

AI时代的数据解决方案

- 百度数据众包产品架构师 沈健
- <http://zhongbao.baidu.com>

人工智能行业现状



人工智能进入公众视野

人工智能的
强大能力已
被证明

● ALPHAGO
00:10:29



● LEE SEDOL
00:01:00

人工智能2017大事记

AlphaGo3:0战胜柯洁，
DeepMind 创始人宣布
AlphaGo “退役”

围棋界再无敌手之后，Alpghgo的下一个目标是“征服”哪里？



国务院印发新一代人工智能
发展规划 中国将人工智能上
升为国家战略

《规划》提出了六方面的重点任务和一系列保障措施，国家层面为AI奠定好的基调。



类人机器人Sophia亮相
《早安英国》

人工安卓智能机器人Sophia与她的发明人David Hanson博士共同做客《早安英国》节目，接受主持人的现场大拷问。



百度All in AI，发布了
DuerOS 和 开源自动驾驶
系统Apollo

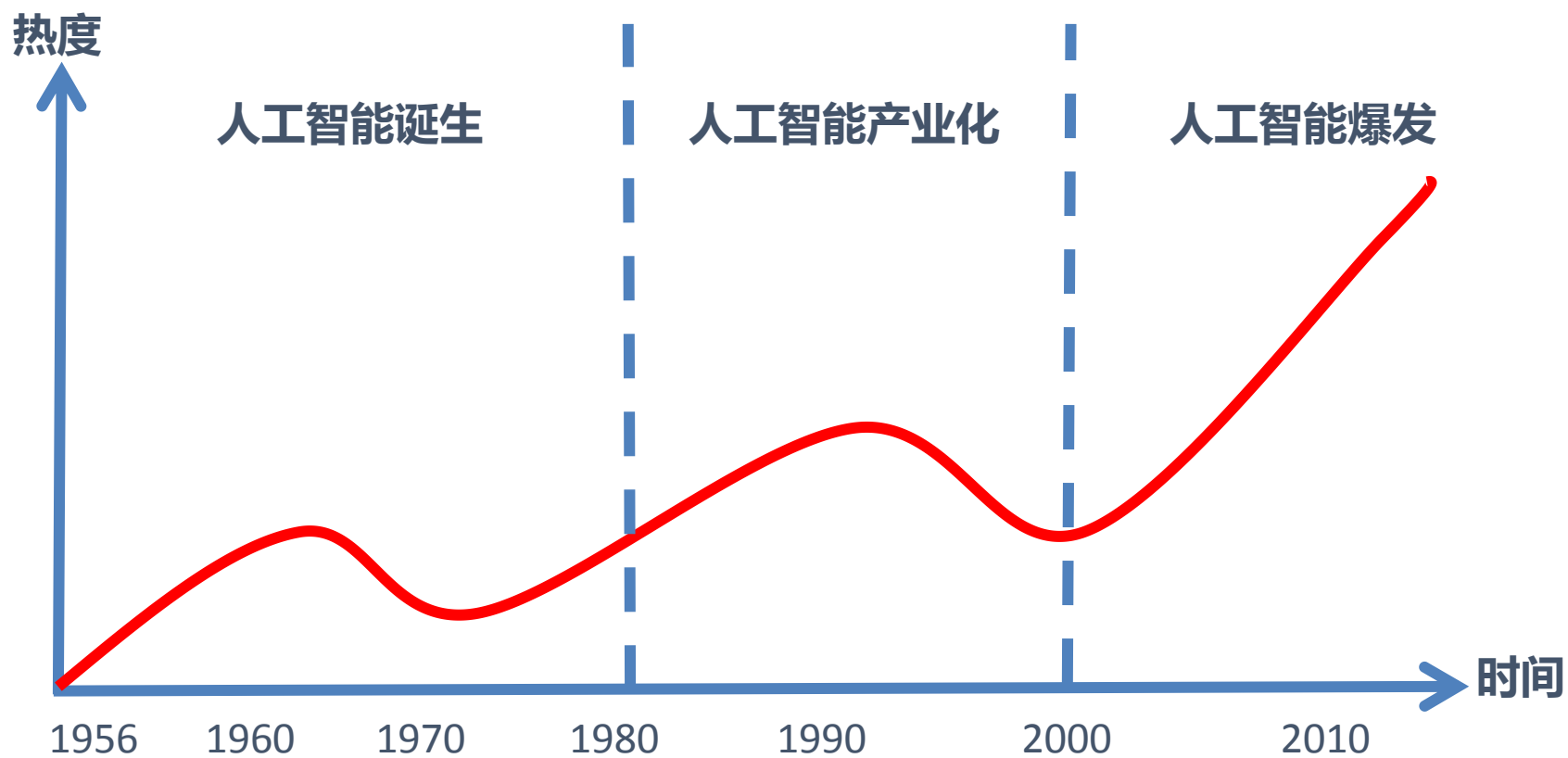
All in AI表明了决心，也为百度贴上了一个新标签——“人工智能公司”



