

# 喜马拉雅语音识别技术和应用介绍

用声音分享人类智慧，用声音服务美好生活

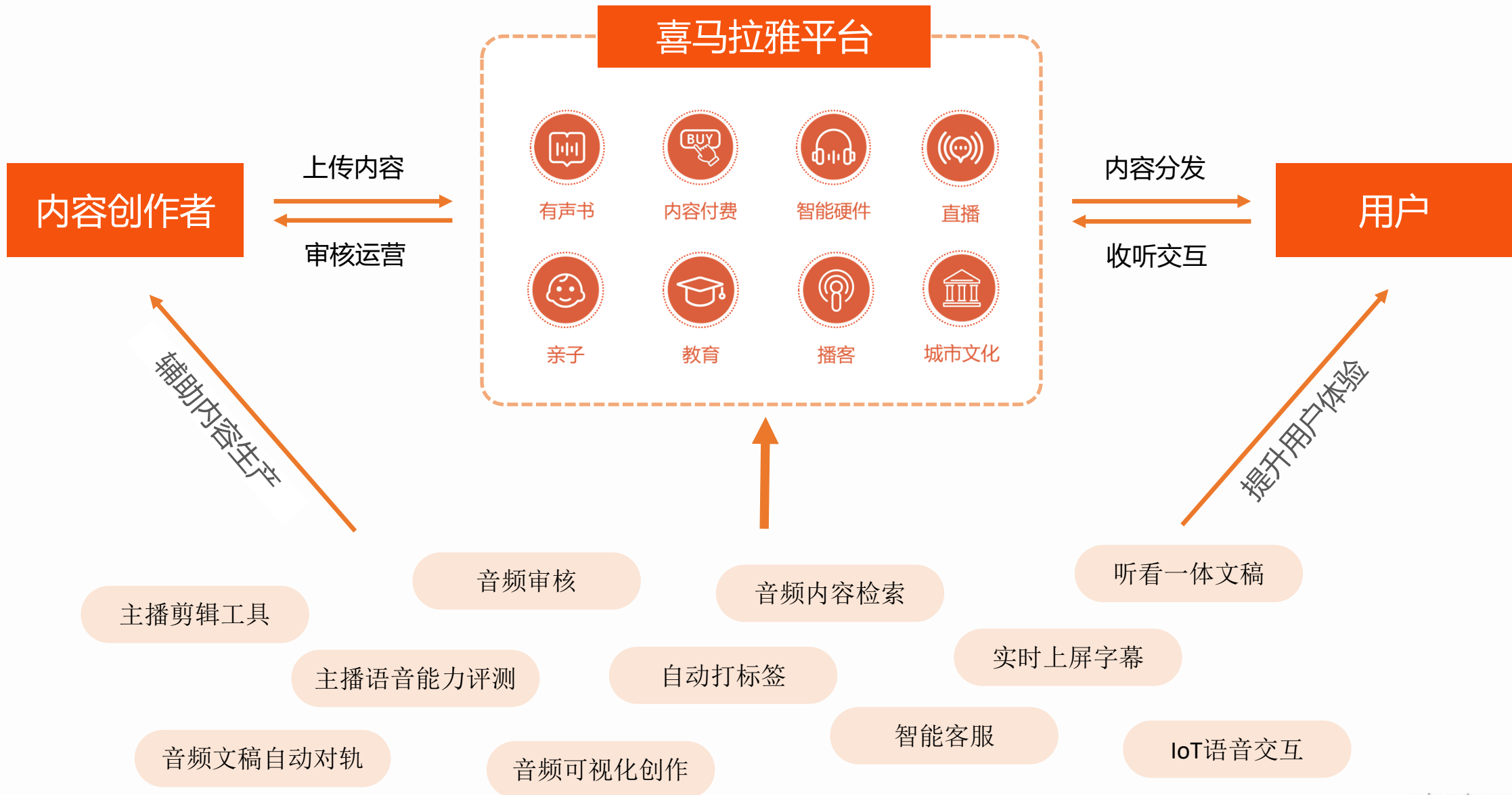


- ▶ 01 | 喜马拉雅ASR应用场景
- ▶ 02 | HiASR算法引擎
- ▶ 03 | HiASR服务架构

主讲人：印晶晶、吕翔

# 01 | 喜马拉雅ASR 应用场景

## 喜马拉雅ASR应用场景



# 喜马拉雅ASR应用场景

声音 AI文稿 评论

这就有了温度计，现在有各种各样的比温度计要高级的多的东西，比如说这个血糖仪呀，**各种各样的各种设备家庭化嘛**，就像商用计算机，像家用电脑，这是所有的高科技产品，它都是这么一个过程嘛。

飞入寻常百姓家，嗯，对，比如说可穿戴设备，它所声称的主张其实是很合理的，就是把体检从一年做一次，做两次体检，变成天天做体检，天天做日日做秒秒做，现在因为它是可以穿戴了嘛。对。

