

# 人工智能产品： 质量保障方案探索

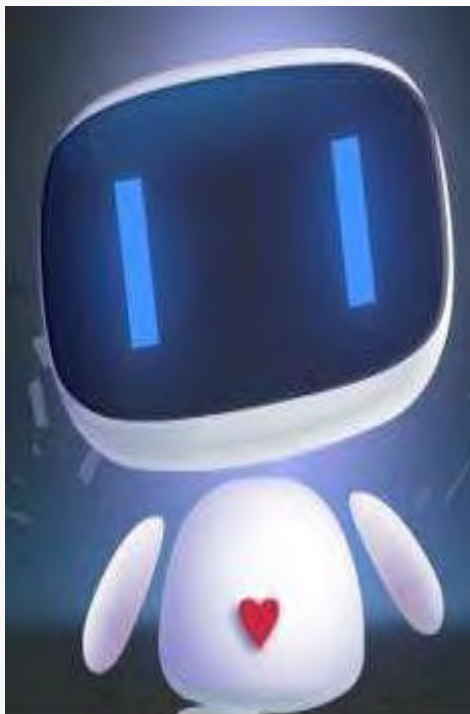
2017.7 钱承君

下一代  
软件研发  
SOFTWARE  
DEVELOPMENT

# 典型人工智能产品：无人车



# 典型人工智能产品：机器人



# 典型人工智能产品：推荐系统



**Baidu 百度** 人工智能 百度一下 百度首页 热搜

网页 新闻 贴吧 知道 音乐 图片 视频 地图 文库 更多

百度为您找到相关结果约100,000,000个 搜索工具

**IBM人工智能: IBM的认知技术平台 - IBM Watson**  
 人工智能-IBM Watson具备理解、学习和推理能力,实现更智能的人机互动,帮助做出更优决策。立即登录IBM官网,了解更多成功案例,下载白皮书!  
[www.ibm.com](http://www.ibm.com) 2017-03 33条评论 - 广告

**国际领先的人工智能技术 云从科技**  
 云从科技,中国科学院背景,国标制定企业,多行业定制化开发案例,全国电话400-151-5992,源自“计算机视觉之父”,拥有中科院与交通大学两大联合实验室,执行新疆安防布控设备用途,智能安防,智能金融,智能商务,设备品牌,识别准确率99%。云从信息  
[www.cloudwalk.cn](http://www.cloudwalk.cn) 2017-03 1条评论 - 广告

**百度云智——基于百度大脑打造的人工智能平台**  
 百度云智提供语音技术、文字识别、人脸识别、深度学习PaddlePaddle和自然语言NLP等人工智能产品及解决方案使用场景:智能客服/智能推荐/身份认证/内容审核/增强现实等  
[cloud.baidu.com](http://cloud.baidu.com) 2017-03 59条评论 - 广告

**英特尔人工智能学院 致力于让人工智能供所有人使用**  
 英特尔致力于实现人工智能创新民主化,通过增加数据、工具、培训,智能机器的可获性,使人工智能更加大众化,更好的改善我们的世界!点击了解  
[www.intel.cn](http://www.intel.cn) 2017-03 39条评论 - 广告

**登录百度账号 交易更有保障**

**人工智能个人助理**

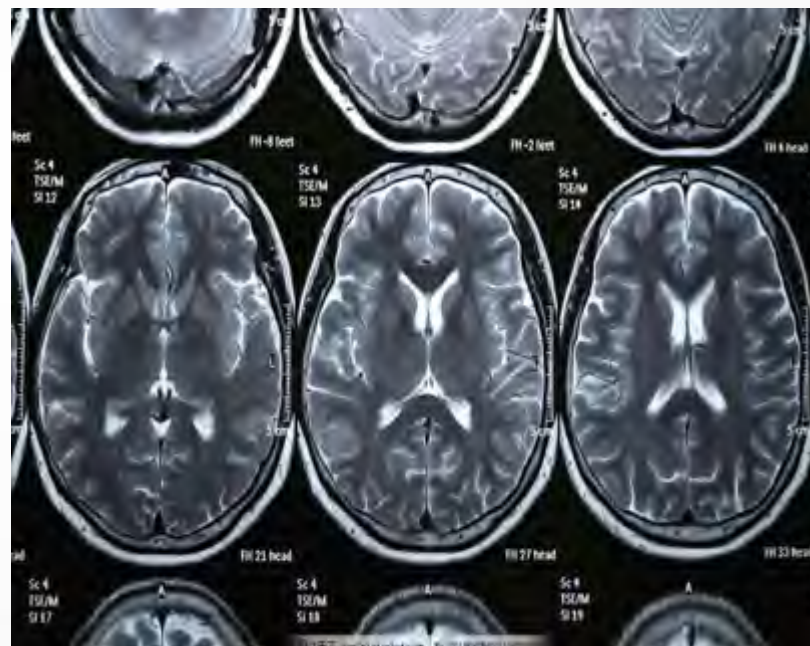
-  [Siri](#)
-  [Google Now](#)
-  [微软Cortana](#)

