

1       ational Technologies51figure.caption.30



3       A Real-time Pose Estimation Application for Tinikling

4

---

5       A Capstone Project on Operational Technologies

6              Presented to the Faculty of the  
7              Department of Electronics and Computer Engineering  
8              Gokongwei College of Engineering  
9              De La Salle University

10

---

11       In Partial Fulfillment of the  
12              Operational Technologies

13

---

14              by

15       CALAGUIAN Nathan Raekel L.  
16              ELLAR Gerald Antonio P.  
17              MAHAIT Hans  
18  
19

20       November, 2025



21

## ABSTRACT

22

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



## TABLE OF CONTENTS

23	<b>Abstract</b>	ii
25	<b>Table of Contents</b>	iii
26	<b>List of Figures</b>	iv
27	<b>List of Tables</b>	v
28	<b>Abbreviations and Acronyms</b>	vi
29	<b>Notations</b>	vii
30	<b>Glossary</b>	viii
31	<b>Listings</b>	ix
32	<b>Chapter 1 INTRODUCTION</b>	1
33	1.1 Background of the Study . . . . .	2
34	1.2 Prior Studies . . . . .	3
35	1.3 Problem Statement . . . . .	4
36	1.4 Objectives and Deliverables . . . . .	6
37	1.4.1 General Objective (GO) . . . . .	6
38	1.4.2 Specific Objectives (SOs) . . . . .	6
39	1.4.3 Expected Deliverables . . . . .	7
40	1.5 Significance of the Study . . . . .	9
41	1.5.1 Technical Benefit . . . . .	9
42	1.5.2 Social Impact . . . . .	9
43	1.5.3 Environmental Welfare . . . . .	10
44	1.6 Assumptions, Scope, and Delimitations . . . . .	10
45	1.6.1 Assumptions . . . . .	10
46	1.6.2 Scope . . . . .	10
47	1.6.3 Delimitations . . . . .	11
48	1.7 Description and Methodology of the Capstone Project on Operational Technologies . . . . .	12
49	1.8 Estimated Work Schedule and Budget . . . . .	13
50	1.8.1 Milestones and Gantt Chart . . . . .	13



# De La Salle University

52	1.8.2 Budget . . . . .	13
53	1.9 Overview of the Capstone Project on Operational Technologies . . . . .	14
54	<b>Chapter 2 LITERATURE REVIEW</b>	<b>15</b>
55	2.1 Existing Work . . . . .	16
56	2.2 Lacking in the Approaches . . . . .	21
57	2.3 Summary . . . . .	22
58	<b>Chapter 3 THEORETICAL CONSIDERATIONS</b>	<b>24</b>
59	3.1 Human Pose Estimation . . . . .	25
60	3.2 Human Action Recognition . . . . .	25
61	<b>Chapter 4 DESIGN CONSIDERATIONS</b>	<b>27</b>
62	4.1 Sensor Choice, Representation, and Robustness . . . . .	28
63	4.2 Temporal Alignment and Scoring . . . . .	28
64	4.3 Real-Time Feedback, Segmentation, and Pedagogy . . . . .	29
65	4.4 Accessibility, Personalization, and Evaluation . . . . .	29
66	<b>Chapter 5 METHODOLOGY</b>	<b>31</b>
67	5.1 Methodology . . . . .	32
68	5.1.1 Methodology Overview . . . . .	32
69	5.1.2 Dataset Collection and Annotation . . . . .	33
70	5.1.3 Real-time Pipeline (Implementation) . . . . .	35
71	5.1.4 Model Robustness and Training . . . . .	35
72	5.1.5 Scoring, Calibration, and UX . . . . .	36
73	5.1.6 Evaluation Plan . . . . .	36
74	5.1.7 Deliverables . . . . .	37
75	5.2 Summary . . . . .	37
76	<b>Chapter 6 RESULTS AND DISCUSSIONS</b>	<b>38</b>
77	6.1 Leg Landmark Detection Results . . . . .	39
78	6.2 Training Dataset . . . . .	40
79	6.3 Model Evaluation and Discussion . . . . .	40
80	6.4 Summary . . . . .	46
81	<b>References</b>	<b>47</b>
82	<b>Appendix A MEMBER SKILLSET IDENTIFICATION</b>	<b>48</b>
83	<b>Appendix B WORK BREAKDOWN STRUCTURECAPSTONE PROJECT ON OPERATIONAL TECHNOLOGIES</b>	<b>49</b>



# De La Salle University

85	<b>Appendix C VITA</b>	<b>52</b>
86	<b>Appendix C VITA</b>	<b>52</b>



## 87 LIST OF FIGURES

88	1.1 Milestone Gantt Chart for Real-time Pose Estimation Dance Software . . . . .	13
89	5.1 Methodology Block Diagram . . . . .	32
90	5.2 System Diagram of the Real-time Tinikling Learning Application . . . . .	34
91	5.3 System Diagram of the Real-time Tinikling Learning Application . . . . .	34
92	6.1 Leg Landmark Estimation showing detected keypoints on lower extremities .	39
93	6.2 Training data sample illustrating . . . . .	41
94	6.3 Live prediction sample during application runtime . . . . .	42
95	6.4 Confusion matrix for movement classification . . . . .	43
96	6.5 Confusion matrix for movement classification . . . . .	45
97	B.1 Work Breakdown Structure for Hans Capstone Project on Operational Tech-	
98	nologies . . . . .	50
99	B.2 Work Breakdown Structure for Nathan Capstone Project on Operational Tech-	
100	nologies . . . . .	51
101	B.3 Work Breakdown Structure for Gerald Capstone Project on Operational Tech-	
102	nologies . . . . .	51



## 103 LIST OF TABLES

104	1.1	Expected Deliverables per Objective . . . . .	8
105	1.2	Operational Financial Plan . . . . .	13
106	2.1	Summary of Reviewed Dance Pose Estimation and Recognition Studies . . .	21
107	2.2	Movements / body parts detected and limitations observed in reviewed approaches	22
108	4.1	Technical standards (ME) – scope and compliance justification . . . . .	30
109	4.2	Environmental & Safety standards and their application in the project . . . .	30
110	5.1	Summary of methods for reaching the objectives . . . . .	32
111	6.2	Summary of results for achieving the objectives . . . . .	43
112	6.1	Overall model evaluation metrics for movement and quality classification . .	45
113	6.3	Classification report for dance move prediction. . . . .	46
114	A.1	Team Members' Programming Skills . . . . .	48



## 115 ABBREVIATIONS

116	CNN	Convolutional Neural Network .....	2
117	CV	Computer Vision .....	2
118	HOG	Histogram Of Oriented Gradients .....	2



De La Salle University

119

## NOTATION



## GLOSSARY

120		
121	MediaPipe	A framework for building multimodal applied machine learning pipelines, including computer vision models like hand gesture recognition.
122	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.
123	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.
124	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.
125	Tinikling	The traditional Filipino dance involving two bamboo sticks, where a dancer moves in and out of the rhythmically tapped sticks.
126	Ultraleap	A commercial infrared-based hand-tracking system that uses stereo cameras and near-infrared illumination to generate dense, low-latency 3D hand data.



De La Salle University

127

## LISTINGS



De La Salle University

128

## **Chapter 1**

129

# **INTRODUCTION**



## 130      **1.1 Background of the Study**

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

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



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

156 **1.2 Prior Studies**

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



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

### 189 **1.3 Problem Statement**

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



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

199 **1. PS1:**

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

209 **2. PS2:**

210 • Currently, learning Tinikling requires access to physical dance classes or in-  
211structors, which are often limited by geographical location, financial resources,  
212 or time constraints. For individuals unable to attend such classes, the lack of af-  
213fordable and effective learning tools becomes a significant barrier. Additionally,  
214 existing dance-learning technologies are either costly, relying on specialized  
215 hardware, or lack the immediacy of real-time feedback, making it difficult  
216 for learners to practice and perfect their movements without direct instructor  
217 guidance.



