



華東師範大學
EAST CHINA NORMAL UNIVERSITY

基于众包的非确定知识图谱清洗

林欣

华东师范大学计算机系



知识图谱构建

- 四种常见方法 [KDD 14]
 - Wiki-based : YAGO、Dbpedia、FreeBase
 - Ontology-based: NELL
 - Open IE: TextRunner, Reverb
 - Taxonomies (is-a): Probase



Wiki-Based方法

- 准确率高
- 知识结构性好

家族成员

父：闔閭

配偶：妃西施（后改嫁，一说被杀）

儿子：太子友（前482年，勾践乘虚而入，大败吴师、太子友被杀）、王子地

孙子：王孙弥庸、王孙寿



Wiki-Based 方法

- 缺点
 - 知识长尾效应
 - 覆盖度有限
 - FreeBase中 71%人物没有“出生地”，75%的人物没有国籍



自动挖掘方法

- 知识覆盖度大
- 噪音大、知识准确率难保证

王菲李亚鹏结婚照首度曝光 天山拍摄简单的甜蜜
[组图]

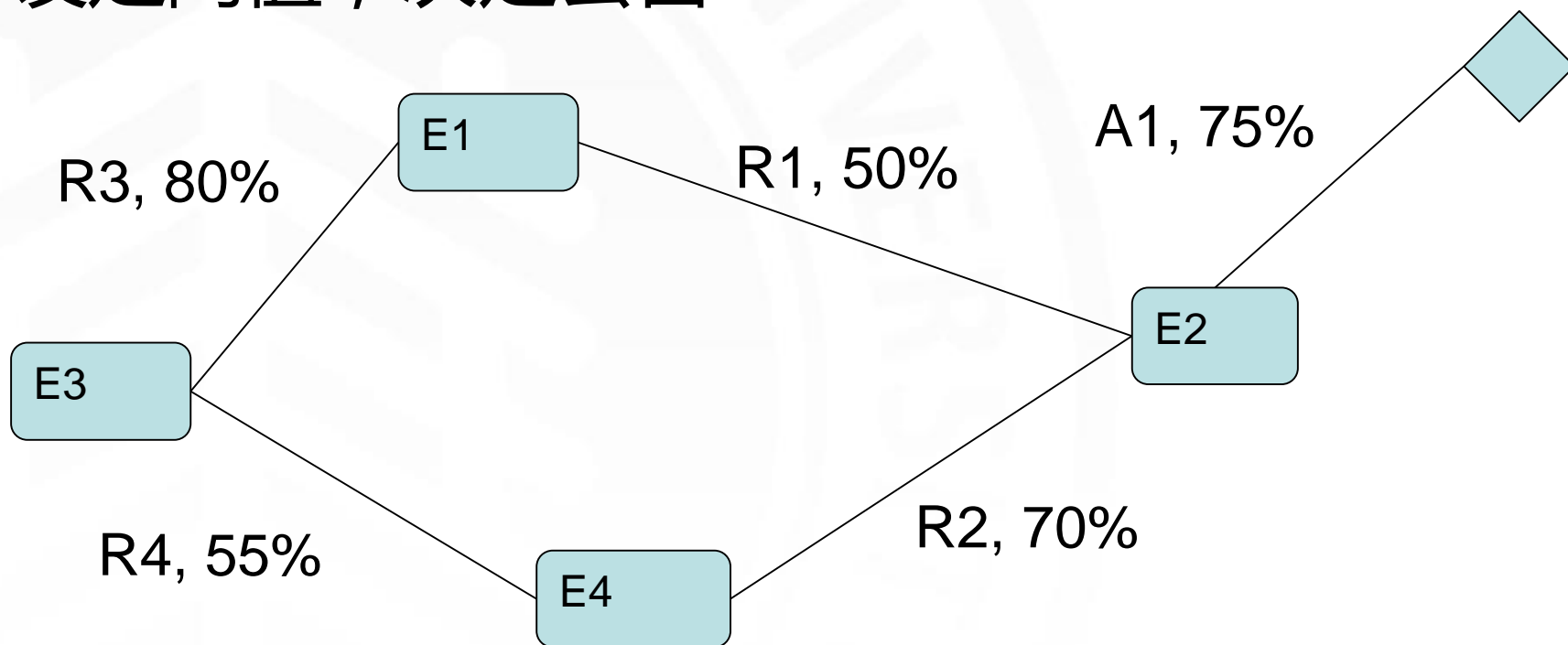
2011-12-07 08:01:00 来源: 中国娱乐网 有0人参与 手机看新闻  转发到微博(0)





