

AN AUTONOMOUS METHOD FOR MEASURING 3D JOINT KINEMATICS FROM
2D XRAY IMAGES

By

ANDREW JAMES JENSEN

A DISSERTATION PRESENTED TO THE GRADUATE SCHOOL
OF THE UNIVERSITY OF FLORIDA IN PARTIAL FULFILLMENT
OF THE REQUIREMENTS FOR THE DEGREE OF
DOCTORATE OF PHILOSOPHY

UNIVERSITY OF FLORIDA

2023

© 2023 Andrew James Jensen

Dedication placeholder.

ACKNOWLEDGEMENTS

First, I would like to thank my parents for their continued support throughout my undergraduate and graduate studies. They fostered my work ethic and drive to get through something as ambitious as a doctorate.

Secondly, I would like to thank my advisor, Scott Banks. He has truly been something of a second dad to me as I navigate through graduate school and my future career. Though most conversations seem to revolve around either beer, music, or hiking, we seem to have gotten quite a lot done.

Third, I would like to thank my wife, Lauren. Though busy with medical school, she still takes the time to understand my research and offer support outside of academics.

TABLE OF CONTENTS

page

LIST OF TABLES

<u>Tables</u>	<u>page</u>
---------------	-------------

LIST OF FIGURES

<u>Figures</u>	<u>page</u>
----------------	-------------

LIST OF ABBREVIATIONS

- TKA Total Knee Arthroplasty. This is the complete or partial resurfacing of the articulating surfaces in the knee.
- ML Machine Learning. This is the process of feeding a computer inputs and outputs in order to determine an algorithm that goes from input → output
- CNN Convolutional Neural Network. This is a type of neural network that uses convolution kernels as the operation between each of the layers
- HRNet High Resolution Convolutional Neural Network [?].
(<https://github.com/HRNet>)

Abstract of Dissertation Presented to the Graduate School
of the University of Florida in Partial Fulfillment of the
Requirements for the Degree of Doctorate of Philosophy

AN AUTONOMOUS METHOD FOR MEASURING 3D JOINT KINEMATICS FROM
2D XRAY IMAGES

By

Andrew James Jensen

March 2023

Chair: Scott Banks

Major: Mechanical Engineering

The primary function of human synovial joints is to support the dynamic motion of the musculoskeletal system. The diseases that affect these systems typically manifest during movement, with mild to severe pain arising during specific activities or during particular ranges of motion. Unsurprisingly, the financial burden of musculoskeletal diseases is roughly USD 300 billion per year in direct and indirect costs [?]. One of the most common conditions affecting human joints is osteoarthritis, which involves the progressive loss of the cartilage between the joint surfaces over time [?]. A highly effective solution for osteoarthritis is arthroplasty, which involves a partial or complete removal and resurfacing of the affected joint with polymeric and metallic components intended to relieve pain and restore a degree of natural function and motion. Despite being highly effective, roughly 20% of patients receiving total knee arthroplasty express some form of dissatisfaction, usually manifested as pain, instability or stiffness during movement [? ? ?]. Surprisingly, standard clinical musculoskeletal diagnostic methods are entirely static. That is, clinicians do not have clinically practical ways to quantify skeletal motion during weight-bearing or dynamic movement when most pain symptoms occur. Unfortunately, most of the tools used to accurately quantify 3D dynamic motion (e.g., 3D motion capture, radiostereometry, fluoroscopic model/image registration) are prohibitively expensive or impractical to use in clinical settings. Methods using single-plane fluoroscopic or flat-panel

imaging with 3D-to-2D model-image registration have been used since the 1990s. They have been shown to provide sufficient accuracy for many clinical joint assessment applications, including natural and replaced knees [? ? ? ?], natural and replaced shoulders [? ? ? ? ?], and extremities [? ? ? ? ?]. One benefit of this approach is that suitable images can be acquired with equipment commonly found in most hospitals. The main impediment for this technology to be used clinically is the time and expense of human operators to supervise the model-image registration process. If the need for human supervision for model-image registration were eliminated, then fluoroscopic imaging could provide a reliable, inexpensive, and accurate method to provide 3D dynamic joint kinematics in a clinical setting. State-of-the-art techniques for generating kinematics using model-image registration involve numerical optimization techniques that iteratively match bone or implant model projections in dynamic x-ray images [? ? ?]. These methods provide accurate 3D bone or implant kinematics when given a rough initial pose estimate for numerical optimization [?]. However, these methods still require human input for an initial pose estimate, making them impractical for clinical use. Recent advancements in computational capabilities and machine learning algorithms provide tools that are well-suited to replace human supervision for a range of time-consuming tasks including model-image registration, implant segmentation, and implant pose initialization. The goal of this dissertation is to provide a framework for the fully autonomous measurement of TKA kinematics from single-plane imaging using a mixture of both recent and historic algorithms.

CHAPTER 1 INTRODUCTION

Total Knee Arthroplasty (TKA) is a standard procedure for alleviating symptoms related to osteoarthritis in the knee. In 2018, orthopaedic surgeons performed more than 715,000 TKA operations in the United States [?]. This number is projected to increase to 3.48 million by 2030 [?] due to an aging population and increased obesity rates. While TKA largely relieves symptomatic osteoarthritis, roughly 20% of TKA patients express postoperative dissatisfaction, citing mechanical limitations, pain, and instability as the leading causes [? ? ?]. Standard methods of musculoskeletal diagnosis cannot quantify the dynamic state of the joint, either pre- or post-operatively; clinicians must rely on static imaging (radiography, MRI, CT) or qualitative mechanical tests to determine the condition of the affected joint, and these tests cannot easily be performed during weight-bearing or dynamic movement when most pain symptoms occur. Unfortunately, most of the tools used to quantify 3D dynamic motion are substantially affected by soft-tissue artifacts [? ? ?], are prohibitively time-consuming or expensive [?], or cannot be performed with equipment available at most hospitals.

Model-image registration is a process where a 3D model is aligned to match an object's projection in an image [?]. Researchers have performed model-image registration using single-plane fluoroscopic or flat-panel imaging since the 1990s. Early methods used pre-computed distance maps [? ?], or shape libraries [? ? ?] to match the projection of a 3D implant model to its projection in a radiographic image. With increasing computational capabilities, methods that iteratively compared implant projections to images were possible [? ? ?]. Most model-image registration methods provide sufficient accuracy for clinical joint assessment applications, including natural and replaced knees [? ? ? ?], natural and replaced shoulders [? ? ? ?], and extremities [? ? ?]. One of the main benefits of this single-plane approach is that suitable images can be acquired with equipment found in most hospitals. The main impediment to implementing this approach into a standard clinical workflow is the time and expense of human operators to supervise the model-image

registration process. These methods require either (1) an initial pose estimate [? ?], (2) a pre-segmented contour of the implant in the image [? ?], or (3) a human operator to assist the optimization routine out of local minima [?]. Each of these requirements makes model-image registration methods impractical for clinical use. Even state-of-the-art model-image registration techniques [?] require human initialization or segmentation to perform adequately.

Machine learning algorithms automate the process of analytical model building, utilizing specific algorithms to fit a series of inputs to their respective outputs. Neural networks are a subset of machine learning algorithms that utilize artificial neurons inspired by the human brain’s connections [?]. These networks have shown a great deal of success in many computer vision tasks, such as segmentation [? ? ?], pose estimation [? ?], and classification [? ? ?]. These capabilities might remove the need for human supervision from TKA model-image registration.

?? gives an overview of some of the fundamental building blocks of computer vision and projective geometry. A cursory understanding of these operations is necessary to the underlying mathematics motivating model-image registration, as well as recreating the images generated by fluoroscopic imaging systems. It also outlines a few key classes of image similarity metrics, which are useful when defining objective functions for optimization routines.

?? discusses neural networks from a bottom-up perspective, and explains how they can be used to automate different tasks in computer vision, with a strong emphasis on image processing. Understanding the “how” behind neural networks and deep learning drives their usage in automating historically tedious and human-supervised tasks.

?? outlines historical methods of utilizing model-image registration for measuring total knee and natural knee kinematics from various imaging modalities. The main emphasis of this chapter is to explain the breadth of research that has been performed in this space, and also explain fundamental limitations with each method that prevent the

adoption of this technology into a clinical setting.

?? utilizes the previous chapters to describe a method for autonomously measuring total knee arthroplasty kinematics from single-plane fluoroscopic imaging. Utilizing the framework built by computer vision and machine learning, and driven by the shortcomings of historic methods, a novel pipeline is presented that removes the need for human supervision for this technology, bringing it one step closer to clinical adoption.

Lastly, ?? discusses the aims of this thesis, beginning with a validation of the pipeline, along with a pilot run with some patients, as well as distribution of the software packages to other researchers as freely distributed libraries and binaries.

1.1 Computer Vision Primer

In order to more fully understand and appreciate the process of extracting TKA kinematics from single-plane digital images, we must first explore the basics of computer vision. Specifically, the elements of image formation of and object projection.

1.1.1 Geometric Transformations

Geometric primitives, such as points, are fundamental building blocks for representing shapes and objects in computer graphics. In this section, we will discuss how points can be represented in 2D and 3D space, along with the different mathematical operations that you can perform on those points.

1.1.1.1 2D and 3D Points

In N-dimensional space, a point is represented as a set of N scalars, each representing a magnitude in a particular direction. This can be represented mathematically as a column vector (??). In two dimensions, a point is represented with two elements, $\mathbf{x} = [x, y]^T$. Similarly, in three dimensions, a point is represented with three elements, $\mathbf{x} = [x, y, z]^T$.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{N-1} \\ x_N \end{bmatrix} \in \mathbb{R}^N \quad (1-1)$$

Homogeneous coordinates provide a way to represent points with an additional scale factor, \tilde{w} . This allows us to perform rotations and translations simultaneously and successively using matrix multiplications. Homogeneous coordinates for a point in N-dimensional space are represented as a column vector with N+1 elements (??). In most model-image registration applications, the scale factor \tilde{w} is set to 1. When dealing with homogeneous coordinates, $\tilde{\mathbf{x}} = \tilde{w}\bar{\mathbf{x}}$, where $\tilde{\mathbf{x}}$ is the scaled version of $\bar{\mathbf{x}}$, and $\bar{\mathbf{x}}$ is simply the original vector, \mathbf{x} , with a 1 appended.

$$\tilde{\mathbf{x}} = \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \\ \vdots \\ \tilde{x}_N \\ \tilde{w} \end{bmatrix} = \tilde{w} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = \tilde{w}\bar{\mathbf{x}} \quad (1-2)$$

1.1.1.2 2D Transformations

Transformations are operations that change the position, orientation, or shape of an object in 2D space. One of the most basic transformations is a translation, which moves an object by adding a displacement vector to its position (??).

$$\mathbf{x}' = \begin{bmatrix} x \\ y \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} = \mathbf{x} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} = \mathbf{x} + \mathbf{t}$$

Or, using homogeneous coordinates and matrix multiplication

(1-3)

$$= \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = [\mathbf{I} \quad \mathbf{t}] \bar{\mathbf{x}}$$

In this equation, \mathbf{x} is the original position of the object, \mathbf{x}' is the transformed position, and \mathbf{t} is the displacement vector that specifies the amount of translation in the x and y directions. Using homogeneous coordinates and matrix multiplication allows for the convenient representation of multiple transformations as a single matrix multiplication, as well as for the composition of multiple transformations (??).

$$\begin{aligned} \bar{\mathbf{x}}'' &= \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & q_x \\ 0 & 1 & q_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} x + t_x + q_x \\ y + t_y + q_y \\ 1 \end{bmatrix} \end{aligned} \quad (1-4)$$

The next type of 2D transformation is a rotation, which changes the orientation of an object, but not its shape (????).

$$\mathbf{x}' = \mathbf{R}\mathbf{x} \quad (1-5)$$

where

$$\mathbf{R} = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \quad (1-6)$$

This will rotate an object by θ in the counter clockwise direction.

It's also possible to perform a rotation and a translation at the same time, by replacing the identity matrix \mathbf{I} in Equation ?? with the rotation matrix \mathbf{R} from Equation ???. This results in a transformation that preserves lengths and angles while rotating and translating the rigid object (??).

$$\mathbf{x}' = \begin{bmatrix} \mathbf{R}_{2 \times 2} & \mathbf{t} \end{bmatrix} \bar{\mathbf{x}} \quad (1-7)$$

A scaled rotation will change the size of the object by some scalar factor, s (??); this transformation preserves angles.

$$\mathbf{x}' = \begin{bmatrix} s\mathbf{R}_{2 \times 2} & \mathbf{t} \end{bmatrix} \bar{\mathbf{x}} \quad (1-8)$$

An affine transformation preserves parallelism, and is simply a pre-multiplication by an arbitrary 2×3 matrix (??).

$$\mathbf{x}' = \mathbf{A}\bar{\mathbf{x}}$$

where

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \quad (1-9)$$

A projection matrix (or perspective transformation) is one that operates on homogeneous coordinates (??).

$$\tilde{\mathbf{x}}' = \tilde{\mathbf{H}}\tilde{\mathbf{x}} \quad (1-10)$$

To obtain inhomogeneous results, the resultant $\tilde{\mathbf{x}}'$ must be normalized (??). A projective transformation preserves straight lines.

$$\begin{aligned} x' &= \frac{\tilde{x}}{\tilde{w}} & y' &= \frac{\tilde{y}}{\tilde{w}} \\ &= \frac{h_{11}x + h_{12}y + h_{13}}{h_{31}x + h_{32}y + h_{33}} & &= \frac{h_{21}x + h_{22}y + h_{23}}{h_{31}x + h_{32}y + h_{33}} \end{aligned} \quad (1-11)$$

1.1.1.3 3D Transformations

Transformations in three dimensions are extremely similar to their two dimensional counterpart, except the dimensionality of each transformation is increased by one. However, three dimensional rotations introduce added complexity, due to multiple axes of rotation, and the non-commutativity of those rotations. We handle this by utilizing Euler angles [?], which decomposes the final orientation of an object as a combination of rotations about the three coordinate axes. Each rotation about an axis is represented with a matrix (??), and the final rotation matrix is represented as the order of the rotations that produced it (e.g. Z-X-Y rotation).

$$R_x = \begin{bmatrix} 1 & 0 & 0 \\ 0 & c_x & -s_x \\ 0 & s_x & c_x \end{bmatrix} \quad R_y = \begin{bmatrix} c_y & 0 & s_y \\ 0 & 1 & 0 \\ -s_y & 0 & c_y \end{bmatrix} \quad R_z = \begin{bmatrix} c_z & -s_z & 0 \\ s_z & c_z & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (1-12)$$

Anatomically, the rotations can be chosen such that they align with the anatomic axes of the joint. This allows for congruence between the mathematical notation used to represent an object in space, and the medical notation used by orthopaedic practitioners (i.e. aligning a Z-X-Y rotation decomposition with the flexion-abduction-external rotation order used by orthopaedic surgeons).

For simplicity, we can also represent 3D rotations as an axis-angle decomposition, which describes the transformation as a single rotation about an arbitrary axis [?] (??). You can use trigonometric identities to determine the axis and angle given an arbitrary rotation matrix.

$$R_{3 \times 3} = \begin{pmatrix} m_x^2 v + c & m_x m_y v - m_z s & m_x m_z v - m_y s \\ m_x m_y v + m_z s & m_y^2 v + c & m_y m_z v - m_x s \\ m_x m_z v - m_y s & m_y m_z v + m_x s & m_z^2 v + c \end{pmatrix}$$

where

$$s = \sin(\theta)$$

$$c = \cos(\theta) \quad (1-13)$$

$$v = (1 - c)$$

and

$$\mathbf{m} = \begin{bmatrix} m_x \\ m_y \\ m_z \end{bmatrix} \text{ is the axis of rotation}$$

3D to 2D Projections Geometric primitives, such as points, lines, and polygons, are the building blocks of computer graphics. These basic shapes can be transformed and combined in various ways to create more complex objects and scenes. In order to represent these 3D primitives and objects in 2D image space, we must apply transformations that manipulate their position, orientation, and other properties. By using these transformations, we can create the illusion of depth and spatial relationships on a flat display. Understanding how to work with primitives and transform them is crucial for creating a wide range of visual effects and graphics in computer graphics.

In computer graphics, one of the most basic methods of projecting three-dimensional objects onto a two-dimensional image plane is an orthographic projection. This projection simply drops the depth component of the object and flattens it onto the image plane (??). This can be thought of as mimicking a camera with a long focal length, or when the depth of the object is shallow compared to its distance from the camera (also known as a weak perspective projection). In this equation, \mathbf{p} represents a point in 3D space and \mathbf{x} represents

the projected point in 2D image space. The $\tilde{\cdot}$ still represents the point in homogeneous coordinates.

$$\tilde{\mathbf{x}} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{p}} \quad (1-14)$$

Perspective Projection The most common 3D-2D projection is *perspective projection*, which accounts for depth perception. This can be done by scaling each point by it's z position relative to the camera (Eq. ??). This necessitates using homogeneous coordinates (Eq. ??).

$$\bar{\mathbf{x}} = \mathcal{P}_z(\mathbf{p}) = \begin{bmatrix} x/z \\ y/z \\ 1 \end{bmatrix} \quad (1-15)$$

$$\tilde{\mathbf{x}} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \tilde{\mathbf{p}} \quad (1-16)$$

The perspective projection is the cornerstone of model-image registration. If a matrix mimics the fluoroscopic projective geometry in the clinic, then this offers a strong possibility of re-creating a virtual scene with the same geometry as the actual scene.

1.1.2 Image Formation and Camera Properties

Using our knowledge of geometric transformations and projective geometries, we can build up a model of a camera step by step. We standardize our reference frame by having the z direction along the focal direction of the camera, the x direction to the right, and the y direction such that the right-hand rule is maintained. The origin is at the center of the camera.

We will describe the object of interest as a collection of points,

$\mathbf{p}_i = [x_i, y_i, z_i, 1]$ for $i = 1, 2, \dots, N$. For a complete picture, any given operation will be performed on all points. For simplicity, the following equations will demonstrate the process on a single point of the object.

First, we need to describe the location and orientation of the object with respect to the camera. This can be done with a 3D homogeneous transformation matrix (translation and rotation) (Eq. ??).

$$\tilde{\mathbf{p}}' = \begin{bmatrix} \mathbf{R}_{3 \times 3} & \mathbf{t}_{3 \times 1} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \tilde{\mathbf{p}} \quad (1-17)$$

Then, we use a projective transformation that determines the perspective scaling between the object's location and the image plane at the focal distance. Geometrically, this relationship can be visualized with similar triangles (??) and quantified using the ratio of lengths (??). In these equations, f' is the focal distance in units of length.

$$\tilde{\mathbf{x}}_i = \begin{bmatrix} f' & 0 & 0 \\ 0 & f' & 0 \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{p}}' \quad (1-18)$$

where

$$x_i = p'_x \frac{f'}{p'_z}$$

$$y_i = p'_y \frac{f'}{p'_z}$$

The standard image reference frame places the origin at the top left corner of the image, with the positive x-direction to the right, and the positive y-direction down. This introduces the idea of the principal point, (c_x, c_y) , which is the location where the optical axis of the camera intersects the image plane perpendicularly. In the camera model, the principal point is a translation starting at the image plane's origin (Eq. ??).

