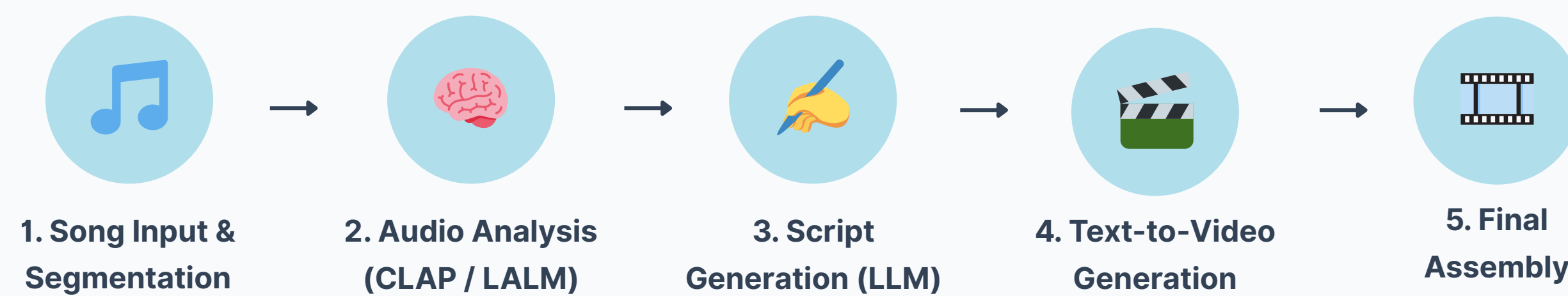


Abstract

Conventional music visualization systems often depend on hand-crafted, abstract transformations with limited expressive power. This research introduces **two novel pipelines** to automatically generate complete **music videos** from any user-specified song, whether instrumental or vocal. By leveraging off-the-shelf deep learning models, the system analyzes audio for musical qualities like **emotion** and **instrumental patterns**, distils them into textual scene descriptions using a language model, and then generates corresponding video clips. A preliminary user study suggests the generated videos show potential for **storytelling**, **visual coherence**, and **emotional alignment** with the music, highlighting the power of deep generative models to expand the art of music visualization.

The Generation Pipeline

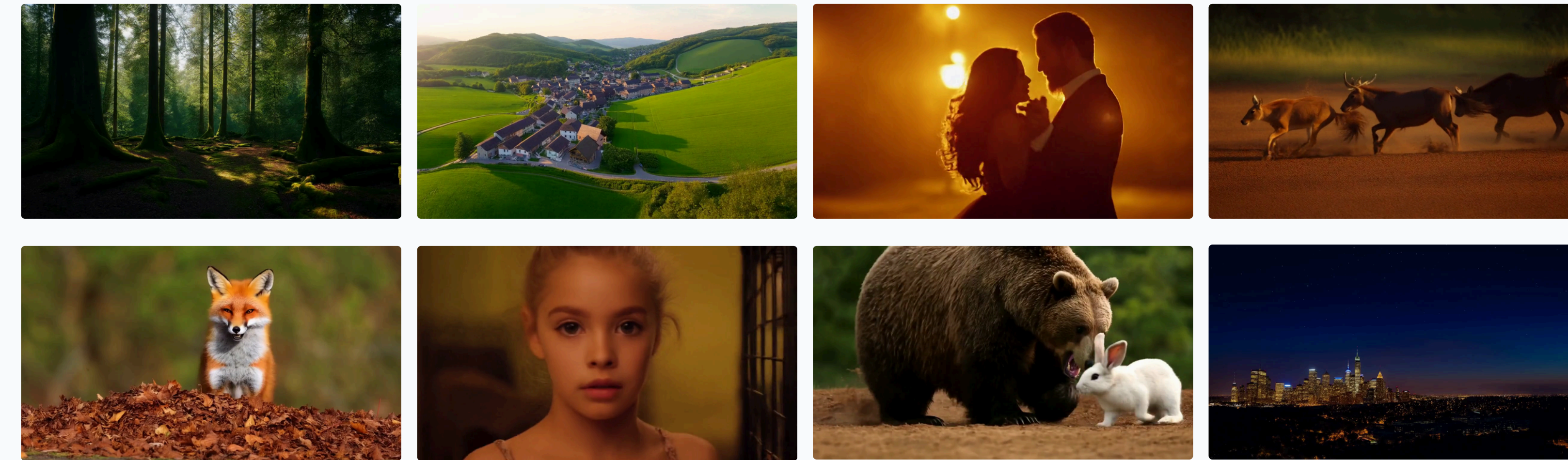
Both proposed pipelines (CLAP and LALM-based) follow a similar multi-step process to generate a music video from any audio, ensuring each stage is **human-readable** and **interpretable**.



Challenges

- **Expressiveness:** Older systems rely on shape/color transformations.
- **Lyric Dependency:** Existing solutions require textual input and guidance.
- **Lack of Cohesion:** Visuals often lack narrative structure and consistency.
- **Quality Assessment:** There are no metrics for music visualization.

