

##### AutoP2C (Score: 0.08 components, 0% trainable)

- **Strengths:** Includes COLMAP data loading
- **Fatal Issues:**
  - Missing all core tetrahedral concepts
  - No Delaunay triangulation or barycentric coordinates
  - Non-functional Python with import errors
- **Result:** Non-functional code that fails at import stage

Table 49. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
Feature dimension	64	✓	✓	✓	×	×
Hidden dimension	128	✓	✓	✓	×	×
MLP layers	3	✓	✓	✓	×	×
Coarse samples	64	✓	×	✓	×	–
Fine samples	128	✓	×	✓	×	–
Learning rate	1e-3	✓	✓	✓	×	–
Optimizer	RAdam	✓	×	×	×	–
<b>W Score</b>	–	1.00	0.60	0.70	0.20	0.00

#### 4.5.8. Hyperparameter Fidelity

##### 4.5.9. Conclusion

All baselines attempt to implement Tetra-NeRF with varying levels of sophistication. GPT-5 and DeepSeek R1 demonstrate understanding of core tetrahedral concepts, correctly implementing Delaunay triangulation and barycentric interpolation. Paper2Code fundamentally misunderstands the paper, implementing a voxel-based approach instead of tetrahedra. AutoP2C fails to implement any meaningful components beyond basic data loading. Despite GPT-5 achieving 58% and DeepSeek R1 achieving 72% component coverage, neither produces trainable code due to missing Nerfstudio integration, inefficient algorithms, and incomplete training pipelines. Only NERFIFY produces immediately trainable code as a complete Nerfstudio plugin with 100% component fidelity. The 0% trainable rate for all baselines versus 100% for NERFIFY demonstrates the critical importance of domain-specific understanding and framework integration in neural rendering research. ““

#### 4.6. E-NeRF: Neural Radiance Fields from a Moving Event Camera

##### 4.6.1. Paper Overview

E-NeRF introduces a groundbreaking approach to neural radiance field reconstruction using event cameras, which capture pixel-level brightness changes rather than absolute intensity. The method replaces traditional RGB photometric losses with event-generation constraints, enabling NeRF reconstruction from asynchronous event streams. This is particularly valuable for scenarios with fast motion, low light, or high dynamic range where conventional cameras fail. The paper demonstrates that event supervision alone can produce high-quality 3D reconstructions, with optional RGB frame integration for enhanced color fidelity.

##### 4.6.2. Implementation Overview

##### 4.6.3. Novel Components

##### 4.6.4. Quantitative Metrics

##### 4.6.5. Component-by-Component Analysis

**Component C1: Event-based Loss Function - All Baselines**



Table 50. E-Nerf Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	487	245	389	1650	156
File Organization	Plugin	Single	Single	Multi-file	Multi-file
Event Loss (Eq. 3)	✓	✓	✓	Partial	×
Linlog Mapping	✓	✓	✓	✓	×
No-Event Loss	✓	Partial	✓	×	×
Framework Integration	✓	×	×	×	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

Table 51. Novel Components in E-Nerf with Importance Weights

ID	Component	Weight $w_i$
C1	Event-based loss function (Eq. 3)	0.25
C2	Linlog brightness mapping (Eq. 2)	0.20
C3	Normalized event loss (Eq. 4)	0.15
C4	No-event loss (Eq. 5)	0.10
C5	Event + RGB combined loss (Eq. 6)	0.10
C6	Event pair sampling strategy	0.08
C7	Camera pose interpolation (slerp/cubic)	0.07
C8	Hash-based encoding (Instant-NGP)	0.05

Table 52. E-Nerf Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	0.95	1.00
GPT-5	0.50	0.25	0.25	0.75	0.60
DeepSeek R1	0.63	0.25	0.13	0.80	0.72
Paper2Code	0.38	0.25	0.38	0.60	0.48
AutoP2C	0.00	0.13	0.88	0.00	0.05

The paper’s core innovation is the event loss enforcing brightness change constraints:  $\mathcal{L}_{\text{evs}}(\Theta) = \|\Delta\hat{\mathbf{L}}(\Theta) - \Delta\mathbf{L}(\Theta)\|_2^2$  where  $\Delta L_k = p_k C$  with contrast threshold  $C = 0.2$ .

```

1 # \nerfify\ : Full integration with event batch
  handling
2 def _event_losses(self, batch):
3     rb_t0 = self._bundle_from_batch_prefix(batch, "
      event_t0")
4     rb_t1 = self._bundle_from_batch_prefix(batch, "
      event_t1")
5
6     I0 = self._render_intensity(rb_t0)
7     I1 = self._render_intensity(rb_t1)
8     L0 = _linlog(I0, B)
9     L1 = _linlog(I1, B)
10    delta_hat = (L1 - L0).reshape(-1)
11
12    delta_gt = (C * batch["event_polarity"]).reshape
      (-1)
13    ev_loss = F.mse_loss(delta_hat, delta_gt)
14    # CORRECT: Proper batch structure and loss
      computation

```

Listing 249. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Mathematically correct but no framework
  integration

```

```

2 def event_loss(deltaL_hat, deltaL_true, normalized=
  False):
3     if normalized:
4         deltaL_hat = deltaL_hat / (torch.norm(
          deltaL_hat) + 1e-8)
5         deltaL_true = deltaL_true / (torch.norm(
          deltaL_true) + 1e-8)
6     return F.mse_loss(deltaL_hat, deltaL_true)
7 # ISSUE: No batch handling, standalone function

```

Listing 250. C1 Implementation: GPT-5 (Score: 0.8)

```

1 # DeepSeek R1: Comprehensive but missing integration
2 def generate_event_loss(self, pred_log_brightness1,
  pred_log_brightness2,
3     measured_polarity, C=0.2):
4     pred_delta_L = pred_log_brightness2 -
      pred_log_brightness1
5     target_delta_L = measured_polarity * C
6     event_loss = F.mse_loss(pred_delta_L,
      target_delta_L)
7     return event_loss
8 # CORRECT: Formula implementation
9 # MISSING: Training loop integration

```

Listing 251. C1 Implementation: DeepSeek R1 (Score: 0.9)

```

1 # Paper2Code: Basic structure without proper batching
2 if self.event_loss_type == "standard":
3     event_loss = torch.mean((predicted_diff -
      event_target_diff_tensor) ** 2)
4 elif self.event_loss_type == "normalized":
5     event_loss = Utils.normalized_event_loss(
      predicted_diff, event_target_diff_tensor)
6 # ERROR: Missing batch structure, incomplete
  integration

```

Listing 252. C1 Implementation: Paper2Code (Score: 0.5)

```

1 # AutoP2C: No implementation
2 # File only contains dataset placeholders
3 # CRITICAL ERROR: Event loss not implemented at all

```

Listing 253. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY correctly implements the event loss with full Nerfstudio integration. GPT-5 and DeepSeek R1 have mathematically correct formulas but lack framework integration. Paper2Code attempts the loss but misses critical batch handling. AutoP2C completely lacks any loss implementation.

### Component C2: Linlog Brightness Mapping - All Baselines

The paper specifies:  $L(\mathbf{u}) = I(\mathbf{u}) \cdot \ln(B)/B$  if  $I < B$ , else  $\ln(I)$  with  $B = 20$ .

```

1 # \nerfify\ : Exact paper formula with proper epsilon
  handling
2 def _linlog(intensity: Tensor, B: float) -> Tensor:
3     eps = 1e-6
4     lnB_over_B = float(torch.log(torch.tensor(B)).
      item() / B)
5     I = torch.clamp(intensity, min=eps)
6     linear = I * lnB_over_B
7     logv = torch.log(I)
8     return torch.where(I < B, linear, logv)
9 # CORRECT: Exact implementation from paper

```

Listing 254. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct implementation with masking
2 def linlog(I, B=20.0):
3     mask = (I < B).float()
4     return mask * (I * torch.log(torch.tensor(B)) / B
5         ) + \
6         (1 - mask) * torch.log(I + 1e-6)
7 # CORRECT: Proper branching logic

```

Listing 255. C2 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Correct with explicit masking
2 def linlog(self, I):
3     linear_mask = I < self.B
4     log_mask = ~linear_mask
5     L = torch.zeros_like(I)
6     L[linear_mask] = I[linear_mask] * torch.log(torch
7         .tensor(self.B)) / self.B
8     L[log_mask] = torch.log(I[log_mask])
9     return L
10 # CORRECT: Accurate implementation

```

Listing 256. C2 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Verbose but correct
2 def lin_log_mapping(self, intensity: torch.Tensor) ->
3     torch.Tensor:
4     B_tensor = torch.tensor(B, device=intensity.
5         device)
6     return torch.where(
7         intensity < B_tensor,
8         intensity * torch.log(B_tensor) / B_tensor,
9         torch.log(intensity + 1e-6)
10    )
11 # CORRECT: Proper implementation

```

Listing 257. C2 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: Not implemented
2 # CRITICAL ERROR: No linlog mapping found in codebase

```

Listing 258. C2 Implementation: AutoP2C (Score: 0.0)

**Analysis:** All baselines except AutoP2C correctly implement linlog mapping. This relatively straightforward component shows that baselines can handle individual mathematical formulas but struggle with system integration.

### Component C3: Normalized Event Loss - All Baselines

For unknown contrast threshold:  $\mathcal{L}_{\text{norm}} = \left\| \frac{\Delta \hat{L}}{\|\Delta \hat{L}\|_2} - \frac{\Delta L}{\|\Delta L\|_2} \right\|_2^2$

```

1 # \nerfify\ : Config-switchable normalized variant
2 if self.config.use_normalized_event_loss:
3     eps = 1e-8
4     nhat = delta_hat / (torch.linalg.norm(delta_hat)
5         + eps)
6     ngt = delta_gt / (torch.linalg.norm(delta_gt) +
7         eps)
8     ev_loss = F.mse_loss(nhat, ngt)
9 else:
10     ev_loss = F.mse_loss(delta_hat, delta_gt)
11 # CORRECT: Proper L2 normalization before MSE

```

Listing 259. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Has normalized variant as parameter
2 def event_loss(deltaL_hat, deltaL_true, normalized=
3     False):
4     if normalized:

```

```

4         deltaL_hat = deltaL_hat / (torch.norm(
5             deltaL_hat) + 1e-8)
6         deltaL_true = deltaL_true / (torch.norm(
7             deltaL_true) + 1e-8)
8         return F.mse_loss(deltaL_hat, deltaL_true)
9 # CORRECT: But not integrated with training

```

Listing 260. C3 Implementation: GPT-5 (Score: 0.8)

```

1 # DeepSeek R1: Separate function for normalized loss
2 def normalized_event_loss(self, pred_delta_L,
3     target_delta_L):
4     pred_norm = pred_delta_L / (torch.norm(
5         pred_delta_L, dim=-1, keepdim=True) + 1e-8)
6     target_norm = target_delta_L / (torch.norm(
7         target_delta_L, dim=-1, keepdim=True) + 1e-8)
8     return F.mse_loss(pred_norm, target_norm)
9 # CORRECT: Proper implementation

```

Listing 261. C3 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: No normalized variant
2 # Only has standard event loss
3 # CRITICAL ERROR: Missing normalized loss entirely

```

Listing 262. C3 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: Not implemented
2 # CRITICAL ERROR: No loss functions at all

```

Listing 263. C3 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY, GPT-5, and DeepSeek R1 correctly implement the normalized variant crucial for real-world data. Paper2Code and AutoP2C completely miss this component.

### Component C4: No-Event Loss - All Baselines

Penalizes brightness changes in areas without events:

$$\mathcal{L}_{\text{noevs}} = \sum_k \text{relu}(|\hat{L}_k - \hat{L}_{k-1}| - C)$$

```

1 # \nerfify\ : Complete no-event handling
2 rb_n0 = self._bundle_from_batch_prefix(batch, "
3     noevent_t0")
4 rb_n1 = self._bundle_from_batch_prefix(batch, "
5     noevent_t1")
6 if rb_n0 is not None and rb_n1 is not None:
7     I0 = self._render_intensity(rb_n0)
8     I1 = self._render_intensity(rb_n1)
9     L0 = _linlog(I0, B)
10    L1 = _linlog(I1, B)
11    diff = torch.abs(L1 - L0) - C
12    noev_loss = torch.clamp(diff, min=0.0).mean()
13 # CORRECT: Full pipeline from sampling to loss

```

Listing 264. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Formula correct but no sampling
2 def no_event_loss(Lhat_k, Lhat_kml, C):
3     return F.relu(torch.abs(Lhat_k - Lhat_kml) - C).
4         mean()
5 # ERROR: Missing no-event location identification

```

Listing 265. C4 Implementation: GPT-5 (Score: 0.4)

```

1 # DeepSeek R1: Class-based implementation
2 class NoEventLoss:
3     def __init__(self, C=0.2):
4         self.C = C
5     def __call__(self, pred_log_brightness1,
6         pred_log_brightness2):
7         brightness_diff = torch.abs(
8             pred_log_brightness2 - pred_log_brightness1)

```

```

7         return F.relu(brightness_diff - self.C).mean
            ()
8 # CORRECT: Formula but missing sampling logic

```

Listing 266. C4 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: Not implemented
2 # Has no-event sampling but no loss computation
3 # CRITICAL ERROR: Missing no-event loss

```

Listing 267. C4 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: Not implemented
2 # CRITICAL ERROR: No loss implementation

```

Listing 268. C4 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY fully implements no-event loss with proper sampling. GPT-5 and DeepSeek R1 have correct formulas but miss the critical sampling component.

### Component C5: Event + RGB Combined Loss - All Baselines

Combines supervision:  $\mathcal{L} = \mathcal{L}_{\text{evs}} + \lambda_{\text{rgb}}\mathcal{L}_{\text{rgb}}$

```

1 # \nerfify\ : Integrated with Nerfacto pipeline
2 def get_loss_dict(self, outputs, batch, metrics_dict=
    None):
3     loss_dict = super().get_loss_dict(outputs, batch,
        metrics_dict)
4     # Scale RGB loss
5     if "rgb_loss" in loss_dict and self.config.
        lambda_rgb != 1.0:
6         loss_dict["rgb_loss"] = self.config.
            lambda_rgb * loss_dict["rgb_loss"]
7     # Add event losses
8     ev_loss, noev_loss = self._event_losses(batch)
9     if ev_loss is not None:
10        loss_dict["event_loss"] = ev_loss
11 # CORRECT: Proper loss combination with weights

```

Listing 269. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Mathematical combination only
2 def combined_loss(Levs, Lrgb=None, lambda_rgb=1.0):
3     if Lrgb is None:
4         return Levs
5     return Levs + lambda_rgb * Lrgb
6 # ERROR: No actual RGB rendering pipeline

```

Listing 270. C5 Implementation: GPT-5 (Score: 0.4)

```

1 # DeepSeek R1: Combined training step
2 def combined_training_step(self, event_batch,
    rgb_batch):
3     event_loss = self.event_training_step(event_batch
        )
4     rgb_loss = 0
5     if rgb_batch is not None:
6         # Simplified RGB loss computation
7         rgb_loss = rgb_loss / len(rgb_batch)
8     total_loss = event_loss + self.config.get('
        lambda_rgb', 1.0) * rgb_loss
9 # ISSUE: Incomplete RGB pipeline

```

Listing 271. C5 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: Not implemented
2 # CRITICAL ERROR: No combined loss logic

```

Listing 272. C5 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: Not implemented
2 # CRITICAL ERROR: No loss functions

```

Listing 273. C5 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY properly integrates event and RGB losses within a working framework. Others either lack RGB pipelines or miss the combination entirely.

### Component C6: Event Pair Sampling Strategy - All Baselines

Samples event pairs within temporal windows (1ms average, 60ms window).

```

1 # \nerfify\ : Integrated with DataManager
2 def _copy_event_tensors(self, src: Dict, dst: Dict)
    -> None:
3     for k, v in src.items():
4         if any(k.startswith(pref) for pref in
            EVENT_PREFIXES):
5             dst[k] = v.to(self.device) if hasattr(v,
                "to") else v
6 # CORRECT: Full event batch handling

```

Listing 274. C6 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Synthetic dummy dataset only
2 class SyntheticEventDataset(Dataset):
3     def __getitem__(self, idx):
4         uk = torch.rand(2)
5         tk, tk1 = torch.rand(1), torch.rand(1)
6         pk = random.choice([-1, 1])
7 # ERROR: No real event sampling implementation

```

Listing 275. C6 Implementation: GPT-5 (Score: 0.2)

```

1 # DeepSeek R1: Placeholder function
2 def sample_event_pairs(self, event_batch,
    max_time_window=60):
3     # Implementation depends on your event data
        structure
4     # This should sample neighboring events or events
        within time windows
5     pass
6 # CRITICAL ERROR: Not implemented

```

Listing 276. C6 Implementation: DeepSeek R1 (Score: 0.3)

```

1 # Paper2Code: Detailed sampling logic
2 def sample_event_pairs(self, n_samples: int) -> List[
    Dict[str, Any]]:
3     for _ in range(n_samples):
4         pixel = random.choice(valid_pixels)
5         events_list = self.events_by_pixel[pixel]
6         i = random.randint(0, len(events_list) - 2)
7         # Accumulate polarities within window
8         polarity_sum = sum([ev[1] for ev in
            events_list[i+1:chosen_index+1]])
9 # CORRECT: Good sampling but not integrated

```

Listing 277. C6 Implementation: Paper2Code (Score: 0.8)

```

1 # AutoP2C: Not implemented
2 # CRITICAL ERROR: No event handling

```

Listing 278. C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY integrates event sampling with Nerf-studio's data pipeline. Paper2Code has detailed sampling

logic but fails to integrate it. Others have minimal or no implementation.

#### Component C7: Camera Pose Interpolation - All Baselines

Uses spherical linear interpolation for rotation and cubic for translation.

```
1 # \nerfify\ : Leverages Nerfstudio's built-in
   interpolation
2 # Handled by Nerfstudio's camera system
3 # CORRECT: Uses framework's optimized implementation
```

Listing 279. C7 Implementation: NERFIFY (Score: 0.9)

```
1 # GPT-5: Correct slerp implementation
2 def slerp(q0, q1, t):
3     dot = torch.sum(q0 * q1, dim=-1, keepdim=True)
4     dot = torch.clamp(dot, -1.0, 1.0)
5     theta_0 = torch.acos(dot)
6     sin_theta_0 = torch.sin(theta_0)
7     s0 = torch.sin(theta_0 - theta_t) / sin_theta_0
8     s1 = torch.sin(theta_t) / sin_theta_0
9     return s0 * q0 + s1 * q1
10 # ISSUE: Missing cubic interpolation for translation
```

Listing 280. C7 Implementation: GPT-5 (Score: 0.6)

```
1 # DeepSeek R1: Not implemented
2 # Uses placeholder compute_rays function
3 # CRITICAL ERROR: No pose interpolation
```

Listing 281. C7 Implementation: DeepSeek R1 (Score: 0.0)

```
1 # Paper2Code: Full implementation
2 def _interpolate_pose(self, timestamp: float) ->
   torch.Tensor:
3     # Slerp for rotation
4     q_interp = Utils.slerp(q0, q1, weight)
5     R_interp = self.quaternion_to_matrix(q_interp)
6     # Cubic for translation
7     t_interp = Utils.cubic_interpolation(t0_vec,
   t1_vec, weight)
8 # CORRECT: Both slerp and cubic interpolation
```

Listing 282. C7 Implementation: Paper2Code (Score: 1.0)

```
1 # AutoP2C: Not implemented
2 # CRITICAL ERROR: No pose handling
```

Listing 283. C7 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Paper2Code surprisingly has the most complete pose interpolation. NERFIFY relies on Nerfstudio's implementation. GPT-5 has partial slerp. Others lack this component.

#### Component C8: Hash-based Encoding - All Baselines

Instant-NGP style multiresolution hash encoding with 16 levels.

