

AiCon

全球人工智能与机器学习技术大会

AI时代的数据解决方案

李鑫

百度资深研发工程师

主办方 **Geekbang** 极客邦科技 **InfoQ**

TABLE OF CONTENTES

人工智能行业现状

数据之于人工智能

人工智能基础数据面临的难题

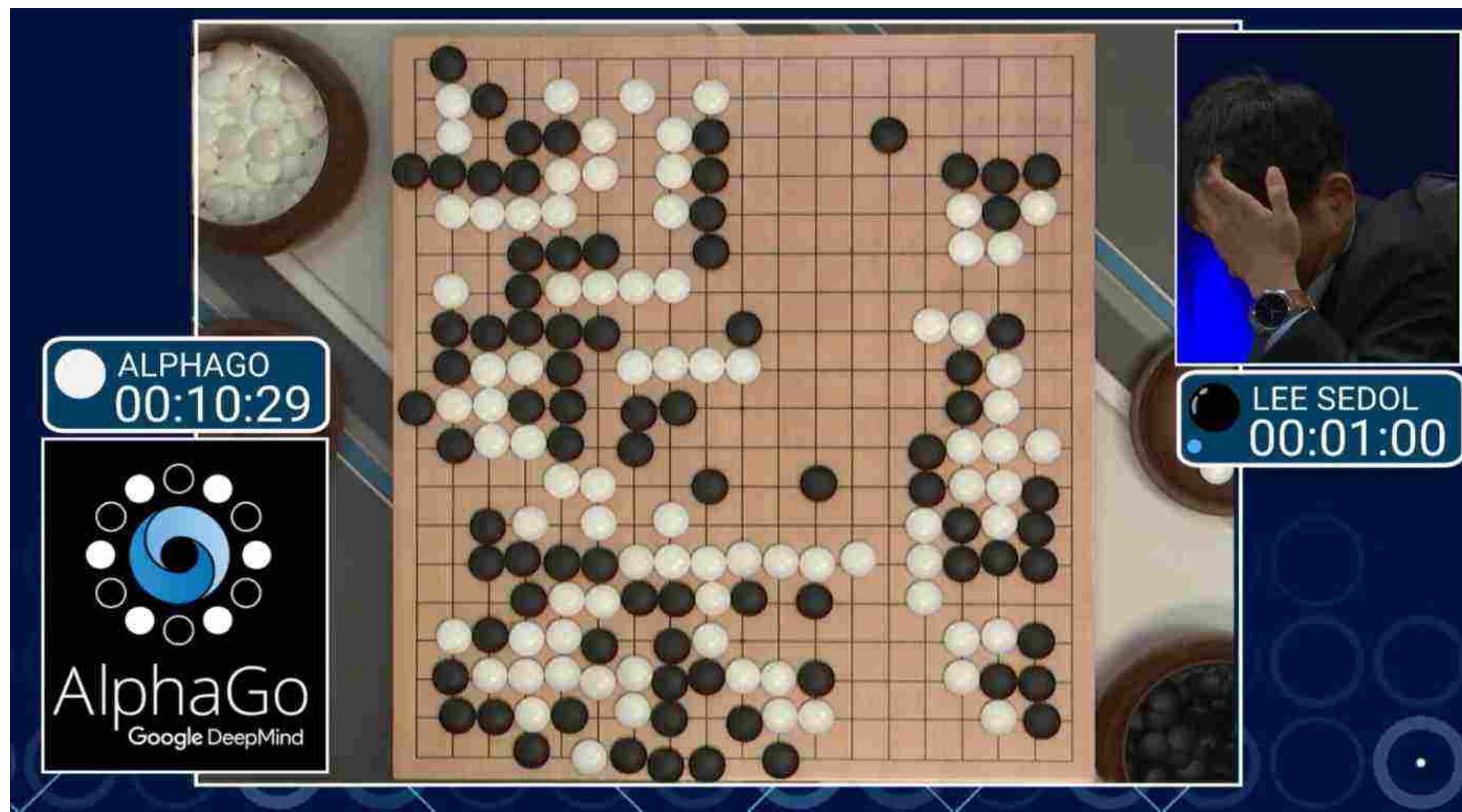
百度是如何应对的

典型人工智能应用场景

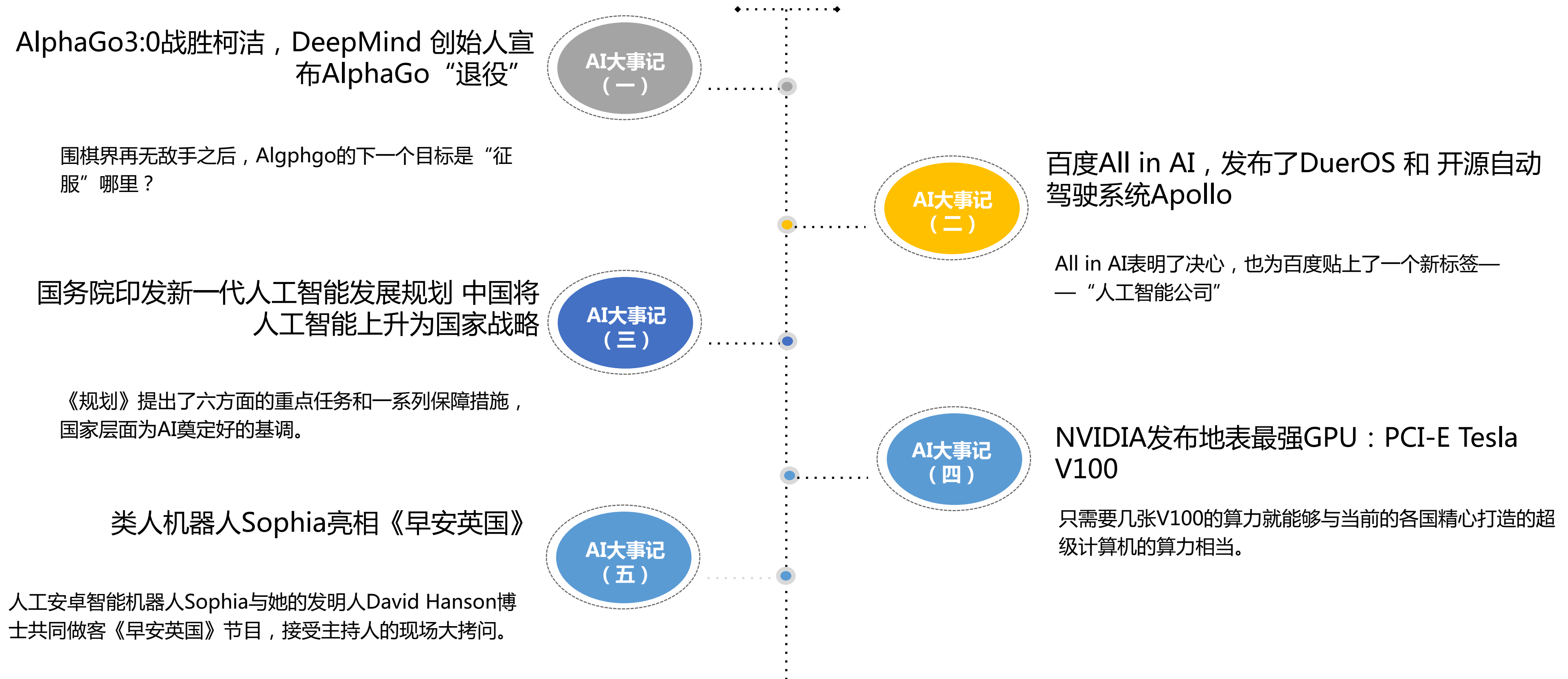
百度数据产品

人工智能进入公众视野

人工智能的强大能力已被证明

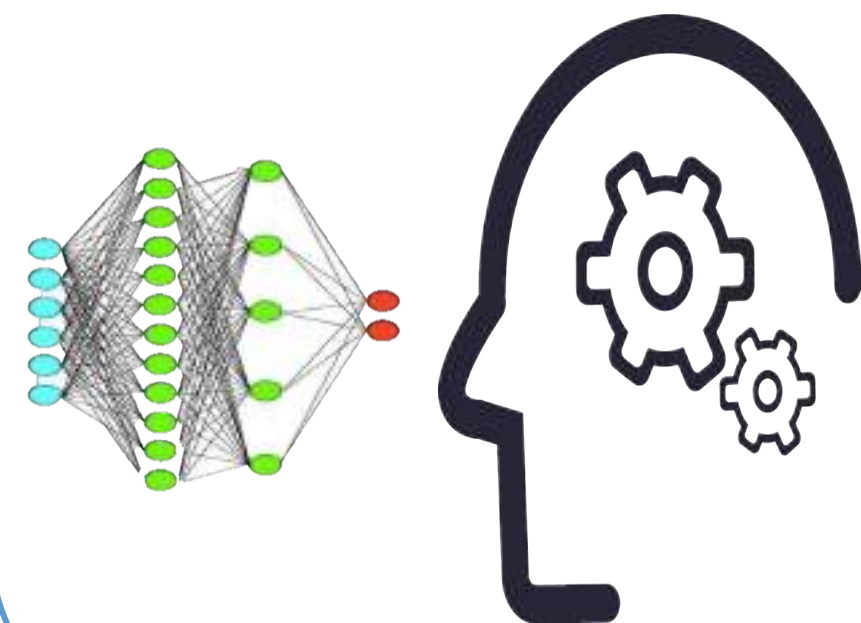


人工智能2017大事记

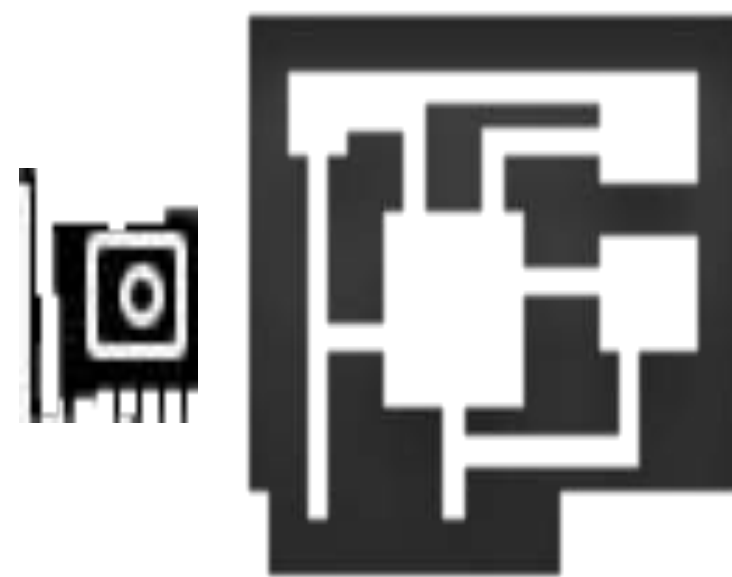


人工智能爆发的三大因素

深度学习



高性能运算



大数据



算法是核心，**计算**、**数据**是基础

TABLE OF CONTENTES

人工智能行业现状

数据之于人工智能

人工智能基础数据面临的难题

