

数据来源：数据库产品上市商用时间



第十三届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2022

数据智能 价值创新



线上直播 | 2022/12/14-16



图数据库构建长城汽车数据一元化 OneID数据底座

长城汽车-产业数智化中心 (IDC)

陈晓 数据中台大数据工程师

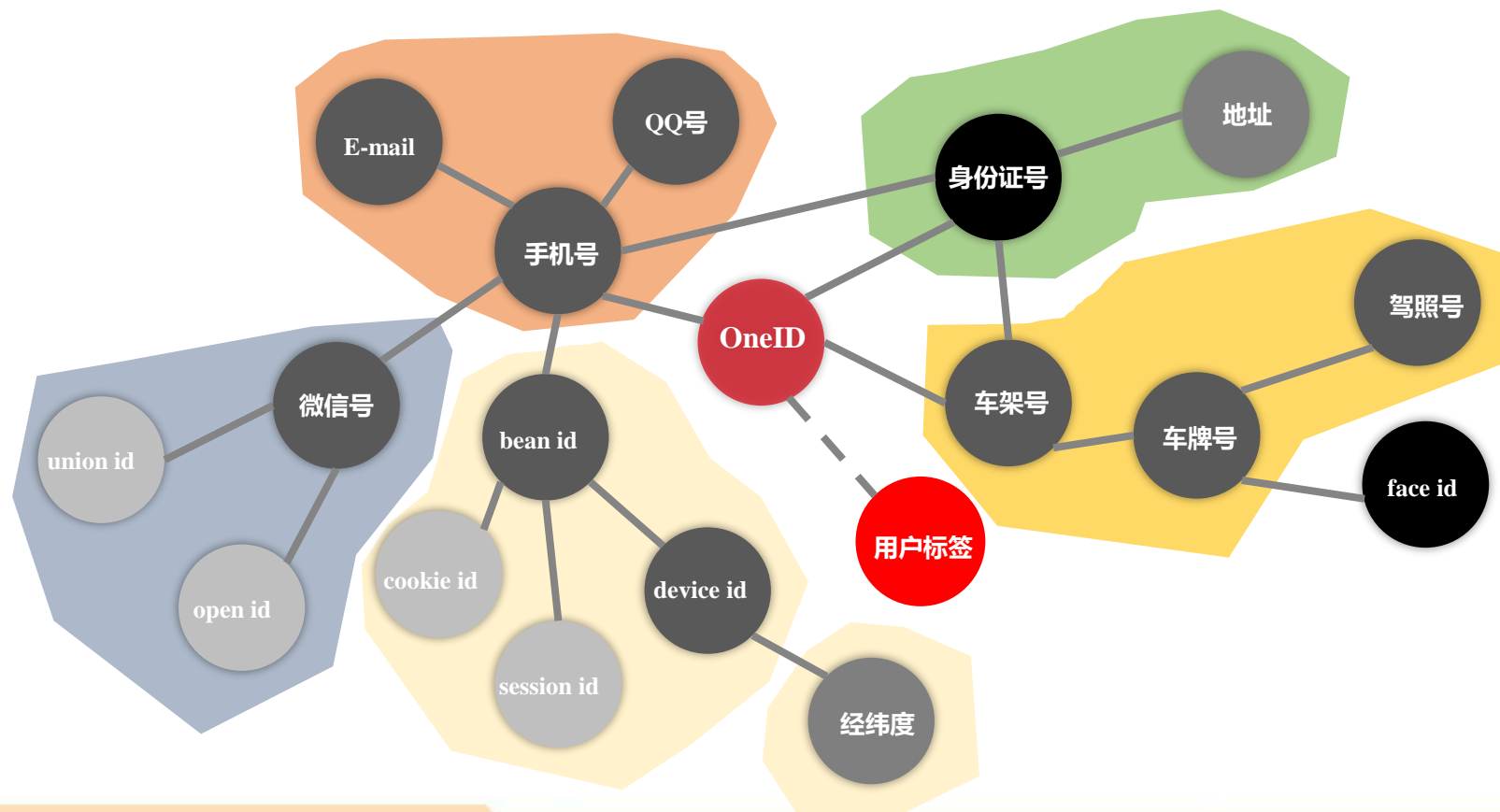
CONTENTS
目录

- 1 图数据库在OneID项目中的应用
- 2 OneID项目中的图挖掘技术
- 3 未来技术方向思考

Part 1 图数据库在OneID项目中的应用

OneID实现数据一元化

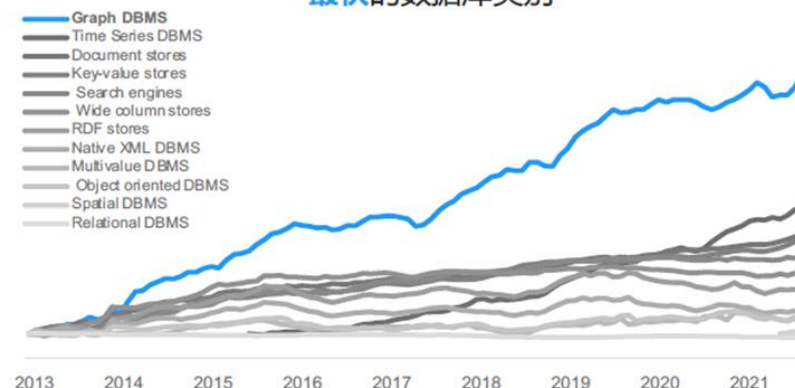
- 数据孤岛问题；例如：PC端、APP端、微信小程序、相关系统单元等.
- 统一用户识别；用ID-mapping技术将多ID归一管理，解决作为用户的唯一身份识别.



图数据库

- 图数据库本身就是基于事物关联关系的模型表达工具，与OneID的拉通数据的理念相契合；
- 图数据库在解决大规模数据实体间复杂关系的查询问题上，具备天然优势，查询效率提升显著；
- 图数据库过去10年增长趋势最快，验证了图数据库成为一种趋势；
- NebulaGraph 在处理超大数据集（千亿节点万亿条边）的查询上保持**毫秒级**查询延时。

