

# Structure-grounded Training Strategies Aid Generalization in Stereo Matching

## Supplementary Material

### A. Depth Score Calculation in Figure 1 of Main Text

For a fixed (model, dataset) and a given training strategy  $s$ , we combine endpoint error (EPE; lower is better) and the  $< 3\text{px}$  accuracy (higher is better) via a normalized composite score:

$$\text{DepthScore}_s = \alpha \left( 1 - \frac{EPE_s - EPE_{\min}}{EPE_{\max} - EPE_{\min}} \right) + (1 - \alpha) \frac{\text{Acc}_s - \text{Acc}_{\min}}{\text{Acc}_{\max} - \text{Acc}_{\min}}, \quad (1)$$

where  $\text{Acc}$  denotes  $< 3\text{px}(\%)$ , and the  $\min / \max$  operators are taken *across strategies within the same (model, dataset) pair*. We use  $\alpha = \frac{1}{2}$  by default (equal weighting). When a denominator is zero (all strategies identical on that metric), we set the corresponding normalized term to 0. For visualization in radar plots, we apply a monotonic remapping  $s' = \varepsilon + (1 - \varepsilon) s$  with  $\varepsilon = 0.15$  to avoid polygon degeneracy at the origin; this does *not* affect ranking. For bar plots, we report the dataset-averaged  $\text{DepthScore}$  per model and annotate  $\Delta$  relative to the model’s Baseline strategy.

### B. Theoretical Analysis of Geometry-oriented Data Augmentation

#### B.1. Pixel-Level Correspondence as Bipartite Matching

We represent the task of matching pixels between two views as finding a maximum-weight matching in a complete bipartite graph. Let  $n$  be the number of feature points (pixels) we extract in each image. Denote by  $\mathcal{I} = \{a_1, \dots, a_n\}$  the set of source-image pixels and by  $\mathcal{J} = \{b_1, \dots, b_n\}$  the set of target-image pixels. We apply a feature extractor  $f(\cdot)$  (e.g., a local descriptor) to each pixel, and a similarity function  $s(\cdot, \cdot)$  (e.g., dot-product or cosine similarity) to pairs of descriptors. The resulting affinity between  $a_i \in \mathcal{I}$  and  $b_j \in \mathcal{J}$  can be denoted by

$$w_{ij} = s(f(a_i), f(b_j)), \quad (2)$$

and we collect these into the weight matrix  $W = [w_{ij}] \in \mathbb{R}^{n \times n}$ .

We note that a matching is a bijection  $\pi : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  that assigns each source pixel  $a_i$  to exactly one target pixel  $b_{\pi(i)}$ . Thus, we denote this permutation by  $\pi$  and score it by the total affinity

$$S(\pi) = \sum_{i=1}^n w_{i, \pi(i)}. \quad (3)$$

The optimal correspondence  $\pi^*$  maximizes this total score:

$$\pi^* = \arg \max_{\pi} S(\pi). \quad (4)$$

To measure how clearly  $\pi^*$  stands out from other candidates, we define the *matching gap* as

$$g = S(\pi^*) - \max_{\pi \neq \pi^*} S(\pi), \quad (5)$$

where the second term is the highest score among all permutations other than  $\pi^*$ . A positive gap  $g > 0$  indicates that  $\pi^*$  is the unique optimum.

#### B.2. Texture Augmentation as Random Perturbation

In regions with weak texture, many pairs  $(a_i, b_j)$  yield similar affinities, leading to small  $g$  and ambiguous matching. We increase surface texture by overlaying random line patterns on the object, which perturbs each affinity by an independent random variable. Let  $\epsilon_{ij}$  be zero-mean noise with variance  $\sigma^2$ , modeling the effect of texture:

$$\tilde{w}_{ij} = w_{ij} + \epsilon_{ij}, \quad \epsilon_{ij} \sim \mathcal{D}(0, \sigma^2), \quad (6)$$

and write  $\tilde{W} = [\tilde{w}_{ij}]$ . Under this perturbation, the score of any matching  $\pi$  becomes

