



# 大规模在线教育知识图谱构建

李涓子  
知识工程研究室  
清华大学计算机科学与技术系



# 主要内容

## □ 智能学术搜索

- 人工智能与知识
- 智能学术搜索

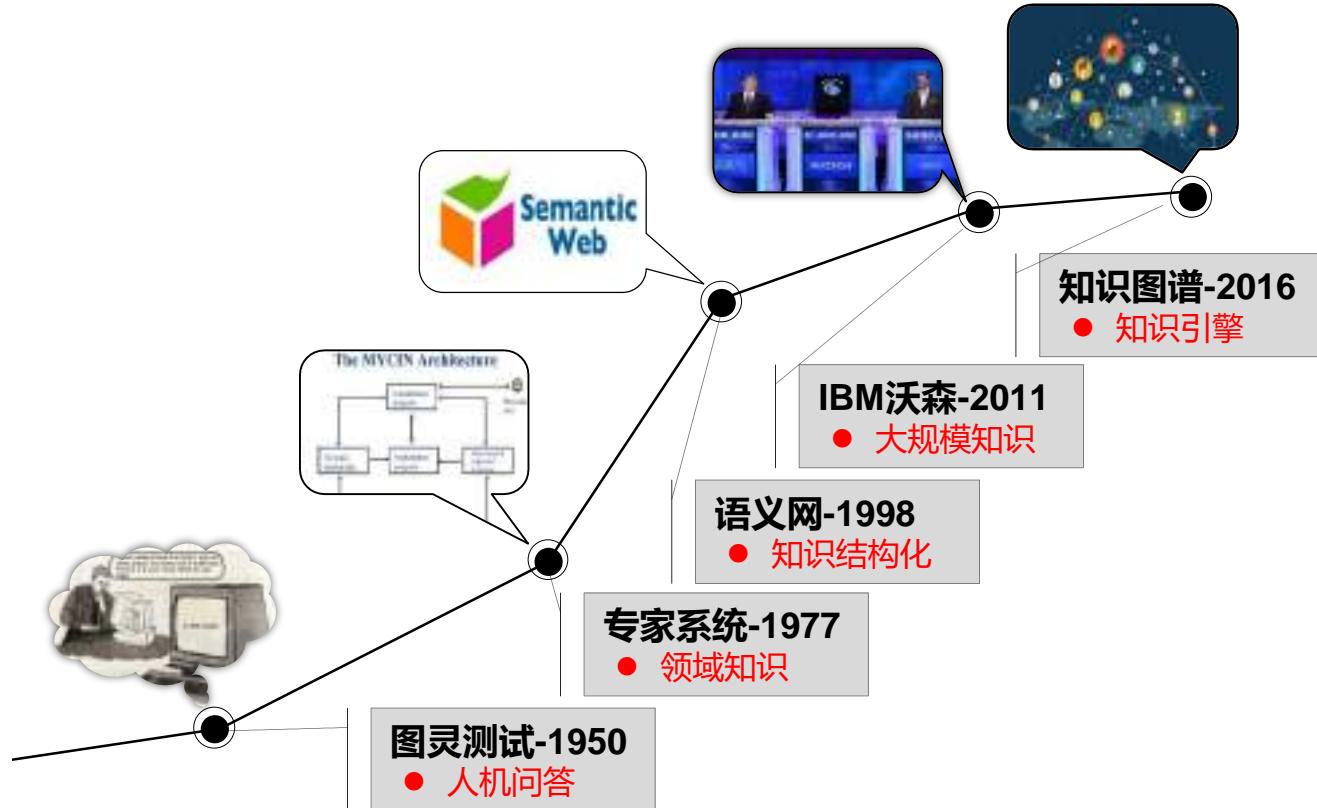
## □ 大规模在教育知识图谱构建

- 以知识为中心的大规模在线教育资源组织
- 概念学习
- 概念扩展
- 概念先后修关系发现

## □ 总结



# 人工智能对知识的需求



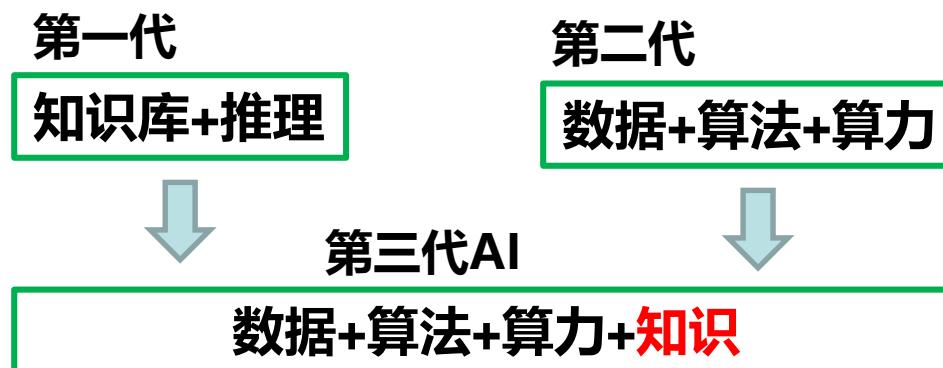
知识是人工智能的核心命题与关键基础



# 第三代人工智能

- Explainable, robust, credible and safe AI  
可解释的、鲁棒的、可信安全的人工智能
- From AI without understanding to AI with understanding  
从不带理解的人工智能成为带理解的人工智能

摘自张钹院士在YOCSEF精英大会上报告



实现从数据智能到知识智能的跃迁是人工智能的必由之路



# 人工智能之学术搜索

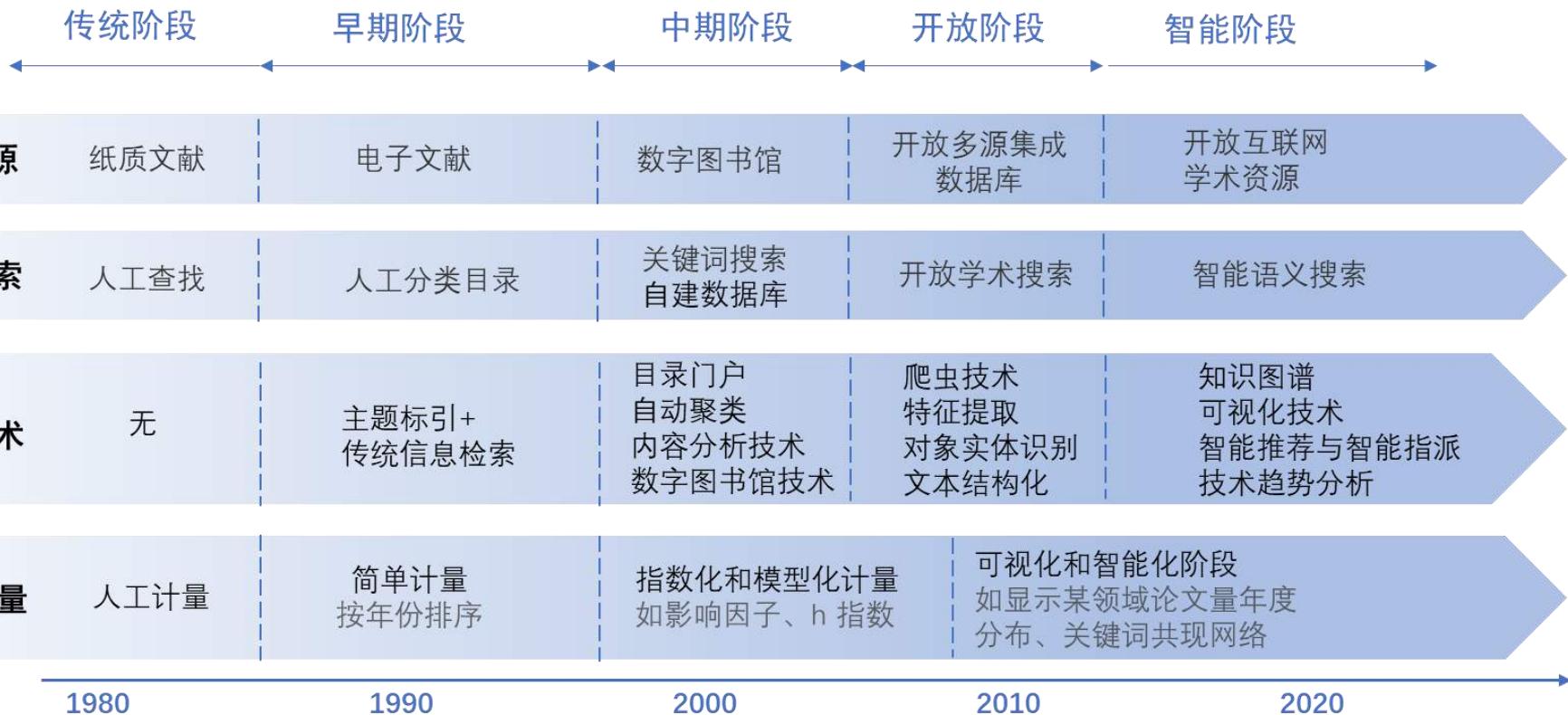


扫码微信阅读



# 人工智能之学术搜索

学术搜索引擎发展历程



# 人工智能之学术搜索

创新应用



双引擎  
自动驱动



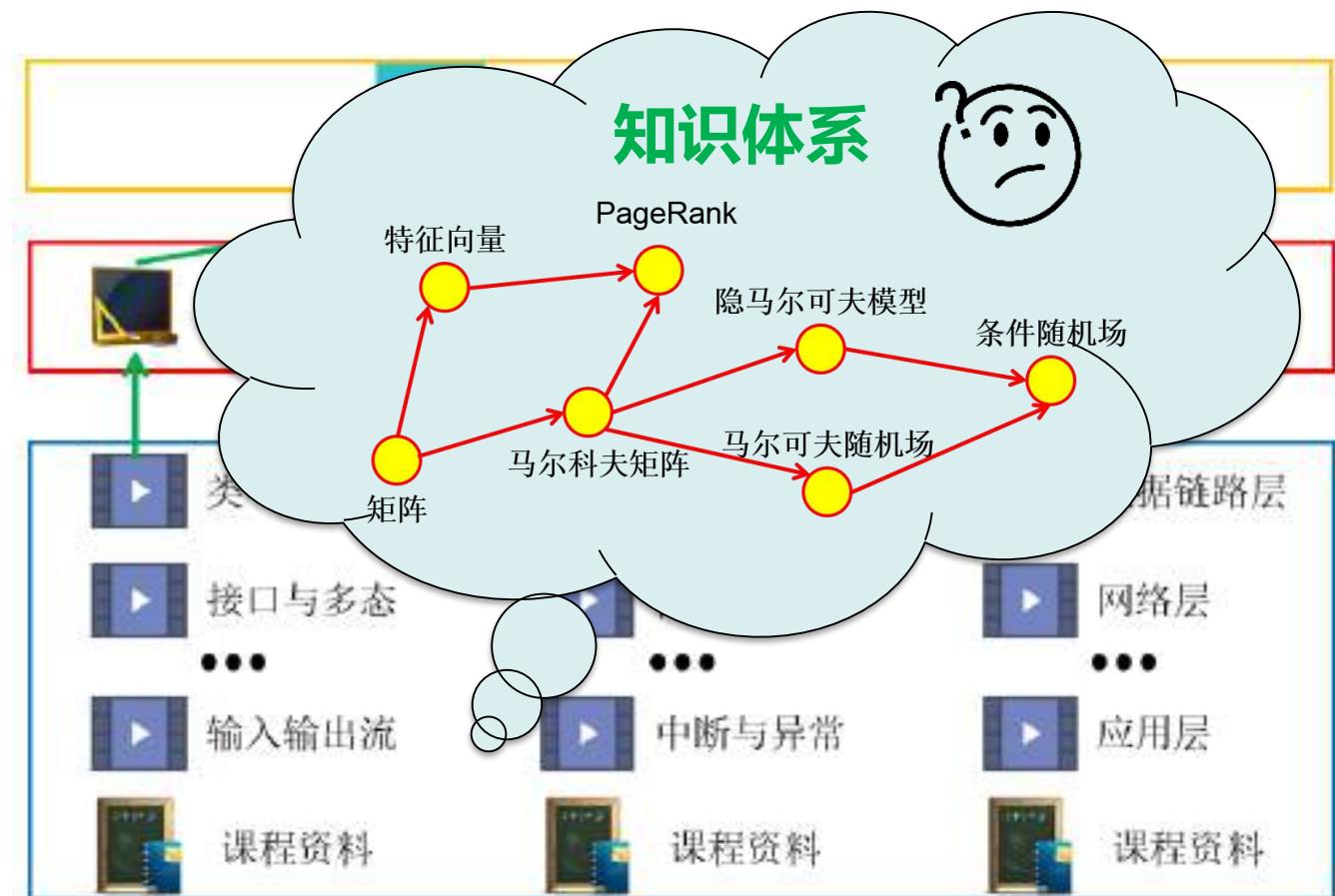
多源异构  
知识融合



# 大规模在线教育知识图谱构建

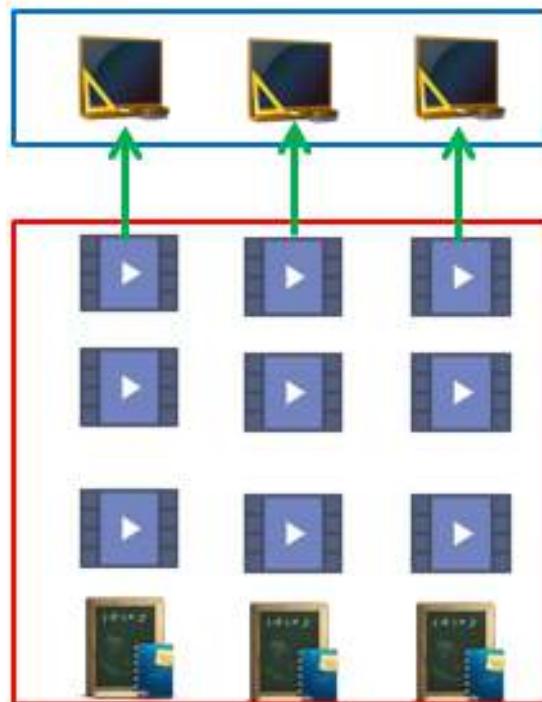
## 领域—课程—视频 三级组织

1 学科  
2 课程  
3 视频

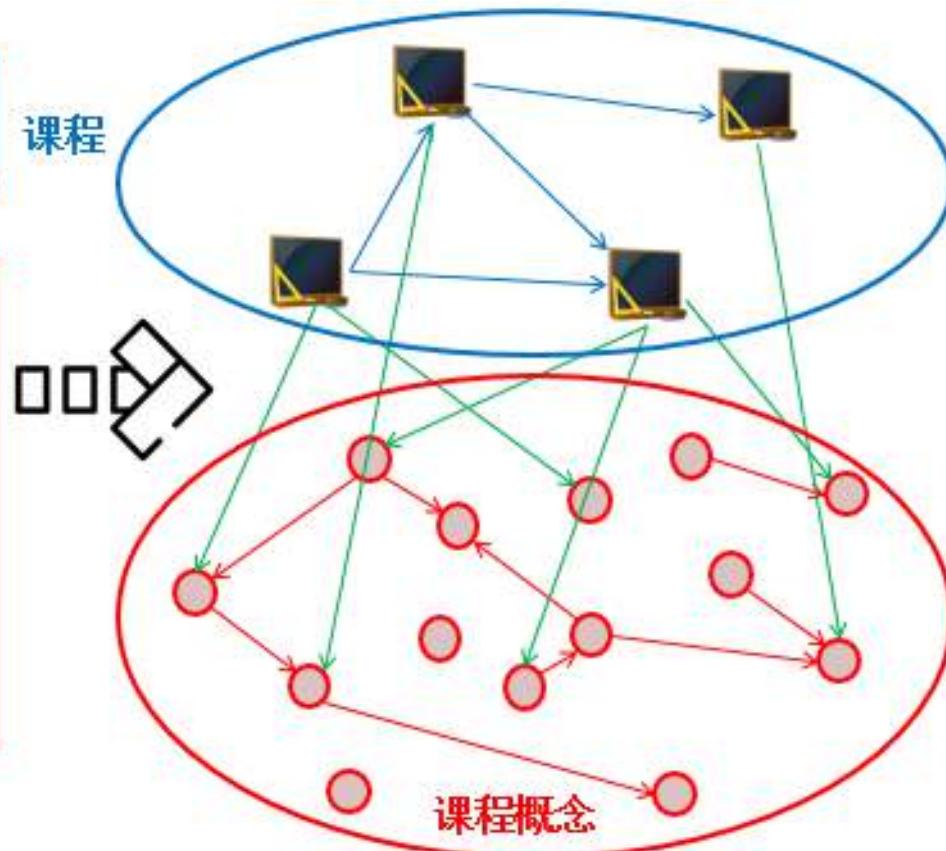


# 大规模在线教育数据：知识概念

视频为中心



知识为中心



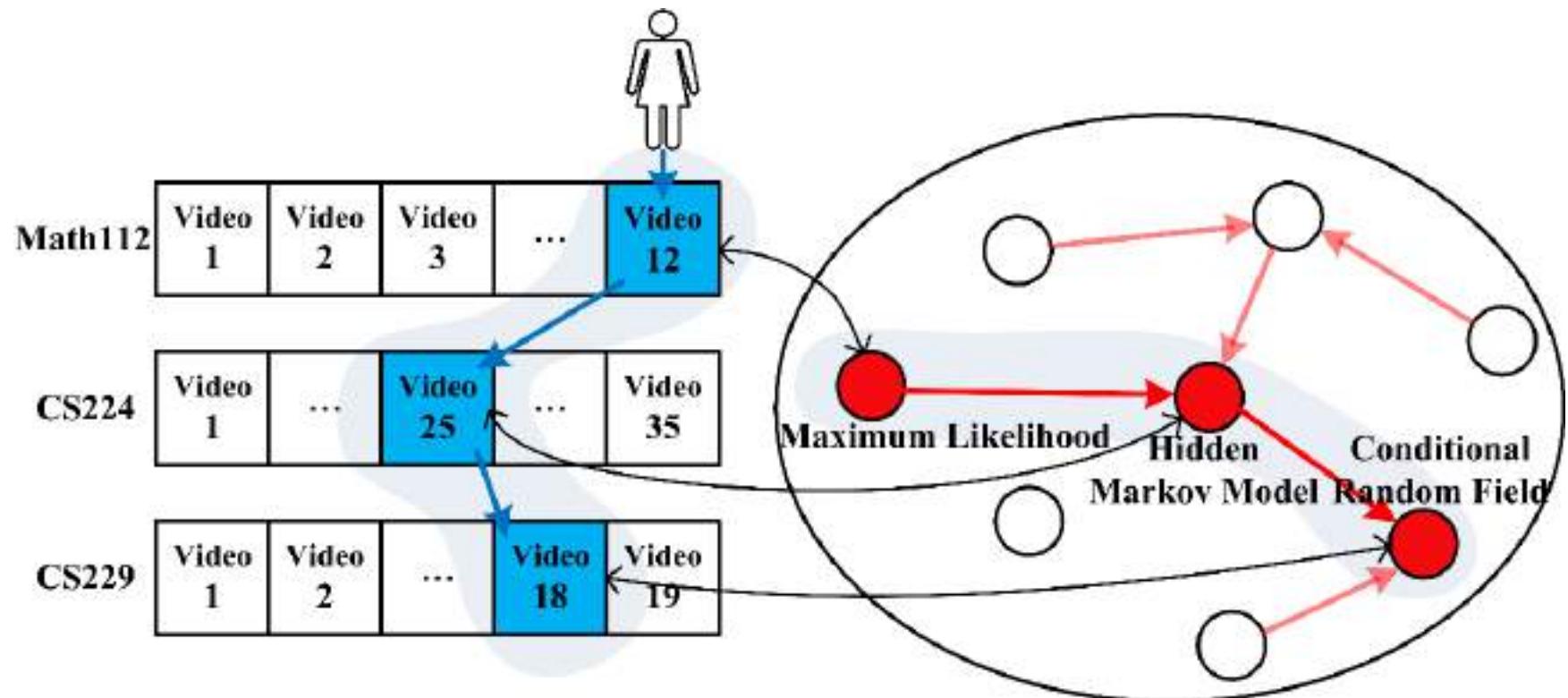
视频

课程概念

# 在线教育课程知识体系构造



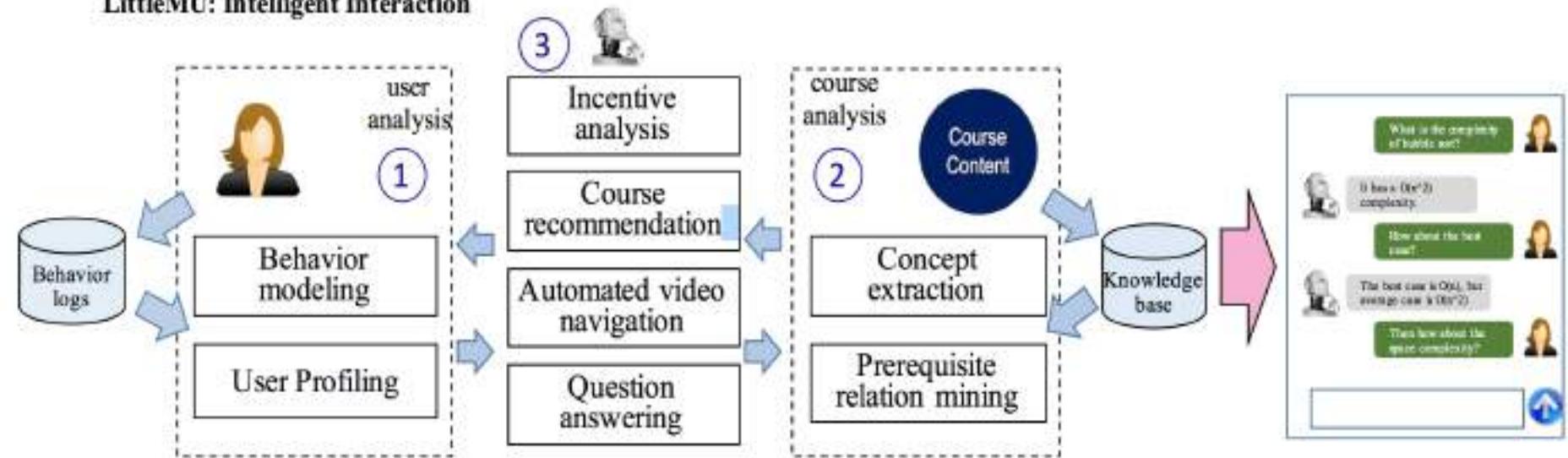
## 应用——智能学习路径规划



# 在线教育智能交互的知识服务

## 应用一智能助教

LittleMU: Intelligent Interaction





# 大规模在线教育知识图谱

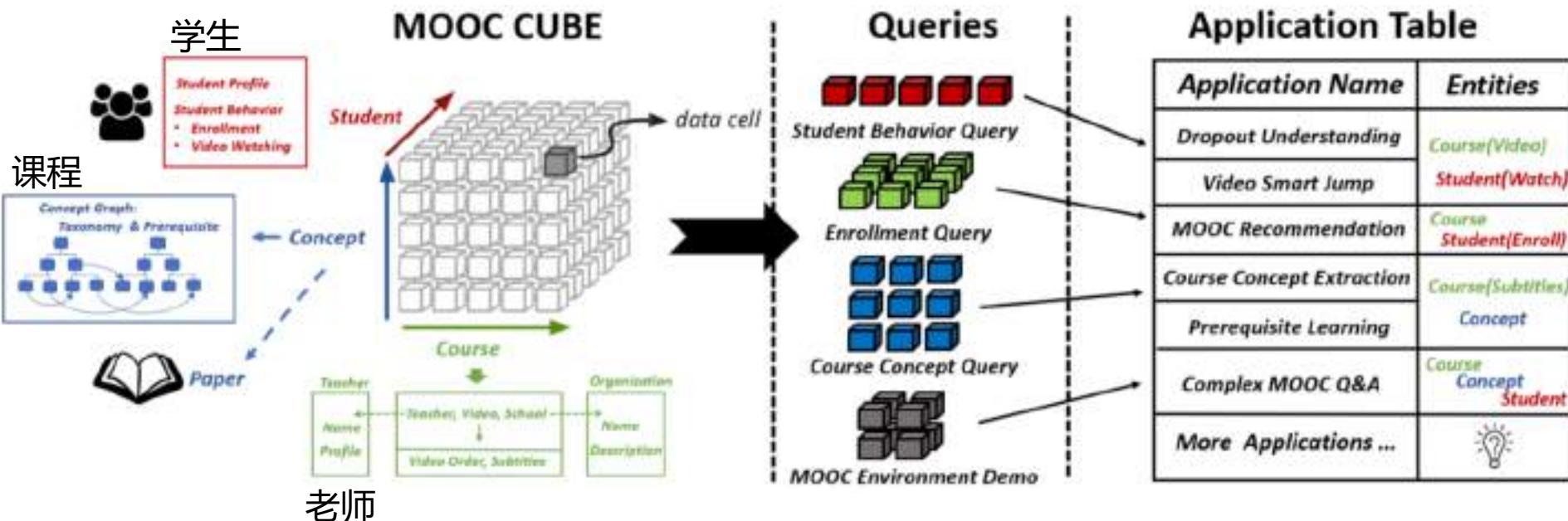


- 以知识为核心大规模在线教育资源组织
- 基于置信度网络传播的课程概念抽取
- 利用外部知识库的概念概念扩展
- 基于多维度特征的课程概念先后序关系学习

- [1] MOOCCube: A Large-scale Data Repository for NLP Applications in MOOCs, ACL2020
- [2] Course Concept Extraction in MOOCs via Embedding-Based Graph Propagation, IJCNLP2017
- [3] Course Concept Expansion in MOOCs with External Knowledge and Interactive Game, ACL2019
- [4] Prerequisite Relation Learning for Concepts in MOOCs. ACL2017

# 大规模在线教育数据立方体

- MOOCCube: 服务于MOOC相关研究的开源大规模数据仓库  
学生行为记录，课程视频资料，知识概念体系  
开源链接：<http://moocdata.cn/data/MOOCCube>



# 大规模在线教育数据：用户行为



## □ MOOCCube 细粒度学生行为记录



- Student Profile
- Student Behavior
  - Enrollment
  - Video Watching

### 用户个人画像

性别，年龄，所在地...

