

AI语音产品测试

演讲人：于海生





于海生

美团技术专家

3年开发经验，9年测试经验。

先后经历过音视频中间件、ERP中间件、搜索、自然语言交互等领域测试

智能语音应用场景有哪些？

智能语音产品如何测试？

主题

- AI简介
- 智能语音技术应用场景
- 如何测试语音产品并评定质量
- QA

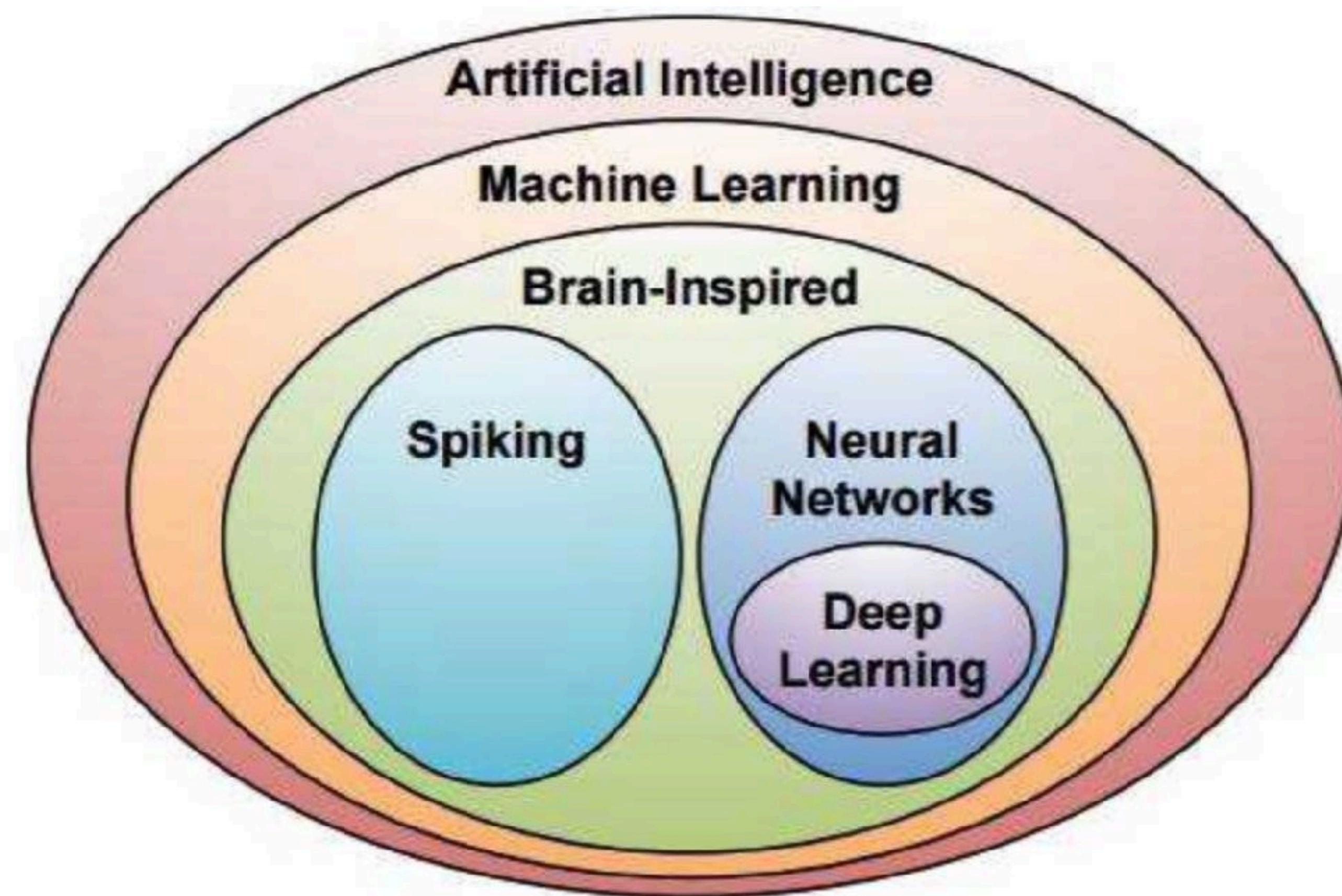
techopedia

Artificial intelligence (AI) is an area of computer science that emphasizes the creation of intelligent machines that work and react like humans.

wikipedia

Artificial intelligence (AI), sometimes called machine intelligence, is intelligence demonstrated by machines, in contrast to the natural intelligence displayed by humans and other animals.