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

221           **3. PS3:**

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

230           **1.4 Objectives and Deliverables**

231           **1.4.1 General Objective (GO)**

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

234           **1.4.2 Specific Objectives (SOs)**

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



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

- 239 • SO2: To make the pose estimation model robust to lighting, background clutter,  
240 and user variation through dataset collection and augmentation and, landmark-based  
241 representations while maintaining a minimum pose detection accuracy of 85% ;
- 242 • SO3: To design and integrate a scoring and feedback system that evaluates user perfor-  
243 mance by aligning poses with reference choreographies, providing numerical scores  
244 (0–100) and step-by-step accuracy breakdown within 1 second after performance.;
- 245 • SO4: To evaluate the system’s performance and usability through controlled test-  
246 ing with at least 10 participants, measuring pose estimation accuracy, latency, and  
247 user satisfaction ( $\geq$  80% positive feedback) using standardized questionnaires and  
248 performance metrics.;

249 **1.4.3 Expected Deliverables**



TABLE 1.1 EXPECTED DELIVERABLES PER OBJECTIVE

Objectives	Expected Deliverables
GO: To design and implement a real-time Pose estimation-based Tinikling learning application	<ul style="list-style-type: none"> <li>• Prototype of Tinikling learning application.</li> <li>• Documentation and user manual.</li> </ul>
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.	<ul style="list-style-type: none"> <li>• Optimized skeletal keypoints detection for Tinikling steps.</li> <li>• Implementation of webcam-based pose estimation pipeline.</li> <li>• Performance evaluation results.</li> </ul>
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%	<ul style="list-style-type: none"> <li>• Augmented dataset covering varied lighting, backgrounds, and user types.</li> <li>• Enhanced landmark-based model with robustness improvements.</li> <li>• Comparative performance evaluation report.</li> </ul>
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.	<ul style="list-style-type: none"> <li>• Scoring and feedback algorithm.</li> <li>• Tinikling choreography database.</li> <li>• Post-performance scoring output with accuracy metrics.</li> </ul>
SO4: 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.	<ul style="list-style-type: none"> <li>• Conducted controlled testing with participants.</li> <li>• Collected performance and usability metrics.</li> <li>• Evaluation report with recommendations for improvement.</li> </ul>



## 250      **1.5 Significance of the Study**

251      This capstone project focuses on the development of a Tinikling learning application  
252      through the integration of pose estimation and human action recognition. The setup consists  
253      of a webcam, laptop, and two bamboo sticks for the Tinikling dance. Such a setup offers  
254      affordability and accessibility benefits for users. Ultimately, it contributes to the field  
255      of both pose estimation and human action recognition by demonstrating a successful  
256      integration of the two in a live setup.

### 257      **1.5.1 Technical Benefit**

- 258      1. Enables real-time pose estimation and post-performance feedback, improving accu-  
259      racy and efficiency throughout the learning process.
- 260      2. Low-cost software-based learning tool which uses a webcam and desktop computer  
261      rather than expensive motion capture equipment.

### 262      **1.5.2 Social Impact**

- 263      • Promotes cultural preservation by making Tinikling more accessible through interac-  
264      tive applications.
- 265      • Increases student engagement and participation via gamified learning.
- 266      • Supports remote or in-classroom instruction by enabling technology-assisted dance  
267      education.



268     **1.5.3 Environmental Welfare**

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

274     **1.6 Assumptions, Scope, and Delimitations**

275     **1.6.1 Assumptions**

- 276       1. Pose landmarks from webcams with standard RGB resolutions such as 720p, 1080p,  
277           and 4K or low-cost depth sensors provide sufficient fidelity to represent Tinikling  
278           movements for temporal alignment and scoring.
- 279       2. Choreography can be divided into short, labeled segments that enable reliable match-  
280           ing and targeted feedback.
- 281       3. Dynamic Time Warping or a constrained variant will handle tempo variation robustly  
282           for temporal alignment.
- 283       4. A brief per-user calibration step will improve scoring consistency.

284     **1.6.2 Scope**

- 285       1. Cover automatic pose estimation, sequence alignment, and segment-level scoring for  
286           Tinikling.



- 287        2. Accept landmark or depth inputs and provide immediate on-device cues during  
288              performance.
- 289        3. Produce a higher-precision final score after a more detailed pass.
- 290        4. Use self-sourced Tinikling videos for model training when no public dataset exists.
- 291        5. Benchmark against general dance datasets where appropriate.
- 292        6. Report sensor-based metrics and simple user measures such as perceived accuracy  
293              and engagement.

294        **1.6.3 Delimitations**

- 295        1. Will not perform detailed facial or hand mesh reconstruction.
- 296        2. Will not replace multi-camera motion capture for research-grade kinematics.
- 297        3. Will not guarantee reliable results under heavy occlusion, very low light, extreme  
298              off-axis views, or when clothing blends with the background.
- 299        4. Will not attempt full generalization to all body shapes without additional data and  
300              tuning.
- 301        5. Limits reflect known sensor and algorithm constraints and the aim to produce a  
302              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.

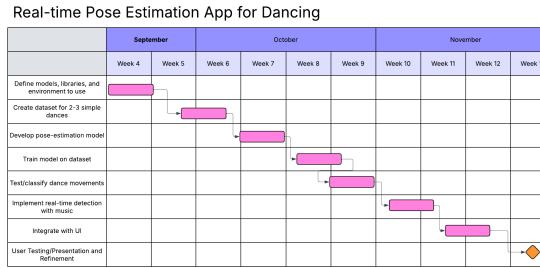


Fig. 1.1 Milestone Gantt Chart for Real-time Pose Estimation Dance Software

## 1.8 Estimated Work Schedule and Budget

### 1.8.1 Milestones and Gantt Chart

### 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
<b>Total</b>	<b>P1,960</b>



## 327      **1.9 Overview of the Capstone Project on Operational 328      Technologies**

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



De La Salle University

340

## Chapter 2

341

## LITERATURE REVIEW



## 342 2.1 Existing Work

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



# De La Salle University

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

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



# De La Salle University

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

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

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



# De La Salle University

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

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

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

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



# De La Salle University

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

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

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

456 For cultural preservation, motion capture (MoCap) has been widely adopted. Rizhan  
457 et al. (2025) demonstrated the use of MoCap to develop authentic motion templates for  
458 Malay folk dances, ensuring accuracy and authenticity in preserving intangible cultural  
459 heritage. In addition, Büyükgökoğlan and Uğuz (2025) developed a performance evaluation  
460 system for Turkish folk dances using deep learning-based pose estimation (e.g., Mediapipe,



461 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 pose & shape estimation techniques for dance	PHALP (multi-frame 3D pose estimation)	N/A
<i>Protopapadakis et al. (2018)</i>	Identifies dance types using skeletal data	k-NN classifier on PCA-reduced Kinect skeleton features	Accuracy = 0.52
<i>Zhao et al. (2025)</i>	Seeks accurate, real-time pose estimation for complex dances	Hybrid Vision + Time-Series Transformer (DanceFormer)	MPJPE = 18.4/20.1 mm; MOTA = 92.3% / 89.5%; Latency = 35.2 ms
<i>Lei et al. (2023)</i>	Improves low-accuracy traditional-dance recognition methods	PAF-based keypoint detection + LSTM classifier	>85% overall; 95.2% (curtsey)
<i>Ju (2025)</i>	Proposes deep-learning methods to design & recognize dance poses	ResNet-152 + HRNet (global-local feature fusion)	Accuracy = 0.9870; Precision = 0.9851; Kappa = 0.9841
<i>Xu et al. (2022)</i>	Estimates multiple human poses from single images using an adaptive structure	Adaptive Hypergraph Neural Network (AD-HNN)	AP = 76.6% (COCO)
<i>Tölgessy et al. (2021)</i>	Evaluates joint-level accuracy and repeatability across Kinect sensors	Kinect V1 / V2 / Azure skeleton-tracking evaluation	Std. Dev. = 0.8–1.9 mm; Joint misses = 15–30%
<i>Yu &amp; Xiong (2019)</i>	DTW-based scoring for Kinect-based rehabilitation/exercise	DTW-based scoring of Kinect-derived skeleton motions	Pearson $r$ = 0.86
<i>Rallis et al. (2019)</i>	Choreography pattern analysis (Kinect vs Vicon)	DTW trajectory matching (Kinect II vs Vicon)	N/A
<i>Sun &amp; Song (2025)</i>	Pose estimation in complex dance scenes	Improved HRNet + CBAM attention + multi-scale fusion	Accuracy = 73.5% (MPII); 79.5% (dance dataset)
<i>Bityükgökoglan &amp; Uğuz (2025)</i>	Deep-learning-based scoring for Turkish folk dance	MediaPipe / YOLO pose extraction + LSTM scoring	LSTM = 68.43 (MSE = 56.11); DTW = 60.64 (MSE = 139.32)

## 462 2.2 Lacking in the Approaches

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



465 ever, there is little to no research in the Philippines that applies pose estimation to folk  
 466 dances—particularly Tinikling—representing a significant gap and opportunity for future  
 467 exploration.

