



Pose Synchronization under Multiple Pair-wise Relative Poses

Yifan Sun Qixing Huang
The University of Texas at Austin
2317 Speedway, Austin, Texas, 78712

yifansun12@utexas.edu

huangqx@cs.utexas.edu

Multiple relative pose candidates -

Abstract

Pose synchronization, which seeks to estimate cons

absolute poses among a collection of orelative poses estimated between pair

applications. This paper studies an extreme setting where multiple relative pose estimates exist between each object pair, and the majority is incorrect. Popular methods that solve pose synchronization via recovering a low-rank matrix that encodes relative poses in block fail under this extreme setting. We introduce a three-step algorithm for pose synchronization under multiple relative pose inputs. The first step performs diffusion and clustering to compute the candidate poses of the input objects. a theoretical result to justify our diffusion formulation. The second step jointly optimizes the best pose for each object. The final step refines the output of the second step. Experimental results on benchmark datasets of structurefrom-motion and scan-based geometry reconstruction show that our approach offers more accurate absolute poses than state-of-the-art pose synchronization techniques.

1. Introduction

Pose synchronization, which seeks to estimate absolute object poses from noisy relative poses estimated between object pairs, is a fundamental problem in many inverse applications in vision and graphics. Examples include multi-view structure from motion [38], 3D reconstruction from RGB-D scans [21], and reassembling fractured objects [14]. This problem has received great process during the past two decades, starting early greedy approaches [14, 21] to recent optimization-based approaches [2, 4, 11, 15, 17, 19, 20, 26-31, 34, 37]. However, existing approaches assume that there is only one relative pose for each object pair, and most relative poses are inliers. This assumption breaks when relative pose estimation is challenging, e.g., in 3D reconstruction from sparse views. The correct poses may differ from the top-ranked relative poses obtained by a pairwise matching method.

In this paper, we study a new pose synchronization

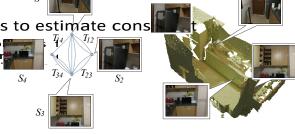


Figure 1. Our approach takes multiple candidate relative poses between pairs of objects as input and outputs absolute poses of the input objects for geometry reconstruction.

Global Poses

setting, where there are multiple relative pose estimates between an object pair, and most of them may be incorrect. This setting is quite popular, e.g., when overlapping ratios are low, or objects possess partial symmetries. approach proceeds in three simple steps. The first step computes for each object a set of candidate poses. This step is based on the fact that a correct relative pose between any object and a root object shall be realized by composing relative poses along multiple paths that connect them. We introduce an iterative procedure that alternates between diffusion and clustering to compute candidate poses. The second step solves a Markov Random Field (MRF) inference problem to jointly selects the best pose for each object so that the induced relative pose agrees with the input relative poses. The third step performs robust optimization to fine-tune the absolute poses of input objects.

Our approach's novelty is a diffusion formulation that synchronizes potentially multiple relative poses between object pairs into candidate poses for each object. The formulation, which utilizes a mixture model, is accurate, robust to noise, and theoretically justified. The resulting candidate poses enable a simple MRF approach via the projected power method [36]. Compared to prior MRF formulations [8,33] that are based on uniform sampling, our approach does not suffer from discretization errors.

We have evaluated our approach on benchmark datasets of multi-view structure-from-motion and geometry reconstruction from depth scans. Experimental results

show that our approach outperforms state-of-the-art pose synchronization approaches.

2. Related Works

Existing approaches fall into three categories. The first category of methods are based on computing a spanning tree among the input graph of relative poses [14, 21]. However, as there are potentially an exponential number of spanning trees in a graph, these approaches are either slow or easily give rise to sub-optimal solutions.

The second category of methods is based on the fact that consistent relative poses can be encoded as a low-rank matrix while the relative poses predicted using pairwise techniques are noisy measurements of the matrix blocks of this low-rank matrix, c.f. [15]. This leads to continuous optimization approaches for synchronization, which recover the underlying low-rank matrix from noisy measurements of its blocks. State-of-the-art techniques include semidefinite programming relaxation [11, 15, 16, 28, 30, 34, 37], spectral techniques [1–3, 17, 20, 23, 26, 27, 29, 31, 32], and non-convex optimization [4, 6, 19, 24, 25, 39]. However, the basis of low-rank matrix recovery is that the spectral gap of the matrix that encodes good input relative poses separates from the matrix that encodes outliers in input relative poses. This assumption is valid mainly when most input relative poses are correct. It does not apply when multiple input relative poses exist between a pair of objects.

The third category of methods solves an inference problem on Markov random field (or MRF) [8, 13, 18]. The key idea is to discretize the rigid pose of each object and then jointly pick the best poses of all objects so that the induced relative poses best align with the input relative poses. The advantage of this approach is that it can incorporate multiple input relative poses between pairs of objects. However, this approach's performance heavily depends on the discretization error when sampling the pose space of each object, which significantly affects the energy landscape of the resulting MRF inference problem. SFM-MRF [8] performs uniform sampling on rotations and translations, which leads to a large-scale MRF problem with many competing solutions that are difficult to solve. In contrast, our approach uses an iterative procedure to generate candidate poses from the input relative poses. This leads to a small-scale MRF problem. Moreover, the resulting MRF is easy to solve, allowing us to adopt an iterative algorithm with guarantees.

Our approach for generating candidate poses alternates between diffusing candidate poses to neighboring objects and then clustering each object's propagated poses from neighboring objects. This paradigm was first proposed in [18] for computing point-based maps among a collection of objects. A recent work [33] extended the idea for pose synchronization among symmetric objects. However, both formulations still use sampling for discretization, i.e., samples on surfaces for point-based maps and partitioning of the pose space into bins for clustering. In contrast, we

introduce a continuous formulation that performs clustering without sampling the pose space. We also provide provable guarantees on the quality of the resulting candidate poses.

Several approaches [4, 8] have studied continuous optimization under robust norms to refine absolute poses from relative poses. A popular paradigm is to employ iteratively reweighted non-linear least squares (IRNLLS) for optimization. Thanks to a highly accurate initial solution for pose refinement, our approach uses the Gauss-Newton method for optimization, which outperforms IRNLLS.

3. Approach

We begin with the problem statement and an overview of the proposed three step approach in Section 3.1. Section 3.2 to Section 3.4 present the technical details of each step.