$$\begin{aligned} S'(\pi) &= \sum_{i=1}^n \tilde{w}_{i, \pi(i)} \\ &= \sum_{i=1}^n (w_{i, \pi(i)} + \epsilon_{i, \pi(i)}) \\ &= S(\pi) + E(\pi), \\ E(\pi) &= \sum_{i=1}^n \epsilon_{i, \pi(i)}, \end{aligned} \quad (7)$$

where  $E(\pi)$  is the total noise contribution for matching  $\pi$ . For each non-optimal permutation  $\pi \neq \pi^*$ , define

$$\Delta(\pi) = S(\pi^*) - S(\pi), \quad X(\pi) = E(\pi^*) - E(\pi). \quad (8)$$

#### Hypothesis: A larger matching gap improves pixel-level matching.

Increasing the matching gap  $g$  through texture augmentation can contribute to reducing matching difficulties in three folds:

**1. Robustness to descriptor errors.** In practice, computed affinities  $w_{ij}$  suffer bounded errors from descriptor quantization, sensor noise, or illumination changes. Suppose each  $w_{ij}$  is perturbed by at most  $\delta$ , yielding observed weights  $\hat{w}_{ij} \in [w_{ij} - \delta, w_{ij} + \delta]$ . Then for the true matching  $\pi^*$ ,

$$\hat{S}(\pi^*) \geq \sum_{i=1}^n (w_{i,\pi^*(i)} - \delta) = S(\pi^*) - n\delta, \quad (9)$$

and for any competitor  $\pi$ ,

$$\hat{S}(\pi) \leq \sum_{i=1}^n (w_{i,\pi(i)} + \delta) = S(\pi) + n\delta. \quad (10)$$

Thus the perturbed gap satisfies

$$\hat{g} = \hat{S}(\pi^*) - \max_{\pi \neq \pi^*} \hat{S}(\pi) \geq g - 2n\delta. \quad (11)$$

If the original gap  $g$  exceeds  $2n\delta$ , the optimum remains  $\pi^*$  despite these errors. Therefore, a larger  $g$  confers greater *tolerance to deterministic weight perturbations*.

**2. Enhanced pruning in matching algorithms.** Combinatorial solvers (e.g. Hungarian method, auction algorithm) build partial matchings of size  $p$  with score  $S_p$  and estimate the maximum possible completion by adding at most  $(n - p) \max_{i,j} w_{ij}$ . A branch is pruned if

$$S_p + (n - p) \max_{i,j} w_{ij} < S(\pi^*) - g. \quad (12)$$

Larger  $g$  allows for *more aggressive pruning* of suboptimal branches, reducing the algorithm’s search space and run-time.

**3. Exponential decay of misassignment probability.** Under random perturbation,  $\pi^*$  remains optimal if for every  $\pi \neq \pi^*$ ,

$$\Delta(\pi) + X(\pi) > 0. \quad (13)$$

Since  $X(\pi)$  has zero mean and variance  $\text{Var}[X(\pi)] = O(n\sigma^2)$ , Hoeffding’s inequality yields

$$\Pr[X(\pi) < -\Delta(\pi)] \leq \exp\left(-\frac{\Delta(\pi)^2}{2n\sigma^2}\right). \quad (14)$$

A union bound over all  $n! - 1$  competitors gives

$$\Pr[\pi^* \text{ remains optimal}] \geq 1 - (n! - 1) \exp\left(-\frac{g^2}{2n\sigma^2}\right). \quad (15)$$

Thus the probability of misassignment decays *exponentially* in  $g^2/(n\sigma^2)$ , showing that a larger gap substantially *reduces the risk of incorrect matches* under random texture noise.

