# Improved Cross-view Completion Pre-training for Stereo Matching

Philippe Weinzaepfel        Vaibhav Arora        Yohann Cabon

Thomas Lucas        Romain Brégier        Vincent Leroy        Gabriela Csurka

Leonid Antsfeld        Boris Chidlovskii        Jérôme Revaud

NAVER LABS Europe

## Abstract

*Despite impressive performance for high-level downstream tasks, self-supervised pre-training methods have not yet fully delivered on dense geometric vision tasks such as stereo matching. The application of self-supervised learning concepts, such as instance discrimination or masked image modeling, to geometric tasks is an active area of research. In this work we build on the recent cross-view completion framework: this variation of masked image modeling leverages a second view from the same scene, which is well suited for binocular downstream tasks. However, the applicability of this concept has so far been limited in at least two ways: (a) by the difficulty of collecting real-world image pairs – in practice only synthetic data had been used – and (b) by the lack of generalization of vanilla transformers to dense downstream tasks for which relative position is more meaningful than absolute position. We explore three avenues of improvement: first, we introduce a method to collect suitable real-world image pairs at large scale. Second, we experiment with relative positional embeddings and demonstrate that they enable vision transformers to perform substantially better. Third, we scale up vision transformer based cross-completion architectures, which is made possible by the use of large amounts of data. With these improvements, we show for the first time that state-of-the-art results on deep stereo matching can be reached without using any standard task-specific techniques like correlation volume, iterative estimation or multi-scale reasoning.*

## 1. Introduction

Self-supervised pre-training methods aim at learning rich representations from large amounts of unannotated data, which can then be finetuned on a variety of downstream tasks. This requires the design of pretext tasks, for which supervision signal can be extracted from the data itself, as well as architectures that can be transferred easily. We hypothesize that successfully pre-training a stereo matching model requires three things together: (a) large-scale real-world data, (b) a well-designed dense pretext task inciting the understanding of 3D scene layout and geometry, and (c) an architecture that can process pairs of images and is robust to cropping and resolution changes.

Early self-supervised methods proceeded by discarding part of the signal (*e.g.*, image color [72], patch ordering [43] or image orientation [21]) and trying to recover it. Later methods based on instance discrimination [9, 10, 13, 24] were the first to surpass supervised pre-training on high-level tasks; they are based on the idea that output features should be invariant to well-designed classes of augmentations. This is typically achieved with a contrastive objective, using positive pairs obtained from two augmentations of the same image. Another recently successful pretext task is masked image modeling (MIM) [2, 18, 23, 63, 66, 74], where part of the input data is masked and an auto-encoder tries to restore the full signal from the remaining visible parts. When pre-trained on balanced and object-centric datasets like ImageNet [49], instance discrimination and MIM methods have achieved excellent performance on semantic tasks such as image classification, in particular with limited amounts of annotated data [2, 14, 56], but have not led to breakthroughs in more geometry-based tasks like stereo matching.

Adapting self-supervised pre-training to geometric vision tasks is an active area of research. Attempts have been made to design contrastive learning objectives at the pixel or patch level [62, 65, 67], but their performance gains have so far been more moderate than for global tasks. Besides, these gains are mainly demonstrated for *semantic* tasks such as semantic segmentation or object detection, rather than for geometric tasks such as depth estimation or stereo matching. Recently, [64] proposed the pretext task of cross-view completion (or CroCo in short), a variant of MIM where a partially masked input image must be reconstructed given visible patches and an additional view of the same scene. This pre-training objective is well suited to geometric stereo
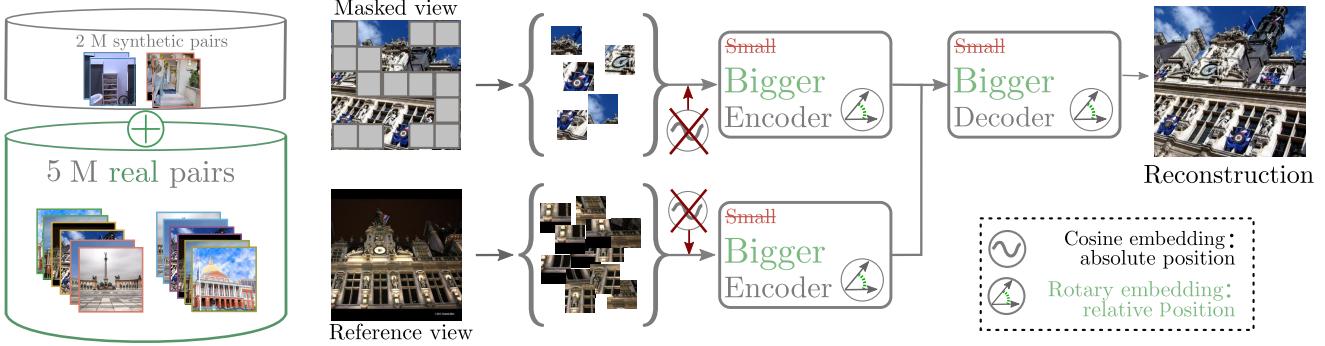
Figure 1. **Overview of the improvements we propose for cross-view completion pre-training:** (a) collecting and using real-world images, (b) using rotary positional embeddings which model *relative* token positions, instead of absolute positions using the standard cosine embedding, (c) increasing network size both in the encoder and the decoder.

downstream tasks as (a) it leverages pairs of images and (b) extracting relevant information from the second view requires geometric understanding of the scene. The CroCo architecture consists of a vision transformer (ViT) [17] encoder to extract features for the non-masked tokens of the first image, as well as for the second reference image, and a transformer to decode the features and reconstruct the masked image, as illustrated in Figure 1.

In spite of these advances, leveraging cross-view completion for geometric vision tasks remains challenging for at least two reasons. First, training with cross-view completion requires image pairs depicting the same scene; this can be hard to acquire at scale, yet scale is the cornerstone of the success of self-supervised pre-training. In practice, the CroCo model of [64] is pre-trained solely with synthetic data, which may thus limit its final performance. Second, most models trained with masking rely on ViTs [17], which typically use absolute positional embeddings. These do not generalize well to new image resolutions when finetuning, and are not always robust to cropping. This limits the applicability of current cross-view completion methods and may explain why the downstream tasks presented in [64] mostly use low-resolution squared images.