```
1 # \nerfify\ : Uses Nerfstudio's optimized hash
   encoding
2 # Configured in NerfactoField with proper parameters
3 # L=16, log2_hashmap=19, N_max=6144
4 # CORRECT: Framework's efficient implementation
```

Listing 284. C8 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Custom hash encoder implementation
2 class HashEncoder(nn.Module):
3     def __init__(self, n_levels=16,
   n_features_per_level=2,
   log2_hashmap_size=19):
4         self.embeddings = nn.ParameterList([
5             nn.Parameter(torch.randn(self.hash_size,
   n_features_per_level) * 1e-4)
6             for _ in range(n_levels)
7         ])
8 # CORRECT: Full implementation from scratch
```

Listing 285. C8 Implementation: GPT-5 (Score: 1.0)

```
1 # DeepSeek R1: Uses tinycudann
2 self.encoding = tcnn.Encoding(
3     n_input_dims=3,
4     encoding_config={
5         "otype": "HashGrid",
6         "n_levels": 16,
7         "n_features_per_level": 2,
8         "log2_hashmap_size": 19
9     }
10 )
11 # CORRECT: Efficient tinycudann implementation
```

Listing 286. C8 Implementation: DeepSeek R1 (Score: 1.0)

```
1 # Paper2Code: Detailed hash grid implementation
2 class HashGridEncoder(nn.Module):
3     def __init__(self, num_levels=16, dims_per_level
   =2,
   hashmap_size_log2=19):
4         for i in range(num_levels):
5             resolution = N_min * ((N_max/N_min)**(i/(
   num_levels-1)))
6             self.resolutions.append(resolution)
7 # CORRECT: Complete custom implementation
```

Listing 287. C8 Implementation: Paper2Code (Score: 1.0)

```
1 # AutoP2C: Not implemented
2 # CRITICAL ERROR: No model at all
```

Listing 288. C8 Implementation: AutoP2C (Score: 0.0)

**Analysis:** All baselines except AutoP2C correctly implement hash encoding, showing this is a well-understood component. The difference lies in integration quality.

#### 4.6.6. Scoring Analysis

Table 53. E-NeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	0.9	1.0	1.00
GPT-5	0.8	1.0	0.8	0.4	0.4	0.2	0.6	1.0	0.60
DeepSeek R1	0.9	1.0	1.0	0.8	0.6	0.3	0.0	1.0	0.72
Paper2Code	0.5	1.0	0.0	0.0	0.0	0.8	1.0	1.0	0.48
AutoP2C	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.05

#### 4.6.7. Why Baselines Fail Despite Component Scores

**GPT-5 (Score: 0.60 components, 0% trainable)**

- **Strengths:** Near-perfect mathematical implementations, clean modular code
- **Fatal Issues:**

- No Nerfstudio integration - standalone PyTorch implementation
- Missing event data pipeline - uses synthetic dummy data
- Training loop runs only 3 demonstration epochs then terminates
- **Result:** Code compiles but cannot train on real event data  
**DeepSeek R1 (Score: 0.72 components, 0% trainable)**
- **Strengths:** Comprehensive component coverage, uses efficient tinycudann
- **Fatal Issues:**
  - Event sampling returns empty - placeholder function not implemented
  - Missing camera pose interpolation entirely
  - No actual training step - only function signatures
- **Result:** Most complete baseline but critical functions are stubs  
**Paper2Code (Score: 0.48 components, 0% trainable)**
- **Strengths:** Extensive 1650-line codebase, detailed event sampling logic
- **Fatal Issues:**
  - Import errors - references non-existent colmap Python module
  - Missing normalized event loss variant
  - Event sampling not connected to training loop
- **Result:** Extensive boilerplate that doesn't execute  
**AutoP2C (Score: 0.05 components, 0% trainable)**
- **Strengths:** Basic dataset structure attempted
- **Fatal Issues:**
  - No model implementation whatsoever
  - No loss functions defined
  - Only contains data loading placeholders
- **Result:** Complete failure - non-functional code

#### 4.6.8. Hyperparameter Fidelity

Table 54. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
Contrast threshold $C$	0.2	✓	✓	✓	✓	×
Linlog threshold $B$	20.0	✓	✓	✓	✓	×
Batch size	30096	4096*	8	Variable	30096	×
Hash levels $L$	16	✓	✓	✓	✓	×
$\log_2$ hashmap	19	✓	✓	✓	✓	×
$N_{\max}$	$B \cdot 2048$	✓	✓	Variable	✓	×
Learning rate	1e-3	1e-2*	✓	✓	✓	×
<b>W Score</b>	–	0.95	0.75	0.80	0.60	0.00

\*NERFIFY adapts parameters for Nerfstudio optimization

#### 4.6.9. Conclusion

All baselines attempt to implement E-NeRF with varying sophistication, from AutoP2C’s complete failure (0.05 score) to DeepSeek R1’s comprehensive but non-functional attempt (0.72 score). GPT-5 achieves mathematically correct components but lacks system integration, while Paper2Code generates extensive boilerplate missing critical training logic.

Only NERFIFY produces immediately trainable code by properly integrating event-specific losses into Nerfstudio’s proven NeRF framework, achieving perfect coverage (1.00 score) and full functionality. This stark contrast—0% trainable rate for all baselines versus 100% for NERFIFY—demonstrates that effective paper-to-code translation for event-based vision demands deep domain specialization over generic approaches.

### 4.7. StyleNeRF: A Style-Based 3D-Aware Generator for High-Resolution Image Synthesis

#### 4.7.1. Paper Overview

StyleNeRF [7] introduces a groundbreaking 3D-aware generative model that integrates neural radiance fields with style-based generation for high-resolution image synthesis. The paper addresses the computational bottleneck in NeRF-based image generation by proposing an efficient approximation that aggregates features in 2D before final color computation, enabling progressive upsampling. This approach achieves both high visual quality and 3D consistency while reducing rendering time from minutes to seconds.

#### 4.7.2. Implementation Overview

Table 55. StyleNeRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	742	385	412	1847	623
File Organization	Plugin	Single	Single	Multi-file	Multi-file
Style Modulation	✓	Simplified	✓	×	×
2D Upsampling Path	✓	Partial	✓	×	×
NeRF-path Regularization	✓	Wrong	×	×	×
Custom Upsampler	✓	~	Partial	×	×
View Direction Removal	✓	×	×	×	×
Nerfstudio Integration	✓	×	×	×	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

#### 4.7.3. Novel Components

Table 56. Novel Components in StyleNeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Style-conditioned NeRF field (Eq. 2-3)	0.20
C2	Fourier positional encoding (Eq. 1)	0.08
C3	Low-res feature rendering + 2D upsampling (Eq. 5-6)	0.22
C4	Custom upsampler with blur kernel (Eq. 7)	0.15
C5	NeRF-path regularization loss (Eq. 9)	0.15
C6	View direction removal for consistency	0.05
C7	Mapping network (StyleGAN2-based)	0.05
C8	NeRF++ foreground/background split	0.05
C9	Progressive training strategy (3-stage)	0.03
C10	Hyperparameter fidelity	0.02



Table 57. StyleNeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	1.00	0.98
GPT-5	0.40	0.30	0.30	0.55	0.52
DeepSeek R1	0.50	0.30	0.20	0.64	0.62
Paper2Code	0.30	0.40	0.30	0.46	0.28
AutoP2C	0.00	0.10	0.90	0.00	0.00

#### 4.7.4. Quantitative Metrics

#### 4.7.5. Component-by-Component Analysis

##### Component C1: Style-Conditioned NeRF Field - All Baselines

StyleNeRF introduces style modulation for NeRF MLPs using FiLM-style conditioning:

$$\phi_w^n(\mathbf{x}) = g_w^n \circ g_w^{n-1} \circ \dots \circ g_w^1 \circ \zeta(\mathbf{x})$$

where  $g_w^i$  are style-modulated layers with weight matrices modulated by style vector  $\mathbf{w}$ .

```

1 # \nerfify\ : Complete FiLM-style modulation with per
  -layer affine transforms
2 class StyleMLP(nn.Module):
3     def __init__(self, in_dim, width, depth,
4         style_dim):
5         super().__init__()
6         self.layers = nn.ModuleList([nn.Linear(...)
7             for _ in range(depth)])
8         self.style_affines = nn.ModuleList([
9             nn.Linear(style_dim, 2 * width) for _ in
10                range(depth)
11            ])
12     def forward(self, x, style):
13         for lin, aff in zip(self.layers, self.
14             style_affines):
15             gb = aff(style)
16             gamma, beta = torch.chunk(gb, 2, dim=-1)
17             x = (1.0 + gamma) * lin(x) + beta # FiLM
18             modulation
19         x = self.act(x)
20         return x

```

Listing 289. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Simplified modulation missing +1 formulation
2 class ModulatedMLP(nn.Module):
3     def forward(self, x, w):
4         s = self.style(w)
5         x = x * (s.unsqueeze(1) + 1e-6) # ERROR:
6         Should be (1 + s)
7         x = self.linear(x)
8         if self.demodulate:
9             d = torch.rsqrt((x ** 2).mean(dim=-1,
10                 keepdim=True) + 1e-8)
11             x = x * d
12         return self.act(x)

```

Listing 290. C1 Implementation: GPT-5 (Score: 0.6)

```

1 # DeepSeek R1: Correct modulation with demodulation
2 class StyleLinear(nn.Module):
3     def forward(self, x, style):
4         style = self.affine(style)
5         weight = self.weight.unsqueeze(0)
6         weight = weight * (style.unsqueeze(1) + 1) #
7         CORRECT: (s + 1)

```

```

7         if self.demodulate:
8             d = torch.rsqrt((weight ** 2).sum(dim=2,
9                 keepdim=True) + 1e-8)
10            weight = weight * d
11            return (x.unsqueeze(1) * weight).sum(dim=2) +
12                self.bias

```

Listing 291. C1 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: No style modulation at all
2 class StyleNerfField(nn.Module):
3     def __init__(self, hidden_units=256):
4         super().__init__()
5         # CRITICAL ERROR: Standard MLP without style
6         # conditioning
7         self.density_network = nn.Sequential(
8             nn.Linear(63, hidden_units),
9             nn.ReLU(),
10            nn.Linear(hidden_units, 1)
11        )
12        # MISSING: Style vector integration entirely

```

Listing 292. C1 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No model implementation at all
2 # CRITICAL ERROR: Only dataset loading code present
3 class AFHQDataset(data.Dataset):
4     def __init__(self, data_path, image_size=512):
5         # MISSING: Entire model, no style modulation
6         # implemented
7         self.dataset = []

```

Listing 293. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and DeepSeek R1 correctly implement FiLM-style modulation with proper  $(1 + \gamma)$  formulation. GPT-5 has the structure but uses incorrect modulation. Paper2Code and AutoP2C completely miss this core innovation.

##### Component C2: Fourier Positional Encoding - All Baselines

The paper uses standard NeRF positional encoding with  $L = 10$  for positions and  $L = 4$  for directions:

$$\zeta^L(x) = [\sin(2^0 x), \cos(2^0 x), \dots, \sin(2^{L-1} x), \cos(2^{L-1} x)]$$

```

1 # \nerfify\ : Correct Fourier encoding with proper
  frequencies
2 self.pos_encoding = NeRFFrequencyEncoding(
3     in_dim=3, num_frequencies=10, # L=10 for
4     positions
5     min_freq_exp=0.0, max_freq_exp=9.0, include_input
6     =True
7 )
8 self.dir_encoding = NeRFFrequencyEncoding(
9     in_dim=3, num_frequencies=4, # L=4 for
10    directions
11    min_freq_exp=0.0, max_freq_exp=3.0, include_input
12    =True
13 )

```

Listing 294. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct implementation
2 class FourierEmbedding(nn.Module):
3     def __init__(self, num_freqs=10):
4         super().__init__()
5         self.freq_bands = 2.0 * torch.linspace(0,
6             num_freqs - 1, num_freqs)

```

```

6     def forward(self, x):
7         out = [x]
8         for f in self.freq_bands:
9             out.append(torch.sin(f * x))
10            out.append(torch.cos(f * x))
11        return torch.cat(out, dim=-1) # CORRECT

```

Listing 295. C2 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Correct encoding
2 class FourierFeatureEncoding(nn.Module):
3     def __init__(self, L=10, include_input=True):
4         super().__init__()
5         self.L = L
6         self.include_input = include_input
7         self.output_dim = 2 * L * 3 * 2 + (3 if
            include_input else 0)

```

Listing 296. C2 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Correct formula, wrong usage
2 def positional_encoding(x: torch.Tensor, L: int) ->
    torch.Tensor:
3     encoding = [x]
4     for i in range(L):
5         for fn in [torch.sin, torch.cos]:
6             encoding.append(fn((2.0 * i) * x))
7     return torch.cat(encoding, dim=-1)
8     # ISSUE: Not integrated into model properly

```

Listing 297. C2 Implementation: Paper2Code (Score: 0.8)

```

1 # AutoP2C: No positional encoding
2 # CRITICAL ERROR: Missing entirely

```

Listing 298. C2 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Most baselines correctly implement Fourier encoding as it's a standard NeRF component. Only AutoP2C completely misses it.

### Component C3: Low-Resolution Feature Rendering + 2D Upsampling - All Baselines

The core efficiency innovation renders features at 32x32 resolution then upsamples:

$$I_w^{\text{Approx}}(\mathbf{r}) \approx h_c \circ [\phi_w^{n_c, n_\sigma}(\mathcal{A}(\mathbf{r})), \zeta(\mathbf{d})]$$

```

1 # \nerfify\ : Complete 2D aggregation for efficient
    upsampling
2 def get_outputs(self, ray_bundle):
3     # Low-res volumetric rendering
4     field_outputs = self.field(ray_samples)
5     weights = ray_samples.get_weights(sigmas)
6     # Aggregate features in 2D (early aggregation)
7     style_feat = field_outputs["style_feat"] # [N, S
        , C]
8     feat_agg = torch.sum(weights * style_feat, dim
        =-2) # [N, C]
9     # Decode once from aggregated features
10    approx_rgb = self.field.
        decode_color_from_features(
11        features=feat_agg, directions=ray_bundle.
        directions
12    )
13    # Then upsample via synthesis network

```

Listing 299. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Partial - has low-res rendering but wrong
    upsampling
2 def forward(self, z, camera_poses, render_full=True):
3     rays_o, rays_d = self.generate_rays(camera_poses,
        res=32)
4     feature_map, _, _ = self.generator.volume_render(
        rays_o, rays_d, styles)
5     # ISSUE: Missing proper 2D aggregation pipeline
6     for block in self.synthesis_blocks:
7         x = block(x, w) # Direct upsampling without
        proper aggregation

```

Listing 300. C3 Implementation: GPT-5 (Score: 0.6)

```

1 # DeepSeek R1: Correct approach with progressive
    upsampling
2 def forward(self, z, camera_poses, render_full=True):
3     if not render_full:
4         rays_o, rays_d = self.generate_rays(
            camera_poses, res=32)
5         feature_map, depth_map, weights = self.
            generator.volume_render(
                rays_o, rays_d, styles)
6         return feature_map, depth_map
7     else:
8         # Progressive upsampling from low-res
            features
9         feature_map = feature_map.reshape(B, 32, 32,
            C)
10        for upsample_layer, block in zip(self.
            upsample_layers, self.blocks):
11            feature_map = upsample_layer(feature_map,
            styles)

```

Listing 301. C3 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: Direct high-res rendering (wrong
    approach)
2 def forward(self, rays_o, rays_d, style):
3     # CRITICAL ERROR: Renders directly at target
        resolution
4     points = sample_points_on_rays(rays_o, rays_d,
        n_samples=64)
5     densities = self.density_network(points) #
        Expensive!
6     colors = self.color_network(points)
7     rgb = volume_render(densities, colors, points)
8     # MISSING: Low-res rendering + upsampling
        strategy entirely

```

Listing 302. C3 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No rendering implementation
2 # CRITICAL ERROR: Missing entirely

```

Listing 303. C3 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY correctly implements 2D feature aggregation for efficient upsampling. DeepSeek R1 has the right approach. GPT-5 partially implements it. Paper2Code and AutoP2C miss this critical efficiency innovation.

### Component C4: Custom Upsampler with Blur Kernel - All Baselines

StyleNeRF's hybrid upsampler combines pixel-shuffle with learnable MLPs and fixed blur:

$$\text{Upsample}(X) = \text{Conv2d}(\text{PixelShuffle}(\text{Repeat}(X, 4) + \psi_\theta(X), 2), K)$$

```

1 # \nerfify\ : Complete hybrid upsampler
2 class HybridUpsampler(nn.Module):
3     def forward(self, x):
4         x_repeat = x.repeat_interleave(4, dim=1) #
5         Repeat 4x
6         psi = self.mlp(x.permute(0,2,3,1)) #
7         Learnable variations
8         x_combined = x_repeat + psi.permute(0,3,1,2)
9         x_shuffled = F.pixel_shuffle(x_combined, 2)
10        # PixelShuffle
11        x_blurred = self.blur_conv(x_shuffled) #
12        Fixed blur kernel
13        return x_blurred

```

Listing 304. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Partial - missing blur kernel
2 class Upsampler(nn.Module):
3     def forward(self, x):
4         x_rep = x.repeat(1, 4, 1, 1)
5         psi = self.mlp(x.permute(0, 2, 3, 1)).permute
6         (0, 3, 1, 2)
7         x = F.pixel_shuffle(x_rep + psi, 2)
8         # MISSING: Fixed blur convolution
9         return x # ERROR: No blur kernel applied

```

Listing 305. C4 Implementation: GPT-5 (Score: 0.4)

```

1 # DeepSeek R1: Has blur but simplified
2 class CustomUpsampler(nn.Module):
3     def __init__(self, in_ch, out_ch, style_dim):
4         super().__init__()
5         self.mlp = nn.Sequential(StyleLinear(...))
6         self.blur = nn.Conv2d(in_ch, out_ch, 3,
7                               padding=1, bias=False)
8         nn.init.dirac_(self.blur.weight) #
9         Initialize as identity
10        def forward(self, x, style):
11            x_repeated = F.pixel_shuffle(x.repeat(1, 4,
12            1, 1), 2)
13            x_mlp = self.mlp(x_repeated, style)
14            x_combined = x_repeated + F.pixel_shuffle(
15            x_mlp, 2)
16            return self.blur(x_combined) # Apply blur

```

Listing 306. C4 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: No custom upsampler
2 # CRITICAL ERROR: Missing entirely, using standard
3 upsampling

```

Listing 307. C4 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No upsampler implementation
2 # CRITICAL ERROR: Missing entirely

```

Listing 308. C4 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY fully implements the hybrid upsampler with all components. DeepSeek R1 has most elements. GPT-5 misses the blur kernel. Paper2Code and AutoP2C don't implement it.

#### Component C5: NeRF-Path Regularization Loss - All Baselines

Novel loss enforcing 3D consistency:

$$\mathcal{L}_{\text{NeRF-path}} = \frac{1}{|S|} \sum_{(i,j) \in S} (I_w^{\text{Approx}}[i,j] - I_w^{\text{NeRF}}[i,j])^2$$

```

1 # \nerfify\ : Correct NeRF-path regularization
2 def get_loss_dict(self, outputs, batch):
3     if self.config.nerf_path_loss_mult > 0.0:
4         rgb_full = outputs["rgb"] # Full volumetric
5         RGB
6         rgb_approx = outputs["approx_rgb"] #
7         Aggregated-feature RGB
8         # Subsample pixels (25% by default)
9         num = int(self.config.nerf_path_sample_frac *
10         rgb_full.shape[0])
11         idx = torch.randperm(rgb_full.shape[0])[:num]
12         nerf_path = F.mse_loss(rgb_approx[idx],
13         rgb_full[idx])
14         loss_dict["nerf_path_reg"] = 0.2 * nerf_path
15         # = 0.2

```

Listing 309. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Wrong formulation entirely
2 def nerf_path_reg(fake_high, fake_low):
3     # CRITICAL ERROR: Comparing high-res vs
4     downsampled
5     # Should be full volumetric vs aggregated
6     features!
7     return F.mse_loss(fake_high, fake_low.detach())
8 # In training:
9 reg = 0.2 * nerf_path_reg(fake_img,
10 F.interpolate(fake_img, scale_factor=0.5)) #
11 WRONG!

```

Listing 310. C5 Implementation: GPT-5 (Score: 0.2)

```

1 # DeepSeek R1: Missing NeRF-path regularization
2 # CRITICAL ERROR: Not implemented
3 class NeRFPathRegularization(nn.Module):
4     def forward(self, approx_output, nerf_output,
5     mask=None):
6     # ERROR: Defined but never used in training
7     loop

```

Listing 311. C5 Implementation: DeepSeek R1 (Score: 0.0)

```

1 # Paper2Code: No NeRF-path regularization
2 def nerf_path_regularization(predictions, targets):
3     # CRITICAL ERROR: Wrong concept - not the paper's
4     formulation
5     squared_differences = (sampled_predictions -
6     sampled_targets) ** 2
7     # MISSING: Comparison between full and approx
8     paths

```

Listing 312. C5 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No loss implementation
2 # CRITICAL ERROR: Missing entirely

```

Listing 313. C5 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY correctly implements the NeRF-path regularization comparing full volumetric vs aggregated-feature RGB. GPT-5 implements a completely wrong loss. Others don't implement it at all.

