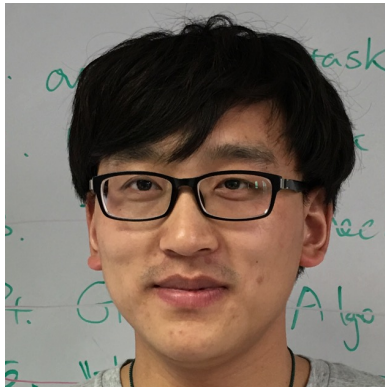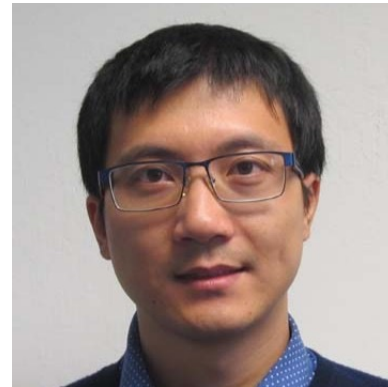# RIAV-MVS: Recurrent-Indexing an Asymmetric Volume for Multi-View Stereo

Changjiang Cai    Pan Ji    Qingan Yan    Yi Xu

OPPO US Research Center, InnoPeak Technology, Inc.
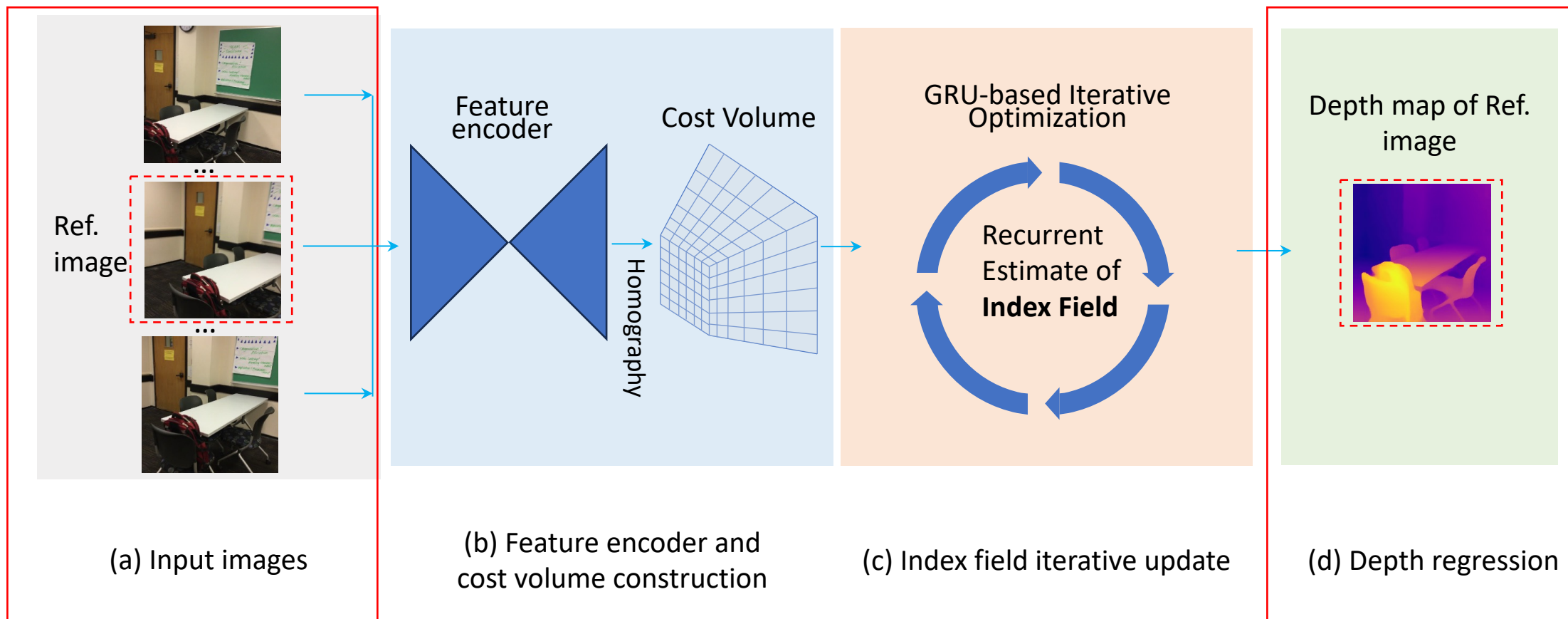
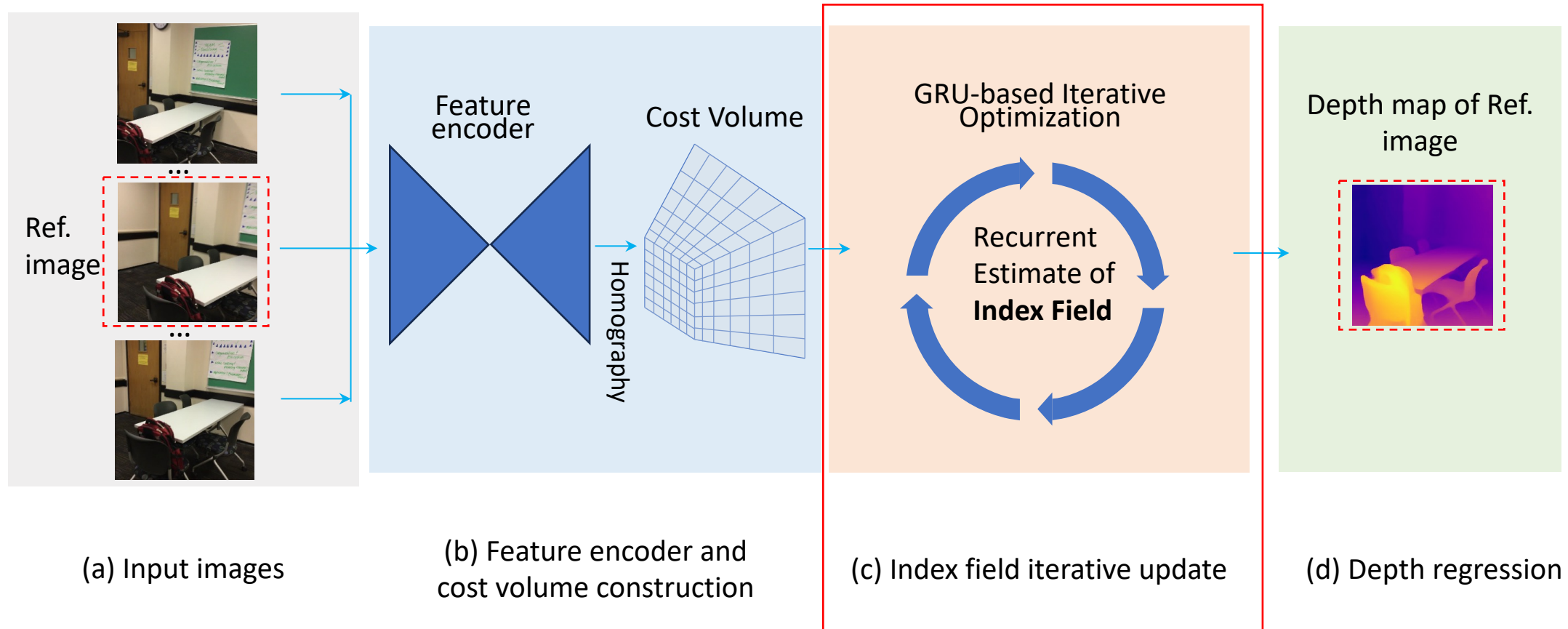Poster Tag: TUE-AM-087    github.com/oppo-us-research/riav-mvs

# Overview

- Our core idea is a "learning-to-optimize" paradigm that iteratively indexes a plane-sweeping cost volume and regresses the depth map via a convolutional GRU.



(a) Input images

(b) Feature encoder and cost volume construction

(c) Index field iterative update

(d) Depth regression

# Overview

- Our core idea is a "learning-to-optimize" paradigm that iteratively indexes a plane-sweeping cost volume and regresses the depth map via a convolutional GRU.
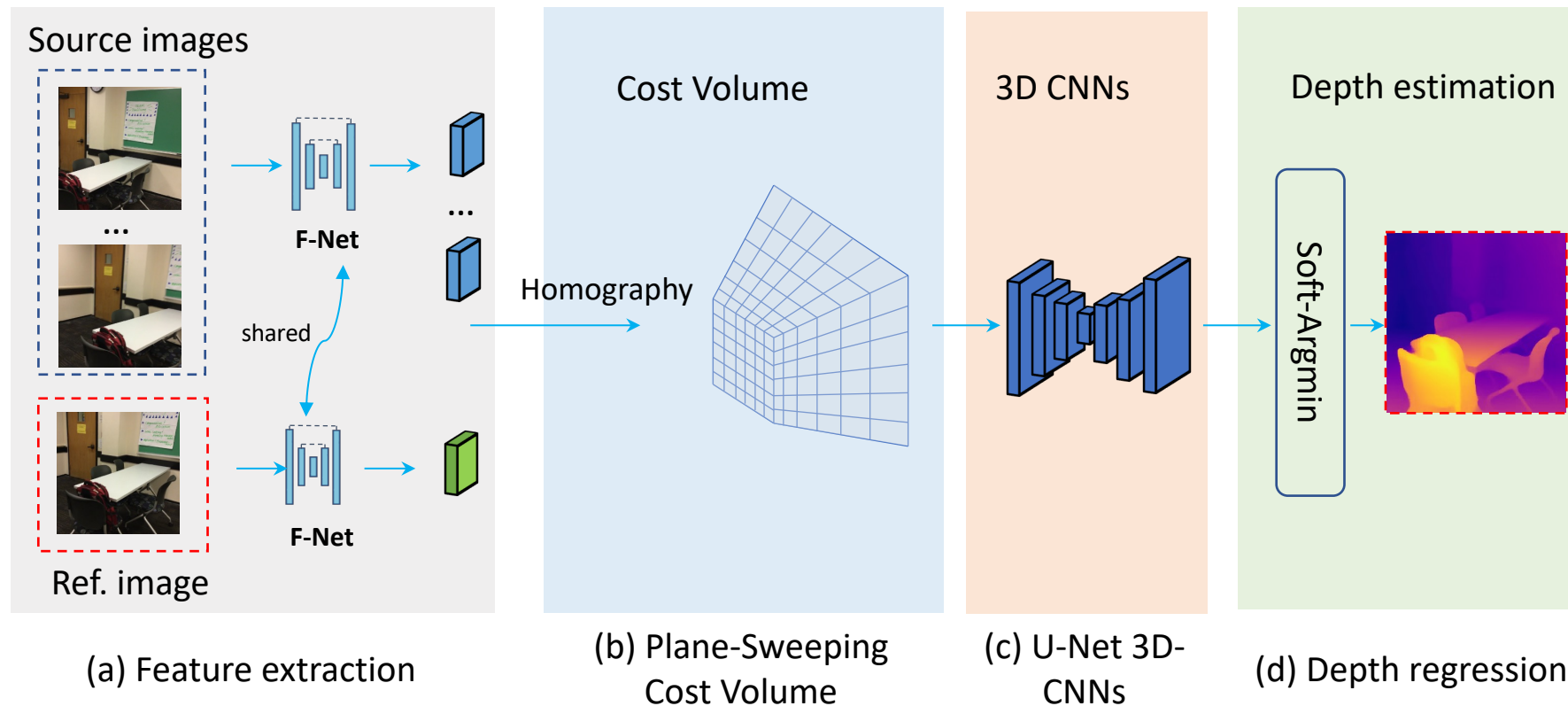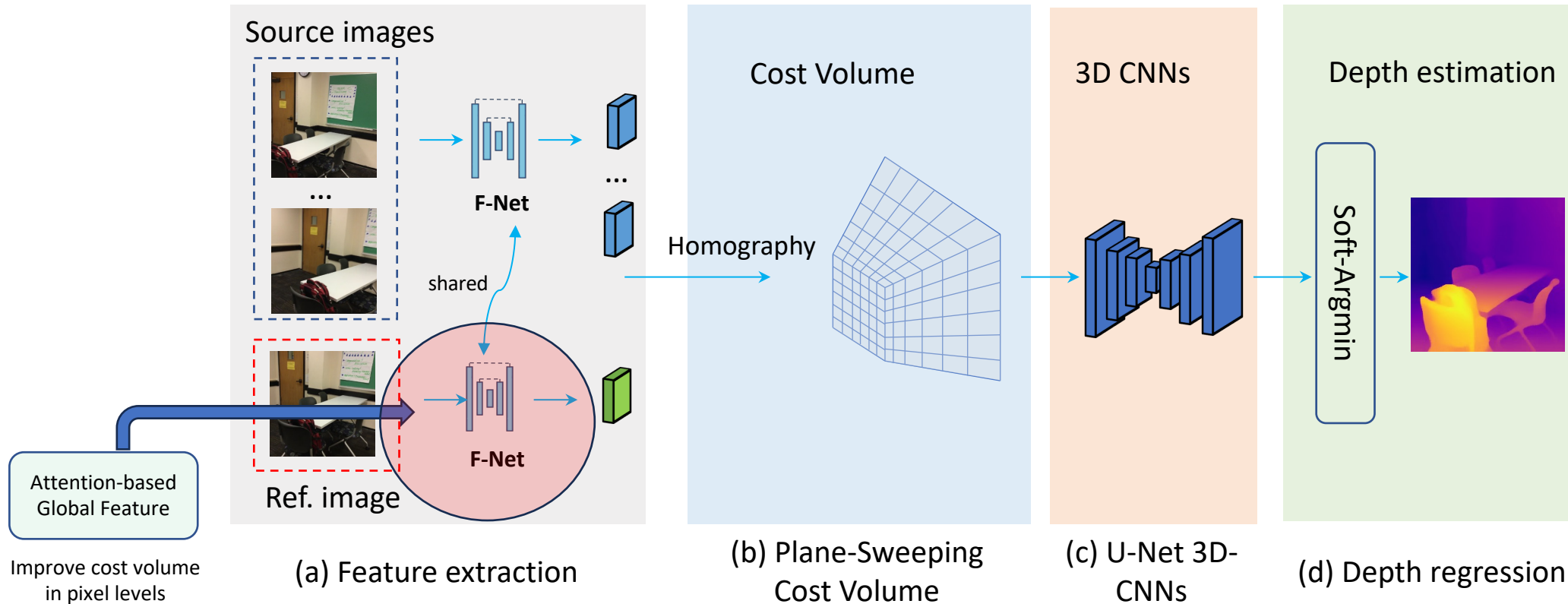