### 用户视频学习记录

(精确到毫秒)

学习时长

学习次数

观看视频的观看范围( $t_1, t_2$ )



199,199名学生，4,874,298人次观看视频记录

# 大规模在线教育数据：课程信息

## □ MOOCCube课程信息：

课程信息：

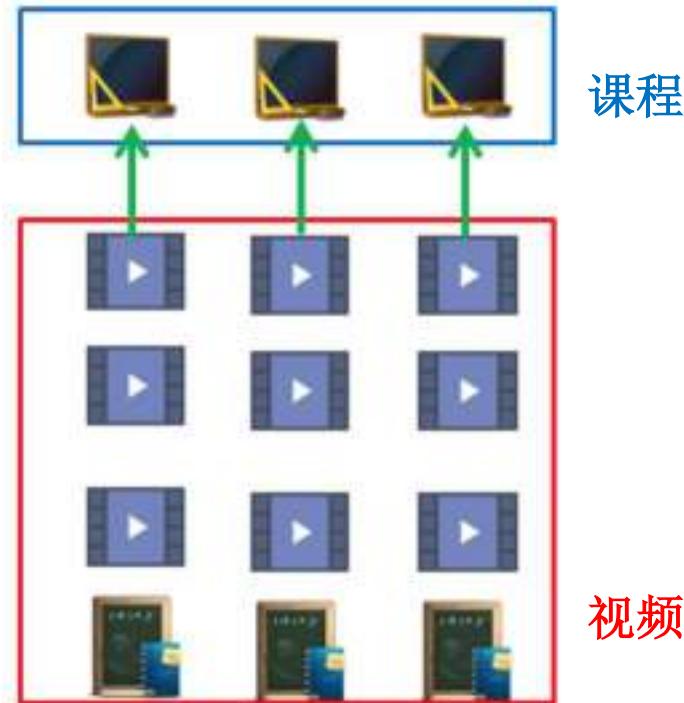
教师

开课机构

视频信息：

字幕文本

视频顺序

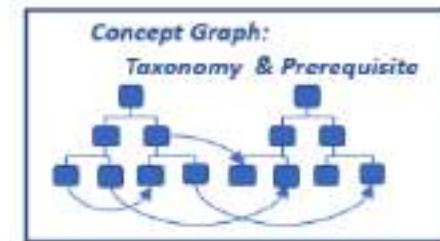


706门课程，38,181个视频

# 大规模在线教育数据：知识概念

## MOOCCube知识概念体系：

- 课程概念集
- 概念先后修关系
- 概念的相关资源(论文)



106,056个知识概念，并被进一步补充



# 在线教育课程知识图谱构建

## 课程概念 + 先后序关系 ☐□△ 知识体系

### 定义

- 课程视频中所教授的知识概念

### 例子

You might learn how to write a **bubble sort** and learn why a **bubble sort** is not as good as a **heapsort**. Next, we are going to talk about the **quick sort** algorithm. **Quicksort** is an algorithm invented in the 1960s by doctor Tony Hoare. It is also called the **partition exchange sort**, and is a typical algorithm based on **divide-and-conquer**.

.....

Now we have the first version of **Q sort**. After we make an analysis on its performance, performance, we will find that **quicksort** is an **unstable sorting algorithm**. Fortunately, the **quick sort** has an average **time complexity** of  $n \log n$ , and in most cases, it can achieve its optimal performance. We first estimate its performance under **independent uniform distribution**.

# 在线教育课程知识体系构造

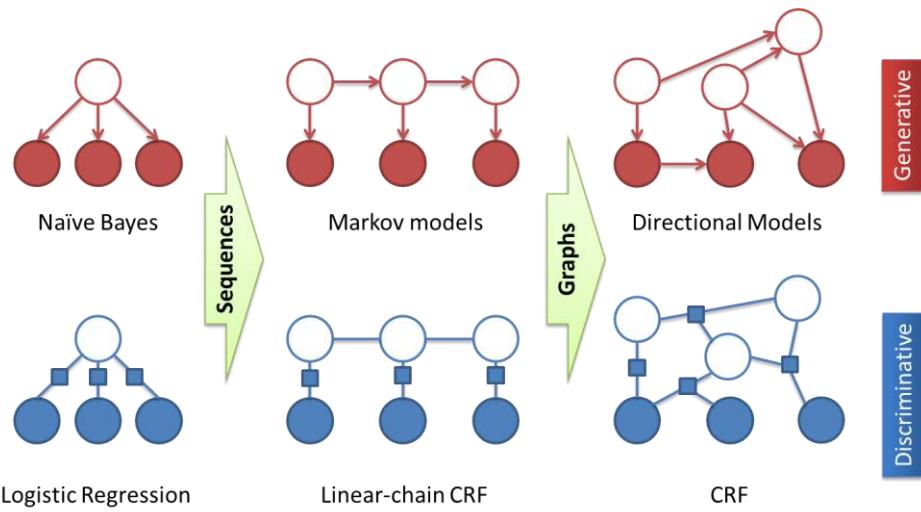


## 课程概念 + 先后序关系 □□□ 知识体系

### 定义

➤ 知识概念间的学习依赖关系

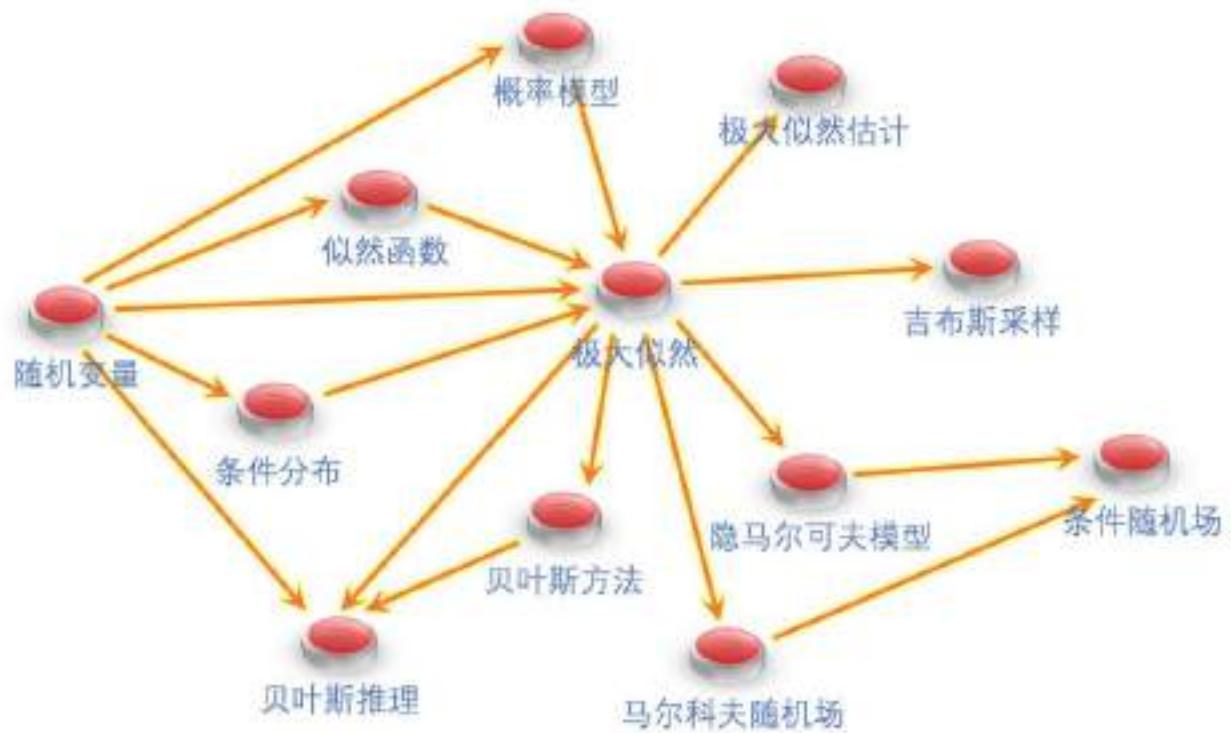
### 例子



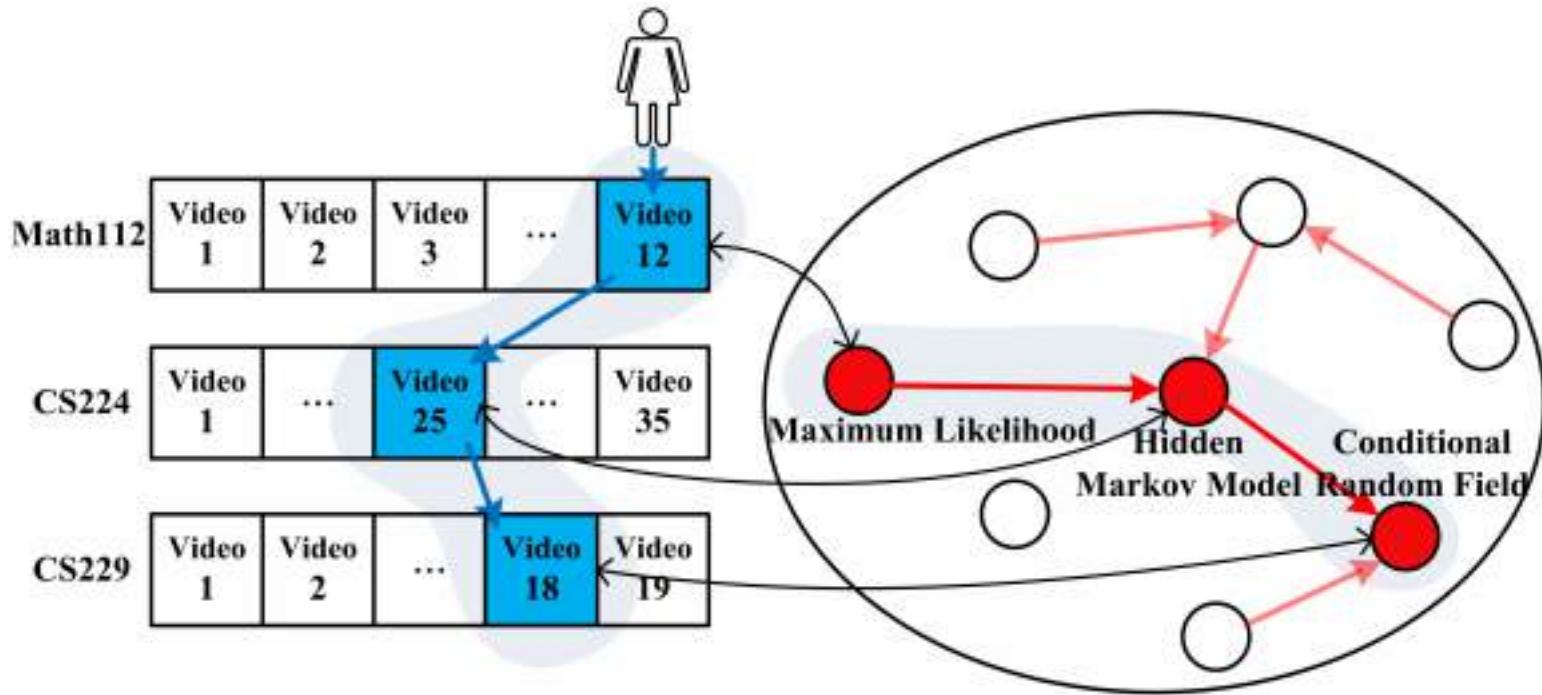
Adapted from C. Sutton, A. McCallum, "An Introduction to Conditional Random Fields", ArXiv, November 2010

# 在线教育课程知识体系构建

## 课程概念 + 先后序关系 知识体系



# 在线教育课程知识体系构建



- How to extract concepts from course scripts?
- How to recognize (prerequisite) relationships between concepts?

[1] Liangming Pan, Chengjiang Li, Juanzi Li, and Jie Tang. Prerequisite Relation Learning for Concepts in MOOCs. **ACL'17**.



# 课程概念学习

Candidate  
Concept  
Extraction

Semantic  
Representation  
Learning

Graph-  
based  
Ranking

In this course, we will teach some basic knowledge about **data mining** and its application in **business intelligence**.

Video script

data mining

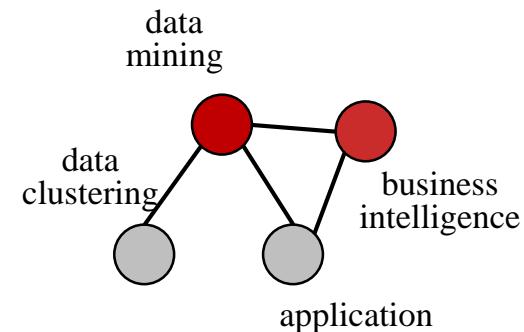
0.8	0.2	0.3	...	0.0	0.0
-----	-----	-----	-----	-----	-----

business intelligence

0.1	0.1	0.2	...	0.8	0.7
-----	-----	-----	-----	-----	-----

Vector representation

Learned via embedding or deep learning





# 课程概念抽取

## 课程概念的三个特性

### 短语性

➤ 课程概念应该是一个**语法正确的词组或短语**

### 信息性

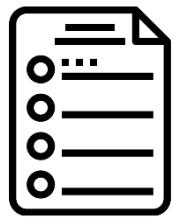
➤ 课程概念在语义上表示一个**科学或技术术语**

### 相关性

➤ 课程概念应该与MOOC课程语料中的**课程相关**

# 课程概念抽取

## 语义相关度

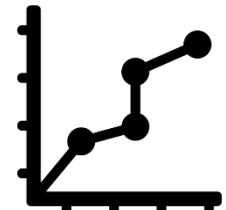


候选课程概念



维基百科语料

- **实体标注**
  - 对维基语料进行**实体标注**
- **词嵌入表示学习**
  - 学习维基语料的**词嵌入表示**
- **候选课程概念表示**
  - 通过**向量加和**得到候选概念向量表示
- **语义相关度度量**
  - **余弦相似度**作为候选概念间语义相关度

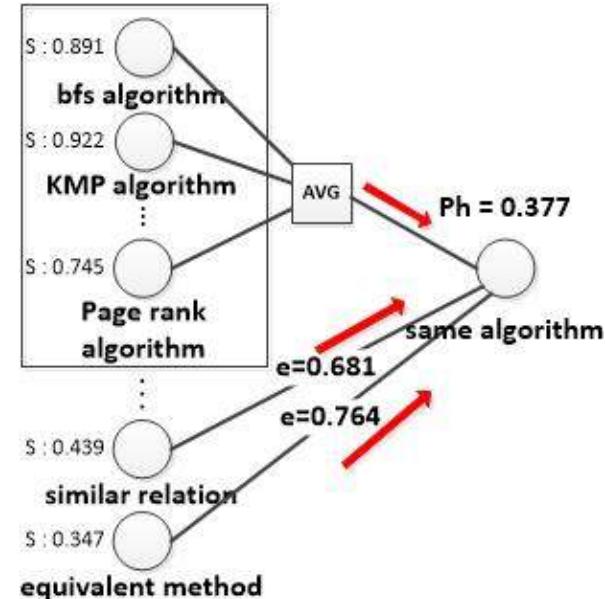
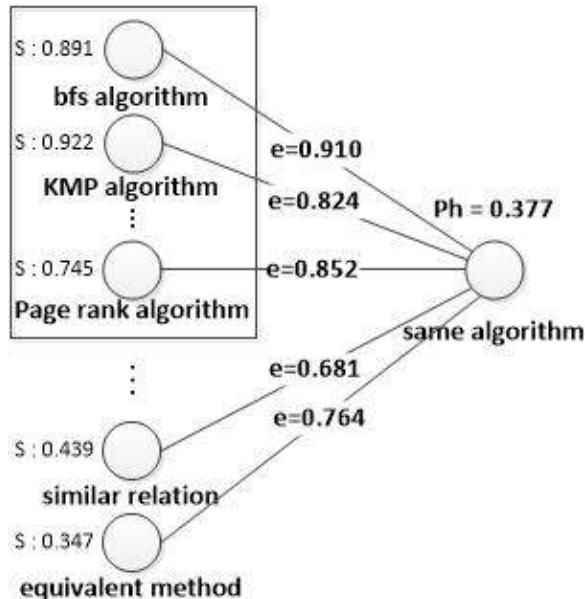
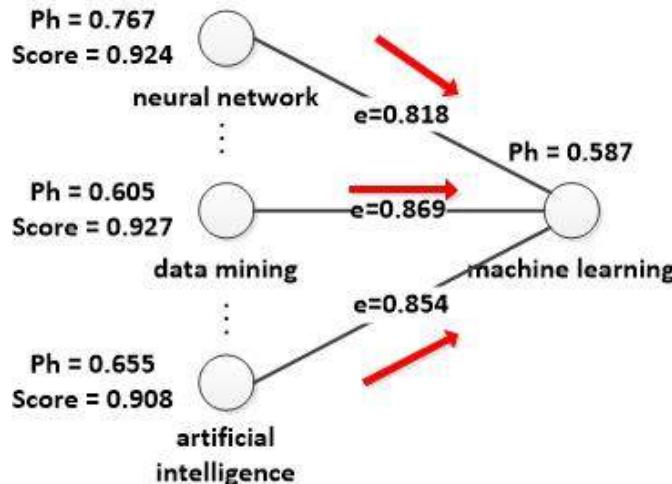


语义相关度

# 课程概念抽取

## 置信度图传播

- 课程概念图中，课程概念与其他课程概念以高语义相关度相连接
- 可以通过少数种子概念通过置信度传播找出更多课程概念





# 课程概念扩展

## 大量课程相关概念并不出现在课程文本中

### *Course subtitles*

... there should be data structures that would give it to you; probably a **binary search tree** is the first thing you'd want to consider...

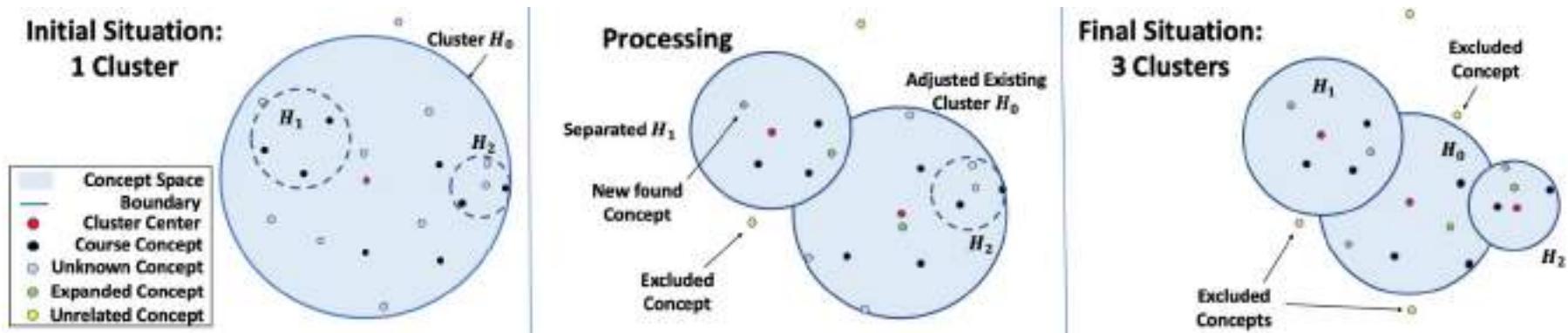
### *Hidden related concepts*

*Heap      Tango Tree  
Sorting Algorithms  
Priority Queue*

这些概念可能由于多种原因能够帮助当前学习

# 课程概念扩展

## 利用外部知识库中的概念，进行课程概念扩展

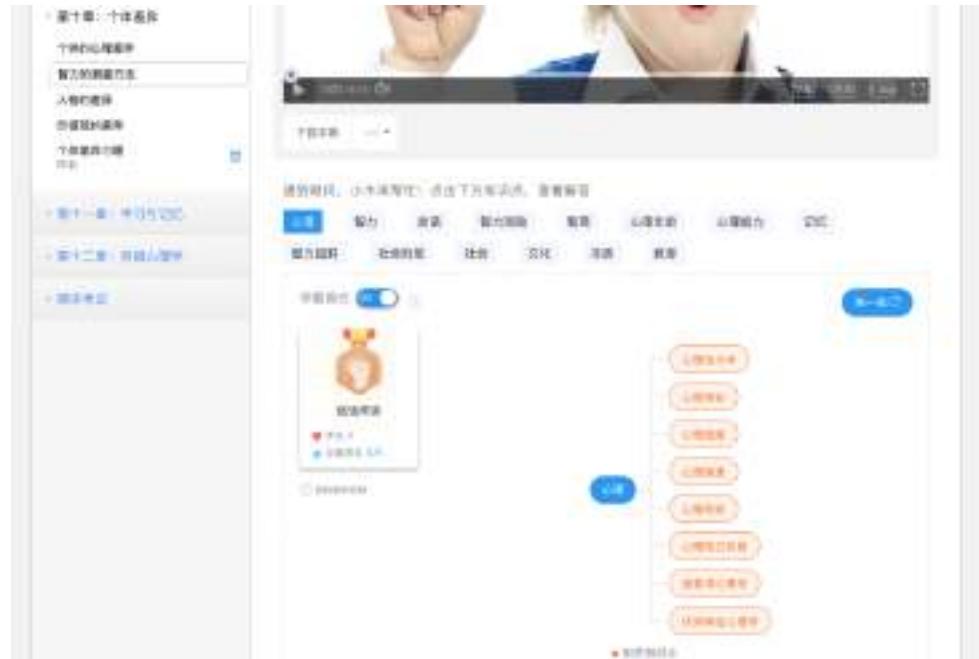


- 语义漂移是课程概念扩展任务中一个重要的挑战，因此使用了基于动态聚类的概念扩展方法



# 课程概念扩展

同时，设计在线互动小游戏，收集用户反馈



- “学霸模式” 小游戏，放置于课程视频下，吸引用户找出相关概念



# 前后序概念关系学习

## 语义特征



## 上下文特征



## 结构特征





# 语义特征

## 语义相关度

- 语义相关度在判断概念对是否具有先后序关系上具有重要作用
- 若两个概念语义相关度很低，则它们之间具有学习依赖性的概率较低
- 采用与课程概念抽取任务中相同的方法度量语义相关度