#### Component C6: View Direction Removal - All Baselines

Paper removes view direction conditioning from color prediction to improve 3D consistency.

```

1 # \nerfify\ : Correctly removes view direction
2 def __init__(self, remove_viewdir: bool = True):
3     self.remove_viewdir = remove_viewdir

```

```

4 def _encode_directions(self, directions):
5     if self.remove_viewdir:
6         # Return empty tensor - no view direction
7         return torch.zeros(*directions.shape[:-1],
8                             0),
9                             device=directions.device)
10    return self.dir_encoding(directions)

```

Listing 314. C6 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Still uses view direction
2 self.color_out = nn.Linear(hidden_dim, 3)
3 # ERROR: Color prediction still depends on view
  direction
4 # MISSING: remove_viewdir flag and implementation

```

Listing 315. C6 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: No view direction removal
2 def forward_color(self, x, style, is_background=False
3 ):
4     # ERROR: Still conditions on view direction
    implicitly
5     color = self.color_head(h) # No removal
    implemented

```

Listing 316. C6 Implementation: DeepSeek R1 (Score: 0.0)

```

1 # Paper2Code: No view direction removal
2 # MISSING: Not implemented

```

Listing 317. C6 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No model implementation
2 # MISSING: Entirely

```

Listing 318. C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY implements view direction removal. All baselines miss this 3D consistency improvement.

**Component C7: Mapping Network - All Baselines**

StyleGAN2-style mapping from  $z$  to  $w$  with 8 FC layers.

```

1 # \nerfify\ : Complete mapping network
2 self.mapping_network = MappingNetwork(
3     latent_dim=512, style_dim=512,
4     num_layers=8, lr_multiplier=0.01 # 100x lower LR
5 )

```

Listing 319. C7 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct mapping network
2 class MappingNetwork(nn.Module):
3     def __init__(self, z_dim=512, w_dim=512, n_layers
4     =8):
5         layers = []
6         for i in range(n_layers):
7             layers.append(nn.Linear(z_dim if i==0
8             else w_dim, w_dim))
9         self.mapping = nn.Sequential(*layers)

```

Listing 320. C7 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Correct implementation
2 class MappingNetwork(nn.Module):
3     def forward(self, z):
4         z = F.normalize(z, dim=-1) # Normalize input
5         for layer in self.layers:
6             z = F.leaky_relu(layer(z), 0.2)
7         return z

```

Listing 321. C7 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Has mapping but wrong integration
2 class MappingNetwork(nn.Module):
3     def __init__(self, latent_dim, style_dim,
4     num_layers):
5         # ISSUE: Correct structure but not properly
    used
6         self.mapping = nn.Sequential(*layers)

```

Listing 322. C7 Implementation: Paper2Code (Score: 0.6)

```

1 # AutoP2C: No mapping network
2 # MISSING: Entirely

```

Listing 323. C7 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Most baselines implement the mapping network correctly as it's a standard StyleGAN2 component.

### Component C8: NeRF++ Foreground/Background Split - All Baselines

NeRF++ architecture with inverted sphere parameterization for background.

```

1 # \nerfify\ : Complete NeRF++ implementation
2 # Foreground in unit sphere, background with inverted
  sphere
3 self.spatial_distortion = None # Use scene box
  normalization
4 # Handles fg/bg split properly in field
  implementation

```

Listing 324. C8 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct fg/bg split
2 # Foreground network
3 self.field = StyleNeRFField(hidden_dim=256)
4 # Background network (smaller)
5 self.bg_field = StyleNeRFField(hidden_dim=128)

```

Listing 325. C8 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Correct implementation
2 def forward_geometry(self, x, style, is_background=
3 False):
4     if is_background:
5         r = torch.sqrt((x ** 2).sum(dim=-1, keepdim=
6         True))
7         x_inv = torch.cat([x / r, 1.0 / r], dim=-1)
8         # Inverted sphere

```

Listing 326. C8 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Mentions but doesn't implement properly
2 # ISSUE: No actual fg/bg separation in rendering

```

Listing 327. C8 Implementation: Paper2Code (Score: 0.4)

```

1 # AutoP2C: No NeRF++ implementation
2 # MISSING: Entirely

```

Listing 328. C8 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY, GPT-5, and DeepSeek R1 correctly implement NeRF++. Paper2Code mentions it but doesn't implement properly.

### Component C9: Progressive Training Strategy - All Baselines

Three-stage training: 500k low-res, 5M progressive, 25M high-res.

```

1 # \nerfify\ : Has progressive structure
2 self.stage1_images = 500000
3 self.stage2_images = 5000000
4 self.stage3_images = 25000000
5 # Implemented through config but simplified

```

Listing 329. C9 Implementation: NERFIFY (Score: 0.8)

```

1 # GPT-5: No progressive training
2 # ERROR: Trains at target resolution from start
3 for iteration in range(num_iterations):
4     # MISSING: Resolution scheduling

```

Listing 330. C9 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: Has stages but incomplete
2 stages = [(32, 500000), (1024, 5000000), (1024,
3     25000000)]
4 # ISSUE: Structure exists but not fully implemented

```

Listing 331. C9 Implementation: DeepSeek R1 (Score: 0.4)

```

1 # Paper2Code: Defines stages but wrong usage
2 class ProgressiveTrainingConfig:
3     stage1_images: int = 500_000
4     # ISSUE: Not integrated with resolution changes

```

Listing 332. C9 Implementation: Paper2Code (Score: 0.6)

```

1 # AutoP2C: No training implementation
2 # MISSING: Entirely

```

Listing 333. C9 Implementation: AutoP2C (Score: 0.0)

**Analysis:** No baseline fully implements progressive training. NERFIFY has the structure. Paper2Code defines it but doesn't use it properly.

#### Component C10: Hyperparameter Fidelity - All Baselines

Matching paper's exact hyperparameters for reproducibility.

```

1 # \nerfify\ : All hyperparameters match paper
2 style_dim=512, trunk_layers=4, trunk_hidden_dim=256,
3 color_layers=2, color_hidden_dim=128,
4 nerf_path_loss_mult=0.2, nerf_path_sample_frac=0.25

```

Listing 334. C10 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Some parameters correct
2 style_dim=512, hidden_dim=256
3 # ISSUE: Missing several key parameters

```

Listing 335. C10 Implementation: GPT-5 (Score: 0.6)

```

1 # DeepSeek R1: Most parameters correct
2 style_dim=512, hidden_dim=256, bg_hidden_dim=128
3 # ISSUE: Missing nerf_path parameters

```

Listing 336. C10 Implementation: DeepSeek R1 (Score: 0.7)

```

1 # Paper2Code: Partial match
2 latent_dim: int = 512
3 # ISSUE: Many incorrect or missing

```

Listing 337. C10 Implementation: Paper2Code (Score: 0.5)

```

1 # AutoP2C: No hyperparameters
2 # MISSING: Entirely

```

Listing 338. C10 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY matches all hyperparameters exactly. Others have partial matches.

#### 4.7.6. Scoring Analysis

Table 58. StyleNeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.8	1.0	0.98
GPT-5	0.6	1.0	0.6	0.4	0.2	0.0	1.0	1.0	0.0	0.6	0.52
DeepSeek R1	1.0	1.0	0.8	0.8	0.0	0.0	1.0	1.0	0.4	0.7	0.62
Paper2Code	0.0	0.8	0.0	0.0	0.0	0.0	0.6	0.4	0.6	0.5	0.28
AutoP2C	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.00

#### 4.7.7. Why Baselines Fail Despite Component Scores

##### GPT-5 (Score: 0.52 components, 0% trainable)

- **Strengths:** Correct Fourier encoding, mapping network, NeRF++ split

- **Fatal Issues:**

- Wrong NeRF-path regularization (compares high-res vs downsampled instead of volumetric vs aggregated)
- No Nerfstudio integration - standalone implementation incompatible with benchmark
- Missing view direction removal and blur kernel in up-sampler

- **Result:** Cannot train on benchmark datasets, incorrect 3D consistency enforcement

##### DeepSeek R1 (Score: 0.62 components, 0% trainable)

- **Strengths:** Perfect style modulation with demodulation, correct NeRF++ architecture

- **Fatal Issues:**

- Missing NeRF-path regularization entirely - loses key 3D consistency benefit
- No Nerfstudio plugin structure - cannot integrate with training pipeline
- Missing view direction removal

- **Result:** Strong mathematical understanding but incomplete system, cannot execute training

##### Paper2Code (Score: 0.28 components, 0% trainable)

- **Strengths:** Configuration management, dataset loading structure

- **Fatal Issues:**

- No style modulation - uses standard MLPs without conditioning
- Missing 2D upsampling path - renders directly at high resolution (computationally infeasible)
- No NeRF-path regularization or custom upsampler

- **Result:** Fundamental misunderstanding of StyleNeRF architecture, cannot scale beyond 256×256

##### AutoP2C (Score: 0.00 components, 0% trainable)



- **Strengths:** Attempted dataset loading (with incorrect URLs)
- **Fatal Issues:**
  - No model implementation whatsoever
  - Only generates dataset and evaluation code
  - Import errors prevent execution
- **Result:** Complete failure to implement any StyleNeRF components

#### 4.7.8. Hyperparameter Fidelity

Table 59. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
Style dim	512	✓	✓	✓	✓	×
Latent dim	512	✓	✓	✓	✓	×
Mapping layers	8	✓	✓	✓	✓	×
FG hidden dim	256	✓	✓	✓	✓	×
BG hidden dim	128	✓	✓	✓	256	×
$L_{\text{pos}}$	10	✓	✓	✓	✓	×
$L_{\text{dir}}$	4	✓	✓	✓	✓	×
Initial resolution	32	✓	✓	✓	64	×
$\beta$ (NeRF-path)	0.2	✓	Wrong	×	×	×
$\lambda$ (R1 reg)	0.5	✓	×	×	0.5	×
<b>W Score</b>	–	1.00	0.55	0.64	0.46	0.00

#### 4.7.9. Conclusion

All baselines attempt to implement StyleNeRF with varying degrees of sophistication. GPT-5 and DeepSeek R1 demonstrate understanding of core architectural concepts but fail on critical details - GPT-5 implements the wrong NeRF-path regularization while DeepSeek R1 omits it entirely. Paper2Code fundamentally misunderstands the method, implementing standard NeRF without style conditioning or the efficiency innovations. AutoP2C completely fails, producing only non-functional dataset code. Despite some baselines achieving reasonable component scores (DeepSeek R1: 0.62, GPT-5: 0.52), none produce trainable code due to missing Nerfstudio integration, incorrect loss formulations, or architectural failures. Only NERFIFY produces immediately trainable code as a complete Nerfstudio plugin, demonstrating the critical importance of domain-specific synthesis with a stark 100% vs 0% trainable rate highlighting the necessity of specialized paper-to-code approaches for complex vision research.

### 4.8. iNeRF: Inverting Neural Radiance Fields for Pose Estimation

#### 4.8.1. Paper Overview

iNeRF [31] introduces a framework for 6DoF pose estimation by inverting a trained NeRF model through gradient-based optimization. The method performs mesh-free pose refinement by backpropagating photometric loss through differentiable volume rendering, optimizing camera poses on the SE(3) manifold. Key innovations include interest-region

ray sampling for computational efficiency, exponential map parameterization for valid pose updates, and self-supervised NeRF training by estimating poses for unlabeled images.

#### 4.8.2. Implementation Overview

Table 60. iNeRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	580	180	650	2100	950
File Organization	Plugin	Single	Single	Multi-file	Multi-file
SE(3) Exp Map	✓	Partial	✓	✓	×
Interest Sampling	✓	Simplified	✓	✓	×
YUV Loss	✓	×	✓	✓	×
LR Schedule	✓	×	Wrong	✓	×
Nerfstudio Integration	✓	×	×	×	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

#### 4.8.3. Novel Components

Table 61. Novel Components in iNeRF with Importance Weights

ID	Component	Weight $w_i$
C1	SE(3) exponential map parameterization	0.20
C2	Interest region sampling with dilation	0.18
C3	Interest point sampling (corner detection)	0.12
C4	Photometric RGB MSE loss	0.10
C5	YUV-UV loss variant (LineMOD)	0.08
C6	Camera LR schedule ( $\alpha_t = \alpha_0 \cdot 0.8^{t/100}$ )	0.10
C7	Small ray batch (2048 rays)	0.08
C8	Gradient-based pose updates	0.07
C9	Volume rendering integration	0.04
C10	Pose-only optimization mode	0.03

#### 4.8.4. Quantitative Metrics

Table 62. iNeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	1.00	0.97
GPT-5	0.50	0.30	0.20	0.60	0.58
DeepSeek R1	0.60	0.30	0.10	0.70	0.68
Paper2Code	0.70	0.20	0.10	0.80	0.75
AutoP2C	0.00	0.10	0.90	0.00	0.05

#### 4.8.5. Component-by-Component Analysis

##### Component C1: SE(3) Exponential Map Parameterization - All Baselines

The paper specifies pose updates via the exponential map from Lie algebra  $\mathfrak{se}(3)$  to SE(3):  $T = \exp([\xi]_{\times})T_0$  where  $\xi = [\omega, v]^T \in \mathbb{R}^6$  and the rotation matrix uses Rodrigues' formula:  $R = I + \frac{\sin \theta}{\theta} [\omega]_{\times} + \frac{1 - \cos \theta}{\theta^2} [\omega]_{\times}^2$ .

```

1 # \nerfify\ : Uses Nerfstudio's camera optimizer with
  correct SE(3)
2 camera_optimizer=CameraOptimizerConfig(mode="SO3xR3")
3 # In callbacks:
4 def set_camera_lr(step: int) -> None:
5     T_update = se3_exponential_map(self.pose_params)
6     T_current = T_update @ self.initial_pose
7 # CORRECT: Full SE(3) manifold parameterization

```

Listing 339. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: SE(3) exp map without small angle handling
2 def se3_expmap(screw):
3     omega = screw[:, :3]; v = screw[:, 3:]
4     theta = torch.norm(omega, dim=1, keepdim=True) +
5         1e-8
6     # MISSING: No small angle approximation
7     R = I + torch.sin(theta) * omega_hat / theta +
8         (1 - torch.cos(theta)) * (omega_hat @
9         omega_hat) / theta**2
10    K = I * theta + (1 - torch.cos(theta)) *
11        omega_hat
12    t = (K @ v[... , None]).squeeze(-1)
13 # ERROR: Numerical instability for small rotations

```

Listing 340. C1 Implementation: GPT-5 (Score: 0.8)

```

1 # DeepSeek R1: Complete SE(3) with small angle
  handling
2 def se3_exp_map(self, screw_axis):
3     omega = screw_axis[:3]; v = screw_axis[3:]
4     theta = torch.norm(omega)
5     if theta < 1e-8: # CORRECT: Small angle
6         approximation
7         R = torch.eye(3) + self.skew_symmetric(omega)
8         t = v
9     else:
10        # CORRECT: Full Rodrigues formula
11        omega_hat = omega / theta
12        omega_hat_skew = self.skew_symmetric(
13            omega_hat)
14        R = torch.eye(3) + torch.sin(theta) *
15            omega_hat_skew +
16            (1 - torch.cos(theta)) * omega_hat_skew @
17            omega_hat_skew
18        V = (torch.sin(theta)/theta) * torch.eye(3) +
19            ...
20        t = V @ v.unsqueeze(-1).squeeze(-1)

```

Listing 341. C1 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Textbook-correct SE(3) implementation
2 def se3_exponential_map(xi: torch.Tensor) -> torch.
  Tensor:
3     omega = xi[0:3]; nu = xi[3:6]
4     theta = torch.norm(omega)
5     eps: float = 1e-6
6     if theta.item() < eps: # CORRECT: Small angle
7         handling
8         R = I + A; t = nu
9     else:
10        sin_theta = torch.sin(theta); cos_theta =
11            torch.cos(theta)
12        R = I + (sin_theta/theta)*A + ((1-cos_theta)/
13            theta**2)*A2
14        t = (I*theta + (1-cos_theta)*A + (theta-
15            sin_theta)*A2) @ (nu/theta)
16        T = torch.eye(4); T[0:3, 0:3] = R; T[0:3, 3] = t
17 # CORRECT: Complete SE(3) with proper
18     parameterization

```

Listing 342. C1 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No SE(3) parameterization at all
2 def optimize_pose(self, image, initial_pose):
3     pose = initial_pose.clone().detach().
4         requires_grad_(True)
5     optimizer = optim.Adam([pose], lr=0.01)
6     # CRITICAL ERROR: Direct matrix optimization
7     for _ in range(100):
8         loss.backward()
9         optimizer.step()
10        with torch.no_grad():
11            pose[:] = self.project_to_se3(pose) #
12        Empty placeholder
13 def project_to_se3(self, pose):
14     return pose # WRONG: No actual projection

```

Listing 343. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY, DeepSeek R1, and Paper2Code correctly implement the full SE(3) exponential map with proper Rodrigues formula and small angle handling. GPT-5 implements the correct formulas but lacks small angle approximation, causing numerical instability. AutoP2C completely fails by optimizing the pose matrix directly without manifold constraints.

### Component C2: Interest Region Sampling with Dilation - All Baselines

Interest region sampling computes gradient magnitudes using Sobel filters, thresholds to create a binary mask, applies morphological dilation (5x5 kernel, 2 iterations), and samples rays from the dilated regions.

```

1 # \nerfify\ : Complete interest region with Sobel and
  dilation
2 def _compute_interest_mask(image_hw3, threshold,
3     dilate_iters):
4     # Sobel filters for edge detection
5     kx = torch.tensor([[-1, 0, 1], [2, 0, -2], [1, 0,
6         -1]])
7     ky = torch.tensor([[-1, 2, 1], [0, 0, 0], [-1, -2,
8         -1]])
9     gx = F.conv2d(gray_, kx, padding=1)
10    gy = F.conv2d(gray_, ky, padding=1)
11    grad = torch.sqrt(gx.pow(2) + gy.pow(2))
12    mask = (grad_norm >= threshold).float()
13    # CORRECT: Morphological dilation via max pooling
14    for _ in range(dilate_iters):
15        mask = F.max_pool2d(mask, kernel_size=5,
16            stride=1, padding=2)

```

Listing 344. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Simplified gradient without Sobel or
  dilation
2 def interest_region_sampling(image, num_rays=2048):
3     # WRONG: Simple difference instead of Sobel
4     grad_x = torch.abs(image[:, :, 1:] - image[:, :,
5         :1])
6     grad_y = torch.abs(image[:, 1:, :] - image[:, :1,
7         :])
8     grad = F.pad(grad_x[... , :-1] + grad_y[... , :-1],
9         (0, 1, 0, 1))
10    # MISSING: No morphological dilation
11    prob = grad.flatten() / grad.sum()
12    idx = torch.multinomial(prob, num_rays,
13        replacement=False)
14 # ERROR: Missing core algorithm components

```

Listing 345. C2 Implementation: GPT-5 (Score: 0.4)

```

1 # DeepSeek R1: Complete with OpenCV Sobel and
  dilation
2 def interest_region_sampling(self, pose, image,
  num_rays):
3     gray = cv2.cvtColor(image, cv2.COLOR_RGB2GRAY)
4     grad_x = cv2.Sobel(gray, cv2.CV_64F, 1, 0, ksize
      =3)
5     grad_y = cv2.Sobel(gray, cv2.CV_64F, 0, 1, ksize
      =3)
6     grad_mag = np.sqrt(grad_x**2 + grad_y**2)
7     mask = (grad_mag > 0.2).astype(np.uint8)
8     # CORRECT: 5x5 kernel, 2 iterations
9     kernel = np.ones((5, 5), np.uint8)
10    mask = cv2.dilate(mask, kernel, iterations=2)

```

Listing 346. C2 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Comprehensive interest region
  implementation
2 def interest_region_sampling(self, observed_img,
  num_rays):
3     sobelx = cv2.Sobel(gray, cv2.CV_64F, 1, 0, ksize
      =3)
4     sobely = cv2.Sobel(gray, cv2.CV_64F, 0, 1, ksize
      =3)
5     gradient_magnitude = np.sqrt(sobelx**2 + sobely
      **2)
6     binary_mask = (gradient_magnitude >= 0.2).astype(
      np.uint8)
7     # CORRECT: Proper dilation parameters
8     kernel = cv2.getStructuringElement(cv2.MORPH_RECT
      , (5, 5))
9     dilated_mask = cv2.dilate(binary_mask, kernel,
      iterations=2)

```

Listing 347. C2 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No interest sampling at all
2 def generate_rays(self):
3     for pose in self.camera_poses:
4         directions = self.compute_ray_directions(self
      .image_size,
5                                                     self
      .intrinsics['focal_length'])
6         rays_world = self.transform_rays(directions,
      pose)
7 # CRITICAL ERROR: Generates rays for ALL pixels
8 # MISSING: No gradient computation, no sampling
  strategy

```

Listing 348. C2 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY , DeepSeek R1, and Paper2Code correctly implement complete interest region sampling with Sobel filters and morphological dilation. GPT-5 uses simplified gradients without proper edge detection or dilation. AutoP2C completely lacks any interest sampling mechanism.

### Component C3: Interest Point Sampling - All Baselines

Interest point sampling detects corners/features and samples rays from high-gradient locations without dilation.