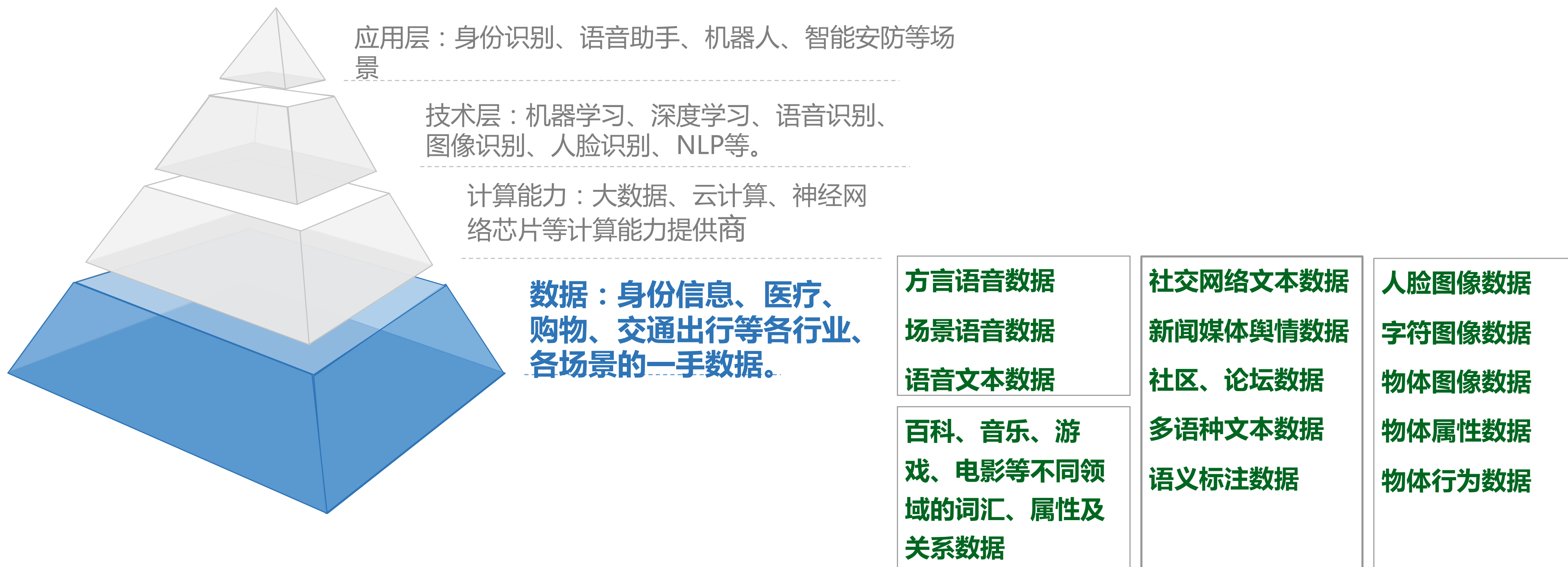
百度是如何应对的

典型人工智能应用场景

百度数据产品

海量、精准、高质量的数据是人工智能的根本

数据是一切人工智能技术和应用实现的基础保障和前提！

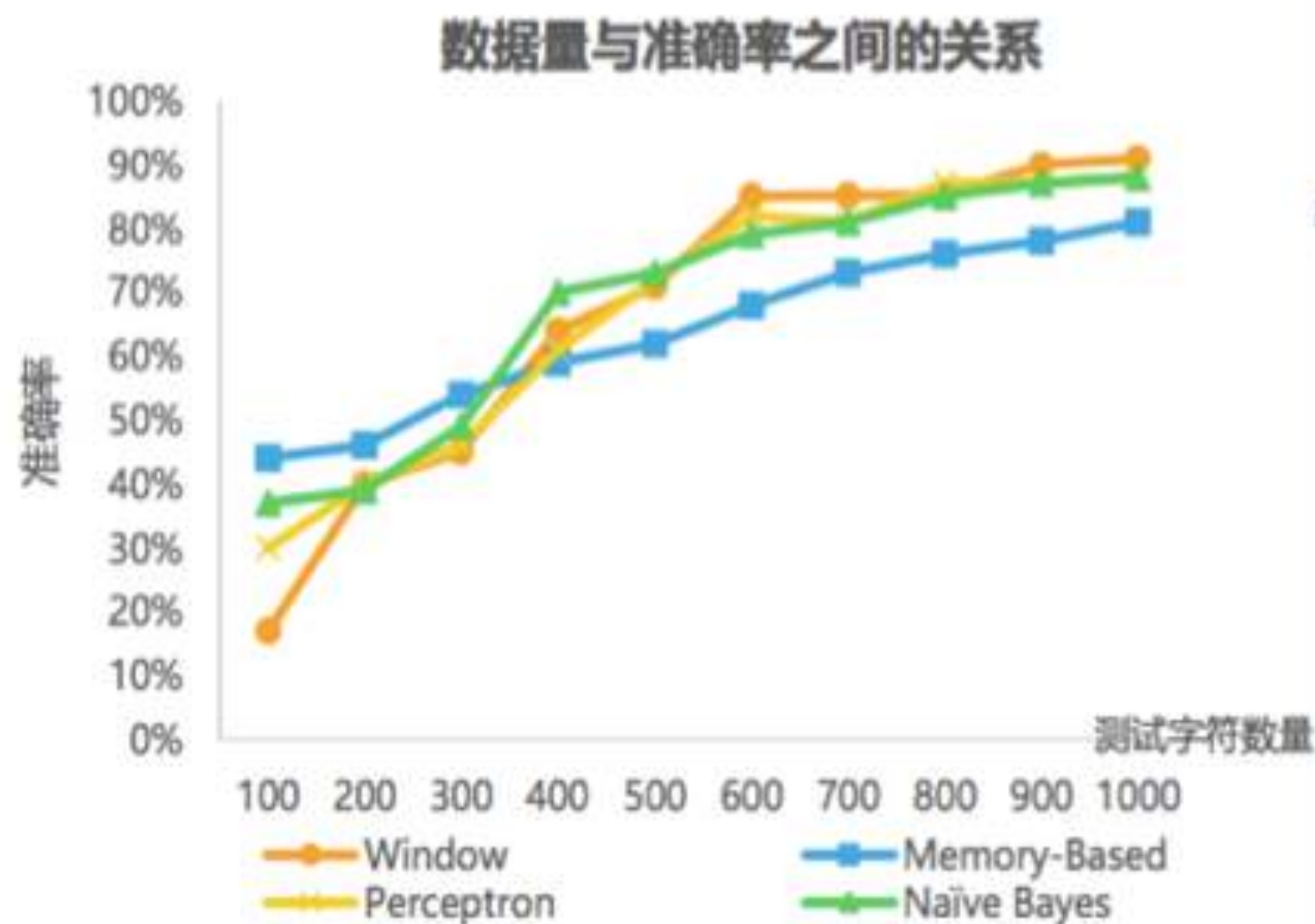


数据样本与算法模型

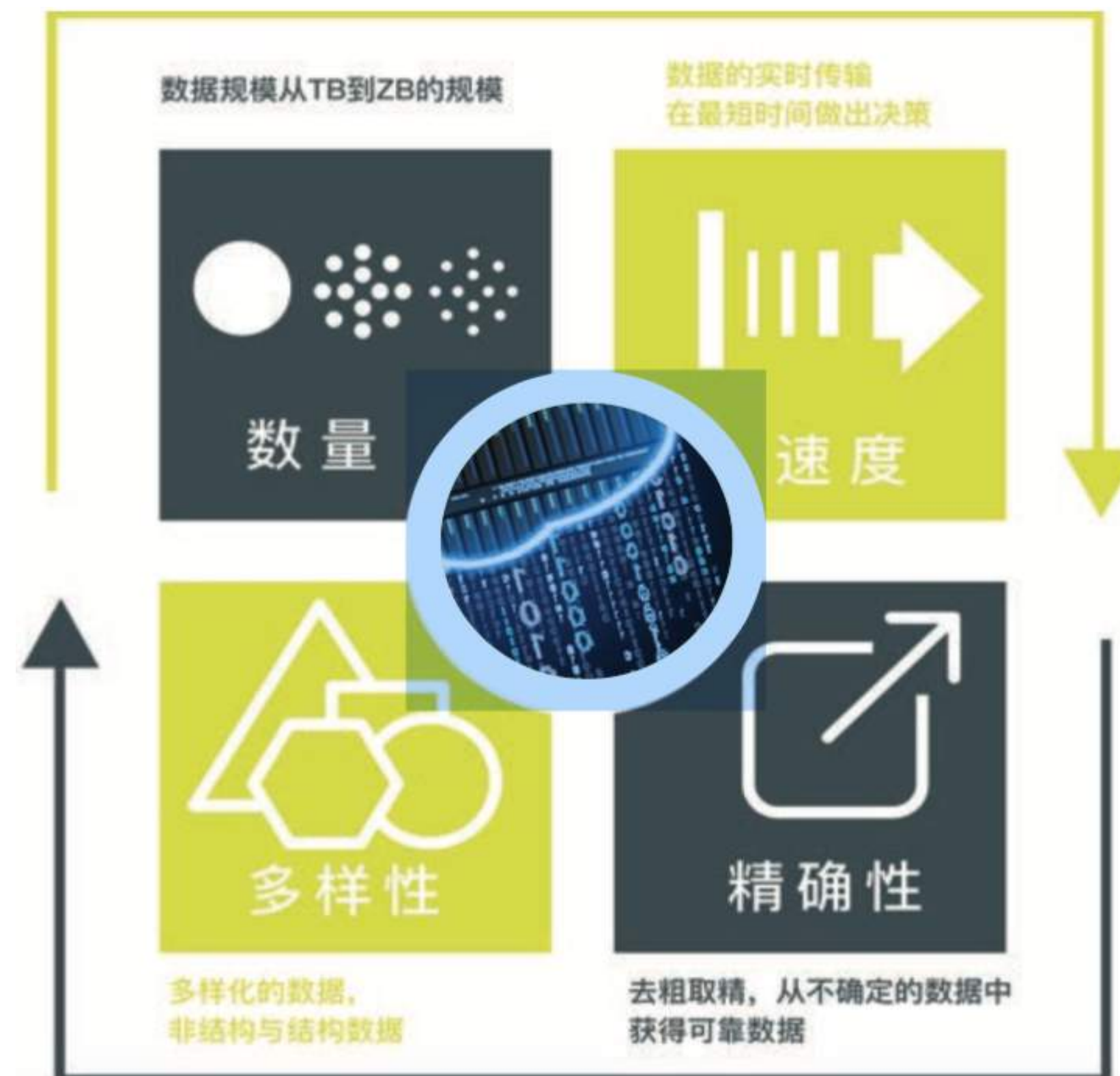


人工智能需要通过**大量的数据样本**来“训练”自己，才能不断提升输出结果的质量。

有时候，数据真的可以秒杀算法



说明：window、memory-based、perceptron、naive bayes 均为不同算法
来源：Stanford机器学习公开课，36氪研究院



快人一步抢占先机，数据竞赛“质&量”取胜

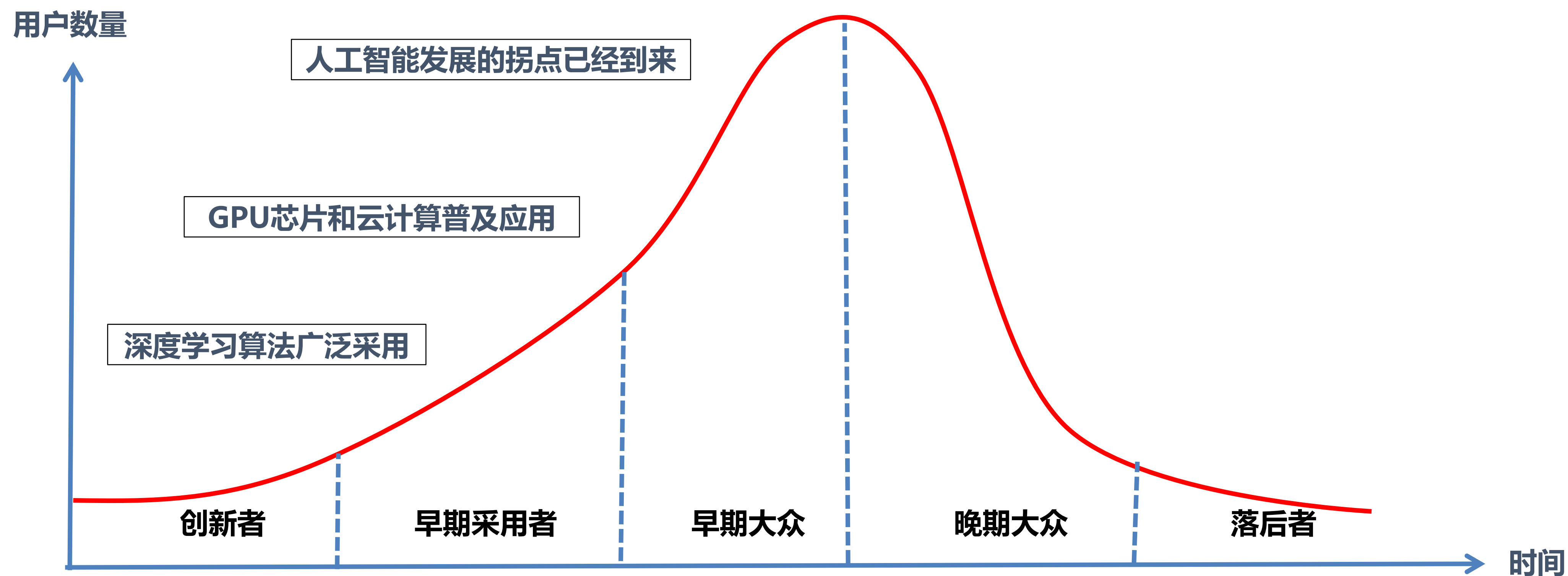


TABLE OF CONTENTES

人工智能行业现状

数据之于人工智能

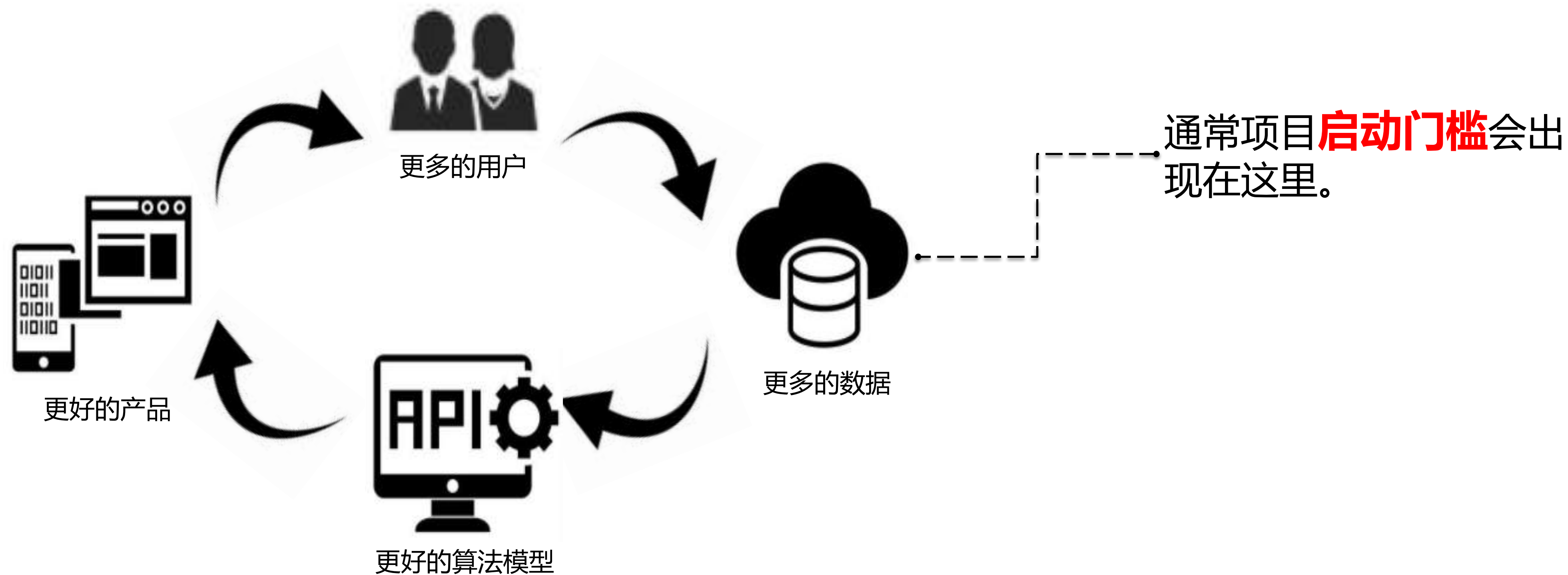
人工智能基础数据面临的难题

百度是如何应对的