AI vs ML vs NN vs DL

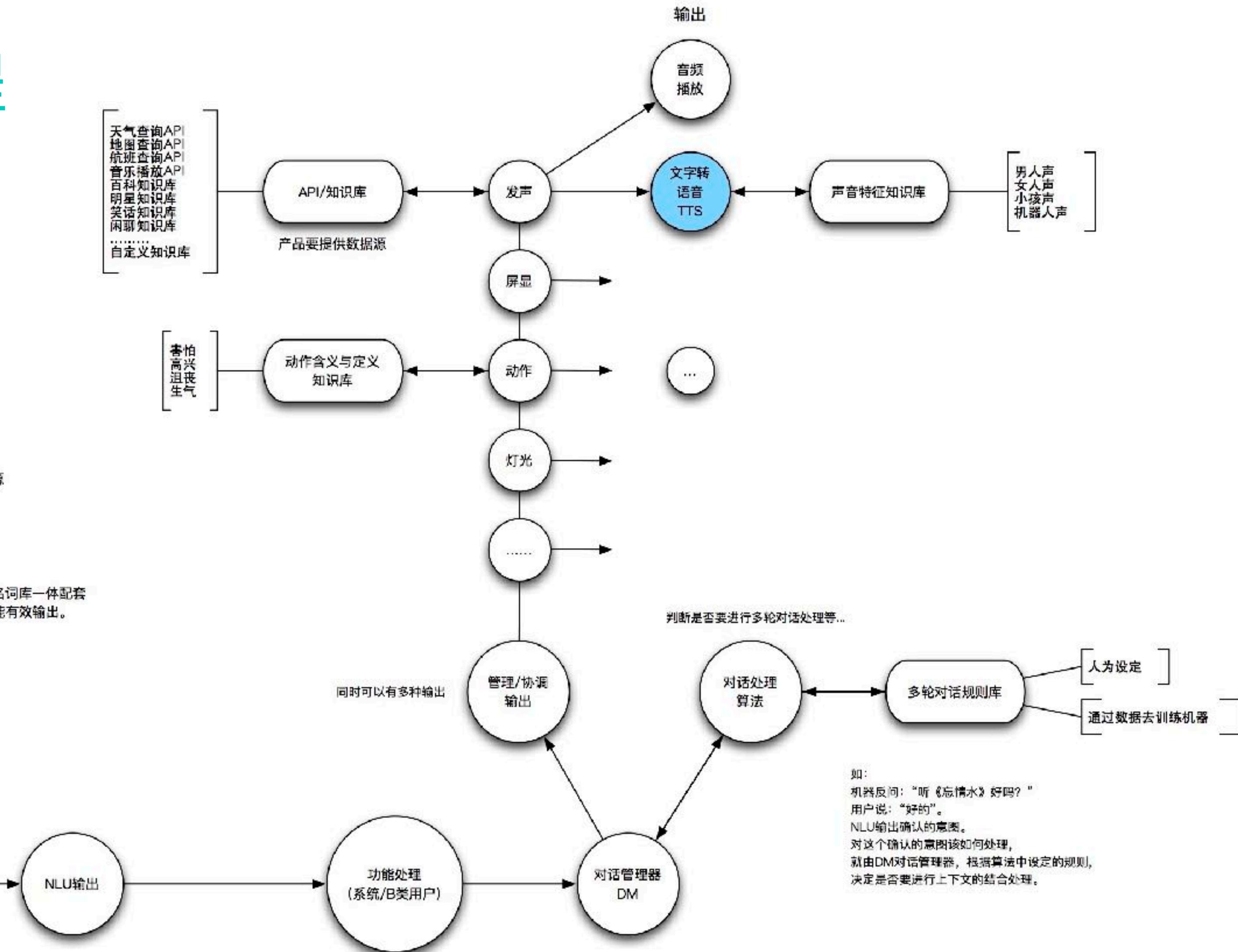
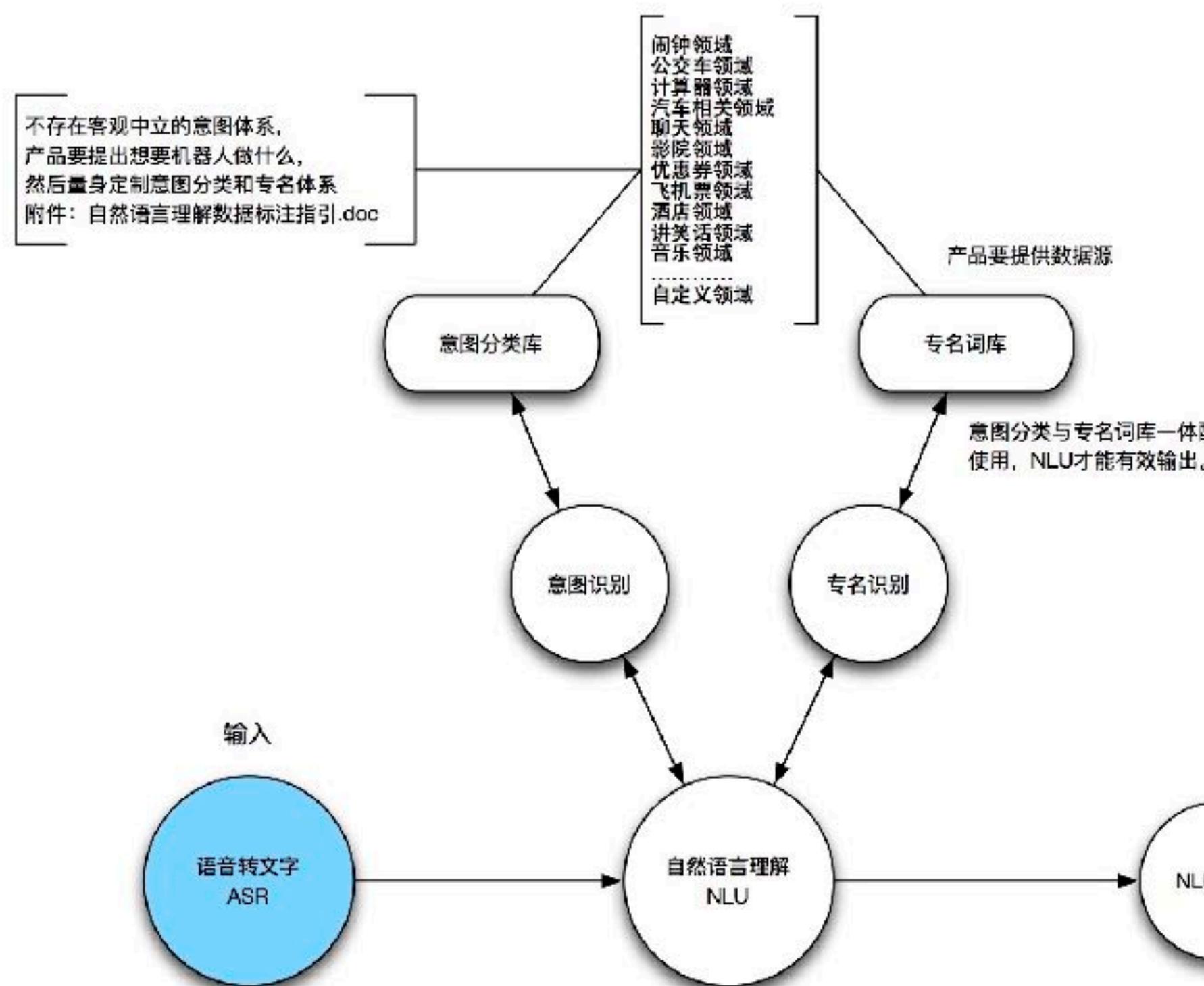


