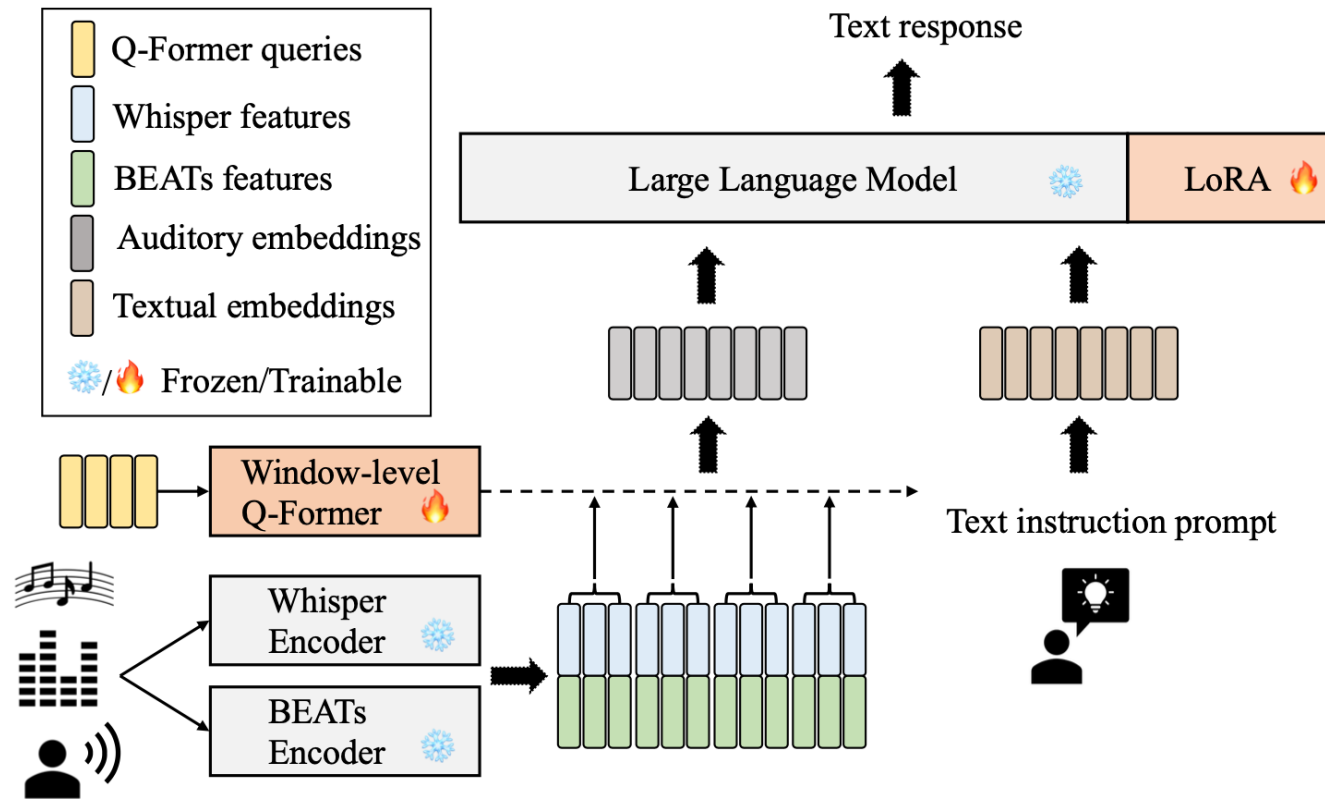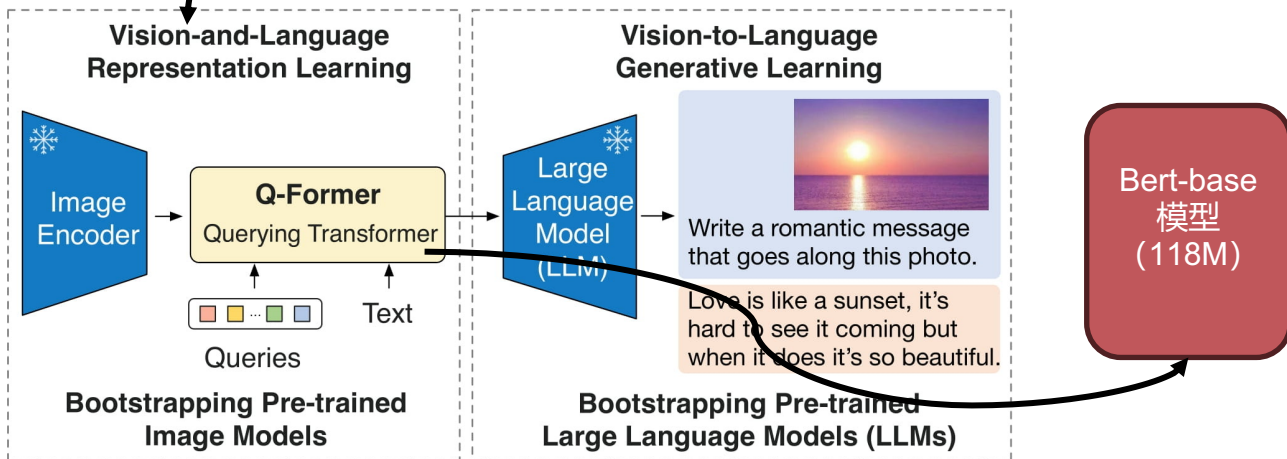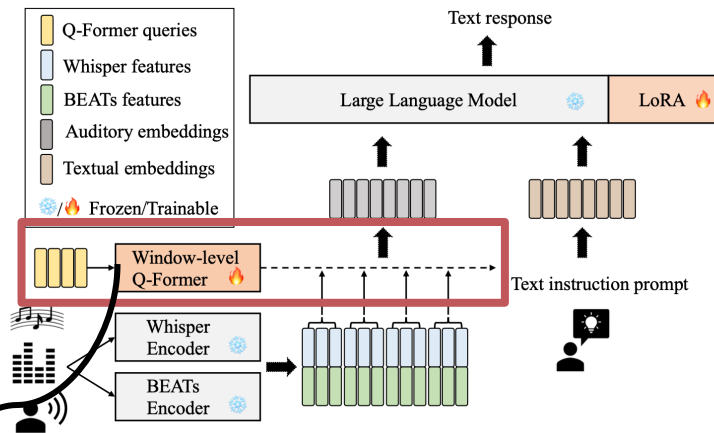# SALMONN: TOWARDS GENERIC <u>HEARING</u> ABILITIES FOR LARGE LANGUAGE MODELS
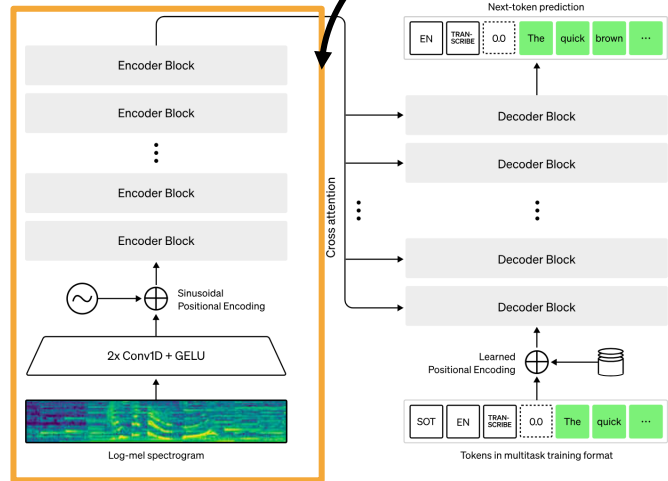
陈安东

