



QCon 全球软件开发大会
INTERNATIONAL SOFTWARE
DEVELOPMENT CONFERENCE

BEIJING 2018

小Q机器人诞生之路

演讲者 / 王辉

自我介绍

- **王辉**，现就职于腾讯，社交平台部创新产品技术负责人，专家工程师；
- 近期主要关注人工智能和在线教育等创新产品的研发；
- 热爱技术，热爱跑步。



小Q机器人

AI技术、内容、硬件三大提升。



小Q机器人一代
2010



小Q机器人二代
2017

小Q机器人二代

Tencent 腾讯 | Qrobot

我们不只是AI智能音箱

时尚家庭新成员



小Q机器人2

具备智能概念的机器人

更好的AI智能音箱

小Q除了具备普通AI语音智能音箱功能外，还有更多出色技能

- 强大的社交应用
- 更好的家居中控平台
- 更强大的后台支持
- 海量的后台内容
- 生活上的小助手
- 更好的声音
- 更适用于家居的新成员



更好的语音识别体验

流畅的人机互动 一秒内响应指令 5M远程唤醒



你好，小Q，帮我把电视关掉，放一首轻松一点的歌

更好的声音体验

大音腔高品质音箱，声音饱满流畅，环绕声场和无损音质俱佳
高配置10W大功率音箱，AI智能音箱里鲜有对手

你好小微，我要听《认真的雪》

薛之谦，《认真的雪》



海量内容 想看就看

小Q接入了腾讯云小微智能服务系统QQ音乐、腾讯视频、企鹅FM、小企鹅乐园，独享2年QQ音乐会员服务，正版高品质音乐，畅享听觉盛宴



持续开通中...

自主强大平台提供更好的体验成就品牌的优势

Agenda





聊天机器人

QQ空间聊天机器人



技术成果

- 峰值每日聊天**2000w**次，用户**130w**（业界领先）
- 单用户会话轮次达到**15**次
- 应用多种业界前沿算法，并有多项创新
- 核心AI模型，算法实践结果做到**70%**的效果提升

技术方案

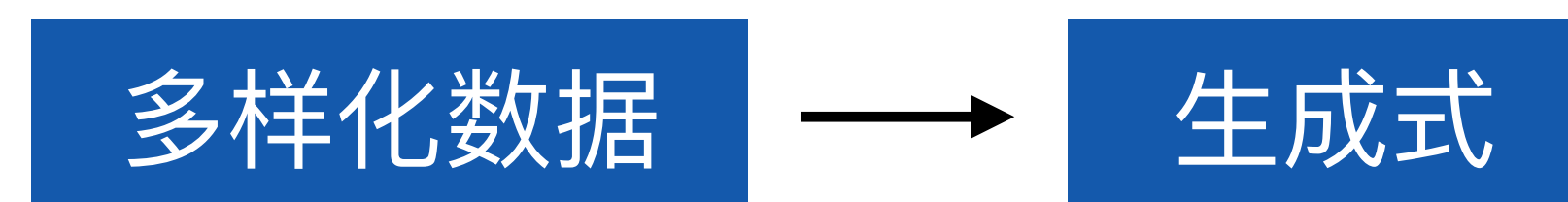
开域



核心场景：

- 聊天
- 推荐语

闭域



核心场景：

- 百科
- 知识问答
- 文字游戏

技术特色——海量数据

数据是人工智能的基石。

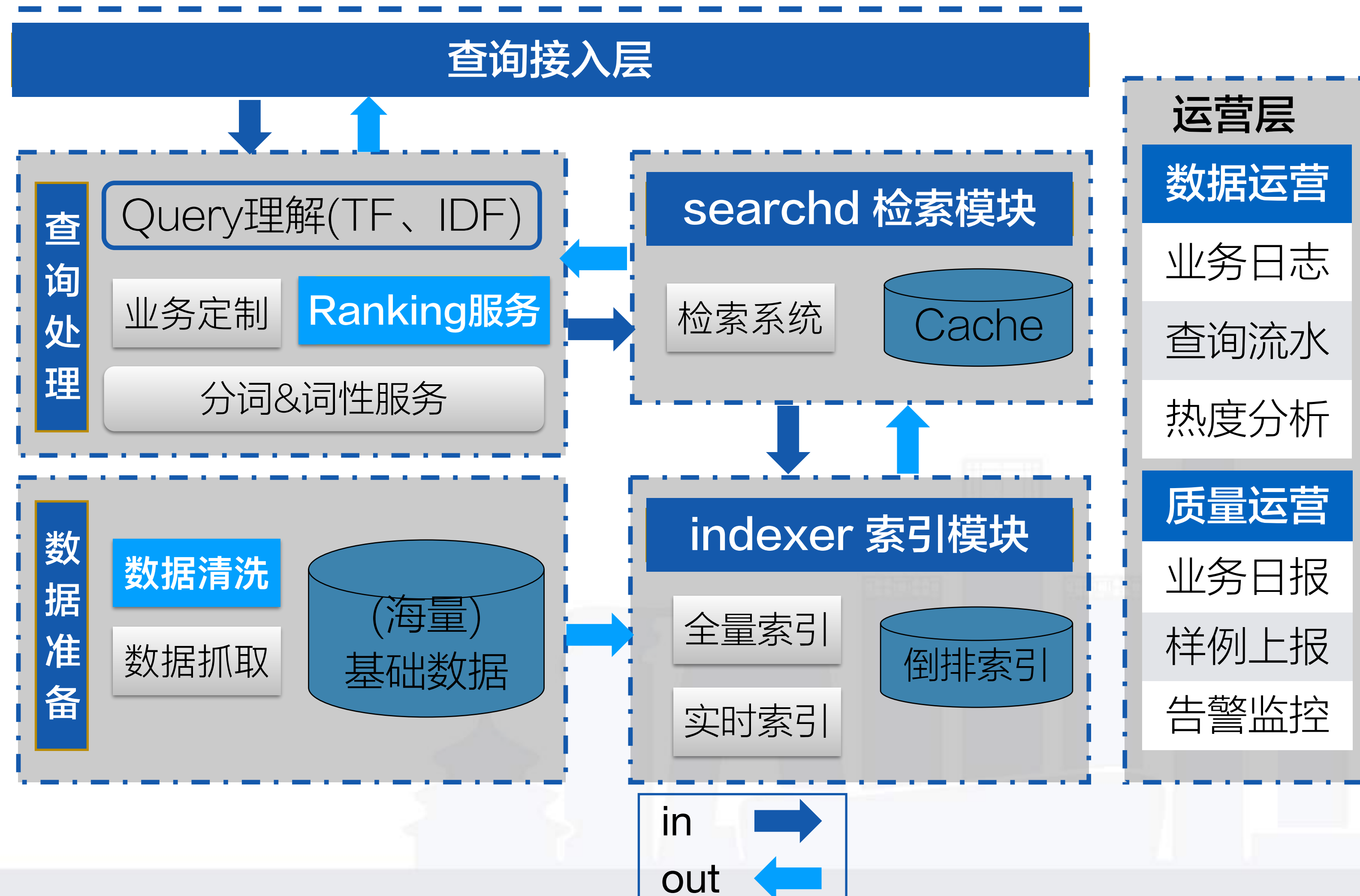


数十亿级公开问答对

技术特色——流动的海量数据

- 自我学习，机器人会“反思”（目前每天可以发现聊天内容约**10%**能继续学习到更好的回答）
- 开域转闭域，深度定制（多个闭域如八卦、新闻、影评等实时知识性数据每日更新，对话过程中的意图领域迁移跟踪）
- 新陈代谢，与时俱进（让聊天随着语料数据得到学习进步，整个数据清洗/检索索引/Rank模型适配流程化）

技术特色——高性能系统



Sphinx

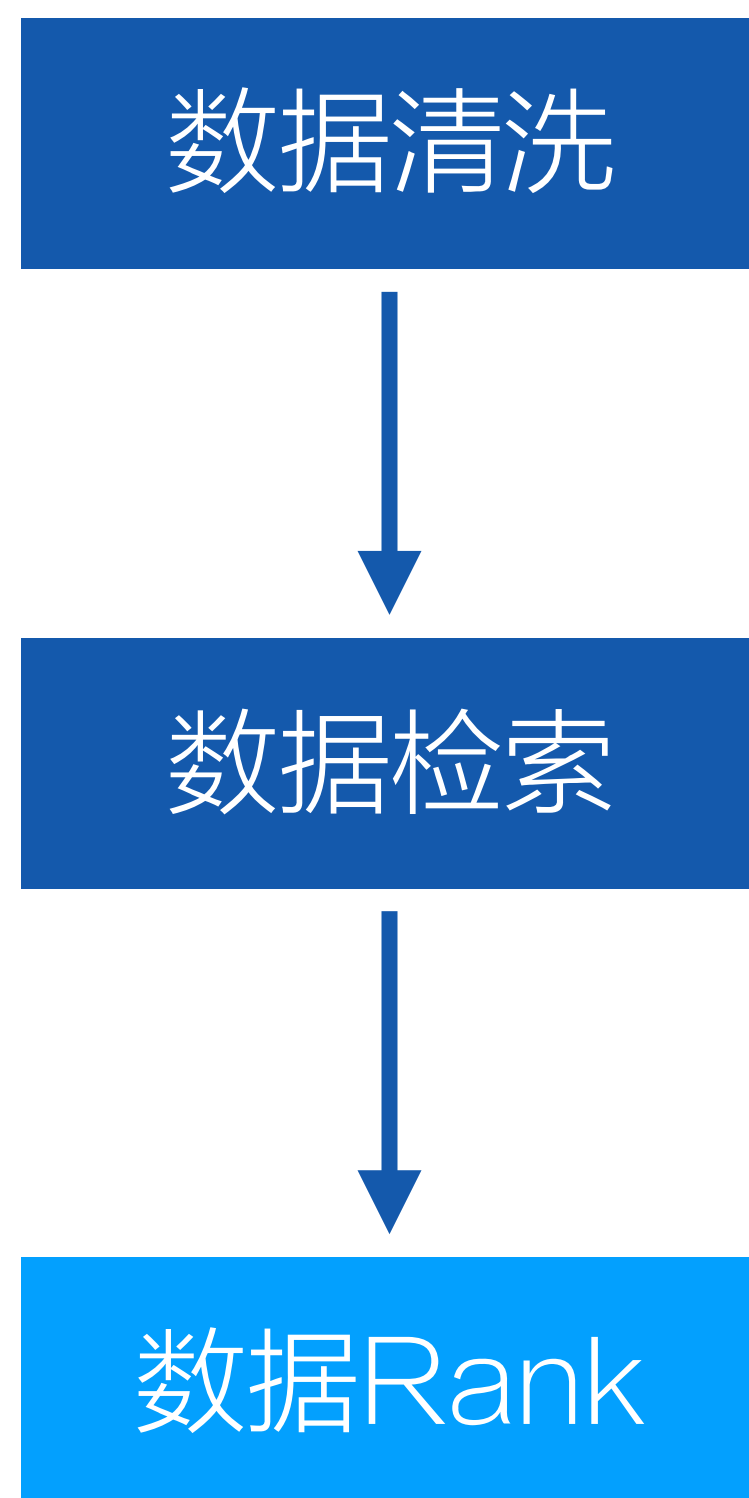
基于Sphinx改造的检索系统

约10亿的中文文档，规模超越知名的coreseek

分布式改造

单机性能优化算法稳定控制在单次计算耗时平均40ms.

