

# Think-Then-Generate: Structural Chain-of-Thought Reasoning for Consistent 3D Generation

## Supplementary Material

411 In this document, we provide a more detailed description  
412 of the model and implementation details to complement the  
413 main content. Additionally, we include extra experiments  
414 and prompt templates.

## 415 6. Implementation Details

### 416 6.1. Inference Pipeline

---

Algorithm 1: Thoughtful3D Generation Framework

---

**Require:** Text prompt  $P$

**Ensure:** 3D model  $\mathcal{M}$

```

1: // Before Generation: Planning Phase
2: semantics  $\leftarrow$  SemanticParser( $P$ ) {Extract objects/at-
   attributes}
3: timeline  $\leftarrow$  Planner( $P$ , semantics) {Plan generation
   stages}
4: // During Generation: Optimization Phase
5: Initialize 3D model  $\mathcal{M}$ 
6: for optimization step  $s \in$  timeline do
7:   Strategy 1: 3DRefine-CoT
8:   Render views  $\mathcal{V} = \{v_i\}$  from  $\mathcal{M}$ 
9:    $\mathcal{L}_{\text{recon}} \leftarrow 0$ 
10:  for each view  $v_i \in \mathcal{V}$  do
11:     $(\text{score}_i, \text{caption}_i) \leftarrow$  Verifier( $v_i, P$ )
12:    if  $\text{score}_i < \text{threshold}$  then
13:      Generate corrections  $\mathcal{C}_i$  using diffusion model
14:       $v'_i \leftarrow$  ConsensusSelector( $\mathcal{C}_i$ )
15:       $\mathcal{L}_{\text{recon}} \leftarrow \mathcal{L}_{\text{recon}} + \text{MSE}(v_i, v'_i)$ 
16:    end if
17:  end for
18:  Strategy 2: Cross-View Alignment
19:   $\mathcal{L}_{\text{align}} \leftarrow$  ComputeAlignmentLoss( $\mathcal{V}$ )
20:  Joint Optimization
21:   $\mathcal{L}_{\text{SDS}} \leftarrow$  ScoreDistillationLoss( $\mathcal{V}$ )
22:   $\mathcal{L}_{\text{total}} \leftarrow \lambda_1 \mathcal{L}_{\text{SDS}} + \lambda_2 \mathcal{L}_{\text{align}} + \lambda_3 \mathcal{L}_{\text{rc}}$ 
23:  Update  $\phi$  via  $\nabla_{\phi} \mathcal{L}_{\text{total}}$ 
24: end for
25: return  $\mathcal{M}$ 

```

---

### 417 6.2. Camera pose setting

418 In 3DRefine-CoT, to ensure that the rendered images  
419  $\{V_i\}_{i=1}^N$  preserve both distinguishable global structures and  
420 local details, we introduce an additional set of  $N$  camera  
421 viewpoints  $\{C_i\}_{i=1}^N$  (with  $N = 4$  in our experiments) ar-  
422 ranged in a spherical coordinate system. The rendered views

generated from these additional cameras are used in the sub-  
sequent CoT reasoning process. They serve as crucial evi-  
dence for identifying inconsistencies or other issues that  
may arise during the 3D generation process.

### 6.3. Hyperparameter Settings

In our experiment, based on the magnitude of each loss  
term, we set the coefficients  $\{\lambda_1, \lambda_2, \lambda_3\}$  defined in Eq.(13)  
as  $\lambda_1 = 1$  for  $\mathcal{L}_{\text{SDS}}$ ,  $\lambda_1 = 0.3$  for  $\mathcal{L}_{\text{rc}}$ ,  $\lambda_3 = 1$  for  $\mathcal{L}_{\text{align}}$ .  
We incorporate the  $\mathcal{L}_{\text{rc}}$  and  $\mathcal{L}_{\text{align}}$  losses only after the  
3D model has developed an initial structural form, typically  
during the mid-to-late stages defined by 3DBlueprint-CoT.

