

MATTHEW KLINGENSMITH

# THESIS PROPOSAL: ARTICULATED 3D SLAM



## *Contents*

*Background* 9

*Framework* 19

*Experiments* 27

*Research Questions* 31

*Timeline* 35

*Bibliography* 37



# *Introduction*

To interact with and understand the world, robots must interpret and fuse noisy and ambiguous signals from their sensors. Consider the case of a multi-jointed robot arm with a 2D or 3D camera attached to its end effector (Fig. 1). The robot receives noisy signals from its joint encoders, the hand-mounted camera, and other sensors. How can the robot make sense of the world given these noisy, heterogeneous signals? How can it then take action based on this information? To successfully complete a task (such as grasping an object), it is critical that the robot have an accurate geometric representation of the world around it.



Figure 1: The ADA robot with hand-mounted camera manipulating a book in a cluttered bookshelf.

Dense 3D reconstruction is a well-studied problem in computer vision and robotics. Given a series of noisy sensor readings a unknown poses, the task is to infer the true geometric structure and color of the world. This is a variant of the Simultaneous Localization and Mapping (SLAM) problem <sup>1</sup>, since it requires simultaneously estimating the state of the sensor and the scene at once. In the context

<sup>1</sup> See page 13

of robotic manipulation, dense 3D reconstruction is useful for scene understanding, planning, grasping, and more.

Typical techniques for solving this problem focus on handheld sensors, and thus assume that little is known about the pose of the sensor aside from weak priors on the magnitude of the motion (though it is sometimes assumed that other sources of motion estimation, such as markers, gyroscopes, or inertial measurement units are available). For this reason, many prior works directly estimate the pose of a handheld sensor using visual (or depth) data alone. Depending on the type of sensor, the lack of additional information can lead to inaccuracies in pose and ultimately in the resulting 3D reconstruction. For instance, using a handheld depth camera, it is impossible to detect motion parallel to a flat plane using the (unchanging) depth data alone<sup>2</sup>, without additional sensors or *a. priori* knowledge of how the sensor moves.

On the other hand, robot linkages provide extremely strong cues of the sensor’s motion, due to their constrained kinematics<sup>3</sup>. As the robot actuates its joints, the sensor moves in a predictable manner. If the joint angles of the robot are perfectly known, and the robot is perfectly rigid, the pose of the sensor is known with absolute certainty through forward kinematics. In this context, dense 3D reconstruction is a simple mapping problem rather than a full SLAM problem.

However, even articulated robotic linkages may have significant noise and systematic inaccuracies<sup>4</sup>. The use of relative joint encoders leads to inaccuracy that compounds over time; and on some systems, joint encoders are separated from the actuation system by a mechanism with unpredictable dynamics. The *Baxter* (Fig. 2) robot, for instance, has series-elastic actuators and deformable plastic gears which significantly increase its kinematic uncertainty. Another robot arm, the *Barrett WAM*, is driven by cables. Stretch in the cables results in systematic error that is configuration and torque dependant. In the extreme case, there may be degrees of freedom that are entirely unknown (such as the pose of the robot’s base, or the extrinsic calibration of the sensor). One way of dealing with systematic error is to calibrate the joint encoders to account for unknown dynamics and other sources of systematic error; but if additional unmodeled dynamics are introduced (such as a grasped object, or exterior forces), the calibration will fail to capture them.

Uncertainty and unmodeled systematic error can be corrected simultaneously by fusing visual data from the camera with the robot’s joint encoder data. If the camera is mounted externally, so that all or part of the robot is visible, the problem becomes the well-studied *articulated tracking*<sup>5</sup> problem. In this case, the robot’s known

<sup>2</sup> See page 15

<sup>3</sup> See page 9

<sup>4</sup> See page 21



Figure 2: A camera mounted on the hand of a *Rethink Robotics* Baxter robot.

<sup>5</sup> See page 11

kinematic model can be used in conjunction with visual sensor data to directly estimate the robot's configuration. In contrast, this work is primarily concerned with *hand-mounted* sensors which, for the most part, cannot see the actuator on which they're mounted. In this case, the localization algorithm can't rely on a known model of the world. Without this information, is it still possible to estimate the robot's configuration?

For my thesis work, I propose a method of simultaneously estimating the joint angles of a robot and a dense 3D map of the world around it by solving a constrained visual SLAM problem<sup>6</sup>. Like in the traditional 3D SLAM problem, this method assumes no *a. priori* knowledge of the scene, and will have to take the sensor's pose uncertainty into account. Like the articulated tracking problem, it leverages the known kinematic structure of the robot to strongly inform the sensor's pose.

In some ways, articulated SLAM is easier than the direct 6 Degree of Freedom (DOF) SLAM problem, because the pose of the sensor is so tightly constrained by the robot's kinematics. However, since the inverse mapping between the sensor's pose and the robot's joint angles is nonlinear and underconstrained, translating 3D SLAM techniques to the articulated case will require work. I believe that the correct state space for this problem is the configuration space of the robot, not the 6 DOF pose of the sensor. Knowing the configuration of the robot allows us to issue commands for servoing and planning, and joint limits, self-collisions and other factors intrinsic to the geometry of the robot prevent the pose estimate from drifting too far from reality.

This proposal includes a brief overview of theory and related works<sup>7</sup>, an outline of potential methods of solving the articulated SLAM problem<sup>8</sup>, 2D and 3D simulations and experiments<sup>9</sup> involving robots with depth sensors, and an outline of future work<sup>10</sup>.

<sup>6</sup> See page [19](#)

<sup>7</sup> See page [9](#)

<sup>8</sup> See page [19](#)

<sup>9</sup> See page [27](#)

<sup>10</sup> See page [31](#)



# Background

Solving the articulated SLAM problem requires an understanding of the robot's kinematics, the mechanics of joint encoders, the characteristics of the sensor, and SLAM techniques.

## Robot Kinematics

A *kinematic linkage*<sup>11</sup> consists of a series of rigid bodies (called *links*) attached to one another through mechanisms (called *joints*) that constrain their motion. A joint usually has between 1 and 6 DOF which define how it constrains the motion of its attached links. For instance, a rotary joint has one degree of freedom that defines a pure rotation around an axis, while a ball joint has two to three degrees of freedom defining a pure rotation around 2 or more axes. The joint which constrains link *A* to link *B*, and which has configuration  $q_i$  has the transformation:

$$T_B^A(q_i) \in SE(3) \quad (1)$$

and as  $q_i$  changes, so does the transformation between links *A* and *B*.

A kinematic linkage can be represented as an undirected graph where vertexes are links, and the edges are joints between pairs of links. When the graph is acyclic, the linkage is said to be a *kinematic tree*. For instance, a robot with two arms, a head, and a fixed base has a kinematic tree with three branches: one for each arm, and one for the head. Any link may be treated arbitrarily as the root of the tree, though usually a *base link* that is rigidly attached to the world is defined as the root.

The transformation of any link  $L_i$  with respect to a fixed reference frame *W* can be calculated by traversing the kinematic tree and appending the transformations implied by each joint from the root of the tree (a process called *forward kinematics*):

$$T_{L_i}^W = T_{L_i}^{L_{i-1}}(q_{i-1}) \dots T_{L_2}^{L_1}(q_1) T_{L_1}^W \quad (2)$$

the path from the root to the link in question is called that link's

<sup>11</sup> Matthew T. Mason . *Mechanics of Robotic Manipulation*. MIT Press, Cambridge, MA, August 2001

*kinematic chain*. The end of the kinematic chain is called the *end effector*. This work will use the reference frame of the camera as the end of the arm's kinematic chain.

A robot's configuration  $\mathbf{q} \in \mathbb{R}^N$  is a vector which concatenates all of its joints' degrees of freedom:

$$\mathbf{q} = \begin{bmatrix} q_1 \\ \vdots \\ q_N \end{bmatrix} \quad (3)$$

The partial derivative of link  $i$ 's reference frame with respect to  $\mathbf{q}$ :

$$\mathbf{J}_i(\mathbf{q}) = \frac{\partial}{\partial \mathbf{q}} T_{L_i}^W \quad (4)$$

is called the link's *kinematic Jacobian*, and for simple kinematic chains it can be computed in closed-form efficiently. The Jacobian is an important object for the articulated SLAM problem, because it allows us to map instantaneous changes of the camera's pose to instantaneous changes in the robot's joints.

The robot can also control its joints. Control signals  $\mathbf{u} \in \mathbb{R}^N$  are sent to the joints, causing them to move based on the robot's physical dynamics <sup>(12)</sup>.

<sup>12</sup> See page [21](#)

## Depth Cameras

My thesis will be primarily concerned with two kinds of sensors: 2D (*i. e monocular*) cameras and depth cameras. In both cases, the output is a 2D image. Call the image  $I_D$ , it is a function with domain  $\Omega \in \mathbb{R}^2$ . The relationship between 3D points in the scene and 2D points on a camera image can be modeled using the simple pinhole camera intrinsic <sup>13</sup> model:

$$\text{Proj}(x, y, z) = \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} f_x x + c_x \\ f_y y + c_y \end{bmatrix} \quad (5)$$

<sup>13</sup> R. I. Hartley and A. Zisserman.  
*Multiple View Geometry in Computer Vision*. Cambridge University Press,  
ISBN: 0521540518, second edition, 2004

where  $u, v$  are the 2D images coordinates,  $x, y, z$  are the 3D point's coordinates in the camera's frame of reference (with  $x$  to the right,  $y$  down, and  $z$  forward), and  $f_x, f_y, c_x, c_y$  are the intrinsic parameters of the camera. The set of all points in 3D space which project onto the camera is called the *field of view* of the camera, and is usually approximated by a truncated pyramid (or *frustum*). I can also define the inverse projection model, which takes a camera coordinate  $u, v$  and depth measurement  $z$ , and converts it into a 3D vector relative to

the camera's focal point:

$$\text{Proj}^{-1}(u, v, z) = z \begin{bmatrix} \frac{u - c_x}{f_x} \\ \frac{v - c_y}{f_y} \\ 1 \end{bmatrix} \quad (6)$$

Depth cameras measure the range ( $z$ ) at each image pixel to points in the scene. Projective depth cameras accomplish this by projecting a pattern of infrared light onto the scene and measuring the disparity of this pattern using an offset infrared camera. Time-of-flight cameras measure depth by directly shooting rays of infrared light at the scene and measuring the time it takes for the rays to reflect back to the sensor. In any case, the depth measured by depth cameras is noisy, incomplete, and contains systematic error.