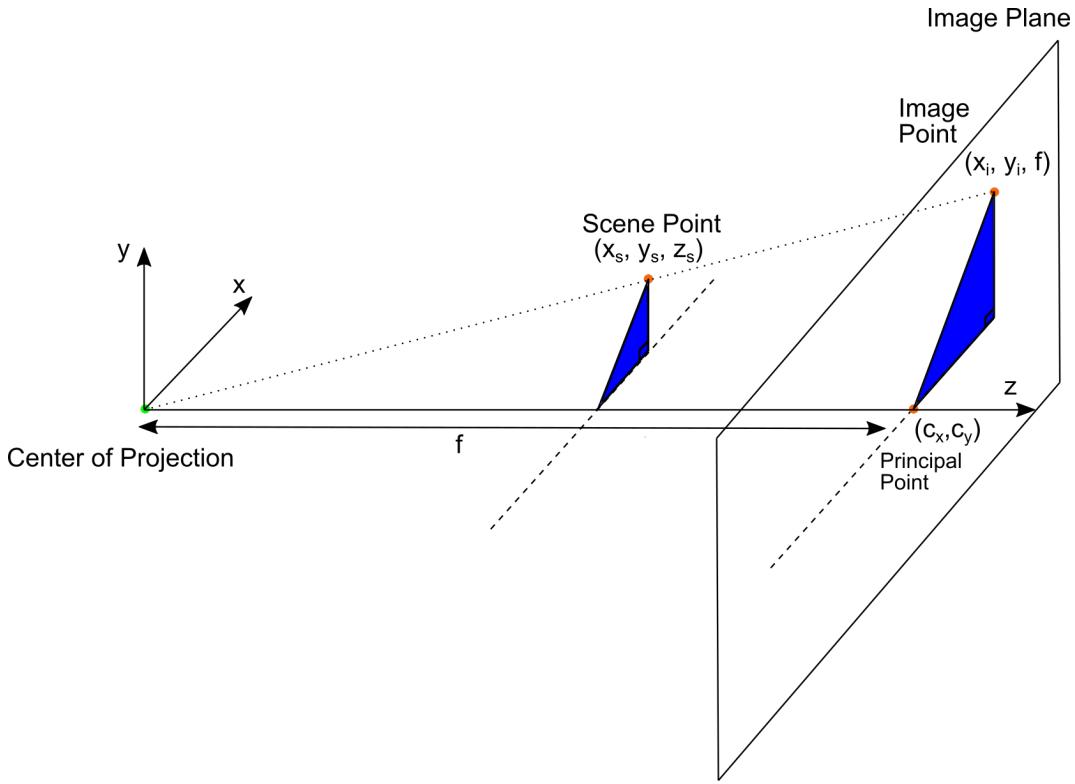


Figure 1-1. The geometry of perspective projection can be visualized by using similar triangles. The overall scaling of the image is based on the ratio of the focal length to the depth of the object.

$$\tilde{\mathbf{x}}_i = \begin{bmatrix} f' & 0 & c_x \\ 0 & f' & c_y \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{p}}' \quad (1-19)$$

where

$$c_x \approx \frac{W_{image}}{2}$$

$$c_y \approx \frac{H_{image}}{2}$$

Finally, we need to convert the image coordinates, $\tilde{\mathbf{x}}_i$, to pixel coordinates using the pixel scale factor (Eq.??). The pixel scale factor is defined by the parameters k_x and k_y , which represent the number of pixels per unit distance in the x and y directions, respectively. By multiplying the perspective projection matrix by the pixel scale factor, we

can obtain a new matrix that maps 3D points in world coordinates directly to pixel coordinates on the image.

$$\begin{aligned}\tilde{\mathbf{x}}_{pix} &= \begin{bmatrix} k_x & 0 & 0 \\ 0 & k_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f' & 0 & c_x \\ 0 & f' & c_y \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{p}}' \\ &= \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{p}}'\end{aligned}\quad (1-20)$$

where

$$f_x = k_x f' \text{ and } f_y = k_y f'$$

Are focal distances in units of pixels

At this point, we have placed an object in front of a camera using 3D transformations, and then projected the points of that object onto the image plane using a camera matrix and projective transform. A simple silhouette rasterization of the image would simply fill in the pixel values interpolated between each of the discrete points obtained from this projection. Generating a photo-realistic image would involve coloring, shading, and ray tracing, but this is unnecessary for this dissertation, and will not be discussed.

1.1.3 Image Processing

Digital image processing is a field of computer vision that deals with the manipulation, analysis, and interpretation of digital images. It focuses on algorithms and techniques that extract meaningful information from images and enhance visual quality. We can leverage these algorithms in order to efficiently and autonomously extract information relevant to model-image registration.

1.1.3.1 Filtering and Convolution

We have already seen how image formation yields a collection of 2D points, \mathbf{x}_{pix} . We can write the intensity values at each pixel location as a digital signal, $f(\mathbf{x}_{pix}) = f(i, j)$,

where (i, j) represents the pixel locations in the image, and the function returns the intensity value. We can then use standard methods of digital signal processing in order to extract meaningful information from images.

The most widely used filter is a linear filter [?], where the output is some linear operation on the neighboring pixels (??). This is known as a *convolution*. In a convolution, the kernel, h , is shifted along the input image, f , and the resultant image, g , is the dot product of those two matrices at that specific location.

$$\begin{aligned} g(i, j) &= \sum_{k,l} f(i - k, j - l)h(k, l) \\ &= \sum_{k,l} f(k, l)h(i - k, j - l) \end{aligned} \tag{1-21}$$

Where we use the following notation

$$g = f * h$$

The convolution operation is *linear shift invariant*, which means that it obeys the superposition principle (??) and the shift invariance principle (??). This is a powerful property, because it will behave the same everywhere on the input signal/image (e.g. an edge detector will detect an edge, no matter where the edge is on the input image)

$$h * (f + g) = h * f + h * g \tag{1-22}$$

$$g(i, j) = f(i + k, j + l) \iff (h * g)(i, j) = (h * f)(i + k, j + l) \tag{1-23}$$

A common filter applied to images is the Gaussian kernel (??). This kernel is shaped as a 2D discrete Gaussian, and has the effect of blurring an image and removing noise.

$$\text{Gaussian filter} = \frac{1}{256} \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix} \quad (1-24)$$

Another is the box kernel, which averages the value of the nearest K pixels (??).

$$\text{Box filter} = \frac{1}{K^2} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & 1 & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} \quad (1-25)$$

Edge filters can also be created to detect vertical (??), horizontal (??), or diagonal edges (??). As each of the filters moves across the feature it is designed for, that region of the output will be more highly activated than other regions, extracting out the desired components. The orientation of each of these filters can be hand-selected to find desirable attributes in images.

$$\text{vertical edge filter} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad (1-26)$$

$$\text{horizontal edge filter} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \quad (1-27)$$

$$\text{diagonal edge filters} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \text{ and } \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \quad (1-28)$$

Lastly, we can use a corner filter to find corners in images (??).

$$\text{Corner filter} = \frac{1}{4} \begin{bmatrix} 1 & -2 & 1 \\ -2 & 4 & -2 \\ 1 & -2 & 1 \end{bmatrix} \quad (1-29)$$

Entire subfields of computer vision and image analysis are devoted to creating increasingly useful and complex filters, or collections of filters. These include steerable filters, which determine information occurring in arbitrary directions [?], recursive filtering [?], and non-linear filtering [?].

1.1.3.2 Edge Detection

Edge detection is a highly motivated sub-field of image processing in computer vision due to the immense usefulness of algorithmically determining the edges in a given image. For a human operator, it can be rather easy to find edges of interest, but how much this be incorporated computationally? The first approach might be in viewing an image topologically, with regions of different colors and intensity represented by different “heights”. Then, an edge simply becomes an area with a steep gradient (??).

$$\mathbf{J}(\mathbf{x}) = \nabla I(\mathbf{x}) = \left(\frac{\partial I}{\partial x}, \frac{\partial I}{\partial y} \right)(\mathbf{x}) \quad (1-30)$$

Finding the direction of the steepest ascent/descent at any given location will give us the normal to the local edge at that point. However, the derivative operator will accentuate and amplify high frequencies in the image, causing noise to overpower the signal. Removing the high-frequency information (a low-pass filter) in the image results in gradient detection that is much more aligned with the salient edges of the image. The

Gaussian kernel is a good option for an isotropic low-pass filter on a 2D signal (image) (??)

$$\begin{aligned}
 \mathbf{J}_\sigma(\mathbf{x}) &= \nabla[G_\sigma(\mathbf{x} * I(\mathbf{x}))] \\
 &= \nabla G_\sigma(\mathbf{x}) * I(\mathbf{x})
 \end{aligned} \tag{1-31}$$

where

$$\nabla G_\sigma(\mathbf{x}) = \left(\frac{\partial G_\sigma}{\partial x}, \frac{\partial G_\sigma}{\partial y} \right) = [-x - y] \frac{1}{\sigma^2} \exp\left(\frac{-(x^2 + y^2)}{2\sigma^2}\right)$$

The ubiquitous edge detection algorithm was proposed by John Canny in 1986 [?], and utilizes a five-step process. First, a Gaussian kernel is applied as a low-pass filter (??). Second, directional filters are used to find the gradients in each direction of the image. Third, a gradient magnitude threshold is applied to remove noise. Fourth, a double threshold is applied to remove both strong and weak edges. Last, edges are determined from hysteresis. The prevailing limitation of this algorithm is the need to set kernel size and edge-intensity.

1.1.3.3 Binary Image Processing

A binary image is a digital image that consists only of black and white pixels. It is often used for labeling or masking an underlying image, where the values of 1 and 0 represent the presence or absence of a particular feature or object. Binary images are often used in computer vision and image processing applications due to low computational overhead and quick analysis. They are also useful for storing and transmitting large amounts of data, as the use of only two values reduces the amount of information that needs to be stored and transmitted.

The primary form of processing on binary images is morphological, that is, changing the shape of the “blob” in order to extract useful information from it.

Dilation and erosion are the two main operations that are used in model-image registration (??). These functions are each two-fold: first, a convolution operation is applied to the existing binary image, then a threshold is applied to the convolution output to determine if the central pixel is a 0 or 1. If f is the input image, s is the convolution

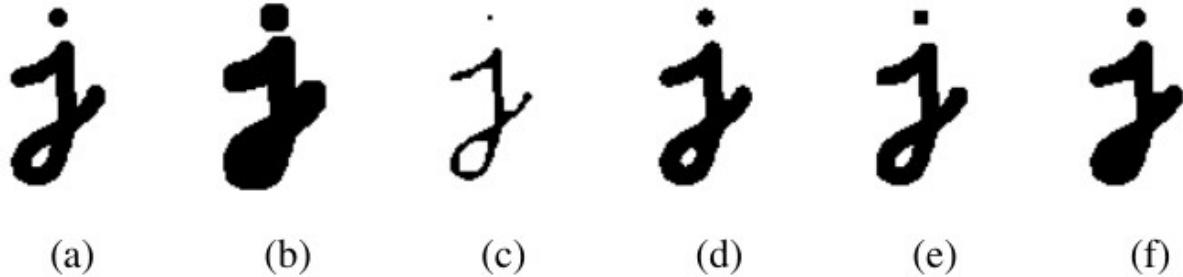


Figure 1-2. A collection of morphological operations on a binary image: (a) original image; (b) dilation; (c) erosion; (d) majority; (e) opening; (f) closing. Image from [?]

kernel of 1s, and $c = f \otimes s$ is the number of 1s in the convolution output, then dilation and erosion can be expressed in the following way.

$$\begin{aligned} \text{dilate}(f, s) &= \theta(c, 1) \\ \text{erode}(f, s) &= \theta(c, S) \end{aligned} \tag{1-32}$$

Where θ represents a thresholding function.

$$\theta(f, t) = \begin{cases} 1 & \text{if } f \geq t \\ 0 & \text{else} \end{cases} \tag{1-33}$$

1.1.4 Image Similarity Metrics

One of the key components in model image registration is image similarity. Fundamentally, this is the method of determining how well the user's synthetic image matches with the actual fluoroscopic image. The choice of similarity metric is going to be determined by many key factors such as the a-priori availability of implant/bone geometry and knowledge of the image quality and contrast. Broadly, there are two classes of image similarity when performing model-image registration: intensity-based and feature-based.

1.1.4.1 Intensity Based

Intensity based measures are those that utilize specific pixel information in order to determine the difference between two images. This can be either a global image similarity

metric, or a measure of the specific regions of interest in the given image.

A canonical difference between two images would be the p-norm separating them (Eq. ??), which iterates through each pixel of the two images and finds the p-norm difference between each pixel pair. Common p-norms are the L_1 norm (*absolute intensity differences* or *mean absolute difference*) [?] ($p = 1$) and the L_2 norm, or Euclidean norm (*squared intensity differences* or *mean squared difference*) [?]($p = 2$).

$$\|A - B\|_p = \left(\sum_{x=0}^w \sum_{y=0}^h |a_{xy} - b_{xy}|^p \right)^{\frac{1}{p}} \quad (1-34)$$

In Equation ??, A and B are the two images being compared, w and h are the width and height of the images, and a_{xy} and b_{xy} are the intensity values at pixel (x, y) in the two images, respectively.

While conceptually easy to use, the main limitation of p-norm measures is their lack of spatial information. For example, an image that has been shifted by a linear transformation would not score well using a p-norm, despite the two images containing only a minor shift, scale, or rotation. One method for overcoming this limitation is to use the cross-correlation, or sliding dot product, between images [? ?] (Eq. ??). When used in conjunction with projective geometry, this can help locate regions of interest for a model-based registration pipeline. The cross-correlation is calculated using the following equation:

$$\begin{aligned} (A \star B)[x, y] &= E[A_{xy} \cdot B_{x+\tau_x, y+\tau_y}] \\ &= \sum_{\tau_x=-\infty}^{\infty} \sum_{\tau_y=-\infty}^{\infty} a_{xy} b_{x+\tau_x, y+\tau_y} \end{aligned} \quad (1-35)$$

This will determine the regions of each image that are similar, causing the correlation function to “light up” at those areas in a similar way to the convolutional operation between two images. The normalized cross-correlation can also be used (??), which removes noise coming from each of the original images.

$$\text{normalized cross correlation}(A, B) = \frac{A \star B}{(A \star A)(B \star B)} \quad (1-36)$$

1.1.4.2 Feature Based

Feature based image similarity metrics involve some method of determining key features in images, and using those notable features for measuring the differences between two images. These types of methods almost always involve some type of feature-extraction step, where the various features of interest are calculated and determined for subsequent use. The two main classes of features are *keypoints* and *edges*. The simplest method of keypoint detection is using a similar method to intensity-based matching, but having one of the “images” as a patch of the desired feature. With keypoints detected in the input image, one could determine the error of the current pose estimate by taking the Euclidean distance between all image keypoints and all projected keypoints: [?] (Eq. ??). With a-priori information about the keypoints, one could attach a weight to every keypoint in order to emphasize specific regions on the image and the model (Eq. ??)

$$\text{Keypoint Error} = \left(\sum_{i=0}^N (KP_{image,i} - KP_{proj,i})^2 \right)^{\frac{1}{2}} \quad (1-37)$$

$$\text{Weighted Keypoint Error} = \left(\sum_{i=0}^N w_i (KP_{image,i} - KP_{proj,i})^2 \right)^{\frac{1}{2}} \quad (1-38)$$

Keypoints are particularly useful when there are invariant features in images and 3D models, like morphological aspects of bones. However, if these features will not, or cannot always be detected, then other measures must be utilized.

Edges as Features Edges are a natural choice of feature when determining image similarity. Similar to intensity-based image similarity, the similarity between the contours of two images offers a promising metric for determining the overall similarity between two images. In model-image registration, the contours of the input image and the projected

model can easily be compared, which presents a reliable scheme for measuring their similarity. When the edges are aligned, you can say that the model is *properly registered* to the image. However, how can we determine when the edges are aligned?

As always, the simplest approach is to take the p-norm between the model and image contours (??), where instead of taking the difference between the two original images, one is taking the difference of the edges of the images. This function will be minimized when there is complete overlap between image and model contours.

1.1.4.2.1. Drawback of Feature Based Similarity Metrics While they can be much more informative and specific, the main drawback of feature based image similarity metrics is the need for feature selection and extraction from the input image. Historically, this required an immense amount of domain-specific knowledge, and the image processing used to extract those features was often relatively limited in scope.

1.2 Deep Learning for Image Processing in Orthopaedics

One of the biggest limitations in historic forms of image processing and filtering is the need to create and tune the filters that are used for feature extraction. While this gives researchers an exceptional amount of control over the filtering techniques, this is often time consuming and difficult to generalize to new data. Deep learning and neural networks offer algorithms for which feature extraction is automated, and techniques have been introduced that make these algorithms robust to new data.

1.2.1 The Purpose of Neural Networks for Image Processing

Neural networks for image processing are an attempt to recreate the visual system in animals, creating algorithmic analogies to the various physical neuronal pathways that are present in the image. Typically, there are many different parts of our visual system that we take for granted, such as the ability to see upside down and extrapolate from prior information. While seemingly intuitive, these are difficult feats to determine algorithmically in a mathematical function.

Using neural networks for image processing is an attempt to re-create the visual system by utilizing algorithmic analogies to the human brain. Subjectively, we take many different aspects of our visual system for granted, like the ability to extrapolate to new information, and understand intuitively the notions of size and scale not affecting the content of an image. However, these are extremely difficult to perform with hand-selected features. Neural networks remove the need to feature selection in images, which has made them a powerful tool in the image processing toolbox.

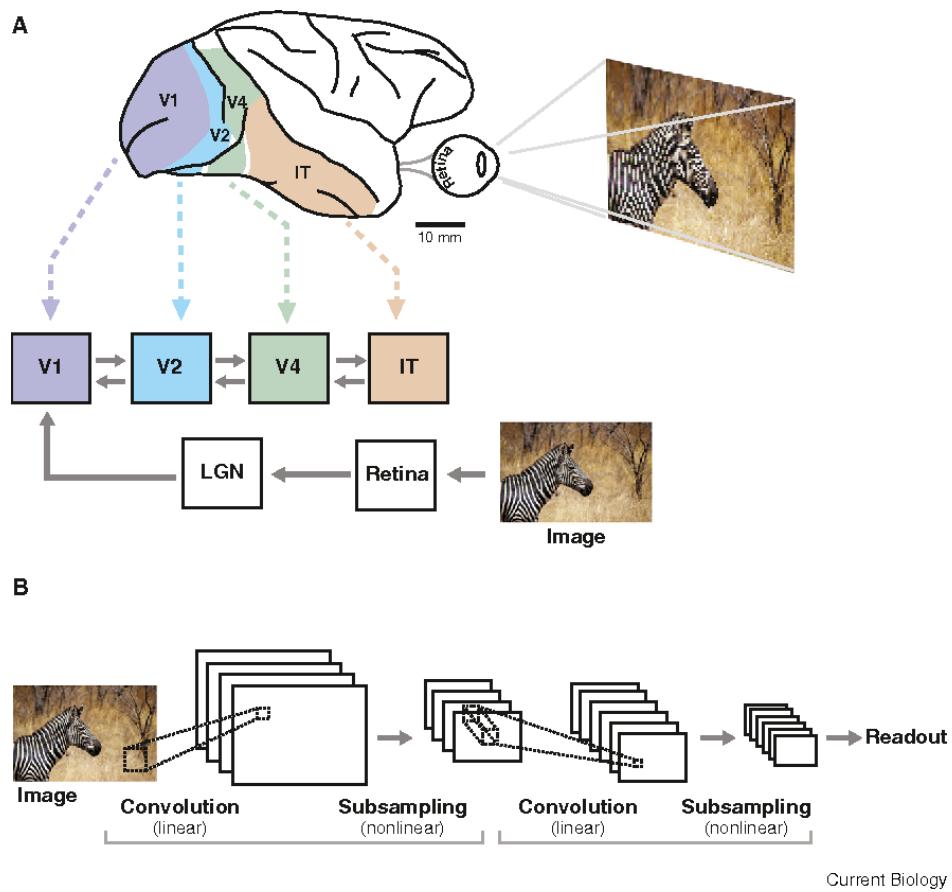


Figure 1-3. An example of a neural network mimicking the human brain.

1.2.2 Neural Network Structure

Neural networks generally have the same constitutive elements, mixed and matched based on the desired performance and complexity of the model that you are trying to build.

1.2.2.1 Neural Network Building Blocks

Generally, neural networks are formed by collections of foundational units, which can generate increasingly complex architectures and yield incredible performance. However, it is always important to start with the fundamental “atoms” of the neural network.

The most basic unit of a neural network is the perceptron (or neuron), which is composed of a summation of inputs multiplied by weights, a bias term, and a (typically non-linear) activation function (??, ??).