(a) Input images

(b) Feature encoder and cost volume construction

(c) Index field iterative update

(d) Depth regression

# Motivation

- Existing CNN-based MVS methods:
  - Concerns in (a), (b), (c) and (d)



(a) Feature extraction

(b) Plane-Sweeping Cost Volume

(c) U-Net 3D-CNNs

(d) Depth regression

# Contributions

- Our contributions: 1) An **asymmetric** cost volume ★★★☆☆



Source images

Attention-based Global Feature

Improve cost volume in pixel levels

Ref. image

**F-Net**

shared

**F-Net**

Cost Volume

Homography

3D CNNs

Soft-Argmin

Depth estimation

(a) Feature extraction

(b) Plane-Sweeping Cost Volume

(c) U-Net 3D-CNNs

(d) Depth regression

# Contributions

- Our contributions: 2) **Residual** pose update

★★★☆



(a) Feature extraction

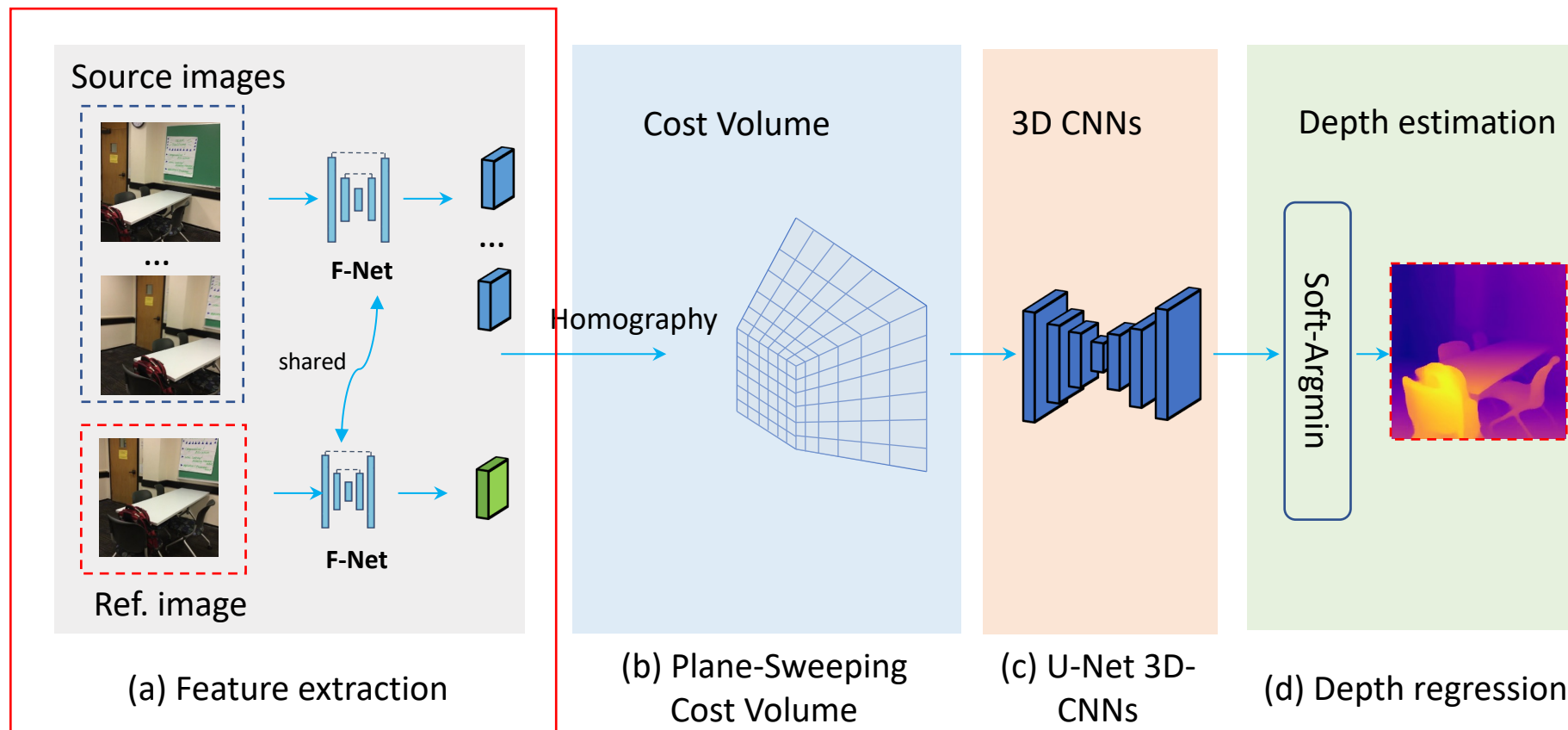(b) Plane-Sweeping Cost Volume

(c) U-Net 3D-CNNs

(d) Depth regression

# Contributions

- Our contributions: 3) A **new paradigm** to predict the depth by ⭐⭐⭐⭐⭐ learning to recurrently index cost volume via GRUs



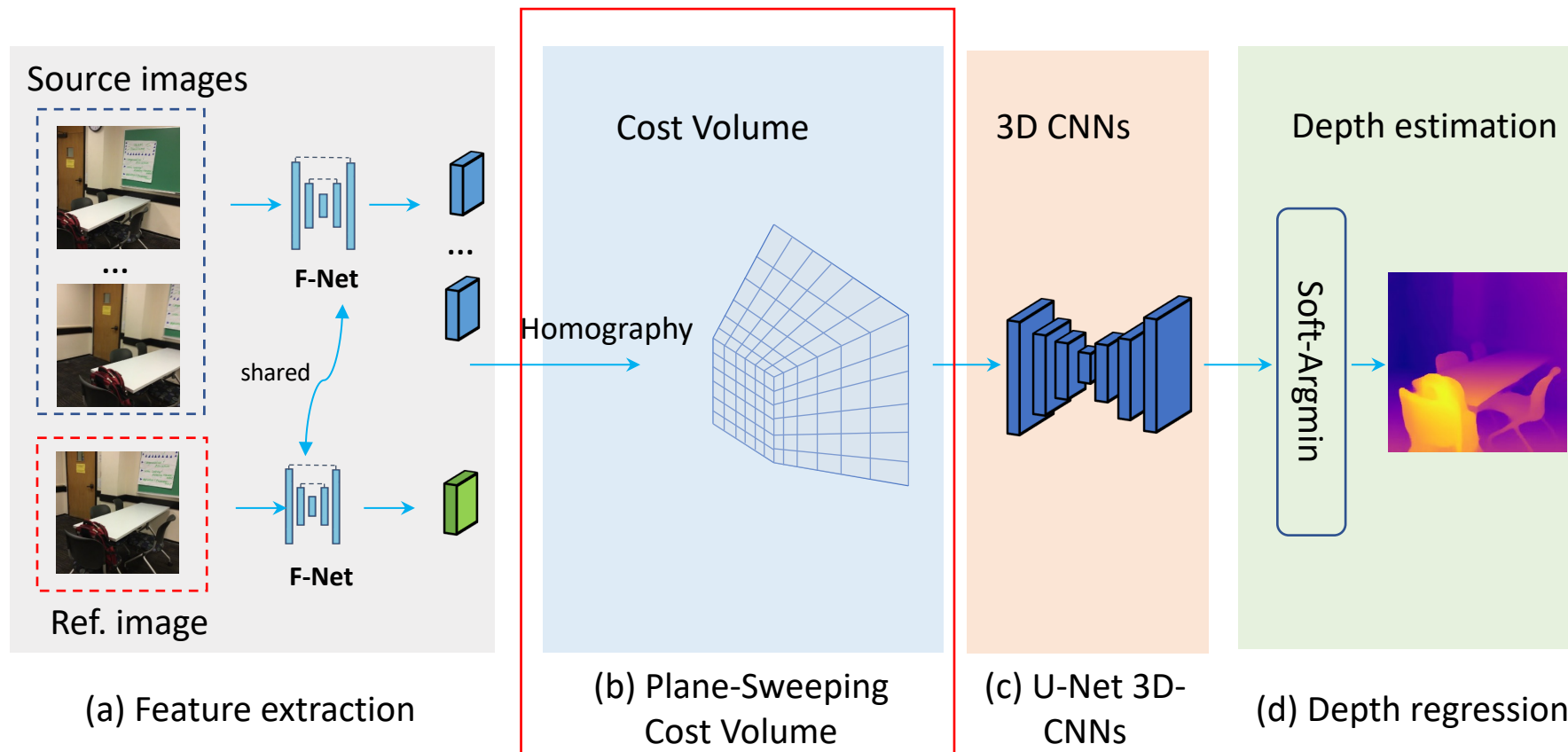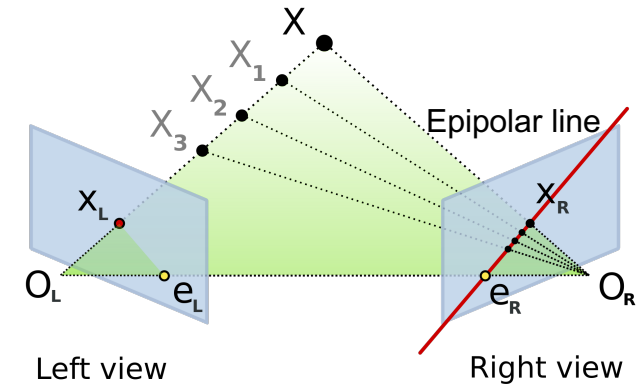Residual Pose

Improve cost volume in frame- levels

Source images

F-Net

shared

F-Net

Ref. image

Attention-based Global Feature

Improve cost volume in pixel levels

(a) Feature extraction

Cost Volume

Homography

(b) Plane-Sweeping Cost Volume

3D CNNs

(c) U-Net 3D-CNNs

Depth estimation

Soft-Argmin

(d) Depth regression

Recurrent Estimate of **Index Field**

# Existing CNN-based MVS Pipeline

- Existing CNN-based MVS methods:
  - Symmetric features, local context



(a) Feature extraction

(b) Plane-Sweeping Cost Volume

(c) U-Net 3D-CNNs

(d) Depth regression

# Existing CNN-based MVS Pipeline

- ## Existing CNN-based MVS methods:
  - ### assuming poses being accurate



(a) Feature extraction

(b) Plane-Sweeping Cost Volume

(c) U-Net 3D-CNNs

(d) Depth regression

# Existing CNN-based MVS Pipeline

- Existing CNN-based MVS methods:
  - 3D CNNs are time and memory consuming
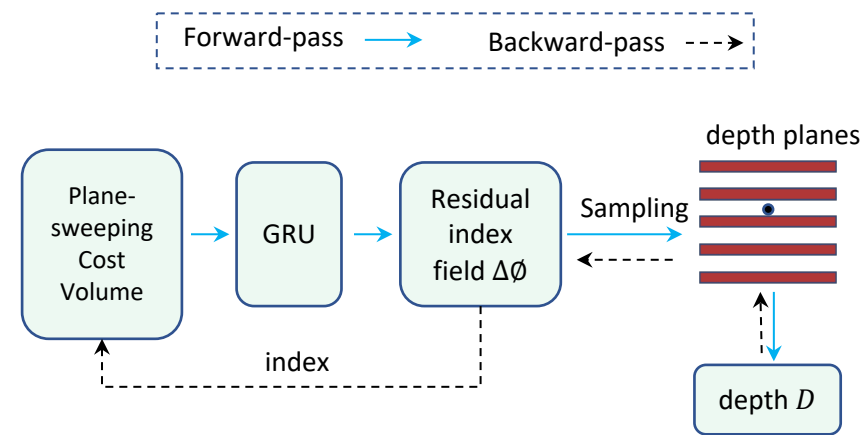  - Soft-argmin is not robust to multi-modal distributions