3.1. Problem Statement and Approach Overview

Problem statement. The input to our approach is given by an observation graph G = (S, E), where S = $\{S_1, \dots, S_n\}$ denotes the input objects, and where E denotes the edges. There may be multiple edges between one object pair, encoding multiple outputs of a pairwise matching algorithm. Let E(i, j) denote the subset of edges from object S_i to S_j . Each edge e $\ensuremath{\mathbb{Z}}$ E(i,j) is associated with a rigid transformation T_e = $\,(R_e,t_e):\Sigma_i\,\rightarrow\Sigma_j\,$ that predicts a relative pose between the local coordinate system Σ_i of object S_i and the local coordinate system Σ_i of object S . We assume that a reverse edge e [2] E (j, i) exists for all e^{J} $\mathbb{E}(i,j)$ and the associated relative pose $T_{e'} = T^{-1}$ is the inverse of T_e. Our goal is to predict the absolute pose $T_i: \Sigma_i \to \Sigma$ of object S_i in a world coordinate system Σ . Approach overview. As shown in Figure 2, our approach proceeds in three steps. The first step computes a candidate set of absolute poses $T_i \supseteq SE(3)$ for each object S_i . Here the absolute poses are expressed with respect to the local coordinate system of the object that has the maximum degree in G. We present a principled approach that alternates between a pose propagation and a pose clustering. The second step joints picks the optimal absolute pose T of each object S₁ by solving a quadratic relaxation of MRF inference. The third step performs continuous optimization to refine the absolute pose T_i of each object S_i .

3.2. Step I: Candidate Pose Generation

Without losing generality, we assume that the object with the maximum degree is S_1 . The underlying principle for computing the candidate pose set T_i is that a candidate pose should be realized by composing relative poses along many paths in G that connects S_1 and S_i . Formally speaking, we first define pose composition along paths.

Definition 1 Consider a path $p = (S_{p_0}S_{p_1}\cdots S_{p_{l(p)}})$ that connects S_{p_0} and $S_{p_{l(p)}}$ where I(p) is the number of edges of p. The composite relative pose $T_p = (R_p, t_p)$ along p is

$$T_p := T_{p_{1(p)-1}p_{1(p)}} \circ \cdots \circ T_{p_0p_1}.$$
 (1)

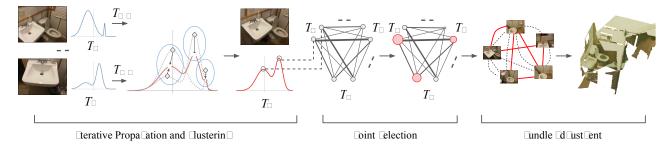


Figure 2. This figure illustrates our three-step approach. Given relative poses predicted between pairs of objects, the first step computes a candidate set of absolute poses for each object. This is done by alternating between pose propagation and pose clustering. The second step jointly selects the best absolute pose for each object. The third step performs bundle adjustment to refine object poses.

Next, we define the induced probability distributions from a rigid pose. Similar to Gaussian mixture models used for mean-shift clustering [7], they are used to aggregate composite relative poses along different paths.

Definition 2 The induced probability distributions over the space of rotations SO(3) and the space of translations R³ from a rigid pose T = (R, t) are given by

$$P^{\sigma_{R},R}(R') = e^{-\frac{\mathbb{E}R'-R\mathbb{E}^{2}_{F}}{2\sigma_{R}^{2}}}, \qquad \mathbb{E}R' \mathbb{E}SO(3) \qquad (2)$$

$$P^{\sigma_{t},t}(t') = e^{-\frac{\mathbb{E}t'-t\mathbb{E}^{2}}{2\sigma_{t}^{2}}}, \qquad \mathbb{E}t' \mathbb{E}R^{3} \qquad (3)$$

$$P^{\sigma_t,t}(t') = e^{-\frac{\alpha t - \tau \Omega^2}{2\sigma_t^2}}, \qquad \exists t' \exists R^3 \qquad (3)$$

where σ_R and σ_t are hyper-parameters of P R and P t, respectively.

Now we are ready to define an induced probability distributions on rotations and translations of a pose graph.

Definition 3 Given a pose graph G, the un-normalized induced probability distributions on rotations and translations of object Si in the local coordinate system of S₁ are defined as

$$Q_{\alpha,I_{\max},I_{0}}^{\{\sigma_{R}^{i}\},R,S_{i}}(R') = \frac{X_{\max}}{I=1} \frac{\alpha^{I}}{(I+I_{0})!} \frac{X}{p^{\mathbb{E}P_{1i}^{(I)}}} P^{\sigma_{R}^{i},R_{p}}(R'), \quad (4)$$

$$Q_{\alpha,l_{\max},l_{0}}^{\{\sigma_{t}^{i}\},t,S_{i}}(t') = \sum_{l=1}^{X^{\max}} \frac{\alpha^{l}}{(l+l_{0})!} \sum_{p \in P_{1i}^{(l)}}^{X^{\max}} P^{\sigma_{t}^{i},t_{p}}(t'), \quad (5)$$

where $P_{1i}^{(1)}$ collects all the paths from S_1 to S_i whose length is 1; I_0 and I_{max} are the minimum and maximum lengths, respectively, of the paths we consider.

There are several hyper-parameters in (4) and (5). α is a decaying parameter that penalizes long paths; σ_R^I and σ_t^I are hyper-parameters of paths of length I. In this paper, we set $\sigma_R^I(\sigma_t^I) = \Gamma_{\sigma_R}^I(\sigma_t)$ where σ_R and σ_t are universal hyper-parameters (to be defined later) shared by all paths. The coefficient \overline{I} is motivated from the fact the variance of composing I independent observations scales as O(1).

