

# Efficient All-Pairs Correlation Volume Sampling for Optical Flow Estimation

## Supplementary Material

Block Size	Row-Major	Patch-Major	Improvement
$1^2$		$1.6 \pm 0.4$	—
$2^2$	$2.5 \pm 0.5$	$2.3 \pm 0.5$	8%
$4^2$	$6.7 \pm 1.1$	$4.2 \pm 0.7$	37%
$8^2$	$20.6 \pm 4.4$	$8.6 \pm 1.4$	58%
$16^2$	$28.2 \pm 8.0$	$27.1 \pm 1.2$	4%

Table 3. Percentage of sampled correlation volume cells depending on the block size and layout, measured over the *SINTEL* [1] training dataset.

## 8. Method Details

In this section, we provide more information on obtaining the block mask, the sampling patterns, as well as complexity analysis and implementation details of our method.

### 8.1. Obtaining the Block Mask

In Figure 7 we provide a conceptual visualization of how the block mask is computed, which is then used for the correlation volume analysis and to determine which blocks need to be computed. Note that in practice it is obtained through splatting into the block mask directly, or computed implicitly through the inter-thread voting.

### 8.2. Correlation Volume Sampling Patterns

In Table 3, we show that averaging the number of sampled values per row-major block increases the ratio of cells that need to be computed, but is still significantly lower than the full matrix for sufficiently small block sizes. In the second column, we show that by averaging rows in our patch-major manner, sparsity is significantly increased without any computational overhead. At the largest measured block size of  $16^2$ , it becomes too large to maintain high level sparsity.

### 8.3. Implementation Details

In this subsection, we provide more details on our implementation of the fused correlation volume sampling algorithm, split in 3 parts according to our algorithmic design: 1) computing implicit computation mask; 2) joint computation of the next correlation block; 3) correlation value sampling. A simplified pseudocode is provided in Algorithm 1, kernel implementation is provided in Listings 2-4. It computes a single-level correlation whereas for multi-level pyramid an average-pooled  $F^2$  is used for each level. We set  $B = 8$  as a trade-off between being able to use efficient warp MMA operations for block computations and limiting the shared memory usage.

**Computing Implicit Block Mask.** To determine the blocks which need to be computed, in each thread we iterate over the  $[-r, r]^2$  neighborhood around the target pixel and compute indices of blocks the value needs to be sampled from. In practice, it is not necessary to iterate over all offsets, as for any 3 consecutive offsets, at least 2 land in the same block. Therefore, we only iterate over offsets with a stride of the block size  $B$ .

Indices of all blocks that a local thread needs to compute are stored in a register array in strictly increasing order, and we store an index pointer to the first block. To determine the next block to compute across all threads in a thread block, we perform a parallel reduction of the local indices with atomic minimum operation to a shared memory scalar, initialized larger than the size of the row. If the current local block matches the next global block, the index pointer is advanced. If the global next block is larger than the total number of blocks, all blocks have been processed and the kernel can exit.

**Computation of the Next Block.** Having agreed on the next block that needs to be computed, all threads participate in its computation, distributing loading of the input values and performing matrix-matrix computations across all threads. At once, we perform  $64 \times 64 \times 32$  product with  $16 \times 8 \times 16$  MMA instructions, `bf16` inputs and `fp32` accumulation. The final correlation values are kept in registers for faster processing.

**Processing of the Correlation Block.** Only threads that voted for the currently computed block sample correlation values by computing the target coordinates for each offset, and splitting them into the block index and local coordinates within the block. Values that can be extracted from the block are added to a global final output buffer, initialized as zeros.

### 8.4. Component Analysis

To analyze the impact of different steps of our algorithm to the total runtime, we approximate the time spent in each step by adding early kernel exits, and observe the aggregated runtime increase over the whole dataset. The runtime breakdown is provided in Tab. 4.

### 8.5. Memory Complexity Analysis

Here we show that our algorithm achieves linear time and space complexity in the number of pixels.

Assuming equal-size feature maps, let  $P = H \times W$  be the number of feature pixels per image,  $B$  the block size,

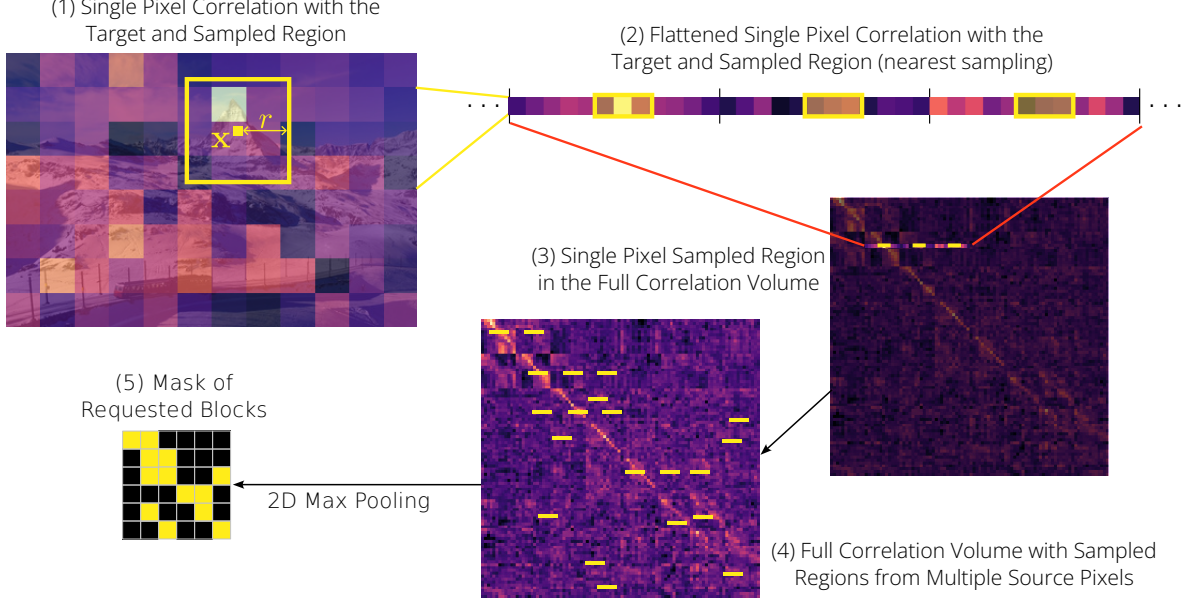


Figure 7. Conceptual visualization of the block mask computation.

<i>Correlation Volume Width</i> =	<b>256px</b>	<b>512px</b>
Preprocessing	11.3%	9.5%
Computing Per-Thread Blocks	<b>32.7%</b>	2.6%
Next-Block Voting	0.2%	0.3%
Matrix-Matrix Multiplication	0.4%	<b>25.5%</b>
Sampling Blocks and Writing Outputs	<b>55.3%</b>	<b>62.2%</b>
Total Runtime, ms	25.9	93.9

Table 4. Approximate runtime breakdown of the steps of our method.

and  $K = 2r + 1$  the lookup region size. We also assume that both input tiles are fully stored in shared memory. In practice, larger channel sizes  $D$  is split into chunks of 32 elements.

As all lookup offsets are continuous, each source pixel can sample from  $(K + 1)^2$  pixels. In the worst case, the first pixel is in a separate block, and, in total, they spread across  $(1 + \lceil \frac{(K + 1) - 1}{B} \rceil)^2 = \lceil \frac{K + 2B - 1}{B} \rceil^2$  blocks. With our hyperparameters  $r = 4, B = 8$ , it results in at most  $3^2 = 9$  blocks per pixel.

The matrix-matrix multiplication for computing each tile requires  $B^4 D$  operations, sampling it requires  $\min(B^2, K^2)$  operations (if the block is smaller than the lookup region, at most all values are sampled). The total time complexity across all steps of the algorithm is there-

fore:

$$\begin{aligned}
 & P \left[ \frac{K + 2B - 1}{B} \right]^2 \left( \underbrace{1}_{\text{Obtain mask}} + \underbrace{B^4 D}_{\text{Compute a tile}} + \underbrace{\min(B^2, K^2)}_{\text{Sample correlation}} \right) \\
 & \leq P(K + 2B - 1)^2 \frac{1}{B^2} (1 + B^4 D + \min(B^2, K^2)) \\
 & = P(K + 2B - 1)^2 \left( \frac{1}{B^2} + B^2 D + \min(1, \frac{K^2}{B^2}) \right) \quad (3) \\
 & \in O(P(K + B)^2 B^2 D),
 \end{aligned}$$

which is linear in the number of input pixels and feature dimensionality, and quadratic in the lookup radius. In practice, the number of requested blocks is lower due to smoothness of the optical flow and multiple source pixels requesting the same block, as shown in Table 3.

In addition to storing the inputs  $O(PD)$  and outputs  $O(PK^2)$ , we also need to use registers and shared memory to store the block indices and inputs and outputs of one computed tile per thread block, requiring:

$$\begin{aligned}
 & P \left[ \frac{K + 2B - 1}{B} \right]^2 + \frac{P}{B^2} (2B^2 D + B^4) \\
 & \leq P \left( \left( \frac{K - 1}{B^2} \right)^2 + 2 \right) + P(2D + B^2) \quad (4) \\
 & \in O(P(K^2 + D + B^2)),
 \end{aligned}$$

which, as inputs, is also linear in the number of pixels and feature dimensionality, quadratic in the lookup radius.

## 9. Extended Evaluation

In this section we provide extended evaluation results.