## 7. Additional Experiments

### 7.1. Comparison with other 3D Consistency Enhancement methods

We demonstrate the superior performance of our Thoughtful3D compared to other approaches aimed at improving 3D consistency. Here, we present quantitative and qualitative comparisons with Perp-Neg, Debias, Hallo3D. As shown in Fig. 8, Thoughtful3D not only effectively addresses the Janus problem, but also significantly improves the quality of generated objects. Specifically, compared to prior methods, Thoughtful3D generates pandas with multi-view consistency and without Janus problem. For storks, it produces structurally complete and realistic shapes, while other methods exhibit varying degrees of structural abnormalities in these cases. A substantial improvement is also observed in the quantitative metric, CLIP-Score.

### 7.2. A Case Study Assessing CoT Reasoning Capabilities in 3D Generation Tasks

Here, by comparing direct LLM image queries with CoT reasoning analysis, we demonstrate that CoT reasoning can effectively locate "hallucination" phenomena in renderings, such as structural errors, repetitive features, and misalignments. Using the prompt "An elegant flamingo standing tall with long legs and pinkish-white feathers." as an example, we selected one rendered view from its training process for analysis. As shown in Fig. 9, when directly queried with an LLM, the model can identify some existing issues. However, it fails to precisely detect a key problem: an abnormally distorted leg structure. In contrast, our structural CoT reasoning accurately identifies this critical issue, as demonstrated in Fig. 10.

In structural CoT reasoning, image analysis is carried out in two stages. In the first stage, the image is carefully

