# Bound and Conquer: Improving Triangulation by Enforcing Consistency

Adam Scholefield, *Member, IEEE,* Alireza Ghasemi, *Student Member, IEEE,*
and Martin Vetterli, *Fellow, IEEE*

## SUPPLEMENTARY MATERIAL

*Abstract*—In this supplementary material, we give the proofs for the two theorems and one proposition of the paper. In addition, we formalize the equivalence between camera localization and triangulation.

— ◆ —

## 1 PROOF OF THEOREM 1

**Theorem 1.** *Consider a multi-camera system of $M$ cameras, each with an $N \times N$ pixel image sensor and define a fixed region of interest, $\mathcal{R}$, with a finite non-zero volume.*

*If we assume that the only source of uncertainty is pixelization, the expected $\ell_2$ reconstruction error of any triangulation algorithm is lower-bounded by a term that is inverse-quadratically dependent on the number of cameras; i.e.,*

$$\mathbb{E}\left(\left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2\right) = \Omega\left(\frac{1}{M^2}\right), \tag{1}$$

*where $\mathbf{U} \in \mathcal{R}$ is any point in the region of interest, and $\hat{\mathbf{U}}$ is the result of reconstructing $\mathbf{U}$, from its images in the multi-camera system, using any triangulation algorithm. Here, the expectation is taken over the location of the point $\mathbf{U}$ in the region of interest.*

*Proof.* A single $N \times N$ pixel camera partitions the world space into $N^2$ regions. Combined with the partitions of other cameras, this leads to a finite number of partitions. Therefore, when a multi-camera system views the region of interest, it splits it into a finite number of partitions. Let $\mathcal{P}$ be the set containing the resulting partitions of $\mathcal{R}$.

We can now consider the the expected $\ell_2$ reconstruction error split over these partitions:

$$\mathbb{E}\left(\left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2\right) = \frac{1}{\mathcal{V}(\mathcal{R})} \iiint_{\mathcal{R}} \left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2 d\mathbf{U}$$

$$= \frac{1}{\mathcal{V}(\mathcal{R})} \sum_{\mathcal{C} \in \mathcal{P}} \iiint_{\mathcal{C}} \left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2 d\mathbf{U}. \tag{2}$$

Here, $\mathcal{V}(\mathcal{R})$ denotes the volume of the region of interest.

The localization error over each partition depends on both its size and shape. Among all partitions with the same volume, the value of this integral would be minimized if the shape was a sphere and the estimate, $\hat{\mathbf{U}}$, was at the centre of that sphere:

$$\iiint_{\mathcal{C}} \left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2 d\mathbf{U} \geq \iiint_{H_r} \|\boldsymbol{c} - \mathbf{U}\|_2^2 d\mathbf{U}, \tag{3}$$

where $H_r$ is a sphere with centre $\boldsymbol{c}$ and radius $r = \sqrt[3]{3\mathcal{V}(\mathcal{C})/(4\pi)}$. Evaluating this integral, we obtain

$$\iiint_{H_r} \|\boldsymbol{c} - \mathbf{U}\|_2^2 d\mathbf{U} = \frac{4\pi}{5} r^5 = K\mathcal{V}(\mathcal{C})^{\frac{5}{3}}, \tag{4}$$

where $K = \frac{4\pi}{5} \sqrt[3]{\frac{3}{4\pi}}$.

Combining (2), (3) and (4) yields

$$\mathbb{E}\left(\left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2\right) > \frac{K}{\mathcal{V}(\mathcal{R})} \sum_{\mathcal{C} \in \mathcal{P}} \mathcal{V}(\mathcal{C})^{\frac{5}{3}}.$$

This lower-bound would be minimized if the available volume, $\mathcal{V}(\mathcal{R})$, was split equally among each of the regions in the sum:

$$\mathbb{E}\left(\left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2\right) > K\frac{1}{\mathcal{V}(\mathcal{R})} \sum_{\mathcal{C} \in \mathcal{P}} \left(\frac{\mathcal{V}(\mathcal{R})}{\#\mathcal{P}}\right)^{\frac{5}{3}}$$

$$= K\left(\frac{\mathcal{V}(\mathcal{R})}{\#\mathcal{P}}\right)^{\frac{2}{3}}. \tag{5}$$

Here, $\#\mathcal{P}$ is the number of partitions (the cardinality of $\mathcal{P}$).

