

Learning to Solve Hard Minimal Problems*

Petr Hruby
ETH Zürich

Department of Computer Science

petr.hruby@inf.ethz.ch

Timothy Duff

University of Washington
Department of Mathematics

timduff@uw.edu

Anton Leykin

Georgia Institute of Technology
School of Mathematics

leykin@math.gatech.edu

Tomas Pajdla

Czech Technical University in Prague
Czech Institute of Informatics, Robotics and Cybernetics

pajdla@cvut.cz

Abstract

We present an approach to solving hard geometric optimization problems in the RANSAC framework. The hard minimal problems arise from relaxing the original geometric optimization problem into a minimal problem with many spurious solutions. Our approach avoids computing large numbers of spurious solutions. We design a learning strategy for selecting a starting problem-solution pair that can be numerically continued to the problem and the solution of interest. We demonstrate our approach by developing a RANSAC solver for the problem of computing the relative pose of three calibrated cameras, via a minimal relaxation using four points in each view. On average, we can solve a single problem in under 70 μ s. We also benchmark and study our engineering choices on the very familiar problem of computing the relative pose of two calibrated cameras, via the minimal case of five points in two views.

1. Introduction

Minimal problems arise from geometrical problems in 3D reconstruction [59, 61, 62], image matching [55], visual odometry, and localization [3, 49, 57, 66]. Many geometrical problems have been successfully formulated and solved as minimal problems [1, 4–6, 11, 12, 18, 27, 33–36, 44, 45, 48, 54, 56, 58, 64, 67]. Technically, minimal problems are systems of polynomial equations which depend on the input data and have a finite number of solutions.

1.1. Motivation

Many geometrical problems are optimization problems that have only one optimal solution. Minimal problems, however, often have many additional spurious solutions. The optimal solution is typically real, satisfies inequality constraints, and fits well all data. Such constraints, however, can not be used by methods of nonlinear algebra [13, 65] which have no ability to bypass finding (or incurring the cost of finding) all solutions of polynomial systems.

RANSAC [23, 53] approximates the optimal solution to a geometrical problem by computing candidate solutions from data samples and picking a solution with maximal data support. This is done by iterating over the samples in an outer loop and over the solutions of a minimal problem for each sample in an inner loop. To find a single solution for a data sample in the inner loop, the state-of-the-art “solve & pick” approach first computes all solutions of a minimal problem and then picks the optimal solutions by removing nonreal solutions, using inequalities, and evaluating the support. Optimization in the inner loop may be very costly when there are many spurious solutions to the minimal problem. Fig. 1 compares the standard “solve & pick” approach with our “pick & solve” approach that learns, for a given data sample, how to first pick a promising starting point and then (ideally) continue it to a meaningful solution.

*This work was partially supported by projects: EU RDF IMPACT No. CZ.02.1.01/0.0/0.0/15 003/0000468, EU H2020 ARTwin No. 856994. The research of AL is partially supported by NSF DMS-2001267. TD acknowledges support from an NSF Mathematical Sciences Postdoctoral Research Fellowship (DMS-2103310.)

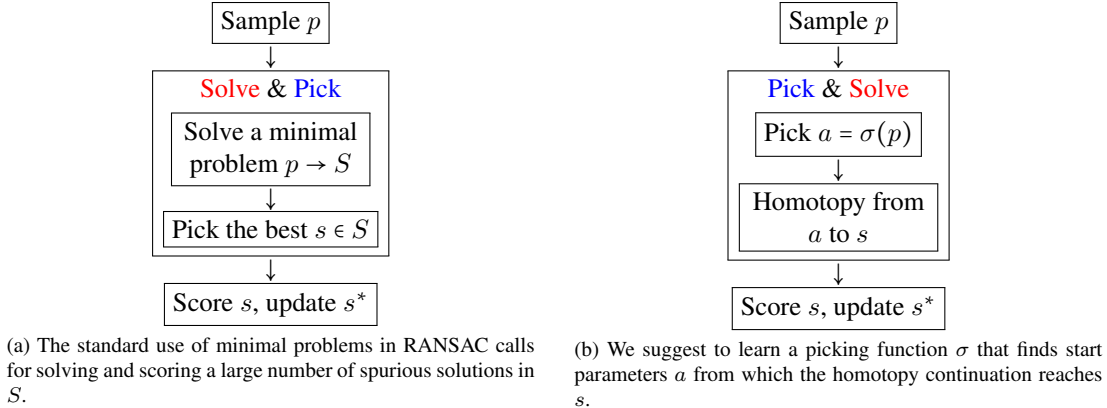


Figure 1. The inner RANSAC loop finds the best solution for a data sample p . This is very expensive when a minimal problem has many spurious solutions. Our efficient homotopy continuation combined with machine learning avoids solving for the spurious solutions. Thus, using minimal problems in RANSAC becomes effectively independent from the number of spurious solutions.

Recent results [15, 16] show that there are many minimal problems in multiview geometry with many spurious solutions which the state-of-the-art polynomial solvers cannot solve efficiently.¹

1.2. Contribution

We present a method for combining optimized homotopy continuation (HC) with machine learning to avoid solving for spurious solutions. The main idea is to learn a single starting point for a real HC path that has a good chance to reach a good solution of the original geometrical problem.

To demonstrate our method on a hard problem, we develop an efficient solver for the “Scranton” minimal problem obtained by relaxing the overconstrained problem of four points in three views (4pt) [50]. We train a model that predicts a starting problem for a single path real HC method to find a good solution. Our solver is implemented efficiently in C++ and evaluated on the state-of-the-art data in computer vision. It successfully solves about 26.3% of inputs in $16.3\mu s$, Tab. 4. In Sec.9 we show that when used in RANSAC, about 4 samples suffice on average to obtain a valid candidate of camera geometry in $61.6\mu s$. No such efficient solver has been known for this problem before. The best-known runtime for a very carefully designed approximation of the problem, reported in [50], was on the order of milliseconds. We thus achieve more than ten times speedup compared to [50]. Most importantly, our approach is general and opens the door to solving other hard minimal problems, e.g., from [15, 16].

We benchmark (Sec. 9) our approach on the classical 5-point problem (5pt) [48] using standard benchmarks [47]. We show that for the 5pt problem, we can solve 29.0% of inputs in about $7.6\mu s$, Tab. 4. Thus, in RANSAC, we can solve it in average in $26.1\mu s$.

Our approach is general. It can be applied even in some cases where the number of spurious solutions is not finite. For instance, our depth formulation of the Scranton problem has an infinite family of solutions where some depths may be zero. Additional polynomial constraints, which do not need to be explicitly enforced, reduce the number of potential solutions to 272—see 15. Thus, by exploiting the “locality” of HC methods, we can guarantee that when starting from a good starting point, we can ignore other spurious solutions with no additional computational cost.

It is important to highlight that, unlike the current symbolic-numeric solvers, our method is coupled with the rest of SfM pipeline, i.e., it uses the real data distribution in a particular vision problem at hand.

1.3. Previous work

The state-of-the-art approach to solving polynomial systems in computer vision is based on symbolic-numeric solvers, which combine elimination (by Gröbner bases [36, 39, 64] or resultants [8, 19, 29]) with eigenvector computation [65] to find all complex solutions. Currently, symbolic-numeric solvers [36, 39–41] provide efficient and stable results for minimal problems with as many as 64 complex solutions [8, 42]. However, these solvers mostly fail to deliver practical results for hard computer vision problems with many solutions. Symbolic/numeric solvers involve two hard computational tasks. First,

¹See Sec. 12 in the SM for more about these problems.

in the symbolic part, large matrices are constructed and triangularized by the Gauss-Jordan elimination. Secondly, in the numeric part, the eigenvectors of $n \times n$ matrices, where n is the number of all complex solutions, are computed. Both steps are prohibitive for generic polynomial systems with many spurious solutions. Methods which deal with real solutions only, e.g. Sturm sequences [48], also generally require expensive manipulations (reduction to a univariate polynomial) and may still need to consider many spurious real solutions.

Global HC methods give an alternative, well-studied approach [7, 10, 14, 68] to finding *all complex solutions* of polynomial systems. Off-the-shelf HC solvers [7, 10, 14, 68] have been proven useful for studying the structure of minimal problems [2, 30, 31, 52]. However, the off-the-shelf solvers are much slower ($10^3 - 10^5$ times) than current symbolic-numeric solvers. For instance, our experiments, Tab. 7, show that solving [48] with complex homotopy continuation in Macaulay2 [43] takes about $10^5 \mu s$ compared to $5 \mu s$ when [38] implementation of [48] is used.

The previous work [20] closest to this work addresses the problem of speeding up minimal HC solvers. This paper took notable steps towards the practical use of homotopy continuation for minimal problem solving. In that work, an off-the-shelf HC implementation [43] has been optimized and efficiently implemented in modern C++. Two hard problems in trifocal geometry have been solved in about $660ms$. Despite not providing practical solvers for RANSAC [23, 53], the results of [20] demonstrated that hard multiview problems involving 312 and 216 solutions can be solved in a stable way and thus could be useful for building practical structure from motion algorithms.

Our paper considers a novel solver for a problem named “Scranton”², which is a minimal relaxation of the overconstrained 4pt problem [50]. Previous work formulated Scranton using camera parameters and found that it has 272 complex solutions [16, 31]. We consider an alternative formulation in terms of 3D point depths, analogous to [52], which has 272 potentially meaningful solutions. The original 4pt problem was solved by numerical search in [50], with about $1ms$ runtime. A depth-formulated 4pt problem was also studied in [52], showing that the overconstrained problem with exact input has a unique solution. Their exact solution does not apply to problems with noisy data.

We also study the classical, well-understood 5-point problem (5pt) of computing the relative pose of two calibrated cameras [48]. Unlike Scranton, this problem has many practical solutions [9, 27, 37, 48, 59]. Currently, the most efficient symbolic-numeric solver [41] of the 5pt problem solves for up to 10 essential matrices in about $5\mu s$. Although the 5pt problem is not as hard as Scranton in terms of the number of spurious solutions, it does provide an important testing ground for us.

2. Our approach

Here we present our approach to solving hard minimal problems. We shall use HC methods to track one real solution of a start problem to obtain one real solution of the target problem. We shall design an algorithm such that this one solution we obtain is a *meaningful* solution with sufficient probability.

2.1. Problem-solution manifold

We operate in the *problem-solution manifold* M of *problem-solution* (p - s) pairs (p, s) , where p is a problem and s is a solution of p . Problem p belongs to a *real* vector space P . Solution s comes from a *real* vector space of solutions. The projection $\pi : M \rightarrow P$ is defined by $(p, s) \mapsto p$. The preimage $\pi^{-1}(p) \in M$ contains all p - s pairs that correspond to a particular problem p .

Example 1. To illustrate the introduced concepts, let us look at one equation $x^3 + ax + b = 0$ in one unknown x with two parameters a and b . Here a problem $p = (a, b)$ has either one or three real solutions depending on whether the discriminant $D = 4a^3 + 27b^2$ is positive or negative; see the corresponding problem-solution manifold M in Fig. 2a.

To be precise, the equation defines an algebraic variety, which is guaranteed to be a smooth manifold when points above the discriminant locus $D = 0$ are removed.

See SM Sec. 13 for the detailed examples of setting up problem-solution manifolds for the 5pt problem and Scranton, a minimal relaxation for the 4pt problem. Below is a condensed version of the 5pt problem.

Example 2. Consider the classical 5pt problem of computing the relative pose of two calibrated cameras which view five world points where the scale is fixed such that the first 3D point lies in the first image plane.

Denote by x the images of 5 points in 2 views: i.e., x is a point in $P = \mathbb{R}^{20}$. Assume the points are in front of both cameras: i.e., their depths $\lambda_{i,j}$, ($i = 1, \dots, 5; j = 1, 2$) are all positive. Our assumptions imply that $\lambda_{1,1} = 1$. The unknown depths vector in the solution space \mathbb{R}^9 determines the relative pose. Therefore, the problem-solution manifold M is contained in the ambient space $\mathbb{R}^{20} \times \mathbb{R}^9$. Equations vanishing on M are given in Sec. 7.1.

²The US telephone area code for Scranton PA is 272.