(a) Feature extraction

(b) Plane-Sweeping Cost Volume

(c) U-Net 3D-CNNs

(d) Depth regression
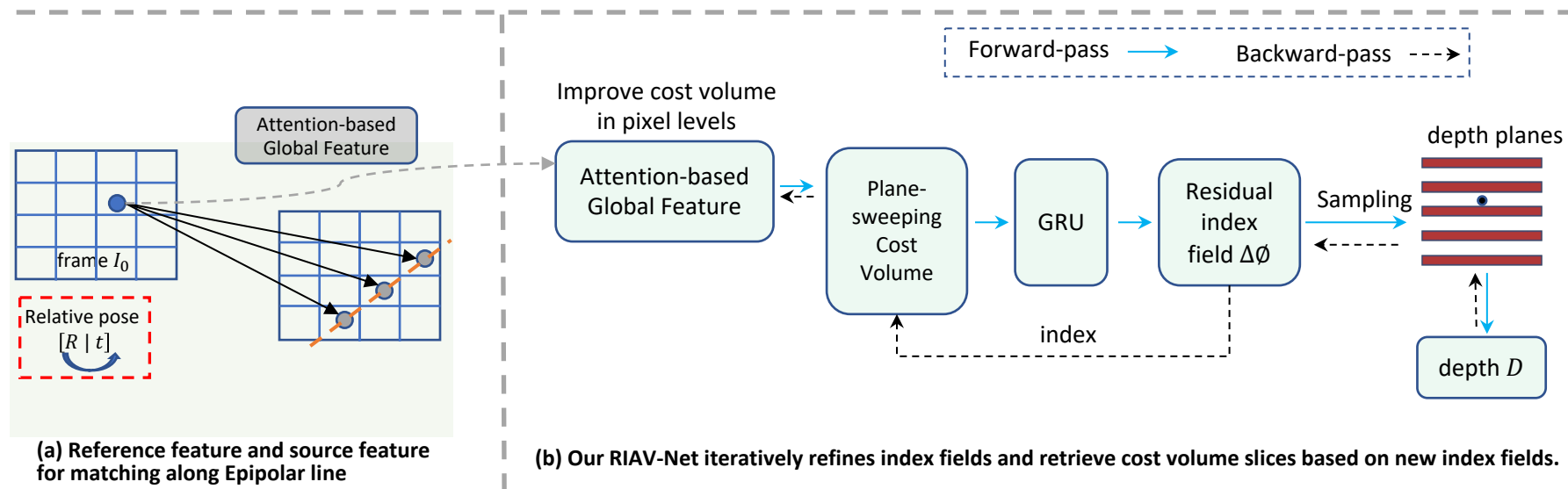
# Our Approach

- Constructing a good cost volume:



(a) Reference feature and source feature for matching along Epipolar line

(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.
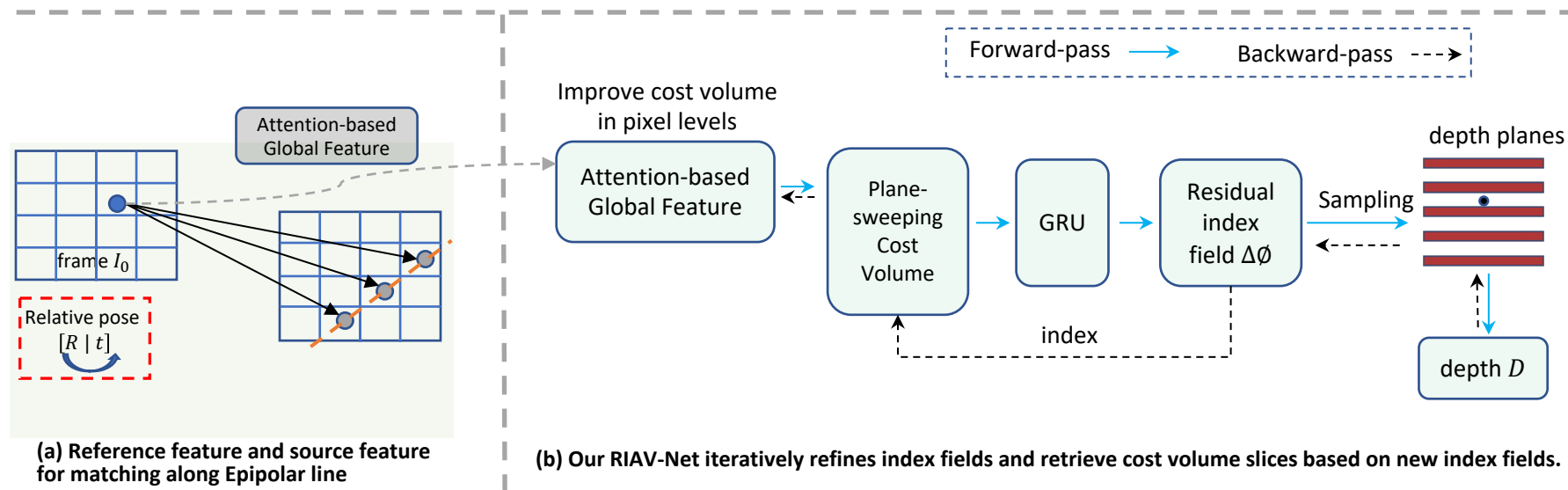
# Our Approach

- Constructing a good cost volume:
  1) To break the symmetry of the Siamese network by introducing a transformer block to the reference image (but not to the source images)
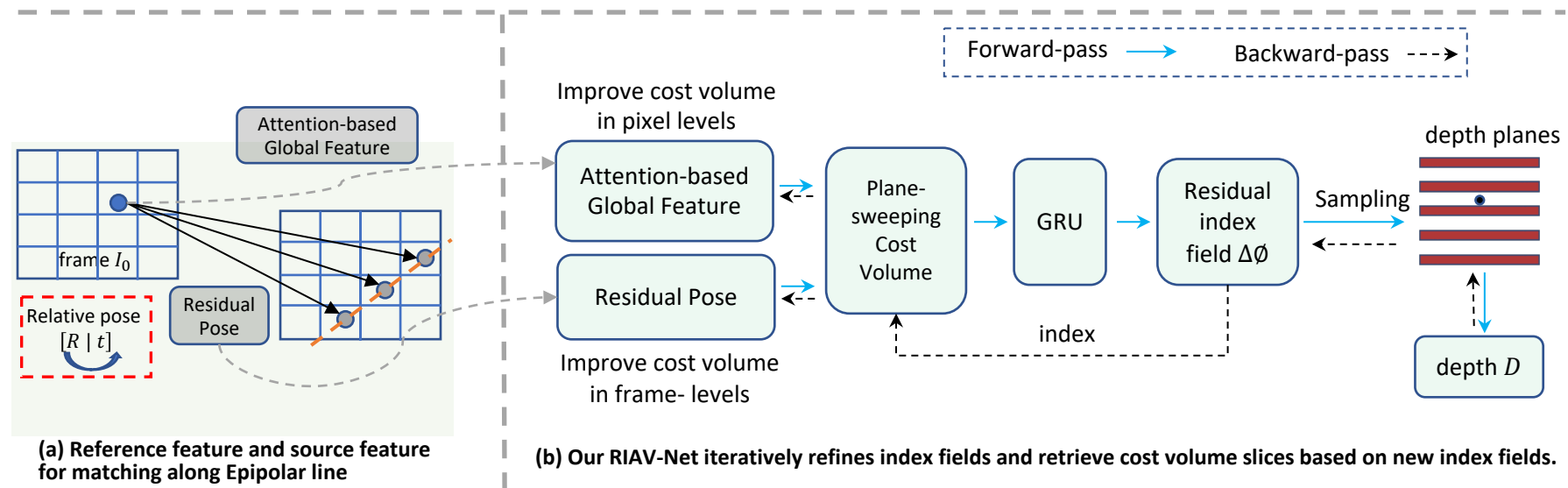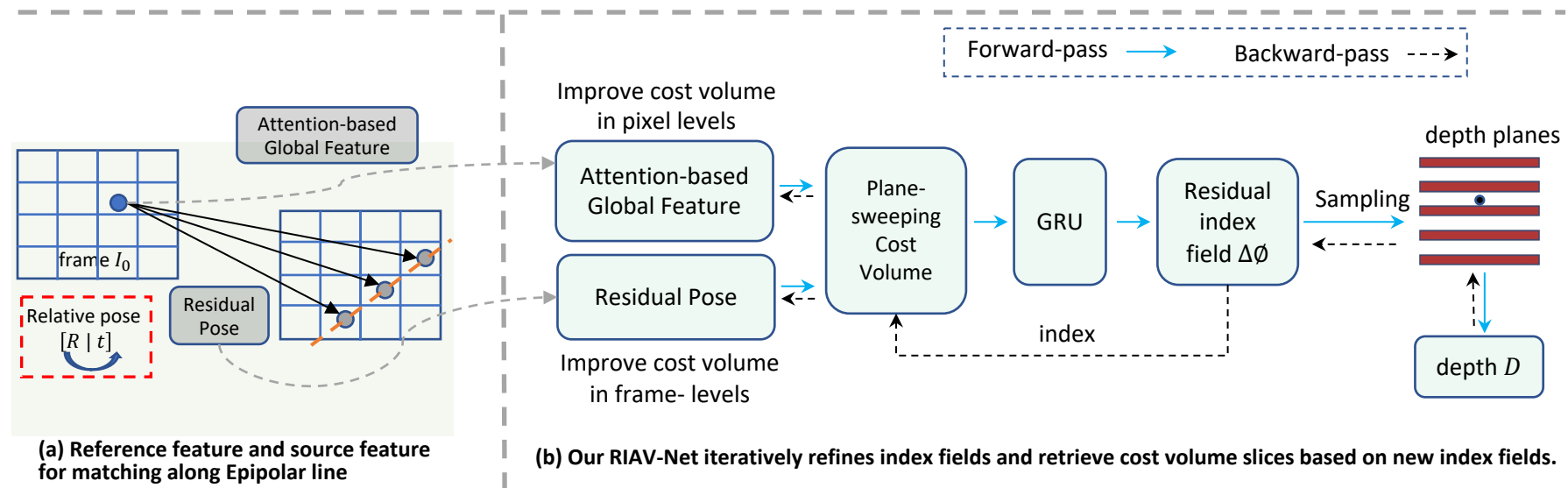


(a) Reference feature and source feature for matching along Epipolar line

(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.

# Our Approach

- Constructing a good cost volume:

  1) To break the symmetry of the Siamese network by introducing a transformer block to the reference image (but not to the source images)



(a) Reference feature and source feature for matching along Epipolar line

(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.

# Our Approach

- Constructing a good cost volume:

  2) To incorporate a residual pose network to correct the relative poses



(a) Reference feature and source feature for matching along Epipolar line

(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.

# Our Approach

- Constructing a good cost volume:
  - 2) To incorporate a residual pose network to correct the relative poses



(a) Reference feature and source feature for matching along Epipolar line

(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.