In this paper, we propose solutions to these limitations that enable to pre-train a large-scale cross-view completion model, see Figure 1, leading to state-of-the-art performance on stereo matching. First, we tackle the problem of scalable pair collection, and gather millions of training pairs from different real-world datasets which cover various scenarios like indoor environments, street view data and landmarks, see Figure 2. In order to generate high-quality pre-training pairs, we carefully control the visual overlap for each pair of images. Indeed, pairs with high overlap make the task trivial, whereas pairs with negligible overlap reduce it to standard MIM [64]. To measure this overlap, we leverage extra information available such as 3D meshes, additional sensors like LIDAR, or Structure-from-Motion (SfM) reconstructions for datasets with sufficient image coverage. Based on these data, we can generate a set of high qual-

ity image pairs – with sufficient overlap and viewpoint difference – while also ensuring high diversity between pairs. Second, these large-scale datasets of pre-training pairs allow to scale up the model: (a) we use a larger encoder to extract better image features and (b) also scale up the decoder, which is responsible for combining information coming from the two views. Third, instead of using the standard cosine positional embedding which encodes absolute positional information, we propose to rely on the Rotary Positional Embedding (RoPE) [58] which efficiently injects relative positional information of token pairs during attention.

We finetune our pre-trained model with this improved cross-view completion scheme on stereo matching using a Dense Prediction Transformer (DPT) [47] head. Our model, termed CroCo-Stereo, is simple and generic: we rely on a plain ViT encoder, followed by a plain transformer decoder which directly predicts the disparity through the DPT head in a single forward pass. For the first time, our experiments demonstrate that a model whose architecture differs significantly from current state-of-the-art approaches, all of which being based on convolutional feature extractors and domain-specific bells and whistles (*e.g.* cost volumes [25,28,29,70], iterative refinement [33,36] and multi-level feature pyramids [15, 33]), can reach state-of-the-art on various stereo benchmarks such as KITTI 2015 [41] or ETH3D [53].

## 2. Related work

**Self-supervised learning.** The success of instance discrimination [9, 10, 13, 22, 24] has drawn a lot of attention to self-supervised learning in computer vision [26]. In that paradigm, variants of an image are obtained by applying different data augmentations. Features extracted from the different variants are trained to be similar, while being pushed away from features obtained from other images. Such self-supervised models are particularly well tailored to image-level tasks, such as image classification, and have led to state-of-the-art performance on various benchmarks.
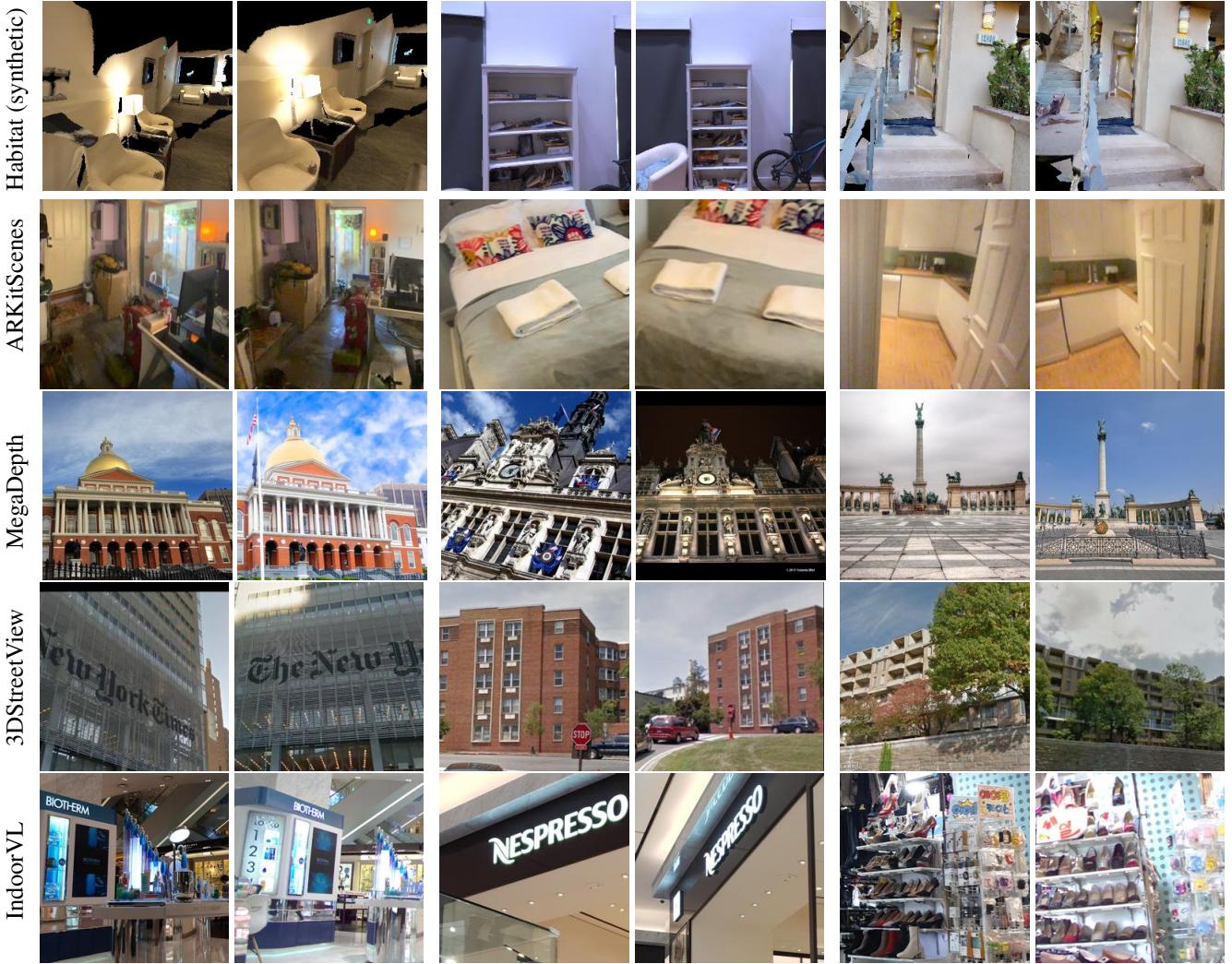
Figure 2. **Example of pre-training cropped image pairs** from Habitat which was the synthetic data used by CroCo [64] on the top row, and from real-world datasets we use in this paper (from ARKitScenes, MegaDepth, 3DStreetView and IndoorVL) below.

Recent studies suggest that this success could be due to the object-centric [44] and the balanced [1] nature of ImageNet [49] that is used for pre-training. Recently, inspired by BERT [16] in natural language processing, different masked modeling methods have been adapted to computer vision. MIM pre-training aims at reconstructing masked information from an input image either in the pixel space [3, 4, 12, 18, 23, 66], or in the feature space [2, 5, 63], and sometimes after quantization [6, 74]. Recent works combine this framework in a teacher-student approach [32, 34] with improved masking strategy [19, 27, 34]. Overall, MIM models perform well on classification tasks. They have obtained some success on denser tasks such as object detection [23] or human pose estimation [69], and have been applied to robotic vision [45] when pre-trained on related datasets. More recently, CroCo [64] introduces the pretext task of cross-view completion, where a second view of the same scene is added to the MIM framework. This is well suited

to geometric downstream tasks: to leverage the second view and improve reconstruction accuracy, the model has to implicitly be aware of the geometry of the scene. The CroCo model outperforms MIM pre-training on an array of geometric downstream tasks; however, it relies on synthetic data only, which may be sub-optimal, and does not reach the performance of task-specific methods.