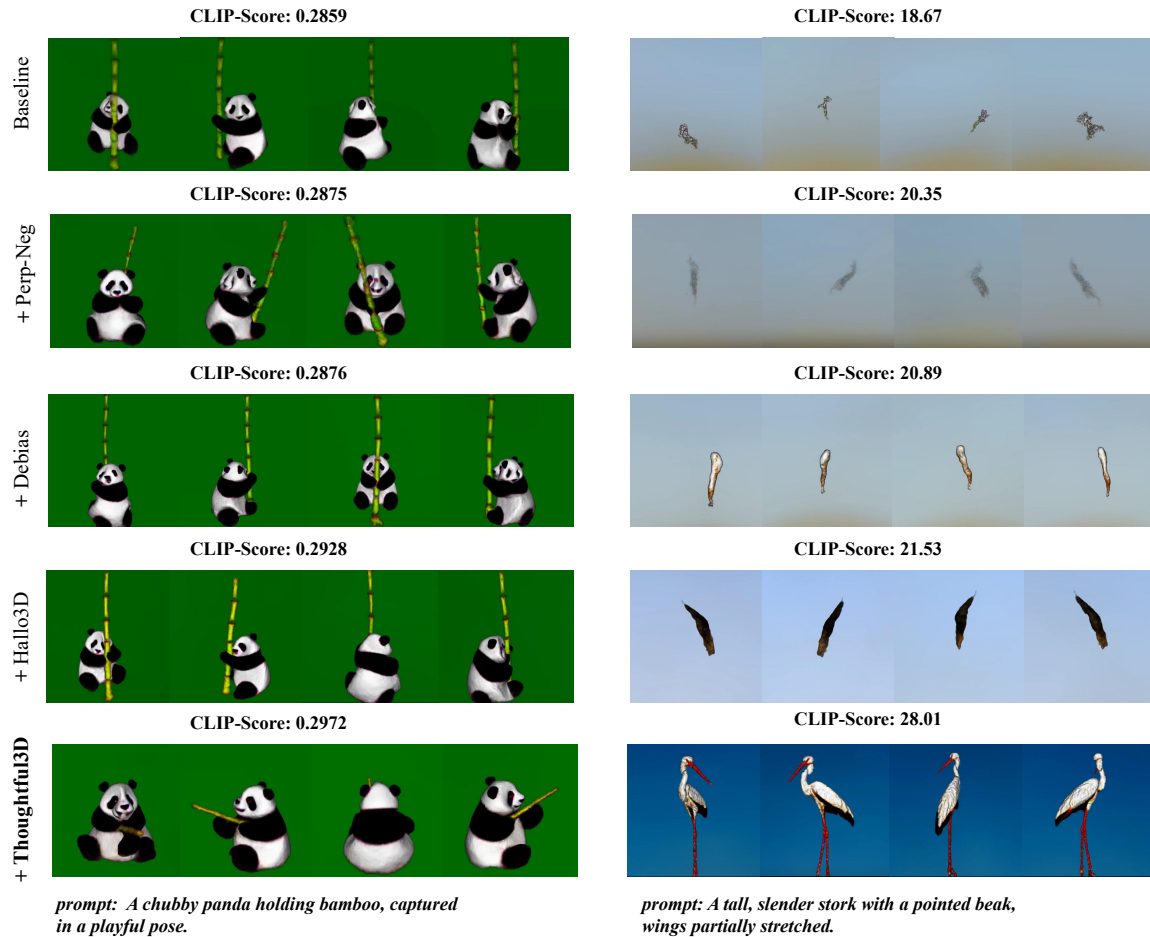


Figure 8. Comparison experiments with Perp-Neg, DeBias and Hallo3D.

467 described and scored. In the second stage, both the image  
 468 and its description are used together to analyze the prob-  
 469 lem from different perspectives. The main issue is then  
 470 identified. By adopting a "description then decision" strat-  
 471 egy, perception is separated from reasoning during analy-  
 472 sis. This improves interpretability, matches human cogni-  
 473 tive processes, and greatly enhances the ability to analyze  
 474 rendered images in 3D generation.

### 475 7.3. Additional Qualitative Results

476 Fig. 11 and Fig. 12 present qualitative results of our method  
 477 integrated with different baseline models. In each image  
 478 pair, the left side shows renderings of 3D models gener-  
 479 ated by the original approach, while the right side dis-  
 480 plays results from our method. Through planning, reflec-  
 481 tion, and correction during 3D model generation, Thought-  
 482 ful3D significantly improves multi-view coherence and en-  
 483 hances output quality. Consequently, it achieves 3D gener-  
