

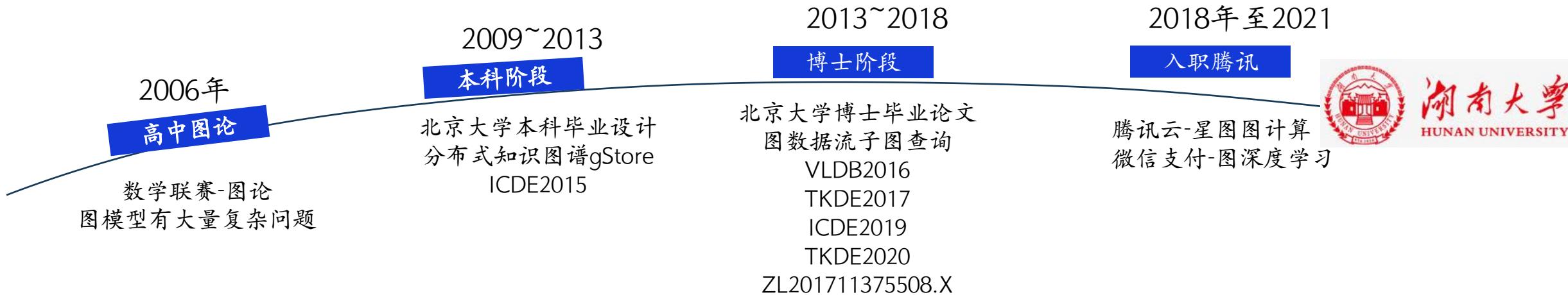


HNU

从图论到图数据模型

李友焕 副教授
liyouhuan@hnu.edu.cn

图相关经历与内容概览



01

图模型

02

图查询

03

图计算

04

图学习

传统的关系数据模型

找出X的所有好友

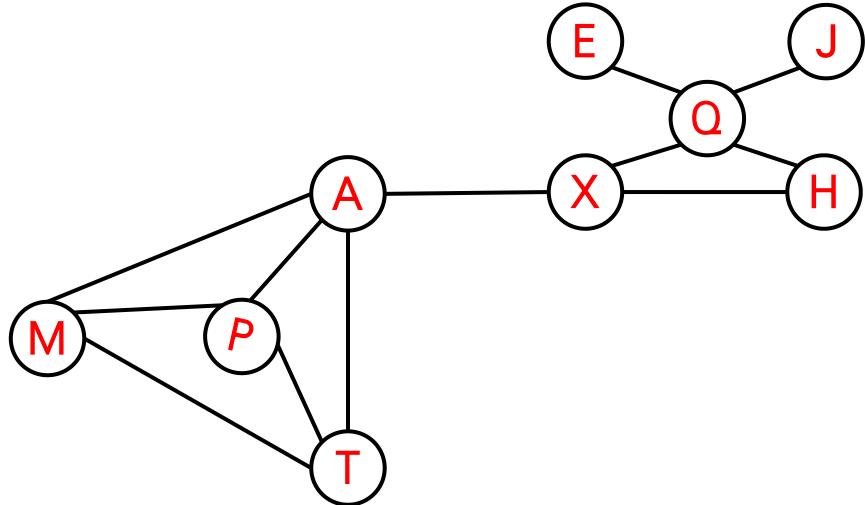
找出朋友不止一个的人

找出所有的三人组，
其中三个人互相认识

传统好友关系表

Name1	Name2
A	M
A	P
A	T
A	X
E	Q
H	X
H	Q
J	Q
M	P
M	T
Q	X
P	T

图谱数据模型



找出朋友不止一个的人

找出所有的三人组，其中三个人互相认识

邻接表

A	M, P, T, X
→	
E	Q
→	
H	→ X, Q
J	Q
→	
M	→ P, T
Q	→ E, H, J, X
P	→ A, T, M
T	→ A, M, P
X	→ A, Q, H

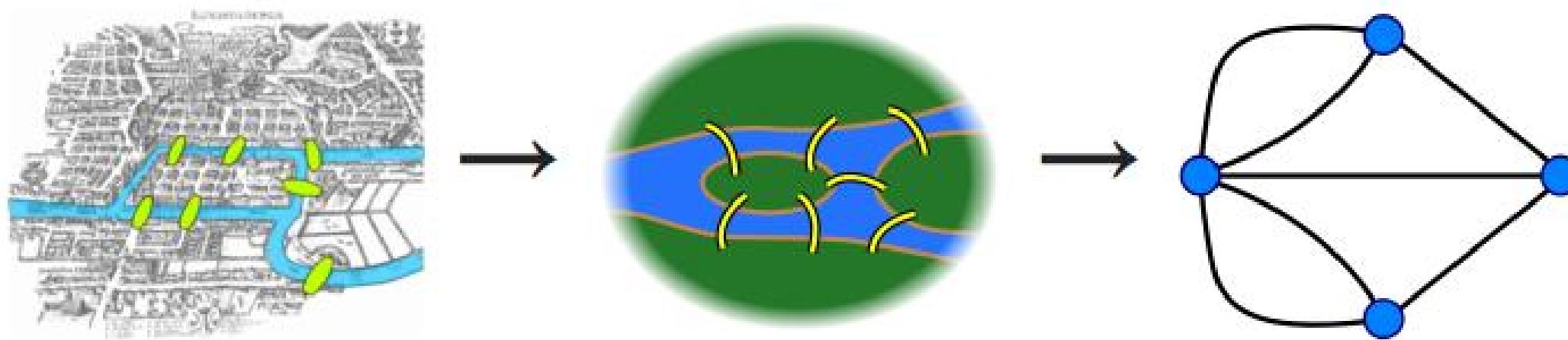
图论起源：欧拉七桥问题

Leonhard Euler



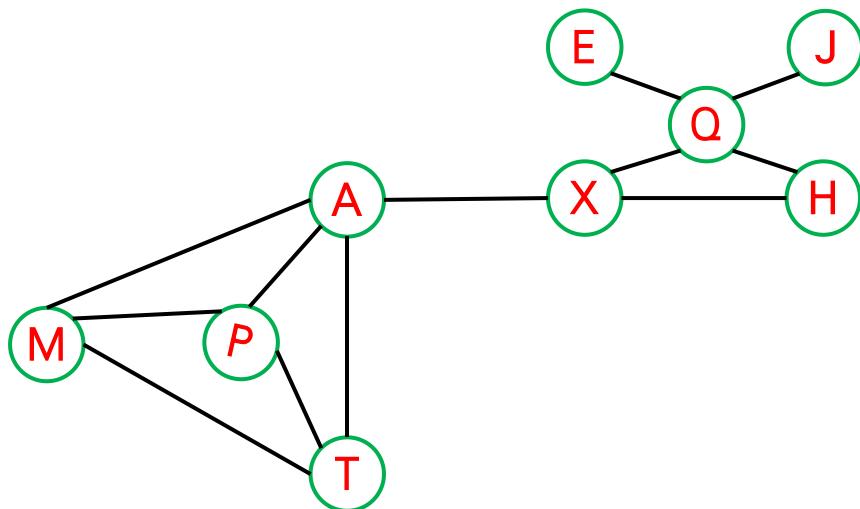
欧拉七桥问题：如何能够不走重复路的情况下走遍哥尼斯堡的七座桥？

图论思路：欧拉通过将七桥问题形式化为点边的一笔画问题来解决。



图形式简单，图技术复杂

点、边、度数



绿色是点，对应现实世界对象

点间连的黑线是边，对应现实世界对象之间的关联

度数是定义在每个点上的一个数值，即跟该点有连边的其它点的数量

请求出所有顶点的度数

一笔画体验图谱拓扑、结论

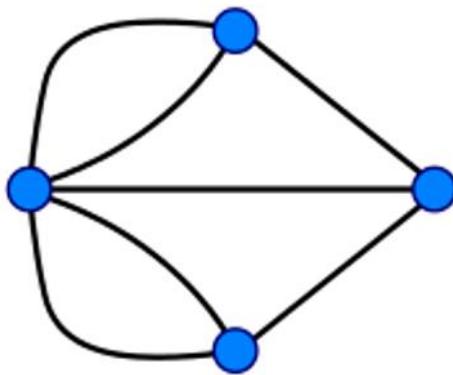


图 A

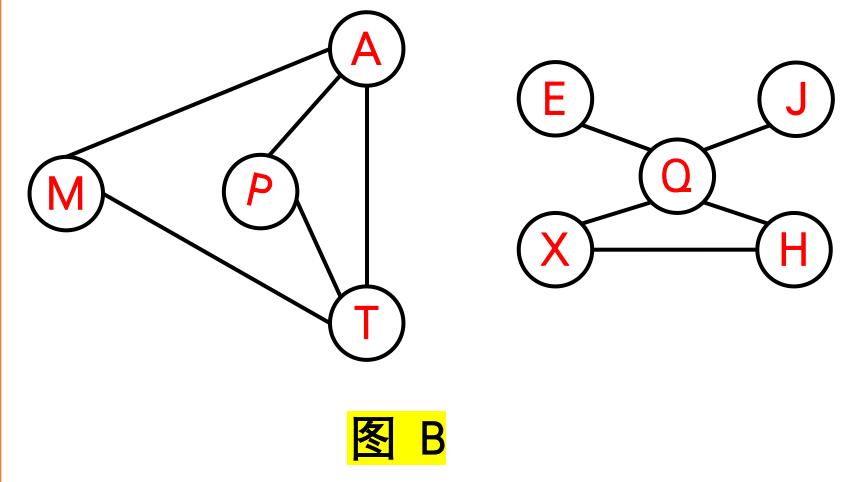


图 B

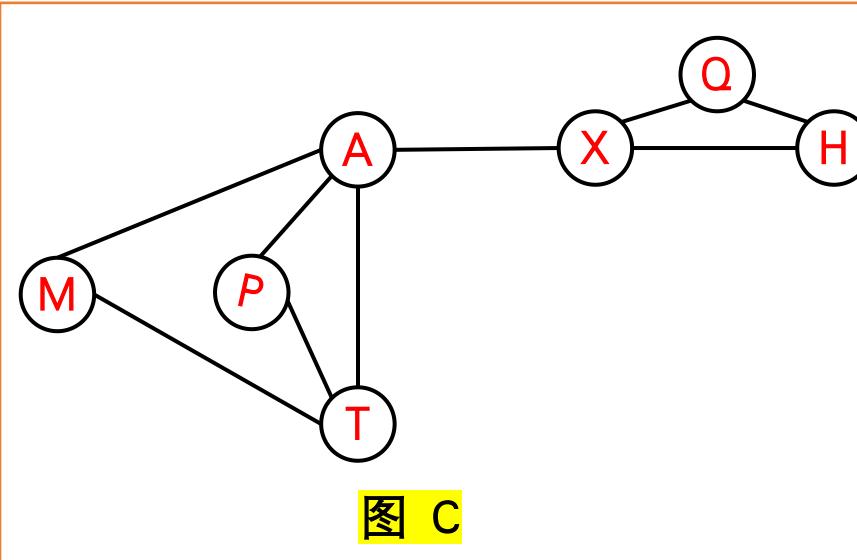


图 C

1. 图形必须是连通的。
2. 图中的“奇度数点”个数是0或2。

图的形式化定义

图论中，图被定义为一个多元组

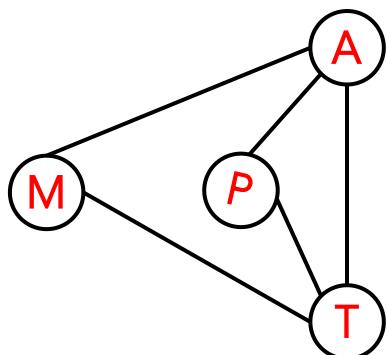
$$G = (V, E, L, P\dots)$$

V 为点 (vertex) 集，
对应对象的集合

$E \subseteq V \times V$ 为边
(edge) 集， 对应
对象间关系、交互
的集合

L 为点/边标签函数
(可选)
 $L(v1)$
= ‘诈骗’， $L(e1)$
= ‘好友’

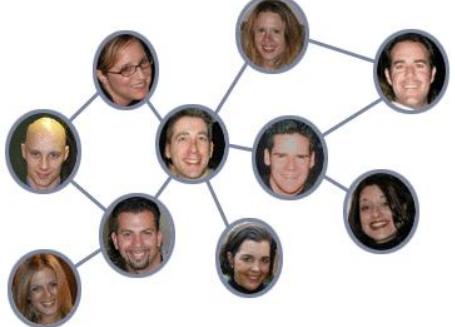
P 为点/边属性函数
(可选)
 $P(v1) = 28$, $P(e1)$
= 20201125



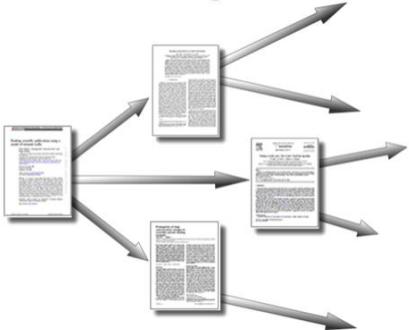
$$G = (V, E), \quad V = \{A, P, T, M\}, \quad E = \{(A, P), (A, T), (A, M), (T, P), (T, M)\}$$