而且呢，结合所谓的大数据技术，对你进行分析，结合所谓社交媒体。

听看一体文稿

《摸金天师》第004章 他是谁？  
摸金天师（活人回避）

00:16 / 00:28

绝大部分人都喜欢叫这里神棍一条街。

付了钱之后，我拖着疲惫的身子下了车。

**也不知为什么，忽然一下子这么疲惫。**

只能在心里安慰自己，是被这一晚上曾出不穷的怪事儿给折腾的筋疲力尽了。

这时候街道上自然是没什么行人。

本文稿由AI生成，可单击修正

生成的片段可在 账号-“咪嚒笔记”中查看

支持60秒

分享 写笔记 字幕视频

音频可视化创作

测试结果

第一题：读单音节字词

女	虐	逛	窘	灭	快	洽	水	润	滑
内	早	寸	河	穷	二	待	老	凑	约
林	酿	听	星	爪	踹	俊	唐	纷	朝

第二题：读多音节词语

诊断	蓬勃	一会儿	牲口	容纳
篡改	小盆儿	蜷缩	佛寺	面条儿
脊髓	班子	然而	语重心长	

第三题：朗读短文

我们知道，水是生物的重要组成部分，许多动物组织的含水量在百分之八十以上，而一些海洋生物的含水量高达百分之九十五。水是新陈代谢的重要媒介，没有它，体内的一系列生理和生物化学反应就无法进行，生命也就停止。

主播语音能力评测

音频切片

只看命中

2021-11-24 14:42:30]

**违禁**

国外的研究所曾经对小学生展开过一。

[2021-11-24 14:42:30-2021-11-24 14:42:40]

**正常**

做早餐实验，他们发现吃了早餐，相比于不吃早餐来上学的孩子，那些吃了早餐的孩子能更快的投入。

[2021-11-24 14:42:40-2021-11-24 14:42:50]

**色情**

学习调皮捣蛋的情况也更上，就说明智是孩子下吃了早饭后意志的一个更强。听到这呢，你应该明白。

音频审核

## 喜马拉雅ASR应用场景

### 特点一：内容生态丰富

- 20+频道，100+细分品类，超2亿条声音
- 泛知识（金融、法律、历史、科技等）：专业性领域多
- 泛娱乐（影视、相声评书、脱口秀等）：背景音、风格差异大
- 信息类（头条、体育、热点资讯等）：时效性要求高
- 其他：英语、小语种、方言、戏曲、音乐等



数据 + 算法

### 特点二：引擎需求多样

- 音频播客：离线长音频转写
- 直播/语音搜索：流式转写
- 智能客服：多说话人识别
- 智能硬件：语音交互、远场识别

### 特点三：调用链路复杂

- 识别的前处理、后处理
- 说话人分割与分离
- 语音评测、副语言信息
- 声音事件检测、文本检测



服务架构



# HiASR (Himalaya ASR) 整体能力建设

## 业务应用

喜马拉雅APP

A+

奇迹文学

主播APP

直播

审核

客服

小雅

.....

## 服务能力

HAS (Himalaya Audio Service)