## C. Additional Results of Geometry-oriented Augmentation

Our augmentation preserves stereo consistency while introducing *two valid correspondences* for pixels within the perturbed regions: (i) the original correspondence  $\mathbf{p}_r^o$ , and (ii) the synthetic correspondence  $\mathbf{p}_r^a$  induced by the warped disparity. As illustrated in Fig. S1, a pixel in the left image may simultaneously admit both  $\mathbf{p}_r^o$  and  $\mathbf{p}_r^a$  on the right view. The original correspondence  $\mathbf{p}_r^o$  aligns with semantic priors and reflects the original disparity, while the synthetic correspondence  $\mathbf{p}_r^a$  stems from the augmented disparity. In training, we supervise using  $\mathbf{p}_r^a$ , encouraging the network to prioritize fine-grained geometric cues over semantic bias. This is particularly beneficial in planar textured regions, where semantic guidance alone may hallucinate depth from texture, while geometric cues enforce the correct planar depth.

Further examples of randomly generated geometry-oriented augmentations are provided in Figures S2 and S3.

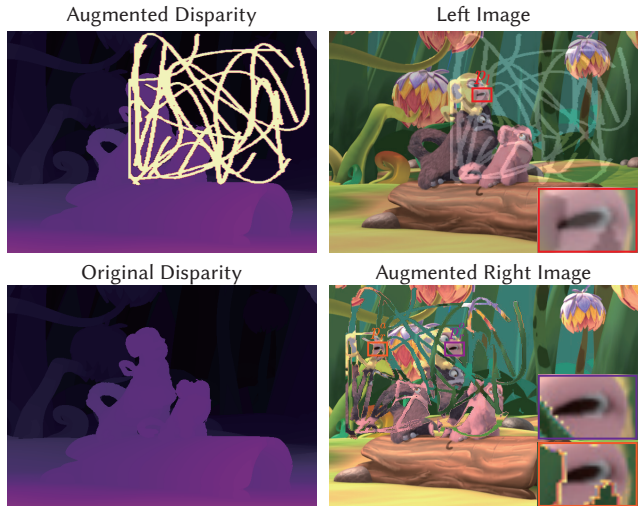


Figure S1. Illustration of dual correspondences induced by augmentation. A selected pixel  $\mathbf{p}_l$  in the left image corresponds to the original match  $\mathbf{p}_r^o$  in the right image (refer to the original disparity) and to the synthetic match  $\mathbf{p}_r^a$  created by the augmented disparity.

## D. Additional Training on Released Models

We further apply our training strategies to the official SceneFlow pretrained models released by RAFT-Stereo, IGEV-Stereo, and Selective-IGEV. Starting from the published checkpoints, we continue training on SceneFlow using our proposed training strategies, with identical hyperparameters across models (batch size 4, 100k steps, learning rate  $1 \times 10^{-5}$ ; for Selective-IGEV, batch size 3 and 135k steps due to memory). For fair comparison, the *Baseline* results in Tables S1a–S1c are obtained by continuing training for the same number of steps using the official training scheme



Figure S2. Examples of geometry-oriented data augmentation with Ribbon mode.