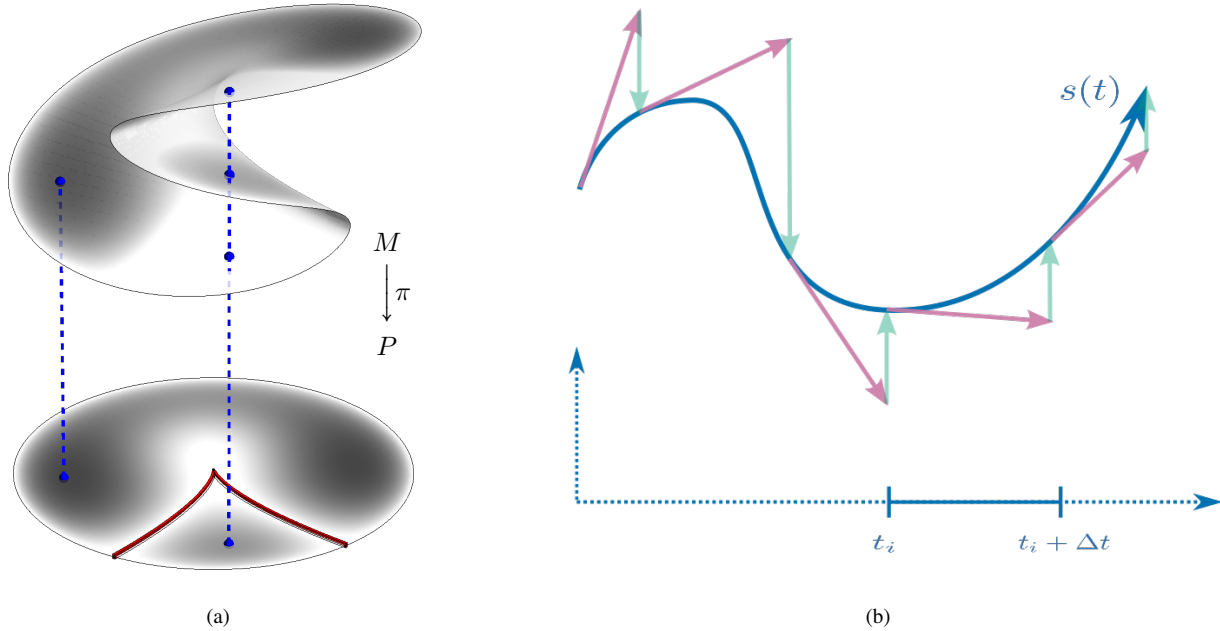


Figure 2. (a) Problem-solution manifold M projected to the problem space P . (b) Numerical HC method.

2.2. Probability distribution on M

We introduce a probability density μ on the problem-solution manifold M that gives the distribution of real-world problem-solution pairs. In Fig. 2a, we give an example of μ depicted with shades of gray (darker is bigger) on the manifold M . Note that the density is such that any problem with a set S of three solutions has only one $s \in S$ with $\mu(s) > 0$. We shall (implicitly) operate under the following assumption:

An input problem p is likely to have one meaningful solution that is dominant, i.e., occurs much more frequently in the real data than other meaningful solutions.

In many problems (e.g., the one in Section 7.2) the number of meaningful solutions is guaranteed to be exactly one generically. The distribution μ is hard to model; in what follows it is represented by training data.

2.3. Pick & solve vs. solve & pick

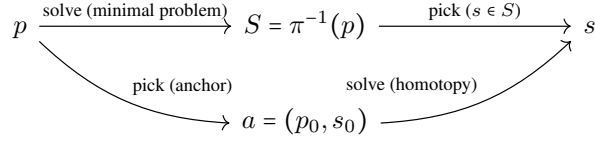
A typical local iterative method (e.g. Newton’s method, gradient descent, etc.) would attempt solving a problem p by obtaining an initial approximate solution s_0 with a hope that it is not too far away from an actual solution and then producing a sequence of its refinements s_0, s_1, s_2, \dots until either a desired quality is reached or some termination-with-failure criterion is satisfied.

Our homotopy approach is a generalization of such local methods. In a nutshell, given a problem p ,

1. we select a suitable start problem-solution pair $(p_0, s_0) \in M$ for p ,
2. we choose a path $p_0 \rightsquigarrow p$ in the problem space P , leading from p_0 to p ,
3. we track the path $(p_0, s_0) \rightsquigarrow (p, s)$ to obtain the target solution s of p .

Selecting a start pair (p_0, s_0) is the key ingredient of our approach. Given a real HC method, one can aim to construct the selection strategy $\sigma(p) = (p_0, s_0)$ in two steps. First, one finds a small set of *anchors* $A \subset M$, such that it is possible to reach (cover) a significant part of M from A by the real homotopy continuation. Secondly, one learns a selection strategy σ such that starting from $(p_0, s_0) \in A$, the meaningful problem-solution pair is reached with sufficiently high frequency to make RANSAC work.

Intuitively, one may perceive a minimal solver employed in RANSAC as an arrow $p \rightarrow S$ in the following diagram:



where an instance p of a minimal problem is “solved” (all solutions in S are found.) Then RANSAC “picks” at most one solution, $S \rightarrow s$, a candidate that maximizes the number of inliers.

Looking for a shortcut that would allow us to go directly $p \rightarrow s$, we reverse this flow: first “picking” an anchor and then tracking *one* HC path to “solve”.

2.4. Structure of our solvers

Our solvers for both 5pt and Scranton minimal problems have a common structure, consisting of an offline training stage and an online evaluation stage. Offline computations may be resource-intensive; however, the online stage must be very efficient to achieve practical sub-millisecond run times. The offline stage consists of:

1. Sampling data D , according to μ , representing the (preprocessed) problem-solution manifold M (Sec 3).
2. Covering a sufficient fraction of the data with anchors $A \subset D$ (Sec. 4).
3. Learning a model σ which selects a starting ps-pair $(p_0, s_0) \in A$ for any given problem p (Sec. 5).

The online stage consists of:

1. Preprocessing the input p to reduce its variability (Sec. 8).
2. Selecting a starting pair from A as $(p_0, s_0) = \sigma(p)$.
3. Constructing polynomial equations of p (Sec. 7).
4. Computing solution s of p by HC from (p_0, s_0) (Sec. 6).
5. Recover solution s of the original problem p (Sec. 8).

The next sections describe these steps in detail.

3. Sampling data representing M and μ

The offline stages of our solvers begin by sampling the data D given by 3D models of various realistic objects. We use models from ETH 3D Dataset to represent μ .

A 3D model consists of 3D points X , cameras C , and relation $I \subset X \times C$ encoding observations ($(X_m, C_i) \in I$ iff C_i observes X_m). For Scranton, we may sample a single p-s pair $(p, s) \in M$ as follows:

1. Select 3 cameras $C_i, C_j, C_k \in C$
2. Select 4 points $X_l, X_m, X_n, X_o \in X : (X_a, C_b) \in I \forall a \in \{l, m, n, o\}, b \in \{i, j, k\}$
3. Project the points to the cameras to get 12 2D points $x_{a,b}$, concatenate them to 24-dim vector $p \in P$
4. Get the depths $\lambda_{a,b}$ of the points in the cameras, concatenate them to 12-dim vector $s \in S$

Sampling for the 5pt problem is similar; in step 1 we select 2 cameras, and in step 2 we select 5 points.

4. Selecting anchors A

We now describe how the starting p-s pairs are obtained. Our goal is to find a small set A of starting p-s pairs (the set of anchors), from which a high portion of p-s pairs from a given distribution can be tracked by HC.

If we limit ourselves to a finite set of p-s pairs and the anchors are selected from the same set, the optimal procedure for the anchor selection consists of building a graph with p-s pairs as vertices. The nodes $(p_i, s_i), (p_j, s_j)$ are connected with an edge, if the correct solution s_j can be obtained by tracking HC from (p_i, s_i) to p_j . A set of anchors covering all problems in this graph is called a *dominating set*. Since computing a minimum-size dominating set is NP-hard, we replace it with a greedy proxy, which is known to perform well³.

³The proposed method is illustrated in SM. Fig. 5

n	5pt problem				Scranton			
	1k	4k	10k	40k	1k	4k	10k	40k
50 %	8	9	8	8	18	18	17	16
75 %	25	28	27	26	47	51	50	50
90 %	59	70	70	70	92	112	120	134
95 %	90	109	115	124	126	168	191	233
100 %	140	235	334	585	176	335	507	1205

Table 1. Number of anchors obtained according to Sec. 4 which are needed to cover 50%, 75%, 90%, 95%, and 100% of n problem-solution pairs. Different values of n are considered. A problem-solution pair is covered by the set of anchors if there is at least one anchor from which the problem-solution pair can be correctly tracked. A track is considered correct if the Euclidean distance from the obtained solution to the ground-truth solution is less than 10^{-5} .

Source	α [%]	Source	α [%]
Courtyard	78.1	Relief 2	77.0
Office	81.1	Off. + Terr.	82.2
Terrains	79.0	Off. + Rel.	80.7
Playground	75.4	O + T + P + R2	79.7

Table 2. Study of sources for anchor selection. Rows correspond to different models or combinations of models, from which the anchors are generated. For each source of anchors, we measure the percentage α of testing p-s problems (generated from models *delivery_area*, *electro*, *facade*, *kicker*, *meadow*, *pipes*) that can be reached from any of the anchors generated from the given source. Anchor sets of 100 anchors are considered for every source.

To show that the selected anchors generalize well to p-s pairs from other scenes, we consider 40000 anchors generated from *office* and *terrains*. For testing data, we generated p-s pairs from the models *delivery_area*, and *facade*.

We have tracked the solutions with HC starting from each of the anchors. The portion of problems correctly tracked from any of the anchors is shown in Tab. 3. The resulting percentage is equivalent to using an oracle that always finds the best anchor to start from.

Next, we show that if the number of vertices in the graph is sufficiently high, a reasonable portion of different p-s pairs from the same scene can be solved by HC starting from one of the anchors A we generate.

We have generated n problem-solution pairs from models *office* and *terrains*⁴. Out of these problem-solution pairs, we have selected anchors which cover 50%, 75%, 90%, 95%, and 100% of data. Tab. 1 shows the number of anchors which cover the data for different values of n for both problems. The number of anchors for 50% and 75% saturates when $n = 10000$. Therefore, we may assume that these anchors will cover a significant portion of problems from the given distribution. Tab. 2 shows a comparison of different sources of anchors. The combination of models *office* and *terrains* gives the best generalizability out of all considered sources.

5. Learning σ to select the starting p-s pair

We formulate the problem of finding the best starting p-s pair as a classification task. Our method relies on a classifier σ , which for a sample problem $p \in P$ assigns a label from $\sigma(p) \in A \cup \{TRASH\}$, where A is the anchor set generated in Sec. 4. The label *TRASH* is included for cases where no problem in A covers p .

Our goal is to minimize the *effective time* $\epsilon_t = \mu_t / \rho$ of the solver, where μ_t is the total time⁵ and ρ is the success rate. Therefore, we must be able to classify p very fast, on the order of $10\mu s$ for a subsequent HC path $\sigma(p) \rightsquigarrow (p, s)$ ($\sigma(p) \neq \{TRASH\}$.) Now, we are going to describe the classifier and its training.

For both problems, we use a Multi-Layer Perceptron (MLP) with 6 hidden layers of 100 neurons with bias.⁶ The input layer has size $\dim P$, and the output layer has size $|A| + 1$. We use the PReLU activation function. During training, we use the dropout before the last layer to prevent overfitting. The classification time of the MLP is about $8\mu s$ for both 5pt and 4pt problems.

⁴<https://www.eth3d.net/datasets>

⁵preprocessing, anchor selection, tracking, and RANSAC scoring

⁶See SM Tab. 11 for a comparison with MLPs with different sizes.

	5 pt problem coverage [%]				
# anchors	8	26	70	124	585
delivery_area	43.3	73.5	86.8	91.4	96.5
facade	51.3	74.9	88.6	92.9	97.0
	4 pt problem coverage [%]				
# anchors	16	50	134	233	1205
delivery_area	47.4	73.6	88.2	92.9	96.6
facade	42.1	71.1	87.9	92.9	96.9

Table 3. Percentage of testing problem-solution pairs which are solvable by the anchors generated from model *Office* and *Terrains*. The anchors are taken from Tab. 1 for $n = 40000$. The problem-solution pair is considered solvable by the anchors, if the correct solution can be obtained by HC starting in any of the anchors. The solution is considered correct if the Euclidean distance from the obtained solution to the ground-truth solution is less than 10^{-5} . This is equivalent to using an oracle classifier that always finds the best anchor to start from.

The input to the MLP is a normalized (Sec. 8) problem $p \in P$. The output is a vector of $|A| + 1$ numbers, which give the score for every starting p-s pair (p_0, s_0) , as well as for *TRASH*. If the score of *TRASH* is higher than the scores of all anchors, we skip the sample. Otherwise, we track from the p-s pair with the highest score.

During training, we normalize the output of the MLP with a softmax layer, and we use a cross-entropy loss, which is a standard method for training classifiers. We use the SGD optimizer, which gives us better results than other optimizers, such as Adam. We generate training and testing data according to Sec. 5.1, and train the MLP by minimizing the loss on the training data for 80 epochs. Then, we select the parameters which maximize the success rate on the validation data. The evaluation of the classifier is shown in Sec. 5.2.

5.1. Training data generation

We use the training and testing data generated from ETH 3D Dataset. Testing data is generated from models *delivery_area* and *facade*, training data from 23 other sequences. First, we generated p-s pairs (p, s) from the models according to Sec. 3. Then, we normalized each problem p (Sec. 8), and tracked the solution to problem p from each anchor $(\check{p}_a, \check{s}_a) \in A$. If the solution to p obtained by HC starting in anchor $(\check{p}_a, \check{s}_a) \in A$ is equal to the expected solution s , then the ID a of the anchor is assigned as the label of problem p . If solution s cannot be reached from any anchor, the label of p is *TRASH*. A problem may have multiple labels. We note that this procedure allows us, in principle, to generate an unlimited amount of training data⁷.

