# Generative models for human motion
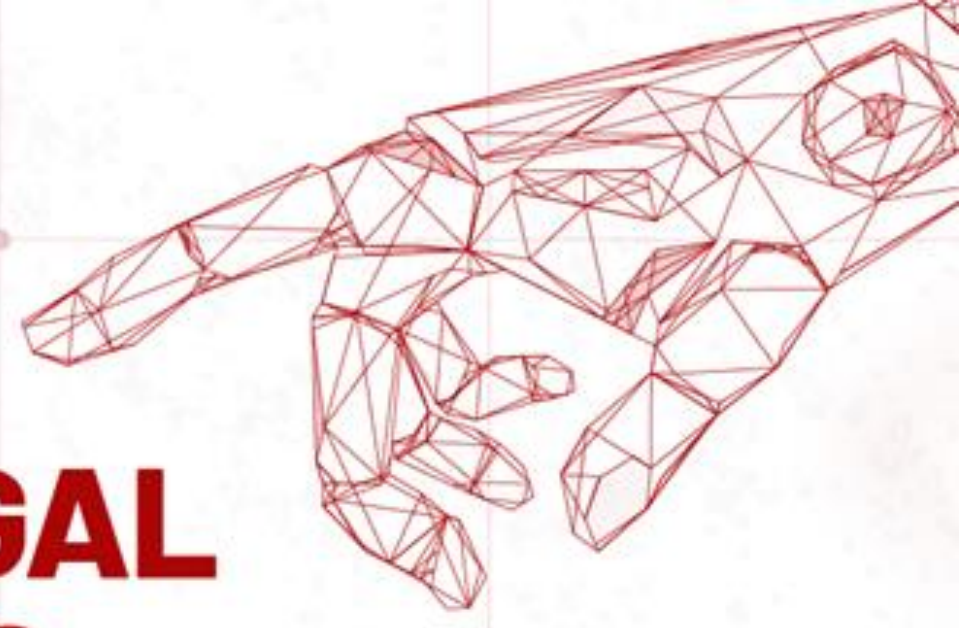
## Jose Ribeiro-Gomes[1], Alexandre Bernardino[1], Fernando de la Torre[2]

1) Instituto Superior Técnico, Lisboa
2) Carnegie Mellon University, Pittsburgh

josepgomes@tecnico.ulisboa.pt, alex@isr.tecnico.ulisboa.pt, ftorre@andrew.cmu.edu

**CMU PORTUGAL SUMMIT 2022**
**NEW FRONTIERS IN TECH**

## Introduction

This work aims to tackle the problem of generating realistic and diverse human motion sequences from highly semantic textual descriptions. We aim to develop a model that can interpret the prompts being fed and generate believable motions that were not part of the training set. We are using an adapted version of the generative variational autoencoder proposed in [1], combined with a language encoder/model (currently CLIP [2]), and training on the highly annotated motion capture dataset BABEL [3]. This generated motion is then fed through a kinematic model to improve the results and realism.

## Contributions

We plan to fully mesh the resulting character, and integrate hand movement in the generated motion, which is often overlooked. Furthermore, we intend to implement physical constraints in our model that will hopefully improve the realism and allow for object and scene interaction.