# Our Approach

- A new paradigm to predict the depth via learning the proposed **index filed** to recurrently index an asymmetric plane-sweeping cost volume via GRUs



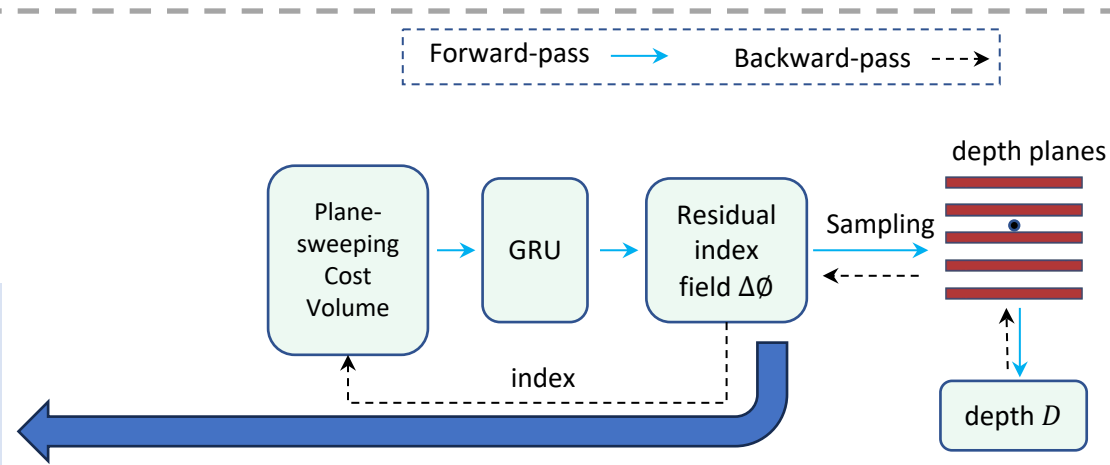(a) Reference feature and source feature for matching along Epipolar line

(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.

# Background - Ours vs RAFT (ECCV'20)
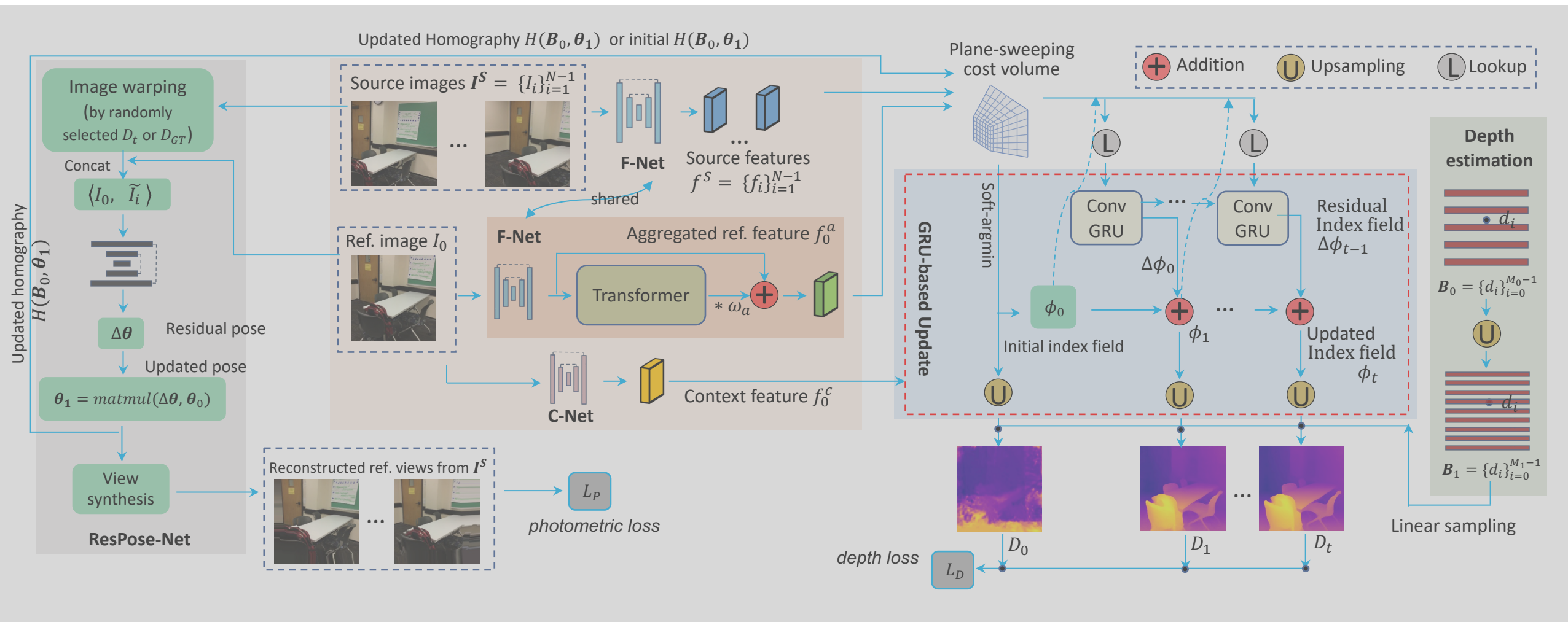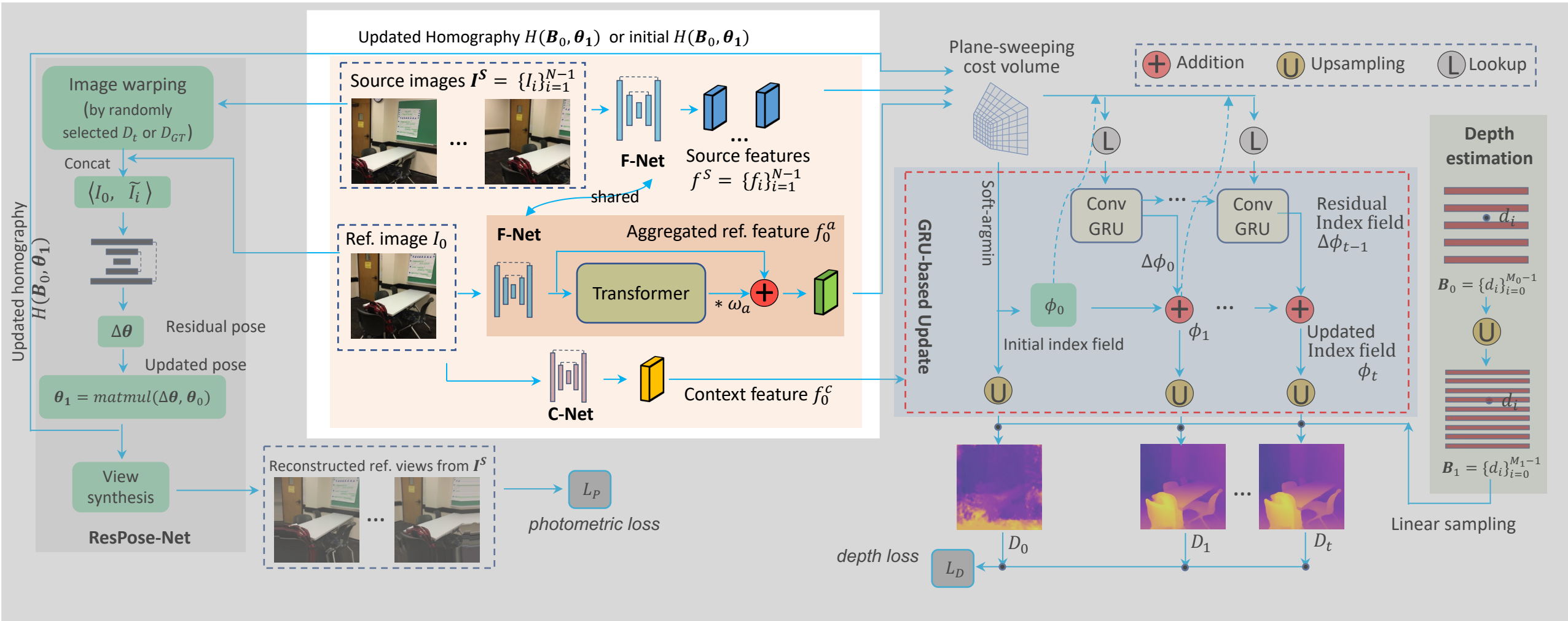
- We borrowed ideas from RAFT for learning to optimize via GRU:

  - RAFT's all-pair correlation for optical flow: NO multi-view geometry constraints → Ours use plane-sweeping cost volume for MVS (Fig. a,b&c)

  - We propose **index field** that serves as a new design to bridge cost volume optimization and depth map estimation (Fig. e)



(a) RAFT (all-pair correlation volume) for Optical Flow.

(b) Ours (plane-sweep cost volume) for MVS depth.

(c) RAFT Iteratively updates optical flow.

(e) Our RIAV-MVS iteratively refines index fields, which are used to retrieve cost volume slices and linearly sample the depth planes to estimate the depth maps.

# Background - Ours vs IterMVS (CVPR'22)

- IterMVS iteratively predicts a depth and reconstructs a new plane-sweeping cost volume using the updated depth planes (Fig. d)

- Ours learns to index the cost volume by approaching the "correct" depth planes per pixel via an index field (Fig. e)



(a) RAFT (all-pair correlation volume) for Optical Flow.

(b) Ours (plane-sweep cost volume) for MVS depth.

(c) RAFT Iteratively updates optical flow.

(d) IterMVS iteratively updates depth and reconstructs a new cost volume using the new depth planes.

(e) Our RIAV-MVS iteratively refines index fields, which are used to retrieve cost volume slices and linearly sample the depth planes to estimate the depth maps.

# Our RIAV-MVS

- Our proposed RIAV-Net iteratively refines index fields and retrieve plane-sweeping cost volume slices based on new index fields.
    - a **residual** index field $\Delta\emptyset$ is predicted as an update direction for next iteration



(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.