NVIDIA发布地表最强
GPU：PCI-E Tesla V100

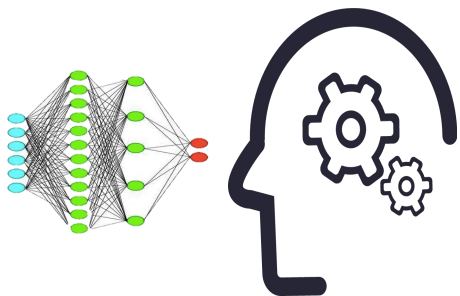
只需要几张V100的算力就能够与当前的各国精心打造的超级计算机的算力相当。

人工智能发展历程



人工智能爆发的三大因素

深度学习



高性能运算



大数据



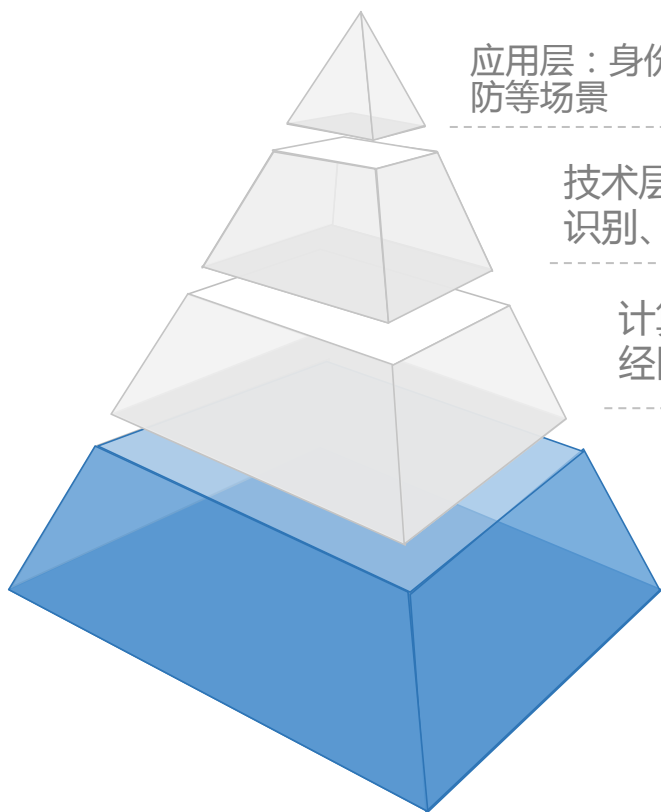
算法是核心，**计算**、**数据**是基础

数据之于人工智能



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！



应用层：身份识别、语音助手、机器人、智能安防等场景

技术层：机器学习、深度学习、语音识别、图像识别、人脸识别、NLP等。

计算能力：大数据、云计算、神经网络芯片等计算能力提供商

数据：身份信息、医疗、购物、交通出行等各行业、各场景的一手数据。

方言语音数据

场景语音数据

语音文本数据

百科、音乐、游戏、电影等不同领域的词汇、属性及关系数据

社交网络文本数据

新闻媒体舆情数据

社区、论坛数据

多语种文本数据

语义标注数据

人脸图像数据

字符图像数据

物体图像数据

物体属性数据

物体行为数据

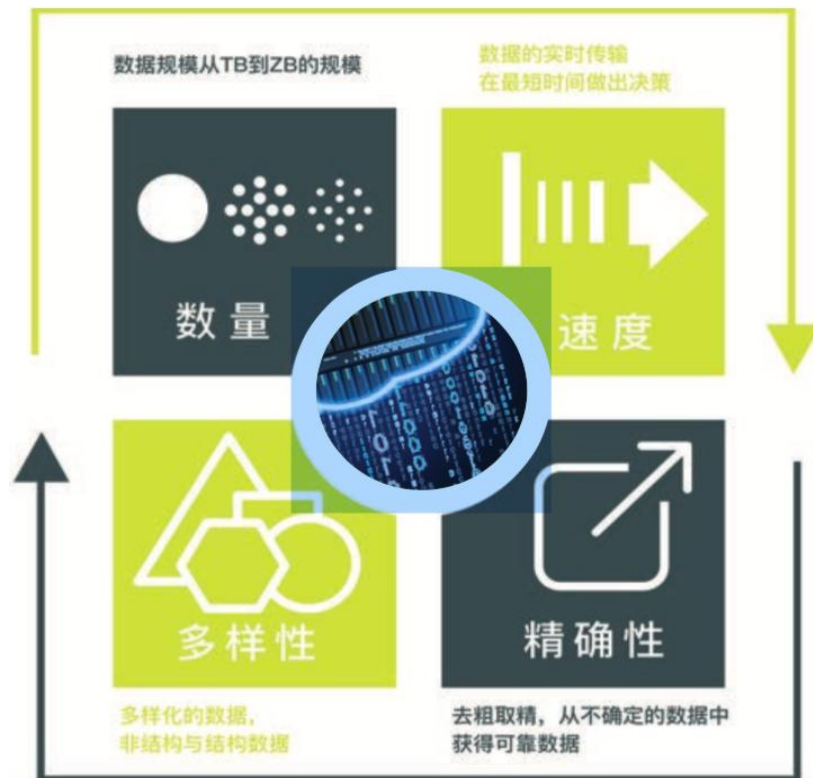
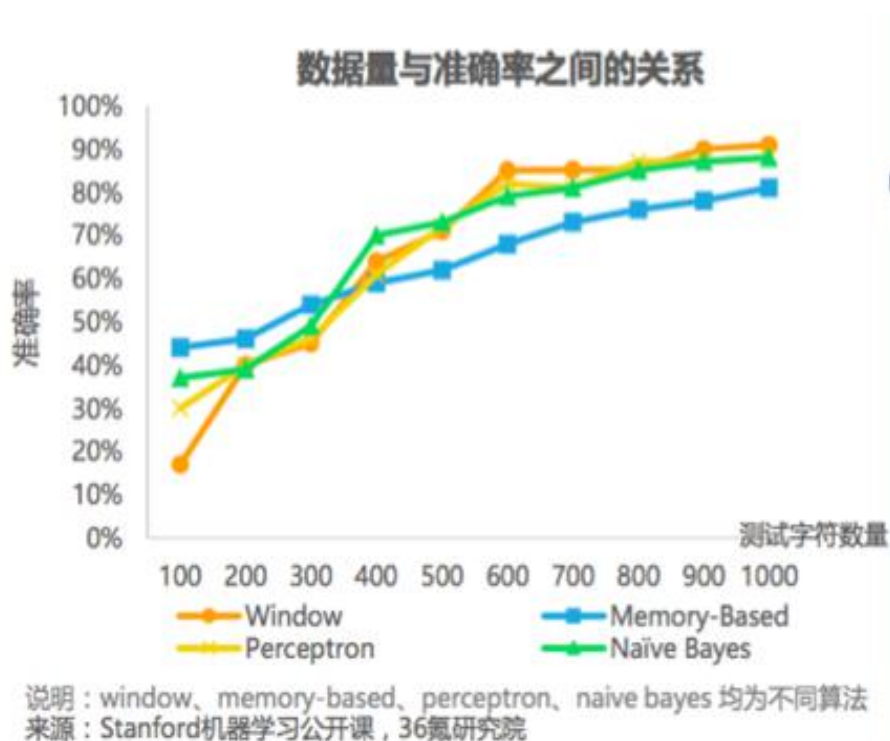
数据样本与算法模型



