

大规模异构网络数据融合

张宇韬

清华大学数据科学研究院



QCon

全球软件开发大会

成为软件技术专家的 必经之路

[北京站] 2018

2018年4月20-22日 北京·国际会议中心

7折 购票中, 每张立减2040元
团购享受更多优惠



识别二维码了解更多

SPEAKER INTRODUCE



张宇韬

清华大学 科技大数据中心 首席研究员

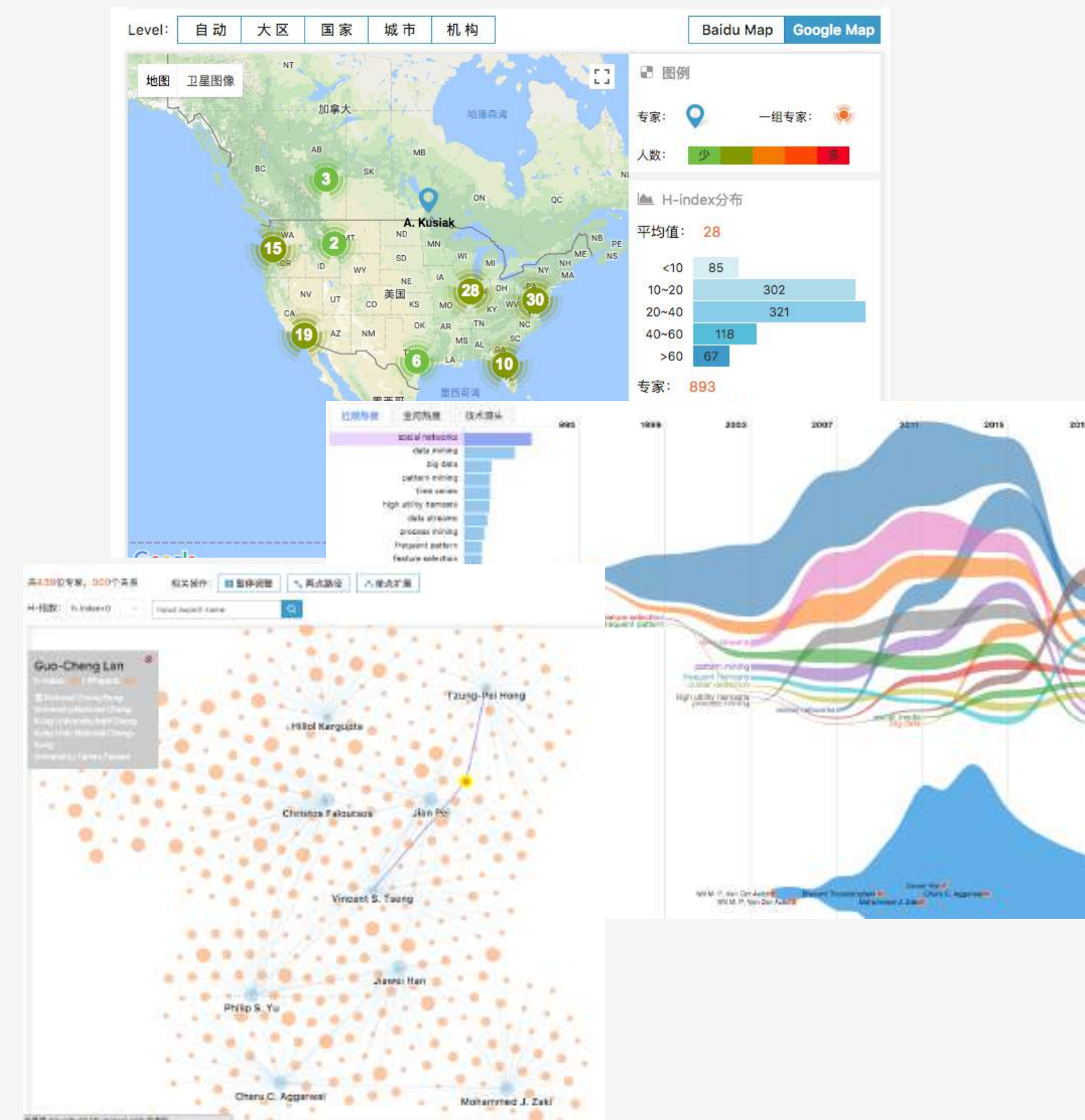
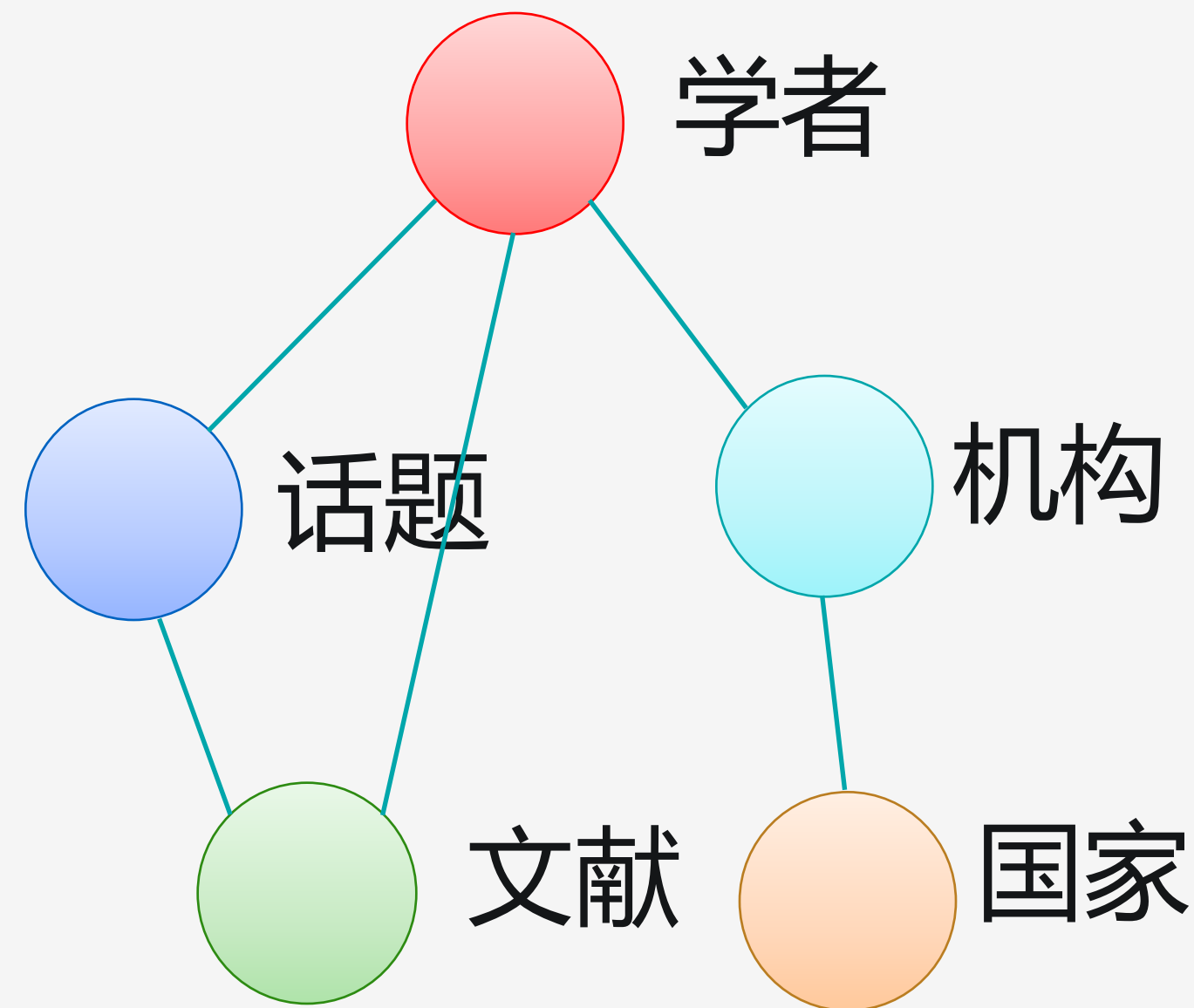
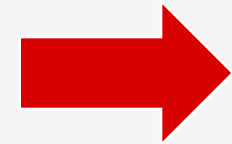
张宇韬，清华大学数据科学研究院科技大数据研究中心首席研究员。研究方向为异构数据融合及知识图谱构建。在KDD、CIKM、VAST等数据挖掘领域国际重要会议上发表多篇论文。作为技术负责人参与研发学术网络分析挖掘系统AMiner，集成上亿的学者、机构、科技文献、专利数据，提供针对科技数据的搜索及可视化分析功能，拥有数百万用户访问量。曾获得吴文俊人工智能科技进步一等奖，IJCAI数据竞赛第二名。

TABLE OF CONTENTS 大纲

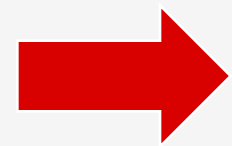
- 异构数据融合与知识图谱构建
- 异构数据融合中的机器学习方法
 - 度量学习
 - 表征学习
- 图嵌入学习与图卷积网络
- 案例：重名实体排歧
- 基于科技知识图谱的智能服务

AMiner科技知识图谱

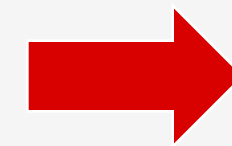
- 文献数据库
- 新闻数据
- 社交网络
- 维基百科
- ...



抽取



融合



挖掘

异构数据融合 (1)

- 相关信息散落在不同的数据源中



Michael I. Jordan
Professor of EECS and Professor of Statistics, [University of California, Berkeley](#)
在 cs.berkeley.edu 的电子邮件经过验证 - 首页
machine learning statistics computational biology artificial intelligence optimization

- 标题**
- Latent dirichlet allocation
DM Blei, AY Ng, MI Jordan
Journal of machine Learning research 3
 - On spectral clustering: Analysis
AY Ng, MI Jordan, Y Weiss
Advances in neural information processing systems 17
 - Adaptive mixtures of local experts
RA Jacobs, MI Jordan, SJ Nowlan, GE Hinton
Neural computation 3 (1), 79-87
 - Sharing clusters among related documents
YW Teh, MI Jordan, MJ Beal, DM Blei
Advances in neural information processing systems 17
 - Hierarchical mixtures of experts
MI Jordan, RA Jacobs
Neural computation 6 (2), 181-214
 - An introduction to variational methods for machine learning
MI Jordan, Z Ghahramani, TS Jaakkola
Machine learning 37 (2), 183-233



Michael I. Jordan
Pehong Chen Distinguished Professor
Department of EECS
Department of Statistics
AMP Lab
Berkeley AI Research Lab
University of California, Berkeley

Emails:
jordan@cs.berkeley.edu
jordan@stat.berkeley.edu

EECS Address:
University of California, Berkeley
EECS Department
387 Soda Hall #1776
Berkeley, CA 94720-1776

Biographical highlights

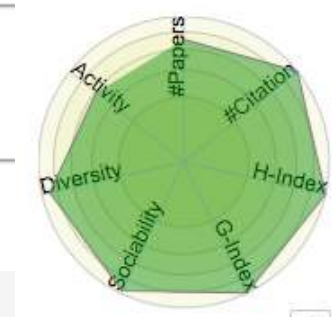
- Professor, University of California, Berkeley
- Professor, MIT, 1988-1991
- Member, National Academy of Sciences
- Member, National Academies of Arts and Sciences
- Fellow, American Association for Artificial Intelligence
- Fellow of the AAAI, ACM, and IEEE
- Elected Member, International Neural Network Society
- Plenary Speaker, International Joint Conference on Artificial Intelligence
- IJCAI Research Excellence Award (2003)
- David E. Rumelhart Prize (2003)
- ACM/AAAI Allen Newell Award (2004)
- SIAM Activity Group on Artificial Intelligence (2004)
- IEEE Neural Networks Award (2004)
- IMS Medallion Lecture (2004)
- Full biography...

Publications

Courses



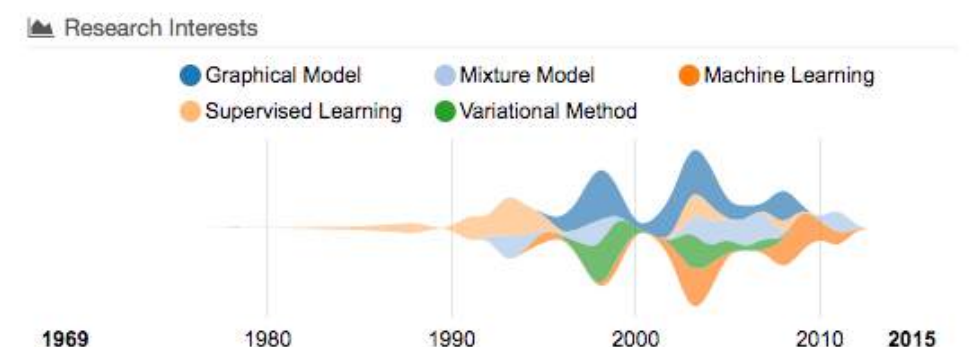
Awards: 



Whatever comes to your mind


Michael I. Jordan  Following | 82

Professor
UC Berkeley Berkeley Engineering
(510) 642-9575
jordan@cs.berkeley.edu, jordan@stat.berkeley.edu
<https://people.eecs.berkeley.edu/~jordan/>



machine learning, statistics, and

Michael I. Jordan



Born February 25, 1956 (age 61)
Residence Berkeley, CA
Alma mater University of California, San Diego
Known for Latent Dirichlet allocation
Fellow of the U.S. National Academy of Sciences^[1]
AAAI Fellow (2002)
Rumelhart Prize (2015) [2]
IJCAI Award for Research Excellence (2016)

异构数据融合 (1)

- 相关信息散落在不同的数据源中

Michael I. Jordan
Professor of EECS and Professor of Statistics, [University of California, Berkeley](http://www.eecs.berkeley.edu)
在 cs.berkeley.edu 的电子邮件经过验证 - 首页
machine learning statistics computational biology artificial intelligence optimization

标题

- Latent dirichlet allocation
DM Blei, AY Ng, MI Jordan
Journal of machine Learning research 3
- On spectral clustering: Analysis
AY Ng, MI Jordan, Y Weiss
Advances in neural information processing systems 14
- Adaptive mixtures of local experts
RA Jacobs, MI Jordan, SJ Nowlan, GE Hinton
Neural computation 3 (1), 79-87
- Sharing clusters among related documents
YW Teh, MI Jordan, MJ Beal, DM Blei
Advances in neural information processing systems 14
- Hierarchical mixtures of experts
MI Jordan, RA Jacobs
Neural computation 6 (2), 181-214
- An introduction to variational methods for machine learning
MI Jordan, Z Ghahramani, TS Jaakkola
Machine learning 37 (2), 183-233

Michael I. Jordan
Pehong Chen Distinguished Professor
Department of EECS
Department of Statistics
AMP Lab
Berkeley AI Research Lab
University of California, Berkeley

Emails:
jordan@cs.berkeley.edu
jordan@stat.berkeley.edu

EECS Address:
University of California, Berkeley
EECS Department
387 Soda Hall #1776
Berkeley, CA 94720-1776

Biographical highlights

- Professor, University of California, Berkeley, 1985-1988
- Professor, MIT, 1988-1991
- Member, National Academy of Sciences, 1991
- Member, National Academies of Arts and Sciences, 1992
- Fellow, American Association for Artificial Intelligence, 1992
- Fellow of the AAAI, ACM, and IEEE
- Elected Member, International Joint Conference on Artificial Intelligence (IJCAI) Research Excellence Award, 1997
- David E. Rumelhart Prize, 1998
- ACM/AAAI Allen Newell Award, 1999
- SIAM Activity Group on Artificial Intelligence, 1999
- IEEE Neural Networks Award, 2000
- IMS Medallion Lecture, 2000
- Full biography...