矩阵  人类学

梯度下降  后向传播算法

# 上下文特征

## 引用距离

- 核心假设：在教师教授概念A时，如果教师经常提及概念B，反过来，在教师教授概念B时，概念A却很少被提及，则概念B很可能 是概念A的先序概念

### 后向传播算法



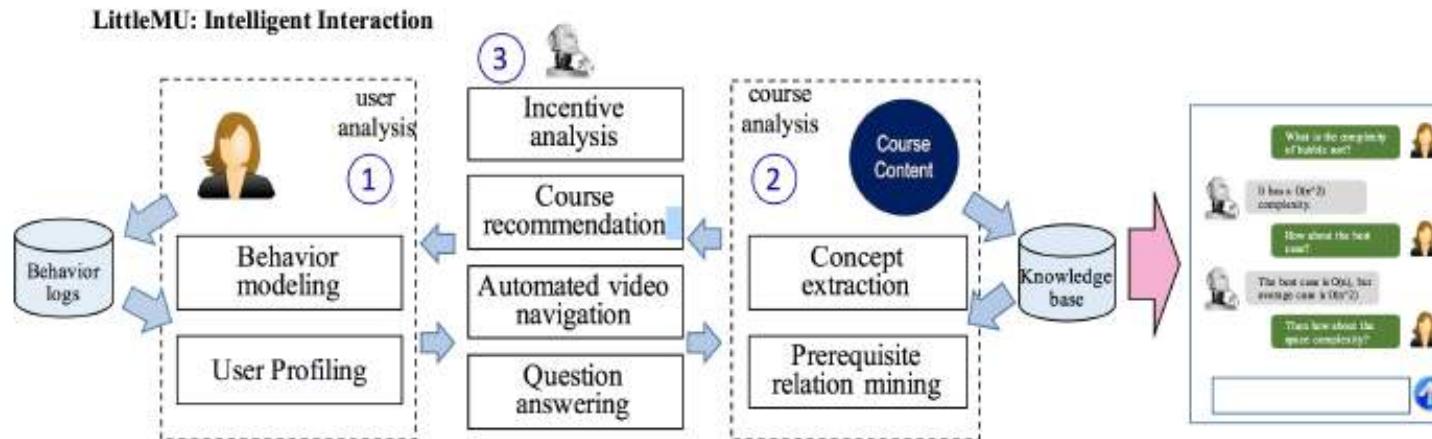
梯度下降

### 梯度下降

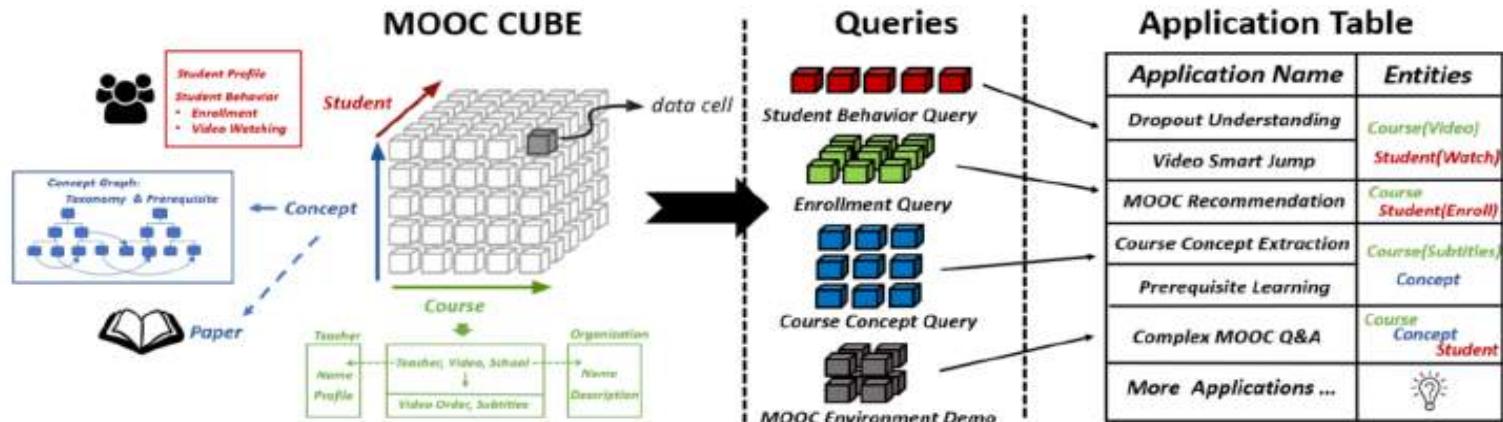


后向传播算法

# 在线教育智能交互的知识服务



## 学堂智能知识机器人“小木”



## 大规模在线教育数据仓库 MOOCCube



# 总 结

- 知识图谱是人工智能的基础设施
- 知识图谱构建需要高质量大规模知识获取技术支持
- 知识驱动的智能学术搜索技术是学术搜索的新型态



# 谢谢大家!

# 大数据知识管理服务平台技术与实践

杜一

中国科学院计算机网络信息中心 大数据知识工程实验室

2020年6月2日

# 学术搜索

Search any topic, author, journal, etc. or any combination of these



Google Scholar



Baidu 学术

# 学术搜索

Search any topic, author, journal, etc. or any combination of these



Google Scholar



Baidu 学术



Academia.edu

AMiner

# 学术搜索

Search any topic, author, journal, etc. or any combination of these



Google Scholar

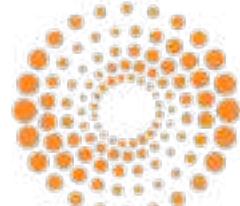
Microsoft Academic Search Beta

Baidu 学术

Semantic Scholar

Academia.edu

AMiner



WEB OF SCIENCE

Scopus®

万方数据  
WANFANG DATA

Cnki 中国知网  
www.cnki.net  
中国知识基础设施工程

# 学术搜索

DOI: 10.1145/3038912.3062558 · Corpus ID: 1644336

Share This Paper

## Explicit Semantic Ranking for Academic Search via Knowledge Graph Embedding

论文结论

Chirantan Gangopadhyay, Russell Power, James P. Callan · Published in WWW 2017 · Computer Science

This paper introduces Explicit Semantic Ranking (ESR), a new ranking technique that leverages knowledge graph embedding. [...] Experiments demonstrate ESR's ability in improving Semantic Scholar's online production system, especially on tailqueries where word-based ranking fails. [\[+Expand Abstract\]](#)

[View On ACM](#)

[Alternative Sources](#)

[e-Library](#)

[Cite Alert](#)

[Cite](#)

[In Feed](#)



Semantic Scholar

ABSTRACT

FIGURES, TABLES, AND TOPICS

CITATIONS

REFERENCES

RELATED PAPERS

Figures, Tables, and Topics from this paper.

论文主题

### Figures and Tables

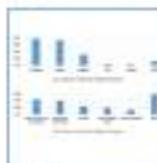


Figure 1

Table 1



Figure 2

Table 2

[SHOW MORE \[N\]](#)

### Explore Further: Topics Discussed in This Paper

- Knowledge Graph
- Graph embedding
- Academic Search
- Semantic Scholar
- Production systems (computer science)
- Web search engine

# 学术搜索

DOI: 10.1145/3039912.3062558 · Corpus ID: 1644336

## Explicit Semantic Ranking for Academic Search via Knowledge Graph Embedding

论文结论

Chenyan Wang, Russell Power, James P. Callan · Published in WWW 2017 · Computer Science

This paper introduces Explicit Semantic Ranking (ESR), a new ranking technique that leverages knowledge graph embedding. [...] Experiments demonstrate ESR's ability in improving Semantic Scholar's online production system, especially on tailqueries where word-based ranking fails. [SI Expand Abstract]

87 Citations

3 Highly Influential Papers

64 Cite Background

30 Cite Methods

2 Cite Results

引文分类



Semantic Scholar

ABSTRACT

FIGURES, TABLES, AND TOPICS

CITATIONS

REFERENCES

RELATED PAPERS

Figures, Tables, and Topics from this paper.

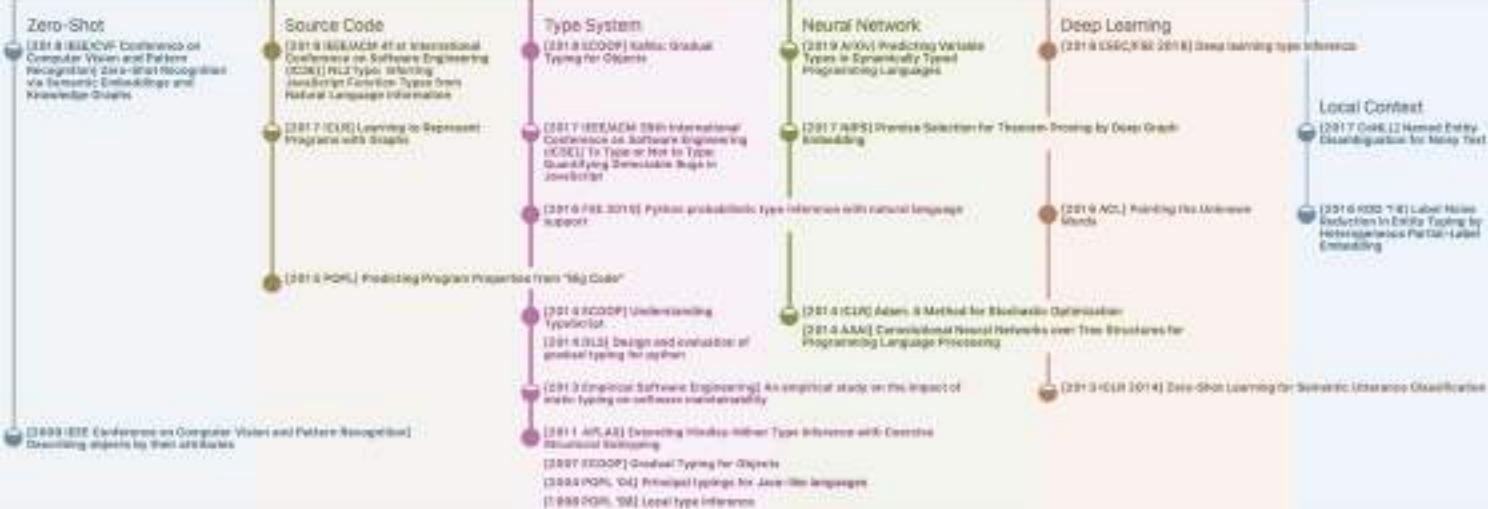
论文主题

### Figures and Tables

[2020 ICSE 2020] LambdaNet: Probabilistic Type Inference using Graph Neural Networks

Explore Further: Topics Discussed in This Paper

- Knowledge Graph
- Graph embedding
- Academic Search
- Semantic Scholar
- Production systems (computer science)
- Web search engine



Source Code

Type System

Neural Network

Deep Learning

Local Context

Aminer 论文溯源树

# 大数据知识管理服务平台



空间科学领域知识图谱与决策  
支持系统



基于知识图谱的科技合作智能  
化管理服务平台



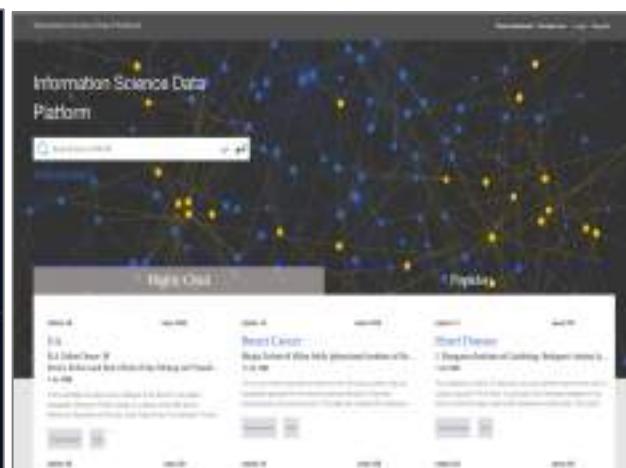
中科院国际合作数据汇聚与分  
析平台



国家自然科学基金大数据知识  
管理服务平台



烟草科技知识图谱服务平台



Information Science Data  
Platform

# 大数据知识管理服务平台

The screenshot shows the homepage of the Big Data Knowledge Management Platform. At the top, there are several logos and language selection buttons (Chinese, English). Below the header, a navigation bar includes links for '空间科学领域知识图谱与决策支持' (Knowledge Graph and Decision Support for the Space Science Field), '基础知识' (Basic Knowledge), '合作网络' (Collaboration Network), '关联路径' (Associated Path), '科研社区' (Research Community), '研究热点' (Research Hotspots), and '态势感知' (Trend Perception). A search bar is located above the main content area. On the left, a sidebar titled '人才信息' (Personnel Information) displays details for a user named '王伟', including their affiliation with '中国科学院国家空间科学中心' and an influence score of '1230'. It also lists 11 research publications. The main right-hand section features a large network graph visualization where nodes represent different entities and connections show relationships. A legend on the far right identifies the colors used for the nodes: orange for '人才' (Personnel), blue for '机构' (Institution), green for '学科' (Discipline), pink for '项目' (Project), red for '会议' (Conference), and yellow for '专著' (Monograph).

空间科学领域知识图谱与决策支持

基础知识 合作网络 关联路径 科研社区 研究热点 态势感知

王伟

人才信息

姓名: 王伟  
机构: 中国科学院国家空间科学中心  
影响力: 1230

项目成果列表

1. 日地物质抛射事件监测日历...【项目】
2. 地球风-磁层相互作用全貌...【项目】
3. 子午工程/空间天气预报系...【论文】
4. 无预警、无自由参数模型在...【论文】
5. 基于中性的ARIMA...【论文】
6. 基于TFIDF模型的极光定位...【论文】
7. 太阳活动与全球气候变化...【论文】
8. 火星空间环境探测设计与实...【论文】
9. 紫光一号火星探测器有效载...【论文】
10. 中国科研信息化蓝皮书20...【专著】
11. 进入太空——日地空间探测【专著】

人才 机构 学科 项目 会议 专著

# 大数据知识管理服务平台



# 大数据知识管理服务平台

来自内部的需求逐渐增多

基金委

科学院

空间中心

烟草总局

出版社

....

交叉学科

学术团队

资助方向

专家画像

热词矩阵

研究热点

人才合作

机构画像

# 大数据知识管理服务平台

来自内部的需求逐渐增多

学术深度挖掘需求迫切

人员关系

科研社区

学术画像

人才引进

学科交叉

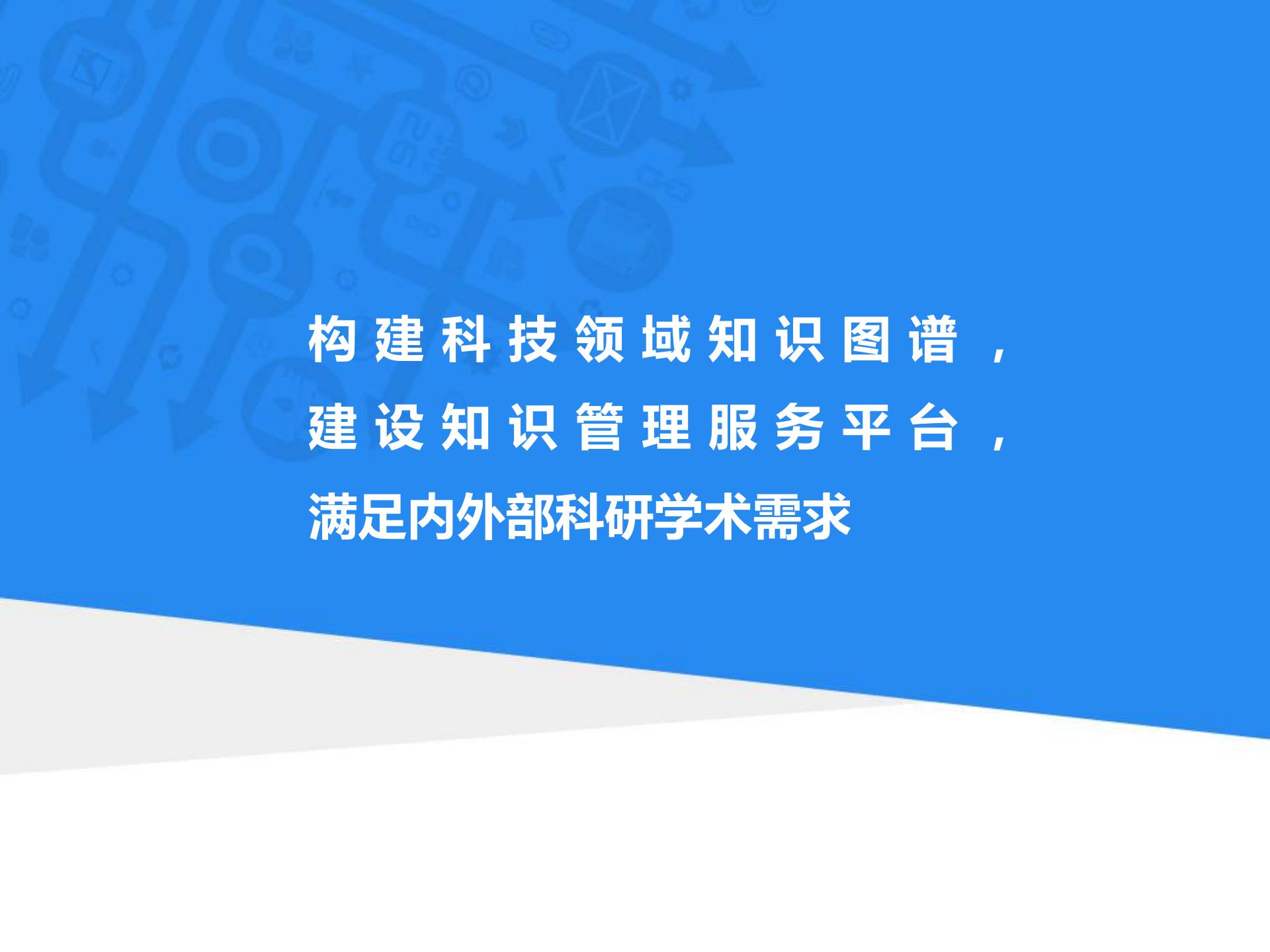
... ...