点、边均只有一种类型的，为同构图，否则为异构图

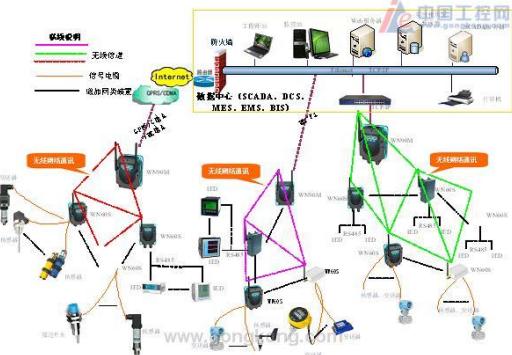
现实中的图数据



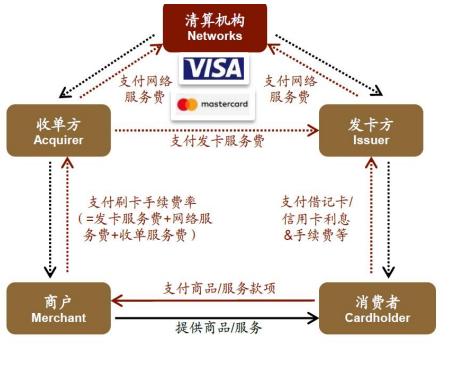
□ 社交网络
□ 好友关系



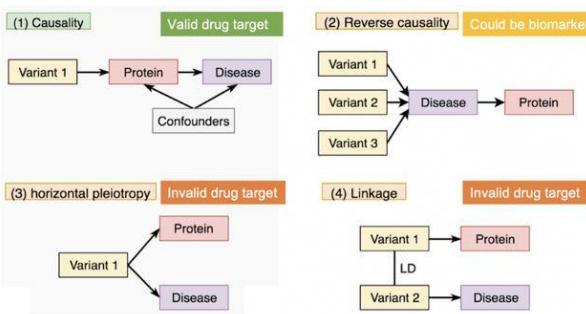
□ 文献引用网络
□ 合作/引用



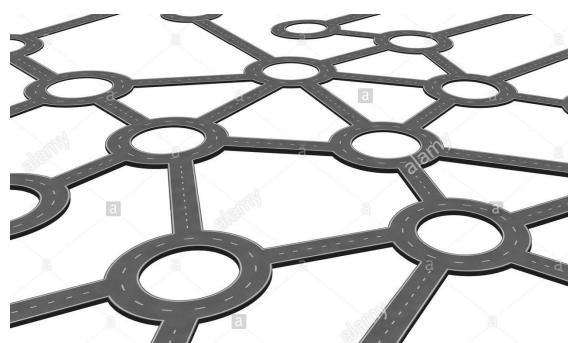
□ 通信网络
□ ip通信



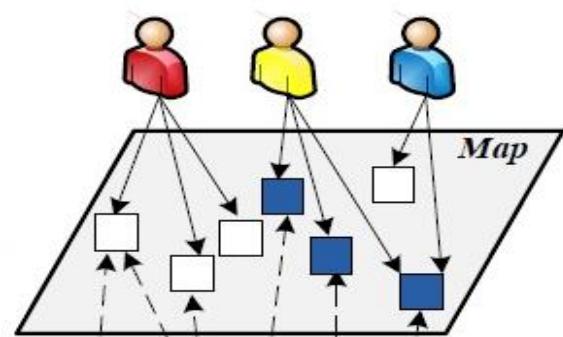
□ 支付网络
□ 资金交易



□ 基因疾病网络
□ 基因-蛋白质-疾病



□ 道路网络
□ 站点通道



□ 移动轨迹网络
□ 人与地理位置

图模型表达能力

信息

图

图模型 $G = (V, E)$

- 对象 → 点
- 对象间关联 → 边



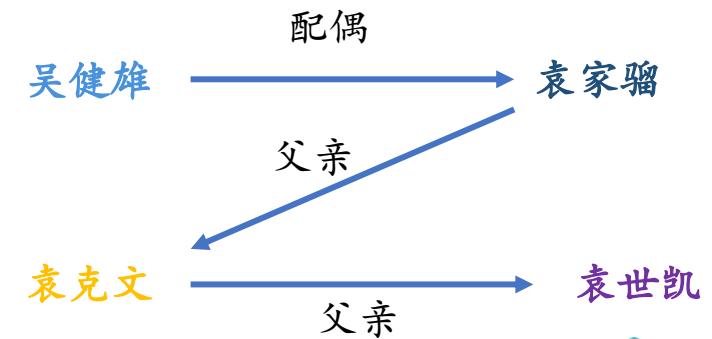
信息表达

- 知识/信息 → 陈述句
- 陈述句 → 主谓宾

吴健雄配偶袁家骝
袁家骝的父亲是袁克文
袁克文的父亲是袁世凯

关联分析

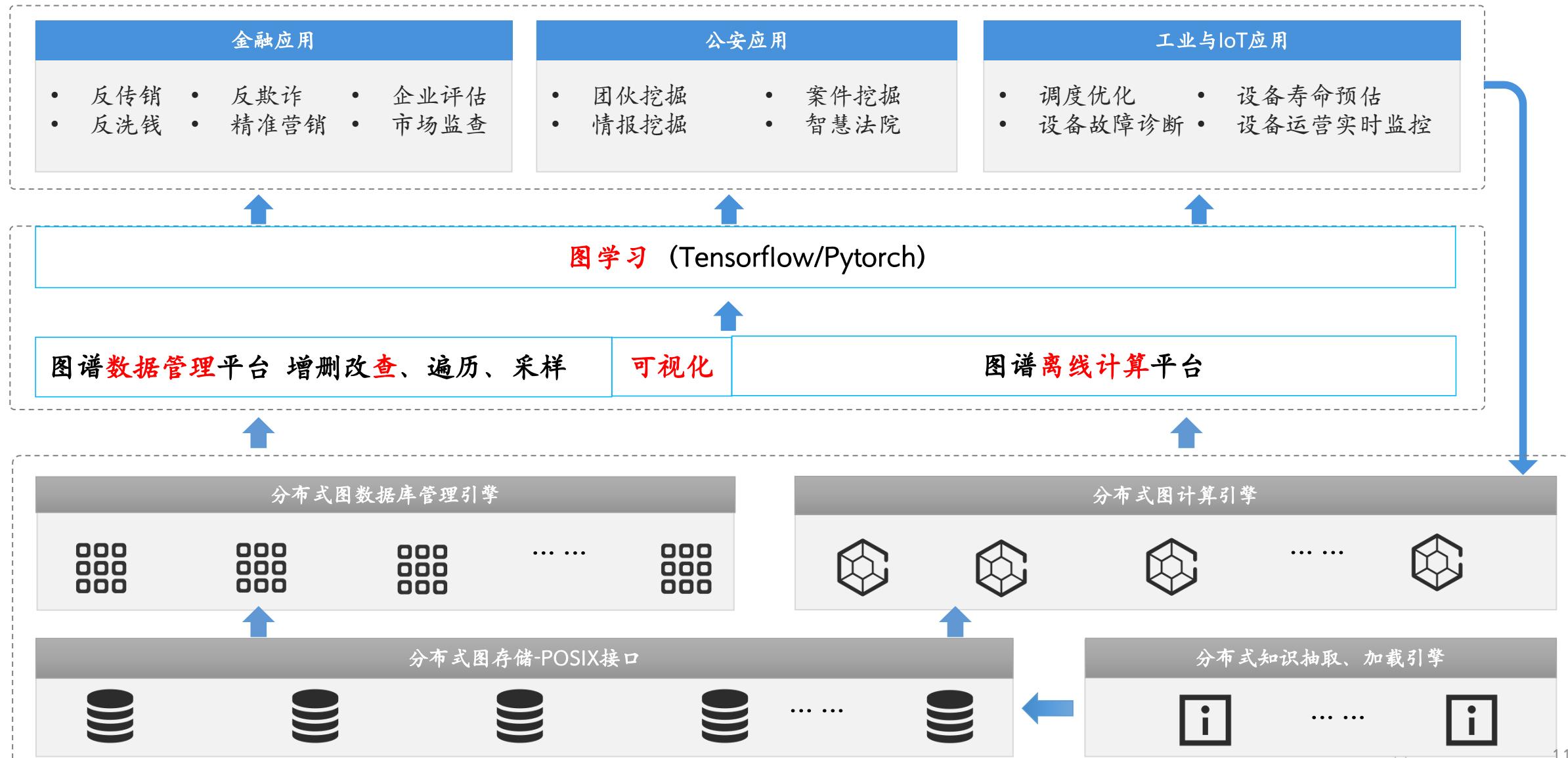
- 对象关联
- 路径



复杂知识/信息的建模、表达与融合

高效的关联发现与分析

图技术四大模块



图信息获取之图查询

图查询
结构关联

邻居信息、路径信息、子图



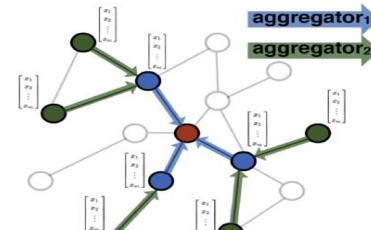
图计算
全局属性

PageRank、K-core

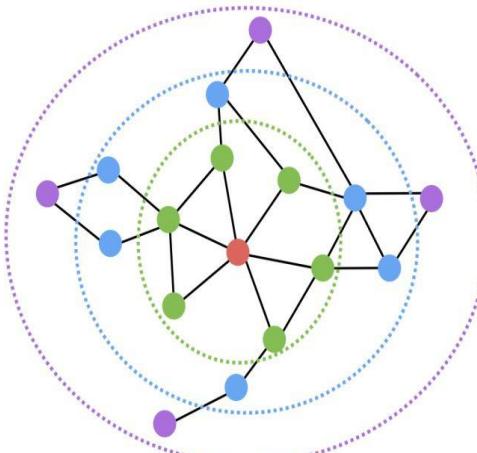


图学习
向量表示

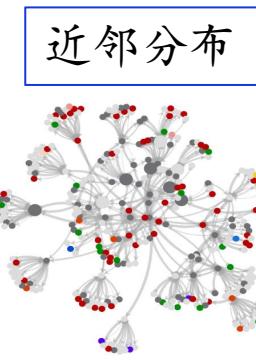
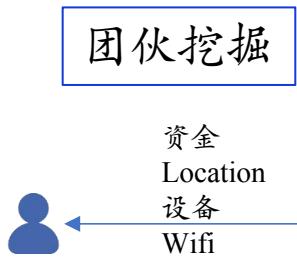
GraphSAGE、LINE



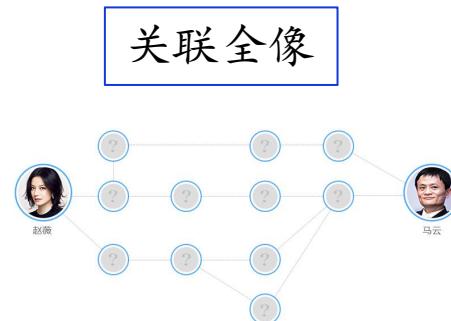
图查询：一阶、二阶邻居查询



- **Red:** Target node
- **Green:** 1-hop neighbors