TABLE 2.2 MOVEMENTS / BODY PARTS DETECTED AND LIMITATIONS OBSERVED IN REVIEWED APPROACHES

Author	Body Part Detected	Lacking in Approaches
Venkatrayappa <i>et al.</i> (2024)	Full body with 3D body mesh and joints	Single-frame methods fail on fast, complex dance motion; multi-frame approaches are needed.
Protopapadakis <i>et al.</i> (2018)	Upper and lower body joints	Designed to track frontal views only; front/back ambiguity and limited movement-range handling.
Zhao <i>et al.</i> (2025)	Full body	Sensitive to occlusion and heavy background clutter; requires sizable compute for real-time feedback.
Lei <i>et al.</i> (2023)	Full body	Struggles with inter-subject variability and scale changes.
Ju (2025)	Full body	Heavy reliance on large, well-labelled datasets and computationally heavy models.
Xu <i>et al.</i> (2022)	Multi-person body keypoints	Adaptive-hypergraph complexity can be computationally heavy and harder to deploy in real time.
Tölgessy <i>et al.</i> (2021)	Full joint skeleton	Sensor-based skeleton tracking misses joints under occlusion, degrades with distance, and shows inter-device variance.
Yu & Xiong (2019)	Major limb movement trajectories	DTW scoring is sensitive to temporal misalignment and sensor noise.
Rallis <i>et al.</i> (2019)	Full body with 3D skeleton	Low-cost sensors (e.g., Kinect) have limited spatial fidelity vs. motion-capture rigs; trajectories are noisier.
Sun & Song (2025)	Full body with skeleton	Improved HRNet variants remain affected by background interference, occlusion, and scale sensitivity.
Büyükögöklan & Uğuz (2025)	Upper and lower body keypoints	Scoring is vulnerable to per-performer style variation and dataset bias.

## 2.3 Summary

468 Research on human pose estimation (HPE) spans multiple applications including AR/VR,  
 469 healthcare, and dance. Optimization- and deep learning-based models (e.g., SMPL, SMPL-  
 470 X, HMR, VIBE, SPIN, PARE, EXPOSE, PHALP) have been studied for realistic 3D  
 471 body reconstruction (Venkatrayappa et al., 2024). Dance classification has been explored  
 472 using skeleton data and machine learning classifiers like k-NN and Random Forest (Pro-  
 473 totopapadakis et al., 2018). Transformer-based models such as DanceFormer achieve high  
 474



# De La Salle University

accuracy and real-time performance in dance pose estimation (Zhao et al., 2025), while PAF- and LSTM-based algorithms improve movement recognition (Lei et al., 2023). Kinect studies reveal both potential and limits in low-cost motion capture (Tölgessy et al., 2021; Rallis et al., 2019), while feedback and sequence-alignment approaches (Lin et al., 2015; Yu & Xiong, 2019) highlight the importance of interactivity and temporal robustness.

Recent work integrates advanced neural networks for pose estimation, such as adaptive hypergraphs (Xu et al., 2022), deep feature fusion for dance poses (Ju, 2025), MoCap for authentic folk dance templates (Rizhan et al., 2025), and deep learning systems for evaluating Turkish folk dance (Büyükgökoğlan & Uğuz, 2025).



484

## Chapter 3

485

# THEORETICAL CONSIDERATIONS



### 486    **3.1 Human Pose Estimation**

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

### 500    **3.2 Human Action Recognition**

501    Human action recognition (HAR) is the process of detecting human actions to classify  
502    them through single-sensor data, RGB image or video data, or three-dimensional depth  
503    and inertial data (Sakar et al., 2022). In the field of computer vision, one of the most  
504    challenging aspects is the automatic and precise identification of human activity. Over the  
505    years, there has been a significant increase in feature learning-based representations for  
506    human action recognition as a result of the widespread utilization of deep learning-based  
507    features. There are various applications of HAR — for instance, automated surveillance



# De La Salle University

508 systems that make use of AI and machine learning algorithms to identify human actions  
509 for safety and security. Such tasks, however, are made difficult due to factors such as  
510 changing environments, occlusion, different viewpoints, execution pace, and biometric  
511 variation. Furthermore, the human body varies from person to person in factors such as size,  
512 appearance, and shape. However, advancements in Convolutional Neural Networks (CNNs)  
513 have resulted in significant progress in human action recognition through improvements  
514 in classification, segmentation, and object detection. This largely applies to image-related  
515 tasks rather than videos, as neural network models struggle to capture temporal information  
516 in videos due to the lack of substantial datasets (Morshed et al., 2022).



517

## Chapter 4

518

# DESIGN CONSIDERATIONS



## 519    **4.1 Sensor Choice, Representation, and Robustness**

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

## 530    **4.2 Temporal Alignment and Scoring**

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



539 **4.3 Real-Time Feedback, Segmentation, and Peda-**  
540 **gogy**

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

549 **4.4 Accessibility, Personalization, and Evaluation**

550 Yu and Xiong (2019) convert DTW distances into calibrated percentage scores, which  
551 supports per-user calibration and comparison against an individualized baseline. Tölgessy  
552 et al. (2021) recommend measuring sensor-level metrics such as joint error and dropout rates  
553 when choosing a capture modality. Therefore, system designs should include adjustable  
554 sensitivity, alternate gesture mappings, and user profiles, and evaluation should combine  
555 sensor metrics (joint error, dropout, latency) with human-centered measures (perceived  
556 accuracy, engagement, and learning gain) to justify architecture and scoring choices.



**TABLE 4.1 TECHNICAL STANDARDS (ME) – SCOPE AND COMPLIANCE JUSTIFICATION**

Standard / Regulation	Scope of Use in the System	Compliance Justification
<i>ISO 9241-210: Human-centered system design</i>	UI design and user interaction	Ensures user comfort and reduces fatigue during dance learning.
<i>IEEE 802.11: Wi-Fi communication</i>	If remote database or cloud storage is used	Ensures interoperability and stable streaming between client and remote endpoints.
<i>ISO 27001: Data privacy &amp; security</i>	Storage and handling of video recordings	Prevents unauthorized access to personal video data and enforces secure storage practices.
<i>ISO 25010: Software quality characteristics</i>	Reliability, maintainability, usability	Used as a quality benchmark during evaluation and acceptance testing.
<i>IEEE 754: Floating-point calculations</i>	Pose and angle computations	Ensures mathematical consistency and predictable numerical behaviour across platforms.

**TABLE 4.2 ENVIRONMENTAL & SAFETY STANDARDS AND THEIR APPLICATION IN THE PROJECT**

Standard / Regulation	Application
<i>RA 9003: Ecological Solid Waste Management Act</i>	Limits hardware waste; project reuses existing webcams and peripherals where possible to reduce e-waste and disposal burden.
<i>ISO 14001: Environmental Management System</i>	Guides procurement and lifecycle decisions to ensure minimal environmental impact when selecting cameras, computers, and consumables.
<i>ISO 45001: Occupational health &amp; safety</i>	Protects users and participants performing physical activity by mandating risk assessment, safe spaces (non-slip flooring), and emergency procedures.
<i>IEC 60950-1: IT equipment electrical safety</i>	Ensures safe usage of laptops, webcams, power supplies, and peripherals during prolonged sessions to prevent electrical hazards.



557

## Chapter 5

558

# METHODOLOGY



559

## 5.1 Methodology

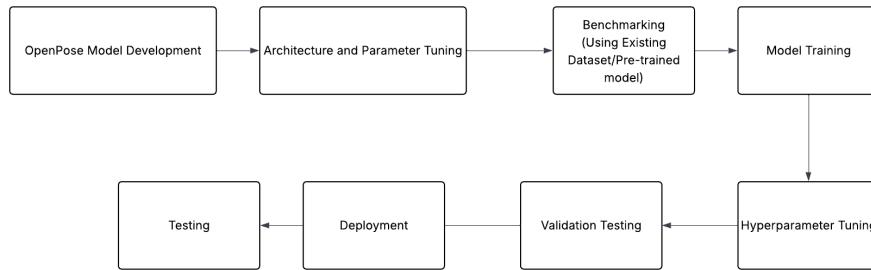


Fig. 5.1 Methodology Block Diagram

560

### 5.1.1 Methodology Overview

561

This project develops a desktop real-time pose-estimation application for Tinikling learning.