典型人工智能应用场景

百度数据产品

项目“冷”启动的数据困扰



获取和加工数据，AI基础数据的两大难题

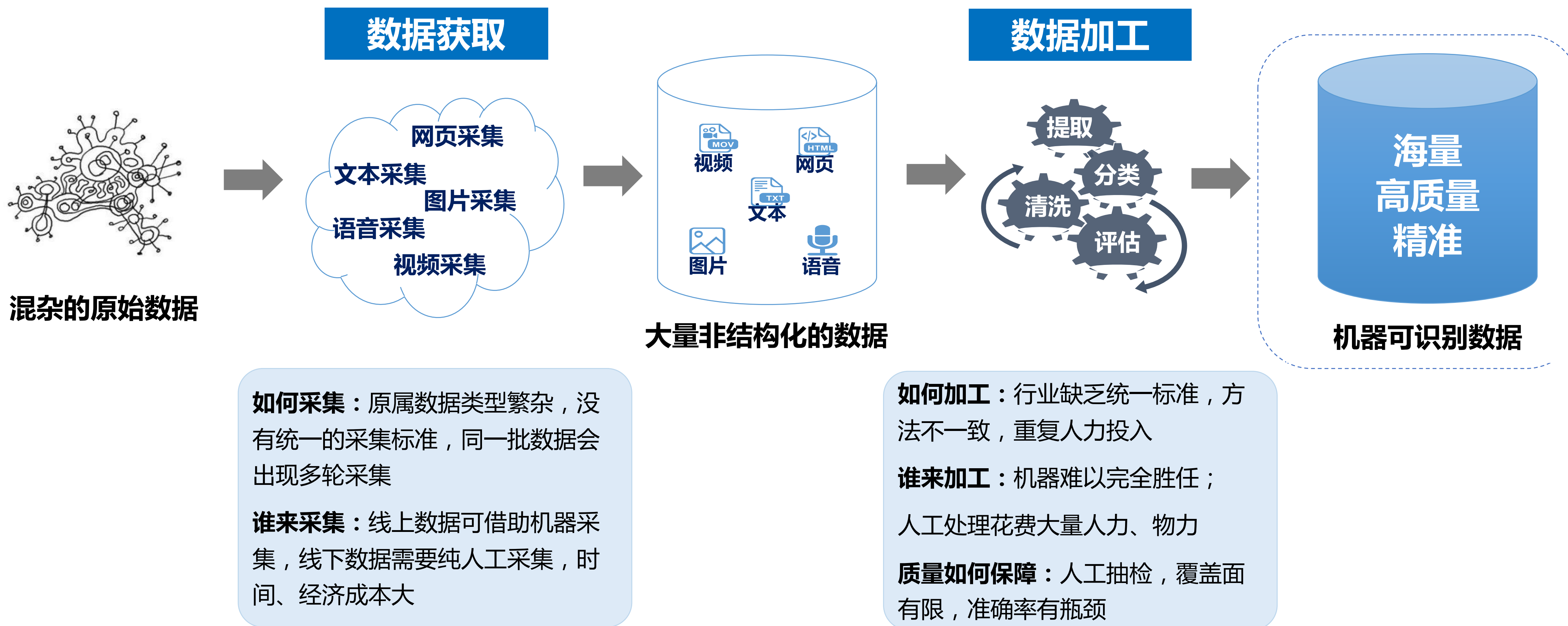


TABLE OF CONTENTES

人工智能行业现状

数据之于人工智能

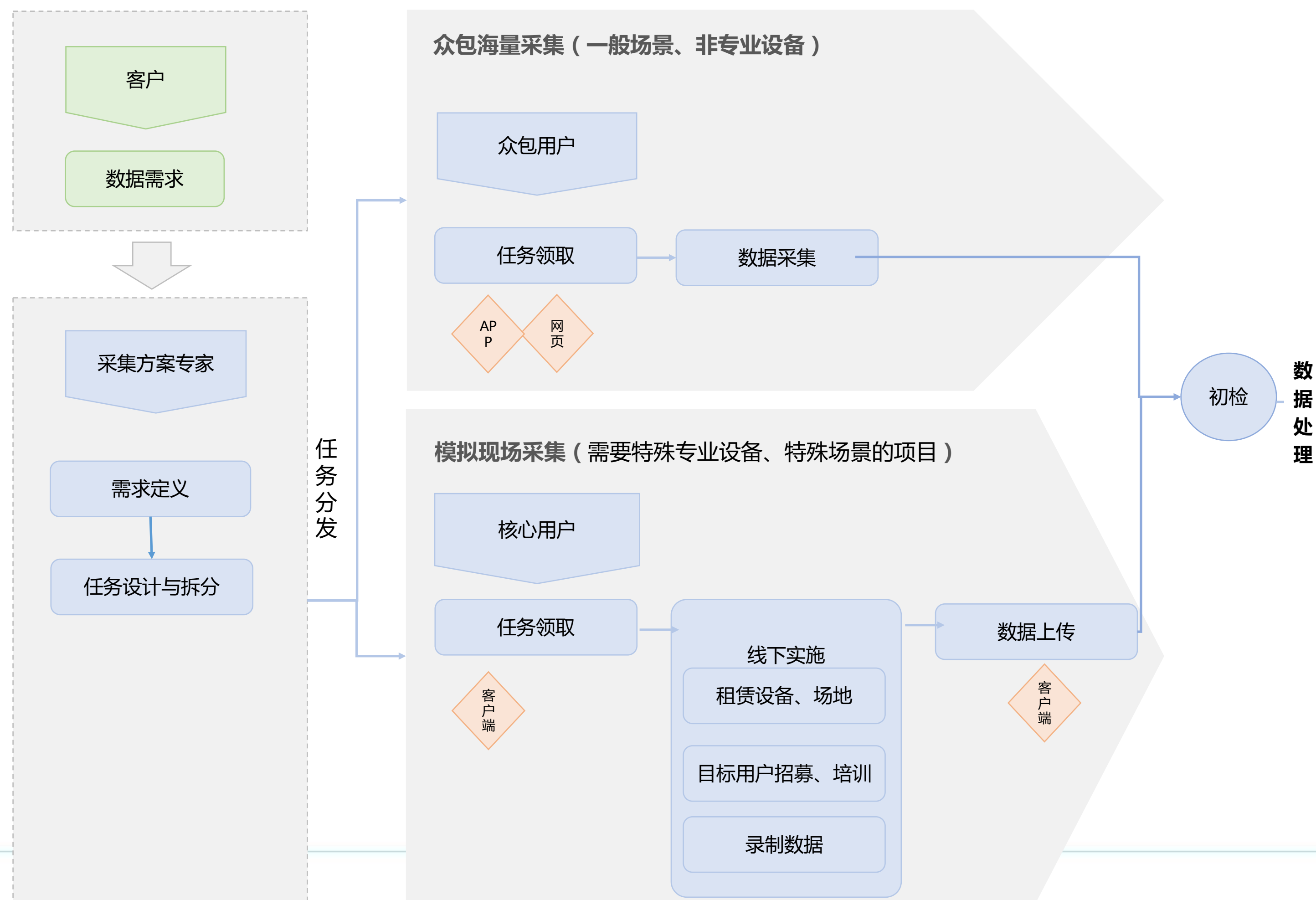
人工智能基础数据面临的难题

百度是如何应对的

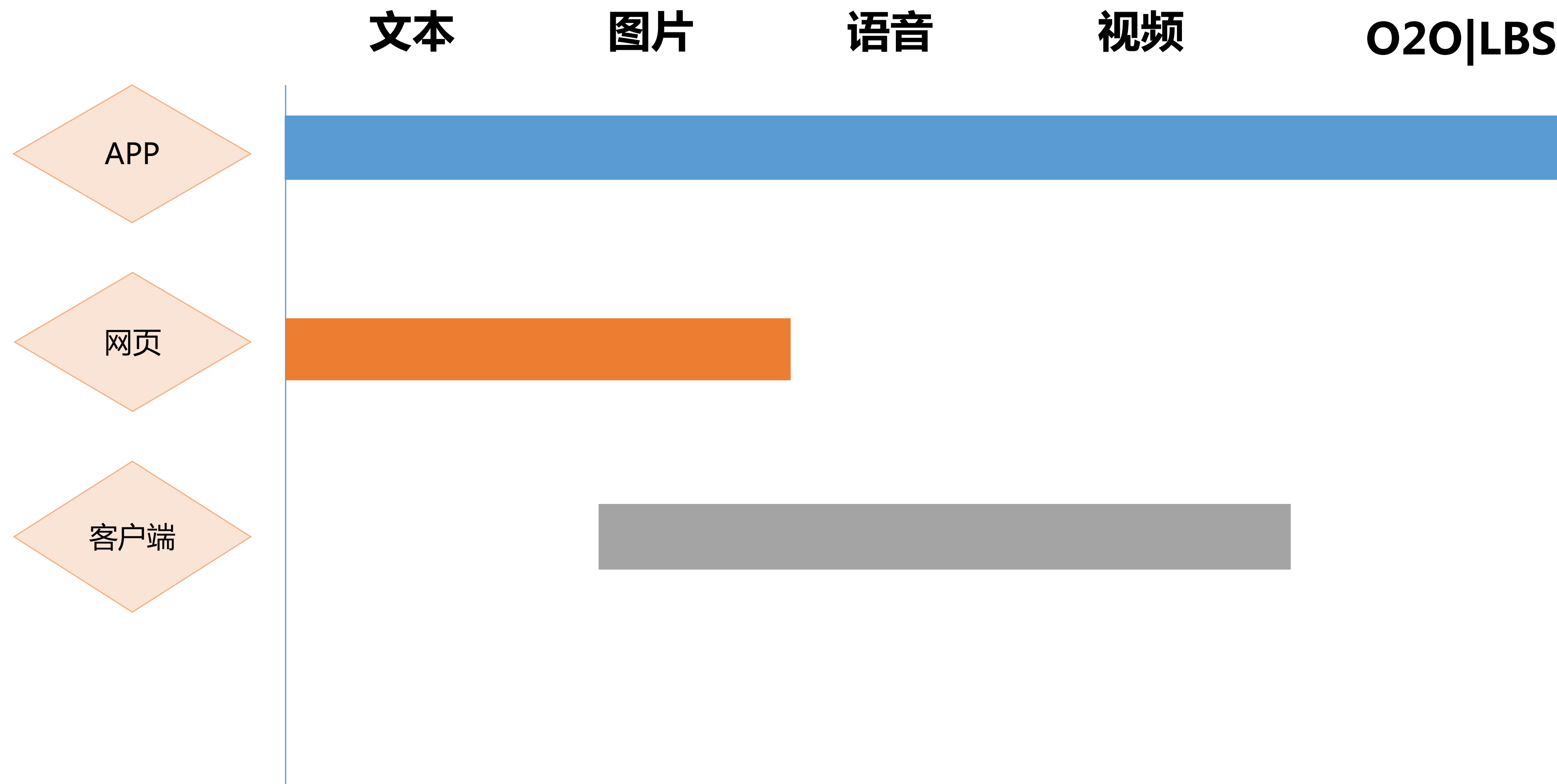
典型人工智能应用场景

百度数据产品

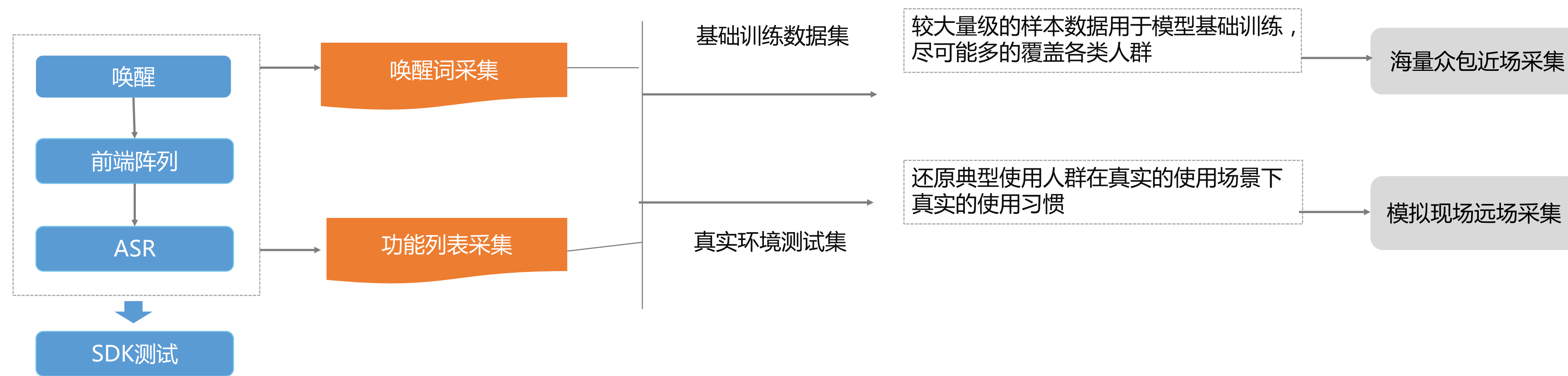
众包数据采集，为人工智能提供海量数据支持



采集场景工具化，全面覆盖各种数据类型



采集场景工具化，全面覆盖各种数据类型



近场数据：

众包用户 手机app自助采集。

采集能力：

累计完成超过5000小时，覆盖10w人的20种唤醒词数据。



近场数据：

核心用户辅助百度官方运营实施。

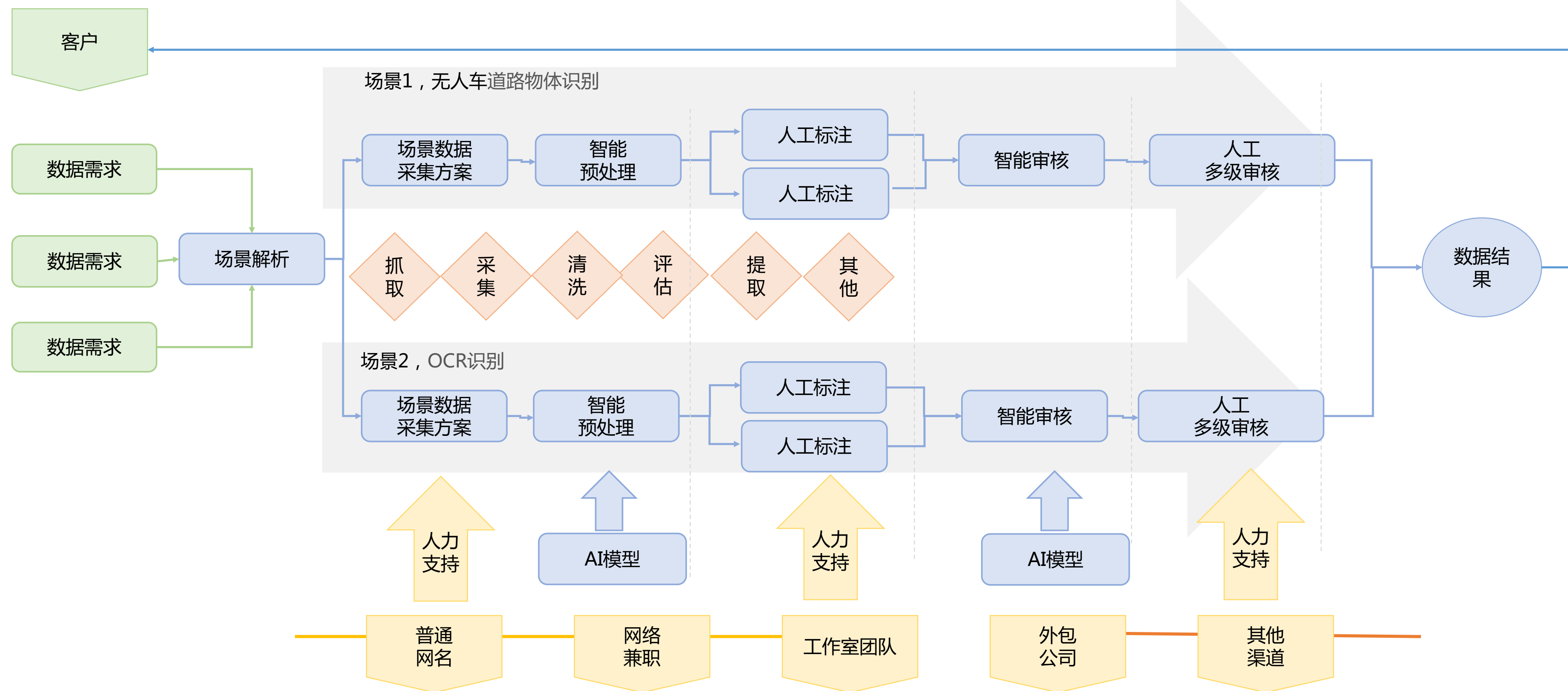