In our experiments, we use about 1 million training p-s pairs per model (23M in total) and 30000 validation p-s pairs per model.

5.2. Classifier evaluation

Now, we are going to show the evaluation of the trained MLPs. During the evaluation, an anchor (p_0, s_0) is selected by the classifier, and HC is tracked from (p_0, s_0) . **Success rate** is the percentage of test p-s pairs (p, s) for which the correct solution s is obtained by HC from (p_0, s_0) . The classification task is difficult because for some problems p , multiple geometrically meaningful solutions (all points in front of cameras, small ratios between the depths, and small baseline) exist. Therefore, we also consider MLP classifiers which return m best anchors. Then, the classification is successful if the correct solution s can be tracked from any of the selected anchors.

To show the benefits of our classifier, we also compare it with the following baselines:

- B1 Start from every anchor in A .
- B2 Start from the closest anchor (p_0, s_0) in terms of Euclidean distance.
- B3 Start from the closest anchor (p_0, s_0) in terms of Mahalanobis distance.

Note that the first baseline gives the upper bound on the success rate for a given anchor set A . The downside of this baseline is that HC paths from all $|A|$ anchors must be tracked. Success rate and total time for different classifiers are shown in Tab. 4. The solution s is considered correct if the squared Euclidean distance from the obtained solution to the ground-truth solution is less than 10^{-5} .

⁷The procedure for generating training data is illustrated in Figure 6.

	5pt problem			4pt problem		
	ρ [%]	μ_t [μs]	ϵ_t [μs]	ρ [%]	μ_t [μs]	ϵ_t [μs]
B1, A_{50}	47.3	∞	∞	44.2	∞	∞
B1, A_{75}	74.2	∞	∞	72.0	∞	∞
B1, A_{90}	87.7	∞	∞	87.9	∞	∞
B2, A_{50}	9.9	11.8	119.5	5.2	16.1	310.5
B2, A_{75}	9.0	12.4	137.8	4.9	16.7	340.9
B2, A_{90}	8.4	12.1	144.5	5.0	16.3	324.5
B2, A	11.2	327.9	2927.7	9.8	150.1	1531.6
B3, A_{50}	14.0	12.2	87.3	5.1	15.9	312.2
B3, A_{75}	13.4	12.8	95.3	4.8	17.0	352.7
B3, A_{90}	4.2	19.5	460.2	4.8	19.9	413.5
MLP, A_{50}	29.3	15.7	53.5	21.6	19.7	91.3
MLP, A_{75}	38.8	15.0	38.7	27.8	20.3	73.0
MLP, A_{90}	39.9	14.3	35.8	29.2	19.6	66.9
MLP _T A_{50}	17.0	4.6	26.9	9.1	8.9	96.8
MLP _T A_{75}	29.0	7.6	26.1	19.0	13.5	71.1
MLP _T A_{90}	36.8	10.8	29.3	26.3	16.2	61.6

Table 4. Classifier evaluation. Rows correspond to start problem selection strategies. The anchors are extracted from datasets *Office* and *Terrains* (Tab. 1). The strategies are evaluated on datasets *Delivery-area* and *Facade*. A_n denotes a set of anchors covering $n\%$ of the training datasets. B1 tracks from all anchors in A_n . The success rate of B1 is equivalent to an “Oracle”, which gives the best possible “retrieval” of the starting problem to reach the target problem, for a given set of anchors. “B2, A_n ” selects the starting problem as the nearest anchor to the target problem (measured by the Euclidean distance in the space of normalized image points) from A_n . “B3, A_n ” selects the starting problem as the nearest anchor to the target problem measured by the Mahalanobis distance. “MLP+T, anchors A_n ” is our method selecting the starting problem as the one from A_n with the highest score given by the MLP constructed in Sec.5. Columns: ρ is the success rate (recall) of retrieving a starting point from which the target problem can be reached, μ_t is the mean solving time, $\epsilon_t = 100 \mu_t / \rho$ is the mean effective solving time, i.e. the average time to obtain one correct solution. If multiple anchors are selected, the classification is considered successful if any of the tracks ends in a correct solution. A solution is considered correct if the squared Euclidean distance from the obtained solution to the ground-truth solution is less than 10^{-5} . We measure total time needed to perform the classification, preprocessing of the input problem and HC from all selected anchors.

6. Homotopy continuation

We now recall the basic principles of HC methods [7, 46, 63] in the framework of our work. Suppose we have a *square system* of n polynomial equations $f(p, s) = (f_1(p, s), \dots, f_n(p, s))$ in n unknowns $s = (s_1, \dots, s_n)$ that vanish on our problem/solution manifold M .

Our task is to numerically continue a known problem/solution pair $(p_0, s_0) \in M$ to a pair $(p, s) \in M$ for some problem of interest $p \in P$. This may be accomplished by introducing a parameter homotopy $H(s, t) = f(p(t), s)$ where $p(t) : [0, 1] \rightarrow P$ is some differentiable function with $p(0) = p_0$ and $p(1) = p$. The goal is to compute a differentiable path $(p(t), s(t)) : [0, 1] \rightarrow M$ such that $s(0) = s_0$ and satisfying the implicit equation $H(p(t), s(t)) = 0$. Note that the homotopy H depends on $p(t)$, and that many choices are possible. We mainly consider **Linear segment HC**: that is, we choose $p(t) = (1-t)p_0 + tp$.

In practice, we compute an approximation of the solution curve $s(t)$ by numerical **predictor/corrector** methods, illustrated in Fig. 2b. In our **predictor** step, the value $s(t_i)^*$ for a given $t_i \in [0, 1)$ is known, and the value $s(t_i + \Delta t)$ for an adaptively-chosen stepsize Δt is approximated using the standard fourth-order Runge-Kutta method. In the **corrector** step, the value $s(t_i + \Delta t)$ is refined by up to 3 steps of Newton’s method.

For both cases of 5pt and Scranton problems, there are additional polynomial constraints which rule out certain spurious solutions. However, one advantage of our HC method is that, although the vanishing set of the f we use is strictly larger than M , these *additional constraints do not need to be explicitly enforced*—See SM Sec. 16

6.1. Efficient HC implementation

Our work builds on the core of an optimized HC solver introduced in [20] that was originally developed for the problem of computing the relative pose of three calibrated cameras from corresponding point-line incidences. The optimized solver

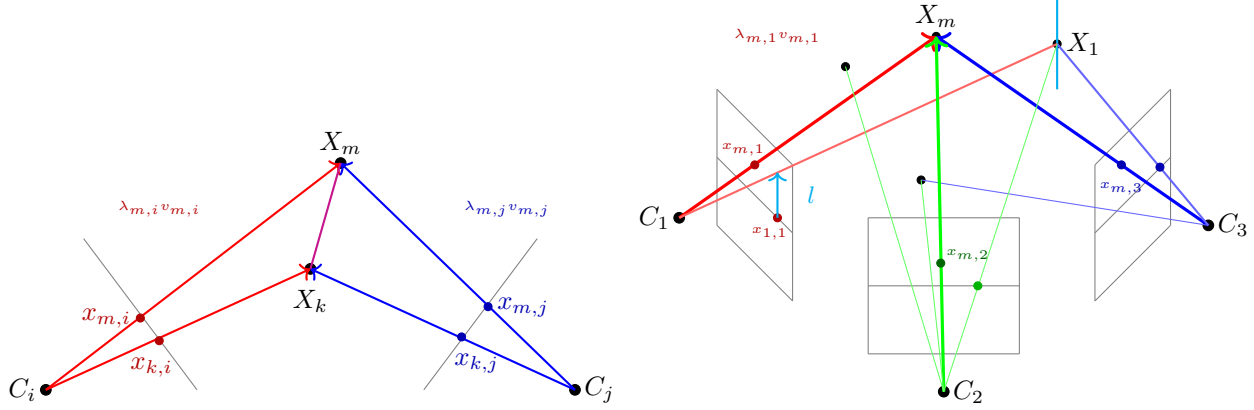


Figure 3. (left) The main geometrical constraint. (right) Scranton formulation. Four points $X_m, m \in \{1, \dots, 4\}$ are projected into three cameras C_1, C_2, C_3 . The first point X_1 projects in the first camera to a point on a vertical line passing through $x_{1,1}$, with l being distance from $x_{1,1}$.

in that work is globally convergent with probability 1, but needs about 500 milliseconds to track 312 complex solution paths to solve a single problem instance.

By contrast, our solver for Scranton tracks a single path in under 10 microseconds, a speedup of more than 1000 \times . As noted in Sec. 1.3, much of this dramatic speedup is because global HC methods must compute all solutions over the complex numbers, whereas our method computes one, allowing now a greater probability of failure for a given data sample. Moreover, the start system in our HC method is tailored to the input by the anchor selection procedure.

There are also significant implementation-specific speedups. For instance, we obtain another $\approx 14\times$ speedup by performing all computations in *real*, instead of complex arithmetic. We also obtain an ≈ 5 speedup by optimizing the linear algebra underlying predictor/corrector steps; the Jacobian matrices of our depth-formulated are sparse, leading to inexpensive closed-form solutions.⁸

7. Minimal problem formulation

We now describe the polynomial systems used by our 5pt and 4pt HC solvers. The unknowns in both systems are the normalized depths in each camera, as formulated in previous works [52, 69]. We choose this formulation over others because (i) due to low degree and sparsity, it leads to fast evaluation of straight-line programs and fast execution of linear algebra subroutines used in homotopy tracking and (ii) it works well in tandem with our normalization procedure described in Sec. 8.

To derive polynomial constraints relating depths and image points, consider two 3D points X_k, X_m labeled by k, m that are projected by two calibrated cameras labeled by i, j into image points $x_{k,i}, x_{m,i}, x_{k,j}, x_{m,j}$ with the corresponding homogeneous coordinates given by $v = [x; 1]$. We see, Fig. 3, that $\|\lambda_{k,i}v_{k,i} - \lambda_{m,i}v_{m,i}\|^2 = \|\lambda_{k,j}v_{k,j} - \lambda_{m,j}v_{m,j}\|^2$ must hold, where the unknown $\lambda_{k,i}$ is the depth of the 3D point X_k in camera i . This constraint means that the distance between every two 3D points, when reconstructed in different cameras, must be the same.

7.1. 5pt minimal problem

The 5pt problem is parametrized by the projection of 5 3D points into 2 calibrated cameras. There are 10 unknown depths $\lambda_{i,j}, i = 1, \dots, 5, j = 1, 2$, and $10 = \binom{5}{2}$ equations

$$\|\lambda_{k,1}v_{k,1} - \lambda_{m,1}v_{m,1}\|^2 = \|\lambda_{k,2}v_{k,2} - \lambda_{m,2}v_{m,2}\|^2 \quad (1)$$

$k, m = 1, \dots, 5, k \neq m$. To dehomogenize this system, we set $\lambda_{1,1} = 1$ to obtain a system of 10 equations in 9 unknowns. It has 80 solutions for generic parameters $v_{i,j}$. There are two isolated singular solutions $\lambda = [1, 0, 0, 0, 0; \pm a, 0, 0, 0, 0]$, with multiplicity 20, and 40 isolated nonsingular solutions. Among the 40 nonsingular solutions, there are at most 10 with all depths positive which extend to a rotation with $\det R_2 = 1$.

⁸See SM Sec. 16 for more details.

For our HC solver, we have to select a square subsystem, i.e., 9 equations for 9 unknowns. For generic parameters, any equation can be dropped to get a square system with 160 solutions, where the two singular solutions have multiplicity 32 and the number of nonsingular solutions rises to 96. As noted in Sec. 6, the dropped equation and other polynomial constraints ($\det R_2 = 1$) need not be explicitly enforced by our HC solver.⁹

7.2. Scranton relaxation of the 4pt problem

The 4pt problem consists of the projection of 4 points into 3 calibrated cameras. Therefore, it involves 12 unknown depths $\lambda_{i,j}$, $i = 1, 2, 3, 4$, $j = 1, 2, 3$, and $18 = \binom{4}{2}\binom{3}{2}$ equations. However, only 2 equations, from each $\binom{3}{2}$ equations involving the same pair of 3D points, are independent. Hence we get 12 equations

$$\|\lambda_{k,i}v_{k,i} - \lambda_{m,i}v_{m,i}\|^2 = \|\lambda_{k,j}v_{k,j} - \lambda_{m,j}v_{m,j}\|^2 \quad (2)$$

$k, m = 1, \dots, 4$, $k \neq m$ and, e.g., $i = 1, 2$, $j = i + 1$. To dehomogenize the system, we set $\lambda_{1,1} = 1$. Unlike for the 5pt problem, here we get an overconstrained system of 12 equations for 11 unknown depths, which has no solution for generic (noisy) parameters $v_{k,i}$. To get a minimal problem, we replace $v_{1,1}$ by $v_{1,1} + l[0; 1; 0]$, where l is a new unknown. This relaxation allows us to “adjust” the second coordinate of the first point in the first camera. Notice that by relaxing the first point in the first view, which has $\lambda_{1,1} = 1$, the equations involving that point

$$\|v_{1,1} + l[0; 1; 0] - \lambda_{m,1}v_{m,1}\|^2 = \|\lambda_{1,2}v_{1,2} - \lambda_{m,2}v_{m,2}\|^2$$

for $m = 2, 3, 4$ remain quadratic in unknowns (λ, l). Solutions to the minimal problem “Scranton”, Fig. 3, are solutions to this square, inhomogeneous system of 12 polynomials used in our HC solver.¹⁰

8. Problem preprocessing

To simplify both the learning the anchor selection strategy σ and the HC tracking, we considered several schemes for *normalizing* the input image correspondences, i.e., the parameters p of the problems. Our chosen normalization yields single representative p for all problems that differ from p up to camera re-orientation or permutation of cameras or correspondences. Once the normalized problem is solved, the original problem may be solved by applying a transformation R_i^{-1} described below.

