**(A) Data Collection and Filtering**

In-the-wild Data    VAD

**(B) Annotation Tool Construction**
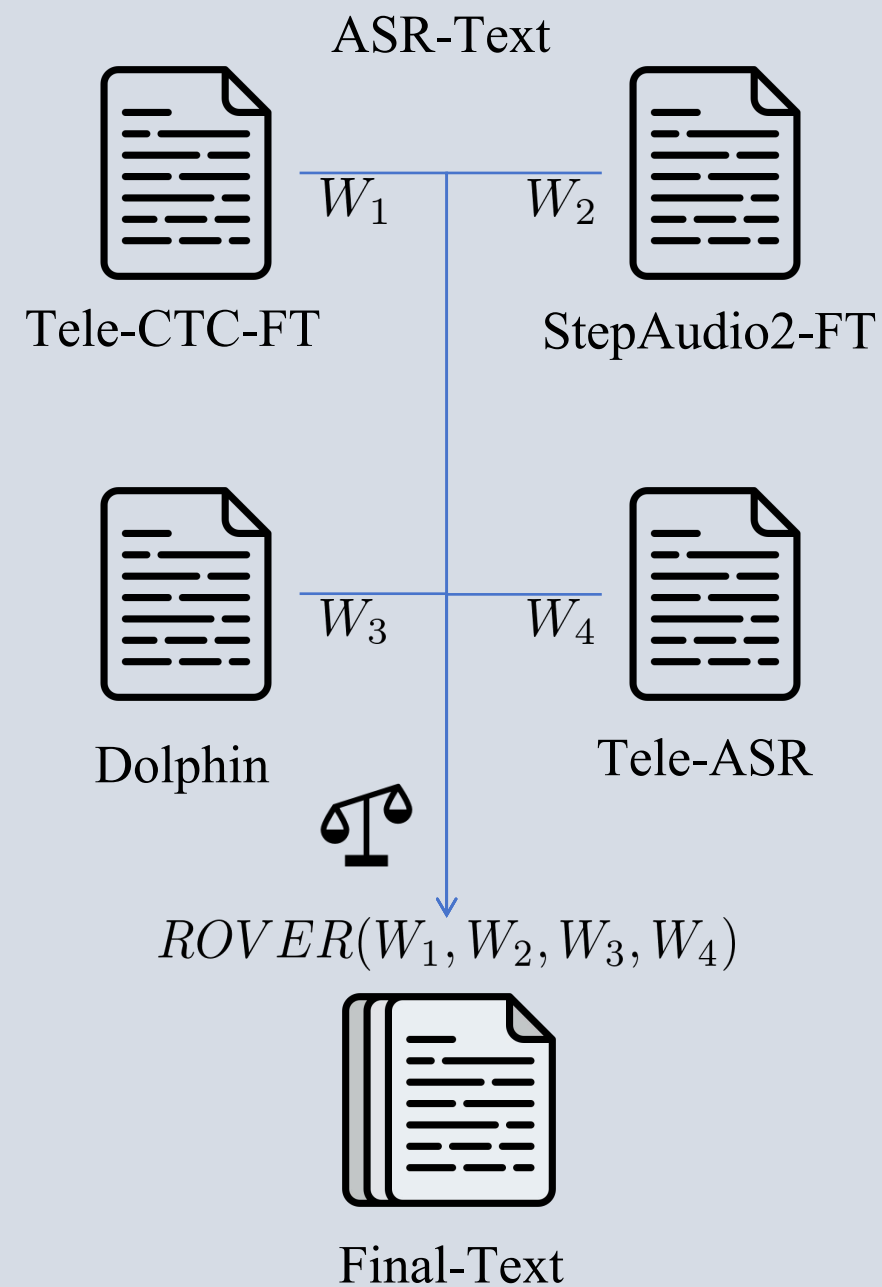
In-house Data → Tele-CTC-FT, Step-Audio2-FT

Fine-Tune

**(C) Automatic Transcription and Fusion**

ASR-Text

Tele-CTC-FT $W_1$    $W_2$ StepAudio2-FT

Dolphin $W_3$    $W_4$ Tele-ASR

$ROVER(W_1, W_2, W_3, W_4)$

Final-Text

**(D) Multi-Dimensional Annotations**

Speaker Attributes

Age    Gender    Multispeaker Detection

Translation

Lexicon-based Translation    Qwen3-8b Refinement

Emotion

(Audio)    (Text)

SenseVoice Emo2vec    Qwen3-8b

Cross Validation Non-neutral

Gemini    Deepseek

Acoustic Features

Energy    Volume

Pitch    Speech Rate

DNSMOS    SNR