# CheapNVS: Real-Time On-Device Narrow-Baseline Novel View Synthesis

Konstantinos Georgiadis[†][*], Mehmet Kerim Yucel[‡][*], and Albert Saà-Garriga[‡][*]

[†]CERTH, Information Technologies Institute, Thessaloniki, Greece

[‡]Samsung R&D Institute UK (SRUK)

[*]Three authors contributed equally.

*Abstract*—Single-view novel view synthesis (NVS) is a notorious problem due to its ill-posed nature, and often requires large, computationally expensive approaches to produce tangible results. In this paper, we propose CheapNVS: a fully end-to-end approach for narrow baseline single-view NVS based on a novel, efficient multiple encoder/decoder design trained in a multi-stage fashion. CheapNVS first approximates the laborious 3D image warping with lightweight learnable modules that are conditioned on the camera pose embeddings of the target view, and then performs inpainting on the occluded regions in parallel to achieve significant performance gains. Once trained on a subset of Open Images dataset, CheapNVS outperforms the state-of-the-art despite being $10 \times$ faster and consuming 6% less memory. Furthermore, CheapNVS runs comfortably in real-time on mobile devices, reaching over 30 FPS on a Samsung Tab 9+.

*Index Terms*—Novel View Synthesis; 3D Photography

## I. INTRODUCTION

The ability to synthesize novel views of a scene/object from a single image is critical in computer vision, with applications in robotics [1], VR [2], multimedia production [3] and biomedicine [4]. Novel view synthesis (NVS) is an ill-posed problem, as it is not only comprised of reverse-projecting an image onto the 3D space, but also of the completion of the missing data in occluded regions. Despite the recent advances in NVS [5], [6], [7], there are several issues with existing methods: they either i) require per-scene training and lack generalization across-scenes [5], [6], ii) overfit to fixed stereo baselines [8], [9], [10], [11], iii) are unfit for on-device, real-time processing [5], [6], [7], [8], [12], [13], [14], [15] or iv) are not end-to-end learnable solutions [8], [16], [11].

In this paper, we aim to take a step towards meeting all the criteria above. To this end, we propose CheapNVS; an end-to-end, device-friendly, real-time solution for narrow baseline NVS. First, CheapNVS introduces an efficient warping module that approximates 3D image warping, conditioned on the camera pose of the target view, which speeds up the NVS pipeline. Second, unlike the contemporary methods that perform inpainting after the warping [16], [13], [12], CheapNVS performs inpainting in parallel, enjoying further speedups. Third, we propose to use Open Images as our training set. Based on a novel architecture that consists of a shared RGBD encoder, an extrinsics encoder, and three decoders, CheapNVS enjoys significant runtime and memory improvements, with a competitive or better accuracy to state-of-the-art. It runs comfortably in real-time on mobile devices, and since it is built on simple network blocks, can leverage hardware accelerators without requiring any low-level optimization expertise. The diagram of CheapNVS is shown in Figure 1.

## II. RELATED WORK

**Novel View Synthesis.** NVS is arguably the ultimate aim of 3D reconstruction, where a plethora of methods are proposed, ranging from earlier methods [17] to neural renderers [5], [6]. Narrow baseline NVS, where the target view is not *distant* from the source view, does not necessarily require exhaustive 3D reconstruction. This has been exploited in many studies, including multi-view NVS [18], [19] and 3D photography [16], [12], [13], to achieve performant methods.

**Single-view Novel View Synthesis.** Single-view NVS aims to do the same but using a single image. Stereo synthesis methods [8], [9], [10], [11] perform NVS for a fixed baseline - mostly including only translation - so they only work for the baseline they are trained on. Modern one-shot neural rendering methods based on diffusion models provide good results, but they are computationally expensive or require per-scene optimization, thus not suitable for device deployment [2], [15].

The most relevant literature to ours is 3D photography [16], [12], [13]. Our work differs from these in the following aspects; CheapNVS i) efficiently approximates 3D warping, whereas others perform 3D warping in the conventional, expensive way, ii) performs inpainting in parallel, not sequentially like others and iii) does away with representations such as Multi-Plane Images (MPI) [13] and Layered Depth Images (LDI) [16], yet shows competitive results with significant efficiency gains and iv) comfortably runs in real-time on mobile devices, whereas others struggle to do so.