微服务架构

AI云平台

CPU/GPU部署

## 算法能力

信号处理

- mp3/m4a/aac等音频解码
- 双通道转单通道
- 降噪或语音分离
- 特征提取、数据增强
- 副语言信息分析

长音频切分

- VAD (cnn/tdnn)
- 说话人分割
- 音频事件检测

E2E-ASR

- LAS、CTC、Transformer、Conformer model
- 中文模型、中英文模型
- 推理优化、剪枝、蒸馏
- 结合TLG

Kaldi-ASR

- Chain model
- ngram领域语言模型
- NNLM rescore
- 热词功能

后处理

- 自动加标点
- 段落划分
- ITN
- 置信度
- 敏感词检测

## 数据能力

喜马拉雅数据

业务场景数据

反馈数据

采购数据

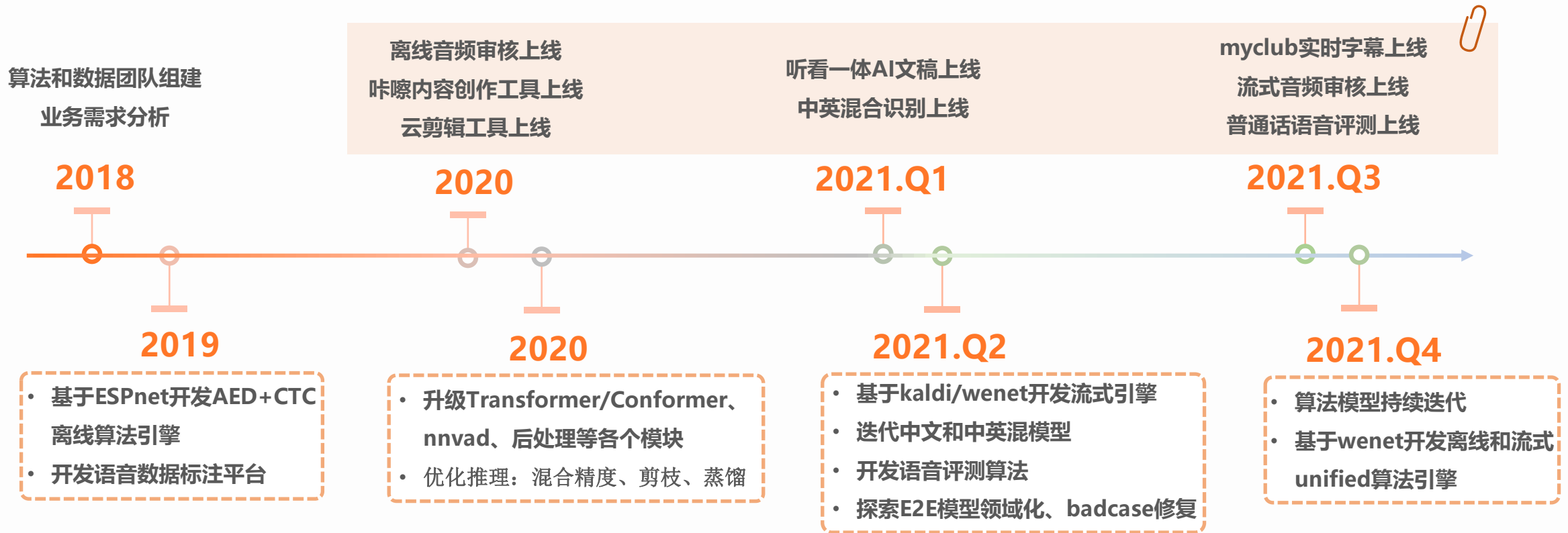
HiCrowd数据标注平台

错误数据反馈链路

# 02 | HiASR 算法引擎

## HiASR算法迭代路径

喜马拉雅以通用转写需求为主 非常适合采用E2E模型





# HiASR-E2E模型

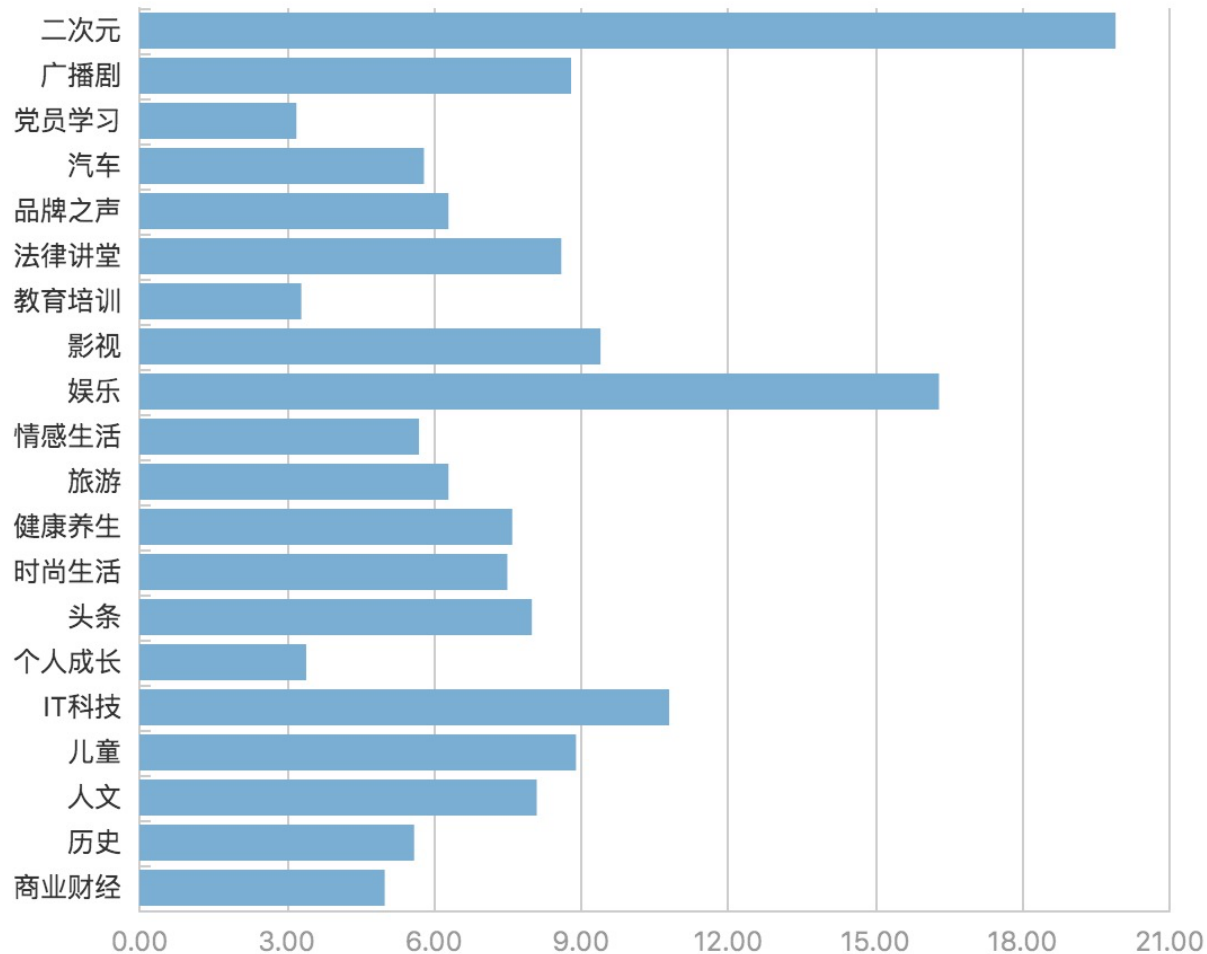
## 数据情况

- 整体：按业务频道需求的优先级以及算法迭代的准确率挑选待处理长音频
- 预处理：音频切分、双模型预识别进行交叉筛选
- 训练数据：数万小时
- 测试数据：20+ 频道130多小时
- 专有名词、人名等特殊词汇，单独收集

## 整体识别率

- 大部分场景准确率>93%
- 部分复杂场景，例如较多背景音乐、太口语化、语速快、口音重、内容领域专业术语多，准确率可能较低，尤其直播娱乐类场景

喜马拉雅主要频道CER (%)



# HiASR-E2E模型

## 框架和模型选择

- 框架：尝试了ESPnet、fairseq和WeNet，前两个不支持流式识别，且需要自行开发推理引擎
- 模型：以AED自回归方案为主，识别率较高
- RNN\_Attention->Transformer->Conformer
- 加入CTC进行联合训练可以有效提升识别率

## 训练策略

- 数据增强：变速、SpecAug、加噪声处理
- 模型迁移和增量训练：基于已经训好的大模型的中间checkpoint进行初始化
- fp16混合精度：可以提升训练和推理速度
- warmup、label smoothing、dropout等有助于收敛和防止过拟合