In the following, we show that under mild conditions, $\begin{array}{ll}Q_{\alpha,l_{\text{max}},l_{0}}^{\{\sigma_{k}^{l}\},R,S_{i}}(\textbf{R}^{'}) & \text{and} & Q_{\alpha,l_{\text{max}},l_{0}}^{\{\sigma_{t}^{l}\},t,S_{i}}(\textbf{t}^{'}) & \text{converge}\\ Q_{\alpha,\infty,0}^{\{\sigma_{k}^{l}\},R,S_{i}}(\textbf{R}^{'}) & \text{and} & Q_{\alpha,\infty,0}^{\{\sigma_{t}^{l}\},t,S_{i}}(\textbf{t}^{'}), & \text{respectively.} \end{array}$

Proposition 1 When $\alpha > 0$, then $Q_{S_i = 0}^{\{\sigma_R^I\}, R, S_i}$ $Q_{S_{+} \infty}^{\{\sigma_{t}^{l}\},t,S_{i}}$ are well-defined and

$$\mathbb{P}Q_{\alpha,l_{\max},0}^{\{\sigma_{R}^{i}\},R,S_{i}} - Q_{\alpha,\infty,0}^{\{\sigma_{R}^{i}\},R,S_{i}}\mathbb{P}_{1} \leq \frac{\frac{e^{\alpha\lambda_{1}}}{|max|}(\alpha\lambda_{1})^{l_{\max}}n, (6)}{s}$$

$$\mathbb{P}Q_{\alpha,l_{\max},0}^{\{\sigma_{t}^{i}\},t,S_{i}} - Q_{\alpha,\infty,0}^{\{\sigma_{t}^{i}\},t,S_{i}}\mathbb{P}_{1} \leq \frac{e^{\alpha\lambda_{1}}}{|max|}(\alpha\lambda_{1})^{l_{\max}}n, (7)$$

$$\mathbb{P}Q_{\alpha,I_{\max},0}^{\{\sigma_{t}^{i}\},t,S_{i}} - Q_{\alpha,\infty,0}^{\{\sigma_{t}^{i}\},t,S_{i}}\mathbb{P}_{1} \leq \frac{e^{\alpha\lambda_{1}}}{I_{\max}!}(\alpha\lambda_{1})^{I_{\max}}n, (7)$$

where λ_1 is spectral norm of the adjacency matrix of G.

To show the usefulness of the induced probability distributions, we next show that under a standard noise model, the global maximums of $Q_{\alpha,\infty,0}^{\{\sigma_R^l\},R,S_i}$ and $Q_{\alpha,\infty,0}^{\{\sigma_l^l\},t,S_i}$ recover the underlying ground-truth pose of S_i .

Proposition 2 (Informal) Suppose all T_{ij} are independent and identically follow a normal distribution with probability w₁ and a uniform distribution with probability w₂. Let P_1 , P_2 be the PDFs of the two distributions. If $(1-w^{\alpha\lambda^1})P_2(t)$ 2 1, then almost surely the global maximums of $Q^{\{\sigma_{l_R}\},R,S}i$ and $Q^{\{\sigma_l^{\dagger}\},t,S_i}$ are with in the Δ -neighborhood of the ground-truth, where

$$\Delta = O(\lambda_1^{-\frac{1}{6}}(\ln \alpha \lambda_1)^{\frac{2}{3}}). \tag{8}$$

Moreover, this neighborhood is almost convex when we choose σ_R , $\sigma_t = O(\Delta)$.

Approximations of induced probabilities. Computing $Q_{\alpha,l_{max},l_0}^{\{\sigma_R^l\},R,S_1}$ and $Q_{\alpha,l_{max},l_0}^{\{\sigma_t^l\},t,S_1}$ directly is infeasible as the size of $P_{11}^{\{1\}}$ quickly explodes when I increases. In the following, we present an iterative algorithm for approximating them. Specifically, let N_i denote the adjacent vertices of S_i in G. It

is clear that the path set
$$P_{1i}^{(l)}$$
 admits the following recursion:
$$P_{1i}^{(l+1)} = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \mathbb{E}[P_{1i}^{(l)}] = [P_{1i}^{(l)}]$$
(9)

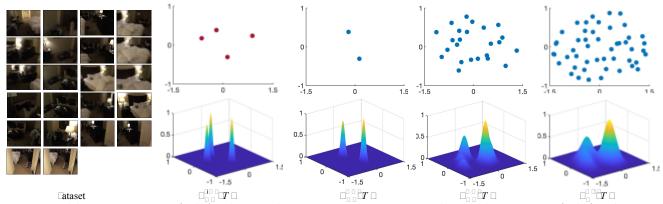


Figure 3. Candidate pose generation. (Col.1) A dataset with 22 scans with two pairwise matches per edge on average. (Col.2) Peak centers and the density function from our results using multi-dimensional scaling. (Col.3-5) Visualizations of $P_{12}^{(4)}(T)$, $P_{12}^{(8)}(T)$, and $P_{12}^{(12)}(T)$.

Applying (9), we obtain the following recursions:

$$Q_{\alpha,l,l'-1}^{\{\sigma_{R}^{i}\},R,S_{i}}(R') = \frac{\alpha}{l'} \sum_{\substack{(1,i) \boxtimes E}} P^{\sigma_{R}^{1},R_{1i}}(R')$$

$$+ \sum_{\substack{j \boxtimes N_{i} \\ j \boxtimes N_{i}}} Q_{\alpha,l-1,l'}^{\sigma_{R}^{i+1}\},R,S_{j}}(R_{ij}R')$$

$$Q_{\alpha,l,l'-1}^{\{\sigma_{t}^{i}\},t,S_{i}}(t') = \frac{\alpha}{l'} \sum_{\substack{(1,i) \boxtimes E}} P^{\sigma_{t}^{1},t_{1i}}(t')$$

$$+ \sum_{\substack{i \boxtimes N_{i} \\ i \boxtimes N_{i}}} Q_{\alpha,l-1,l'}^{\{\sigma_{t}^{i+1}\},t,S_{j}}(R_{ij}R'+t_{ij})$$

$$(10)$$

It remains to control the complexity of $Q_{\alpha, \tilde{l}, \tilde{l}'}^{\{\sigma_{\tilde{l}}^{i}\}}, ^{R,S}i$ and $Q^{\{\sigma_t^i\},t,S_i}_{-\cdots}$ so that the number of mixture components in them do not explode. The key idea is to approximate each of them using a low complexity Gaussian mixture function (GMM) after applying (10) and (11). It is easy to see that both recursions (10) and (11) apply whenever induced probabilities are given by GMMs.

To simplify a GMM, we propose to detect its local maximums and fit a Gaussian function around each local minimum. Consider a GMM on SE(3):

$$Q^{T}(T') = w_{k} e^{-d_{\sigma_{k,R},\sigma_{k,t}}^{2}(T',T_{k})}$$

$$d_{\sigma_{k,R},\sigma_{k,t}}^{2}(T',T_{k}) := \frac{\mathbb{E}R' - R_{k}\mathbb{E}_{F}^{2}}{2\sigma_{\nu_{R},F}^{2}} + \frac{\mathbb{E}t' - k\mathbb{E}^{2}}{2\sigma_{\nu_{L}}^{2}}$$
(12)

where w_k and $(\sigma_{k,R}, \sigma_{k,t})$ are the magnitude and variance parameters of the k-th component.