# 模型架构及方法

Q-Former queries
Whisper features
BEATs features
Auditory embeddings
Textual embeddings
❄/🔥 Frozen/Trainable

Text response

Large Language Model ❄ | LoRA 🔥

Window-level Q-Former 🔥

Whisper Encoder ❄

BEATs Encoder ❄

Text instruction prompt

**Vision-and-Language Representation Learning**

Image Encoder ❄

**Q-Former** Querying Transformer

Queries | Text

**Bootstrapping Pre-trained Image Models**

**Vision-to-Language Generative Learning**

Large Language Model (LLM) ❄

Write a romantic message that goes along this photo.

Love is like a sunset, it's hard to see it coming but when it does it's so beautiful.

**Bootstrapping Pre-trained Large Language Models (LLMs)**

Bert-base 模型 (118M)

BLIP-2 architecture.

**Legend:**
- Q-Former queries
- Whisper features
- BEATs features
- Auditory embeddings
- Textual embeddings
- ❄️/🔥 Frozen/Trainable

Text response

Large Language Model ❄️   LoRA 🔥

Window-level Q-Former 🔥

Whisper Encoder ❄️

BEATs Encoder ❄️

Text instruction prompt

**Left module (orange box):**

Next-token prediction

EN | TRAN-SCRIBE | 0.0 | The | quick | brown | ...

Encoder Block
Encoder Block
⋮
Encoder Block
Encoder Block

Cross attention

Decoder Block
Decoder Block
⋮
Decoder Block
Decoder Block

Sinusoidal Positional Encoding

2x Conv1D + GELU

Log-mel spectrogram

Learned Positional Encoding

SOT | EN | TRAN-SCRIBE | 0.0 | The | quick | ...

