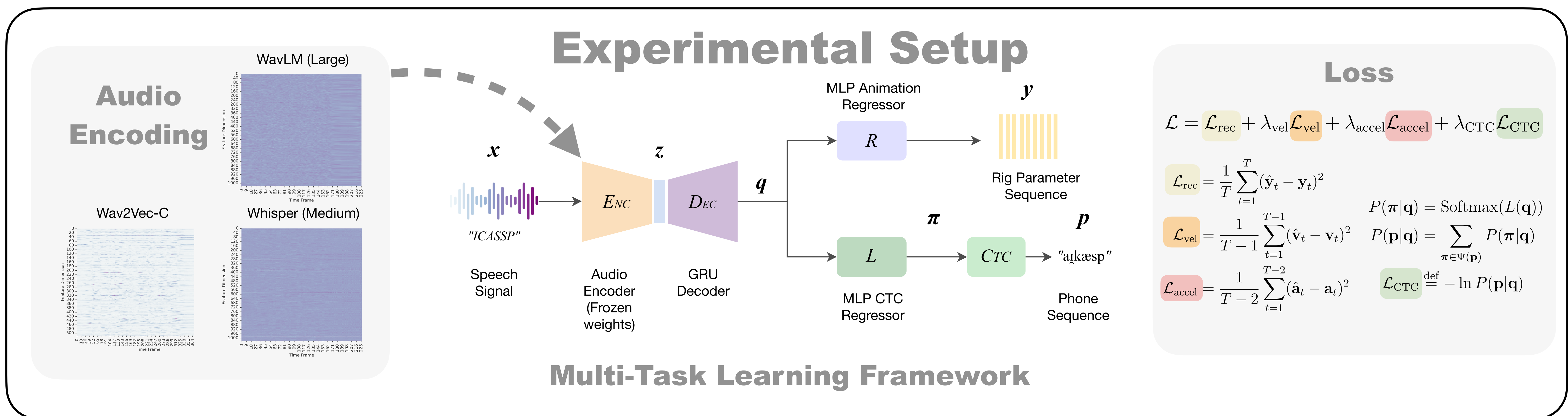
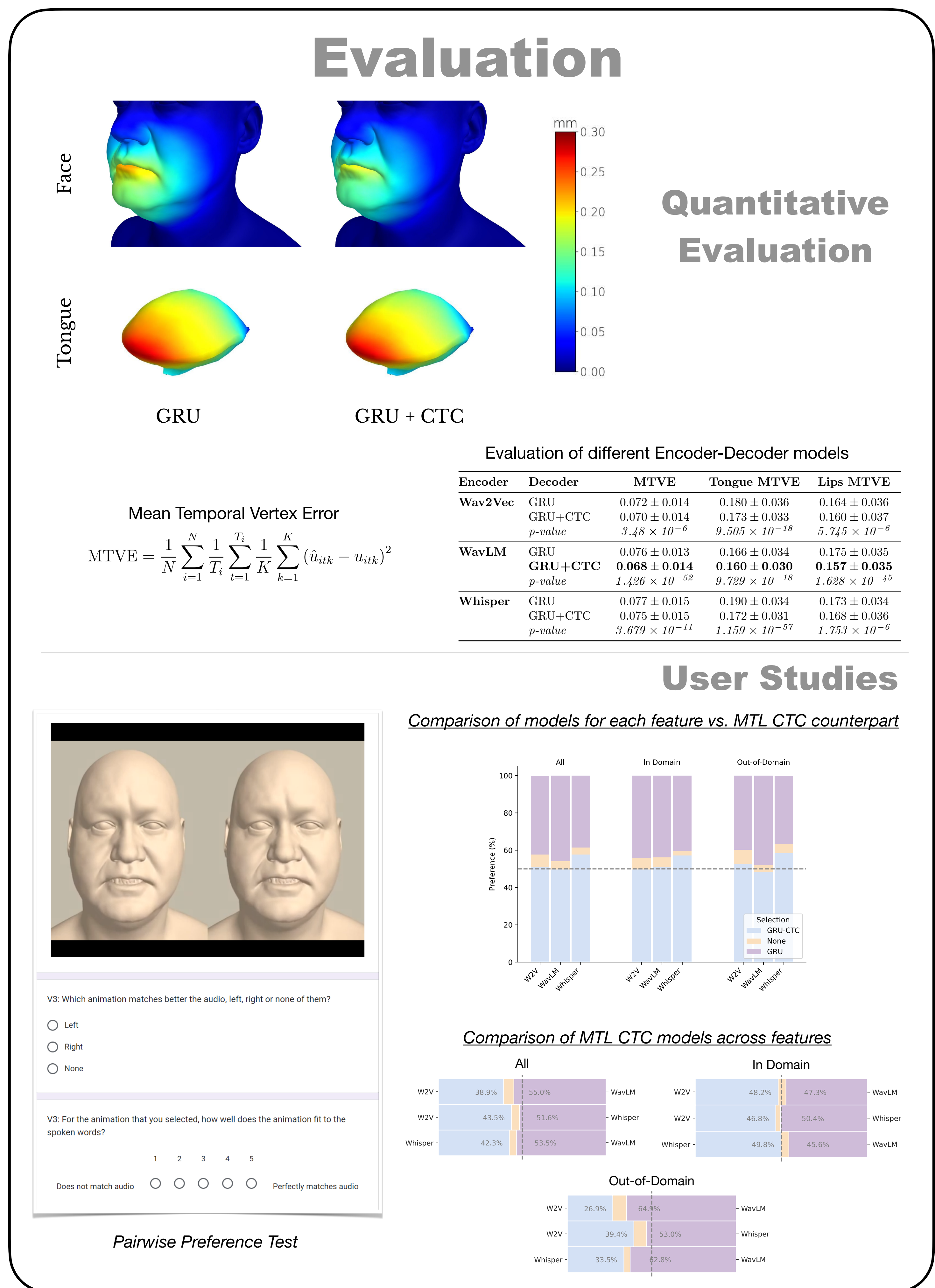