---

**Algorithm 1** Single-Level Correlation Volume Sampling

---

**Require:** Flattened input features  $F^{1,2} \in \mathbb{R}^{H \times W \times D}$ , lookup centroids  $\mathbf{X} \in \mathbb{R}^{N \times H \times W \times 2}$ , block size  $B$ , lookup radius  $r$ , output buffer  $\mathbf{O} \in \mathbb{R}^{N \times H \times W \times (2r+1)^2}$ .

```
 $bH, bW \leftarrow \lceil H/B \rceil, \lceil W/B \rceil$  ▷ Setting the number of blocks
 $\bar{F}^{1,2} \leftarrow \text{rearrange}(\bar{F}^{1,2}, '(bH \times B_h) \times (bW \times B_w) \times D \rightarrow (bH \times bW) \times (B_h \times B_w) \times D')$  ▷ Patch-major reshaping
 $nBlocks \leftarrow bH \cdot bW$  ▷ The total number of target blocks
for  $i = 0, 1, \dots, N - 1$  do ▷ For every lookup iteration
  On Device block  $by \in [bH \cdot bW]$ , thread  $t \in [B^2]$ 
   $\mathbf{x} \leftarrow \mathbf{X}[i][by \cdot B^2 + t]$  ▷ Get the respective target coordinates
  blocks  $\leftarrow [ ]$  ▷ Empty list of blocks to compute
  for all  $\mathbf{dx}' \in \{-r, -r+1, \dots, r, r+1\}^2$  do ▷ Find blocks for all offsets
    blockId  $\leftarrow \lfloor (\mathbf{x} + \mathbf{dx})/B \rfloor$ 
    if  $0 \leq \text{blockId} < \{bH, bW\}$  and (empty(blocks) or blockId > blocks[-1]) then
      blocks.append(blockId)
    end if
  end for
  blocks.append(nBlocks) ▷ Add end marker

  while True do
    nextBlock  $\leftarrow \text{inter-thread minimum of blocks.front()}$  ▷ Vote for the next block to compute
    if nextBlock  $\geq nBlocks$  then
      break ▷ All blocks have been processed
    end if

    corrTile  $\in \mathbb{R}^{B^2 \times B^2} \leftarrow (\bar{F}^1[by])(\bar{F}^2[\text{nextBlock}])^T$  ▷ Matrix-matrix multiplication of the block

    if nextBlock = blocks.front() then ▷ This thread requested the current block
      blocks.pop(0) ▷ Advance local block pointer
      for all  $\mathbf{dx} \in \{-r, -r+1, \dots, r\}^2$  do ▷ Sample from the computed block
         $\mathbf{x}' \leftarrow \mathbf{x} + \mathbf{dx} - \text{nextBlock} \cdot B$  ▷ Local coordinates within the block
         $\mathbf{O}[i][by \cdot B^2 + t][\mathbf{dx}] \leftarrow \text{sample corrTile}[t][\mathbf{x}']$ 
      end for
    end if
  end while
End Device
end for
```

---

## 9.1. Correlation Sampling Measurements

In Table 10, we provide full benchmarking results, corresponding to all runtime and memory usage plots. We further analyze the impact of the number of iterations, input feature dimensions, and GPU model. By default,  $(512 \times 224)^2$  correlation volume, 256 feature channels, and 32 flow update iterations are used.

### 9.1.1. Number of Iterations

In Figure 8, we show the impact of number of flow update iterations on the runtime and memory usage. All methods show approximately linear runtime increase and constant peak memory consumption. The gap between our method and the default implementation slightly narrows with in-

creasing number of iterations, as it has a fixed correlation volume precomputation cost, while in further steps only sampling is performed.

### 9.1.2. Feature Dimensionality

In Figure 9, we report the results at different feature  $F^{1,2}$  dimensionality. The runtime of all methods increases approximately linearly with the feature dimension, with slower rate for our method and the default implementation.

As both our method and on-demand sampling stores down-sampled replicas of the target features for multi-level computations, they show an increase of memory with larger dimensions. However, the absolute increase is small and typically negligible.

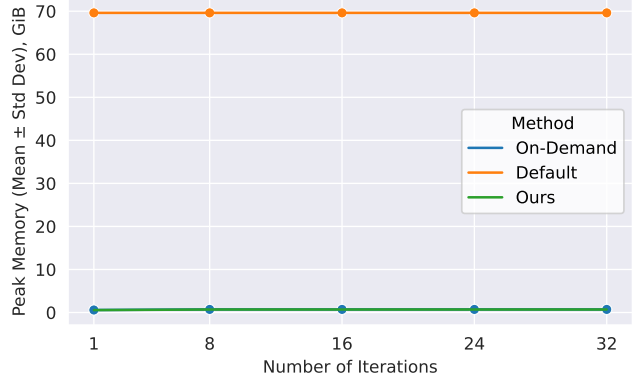
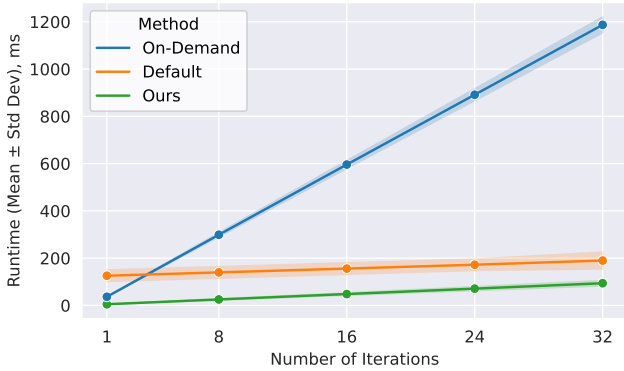


Figure 8. Runtime and peak memory consumption by the number of flow update iterations.

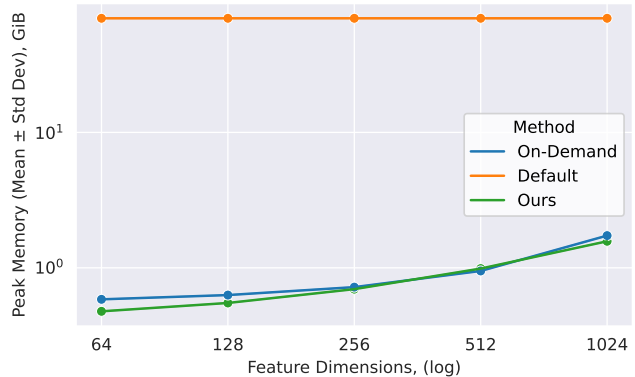
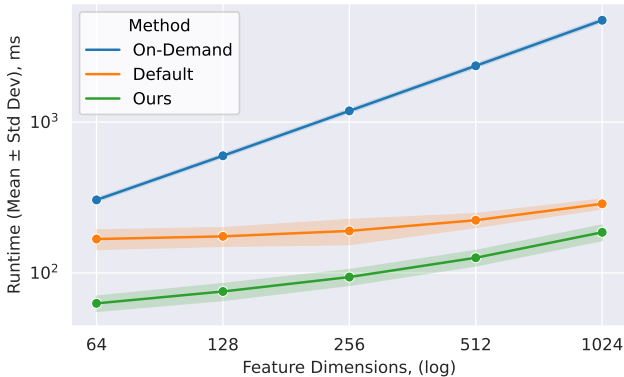


Figure 9. Runtime and peak memory consumption by input feature dimensionality.

Shot	Start Frame	End Frame
010_0050	101	196
040_0040	101	264
060_0130	101	175

Table 5. List of rendered CHARGE sequences.

### 9.1.3. Hardware

We perform additional experiments with different GPU models - *NVIDIA RTX 3090*, *NVIDIA RTX 4090*, and *NVIDIA A100* - and show results in Figure 10. It can be observed that hardware has little impact on the relative performance between different methods.

Note that for some GPU models, such as the *A100*, at  $128px$  correlation volume width, the overhead of pre-computing the full cost matrix becomes less significant, slightly outperforming our method. However, it still obtained at the expense of much higher memory requirements.

## 9.2. Dataset Generation Details

For the data generation, we use the *CYCLES* renderer to generate 335 frames from 3 sequences of the *BLENDER* movie *CHARGE* as listed in Table 5. Sample frames from each of the sequences can be seen in Figures 13-15.

To improve the quality of the rendered motion vectors, we remove lens flares and volumes, and disable motion blur to decrease noise levels. Following Mehl *et al.* [7], we render super-resolved motion vectors at 16K resolution with four ground truth values for each pixel. The *Python* script, used to apply scene adjustments, is provided in Listing 1.

## 9.3. Extended Results

In Table 9, we further analyze the improvements from our method by considering alternative strategies to reduce computational costs. First, we consider *SEA-RAFT* cascaded variants with more iterations at the lower resolution, but observe significant reduction in flow quality. Second, instead of using our correlation sampler, we run the default implementation on internal processing resolutions it is feasible on 80GB GPU, before switching to the on-demand sampler. While small improvement over on-demand sampler is achieved, it is significantly outperformed by our algorithm.

In Table 11, we provide the full results of all evaluated methods on the *Charge 8K* dataset. The first part of the table contains methods that do not employ correlation volume sampling and do not benefit from performance improvements using our efficient volume sampling algorithm.

In Figure 11, we visually plot the 1px error results and

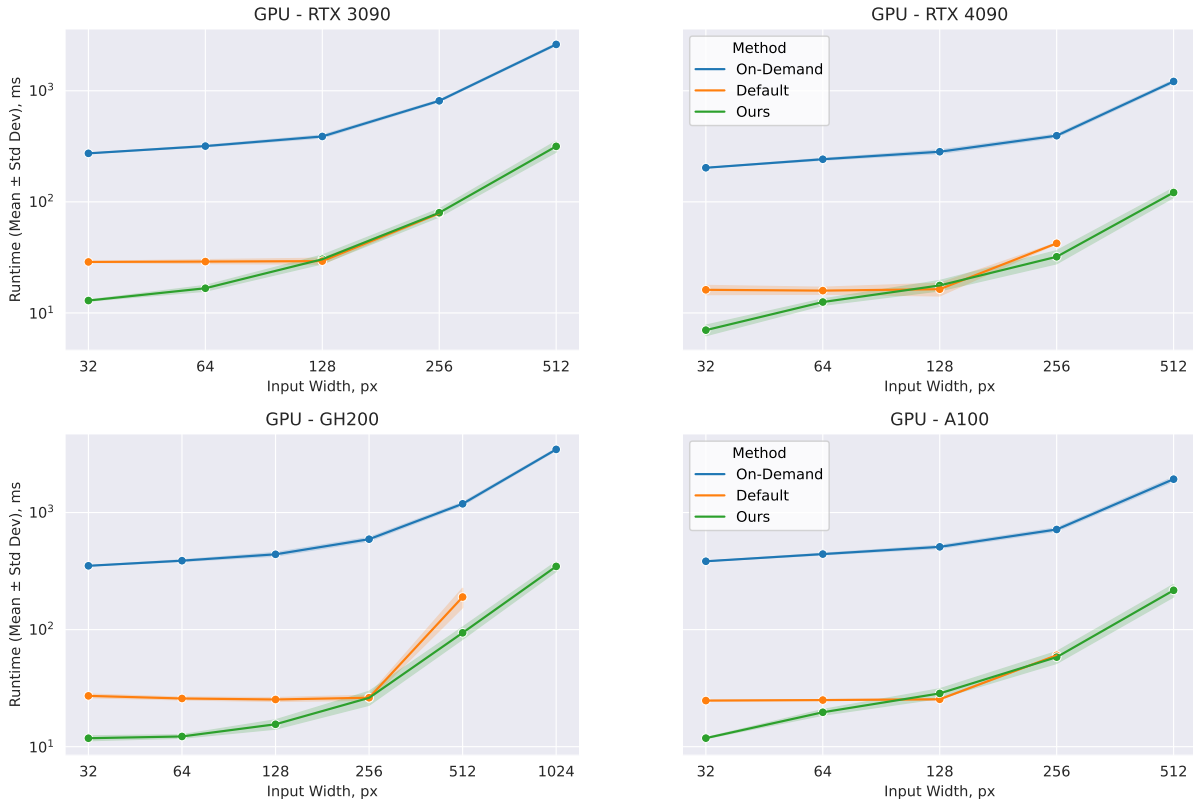


Figure 10. Runtime depending on input image size and different hardware.

runtime depending on the input scaling resolution. It can be observed that several methods degrade in quality when used with the highest resolution inputs.