Proposition 3 T' = (R', t') is a critical point of Q^T if

$$R' = UV^{T}, U\Sigma V^{T} = \frac{X^{K}}{\kappa_{e=1}} \frac{w_{k}}{\sigma_{k,R}^{2}} e^{-d_{\sigma_{k,R},\sigma_{k,t}}^{2}(T',T_{k})}$$

$$t' = \frac{X^{K}}{\kappa_{e=1}} \frac{w_{k}}{\sigma_{k,t}^{2}} e^{-d_{\sigma_{k},R},\sigma_{k}}^{2}(T',T_{k})$$
(13)

where U ΣV^T is the singular value decomposition (SVD). Applying (13), we arrive at the following iterative procedure for computing local maximums of Q^T :

$$R^{n} \leftarrow UV^{T}, U\Sigma V^{T} = \frac{X^{K}}{k=1} \frac{w_{k}}{\sigma_{k,R}^{2}} e^{-d^{2}(T^{c},T_{k})} \cdot R_{k},$$

$$t^{n} \leftarrow \frac{X^{K}}{k=1} \frac{w_{k}}{\sigma_{k,t}^{2}} e^{-d^{2}(T^{c},T_{k})} \cdot t_{k} \frac{X^{K}}{k=1} \frac{w_{k}}{\sigma_{k,t}^{2}} e^{-d^{2}(T^{c},T_{k})},$$
(14)

where $T^n(T^c)$ denote the pose at the next(current) iteration. The starting solutions for (14) are $T_k = (R_k, t_k)$.

After convergence, we detect distinctive local maximums. Decompose $\{1, \dots, K\} = C_1 \ \boxdot \cdots C_n$ where each C_j collects initial poses T_k that converge to the same local maximum T_{i}^{\square} . We then compute the optimal variance parameters and an optimal mixture weight associated with Tp via least square fitting.

Proposition 4 The optimal solution to
$$\min_{\substack{\sigma \\ R, \sigma_t, \\ R}} \frac{Z}{\sigma} w_{K} e^{-\frac{d}{2}\sigma_{R}, \sigma_t} (T^{'}, T^{\boxtimes})$$

$$- X w_{K} e^{-\frac{d^2}{\kappa^{'}, \sigma_{K}^{'}, R^{'}, \sigma_{K}^{'}, t}} (T, T_{K})^{-2}$$

$$(15)$$

is given by

$$w_{j}^{\mathbb{Z}} = \frac{X}{k_{\mathbb{Z}}^{\mathbb{Z}} C_{j}} w_{k} \frac{2\sigma_{k,R}^{2}}{\sigma_{j,R}^{\mathbb{Z}}^{2} + \sigma_{k,R}^{2}}^{2} \frac{2\sigma_{k,t}^{2}}{\sigma_{j,t}^{\mathbb{Z}}^{2} + \sigma_{k,t}^{2}}^{2} \frac{3}{\sigma_{j,t}^{\mathbb{Z}}^{2} + \sigma_{k,t}^{2}}^{2} \frac{3}{\sigma_{k,t}^{2}}^{2} \frac$$

We compute $(\sigma_{j,R}^{\mathbb{B}}, \sigma_{j,t}^{\mathbb{B}})$ in (16) via alternating maximization, starting from $\sigma_{R}' = \text{median}(\sigma_{k,R})$ and $\sigma_{R}' = \text{median}(\sigma_{k,R})$. See the supplementaries for details

 $\sigma_t^{'} = \underset{k \boxtimes C}{\text{median}} (\sigma_{k,t}).$ See the supp. material for details.

After compressing the mixture components of rotations and translations, we keep the top-k (k = 4 in all of our experiments) components with maximum weights. Figure 3 illustrates the diffusion and clustering procedure on a dataset with 22 scans.

Hyper-parameters σ_R and σ_t . We compute them by detecting consistent cycles in the input graph. Specifically, we first use [12] to compute a collection of cycles C. For each cycle c $\[\]$ C, we say it is consistent if the composite rotation R_c and translation t_c satisfy $\[\]$ $R_c - I_3 \[\]$ $\[\]$

3.3. Step II: Optimal Pose Selection

$$\max_{\substack{x_1 \boxtimes [n_1], \cdots, x_n \boxtimes [n_{|S|}] \\ i = 1 \ j = 1 \ e \boxtimes E(i,j)}} \phi_e(T_{i_{X_i}}, T_{j_{X_j}}) \quad (17)$$

where

$$\begin{split} \phi_{e} \big(T_{ix_{i}}, T_{jx_{j}} \big) = & \; exp \; - \frac{ \mathbb{R} R_{e}^{in} R_{ix_{i}} - R_{x_{j}} \mathbb{P}_{F}^{2}}{ 2 \sigma_{R}^{2} \ \, j} \\ & - \frac{ \mathbb{R} R_{e}^{in} t_{ix_{i}} + t_{e}^{in} - t_{jx_{j}} \mathbb{P}^{2}}{ 2 \sigma_{e}^{2}} \end{split}$$

where 2 · 2 F is the matrix Frobenius norm.

Since the candidate poses are computed via diffusion, they are more accurate than those obtained from sampling the space of rotations and translations (e.g., SFM-MRF [8]), which are subject to discretization errors. We find that (17) can be solved effectively using projected power method [36], which is a simple iterative procedure in contrast to more complex MRF_pinference techniques [8,22].