在真实使用场地，模拟真实使用情况，寻找设备真实受众人群的代表按照要求进行统一采集。

项目执行方案：

5种真实环境、3种真实距离

200人/天

链条化AI数据加工厂，为AI发展保驾护航



沉淀数据处理方法，建立数据处理规则

- 不完整数据
- 错误数据
- 冗余数据
- 数据标签化
- 垂类数据

数据清洗

1

数据评估

2

- 相关性评估
- 时效性评估
- 竞品评估
- 互联网，社交网络舆情
- 电子商务评论

- 关键词提取
- 网页内容提取
- 图片内容提取（OCR识别，人脸识别，物体识别等）

数据内容
获取

3

特殊信息
处理

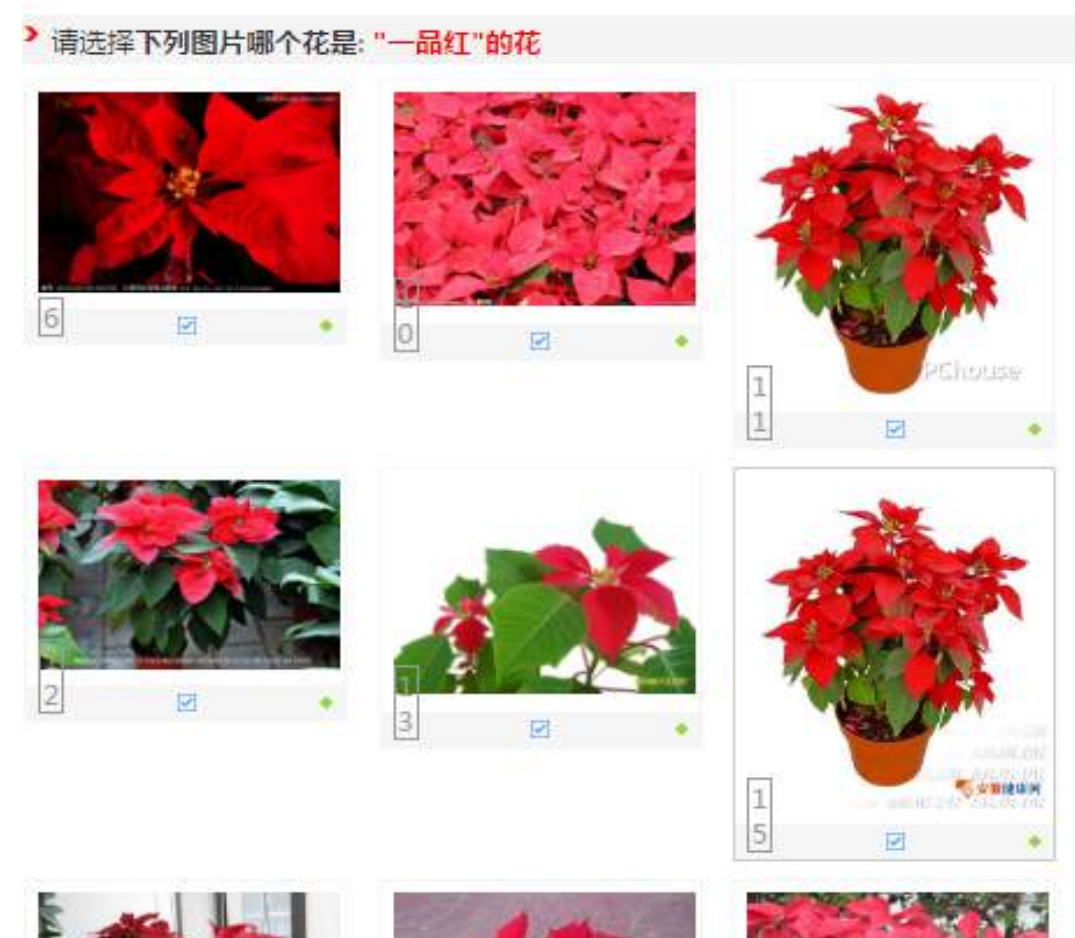
4

- 地图信息制作
- 语音转写
- 其他数据标注

固化数据处理工具：通用图片检测

通用图片检测类型涵盖商品、动物、植物、菜品、服装搭配、黄反、暴恐、建筑、素材等多种垂类。

1. 多图 vs. 单图；
2. 图+参考文字/参考图/搜索页面/参考链接/预识别结果/特定内部参考页面；
3. 多选题 vs. 单选题；
4. 题目类型：单选/多选/多级菜单选择/填写