技术特色——独家配方Rank模型



关于：最近加班频繁啊 的相关回答：

共检索到匹配答案：20

答	原问题
还在加？	最近加班比较频繁啊
有加班费吗？	最近加班比较频繁啊
还在加？	最近加班比较频繁啊
有加班费吗？	最近加班比较频繁啊
是滴，有点小忙	最近星期天加班比较频繁嘛。
🙄 恩最近加班是正常的，不加班不正常	加班挺频繁的啊
，不是你的错，飞不起来可就是你的不对了	晚上加班太频繁了
还在加？	最近加班比较频繁啊
有加班费吗？	最近加班比较频繁啊

原创baseline基准算法

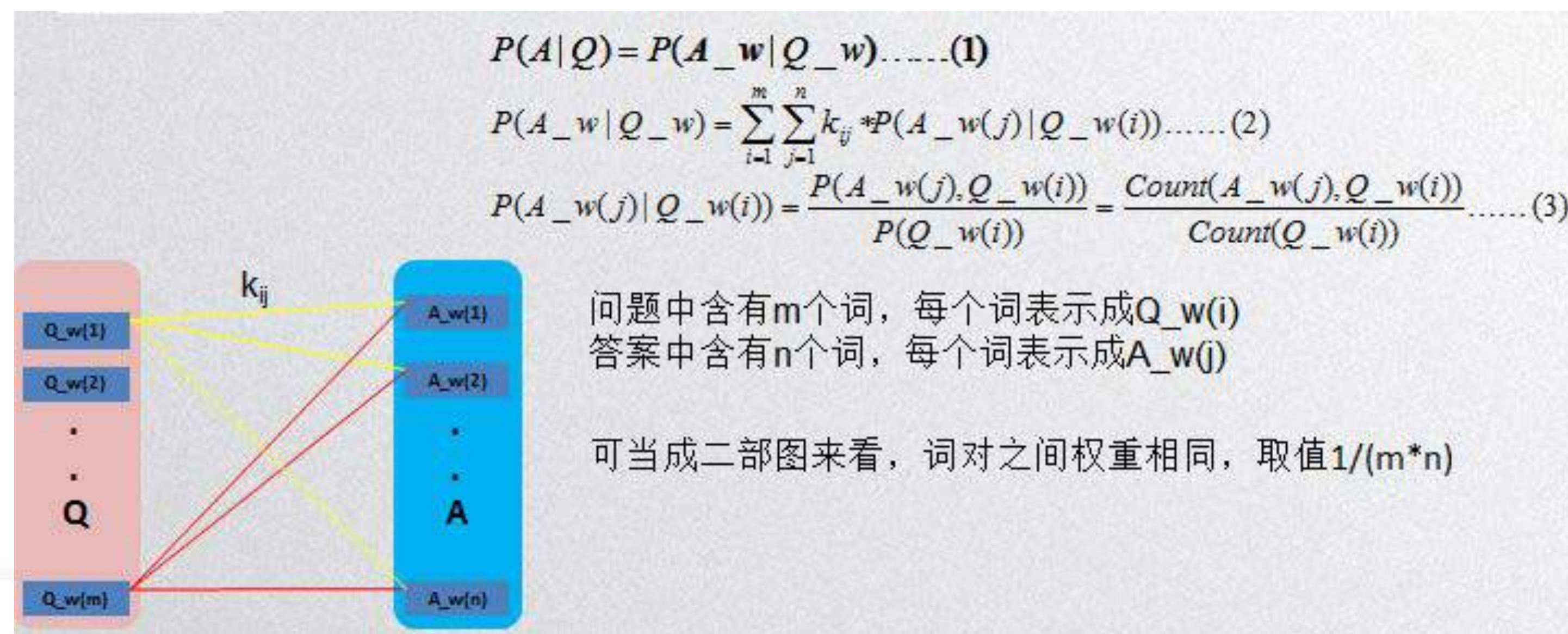
理论基础是条件概率，通过统计海量数据，当人问出一个什么词的条件下，人类回答什么词汇的概率比较高，得出聊天对话中语义上共现的规律。

Q: 今天踢球踢得很累

A1: 休息下

A2: 你还会足球啊?

A3: 下次叫我一起，球场见



算法	P准确率	R召回率	F-measure	效果提升
词共现算法	0.532	0.265	0.354	-----

POS-IDF词向量权重模型

很多词语对确实在QA里共同出现了，但是他们并不是语义上有对话关系，所以需要修正这个条件概率的公式。

Q: 今天踢球踢得很累

A1: 休息下

A2: 你还会足球啊?

A3: 下次叫我一起，球场见

原理:

IDF文档倒频率体现词语关键性

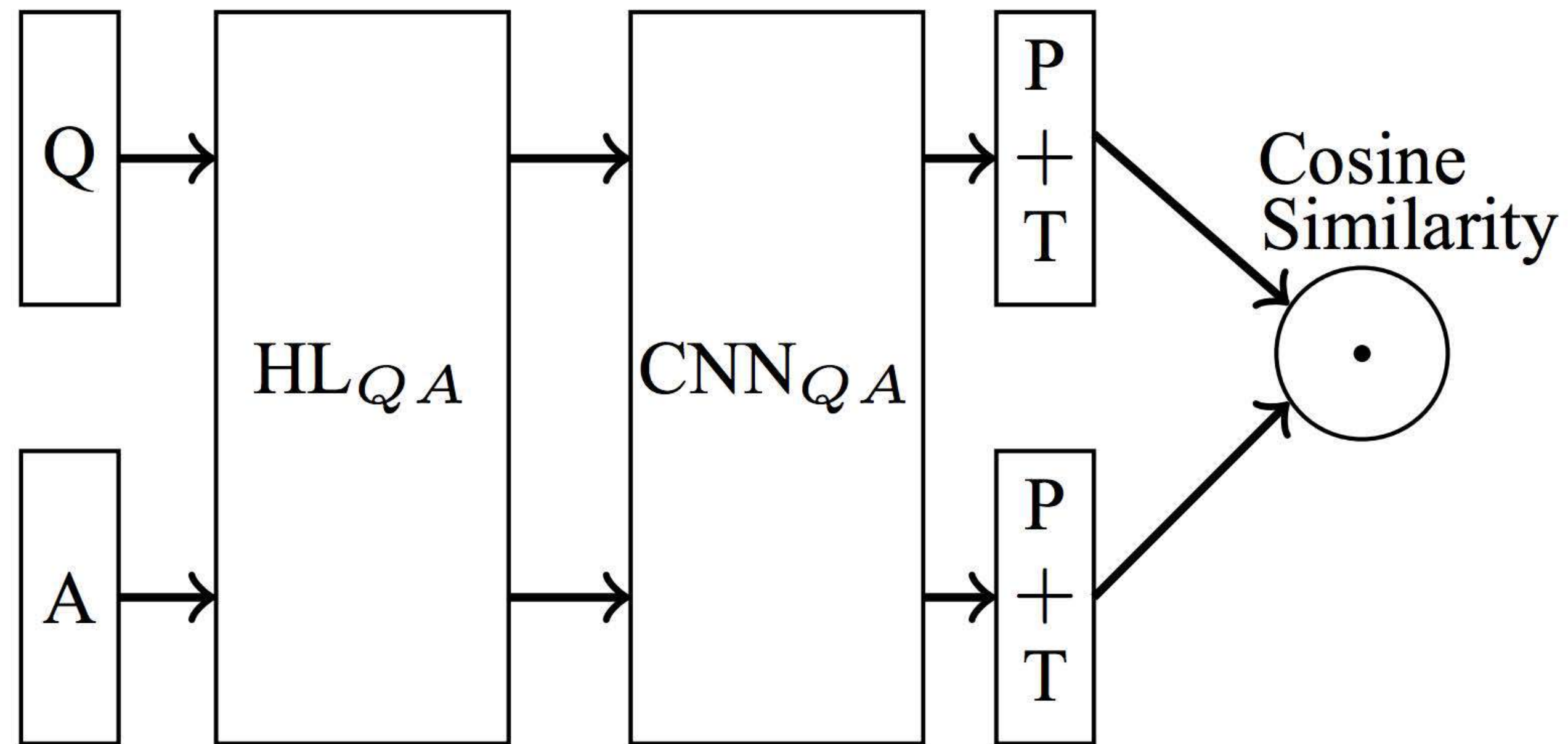
POS词性和句子主干有关

Word2vec的词向量本身具有语义上下文信息

算法	P准确率	R召回率	F-measure	效果提升
POS-IDF-词共现 算法	0.624	0.345	0.445	25%

深度学习CNN

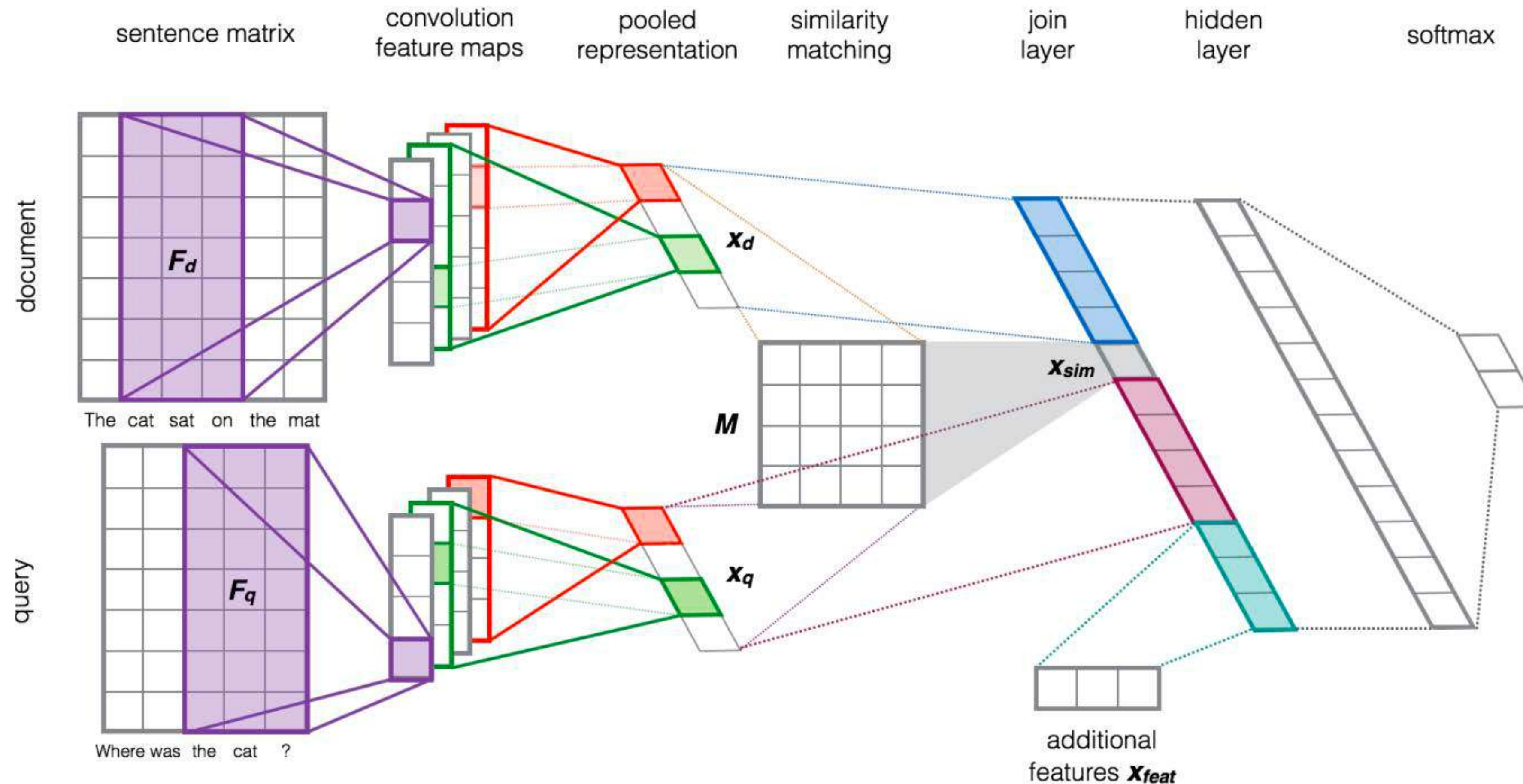
引入深度学习CNN模型，基于Tensorflow开发\GPU训练\评测\部署，QA数据做了拼接共用一个卷积池化的参数，然后求余弦相似度。



Applying Deep Learning to Answer Selection: A Study and An Open Task: <https://arxiv.org/abs/1508.01585>

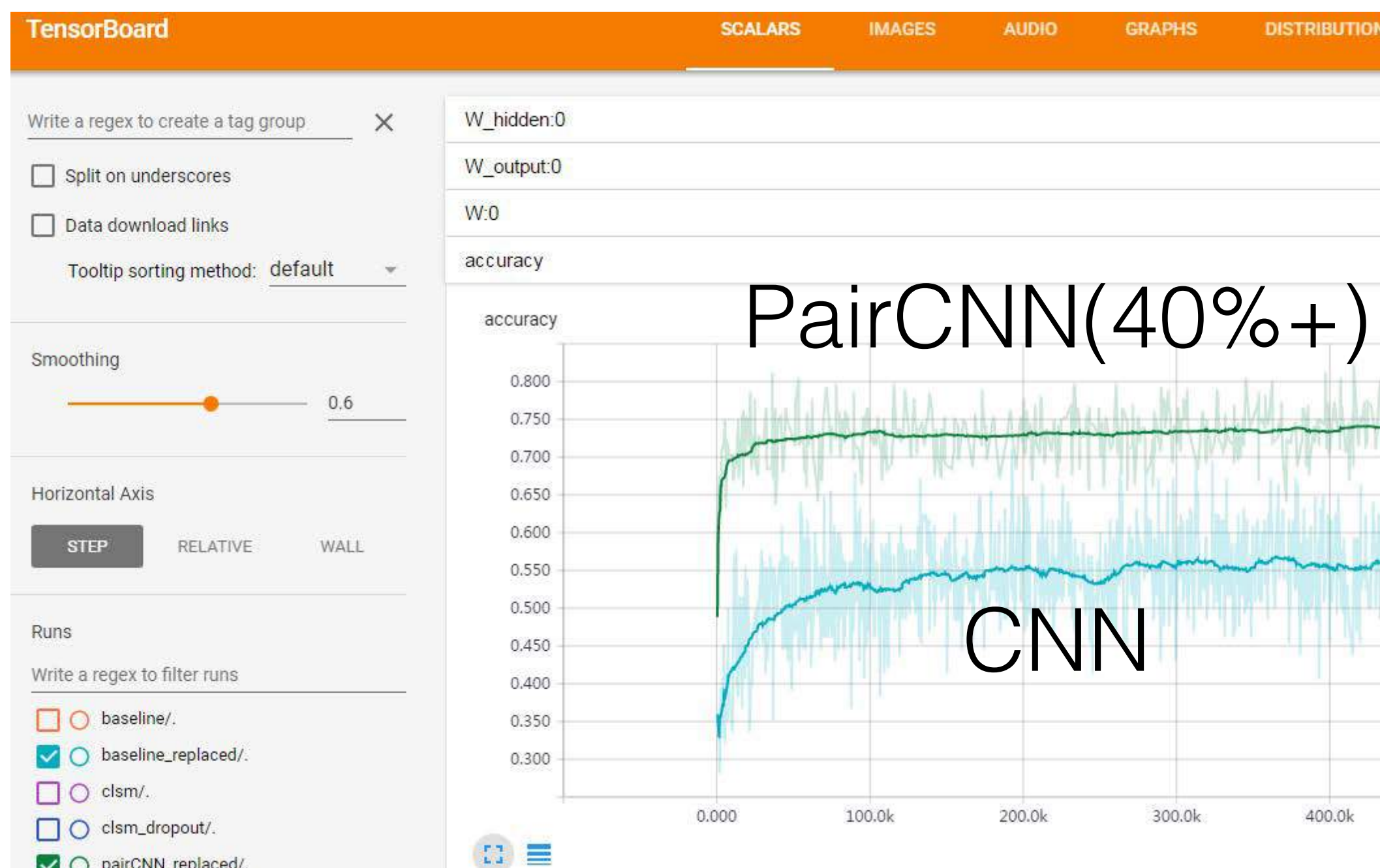
深度学习PairCNN

我们引入PairCNN改进模型对原有的CNN模型对比优化。



Learning to Rank Short Text Pairs with Convolutional Deep Neural Networks: <https://dl.acm.org/citation.cfm?id=2767738>

模型对比



模型对比

CNN

- 1、共享所有卷积池化的参数，且卷积核数量较大，计算量大，对于Q和A的文本的区分度 (diversity)无法分别建模；
- 2、Q和A的交互发生在最后MLP层中，表达相关性的能力有限

PairCNN

- 1、Q和A分别用不同的卷积池化参数进行建模，且卷积核数量小于原有模型，复杂度更低；
- 2、在MLP层之前增加了Q与A相关性抽取的交互，从 high-level的语义向量中进一步抽取了相关性特征，能更加丰富语义表征