- AI简介
- 智能语音技术应用场景
- 如何测试语音产品并评定质量
- QA



- AI简介
- 智能语音技术应用场景
- **如何测试语音产品并评定质量**
- QA

语音硬件自然对话过程



我想听刘德华的歌

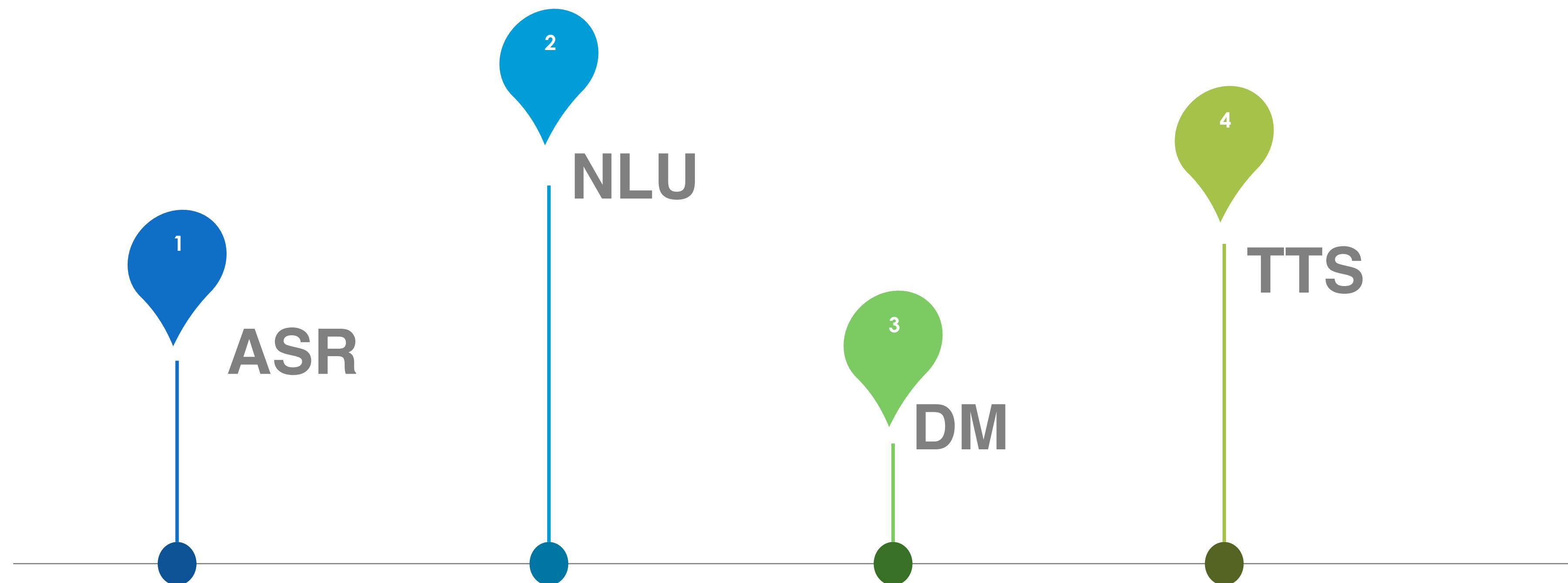
我/可忽略 想听/关键字 刘德华/歌手 的/可忽略 歌/关键字

意图="music_play-music"
专名={"歌手": "刘德华"}

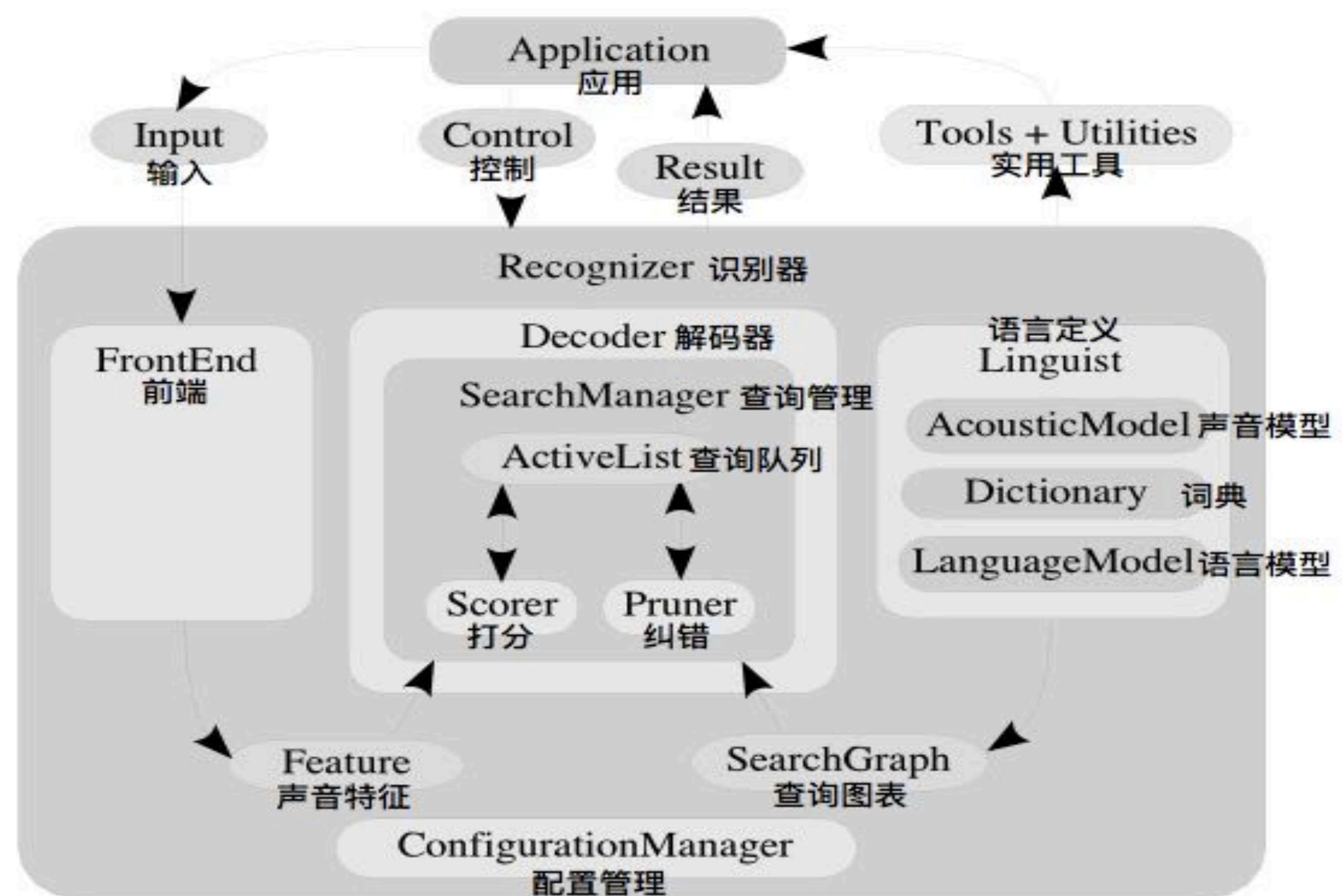
以上是数据标注的示例
绿色是意图识别的关键词(特征)
蓝色是抽取出来的专名(实体)

其中的意图输出叫意图标签

语音产品核心架构图



ASR语音识别处理过程



ASR
质量指标

A 识别率

- 客户体验角度

识别率 : 100-SER (句错率)

SER: 句子识别错误的个数 / 总的句子个数

- 工程角度

识别率 : 100-WER (word error rate , 词错率)

分为男女、快慢、口音、数字、英文、中文。

B 语音唤醒
VOICE TRIGGER , VT :

- 唤醒率
- 误唤醒率 : 没叫AI的时候 , ta自己跳出来讲话的比率
- 唤醒词的音节长度 : 一般技术上要求 , 最少3个音节
- 唤醒响应的时间
- 功耗