**Positional embeddings.** Since a ViT treats its input as an orderless set of image patches, or tokens, positional embeddings are a necessary tool to keep track of the position of each patch token from the original image. They can be either learned [10, 17] or handcrafted, such as the cosine positional embeddings from the original transformer [61]. Cosine embeddings performs similarly in practice and are now more commonly used for ViTs, since learned embeddings do not generalize well when varying the image resolutions (*i.e.*, the number or tokens). Both learned and cosine embeddings are added explicitly to the signal and con-

tain absolute positional information. However, pixel-level dense computer vision tasks should be able to process various image resolutions and be robust to cropping. Thus, relative positional embeddings, *e.g.* [54], that consider distances between tokens, might be preferable. For instance, Bello *et al.* [8] achieve better object detection using relative self-attention. Similarly, Swin Transformers [38] and Swin v2 [37] observed improved performance using relative positional embeddings. More recently, [58] introduced the Rotary Positional Embedding (RoPE): a transformation to each key and query features is applied according to their absolute position, but in such a way that the pairwise similarity scores used in the attention computation only depend on the relative positions of the token pairs and on their feature similarity. Thus, RoPE is well suited to our purpose: it models relative positions, which can be computed at any resolution, and avoids polluting the signal with position information.

**Deep stereo matching.** The seminal work on end-to-end networks of Mayer *et al.* [40] has shown the benefit of leveraging correlation for disparity estimation, compared to simply applying convolutions on concatenated frames. Since then, most methods are based on cost/correlation volumes, see [30] for a survey. Image features are extracted from both images before a raw cost volume is built. The cost volumes typically have three dimensions [25, 29, 70], the third dimension representing a discretization of the disparity level. Some recent methods use a fourth dimension in the cost volume [11, 28, 42], by keeping the feature dimensionality for later stages. This cost volume can be further regularized with different priors before disparity is estimated and optionally some post-processing steps can be applied. Recent trends in deep stereo matching include iterative refinement and multi-level interaction. RAFT-Stereo [36] adapted the paradigm proposed by RAFT [60] for optical flow estimation with an iterative refinement scheme for the stereo matching task. Hierarchical neural architecture search has been introduced in LEAStereo [15]. Lately, CREStereo [33] proposed a hierarchical network with recurrent refinement to update disparities in a coarse-to-fine manner.

## 3. Cross-view completion pre-training at scale

Our proposed pre-training method is based on the recently introduced cross-view completion (CroCo) framework [64]. It extends masked image modeling to pairs of images. Given two different images depicting a given scene, the two images are divided into sets of non-overlapping patches, denoted as tokens, and 90% of the tokens from the first image are masked. The remaining ones are fed to a ViT [17] encoder to extract features for the first image. Similarly, tokens from the second image are fed to the same encoder with shared weights, and a ViT decoder processes the two sets of features together to reconstruct the target. Fig-

ure 1 provides an overview of the pre-training stage. Compared to standard masked image modeling methods, this approach can leverage the information in the second view to resolve some of the ambiguities about the masked context. To leverage this information, the model has to implicitly reason about the scene geometry and spatial relationship between the two views, which primes it well for geometric downstream tasks and stereo matching in particular.

**Training data.** Collecting pairs of images that are suitable for this approach is highly non-trivial. First, images have to be paired together without manual annotation; second, their visual overlap has to be carefully controlled for the pairs to be useful. In [64], only synthetic data generated with the Habitat simulator [50] is used, which restricts the variety of the pre-training data. In contrast, we propose an approach to this real-world image pairing problem, necessary to use cross-view completion at scale, as detailed in Section 3.1.

**Positional embeddings.** The architecture used in [64] adapts vision transformers to process pairs of images, by using cross-attention layers inside the decoder. Following standard practice, in their work cosine positional embedding is added to the token features prior to the encoder and the decoder. However, this models absolute position while dense tasks such as stereo matching must be robust to cropping or images of various resolutions. In Section 3.2, we describe how relative positional embeddings can be adapted to cross-view completion.

**Large-scale models.** Finally, we discuss scaling-up the model in Section 3.3. CroCo [64] uses a ViT-Base encoder (12 blocks, 768-dimension features, 12 attention heads) and a decoder composed of 8 blocks with 512-dimensional features and 16 heads. Using our large-scale dataset of real-world image pairs, we are able to scale to larger ViT architectures and demonstrate consistent gains.

### 3.1. Collecting real-world image pairs

We now present our approach to automatically select image pairs from real-world datasets that are suitable for pre-training. To be useful, pairs need to depict the same scene with some partial overlap. The overlap should not be small to the point where the task boils down to auto-completion. It should not be high either to the point where the task becomes a trivial 'copy and paste', *i.e.*, without requiring any understanding of scene geometry. On top of that, diversity should be as high as possible among pairs. We propose to use datasets which offer some way to get information about the geometry of the scene and the camera poses. This information can be captured using additional sensors like LIDAR, or it can be extracted using structure-from-motion (SfM) techniques if the images offer enough coverage of the scene. We use this information to build an image pair quality score based on overlap and difference in viewpoint angle. We then use a greedy algorithm to select a diverse
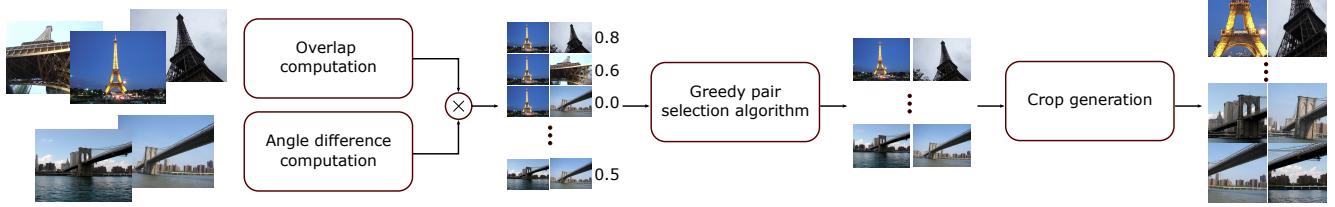
Figure 3. **Overview of our pre-training cropped image pair collection method.** Given a dataset of posed images, optionally with point clouds (*e.g.* from SfM) or meshes of the scene, we first measure the visual overlap between pairs and the viewpoint angle difference. Based on these scores, we use a greedy algorithm to select diverse image pairs and finally generate crops from them.

set of image pairs. Finally, we generate overlapping image crops by leveraging image matching. Figure 3 gives an overview of our approach and we detail each step below.

**Computing overlap scores.** The first step is to compute overlap scores for candidate pairs. We develop several approaches depending on the available information.