### *Articulated Tracking*

In practice, robots cannot directly measure their configuration, instead they must rely on sensors (called *joint encoders*) which indirectly measure the configuration of their joints. Between the robot's joint encoder and the actual degree of freedom of the joint, there may be mechanisms with unmodeled dynamics (such as non-rigid links, springs, gear trains, or pulley systems) which prevent the robot from knowing its own configuration with certainty. For example, each joint of the *Baxter* robot consists of a motor with plastic gearing, followed by a spring, followed by a rotating mechanism. In the extreme case (*e. g.* human pose tracking or hand tracking), some degrees of freedom may be completely unknown, and can only be inferred through external sensing.

Tracking articulated bodies with known kinematic structure using externally-mounted visual sensors is a well-studied topic in computer vision. For instance, commercial motion capture systems use markers (such as fiducials or reflective spheres) attached to articulated bodies along with an external camera to track the pose of human actors and objects in the scene.

When only the kinematic structure of the body is known, but no markers are available, the problem becomes difficult due to the unknown correspondences between sensor measurements and the body itself (*i. e* the *segmentation problem*). But even in this case, efficient solutions for tracking articulated bodies exist. In general, given a series of sensor measurements from timestep 1 to  $(t - 1)$ , and control inputs

$$Z_{1\dots t} = \{Z_1, \dots, Z_t\} \quad (7)$$

$$u_{1\dots t} = \{\mathbf{u}_1, \dots, \mathbf{u}_t\} \quad (8)$$

the task is to estimate the configuration of the body at time  $t$ :

$$\mathbf{q}_t = \operatorname{argmax}_{\mathbf{q}} P(\mathbf{q}|Z_{1\dots t}, \mathbf{q}_{1\dots(t-1)}, \mathbf{u}_{1\dots(t-1)}) \quad (9)$$

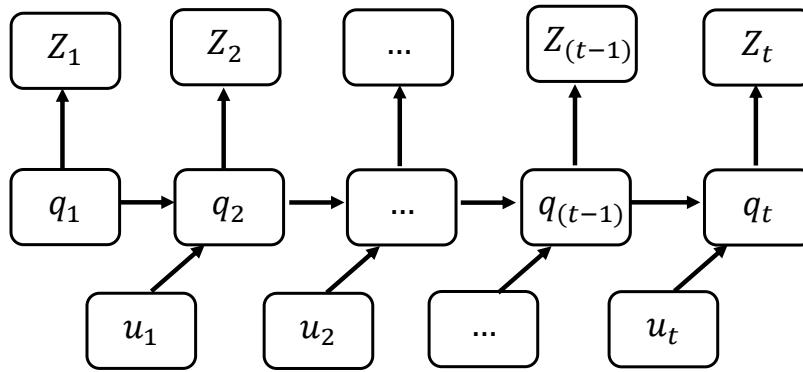
If the system is assumed to have the Markov property (*i.e.* the current state is conditionally independent of everything given the prior state), the articulated tracking problem reduces to merely finding the next configuration of the body given its previous configuration and sensor measurement:

$$\mathbf{q}_t = \operatorname{argmax}_{\mathbf{q}} P(\mathbf{q}|Z_{(t-1)}, \mathbf{q}_{(t-1)}, \mathbf{u}_{(t-1)}) \quad (10)$$

$$= \operatorname{argmax}_{\mathbf{q}} \frac{P(Z_{(t-1)}|\mathbf{q}, \mathbf{q}_{(t-1)}, \mathbf{u}_{(t-1)})P(\mathbf{q}|\mathbf{q}_{(t-1)}, \mathbf{u}_{(t-1)})}{P(Z_{(t-1)}, \mathbf{q}_{(t-1)}, \mathbf{u}_{(t-1)})} \quad (11)$$

$$= \operatorname{argmax}_{\mathbf{q}} \left[ \log P(Z_{(t-1)}|\mathbf{q}, \mathbf{q}_{(t-1)}, \mathbf{u}_{(t-1)}) + \log P(\mathbf{q}|\mathbf{q}_{(t-1)}, \mathbf{u}_{(t-1)}) \right] \quad (12)$$

the term  $P(Z_{(t-1)}|\mathbf{q}, \mathbf{q}_{(t-1)}, \mathbf{u}_{(t-1)})$  represents the posterior probability that a sensor measurement was observed given some body configuration, while the term  $P(\mathbf{q}|\mathbf{q}_{(t-1)}, \mathbf{u}_{(t-1)})$  represents the prior probability of the body's configuration given its previous configuration (called the *motion model*). Articulated tracking approaches differ mainly in how they compute these two quantities, as well as how they maximize them.



For instance, one method intended for tracking humans in 2D images, Articulated Iterative Closest Point (ICP)<sup>14</sup>, computes the pos-

Figure 3: A graphical model of the articulated tracking problem. Notice that the graph has the Markov property.

<sup>14</sup> L. Mundermann, S. Corazza, and T.P. Andriacchi. Accurately measuring human movement using articulated icp with soft-joint constraints and a repository of articulated models. In *Computer Vision and Pattern Recognition, 2007. CVPR '07. IEEE Conference on*, pages 1–6, June 2007. DOI: 10.1109/CVPR.2007.383302

terior using image-based edge features, and maximizes the likelihood of a configuration by coordinate descent in configuration space.

One of my earlier works, Real-time Markerless Articulated Tracking (RMAT)<sup>15</sup> tracks robots using a depth camera, a simple posterior model that sensor measurements are near the robot's surface as projected onto the depth image, and a motion model which assumes the robot's configuration changes slowly between timesteps. RMAT matches sensor measurements to corresponding points on the robot's body using an octree, and gradient descent is used to maximize the likelihood of the robot's configuration given its previous configuration and depth image. RMAT (though simple) is fast enough to run in a control loop; using articulated tracking, the robot is able to perform delicate manipulation tasks even with severe joint encoder error.

A related but separate work, Dense Articulated Real-time Tracking (DART)<sup>16</sup>, improves upon RMAT by using a signed distance field representation of the robot's body rather than an octree, a more complex motion model that considers joint velocity, and an extended Kalman filter for tracking. DART has been shown to effectively track robots, human hands, and rigid bodies in real-time with commercial depth cameras.

The articulated tracking problem is a subset of the problem my thesis will address. The requirement that the robot's sensor be able to see its own links greatly reduces the complexity of the problem, turning it into a tracking problem rather than a SLAM problem – but the same techniques used to compute the robot's configuration from sensor measurements can be used even when the scene is unknown. For one thing, the insight that tracking ought to be done in the configuration space of the robot, rather than in  $SE(3)$ , is critical to my thesis topic.

## *Simultaneous Localization and Mapping*

The Simultaneous Localization and Mapping (SLAM) problem<sup>17</sup> involves concurrently estimating the pose of a moving sensor and the structure of the world around it when both are unknown. That is, given a sequence of noisy sensor readings

$$Z_{1 \dots t} = \{Z_1, Z_2, Z_3, \dots, Z_t\} \quad (13)$$

estimate the pose of the sensor at all times,  $H_{1 \dots t} \in SE(3)$  and the scene  $M$ :

$$H_{1 \dots t}, M = \underset{H, M}{\operatorname{argmax}} P(H, M | Z_{1 \dots t}) \quad (14)$$

<sup>15</sup> Matthew Klingensmith, Thomas Galluzzo, Christopher Dellin, Moslem Kazemi, J Andrew Bagnell, and Nancy Pollard. Closed-loop Servoing using Real-time Markerless Arm Tracking. In *International Conference on Robotics And Automation (Humanoids Workshop)*, May 2013

<sup>16</sup> Tanner Schmidt, Richard Newcombe, and Dieter Fox. DART: Dense Articulated Real-Time Tracking. In *Robotics: Science and Systems*, number 1, 2014. URL <http://www.roboticsproceedings.org/rss10/p30.pdf>

<sup>17</sup> Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic Robotics (Intelligent Robotics and Autonomous Agents)*. The MIT Press, 2005. ISBN 0262201623

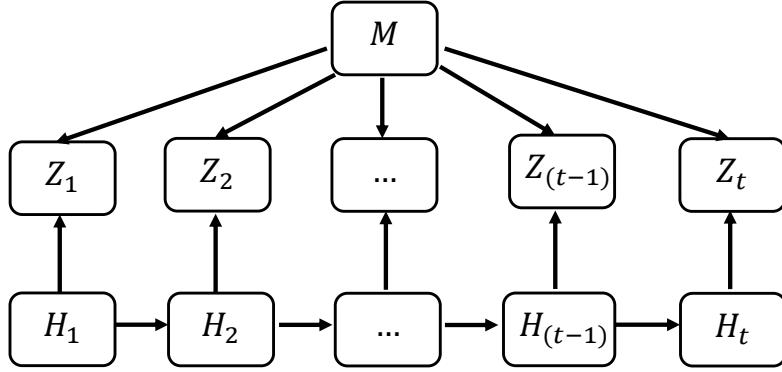


Figure 4: A graphical model of the SLAM problem. Note that given a known map, it has the Markov property, and is just the tracking problem. Without a known map, all measurements are correlated.

