

Building Rome on a Cloudless Day (submitted to ECCV 2010)

Jan-Michael Frahm¹, Pierre Georgel¹, David Gallup¹, Tim Johnson¹, Rahul Raguram¹, Changchang Wu¹, Yi-Hung Jen¹, Enrique Dunn¹, Brian Clipp¹, Svetlana Lazebnik¹, Marc Pollefeys^{1,2}

¹University of North Carolina at Chapel Hill, Department of Computer Science

²ETH Zürich, Department of Computer Science

Abstract. This paper introduces an approach for dense 3D reconstruction from unregistered Internet-scale photo collections with about 3 million of images within the span of a day on a single PC (“cloudless”). Our method advances image clustering, stereo, stereo fusion and structure from motion to achieve high computational performance. We leverage geometric and appearance constraints to obtain a highly parallel implementation on modern graphics processors and multi-core architectures. This leads to two orders of magnitude higher performance on an order of magnitude larger dataset than competing state-of-the-art approaches.

1 Introduction



Fig. 1. Example models of our method from Rome (left) and Berlin (right) computed in less than 24 hrs from subsets of photo collections of 2.9 million and 2.8 million images respectively.

Recent years have seen an explosion in consumer digital photography and a phenomenal growth of community photo-sharing websites. More than 80 million photos are uploaded to the web every day,¹ and this number shows no signs of

¹ <http://royal.pingdom.com/2010/01/22/internet-2009-in-numbers>

21 slowing down. More and more of the Earth’s cities and sights are photographed 21
 22 each day from a variety of cameras, viewing positions, and angles. This has cre- 22
 23 ated a growing need for computer vision techniques that can provide intuitive 23
 24 and compelling visual representations of landmarks and geographic locations. In 24
 25 response to this challenge, the field has progressed quite impressively. Snavely et 25
 26 al. [1] were the first to demonstrate successful structure from motion (SfM) from 26
 27 Internet photo collections. Agarwal et al. [2] have performed camera registration 27
 28 and sparse 3D reconstruction starting with 150,000 images in a day on 62 cloud 28
 29 computers with 500 cores. Li et al. [3] have presented a system that combines 29
 30 appearance and multi-view geometry constraints to process tens of thousands 30
 31 of images in little more than a day on a single computer. There also exist tech- 31
 32 niques for accurate reconstruction of dense 3D models from community photo 32
 33 collections [4, 5], but they are currently much slower and more computationally 33
 34 intensive than the SfM approaches. Overall, existing systems do not measure up 34
 35 to the needs for reconstruction at city scale as, for example, a query for “Rome” 35
 36 on Flickr.com returns about 3 million images. This paper proposes a highly effi- 36
 37 cient system for camera registration combined with *dense* geometry estimation 37
 38 for city-scale reconstruction from millions of images on a single PC (no cloud 38
 39 computers = “cloudless”). The proposed system brings the computation of mod- 39
 40 els from Internet photo collections on par with state-of-the-art performance for 40
 41 reconstruction from video [6] by extending the capabilities of each step of the 41
 42 reconstruction pipeline to efficiently handle the variability and complexity of 42
 43 large-scale, unorganized, heavily contaminated datasets. 43

44 Our method efficiently combines 2D appearance and color constraints with 44
 45 3D multi-view geometry constraints to estimate the geometric relationships be- 45
 46 tween millions of images. The resulting registration serves as a basis for dense 46
 47 geometry computation using fast plane sweeping stereo [7] and a new method 47
 48 for robust and efficient depth map fusion. We take advantage of the appearance 48
 49 and geometry constraints to achieve parallelization on graphics processors and 49
 50 multi-core architectures. All timings in the paper are obtained on a PC with 50
 51 dual Intel quadcore Xeon 3.33 Ghz processors, four NVidia 295GTX commodity 51
 52 graphics cards,² 48 GB RAM and a 1 TB solid state hard drive for data storage. 52
 53 The major steps of our method are: 53

54 **1) Appearance-based clustering with small codes** (Sec. 3.1): Similarly to 54
 55 Li et al. [3] we use the *gist* feature [8] to capture global image appearance. The 55
 56 complexity of the subsequent geometric registration is reduced by clustering the 56
 57 *gist* features to obtain a set of *canonical* or *iconic* views [3]. In order to be able to 57
 58 fit several million *gist* features in GPU-memory, we compress them to compact 58
 59 binary strings using a locality sensitive scheme [9–11]. We then cluster them 59
 60 based on Hamming distance with the *k*-medoids algorithm [12] implemented 60
 61 on the GPU. To our knowledge, this is the first application of small codes in 61
 62 the style of [11] outside of proof-of-concept recognition settings, and the first 62
 63 demonstration of their effectiveness for large-scale clustering problems. 63

² By the time of ECCV this will correspond to two graphics cards of the next gener-
 ation that will then be available, making this a state-of-the-art gaming computer.