In Figures 13-15, we provide additional qualitative comparisons.

To obtain results of all methods, we extend PTLFlow [9], and use checkpoints fine-tuned on the Sintel dataset, unless otherwise noted in Table 11. *FlowFormer* results were obtained with their default tiling technique. We do not run *WAF*T on the full resolution due to its very slow runtime ( $\approx 2$  min / image).

#### 9.4. Additional Datasets

We run similar evaluations on the *Kubric-8K* [10] dataset, which also contains ultra-high-resolution synthetic ground truth, but is obtained from generated scenes with moving objects. The respective results are provided in Figure 12 and Tables 6 and 7. Unlike the *Charge* dataset, where several methods benefit from fine-grained detail reconstruction at the highest resolution, on *Kubric-8K* most methods reach their lowest 1px error at quarter resolution. This could be explained from the use of less-detailed objects in the dataset. However, our algorithm still provides runtime and memory improvements across different evaluations scales (Table 7).

In addition to the *8K* datasets used in most experiments, we perform end-to-end runtime measurements on *Spring* [7], *Sintel* [1] and *KITTI* [8] datasets and report results in Table 8. Despite these datasets having a much lower resolution than our target, making the fast default implementation feasible, our method can still provide significant runtime or memory improvements. As the cascaded inference is not applicable on low-resolution inputs of these datasets, and we focus on performance improvements, we refer to their official benchmarks for full numerical results.

#### References

- [1] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In *European Conf. on Computer Vision (ECCV)*, pages 611–625. Springer-Verlag, 2012. 1, 5
- [2] Zhaoyang Huang, Xiaoyu Shi, Chao Zhang, Qiang Wang, Ka Chun Cheung, Hongwei Qin, Jifeng Dai, and Hongsheng Li. FlowFormer: A Transformer Architecture for Optical Flow. In *Computer Vision – ECCV 2022*, pages 668–685. Springer Nature Switzerland, Cham, 2022. 9, 16
- [3] Azin Jahedi, Maximilian Luz, Marc Rivinius, and Andrés Bruhn. CCMR: High Resolution Optical Flow Estimation via Coarse-To-Fine Context-Guided Motion Reasoning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 6899–6908, 2024. 9, 16

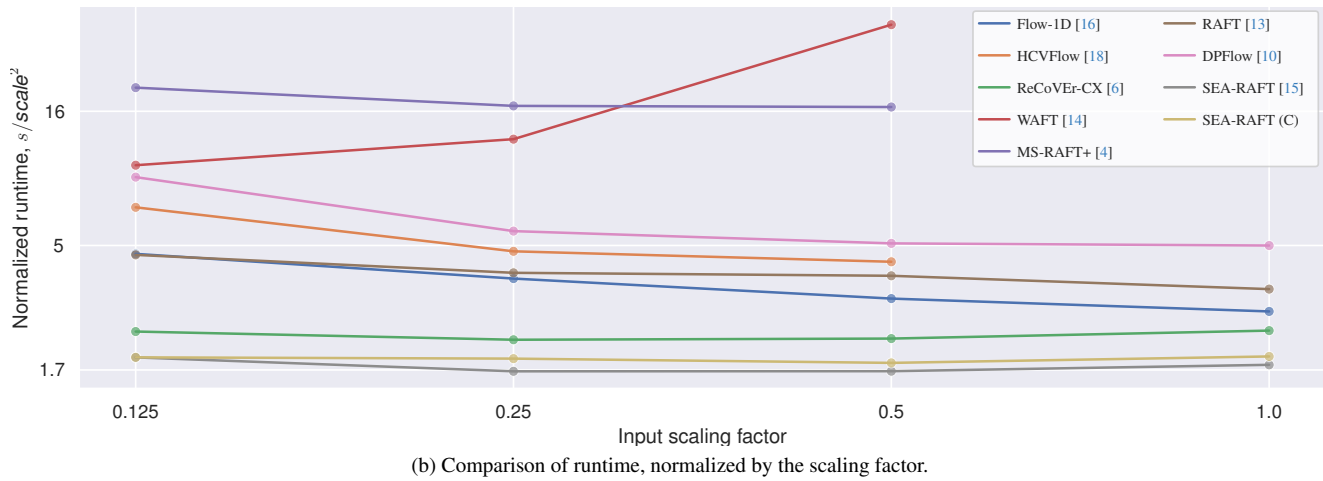
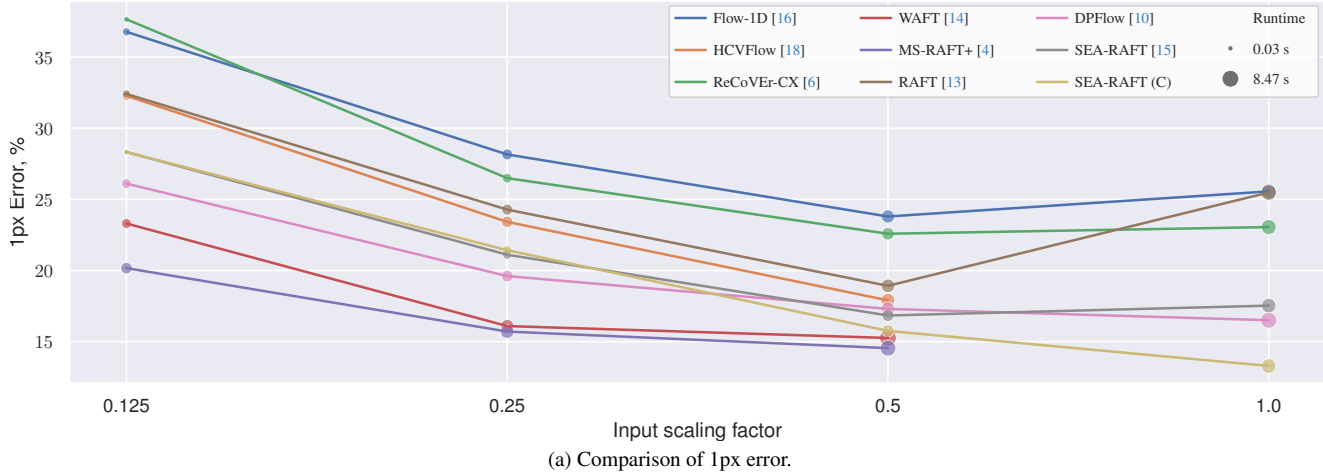


Figure 11. Quality and runtime comparison of selected methods on the CHARGE dataset depending on the input scaling resolution.

[4] Azin Jahedi, Maximilian Luz, Marc Rivinius, Lukas Mehl, and Andrés Bruhn. MS-RAFT+: High Resolution Multi-Scale RAFT. *International Journal of Computer Vision*, 132 (5):1835–1856, 2024. 6, 8, 9, 10, 16

[5] Shihao Jiang, Yao Lu, Hongdong Li, and Richard Hartley. Learning Optical Flow From a Few Matches. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16592–16600, 2021. 9, 16

[6] Simon Kiefhaber, Stefan Roth, and Simone Schaub-Meyer. Removing cost volumes from optical flow estimators. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 79–89, 2025. 6, 8, 9, 16

[7] Lukas Mehl, Jenny Schmalfluss, Azin Jahedi, Yaroslava Nalivayko, and Andrés Bruhn. Spring: A high-resolution high-detail dataset and benchmark for scene flow, optical flow and stereo. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 4, 5

[8] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3061–3070, 2015. 5

[9] Henrique Morimitsu. Pytorch lightning optical flow. <https://github.com/hmorimitsu/ptlflow>, 2021. 5

[10] Henrique Morimitsu, Xiaobin Zhu, Roberto M. Cesar, Xiangyang Ji, and Xu-Cheng Yin. DPFlow: Adaptive Optical Flow Estimation with a Dual-Pyramid Framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17810–17820, 2025. 5, 6, 8, 9, 16

[11] Xiaoyu Shi, Zhaoyang Huang, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. FlowFormer++: Masked Cost Volume Auto-encoding for Pretraining Optical Flow Estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1599–1610, 2023. 9, 16

[12] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8934–8943, 2018. 9, 16

[13] Zachary Teed and Jia Deng. RAFT: Recurrent All-Pairs

- Field Transforms for Optical Flow. In *Computer Vision – ECCV 2020*, pages 402–419, Cham, 2020. Springer International Publishing. [6](#), [8](#), [9](#), [16](#)
- [14] Yihan Wang and Jia Deng. WAFT: Warping-Alone Field Transforms for Optical Flow, 2025. [6](#), [8](#), [9](#), [16](#)
- [15] Yihan Wang, Lahav Lipson, and Jia Deng. SEA-RAFT: Simple, Efficient, Accurate RAFT for Optical Flow. In *Computer Vision – ECCV 2024*, pages 36–54. Springer Nature Switzerland, Cham, 2025. [6](#), [8](#), [9](#), [10](#), [16](#)
- [16] Haofei Xu, Jiaolong Yang, Jianfei Cai, Juyong Zhang, and Xin Tong. High-Resolution Optical Flow From 1D Attention and Correlation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10498–10507, 2021. [6](#), [8](#), [9](#), [16](#)
- [17] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezaatofghi, and Dacheng Tao. GMFlow: Learning Optical Flow via Global Matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8121–8130, 2022. [9](#), [16](#)
- [18] Yang Zhao, Gangwei Xu, and Gang Wu. Hybrid Cost Volume for Memory-Efficient Optical Flow. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 8740–8749, New York, NY, USA, 2024. Association for Computing Machinery. [6](#), [8](#), [9](#), [16](#)
- [19] Zihua Zheng, Ni Nie, Zhi Ling, Pengfei Xiong, Jiangyu Liu, Hao Wang, and Jiankun Li. DIP: Deep Inverse Patch-match for High-Resolution Optical Flow. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8925–8934, 2022. [9](#), [16](#)