传统非确定性数据解决办法

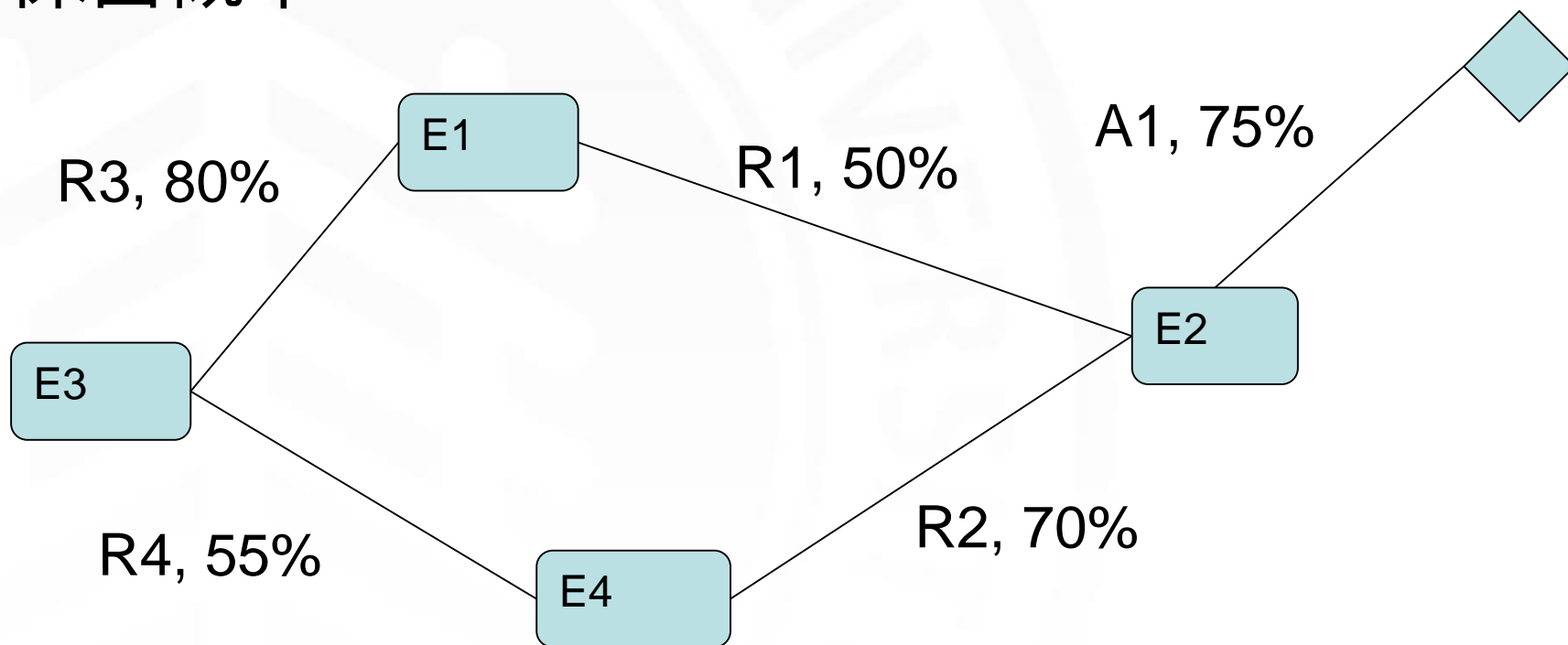
- 设定阈值，决定去留





非确定知识图谱

- 保留概率





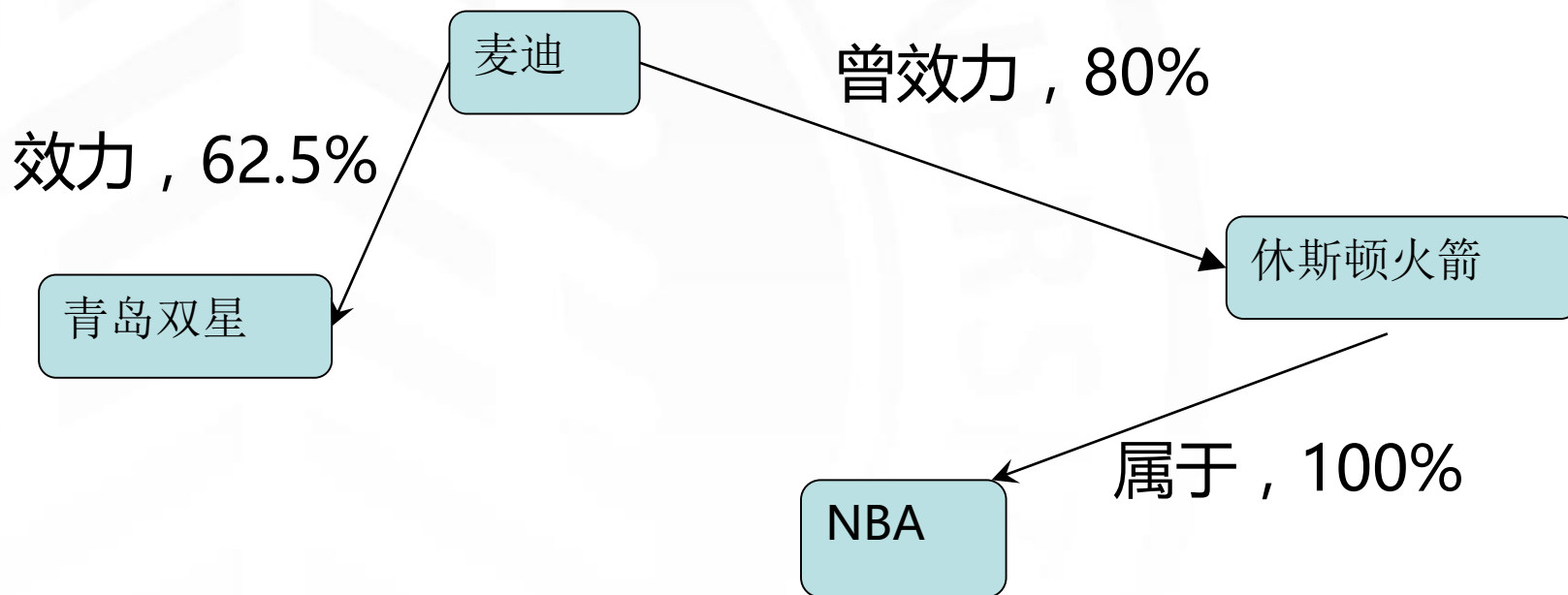
非确定知识图谱