Publications

Courses

Michael Jordan
Machine learning, statistics, and optimization
Microsoft Research

Michael I. Jordan
Born February 25, 1956 (age 61)
Residence Berkeley, CA
Alma mater University of California, San Diego
Latent Dirichlet allocation
Fellow of the U.S. National Academy of Sciences^[1]
AAAI Fellow (2002)
Rumelhart Prize (2015) ^[2]
IJCAI Award for Research Excellence (2016)

Michael I. Jordan
Professor
UC Berkeley Berkeley Engineering
(510) 642-9575
jordan@cs.berkeley.edu, jordan@stat.berkeley.edu
<https://people.eecs.berkeley.edu/~jordan/>

Research Interests

- Graphical Model
- Mixture Model
- Machine Learning
- Supervised Learning
- Variational Method

Similar Authors

Ego Network

Awards:

Activity

- #Papers
- #Citation
- H-Index
- G-Index
- Diversity
- Sociability

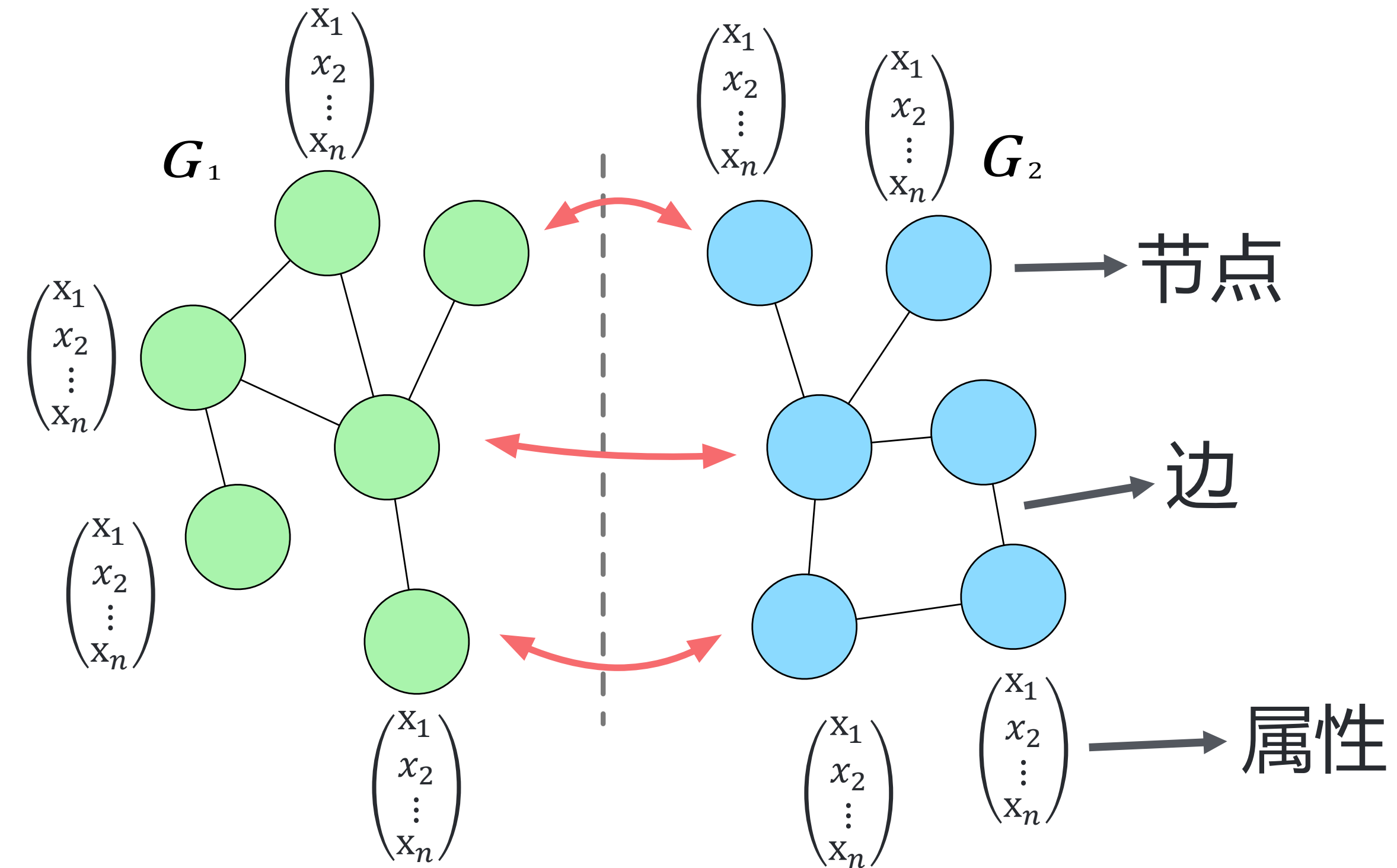
异构数据融合 (2)

- 不同数据源提供了对于同一实体不同视角的描述
- 如何回答：
 - 在NIPS上发表过一作论的应届博士生？
 - 中国2017年获得发明专利最多的机构是那些？
 - 机器学习领域引用数最多的女性学者是谁？

异构数据融合 (3)

- 不同数据源间缺乏对于同一实体的统一识别
 - 自然语言的二义性
 - 同名异义(hyponym)
 - 异名同义(synonym)

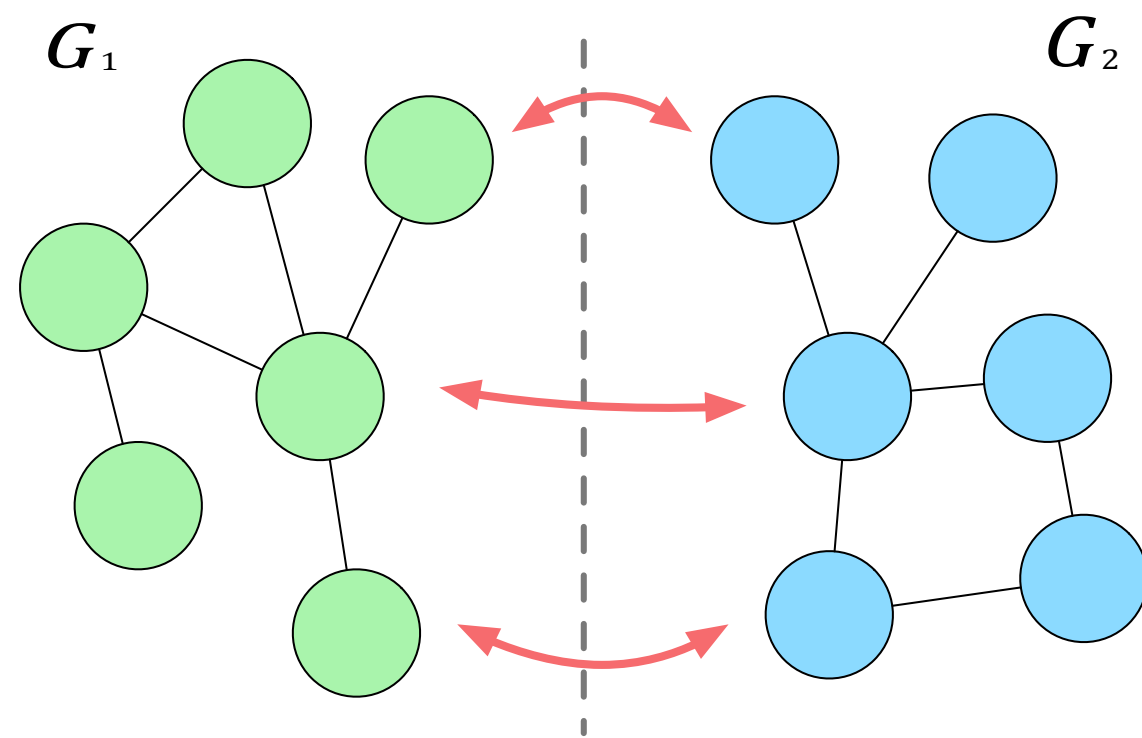
网络数据融合



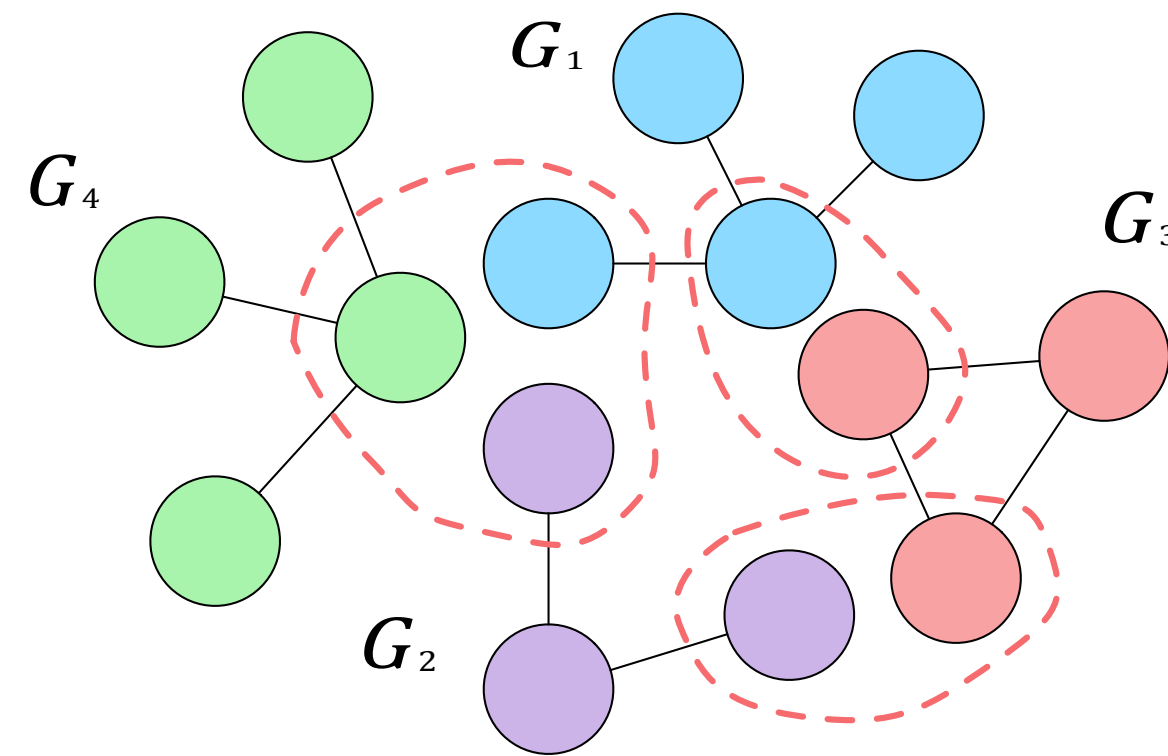
- 社交账户匹配、蛋白质网络(PPI)对齐、重名实体排歧、知识库实体链接、关系数据库集成、视觉图像匹配.....

网络数据融合的场景

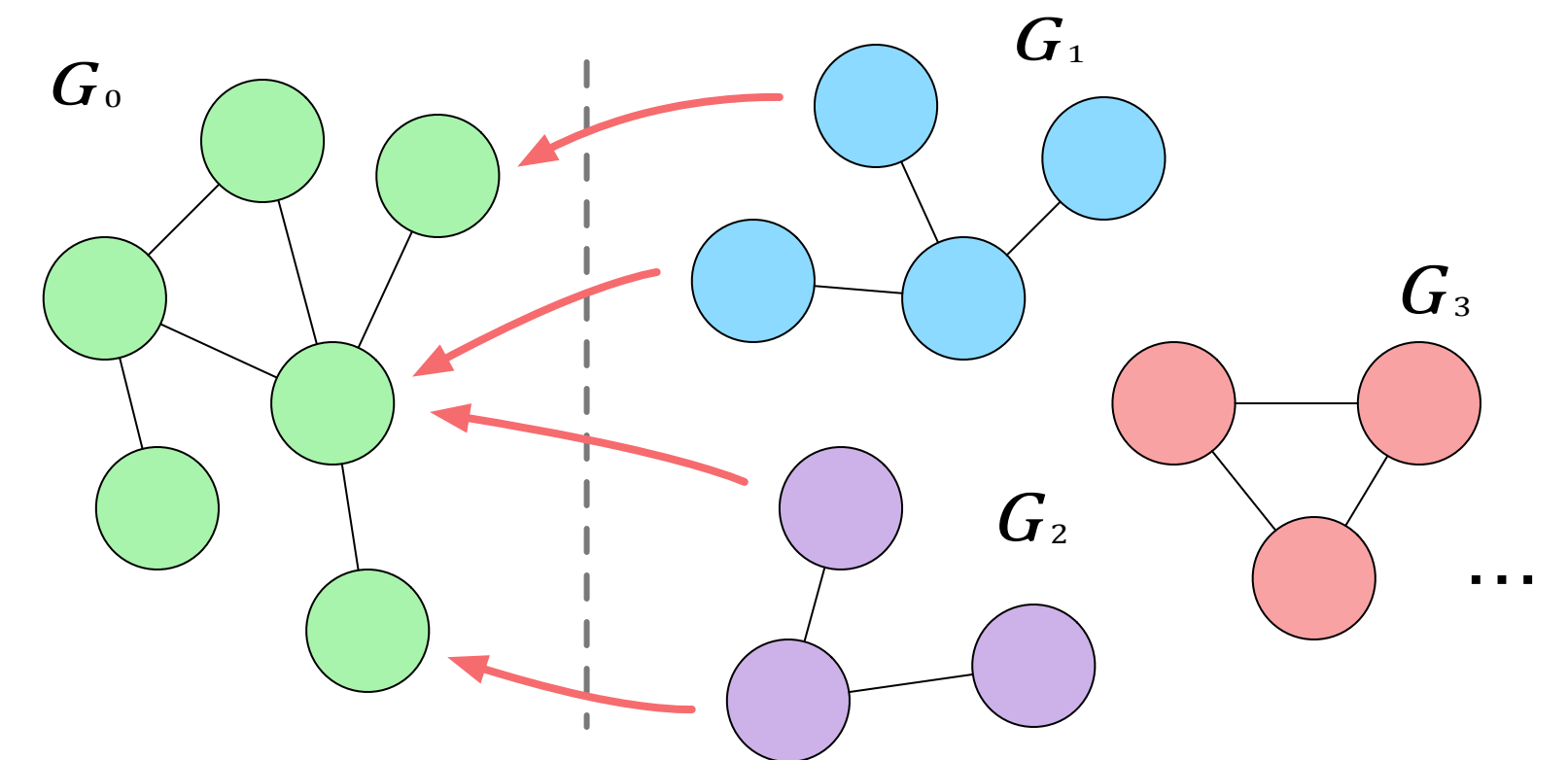
匹配问题 (Matching)



聚类问题 (Clustering)



分配问题 (Assignment)



核心问题：度量节点间的相似度

- 传统方法：定义相似度度量算法
 - Euclidean Distance, Jaccard Similarity, Edit Distance...
- 机器学习方法
 - 度量学习 (Metric Learning)
 - 表征学习 (Representation Learning)

度量学习 (Metric Learning)