Since the volume of the region of interest, $\mathcal{V}(\mathcal{R})$, is fixed, we just need to consider how the number of regions, $\#\mathcal{P}$, grows as we add more cameras to the system. To do so, we first consider how many regions can be created from $L$ planes in $\mathbb{R}^3$. In computational geometry, this quantity is known as the number of cells in an arrangement of hyperplanes (see for example [Goodman and Pollack 1986]). It can be shown that, with

$L$ planes, the 3-D space $\mathbb{R}^3$ is partitioned into at most $k$ regions and $k$ grows cubically with $L$, i.e. $k = \mathcal{O}(L^3)$.

In our case, partitions are created by the boundaries of the pixels. We can see that each camera in a multi-camera system partitions the space with at most $2(N+1)$ planes intersected by rays starting from the camera centre and passing through pixel boundaries[1] (we have an upper bound since some or all of these planes may not pass through the region of interest). Therefore, for $M$ cameras, we have at most $2M(N+1)$ such planes passing through the region of interest and thus we can conclude that the number of regions $(\#\mathcal{P})$ satisfies

$$\#\mathcal{P} = \mathcal{O}\left(M^3 N^3\right). \tag{6}$$

Substituting (6) into (5) gives

$$\mathbb{E}\left(\left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2\right) = \Omega\left(\frac{\mathcal{V}(\mathcal{R})}{M^2 N^2}\right),$$

which proves that $\mathbb{E}\left(\left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2\right) = \Omega\left(1/M^2\right)$ for fixed $N$ and $\mathcal{R}$, hence the fact that best possible decay rate for a geometric reconstruction algorithm is quadratic. $\quad\square$

## 2 PROOF OF PROPOSITION 1

**Proposition 1.** *Consider a multi-camera system viewing a point and assume that the image points are subjected to $\ell_q$-norm bounded noise:*

$$\|\mathbf{u}_i - \mathcal{P}_i(\mathbf{U})\|_q \leq \delta \quad \text{for } i = 1...M.$$

*Then, any algorithm that minimizes the $(\ell_q, \ell_\infty)$-norm of the reprojection error is a q-consistent triangulation algorithm.*

*Proof.* The proof will be by contradiction. Let $\hat{\mathbf{U}}$ be the minimum $(\ell_q, \ell_\infty)$-norm solution:

$$\hat{\mathbf{U}} = \arg\min_{\mathbf{U}} \max_{i=1..M} \|\mathbf{u}_i - \mathcal{P}_i(\mathbf{U})\|_q. \tag{7}$$

Assume that $\hat{\mathbf{U}}$ is not $q$-consistent. Then, there exists an $i$ such that

$$\left\|\mathbf{u}_i - \mathcal{P}_i(\hat{\mathbf{U}})\right\|_q > \delta. \tag{8}$$

Alternatively, let $\boldsymbol{U}_c$ be a $q$-consistent estimate. By definition,

$$\|\mathbf{u}_i - \mathcal{P}_i(\boldsymbol{U}_c)\|_q \leq \delta \quad \text{for all } i = 1...M. \tag{9}$$

Therefore,

$$\max_{i=1..M} \left\|\mathbf{u}_i - \mathcal{P}_i(\hat{\mathbf{U}})\right\|_q > \max_{i=1..M} \|\mathbf{u}_i - \mathcal{P}_i(\boldsymbol{U}_c)\|_q. \tag{10}$$

But, this contradicts (7) and thus $\hat{\mathbf{U}}$ must be $q$-consistent. $\quad\square$

1. In the case of orthogonal projection, rays do not originate from the centre of the camera, but their cardinality and hence the rest of the proof remain unchanged.

## 3 PROOF OF THEOREM 2

The proof makes use of the following corollary.

**Corollary 1** (Powell and Whitehouse 2016). *Assume random vectors $\{\phi_i\}_{i=1}^M \subset \mathbb{R}^d$ are i.i.d. and uniformly distributed on the unit d-dimensional sphere. Suppose a point in $\mathbb{R}^d$ is orthogonal projected onto the random vectors and subjected to zero-mean uniform bounded noise with bandwidth $\delta$. Then, constants $c_1, c_2 > 0$ exist such that*

$$\mathbb{E}\{(W_M)^2\} \leq \frac{c_2 d^3 \delta^3}{M^2}, \qquad \forall M \geq c_1 d \ln d. \tag{11}$$

*Here, $W_M$ is the radius of the smallest d-dimensional sphere containing the consistency region formed from the M samples.*

*Proof.* See [Powell and Whitehouse 2016, Corollary 6.2]. $\quad\square$

**Theorem 2.** *Place $M$ cameras in a plane, i.i.d. uniformly at random on a finite radius circle oriented towards the centre of the circle. Define the region of interest, $\mathcal{R}$, to be the intersection of the field of view of all cameras as $M \to \infty$ and place a point anywhere in this region.*