- 优点：
 - 知识覆盖度
 - 信息更丰富 (Provenance)
- 缺点
 - 查询质量 (Ambiguity)
 - 0.5 问题



0.5问题

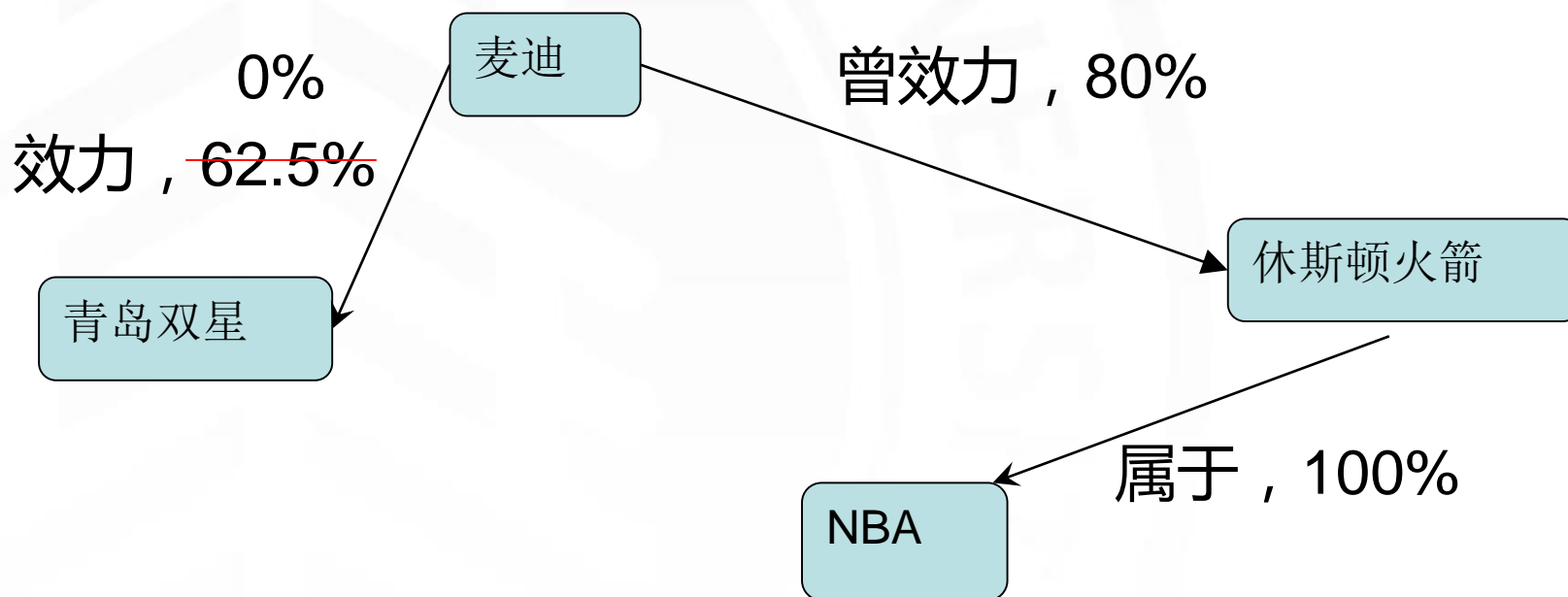
- “现在在中国打球的前NBA球员？”