## III. CHEAPNVS

### A. Preliminaries

Given an input image $\mathbb{I}_s \in \mathbb{R}^{H \times W \times 3}$ and its depth map $\mathbb{D} \in \mathbb{R}^{H \times W}$, single-view NVS aims to synthesize a new image $\mathbb{I}_t \in \mathbb{R}^{H \times W \times 3}$ of the same scene from a given target camera pose $\mathbb{T} \in \mathbb{R}^{3 \times 4}$. We can write single-view NVS as

$$\mathbb{I}_t = f(w(\mathbb{I}_s, \mathbb{T}, \mathbb{D}; \theta_w), \mathbb{M}; \theta_f) \tag{1}$$

where $w(\cdot; \theta_w)$ is the function that implements image warping, $\mathbb{M} \in \mathbb{R}^{H \times W}$ is the occlusion mask which indicates the areas to be *filled* in the warped image, and $f(\cdot; \theta_f)$ is the function that performs the filling to synthesize the novel view.

### B. A better NVS formulation.

The primary issue of Equation 1 is that $f(\cdot; \theta_f)$ requires the output of $w(\cdot; \theta_w)$, which makes the process inherently sequential. We hypothesize that we can perform these in parallel, and have a modular and an interpretable architecture. An alternative to Equation 1 is

$$\mathbb{I}_t = w(\mathbb{I}_s, \mathbb{T}, \mathbb{D}; \theta_w) \cdot \phi(\mathbb{I}_s, \mathbb{T}, \mathbb{D}; \theta_\phi) +$$
$$f(\mathbb{I}_s, \mathbb{T}, \mathbb{D}; \theta_f) \cdot (1 - \phi(\mathbb{I}_s, \mathbb{T}, \mathbb{D}; \theta_\phi)) \tag{2}$$

where $\phi(\cdot; \theta_\phi)$ is the function that outputs the occlusion mask. Reminiscent of the image blending formulation [20], [21], note that both $w(\cdot; \theta_w)$, $f(\cdot; \theta_f)$ and $\phi(\cdot; \theta_\phi)$ take the same inputs, which shows the possibility of a shared backbone, which is better for both the knowledge sharing between all functions and computational
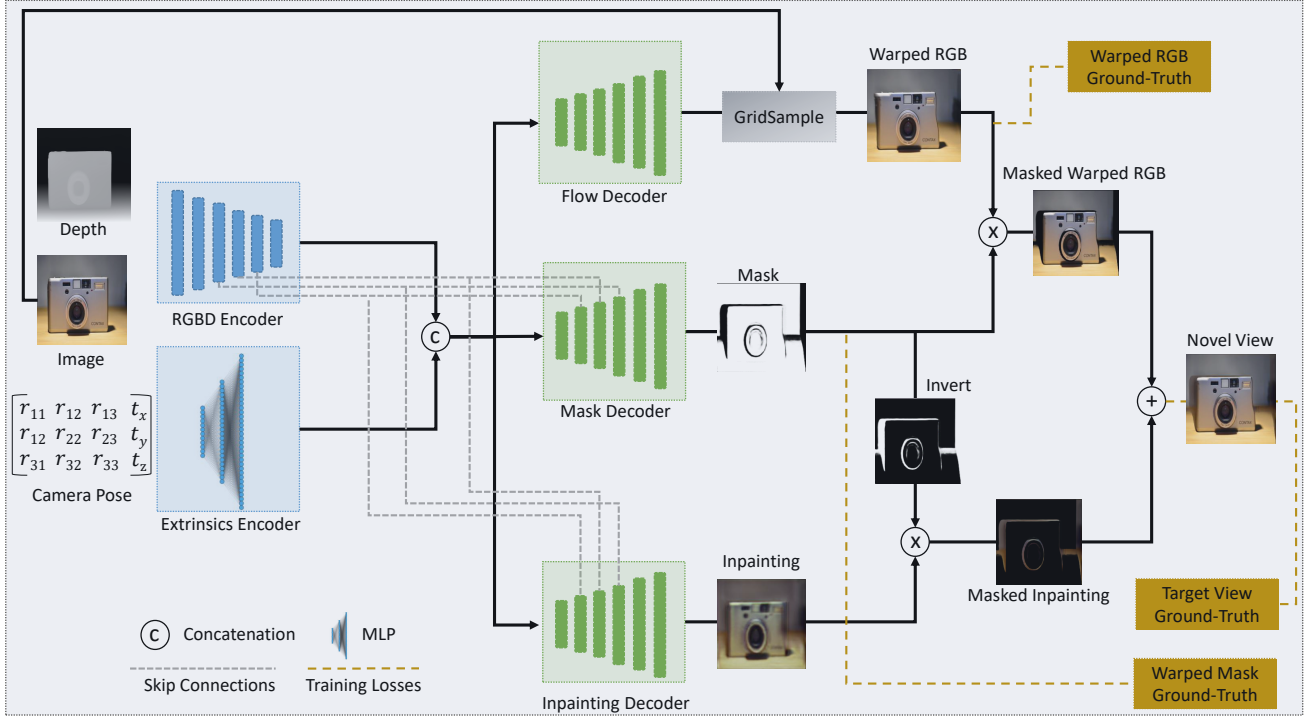
Fig. 1. Our CheapNVS architecture. CheapNVS embeds target camera pose and RGBD input information into a shared latent space, which is then used by flow and mask decoders to perform learnable warping, and inpainting decoder to fill in occluded areas. Performing inpainting and warping in parallel, as well as approximating 3D warping via flow and mask decoders make CheapNVS more computationally friendly than existing methods.

