



A Real-time Pose Estimation Application for Tinikling

A Capstone Project on Operational Technologies
Presented to the Faculty of the
Department of Electronics and Computer Engineering
Gokongwei College of Engineering
De La Salle University

In Partial Fulfillment of the
Operational Technologies

by
CALAGUIAN Nathan Raekel L.
ELLAR Gerald Antonio P.
MAHAIT Hans

October, 2025



20

ABSTRACT

21

Index Terms—Dance, Pose Estimation, Real-time, OpenPose .



22

TABLE OF CONTENTS

23

Abstract	ii
-----------------	-----------

24

Table of Contents	iii
--------------------------	------------

25

List of Figures	v
------------------------	----------

26

List of Tables	vi
-----------------------	-----------

27

Abbreviations and Acronyms	vii
-----------------------------------	------------

28

Notations	viii
------------------	-------------

29

Glossary	ix
-----------------	-----------

30

Listings	x
-----------------	----------

31

Chapter 1 INTRODUCTION	1
-------------------------------	----------

32

1.1 Background of the Study	2
---------------------------------------	---

33

1.2 Prior Studies	3
-----------------------------	---

34

1.3 Problem Statement	4
---------------------------------	---

35

1.4 Objectives and Deliverables	6
---	---

36

1.4.1 General Objective (GO)	6
--	---

37

1.4.2 Specific Objectives (SOs)	6
---	---

38

1.4.3 Expected Deliverables	7
---------------------------------------	---

39

1.5 Significance of the Study	7
---	---

40

1.5.1 Technical Benefit	8
-----------------------------------	---

41

1.5.2 Social Impact	9
-------------------------------	---

42

1.5.3 Environmental Welfare	9
---------------------------------------	---

43

1.6 Assumptions, Scope, and Delimitations	9
---	---

44

1.6.1 Assumptions	9
-----------------------------	---

45

1.6.2 Scope	10
-----------------------	----

46

1.6.3 Delimitations	10
-------------------------------	----

47

1.7 Description and Methodology of the Capstone Project on Operational	
--	--

48

Technologies	11
------------------------	----

49

1.8 Estimated Work Schedule and Budget	12
--	----

50

1.8.1 Milestones and Gantt Chart	12
--	----



51	1.8.2 Budget	12
52	1.9 Overview of the Capstone Project on Operational Technologies	13
53	Chapter 2 LITERATURE REVIEW	14
54	2.1 Existing Work	15
55	2.2 Lacking in the Approaches	22
56	2.3 Summary	23
57	Chapter 3 THEORETICAL CONSIDERATIONS	24
58	Chapter 4 DESIGN CONSIDERATIONS	25
59	Chapter 5 METHODOLOGY	26
60	5.1 Methodology	27
61	5.2 Design Considerations	27
62	5.2.1 Sensor choice, representation, and robustness	27
63	5.2.2 Temporal alignment and scoring	28
64	5.2.3 Real-time feedback, segmentation, and pedagogy	28
65	5.2.4 Accessibility, personalization, and evaluation	28
66	5.3 Theoretical Considerations	29
67	5.3.1 Human Pose Estimation	29
68	5.3.2 Human Action Recognition	30
69	5.4 Summary	30
70	References	31
71	Appendix A MEMBER SKILLSET IDENTIFICATION	33
72	Appendix B WORK BREAKDOWN STRUCTURECAPSTONE PROJECT	
73	ON OPERATIONAL TECHNOLOGIES	34



74

LIST OF FIGURES

75

1.1 Milestone Gantt Chart for Real-time Pose Estimation Dance Software 12

76

5.1 Methodology Flowchart 27

77

B.1 Work Breakdown Structure for Hans Capstone Project on Operational Technologies 35

78

79

B.2 Work Breakdown Structure for Nathan Capstone Project on Operational Technologies 36

80

81

B.3 Work Breakdown Structure for Gerald Capstone Project on Operational Technologies 36

82



83

LIST OF TABLES

84

1.1 Expected Deliverables per Objective 8

85

1.2 Operational Financial Plan 12

86

2.1 Summary of Reviewed Dance Pose Estimation and Recognition Studies . . . 20

87

A.1 Team Members' Programming Skills 33



88

ABBREVIATIONS

89

CV Computer Vision 2

90

HOG Histogram Of Oriented Gradients 2

91

CNN Convolutional Neural Network 2



De La Salle University

92

NOTATION



93

GLOSSARY

94

Tinikling

The traditional Filipino dance involving two bamboo sticks, where a dancer moves in and out of the rhythmically tapped sticks.

95

OpenCV

An open-source computer vision library widely used for real-time image capture and processing, including camera I/O, preprocessing, filtering, and contour extraction.

96

Ultraleap

A commercial infrared-based hand-tracking system that uses stereo cameras and near-infrared illumination to generate dense, low-latency 3D hand data.

97

MediaPipe

A framework for building multimodal applied machine learning pipelines, including computer vision models like hand gesture recognition.

98

Pose estimation

A computer vision technique used to detect human poses (such as hand or body positions) from images or videos, often used for gesture and movement analysis.

99

Operational Technologies

Programmable systems or devices that interact with the physical environment (or manage devices that interact with the physical environment). These systems/devices detect or cause a direct change through the monitoring and/or control of devices, processes, and events. Examples include industrial control systems, building management systems, fire control systems, and physical access control mechanisms.



De La Salle University

100

LISTINGS



De La Salle University

101

Chapter 1

102

INTRODUCTION



1.1 Background of the Study

Classical Computer Vision (CV) approaches used skin color segmentation, contour analysis, optical flow, and handcrafted descriptors (Histogram of Oriented Gradients (HOG), motion history images) to detect and classify gestures. Despite being simple and interpretable, those methods struggle with background variation and scale. The deep-learning era replaced handcrafted features with Convolutional Neural Network (CNN)s that learn hierarchical visual features directly from image data, yielding much higher accuracy for static hand-pose and short-sequence recognition tasks. Many recent capstone and journal implementations pair OpenCV (for capture/preprocessing) with CNN built and trained in TensorFlow/PyTorch to recognize a fixed vocabulary of gestures in real time. These hybrid pipelines are practical for capstone projects because OpenCV handles efficient frame processing while CNNs provide generalization across users and backgrounds. Furthermore, Operational Technologies plays a crucial role in deploying these systems in real-world applications where physical devices and processes are monitored and controlled, such as in industrial automation or building management systems, which benefit from enhanced gesture recognition. (Oudah et al., 2020)