Specifically, denote N = $\bigcap_{j=1}^{n} n_i$. Let A \bigcirc R $\bigcap_{j=1}^{N} n_j$ be a block matrix that stores the values $\varphi_e(T_{ix_i}, T_{jx_j})$, i.e., the elemepts of the (i,j)-th block A_{ij} are given by $A_{ij}(x_i, x_j) = \bigoplus_{e \supseteq E(i,j)} \varphi_e(T_{ik}, T_{jl}), x_i \bigcirc [n_i], x_j \bigcirc [n_j]$. Introduce the corresponding $n \times 1$ block indicator vector $\bigcap_{j=1}^{N} R^{N_j}$ whose element V_{ix_j} , $I_i \bigcirc [|S|], x_i \bigcirc [n_i]$ specifies the confidence on selecting I_{ix_j} as the pose for the i-th object. Initialize $V_{ix_j}^{(0)} = \bigvee_{j=1}^{N} \frac{1}{2} I_j x_j \bigcirc [n_i], i \bigcirc [|S|]$. At each iteration $t \ge 1$, we alternate between a matrix-vector

multiplication operation and projection operation:

$$\hat{v}^{(t)} = Av^{(t-1)}, \quad v_i^{(t)} = P_{\Delta_{n_i}}(\mu_t \hat{v}_i).$$
 (18)

where $P_{\Delta_{n_i}}: R^{n_i} \to \Delta_{n_i}$ is the projection operation on the simplex $\Delta_{n_i}:=\{x\,|\,x\,\boxtimes\,R^{n_i},\mathbf{1}^Tx=1,x\geq0\}$:

$$P_{\Delta_{n_i}}(\hat{v}_i) := \underset{y \boxtimes \Delta_{i_i}}{\operatorname{argmin}} \hat{v}_i - y \mathbb{Z}^2.$$

 $\mu_t \to \infty$ is a scaling parameter. We set $\mu_t = 2A2_{1,\infty}t$ for all of our experiments. (18) is run $t_{max} = 200$ steps. After that, the selected pose is given by

$$T_i^g = T_{ix_i^g}, \quad x_i^g = \arg \max_{x_i} v_{ix_i}^{(t)}.$$

3.4. Step III: Pose Refinement

The third step performs continuous optimization to refine the object poses ($R_{_{i}}$, $t_{_{i}}$). This is achieved by using robust norms to extract inliers in the input pairwise rotations R_{e}^{in} and t_{e}^{in} translations for e in each set E(i,j) that satisfy $R_{e}^{in}\,R_{i}=R_{j}$ and $R_{e}^{in}\,t_{i}+t_{e}^{in}=t_{j}$ approximately:

$$\min_{\substack{R \text{ i} \\ e \supseteq E(i,j)}} \sum_{\substack{e \supseteq E(i,j) \\ t_{i} \\ e \supseteq E(i,j)}} \sum_{\substack{R \text{ in} \\ e \supseteq E(i,$$

where ρ_R and ρ_t are Geman-Mcclure robust functions:

$$\rho_R(x) = \frac{s_R^2 x^2}{s_R^2 + x^2}, \quad \rho_t(x) = \frac{s_t^2 x^2}{s_t^2 + x^2}$$
 (21)

where the hyper-parameters s_R and s_t are given by

$$\begin{split} s_R &= 2 \cdot \underset{1 \leq i \leq N}{\text{median}} \quad \underset{j:E(i,j)=\mathbb{Z}}{\text{min}} \quad \mathbb{R}_e^{\text{in}} - R_j^g R_i^{gT} \mathbb{Z}_F \\ s_t &= 2 \cdot \underset{1 \leq i \leq N}{\text{median}} \quad \underset{j:E(i,j)=\mathbb{Z}}{\text{min}} \quad \mathbb{R}_e^{\text{in}} t_i^g + t_{ij}^{\text{in}} - t_j^g \mathbb{Z}_F \end{split}$$

Here (R_i^g, t_i^g) are the optimal scan poses from step II.

Unlike [4] that solves (19) and (20) using iteratively reweighted non-linear least squares (IRNLLS), we solve both using Gauss-Newton method (GNM). This is motivated from the fact that the output of step II is already close to the optimal solution. Specifically, starting from $R_i^{(1)} = R_i^g$, at iteration k we optimize the velocity vector $c_i \ \mathbb{Z} \ R^3$ for R_i by minimizing

$$F = \frac{X^{S|} X^{S|} X}{\sum_{i=1}^{K} j=1} \frac{X}{e^{\mathbb{Z}E(i,j)}} \frac{S_R^2}{S_R^2 + \mathbb{Z}E_e^{(k)}} \mathbb{Z}^2 E_e^{(k)} - (c_j \times) R_{ij}^{(k)} x + C_{ij}^{(k)} + C_$$

$$R_{ij}^{(k)}(c_i \times) - \frac{E_e^{(k)} \langle E_e^{(k)}, -(c_j \times) R_{ij}^{(k)} + R_{ij}^{(k)}(c_i \times) \rangle}{s_e^2 + \mathbb{P} E_e^{(k)} \mathbb{P}_e^2},$$

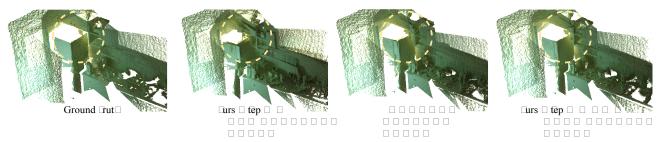


Figure 4. Comparison between reconstruction results before and after GNM for optimizing (19) on a dataset with 18 scans, with mean rotation and translation errors below the visualization.

s.t $c_1=0$, where $R_{ij}^{(k)}=R_j^{(k)}R_i^{(k)^T}$ and $E_e^{(k)}=R_e^{in}-R_i^{(k)}$. F is quadratic in c_i , therefore, the optimal values $c_i^{(k)}$ can be obtained by solving a linear system. After that, we update the rotations as $R_i^{(k+1)}=\exp(\alpha^{(k)}c^{(k)})R_i^{(k)}$, where $\alpha^{(k)}$ is the optimal stepsize obtained through line search.

Similarly, starting from $t_i^{(1)} = t_i^g$, at iteration k we optimize the displacement vector \mathbf{c}_i \mathbf{R}^3 for t_i via

$$\underline{\min}_{c_{i}} X^{|S|} X^{|S|} X \qquad \underline{\sigma_{t}^{2}}_{e^{-|S|}} \mathbb{P}^{(k)} r_{e}^{(k)} + R_{ij}^{in} c_{i} - c_{j}$$

$$- \frac{r_{e}^{(k)} r_{e}^{(k)}}{\sigma_{t}^{2}} \mathbb{P}^{(k)} \mathbb{P}^{2} \qquad s.t. \overline{c}_{1} = 0 \qquad (22)$$

where $r_e^{(k)} = R_e^{in} t_i^{(k)} + t_e^{in} - t_j^{(k)}$. Figure 4 shows the significant improvement in synchronization accuracy with GNM on a dataset with 18 scans. Please refer to the supp. material for more comparisons.

