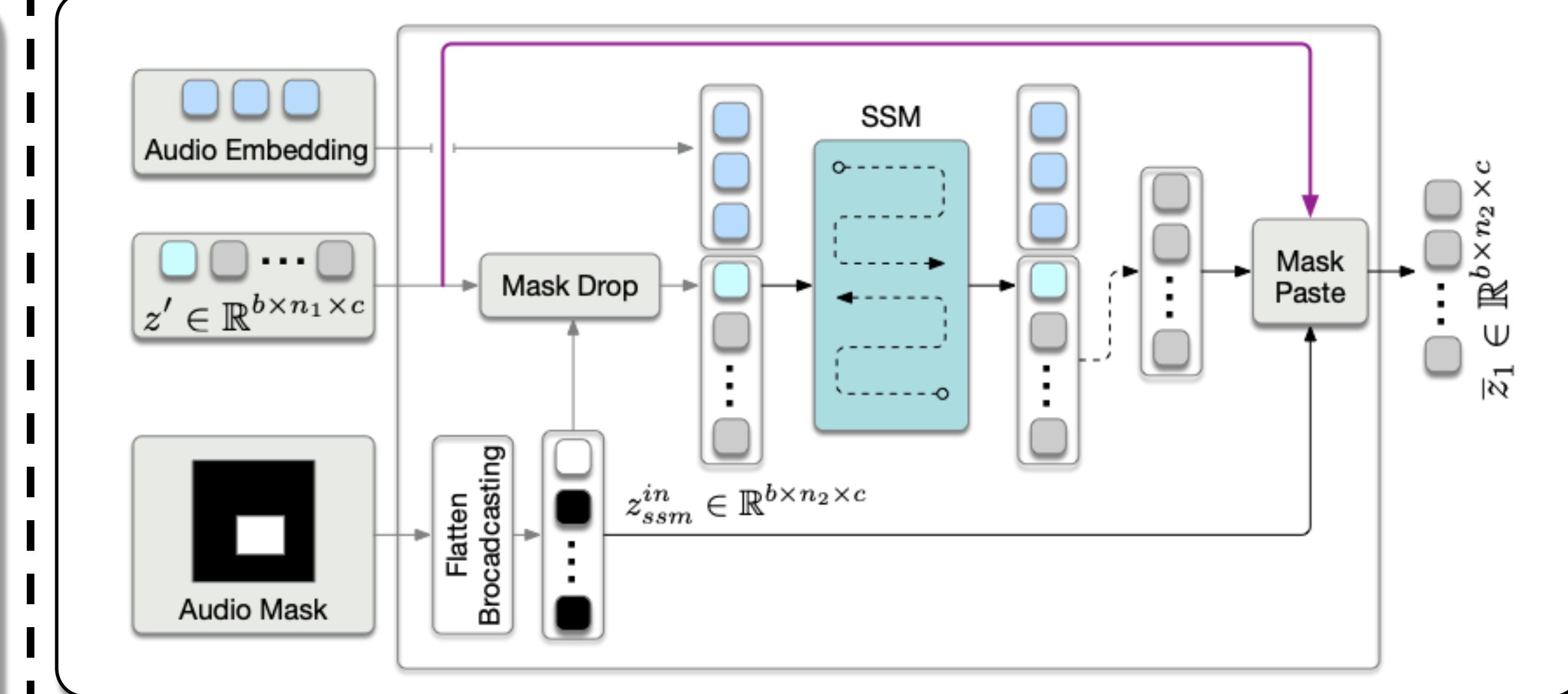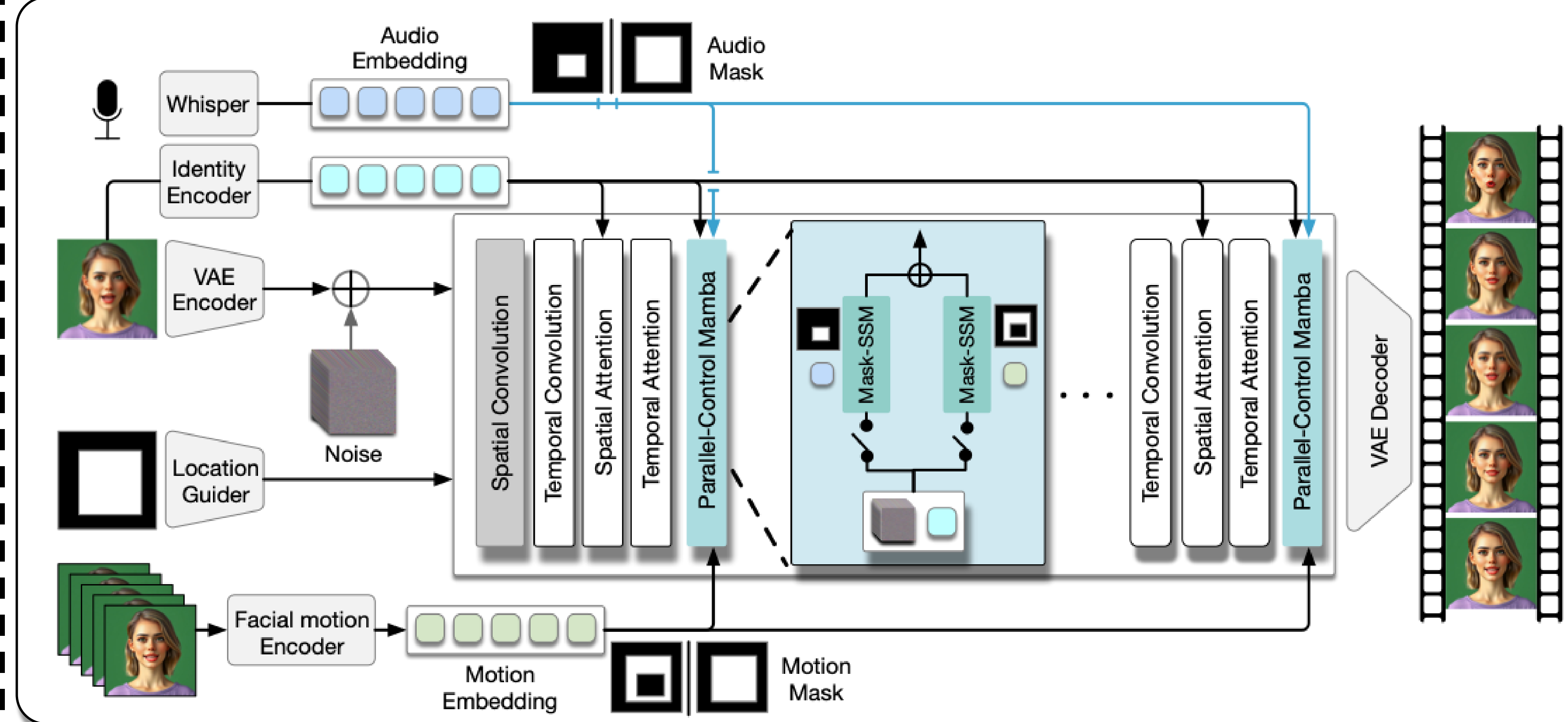# Audio-visual Controlled Video Diffusion with Masked Selective State Spaces Modeling for Natural Talking Head Generation