**相关电影** 展开

-  [机器人启示录](#)  
2017年上映  
科幻片
-  [环太平洋2](#)  
2017年上映  
科幻片
-  [独立日](#)  
打败外星人  
拯救地球
-  [猩球崛起](#)  
Charles Darwin  
科幻片



# 典型人工智能产品：图像识别



# 通常被谈及的质量保障范畴

## 流程控制

- 甘特图 / 看板
- 敏捷 / 持续集成

## 测试设计

- 等价类 / 边界值 / 因果图
- 异常 / 容错处理 / 安全

## 测试执行

- 探索性测试
- 自动化

## 测试分析

- 缺陷根因 / 收敛分析
- 性能瓶颈分析

传统测试方法的外延

众包、监控、用户体验评测



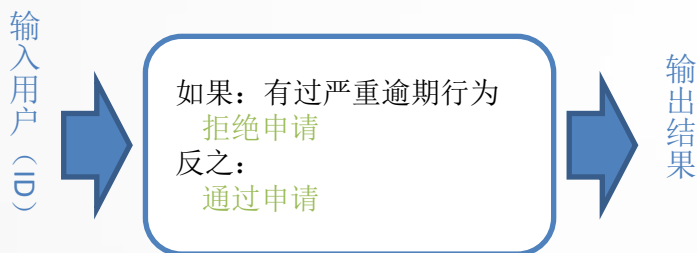
更完整、更高效、更低成本地

戳一戳，怼一怼，看看对不对

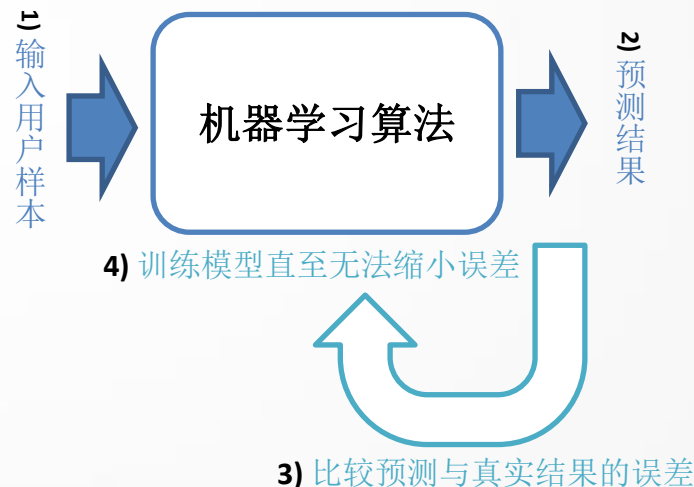
# 编外：简易版机器学习介绍

案例：判断是否准许特定用户的信用卡申请

传统思路（规则化思路）



机器学习（持续自优化的思路）



“A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.”

-- Tom Mitchell, Carnegie Mellon University

# 传统软件 vs 人工智能



后端工程测试

页面测试

手机端测试

性能 / 灾备 / 安全性



找到所需信息

更多样的内容

更新鲜的内容

更少干扰和无效信息



# 人工智能产品：质量保障思路初探

验收标准

评测数据

评测手段

结果分析

语音识别

近场识别

字准确率 > 95%

说明：识别结果与标注答案对齐后，去掉插入、删除、替换三种错误后，字级别正确的比例

- 1) 手百多模线上采样数据
- 2) 输入法线上采样数据

远场识别

五米，字准 > 92%

说明：距离麦克风一定距离场景下，模拟家居场景摆放，综合得出评估结论

- 1) 现场人工录制数据
- 2) 软件合成数据
- 3) 加噪数据，噪音源包含聊天、电视节目、音箱播放音乐、环境噪音

1) 设备按实际使用场景摆放，正交人工做测试输入获得结果数据

2) 通过软件手段合成与回放音频，获得更大批量结果数据

1) 多测试方法间结果趋于一致、或有合理解读

2) 与历史版本比，结果趋于一致、或有合理解读

3) 测试数据分布吻合用户实际使用分布

4) 修正和解读指标计算、样本选择带来的结果偏差

# 基础验收指标

Cat detector: return TRUE when there is a cat.