# Our RIAV-MVS
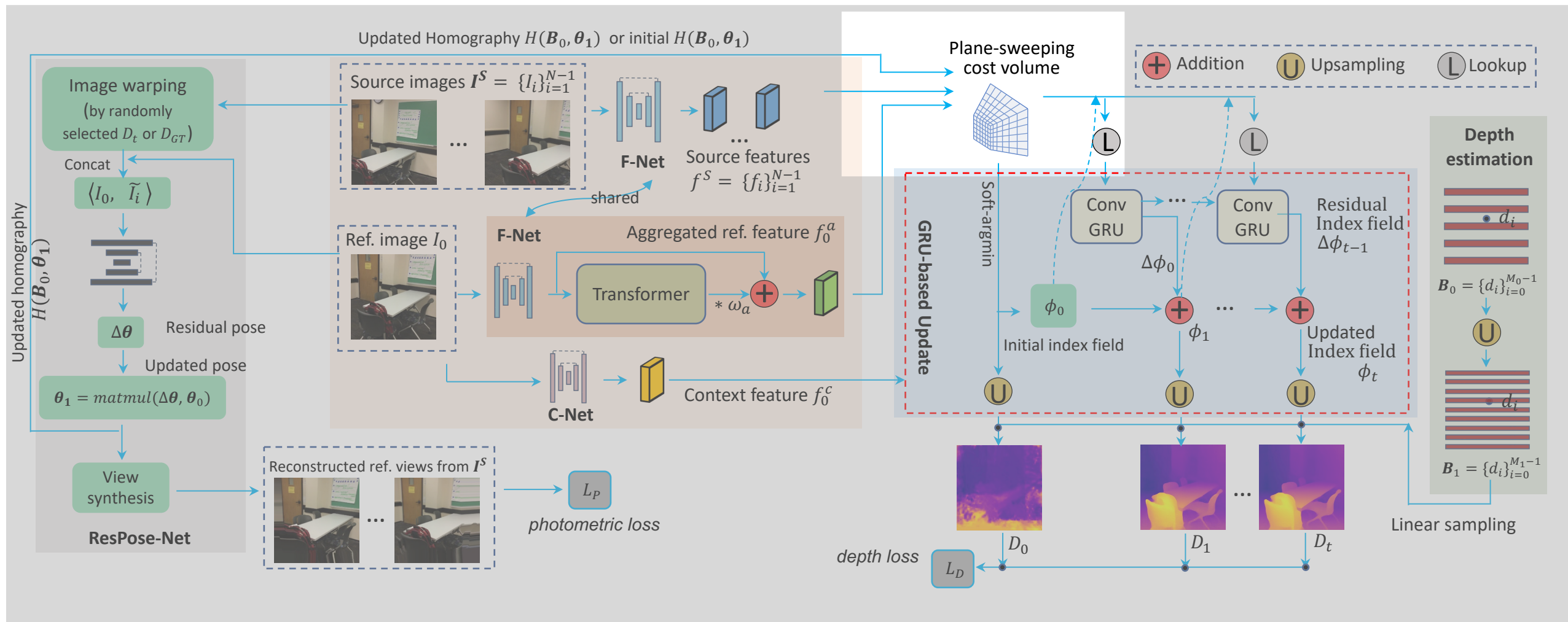
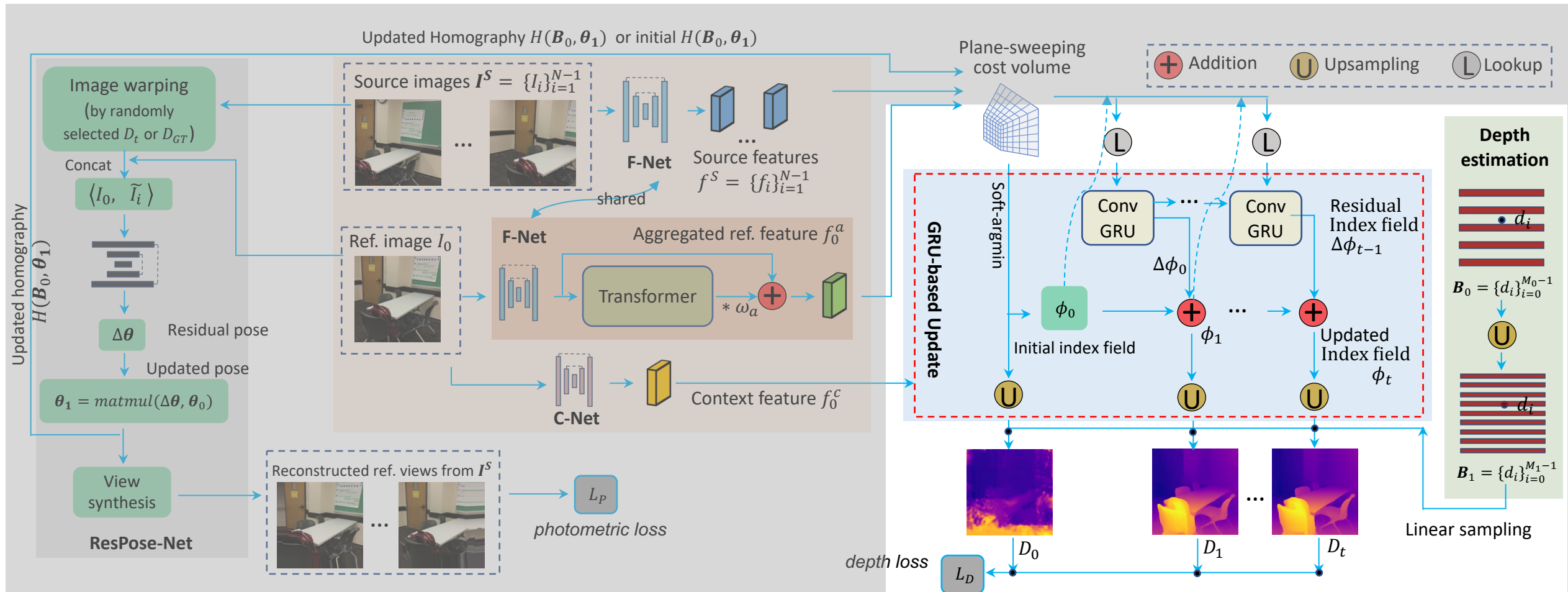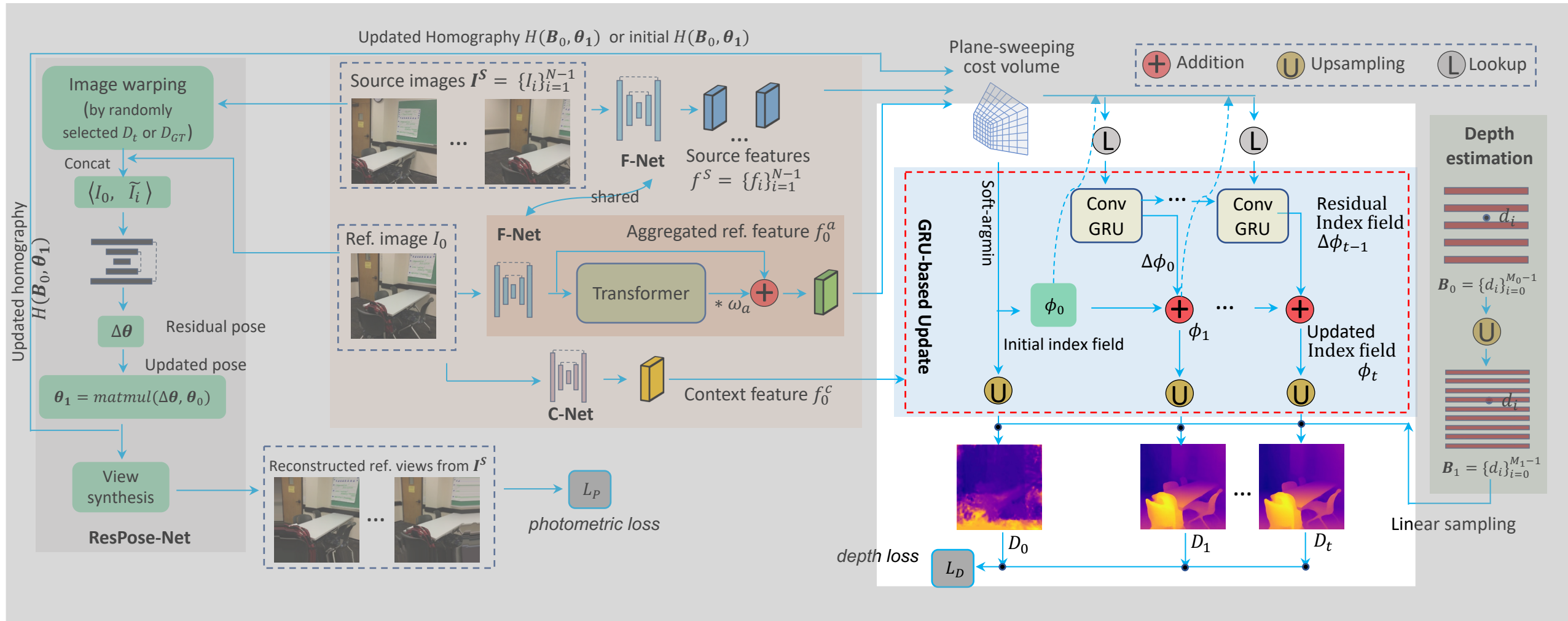- Our proposed RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.
  - a depth map is estimated by sampling depth plane hypotheses via index field



(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.

# Our RIAV-MVS

- Our proposed RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.
    - a depth map is estimated by sampling depth plane hypotheses via **index field**



Our *index field* serves as a new design to bridge cost volume optimization and depth map estimation

(b) Our RIAV-Net iteratively refines index fields and retrieve cost volume slices based on new index fields.

# Architecture

# Architecture

# Architecture

# Architecture

# Architecture

# Architecture

# Architecture

# Architecture

# Architecture

# Architecture: Feature & Cost Volume



Our **asymmetric** employment of this transformer layer provides the capability to better balance the **high-frequency** (by high-pass CNNs) and **low-frequency** features (by self-attention).

# Architecture: GRU-based Iterative Updates

# Architecture: GRU-based Iterative Updates



The proposed recurrent estimate of **index field** (i.e., a grid of indices to identify the depth hypotheses) enables the learning to be anchored at the cost volume domain.

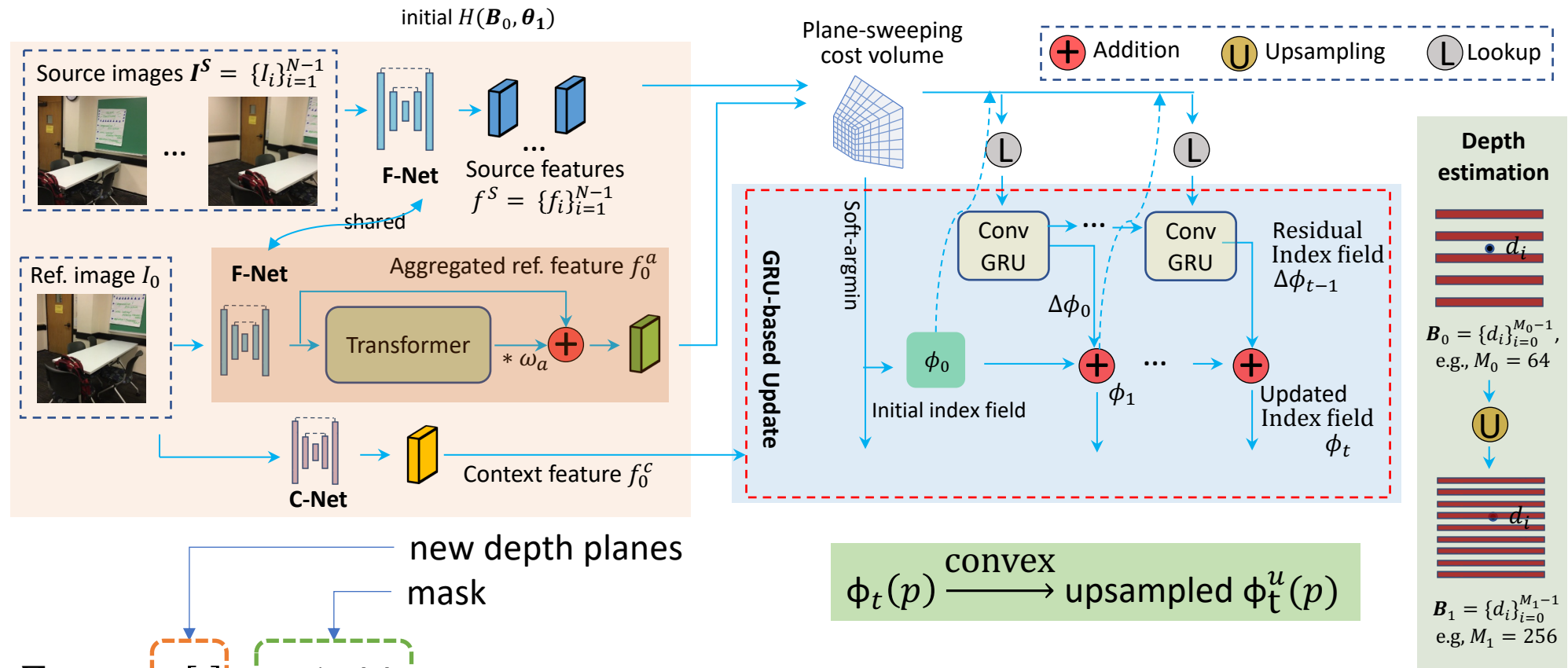# Architecture: Upsampling & Depth Estimation

# Architecture: Upsampling & Depth Estimation



$$D_t(p) = \frac{\sum_{i \in \Omega(p)} B[i] \cdot W_1(p, \lfloor i \rfloor)}{\sum_{j \in \Omega(p)} W_1(p, \lfloor j \rfloor)}$$