For the 5pt problem, a problem p is given by image coordinates $x_{i,j} \in \mathbb{R}^2$ with cameras indexed by $i = 1, 2$ and points indexed by $j = 1, \dots, 5$. First, we construct unit 3D vectors representing the rays of the image points as $v_{i,j} = [x_{i,j}; 1] / \|[x_{i,j}; 1]\|$. Next, we compute the mean ray for each camera $m_i = \text{mean}_j(v_{i,j})$. Then, we find the ray $v_{i^*j^*}$ that contains the largest angle with the mean ray m_i of its camera, i.e., $(i^*, j^*) = \text{argmax}_{(i,j)} \angle(x_{i,j}, m_i)$. Next, we compute $w_{i,j} = R_i v_{i,j}$ such that $R_i m_i = [0; 0; 1]$ and $y_{i^*j^*}$, as well as the corresponding y_{1-i^*,j^*} , have the second coordinate equal to 0, i.e., we put them on the “ x axis”. Finally, we swap the cameras to make the camera i^* the first one, project 3D rays $w_{i,j}$ back to the image points $x_{i,j} = w_{i,j} / w_{i,j}^{(3)}$, and reorder the image correspondences counterclockwise starting with j^* ¹¹.

For the 4pt problem, cameras are ordered according to angles of the first point; $\angle([x_{1,1}; 1], [0, 0, 1]) \geq \angle([x_{2,1}; 1], [0, 0, 1]) \geq \angle([x_{3,1}; 1], [0, 0, 1])$.

9. Experiments - RANSAC evaluation

To show how our method generalizes to different real scenes and to data contaminated by noise and wrong matches, we evaluate our approach for the 5pt problem on the dataset from the CVPR 2020 RANSAC Tutorial [47] consisting of 2 validation scenes and 11 test scenes, each comprising 4950 camera pairs. For every camera pair, a set of matched 2D points is known. The points are contaminated with noise and mismatches. We evaluate solvers by plugging them into a RANSAC scheme [53] and computing relative poses for camera pairs in each scene. We evaluate the rotation error and translation error separately. Our evaluation metric is the percentage of relative poses whose angular distance from the ground truth is less than 10° . We believe that this metric is justified, because the main purpose of the RANSAC procedure is to separate the correct matches from the mismatches, and a more precise relative pose can be obtained by local optimization on the inliers.

⁹See SM Sec. 14 for additional details about the problem formulation.

¹⁰See SM Sec. 15 for additional details about Scranton.

¹¹See SM Fig. 8 for an example of image correspondences after the normalization, and SM Sec. 17 for a comparison with other normalizations which delivered worse results in our evaluation.

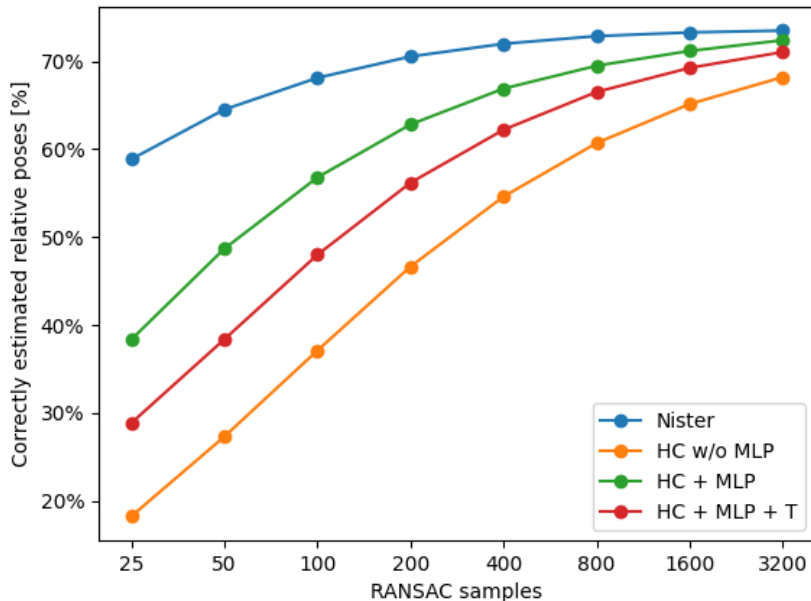


Figure 4. Percentage of camera pairs from the 2020 RANSAC Tutorial [47] for which the relative pose obtained by RANSAC has rotation and translation error less than 10° .

We consider our HC solver with a single anchor, our HC solver with MLP without trash, and our HC solver with MLP and trash. To estimate the success rate of our solver on this data, we compare it with the Nistér 5 point solver [48]. The success rate of the Nistér solver is close to 100%, and its errors are only due to the noise and mismatches in the data. Therefore, if the success rate of a solver on the given data is, e.g., 25%, we expect it to need 4 times the number of samples used for the Nistér solver to get the same results. We have considered RANSAC with 25, 50, 100, 200, 400, 800, 1600, and 3200 samples. The inlier ratio is 3px. The relation between the number of samples and the percentage of correctly estimated cameras is shown in Fig. 4. The graph shows that the lower success rate of our method can be compensated by running RANSAC for more samples. Our method with MLP requires about 4 times more samples than the Nistér solver. Therefore, the success rate on the data from [47] is around 25%, which is about 1.6 times lower than the success rate on the testing data from ETH 3D dataset ¹².

10. Conclusion

Our approach to solving hard minimal problems for RANSAC framework, which uses efficient homotopy continuation and machine learning to avoid solving for many spurious solutions, is fast and delivers correct results. Supplementary Material (SM) presents more details and experiments. Our code and data are available at https://github.com/petrruby97/learning_minimal
Limitations of our approach: First, we sacrifice the high success rate of a complex HC method for a fast, real HC method that fails more frequently. Nevertheless, when combined with trained models, our method succeeds in computing real solutions often enough to be useful in RANSAC. Secondly, our MLP model represents only what it is trained for. Still, we saw that it was able to represent real data distributions while keeping small size and fast evaluation. Fitting to a particular data distribution may also be useful in special situations, e.g., when cameras are mounted on a vehicle, hence having special motions.

References

- [1] Sameer Agarwal, Hon-leung Lee, Bernd Sturmfels, and Rekha R. Thomas. On the existence of epipolar matrices. *International Journal of Computer Vision*, 121(3):403–415, 2017. 1

¹²See SM Sec. 18 for the study of our engineering choices.

- [2] Chris Aholt and Luke Oeding. The ideal of the trifocal variety. *Math. Comput.*, 83(289):2553–2574, 2014. [3](#)
- [3] Hatem Said Alismail, Brett Browning, and M Bernardine Dias. Evaluating pose estimation methods for stereo visual odometry on robots. In *the 11th International Conference on Intelligent Autonomous Systems (IAS-11)*, January 2011. [1](#)
- [4] Daniel Barath. Five-point fundamental matrix estimation for uncalibrated cameras. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 235–243, 2018. [1](#)
- [5] Daniel Barath and Levente Hajder. Efficient recovery of essential matrix from two affine correspondences. *IEEE Trans. Image Processing*, 27(11):5328–5337, 2018. [1](#)
- [6] Daniel Barath, Tekla Toth, and Levente Hajder. A minimal solution for two-view focal-length estimation using two affine correspondences. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2557–2565, 2017. [1](#)
- [7] Daniel J. Bates, Andrew J. Sommese, Jonathan D. Hauenstein, and Charles W. Wampler. *Numerically Solving Polynomial Systems with Bertini*, volume 25 of *Software, environments, tools*. SIAM, 2013. [3](#), [8](#)
- [8] Snehal Bhayani, Zuzana Kukelova, and Janne Heikkilä. A sparse resultant based method for efficient minimal solvers. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 1767–1776. IEEE, 2020. [2](#)
- [9] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000. [3](#)
- [10] Paul Breiding and Sascha Timme. Homotopycontinuation.jl: A package for homotopy continuation in julia. In James H. Davenport, Manuel Kauers, George Labahn, and Josef Urban, editors, *Mathematical Software - ICMS 2018 - 6th International Conference, South Bend, IN, USA, July 24-27, 2018, Proceedings*, volume 10931 of *Lecture Notes in Computer Science*, pages 458–465. Springer, 2018. [3](#)
- [11] Martin Byröd, Klas Josephson, and Kalle Åström. A column-pivoting based strategy for monomial ordering in numerical Gröbner basis calculations. In *European Conference on Computer Vision (ECCV)*, volume 5305, pages 130–143. Springer, 2008. [1](#)
- [12] F. Camposeco, T. Sattler, and M. Pollefeys. Minimal solvers for generalized pose and scale estimation from two rays and one point. In *ECCV – European Conference on Computer Vision*, pages 202–218, 2016. [1](#)
- [13] David Cox, John Little, and Donald O’Shea. *Using Algebraic Geometry*. Springer, 1998. [1](#)
- [14] Timothy Duff, Cvetelina Hill, Anders Jensen, Kisun Lee, Anton Leykin, and Jeff Sommars. Solving polynomial systems via homotopy continuation and monodromy. *IMA Journal of Numerical Analysis*, 2018. [3](#)
- [15] T. Duff, K. Kohn, A. Leykin, and T. Pajdla. PLMP - point-line minimal problems in complete multi-view visibility. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1675–1684, 2019. [2](#), [15](#)
- [16] Timothy Duff, Kathlén Kohn, Anton Leykin, and Tomás Pajdla. PL₁P - point-line minimal problems under partial visibility in three views. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXVI*, volume 12371 of *Lecture Notes in Computer Science*, pages 175–192. Springer, 2020. [2](#), [3](#), [15](#), [17](#)
- [17] Timothy Duff, Viktor Korotynskiy, Tomas Pajdla, and Margaret H Regan. Galois/monodromy groups for decomposing minimal problems in 3d reconstruction. *arXiv preprint arXiv:2105.04460*, 2021. [18](#)
- [18] Ali Elqursh and Ahmed M. Elgammal. Line-based relative pose estimation. In *The 24th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2011, Colorado Springs, CO, USA, 20-25 June 2011*, pages 3049–3056. IEEE Computer Society, 2011. [1](#)
- [19] Ioannis Z. Emiris. A general solver based on sparse resultants. *CoRR*, abs/1201.5810, 2012. [2](#)
- [20] Ricardo Fabbri, Timothy Duff, Hongyi Fan, Margaret H. Regan, David da Costa de Pinho, Elias P. Tsigaridas, Charles W. Wampler, Jonathan D. Hauenstein, Peter J. Giblin, Benjamin B. Kimia, Anton Leykin, and Tomás Pajdla. TRPLP - Trifocal relative pose from lines at points. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 12070–12080. IEEE, 2020. [3](#), [8](#), [21](#)
- [21] R. Fabbri, P. Giblin, and B. Kimia. Camera pose estimation using first-order curve differential geometry. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. [19](#)
- [22] Jean-Charles Faugère, Guillaume Moroz, Fabrice Rouillier, and Mohab Safey El Din. Classification of the perspective-three-point problem, discriminant variety and real solving polynomial systems of inequalities. In J. Rafael Sendra and Laureano González-Vega, editors, *Symbolic and Algebraic Computation, International Symposium, ISSAC 2008, Linz/Hagenberg, Austria, July 20-23, 2008, Proceedings*, pages 79–86. ACM, 2008. [15](#)
- [23] M. A. Fischler and R. C. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, 1981. [1](#), [3](#)
- [24] Phillip Griffiths and Joseph Harris. *Principles of algebraic geometry*. Wiley Classics Library. John Wiley & Sons, Inc., New York, 1994. Reprint of the 1978 original. [18](#)
- [25] J. A. Grunert. Das pothenotische Problem in erweiterter Gestalt nebst über seine Anwendungen in Geodäsie. In *Grunerts Archiv für Mathematik und Physik*, 1841. [15](#)
- [26] Robert M. Haralick, Chung-Nan Lee, Karsten Ottenberg, and Michael Nölle. Review and analysis of solutions of the three point perspective pose estimation problem. *Int. J. Comput. Vis.*, 13(3):331–356, 1994. [15](#)