资助方向

专家画像

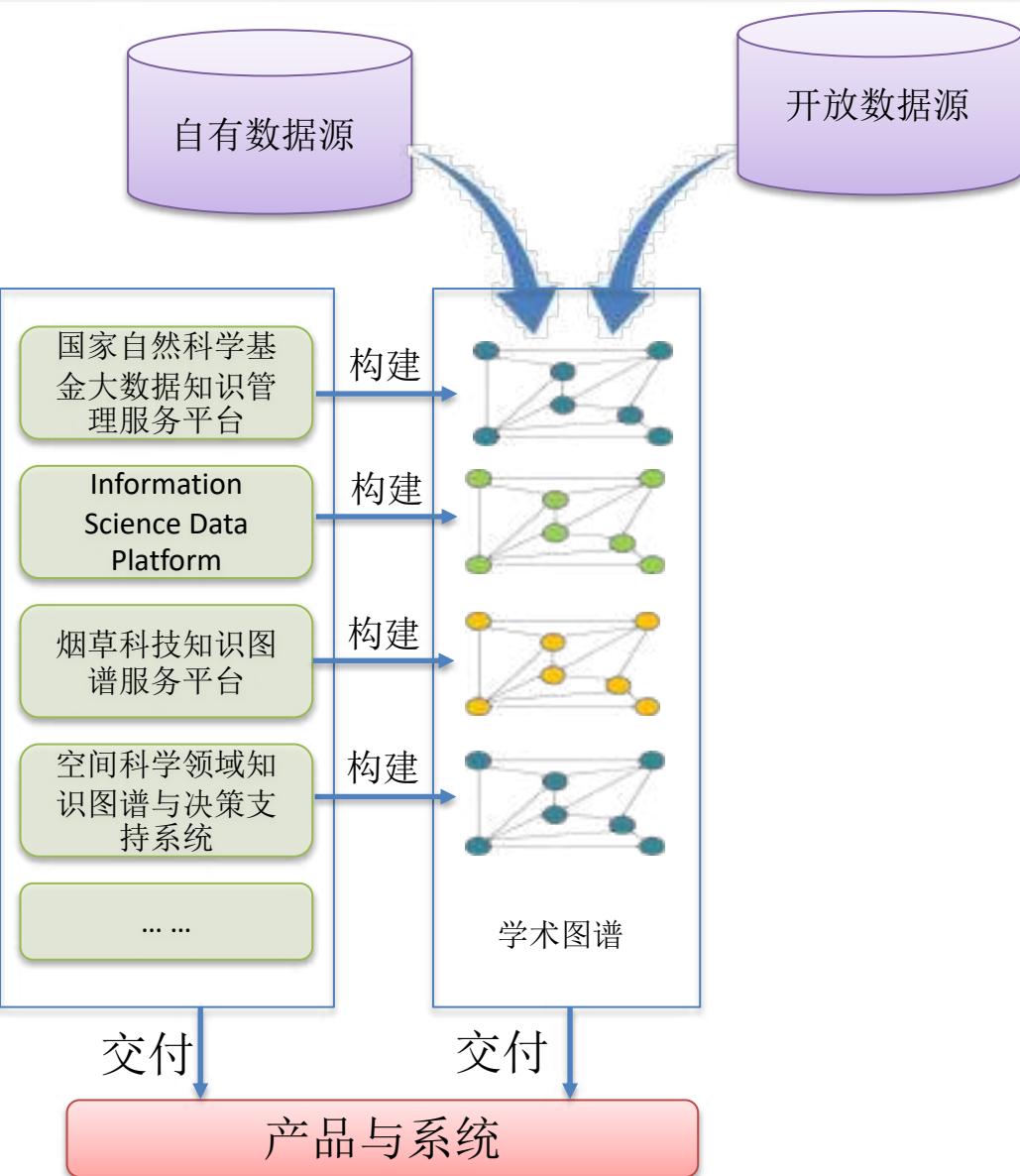
人才合作

机构画像

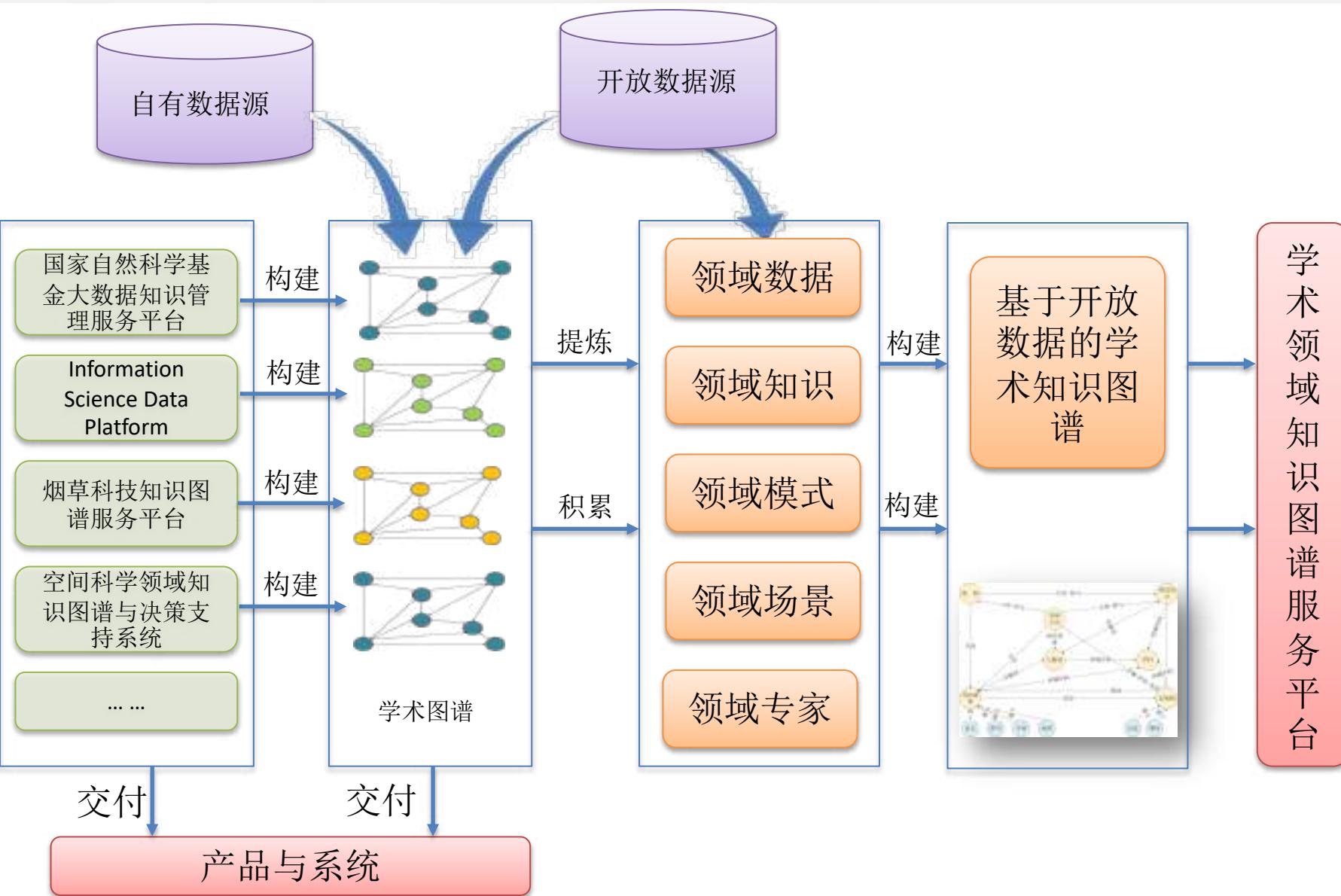


构建科技领域知识图谱，  
建设知识管理服务平台，  
满足内外部科研学术需求

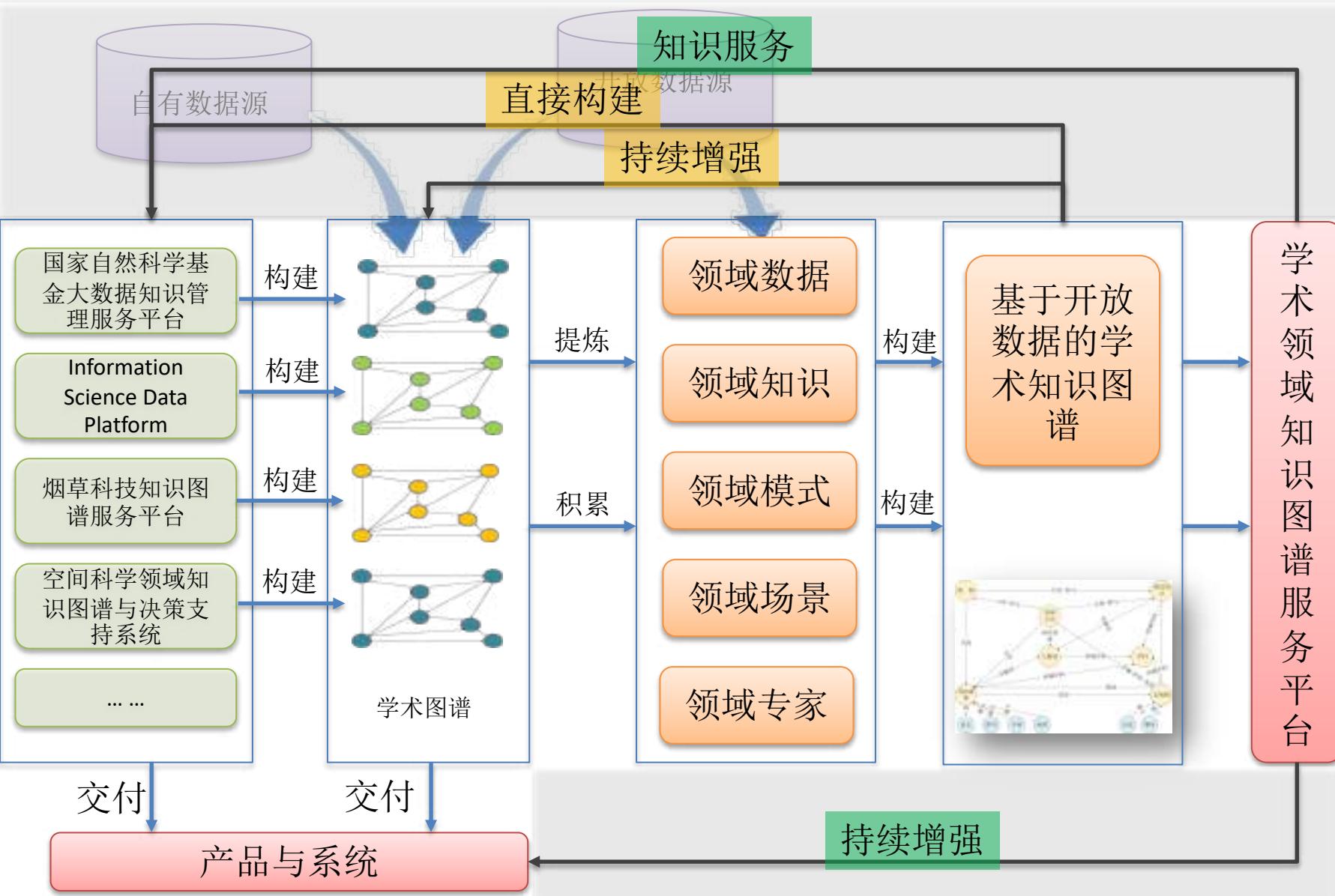
# 科技领域知识图谱



# 科技领域知识图谱

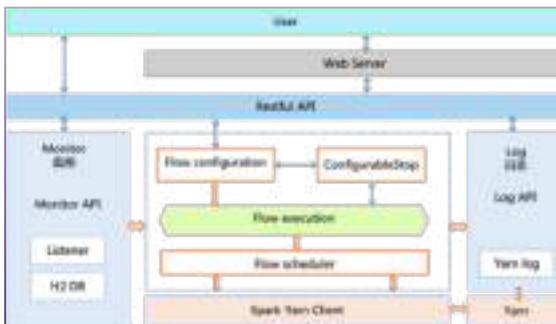


# 科技领域知识图谱



# 科技领域知识图谱：关键技术

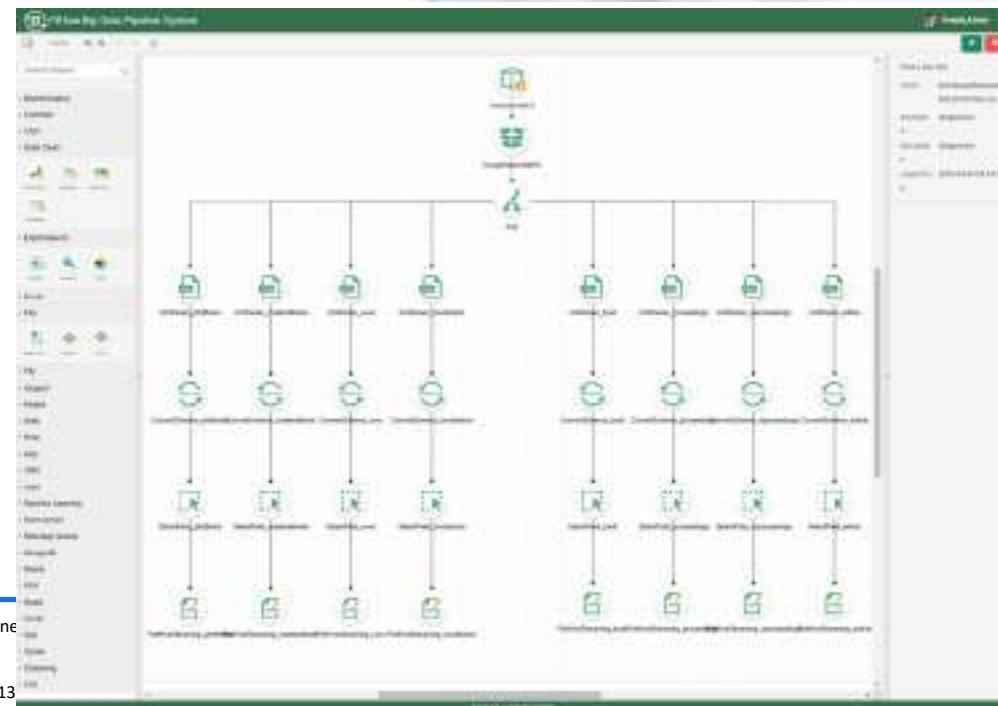
## 自动及半自动的数据采集、处理流水线机制



PiFlow 架构



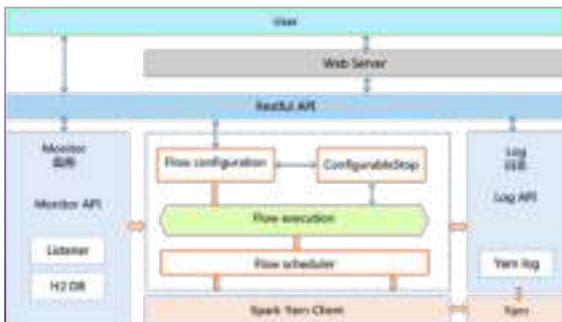
流水线实例



1. 乔子越,杜一\*,傅衍杰,王鹏飞,周园春,Unsupervised Author Disambiguation using Heterogeneous 2019年“智源-AMiner姓名排歧大赛”第一名
2. 乔子越,杜一,王寒雪,周园春,一种基于网络表征和语义表征的同名作者消歧方法,20191113
3. PiFlow: 阿云2019年最有价值开源项目;中国科技云首届开源大赛二等奖
4. 杜一,王寒雪,乔子越,周园春,一种基于异质图卷积神经网络嵌入的作者名字消歧方法,2019106357954,2019
5. 杜一,乔子越,周园春,一种基于异质网络嵌入的学者名字消歧方法,2018112671819,2018
6. 乔子越,周园春,一种基于异质网络嵌入的学者名字消歧方法,2018112671819,2018
7. 中国专利: 一种大数据ETL任务的编排方法与系统. 完成人: 朱小杰,沈志宏,杜一,赵子豪,周园春,2019
8. 中国专利: 一种大数据ETL任务的调度方法. 完成人: 朱小杰,沈志宏,杜一,赵子豪,周园春,2018

# 科技领域知识图谱：关键技术

## 自动及半自动的数据采集、处理流水线机制

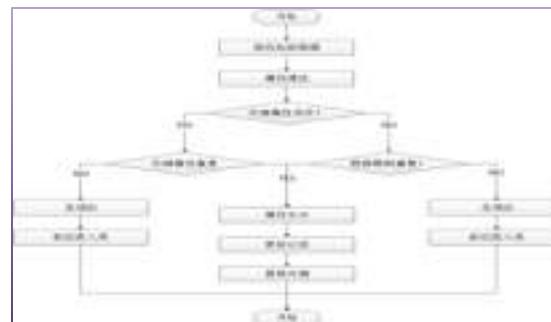


PiFlow 架构

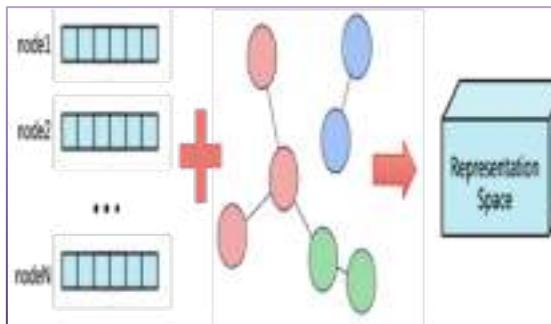


流水线实例

## 基于统计规则及深度学习方法的实体融合方法

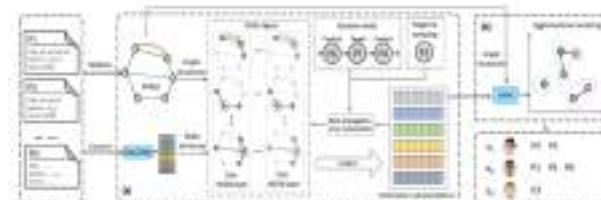


基于规则组合的实体融合



基于GraphEmbedding实体融合

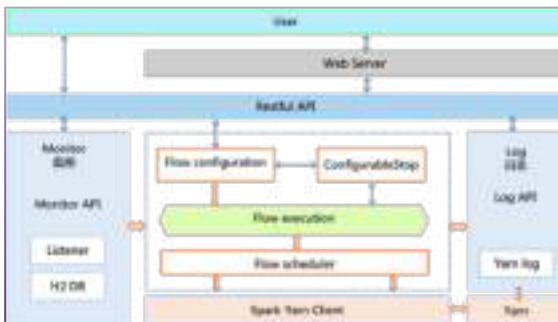
Name	Our method	Components	Zhang et al. 2018	Xu et al. 2018	Zhang et al. 2017
Ajay Gupta	<b>0.750</b>	0.329	0.348	0.352	0.638
Alok Gupta	1	0.490	0.689	0.892	0.540
Ben Yu	<b>0.696</b>	0.295	0.434	0.583	0.616
David Cooper	0.900	0.327	0.737	0.884	<b>0.931</b>
David Nelson	<b>0.944</b>	0.219	0.750	0.735	0.556
Esi Su	1	0.648	0.952	0.630	0.941
Fan Wang	<b>0.604</b>	0.096	0.403	0.557	0.545
Jie Tang	<b>0.982</b>	0.883	0.857	0.522	0.910
Thomas Wolf	<b>0.860</b>	0.502	0.709	0.522	0.352
Yang Wang	0.528	0.118	0.273	<b>0.574</b>	0.409
Avg.	<b>0.786</b>	0.307	0.715	0.681	0.680



- 乔子越,杜一\*,傅衍杰,王鹏飞,周园春,Unsupervised Author Disambiguation using Heterogeneous Graph Convolutional Network Embedding, IEEE Big Data, 2019
- 2019年“智源-AMiner姓名排歧大赛”第一名
- PiFlow: 码云2019年最有价值开源项目;中国科技云首届开源大赛二等奖
- 杜一,王寒雪,乔子越,周园春,一种基于网络表征和语义表征的同名作者消歧方法,2019113223833 | PCT/CN2019/128642, 2019
- 杜一,乔子越,周园春,一种基于异质图卷积神经网络嵌入的作者名字消歧方法,2019106357994, 2019
- 杜一,乔子越,周园春,一种基于异质网络嵌入的学者名字消歧方法,2018112671819, 2018
- 中国专利: 一种大数据ETL任务的编排方法与系统. 完成人: 朱小杰,沈志宏,杜一,赵子豪,周园春,2019
- 中国专利: 一种大数据ETL任务的调度方法. 完成人: 朱小杰,沈志宏,杜一,赵子豪,周园春,2018

# 科技领域知识图谱：关键技术

## 自动及半自动的数据采集、处理流水线机制

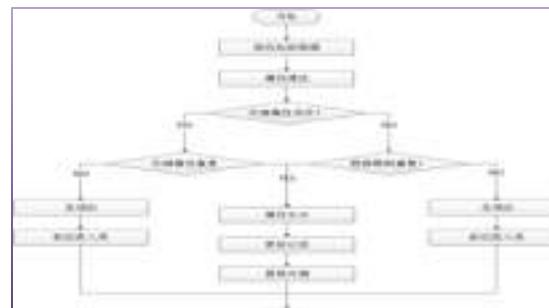


PiFlow 架构

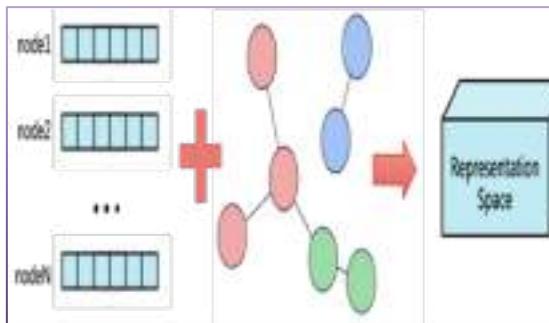


流水线实例

## 基于统计规则及深度学习方法的实体融合方法

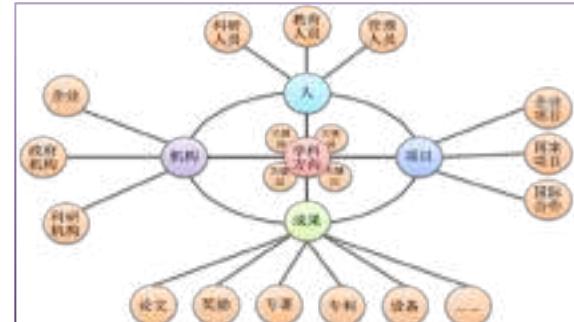


基于规则组合的实体融合

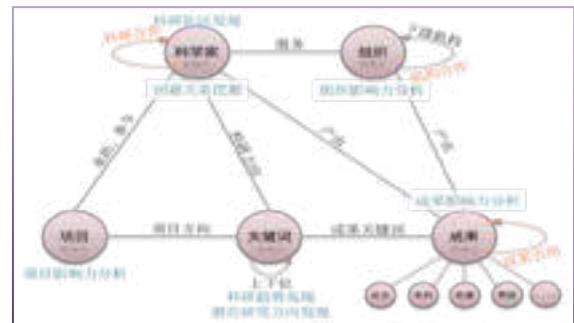


基于GraphEmbedding实体融合

## 基于属性图的科技领域知识图谱构建方法



基于规则的学术知识图谱



基于推理学习的学术知识图谱

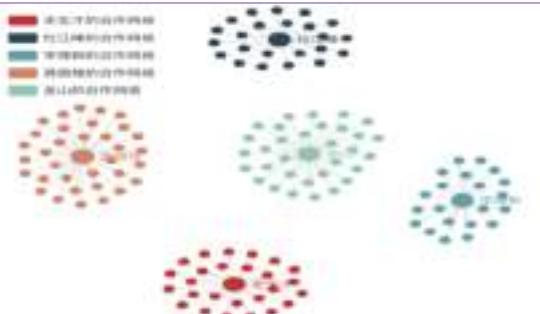
- 乔子越,杜一\*,傅衍杰,王鹏飞,周园春,Unsupervised Author Disambiguation using Heterogeneous Graph Convolutional Network Embedding, IEEE Big Data, 2019
- 2019年“智源-AMiner姓名排歧大赛”第一名
- PiFlow: 码云2019年最有价值开源项目;中国科技云首届开源大赛二等奖
- 杜一,王寒雪,乔子越,周园春,一种基于网络表征和语义表征的同名作者消歧方法,2019113223833 | PCT/CN2019/128642, 2019
- 杜一,乔子越,周园春,一种基于异质图卷积神经网络嵌入的作者名字消歧方法,2019106357994, 2019
- 杜一,乔子越,周园春,一种基于异质网络嵌入的学者名字消歧方法,2018112671819, 2018
- 中国专利: 一种大数据ETL任务的编排方法与系统. 完成人: 朱小杰, 沈志宏, 杜一, 赵子豪, 周园春, 2019
- 中国专利: 一种大数据ETL任务的调度方法. 完成人: 朱小杰, 沈志宏, 杜一, 赵子豪, 周园春, 2018

# 科技领域知识图谱：关键技术

## 基于学科的科研社区发现算法



基于PageRank+SVD的高影响力人员关系网络挖掘



基于网络度量指标+标签传播的科研社区发现

## 学科交叉性评估及学科趋势预测方法

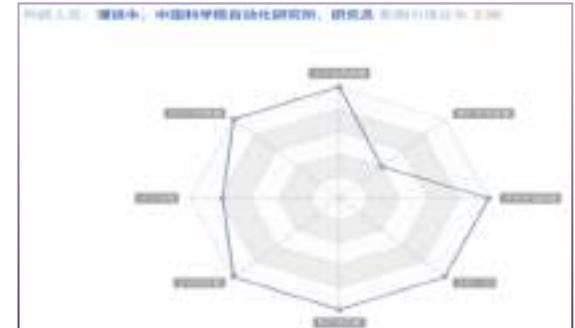


基于统计指标的学科交叉性评价

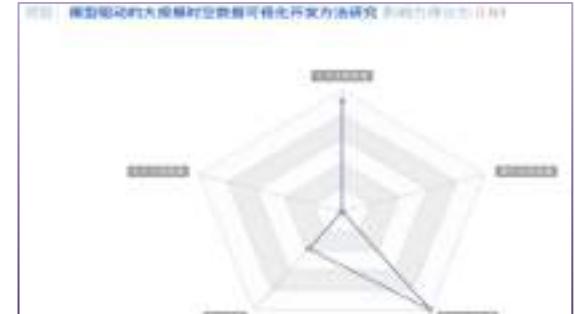


基于时序+图数据库的学科趋势评估

## 学术影响力评价方法



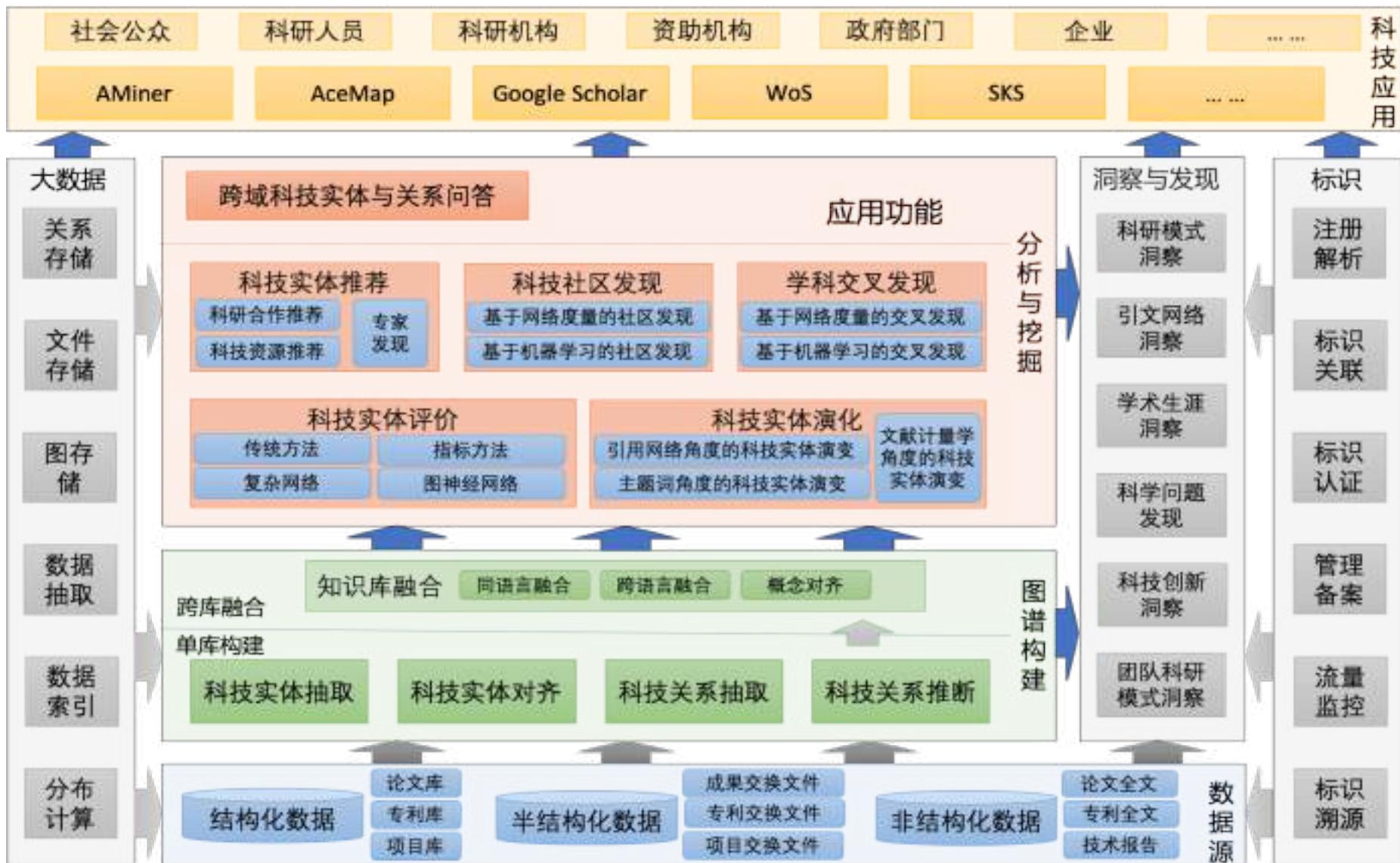
基于网络度量指标+计量学的科研人员影响力评价



基于网络度量指标+计量学的科研项目影响力评价

- 王卫军, 崔文娟, 杜一\*, 周园春, 基于词嵌入的国家自然科学基金学科交叉知识点分析方法, 《情报学报》, 2020. (in submission)
- 姚畅, 王晓帆, 杜一, 张兆田, 李建军, 郝艳妮, 国家自然科学基金大数据知识管理服务平台总体方案及关键技术研究, 中国科学基金, 2019.