Final depth

new depth planes

mask

Neighbor with a radius r=4 centered at upsampled index $\phi_t^u(p)$ for a given pixel p
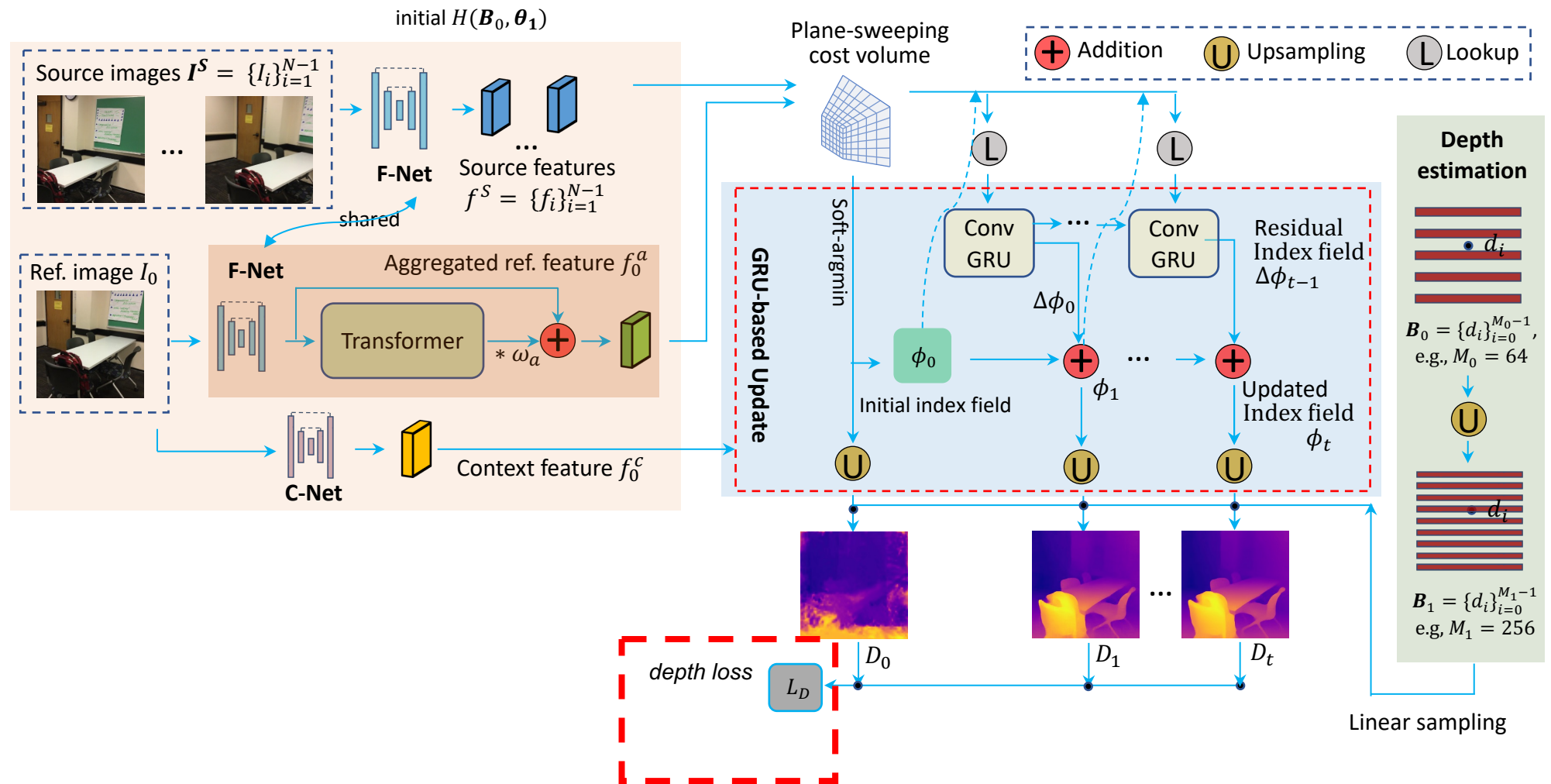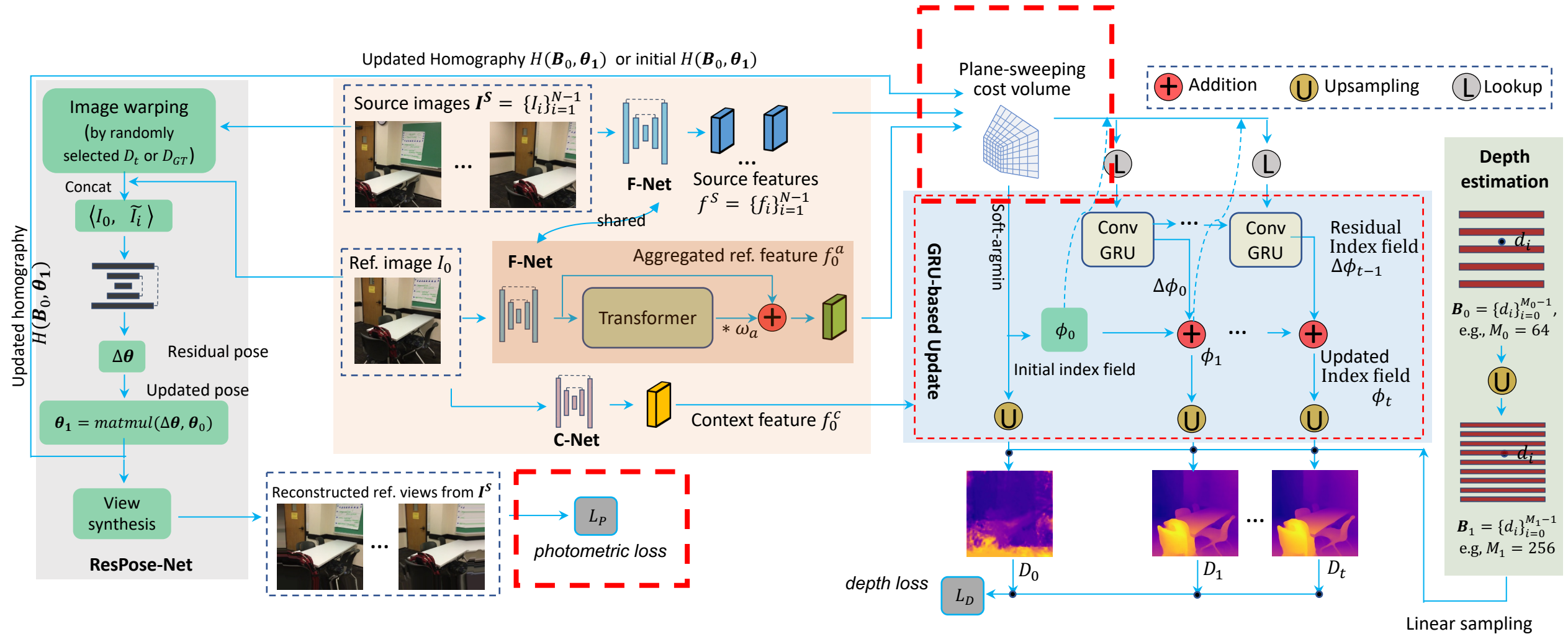
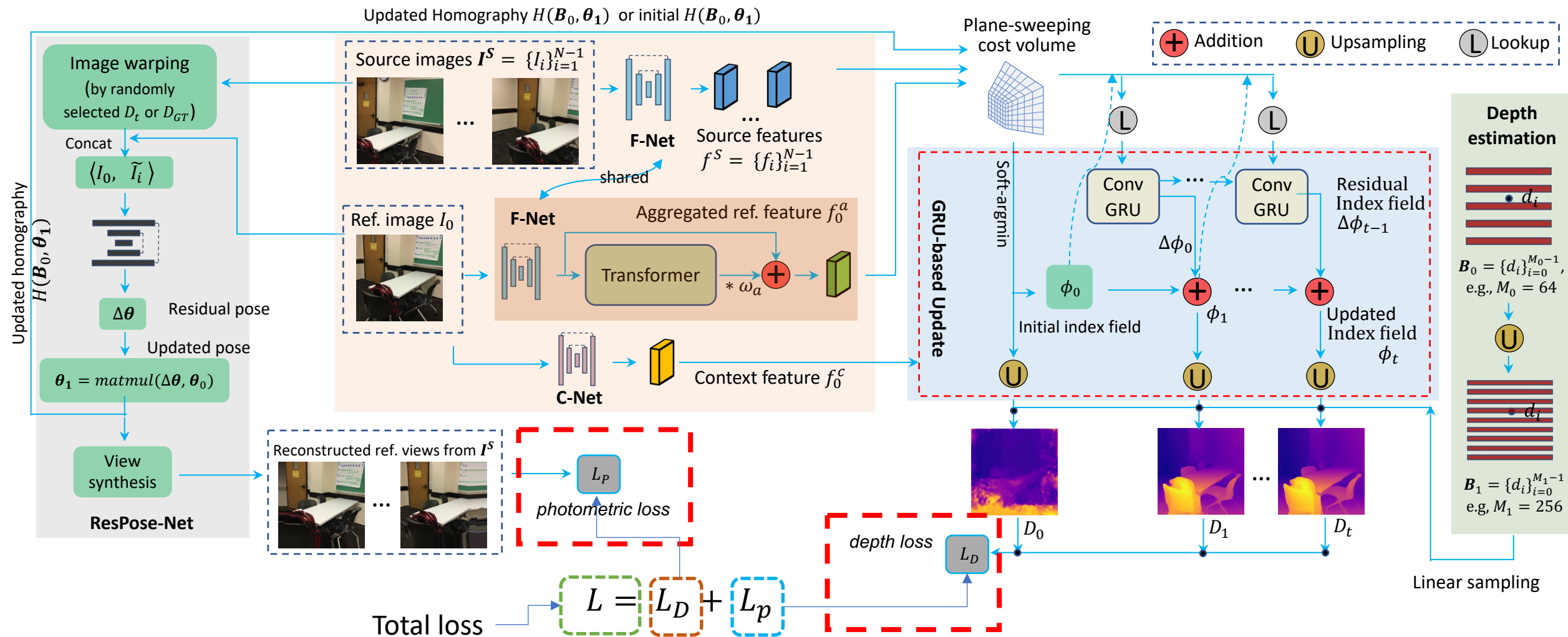# Architecture: Upsampling & Depth Estimation