○ *ARKitScenes* [7] provides 450,000 frames from 1,667 different indoor environments. The availability of the corresponding mesh for each frame enables the computation of the overlap between every pair of images. For each image $I$, we retrieve the set of mesh vertices $\mathcal{P}(I)$ that are visible. We then measure the intersection-over-union (IoU) of the vertices (3D points) for each pair of images $(I_1, I_2)$ as:

$$IoU(I_1, I_2) = \frac{|\mathcal{P}(I_1) \cap \mathcal{P}(I_2)|}{|\mathcal{P}(I_1) \cup \mathcal{P}(I_2)|}. \quad (1)$$

○ *MegaDepth* [35] consists of around 300,000 images downloaded from the web corresponding to 200 different landmarks. For the images of the MegaDepth dataset, a point cloud model for each landmark obtained using structure-form-motion (SfM) with COLMAP [52] is also provided. As above, it is possible to measure the vertex-based IoU between pairs of images, where each vertex is in this case a 3D point from the point cloud. Unfortunately, occlusions cannot be taken into account due to the absence of 3D mesh, which greatly degrades the overlap estimation. We propose a simple yet effective solution: we create an artificial occlusion model by attaching a ball of fixed radius to each 3D point, which occludes the vertices placed behind it. This way, we can compute a set of visible vertices for each image and evaluate the IoU as done previously.

○ *3D Street View* [71] contains 25 million street view images from 8 cities. In addition to the camera pose, the 3D location and orientation (normal vector) of the target buildings are provided. To compute the overlap score, we create a pseudo 3D point cloud and apply the same technique as for MegaDepth. We start from an empty point cloud and append, for each target building, a of $10 \times 6$ meters grid of $7 \times 11$ balls oriented according to the provided annotation.

○ *Indoor Visual Localization datasets* (IndoorVL) [31] contains over 135,000 images from a large shopping mall and a large metro station in Seoul, South Korea, captured regularly with several months interval with 10 cameras and

2 laser scanners. The data is provided with accurate camera poses obtained via LiDAR SLAM refined by SfM based optimization. Due to the LIDAR noise caused by heavy reflections on the floor or on windows, and thanks to the accurate camera poses, we directly measure the overlap between images using the intersection between the camera frustrums. To encourage further diversity, we multiply this score by a factor $0.8$ if both images come from the same capture session, thus favoring pairs that are taken with several months interval.

**Greedy image pair selection.** We rely on the overlap scores described above to select high quality pairs. This is however not sufficient: we also need pairs to be diverse, which would not be the case when randomly selecting good pairs, as images in the dataset can be very correlated. Therefore, we use a greedy algorithm to select non-redundant image pairs for pre-training. First, for each image pair $(I_1, I_2)$ we use a quality pair score $s$ given by:

$$s(I_1, I_2) = IoU(I_1, I_2) \times 4\cos(\alpha)\big(1 - \cos(\alpha)\big), \quad (2)$$

where $\alpha$ denotes the viewpoint angle difference between the two images (note that all the datasets above provide camera poses). The function $4\cos(x)(1 - \cos(x))$ has a maximum value of 1 for $x = 60°$, 0 value for $x = 0°$ and $x = 90°$, and it is negative for angles larger than $90°$. Thus this score favors image pairs that have different viewpoints while still having large overlaps. Given the score for every pair, we aim at building a large number of image pairs while ensuring diversity, *i.e.*, avoiding content redundancy. To do this, we use a greedy algorithm, where each time we select a pair of images with maximum score, we discard the two images forming it, as well as images that have too large IoU (above 0.75) with any of the two. We then repeat this process iteratively until we have no image pairs with a score above a certain threshold.

**Crop generation per pair.** For pre-training, we use fixed size crops of size $224 \times 224$, as considering higher resolution images would significantly increase the pre-training cost. To generate crops on pairs of images while maintaining overlaps, we rely on quasi-dense keypoint matching. More precisely, we use DeepMatching [48], except for pairs from ARKitScenes where we directly use matches obtained from the mesh. Given the matches, we consider a grid of

5

crops in the first image, estimate the correspond matching crop in the second image and keep the ones that do not overlap and have the most consistent matches.

**Overall statistics.** In total, we collected about 5.3 million real-world pairs of crops with the process described above, with respectively 1,070,414 pairs from ARKitScenes [7], 2,014,789 pairs from MegaDepth [35], 655,464 from 3DStreetView [71], and 1,593,689 pairs from IndoorVL [31]. We added this to 1,821,391 synthetic pairs generated with the Habitat simulator [50], following the approach of [64]. Example pairs for each dataset are shown in Figure 2. They cover various scenarios, from indoor rooms – synthetic with Habitat or real with ARKitScenes – to larger crowded indoor environment (IndoorVL), landmarks (MegaDepth) and outdoor streets (3DStreetView).

### 3.2. Positional embeddings

We replace the cosine embeddings, which inject absolute positional information, by Rotary Positional Embedding (or RoPE) [58]. RoPE efficiently injects information about the *relative* positioning of feature pairs when computing attention. Formally, let $q$ and $k$ represent a query and a key feature, at absolute positions $m$ and $n$ respectively. The main idea of RoPE is to design an efficient function $f(x, p)$ that transforms a feature $x$ according to its absolute position $p$ such that the similarity between the transformed query and the transformed key $\langle f(q, m), f(k, n) \rangle$ is a function of $q$, $k$ and $m - n$ only. [58] showed that a simple transformation such as applying rotations on pairs of dimensions according to a series of rotation matrices at different frequencies satisfy this desirable property. To deal with 2D signals such as images, we split the features into 2 parts, we apply the positional embedding of the x-dimension on the first part, and the embedding of the y-dimension on the second part.

### 3.3. Scaling up the model

The combination of information extracted from the two images only occurs in the decoder. Following MAE [23], CroCo [64] uses a small decoder of 8 blocks consisting of self-attention, cross-attention and an MLP, with 512 dimensions and 16 attention heads. As the decoder is crucial for binocular tasks such as stereo matching, we scale up the decoder and follow the ViT-Base hyper-parameters with 12 blocks, 768-dimensional features and 12 heads. We also scale up the image encoder from ViT-Base to ViT-Large, *i.e.*, increase the depth from 12 to 24, the feature dimension from 768 to 1024 and the number of heads from 12 to 16.

**Pre-training detailed setting.** We pre-train the network for 100 epochs with the AdamW optimizer [39], a weight decay of 0.05, a cosine learning rate schedule at a base learning rate of $3.10^{-4}$ with a linear warmup in the first 10 epochs, a batch size of 512 spread on 8 GPUs. During pre-training, we simply use color jittering and random swapping between the first and the second image as data augmentation. We mask 90% of the tokens from the first image.