Instead of classifying raw images, several high-performance systems first extract skeletal landmarks (e.g., MediaPipe's 21-point hand model) and feed those coordinates to a classifier (small CNN, MLP, or temporal model like LSTM). Landmark-based pipelines reduce sensitivity to background and scale and make models smaller and faster, which is ideal for mobile or AR deployment. Markerless commercial devices such as the Leap Motion Controller and Ultraleap cameras provide very accurate 3D joint data using IR illumination and multi-camera setups; those give superior fidelity but add hardware cost and integration



work. For a capstone aiming at broad deployability, a practical approach is to prototype with MediaPipe + OpenCV + CNN (or lightweight temporal model) and consider Ultraleap integration later for high-precision installations. (Zhang et al., 2020)

1.2 Prior Studies

Prior research on the topic at hand has shown substantial progress in the integration of pose estimation, computer vision, and interactive technologies for the sake of movement-based learning. For instance, a study by Kim et al. (2023) presents a human pose estimation method which integrates MediaPipe Pose with additional optimization techniques in order to improve its accuracy and robustness. The designed framework is capable of real-time landmark detection through the use of only a single RGB camera, while optimization methods such as smoothing filters and Kalman filtering are used to reduce jitter and improve the temporal consistency. Results depicted a high detection accuracy for various body parts, with its performance remaining stable under varying lighting and background. This shows MediaPipe’s suitability for real-time applications where both speed and stability is crucial, especially in aspects such as gesture recognition, sports monitoring, and motion analysis. Tharatipyakul et al. (2024) explores various deep learning-based human pose estimation techniques and their applications in health, rehabilitation, and human motion analysis. The paper looks into both 2D and 3D pose estimation. It is noted that 2D methods are widely used for real-time applications as they have much lower computational requirements in comparison to 3D. Deep convolutional neural networks and transformer-based models proved to significantly improve the landmark localization accuracy in comparison to classical approaches. Ultimately, the paper emphasized that integrating temporal information



enhances performance in sequential movement tasks, making these methods highly relevant for motion learning, sports training, and interactive systems. El Raheb et al. (2019) focuses on interactive dance learning systems and how such technology has the potential to support dance pedagogy through utilizing real-time feedback and structured interaction workflows. Multiple systems were analyzed and, afterwards, a framework was perfected which made use of motion capture, real-time analysis, and visual feedback in order to support users, who are both learners and instructors. Key interaction patterns were identified such as mirroring, guidance, and correction, which enhances the overall learning experience and, in turn, effectiveness. It also looks into usability considerations such as responsiveness, clarity of feedback, and alignment with existing teaching approaches, which is relevant to the creation of dance learning systems. Ultimately, such studies depict the intersection of pose estimation, feedback systems, and immersive interfaces, which lays a strong groundwork for future developments in digital dance education and interactive movement learning systems.

1.3 Problem Statement

To this day, the national dance of the Philippines known as ‘Tinikling’ continues to hold cultural significance among students, educators, and dance enthusiasts. However, despite its importance, those that aspire to learn the dance lack access to physical classes or qualified instructors be it due to geographical or time constraints. Existing methods of learning may be costly or unable to provide feedback to the student in real-time, which makes the learning process difficult for individuals in terms of practicing effectively on their own. Such a gap highlights the need for a much more accessible, interactive, and accurate tool



170 which would be able to guide learners remotely in an efficient manner and, ultimately,
171 ensuring that tradition is preserved and passed on to future generations.

172 **1. PS1:**

- 173 • The ideal scenario for our intended audience (students, educators, and dance
174 enthusiasts) is to have an intuitive and interactive learning tool that facilitates
175 the practice of Tinikling, the traditional Filipino dance. This tool should provide
176 real-time feedback on users' dance movements, enabling them to learn and
177 improve their technique. The desired state includes accessibility to the tool on
178 various devices (e.g., desktop, mobile) with a user-friendly interface and a high
179 level of accuracy in tracking the dance steps. Additionally, it should support
180 personalized feedback, enabling users of all skill levels to progress and feel
181 engaged in learning this cultural heritage.

182 **2. PS2:**

- 183 • Currently, learning Tinikling requires access to physical dance classes or in-
184 structors, which are often limited by geographical location, financial resources,
185 or time constraints. For individuals unable to attend such classes, the lack of af-
186 fordable and effective learning tools becomes a significant barrier. Additionally,
187 existing dance-learning technologies are either costly, relying on specialized
188 hardware, or lack the immediacy of real-time feedback, making it difficult
189 for learners to practice and perfect their movements without direct instructor
190 guidance.



- The pain point is that students who want to practice Tinikling at home or in remote areas are unable to receive real-time guidance or feedback, leading to slower progress, incorrect technique, and a loss of motivation.

3. PS3:

- Without a tool that offers immediate feedback and a clear learning path, students practicing Tinikling on their own are likely to struggle with incorrect movements, which may lead to frustration. Over time, this lack of progress could result in a lack of confidence, disengagement from the learning process, and ultimately, the inability to learn the dance correctly. Furthermore, the absence of accessible learning tools risks the loss of cultural knowledge and the fading of the Tinikling tradition, especially among younger generations who may not have easy access to traditional learning methods.

1.4 Objectives and Deliverables

1.4.1 General Objective (GO)

- GO: To design and implement a real-time Pose estimation-based Tinikling learning application;

1.4.2 Specific Objectives (SOs)

- SO1: To develop a real-time pose estimation pipeline that captures dancers' movements using a webcam, detects key skeletal landmarks, and analyzes Tinikling steps



with at least 30 frames per second (fps) processing speed and $\geq 90\%$ detection accuracy.;;

- SO2: To make the pose estimation model robust to lighting, background clutter, and user variation through dataset collection and augmentation and, landmark-based representations while maintaining a minimum pose detection accuracy of 85% ;
- SO3: To design and integrate a scoring and feedback system that evaluates user performance by aligning poses with reference choreographies, providing numerical scores (0–100) and step-by-step accuracy breakdown within 1 second after performance.;
- SO4: To develop a desktop-based user interface that provides real-time visual cues, instant audio or visual feedback within 200 ms, and performance scores to guide learners effectively during practice sessions.;
- SO5: To evaluate the system's performance and usability through controlled testing with at least 10 participants, measuring pose estimation accuracy, latency, and user satisfaction ($\geq 80\%$ positive feedback) using standardized questionnaires and performance metrics.;