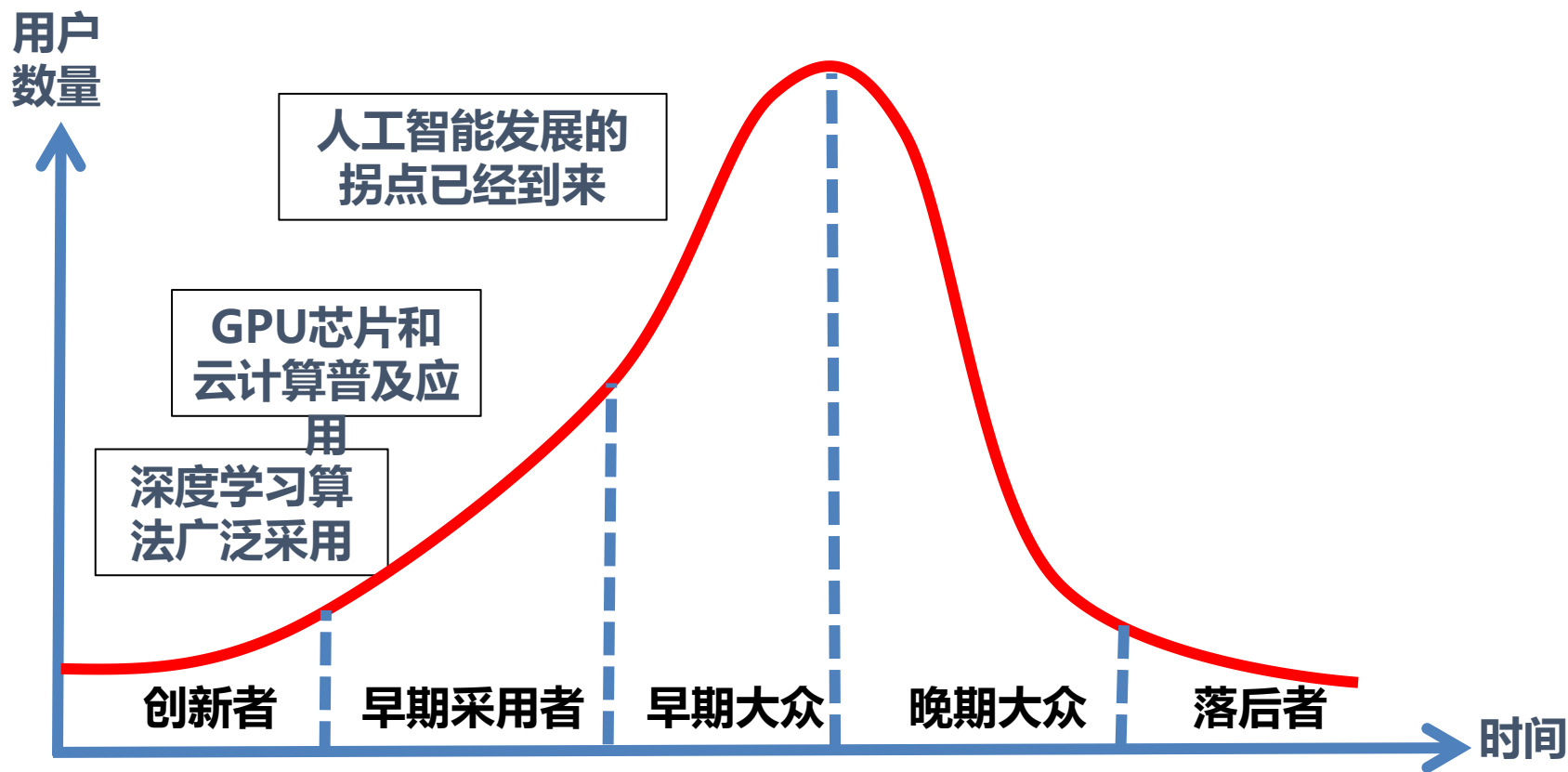
人工智能需要通过**大量的数据样本**来“训练”自己，才能不断提升输出结果的质量。