 484 ation that faithfully aligns with diverse input prompts.

### 7.4. Robustness Analysis

485 We address MLLM stochasticity by employing Consensus  
 486 Selector as a *robust filter*. 1) *High Reliability*. A user study  
 487 (300 pairs) shows 87% feedback accuracy and only 5% hal-  
 488 lucination (Tab. 4). 2) *Error Filtering*. To test sensitivity to  
 489 errors, we injected 20% random/incorrect negative prompts.  
 490 The Consensus Selector effectively *intercepted these er-*  
 491 *rors*, raising the discard rate to 28% (Tab. 5) and prevent-  
 492 ing model degradation. 3) *Stability*. By aggregating inde-  
 493 pendent MLLM votes, our method reduces cross-run vari-  
 494 ance by 44% compared to a single-model setup (Tab. 6),  
 495 ensuring stable 3D generation despite MLLM stochasticity.  
 496

Table 4. User Study of MLLM feedback.

Neg. Prompt	Count	Percentage
Accurate	261	87%
Subjective	24	8%
Hallucination	15	5%

Table 5. MLLM Consensus Selector discard rate.

Neg. Prompt	CLIP	Discard Rate
<b>Standard (Ours)</b>	<b>30.11</b>	<b>10%</b>
Noise (Random)	29.45	21%
Adversarial (Wrong)	29.02	28%

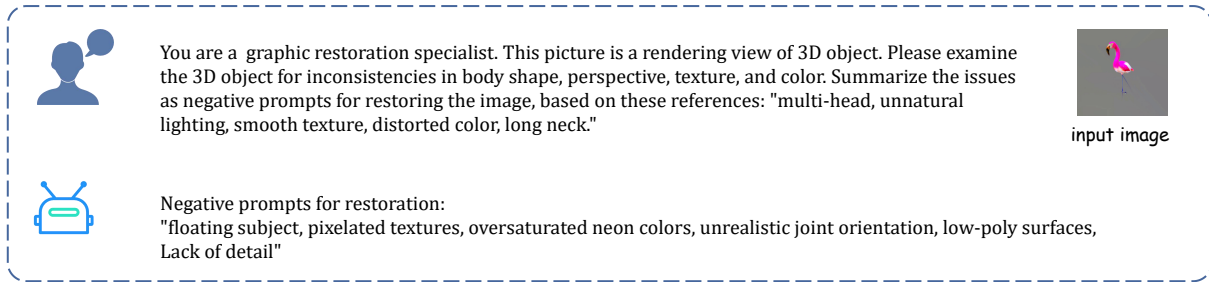


Figure 9. Analyze rendered images directly using MLLMs.

