

AliGraph：大规模图神经网络平台

李永（九丰）

阿里巴巴资深技术专家
计算平台事业部

PAI机器学习平台介绍

阿里云 2019阿里云峰会·上海
开发者大会

机器学习PAI

▶ 视频简介

阿里云机器学习平台PAI（Platform of Artificial Intelligence），为传统机器学习和深度学习提供了从数据处理、模型训练、服务部署到预测的一站式服务。

立即购买

产品定价

使用文档

<https://data.aliyun.com/product/learn>

简单易用

功能强大

PAI机器学习平台介绍

首页

实验

数据源

组件

模型

设置

模版列表

新建实验

实时热点新闻挖掘

PAI OnlineLearning挖掘实时热点新闻

370 位用户

从模版创建 查看文档

【推荐算法】商品推荐

通过协同过滤算法实现商品推荐。

3390 位用户

从模版创建 查看文档

【文本分析】新闻分类

通过主题模型实现了整个文本分类的流程。

2897 位用户

从模版创建 查看文档

【图算法】金融风控实验

利用图算法，针对个人信用，解决金融行业的风控问题。

1919 位用户

从模版创建 查看文档

雾霾天气预测

机器学习算法计算出二氧化氮对于雾霾影响最大。

1402 位用户

从模版创建 查看文档

心脏病预测案例

包括数据预处理、特征工程、模型训练和预测等一套机器学习流程。

3707 位用户

从模版创建 查看文档

文本聚类舆情分析

本实验提炼了芸聆平台在用户反馈舆情、客满舆情、社会舆情等文本数据上进行聚类分析的流程

368 位用户

从模版创建 查看文档