1.4.3 Expected Deliverables

1.5 Significance of the Study

This capstone project focuses on the development of a Tinikling learning application through the integration of pose estimation and human action recognition. The setup consists of a webcam, laptop, and two bamboo sticks for the Tinikling dance. Such a setup offers



TABLE 1.1 EXPECTED DELIVERABLES PER OBJECTIVE

Objectives	Expected Deliverables
GO: To develop a real-time pose estimation-based Tinikling learning application	<ul style="list-style-type: none"> • Prototype of Tinikling learning application. • Documentation and user manual.
SO1: To develop a real-time pose estimation pipeline that captures the movement of dancers through a webcam, detects skeletal keypoints, and analyzes poses for Tinikling steps with low latency and high accuracy.	<ul style="list-style-type: none"> • Optimized skeletal keypoints detection for Tinikling steps. • Implementation of webcam-based pose estimation pipeline. • Performance evaluation results.
SO2: To make the pose estimation model robust to lighting, background clutter, and user variation through dataset collection, augmentation, and landmark-based representations.	<ul style="list-style-type: none"> • Augmented dataset covering varied lighting, backgrounds, and user types. • Enhanced landmark-based model with robustness improvements. • Comparative performance evaluation report.
SO3: To design and integrate a scoring and feedback system that evaluates users' dance accuracy in a post-performance review by aligning user poses with reference choreographies.	<ul style="list-style-type: none"> • Scoring and feedback algorithm. • Tinikling choreography database. • Post-performance scoring output with accuracy metrics.
SO4: To evaluate the system's performance and usability through controlled testing with dancers or students, measuring accuracy, latency, and user experience for future refinement and educational deployment.	<ul style="list-style-type: none"> • Conducted controlled testing with participants. • Collected performance and usability metrics. • Evaluation report with recommendations for improvement.

affordability and accessibility benefits for users. Ultimately, it contributes to the field of both pose estimation and human action recognition by demonstrating a successful integration of the two in a live setup.

1.5.1 Technical Benefit

1. Enables real-time pose estimation and post-performance feedback, improving accuracy and efficiency throughout the learning process.



2. Low-cost software-based learning tool which uses a webcam and desktop computer rather than expensive motion capture equipment.

1.5.2 Social Impact

- Promotes cultural preservation by making Tinikling more accessible through interactive applications.
- Increases student engagement and participation via gamified learning.
- Supports remote or in-classroom instruction by enabling technology-assisted dance education.

1.5.3 Environmental Welfare

- Utilizes existing and widely available hardware such as webcams and desktop computers rather than new specialized equipment, which ultimately lessens electronic waste.
- Encourages digital preservation of cultural heritage, lessening reliance on physical materials or infrastructure.

1.6 Assumptions, Scope, and Delimitations

1.6.1 Assumptions

1. Pose landmarks from webcams with standard RGB resolutions such as 720p, 1080p, and 4K or low-cost depth sensors provide sufficient fidelity to represent Tinikling



- 254 movements for temporal alignment and scoring.
- 255 2. Choreography can be divided into short, labeled segments that enable reliable match-
- 256 ing and targeted feedback.
- 257 3. Dynamic Time Warping or a constrained variant will handle tempo variation robustly
- 258 for temporal alignment.
- 259 4. A brief per-user calibration step will improve scoring consistency.

260 1.6.2 Scope

- 261 1. Cover automatic pose estimation, sequence alignment, and segment-level scoring for
- 262 Tinikling.
- 263 2. Accept landmark or depth inputs and provide immediate on-device cues during
- 264 performance.
- 265 3. Produce a higher-precision final score after a more detailed pass.
- 266 4. Use self-sourced Tinikling videos for model training when no public dataset exists.
- 267 5. Benchmark against general dance datasets where appropriate.
- 268 6. Report sensor-based metrics and simple user measures such as perceived accuracy
- 269 and engagement.

270 1.6.3 Delimitations

- 271 1. Will not perform detailed facial or hand mesh reconstruction.



2. Will not replace multi-camera motion capture for research-grade kinematics.

3. Will not guarantee reliable results under heavy occlusion, very low light, extreme off-axis views, or when clothing blends with the background.

4. Will not attempt full generalization to all body shapes without additional data and tuning.

5. Limits reflect known sensor and algorithm constraints and the aim to produce a practical, lightweight prototype.

1.7 Description and Methodology of the Capstone Project on Operational Technologies

1. Phase 1: Model Development serves as a precursor for Phase 2 wherein the specifics of the model, libraries, and environment to use are defined. In total, Phase 1 would last 4 weeks spanning from week 4 to 7. The bulk of the research for the project would be carried out during this phase. The dataset to be used for training would be collected during this phase as well.

2. Phase 2: Model Training consists of training the model using the dataset collected in the previous phase. This phase will largely consist of testing and improving the resulting model. Tests would be conducted using the group members as dancers. This phase also includes the optimization of the model for real-time detection simultaneously with the music. In total, this phase would last 4 weeks spanning from week 8 to 11.



3. Phase 3: UI/UX Development consists of the integration of the trained model with a user interface. Once integrated final testing and refinement of the final program would be carried out. The final output would be presented as well during this phase along with the finalization of the documentation. This phase would last for 3 weeks spanning from week 11 to 13.

1.8 Estimated Work Schedule and Budget

1.8.1 Milestones and Gantt Chart

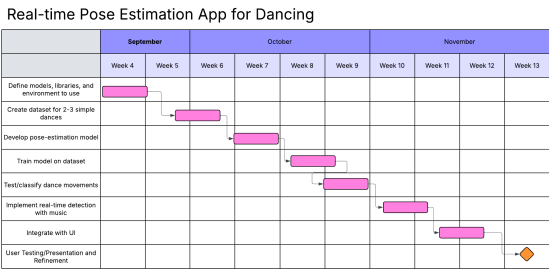


Fig. 1.1 Milestone Gantt Chart for Real-time Pose Estimation Dance Software

1.8.2 Budget

Given that the capstone project largely consists of software, apart from the use of a laptop for both programming, as well as actual implementation and usage of the dance program, the only expense to consider would be for that of a Webcam, which is already owned.