64 **2) Geometric cluster verification** (Sec. 3.2) is used to identify in each cluster 64
 65 a “core” set images with mutually consistent epipolar geometry using a fast 65
 66 RANSAC method [13]. All other cluster images are verified to match to either of 66
 67 the “core” images, and the ones found to be inconsistent are removed. Finally, 67
 68 we select a single iconic view as the best representative of the cluster. Given that 68
 69 geo-location is available for many images in the Internet photo collection, we are 69
 70 typically able to geo-locate a large fraction of our clusters ($> 50\%$). 70

71 **3) Local iconic scene graph reconstruction** (Sec. 3.3) establishes the skele- 71
 72 ton registration of the iconic images in the different locations. We use vocabulary 72
 73 tree search [14] and clustering based on geo-location and image appearance to 73
 74 identify neighboring iconics. Both of these strategies typically lead to sets of lo- 74
 75 cally connected images corresponding to the different geographically separated 75
 76 sites of the city. We dub these sets *local* iconic scene graphs. These graphs are 76
 77 extended by registering additional views from the iconic clusters. Our registra- 77
 78 tion method uses incremental SfM combined with periodic bundle adjustment 78
 79 to mitigate errors. 79

80 **5) Dense model computation** (Sec. 3.4) uses all registered views in the local 80
 81 iconic scene graphs to obtain dense scene geometry for the captured sites. Tak- 81
 82 ing advantage of the initial appearance-based image grouping, we deploy fast 82
 83 plane sweeping stereo to obtain depth maps from each iconic cluster. To mini- 83
 84 mize the computational load we perform visibility-based view selection for the 84
 85 dense depth map computation. Then we apply a novel extension to a depthmap 85
 86 fusion method to obtain a watertight scene representation from the noisy but 86
 87 redundant depth maps. 87
 88 88

89 2 Previous Work 89

90 Our method is the first system performing dense modeling from Internet photo 90
 91 collections consisting of millions of images. Systems for urban reconstruction 91
 92 from video have been proposed in [15, 6], with [6] achieving real-time dense 3D 92
 93 reconstruction. However, modeling from video is inherently much more efficient 93
 94 as it takes advantage of spatial proximity between the camera positions of suc- 94
 95 cessive frames, whereas the spatial relationships between images in a community 95
 96 photo collection are unknown a priori, and in fact, 40% to 60% of images in such 96
 97 collections turn out to be irrelevant clutter [3]. 97

98 The first approach for organizing unordered image collections was proposed 98
 99 by Schaffalitzky and Zisserman [16]. Sparse 3D reconstruction of landmarks from 99
 100 Internet photo collections was first addressed by the *Photo Tourism* system [17], 100
 101 which achieves high-quality results through exhaustive pairwise image matching 101
 102 and frequent global bundle adjustment. Neither one of these steps is very scal- 102
 103 able, so in practice, the Photo Tourism system can be applied to a few thousand 103
 104 images at most. Aiming at scalability, Snavely et al. [18] construct *skeletal sets* of 104
 105 images whose reconstruction approximates the full reconstruction of the whole 105

106 dataset. However, computing these sets still requires initial exhaustive pairwise 106
 107 image matching. Agarwal et al. [2] parallelize the matching process and use 107
 108 approximate nearest neighbor search and query expansion [19] on a cluster of 62 108
 109 machines each one comparable to our single PC. With that single PC, we tackle 109
 110 an order of magnitude more data in the same amount of time. 110

111 The speed of our approach is a result of efficient early application of 2D 111
 112 appearance-based constraints, similarly to the approach of Li et al. [3]. But 112
 113 our system extends [3] to successfully process two orders of magnitude more 113
 114 data by parallelizing the computation on graphics processors and multi-core 114
 115 architectures. We summarize the dataset and select *iconic images* using 2D image 115
 116 appearance as a prerequisite for efficient camera registration. This is the opposite 116
 117 of the approach of Simon et al. [20], who treat scene summarization as a by- 117
 118 product of 3D reconstruction and select canonical views through clustering the 118
 119 3D camera poses. While our method of image organization is initially looser than 119
 120 that of [20], it provides a powerful pre-selection mechanism for advancing the 120
 121 reconstruction efficiency significantly. 121

122 After selecting the iconic images, the next step of our system is to discover the 122
 123 geometric relationships between them and register them together through SfM. 123
 124 Li et al. [3] deployed a vocabulary tree [14] to rapidly find related iconics. Our 124
 125 system can use a vocabulary tree in the absence of geo-location information for 125
 126 the iconics. If this information is available, we use it to help identify possible links 126
 127 between iconics. The latter approach is more efficient since it avoids building 127
 128 the vocabulary tree. Note that the methods of [21, 22] are also applicable to 128
 129 discovering spatial relationships in large collections of data. 129

130 To perform SfM on the set of iconic images, Li et al. [3] partitioned the 130
 131 iconic scene graph into multiple connected components and performed SfM on 131
 132 each component. In contrast, we do not cut the iconic scene graph, as such 132
 133 an approach is prone to excessive fragmentation of scene models. Instead, we 133
 134 use a growing strategy combined with efficient merging and periodic bundle 134
 135 adjustment to obtain higher-quality, more complete models. Our method is open 135
 136 to use techniques for out-of-core bundle-adjustment [23], which take advantage 136
 137 of the uneven viewpoint distribution in photo collections. 137

138 Given the registered viewpoints recovered by SfM, we next perform multi- 138
 139 view stereo to get dense 3D models. The first approach demonstrating dense 139
 140 modeling from photo collections was proposed by Goesele et al. [4]. It uses per- 140
 141 pixel view selection and patch growing to obtain a set of surface elements, which 141
 142 are then regularized into a Poisson surface model. However, this approach does 142
 143 not make it easy to provide textures for the resulting models. Recently, Furukawa 143
 144 et al. [24] proposed a dense reconstruction method from large-scale photo collec- 144
 145 tions using view clustering to initialize the PMVS approach [25]. This method 145
 146 computes a dense model from approximately 13,000 images in about two days on 146
 147 a single computer assuming known camera registration. Our proposed method 147
 148 uses an extended version of Yang and Pollefeys stereo [7] combined with novel 148
 149 multi-layer depth map fusion [26]. While achieving comparable quality, it com- 149

putationally outperforms [4, 24] by achieving modeling on a single PC within less than an hour instead of multiple days.