4. Evaluation

Section 4.1 introduces the experimental setup. Section 4.2, Section 4.3, and Section 4.4 present the results on synthetic datasets, RGB-D scans, and RGB images, respectively.

4.1. Experimental Setup

Baseline approaches. We consider four types of baseline approaches.

- Convex programming relaxation. We evaluate against SE-Sync [28], a state-of-the-art approach for pose synchronization. We also compare the performance of rotation synchronization against two other approaches: Shonan-Rot [11] and SDP-Rot [34].
- Spectral synchronization. The second category of baseline employs spectral approaches to perform synchronization. We compare against TransSync [20], a state-of-the-art approach in this category.
- Non-convex optimization. The third category of baselines minimizes non-convex objective functions.
 We compare against SFM-MRF [9] and K-Best [33],

two state-of-the-art approaches, and IRLS-L0 [5], a state-of-the-art rotation synchronization approach.

 Cycle-pruning. The fourth category leverages the cycle-consistency constraint to filter out incorrect relative poses. We compare against RobustRecons [6], a state-of-the-art approach in this category.

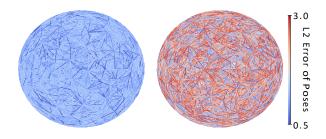


Figure 5. Synthetic datasets where G are geometric graphs on 1000 random samples on the unit sphere. (Left) Sync-Easy. The average vertex degree is 30. (Right) Sync-Hard. The average degree is 20.

Evaluation metrics. Following [9], we report the percentage of absolute angular rotation errors and absolute translation errors that fall within a varying threshold.

4.2. Results on Synthetic Datasets

As illustrated in Figure 5, we generate synthetic datasets by constructing G as geometric graphs. This is motivated from the trial application of pose synchronization in 3D reconstruction from RGB-D or RGB images where each object is matched with adjacent objects. Specifically, we consider a Sync-Easy and a Sync-Hard dataset, where both of them have 1000 objects distributed on a unit sphere. Each dataset is parameterized by five parameters k_n, n_g, p, q, and δ . k_n denotes the average number of adjacent objects in G. For Sync-Easy, we connect each object with 30 nearest neighbors. For Sync-Hard, we connect each object with 10 nearest neighbors. $n_{\rm g}$ denotes the number of absolute poses associated with each object, i.e., (R^i, t^i) , $1 \le i \le n, 1 \le n$ $I \le n_g$. Here the ground-truth is given by $R_i^1 = I_3$, $t_i^1 = 0$. $R_i^{l}, 2 \le l \le n_g$ is a random rotation and $t, 2 \le l \le n_g$ is a random translation in the cube [-1, 1] . p, \dot{q} , and δ control how the relative poses are generated. We generate

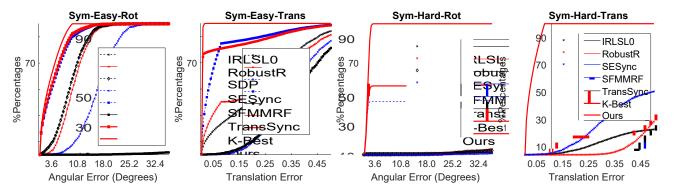


Figure 6. Results on synthetic datasets shown in Figure 5. (a) Angular rotation errors on the Sync-Easy dataset. (b) Translation errors on the Sync-Easy dataset. (c) Angular rotation errors on the Sync-Hard dataset. (d) Translation errors on the Sync-Hard dataset. Baselines: IRLSLO [5]; RobustR [6]; SDP [34]; SESync [28]; SFMMRF [9]; TransSync [20]; K-Best [33].

 n_g relative poses along each edge: $(R_{ij}^I, t_{ij}^I), 1 \le I \le n_g$. With probability p_I ,

$$R_{ij}^{l} = exp(c_{R}^{ijl} \times) R_{j}^{l} R_{i}^{lT}, t_{ij}^{l} = t_{j}^{l} - R_{j}^{l} R_{i}^{lT} t_{i}^{l} + c_{t}^{ijl}$$

where $c_R^{ij\,l}$ and $c_t^{ij\,l}$ are independent random vectors in cube $[-\delta,\delta]^3$. Likewise, with probability $1-p_l$, R_{ij}^l is a random rotation and t_i^l is a random vector in cube $[-1,1]^3$. p_l are chosen so that $p_1=p$ and $p_l=q,2\leq l\leq n_g$. The parameters for Sync-Easy and Sync-Hard are $k=30,n_g=2,p=1,q=0.5,\delta=4\cdot10^{-3}$ and $k_n=20,n_g=3,p=0.8,q=0.5,\delta=2\cdot10-2$.

Figure 6 compares our approach and baseline approaches. First of all, we can see that on both Sync-Easy and Sync-Hard, convex optimization formulations (SE-Sync, Shonan-Rot, and SDP-Rot) and spectral formulations (TransSync) fail to recover the underlying ground-truth. This is expected as both formulations are based on the assumption that there is at most one relative pose estimation between each pair of objects. RobustRecons and IRLS-LO deliver reasonable results on Sync-Easy. We can understand this from the fact that pruning inconsistent cycles applies when the noise-level among inliers is small. while optimization under robust norm may work when the fraction of inliers is significant (50% for Sync-Easy). However, both RobustRecons and IRLS-LO fail on Sync-Hard, in which cycle-consistency is not exactly satisfied among inliers and the fraction of outliers is significant (80% for Sync-Hard).

The top performing baselines are SFM-MRF and K-Best, which can take as input multiple relative poses between a pair of objects. However, the accuracy of our approach is much higher than both. This is because both baselines sample the space of rotations and translations, which is subject to discretization errors. Such discretization errors affect the optimal solution greatly when the fraction of outliers is large. Therefore, our approach outperforms SFM-MRF and K-Best significantly on Sym-Hard.

4.3. Results on RGB-D Scans

Next, we compare our approach and baseline approaches on pose synchronization of RGB-D scans. To this end, we use ScanNet [10], which consists of 1513 3D scenes. Our goal is to test the performance of 3D reconstruction from sparse views. To generate the input data, we randomly sample 10 to 25 scans from each scene. The relative poses between each pair of scans are computed using spectral matching [35]. We consider two settings. In the first setting, we consider one relative pose between each object pair whenever the number of matched features is bigger than 20% times the minimum number of features on each object. In the second setting, we use up to four relative poses per object pair whenever the number of matched features is bigger than 20% times the minimum number of features on each object. We report statistics over the entire ScanNet dataset

Figure 7(Left) shows the rotation and translation errors of our approach and baseline approaches under the first setting. In this setting, the major challenge is that most edges are outliers that do not have a reasonable relative pose estimate. In this case, all methods, including our approach, do not produce accurate synchronization results. However, our approach outperforms all baseline approaches when considering the relative performance with salient margins. Specifically, regarding the percentage of rotation estimates whose errors fall under 10°, our approach achieves 56.4%, which is 17.3% higher than K-Best, the top-performing baseline. Regarding the percentage of translation estimates whose errors fall under 0.25m, our approach achieves 58.9%, which is again 4.6% higher than K-Best.