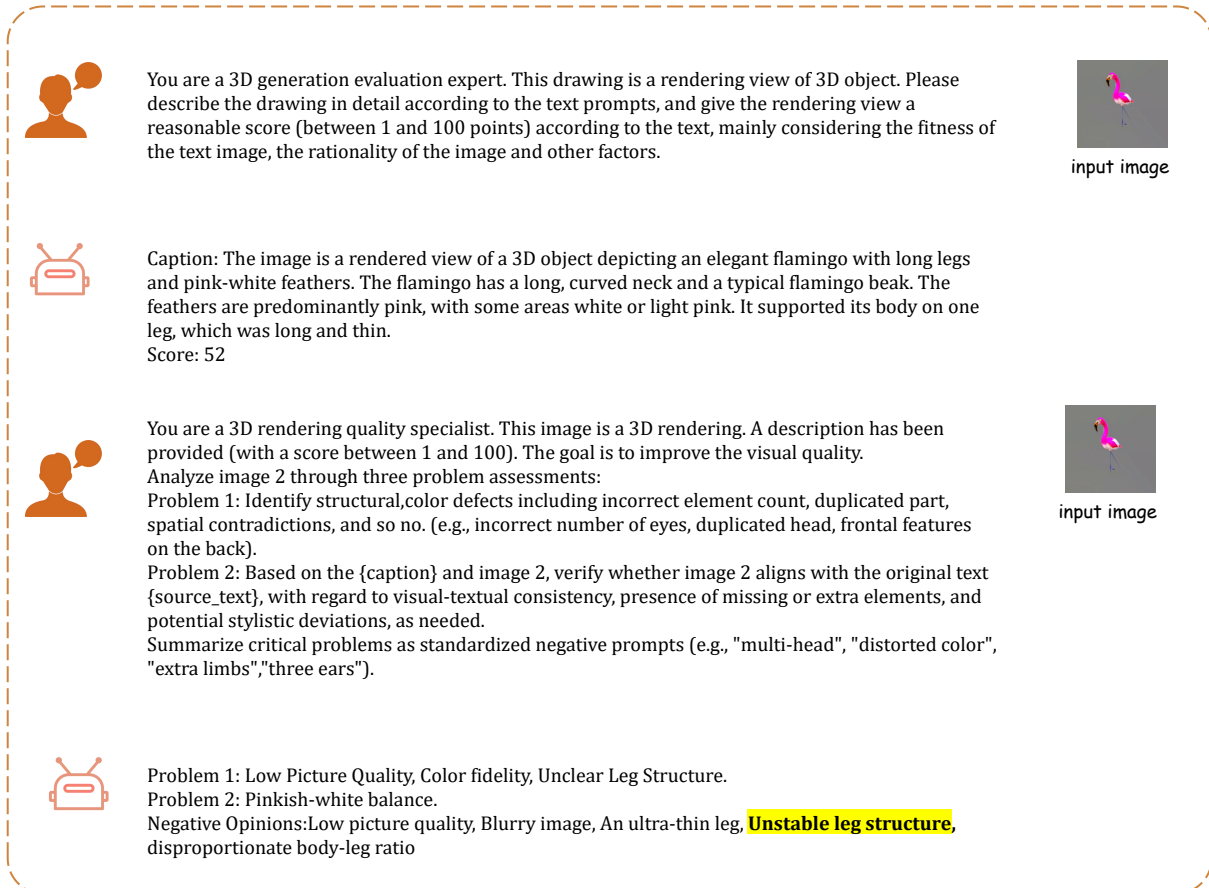


Figure 10. Analyze rendered images through structural CoT Reasoning.

