

Pixel-Perfect Structure-from-Motion with Featuremetric Refinement

Philipp Lindenberger^{1*} Paul-Edouard Sarlin^{2*} Viktor Larsson² Marc Pollefeys^{2,3}

Departments of ¹Mathematics and ²Computer Science, ETH Zurich ³Microsoft

Abstract

Finding local features that are repeatable across multiple views is a cornerstone of sparse 3D reconstruction. The classical image matching paradigm detects keypoints per-image once and for all, which can yield poorly-localized features and propagate large errors to the final geometry. In this paper, we refine two key steps of structure-from-motion by a direct alignment of low-level image information from multiple views: we first adjust the initial keypoint locations prior to any geometric estimation, and subsequently refine points and camera poses as a post-processing. This refinement is robust to large detection noise and appearance changes, as it optimizes a featuremetric error based on dense features predicted by a neural network. This significantly improves the accuracy of camera poses and scene geometry for a wide range of keypoint detectors, challenging viewing conditions, and off-the-shelf deep features. Our system easily scales to large image collections, enabling pixel-perfect crowdsourced localization at scale. Our code is publicly available at github.com/cvg/pixel-perfect-sfm as an add-on to the popular SfM software COLMAP.

1. Introduction

Mapping the world is an important requirement for spatial intelligence applications in augmented reality or robotics. Tasks like visual localization or path planning can benefit from accurate sparse or dense 3D reconstructions of the environment. These can be built from images using Structure-from-Motion (SfM), which associates observations across views to estimate camera parameters and 3D scene geometry. Sparse reconstruction based on matching local image features [10, 21, 23, 34, 51, 57, 59, 65] is the most common due to its scalability and its robustness to appearance changes introduced by varying devices, viewpoints, and temporal conditions found in crowdsourced scenarios [2, 29, 35, 41, 47, 50, 58].

SfM assumes that sparse interest points [10, 21, 23, 34, 51, 59, 62, 84, 92] can be reliably detected across views. It typi-

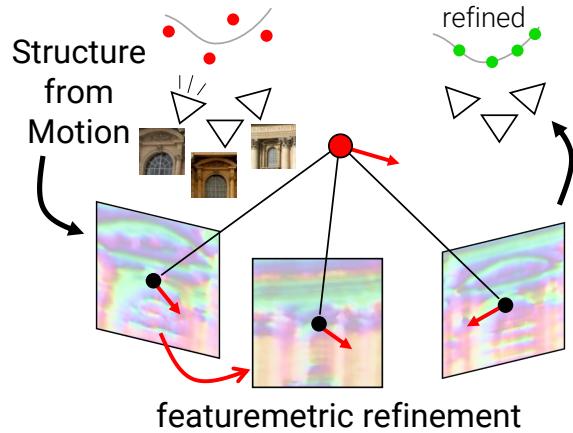


Figure 1: **From sparse to dense.** We improve the accuracy of sparse Structure-from-Motion by refining 2D keypoints, camera poses, and 3D points using the direct alignment of deep features. This featuremetric optimization leverages dense image information but can scale to scenes with thousands of images. Such refinement results in subpixel-accurate reconstructions, even in challenging conditions.

cally selects such points for each image independently and relies on these initial detections for the remainder of the reconstruction process. However, detecting keypoints from a single view is inherently inaccurate due to appearance changes and discrete image sampling [31]. The advent of convolutional neural network (CNNs) for detection has magnified this issue, as they generally do not retain local image information and instead favor global context.

Multi-view geometric optimization with bundle adjustment [4, 42, 82] is commonly used to refine cameras and points using reprojection errors. Dusmanu *et al.* [24] proposed to refine keypoint locations prior to SfM via an analogous geometric cost constrained with local optical flow. This can improve SfM, but has limited accuracy and scalability.

In this work, we argue that local image information is valuable throughout the SfM process to improve its accuracy. We adjust both keypoints and bundles, before and after reconstruction, by direct image alignment [18, 26, 52] in a learned feature space. Exploiting this locally-dense informa-

*indicates equal contributions

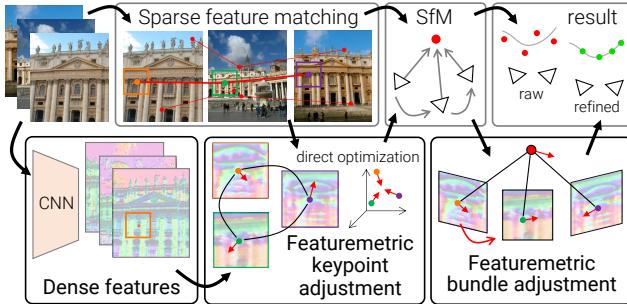


Figure 2: Refinement pipeline. Our refinement works on top of any SfM pipeline that is based on local features. We perform a two-stage adjustment of keypoints and bundles. The approach first refines the 2D keypoints only from tentative matches by optimizing a direct cost over dense feature maps. The second stage operates after SfM and refines 3D points and poses with a similar featuremetric cost.

tion is significantly more accurate than geometric optimization, while deep, high-dimensional features extracted by a CNN ensure wider convergence in challenging conditions. This formulation elegantly combines globally-discriminative sparse matching with locally-accurate dense details. It is applicable to both incremental [70, 75] and global [9, 12, 54] SfM irrespective of the types of sparse or dense features.

We validate our approach in experiments evaluating the accuracy of both 3D structure and camera poses in various conditions. We demonstrate drastic improvements for multiple hand-crafted and learned local features using off-the-shelf CNNs. The resulting system produces accurate reconstructions and scales well to large scenes with thousands of images. In the context of visual localization, it can, in addition to providing a more accurate map, also refine poses of single query images with minimal overhead.

For the benefit of the research community, we will release our code as an extension to COLMAP [70, 71] and to the popular localization toolbox hloc [63, 64]. We believe that our featuremetric refinement can significantly improve the accuracy of existing datasets [67] and push the community towards sub-pixel accurate localization at large scale.

2. Related work

Image matching is at the core of SfM and visual SLAM, which typically rely on sparse local features for their efficiency and robustness. The process i) detects a small number of interest points, ii) computes their visual descriptors, iii) matches them with a nearest neighbor search, and iv) verifies the matches with two-view epipolar estimation and RANSAC. The correspondences then serve for relative or absolute pose estimation and 3D triangulation. As keypoints are sparse, small inaccuracies in their locations can result in large errors for the estimated geometric quantities.

Differently, dense matching [13, 49, 61, 74, 77, 81, 83] considers all pixels in each image, resulting in denser and more accurate correspondences. It has been successful for constrained settings like optical flow [40, 76] or stereo depth estimation [90], but is not suitable for large-scale SfM due to its high computational cost due to many redundant correspondences. Several recent works [46, 60, 78, 96] improve the matching efficiency by first matching coarsely and subsequently refining correspondences using a local search. This is however limited to image pairs and thus cannot create point tracks required by SfM.

Our work combines the best of both paradigms by leveraging dense local information to refine sparse observations. It is inherently amenable to SfM as it can optimize all locations over multiple views in a track simultaneously.

Subpixel estimation is a well-studied problem in correspondence search. Common approaches either upsample the input images or fit polynomials or Gaussian distributions to local image neighborhoods [28, 36, 39, 51, 69]. With the widespread interest in CNNs for local features, solutions tailored to 2D heatmaps have been recently developed, such as learning fine local sub-heatmaps [38] or estimating subpixel corrections with regression [14, 80] or the soft-argmax [55, 93]. Cleaner heatmaps can also arise from aggregating predictions over multiple virtual views using data augmentation [21].

Detections or local affine frames can be combined across multiple views with known poses in a least-squares geometric optimization [25, 82]. Dusmanu *et al.* [24] instead refine keypoints solely based on tentative matches, without assuming known geometry. This geometric formulation exhibits remarkable robustness, but is based on a local optical flow whose estimation for each correspondence is expensive and approximate. We unify both keypoint and bundle optimizations into a joint framework that optimizes a featuremetric cost, resulting in more accurate geometries and a more efficient keypoint refinement.

Direct alignment optimizes differences in pixel intensities by implicitly defining correspondences through the motion and geometry. It therefore does not suffer from geometric noise and is naturally subpixel accurate via image interpolation. Direct photometric optimization has been successfully applied to optical flow [8, 52], visual odometry [18, 26, 27, 44], SLAM [5, 72], multi-view stereo (MVS) [19, 22, 91], and pose refinement [73]. It generally fails for moderate displacements or appearances changes, and is thus not suitable for large-baseline SfM. One notable work by Woodford & Rossten [88] refines dense SfM+MVS models with a robust image normalization. It focuses on dense mapping with accurate initial poses and moderate appearance changes. Georgel *et al.* [30] instead estimate more accurate relative poses by elegantly combining photometric and geometric costs. They show that dense information can improve sparse estimation

but their approach ignores appearance changes. Differently, our work improves the entire SfM pipeline starting with tentative matches and addresses larger, challenging changes.

To improve on the weaknesses of photometric optimization, numerous recent works align multi-dimensional image representations. Examples of this *featuremetric* optimization include frame tracking with handcrafted [6, 56] or learned descriptors [17, 53, 86, 87, 89], optical flow [7, 11], MVS [94], and dense SfM in small scenes [79]. Closer to our work, PixLoc [66] learns deep features with a large basin of convergence for wide-baseline pose refinement. It improves the accuracy of sparse matching but is designed for single images and disregards the scalability to multiple images or large scenes. Here we extend this paradigm to other steps of SfM and propose an efficient algorithm that scales to thousands of images. We show that learning task-specific wide-context features is not necessary and demonstrate highly accurate refinements with off-the-shelf features.

In conclusion, our work is the first to apply robust featuremetric optimization to a large-scale sparse reconstruction problem and show significant benefits for visual localization.

3. Background

Given N images $\{\mathbf{I}_i\}$ observing a scene, we are interested in accurately estimating its 3D structure, represented as sparse points $\{\mathbf{P}_j \in \mathbb{R}^3\}$, intrinsic parameters $\{\mathbf{C}_i\}$ of the cameras, and the poses $\{(\mathbf{R}_i, \mathbf{t}_i) \in \text{SE}(3)\}$ of the images, represented as rotation matrices and translation vectors.