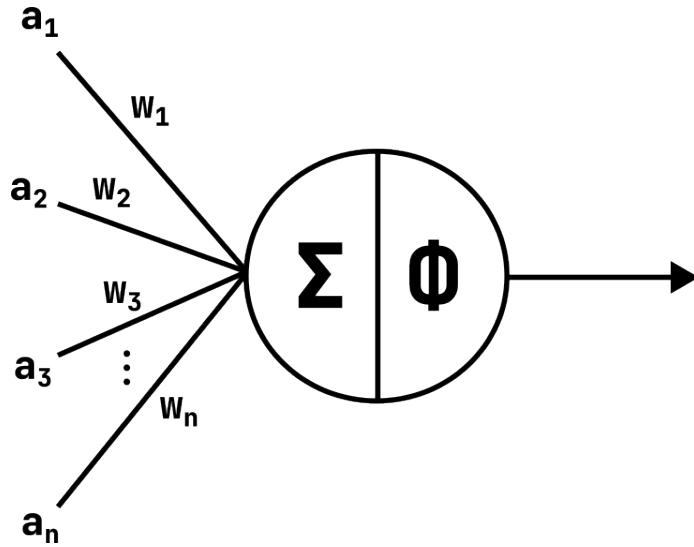


Figure 1-4. A schematic representing a single neuron that receives n inputs and applies ϕ as an activation function.

$$y = \phi\left(\sum_{i=1}^n a_i w_i + b\right) \quad (1-39)$$

Activation functions are a crucial component in neural networks as they allow for the introduction of non-linearity, which is essential for the network’s ability to learn complex representations of the input data. Without activation functions, a neural network would only be able to learn linear relationships between the input and output. However, many real-world problems involve non-linear relationships that cannot be captured by a linear model alone. Activation functions provide a way to move beyond a linear relationship,

allowing the neural network to learn more nuanced mappings between the input and output. The choice of activation function depends on the specific problem you are trying to solve and the architecture of your network. Some activation functions introduce more non-linearity than others and some trade-off between non-linearity and computational efficiency. Experimenting with different activation functions and observing the impact on the network's performance can be a useful technique for optimizing the performance of a neural network. A list of common activation functions and their equations is shown in Table ??.

Table 1-1. A list of activation functions and their corresponding mathematical formula

Activation Function	Equation
Linear	$\phi(x) = x$
Sigmoid	$\phi(x) = \frac{e^x}{1+e^x}$
Hyperbolic Tangent	$\phi(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
Rectified Linear Unit (ReLU)	$\phi(x) = \max(0, x)$
Leaky ReLU	$\phi(x) = \max(0.1x, x)$

1.2.2.2 Fully Connected Network

The fully connected network, also known as the multi-layer perceptron, is a basic type of neural network that utilizes the neurons previously discussed as building blocks. Its schematic representation is a familiar image to many when considering neural networks. By stacking the summations and multiplications of each neuron, we can derive the equation for a single layer of a fully connected network, which is simply a matrix multiplication (??). In this equation, W represents the collection of weights for each neuron, a represents the input, b represents the bias, and ϕ represents the activation function.

One of the key strengths of this type of network is the utilization of non-linear activation functions. A well-chosen activation function can greatly impact the network's performance and ability to achieve specific tasks. For example, in a binary classification task, the sigmoid activation function can be used at the output layer, constraining the

output between 0 (false) and 1 (true). The output can then represent the probability of the given input being classified as 'true'. The ReLU activation function is another popular choice because it is computationally efficient and often used in hidden layers of deep networks. Additionally, the hyperbolic tangent (tanh) activation function is used in the context of a model where data follows a Gaussian distribution.

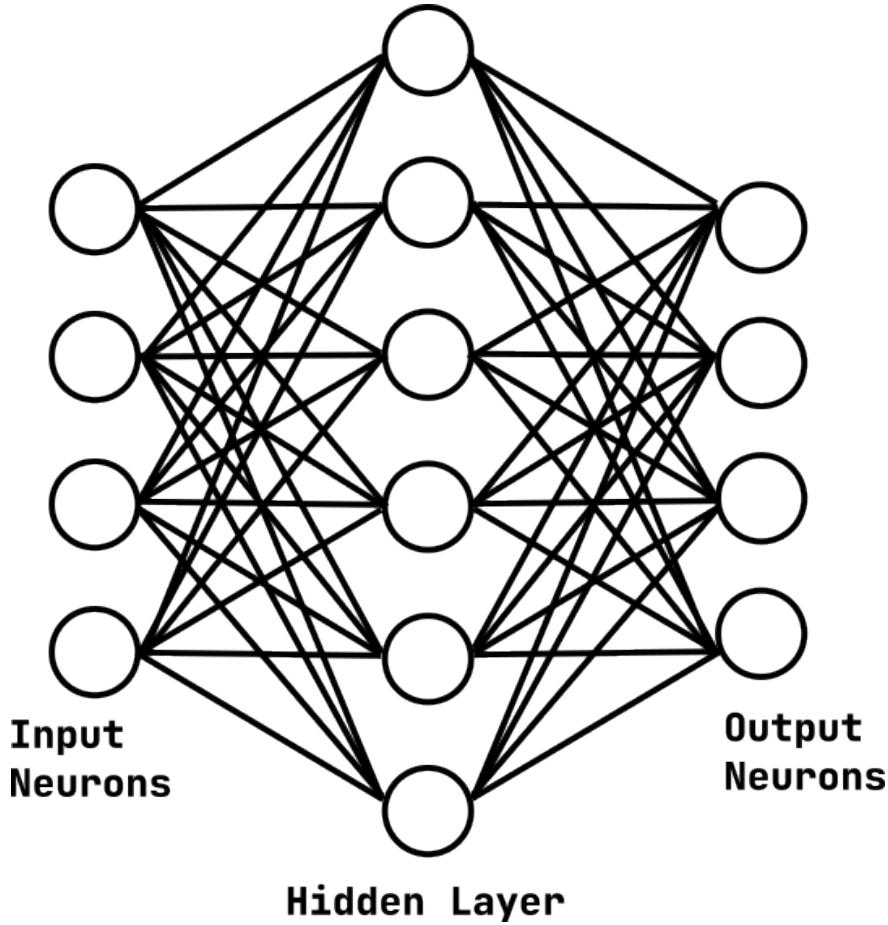


Figure 1-5. A basic fully connected network with a single hidden layer. Each of the neurons are exactly the same as the neurons shown in figure ??

$$y = \phi(W^T a + b) \quad (1-40)$$

Additional layers can be added by taking the output of one layer and sending that into the input of the next. Increasingly complex mappings can be generated by increasing either the size or the number of hidden layers in a network.

$$y = \phi_2(W_2^T(\phi_1(W_1^T a + b_1) + b_2)) \quad (1-41)$$

One of the main limitations of fully-connected networks is exponential increase in computational complexity as your input grows. For a standard 1024×1024 image, you have roughly 1 million input nodes, which can lead to hundreds of millions of parameters that need to be learned depending on the depth of the network. We can use standard image processing techniques to overcome some of these limitations as we explore different network architectures.

1.2.2.3 Convolutional Neural Networks

Alex Krizhevsky sparked renewed interest in deep learning in 2012 by utilizing a convolutional neural network to win the ImageNet challenge [?] by more than 10 points over second place. Since then, neural networks have found their way into many different computer vision tasks, including those in the medical field. The power of convolutional neural networks lay in their ability to autonomously extract latent features from images useful for many processing and analysis tasks. The medical field has used them to segment and classify different bones, structures, and pathologies from a wide array of imaging modalities (CT, x-ray, MRI, etc). In most cases, it can completely remove the need for human supervision in many repetitive image processing tasks.

Convolutional neural networks (CNN) utilize the convolution operation (??) in order to both reduce the size and complexity of the network and capture spatial information present in images. In the same way that the matrix W in the fully connected network is a learnable parameter in a FCN, the individual kernels are the learnable parameters in a convolution operation. In practice, with the correct cost function and optimization routine, this allows each of the kernels to learn the latent structure in the image and make connections between those structures. At a high level, these kernels often represent feature extractors for edges, lines, corners, curves, and other geometric primitives. However, as one traverses deeper into the network, the combinations of these features are often

incomprehensible to a human (??).

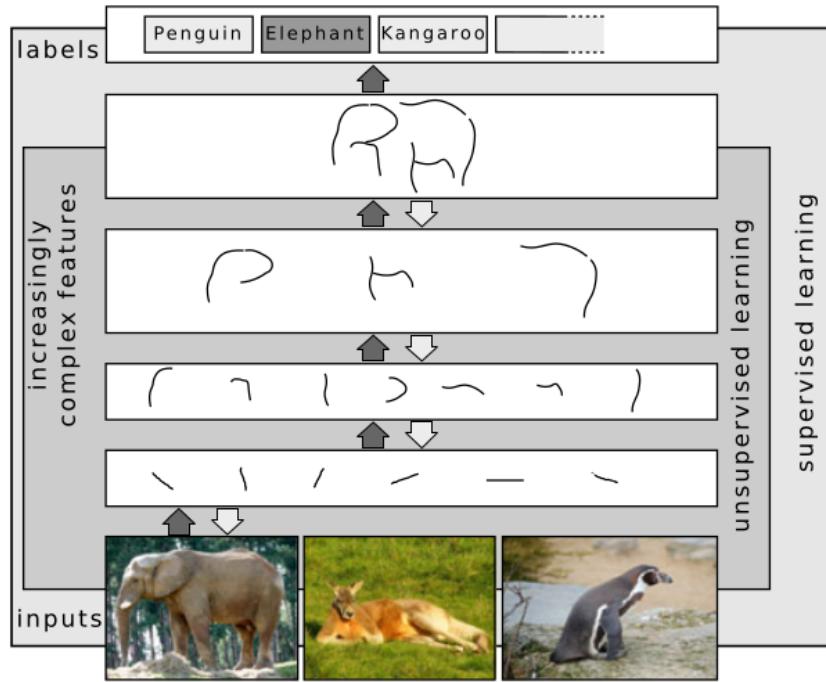


Figure 1-6. An example of extracted features from a convolutional neural network. As shown, the deeper into the network one explores, the combination of core features start to create higher level features that might represent the shape of a specific animal [?].

Typically, these networks will downsample the image to capture the most salient information, then upsample to regenerate the features from the underlying latent representations. This network architecture has been popularized by the U-Net, which is known as a standard autoencoder (Fig. ??). Each of the layers is composed of a collection of convolution kernels, and the downsampling occurs based on the stride or size of the kernel.

CNN architecture can be altered by changing the behavior and structure of the underlying kernels, namely size, stride, and shape. Stride controls the discrete steps the kernel takes as it is moving across the input image and can be used to downsample more aggressively. An atrous convolution involves internally padding the convolution entries with zeros. This also has the effect of more aggressively downsizing an image and capturing

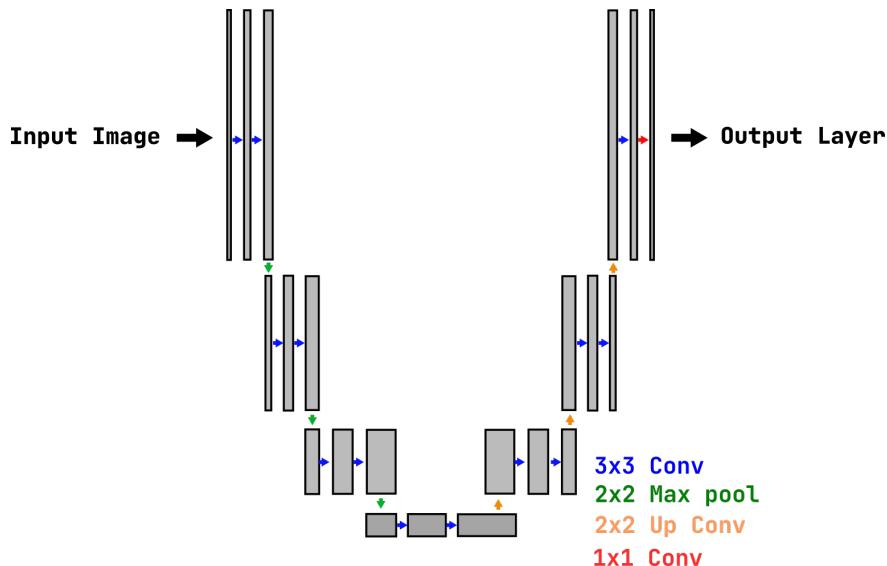


Figure 1-7. A standard U-Net architecture for a convolutional neural network. This architecture is the most common form of network used for biomedical image analysis and processing.

a larger region around the center pixel.

A convolutional neural network can have an additional fully connected network appended to the output in order to learn a mapping that involves a linear combination of the extracted features from the convolution operations. This is often used when a CNN is applied to a classification task, where the total number of outputs of the final FCN is the number of classes.

1.2.3 Training Neural Networks

The power of neural networks is that they are able to “learn” incredibly complex mappings from a set of training data. But how *exactly* do they learn?

1.2.3.1 Neural Network Cost Functions

The first step in updating a machine learning pipeline is to understand the metric you are trying to minimize (or maximize). This can encode a specific understanding about the task at hand, and it can also be used to promote a specific structure in the network itself (parsimony, sparsity, etc).

The **cost function** is exactly the metric that the learning algorithm is attempting to

minimize. The application of different cost functions for different tasks in neural networks has been well studied. For regression, typically some type of distance metric is employed (mean absolute error, mean squared error, Kullback-Liebler distance), while classification tasks mostly utilize some type of cross-entropy or log likelihood metric [?].

1.2.3.2 Optimizing and Updating Weights

For convex problems, there often exists a closed-form solution for the parameters that minimize a given objective function. However, neural networks are highly non-linear and non-convex, forcing researchers to employ different methods of updating parameter values in a direction of minimizing the objective function. Backpropagation has become the ubiquitous choice for updating these highly non-linear systems. [?]. This utilizes gradient descent (??) and the chain rule (??) in order to calculate the affect that each node has on the final output, and update it in the direction that maximally minimizes the current error. Researchers can control how aggressive the update is by changing the learning rate, η .

$$w^{(j+1)} = w^{(j)} + \Delta w \quad \text{where} \quad (1-42)$$

$$\Delta w = -\eta \nabla J(w^{(j)})$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial e} \frac{\partial e}{\partial \phi} \frac{\partial \phi}{\partial v} \frac{\partial v}{\partial w} \quad (1-43)$$

One of the main limitations of gradient descent and backpropagation is the need to tune the learning rate. Too small a value will cause a network to train extremely slowly and almost always get stuck in local minima. Too large a learning rate can cause the network to “bounce out” of the global minima due to the rate of change being too large. Hyperparameter tuning can be a difficult process without a refined methodology beyond trial and error. Some groups have proposed different methods of incorporating dynamically changing update rules in order to incorporate physical properties into the network training stage. The most common is Adaptive Moment Estimation [?], which incorporates Root

Mean Squared Propagation [?] and momentum learning.

1.3 Model-Image Registration for Measuring Kinematics from Fluoroscopic Images

3D/2D model-image registration utilizes the computer vision principles discussed previously to determine the position and orientation of the model, given an image containing that model. Relying on the idea that the projective scheme of the real-world camera can be emulated, then the goal is to find some type of similarity metric (??) between a projected 3D model (??, ??) and the existing image such that minimizing this metric indicates that our object has been “placed” correctly. Many different branches of computer vision and optimization have been explored for determining the kinematics of total knee arthroplasty components, each with varying levels of computational intensity, objective function (via image similarity), and optimization routine. General groupings of each will be discussed here, along with pitfalls and limitations that prevent the technique from being applicable in a clinical setting.

1.3.1 Pre-Computed Projective Geometries

The earliest methods of model-image registration had to deal with many of the limiting factors of computational availability at the time. This forced researchers to find clever ways to determine image similarity metrics without the need to iteratively compute the projective geometry of the model thousands of times per second.

First established in the early 80s [? ?], normalized Fourier descriptors provide a way to normalize 2-dimensional images using information from the latent 3D characteristics (position and orientation). This was used to determine TKA kinematics to high levels of accuracy [? ?], so long as 3D shape information was known, and the camera matrix could be deduced. The image similarity metric utilized the L_2 difference between the normalized fourier descriptors of the input image and the precomputed shape library at known rotations. Further interpolation was used to increase accuracy significantly.

In parallel, another group utilized pre-computed distance maps intrinsic to the 3D model [? ?]. These distance maps could then be used to quickly determine the Euclidean

distance between any node in the model and an arbitrary line in 3D space. Then, 3D vectors were created starting at sampled points along the contour of the implant, and concluding at the origin of the camera. Assuming an accurate focal distance, then the objective is simply minimizing the distance between all the vectors and their corresponding nearest node on the 3D model. Once minimized, the 3D object fell into the “slot” carved out for it by the pseudo-conical shape.

The main issue with these pre-computed geometries is need for the researchers to hand-select the contours belonging to the implant. Despite the recently available Canny edge-detector [?], one still needed to hand-label the specific edges of interest. This would be far too time consuming in a clinical setting.

1.3.2 Motion Capture

Motion capture is a common method in computer vision involving the tracking of reflective markers using multiple cameras to over-determine the location of the markers. Then, if the position of the markers relative to anatomic features is known, a series of transformation matrices can be used to determine the location of the anatomic features, given the measurable location of the markers. Historically, this has taken two different approaches, skin-mounted and bone-mounted.

Skin mounted markers rely on a set of fiducial points on the body that are closely related to joint centers, such that relative translations and rotations can be resolved. However, there have been numerous studies showing that skin-motion during activity leads to high levels of inaccuracy at the joint [? ? ? ?]. These inaccuracies preclude this method from being a clinically viable option for measuring joint kinematics.

Other groups have taken a much more accurate approach: drilling bone pins into the joints of interest and attaching motion-capture reflectors to the end of the pins [?]. While extremely accurate, this method is far too invasive and time consuming for a clinical setting.

1.3.3 Iterative Projections

Increased computational performance and widespread availability allowed for different types of methods of model-image registration that were previously untenable due to computational complexity. One of the main areas that saw increased speed-ups were in parallel computing, which drastically increased computer graphics performance. By leverage improved computer graphics pipelines, researchers were able to project and render 3D implants in near real-time and perform image similarity metrics quickly. This seemingly overcame the need for pre-computed projection schemes.

Different groups have achieved strong performance using simulated annealing with direct image-to-image similarity metrics [?], and Lipschitzian optimization using contour matching [?]. While powerful and quick, these methods still rely on human supervision to escape local minima and set an initial pose estimate. Noisy images and the need to set hyperparameters for the edge detection scheme also reduce the scope of generalization to new data in a clinically practical way.

1.3.4 Fully Human Supervised

One of the more dominant softwares in measuring joint kinematics is JointTrack, which utilizes the strongest neural network available (the human brain) connected to one of the most dexterous manipulators available (the human hand) to make accurate measurements of joint kinematics from single-plane images [?]. This software works by accurately recreating the fluoroscopic system's projective matrix and allowing users to manipulate a 3D model of the desired bone or implant to align it with the provided image.

There were additional views that allowed the user to see the alignment from the coronal plane, as well as graphs demonstrating the kinematics throughout the movement, which allowed abnormal frames to be identified and dealt with quickly.

Hundreds of papers have been published using this software as the method for determining kinematics, and it has been extensively validated using many different methods offering ground-truth solutions to kinematics. The main issue with this method is

both the upfront time to train a user, and the necessity for human-supervision during the entire measurement process.

1.3.5 Biplane Kinematics Measurements

One of the most effective methods for resolving single-camera limitations in measuring kinematics is to add a second camera! This offers much greater accuracy and resolution, especially when determining out-of-plane translations, because the out of plane translation for one camera is an in-plane translation for the second.

The groups that have used this have cited sub-mm and sub-degree accuracy for all translations and rotations [? ? ?], and have used a wide variety of optimization routines and image similarity metrics.

While this seems extremely promising, the general cost and unavailability of bi-plane fluoroscopic imaging systems presents a problem for integrating this technology into a clinical setting. If kinematic analysis is going to be have widespread clinical adoption, it **must** be integrating into the imaging systems present at most hospitals and clinics, which is single-plane systems. However, the accuracy of bi-plane systems can be used to validation of performance for various single-plane pipelines [?].

1.3.6 Roentgen Stereophotogrammetric Analysis

Roentgen Stereophotogrammetry is one of the most accurate methods of measuring implant kinematics, popularized in the early 1970s [?]. Most often, it is used to determine implant micromotion in post-operative followups due to the extremely high resolution at which it measures pose and orientation. This method utilizes highly radio-opaque tantalum beads of known shape and size implanted on and adjacent to the implants. Utilizing similar methods to motion capture, algorithms can resolve the location of the tantalum beads extremely accurately, providing a good estimate for the location of the implant to which they are attached.