NLU

自然语言理解

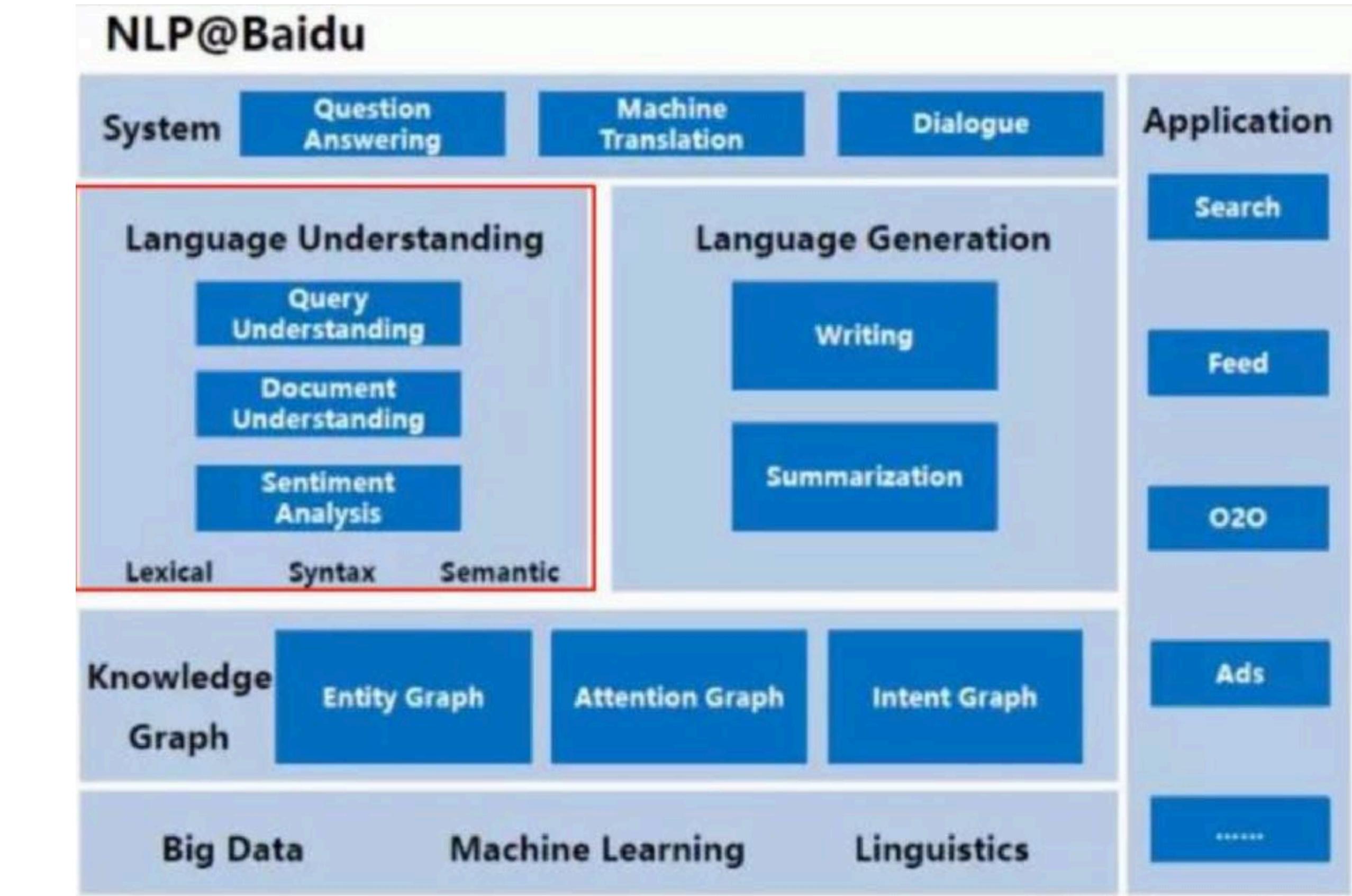
Natural Language Understanding, NLU

- ◆ 句子检测
- ◆ 分词
- ◆ 词性标注
- ◆ 句法分析
- ◆ 文本分类/聚类