Experiments - Selected Frames from Generated Videos



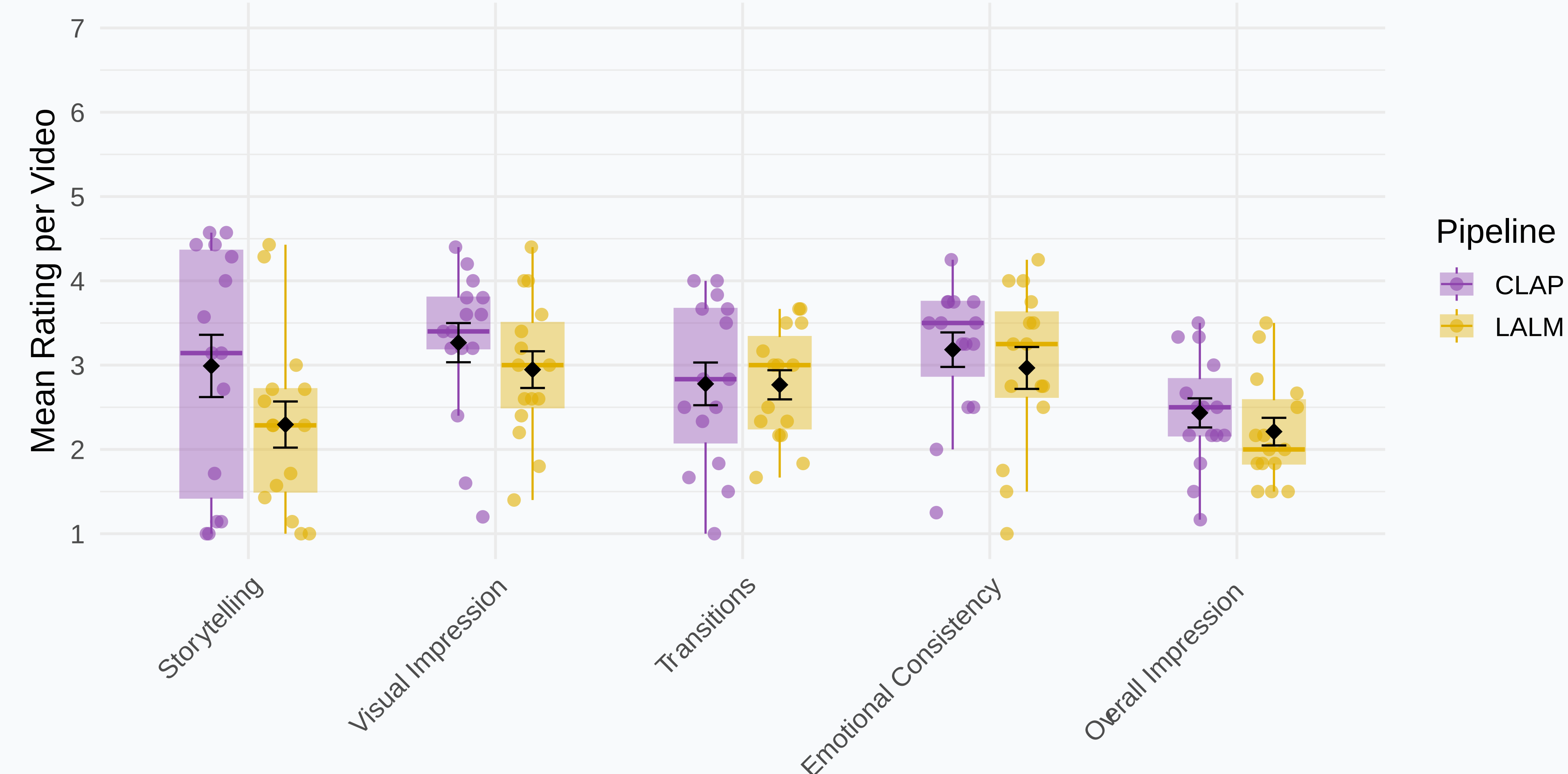
Example music video

Conclusion & Future Work

The proposed pipelines demonstrate significant potential for creating conceptually coherent music videos by integrating audio analysis with LLM-driven text-to-video generation. Next steps will focus on improving **visual consistency**, incorporating **lyrics**, expanding **user studies** and exploring an end-to-end, **music-to-video** based approach without relying on text as intermediary.

Results - Human Evaluation

Distribution of Participant Ratings by Evaluation Dimension and Generation Pipeline

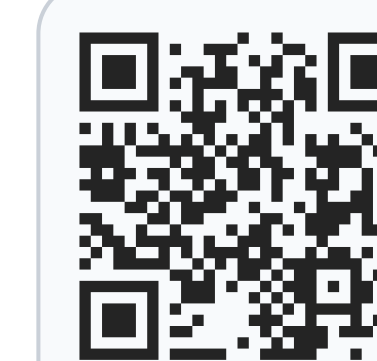


Acknowledgment

This work is funded by the European Union within the Horizon Europe research and innovation program under grant agreement No. 101136006 – XTREME project, coordinated by S.S. Brandt/IT University of Copenhagen, Denmark. Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the European Union, European Union can not be held responsible for them. This project was partially funded by the Pioneer Centre for AI, DNRF grant number P1.

Affiliations

¹IT University of Copenhagen ²Pioneer Centre for Artificial Intelligence ³Aalto University ⁴University of Twente ⁵University of Limerick ⁶University of Nottingham



Paper