实际是猫 (P)

实际不是猫 (N)

TP

FP

判断是猫  
(T)



判断不是猫  
(F)



FN

TN

$$\text{precision} = \text{TP} / (\text{TP} + \text{FP})$$

判定为猫的图片中，有多少比例是真的猫

$$\text{recall} = \text{TP} / (\text{TP} + \text{FN})$$

实际为猫的图片中，有多少比例被判定为猫

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

$$\text{accuracy} = (\text{TP} + \text{TN}) / (\text{P} + \text{N})$$

是猫的判定为猫，不是猫的判定不是猫

$$\text{Specificity} = \text{TN} / (\text{TN} + \text{FP})$$

$$\text{Missing Alarm} = \text{FN} / (\text{TP} + \text{FN})$$

$$\text{False Alarm} = \text{FP} / (\text{TP} + \text{FP})$$

# 案例：色情图片检测

场景	指标	数值
色情	召回率	96%
	准确率	98%
	F-Score	0.97
正常	召回率	98%
	准确率	84%
	F-Score	90%

1000张色情图，混杂正常样本若干

总共有1000张色情图，判断对了960张，占比96%，漏判了40张  
所有被判为色情的图有980张，判对960张，判错20张，判对占比98%

1000张正常图，混杂色情样本若干

总共有1000张正常图，判断对了980张，占比98%，漏判了20张  
所有被判为正常的图有1167张，判对980张，判错187张，判对占比84%

	实际色情 (P)	实际正常 (N)	
判断色情 (T)	960	20	判断色情 (T)
判断正常 (F)	40	?	判断正常 (F)
			实际色情 (P)      实际正常 (N)
			判断色情 (T)      判断正常 (F)
			187      980