*Furthermore, assume that the images of the world point in the cameras are perturbed with $\ell_\infty$ uniform bounded noise; i.e., for the world point $\mathbf{U}$, the image $\mathbf{u}_i$ in the $i$-th camera is*

$$\mathbf{u}_i = \mathcal{P}_i(\mathbf{U}) + \boldsymbol{\epsilon}_i, \tag{12}$$

*where $\boldsymbol{\epsilon}_i$ is zero-mean uniform bounded random satisfying $\|\boldsymbol{\epsilon}_i\|_\infty \leq \delta$.*

*In this situation, the expected $\ell_2$ reconstruction error of any $\infty$-consistent triangulation algorithm is upper-bounded by a term which decreases quadratically with the number of cameras; i.e.,*

$$\mathbb{E}\left(\left\|\hat{\mathbf{U}} - \mathbf{U}\right\|_2^2\right) = \mathcal{O}\left(\frac{1}{M^2}\right), \tag{13}$$

*where $\mathbf{U} \in \mathcal{R}$ is any point in the region of interest, and $\hat{\mathbf{U}}$ is the result of reconstructing $\mathbf{U}$, from its images in the multi-camera system, using a $\infty$-consistent triangulation algorithm. Here, the expectation is taken over both the noise and the camera locations.*

*Proof.* Let $\mathbf{U} = [U_X, U_Y, U_Z]^T$, $\boldsymbol{\epsilon_i} = [\epsilon_{i,x}, \epsilon_{i,y}]^T$ and assume, without loss of generality, that the circle lies in the $X$-$Z$ plane.

Before considering the central projection case, we assume the cameras are orthographic. In this case, the vertical coordinate of the image points are given by

$$u_{i,y} = U_Y + \epsilon_{i,y}, \quad i \in [1, M]. \tag{14}$$

The $y$-coordinate of the world-space $\infty$-consistency region is the following 1-D interval:

$$\left\{\hat{U}_Y : \max_i \epsilon_{i,y} - \delta \leq \hat{U}_Y - U_Y \leq \min_i \epsilon_{i,y} + \delta\right\}.$$

Therefore, the maximum reconstruction error in this coordinate is

$$\mathcal{E} := \max_{\hat{U}_Y \in \mathcal{C}_y} \left| \hat{U}_Y - U_Y \right|$$
$$= \max \left\{ \left| \max_i \epsilon_{i,y} - \delta \right|, \left| \min_i \epsilon_{i,y} + \delta \right| \right\}$$
$$= \max \left\{ \mathcal{E}_l, \mathcal{E}_u \right\}, \tag{15}$$

where $\mathcal{E}_l := |\max_i \epsilon_{i,y} - \delta| = \delta - \max_i \epsilon_{i,y}$ and $\mathcal{E}_u := |\min_i \epsilon_{i,y} + \delta| = \min_i \epsilon_{i,y} + \delta$ are the absolute values of the lower and upper bounds, respectively.

The expected maximum squared error can be computed as

$$\mathbb{E}\left(\mathcal{E}^2\right) = \int_0^\infty \lambda^2 \frac{d\mathbb{P}\left(\mathcal{E} \leq \lambda\right)}{d\lambda} d\lambda = 2 \int_0^\infty \lambda \mathbb{P}\left(\mathcal{E} \geq \lambda\right) d\lambda,$$

where $\mathbb{P}(\cdot)$ represents the probability. Furthermore, from (15), we have

$$\mathbb{P}\left(\mathcal{E} \geq \lambda\right) = \mathbb{P}\left(\mathcal{E}_l \geq \lambda \cup \mathcal{E}_u \geq \lambda\right)$$
$$= \mathbb{P}\left(\mathcal{E}_l \geq \lambda\right) + \mathbb{P}\left(\mathcal{E}_u \geq \lambda\right) - \mathbb{P}\left(\mathcal{E}_l \geq \lambda \cup \mathcal{E}_u \geq \lambda\right).$$