### Sparse SLAM

Depending on the scene representation and sensor type, different approaches to solving this problem are appropriate. If the measurement model of the sensor consists of a sparse set of features (such as SIFT features, corner features, or related), a *sparse* landmark-based scene representation is most appropriate. Sparse SLAM techniques build a globally-consistent graph of tracked landmarks (called a *pose graph*) with pairwise edges on sensor motion between frames; and then optimize this pose graph based on a consistency metric. Since only a very small subset of the sensor data is used, this method can be made to work efficiently over very long datasets featuring loop closures. Specifically, the typical sparse SLAM pose graph has the form:

$$M = \{V, E\} \quad (15)$$

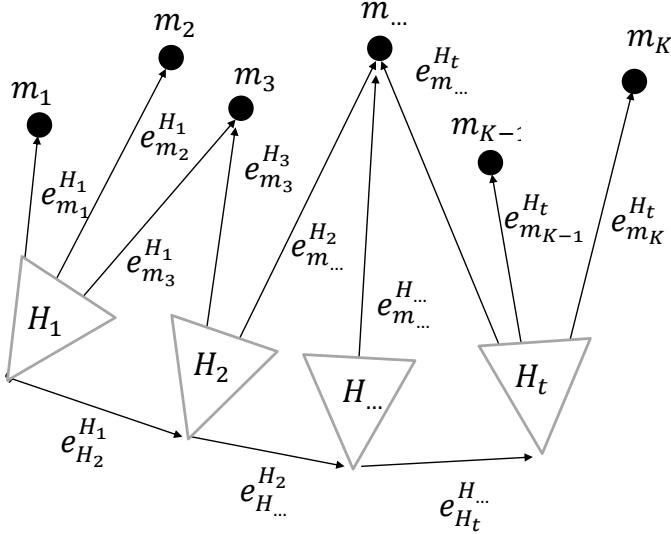
where  $V = \{v_1, \dots, v_N\}$  is a collection of vertices, with  $v_i \in SE(3)$  as either a landmark pose or a sensor pose, and  $E = \{e_1, \dots, e_M\}$  is a collection of directed edges, with each edge defining a transformation  $e_i = T_{v_a}^{v_b}$  from one vertex  $v_a$  to another vertex  $v_b$ . A pose graph is said to be *consistent* whenever each edge is consistent (*i.e.*, whenever  $v_a e_i = v_b$ ). A *loop closure* happens whenever two edges share the same outgoing vertex.

This breaks up the graphical model (Fig. 4 ) so that sensor measurements are only corelated by shared landmarks, rather than all sensor measurements being corelated.

More generally, a pose graph can contain arbitrary nodes as states, and arbitrary links between nodes which represent factors in a factor graph<sup>18</sup>. In my thesis work, the nodes in the pose graph will be robot configurations, links between nodes will come from an observed motion derived from the sensor, as well as additional factors from the

<sup>18</sup> M. Kaess, A. Ranganathan, and F. Dellaert. iSAM: Incremental smoothing and mapping. *IEEE Trans. on Robotics (TRO)*, 24(6):1365–1378, December 2008

Figure 5: An example of a pose graph.



robot’s joint encoders.

### Dense SLAM

However, our use case calls for a robot to localize itself only in a very small workspace, and to use the reconstructed scene for geometric reasoning about objects for the purpose of manipulating them. For this use case, it is much more important to have a detailed geometric representation of the scene than it is to have consistent pose over very long distances.

Dense SLAM techniques, in contrast to sparse SLAM techniques, use all or most of the sensor data to reconstruct a geometric model of the scene. This technique ensures that at all timesteps, a plausible geometric model taking into account all of the sensor data exists – but consequently, the memory required to store the map at all timesteps is much larger than in landmark-based SLAM. Some dense methods (such as *point fusion*<sup>19</sup>) still manage to store a globally consistent pose graph as well as a geometric model by sacrificing some model quality, while most others (such as *kinect fusion*<sup>20</sup> and variants) store the map and pose only for time  $t - 1$ , and then maximize:

$$M_{t-1} = \underset{M}{\operatorname{argmax}} P(M|Z_{t-1}, H_{t-1}) \quad (16)$$

$$H_t = \underset{H}{\operatorname{argmax}} P(H|H_{t-1}, M_{t-1}) \quad (17)$$

before generating the current model  $M_t$  from the previous model

<sup>19</sup> M. Keller, D. Lefloch, M. Lambers, S. Izadi, T. Weyrich, and A. Kolb. Real-time 3d reconstruction in dynamic scenes using point-based fusion. In *3D Vision - 3DV 2013, 2013 International Conference on*, pages 1–8, June 2013. DOI: 10.1109/3DV.2013.9

<sup>20</sup> Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, and Andrew Fitzgibbon. Kinectfusion: Real-time 3d reconstruction and interaction using a moving depth camera. In *UIST 2011*, 2011

$M_{t-1}$  and current pose  $H_t$ , and so on. By only storing a single map, these dense techniques avoid having to store dense models for every frame. This simplification ignores the fact that the unknown map correlates all of the sensor readings together (Fig. 4), but results in systemic errors in pose/map estimation building up over time.

A dense map representation of particular interest to our problem is the Signed Distance Field (SDF)<sup>21</sup>. The SDF stores a function that gives the distance to the nearest surface in the scene:

$$\Phi(\mathbf{x}) = \min_{\mathbf{s} \in S} \|\mathbf{x} - \mathbf{s}\| I(\mathbf{s}) \quad (18)$$

where  $S \subseteq \mathbb{R}^3$  is the set of all occupied points in the scene, and

$$I(\mathbf{s}) = \begin{cases} -1 & \text{if } \mathbf{s} \text{ is inside an object} \\ +1 & \text{otherwise} \end{cases} \quad (19)$$

$\Phi$  is negative inside objects, positive outside objects, and zero along the surface. The gradient  $\nabla \Phi$  can be easily calculated, and the representation very flexible, allowing  $\Phi$  to represent many kinds of scenes. On the other hand, storing the SDF requires memory on the order of the volume of the scene.

In practice, the SDF is computed only for a small region around the surfaces of objects (called the *truncation region*) for which the locally-linear approximation is good. The result is the Truncated Signed Distance Field (TSDF)<sup>22</sup> (Alg. 1). TSDF-based dense SLAM techniques (such as *kinect fusion* and variants) build  $\Phi$  at time  $t$  by fusing the depth image produced by the camera at time  $t - 1$ . A depth image can be fused into an SDF by computing a locally-linear approximation of the SDF for the depth image, and averaging it into the SDF computed at  $t - 1$ . Since only one dense model is maintained at each timestep, errors in the model can compound over time.

One of my earlier works, *CHISEL*<sup>23</sup>, explores techniques for efficiently building and storing the SDF used in dense SLAM, and extends the technique to mobile phones. My thesis will be using *CHISEL* as a backend dense mapping technique.

### Direct Monocular SLAM

For monocular cameras, techniques exist which combine the benefits of sparse and dense SLAM. These techniques (called direct, or direct semi-dense techniques) track the pose of a moving camera by directly minimizing the *photometric error* between subsequent camera frames. This involves directly estimating the depth of pixels in the 2D camera image and reprojecting pixels from one frame onto another, and then

<sup>21</sup> Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '96, 1996

<sup>22</sup> Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '96, 1996

<sup>23</sup> Matthew Klingensmith, Ivan Dryanovski, Siddhartha Srinivasa, and Jizhong Xiao. Chisel: Real time large scale 3d reconstruction onboard a mobile device. In *RSS*, 2015

---

**Algorithm 1:** FUSETSDF
 

---

```

1 // Given a depth image, sensor pose, previous SDF, previous
   voxel weights, a weighting function, a volume of interest, and
   a truncation distance
2 Input:  $I_D, H_t, \Phi_{t-1}, W_{t-1}, w, V, \tau$ 
3  $H' \leftarrow H_t^{-1}$ 
4 // Initialize the new distance field and voxel weights
5  $\Phi_t \leftarrow \Phi_{t-1}$ 
6  $W_t \leftarrow W_{t-1}$ 
7 // For each voxel center in a volume of interest
8 for  $\mathbf{v} \in V$  do
9   // Transform the voxel to the camera frame and project it on
     the depth image.
10   $\mathbf{v}_h \leftarrow H'\mathbf{v}$ 
11   $d_i \leftarrow I_D [\text{Proj}(\mathbf{v}_h)]$ 
12  // The geometric depth of the voxel to the camera plane.
13   $d_v \leftarrow \mathbf{v}_h(z)$ 
14  // The difference between the geometric depth and the
     reported depth is a local approximation of the SDF
15   $u \leftarrow d_v - d_i$ 
16  // If the local approximation is within the truncation
     distance...
17  if  $|u| < \tau$  then
18    // Compute the weighted average with the previous SDF.
19     $\Phi_t(\mathbf{v}) \leftarrow \frac{W_t(\mathbf{v})\Phi_t(\mathbf{v}) + w(u)u}{W_t(\mathbf{v}) + w(u)}$ 
20     $W_t(\mathbf{v}) \leftarrow W_t(\mathbf{v}) + w(u)$ 

```

---

**Output:**  $\Phi_t, W_t$ 


---

maximizing the likelihood of the subsequent camera poses:

$$H_t = \underset{H}{\operatorname{argmin}} \| I_{t-1}[Z_{t-1}] - I_t[\operatorname{Proj}(HZ_{t-1})] \| \quad (20)$$

where  $I$  is the camera image at a given time, and  $Z_t$  is the point cloud as inferred by stereo matching between pose  $H_t$  and  $H_{t-1}$ . Examples of direct techniques include ORB-SLAM<sup>24</sup>, LSD-SLAM<sup>25</sup>, and DTAM<sup>26</sup>.