算法

P准确率

R召回率

F-measure

效果提升

PairCNN算法

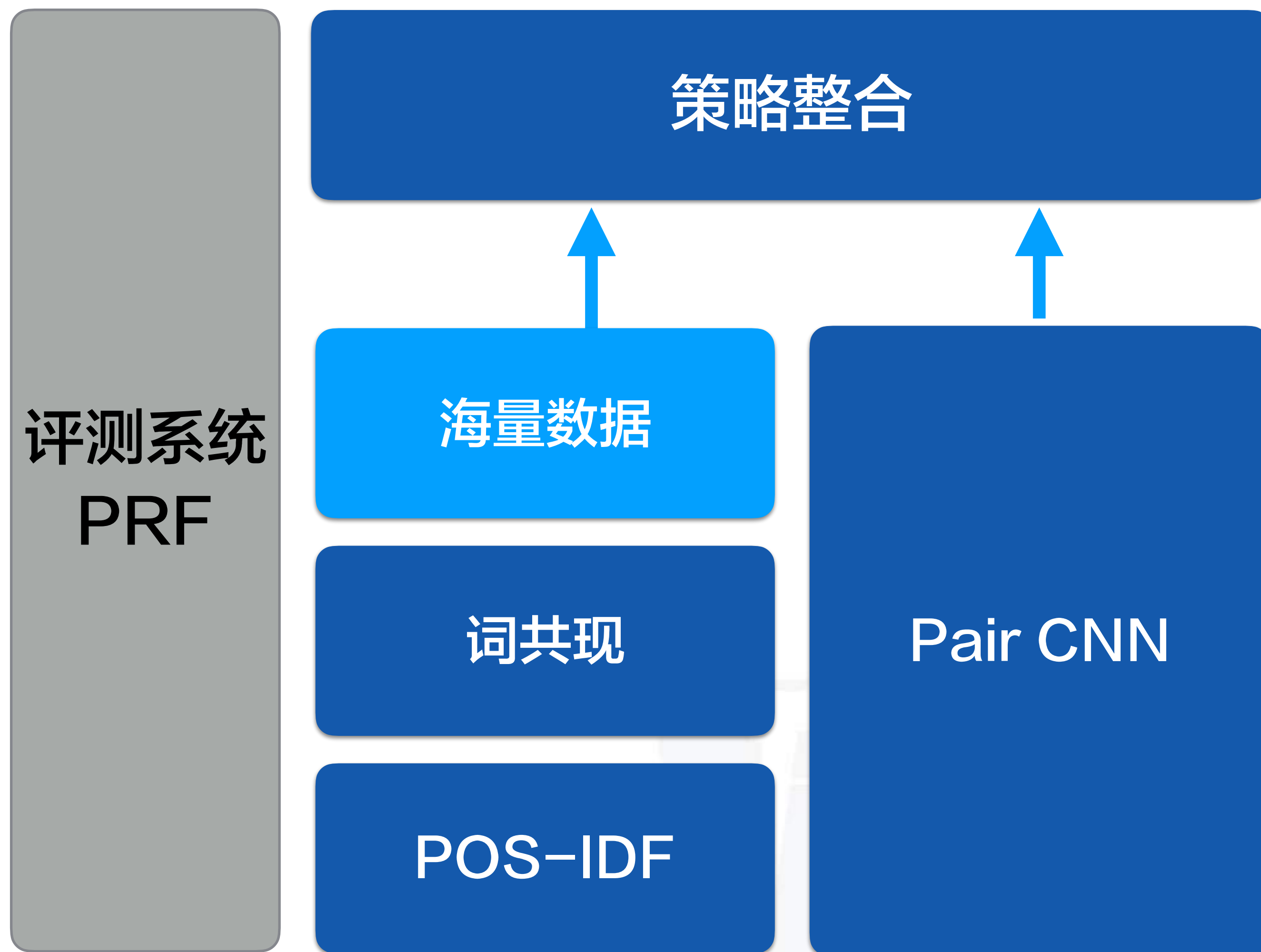
0.816

0.415

0.551

55.6%

技术特色——独家配方Rank模型



- 1、各自评测
- 2、综合配比

技术特色——独家配方Rank模型

综合词共现 / IDF词向量 / PairCNN取得比baseline在F值上约70%的提升。

$$F=2*P*R/(P+R)$$

算法	P准确率	R召回率	F-measure	效果提升
词共现算法	0.532	0.265	0.354	-----
POS-IDF-词共现算法	0.624	0.345	0.445	25%
PairCNN	0.816	0.415	0.551	55.6%
词共现+idf-word2vec	0.809	0.431	0.563	59%
PairCNN+词共现	0.812	0.448	0.577	63%
三种算法结合起来	0.804	0.479	0.601	69.77%

技术特色——文本情感分析

- 基于UGC中文和Emoji表情label混杂得到训练数据
- 利于深度学习LSTM模型进行文本情感分类
- 目前6分类的成功率大于80%，业界第一梯队水平标准了
- 结合宠物的情感分析商业价值也开始显现(KFC &宠物)

只有。。。好吧，保持微笑 😊

两个女汉子勇闯魔都，体育界的老法师纷纷伸出友爱之手，今天是收获的一天 🤔 🤔



Agenda

1 聊天机器人

2 腾讯云小微

3 小Q机器人





腾讯云小微

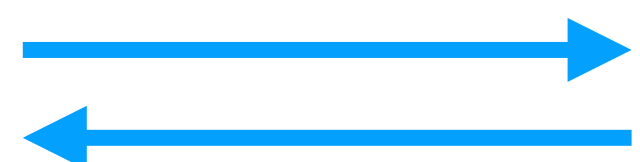


产品定位



小Q机器人

整合+开放+赋能



能力展示

腾讯云小微

"小微，你好！"

小微，是一套腾讯云的智能服务系统，也是一个智能服务开放平台，接入小微的硬件可以快速具备听觉和视觉感知能力，帮助智能硬件厂商实现语音人机互动和音视频服务能力。

在使用小微的时候，只需要说一声“小微”，就可以开始播放音乐和视频、听有声故事和新闻、查询天气、学习英语、与朋友聊聊天、创建任务提醒、设定闹钟时间等，小微还可以和各种智能设备进行交互，用来控制调节灯光空调和电视，小微还能通过图像识别技术认识很多东西，这看起来很酷，是吧！来吧，让我们迎接AI时代的到来，用声音连接物理世界！

云服务
(腾讯云小微)