## Methods

- Motion and text are encoded onto latent space
- Both used for training; only text is used for testing
- Generated motion is sampled from latent space
- Generated motion is fine-tuned using physical constraints, adding kinematic models to latent spaces
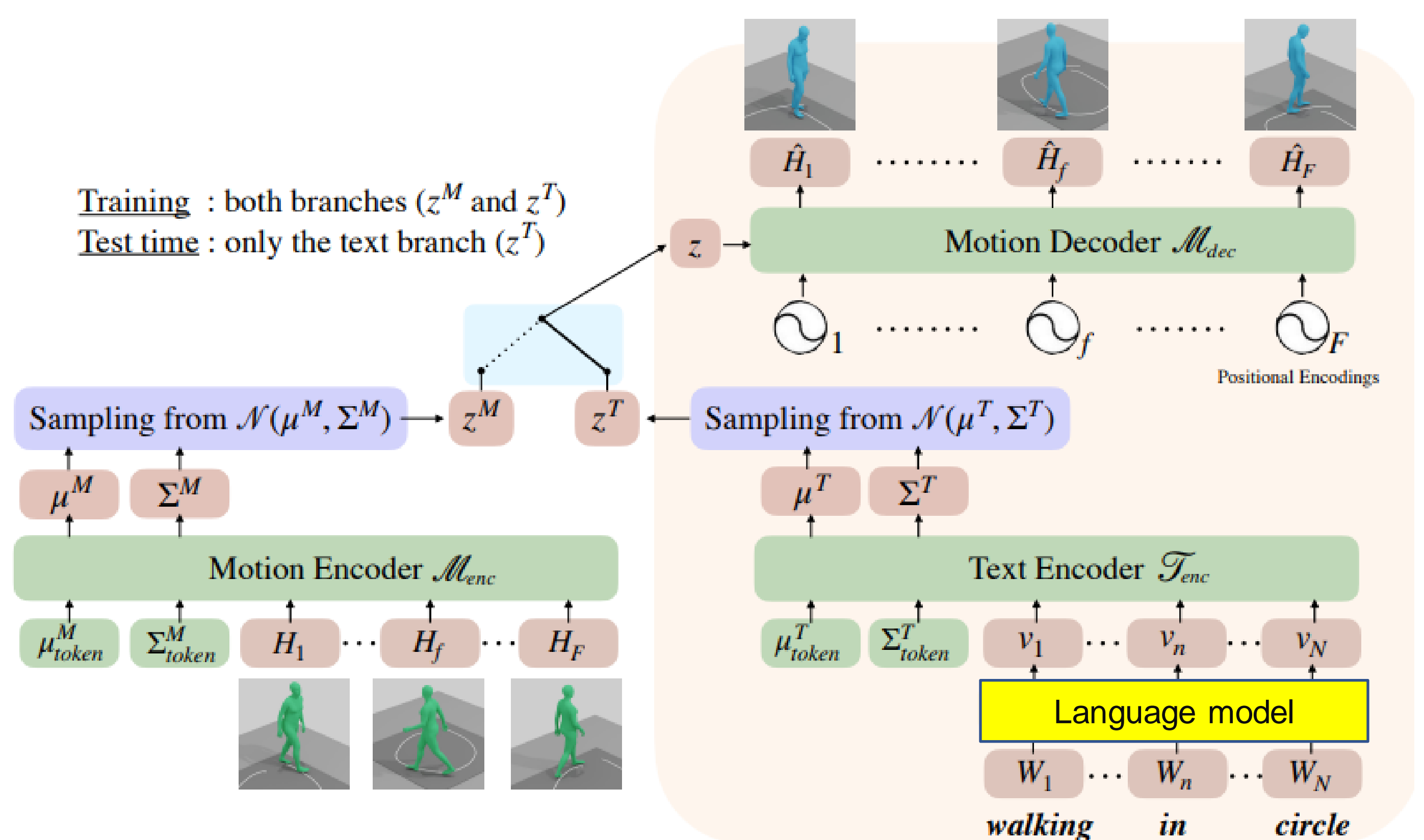


Fig. 1: Overview of the current model under analysis, where motion and text are encoded into the latent space at training time. At test time, only the language branch is used to sample from the latent space and generate the motion. Adapted from [4].

## References

1) Petrovich, Mathis, Michael J. Black, and Gül Varol. "Action-conditioned 3d human motion synthesis with transformer vae." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.
2) Radford, Alec, et al. "Learning transferable visual models from natural language supervision." International Conference on Machine Learning. PMLR, 2021.
3) Punnakkal, Abhinanda R., et al. "BABEL: Bodies, action and behavior with english labels." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.
4) Petrovich, Mathis, Michael J. Black, and Gül Varol. "TEMOS: Generating diverse human motions from textual descriptions." arXiv preprint arXiv:2204.14109 (2022).

## Acknowledgements