

1 Building Rome on a Cloudless Day (ECCV 2010) 1

2 Jan-Michael Frahm¹, Pierre Georgel¹, David Gallup¹, Tim Johnson¹, Rahul 2
3 Raguram¹, Changchang Wu¹, Yi-Hung Jen¹, Enrique Dunn¹, Brian Clipp¹, 3
4 Svetlana Lazebnik¹, Marc Pollefeys^{1,2} 4

5 ¹University of North Carolina at Chapel Hill, Department of Computer Science 5
6 ²ETH Zürich, Department of Computer Science 6

7 **Abstract.** This paper introduces an approach for dense 3D reconstruction 7
8 from unregistered Internet-scale photo collections with about 3 mil- 8
9 lion of images within the span of a day on a single PC (“cloudless”). Our 9
10 method advances image clustering, stereo, stereo fusion and structure 10
11 from motion to achieve high computational performance. We leverage 11
12 geometric and appearance constraints to obtain a highly parallel imple- 12
13 mentation on modern graphics processors and multi-core architectures. 13
14 This leads to two orders of magnitude higher performance on an order 14
15 of magnitude larger dataset than competing state-of-the-art approaches. 15

16 1 Introduction 16



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.

17 Recent years have seen an explosion in consumer digital photography and a 17
18 phenomenal growth of community photo-sharing websites. More than 80 million 18
19 photos are uploaded to the web every day,¹ and this number shows no signs of 19
20 slowing down. More and more of the Earth’s cities and sights are photographed 20

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

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

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

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

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

63 **2) Geometric cluster verification** (Sec. 3.2) is used to identify in each cluster 63
 64 a “core” set images with mutually consistent epipolar geometry using a fast 64
 65 RANSAC method [13]. All other cluster images are verified to match to either of 65
 66 the “core” images, and the ones found to be inconsistent are removed. Finally, 66
 67 we select a single iconic view as the best representative of the cluster. Given that 67
 68 geo-location is available for many images in the Internet photo collection, we are 68
 69 typically able to geo-locate a large fraction of our clusters ($> 50\%$). 69

70 **3) Local iconic scene graph reconstruction** (Sec. 3.3) establishes the skele- 70
 71 ton registration of the iconic images in the different locations. We use vocabulary 71
 72 tree search [14] and clustering based on geo-location and image appearance to 72
 73 identify neighboring iconics. Both of these strategies typically lead to sets of lo- 73
 74 cally connected images corresponding to the different geographically separated 74
 75 sites of the city. We dub these sets *local* iconic scene graphs. These graphs are 75
 76 extended by registering additional views from the iconic clusters. Our registra- 76
 77 tion method uses incremental SfM combined with periodic bundle adjustment 77
 78 to mitigate errors. 78

79 **5) Dense model computation** (Sec. 3.4) uses all registered views in the local 79
 80 iconic scene graphs to obtain dense scene geometry for the captured sites. Tak- 80
 81 ing advantage of the initial appearance-based image grouping, we deploy fast 81
 82 plane sweeping stereo to obtain depth maps from each iconic cluster. To mini- 82
 83 mize the computational load we perform visibility-based view selection for the 83
 84 dense depth map computation. Then we apply a novel extension to a depthmap 84
 85 fusion method to obtain a watertight scene representation from the noisy but 85
 86 redundant depth maps. 86
 87 87

88 2 Previous Work 88

89 Our method is the first system performing dense modeling from Internet photo 89
 90 collections consisting of millions of images. Systems for urban reconstruction 90
 91 from video have been proposed in [15, 6], with [6] achieving real-time dense 3D 91
 92 reconstruction. However, modeling from video is inherently much more efficient 92
 93 as it takes advantage of spatial proximity between the camera positions of suc- 93
 94 cessive frames, whereas the spatial relationships between images in a community 94
 95 photo collection are unknown a priori, and in fact, 40% to 60% of images in such 95
 96 collections turn out to be irrelevant clutter [3]. 96