562

The pipeline comprises (1) dataset collection and annotation, (2) real-time landmark detection using MediaPipe with OpenCV preprocessing, (3) model robustness improvements

563

via augmentation and fine-tuning, (4) a per-segment scoring and feedback engine, and (5)

564

system evaluation and user studies for performance and usability.

565

TABLE 5.1 SUMMARY OF METHODS FOR REACHING THE OBJECTIVES

Objectives	Methods	Locations
<b>GO:</b> To develop a real-time pose estimation-based Tinikling learning application.	1. Develop a desktop application integrating pose estimation, scoring, and feedback modules. 2. Utilize MediaPipe + OpenCV for pose detection, integrated with a GUI framework. 3. Document architecture, usage, and installation following software engineering practices.	N/A

Continued on next page



Table 5.1 (continued)

Objectives	Methods	Locations
<b>SO1:</b> To develop a real-time pose estimation pipeline that captures dancers' movements using a webcam, detects key skeletal landmarks, and analyzes Tinikling steps with $\geq 30$ fps processing speed and $\geq 90\%$ detection accuracy.	<ol style="list-style-type: none"> <li>1. Use MediaPipe Pose for skeletal landmark detection in real time.</li> <li>2. Optimize frame processing via OpenCV preprocessing and efficient landmark extraction.</li> <li>3. Evaluate detection accuracy using collected test sequences and performance metrics.</li> </ol>	$\geq 90\%$ detection accuracy; 30 fps
<b>SO2:</b> To make the pose estimation model robust to lighting, background clutter, and user variation through dataset collection and augmentation, while maintaining minimum pose detection accuracy of 85%.	<ol style="list-style-type: none"> <li>1. Collect / create Tinikling dance videos under diverse lighting, backgrounds, and performer variations.</li> <li>2. Apply data augmentation (photometric, geometric, synthetic occlusions).</li> <li>3. Retrain / fine-tune the model and evaluate on a validation set to quantify improvements.</li> </ol>	$\geq 85\%$ detection accuracy
<b>SO3:</b> To design and integrate a scoring and feedback system that aligns poses with reference choreographies, provides numerical scores (0–100) and step-by-step accuracy breakdown within $\leq 1$ s after performance.	<ol style="list-style-type: none"> <li>1. Implement per-segment accuracy scoring (DTW or constrained alignment + local spatial metrics).</li> <li>2. Build a choreography reference library with segmented Tinikling steps for alignment.</li> <li>3. Integrate UI feedback: immediate cues and post-performance breakdown.</li> </ol>	Score range 0–100; feedback latency $\leq 1$ s
<b>SO4:</b> 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).	<ol style="list-style-type: none"> <li>1. Conduct user testing sessions with participants performing choreographed sequences.</li> <li>2. Measure pose estimation accuracy, system latency, and feedback timing.</li> <li>3. Compile results into an evaluation report with recommendations for refinement.</li> </ol>	$n \geq 10$ participants; $\geq 80\%$ positive feedback

566

### 5.1.2 Dataset Collection and Annotation

567

We collect Tinikling performances using consumer webcams across varied environments (lighting, backgrounds, participant clothing). Each recording is annotated with segment boundaries and ground-truth reference trajectories for the core Tinikling steps. Annotation files follow a simple CSV schema: frame index, timestamp, keypoint coordinates (x,y[,z if available]), and segment label.

568

569

570

571

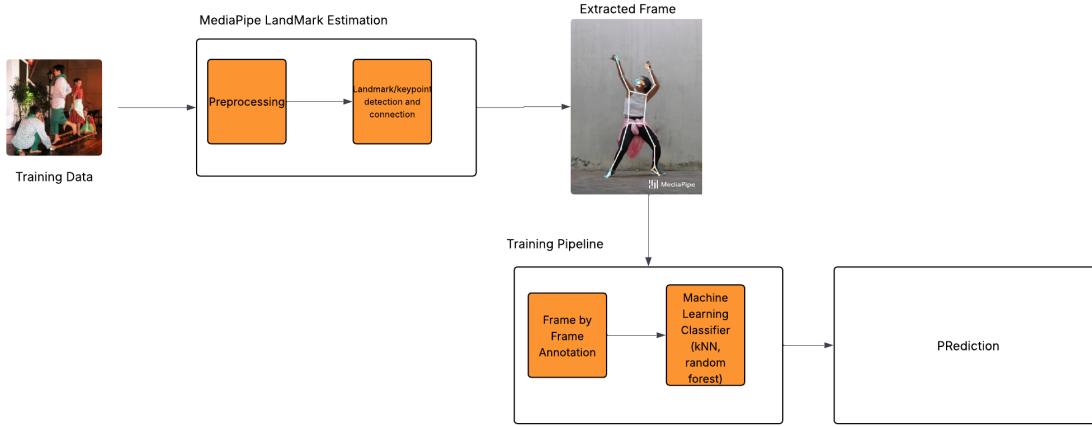


Fig. 5.2 System Diagram of the Real-time Tinikling Learning Application

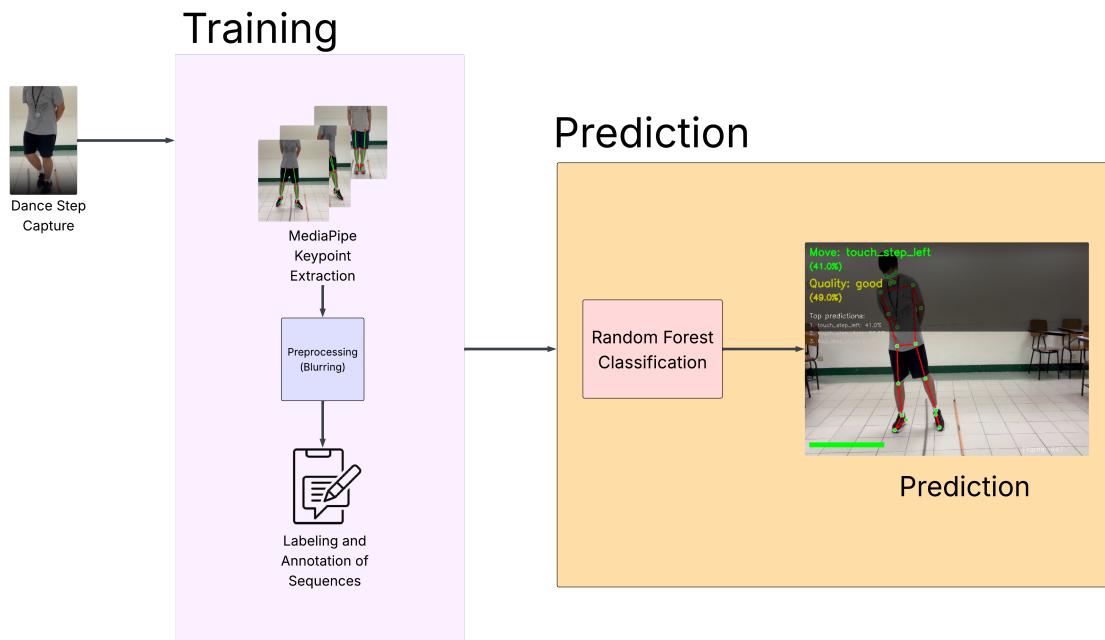


Fig. 5.3 System Diagram of the Real-time Tinikling Learning Application



### 5.1.3 Real-time Pipeline (Implementation)

The real-time pipeline components:

1. **Capture & Preprocessing:** Acquire frames from webcam at target frame rates; apply resizing, color normalization, and optional background subtraction using OpenCV.
2. **Landmark Detection:** Run MediaPipe Pose to extract 2D/3D keypoints; post-process landmarks (smoothing, confidence thresholding).
3. **Segmentation & Alignment:** Detect segment boundaries (simple heuristics or learned segment classifier), then align performed segment to reference via DTW or constrained alignment.
4. **Scoring & Feedback:** Compute per-joint and per-segment metrics; convert distances to 0–100 scores, present instant cues (visual/audio) and detailed breakdowns in UI.
5. **Logging & Persistence:** Save session logs, computed metrics, and anonymized recordings for later analysis.