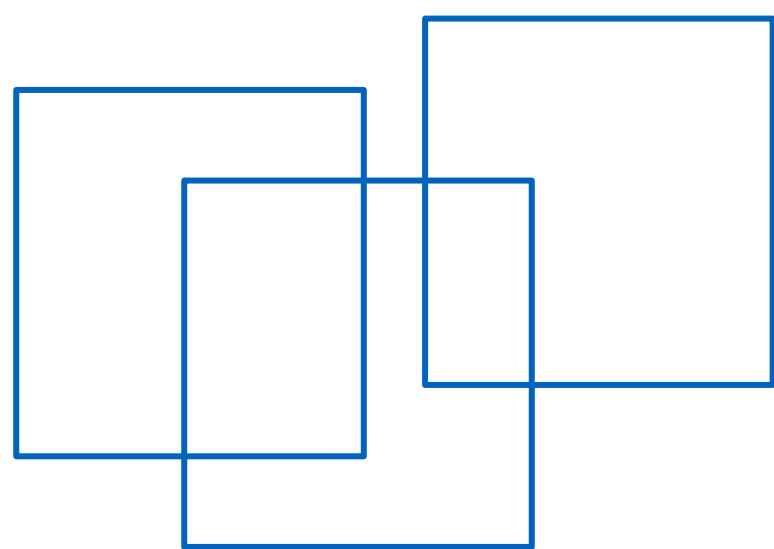


标注工具：目标框选类

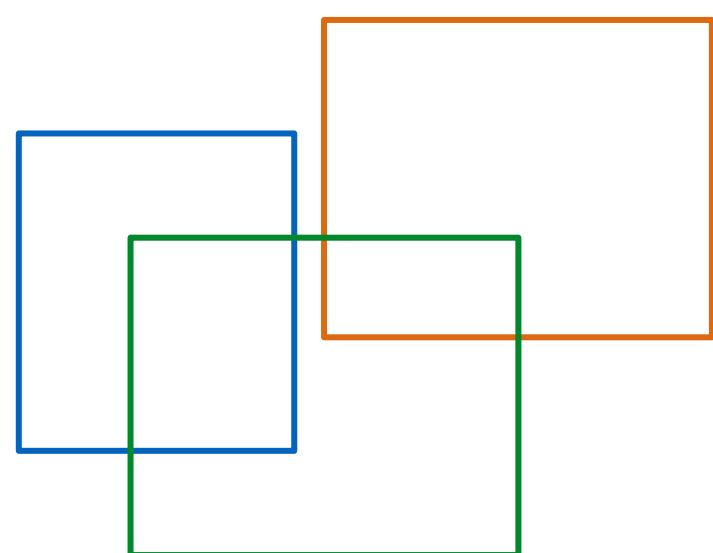
框选类能力涵盖：

普通矩形、分类矩形、普通多边形、分类多边形、区域填色、多级属性多边形、Parsing、点+线+区域复合检测

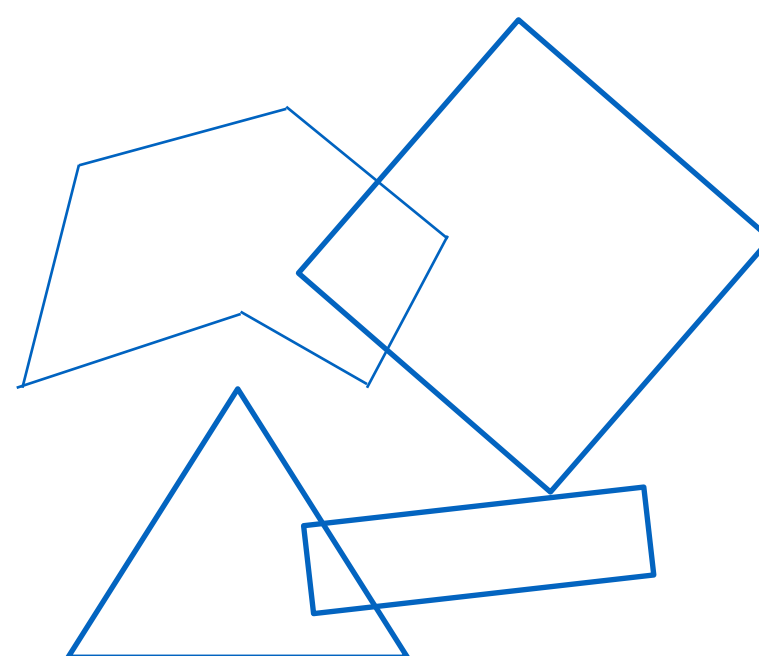
普通矩形框



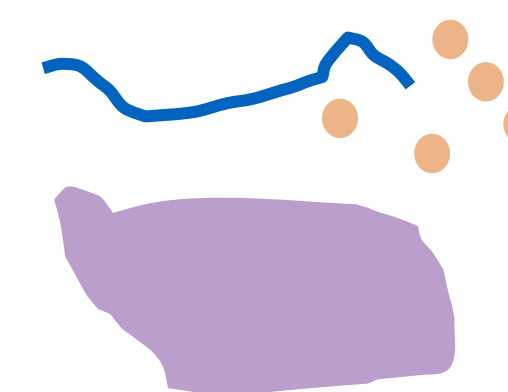
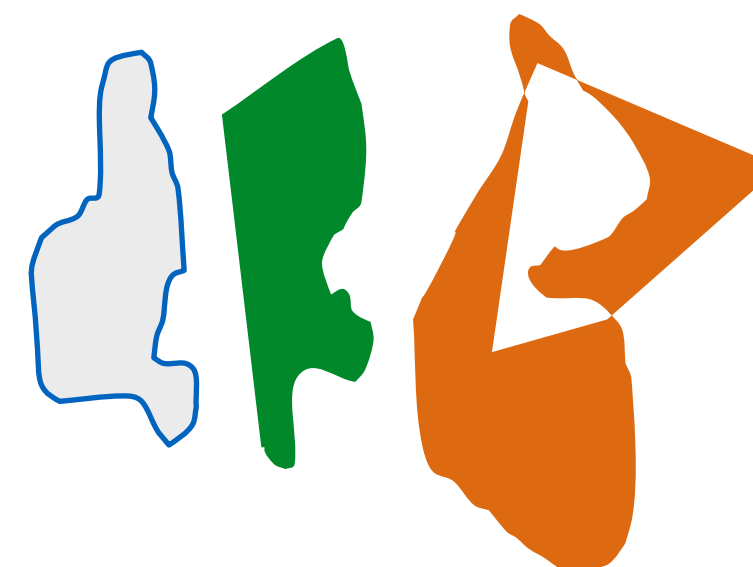
分类矩形



普通多边形



区域填色|多级属性多边形|Parsing 点+线+区域复合



标注工具：内容评估

用户行为画像

对“兴趣偏好”属性进行策略优化，通过第三方人工标注，通过用户人工贡献评价，评估策略优化后的标签准确率



蒜香花甲

2013-11-05 yutou

中国菜谱网
ChinaCaipu.com

用户搜索词：**蒜香花甲

分类标注结果：餐饮美食

分类标注辅助判断信息：饮食相关内容，如餐饮类别（如，火锅、茶酒饮料、烘焙等）、饮食口味（如，麻、辣、甜、生鲜等）、餐饮制作（如，烹饪、烹饪、烹饪）、餐饮活动（如，家庭聚会、情侣约会、朋友聚餐等）、餐饮地点（餐馆商圈、档次等）等内容。

我们给出的是一个人的网页搜索或浏览行为，请根据网页判断这个人的搜索是否符合我们的分类标注结果（例如用户搜索“十大最佳旅游地点”，我们的分类是“旅游出行”。那么这道题就选“是”）

是 不是



海贼VS火影 此游戏文件较大 (10.04 MB) 加载时间可能较长,请耐心等待... 游戏说明

4399小游戏

**海贼vs火影小游戏,在线玩,4399小游戏

我们给出的是一个人的网页搜索或浏览行为，请推测进行这些搜索和浏览的人，是不是对【游戏】相关内容感兴趣？

【游戏】是指：手机游戏，网页游戏，PC游戏，游戏机，游戏直播平台等

是 不是

我们给出的是一个人的网页搜索或浏览行为，请推测进行这些搜索和浏览的人，是不是对【网页游戏】相关内容感兴趣？

【网页游戏】是指：网页游戏类型、玩法、题材等

是 不是

要素提取

依据客户要求对文字内容或槽位进行提取并定位具体属性。



例句： 我是要成为海贼王队友的人！！

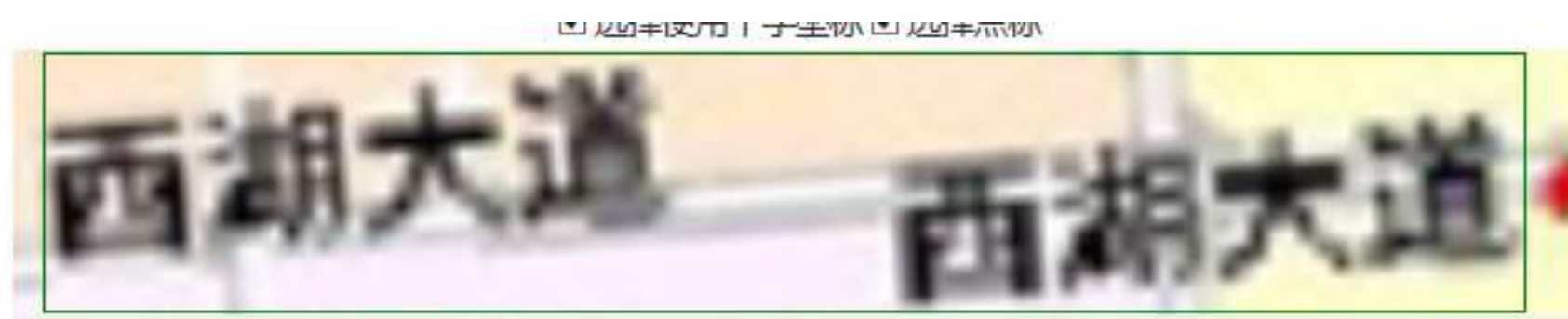
	提取内容	分类
* 提取分类1	海贼王	人名-人名 (PERSON.PERSON_NAME)

添加

[点击此处跳转对应搜索页面](#)