A typical SfM pipeline performs geometric estimation from correspondences between sparse 2D keypoints $\{\mathbf{p}_u\}$ observing the same 3D point from different views, collectively called a track. Association between observations is based on matching local image descriptors $\{\mathbf{d}_u \in \mathbb{R}^D\}$, but the estimated geometry relies solely on the location of the keypoints, whose accuracy is thus critical. Keypoints are detected from local image information for each image individually, without considering multiple views simultaneously. Subsequent steps of the pipeline discover additional information about the scene, such as its geometry or its multi-view appearance. Two approaches leverage this information to reduce the detection noise and refine the keypoints.

Global refinement: Bundle adjustment [82] is the gold standard for refining structure and poses given initial estimates. It minimizes the total geometric error

$$E_{\text{BA}} = \sum_j \sum_{(i,u) \in \mathcal{T}(j)} \|\Pi(\mathbf{R}_i \mathbf{P}_j + \mathbf{t}_i, \mathbf{C}_i) - \mathbf{p}_u\|_\gamma , \quad (1)$$

where $\mathcal{T}(j)$ is the set the images and keypoints in track j , $\Pi(\cdot)$ projects to the image plane, and $\|\cdot\|_\gamma$ is a robust norm [33]. This formulation implicitly refines the keypoints while ensuring their geometric consistency. It however ignores the uncertainty of the initial detections and thus re-

quires many observations to reduce the geometric noise. Operating on an existing reconstruction, it cannot recover observations arising from noisy keypoints that are matched correctly but discarded by the geometric verification.

Track refinement: To improve the accuracy of the keypoints prior to any geometric 3D estimation, Dusmanu *et al.* [24] optimize their locations over tentative tracks formed by raw, unverified matches. They exploit the inherent structure of the matching graph to discard incorrect matches without relying on geometric constraints. Given two-view dense flow fields $\{\mathbf{T}_{v \rightarrow u}\}$ between the neighborhoods of matching keypoints u and v , this *keypoint adjustment* optimizes, for each tentative track j , the multi-view cost

$$E_{\text{KA}}^j = \sum_{(u,v) \in \mathcal{M}(j)} \|\mathbf{p}_v + \mathbf{T}_{v \rightarrow u}[\mathbf{p}_v] - \mathbf{p}_u\|_\gamma , \quad (2)$$

where $\mathcal{M}(i)$ denotes the set of matches that forms the track and $[\cdot]$ is a lookup with subpixel interpolation. A deep neural network is trained to regress the flow of a single point from two input patches and the flow field is interpolated from a sparse grid. This dramatically improves the keypoint accuracy, but some errors remain as the regression and interpolation are only approximate.

Both bundle and keypoint adjustments are based on geometric observations, namely keypoint locations and flow, but do not account for their respective uncertainties. They thus require a large number of observations to average out the geometric noise and their accuracy is in practice limited.

4. Approach

Summarizing dense image information into sparse points is necessary to perform global data association and optimization at scale. However, refining geometry is an inherently local operation, which, we show, can efficiently benefit from locally-dense pixels. Given constraints provided by coarse but global correspondences or initial 3D geometry, the dense information only needs to be locally accurate and invariant but not globally discriminative. While SfM typically discards image information as early as possible, we instead exploit it in several steps of the process thanks to direct alignment. Leveraging the power of deep features, this translates into featuremetric keypoint and bundle adjustments that elegantly integrate into any SfM pipeline by replacing their geometric counterparts. Figure 2 shows an overview.

We first introduce the featuremetric optimization in Section 4.1. We then describe our formulations of keypoint adjustment, in Section 4.2, and bundle adjustment, in Section 4.3, and analyze their efficiency.

4.1. Featuremetric optimization

Direct alignment: We consider the error between image intensities at two sparse observations: $\mathbf{r} = \mathbf{I}_i[\mathbf{p}_u] - \mathbf{I}_j[\mathbf{p}_v]$.

Local image derivatives implicitly define a flow from one point to the other through a gradient descent update:

$$\mathbf{T}_{v \rightarrow u}[\mathbf{p}_v] \propto -\frac{\partial \mathbf{I}_j}{\partial \mathbf{p}} [\mathbf{p}_v]^\top \mathbf{r} . \quad (3)$$

This flow can be efficiently computed at any location in a neighborhood around v , without approximate interpolation nor descriptor matching. It naturally emerges from the direct optimization of the photometric error, which can be minimized with second-order methods in the same way as the aforementioned geometric costs. Unlike the flow regressed from a black-box neural network [24], this flow can be made consistent across multiple view by jointly optimizing the cost over all pairs of observations.

Learned representation: SfM can handle image collections with unconstrained viewing conditions exhibiting large changes in terms of illumination, resolution, or camera models. The image representation used should be robust to such changes and ensure an accurate refinement in any condition. We thus turn to features computed by deep CNNs, which can exhibit high invariance by capturing a large context, yet retain fine local details. For each image \mathbf{I}_i , we compute a D -dimensional, L2-normalized feature map $\mathbf{F}_i \in \mathbb{R}^{W \times H \times D}$ at identical resolution. We use the same representations for keypoint and bundle adjustments, requiring a single forward pass per image. Our experiments show that multiple off-the-shelf dense local descriptors can result in highly accurate refinements. However, our formulation can also be applied to robust intensity representations, such as the normalized cross-correlation (NCC) over local image patches [88].

4.2. Keypoint adjustment

Once local features are detected, described, and matched, we refine the keypoint locations before geometrically verifying the tentative matches.

Track separation: Connected components in the matching graph define tentative tracks – sets of keypoints that are likely to observe the same 3D point, but whose observations have not yet been geometrically verified. Because a 3D point has a single projection on a given image plane, valid tracks cannot contain multiple keypoints detected in the same image. We can leverage this property to efficiently prune out most incorrect matches using the track separation algorithm introduced in [24]. This speeds up the subsequent optimization and reduces the noise in the estimation.

Objective: We then adjust the locations of 2D keypoints belonging to the same track j by optimizing its featuremetric consistency along tentative matches with the cost

$$E_{\text{FKA}}^j = \sum_{(u,v) \in \mathcal{M}(j)} w_{uv} \left\| \mathbf{F}_{i(u)}[\mathbf{p}_u] - \mathbf{F}_{i(v)}[\mathbf{p}_v] \right\|_\gamma , \quad (4)$$

where w_{uv} is the confidence of the correspondence (u, v) , such as the similarity of its local feature descriptors $\mathbf{d}_u^\top \mathbf{d}_v$.

This allows the optimization to split tracks connected by weak correspondences, providing robustness to mismatches. The confidence is not based on the dense features since these are not expected to disambiguate correspondences at the global image level.

Efficiency : This direct formulation simply compares pre-computed features on sparse points and is thus much more scalable than patch flow regression (Eq. 2), which performs a dense local correlation for each correspondence. All tracks are optimized independently, which is very fast in practice despite the sheer number of tentative matches.

Drift: Because of the lack of geometric constraints, the points are free to move anywhere on the underlying 3D surface of the scene. The featuremetric cost biases the updates towards areas with low spatial feature gradients and with better-defined features. This can result in a large drift if not accounted for. Keypoints should however remain repeatable w.r.t. unrefined detections to ensure the matchability of new images, such as for visual localization. It is thus critical to limit the drift, while allowing the refinement of noisier keypoints. For each track, we freeze the location of the keypoint \bar{u} with highest connectivity, as in [24], and constrain the location \mathbf{p}_u of each keypoint w.r.t. to its initial detection \mathbf{p}_u^0 , such that $\|\mathbf{p}_u - \mathbf{p}_u^0\| \leq K$.

Once all tracks are refined, the geometric estimation proceeds, typically using two-view epipolar geometric verification followed by incremental or global SfM.

4.3. Bundle adjustment

The estimated structure and motion can then be refined with a similar featuremetric cost. Here keypoints are implicitly defined by the projections of the 3D points into the 2D image planes, and only poses and 3D points are optimized.

Objective: We minimize for each track j the error between its observations and a reference appearance \mathbf{f}^j :

$$E_{\text{FBA}} = \sum_j \sum_{(i,u) \in \mathcal{T}(j)} \left\| \mathbf{F}_i [\Pi(\mathbf{R}_i \mathbf{P}_j + \mathbf{t}_i, \mathbf{C}_i)] - \mathbf{f}^j \right\|_\gamma . \quad (5)$$

The reference is selected at the beginning of the optimization and kept fixed from then on. This reduces the drift of the points significantly, as also noted in [5], but is more flexible than the common ray-based parametrization [26, 44, 88].

The reference is defined as the observation closest to the robust mean μ over all initial observations \mathbf{f}_u^j of the track:

$$\mathbf{f}^j = \operatorname{argmin}_{\mathbf{f} \in \{\mathbf{f}_u^j\}} \|\mu^j - \mathbf{f}\| \quad (6)$$

$$\text{with } \mu^j = \operatorname{argmin}_{\mu \in \mathbb{R}^D} \sum_{\mathbf{f} \in \{\mathbf{f}_u^j\}} \|\mathbf{f} - \mu\|_\gamma . \quad (7)$$

This ensures robustness to outlier observations and accounts for the unknown topology of the feature space.