图数据库年增长率**100%**，是过去十年**采用率增长最快**的数据库类别¹



75%世界百强企业采用图数据库，起到标杆作用，未来**快速辐射腰部企业**²

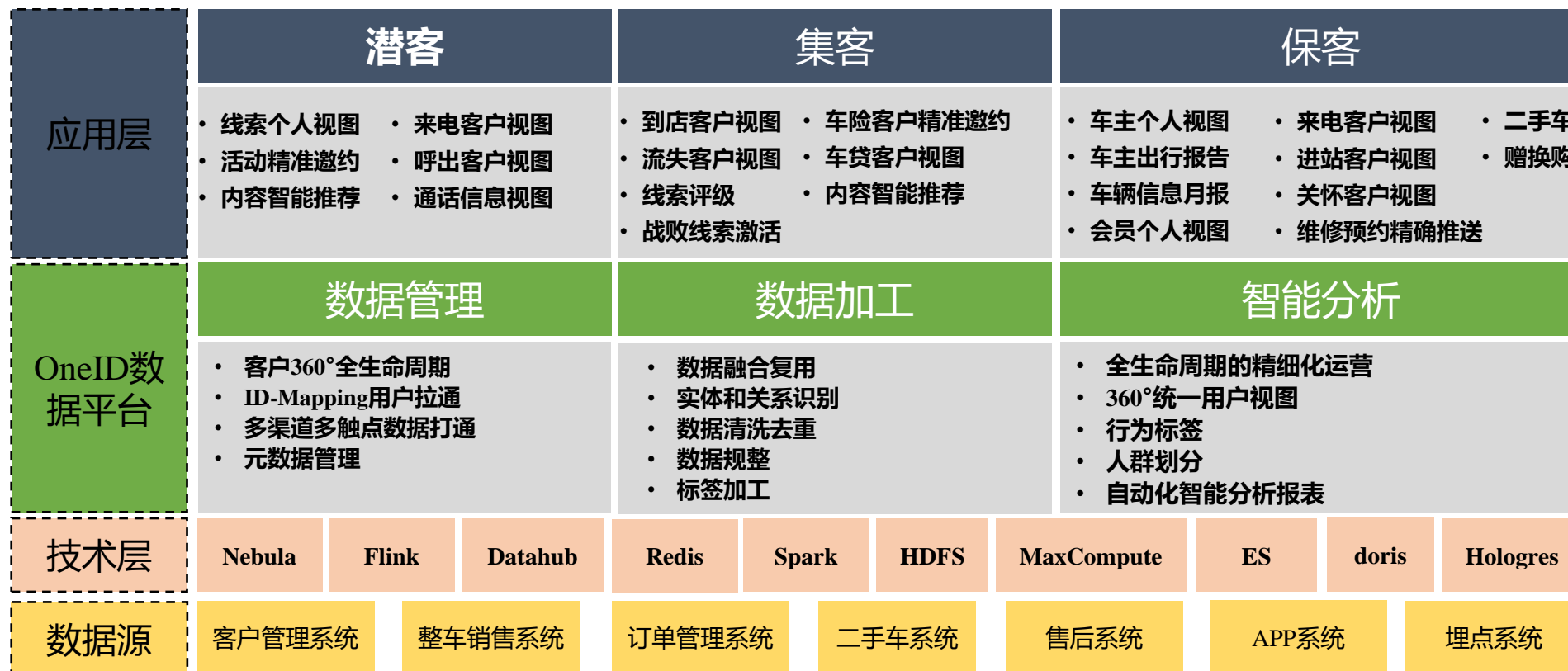
20/25
金融企业

7/10
零售企业

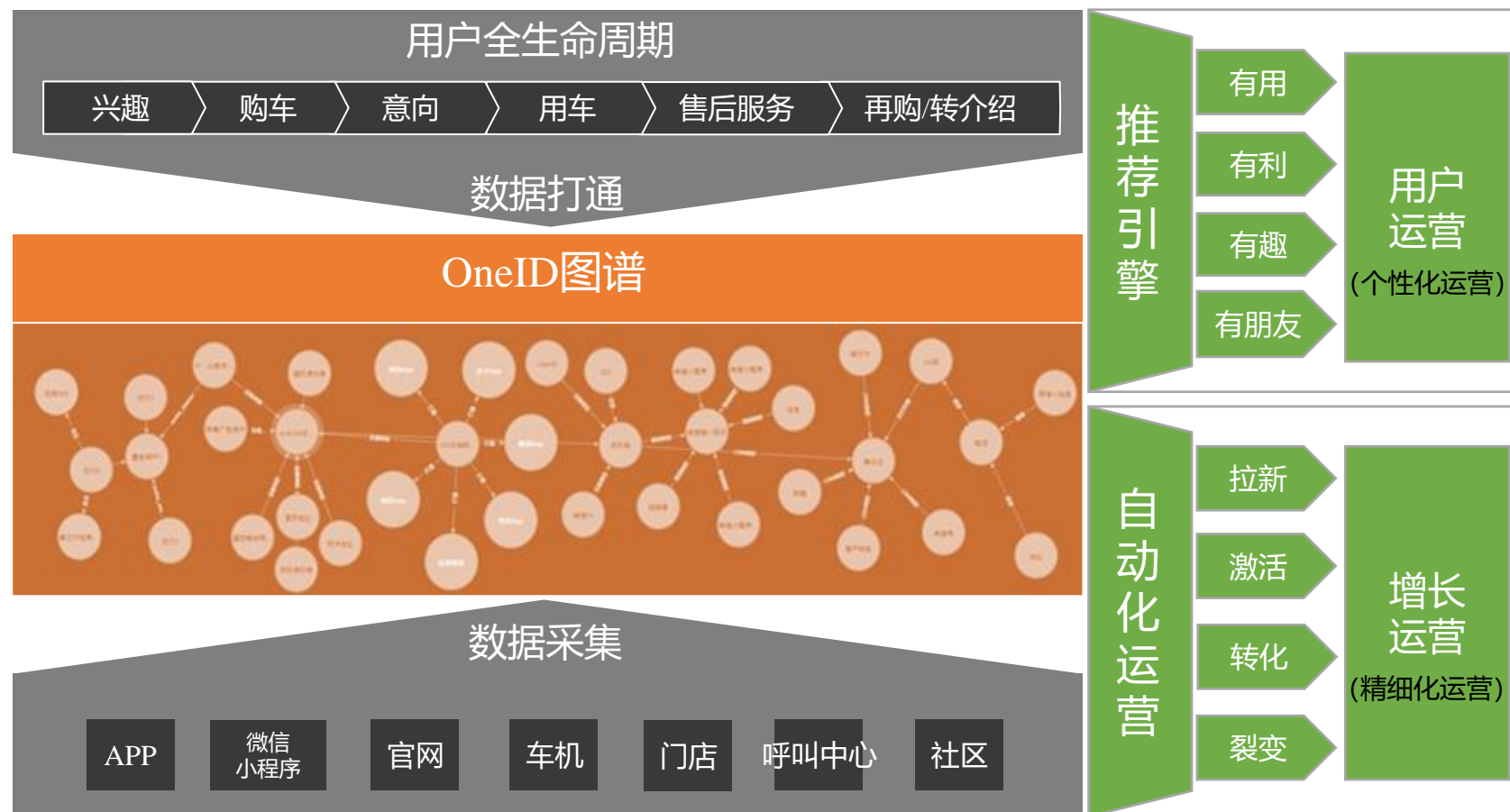
7/10
互联网企业

4/5
电信企业

OneID平台架构



OneID标签图谱赋能用户增长和个性化运营



各大推荐场景:

- ① 首页个性化推荐
- ② 频道首页信息流推荐
- ③ 相关推荐
- ④ 搜索推荐
- ⑤ 热门推荐
- ⑥ 猜你喜欢
- ⑦ 社交推荐
- ⑧ 商品图谱推荐

OneID赋能业务场景

OneID为客服人员、销售顾问、运营人员，在流量、线索、场景等价值流转化上提供落地服务。

价值流

潜客(流量转化)

集客(线索转化)

保客(场景转化)

客服人员

销售顾问

运营人员



呼叫中心 – 来电弹屏页
车辆查询页



登录明细页



用户增长模型页

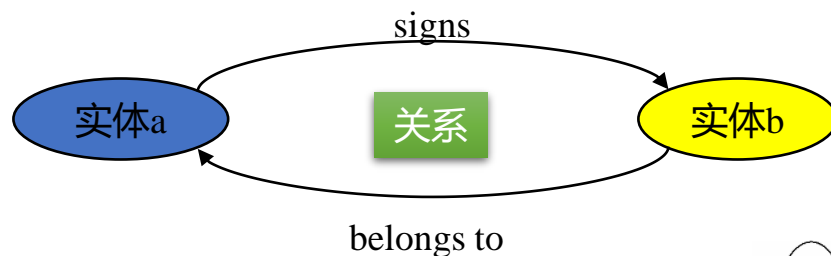
OneID实现了用户360°视图，应该将OneID的建设经验和能力复用到车辆360°、零部件360°、供应商360°上，为更多领域业务提供数据服务。

Part 2 OneID项目中的图挖掘技术

什么是图

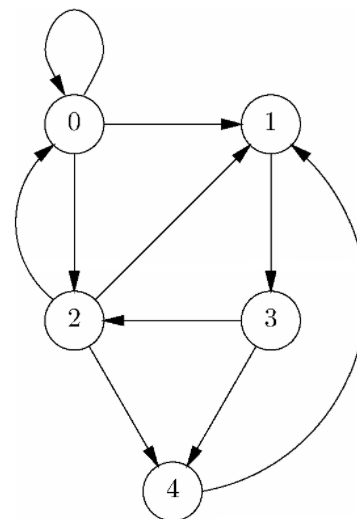
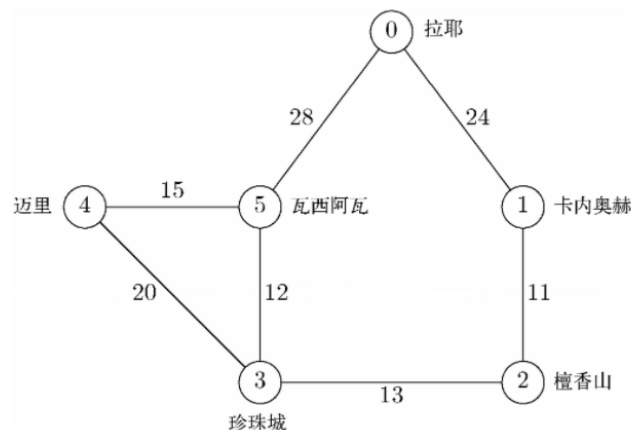
图就是二元关系，它利用一系列由线（称为边）或箭头（称为弧）连接的点（称为节点）提供了强大的视觉效果。