# 科技领域知识图谱：研究架构



请批评指正！



# AMiner

## —Author-centric academic search and mining

Peng Zhang  
Tsinghua University

# Academic search is far from sufficient...



# Academic search is far from sufficient...

Google Scholar search results for "data mining". The results page shows various academic papers and books. A specific result is highlighted: "Advances in knowledge discovery and data mining" by Jiawei Han, Micheline Kamber, and Jian Pei. The page includes filters like "Exclude self-cited papers" and "Exclude patents".

Springer Academic Search results for "data mining". The interface shows a search bar, filters, and a results summary. One result is highlighted: "Data Mining - DM" by Jiawei Han, Micheline Kamber, and Jian Pei. The page includes a timeline chart showing the growth of publications from 1981 to 2010.

EBSCOhost search results for "data mining". The results page shows various academic papers and books. A specific result is highlighted: "The WEKA data mining software: an update" by Mark Hall, Eibe Frank, and Bernhard Pfahringer. The page includes filters like "Exclude self-cited papers" and "Exclude patents".

ResearchGate search results for "data mining". The results page shows various academic profiles and publications. A specific profile is highlighted: "Jiawei Han" (University of Illinois at Urbana-Champaign). The page includes a sidebar with research interests like "Data Mining", "Machine Learning", and "Knowledge Discovery".

# AI Powered Research Career

133,208,365  
RESEARCHERS

272,411,089  
PUBLICATIONS

8,796,476  
CONCEPTS

754,201,878  
CITATIONS

An **author-centric** academic search  
and mining system...



COVID-19 Graph

Open Datasets

Knowledge Dashboard

Knowledge Graph



Experts



Academic Rankings



THU AI TR



Topic Trends



Open



Open Data



TOPIC Must Reading Papers  
COVID-19, Knowledge Graph



AI 2000  
AI 2000 Most Influential  
Scholars



Master Reading Tree  
Tools to help scholars study the  
evolution of papers



Xiaomai's Star Talents  
Predict Stars of Tomorrow with  
AI

Experts

A.M. Turing Award Winners

Chinese Academy of Sciences Talents

Academic Rankings

Want access to check cited publications

Women in AI new

THU AI TR

THU AI TR: 人工智能之机器学习 Hot

ICLR 2020反事实因果理论如何帮助深度学习? new

# Expert Search

Finding experts,  
for “data mining”

with more semantic  
constraints, e.g., from  
USA, speak Chinese,  
gender, and age...

data mining

Home |

Expert Paper

Search Results: 17

H-index: 17

Gender: Male (8) Female (4)

Language: Chinese (298) English (294) Greek (21) German (87) French (22) Japanese (20) Korean (14) Italian (13)

Location: USA (218) China (141) Taiwan (36) Australia (36) Canada (26) Japan (36) Germany (36) Italy (20) Hong Kong (26) Singapore (26)

Relevance (17) H-index Activity Diversity Rating Star Publication #Paper

Jiawei Han (汉江伟) ⓘ Follow  
H-Index: 125 | #Paper: 790 | #Citation: 9049  
Department of Computer Sciences, University of Illinois at Urbana-Champaign  
Professor  
Data Mining, Information Retrieval, Data Mining  
Similar Authors ⓘ 2567 views

Philip S. Yu ⓘ Follow  
H-Index: 134 | #Paper: 338 | #Citation: 6648  
Department of Computer Sciences, University of Illinois Chicago  
Professor and Walter Chair in Information Technology  
Distributed Systems, Query Optimization, Query Processing, Database Systems  
Data Mining  
Similar Authors ⓘ 323 views

Hiloi Kargupta ⓘ Follow  
H-Index: 40 | #Paper: 141 | #Citation: 6182  
Department of Computer Science and Electrical Engineering, University of Maryland Baltimore County  
Associate Professor  
Data Mining, Machine Learning, Data Analysis, Knowledge Discovery  
Genetic Algorithms  
Similar Authors ⓘ 338 views

Xindong Wu ⓘ Follow  
H-Index: 45 | #Paper: 231 | #Citation: 9544  
Department of Computer Sciences, University of Vermont  
Professor  
Machine Learning, Information Extraction, Bayesian Networks, Data Mining  
Supervised Learning  
Similar Authors ⓘ 25 views

Demographics: gender, language, location, etc.

Similar Authors

Data mining

Diagram illustrating the process of "Knowledge Discovery in Databases" (KDD):  
Raw Data → Data Cleaning → Data Integration → Data Selection → Data Transformation → Data Mining → Data Interpretation → Knowledge Representation.

Knowledge about “data mining”

Data mining (the analysis step of the “Knowledge Discovery in Databases” process, or KDD), a relatively young and interdisciplinary field of computer science, is the process that attempts to discover patterns in large data sets. It utilizes methods at the intersection of artificial intelligence, machine learning, statistics, and database systems. The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. Aside from the raw analysis step, it involves database and data management aspects, data preprocessing, model and inference considerations, interestingness metrics, complexity considerations, post-processing of discovered structures, visualization, and online updating.

Super Concepts:  
Data analysis, Data mining, Formal sciences, Applied sciences, Networks, Artificial intelligence

Related Concepts:  
Data compression, Data visualization, Natural language processing, Data cleansing, Distributed computing, Information retrieval, Speech recognition, Business intelligence, Pattern recognition, Spatial databases, Full-text search, Metadata, Computer vision, ISAM, Biological neural network, Database, Grid computing, Database marketing, Parallel computing

数据挖掘

数据挖掘(Data Mining)是通过分析每个数据，从大量数据中寻找其规律的技术。主要有数据准备、特征寻找和规则表达3个步骤。数据挖掘的任务有关联分析、聚类分析、分类分析、异常检测分析和集成分析等。

上位词:  
数据挖掘 - 计算机科学基础理论 - 决策支持系统 - 信息管理术语 - 数据挖掘 - 演示科学

# Researcher Profile

Whatever comes to your mind

Home | Log In



Jiawei Han (韩家炜) [Follow](#)

Department of Computer Science, University of Illinois at Urbana-Champaign  
Professor  
(217) 333-6903  
hanj@cs.uiuc.edu  
<http://www.cs.uiuc.edu/~hanj/>  
External Links [\[Edit\]](#) [\[Update\]](#)

**Basic Info.**

**Research Interests**

Data Mining, Information Extraction, Machine Learning, Text Mining, Data Analysis

**Ego Network**

**Citation Statistics**

Bringing Structure to Text: Mining Phrases, Entity Concepts, Topics, and Hierarchies  
Mining Massive RFID, Trajectory, and Traffic Data Sets  
20th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD), New York 2014  
11th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), Las Vegas 2008

About

Papers 170  
Lectures 11  
Patents 1

Entity Linking with a Knowledge Base  
Hai Shen, Jinyong Wang, Jiawei Han  
Knowledge and Data Engineering, IEEE Transactions on, Volume: 27, Issue: 1, January 2015, pp. 16-30  
<https://doi.org/10.1109/TKDE.2014.2350963>

Patterns  
Jing Yan, Jiawei Han  
Mining (2016)  
[https://doi.org/10.1007/978-3-319-27070-5\\_10](https://doi.org/10.1007/978-3-319-27070-5_10)

Troubleshooting Interactive components  
Mohammad Madi Hassan Khan, Hieu Khanh Le, Hossain Ahmed, Tarek F. Abdelsaleh, Jiawei Han  
ACM Transactions on Sensor Networks (TOSN) (2016)  
<https://doi.org/10.1145/2850290>

500+ connections  
Cited by 4  
Cited by 20  
Cited by 5

Search for people, jobs, companies, and more... [\[Search\]](#)

Home Profile My Network Jobs Interests Business Year By Citation

Jiawei Han  
Professor at UIUC  
Urbana-Champaign, Illinois, USA | Computer Software  
Education: University of Wisconsin-Madison  
[Send a message](#) [500+ connections](#)  
[https://www.linkedin.com/in/jiawehan](#) [See Contact Info](#)

Background

Experience

Professor  
UIUC  
August 2001 – Present (14 years 7 months)

Similar Authors | Ego Network

Watch +

Watch +

Bringing Structure to Text: Mining Phrases, Entity Concepts, Topics, and Hierarchies

20th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD), New York 2014

Mining Massive RFID, Trajectory, and Traffic Data Sets

11th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), Las Vegas 2008

Citation Statistics

1

Systems and Methods for Detecting a Novel Data Class

Mohammad Mehedy Masud, Latifur Rahman Khan, Bhavani Maronne Thuraisingham, Qing Chen, Jing Gao, Jiawei Han

Publication-date: 2012-03-01 Application-date: 2011-08-22

187

6

# AMiner

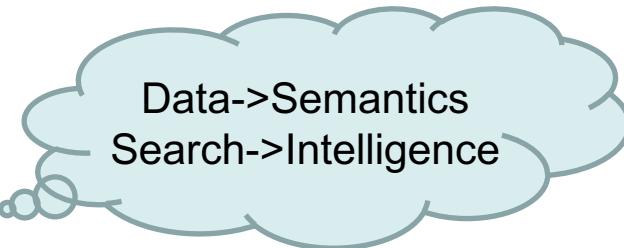
- Academic search and mining system—<https://aminer.cn>
  - Online since 2006
  - >136 million researcher profiles
  - >267 million publication papers
  - >341 million requests
  - 20M IP access from 220 countries per month
- Deep analysis, mining, and search

The image displays two screenshots of the AMiner platform. The top screenshot shows a detailed researcher profile for 'Jiawei Han'. It includes a photo, basic information like address and phone number, and sections for 'Publications' and 'Research Interests'. The bottom screenshot shows a search results page for 'Jiawei Han', displaying a list of related profiles and a network visualization.

# Technology Overview



Tim Berners Lee  
WWW Inventor  
Turing Award Winner

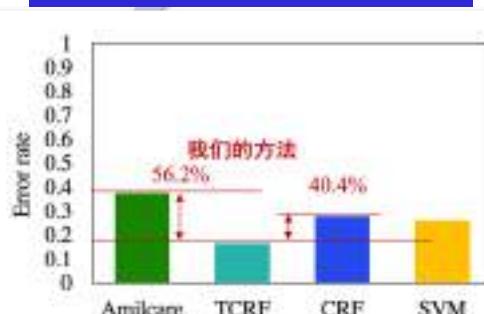


Intelligence

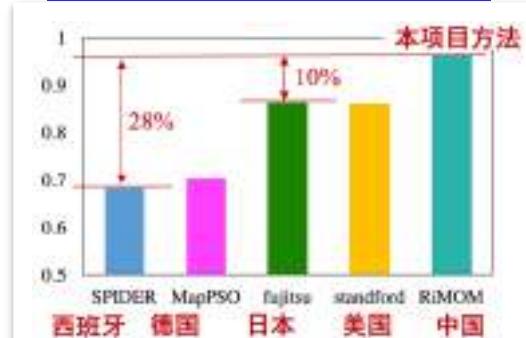
Linking

Extraction

Error rate reduced  
40-56%



15 champions in  
the past 7 years



Recommendation  
accuracy +161%



Reported by UN

UNITED NATIONS GLOBAL PULSE  
Harnessing big data for development and humanitarian action

PULSE LAB DIARIES  
Research Bites: "Inferring User Demographics and Social Strategies in Mobile Social Networks"  
August 2014

#citation>8,000, published ~20 papers on KDD

# Extracting Profile Semantics from the Web

(ACM TKDD, WWW'12, ISWC'06, ICDM'07, ACL'07)

**Ruud Bolle**

Office: 1S-D58

**Contact Information**

Office: 1S-D58  
Letters: IBM T.J. Watson Research Center  
P.O. Box 704  
Yorktown Heights, NY 10598 USA  
Packages: IBM T.J. Watson Research Center  
19 Skyline Drive  
Hawthorne, NY 10532 USA  
Email: bolle@us.ibm.com

**Ruud M. E.**  
Degree in Engineering, Netherlands, 1984 the P. Island. In Research Department Vision Group  
Currently, processing  
**Ruud M. E.**  
Vision and M. Bolle is

**DBLP: R**

2006		
50	EE	Nalini K. Ratha, Jonathan Connell, Ruud M. Bolle, Sharat Chikkerur: Cancelable Biometrics: A Case Study in Fingerprints. ICPR (4) 2006: 370-373
49	EE	Sharat Chikkerur, Sharath Pankanti, Alan Jea, Nalini K. Ratha, Ruud M. Bolle: Fingerprint Representation Using Localized Texture Features. ICPR (4) 2006: 521-524
48	EE	Andrew Senior, Arun Hampapur, Ying-li Tian, Lisa Brown, Sharath Pankanti, Ruud M. Bolle: Appearance models for occlusion handling. Image Vision Comput. 24(11): 1233-1243 (2006)
<b>2005</b>		
47	EE	Ruud M. Bolle, Jonathan H. Connell, Sharath Pankanti, Nalini K. Ratha, Andrew W. Senior: The Relation between the ROC Curve and the CMC. AutoID 2005: 15-20
46	EE	Sharat Chikkerur, Venu Govindaraju, Sharath Pankanti, Ruud M. Bolle, Nalini K. Ratha: Novel Approaches for Minutiae Verification in Fingerprint Images. WACV. 2005: 111-116

...

The diagram illustrates the extraction of semantic-based profile information from the Web. It shows how publications, research interests, and affiliations are linked to a profile.

**Publications:**

- Cancelable Biometrics: A Case Study in Fingerprints (ICPR 2006, pages 370-373)
- Fingerprint Representation Using Localized Texture Features (ICPR 2006, pages 521-524)
- Appearance models for occlusion handling (Image Vision Comput. 2006, pages 1233-1243)
- The Relation between the ROC Curve and the CMC (AutoID 2005, pages 15-20)
- Novel Approaches for Minutiae Verification in Fingerprint Images (WACV. 2005, pages 111-116)

**Research Interest:**

- video database indexing
- video processing
- visual human-computer interaction
- biometrics applications

**Affiliation:**

- IBM T.J. Watson Research Center
- Research Staff
- IBM T.J. Watson Research Center
- 19 Skyline Drive Hawthorne, NY 10532 USA
- IBM T.J. Watson Research

**Address:**

- IBM T.J. Watson Research Center  
19 Skyline Drive  
Hawthorne, NY 10532 USA

**Profile:**

- Ruud Bolle
- Office: 1S-D58
- Research Interest: video database indexing, video processing, visual human-computer interaction, biometrics applications
- Affiliation: IBM T.J. Watson Research Center
- Address: IBM T.J. Watson Research Center, 19 Skyline Drive, Hawthorne, NY 10532 USA
- Position: Professor
- UIUC

**Two questions:**

- How to accurately extract the semantic-based profile information from the Web?
- How to deal with the name ambiguity problem?

# A Big Data Approach

- Search Engine as Data Source

The image shows two Google search results side-by-side. The left search is for "philip s. yu email" and the right search is for "philip s. yu affiliation". Both searches yield approximately 951,000 and 616,000 results respectively. The results are identical, showing links to Philip S. Yu's profile at UIC, his Google Scholar page, Wikipedia entry, dbpedia page, and several of his publications.

**philip s. yu email**

About 951,000 results (1.05 seconds)

**Philip S. Yu - UIC - Computer Science - University of Illinois ...**  
https://www.cs.uic.edu/~PSYu/ • University of Illinois at Chicago  
May 20, 2009 · Philip S. Yu's main research interests include data mining especially ... Dr. Yu is a Fellow of the ACM and the IEEE. ... e-mail: psyu@cs.uic.edu.  
Lab - Research - Teaching - Awards

**Philip S. Yu - Google Scholar Citations**  
scholar.google.com/citations?user=DBL1rQAAAAJ • Google Scholar  
Professor of Computer Science, University of Illinois at Chicago - cs.uic.edu  
Philip S. Yu, Professor of Computer ... Verified email at cs.uic.edu Homepage ...  
Scholar ... CC Aggarwal, JL Wolf, PS Yu, C Procopius, JS Park, ACM SIGKDD ...

**Philip S. Yu - Wikipedia, the free encyclopedia**  
https://en.wikipedia.org/wiki/Philip\_S.\_Yu • Wikipedia  
Philip S. Yu (born ca 1962) is an American computer scientist and Professor in Information Technology at the University of Illinois at Chicago, known for his work ... Missing: email

**dblp: Philip S. Yu**  
dblp.uni-trier.de • Home > Persons • University of Trier  
List of computer science publications by Philip S. Yu.  
Missing: email

**Jiawei Han**  
han@illinois.edu •  
E-mail: han@cs.uiuc.edu Ph.D. ... Philip S. Yu, Jiawei Han, and Christos Faloutsos (eds.), *Link Mining: Models, Algorithms, and Applications*, Springer, 2010. You've visited this page 2 times. Last visit: 6/14/15

**Philip Yu | Computer Science and Computer Engineering ...**  
computer-science-and-computer-engineering.uark.edu/yu.php •  
Philip Yu distinguished speaker: ... Dr. Philip S. Yu, University of Illinois, Chicago Friday, November 6, 2015 3:05pm - 3:55pm, JBLH 144. Abstract: Philip Yu,

**Philip S. Yu - Journals, Conferences, Proceedings, Open ...**  
www.acip.org/JournalDetailedInfoOffEditorCont.aspx?personID... •  
Philip S. Yu, Computer Science Department, University of Illinois at Chicago, USA. Email: psyu@uic.edu Qualifications: 1978 Ph.D., New York University, USA ...

**philip s. yu affiliation**

About 616,000 results (0.72 seconds)

**dblp: Philip S. Yu**  
dblp.uni-trier.de • Home > Persons • University of Trier  
List of computer science publications by Philip S. Yu ... Person information ... affiliation: University of Illinois at Chicago [+] [-] ... Charu C. Aggarwal, Philip S. Yu

**Philip S. Yu - Google Scholar Citations**  
scholar.google.com/citations?user=DBL1rQAAAAJ • Google Scholar  
Professor of Computer Science, University of Illinois at Chicago - cs.uic.edu  
Mining concept-drifting data streams using ensemble classifiers. H Wang, W Fan, PS Yu, J Han. Proceedings of the ninth ACM SIGKDD International conference ...

**Philip S. Yu - Wikipedia, the free encyclopedia**  
https://en.wikipedia.org/wiki/Philip\_S.\_Yu • Wikipedia  
Philip S. Yu (born ca 1962) is an American computer scientist and Professor in Information Technology at the University of Illinois at Chicago, known for his work ...

**Link Mining: Models, Algorithms, and Applications**  
https://books.google.com/books?isbn=1441885157  
Philip S. Yu, Jiawei Han, Christos Faloutsos - 2010 - Science  
Philip S. Yu, Jiawei Han, Christos Faloutsos ... Affiliation is a transitive relationship; therefore, all individuals sharing an affiliation form a clique. 2 An affiliation ...

**Discovering High-Order Periodic Patterns - Springer**  
link.springer.com/10.1007%2Fst101... • Springer Science+Business Media ... by J Yang - 2004 - Cited by 20 - Related articles  
Authors: Jong Yang - jongyang@cs.umn.edu (1), Wei Wang (2), Philip S. Yu (3).  
Author Affiliations: (1) Computer Science Department, UIUC, Urbana, IL, USA; (2)

**On clustering massive text and categorical data streams ...**  
link.springer.com/10.1007%2Fst101... • Springer Science+Business Media ... by CC Aggarwal - 2010 - Cited by 47 - Related articles  
Aug 6, 2009 · Authors: Charu C. Aggarwal - charu@us.ibm.com (1), Philip S. Yu (2). Author Affiliations: (1) IBM T. J. Watson Research Center, 19 Skyline Drive,

**An Introduction to Privacy-Preserving Data Mining - Springer**  
link.springer.com/10.1007%2F978... • Springer Science+Business Media ... by CC Aggarwal - 2008 - Cited by 7 - Related articles  
(1) Philip S. Yu - psyu@cs.uic.edu (4). Editor Affiliations: (1) IBM Thomas J. Watson Research Center; (4) Department of Computer Science, University of Illinois at ...

# MagicFG: Markov Logic Factor Graph

- ✓ Depict and utilize correlations between possible candidates from redundant data.
- ✓ Incorporate human knowledge to guide and amend the classification model.