场景一：图搜下架色情图片

场景二：贴吧禁发色情图片

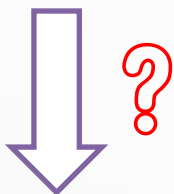
# 线上曾经存在的问题（图搜色情）



# TiD2017 用户向的验收指标 (User Side Matrix)



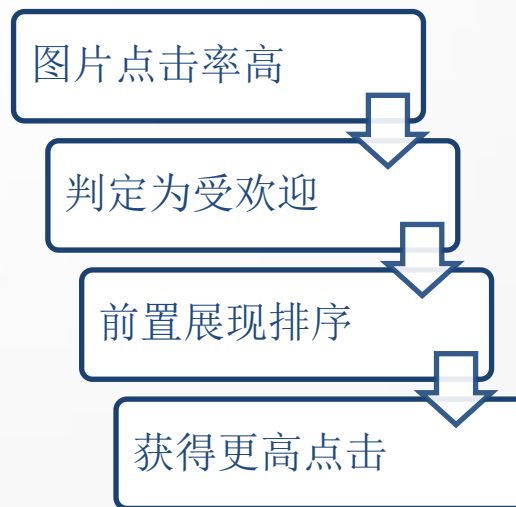
色情图片召回率  
**80% → 95%**  
漏检数量降低四倍



图搜色情图片流量  
**2% → 0.5%**  
同比降低四倍

备注：非真实数据

即使只有少量色情图片存在，会因以下搜索策略被推高，实际展现量远高于推算预期。





# 指标制定过程常见错误

## 指标与用户需求偏离

- 色情图片检出率 vs 色情图片误检率

## 指标不对终端用户负责

- 色情图片检出率 vs 线上色情图展示率

## 忽视指标间联动与完整

- 准确与召回不可兼得
- 复杂模型带来更好的效果，但有更大的功耗

# 评测数据集举例

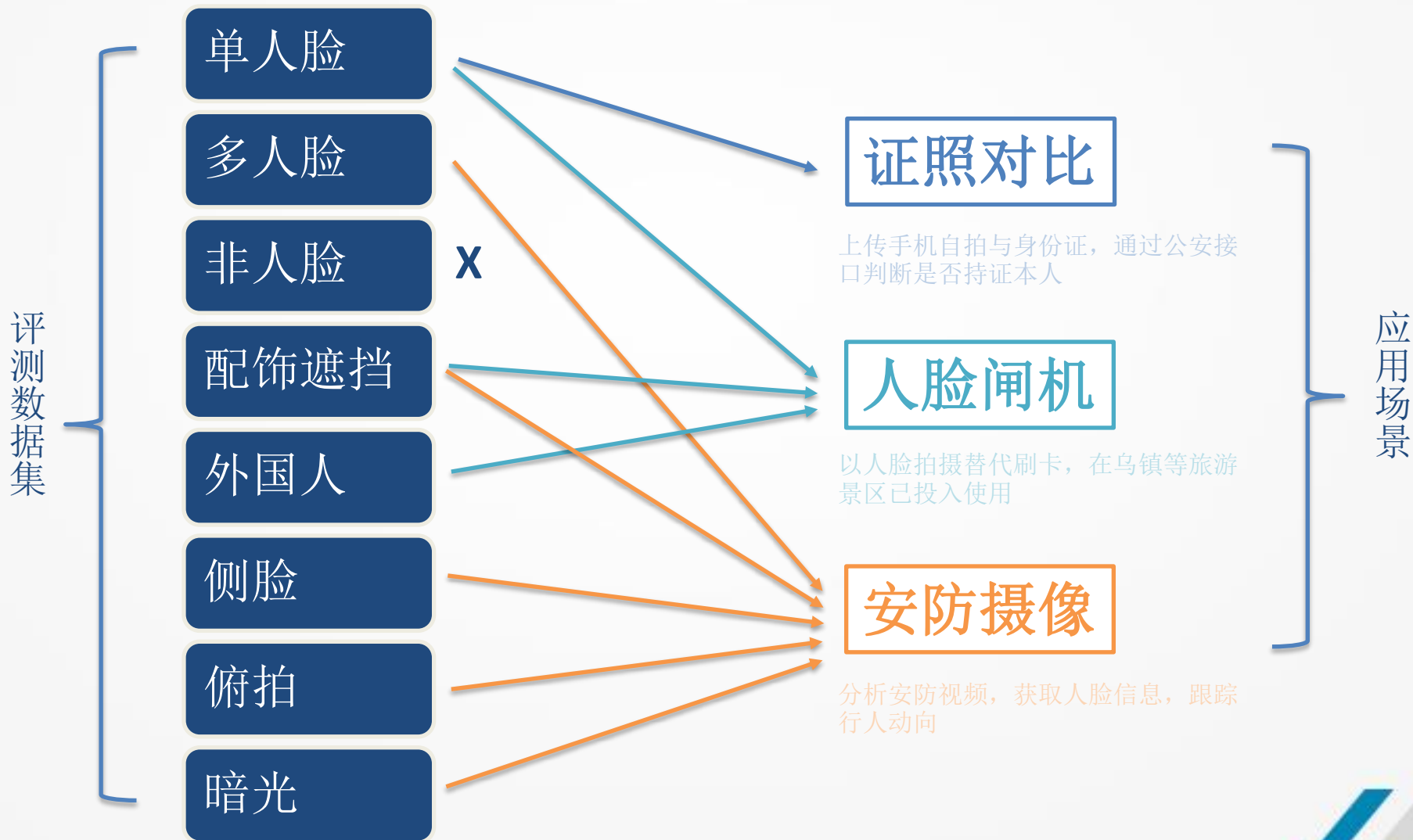
## 语音识别

- 发音特质：男女老少，口音，声调，语速
- 语言内容：语种，长句，中英混合，专名
- 噪音环境：人声，交通，空调，电视，音乐，通讯设备
- 应用场景：车载离线，远场识别，麦克风阵列

## 人脸检索

- 基础：多人脸，人脸尺寸
- 干扰：遮挡，侧脸，帽子，眼镜
- 光源：暗光，反光，曝光过度，隔玻璃反光
- 图像：模糊，失焦，黑白
- 细分：老人，外国人，卡通人脸

# 针对场景构建评测集



# 评测数据获取

## 人工标注

- 自标注
- 自建标注团队
- 众测模式

## 数据复用

- 行标数据，采购数据
- 线上数据 / 用户数据

## 数据合成

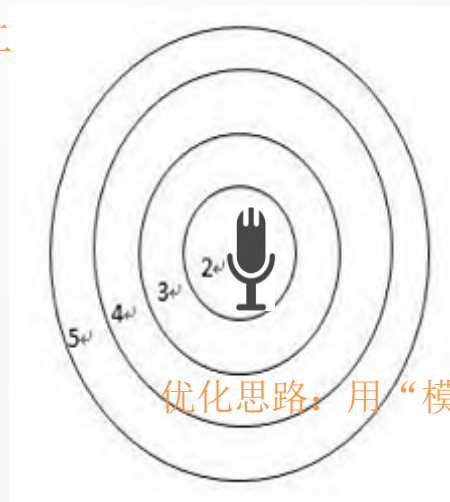
- 生成不同噪音背景下的语料
- 生成同一人脸不同曝光下的图像

案例一



优化思路：用“验证”代替“标注”

案例二



优化思路：用“模拟”代替“采集”