有时候，数据真的可以秒杀算法

有时候谁能够取胜，并不取决于谁拥有更好的算法模型，而是看谁掌握着**更多、更好**的数据资源。



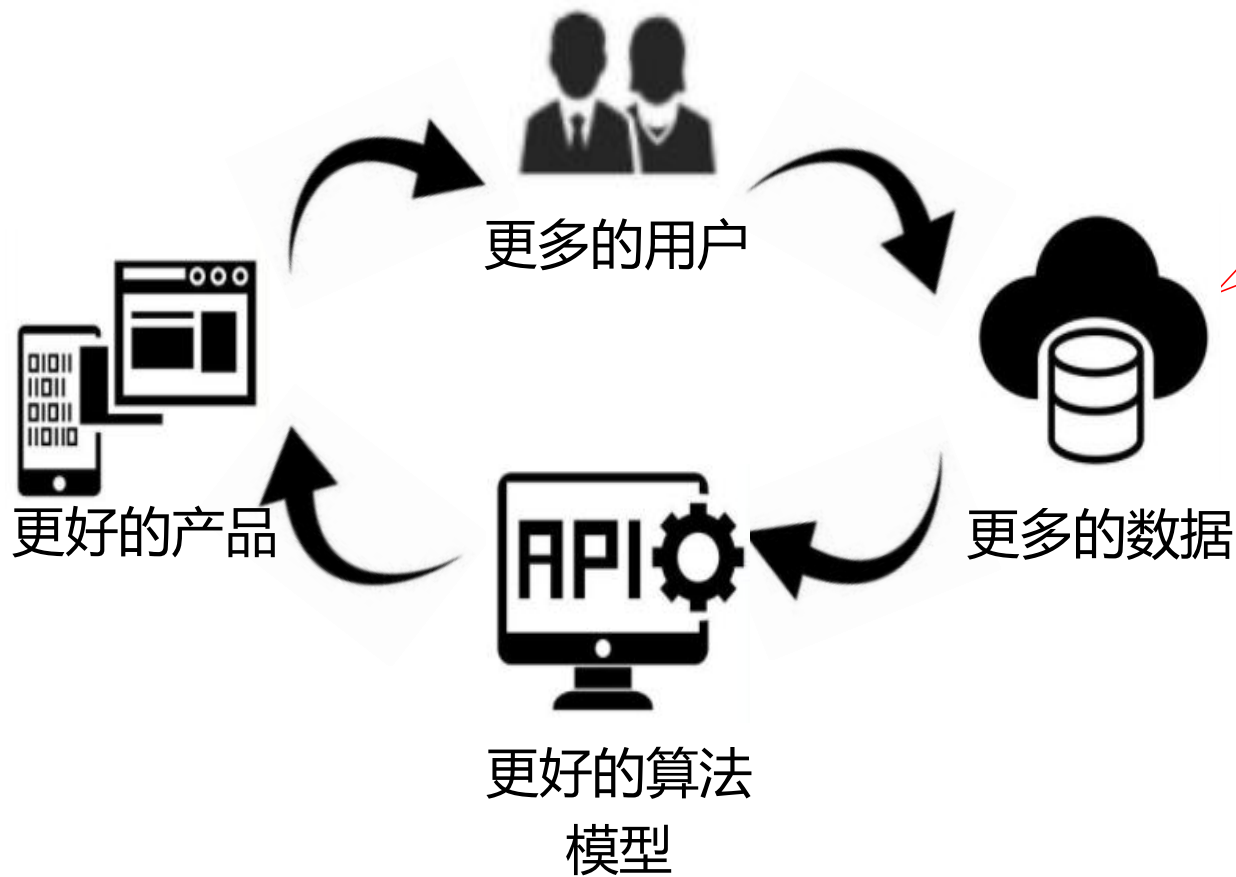
快人一步抢占先机，数据竞赛“质&量”取胜



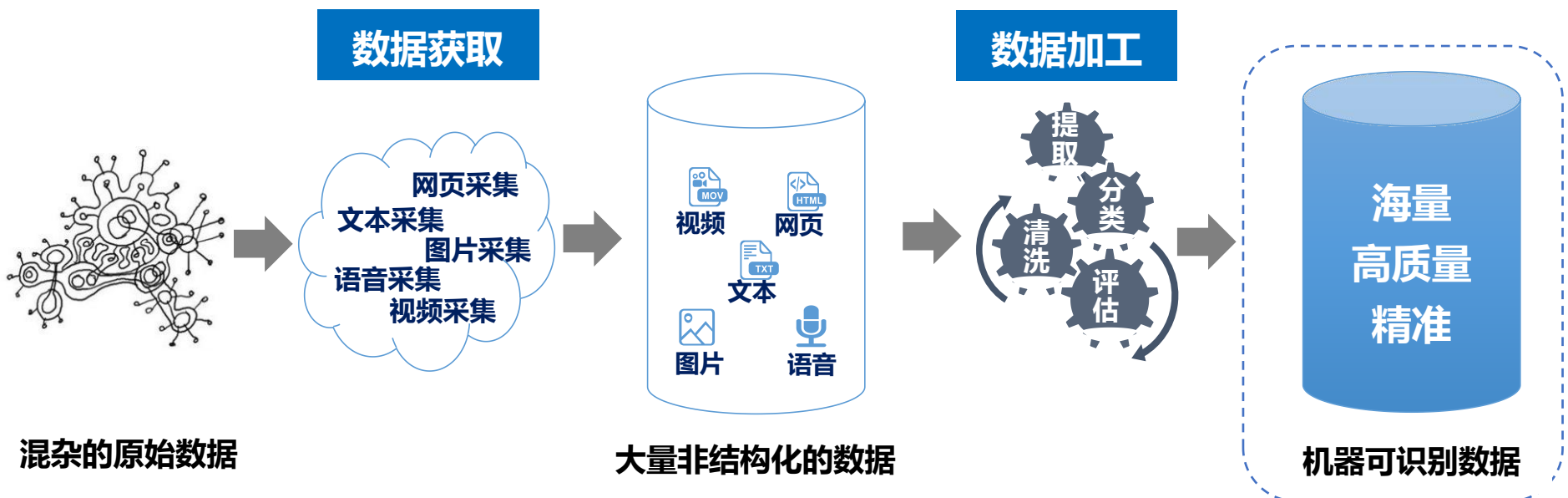
人工智能基础数据面临的难题



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



如何采集：原属数据类型繁杂，没有统一的采集标准，同一批数据会出现多轮采集

谁来采集：线上数据可借助机器采集，线下数据需要纯人工采集，时间、经济成本大

