

DanceVision Proposal

Anna Guo Jenna Jacob Katie Wang Natalie Leal Blanco
UT Austin: Computer Vision Course
annayuab@icloud.com

Abstract

Dance is a universal form of expression and entertainment, yet mastering its intricacies can be daunting for beginners. To democratize dance learning and make it more accessible to enthusiasts, we propose a novel approach utilizing computer vision techniques. By breaking down dance videos into vectorized and text-based representations with counts, the complexities of dance routines are simplified, enabling easier comprehension and replication by average fans.

1. Introduction

In recent years, dance has surged in popularity, whether through TikTok trends in younger generations, the growing fanbase for elaborate K-pop dance videos, or the performances in reality competition shows (which themselves cover wide genres from ballet and contemporary to ballroom dances). As a result of this growing interest, many eager dance novices have attempted to self-teach from the abundant dance videos online, available through a variety of platforms. However, such individuals often struggle with breaking down movements due to the speed or intricacy with which they were executed.

We aim to use Computer Vision to partition dance movements into more easily digestible chunks by creating a model that can learn (vectorize) dance choreography from a given dance video, and output it in a more human-readable form, such as by describing multiple steps in a motion and adding counts or timings for context. By developing a tool that provides detailed descriptions of human positions and gestures alongside counts or musical beats, learning dance choreographies could become much more accessible and approachable for a broader audience.

1.1. Existing Attempts

Several existing attempts have been made to tackle the challenge of learning dance from online videos, each with its own approach and limitations. One such approach involves platforms that generate reskinned versions of dance

performances, offering visual references for users. These platforms overlay graphics or animations onto the original video, highlighting key movements or positions to provide guidance. While these reskinned versions serve as visual aids, they typically do not offer active instruction or detailed breakdowns of the dance steps. As a result, users may struggle to fully grasp the nuances and intricacies of the movements, limiting the effectiveness of this approach as a learning tool.

Another example is a DancePal Hackathon project, which adopts a more technical approach by comparing an ideal video with the end user's video skeletons[5]. This method involves analyzing the skeletal poses extracted from both videos and identifying discrepancies between them. While this approach has the potential to provide personalized feedback to users based on their own performances, it comes with several challenges. Firstly, it requires an ideal video with specific lighting, background, and perspective conditions to serve as a reference, which may not always be feasible to obtain. Additionally, accurately mapping movements onto predefined definitions can be difficult, especially for complex dance sequences with varying styles and interpretations. As a result, the effectiveness of this approach may be limited by the availability of suitable reference videos and the accuracy of the mapping process.

Unlike reskinned versions of dance performances, our proposed method provides detailed instructional breakdowns of dance movements, guiding learners through each step with clarity and precision. This active instruction facilitates better understanding and retention of choreography. Additionally, our method could be expanded to offer personalized feedback based on the analysis of the learner's own video and that of an original dance video, ideally without the requirement for perfect lighting or camera angle conditions.

1.2. Useful Techniques

Several computer vision techniques offer promising avenues for addressing the challenge of learning dance from online videos by providing deeper insights into the dynamics and nuances of dance movements. One such tech-

nique involves exploring 3D models instead of traditional 2D skeletonization methods to create more comprehensive representations of dance movements. Unlike 2D skeletonization, which captures only the spatial coordinates of key body joints in a single plane, 3D modeling techniques can capture the full spatial extent of movements, including depth information[1]. This allows for a more accurate and detailed representation of the dancer's body and movements in three-dimensional space, which is particularly valuable for capturing complex and multi-dimensional dance choreographies[12].

In addition to 3D modeling, other computer vision techniques such as action recognition, pose estimation, trajectory analysis, and video segmentation and tracking can provide valuable insights into the dynamics of dance movements. Action recognition algorithms can automatically identify and classify different dance gestures and sequences, allowing for automated analysis of choreographic patterns and styles. Pose estimation techniques can accurately estimate the skeletal poses of dancers from video frames, providing a structural representation of their movements that can be used for further analysis and interpretation.

Trajectory analysis techniques can track the trajectories of specific body parts or movements over time, allowing for the analysis of movement patterns, velocities, and accelerations[2]. This can provide valuable insights into the rhythmic and dynamic aspects of dance performances, helping learners to better understand the timing and pacing of movements. Video segmentation and tracking methods can segment dance videos into individual sequences or movements, allowing for more focused analysis and comparison of specific choreographic elements.

By leveraging these advanced computer vision techniques, learners can gain a deeper understanding of the dynamics and nuances of dance movements, facilitating more effective learning experiences. These techniques can help learners to break down complex choreographies into smaller, more manageable components, identify areas for improvement, and refine their technique with greater precision and accuracy. Overall, the integration of computer vision into dance learning has the potential to revolutionize the way dance is taught and learned, making it more accessible, engaging, and effective for dancers of all levels.