3 The Approach

In this section we will describe our proposed method. Section 3.1 discusses the initial appearance-based clustering, and Section 3.2 discusses our method for efficient geometric verification of the resulting clusters. Section 3.3 explains our SfM scheme, and Section 3.4 explains our stereo fusion that leverages the appearance-based clustering for dense 3D model generation.

3.1 Appearance-Based Clustering with Small Codes

Similarly to Li et al. [3], we begin by computing a global appearance descriptor for every image in the dataset. We generate a gist feature [8] for each image by computing oriented edge responses at three scales (with 8, 8 and 4 orientations, respectively), aggregated to a 4×4 spatial resolution. To ensure better grouping of views for our dense reconstruction method, we concatenate the gist with a subsampled RGB image at 4×4 spatial resolution. Both the gist and the color parts of the descriptor are rescaled to have unit norm. The combined descriptor has 368 dimensions, and it is computed on the 8 GPU cores at a rate of 781Hz³.

The next step is to cluster the gist descriptors to obtain groups of images consistent in appearance. For efficiency in the clustering we aim at a GPU-based implementation, given the inherent parallelism in the distance computation of clustering algorithms like k -means and k -medoids. Since it is impossible to cluster up to 2.8 million 368-dimensional double-precision vectors in the GPU memory of 768 MB, we have chosen to compress the descriptors to much shorter binary strings, such that the Hamming distances between the compressed strings approximate the distances between the original descriptors. To this end, we have implemented on the GPU the locality sensitive binary code (LSBC) scheme of Raginsky and Lazebnik [10], in which the i th bit of the code for a descriptor vector \mathbf{x} is given by $\varphi_i(\mathbf{x}) = \text{sgn}[\cos(\mathbf{x} \cdot \mathbf{r}_i + b_i) + t_i]$, where $\mathbf{r} \sim \text{Normal}(0, \gamma I)$, $b_i \sim \text{Unif}[0, 2\pi]$, and $t_i \sim \text{Unif}[-1, 1]$ are randomly chosen code parameters. As shown in [10], as the number of bits in the code increases, the *normalized* Hamming distance (i.e., Hamming distance divided by code length) between two binary strings $\varphi(\mathbf{x})$ and $\varphi(\mathbf{y})$ approximates $(1 - K(\mathbf{x}, \mathbf{y}))/2$, where $K(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2 / 2}$ is a Gaussian kernel between \mathbf{x} and \mathbf{y} . We have compared the LSBC scheme with a simple locality sensitive hashing (LSH) scheme for unit norm vectors where the i th bit of the code is given by $\text{sgn}(\mathbf{x} \cdot \mathbf{r}_i)$ [9]. As shown in the recall-precision plots in Figure 2, LSBC does a better job of preserving the distance relationships of our descriptors.

We have found that $\gamma = 4.0$ works well for our data, and that the code length of 512 offers the best tradeoff between approximation accuracy and memory

³ code in preparation for release

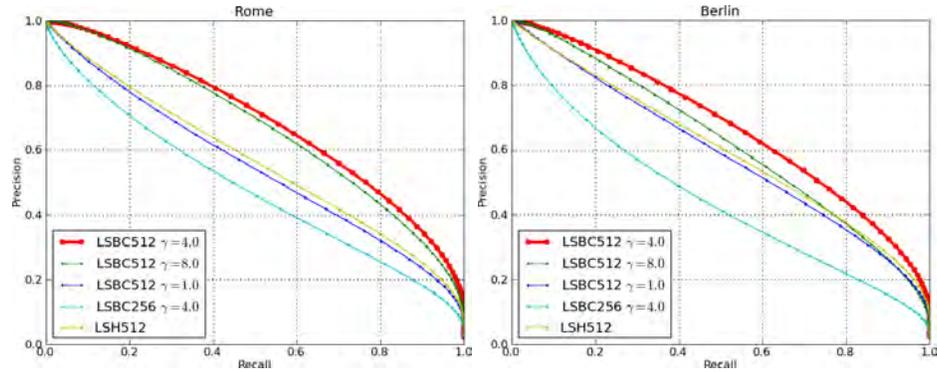


Fig. 2. Comparison of LSH coding scheme [9] and LSBC [10] scheme with different settings for γ and code size on Rome data (left) and Berlin data (right). These plots show the recall and precision of nearest-neighbor search with Hamming distance on binary codes for retrieving the “true” k nearest neighbors according to Euclidean distance on the original gist features (k is our average cluster size, 28 for Rome and 26 for Berlin). For our chosen code size of 512, the LSBC scheme with $\gamma = 4$ outperforms LSH.