Table 6. Consistency across multiple generations.

Method	CLIP $\uparrow$	Avg. Std $\downarrow$	Std@95% $\downarrow$	Run-to-Run Sim $\uparrow$
w/o consensus (Single)	0.283 $\pm$ 0.018	0.018	0.027	93.75
Ours (K=3 + Consensus)	<b>0.302 <math>\pm</math> 0.010</b>	<b>0.010</b>	<b>0.018</b>	<b>96.48</b>

head, we consider it justified by the quality gains. We ensure efficiency by applying 3DRefine-CoT exclusively in later stages to refine established geometry, striking an optimal balance.

501  
502  
503

497

## 7.5. Efficiency Comparison

498

499

500

As shown in Tab. 7, we report the runtime on an NVIDIA A100 GPU using *GaussianDreamer* and *DreamFusion* as baselines. Although our method adds computational over-

Table 7. Analysis of runtime and memory usage.

Model	CLIP (Baseline)	CLIP (Ours)	Iteration	Refine Start	Org. Time	Extra Time
GaussianDreamer	27.33	<b>30.39</b>	1200	900	~15 min	~5 min
DreamFusion	23.22	<b>28.67</b>	2500	2100	~27 min	~6 min

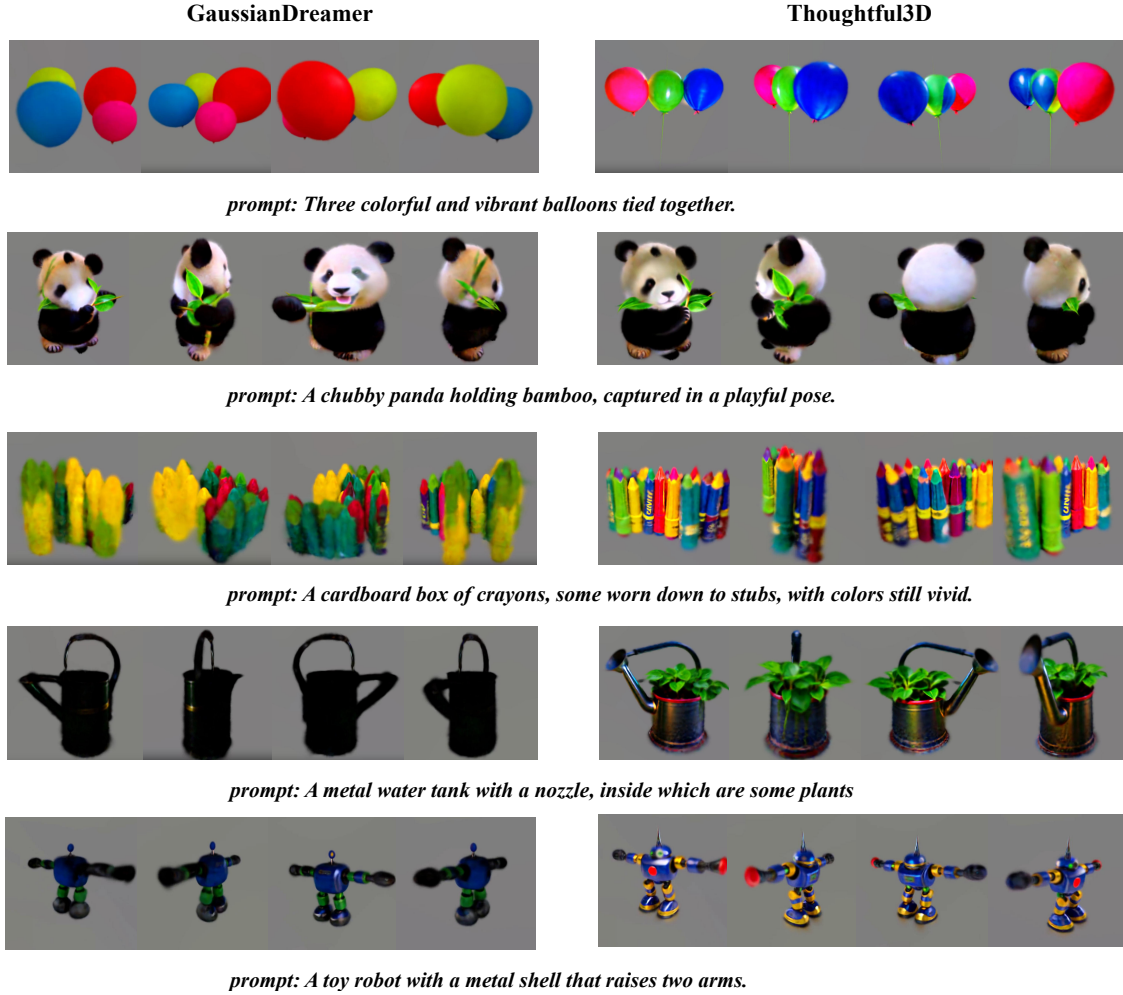


Figure 11. Analyze rendered images directly using MLLMs.

## 504 7.6. MLLM Sensitivity

505 As shown in Tab. 8, we evaluate our method across multiple open- and closed-source MLLM. performance differences among other high-capacity models like GPT-5.1 and Qwen2.5-VL-72B are minimal, suggesting low model sensitivity. Overall, our method remains stable across strong MLLMs and benefits from more capable models.

Table 8. Performance comparison across MLLMs.

Model	GaussianDreamer			DreamFusion		
	B/32	B/16	L/14	B/32	B/16	L/14
Baseline	21.45	26.71	27.33	16.40	22.98	23.22
<b>GPT-4o (Ours)</b>	<b>24.65</b>	<b>29.15</b>	<b>30.39</b>	<b>21.98</b>	<b>27.87</b>	<b>28.67</b>
GPT-5.1	25.17	30.21	30.87	22.36	29.01	29.94
Qwen2.5-VL-72B-Instruct	22.87	28.79	29.42	21.44	27.13	28.52

Table 9. Ablation on the number of candidate negative prompts.

	K=1	K=2	K=3	K=4	K=6	K=8
CLIP	28.43	29.62	<b>30.15</b>	30.21	30.17	30.20
Time (min)	16.7	18.3	<b>20.4</b>	23.5	30.1	36.9
$\Delta\text{CLIP}/\Delta\text{Time}$	-	0.744	<b>0.252</b>	0.019	0.006	0.004

## 511 7.7. Discussion on the test-time scaling

512 As shown in Tab. 9, we study the effect of the number of candidate negative prompts in 3DRefine-CoT. The metric  $\Delta\text{CLIP}/\Delta\text{Time}$  measures the marginal CLIP gain per additional minute relative to the previous  $K$ .  $K = 3$  yields the highest gain per minute, while further increasing  $K$  (e.g.,  $K = 6$ ) incurs higher computational cost but leads to degraded performance, indicating that additional candidates may introduce noise or negative effects. We therefore set  $K = 3$  as the best efficiency–performance trade-off.