 - A (i.e., adjacency matrix)
- **Blue:** 2-hop neighbors
 - A^2
- **Purple:** 3-hop neighbors
 - A^3



招商银行
CHINA MERCHANTS BANK

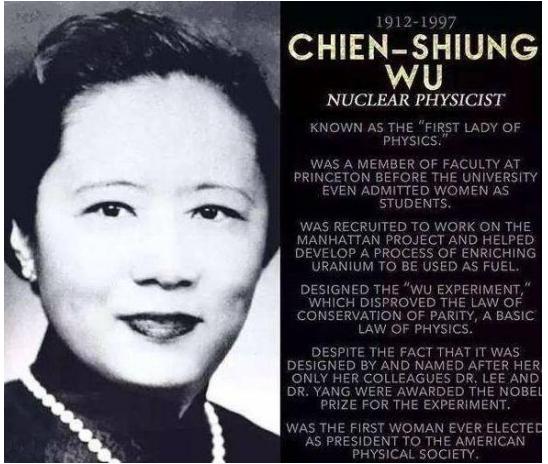


天眼查
www.tianyancha.com



朋友圈

图查询：路径查询



1912-1997
CHIEN-SHIUNG WU
NUCLEAR PHYSICIST

KNOWN AS THE "FIRST LADY OF PHYSICS."

WAS A MEMBER OF FACULTY AT PRINCETON BEFORE THE UNIVERSITY EVEN ADMITTED WOMEN AS STUDENTS.

WAS RECRUITED TO WORK ON THE MANHATTAN PROJECT AND HELPED DEVELOP A PROCESS OF ENRICHING URANIUM TO BE USED AS FUEL.

DESIGNED THE "WU EXPERIMENT," WHICH DISPROVED THE LAW OF CONSERVATION OF PARITY, A BASIC LAW OF PHYSICS.

DESPITE THE FACT THAT IT WAS DESIGNED BY AND NAMED AFTER HER ONLY HER COLLEAGUES DR. LEE AND DR. YANG WERE AWARDED THE NOBEL PRIZE FOR THE EXPERIMENT.

WAS THE FIRST WOMAN EVER ELECTED AS PRESIDENT TO THE AMERICAN PHYSICAL SOCIETY.

....

吴健雄配偶袁家骝

袁家骝的父亲是袁克文

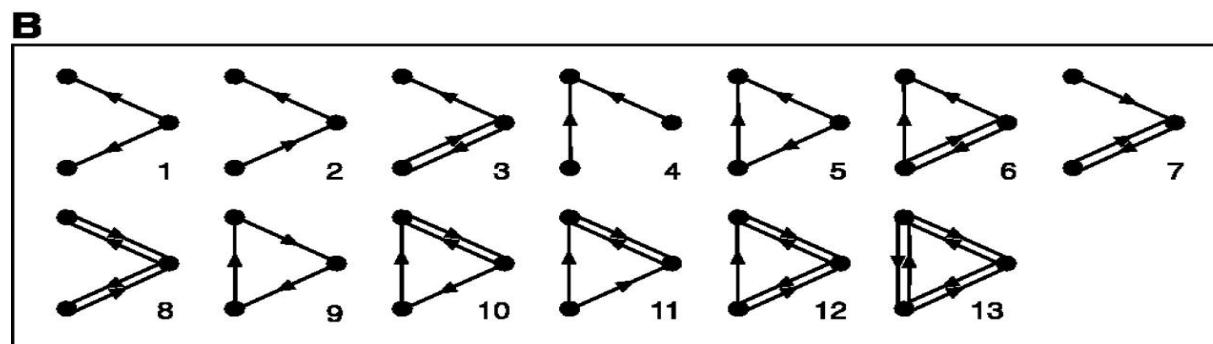
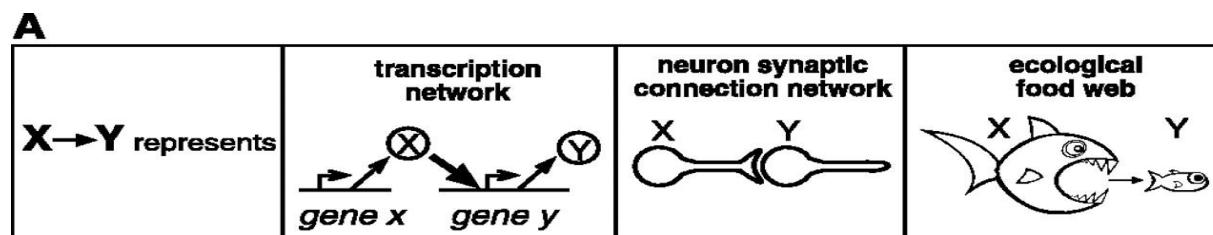
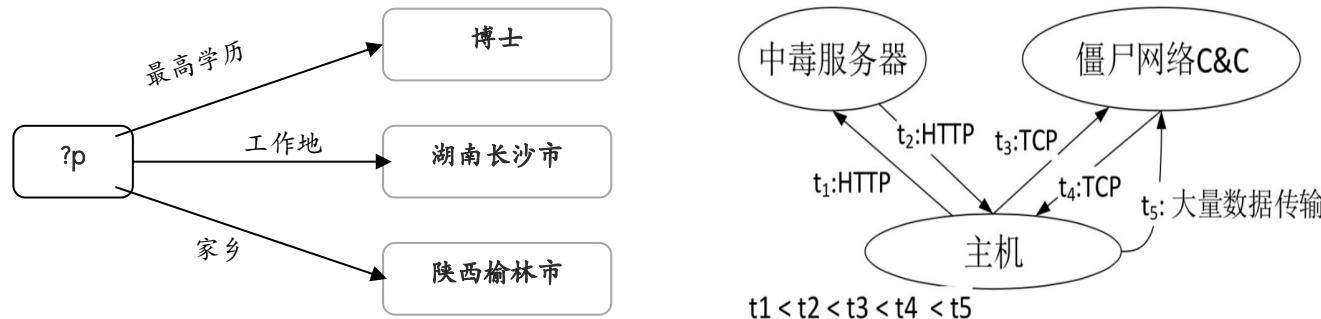
袁克文的父亲是袁世凯

....