97 The first approach for organizing unordered image collections was proposed 97
 98 by Schaffalitzky and Zisserman [16]. Sparse 3D reconstruction of landmarks from 98
 99 Internet photo collections was first addressed by the *Photo Tourism* system [17], 99
 100 which achieves high-quality results through exhaustive pairwise image matching 100
 101 and frequent global bundle adjustment. Neither one of these steps is very scal- 101
 102 able, so in practice, the Photo Tourism system can be applied to a few thousand 102
 103 images at most. Aiming at scalability, Snavely et al. [18] construct *skeletal sets* of 103
 104 images whose reconstruction approximates the full reconstruction of the whole 104
 105 dataset. However, computing these sets still requires initial exhaustive pairwise 105

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

110 The speed of our approach is a result of efficient early application of 2D 110
111 appearance-based constraints, similarly to the approach of Li et al. [3]. But 111
112 our system extends [3] to successfully process two orders of magnitude more 112
113 data by parallelizing the computation on graphics processors and multi-core 113
114 architectures. We summarize the dataset and select *iconic images* using 2D image 114
115 appearance as a prerequisite for efficient camera registration. This is the opposite 115
116 of the approach of Simon et al. [20], who treat scene summarization as a by- 116
117 product of 3D reconstruction and select canonical views through clustering the 117
118 3D camera poses. While our method of image organization is initially looser than 118
119 that of [20], it provides a powerful pre-selection mechanism for advancing the 119
120 reconstruction efficiency significantly. 120

121 After selecting the iconic images, the next step of our system is to discover the 121
122 geometric relationships between them and register them together through SfM. 122
123 Li et al. [3] deployed a vocabulary tree [14] to rapidly find related iconics. Our 123
124 system can use a vocabulary tree in the absence of geo-location information for 124
125 the iconics. If this information is available, we use it to help identify possible links 125
126 between iconics. The latter approach is more efficient since it avoids building 126
127 the vocabulary tree. Note that the methods of [21, 22] are also applicable to 127
128 discovering spatial relationships in large collections of data. 128

129 To perform SfM on the set of iconic images, Li et al. [3] partitioned the 129
130 iconic scene graph into multiple connected components and performed SfM on 130
131 each component. In contrast, we do not cut the iconic scene graph, as such 131
132 an approach is prone to excessive fragmentation of scene models. Instead, we 132
133 use a growing strategy combined with efficient merging and periodic bundle 133
134 adjustment to obtain higher-quality, more complete models. Our method is open 134
135 to use techniques for out-of-core bundle-adjustment [23], which take advantage 135
136 of the uneven viewpoint distribution in photo collections. 136

137 Given the registered viewpoints recovered by SfM, we next perform multi- 137
138 view stereo to get dense 3D models. The first approach demonstrating dense 138
139 modeling from photo collections was proposed by Goesele et al. [4]. It uses per- 139
140 pixel view selection and patch growing to obtain a set of surface elements, which 140
141 are then regularized into a Poisson surface model. However, this approach does 141
142 not make it easy to provide textures for the resulting models. Recently, Furukawa 142
143 et al. [24] proposed a dense reconstruction method from large-scale photo collec- 143
144 tions using view clustering to initialize the PMVS approach [25]. This method 144
145 computes a dense model from approximately 13,000 images in about two days on 145
146 a single computer assuming known camera registration. Our proposed method 146
147 uses an extended version of Yang and Pollefeys stereo [7] combined with novel 147
148 multi-layer depth map fusion [26]. While achieving comparable quality, it com- 148
149 putationally outperforms [4, 24] by achieving modeling on a single PC within 149
150 less than an hour instead of multiple days. 150

151 3 The Approach 151

152 In this section we will describe our proposed method. Section 3.1 discusses the 152
 153 initial appearance-based clustering, and Section 3.2 discusses our method for effi- 153
 154 cient geometric verification of the resulting clusters. Section 3.3 explains our SfM 154
 155 scheme, and Section 3.4 explains our stereo fusion that leverages the appearance- 155
 156 based clustering for dense 3D model generation. 156