```

1 # \nerfify\ : Interest point via top-k gradient
  magnitudes
2 elif strategy == "interest_point":
3     # Sample from top-k gradient magnitudes (no
      dilation)
4     mag = torch.sqrt(gx.pow(2) + gy.pow(2))
5     flat = mag.view(-1)
6     k = min(N * 10, flat.numel())
7     topk = torch.topk(flat, k=k, largest=True)
8     idxs = topk.indices

```

```

9     ys, xs = idxs // W, idxs % W
10 # CORRECT: Distinct from region sampling

```

Listing 349. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: No separate interest point implementation
2 # Uses same simplified gradient approach for all
  sampling
3 def interest_region_sampling(image, num_rays=2048):
4     grad = F.pad(grad_x[... , :-1] + grad_y[... , :-1])
5     prob = grad.flatten() / grad.sum()
6 # ISSUE: No distinction between point and region
  sampling

```

Listing 350. C3 Implementation: GPT-5 (Score: 0.5)

```

1 # DeepSeek R1: Harris corner detection for interest
  points
2 def interest_point_sampling(self, pose,
  observed_image, batch_size):
3     gray = cv2.cvtColor(observed_image, cv2.
      COLOR_RGB2GRAY)
4     corners = cv2.cornerHarris(gray.astype(np.float32
      ), 2, 3, 0.04)
5     corners = cv2.dilate(corners, None) # Local
      maxima
6     threshold = 0.01 * corners.max()
7     interest_mask = corners > threshold
8     interest_coords = np.argwhere(interest_mask)
9 # CORRECT: Proper corner detection

```

Listing 351. C3 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: goodFeaturesToTrack for interest points
2 def _interest_point_sampling(self, image_np, H, W):
3     corners = cv2.goodFeaturesToTrack(gray,
4     maxCorners=self.batch_size, qualityLevel
      =0.01, minDistance=10)
5     if corners is not None:
6         keypoints = corners.reshape(-1, 2)
7 # CORRECT: Standard feature detection

```

Listing 352. C3 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No interest point sampling
2 # MISSING: No feature detection or corner detection
3 # MISSING: No sampling strategy differentiation

```

Listing 353. C3 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY , DeepSeek R1, and Paper2Code implement distinct interest point sampling methods. GPT-5 doesn't differentiate between point and region sampling. AutoP2C lacks any interest-based sampling.

### Component C4: Photometric RGB MSE Loss - All Baselines

The paper uses standard RGB mean squared error loss:  $\mathcal{L} = \sum_{r \in \mathcal{R}} \|\hat{C}(r) - C(r)\|_2^2$ .

```

1 # \nerfify\ : Standard RGB MSE via Nerfacto base
  class
2 loss_dict = super().get_loss_dict(outputs, batch,
  metrics_dict)
3 # Inherits correct RGB loss: torch.mean((pred - gt)
  **2)

```

Listing 354. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct RGB MSE
2 loss = F.mse_loss(rendered, observed)
3 # CORRECT: Standard photometric loss

```

Listing 355. C4 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: RGB MSE implementation
2 if loss_type == 'l2':
3     loss = torch.mean((rendered_rgb - observed_rgb)
4         ** 2)
5 # CORRECT: Standard MSE loss

```

Listing 356. C4 Implementation: DeepSeek R1 (Score: 1.0)

```

1 # Paper2Code: Photometric loss with proper
  normalization
2 def photometric_loss(predicted, observed):
3     difference = predicted - observed
4     squared_difference = difference ** 2
5     photometric_loss = total_error / (H * W * 3)
6 # CORRECT: Normalized MSE

```

Listing 357. C4 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: Generic loss without context
2 def photometric_loss(predicted, observed):
3     difference = predicted - observed
4     squared_difference = difference ** 2
5 # ISSUE: Oversimplified, no ray context

```

Listing 358. C4 Implementation: AutoP2C (Score: 0.5)

**Analysis:** All methods implement RGB MSE loss, but AutoP2C's implementation is overly generic without ray-based context.

#### Component C5: YUV-UV Loss Variant - All Baselines

LineMOD variant uses only U and V channels in YUV space with BT.709 conversion matrix.

```

1 # \nerfify\ : YUV-UV loss with BT.709 standard
2 if self.config.use_yuv_uv_loss:
3     yuv_gt = rgb_to_yuv_bt709(image) # BT.709
4     coefficients
5     yuv_pr = rgb_to_yuv_bt709(pred)
6     gt_uv = yuv_gt[:, :, 1:] # U and V only
7     pr_uv = yuv_pr[:, :, 1:]
8     uv_loss = torch.mean((gt_uv - pr_uv) ** 2)
9 # CORRECT: Exact paper specification

```

Listing 359. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: No YUV loss implementation
2 # MISSING: Only RGB MSE, no YUV variant

```

Listing 360. C5 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: YUV with different coefficients
2 elif loss_type == 'yuv':
3     rgb_to_yuv = torch.tensor([
4         [0.299, 0.587, 0.114], # ISSUE: BT.601, not
5         BT.709
6         [-0.14713, -0.28886, 0.436],
7         [0.615, -0.51499, -0.10001]
8     ])
9     loss = torch.mean((rendered_yuv[:, 1:] -
10         observed_yuv[:, 1:]) ** 2)
11 # ERROR: Wrong color space standard

```

Listing 361. C5 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: Correct BT.709 YUV conversion
2 conversion_matrix = np.array(
3     [[0.2126, 0.7152, 0.0722], # CORRECT: BT.709
4      [-0.09991, -0.33609, 0.436],
5      [0.615, -0.55861, -0.05639]])
6 yuv_image = torch.tensordot(image, conversion_matrix,
7     T, dims=[2], [1]))

```

Listing 362. C5 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No YUV implementation
2 # MISSING: No color space conversion

```

Listing 363. C5 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and Paper2Code correctly implement BT.709 YUV conversion. DeepSeek R1 uses wrong coefficients (BT.601). GPT-5 and AutoP2C lack YUV loss entirely.

#### Component C6: Camera LR Schedule - All Baselines

Paper specifies exponential decay:  $\alpha_t = \alpha_0 \cdot 0.8^{t/100}$ .

```

1 # \nerfify\ : Exact formula via training callback
2 def set_camera_lr(step: int) -> None:
3     lr = self._camera_lr * (0.8 ** (step / 100.0))
4     for g in optim.param_groups:
5         g["lr"] = lr
6 # CORRECT: Precise paper formula

```

Listing 364. C6 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Fixed learning rate
2 optimizer = torch.optim.Adam([pose], lr=lr)
3 # MISSING: No learning rate schedule

```

Listing 365. C6 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: Wrong exponential schedule
2 self.scheduler = torch.optim.lr_scheduler.
3     ExponentialLR(
4         self.optimizer, gamma=0.8)
5 # ERROR: gamma^t instead of gamma^(t/100)

```

Listing 366. C6 Implementation: DeepSeek R1 (Score: 0.5)

```

1 # Paper2Code: Correct custom schedule
2 def get_learning_rate(initial_lr, decay_factor,
3     decay_steps, iteration):
4     lr = initial_lr * (decay_factor ** (iteration /
5         decay_steps))
6     return lr
7 # CORRECT: Matches paper formula

```

Listing 367. C6 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No learning rate schedule
2 optimizer = optim.Adam([pose], lr=0.01)
3 # MISSING: Fixed LR throughout

```

Listing 368. C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY and Paper2Code implement the exact LR schedule. DeepSeek R1 uses wrong decay rate. GPT-5 and AutoP2C have no scheduling.

#### Component C7: Small Ray Batch Size - All Baselines

Paper uses 2048 rays per optimization step for efficiency.

```

1 # \nerfify\ : Exact batch size
2 datamanager=TemplateDataManagerConfig(
3     train_num_rays_per_batch=2048, # CORRECT: Paper
4     value
5     eval_num_rays_per_batch=4096,
6 )

```

Listing 369. C7 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct batch size
2 def interest_region_sampling(image, num_rays=2048):
3 # CORRECT: 2048 rays

```

Listing 370. C7 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Wrong default
2 def __init__(self, batch_size=512): # ERROR: Default
3     512
4     self.batch_size = batch_size

```

Listing 371. C7 Implementation: DeepSeek R1 (Score: 0.5)

```

1 # Paper2Code: Correct from config
2 self.batch_size = config["ray_sampling"].get("
3     batch_size", 2048)
4 # CORRECT: 2048 default

```

Listing 372. C7 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No batch control
2 rays_world = self.transform_rays(directions, pose)
3 # CRITICAL ERROR: Processes full image

```

Listing 373. C7 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY, GPT-5, and Paper2Code use correct batch size. DeepSeek R1 defaults to 512. AutoP2C processes entire images.

### Component C8: Gradient-based Pose Updates - All Baselines

All methods use gradient descent, but implementation quality varies.

```

1 # \nerfify\ : Proper gradient flow through volume
2 rendering
3 loss = self.get_loss_dict(outputs, batch)
4 loss.backward()
5 optimizer.step()
6 # CORRECT: Full differentiable pipeline

```

Listing 374. C8 Implementation: NERFIFY (Score: 1.0)

```

1 # Standard gradient descent present in all baselines
2 loss.backward()
3 optimizer.step()

```

Listing 375. C8 Implementation: All Others (Score: 0.8-1.0)

**Analysis:** All methods implement basic gradient descent, with varying quality in the overall pipeline.

### Component C9: Volume Rendering Integration - All Baselines

Standard NeRF volume rendering with alpha compositing.

```

1 # \nerfify\ : Uses Nerfacto's volume renderer
2 # Inherits complete volume rendering pipeline

```

Listing 376. C9 Implementation: NERFIFY (Score: 1.0)

```

1 # All baselines have some volume rendering
2 alpha = 1.0 - torch.exp(-sigma * deltas)
3 weights = alpha * transmittance
4 final_color = torch.sum(weights * rgb, dim=1)

```

Listing 377. C9 Implementation: Others (Score: 0.6-1.0)

**Analysis:** All methods include volume rendering with varying completeness.

### Component C10: Pose-only Optimization Mode - All Baselines

Freezing NeRF weights and only optimizing camera poses.

```

1 # \nerfify\ : Explicit pose-only mode
2 if self.config.invert_poses_only:
3     for p in self.parameters():
4         p.requires_grad = False
5     for p in self.camera_optimizer.parameters():
6         p.requires_grad = True
7 # CORRECT: Clear separation

```

Listing 378. C10 Implementation: NERFIFY (Score: 1.0)

```

1 # No explicit pose-only mode in other baselines
2 # MISSING: All assume fixed NeRF implicitly

```

Listing 379. C10 Implementation: Others (Score: 0.0)

**Analysis:** Only NERFIFY has explicit pose-only optimization control.

## 4.8.6. Scoring Analysis

Table 63. iNeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.7	0.97
GPT-5	0.8	0.4	0.5	1.0	0.0	0.0	1.0	1.0	0.8	0.0	0.58
DeepSeek R1	1.0	1.0	1.0	1.0	0.8	0.5	0.5	1.0	1.0	0.0	0.68
Paper2Code	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.8	0.0	0.75
AutoP2C	0.0	0.0	0.0	0.5	0.0	0.0	0.0	0.3	0.6	0.0	0.05

### 4.8.7. Why Baselines Fail Despite Component Scores

#### GPT-5 (Score: 0.58 components, 0% trainable)

- Strengths:** Implements SE(3) structure, basic volume rendering, correct batch size

- Fatal Issues:**

- No Nerfstudio integration - standalone script cannot be deployed
- Missing small angle handling causes numerical instability
- Simplified gradient computation without proper edge detection
- No learning rate schedule leads to poor convergence

- Result:** Code runs but produces unstable optimization with divergent poses



#### DeepSeek R1 (Score: 0.68 components, 0% trainable)

- **Strengths:** Complete SE(3) implementation, proper interest sampling, YUV loss variant
- **Fatal Issues:**
  - No framework integration - requires complete rewrite for Nerfstudio
  - Wrong LR schedule causes 100× faster decay than intended
  - Default batch size mismatch reduces gradient quality
- **Result:** Mathematically sound but practically unusable without major refactoring

#### Paper2Code (Score: 0.75 components, 0% trainable)

- **Strengths:** Most complete component implementation, correct hyperparameters
- **Fatal Issues:**
  - 2100+ lines of framework-agnostic code incompatible with Nerfstudio
  - Over-engineered module structure prevents simple integration
  - Requires specific config.yaml format not matching Nerfstudio conventions
- **Result:** High-quality standalone system that cannot be deployed as intended plugin

#### AutoP2C (Score: 0.05 components, 0% trainable)

- **Strengths:** None - fundamental algorithmic failures throughout
- **Fatal Issues:**
  - No SE(3) manifold constraints - optimizes invalid rotation matrices
  - Missing all sampling strategies - processes entire images
  - Import errors prevent compilation (non-existent colormap\_wrapper)
  - Empty placeholder functions throughout codebase
- **Result:** Code neither compiles nor demonstrates any domain understanding

#### 4.8.8. Hyperparameter Fidelity

Table 64. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
Ray batch size	2048	✓	✓	×	✓	×
Initial LR ( $\alpha_0$ )	0.01	✓	×	✓	✓	✓
LR decay factor	0.8	✓	×	×	✓	×
LR decay steps	100	✓	×	×	✓	×
Optimizer	Adam	✓	✓	✓	✓	✓
Adam $\beta_1$	0.9	✓	×	✓	×	×
Adam $\beta_2$	0.999	✓	×	✓	×	×
Edge threshold	0.2	✓	×	✓	✓	×
Dilation iterations	2	✓	×	✓	✓	×
Kernel size	5×5	✓	×	✓	✓	×
<b>W Score</b>	–	1.00	0.20	0.60	0.80	0.10

#### 4.8.9. Conclusion

All baselines attempt to implement iNeRF with varying sophistication. Paper2Code achieves the highest component

coverage (75%) with correct SE(3) parameterization and complete interest sampling but produces framework-agnostic code. DeepSeek R1 implements core algorithms correctly (68%) but uses wrong hyperparameters that prevent convergence. GPT-5 provides a functional skeleton (58%) missing critical details like small angle handling and proper edge detection. AutoP2C completely fails (5%) with no SE(3) constraints, missing sampling strategies, and import errors. Only NERFIFY produces immediately trainable code as a complete Nerfstudio plugin with 97% component fidelity and exact hyperparameter matching, demonstrating that the 0% trainable rate for baselines versus 100% for NERFIFY validates the necessity of domain-specific synthesis for complex vision research.

### 4.9. SIGNeRF: Scene Integrated Generation for Neural Radiance Fields

#### 4.9.1. Paper Overview

SIGNeRF introduces a framework for editing neural radiance fields by generating multi-view consistent reference sheets through depth-conditioned ControlNet inpainting. The method enables scene-integrated generation by first creating a reference sheet with multiple views arranged in a grid, processing them through ControlNet with depth conditioning, and then propagating the edits to all training views. This addresses the critical challenge of maintaining 3D consistency when editing NeRF scenes with diffusion models.

#### 4.9.2. Implementation Overview

Table 65. SIGNeRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	482	876	624	1842	658
File Organization	Plugin	Single	Single	Multi-file	Multi-file
Reference Sheet	✓	✓	Partial	✓	×
ControlNet Integration	✓	✓	Wrong Type	✓	×
Dataset Updating	✓	Partial	Simplified	✓	×
Two-Stage Training	✓	✓	Simplified	Partial	×
Selection Modes	✓	Partial	Partial	Partial	×
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

#### 4.9.3. Novel Components

#### 4.9.4. Quantitative Metrics

#### 4.9.5. Component-by-Component Analysis

##### Component C1: Reference Sheet Generation - All Baselines

The paper specifies creating a reference sheet by arranging M views in a grid with one empty slot:

$$\mathcal{R} \leftarrow \text{ControlNet}(\bar{\mathcal{I}}_R, \bar{\mathcal{D}}_R, \bar{\mathcal{M}}_R, y) \quad (1)$$

Table 66. Novel Components in SIGNeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Reference sheet generation with grid layout	0.20
C2	ControlNet depth-conditioned inpainting	0.20
C3	Multi-view dataset updating	0.15
C4	Two-stage NeRF training (initial + fine-tune)	0.10
C5	Selection modes (bounding box and proxy)	0.10
C6	Iterative refinement with parameter adjustment	0.10
C7	ControlNet hyperparameters (scale, guidance)	0.08
C8	LPIPS loss for consistency	0.07

Table 67. SIGNeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	1.00	1.00
GPT-5	0.50	0.25	0.25	0.63	0.58
DeepSeek R1	0.63	0.25	0.12	0.75	0.72
Paper2Code	0.38	0.38	0.24	0.50	0.52
AutoP2C	0.00	0.13	0.87	0.00	0.08

```

1 # \nerfify\ : Complete grid generation with
  Nerfstudio integration
2 class ControlNetEditor:
3     def create_reference_sheet(self, rgb_images,
4                               depth_images,
5                               masks, prompt,
6                               grid_size=(2, 3)):
7         """Create edited reference sheet with M views
8         in grid."""
9         # Arrange images into grid with empty slot
10        reference_sheet = self._arrange_grid(
11            rgb_images, depth_images, masks,
12            grid_size
13        )
14        # Process with ControlNet
15        edited_sheet = self.
16        _run_controlnet_inpainting(
17            reference_sheet['rgb_grid'],
18            reference_sheet['depth_grid'],
19            reference_sheet['mask_grid'], prompt
20        )
21        return edited_sheet
22 # CORRECT: Proper grid layout with empty slot

```

Listing 380. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Standalone implementation
2 def generate_reference_sheet(color_imgs, depth_imgs,
3                               mask_imgs, prompt, ncols
4                               =3):
5     """Combine M views into tiled reference sheet."""
6     grid_color = make_grid(colors, ncols)
7     grid_depth = make_grid(depths, ncols)
8     grid_mask = make_grid([cv2.merge([m,m,m])
9                             for m in masks], ncols)
10    gen = pipe(prompt=prompt, image=grid_color,
11              control_image=grid_depth, mask_image=
12              grid_mask)
13    return gen.images[0]
14 # ERROR: Not integrated with Nerfstudio framework

```

Listing 381. C1 Implementation: GPT-5 (Score: 0.8)

```

1 # DeepSeek R1: Simplified grid generation
2 def generate_reference_sheet(self, prompt,
3                               selection_mode,
4                               num_reference_views=5):
5     reference_cameras = self.
6     _select_reference_cameras(
7         num_reference_views, selection_params)
8     color_grid, depth_grid, mask_grid = [], [], []
9     for cam in reference_cameras:
10        color, depth = self.nerf_model.render(cam)
11        mask = self._generate_mask(cam,
12                                   selection_mode,
13                                   selection_params,
14                                   depth)
15        color_grid.append(color)
16    grid_size = int(np.ceil(np.sqrt(
17        num_reference_views)))
18    # ISSUE: Missing actual ControlNet processing
19    reference_sheet = self._create_reference_sheet(
20        color_grid, depth_grid, mask_grid, grid_size,
21        prompt)

```

Listing 382. C1 Implementation: DeepSeek R1 (Score: 0.7)

```

1 # Paper2Code: Complete but standalone
2 def generate_reference_sheet(self, image_grid,
3                               depth_grid,
4                               mask_grid, prompt=None):
5     """Generate reference sheet using diffusion
6     pipeline."""
7     image_pil = self._np_to_pil(image_grid)
8     depth_pil = self._np_to_pil(depth_grid)
9     mask_pil = self._np_to_pil(mask_grid)
10    generator = torch.Generator(device=self.device).
11    manual_seed(
12        self.seed)
13    result = self.pipeline(prompt=prompt, image=
14        image_pil,
15        mask_image=mask_pil,
16        controlnet_conditioning_image=depth_pil)
17    return self._pil_to_np(result.images[0])
18 # ERROR: Not a Nerfstudio plugin

```

Listing 383. C1 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No reference sheet generation
2 class BlenderDataset:
3     def load_data(self):
4         """Loads images and camera poses."""
5         with open('transforms.json', 'r') as f:
6             metadata = json.load(f)
7             self.camera_poses = [np.array(frame['
8                 transform_matrix'])
9                 for frame in metadata[
10                 'frames']]
11 # MISSING: No reference sheet or ControlNet
12 # implementation

```

Listing 384. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and Paper2Code correctly implement grid-based reference sheet generation, but only NERFIFY integrates it with Nerfstudio. GPT-5 has correct logic but lacks framework integration. DeepSeek R1 has the structure but misses actual ControlNet processing. AutoP2C completely lacks this component.

### Component C2: ControlNet Depth-Conditioned Inpainting - All Baselines

The paper specifies using SDXL with depth ControlNet for multi-view consistency.

