

# 声音技术的未来 大模型带来的音频算法革新

张俊博

小米AI实验室 语音技术专家

## 讲师简介



小米语音技术专家。



博士毕业于中国科学院声学研究所，多年从事智能语音技术的研究和应用，在语音识别、发音评测、语音合成、音频标记等领域都做过深入的工作，在顶级会议和期刊发表论文 30 余篇，著有出版物《Kaldi 语音识别实战》。

目前在小米负责若干项声学语音新技术的研发。



## 内容提要

- 对大模型的思考
- 小米的音频大模型探索

# 对大模型的思考

# 是“发现”，而不是“发明”

```
LlamaModel(  
  (embed_tokens): Embedding(32000, 256)  
  (layers): ModuleList(  
    (0-7): 8 x LlamaDecoderLayer(  
      (self_attn): LlamaAttention(  
        (q_proj): Linear(in_features=256, out_features=256, bias=False)  
        (k_proj): Linear(in_features=256, out_features=1024, bias=False)  
        (v_proj): Linear(in_features=256, out_features=1024, bias=False)  
        (o_proj): Linear(in_features=256, out_features=256, bias=False)  
        (rotary_emb): LlamaRotaryEmbedding()  
      )  
      (mlp): LlamaMLP(  
        (gate_proj): Linear(in_features=256, out_features=11008, bias=False)  
        (up_proj): Linear(in_features=256, out_features=11008, bias=False)  
        (down_proj): Linear(in_features=11008, out_features=256, bias=False)  
        (act_fn): SiLUActivation()  
      )  
      (input_layernorm): LlamaRMSNorm()  
      (post_attention_layernorm): LlamaRMSNorm()  
    )  
  )  
  (norm): LlamaRMSNorm()  
)
```

Llama2 模型：没有任何模型结构上的创新

原理上是量变，效果上是质变

无法解释，只好说“涌现”

大模型的成功，证明了这样的路线是可行的

**为 AI 研究指明了方向**

# 为什么大模型具备如此神奇的能力?

## 不知道

人类对它的原理还远远称不上理解

但大模型研发并没有技术原理上的门槛

虽然不知道麦克斯韦方程组  
不妨碍古人发明指南针

虽然暂时未能全面理解大模型  
不妨碍我们做出更强的大模型



# 小米自研大语言模型

## 本地化、轻量部署

	平均分	STEM	人文学科	社会科学	其他	中国特定主题
MiLM-1.3B	50.79	40.51	54.82	54.15	53.99	52.26
Baichuan-13B	54.63	42.04	60.49	59.55	56.6	55.72
ChatGLM 2-6B	49.95	41.28	52.85	53.37	52.24	50.58

CMMLU 中文多任务语音理解评估, 2023年8月10日数据

手机端侧大模型部分场景媲美云端



# 大模型 == 大语言模型？

文本形式训练数据相对更易获取和处理  
大模型首先以文本模态出现

但人类更倾向于使用视觉和声音交互



WIKIPEDIA  
The Free Encyclopedia

Search Wikipedia

### Search results

large model

Advanced search: Sort by relevance

Search in:

*The page "**Large model**" does not exist. You can [create a draft](#) and submit it for review or request that a [redirect](#) be created, but consider checking the search results below to see whether the topic is already covered.*

View (previous 20 | [next 20](#)) (20 | 50 | 100 | 250 | 500)

**Large language model**

A **large language model** (LLM) is a type of language **model** notable for its ability to achieve general-purpose language understanding and generation. LLMs...

103 KB (9,271 words) - 01:13, 3 December 2023

# GPT4-V(ision)



User

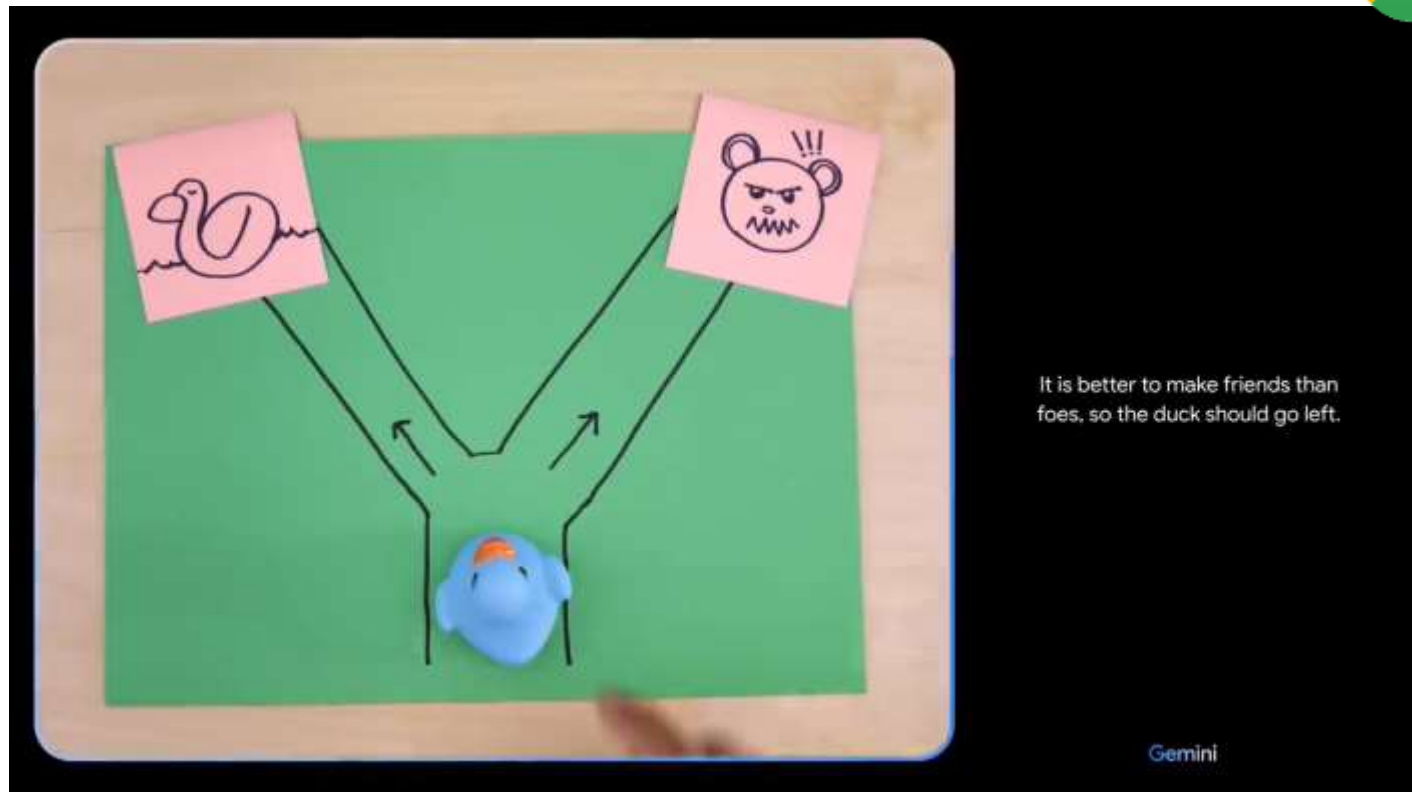
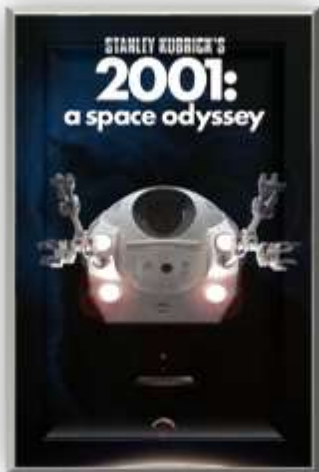
What is unusual about this image?