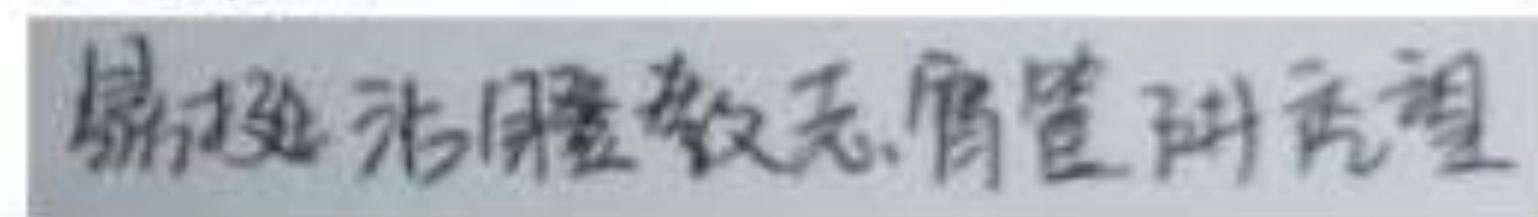
标注工具：图片&语音转写

- 1.进行多种语言OCR文字转写
- 2.进行多种口音的语音文字转写



选框 2 西湖大道西湖大道

参考图片：



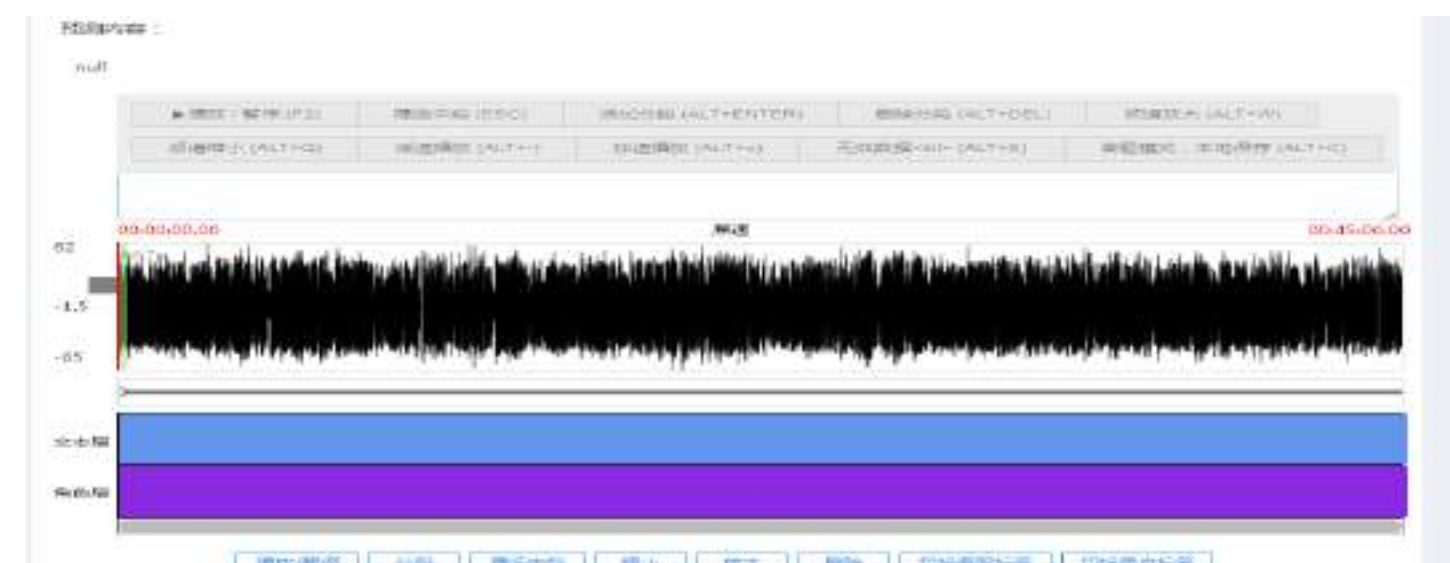
预识别文本：

鼎挺沾瞪教志宵筐阱秃望

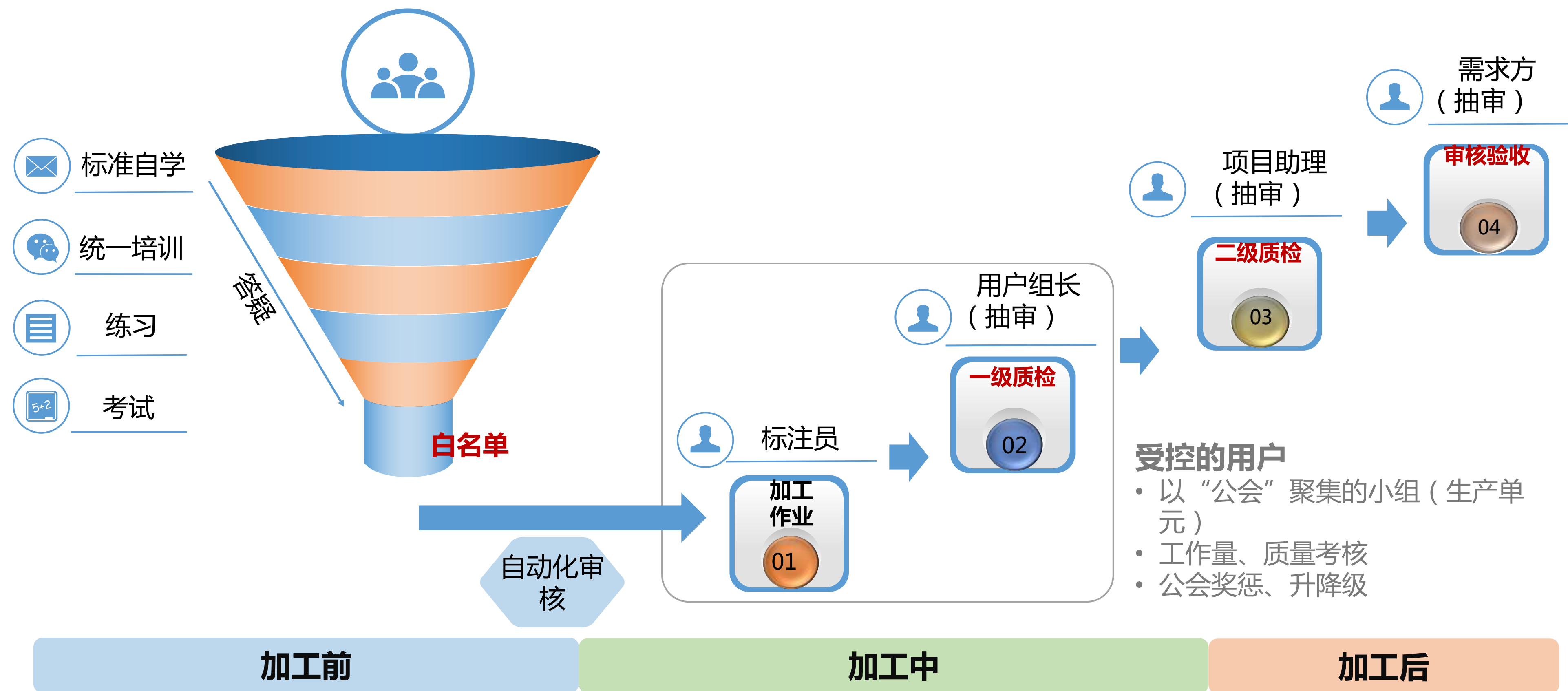
请仔细听下面的音频，按照规则将语音内容转写下来：

▶ 0:00

- > 1. 当前语音是否包含有效语音 包含有效语音且语言情况确定 包含有效语音但语言情况不确定 不包含有效语音
- > 2. 当前语音的噪声情况 安静 含噪音
- > 3. 语音内容
- > 4. 说话人类型 男声 女声 儿童
- > 5. 是否包含口音 否 是