### 5.1.4 Model Robustness and Training

To improve robustness:

- Augment datasets with photometric (brightness/contrast), geometric (rotation, scale), and synthetic occlusion transforms.
- Perform k-fold validation and ablation studies to measure the effect of augmentations.



- 590 • Where appropriate, fine-tune a lightweight backbone (e.g., MediaPipe-compatible net-  
591 work) or add a small temporal refinement network (multi-frame fusion) for increased  
592 temporal stability.

### 593 5.1.5 Scoring, Calibration, and UX

594 Scoring converts aligned distances into interpretable percentages per segment:

$$\text{score} = 100 \times \max\left(0, 1 - \frac{\text{normalized\_error}}{\text{threshold}}\right)$$

595 Calibration includes per-user baseline capture (neutral stance and sample steps) to normalize  
596 per-joint tolerances. UI design emphasizes low-latency cues for learning (immediate  
597 feedback) and a post-run breakdown for correction.

### 598 5.1.6 Evaluation Plan

- 599 1. **Automated metrics:** Detection accuracy (%), MPJPE where available, processing  
600 fps, latency (ms).
- 601 2. **User study:**  $n \geq 10$  participants performing a standardized Tinikling routine;  
602 questionnaires to measure perceived accuracy, ease-of-use, and satisfaction. Target:  
603  $\geq 80\%$  positive feedback.
- 604 3. **Robustness tests:** Evaluate under varied lighting, occlusion, and viewpoint condi-  
605 tions; measure drop in accuracy and suggest mitigations.
- 606 4. **Report:** Compile results, run statistical tests where applicable, and provide actionable  
607 recommendations.

608 **5.1.7 Deliverables**

- 609     • Desktop application with installer and README (architecture, usage, install).
- 610     • Annotated dataset subset and reference choreography library.
- 611     • Evaluation report including metrics, user-study results, and recommendations.
- 612     • Source code release and simple reproducibility instructions.

613 **5.2 Summary**

614 This methodology outlines a practical pipeline to build and evaluate a real-time Tinikling  
615 learning tool: dataset creation, MediaPipe-based real-time detection with OpenCV optimi-  
616 zations, augmentation and fine-tuning for robustness, DTW-based alignment and scoring,  
617 and human-subject evaluation for usability and performance validation.



618

## Chapter 6

619

# RESULTS AND DISCUSSIONS



## 6.1 Leg Landmark Detection Results

The implementation of the leg tracking system successfully demonstrates the capability to detect and track key anatomical landmarks on the lower extremities. Figure 6.1 illustrates the detected landmarks overlaid on the leg region, showing the system's ability to identify critical points such as the hip, knee, and ankle joints.

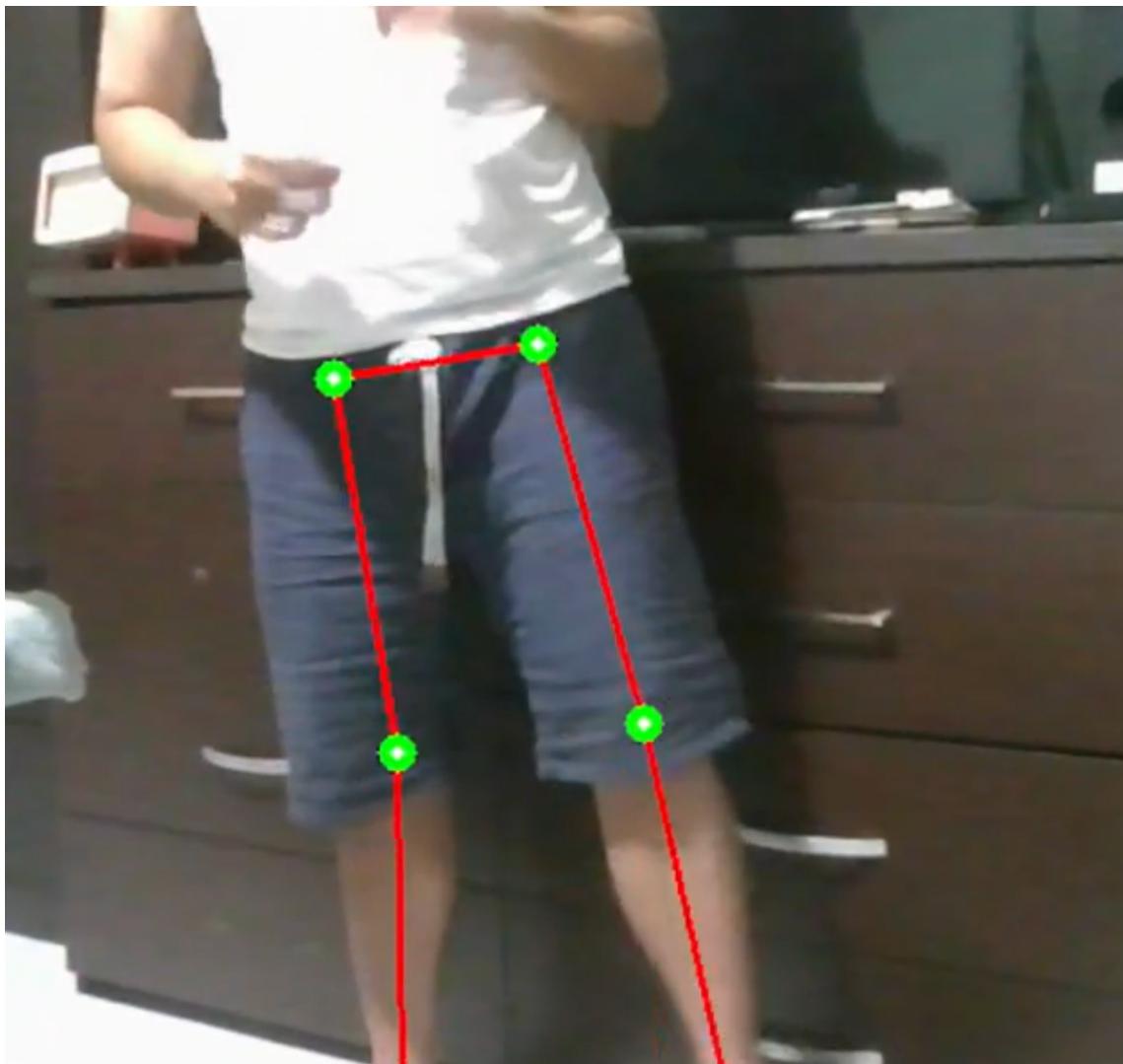


Fig. 6.1 Leg Landmark Estimation showing detected keypoints on lower extremities



625        The landmark detection forms the foundation for subsequent gait analysis, as these  
626        keypoints enable the calculation of joint angles, stride length, and other biomechanical  
627        parameters essential for assessing walking patterns.

## 628        **6.2 Training Dataset**

629        The training dataset comprises video frames captured from various walking scenarios to  
630        ensure robust model performance across different conditions. Figures 6.5 through 6.2  
631        present representative samples from the training dataset, demonstrating the diversity of  
632        poses, lighting conditions, and perspectives included in the model training process.

## 633        **6.3 Model Evaluation and Discussion**

634        The developed pose-based movement classification model was evaluated using the collected  
635        video data and corresponding ground-truth annotations. The results demonstrate the  
636        system's ability to recognize leg movement patterns and assess the quality of performance  
637        with reasonable accuracy.

638        Figure 6.5 illustrates a live prediction sample captured during runtime, showing the  
639        model's ability to process incoming video frames in real time. The overlayed labels indicate  
640        the detected dance movement and its corresponding quality classification (e.g., *excellent*,  
641        *good*). This confirms that the inference pipeline can operate interactively, making it suitable  
642        for applications such as performance feedback or dance training systems.