Fig. 3. Images closest to the center of one cluster from Rome.

189 usage. To give an idea of the memory savings afforded by this scheme, at 32 189
 190 bytes per dimension, each original descriptor takes up 11,776 bytes, while the 190
 191 corresponding binary vector takes up only 64 bytes, thus achieving a compression 191
 192 factor of 184. With this amount of compression, we can cluster up to about 192
 193 4 million images on our memory budget of 768 MB, vs. only a few hundred 193
 194 thousand images in the original GIST representation. An example of a gist cluster 194
 195 is shown in Figure 3. 195

196 For clustering the binary codevectors with the Hamming distance, we have 196
 197 implemented the k -medoids algorithm [12] on the GPU. Like k -means, k -medoids 197
 198 alternates between updating cluster centers and cluster assignments, but unlike 198
 199 k -means, it forces each cluster center to be an element of the dataset. For every 199
 200 iteration, we compute the Hamming distance matrix between the binary codes 200
 201 of all images and those that correspond to the medoids. Due to the size of the 201

dataset and number of cluster centers, this distance matrix must be computed piecewise, as it would require roughly 1050 GB to store on the GPU.

A generally open problem for clustering in general is how to initialize the cluster centers, as the initialization can have a big effect on the end results. We found that images with available geo-location information (typically 10 – 15% of our city-scale datasets) provide a good sampling of the points of interest (see Figure 4). Thus, we first cluster the codevectors of images with available geo-location into k_{geo} clusters initialized randomly. Then we use the resulting centers together with additional k_{rand} random centers to initialize the clustering of the complete dataset (in all our experiments $k_{geo} = k_{rand}$). From Table 2 it can be seen that we gain about 20% more geometrically consistent images by this initialization strategy.

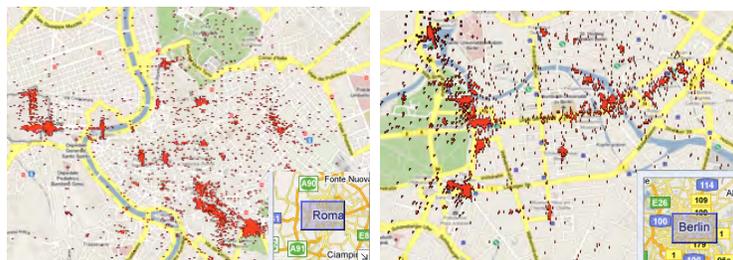


Fig. 4. Geo-tag density map for Rome (left) and Berlin (right).

3.2 Geometric Verification

The clusters obtained in the previous step consist of images that are visually similar but may be geometrically and semantically inconsistent. Since our goal is to reconstruct scenes with stable 3D structure, we next enforce geometric consistency for images within a cluster. A cluster is deemed to be consistent if it has at least n images with a valid pairwise epipolar geometry. This is determined by selecting an initial subset of n images (those closest to the cluster medoid) and estimating the two-view geometry of all the pairs in this subset while requiring at least m inliers (in all our experiments we use $n = 4$, $m = 18$). Inconsistent images within the subset are replaced by others until n valid images are found, or all cluster images are exhausted and the cluster is rejected.

The computation of two-view epipolar geometry is performed as follows. We extract SIFT features [27] using an efficient GPU implementation,⁴ processing 1024 × 768 images at up to 16.8 Hz on a single GPU. In the interest of computational efficiency and memory bandwidth, we limit the number of features

⁴ <http://www.cs.unc.edu/ccwu/siftgpu>



Fig. 5. The geometrically verified cluster showing the Coliseum in Rome.

229 extracted to 4000 per image. Next, we calculate the putative SIFT matches for 229
 230 each image pair. This computationally demanding process (which could take a 230
 231 few seconds per pair on the CPU) is cast as a matrix multiplication problem 231
 232 on multiple GPUs (with a speedup of three orders of magnitude to 740 Hz), 232
 233 followed a subsequent distance ratio test [27] to identify likely correspondences. 233

234 The putative matches are verified by estimation of the fundamental matrix 234
 235 with the 7-point algorithm [28] and ARRSAC [13], which is a robust estimation 235
 236 framework designed for efficient real-time operation. For small inlier ratios, even 236
 237 ARRSAC significantly degrades in performance. However, we have observed that 237
 238 of all registered images in the three datasets a significant fraction had inlier 238
 239 ratios above 50% (e.g., for San Marco, this fraction is 72%). We use this to 239
 240 our advantage by limiting the maximum number of tested hypotheses to 400 in 240
 241 ARRSAC, which corresponds to inlier ratio of approximately 50%. To improve 241
 242 registration performance, we take the best solution deemed promising by the 242
 243 SPRT test of ARRSAC, and perform a *post hoc* refinement procedure. The latter 243
 244 enables us to recover a significant fraction of solutions with less than 50% inlier 244
 245 ratio. Comparing the number of registered images by the standard ARRSAC and 245
 246 the number of images registered by our modified procedure shows a loss of less 246
 247 than 3% for Rome and less than 5% for Berlin of registered images while having 247
 248 an approximately two- to five-fold gain in speed. This result makes intuitive 248
 249 sense: it has been observed [18, 3] that community photo collections contain a 249
 250 tremendous amount of viewpoint overlap and redundancy, which is particularly 250
 251 pronounced at the scale at which we operate. 251