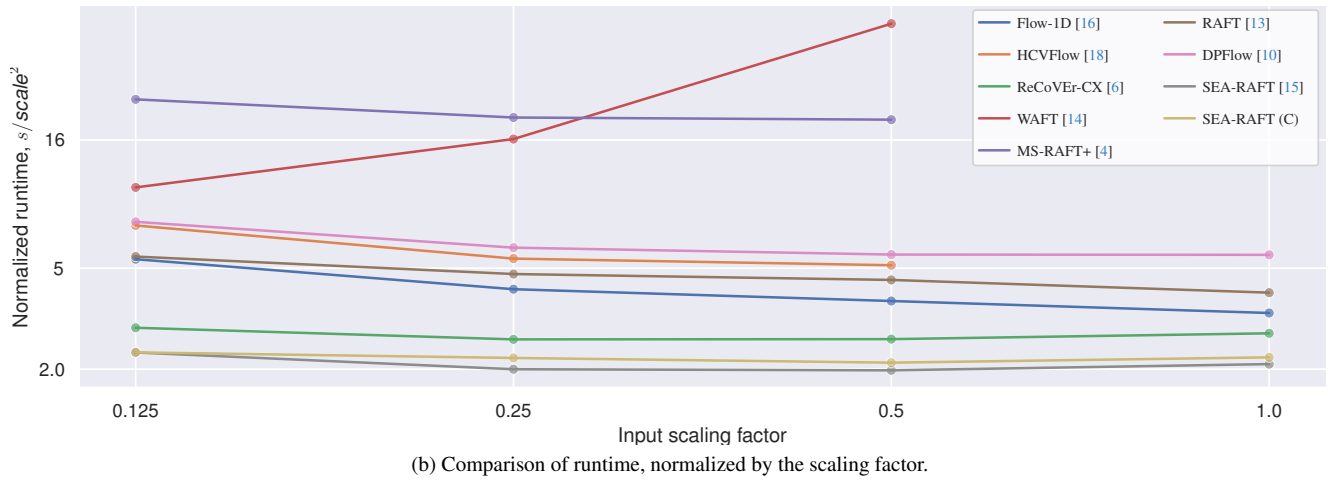
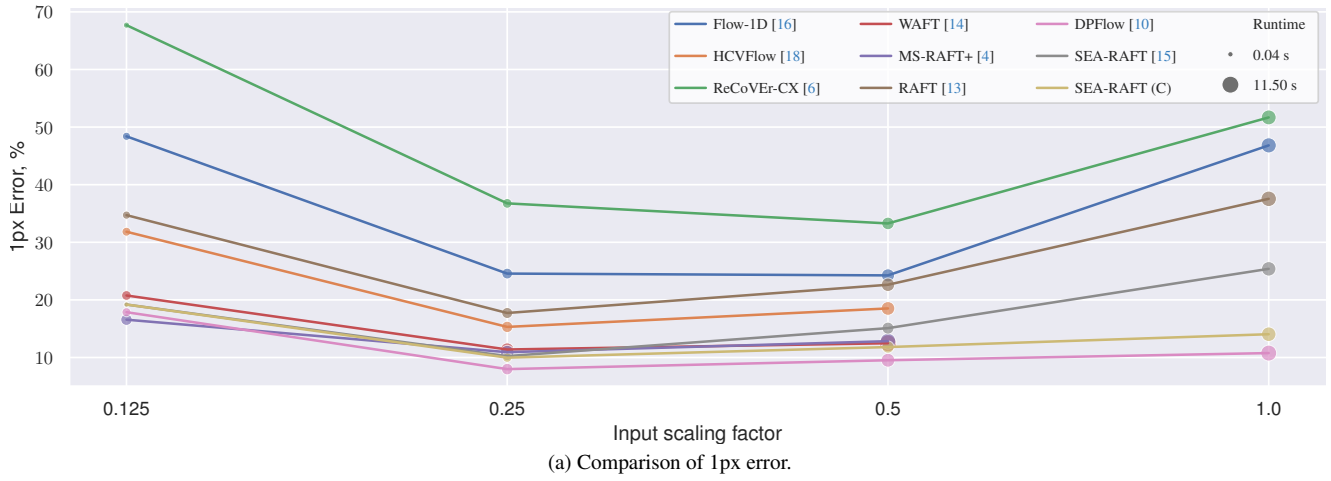


Figure 12. Quality and runtime comparison of selected methods on the *Kubric-8K* dataset depending on the input scaling resolution.

	Best-Accuracy Input Width	1px error % ↓	EPE px ↓	LM - 1px error % ↓	LM-EPE px ↓	Best Runtime, s	Without Ours, s	Our Improvement
GMFlow [17]	2048 - <sup>1</sup> / <sub>4</sub>	73.4	5.08	84.4	15.36	0.94	n/a	
PWC-Net [12]	4096 - <sup>1</sup> / <sub>2</sub>	19.7	21.79	44.5	98.94	0.22	n/a	
Flow-1D [16]	4096 - <sup>1</sup> / <sub>2</sub>	24.3	7.80	49.3	31.45	0.93	n/a	
DIP [19]	2048 - <sup>1</sup> / <sub>4</sub>	14.9	6.44	33.5	27.09	0.96	n/a	
FlowFormer [2]	2048 - <sup>1</sup> / <sub>4</sub>	15.6	4.10	33.7	15.92	1.79	n/a	
FlowFormer++ [11]	2048 - <sup>1</sup> / <sub>4</sub>	14.7	4.06	32.5	14.87	1.80	n/a	
SCV [5]	2048 - <sup>1</sup> / <sub>4</sub>	22.8	4.41	39.8	17.29	3.48	n/a	
HCVFlow [18]	2048 - <sup>1</sup> / <sub>4</sub>	15.3	4.34	36.2	17.61	0.34	n/a	
ReCoVer-CX [6]	4096 - <sup>1</sup> / <sub>2</sub>	33.3	20.33	83.2	93.02	0.66	n/a	
WAFT [14]	2048 - <sup>1</sup> / <sub>4</sub>	11.4	4.81	30.8	20.17	1.01	n/a	
MS-RAFT+ [4]	2048 - <sup>1</sup> / <sub>4</sub>	10.9	4.46	22.0	18.95	1.23	1.65	-26%
RAFT [13]	2048 - <sup>1</sup> / <sub>4</sub>	17.7	4.81	38.8	19.73	0.30	0.31	-4%
CCMR [3]	2048 - <sup>1</sup> / <sub>4</sub>	13.5	4.79	26.2	20.07	1.29	1.60	-20%
DPFlow [10]	2048 - <sup>1</sup> / <sub>4</sub>	<b>8.0</b>	<u>3.16</u>	<b>21.0</b>	<b>11.22</b>	0.38	0.42	-11%
SEA-RAFT [15]	2048 - <sup>1</sup> / <sub>4</sub>	10.2	4.25	23.7	17.92	<b>0.12</b>	0.17	-25%
SEA-RAFT (Cascaded) 1/2	4096 - <sup>1</sup> / <sub>2</sub>	11.8	3.95	<u>21.3</u>	14.75	0.53	0.87	-39%
SEA-RAFT (Cascaded)	2048 - <sup>1</sup> / <sub>4</sub>	<u>10.0</u>	<b>3.04</b>	22.5	<u>11.78</u>	<u>0.14</u>	0.18	-22%

Table 6. Quantitative evaluation of optical flow estimation methods on the *Kubric-8K* dataset. For each method, we list the resolution that can be run with 80GB of GPU memory and obtains the best 1px error. For a full comparison with other methods we also list our <sup>1</sup>/<sub>2</sub> results. We report the 1px outlier rate, endpoint-error (EPE), both metrics for pixels with large motion (LM, magnitude over 128px) and the best runtime across all variants with and without our improvements. We highlight the best (in **bold**) and second-best (underlined) method.

Method	Input Width	Default		On-Demand Sampling		Ours	
		Runtime	Memory	Runtime	Memory	Runtime	Memory
RAFT	960 - <sup>1</sup> / <sub>8</sub>	0.09 (-13%)	2.54 (+9%)	0.61 (+510%)	2.33 (=)	0.10	2.33
	1920 - <sup>1</sup> / <sub>4</sub>	0.31 (+4%)	9.98 (+200%)	0.99 (+235%)	3.33 (=)	0.30	3.33
	3840 - <sup>1</sup> / <sub>2</sub>		OOM	3.34 (+197%)	7.31 (=)	1.12	7.31
	7680 - <sup>1</sup> / <sub>1</sub>		OOM	12.49 (+212%)	23.26 (=)	4.00	23.26
MS-RAFT+	960 - <sup>1</sup> / <sub>8</sub>		OOM	0.52 (+43%)	3.73 (-1%)	0.36	3.76
	1920 - <sup>1</sup> / <sub>4</sub>		OOM	1.65 (+34%)	8.77 (-1%)	1.23	8.89
	3840 - <sup>1</sup> / <sub>2</sub>		OOM	6.46 (+34%)	28.75 (-2%)	4.81	29.25
DPFlow	960 - <sup>1</sup> / <sub>8</sub>	0.13 (+8%)	2.61 (+13%)	0.16 (+33%)	2.31 (=)	0.12	2.31
	1920 - <sup>1</sup> / <sub>4</sub>	0.42 (+12%)	10.48 (+244%)	0.45 (+19%)	3.05 (=)	0.38	3.05
	3840 - <sup>1</sup> / <sub>2</sub>		OOM	1.59 (+13%)	6.12 (=)	1.41	6.12
	7680 - <sup>1</sup> / <sub>1</sub>		OOM	6.25 (+11%)	18.45 (=)	5.65	18.45
SEA-RAFT	960 - <sup>1</sup> / <sub>8</sub>	0.04 (+2%)	2.60 (+17%)	0.10 (+175%)	2.23 (=)	0.04	2.23
	1920 - <sup>1</sup> / <sub>4</sub>	0.17 (+33%)	10.07 (+271%)	0.21 (+70%)	2.71 (=)	0.12	2.71
	3840 - <sup>1</sup> / <sub>2</sub>		OOM	0.77 (+55%)	4.66 (=)	0.49	4.66
	7680 - <sup>1</sup> / <sub>1</sub>		OOM	3.13 (+50%)	12.46 (=)	2.09	12.46

Table 7. Runtime (s) and peak memory usage (GiB) of the full optical flow method end-to-end evaluation on *Kubric-8K* dataset depending on the correlation computation variant at different scales of the inputs. We report the relative difference compared to our method in parentheses. OOM indicates that the method requires more than 80GB of memory and fails with an out-of-memory error.