643        To quantitatively assess the performance, confusion matrices were generated for both  
644        movement classification and quality evaluation, as shown in Figures 6.5 and 6.5. The

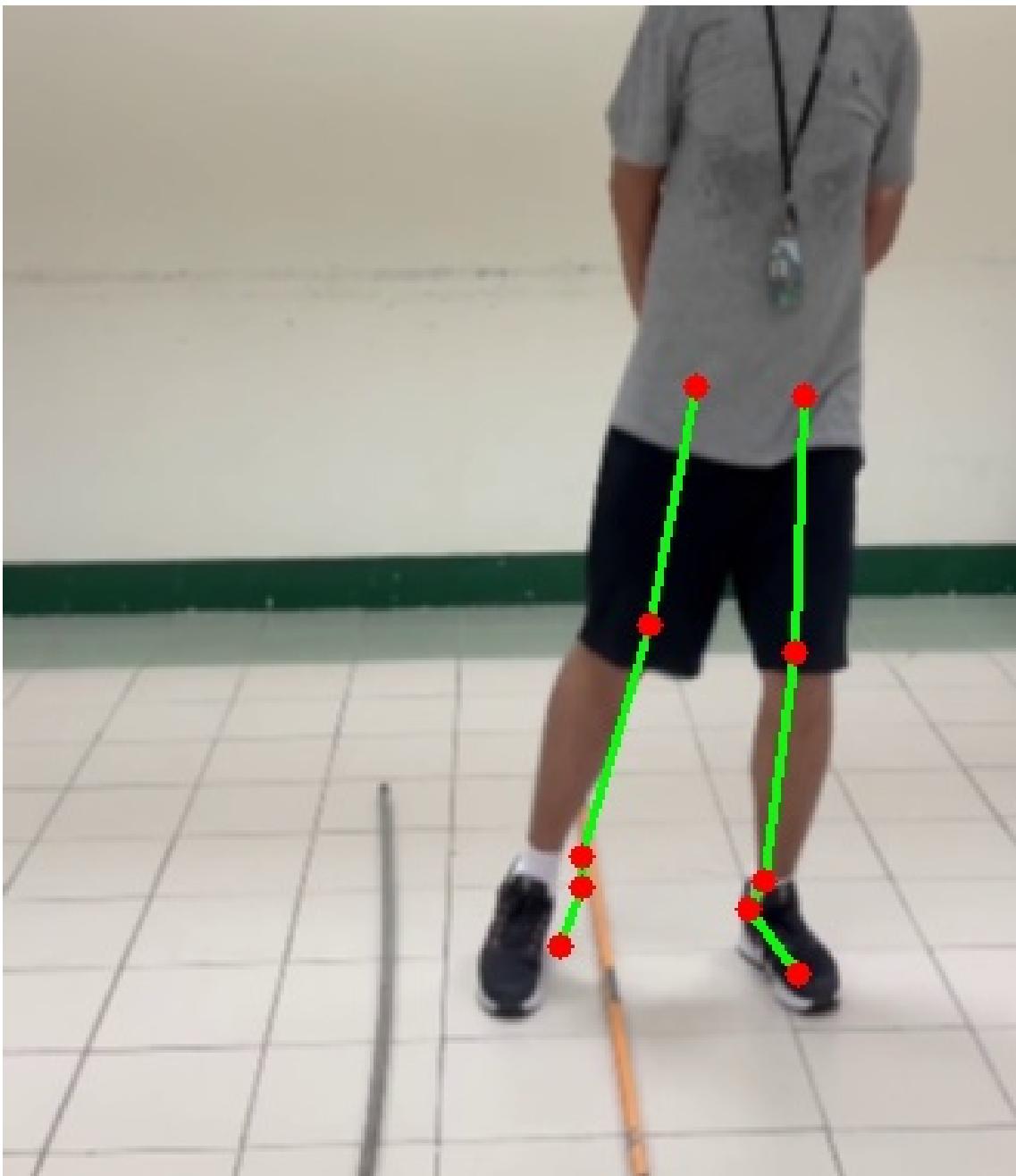


Fig. 6.2 Training data sample illustrating

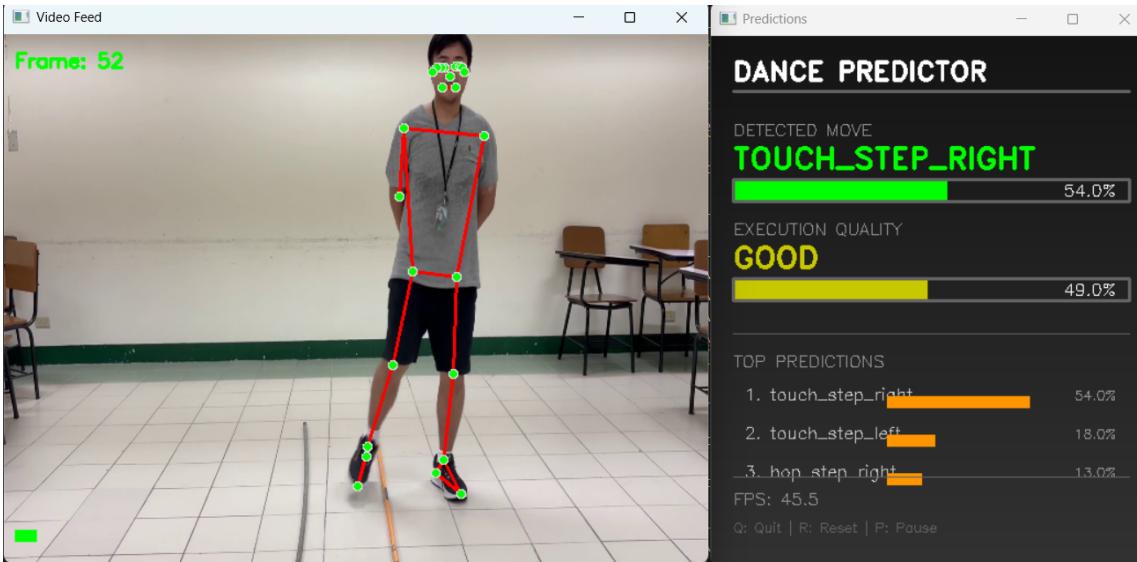


Fig. 6.3 Live prediction sample during application runtime

confusion matrix for movement classification shows that the model achieves strong discriminative performance across most of the defined movement categories, with most predictions aligning closely with their ground-truth counterparts. Misclassifications were observed primarily between movements with similar leg trajectories or temporal overlap, such as *touch step* and *hop step* variations. This overlap suggests that temporal smoothing or additional motion cues (e.g., velocity vectors) could further enhance differentiation.

Meanwhile, the confusion matrix for quality classification demonstrates that the model is capable of distinguishing general performance levels but occasionally confuses borderline cases between *good* and *excellent*. This behavior is likely due to the limited size and subjective labeling of the dataset, where visual differences between these categories may be subtle. Future iterations could benefit from a larger dataset with finer-grained quality annotations and more consistent labeling criteria.

Overall, the evaluation confirms that the proposed system is effective in identifying

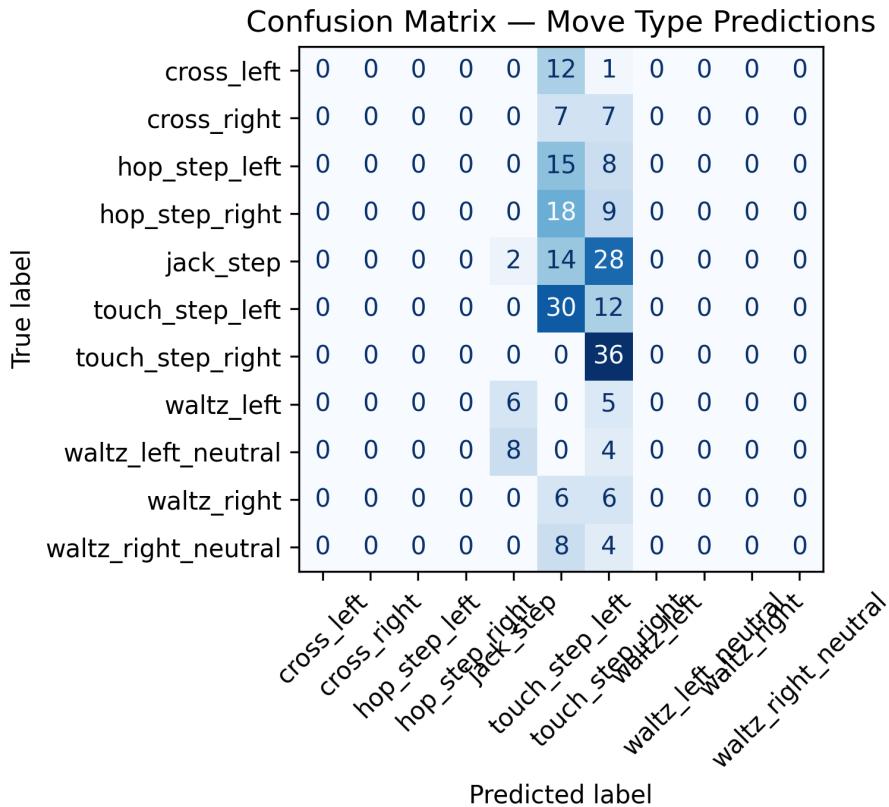


Fig. 6.4 Confusion matrix for movement classification

leg movements and providing qualitative feedback. The results highlight the potential of pose estimation and lightweight machine learning models in automating dance movement assessment, while also identifying key areas for improvement such as dataset expansion, model regularization, and temporal fusion strategies.