如何加工：行业缺乏统一标准，方法不一致，重复人力投入

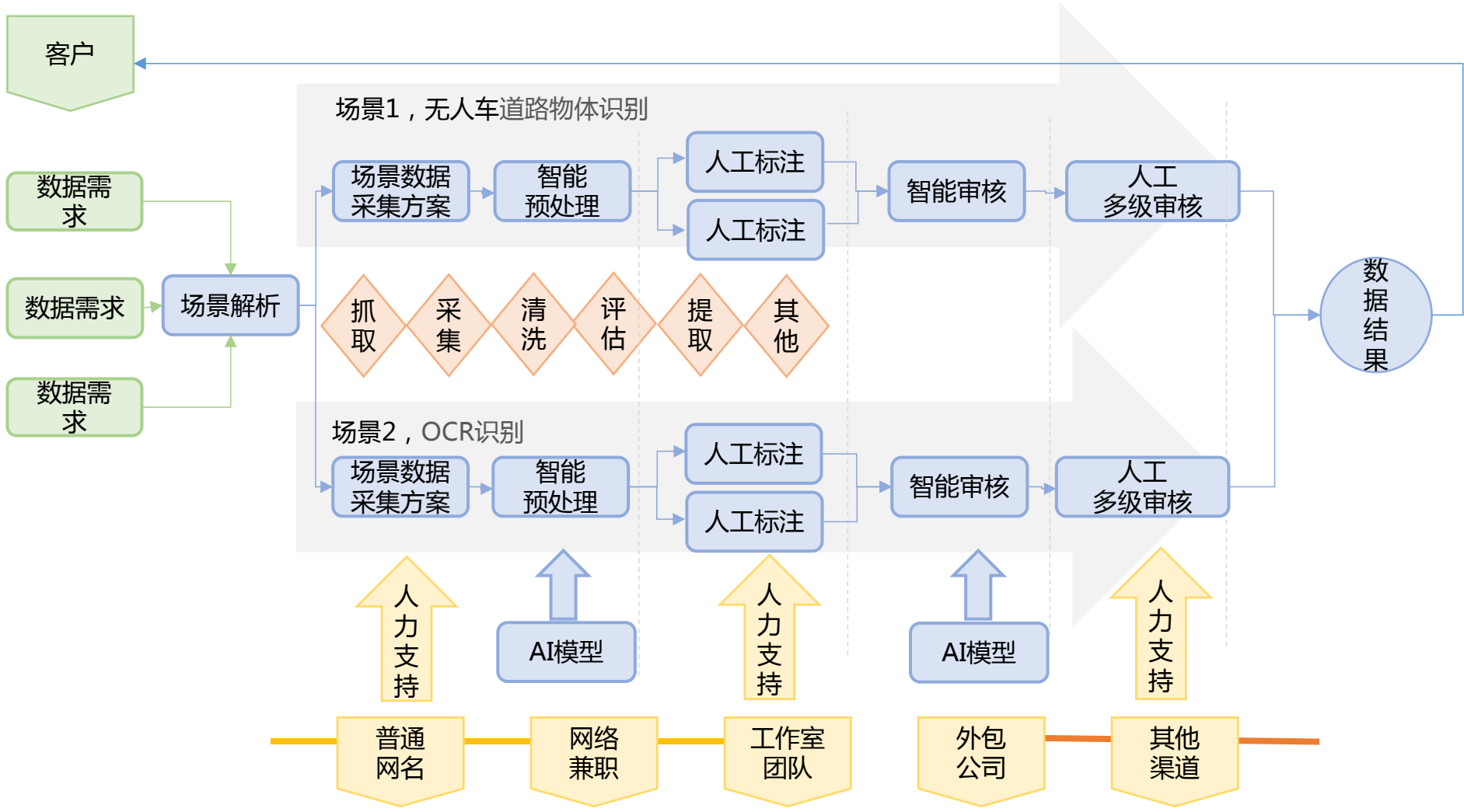
谁来加工：机器难以完全胜任；人工处理花费大量人力、物力

质量如何保障：人工抽检，覆盖面有限，准确率有瓶颈

百度是如何应对的



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

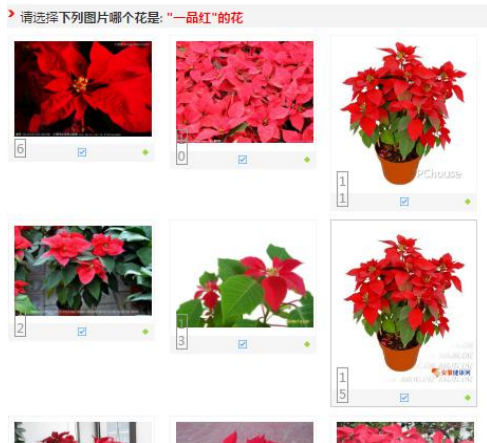
4

- 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. 多题 vs. 单题；
4. 题目类型：单选/多选/多级菜单选择/填写