1.3. Challenges

Several challenges arise from attempting this project. One is the quality of input data, which encompasses various factors such as lighting, angles, occlusions, frame rates, and blurring. These factors can significantly impact the accuracy of movement analysis, as they can obscure or distort key visual information needed for understanding dance movements. Addressing these challenges may require the

development of robust computer vision algorithms capable of handling noisy or low-quality video data, as well as techniques for preprocessing and enhancing the visual information to improve analysis accuracy.

Moreover, handling interactions between multiple dancers in an ideal video presents additional complexities in accurately capturing and representing dance movements. In many dance performances, especially group or partner dances, dancers may interact with each other in complex ways, such as lifts, partner work, or coordinated movements. Capturing these interactions and accurately representing them in the analysis poses challenges for computer vision algorithms, as they must be able to differentiate between individual dancers and track their movements independently while also accounting for their interactions with others. Addressing this challenge may require the development of advanced tracking and segmentation techniques capable of accurately identifying and tracking multiple dancers in a crowded or dynamic environment.

Overall, addressing these challenges requires a multidisciplinary approach that combines expertise in computer vision, dance analysis, and human-computer interaction. By developing robust algorithms, incorporating domain-specific knowledge, and leveraging advanced techniques for handling complex visual data, it is possible to create a tool that effectively supports learning dance from online videos using computer vision. However, continued research and development efforts are needed to overcome these challenges and realize the full potential of computer vision technology in the field of dance education.

1.4. Technical Gap

The primary aim of this project is to reduce the technical gap between online dance learning resources and effective dance instruction. By enabling the correspondence between individuals with different body proportions captured from various angles, this tool seeks to enhance the accessibility and inclusivity of dance education. This interdisciplinary endeavor involves collaboration between experts in dance, biomechanics, and computer vision, with computer vision providing quantitative analysis capabilities to enhance the understanding and teaching of dance movements.

In conclusion, developing a tool for learning dance from online videos using computer vision holds great potential for revolutionizing dance education and making it more accessible to a wider audience. By leveraging advanced computer vision techniques and interdisciplinary collaboration, this project aims to bridge the gap between online dance resources and effective dance instruction, ultimately empowering individuals to learn and appreciate dance in new and innovative ways.

2. Pose Estimation Models

The first step in developing this project was to define the parameters and requirements in order to choose models and algorithms that satisfy our expectations. The first model that will be implemented on our project is a Pose Estimation Model. The requirements we defined were: limit GPU usage, multiple persons detection, and output in the form of JSON files, in order to be transferable and readable for the subsequent sections of our project. We also set a high priority on accuracy, due to the requirements for the other portions of our project[10].

2.1. Model Choices

Options that we considered for use in this project include OpenPose, MediaPipe, and MoveNet. All of which will be detailed below.

The first model we tested was MediaPipe, which is a Framework with many applications that has been widely used by Google for numerous products and services[11]. Although, at first glance MediaPipe provided satisfactory results, after further inspection and testing, the model did not satisfy our requirements. The first limitation is that MediaPipe only performs single-person analysis. Other limitations we encountered after testing were lower accuracy than expected, as well as difficulty converting the output into JSON files. Thus, we continued our research for other models.

Our second option was MoveNet, which is a pose estimation model developed by Google, which uses deep learning techniques. It is an improved version of the previous generation PoseNet model, released in 2017[9]. Though this worked quickly, it did not have the desired accuracy for the video inputs. Similar to MediaPipe, there was significant difficulty converting the output to JSON files. Although this proved more successful than MediaPipe, it still was not satisfying our requirements.

Finally, we tested OpenPose, a software for pose estimation, which has proved to work well when evaluating multiple persons, even for more crowded scenes. It uses a bottom-up approach to produce higher quality results, though at a higher computation cost compared to MediaPipe and MoveNet[6]. Additionally, it outputs JSON files, one per video frame, containing data for keypoints (e.g., left knee, right knee, etc.). Because of its high accuracy and usable output, OpenPose was selected for this project over MediaPipe and MoveNet.

2.2. Challenges with OpenPose

Since OpenPose was a rather new concept to our team, we initially had some trouble setting it up and running on our own devices. At first, it appeared to require more significant computational resources that were readily available. However, after extraneous research we were able to set it up and

get it running properly. Due to our lack of experience with OpenPose, there was a learning curve in figuring out how exactly the model operates and what we could do with it, but over time, we were able to learn how to use it well.