- 输入：节点配对 (x_i^1, x_i^2) , 相似度 $y_i \in R$ 或标注 $y_i \in \{0,1\}$
- 目标：学习一个函数 $f(\cdot, \cdot)$, 使得 $f(x_i^1, x_i^2) \sim y_i$
- 马氏距离(Mahalanobis Distances):

$$d_M(x_1, x_2) = \sqrt{(x_1 - x_2)^T M (x_1 - x_2)}$$

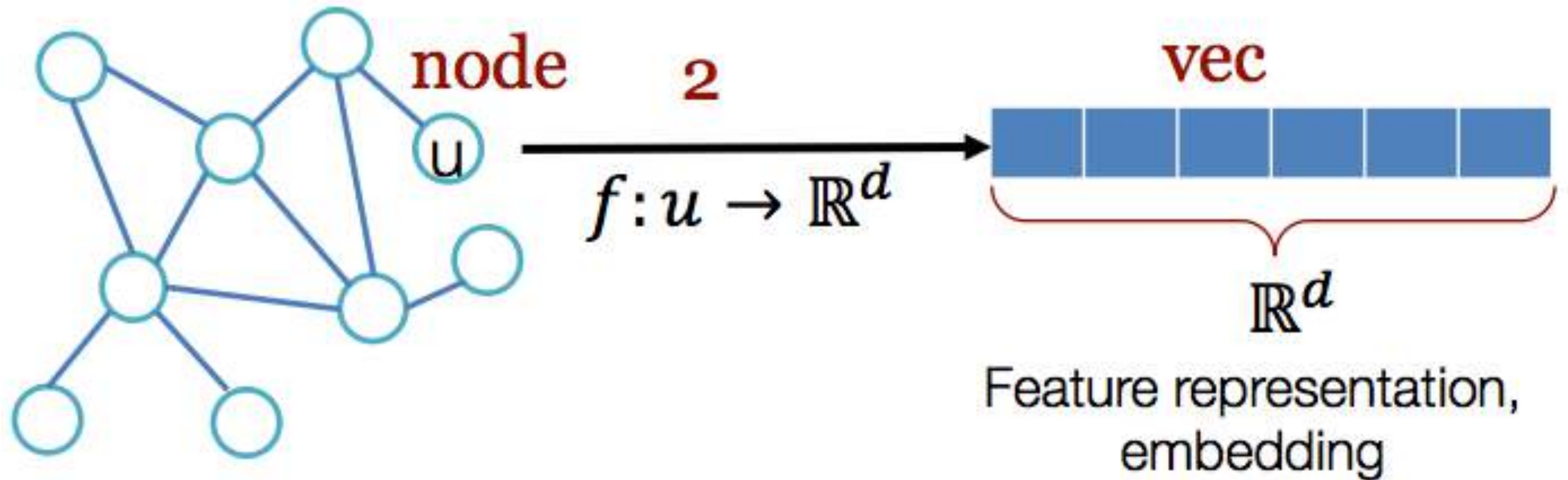
表征学习 (Representation Learning)

- Embedding : $f: X \rightarrow Y$ 即一个映射函数，将数据从原特征空间 X 投影到表征空间 Y
 - f 为单射函数，即任意 $x \in X$ 有唯一的 $y \in Y$ 与之对应
 - Structure preserving性质，例如
$$d(x_1, x_2) < d(x_1, x_3) \rightarrow d(y_1, y_2) < d(y_1, y_3)$$

表征学习在实际应用中的优势

- 度量学习：
 - 直接预测两点间的相似度
- 表征学习
 - 将节点投影到低维连续空间，可容易计算相似度
 - Embedding可作为实体的压缩表示
 - Embedding可作为特征被用于其他数据挖掘任务

网络嵌入学习 (Network Embedding)



图嵌入的难点

- 文本是线性序列 (word2vec, RNN)
- 图片是固定大小的二维网格 (CNN)
- 网络：
 - 节点间的连接是随机的
 - 有着复杂的子结构

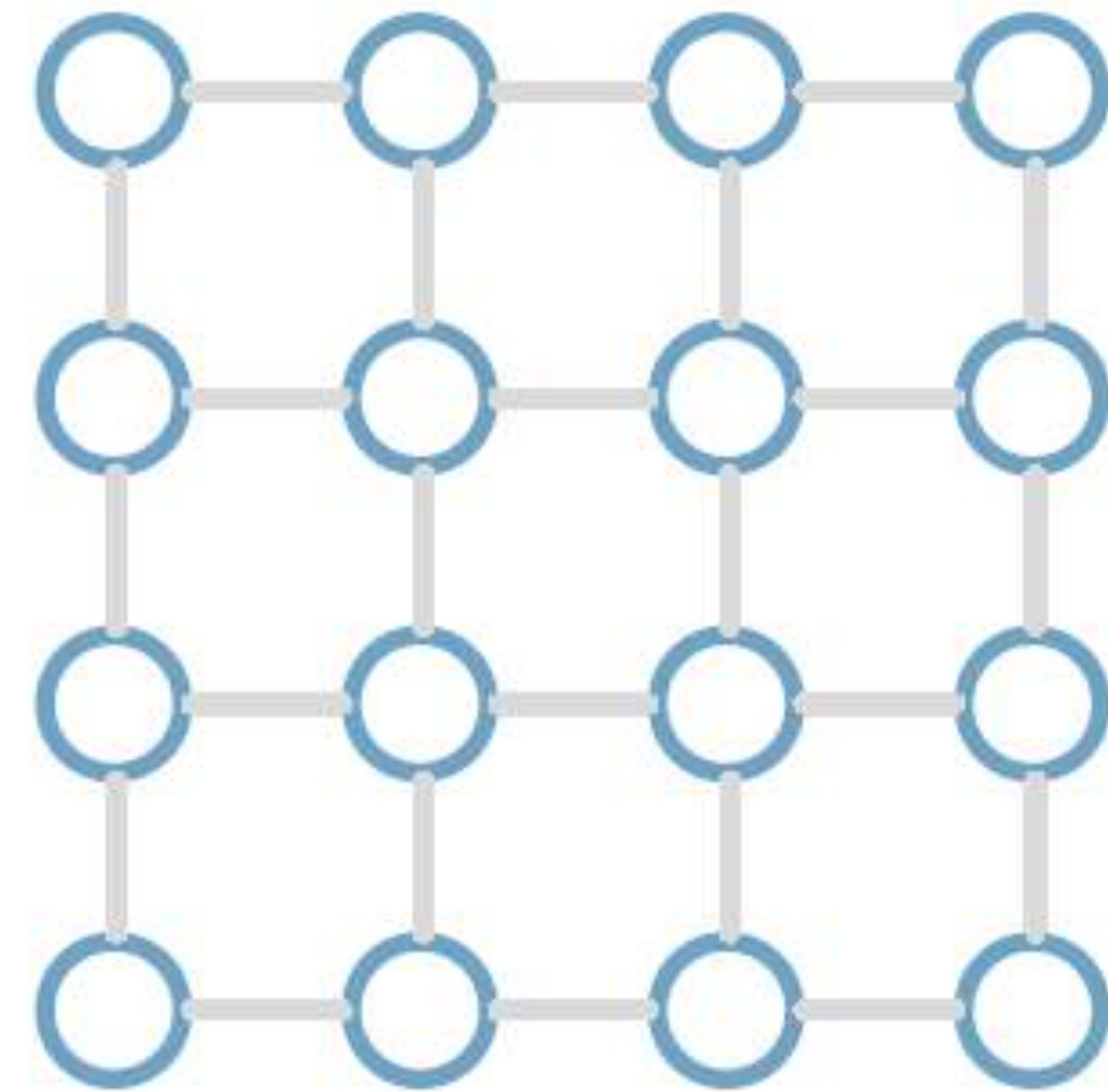
图嵌入的难点

- 文本是线性序列 (word2vec, RNN)
- 图片是固定大小的二维网格 (CNN)
- 网络：
 - 节点间的连接是随机的
 - 有着复杂的子结构



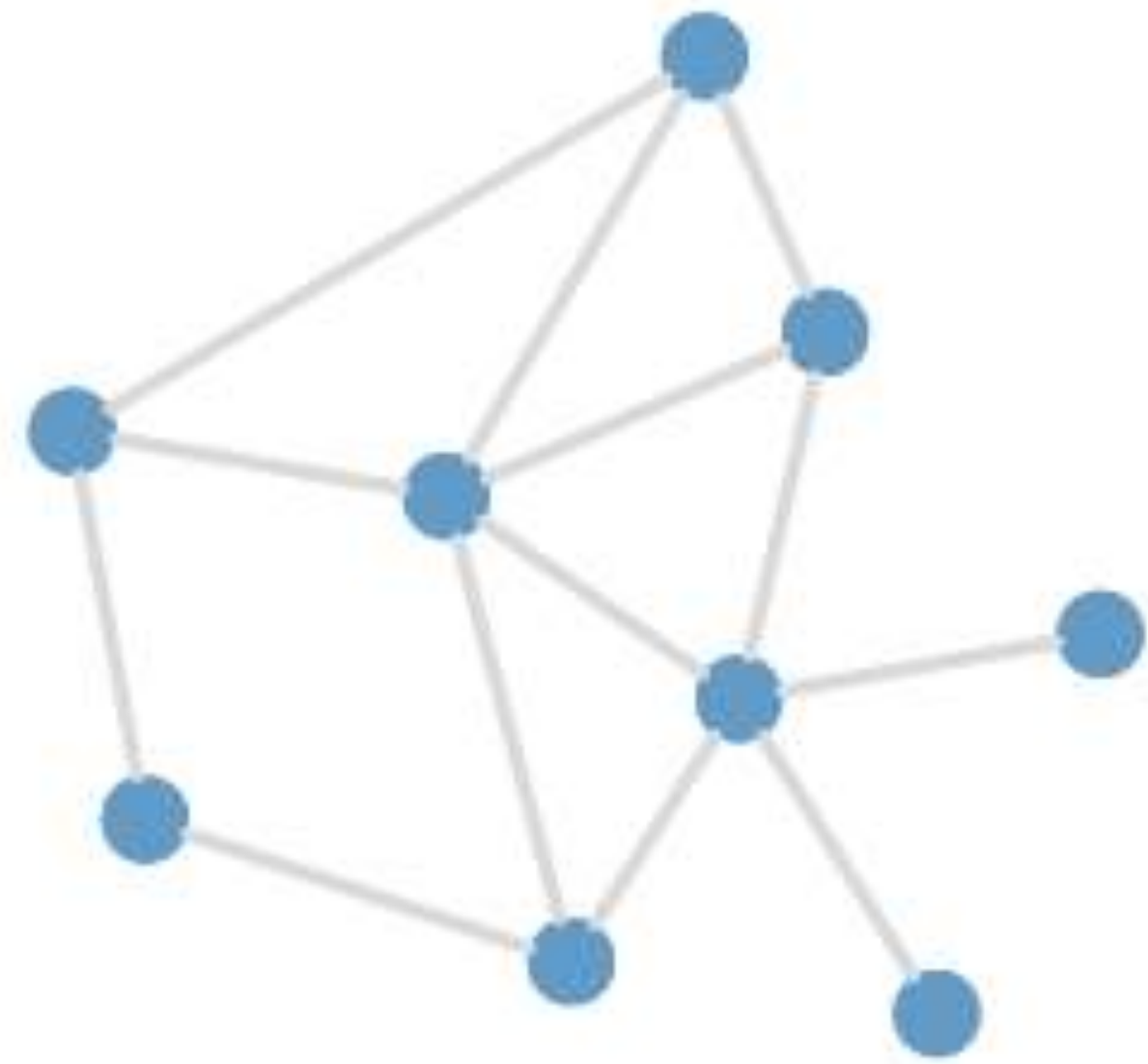
图嵌入的难点

- 文本是线性序列 (word2vec, RNN)
- 图片是固定大小的二维网格 (CNN)
- 网络：
 - 节点间的连接是随机的
 - 有着复杂的子结构



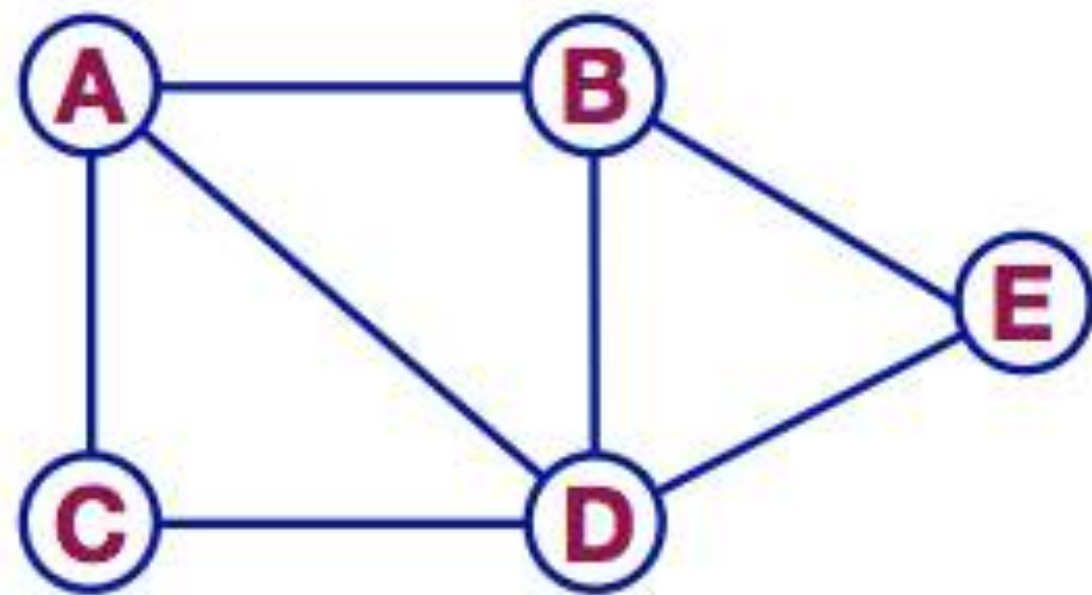
图嵌入的难点

- 文本是线性序列 (word2vec, RNN)
- 图片是固定大小的二维网格 (CNN)
- 网络：
 - 节点间的连接是随机的
 - 有着复杂的子结构



图嵌入的难点

Graph: $G = (\mathcal{V}, \mathcal{E})$



邻接矩阵 A

特征向量 X

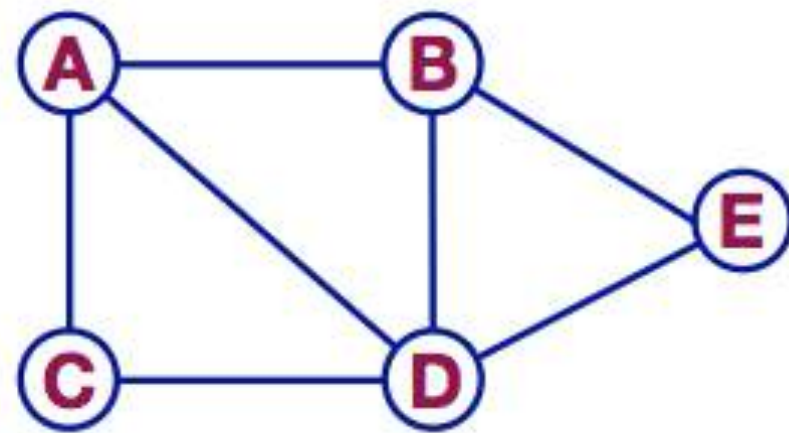
	A	B	C	D	E	Feat	
A	0	1	1	1	0	1	0
B	1	0	0	1	1	0	0
C	1	0	0	1	0	0	1
D	1	1	1	0	1	1	1
E	0	1	0	1	0	1	0

[Thomas Kipf]