The largest limitation of this method is two-fold (1) the need for additional surgical steps precludes this from being a widely applied method, and (2) the need for biplane

imaging system poses a problem to availability and portability.

1.4 Automating Measurements of Kinematics from Fluoroscopic Images

The proposed method overcomes the various limitations of previous attempts to autonomously measure kinematics from single plane fluoroscopy. The key feature of all these limitations is that they prevent the adoption of this technique in a clinical setting, due to the labor overhead or equipment required in order to generate an accurate kinematic report. The proposed combination of methods will alleviate these requirements, and the extensibility of the software will allow this technology to be used in a clinical setting without disrupting the standard clinical workflow.

1.4.1 Autonomous Implant Detection

Determining the location of an object in an image is a historically intractable problem. These are one of those tasks that are often relegated to the corner of “easy for humans, extremely difficult for computers”, especially when there is very little a-priori information available. However, as discussed in ??, deep learning has paved the way for computer vision programs to be able to perform tasks that were once only possible by humans. Two convolutional neural networks were trained to segment the tibial and femoral components from the single plane fluoroscopic images. The network architecture used was the High-Resolution Net [?], which leverages low-level features with higher resolution parallel processing in order to better determine the spatial characteristics of the image and produce a better output. At the time of writing, this network sets the state-of-the-art standard for performance on the COCO and ImageNet datasets, demonstrating robustness for use in many different arenas.

Many of the historic methods of determining kinematics were limited by the researchers ability to quickly and reliably determine the location in the implant. The contours were either hand-segmented [? ?], or a normal edge detector was used, which introduces extra tuning parameters for any given study due to variations in image quality.

1.4.1.1 Neural Network Robustness

One of the main problems with neural networks is overfitting. With millions of parameters to tune, it can be extremely easy to “overfit” on your training set, leading to the network’s inability to perform well on any image that was not directly in the training set. When dealing with fluoroscopic images, this might look like a neural network that can perform extremely well on high-quality, high frame rate, low blur images from a hospital that has a budget to support such a machine, while failing to segment images from a machine more than a decade old. We overcame this challenge through a mixture of additional image augmentations [?] and using a wide range of training data. The neural networks were trained on roughly 8000 images from 7 different human-supervised total knee kinematics studies spanning the last two decades. The image qualities range from extremely clear and high quality to nearly indiscernible without intense human supervision. The goal of this two-pronged approach was to have both artificial and real “low quality” images for the network to train on so that any hospital or researcher, regardless of the available equipment, might be able to leverage this technology in their practice. The authors hope that this approach provides equal access to this informative measurement.

1.4.2 Initial Pose Estimation

Hand-in-hand with implant detection is the initial pose estimate that very often needs to be input into the optimization routine. Once the contour of the implant was determined, many methods required a human operator to place the implant in the “capture region” of their optimizer in order for it to eventually find the correct solution. This takes human supervision to get correct, thus adding another impediment toward getting this technology into the clinic.

To determine an initial estimate, we must rely on those methods that can leverage information present in the camera projection (??) and the 2D CNN output (??). The primary method that takes advantage of this information utilizing normalized fourier descriptor shape libraries, and extracting each of the 6 degrees-of-freedom from either the

normalized coefficients or matching with the closest entry in the library [? ? ? ?]. The key feature that makes this method tractable for autonomous measurements is the availability of the implant pixels from the convolutional neural network.

1.4.2.1 Generating Normalized Fourier Descriptors

The Fourier Transform is one that takes in a continuous and repeating function and represents it as a sum of sinusoidal signals. Because the implant contour is self intersection, we can view it as a continuous function that has period 2π radians, as it will loop back onto itself and start over. This allows us to take advantage of the Fast Fourier Transform (FFT) [?], which operates on a discrete set of points rather than a continuous function. First, the contour of the segmented implant is taken then resampled into 128 equi-spaced points. The choice of 2^n points allows the FFT algorithm to perform much more quickly than another choice of points. Each contour point on the image (x, y) is then represented as a single complex point, $x = jy$, such that the 1D FFT algorithm can be used. If $s(n)$ represents the complex sequence of equi-spaced contour points in the implant, and $S(i)$ represents the frequency-domain representation of those points after the FFT is applied, then we can represent these functions as shown in ??.

$$S(i) = \sum_{n=-\frac{N}{2}+1}^{\frac{N}{2}} s(n)e^{-j(2\pi in)/N}$$

$$\text{for } -\frac{N}{2} + 1 \leq i \leq \frac{N}{2} \quad (1-44)$$

$$s(n) = \frac{1}{N} \sum_{i=\frac{N}{2}+1}^{\frac{N}{2}} S(i)e^{j(2\pi in)/N}$$

$$\text{for } -\frac{N}{2} + 1 \leq n \leq \frac{N}{2}$$

As discussed in ??, spatial information in images is difficult for computers to understand without explicit calculation. So, we use the properties of the FFT to normalize each of the shapes based on location, rotation, and size in order to accurately compare the segmented mask to a generated library.

Normalize Position By the properties of the FFT, we know that $S(0)$ is the geometric centroid of the contour. Thus, we can set this to zero for all contours to give a consistent reference point that is independent of the location of the input contour. We can save this value as the “position normalization coefficient” to determine the (x, y) location of the implant later in the process. And, because of our usage of the 1D FFT, we know that the real portion of this coefficient is the x-value and the imaginary portion is the y-value.

$$\begin{aligned} \text{Position Coefficient} &\leftarrow S(0) \\ S(0) &= 0 \end{aligned} \tag{1-45}$$

Normalize Size Because the implant is not self-intersecting, we know that $S(1)$ is the coefficient with the largest size, and it also represent the scale of the shape of the implant. So, we can normalize each shape by dividing each coefficient by the overall size of the contour, shown below.

$$\begin{aligned} \text{Size Coefficient} &\leftarrow |S(1)| \\ S(i) &= \frac{S(i)}{|S(1)|} \\ \text{for } -\frac{N}{2} + 1 < i < \frac{N}{2} \end{aligned} \tag{1-46}$$

Normalize In-Plane Rotation and Contour Starting Point One of the most difficult parts of measuring the similarity between two contours, especially when they are composed of a set of discrete points, is that any distance measurement necessarily takes into account those discrete points in the order that they are presented. To illustrate this example, imagine two squares, each defined by the location of the corners A_i, B_i, C_i, D_i . If these two squares are perfectly overlapping, one might imagine that the Euclidean distance between each of the points, $(\sum_{P \in A, B, C, D} (P_1 - P_2)^2)^{\frac{1}{2}}$ would be equal to zero. This would only be true if the starting point of the contour was the same for each square (e.g. Corner A was always the top left corner). If each square had A starting in different locations, then

even when the contours are perfectly aligned, the distance metric would be non-zero and uninformative.

We can use properties of the FFT to normalize the starting position of the contour in each of the segmentations as well as the general orientation of the contour. We can leverage the starting point shift property of the fourier transform (??) and the rotation property (??) in order to normalize both of these factors and ensure that similar shapes have the same orientation and starting point.

$$s(n - T) \xleftrightarrow{DFT} S(i)e^{-jiT} \quad (1-47)$$

$$s(n)e^{j\theta} \xleftrightarrow{DFT} S(i)e^{j\theta} \quad (1-48)$$

To normalize the starting point and rotation, we find k such that $S(k)$ is the coefficient of second largest magnitude. We then apply a mixture of the starting point shift and the rotation property of the FFT to orient each contour (??). We also want to find the “rotation normalization coefficient”, which is the angle through which the contour needs to rotate to reach the normalized orientation. This is done by determining the phase of the normalization equation at $i = 0$, which controls the overall orientation of the contour. If u is the phase of $S(1)$, and v is the phase of $S(k)$, then we can find the normalized orientation.

$$\begin{aligned} \text{Rotation Normalization Coefficient} &\leftarrow \frac{v - ku}{k - 1} \\ S(i)_{norm} &= S(i)e^{j\frac{(i-k)u + (1-i)v}{k-1}} \\ \text{for } -\frac{N}{2} + 1 \leq i \leq \frac{N}{2} \end{aligned} \quad (1-49)$$

Once the in-plane rotations have been normalized, the contour has been completely normalized for x, y, z translations and z rotations. And, by storing these values, we are

able to utilize them in determining an initial pose estimate. Then, we can use a library made up of known x, y rotations, and compare the segmentation to this library to determine all 6 degrees-of-freedom.

1.4.2.2 Shape Library

A shape library is created using a flat panel projection (??) of the implant at known x and y rotations, while holding all positions and orientations constant, then applying the normalization protocol described above to determine the normalized coefficients of each library entry. If $s_{x,y}(n)$ is the flat-panel projection of the implant with x and y rotations, then we can generate a library using the following procedure, where FFT is the fast fourier transform (??) and NFD is the process of normalizing the contour and extracting the relevant coefficients (??).

$$S_{x,y}^{lib}(i) = NFD(FFT(s_{x,y}(n))) \quad (1-50)$$

Once the shape library is generated for a specific implant, we can start to determine the pose estimates for each degree of freedom.

1.4.2.3 Generating a Pose Estimate

First, we compare the Euclidean distance of the normalized segmentation contour to each value of the shape library; the x and y rotations that minimize this function are taken as the x and y rotations of the implant (??).

$$(\theta_{x,est}, \theta_{y,est}) = \operatorname{argmin}_{x,y} \left(\sum_{i=-\frac{N}{2}+1}^{\frac{N}{2}} (S^{seg}(i) - S_{x,y}^{lib}(i))^2 \right)^{\frac{1}{2}} \quad (1-51)$$

Then, we can determine the z rotation estimate by comparing the values of the normalized θ values that were needed in ???. This process is shown in ??.

$$\theta_{z,est} = \theta^{seg} - \theta_{x,y}^{lib} \quad (1-52)$$

where $\theta^{seg,lib} \equiv$ Rotation Normalized Coefficient

We can then use the principals of projective geometry (??) and similar triangles to determine the out-of-plane translation of the implant, given that the library was projected at a known depth. Using similar triangles, we are able to determine that the depth is inversely proportional to the overall magnitude of the projection, captured by the “Size Coefficient” (??). We also make the assumption that we have a weak perspective projection, meaning that the out-of-plane translations are small compared to the focal length of the fluoroscopy imaging setup.

$$\begin{aligned} \frac{M^{seg}}{f} &= \frac{h}{z_{est}} \\ \frac{M^{lib}}{f} &= \frac{h}{z_{lib}} \\ \rightarrow \\ z_{est} &= \frac{M^{lib}}{M^{seg}} z_{lib} \end{aligned} \tag{1-53}$$

where

$$M^{seg,lib} \equiv \text{Size Coefficients}$$

$$h \equiv \text{Implant Size}$$

The x and y translations can be determined using the value of the z translation estimate, along with the location of the centroid, saved as the “Position Coefficient”. This is then refined further to account for the rotation of the implant and the distance of the implant’s centroid to its origin. We can express this with a single matrix multiplication multiplied by a scale factor (??).

$$\begin{pmatrix} x_{est} \\ y_{est} \end{pmatrix} = \begin{pmatrix} Re(S^{seg}(0)) \\ Imag(S^{seg}(0)) \end{pmatrix} - \begin{pmatrix} \cos(\theta_z) & -\sin(\theta_z) \\ \sin(\theta_z) & \cos(\theta_z) \end{pmatrix} \begin{pmatrix} Re(S_{x,y}^{lib}(0)) \\ Imag(S_{x,y}^{lib}(0)) \end{pmatrix} \times \begin{pmatrix} z_{est} \\ z_{lib} \end{pmatrix} \tag{1-54}$$

Once this step is complete, we have accounted for x and y rotations (??), z rotations (??), x and y translations (??), and z translations (??). This provides a robust initial

estimate when the only information available is the implant geometry and the segmentation output from the neural network. Furthermore, this can be done without any human supervision, providing a fully autonomous initialization for any pose refinement strategy, so long as the estimate is within the convergence region.

1.4.3 Objective Function

In a perfect situation, our objective function would directly measure the error between our 3D model’s current pose and the true pose of the object. However, if we had a-priori access to the true pose of the object, then this entire pipeline would be worthless. Thus, we must find an objective function that can act as a heuristic for the difference between true pose of our model and the current pose of our model. Our access to the segmentation output from the CNN and the ability to project the silhouette of our model quickly makes contour comparison a natural choice for an objective function. The only assumptions that we make are (1) our projective algorithm and camera definition are the same as the camera that was used to take the original fluoroscopic image and (2) our 3D model is the same 3D model that is present in the image. If these two assumptions are correct, then the alignment of the image contour and the projected contour necessarily means that our pose is correct, unless you have a symmetric object (??).

First, we apply a Canny edge detector [?] to extract the edges from our segmentation contour, S , and our projected 3D model, P , where edges are 1, and every other pixel is 0. We can then iterate over each pixel and take the absolute values of the L_1 norm between our segmented and projected contours (??).

$$J = \sum_i^{\text{Height}} \sum_j^{\text{Width}} |S_{ij} - P_{ij}| \quad (1-55)$$

Unfortunately, the contours of the projected model are extremely sensitive to pose, especially when representing angles using Euler decomposition. This results in a chaotic similarity function that has an extensive amount of local minima. Past methods have overcome this by dilating the contour of the projected image (??) and performing the same

L_1 optimization routine. However, a lack of an isolated contour for the fluoroscopic image still lead to a slightly noisy objective function. Our proposed method takes advantage of the segmentation output for the neural network and dilates both the segmentation contour and the projected contour for a much smoother objective function allowing for a wider search range. As our objective function is minimized, we can decrease the level of dilation to return the metric back into its original form, which most accurately describes the difference between the projection and image.

1.4.4 Optimization Routine

Broadly, optimization is the process of minimizing or maximizing an objective function, $f(\mathbb{R}^n) \rightarrow \mathbb{R}$, potentially subject to some constraints (e.g. $x \in \Omega$) [?]. We formalize this below (??). The simplest optimization problems have an analytic form of the gradient of f that can be solved directly, typically by setting the first derivative to zero (e.g. least squares).

$$\operatorname{argmin}_x \{f(x) : x \in \Omega\} \quad (1-56)$$

A drawback to our pipeline is that there is no analytic form of the objective function between each segmentation and projected contour; they must be resampled over a specified range in order to approximate objective function values and gradients. This defines a *black box* optimization routine, which is well studied in the literature [?]. In this type of function, it is not possible to use gradient-based methods to determine a minimum value for the objective function, one must use heuristics or ad-hoc methods to find the minimum location. Lipschitzian optimization offers an appealing black box optimization approach because it satisfies our need for a global search algorithm with provable convergence. First, we start with the definition of Lipschitz continuous, which places bounds on the rate of change of a function specified by some constant, called the Lipschitz constant. With a known Lipshitz constant, is is possible to find the value for the global minimum of optimization function [?].

Definition 1-1 (Lipschitz Continuous). *The function g is said to be Lipschitz Continuous on the set $\mathbf{X} \in \mathbb{R}^n$ if and only if there exists a scalar $K > 0$ for which*

$$\|g(x) - g(y)\| \leq K\|x - y\|$$

for all $x, y \in \mathbf{X}$

The scalar K is called the Lipschitz Constant of g relative to the set \mathbf{X} .

We can illustrate this convergence with a simple example: consider the function $f(x) \rightarrow \mathbb{R}, x \in [a, b]$. If we know that the function is Lipschitz continuous, the following conditions are true on the domain $x \in [a, b]$. This corresponds to the positive and negative slopes, K , applied to the extrema of the domain, and the intersection, x is selected as the choice for subdivision. This process is repeated, where the region is further subdivided based on the lowest value of $f(x_i)$. The iterative process is stopped once the difference between successive domain splitting is below a pre-specified global tolerance.

$$\begin{aligned} f(x) &\geq f(a) - K(x - a) \\ f(x) &\geq f(b) + K(x - b) \end{aligned} \tag{1-57}$$

While powerful, a-priori knowledge of the Lipschitz constant is needed for determining the global minima. Without it, there is no way of determining intersection points, and no way of selected new regions for subdivision and sampling. Shubert's Algorithm also has slow convergence, due to the inability to define parameters for when to explore local vs. global search. The Lipschitz constant, K , acts as a weight that places larger emphasis on global serach when high, and local search when low.

Fortunately, methods exist for utilizing the power of Lipschitzian optimization without the need for explicit knowledge of the Lipschitz constant [?]. These can both overcome the need for an a-priori knowledge of the Lipschitz constant, as well as offer some

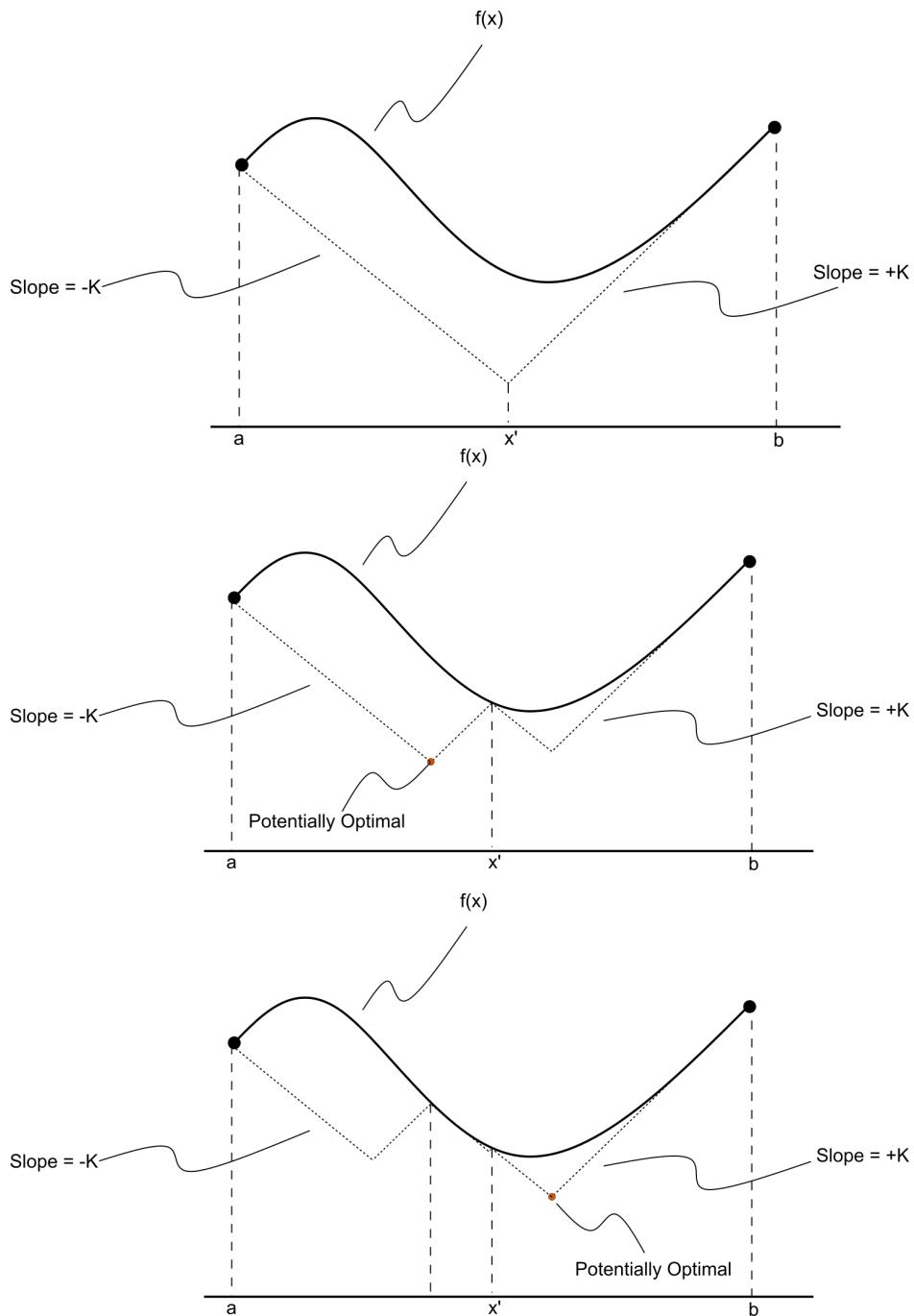


Figure 1-8. A visual representation of Shubert's algorithm, which finds the global minima iteratively through a repeated calculation of the intersection of two lines with slope $\pm K$. If K is known, it will always find the global minimum.

solutions of the slow convergence by providing methods for exploiting both local and global search simultaneously to find the minimum function values.