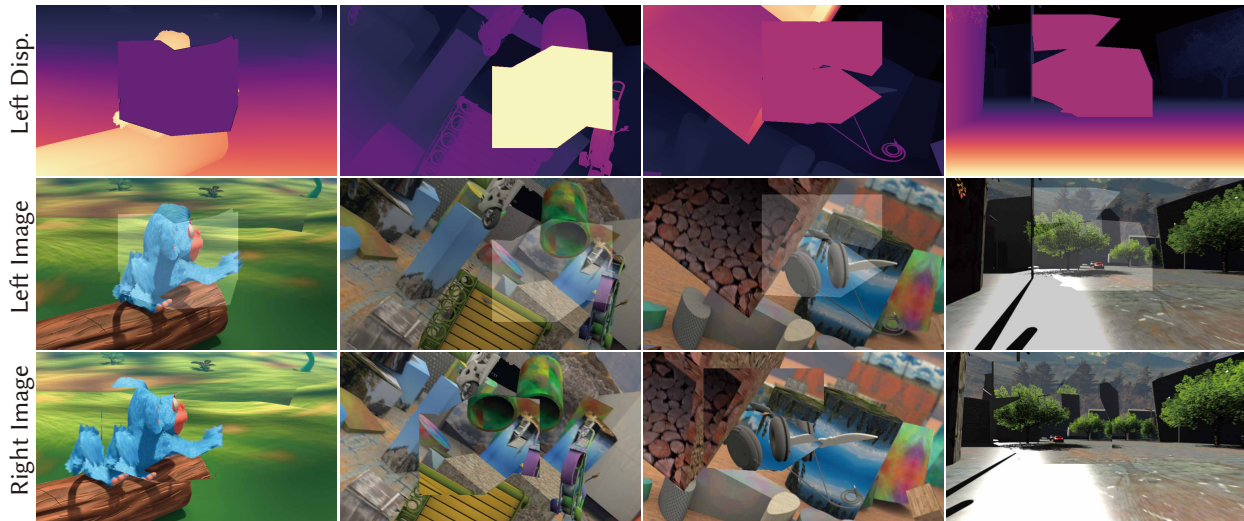


Figure S3. Examples of geometry-oriented data augmentation with Blob mode.

from each paper. We then evaluate the resulting models on KITTI 2012/2015, Middlebury, and ETH3D.

As shown in Tables S1a–S1c, our strategies consistently improve or maintain accuracy compared with the baselines. These results demonstrate that our training paradigm can be applied to existing SceneFlow-trained models and leads to measurable performance improvements across diverse architectures.

## E. Additional Qualitative Results

Figures S4, S6, S5, and S7 extend the visual result of Fig. 8 in the main paper to all four models discussed there (RAFT-Stereo, IGEV-Stereo, Selective-IGEV, and DLNR). For each model, we show disparity predictions from the

baseline and from our training strategies side by side.

Table S1. Results obtained by continuing training from the official SceneFlow-pretrained models of RAFT-Stereo, IGEV-Stereo, and Selective-IGEV. **Bold** = best, underline = second best, **yellow background** = improved over Baseline.

(a) RAFT-Stereo

Method	KITTI-12		KITTI-15		Middlebury		ETH3D	
	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)
Baseline	0.929	95.1	1.173	94.1	<b>1.015</b>	<b>94.1</b>	0.275	98.7
GeomAug	<u>0.923</u>	<u>95.2</u>	<u>1.150</u>	<u>94.3</u>	<u>1.099</u>	93.2	0.286	98.8
AuxTask	<b>0.905</b>	<b>95.3</b>	<b>1.136</b>	<b>94.4</b>	1.172	93.6	<b>0.256</b>	<b>99.0</b>
UpdReg	0.966	94.8	1.195	94.0	1.275	93.0	0.270	98.8
Joint	<u>0.908</u>	<b>95.3</b>	1.154	94.2	1.332	<u>93.9</u>	<u>0.267</u>	<u>98.9</u>

(b) IGEV-Stereo

Method	KITTI-12		KITTI-15		Middlebury		ETH3D	
	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)
Baseline	1.076	94.2	1.199	<u>94.0</u>	<b>0.891</b>	94.0	0.529	<u>98.3</u>
GeomAug	<u>1.048</u>	<u>94.4</u>	<u>1.192</u>	93.9	1.006	<b>94.6</b>	<u>0.335</u>	<b>98.5</b>
AuxTask	1.109	94.0	1.202	93.9	1.044	94.1	<b>0.318</b>	<b>98.5</b>
UpdReg	1.150	93.5	1.252	93.5	<u>0.989</u>	<u>94.5</u>	0.540	98.2
Joint	<b>1.002</b>	<b>94.8</b>	<b>1.182</b>	<b>94.2</b>	1.033	<u>94.5</u>	0.383	<u>98.3</u>