## 4. Application to stereo matching

We now present CroCo-Stereo, our ViT-based correlation-free architecture for stereo matching that leverages cross-view completion pre-training. This is much in contrast to current state-of-the-art methods which rely on task-specific design in the form of cost volumes [25, 28, 29, 57, 68, 70, 73], iterative refinement [33, 36] and multi-level feature pyramids [15, 33, 59].

**Architecture.** When finetuning the model for stereo matching, both images are fed to the encoder as during pre-training (but without masking), and the decoder processes the tokens of both images together. To output a pixel-wise prediction, we rely on the DPT module [47], which adapts the standard up-convolutions and fusions from multiple layers used in fully-convolutional approaches for dense tasks, to vision transformers. This allows to combine features from different blocks of the transformers by reshaping them to different resolutions and fusing them with convolutional layers. In practice, we use the features from the penultimate layer of the encoder of the first image, as well as features from 3 blocks regularly spread in the decoder, *e.g.* the output of blocks 4, 8 and 12 for a decoder with 12 blocks.

**Stereo matching training data.** For training, we use $352 \times 704$ crops from various stereo datasets. We use 196,000 pairs from the training set of the synthetic CREStereo dataset [33], a subset of 13,014 synthetic training pairs from SceneFlow [40] (SF) in its final rendering version, 24 training pairs from ETH3D [53], 223 pairs from Booster [46], and 763 image pairs from different versions of Middlebury [51] (51 pairs from the 2005 version, 186 pairs from 2006, 150 pairs from 2014 and 376 pairs from 2021). To better balance different datasets, we duplicate the pairs from ETH3D 30 times, and the pairs from Middlebury 50 times. During training, prior to cropping, we rescale the images by a factor randomly chosen in the range of $[2^{-0.2}, 2^{0.4}]$, we apply asymmetric color jittering, add some grey rectangles and some shift and rotation to the right image only.

**Training.** We finetune a model pre-trained with our improved cross-view completion. We use an L1 loss between the network prediction and the ground-truth disparity. We train the model for 16 epochs on the data listed above using batches of 6 pairs. We use the AdamW optimizer [39] with a weight decay of 0.05, a cosine learning rate schedule with a single warm-up epoch and a learning rate of $3.10^{-5}$.

**Inference.** We use a tiling-based approach. We sample overlapping tiles of $352 \times 704$ pixels in the first image, with 70% overlap along each dimension, *i.e.*, a stride of 105 and 211 along the y- and the x-dimension respectively. For each tile, we create a pair by sampling a corresponding tile at the same position from the second image. We then predict
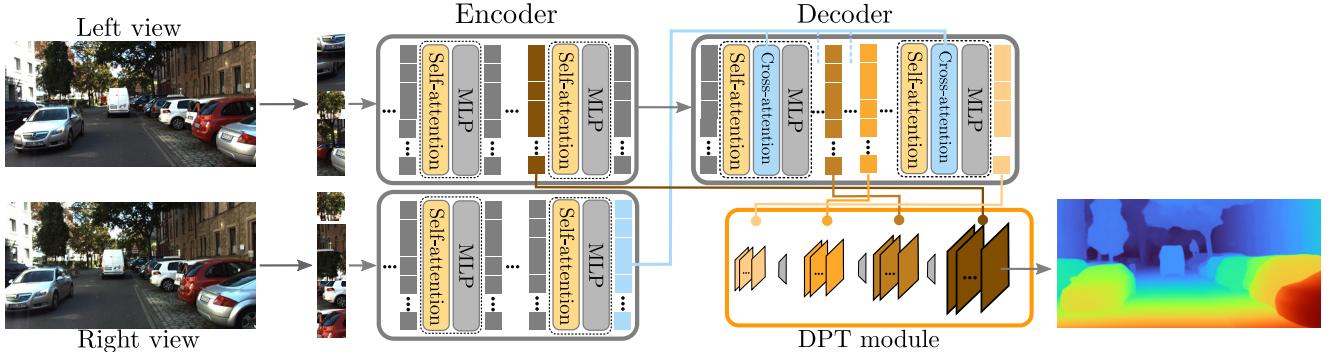
Figure 4. **Overview of the downstream CroCo-Stereo architecture.** The two images are split into patches and encoded with a series of transformer blocks that use the RoPE positional embeddings. The decoder consists in a series of transformer decoder blocks (self-attention among token features from the first image, cross-attention with the token features from the second image, and an MLP). Token features from different intermediate blocks are fed to the DPT module [47] which outputs the final disparity prediction.
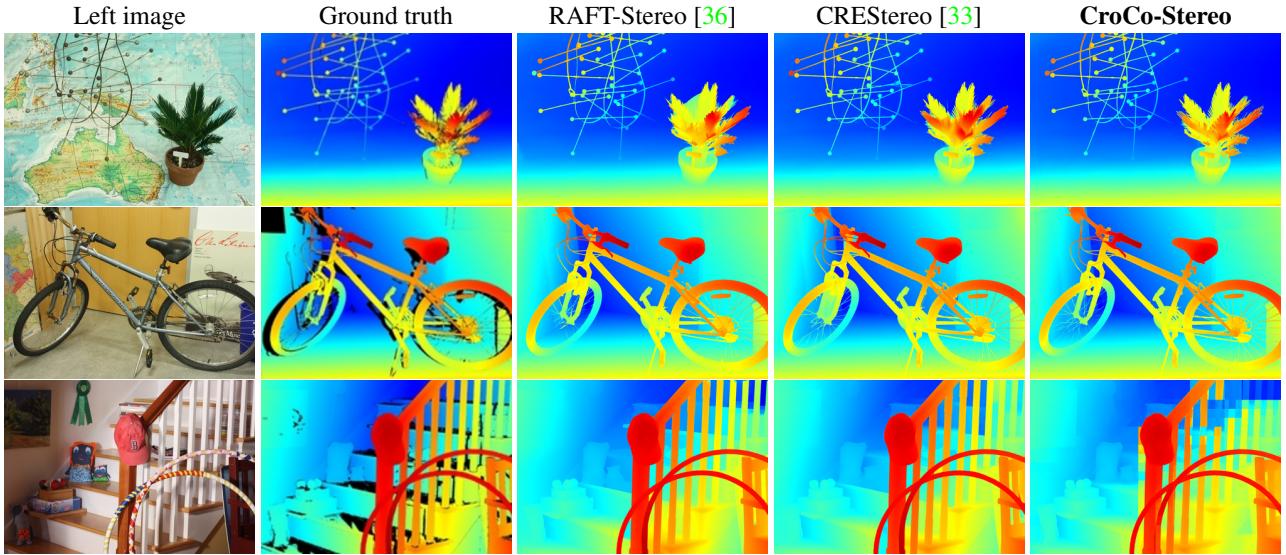


Figure 5. **Three example results from the Middlebury test set** (Australia, Bicycle2 and Hoops) with from left to right: the left image, the ground-truth disparity, RAFTStereo [36], CREStereo [33] and CroCo-Stereo.