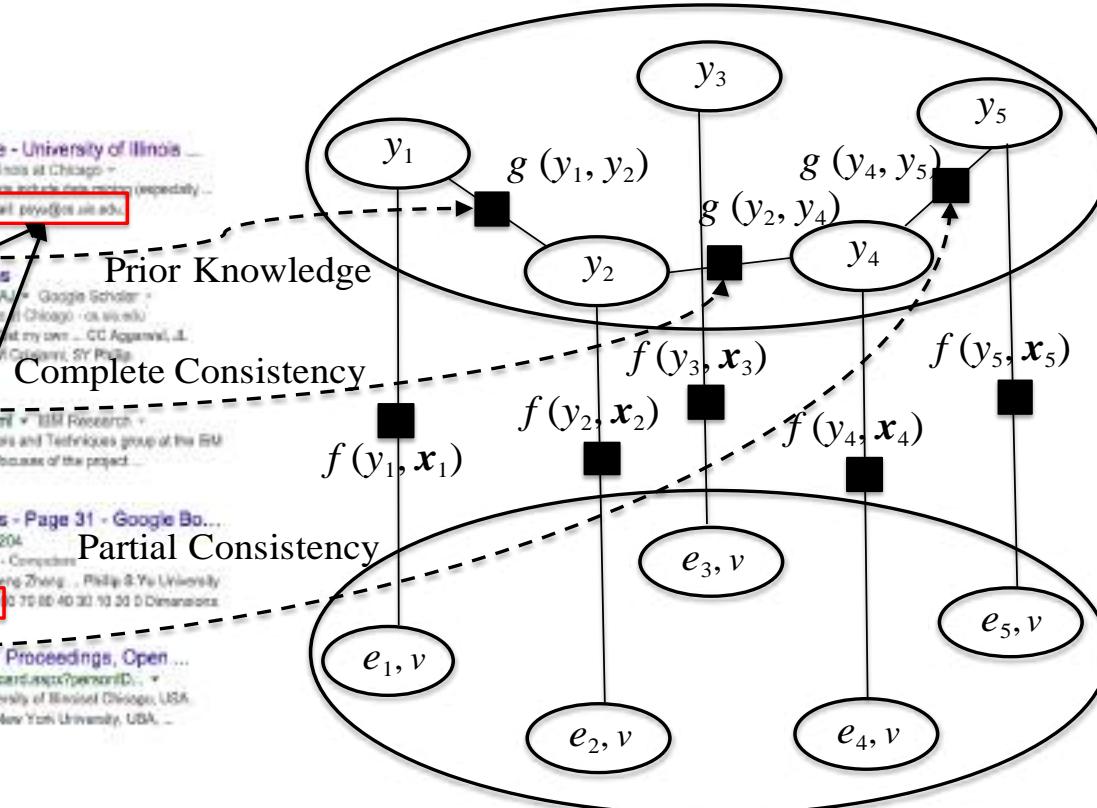
Philip S. Yu - UIC - Computer Science - University of Illinois  
<https://www.cs.uic.edu/~BYU/> - University of Illinois at Chicago -  
May 20, 2018 - Philip S. Yu's main research interests include data mining (especially  
Dr. Yu is a Fellow of the ACM and the IEEE. E-mail: [psyu@cs.uic.edu](mailto:psyu@cs.uic.edu)  
Labs - Research - Teaching - Awards

Philip S. Yu - Google Scholar Citations  
[scholar.google.com/citations?user=00L1HQAQAAQ](https://scholar.google.com/citations?user=00L1HQAQAAQ) - Google Scholar -  
Professor of Computer Science, University of Illinois at Chicago - cs.uic.edu  
Verified email at cs.uic.edu - Homepage - Scholar - End my own - CG Aggarwal, J.  
Wu, P. Yu, C. Principe, S. Park ... V. Carvalho, M. Corrent, S.Y. Philip

IBM Research - Dr. Philip S. Yu  
<http://www.research.ibm.com/people/b/yulyu/> - IBM Research -  
Dr. Philip S. Yu is the manager of the Business Tools and Techniques group of the IBM Thomas J. Watson Research Center. This project focuses on the project:  
Missing email

Data Mining for Business Applications - Page 31 - Google Bo...  
<https://books.google.com/books?id=e0387794204>  
Longbing Cao, Philip S. Yu, Chengqi Zhang - 2008 - Computers  
Longbing Cao, Philip S. Yu, Chengqi Zhang, Huafeng Zheng, ... Philip S. Yu, University  
of Illinois at Chicago. E-mail: [psyu@cs.uic.edu](mailto:psyu@cs.uic.edu) 50 0 70 80 40 30 10 30 0 Dimensions  
(1 X 30)

Philip S. Yu - Journals, Conferences, Proceedings, Open...  
[www.acnp.org/journals/DelectronOffshoreBoard.aspx?personID=...](http://www.acnp.org/journals/DelectronOffshoreBoard.aspx?personID=...)  
Philip S. Yu, Computer Science Department, University of Illinois Chicago, USA.  
Email: [psyu@uic.edu](mailto:psyu@uic.edu) Publications: 1978 Ph.D., New York University, USA, ...



# Is this Enough?

Whatever comes to your mind.

jeannette Wing

Computer Science Department Carnegie Mellon University  
Professor of Computer Science  
412-268-6926 (cell)  
wing@cs.cmu.edu  
<http://www.cs.cmu.edu/~wing/>

External Links

Update

Research Interests

- Formal Methods
- Software Engineering
- Formal Verification
- Embedded Systems
- Specification Languages

1982 1985 1990 1995 2000 2005 2009

Similar Authors

Ego Network

Overview

Papers 24

Merge 24

Computational thinking

All (24) Recent 30 2009 2008 2007 2006 2005 2004 2003 2002 2001 2000 1999 1998 1997 1996 1995 1994 1993 1992 1991 1990 1989 1988 1987 1986 1985 1984 1983 1982

Add Paper Remove Paper By Year By Citation

Cited by 1377

# Linking Semantics across Networks

- Identifying users from multiple heterogeneous networks and integrating semantics from the different networks together.

LinkedIn



A screenshot of Jeannette Wing's LinkedIn profile. It shows her title as Corporate Vice President at Microsoft, her location as Redmond, Washington, and her education from Carnegie Mellon University. Her LinkedIn ID is 674.

WikiPedia



A screenshot of Jeannette Wing's Wikipedia page. It includes a photo, her title as Corporate Vice President, and her education from Carnegie Mellon University.

Same Person



A screenshot of Jeannette Wing's Google Scholar profile. It shows her publications, including "Essential methods in machine learning" and "Solving multi-label learning problems via efficient covering codes".

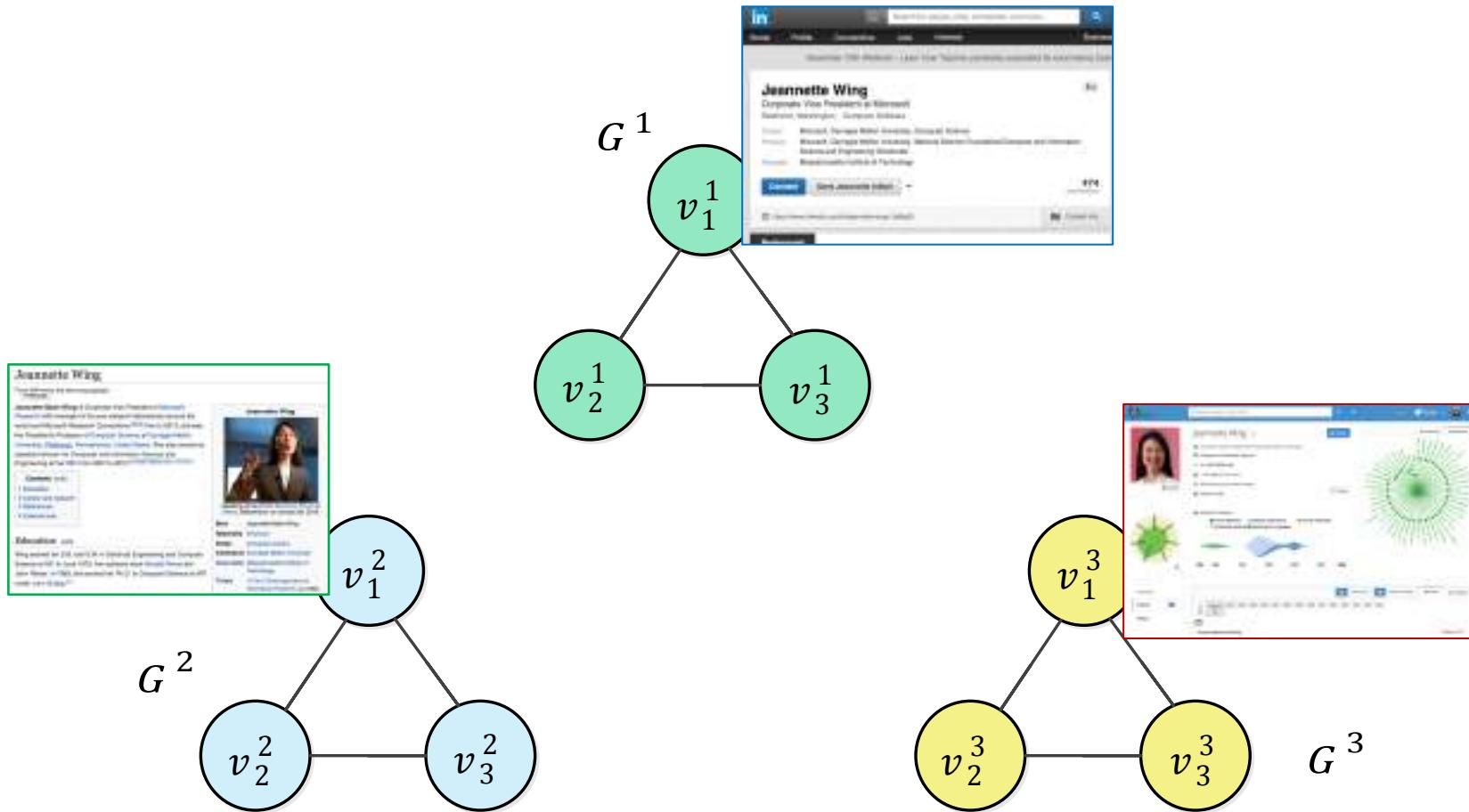
Google Scholar



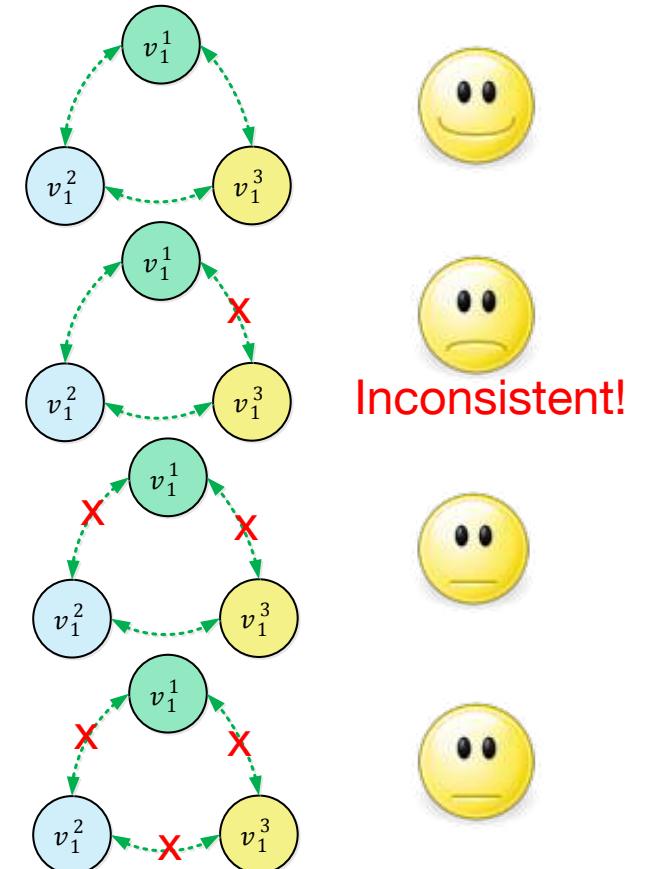
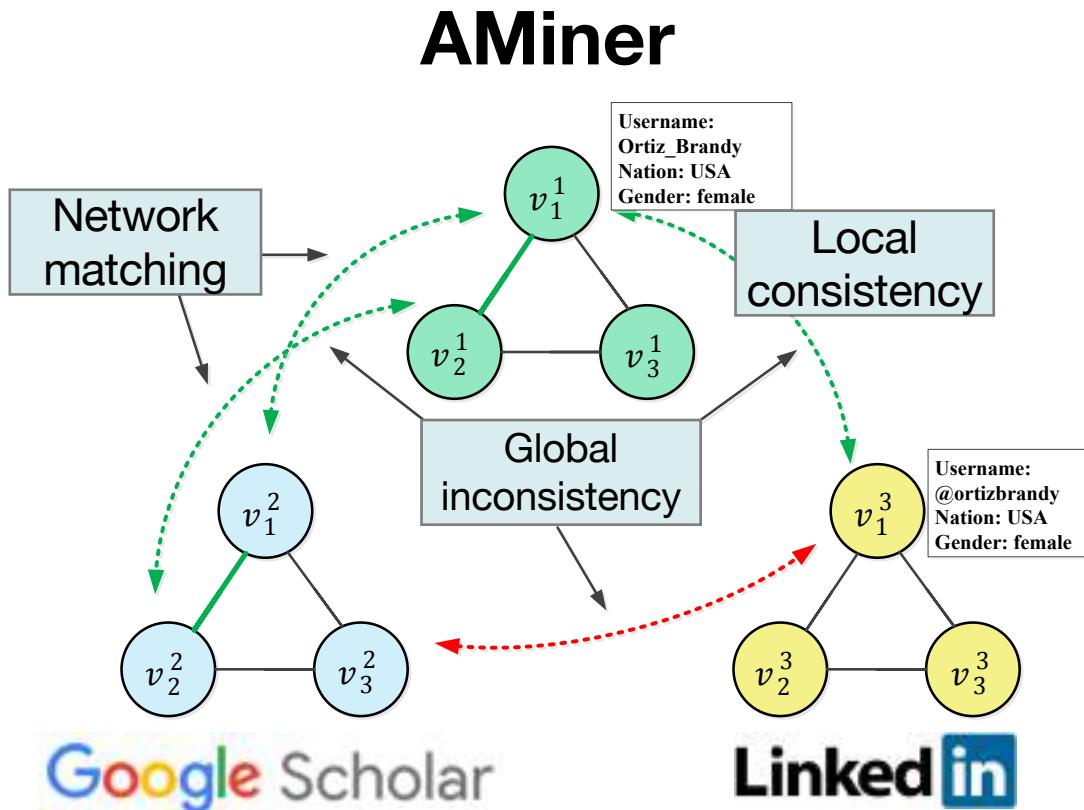
A screenshot of Jeannette Wing's AMiner profile. It features a circular network visualization showing her research connections and a timeline of her publications.

AMiner

# Considering the networks...



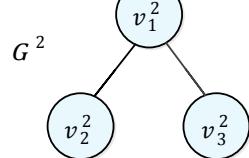
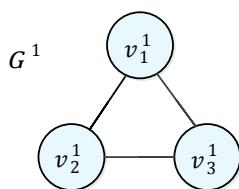
# Local + Global consistency



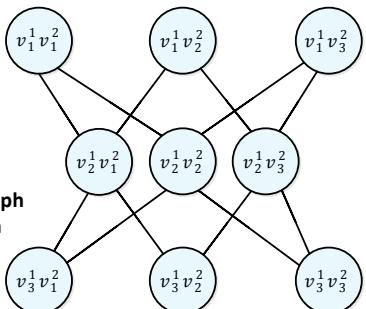
# COSNET: Connecting Social Networks with Local and Global Consistency

- **Input:**  $\mathbf{G}=\{G^1, G^2, \dots, G^m\}$ , with  $G^k=(V^k, E^k, R^k)$
- **Formalization:**  $\mathbf{X}=\{x_i\}$ , all possible pairwise matchings and each corresponds to  $y_i \in \{1, 0\}$
- **COSNET:** an energy-based model
$$Y^* = \arg \max E(Y, X)$$

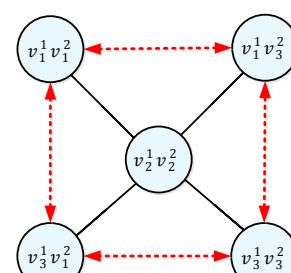
# Model Construction



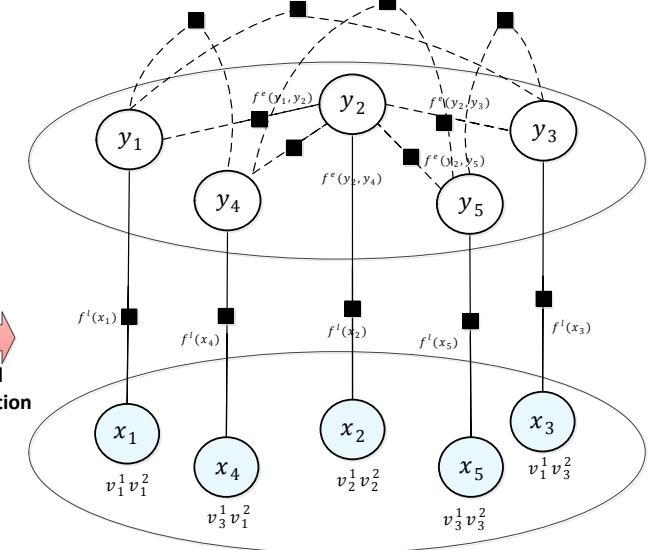
Matching Graph Generation



Candidate Pruning



Model Construction



(a) Two input networks

(b) The generated matching graph

(c) Matching graph after pruning

(d) The constructed model

Objective function by combining all the energy functions

$$\begin{aligned}
 E(Y, X) &= \sum_{\mathbf{x}_i \in V_{MG}} \mathbf{w}_i^\top \mathbf{g}_i(\mathbf{x}_i, y_i) + \sum_{\langle \mathbf{x}_i, \mathbf{x}_j \rangle \in E_{MG}} \mathbf{w}_e^\top \mathbf{f}_e(y_i, y_j) \\
 &\quad + \sum_{c \in T_{MG}} \mathbf{w}_t^\top \mathbf{f}_t(Y_c)
 \end{aligned} \tag{2}$$

# AMiner

Whatever comes to your mind

Jiawei Han (韩家炜)

Department of Computer Science, University of Illinois at Urbana-Champaign, Illinois, USA | Computer Software

Professor

(217) 333-6903

han@cs.uiuc.edu

<http://www.cs.uiuc.edu/~han/>

External Links

Google Scholar LinkedIn YouTube

Research Interests

Data Mining, Information Extraction, Machine Learning, Text Mining

Background

Experience

Professor

UIUC  
August 2001 – Present (14 years 7 months)

Add Paper Remove Paper By Year By Citation

Entity Linking with a Knowledge Graph Watch

Bringing Structure to Text: Mining Phrases, Entity Concepts, Topics, and Hierarchies Watch

20th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD), New York 2014 Watch

Troubleshooting interactive complexity bugs in wireless sensor networks using data mining techniques Watch

Mohammad Mehedy Masud, Latifur Rahman Khan, Bhavani Maronne Thuraisingham, Qing Chen, Jing Gao, Jiawei Han

Publication date: 2012-03-01 Application date: 2011-08-22

18

# Discovering Expertise Semantics

- Quantifying researchers' expertise using publications

818 Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions  
Wei Shan, Jinyong Wang, **Jiawei Han**  
Knowledge and Data Engineering, IEEE Transactions (2015)  
# Biblio: 0 <http://dx.doi.org/10.1109/TKDE.2014.2327028>

817 Power-Based Diagnosis of Node Silence in Remote High-End Sensing Systems.  
Ying Ying, Lu Su, Mohammad Nafeel Hassan Khan, Michael Lebtah, Tarek F. Abdelzaher, **Jiawei Han**  
TOEN (2015)  
# Biblio: 0 <http://doi.acm.org/10.1145/2661638>

816 A Framework of Mining Trajectories from Untrustworthy Data in Cyber-Physical Systems.  
Lu An Tang, Xiao Yu, Quenquan Gu, **Jiawei Han**, Guohui Jiang, Alice Leung, Thomas H. La Porta  
TKDD (2015)  
# Biblio: 0 <http://doi.acm.org/10.1145/2700394>

815 A Unifying Framework of Mining Trajectory Patterns of Various Temporal Tightness.  
Jae-Gil Lee, **Jiawei Han**, Xiaolei Li  
IEEE Trans. Knowl. Data Eng. (2015)  
# Biblio: 0 <http://dx.doi.org/10.1109/TKDE.2014.2377742>

814 ePeriodicity: Mining Event Periodicity from Incomplete Observations  
Zhenwei Li, Jingjing Wang, **Jiawei Han**  
Knowledge and Data Engineering, IEEE Transactions (2015)  
# Biblio: 0 <http://dx.doi.org/10.1109/TKDE.2014.2366601>



Jiawei Han (韩家炜)

H-Index: 133 | #Paper: 818 | #Citation: 111184

9 Department of Computer Science, University of Illinois at Urbana-Champaign  
Professor

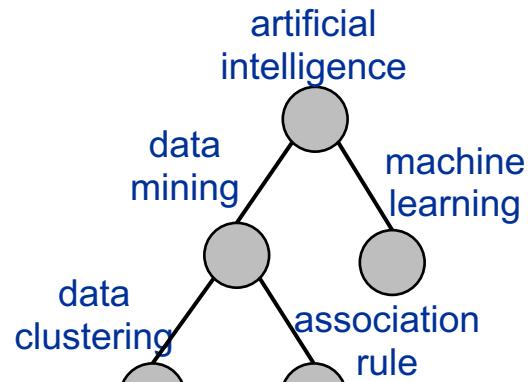
Data Mining Information Extraction Data Analysis Machine Learning Text Mining

The straightforward method is to extract high frequent terms as expertise; however, the extracted terms may be not good terms to represent expertise.

# Knowledge-driven Semantic Mining

The screenshot displays a list of four research papers from a database:

- 1. Efficient Linking with a Knowledge Base: Resources, Techniques, and Solutions. Author: Jiawei Han. DOI: 10.1109/TKDE.2010.2020728.
- 2. Power-Based Diagnosis of Node Silence in Remote High-End Sensing Systems. Authors: Yang Tang, Li Su, Mohammad Reza Hassan Ravan, Michael Ladday, Tarek F. Abdellatif. DOI: 10.1109/TNSM.2013.2273038.
- 3. A Framework of Mining Trajectories from Untrustworthy Data in Cyber-Physical Systems. Authors: Li An Tang, Ren Yu, Guozeng Bi, Jiawei Han, Daniel Jiang, Alvin Leong, Romeo F. Lai, Po-Yen. DOI: 10.1109/TNSM.2013.2273039.
- 4. A Unifying Framework of Mining Trajectory Patterns of Various Temporal Tightness. Authors: Jun-Ge Lee, Jiawei Han, Daode Li. DOI: 10.1109/TNSM.2013.2273040.