(c) Selective-IGEV

Method	KITTI-12		KITTI-15		Middlebury		ETH3D	
	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)	EPE (px)	< 3 px (%)
Baseline	1.069	<u>93.9</u>	1.265	<u>93.9</u>	0.879	<b>95.6</b>	0.354	<u>98.3</u>
GeomAug	<u>1.052</u>	<b>94.2</b>	1.253	<b>94.0</b>	<b>0.816</b>	<u>94.9</u>	0.373	<u>98.3</u>
AuxTask	1.095	93.7	<u>1.241</u>	<u>93.9</u>	0.844	94.5	0.384	98.2
UpdReg	1.076	93.8	1.255	93.8	<u>0.860</u>	<b>95.6</b>	<u>0.353</u>	<u>98.3</u>
Joint	<b>1.042</b>	<b>94.2</b>	<b>1.220</b>	<b>94.0</b>	<u>0.820</u>	<b>95.6</b>	<b>0.340</b>	<b>98.5</b>

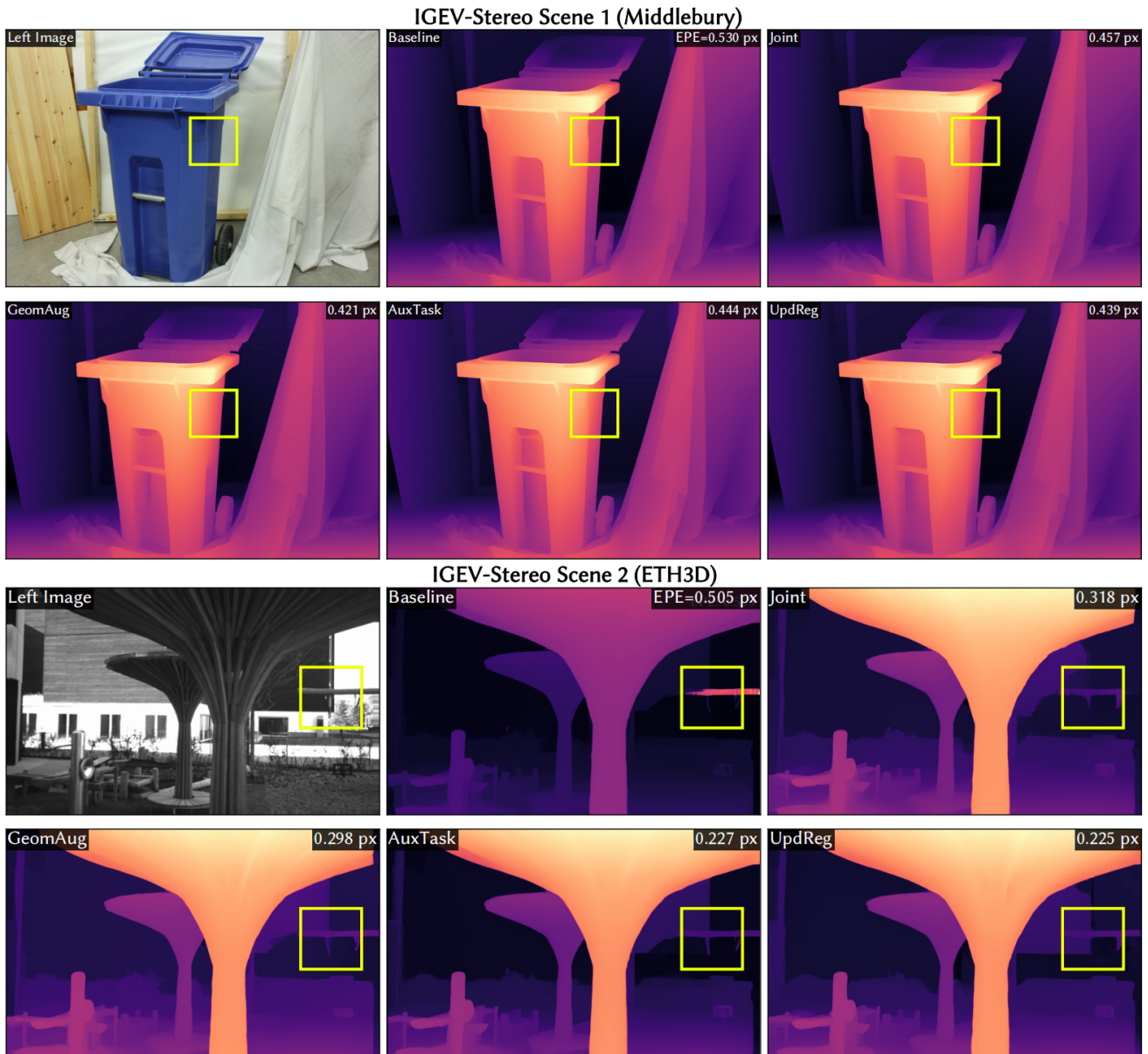