图查询：子图查询

子图：以月度转账图为例，当前听众之间互相转账形成的图就是全国月度转账网络的子图



数据分析业务增效利器

图数据库

工业界图数据库的定位很多时候是分析平台，图数据库并不是作为关系数据库的替代，而是补充。

工业界需要图模型方法论来保证相应效率

<https://db-engines.com/en/ranking/graph+dbms>

Rank			DBMS	Database Model	Score		
Feb 2025	Jan 2025	Feb 2024			Feb 2025	Jan 2025	Feb 2024
1.	1.	1.	Neo4j	Graph	45.34	+1.65	-1.27
2.	2.	2.	Microsoft Azure Cosmos DB	Multi-model	22.13	-0.83	-9.86
3.	3.	3.	Aerospike	Multi-model	5.09	+0.05	-2.00
4.	4.	4.	Virtuoso	Multi-model	3.14	-0.18	-2.18
5.	5.	5.	ArangoDB	Multi-model	2.91	-0.01	-1.43
6.	6.	6.	OrientDB	Multi-model	2.73	-0.06	-0.92
7.	7.	↑ 9.	GraphDB	Multi-model	2.63	-0.09	-0.12
8.	8.	↓ 7.	Memgraph	Graph	2.54	-0.08	-0.63
9.	9.	↓ 8.	Amazon Neptune	Multi-model	2.12	-0.07	-0.81
10.	10.	↑ 11.	JanusGraph	Graph	1.69	-0.03	-0.61
11.	11.	↑ 12.	Stardog	Multi-model	1.61	-0.07	-0.52
12.	12.	↓ 10.	NebulaGraph	Graph	1.58	-0.08	-0.96
13.	13.	↑ 14.	Fauna	Multi-model	1.49	+0.07	-0.30
14.	14.	↓ 13.	TigerGraph	Graph	1.45	+0.03	-0.52
15.	15.	15.	Dgraph	Graph	1.20	-0.05	-0.40
16.	16.	16.	Giraph	Graph	1.07	-0.01	-0.23
17.	17.	17.	SurrealDB	Multi-model	0.95	-0.02	-0.17
18.	18.	↑ 20.	Blazegraph	Multi-model	0.73	+0.02	-0.02
19.	19.	↓ 18.	AllegroGraph	Multi-model	0.70	-0.02	-0.39
20.	20.	↓ 19.	TypeDB	Multi-model	0.59	-0.02	-0.47

图系统是数据分析利器之一



Apache Flink

图数据库易用性的关键



足够简易就能可视化

产品除了自身满足功能需求之外，还要重视使用者的体验，如果产品使用对使用者造成大的时间成本，则容易失去市场

关联查询性能对比

异常用户
24*****30

SQL **VS** 模版UI

```
1 SELECT DISTINCT dst FROM (
2     SELECT t2.dst FROM (
3         SELECT dst
4             FROM [REDACTED]
5             PARTITION(p_201904)
6     ) t1 JOIN [REDACTED]
7             PARTITION(p_201904) t2
8         ON (t1.dst=t2.src)
9     UNION ALL
10    SELECT dst
11        FROM [REDACTED]
12        PARTITION(p_201904) t1 WHERE src=24*****30
13 );
```

submit time | sql string | processing | finish time
2019/05/27 16:47:18.1 | SELECT t1.uin,N... | 100% | 2019-05-27 17:04:29.

→ 17min ←

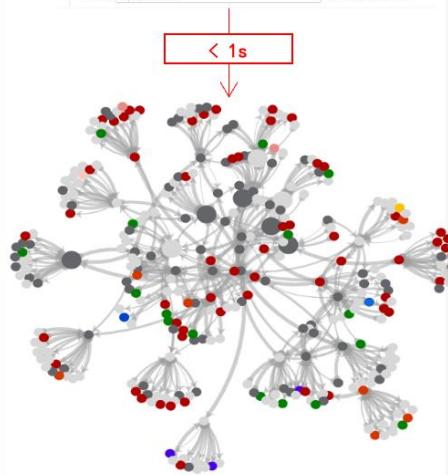
1	1C ** 63
2	12 ** 15
3	25 ** 45
4	34 ** 55 轻度返利
5	56 ** 40 养号
6	61 ** .78
7	64 ** 82
8	11 ** 100 全量游戏
9	12 ** 900
10	12 ** 420 涉赌
11	14 ** 865
12	15 ** .564
13	16 ** 061 红包赌博
14	16 ** 240 涉赌
15	16 ** .721
16	18 ** 417 外挂/ [REDACTED] 赌账号被...

设置属性颜色

属性名称: abnormal_type

属性值:	颜色值: #0000ff
属性值:	颜色值: #cccccc
属性值:	颜色值: #ffff00
属性值:	颜色值: #ff9999
属性值:	颜色值: #000080
属性值:	颜色值: #008000
属性值:	颜色值: #800000
属性值:	颜色值: #800080
属性值:	颜色值: #ff0000
属性值:	颜色值: #ff00ff
属性值:	颜色值: #00ffff
属性值:	颜色值: #0000ff

< 1s



表格 **VS** 网络图形

关系模型在大数据场景的缺点

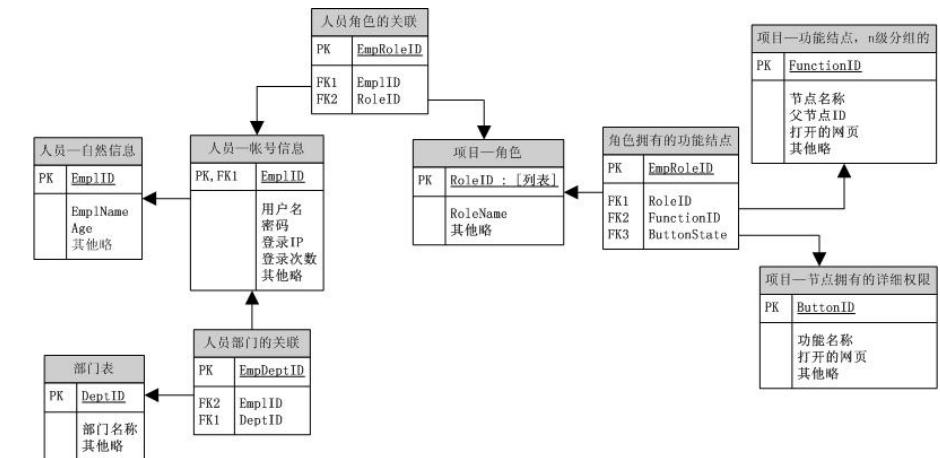
大数据场景下的数据

- 数据来源复杂
- 数据形式多样
- 数据规模各异
- 数据动态变化

关系模型

- 关联模型处理结构化数据，管理成本高

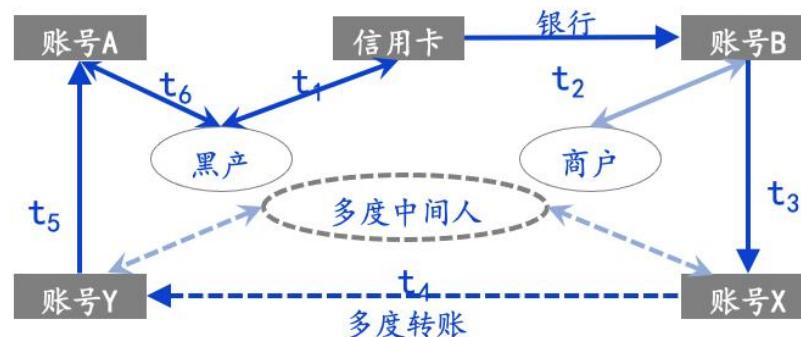
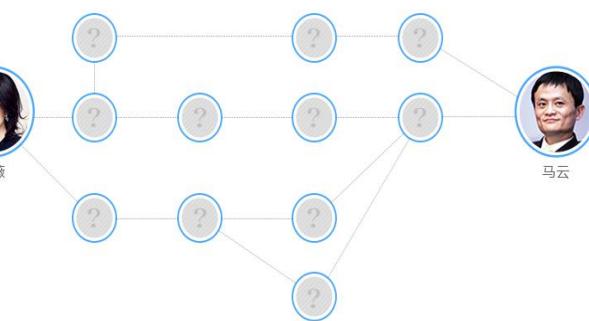
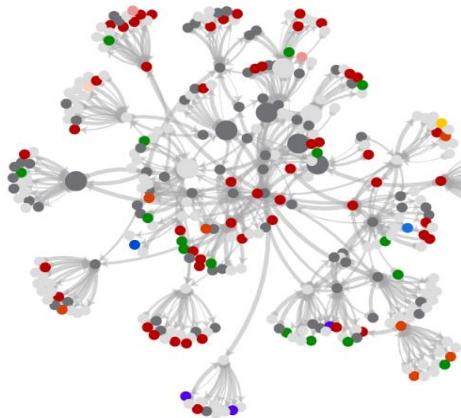
BIG DATA



关系模型在大数据场景的缺点

大数据场景下的查询

- 查询关联多样：直接关联、间接关联、多



关系模型

- 关系模型查询效率低，尤其在常见的多阶关

Meals	
Omlet	🔍
Fried Egg	🔍
Sausage	🔍

Drinks	
Orange Juice	🔍
Tea	🔍
Coffee	🔍

CROSS JOIN

Menu Combination	
🔍	🥤
🔍	🥤
🌭	🥤
🔍	🥤
🍳	🥤
🌭	🥤
🔍	☕️
🍳	☕️
🌭	☕️

关联查询的优化方向：关联聚集

毕业

Name	Name
史玉柱	深圳大学
史玉柱	浙江大学
张一鸣	南开大学

创办

Name	Name
史玉柱	巨人网络
张维	深圳大学
陈嘉庚	厦门大学

外公

Name	Name
高晓松	张维
溥仪	荣禄
....