157 3.1 Appearance-Based Clustering with Small Codes 157

158 Similarly to Li et al. [3], we begin by computing a global appearance descriptor 158
 159 for every image in the dataset. We generate a gist feature [8] for each image by 159
 160 computing oriented edge responses at three scales (with 8, 8 and 4 orientations, 160
 161 respectively), aggregated to a 4×4 spatial resolution. To ensure better grouping 161
 162 of views for our dense reconstruction method, we concatenate the gist with a 162
 163 subsampled RGB image at 4×4 spatial resolution. Both the gist and the color 163
 164 parts of the descriptor are rescaled to have unit norm. The combined descriptor 164
 165 has 368 dimensions, and it is computed on the 8 GPU cores at a rate of 781Hz³. 165

166 The next step is to cluster the gist descriptors to obtain groups of images 166
 167 consistent in appearance. For efficiency in the clustering we aim at a GPU-based 167
 168 implementation, given the inherent parallelism in the distance computation of 168
 169 clustering algorithms like k -means and k -medoids. Since it is impossible to cluster 169
 170 up to 2.8 million 368-dimensional double-precision vectors in the GPU memory 170
 171 of 768 MB, we have chosen to compress the descriptors to much shorter binary 171
 172 strings, such that the Hamming distances between the compressed strings ap- 172
 173 proximate the distances between the original descriptors. To this end, we have 173
 174 implemented on the GPU the locality sensitive binary code (LSBC) scheme of 174
 175 Raginsky and Lazebnik [10], in which the i th bit of the code for a descriptor 175
 176 vector \mathbf{x} is given by $\varphi_i(\mathbf{x}) = \text{sgn}[\cos(\mathbf{x} \cdot \mathbf{r}_i + b_i) + t_i]$, where $\mathbf{r} \sim \text{Normal}(0, \gamma I)$, 176
 177 $b_i \sim \text{Unif}[0, 2\pi]$, and $t_i \sim \text{Unif}[-1, 1]$ are randomly chosen code parameters. 177
 178 As shown in [10], as the number of bits in the code increases, the *normal-* 178
 179 *ized* Hamming distance (i.e., Hamming distance divided by code length) be- 179
 180 tween two binary strings $\varphi(\mathbf{x})$ and $\varphi(\mathbf{y})$ approximates $(1 - K(\mathbf{x}, \mathbf{y}))/2$, where 180
 181 $K(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2 / 2}$ is a Gaussian kernel between \mathbf{x} and \mathbf{y} . We have compared 181
 182 the LSBC scheme with a simple locality sensitive hashing (LSH) scheme for unit 182
 183 norm vectors where the i th bit of the code is given by $\text{sgn}(\mathbf{x} \cdot \mathbf{r}_i)$ [9]. As shown 183
 184 in the recall-precision plots in Figure 2, LSBC does a better job of preserving 184
 185 the distance relationships of our descriptors. 185

186 We have found that $\gamma = 4.0$ works well for our data, and that the code length 186
 187 of 512 offers the best tradeoff between approximation accuracy and memory 187
 188 usage. To give an idea of the memory savings afforded by this scheme, at 32 188
 189 bytes per dimension, each original descriptor takes up 11,776 bytes, while the 189
 190 corresponding binary vector takes up only 64 bytes, thus achieving a compression 190
 191 factor of 184. With this amount of compression, we can cluster up to about 191

³ code in preparation for release

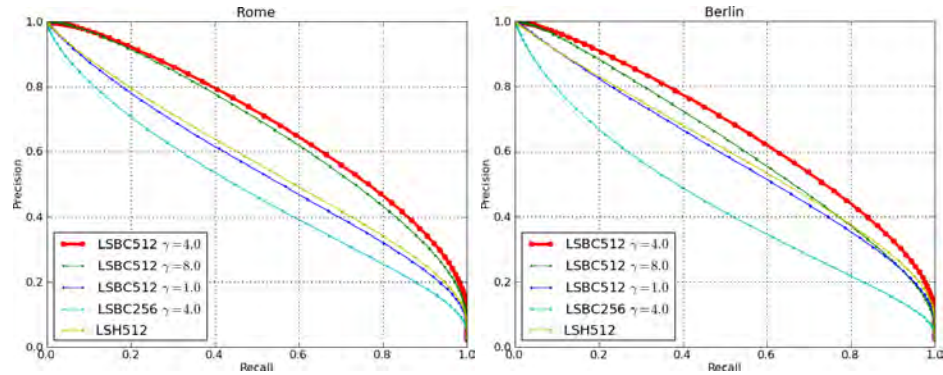


Fig. 2. Comparison of LSH coding scheme [9] and LSBC [10] scheme with different settings for γ and code size on Rome data (left) and Berlin data (right). These plots show the recall and precision of nearest-neighbor search with Hamming distance on binary codes for retrieving the “true” k nearest neighbors 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.