```

1 # \nerfify\ : Correct depth ControlNet setup
2 from diffusers import (
3     StableDiffusionXLControlNetInpaintPipeline,
4     ControlNetModel)
5 self.controlnet_editor = ControlNetEditor(
6     controlnet_scale=self.config.controlnet_scale,
7     # 0.7
8     guidance_scale=self.config.guidance_scale,
9     # 7.5
10    denoising_strength=self.config.denoising_strength
11    # 0.8
12 )
13 # CORRECT: Uses depth ControlNet as specified

```

Listing 385. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct depth ControlNet
2 controlnet = ControlNetModel.from_pretrained(
3     "diffusers/controlnet-depth-sdxl-1.0",
4     torch_dtype=torch.float16)
5 self.pipe =
6     StableDiffusionXLControlNetInpaintPipeline\
7     .from_pretrained("stabilityai/stable-diffusion-xl-
8     -base-1.0",
9     controlnet=controlnet,
10    torch_dtype=torch.float16).to(
11    device)
12 # CORRECT: Proper depth conditioning

```

Listing 386. C2 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Wrong ControlNet type
2 if self.config['use_sdxl']:
3     controlnet = ControlNetModel.from_pretrained(
4         "diffusers/controlnet-canny-sdxl-1.0", #
5         WRONG!
6         torch_dtype=torch.float16)
7 # CRITICAL ERROR: Uses Canny instead of Depth

```

Listing 387. C2 Implementation: DeepSeek R1 (Score: 0.5)

```

1 # Paper2Code: Correct but hardcoded parameters
2 self.pipeline =
3     StableDiffusionControlNetInpaintPipeline\
4     .from_pretrained("stabilityai/stable-diffusion-xl-
5     -base-1.0",
6     torch_dtype=torch.float16,
7     revision="fp16")
8 self.controlnet_scale = float(
9     diffusion_config.get("controlnet_scale", [0.4,
10    1.0])[0])
11 # Uses 0.4 instead of paper's 0.7-0.8
12 # ISSUE: Wrong hyperparameters

```

Listing 388. C2 Implementation: Paper2Code (Score: 0.9)

```

1 # AutoP2C: No ControlNet implementation
2 class ControlNet(nn.Module):
3     def __init__(self, latent_dim=512, control_scale
4     =0.7):
5         super(ControlNet, self).__init__()
6         self.encoder = self.create_encoder()
7         # Just a placeholder NN module
8 # MISSING: No actual diffusion model integration

```

Listing 389. C2 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 correctly use depth ControlNet. DeepSeek R1 critically uses the wrong ControlNet type (Canny), breaking 3D consistency. Paper2Code has correct implementation but wrong hyperparameters. AutoP2C has no real ControlNet.

## Component C3: Multi-View Dataset Updating - All Baselines

The paper describes propagating edits to all views by compositing each into the reference sheet.

```

1 # \nerfify\ : Complete view propagation
2 def propagate_to_views(self, reference_sheet,
3     view_rgbs,
4     view_depths, view_masks, prompt
5 ):
6     edited_views = []
7     for rgb, depth, mask in zip(view_rgbs,
8     view_depths, view_masks):
9         composite_sheet = self._replace_grid_slot(
10            reference_sheet, rgb, depth, mask,
11            slot_index=-1)
12         edited_view = self._run_controlnet_inpainting(
13             composite_sheet['rgb_grid'],
14             composite_sheet['depth_grid'],
15             composite_sheet['mask_grid'], prompt)
16         edited_slot = self._extract_grid_slot(
17             edited_view,
18             slot_index=-1)
19         edited_views.append(edited_slot)
20     return edited_views
21 # CORRECT: Proper view compositing and extraction

```

Listing 390. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Dataset update but saves to disk
2 def update_dataset(ref_sheet_path, color_imgs,
3     depth_imgs,
4     mask_imgs, prompt, out_dir):
5     ref_sheet = Image.open(ref_sheet_path).convert("
6     RGB")
7     for idx, (ci, di, mi) in enumerate(zip(color_imgs
8     ,
9     depth_imgs,
10    mask_imgs)):
11         grid = make_grid([np.array(ref_sheet), np.
12         array(color)],
13             ncols=2)
14         gen = pipe(prompt=prompt, image=grid,
15         control_image=depth,
16         mask_image=mask).images[0]
17         gen.save(os.path.join(out_dir, f"view_{idx:04
18         d}.png"))
19 # ERROR: Saves to disk instead of returning for
20 training

```

Listing 391. C3 Implementation: GPT-5 (Score: 0.7)

```

1 # DeepSeek R1: Incomplete update logic
2 def update_image_set(self, reference_sheet):
3     for i, cam in self.camera_poses.items():
4         color, depth = self.nerf_model.render(cam)
5         mask = self._get_camera_mask(cam)
6         # MISSING: No actual ControlNet propagation
7         edited_view = reference_sheet # WRONG!
8         self.edited_images[i] = edited_view

```

Listing 392. C3 Implementation: DeepSeek R1 (Score: 0.4)

```

1 # Paper2Code: Complete implementation
2 def update_view(self, reference_sheet, view_rgb,
3     view_depth,
4     view_mask, prompt):
5     composite_rgb = self._inject_view_into_grid(
6         self.ref_sheet, view_rgb, self.
7         empty_slot_position)
8     result = self.pipeline(prompt=prompt, image=
9     composite_rgb_pil,

```

```

7         mask_image=
        composite_mask_pil,
8
        controlnet_conditioning_image=
        composite_depth_pil)
9        updated_view = updated_composite_np[y_start:y_end
        ,
10                                           x_start:x_end,
        :]
11    return updated_view
12 # ERROR: Not integrated with Nerfstudio training

```

Listing 393. C3 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No dataset updating
2 def update_dataset(ref_sheet_path, view_list, prompt)
3 :
4     # Function defined but no implementation
5     pass
6 # MISSING: No implementation at all

```

Listing 394. C3 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and Paper2Code correctly implement view propagation with compositing. GPT-5 has the logic but saves to disk. DeepSeek R1 lacks actual ControlNet propagation. AutoP2C has no implementation.

#### Component C4: Two-Stage NeRF Training - All Baselines

The paper specifies training initial NeRF then fine-tuning on edited data.

```

1 # \nerfify\ : Integrated two-stage with Nerfstudio
2 class SIGNeRFPipeline (VanillaPipeline):
3     def get_train_loss_dict(self, step: int):
4         ray_bundle, batch = self.datamanager.
        next_train(step)
5         model_outputs = self.model(ray_bundle)
6         loss_dict = self.model.get_loss_dict(
        model_outputs, batch, metrics_dict)
7         return model_outputs, loss_dict, metrics_dict
8 # Config: initial_training: 30k iterations
9 # fine_tuning: load checkpoint, continue
10 # CORRECT: Proper two-stage pipeline

```

Listing 395. C4 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Explicit two stages but standalone
2 def train(self, data_dir, output_name="nerf", steps
    =30000):
3     config = TrainerConfig(method_name="nerfacto",
4                             max_num_iterations=steps,
5                             data=data_dir)
6     trainer = Trainer(config)
7     trainer.train()
8
9 def finetune(self, pretrained_dir, edited_data_dir,
    steps=10000):
10    config = TrainerConfig(load_dir=pretrained_dir,
11                            max_num_iterations=steps,
12                            data=edited_data_dir)
13 # ERROR: Not a Nerfstudio plugin

```

Listing 396. C4 Implementation: GPT-5 (Score: 0.9)

```

1 # DeepSeek R1: Simplified placeholder
2 def train_nerf(self, images, camera_poses):
3     self.nerf_model.train_model(images, camera_poses)
4
5 def refine_nerf(self, num_iterations=1000):
6     self.nerf_model.fine_tune(edited_images_list,

```

```

7         camera_poses_list,
        num_iterations)
8 # ISSUE: SimplifiedNeRF is dummy implementation

```

Listing 397. C4 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: Generic training only
2 class NeRFTrainer:
3     def train(self):
4         for iter_idx in range(1, self.iterations + 1)
5         :
6             batch = self.training_data[iter_idx %
        num_batches]
7             prediction = self.model(batch)
8             loss = F.mse_loss(prediction, batch)
9 # MISSING: No fine-tuning stage

```

Listing 398. C4 Implementation: Paper2Code (Score: 0.4)

```

1 # AutoP2C: Generic NeRF training
2 class Trainer:
3     def train(self):
4         for epoch in range(self.num_epochs):
5             predictions = self.model(rays['coords'],
6                                     rays['view_dirs'])
7             loss = self.compute_loss(predictions,
8                                     batch['images'])
9 # MISSING: No two-stage process at all

```

Listing 399. C4 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY properly integrates two-stage training with Nerfstudio. GPT-5 has explicit stages but standalone. DeepSeek R1 has placeholder implementation. Paper2Code and AutoP2C lack fine-tuning.

#### Component C5: Selection Modes - All Baselines

The paper specifies bounding box and proxy mesh selection modes.

```

1 # \nerfify\ : Complete selection modes
2 def generate_mask(self, depth_map, selection):
3     mode = selection.get("mode", self.default_mode)
4     if mode == "bounding_box":
5         bbox = selection["bbox"]
6         region = depth_map[bbox["y_min"]:bbox["y_max"]
7                             ,
8                             bbox["x_min"]:bbox["x_max"]
9                             ]
10    mask = (depth_map >= th_min) & (depth_map <=
        th_max)
11    elif mode == "proxy":
12        proxy_mask = selection.get("proxy_mask")
13        proxy_depth = selection.get("proxy_depth")
14        # Full implementation for both options
15 # CORRECT: Both modes fully implemented

```

Listing 400. C5 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Basic selection modes
2 def generate_masks_from_bbox(depths, bbox):
3     x1,y1,x2,y2 = bbox
4     masks=[]
5     for d in depths:
6         m=np.zeros_like(d,dtype=np.uint8)
7         m[y1:y2,x1:x2]=255
8         masks.append(m)
9 def proxy_mask_from_mesh(depths, proxy_depths):
10    masks.append(((p>0)&(p<d)).astype(np.uint8)*255)
11 # ISSUE: Simplified, not integrated

```

Listing 401. C5 Implementation: GPT-5 (Score: 0.7)

```

1 # DeepSeek R1: Both modes but simplified
2 def _shape_selection_mask(self, camera,
    selection_params,
3
4     depth_map, h, w):
5     bbox = selection_params['bounding_box']
6     # Project bounding box to image space
7 def _proxy_selection_mask(self, camera,
    selection_params,
8
9     depth_map, h, w):
10    proxy_mesh = selection_params['proxy_mesh']
11    # Render proxy mesh from camera view
12 # ISSUE: Incomplete implementations

```

Listing 402. C5 Implementation: DeepSeek R1 (Score: 0.7)

```

1 # Paper2Code: Only bounding box mode
2 if mode == "bounding_box":
3     bbox = selection["bbox"]
4     region = depth_map[y_min:y_max, x_min:x_max]
5     mask_bool = (depth_map >= th_min) & (depth_map <=
        th_max)
6 # MISSING: No proxy mode implementation

```

Listing 403. C5 Implementation: Paper2Code (Score: 0.3)

```

1 # AutoP2C: No selection modes
2 # No implementation found

```

Listing 404. C5 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY fully implements both selection modes. GPT-5 and DeepSeek R1 have both but simplified. Paper2Code only has bounding box. AutoP2C has none.

#### Component C6: Iterative Refinement - All Baselines

The paper describes iterative refinement with parameter adjustment.

```

1 # \nerfify\ : Full iterative refinement support
2 # Config allows multiple iterations:
3 trainer_config = TrainerConfig(
4     pipeline=SIGNeRFModelConfig(
5         model=SIGNeRFModelConfig(
6             use_lpips_loss=True,
7             lpips_loss_mult=0.01,
8             depth_supervision_mult=0.1
9         )
10     )
11 )
12 # CORRECT: Supports iterative refinement

```

Listing 405. C6 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Optional second iteration
2 if second_iter:
3     self.ref_gen.controlnet_scale=1.0
4     self.ref_gen.denoising_strength=0.5
5     ref_sheet2=os.path.join(base_dir, "
        reference_sheet_iter2.png")
6     self.ref_gen.generate_reference_sheet(...)
7     self.trainer.finetune(...)
8 # ISSUE: Only supports one refinement

```

Listing 406. C6 Implementation: GPT-5 (Score: 0.6)

```

1 # DeepSeek R1: Iterative refinement method
2 def iterative_refinement(self, num_iterations=2):
3     for iteration in range(num_iterations):
4         if iteration > 0:
5             self.config['controlnet_scale'] = min
6                 (1.0,

```

```

6         self.config.get('controlnet_scale',
7             0.8) + 0.2)
8         reference_sheet = self.
9         generate_reference_sheet(...)
10        self.update_image_set(reference_sheet)
11        self.refine_nerf(500)
12 # CORRECT: Good iterative structure

```

Listing 407. C6 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: No iterative refinement
2 # Single pass only, no iteration support

```

Listing 408. C6 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No iterative refinement
2 # No implementation found

```

Listing 409. C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY supports full iterative refinement through config. DeepSeek R1 has good iterative structure. GPT-5 only supports one refinement. Paper2Code and AutoP2C lack this feature.

#### Component C7: ControlNet Hyperparameters - All Baselines

The paper specifies  $\text{controlnet\_scale} 0.4 - 1.0$ ,  $\text{guidance\_scale} 6 - 10$ ,  $\text{denoising\_strength} 0.5 - 0.95$ .

```

1 # \nerfify\ : Correct hyperparameters
2 self.controlnet_editor = ControlNetEditor(
3     controlnet_scale=self.config.controlnet_scale,
4     # 0.7
5     guidance_scale=self.config.guidance_scale,
6     # 7.5
7     denoising_strength=self.config.denoising_strength
8     # 0.8
9 )
10 # CORRECT: All within paper ranges

```

Listing 410. C7 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct parameters
2 self.controlnet_scale = controlnet_scale # 0.8
3 self.guidance_scale = guidance_scale # 7.0
4 self.denoising_strength = denoising_strength # 0.8
5 # CORRECT: All within specified ranges

```

Listing 411. C7 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Mostly correct
2 controlnet_conditioning_scale=self.config.get('
    controlnet_scale', 0.8)
3 guidance_scale=self.config.get('guidance_scale', 7.5)
4 num_inference_steps=15 # Too few, paper uses 40-50
5 # ISSUE: Inference steps too low

```

Listing 412. C7 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: Wrong values from config
2 self.controlnet_scale = float(controlnet_scale_range
    [0]) # 0.4
3 self.guidance_scale = float(guidance_scale_range[0])
4 # 6.0
5 self.denoising_strength = float(denoising_range[0])
6 # 0.5
7 # ERROR: Takes first value from ranges, suboptimal

```

Listing 413. C7 Implementation: Paper2Code (Score: 0.5)



```

1 # AutoP2C: No ControlNet parameters
2 control_scale = 0.7 # Defined but never used
3 # MISSING: No actual ControlNet usage

```

Listing 414. C7 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 use correct hyperparameters. DeepSeek R1 mostly correct but inference steps too low. Paper2Code uses suboptimal values. AutoP2C has no real parameters.

#### Component C8: LPIPS Loss for Consistency - All Baselines

The paper specifies using LPIPS loss for enhanced texture consistency during fine-tuning.

```

1 # \nerfify\ : LPIPS loss integrated
2 from torchmetrics.image.lpip import
   LearnedPerceptualImagePatchSimilarity
3 if self.config.use_lpips_loss:
4     self.lpips =
       LearnedPerceptualImagePatchSimilarity(normalize=
         True)
5     lpips_loss = self.lpips(pred_rgb_permuted,
       gt_rgb_permuted)
6     loss_dict["lpips_loss"] = self.config.
       lpips_loss_mult * lpips_loss
7 # CORRECT: Proper LPIPS integration

```

Listing 415. C8 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: No LPIPS loss
2 # Uses standard MSE loss only
3 # MISSING: No LPIPS implementation

```

Listing 416. C8 Implementation: GPT-5 (Score: 0.0)

```

1 # DeepSeek R1: No LPIPS loss
2 # Standard NeRF losses only
3 # MISSING: No LPIPS implementation

```

Listing 417. C8 Implementation: DeepSeek R1 (Score: 0.0)

```

1 # Paper2Code: No LPIPS in training
2 # Only mentions LPIPS in evaluation metrics
3 # MISSING: Not used in training loss

```

Listing 418. C8 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: LPIPS defined but misused
2 def lpips_loss(predictions, targets, lpips_model):
3     lpips_value = lpips_model(predictions, targets)
4     return lpips_value.mean()
5 # ERROR: Defined but never integrated in training

```

Listing 419. C8 Implementation: AutoP2C (Score: 0.5)

**Analysis:** Only NERFIFY properly integrates LPIPS loss in training. AutoP2C defines it but doesn't use it. Others completely lack LPIPS implementation.

#### 4.9.6. Scoring Analysis

#### 4.9.7. Why Baselines Fail Despite Component Scores

##### GPT-5 (Score: 0.58 overall, 0% trainable)

- **Strengths:** Correct depth ControlNet, proper hyperparameters, explicit two-stage training

Table 68. SIGNeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	Weighted Avg <sub>LLM</sub>
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.00
GPT-5	0.8	1.0	0.7	0.9	0.7	0.6	1.0	0.0	0.71
DeepSeek R1	0.7	0.5	0.4	0.6	0.7	0.8	0.8	0.0	0.56
Paper2Code	1.0	0.9	1.0	0.4	0.3	0.0	0.5	0.0	0.51
AutoP2C	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.5	0.06

#### • Fatal Issues:

- No Nerfstudio integration - standalone script only
- Saves edited images to disk instead of training pipeline
- Missing LPIPS loss for consistency

- **Result:** Cannot be loaded as 'ns-train signerf', requires 3-5 hours manual integration

##### DeepSeek R1 (Score: 0.72 overall, 0% trainable)

- **Strengths:** Good iterative refinement structure, both selection modes

#### • Fatal Issues:

- Uses Canny ControlNet instead of Depth - fundamentally breaks 3D consistency
- SimplifiedNeRF is placeholder with no real implementation
- Missing actual ControlNet propagation in dataset update

- **Result:** Wrong ControlNet type makes method unusable for 3D-consistent editing

##### Paper2Code (Score: 0.52 overall, 0% trainable)

- **Strengths:** Complete reference sheet generation, correct depth ControlNet, proper view updating

#### • Fatal Issues:

- Not a Nerfstudio plugin - multi-file standalone code
- Missing fine-tuning stage in training
- Wrong hyperparameter values from config

- **Result:** Core algorithm correct but cannot integrate with Nerfstudio framework

##### AutoP2C (Score: 0.08 overall, 0% trainable)

- **Strengths:** Basic dataset loading structure

#### • Fatal Issues:

- No reference sheet generation or ControlNet implementation
- No dataset updating or editing pipeline
- Generic NeRF code without any SIGNeRF-specific features

- **Result:** Complete failure to understand paper, produces generic NeRF only

#### 4.9.8. Hyperparameter Fidelity

#### 4.9.9. Conclusion

All baselines attempt to implement SIGNeRF with varying levels of sophistication. GPT-5 correctly implements the core algorithm including depth ControlNet and two-stage training but lacks Nerfstudio integration. DeepSeek R1 has good structure and iterative refinement but critically uses the wrong ControlNet type (Canny instead of Depth). Pa-

Table 69. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
ControlNet scale	0.4-1.0	✓	✓	✓	×	×
Guidance scale	6-10	✓	✓	✓	✓	×
Denoising strength	0.5-0.95	✓	✓	✓	✓	×
Inference steps	40-50	✓	✓	×	✓	×
Initial iterations	30000	✓	✓	×	×	×
Fine-tune iterations	10000	✓	✓	×	×	×
LPIPS weight	0.01	✓	×	×	×	×
<b>W Score</b>	–	1.00	0.71	0.43	0.43	0.00

per2Code implements reference sheet generation and view updating correctly but misses fine-tuning and uses wrong hyperparameters. AutoP2C completely fails to understand the paper, producing only generic NeRF code without any editing capabilities. Only NERFIFY produces immediately trainable code.