0.5问题

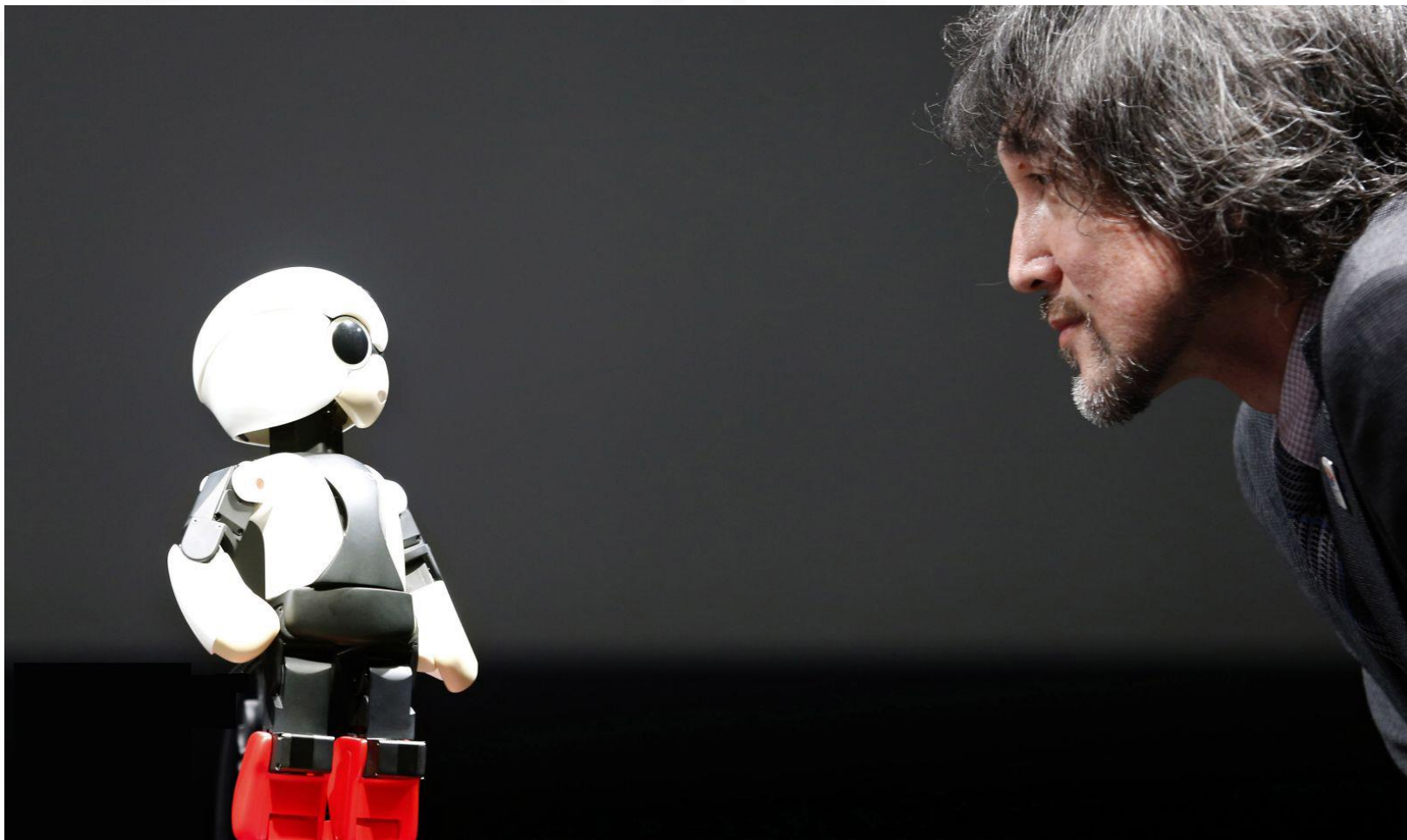
- 怎么办？概率清洗





華東師範大學
EAST CHINA NORMAL UNIVERSITY

机器还不够聪明，且还会不够聪明很久





众包清洗，用外脑的清洗架构



双低

机器+人脑

双高

机器 + 机器

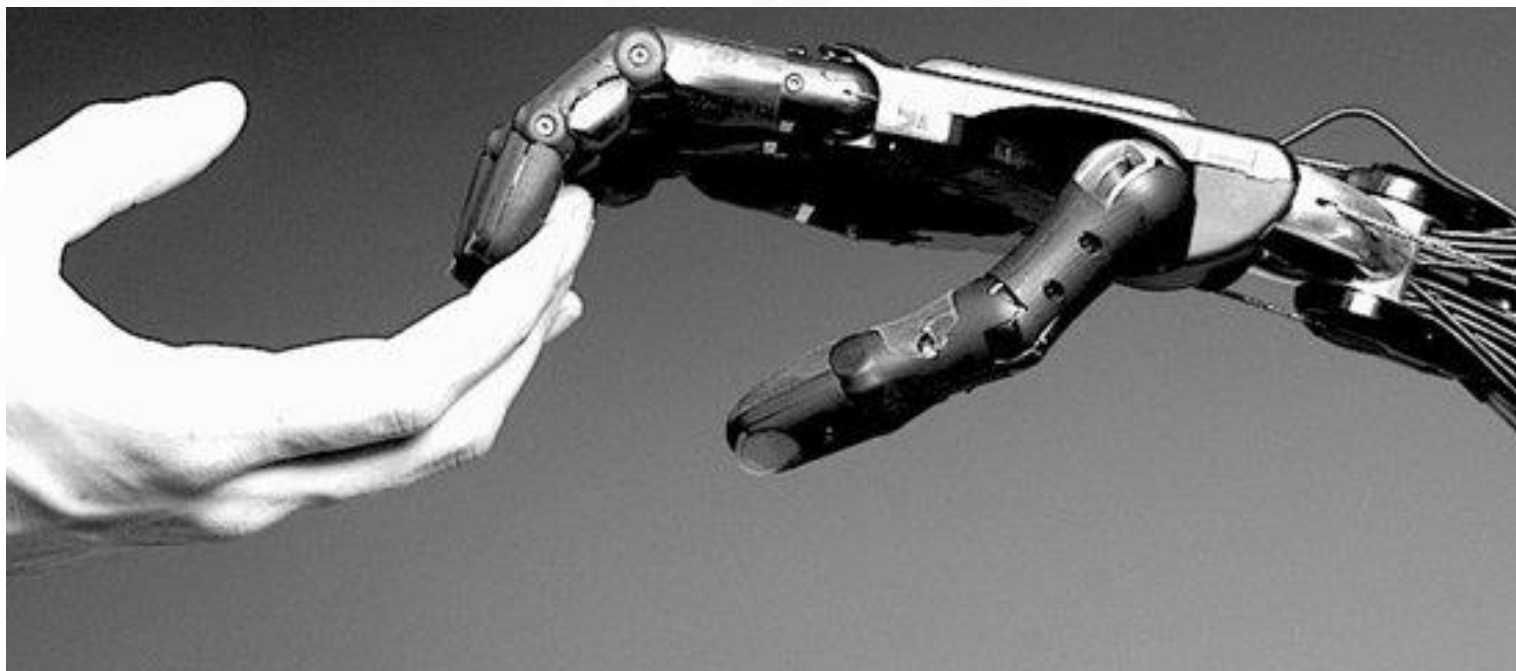
人脑 + 人脑

代价

质量

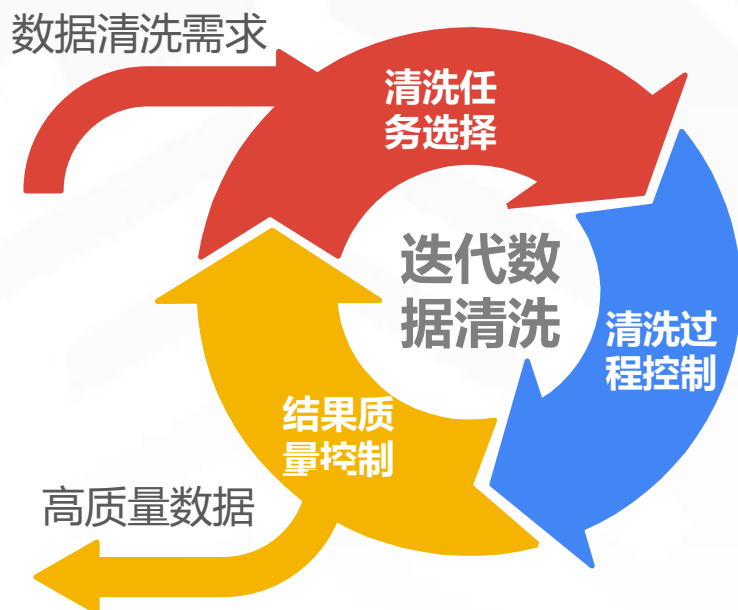


“众包清洗的科学问题”





问题1：基于众包的迭代清洗机制



问题描述：

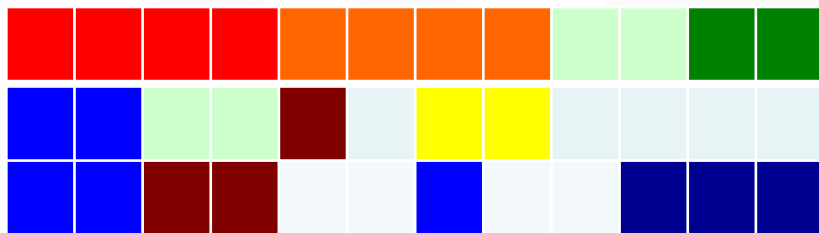
数据清洗的过程是不断提高数据质量的过程，而这种提升往往不能通过一次清洗而到位。正确的过程是通过若干次群智计算与机器学习的交织迭代，不断提高数据质量。

研究重点：

1. 迭代终止条件（收敛程度监测？）
2. 迭代过程自动化（包括子步骤间的衔接）



问题2：清洗收益的高效算法



问题描述：

众包更擅长清洗较小的数据，如识别两张图片，识别一个单词。而小数据之间存在关联，清洗每份小数据对全局质量提升量不同（如左图），需要将有限的资源投放给清洗收益最高的小数据。