- 人工智能探索性产品
- 拥有多项AI能力的智能机器人
- 开放成为“腾讯云小微”云服务
- “腾讯云小微”的行业标杆产品

公司战略

人工智能，令人振奋的长远投资。

整合人工智能基础能力，赋能生态系统合作伙伴。——腾讯2017年Q2财报

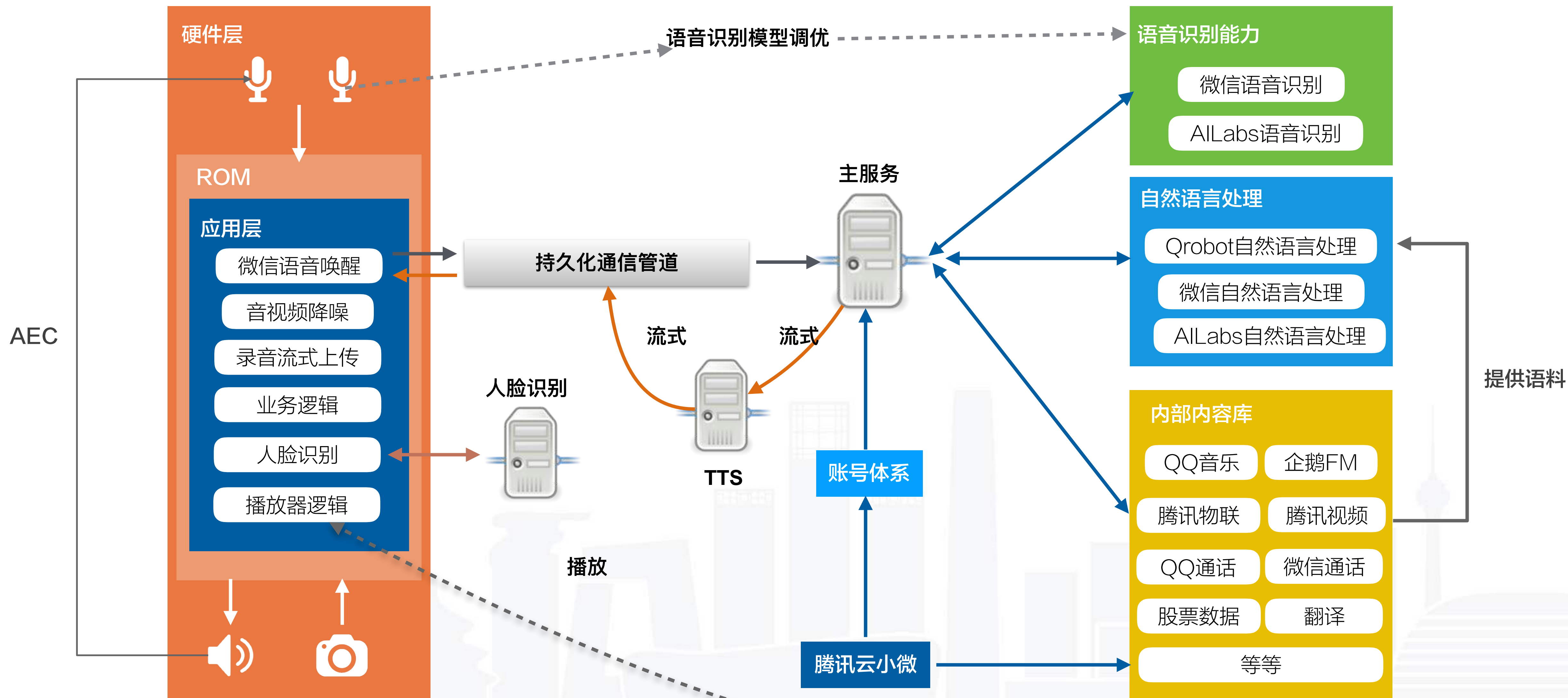


能力整合

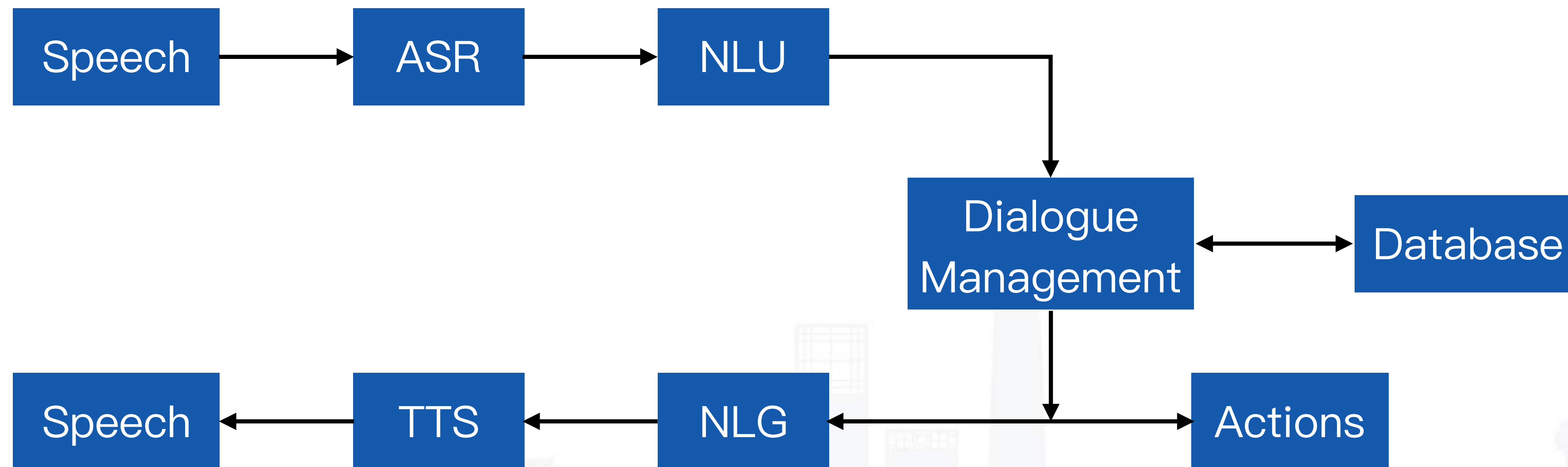
整合腾讯内多项AI能力。



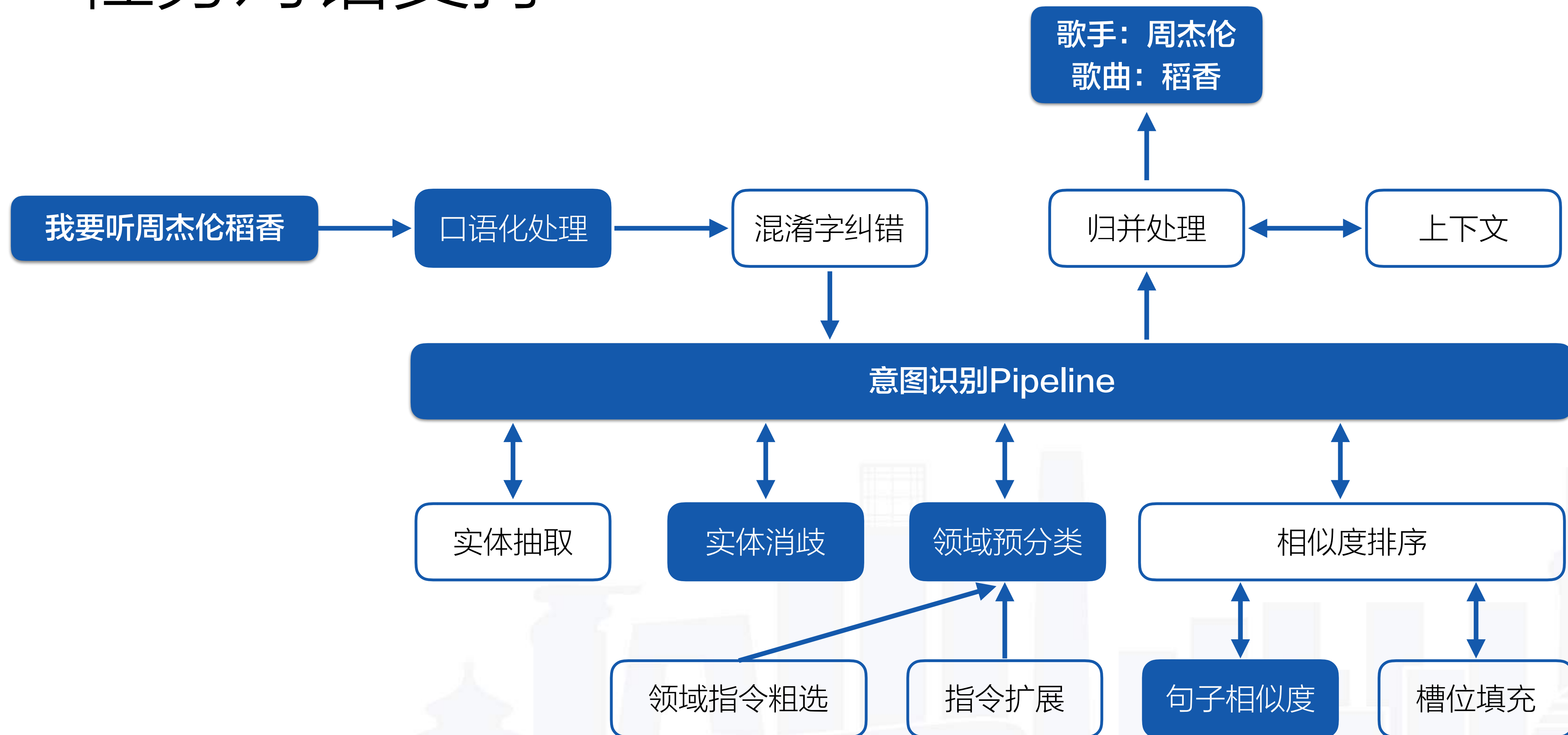
系统架构



任务对话支持

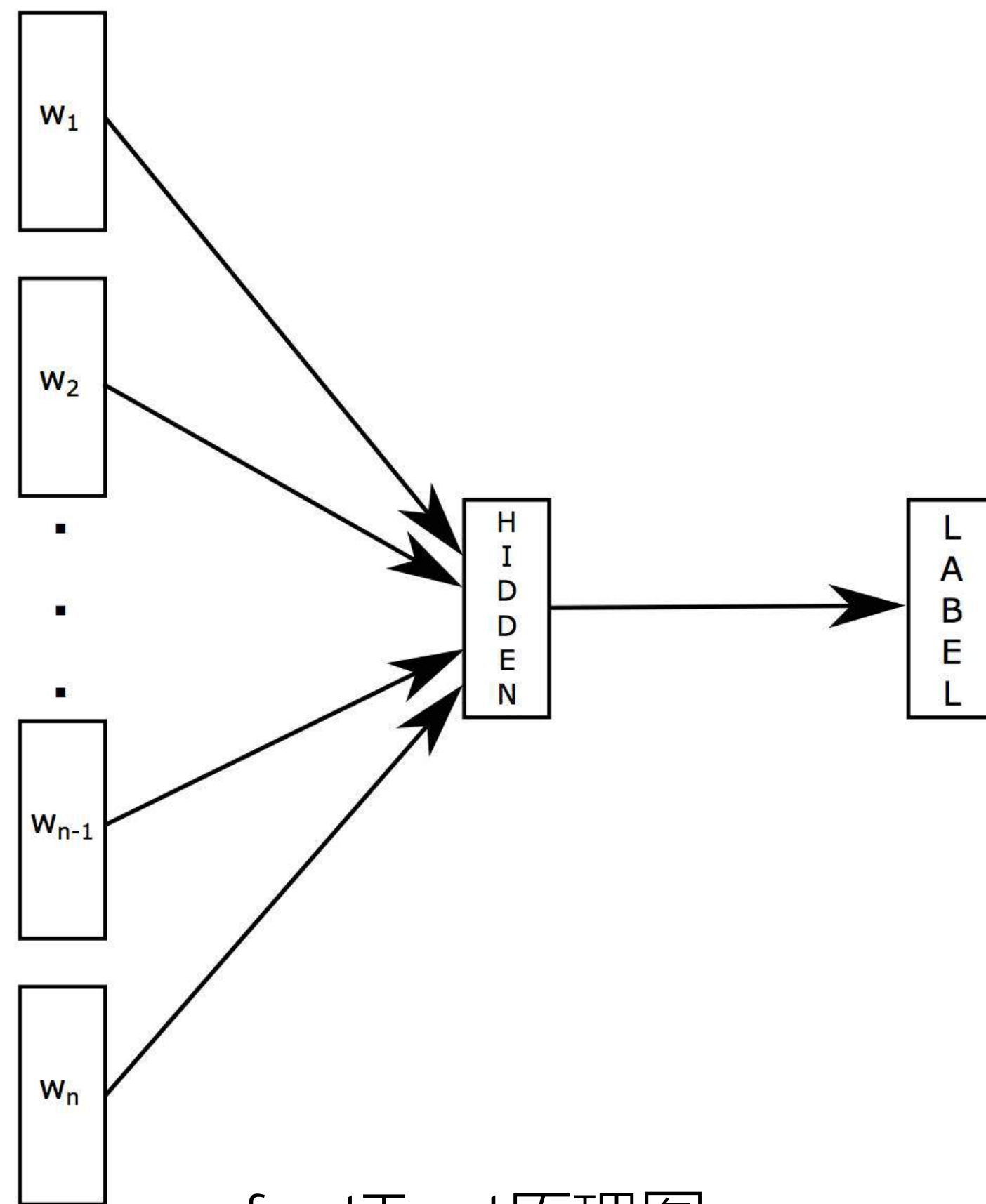


任务对话支持



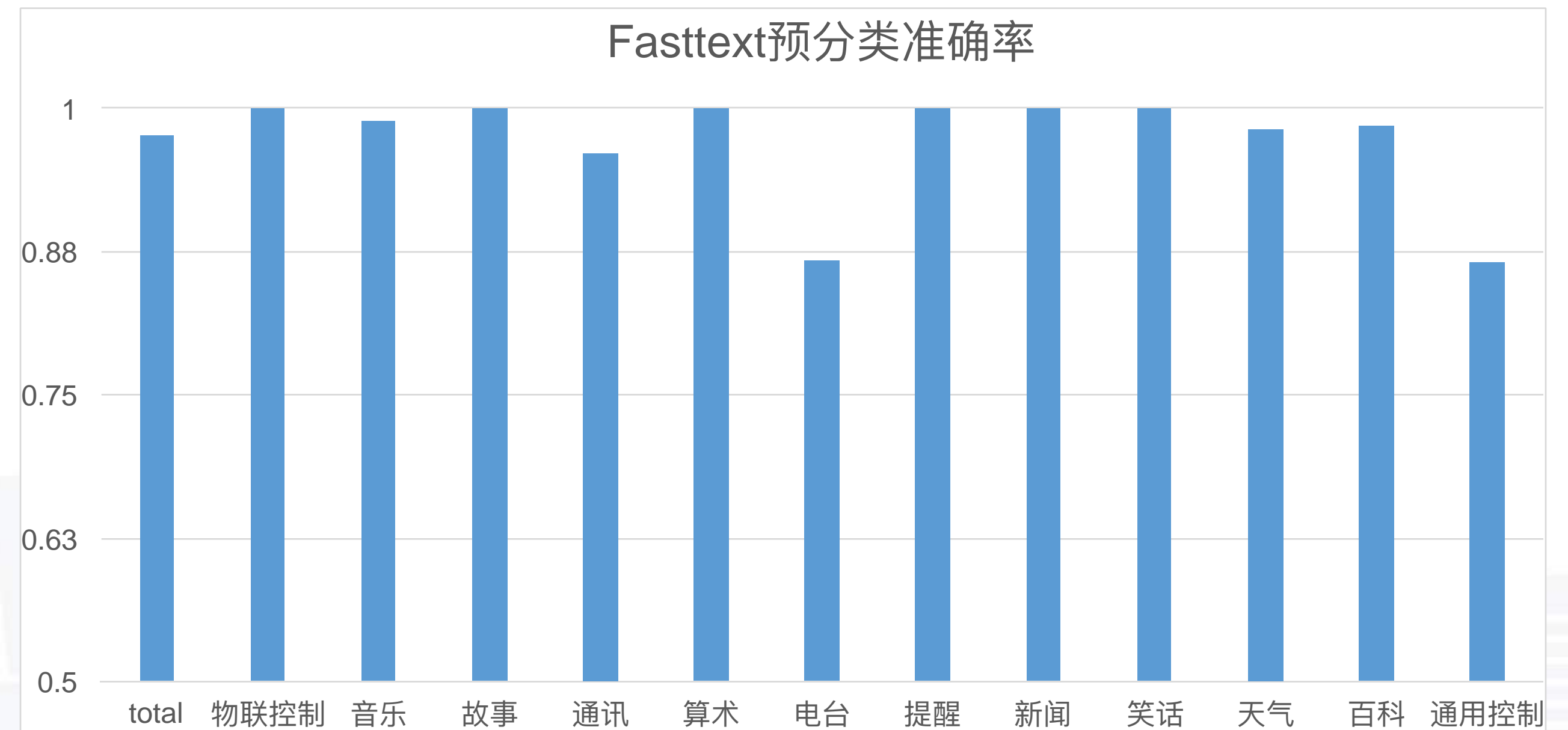
技术特色——领域预分类

提高响应速率策略：fastText预分类模型。



fastText原理图

优点是它是浅层神经网络，计算简单；
在top3的领域上选中目标准确率达到了97.6%。



平均准确率：0.976

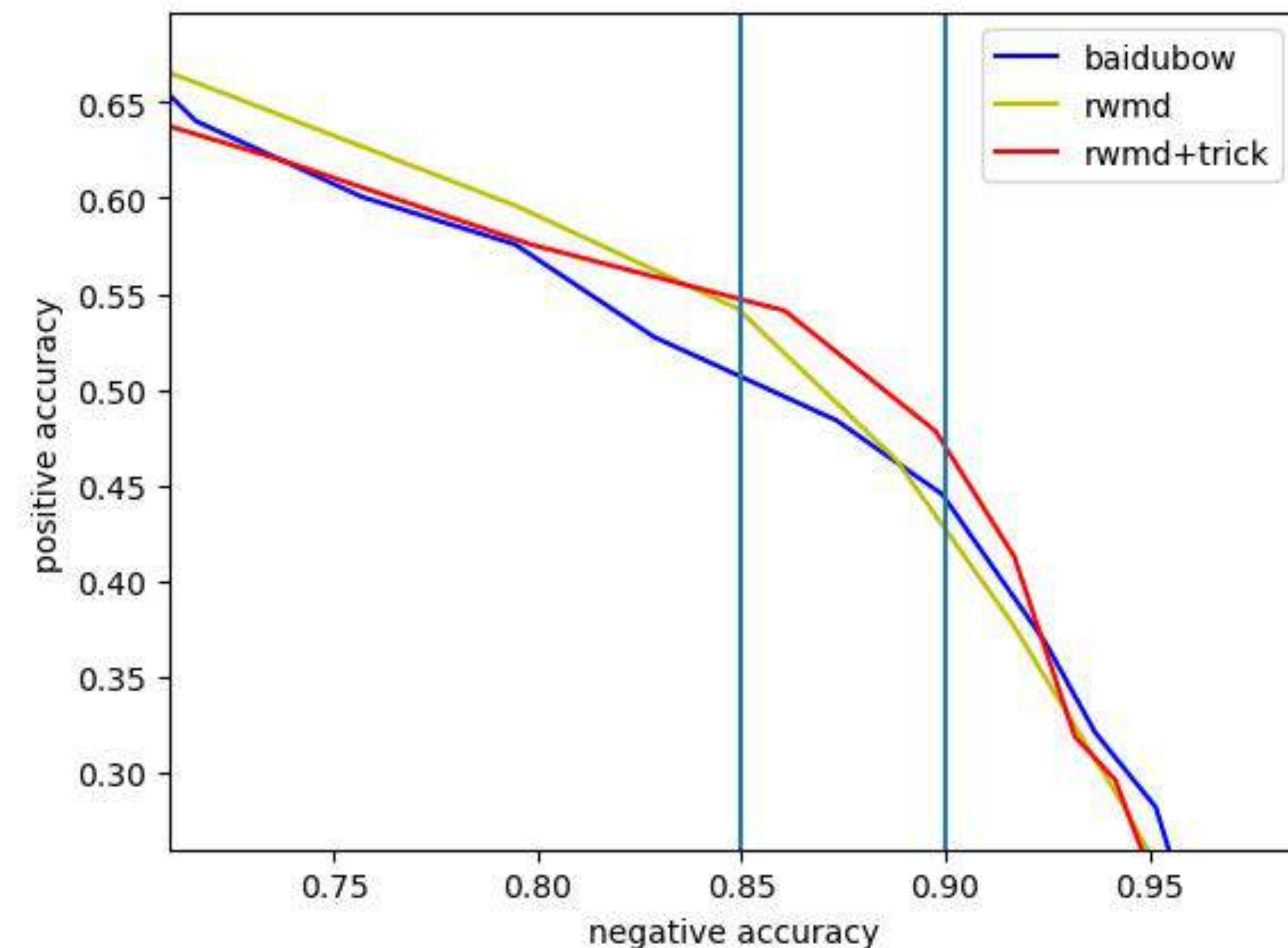
技术特色——句子相似度

提高正确率策略：改进语义相似度算法

主流语义相似度算法：CNN, wmd

Qrobot语义相似度算法：rwmd+调整策略