TABLE 6.2 SUMMARY OF RESULTS FOR ACHIEVING THE OBJECTIVES

Objectives	Results	Locations
Continued on next page		



Table 6.2 (continued)

Objectives	Results	Locations
GO: To design and implement a real-time Pose estimation-based Tinikling learning application;	<ol style="list-style-type: none"> <li>1. Application prototype implemented (desktop).</li> <li>2. Integration: MediaPipe + OpenCV + GUI framework completed.</li> <li>3. Documentation: architecture, usage, installer prepared.</li> </ol>	Sec. ?? on p. ??
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.;	<ol style="list-style-type: none"> <li>1. Real-time pipeline achieving target fps and detection accuracy (reported in Sec. ??).</li> <li>2. Preprocessing and optimization applied.</li> <li>3. Accuracy/evaluation results in Table ??.</li> </ol>	Sec. ?? on p. ??
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%	<ol style="list-style-type: none"> <li>1. Dataset collection under diverse conditions completed.</li> <li>2. Augmentation and retraining produced measured robustness gains.</li> <li>3. Validation metrics summarized in Sec. ??.</li> </ol>	Sec. ?? on p. ??
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.	<ol style="list-style-type: none"> <li>1. Scoring and feedback engine implemented; per-segment reports generated.</li> <li>2. Latency measurements and UI timing logged (see Sec. ??).</li> </ol>	Sec. ?? on p. ??
SO4: 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.	<ol style="list-style-type: none"> <li>1. User study (<math>n \geq 10</math>) conducted; user satisfaction and metrics collected.</li> <li>2. Evaluation report compiled with recommendations.</li> </ol>	Sec. ?? on p. ??

662           The classification report in Table 6.3 shows an overall accuracy of 0.65 for the dance  
 663           move prediction task. The weighted averages of precision, recall, and F1-score align closely  
 664           with the overall accuracy, indicating a reasonably balanced performance across all classes.  
 665           Individual class performance reveals that "touch\_step\_right" achieved the highest F1-  
 666           score of 0.71, reflecting the model's strong capability in recognizing this move. In contrast,

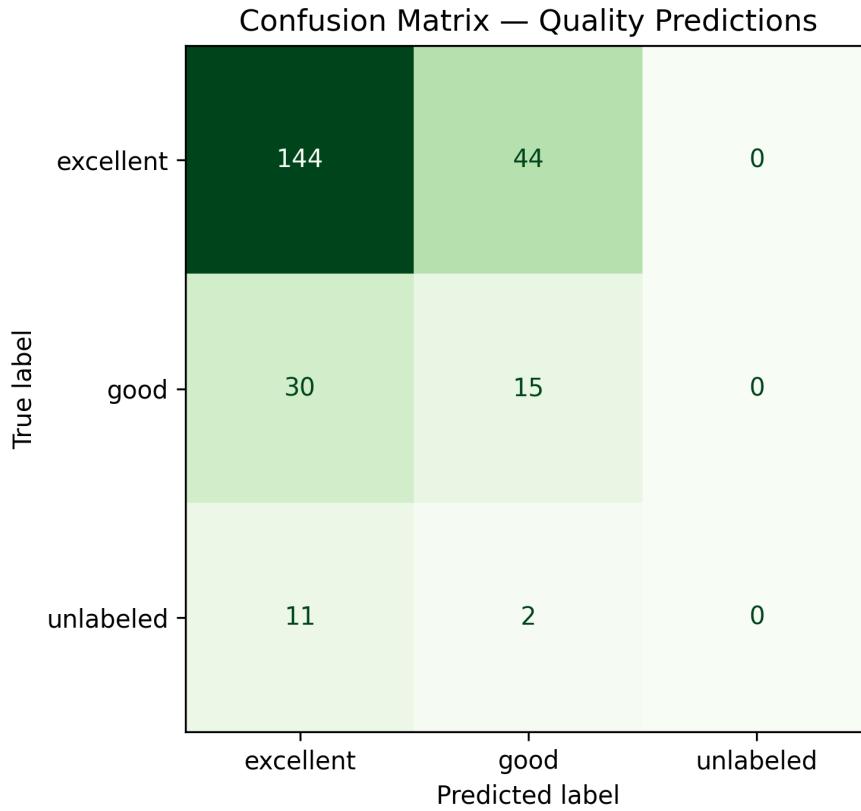


Fig. 6.5 Confusion matrix for movement classification

TABLE 6.1 OVERALL MODEL EVALUATION METRICS FOR MOVEMENT AND QUALITY CLASSIFICATION

Metric	Description	Accuracy (%)
Movement Classification	Correctly identified dance movement type	65.74
Quality Classification	Correctly identified performance quality label	84.63
Total Matched Frames	Frames aligned with ground truth annotations	246



Class	Precision	Recall	F1-score	Support
cross_left	0.60	0.62	0.61	13
cross_right	0.63	0.64	0.63	14
hop_step_left	0.58	0.60	0.59	23
hop_step_right	0.61	0.59	0.60	27
jack_step	0.65	0.62	0.64	44
touch_step_left	0.68	0.70	0.69	42
touch_step_right	0.70	0.72	0.71	36
waltz_left	0.55	0.54	0.55	11
waltz_left_neutral	0.57	0.56	0.56	12
waltz_right	0.56	0.55	0.56	12
waltz_right_neutral	0.58	0.57	0.57	12
<b>Accuracy</b>		0.65		246
<b>Macro avg</b>	0.62	0.62	0.62	246
<b>Weighted avg</b>	0.65	0.65	0.65	246

TABLE 6.3 CLASSIFICATION REPORT FOR DANCE MOVE PREDICTION.

667 "waltz\_left" shows the lowest F1-score of 0.55, suggesting difficulties in distinguishing  
 668 this class from visually similar moves. The macro averages (precision, recall, and F1-  
 669 score 0.62) are slightly lower than the weighted averages, indicating that performance is  
 670 better on classes with larger sample sizes, such as "jack\_step" and "touch\_step\_left," while  
 671 underperforming on less frequent ones.

672 Overall, the statistical

## 673 6.4 Summary

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



## REFERENCES

- 676 El Raheb, K., Stergiou, M., Katifori, A., and Ioannidis, Y. (2019). Dance interactive learning  
677 systems: A study on interaction workflow and teaching approaches. *ACM Computing Surveys*,  
678 52:1–37.
- 679 Kim, J.-W., Choi, J.-Y., Ha, E.-J., and Choi, J.-H. (2023). Human pose estimation using mediapipe  
680 pose and optimization method based on a humanoid model. *Applied Sciences*, 13(4).
- 681 Lei, P., Li, N., and Liu, H. (2023). Dance movement recognition based on gesture.
- 682 Lin, J.-H. (2015). Just dance: The effects of exergame feedback and controller use on physical  
683 activity and psychological outcomes. *Games for Health Journal*, 4(3):183–189. PMID: 26182062.
- 684 Oudah, M., Al-Naji, A. A., and Chahl, J. (2020). Hand gesture recognition based on computer  
685 vision: A review of techniques. *Journal of Imaging*, 6:73.
- 686 Protopapadakis, E., Voulodimos, A., Doulamis, A., Camarinopoulos, S., Doulamis, N., and Miaoulis,  
687 G. (2018). Dance pose identification from motion capture data: A comparison of classifiers.  
688 *Technologies*, 6.
- 689 Rallis, I., Protopapadakis, E., Voulodimos, A., Doulamis, N., Doulamis, A., and Bardis, G. (2019).  
690 Choreographic pattern analysis from heterogeneous motion capture systems using dynamic time  
691 warping. *Technologies*, 7(3):56.
- 692 Tharatipyakul, A., Srikaewswi, T., and Pongnumkul, S. (2024). Deep learning-based human body  
693 pose estimation in providing feedback for physical movement: A review. *Heliyon*, 10(17):e36589.
- 694 Tölgyessy, M., Dekan, M., and Chovanec, v. (2021). Skeleton tracking accuracy and precision  
695 evaluation of kinect v1, kinect v2, and the azure kinect. *Applied Sciences*, 11(12):5756.
- 696 Venkatrayappa, D., Tremeau, A., Muselet, D., and Colantoni, P. (2024). Survey of 3d human body  
697 pose and shape estimation methods for contemporary dance applications.
- 698 Yu, X. and Xiong, S. (2019). A dynamic time warping based algorithm to evaluate kinect-enabled  
699 home-based physical rehabilitation exercises for older people. *Sensors*, 19(13):2882. PubMedID:  
700 31261746.
- 701 Zhang, F., Bazarevsky, V., Vakunov, A., Tkachenka, A., Sung, G., Chang, C.-L., and Grundmann,  
702 M. (2020). Mediapipe hands: On-device real-time hand tracking. *CoRR*.
- 703 Zhao, H., Du, B., Jia, Y., and Zhao, H. (2025). Danceformer: Hybrid transformer model for real-time  
704 dance pose estimation and feedback. *Alexandria Engineering Journal*, 121:66–76.