#### 4.10. MCNeRF: Monte Carlo Rendering and Denoising for Real-Time NeRFs

##### 4.10.1. Paper Overview

MCNeRF [8] introduces a general-purpose acceleration technique for volumetric NeRF rendering through Monte Carlo importance sampling. The method replaces dense ray marching with stochastic sampling based on transmittance-weighted probability distributions, achieving 7× reduction in color MLP evaluations while maintaining visual quality through a lightweight image-space denoiser. This work addresses the critical bottleneck of real-time NeRF rendering by demonstrating that careful importance sampling with only 5 samples per pixel can match the quality of traditional 128-sample quadrature.

##### 4.10.2. Implementation Overview

Table 70. MCNeRF Implementation Summary Across All Baselines

Aspect	NERFIFY	GPT-5	DeepSeek R1	Paper2Code	AutoP2C
Lines of Code	476	287	892	1243	892
File Organization	Plugin	Single	Single	Multi-file	Multi-file
Monte Carlo Sampling	✓	✓	Partial	×	×
Two-Pass Rendering	✓	✓	Partial	×	×
Importance Sampling	✓	✓	✓	×	×
Denoiser Network	✓	✓	✓	Simplified	Attempted
Trainable	✓	×	×	×	×

Note: NERFIFY always produces trainable code as a complete Nerfstudio plugin

##### 4.10.3. Novel Components

##### 4.10.4. Quantitative Metrics

##### 4.10.5. Component-by-Component Analysis

###### Component C1: Monte Carlo Integration - All Baselines

The paper replaces standard quadrature  $\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i \alpha_i \mathbf{c}_i$  with Monte Carlo estimation:  $\hat{C}(\mathbf{r}) = \mathbb{E}_{i \sim p_i}[\mathbf{c}_i W]$  where  $p_i = w_i/W$  and  $W = \sum_i T_i \alpha_i$ .

Table 71. Novel Components in MCNeRF with Importance Weights

ID	Component	Weight $w_i$
C1	Monte Carlo integration replacing dense quadrature	0.25
C2	Importance sampling via transmittance CDF	0.20
C3	Stratified sampling for variance reduction	0.15
C4	Two-pass rendering (density then color)	0.15
C5	Lightweight 3-layer CNN denoiser	0.10
C6	Combined MSE + SSIM loss	0.05
C7	Per-scene denoiser training	0.05
C8	Real-time performance (5 spp target)	0.05

Table 72. MCNeRF Implementation Coverage Metrics

Method	C	I	M	W	Score <sub>LLM</sub>
NERFIFY (Ours)	1.00	0.00	0.00	1.00	0.95
GPT-5	0.75	0.25	0.00	0.85	0.95
DeepSeek R1	0.50	0.38	0.13	0.80	0.74
Paper2Code	0.00	0.13	0.88	0.20	0.15
AutoP2C	0.00	0.25	0.75	0.10	0.08

```

1 # \nerfify\ : Unbiased MC estimator with proper
   normalization
2 rgb_sel = field_outputs_sel[FieldHeadNames.RGB] # [R
   ,M,3]
3 # Unbiased estimator: mean_j( c_{i_j} ) * W
4 rgb_est = rgb_sel.mean(dim=1) * W # [R,3]
5 # Composite with background using (1 - W)
6 bg = self._background_color(R, device)
7 rgb = rgb_est + (1.0 - W_clamped) * bg # [R,3]
8 # CORRECT: Proper MC formula W * mean(c_samples)

```

Listing 420. C1 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct MC integration with weight
   normalization
2 for ep in tqdm(range(epochs)):
3     rays_o, rays_d = rays
4     preds = monte_carlo_render(rays_o, rays_d, model)
5     # Proper MC estimator in monte_carlo_render
6     rgb_out = rgb_out / n_samples
7     return rgb_out
8 # CORRECT: Proper averaging over MC samples

```

Listing 421. C1 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: MC sampling but simplified averaging
2 noisy_colors = torch.mean(colors, dim=1) # (B, 3)
3 alpha = torch.mean(torch.sigmoid(densities), dim=1,
   keepdim=True)
4 # ISSUE: Missing proper weight normalization

```

Listing 422. C1 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: No MC implementation, only metrics
2 def calculate_psnr(predictions, ground_truth):
3     mse = F.mse_loss(predictions, ground_truth)
4     psnr = 20 * torch.log10(max_pixel_value) - 10 *
   torch.log10(mse)
5 # MISSING: No Monte Carlo sampling implementation

```

Listing 423. C1 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: Dataset loading only, no rendering
2 def generate_rays(self):
3     for pose in self.camera_poses:
4         directions = np.stack([(i - cx) / fx, (j - cy) / fy])
5 # MISSING: No MC integration, only ray generation

```

Listing 424. C1 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 correctly implement the unbiased MC estimator with  $W \cdot \text{mean}(c)$ . DeepSeek R1 attempts MC but lacks proper weight normalization. Paper2Code and AutoP2C completely miss the MC algorithm.

### Component C2: Importance Sampling via Transmittance CDF - All Baselines

The paper samples according to  $p_i = w_i/W$  using inverse CDF transform on cumulative weights.

```

1 # \nerfify\ : Complete CDF-based importance sampling
2 W = weights.sum(dim=1) # [R,1]
3 # CDF over samples
4 cdf = torch.cumsum(weights[..., 0], dim=1) # [R,S]
5 # searchsorted returns indices in [0,S]
6 idx = torch.searchsorted(cdf, thresholds, right=True).clamp(max=S - 1)
7 # CORRECT: Proper inverse CDF sampling

```

Listing 425. C2 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Proper cumulative weight building
2 W = weights.sum(dim=-1, keepdim=True)
3 W[W == 0] = 1e-6
4 pdf = weights / W
5 cdfs = torch.cumsum(pdf, dim=-1)
6 idx = torch.searchsorted(cdfs, us, right=True).clamp(max=n_coarse - 1)
7 # CORRECT: Textbook inverse CDF implementation

```

Listing 426. C2 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: CDF sampling with minor issues
2 cdf = torch.cumsum(probs, dim=-1)
3 indices = torch.searchsorted(cdf, u, right=True)
4 sampled_depths = depth_below + t * (depth_above - depth_below)
5 # ISSUE: Overcomplicated interpolation logic

```

Listing 427. C2 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: No importance sampling
2 def render_view(self, camera_pose, image_size=[800, 800]):
3     rays = self.generate_rays(camera_pose, image_size)
4     rendered_image = self.model.render(rays)
5 # MISSING: No CDF-based sampling

```

Listing 428. C2 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: Basic ray generation only
2 ray_directions = directions @ pose_matrix[:3, :3].T
3 ray_origins = np.broadcast_to(pose_matrix[:3, 3], ray_directions.shape)
4 # MISSING: No importance sampling logic

```

Listing 429. C2 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 correctly implement inverse CDF sampling. DeepSeek R1 has the right idea but overcomplicates interpolation. Paper2Code and AutoP2C lack any importance sampling.

### Component C3: Stratified Sampling for Variance Reduction - All Baselines

The paper uses stratified random sampling:  $t_k = \frac{k+u_k}{M}$  where  $u_k \sim U(0, 1)$ .

```

1 # \nerfify\ : Perfect stratified sampling
2 if self.config.mc_stratified and M > 1:
3     u = torch.rand(R, M, device=device)
4     k = torch.arange(M, device=device)[None, :].expand(R, -1).float()
5     thresholds = ((k + u) / float(M)) * W # [R,M]
6 # CORRECT: Proper jittered uniform sampling

```

Listing 430. C3 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Correct stratification
2 us = torch.rand((N_rays, n_samples), device=device)
3 # Stratified MC sampling
4 cdfs = torch.cumsum(pdf, dim=-1)
5 idx = torch.searchsorted(cdfs, us, right=True).clamp(max=n_coarse - 1)
6 # CORRECT: Proper stratified sampling implementation

```

Listing 431. C3 Implementation: GPT-5 (Score: 1.0)

```

1 # DeepSeek R1: Partial stratification
2 u = torch.rand(B, num_samples, device=depths.device)
3 u = (u + torch.arange(num_samples).float()).view(1, -1) / num_samples
4 # ISSUE: Missing device specification for arange

```

Listing 432. C3 Implementation: DeepSeek R1 (Score: 0.6)

```

1 # Paper2Code: No stratified sampling
2 # Placeholder for ray generation logic
3 return []
4 # MISSING: No stratification implementation

```

Listing 433. C3 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: Uniform sampling only
2 i, j = np.meshgrid(np.arange(w), np.arange(h), indexing='xy')
3 # MISSING: No stratified sampling

```

Listing 434. C3 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 correctly implement stratified sampling with jittered uniforms. DeepSeek R1 partially implements but has device handling issues. Paper2Code and AutoP2C lack stratification.

### Component C4: Two-Pass Rendering - All Baselines

The paper specifies Pass 1 for density-only queries to build weights, Pass 2 for selective color evaluation.

```

1 # \nerfify\ : Clean two-pass separation
2 # First pass: densities, weights, accumulations
3 density, density_embedding = self.field.get_density(ray_samples)
4 weights = ray_samples.get_weights(density)
5 # Second pass: MC sample color indices
6 idx = torch.searchsorted(cdf, thresholds, right=True)
7 field_outputs_sel = self.field.get_outputs(selected_ray_samples, de_sel)

```

```
8 # CORRECT: Clear separation of passes
```

Listing 435. C4 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Explicit two-pass with comments
2 # First pass: compute opacity weights
3 for i in range(n_coarse):
4     _, sigma = model(pts, rays_d)
5     weights[:, i] = T * alpha
6 # Second pass: evaluate radiance at sampled indices
7 for j in range(n_samples):
8     rgb, _ = model(pts, rays_d)
9 # CORRECT: Well-documented two-pass structure
```

Listing 436. C4 Implementation: GPT-5 (Score: 1.0)

```
1 # DeepSeek R1: Two-pass but inefficient
2 with torch.no_grad():
3     coarse_densities = scene_representation.
4     compute_density(points_flat)
5 # Later: compute radiance
6 colors_flat = scene_representation.compute_radiance(
7     points_flat, rays_d_flat)
8 # ISSUE: Not clearly separated passes
```

Listing 437. C4 Implementation: DeepSeek R1 (Score: 0.8)

```
1 # Paper2Code: Single-pass rendering assumption
2 rendered_image = self.model.render(rays)
3 # MISSING: No two-pass implementation
```

Listing 438. C4 Implementation: Paper2Code (Score: 0.0)

```
1 # AutoP2C: No rendering implementation
2 self.rays.append((ray_origins, ray_directions))
3 # MISSING: No pass structure at all
```

Listing 439. C4 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 correctly separate density and color evaluation. DeepSeek R1 has two passes but lacks clear separation. Paper2Code and AutoP2C miss the two-pass structure entirely.

#### Component C5: Lightweight 3-Layer CNN Denoiser - All Baselines

The paper specifies a 3-layer CNN with 8 output channels, ReLU activations, and 4 input channels (RGB+ $\alpha$ ).

```
1 # \nerfify\ : Complete configurable denoiser
2 class SimpleDenoiser(nn.Module):
3     def __init__(self, in_ch: int = 4, feat: int = 8):
4         :
5         self.net = nn.Sequential(
6             nn.Conv2d(in_ch, feat, 3, padding=1),
7             nn.ReLU(inplace=True),
8             nn.Conv2d(feat, feat, 3, padding=1),
9             nn.ReLU(inplace=True),
10            nn.Conv2d(feat, 3, 3, padding=1),
11        )
12 # CORRECT: Exact paper specification
```

Listing 440. C5 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Proper 3-layer denoiser
2 class DenoiseNet(nn.Module):
3     def __init__(self):
4         self.fc1 = nn.Linear(3, hidden_dim)
5         self.fc2 = nn.Linear(hidden_dim, hidden_dim)
6         self.sigma_out = nn.Linear(hidden_dim, 1)
```

```
7 # CORRECT: 3-layer architecture with proper channels
```

Listing 441. C5 Implementation: GPT-5 (Score: 1.0)

```
1 # DeepSeek R1: Denoiser with feature computation
2 self.net = nn.Sequential(
3     nn.Conv2d(4, 8, kernel_size=3, padding=1),
4     nn.ReLU(), nn.Conv2d(8, 8, kernel_size=3, padding=
5     =1),
6     nn.ReLU(), nn.Conv2d(8, self.feature_dim + 2,
7     kernel_size=3, padding=1))
8 # ISSUE: Extra complexity with feature branches
```

Listing 442. C5 Implementation: DeepSeek R1 (Score: 0.8)

```
1 # Paper2Code: Simplified denoiser attempt
2 class DenoiseNet(nn.Module):
3     def forward(self, noisy_image, alpha_channel):
4         x = self.activation(self.conv1(input_tensor))
5         x = self.activation(self.conv2(x))
6 # ERROR: Missing proper initialization
```

Listing 443. C5 Implementation: Paper2Code (Score: 0.4)

```
1 # AutoP2C: Placeholder denoiser
2 denoised_output = F.conv2d(noisy_image, x, padding=1)
3 # CRITICAL ERROR: No actual network defined
```

Listing 444. C5 Implementation: AutoP2C (Score: 0.2)

**Analysis:** NERFIFY and GPT-5 implement the exact 3-layer specification. DeepSeek R1 adds unnecessary complexity. Paper2Code has structure but missing initialization. AutoP2C only has placeholder code.

#### Component C6: Combined MSE + SSIM Loss - All Baselines

The paper uses  $\mathcal{L} = \mathcal{L}_{\text{MSE}} + \lambda \cdot \mathcal{L}_{\text{SSIM}}$  with  $\lambda = 0.1$ .

```
1 # \nerfify\ : Proper combined loss
2 total_loss = mse_loss + 0.1 * ssim_loss # lambda =
3     0.1 as in paper
4 # CORRECT: Exact paper specification
```

Listing 445. C6 Implementation: NERFIFY (Score: 1.0)

```
1 # GPT-5: Correct SSIM implementation
2 loss_recon = F.mse_loss(denoised, gt_colors)
3 loss_ssim = ssim_loss(denoised, gt_colors)
4 loss = loss_recon + 0.1 * loss_ssim
5 # CORRECT: Proper weighting with lambda=0.1
```

Listing 446. C6 Implementation: GPT-5 (Score: 1.0)

```
1 # DeepSeek R1: SSIM but simplified
2 ssim_loss = 1 - self._ssim(denoised_images,
3     target_images)
4 total_loss = mse_loss + 0.1 * ssim_loss
5 # ISSUE: Simplified SSIM computation
```

Listing 447. C6 Implementation: DeepSeek R1 (Score: 0.6)

```
1 # Paper2Code: Has SSIM but no MC
2 ssim_val = compare_ssim(ground_truth, prediction,
3     multichannel=True)
4 # CORRECT: SSIM computation but missing integration
```

Listing 448. C6 Implementation: Paper2Code (Score: 1.0)

```

1 # AutoP2C: No loss implementation
2 # MISSING: No loss functions defined

```

Listing 449. C6 Implementation: AutoP2C (Score: 0.0)

**Analysis:** NERFIFY and GPT-5 implement the exact loss combination. DeepSeek R1 has correct weighting but simplified SSIM. Paper2Code has SSIM but not integrated. AutoP2C missing entirely.

#### Component C7: Per-Scene Denoiser Training - All Baselines

The paper trains a separate denoiser for each scene with 100k steps.

```

1 # \nerfify\ : Full per-scene pipeline
2 # Per scene, render training views with MC sampling
3 # Optimize DenoiseNet with Adam, 100k steps, batch
  size 32
4 self.denoiser_optimizer = torch.optim.Adam(
5     model.denoiser.parameters(), lr=config.get('
      denoiser_lr', 1e-3))
6 # CORRECT: Complete per-scene training

```

Listing 450. C7 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Per-scene training present
2 def train_mcnrf(model, denoiser, optimizer, rays,
  gt_colors, epochs=200):
3     for ep in tqdm(range(epochs)):
4         denoised = denoiser(inputs).view(-1, 3)
5 # ISSUE: Simplified training loop

```

Listing 451. C7 Implementation: GPT-5 (Score: 0.6)

```

1 # DeepSeek R1: Has per-scene training
2 self.denoiser_optimizer = torch.optim.Adam(
3     model.denoiser.parameters(), lr=config.get('
      denoiser_lr', 1e-3))
4 # CORRECT: Per-scene optimization setup

```

Listing 452. C7 Implementation: DeepSeek R1 (Score: 0.8)

```

1 # Paper2Code: Training loop exists
2 self.optimizer = optim.Adam(self.model.parameters(),
  lr=learning_rate)
3 for global_step in range(self.total_steps):
4     self.optimizer.step()
5 # ISSUE: Not scene-specific

```

Listing 453. C7 Implementation: Paper2Code (Score: 0.6)

```

1 # AutoP2C: Generic training attempt
2 trainer = Trainer(model=mcnerf_model, optimizer=
  optimizer)
3 trainer.train(num_epochs=config['training']['
  num_epochs'])
4 # ERROR: No actual implementation

```

Listing 454. C7 Implementation: AutoP2C (Score: 0.4)

**Analysis:** NERFIFY has complete per-scene training. GPT-5 and DeepSeek R1 have training but simplified. Paper2Code has training but not scene-specific. AutoP2C only has placeholders.

#### Component C8: Real-time Performance Target - All Baselines

The paper targets 5 samples per pixel for real-time rendering at 20 fps.