## 5. Experiments

**Ablations.** We perform our ablations on 3 validation pairs from Middlebury 2006, 2014 and 2021 and 50 validation pairs from SceneFlow in its clean rendering, that we have removed from our training set. We first ablate in Table 1 (first five rows) the changes we propose to improve CroCo (pre-training data, positional embedding, larger encoder and decoder) by reporting the average disparity error for different models. We observe that they all lead to consistent improvements: starting from CroCo [64], replacing the cosine absolute positional embedding by RoPE reduces the mean error from 1.444 to 1.297, scaling up the decoder allows to go down to 1.201, using larger-scale pre-training data and

a larger encoder leads to the best performance with 1.086. Finally, we measure the impact of pre-training by training the model directly for stereo matching (last row); this leads to poor performance with a mean average error of 2.378.

**Comparison to the state of the art.** We now evaluate CroCo-Stereo on the official leaderboards of Middlebury, KITTI 2015 and ETH3D. For KITTI, we additionally fine-tuned the model on the 394 KITTI training pairs (194 from KITTI 2012 [20] and 200 from KITTI 2015) for 500 epochs.

On Middlebury (Table 2), CroCo-Stereo sets a new state of the art on 7 out of 15 sequences, in spite of using a generic patch-based transformer without any of the usual apparatus for stereo matching (*e.g.* cost-volume, coarse-to-scale processing, iterative refinement). CREStereo [33], in comparison, performs best on 4 sequences, but has a lower average error due to the fact that CroCo-Stereo produces really large errors for a few image pairs. In fact, these er-

| pos. embed | encoder | decoder | pre-training data | | MD14 | MD21 | MD06 | SFclean | mean |
|---|---|---|---|---|---|---|---|---|---|
| cosine | ViT-Base | Small | 2M habitat | (= CroCo [64]) | 0.723 | 1.629 | 1.530 | 1.893 | 1.444 |
| RoPE | ViT-Base | Small | 2M habitat | | 0.607 | 1.371 | 1.404 | 1.804 | 1.297 |
| RoPE | ViT-Base | Base | 2M habitat | | 0.532 | 1.265 | 1.327 | 1.681 | 1.201 |
| RoPE | ViT-Base | Base | 2M habitat + 5.3M real | | 0.495 | 1.202 | **1.300** | **1.606** | 1.151 |
| RoPE | ViT-Large | Base | 2M habitat + 5.3M real | (= **CroCo-Stereo**) | **0.435** | **0.977** | 1.306 | 1.625 | **1.086** |
| RoPE | ViT-Large | Base | none | | 1.369 | 3.190 | 2.109 | 2.937 | 2.378 |

Table 1. **Ablative study** of each change between CroCo and CroCo-Stereo measured as the average disparity error on the low-resolution validation sets from Middlebury 2014 (MD14), 2021 (MD21), 2006 (MD06), as well as SceneFlow (SFclean) in its clean rendering, and the mean over these validation sets. A *Small* decoder consists of 8 decoder blocks with 16 attention heads on 512-dimensional features, while the *Base* one has 12 blocks with 12 heads on 768-dimensional features. The last row represents CroCo-Stereo without any pre-training.

| Method nd < 400px | Bicyc2 ✓ | Compu ✓ | Austr ✓ | AustrP ✓ | Djemb ✓ | DjembL ✓ | Livgrm ✓ | Plants ✓ | Hoops | Stairs | Nkuba | Class | ClassE | Crusa | CrusaP | **avg** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AdaStereo [57] | 7.20 | 3.37 | 16.7 | 13.4 | 3.83 | 4.14 | 7.82 | 19.4 | 11.0 | 10.5 | 13.8 | 8.73 | 8.75 | 13.2 | 11.2 | 10.2 |
| HITNet [59] | 6.18 | 4.33 | 15.7 | 15.1 | 2.63 | 21.5 | 6.12 | 10.4 | 8.91 | 11.6 | 16.2 | 6.72 | 8.88 | 11.3 | 12.3 | 9.97 |
| RAFT-Stereo [36] | 4.66 | 2.91 | 12.1 | 12.8 | **1.75** | 3.14 | 5.83 | 11.1 | 8.79 | 9.11 | 17.9 | **4.96** | **5.97** | 10.6 | 13.0 | 8.41 |
| LEAStereo [15] | 6.42 | 4.19 | 13.1 | 12.7 | 2.28 | 3.32 | 7.27 | 11.7 | 10.2 | 7.92 | **12.7** | 7.60 | 7.37 | 7.93 | 7.75 | 8.11 |
| CREStereo [33] | 5.79 | 3.15 | 13.4 | 13.5 | 1.87 | **2.93** | 4.88 | 13.4 | **8.77** | 8.47 | 14.6 | 5.92 | 6.34 | **6.58** | **6.62** | **7.70** |
| **CroCo-Stereo** | **4.50** | **2.60** | **9.78** | **10.1** | 1.81 | 6.25 | **4.31** | **7.72** | 28.9 | **6.08** | 21.5 | 19.6 | 20.4 | 36.0 | 39.1 | 14.6 |

Table 2. **Evaluation on Middlebury** with the RMSE over non-occluded pixels for each sequence and the average (last column). Sequences are ordered according to their 'nd' value, which is the official threshold of maximum disparity used to clip predictions before evaluation.

| Method | D1-bg | D1-fg | D1-all |
|---|---|---|---|
| AdaStereo [57] | 2.59 | 5.55 | 3.08 |
| HITNet [59] | 1.74 | 3.20 | 1.98 |
| ACVNet [68] | **1.37** | 3.07 | **1.65** |
| LEAStereo [15] | 1.40 | 2.91 | **1.65** |
| CREStreo [33] | 1.45 | 2.86 | 1.69 |
| PCWNet [55] | **1.37** | 3.16 | 1.67 |
| **CroCo-Stereo** | 1.89 | **2.72** | 2.03 |

Table 3. **Evaluation on the KITTI 2015 benchmark** with the percentage of outliers (*i.e.*, error above 3 pixels) for background (D1-bg), foreground (D1-fg) and all (D1-all) pixels.

| Method | Bad 1.0 | | RMSE | | 99% quantile | |
|---|---|---|---|---|---|---|
| | noc | all | noc | all | noc | all |
| AdaStereo [57] | 3.09 | 3.34 | 0.44 | 0.48 | 1.45 | 2.22 |
| HITNet [59] | 2.79 | 3.11 | 0.46 | 0.55 | 2.07 | 2.49 |
| RAFT-Stereo [36] | 2.44 | 2.60 | 0.36 | 0.42 | 1.20 | 2.13 |
| DIP-Stereo [73] | 1.97 | 2.12 | 0.37 | 0.43 | 1.24 | 2.03 |
| CREStereo [33] | **0.98** | **1.09** | **0.28** | **0.31** | **0.92** | **1.00** |
| **CroCo-Stereo** | 1.05 | 1.18 | 0.32 | 0.37 | 1.05 | 1.14 |

Table 4. **Evaluation on ETH3D** with the percentage of pixels with an error over 1px (Bad 1.0), the RMSE and the error at 99% quantile evaluated over non-occluded (noc) or all pixels.

rors correspond to sequences with large maximum disparities (based on the maximum threshold value applied before evaluation), which is harmful for our simple tiling-based inference approach. This effect is visible in the prediction of the bottom example of Figure 5 where one can observe tiling artefacts, *e.g.* next to the stair pillars. In general, however, our method remains very accurate, especially on thin structures like the pins on the map or the radius of the bicycle wheels in Figure 5.