CER	cn7k model			某第三方
	rnn_att	transformer	conformer	
商业财经	8.1%	5.59%	5.08%	5.4%
历史	7.6%	5.94%	5.46%	6.2%
人文	9.7%	7.94%	7.54%	8.2%
儿童	11.2%	9.28%	8.51%	7.0%
IT科技	13.0%	11.11%	10.48%	8.9%
个人成长	4.9%	3.88%	3.91%	4.9%
头条	7.8%	7.90%	7.41%	6.1%
娱乐	19.9%	17.31%	16.57%	21.2%
教育培训	5.0%	3.51%	3.28%	3.3%
广播剧	10.2%	9.62%	8.20%	8.4%

# HiASR-E2E模型

## 中英混合模型优化

- 混合建模单元：中文汉字、英文子词（BPE）
- 训练语料：大量两种单语数据 + 部分句内code-switching数据
- 采用bi-encoder：可以有效降低整体MER，并且中文CER和英文WER都会降低
  - encoder分别用单语模型初始化
  - 卷积层共享，encoder层部分拆分部分共享

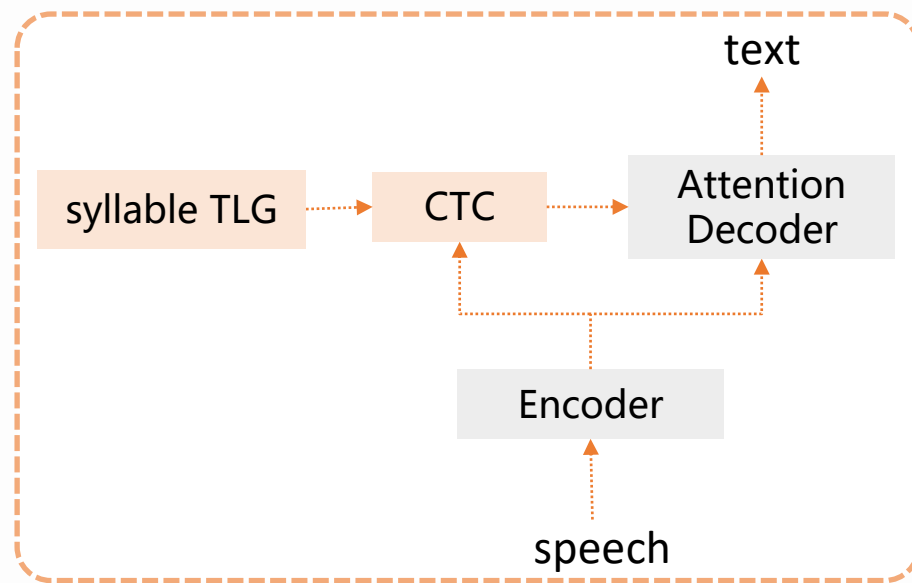
cn15k + en8k model	中英混 MER (CER/WER)
baseline transformer	14.8 (10.49/33.58)
+ cnen0.8k	12.4 (9.18/26.15)
+ cnen0.9k + cn3k	11.0 (8.31/20.93)
+ biencoder	10.4 (7.64/19.86)



## HiASR-E2E模型

### 模型领域化

- shallow fusion: 加入语言模型的方式，小数据集上可以获得一定的效果，在大数据量上效果不明显
- 语音合成+语音转换：可用于修复badcase
- 探索syllable TLG:
  - 基于multi-cn数据集进行实验验证，syllable TLG相比char TLG获得CER相对10.7%的提升
  - syllable TLG在喜马测试集进行诗词类数据的领域化，也获得了相对11.6%的提升



multi-cn	aishell	aidatang	magicdata	thchs	avg
conformer	4.85%	4.75%	3.08%	13.90%	6.65%
+ char TLG	4.26%	3.91%	2.44%	14.43%	6.26%
+ syllable TLG	4.22%	3.72%	2.31%	12.12%	5.59%

# HiASR-推理引擎

## 模型压缩

- 模型剪枝：利用剪枝削减模型的参数量，减少显存占用和推理耗时，其中CNN部分效果较好
- 模型蒸馏：用conformer蒸馏transformer，最终达到与conformer接近的性能，而推理速度更快，但训练时间也会比较长

$$L = \alpha \cdot CE(y, p) + \beta \cdot KL(p, q) \quad q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

cn7k model	params	aishell_test	cctv_news
conformer	93M	3.72%	1.62%
transformer	49M	3.94%	1.72%
probs_t2_alpha0.5	49M	5.38%	2.44%
probs_t2_alpha0.8	49M	5.36%	2.47%
probs_t1_alpha0.5	49M	3.77%	1.62%

## 推理优化

- 相对阈值裁剪：进入下个step时，先对新路径做相对裁剪，再全局beam裁剪
- 动态batchsize：多个长音频分句后，根据句子时长排序重新组batch，充分利用显存
- 单GPU多进程部署：提升显卡利用率
- python/C++混合开发：line\_profile工具分析，将推理中耗时较多的模块用C++改写
- 基于TensorRT/CUDA：模型网络改写

离线 rtf	fbank	vad	transf.	单进程 all	多进程 all
2080ti	2332	550	783	290	>400
3090	2636	763	882	360	>1000

# HiASR-推理引擎

## 原始方案

### E2E离线引擎

- transformer model, 识别率高
- 无法支持流式识别
- GPU推理: 强调高吞吐量
- 优化后2080ti rtf>400, 3090 rtf>1000
- 少量GPU服务器一天可处理几十万小时
- 热词和领域化需要自己开发