252 We choose a representative or “iconic” image for each verified cluster as the 252
 253 image with the most inliers to the other $n - 1$ top images. Afterwards all other 253
 254 cluster images are only verified with respect to the iconic image. Our system 254
 255 processes all the appearance-based clusters independently using 16 threads on 8 255
 256 CPU cores and 8 GPU cores. In particular, the process of putative matching is 256
 257 distributed over multiple GPUs, while the robust estimation of the fundamental 257
 258 matrix utilizes the CPU cores. This enables effective utilization of all available 258
 259 computing resources and gives a significant speedup to about 480 Hz verification 259
 260 rate an example is shown in Figure 5 260

261 If user provided geo-tags are available (all our city datasets have between 10% 261
 262 and 15% geo-tagged images) we use them to geo-locate the clusters. Our geo- 262
 263 location evaluates the pairwise distances of all geo-tagged image in the iconic 263
 264 cluster. Then it performs a weighted voting on the locations of all images within a 264

265 spatial proximity of the most central image as defined by the pairwise distances. 265
 266 This typically provides a geo-location for about two thirds of the iconic clusters. 266

267 3.3 Local Iconic Scene Graph Reconstruction 267

268 After identifying the geometrically consistent clusters, we need to establish pair- 268
 269 wise relationships between the iconics. Li et al. [3] introduced the *iconic scene* 269
 270 *graph* to encode these relationships. We use the same concept but identify mul- 270
 271 tiple *local* iconic scene graphs corresponding to the multiple geographic sites 271
 272 within each dataset. This keeps the complexity low despite the fact that our sets 272
 273 of iconics are comparable in size to the entire datasets of [3]. 273

274 We experimented with two different schemes for efficiently obtaining candi- 274
 275 date iconic pairs for geometric verification. The first scheme is applicable in 275
 276 the absence of any geo-location. It is based on building a vocabulary tree index 276
 277 for the SIFT features of the iconics, and using each iconic to query for related 277
 278 images. The drawback of this scheme is that the mapping of the vocabulary 278
 279 tree has to be rebuilt specifically for each set of iconics, imposing a significant 279
 280 overhead on the computation. The second scheme avoids this overhead by using 280
 281 geo-location of iconic clusters. In this scheme, the candidate pairs are defined 281
 282 as all pairs within a certain distance s of each other (in all our experiments set 282
 283 to $s = 150$ m). As for the iconics lacking geo-location, they are linked to their 283
 284 l -nearest neighbors ($l = 10$ in all experiments) in the binarized gist descriptor 284
 285 space (the distance computation uses GPU-based nearest-neighbor search as in 285
 286 the k -medoids clustering). We have found this second scheme to be more effi- 286
 287 cient whenever geo-location is available for a sufficient fraction of the iconics (as 287
 288 in our Rome and Berlin datasets). For both schemes, all the candidate iconic 288
 289 pairs are geometrically verified as described in Section 3.2, and the pairs with a 289
 290 valid epipolar geometry are connected by an edge. Each connected set of iconics 290
 291 obtained in this way is a *local iconic scene graph*, usually corresponding to a 291
 292 distinct geographic site in a city. 292

293 Next, each local iconic scene graph is processed independently to obtain a 293
 294 camera registration and a sparse 3D point cloud using an incremental approach. 294
 295 The algorithm picks the pair of iconic images whose epipolar geometry given by 295
 296 the essential matrix (computed as similarly to Section 3.2) has the highest inlier 296
 297 number and delivers a sufficiently low reconstruction uncertainty, as computed 297
 298 by the criterion of [29]. Obtaining a metric two-view reconstruction requires a 298
 299 known camera calibration, which we either obtain from the EXIF-data of the 299
 300 iconics (there are 34% EXIF based calibrations for the Berlin dataset and 40% 300
 301 for Rome), or alternatively we approximate the calibration by assuming a popular 301
 302 viewing angle for the camera model. The latter estimate typically approximates 302
 303 the true focal length within the error bounds of successfully executing the five- 303
 304 point method [30]. To limit drift after inserting i new iconics, the 3D sub-model 304
 305 and camera parameters are optimized by a sparse bundle adjustment [31]. The 305
 306 particular choice of i is not critical and in all our experiments we use $i = 50$. If 306
 307 no new images can be registered into the current sub-model, the process starts 307

afresh by picking the next best pair of iconics not yet registered to any sub-model. Note that we intentionally construct multiple sub-models that may share some images. We use these images to merge newly completed sub-models with existing ones whenever sufficient 3D matches exist. The merging step again uses ARRISAC [32] to robustly estimate a similarity transformation based on the identified 3D matches.

In the last stage of the incremental reconstruction algorithm, we complete the model by incorporating non-iconic images from iconic clusters of the registered iconics. This process takes advantage of the feature matches between the non-iconic images and their respective iconics known from the geometric verification (Section 3.2). The 2D matches between the image and its iconic determine 2D-3D correspondences between the image and the 3D model into which the iconic is registered, and ARRISAC is once again used to determine the camera pose. Detailed results of our 3D reconstruction algorithm are shown in Figure 6, and timings in Table 1.