SfM features ↓ Refinement	ETH3D indoor						ETH3D outdoor					
	Accuracy (%)			Completeness (%)			Accuracy (%)			Completeness (%)		
	1cm	2cm	5cm	1cm	2cm	5cm	1cm	2cm	5cm	1cm	2cm	5cm
SIFT [51]	75.62	85.04	92.45	0.21	0.87	3.61	57.64	71.92	85.23	0.06	0.34	2.45
↳ Patch Flow	80.99	89.06	95.06	0.24	0.97	3.88	64.79	78.90	90.04	0.08	0.41	2.76
↳ ours	82.82	89.77	94.77	0.25	0.96	3.75	68.43	80.73	91.28	0.08	0.42	2.75
SuperPoint [21]	75.76	85.61	93.38	0.59	2.21	8.89	50.45	65.07	80.26	0.10	0.55	3.92
↳ Patch Flow	85.77	91.57	95.85	0.72	2.51	9.59	64.94	77.65	88.86	0.15	0.77	4.93
↳ ours	89.33	93.58	96.58	0.74	2.53	9.51	71.27	82.58	92.08	0.16	0.83	5.06
D2-Net [23]	47.18	64.94	83.37	0.47	1.87	7.07	20.87	34.55	56.53	0.03	0.19	1.78
↳ Patch Flow	79.10	86.64	93.26	1.45	4.53	12.95	57.34	70.71	84.12	0.21	1.06	6.02
↳ ours	82.49	88.83	94.35	1.36	4.13	11.80	65.71	77.95	89.22	0.21	1.01	5.63
R2D2 [59]	66.30	79.21	90.00	0.53	2.06	8.62	49.32	66.10	83.10	0.11	0.55	3.63
↳ Patch Flow	77.94	85.82	92.48	0.66	2.32	9.07	64.14	78.10	90.18	0.16	0.71	4.09
↳ ours	80.67	87.61	93.42	0.67	2.31	8.95	67.77	80.85	91.91	0.16	0.73	4.09

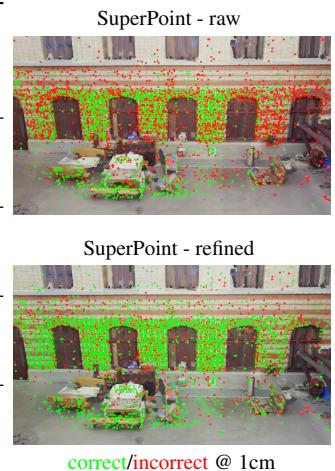


Table 1: **3D sparse triangulation.** Our refinement yields significantly more accurate and complete point clouds than the common geometric SfM pipeline. It is more effective than the existing Patch Flow [24], especially at 1cm or with SIFT.

Efficiency: Compared to the keypoint adjustment (Eq. 4), using a reference feature reduces the number of residuals from $\mathcal{O}(N^2)$ to $\mathcal{O}(N)$. On the other hand, all tracks need to be updated simultaneously because of the interdependency caused by the camera poses. To accelerate the convergence, we form a reduced camera system based on the Schur complement and use embedded point iterations [42]. The refinement generally converges within a few camera updates.

4.4. Implementation

Dense extractor: Our refinement can work with any off-the-shelf CNN that produces feature maps that are locally discriminative. These should be of the same resolution as the input (stride 1) to enable subpixel accuracy. The radius of convergence, or context, of such features depends on the amount of noise in the keypoints. Most detectors like SIFT have at most a few pixels of error, while others like D2-Net exhibit a much larger detection noise. In our experiments, we use S2DNet [31] for dense feature extraction, as it computes fine features very efficiently in only 4 convolutions, but also produce, if required, deeper features with a larger context. These can then be combined into a multi-level optimization scheme [26, 66, 86] that sequentially refines based on coarse to fine features. The convergence can thus be adjusted depending on the detector and on the image resolution. We show in Section 5.4 that other dense features work well too.

Optimization: The optimization problems of both keypoint and bundle adjustments are solved with the Levenberg-Marquardt [45] algorithm implemented using Ceres [3]. Feature maps are stored as collections of 16×16 patches centered around the initial keypoint detections. We thus constrain points to move at most $K=8$ pixels. The feature lookup is implemented as bicubic interpolation. We use the Cauchy loss γ with a scale of 0.25. The robust mean in Eq. 7

is computed with iteratively reweighted least squares [37].

Simultaneously storing all high-dimensional feature patches incurs high memory requirements during BA. We dramatically increase its efficiency by exhaustively precomputing patches of feature distances and directly interpolate an approximate cost $\bar{E}_{ij} = \|\mathbf{F}_i - \mathbf{f}^j\|_\gamma [\mathbf{p}_{ij}]$. To improve the convergence, we store and optimize its spatial derivatives $\partial \bar{E}_{ij} / \partial \mathbf{p}_{ij}$. This reduces the residual size from D to 3 with no loss of accuracy. See Appendix C for more details.

Run time and memory: S2DNet can extract 3-5 dense feature maps per second and both featuremetric adjustments run in less than 5 minutes for 100 images. As these features are 128-dimensional, the memory consumption can be a bottleneck. We believe that much fewer dimensions are actually required for refinement, and retraining a compact feature extractor would improve the efficiency of the optimization.

5. Experiments

We evaluate our featuremetric refinement on various SfM tasks with several handcrafted and learned local features and show substantial improvements for all of them. We first evaluate its accuracy on the tasks of triangulation and camera pose estimation in Sections 5.1 and 5.2, respectively. We then assess in Section 5.3 the impact of the refinement on two-view and multi-view pose estimation for end-to-end reconstruction in challenging conditions. Lastly, Section 5.4 analyzes the validity and scalability of our design decisions through an ablation study.

5.1. 3D triangulation

We first evaluate the accuracy of the refined 3D structure given known camera poses and intrinsics.

Evaluation: We use the ETH3D benchmark [73], which

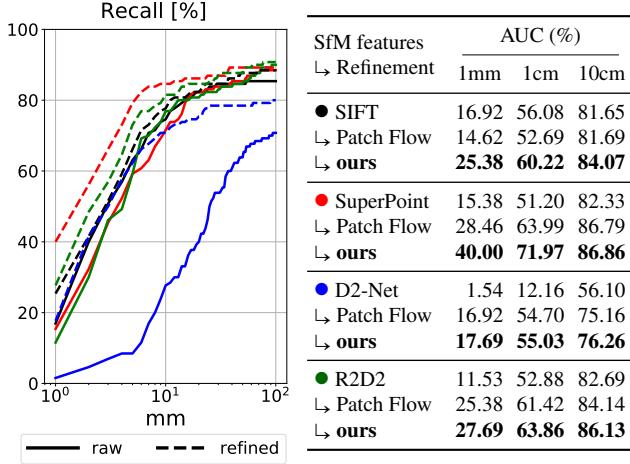


Table 2: **Camera pose estimation.** We plot the cumulative translation error and report its AUC. Our refinement improves the accuracy of the query camera poses for all local features, even when for SIFT, whose detections are already well-localized. It is generally more accurate than Patch Flow.

is composed of 13 indoor and outdoor scenes and provides images with millimeter-accurate camera poses and highly-accurate ground truth dense reconstructions obtained with a laser scanner. We follow the protocol introduced in [24], in which a sparse 3D model is triangulated for each scene using COLMAP [70] with fixed camera poses and intrinsics. Following the original benchmark setup, we report the accuracy and completeness of the reconstruction, in %, as the ratio of triangulated and ground-truth dense points that are within a given distance of each other.

Baselines: We evaluate our featuremetric refinement with the hand-crafted local features SIFT [51] and the learned ones SuperPoint [21], D2-Net [23], and R2D2 [59], using the associated publicly available code repositories. We compare our approach to the geometric optimization of [24], referred here as Patch Flow. We re-compute the numbers provided in the original paper using the code provided by the authors.

Results: Table 1 shows that our approach results in significantly more accurate and complete 3D reconstructions compared to the traditional geometric SfM. It is more accurate than Patch Flow, especially at the strict threshold of 1cm, and exhibits similar completeness. The improvements are consistent across all local features, both indoors and outdoors. The gap with Patch Flow is especially large for SIFT, which already detects well-localized keypoints. This confirms that our featuremetric optimization better captures low-level image information and yields a finer alignment. Patch Flow is more complete for larger thresholds as it partly solves a different problem by increasing the keypoint repeatability with its large receptive field, while we focus on their localization.

SfM features (# keypoints)	Task 1: Stereo		Task 2: Multiview		
	AUC@K°		AUC@5° @N		
	5°	10°	5	10	25
SuperPoint+SuperGlue (2k)	58.78	71.01	63.02	77.36	86.76
ours	65.89	76.51	68.87	82.09	89.73
SIFT (2k)	38.09	48.05	25.12	50.82	77.28
ours	40.59	50.87	28.01	53.59	79.49
D2-Net (4k)	16.83	22.40	16.52	33.07	49.35
ours	25.89	33.32	21.33	40.69	57.93

Table 3: **End-to-end SfM.** The proposed refinement improves the accuracy of poses estimated by epipolar geometry (stereo) or a complete SfM pipeline (multiview) with crowd-sourced imagery. Improvements are substantial for both standard (SIFT) and recent (SuperGlue) matching configurations, especially when few images N observe the scene.

5.2. Camera pose estimation

We now evaluate the impact of our refinement on the task of camera pose estimation from a single image.

Evaluation: We again follow the setup of [24] based on the ETH3D benchmark. For each scene, 10 images are randomly selected as queries. For each of them, the remaining images, excluding the 2 most covisible ones, are used to triangulate a sparse 3D partial model. Each query is then matched against its corresponding partial model and the resulting 2D-3D matches serve to estimate its absolute pose using LO-RANSAC+PnP [15] followed by geometric refinement. We compare the 130 estimated query poses to their ground truth and report the area under the cumulative translation error curve (AUC) up to 1mm, 1cm, and 10cm.

Baselines: Patch Flow performs multi-view optimization over each partial model independently as well as over the matches between each query and its partial model. Similarly, we first refine each partial model as in Section 5.1. We then adjust the query keypoints using its tentative matches, estimate an initial pose, and refine it with featuremetric BA.

Results: The AUC and its cumulative plot are shown in Table 2. Our refinement substantially improves the localization accuracy for all local features, including SIFT, for which Patch Flow does not show any benefit. At all error thresholds, featuremetric optimization is consistently more accurate than its geometric counterparts. The accuracy of SuperPoint is raised far higher than other detectors, despite the high sparsity of the 3D models that it produces. This shows how more accurate keypoint detections can result in much more accurate visual localization.

5.3. End-to-end Structure-from-Motion

While the previous experiments precisely quantify the accuracy of the refinement, they do not contain any variations

of appearance or camera models. We thus turn to crowd-sourced imagery and evaluate the benefits of our featuremetric optimization in an end-to-end reconstruction pipeline.