As we became more familiar and comfortable with the OpenPose model, we were able to test the effects of different types of videos as inputs. We learned that lowering the resolution of a video and then inputting it to the program resulted in a significantly faster processing time. However, this also consistently resulted in decreased accuracy in the pose estimation. Thus, after several rounds of testing, we ultimately decided against reducing the resolution in our input videos since accuracy was crucial to this project.

3. Algorithm to Prepare OpenPose outputs for Feature Extraction

For a given video, OpenPose outputs a directory of json files, with one for each frame. Before passing these outputs to a second machine learning model, we want to write an algorithm that maps similar movements to the same values, regardless of camera angle or proximity of the person to the camera. (Due to the nature of OpenPose outputting x, y values for keypoints, we have already circumvented the issue of changes in illumination.) Our algorithm aims to:

1. Remove absolute position by calculating the change in positions of key points between consecutive frames.
2. Normalize position gradient using the maximum distance (in the 32 most recent frames) from right shoulder to right hip.
3. Accommodate consistent elliptical motion by calculating angles at body joints.

We choose to transition away from absolute positions to position gradients due to inspiration from SIFTs usage of illumination gradients. Absolute positions could cause our model to unintentionally train on a dancers location in the camera view rather than their movements, whereas position gradient should make our model invariant to the dancers x, y-positions.

Similarly, by attempting to normalize the position gradients to the length of the dancers body part, we hope to make our model invariant to the dancers height and their distance from the camera (people in the foreground tend to appear taller than people in the background). We choose the right shoulder and right hip specifically because we believe the distance in between to represent a relatively rigid part of the body with a length correlated to a persons height. Furthermore, we choose the longest length found to account for dance movements that may cause the shoulder and hip to be closer together than usual, while limiting the range to 2 seconds (32 frames) so that were still sensitive to global changes (changes throughout the video) in z position (the distance between the dancer and camera center).

Finally, we want to pass in angles calculated from the keypoints. Firstly, this will likely aid the model in considering elliptical motions, which are prevalent in dances, as a continuous motion, as opposed to a single elliptical motion being calculated as multiple straight-line movements. Additionally, this encodes information about body parts relative to each other. For instance, moving ones right hand downwards on the left side of ones body could be differentiated from moving ones right hand downwards on the right side of ones body through the angle of the right elbow relative to the right shoulder.

4. Applications: Physical Therapy

Although this project originally had a focus on dance, we believe that it can be applied more widely in medical care, specifically in physical therapy. Since our project is able to read video inputs and detect the poses in videos, it can do the same for physical therapy. This would greatly serve patients in their PT journeys, whether healing from a sports injury, recovering from a surgical procedure, or relieving pain from a specific condition. Patients undergoing this type of treatment will typically receive a plan outlining a set of exercises and stretches to be completed and aid in their process to restoration. With the scope of our project, it can similarly analyze the data from these exercises and stretches being done, then evaluate the patients progress over time. This could help motivate the patient in their recovery and provide further insight for both the patient and their physical therapist as to how their recovery journey is progressing, also allowing for a better understanding of how ones plan might be adjusted accordingly. We aim for our project to have a level of flexibility that will allow for this to be easily adapted[4].

References

- [1] A.Mumuni and F. Mumuni. Data augmentation: A comprehensive survey of modern approaches. *Elsevier*, 16(1), 2022.
- [2] B.B. Armor, J. Su, and A. Srivastava. Action recognition using rate-invariant analysis of skeletal shape trajectories. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(1):1–13, 2016.
- [3] F. Baradel. *Structured deep learning for video analysis*. PhD thesis, 2004.
- [4] S. Bernstein. What is physical therapy? *WebMD*, 2023.
- [5] V. Buwaneka, A. Samadder, and K. Sharma. Dancevideo, 2023.
- [6] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. Openpose: Realtime multi-person 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [7] C. Chan, S. Ginosar, T. Zhou, and A. Efros. Everybody dance now, 2019.
- [8] Y. Jafarian and H. Park. Learning high fidelity depths of dressed humans by watching social media dance videos, 2021.
- [9] K. LeViet. Pose estimation and classification on edge devices with movenet and tensorflow lite, 2021.
- [10] Maureentkt. Selecting your real-time pose estimation models. *Medium*, 2021.
- [11] MediaPipe. Pose landmark detection guide, 2023.
- [12] OpenCV. Open source computer vision library, 2024.
- [13] Y. Pang and Y. Niu. Dance video motion recognition based on computer vision and image processing. *Applied Artificial Intelligence*, 37(1), 2023.
- [14] M. Yang and Z. He. Dance action recognition model using deep learning network in streaming media environment. *PMC*, 2022.