3.4 Dense geometry estimation

Once the camera poses have been recovered, the next step is to recover the surface of the scene, represented as a polygonal mesh, and to reconstruct the surface color represented as a texture map. We use a two-phase approach for surface reconstruction: first, recover depthmaps for a select number of images, and second, fuse the depthmaps into a final surface model.

One of the major challenges of stereo from Internet photo collections is appearance variation. Previous approaches [4, 33] take great care to select compatible views for stereo matching. We use the clustering approach from Section 3.1 to cluster all images registered in the local iconic scene graph. Since our gist descriptor encodes color, the resulting clusters are color-consistent. The availability of color-consistent images within a spatially confined area enables us to use traditional stereo methods and makes dense reconstruction a simpler task than might otherwise be thought. We use a GPU-accelerated plane sweep stereo [34] with a 3×3 normalized cross-correlation matching kernel. Our stereo deploys 20 matching views, and handles occlusions (and other outliers) through taking the best 50% of views per pixel as suggested in [35]. We have found that within a set of 20 views, non-identical views provide a sufficient baseline for accurate depth computation.

We adapted the vertical heightmap approach of [36] for depthmap fusion to handle geometrically more complex scenes. This method is intended to compute a watertight approximate surface model. The approach assumes that the vertical direction of the scene is known beforehand. For community photo collections, this direction can be easily obtained using the approach of [37] based on the assumption that most photographers will keep the camera's x -axis perpendicular the vertical direction. The heightmap is computed by constructing an occupancy grid over a volume of interest. All points below the heightmap surface are considered full and all points above are considered empty. Each vertical column of the grid is computed independently. For each vertical column, occupancy votes

Dataset	Gist & Clustering	SIFT & Geom. verification	Local iconic scene graph	Dense	total time
Rome & geo	1:35 hrs	11:36 hrs	8:35 hrs	1:58 hrs	23:53 hrs
Berlin & geo	1:30 hrs	11:46 hrs	7:03 hrs	0:58 hrs	21:58 hrs
San Marco	0:03 hrs	0:24 hrs	0:32 hrs	0:07 hrs	1:06 hrs

Table 1. Computation times (hh:mm hrs) for the photo collection reconstruction for the Rome dataset using geo-tags, the Berlin dataset with geo-tags, and the San Marco dataset without geo-tags.

Dataset	total	LSBC clusters	#images			
			iconics	verified	3D models	largest model
Rome & geo	2,884,653	100, 000	21,651	306788	63905	5671
Rome	2,884,653	100, 000	17874	249689	-	-
Berlin & geo	2,771,966	100, 000	14664	124317	31190	3158
San Marco	44, 229	4,429	890	13604	1488	721

Table 2. Image sizes for the the Rome dataset, the Berlin dataset, and the San Marco dataset.

352 are accumulated from the depthmaps. Points between the camera center and the 352
353 depth value receive empty votes, and points beyond the depth value receive a 353
354 full vote with a weight that falls off with distance. Then a height value is deter- 354
355 mined that minimizes the number of empty votes above and the number of full 355
356 votes below. Our extension is to allow the approach to have multiple connected 356
357 “segments” within the column, which provides higher quality mesh models while 357
358 maintaining the regularization properties of the original approach. A polygonal 358
359 mesh is then extracted from the heightmap and texture maps are generated from 359
360 the color images. The heightmap model is highly robust to noise and it can be 360
361 computed very efficiently on the GPU. 361

362 The resolution of the height map is determined by the median camera-to- 362
363 point distance which is representative of the scale of the scene and the accuracy 363
364 of the depth measurements. The texture of the mesh models is then computed as 364
365 the mean of all images observing the geometry. Runtimes are provided in Table 365
366 1. 366

367 4 Conclusions 367

368 This paper demonstrated the first system able to deliver dense geometry for 368
369 Internet scale photo collections with millions of images of an entire city within 369
370 the span of a day on a single PC. Our novel methods extend to the scale of 370
371 millions of images state-of-the-art methods for appearance-based clustering [3], 371
372 robust estimation [32], and stereo fusion [36]. To successfully handle reconstruction 372
373 problems of this magnitude, we have incorporated novel system components 373
374 for clustering of small codes, geo-location of iconic images through their clusters 374
375 for clustering of small codes, geo-location of iconic images through their clusters 375

376 through appearance-based clustering. Beyond making algorithmic changes, we 376
 377 significantly improve performance by leveraging the constraints from appearance 377
 378 clustering and location independence to parallelize the processing on modern 378
 379 multi-core CPUs and commodity graphics cards. 379