# 评测集的拟合

## 模拟数据与实测数据拟合

- 通过模拟手段扩容样本
- 该场景下的实测数据
- 对比关键指标分布一致性
- 修正数据合成方案

## 评测样本集间拟合

- 多样本集间获取交集
- 比对交集部分的偏差，标识相对置信度
- 对部分评测集进行修正取舍

## 案例一

实测方法：距离设备0-5米距离，固定摆放，每个关键点不同录制人多次朗读语料。

模拟方法：通过近场（0米）录制语料，计算得出不同距离的语料情况。

当上述结论趋于一致，认为模拟得出的数据有较高的仿真度，可用于验收。

## 案例二

样本一：公司 ERP 系统导出信息。（性别/年龄）

样本二：合作友商导入CRM数据。

样本三：用户注册时自填信息。

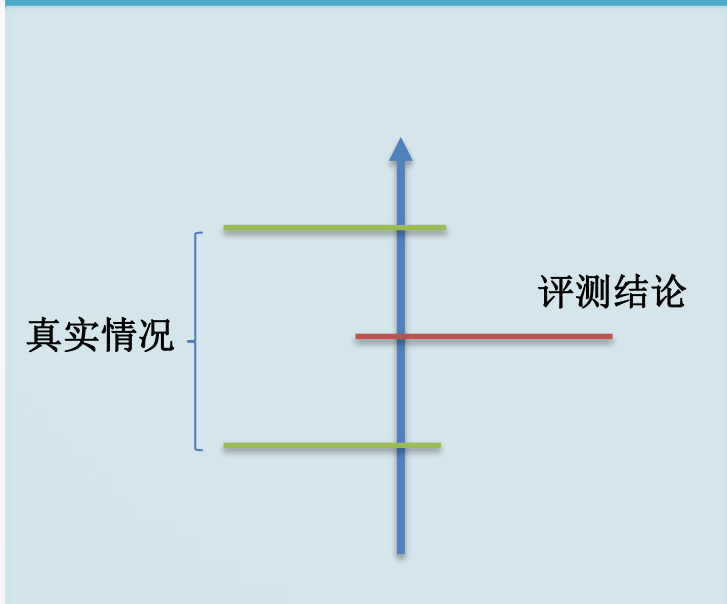
样本四：依据用户行为标注推导出信息。

各用于评测的样本来源不一致，置信度不一致。依据样本间相对关系，做为样本置信分析的依据。

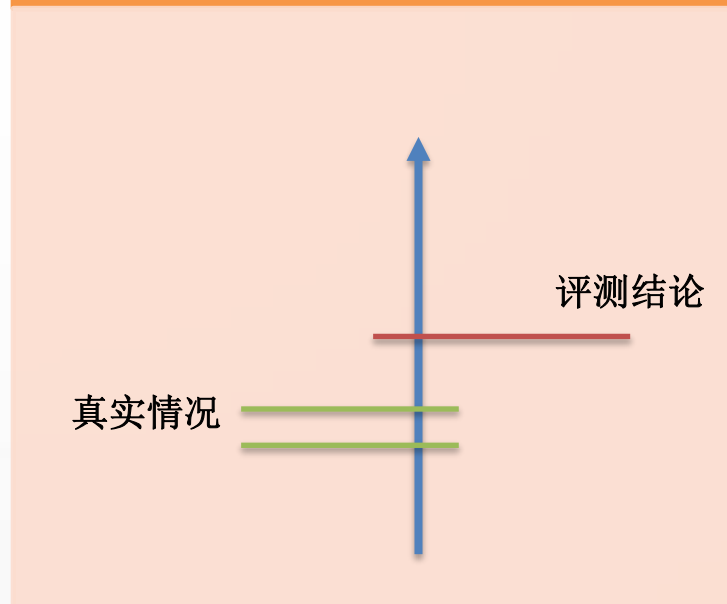


# 评测样本选择

## 样本数量



## 样本分布



# 评测数据选取常见错误

## 样本数量

- 数量过小，例：拿个位数样本判断准确率
- 数据偏差，例：一千张图片中只有五张色情图

## 样本分布

- 样本泛化，例：导航和智能音箱场景下的中英混合识别
- 场景偏差，例：部署人脸识别时隔着玻璃
- 正负例数量不对等，例：人脸闸机需要彩照等负例

## 评测集选择

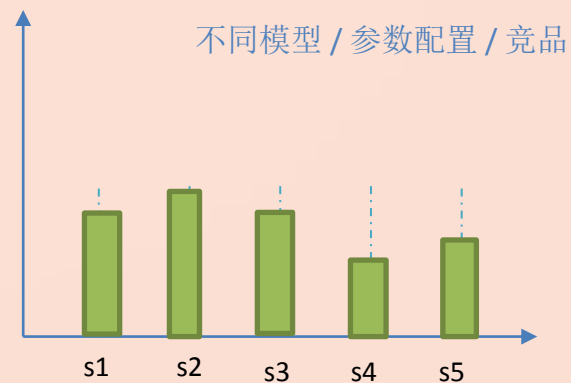
- 过拟合，测试数据直接用于模型训练
- 把自身优势场景用于竞品对标，不做交叉验证