GPT-4

The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi.

# Gemini

令人震惊的多模态能力  
强人工智能已实现？





# 小米的音频大模型探索

# AI 时代的小米

## 全球最大消费级 IoT 平台

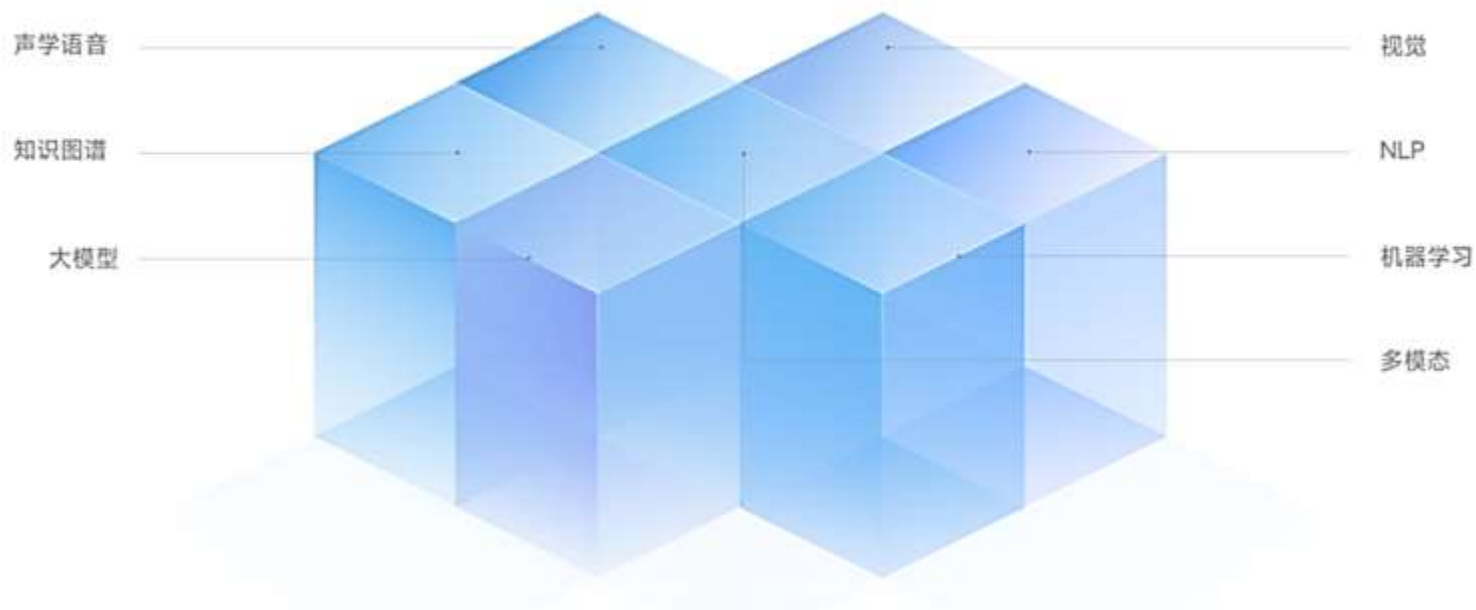
6.99 亿

IoT 平台已连接设备数

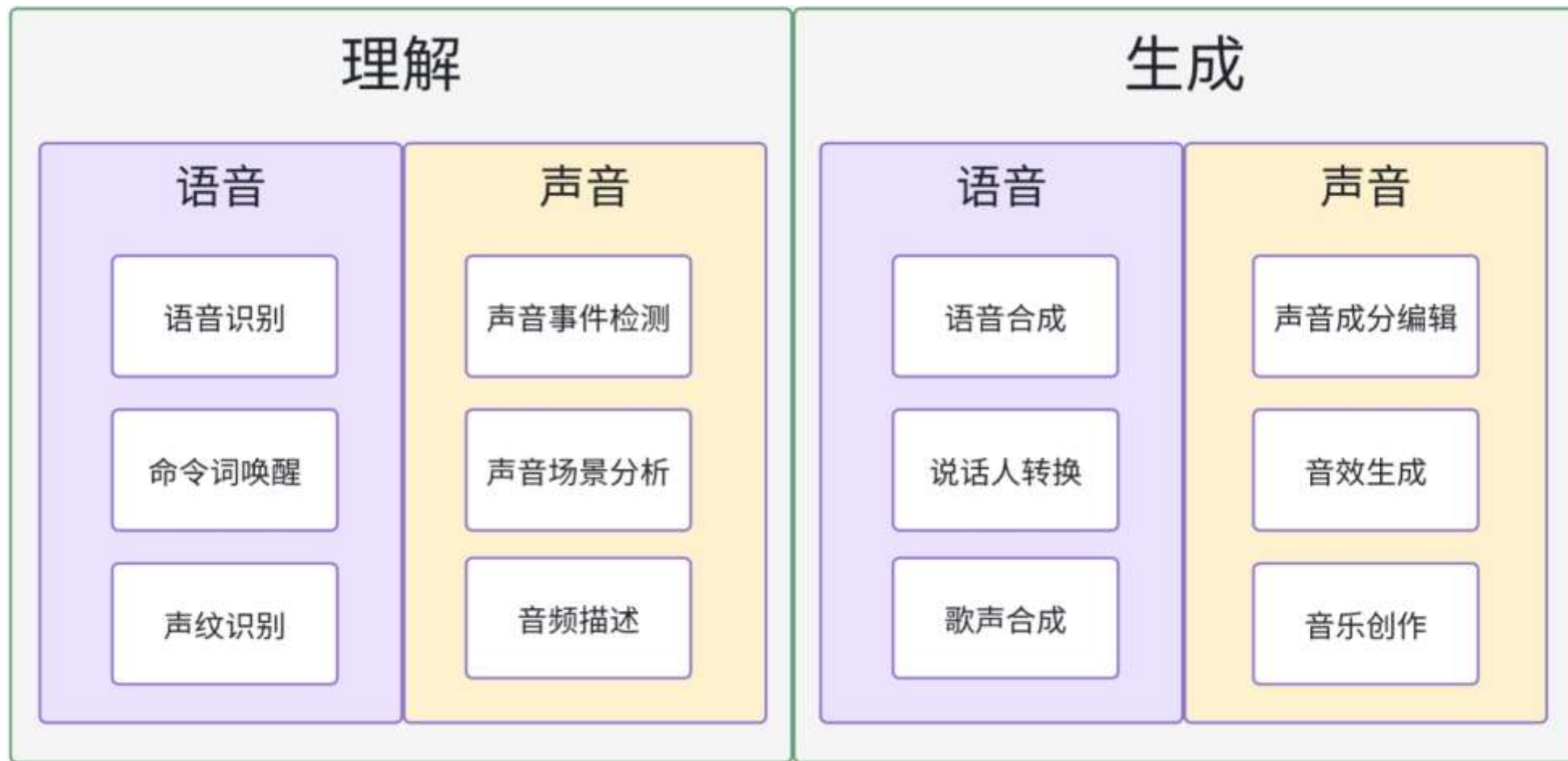
1370 万

拥有 5 个及以上小米 IoT 设备的用户数

视觉 | 声学语音 | NLP | 知识图谱 | 机器学习 | 大模型 | 多模态



# 小米声学语音技术







# Whisper: 大模型语音识别



Whisper examples:

Speed talking ▾



This is the Micro Machine Man presenting the most midget miniature motorcade of Micro Machines. Each one has dramatic details, terrific trim, precision paint jobs, plus incredible Micro Machine Pocket Play Sets. There's a police station, fire station, restaurant, service station, and more. Perfect pocket portables to take any place. And there are many miniature play sets to play with, and each one comes with its own special edition Micro Machine vehicle and fun, fantastic features that miraculously move. Raise the boatlift at the airport marina. Man the gun turret at the army base. Clean your car at the car wash. Raise the toll bridge. And these play sets fit together to form a Micro Machine world. Micro Machine Pocket Play Sets, so tremendously tiny, so perfectly precise, so dazzlingly detailed, you'll want to pocket them all. Micro Machines are Micro Machine Pocket Play Sets sold separately from Galoob. The smaller they are, the better they are.

Whisper examples:

French ▾



Whisper is an automatic speech recognition system based on 680,000 hours of multilingual and multitasking data collected on the Internet. We establish that the use of such a number of data is such a diversity and the reason why our system is able to understand many accents, regardless of the background noise, to understand technical vocabulary and to successfully translate from various languages into English. We distribute as a free software the source code for our models and for the inference, so that it can serve as a starting point to build useful applications and to help progress research in speech processing.

Whisper examples:

K-Pop ▾



While darkness was my everything  
I ran so hard that I ran out of breath  
Never say time's up  
Like the end of the boundary  
Because my end is not the end

Whisper examples:

Accent ▾



One of the most famous landmarks on the Borders, it's three hills and the myth is that Merlin, the magician, split one hill into three and left the two hills at the back of us which you can see. The weather's never good though, we stay on the Borders with the mists on the Yildens, we never get the good weather and as you can see today there's no sunshine, it's a typical Scottish Borders day.