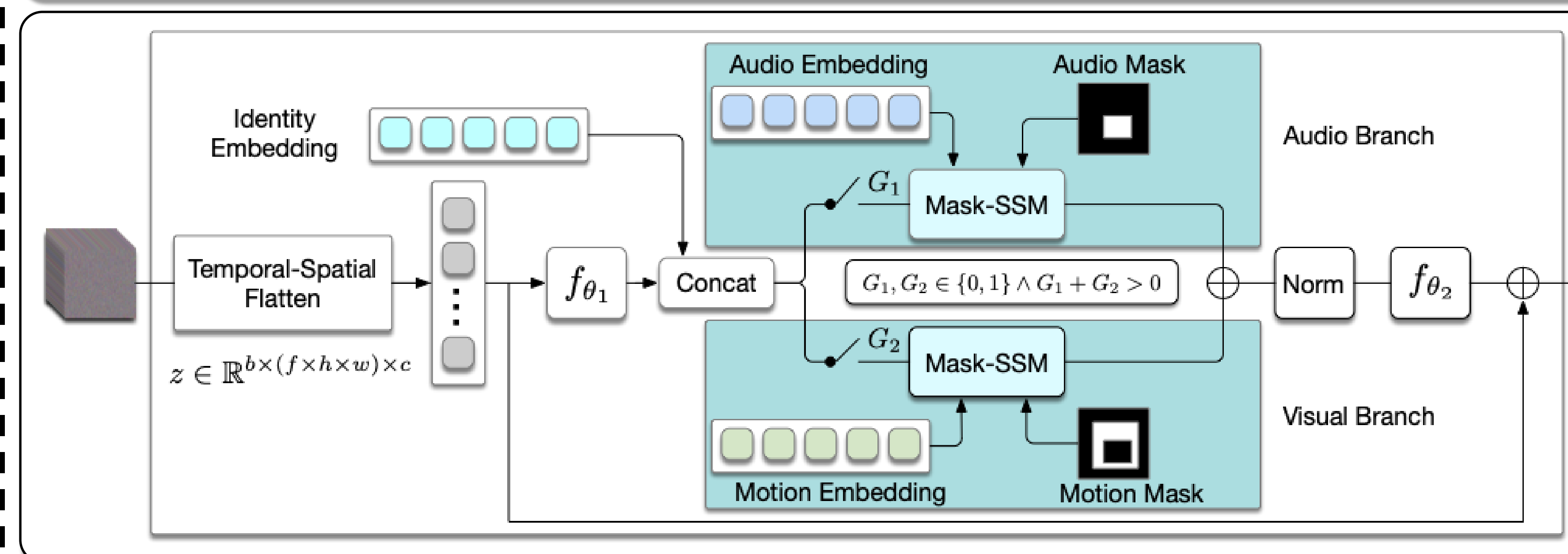Fa-Ting Hong, Zunnan Xu, Zixiang Zhou, Jun Zhou, Xiu Li, Qin Lin, Qinglin Lu, Dan Xu

THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY · 香港科技大學 · 腾讯混元 · Tsinghua University 清华大学 · ICCV OCT 19-23, 2025 HONOLULU HAWAII

**Mask-Drop:** Drops irrelevant tokens to focus on specific control tokens.

**Mask-Paste:** Connects unmodified tokens to generate a complete image.

## Motivation:

- **Overcoming Single-Modality Limitations:** Existing systems are limited to a single control input which restricts realism.
- **Resolving Conflicting Control Signals:** Combining multiple inputs often leads to conflicting and unnatural facial animations.
- **Enabling Multi-Modal Expressiveness:** We aim to seamlessly integrate audio and motion signals for natural and realistically controlled results.

## We Propose:

- **A Novel End-to-End Framework:** An end-to-end framework enables seamless and simultaneous control of generated videos using both audio and fine-grained facial motion signals, leading to more realistic and expressive outputs.
- **A Core Mask-SSM :** It coordinates specific driving signals with their relevant facial regions to resolve signal conflicts.
- **A Parallel-control Mamba Layer:** Effectively coordinates multiple driving signals without conflicts, ensuring smooth integration of audio and facial motion signals.

Project Page

Personal Page