- [27] R. Hartley and Hongdong Li. An efficient hidden variable approach to minimal-case camera motion estimation. *IEEE PAMI*, 34(12):2303–2314, 2012. [1](#), [3](#)
- [28] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge, 2nd edition, 2003. [15](#), [17](#)
- [29] Janne Heikkilä. Using sparse elimination for solving minimal problems in computer vision. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 76–84. IEEE Computer Society, 2017. [2](#)
- [30] Robert J. Holt and Arun N. Netravali. Uniqueness of solutions to three perspective views of four points. *IEEE Trans. Pattern Anal. Mach. Intell.*, 17(3):303–307, 1995. [3](#)
- [31] Joe Kileel. Minimal problems for the calibrated trifocal variety. *SIAM Journal on Applied Algebra and Geometry*, 1(1):575–598, 2017. [3](#), [17](#)
- [32] L. Kneip, D. Scaramuzza, and R. Siegwart. A novel parametrization of the perspective-three-point problem for a direct computation of absolute camera position and orientation. In *CVPR – IEEE Conference on Computer Vision and Pattern Recognition*, pages 2969–2976, 2011. [15](#)
- [33] L. Kneip, R. Siegwart, and M. Pollefeys. Finding the exact rotation between two images independently of the translation. In *ECCV – European Conference on Computer Vision*, pages 696–709, 2012. [1](#)
- [34] Yubin Kuang and Kalle Åström. Pose estimation with unknown focal length using points, directions and lines. In *IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013*, pages 529–536, 2013. [1](#)
- [35] Yubin Kuang and Kalle Åström. Stratified sensor network self-calibration from TDOA measurements. In *21st European Signal Processing Conference*, 2013. [1](#)
- [36] Zuzana Kukelova, Martin Bujnak, and Tomas Pajdla. Automatic generator of minimal problem solvers. In *European Conference on Computer Vision (ECCV)*, 2008. [1](#), [2](#)
- [37] Zuzana Kukelova, Martin Bujnak, and Tomas Pajdla. Polynomial eigenvalue solutions to minimal problems in computer vision. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2012. [3](#)
- [38] Viktor Larsson. PoseLib - Minimal Solvers for Camera Pose Estimation, 2020. [3](#)
- [39] Viktor Larsson, Kalle Åström, and Magnus Oskarsson. Efficient solvers for minimal problems by syzygy-based reduction. In *Computer Vision and Pattern Recognition (CVPR)*, 2017. [2](#)
- [40] Viktor Larsson, Kalle Åström, and Magnus Oskarsson. Polynomial solvers for saturated ideals. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2307–2316, 2017. [2](#)
- [41] Viktor Larsson, Zuzana Kukelova, and Yinqiang Zheng. Making minimal solvers for absolute pose estimation compact and robust. In *International Conference on Computer Vision (ICCV)*, 2017. [2](#), [3](#), [15](#)
- [42] Viktor Larsson, Magnus Oskarsson, Kalle Åström, Alge Wallis, Zuzana Kukelova, and Tomás Pajdla. Beyond grobner bases: Basis selection for minimal solvers. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 3945–3954, 2018. [2](#)
- [43] Anton Leykin. Numerical algebraic geometry. *Journal of Software for Algebra and Geometry*, 3(1):5–10, 2011. [3](#), [21](#)
- [44] Pedro Miraldo, Tiago Dias, and Srikumar Ramalingam. A minimal closed-form solution for multi-perspective pose estimation using points and lines. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XVI*, pages 490–507, 2018. [1](#)
- [45] Faraz M Mirzaei and Stergios I Roumeliotis. Optimal estimation of vanishing points in a manhattan world. In *International Conference on Computer Vision (ICCV)*, 2011. [1](#)
- [46] A. Morgan. *Solving Polynomial Systems Using Continuation for Engineering and Scientific Problems*. Classics in Applied Mathematics. Society for Industrial and Applied Mathematics (SIAM, 3600 Market Street, Floor 6, Philadelphia, PA 19104), 2009. [8](#)
- [47] Dmytro Myshkin. Benchmarking robust estimation methods. *Tutorial “RANSAC in 2020”, CVPR*, 2020. [2](#), [10](#), [11](#)
- [48] D. Nistér. An efficient solution to the five-point relative pose problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(6):756–770, June 2004. [1](#), [2](#), [3](#), [11](#)
- [49] David Nistér, Oleg Naroditsky, and James Bergen. Visual odometry. In *Computer Vision and Pattern Recognition (CVPR)*, pages 652–659, 2004. [1](#)
- [50] David Nistér and Frederik Schaffalitzky. Four points in two or three calibrated views: Theory and practice. *International Journal of Computer Vision*, 67(2):211–231, 2006. [2](#), [3](#)
- [51] Mikael Persson and Klas Nordberg. Lambda twist: An accurate fast robust perspective three point (P3P) solver. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part IV*, volume 11208 of *Lecture Notes in Computer Science*, pages 334–349. Springer, 2018. [15](#)
- [52] Long Quan, Bill Triggs, and Bernard Mourrain. Some results on minimal euclidean reconstruction from four points. *Journal of Mathematical Imaging and Vision*, 24(3):341–348, 2006. [3](#), [9](#), [23](#)
- [53] R. Raguram, O. Chum, M. Pollefeys, J. Matas, and J.-M. Frahm. USAC: A universal framework for random sample consensus. *IEEE Transactions on Pattern Analysis Machine Intelligence*, 35(8):2022–2038, 2013. [1](#), [3](#), [10](#)
- [54] S. Ramalingam and P. F. Sturm. Minimal solutions for generic imaging models. In *CVPR – IEEE Conference on Computer Vision and Pattern Recognition*, 2008. [1](#)

- [55] Ignacio Rocco, Mircea Cimpoi, Relja Arandjelović, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Neighbourhood consensus networks, 2018. [1](#)
- [56] Johann Salaün, Renaud Marlet, and Pascal Monasse. Robust and accurate line- and/or point-based pose estimation without manhattan assumptions. In *European Conference on Computer Vision (ECCV)*, 2016. [1](#)
- [57] Torsten Sattler, Bastian Leibe, and Leif Kobbelt. Efficient & effective prioritized matching for large-scale image-based localization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(9):1744–1756, 2017. [1](#)
- [58] Olivier Saurer, Marc Pollefeys, and Gim Hee Lee. A minimal solution to the rolling shutter pose estimation problem. In *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*, pages 1328–1334. IEEE, 2015. [1](#)
- [59] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [1](#), [3](#)
- [60] Thomas Schöps, Johannes L. Schönberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. [20](#)
- [61] N. Snavely, S. M. Seitz, and R. Szeliski. Photo tourism: exploring photo collections in 3D. In *ACM SIGGRAPH*, 2006. [1](#)
- [62] Noah Snavely, Steven M Seitz, and Richard Szeliski. Modeling the world from internet photo collections. *International Journal of Computer Vision (IJCV)*, 80(2):189–210, 2008. [1](#)
- [63] Andrew J. Sommese and Charles W. Wampler II. *The numerical solution of systems of polynomials - arising in engineering and science*. World Scientific, 2005. [8](#), [18](#)
- [64] H. Stewenius, C. Engels, and D. Nistér. Recent developments on direct relative orientation. *ISPRS J. of Photogrammetry and Remote Sensing*, 60:284–294, 2006. [1](#), [2](#)
- [65] Bernd Sturmfels. *Solving Systems of Polynomial Equations*, volume 97 of *CBMS Regional Conferences Series*. Amer.Math.Soc., Providence, Rhode Island, 2002. [1](#), [2](#)
- [66] Hajime Taira, Masatoshi Okutomi, Torsten Sattler, Mircea Cimpoi, Marc Pollefeys, Josef Sivic, Tomas Pajdla, and Akihiko Torii. InLoc: Indoor visual localization with dense matching and view synthesis. In *CVPR*, 2018. [1](#)
- [67] Jonathan Ventura, Clemens Arth, and Vincent Lepetit. An efficient minimal solution for multi-camera motion. In *International Conference on Computer Vision (ICCV)*, pages 747–755, 2015. [1](#)
- [68] Jan Verschelde. Polynomial homotopy continuation with phpack. *ACM Commun. Comput. Algebra*, 44(3/4):217–220, 2010. [3](#)
- [69] Ji Zhang, Mireille Boutin, and Daniel G Aliaga. Pose-free structure from motion using depth from motion constraints. *IEEE transactions on image processing*, 20(10):2937–2953, 2011. [9](#)

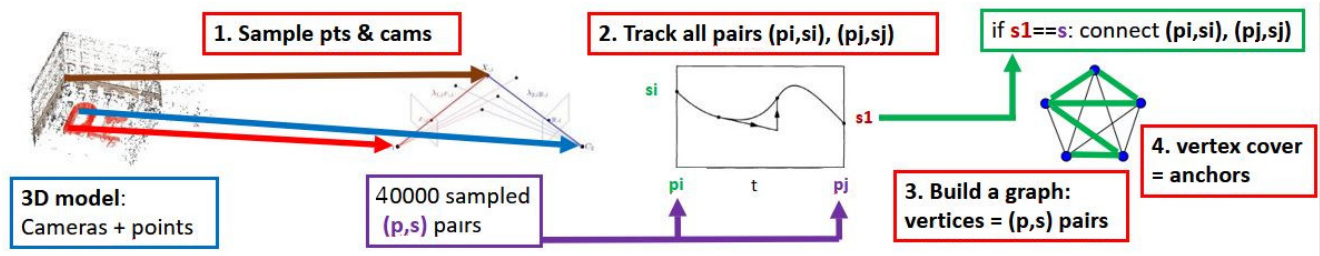


Figure 5. Illustration of generating anchors. A minimal sample of cameras and points is sampled from an existing 3D model. Then, the sampled geometry is converted to the problem-solution pairs. A graph is built whose nodes are the sampled problem-solution pairs. HC is tracked from every p-s pair to every other p-s pair. If the obtained solution is equal to the sampled solution, the p-s pairs are connected with an edge. The selected anchors are a vertex cover of the graph obtained with a greedy algorithm.

Supplementary Material

Here we give additional details, including an analysis of the 5pt and Scranton minimal problem solutions, details of our efficient homotopy continuation implementation, and experiments justifying our engineering choices. Our code is available at https://github.com/petrhruby97/learning_minimal

11. A classical example of a minimal problem

A classical, easy, but still essential, minimal problem in computer vision is computing the pose of a calibrated perspective camera [28] from three points in space and their image projections [25,26,32,41,51]. In one of its classical formulations [25], it leads to a polynomial system of three equations

$$\begin{aligned} \|X_1 - X_2\|^2 &= \|\lambda_1 u_1 - \lambda_2 u_2\|^2 \\ \|X_2 - X_3\|^2 &= \|\lambda_2 u_2 - \lambda_3 u_3\|^2 \\ \|X_3 - X_1\|^2 &= \|\lambda_3 u_3 - \lambda_1 u_1\|^2 \end{aligned}$$

of degree two in three unknown depths $\lambda_1, \lambda_2, \lambda_3$. Parameters of the problem are three 3D points $X_i \in \mathbb{R}^3$ and homogeneous coordinates $u_i \in \mathbb{R}^2$ of three image projections, altogether on $3 \times 3 + 3 \times 2 = 15$ parameters. For generic data, the system has eight complex solutions for λ 's with up to eight real solutions [22]. However, often, there are only zero, two, or four real solutions with positive λ 's.

This example illustrates a typical situation occurring in minimal problem solving. The minimal problem obtained by relaxing a geometrical optimization problem, which has one optimal solution, brings in seven additional (spurious) solutions. In this case, there are always at least two real solutions corresponding to seeing the three points from the opposite sides of the plain they span.

12. Interesting hard minimal problems

Recent results [15, 16] suggest that solving minimal problems with many complex solutions is interesting. A complete classification of minimal problems for points, lines, and their incidences in the calibrated multi-view geometry appeared for the case of complete multi-view visibility [15]. It has been found that there are only 30 minimal problems in that setting, but it also became clear that problems involving more than two cameras are hard for the current symbolic-numeric and homotopy continuation solvers. The number of solutions for three views starts with 64, interesting cases have 200+ solutions, and 5-view cases have 10000+ solutions. Allowing for occlusion or missed detection in images leads to even harder problems. The follow-up work [16] developed a complete classification of minimal problems for generic arrangements of points and lines in space, observed partially by three calibrated perspective cameras when each line is incident to at most one point. It has been found that there is an infinite number of such minimal problems arranged into 74575 equivalence classes when caring only about camera configurations. Interestingly, this classification involves all calibrated trifocal geometry of computer vision for nonincident points and lines in space. Out of 74575 classes, only 759 classes have less than 300 solutions. The rest have (many) more solutions. Thus, for many interesting and potentially practical problems, computing *all* solutions in a reasonable time is a task beyond the reach of current symbolic-numeric and homotopy continuation solving methods.

13. Examples of problem-solution manifolds worked out in detail

Let us now provide a detailed explanation of the problem-solution manifold concept introduced in Sec. 2.1 for the 5pt and 4pt problems solved in this work.