Jones et al. [?] propose a method by which the center, c , of a domain is sampled, rather than the endpoints. This produces the equations below (??). The inequalities represent slopes $+K$ and $-K$, respectively, and provide a maximum value for the lower bound of the function at the endpoints, a and b . The midpoints of $[a, c]$, and $[c, b]$ are then calculated and the process is then repeated (??). The power of this method is that you can determine potentially optimal regions of the domain by choosing those points along the lower convex hull of the graph plotting sub-domain size vs center point function value. The points along this hull are those that could potentially include the function minimum, and each is chosen for further sub-sampling (??). Determining the convex hull is a problem well studied in the literature [? ? ? ?]. This elegantly mixes local vs. global search, and drastically increases the speed of convergence.

$$f(x) \geq \begin{cases} f(c) + K(x - c) & \text{if } x \leq c \\ f(c) - K(x - c) & \text{if } x \geq c \end{cases} \quad (1-58)$$

This can be extended into multiple dimensions without loss of generality. First, each dimension in the domain is normalized to the range $[0, 1]$, and a hypercube is created in \mathbb{R}^D , where D is the dimension of the domain you are searching. The first iteration trisects this hypercube along an arbitrary dimension, and further iterations trisect along the largest dimension of the hypercube. We select potentially optimal hypercubes by identifying points along the lower convex hull of the graph plotting hypercube size vs center point function value.

For our purposes, we construct a hypercube along each of the 6 degrees-of-freedom that describe the pose of the implant in space, using bounds set by the user. The first epoch involves minimizing the objective function with larger bounds and a larger dilation. After all iterations have been used up, the algorithm is re-started with a smaller dilation

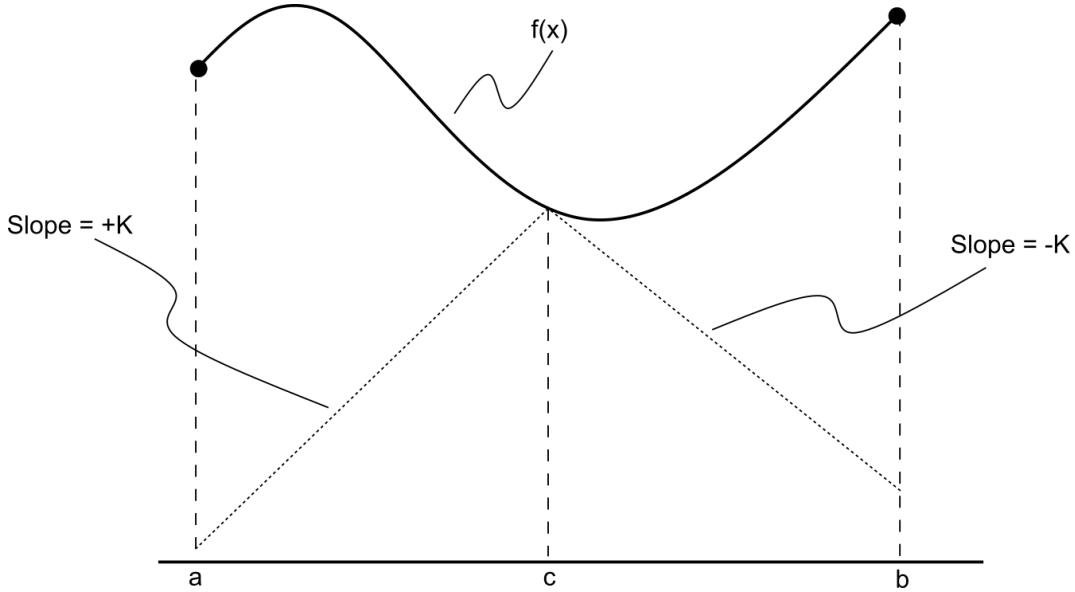


Figure 1-9. The DIviding RECTangles (DIRECT) algorithm in one dimension. It can find the global minimum of a function without a-priori knowledge of the Lipschitz constant. The value of the line with slope $\pm K$ at each of the end-points represents the theoretical minimum for the value of the function in that region. Thus, the size of the region and the value of the function at the center help the algorithm determine potentially optimal sub-regions.

and tighter bounds. The last epoch typically has the tightest bounds and no dilation. This pyramidal scheme offers a smooth objective function when the bounds are largest, which assists in escaping local minima, and an aggressive objective when the bounds are tight, and fewer local minima are present.

1.4.5 Overcoming Single-Plane Limitations

Despite a fully autonomous solution for measuring total knee arthroplasty kinematics, there are some fundamental limitations when using monocular vision to determine the three dimensional position and orientation of an object. ?? discusses methods for overcoming these limitations. So far, incorporating digital ligaments into the cost function has been used to a great deal of success.

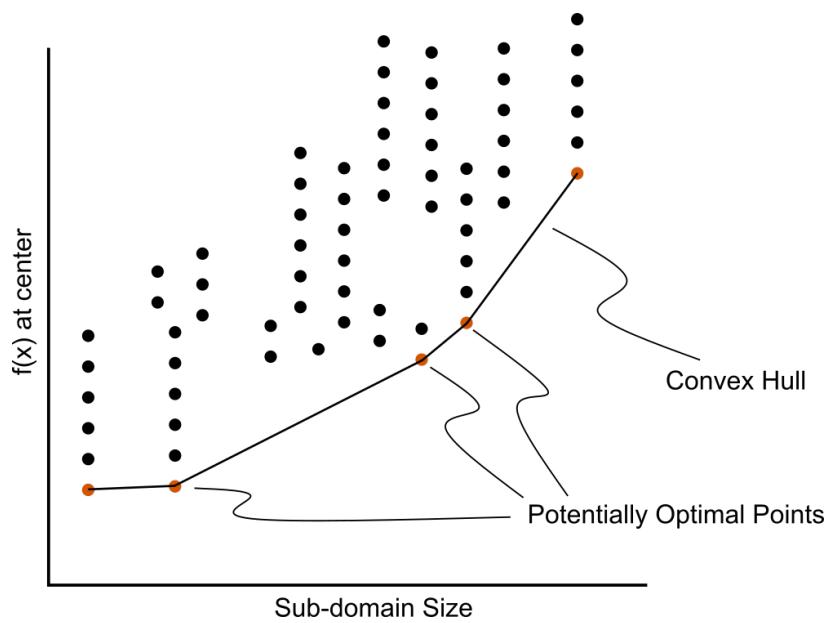


Figure 1-10. The potentially optimal regions of the DIRECT algorithm are those points that lay along the lower convex hull of the scatter plot of sub-domain size vs function value at center. This is due to those regions being the locations where the maximum possible rate of change in the function might be a minimum, without any prior knowledge of the Lipschitz constant.

CHAPTER 2 THE AIMS OF THIS DISSERTATION

2.1 Aim 1: Validating an Autonomous Pipeline for Measuring Total Knee Arthroplasty Kinematics from Single Plane Fluoroscopic Images

The primary aim of this dissertation is to test the feasibility of the fully autonomous pipeline, and validate the results against known gold-standard kinematics measurements to test reproducibility and accuracy. These measures are discussed and quantified in the attached paper (??) The work presented in this paper won the HAP Paul Award for Best Paper at the International Society for Technology in Arthroplasty's 2022 Annual Meeting, and is currently under review for publication in the Journal of Arthroplasty.

2.2 Aim 2: Overcoming Inherent Single Plane Limitations When Measuring Total Knee Arthroplasty Kinematics

While establishing a pipeline for the fully autonomous measurements of TKA kinematics, we encountered many of the different limitations present in using single-plane fluoroscopy. Fundamentally, this is a problem that exists inherently in the system, as you have a severely underconstrained problem, leading us to the inverse problem of computer vision (??).

Definition 2-1 (Inverse Problem). *The inverse problem in computer vision is the process of calculating the causal factors (kinematics) that produced a set of observations (fluoroscopic images).*

However, we are equipped with a-priori information about human anatomy that can dictate a set of rules and procedures to follow in order to overcome some of the different limitations present in using single-plane fluoroscopy.

2.2.1 Depth Perception

One of the most apparent limitations is depth perception. When you only have a single camera to resolve the pose of your object, sensitivity parallel to the focal ray becomes increasingly difficult. However, when dealing with anatomic structures, we know

that there are specific poses that are, at the very least, pathological, and at most, downright impossible. In the objective function used during our black-box optimization, we add linear constraints to the relative mediolateral translation between the two implants, simulating the role of ligaments and soft tissue structures.

2.2.2 Projection Ambiguities and Symmetry Traps

One of the more pernicious limitations in single-plane fluoroscopy is an issue that we've dubbed the "symmetry trap", which causes multiple global minima when using a strictly contour-based objective function. The major contributor to these issues is symmetric tibial implants, which are mediolaterally symmetric (i.e. no different between right and left implants).

Definition 2-2 (Symmetry Trap). *A symmetry trap occurs when a symmetric object has a projective geometry with more than one unique pose that can produce it. The simplest case is a sphere, where all poses produce the same circular projective geometry.*

Based on a weak perspective projection, we can take advantage of the axis-angle rotations (??) to determine the "symmetric pose" for any given orientation of a symmetric object.

Algorithm for Determining the Dual-Pose of a Symmetric Object

1. Determine the viewing ray from camera → object (Eq. ??).
2. Determine the axis-angle (m, θ) rotation between the viewing ray and the symmetric axis of the object (Eq. ??, ??).
3. Rotate the object -2θ about the same axis, reflecting the rotation about the viewing ray (??).
4. The final orientation of the object is exactly the "dual pose", producing indistin-

guishable projective geometry (Eq. ??).

If T is the homogeneous transformation matrix describing the object

$$\vec{v}' = T_{1:3,4} \quad (2-1)$$

$$\vec{v} = \frac{\vec{v}'}{\|\vec{v}'\|}$$

We can use trigonometry to determine the angle (Eq. ??) and perpendicular axis (Eq. ??) between two vectors. For our example, we use the normalized viewing ray and the z-axis (symmetric axis) as the two vectors.

$$\cos(\theta) = \vec{v} \cdot \vec{z} \rightarrow \theta = \arccos(\vec{v} \cdot \vec{z}) \quad (2-2)$$

$$\vec{m} = \frac{\vec{v} \times \vec{z}}{\|\vec{v} \times \vec{z}\|} \quad (2-3)$$

Then, we can build a rotation matrix using an axis and an angle [?], (Eq. ??).

Then, we obtain the final transformation matrix describing the dual pose of the object by a post-multiplication of this rotation matrix.

$$T_{dual} = T_{orig} * \begin{pmatrix} R_{3 \times 3} & \vec{0}_{3 \times 1} \\ \vec{0}_{1 \times 3} & 1 \end{pmatrix} \quad (2-4)$$

In order to solve this problem, we must take advantage of a few key pieces of information:

- We know how humans think when they are performing model-image registration manually
- We know the mechanical properties of the soft tissue surrounding the knee, and can emulate these problems algebraically
- We can measure the relative performance between single-plane and bi-plane imaging, and create a “correction factor” to adjust single-plane measurements.

Each of these understanding motivates a different method by which we attempt to correct the symmetry trap. This aim serves to expand our understanding of how to overcome single-plane limitations in a rigorous manner.

2.2.2.1 Method 1: Artificial Ligaments Using Linear Springs

The first method involves modifying the cost function in order to incorporate linear springs that adjust for varus/valgus position, in much the same way that medial and lateral collateral ligaments do in human anatomy. First, we measure the relative kinematics between the femoral and tibial components in the current orientation, then use Euler angle decomposition to find the relative varus/valgus between them, and add a linear cost to that angle (??). This has the effect of choosing the choice of symmetric pose with a smaller varus/valgus angle, which is how human operators typically distinguish between the two options. The user must set the hyperparameter λ such that the function does not prioritize the angle over the contour matching. Once selected, this value does not change.

$$J = L_1 + \lambda |\theta_{VV}| \quad (2-5)$$

where $L_1 = ??$

2.2.2.2 Method 2: Binary Selection Between Symmetric Poses

This method relies on the fact that, given the pose of an object and the axis of symmetry, we can determine the symmetric pose of that object that will produce the same projective geometry.

The next method, rather than incorporating varus/valgus information in the objective function, simply calculates the angle for the current pose and it's symmetric pose, then picks the orientation that has the smaller absolute value. This can be more applicable than incorporating into the objective, because no hyperparameters need to be set, and the only information determining the correct orientation is the contour matching.

2.2.2.3 Method 3: Bi-plane Calibration

This method utilizes the validation data that we used from Aim 1 (??) in order to create post-hoc calibration for the position of the implant. First, we compare our single-plane data to gold standard bi-plane data, which represents the ground truth for the position of the implant. Then, using Bland-Altman plot, we can determine a calibration constant for our implant’s pose based on ground truth data. These calibration corrections will be applied after model-image registration has already been completed.

2.2.2.4 Method 4: Decision Tree for Binary Selection

Next, we turn to machine learning in order to determine the correct pose of the object. A decision tree is a type of data structure that traverses through a series of binary choices in order to arrive at some classification. One might imagine that Method 2 is a decision tree with a single query that is used to select the correct pose from the incorrect pose.

Given the 8000 frames of human-supervised measuring TKA kinematics, there are plenty of samples that can be used to train this decision tree. The driving force behind this method is capturing the latent human intuition used to make decisions for the correct pose when given two symmetric poses. Though we claim that varus/valgus is typically used, perhaps there are edge cases using different criteria that this formalism might elucidate.

2.2.2.5 Method 5: Neural Network for Binary Selection

This approach is very similar to the decision tree, but a fully connected network (??) is used to select between the two poses. Rather than a series of binary queries, a densely connected feature space will extract latent decision making for choosing the correct pose.

In order to perform this, first, the process of calculating the symmetric pose will be encoded into the neural network, and those layers will be frozen (i.e. the weights will not be updated during training). Then, each tibiofemoral pose will be randomized between “true” and “symmetric” for the training set, and represented as a quaternion. The pipeline will then fuse both the current (either true or symmetric) and the secondary pose at the

initial layer for a neural network, with the final output as a binary to “keep” or “switch” the pose. Different network weights will be tried, with a final emphasis on parsimony.

2.3 Aim 3: Pilot Trial of Application on Series of Patients

With the feasibility of an autonomous pipeline established, and single-plane limitations corrected, the final step is determining whether this pipeline is actually possible to use in a clinical workflow.

The third aim of this thesis hopes to test these capabilities on a pilot study measuring the kinematics of knee implants in a series of patients that have received total knee arthroplasty. Each patient will perform a series of tasks, and the time it takes to return a full kinematic evaluation will be reported, along with any areas that necessitated human supervision. The hypothesis of this aim is that the only areas that will require human supervision will be data-transfer to Joint Track Machine Learning, and post-processing in order to get useful results from the output kinematics. We will also use a human supervised method of measuring TKA kinematics to determine the difference between kinematics and the difference in time to examination report.

By establishing and benchmarking the difference in both time and levels of supervision required to generate a kinematic report, we hope to alleviate concerns over using this in the clinic.

2.4 Aim 4: Standardizing Kinematic Examination Protocols Using Novel Machine Learning Methods

After the pipeline has been utilized in a pilot study determining the efficacy of utilizing measurements in a clinical setting, we hope to refine the kinematic examination protocol in order to establish a standardized set of measurements that can give the most information to the surgeons. The overarching goal is to find the movements and stances that provide the most information about post-operative outcome, and to offer a robust series of movements that can be measured in any hospital with any imaging equipment. Historically, these movements has been decided by the researchers somewhat arbitrarily; we

hope to make an improvement to the standard of care by rigorously determining which movements provide the best “bang for your buck”.

2.4.1 The Movements

First, we are going to measure movements spanning the entire spectrum of imaging possibilities. Normal dynamic movements like walking, lunging, sit-to-stand, stair rise, stair descent, squatting, open-chain extension with the hip in flexion, neural, or extension positions, among others. Similarly, a wide array of static images will be taken, including those at maximum flexion and extension, weight and non-weight bearing, and during other static positions. We will cover the array of movements that patients cite as “unstable” or “uncomfortable”, and hopefully span every possible pathology with multiple static and dynamic images.

2.4.2 Statistical Analysis

One of the most recently developments in machine learning is Transformers. These new foundational networks utilize attention blocks in place of recurrent neural networks or convolutional neural networks in order to process time-series data [?]. One of the most interesting areas where these networks are highly performant is translation; both between languages, as well as image-to-text [?]. Key to the idea of translation is the notion of the “essence” of an “idea” (in the Platonic form) having multiple real-world instantiations. Specifically, there is an underlying meaning to a sentence or image that can be represented by words in another language. So, using this paradigm, is there an underlying “idea” present in a parametrized kinematic measurement that might be represented in either pathology, outcome or other measurement?

This aim seeks to answer the following questions: (1) can a transformer-based architecture be used to “translate” between parametrized representations of kinematic motion into another movement by the same patient, and can statistical salience between movements be inferred from an exploration into these networks? (2) Can these kinematic examinations be used to statistically infer post-operative outcome? (3) Can the total

number of measurements necessary to generate an accurate post-operative prediction be minimized by clustering different movements?

2.5 Aim 5: Joint Track Auto Toolkit: An Open Source Framework for Model-Image Registration

Parallel to clinical adoption of the Joint Track Machine Learning framework, we hope to expand the scope of use to research personnel studying joint kinematics around the world. Unfortunately, the current system is a system containing an exorbitant amount of legacy C++ code that requires extensive time and understanding to make changes. As such, the number of possibilities and ideas that can be tried and implemented remains in the hands of those students in the Gary J. Miller Ph.D. Orthopaedic Biomechanics Laboratory.

Thankfully, with the widespread adoption of Python as an easy-to-use language for researchers, and the ability to wrap back-end C++ subroutines with Python functions, it is possible to extend development of new and novel model-image registration pipelines to those without much software development expertise.

While working on specific research questions addressed above, current legacy code will slowly be replaced with dedicated and well documented libraries and subroutines, all of which will get Python wrappings. We hope that the introduction of this will allow other research groups to develop novel and unique model-image registration pipelines using the speed and accuracy of Joint Track Machine Learning as a foundation.

CHAPTER 3

JOINT TRACK MACHINE LEARNING: AN AUTONOMOUS METHOD OF MEASURING 6-DOF TKA KINEMATICS FROM SINGLE-PLANE FLUOROSCOPIC IMAGES

3.1 Introduction

Total Knee Arthroplasty (TKA) is a standard procedure for alleviating symptoms related to osteoarthritis in the knee. In 2018, orthopaedic surgeons performed more than 715,000 TKA operations in the United States [?]. This number is projected to increase to 3.48 million by 2030 [?] due to an aging population and increased obesity rates. While TKA largely relieves symptomatic osteoarthritis, roughly 20% of TKA patients express postoperative dissatisfaction, citing mechanical limitations, pain, and instability as the leading causes [? ? ?]. Standard methods of musculoskeletal diagnosis cannot quantify the dynamic state of the joint, either pre- or post-operatively; clinicians must rely on static imaging (radiography, MRI, CT) or qualitative mechanical tests to determine the condition of the affected joint, and these tests cannot easily be performed during weight-bearing or dynamic movement when most pain symptoms occur. Unfortunately, most of the tools used to quantify 3D dynamic motion are substantially affected by soft-tissue artifacts [? ? ?], are prohibitively time-consuming or expensive [?], or cannot be performed with equipment available at most hospitals.