- 1.由于停用词影响正确率（你、我、他等），这里设置停用词去掉后剩余词多于3个（谓宾短语），则去掉停用词。否则保留停用词。
- 2.另外对于slot的实体（歌曲、歌手等）给予较低权重。



添加调整策略后正负样本准确率曲线对比

结论：
采用rwmd+调整策略方法后系统语义相似度的正确率拥有明显提升。

技术特色——实体消歧

听张学友 心如刀割 →



→ 张学友 心如刀割



看深南大道 堵车 →

→ 看深南大道 堵车 (无实体)

基于29个句子/词性特征，实体消歧采用**GBDT+Boosting** 准确率达到**95%**，多个模型多轮提升效果如下图：

使用十折交叉验证的方法，使用 6 种常用分类模型对全量数据进行了拟合和效果验证，


	Native Bayes	LR	SVM	决策树	GBDT	<u>GBDT+Boost</u>
第一轮特征	0.661	0.850	0.878	0.852	0.906	0.893
第二轮特征	0.716	0.858	0.882	0.862	0.908	0.898
第三轮特征	0.676	0.867	0.884	0.877	0.912	0.908

按照 7:3 的比例进行训练测试集划分之后，各个模型的表现如下表所示：

	Native Bayes	LR	SVM	决策树	GBDT	<u>GBDT+Boost</u>
第三轮特征	0.60	0.87	0.89	0.88	0.91	0.95