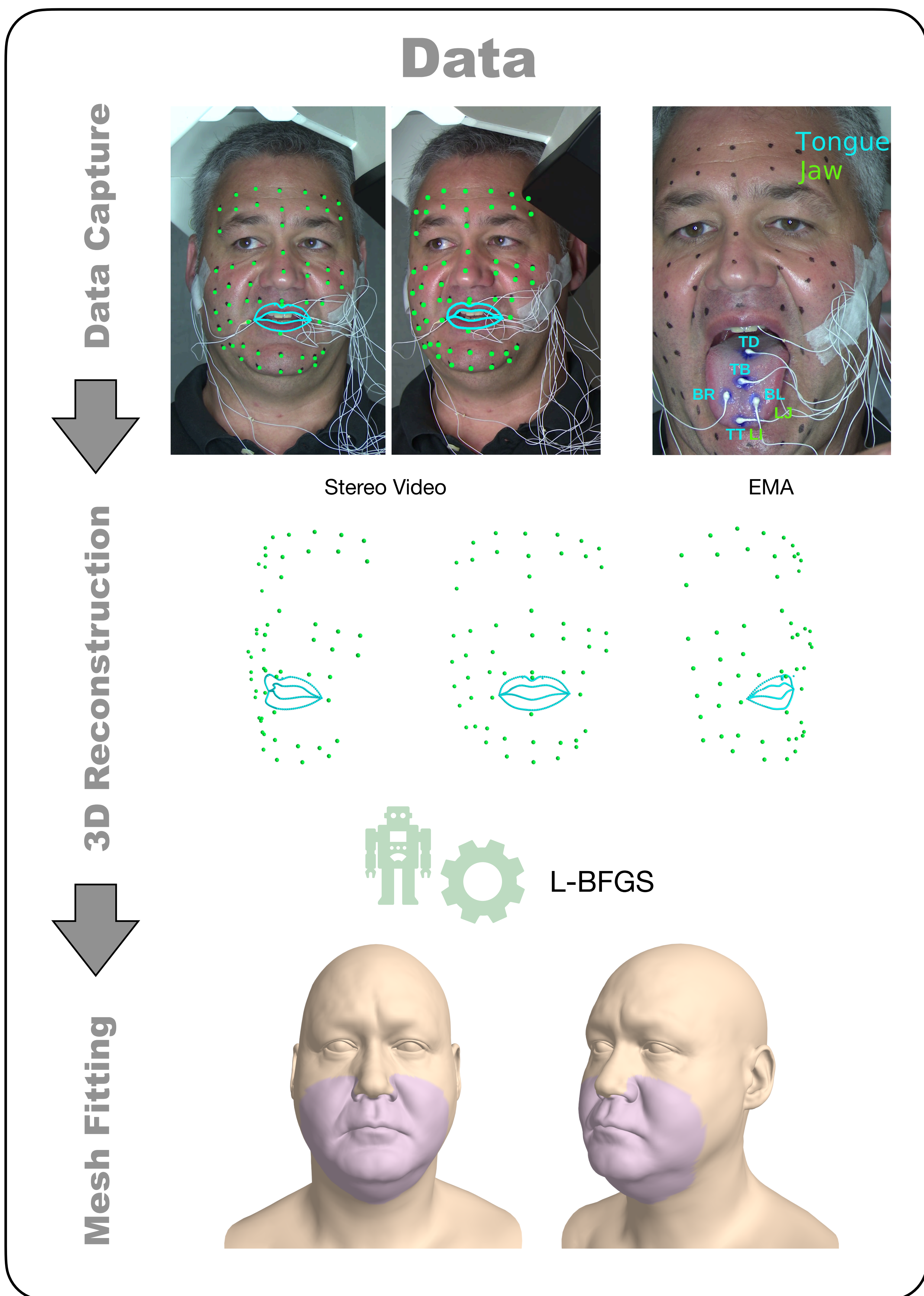