Example 3. Consider the 5pt problem of computing the relative pose of two calibrated cameras from 5 correspondences in two images, i.e., two 3×4 matrices $C_1 = [R_1 \ t_1] = [I \ 0]$, $C_2 = [R_2 \ t_2]$ in the special Euclidean group $\text{SE}_{\mathbb{R}}(3)$, which view five world points $X_1, \dots, X_5 \in \mathbb{R}^3$ where X_1 is normalized to lie in the first image plane: $\{X_1 \mid X_1^{(3)} = 1\} \cong \mathbb{R}^2$. The points are in front of both cameras iff their depths

$$\lambda_{i,j} = \lambda_{i,j}(X, C) = R_j^{(3,\cdot)} X_i + t_j^{(3)}$$

are all positive ($R_j^{(3,\cdot)}$ is the third row of R_j and $t_j^{(3)}$ the third entry of t_j .) Consider

$$\Psi^{5pt} : (\mathbb{R}^2 \times (\mathbb{R}^3)^4) \times \text{SE}_{\mathbb{R}}(3) \rightarrow (\mathbb{R}^2)^{10} \times \mathbb{R}^9$$

$$(X, C) \mapsto (x, \lambda)$$

where

$$x = \left(\lambda_{i,j}^{-1} (R_j X_i + t_j)^{(1:2)} \right)_{i=1, \dots, 5; j=1, 2; (i,j) \neq (1,1)}. \quad (3)$$

Here the problem space is $P = (\mathbb{R}^2)^{10}$ the solution space is $\mathcal{S} = \mathbb{R}^9$, $\pi(x, \lambda) = x$, and our problem-solution manifold $M = M^{5pt}$ is the set of smooth points in the semialgebraic set $\text{im}(\Psi^{5pt}) \cap ((\mathbb{R}_{>0})^9 \times (\mathbb{R}^2)^{10})$.

Remark 1. For a generic problem $x \in M^{5pt}$, the fiber $\pi^{-1}(x)$ consists of at most 10 solutions $\lambda > 0$, and every such λ can be extended uniquely to a pair $(X, C) \mapsto (x, \lambda)$.

Remark 2. Our assumptions in Example 3 imply that $\lambda_{1,1} = 1$. In subsequent sections, we treat the five-point problem as a system of equations in the nine remaining unknown depths. More generally, we could dehomogenize our system by setting any linear form in λ 's equal to 1.

Example 4. Consider the Scranton relaxation of the 4pt problem computing the relative pose of three calibrated cameras from 4 correspondences in 3 images. We have 4 world points X_1, \dots, X_4 with $X_1^{(3)} = 1$ and three cameras $C_1 = [I \ 0]$, $C_2 = [R_2 \ t_2]$, $C_3 = [R_3 \ t_3]$ such that cameras C_1, C_2, C_3 view X_2, \dots, X_4 , cameras C_2, C_3 view X_1 , and camera C_1 views the line ℓ in the direction $e_2 = [0 \ 1 \ 0]^T$ that passes through X_1 , parametrized as $\ell(l) = X_1 + le_2$. Now, we have a map

$$\Psi^{Scr} : (\mathbb{R}^2 \times (\mathbb{R}^3)^3 \times \mathbb{R}) \times (\text{SE}_{\mathbb{R}}(3))^2 \rightarrow (\mathbb{R}^2)^{12} \times \mathbb{R}^{12}$$

$$((X, l), C) \mapsto (x, (\lambda, l))$$

where $x_{i,j}$ are as in Eq. (3) except that

$$x_{1,1} = (X_i + t_j - le_2)^{(1:2)}.$$

Here the problem space is $P = (\mathbb{R}^2)^{12}$ the solution space is $\mathcal{S} = \mathbb{R}^{12}$, $\pi(x, \lambda) = x$, and our problem-solution manifold $M = M^{Scr}$ is the set of smooth points in the semialgebraic set $\text{im}(\Psi^{Scr}) \cap ((\mathbb{R}_{>0})^{11} \times \mathbb{R} \times (\mathbb{R}^2)^{12})$.

14. Additional details for 5pt formulation

Here we provide additional details concerning the solutions of the depth-formulated 5pt problem developed in Section 7.1. Given $(x, \lambda) \in M^{5pt}$, we first note that a valid rotation matrix $R(x, \lambda)$ may be estimated by computing certain auxiliary quantities: for $i \in \{1, \dots, 5\}$ and $v \in \{1, 2\}$, we define $X_{i,v} = \lambda_{i,v} x_{i,v}$, and, for distinct $i, j, k \in \{2, 3, 4, 5\}$,

$$A_{i,j,k}^{(v)} = \begin{pmatrix} X_{i,v} - X_{1,v} & X_{j,v} - X_{1,v} & X_{k,v} - X_{1,v} \end{pmatrix}.$$

Thus, $\det A_{i,j,k}^{(v)}$ gives the oriented volume of a tetrahedron whose vertices are $X_{1,v}, X_{i,v}, X_{j,v}, X_{k,v}$. To estimate the rotation from a geometrically meaningful solution, one may compute either

$$R_2(x, \lambda) = A_{2,3,4}^{(2)} \left(A_{2,3,4}^{(1)} \right)^{-1} \quad (4)$$

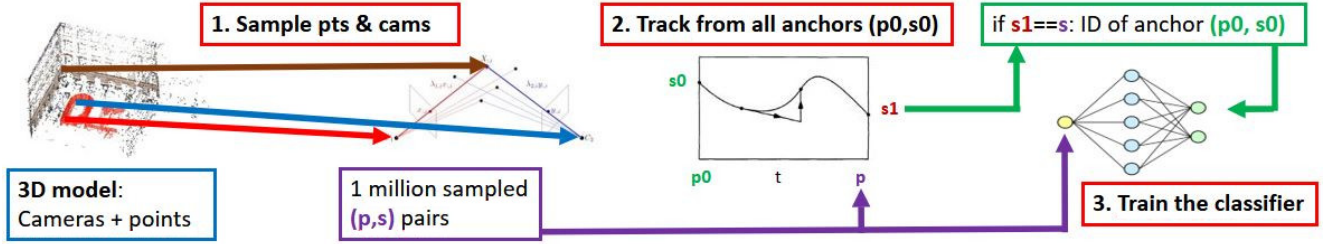


Figure 6. Illustration of generating training data and training the classifier. A minimal sample of cameras and points is sampled from an existing 3D model. Then, the sampled geometry is converted to the problem-solution pairs. Every generated problem-solution pair (p, s) is tracked from every anchor (p_0, s_0) . If the correct solution is obtained, problem p is added to the training data, whereby the ID of the anchor (p_0, s_0) is used as the expected label. Then, the MLP is trained on the generated training data.

or

$$R_2(x, \lambda) = A_{2,3,5}^{(2)} \left(A_{2,3,5}^{(1)} \right)^{-1}. \quad (5)$$

The solutions to Equation (1) need not satisfy the additional constraint $\det R_2(x, \lambda) = 1$, since there is a sign-symmetry $\lambda_{i,v} \mapsto (-1)^{v+1} \lambda_{i,v}$ which changes the sign of $\det R_2(x, \lambda)$ but leaves Equation (1) invariant. Moreover, there are 76 nonsingular, spurious solutions to the square subsystem obtained by dropping one equation, plus an additional 2 of higher multiplicity. For these 78 spurious solutions, either of the matrices in Eq. (4) or Eq. (5) may have determinant -1 . These spurious solutions may be ruled out by enforcing $\det = 1$ for either both of these matrices, or for one of these matrices in addition to all original depth constraints. With these constraints enforced, the generic number of solutions drops to 20. Moreover, there is an additional symmetry given by the “twisted pair” [28]: letting

$$t(x, \lambda) = \lambda_{1,2} v_{1,2} - R(x, \lambda) \lambda_{1,1} v_{1,1},$$

we define $tw(x, \lambda)$ coordinate-wise fixing x and

$$tw(\lambda_{i,j}) = \frac{(-1)^{j+1} \|t(x, \lambda)\|^2 \lambda_{i,j}}{\|\lambda_{i,2} v_{i,2}\|^2 - \|\lambda_{i,1} v_{i,1}\|^2}$$

The map on p-s pairs $(x, \lambda) \mapsto (x, tw(x, \lambda))$ reverses the signs of depths in the second view. This justifies our claim that there are at most 10 geometrically meaningful solutions to Eq. (1). We remark that the partition of 20 solutions into twisted pairs is preserved along non-singular solution curves computed by our HC method.

15. Additional details for Scranton formulation

Unlike the system used for the 5pt problem, the square system for Scranton has infinitely many solutions. Recall that this system is given by the relaxed depth constraint

$$\|v_{1,1} + l[0; 1; 0] - \lambda_{m,1} v_{m,1}\|^2 = \|\lambda_{1,2} v_{1,2} - \lambda_{m,2} v_{m,2}\|^2$$

and Eq. (2) for remaining points and cameras. A one-dimensional family of solutions may be obtained by setting all depths except $\lambda_{1,1}, \lambda_{1,2}, \lambda_{1,3}$ to 0, resulting in 2 nontrivial equations in the remaining 3 unknowns: namely,

$$\begin{aligned} \|v_{1,1} + l[0; 1; 0]\|^2 &= \|\lambda_{1,2} v_{1,2}\|^2 \\ \|\lambda_{1,2} v_{1,2}\|^2 &= \|\lambda_{1,3} v_{1,3}\|^2 \end{aligned}$$

The square system for Scranton also has several families of isolated singular solutions where certain depths equal 0. However, for generic data, the number of *nonsingular* solutions equals the number of solutions with *nonzero depths*, which is 1408. Among these, there is a four-fold sign symmetry where $\lambda_{i,2} \mapsto \pm \lambda_{i,2}, \lambda_{i,3} \mapsto \pm \lambda_{i,3}$, and $320 = 4 \times 80$ cannot be lifted to a valid pair of rotations $(R_2(x, \lambda), R_3(x, \lambda))$.

Taking these facts into account, there are at most $1408 - 3 \times 272 - 4 \times 80 = 272$ geometrically relevant solutions on the problem-solution manifold. This agrees with the number of solutions reported in both [31] and [16], where formulations in

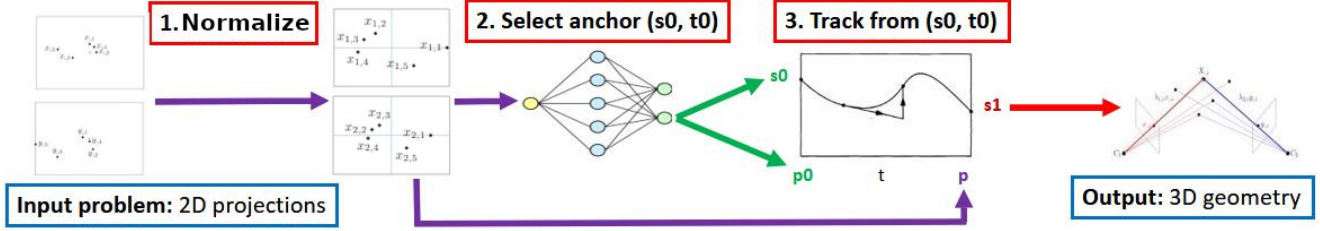


Figure 7. Illustration of testing the classifier. A minimal sample p is obtained in the RANSAC scheme. Then, the sample is normalized. The normalized sample is used as the input to the trained MLP, which selects the starting p-s pair (p_0, s_0) . Then, HC is tracked from (p_0, s_0) to p . If a solution s is obtained, it is converted to a relative pose and the RANSAC score of it is evaluated.

terms of trifocal tensors and camera matrices, respectively, were employed. We note that, unlike the five-point problem, there is no further reduction in the number of solutions implied by a symmetry such as the twisted pair; this follows by numerically computing the Galois group associated to Scranton, using techniques described in [17], which turns out to be the symmetric group on 272 letters.

16. Additional details on HC methods

As noted in Sec. 6, our homotopy H depends on the choice of a path $p(t)$ connecting p_0 to p , where (p_0, s_0) is a known p-s pair and p is the problem to be solved. In all of our experiments, we consider one of two choices. Mostly, we use **1) Linear segment HC**: that is, we choose $p(t) = (1-t)p_0 + tp$. This linear segment homotopy has several advantages; among them, the straight-line programs needed to evaluate H and its derivatives are much simpler than for other paths, and the fact that $p(t)$ is real-valued for all t . However, under Linear segment HC, a differentiable solution curve $s(t)$ satisfying $H(p(t), s(t)) = 0$ need not exist for all $t \in [0, 1]$. For instance, a problem with singular solutions may exist somewhere along the segment connected in P connecting p_0 and p . However, the solution curve $s(t)$ will exist for all $t \in [0, 1]$ if p_0 and p are “close enough”—more precisely, if $p(t)$ avoids the lower-dimensional set of critical values of π in P for all $t \in [0, 1]$.

