

Multi-View Stereo Problems

- Depth map-based MVS algorithms estimate the reference view depth maps using multiple RGB inputs (Reference + Source views)
- A consistent scene requires geometric consistency of depth estimates across multiple views

Two broader approaches are undertaken to ensure geometric consistency in estimated depth maps:

- Repeated application of geometric constraints during the depth estimation process → Traditional MVS Algorithms
- Geometric constraints applied as a post-processing step → Learning-based MVS Algorithms

GC-MVSNet is a learning-based algorithm with geometric constraints applied during the learning process.

Learning-Based MVS Algorithms

A learning-based MVS method:

- Extracts multi-level features using CNNs
- Creates a matching 3D cost volume using features
- Regularizes cost volume using 3D-CNN
- Filters geometrically consistent points to generate 3D point-cloud

They only use Geometric Constraints as a post-processing step for filtering multi-view consistent points. This leads to:

- Limited geometric cues during the learning process
- Requiring more training iterations to learn to reason about geometry

Hypothesis

GC-MVSNet:

- Explicitly models cross-view geometric constraints during learning
- Penalizes geometrically inconsistent estimates during learning

With such explicit geometric constraint modeling, GC-MVSNet should:

- Develop a better understanding of multi-view geometry → Improved quantitative results
- Learn quickly to reason about scene geometry → Require less training iterations

Forward-Backward-Reprojection

Inputs: $D_0, c_0, D_i^{gt}, c_i^{gt}$
Output: $D_{P_0}^i, P_0^i$

$K_R, E_R \leftarrow c_0; K_S, E_S \leftarrow c_i^{gt}$
 $D_{(R \rightarrow S)} \leftarrow K_S \cdot E_S \cdot E_R^{-1} \cdot K_R^{-1} \cdot D_0$ ▷ Project
 $X_{D_{(R \rightarrow S)}}, Y_{D_{(R \rightarrow S)}} \leftarrow D_{(R \rightarrow S)}$
 $D_{S_{remap}} \leftarrow REMAP(D_i^{gt}, X_{D_{(R \rightarrow S)}}, Y_{D_{(R \rightarrow S)}})$ ▷ Remap
 $D_{P_0}^i \leftarrow K_R \cdot E_R \cdot E_S^{-1} \cdot K_S^{-1} \cdot D_{S_{remap}}$ ▷ Back project
 $P_0^i \leftarrow (X_{D_{P_0}^i}, Y_{D_{P_0}^i})$

Other Modifications

Two additional modifications were to stabilize the model's performance.

- Keep the feature-extraction network as a Feature Pyramid Network, replaced the regular conv-layers with deformable conv-layers
- Replaced BatchNorm-layers with GroupNorm-layers as BatchNorm is not well suited for small batch-sizes

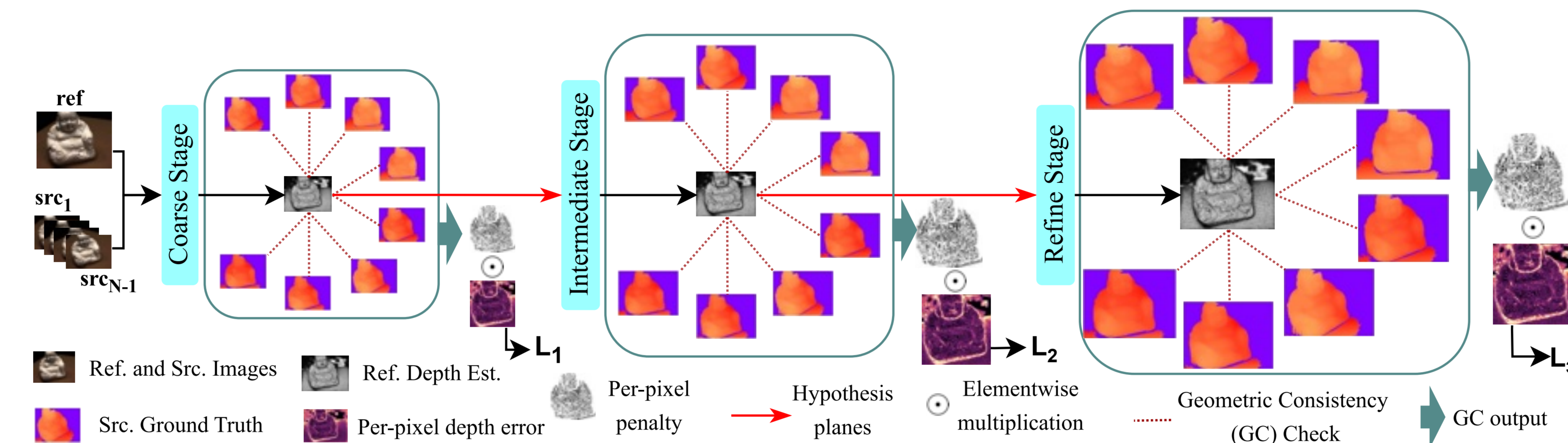
Reconstructed Scene Point Clouds



Method

Geometric-Consistency (GC) Module:

- Applied at the end of each stage to check cross-view consistency of the reference view depth maps
- Generates penalty for geometrically inconsistent estimates for each stage



Geometric-Consistency Module

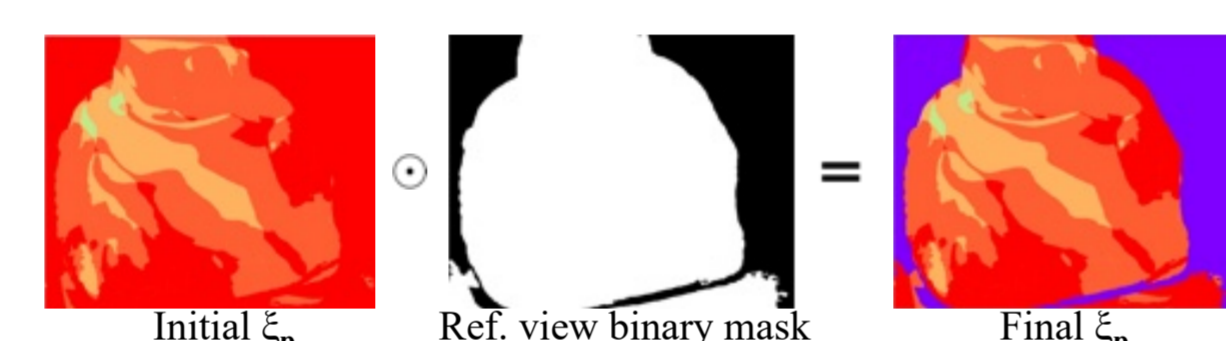
Complete GC-Algorithm

Initialize Mask-Sum $\rightarrow 0$

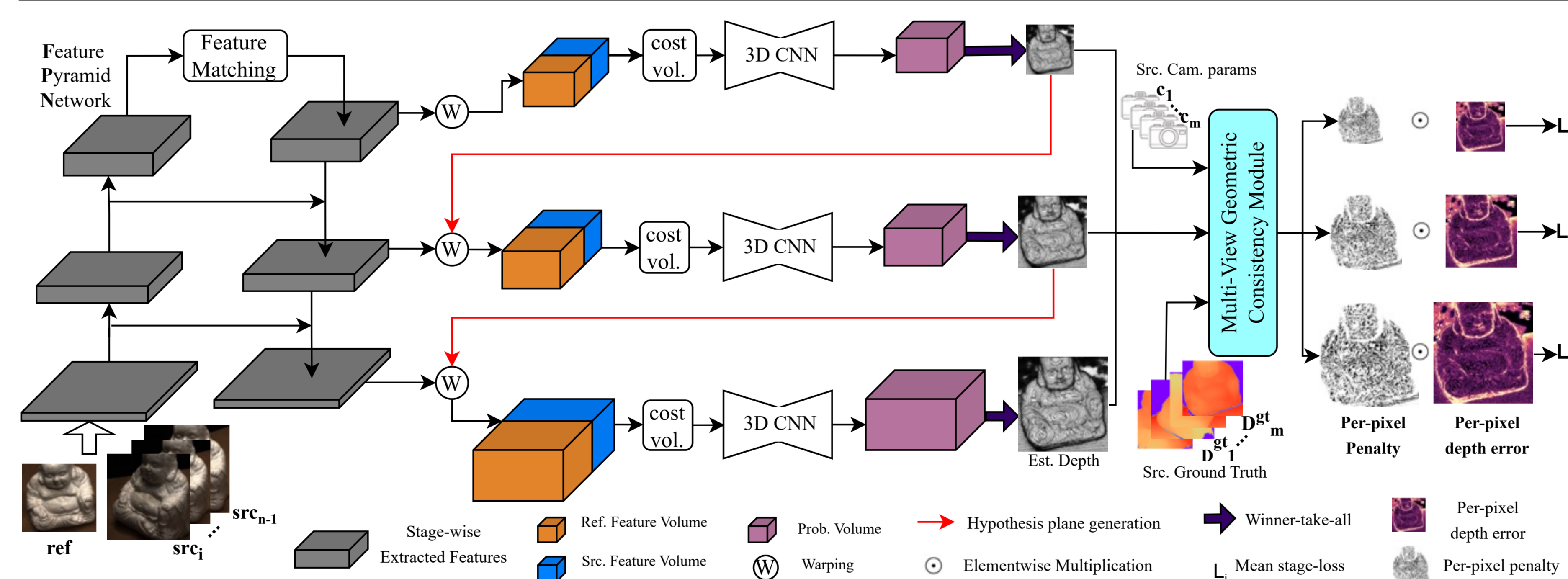
For each Src. depth map:

- forward-backward-reprojection to get PDE and RDD
 - $PDE \leftarrow \|P_0 - P_0^i\|_2$
 - $RDD \leftarrow 1/d_0 \|D_{P_0}^i - D_0\|_1$
- Select geometrically inconsistent pixels
 - $PDE_{mask} > D_{pixel}$
 - $RDD_{mask} > D_{depth}$
- Combine inconsistent pixels from both masks
 - Logical-OR (PDE_{mask}, RDD_{mask})
- Current-Mask \leftarrow Assign penalty to each pixel
 - Inconsistent pixels $\rightarrow 1$
 - All other pixels $\rightarrow 0$
- Add Current-Mask to initial Mask-Sum

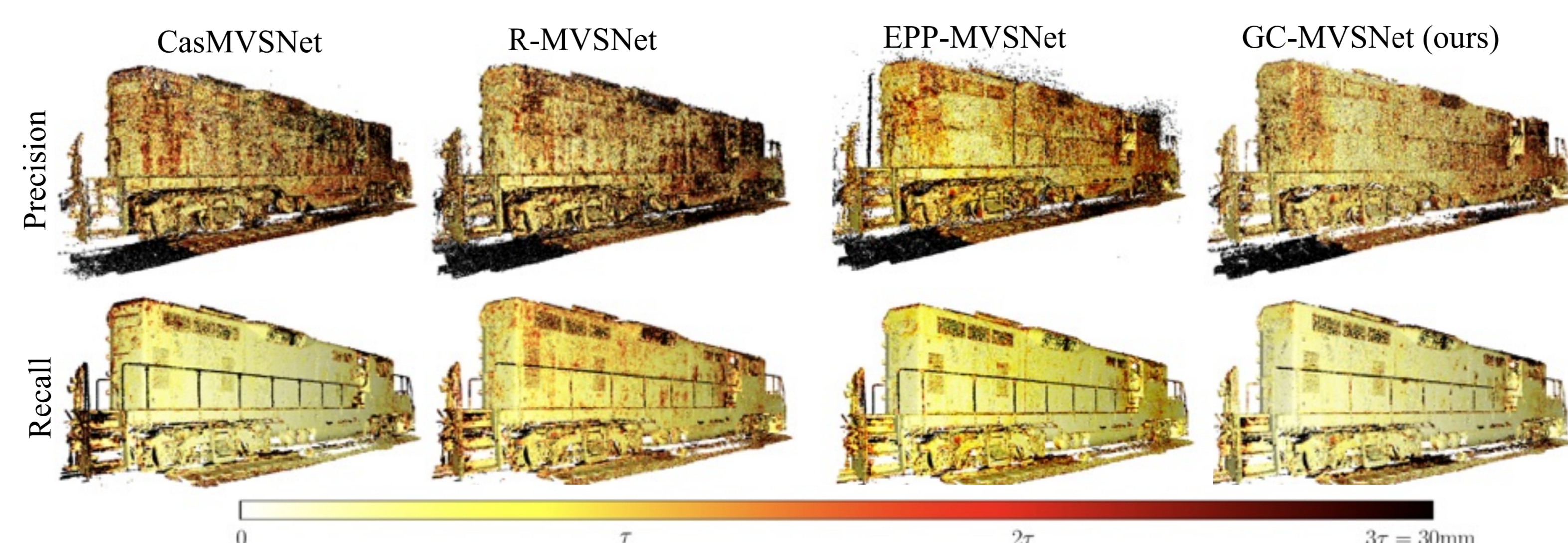
Geometric penalty (ξ_p) \leftarrow average Mask-Sum
Apply reference view binary mask to generate final ξ_p