192 4 million images on our memory budget of 768 MB, vs. only a few hundred 192
 193 thousand images in the original GIST representation. An example of a gist cluster 193
 194 is shown in Figure 3. 194

195 For clustering the binary codevectors with the Hamming distance, we have 195
 196 implemented the k -medoids algorithm [12] on the GPU. Like k -means, k -medoids 196
 197 alternates between updating cluster centers and cluster assignments, but unlike 197
 198 k -means, it forces each cluster center to be an element of the dataset. For every 198
 199 iteration, we compute the Hamming distance matrix between the binary codes 199
 200 of all images and those that correspond to the medoids. Due to the size of the 200
 201 dataset and number of cluster centers, this distance matrix must be computed 201
 202 piecewise, as it would require roughly 1050 GB to store on the GPU. 202

203 A generally open problem for clustering in general is how to initialize the 203
 204 cluster centers, as the initialization can have a big effect on the end results. We 204

205 found that images with available geo-location information (typically 10 – 15% 205
 206 of our city-scale datasets) provide a good sampling of the points of interest (see 206
 207 Figure 4). Thus, we first cluster the codevectors of images with available geo- 207
 208 location into k_{geo} clusters initialized randomly. Then we use the resulting centers 208
 209 together with additional k_{rand} random centers to initialize the clustering of the 209
 210 complete dataset (in all our experiments $k_{geo} = k_{rand}$). From Table 2 it can 210
 211 be seen that we gain about 20% more geometrically consistent images by this 211
 212 initialization strategy. 212

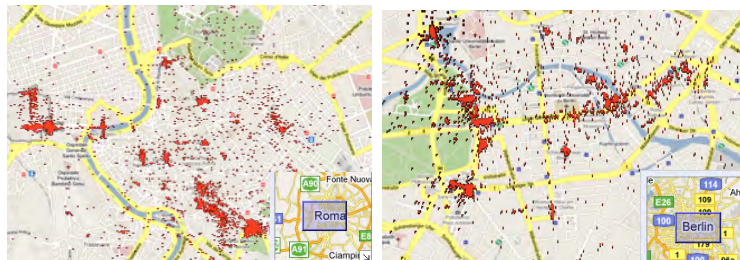


Fig. 4. Geo-tag density map for Rome (left) and Berlin (right).

213 3.2 Geometric Verification 213

214 The clusters obtained in the previous step consist of images that are visually 214
 215 similar but may be geometrically and semantically inconsistent. Since our goal 215
 216 is to reconstruct scenes with stable 3D structure, we next enforce geometric 216
 217 consistency for images within a cluster. A cluster is deemed to be consistent if it 217
 218 has at least n images with a valid pairwise epipolar geometry. This is determined 218
 219 by selecting an initial subset of n images (those closest to the cluster medoid) and 219
 220 estimating the two-view geometry of all the pairs in this subset while requiring 220
 221 at least m inliers (in all our experiments we use $n = 4$, $m = 18$). Inconsistent 221
 222 images within the subset are replaced by others until n valid images are found, 222
 223 or all cluster images are exhausted and the cluster is rejected. 223

224 The computation of two-view epipolar geometry is performed as follows. We 224
 225 extract SIFT features [27] using an efficient GPU implementation,⁴ processing 225
 226 1024×768 images at up to 16.8 Hz on a single GPU. In the interest of com- 226
 227 putational efficiency and memory bandwidth, we limit the number of features 227
 228 extracted to 4000 per image. Next, we calculate the putative SIFT matches for 228
 229 each image pair. This computationally demanding process (which could take a 229
 230 few seconds per pair on the CPU) is cast as a matrix multiplication problem 230
 231 on multiple GPUs (with a speedup of three orders of magnitude to 740 Hz), 231
 232 followed a subsequent distance ratio test [27] to identify likely correspondences. 232