Tokens in multitask training format

**Center boxes:**

多语言
多任务
语音
预训练模型

Bidirectional
语音特征
编码器

**Right module:**

Masked Audio Prediction Loss

Target Labels

$\hat{z}_1$ $\hat{z}_2$ $\hat{z}_3$ $\hat{z}_4$ $\hat{z}_5$
$\hat{z}_6$ $\hat{z}_7$ $\hat{z}_8$ $\hat{z}_9$ $\hat{z}_{10}$
$\hat{z}_{11}$ $\hat{z}_{12}$ $\hat{z}_{13}$ $\hat{z}_{14}$ $\hat{z}_{15}$
$\hat{z}_{16}$ $\hat{z}_{17}$ $\hat{z}_{18}$ $\hat{z}_{19}$ $\hat{z}_{20}$

$Z_1$ $Z_2$ $Z_3$ $\cdots$ $Z_{18}$ $Z_{19}$ $Z_{20}$

Tokenizer

Label Predictor

$\mathbf{r}_1$ [M] $\mathbf{r}_3$ $\cdots$ [M] $\mathbf{r}_{19}$ [M]

Input Feature

$\mathbf{r}_1$ $\mathbf{r}_3$ $\mathbf{r}_{10}$ $\mathbf{r}_{12}$ $\mathbf{r}_{19}$

Transformer Encoder

$\mathbf{e}_1$ $\mathbf{e}_3$ $\mathbf{e}_{10}$ $\mathbf{e}_{12}$ $\mathbf{e}_{19}$

Masked Feature

Projection Layer

# 训练方案

Q-Former queries
Whisper features
BEATs features
Auditory embeddings
Textual embeddings
❄️/🔥 Frozen/Trainable

Text response

Large Language Model ❄️    LoRA 🔥

Window-level
Q-Former 🔥

Text instruction prompt

Whisper
Encoder ❄️

BEATs
Encoder ❄️

Pre-training

Instruction Tuning

Activation Tuning

Pre-training

Instruction Tuning

Activation Tuning

ASR

Audio Caption

大数据量

| Task | Data Source | #Hours | #Samples |
|------|-------------|--------|----------|
| ASR | LibriSpeech + GigaSpeech | 960 + 220 | 280K + 200K |
| En2Zh | CoVoST2-En2Zh (Wang et al., 2021) | 430 | 290K |
| AAC | AudioCaps + Clotho | 130 + 24 | 48K + 4K |
| PR | LibriSpeech | 960 | 280K |
| ER | IEMOCAP Session 1-4 (Busso et al., 2008) | 5 | 4K |
| MC | MusicCaps (Agostinelli et al., 2023) | 14 | 3K |
| OSR | LibriMix (Cosentino et al., 2020) | 260 | 64K |
| SV | VoxCeleb1 (Nagrani et al., 2019) | 1200 | 520K |
| GR | LibriSpeech | 100 | 28K |
| SQA | LibriSpeech | 960 | 280K |
| AQA | WavCaps + AudioCaps | 760 + 130 | 270K + 48K |
| MQA | MillionSong[5] + MusicNet (Thickstun et al., 2017) | 400 + 3 | 48K + 0.3K |
| **Total** | – | $\sim$4400 | $\sim$2.3M |

多任务

Pre-training

Instruction Tuning

Activation Tuning

| Method | En2De↑ | En2Ja↑ | KE↑ | SQQA↑ | SF↑ | Story↑ | SAC↑ |
|---|---|---|---|---|---|---|---|
| w/o Activation | 19.7 | 22.0 | 0.30 | 0.19 (0.29) | 0.33 (0.77) | 7.77 (0.00) | 0.02 (0.04) |
| w/ Activation | 18.6 | 22.7 | 0.32 | 0.41 (0.98) | 0.41 (0.99) | 82.57 (1.00) | 0.50 (0.73) |
| Reference Value | 16.5 | 15.6 | 0.31 | 0.77 (1.00) | 0.46 (1.00) | - | - |

(b) Results of the level 2 and level 3 tasks.

Pre-training

Instruction Tuning

Activation Tuning



h

$B = 0$

Pretrained Weights

$W \in \mathbb{R}^{d \times d}$

$r$

$A = \mathcal{N}(0, \sigma^2)$

$d$

x

| Method | En2De↑ | En2Ja↑ | KE↑ | SQQA↑ | SF↑ | Story↑ | SAC↑ |
|---|---|---|---|---|---|---|---|
| w/o Activation | 19.7 | 22.0 | 0.30 | 0.19 (0.29) | 0.33 (0.77) | 7.77 (0.00) | 0.02 (0.04) |
| Reference Value | 16.5 | 15.6 | 0.31 | 0.77 (1.00) | 0.46 (1.00) | - | - |

(b) Results of the level 2 and level 3 tasks.