Evaluation: We use the data, protocol, and code of the 2020 Image Matching Challenge [1, 43]. It is based on large collections of crowd-sourced images depicting popular landmarks around the world. Pseudo ground truth poses are obtained with SfM [70] and used for two tasks. The stereo task evaluates relative poses estimated from image pairs by decomposing their epipolar geometry. This is a critical step of global SfM as it initializes its global optimization. The multiview task runs incremental SfM for small subsets of images, making the SfM problem much harder, and evaluates the final relative poses within each subset. For each task, we report the AUC of the pose error at the threshold of 5° , where the pose error is the maximum of the angular errors in rotation and translation. As the evaluation server accepts at most correspondences, we cannot evaluate our method using the test data. We instead test on a subset of the publicly available validation scenes, and tune the RANSAC and matching parameters on the remaining scenes. More details on this setup are provided in the Appendix.

Baselines: We evaluate our refinement in combination with SIFT [51], D2-Net [23], and SuperPoint+SuperGlue [21, 65]. We limit the number of detected keypoints to 2k for computational reasons, but increase this number to 4k for D2-Net as it otherwise performs poorly. In the stereo task, we adjust the keypoints using the entire exhaustive tentative match graph (4950 pairs per scene). We use LO-DEGENSAC [15, 16] for match verification, the ratio test for SIFT, and the mutual check for SIFT and D2-Net. In the multiview task, we adjust keypoints for each subset independently, considering only the matches between images in the subset, and run our bundle adjustment after SfM.

Results: Table 3 summarizes the results. For stereo, our featuremetric keypoint adjustment significantly improves the accuracy of the two-view epipolar geometries across all local features and despite the challenging conditions. In multiview setting, it also improves the accuracy of the SfM poses, especially for small sets of images. Featuremetric optimization is particularly effective in this situation, as geometric optimization cannot fully suppress the detection noise due to the small number of observations. We visualize tracks of a 5-image reconstruction in Figure 4 and highlight the accuracy of the refined SfM model.

5.4. Additional insights

Ablation study: Table 4 shows the performance of several variants of our featuremetric optimization on ETH3D in terms of triangulation (scene *Facade* only) and localization (all scenes). We compare both types of adjustments, minor tweaks, and different image representations, including

	SuperPoint ↓ Refinement	Acc. (%)		Compl. (%)		track length	AUC 1cm
		1cm	2cm	1cm	2cm		
KA vs. BA	unrefined	18.42	32.23	0.06	0.49	4.17	51.20
	↓ Patch Flow [24]	37.00	55.18	0.15	0.93	5.24	63.53
	↓ F-KA	36.85	54.48	0.15	0.90	5.02	69.84
	↓ F-BA	43.65	62.44	0.18	1.06	4.17	67.61
KA vs. BA (full)	↓ F-KA+BA (full)	46.46	65.41	0.19	1.14	5.02	71.97
	w/ F-BA drift	47.93	66.52	0.20	1.17	5.02	64.51
	Patch Flow + F-BA	46.30	65.22	0.19	1.13	5.24	-
	higher resolution	47.67	65.39	0.21	1.21	5.12	-
dense feats	photometric BA [88]	28.43	45.87	0.11	0.72	4.17	-
	VGG-16 ImageNet	36.86	54.99	0.15	0.90	4.61	-
	DSIFT [49]	38.78	56.46	0.16	0.96	4.73	-
	PixLoc [66]	29.49	46.60	0.12	0.74	4.48	-

Table 4: **Ablation study on ETH3D.** i) Featuremetric keypoint and bundle adjustments (KA and BA) both largely improve the **triangulation** and **localization** accuracy. Patch Flow produces a longer track length because of its larger receptive field but is less accurate. ii) Letting the BA drift by updating reference features or increasing the image resolution both improve the triangulation, at the expense of poorer localization and increased run time, respectively. iii) Different image representations are better than the unrefined detections but S2DNet (our default) works best.

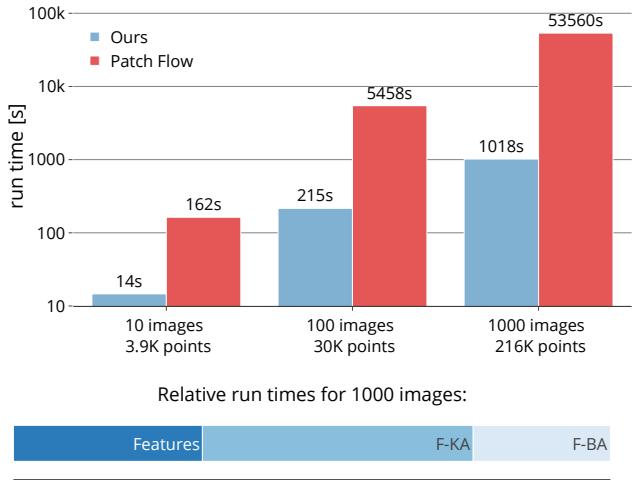


Figure 3: **Run-times.** We show the duration, in logarithmic scale, of the refinement for varying numbers of images. Our refinement is more than ten times faster than Patch Flow [24], whose run-time is dominated by the computation of the pairwise flow, which scales quadratically. Thanks to our precomputed cost patches, the featuremetric BA is fast. The KA amounts for the majority of the refinement time.

NCC-normalized intensity patches with fronto-parallel warping. Our final configuration, based on the dense features of S2DNet [31], performs best across all metrics. We will now show that it is also fairly efficient.



Figure 4: Refined SfM tracks. We show patches centered around reprojections of 3x 3D points observed in 4 images of the *St. Peter’s Square* scene. Deep features and their correlation maps with a reference are robust to scale or illumination changes, yet preserve local details required for fine alignment. Points refined with our approach (in green) are consistent across multiple views while those of a standard SfM pipeline (in red) are misaligned because the initial keypoint detections (in blue) are noisy.

Scalability: We run SfM on subsets of images of the Aachen Day-Night dataset [67, 68, 95]. Figure 3 shows the run times of the refinement for subsets of 10, 100 and 1000 images. The featuremetric refinement is an order of magnitude faster than Patch-Flow [24]. Precomputing distance maps reduces the peak memory requirement of the bundle adjustment from 80 GB to less than 10GB for 1000 images. As storing feature maps only requires 50 GB of disk space, this refinement can easily run on a desktop PC. We thus refined the entire Aachen Day-Night v1.1 model, composed of 7k images, in less than 2 hours. Scene partitioning [70] could further reduce the peak memory. See Appendix D for more details.

6. Conclusion

In this paper we argue that the recipe for accurate large-scale Structure-from-Motion is to perform an initial coarse estimation using sparse local features, which are by necessity globally-discriminative, followed by a refinement using locally-accurate dense features. Since the dense feature only need to be locally-discriminative, they can afford to capture much lower-level texture, leading to more accurate corre-

spondences. Through extensive experiments we show that this results in more accurate camera poses and structure; in challenging conditions and for different local features.

While we optimize against dense feature maps, we keep the sparse scene representation of SfM. This ensures not only that the approach is scalable but also that the resulting 3D model is compatible with downstream applications, e.g. mapping for visual localization. Since our refinement works well even with few observations, as it does not need to average out the keypoint detection noise, it has the potential to achieve more accurate results using fewer images.

We thus believe that our approach can have a large impact in the localization community as it can improve the accuracy of the ground truth poses of standard benchmark datasets, of which many are currently saturated. Since this refinement is less sensitive to under-sampling, it enables benchmarking for crowd-sourced scenarios beyond densely-photographed tourism landmarks.

Acknowledgements: The authors thank Mihai Dusmanu, Rémi Pautrat, Marcel Geppert, and the anonymous reviewers for their thoughtful comments. Paul-Edouard Sarlin was supported by gift funding from Huawei, and Viktor Larsson by an ETH Zurich Postdoctoral Fellowship.

Appendix

A. Additional results on ETH3D

A.1. Triangulation

We refine the triangulation of SuperPoint [21] keypoints for the ETH3D *Courtyard* scene and show in Figure 5 the distribution of triangulation errors for points observed by different numbers of images (track length). Our featuremetric refinement provides the largest improvement for points with low track length, for which the estimates of the traditional geometric BA are dominated by the noise of the keypoint detection. For larger track lengths, the refined point cloud has an accuracy close to the Faro Focus X 330 laser scanner from which the ground truth is computed.

We show in Figure 10 the raw and refined point clouds for SuperPoint and D2-Net. The benefits of our refinement are easily visible in 3D. Planar walls exhibit fewer noisy keypoints and the refined point clouds are more complete.

A.2. Camera pose estimation

We analyze in Table 5 how the different kinds of adjustments impact the accuracy of camera localization. The full method presented in the main paper first refines the 3D SfM model with featuremetric keypoint and bundle adjustments. It then refines each keypoint in the query image using its tentative 2D-3D correspondences by minimizing the featuremetric error between its observation in the query and the most similar observation of the respective 3D points. Refining the query keypoints before RANSAC increases the number of inlier matches and stabilizes the pose estimation in challenging scenarios where few 3D points are matches.

Once an initial pose is estimated with PnP+RANSAC, we refine it via a small featuremetric bundle adjustment over the inlier correspondences. This optimizes each query keypoint against the closest descriptor within the matched track. As opposed to refining each query keypoint against all observations in the matched track, this has the benefit of scaling linearly in the number of query keypoints and yields a similar accuracy.

B. Impact of various parameters

B.1. Patch size

Figure 6 shows how much our refinement displaces the detected keypoints during the triangulation of SuperPoint on *Courtyard* using dense features extracted from 1600x1066-pixel images. When using full feature maps without any constraints in keypoint adjustment, most points are moved by more than 1 pixel, but most often by less than 8 pixels. This confirms that storing the feature maps as 16×16 patches is sufficient and rather conservative.

