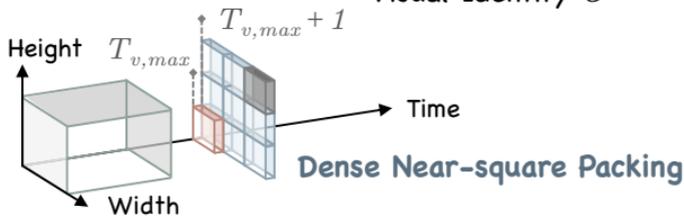
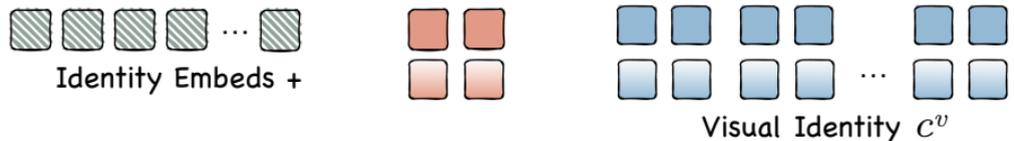
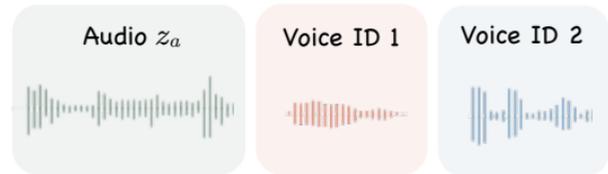




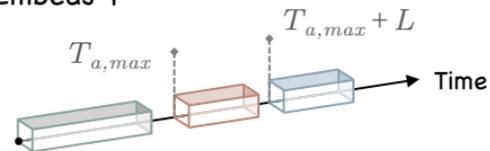
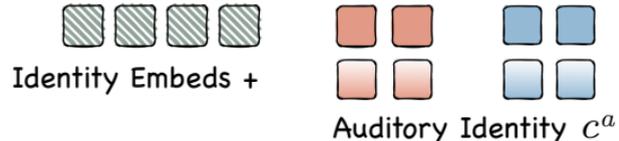
Video Encoder & Patchify



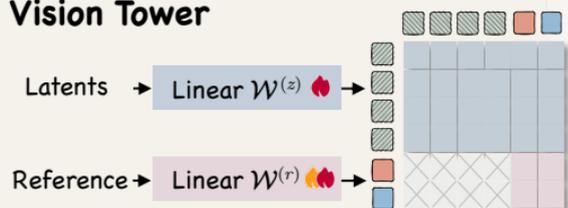
- Trainable in Stage 1
- Trainable in Stage 2&3
- Frozen



Audio Encoder & Patchify



Vision Tower



$\langle \text{REF}_1_S \rangle$ a man with short curly brown hair...
 $\langle \text{REF}_1_E \rangle$ says calmly, $\langle S1 \rangle$ The story has been told. $\langle E1 \rangle$... $\langle \text{REF}_2_S \rangle$ a man with short black hair, wears... $\langle \text{REF}_2_E \rangle$ responds gently, $\langle S2 \rangle$ But the touching part has just begun. $\langle E2 \rangle$

Asymmetric Self-Attention

Cross-Attn

Cross-Attn V2A

Norm + FFN

$\times N$

Audio Tower

Asymmetric Self-Attention

Structured Caption

Cross-Attn

Cross-Attn A2V

Norm + FFN

$\times N$