On KITTI 2015 (Table 3), CroCo-Stereo performs the best on foreground pixels according to the benchmark metric (outliers ratio at a 3px error threshold). Foreground areas correspond to pixels where cross-attention of the decoder can be leveraged effectively. CroCo-Stereo ranks slightly lower on background pixels and thus on average. Again, we do not benefit on that aspect from any task-specific design, *e.g.* iterative refinement or multi-scale reasoning.

Finally, on ETH3D, CroCo-Stereo performs second, just below CREStereo [33], see Table 4. It outperforms recent approaches like RAFT-Stereo [36], DIP-Stereo [73] or HIT-Net [59] by a large margin, *e.g.* decreasing the ratio of outliers at 1px to 1.05%, 1 to 2% below these approaches.

**Limitations.** CroCo-Stereo is a large model with about 447M parameters (320M for the ViT-Large encoder, 115M for the Base decoder and 30M for the DPT head), *i.e.* significantly more weights than state-of-the-art models. Note however that in this work we did not aim at efficiency in terms of model size, but at showing that CroCo-Stereo pre-training can benefit to stereo matching. A second limitation is our inference strategy with tiling. The boundary of the tiles can be clearly visible in some cases, especially with high disparity values where the prediction might become totally wrong as the matching pixel in the second image can be outside of the tile. A smarter tiling strategy than taking the same cropping coordinates in a pair of images could be considered, and we plan to improve this in future work.

# 6. Conclusion

For the first time, we have shown that large-scale pre-training can be successful for a dense 3D geometric task like stereo matching thanks to a well-adapted pretext task and real-world data at scale. This enables to reach state-of-the-art performance with a ViT-based architecture without using task-specific designs, thereby opening novel routes to tackle the problem.

# References

[1] Mahmoud Assran, Randall Balestriero, Quentin Duval, Florian Bordes, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, and Nicolas Ballas. The hidden uniform cluster prior in self-supervised learning. *arXiv preprint arXiv:2210.07277*, 2022. 3

[2] Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Florian Bordes, Pascal Vincent, Armand Joulin, Michael Rabbat, and Nicolas Ballas. Masked siamese networks for label-efficient learning. In *ECCV*, 2022. 1, 3

[3] Sara Atito, Muhammad Awais, and Josef Kittler. SiT: Self-supervised vIsion Transformer. *arXiv preprint arXiv:2104.03602*, 2021. 3

[4] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. MultiMAE: Multi-modal Multi-task Masked Autoencoders. In *ECCV*, 2022. 3

[5] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. data2vec: A General Framework for Self-supervised Learning in Speech, Vision and Language. *arXiv preprint arXiv:2202.03555*, 2022. 3

[6] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEiT: BERT Pre-Training of Image Transformers. In *ICLR*, 2022. 3

[7] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, et al. Arkitscenes–a diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data. In *NeurIPS*, 2021. 5, 6

[8] Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V. Le. Attention Augmented Convolutional Networks. In *ICCV*, 2019. 4

[9] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments. In *NeurIPS*, 2020. 1, 2

[10] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging Properties in Self-Supervised Vision Transformers. In *ICCV*, 2021. 1, 2, 3

[11] J. Chang and Y. Chen. Pyramid stereo matching network. In *CVPR*, 2018. 4

[12] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative Pretraining From Pixels. In *ICML*, 2020. 3

[13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 1, 2

[14] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E. Hinton. Big self-supervised models are strong semi-supervised learners. In *NeurIPS*, 2020. 1

[15] Xuelian Cheng, Yiran Zhong, Mehrtash Harandi, Yuchao Dai, Xiaojun Chang, Hongdong Li, Tom Drummond, and Zongyuan Ge. Hierarchical neural architecture search for deep stereo matching. In *NeurIPS*, 2020. 2, 4, 6, 8

[16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL HLT*, 2019. 3

[17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *ICLR*, 2021. 2, 3, 4

[18] Alaaeldin El-Nouby, Gautier Izacard, Hugo Touvron, Ivan Laptev, Hervé Jegou, and Edouard Grave. Are Large-scale Datasets Necessary for Self-Supervised Pre-training? *arXiv preprint arXiv:2112.10740*, 2021. 1, 3

[19] Yuxin Fang, Li Dong, Hangbo Bao, Xinggang Wang, and Furu Wei. Corrupted Image Modeling for Self-Supervised Visual Pre-Training. *arXiv preprint arXiv:2202.03382*, 2022. 3

[20] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 7

[21] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised Representation Learning by Predicting Image Rotations. In *ICLR*, 2018. 1

[22] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Ávila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap Your Own Latent - A New Approach to Self-Supervised Learning. In *NeurIPS*, 2020. 2

[23] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked Autoencoders are Scalable Vision Learners. In *CVPR*, 2022. 1, 3, 6

[24] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum Contrast for Unsupervised Visual Representation Learning. In *CVPR*, 2020. 1, 2

[25] Zequn Jie, Pengfei Wang, Yonggen Ling, Bo Zhao, Yunchao Wei, Jiashi Feng, and Wei Liu. Left-right comparative recurrent model for stereo matching. In *CVPR*, 2018. 2, 4, 6

[26] Longlong Jing and Yingli Tian. Self-supervised Visual Feature Learning with Deep Neural Networks: A Survey. *IEEE Trans. PAMI*, 2021. 2

[27] Ioannis Kakogeorgiou, Spyros Gidaris, Bill Psomas, Yannis Avrithis, Andrei Bursuc, Konstantinos Karantzalos, and Nikos Komodakis. What to Hide from Your Students: Attention-Guided Masked Image Modeling. In *ECCV*, 2022. 3

[28] A. Kendall, H. Martirosyan, S. Dasgupta, P. Henry, R. Kennedy, A. Bachrach, and A. Bry. End-to-end learning of geometry and context for deep stereo regression. In *ICCV*, 2017. 2, 4, 6

[29] Sameh Khamis, Sean Fanello, Christoph Rhemann, Adarsh Kowdle, Julien Valentin, and Shahram Izadi. Stereonet: Guided hierarchical refinement for real-time edge-aware depth prediction. In *ECCV*, 2018. 2, 4, 6