图的本质是由二元关系组构成，实体-关系-实体模型：



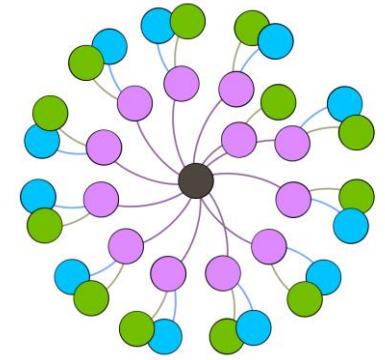
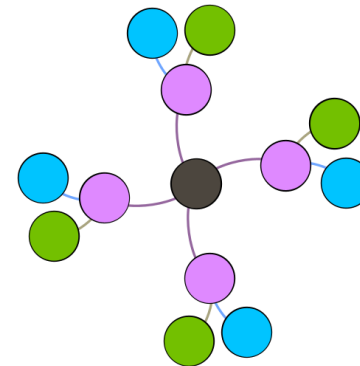
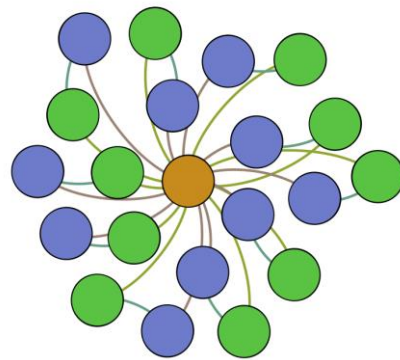
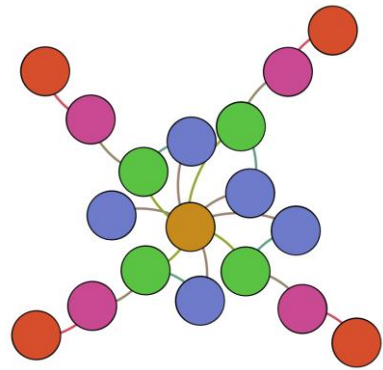
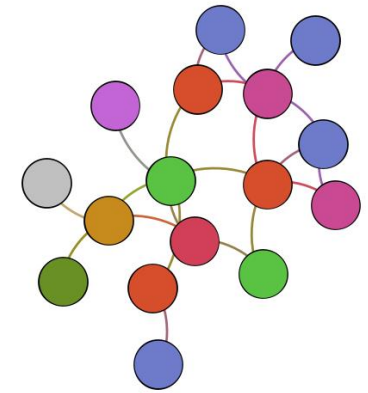
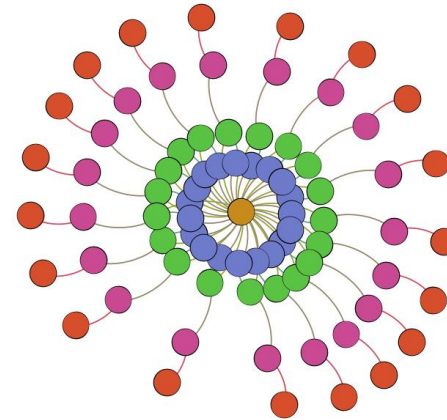
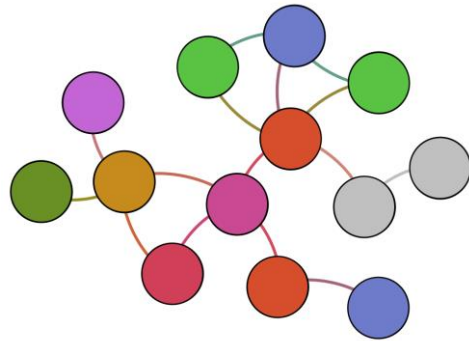
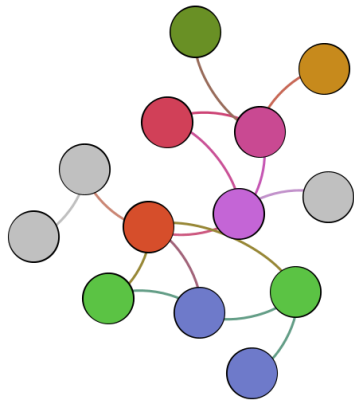
图有多种形式：

- 无向图、有向图
- 加权、未加权
- 同构、异构
- 单边、多重边
- 静态图、动态图



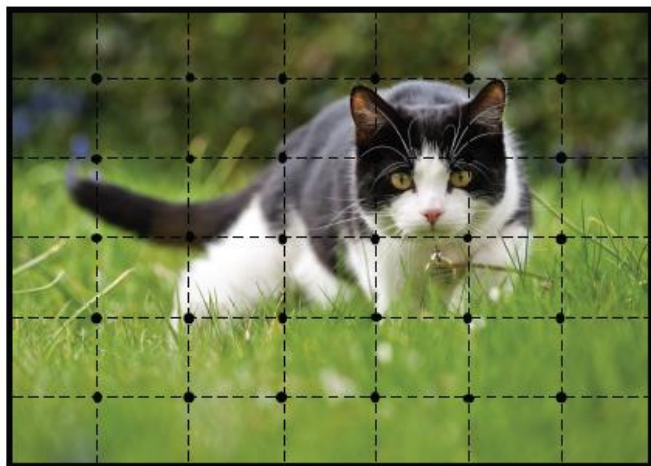
OneID图库中不同形态的图

形态各异的图网络：

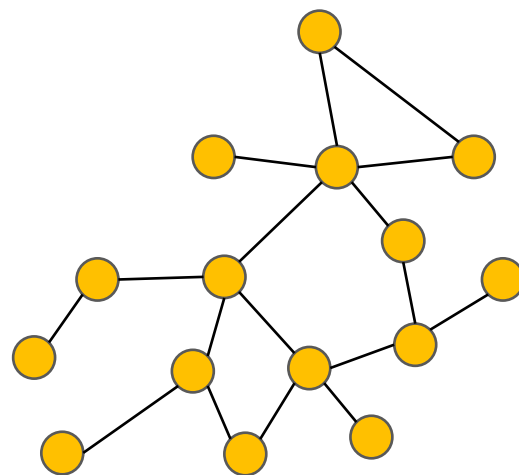


如何进行图建模

图像和拓扑图对比：



左:欧氏空间中的图像



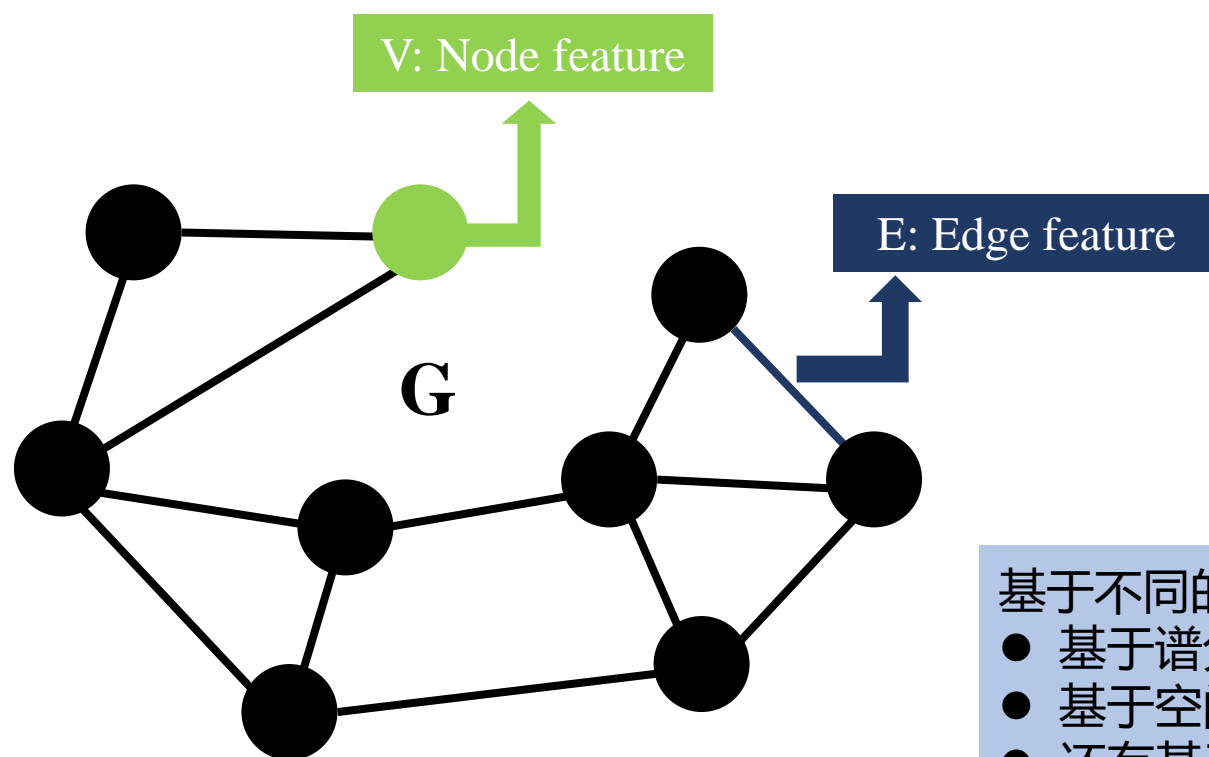
右:非欧氏空间中的图

CNN在非欧氏空间这里不适用.