研究重点：

1. 单个小数据清洗收益算法；
2. 多个关联小数据清洗综合收益算法



问题3：众包决策的正确性判决

最终结果

正确性判断

问题描述：

系统无法控制众包返回的结果，因此，每个工人的结果可能差异性较大，而上层应用只需要众包平台通过正确性判决返回统一的结果。

结果1

结果2

结果3

结果4

..... 研究重点：



工人1

工人2

工人3

工人4

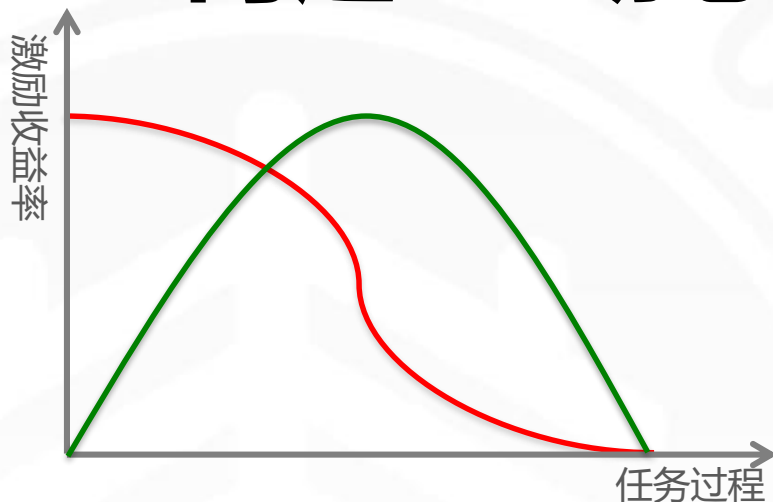
.....

1. 工人建模分析：包括使用工人背景、众包提交历史记录、互评价等数据，为清洗筛选工人或为工人建模；

2. 结果正确性判断：综合所有返回的结果、工人错误率等模型，获得最终结果



问题4：动态激励机制研究



- 独立数据
- 关联数据

问题描述：

众包的激励机制（如金钱、积分）是工人完成任务的动力，越多激励投入，任务收益也越大。然而在任务进行过程中，激励的收益率是动态的，也和具体数据相关，如独立性数据收益率会越来越小，关联性数据收益率呈抛物线。因此，需要研究如何利用最少的代价获得最多的清洗收益。