### kaldi流式引擎

- chain model, 支持流式识别
- 可以额外训练超大规模语言模型
- 方便进行热词更新和领域化定制
- CPU推理: 强调高并发量低延迟
- 实测单核并发>2路, 平均尾包延迟0.014s

## 基于 WeNet Unified 方案

- unified流式和非流式, 可以迭代一套模型复用到不同场景
- 支持多种解码方式, 可以根据准确率或速度要求灵活配置
- 支持GPU推理和CPU推理, 可以根据吞吐量或并发量需求进行部署
- 支持语言模型、热词功能、字级时间戳, 可以实现功能的快速定制
- 具有完善的引擎和服务的框架, 可以实现更加快速的工业化落地

cn15k model	中文 test1	中文 test2	中文 test3	avg
transformer (non-streaming)	6.81%	8.89%	7.96%	7.5%
unified_conformer (non-streaming)	6.52%	7.96%	7.42%	7.0%
unified_conformer (streaming)	7.10%	8.72%	8.09%	7.7%

# 03 | HiASR 服务架构

## 典型流式语音识别服务架构

### 纯C++ 服务

接口实现

VAD 判断

ASR 识别

后处理

#### Pros

- ✓ 高效率
- ✓ 健壮性

#### Cons

- ✗ 应用库少
- ✗ 开发难度大

### 动态库封装

VAD判断

ASR识别

后处理

JNI调用

Java  
服务

#### Pros

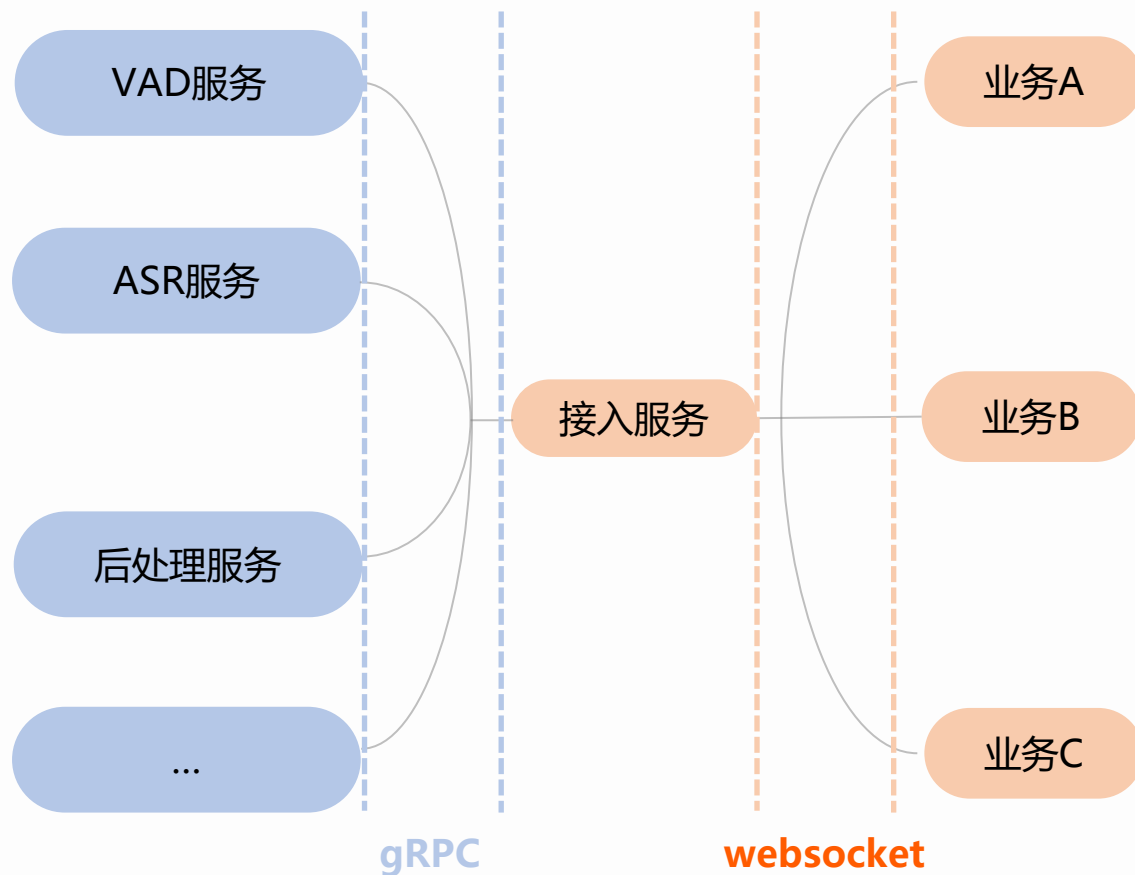
- ✓ 可用库丰富
- ✓ 开发难度/效率均衡

#### Cons

- ✗ JNI调用问题排查
- ✗ 应用库少
- ✗ 开发难度大



# 基于微服务/gRPC 的云原生架构



```
service ASR {  
  rpc Recognize (stream Request) returns (stream Response) {}  
}  
  
message Request {  
  message DecodeConfig {  
    int32 nbest_config = 1;  
    bool continuous_decoding_config = 2;  
  }  
  oneof RequestPayload {  
    DecodeConfig decode_config = 1;  
    bytes audio_data = 2;  
  }  
}  
  
message Response {  
  message OneBest {  
    string sentence = 1;  
    repeated OnePiece wordpieces = 2;  
  }  
  message OnePiece {  
    string word = 1;  
    int32 start = 2;  
    int32 end = 3;  
  }  
  enum Status {  
    ok = 0;  
    failed = 1;  
  }  
  enum Type {  
    server_ready = 0;  
    partial_result = 1;  
    final_result = 2;  
    speech_end = 3;  
  }  
  Status status = 1;  
  Type type = 2;  
  repeated OneBest nbest = 3;  
}
```