技术特色——口语化处理

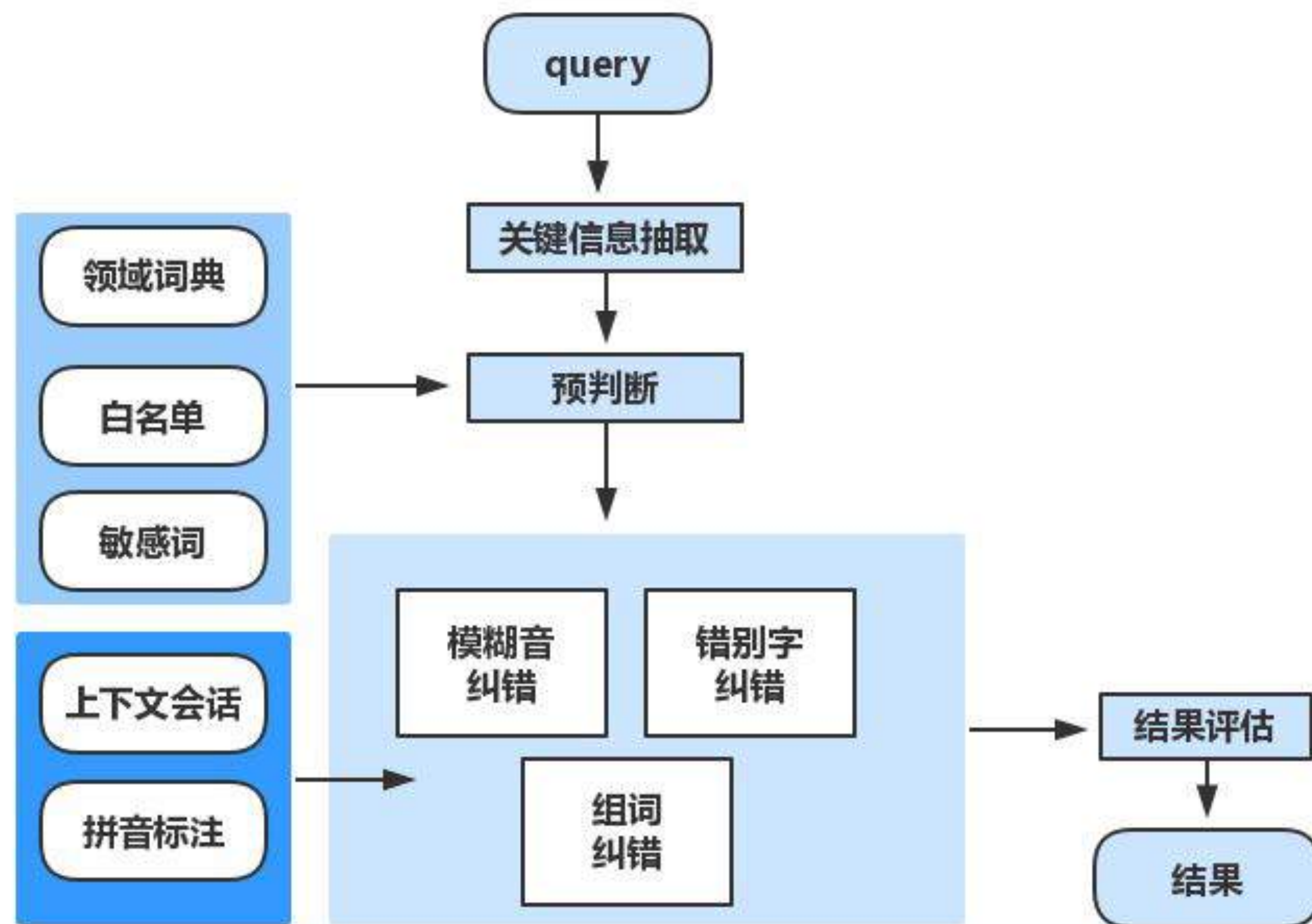
提高容错能力策略：Query Correction(QC)实现纠错

来一首周杰伦的稻香 

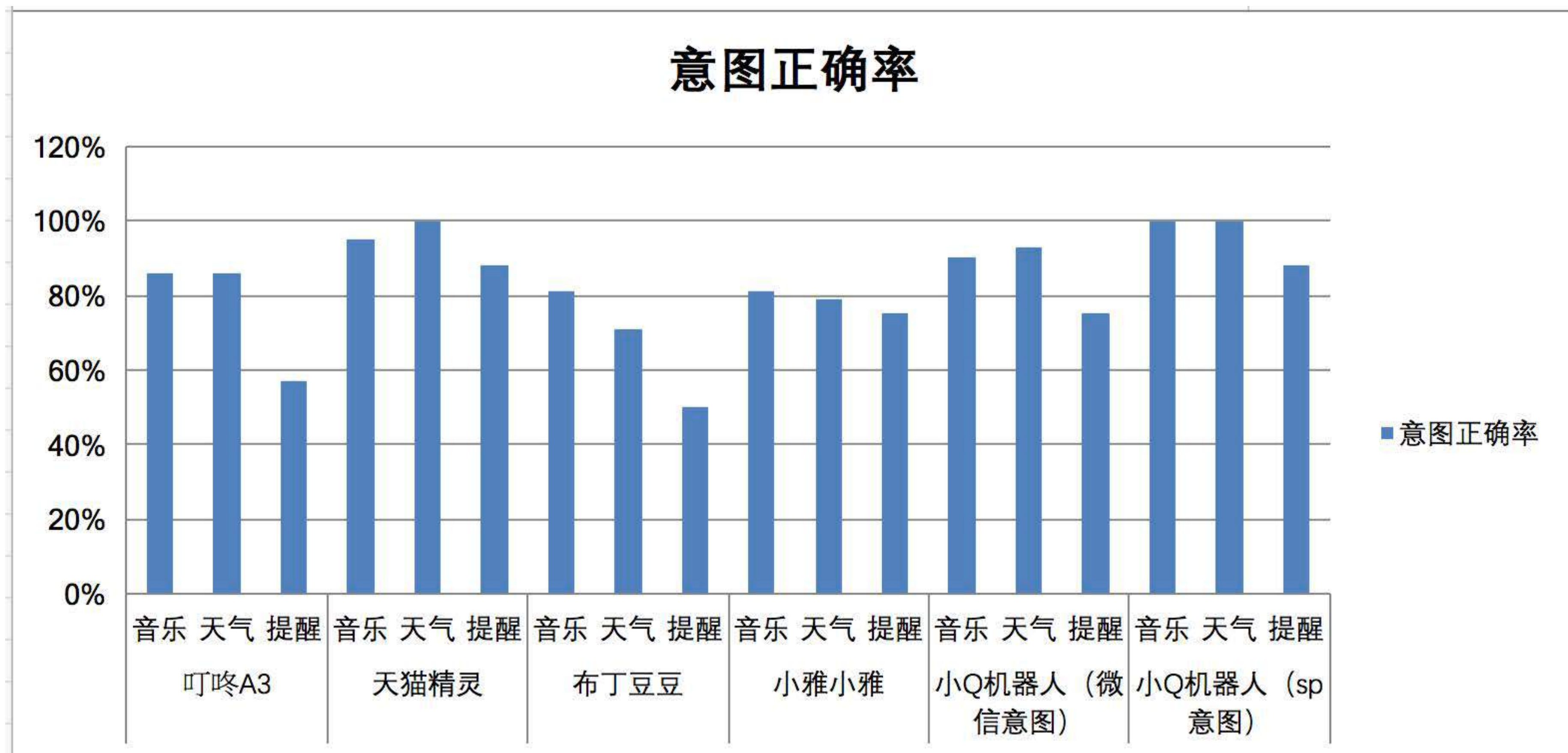
意图:音乐/点播,歌曲名:到香,歌手名:周杰伦,
电台类型:[电台类型],歌曲类型:[歌曲类型],歌
曲语言:[歌曲语言],专辑名:[专辑],影视剧名:[



```
{  
  "is_merge_second_cut" : 1,  
  "qc CGI" : [  
    {  
      "qc_cid" : 0,  
      "qc_ctr" : 0.50,  
      "qc_hotness" : 0.0,  
      "qc_query" : "周杰伦的稻香",  
      "qc_query_annotate" : "周杰伦的<em>稻</em>香",  
      "qc_score" : 2000,  
    }  
  ]  
}
```



技术成果



Agenda

1 聊天机器人

2 腾讯云小微

3 小Q机器人



小Q机器人



面临的挑战

项目挑战

团队缺乏硬件基因，硬件如何选型，**硬件研发**如何开展？

硬件项目周期长，如何做到快速验证，快速试错，快速决策

消费级产品，软硬件调优，保障基础体验，保障质量

硬件拆解



落地方案

2017.03



成熟硬件产品上快速打磨软件 →



硬件研发，原型机快速完善软件 →

软硬结合优化 →

2017.10



硬件问题



硬件研发——填坑之旅!

响应速度优化

第一个版本情况:

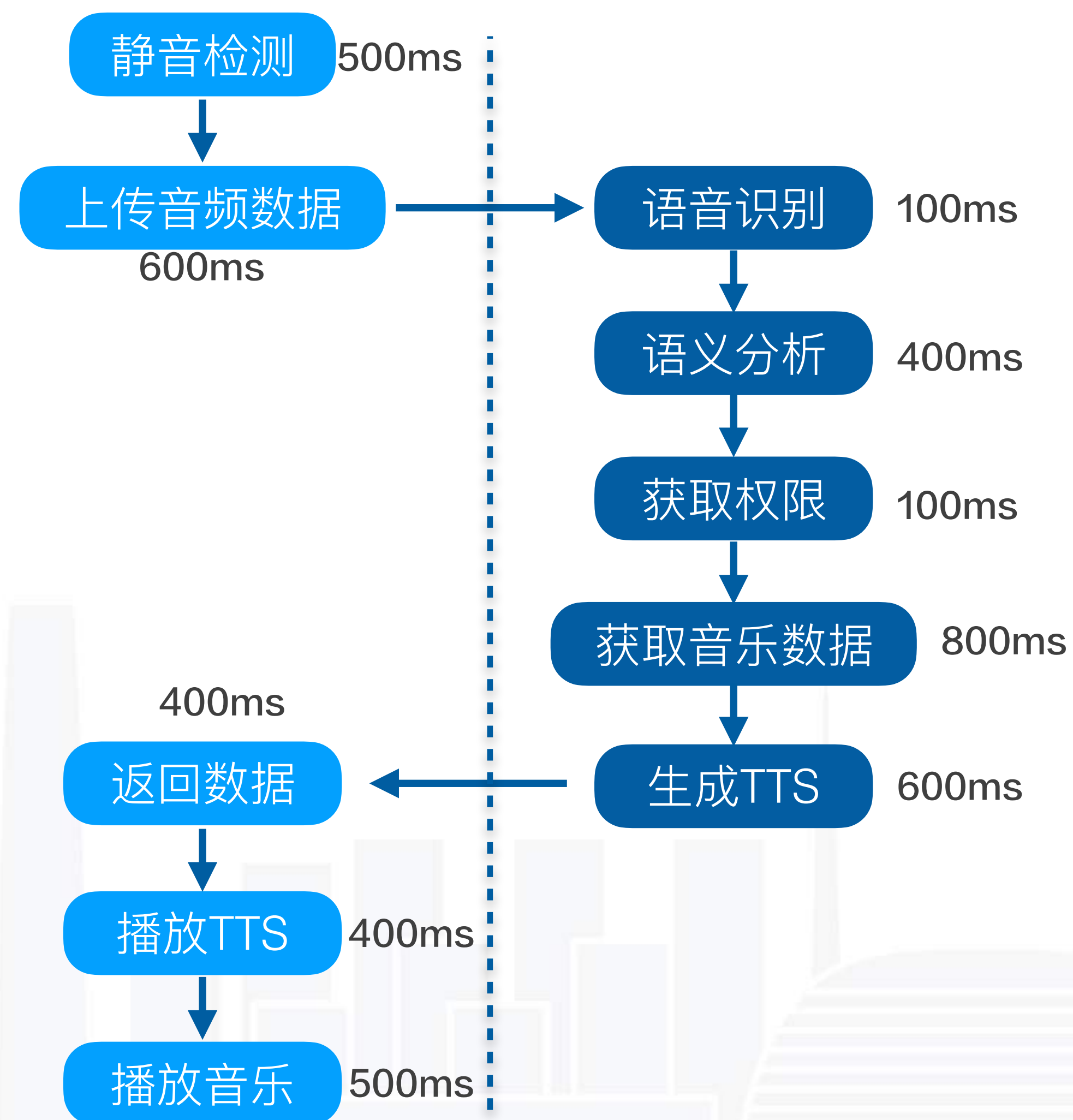
- 1、总耗时达4.4秒!
- 2、大部分步骤都必须串行;
- 3、每个流程都在不同的部门,有些功能还是刚走出实验室;
- 4、TTS耗时过长,需要1秒多;
- 5、静音检测每次都要消耗500ms;
- 6、后台耗时过长,接近2秒。