Direct monocular methods, like sparse keypoint methods, are scale ambiguous, and cannot distinguish between two scenes that differ only in scale. Like sparse keypoint methods, they generally build a pose graph on top of the dense or semi-dense localization system.

Direct methods generally use Gauss-Newton nonlinear optimization to solve for Eq. 20. They apply a small motion  $\Delta\xi$  to the camera's pose by solving:

$$\Delta\xi = (J_I^T J_I)^{-1} J_I r_I(\xi) \quad (21)$$

where  $r_I$  is the *image residual* (the inner term of (20)), and  $J_I = \frac{\partial r_I}{\partial \Delta\xi}$  is the image Jacobian. The image Jacobian is an interesting term, as it captures the incremental change in the pixel intensities of an image with respect to an incremental motion.

Unlike sparse keyframe-based algorithms, semi-dense direct techniques are capable of creating high-quality, dense 3D reconstructions. Unlike purely dense algorithms, they do not require a depth image to work. These qualities make semi-dense techniques attractive for my thesis.

<sup>24</sup> Raúl Mur-Artal and Juan Tardós. Probabilistic semi-dense mapping from highly accurate feature-based monocular slam. In *RSS*, 2015

<sup>25</sup> J. Engel, T. Schöps, and D. Cremers. LSD-SLAM: Large-scale direct monocular SLAM. In *ECCV*, 2014

<sup>26</sup> Richard A. Newcombe, Steven J. Lovegrove, and Andrew J. Davison. Dtam: Dense tracking and mapping in real-time. In *ICCV*, 2011

# Framework

## Articulated SLAM

When the sensor mounted to the robot is outward-facing, and nothing is known about the world around the robot, how can the robot simultaneously estimate its configuration, and reconstruct the 3D geometry of the world? This is equivalent to the articulated tracking problem with an unknown world model, and also equivalent to the visual SLAM problem with additional constraints on the pose of the sensor.

Given a series of sensor measurements  $Z_1, \dots, Z_t$  and control inputs  $u_1, \dots, u_t$ , the algorithm must simultaneously estimate the trajectory of the robot's configuration  $q_1, \dots, q_t$ , and a dense model of the world  $M$ :

$$q_1, \dots, q_t, M = \underset{q_{\dots}, M}{\operatorname{argmax}} P(q_{\dots}, M | Z_{1\dots t}, u_{1\dots t}) \quad (22)$$

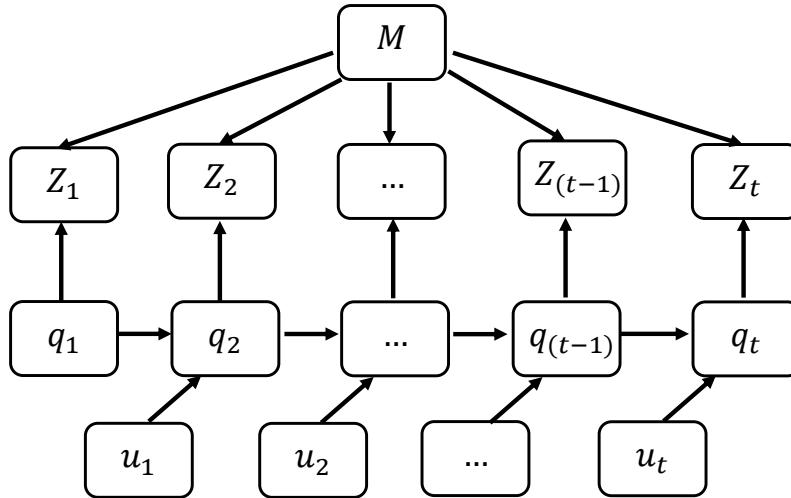


Figure 6: A graphical model of the articulated SLAM problem. It combines aspects of the articulated tracking problem and the SLAM problem..

This problem is made difficult by the fact that the unknown map

$M$  correlates all of the sensor measurements together (Fig. 6). First, consider the articulated tracking problem with a *known* map. In this case, Eq. 22 reduces to

$$\mathbf{q}_1, \dots, \mathbf{q}_t = \operatorname{argmax}_{\mathbf{q}_{\dots}} P(\mathbf{q}_{\dots} | M, Z_{1\dots t}, \mathbf{u}_{1\dots t}) \quad (23)$$

This becomes a state estimation problem rather than a full SLAM problem, and therefore acquires the Markov property. So it suffices to estimate the current state given the previous state estimate only:

$$\mathbf{q}_t = \operatorname{argmax}_{\mathbf{q}} P(\mathbf{q} | \mathbf{q}_{t-1}, \mathbf{u}_{t-1}, M, Z_{1\dots t}) \quad (24)$$

$$= \operatorname{argmax}_{\mathbf{q}} \frac{P(Z_{1\dots t} | \mathbf{q}, \mathbf{q}_{t-1}, \mathbf{u}_{t-1}, M) P(\mathbf{q} | \mathbf{q}_{t-1}, \mathbf{u}_{t-1}, M)}{P(Z_{1\dots t}, \mathbf{q}_{t-1}, \mathbf{u}_{t-1}, M)} \quad (25)$$

$$= \operatorname{argmax}_{\mathbf{q}} \log P(Z_{1\dots t} | \mathbf{q}, \mathbf{q}_{t-1}, \mathbf{u}_{t-1}, M) + \log P(\mathbf{q} | \mathbf{q}_{t-1}, \mathbf{u}_{t-1}, M) \quad (26)$$

$$= \operatorname{argmax}_{\mathbf{q}} \log P(Z_t | \mathbf{q}, M) + \log P(\mathbf{q} | \mathbf{q}_{t-1}, \mathbf{u}_{t-1}) \quad (27)$$

yielding two terms which must be maximized: a posterior term  $P(Z_t | \mathbf{q}, M)$ , representing the likelihood of the sensor measurements given a robot configuration and known map, and a prior term  $P(\mathbf{q} | \mathbf{q}_{t-1}, \mathbf{u}_{t-1})$ , representing the likelihood of a robot configuration given its past configuration and control input.

### Prior Motion Model

The prior  $P(\mathbf{q} | \mathbf{q}_{t-1}, \mathbf{u}_{t-1})$  depends on the robot's dynamics, and its relationship to gravity and other external forces. In general, the forward dynamics model of a robot arm can be expressed using the well-known differential equation based on its Lagrangian:

$$\tau(\mathbf{u}) = \mathbf{M}(\mathbf{q})\ddot{\mathbf{q}} + \mathbf{C}(\mathbf{q}, \dot{\mathbf{q}})\dot{\mathbf{q}} + \mathbf{G}(\mathbf{q}, \dot{\mathbf{q}}, F_{\text{ext}}) \quad (28)$$

where  $\tau$  is the torque the robot produces at its joints given  $\mathbf{q}$ ,  $\mathbf{M}$  is a function representing the robot's mass projected into joint space,  $\mathbf{C}$  is a function representing the internal coriolis forces that the robot experiences, and  $\mathbf{G}$  is a function representing how the robot responds to external forces  $F_{\text{ext}}$  (usually gravity).

If this dynamics model were known perfectly, then the prior motion model would be completely deterministic, and indeed, no external sensors of any kind (not even joint encoders) would be needed to track its state over time. Unfortunately, the dynamics model of a robot manipulator is not easy to calculate, and depends

on having a full account of external and internal forces the robot experiences.

These unmodeled dynamics can be mitigated in two ways: first, by treating them like noise in the motion model, and second, by using external sensors to account for them. At the very least, the dynamic model of the robot provides a narrow region of configuration space that the robot's configuration can be in, given its previous configuration and control input.

### *Motor Encoder Posterior*

Robots invariably have sensors (called *encoders*) mounted on their motors and/or joints. These sensors report a value  $\mathbf{q}^{(e)}$  that correlates to the actual geometric joint angle  $\mathbf{q}^{(j)}$ , but, due to the complications of an unknown intervening mechanism (Fig. 7), the encoder's reported value may have noise and systematic error (Fig. 8). For instance, the Barrett WAM robot arm is cable-driven. Its' encoders are attached to its' motors, and a complicated system of cables and pulleys transmits torque from the motor to the joint. Unmodeled dynamics from stretch in the cable system leads to a discrepancy between  $\mathbf{q}^{(e)}$  and  $\mathbf{q}^{(j)}$ .

Define a simple probabilistic model:

$$\mathbf{q}^{(j)} = g(\mathbf{q}^{(e)}) + \epsilon(\mathbf{q}^{(e)}) \quad (29)$$

where  $g : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is a static calibration function that maps motor encoders to joint angles, and  $\epsilon(\mathbf{q}^{(e)})$  is a random variable capturing the variance of the mapping from the motor encoders to joint angles (Fig. 8). This model may be learned by Gaussian Process Regression, or other regression techniques, when ground truth data is known.

Initial tests I have carried out prior to this proposal on a Barret WAM arm using optical joint encoders as ground truth have revealed that the mapping (Eq. 29) between the motor encoders and joint angles depends not only on the configuration of the robot, but also on the applied torque, the path that the robot has taken to get to a particular configuration, and other variables that are not easy to capture using a simple Gaussian model. More work is needed to correctly model joint encoder error.

