



- Frozen module
- Training module
- Slow AR output latent
- Quantized linguistic content feature
- Speaker embedding
- Acoustic codes of different codebooks
- Wait-for-start embedding
- Wait-for-end embedding