Figure S4. Qualitative comparison on two scenes (from KITTI 2012 and Middlebury) based on IGEV-Stereo. Two scenes are shown to compare our four training strategies with the baseline, original IGEV-Stereo. The highlighted boxes emphasize difficult areas (e.g., textureless regions, occlusions, and thin structures) where disparities diverge across methods. Per-panel EPE (top-right) in disparity map quantifies the error and complements the visual comparison, lower which is better. (top row: Left Image, Baseline, Joint; bottom row: GeomAug, AuxTask, UpdReg).

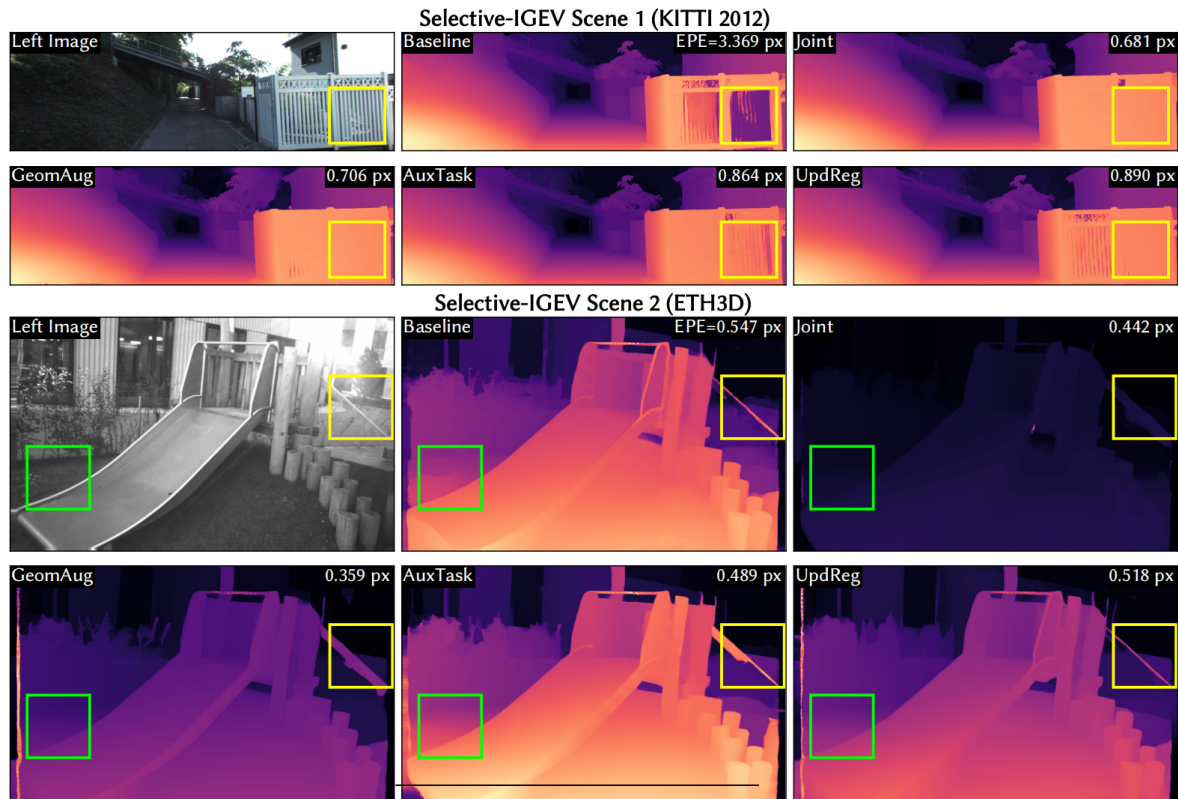


Figure S5. Qualitative comparison on two scenes (from KITTI 2012 and ETH3D) based on Selective-IGEV. Two scenes are shown to compare our four training strategies with the baseline, original Selective-IGEV.

