

Test for Prospective HumanAI GSoC 2025 Applicants:

AI-Enabled Choreography: Dance Beyond Music

Part 1: Why this project?

Out of the many other GSoC organizations and projects, this project caught my attention and I was instantly interested in being a part of it. As a 19 year old girl, I was quite enthusiastic in dance, music and arts and crafts. As I'm growing older, I seem to have lost the touch of the artforms that there are due to the increased academic workload. Even after this, after entering college, I did try to pursue my hobbies by being a part of the art, music and dance club and participating in their activities. I have been very active in graphic designing since the past year and this project seems like it would bring me closer to the things I love doing while also developing skills and gaining some real world experience.

As far as I have thought on this, potential benefits in incorporating AI into the creative process might be the fact that: AI can automate repetitive tasks like color correction, formatting, and rendering, allowing creators to focus on higher-level creative work; AI tools can suggest themes, generate images, or assist with brainstorming by providing new perspectives; AI can reduce the need for large teams, cutting down expenses in creative industries like design, marketing, and filmmaking; etc. Whereas there are also some drawbacks to this: AI-generated content may lack emotional depth, originality, and unique artistic expression(loss of human touch); Automation in creative industries could reduce job opportunities for artists, writers, and designers; AI tools often rely on user data, raising concerns about data security and ownership; excessive use of AI has also raised some environmental issues that the GPUs used for these produce excessive heat and also too much of water is used to cool it down. These issues can be seen from the recent "ChatGPT trend" of turning normal images into Studio Ghibli art.

Since I'm a student, and during the project timeline, I would have ongoing summer vacations, I plan on collaborating with the mentors to work on this project and get their guidance on the next steps. If I do get selected for this project, I plan on learning from the mentors, about the workflow on such big projects, how to research properly, how to implement the knowledge I have, professional work ethics and many more. I'm looking forward to being assigned this project and getting started ASAP.

Furthermore, there are some ideas that I came up with, on the subject of multimodal representations of dance: extract dance movements and synthesize new sequences and choreography(using GANs, Transformers); studying complex dance movements and learning the physics behind the movements to replicate that(there was a movie named "Ice Princess" where a physics student used her knowledge to study the movement of the figure skaters and their dance movements and came up with the solutions to correct their mistakes.); using dance as therapy and making curated dance choreos for specific ailments according to the patients(using IMU sensors in smartwatches, motion trackers to analyze dance movements for physical therapy); a system where

users input descriptions (e.g., "graceful ballet spin"), and an AI generates corresponding dance movements (LLMs paired with motion capture datasets to convert text into 3D animated dance sequences); suggesting songs based on the dance movements.

These are some of the additional project ideas that can be further developed apart from the initial aim of this project.

Part 2: Coding

Before starting on the 2 tasks, we load the data as follows:

Code: [load_data](#)

→ **Animate some dance data**

Code: [code1](#)

First Task: Visualizing a 3D Dance Sequence

Alright, so the goal was to take some motion capture data (which is in .npy format: the Mariel dataset) and create a 3D animated plot that shows the dancer moving.

Here's how I worked through the given problem:

1. Understanding the Data

- First, I loaded the data file and checked its shape. It turned out to be (T, J, 3), where:
 - T = Number of timesteps (frames in the animation).
 - J = Number of joints (points on the dancer's body).
 - 3 = X, Y, Z coordinates for each joint.

2. Plotting a Single Frame

- Before jumping into animations, I wanted to make sure I could plot one frame of the dance.
- I used Matplotlib's 3D scatter plot to visualize the joints at $t=0$.
- After running the code, the result wasn't as it was supposed to be; rather it was a plotting and movement of random points.

3. Connecting the Joints

- To make it actually look like a human, I needed to connect certain joints with lines.
- I looked at common motion capture skeleton structures and defined some connections between key joints (e.g., head to spine, arms to shoulders, legs to hips, etc.).
- Once I plotted these connections, I was expecting to see the shape of a human- dancing. But unfortunately, the result was not achieved due to some issues.

4. Animating the Dance

- Now, I needed to loop through all frames and update the positions of the joints in real time.
- I used Matplotlib's animation module to update the plot at each timestep.
- The tricky part was making sure the axes stayed fixed so the dancer didn't keep shifting around weirdly.

5. Issues

While the motion capture data did animate over time, several challenges were observed:

- **Lack of Joint Mapping Information:**

- The dataset did not provide explicit joint labels, making it difficult to correctly connect points into a human-like skeleton.
- This could be fixed by obtaining a joint index mapping from the dataset creators or aligning it with a standard skeleton format (like SMPL or CMU MoCap).

- **Possible Incorrect Joint Order:**

- If the joint indices are not mapped correctly, the plotted connections may be linking the wrong points, causing the figure to appear disorganized.

- **Normalization Issues:**

- The data may have required further normalization or transformation to align all joints correctly.
- A better approach would be to center the body around the origin and ensure consistent scale across frames.

Final Thought Process

At first, the data just looked like a big mess of numbers. But step by step, by visualizing, connecting joints, and animating frame by frame, it started to actually

look like a dancing human! The final result was a 3D animation where we could clearly see the dancer moving through space.

Would I do anything differently?