史玉柱

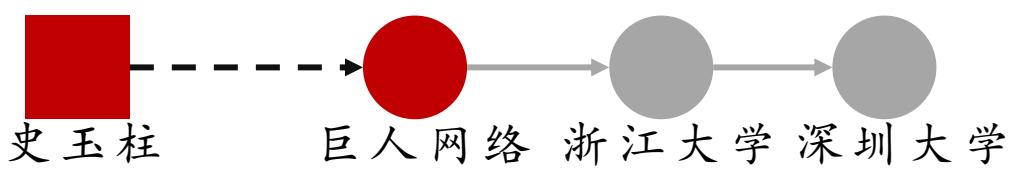
张维

Name	Name	关系
史玉柱	深圳大学	毕业
史玉柱	浙江大学	毕业
史玉柱	巨人网络	创办

Name	Name	关系
深圳大学	张维	创办人
张维	高晓松	外孙
张维	深圳大学	创办

关联查询的优化方向：去除行冗余

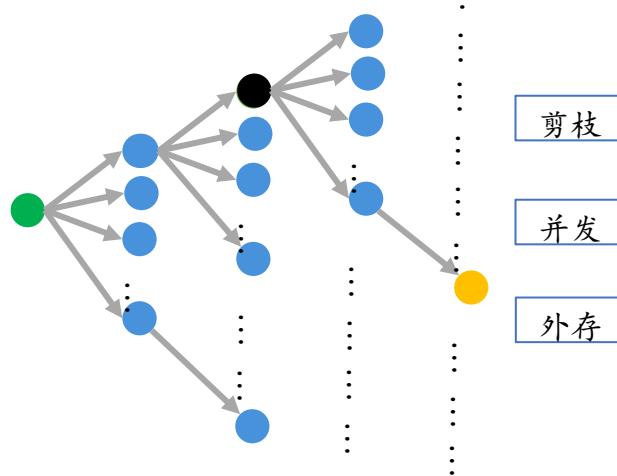
Name	Name	关系
史玉柱	深圳大学	毕业
史玉柱	浙江大学	毕业
史玉柱	巨人网络	创办



关联查询的优化方向：KV之上的图方法论

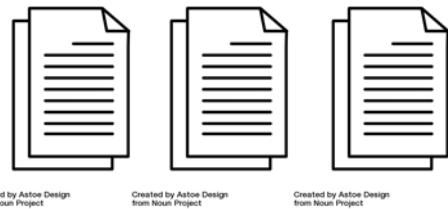
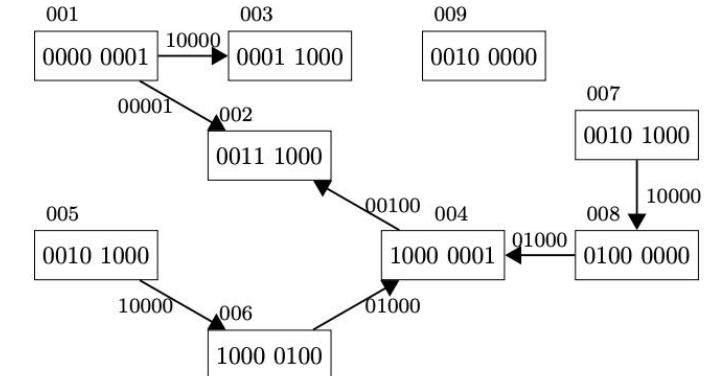
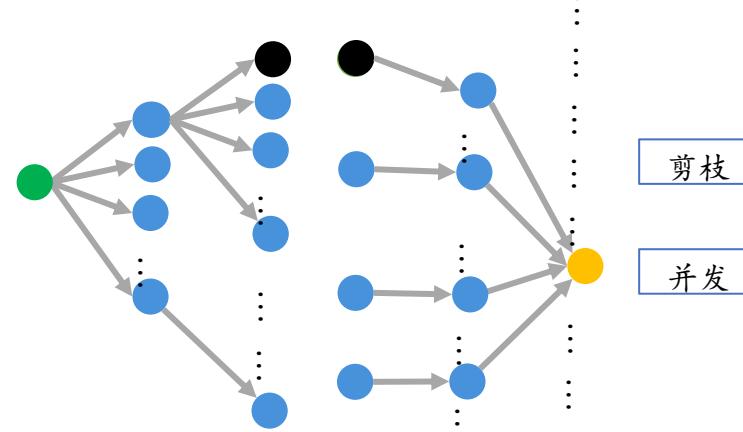
暴力宽度搜索

- 深度 $d \rightarrow N^d$
- $100^4 = \text{亿级}$

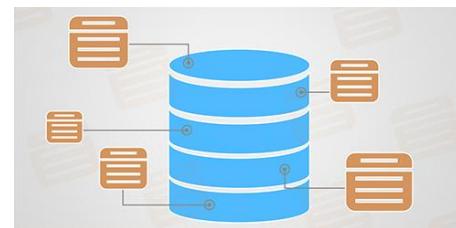


双向搜索

- 深度 $d \rightarrow N^{d/2}$
- $100^2 = \text{万级}$



文件系统



Database

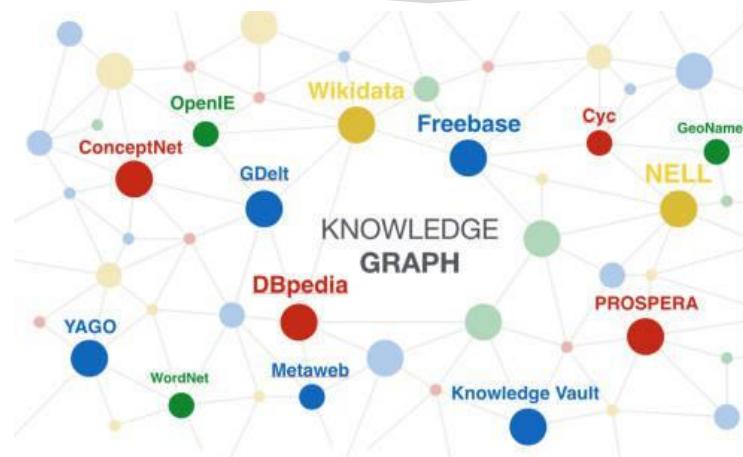
知识图谱

知识和图谱的融合，基于图的模型与方法，对知识进行建模、管理、查询、分析进而对实际应用输出解决方案的一整套技术体系。

图 谱



知 识



图信息获取之图计算

图查询
结构关联

邻居信息、路径信息、子图



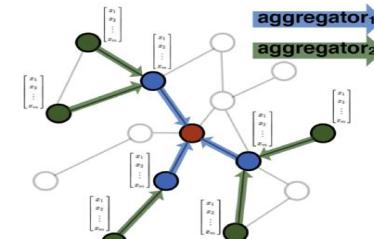
图计算
全局属性

PageRank、K-core

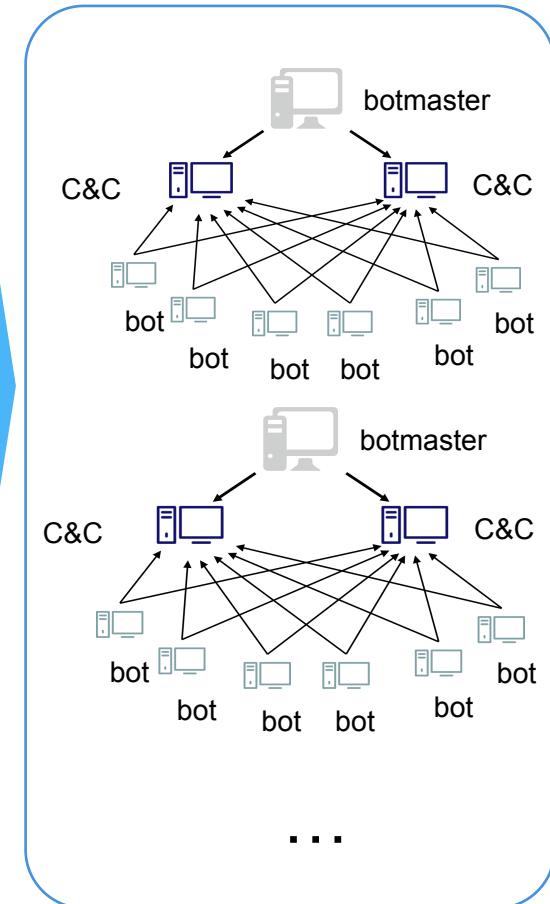
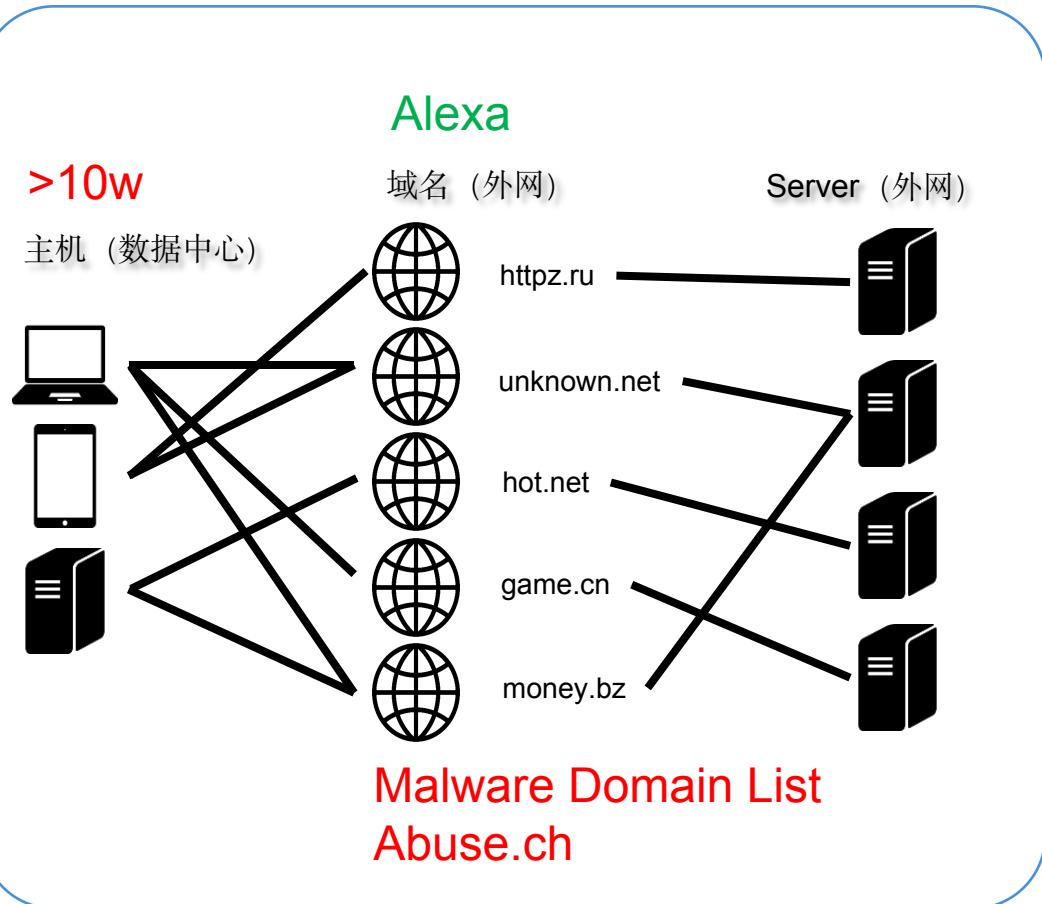
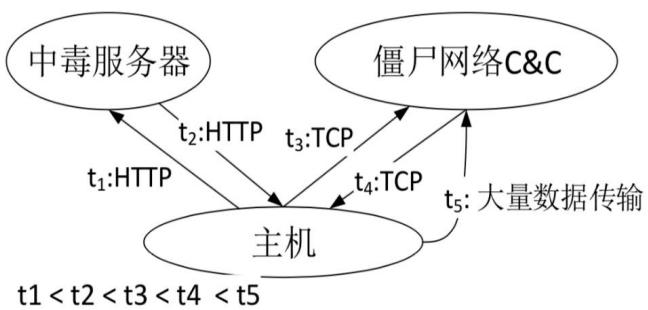