Figure 7(Right) shows the rotation and translation errors of our approach and baseline approaches under the second setting. Due to representation limitations of convex programming and spectral approaches, i.e., SE-Sync, Shanon-Rot, SDP-Rot, and TransSync, their performance is worse than the first setting. Thanks to certain abilities for handling multiple relative poses per object pair, IRLS-LO, RobustRecons, SFM-MRF, and K-

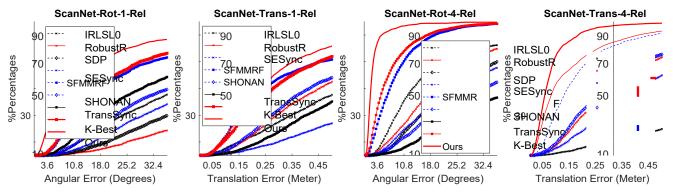


Figure 7. Results on ScanNet [10]. (Left) Average angular rotation errors (a) and average translation errors given one relative pose per object pair. (Right) Average angular rotation errors (c) and average translation errors (d) given four relative poses per object pair.

Best result in performance improvements compared to the first setting. In particular, K-Best is still the top-performing baseline with 67.8% percent of rotation predictions under 10° and 73.2% percent of translation predictions under 0.25m. Our approach is significantly better than all baseline approaches, improving the rotation and translation estimates by 22.1% and 12.5% under the same thresholds. Moreover, the performance gains of our approach in the second setting are larger than in the first setting. In other words, our approach can fully utilize the additional inputs.

4.4. Results on RGB Images

Finally, we evaluate rotation synchronization on RGB images. To this end, we consider two large-scale datasets Cornell-Artquad and San-Francisco from [9]. Cornell-Artquad has 5530 images with 222044 edges. San-Francisco has 7866 images and 101512 edges. Each edge has one relative rotation as input.

Figure 8 plots the angular errors of our approach and baseline approaches. Since G is sparse, approaches based on low-rank matrix recovery that require noise-level smaller than the spectral gap of G, i.e., SE-Sync, Shonan-Rot, SDP-Rot, and Trans-Sync, only exhibit limited performance. For the same reason, the performance of these approaches on San-Francisco is lower than that on Cornell-Artquad. RobustRecons also shows limited performance on both datasets due to difficulty in sampling cycles on sparse graphs. The top-performing baselines are SFM-MRF and K-Best. However, our approach still outperforms both of them. The running time of our method is 263s over San-Francisco and 570s over Cornell-Artsquad, both of which are competitive against baseline methods (see Appendix C for a full comparison of time efficiency).

5. Limitations

One limitation of our approach is that candidate pose generation employs pose clustering at each iteration for each object. When the number of adjacent edges is large, this step is time-consuming. We plan to address this challenge by developing more efficient clustering

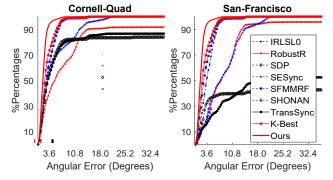


Figure 8. Results of rotation synchronization on (Left) Cornell-Artquad and (Right) San-Francisco.

approaches or using a subset of propagated poses for clustering. Another limitation is that step I which computes candidate poses, and step II, which jointly selects candidate poses, are decoupled. We plan to study how to unify these two steps as both steps employ a diffusion operation that aggregates information from neighboring objects.

6. Conclusion

This paper introduced a new pose synchronization approach that addresses the open challenge where there are multiple relative poses between each pair of objects, in which the correct one may not be the top ranked relative pose. Our approach employs a novel iterative diffusion and clustering step for computing the candidate poses, a joint optimization strategy for selecting consistent poses, and a pose refinement step which optimizes the relative poses under a robust norm. Empirical results on real and synthetic datasets demonstrate the advantage of our approach against state-of-the-art pose synchronization approaches.

Acknowledgement. We would like to acknowledge NSF IIS-2047677, HDR-1934932, and CCF-2019844.

References

- [1] Federica Arrigoni, Andrea Fusiello, and Beatrice Rossi. Camera motion from group synchronization. In 3D Vision (3DV), 2016 Fourth International Conference on, pages 546– 555. IEEE, 2016.
- [2] Federica Arrigoni, Beatrice Rossi, and Andrea Fusiello. Spectral synchronization of multiple views in se (3). SIAM Journal on Imaging Sciences, 9(4):1963–1990, 2016.
- [3] Florian Bernard, Johan Thunberg, Peter Gemmar, Frank Hertel, Andreas Husch, and Jorge Goncalves. A solution for multi-alignment by transformation synchronisation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 2161–2169, 2015.
- [4] Avishek Chatterjee and Venu Madhav Govindu. Efficient and robust large-scale rotation averaging. In ICCV, pages 521– 528. IEEE Computer Society, 2013.
- [5] Avishek Chatterjee and Venu Madhav Govindu. Robust relative rotation averaging. IEEE transactions on pattern analysis and machine intelligence, 40(4):958–972, 2017.
- [6] Sungjoon Choi, Qian-Yi Zhou, and Vladlen Koltun. Robust reconstruction of indoor scenes. In CVPR, pages 5556–5565. IEEE Computer Society, 2015.
- [7] Trial Comaniciu and Peter Meer. Mean shift: A robust approach toward feature space analysis. IEEE Trans. Pattern Anal. Mach. Intell., 24(5):603–619, may 2002.
- [8] David Crandall, Andrew Owens, Noah Snavely, and Daniel Huttenlocher. SfM with MRFs: Discrete-continuous optimization for large-scale structure from motion. IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI), 35(12):2841–2853, December 2013.
- [9] David J. Crandall, Andrew Owens, Noah Snavely, and Daniel P. Huttenlocher. Sfm with mrfs: Discrete-continuous optimization for large-scale structure from motion. IEEE Trans. Pattern Anal. Mach. Intell., 35(12):2841–2853, 2013.
- [10] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes, 2017. cite arxiv:1702.04405.
- [11] Frank Dellaert, David M Rosen, Jing Wu, Robert Mahony, and Luca Carlone. Shonan rotation averaging: Global optimality by surfing so(n). In European Conference on Computer Vision, pages 292–308. Springer, 2020.
- [12] Leonidas J. Guibas, Qixing Huang, and Zhenxiao Liang. A condition number for joint optimization of cycle-consistent networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 1005—1015, 2019.
- [13] Qi-Xing Huang, Yuxin Chen, and Leonidas J. Guibas. Scalable semidefinite relaxation for maximum A posterior estimation. In Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014, pages 64–72, 2014.
- [14] Qixing Huang, Simon Flöry, Natasha Gelfand, Michael Hofer, and Helmut Pottmann. Reassembling fractured objects by geometric matching. In ACM SIGGRAPH 2006 Papers, SIGGRAPH '06, pages 569–578, New York, NY, USA, 2006. ACM.