研究重点：

1. 动态激励收益率计算方法；
2. 最佳激励分配机制的近似算法；



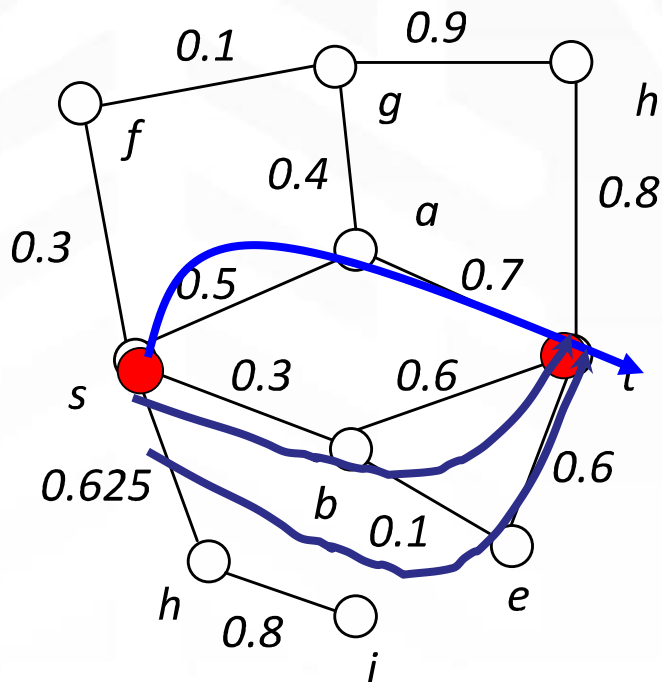
我们的工作

- 1. 众包关系选取
- 2. 众包平台开发



众包关系选取

- 特定查询优化（可达性查询）



Q: 点s和t是否3跳可达？

A: $P = 0.47$

目标：选一条边清洗，
使得清洗后的答案的不
确定性最低（靠近0或者
1）



问题定义

- 不确定性度量
 - 熵理论
 - $H(A) = -P(A)\log P(A) - (1-P(A))\log(1-P(A))$
 - 清洗后的熵 $H_{e=1}(A)$, $H_{e=0}(A)$
- 期望熵
 - $E = P(A) H_{e=1}(A) + (1-P(A)) H_{e=0}(A)$



基本思路

- P^* 因子
 - 所有过 e 的路径连通，而不过 e 的路径都不连通的概率
 - 证明： P^* 与 E 正相关
- P^* 快速排序算法

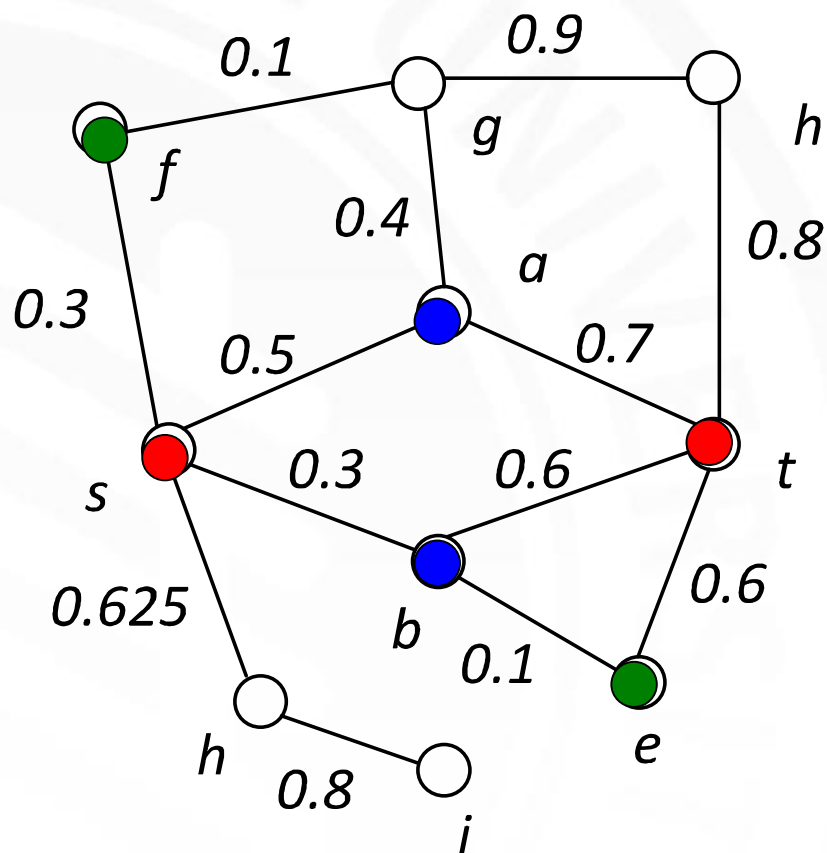


问题扩展

- 全图清洗
- 多边清洗
- Noisy Crowd
- 其他布尔查询 (Connectivity、K-core)
- 其他枚举型查询 (Subgraph Matching)