图嵌入的难点

- 将邻接矩阵和特征向量直接作为输入构建三层神经网络

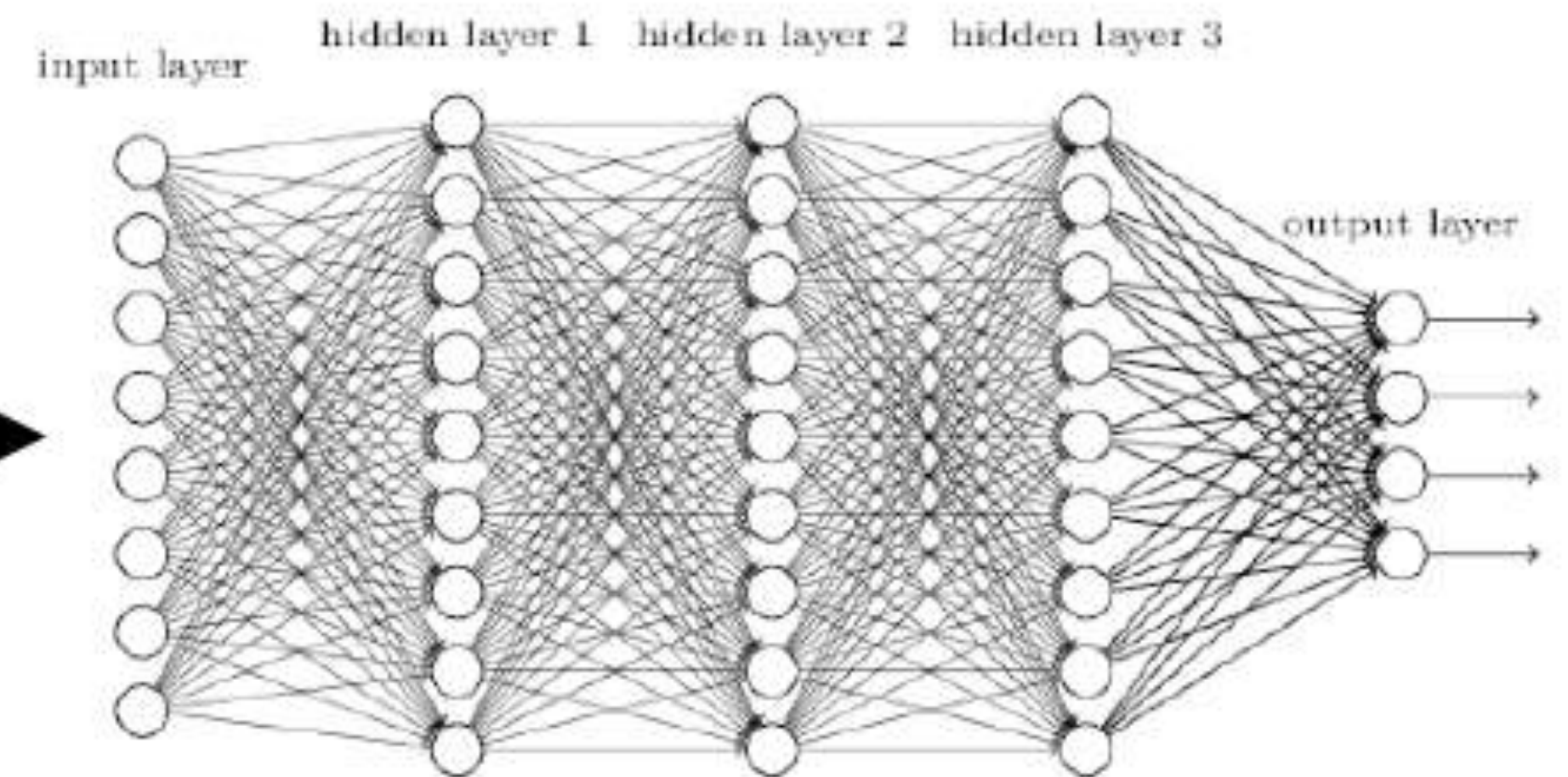
Graph: $G = (\mathcal{V}, \mathcal{E})$



邻接矩阵 A

特征向量 X

	A	B	C	D	E	Feat	
A	0	1	1	1	0	1	0
B	1	0	0	1	1	0	0
C	1	0	0	1	0	0	1
D	1	1	1	0	1	1	1
E	0	1	0	1	0	1	0

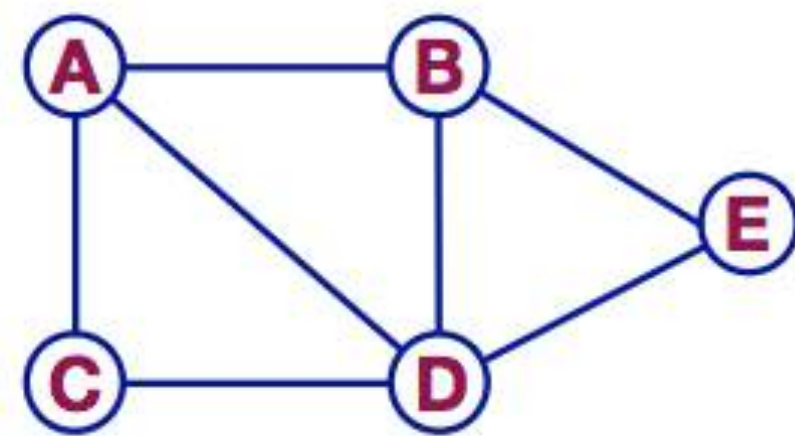


[Thomas Kipf]

图嵌入的难点

- 将邻接矩阵和特征向量直接作为输入构建三层神经网络

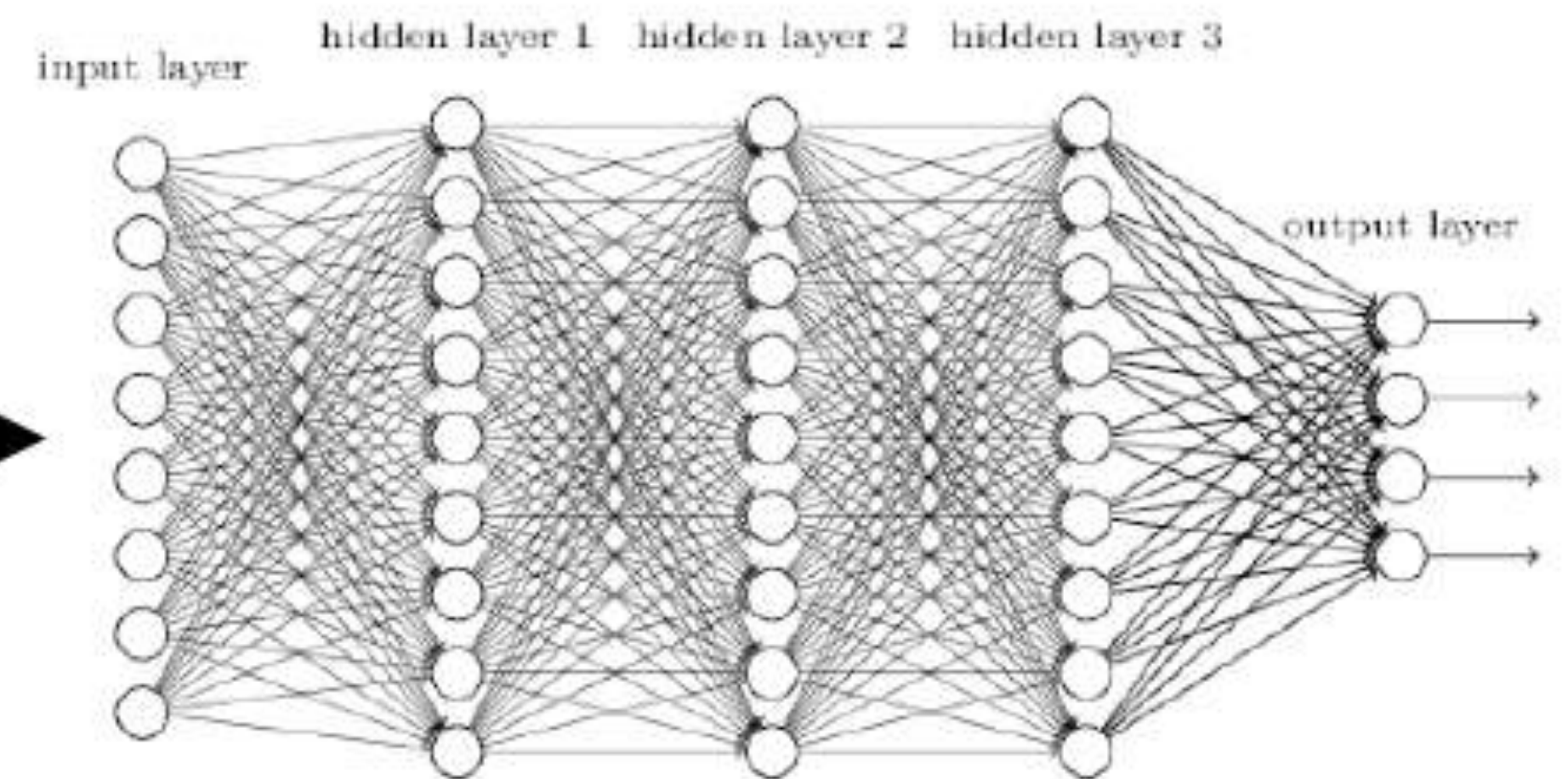
Graph: $G = (\mathcal{V}, \mathcal{E})$



邻接矩阵 A

特征向量 X

	A	B	C	D	E	Feat	
A	0	1	1	1	0	1	0
B	1	0	0	1	1	0	0
C	1	0	0	1	0	0	1
D	1	1	1	0	1	1	1
E	0	1	0	1	0	1	0



问题：

- 参数数量巨大
- 当图结构变化时需要重新训练

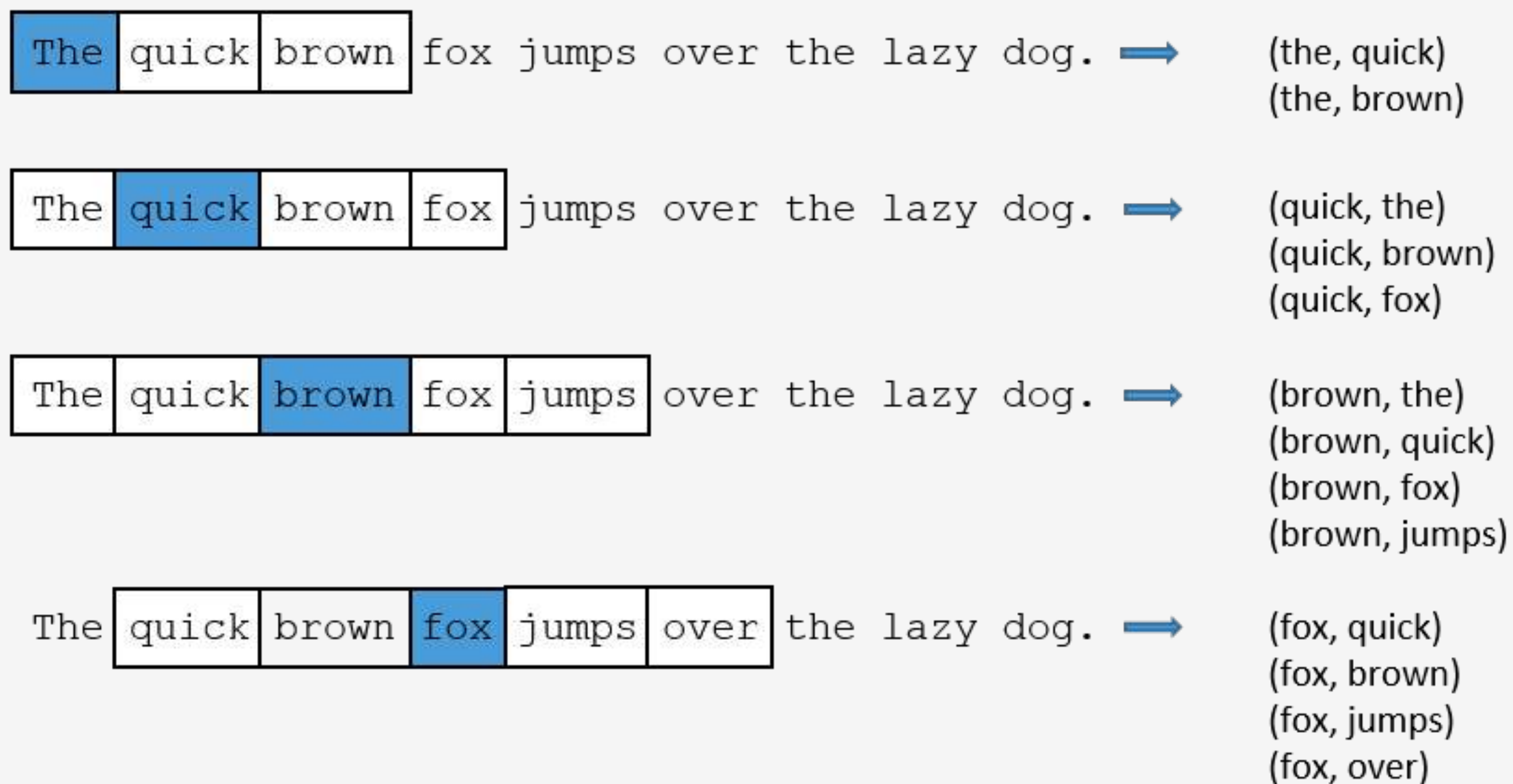
[Thomas Kipf]

常用的图嵌入方法

- 矩阵分解 (Unsupervised, Transductive, Pure Structure)
 - 谱方法
- 图的“线性化” (Unsupervised, Transductive, Pure Structure)
 - 使用RandomWalk等方法将图转化为类似句子的点序列
 - DeepWalk, node2vec, LINE
- 图卷积网络 (Semi-Supervised, Inductive, Structure & Attributes)

图的“线性化”