Alternatively, one may consider **2) Circular arc HC**: here, we reparametrize the segment $p(t)$ via a circular arc $t(\tau) : [0, 1] \rightarrow [0, 1]$ obtained by fixing a random $\gamma \in \mathbb{C}$ (typically of modulus 1) and $t(\tau) = \frac{\gamma\tau}{1+(\gamma-1)\tau}$. This is the γ -trick of [63, Lemma 7.1.3]. Numerical continuation with the resulting homotopy $H(s, \tau)$ is *globally convergent with probability one*: for almost all choices of γ , the solution curves $s(t)$ are defined for all $t \in [0, 1]$, and any isolated target solution (p, s) is the endpoint of some solution curve. However, this necessitates computing many spurious solutions. An experimental comparison of Linear segment and Circular arc HC may be found in Tab. 7.

We now explain why the square systems used in our homotopies are sufficient when starting from a p-s pair on the problem-solution manifold. Consider a path $t \mapsto (p(t), s(t)) \in P \times \mathcal{S}$ where $(p(0), s(0)) = (p_0, s_0) \in M$ and such that that the $n \times n$ Jacobian matrix $\frac{dH}{ds}(p(t), s(t))$ has rank n for all $t \in [0, 1]$. Thus, the path $t \mapsto (p(t), s(t))$ is contained in a single connected component of the set of nonsingular points in the complex vanishing set $\{(p, s) \in P_{\mathbb{C}} \times \mathcal{S}_{\mathbb{C}} \mid f(p, s) = 0\}$. Among these connected components is the set of smooth points in the Zariski closure of M . Indeed, the complex Zariski closure of M has a rational parametrization (given by one of the maps Ψ^{5pt}, Ψ^{Scr} defined in Sec. 13), so it is irreducible, and the connectedness of its smooth points follows by [24, pp. 21–22]. Since the point (p_0, s_0) is contained in this connected component, so also must $(p(t), s(t))$ for all $t \in [0, 1]$. Thus, any polynomial $g(p, s)$ vanishing on M satisfies $g(p(t), s(t)) = 0$ for all $t \in [0, 1]$. This means that, if HC tracking from $(p_0, s_0) \in M$ succeeds using our square system of constraints, then *all* additional constraints which are polynomial equalities are automatically satisfied. This justifies the fact that we do not explicitly enforce constraints like $\det R_2(x, \lambda) = 1$, since it is enough to enforce them for the initial p-s pair.

16.1. Efficient evaluation of predictor/corrector

The Runge-Kutta method used for the predictor step and Newton’s method used for the corrector step in our HC implementation both require solving systems of linear equations. In either step, the coefficient matrix is given by the Jacobian $\frac{\partial H(s, t)}{\partial s}$ (Sec. 6). In the case of the depth formulation of the Five-Point problem (1) and the Four-Point problem (2), the associated Jacobian matrix is sparse. The sparsity pattern of the Jacobian matrix $\frac{\partial H(s, t)}{\partial s}$ is shown in (6) for the Five-Point problem, and in (7) for the Four-Point problem.

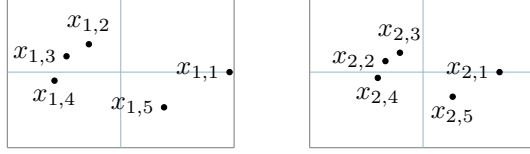


Figure 8. An example of a 5pt problem after normalization of image coordinates.

$$\begin{bmatrix} A_{0,3} & 0 & 0 & 0 & A_{0,4} & A_{0,5} & 0 & 0 & 0 \\ 0 & A_{1,1} & 0 & 0 & A_{1,4} & 0 & A_{1,6} & 0 & 0 \\ A_{2,0} & A_{2,1} & 0 & A_{2,3} & 0 & 0 & A_{2,6} & 0 & 0 \\ 0 & 0 & A_{3,2} & 0 & A_{3,4} & 0 & 0 & A_{3,7} & 0 \\ A_{4,0} & 0 & A_{4,2} & 0 & 0 & A_{4,5} & 0 & A_{4,7} & 0 \\ 0 & A_{5,1} & A_{5,2} & 0 & 0 & 0 & A_{5,6} & A_{5,7} & 0 \\ 0 & 0 & 0 & A_{6,3} & A_{6,4} & 0 & 0 & 0 & A_{6,8} \\ A_{7,0} & 0 & 0 & A_{7,3} & 0 & A_{7,5} & 0 & 0 & A_{7,8} \\ 0 & A_{8,1} & 0 & A_{8,3} & 0 & 0 & A_{8,6} & 0 & A_{8,8} \end{bmatrix} \quad (6)$$

$$\begin{bmatrix} A_{0,0} & 0 & 0 & A_{0,3} & A_{0,4} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & A_{1,1} & 0 & A_{1,3} & 0 & A_{1,5} & 0 & 0 & 0 & 0 & 0 \\ A_{2,0} & A_{2,1} & 0 & 0 & A_{2,4} & A_{2,5} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & A_{3,2} & A_{3,3} & 0 & 0 & A_{3,6} & 0 & 0 & 0 & 0 \\ A_{4,0} & 0 & A_{4,2} & 0 & A_{4,4} & 0 & A_{4,6} & 0 & 0 & 0 & 0 \\ 0 & A_{5,1} & A_{5,2} & 0 & 0 & A_{5,5} & A_{5,6} & 0 & 0 & 0 & 0 \\ A_{6,0} & 0 & 0 & 0 & 0 & 0 & 0 & A_{6,7} & A_{6,8} & 0 & 0 \\ 0 & A_{7,1} & 0 & 0 & 0 & 0 & 0 & A_{7,7} & 0 & A_{7,9} & 0 \\ A_{8,0} & A_{8,1} & 0 & 0 & 0 & 0 & 0 & 0 & A_{8,8} & A_{8,9} & 0 \\ 0 & 0 & A_{9,2} & 0 & 0 & 0 & 0 & A_{9,7} & 0 & 0 & A_{9,10} & A_{9,11} \\ A_{10,0} & 0 & A_{10,2} & 0 & 0 & 0 & 0 & 0 & A_{10,8} & 0 & A_{10,10} & A_{10,11} \\ 0 & A_{11,1} & A_{11,2} & 0 & 0 & 0 & 0 & 0 & A_{11,9} & A_{11,10} & A_{11,11} \end{bmatrix} \quad (7)$$

Using generic methods such as LU decomposition, as in previous work [21], numerical linear algebra becomes a significant bottleneck in both the predictor and corrector stages. To overcome this bottleneck, we replace the generic numerical linear algebra with closed-form solutions to the systems of linear equations with coefficient matrices (6) and (7). This replacement results in about 5x speedup for both problems.

17. Variations of the normalization

Let us provide additional details about our normalization of problems to simplify their variability and thus to make learning of the picking function σ easier.

Fig. 8 shows an example of the normalized 5pt problem. The mean direction vectors m_1, m_2 in both images are at $[0; 0; 1]$. The first correspondence, $x_{1,1}, x_{2,1}$ is on the x axis. The first image is chosen such that it contains the larger angle with its corresponding m_u . This also means that the first correspondence point has a larger x image coordinate: $x_{1,1}^{(1)} > x_{2,1}^{(1)}$. To make the normalized problems independent on the ordering of the correspondences, we sort them by their polar angles in the coordinate system in the the first image. Notice that their order may be swapped in the second image, e.g., as for $x_{2,2}, x_{3,2}$. Such a swap is mainly due to a large change of the order of depth of the corresponding points in the scene, as seen from different view points, which is in practice much less frequent than keeping the order [70].

Our normalization is chosen as the best one among several meaningful alternatives. The evaluation of the alternative normalization methods is shown in Tab. 5 for the 5pt problem and in Tab. 6 for the Scranton problem. The tables show that our strategy, labeled by A, which rotates the center of mass to zero and the farthest point on the x -axis, has the best success rate for both problems. Note that every normalization strategy performs better than when tracking without normalization.

The normalization strategies are as follows:

- Rotate the center of mass to zero, rotate the point farthest from zero to x -axis.
- Rotate the center of mass to zero with an iterative procedure, rotate the point farthest from zero to x -axis.
- Rotate the closest point to center of mass to zero and the point farthest from zero to x -axis.
- Rotate the center of mass to zero and the maximal variance to x -axis.
- Rotate the closest point to center of mass to zero and the maximal variance to x -axis.

In the case of the Scranton problem, we also have to decide which point in which view to relax on the line. Here, we consider:

Strategy	Succ. rate	Time inv.	Time HC
No inv.	0.53%	0	14.25 μs
A.	3.73%	0.43 μs	11.80 μs
B.	3.70%	0.77 μs	11.88 μs
C.	2.29%	0.36 μs	11.19 μs
D.	1.56%	1.03 μs	12.18 μs
E.	0.71%	0.68 μs	10.44 μs

Table 5. Evaluation of the normalization for the 5 pt problem. We have generated 4000 problem-solution pairs, normalized them with a given strategy and tracked HC from every p-s pair to every other. We consider strategies from Sec. 17. We measure the success rate, average time of the normalization and of HC. The track is considered successful if the squared Euclidean distance from the obtained solution to the ground-truth is less than 10^{-5} .

Strategy	Line strategy	Succ. rate	Time inv.	Time HC
No inv.	-	0.21%	0	22.11 μs
A.	a)	1.44%	0.50 μs	17.82 μs
B.	a)	1.44%	1.16 μs	17.40 μs
C.	a)	0.80%	0.82 μs	19.57 μs
C.	b)	0.32%	0.78 μs	17.37 μs
D.	a)	0.61%	1.39 μs	20.13 μs
E.	a)	0.37%	1.01 μs	20.43 μs
E.	b)	0.27%	0.93 μs	18.46 μs

Table 6. Evaluation of the normalization for the Scranton problem. We have generated 4000 problem-solution pairs, normalized them with a given strategy and tracked HC from every p-s pair to every other. We consider strategies from Sec. 17. We measure the success rate, average time of the normalization and of HC. The track is considered successful if the squared Euclidean distance from the obtained solution to the ground-truth is less than 10^{-5} .

- a) The farthest point and view.
- b) The point rotated to zero (if possible).

Our normalization induces three linear constraints for every view. The instance p of the 5pt problem consists of 2D projections of 5 points into two views, therefore $p \in \mathbb{R}^{20}$. The normalized instances live in a $20 - 2 \times 3 = 14$ dimensional subspace of \mathbb{R}^{20} . The instance p of the 4pt problem consists of 2D projections of 4 points into 3 views, $p \in \mathbb{R}^{24}$. The normalized instances live in a $24 - 3 \times 3 = 15$ dimensional subspace of \mathbb{R}^{24} .

The values of the success rate in this experiment are low because we use randomly sampled data and track from every p-s pair to every other p-s pair. Tab. 9 shows that the success rate significantly increases if we track from preselected anchors which are chosen to perform well and we select the best starting anchor with a trained classifier.

18. Study to justify engineering choices

Let us describe the data sets we use to study our engineering choices.

Training data set D^{5pt} consists of 40000 p-s pairs. We randomly sample pairs of cameras and 5-tuples of 3D points from the ETH 3D dataset [60] ‘‘Office’’ and ‘‘Terrains’’. Problem parameters p are 10D vectors of 2D image coordinates obtained by projecting the sampled 5-tuples of 3D points by the camera pairs. The corresponding solutions s are 10D vectors of the depths of the 3D points in the two cameras normalized to have the first depth equal to 1. Data set D^{Scr} , consisting of 40000 p-s pairs, is constructed analogously by sampling 4-tuples of 3D points and triplets of cameras. Data sets are used to select anchor sets and to train our model for selecting the starting problem-solution pair.

Anchor sets $A_{50}^{5pt}, A_{75}^{5pt}, A_{90}^{5pt}, A_{100}^{5pt}$ are selected by the procedure described in Sec. 4 such that $A_{50}^{5pt} \subset A_{75}^{5pt} \subset A_{90}^{5pt} \subset A_{100}^{5pt} \subset D^{5pt}$ where A_{50}^{5pt} of 8 anchors covers 50% of problems in D^{5pt} , A_{75}^{5pt} of 26 anchors covers 75% of problems in D^{5pt} , A_{90}^{5pt} of 70 anchors covers 90% of problems in D^{5pt} and A_{100}^{5pt} of 465 anchors covers 100% of problems in D^{5pt} . Data sets, A_{50}^{Scr} , of 16 anchors, A_{75}^{Scr} , of 50 anchors, A_{90}^{Scr} , of 134 anchors and A_{100}^{Scr} , of 1205 anchors, are constructed analogously from A^{Scr} .

Test data set V^{5pt} consists of 60000 p-s pairs constructed as above from the ETH 3D dataset ‘‘Delivery area’’ and ‘‘Facade’’.