- ◆ 信息抽取/自动摘要等

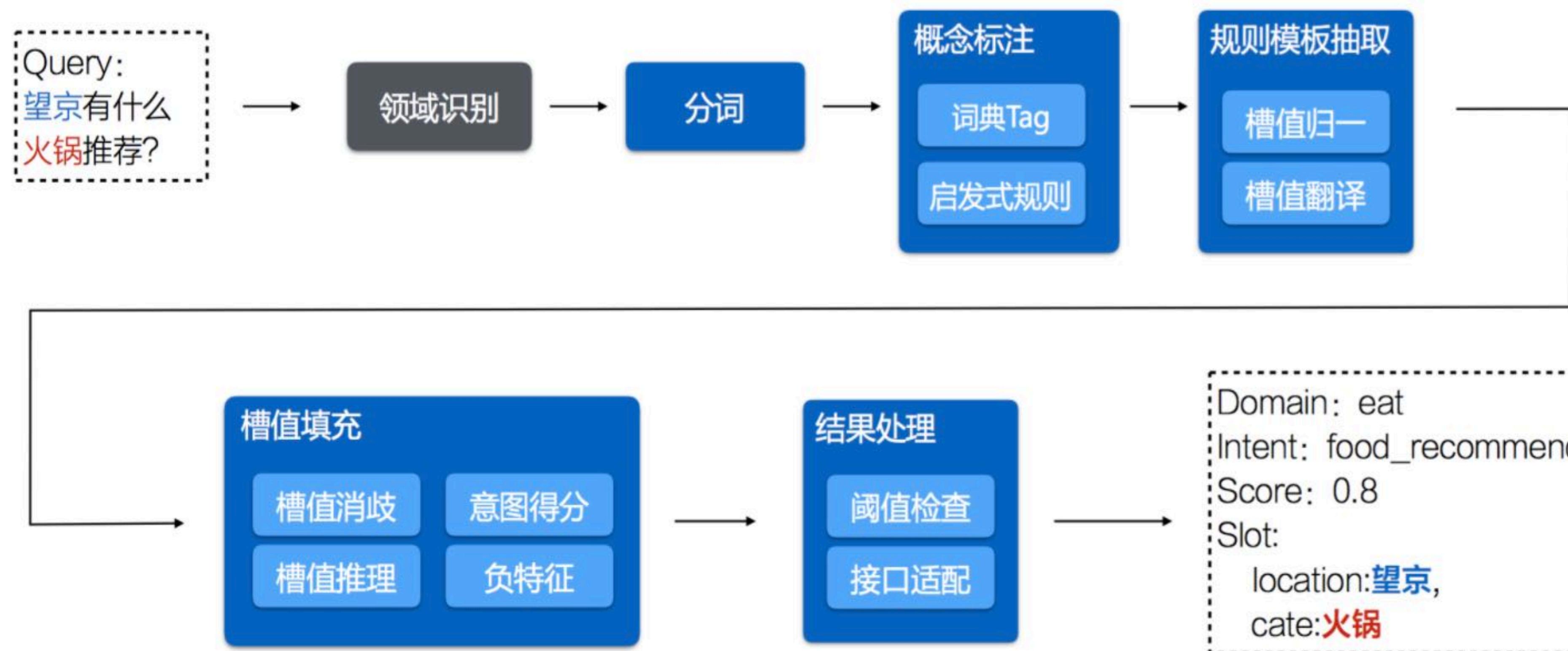
自然语言处理

Natural Language Processing, NLP

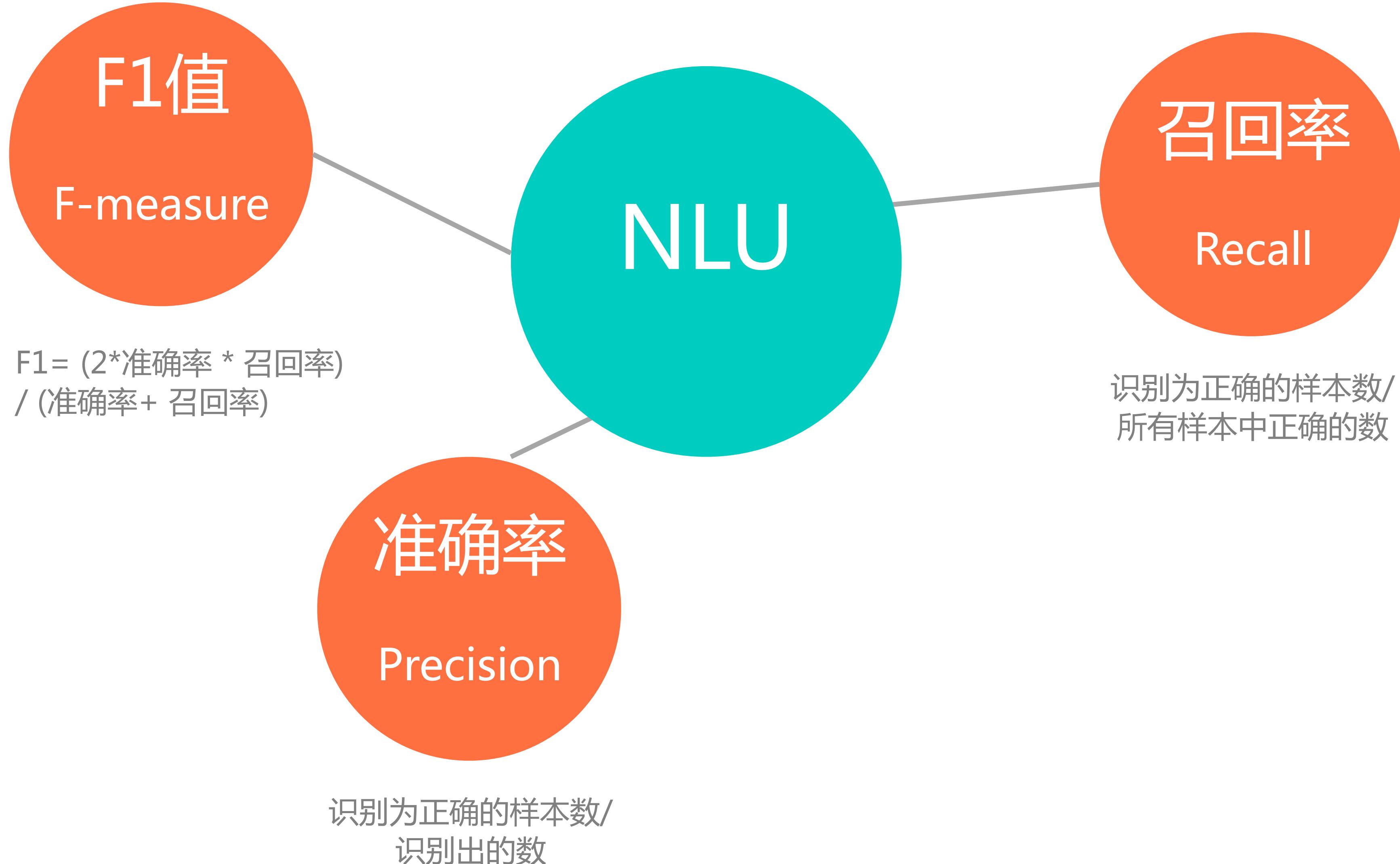