savings. Furthermore, this formulation is modular, as one can obtain the occlusion mask, warped image and inpainted textures separately, unlike the recent diffusion-based NVS solutions [2], [15].

### C. CheapNVS

CheapNVS performs single-view NVS by implementing the formulation of Equation 2. More specifically, we exploit the inherent parallelism of Equation 2 and propose to implement

$$\mathbb{M} = \phi(\mathbf{F}; \theta_\phi)$$
$$\mathbb{S} = \overline{w}(\mathbf{F}; \theta_{\overline{w}})$$
$$\mathbb{P} = f(\mathbf{F}; \theta_f)$$
$$\mathbf{F} = concat(\delta(\mathbb{I}_s, \mathbb{D}; \theta_\delta), \nabla(\mathbb{T}; \theta_\nabla))$$

(3)

where $\mathbb{S} \in \mathbb{R}^{H \times W}$, $\mathbb{P} \in \mathbb{R}^{H \times W}$, $\mathbf{F}$ and $\overline{w}(\cdot; \theta_{\overline{w}})$ are the predicted flow, inpainting, the shared embedding and the flow predictor, respectively. $\mathbf{F}$ is obtained by concatenating the outputs of $\delta(\cdot; \theta_\delta)$ and $\nabla(\cdot; \theta_\nabla)$, which encode input RGBD and the camera transformation matrix, respectively. We then synthesize the target view as

$$\mathbb{I}_t = gs(\mathbb{S}, \mathbb{I}_s) \cdot \mathbb{M} + \mathbb{P} \cdot (1 - \mathbb{M})$$

(4)

where $gs(\cdot)$ is the grid sampling function (we use PyTorch in practice) that warps the image $\mathbb{I}_s$ via the shift map $\mathbb{S}$. We implement $\delta(\cdot), \nabla(\cdot), \overline{w}(\cdot), f(\cdot)$ and $\phi(\cdot)$ with neural networks by learning their parameters $\theta_{\delta,\phi,\nabla,f,\overline{w}}$ in an end-to-end manner.

### D. The architecture

**RGBD encoder.** RGBD encoder takes in an input image and the its corresponding depth map, and maps these into a shared latent space. This latent space is shared between all three decoders, which facilitates an implicit knowledge sharing between the decoders and also saves computational budget. We use a MobileNetv2 [22] to implement the RGBD encoder.

**Extrinsics encoder.** The image warping process is conditioned on the camera pose of the target view. Although the conventional warping methods use this information [13], learnable warping methods discard this as they work for fixed baselines [8], [9], [10], [11]. Our extrinsics encoder fixes that issue by learning an embedding for the target camera view. Specifically, it takes in the transformation matrix (between $\mathbb{I}_s$ and $\mathbb{I}_t$) and produces an embedding, which is then concatenated with the output of the RGBD encoder. We implement the extrinsics encoder with a cheap, two-layer MLP, which progressively maps the 12-D input to 256-D output.

**Discussion.** The two-encoder approach facilitates parallelism, and saves compute compared to having a single, large encoder for both RGBD and camera pose inputs. Additionally, we use skip connections from the RGBD encoder to mask and inpainting decoders to facilitate better knowledge sharing. We get negative returns when we have skip connections to the flow decoder, which we discuss in Section IV-C.

**Flow decoder.** Having encoded the necessary inputs, we first focus on our flow decoder. The aim of this decoder is to take the shared embedding $\mathbf{F}$ as input and output the shift-map $\mathbb{S}$, which calculates the offset values to be applied on each pixel to perform the warping. Once the shift map is applied to the input image $\mathbb{I}_s$, we get the warped RGB image $\mathbb{I}_{warped}$, which is then multiplied by $\mathbb{M}$ to only contain the pixels coming from the source image. Note that our shift map network is conditioned on the target camera pose as the input $\mathbf{F}$ is conditioned on $\nabla(\cdot)$, unlike the stereo flow networks [8], [9], [10], [11] that completely discard it.