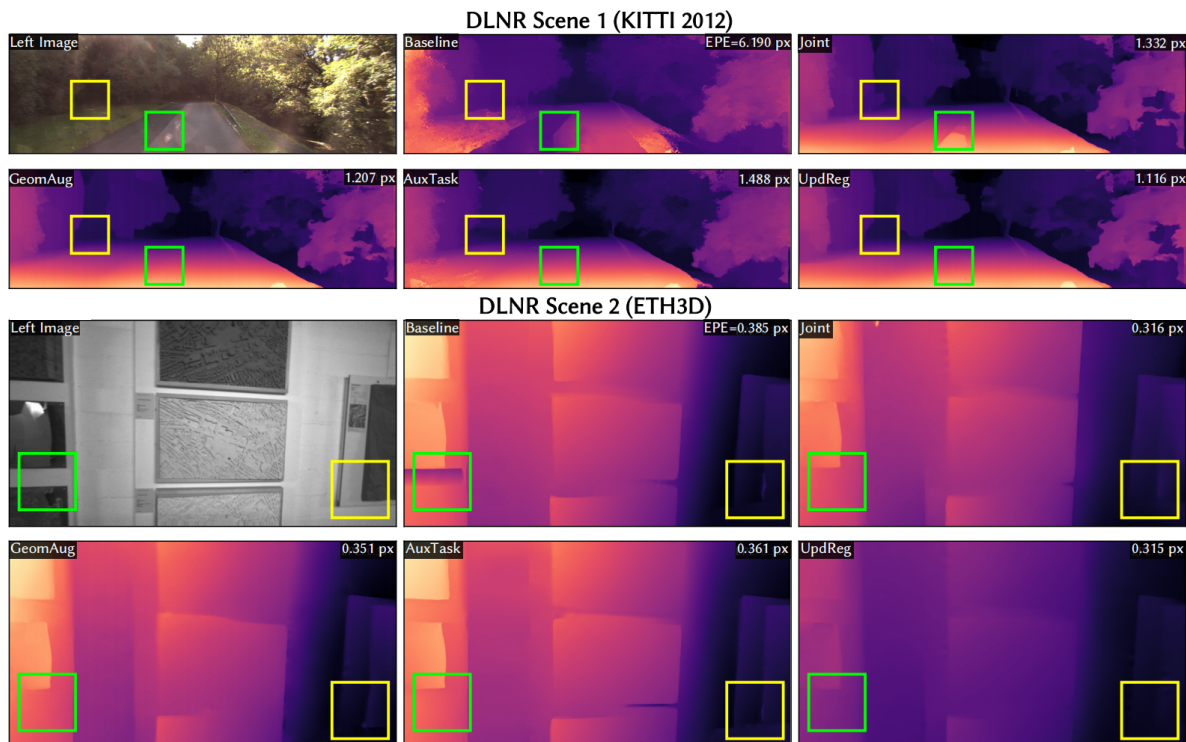


Figure S6. Qualitative comparison on two scenes (from KITTI 2012 and ETH3D) based on DLNR model. Two scenes are shown to compare our four training strategies with the baseline, original DLNR.