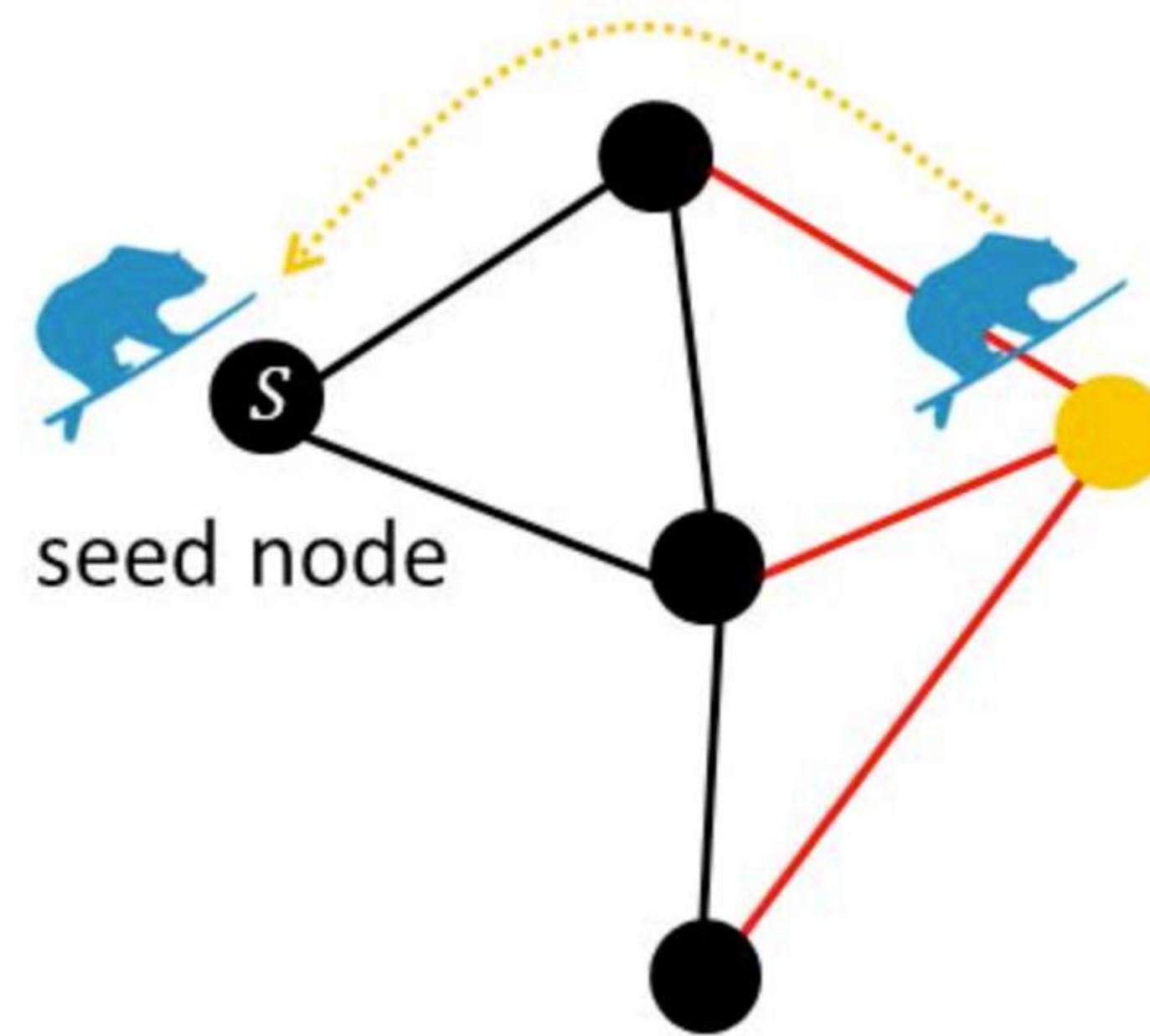
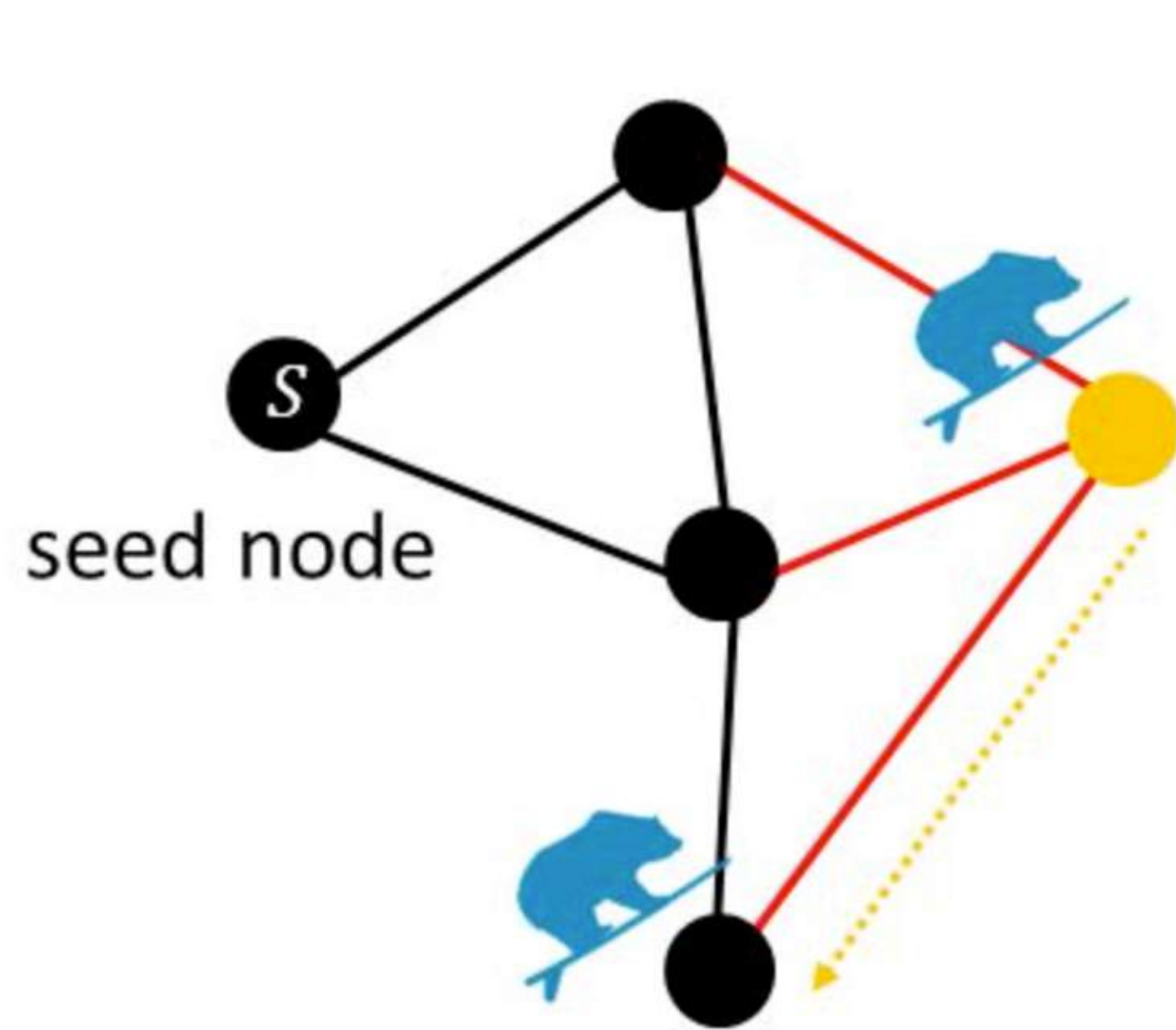
- 词向量模型 Word2vec



<http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/>

图的“线性化”

- 根据随机游走（以概率 p 跳回源点）采样节点序列



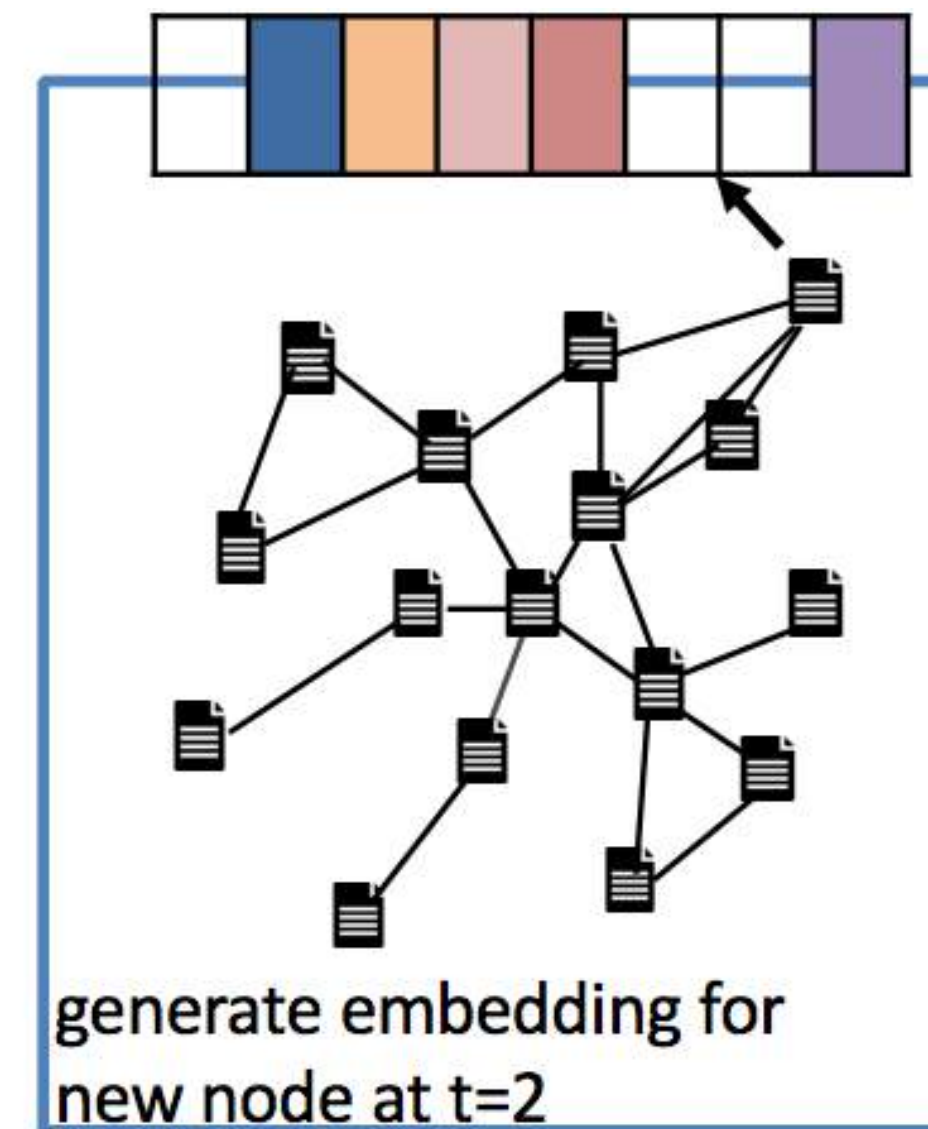
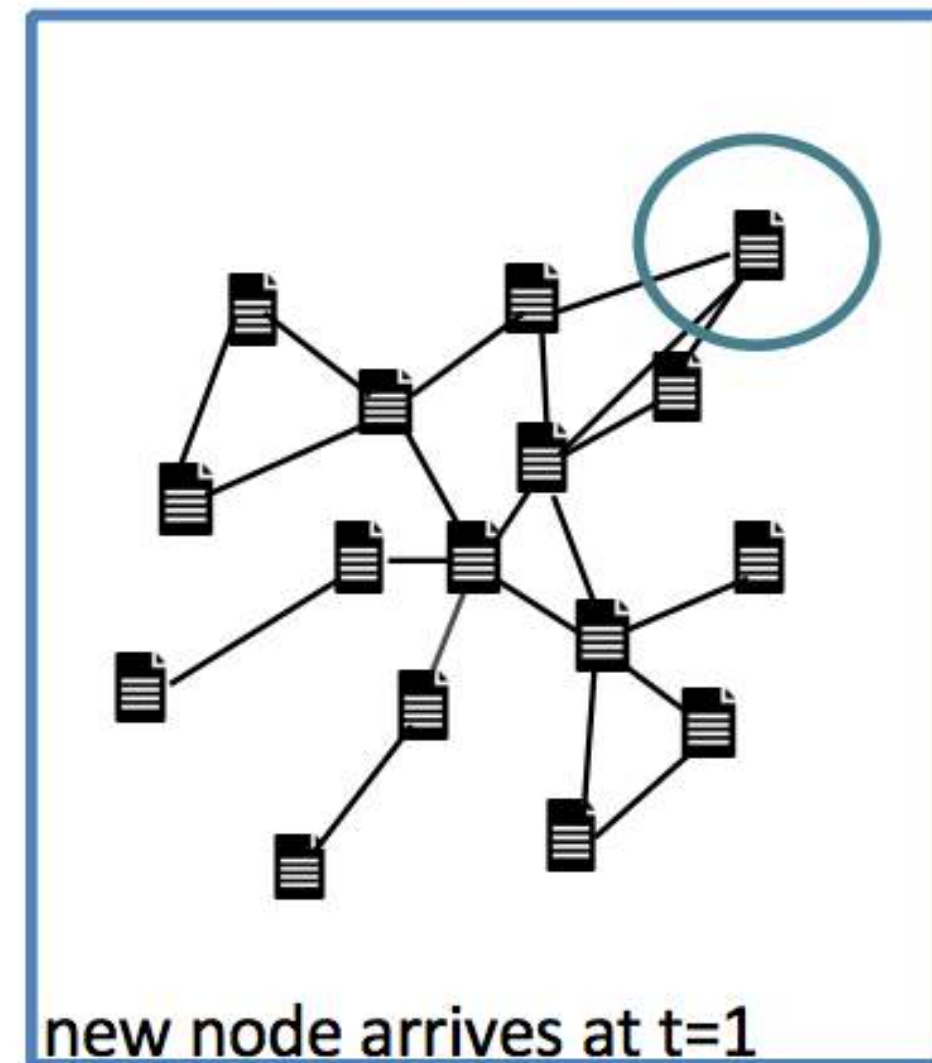
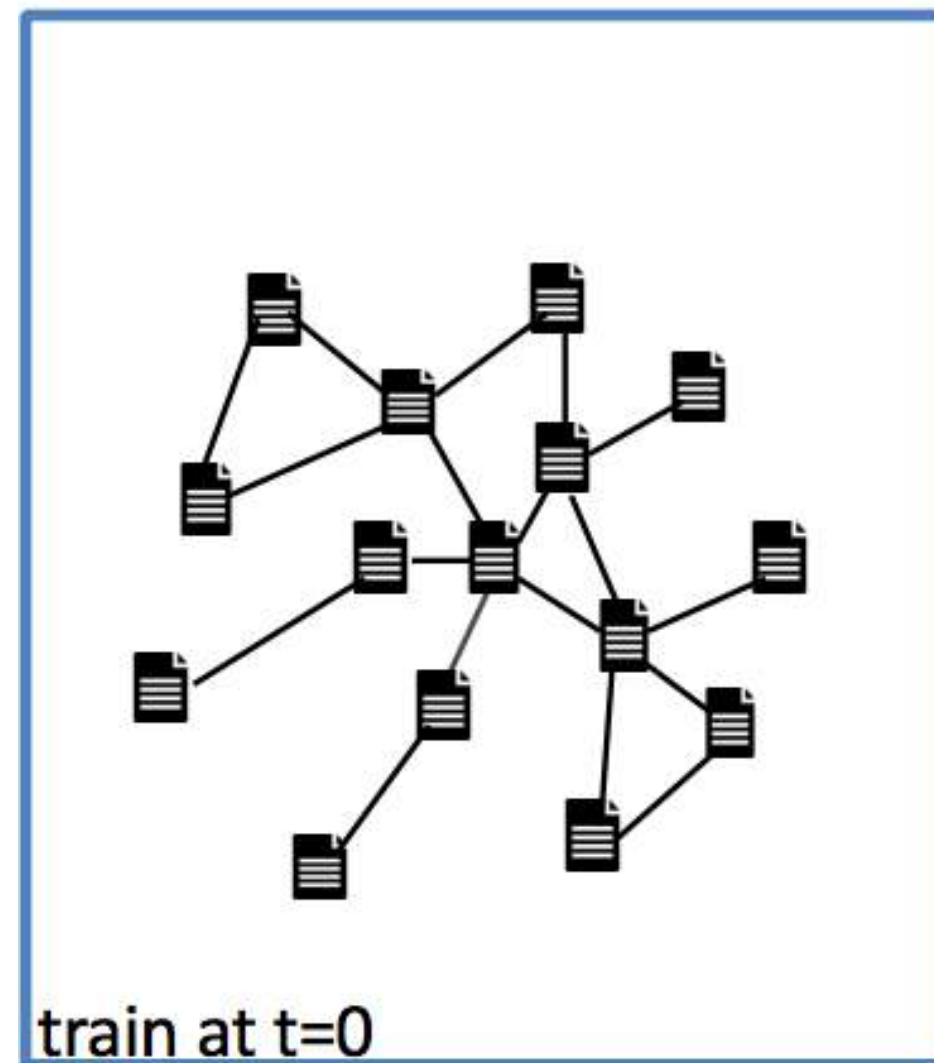
Transductive VS Inductive

- Transductive Learning:

当网络变化时需要重新训练

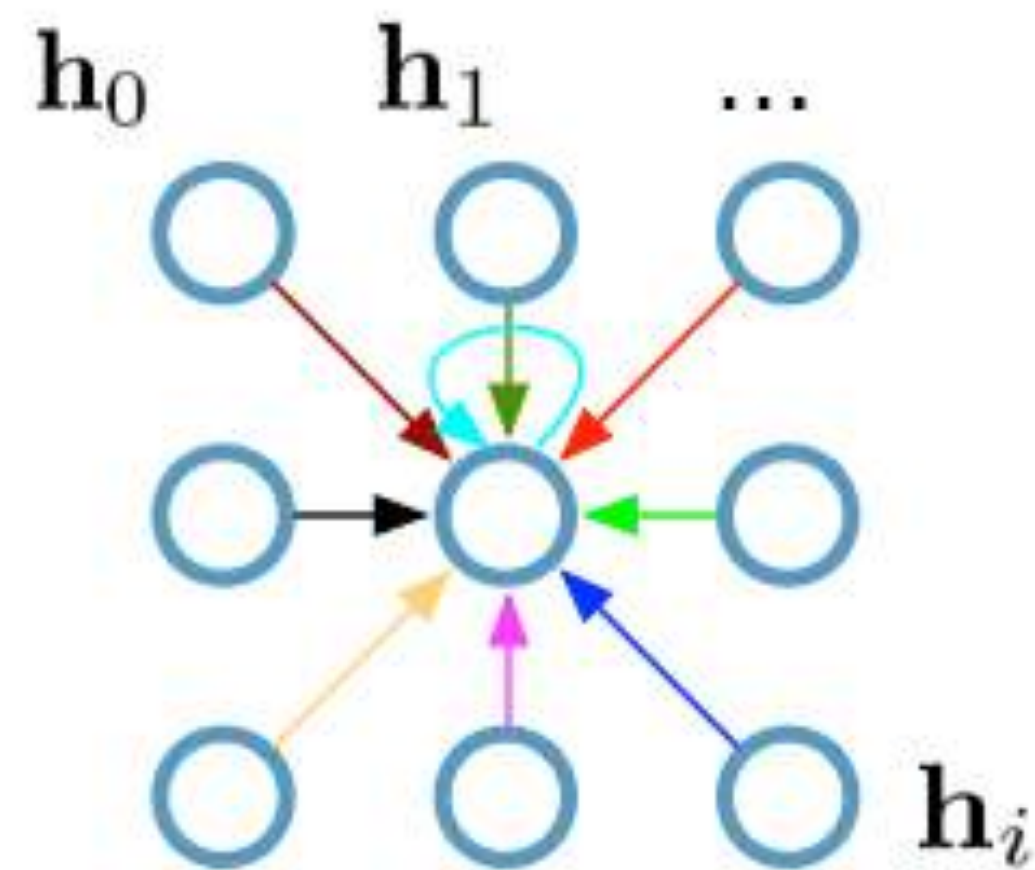
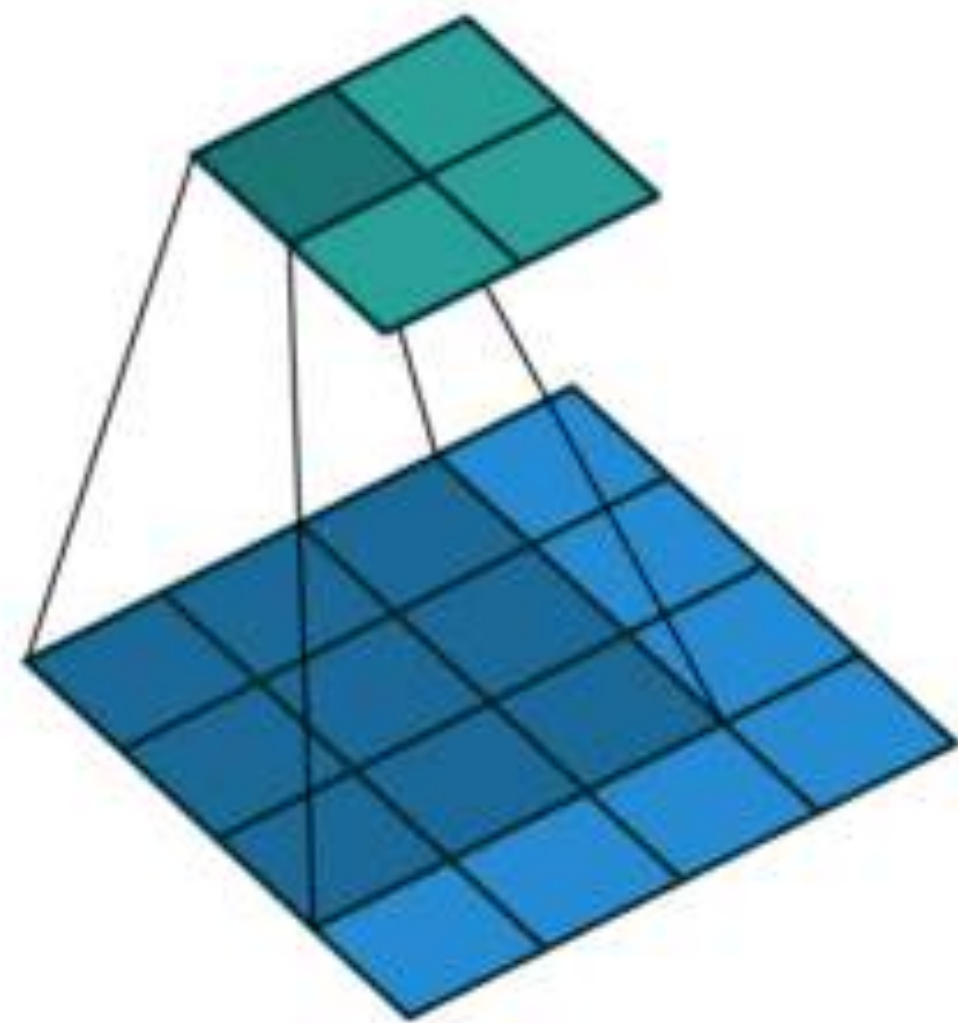
- Inductive Learning:

可预测训练时未见过的新节点的embedding



二维网格上的卷积网络

- 3x3的卷积神经网络

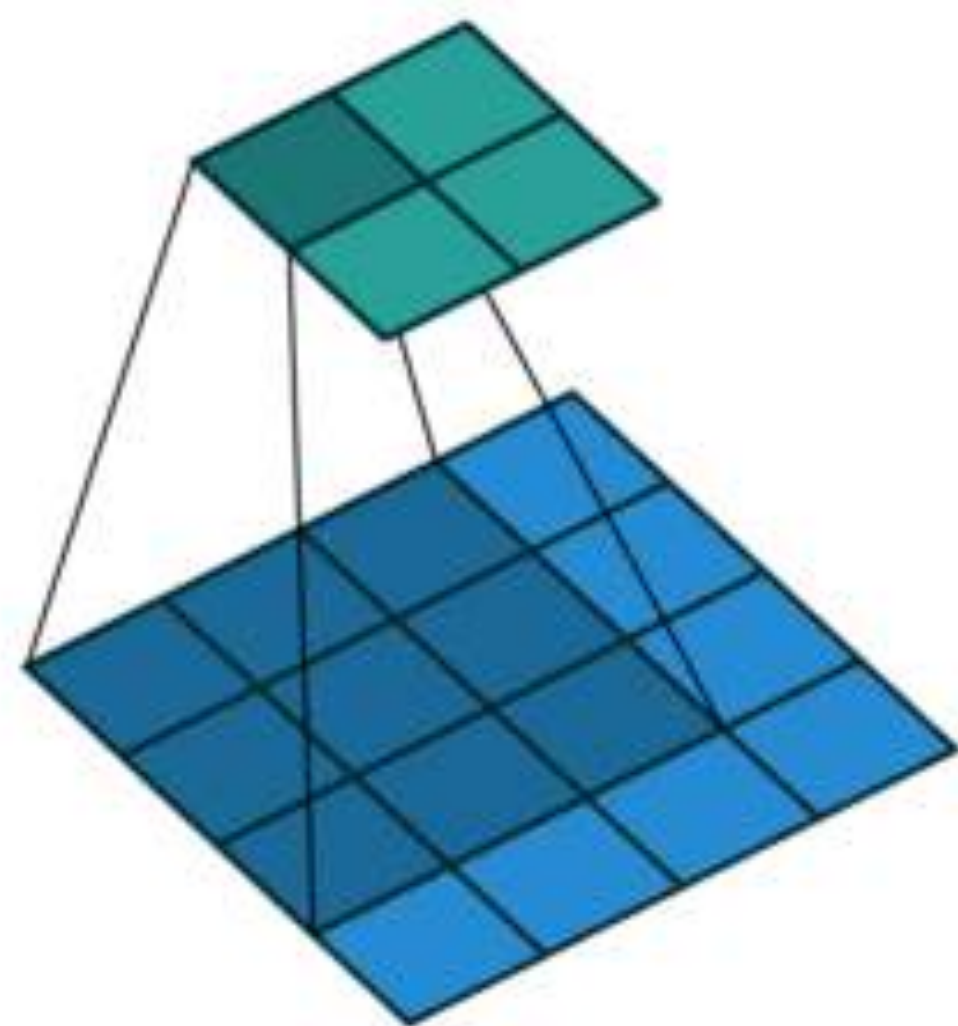


感受域
Receptive field

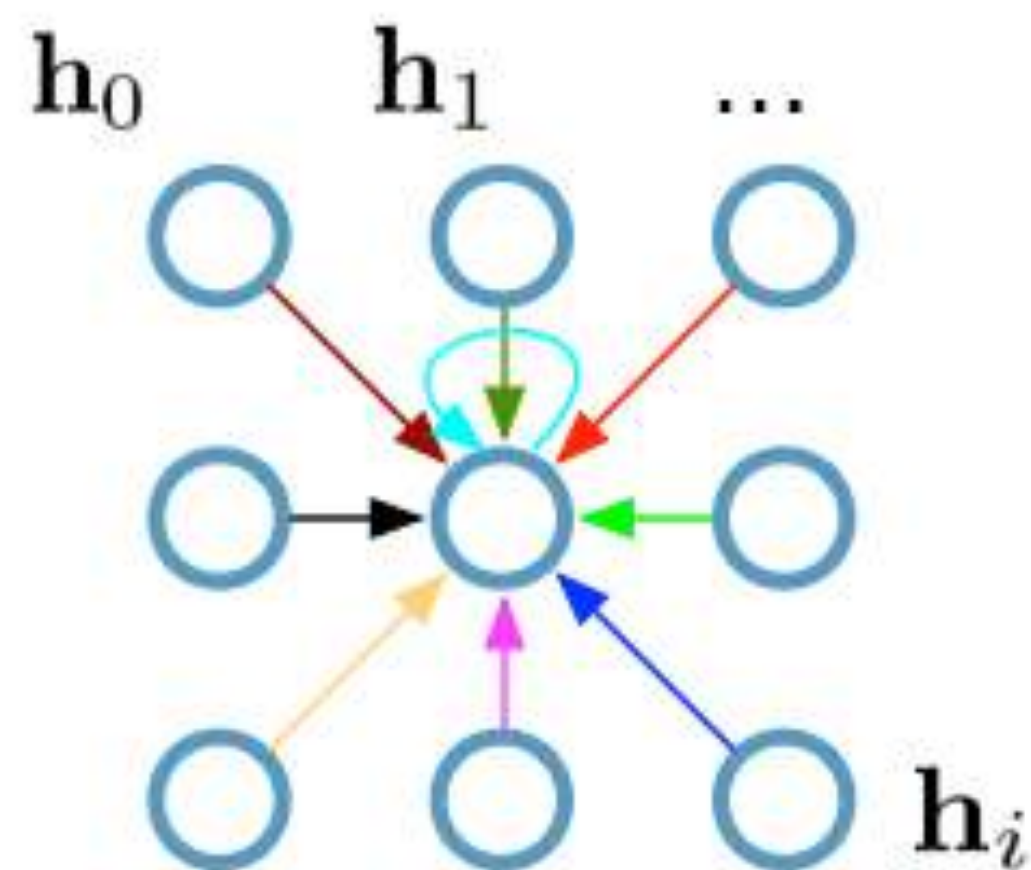
[Thomas Kipf]

二维网格上的卷积网络

- 3x3的卷积神经网络



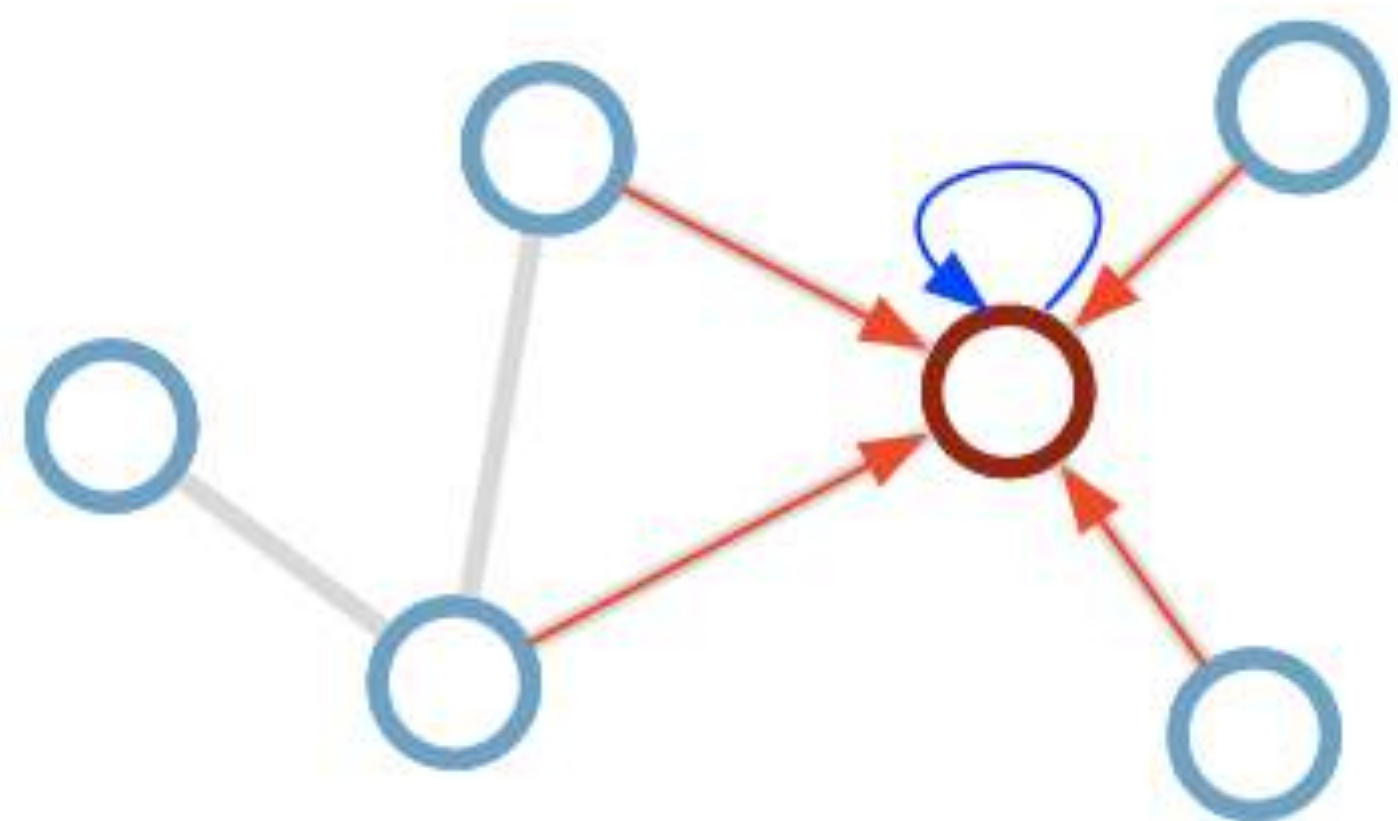
图上的消息传递



$$\mathbf{h}_4^{(l+1)} = \sigma \left(\mathbf{W}_0^{(l)} \mathbf{h}_0^{(l)} + \mathbf{W}_1^{(l)} \mathbf{h}_1^{(l)} + \dots + \mathbf{W}_8^{(l)} \mathbf{h}_8^{(l)} \right)$$