TABLE 1.2 OPERATIONAL FINANCIAL PLAN

Item	Price
Webcam	P1,850
4pc. Tinikling Sticks	P110
Total	P1,960

1.9 Overview of the Capstone Project on Operational Technologies

This capstone project focuses on developing a real-time pose estimation-based learning application for Tinikling, the Philippine national dance. It integrates computer vision and machine learning techniques in order to create an interactive learning platform that provides performance scoring to users. The project utilizes webcams and MediaPipe-based skeletal landmark extraction to analyze users' movements relative to reference choreography. Unlike expensive motion capture systems, this setup uses low-cost and accessible hardware, making the system practical for classroom, cultural, and home use. The system emphasizes cultural preservation by modernizing Tinikling education through technology. It enables students to learn and practice the dance interactively, provides technical benefits such as real-time feedback without costly sensors, and supports social and environmental goals through cultural engagement and sustainable use of existing hardware.



316

Chapter 2

317

LITERATURE REVIEW



2.1 Existing Work

A study by Venkatrayappa et al. (2024) focused on surveying the various existing 3D human body pose and shape estimation techniques, given its crucial nature in fields such as augmented or virtual reality, healthcare and fitness technology, and virtual retail. The solutions explored consisted of mainly three types of inputs, which were single images, multi-view images, and videos. Various issues pertaining to dance, such as fast motion, occlusion, and unusual poses were analyzed to see how each affected the performance of each method. The specific models consisted of SMPL -A, SMPL -X, MANO, STAR, FLAME, which are optimization-based models, as well as HMR, VIBE, SPIN, PARE, EXPOSE, and PHALP, which are deep learning-based models. SMPL was found to be beneficial in terms of realistic body representation, efficiency for real time applications, and wide availability, however it has limitations in areas pertaining to facial and hand modeling, as well as representation of ethnic diversity. SMPL -X proved to provide several advantages such as facial expressions, hand gestures, and improved expressiveness. Its limitations, however, consisted of simplified hand modeling and its limited pose variability. MANO offers detailed hand gesture modeling and realistic hand deformations, but has limitations due to its focus being exclusively on the modeling of hands, as well as computational challenges. STAR leverages sparse coding and temporal modeling, which allowed for a much more powerful framework for pose estimation., depicting state-of-the-art results throughout various benchmarks and practical implementations in sports analysis, human-computer interaction, and VR. FLAME was advantageous when it comes to computational efficiency, which made it suitable for real-time applications of pose estimation. As for its limitations, it primarily focuses on facial and lip modeling, which introduces complexity



341 and potential computational challenges. MANO. HMR produces richer and more useful
 342 mesh representation, which is parameterized by shape and 3D joint angles. The network
 343 implicitly learns the angle limits of each joint. As such its use is discouraged for people
 344 with unusual body shapes. Its re-projection loss is highly under-constrained and it needs
 345 adversarial supervision in order to avoid unrealistic outputs. VIBE makes use of CNNs,
 346 RNNs and GANs, as well as a self-attention layer in order to achieve state-of-the-art
 347 results. A motion discriminator is used to help produce more realistic motion. Ultimately,
 348 the model is a standard SMPL body model format with sequences of poses and shape
 349 parameters. SPIN makes use of a self improving loop wherein better fits allow the network
 350 to train in a much more efficient manner while better initial estimates from the network
 351 aids the optimization routine in order to result in better fits. PARE consists of a guided
 352 attention mechanism which exploits information on visibility of individual body parts all
 353 the while leveraging information from neighboring body parts in order to predict parts
 354 which are occluded. EXPOSE includes body, face, and hand estimation. It is able to
 355 estimate expressive 3D humans in a much more accurate manner in comparison to existing
 356 optimization methods at only a fraction of the computational costs. PHALP out performs
 357 all of the aforementioned methods. Despite this, it still has its limitations as well such as its
 358 reliance on a single camera, which may lead to issues such as occlusion and motion blur. It
 359 may also not work well in low-light conditions or when a person's clothes is a similar color
 360 to that of the background. Lastly, it also requires a significant amount of computational
 361 resources, which may make it not suitable for real-time applications.

362 A study by Protopapadakis et al. (2018), analyzes the effectiveness of various classifica-
 363 tion techniques in recognizing different dance types based on motion-capture skeleton data.
 364 Classifiers explored consisted of k-Nearest Neighbors (k-NN), Naïve Bayes, Discriminant



365 Analysis, Classification Trees, Random Forests (TreeBagger), Support Vector Machines
366 (SVMs), and Ensemble Classifiers. Poses are identified through the use of body joints via
367 Kinect sensor. The data set used consisted of various dances such as Enteka, Kalamatianos,
368 Syrtos (Two-beat), Syrtos (Three-beat). The kinect was used to capture skeletal joining
369 data, to which feature extraction techniques such as principal component analysis and frame
370 differencing were used in order to improve the classification accuracy. Ultimately, results
371 showed that k-nearest neighbors and random forests are the best-performing classifiers
372 among those that were explored. It was also proposed that the use of mulit-sensor or
373 multimodal data may serve as a potential solution for issues specific to pose recognition in
374 dance such as occlusion and complex movement patterns.

375 A study by Zhao et al. (2025), looks into dance pose estimation and introduces the model
376 DanceFormer. DanceFormer is a transformer-based model for dance pose estimation which
377 makes use of the Vision Transformer, Time Series Transformer, and an edge computation
378 layer in order to achieve a deep fusion of multimodal features and to overall increase
379 its accuracy and real-time performance. The AIST and DanceTrack datasets were used
380 throughout the experimentation. Results showed that DanceFormer out performs other
381 models, with it achieving a pose estimation accuracy or MPJPE of 18.4mm and 20.1mm,
382 as well as a multi-object tracking accuracy or MOTA of 92.3% and 89.5%. It is also
383 suitable for real-time processing in even low-resource with an average latency of 35.2ms.
384 Ultimately, it serves as an efficient, precise and real time solution for rather complex dance
385 scenarios. It also has applications in a much more broad sense be it in dance education or
386 in real-time motion analysis.

387 A study by Lei et al. (2023) discusses dance movement recognition based on gesture. A
388 low accuracy traditional dance movement recognition algorithm based on human posture



estimation was proposed. PAFs algorithm was used in order to recognize the spatial skeleton nodes and connections of joints in the human body. The pose of the body is estimated based on the movement of the spatial skeleton. Once the information on the detected posture is preprocessed and its features are extracted, LSTM time series algorithm was used in order to classify and recognize certain dance movements. Ultimately, results showed that the proposed algorithm has the capacity to reliably identify dance movements based on the skeleton nodes. It was able to achieve a recognition accuracy and recall rate upwards of 85% for the different movement categories. As for its recognition accuracy of curtsy movement, it achieved upwards of 95.2%.