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



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

233 The putative matches are verified by estimation of the fundamental matrix 233
 234 with the 7-point algorithm [28] and ARRSAC [13], which is a robust estimation 234
 235 framework designed for efficient real-time operation. For small inlier ratios, even 235
 236 ARRSAC significantly degrades in performance. However, we have observed that 236
 237 of all registered images in the three datasets a significant fraction had inlier 237
 238 ratios above 50% (e.g., for San Marco, this fraction is 72%). We use this to 238
 239 our advantage by limiting the maximum number of tested hypotheses to 400 in 239
 240 ARRSAC, which corresponds to inlier ratio of approximately 50%. To improve 240
 241 registration performance, we take the best solution deemed promising by the 241
 242 SPRT test of ARRSAC, and perform a *post hoc* refinement procedure. The latter 242
 243 enables us to recover a significant fraction of solutions with less than 50% inlier 243
 244 ratio. Comparing the number of registered images by the standard ARRSAC and 244
 245 the number of images registered by our modified procedure shows a loss of less 245
 246 than 3% for Rome and less than 5% for Berlin of registered images while having 246
 247 an approximately two- to five-fold gain in speed. This result makes intuitive 247
 248 sense: it has been observed [18, 3] that community photo collections contain a 248
 249 tremendous amount of viewpoint overlap and redundancy, which is particularly 249
 250 pronounced at the scale at which we operate. 250

251 We choose a representative or “iconic” image for each verified cluster as the 251
 252 image with the most inliers to the other $n - 1$ top images. Afterwards all other 252
 253 cluster images are only verified with respect to the iconic image. Our system 253
 254 processes all the appearance-based clusters independently using 16 threads on 8 254
 255 CPU cores and 8 GPU cores. In particular, the process of putative matching is 255
 256 distributed over multiple GPUs, while the robust estimation of the fundamental 256
 257 matrix utilizes the CPU cores. This enables effective utilization of all available 257
 258 computing resources and gives a significant speedup to about 480 Hz verification 258
 259 rate an example is shown in Figure 5 259

260 If user provided geo-tags are available (all our city datasets have between 10% 260
 261 and 15% geo-tagged images) we use them to geo-locate the clusters. Our geo- 261
 262 location evaluates the pairwise distances of all geo-tagged image in the iconic 262
 263 cluster. Then it performs a weighted voting on the locations of all images within a 263
 264 spatial proximity of the most central image as defined by the pairwise distances. 264
 265 This typically provides a geo-location for about two thirds of the iconic clusters. 265

3.3 Local Iconic Scene Graph Reconstruction

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

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

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

ARRSAC [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 ARRSAC 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 are accumulated from the depthmaps. Points between the camera center and the depth value receive empty votes, and points beyond the depth value receive a full vote with a weight that falls off with distance. Then a height value is determined that minimizes the number of empty votes above and the number of full

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.

355 votes below. Our extension is to allow the approach to have multiple connected 355
356 “segments” within the column, which provides higher quality mesh models while 356
357 maintaining the regularization properties of the original approach. A polygonal 357
358 mesh is then extracted from the heightmap and texture maps are generated from 358
359 the color images. The heightmap model is highly robust to noise and it can be 359
360 computed very efficiently on the GPU. 360

361 The resolution of the height map is determined by the median camera-to- 361
362 point distance which is representative of the scale of the scene and the accuracy 362
363 of the depth measurements. The texture of the mesh models is then computed as 363
364 the mean of all images observing the geometry. Runtimes are provided in Table 364
365 1. 365

366 4 Conclusions 366

367 This paper demonstrated the first system able to deliver dense geometry for 367
368 Internet scale photo collections with millions of images of an entire city within 368
369 the span of a day on a single PC. Our novel methods extend to the scale of 369
370 millions of images state-of-the-art methods for appearance-based clustering [3], 370
371 robust estimation [32], and stereo fusion [36]. To successfully handle reconstruction 371
372 problems of this magnitude, we have incorporated novel system components 372
373 for clustering of small codes, geo-location of iconic images through their clusters 373
374 clusters, efficient incremental model merging, and enhanced stereo view selection 374
375 through appearance-based clustering. Beyond making algorithmic changes, we 375
376 significantly improve performance by leveraging the constraints from appearance 376
377 clustering and location independence to parallelize the processing on modern 377
378 multi-core CPUs and commodity graphics cards. 378

