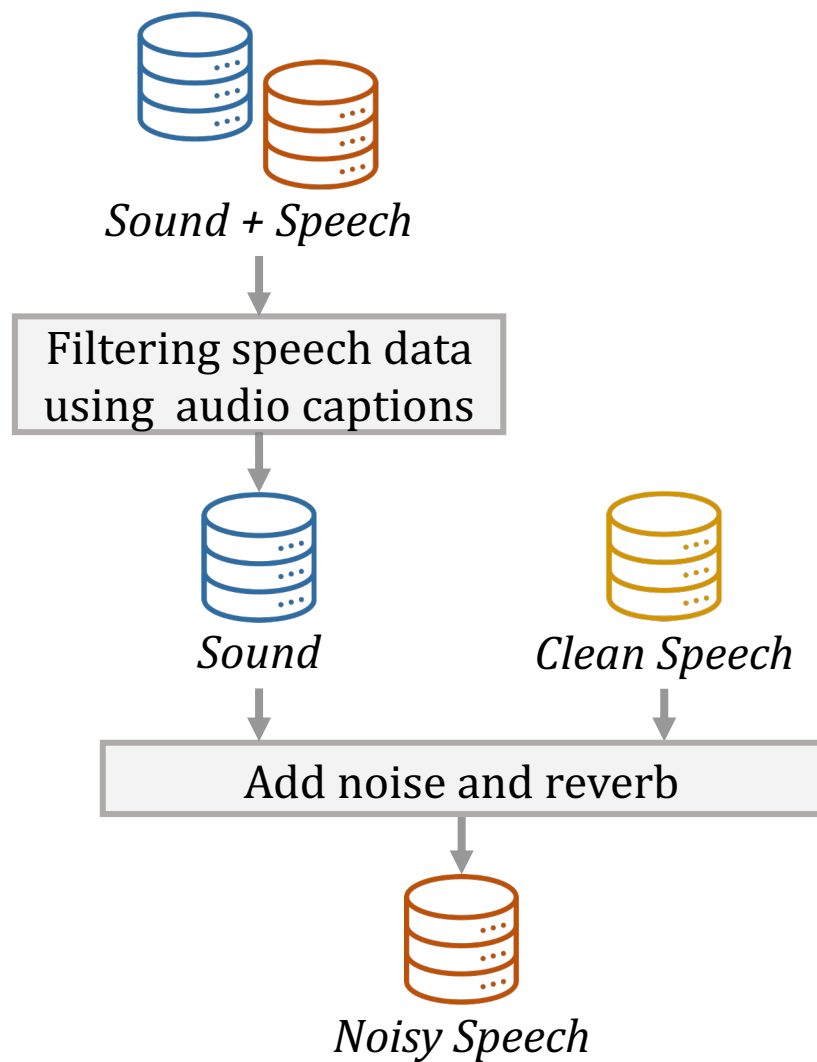
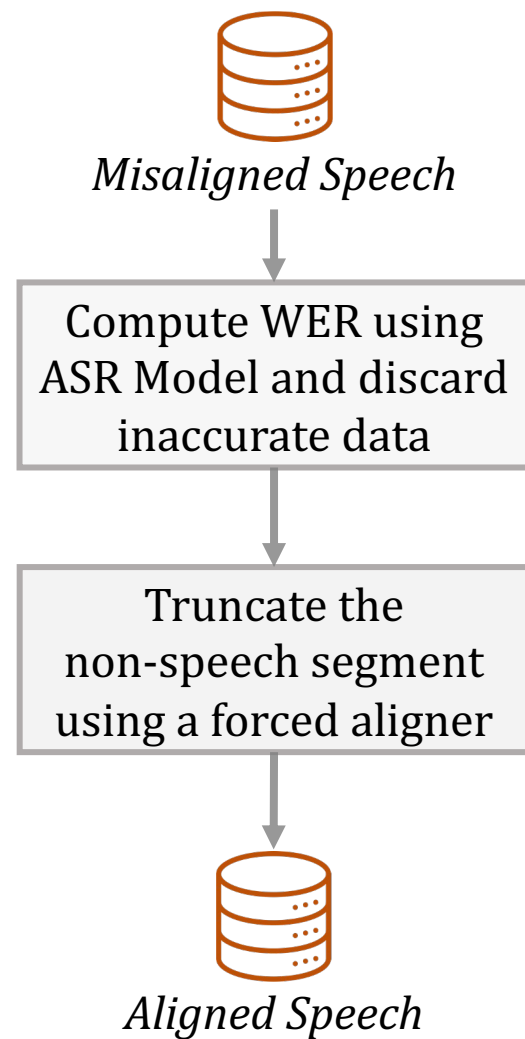


Data Pre-processing

(a) Synthetic Dataset



(b) Filtering



Training

"Good morning!"



Audio

TTS Module

CLAP

*: Fixed parameters

Dual-DiT



Inference

"Good morning!"



Image



Text



Audio

TTS Module

CLIP

I2A
Translator

CLAP

Dual-DiT