### *Depth Camera Posterior*

The sensor posterior term,  $P(Z_t | \mathbf{q}, M)$ , represents the likelihood of sensor measurements given a particular robot configuration  $\mathbf{q}$  and dense map  $M$ . Here, the relevant work from dense 3D SLAM can

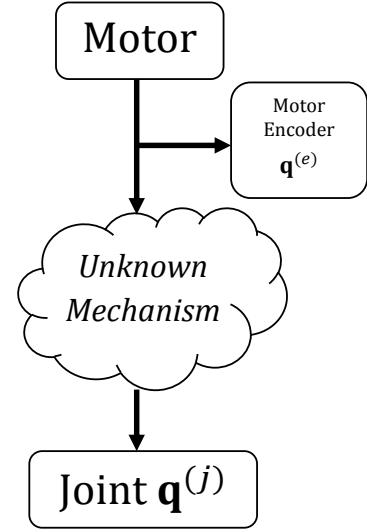


Figure 7: There is always some unknown mechanism between the method of actuation (a motor), and the physical robot degree of freedom.

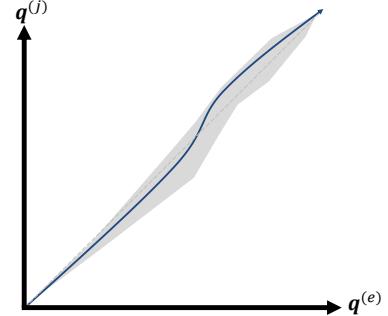


Figure 8: A sketch of a hypothetical relationship between the motor encoders  $\mathbf{q}^{(e)}$  and the physical joint angles  $\mathbf{q}^{(j)}$ , note that the relationship is multivariate, even though it is shown as one-dimensional here.

be taken into account. Consider a depth camera mounted on the robot's end effector which can see portions of the dense map. The depth data can be modeled as cones, rays, or points. Since each depth pixel measurement is independent given the map and robot configuration, the likelihood of each ray, cone, or point, can be calculated independently.

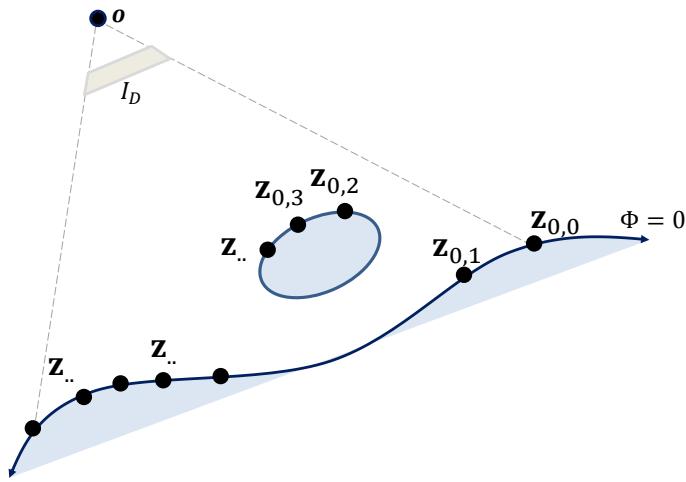


Figure 9: A depth camera with focal point  $\mathbf{o}$  and depth image  $I_D$  projects points  $\mathbf{z}_{u,v}$  into the scene. All the points are expected to lie along surfaces where the signed distance to surfaces in the scene  $\Phi = 0$ .

**Point Cloud Model** One of the simplest means of modeling the depth camera is to unproject each depth pixel into the scene as a point using the camera intrinsics. The result is the unstructured *point cloud* of the depth image (Fig. 9). So the sensor measurement is:

$$Z_t = \{H_t \text{ Proj}^{-1}(I_D[u, v]) = \mathbf{z}_{u,v} \forall [u, v] \in \Omega\} \quad (30)$$

and due to the conditional independence of the depth image given a robot configuration and map, the posterior is:

$$P(Z_t | \mathbf{q}, M) = \prod_{u,v \in \Omega} P(\mathbf{z}_{u,v} | \mathbf{q}, M) \quad (31)$$

Assume that the points in the point cloud are corrupted by simple isotropic Gaussian noise. If this is the case,

$$P(\mathbf{z}_{u,v} | \mathbf{q}, M) \propto \exp(-\Phi(\mathbf{z}_{u,v})^2) \quad (32)$$

where  $\Phi$  is the signed distance field representation of the scene  $M$ .

This model, while simple, ignores occlusions. It makes no distinction between points behind obstacles and points in front of obstacles

from the perspective of the camera. It also ignores the anisotropic (direction-dependant) nature of depth sensor noise. A more complicated model may be needed to correctly capture the depth camera posterior. The generalized ICP<sup>27</sup> algorithm directly makes use of this model to align ray clouds.

<sup>27</sup> A. Segal, D. Haehnel, and S. Thrun. Generalized-icp. In *Proceedings of Robotics: Science and Systems*, Seattle, USA, June 2009

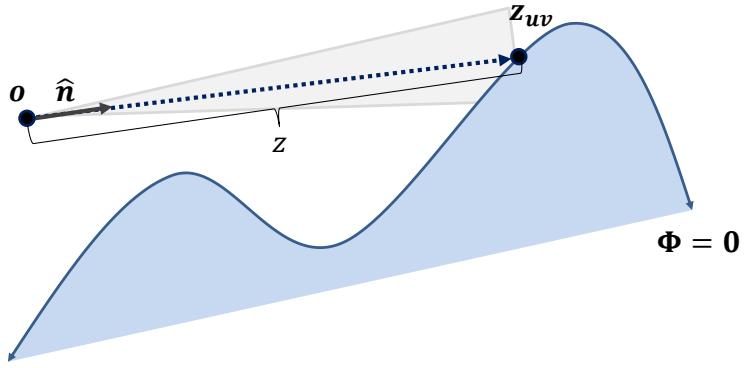


Figure 10: A cone/ray model of the depth camera. The cone emanates from the sensor's focal point  $\mathbf{o}$ , with direction  $\hat{\mathbf{n}}$ . It hits the surface of the scene at point  $\mathbf{z}_{uv}$ , which is distance  $z$  along its center.

**Ray/Cone Cloud Model** The depth camera can be modeled as a set of cones emanating from the camera's focal point – one cone for each depth pixel. The cone is the swept volume of the depth pixel projected into the scene, and has a length equal to the depth pixel's measurement (Fig. 10). The cone may be seen as the union of all rays that pass through the depth pixel from the focal point of the camera.

This is a much more complex model, incorporating information about which points in the scene have been passed through by rays as well as which points lie on surfaces (as well as their surface normal).

### Dense Sensor-To-Map Constrained SLAM

Given a model of the sensor posterior and motion model, and a dense signed distance field representation of the map, it is possible to formulate a simple dense<sup>28</sup> sensor-to-map algorithm for articulated SLAM. Using the point cloud posterior model<sup>29</sup>, and a simple motion model which just assumes that the difference between the robot's configuration and the reported motor encoder readings<sup>30</sup>  $\|\mathbf{q}_t - \mathbf{q}_t^e\|$  is small, the problem can be treated as a simple optimization

<sup>28</sup> See page 15

<sup>29</sup> See page 22

<sup>30</sup> See page 20

problem at each timestep.

$$\mathbf{q}_t \leftarrow \operatorname{argmax}_{\mathbf{q}} P(\mathbf{q}|Z_t, M_{t-1}, \mathbf{q}_{t-1}) \quad (33)$$

$$M_t \leftarrow \operatorname{argmax}_M P(M|\mathbf{q}_t, Z_t, M_{t-1}) \quad (34)$$

The configuration update (Eq. 34) can be maximized by breaking it into its component parts

$$\operatorname{argmax}_{\mathbf{q}} P(\mathbf{q}|Z_t, M_{t-1}, \mathbf{q}_{t-1}) = \operatorname{argmax}_{\mathbf{q}} \log P(Z_t|\mathbf{q}, M_{t-1}) + \log P(\mathbf{q}|\mathbf{q}_{t-1}) \quad (35)$$

and the sensor posterior from the point cloud model can be computed as:

$$\log P(Z_t|\mathbf{q}, M_{t-1}) = \log \prod_{[u,v] \in \Omega} P(\mathbf{z}_{u,v}|\mathbf{q}, M_{t-1}) \quad (36)$$

$$\propto \sum_{[u,v] \in \Omega} \log \exp(-\Phi(\mathbf{z}_{u,v})^2) \quad (37)$$

$$= \sum_{[u,v] \in \Omega} -\Phi(\mathbf{z}_{u,v})^2 \quad (38)$$

where  $\mathbf{z}_{u,v} = H_t \operatorname{Proj}^{-1}(I_D[u, v])$  is a point from the point cloud of the depth image at time  $t$ , and  $\Phi$  is the signed distance field of the map at time  $t-1$ .

So, maximizing the sensor posterior is equivalent to minimizing the simple cost function:

$$c(\mathbf{q}) = \frac{1}{2} \sum_{[u,v] \in \Omega} \Phi(\mathbf{z}_{u,v})^2 \quad (39)$$

and  $c$  is minimized whenever  $\frac{\partial c}{\partial \mathbf{q}} = 0$ . That is:

$$\nabla c = \frac{1}{2} \sum_{[u,v] \in \Omega} \left[ \frac{\partial}{\partial \mathbf{q}} \Phi(\mathbf{z}_{u,v})^2 \right] \quad (40)$$

$$= \sum_{[u,v] \in \Omega} \left[ \Phi(\mathbf{z}_{u,v}) \frac{\partial}{\partial \mathbf{q}} \nabla \Phi(\mathbf{z}_{u,v}) \right] \quad (41)$$

$$= \sum_{[u,v] \in \Omega} \left[ \Phi(\mathbf{z}_{u,v}) \mathbf{J}(\mathbf{q}, \mathbf{z}_{u,v})^T \nabla \Phi(\mathbf{z}_{u,v}) \right] \quad (42)$$

where  $\mathbf{J}(\mathbf{q}, \mathbf{z}_{u,v}) \in \mathbb{R}^{3 \times N}$  is the manipulator's kinematic jacobian<sup>31</sup>

<sup>31</sup> See page 9

evaluated for the point  $\mathbf{z}_{u,v}$  as rigidly attached to the end effector.