Model-image registration is a process where a 3D model is aligned to match an object's projection in an image [?]. Researchers have performed model-image registration using single-plane fluoroscopic or flat-panel imaging since the 1990s. Early methods used pre-computed distance maps [? ?], or shape libraries [? ? ?] to match the projection of a 3D implant model to its projection in a radiographic image. With increasing computational capabilities, methods that iteratively compared implant projections to images were possible [? ? ?]. Most model-image registration methods provide sufficient accuracy for clinical joint assessment applications, including natural and replaced knees [? ? ? ?], natural and replaced shoulders [? ? ? ?], and extremities [? ? ?]. One of the main benefits of this

single-plane approach is that suitable images can be acquired with equipment found in most hospitals. The main impediment to implementing this approach into a standard clinical workflow is the time and expense of human operators to supervise the model-image registration process. These methods require either (1) an initial pose estimate [? ?], (2) a pre-segmented contour of the implant in the image [? ?], or (3) a human operator to assist the optimization routine out of local minima [?]. Each of these requirements makes model-image registration methods impractical for clinical use. Even state-of-the-art model-image registration techniques [?] require human initialization or segmentation to perform adequately.

Machine learning algorithms automate the process of analytical model building, utilizing specific algorithms to fit a series of inputs to their respective outputs. Neural networks are a subset of machine learning algorithms that utilize artificial neurons inspired by the human brain’s connections [?]. These networks have shown a great deal of success in many computer vision tasks, such as segmentation [? ? ?], pose estimation [? ?], and classification [? ? ?]. These capabilities might remove the need for human supervision from TKA model-image registration. Therefore, we propose a three-stage data analysis pipeline (Fig. ??) where a convolutional neural network (CNN) is used to segment, or identify, the pixels belonging to either a femoral or tibial component. Then, an initial pose estimate is generated comparing the segmented implant contour to a pre-computed shape library. Lastly, the initial pose estimate serves as the starting point for a Lipschitzian optimizer that aligns the contours of a 3D implant model to the contour of the CNN-segmented image.

This paper seeks to answer the following three questions: (1) How well does a convolutional neural network segment the femoral and tibial implants from fluoroscopic and flat-panel images? (2) Can a Fourier descriptor-based pose estimation method produce useful initial guesses of 3D implant pose from the CNN-segmented images? (3) Can the Lipschitzian optimizer, given reasonable initial guesses, replicate human-supervised TKA

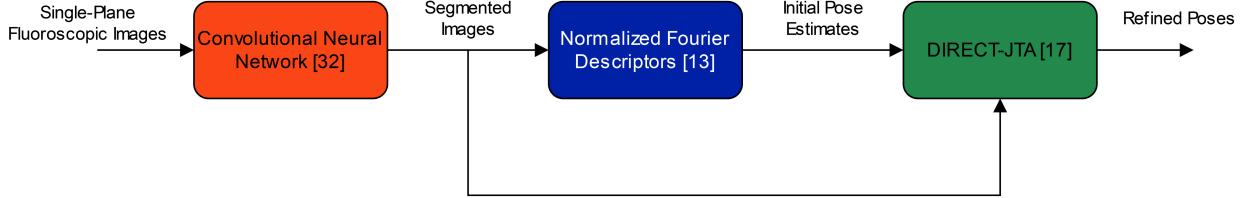


Figure 3-1. An overview of the pipeline for autonomous measurements of total knee arthroplasty kinematics. First, the data is processed through a convolutional neural network to locate the pixels belonging to the femoral and tibial implants [?], then, Normalized Fourier Descriptor shape libraries are used to determine and initial pose estimate [?], and lastly, DIRECT-JTA [?] is run on those segmented images using the NFD estimates as initializations for pose.

kinematic measurements?

3.2 Methods

Data from seven previously reported TKA kinematics studies were used for this study [? ? ? ? ? ? ?]. These studies utilized single-plane fluoroscopy or flat-panel imaging to measure tibiofemoral implant kinematics during lunge, squat, kneel, and stair climbing movements from 8248 images in 71 patients with implants from 7 manufacturers, including 36 distinct implants. From each of these studies, the following information was collected: (1) deidentified radiographic images, (2) x-ray calibration files, (3) manufacturer-supplied tibial and femoral implant surface geometry files (STL format), and (4) human supervised kinematics for the tibial and femoral components in each of the images. CNNs were trained with images from six of the studies using a transfer-learning paradigm with an open-source network [?]. CNN performance was tested using two image collections: a standard test set including images from the six studies used for training and a wholly naïve test set using images from the seventh study, where the imaging equipment and implants were different from anything used in training (Fig. ??). We used both test image sets to compare human-supervised kinematics with autonomously measured kinematics. Separately, two independent groups utilized our software to assess the accuracy of TKA kinematics measurements compared to their previously reported reference standard systems using RSA [?] or motion capture [?].

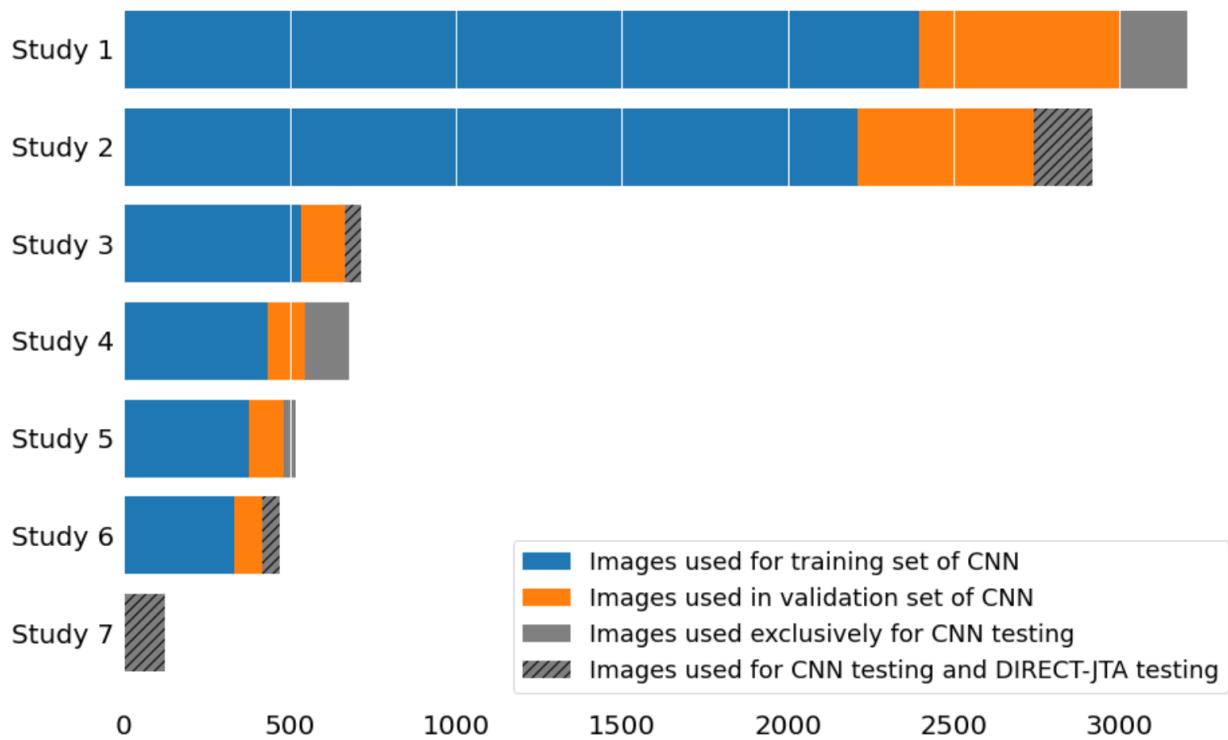


Figure 3-2. Data from seven studies were used to train and test the TKA kinematics measurement pipeline. Color coding in the figure identifies how many images were used for the training, validation, and testing functions. Images from the seventh study were used exclusively for testing the measurement pipeline that was trained using images from the other six studies.

3.2.1 Image Segmentation

Images were resized and padded to 1024x1024 pixels. Images containing bilateral implants had the contralateral knee cropped from the image. Segmentation labels were created by taking the human-supervised kinematics for each implant and generating a flat-shaded ground-truth projection image (Fig. ??). Two neural networks [?] were trained to segment the tibial and femoral implants, respectively, from the x-ray images. Each network was trained using a random 6284/1572 (80/20) training/validation split. Augmentations were introduced in the training pipeline to improve the network's generalization to new implants and implant types [?]. Each neural network was trained on an NVIDIA A100 GPU for 30 epochs. The performance of the segmentation networks was measured using the Jaccard Index [?]. This calculates the intersection between the

estimated and ground-truth pixels over the union of both sets of pixels. The ideal Jaccard index is 1.

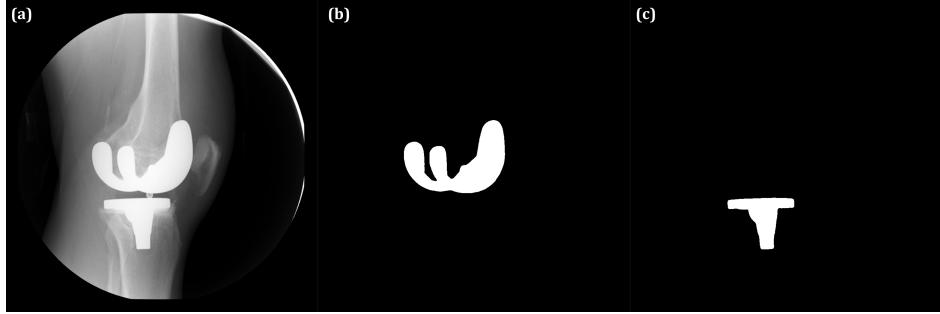


Figure 3-3. A representative fluoroscopic image is shown (a) with corresponding femoral (b) and tibial (c) ground-truth images created by flat-shaded projections of registered implant models.

3.2.2 Initial Pose Estimates

Initial pose estimates were generated from bounding contours of the CNN-segmented implant regions using Normalized Fourier Descriptor (NFD) shape libraries [? ? ?]. Shape libraries were created by projecting 3D implant models using the corresponding x-ray calibration parameters with $\pm 30^\circ$ ranges for the out-of-plane rotations at 3° increments (Fig. ??). Pose estimates were determined as previously described [?] NFD-derived femoral and tibial implant poses were transformed to anatomic joint angles and translations [?] and compared to the human-supervised kinematics for the same images using RMS differences for each joint pose parameter. The performance of this method was also assessed using flat-shaded projection images with perfect segmentation as a ground-truth reference standard.

3.2.3 Pose Refinement

A modified Dividing Rectangles (DIRECT) algorithm called DIRECT-JTA [?] generated the final pose estimates. This method of Lipschitzian optimization divides the search into three stages, the “trunk,” “branch,” and “leaf.” Each of the three stages was assigned distinct cost function parameters and search regions. The cost function used a computationally efficient L1-norm between the dilated contour from the segmentation label

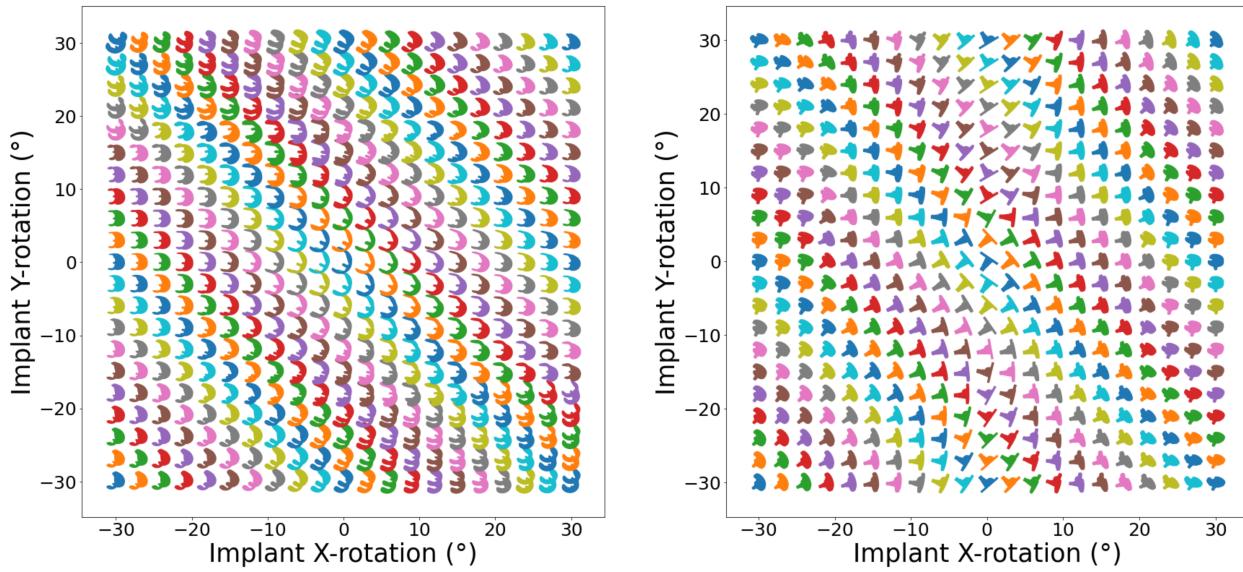


Figure 3-4. Femoral (left) and tibial (right) NFD shape libraries were generated to capture the variation in projection silhouette geometry with out-of-plane rotation [?]. Initial pose estimates were generated by comparing the NFD contour from the x-ray image to the shape library.

and the projected implant. Successively decreasing the dilation coefficient allowed the optimization routine to escape local minima, and the leaf branch served to find the optimal out-of-plane translation. Transversely symmetric tibial implants posed problems during registration because two distinct poses produced roughly identical projections [?]. Because of this pose ambiguity, the tibial implant was always optimized after the non-symmetric femoral implant. In addition to the dilation metric, the tibial mediolateral translation and varus/valgus rotations relative to the femur were penalized. Final implant poses were transformed into knee joint rotations and translations [?] and compared to the human-supervised kinematics for the same images using RMS differences for each joint pose parameter. Squared differences between data sets were compared using one-way MANOVA with post-hoc multiple pair-wise comparisons using the Games-Howell test (R v4.2.0 using R Studio, rstatix, and stats).

3.2.4 Pose Ambiguities and Registration Blunders

A blunder was defined as an image frame with the squared sum of rotation differences greater than 5° between autonomous and human-supervised measures. These blunder frames contain errors considerably larger than would be clinically acceptable and warrant further exploration. Blunders were analyzed with respect to the tibial implant's apparent varus/valgus rotation relative to the viewing ray (Fig. ??). A probability density function and cumulative density function were calculated for the blunder likelihood. Due to the high likelihood of blunders in this region, an ambiguous zone was defined for all apparent tibial varus/valgus-rotation less than 3.6° degrees, which is the mean + 1std of the blunder distribution (Fig. ??). Squared measurement differences between images inside and outside the ambiguous zone were also compared using one-way MANOVA with post-hoc multiple pair-wise comparisons using the Games-Howell test.

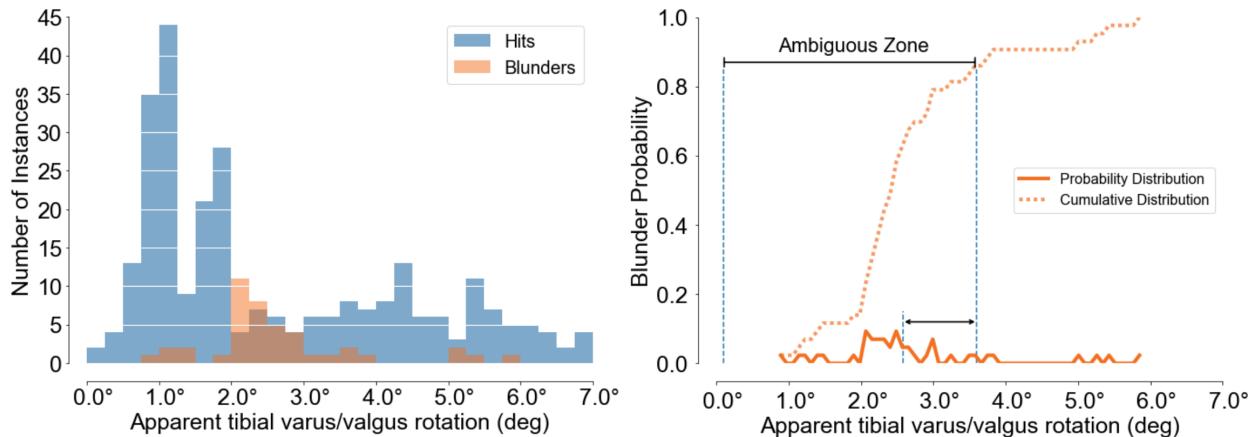


Figure 3-5. The histogram (left) shows the correctly registered frames (Hits, blue) and incorrectly registered frames (Blunders, orange) plotted as a function of the apparent tibial varus/valgus angle relative to the viewing raw. The probability plot (right) shows the distribution of blunders (solid orange) and the cumulative probability of blunders (dotted orange). The Ambiguous Zone is defined as apparent tibial varus/valgus rotations less than the mean + one standard deviation of the blunder probability distribution, capturing approximately 85 % of the blunders.

3.3 Results

CNN segmentation of standard test set images produced Jaccard indices of 0.936 for the femoral and 0.883 for the tibial components. CNN segmentation performance on the completely naïve test set was lower, 0.715 and 0.753, respectively.

The initial pose estimates were within the range of convergence for the DIRECT-JTA optimizer and offered a robust initialization for optimization (Table 1). The RMS differences for initial pose estimates on ground-truth images were smaller (better) than for CNN-segmented images, but the differences were mostly within a few millimeters or degrees. Due to poor sensitivity for measuring out-of-plane translation with monocular vision, the mediolateral translation had the largest RMS differences for both image types.

Table I
RMS Differences Between NFD Initial Estimates and Human-Supervised Kinematics

Implant	Images	Translation (mm)			Rotation (deg)		
		x	y	z	z	x	y
Femoral	CNN-Segmented Images	2.37	0.71	17.59	2.54	2.45	4.75
	Ground-Truth Projections	2.06	0.57	13.53	0.85	1.42	4.00
Tibial	CNN-Segmented Images	2.06	1.49	29.93	0.94	5.59	9.47
	Ground-Truth Projections	2.05	0.87	14.60	0.55	4.73	6.23

RMS differences between DIRECT-JTA optimized kinematics and human-supervised kinematics were sub-millimeters for all in-plane translations (Table II). Mediolateral translations and out-of-plane rotation differences were smaller when the pose of the tibia was outside the ambiguous zone. The RMS differences for the completely naïve test set were within 0.5 mm or 0.5 deg compared to the standard test set, indicating similar performance on the entirely novel dataset.

There was one femoral blunder and 43 tibial blunders out of 392 test images. Using the definition of the ambiguous zone as apparent tibial varus/valgus rotation less than 3.6 deg, 11% of images have a tibial blunder within this zone, compared to 3.2% outside. Sixty-six percent of tibial blunders were due to symmetry ambiguities (Fig ??).

One-hundred thirteen image pairs from an RSA study of TKA were used to

independently assess the accuracy of the autonomous kinematics measurement for single-plane lateral TKA images. RMS errors were 0.8mm for AP translation, 0.5mm for SI translation, 2.6mm for ML translation, 1.0° for flexion-extension, 1.2° for abduction-adduction, and 1.7° for internal-external rotation. At a different institution, 45 single-plane radiographic images were acquired with an instrumented sawbones phantom that was independently tracked using motion capture. Comparing the motion capture and autonomously measured radiographic kinematics, the RMS errors were 0.72mm for AP translation, 0.31mm for SI translation, 1.82mm for ML translation, 0.56° for flexion-extension, 0.63° for abduction-adduction, and 0.84° for internal-external rotation.