380 References 380

- 381 1. Snavely, N., Seitz, S.M., Szeliski, R.: Photo tourism: Exploring photo collections 381
 382 in 3d. In: SIGGRAPH. (2006) 835–846 382
- 383 2. Agarwal, S., Snavely, N., Simon, I., Seitz, S.M., Szeliski, R.: Building Rome in a 383
 384 day. In: ICCV. (2009) 384
- 385 3. Li, X., Wu, C., Zach, C., Lazebnik, S., Frahm, J.M.: Modeling and recognition of 385
 386 landmark image collections using iconic scene graphs. In: ECCV. (2008) 386
- 387 4. Goesele, M., Snavely, N., Curless, B., Hoppe, H., Seitz, S.M.: Multi-view stereo 387
 388 for community photo collections. In: ICCV. (2007) 388
- 389 5. Furukawa, Y., Curless, B., Seitz, S.M., Szeliski, R.: Towards internet-scale multi- 389
 390 view stereo. In: CVPR. (2010) 390
- 391 6. Pollefeys, M., Nister, D., Frahm, J.M., Akbarzadeh, A., Mordohai, P., Clipp, B., 391
 392 Engels, C., Gallup, D., Kim, S.J., Merrell, P., Salmi, C., Sinha, S., Talton, B., 392
 393 Wang, L., Yang, Q., Stewenius, H., Yang, R., Welch, G., Towles, H.: Detailed 393
 394 real-time urban 3d reconstruction from video. *International Journal of Computer 394*
 395 *Vision special issue on Modeling Large-Scale 3D Scenes* (2008) 395
- 396 7. Yang, R., Pollefeys, M.: Multi-resolution real-time stereo on commodity graphics 396
 397 hardware. In: *Int. Conf. on Computer Vision and Pattern Recognition*. (2003) 397
 398 211–217 398
- 399 8. Oliva, A., Torralba, A.: Modeling the shape of the scene: a holistic representation 399
 400 of the spatial envelope. *IJCV* **42** (2001) 145–175 400
- 401 9. Andoni, A., Indyk, P.: Near-optimal hashing algorithms for approximate nearest 401
 402 neighbor in high dimensions. *Communications of the ACM* **51** (2008) 117–122 402
- 403 10. Raginsky, M., Lazebnik, S.: Locality sensitive binary codes from shift-invariant 403
 404 kernels. In: NIPS. (2009) 404
- 405 11. Torralba, A., Fergus, R., Weiss, Y.: Small codes and large databases for recognition. 405
 406 In: CVPR. (2008) 406
- 407 12. Kaufman, L., Rousseeuw, P.J.: *Finding Groups in Data: An Introduction to Cluster 407*
 408 *Analysis*. Wiley (1990) 408
- 409 13. Raguram, R., Frahm, J.M., Pollefeys, M.: A comparative analysis of RANSAC 409
 410 techniques leading to adaptive real-time random sample consensus. In: ECCV. 410
 411 (2008) 411
- 412 14. Nister, D., Stewenius, H.: Scalable recognition with a vocabulary tree. In: CVPR. 412
 413 (2006) 413
- 414 15. Cornelis, N., Cornelis, K., Van Gool, L.: Fast compact city modeling for navigation 414
 415 pre-visualization. In: *Int. Conf. on Computer Vision and Pattern Recognition*. 415
 416 (2006) 416
- 417 16. Schaffalitzky, F., Zisserman, A.: Multi-view matching for unordered image sets, 417
 418 or "how do i organize my holiday snaps?". In: *ECCV '02: Proceedings of the 7th 418*
 419 *European Conference on Computer Vision-Part I*, London, UK, Springer-Verlag 419
 420 (2002) 414–431 420
- 421 17. Snavely, N., Seitz, S.M., Szeliski, R.: Modeling the world from Internet photo 421
 422 collections. *International Journal of Computer Vision* **80** (2008) 189–210 422

- 423 18. Snavely, N., Seitz, S.M., Szeliski, R.: Skeletal sets for efficient structure from 423
424 motion. In: CVPR. (2008) 424
- 425 19. Chum, O., Philbin, J., Sivic, J., Isard, M., Zisserman, A.: Total recall: Automatic 425
426 query expansion with a generative feature model for object retrieval. In: ICCV. 426
427 (2007) 427
- 428 20. Simon, I., Snavely, N., Seitz, S.M.: Scene summarization for online image collec- 428
429 tions. In: ICCV. (2007) 429
- 430 21. Chum, O., Matas, J.: Web scale image clustering: Large scale discovery of spatially 430
431 related images. Technical report, CTU-CMP-2008-15 (2008) 431
- 432 22. Philbin, J., Zisserman, A.: Object mining using a matching graph on very large 432
433 image collections. In: Proceedings of the Indian Conference on Computer Vision, 433
434 Graphics and Image Processing. (2008) 434
- 435 23. Ni, K., Steedly, D., Dellaert, F.: Out-of-core bundle adjustment for large-scale 3d 435
436 reconstruction. In: ICCV. (2007) 436
- 437 24. Furukawa, Y., Curless, B., Seitz, S.M., Szeliski, R.: Towards internet-scale multi- 437
438 view stereo. In: In Proceedings of IEEE CVPR. (2010) 438
- 439 25. Furukawa, Y., Ponce, J.: Accurate, dense, and robust multi-view stereopsis. In: 439
440 PAMI. (2009) 440
- 441 26. Gallup, D., Frahm, J.M., Pollefeys, M.: A heightmap model for efficient 3d recon- 441
442 struction from street-level video. In: In Proceedings of 3DPVT. (2010) 442
- 443 27. Lowe, D.: Distinctive image features from scale-invariant keypoints. *IJCV* **60** 443
444 (2004) 91–110 444
- 445 28. Hartley, R.I., Zisserman, A.: Multiple View Geometry in Computer Vision. Second 445
446 edn. Cambridge University Press (2004) 446
- 447 29. Beder, C., Steffen, R.: Determining an initial image pair for fixing the scale of a 447
448 3d reconstruction from an image sequence. In: Proc. DAGM. (2006) 657–666 448
- 449 30. Nistér, D.: An efficient solution to the five-point relative pose problem. *PAMI* **26** 449
450 (2004) 756–770 450
- 451 31. Lourakis, M., Argyros, A.: The design and implementation of a generic sparse 451
452 bundle adjustment software package based on the Levenberg-Marquardt algorithm. 452
453 Technical Report 340, Institute of Computer Science - FORTH (2004) 453
- 454 32. Raguram, R., Lazebnik, S.: Computing iconic summaries of general visual concepts. 454
455 In: Workshop on Internet Vision CVPR. (2008) 455
- 456 33. Furukawa, Y., Seitz, S.: Towards internet-scale multi-view stereo. In: In Proceed- 456
457 ings of IEEE CVPR. (2010) 457
- 458 34. Kim, S., Gallup, D., Frahm, J., Akbarzadeh, A., Yang, Q., Yang, R., Nister, D., 458
459 Pollefeys, M.: Gain adaptive real-time stereo streaming. In: International Confer- 459
460 ence on Computer Vision Systems (ICVS). (2007) 460
- 461 35. Kang, S., Szeliski, R., Chai, J.: Handling occlusions in dense multi-view stereo. In: 461
462 CVPR. (2001) 462
- 463 36. NN: blank for double blind review. In: NN. (2010) 463
- 464 37. Szeliski, R.: Image alignment and stitching: A tutorial. In: Microsoft Research 464
465 Technical Report. (2004) 465