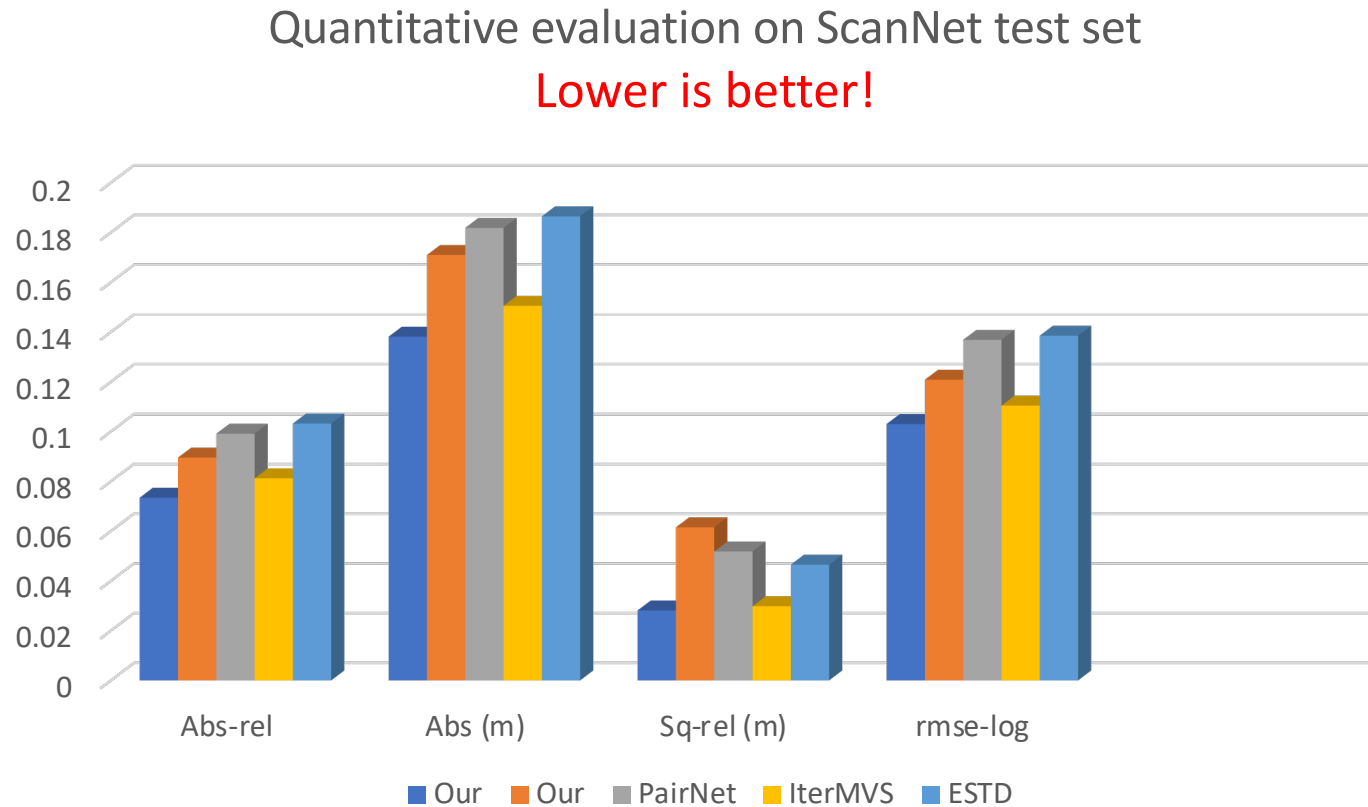# Architecture: Residual Pose Net

# Architecture: Residual Pose Net

# Architecture: Loss Function

# Experimental Results

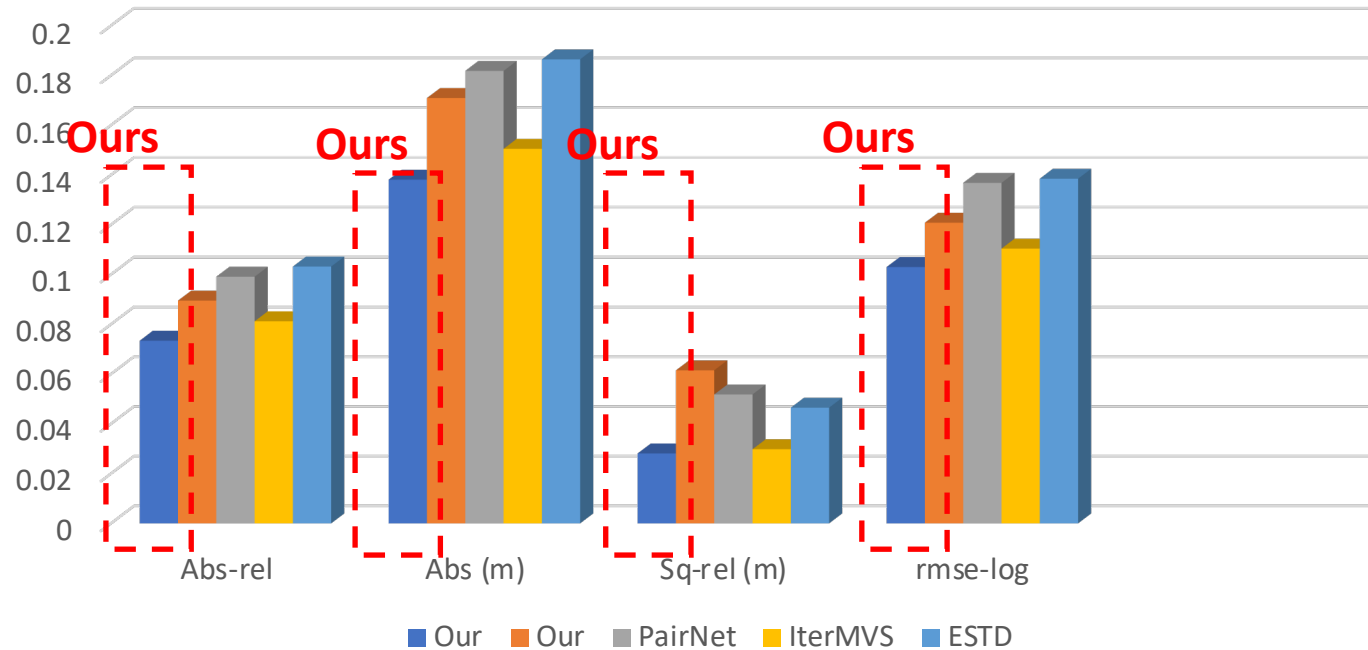- Depth map evaluation on ScanNet Testset



Quantitative evaluation on ScanNet test set
Lower is better!

# Experimental Results

- Depth map evaluation on ScanNet Testset



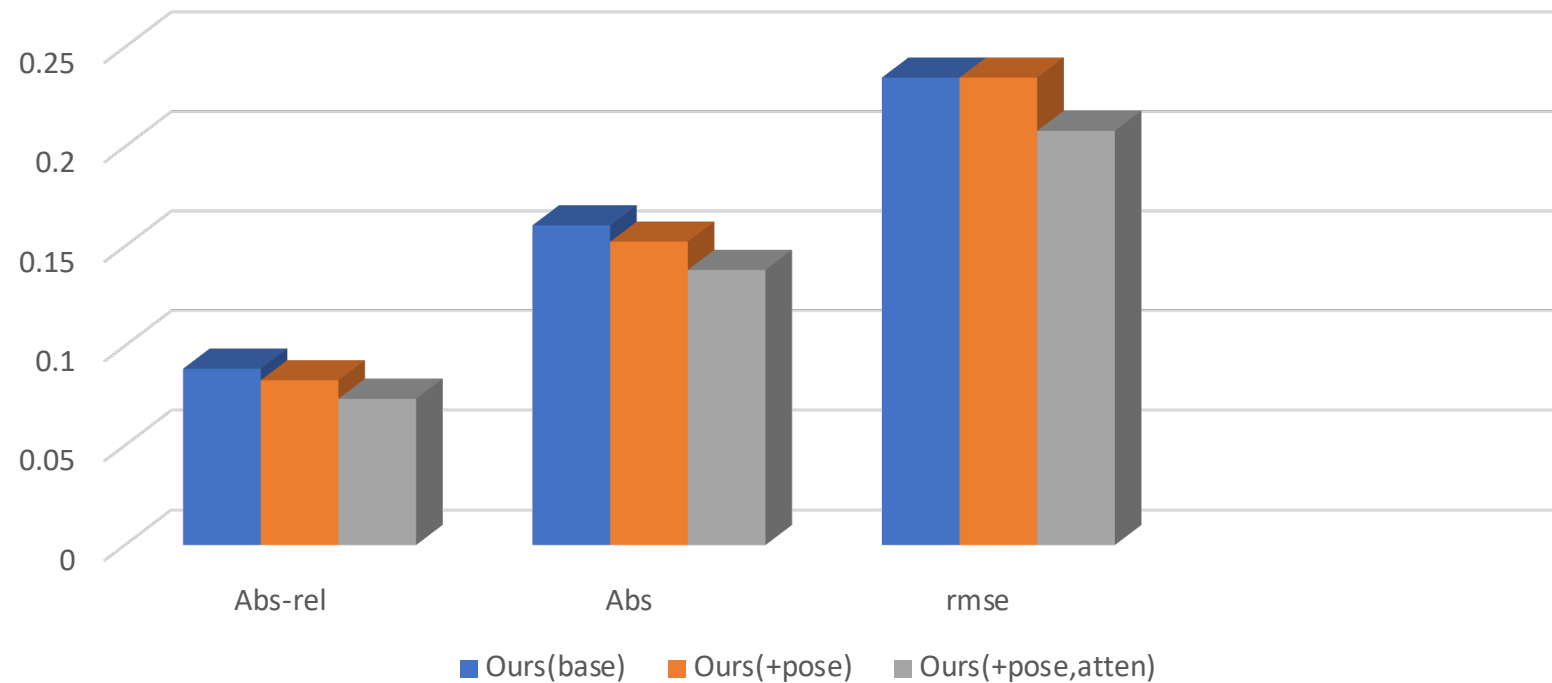Quantitative evaluation on ScanNet test set
Lower is better!

# Experimental Results

- Ablation study: three variants of our method

Comparison of three variants of our models on ScanNet test set
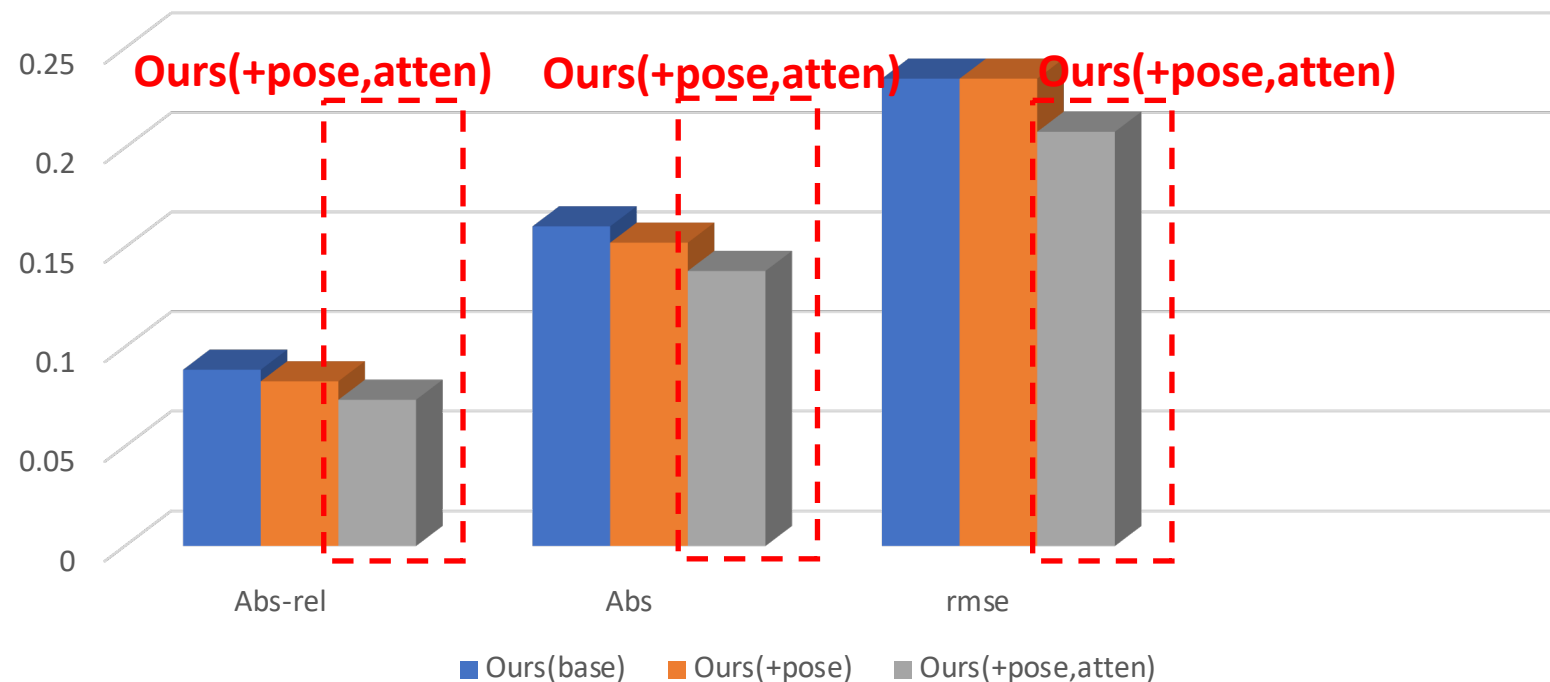Lower is better!

# Experimental Results

- Ablation study: three variants of our method



Comparison of three variants of our models on ScanNet test set
Lower is better!
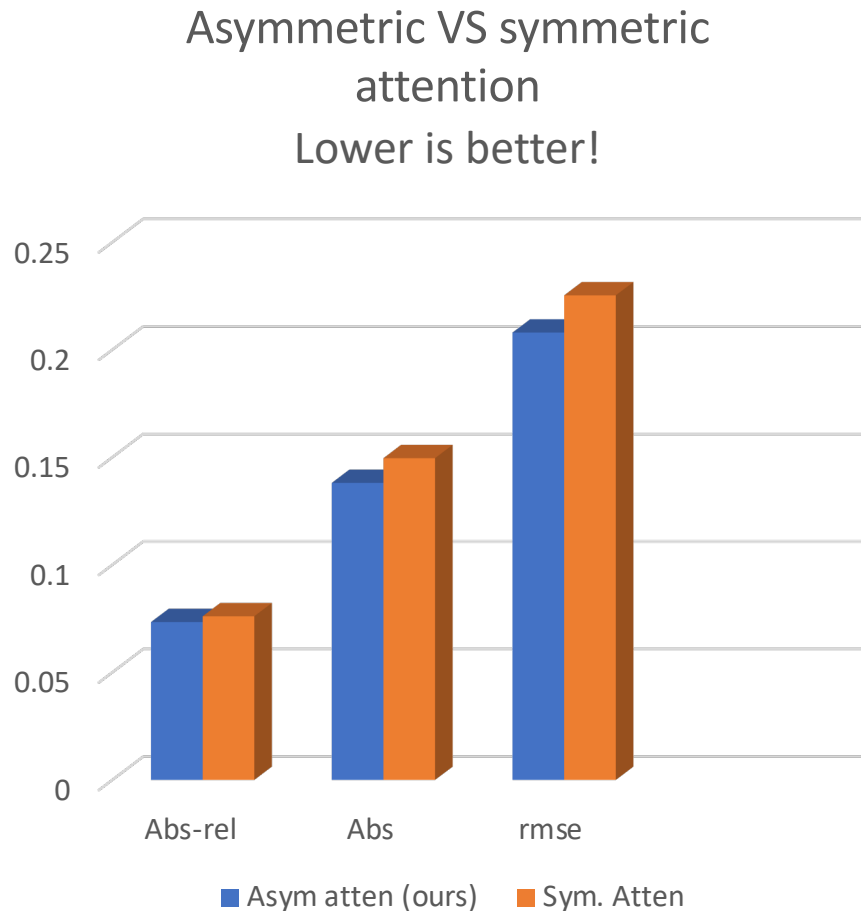
# Experimental Results

- ## Asymmetric attention

Asymmetric VS symmetric attention
Lower is better!