常用优化措施:

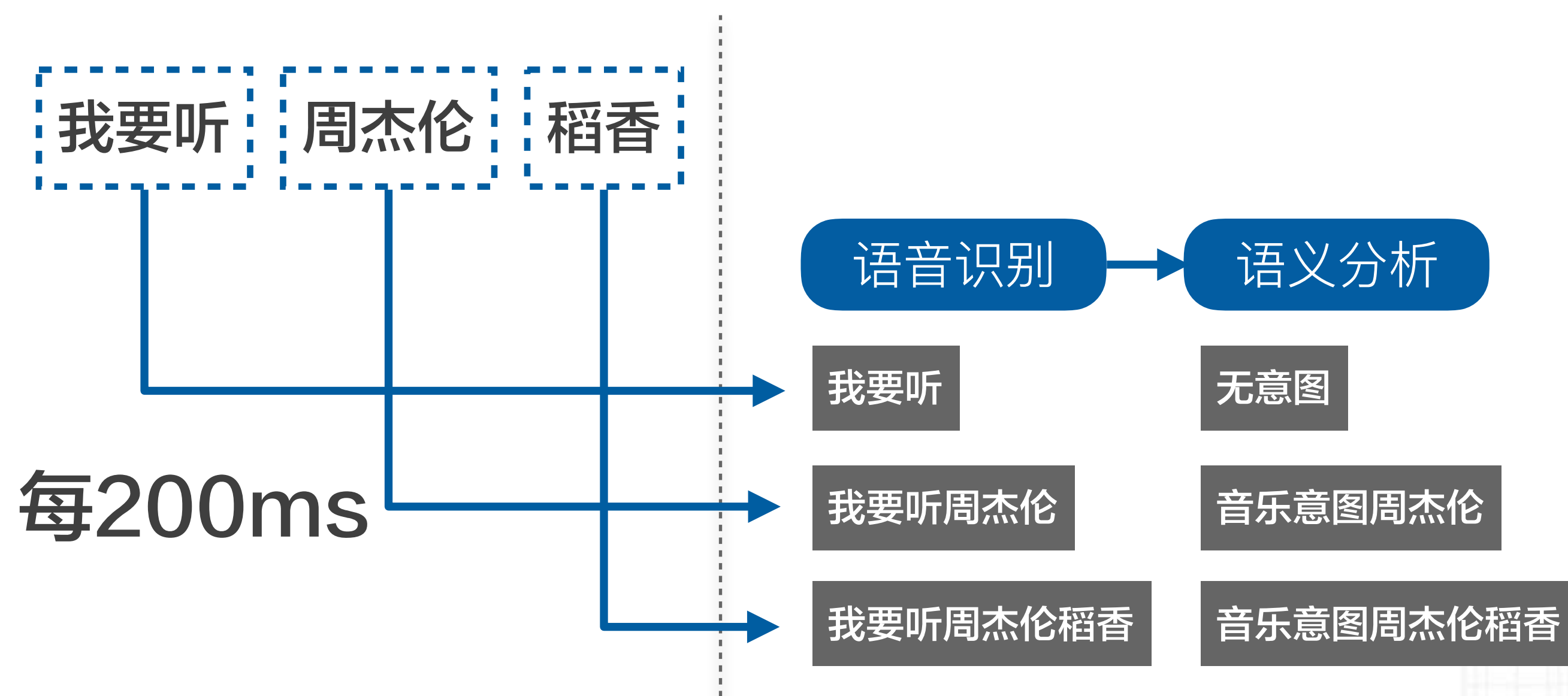
- 1、后台服务就近部署,并分离研发和线上环境,减少不必要逻辑;
- 2、优化各个步骤耗时,比如音乐从800ms优化到400ms;
- 3、一些常用指令做数据缓存。

目标是秒开,追赶Amazon Echo的应答速度。

耗时仍然无法满足要求。

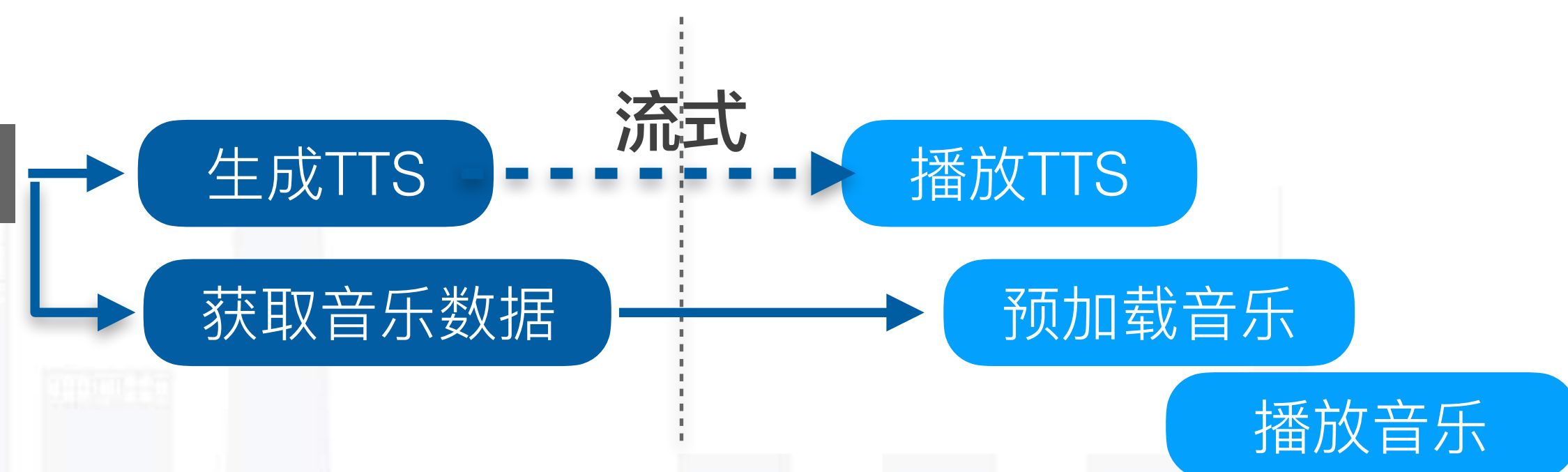


响应速度优化—全流程流式



- 1、语音识别和语义分析改为流式计算，检测到静音的同时，就能拿到语义的结果，节省耗时1000ms以上；
- 2、静音检测VAD逻辑迁移到后台，算法调整更灵活，节省约200ms。

- 3、TTS不用等待音乐查询结果，节省约400ms；
- 4、TTS改为流式生成，流式传输，流式播放，节省约800ms；
- 5、播放TTS时，预加载音乐内容，TTS播放结束，音乐可以立即播放，节省约200ms。



响应速度优化——TTS流式优化

传统TTS模式：

一次性生成音频文件，返回地址客户端下载播放。



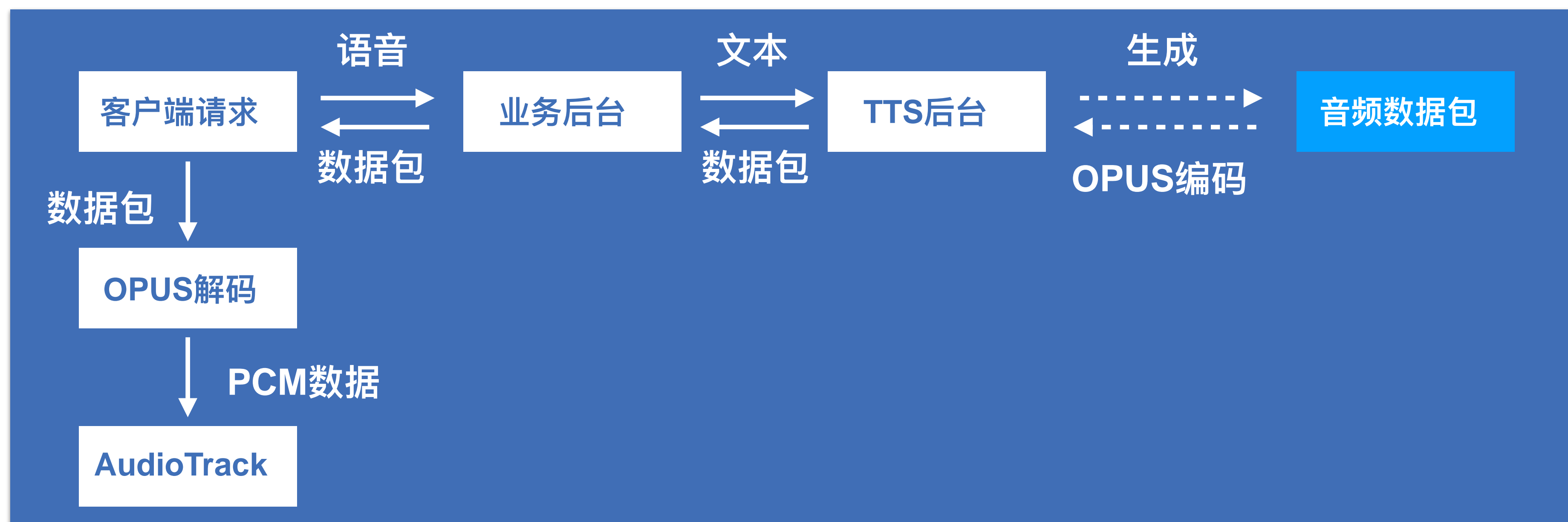
痛点：耗时！（在百科类长内容上非常明显）

- 1、TTS后台生成完整音频文件时间
- 2、客户端拿到音频地址后时间到可以播放的时间

响应速度优化——TTS流式优化

流式TTS模式：

根据文本实时生成语音包数据，借助TCP长链接，OPUS将语音包编码后，有序送往客户端解码播放，OPUS可将语音包从2880b减至300b。



与传统模式对比优势：

- 1、流式TTS的播放时间不依赖于文本的长度，在长文本和短文本上时间无差距
- 2、客户端无需要下载，长链接会以PUSH的形式有序传递给客户端解码播放