# Whisper 原理有何不同?

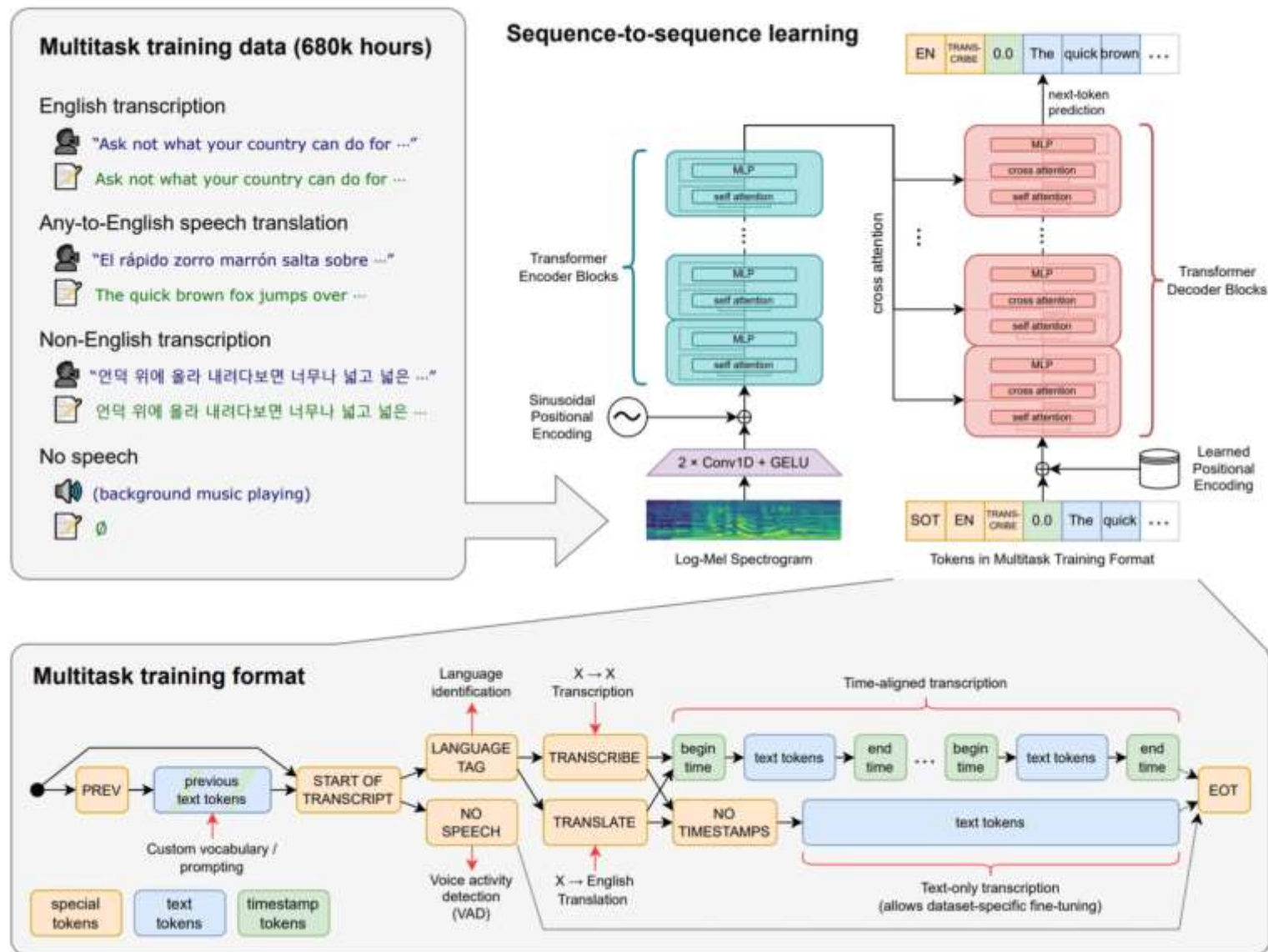
更先进的模型结构? No

模型结构并无不同

多语种训练数据

带有多任务标签

680,000 hours of multilingual and multitask supervised data collected from the web

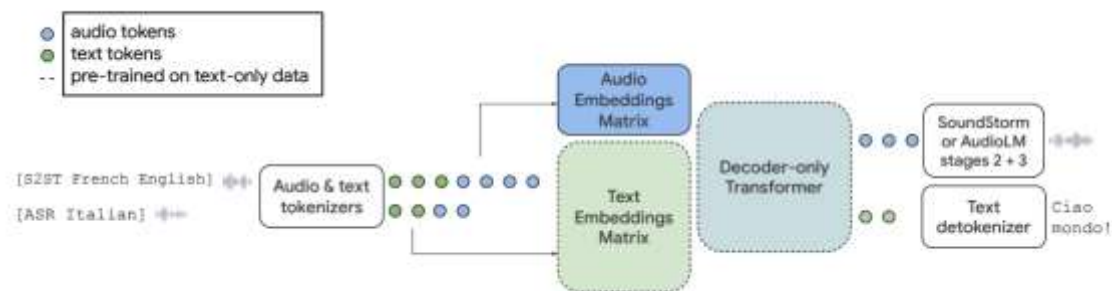


# AudioPaLM: 多语种语音直译



大语言模型作为模型骨架和初始化参数

多语种音频和文本数据迭代训练



Mandarin / 中文

AudioPaLM

Original

我非常开心和你合作

Translation with  
AudioPaLM

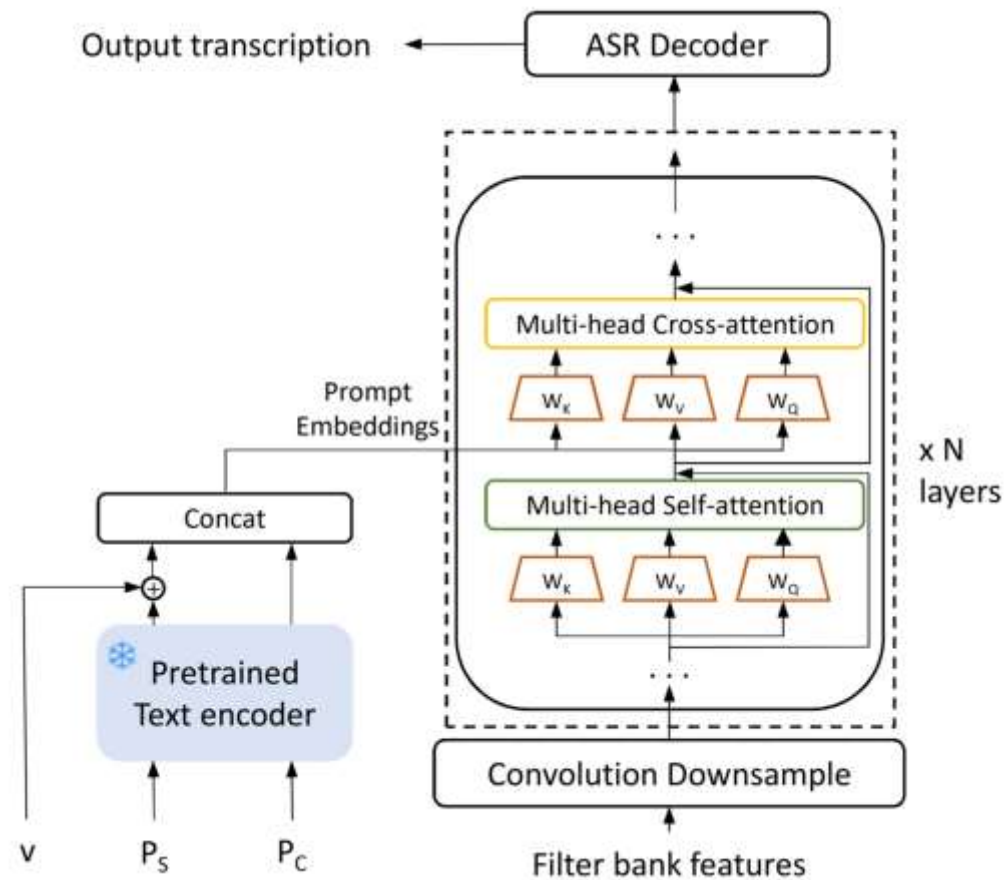
I am very happy to work with you

# 小米 Prompt-ASR

用 prompt 约束语音识别领域，提升识别率

把大语言模型输出通过 cross-attention 联入 encoder

Style Prompt	WITHOUT CASING OR PUNCTUATION
Content Prompt	Welcome to the UEFA Champions League final!
Reference text	TODAY'S MATCH IS BETWEEN REAL MADRID AND LIVERPOOL
Style Prompt	Mixed-cased English with punctuation
Content Prompt	Welcome to the UEFA Champions League final!
Reference text	Today's match is between Real Madrid and Liverpool.