```

1 # \nerfify\ : Optimized for real-time
2 mc_spp: int = 5 # Samples per pixel
3 # WebGL shader implementing two-pass MC importance
  sampling
4 # Achieves interactive frame rates on 2021 Apple M1
  laptop
5 # CORRECT: Full real-time optimization

```

Listing 455. C8 Implementation: NERFIFY (Score: 1.0)

```

1 # GPT-5: Targets 5 spp but not optimized
2 M = 5 # use M=5 for real-time
3 # Basic implementation without GPU optimizations
4 # ISSUE: No real-time optimizations

```

Listing 456. C8 Implementation: GPT-5 (Score: 0.4)

```

1 # DeepSeek R1: Has 5 spp setting
2 self.num_importance_samples = config.get('
  num_importance_samples', 5)
3 # ISSUE: No performance optimizations

```

Listing 457. C8 Implementation: DeepSeek R1 (Score: 0.4)

```

1 # Paper2Code: No performance considerations
2 # MISSING: No real-time optimizations

```

Listing 458. C8 Implementation: Paper2Code (Score: 0.0)

```

1 # AutoP2C: No performance implementation
2 # MISSING: No optimization for speed

```

Listing 459. C8 Implementation: AutoP2C (Score: 0.0)

**Analysis:** Only NERFIFY fully optimizes for real-time performance. GPT-5 and DeepSeek R1 set 5 spp but lack optimizations. Paper2Code and AutoP2C ignore performance entirely.

#### 4.10.6. Scoring Analysis

Table 73. MCNeRF Implementation Scores

Method	C1	C2	C3	C4	C5	C6	C7	C8	Weighted AvgLLM
NERFIFY	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.00
GPT-5	1.0	1.0	1.0	1.0	1.0	1.0	0.6	0.4	0.95
DeepSeek R1	0.8	0.8	0.6	0.8	0.8	0.6	0.8	0.4	0.74
Paper2Code	0.0	0.0	0.0	0.0	0.4	1.0	0.6	0.0	0.15
AutoP2C	0.0	0.0	0.0	0.0	0.2	0.0	0.4	0.0	0.08

#### 4.10.7. Why Baselines Fail Despite Component Scores

##### GPT-5 (Score: 0.95 components, 0% trainable)

- **Strengths:** Perfect Monte Carlo mathematics, correct stratified sampling, proper SSIM implementation
- **Fatal Issues:**
  - Standalone implementation with no framework integration
  - Missing training loop orchestration
  - No data loading pipeline
- **Result:** Algorithmically correct but requires manual integration to train

**DeepSeek R1 (Score: 0.74 components, 0% trainable)**



- **Strengths:** Multiple backend support (TensorRF, Instant-NGP), comprehensive architecture
- **Fatal Issues:**
  - No framework integration or training pipeline
  - Device handling errors in stratified sampling
  - Overcomplicated interpolation logic
- **Result:** Good coverage but execution fails due to missing infrastructure

**Paper2Code (Score: 0.15 components, 0% trainable)**

- **Strengths:** Has evaluation metrics and SSIM implementation
- **Fatal Issues:**
  - Completely missing Monte Carlo sampling (core algorithm)
  - No importance sampling or two-pass rendering
  - Placeholder functions throughout
- **Result:** Fundamental misunderstanding of paper’s core contribution

**AutoP2C (Score: 0.08 components, 0% trainable)**

- **Strengths:** Basic dataset loading structure
- **Fatal Issues:**
  - No Monte Carlo implementation whatsoever
  - Missing all core algorithms
  - Placeholder denoiser with no actual network
- **Result:** Only implements data loading, misses entire algorithm

only peripheral components like metrics and data loading. The stark contrast between NERFIFY’s 100% trainable implementation and the 0% trainable rate across all baselines underscores the critical importance of domain-specific synthesis with proper framework integration for reproducible research code generation.

#### 4.10.8. Hyperparameter Fidelity

Table 74. Hyperparameter Implementation Accuracy

Parameter	Paper	NERFIFY	GPT-5	R1	P2C	AutoP2C
Samples per pixel	5	✓	✓	✓	×	×
Coarse samples	128	✓	✓	✓	×	×
Denoiser layers	3	✓	✓	✓	~	×
Denoiser channels	8	✓	✓	✓	~	×
SSIM weight $\lambda$	0.1	✓	✓	✓	✓	×
Stratified sampling	Yes	✓	✓	~	×	×
Learning rate	$10^{-3}$	✓	✓	✓	✓	✓
Batch size	32	✓	~	~	✓	✓
Training steps	100k	✓	~	~	✓	~
<b>W Score</b>	–	1.00	0.85	0.80	0.20	0.10

#### 4.10.9. Conclusion

All baselines attempt to implement MCNeRF with varying degrees of sophistication. NERFIFY achieves perfect implementation (1.00 correct) as a complete Nerfstudio plugin with all components properly integrated. GPT-5 demonstrates exceptional algorithmic understanding (0.95 score) with correct Monte Carlo mathematics and stratified sampling, but lacks framework integration. DeepSeek R1 provides good coverage (0.74) with multiple backend support but has critical execution issues. Paper2Code (0.15) and AutoP2C (0.08) fundamentally fail by completely missing the core Monte Carlo sampling algorithm, implementing

## 5. Implementation Details

### 5.1. Context-Free Grammar

```
1
2 Context-Free Grammar
3 // Notation:
4 // - Terminal symbols are in double quotes.
5 // - Non-terminals are in angle brackets.
6 // - \epsilon denotes the empty string.
7 // - [] = optional, {} = zero or more repetitions.
8 // - <PY_EXPR>, <PY_STMT>, <MARKDOWN>, <TOML> are lexical non-terminals
9 //   handled by a separate Python / Markdown / TOML grammar or left unconstrained.
10
11 // =====
12 // Top-level \nerfify\ repository & multi-file protocol for the LLM
13 // =====
14
15 <NeRFifyRepo> ::= <ConfigFile> <PipelineFile> <ModelFile> <DataManagerFile>
16                  [<FieldFile>]
17                  {<AuxFile>}
18                  <InitFile>
19                  <PyProjectFile>
20                  [<ReadmeFile>]
21
22 // Each file is wrapped in explicit markers that the LLM must output.
23 // In the paper we show them as terminals; in code they are literal strings.
24
25 <ConfigFile>      ::= "###FILE: method_template/template_config.py\n"
26                      <ConfigModule>
27                      "\n###END_FILE\n"
28
29 <PipelineFile>    ::= "###FILE: method_template/template_pipeline.py\n"
30                      <PipelineModule>
31                      "\n###END_FILE\n"
32
33 <ModelFile>       ::= "###FILE: method_template/template_model.py\n"
34                      <ModelModule>
35                      "\n###END_FILE\n"
36
37 <DataManagerFile> ::= "###FILE: method_template/template_datamanager.py\n"
38                      <DataManagerModule>
39                      "\n###END_FILE\n"
40
41 <FieldFile>       ::= "###FILE: method_template/template_field.py\n"
42                      <FieldModule>
43                      "\n###END_FILE\n"
44
45 // Optional auxiliary modules that the LLM is allowed to generate.
46 // They share the generic <AuxModuleBody> grammar (Python module).
47
48 <AuxFile> ::= "###FILE: method_template/" <AuxFileName> ".py\n"
49             <AuxModuleBody>
50             "\n###END_FILE\n"
51
52 <AuxFileName> ::= "utils" | "losses" | "encoders" | "render_utils"
53                 | "metrics" | "diffusion" | "samplers" | "scripts"
54
55 // Package initialiser and entry point in pyproject.toml
56
57 <InitFile> ::= "###FILE: method_template/__init__.py\n"
58               [<ModuleDocstring>]
59               "from .template_config import " <MethodVar> " \n"
60               "__all__ = [\" " <MethodVar> "\"]\n"
61               "###END_FILE\n"
62
63 <PyProjectFile> ::= "###FILE: pyproject.toml\n"
64                   "[project.entry-points.'nerfstudio.method_configs']\n"
65                   <MethodName> " = \"method_template.template_config:\" <MethodVar> "\"\n"
66                   {<TomlLine>}
67                   "###END_FILE\n"
68
69
70 <ReadmeFile> ::= "###FILE: README.md\n"
71                 <MARKDOWN>
72                 "###END_FILE\n"
73
74
75 <TomlLine> ::= <TOML> // any additional TOML lines
```

```

76
77 // Method identity used consistently across files
78 <MethodName> ::= <LowerIdent> // snake_case method name
79 <MethodVar> ::= <LowerIdent> // Python variable (e.g. method_template)
80
81 // =====
82 // Configuration module: MethodSpecification + TrainerConfig wiring
83 // =====
84
85 <ConfigModule> ::= [<ModuleDocstring>]
86                   <ConfigImports>
87                   <MethodSpecDef>
88
89 <ModuleDocstring> ::= '""' {<Char>} '""' "\n"
90
91 <ConfigImports> ::= <RequiredConfigImports> {<ExtraImport>}
92
93 <RequiredConfigImports> ::=
94     "from __future__ import annotations\n"
95     "from method_template.template_datamanager import TemplateDataManagerConfig\n"
96     "from method_template.template_model import TemplateModelConfig\n"
97     "from method_template.template_pipeline import TemplatePipelineConfig\n"
98     "from nerfstudio.configs.base_config import ViewerConfig\n"
99     "from nerfstudio.data.dataparsers.nerfstudio_dataparser import NerfstudioDataParserConfig\n"
100    "from nerfstudio.engine.optimizers import AdamOptimizerConfig\n"
101    "from nerfstudio.engine.schedulers import ExponentialDecaySchedulerConfig\n"
102    "from nerfstudio.engine.trainer import TrainerConfig\n"
103    "from nerfstudio.plugins.types import MethodSpecification\n"
104
105 <ExtraImport> ::=
106     "from" <ModulePath> "import" <IdentList> "\n"
107     | "import" <ModulePath> "\n"
108
109 <ModulePath> ::= <LowerIdent> { "." <LowerIdent> }
110 <IdentList> ::= <Ident> { "," <Ident> }
111 <Ident> ::= <LowerIdent> | <UpperIdent>
112
113 <MethodSpecDef> ::= <MethodVar> " = MethodSpecification("
114                   <TrainerConfig> ","
115                   <DescriptionArg>
116                   ")\n"
117
118 <TrainerConfig> ::= "config=TrainerConfig(" <TrainerKwargs> ")"
119
120 <TrainerKwargs> ::= <TrainerKwarg> { "," <TrainerKwarg> }
121
122 <TrainerKwarg> ::= "method_name=" <String>
123                 | "pipeline=" <PipelineConfig>
124                 | "datamanager=" <DataManagerConfig>
125                 | "model=" <ModelConfig>
126                 | "viewer=" <ViewerConfigExpr>
127                 | "optimizers=" <OptimizerDict>
128                 | "vis=" <String>
129                 | <Ident> "=" <PY_EXPR> // catch-all extra kwargs
130
131 <PipelineConfig> ::= "TemplatePipelineConfig(" <PyArgs> ")"
132 <DataManagerConfig> ::= "TemplateDataManagerConfig(" <PyArgs> ")"
133 <ModelConfig> ::= "TemplateModelConfig(" <PyArgs> ")"
134
135 <your entire grammar goes here>
136
137 <ViewerConfigExpr> ::= "ViewerConfig(" <PyArgs> ")"
138
139 <OptimizerDict> ::= "{" <OptimizerEntry> { "," <OptimizerEntry> } "}"
140
141 <OptimizerEntry> ::= <String> ": {"
142                   "'optimizer': ' <OptimizerConfig> ','
143                   "'scheduler': ' <SchedulerConfig>
144                   "}"
145
146 <OptimizerConfig> ::= <Ident> "(" <PyArgs> ")"
147 <SchedulerConfig> ::= <Ident> "(" <PyArgs> ")"
148
149 <PyArgs> ::= [<PY_EXPR> { "," <PY_EXPR> }]
150
151 <DescriptionArg> ::= "description=" <String>
152
153 <String> ::= '""' {<Char>} '""'
154

```

```

155 // =====
156 // Pipeline module: VanillaPipelineConfig + VanillaPipeline subclass
157 // =====
158
159 <PipelineModule> ::= [<ModuleDocstring>]
160                     <PipelineImports>
161                     <PipelineConfigClass>
162                     <PipelineClass>
163
164 <PipelineImports> ::= {<ImportStmt>}
165
166 <ImportStmt> ::=
167     "import" <ModulePath> "\n"
168 | "from" <ModulePath> "import" <IdentList> "\n"
169
170 <PipelineConfigClass> ::=
171     "@dataclass\n"
172     "class " <PipelineConfigName> "(VanillaPipelineConfig):" <IndentedSuite>
173
174 <PipelineClass> ::=
175     "class " <PipelineName> "(VanillaPipeline):" <PipelineSuite>
176
177 <PipelineSuite> ::= <InitMethod> {<PipelineMethod>}
178
179 <InitMethod> ::=
180     "def __init__(self, " <PipelineParams> "):" <IndentedSuite>
181     // Inside the suite we require core wiring patterns:
182     // - super().__init__(...)
183     // - datamanager setup
184     // - model setup
185     // - optional DDP setup
186
187 <PipelineParams> ::= <Param> {"", " <Param>}
188 <Param> ::= <Ident> [": " <Type>] [" = " <PY_EXPR>]
189
190 <Type> ::= "int" | "float" | "bool" | "Literal" | "Optional" | "Type" | "DataManager" | "Model" | <Ident>
191
192 <PipelineMethod> ::= <DefHeader> <IndentedSuite>
193
194 <DefHeader> ::= "def " <Ident> "(" [<Param> {"", " <Param>}] ")" ":"
195
196 <IndentedSuite> ::= "\n" <Indent> {<PY_STMT>} <Dedent>
197
198 // We do not expand <PY_STMT> further; it's a lexical non-terminal for valid Python code.
199
200 // =====
201 // Model module: NerfactoModelConfig + NerfactoModel subclass
202 // =====
203
204 <ModelModule> ::= [<ModuleDocstring>]
205                 <ModelImports>
206                 <ModelConfigClass>
207                 <ModelClass>
208
209 <ModelImports> ::= {<ImportStmt>}
210
211 <ModelConfigClass> ::=
212     "@dataclass\n"
213     "class " <ModelConfigName> "(" <BaseModelConfig> "):" <IndentedSuite>
214
215 <BaseModelConfig> ::= "NerfactoModelConfig" | "ModelConfig" | <Ident>
216
217 <ModelClass> ::=
218     "class " <ModelName> "(" <BaseModel> "):" <ModelSuite>
219
220 <BaseModel> ::= "NerfactoModel" | "Model" | <Ident>
221
222 <ModelSuite> ::= <PopulateModules>
223                 [<GetLossDict>]
224                 {<ModelMethod>}
225
226 <PopulateModules> ::=
227     "def populate_modules(self):" <IndentedSuite>
228
229 <GetLossDict> ::=
230     "def get_loss_dict(self, outputs, batch, metrics_dict=None):" <IndentedSuite>
231
232 <ModelMethod> ::= <DefHeader> <IndentedSuite>
233

```

```

234 // =====
235 // DataManager module: VanillaDataManagerConfig + VanillaDataManager subclass
236 // =====
237
238 <DataManagerModule> ::= [<ModuleDocstring>]
239                        <DataManagerImports>
240                        <DataManagerConfigClass>
241                        <DataManagerClass>
242
243 <DataManagerImports> ::= {<ImportStmt>}
244
245 <DataManagerConfigClass> ::=
246     "@dataclass\n"
247     "class " <DataManagerConfigName> "(VanillaDataManagerConfig):" <IndentedSuite>
248
249 <DataManagerClass> ::=
250     "class " <DataManagerName> "(VanillaDataManager):" <DataManagerSuite>
251
252 <DataManagerSuite> ::= <DataManagerInit>
253                        [<NextTrainMethod>]
254                        [<NextEvalMethod>]
255                        {<DataManagerMethod>}
256
257 <DataManagerInit> ::=
258     "def __init__(self, " <Param> {"", " <Param>}" "):" <IndentedSuite>
259
260 <NextTrainMethod> ::=
261     "def next_train(self, step: int) -> Tuple[RayBundle, Dict]:" <IndentedSuite>
262
263 <NextEvalMethod> ::=
264     "def next_eval(self, step: int) -> Tuple[RayBundle, Dict]:" <IndentedSuite>
265
266 <DataManagerMethod> ::= <DefHeader> <IndentedSuite>
267
268 // =====
269 // Field module (optional): NerfactoField subclass
270 // =====
271
272 <FieldModule> ::= [<ModuleDocstring>]
273                <FieldImports>
274                <FieldClass>
275
276 <FieldImports> ::= {<ImportStmt>}
277
278 <FieldClass> ::=
279     "class " <FieldName> "(" <BaseField> ")" <FieldSuite>
280
281 <BaseField> ::= "NerfactoField" | "Field" | <Ident>
282
283 <FieldSuite> ::= <FieldInit>
284                [<ForwardMethod>]
285                {<FieldMethod>}
286
287 <FieldInit> ::=
288     "def __init__(self, " <Param> {"", " <Param>}" "):" <IndentedSuite>
289
290 <ForwardMethod> ::=
291     "def forward(self, ray_samples: RaySamples) -> Dict[str, Tensor]:" <IndentedSuite>
292
293 <FieldMethod> ::= <DefHeader> <IndentedSuite>
294
295 // =====
296 // Generic auxiliary module body (for utils.py, diffusion.py, etc.)
297 // =====
298
299 <AuxModuleBody> ::= [<ModuleDocstring>]
300                  {<ImportStmt>}
301                  {<TopLevelDef>}
302
303 <TopLevelDef> ::= <DefHeader> <IndentedSuite>
304                | "@dataclass\n" "class " <UpperIdent> "(" [<Ident>] ")" <IndentedSuite>
305                | "class " <UpperIdent> "(" [<Ident>] ")" <IndentedSuite>
306
307 // =====
308 // Lexical categories (sketched)
309 // =====
310
311 <LowerIdent> ::= <LowerLetter> {<Letter> | "_" | <Digit>}
312 <UpperIdent> ::= <UpperLetter> {<Letter> | "_" | <Digit>}

```



```

313
314 <LowerLetter> ::= "a" | "b" | ... | "z"
315 <UpperLetter> ::= "A" | "B" | ... | "Z"
316 <Letter>      ::= <LowerLetter> | <UpperLetter>
317 <Digit>       ::= "0" | "1" | ... | "9"
318 <Char>        ::= any printable character except "'" and newline

```

### 5.1.1. Code Generation Templates

## 5.2. Multi-Agent Architecture and Agent Specifications

The system is built on LangChain for agent orchestration and LangGraph for workflow management. We use DeepAgent as the base agent framework and GPT-5 API for code generation. Paper parsing employs MinerU [19] to extract text, equations, and figures from PDFs. External knowledge retrieval uses Tavily for web search when agents encounter unfamiliar techniques or require implementation details absent from the paper. For visual feedback, NeRFify uses Qwen3-VL 8B model.

## 5.3. LLM Prompts Used in NeRFify

This section lists the main prompts used by the agents in our NERFIFY framework. For brevity, we show the system message and the static part of the user template for each agent. Placeholders such as {template\_tree} or {image\_paths} are filled programmatically at runtime.

### 5.3.1. Prompt 1: Markdown Cleaning Agent (agentic\_clean.py)

#### System prompt

```

1 You are a meticulous Markdown cleaner for research papers.
2 Apply the following strictly. Output ONLY valid GitHub-flavored
3 Markdownno commentary, no code fences, no extra prose.
4
5 Rules:
6 1) Text hygiene: de-hyphenate wrapped words, fix OCR ligatures,
7    normalize quotes/dashes/units; remove duplicate lines/sections;
8    fix spacing/formatting.
9
10 2) Equations: NEVER delete or alter equations. Preserve inline
11     $...$  and display  $...$  math exactly as authored, including
12    numbering/labels if present.
13
14 3) Scope pruning: remove generic narrative sections
15    (Introduction/background, Related Work surveys,
16    qualitative/marketing-style Results). Keep ONLY content needed
17    to implement and reproduce experiments: problem setup,
18    assumptions, notation, model architecture, objectives/losses,
19    algorithms/pseudocode, training schedule, datasets,
20    preprocessing, hyperparameters, ablations that affect
21    implementation, and evaluation protocol/metrics definitions.
22
23 4) Tables: delete benchmark/comparison tables. If a single clear
24    takeaway is obvious, replace with a one-line textual takeaway.
25    Retain implementation-critical tables (hyperparameters, layer
26    configs, dataset splits) but convert them into concise bullet
27    lists; no raw HTML or Markdown tables.
28
29 5) Figures: remove images/captions/links/placeholders. If nearby
30    text contains implementation-relevant details, keep that text
31    as plain prose. Do NOT invent content.
32
33 6) Strip raw HTML artifacts entirely (<td>, <tr>, <table>,
34    inline styles). Keep only pure Markdown and math.
35
36 7) Citations: remove unnecessary/dangling citations and citation

```

```

37     dumps. Keep citations only when required to identify
38     datasets/codebases/definitions essential for reproduction.
39
40 8) Summaries: when paragraphs are verbose or narrative, compress
41     to 3 6 bullets emphasizing actionable implementation details
42     and experimental setup (no fixed word threshold).
43
44 9) Final check: output must parse as clean, minimal Markdown with
45     intact math, no images, no HTML table tags, no benchmark
46     tables, and no broken anchors.
47
48 10) Brevity: make the output as short as possible while
49     preserving all information required for implementation and
50     experiments.

```

### User prompt template

```

1 Clean the following Markdown. Output ONLY the cleaned Markdown.
2
3 <MARKDOWN CONTENT HERE>

```

### 5.3.2. Prompt 2: Citation Search Agent (agentic\_citation\_recovery.py)

#### System prompt (SEARCH\_AGENT\_SYSTEM\_PROMPT)

```

1 You are a NeRF and 3D vision research assistant focused on
2 citation discovery and dependency analysis.
3
4 For each user query:
5 - Use 'internet_search' to find the canonical page(s) for the
6   requested NeRF-related paper (arXiv, project page, or main
7   publication page).
8 - Skim titles, abstracts, and key metadata from the results.
9 - Produce a concise textual bundle that will later be parsed by
10  another LLM.
11
12 Your answer MUST:
13 - Clearly state the best-guess canonical title, arXiv ID (if any),
14   year, venue (if obvious), and 3 6 bullet points summarizing the
15   method.
16 - List 3 8 likely upstream dependency / baseline papers, each with
17   a one-line reason describing what is reused or extended.
18 - Include the strongest URLs (arXiv / project page / PDF) inline.
19
20 Keep the answer compact but information-dense. Do not write
21 generic advice; only report what you found in the search results.

```

### User prompt template for a single query

```

1 Search for the NeRF-related paper and its canonical page or arXiv:
2 {query}
3
4 Use the 'internet_search' tool as needed. Return a concise
5 textual bundle with:
6 - Canonical title and year
7 - arXiv ID (if any)
8 - Venue (if obvious)
9 - 3 6 bullet points summarizing the method

```

```
10 - 3 8 likely upstream dependency papers with one-line reasons each
11 - Best URLs (arXiv / project / PDF)
```

### 5.3.3. Prompt 3: Dependency Extraction Prompt

#### System prompt (DEPENDENCY\_EXTRACTION\_SYSTEM)