	Succ. rate $\mu_s \pm \delta_s$ [%] / Time μ_t [μs]		
	5pt problem		
Solving technique	M2 [43]	MINUS [20]	OUR
\mathbb{C} -HC, All Sols	$98.9 \pm 0.2 / 1.9 \times 10^5$	$97.7 \pm 0.2 / 15197.1$	$97.7 \pm 0.1 / 5133.7$
\mathbb{C} -HC, \mathbb{R} Sols	$56.1 \pm 3.3 / 1.2 \times 10^5$	$55.3 \pm 2.6 / 5704.7$	$54.7 \pm 3.2 / 1895.1$
\mathbb{C} -HC, Fab Sol	$12.9 \pm 0.9 / 9.8 \times 10^4$	$12.0 \pm 1.5 / 638.0$	$13.1 \pm 1.4 / 165.8$
\mathbb{R} -HC, \mathbb{R} Sols	$9.9 \pm 1.7 / 5.7 \times 10^4$	$9.7 \pm 1.5 / 647.9$	$9.7 \pm 1.5 / 106.8$
\mathbb{R} -HC, Fab Sol	$3.4 \pm 1.3 / 4.4 \times 10^4$	$2.7 \pm 0.8 / 67.6$	$2.7 \pm 0.8 / 11.1$
Newton, Fab Sol	$4.0 \pm 0.6 / 1.4 \times 10^4$	$4.0 \pm 0.7 / 8.5$	$4.0 \pm 0.7 / 1.3$

	Succ. rate $\mu_s \pm \delta_s$ [%] / Time μ_t [μs]	
	Scranton	
Solving technique	MINUS [20]	OUR
\mathbb{C} -HC, All Sols	$95.5 \pm 2.6 / 608332.0$	$95.7 \pm 2.6 / 187364.5$
\mathbb{C} -HC, \mathbb{R} Sols	$22.7 \pm 1.6 / 47961.4$	$22.5 \pm 1.7 / 14905.1$
\mathbb{C} -HC, Fab Sol	$3.5 \pm 0.7 / 1280.7$	$3.3 \pm 0.9 / 405.7$
\mathbb{R} -HC, \mathbb{R} Sols	$7.5 \pm 1.2 / 2548.4$	$7.6 \pm 1.2 / 484.4$
\mathbb{R} -HC, Fab Sol	$1.2 \pm 0.3 / 70.2$	$1.2 \pm 0.3 / 14.1$
Newton, Fab Sol	$1.9 \pm 0.6 / 8.6$	$2.0 \pm 0.4 / 1.4$

Table 7. Homotopy continuation study. The rows represent variations mixing the solving technique of complex (\mathbb{C} -HC) and real (\mathbb{R} -HC) homotopy continuation, the Newton’s local method (Newton) with starting from all solutions (All Sols), real solutions only (\mathbb{R} Sols), and the fabricated solution only (Fab Sol) of a problem. The columns represent different implementations of homotopy continuation. M2 denotes the off-the-shelf implementation in Macaulay2 [43]. MINUS denotes the implementation based on [20]. OUR denotes our efficient implementation. To compare MINUS and OUR, we selected 20 subsets P_i , $i = 1, \dots, 20$, each containing 50 random problems from P^{5pt} for 5pt problem and from P^{Scr} for Scranton. All problems were normalized. For each P_i , we compute the success rate μ_{s_i} of $50^2 - 50$ of homotopy continuations from each start problem $p_{i,j} \in P_i$ to each different target problem $p_{i,k} \in P_i$. We consider a homotopy continuation successful if the fabricated solution of the target problem $p_{i,k}$ is among the solutions reached by the homotopy continuation within 10^{-5} Euclidean distance in the solution space of depths. We report the mean success rate $\mu_s = \text{mean}(\mu_{s_i})$ and the standard deviation $\delta_s = \text{std}(\mu_{s_i})$ over all P_i ’s for each implementation and mean computation times μ_t .

Data set V^{Scr} , consisting of 60000 p-s pairs, is constructed analogously. We use V^{5pt} and V^{Scr} in the experiments studying anchor set selection methods reported in Tab. 4.

18.1. Comparison of different tracking approaches

Data set P^{5pt} consists of 3751 p-s pairs sampled from the ETH 3D dataset ‘‘Courtyard’’. A^{5pt} and P^{5pt} are disjoint. All problems in P^{5pt} are checked to be generic and can be used as good starting problem-solution pairs. We select 20 random pairwise disjoint $P_i^{5pt} \subset P^{5pt}$ consisting of 50 problem-solution pairs. Data sets P^{Scr} , consisting of 5727 problem-solutions pairs, and P_i^{Scr} , consisting of 50 p-s pairs, are constructed analogously. We use P_i^{5pt} ’s and P_i^{Scr} in the experiments studying variations of the homotopy continuation methods reported in Tab. 7. In this table, we measure the success rate for a given subset P_i as a percentage of different pairs $p_{i,j} \in P_i, p_{i,k} \in P_i$ for which the fabricated solution to the target problem $p_{i,k}$ can be recovered when tracking from $p_{i,j}$ to $p_{i,k}$. Then, we find the mean success rate μ_s , and the standard deviation δ_s over all subsets $P_i, i \in \{1, \dots, 20\}$.

Tab. 7 shows that for every setting, our evaluation (Sec. 16.1) brings about 5x speedup over the previous work [20] without any impact on the success rate.

The table also shows that Homotopy Continuation in the complex domain tracked from every solution has almost 100% success rate, but the running time of the solver is prohibitively slow to be used in the RANSAC scheme. We can see that reducing the number of tracks, as well as tracking in \mathbb{R} instead of \mathbb{C} can significantly reduce the running time (about 10000x for the Scranton problem) at the cost of a lower success rate. Note (Tab. 9), that the issue with the low success rate may be remedied by selecting an appropriate starting problem-solution pair (p_0, s_0) for every target problem p .

Tab. 7 also shows that using the Newton method may be a promising approach for solving the minimal problems, as it has a higher success rate and lower running time than Homotopy Continuation with the same setting. Therefore, we have found

Method	ρ [%]	μ_t [\mu s]	ϵ_t [\mu s]
N3 + B1	0.01	1.73	15159.7
N3 + MLP	0.12	8.2	6811.2
N15 + B1	2.1	2.6	119.7
N15 + MLP	10.6	10.6	100.0
HC + B1	4.2	18.7	442.1
HC + MLP	26.3	16.2	61.6

Table 8. Study of methods for starting problem selection and tracking for Scranton problem. The strategies are evaluated on datasets *Deilvery_area* and *Facade*. Tracking methods: ‘N3’: Newton method with 3 steps, ‘N15’: Newton method with 15 steps (this number of steps maximizes the efficient time), ‘HC’: Homotopy Continuation as described in Sec 6. Problem selection methods: ‘B1’: Tracks always from the same anchor. ‘MLP’: selecting the starting problem as the one with the highest score given by the MLP.

Result [%]	All pairs		MLP	
	5 pt	4 pt	5pt	4pt
Fabricated sol.	3.64	1.49	38.8	29.2
1 rel. pose correct	-	0.74	-	4.53
Other meaningful	9.50	13.43	31.6	34.2
Sol. with det -1	0.00	0.00	0.00	0.00
Sol. with zeros	0.03	1.53	0.01	0.52
Negative sol.	0.52	10.86	0.83	8.54
Failed track	86.31	71.94	28.8	23.0

Table 9. Percentage of different results of real homotopy continuation. In “All pairs”, we track from each $p_i \in P$ to each other $p_j \in P$. In “MLP”, we solve each $p \in V$ by selecting a starting p-s pair (p_0, s_0) from A_{90} , and by tracking HC from p_0 to p . “Fabricated sol.” means that the fabricated solution was reached, “1 rel. pose correct” means that if we convert the obtained solution of the Scranton problem to the relative poses, then at least one of three relative poses is correct. “Other meaningful” means that a non-fabricated solution with positive depths and valid rotation matrix was reached, “Sol. with det -1” means that the matrix R is not a valid rotation. “Sol. with zeros” means that some depths are equal to 0, “Negative sol.” means that some depths are negative, and “Failed track” means that the HC track has failed and, thus, the solution has not been found.

starting p-s pairs (Sec. 4) and trained the MLP for the Newton method. The comparison of the solvers using the Homotopy Continuation and the Newton method is shown in Tab. 8. We can see that the effective time (the average time needed to obtain one correct solution in the RANSAC scheme) of the solver using Homotopy Continuation is about 2x lower than the effective time of the solver using Newton method. While the solvers using Newton method are faster, the success rate of the solver, which combines Homotopy Continuation and MLP classifier has a higher success rate. The possible explanation for this is that the Newton method behaves “more randomly” than the Homotopy Continuation, and therefore, it is more difficult to train.

Tab. 9 provides a more detailed analysis of the frequency of different results of real Homotopy Continuation tracked from the fabricated solution. We consider two different settings. In the “All pairs” setting, we track from each $p_i \in P$ to each other $p_j \in P$. Then, in the MLP setting, we select a starting p-s pair (p_0, s_0) from A_{90} , and track from p_0 to p . We measure how often we reach the fabricated solution, non-fabricated meaningful solution, a non-valid solution, and how often HC fails and does not deliver any solution. The table shows that the MLP increases the probability of reaching the fabricated solution about 10x for 5pt and 20x for Scranton. The probability of reaching another meaningful solution increases about 3x for both problems, while the probability of reaching a non-valid solution and the probability of failing decreases.

18.2. Comparison of different settings our solver

Here, we show how different settings of the solver influence the resulting success rate ρ , mean running time μ_t , and efficient time ϵ_t . We perform this study on our solver for the Scranton problem. For every experiment, we use the settings from the main paper, except that one parameter is varied. The solver uses anchors A_{90}^{Scr} and it is evaluated on data V^{Scr} . Our goal is to show that we have selected the optimal settings for our solver.

Formulation	Succ. rate
First depth fixed	2.30 %
Quan symmetrical [52]	1.47 %
Quan asymmetrical [52]	1.13 %

Table 10. Scranton dehomogenization study. Rows correspond to different formulations of the problems. For each method, we compute the success rate of $4000^2 - 4000$ HC calls from each starting p-s pair to each other target p-s pair. We consider the result successful if the fabricated solution of the target problem is and the result computed by HC are sufficiently close ($\leq 10^{-5}$ Euclidean distance in the solution space of depths.)

Layer size	ρ [%]	μ_t [μs]	ϵ_t [μs]
100	27.8	20.3	73.1
200	31.3	30.8	98.3
500	34.7	79.0	227.6

Table 11. Study of different MLP sizes. The strategies are evaluated on datasets *Deilvery_area* and *Facade*. Scranton problem, MLP+HC, A_{75} . Rows correspond to different sizes of hidden MLP layers.

# Tracks	ρ [%]	μ_t [μs]	ϵ_t [μs]
1	29.2	19.6	67.0
2	37.2	33.3	89.6
3	42.2	45.0	106.8
4	45.9	60.8	132.6
8	56.4	118.1	209.5
16	67.9	245.3	361.5

Table 12. Number of tracks study. The strategies are evaluated on datasets *Deilvery_area* and *Facade*. Scranton problem, MLP+HC, A_{90} . Rows correspond to different numbers of tracks conducted after the MLP is evaluated.

We compare different methods of dehomogenizing Scranton problem in Tab. 10. This justifies our choice of fixing the first depth $\lambda_{1,1} = 1$, which gives a superior success rate compared to “symmetric” and “asymmetric” dehomogenization proposed in [52].

See Fig. 9 for a code snippet showing the structure of our MLP model. In Tab. 11, we show how the success rate ρ and running time μ_t depends on the size of the MLP that is used for the classification of the anchors. Larger networks have a higher success rate. However, the efficient time of the solver with the smaller MLP is better because the time needed for the evaluation of the MLP grows faster than the success rate.

Finally, in Tab. 12, we show how the number of HC tracks per problem influences the success rate ρ , running time μ_t , and efficient time ϵ_t . In this experiment, we use the MLP trained in 5, and we perform n tracks from n anchors with the highest score. We consider the solution successful if at least one track reaches the fabricated solution of the target problem. Much like using larger MLP as Tab. 11, such a strategy involving multiple anchors suggests a future approach to improving our solvers’ success rates. However, we note that the efficient time ϵ_t in Tab. 12 grows with the number of tracks, since the evaluation time grows faster than the success rate.

References

[70] A.L. Yuille and T. Poggio. A generalized ordering constraint for stereo correspondence. A.I. Memo 777, AI Lab, MIT, 1994.

```

class Net(nn.Module):
    def __init__(self, anchors):
        super(Net, self).__init__()
        self.fc1 = nn.Linear(20,100)
        self.relu1 = nn.PReLU(100, 0.25)
        self.fc2 = nn.Linear(100,100)
        self.relu2 = nn.PReLU(100, 0.25)
        self.fc4 = nn.Linear(100,100)
        self.relu4 = nn.PReLU(100, 0.25)
        self.fc5 = nn.Linear(100,100)
        self.relu5 = nn.PReLU(100, 0.25)
        self.fc6 = nn.Linear(100,100)
        self.relu6 = nn.PReLU(100, 0.25)
        self.fc7 = nn.Linear(100,100)
        self.relu7 = nn.PReLU(100, 0.25)
        self.drop3 = nn.Dropout(0.5)
        self.fc3 = nn.Linear(100,anchors+1)
    def forward(self, x):
        x = self.relu1(self.fc1(x))
        x = self.relu2(self.fc2(x))
        x = self.relu4(self.fc4(x))
        x = self.relu5(self.fc5(x))
        x = self.relu6(self.fc6(x))
        x = self.relu7(self.fc7(x))
        x = self.drop3(x)
        return self.fc3(x)

```

Figure 9. Code snippet describing our MLP.