# 基于大模型的语音合成

更加自然



视频来源

<https://www.bilibili.com/video/BV1e84y1U7j4>

支持 Prompt 定制

Microsoft VALL-E X

Look a little closer while our guide lets the light of his lamp fall upon the black wall at your side.

baseline



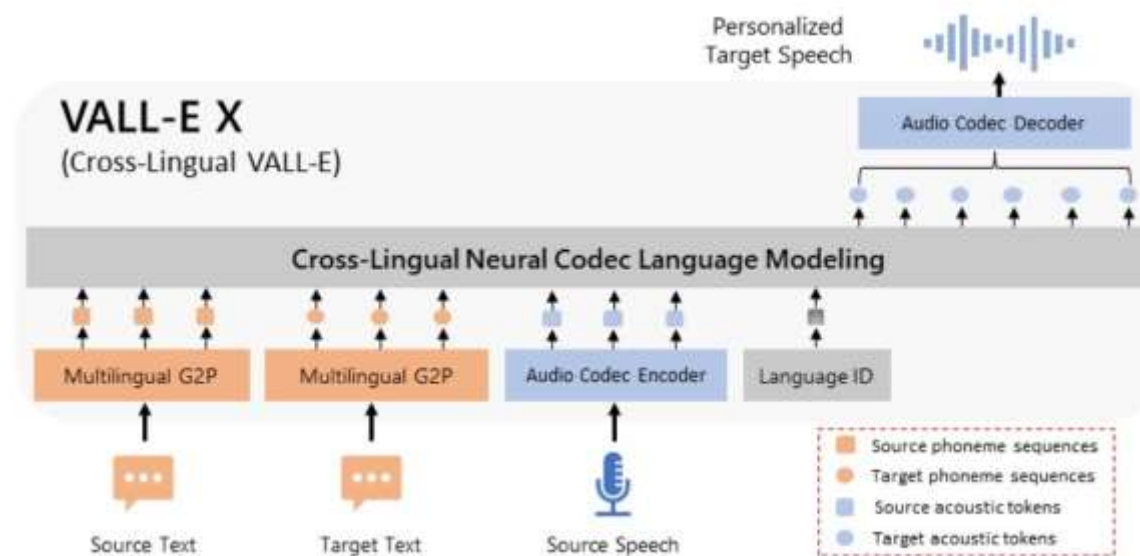
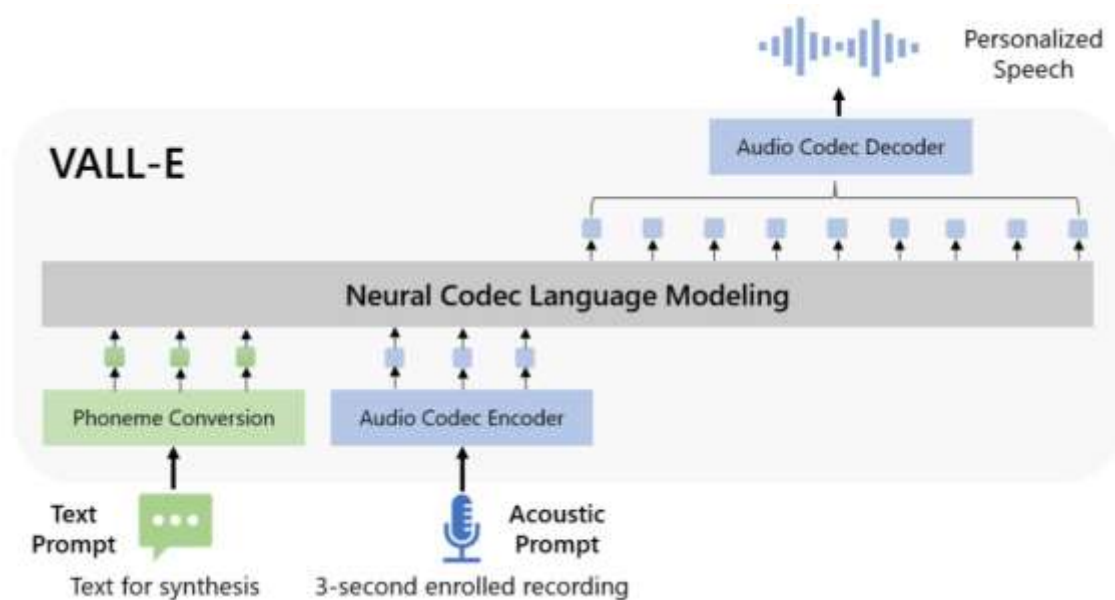
中文说话人