We show in Figure 7 the accuracy of the triangulation for

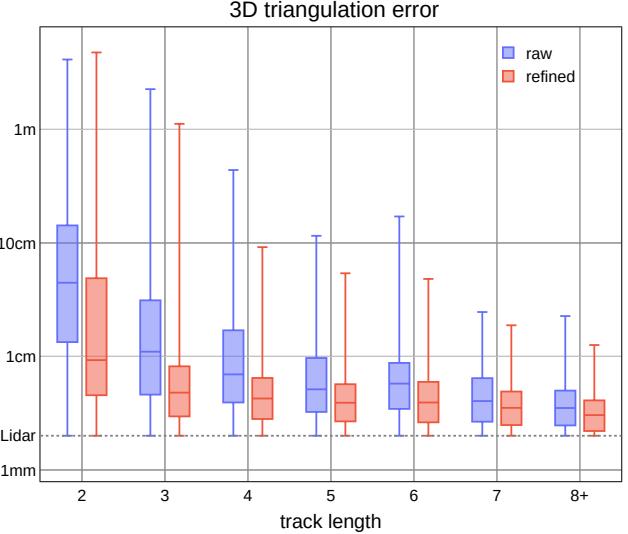


Figure 5: **Triangulation errors vs. track length.** The initial, unrefined output, based on geometric BA, exhibits high errors for 3D points that are observed by few images (low track length). Our refinement significantly reduces these errors and brings the accuracy of the sparse point cloud close to the ground truth acquired by Lidar (2mm accuracy).

SuperPoint ↳ Refinement	KA	BA	qKA	qBA	AUC (%)		
					1mm	1cm	10cm
unrefined					15.38	51.20	82.33
↳ refined	✓				16.15	53.34	82.49
↳ refined	✓	✓			16.92	54.71	84.08
↳ refined	✓	✓	✓		38.46	70.44	85.28
↳ refined (full)	✓	✓	✓	✓	40.00	71.97	86.86
↳ Patch Flow	✓		✓		28.46	63.04	86.65

Table 5: **Ablation study for pose estimation.** The accuracy of the camera pose is improved by refining the map (KA and BA) and by refining the query keypoints before (qKA) and after (qBA) pose estimation. The largest improvement is brought by qKA. It increases the number of inlier matches and the likelihood of finding the correct pose with RANSAC.

various patch sizes. Smaller 10×10 patches achieve sufficient accuracys and require significantly less memory.

B.2. Image resolution

The image resolution at which the dense features are extracted has a large impact on the accuracy of the refinement. In Figure 8 we quantify in the impact on both triangulation accuracy and run time for the ETH3D *Courtyard* scene (38 images). The accuracy drops significantly when the resolution is smaller than 1600×1066px, which amounts to 25% of the full image resolution. Doubling the resolution to 3200×2132px yields noticeable improvements, albeit sig-

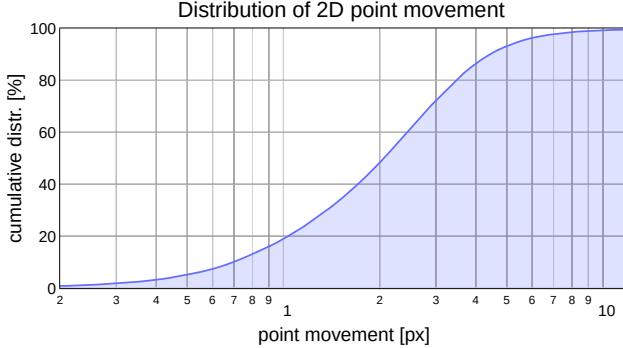


Figure 6: Distribution of point movements. We show the cumulative distribution of the distance traveled by the 2D keypoints during the featuremetric refinement of SuperPoint with KA and BA. 60% of the points move by fewer than 2 pixels and 99% remain within 8 pixels of the initial detections.

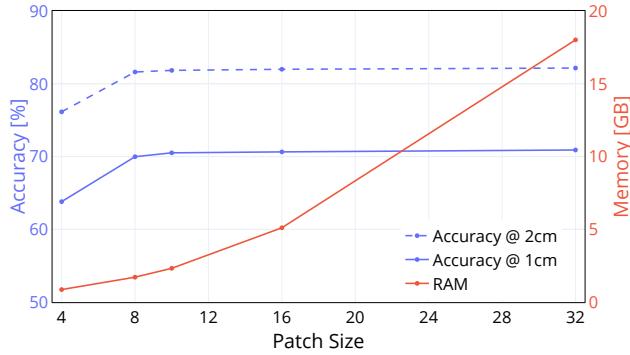


Figure 7: Impact of the patch size. Smaller patches for each observation significantly reduce memory requirements but can impair the accuracy of the refinement. Patches of size 10×10 offer a good trade-off with high accuracy and moderate memory consumption.

nificantly increases the extraction time and the consumption of GPU VRAM. As a reference, extracting only fine-level S2DNet features (4 convolutions) from 3200×2132 px images requires around 10GB of GPU VRAM.

B.3. Reference selection for keypoint adjustment

Selecting some observations as references is necessary to avoid the drift. In a given track, the keypoint adjustment selects the point that is the most connected (topological center), while the bundle adjustment selects the point closest to the robust mean in feature space (feature center). Could we use the feature center for selecting the reference of the keypoint adjustment? By minimizing the feature distance to this unique reference, we could reduce the number of residuals from quadratic (pairwise constraints) to linear (unary constraints) and thus accelerate the optimization.

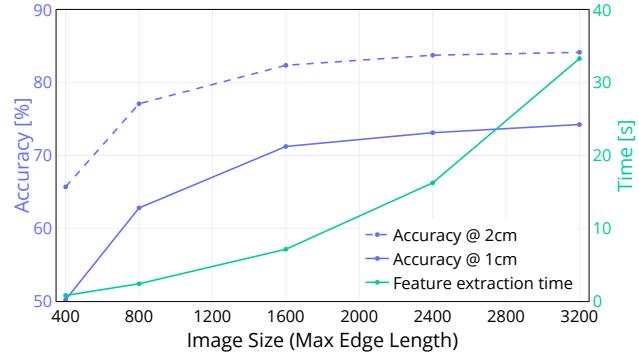


Figure 8: Impact of the image resolution. Increasing the image resolution increases the accuracy, but at the cost of longer feature extraction time and higher VRAM requirements. For all experiments on ETH3D, we used a maximum edge length of 1600px, which is very close to saturating the accuracy while providing low run times.

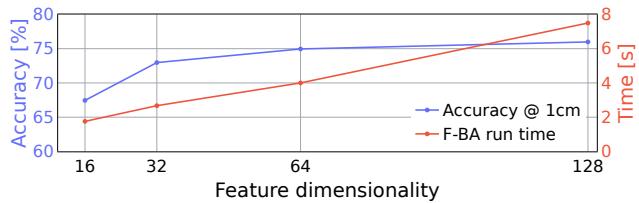


Figure 9: Impact of the feature dimensionality. Dense features computed by S2DNet can be naively reduced to accelerate the featuremetric bundle adjustment by 2 while incurring only a minor drop of triangulation accuracy.

Retaining pairwise constraints however allows the optimization to separate tracks that were incorrectly merged by the track separation algorithm. This is not necessary in the bundle adjustment, as tracks are already filtered by the robust geometric estimation and can thus be assumed to be correct, but is common for unverified track. We evaluate the impact of the reference selection in the keypoint adjustment and report the results in Table 6. For both SuperPoint and D2-Net, using the feature center results in lower completeness and accuracy than the topological center. It also results in a lower track length, which confirms that the topological reference allows to retain incorrectly-merged tracks. Since the feature center still performs relatively well, it could be considered in case of tighter computational constraints.

Furthermore, Table 6 highlights the importance of the featuremetric keypoint adjustment. The benefits are larger for D2-Net, which detects very noisy keypoints. As a consequence, many correct albeit noisy matches are rejected by the geometric verification. Our keypoint adjustment not only allows more points to be triangulated, thus increasing the completeness of the model, but also increases the accuracy of the triangulated points.

	Triangulation	Acc. (%)		Compl. (%)			track length
		1cm	2cm	1cm	2cm	5cm	
SuperPoint	unrefined	18.03	31.97	0.07	0.49	5.03	4.17
	↪ Patch Flow [24]	37.00	55.18	0.15	0.93	7.44	5.24
	↪ F-BA	43.65	62.44	0.18	1.06	7.70	4.17
	↪ +F-KA (feat-ref)	45.05	64.84	0.18	1.12	7.76	4.88
	↪ +F-KA (topol-ref)	46.46	65.41	0.19	1.14	8.19	5.02
D2-Net	unrefined	7.68	13.98	0.02	0.17	2.19	3.29
	↪ Patch Flow [24]	34.64	52.36	0.16	1.00	8.10	4.99
	↪ F-BA	39.30	58.59	0.15	0.94	6.99	3.29
	↪ +F-KA (feat-ref)	43.35	62.54	0.19	1.18	8.36	4.49
	↪ +F-KA (topol-ref)	44.21	64.22	0.20	1.20	8.72	4.63

Table 6: **Additional ablation study on ETH3D Facade.** i) Featuremetric keypoint adjustment significantly improves the completeness, especially for noisy keypoints as in D2-Net. ii) Keypoint adjustment against the topological center in each tentative track (topol-ref) improves the point cloud in accuracy and completeness over KA towards the robust feature center (feat-ref) because it allows to merge tracks.

B.4. Number of feature levels

Using multiple feature levels enlarges the basin of convergence but increases the computational requirements. The radius of convergence that is required depends on the noise of the keypoint detector and on the resolution of the image from which keypoints are detected. When performing detection and refinement at identical image resolutions, the optimal displacement is at most a few pixels for most keypoint detectors. In this case, the fine level of S2DNet feature maps is sufficient. We empirically measured that its radius of convergence is approximately 3 pixels, although the multiview constraints enable to refine over much larger distances.

We thus use a single feature level for all experiments involving SIFT, SuperPoint, and R2D2. D2-Net require a different treatment, as its detection noise is significantly larger. This is partly due to the aggressive downsampling of its CNN backbone and to the low resolution of its output heatmap. As a consequence, we employ both fine and medium feature levels for D2-Net. Both keypoint and bundle adjustments run the optimization successively at the coarser and finer levels.

B.5. Dimensionality of the features

Throughout this paper, we used 128-dimensional dense features extracted by S2DNet [31]. Relying on compact features would easily reduce the memory footprint and the run time of the refinement. To demonstrate these benefits, we show in Figure 9 the relationship between the dimension, the run time of the BA, and the triangulation accuracy when retaining only the first k channels of the S2DNet features. Features with fewer dimensions yield a faster refinement. The accuracy drops moderately but we expect a smaller reduction with features explicitly trained for smaller dimensions.