3.4 Discussion

Dynamic radiographic measurement of 3D TKA kinematics has provided important information for implant design and surgical technique for over 30 years. Many surgeons have expressed an interest in utilizing this type of measurement in their clinical practices; however, current methods are impractical. We developed a completely autonomous TKA kinematics measurement pipeline that can potentially provide a practical method for clinical implementation. This study sought to answer three questions, (1) How well does a neural network segment TKA implants from fluoroscopic and flat-panel images? (2) How well can an NFD shape library estimate the pose of a TKA implant given a CNN-segmented image? And (3) How well does a Lipschitzian optimization routine replicate human-supervised kinematics for TKA implants given an approximate initial

Table II
RMS Differences Between DIRECT-JTA Optimized and Human-Supervised Kinematics

Test Set	Image Group	Number of Images	A/P (mm)	S/I (mm)	M/L (mm)	Flx/Ext (°)	I/E (°)	V/V (°)
Standard	Inside AZ	187	0.694	0.523 ^b	1.752 ^a	0.730 ^a	3.380	1.938 ^a
	Outside AZ	83	0.685	0.466 ^c	0.917	1.029	1.811	0.605
Naïve	Inside AZ	47	0.802	0.739	1.715 ^d	1.388	4.044	2.480 ^d
	Outside AZ	75	0.692	0.644	0.691	1.031	1.154	0.846

AZ = Ambiguous Zone

Superscripts denote pairwise differences ($p < 0.05$) in squared errors for:

a. Standard Inside AZ vs Standard Outside AZ
b. Standard Inside AZ vs Naïve Inside AZ

c. Standard Outside AZ vs Naïve Outside AZ
d. Naïve Inside AZ vs Naïve Outside AZ

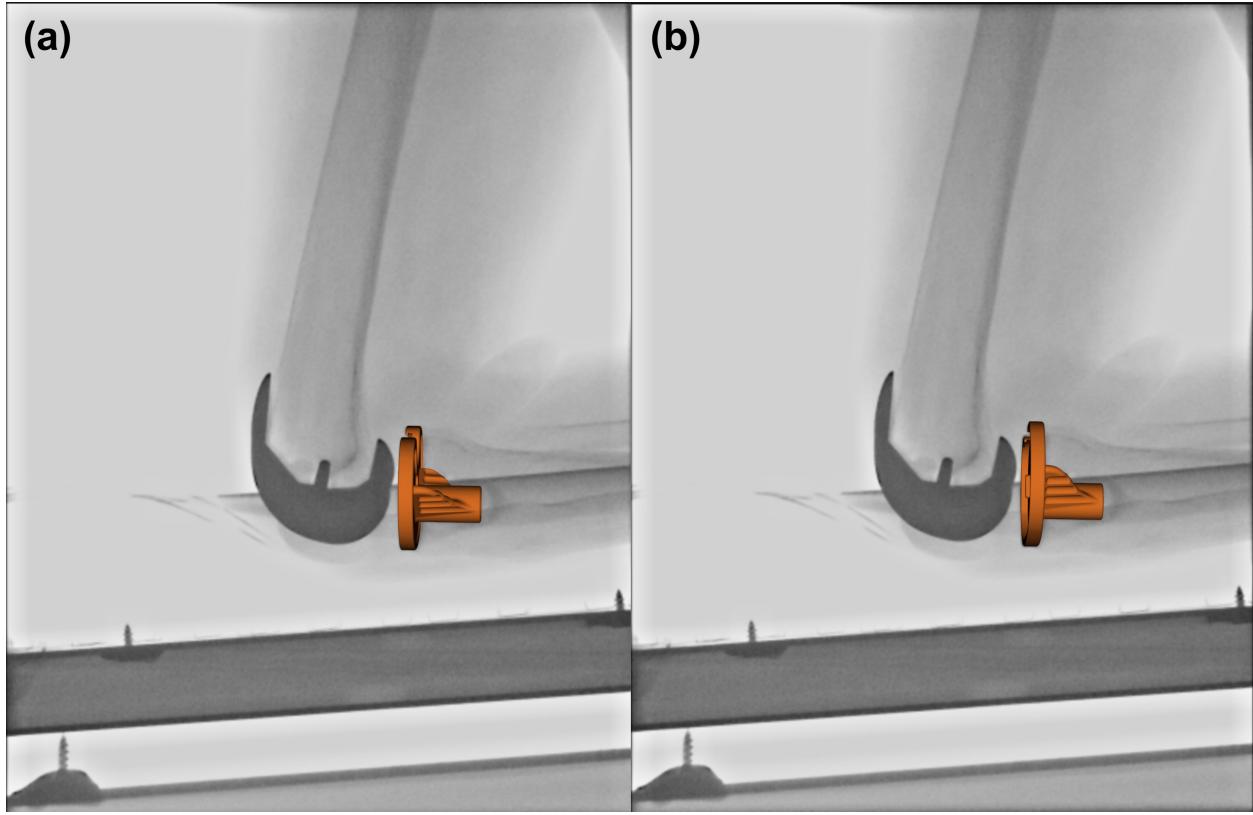


Figure 3-6. The figure shows the same radiographic image with two registered tibial implant poses: (a) shows a correctly registered tibial implant, while (b) shows an implant caught in a local cost function minimum corresponding to a nearly symmetric pose.

guess?

CNN image segmentation of TKA implants worked well, with Jaccard indices greater than 0.88 for the standard test set, and greater than 0.71 for the naïve test set.

Segmentation performance for the standard test set outperformed published examples by 0.05-0.1 Jaccard points [? ?], with the naïve test set on par with other segmentation examples. The most notable decrease in segmentation performance occurred along the perimeter of the segmented pixel region, especially in areas where implant projections occluded each other. These imperfectly segmented perimeter regions likely affect the initial pose estimate and the DIRECT-JTA optimization solution since both methods rely heavily on the segmented implant boundary. Further improvements can be made for the perimeter segmentation results by introducing intelligent augmentations during training using

generative models [?] and performing neural network bolstered contour improvement strategies [?].

Our initial pose estimates were satisfactory as an initialization for the DIRECT-JTA optimization, falling within the convergence region of $\pm 30^\circ$ [?]. However, the performance for the ground-truth projections was not as good as the cited method [?], which achieved errors of less than 1mm for in-plane translation and 2° for rotation. The cited method utilized an additional refinement step for the NFD estimation, interpolating the apparent out-of-plane angles between nearest shapes in the library. This extra step was not done because only approximate initial pose estimates were needed. In addition, the current study incorporated a vastly larger set of implant shapes (36 vs. 2) and image quality and calibration variations. Distinct implant shapes manifest unique normalization maps, where there can be discontinuities or jumps in normalization angles which affect the best-fitting library entry (Fig. ??) [? ?]. These details are easily upgraded with additional code using previously reported methods but were not pursued because the initial pose results were well within the DIRECT-JTA convergence region. The initial pose estimates for the CNN-segmented images were not as good as for the ground-truth projections. This follows directly from the fact that the perimeter of the segmented implants was not as accurately rendered, leading to poorer results with the edge-based NFD method. Finally, the out-of-plane translation estimates were relatively poor for both ground-truth projects and CNN-segmented images. This translation estimate is extremely sensitive to model projection and edge detection details and can be adjusted for better results if required.

RMS differences between human-supervised and DIRECT-JTA optimized kinematics demonstrate the two methods provide similar results. In-plane translation differences of less than 0.8mm and out-of-plane less than 1.8 mm, indicate good consistency in determining the relative locations of TKA implants. Rotation differences of 4° or less for frames within the ambiguous zone, and less than 1.7° for frames outside the ambiguous zone, indicate joint rotation measures with sufficient resolution to be clinically useful. We

observed two important characteristics in the measurement comparisons that will affect future implementations and use. First, we identified an ambiguous zone of apparent tibial rotations wherein there is a higher incidence of registration errors. These errors resulted in significant differences in measurement performance for the out-of-plane translations and rotations. This phenomenon, resulting from the nearly symmetric nature of most tibial implants [? ? ? ? ?] prompts either practical modification to imaging protocols to bias the tibial view outside the ambiguous zone or modifications of the model-image registration code to enforce smooth kinematic continuity across image frames and/or to impose joint penetration/separation penalties [?]. Second, we observed similar measurement performance for the standard and naïve test sets, which differed only in the superior/inferior joint translation. This suggests that the autonomous kinematic processing pipeline can provide reliable measures for implants and imaging systems that were not part of the training set, which will be important for application in novel clinical environments.

Two independent research teams utilized our software to evaluate the accuracy of our autonomous measurement pipeline compared to their reference standard methods using implants and image detectors that were not part of our training sets. In both cases, the accuracy results were comparable to results reported for contemporary human-supervised single-plane model-image registration methods for TKA kinematics [? ? ? ? ?]. Interestingly, the independent accuracy results appeared superior to our assessment of differences between autonomous and human-supervised measures of TKA kinematics. In both cases, the independent centers used high-resolution flat-panel detectors that provided better spatial resolution and grayscale contrast than most of the imaging systems included in our datasets. With images of similar quality, it is reasonable to expect similar measurement accuracy.

This work has several limitations. First, the image data sets resulted from previous studies in our labs, so there was no prospective design of which implant systems and image detectors should be included for a pipeline that generalizes well to other implants and

detectors. Nevertheless, the naïve data set and the independent assessments, all involving implants and detectors not used for training, performed well and suggest that the method can usefully generalize to measurements of traditionally configured TKA implants. Future work is required to evaluate measurement performance with partial knee arthroplasty or revision implants. Second, many methodologic and configuration options and alternatives remain to be explored, and the current pipeline implementation should not be considered optimal. How best to disambiguate tibial poses and determine the most effective and robust optimization cost functions are areas of current effort.

We present an autonomous pipeline for measuring 3D TKA kinematics from single-plane radiographic images. Measurement reproducibility and accuracy are comparable to contemporary human-supervised methods. We believe capabilities like this will soon make it practical to perform dynamic TKA kinematic analysis in a clinical workflow, where these measures can help surgeons objectively determine the best course of treatment for their patients.

CHAPTER 4

AIM 2: CORRECTING SYMMETRIC IMPLANT AMBIGUITY IN MEASURING
TOTAL KNEE ARTHROPLASTY KINEMATICS FROM SINGLE-PLANE
FLUOROSCOPY

4.1 Abstract

Recent advancements in computer vision and machine learning enable autonomous measurement of total knee arthroplasty kinematics through single-plane fluoroscopy. However, symmetric implants present challenges in optimization routines, causing 'symmetry traps' and ambiguous poses. Achieving clinically robust kinematics measurement requires addressing this issue. We devised an algorithm that converts a 'true' pose to its corresponding 'symmetry trap' orientation. From a dataset of nearly 13,000 human supervised kinematics, this algorithm constructs an augmented dataset of 'true' and 'symmetry trap' kinematics, used to train eight classification machine learning algorithms. The outputs from the highest-performing algorithm classify kinematics sequences as 'obviously true' or 'potentially ambiguous.' We construct a spline through 'obviously true' poses, and 'ambiguous' poses are compared to the spline to determine correct orientation. The machine learning algorithms achieved 88-94% accuracy on our internal test set and 91-93% on our external test set. Applying our spline algorithm to kinematics sequences yielded 91.1% accuracy, 94% specificity, but 67% sensitivity. The accuracy of standard ML algorithms for implants within 5 degrees of a pure-lateral view was 71%, rising to 88% beyond 5 degrees. This pioneering study systematizes addressing model-image registration issues with symmetric tibial implants. High accuracy suggests potential use of ML algorithms to mitigate shape-ambiguity errors in pose measurements from single-plane fluoroscopy. Our results also suggest an imaging protocol for measuring kinematics that favors more oblique viewing angles, which could further disambiguate 'true' and 'symmetry trap' poses.

4.2 Introduction

Measuring total knee arthroplasty (TKA) kinematics from fluoroscopic images has been an important contributor to knee implant design, post-operative assessment, and

predictive modeling for wear and failure patterns for nearly 30 years [? ? ?]. However, the application of this technology has been limited to research use by the challenges of performing these measurements quickly and reliably, as they often require expensive equipment and time-consuming processes [? ? ? ?]. Recent advancements in computer vision and machine learning have opened up the possibility of using a single fluoro-camera for fully autonomous TKA kinematic measurement[? ?], making it more accessible and cost-effective for hospitals and clinics. Nonetheless, this approach faces inherent limitations in accurately resolving orientation and location information using only one fluoro-camera view [? ? ? ? ?].

One of the key issues researchers encounter when dealing with imaging from a single camera arises from symmetric implant geometries. Under a weak perspective/nearly orthographic projection paradigm [?], symmetric implants have distinct 3D orientations that produce nearly indistinguishable 2-dimensional projective geometry (??). Due to the model-image registration process relying on the information present in the 2D fluoroscopic image, this causes multiple local minima for many optimization algorithms.

Human-supervised kinematics measurements for TKA with symmetric tibial implants frequently rely on the relative location of the fibula in the fluoroscopic image to disambiguate the pose. Unfortunately, fully autonomous solutions can't utilize this reference point, causing difficulty in algorithmic implementation.

In this paper, we utilize 12,592 images from seven studies that utilized human-supervised TKA kinematics [? ? ? ? ? ? ?] to explore the potential of a novel method to classify and correct 'symmetry trap' poses in single-plane TKA kinematic measurements. The goal is to improve the accuracy and resolution of single-plane TKA measurements and enable more accessible and reliable kinematic data for orthopaedic applications.

This research paper will answer two key questions: (1) How effectively can machine learning methods distinguish between 'true' and 'symmetry trap' poses in cases involving

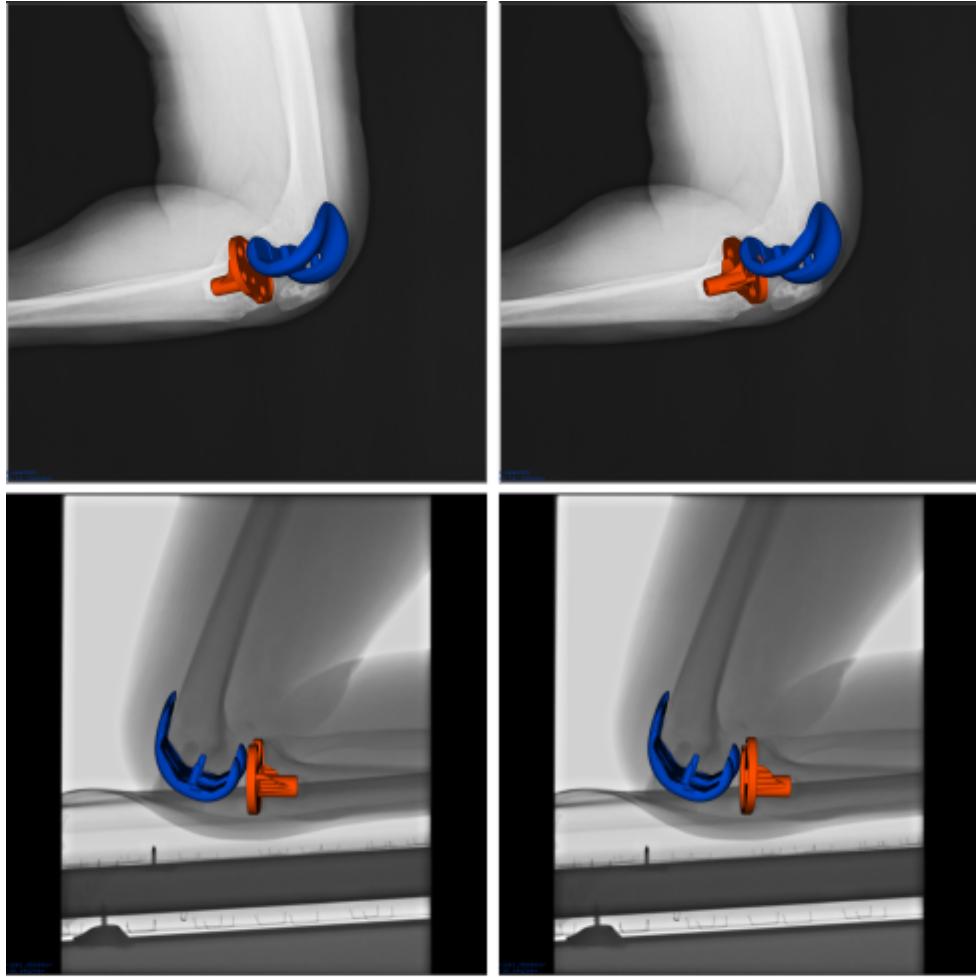


Figure 1. Two examples of symmetric pose dyads with correct (left) and symmetry trap (right) kinematics. For some orientations, the symmetry trap is obviously incorrect (top), and for others, it is more ambiguous (bottom).

Figure 4-1. Two examples of symmetry pose dyads with the correct (left) and symmetry trap (right) kinematics. For some orientations, the symmetry trap is obviously incorrect (top), and for others, it is more ambiguous (bottom).

symmetric implants? (2) How well can we systematically correct 'symmetry traps' arising in kinematics sequences obtained from fluoroscopic images?

4.3 Methods

4.3.1 Study Design and Sample

In this paper, we are addressing the issue of 'symmetry traps' in TKA kinematics measurements from single-plane fluoroscopy. In order to do so, we take a data-driven approach. We start with a collection of human-supervised 'true' kinematic measurements,

then determine the respective ‘symmetry trap’ orientation using a novel algorithm. Using this collection of both ‘true’ and ‘symmetric’ poses, we train eight different machine learning classifiers and evaluate their performance on both an internal and external test set. We then take the highest performing classifier and use it to establish an algorithm that takes as input a contiguous sequence of kinematics measurements, and imposes continuity constraints to fix ‘symmetry traps’ in real-world data. We report the performance of our cubic spline correction algorithm on a dataset of autonomously measured kinematics [?]. All programs were written in Python 3.10.

The dataset used for development of the methodology comprises seven total studies [? ? ? ? ? ? ?] comprising 12,592 frames of radiographic data. We completely withheld data from one study [?] during training to use later as an external test set.

4.3.2 Determining Symmetric Orientations

First, ‘symmetry trap’ poses were identified from ‘true’ poses through a novel algorithm devised to “flip” any given pose into its symmetric counterpart (??). The algorithm proceeds as follows:

1. Determine the viewing ray from the object origin to the camera origin. Denote this as \vec{v} , where $||\vec{v}|| = 1$.
2. Determine the symmetric axis of the object, \vec{s} , where $||\vec{s}|| = 1$. This symmetric axis is equivalent to the normal vector of the “mirror plane” of symmetry for the object.
3. Determine the axis and angle of rotation between \vec{s} and \vec{v} and construct the equivalent rotation matrix.
 - (a) The axis is the cross product, $\vec{m} = \vec{v} \times \vec{s}$.
 - (b) The angle is the normalized dot product, $\psi = 2 \times \cos^{-1} \frac{\vec{v} \cdot \vec{s}}{||\vec{v}|| ||\vec{s}||}$
4. Apply body-centered rotation to the symmetric obhject about \vec{m} by ψ .

For each projective geometry, there exist exactly two poses (i.e., a dyad) that produce it, thus applying this algorithm twice will return the original pose.

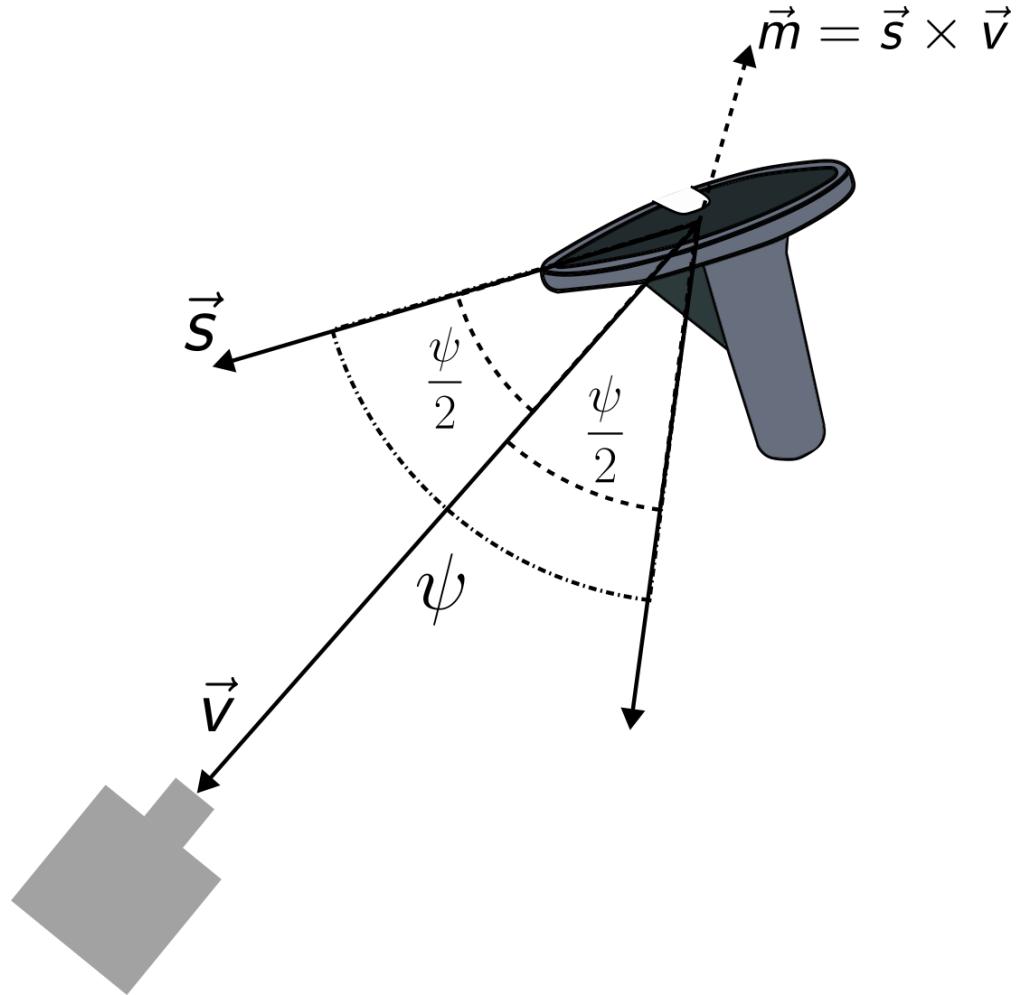


Figure 4-2. A visualization of the “symmetry flipper” algorithm that converts 3D orientations into their symmetric counterpart.