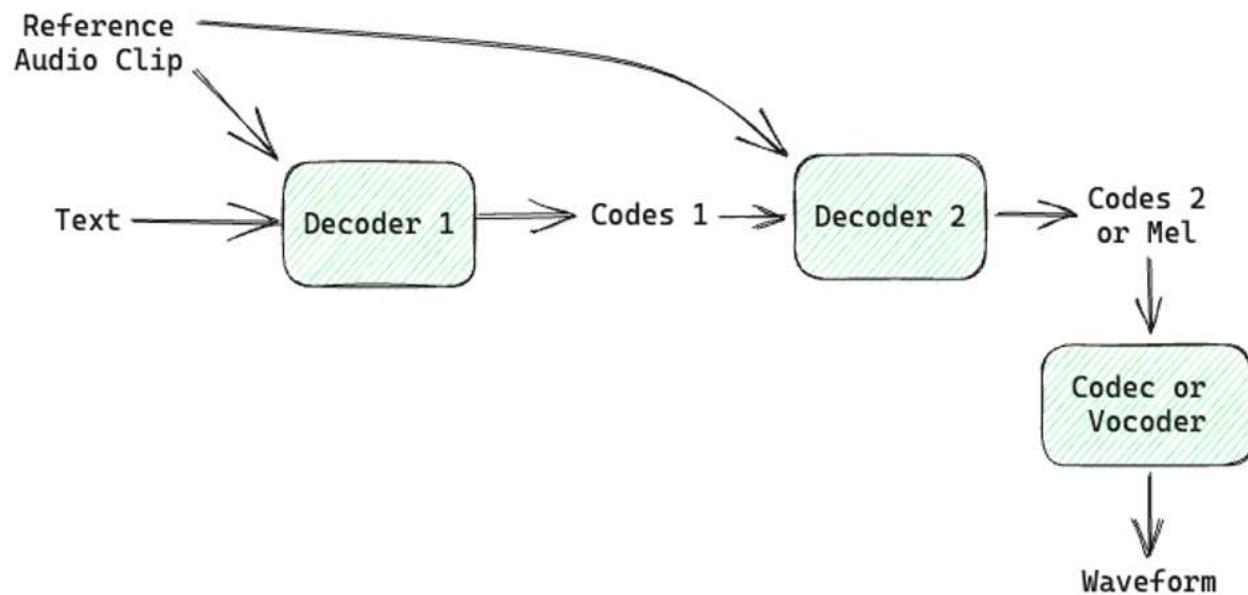
合成效果



# VALL-E (X) 算法框架



# 基于大模型的小米自然语音 TTS



使用小爱默认音色



说话人迁移

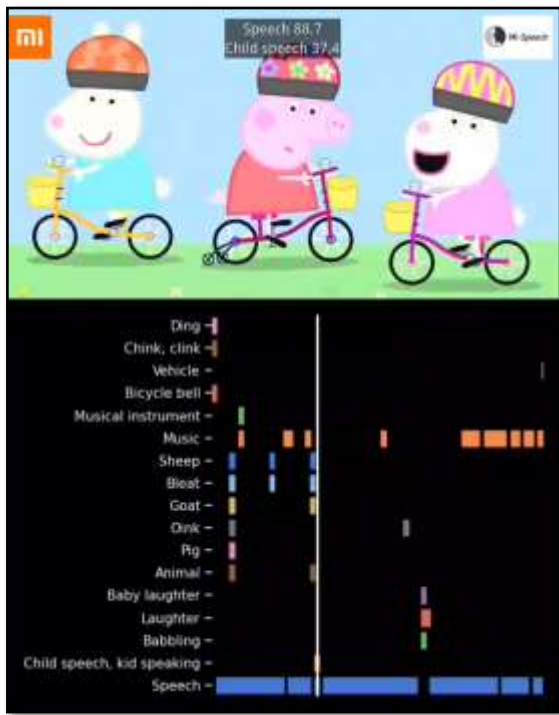


(prompt)





# 小米声音识别技术



目前支持 85 种声音事件



# 大模型时代的语音理解



▶ 0:00 / 0:06 ———▶ 🔊 ⋮

请描述这段声音的主要内容？

▶ 0:00 / 0:04 ———▶ 🔊 ⋮

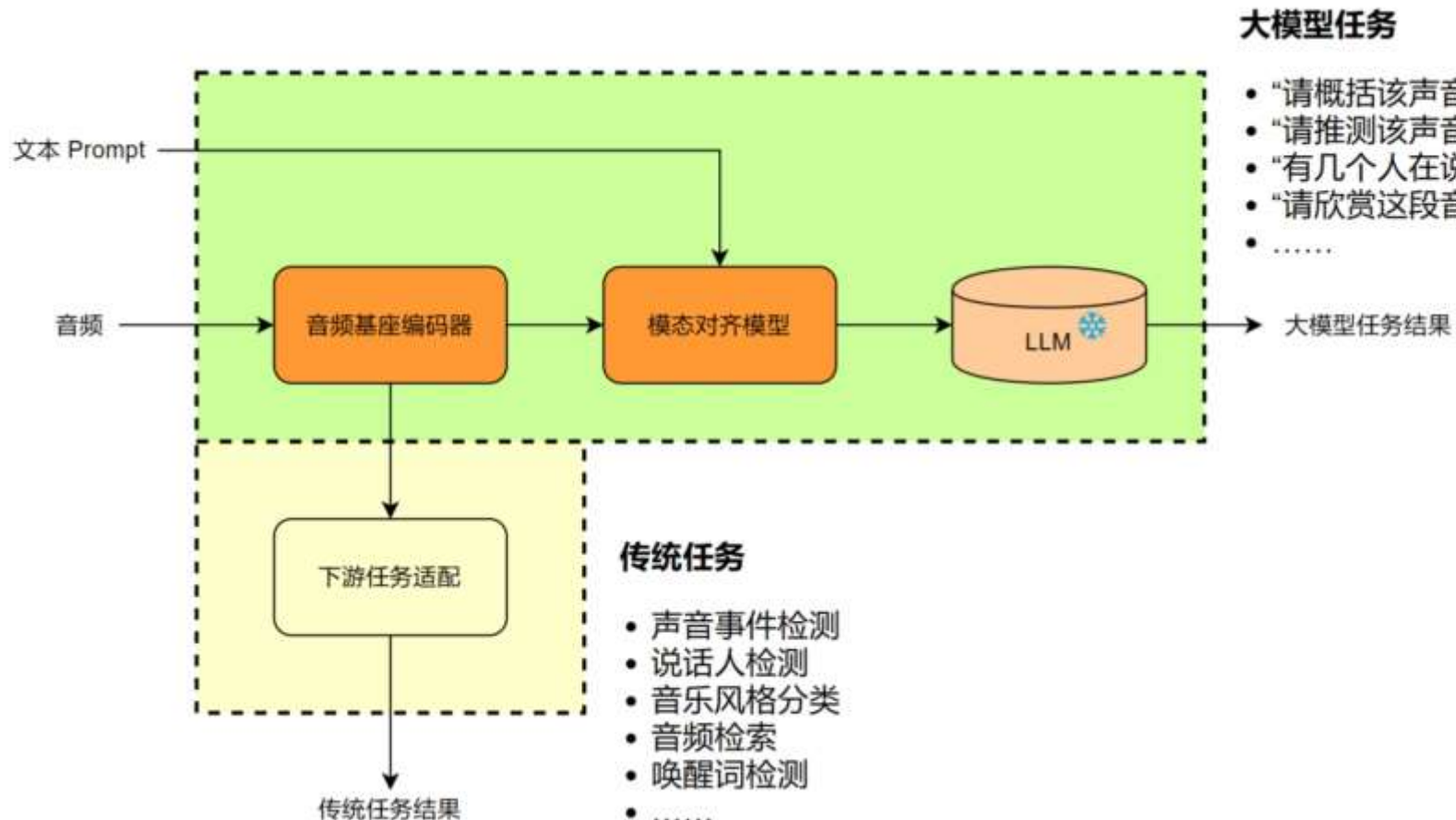
他们在讨论什么问题？结论是什么？

▶ 0:00 / 0:05 ———▶ 🔊 ⋮

现场的气氛怎样？



# 我们的算法框架



- 大模型任务**
- “请概括该声音的内容？”
  - “请推测该声音发生的地点？”
  - “有几个人在说话、他们在争论什么问题、情绪如何？”
  - “请欣赏这段音乐，并针对其艺术水平写一段300字的评论。”
  - .....