因此，不是所有的事情
可以表示为序列或网格。
我们将如何应用神经网络？

如何进行图建模

图卷积的应用：



图表示: $G = (V, E)$

邻域矩阵 A_{ij} :

$$A_{ij} = \begin{cases} 1 & \text{If } \{v_i, v_j\} \in E \text{ and } i \neq j \\ 0 & \text{otherwise} \end{cases}$$

度矩阵 D : $D_{ii} = d(v_i)$

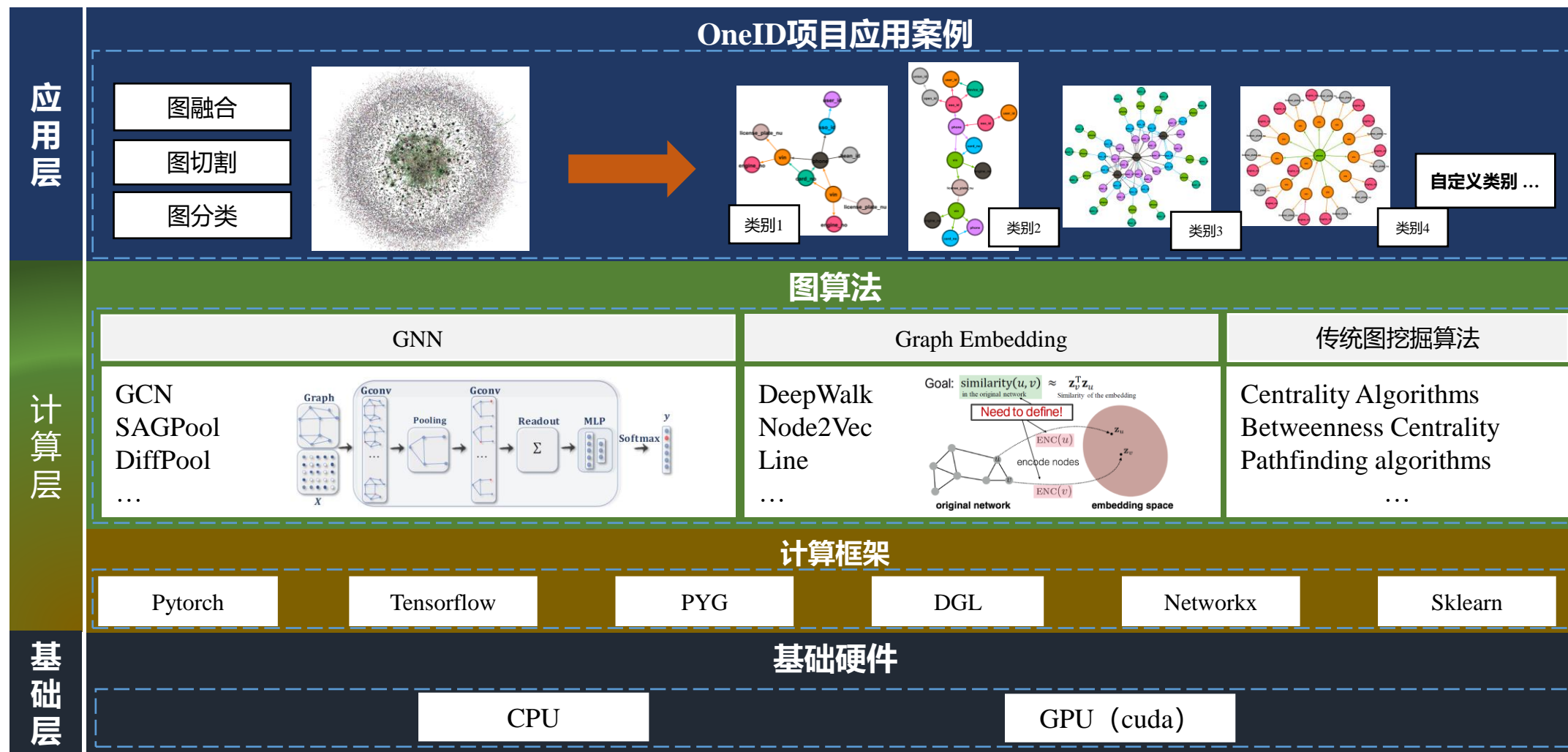
图卷积神经网络卷积层的基本传播规则:

$$H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)})$$

基于不同的传播规则和采样技术:

- 基于谱分解的方法: GCN、ChebNet、Spectral Network
- 基于空间结构的方法: GraphSAGE、GAT
- 还有基于递归运算和跳链接的其他传播模块