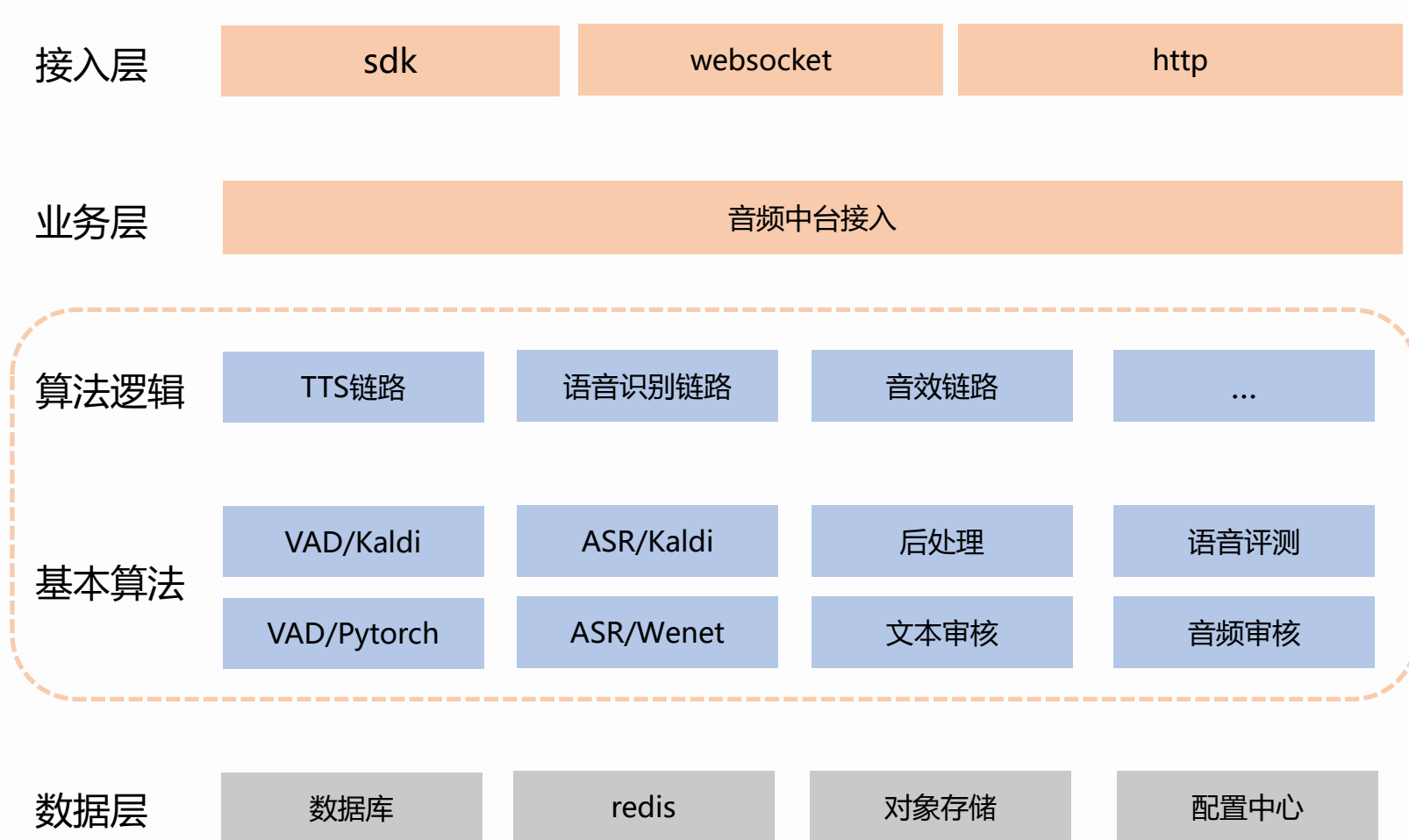
## Pros

- ✓ 模块解耦，独立迭代
- ✓ 多路复用
- ✓ 双向流式
- ✓ 自带队列机制

## Cons

- x 联调难度增大

# 喜马拉雅流式云原生架构



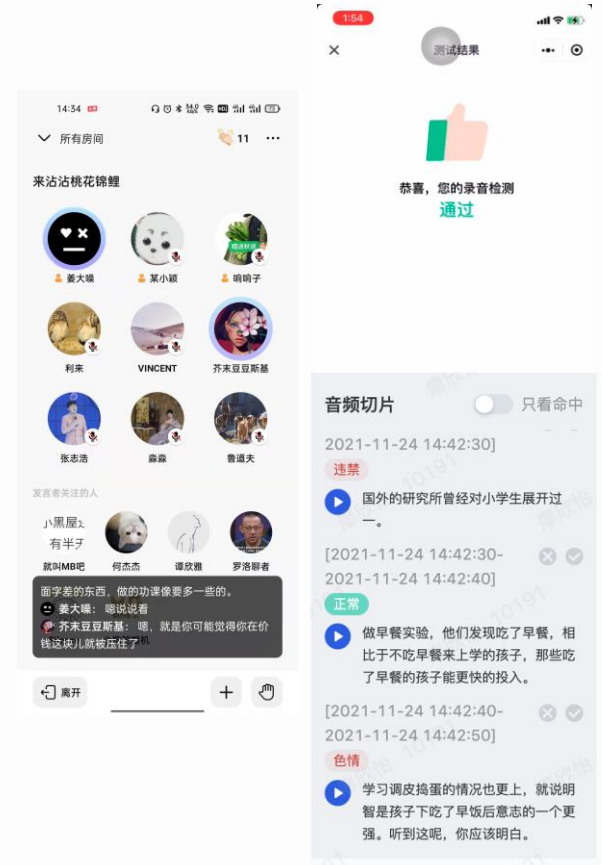
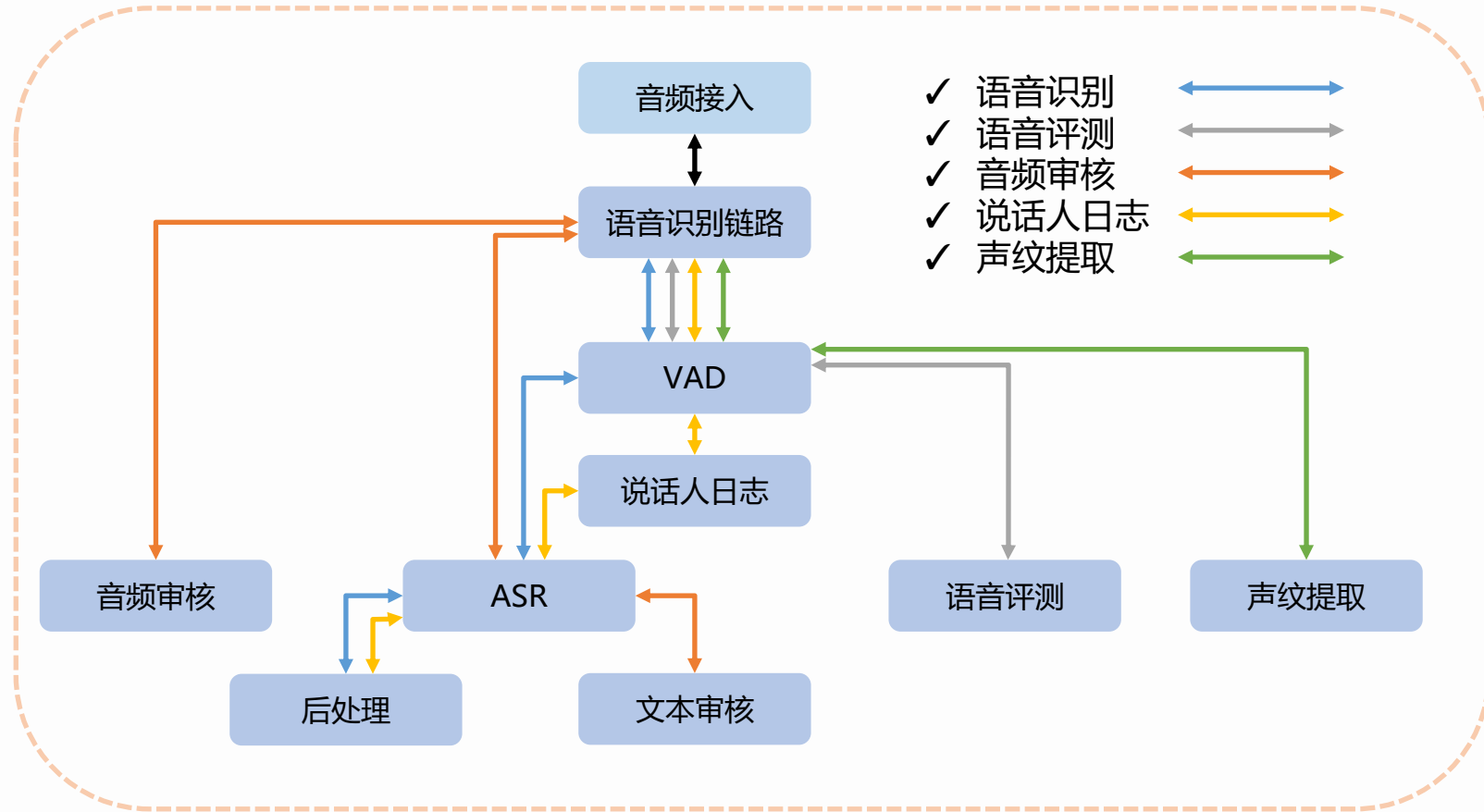
## Pros

- ✓ 鉴权/流量控制
- ✓ 接口统一
- ✓ 算法独立迭代

## Cons

- ✗ 联调难度增加
- ✗ 服务数量增加

# 喜马拉雅流式云原生架构



Wenet service

gRPC server

# client

```
void TorchAsrDecoder::Rescoring() {
    // Do attention rescoring
    Timer timer;
    AttentionRescoring();
    LOG(INFO) << "Rescoring cost latency: " << timer.Elapsed() << "ms.";
    mkl_free_buffers();
}
```