## 传统任务

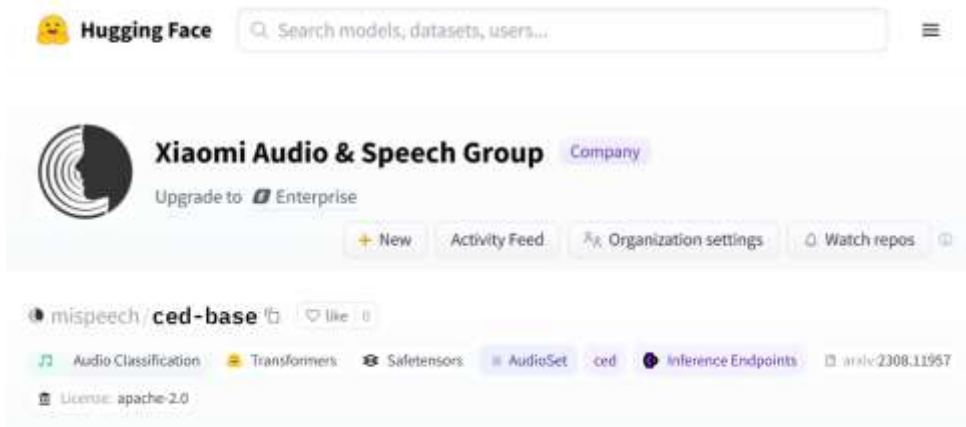
- 声音事件检测
- 说话人检测
- 音乐风格分类
- 音频检索
- 唤醒词检测
- .....

# 基座音频编码器

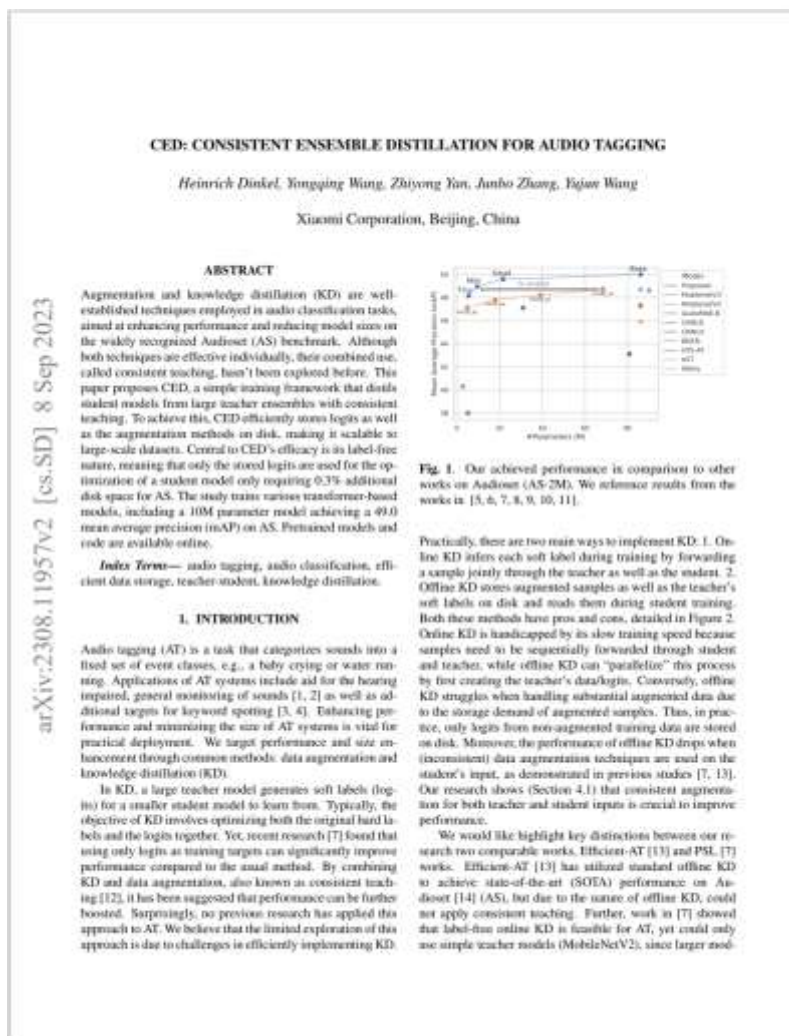
训练数据时长超过30年

参数量超过10亿  
正在探索百亿参数量的模型

独创的一致性集成蒸馏技术  
论文已被 ICASSP 2024 接收



模型开源可下载



## 小米声音识别算法 性能在国际上 排名第一

模型名称	机构	mAP	参数量 (M)
CED	小米	50.0	86
CED-mini	小米	49.0	10
Efficient-AT	JKU	48.7	68
BEATs	微软	48.6	90
Audio-MAE	Meta, CMU	47.3	86
ConcFetS	ANITI 等	47.1	28
HTS-AT	UCSD, 字节跳动	47.1	31
MayaSpec	北大, 浙大	47.1	86
AST	MIT	45.9	86
PANNs-CNN14	萨里大学, 字节跳动	45.1	81
PSL-Mbv2	小米	40.3	3