The gradient of the SDF  $\nabla \Phi$  is easy to approximate using finite differencing.

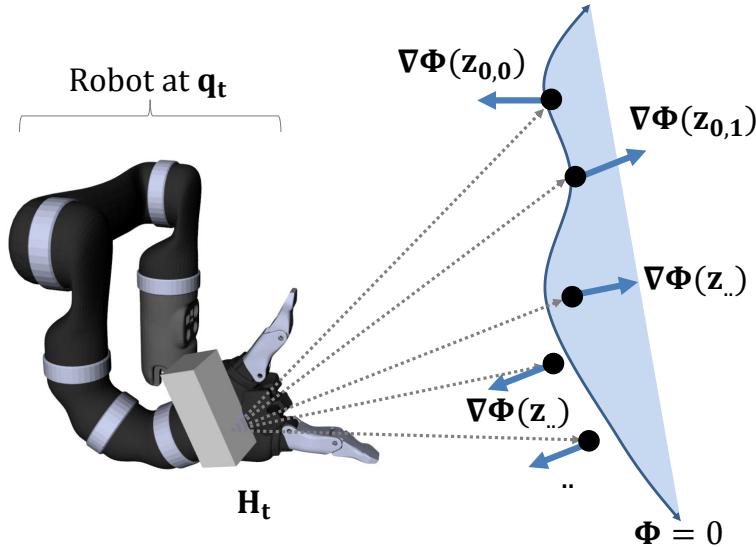


Figure 11: A rendering of a robot manipulator with attached depth sensor. Points  $z_{u,v}$  are projected into the scene, which is compared with the map's distance field from the previous timestep  $\Phi$ . The gradient of  $\Phi$  at each point in the point cloud is shown as an arrow. The end effector transform at time  $t$  is given by  $H_t$ , the robot's configuration at that time is  $q_t$ .

This cost function and gradient has an interesting physical interpretation (Fig. 11). If all of the points in the point cloud are interpreted as being rigidly attached to the camera frame, and a force is applied at each of the points  $z_{u,v}$  with a force in the direction  $\nabla\Phi(z_{u,v})$  with magnitude  $\Phi(z_{u,v})$ , the gradient  $\nabla c$  is the same as the instantaneous acceleration the joints would experience from the sum of such forces.

The likelihood of  $q_t$  can be maximized using gradient descent of the cost function  $c$  regularized by the motor encoders (Alg. 2). In this way, the robot's configuration stays near what is reported by its motor encoders, but is allowed to change in response to the sensor data.

---

**Algorithm 2:** Regularized Dense SDF Gradient Descent.

---

```

1 // Where  $\mathbf{q}_t^{(e)}$  are the motor encoders at time  $t$ ,  $\lambda$  is a learning
   rate, and  $\gamma$  is a regularization parameter.
Input:  $Z_t, \mathbf{q}_t^{(e)}, \mathbf{q}_{t-1}, \Phi_{t-1}, \lambda, \gamma$ 
2 // Initialize configuration offset  $\epsilon \in \mathbb{R}^N$  to zero.
3  $\epsilon \leftarrow \mathbf{0}$ 
4 repeat
5   // Current estimate for  $\mathbf{q}_t$  is the joint encoders plus the offset
6    $\mathbf{q}_t \leftarrow \mathbf{q}_t^{(e)} + \epsilon$ 
7   // Compute the gradient of the sensor measurement
      posterior
8    $\nabla c(\mathbf{q}_t) \leftarrow \sum_{[u,v] \in \Omega} [\Phi(\mathbf{z}_{u,v}) \mathbf{J}(\mathbf{q}_t, \mathbf{z}_{u,v})^T \nabla \Phi(\mathbf{z}_{u,v})]$ 
9   // Descend the gradient regularized by  $\gamma$ 
10   $\epsilon \leftarrow \epsilon - \lambda \nabla c(\mathbf{q}_t) - \gamma \epsilon$ 
11 until convergence;
12 // Now that there is a new estimate of  $\mathbf{q}_t$ , fuze the sensor data
    into the SDF.
13  $\Phi_t \leftarrow \text{FUSETSDF}(Z_t, F(\mathbf{q}_t), \Phi_{t-1}, \dots)$ 
Output:  $\mathbf{q}_t, \Phi_t$ 

```

---

# *Experiments*

## *2D Simulation Experiments*

I constructed a simple 2D simulation to demonstrate the constrained descent algorithm (Alg. 2) (Table 13). In the simulation, a 3-link serial robot manipulator with a simulated 1D depth sensor scans a scene. I added zero-centered Perlin<sup>32</sup> noise to its joint encoder readings. That is,

$$\mathbf{q}_t^{(e)} = \beta_n \text{PERLIN}(s_n \mathbf{q}_t^{(j)}) \quad (43)$$

where  $s_n$  is a noise scale parameter, and  $\beta_n$  is a noise magnitude parameter.

For the world model, I constructed a simple 2D TSDF. The performance of the constrained gradient descent algorithm (Alg. 2) is compared against a simple unconstrained descent algorithm which assumes the sensor can move and rotate freely, without considering the robot kinematics.

The robot scans the scene, and makes a single loop closure. When using the (noisy) joint encoders for pose only, much of the detail of the scene is lost, since the TSDF ( $\Phi$ ) smooths out errors through averaging. Constrained gradient descent allows the robot to correct its joint encoder error and recapture fine details in the scene (though drift is still made apparent by the loop closure). Unconstrained gradient descent, on the other hand, causes the robot's estimate of the sensor pose to drift; and the 3D reconstruction is quickly destroyed as error accumulates (Fig. 12).

<sup>32</sup> Ken Perlin. Improving noise. In *SIGGRAPH*, 2002

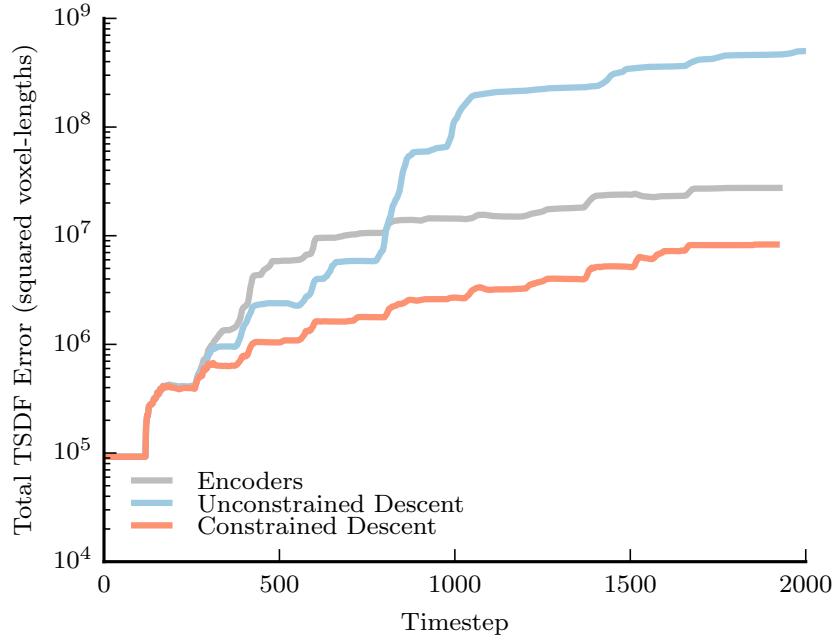
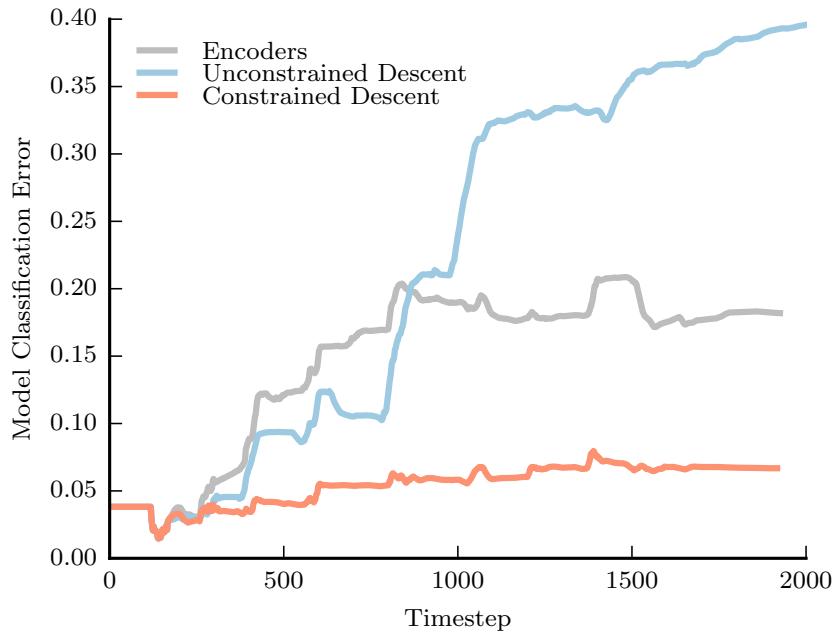
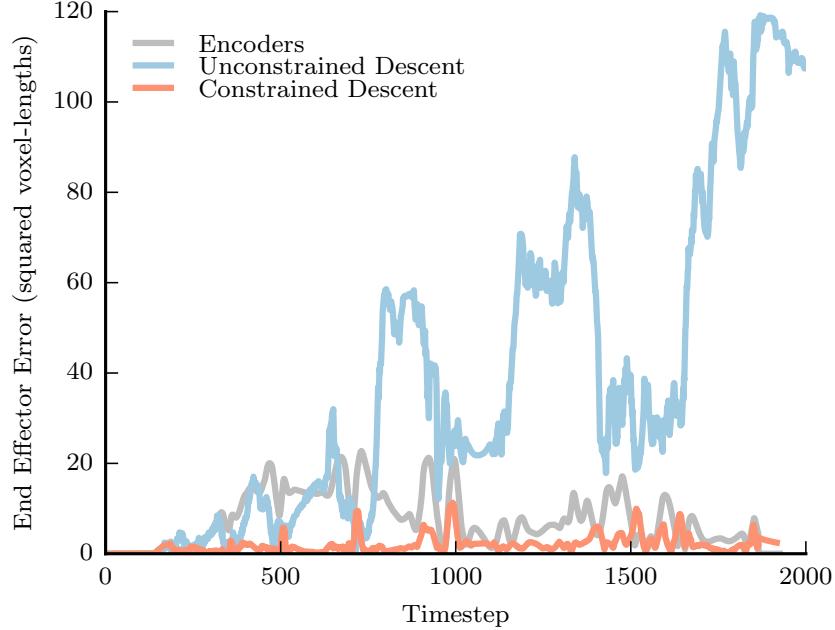