Recalling that $\epsilon_{i,y}$ is uniformly distributed on $[-\delta, \delta]$, we can calculate each term as

$$\mathbb{P}\left(\mathcal{E}_l \geq \lambda\right) = \mathbb{P}\left(\epsilon_{i,y} \leq \delta - \lambda, \, i \in [1, M]\right)$$
$$= \left(1 - \frac{\lambda}{2\delta}\right)^M \quad \text{for } 0 \leq \lambda \leq 2\delta,$$

$$\mathbb{P}\left(\mathcal{E}_u \geq \lambda\right) = \mathbb{P}\left(\epsilon_{i,y} \geq \lambda - \delta, \, i \in [1, M]\right)$$
$$= \left(1 - \frac{\lambda}{2\delta}\right)^M \quad \text{for } 0 \leq \lambda \leq 2\delta,$$

and

$$\mathbb{P}\left(\mathcal{E}_l \geq \lambda \cup \mathcal{E}_u \geq \lambda\right) = \mathbb{P}\left(\lambda - \delta \leq \epsilon_{i,y} \leq \delta - \lambda, i \in [1, M]\right)$$
$$= \left(1 - \frac{\lambda}{\delta}\right)^M \quad \text{for } 0 \leq \lambda \leq \delta.$$

Therefore,

$$\mathbb{E}\left(\mathcal{E}^2\right) = 4 \int_0^{2\delta} \lambda \left(1 - \frac{\lambda}{2\delta}\right)^M d\lambda - 2 \int_0^\delta \lambda \left(1 - \frac{\lambda}{\delta}\right)^M d\lambda$$
$$= \frac{14\delta^2}{(M+1)(M+2)}$$

and so

$$\mathbb{E}\left(\left|\hat{U}_Y - U_Y\right|^2\right) \leq \frac{14\delta^2}{(M+1)(M+2)} < \frac{14\delta^2}{M^2}, \tag{16}$$

for any $\infty$-consistent estimate $\hat{U}_Y$ of $U_Y$.

Let's now consider the horizontal coordinate of the image points. If we continue to assume orthographic projection, we have

$$\begin{bmatrix} u_{1,x} \\ u_{2,x} \\ \vdots \\ u_{M,x} \end{bmatrix} = \begin{bmatrix} -\sin\theta_1 & \cos\theta_1 \\ -\sin\theta_2 & \cos\theta_2 \\ \vdots & \vdots \\ -\sin\theta_M & \cos\theta_M \end{bmatrix} \begin{bmatrix} U_X \\ U_Z \end{bmatrix} + \begin{bmatrix} \epsilon_{1,x} \\ \epsilon_{2,x} \\ \vdots \\ \epsilon_{M,x} \end{bmatrix}.$$

This is a linear inverse problem in two dimensions, seeking unknowns $U_X$ and $U_Z$. The solution defines the $x$ and $z$ coordinates of the world-space $\infty$-consistency region. Since the $x$ and $z$ coordinates cannot be split into 1-D intervals, the geometry of the resulting 2-D consistency region is more complex than for the $y$ coordinate; however, the assumption that the cameras are uniformly distributed on the circle simplifies this geometrical dependence. This is exploited in [Powell and Whitehouse 2016] to prove various bounds including Corollary 1. Directly applying this corollary yields

$$\mathbb{E}\left( \left\| \begin{bmatrix} \hat{U}_X \\ \hat{U}_Z \end{bmatrix} - \begin{bmatrix} U_X \\ U_Z \end{bmatrix} \right\|_2^2 \right) \leq \frac{K_1 \delta^2}{M^2}, \tag{17}$$

for any $\infty$-consistent estimate $[\hat{U}_X, \hat{U}_Z]^T$ of $[U_X, U_Z]^T$. Here, $K_1$ is a constant independent of the number of cameras and the support of the bounded noise.

Combining (16) and (17) yields

$$\mathbb{E}\left( \left\| \hat{\mathbf{U}} - \mathbf{U} \right\|_2^2 \right) \leq \frac{K_2 \delta^2}{M^2},$$

for the orthographic case. Here $K_2$ is a constant independent of the number of cameras and the support of the bounded noise.

Now, to extend this result to the pinhole camera case, let $r$ be the radius of the circle and $f$ be the focal length of all cameras. Then, the pinhole projection $\infty$-consistency region corresponding to an image point measurement with a noise bandwidth of $\delta$ has a smaller volume than the $\infty$-consistency region of an orthogonal projection, with larger bandwidth and a circle of interest of radius $r - f$. The bandwidth $\delta_{equiv}$ of this corresponding parallel projection camera is computed as

$$\delta_{equiv} = \delta \left(1 + \frac{r-f}{f}\right) = \delta \left(\frac{r}{f}\right). \tag{18}$$