Tölgyessy et al. (2021) present a detailed evaluation of Kinect v1, Kinect v2, and Azure Kinect skeleton tracking, analyzing joint-level error distributions and repeatability across distances and orientations. Their results highlight degradation in accuracy under occlusion, off-axis angles, and larger working distances, conditions typical of casual living-room dance setups. The findings underline both the potential and the limits of Kinect-class sensors, suggesting that practical applications often require either sensor fusion and smoothing to handle jitter or a focus on more reliable joints for robust real-time scoring.

Lin (2015) investigate how interactive feedback design influences user motivation in the context of Just Dance. Their study demonstrates that timely, clear cues significantly improve engagement, perceived competence, and sustained participation, with direct effects on physical activity outcomes. These findings show that feedback modalities and latency are as critical as recognition accuracy in shaping the player experience, emphasizing the importance of immediate, multimodal responses in dance or pose-based teaching applications.

Yu and Xiong (2019) propose and validate a Dynamic Time Warping method for



evaluating rehabilitation exercises tracked with Kinect. Their algorithm successfully aligns noisy, tempo-varying motion with reference trajectories, producing reliable correctness scores even with partial occlusion. Applied to dance or short choreographies, DTW offers a robust foundation for handling tempo shifts and timing variation, supporting sequence-based scoring that is more forgiving than strict frame-to-frame comparison.

Rallis et al. (2019) compare Kinect II with the high-precision Vicon system in the context of choreography retrieval and analysis, using trajectory similarity measures such as DTW. While Kinect data contain noise and smoothing artifacts, the study shows that trajectory-level patterns remain useful when algorithms are designed to tolerate sensor bias. Their results support the use of low-cost consumer sensors, including RGB landmark pipelines, in applications where robust temporal alignment and trajectory modeling can offset hardware limitations.

Human pose estimation (HPE) has become an important area of study due to its applications in action recognition, sports, and performing arts. Xu, Zou, and Lin (2022) introduced the Adaptive Hypergraph Neural Network (AD-HNN), which captures high-order semantic dependencies among joints to improve multi-person pose estimation, particularly in handling occlusion and pose variability. In dance analysis, Ju (2025) applied deep learning with ResNet-152 and HR-Net to enhance dance pose recognition, addressing class imbalance and improving classification accuracy through global–local feature fusion.

For cultural preservation, motion capture (MoCap) has been widely adopted. Rizhan et al. (2025) demonstrated the use of MoCap to develop authentic motion templates for Malay folk dances, ensuring accuracy and authenticity in preserving intangible cultural heritage. In addition, Büyükgökoğlu and Uğuz (2025) developed a performance evaluation system for Turkish folk dances using deep learning–based pose estimation (e.g., Mediapipe, YOLO, LSTM), enabling objective assessment compared to traditional jury scoring.



TABLE 2.1 SUMMARY OF REVIEWED DANCE POSE ESTIMATION AND RECOGNITION STUDIES

Paper	Focus	Methodology	Results
<i>Venkatrayappa et al. (2024)</i>	Evaluates 3D human body pose and shape estimation methods for contemporary dance	Comparative survey of model families: SMPL(-A/X), MANO, STAR, FLAME (optimization-based) and HMR, VIBE, SPIN, PARE, EX-POSE, PHALP (learning-based); analysis by input modality (single-image, multi-view, video).	PHALP strong overall; SMPL-X improves expressiveness; STAR excels in temporal modeling. Common limits: occlusion, lighting, and compute needs.
<i>Protopapadakis et al. (2018)</i>	Identifies dance types from motion-capture skeletal data	Kinect skeletal features with PCA and frame differencing; compared classifiers (k-NN, Naïve Bayes, LDA/DA, decision trees, Random Forest, SVM, ensembles).	k-NN and Random Forest performed best; multimodal or multi-sensor data recommended to handle occlusion.
<i>Zhao et al. (2025)</i>	Real-time pose estimation for complex dances	Hybrid architecture: Vision Transformer + Time-Series Transformer trained on AIST and DanceTrack datasets.	MPJPE: 18.4 mm / 20.1 mm; MOTA: 92.3% / 89.5%; latency \approx 35.2 ms (real-time capable).

Continued on next page



Table 2.1 (continued)

Paper	Focus	Methodology	Results
<i>Lei et al. (2023)</i>	Improves recognition accuracy for traditional dance movements	Keypoint detection via Part Affinity Fields (PAFs); temporal modeling with LSTM classifiers.	>85% overall accuracy; 95.2% for curtsey movements.
<i>Zheng et al. (2023)</i>	Deep-learning approaches for pose design and recognition	Backbone fusion (ResNet-152 + HR-Net) with global-local feature fusion and class-imbalance handling strategies.	Reported metrics: accuracy 0.9870; precision 0.9851; sensitivity 0.9873; F1 0.9861; Kappa 0.9841.
<i>Xu et al. (2022)</i>	Multi-person pose estimation from single images	Two-stage Adaptive Hypergraph Neural Network (keypoint localization + adaptive hypergraph) with SIC module; end-to-end training.	Achieves state-of-the-art performance on MS-COCO, MPII, and CrowdPose benchmarks.
<i>Tölgyessy et al. (2021)</i>	Quantifies joint-level accuracy and repeatability across Kinect sensors	Controlled robotic-manipulator and figurine measurements across positions; compared Kinect v1, v2, Azure Kinect (NFOV/WFOV).	Azure NFOV shows highest accuracy (0.8–1.9 mm SD); joint failures 15–30% under occlusion; performance declines at long ranges.
<i>Lin (2015)</i>	Effects of feedback and controller use in dance exergames	2×2×2 factorial experimental design (feedback × controller × sex); 129 participants; 12-minute sessions.	Mean HR ≈ 109 bpm; immediate/clear feedback increased engagement and perceived competence.