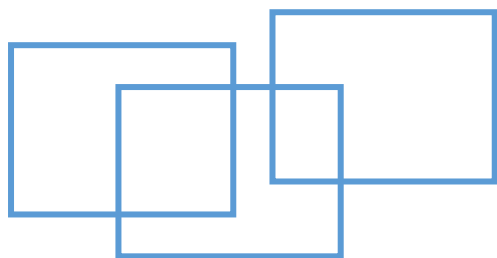


标注工具——目标框选类

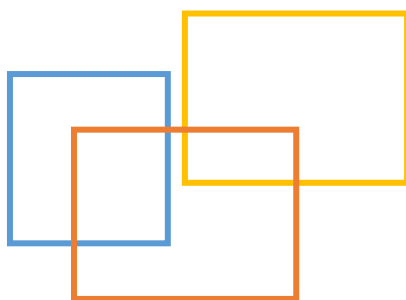
框选类能力涵盖：

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

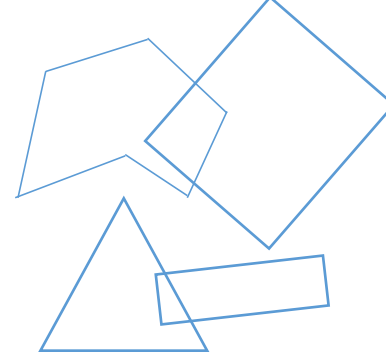
普通矩形框



分类矩形



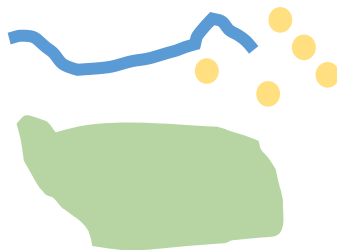
普通多边形



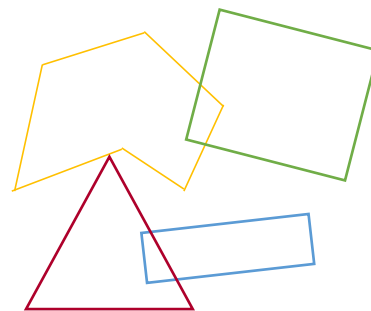
区域填色|多级属性多边形|Parsing



点+线+区域复合



分类多边形



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



◆ 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)

[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：

鼎挺沾證教志宵筐阱秃望

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

您好女士，我是联通的客服代表三三四，然后这边我刚帮您核实了一下，然后您之前那个在微信上购买的那个流量包，然后您对这个核销的顺序不认可，这边我能帮您然后再帮您看一下，然后这个也是建议您得问一下微信的那个客服，然后这边我也把这个微信客服这个联系电话告诉您，这边方便记一下吗

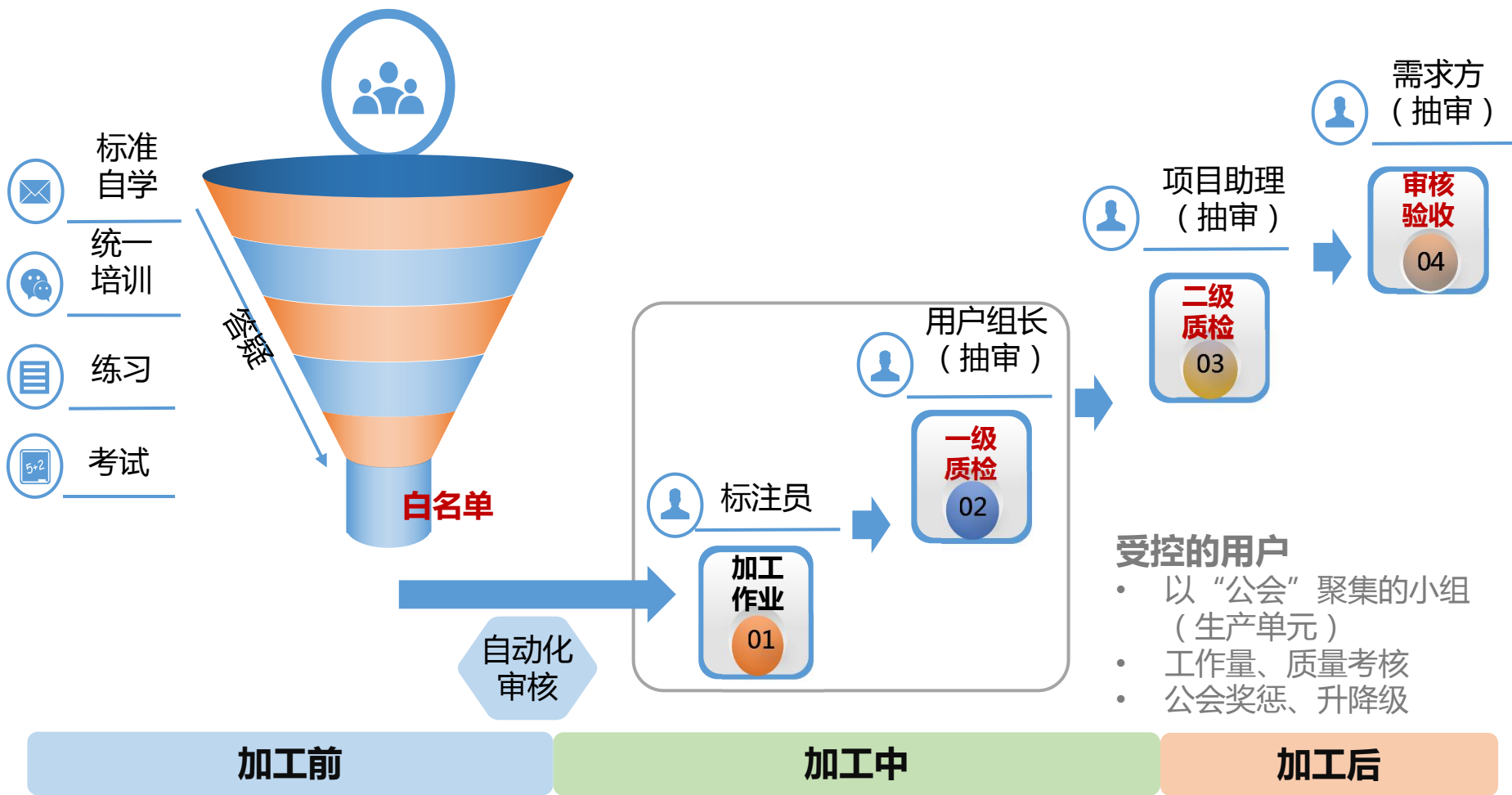
A screenshot of an audio transcription tool. It shows a waveform of the audio file. Below the waveform, there is a transcription interface with a text area containing the transcribed text: '您好女士，我是联通的客服代表三三四，然后这边我刚帮您核实了一下，然后您之前那个在微信上购买的那个流量包，然后您对这个核销的顺序不认可，这边我能帮您然后再帮您看一下，然后这个也是建议您得问一下微信的那个客服，然后这边我也把这个微信客服这个联系电话告诉您，这边方便记一下吗'. Below the text area, there is a table with columns for '文本层' (Text Layer) and '角色层' (Character Layer). The text layer contains the transcribed text, and the character layer contains the corresponding phonetic symbols for each character.