We present **PHISANet**, an end-to-end model that fully animates the lower face, including the jaw and tongue from speech. In this work we:

- Compare the use of **WavLM** [1], **Whisper** [2], and **Wav2Vec** [3] audio features for speech-to-animation.
- Improve articulation animation by regularizing the speech-animation model through **multi-task learning (MTL)** with a **Connectionist Temporal Classification (CTC)** [4] task.



Conclusions

- PHISANet is an encoder-decoder model that generates realistic speech animation by training on data derived from high-quality 3D capture of facial and tongue movements.
- PHISANet delivers realistic and high-quality animations regardless of the gender, age, or language, leveraging on robust audio encodings.
- Incorporating Connectionist Temporal Classification multi-task learning, enhances the realism of the generated speech animations.
- State-of-the-art speech audio encoders, such as Wav2Vec, WavLM, or Whisper, can effectively drive plausible speech animation generation.
- Animations generated by the WavLM-based model were preferred by users due to their *“natural and lifelike motion”*.

[1] S. Chen, C. Wang, Z. Chen, Y. Wu, et al., “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” IEEE Journal of Selected Topics in Signal Processing, vol. 16, pp. 1505–1518, 2021.

[2] A. Radford, J.W. Kim, T. Xu, et al., “Robust speech recognition via large-scale weak supervision,” ArXiv preprint, vol. abs/2212.04356, 2022.

[3] S. Schneider, A. Baevski, R. Collobert, M. Auli, “Wav2Vec: Unsupervised pre-training for speech recognition,” in INTERSPEECH, pp 1–9, 2019.

[4] A. Graves, S. Fernandez, F.J. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in Proc. of ICML, vol. 148, pp. 369–376, 2006.