Method	Dataset	Default		On-Demand Sampling		Ours	
		Runtime	Memory	Runtime	Memory	Runtime	Memory
RAFT	<i>KITTI 2015</i>	0.08 (-29%)	0.50 (+41%)	0.60 (+425%)	0.35 (=)	0.11	0.35
	<i>Sintel</i>	0.08 (-16%)	0.53 (+57%)	0.58 (+518%)	0.34 (=)	0.09	0.34
	<i>Spring</i>	0.32 (+7%)	8.22 (+423%)	0.98 (+227%)	1.57 (=)	0.30	1.57
MS-RAFT+	<i>KITTI 2015</i>		OOM	0.48 (+34%)	1.64 (-2%)	0.36	1.67
	<i>Sintel</i>		OOM	0.46 (+39%)	1.57 (-2%)	0.33	1.60
	<i>Spring</i>		OOM	1.68 (+34%)	7.01 (-2%)	1.25	7.14
DPFLow	<i>KITTI 2015</i>	0.12 (+11%)	0.57 (+63%)	0.15 (+39%)	0.35 (=)	0.11	0.35
	<i>Sintel</i>	0.12 (+9%)	0.53 (+56%)	0.14 (+33%)	0.34 (=)	0.11	0.34
	<i>Spring</i>	0.43 (+13%)	8.72 (+579%)	0.45 (+17%)	1.29 (=)	0.39	1.29
SEA-RAFT	<i>KITTI 2015</i>	0.04 (-33%)	0.55 (+115%)	0.10 (+81%)	0.26 (=)	0.05	0.26
	<i>Sintel</i>	0.03 (=)	0.58 (+64%)	0.09 (+184%)	0.34 (-4%)	0.03	0.35
	<i>Spring</i>	0.17 (+34%)	8.31 (+773%)	0.21 (+66%)	0.95 (=)	0.13	0.95

Table 8. Runtime (s) and peak memory usage (GiB) of the full optical flow end-to-end evaluation depending on the correlation computation method on three additional datasets - *KITTI 2015*, *Sintel*, and *Spring*.

		Metrics				Runtime (s) / Memory (GiB)			
		1px error	EPE	LM-1px	LM-EPE	Default	On-Demand	Mixed	Ours
MS-RAFT+ [4]	1/8	20.2	1.94	64.1	22.57	OOM	0.44 / 4.4	0.38 / 9.2	0.30 / 4.4
	1/4	15.7	<b>1.64</b>	41.1	<b>19.85</b>	OOM	1.40 / 8.7	1.35 / 9.9	1.04 / 8.8
	1/2	<u>14.5</u>	1.92	<u>34.6</u>	32.03	OOM	5.47 / 25.8	5.45 / 25.8	4.11 / 26.2
SEA-RAFT [15] 4LR + 0HR (default)	1/8	28.3	2.30	81.5	21.51	0.03 / 3.5	0.09 / 3.2	0.03 / 3.5	0.03 / 3.3
	1/4	21.1	<u>2.01</u>	54.5	<u>19.18</u>	0.13 / 8.9	0.19 / 3.6	0.13 / 8.9	0.10 / 3.6
	1/2	<u>16.8</u>	4.94	<u>39.8</u>	39.83	OOM	0.63 / 5.2	0.62 / 5.2	0.42 / 5.2
	1/1	17.5	17.65	49.2	107.43	OOM	2.63 / 11.8	2.63 / 11.8	1.78 / 11.8
SEA-RAFT 3LR + 1HR	1/8	31.3	2.67	84.8	27.48	0.02 / 3.5	0.04 / 3.2	0.02 / 3.5	0.02 / 3.2
	1/4	25.9	2.25	62.9	20.29	0.12 / 8.9	0.15 / 3.6	0.12 / 8.9	0.09 / 3.6
	1/2	20.7	1.93	44.3	18.41	OOM	0.44 / 5.2	0.40 / 5.2	0.35 / 5.2
	1/1	<u>17.9</u>	<u>1.83</u>	<u>37.0</u>	<u>17.04</u>	OOM	1.71 / 11.8	1.64 / 11.8	1.42 / 11.8
SEA-RAFT 6LR + 2HR	1/8	29.6	2.42	82.3	22.90	0.03 / 3.5	0.05 / 3.2	0.02 / 3.5	0.02 / 3.3
	1/4	23.1	2.03	56.1	17.82	0.13 / 8.9	0.22 / 3.6	0.13 / 8.9	0.10 / 3.6
	1/2	18.0	1.84	39.0	17.31	OOM	0.58 / 5.2	0.49 / 5.2	0.39 / 5.2
	1/1	<u>15.0</u>	<u>1.82</u>	<u>33.2</u>	<u>17.05</u>	OOM	2.19 / 11.8	2.02 / 11.8	1.58 / 11.8
SEA-RAFT 4LR + 4HR (Cascaded)	1/8	28.3	2.30	81.5	21.51	0.03 / 3.5	0.09 / 3.2	0.03 / 3.5	0.03 / 3.3
	1/4	21.4	<u>1.86</u>	54.1	<b>16.15</b>	0.15 / 8.9	0.25 / 3.6	0.15 / 8.9	0.12 / 3.6
	1/2	15.8	1.89	36.8	18.55	OOM	0.71 / 5.2	0.65 / 5.2	0.45 / 5.2
	1/1	<b>13.3</b>	2.70	<b>31.6</b>	21.53	OOM	2.88 / 11.8	2.77 / 11.8	1.89 / 11.8

Table 9. Analysis of alternative inference strategies on CHARGE8K dataset, with larger number of iterations at low 1/4 resolution (LR) than high resolution (HR), and mixed correlation sampler, running default implementation up to the highest feasible scale before switching to the on-demand sampler.

```

import bpy

# Choose to render RGB or Flow pass
render_flow = False

for scene in bpy.data.scenes:

    width = 1024
    aspect_ratio = width / scene.render.resolution_x
    height = int(round(scene.render.resolution_y * aspect_ratio))

    scale = 16 if render_flow else 8
    scene.render.resolution_x = width * scale
    scene.render.resolution_y = height * scale

    scene.render.use_motion_blur = False
    scene.render.engine = "CYCLES"
    scene.cycles.device = "CPU"

    n_samples = 1 if render_flow else 1024
    scene.cycles.samples = n_samples
    scene.cycles.adaptive_min_samples = n_samples > 1
    scene.cycles.adaptive_threshold = 0.01
    scene.cycles.use_adaptive_sampling = True
    scene.cycles.denoiser = "OPENIMAGEDENOISE"
    scene.cycles.use_denoising = (not render_flow)

    # Removing very rare but strong fireflies
    scene.cycles.sample_clamp_direct = 50.0

    # Removing lens flares and overlays
    scene.use_nodes = False

    collection = bpy.data.collections.get("flares")
    if collection:
        bpy.data.collections.remove(collection)

    for view_layer in scene.view_layers:
        view_layer.use_pass_combined = True

# Removing volumes
for mat in bpy.data.materials:
    if mat.use_nodes:
        for node in mat.node_tree.nodes:
            if node.bl_static_type == "OUTPUT_MATERIAL" and node.is_active_output:
                for link in node.inputs["Volume"].links:
                    mat.node_tree.links.remove(link)

```

Listing 1. Blender scene setup script.

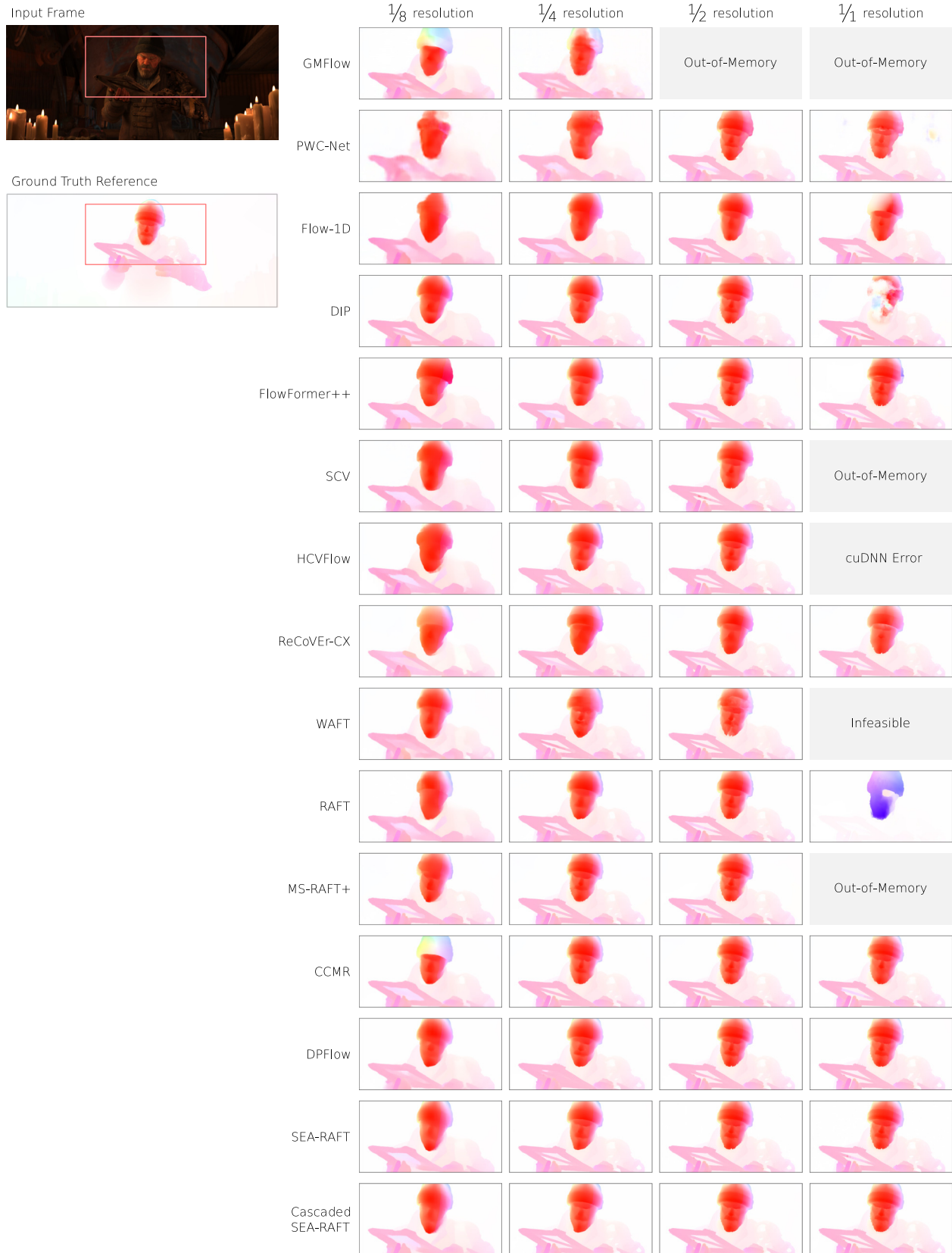


Figure 13. Qualitative comparison across different evaluation scales on a frame from 010\_0050 shot. Image from Charge by Blender Studio.