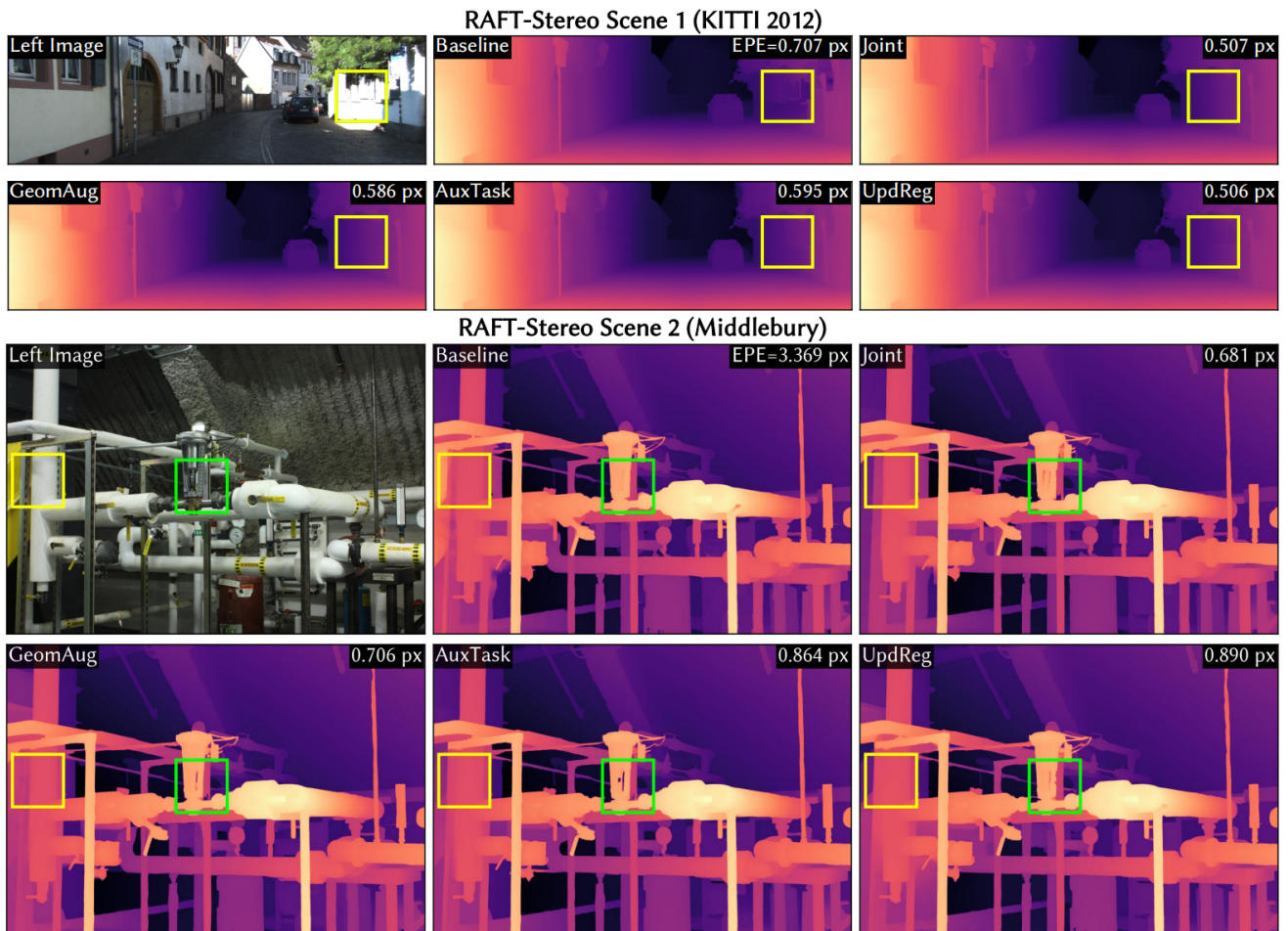


Figure S7. Qualitative comparison on two scenes (from KITTI 2012 and Middlebury) based on RAFT-Stereo. Two scenes are shown to compare our four training strategies with the baseline, original RAFT-Stereo. (top row: Left Image, Baseline, Joint; bottom row: GeomAug, AuxTask, UpdReg).