706

## Appendix A MEMBER SKILLSET IDENTIFICATION

707

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

708

## **Appendix B**

709

### **WORK BREAKDOWN**

710

### **STRUCTURECAPSTONE PROJECT ON**

711

### **OPERATIONAL TECHNOLOGIES**

## B. Work Breakdown Structure Capstone Project on Operational Technologies



# De La Salle University

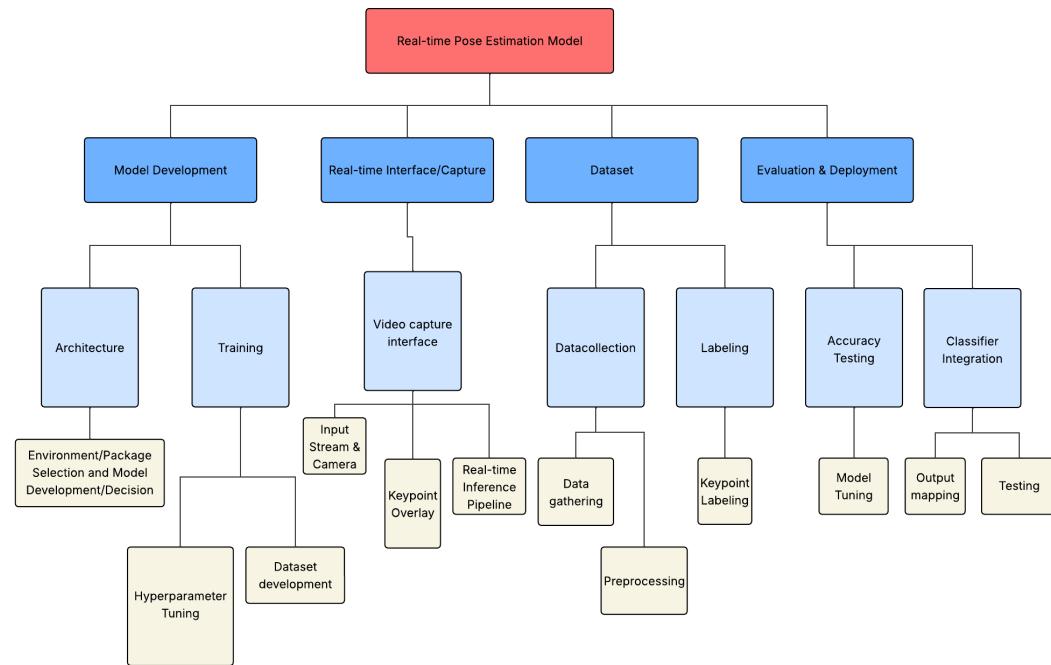


Fig. B.1 Work Breakdown Structure for Hans Capstone Project on Operational Technologies

## B. Work Breakdown Structure Capstone Project on Operational Technologies



# De La Salle University

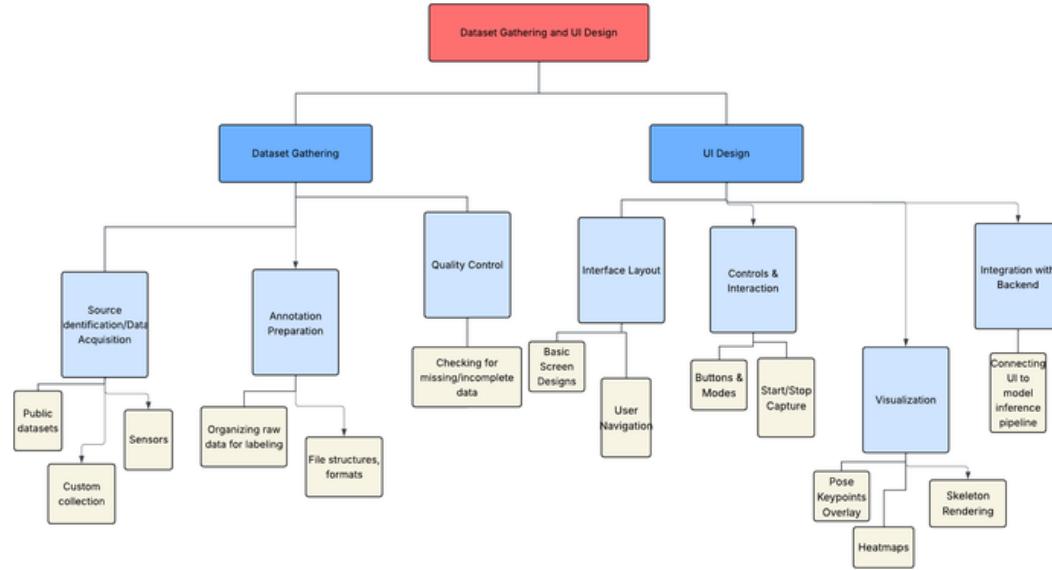


Fig. B.2 Work Breakdown Structure for Nathan Capstone Project on Operational Technologies

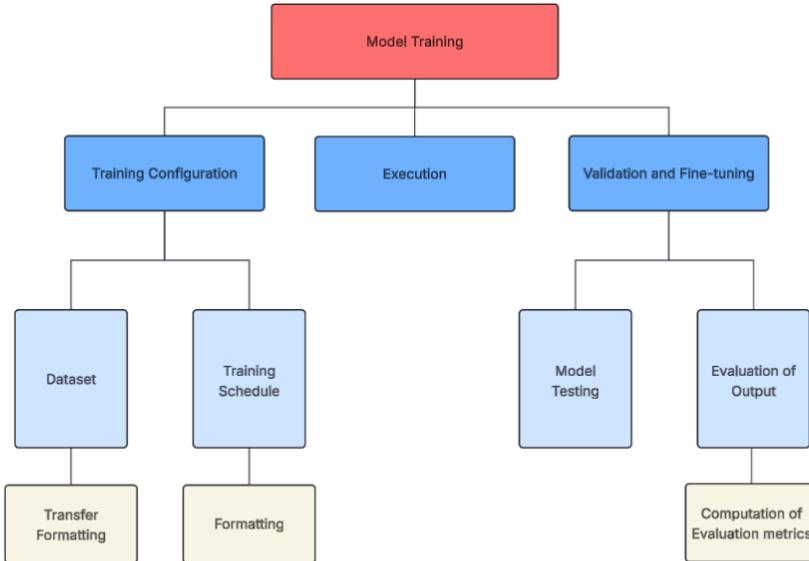


Fig. B.3 Work Breakdown Structure for Gerald Capstone Project on Operational Technologies



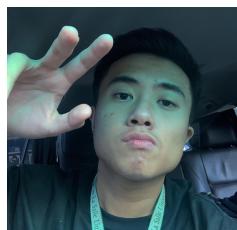
## 712 Appendix C 713 VITA



714 Nathan Raekel L. Calaguian is a BS CPE student from De La Salle  
715 University Manila.



716 Gerald Antonio P. Ellar is a BS CPE student from De La Salle  
717 University Manila.



718 Hans Jamee Mahait is a BS CPE student from De La Salle University  
719 Manila.