Knowledge graph



WIKIPEDIA  
The Free Encyclopedia

Freebase



Jiawei Han (叶家伟)

H-Index: 133 | #Paper: 818 | #Citation: 111184

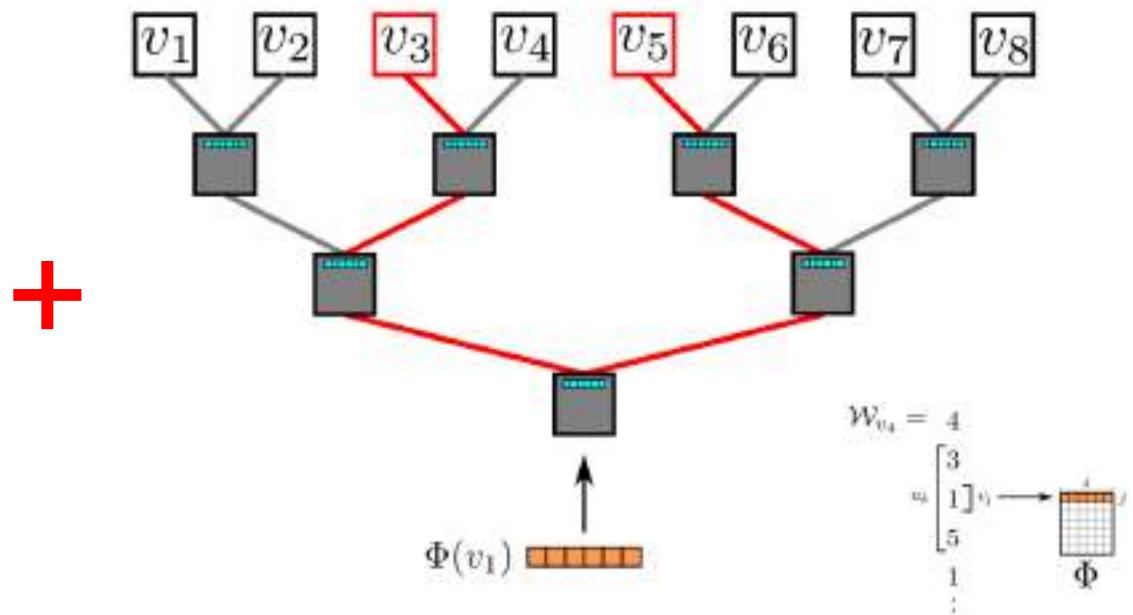
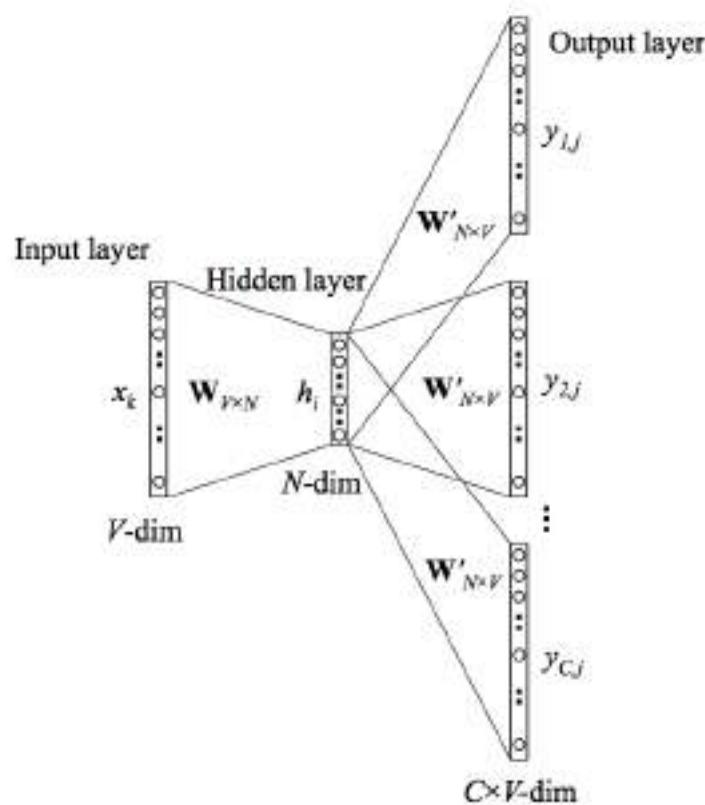
Department of Computer Science, University of Illinois at Urbana-Champaign  
Professor



Similar

Data Mining Information Extraction Data Analysis Machine Learning Text Mining

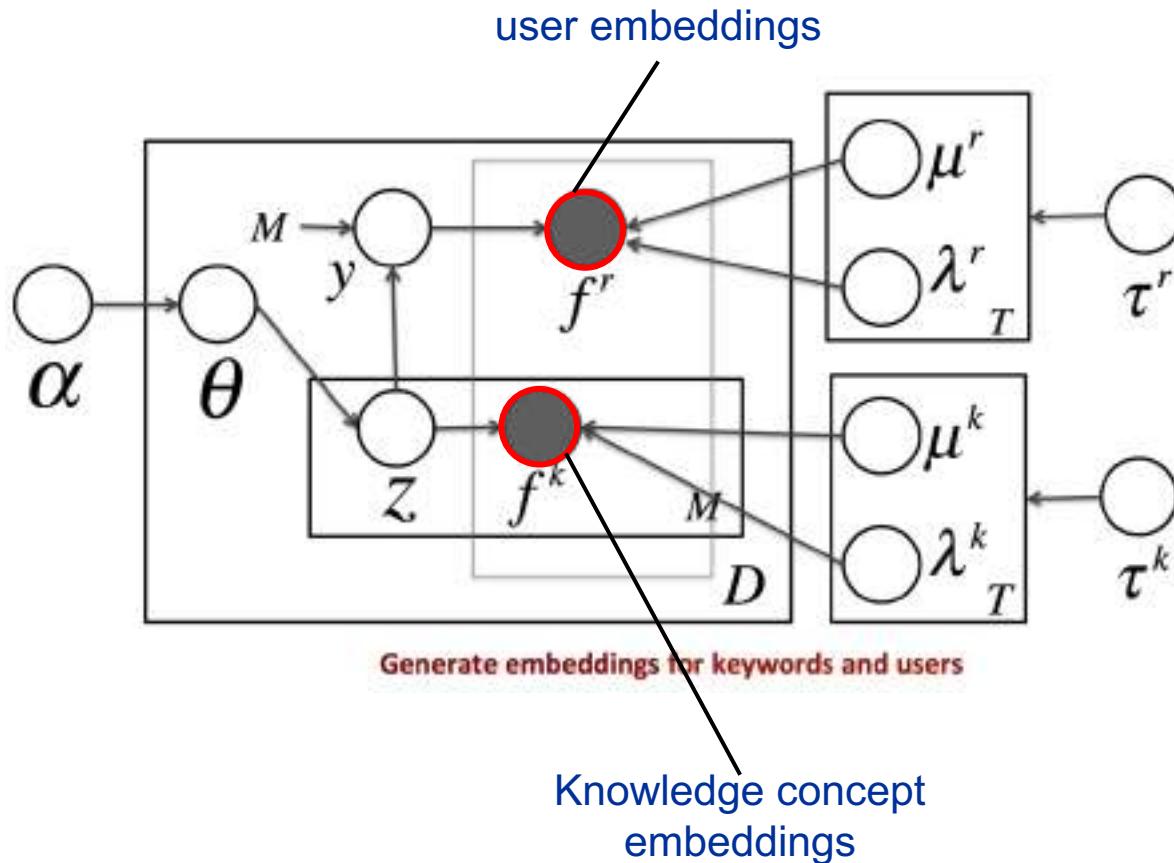
# Embedding Knowledge and Networks



Knowledge concept embedding

Network-based researcher embedding

# GenVector—bridging researcher network and knowledge graph



# Examples

## Case study: Dan Klein

GenVector	AM-base
Language models	Machine translation
Markov models	Word alignment
Probabilistic models	Bleu score
Natural language	Best result
CorefERENCE resolution	Language model

## Case study: Xiaou Tang

GenVector	AM-base
Feature extraction	Face recognition
Image segmentation	Face image
Image matching	Novel approach
Image classification	Line drawing
Face recognition	Discriminant analysis

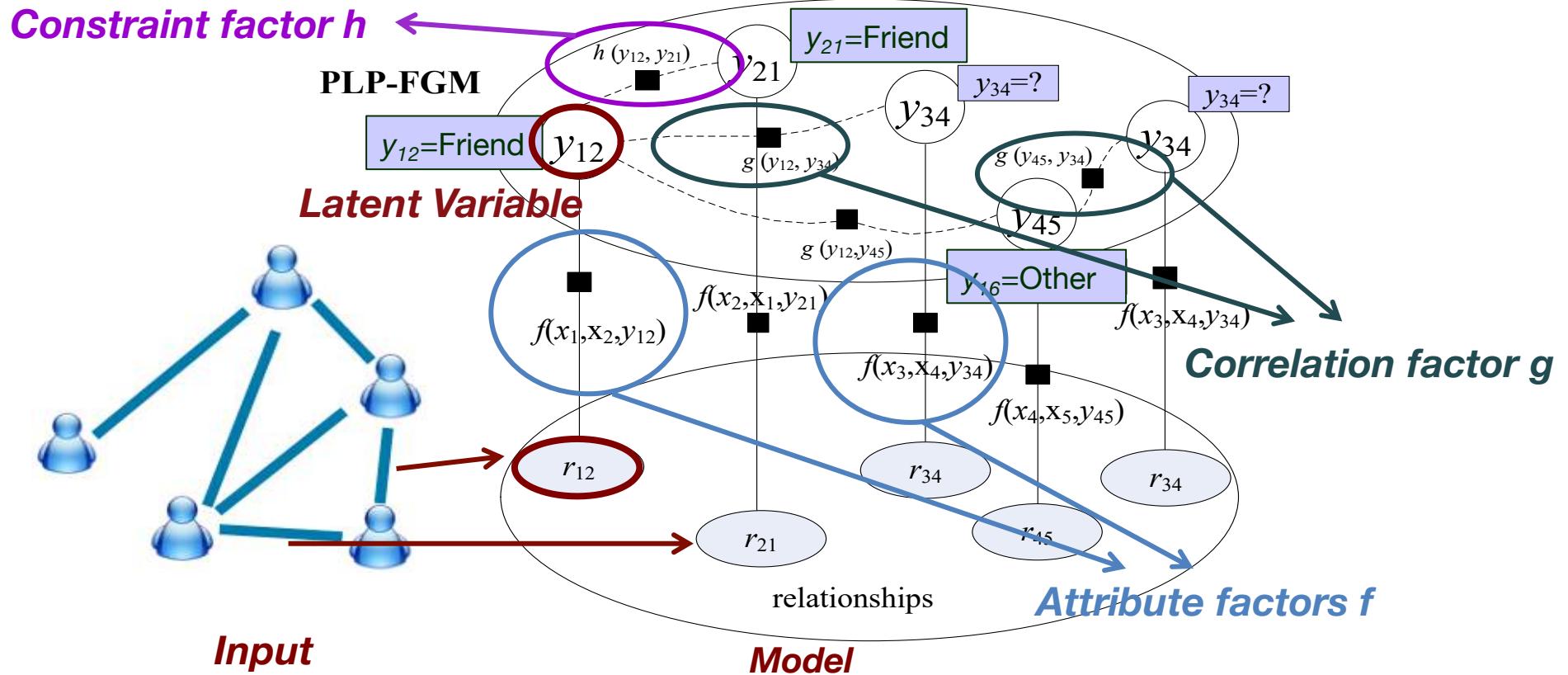
GenVector—the proposed Bayesian embedding method;  
AM-base—previous language modeling method.

# Discovering Network Semantics

- How to mine semantics within networks?
    - Who are Jiawei Han's students and advisor?



# Partially Labeled Pairwise Factor Graph Model (PLP-FGM)



Map relationship to nodes in model

Example

#Coaut

Example:

Author A has a longer publication history than both B and C.

# AMiner

Whatever comes to your mind

Jiawel Han (韩家炜)

Department of Computer Science, University of Illinois at Urbana-Champaign  
Professor  
1 (217) 333-6903  
hanj@cs.uiuc.edu  
<http://www.cs.uiuc.edu/~han/>

External Links:

Research Interests: Data Mining (blue), Machine Learning (orange), Information Extraction (light blue), Text Mining (green), Data Analysis (yellow)

Timeline: 1985, 1990, 1995, 2000, 2005, 2010, 2015



Ego Network

Papers	790
Lectures	11
Patents	1

Add Paper | Remove Paper | By Year | By Citation

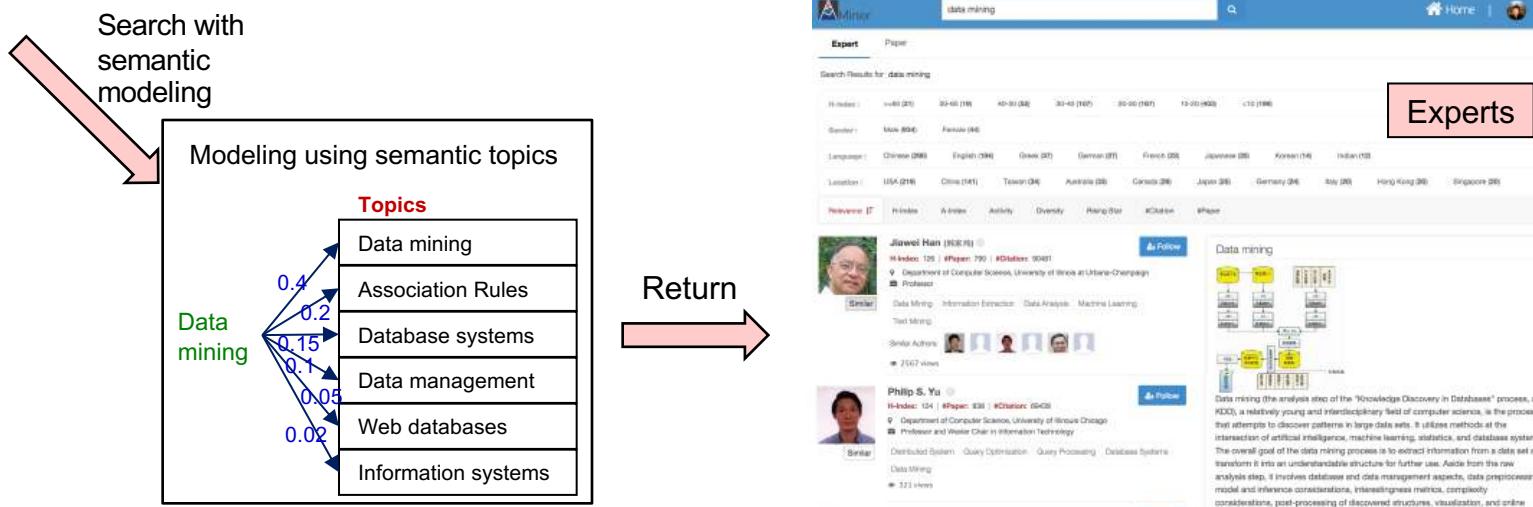
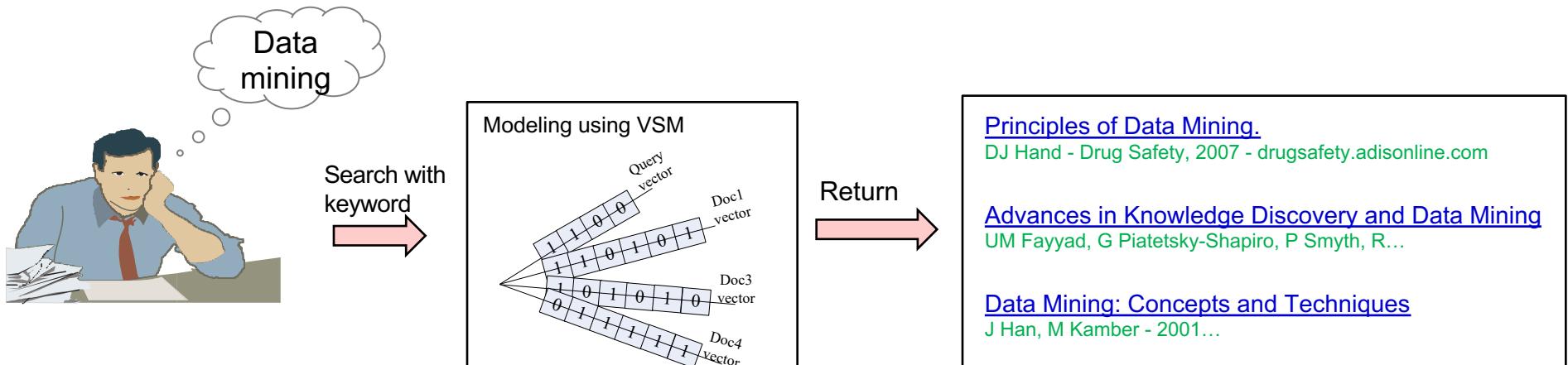
790 Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions Cited by 4  
Hai Shan, Jinyong Wang, **Jiawel Han**  
Knowledge and Data Engineering, IEEE Transactions (2016)  
[View](#) [DOI](#) <https://doi.org/10.1109/KDE.2014.2327528>

789 Mining Graph Patterns Cited by 20  
Hong Chang, Xileng Yan, **Jiawel Han**  
Frequent Pattern Mining (2016)  
[View](#) [DOI](#) [https://doi.org/10.1007/978-1-4614-6044-0\\_13](https://doi.org/10.1007/978-1-4614-6044-0_13)

788 Troubleshooting interactive complexity bugs in wireless sensor networks using data mining techniques Cited by 3  
Mohammad Madi Hassan Khan, Hieu Khan Lu, Hossain Ahmed, Tarek F. Abdessalem, **Jiawel Han**  
ACM Transactions on Sensor Networks (TOSN) (2016)  
[View](#) [DOI](#) <https://doi.org/10.1145/2550290>

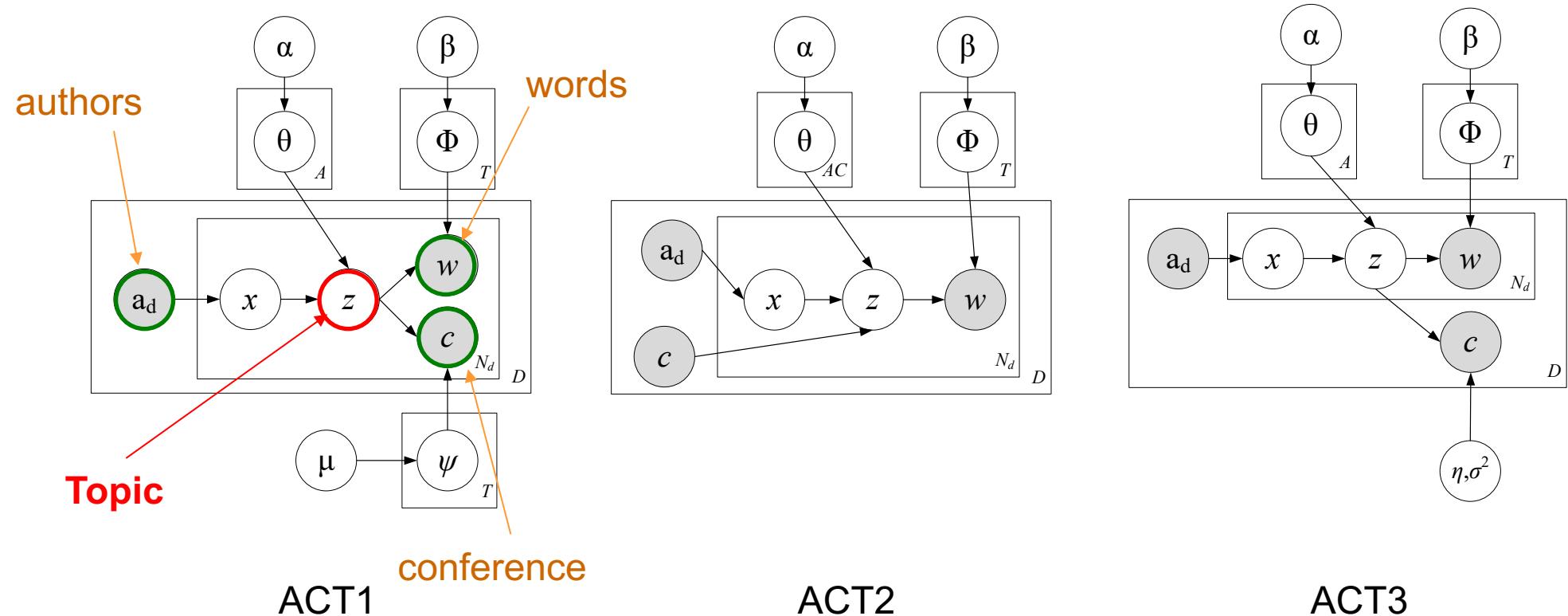
787

# Search in AMiner



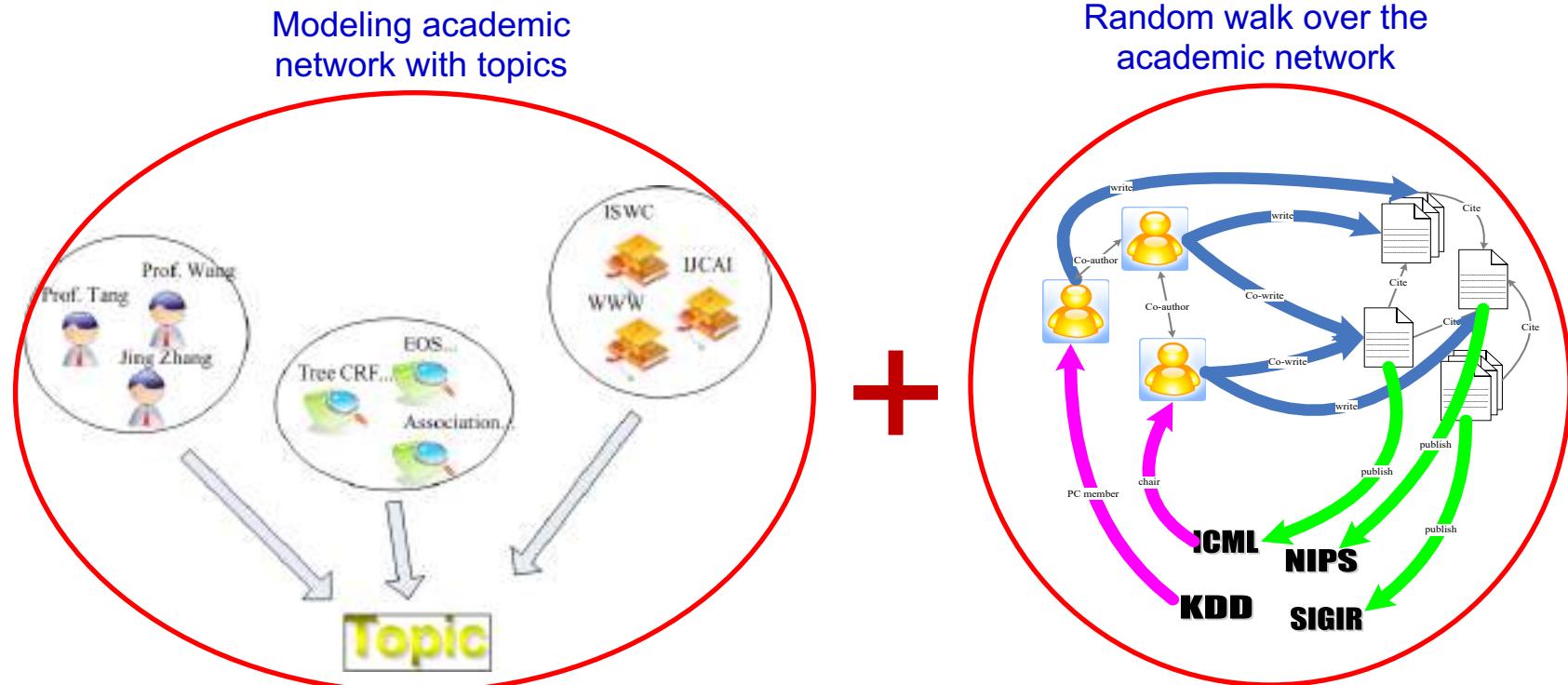
The screenshot shows the AMiner search interface with the query "data mining" entered. The results are categorized under "Expert" and "Paper". The "Expert" tab is selected, displaying profiles for three experts: Jiliwei Han, Philip S. Yu, and Hilito Kargupta. Each profile includes basic information like H-index, number of papers, and citations, along with a list of research interests. The "Paper" tab is also visible, showing a grid of paper thumbnails. A pink box highlights the "Experts" tab. At the bottom right, there is a sidebar with sections for "Super Concepts", "Related Concepts", and "Database".