## 521 7.8. Transition window for prompt switching.

522 As mentioned earlier, the transition window during prompt switching is utilized in our 3DBlueprint-CoT process. The

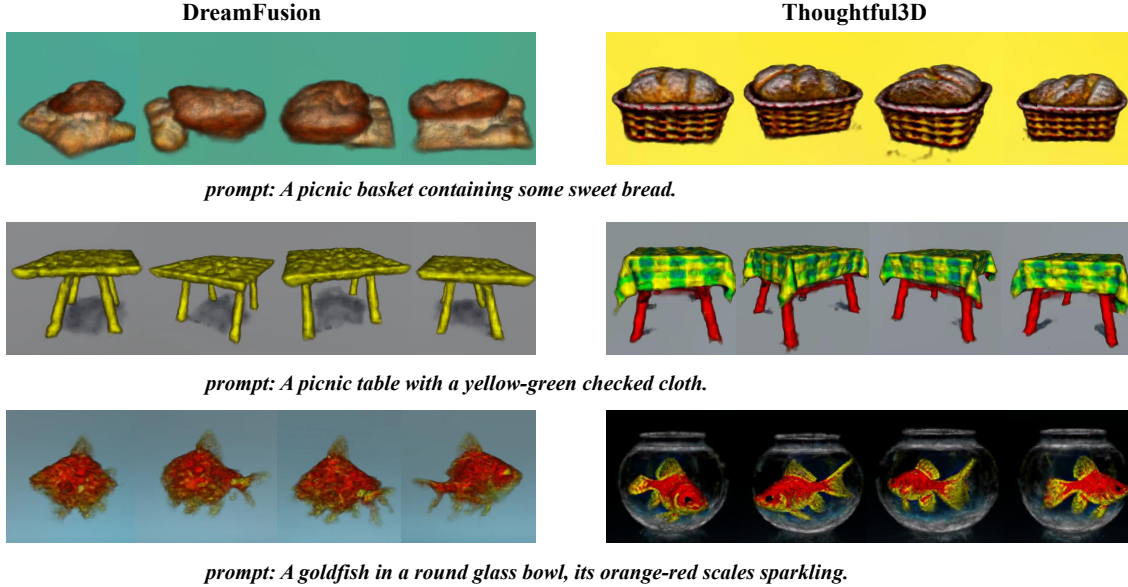


Figure 12. Analyze rendered images directly using MLLMs.

524 switching procedure for dynamic prompts can be expressed  
525 as:

$$526 \quad p_k(s) = \begin{cases} S_o, & s \in \tau_1, \\ S_c + \lambda_a(s)S_a, & s \in \tau_2, \\ S_c + S_a + \lambda_h(s)S_h, & s \in \tau_3, \end{cases} \quad (14)$$

527 The weighting  $\lambda_a(s)$  and  $\lambda_h(s)$  coefficients are defined as:

$$528 \quad \lambda_a(s) = \begin{cases} 0, & s \leq T_1, \\ \frac{s-T_1}{\Delta T}, & T_1 < s \leq T_1 + \Delta T, \\ 1, & s > T_1 + \Delta T, \end{cases} \quad (15)$$

$$529 \quad \lambda_h(s) = \begin{cases} 0, & s \leq T_2, \\ \frac{s-T_2}{\Delta T}, & T_2 < s \leq T_2 + \Delta T, \\ 1, & s > T_2 + \Delta T. \end{cases} \quad (16)$$

530 where  $\Delta T$  represents the width of the transition window  
531 during prompt switching, defined as the number of training  
532 iterations required for smooth transitions between stages. In  
533 our experiments, this parameter is typically set to 100.

## 534 8. Prompt Templates

535 We present the core prompts used in this work, specifically  
536 those applied in 3DBlueprint-CoT and 3DRefine-CoT, re-  
537 spectively. Fig. 13 and Fig. 15 belong to the 3DBlueprint-  
538 CoT. Fig. 13 parses the input prompt, while Fig. 15 plans  
539 the specific generation process. Fig. 16 and Fig. 14 belong  
540 to the 3DRefine-CoT. Fig. 16 provides a detailed descrip-  
541 tion and scoring of the rendered image. Fig. 14 analyzes  
542 issues and summarizes findings for the rendered image.

You are a prompt semantic expert.  
Analyze the 3D generation prompt through systematic linguistic  
decomposition.  
The input query is: "{query}". Please analyze which words  
belong to the main structure of the sentence, which words are  
adjectives, and which words are abstract terms.

Reference Analysis:  
Here is a reference example: For the sentence "An elegant  
flamingo standing tall with long legs and pinkish-white  
feathers." The main structure is: "A flamingo standing with legs  
and feathers.". The adjectives are: "tall, long, pinkish-white".  
The abstract word is: "elegant".  
then the returned result is:  
Structure: "A flamingo standing with legs and feathers."  
Adjectives: "tall, long, pinkish-white"  
Abstract: "elegant"

Output Format:  
Structure: "[Core structural sentence]"  
Adjectives: "[Comma-separated modifiers]"  
Abstract: "[Comma-separated abstract terms]"  
Please think carefully and provide your response.