- ◆ 是自然语言和AI交集的学科，包括语音处理相关的所有技术



NLU处理过程

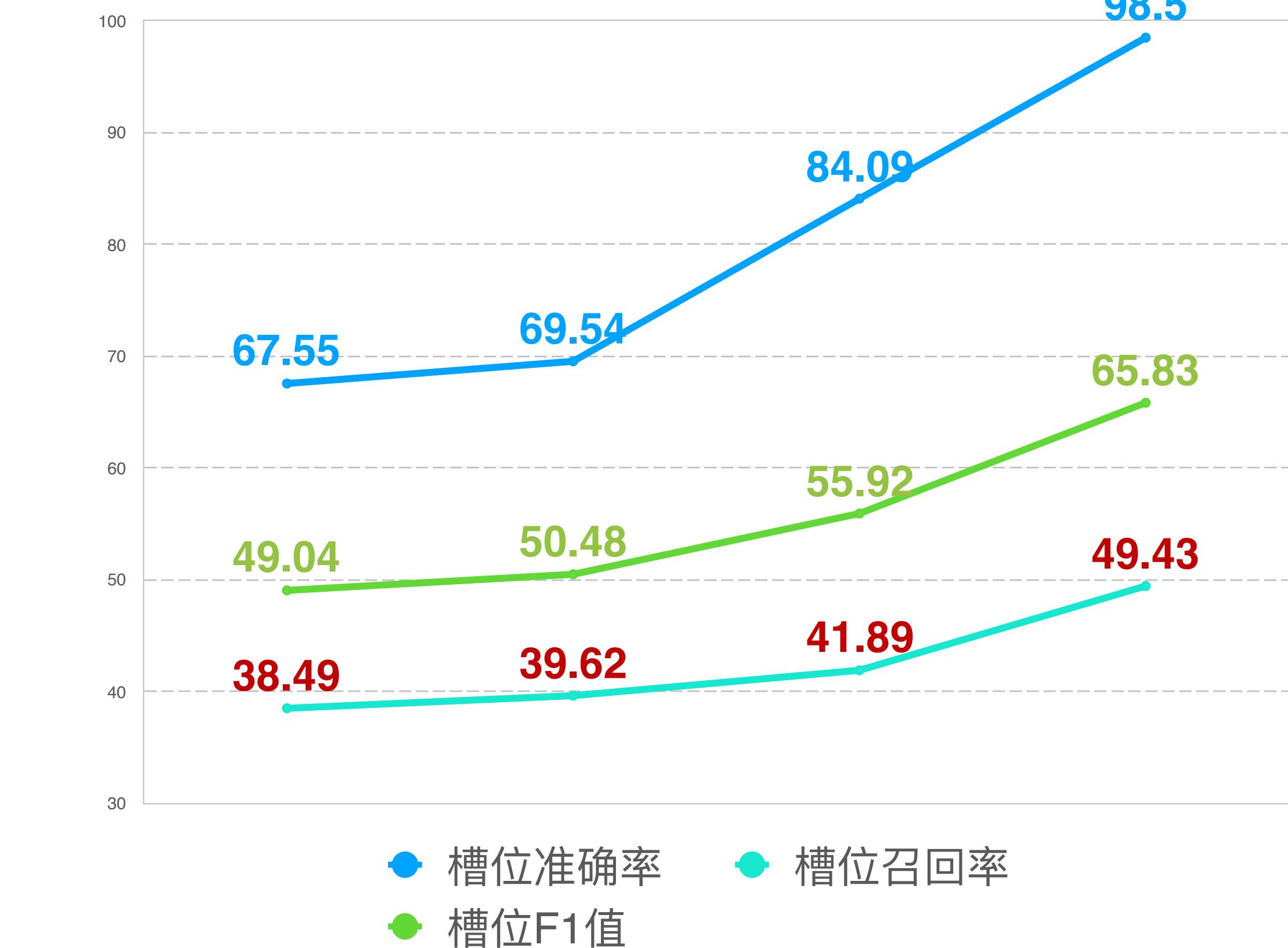
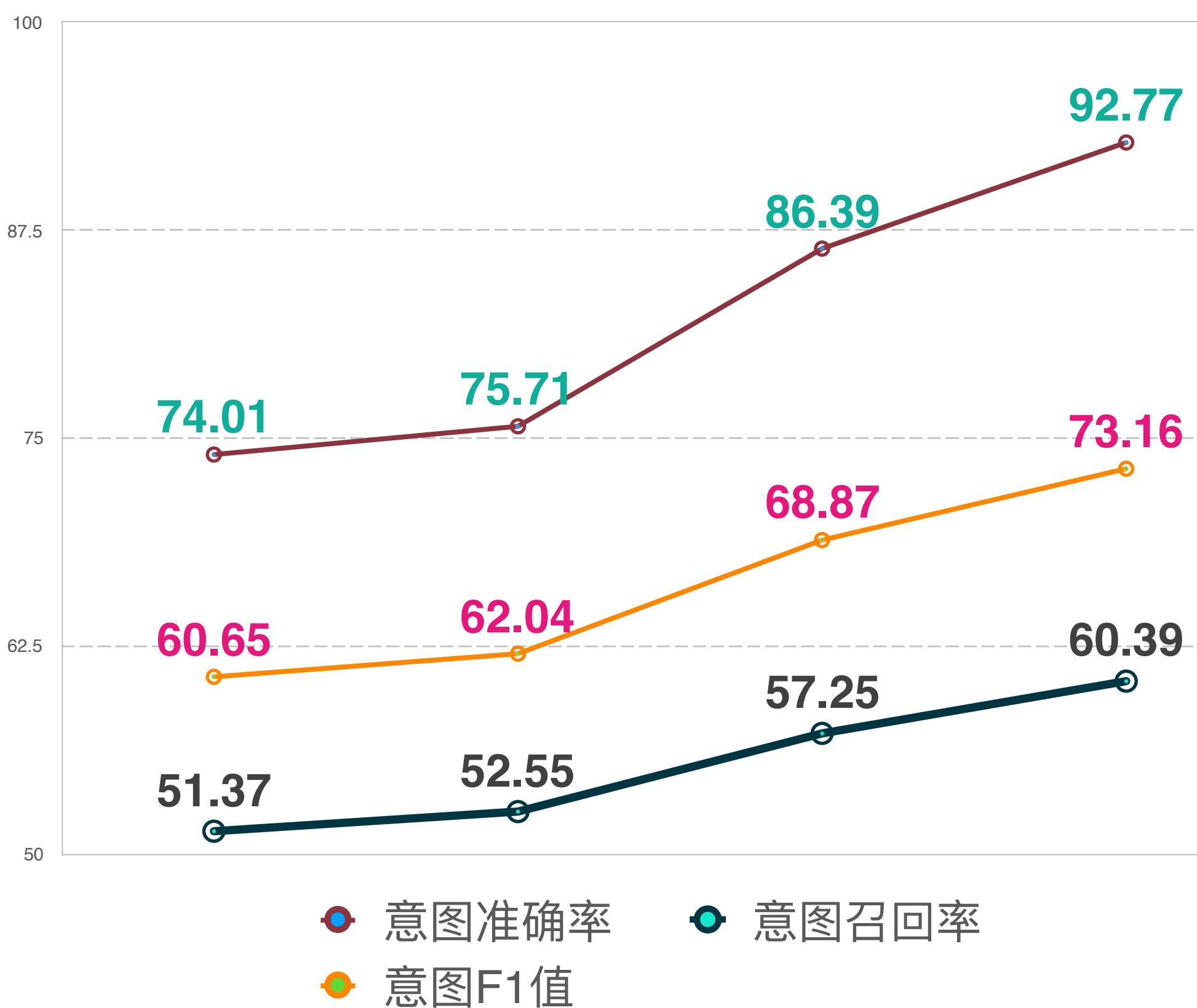