[30] Hamid Laga, Laurent Valentin Jospin, Farid Boussaid, and Mohammed Bennamoun. A survey on deep learning techniques for stereo-based depth estimation. *IEEE Trans. PAMI*, 2020. 4

[31] Donghwan Lee, Soohyun Ryu, Suyong Yeon, Yonghan Lee, Deokhwa Kim, Cheolho Han, Yohann Cabon, Philippe Weinzaepfel, Nicolas Guérin, Gabriela Csurka, et al. Large-scale localization datasets in crowded indoor spaces. In *CVPR*, 2021. 5, 6

[32] Youngwan Lee, Jeffrey Willette, Jonghee Kim, Juho Lee, and Sung Ju Hwang. Exploring the role of mean teachers in self-supervised masked auto-encoders. *arXiv preprint arXiv:2210.02077*, 2022. 3

[33] Jiankun Li, Peisen Wang, Pengfei Xiong, Tao Cai, Ziwei Yan, Lei Yang, Jiangyu Liu, Haoqiang Fan, and Shuaicheng Liu. Practical stereo matching via cascaded recurrent network with adaptive correlation. In *CVPR*, 2022. 2, 4, 6, 7, 8

[34] Zhaowen Li, Zhiyang Chen, Fan Yang, Wei Li, Yousong Zhu, Chaoyang Zhao, Rui Deng, Liwei Wu, Rui Zhao, Ming Tang, and Jinqiao Wang. MST: Masked Self-Supervised Transformer for Visual Representation. In *NeurIPS*, 2021. 3

[35] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *CVPR*, 2018. 5, 6

[36] Lahav Lipson, Zachary Teed, and Jia Deng. Raft-stereo: Multilevel recurrent field transforms for stereo matching. In *3DV*, 2021. 2, 4, 6, 7, 8

[37] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *CVPR*, 2022. 4

[38] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 4

[39] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *ICLR*, 2019. 6

[40] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, 2016. 4, 6

[41] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *CVPR*, 2015. 2

[42] Guang-Yu Nie, Ming-Ming Cheng, Yun Liu, Zhengfa Liang, Deng-Ping Fan, Yue Liu, and Yongtian Wang. Multi-level context ultra-aggregation for stereo matching. In *CVPR*, 2019. 4

[43] Mehdi Noroozi and Paolo Favaro. Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles. In *ECCV*, 2016. 1

[44] Senthil Purushwalkam and Abhinav Gupta. Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases. In *NeurIPS*, 2020. 3

[45] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. Real-world robot learning with masked visual pre-training. *CoRL*, 2022. 3

[46] Pierluigi Zama Ramirez, Fabio Tosi, Matteo Poggi, Samuele Salti, Stefano Mattoccia, and Luigi Di Stefano. Open challenges in deep stereo: the booster dataset. In *CVPR*, 2022. 6

[47] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *ICCV*, 2021. 2, 6, 7

[48] Jerome Revaud, Philippe Weinzaepfel, Zaid Harchaoui, and Cordelia Schmid. Deepmatching: Hierarchical deformable dense matching. *IJCV*, 2016. 5

[49] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. 1, 3

[50] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A Platform for Embodied AI Research. In *ICCV*, 2019. 4, 6

[51] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nešić, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *GCPR*, 2014. 6

[52] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 5

[53] Thomas Schops, Johannes L Schonberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *CVPR*, 2017. 2, 6

[54] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with Relative Position Representations. In *NAACL HLT*, 2018. 4

[55] Zhelun Shen, Yuchao Dai21, Xibin Song11, Zhibo Rao, Dingfu Zhou, and Liangjun Zhang. Pcw-net: Pyramid combination and warping cost volume for stereo matching. In *ECCV*, 2022. 8

[56] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. In *NeurIPS*, 2020. 1

[57] Xiao Song, Guorun Yang, Xinge Zhu, Hui Zhou, Zhe Wang, and Jianping Shi. Adastereo: a simple and efficient approach for adaptive stereo matching. In *CVPR*, 2021. 6, 8

[58] Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021. 2, 4, 6

[59] Vladimir Tankovich, Christian Hane, Yinda Zhang, Adarsh Kowdle, Sean Fanello, and Sofien Bouaziz. Hitnet: Hierarchical iterative tile refinement network for real-time stereo matching. In *CVPR*, 2021. 6, 8

[60] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020. 4

[61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017. 3

[62] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense Contrastive Learning for Self-Supervised Visual Pre-Training. In *CVPR*, 2021. 1

[63] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked Feature Prediction for Self-Supervised Visual Pre-Training. In *CVPR*, 2022. 1, 3

[64] Weinzaepfel, Philippe and Leroy, Vincent and Lucas, Thomas and Brégier, Romain and Cabon, Yohann and Arora, Vaibhav and Antsfeld, Leonid and Chidlovskii, Boris and Csurka, Gabriela and Revaud Jérôme. CroCo: Self-Supervised Pre-training for 3D Vision Tasks by Cross-View Completion. In *NeurIPS*, 2022. 1, 2, 3, 4, 6, 7, 8

[65] Zhenda Xie, Yutong Lin, Zheng Zhang, Yue Cao, Stephen Lin, and Han Hu. Propagate Yourself: Exploring Pixel-Level Consistency for Unsupervised Visual Representation Learning. In *CVPR*, 2021. 1

[66] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. SimMIM: A Simple Framework for Masked Image Modeling. In *CVPR*, 2022. 1, 3

[67] Yuwen Xiong, Mengye Ren, and Raquel Urtasun. Loco: Local contrastive representation learning. In *NeurIPS*, 2020. 1

[68] Gangwei Xu, Junda Cheng, Peng Guo, and Xin Yang. Attention concatenation volume for accurate and efficient stereo matching. In *CVPR*, 2022. 6, 8

[69] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vitpose: Simple vision transformer baselines for human pose estimation. *arXiv preprint arXiv:2204.12484*, 2022. 3

[70] Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for match density estimation. In *CVPR*, 2019. 2, 4, 6

[71] Amir R Zamir, Tilman Wekel, Pulkit Agrawal, Colin Wei, Jitendra Malik, and Silvio Savarese. Generic 3D representation via pose estimation and matching. In *ECCV*, 2016. 5, 6

[72] Richard Zhang, Phillip Isola, and Alexei A. Efros. Colorful Image Colorization. In *ECCV*, 2016. 1

[73] Zihua Zheng, Ni Nie, Zhi Ling, Pengfei Xiong, Jiangyu Liu, Hao Wang, and Jiankun Li. Dip: Deep inverse patchmatch for high-resolution optical flow. In *CVPR*, 2022. 6, 8

[74] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. iBOT: Image BERT Pre-training with Online Tokenizer. In *ICLR*, 2022. 1, 3