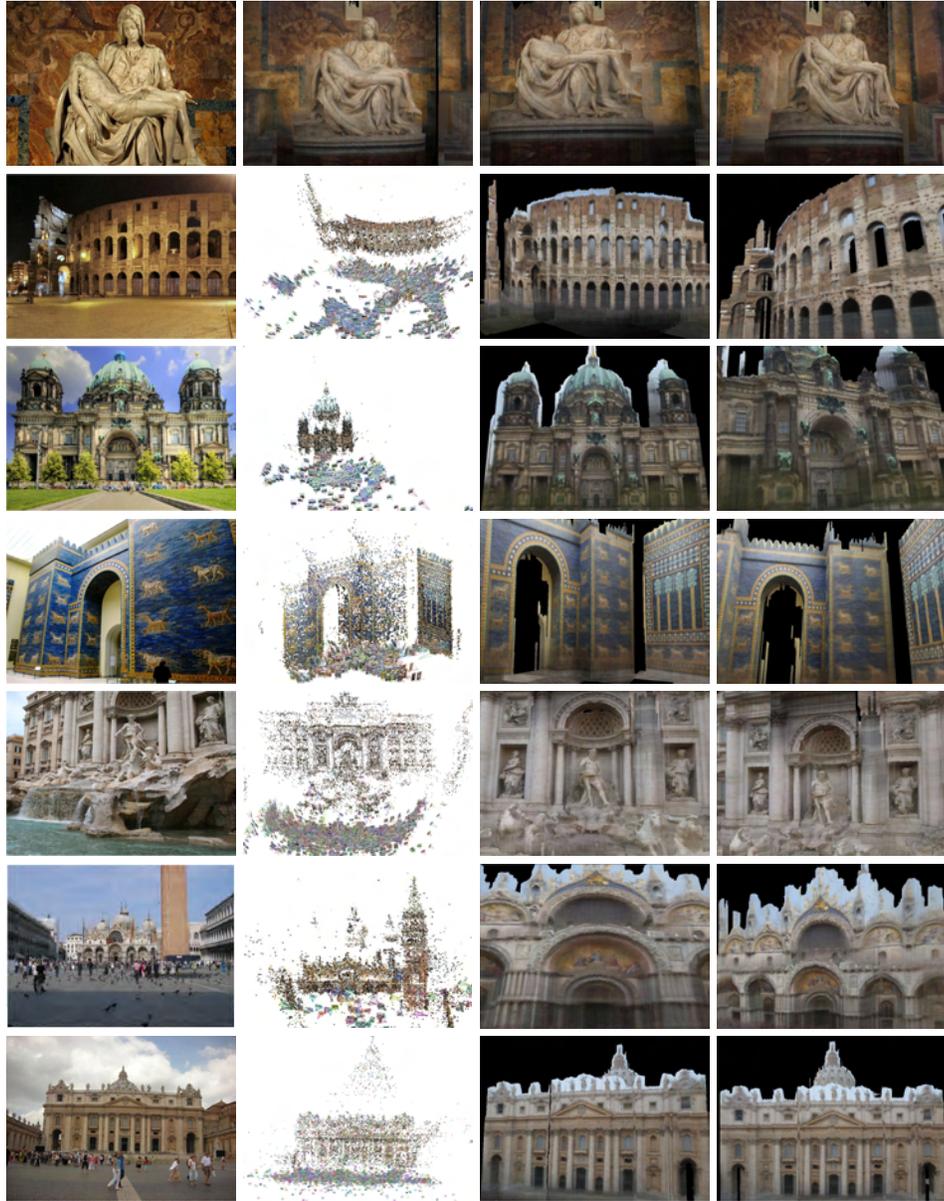


Fig. 6. Original images, local iconic scene graph and 3D model.