4.3.3 Constructing Dataset

We collected data from seven studies utilizing human supervised total knee arthroplasty kinematics, totaling approximately 12,592 individual frames. Anatomic tibiofemoral kinematics were stored as ‘true’ poses. The proposed algorithm was then applied to each frame’s tibial implant, generating corresponding “symmetric” poses, for a total of 25,184 samples (??). Solid angle distances (ψ) obtained from the symmetry flipper algorithm were also stored, and the dataset was stratified based on the solid-angle value

and split into 67% training and 33% testing sets for unbiased evaluation. Data from one study [?] was completely withheld from training to be used as an external test set.

Thus, for each sample, the input was $\{\theta_{Int/Ext}, \theta_{Flx/Ext}, \theta_{Var/Val}, \psi\}$ with $\theta_{Int/Ext}, \theta_{Flx/Ext}, \theta_{Var/Val}$ representing anatomic internal/external rotation, flexion/extension, and varus/valgus respectively. The target was either ‘true’ or ‘symmetry trap’.

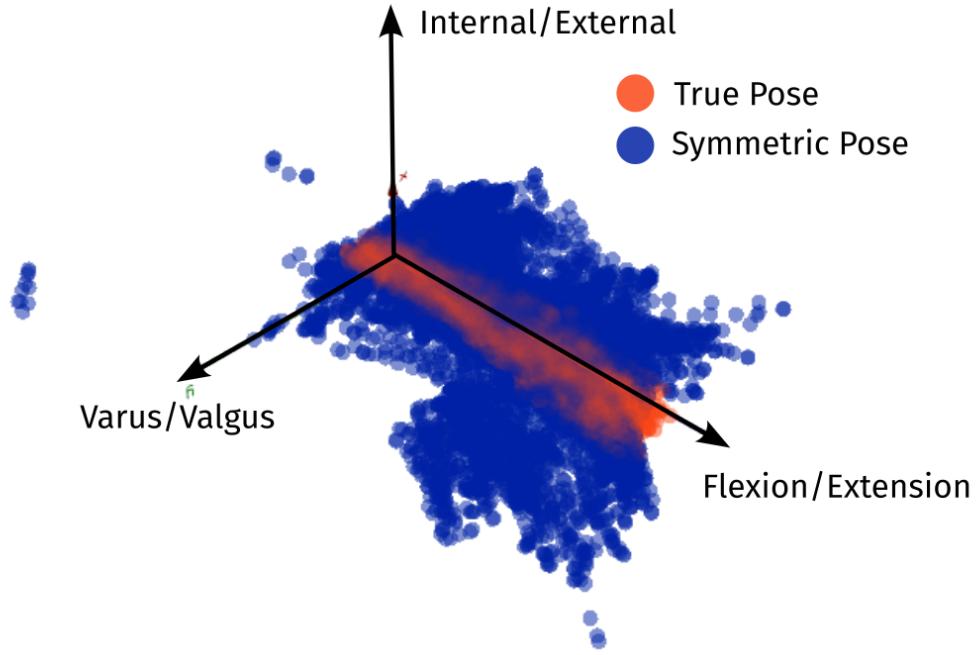


Figure 4-3. A 3D scatter plot of our training data with ‘true’ poses in orange and ‘symmetry trap’ poses in blue. Of note, there are distinct regions of exclusively ‘symmetry trap’ poses, while the region of predominately ‘true’ poses (orange cylinder along Flexion/Extension axis) also has many ‘symmetry trap’ poses.

4.3.4 Classification Algorithms

In this study, we implemented a selection of classification algorithms (Table 1) from Scikit-learn [?] to distinguish between ‘true’ and ‘symmetry trap’ poses. The chosen algorithms were K-Nearest-Neighbors (KNN) [?], Support Vector Machines (SVM) [?], AdaBoost [?], Histogram Gradient Boosting [?], Bagging Meta-Estimator [?], Stacked Generalization [? ?], and (majority) Voting Classifier. The hyperparameters for each method were tuned using stratified 5-fold cross validation using grid-search to test all

possible combinations. Stacked generalization and the voting classifier used the tuned hyperparameters from each of the other estimators in their predictions. We recorded sensitivity [?], specificity [?], accuracy [?], and the F1 score [?] for each algorithm on our test set. We also report these metrics of the highest performing individual classifier on stratified partitions of the test set to determine relative performance for different ψ values.

4.3.5 Fixing Incorrect Symmetry Trap Poses

The primary objective of this study is to accurately measure the kinematics of all frames in a TKA kinematics sequence and to robustly fix incorrect symmetry trap poses. However, due to the inability to know a-priori which poses are ‘true’ and in a ‘symmetry trap’, we must employ a separate procedure for systematically correcting individual frames. Thus, for every kinematics sequence we employ the following procedure:

1. Determining Ambiguous and Non-ambiguous Poses for each frame
 - (a) For each frame, we generate the ‘symmetric’ pose corresponding to the input pose.
 - (b) We run both the input and its symmetric pose through the highest performing classification algorithm.
 - (c) If the input pose and its symmetric pose are differently labeled by the classifier (i.e., one is labeled ‘true’, and the other ‘symmetry trap’), take the pose labeled ‘true’ as the actual pose.
 - (d) If both poses return the same label (i.e., both ‘true’ or both ‘symmetry trap’), label the pose as “ambiguous” and move to step 2.
 - (e) If no poses were labeled “ambiguous”, the procedure is finished.
2. Construction of 3-Dimensional cubic b-spline of the movement
 - (a) A spline is individually fit through the flexion/extension, internal/external rotation, and adduction/abduction angles for all frames that were **not** labeled

ambiguous.

3. Correcting Ambiguous Poses

- (a) We sample the cubic spline at frame locations of ambiguous poses.
- (b) We compare the solid-angle difference between the input and symmetric pose to the pose at the sampled spline location.
- (c) Whichever pose is closer to the spline, we take this as the ‘true’ pose.

To evaluate this procedure, we used a dataset of kinematics generated from a fully autonomous method [?] to emulate the real-world use case of this algorithm as a post-processing operation. This test set has robustly quantified the presence of ‘symmetry traps’, and so allows us to report sensitivity, specificity, accuracy, and the F1 measure of its ability to correct symmetry traps in a clinical setting.

4.4 Results

4.4.1 Machine Learning Classification Performance

Classification accuracy ranged from 87.8% to 94.0%, sensitivity ranged from 92.9% to 94.7%, and specificity ranged from 83.5% to 93.6% for all eight methods evaluated (Table 1). Stacked generalization and the voting classifier tended to outperform the other classifiers, as they both incorporated a combination of the other models in their decision-making. For the external test set (Table 1), the Support Vector Machine with Radial Basis Function kernel had the highest accuracy and F1 score, at 94% and 0.94 respectively. Stacked generalization performed slightly worse in the external test set. Overall, the performance in the external test set was comparable to the performance on the internal test set, with some classifiers seeing decreased performance (K-Nearest-Neighbors, Stacked Generalization), and others seeing slightly increased performance (Support Vector Machines, AdaBoost, Voting Classifier).

Table 1: Machine Learning Classifier Performance						
Classifier	Tuned Hyperparameters	Test Set	Accuracy (%)	Sensitivity (%)	Specificity (%)	F1-Score
Support Vector Machines (Radial Basis Function)	C = 1000	Internal	92.08	94.75	89.73	0.92
		External	94.21	97.12	91.64	0.94
Support Vector Machine (Polynomial Kernel)	C = 1000 Polynomial Degree = 2	Internal	87.65	92.54	83.8	0.87
		External	92.08	96.58	88.38	0.92
K-Nearest-Neighbors	Neighbors = 4 Distance Metric = Minkowski	Internal	93.12	93.96	92.32	0.93
		External	90.93	93.62	88.55	0.91
AdaBoost	Num. Estimators = 200	Internal	88.78	91.08	86.74	0.88
		External	92.86	97.23	89.22	0.93
Histogram Gradient Boosting	Learning Rate = 0.1	Internal	93.11	94.98	91.4	0.93
		External	93.24	96.67	90.29	0.93
Bagging Estimators	Num. Estimators = 500	Internal	93.32	94.28	92.41	0.93
		External	93.82	95.95	91.88	0.94
Stacked Generalization	N/A	Internal	94.3	94.78	93.84	0.94
		External	92.86	94.94	90.96	0.93
Majority Voting Classifier	N/A	Internal	92.63	95.85	89.86	0.92
		External	93.34	96.87	90.31	0.93

Table 2: Stratified ψ Test Set Stacked Generalization Classification Performance					
Psi Range	Sample Size	Accuracy (%)	Sensitivity (%)	Specificity (%)	F1-Score
0 – 5°	488	71.04	71.43	70.7	0.69
5 – 10°	1132	88.21	90.53	85.98	0.88
10 – 15°	1224	93.0	92.78	93.2	0.93
15 – 20°	1107	96.13	96.97	95.31	0.96
> 20°	3568	98.28	98.32	98.24	0.98

4.4.2 Stratified ψ Test Set Performance

As the value of ψ increases, the performance of stacked generalization classification increases (Table 2). At values closer to a pure-lateral view (0), accuracy is roughly 71%, and once the implants exceed 5 off-lateral, accuracy increases to above 88%. See Appendix A for visualization of different test sets.

4.4.3 Symmetry Trap Correction Performance

When applying the procedure to correct symmetry traps, we achieve an accuracy of 91.1%, a sensitivity of 67.4% and a specificity of 94% on an external test set. The average “distance to symmetric pose” (ψ) was for the frames correctly classified, and for the frames incorrectly classified. For our initial internal test set, the distance to the symmetric pose was for frames correctly classified, and for incorrectly classified frames. From this, we can see that the incorrect frames are much closer to a true lateral single-plane image, meaning the difference between the symmetric dyad of poses is smaller.

4.5 Discussion

Orthopaedic surgeons would value a practical method to quantify TKA kinematics clinically for post-operative assessment or exploration of unsatisfactory outcomes [?]. To date, research methods for quantifying TKA kinematics from single-plane radiographic images have been too time-consuming for clinical use, but the application of machine learning methods has been rapidly improving their convenience [?]. A critical step in making this technology ready for clinical use is to address the pose ambiguity issue that arises from single-camera views of symmetric tibial implant components in a time-efficient manner. We found that several machine learning algorithms were able to robustly identify and correct symmetric erroneous poses with greater than 90% accuracy, making these measurement methods sufficiently autonomous and robust to be considered for clinical use.

We demonstrated that classical machine learning techniques are extremely adept at handling the classification of ‘true’ and ‘symmetry trap’ TKA kinematics measurements from single-plane fluoroscopy. As a first study in this domain, benchmark comparisons with

other metrics, datasets, or algorithms are not available. However, an accuracy level of at least 87% establishes a useful performance benchmark in terms of clinical applications of these methods and for future research in this field. As expected, stacked generalization [? ?], where “the worst possible performance is the highest of any individual estimator” held true for our internal test set, and outperformed all other algorithms. Interestingly, the same trend did not hold for our external test set, which contained kinematics from patients, imaging machines, and surgeons that were completely withheld from training. In our external test set, SVMs with radial basis function kernel seems to outperform the other classifiers.

Our symmetry trap correction procedure has strong performance in the accuracy ($> 91\%$) and specificity ($> 94\%$). We note decreased sensitivity related to decreasing values of our term, consistent with intuition: when two possible poses are closer together, they are harder to disambiguate. When performance is measured based on ψ values greater than 5 or 10 degrees, we see that our symmetry trap correction algorithm boasts performance similar to that of the stacked generalization estimator (Table 2). Previous work has also shown an increase in the prevalence of symmetry traps in low- poses [?]. We also note that in real-world data, we have imbalanced class representation because most symmetry traps are caught in the initial kinematics measurement, leaving only the most difficult cases for this algorithm to correct. These reasons both culminate with an expected decreased sensitivity when trying to correct symmetry trap poses. Future work might be able to address this problem using time-series machine learning algorithms like recurrent neural networks [?] and transformers [?] to isolate specific frames that have fallen into symmetry traps based on local and global relationships with other kinematics in that particular sequence. Additionally, 3D geometries could be used to impose penetration and separation penalties that are further used to refine ambiguous poses.

The performance of autonomous measurements of single-plane fluoroscopy kinematics measurements already achieve clinically acceptable levels of error. Despite this,

measurement errors still occur and roughly 67% of frames that are incorrectly optimized fall into symmetry traps [? ?]. With the techniques proposed in this work, we can drastically improve the overall performance of autonomous single-plane TKA kinematics measurements and remove the largest source of inaccuracy.

Stratifying measurement performance by (Table 3), it is clear that the joint viewing perspective is a dominant contributor to pose ambiguities. This suggests a modification to imaging protocols for measuring TKA kinematics to minimize the number of frames that fall into the “ambiguous zone”. Instead of “pure lateral” imaging that attempts to perfectly center the knee in the image frame, we suggest positioning the center of the knee 1-2 inches away from the center of projection, and to rotate the plane of the observed leg by 5 degrees from parallel with the image detector. These small adjustments to the imager/patient alignment would cause all kinematics to have $\angle 5$, which improves classification accuracy by at least 20%. Similar imaging protocol improvements have been proposed for bi-plane RSA [?]. This study has some limitations. In this study, we use human-supervised kinematics from single-plane images, where there is always the possibility that the human operator fell into error when distinguishing between ‘true’ and ‘symmetry trap’ poses in a sequence of images. This can occur due to a myriad of factors, including low image quality and near-lateral imaging angles. To mitigate this, validated ground truth kinematics are necessary to create a robust algorithm without any errors in the training data.

The dataset constructed in this study assumes that the femoral implant has been properly registered with correct kinematics. Often, this is not an issue because most femoral implants are asymmetric, and have an extremely low likelihood of falling into their own ‘symmetry traps’ [?]. However, improper registration of both a symmetric tibial implant and symmetric femoral implant would decrease the feasibility of the proposed procedure.

This study presents a novel method of overcoming one of the most pernicious issues facing researchers measuring TKA kinematics from single-plane fluoroscopy. Our algorithm

combines 3 elements (1) calculating 'symmetry trap' poses from 'true' poses, (2) classical machine learning algorithms to classify a single pose as 'true' or 'symmetry trap' and (3) a procedure that applies those machine learning outputs to a contiguous sequence of kinematics to construct a spline and correct symmetry traps in real-world data. Further, our results suggest slight alterations to imaging angles could reduce errors due to ambiguous poses resulting from symmetrical tibial implants. We believe that this method, along with recent advancements in autonomous kinematics measurements, will soon make it practical to perform dynamic TKA kinematics measurements in a clinical setting without the inherent limitations of single-plane radiographic imaging.

4.6 Conclusion

Our three-stage approach to correcting 'symmetry traps' that arise from symmetric tibial implants when performing kinematics measurements on single-plane imaging provides a groundbreaking solution to a problem facing researchers in the field for nearly three decades. Further, the "symmetric flipper" algorithm provides a strong foundational approach for any symmetric object in a weak-perspective paradigm and can drastically reduce the complexity of many optimization algorithms by reducing the search space. Because this algorithm is model-agnostic, it could be used in many applications where orientation ambiguities arise. This approach further refines the ability for a fully autonomous kinematics measurement platform, bringing it one step closer to being a clinically viable technology, where these measures can help surgeons objectively determine the best course of treatment for their patients.

CHAPTER 5

MUSINGS ON LATENT KINEMATICS SPACE AND SYNTHETIC BIOMECHANICS DATA

5.1 Abstract

5.2 Introduction

In order to make autonomously measuring joint kinematics a clinically viable tool, the movements and motions being measured must be standardized. Similar efforts have been put into standardized post-operative outcome metrics like the Knee Society Score (KSS), Knee Injury and Osteoarthritis Outcome Score (KOOS), the Forgotten Joint Score, the Western Ontario and McMaster Universities Osteoarthritis Index (WOMAC), among others. Efforts put toward standardizing these metrics allow different groups, different implants, and different surgical techniques to be reliably compared to one another as objectively as possible. This same need for standardized metrics exists when quantifying joint kinematics for research and clinical purposes. Unfortunately, no such standardization exists for quantifying joint kinematics. Most groups will perform roughly the same activities of daily living (Some form of walking, stair rise, chair rise, kneel, lunge, and squat), but the differences in how each group performs these activities prevent highly accurate direct comparison.

So, this chapter describes a data-driven approach toward being able to select the highest yield activities to be included in a clinical kinematics examination. In order to be a clinically practical method, the time to collect the data must not be excessive, and so a brief and succinct set of activities is preferred. From these constraints, the first approach stems from the thought, “All kinematics data is coming from the same patient. Are there underlying characteristics from one activity that translate to all activities?” An example might arise: “Can we predict the patient’s stair rise kinematics from their walking kinematics?” The motivating intuition behind these ideas is that they all come from the same patient, and are being driven by the same underlying musculature. The following chapter will detail the beginning of an experimental procedure, and why the data was

insufficient for answering the above questions. However, another question arises: “Given the statistical parameters surrounding a “good outcome” or “bad outcome”, might we be able to generate synthetic patient data to overcome the lack of real-world kinematics data?” This is an extremely pertinent question in the world of “big data”. Many modern algorithms leverage an enormous amount of data, often more than traditional methods of quantified biomechanics are able to generate. Additionally, HIPAA requirements make it incredibly difficult to share data, slowing down the rapid progress in curated algorithms for biomechanics data that we see with open datasets like ImageNet [?].

And so, the following chapter will contain some background information toward solving these problems.

5.3 Latent Kinematics Space and Reducing the Number of Movements Required for a Full Kinematics Evaluation

The concept of latent space in kinematics is analogous to the latent spaces found in other domains, such as natural language processing (NLP). In NLP, latent space refers to the underlying, hidden structure in language data that algorithms attempt to uncover [?]. This is where the power of transformers [?] in NLP becomes relevant to our study of kinematics.

In kinematics, each movement pattern can be thought of as a unique “expression” of the underlying musculoskeletal system, akin to how sentences are expressions of underlying linguistic rules and meanings in language. By exploring the latent kinematics space, we aim to reduce the number of movements required for a comprehensive kinematics evaluation without losing the predictive power or statistical significance of the evaluation. This approach is based on the hypothesis that certain movement patterns contain redundant information, which could be inferred from other, more distinct movements. For instance, we might predict the kinematics of a “stair rise” from the “walking” pattern, assuming that the essential biomechanical information is embedded within the walking movement.

Test

5.4 Toward Generative Synthetic Biomechanics Data

BIOGRAPHICAL SKETCH

Andrew Jensen is a Florida native from Sarasota, Florida. He attended the University of Florida for his undergraduate degree in Mechanical Engineering, for which he received high honors. He took a brief hiatus from school to work at an orthopaedic solutions company, Exactech. The COVID-19 pandemic cut his time at Exactech short, so he joined the Gary J Miller Orthopaedic Biomechanics Laboratory as a part-time researcher during the summer leading up to his first official semester of graduate school.

Andrew enjoys reading (mostly philosophy and science fiction), working out, and hanging out with his dog.