C. Cost map approximation

We mention in Section 4.4 that the memory efficiency of the bundle adjustment can be improved by precomputing the featuremetric cost. We provide here more details.

Description: Given D -dimensional features, the featuremetric bundle adjustment (Eq. 5) involves residuals and Jacobian matrices of dimension D . Unlike the keypoint adjustment, which can optimize tracks independently, all bundle parameters are updated simultaneously and the memory requirements are thus prohibitive. Given the 2D reprojection $\mathbf{p}_{ij} = \Pi(\mathbf{R}_i \mathbf{P}_j + \mathbf{t}_i, \mathbf{C}_i)$, this formulation loads in memory the dense features \mathbf{F}_i , interpolates them at \mathbf{p}_{ij} , and compute the residuals $\mathbf{r}_{ij} = \mathbf{F}_i[\mathbf{p}_{ij}] - \mathbf{f}^j$ for the cost $E_{ij} = \|\mathbf{r}_{ij}\|_\gamma$.

To reduce the memory footprint, we can exhaustively precompute patches of feature distances and treat them as one-dimensional residuals $\bar{\mathbf{r}}_{ij} = \|\mathbf{F}_i - \mathbf{f}^j\|[\mathbf{p}_{ij}]$. The cost then becomes $\bar{E}_{ij} = \gamma(\bar{\mathbf{r}}_{ij})$. Such distances only need to be computed once since the reference \mathbf{f}^j is kept fixed throughout the optimization. This precomputed cost reduces the peak memory by a factor D , with often $D=128$. It is similar to the Neural Reprojection Error recently introduced by Germain *et al.* [32] for camera localization.

Analysis: Swapping the distance computation and the sparse interpolation introduces an approximation error. We first write the bilinear or bicubic interpolation as a sum over features \mathbf{F}_k on the discrete grid:

$$\mathbf{F}[\mathbf{p}] = \sum_k w_k \mathbf{F}_k \quad \text{with} \quad \sum_k w_k = 1 . \quad (8)$$

We assume that the features are L2-normalized $\|\mathbf{F}_k\| = 1$, such that $\|\mathbf{F}[\mathbf{p}]\| \approx 1$. For a squared loss function, the approximation error can then be written as:

$$\begin{aligned} & \|\mathbf{F} - \mathbf{f}\|^2[\mathbf{p}] - \|\mathbf{F}[\mathbf{p}] - \mathbf{f}\|^2 \\ & \approx 1 - \|\mathbf{F}[\mathbf{p}]\|^2 = \frac{1}{2} \sum_k \sum_l w_k w_l \|\mathbf{F}_k - \mathbf{F}_l\|^2 . \end{aligned} \quad (9)$$

This error is zero at points on the discrete grid and increases with the roughness of the feature space. This approximation thus displaces the local minimum of the cost by at most 1 pixel but most often by much less.

Improvement: This approximation however degrades the correctness of the approximate Hessian matrix that the Levenberg-Marquardt algorithm [45] relies on for fast convergence. We found that also optimizing the squared spatial derivatives of this cost significantly improves the convergence. This simply amounts to augmenting the scalar residual map with dense derivative maps:

$$\tilde{\mathbf{r}}_{ij} = \begin{pmatrix} \|\mathbf{F}_i - \mathbf{f}^j\| \\ \frac{\partial \|\mathbf{F}_i - \mathbf{f}^j\|}{\partial x} \\ \frac{\partial \|\mathbf{F}_i - \mathbf{f}^j\|}{\partial y} \end{pmatrix} [\mathbf{p}_{ij}] . \quad (10)$$

SuperPoint ↳ Refinement	Acc. (%)		Compl. (%)		Time (s)	Memory (GB)
	1cm	2cm	1cm	2cm		
unrefined	64.27	76.47	0.37	1.44	-	-
↳ ours (exact)	81.31	88.50	0.47	1.74	42.22	7.3
↳ ours (cost maps)	80.27	87.81	0.47	1.72	29.86	0.15

Table 7: **Triangulation with cost map approximations.** Using precomputed cost maps increase the efficiency of the bundle adjustment with a marginal loss of accuracy.

This improvement results in three-dimensional residuals, which is still smaller than D when $D=128$. Using the spatial derivatives, we can also compute an exact, more accurate bicubic spline interpolation of the cost landscape.

Evaluation: We now show experimentally that this approximation often does not, or only minimally, impairs the accuracy of the refinement. Table 7 reports the results of the triangulation of SuperPoint features on the ETH3D dataset. The approximation reduces the accuracy by less than 1% and does not alter the completeness. It however significantly reduces the memory consumption of the bundle adjustment, allowing it to scale to thousands of images. Note that all experiments in Sections 5.1, 5.2, and 5.3 do not use the cost map approximation as the corresponding scenes are relatively small.

D. Experimental details

D.1. ETH3D - Sections 5.1 and 5.2

For the experiments on ETH3D, we use the evaluation code provided by Dusmanu *et al.* [24]. We use the original implementations of SuperPoint [21], D2-Net [23], and R2D2 [59], and extract root-normalized SIFT [51] features using COLMAP [70]. For both sparse and dense feature extraction, the images are resized so that their longest dimension is equal to 1600 pixels. The tentative matches are filtered according to the recipe described in [24].

D.2. Structure-from-Motion - Section 5.3

We tune the hyperparameters on the training scenes *Temple Nara Japan*, *Trevi Fountain*, and *Brandenburg Gate*. The results in the main paper are computed on the test scenes *Sacre Coeur*, *Saint Peter’s Square*, and *Reichstag*, using the data and code provided by the challenge organizers.

For SIFT [51], we use the mutual check, a ratio test with threshold 0.85 for the multi-view and 0.9 for the stereo tasks, and DEGENSAC with an inlier threshold of 0.5px. For D2-Net [23], we use the mutual check and inlier thresholds of 2px and 0.5px for raw and refined keypoints, respectively. For SuperPoint+SuperGlue [21, 65], we do not use additional match filtering and we select an inlier thresholds of 1.1px and 0.5px for raw and refined keypoints, respectively. All

sparse local and dense features are extracted at full image resolution, which is generally not larger than 1024px.

D.3. Ablation study - Section 5.4

The triangulation metrics are reported for the ETH3D scene *Facade*, which is the largest with 76 images. We use SuperPoint local features as they perform best in all earlier experiments and we store dense feature maps in every experiment. The localization AUC is measured over all 13 scenes in ETH3D with 10 holdout images per scene. We now detail the different baselines.

Localization is achieved in “F-KA” by first refining the keypoints, triangulating the map and finally performing query keypoint adjustment as described in section A.2. For localization with “F-BA”, we refined the triangulated model using featuremetric bundle adjustment and then refined the pose from PnP+RANSAC using qBA.

In the entry “w/ F-BA drift”, we use the robust reference (Eq. 7) to select the observation in each track which is most similar to the robust reference as the source frame. The optimizer then minimizes the error between each other observation and the current, moving reference of the source frame. Since only the index of the source frame is fixed during the optimization, this method does not account for drift, which appears to yield higher accuracy but suffers from repeatability problems during localization.

The baseline “PatchFlow + F-BA” uses the keypoint refinement from Dusmanu *et al.* [24] as initialization, and runs our featuremetric bundle adjustment on top of it. We used the exact same parameters for PatchFlow as presented in [24].

The entry “higher resolution” corresponds to input images at double the resolution than all the other experiments, i.e. 3200 pixels in the longest dimension.

For the “photometric” baseline, we use RGB images (while Woodford *et al.* [88] use grayscale images), we warp patches of 4×4 pixels at the featuremap resolution (1600 pixels in the longest dimension) with fronto-parallel assumption, and apply normalized cross correlation (NCC). Identically to our featuremetric BA and to LSPBA [88], the source frame is selected as the observation closest to the robust mean.

We report results for dense features extracted from a VGG-16 CNN, trained on ImageNet [20], at the layer conv1_2 (64 channels) and for the fine feature map predicted by PixLoc [66] (32 channels). The model of PixLoc, trained on MegaDepth [48], was kindly provided by its authors. In DSIFT [49] (128 channels), we apply a bin size of 4 and a step size of 1 and refer to the VLFeat implementation [85] for more details.

D.4. Scalability

All experiments were conducted on 8 CPU cores (Intel Xeon E5-2630v4) and one NVIDIA RTX 1080 Ti. The subsets from the Aachen Day-Night v1.1 model [67, 68, 95]

were selected as the images with the largest visibility overlap, in descending order. To accelerate the feature matching, each image was matched only to its top 20 most covisible reference images in the original Aachen SfM model. We use SuperPoint [21] features and match image pairs with the mutual check and distance thresholding at 0.7. During BA, we apply the sparse Schur solver from Ceres for each linear system in LM, while we use sparse Cholesky in KA, similar to [24]. Featuremetric bundle adjustment is stopped after 30 iterations while KA runs for at most 100 iterations and stops when parameters change by less than 10^{-4} .

To refine the full Aachen Day-Night model, we use SuperPoint features matched with SuperGlue [65] from the Hierarchical Localization toolbox [63, 64]. We refine the keypoints with KA, then triangulate the points with fixed poses from the reference model. Finally, we run a full bundle adjustment of the model with the proposed approximation by cost maps.