NLU
质量指标



AI

NLU效果评估



AI

对话系统

特定领域



通用领域

任务型对话

- 机器人
- 智能音箱

QA问答

- 手机语音助手
- 车载语音系统

开放域聊天

- 小冰
- Siri
- 机器人
- 智能音箱

AI
对话系统



用户Query+Context

意图识别

中控系统

Chat-Bot

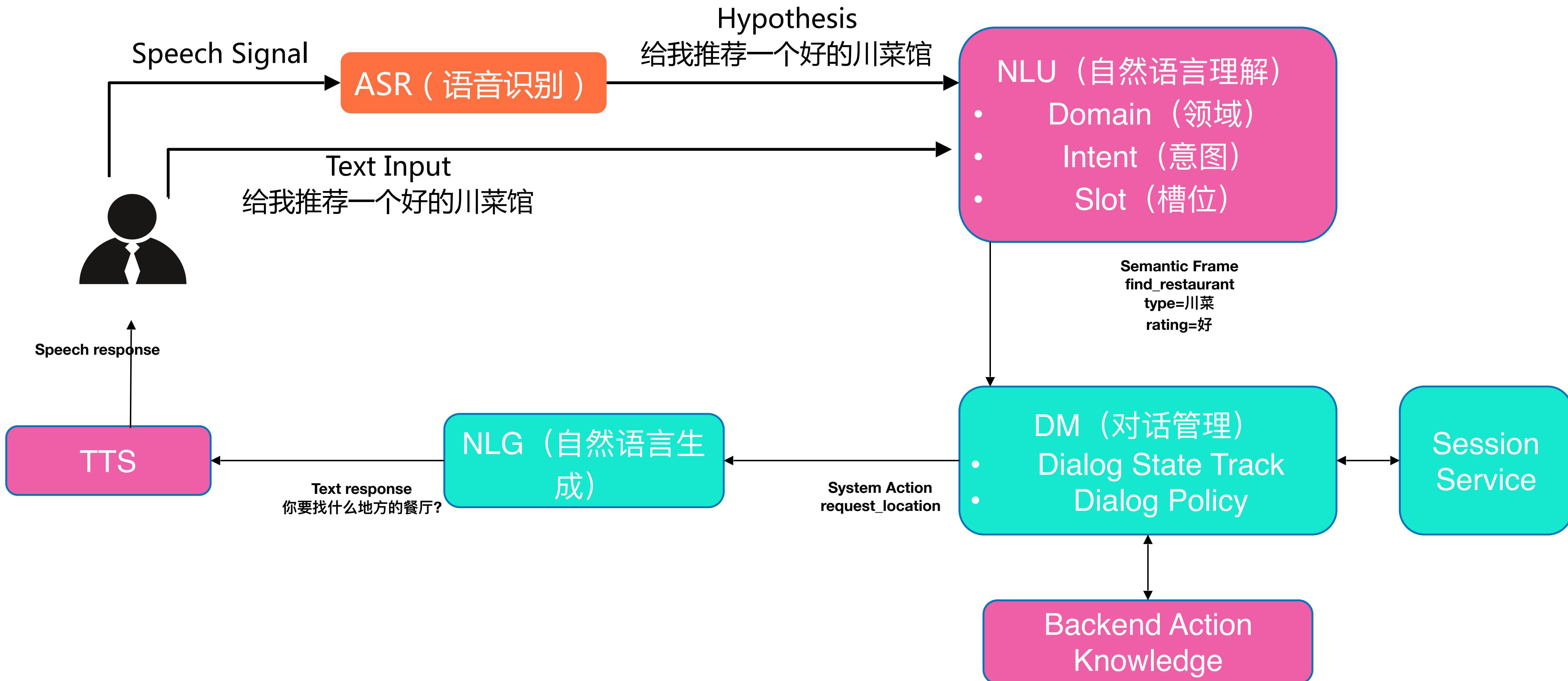
QA-Bot

Task-Bot

Other-Bots

AI

任务型对话系统



DM

质量指标

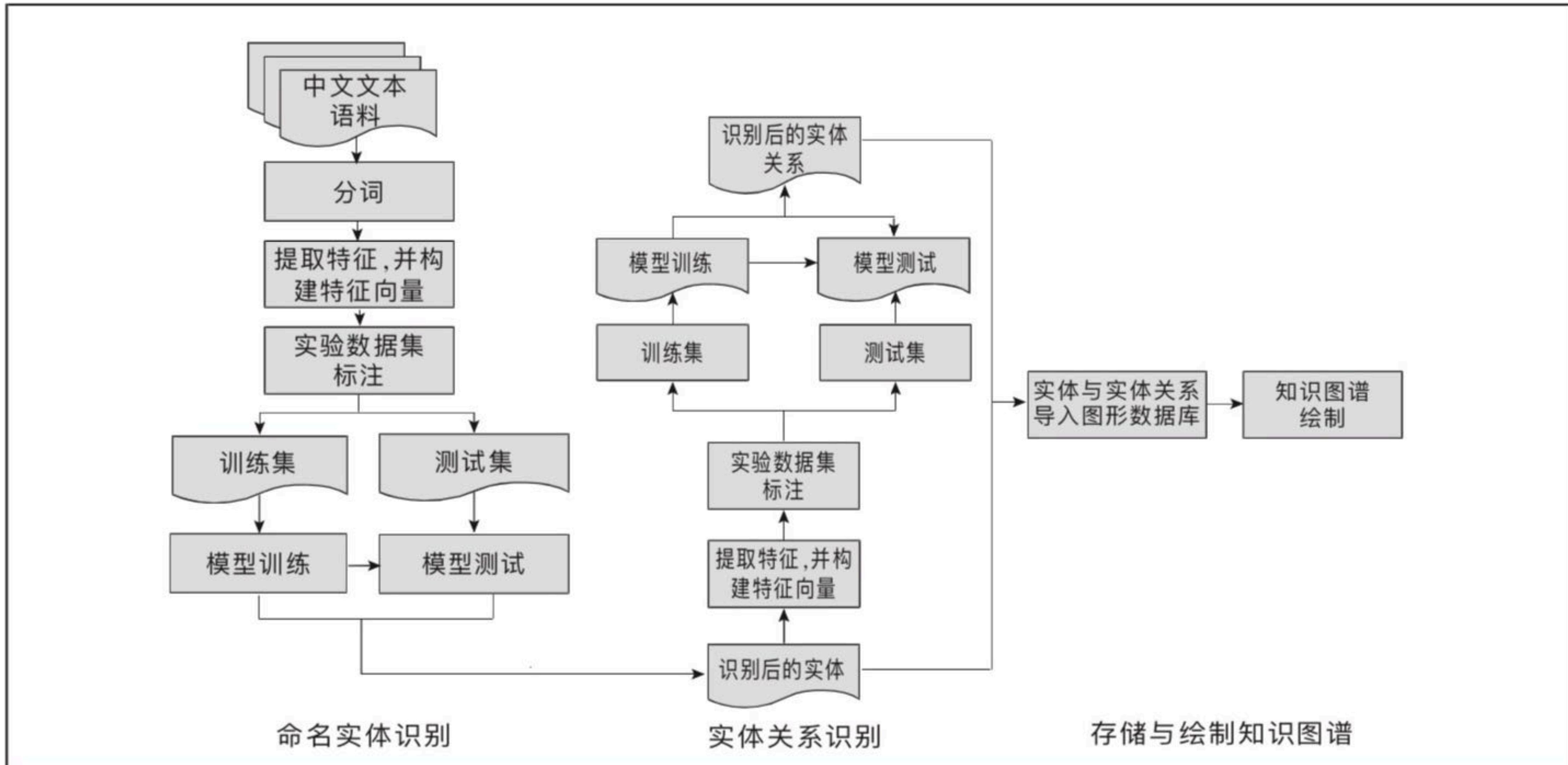
闲聊型

任务型

问答型

- | | | |
|---|-------------------------------------|------------------|
| ① CPS (Conversation
s Per Session, 平均
单次对话轮数) | ① 留存率：可发现用
户有没有形成使用
习惯 | ① 最终求助人工比例 |
| ② 相关性和新颖性 | ② 完成率：用户任务
达成率 | ② 重复问同样问题的
比例 |
| ③ 话题终结者 | ③ 相关指标：每个任
务评价slot填入轮数
或填充完整度 | ③ “没答案”之类的比
例 |