\*2025年9月12日结果

# 基座音频编码器的多任务应用



Holistic Evaluation of Audio Representations

## HEAR Leaderboard

Model	URL	Submission Date	Redrive	Beijing Opera	CREMA-D	DCASE 2016	ESC-50	FSD50K	GTZAN Genre	GTZAN Music/Speech	Gunshot	Libelcount	Maestro 5h	Mridangam Stroke	Mridangam Tune	NSynth Pitch 50h	NSynth Pitch 5h	Speech commands 5h	Speech commands full	Vocal Imitation	VoxLingua107 top 10
RedRice cod_large	<a href="#">🔗</a>	2023-09-26	0.4835	0.9660	0.0010	0.9218	0.9665	0.6548	0.8995	0.9436	0.9929	0.8785	0.1476	0.9743	0.9655	0.8281	0.6820	0.9683	0.9602	0.2269	0.3852
GURA Fuse Cat H+w+C	<a href="#">🔗</a>	HEAR 2021		0.9660	0.7474	0.8260	0.7335	0.4197	0.8050	0.9282	0.9345	0.6872	0.4413	0.9725	0.9235	0.8852	0.8480	0.9665	0.9677	0.1069	0.7262
RedRice cod_small	<a href="#">🔗</a>	2023-09-26	0.5170	0.9660	0.0004	0.9163	0.9595	0.6433	0.8950	0.9122	0.9345	0.6559	0.1006	0.9692	0.9304	0.7995	0.6020	0.8992	0.8519	0.2192	0.3633
GURA Fuse Cat H+w+C (time)	<a href="#">🔗</a>	HEAR 2021		0.9618	0.7427	0.8260	0.6535	0.3742	0.7089	0.9436	0.9949	0.6591	0.4413	0.9703	0.9243	0.8989	0.8540	0.9508	0.9681	0.2152	0.8286
RedRice cod_mini	<a href="#">🔗</a>	2023-09-26	0.5917	0.9618	0.6520	0.9066	0.9535	0.6388	0.9030	0.9449	0.9601	0.6402	0.08289	0.9626	0.9332	0.7520	0.5500	0.7738	0.8196	0.2037	0.3467
Logitech AI SERLAB BYOL-S	<a href="#">🔗</a>	HEAR 2021	0.5487	0.9533	0.6560	0.6121	0.8050	0.3068	0.8370	0.9385	0.8171	0.7853	0.087860	0.9726	0.9285	0.7116	0.3960	0.9137	0.9481	0.1598	0.4378
CP-JKU mn46_se (x1+1+all_se)	<a href="#">🔗</a>	2023-01-05	0.5027	0.9534	0.6480	0.8130	0.9615	0.6212	0.8980	0.9840	0.9107	0.7253	0.01862	0.9713	0.9047	0.6355	0.3040	0.7767	0.8472	0.2019	0.3179
GURA Fuse Hubert	<a href="#">🔗</a>	HEAR 2021		0.9483	0.7521	0.8259	0.7435	0.4132	0.7960	0.9350	0.9286	0.6834	0.1657	0.9738	0.9096	0.8882	0.3820	0.9468	0.9571	0.1848	0.7140
GURA Cat H+w+C	<a href="#">🔗</a>	HEAR 2021		0.9363	0.6294	0.6808	0.5130	0.3144	0.7220	0.9669	0.8810	0.6390	0.4091	0.9378	0.8591	0.8968	0.8680	0.9270	0.9429	0.1114	0.4389
OpenJ3	<a href="#">🔗</a>	HEAR 2021	0.6040	0.9746	0.5497	0.8328	0.7565	0.4470	0.8790	0.9696	0.9494	0.6414	0.01650	0.9688	0.9669	0.7310	0.5680	0.6796	0.7634	0.07812	0.3313
GURA Fuse war2vec2	<a href="#">🔗</a>	HEAR 2021		0.9449	0.6024	0.7983	0.6050	0.4028	0.7930	0.9532	0.9673	0.6026	0.1110	0.9623	0.8376	0.6059	0.3288	0.9571	0.9687	0.1740	0.7058
GURA Avg H+w+C	<a href="#">🔗</a>	HEAR 2021		0.9448	0.5473	0.6238	0.4480	0.2636	0.7060	0.9365	0.8571	0.6289	0.4599	0.9144	0.8373	0.8963	0.8020	0.8225	0.8808	0.08415	0.3210
CP-JKU PaSST 2b1+mal	<a href="#">🔗</a>	HEAR 2021		0.9669	0.6104	0.9254	0.9475	0.6400	0.8830	0.9769	0.9403	0.6001		0.9650	0.8194	0.5409	0.2500	0.6810	0.6387	0.1820	0.2593
CP-JKU PaSST 2b1	<a href="#">🔗</a>	HEAR 2021		0.9660	0.6104	0.9132	0.9475	0.6409	0.8830	0.9769	0.9405	0.6001		0.9650	0.8194	0.5409	0.2500	0.6810	0.6387	0.1820	0.2593
GURA Avg Hubert+Crpe	<a href="#">🔗</a>	HEAR 2021		0.9323	0.5387	0.6104	0.4365	0.2528	0.6980	0.9455	0.8452	0.6224	0.4634	0.9170	0.8389	0.8966	0.8790	0.7372	0.8233	0.07805	0.2832
RedRice EfficientNet-B2	<a href="#">🔗</a>	HEAR 2021	0.5332	0.8535	0.5746	0.7903	0.8045	0.6071	0.8780	0.9670	0.8780	0.6509	0.0001910	0.9485	0.8432	0.3914	0.1680	0.5734	0.6757	0.1385	0.2551
CP-JKU PaSST base	<a href="#">🔗</a>	HEAR 2021		0.9660	0.6104	0.7879	0.9475	0.6400	0.8830	0.9769	0.9405	0.6001		0.9650	0.8194	0.5409	0.2500	0.6810	0.6387	0.1820	0.2593
GURA Cat war2vec2+crpe	<a href="#">🔗</a>	HEAR 2021		0.9108	0.4598	0.5855	0.3430	0.2341	0.6810	0.9378	0.8333	0.5694	0.4626	0.8977	0.8226	0.8990	0.8680	0.8853	0.9188	0.07640	0.3097
war2vec2	<a href="#">🔗</a>	HEAR 2021		0.9667	0.6562	0.6530	0.5610	0.3417	0.7880	0.9462	0.8482	0.6821	0.02280	0.9432	0.8283	0.6530	0.4039	0.8382	0.8795	0.08606	0.4926



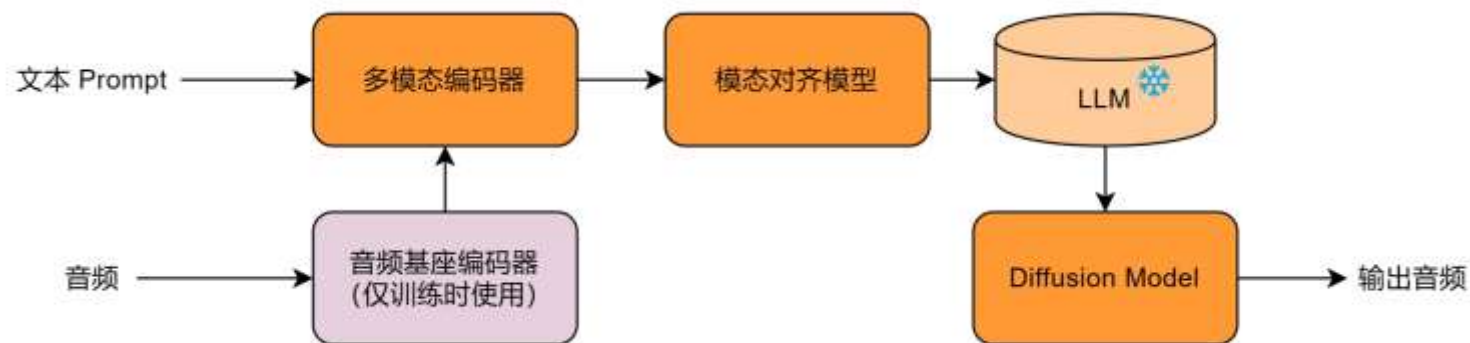
## 声音增强/编辑/生成

已有成果其实已经具备了部分大模型的能力

需要进一步整合



# 基于 Prompt 的声音生成



Prompt:  
风呼啸，树叶沙沙作响

Prompt:  
在雪地中行走

Prompt:  
一辆汽车正在驶过并加速

Prompt:  
猫和狗在打架

## 结语

- 大模型的成功为 AI 研究指明了方向
- 多任务统一学习可以带来真正的理解能力和强大的任务自推广能力
- 各任务的统一、各模态的统一是大势所趋



微信官方公众号：壹佰案例  
关注查看更多年度实践案例