**Mask decoder.** The aim of the mask decoder is to generate the mask $\mathbb{M}$ necessary to perform the blending of Equation 2. Using the shared embedding $\mathbf{F}$ as input, the mask decoder predicts a binary mask that contains the spatial information to blend the warped RGB image and the result of inpainting.

**Inpainting decoder.** The aim of the inpainting decoder is to map the input features $\mathbf{F}$ to the dense inpainting output $\mathbb{P}$. This inpainting output is then multiplied with the inverse of $\mathbb{M}$, which filters out all

| | Open Images | | | | | | COCO | | | | | | Performance | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Warping | | | Inpainting | | | Warping | | | Inpainting | | | Runtime | | Memory |
| Method | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | GPU (ms)↓ | Mobile (ms)↓ | GB ↓ |
| AdaMPI [13] ‡ | - | - | - | 0.303 | 20.46 | 0.583 | - | - | - | 0.306 | 20.00 | 0.571 | 254 | N/A | 0.15 |
| Ours | **0.143** | **26.87** | **0.877** | **0.088** | **29.33** | 0.883 | **0.095** | **26.87** | **0.902** | **0.069** | 29.03 | 0.895 | 26 | 33 | **0.14** |
| AdaMPI [13] | - | - | - | 0.100 | 26.23 | **0.888** | - | - | - | 0.078 | 26.38 | **0.909** | 254 | N/A | 0.15 |
| Ours | 0.158 | 25.90 | 0.848 | 0.099 | 28.57 | 0.865 | 0.100 | 26.82 | 0.893 | 0.071 | **29.31** | 0.892 | **26** | **33** | **0.14** |

TABLE I

RESULTS ON COCO AND OPEN IMAGES. ROWS 1 AND 2 ARE MODELS TRAINED ON OPEN IMAGES, WHEREAS THE OTHERS ARE TRAINED ON COCO. OPENIMAGES AND COCO COLUMNS INDICATE EVALUATION ON OPEN IMAGES AND COCO TEST SETS. ADAMPI [13] PERFORMS THE CONVENTIONAL WARPING WE USE AS GROUND-TRUTH, SO WE DO NOT REPORT ITS WARPING RESULTS. ‡ IS ADAMPI TRAINED ON OPEN IMAGES BY US.

the shifted pixels from the input image and adds only the inpainted colours to the final output.

**Discussion.** These three decoders run concurrently and share the same architecture, which is formed of decoder blocks having bilinear upsampling, convolution and ELU activation steps. Conceptually, the warping is implemented by the flow and mask decoders, whereas the occlusion-filling is implemented by the inpainting decoder. Note that each of these decoders are working independently expect their shared encoders, which means they can also be used as isolated modules.

*E. Multi-stage training*

Despite the parallelism-friendly formulation of Equation 2, in single-view NVS, inpainting is conceptually dependant on warping as the warping supplies inpainting inputs, namely the masked warped RGB and the mask itself. From a practical point of view, that necessitates a functional warping strategy before starting to optimize the inpainting module to specialise in narrow baseline occlusion masks.

To this end, we devise a simple multi-stage training strategy, where we gradually activate new decoders in training. We first train the flow and mask decoders on their respective ground-truths, which helps the shared encoders and these two decoder supply rich information to the inpainting module. After a few epochs, we activate the training of the inpainting module, and train the entire pipeline until convergence.

## IV. EXPERIMENTAL RESULTS

*A. Experimental Details*

**Training.** We train our models with a combination of several losses. We use L1 loss for the flow decoder, cross-entropy loss for the mask decoder. For the inpainting decoder we use L1, as we observe using the losses of [13] degrade our results. The overall loss is defined as

$$\mathcal{L}_{total} = \underbrace{\lambda_1 \mathcal{L}_{L_1}}_{\mathcal{L}_{inpaint}} + \underbrace{\lambda_2 \mathcal{L}_{CE}}_{\mathcal{L}_{mask}} + \underbrace{\lambda_3 \mathcal{L}_{L_1}}_{\mathcal{L}_{flow}} \quad (5)$$

where at epoch 0, we set $\lambda_1 = 0$, $\lambda_2 = \lambda_3 = 1$. Starting from the 5th epoch, we set $\lambda_1 = \lambda_2 = \lambda_3 = 1$. This setting facilitates our multi-stage training. We train our model for 20 epochs with a learning rate of $1e-4$ and a batch size of 32. We train on random crops of size 224x224 with horizontal flip augmentation, and use Adam optimizer [23].