There were some time constraints, but I would definitely fidget with other datasets to get better optimised results. This was my first time working with this much of a tech stack, reading the [paper](#) to write the code and visualize the data in 3D. But overall, this was a solid first step in visualizing dance movements! This was a great starting point and I would love to work more on this topic under the guidance of the mentors and contribute to this project.

Second Task: Train a multimodal model of dance & text

Code: [code2](#)

For the second part of the project, the goal was to train a neural network that maps both dance movements and their corresponding text descriptions into a shared embedding space using contrastive learning. The idea was that similar dances should have embeddings close to their textual descriptions, while different dances should be far apart.

This task was significantly more complex than the first one, since it required deep learning, natural language processing (NLP), and sequence modeling for dance movements.

What I Used in the Code for the Second Task

For the contrastive learning model, I used LSTMs instead of GNNs. Here's why:

1. The goal was to embed dance sequences, not predict motion – Since the task was to embed short dance phrases into a shared space with text descriptions, an LSTM was sufficient to encode motion temporally.
2. Simpler to Train – Training a GNN requires designing a skeleton-based graph and ensuring the model learns meaningful relationships between joints. For this task, an LSTM was easier to optimize and train.
3. Dataset Constraints – The dataset likely contains multiple sequences but may not be structured optimally for a GNN. Using an LSTM meant I could process the data without needing to explicitly model joint connections.
4. Focus on Time-Series Encoding – Since I needed a fixed-length representation of movement sequences, an LSTM was a natural choice because it compresses time-dependent features effectively.

In contrast, if the task required learning structural dependencies between joints (e.g., generating realistic motion), I would have considered using a GNN-based approach. However, since the focus was embedding dance movements and text into a shared space, LSTMs were the better fit for the contrastive learning framework.

Approach & Thought Process

1. Encoding Dance Sequences

The dance data consists of sequences of joint positions over time. To embed this into a meaningful representation, I:

- Used an LSTM network to process the sequence, since LSTMs are designed for time-series data.
- Made the LSTM bidirectional, allowing it to capture both past and future movements.
- Took the final hidden state from the LSTM and projected it into a 256-dimensional embedding space.

2. Encoding Text Descriptions

Since each dance phrase comes with a natural language description, I needed a way to convert text into a numerical embedding. For this:

- I used DistilBERT, which transforms sentences into a 768-dimensional vector.
- Reduced the dimensionality using a fully connected layer, mapping it to the same 256-dimensional space as the dance embeddings.
- This ensures that dance sequences and their descriptions are in a shared vector space, making similarity comparisons possible.

3. Training with Contrastive Learning

Once both dance and text were converted into embeddings, the next step was training the model. The goal was:

- If a dance matches a text description → Their embeddings should be close together.
- If a dance doesn't match a text description → Their embeddings should be far apart.

I used NT-Xent Loss (Normalized Temperature-scaled Cross-Entropy Loss) for this. It works by:

1. Computing cosine similarity between all dance-text pairs in a batch.
2. Ensuring that correct pairs have high similarity and incorrect pairs have low similarity.

3. Using a temperature scaling factor to control how aggressively the model pulls embeddings together/apart.

{

LSTM:

1. Designed for Sequential Data – Since dance movements are inherently time-dependent, LSTMs handle temporal dependencies well.
 2. Good at Capturing Long-Term Dependencies – LSTMs have gating mechanisms (forget, input, and output gates) that help retain relevant motion information over long sequences.
 3. Computationally Simpler than GNNs – LSTMs operate on linear sequences, which makes them easier to train than GNNs that require graph processing.
 4. Works Well with Small Datasets – If data is limited, LSTMs can still learn useful representations, whereas GNNs often require larger datasets for effective learning.
 5. Ignores Spatial Structure – LSTMs treat joint positions as independent time-series data, ignoring the fact that joints are connected in a skeletal structure.
 6. Not Ideal for Multi-Joint Coordination – Dance movements involve complex joint interactions (e.g., arm movement affecting torso), which LSTMs do not explicitly model.
 7. Struggles with Complex, Multi-Dimensional Data – Motion capture data includes dozens of joints moving in 3D space, making it harder for LSTMs to learn meaningful features.
-

GNN

1. Captures Spatial Relationships – GNNs process joint data as a graph, preserving the physical connectivity between body parts.
2. Models Both Temporal and Structural Information – When combined with RNNs or Transformers, GNNs can handle both joint relationships and motion over time effectively.
3. More Robust to Variations – GNNs can generalize better to different body types, styles, or missing data since they rely on neighboring joints for context.
4. Scalable to Different Skeletons – If the dataset includes multiple skeleton structures (e.g., human vs. animal motion), GNNs adapt naturally.
5. Computationally Expensive – GNNs require graph construction and matrix multiplications, which increase memory and training time.
6. Requires More Data – Since GNNs learn joint relationships, they perform best with larger datasets. Small datasets may lead to overfitting.
7. More Complex to Implement – Unlike LSTMs, which process sequences directly, GNNs require defining a graph structure (i.e., which joints are connected).

}

Output:

Even though the second task was unsuccessful, I tried my hands on these new problems which was my first time tackling such problem statements. I'm looking forward to working on this project and get to learn from the mentors and peers.