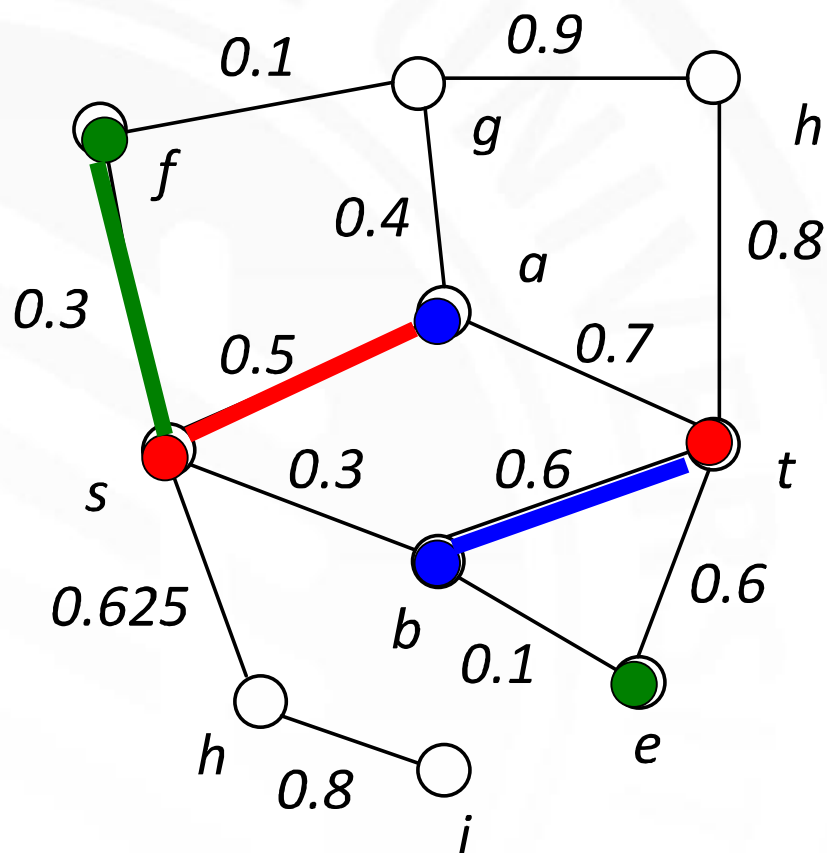


全图清洗





多边清洗





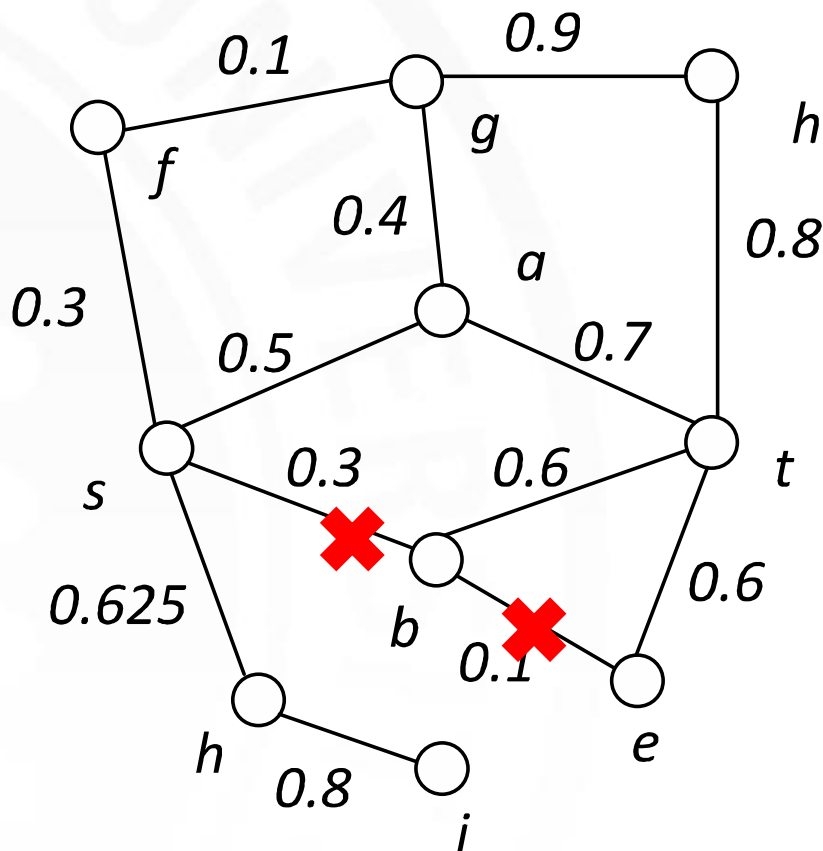
Noisy Crowd

- 众包并不能保证完全正确 [ICDE15]
- 假设Crowd返回的结果和真实值有一定偏差
- 在此假设下选取最优边



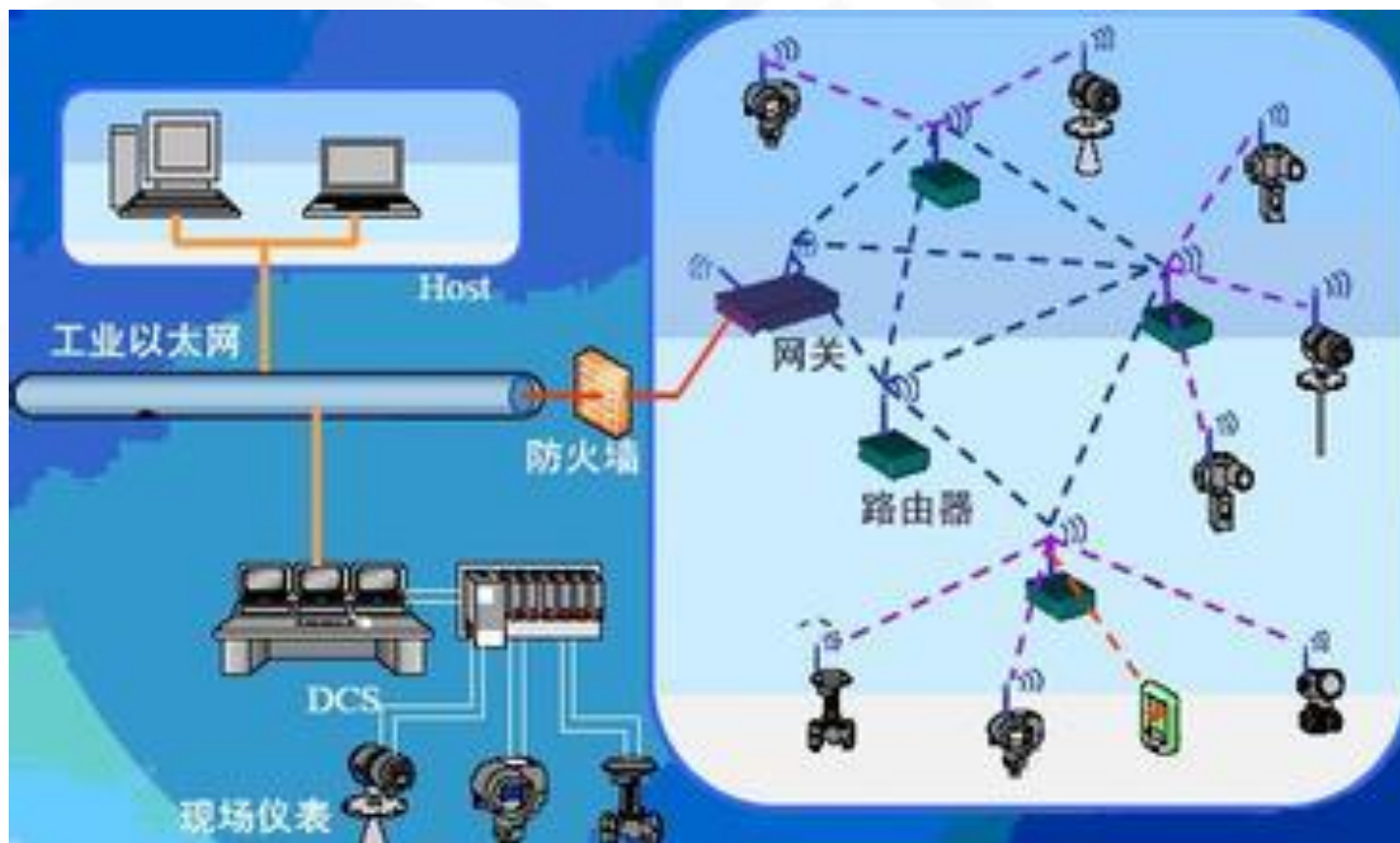
其他布尔查询清洗

- 连通性查询
– 2-Connectivity