**Training datasets.** Following [13], we use COCO to train and use 118K images for training. We also use Open Images [24] as another training set, and randomly sample 174K images for training. We generate our ground-truths on-the-fly as [13]. For mask and flow decoders, we generate their labels through conventional 3D warping. We use the random transformation matrix generation approach of [13], which effectively defines the narrow baseline definition throughout this paper. For inpainting, we create our own pseudo-labeler; we finetune a pretrained MI-GAN [25] on warp-back data using OpenImages. This model is then used to generate inpainted images which are used as ground-truth for the inpainting decoder.

**Evaluation datasets.** In ablations, we use OpenImages to train and evaluate, and randomly select 12.5K images for testing. For

| | Warping | | | Inpainting | | |
|---|---|---|---|---|---|---|
| | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ |
| $\mathcal{L}_{inpaint}$ | 0.211 | 22.43 | 0.778 | 0.095 | **29.56** | 0.877 |
| $\mathcal{L}_{inpaint} + \mathcal{L}_{mask}$ | 0.147 | 26.55 | 0.864 | 0.090 | 29.31 | 0.878 |
| $\mathcal{L}_{total}$ | **0.143** | **26.87** | **0.877** | **0.088** | 29.33 | **0.883** |
| No SC | 0.198 | 24.91 | 0.820 | 0.109 | 28.64 | 0.869 |
| SC to all decoders | 0.153 | 25.93 | 0.820 | 0.102 | 28.10 | 0.828 |
| SC to $f(\cdot)$ and $\phi(\cdot)$ | 0.143 | 26.87 | 0.877 | 0.088 | 29.33 | 0.883 |
| L1 + FFL + Perc + SSIM | 0.154 | 25.73 | 0.827 | 0.093 | 28.54 | 0.853 |
| L1 + SSIM | 0.152 | 26.56 | 0.866 | 0.091 | 29.32 | **0.888** |
| L1 + Perc | **0.142** | 26.69 | 0.872 | 0.089 | 29.21 | 0.879 |
| L1 + FFL | 0.147 | 26.46 | 0.858 | 0.093 | 28.54 | 0.853 |
| L1 | 0.143 | **26.87** | **0.877** | **0.088** | **29.33** | 0.883 |

TABLE II

ABLATION RESULTS ON OPENIMAGES TEST SET. **SC** STANDS FOR SKIP CONNECTIONS FROM RGBD ENCODER. **FFL** AND **PERC** INDICATE THE USAGE OF FOCAL FREQUENCY [26] AND PERCEPTUAL [27] LOSSES.

comparison with state-of-the-art, we use the same randomly selected 12.5K images from OpenImages and the test set of COCO. We generate the ground-truths the same way as we do in training.

**Metrics.** We evaluate CheapNVS in warping as well as inpainting, and use SSIM, PSNR and LPIPS [27] metrics. For warping, the ground-truth warping is the conventional depth-based 3D warping implemented in PyTorch. We use DPT [28] depth maps for COCO training to be fair against AdaMPI, and Marigold [29] depth maps for Open Images training.

*B. Results*

We compare CheapNVS against AdaMPI [13], the current state-of-the-art method for 3D photography. We use its original weights, but we also train it using the original code on OpenImages for a fair comparison against our method. We omit other older methods from our comparison, as they are already inferior to AdaMPI in accuracy.

**Quantitative.** The results shown in Table I indicate that CheapNVS, except SSIM scores, is in fact better than AdaMPI in inpainting on both Open Images and COCO test sets. This shows that our lightweight decoder manages to do as well as the heavy multi-decoder approach of AdaMPI. We credit this to CheapNVS' shared latent space between multiple decoders that facilitate rich interaction between them. As AdaMPI uses the conventional 3D warping pipeline that we use to generate our ground-truths, the warping results in Table I shows how far off we are from the ground-truth. The results show that we are doing fine in approximating the laborious 3D warping with our cheap double encoders. The inpainting results also verify that, as inpainting inherently relies on the warping accuracy.

Finally, our model trained on Open Images beat our model trained on COCO and the original AdaMPI trained on COCO. We also train AdaMPI on Open Images with its original settings, but this training fails to provide a tangible result. Ultimately, the results show that Open Images as a training dataset can offer better scale, and manages to outperform COCO as a training set for narrow baseline NVS.