预览内容：

null

A screenshot of an audio transcription tool. It shows a waveform of the audio file. Below the waveform, there is a transcription interface with a text area containing the transcribed text: '您好女士，我是联通的客服代表三三四，然后这边我刚帮您核实了一下，然后您之前那个在微信上购买的那个流量包，然后您对这个核销的顺序不认可，这边我能帮您然后再帮您看一下，然后这个也是建议您得问一下微信的那个客服，然后这边我也把这个微信客服这个联系电话告诉您，这边方便记一下吗'. Below the text area, there is a table with columns for '文本层' (Text Layer) and '角色层' (Character Layer). The text layer contains the transcribed text, and the character layer contains the corresponding phonetic symbols for each character.

多级质量管控，突破准确率瓶颈



根据数据需求类型，覆盖更多实际应用场景

经过多年的数据积累，目前百度的人工智能数据仓库已经覆盖了超过5个大类，50多个小类别的实际人工智能模型数据集。



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

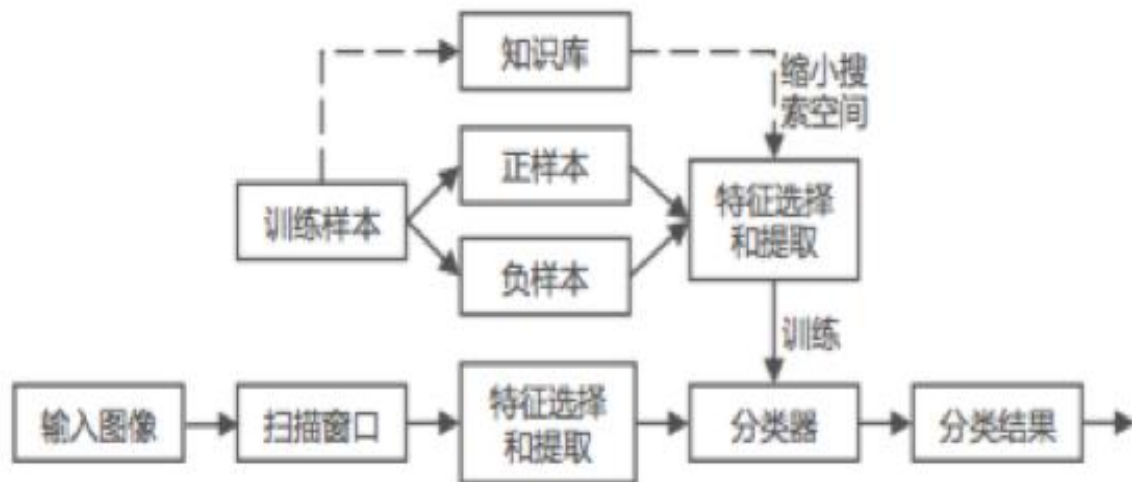
模式识别

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

典型人工智能应用场景



计算机视觉数据解决方案



计算机视觉的一般识别流程

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



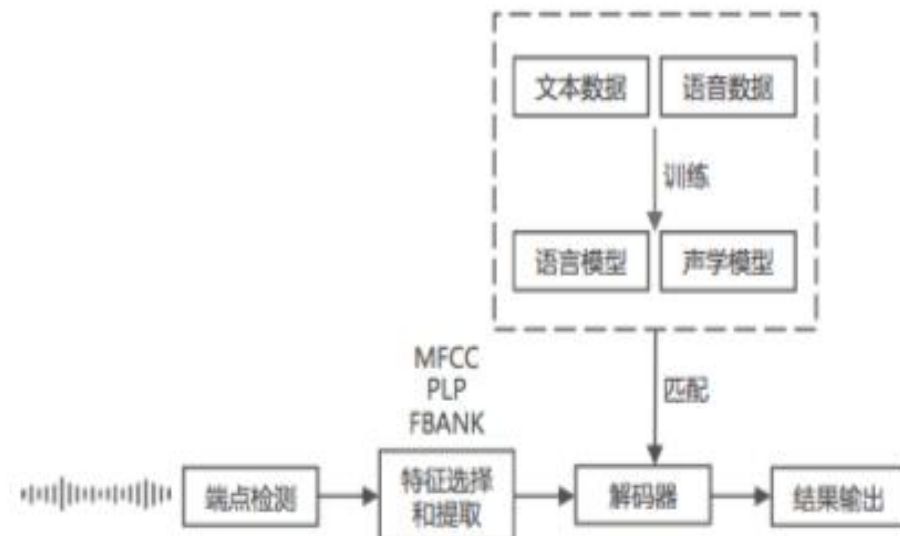
特殊要求人脸图像采集

- 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



语音识别的一般识别流程

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

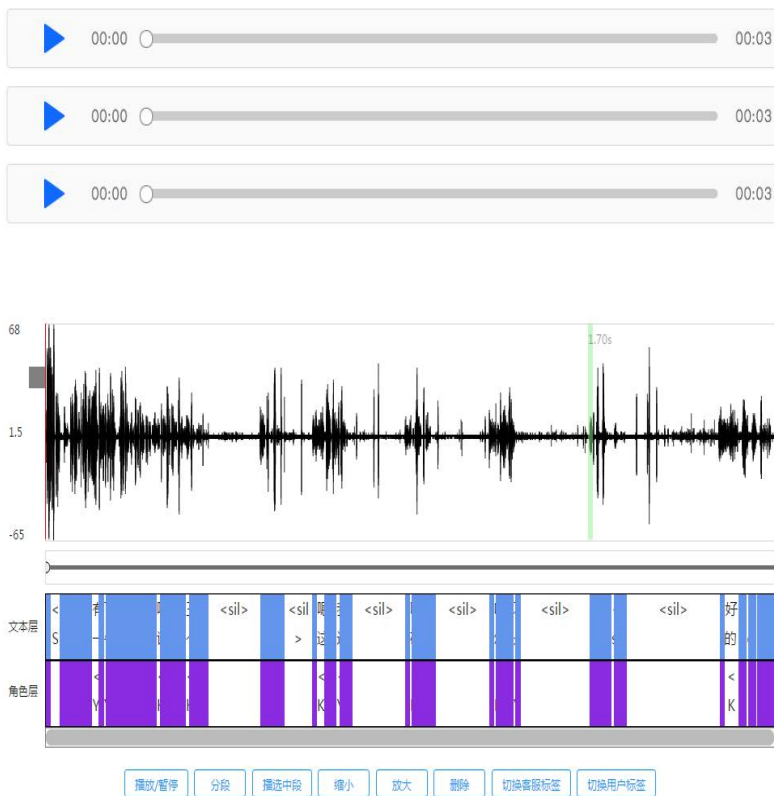
自然语言
生成

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

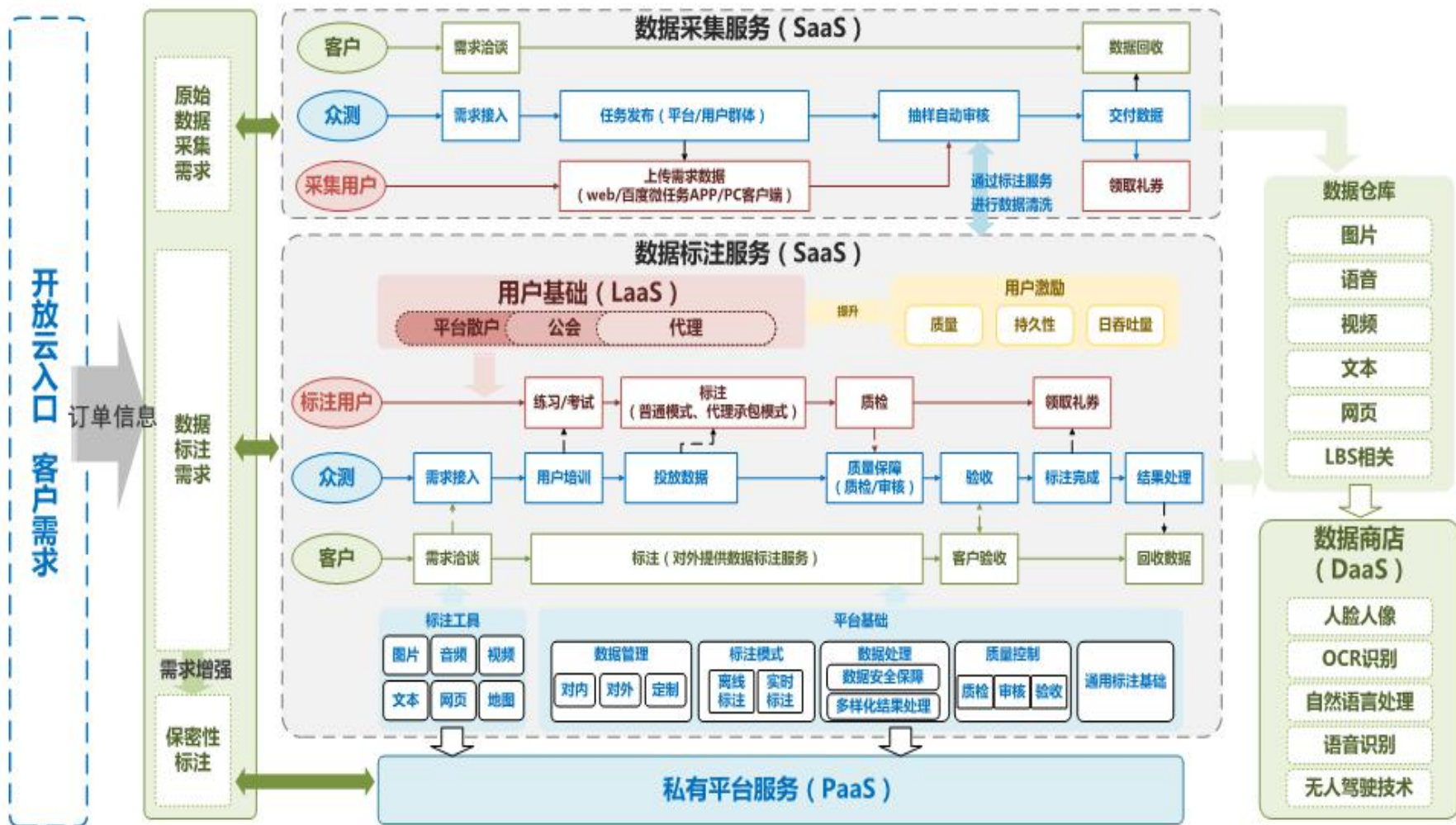
语音数据转写标注

- 中文方言、普通话
- 转写准确率98%

百度数据产品



百度数据产品矩阵



Thank

合作联系

zhongbao.baidu.com