# 实验设计

| Task | Test Data | Eval Metrics | Reference Value |
|------|-----------|--------------|-----------------|
| ASR | LibriSpeech test-clean/-other, | %WER | Whisper |
| ASR | GigaSpeech test | %WER | Whisper |
| En2Zh | CoVoST2-En2Zh | BLEU4 | Whisper + Vicuna |
| AAC | AudioCaps | METEOR \| SPIDEr | SOTA (Mei et al., 2023) |
| PR | LibriSpeech test-clean | %PER | WavLM (Chen et al., 2022) |
| ER | IEMOCAP Session 5 | Accuracy | (Wu et al., 2021) |
| MC | MusicCaps | BLEU4, RougeL | SOTA (Doh et al., 2023) |
| OSR | LibriMix | %WER | (Huang et al., 2023c) |
| SV | Voxceleb1 | Accuracy | - |
| En2De | CoVoST2-En2De | BLEU4 | Whisper + Vicuna |
| En2Ja | CoVoST2-En2Ja | BLEU4 | Whisper + Vicuna |
| KE | Inspec (Hulth, 2003) | Accuracy | Whisper + Vicuna |
| SQQA | WikiQA (Yang et al., 2015) | Accuracy (FR) | Whisper + Vicuna |
| SF | SLURP (Bastianelli et al., 2020) | Accuracy (FR) | Whisper + Vicuna |
| Story | AudioCaps | Diversity (FR) | – |
| SAC | In-house Data | Accuracy (FR) | – |

Level 1 微调的任务内

Level 2 分布外任务

Level 3 逻辑推理题
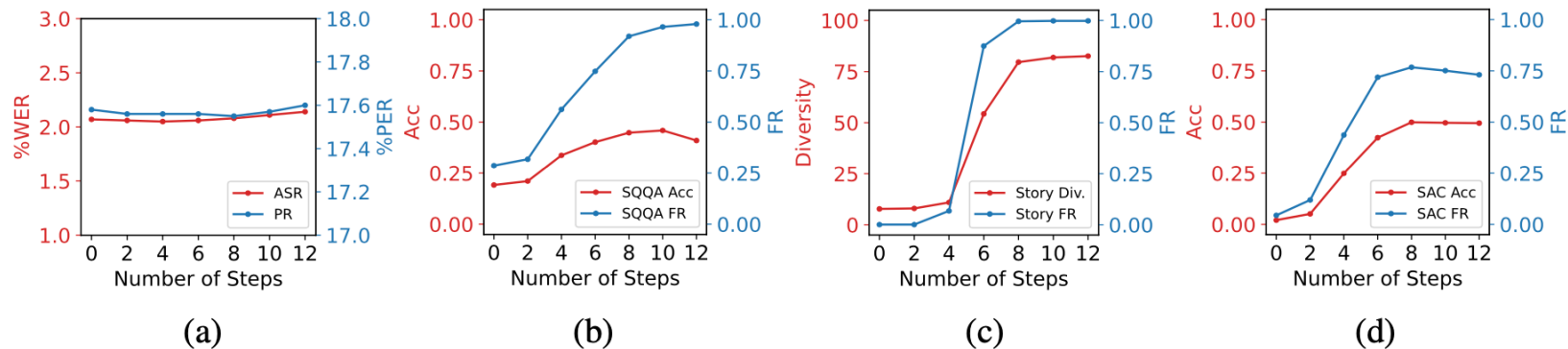
Figure 2: Performance changes on ASR & PR (a), SQQA (b), Story (c) and SAC (d) along with the FR of the emergent abilities against the number of training steps during activation tuning.

Level 1微调的任务内

Level 2 分布外任务

Level 3 逻辑推理题

# SALMONN: Towards Generic Hearing Abilities for Large Language Models

① 当前语音大模型的基本架构变化较为统一

② 机器同传也会面临机器翻译领域相同的问题（多语言/低资源）

③ 机器同传具有独有的特点（单调性/源端不完整），具体较大的挖坑空间

④ 机器同传具有多模态的属性，多模态的同传场景工作不多，但每年都有

# SALMONN: Towards Generic Hearing Abilities for Large Language Models

① 当前语音大模型的基本架构变化较为统一

② 发展趋势使得我们对多模态的组件需要有进一步的认识

③ 机器同传具有独有的特点（单调性/源端不完整），具体较大的挖坑空间

④ 机器同传具有多模态的属性，多模态的同传场景工作不多，但每年都有

# SALMONN: Towards Generic Hearing Abilities for Large Language Models

① 当前语音大模型的基本架构变化较为统一

② 发展趋势使得我们对多模态的组件需要有进一步的认识

③ 大模型在多模态应用使得模态对齐工作成为一些工作的切入点