# Modeling the Academic Network



Author-Conference-Topic Model [Tang et al., 08]

# Integrating Topic Model into Random Walk



Author-Conference-Topic  
Model [Tang et al., 08]

# Heterogeneous Cross-domain Ranking



$$\min_{w_S, w_T} \left\{ \sum_{i=1}^{n_1} \left[ 1 - z_{S_i} \langle w_S, \underbrace{x_{S_i}^a - x_{S_i}^b}_{\text{Loss in one domain}} \rangle \right]_+ + C \sum_{i=1}^{n_2} \left[ 1 - z_{T_i} \langle w_T, \underbrace{x_{T_i}^a - x_{T_i}^b}_{\text{Loss in another domain}} \rangle \right]_+ + \lambda \|W\|_{2,1} \right\}$$



$$\min_{w_S, w_T, U} \left\{ \sum_{i=1}^{n_1} \left[ 1 - z_{S_i} \langle w_S, \underbrace{U^T(x_{S_i}^a - x_{S_i}^b)}_{\text{Common feature space}} \rangle \right]_+ + C \sum_{i=1}^{n_2} \left[ 1 - z_{T_i} \langle w_T, \underbrace{U^T(x_{T_i}^a - x_{T_i}^b)}_{\text{Common feature space}} \rangle \right]_+ + \lambda \|W\|_{2,1} \right\}$$

# Results on Heterogeneous Tasks

- Expert finding verse Bole search (finding best supervisor)
- To obtain ground truth of bole for each query
  - We sent emails to 50 senior researchers and 50 junior researchers (91.6% are post doc or graduates)
  - Average their feedbacks

**Table Results on heterogeneous tasks.**

Approach	P@5	P@10	P@15	MAP	N@5	N@10
RSVM	.7714	<b>.8429</b>	.8285	.7756	.5545	.5947
RSVMt	.8000	.8286	.8476	.7837	.5923	.5999
MTRSVM	.8000	.8286	.8476	.7875	.6140	.6075
HCDRank	<b>.8285</b>	.7857	<b>.8571</b>	<b>.7971</b>	<b>.6189</b>	<b>.6112</b>
Language model	.6250	.6875	.6500	.6726	.3343	.3809

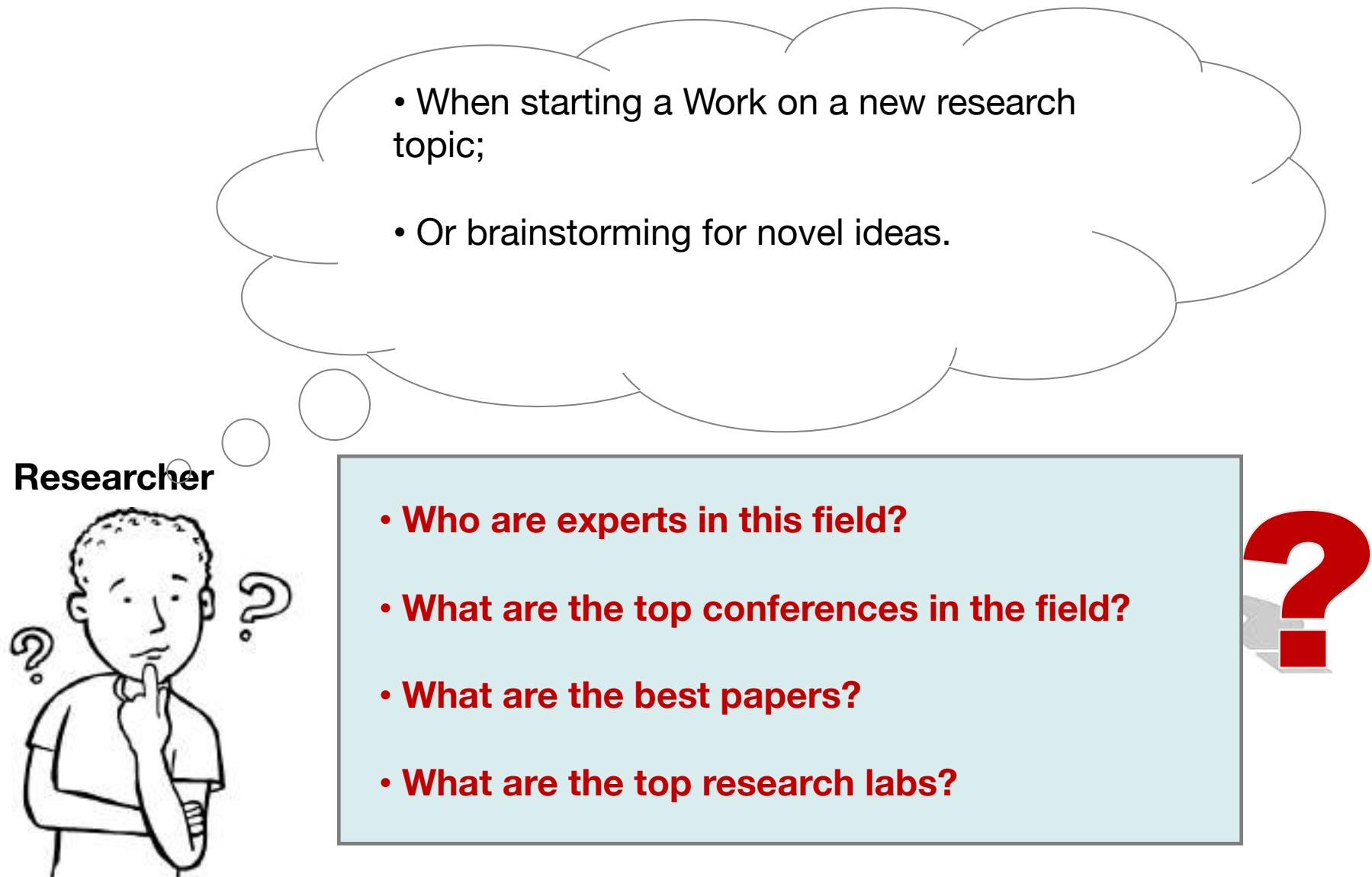


# Why you will love AMiner?

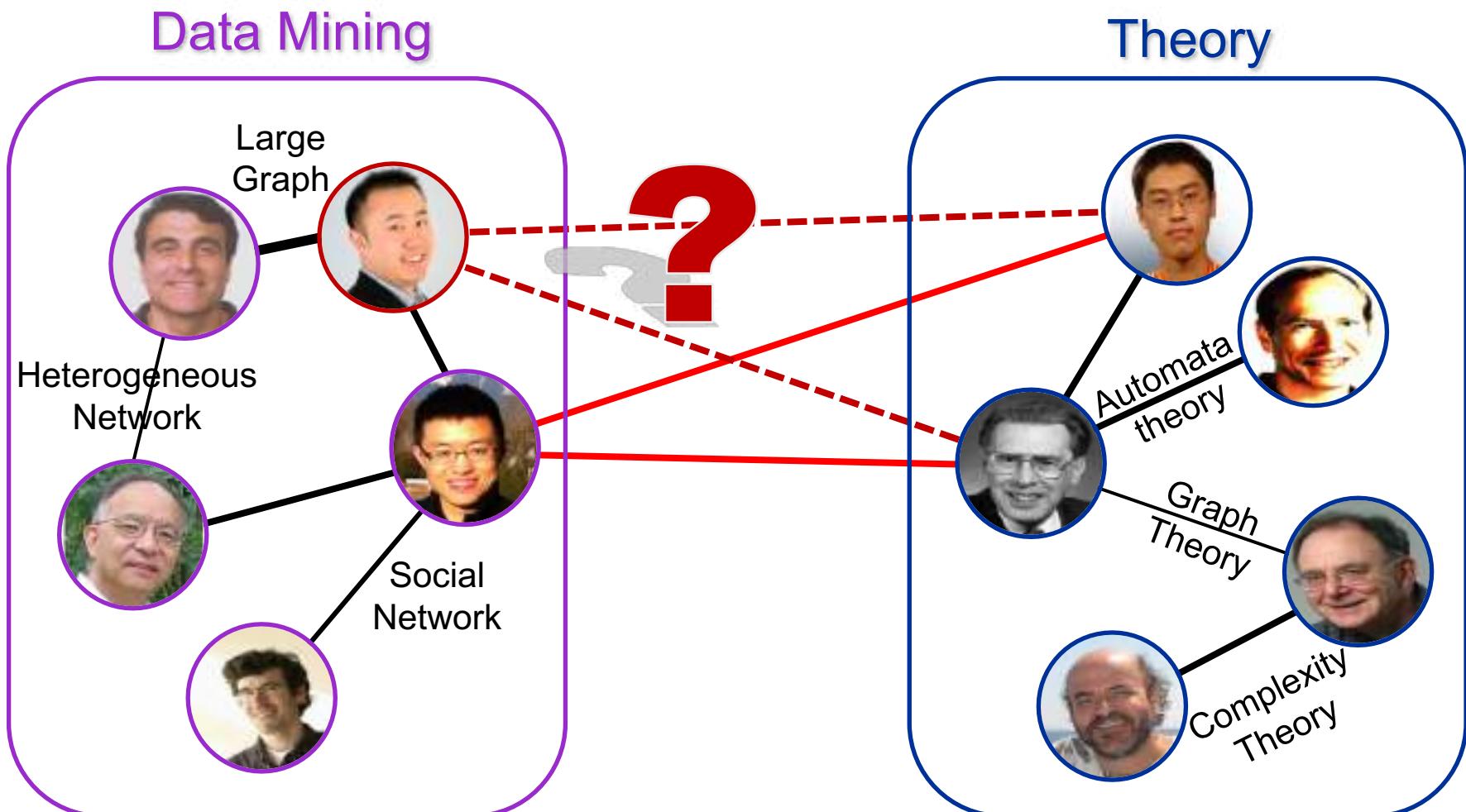
—finding collaborations, applying  
for PhD, proposal reviewing,  
recruiting...

Related publications: ACM TKDD (2), IEEE TKDE (3), KDD'08-15, WWW'12-15, J. Informetrics, SIGMOD'09, IJCAI'09-15, ICML'14

# Examples – Expert search



# Examples – Collaboration Recommendation



# Examples – More Challenging Questions...

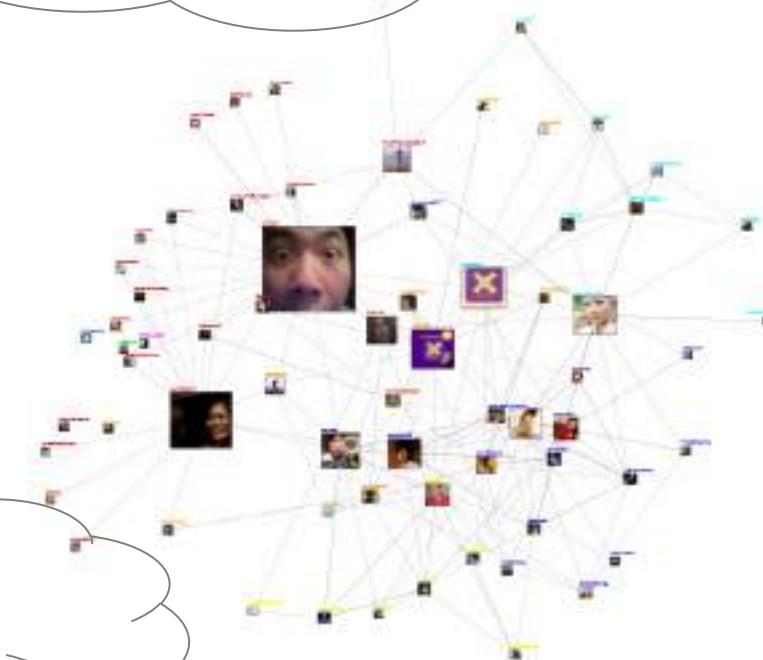
Organization  
or Agency



- Looking for researchers  
**from the U.S.** who can  
speak Chinese...



- Want to recruit **female**  
**researchers** working on  
WSDM



# Data Mining Experts

We use this tool to help  
**MOST (Ministry of Science and Technology of the People's Republic of China)** to identify experts in different research fields.

Global Data Mining Experts 

This collection includes experts in data mining field. The basic metric is that if a researcher served as SPC at KDD (the best conference in Data Mining) for more than once, she/he will be considered to be included in this collection.  
1. Jim Yang - 02/18-02-18

0 ~ 20 of 310 results | 1 | Go

Statistics Extract emails Export Members

Keywords	Name	Organization	Find All Persons					
H-Index :	>=60 (10)	50-59 (26)	40-49 (28)	30-39 (82)	20-29 (71)	10-19 (106)	<10 (86)	
Gender :	Male (302)	Female (85)						
Language :	Chinese (113)	English (88)	Indian (21)	Greek (16)	German (12)	Japanese (8)	French (8)	Korean (4)
Location :	USA (170)	China (26)	Hong Kong (14)	Germany (11)	Canada (11)	Australia (10)	Singapore (8)	Belgium (7)

Relevance  H-Index Activity Diversity Publishing #Citation #Paper AB



Jiawei Han [韩家炜]  
✉ Department of Computer Science, University of Illinois at Urbana-Champaign  
☞ Professor

 Tags Male Chinese



Philip S. Yu  
✉ Department of Computer Science, University of Illinois Chicago  
☞ Professor and Walter Chair in Information Technology

 Tags Male Chinese



Christos Faloutsos  
✉ Dept. of Computer Science Carnegie Mellon University  
☞ Professor

 Tags English



Vipin Kumar  
✉ University of Minnesota  
☞ Professor

 Tags Indian



Rajeev Motwani  
✉ Stanford University  
☞ Professor and Director of Graduate Studies

 Tags English



George Karypis



Jian Pei



Gerhard Weikum



Charu C. Aggarwal

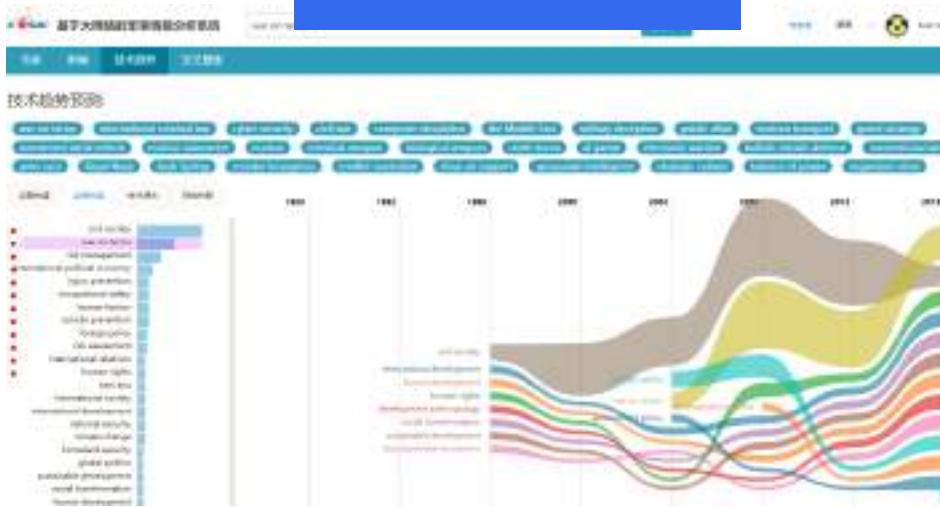


Zhi-Hua Zhou [周志华]

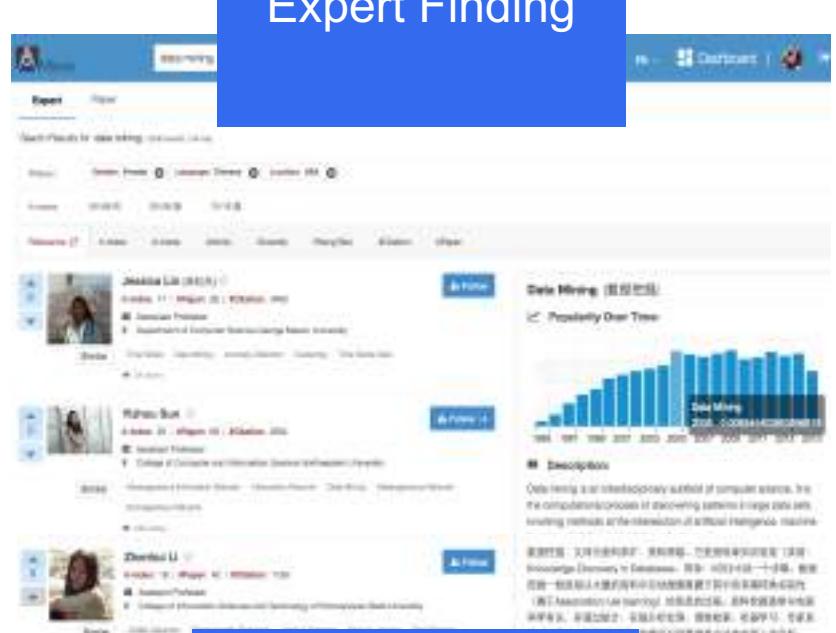
<https://aminer.org/datamining-experts/>

# AMiner Services

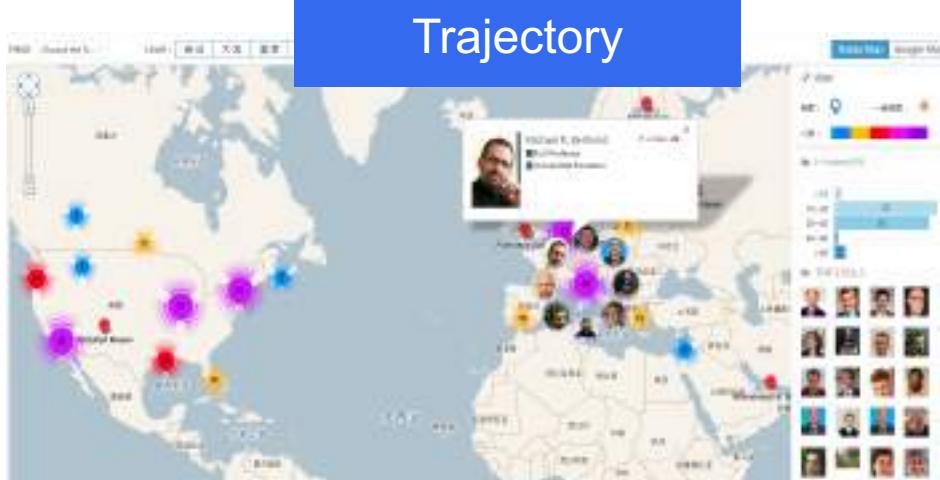
# Trend



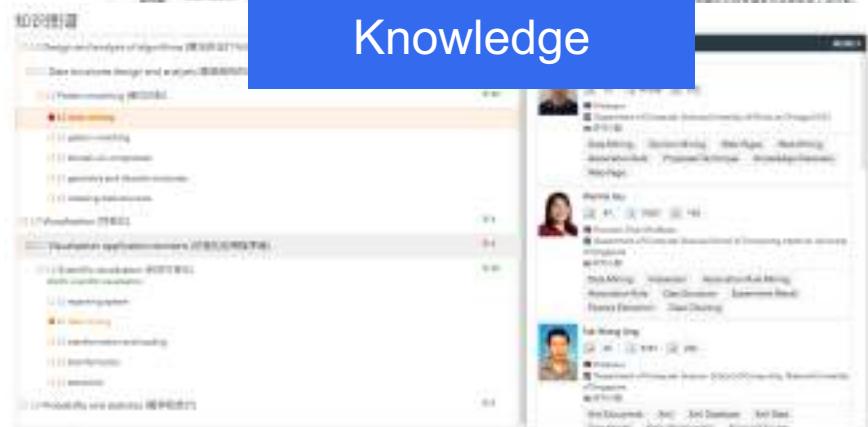
## Expert Finding



# Trajectory



# Knowledge



# Talent Map

阿里巴巴 学术资源地图 Data Mining 搜索 首页 地图

# Career Trajectory of Top Experts



# AMiner

## —AI-Powered Academic Search

<https://aminer.cn>

### Thanks to our partners

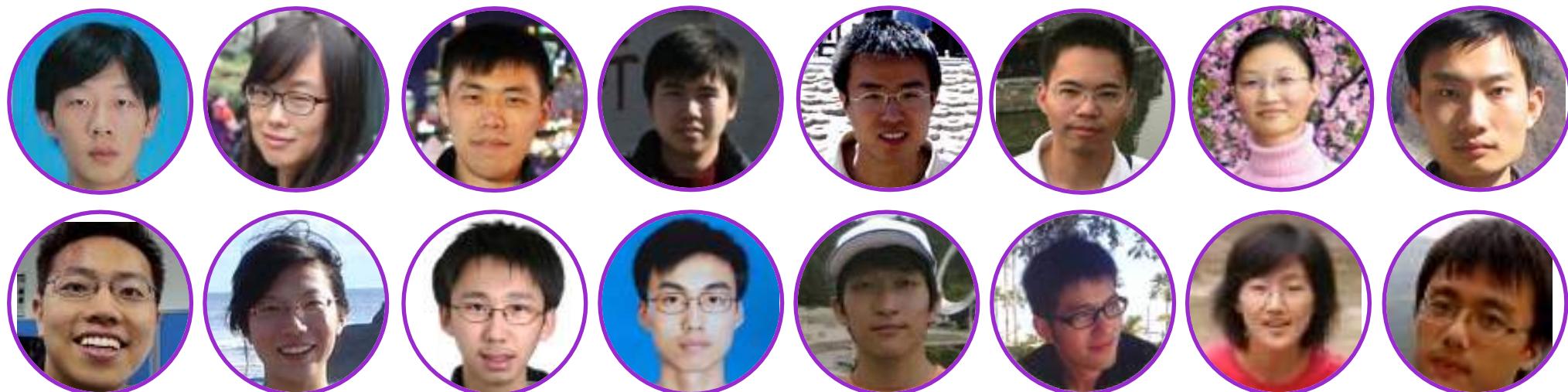




# AMiner

## —AI-Powered Academic Search

<https://aminer.cn>





谢谢观看！



# 人工智能是学术搜索的未来吗？

AI Time 第十二期

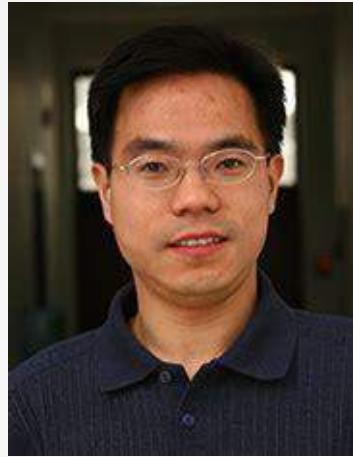
AI Time是一群关注人工智能发展，并有思想情怀的青年人创办的圈子。AI Time旨在发扬科学思辨精神，邀请各界人士对人工智能理论、算法、场景、应用的本质问题进行探索，加强思想碰撞，打造成为北京乃至全国知识分享的聚集地。



## ■ 嘉宾介绍



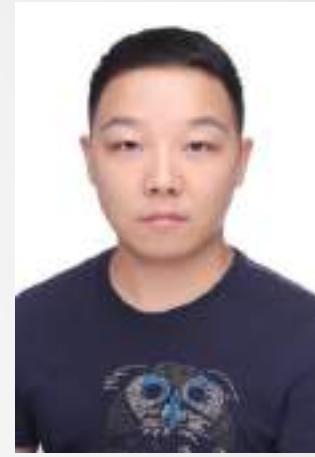
唐杰  
清华大学计算机系教授  
副系主任



周园春  
中科院网络信息中心研究员，  
大数据技术与应用发展部主任



刘筱敏  
中科院文献情报中心研究馆员，  
知识服务创业中心数据产品与  
管理服务事业部主任



刘雪峰  
百度学术搜索产品负责人



李文珏  
中科创星投资总监  
CCF YOCSEF学术秘书福  
布斯30Under30



杜一  
中科院网络信息中心大数  
据知识工程实验室主任  
北京市科技新星

## ■ 主持人



## ■ 目录

- 1. 学术搜索的内涵与外延**
- 2. 思辨及讨论**
  - a. 技术篇**
  - b. 应用篇**
  - c. 产业篇**
- 3. 总结与畅想**



## ■ 技术篇

- 学术搜索与普通搜索有什么区别，是更难还是更容易？
- 在学术搜索领域，我国现状是处在领跑、并跑还是跟跑？
- 关于学术搜索中的实体、知识融合，学术界已经有很多方法，但是在实践上并不能完全满足要求。如何应对这类问题？
- 在下列五个场景中：自动摘要生成、交叉学科发现、趋势分析与预测、学术影响力评价、专家&机构画像，还有哪些技术可以赋能？





## ■ 应用篇

- 从用户角度（例如学者、期刊/会议、高校、基金资助机构等），AI在哪些场景已经发挥了作用？哪些场景还需要改进？还有哪些需求没有被满足？
- 在学术搜索领域还有哪些场景可以用AI赋能？
- 学术搜索+AI对国内外学术环境有哪些影响？





## ■ 产业篇

- 目前谷歌、百度、微软等大型企业都持续在学术搜索上发力，作为Semantic Scholar、Aminer这样的学术搜索平台，在哪些方向会有所突破？
- 面对Aminer这样的平台的冲击，百度学术这样的传统学术搜索引擎，应该如何响应？
- 学术搜索对学术出版，科技情报产业的意义在哪里，能否突破产业局限？还有哪些商业上的想象空间？
- 自内而外 VS 自外而内的模式，这两种模式在应用场景、构建模式上有何不同？



## ■ 畅想未来

- 作为计算机、文献情报界以及产业界的代表，各位对学术搜索的未来怎么看？
- 为学术搜索领域从业者（科研人员、项目研发、产品设计等）提供一些建议？





# ■ It Is Your Turn

## Q&A



谢谢大家



AI Time欢迎每一位有情怀的AI爱好者加入！我们需要您的思辨和碰撞！