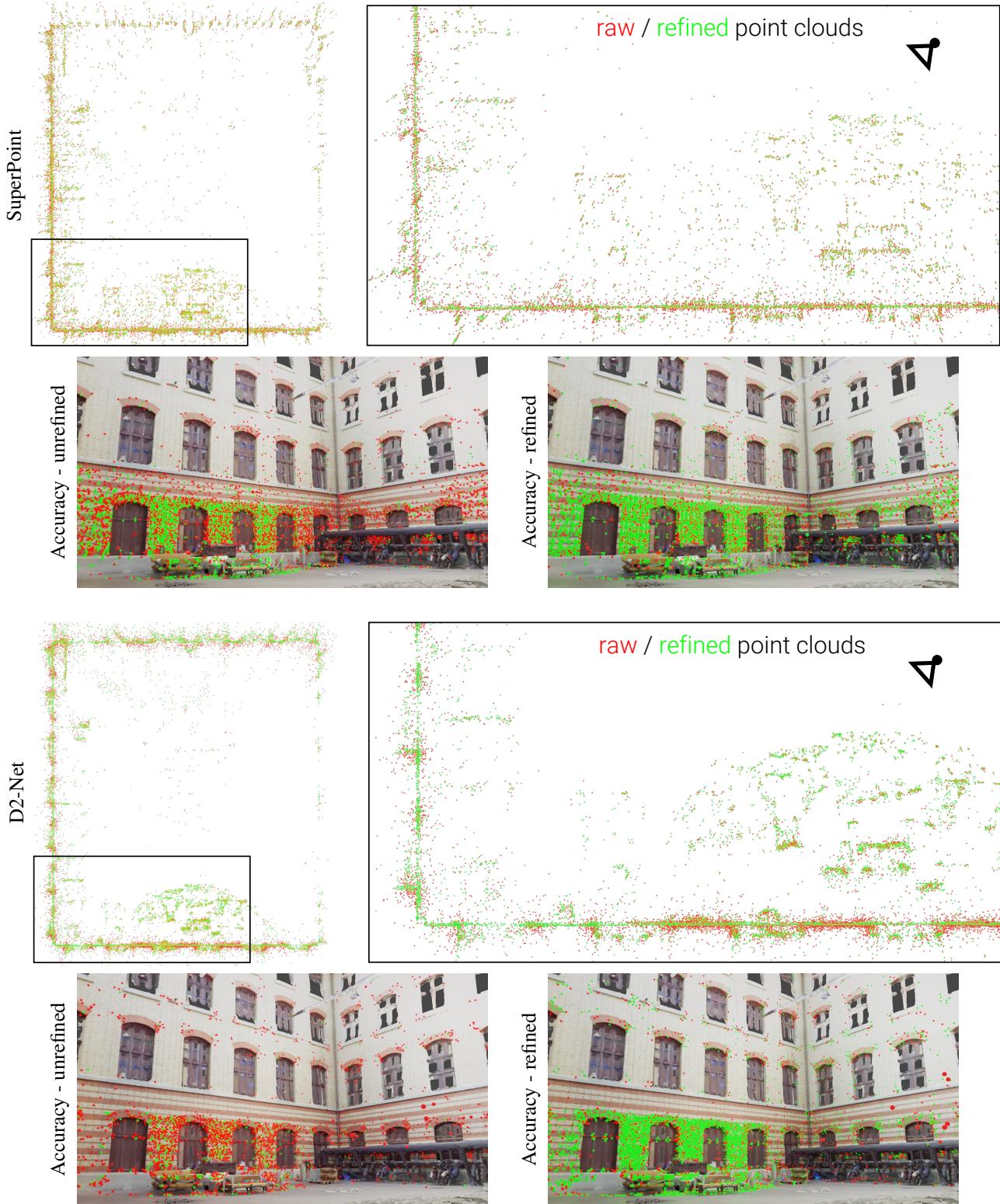


Figure 10: **Refinement on ETH3D Courtyard.** In the top parts, we show for both SuperPoint (top) and D2-Net (bottom) top-down views of the sparse point clouds triangulated with raw (in red) and refined (in green) keypoints. The refined point clouds better fit the geometry of the scene, especially on planar walls. In the lower parts, we also show images in which points are colored as accurate (in green) or inaccurate (in red) at 1cm for raw (left) and refined (right) point clouds.

References

- [1] CVPR 2020 Image Matching Challenge. [https://www.cs.ubc.ca/research/
image-matching-challenge/](https://www.cs.ubc.ca/research/image-matching-challenge/). Accessed March 1, 2021. 7
- [2] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M Seitz, and Richard Szeliski. Building Rome in a day. *Communications of the ACM*, 54(10):105–112, 2011. 1
- [3] Sameer Agarwal, Keir Mierle, and Others. Ceres solver. <http://ceres-solver.org>. 5
- [4] Sameer Agarwal, Noah Snavely, Steven M Seitz, and Richard Szeliski. Bundle adjustment in the large. In *ECCV*, 2010. 1
- [5] Hatem Alismail, Brett Browning, and Simon Lucey. Photometric bundle adjustment for vision-based SLAM. In *ACCV*, 2016. 2, 4
- [6] Hatem Alismail, Brett Browning, and Simon Lucey. Robust tracking in low light and sudden illumination changes. In *3DV*, 2016. 3
- [7] Epameinondas Antonakos, Joan Alabert-i Medina, Georgios Tzimiropoulos, and Stefanos P Zafeiriou. Feature-based lucas-kanade and active appearance models. *IEEE Transactions on Image Processing*, 2015. 3
- [8] Simon Baker, Ralph Gross, and Iain Matthews. Lucas-kanade 20 years on: A unifying framework. *IJCV*, 56, 2003. 2
- [9] Daniel Barath, Dmytro Mishkin, Ivan Eichhardt, Ilia Shipachev, and Jiri Matas. Efficient initial pose-graph generation for global sfm. In *CVPR*, 2021. 2
- [10] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. SURF: Speeded up robust features. In *ECCV*, 2006. 1
- [11] Che-Han Chang, Chun-Nan Chou, and Edward Y Chang. CLKN: Cascaded Lucas-Lanade networks for image alignment. In *CVPR*, 2017. 3
- [12] Avishek Chatterjee and Venu Madhav Govindu. Efficient and robust large-scale rotation averaging. In *CVPR*, 2013. 2
- [13] Christopher B Choy, JunYoung Gwak, Silvio Savarese, and Manmohan Chandraker. Universal correspondence network. In *NIPS*, 2016. 2
- [14] Peter Hviid Christiansen, Mikkel Fly Kragh, Yury Brodskiy, and Henrik Karstoft. UnsuperPoint: End-to-end unsupervised interest point detector and descriptor. *arXiv:1907.04011*, 2019. 2
- [15] Ondřej Chum, Jiří Matas, and Josef Kittler. Locally optimized RANSAC. In *Joint Pattern Recognition Symposium*, pages 236–243. Springer, 2003. 6, 7
- [16] Ondrej Chum, Tomas Werner, and Jiri Matas. Two-view geometry estimation unaffected by a dominant plane. In *CVPR*, 2005. 7
- [17] Ronald Clark, Michael Bloesch, Jan Czarnowski, Stefan Leutenegger, and Andrew J. Davison. LS-Net: Learning to solve nonlinear least squares for monocular stereo. In *ECCV*, 2018. 3
- [18] Jan Czarnowski, Stefan Leutenegger, and Andrew J. Davison. Semantic texture for robust dense tracking. In *ICCV Workshops*, 2017. 1, 2
- [19] Amaël Delaunoy and Marc Pollefeys. Photometric bundle adjustment for dense multi-view 3d modeling. In *CVPR*, 2014. 2
- [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009. 12
- [21] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabenovich. SuperPoint: Self-supervised interest point detection and description. In *CVPR Workshop on Deep Learning for Visual SLAM*, 2018. 1, 2, 5, 6, 7, 9, 12, 13
- [22] Frédéric Devernay and Olivier D Faugeras. Computing differential properties of 3-D shapes from stereoscopic images without 3-D models. 1994. 2
- [23] Mihai Dusmanu, Ignacio Rocco, Tomas Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-Net: A trainable CNN for joint detection and description of local features. In *CVPR*, 2019. 1, 5, 6, 7, 12
- [24] Mihai Dusmanu, Johannes L. Schönberger, and Marc Pollefeys. Multi-View Optimization of Local Feature Geometry. In *ECCV*, 2020. 1, 2, 3, 4, 5, 6, 7, 8, 11, 12, 13
- [25] Ivan Eichhardt and Daniel Barath. Optimal multi-view correction of local affine frames. In *BMVC*, 2019. 2
- [26] Jakob Engel, Vladlen Koltun, and Daniel Cremers. Direct sparse odometry. *TPAMI*, 2017. 1, 2, 4, 5
- [27] Jakob Engel, Thomas Schöps, and Daniel Cremers. LSD-SLAM: Large-scale direct monocular SLAM. In *ECCV*, 2014. 2
- [28] Wolfgang Förstner and Eberhard Gülich. A fast operator for detection and precise location of distinct points, corners and centres of circular features. In *Proc. ISPRS intercommission conference on fast processing of photogrammetric data*, 1987. 2
- [29] Jan-Michael Frahm, Pierre Fite-Georgel, David Gallup, Tim Johnson, Rahul Raguram, Changchang Wu, Yi-Hung Jen, Enrique Dunn, Brian Clipp, Svetlana Lazebnik, et al. Building Rome on a cloudless day. In *ECCV*, 2010. 1
- [30] P. Georgel, Selim Benhimane, and Nassir Navab. A unified approach combining photometric and geometric information for pose estimation. In *BMVC*, 2008. 2
- [31] Hugo Germain, Guillaume Bourmaud, and Vincent Lepetit. S2DNet: Learning accurate correspondences for sparse-to-dense feature matching. In *ECCV*, 2020. 1, 5, 7, 11
- [32] Hugo Germain, Vincent Lepetit, and Guillaume Bourmaud. Neural Reprojection Error: Merging feature learning and camera pose estimation. In *CVPR*, 2021. 11
- [33] Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. *Robust statistics: the approach based on influence functions*. Wiley, 1986. 3
- [34] Christopher G Harris, Mike Stephens, et al. A combined corner and edge detector. In *Alvey vision conference*, 1988. 1
- [35] Jared Heinly, Johannes L Schonberger, Enrique Dunn, and Jan-Michael Frahm. Reconstructing the World* in Six Days *(as Captured by the Yahoo 100 Million Image Dataset). In *CVPR*, 2015. 1
- [36] Heiko Hirschmüller, Peter R Innocent, and Jon Garibaldi. Real-time correlation-based stereo vision with reduced border errors. *IJCV*, 2002. 2