- K-Core





应用（无线传感器网络）



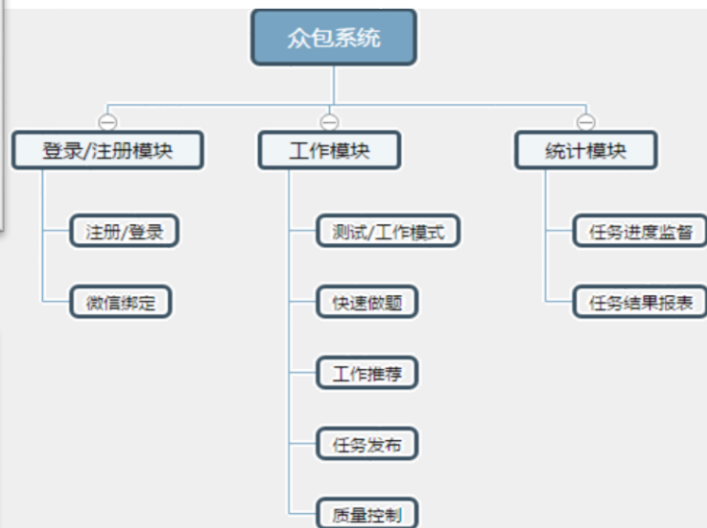


枚举型查询的清洗

- 不仅仅只有Yes & NO两种答案
- 每种答案都有概率，组成了信息熵
- SPARQL 查询：找出所有在中国打球的前NBA球员

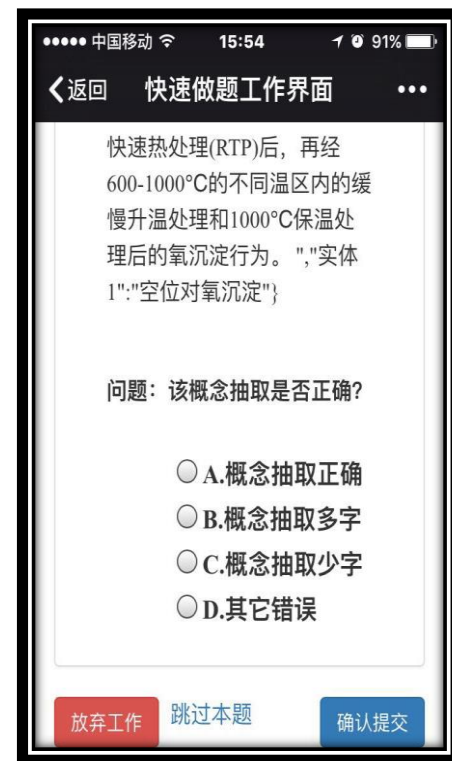


众包平台搭建





移动端界面





華東師範大學
EAST CHINA NORMAL UNIVERSITY

Thank you!