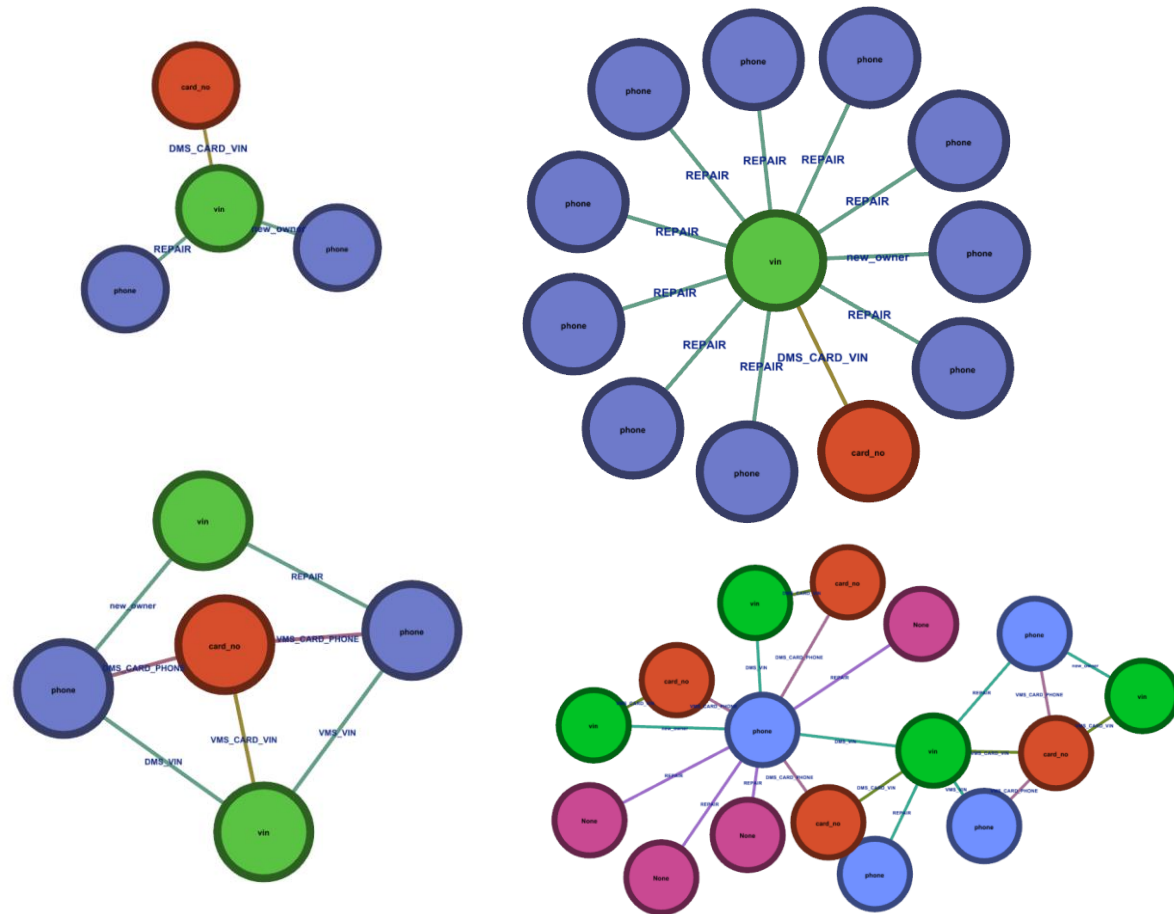
图算法在OneID项目中的应用



图算法应用1—车主信息融合

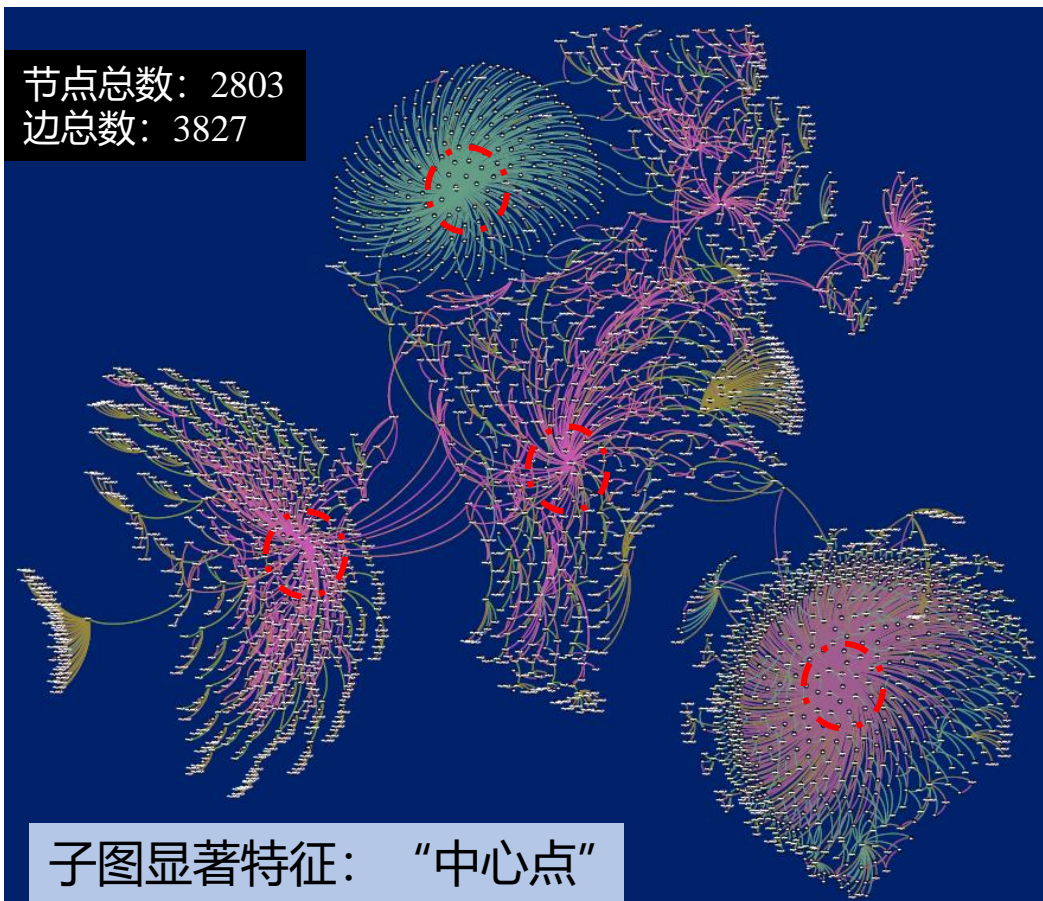
例如：利用图技术，对车和人是一对多的关系进行融合、优化。

补全：利用融合信息，对缺失信息推断；
去重：重复、真伪信息过滤；
推断：先验知识辅助推断。



图算法应用2—图切割提纯

信息融合过程中存在的“信息冗余”问题：



“信息冗余”问题：

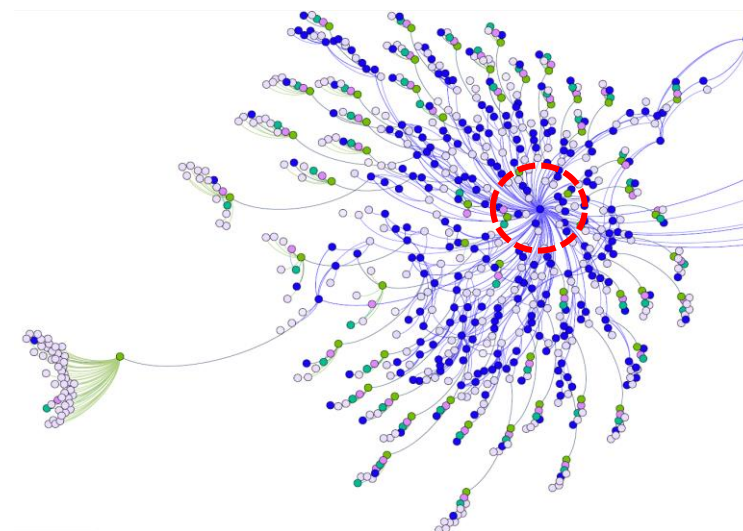
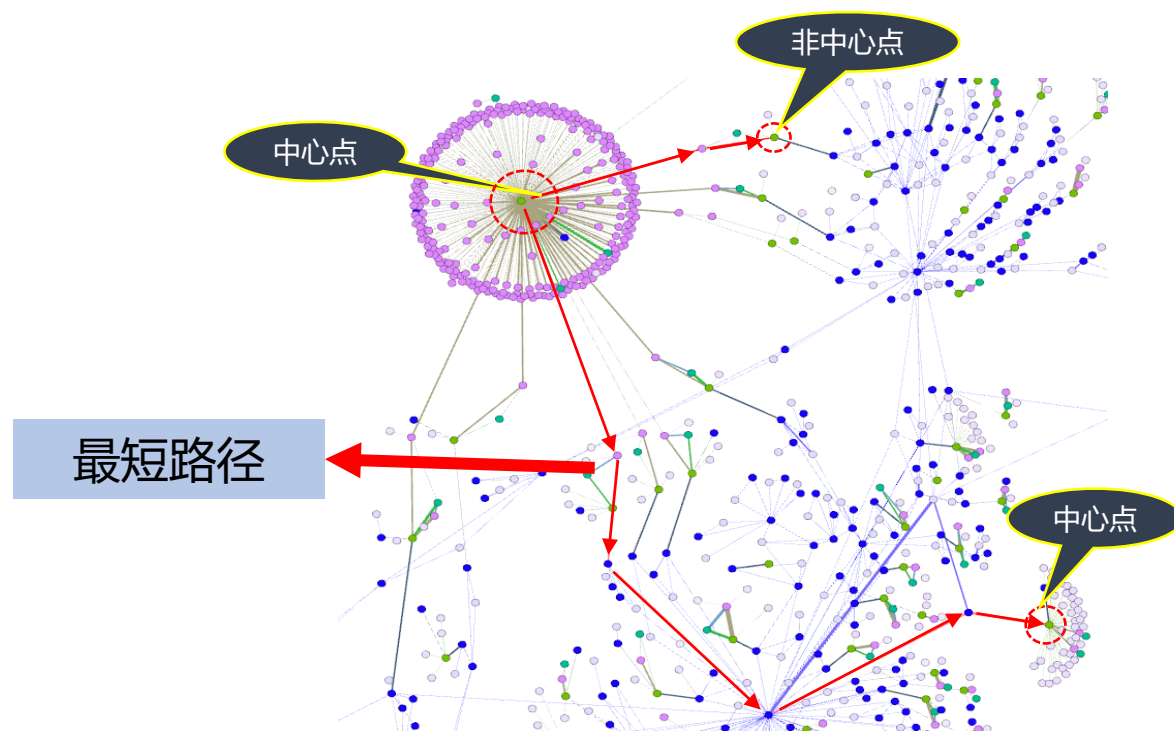
- 销售、企业 通过手机、身份证将多人信息糅在同一张子图中；
- “羊毛党” 刷单，将多人信息绑定到同一张子图中；
- “脏数据” 流入绑定多人信息。