Figure 13. Prompt of Semantic Parsing

## 543 9. limitations

544 Our approach functions as an optimization strategy built  
545 upon a baseline model, wherein the final output quality is  
546 directly contingent upon the baseline’s performance. Criti-  
547 cally, when the baseline exhibits suboptimal performance,  
548 our method demonstrates limited efficacy in substantially  
549 enhancing the quality and consistency of the generated 3D

You are a 3D rendering quality specialist. Two images are provided: image 1 is a composite of high-quality renderings of "{source\_text}", and image 2 is the target view described as "{caption}" with a current score of {score}/100 requiring improvement. Analyze image 2 through three problem assessments:

Problem 1: Identify structural,color defects including incorrect element count, duplicated part, spatial contradictions, and so no. (e.g., incorrect number of eyes, duplicated head, frontal features on the back).

Problem 2: Based on the {caption} and image 2, verify whether image 2 aligns with the original text {source\_text}, with regard to visual-textual consistency, presence of missing or extra elements, and potential stylistic deviations, as needed.

Problem 3: Compare image 2 with the reference rendering in image 1 to identify structural inconsistencies or abnormalities, such as additional ears, distorted limbs, or duplicated facial features.

Summarize critical problems as standardized negative prompts (e.g., "multi-head", "distorted color", "extra limbs","three ears").

Output Format:

Problem 1: [findings], Problem 2: [findings], Problem 3: [findings], Negative Prompts: [comma-separated descriptors]  
Please think carefully and provide your response.

Figure 14. Prompt of Summarize the negative prompts

550 content.

You are a linguistic reconstruction expert. Reconstruct the input prompt through three progressive enrichment stages using the provided structural analysis.

Input Components: Original prompt "{query}" and Structural analysis "{structure}".

Reconstruction Protocol:

Stage 1 (Core Skeleton): Extract essential syntactic elements (primary subject nouns, key action verbs, critical spatial prepositions) to form a minimal grammatical sentence.

Stage 2 (Concrete Enhancement): Integrate concrete descriptors (physical attributes, quantitative specifications, material indicators) while maintaining syntax.

Stage 3 (Abstract Enrichment): Incorporate non-tangible elements (subjective qualities, stylistic indicators, atmospheric cues) to achieve near-perfect semantic equivalence.

Reference Implementation:

Input Query: "An elegant flamingo standing tall with long legs and pinkish-white feathers."

Structural Analysis:

Structure: "A flamingo standing with legs and feathers"

Adjectives: tall, long, pinkish-white

Abstract: elegant

Output:

Stage 1: A flamingo standing with legs and feathers.

Stage 2: A flamingo standing with long legs and pinkish-white feathers.

Stage 3: An elegant flamingo standing tall with long legs and pinkish-white feathers.

Output Format Requirements:

Strictly use this sequence format:

Stage 1: [Core skeleton reconstruction]

Stage 2: [Concrete-enhanced reconstruction]

Stage 3: [Abstract-enriched reconstruction]

Please think carefully and provide your response.

Figure 15. Prompt of Generation Planning

You are a 3D rendering evaluation expert.

Analyze the image of the 3D object "{text}" by providing a comprehensive description covering geometric structure, proportions, material properties, surface textures, camera perspective, spatial relationships, and fidelity to textual specifications.

Then perform quantitative evaluation (1-100 scale) using weighted criteria: Text-Image Alignment (50% weight, evaluating accuracy of entities and attributes, spatial consistency, and style adherence) and Visual Plausibility (50% weight, assessing structural coherence, absence of artifacts, and perspective accuracy).

Calculate the final score as (Alignment Score  $\times$  0.5) + (Plausibility Score  $\times$  0.5).

Output the results as: Description: [your analysis] Keep the description within 100 words.

Score: [final score to one decimal place].

Please think carefully and provide your response.

Figure 16. Prompt of Reflection