# 评测结果：可连续对标

## 纵向对标

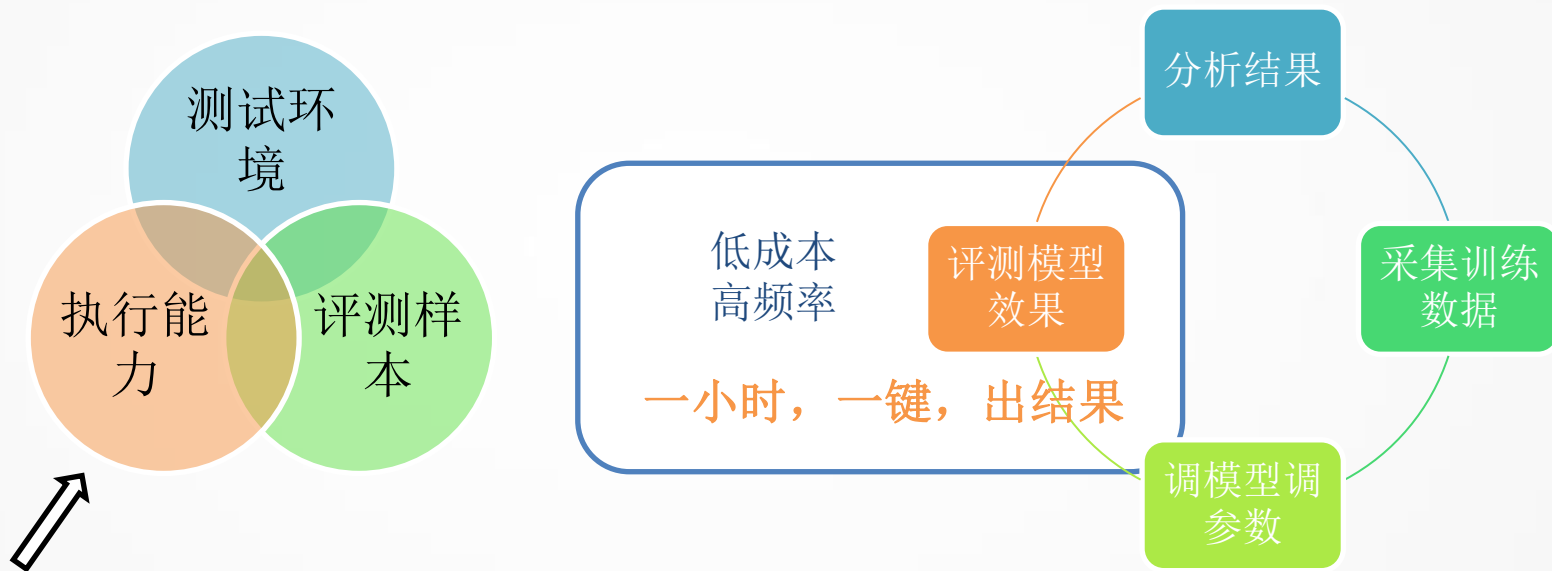


## 横向对标



常见问题：前后测试报告给出差异很大的结论，数据不连续

# 评测过程：一键执行



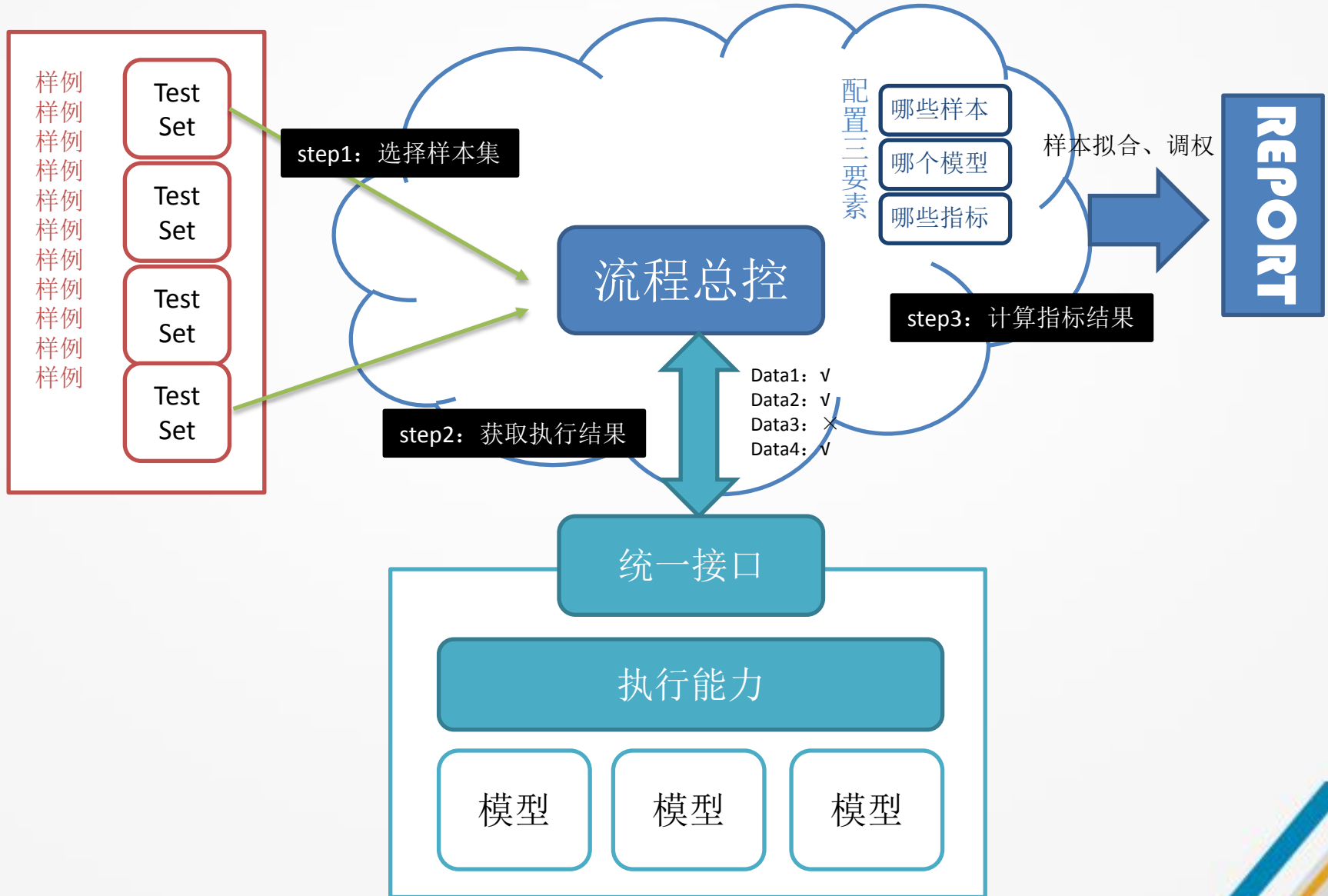
案例：语音语料复用