Figure 12: 2D experiment (Table 13) data. Here, end effector error is the squared distance between the reported end effector position and ground truth. The model classification error is the proportion of TSDFvoxels that are misclassified (full vs. empty), and the total TSDFerror is the total squared difference between the ground truth distances and the reported distances.

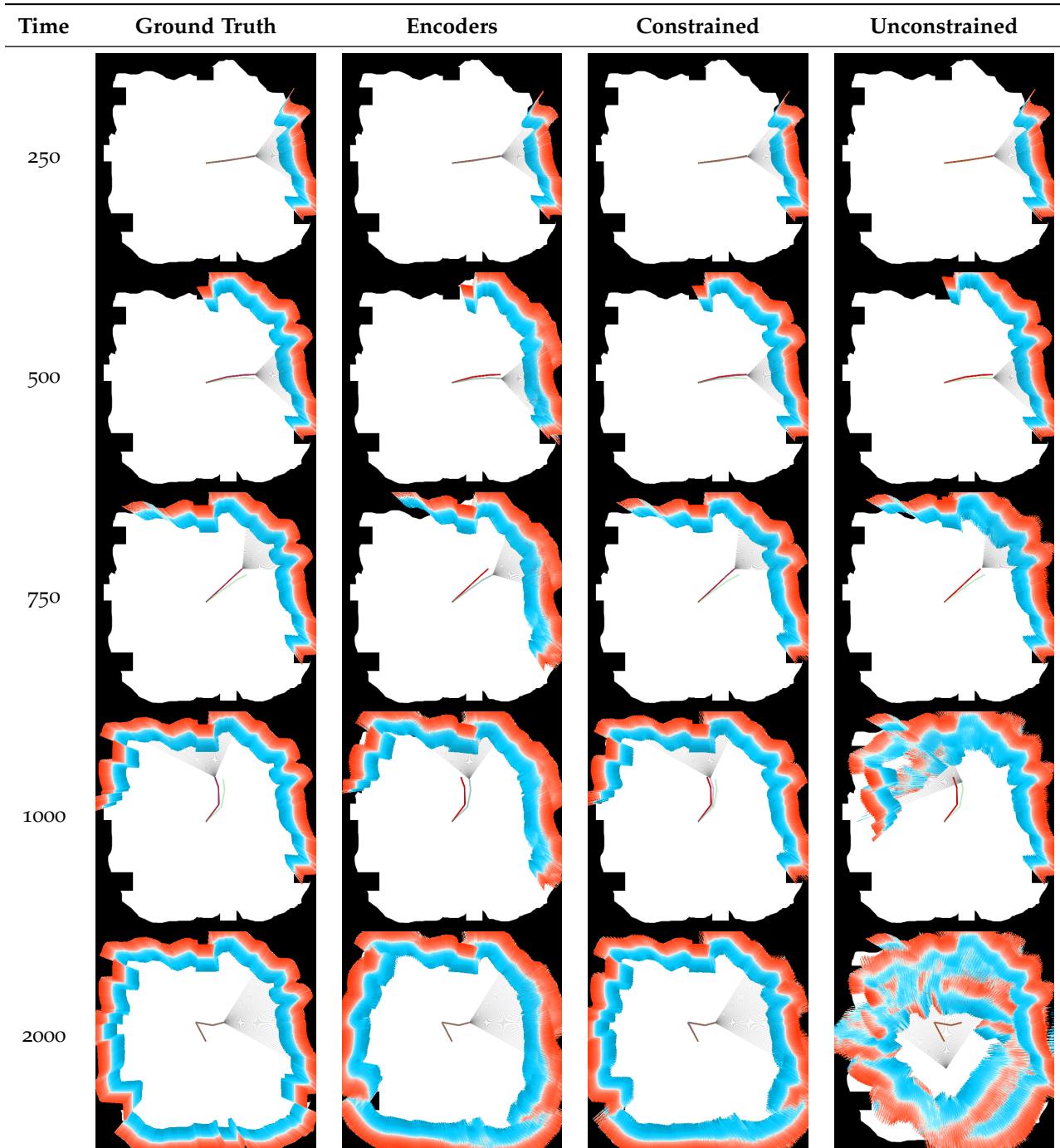


Table 13 : 2D experiments. A 3-joint robot arm (dark red) scans the scene with a simulated depth camera (light grey lines). Its joint encoders are corrupted by noise (green). It builds up a TSDF,  $\Phi$  (red and blue pixels) of the scene (black pixels). The ground truth is compared to just using the robot encoders, constrained descent (Alg. 2), and unconstrained descent over time.

### 3D Simulation Experiments

The same concept can be extended to 3D scenes (Fig. 13). I built a framework for doing dense sensor-to-model articulated tracking (Alg. 2) in 3D, as well as a depth sensor simulator. Like in the 2D experiments, the robot follows a fixed trajectory. I again add Perlin noise to the joint encoders. The robot builds a 3D TSDF using the simulated depth data (Fig. 15). Fig. 16 shows the result of the same robot trajectory and depth data being used to reconstruct a bookshelf without any joint encoder noise, with noise and dense articulated 3D tracking, and without any tracking. Even under conditions of extreme encoder noise, dense frame-to-model SLAM is able to create a high-quality 3D reconstructions of the scene. On the other hand, using the noisy encoders alone, the 3D reconstruction is almost totally unrecognizable.

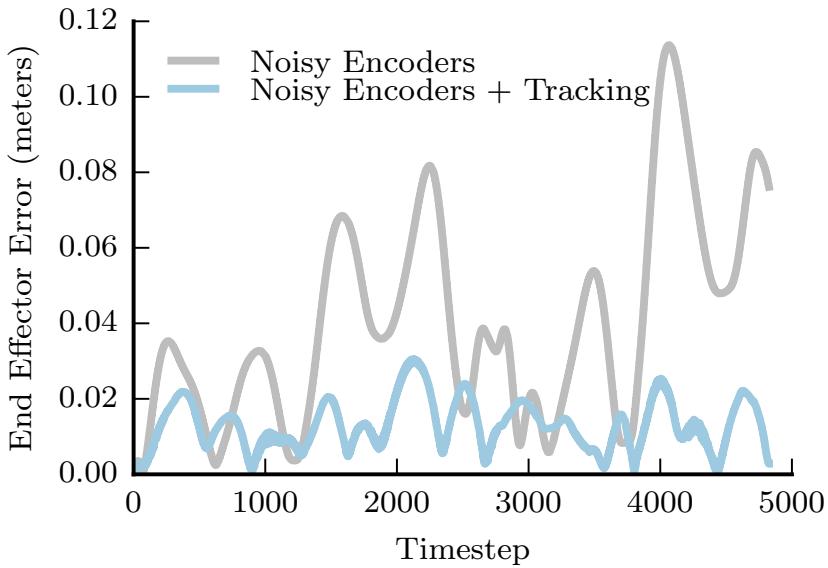


Figure 14: The end effector error in the 3D simulation experiment.

Fig. 14 shows the translation error of the end effector vs. ground truth during tracking. The end effector error consistently converges to be beneath 2 cm even when the error due to noisy joint encoders is more than 11 cm.

The deficiencies of the simple gradient descent algorithm can be seen around timestep 500, 3000, etc. where the noisy encoders have better end effector accuracy than the tracking. This is because the prior motion model is too simple to account for high-frequency noise.

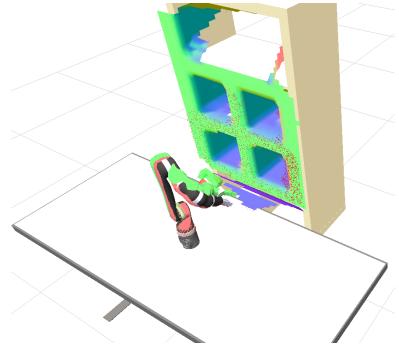


Figure 13: A 3D simulation environment featuring a robot mounted to a table, looking at a bookshelf.