基于深度学习算法的商业知识图谱构建流程



KG
质量指标

KG

knowledge Graph

准确率

覆盖率

AI应用孵化器

效果指标-知识图谱效果

北京	属性名	属性值	覆盖率	全量	属性名	属性值	覆盖率
商品总数: 51500 抽取总数: 51500 upc缺失: 5116	品牌	蒙牛等82个	96.12%(49504)	商品总数: 2246104 抽取总数: 2246104 upc缺失: 515597	品牌	蒙牛等82个	87.28%(1960508)
	规格	ml	86.94%(44775)		规格	ml	78.98%(1773886)
	包装	盒, 瓶, 袋, 罐, 包	54.67%(28157)		包装	盒, 瓶, 袋, 罐, 包	69.02%(1550247)
	口味 (多)	草莓等19个	26.39%(13592)		口味 (多)	草莓等19个	31.44%(706192)
	营养成分	全脂, 低脂, 脱脂	10.59%(5456)		营养成分	全脂, 低脂, 脱脂	4.55%(102207)

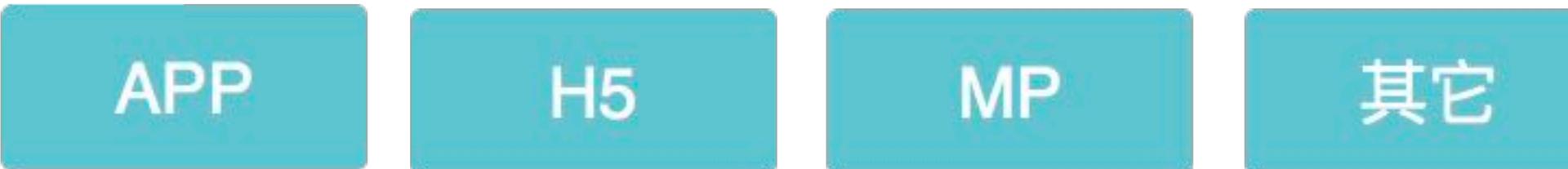
	品类	品牌	规格	单重	数量	包装	口味	营养成分
0	0	0	31	31	31	64	137	185
1	198	181	160	161	160	135	56	15
2	2	0	1	0	1	1	0	0
3	0	19	8	8	8	0	7	0
准确率	99.00%	90.50%	95.50%	96.00%	95.50%	99.50%	96.50%	100.00%

商超领域，知识图谱完成初次数据覆盖度、
数据准确率评估，覆盖28个品类，10个
属性，图例为牛奶品类，数据覆盖度、数
据准确率情况

AI

语音产品核心

用户层



接入层



业务层



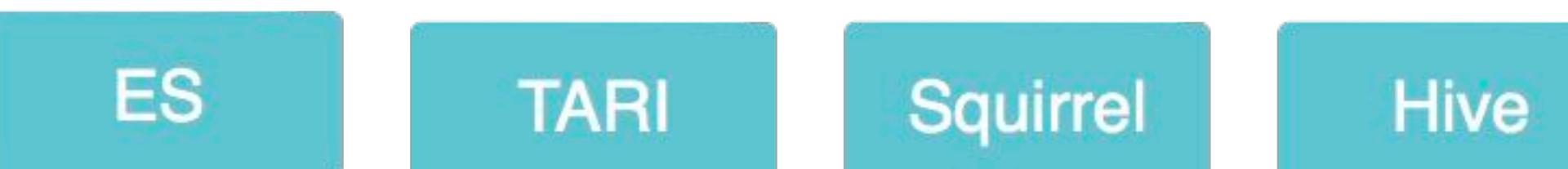
基础服务层



基础数据层



数据源



01

ASR

将声音转化为文字的过程,相当于人类的耳朵

02

NLU

让计算机能够理解和生成人类语言,包括分词、词性分析、信息抽取等, 相当于人类的大脑

03

DM

负责聊天管理和体验, 即维护和更新对话的状态, 以及对话状态, 做出决策, 相当于人类的大脑

04

TTS

将文字转化为声音, 相当于人类的嘴巴

AI

效果指标-实现策略

01 按优先级

按重要程度

优先实现核心业务的效果指标

02 按实现成本

按实现难度、时间成本

优先实现成本低收益高的效果指标

03 采集数据策略

新项目PM 提供输入数据，QA补充
随着产品迭代，扩充真实用户输入

指标	评价公式	优先级	成本
平均对话轮数	完成任务的平均交互次数	★★	★★
任务完成率	达成用户意图数/用户意图数	★★	★★
没答案比例	没答案数/样本数	★★	★★
唤醒率	呼唤次数/被成功唤醒次数	★	★
识别率	(插入词数+替换词数+删除词数) / 标准词	★	★★
召回率	正确数/样本数	★★★	★
准确率	正确数/识别数	★★★	★
F1值	(2*准确率 * 召回率) / (准确率+召回率)	★★★	★
用户任务达成率	同任务完成率	★★	★★
对话交互效率	跳出DM现有交互流程次数/总交互次数	★★	★★
数据准确率	正确数/识别数	★★	★★
数据覆盖度	分维度已有数据属性/全网分维度数据属性	★★★	★★

Q&A

二分类问题常用指标：

<https://blog.csdn.net/quincuntial/article/details/69596456>

https://blog.csdn.net/zhihua_oba/article/details/78677469

书单

- 科学的极致：漫谈人工智能
- 走进2050，注意力、互联网与人
工智能
- 与机器人共舞
- 智能时代
- 人工智能：一种现代的方法

公众号

- 机器之心
- 全球人工智能
- 哈工大SCIR
- Ai科技评论

CODE A BETTER LIFE

一 行 代 码 亿 万 生 活



更多技术干货
欢迎关注“美团技术团队”

招聘：测试开发岗位
邮箱：yuhaiseng@meituan.com