Learning Loop

- 软件注入
- 耳机线导入
- 人工嘴

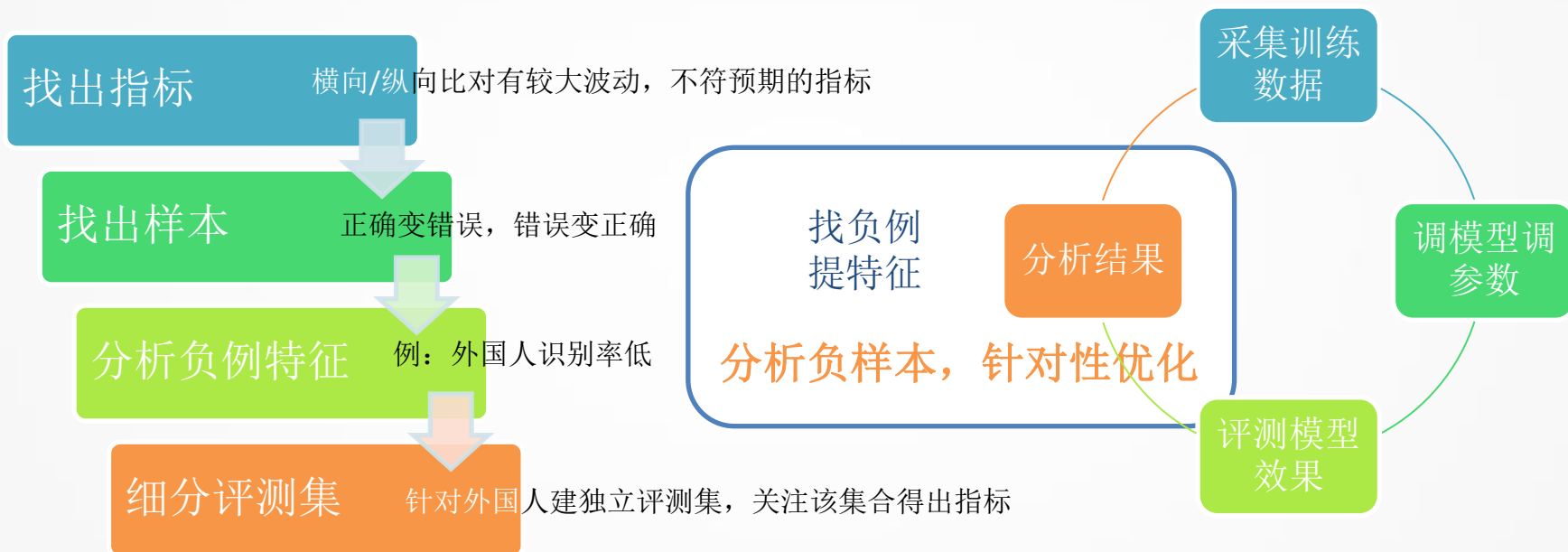


# 评测系统的实现思路



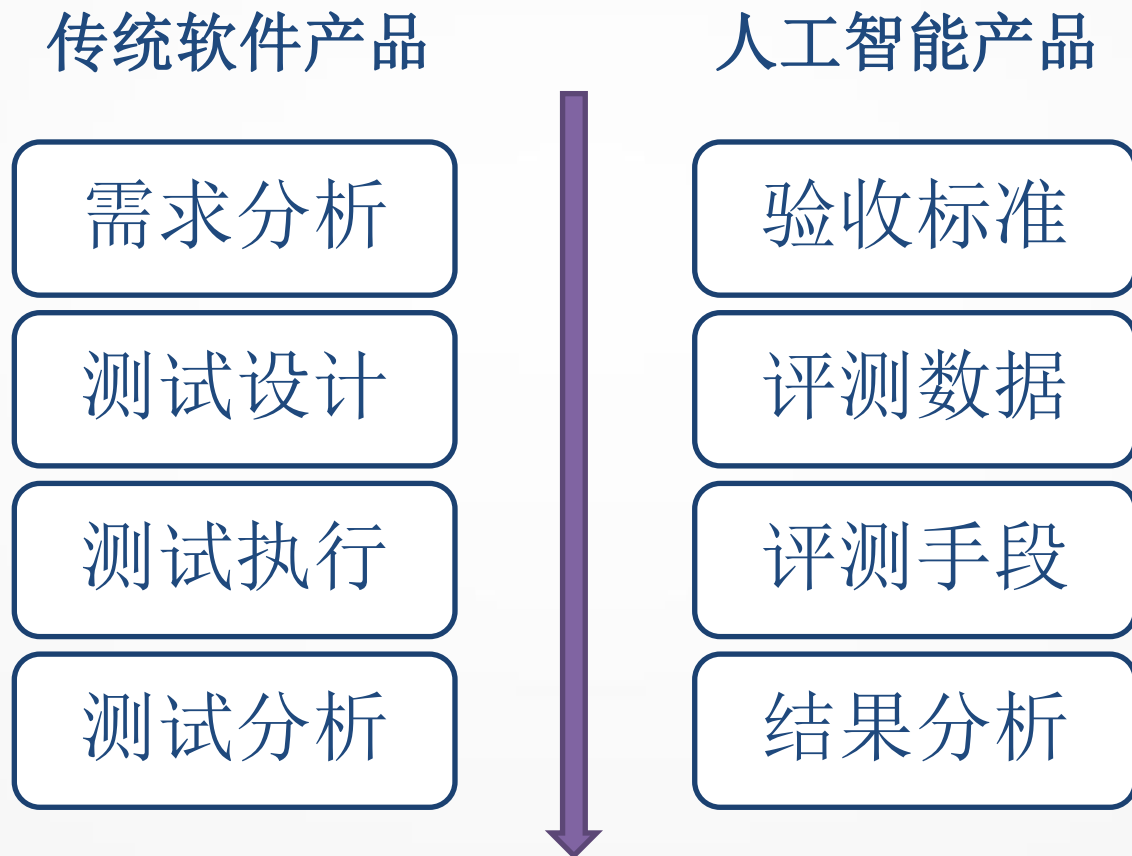


# 结果分析



执行过程变体：针对邻接版本差集，抽样标注评估

# 质量保障过程类比



关键点：场景理解 + 测试设计

# 务虚谈：行业走向的个人观点

真正传统的领域  
例如：航天所

软件工程化，外包盛行  
例如：CMMI

互联网谈敏捷，自动化盛行  
例如：Selenium

本质未变，基调收敛  
承认手工测试，不同技能组合

测试设计

研发流程

自动化

移动应用

下一幕是什么？

感谢

THE END