- [15] Qixing Huang and Leonidas Guibas. Consistent shape maps via semidefinite programming. In Proceedings of the Eleventh Eurographics/ACMSIGGRAPH Symposium on Geometry Processing, pages 177–186, 2013.
- [16] Qixing Huang, Fan Wang, and Leonidas J. Guibas. Functional map networks for analyzing and exploring large shape collections. ACM Trans. Graph., 33(4):36:1–36:11, 2014.
- [17] Qi-Xing Huang, Guo-Xin Zhang, Lin Gao, Shi-Min Hu, Adrian Butscher, and Leonidas J. Guibas. An optimization approach for extracting and encoding consistent maps in a shape collection. ACM Trans. Graph., 31(6):167:1–167:11, 2012.
- [18] Qi-Xing Huang, Guo-Xin Zhang, Lin Gao, Shi-Min Hu, Adrian Butscher, and Leonidas Guibas. An optimization approach for extracting and encoding consistent maps in a shape collection. ACM Trans. Graph., 31(6), nov 2012.
- [19] Xiangru Huang, Zhenxiao Liang, Chandrajit Bajaj, and Qixing Huang. Translation synchronization via truncated least squares. In NIPS, 2017.
- [20] Xiangru Huang, Zhenxiao Liang, Xiaowei Zhou, Yao Xie, Leonidas J Guibas, and Qixing Huang. Learning transformation synchronization. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 8082–8091, 2019.
- [21] Daniel Huber. Automatic Three-dimensional Modeling from Reality. PhD thesis, Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, December 2002.
- [22] Jörg H. Kappes, Björn Andres, Fred A. Hamprecht, Christoph Schnörr, Sebastian Nowozin, Dhruv Batra, Sungwoong Kim, Bernhard X. Kausler, Jan Lellmann, Nikos Komodakis, and Carsten Rother. A comparative study of modern inference techniques for discrete energy minimization problems. In 2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013, pages 1328–1335. IEEE Computer Society, 2013.
- [23] Vladimir G. Kim, Wilmot Li, Niloy J. Mitra, Stephen DiVerdi, and Thomas Funkhouser. Exploring collections of 3d models using fuzzy correspondences. ACM Trans. Graph., 31(4):54:1–54:11, July 2012.
- [24] Spyridon Leonardos, Xiaowei Zhou, and Kostas Daniilidis. Distributed consistent data association via permutation synchronization. In ICRA, pages 2645–2652. IEEE, 2017.
- [25] Shaohan Li, Yunpeng Shi, and Gilad Lerman. Fast, accurate and memory-efficient partial permutation synchronization. In IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pages 15714–15722. IEEE, 2022.
- [26] Deepti Pachauri, Risi Kondor, Gautam Sargur, and Vikas Singh. Permutation diffusion maps (PDM) with application to the image association problem in computer vision. In NIPS, pages 541–549, 2014.
- [27] Deepti Pachauri, Risi Kondor, and Vikas Singh. Solving the multi-way matching problem by permutation synchronization. In NIPS, pages 1860–1868, 2013.
- [28] D.M. Rosen, L. Carlone, A.S. Bandeira, and J.J. Leonard. SE-Sync: A certifiably correct algorithm for synchronization over the special Euclidean group. Intl. J. of Robotics Research, 38(2–3):95–125, Mar. 2019.

- [29] Yanyao Shen, Qixing Huang, Nati Srebro, and Sujay Sanghavi. Normalized spectral map synchronization. In NIPS, pages 4925–4933, 2016.
- [30] Amit Singer. Angular synchronization by eigenvectors and semidefinite programming. Applied and computational harmonic analysis, 30:20–36, 2011.
- [31] Amit Singer and Hau tieng Wu. Vector diffusion maps and the connection laplacian. Communications on Pure and Applied Mathematics, 65:1067–1144, March 2012.
- [32] Yifan Sun, Zhenxiao Liang, Xiangru Huang, and Qixing Huang. Joint map and symmetry synchronization. In Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part V, pages 257–275, 2018.
- [33] Yifan Sun, Jiacheng Zhuo, Arnav Mohan, and Qixing Huang. K-best transformation synchronization. In 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019, pages 10251–10260, 2019.
- [34] Lanhui Wang and Amit Singer. Exact and stable recovery of rotations for robust synchronization. Information and Inference: A Journal of the IMA, 2:145–193, December 2013.
- [35] Zhenpei Yang, Siming Yan, and Qixing Huang. Extreme relative pose network under hybrid representations. In IEEE Conference on Computer Vision and Pattern Recognition, pages 2452–2461, 2020.
- [36] Yuxin Chen and Emmanuel Candes. The projected power method: An efficient algorithm for joint alignment from pairwise differences. Communications on Pure and Applied Mathematics, 71(8), August 2018.
- [37] Yuxin Chen, Leonidas J. Guibas, and Qi-Xing Huang. Near-optimal joint object matching via convex relaxation. In Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014, pages 100–108, 2014.
- [38] Christopher Zach, Manfred Klopschitz, and Manfred Pollefeys. Disambiguating visual relations using loop constraints. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 1426–1433, 2010.
- [39] Xiaowei Zhou, Menglong Zhu, and Kostas Daniilidis. Multi-image matching via fast alternating minimization. In Proceedings of the IEEE International Conference on Computer Vision, pages 4032–4040, 2015.