This means that we can upper-bound the reconstruction error of a circular array of $M$ pinhole cameras with a measurement error bandwidth of $\delta$, with the reconstruction error of a circular array of *parallel* cameras, with the bandwidth $\delta_{equiv}$ as defined above. Using this fact, we have the following bound:

$$\mathbb{E}\left( \left\| \hat{\mathbf{U}} - \mathbf{U} \right\|_2^2 \right) \leq \frac{K_2 \delta^2 r^2}{M^2 f^2}. \tag{19}$$

$\square$

# 4 EQUIVALENCE BETWEEN TRIANGULATION AND CAMERA LOCALIZATION

As briefly mentioned in the main text, triangulation is equivalent to a restricted version of camera localization, where one assume that the orientation of the camera is known. Therefore, the scaling laws derived for triangulation also apply to this restricted version of camera localization. To formalize this, we now define the camera localization problem and prove the equivalence.

**Definition 4.** A *camera localization problem* takes as input

$$L = \{\mathbf{K}, \mathbf{R}\} \cup \{(\mathbf{u}_i, \mathbf{U}_i) | 1 \le i \le M\}, \tag{20}$$

and estimates the unknown camera centre $\mathbf{C}$ as follows:

$$\hat{\mathbf{C}} = \arg\min_{\tilde{\mathbf{C}}} \sum_{i=1}^{M} \left\| \mathbf{u}_i - \tilde{\mathcal{P}}(\mathbf{U}_i) \right\|_{p'}^{p}. \tag{21}$$

Here, $\tilde{\mathcal{P}}(\cdot)$ denotes denotes the projection operator corresponding to the camera matrix $\tilde{\mathbf{P}} = \mathbf{K}\mathbf{R}[\mathbf{I}|-\tilde{\mathbf{C}}]$, where $\mathbf{K}$ and $\mathbf{R}$ are the known intrinsic and orientation matrices of the camera, respectively. In addition, the $(\mathbf{u}_i, \mathbf{U}_i)$ pairs denote $M$ world points along with their projections onto the image plane of the camera. Finally, $p'$ and $p$ are the image-space and residual-space norms, respectively.

**Proposition 2.** *The solution to a camera localization problem* $L = \{\mathbf{K}, \mathbf{R}\} \cup \{(\mathbf{u}_i, \mathbf{U}_i) | 1 \le i \le M\}$ *is the same as the solution to a triangulation problem* $T = \{(\mathbf{u}_i, \mathbf{P}_i) | 1 \le i \le M\}$, *where* $\mathbf{P}_i = \mathbf{K}\mathbf{R}[\mathbf{I}| - \mathbf{U}_i]$.

*Proof.* From Definition 4, the solution to $L$ is given by

$$\hat{\mathbf{C}} = \arg\min_{\tilde{\mathbf{C}}} \sum_{i=1}^{M} \left\| \mathbf{u}_i - \tilde{\mathcal{P}}(\mathbf{U}_i) \right\|_{p'}^{p}, \tag{22}$$

where $\tilde{\mathcal{P}}(.)$ is the projection operator corresponding to the camera matrix $\tilde{\mathbf{P}} = \mathbf{K}\mathbf{R}[\mathbf{I} - \tilde{\mathbf{C}}]$. Now,

$$\mathbf{K}\mathbf{R}[\mathbf{I}| - \tilde{\mathbf{C}}] \begin{bmatrix} \mathbf{U}_i \\ 1 \end{bmatrix} = -\mathbf{K}\mathbf{R}[\mathbf{I}| - \mathbf{U}_i] \begin{bmatrix} \tilde{\mathbf{C}} \\ 1 \end{bmatrix}. \tag{23}$$

Therefore, denoting the projection operator corresponding to the camera matrix $\mathbf{P}_i = \mathbf{K}\mathbf{R}[\mathbf{I}| - \mathbf{U}_i]$ by $\mathcal{P}_i(.)$, we have

$$\tilde{\mathcal{P}}(\mathbf{U}_i) = \mathcal{P}_i(\tilde{\mathbf{C}}). \tag{24}$$

Note that there is no sign difference here because the sign difference in (23) cancels when converting from homogeneous to Cartesian coordinates.

Substituting (24) into (22) yields

$$\hat{\mathbf{C}} = \arg\min_{\tilde{\mathbf{C}}} \sum_{i=1}^{M} \left\| \mathbf{u}_i - \mathcal{P}_i(\tilde{\mathbf{C}}) \right\|_{p'}^{p}. \tag{25}$$

By replacing $\tilde{\mathbf{C}}$ with $\mathbf{U}$ and comparing to Definition 1 of the main text, we see that $\hat{\mathbf{C}} = \hat{\mathbf{U}}$, hence proving the proposition. $\qquad \square$