379 **References**

- 380 1. Snavely, N., Seitz, S.M., Szeliski, R.: Photo tourism: Exploring photo collections 380
381 in 3d. In: SIGGRAPH. (2006) 835–846 381
- 382 2. Agarwal, S., Snavely, N., Simon, I., Seitz, S.M., Szeliski, R.: Building Rome in a 382
383 day. In: ICCV. (2009) 383
- 384 3. Li, X., Wu, C., Zach, C., Lazebnik, S., Frahm, J.M.: Modeling and recognition of 384
385 landmark image collections using iconic scene graphs. In: ECCV. (2008) 385
- 386 4. Goesele, M., Snavely, N., Curless, B., Hoppe, H., Seitz, S.M.: Multi-view stereo 386
387 for community photo collections. In: ICCV. (2007) 387
- 388 5. Furukawa, Y., Curless, B., Seitz, S.M., Szeliski, R.: Towards internet-scale multi- 388
389 view stereo. In: CVPR. (2010) 389
- 390 6. Pollefeys, M., Nister, D., Frahm, J.M., Akbarzadeh, A., Mordohai, P., Clipp, B., 390
391 Engels, C., Gallup, D., Kim, S.J., Merrell, P., Salmi, C., Sinha, S., Talton, B.,
392 Wang, L., Yang, Q., Stewenius, H., Yang, R., Welch, G., Towles, H.: Detailed 392
393 real-time urban 3d reconstruction from video. *International Journal of Computer*
394 *Vision special issue on Modeling Large-Scale 3D Scenes* (2008) 394
- 395 7. Yang, R., Pollefeys, M.: Multi-resolution real-time stereo on commodity graphics 395
396 hardware. In: *Int. Conf. on Computer Vision and Pattern Recognition*. (2003) 396
397 211–217 397
- 398 8. Oliva, A., Torralba, A.: Modeling the shape of the scene: a holistic representation 398
399 of the spatial envelope. *IJCV* **42** (2001) 145–175 399
- 400 9. Andoni, A., Indyk, P.: Near-optimal hashing algorithms for approximate nearest 400
401 neighbor in high dimensions. *Communications of the ACM* **51** (2008) 117–122 401
- 402 10. Raginsky, M., Lazebnik, S.: Locality sensitive binary codes from shift-invariant 402
403 kernels. In: NIPS. (2009) 403
- 404 11. Torralba, A., Fergus, R., Weiss, Y.: Small codes and large databases for recognition. 404
405 In: CVPR. (2008) 405
- 406 12. Kaufman, L., Rousseeuw, P.J.: *Finding Groups in Data: An Introduction to Cluster*
407 *Analysis*. Wiley (1990) 407
- 408 13. Raguram, R., Frahm, J.M., Pollefeys, M.: A comparative analysis of RANSAC 408
409 techniques leading to adaptive real-time random sample consensus. In: ECCV.
410 (2008) 410
- 411 14. Nister, D., Stewenius, H.: Scalable recognition with a vocabulary tree. In: CVPR.
412 (2006) 412
- 413 15. Cornelis, N., Cornelis, K., Van Gool, L.: Fast compact city modeling for navigation
414 pre-visualization. In: *Int. Conf. on Computer Vision and Pattern Recognition*.
415 (2006) 415
- 416 16. Schaffalitzky, F., Zisserman, A.: Multi-view matching for unordered image sets,
417 or "how do i organize my holiday snaps?". In: *ECCV '02: Proceedings of the 7th*
418 *European Conference on Computer Vision-Part I*, London, UK, Springer-Verlag
419 (2002) 414–431 419
- 420 17. Snavely, N., Seitz, S.M., Szeliski, R.: Modeling the world from Internet photo
421 collections. *International Journal of Computer Vision* **80** (2008) 189–210 421
- 422 18. Snavely, N., Seitz, S.M., Szeliski, R.: Skeletal sets for efficient structure from
423 motion. In: CVPR. (2008) 423
- 424 19. Chum, O., Philbin, J., Sivic, J., Isard, M., Zisserman, A.: Total recall: Automatic
425 query expansion with a generative feature model for object retrieval. In: ICCV.
426 (2007) 426