GC-MVSNet Architecture



Error Plot - Train (Tanks & Temples)



Quantitative Results

Our method achieves state-of-the-art results on two datasets: DTU and BlendedMVS

	Method	Acc ↓	Comp ↓	Overall ↓
Traditional	Furu [9]	0.613	0.941	0.777
	Tola [36]	0.342	1.190	0.766
	Gipuma [10]	0.283	0.873	0.578
	COLMAP [33]	0.400	0.664	0.532
Learning-based	SurfaceNet [16]	0.450	1.040	0.745
	MVSNet [48]	0.396	0.527	0.462
	P-MVSNet [25]	0.406	0.434	0.420
	R-MVSNet [49]	0.383	0.452	0.417
	Point-MVSNet [2]	0.342	0.411	0.376
	CasMVSNet [12]	0.325	0.385	0.355
	CVP-MVSNet [47]	0.296	0.406	0.351
	UCS-Net [3]	0.338	0.349	0.344
	AA-RMVSNet [41]	0.376	0.339	0.357
	UniMVSNet [30]	0.352	0.278	0.315
	TransMVSNet [6]	0.321	0.289	0.305
	GBi-Net* [28]	0.312	0.293	0.303
	MVSTER [39]	0.350	0.276	0.313
	GC-MVSNet (ours)	0.330	0.260	0.295
GBi-Net [28]	0.315	0.262	0.289	
GC-MVSNet (ours)	0.323	0.255	0.289	

DTU Dataset

Method	EPE ↓	ϵ_1 ↓	ϵ_3 ↓
MVSNet [48]	1.49	21.98	8.32
CasMVSNet [12]	1.43	19.01	9.77
CVP-MVSNet [47]	1.90	19.73	10.24
Vis-MVSNet [54]	1.47	15.14	5.13
EPP-MVSNet [26]	1.17	12.66	6.20
TransMVSNet [6]	0.73	8.32	3.62
GC-MVSNet (ours)	0.48	0.89	0.97

BlendedMVS Dataset

GC: A Plug-in Module

The GC module is designed as a plug-in module

- Plug into any depth map-based MVS method
- Retraining the network with GC module provides:
 - Improved quantitative results
 - Requires less training iterations to achieve optimal performance

We demonstrate this on two different methods:

- CasMVSNet and TransMVSNet

Methods	Loss	Other	GC	Overall ↓	Epoch
CasMVSNet [2]	L_1	×	×	0.355	16
	L_1	✓	×	0.357	16
	L_1	×	✓	0.335	11
TransMVSNet [1]	FL	×	×	0.305	16
	FL	✓	×	0.322	16
	FL	×	✓	0.303	8

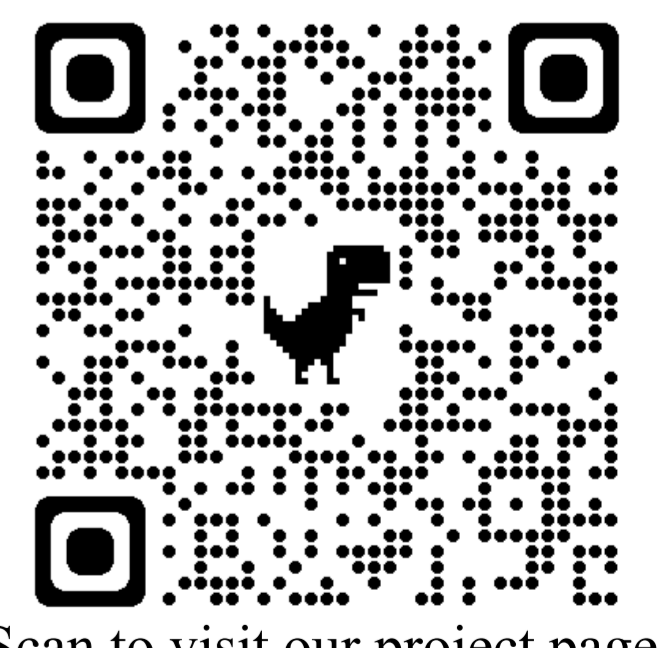
Table 1. GC-module as a plug-in in TransMVSNet and CasMVSNet

References

- Yikang Ding, Wentao Yuan, Qingtian Zhu, Haotian Zhang, Xiangyue Liu, Yuanjiang Wang, and Xiao Liu. Transmvsnet: Global context-aware multi-view stereo network with transformers. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022.
- Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition, 2020.

Connect with us

Provide feedback:



Scan to visit our project page