Continued on next page




Table 2.1 (continued)

Paper	Focus	Methodology	Results
<i>Yu and Xiong (2019)</i>	DTW-based scoring for rehabilitation/exercise movements	Dynamic Time Warping on 8 bone vectors + body orientation; algorithm converting DTW distance to a 0–100% performance score.	Scores correlated strongly with expert ratings ($r \approx 0.86$); robust to tempo variation and some occlusion.
<i>Rallis et al. (2019)</i>	Choreographic pattern analysis from heterogeneous capture systems	Trajectory-based DTW similarity for choreography; sensor comparison (VICON vs Kinect); experiments on smoothing and joint selection.	Kinect is noisier but DTW reduces sensor bias; smoothing and selective joint use improve retrieval accuracy.
<i>Sun and Song (2025)</i>	Pose estimation in complex dance scenes	Enhanced HRNet backbone with improved feature extraction and robustness modules for cluttered scenes.	Improved keypoint accuracy and robustness under occlusion/clutter.
<i>Büyükgökoğlu and Uğuz (2025)</i>	Deep-learning-based scoring for Turkish folk dance	Webcam capture; pose extraction via MediaPipe/YOLO; comparative models: DTW, TLCC, LSTM, Siamese networks.	LSTM produced higher scores (≈ 68.43 , $MSE \approx 56.11$) vs DTW (≈ 60.64 , $MSE \approx 139.32$); system is sensitive to camera angle.

2.2 Lacking in the Approaches

These studies show the potential of pose estimation and deep learning for advancing both modern dance movement design and traditional folk dance preservation. How-

		2. Literature Review	
	 <div>De La Salle University</div>		
441	ever, there is little to no research in the Philippines that applies pose estimation to folk		
442	dances—particularly Tinikling—representing a significant gap and opportunity for future		
443	exploration.		
444	2.3 Summary		
445	Research on human pose estimation (HPE) spans multiple applications including AR/VR,		
446	healthcare, and dance. Optimization- and deep learning–based models (e.g., SMPL, SMPL-		
447	X, HMR, VIBE, SPIN, PARE, EXPOSE, PHALP) have been studied for realistic 3D		
448	body reconstruction (Venkatrayappa et al., 2024). Dance classification has been explored		
449	using skeleton data and machine learning classifiers like k-NN and Random Forest (Pro-		
450	topapadakis et al., 2018). Transformer-based models such as DanceFormer achieve high		
451	accuracy and real-time performance in dance pose estimation (Zhao et al., 2025), while		
452	PAF- and LSTM-based algorithms improve movement recognition (Lei et al., 2023). Kinect		
453	studies reveal both potential and limits in low-cost motion capture (Tölgyessy et al., 2021;		
454	Rallis et al., 2019), while feedback and sequence-alignment approaches (Lin et al., 2015;		
455	Yu & Xiong, 2019) highlight the importance of interactivity and temporal robustness.		
456	Recent work integrates advanced neural networks for pose estimation, such as adaptive		
457	hypergraphs (Xu et al., 2022), deep feature fusion for dance poses (Ju, 2025), MoCap		
458	for authentic folk dance templates (Rizhan et al., 2025), and deep learning systems for		
459	evaluating Turkish folk dance (Büyükgökoğlu & Uğuz, 2025).		
			23



460

Chapter 3

461

THEORETICAL CONSIDERATIONS



462

Chapter 4

463

DESIGN CONSIDERATIONS



De La Salle University

464

Chapter 5

465

METHODOLOGY

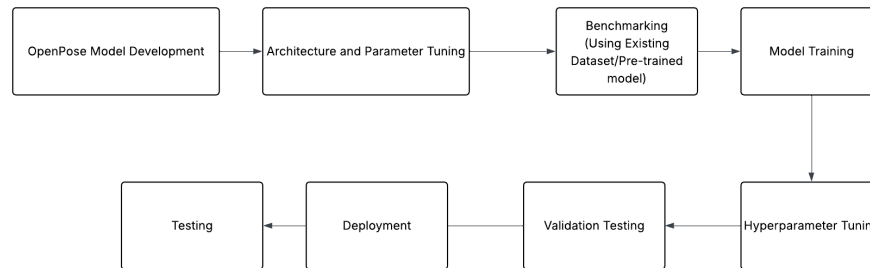


Fig. 5.1 Methodology Flowchart

5.1 Methodology

5.2 Design Considerations

5.2.1 Sensor choice, representation, and robustness

A study by Tölgyessy, Dekan, and Chovanec (2021) demonstrated that Kinect-family depth sensors produce explicit 3D skeletons and give higher joint fidelity in controlled settings, but the accuracy falls with occlusion, off-axis views, and increased distance. Zhang et al. (2020) described MediaPipe, which yields compact 2D/3D landmark coordinates from ordinary RGB cameras and runs in real time on mobile devices. Therefore, designers often choose landmarks for rapid, lightweight prototypes and mobile deployment, and reserve depth or IR systems for installation-grade fidelity when hardware is available. To reduce real-world failure modes, practitioners apply photometric and background augmentation and synthetic occlusions during training, and they add a short calibration step so system metrics align with an individual user's range of motion.



5.2.2 Temporal alignment and scoring

Dance is a temporal activity and should be compared as a sequence rather than as isolated frames. Yu and Xiong (2019) demonstrate that Dynamic Time Warping (DTW) can align noisy, tempo-varying Kinect skeleton sequences and convert DTW distances into meaningful performance scores. Rallis et al. (2019) apply DTW to choreographic trajectories and show it can match patterns across high-precision (VICON) and low-cost (Kinect) capture systems. Thus, a practical scoring pipeline first aligns sequences with DTW (or a constrained variant) and then evaluates local spatial metrics such as joint-angle differences or normalized trajectory distances to produce interpretable, per-segment correctness scores.

5.2.3 Real-time feedback, segmentation, and pedagogy

Lin (2015) finds that immediate, clear feedback in dance exergames improves engagement and supports learning. Zhang et al. (2020) show that on-device landmark extraction can run at real-time rates suitable for low-latency feedback. Combining these results suggests a two-tier runtime design: use a fast, coarse matcher (enabled by on-device landmarks) for instant cues, and run a slower, higher-precision alignment and scoring pass for final grading. Breaking choreography into short labeled segments also simplifies alignment and reduces error accumulation; Rallis et al. (2019) illustrate that segment- or trajectory-level matching better supports choreographic retrieval and per-segment feedback.

5.2.4 Accessibility, personalization, and evaluation

Yu and Xiong (2019) convert DTW distances into calibrated percentage scores, which supports per-user calibration and comparison against an individualized baseline. Tölgyessy



et al. (2021) recommend measuring sensor-level metrics such as joint error and dropout rates when choosing a capture modality. Therefore, system designs should include adjustable sensitivity, alternate gesture mappings, and user profiles, and evaluation should combine sensor metrics (joint error, dropout, latency) with human-centered measures (perceived accuracy, engagement, and learning gain) to justify architecture and scoring choices.