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

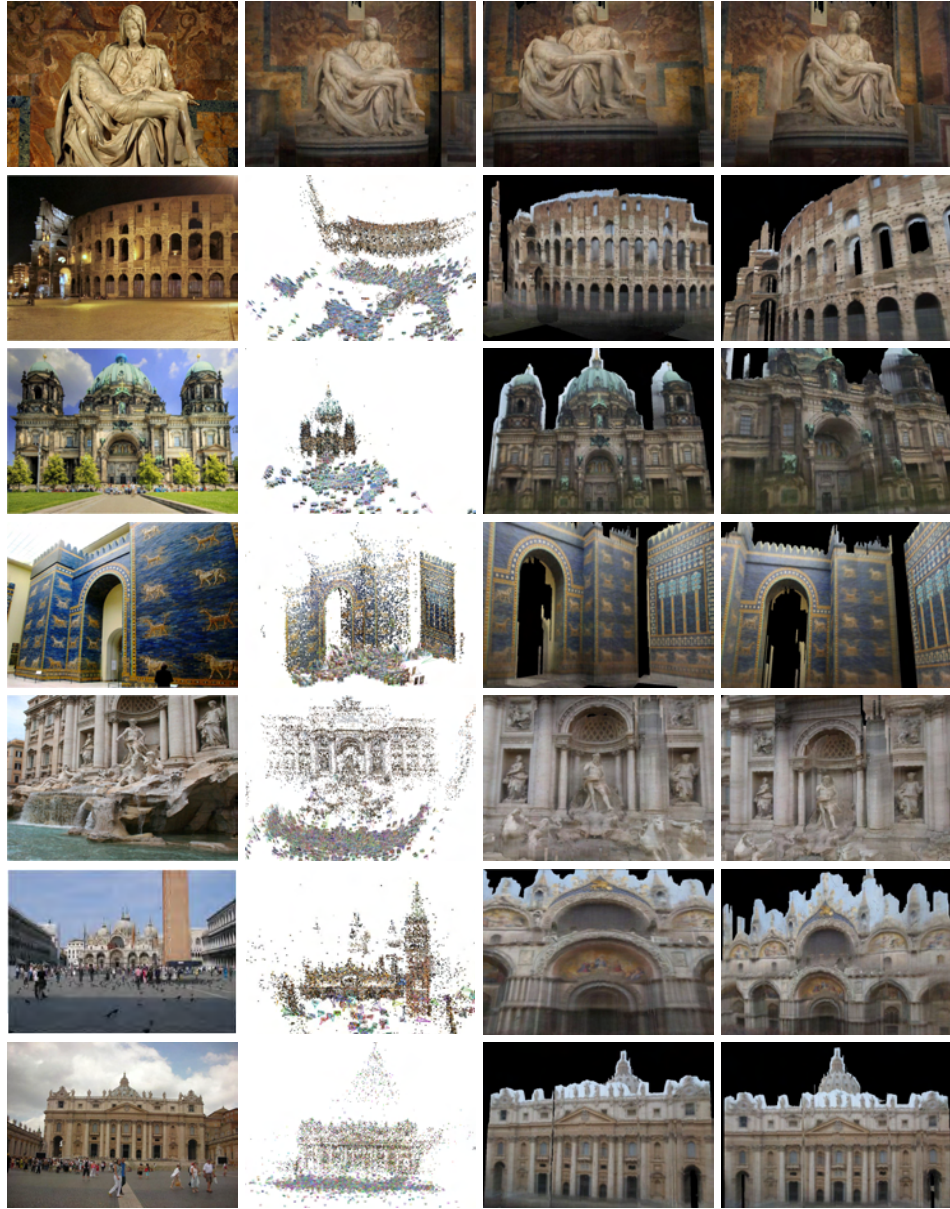


Fig. 6. Original images, local iconic scene graph and 3D model.