思路：在同一子图中,存在出入度比较高的一类节点，结合业务将此类节点作为“中心节点”，探索它与临近节点的关系、与他人信息之间的最短路径，进而将子图切割，实现无关信息的隔离。

图算法应用2—图切割提纯

利用最短路径算法、中心点算法等实现图切割：

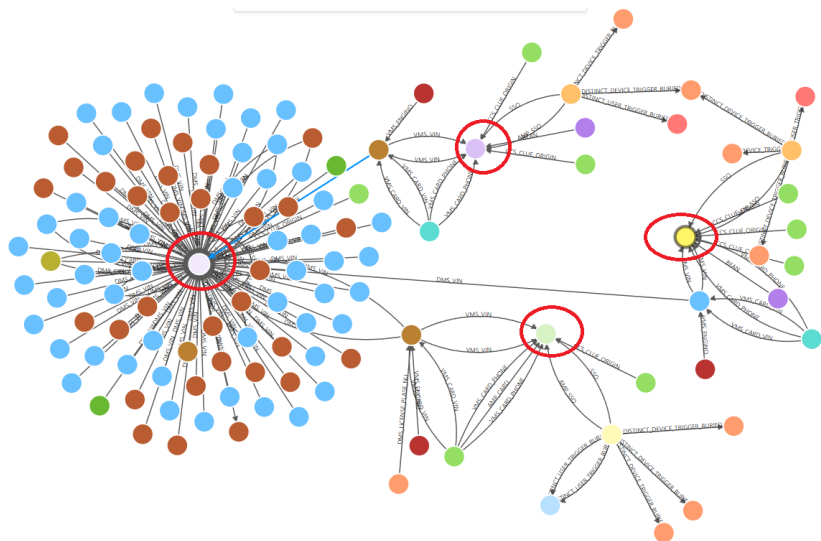
- 识别出中心点节点和非中心点节点
- 使用最短路径算法完成切割
- 识别出问题设备“羊毛党”



图算法应用2—图切割提纯

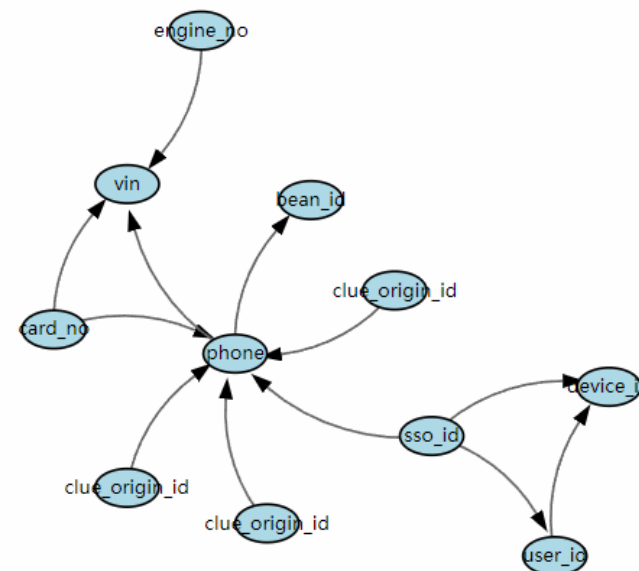
切割后效果展示：

案例1：在保证图信息的完整、独立的情况下，提取出某公司的相关信息



- 提升图的数据质量
- 支撑推荐系统、用户画像、图打分等应用服务体系建设

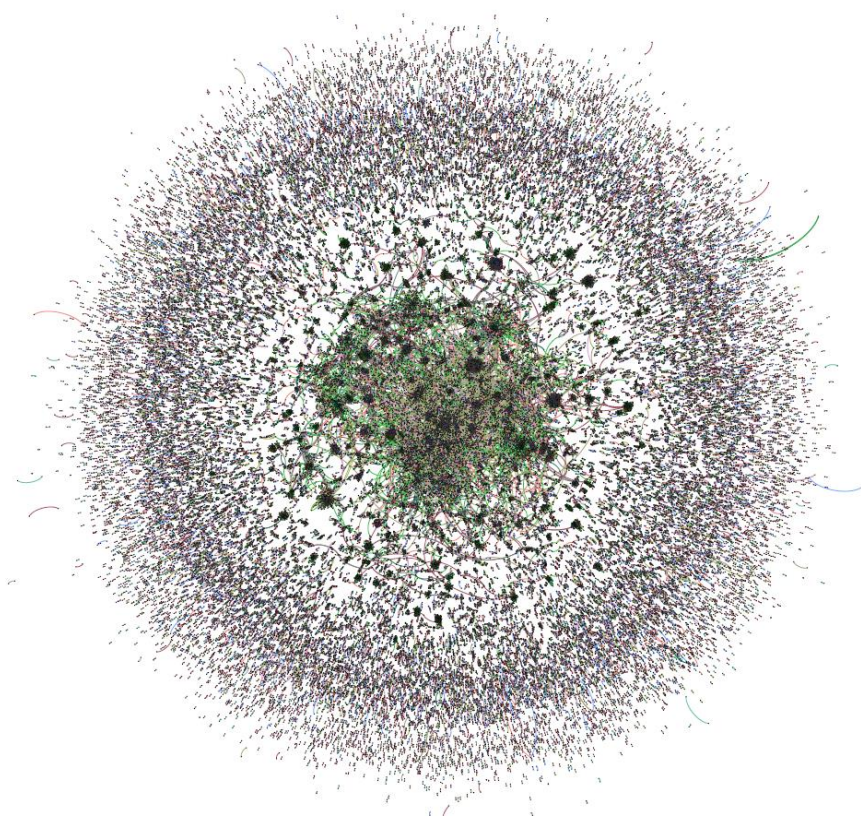
案例2：某个人的信息提纯



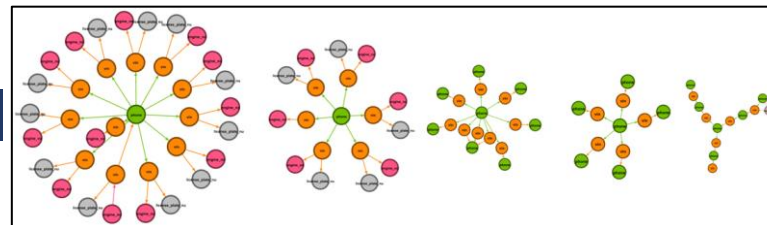
图算法应用3—图分类

子图分类：物以类聚人以群分

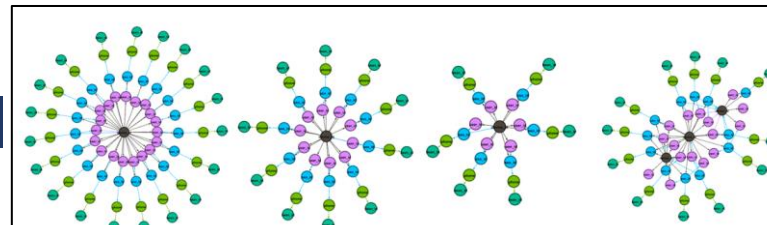
例如：家庭用户、销售人员、公司大客户等.