5.3 Theoretical Considerations

5.3.1 Human Pose Estimation

Human pose estimation is the process of predicting the pose of human body parts. The data are typically stemming from RGB images or videos. Given that certain motions are motivated by human actions, detecting poses is a critical aspect of human action recognition (Song et al., 2021). It has a wide range of applications such as human-computer interaction, motion analysis, augmented reality, and virtual reality. The resulting output of human pose estimation is a skeleton-like representation of the human body consisting of nodes and limbs (Zheng et al, 2020)). There are 2 main types of human pose estimation, namely 2D and 3D. 2D pose estimation consists of predicting the posture of each of the body's key points in a 2D plane, considering the X and Y axis. As for 3D pose estimation, it considers the Z axis, situating each point in a 3D space. It goes without saying that the 3D estimation would be much more difficult in comparison to 2D estimation in process or complexity due to underlying issues which may manifest such as noisy backgrounds, clothing, lighting, undetected joints, or occlusion (Ben Gamra & Akhloufi, 2021).



5.3.2 Human Action Recognition

Human action recognition, otherwise known as HAR, is the process of detecting human actions in order to classify them through single sensor data, RGB image or video data, or three-dimensional depth and inertial data (Sakar et al., 2022). In the field of computer vision, one of the most challenging aspects of it is the automatic and precise identification of human activity. Over the years, there has been a significant increase in feature learning-based representations for human action recognition as a result of the widespread utilization of deep learning-based features. There are various applications of Human action recognition. For instance, automated surveillance systems make use of AI and machine learning algorithms in order to identify human actions for the sake of safety and security. Such a task, however, is made difficult due to various factors such as changing online environments, occlusion, different viewpoints, execution pace and biometric change. Not only this, but the human body also varies from person to person in factors such as size, appearances, and shapes. However, advancements in Convolutional Neural Networks, otherwise known as CNNs, resulted in significant progress for human action recognition through improvements on classification, segmentation and object detection. This largely applies more on image-related tasks rather than videos as neural network models struggle to capture temporal information in videos due to a lack of substantial datasets (Morshed et al., 2022).

5.4 Summary

Provide the gist of this chapter such that it reflects the contents and the message.



REFERENCES

- Büyükökoğlu, E. and Uğuz, S. (2025). Development of a performance evaluation system in turkish folk dance using deep learning-based pose estimation. *Tehnički vjesnik – Technical Gazette*, 32:1817–1824.
- El Raheb, K., Stergiou, M., Katifori, A., and Ioannidis, Y. (2019). Dance interactive learning systems: A study on interaction workflow and teaching approaches. *ACM Computing Surveys*, 52:1–37.
- Kim, J.-W., Choi, J.-Y., Ha, E.-J., and Choi, J.-H. (2023). Human pose estimation using mediapipe pose and optimization method based on a humanoid model. *Applied Sciences*, 13(4).
- Lei, P., Li, N., and Liu, H. (2023). Dance movement recognition based on gesture.
- Lin, J.-H. (2015). Just dance: The effects of exergame feedback and controller use on physical activity and psychological outcomes. *Games for Health Journal*, 4(3):183–189. PMID: 26182062.
- Oudah, M., Al-Naji, A. A., and Chahl, J. (2020). Hand gesture recognition based on computer vision: A review of techniques. *Journal of Imaging*, 6:73.
- Protopapadakis, E., Voulodimos, A., Doulamis, A., Camarinopoulos, S., Doulamis, N., and Miaoulis, G. (2018). Dance pose identification from motion capture data: A comparison of classifiers. *Technologies*, 6.
- Rallis, I., Protopapadakis, E., Voulodimos, A., Doulamis, N., Doulamis, A., and Bardis, G. (2019). Choreographic pattern analysis from heterogeneous motion capture systems using dynamic time warping. *Technologies*, 7(3):56.
- Sun, J. and Song, L. (2025). Dance movement pose estimation in complex scenes based on improved high-resolution networks. *Taiwan Ubiquitous Information*, 10(1). Journal of Network Intelligence.
- Tharatipyakul, A., Srikaewsiew, T., and Pongnumkul, S. (2024). Deep learning-based human body pose estimation in providing feedback for physical movement: A review. *Heliyon*, 10(17):e36589.
- Tölgyessy, M., Dekan, M., and Chovanec, v. (2021). Skeleton tracking accuracy and precision evaluation of kinect v1, kinect v2, and the azure kinect. *Applied Sciences*, 11(12):5756.
- Venkatrayappa, D., Treméau, A., Muselet, D., and Colantoni, P. (2024). Survey of 3d human body pose and shape estimation methods for contemporary dance applications.
- Xu, X., Zou, Q., and Lin, X. (2022). Adaptive hypergraph neural network for multi-person pose estimation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(3):2955–2963.
- Yu, X. and Xiong, S. (2019). A dynamic time warping based algorithm to evaluate kinect-enabled home-based physical rehabilitation exercises for older people. *Sensors*, 19(13):2882. PubMedID: 31261746.



- 574 Zhang, F., Bazarevsky, V., Vakunov, A., Tkachenka, A., Sung, G., Chang, C.-L., and Grundmann,
575 M. (2020). Mediapipe hands: On-device real-time hand tracking. *CoRR*.
- 576 Zhao, H., Du, B., Jia, Y., and Zhao, H. (2025). Danceformer: Hybrid transformer model for real-time
577 dance pose estimation and feedback. *Alexandria Engineering Journal*, 121:66–76.
- 578 Zheng, C., Wu, W., Chen, C., Yang, T., Zhu, S., Shen, J., Kehtarnavaz, N., and Shah, M. (2023).
579 Deep learning-based human pose estimation: A survey.