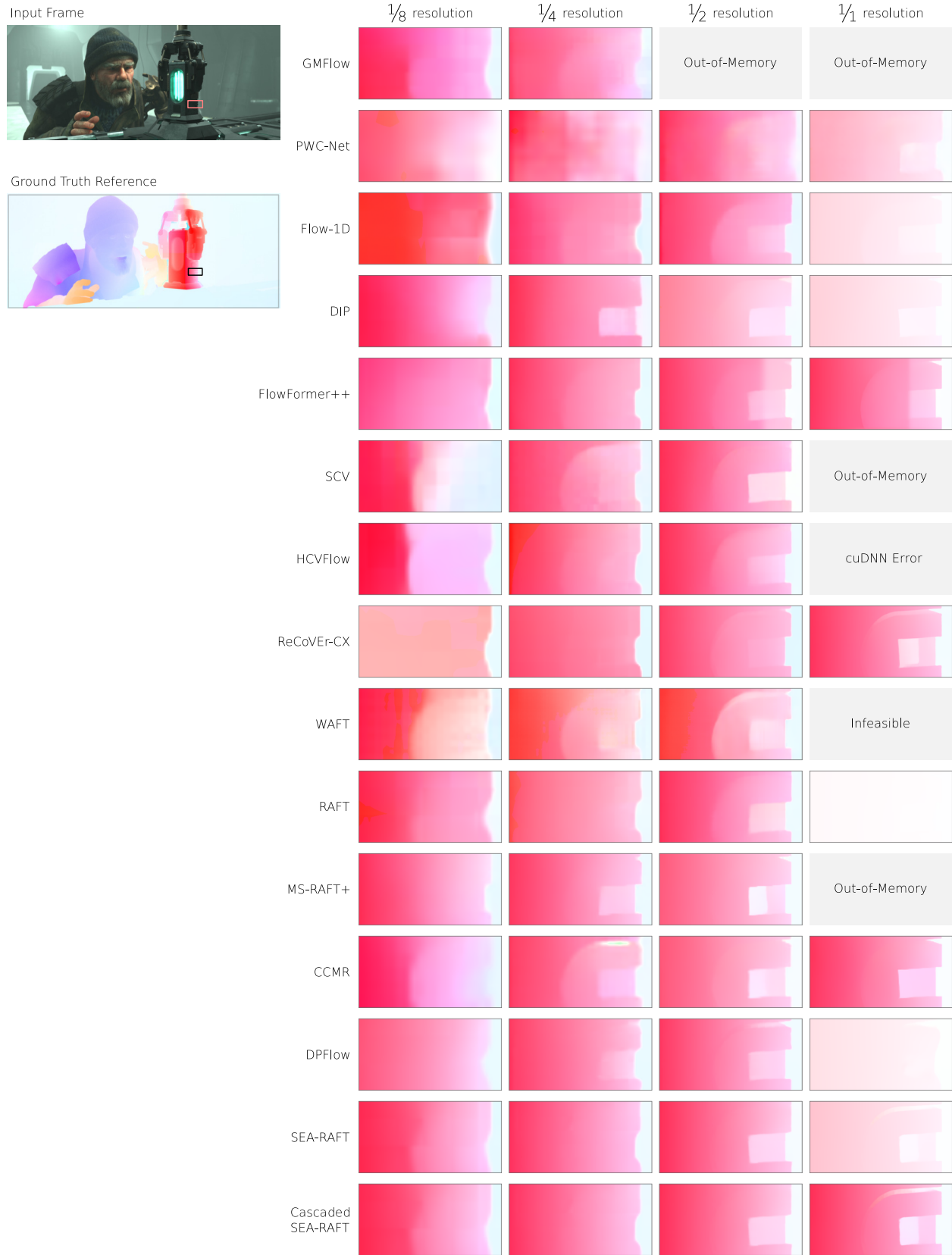


Figure 14. Qualitative comparison across different evaluation scales on a frame from 040\_0040 shot. Image from Charge by Blender Studio.



Figure 15. Qualitative comparison across different evaluation scales on a frame from 060\_0130 shot. Image from Charge by Blender Studio.

Variable	Value	Runtime (mean $\pm$ std), ms				Peak Memory (mean $\pm$ std), MB			
		Default	On-Demand	Ours	Improvement	Default	On-Demand	Ours	Improvement
Input Width GH200	32	27 $\pm$ 1	351 $\pm$ 0	12 $\pm$ 1	96.6%	4 $\pm$ 0	3 $\pm$ 0	3 $\pm$ 0	13.3%
	64	26 $\pm$ 1	388 $\pm$ 6	12 $\pm$ 0	96.8%	27 $\pm$ 1	12 $\pm$ 0	12 $\pm$ 0	55.3%
	128	25 $\pm$ 1	440 $\pm$ 15	16 $\pm$ 1	96.5%	291 $\pm$ 0	47 $\pm$ 0	46 $\pm$ 0	84.3%
	256	26 $\pm$ 1	592 $\pm$ 20	26 $\pm$ 4	95.6%	4091 $\pm$ 0	184 $\pm$ 0	178 $\pm$ 0	95.6%
	512	190 $\pm$ 35	1187 $\pm$ 33	94 $\pm$ 11	92.1%	71288 $\pm$ 0	737 $\pm$ 0	712 $\pm$ 0	99.0%
	1024	OOM	3462 $\pm$ 62	347 $\pm$ 32	90.0%	OOM	2944 $\pm$ 0	2829 $\pm$ 0	-
Input Width A100	2048	OOM	13522 $\pm$ 288	1334 $\pm$ 133	90.1%	OOM	11764 $\pm$ 0	11274 $\pm$ 0	-
	32	25 $\pm$ 0	384 $\pm$ 0	12 $\pm$ 0	96.9%	4 $\pm$ 0	3 $\pm$ 0	3 $\pm$ 0	13.3%
	64	25 $\pm$ 0	442 $\pm$ 6	20 $\pm$ 1	95.5%	27 $\pm$ 0	12 $\pm$ 0	12 $\pm$ 0	55.3%
	128	25 $\pm$ 0	509 $\pm$ 15	29 $\pm$ 3	94.4%	291 $\pm$ 0	47 $\pm$ 0	46 $\pm$ 0	84.3%
	256	60 $\pm$ 0	717 $\pm$ 23	58 $\pm$ 7	91.9%	4090 $\pm$ 0	184 $\pm$ 0	178 $\pm$ 0	95.6%
	512	OOM	1928 $\pm$ 73	217 $\pm$ 26	88.8%	OOM	737 $\pm$ 0	712 $\pm$ 0	-
Input Width RTX 4090	32	16 $\pm$ 2	203 $\pm$ 0	7 $\pm$ 1	96.5%	4 $\pm$ 0	3 $\pm$ 0	3 $\pm$ 0	13.3%
	64	16 $\pm$ 1	242 $\pm$ 4	13 $\pm$ 1	94.8%	27 $\pm$ 0	12 $\pm$ 0	12 $\pm$ 0	55.3%
	128	16 $\pm$ 2	283 $\pm$ 11	18 $\pm$ 2	93.7%	291 $\pm$ 0	47 $\pm$ 0	46 $\pm$ 0	84.3%
	256	42 $\pm$ 0	394 $\pm$ 12	32 $\pm$ 4	91.9%	4090 $\pm$ 0	184 $\pm$ 0	178 $\pm$ 0	95.6%
	512	OOM	1213 $\pm$ 34	121 $\pm$ 12	90.0%	OOM	737 $\pm$ 0	712 $\pm$ 0	-
Input Width RTX 3090	32	29 $\pm$ 0	273 $\pm$ 1	13 $\pm$ 0	95.3%	4 $\pm$ 0	3 $\pm$ 0	3 $\pm$ 0	13.3%
	64	29 $\pm$ 1	318 $\pm$ 5	17 $\pm$ 1	94.8%	27 $\pm$ 0	12 $\pm$ 0	12 $\pm$ 0	55.3%
	128	29 $\pm$ 2	388 $\pm$ 11	30 $\pm$ 3	92.2%	291 $\pm$ 0	47 $\pm$ 0	46 $\pm$ 0	84.3%
	256	79 $\pm$ 0	812 $\pm$ 15	80 $\pm$ 6	90.2%	4090 $\pm$ 0	184 $\pm$ 0	178 $\pm$ 0	95.6%
	512	OOM	2615 $\pm$ 49	316 $\pm$ 35	87.9%	OOM	737 $\pm$ 0	712 $\pm$ 0	-
Number of Iterations GH200	1	126 $\pm$ 24	37 $\pm$ 0	5 $\pm$ 0	86.3%	71288 $\pm$ 0	599 $\pm$ 0	574 $\pm$ 0	99.2%
	8	140 $\pm$ 24	299 $\pm$ 8	25 $\pm$ 2	91.5%	71288 $\pm$ 0	737 $\pm$ 0	712 $\pm$ 0	99.0%
	16	156 $\pm$ 24	596 $\pm$ 16	48 $\pm$ 5	91.9%	71288 $\pm$ 0	737 $\pm$ 0	712 $\pm$ 0	99.0%
	24	172 $\pm$ 24	891 $\pm$ 25	71 $\pm$ 8	92.0%	71288 $\pm$ 0	737 $\pm$ 0	712 $\pm$ 0	99.0%
	32	190 $\pm$ 35	1187 $\pm$ 33	94 $\pm$ 11	92.1%	71288 $\pm$ 0	737 $\pm$ 0	712 $\pm$ 0	99.0%
Feature Dimensionality GH200	64	168 $\pm$ 25	305 $\pm$ 8	63 $\pm$ 7	79.4%	71288 $\pm$ 0	599 $\pm$ 0	489 $\pm$ 0	99.3%
	128	175 $\pm$ 24	598 $\pm$ 16	75 $\pm$ 10	87.4%	71288 $\pm$ 0	644 $\pm$ 0	563 $\pm$ 0	99.2%
	256	190 $\pm$ 35	1187 $\pm$ 33	94 $\pm$ 11	92.1%	71288 $\pm$ 0	737 $\pm$ 0	712 $\pm$ 0	99.0%
	512	224 $\pm$ 23	2362 $\pm$ 68	126 $\pm$ 14	94.7%	71288 $\pm$ 0	971 $\pm$ 0	1011 $\pm$ 0	98.6%
	1024	287 $\pm$ 21	4734 $\pm$ 145	186 $\pm$ 22	96.1%	71288 $\pm$ 0	1771 $\pm$ 0	1608 $\pm$ 0	97.7%

Table 10. Full correlation volume sampling isolated benchmarking results depending on one variable. We report the runtime and peak memory usage of each implementation, as well as memory improvement over default implementation and runtime improvement over the on-demand sampling method. OOM indicates that the method fails with an out-of-memory error.