```
1 You are a research assistant building a citation dependency graph
2 for NeRF papers.
3
4 Given text that may include scraped search results, references,
5 abstracts, or markdown, identify:
6 (a) the paper's canonical title,
7 (b) arXiv id if any,
8 (c) year,
9 (d) core components it introduces or uses as architectural modules,
10 loss functions, training protocols, or datasets,
11 (e) explicit dependencies on prior NeRF-style methods.
12
13 Return a single JSON object with keys:
14 - title: string
15 - arxiv_id: string or null
16 - year: integer or null
17 - venue: string or null
18 - url: string or null
19 - summary: short string
20 - components: {modules:[], losses:[], protocols:[]}
21 - dependencies: [{title_or_id, reason, components_borrowed}]
22 - citations: [{title_or_id, url?}]
23
24 Return STRICT JSON only, with no extra commentary.
```

#### User prompt template

```
1 Use the following bundle to extract paper facts and dependencies.
2
3 <BUNDLE AS JSON OR TEXT HERE>
4
5 Return strict JSON only.
```

### 5.3.4. Prompt 4: Target Analysis Prompt

#### System prompt (TARGET\_ANALYSIS\_SYSTEM)

```
1 You extract novelties, formulas, and implementation details from
2 NeRF papers and their dependencies.
3
4 Given a JSON citation/dependency graph plus metadata for a target
5 paper, produce a compact analysis describing:
6 - core novelties of the target method (architecture, losses,
7 training/inference protocols),
8 - key equations (with informal names, LaTeX-like forms, and short
9 descriptions),
10 - a high-level but implementation-ready plan of components.
11
12 Return STRICT JSON with:
13 - novelties: [{aspect, description}]
14 - formulas: [{name, equation, description}]
```

```

15 - implementation_plan: {
16     modules:[...],
17     losses:[...],
18     protocols:[...],
19     data_requirements:[...],
20     pseudo_steps:[...]
21 }
22
23 No extra commentary outside this JSON.

```

## User prompt template

```

1 Analyze the following target paper text to extract novelties,
2 formulas, and an implementation plan.
3
4 <TARGET PAPER TEXT OR DERIVED MARKDOWN HERE>
5
6 Return strict JSON only.

```

### 5.3.5. Prompt 5: DAG + File Plan with Repo Snapshot (agentic\_dag.py)

#### User prompt (DAG + file-plan generator)

```

1 You are a senior Nerfstudio engineer.
2
3 TASK A      IMPORT DAG:
4 - Build a file-level DAG where each node is a relative file path
5   (e.g., "method_template/template_model.py").
6 - Add a directed edge for each internal Python import:
7   { "from": "<file>", "to": "<file>", "relation": "imports" }.
8 - The result must be a valid DAG (no cycles). If cycles are
9   detected, break them by removing the least number of edges.
10 - Use the REPO SNAPSHOT and OBSERVED IMPORTS as ground truth.
11   If there is disagreement, prefer OBSERVED IMPORTS.
12
13 TASK B      FILE PLAN (optional for codegen):
14 - Propose the minimal set of files under "method_template/"
15   (subset of TEMPLATE TREE) that must be implemented for this
16   paper.
17 - For each file, provide:
18   path,
19   purpose,
20   depends_on (relative file paths),
21   key_classes,
22   key_functions.
23
24 Return STRICT JSON of the form:
25 {
26   "nodes": [
27     {"id": "<file>", "label": "<short description>"}
28   ],
29   "edges": [
30     {"from": "<src_file>", "to": "<dst_file>", "relation": "imports"}
31   ],
32   "files": [
33     {
34       "path": "method_template/template_model.py",
35       "purpose": "<short description>",

```

```

36     "depends_on": [
37         "method_template/template_field.py",
38         ...
39     ],
40     "key_classes": ["METHOD_Model"],
41     "key_functions": ["forward"]
42 },
43 ...
44 ]
45 }
46
47 Constraints:
48 - Only include paths that exist in the TEMPLATE TREE or REPO
49   SNAPSHOT.
50 - Keep depends_on lists short and meaningful.
51 - Do NOT invent new directories or top-level packages beyond what
52   you see.
53 - JSON must be syntactically valid and parseable.
54
55 REPO SNAPSHOT (selected files with excerpts):
56 <snapshot_text>
57
58 OBSERVED INTERNAL IMPORTS (static analysis):
59 <imports_text>

```

### System prompt used for this call

```

1 You are a NeRF planning agent. Return JSON only.

```

### 5.3.6. Prompt 6: DAG + File Plan from Template Tree Only (build\_dag.py)

#### User prompt

```

1 You are a senior Nerfstudio engineer. Using ONLY the provided
2 TEMPLATE TREE and TEMPLATE FILE CONTENTS, design a complete DAG
3 of the method and a file-generation plan.
4 - Do NOT invent files outside TEMPLATE TREE.
5 - Do NOT include any in-context examples.
6 - Be concise but complete.
7
8 Return strict JSON with:
9 - "nodes": [
10     { "id": "snake_case_id",
11       "label": "Short title",
12       "methods": ["method1", "method2"] }
13   ]
14 - "edges": [
15     { "from": "node_id",
16       "to": "node_id",
17       "relation": "feeds|queries|supervises|produces|writes" }
18   ]
19 - "files": [
20     {
21       "path": "relative/path.py or README.md",
22       "purpose": "short description",
23       "depends_on": ["other/file.py", ...],
24       "key_classes": ["ClassA"],
25       "key_functions": ["fn_a", "fn_b"]

```



```

26     }
27 ]
28
29 TEMPLATE TREE (authoritative list of allowed files):
30 {template_tree}
31
32 TEMPLATE FILE CONTENTS (use to match APIs/imports):
33 {template_files}

```

## System prompt

```

1 You are a NeRF planning agent. Return JSON only.

```

### 5.3.7. Prompt 7: NeRFify Code Generation Prompt (agentic\_coding.py)

#### Base full-code prompt (\_build\_full\_prompt)

```

1 You are a senior NeRFStudio engineer. Synthesize a complete
2 implementation for the TARGET PAPER by learning patterns from
3 in-context examples and by adhering EXACTLY to the provided
4 NeRFStudio method template.
5
6 ### WHAT YOU GET
7 A) TEMPLATE TREE (authoritative list of output files & their
8     relative paths under method_template):
9 {template_tree}
10
11 B) TEMPLATE FILE CONTENTS (current stubs you must overwrite; use
12     for APIs/imports naming):
13 {template_files}
14
15 C) IN-CONTEXT EXAMPLES (order matters; learn structure, naming,
16     losses, schedulers, data flows):
17 {examples}
18
19 D) TARGET PAPER - FINAL to IMPLEMENT
20 {final_paper}
21 {optional_file_plan_block}
22
23 ### IMPLEMENTATION REQUIREMENTS
24 1) OUTPUT FILES:
25     Generate code ONLY for the files that exist in TEMPLATE TREE
26     under method_template. Do NOT invent new files or paths
27     beyond this list.
28
29 2) METHOD NAME:
30     Infer a concise snake_case METHOD_NAME from the paper title
31     (e.g., "mip_nerf", "seathru_nerf"). Use this consistently
32     across files (config, model, field, datamanager, pipeline).
33
34 3) COMPATIBILITY & STYLE:
35     - Match Nerfstudio APIs and signatures shown in TEMPLATE
36       FILE CONTENTS.
37     - Use type hints and relative imports.
38     - Use defaults from the TARGET PAPER; encode losses exactly.
39     - Put data requirements/transforms in template_datamanager.py.
40     - Wire metrics in template_pipeline.py; compute them in model/
41       field as needed.

```

```

42 - Avoid unresolved imports or circular dependencies.
43
44 4) SELF-CONSISTENCY:
45 - The project should install with 'pip install -e .'.
46 - Imports like
47     `from method_template.template_model import METHOD_NAME_Model `
48     must succeed.
49 - All files should form a coherent, runnable Nerfstudio method.
50
51 5) OUTPUT FORMAT:
52 - Return your answer as a sequence of file blocks:
53     @@@FILE: relative/path.py
54     <file contents>
55     @@@END_FILE
56 - One block per file. No extra commentary, no Markdown fences,
57     no explanations outside the file blocks.

```

### Incremental / consolidation variants

The incremental and consolidation prompts (`_build_incremental_prompt` and `_build_consolidation_prompt`) reuse the same structure as above, but:

- Restrict the allowed output paths to a subset `{allowed_files}` (plus optional `{extra_allowed_files}`).
- Add an extra section E) CURRENT IMPLEMENTATION (YOU MAY MODIFY AS NEEDED) containing the current generated files:

```

1 E) CURRENT IMPLEMENTATION (YOU MAY MODIFY AS NEEDED)
2 The following files were generated incrementally. You may change
3 ANY of them to ensure correctness and consistency:
4 {current_text}

```

### 5.3.8. Prompt 8: ReAct Wrapper for Code Generation (agentic\_coding.py)

#### System prompt for ReAct agent

```

1 You are a senior NerFStudio engineer. You may use tools
2 (web_search, http_get) to verify facts. Your FINAL answer must
3 be ONLY the multi-file output in the exact format specified by
4 the user using @@@FILE blocks and @@@END_FILE, with NO extra
5 commentary, NO code fences, NO prefixes, and NO explanations.

```

The human message for this agent is simply:

```

1 {input}

```

where `{input}` is the full-code prompt from Prompt 7 (or its incremental variant).

### 5.3.9. Prompt 9: NeRF Artifact Detection Prompt (nerf\_qwen\_vl\_artifacts.py)

#### Image-level Qwen3-VL prompt (build\_prompt)

```

1 You are an expert in Neural Radiance Field (NeRF) rendering
2 diagnostics.
3 The image you see is a NeRF render named '{image_name}' with
4 resolution {width}x{height} pixels.
5
6 Your tasks:
7 1. Inspect the image for visible NeRF-specific rendering defects.
8 2. For each defect, provide a tight bounding box and classify it.
9 3. For each defect, propose concrete code-level changes to a
10    typical NeRF implementation that would likely reduce or

```

```

11     remove this artifact.
12
13 Consider the following defect types and use their keys exactly as
14 written:
15 {defect_list_text}
16
17 Examples of code-level suggestions you may use (adapt as
18 appropriate):
19 - Increase samples per ray (coarse and/or fine); increase
20   proposal network steps.
21 - Enable mip-NeRF-style cone tracing / anti-aliasing to reduce
22   aliasing.
23 - Increase distortion loss weight or density regularization to
24   remove floaters.
25 - Adjust near/far bounds or background handling to fix background
26   leakage.
27 - Add depth/normal smoothness loss to fix tearing or missing
28   geometry.
29 - Improve pose estimation (re-run COLMAP with higher quality) if
30   multi-view misalignment is visible.
31 - Increase training iterations or reduce learning rate if
32   underfit / noisy.
33 - Use higher input image resolution or improve dataloader
34   resizing.
35
36 IMPORTANT:
37 - Set width to {width} and height to {height} in the JSON.
38 - Use as few defects as needed but up to 10 per image.
39 - Be specific in code_suggestions: reference knobs like samples
40   per ray, proposal networks, loss weights, occupancy grids,
41   background color handling, pose optimization, etc.
42
43 Return ONLY a single valid JSON object, with this structure:
44
45 {
46   "image_name": "<exact image file name>",
47   "width": <integer width in pixels>,
48   "height": <integer height in pixels>,
49   "defects": [
50     {
51       "type": "<one of the defect type keys above>",
52       "bbox": [x_min, y_min, x_max, y_max],
53       "description": "<short natural-language description>",
54       "code_suggestions": [
55         "<one concrete change to the NeRF code or hyperparameters>",
56         "<another concrete suggestion if useful>"
57     ]
58   ]
59 }
60
61
62 Constraints:
63 - bbox coordinates must be integers in pixel space, with:
64   0 <= x_min < x_max <= width,
65   0 <= y_min < y_max <= height.
66 - If no defects are present, return "defects": [] but keep the
67   rest of the JSON.
68 - Do NOT output any extra commentary, markdown, or text outside

```

### 5.3.10. Prompt 10: ReAct Aggregation Prompt for NeRF Artifacts

#### System prompt for LangChain ReAct agent

```

1 You are orchestrating an automated inspection pipeline for NeRF
2 renderings.
3 You have access to a tool called 'inspect_nerf_render' which,
4 given an image path string, returns a JSON string describing
5 defects in that image.
6
7 Your job:
8 1. For each image path provided in the user input, call
9   'inspect_nerf_render'.
10 2. Parse each JSON string returned by the tool.
11 3. Aggregate them into a single JSON object with this structure:
12
13 {
14   "images": [
15     {
16       "image_name": "<file name>",
17       "width": <integer>,
18       "height": <integer>,
19       "defects": [
20         {
21           "type": "<defect type>",
22           "bbox": [x_min, y_min, x_max, y_max],
23           "description": "<description>",
24           "code_suggestions": ["<suggestion 1>", "<suggestion 2>", "..."]
25         }
26       ]
27     }
28   ]
29 }
30
31 Rules:
32 - Call 'inspect_nerf_render' on EVERY image path and do not skip
33   any.
34 - Preserve the order of images as given.
35 - If the tool reports an error or returns no defects, still
36   include the image with "defects": [] (you may carry over an
37   error field if present).
38 - Your FINAL answer must be ONLY that JSON object, with no extra
39   text.

```

#### User prompt template

```

1 You are given the following list of NeRF render image paths:
2
3 {JSON-ENCODED LIST OF IMAGE PATHS}
4
5 For EACH path in this list, you must:
6 - Call the tool 'inspect_nerf_render' with the exact path string.
7 - Afterwards, aggregate all tool outputs into the single JSON
8   object described in the system message. Output ONLY that JSON
9   object.

```

=

**Public Release.** All resources associated with NERFIFY are publicly available at <https://seemandhar.github.io/NERFIFY/>. This includes the complete source code for the NERFIFY framework, all generated implementations for the 30 papers in NERFIFY-Bench, the benchmark dataset with evaluation scripts, processed markdown versions of all papers, input/output token usage and cost analysis for each stage of the pipeline, wall-clock time breakdowns, in-context example paper-code pairs, and the full configuration files used in our experiments. We release these resources to support reproducibility and to enable the community to build upon our work.

## References

- [1] Relja Arandjelović and Andrew Zisserman. NeRF in detail: Learning to sample for view synthesis. *arXiv preprint arXiv:2106.05264*, 2021. 1, 2, 5
- [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5855–5864, 2021. 4, 21
- [3] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensoRF: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350, 2022. 4, 29
- [4] Congyue Deng, Jiawei Yang, Leonidas Guibas, and Yue Wang. Rethinking directional integration in neural radiance fields. *arXiv preprint arXiv:2311.16504*, 2023. 1, 2, 5
- [5] Jan-Niklas Dihlmann, Andreas Engelhardt, and Hendrik Lensch. SigNeRF: Scene integrated generation for neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6679–6688, 2024. 4
- [6] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. FastNeRF: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14346–14355, 2021. 5
- [7] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. StyleNeRF: A style-based 3D-aware generator for high-resolution image synthesis. *arXiv preprint arXiv:2110.08985*, 2021. 4, 45
- [8] Kunal Gupta, Milos Hasan, Zexiang Xu, Fujun Luan, Kalyan Sunkavalli, Xin Sun, Manmohan Chandraker, and Sai Bi. MCNeRF: Monte Carlo rendering and denoising for real-time NeRFs. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–11, 2023. 4, 62
- [9] HKUDS. DeepCode: Open agentic coding. <https://github.com/HKUDS/DeepCode>, 2025. GitHub repository. 5
- [10] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2023. 5
- [11] Mijeong Kim, Seonguk Seo, and Bohyung Han. InfoNeRF: Ray entropy minimization for few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12912–12921, 2022. 1, 3, 5
- [12] Jonas Kulhanek and Torsten Sattler. Tetra-NeRF: Representing neural radiance fields using tetrahedra. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 18458–18469, 2023. 4, 35
- [13] Deborah Levy, Amit Peleg, Naama Pearl, Dan Rosenbaum, Derya Akkaynak, Simon Korman, and Tali Treibitz. SeaThru-NeRF: Neural radiance fields in scattering media. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 56–65, 2023. 3, 5
- [14] Liangchen Li and Juyong Zhang. L0-Sampler: An L0 model guided volume sampling for NeRF. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21390–21400, 2024. 1, 3, 5
- [15] Li Ma, Xiaoyu Li, Jing Liao, Qi Zhang, Xuan Wang, Jue Wang, and Pedro V Sander. Deblur-NeRF: Neural radiance fields from blurry images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12861–12870, 2022. 3, 5
- [16] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision (ECCV)*, pages 405–421, 2020. 1, 4, 5
- [17] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics*, 41 (4):1–15, 2022. 4, 5
- [18] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H Mueller, Chakravarty R Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. DOnERF: Towards real-time rendering of compact neural radiance fields using depth oracle networks. In *Computer Graphics Forum*, pages 45–59. Wiley Online Library, 2021. 5
- [19] Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, et al. MinerU2.5: A decoupled vision-language model for efficient high-resolution document parsing. *arXiv preprint arXiv:2509.22186*, 2025. 71
- [20] OpenAI. GPT-5 technical report. <https://openai.com>, 2025. 5
- [21] Marco Orsingher, Anthony Dell’Eva, Paolo Zani, Paolo Medici, and Massimo Bertozzi. Informative rays selection for few-shot neural radiance fields. *arXiv preprint arXiv:2312.17561*, 2023. 1, 2, 5
- [22] Leandro A Passos, Douglas Rodrigues, Danilo Jodas, Kelton AP Costa, Ahsan Adeel, and Joao Paulo Papa. BioNeRF: Biologically plausible neural radiance fields for view synthesis. *arXiv preprint arXiv:2402.07310*, 2024. 3, 4, 5
- [23] Chen Qian, Xin Cai, Cheng Liu, Yang Liu, Juyuan Dang, Lin Wang, et al. ChatDev: Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023. 5



- [24] Shilei Sun, Ming Liu, Zhongyi Fan, Qingliang Jiao, Yuxue Liu, Liquan Dong, and Lingqin Kong. Efficient ray sampling for radiance fields reconstruction. *Computers & Graphics*, 118:48–59, 2024. [1](#), [2](#), [5](#)
- [25] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, et al. Nerfstudio: A modular framework for neural radiance field development. In *ACM SIGGRAPH Conference Proceedings*, 2023. [1](#), [4](#), [5](#)
- [26] Haithem Turki, Michael Zollhöfer, Christian Richardt, and Deva Ramanan. PyNeRF: Pyramidal neural radiance fields. *Advances in neural information processing systems*, 36:37670–37681, 2023. [4](#)
- [27] Yifan Wang, Jun Xu, Y Zeng, and Y Gong. Anisotropic neural representation learning for high-quality neural rendering. *arXiv preprint arXiv:2311.18311*, 2023. [1](#)
- [28] Yifan Wang, Yi Gong, and Yuan Zeng. Hyb-NeRF: A multiresolution hybrid encoding for neural radiance fields. In *2024 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 3677–3686. IEEE, 2024. [1](#), [2](#), [5](#)
- [29] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. NerfingMVS: Guided optimization of neural radiance fields for indoor multi-view stereo, 2021. [3](#), [5](#)
- [30] Qingshan Xu, Xuanyu Yi, Jianyao Xu, Wenbing Tao, Yew-Soon Ong, and Hanwang Zhang. Few-shot NeRF by adaptive rendering loss regularization. In *European Conference on Computer Vision*, pages 125–142. Springer, 2024. [2](#), [5](#)
- [31] Lin Yen-Chen, Pete Florence, Jonathan T Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. iNeRF: Inverting neural radiance fields for pose estimation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1323–1330. IEEE, 2021. [4](#), [51](#)
- [32] Hye Bin Yoo, Hyun Min Han, Sung Soo Hwang, and Il Yong Chun. Improving neural radiance fields using near-surface sampling with point cloud generation. *Neural Processing Letters*, 56(4):214, 2024. [1](#), [2](#), [5](#)
- [33] Yao Zhang, Jiangshu Wei, Bei Zhou, Fang Li, Yuxin Xie, and Jiajun Liu. TVNeRF: Improving few-view neural volume rendering with total variation maximization. *Knowledge-Based Systems*, 301:112273, 2024. [1](#), [2](#), [5](#)
- [34] Hanxin Zhu, Tianyu He, Xin Li, Bingchen Li, and Zhibo Chen. Is vanilla MLP in neural radiance field enough for few-shot view synthesis? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20288–20298, 2024. [1](#), [2](#), [5](#)