- [37] Paul W Holland and Roy E Welsch. Robust regression using iteratively reweighted least-squares. *Communications in Statistics-theory and Methods*, 6(9):813–827, 1977. 5
- [38] Danying Hu, Daniel DeTone, and Tomasz Malisiewicz. Deep ChArUco: Dark ChArUco Marker Pose Estimation. In *CVPR*, 2019. 2
- [39] Andres Huertas and Gerard Medioni. Detection of intensity changes with subpixel accuracy using laplacian-gaussian masks. *TPAMI*, 1986. 2
- [40] Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *CVPR*, 2017. 2
- [41] Arnold Irschara, Christopher Zach, Jan-Michael Frahm, and Horst Bischof. From structure-from-motion point clouds to fast location recognition. In *CVPR*, 2009. 1
- [42] Yekeun Jeong, David Nister, Drew Steedly, Richard Szeliski, and In-So Kweon. Pushing the envelope of modern methods for bundle adjustment. *TPAMI*, 2011. 1, 5
- [43] Yuhe Jin, Dmytro Mishkin, Anastasiia Mishchuk, Jiri Matas, Pascal Fua, Kwang Moo Yi, and Eduard Trulls. Image Matching across Wide Baselines: From Paper to Practice. *IJCV*, 2020. 7
- [44] Christian Kerl, Jürgen Sturm, and Daniel Cremers. Dense visual slam for RGB-D cameras. In *IROS*, 2013. 2, 4
- [45] Kenneth Levenberg. A method for the solution of certain non-linear problems in least squares. *Quarterly of applied mathematics*, 2(2):164–168, 1944. 5, 11
- [46] Xinghui Li, Kai Han, Shuda Li, and Victor Prisacariu. Dual-resolution correspondence networks. In *NeurIPS*, 2020. 2
- [47] Yunpeng Li, Noah Snavely, Dan Huttenlocher, and Pascal Fua. Worldwide pose estimation using 3D point clouds. In *ECCV*, 2012. 1
- [48] Zhengqi Li and Noah Snavely. MegaDepth: Learning single-view depth prediction from internet photos. In *CVPR*, 2018. 12
- [49] Ce Liu, Jenny Yuen, and Antonio Torralba. SIFT Flow: Dense correspondence across scenes and its applications. *TPAMI*, 2010. 2, 7, 12
- [50] Liu Liu, Hongdong Li, and Yuchao Dai. Efficient global 2D-3D matching for camera localization in a large-scale 3D map. In *ICCV*, 2017. 1
- [51] David G Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, 2004. 1, 2, 5, 6, 7, 12
- [52] Bruce D. Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. In *IJCAI*, 1981. 1, 2
- [53] Zhaoyang Lv, Frank Dellaert, James M. Rehg, and Andreas Geiger. Taking a deeper look at the inverse compositional algorithm. In *CVPR*, 2019. 3
- [54] Daniel Martinec and Tomas Pajdla. Robust rotation and translation estimation in multiview reconstruction. In *CVPR*, 2007. 2
- [55] Yuki Ono, Eduard Trulls, Pascal Fua, and Kwang Moo Yi. LF-Net: Learning local features from images. In *NeurIPS*, 2018. 2
- [56] Seonwook Park, Thomas Schöps, and Marc Pollefeys. Illumination Change Robustness in Direct Visual SLAM. In *ICRA*, 2017. 3
- [57] Rémi Pautrat, Viktor Larsson, Martin R Oswald, and Marc Pollefeys. Online invariance selection for local feature descriptors. In *ECCV*, 2020. 1
- [58] Filip Radenovic, Johannes L Schonberger, Dinghuang Ji, Jan-Michael Frahm, Ondrej Chum, and Jiri Matas. From dusk till dawn: Modeling in the dark. In *CVPR*, 2016. 1
- [59] Jerome Revaud, Philippe Weinzaepfel, César De Souza, Noe Pion, Gabriela Csurka, Yohann Cabon, and Martin Humenberger. R2D2: Repeatable and reliable detector and descriptor. In *NeurIPS*, 2019. 1, 5, 6, 12
- [60] Ignacio Rocco, Relja Arandjelović, and Josef Sivic. Efficient neighbourhood consensus networks via submanifold sparse convolutions. In *ECCV*, 2020. 2
- [61] Ignacio Rocco, Mircea Cimpoi, Relja Arandjelović, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Neighbourhood consensus networks. In *NeurIPS*, 2018. 2
- [62] Edward Rosten and Tom Drummond. Machine learning for high-speed corner detection. In *ECCV*, 2006. 1
- [63] Paul-Edouard Sarlin. Visual localization made easy with hloc. <https://github.com/cvg/Hierarchical-Localization/>. 2, 13
- [64] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *CVPR*, 2019. 2, 13
- [65] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. In *CVPR*, 2020. 1, 7, 12, 13
- [66] Paul-Edouard Sarlin, Ajaykumar Unagar, Måns Larsson, Hugo Germain, Carl Toft, Viktor Larsson, Marc Pollefeys, Vincent Lepetit, Lars Hammarstrand, Fredrik Kahl, and Torsten Sattler. Back to the Feature: Learning robust camera localization from pixels to pose. In *CVPR*, 2021. 3, 5, 7, 12
- [67] Torsten Sattler, Will Maddern, Carl Toft, Akihiko Torii, Lars Hammarstrand, Erik Stenborg, Daniel Safari, Masatoshi Okutomi, Marc Pollefeys, Josef Sivic, Fredrik Kahl, and Tomas Pajdla. Benchmarking 6DOF outdoor visual localization in changing conditions. In *CVPR*, 2018. 2, 8, 12
- [68] Torsten Sattler, Tobias Weyand, Bastian Leibe, and Leif Kobbelt. Image retrieval for image-based localization revisited. In *BMVC*, 2012. 8, 12
- [69] Daniel Scharstein and Richard Szeliski. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *IJCV*, 2002. 2
- [70] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 2, 6, 7, 8, 12
- [71] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, 2016. 2
- [72] Thomas Schops, Torsten Sattler, and Marc Pollefeys. BAD SLAM: Bundle Adjusted Direct RGB-D SLAM. In *CVPR*, June 2019. 2

- [73] 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](#), [5](#)
- [74] Xi Shen, Fran ois Darmon, Alexei A Efros, and Mathieu Aubry. RANSAC-Flow: generic two-stage image alignment. In *ECCV*, 2020. [2](#)
- [75] Noah Snavely, Steven M Seitz, and Richard Szeliski. Modeling the world from internet photo collections. *IJCV*, 2008. [2](#)
- [76] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *CVPR*, 2018. [2](#)
- [77] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. LoFTR: Detector-free local feature matching with Transformers. *CVPR*, 2021. [2](#)
- [78] Hajime Taira, Masatoshi Okutomi, Torsten Sattler, Mircea Cimpoi, Marc Pollefeys, Josef Sivic, Tomas Pajdla, and Akihiko Torii. InLoc: Indoor Visual Localization with Dense Matching and View Synthesis. *TPAMI*, 2019. [2](#)
- [79] Chengzhou Tang and Ping Tan. BA-Net: Dense bundle adjustment network. In *ICLR*, 2019. [3](#)
- [80] Jiexiong Tang, Hanme Kim, Vitor Guizilini, Sudeep Pillai, and Rares Ambrus. Neural outlier rejection for self-supervised keypoint learning. In *ICLR*, 2020. [2](#)
- [81] Engin Tola, Vincent Lepetit, and Pascal Fua. DAISY: An efficient dense descriptor applied to wide-baseline stereo. *TPAMI*, 2009. [2](#)
- [82] Bill Triggs, Philip F McLauchlan, Richard I Hartley, and Andrew W Fitzgibbon. Bundle adjustment — a modern synthesis. In *International workshop on vision algorithms*, 1999. [1](#), [2](#), [3](#)
- [83] Prune Truong, Martin Danelljan, and Radu Timofte. GLU-Net: Global-local universal network for dense flow and correspondences. In *CVPR*, 2020. [2](#)
- [84] Michał J Tyszkiewicz, Pascal Fua, and Eduard Trulls. Disk: Learning local features with policy gradient. In *NeurIPS*, 2020. [1](#)
- [85] Andrea Vedaldi and Brian Fulkerson. VLFeat: An open and portable library of computer vision algorithms. In *ACM international conference on Multimedia*, 2010. [12](#)
- [86] Lukas Von Stumberg, Patrick Wenzel, Qadeer Khan, and Daniel Cremers. GN-Net: The Gauss-Newton loss for multi-weather relocalization. *RA-L*, 5(2):890–897, 2020. [3](#), [5](#)
- [87] Lukas Von Stumberg, Patrick Wenzel, Nan Yang, and Daniel Cremers. LM-Reloc: Levenberg-Marquardt based direct visual relocalization. In *3DV*, 2020. [3](#)
- [88] Oliver J Woodford and Edward Rosten. Large scale photometric bundle adjustment. In *BMVC*, 2020. [2](#), [4](#), [7](#), [12](#)
- [89] Binbin Xu, Andrew J. Davison, and Stefan Leutenegger. Deep probabilistic feature-metric tracking. *RA-L*, 6(1):223–230, 2021. [3](#)
- [90] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. MVSNet: Depth inference for unstructured multi-view stereo. *ECCV*, 2018. [2](#)
- [91] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent MVSNet for high-resolution multi-view stereo depth inference. In *CVPR*, 2019. [2](#)
- [92] Kwang Moo Yi, Eduard Trulls, Vincent Lepetit, and Pascal Fua. LIFT: Learned invariant feature transform. In *ECCV*, 2016. [1](#)
- [93] Baosheng Yu and Dacheng Tao. Heatmap regression via randomized rounding. *arXiv:2009.00225*, 2020. [2](#)
- [94] Zehao Yu and Shenghua Gao. Fast-MVSNet: Sparse-to-dense multi-view stereo with learned propagation and gauss-newton refinement. In *CVPR*, 2020. [3](#)
- [95] Zichao Zhang, Torsten Sattler, and Davide Scaramuzza. Reference Pose Generation for Long-term Visual Localization via Learned Features and View Synthesis. *IJCV*, 2020. [8](#), [12](#)
- [96] Qunjie Zhou, Torsten Sattler, and Laura Leal-Taixe. Patch2Pix: Epipolar-guided pixel-level correspondences. In *CVPR*, 2021. [2](#)