		Metrics				Runtime		
		lpx error	EPE	LM-lpx	LM-EPE	Default	On-Demand	Ours
GMFlow [17]	1/8	<u>43.9</u>	2.95	94.3	<u>24.46</u>	0.09 ± 0.00		n/a
	1/4	47.0	<u>2.74</u>	<u>86.3</u>	28.61	0.71 ± 0.00		n/a
PWC-Net [12]	1/8	42.8	3.89	93.6	47.94	0.07 ± 0.04		n/a
	1/4	30.3	3.32	78.6	<u>47.27</u>	0.10 ± 0.05		n/a
	1/2	<u>24.4</u>	<u>3.26</u>	<u>61.4</u>	57.14	0.21 ± 0.07		n/a
	1/1	26.0	6.65	76.1	128.30	0.64 ± 0.10		n/a
Flow-1D [16] Highres	1/8	36.8	3.14	95.2	37.06	0.07 ± 0.01		n/a
	1/4	28.2	2.39	83.7	<u>29.80</u>	0.23 ± 0.05		n/a
	1/2	<u>23.8</u>	<u>2.23</u>	71.3	31.58	0.79 ± 0.02		n/a
	1/1	25.6	5.86	<u>69.2</u>	57.55	2.82 ± 0.01		n/a
DIP [19]	1/8	28.3	2.44	84.4	<u>28.92</u>	0.28 ± 0.01		n/a
	1/4	22.5	<u>2.23</u>	60.6	33.65	0.84 ± 0.00		n/a
	1/2	20.8	2.76	<u>46.0</u>	56.97	2.93 ± 0.01		n/a
	1/1	<u>19.8</u>	4.00	51.2	85.62	11.57 ± 0.06		n/a
FlowFormer [2]	1/8	33.2	2.62	89.0	26.13	0.79 ± 0.00		n/a
	1/4	25.0	<u>2.19</u>	68.0	<u>22.20</u>	1.81 ± 0.01		n/a
	1/2	19.2	2.31	<u>52.1</u>	33.86	5.06 ± 0.03		n/a
	1/1	<u>16.9</u>	3.22	56.4	58.16	16.57 ± 0.06		n/a
FlowFormer++ [11]	1/8	32.0	2.53	88.6	22.44	0.80 ± 0.00		n/a
	1/4	24.1	<u>2.00</u>	66.5	<u>20.04</u>	1.82 ± 0.01		n/a
	1/2	18.9	2.46	<u>51.9</u>	32.26	5.09 ± 0.03		n/a
	1/1	<u>16.8</u>	3.41	55.3	55.70	16.68 ± 0.05		n/a
SCV [5]	1/8	35.3	2.60	85.8	29.57	0.94 ± 0.08		n/a
	1/4	23.7	<u>2.01</u>	61.0	<u>23.19</u>	3.08 ± 0.21		n/a
	1/2	<u>19.1</u>	2.83	<u>46.1</u>	46.42	14.66 ± 0.88		n/a
HCVFlow [18]	1/8	32.3	2.61	91.0	28.68	0.11 ± 0.00		n/a
	1/4	23.4	<u>1.99</u>	70.9	<u>25.37</u>	0.30 ± 0.03		n/a
	1/2	<u>17.9</u>	2.08	<u>54.4</u>	32.92	1.09 ± 0.01		n/a
ReCoVer-CX [6] Mixed	1/8	37.7	3.20	97.3	43.77	0.04 ± 0.00		n/a
	1/4	26.5	<u>2.67</u>	95.1	<u>40.80</u>	0.14 ± 0.00		n/a
	1/2	<u>22.6</u>	3.01	<u>94.6</u>	58.71	0.56 ± 0.00		n/a
	1/1	23.1	4.98	99.6	114.81	2.39 ± 0.01		n/a
WAFT [14]	1/8	23.3	2.05	83.1	<u>21.35</u>	0.16 ± 0.03		n/a
	1/4	16.1	<u>1.93</u>	<u>62.0</u>	27.44	0.79 ± 0.01		n/a
	1/2	<u>15.3</u>	2.74	62.6	65.14	8.47 ± 0.07		n/a
MS-RAFT+ [4]	1/8	20.2	1.94	64.1	22.59	OOM	0.44 ± 0.00	0.31 ± 0.00
	1/4	15.7	<u>1.64</u>	41.1	<u>19.85</u>	OOM	1.41 ± 0.01	1.05 ± 0.00
	1/2	<u>14.5</u>	1.92	<u>34.6</u>	32.03	OOM	5.51 ± 0.02	4.15 ± 0.02
RAFT [13]	1/8	32.4	2.63	90.3	26.73	0.07 ± 0.00	0.56 ± 0.01	0.09 ± 0.00
	1/4	24.3	<u>2.01</u>	69.8	<u>22.16</u>	0.25 ± 0.00	0.89 ± 0.01	0.25 ± 0.00
	1/2	<u>18.9</u>	5.90	<u>49.8</u>	36.41	OOM	2.62 ± 0.03	0.96 ± 0.01
	1/1	25.5	59.94	71.3	266.20	OOM	10.47 ± 0.13	3.43 ± 0.04
CCMR [3]	1/8	25.9	2.14	76.1	<u>24.40</u>	0.34 ± 0.02	0.47 ± 0.02	0.33 ± 0.01
	1/4	17.5	<u>1.88</u>	50.5	27.71	OOM	1.35 ± 0.04	1.09 ± 0.02
	1/2	<u>15.8</u>	2.16	<u>39.7</u>	42.86	OOM	5.23 ± 0.02	4.39 ± 0.01
	1/1	18.0	4.08	45.7	76.96	OOM	21.09 ± 0.15	17.58 ± 0.02
DPFlow [10]	1/8	26.1	2.00	82.4	18.88	0.14 ± 0.01	0.15 ± 0.01	0.14 ± 0.02
	1/4	19.6	<u>1.60</u>	57.8	<u>13.83</u>	0.40 ± 0.01	0.41 ± 0.01	0.35 ± 0.01
	1/2	17.3	1.71	40.3	14.29	OOM	1.41 ± 0.00	1.27 ± 0.00
	1/1	<u>16.5</u>	1.92	<u>34.4</u>	18.03	OOM	5.48 ± 0.02	5.00 ± 0.02
SEA-RAFT [15]	1/8	28.3	2.30	81.5	21.51	0.03 ± 0.00	0.09 ± 0.00	0.03 ± 0.00
	1/4	21.1	<u>2.01</u>	54.5	<u>19.18</u>	0.13 ± 0.00	0.18 ± 0.00	0.11 ± 0.00
	1/2	<u>16.8</u>	4.94	<u>39.8</u>	39.83	OOM	0.62 ± 0.00	0.42 ± 0.00
	1/1	17.5	17.65	49.2	107.43	OOM	2.63 ± 0.02	1.78 ± 0.01
SEA-RAFT (Cascaded)	1/8	28.3	2.30	81.5	21.51	0.03 ± 0.00	0.09 ± 0.00	0.03 ± 0.00
	1/4	21.4	<u>1.86</u>	54.1	<u>16.15</u>	0.15 ± 0.00	0.25 ± 0.00	0.12 ± 0.00
	1/2	15.8	1.90	36.8	18.58	OOM	0.72 ± 0.00	0.45 ± 0.00
	1/1	<u>13.3</u>	2.70	<u>31.6</u>	21.53	OOM	2.89 ± 0.02	1.91 ± 0.01