581

582

Appendix A MEMBER SKILLSET IDENTIFICATION

TABLE A.1 TEAM MEMBERS' PROGRAMMING SKILLS

Member	Model Dev.	UI Design	Source Control (GitHub)	Problem Solving & Opt.	Python
Hans	Intermediate	Novice	Expert	Intermediate	Intermediate
Gerald	Intermediate	Basic	Novice	Intermediate	Intermediate
Nathan	Intermediate	Novice	Novice	Intermediate	Intermediate



De La Salle University

583

Appendix B

584

WORK BREAKDOWN

585

STRUCTURECAPSTONE PROJECT ON

586

OPERATIONAL TECHNOLOGIES

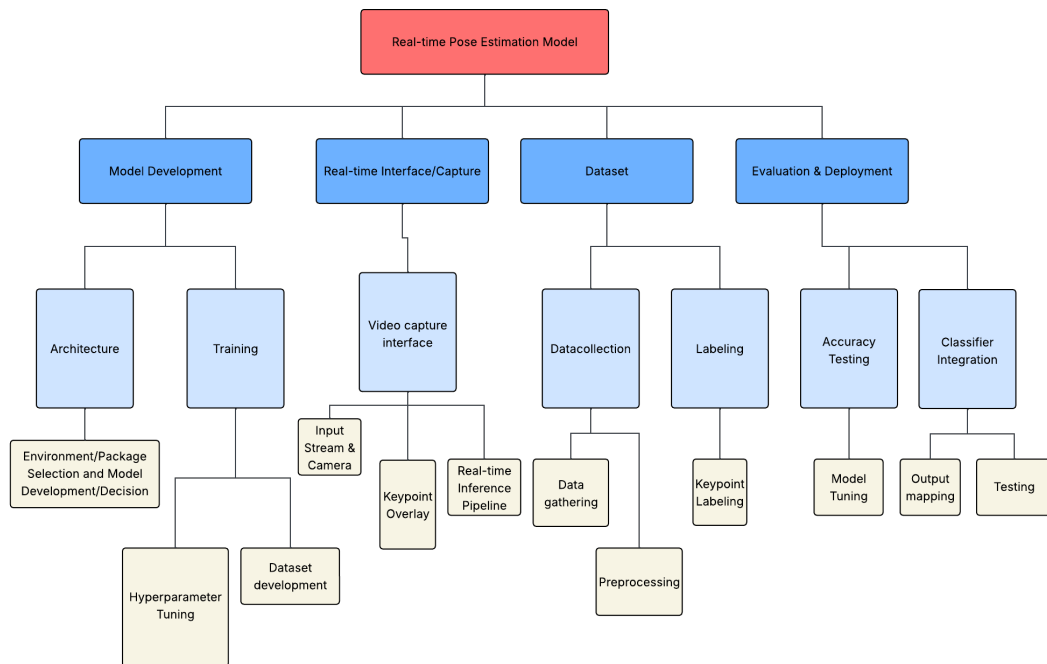


Fig. B.1 Work Breakdown Structure for Hans Capstone Project on Operational Technologies

B. Work Breakdown StructureCapstone Project on Operational Technologies



De La Salle University

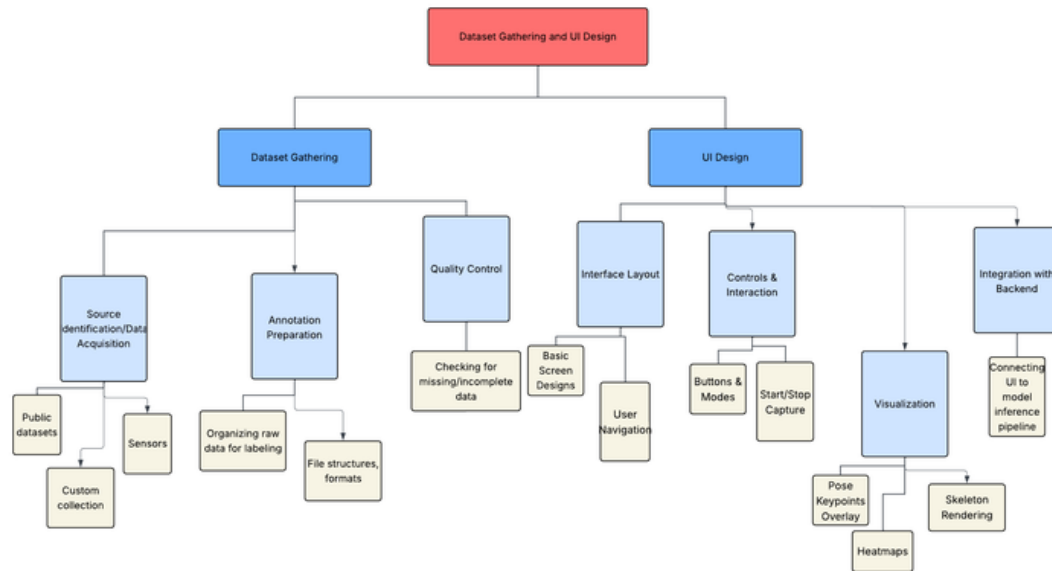


Fig. B.2 Work Breakdown Structure for Nathan Capstone Project on Operational Technologies

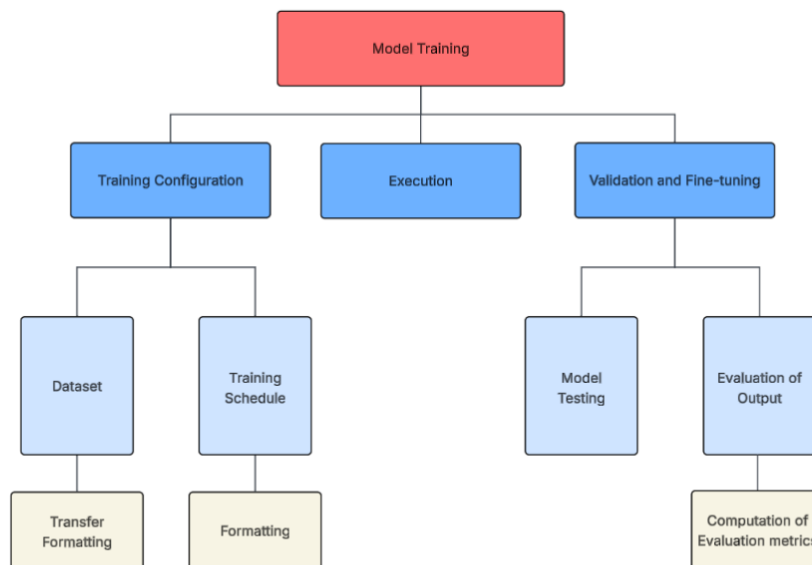


Fig. B.3 Work Breakdown Structure for Gerald Capstone Project on Operational Technologies