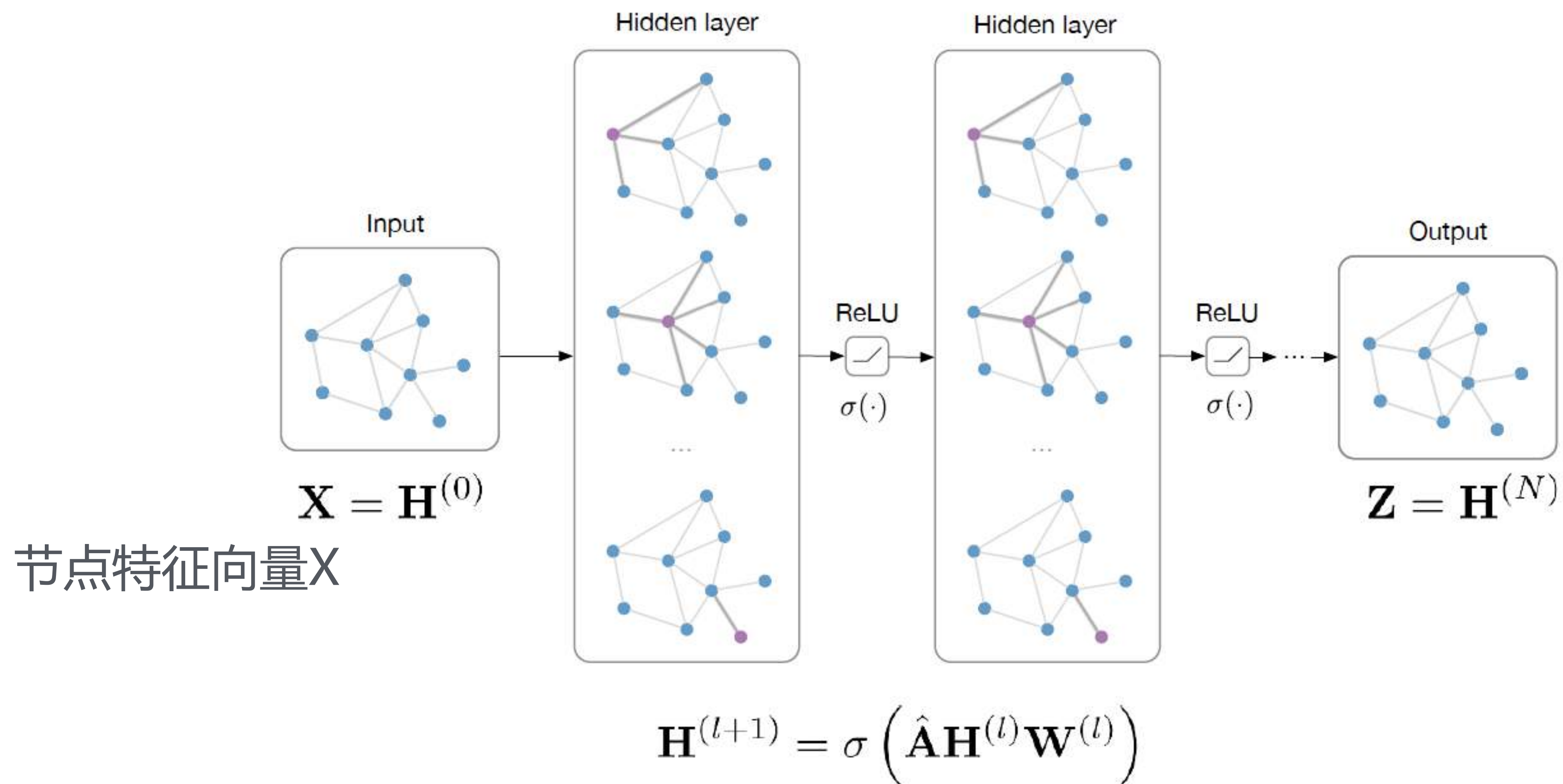
图卷积网络

- 拓展到广义的图上，定义图上的卷积操作



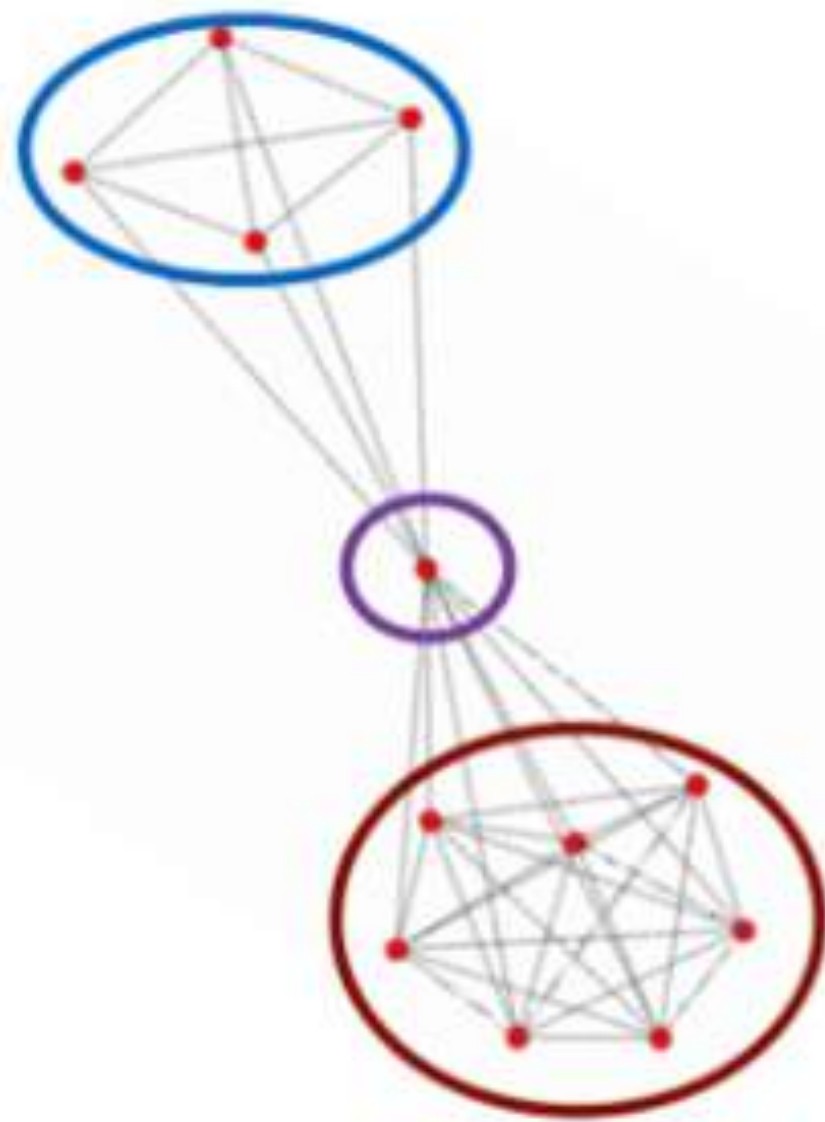
$$\mathbf{h}_i^{(l+1)} = \sigma \left(\mathbf{h}_i^{(l)} \mathbf{W}_0^{(l)} + \sum_{j \in \mathcal{N}_i} \frac{1}{c_{ij}} \mathbf{h}_j^{(l)} \mathbf{W}_1^{(l)} \right)$$

图卷积网络



与图同构问题的联系

- Weisfeiler-Lehman算法



Algorithm 1: WL-1 algorithm (Weisfeiler & Lehmann, 1968)

Input: Initial node coloring $(h_1^{(0)}, h_2^{(0)}, \dots, h_N^{(0)})$

Output: Final node coloring $(h_1^{(T)}, h_2^{(T)}, \dots, h_N^{(T)})$

$t \leftarrow 0$;

repeat

for $v_i \in \mathcal{V}$ **do**

$h_i^{(t+1)} \leftarrow \text{hash} \left(\sum_{j \in \mathcal{N}_i} h_j^{(t)} \right)$;

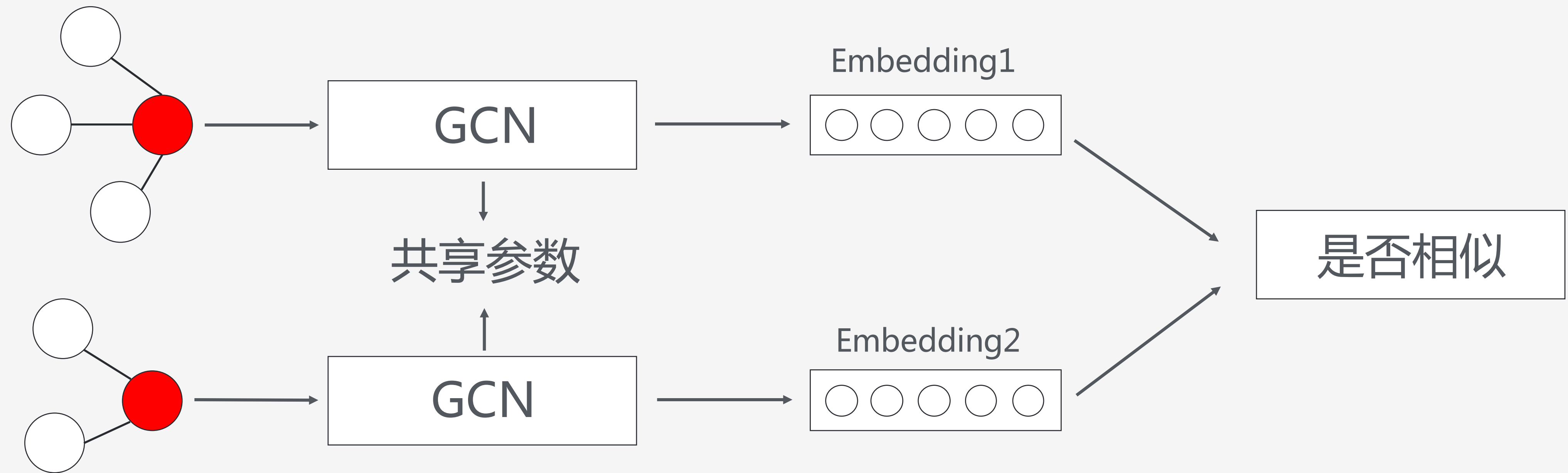
$t \leftarrow t + 1$;

until *stable node coloring is reached*;

[Thomas Kipf]

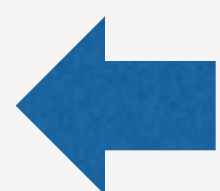
网络训练

- 优化目标：表示同一实体节点的embedding尽可能相近
- Siamese网络



案例: 重名实体排歧

Name	Affiliation
Jing Zhang	Shanghai Jiao Tong Univ.
	Yunnan Univ.
	Tsinghua Univ.
	Alabama Univ.
	Univ. of California, Davis
	Carnegie Mellon University
	Henan Institute of Education



Jing Zhang

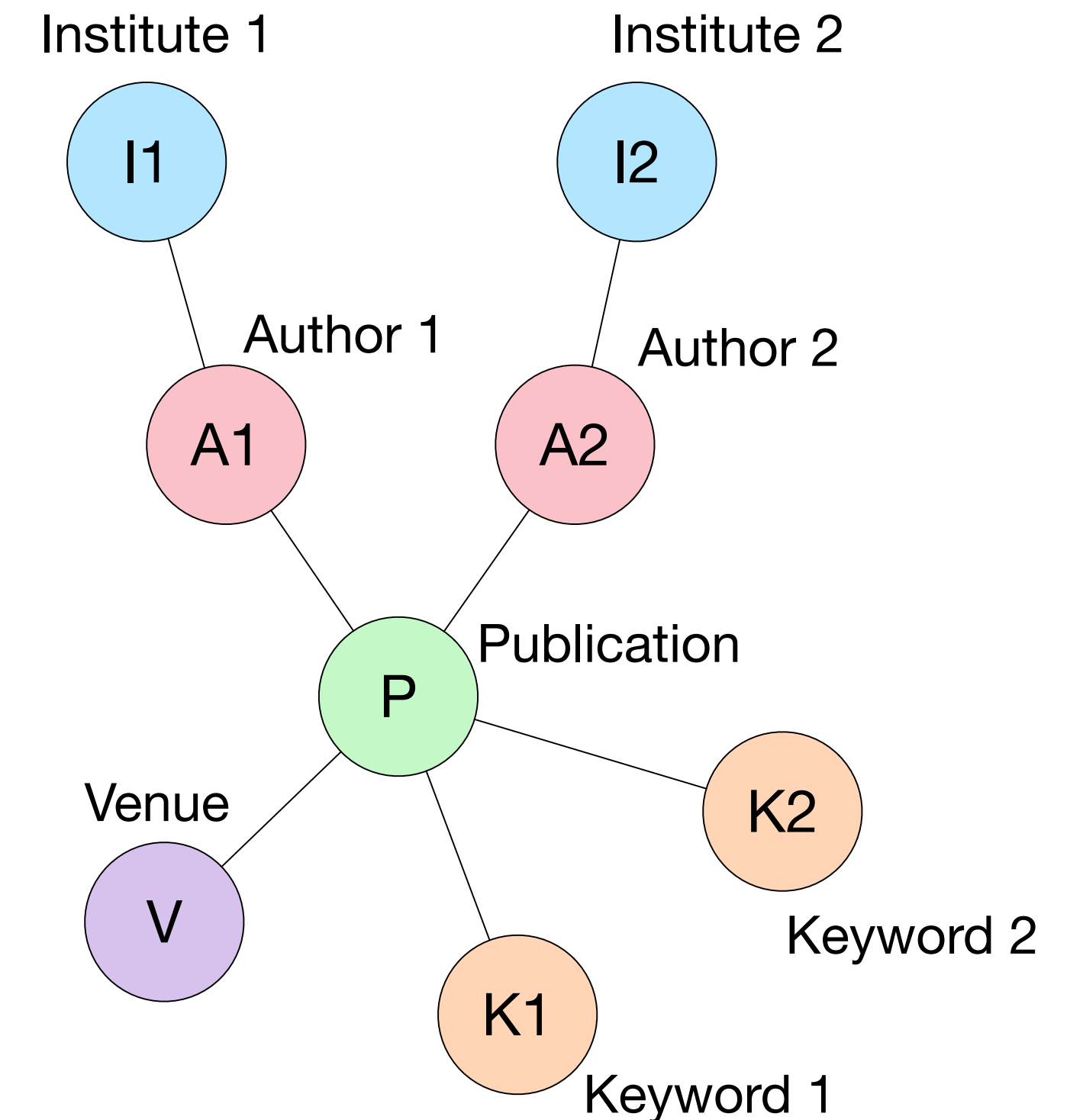
List of publications from the [DBLP Bibliography Server](#) - [FAQ](#)

[Coauthor Index](#) - Ask others: [ACM DL](#) - [ACM Guide](#) - [CiteSeer](#) - [CSB](#) - [Google](#)

		2007
83	EE	Jing Zhang, Guizhong Liu : Hyperspectral images lossless compression by a novel three-dimensional wavelet coding
82	EE	Jing Peng, Dongqing Yang , Changjie Tang , Jing Zhang, Jianjun Hu : CACS: A Novel Classification Algorithm Base
81	EE	Jing Zhang, Xi Chen , Ming Li : Computing Exact p-Value for Structured Motif. <i>CPM 2007</i> : 162-172
80	EE	Jing Zhang, Jie Tang , Juan-Zi Li : Expert Finding in a Social Network. <i>DASFAA 2007</i> : 1066-1069
79	EE	Guojun Chen , Jing Zhang, Xiaoli Xu , Yuan Yin : Real-Time Visualization of Tire Tracks in Dynamic Terrain with LC
78	EE	Jing Zhang, Hai Huang : Federate Job Mapping Strategy in Grid-Based Virtual Wargame Collaborative Environmen
77	EE	Maria Wimmer , Michael Goul , Jing Zhang: Minitrack: E-Government Information and Knowledge Management. <i>H</i>
76	EE	Kai Kang , Jing Zhang, Baoshan Xu : Optimizing the Selection of Partners in Collaborative Operation Networks. <i>IC</i>
75	EE	Lingshuang Shao , Jing Zhang, Yong Wei , Junfeng Zhao , Bing Xie , Hong Mei : Personalized QoS Prediction for We
74	EE	Benyong Liu , Jing Zhang, Xiaowei Chen : Adaptive Training of a Kernel-Based Representative and Discriminative I
73	EE	Jilong Wang , Jing Zhang: Federation Based Solution for Peer-to-Peer Network Management. <i>International Confer</i>
72	EE	Jing Zhang, Fanhuai Shi , Jianhua Wang , Yuncaai Liu : 3D Motion Segmentation from Straight-Line Optical Flow. <i>M</i>

问题形式化

- 每篇文献可表示成一个异构子图
- 包含文献、作者、机构、会议、关键词等节点
- 需要匹配同类型且表示同一实体的节点
- 关键在于度量节点相似度



两个应用场景

- 知识图谱初始建立
 - 根据Embedding将表示同一实体的节点聚类在一起 (Clustering)
- 增量数据对齐将Embedding扩展到新来的子图 (Inductive Learning)
 - 将新的节点分配给知识库中的已有节点 (Assignment)

知识图谱初始建立

- 问题一：缺乏人工标注训练样本
- 解决方案：弱监督学习
- 基于规则生成正例
 - 例如：将同名且名字不常见的作者认为是同一个人
- 随机采样负例

知识图谱初始建立

- 问题二：聚类个数未知
- 解决方案：使用贝叶斯信息准则(BIC)估计聚类个数

增量数据更新与维护

- 根据模型预测新节点的embedding
- 将新进入的子图中的节点与知识图谱对齐
 - 可能需要新建节点
- 构建二分类模型
 - 输入：子图节点embedding，知识图谱中候选节点的embedding
 - 输出：是否为同一实体