- ## Our attention applied to MVSNet
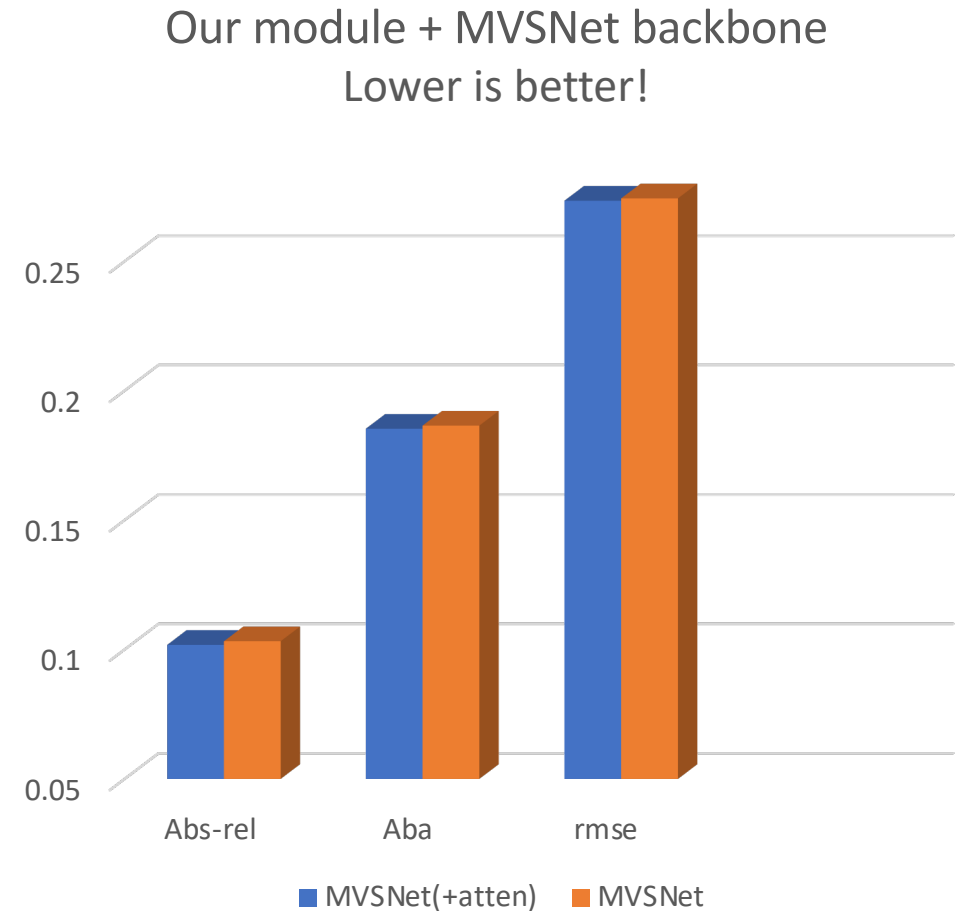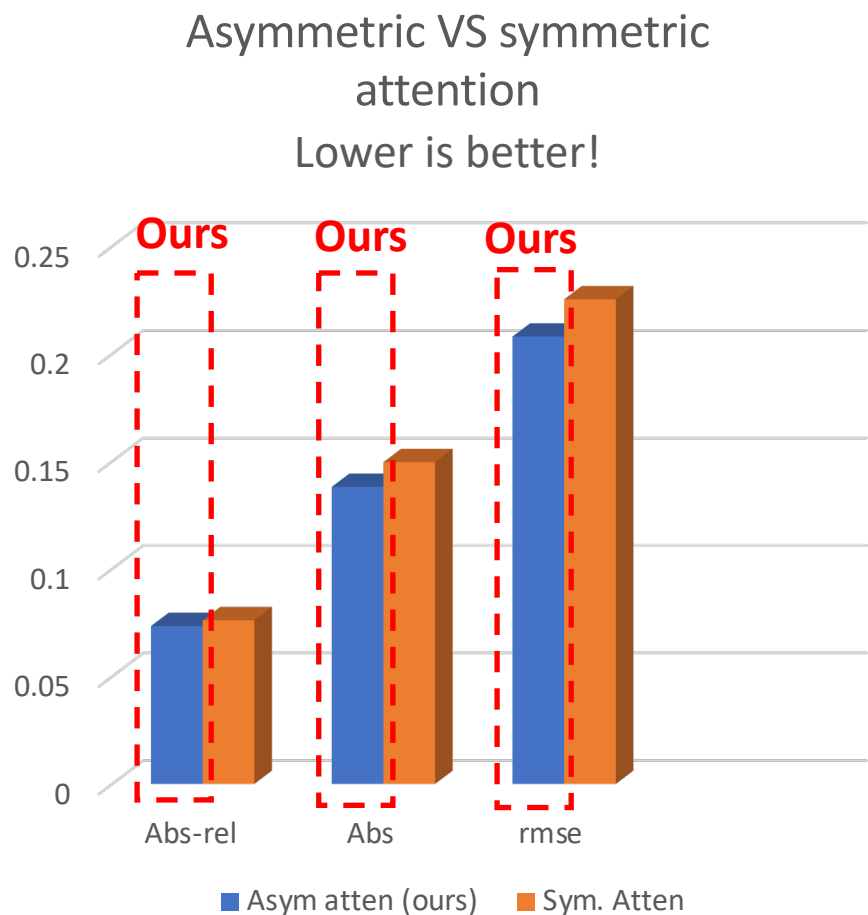
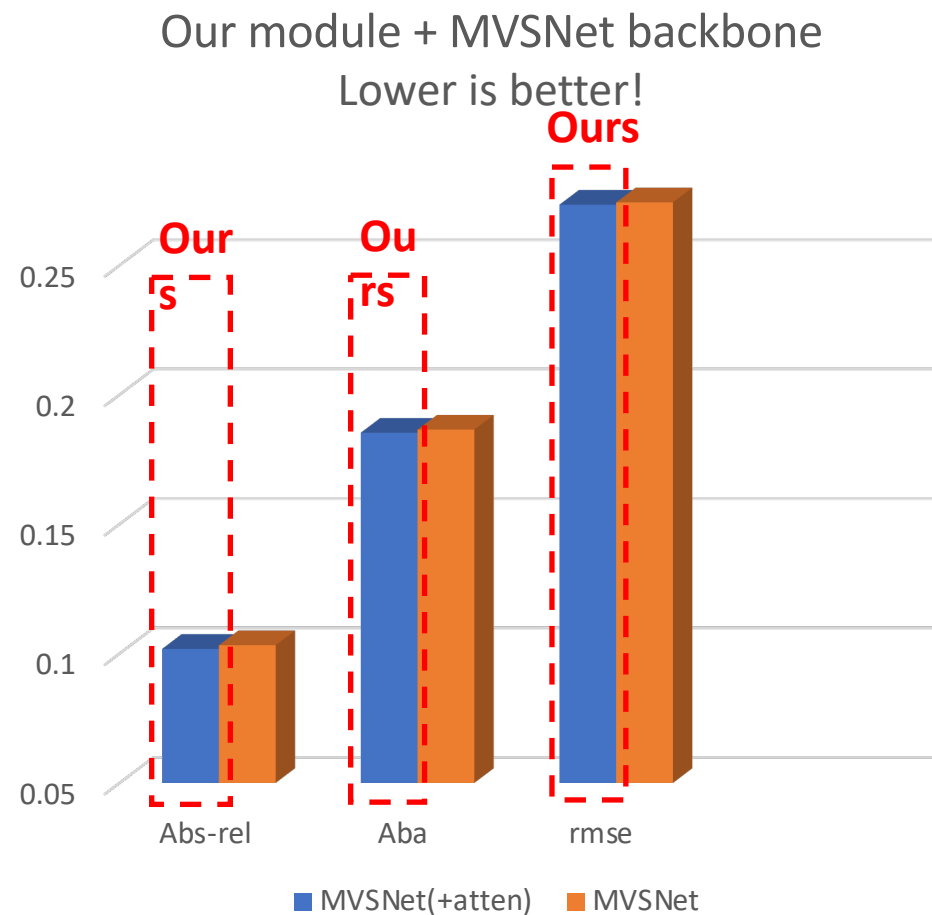Our module + MVSNet backbone
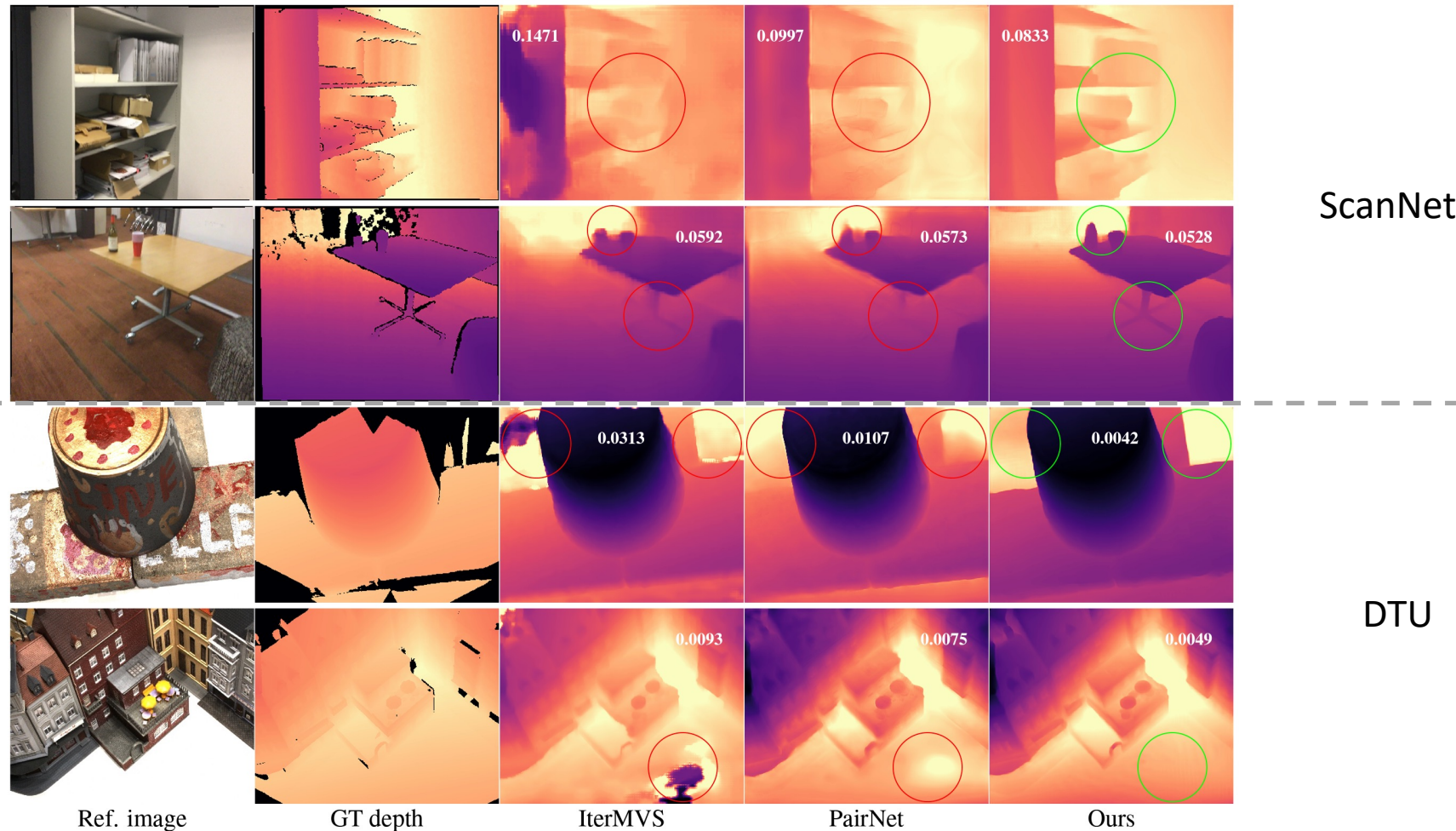Lower is better!

# Experimental Results

- Asymmetric attention



- Our attention applied to MVSNet

# Experimental Results

- Qualitative results on ScanNet (top two rows) and DTU test set



ScanNet

DTU

| Ref. image | GT depth | IterMVS | PairNet | Ours |

# Experimental Results

- More depth results and 3D point clouds on ScanNet

# Conclusion

- RIAV-MVS, as a new paradigm to predict depth by learning to recurrently index cost volume via GRUs

★★★★★

- An asymmetric cost volume by a transformer block applied to the reference image

★★★☆☆

- A Residual pose network to update the relative poses to improve cost volume

★★★★☆

# Conclusion

- RIAV-MVS, as a new paradigm to predict depth by learning to recurrently index cost volume via GRUs

  ★★★★★

- An asymmetric cost volume by a transformer block applied to the reference image

  ★★★☆☆

- A Residual pose network to update the relative poses to improve cost volume

  ★★★★☆

# Conclusion

- RIAV-MVS, as a new paradigm to predict depth by learning to recurrently index cost volume via GRUs ★★★★★

- An asymmetric cost volume by a transformer block applied to the reference image ★★★☆☆

- A Residual pose network to update the relative poses to improve cost volume ★★★★☆

# Thank You!

Code coming soon
https://github.com/oppo-us-research/riav-mvs