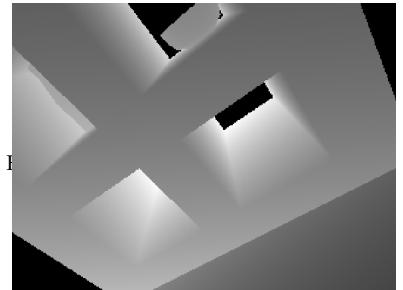
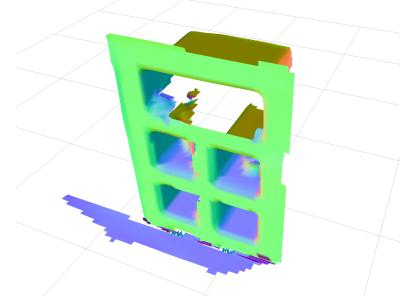
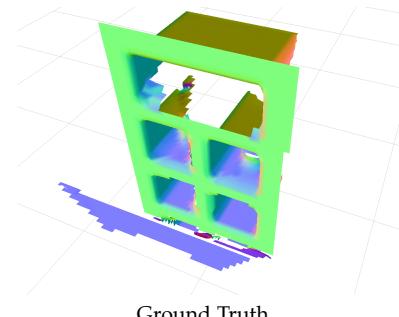
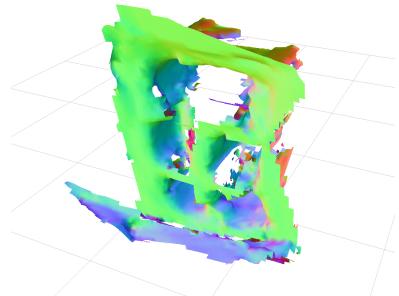


Figure 15: Simulated depth camera.



Noisy Encoders + Tracking



Noisy Encoders Only  
Figure 16: Simulated 3D TSDF reconstructions.

# *Research Questions*

This thesis proposal has introduced a novel framework for doing 3D visual SLAM in articulated systems; but much further research is needed to investigate the theoretical grounding and concrete properties of this framework. Here are a few questions which must be answered by the thesis:

## ***What applications are enabled by articulated SLAM?***

Building high quality 3D reconstructions is one thing, but what can be done with the output of the articulated SLAM system? One possibility is to use the model to plan manipulation strategies (such as grasping and reaching into tight spaces), or to control the arm (*i.e.* visual servoing). When the robot is no longer passively observing the model of the world, but instead can control where it places the sensor, the whole field of active perception opens up. How should the robot move to best reconstruct the scene?

As interesting as these questions are, they will have to be secondary to the fundamental question of how to get the SLAM system to work in the first place.

## ***How can the encoder noise be modeled on real robot arms?***

So far, I have only considered vague properties of encoder noise<sup>33</sup> – that it is really unmodeled dynamics causing the discrepancy between joint encoders and the actual joint angles, and that it is likely configuration-dependant. In simulated experiments, I have modeled the noise using a random Perlin function<sup>34</sup>. But what exactly *are* the unmodeled dynamics I am trying to capture? What physical processes cause them, and how can these physical processes be modeled?

Answering this question will require a literature review of existing methods, data collection on real robots, and a framework for testing different dynamics models.

<sup>33</sup> See page [21](#)

<sup>34</sup> See page [27](#)

## *What are the limitations of dense camera-to-model SLAM using real data/robots?*

So far, my experiments have only covered 2D and 3D simulated datasets. Obviously, real robots have unique challenges not addressed by the simulations: extrinsic/intrinsic calibration error, sensor noise, unmodeled dynamics, etc. I am eager to test articulated SLAM on a real robot. Experiments must be devised to test the convergence basin of the SLAM system. One particularly challenging component will be getting ground-truth data for the pose of the robot and geometric structure of the world (which are easy to get in simulation).

## *How can a pose graph be incorporated into the articulated SLAM framework?*

So far, I have only derived dense camera-to-model 3D SLAM<sup>35</sup>. This approach works well for small workspaces when the displacement between camera poses is small, but fails over long trajectories and large workspaces. The camera-to-keyframe or pose graph methods<sup>36</sup> correct for drift by globally optimizing the entire trajectory with respect to every keyframe. Incorporating a pose graph into the articulated SLAM problem would be an ideal way of correcting for drift.

Answering this question will require re-deriving pose graph optimization in terms of *robot configurations* rather than *camera poses*.

<sup>35</sup> See page [23](#)

<sup>36</sup> See page [14](#)

## *How can direct monocular articulated SLAM be accomplished?*

So far, I have only discussed solutions for dense camera-to-model 3D SLAM, which requires knowledge of a depth image at all time steps. How can I reconstruct a scene and estimate the robot's configuration when only a *monocular* camera is available? Direct or semi-direct<sup>37</sup> SLAM is an attractive set of solutions for 3D reconstruction using monocular cameras. Such methods rely on direct optimization in  $SE(3)$  of camera poses given reprojected pixel intensities.

<sup>37</sup> See page [16](#)

Answering this question will require constructing a direct image residual optimization technique in the configuration space of the robot, instead of in  $SE(3)$ . That is, I will have to derive the direct relationship between *pixel intensities* and *joint angle displacements* (an equivalent to equation Eq. [21](#)).

I have done some back-of-the-envelope calculations which seem to suggest that the image Jacobian with respect to the manipulator can

be computed as:

$$J_q|_{p_i} = \frac{\partial r_i}{\partial q} = \nabla I|_{p_i} \frac{\partial \pi}{\partial p_i}|_{F(q)} \frac{\partial x_i}{\partial q}. \quad (44)$$

where  $\nabla I$  is an image gradient,  $p_i$  is an image pixel,  $\pi$  is the pinhole projection function, and  $x_i$  is a pixel as projected into the scene.

$J_q$  will take the form of a set of  $N$  images, one for each robot joint. Each image will correspond to the change in pixel intensities with respect to the incremental movement of that joint. Notice that the last term of (44) is just the manipulator Jacobian evaluated for a particular point in the point cloud.

### *Why articulated SLAM?*

This is perhaps the biggest research question. What advantages does an articulated SLAM system have over a traditional unconstrained 6DOF SLAM system? I believe that doing SLAM directly in the natural space of the robot arm (its configuration space) is a more accurate and theoretically grounded way of solving the problem – but data will have to be collected to verify this, comparing existing SLAM systems with their articulated counterparts.



## *Timeline*

I want to complete this thesis work in one year. I will be focusing my efforts on two major conference papers: one for ICRA 2016, and one for IROS 2016. Other conferences may be selected depending on how my research is going.

Date	Conference
July 11, 2015	Present CHISEL at RSS.
August 2015	Real experiments with dense frame-to-model tracking.
Sep 15, 2015	Write paper on dense frame-to-model tracking.
Fall 2015	Begin implementing direct monocular system.
Mar 1, 2016	Write paper on direct monocular system.
Spring 2016	Implement pose graph.
June 8, 2016	Finish and defend this thesis

The first paper will focus on dense camera-to-model SLAM <sup>38</sup>, and will theoretical ground it, and provide results against other (non-articulated) SLAM algorithms. The second paper will focus on a direct monocular method <sup>39</sup>. Finally, I will look at incorporating a pose graph <sup>40</sup>. Work from these conference papers will be incorporated into my thesis.

<sup>38</sup> See page [23](#)

<sup>39</sup> See page [16](#)

<sup>40</sup> See page [14](#)



## Bibliography

Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '96*, 1996.

J. Engel, T. Schöps, and D. Cremers. LSD-SLAM: Large-scale direct monocular SLAM. In *EECV*, 2014.

R. I. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, ISBN: 0521540518, second edition, 2004.

Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, and Andrew Fitzgibbon. Kinect-fusion: Real-time 3d reconstruction and interaction using a moving depth camera. In *UIST 2011*, 2011.

M. Kaess, A. Ranganathan, and F. Dellaert. iSAM: Incremental smoothing and mapping. *IEEE Trans. on Robotics (TRO)*, 24(6):1365–1378, December 2008.

M. Keller, D. Lefloch, M. Lambers, S. Izadi, T. Weyrich, and A. Kolb. Real-time 3d reconstruction in dynamic scenes using point-based fusion. In *3D Vision - 3DV 2013, 2013 International Conference on*, pages 1–8, June 2013. DOI: [10.1109/3DV.2013.9](https://doi.org/10.1109/3DV.2013.9).

Matthew Klingensmith, Thomas Galluzzo, Christopher Dellin, Moslem Kazemi, J Andrew Bagnell, and Nancy Pollard. Closed-loop Servoing using Real-time Markerless Arm Tracking. In *International Conference on Robotics And Automation (Humanoids Workshop)*, May 2013.

Matthew Klingensmith, Ivan Dryanovski, Siddhartha Srinivasa, and Jizhong Xiao. Chisel: Real time large scale 3d reconstruction onboard a mobile device. In *RSS*, 2015.

Matthew T. Mason . *Mechanics of Robotic Manipulation*. MIT Press, Cambridge, MA, August 2001.

L. Mundermann, S. Corazza, and T.P. Andriacchi. Accurately measuring human movement using articulated icp with soft-joint constraints and a repository of articulated models. In *Computer Vision and Pattern Recognition, 2007. CVPR '07. IEEE Conference on*, pages 1–6, June 2007. DOI: 10.1109/CVPR.2007.383302.

Raúl Mur-Artal and Juan Tardós. Probabilistic semi-dense mapping from highly accurate feature-based monocular slam. In *RSS*, 2015.

Richard A. Newcombe, Steven J. Lovegrove, and Andrew J. Davison. Dtam: Dense tracking and mapping in real-time. In *ICCV*, 2011.

Ken Perlin. Improving noise. In *SIGGRAPH*, 2002.

Tanner Schmidt, Richard Newcombe, and Dieter Fox. DART: Dense Articulated Real-Time Tracking. In *Robotics: Science and Systems*, number 1, 2014. URL <http://www.roboticsproceedings.org/rss10/p30.pdf>.

A. Segal, D. Haehnel, and S. Thrun. Generalized-icp. In *Proceedings of Robotics: Science and Systems*, Seattle, USA, June 2009.

Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic Robotics (Intelligent Robotics and Autonomous Agents)*. The MIT Press, 2005. ISBN 0262201623.