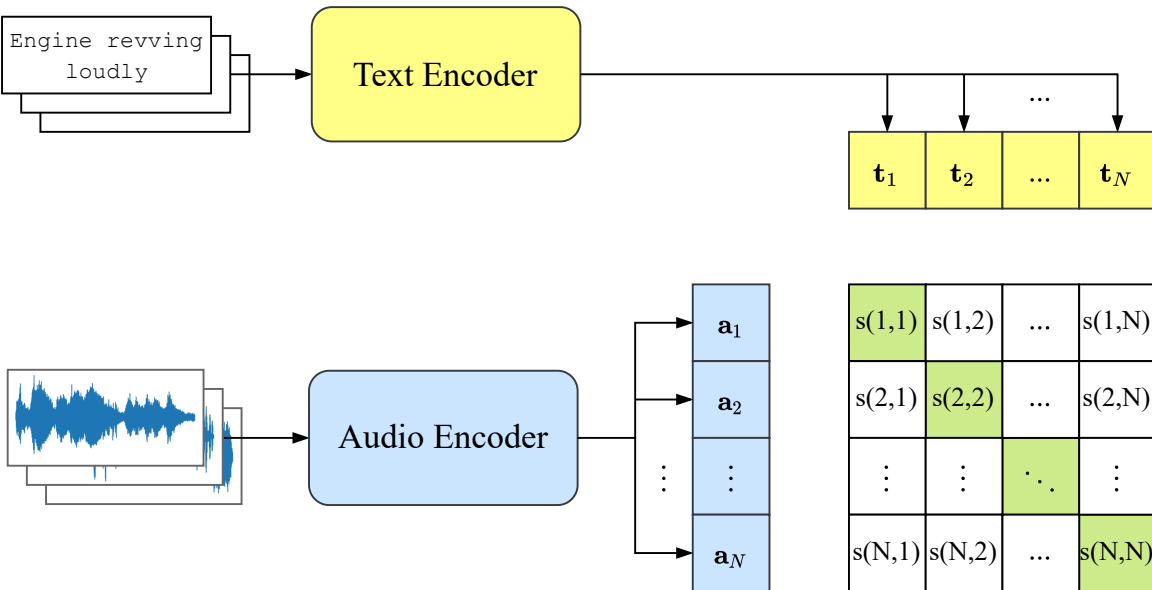
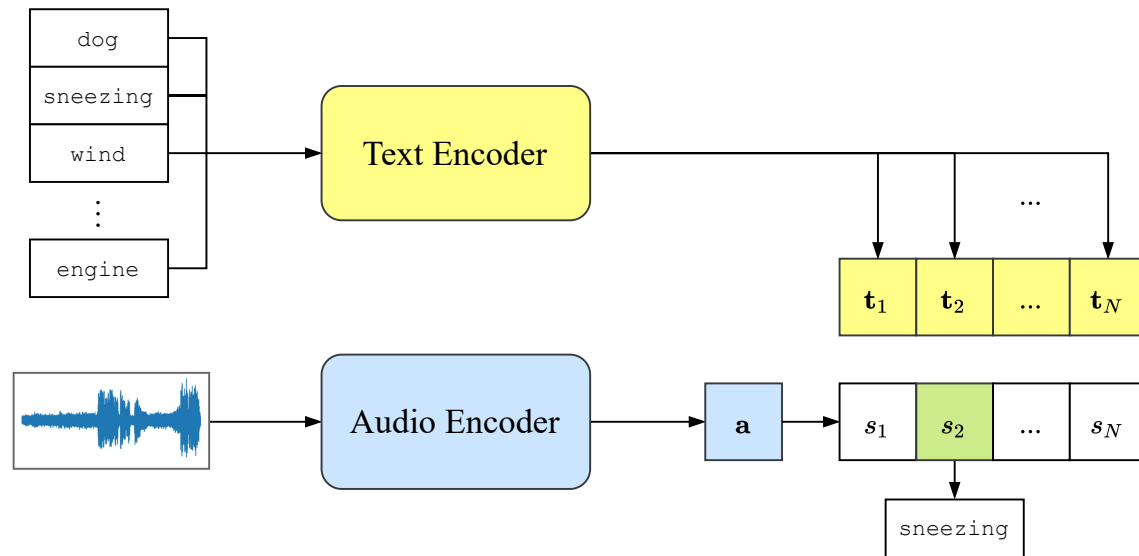


Contrastive pre-training



Zero-shot inference



Audio encoder fine-tuning