图学习
向量表示

GraphSAGE、LINE



主机入侵检测



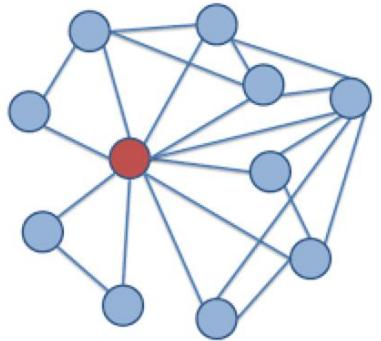
特征?

行为!

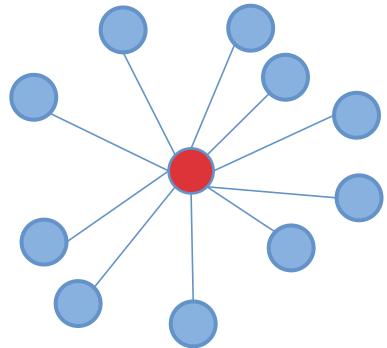
基于全局关联结构
放大 细微的、隐秘的 关联

Motif

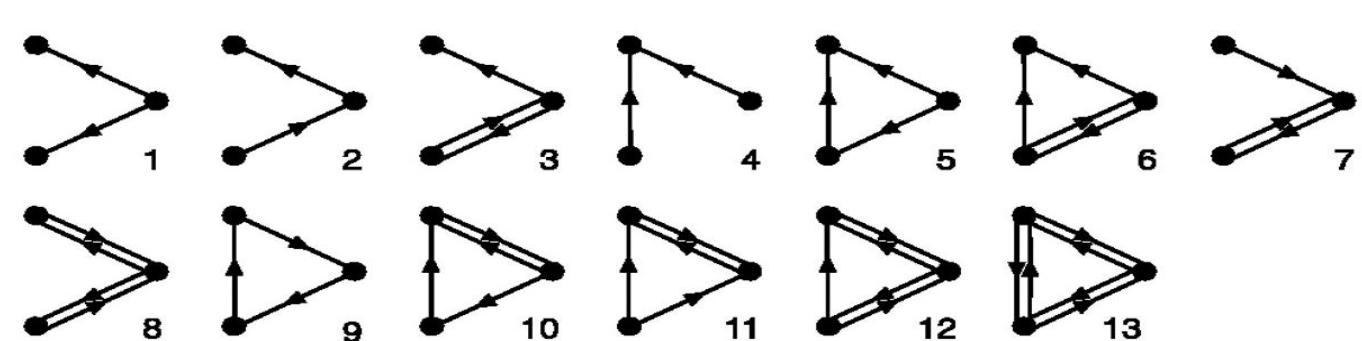
EGO网络1-同学群



EGO网络2-赌博群



子结构示意



案例分享：刷单团队人员角色的子结构分布



“正常商家”



“恶意商家”



“中介”



“刷手”



“普通用户”

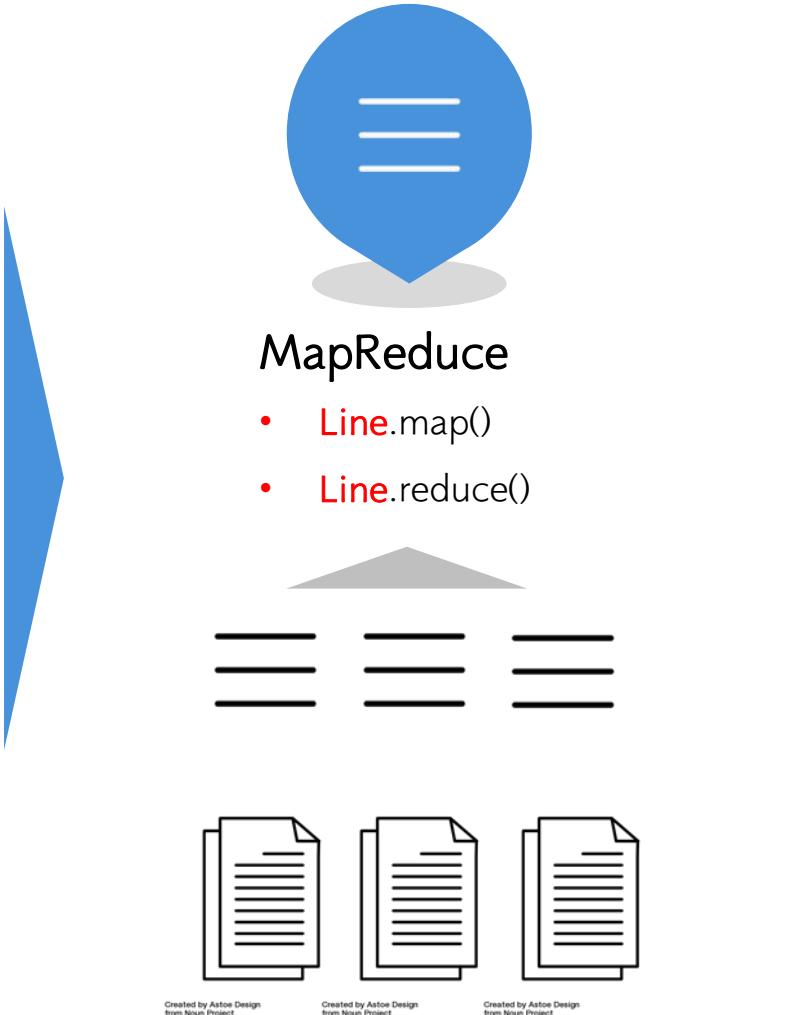
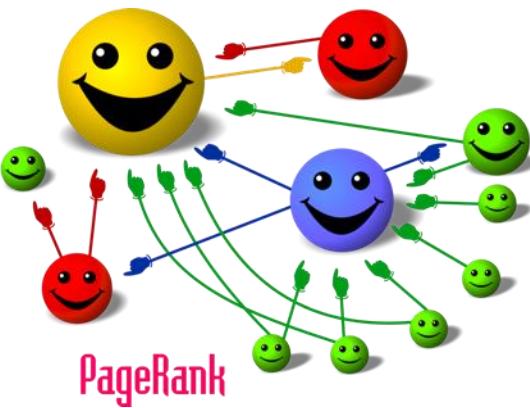
点中心图计算引擎原理

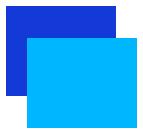
Pagerank/LPA

- Pagerank: 邻居权重 → 点权重
- LPA: 邻居标签 → 点标签

点中心编程

- 将大规模图计算抽象为点操作
- 调度集群更新点状态至全图收敛



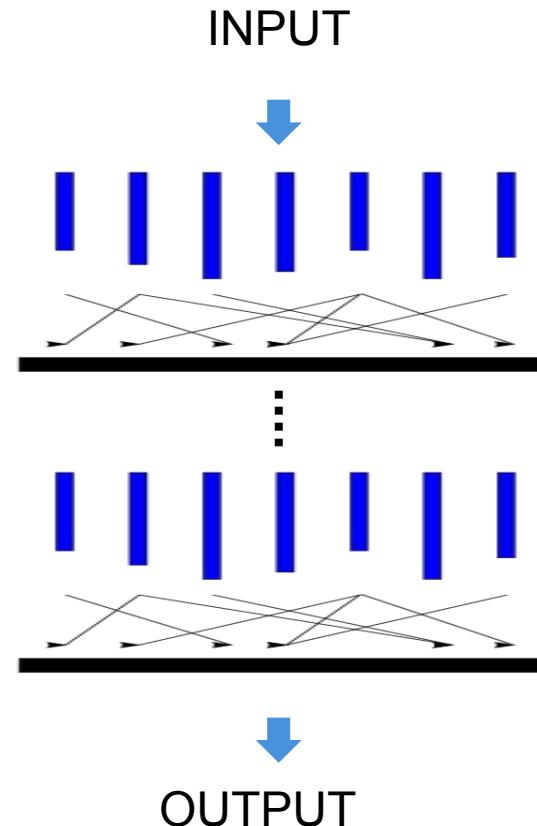


SIGMOD2010 (Google)