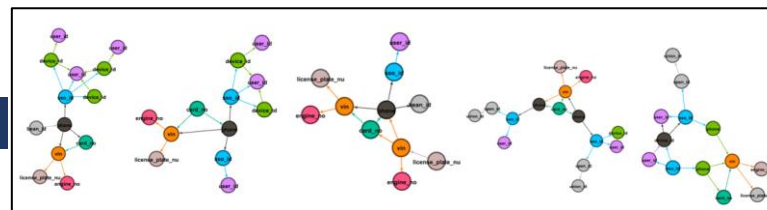
类别1



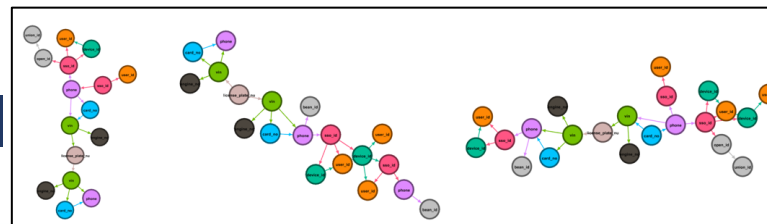
类别2



类别3

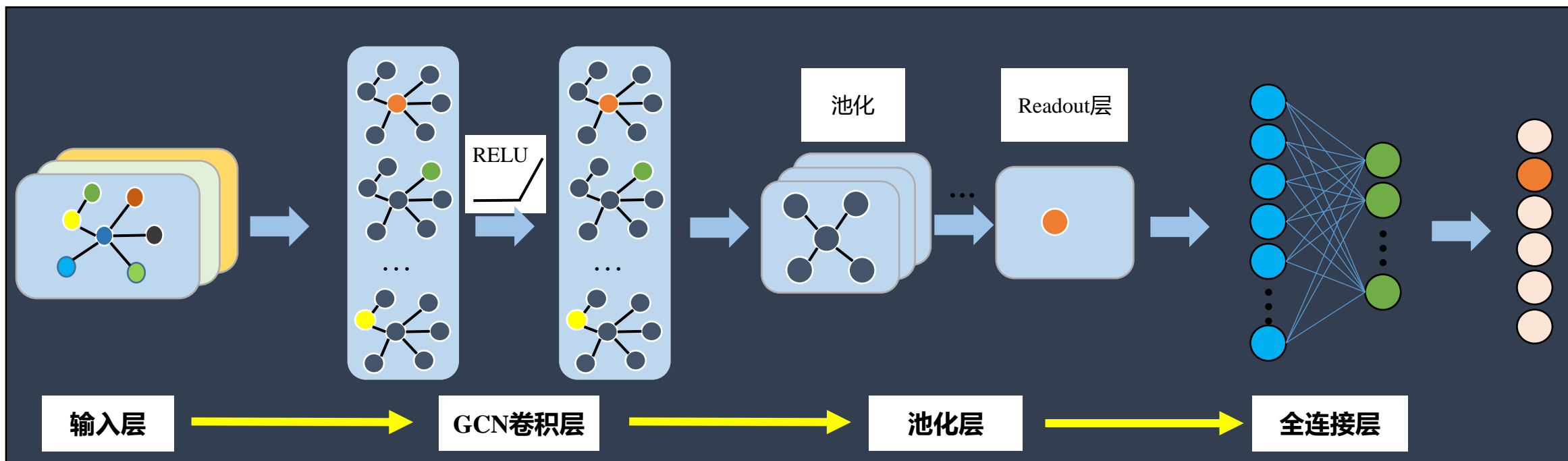


类别4



图算法应用3—图分类

图神经网络子图分类模型架构：



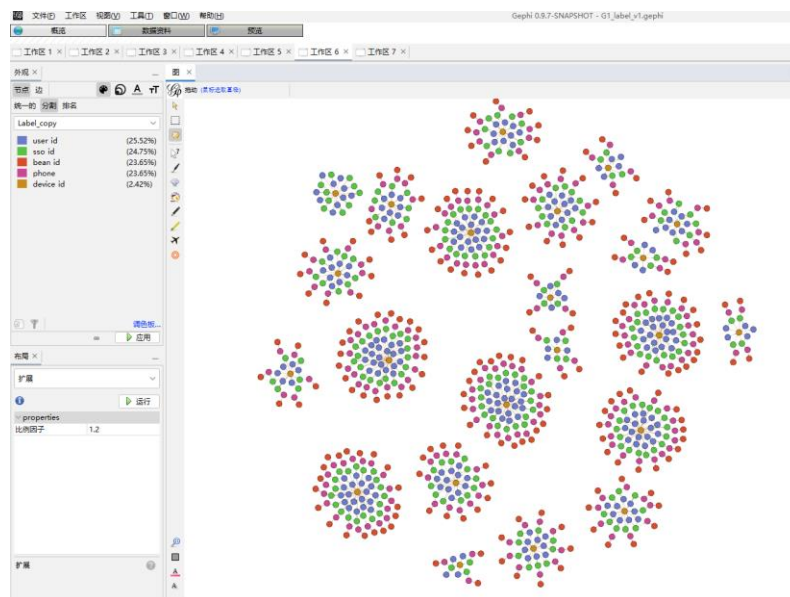
- 输入层：度矩阵 D 、邻域矩阵 A 、节点特征 $H_{(0)}$
- 卷积层： $H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)})$
- 池化层： $idx = \text{top-rank}(Z, \lceil kN \rceil), \quad Z_{mask} = Z_{idx}$
- 全连接层

图算法应用3—图分类

数据准备工作:

图采样模块设计

- **动态图**: 对采样数据要做离线存储, 对子图进行随机采样;
- **信息冗余**: 对子图进行预剪枝, 过滤掉无关节点;
- **信息缺失**: 保留了缺失属性的节点.



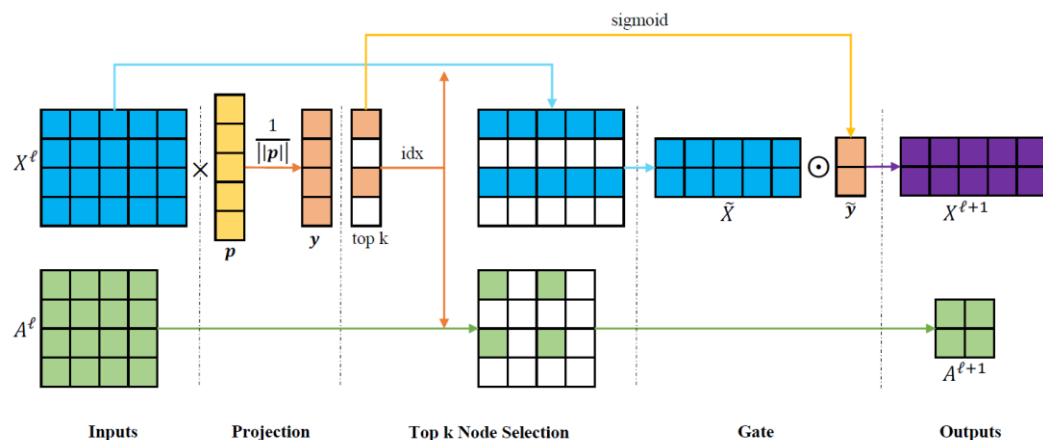
图数据加工模块: Dataset

利用nebula3、networkx、Gephi完成图内容的解析、图重构、图可视化展示等:

- 邻接矩阵
- 节点的图标识
- 图类型
- 节点特征

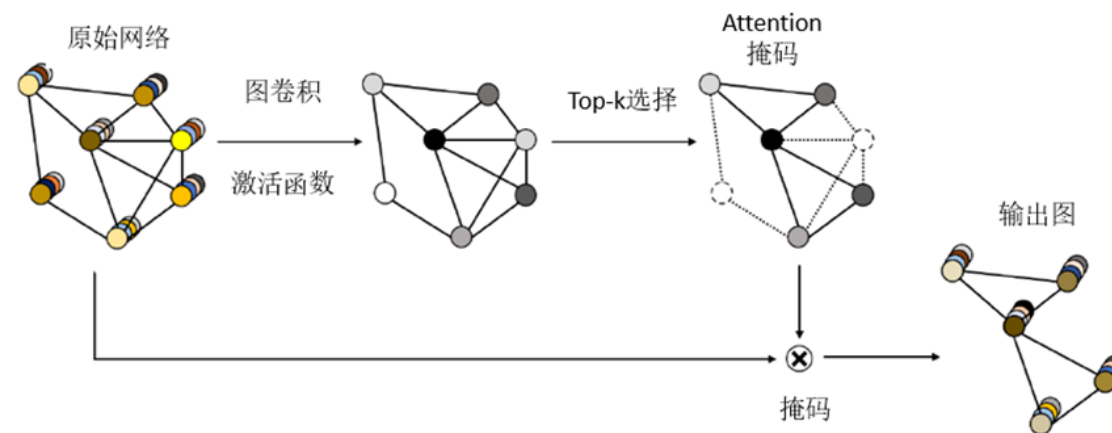
图算法应用3—图分类

Pooling常用于图分类任务中，常用的池化方式：



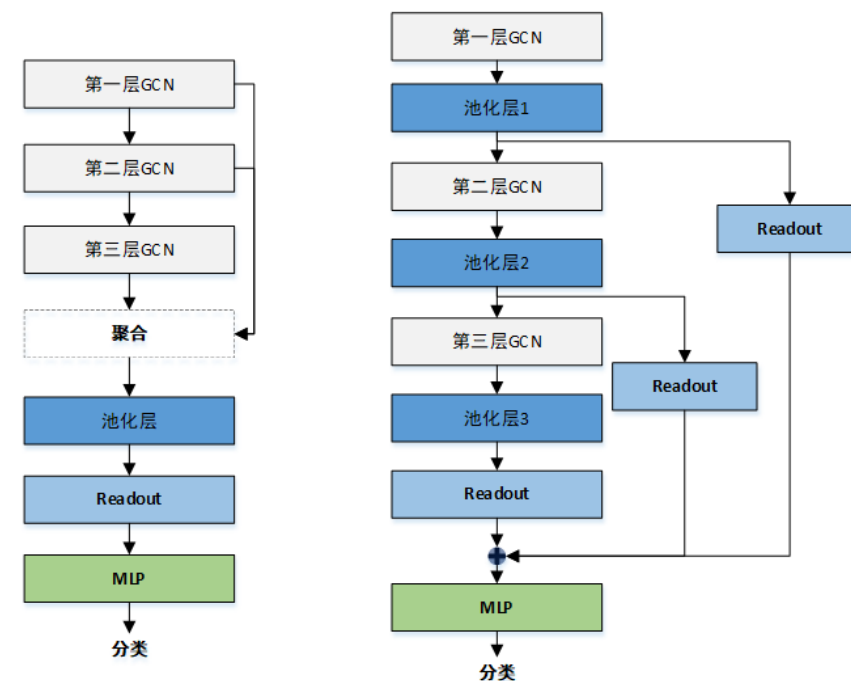
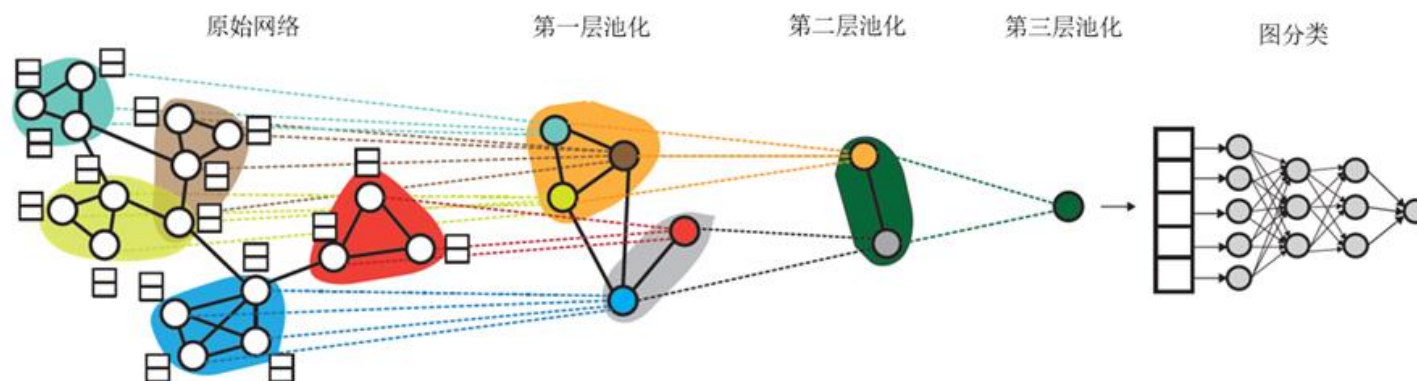
TopkPool：全局筛选，难以捕捉到局部信息。

SAGPool：采用GNN获取节点的self-attention得分，融合图的特征和拓扑结构。



图算法应用3—图分类

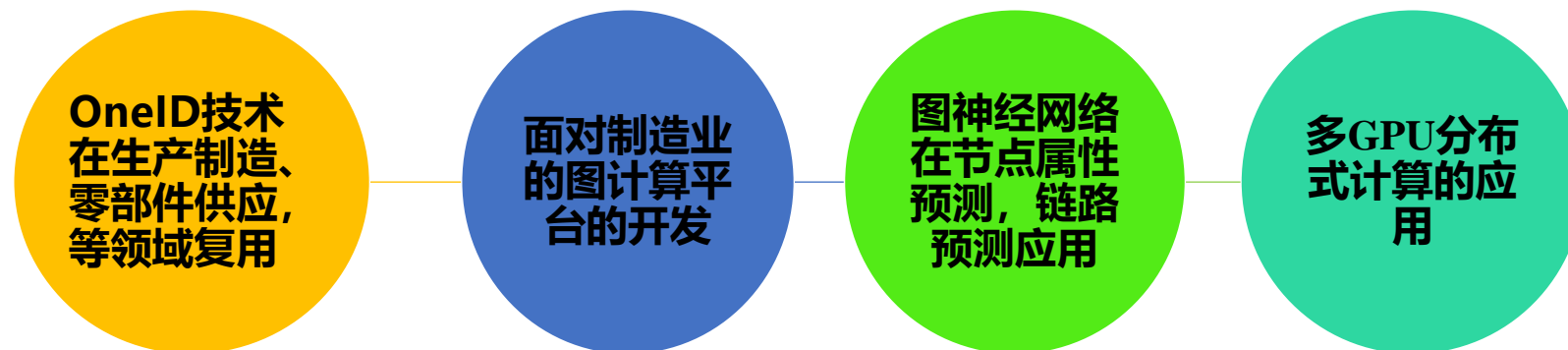
DiffPool: 采用GNN获取节点embedding表达, 然后对节点进行聚类, 生成一个粗粒度的图, 重复以上过程.



Part 3 / 未来技术方向思考

未来技术方向思考

图技术能力建设的四个方向拓展：



文献综述

图神经网络综述类:

- 【1】 Jie Zhou et al. [Graph neural networks: A review of methods and applications.](#)
- 【2】 Ziwei Zhang et al. [Deep Learning on Graphs: A Survey.](#)
- 【3】 Jie Zhou et al. [Graph Neural Networks: A Review of Methods and Applications.](#)
- 【4】 Zonghan Wu et al. [A Comprehensive Survey on Graph Neural Networks.](#)

图表示学习Embeddings:

- 【5】 **DeepWalk**: Perozzi et al. [Online Learning of Social Representations.](#)
- 【6】 **LINE**: Jian Tang et al. [Large-scale Information Network Embedding.](#)
- 【7】 **node2vec**: Aditya Grover et al. [Scalable Feature Learning for Networks.](#)

不同卷积核:

- 【8】 **GCNConv**: Kipf and Welling et al. [Semi-Supervised Classification with Graph Convolutional Networks.](#)
- 【9】 **ChebConv**: Defferrard et al. [Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering.](#)
- 【10】 **GATConv**: Veličković et al. [Graph Attention Networks.](#)
- 【11】 **SAGEConv**: Hamilton et al. [Inductive Representation Learning on Large Graphs.](#)

文献综述

不同池化方式:

- 【12】 **Top-K Pooling**: Gao and Ji et al. [Graph U-Nets](#), Cangea et al. [Towards Sparse Hierarchical Graph Classifiers](#) and Knyazev et al. [Understanding Attention and Generalization in Graph Neural Networks](#).
- 【13】 **DiffPool**: Ying et al. [Hierarchical Graph Representation Learning with Differentiable Pooling](#).
- 【14】 **Set2Set**: Vinyals et al. Order Matters: [Sequence to Sequence for Sets](#).
- 【15】 **SAG Pooling**: Lee et al. [Self-Attention Graph Pooling](#) and Knyazev et al. [Understanding Attention and Generalization in Graph Neural Networks](#).
- 【16】 **ASAPooling**: Ranjan et al. ASAP: [Adaptive Structure Aware Pooling for Learning Hierarchical Graph Representations](#).

THANKS

SQL Server
vertica
D B 2
G B a s e
O r a c l e
达梦数据库
神舟通用
KingbaseES

2010

2014

2018

openGauss
OceanBase
ArkDB
RASESQL
HotDB
StellarDB
QianBase xTP
云树Shard
GoldenDB
DolphinDB
MatrixDB
DynamoDB
SinoDB
FastData
Galaxybase
KunDB
GDB
GaussDB
PolarDB
KunDB
Spacture
SequoiaDB
OushuDB
ArgoDB
开务数据库
GreatDB
MongoDB
TDSQL
TiDB
Tapdata
StarRocks
UbiSQL