Table 11. Full qualitative evaluation of optical flow estimation methods on the CHARGE-8K dataset. We split the results depending on the evaluation scale and report endpoint-error (EPE), lpx outlier rate, as well as metrics for pixels with large motion. Best result of each metric is highlighted in **bold**, best performing scale for each method is underlined. OOM indicates that the method fails with an out-of-memory error, while n/a - not applicable. We also underline the scale of the row corresponding to the results in Table 1 in the main paper.

```

1  from dataclasses import dataclass
2  import cutlass
3  from cutlass import cute
4  from cutlass.cute.nvgpu import warp
5
6  @dataclass
7  class Gemm:
8      """Helper class to perform tiled GEMM"""
9      thr_copy: cute.TiledCopy
10     rMat1: cute.Tensor
11     rMat2: cute.Tensor
12     rC: cute.Tensor
13     tCrC: cute.Tensor
14     gmem_copy: cute.TiledCopy
15     smem_copy1: cute.TiledCopy
16     smem_copy2: cute.TiledCopy
17     tM1gM1: cute.Tensor
18     tM1sM1: cute.Tensor
19     tM2sM2: cute.Tensor
20     mma: cute.TiledMma
21     tSsMat1: cute.Tensor
22     tSrMat1: cute.Tensor
23     tSsMat2: cute.Tensor
24     tSrMat2: cute.Tensor
25
26     @classmethod
27     def make(cls, gMat1: cute.Tensor, sMat1: cute.Tensor,
28             sMat2: cute.Tensor, sCorr: cute.Tensor,
29             gmem_copy: cute.TiledCopy,
30             mma: cute.TiledMma):
31         tid, _, _ = cute.arch.thread_idx()
32         thr_copy = gmem_copy.get_slice(tid)
33         tM1gM1 = thr_copy.partition_S(gMat1)
34         tM1sM1 = thr_copy.partition_D(sMat1)
35         tM2sM2 = thr_copy.partition_D(sMat2)
36         thr_mma = mma.get_slice(tid)
37         rMat1 = thr_mma.make_fragment_A(thr_mma.partition_A(sMat1))
38         rMat2 = thr_mma.make_fragment_B(thr_mma.partition_B(sMat2))
39         rC = cute.make_fragment(thr_mma.partition_shape_C(sCorr.shape), sCorr.element_type)
40         tCrC = thr_mma.partition_C(sCorr)
41         smem_atom = cute.make_copy_atom(warp.LdMatrix8x8x16bOp(num_matrices=4), gMat1.element_type)
42         smem_copy1 = cute.make_tiled_copy_A(smem_atom, mma)
43         smem_copy2 = cute.make_tiled_copy_B(smem_atom, mma)
44         s1, s2 = smem_copy1.get_slice(tid), smem_copy2.get_slice(tid)
45         tSsMat1, tSrMat1 = s1.partition_S(sMat1), s1.retile(rMat1)
46         tSsMat2, tSrMat2 = s2.partition_S(sMat2), s2.retile(rMat2)
47         return cls(thr_copy, rMat1, rMat2, rC, tCrC, gmem_copy, smem_copy1, smem_copy2,
48                 tM1gM1, tM1sM1, tM2sM2, mma, tSsMat1, tSrMat1, tSsMat2, tSrMat2)
49
50     def compute_tile(self, gMat2: cute.Tensor, n_steps: int):
51         tM1gM2 = self.thr_copy.partition_S(gMat2)
52
53         self.rC.fill(0.0)
54
55         for ko in range(n_steps):
56             cute.copy(self.gmem_copy, self.tM1gM1[None, None, ko], self.tM1sM1[None, None, 0])
57             cute.copy(self.gmem_copy, tM1gM2[None, None, ko], self.tM2sM2[None, None, 0])
58             cute.arch.cp_async_commit_group()
59             cute.arch.cp_async_wait_group(0)
60             cute.arch.sync_threads()
61
62             cute.copy(self.smem_copy1, self.tSsMat1, self.tSrMat1)
63             cute.copy(self.smem_copy2, self.tSsMat2, self.tSrMat2)
64
65             cute.gemm(self.mma, self.rC, self.rMat1, self.rMat2, self.rC)
66
67             cute.arch.sync_threads()
68
69         cute.autovec_copy(self.rC, self.tCrC)
70         cute.arch.sync_threads()

```

Listing 2. CuTe DSL GEMM helper class.

```

1  def _floor(v: cutlass.Numeric):
2      return cute.Int32(cutlass._mlir.dialects.math.floor(v.ir_value()))
3
4
5  def _atomicMin(p: cutlass.Pointer, v: cutlass.Int32):
6      return cutlass._mlir.dialects.nvvm.atomicrmw(
7          cutlass.cutlass_dsl.T.i32(),
8          cutlass._mlir.dialects.nvvm.AtomicOpKind.MIN,
9          p.llvm_ptr, v.ir_value(),
10     )
11
12
13 def _map_index(i: int, h: cutlass.Constexpr, bh: cutlass.Constexpr, w: cutlass.Constexpr, bw: cutlass.Constexpr):
14     bw_i = i % bw
15     bh_i = (i // bw) % bh
16     w_i = (i // (bw * bh)) % w
17     h_i = i // (bw * bh * w)
18     return h_i * bh + bh_i, w_i * bw + bw_i
19
20
21 @cute.jit
22 def compute_coefficients(coords, bi, by, tidx, scale, h_block, cs, r, dtype):
23     ch, cw = coords.shape[2:]
24     cy, cx = _map_index(by * h_block + tidx, cute.ceil_div(ch, cs), cs, cute.ceil_div(cw, cs), cs)
25     is_valid_row = (cy < ch) & (cx < cw)
26     x_floor, y_floor = cute.Int32(0), cute.Int32(0)
27     nw, ne, sw, se = dtype(0.0), dtype(0.0), dtype(0.0), dtype(0.0)
28     if is_valid_row:
29         x_start = coords[bi, 0, cy, cx] * scale - r // 2
30         y_start = coords[bi, 1, cy, cx] * scale - r // 2
31         x_floor, y_floor = _floor(x_start), _floor(y_start)
32         xc, yc = x_start - x_floor, y_start - y_floor
33         nw, ne = (1 - xc) * (1 - yc), xc * (1 - yc)
34         sw, se = (1 - xc) * yc, xc * yc
35     return (cy, cx), is_valid_row, (x_floor, y_floor), (nw, ne, sw, se)
36
37
38 @cute.jit
39 def sample_block(sCorr, tidx, rCorr, mat2x, mat2y,
40                 x_floor, y_floor, nw, ne, sw, se,
41                 output, bi, cy, cx, cs, r, dtype):
42     cute.autovec_copy(sCorr[tidx, None], rCorr)
43
44     # Only iterate over parts that are in the current tile
45     rx_start = max(mat2x * cs - x_floor - 1, 0)
46     rx_end = min(rx_start + cs + 1, r)
47     ry_start = max(mat2y * cs - y_floor - 1, 0)
48     ry_end = min(ry_start + cs + 1, r)
49
50     for ry in range(ry_start, ry_end):
51         y = y_floor - mat2y * cs + ry
52         for rx in range(rx_start, rx_end):
53             x = x_floor - mat2x * cs + rx
54             value = dtype(0.0)
55             if (y >= 0) & (y < cs):
56                 if (x >= 0) & (x < cs):
57                     value += rCorr[y * cs + x] * nw
58                 if (x + 1 >= 0) & (x + 1 < cs):
59                     value += rCorr[y * cs + x + 1] * ne
60             if (y + 1 >= 0) & (y + 1 < cs):
61                 if (x >= 0) & (x < cs):
62                     value += rCorr[(y + 1) * cs + x] * sw
63                 if (x + 1 >= 0) & (x + 1 < cs):
64                     value += rCorr[(y + 1) * cs + x + 1] * se
65             output[bi, ry + rx * r, cy, cx] += value

```

Listing 3. Kernel helper functions.

```

1  @cute.kernel
2  def kernel(
3      coords: cute.Tensor, # [B 2 H W]
4      mat1: cute.Tensor, # patch-major [B (h w bh bw) C]
5      mat2: cute.Tensor, # patch-major [B h w (bh bw) C]
6      output: cute.Tensor, # output [B (2r+1)^2 H H]
7      sMat1_layout: cute.ComposedLayout, # with Swizzle<2,3,3> o 0 o (8,32):(32,1)
8      sMat2_layout: cute.ComposedLayout, # with Swizzle<2,3,3> o 0 o (8,32):(32,1)
9      corr_layout: cute.Layout, # [bh, bw]
10     gmem_tiled_copy_mat: cute.TiledCopy, # with cute.nvgpu.cpasync.CopyG2SOp
11     tiled_mma: cute.TiledMma, # MmaF16BF16Op
12     coord_scaler: cutlass.Numeric, # float, scaler for coords on lower resolutions
13     cell_size: cutlass.Constexpr, radius: cutlass.Constexpr, k_block: cutlass.Constexpr,
14 ):
15     (by, bi, _) , (tidx, _, _) = cute.arch.block_idx(), cute.arch.thread_idx()
16     h_block, w_block = corr_layout.shape[:2]
17     dtype = output.element_type
18     _, targetH, targetW, _, channels = mat2.shape
19
20     smem = cutlass.utils.SmemAllocator()
21     next_block = smem.allocate_tensor(cutlass.Int32, cute.make_layout((1,)))
22     sMat1 = smem.allocate_tensor(mat1.element_type, sMat1_layout, byte_alignment=16)
23     sMat2 = smem.allocate_tensor(mat2.element_type, sMat2_layout, byte_alignment=16)
24     sCorr = cute.make_tensor(cute.recast_ptr(sMat1.iterator, dtype=dtype), corr_layout)
25
26     gMat1 = cute.local_tile(mat1[bi, None, None], (h_block, channels), (by, 0))
27     gMat1 = cute.make_tensor(gMat1.iterator.align(16), gMat1.layout)
28     gemm = Gemm.make(gMat1, sMat1, sMat2, sCorr, gmem_tiled_copy_mat, tiled_mma)
29     rCorr = cute.make_fragment(cute.make_layout((w_block,)), dtype)
30     cute.arch.sync_threads()
31
32     (cy, cx), is_valid_row, (x_floor, y_floor), (nw, ne, sw, se) = compute_coefficients(
33         coords, bi, by, tidx, coord_scaler, h_block, cell_size, radius, dtype)
34
35     # Find the blocks this thread requires and prepare the vote
36     rpl = radius + 1
37     nr = rpl // cell_size + 2 # Maximum number of blocks this thread can request
38     blocks = cute.make_fragment(cute.make_layout((nr * nr + 1,)), cutlass.Int32)
39     max_block = targetH * targetW
40     blocks.fill(max_block)
41     block_i = 0
42     if is_valid_row:
43         for i in range(nr * nr):
44             rx, ry = i % nr, (i // nr) % nr
45             x = (x_floor + min(rx * cell_size, rpl - 1)) // cell_size
46             y = (y_floor + min(ry * cell_size, rpl - 1)) // cell_size
47             if (x >= 0) & (x < targetW) & (y >= 0) & (y < targetH):
48                 val = y * targetW + x
49                 if block_i == 0 or blocks[block_i - 1] < val:
50                     blocks[block_i] = val
51                     block_i += 1
52
53     if tidx == 0:
54         next_block[0] = max_block
55     cute.arch.sync_threads()
56
57     block_i = 0
58     next_local = 0
59     while next_local < max_block:
60         # Vote for the next block
61         if is_valid_row and blocks[block_i] < max_block:
62             _atomicMin(next_block.iterator, blocks[block_i])
63             cute.arch.sync_threads()
64             next_local = next_block[0]
65
66         if next_local < max_block:
67             mat2y, mat2x = next_local // targetW, next_local % targetW
68             gMat2 = cute.local_tile(mat2[bi, mat2y, mat2x, None, None], (w_block, channels), (0, 0))
69             gMat2 = cute.make_tensor(gMat2.iterator.align(16), gMat2.layout)
70             cute.arch.sync_threads()
71             if tidx == 0: # Resetting the vote
72                 next_block[0] = max_block
73
74             gemm.compute_tile(gMat2, n_steps=(channels + k_block - 1) // k_block)
75
76         if next_local == blocks[block_i]: # If this thread requested the block, sample it and advance to the next
77             sample_block(sCorr, tidx, rCorr, mat2x, mat2y, x_floor, y_floor,
78                 nw, ne, sw, se, output, bi, cy, cx, cell_size, radius, dtype)
79             block_i += 1

```

Listing 4. Kernel implementation using CuTe DSL, reformatted for space constraints.