**Qualitative.** We provide visual results in Figure 2. CheapNVS manages to perform warping and occlusion mask prediction successfully. Specifically, it not only approximates the ground-truth predictions quite well, but manages to *smooth* the occlusion masks by itself
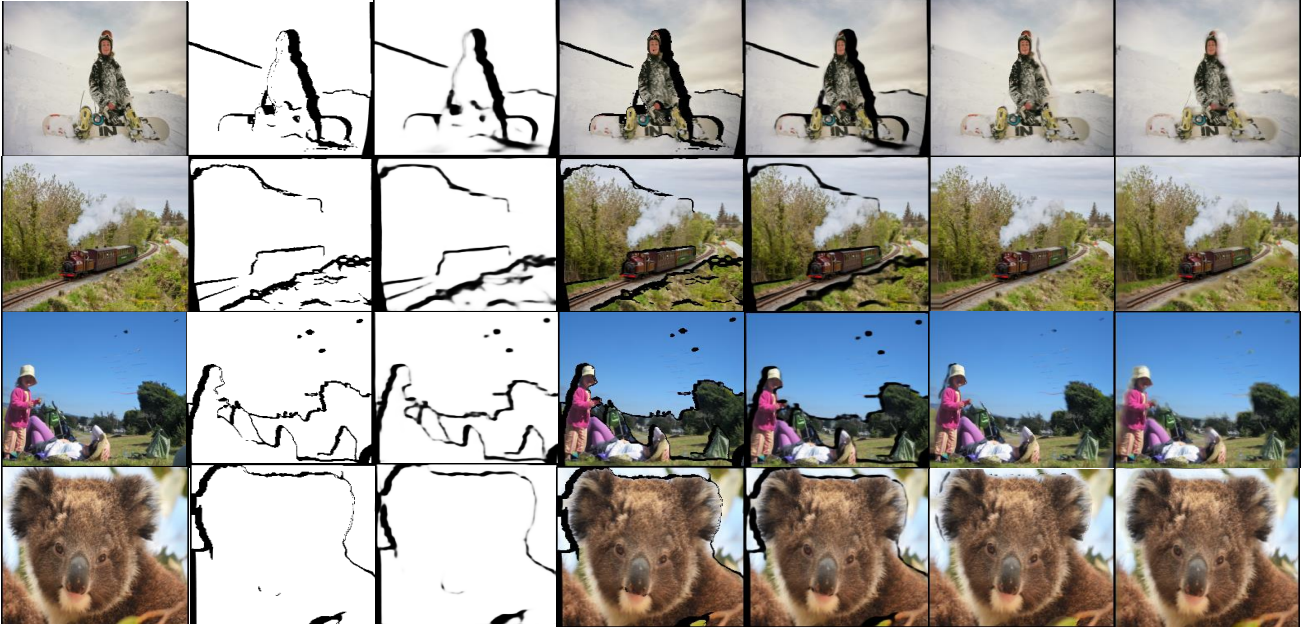
Fig. 2. Left to right: Input, ground-truth occlusion mask, CheapNVS occlusion mask, ground-truth warped RGB, CheapNVS warped RGB, AdaMPI inpainting, CheapNVS inpainting. CheapNVS manages to handle object border artefacts in inpainting (1st, 3rd and 4th rows) and approximates warping successfully.

(see first and third rows of Figure 2). Note that occlusion masks are often post-processed with operations such as smoothing, connected components and dilation/erosion in conventional methods, as these are necessary to perform an accurate blending. CheapNVS learns to do these automatically, as it jointly learns inpainting as well.

In inpainting, against much more complex competitors, CheapNVS still performs competitively. CheapNVS actually does a better job in object boundary artefact removal compared to AdaMPI (see the person on 1st and 3rd row, and Koala on the last row). In cases where there are large translations on the image boundaries, such as left edge of image on the second row, AdaMPI exhibits stretching artefacts, whereas CheapNVS avoids this behaviour.

**Runtime.** Table I shows that CheapNVS is worthy of its name and runs at 26 ms, which is 10 × faster than AdaMPI on an RTX 3090 GPU. Furthermore, CheapNVS has lower memory consumption during inference (a batch size of 1). CheapNVS runs at 33ms (via TFLite GPU delegate) on a Samsung Tab 9+, whereas AdaMPI can not be ported naively as it uses external libraries for 3D warping. Note that porting the conventional image warping efficiently would require domain expertise, yet CheapNVS can be seamlessly run on various hardware accelerators without requiring such expertise. Other methods [16] report 800ms on an IPhone with resolution 1152x1536 or 300ms on a desktop P100 GPU with 762x1008 resolution [30]. Once lifted to these resolutions, CheapNVS takes 618 and 275ms, which is still much faster than [30], [16].