- 图数据访问的局部性很差
- 每个点计算量很少
- 图系统开发的复杂性

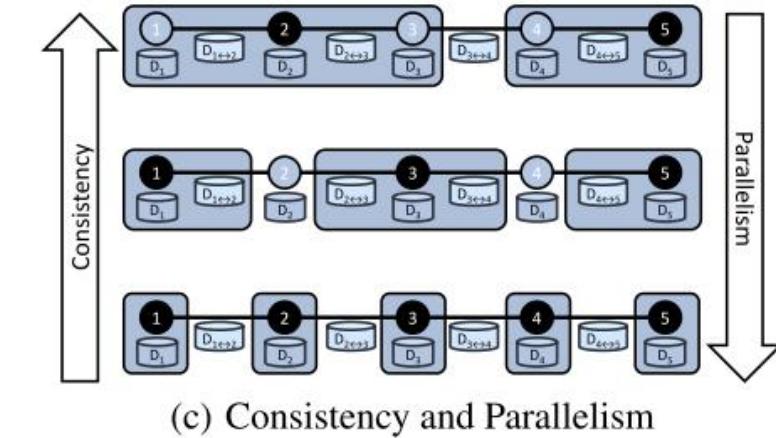
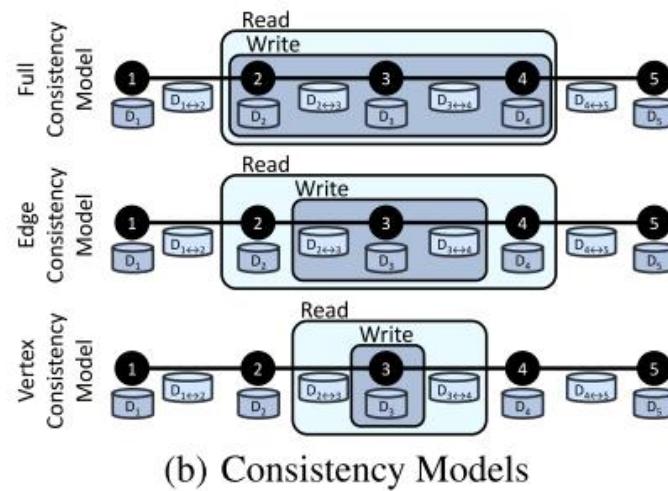
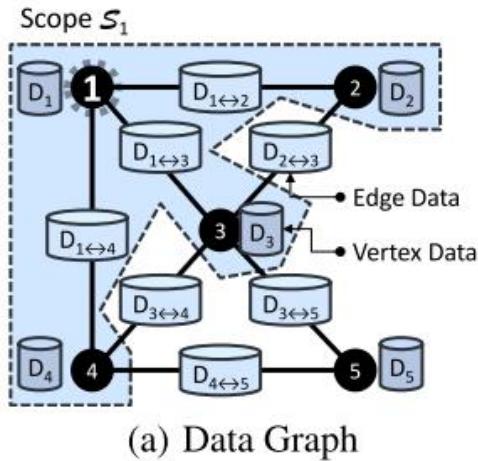
点中心

- `Vertex_update(v, iterator<message>)`
 - 更新点v的状态
 - 从点v向外发送消息
- 缺点
 - 强同步：水桶效应
 - 大量消息通信：边 >> 点

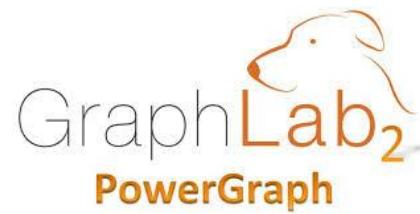


VLDB2012 (UCB AMPlab)

- 2~3X Pregel
- 分布式共享内存
- 随时读/写邻居状态
- 弱同步，本轮结束的点，直接开启下一轮计算



PowerGraph

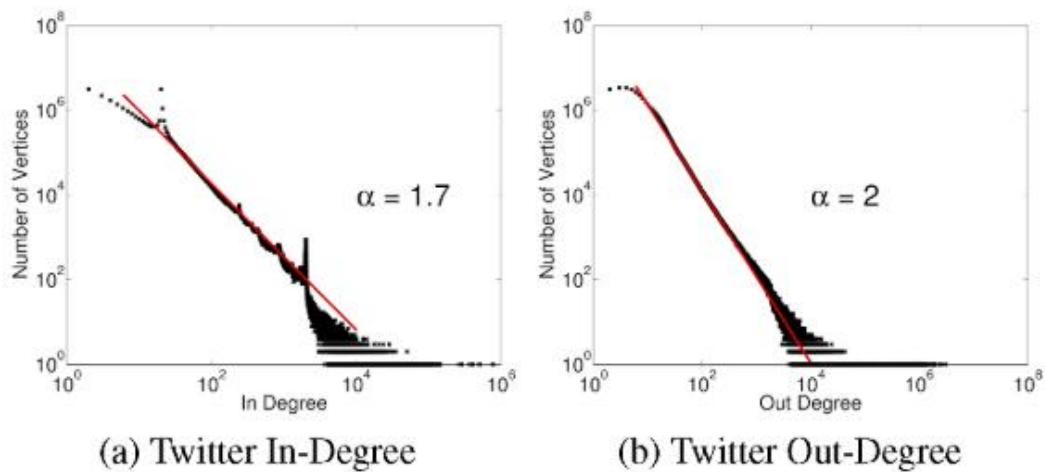


OSDI2012 (UCB AMPLab)

- 5X GraphLab

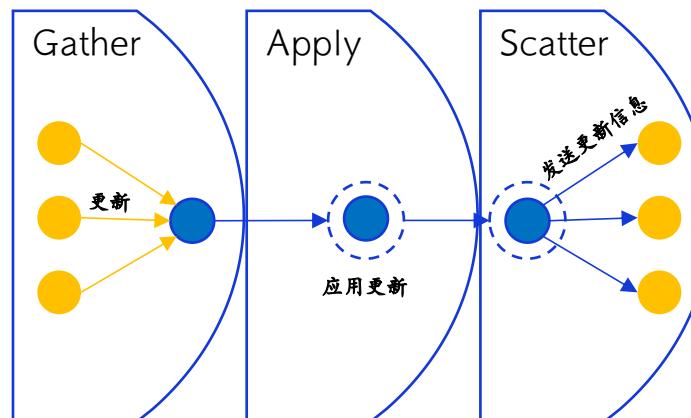
度数的指数分布：基于度数的分割 (twitter)

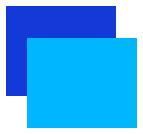
- 2~3X Random



GAS

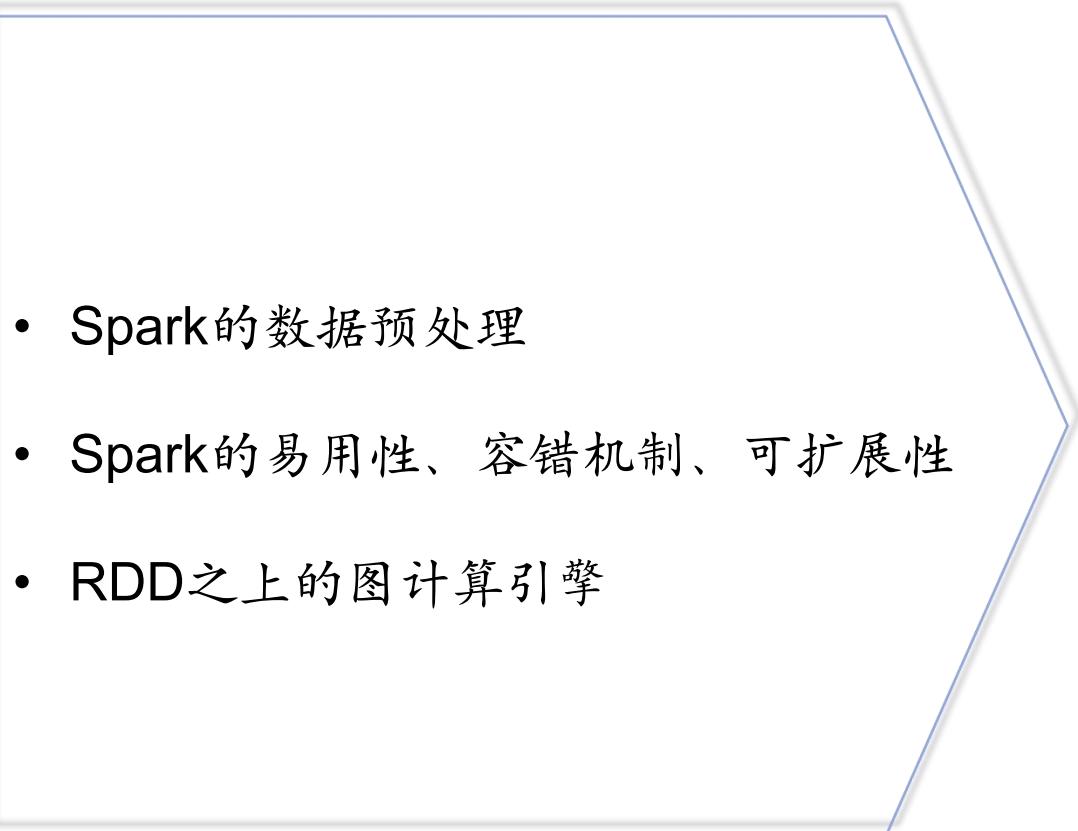
- 统一Pregel、GraphLab计算模型



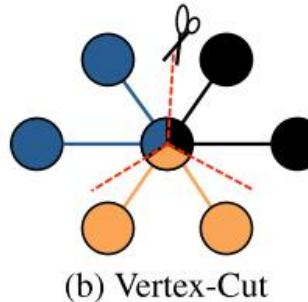
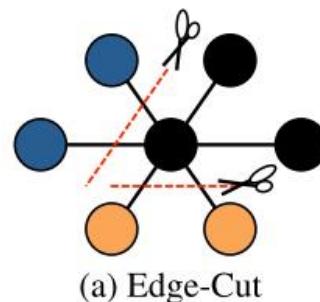