多级质量管控，突破准确率瓶颈



智能标注：提升标注效率&降低标注成本

- 通过自动化审核规则&算法，抽样概率较大的badcase，降低审核成本

成本

= 题目数 × 拟合人数 × 标注单价 + 审核成本

标注前：数据预处理

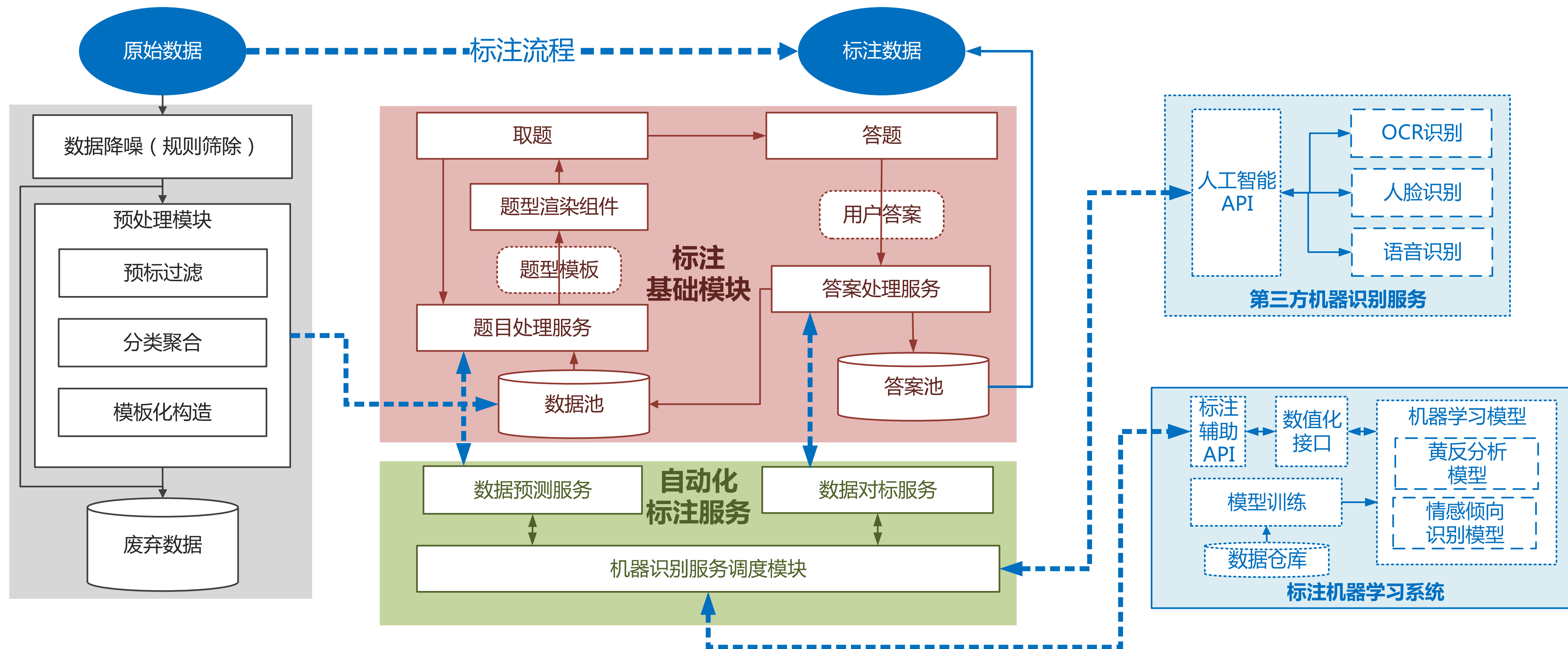
- 按标注难度区分任务
- 分类原始数据
- 去除无效的空白数据

- 机器实时对标人工标注结果，减少需要拟合的人工标注数量
- 机器标注结果与普通用户标注结果一起参与审核和验收

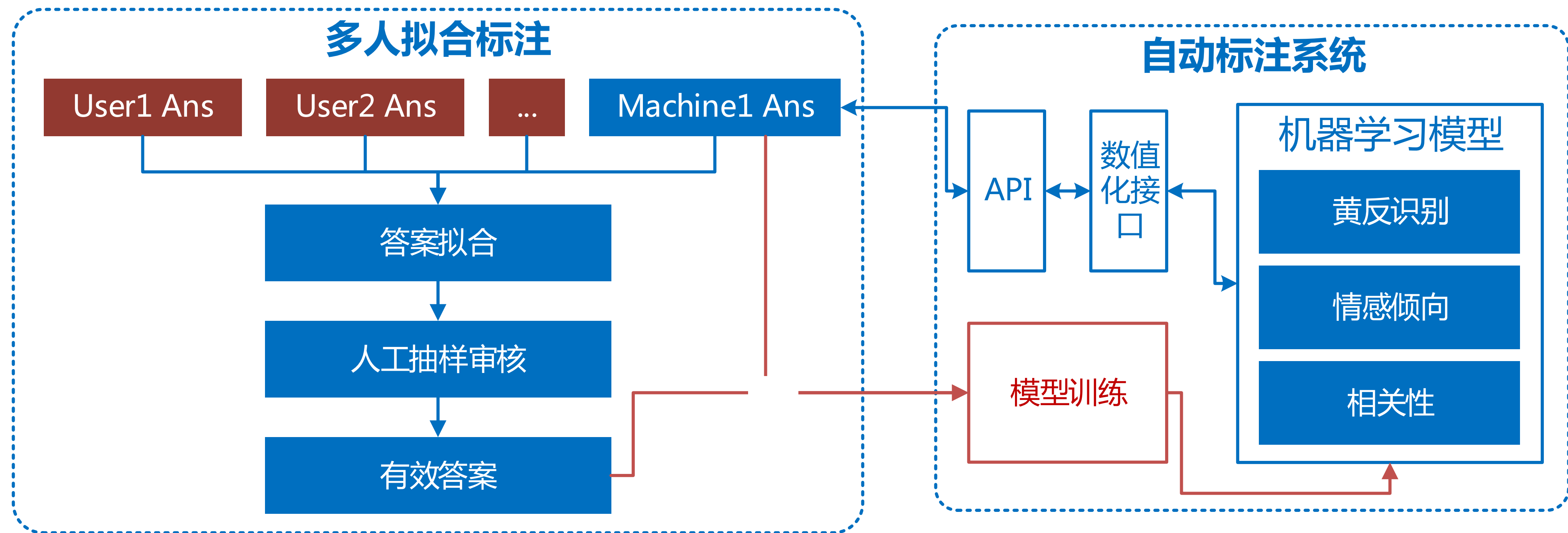
标注中：辅助标注

- 自动化生成预标注结果
- 辅助用户标注，简化操作步骤
- 尽量进行内容修改，降低标注难度

智能标注：架构总体设计



智能标注：标注中（人工+机器对标）



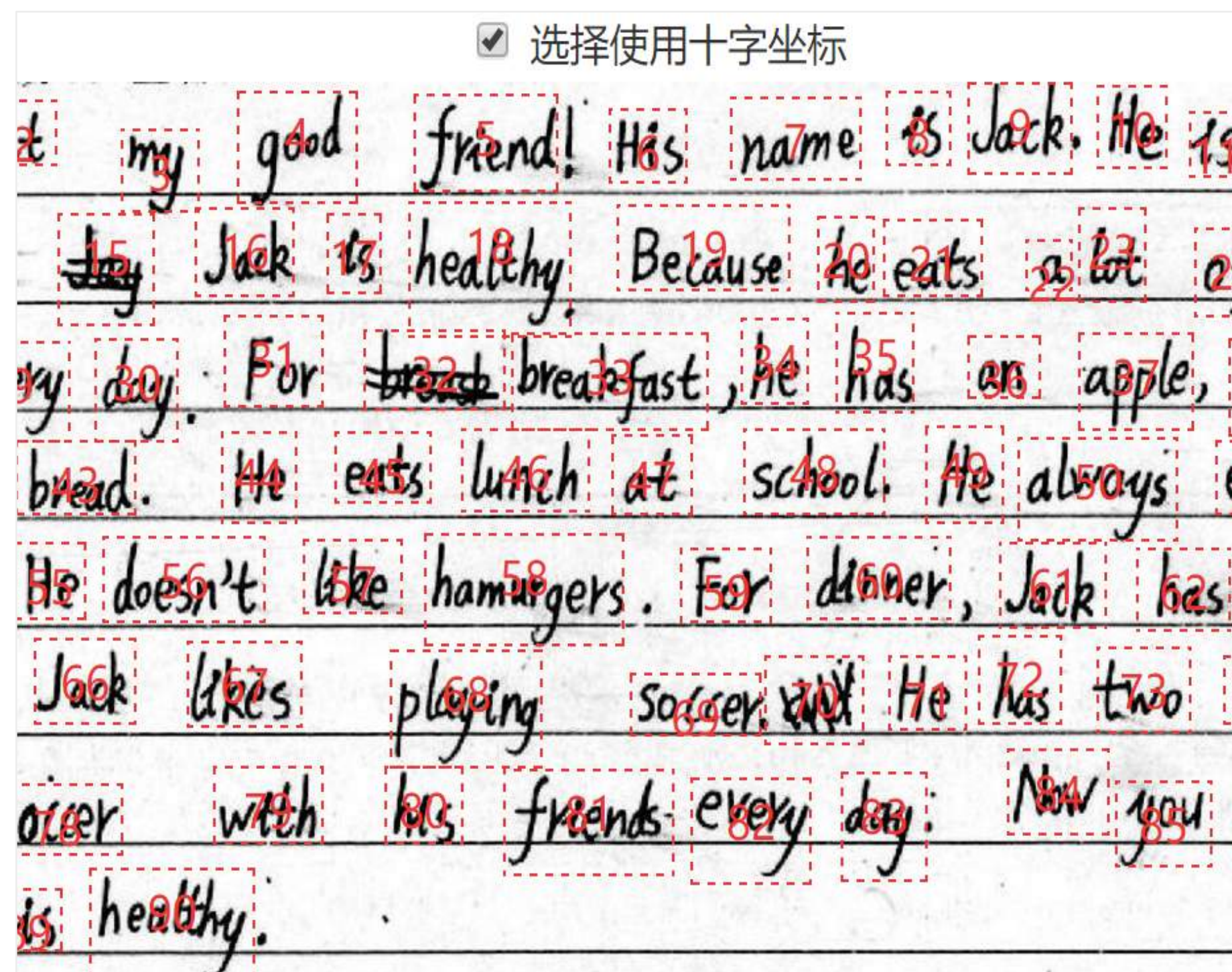
智能标注：标注中（辅助标注）

标注特点

- OCR模式识别技术已经比较成熟，进一步优化模型需要更多的bad case，因此需要更好地获取算法与人工标注存在较大diff的数据反馈到模型训练中
- OCR图片识别产品化，可直接调用接口服务对图片进行预识别

实时预测校正标注

- OCR识别接口的预标注准确率较高，用户只需要对少量bad case进行调整
- 单价构成： $P_m * N_m$
- 记录预测数据修改，大量修改的数据可作为训练数据



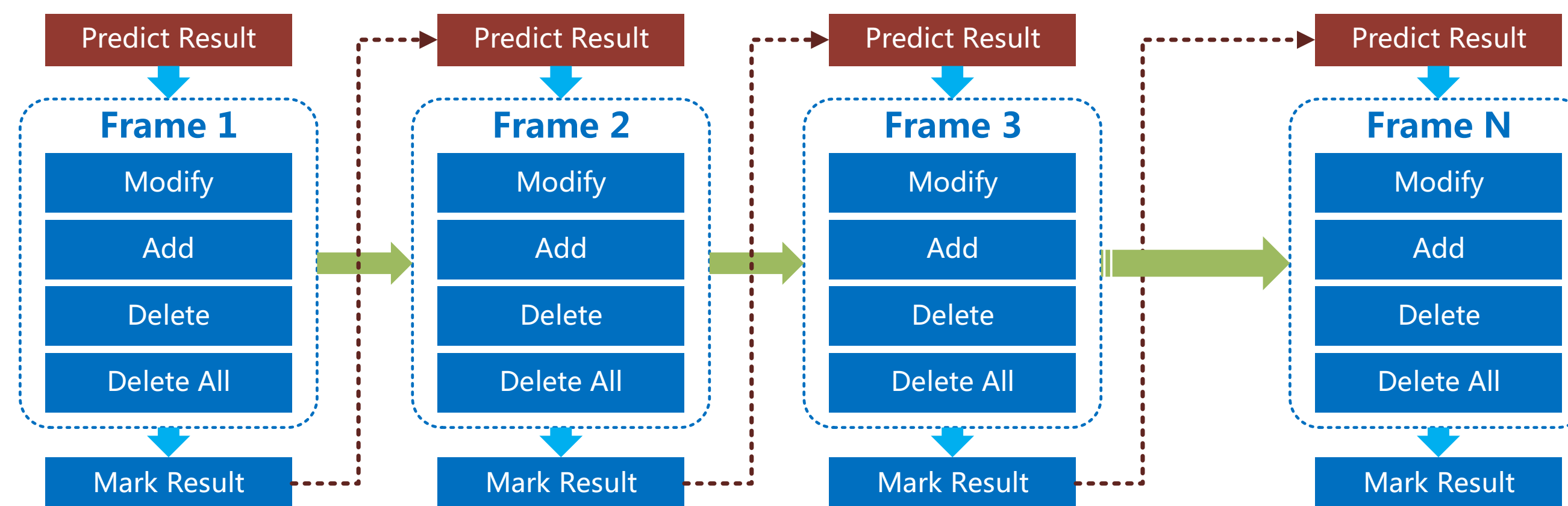
智能标注：标注中（辅助标注）

标注特点

- 投放的标注图片一般为连续帧，前后相邻帧中需要标注的语义类型、位置、大小等属性基本一致
- 通过算法进行预测数据生成准确性低，且无法达到实时获取

串行连续预测标注

- 连续帧图片相似度达到80%以上



智能标注：标注后（自动校验）



图片连续帧自动化审核

- 结果跳变检测：图片连续帧来说，每一帧的数据结果之间应该不存在跳变，如在第三帧的结果中突然少了很多第二帧的车辆，这种应该就需要被标记为异常结果，进而需要着重来进行查看