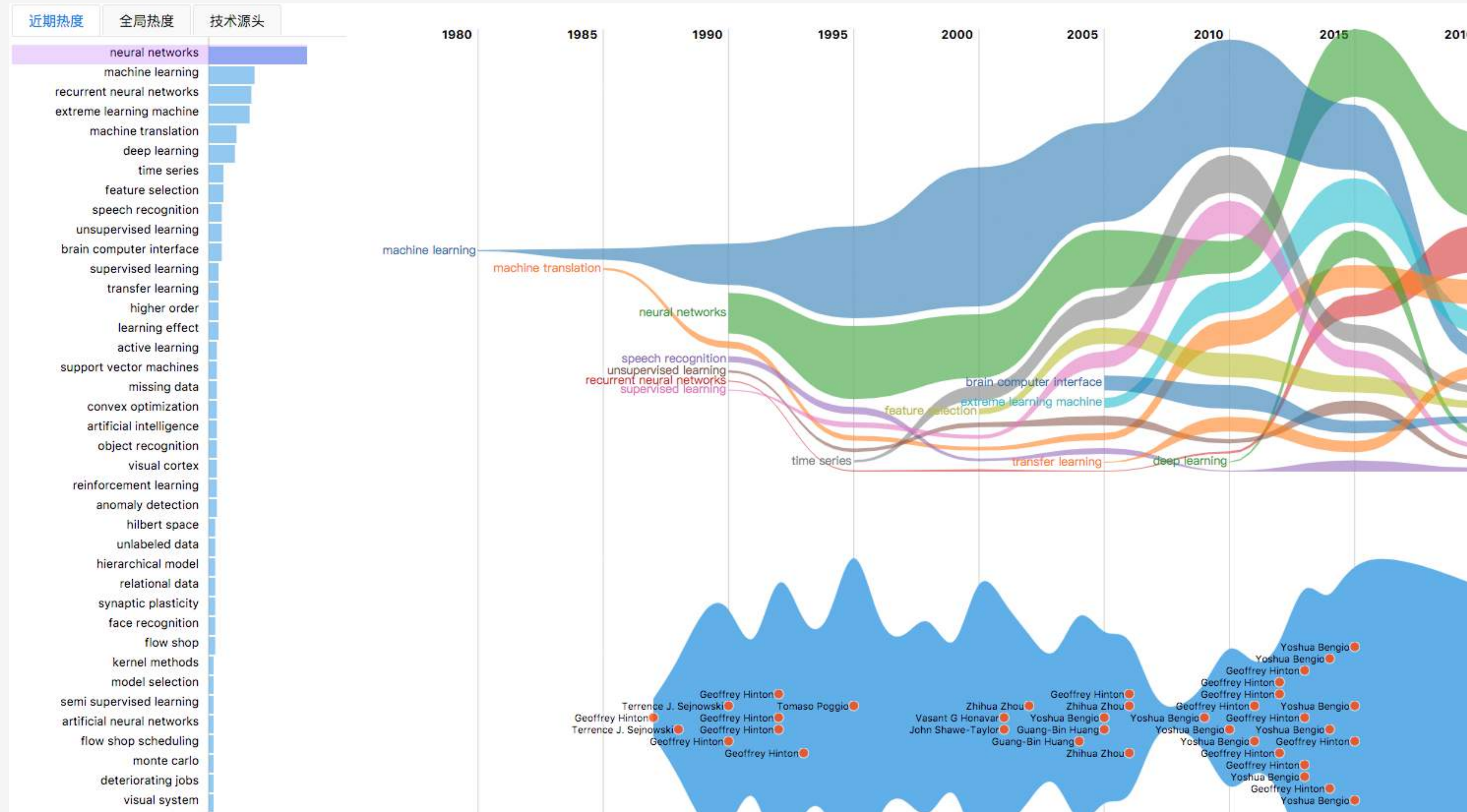
增量数据更新与维护

- 候选配对过滤 (Blocking)
 - 使用粗略的快速算法（例如基于规则），对候选配对空间进行过滤
 - 降低候选匹配个数，避免 $O(N^2)$ 次相似度计算

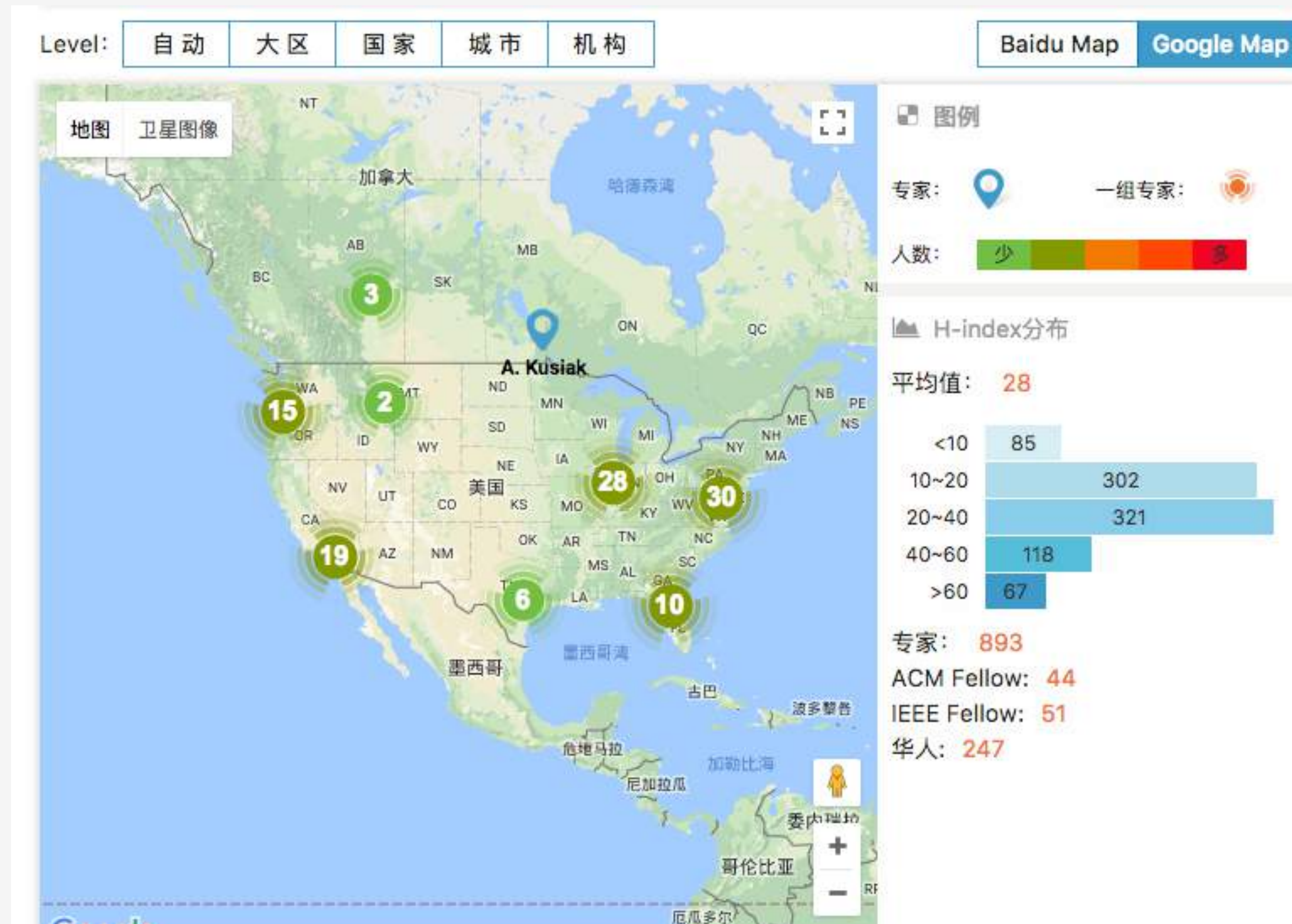
TABLE OF CONTENTS 大纲

- 异构数据融合与知识图谱构建
- 异构数据融合中的机器学习方法
 - 度量学习
 - 表征学习
- 图嵌入学习与图卷积网络
- 案例：重名实体排歧
- 基于科技知识图谱的智能服务

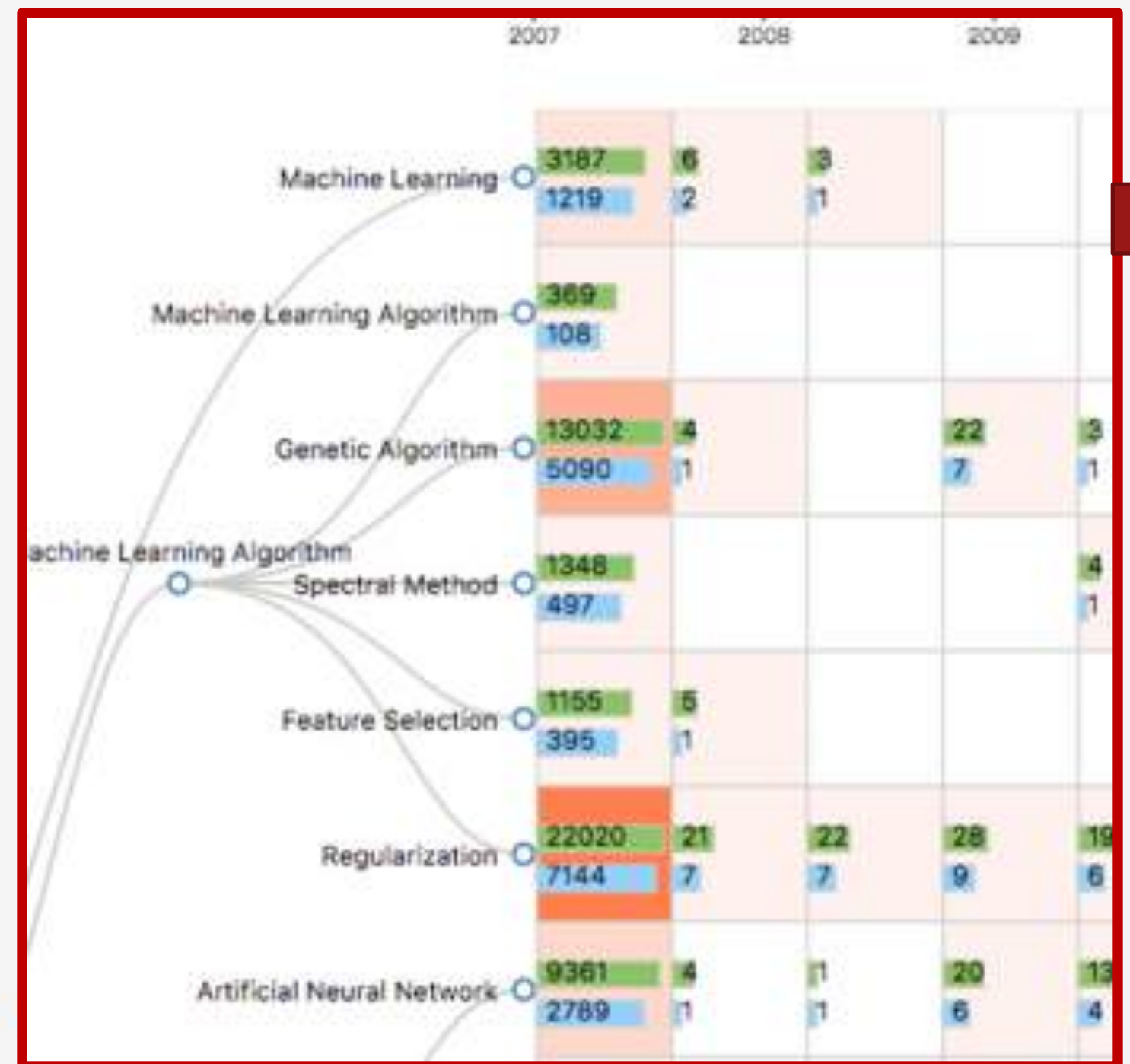
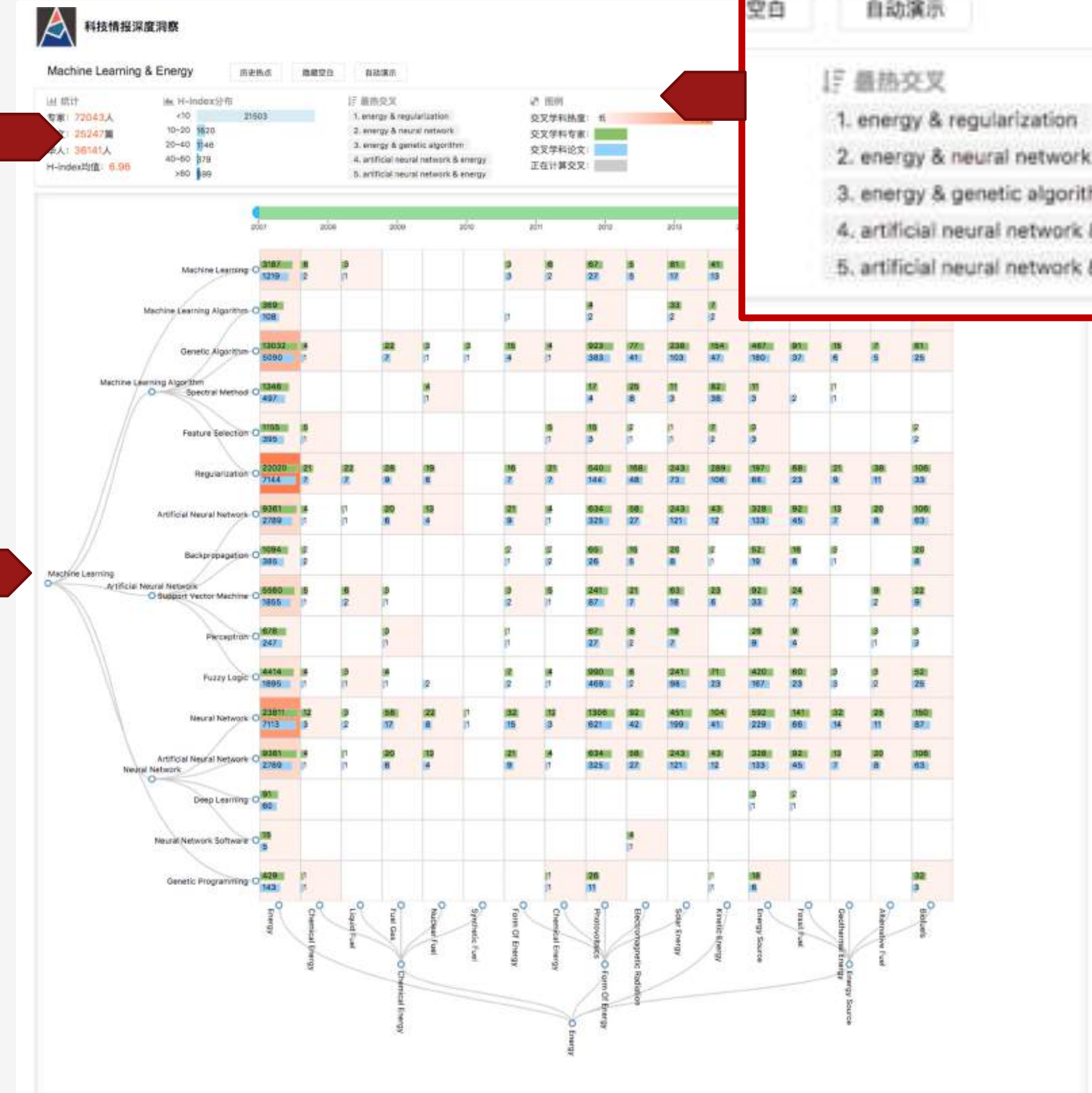
技术趋势分析



基于地理位置的专家发现



学科交叉创新分析



THANK YOU

如有需求，欢迎至 [讲师交流会议室] 与我们的讲师进一步交流