OSDI2014 (UCB AMPlab)

- 1/7 X PowerGraph
- 生产力弥补执行性能



VertexCut



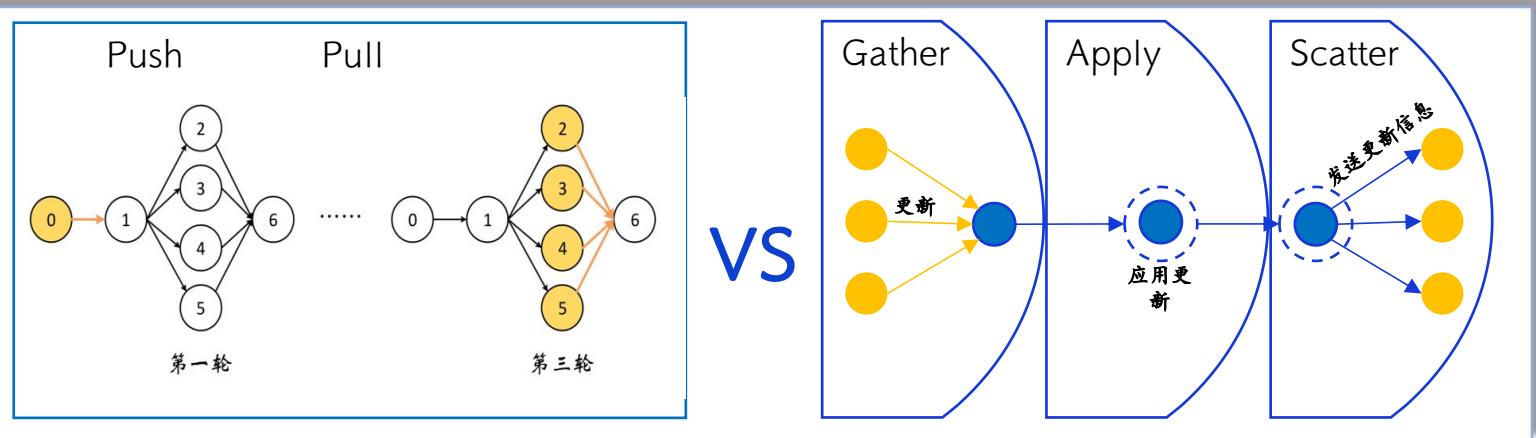


Gemini/Plato

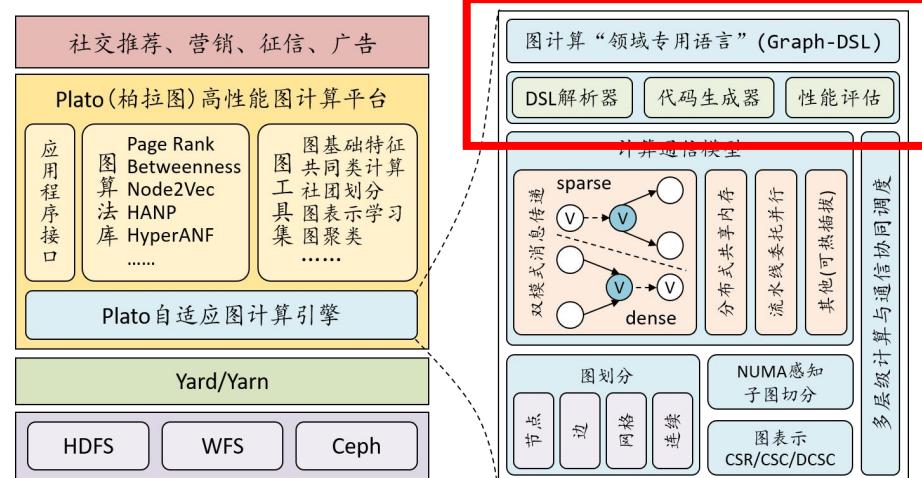
OSDI2016 (THU)

- 5~10X Powergraph

- 双消息传播
- 拉取(Pull)与推送(Push)切换
- 计算为中心



- 柏拉图
- Gemini内核
- DSL语言



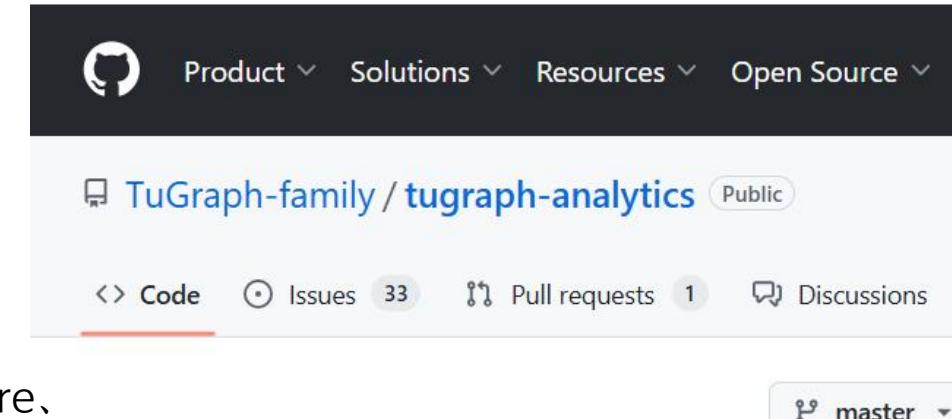
点中心图计算系统

工业界图计算系统大幅加速了相关目标的计算，对异构图支持有待提升

大量图计算算法LPA、FastUnfolding、Betweenness、K-Core、TriangleCounting…等需要在特定业务中起到对应的作用

图计算主要支持同构图，无论在方法论还是系统实现上，对异构图的支持均不足，实际业务场景中更关注异构图

图计算任务同时属于计算密集型与资源密集型，计算代价高昂，需要高效的计算系统，如angel



<https://github.com/TuGraph-family/tugraph-analytics>

图信息获取之图学习

图查询
结构关联

邻居信息、路径信息、子图



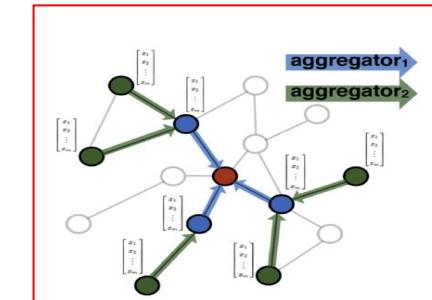
图计算
全局属性

PageRank、K-core



图学习
向量表示

GraphSAGE、LINE



图表示学习

图新视角

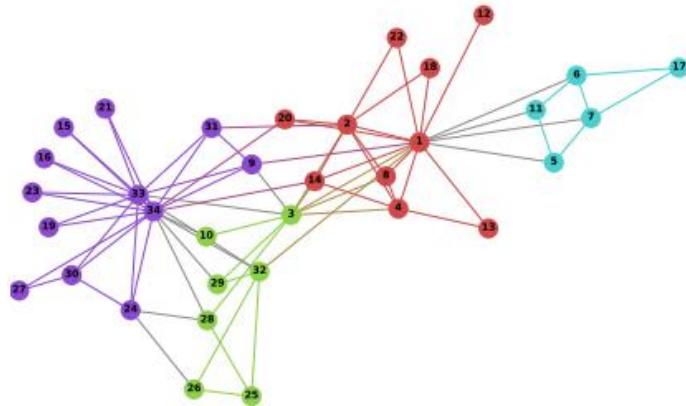
复杂信息的建模、表达

高效、直接的关联查询

图中潜在的重要信息

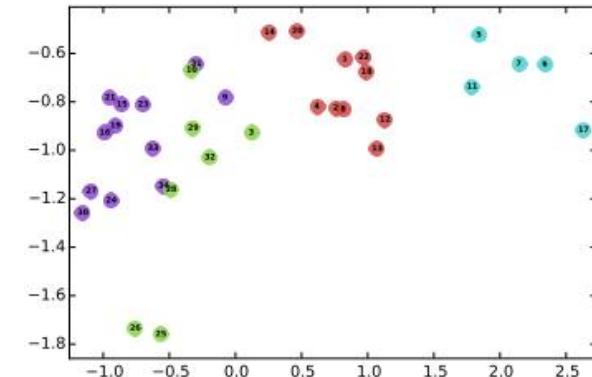
结构化图与深度学习的有机结合

非欧特性



复杂非欧数据向低维数据的降维

向量空间



深层关联关系的动态刻画与度量

“

An embedding maps each node to a **low** dimensional **feature** vector and tries to **preserve** the **connection strengths** between vertices.

”

经验与思考

图查询的关键在于可视化与即时关联分析的高效

图计算的核心作用再全局关联计算中的性能加速

图学习同目前业务需求关系最为紧密，作用最为明显

图的运用应该在遇到业务瓶颈之后

图的产品应该聚焦业务需求、使用体验而非图技术本身

目前的图数据库大都不是合格的数据库，也不被定位为数据库

图很多‘坑’是源自少数工业界的浮躁人群，而非图本身

把Graph的优势融入Hive

TuGraph

图数据库和图计算部分内容共建

图数据库概述：概念、历史、复杂查询及其应用、图数据库系统

图计算概述：概念、常见算法及其应用

图计算系统：基于大数据工具实现、基于点中心编程实现、基于图数据库实现、流式图计算

图计算的工业应用实践（蚂蚁TuGraph团队讲授）：风控、社交等

- 大作业（课程实践，1~2人一组）：同蚂蚁公司TuGraph团队共建，待发布
- 平时作业（仅示例参考，最终发布见后续通知）：
 1. ① TuGraph的部署和试运行GQL；
 2. ② 基于TuGraph图计算实现PageRank，以wiki-Talk为数据集
(<https://snap.stanford.edu/data/wiki-Talk.html>)

Thank you