根据数据需求类型，覆盖更多实际应用场景



人像识别

多角度自拍
跨年龄段
暗光人脸
亲子全家福
人脸打点

语音识别

唤醒词语料
客服语音
普通话文本转录
中英文混读
方言 - 粤语
方言 - 四川话

OCR识别

驾驶证图片
名片图片
商标LOGO
彩票图片
医疗单据图片

无人驾驶

红绿灯图片
道路障碍物
交通行驶区域
道路分界线
交通路面边界
泰国车牌

模式识别

手部图片
时尚服装
推荐菜品图片
汽车外观图片
动物图像
花卉图像

TABLE OF CONTENTES

人工智能行业现状

数据之于人工智能

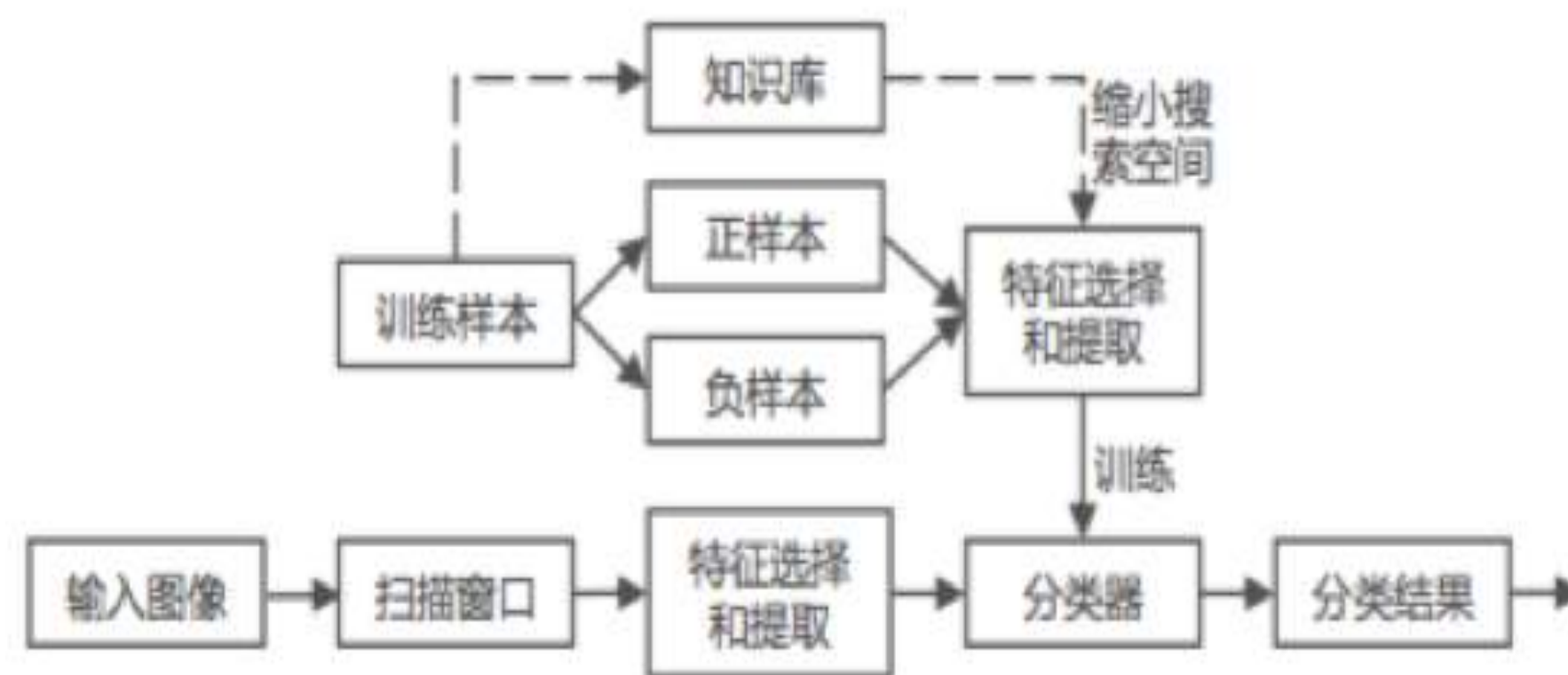
人工智能基础数据面临的难题

百度是如何应对的

典型人工智能应用场景

百度数据产品

计算机视觉数据解决方案



数据采集

根据实际计算机识别模型的要求，采集相应的图片、视频内容。

数据加工

将采集内容加工处理：标注关键点定位、提取特征信息打标签。

模型训练

将原始数据和特征标签数据提交到学习平台进行训练，提高识别精度

识别反馈

进行多次的迭代训练，最终计算机给予相应的识别反馈信息。

计算机视觉应用下的数据方案

特殊场景人脸图像数据



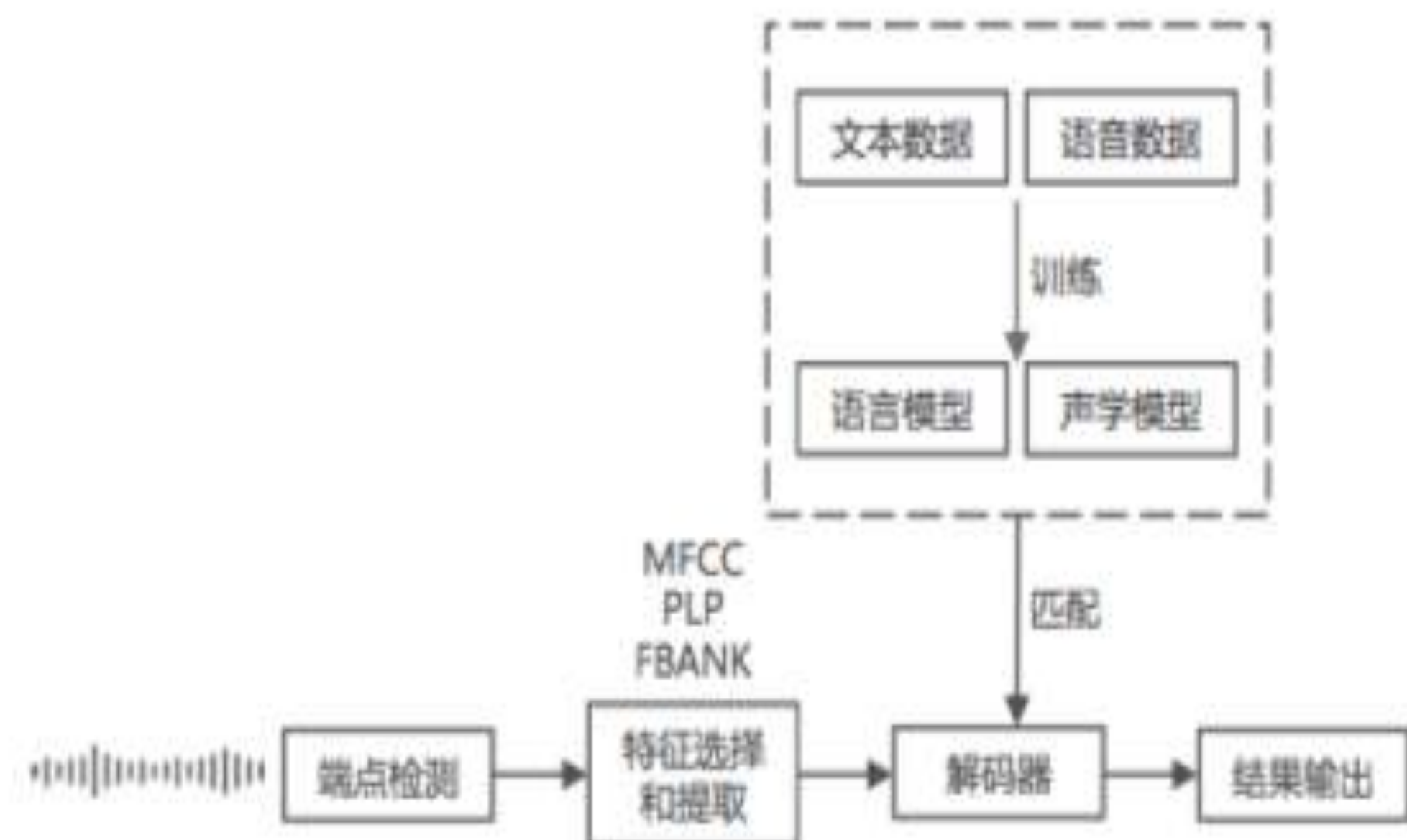
特殊要求人脸图像采集

- 采集指定条件下的人脸图像照片
- 通过手机自带相机拍摄
- 正常、暗光、微光多条件拍摄
- 口罩、墨镜、帽子多遮挡条件拍摄

人脸图像标注

- 人脸检测标注：人脸位置框选
- 人脸关键点标注：人脸5点-72点标注

语音识别数据解决方案



唤醒词、中英文语料、
方言语音识别。

语音
识别

语义
理解

多轮对话：上下文可随时
打断,加入语境分析功能

机器翻译、实时同声传译

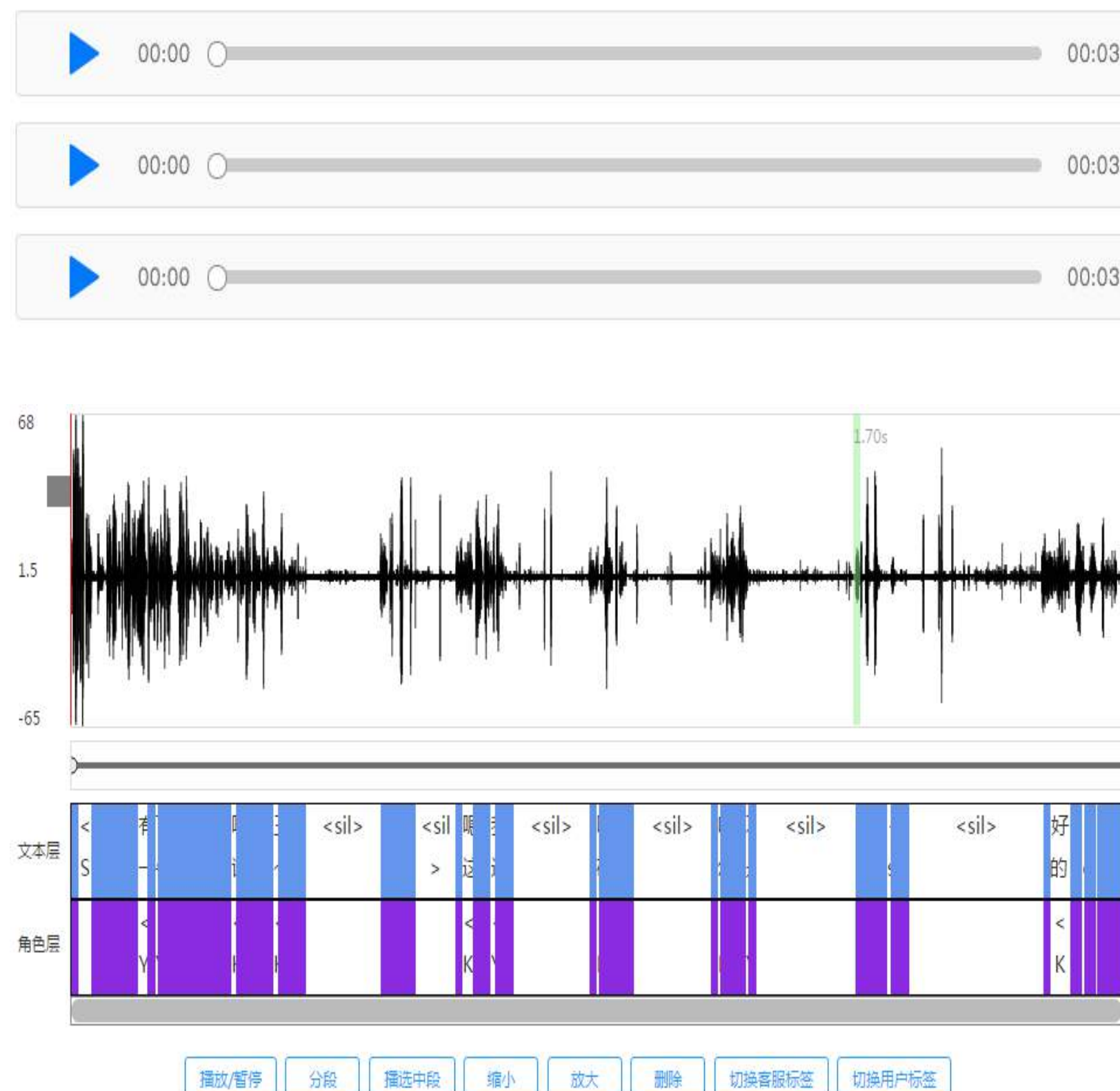
自然语言
生成

语音
合成

中文语音合成、中英文
混合语音合成

语音识别数据解决方案

汉语方言语音数据



汉语方言语音数据采集

- 采集指定地区的汉语方言数据
- 通过手机自带麦克录制
- 四川话 / 上海话 / 湖南话等8种方言
- 安静 / 吵闹环境录制

语音数据转写标注

- 中文方言、普通话
- 转写准确率98%

TABLE OF CONTENTES

人工智能行业现状

数据之于人工智能

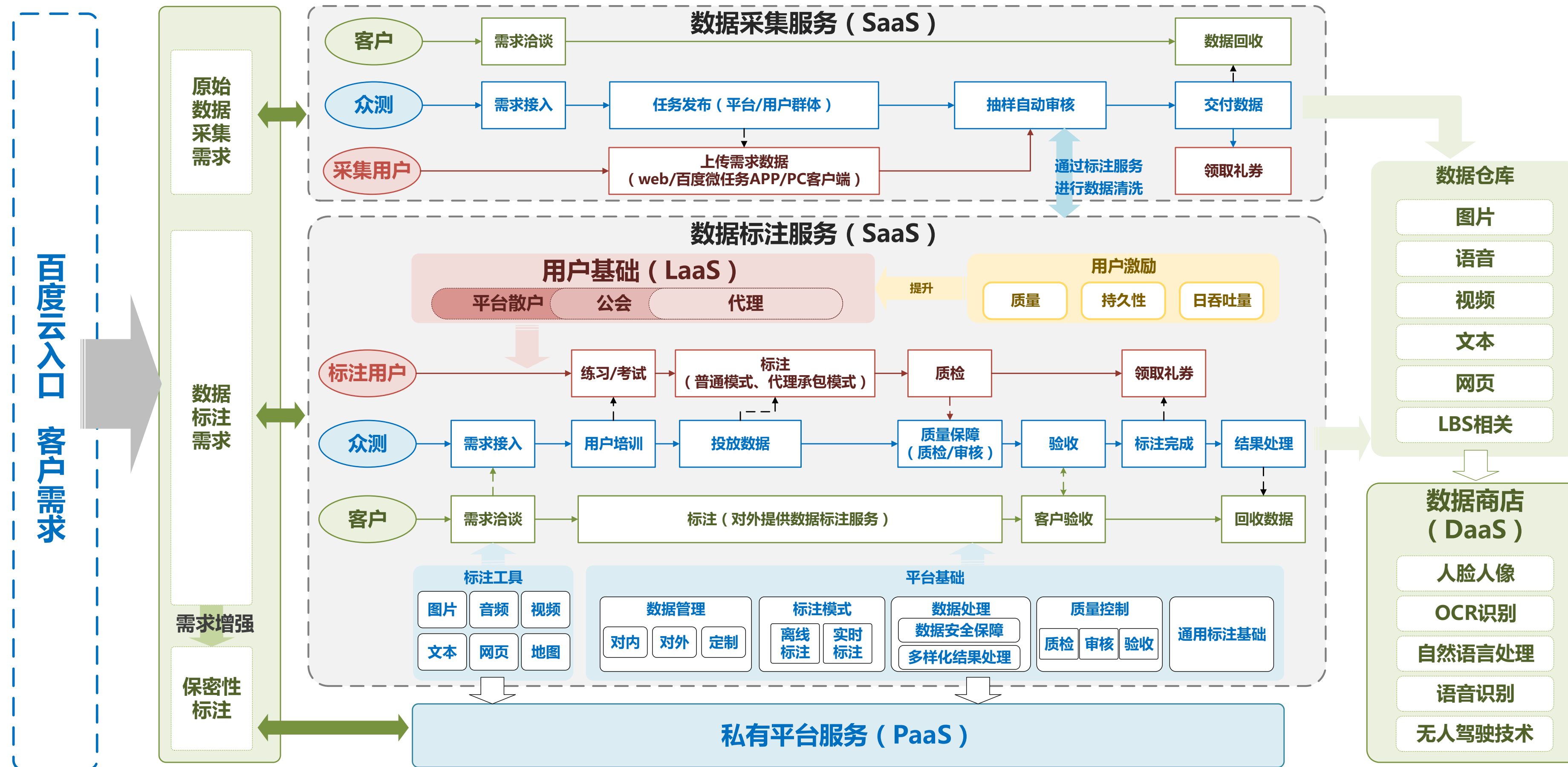
人工智能基础数据面临的难题

百度是如何应对的

典型人工智能应用场景

百度数据产品

百度数据产品矩阵



Thanks!