### C. Ablations

**Supervision.** Our method leverages three distinct supervisory losses, one for each decoder, as shown in Equation 5. The first three rows of Table II shows that addition of each loss term helps the final warping and inpainting accuracy. Furthermore, the first row shows the possibility of learning warping via purely inpainting supervision, paving the way for weak supervised methods for learnable warping.

**Skip connections.** As Figure 1 shows, we have skip connections from RGBD encoder to inpainting and mask decoders. The second part (rows 3 to 6) of Table II shows the effect of adding such skip connections. Adding skip connections definitely help the final results both for warping and inpainting. Adding a skip connection to flow

decoder does the opposite and degrades the results. We believe the reason is the nature of each decoders' output; inpainting and mask decoders predict outputs that are spatially-meaningful and relevant to the content of the input image. The flow decoder, however, predicts a shift-map that has no relation with the other two decoders' output. This discrepancy, we believe, is the reason why skip connections to the flow decoder harm the results.

**Losses.** AdaMPI uses L1, focal frequency [26], perceptual [27] and SSIM losses for training. We explore the loss design space, and observe that we obtain better quantitative results by using only the L1 loss. As shown in the last rows of Table II, L1 loss training obtains best results in warping and inpainting, with close second being L1 + perceptual loss. We do not observe tangible qualitative differences between the results of training with or without the perceptual loss.

### D. Limitations and Future Work

A common failure mode of CheapNVS is in warping, which ultimately effects the quality of the synthesized view. Similarly, since we learn the warping on a narrow baseline specified during the training, CheapNVS might struggle to warp in baselines that are larger than what it is trained on, whereas AdaMPI might do well since it performs the highly expensive, handcrafted warping. We believe this to be the reason why we outperform AdaMPI in narrow baselines despite its use of complex multi-plane images and multiple inpainting networks, whereas they might do better than us in larger baselines (such as some stereo datasets). We note, however, that these larger baselines are not within the scope of our work. In future work, we aim to alleviate this warping domain gap, and also improve synthesis quality by using a more powerful inpainting teacher.

### V. CONCLUSION

We propose CheapNVS, a narrow-baseline single-view novel view synthesis method that comfortably runs on real-time on mobile devices. CheapNVS approximates the laborious depth-based 3D image warping by efficient modules that learn to perform warping conditioned on the target view. Furthermore, we cast single-view NVS as an image blending problem, and perform warping and inpainting concurrently to save computational resources. Trained in multiple

stages and on a diverse, large-scale dataset, CheapNVS outperforms the existing state-of-the-art method while performing up to 10 × faster and consuming 6% less memory.

## REFERENCES

[1] Michal Adamkiewicz, Timothy Chen, Adam Caccavale, Rachel Gardner, Preston Culbertson, Jeannette Bohg, and Mac Schwager, "Vision-only robot navigation in a neural radiance world," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4606–4613, 2022.

[2] Sixu Li, Chaojian Li, Wenbo Zhu, Boyang Yu, Yang Zhao, Cheng Wan, Haoran You, Huihong Shi, and Yingyan Lin, "Instant-3d: Instant neural radiance field training towards on-device ar/vr 3d reconstruction," in *Proceedings of the 50th Annual International Symposium on Computer Architecture*, 2023, pp. 1–13.

[3] Evangelos Skartados, Mehmet Kerim Yucel, Bruno Manganelli, Anastasios Drosou, and Albert Saà-Garriga, "Finding waldo: Towards efficient exploration of nerf scene spaces," in *Proceedings of the 15th ACM Multimedia Systems Conference*, 2024, pp. 155–165.

[4] Abril Corona-Figueroa, Jonathan Frawley, Sam Bond-Taylor, Sarath Bethapudi, Hubert PH Shum, and Chris G Willcocks, "Mednerf: Medical neural radiance fields for reconstructing 3d-aware ct-projections from a single x-ray," in *2022 44th annual international conference of the IEEE engineering in medicine & Biology society (EMBC)*. IEEE, 2022, pp. 3843–3848.

[5] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[6] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis, "3d gaussian splatting for real-time radiance field rendering.," *ACM Trans. Graph.*, vol. 42, no. 4, pp. 139–1, 2023.

[7] Zhiwen Fan, Panwang Pan, Peihao Wang, Yifan Jiang, Hanwen Jiang, Dejia Xu, Zehao Zhu, Dilin Wang, and Zhangyang Wang, "Pose-free generalizable rendering transformer," *arXiv e-prints*, pp. arXiv–2310, 2023.