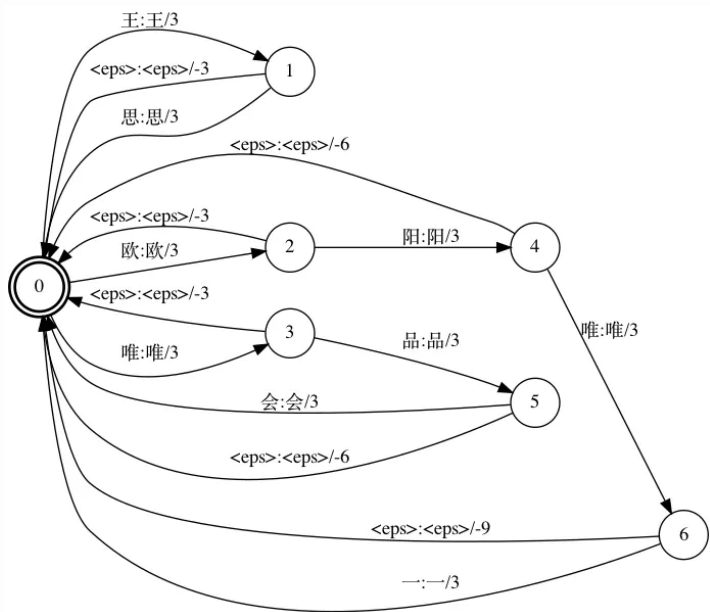
```
==111961== LEAK SUMMARY:
==111961==     definitely lost: 0 bytes in 0 blocks
==111961==     indirectly lost: 0 bytes in 0 blocks
==111961==     possibly lost: 0 bytes in 0 blocks
==111961==     still reachable: 792,039 bytes in 12,527 blocks
==111961==             suppressed: 0 bytes in 0 blocks

==28566== LEAK SUMMARY:
==28566==     definitely lost: 0 bytes in 0 blocks
==28566==     indirectly lost: 0 bytes in 0 blocks
==28566==     possibly lost: 3,840 bytes in 10 blocks
==28566==     still reachable: 303,244,294 bytes in 714,015 blocks
==28566==                   of which reachable via heuristic:
==28566==                       newarray          : 1,088 bytes in 2 blocks
==28566==             suppressed: 0 bytes in 0 blocks
```

每次识别结束释放缓存结果

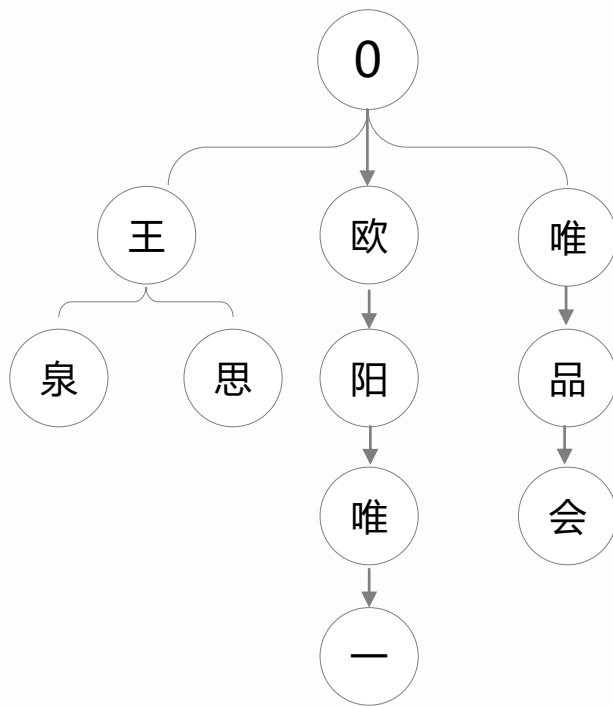
# Wenet 算法服务具体实现

## 基于fst的热词方案



```
if (context_graph_) {
    if (arc.olabel == 0) {
        e_next->val->context_state = tok->context_state;
    } else {
        e_next->val->context_state = context_graph_->GetNextState(
            tok->context_state, arc.olabel, &context_score,
            &is_start_boundary, &is_end_boundary);
        graph_cost -= context_score;
    }
}
```

## 基于字典树的热词方案 (不依赖fst)



```
if (config.trie != NULL)
{
    std::tie(next_trie_, history_hot_words_) = UpdateByHotword(tok, arc.olabel, tot_cost, ac_cost, graph_cost);
}
```

```
{"nbest": [{"sentence": "神智出现交易几乎停滞的情况"}], "status": "ok",
```

# Kaldi/Python 算法服务具体实现

Kaldi/Python service  
gRPC server

Proto request/response

client

## Kaldi VAD

后验概率阈值  
有效语音帧  
有效静音帧

## Kaldi ASR

尾包延迟低于wenet  
领域语言模型  
领域热词

## Kaldi 评测

限定解码空间  
考虑跳读/漏读  
考虑最优/次优phone区分性

```
export PKG_CONFIG_PATH=/usr/local/lib/pkgconfig/  
LDLIBS += -Wl,--as-needed -L/usr/local/lib -lgrpc++_reflection \  
'pkg-config --libs grpc++ libglog protobuf'
```

```
message Recognize_Request{  
  string biz_id = 1;  
  string session_id = 2;  
  string config_json = 3;  
  bytes audio_data = 4;  
}  
  
message Recognize_Response{  
  string response_json = 1;  
  int32 error_code = 2;  
}
```

## Python后处理

正则ITN  
循环神经网络加标点

## Python文本检测

构建敏感词库字典树，耗时80ms->5ms  
bert蒸馏textcnn，速度提升5倍，精度保持99%

## Python声纹

pytorch tdnn->onnx推理，速度提升6倍

## Python说话人日志

cosine+sc+vb

## ■ 优劣分析

经济学十大原理第一条

### 人们面临权衡取舍

无意于强调架构的优劣，合适的才是最好的

#### Pros

- ✓ 算法/代码，自由选择
- ✓ 离线/实时统一，业务逻辑层兼容处理
- ✓ 各个模块单独迭代，增加服务可靠性
- ✓ 提升工作效率

#### Cons

- x 尾包延时增加，100ms以内尚可接受
- x 机器利用效率降低
- x 负载均衡处理，需要L7级负载均衡
- x 无状态处理

# Q & A