[8] Yang Zhou, Hanjie Wu, Wenxi Liu, Zheng Xiong, Jing Qin, and Shengfeng He, "Single-view view synthesis with self-rectified pseudo-stereo," *International Journal of Computer Vision*, vol. 131, no. 8, pp. 2032–2043, 2023.

[9] Yi-Nan Chen, Hang Dai, and Yong Ding, "Pseudo-stereo for monocular 3d object detection in autonomous driving," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 887–897.

[10] Simon Evain and Christine Guillemot, "A lightweight neural network for monocular view generation with occlusion handling," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 6, pp. 1832–1844, 2019.

[11] Gaurav Chaurasia, Arthur Nieuwoudt, Alexandru-Eugen Ichim, Richard Szeliski, and Alexander Sorkine-Hornung, "Passthrough+ real-time stereoscopic view synthesis for mobile mixed reality," *Proceedings of the ACM on Computer Graphics and Interactive Techniques*, vol. 3, no. 1, pp. 1–17, 2020.

[12] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang, "3d photography using context-aware layered depth inpainting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8028–8038.

[13] Yuxuan Han, Ruicheng Wang, and Jiaolong Yang, "Single-view view synthesis in the wild with learned adaptive multiplane images," in *ACM SIGGRAPH 2022 Conference Proceedings*, 2022, pp. 1–8.

[14] Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen, Chong Zeng, and Hao Su, "Zero123++: a single image to consistent multi-view diffusion base model," *arXiv preprint arXiv:2310.15110*, 2023.

[15] Shitao Tang, Jiacheng Chen, Dilin Wang, Chengzhou Tang, Fuyang Zhang, Yuchen Fan, Vikas Chandra, Yasutaka Furukawa, and Rakesh Ranjan, "Mvdiffusion++: A dense high-resolution multi-view diffusion model for single or sparse-view 3d object reconstruction," *arXiv preprint arXiv:2402.12712*, 2024.

[16] Johannes Kopf, Kevin Matzen, Suhib Alsisan, Ocean Quigley, Francis Ge, Yangming Chong, Josh Patterson, Jan-Michael Frahm, Shu Wu, Matthew Yu, et al., "One shot 3d photography," *ACM Transactions on Graphics (TOG)*, vol. 39, no. 4, pp. 76–1, 2020.

[17] Johannes Lutz Schönberger and Jan-Michael Frahm, "Structure-from-motion revisited," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[18] Vincent Jantet, Luce Morin, and Christine Guillemot, "Incremental-ldi for multi-view coding," in *2009 3DTV Conference: The True Vision-Capture, Transmission and Display of 3D Video*. IEEE, 2009, pp. 1–4.

[19] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker, "Deepview: View synthesis with learned gradient descent," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2367–2376.

[20] Konstantinos Georgiadis, Albert Saà-Garriga, Mehmet Kerim Yucel, Anastasios Drosou, and Bruno Manganelli, "Adaptive mask-based pyramid network for realistic bokeh rendering," in *European Conference on Computer Vision*. Springer, 2022, pp. 429–444.

[21] Mehmet Kerim Yücel, Valia Dimaridou, Bruno Manganelli, Mete Ozay, Anastasios Drosou, and Albert Saa-Garriga, "Lra&ldra: Rethinking residual predictions for efficient shadow detection and removal," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 4925–4935.

[22] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.

[23] Diederik Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, San Diega, CA, USA, 2015.

[24] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, et al., "The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale," *International journal of computer vision*, vol. 128, no. 7, pp. 1956–1981, 2020.

[25] Andranik Sargsyan, Shant Navasardyan, Xingqian Xu, and Humphrey Shi, "Mi-gan: A simple baseline for image inpainting on mobile devices," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 7335–7345.

[26] Liming Jiang, Bo Dai, Wayne Wu, and Chen Change Loy, "Focal frequency loss for image reconstruction and synthesis," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 13919–13929.

[27] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.

[28] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun, "Vision transformers for dense prediction," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 12179–12188.

[29] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Konrad Schindler, "Repurposing diffusion-based image generators for monocular depth estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9492–9502.

[30] Varun Jampani, Huiwen Chang, Kyle Sargent, Abhishek Kar, Richard Tucker, Michael Krainin, Dominik Kaeser, William T Freeman, David Salesin, Brian Curless, et al., "Slide: Single image 3d photography with soft layering and depth-aware inpainting," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 12518–12527.