【NLP.AI】机器人工厂问答

自然语言理解算法在问答场景的应用，机器人工厂帮你定制云端智能问答服务。

229 位用户

从模版创建 查看文档

【推荐...390 X

运行 部署 Auto ML DataStudio

cf_训练_data

SQL取购买行为

协同过滤...ec-1

SQL...出格式

JOIN-1

JOIN-2

SQL去重

全表统计-1

cf_结果_data

过滤与映射-1

SQL去重2

全表统计-2

一站式

机器学习平台

数据处理，特征加工
模型训练，在线预测

STOA

机器学习算法

经典机器学习
阿里自研算法
兼顾离线&实时

精心优化的 深度学习引擎

阿里内部大规模任务打磨，
PAI-Tensorflow引擎

弹性

模型预测服务

简化模型部署负担
支持扩缩容，降低成本



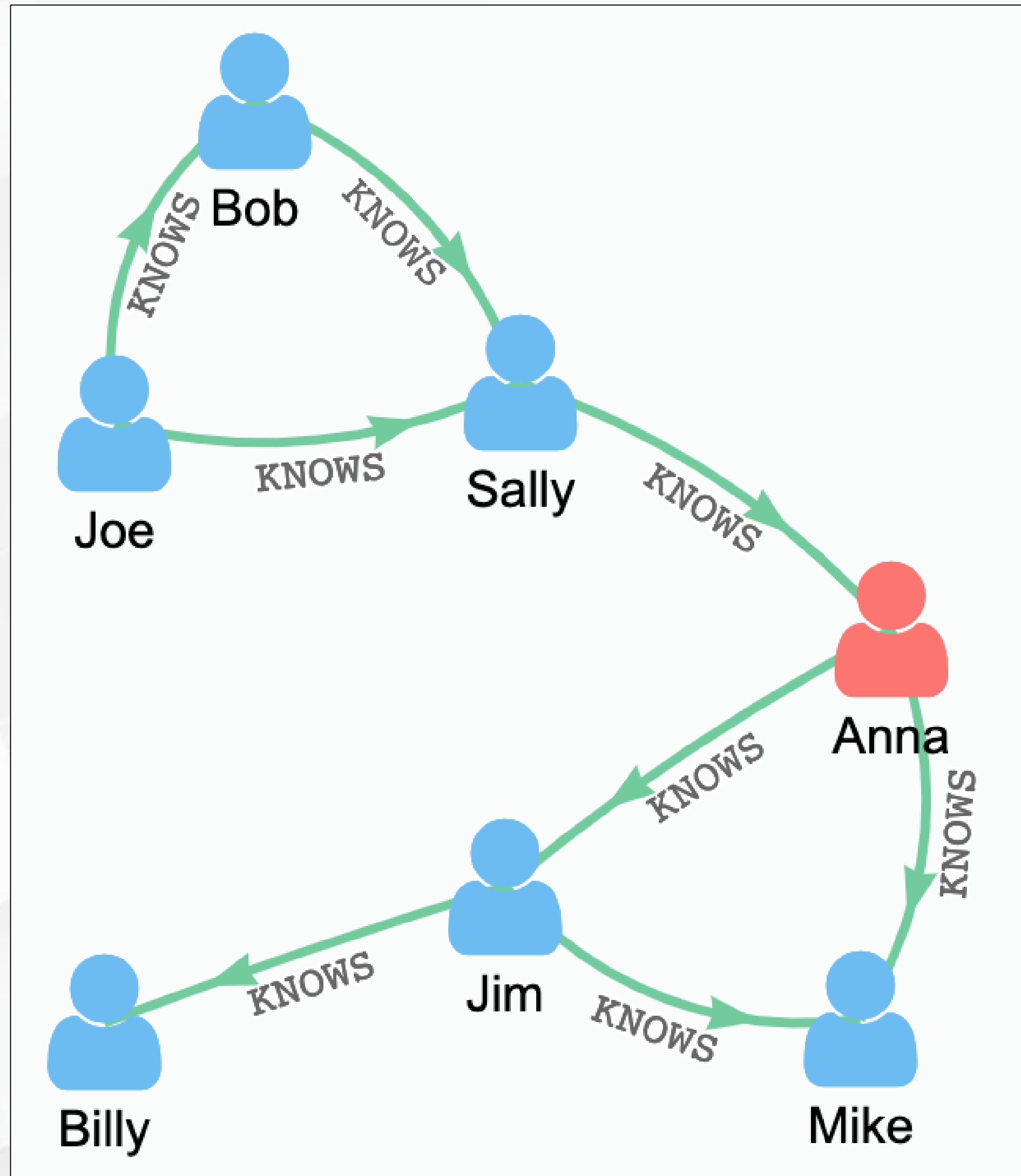
模型可解释性

难以理解模型特征以及决策逻辑
缺少数学工具来评测网络的表达能力



因果关系推理

美国增加关税会影响中国GDP吗
博士毕业的薪水会比硕士要高吗？



$$G = (V, E)$$

$V = \text{Vertex}$ $E = \text{Edge}$

同构图

异构图

有向图

无向图

图数据库

Neo4J, Titan

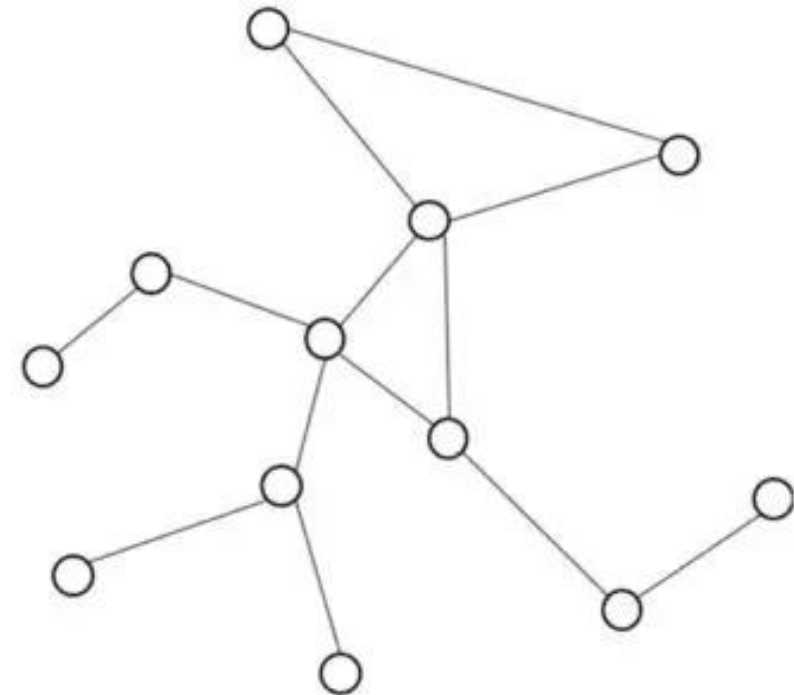
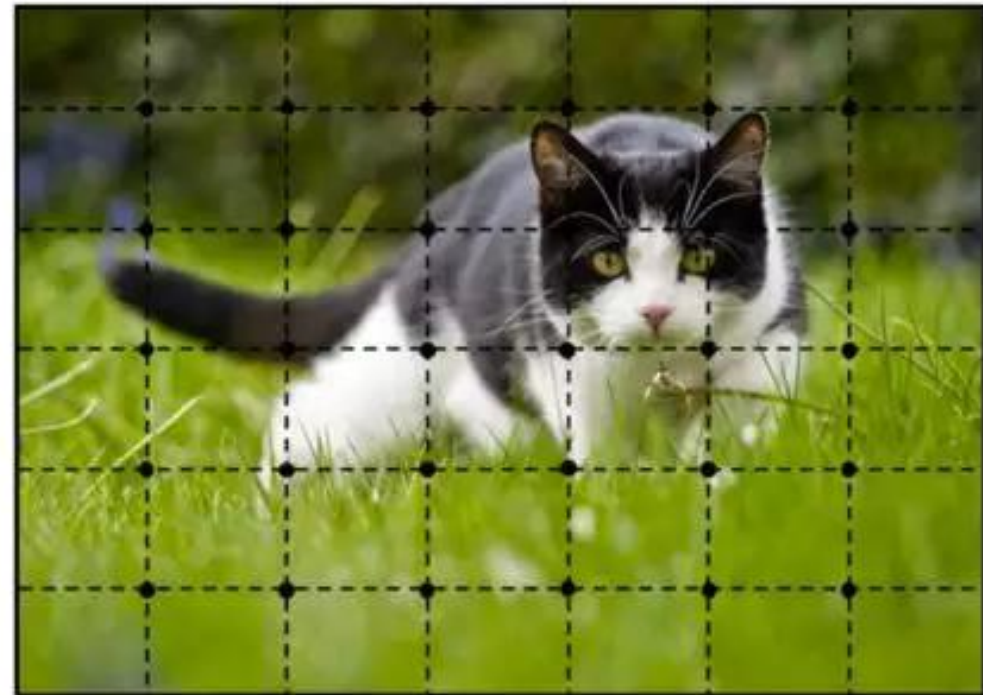
实时图分析查询

Aliyun GraphCompute, Amazon Neptune

离线图处理

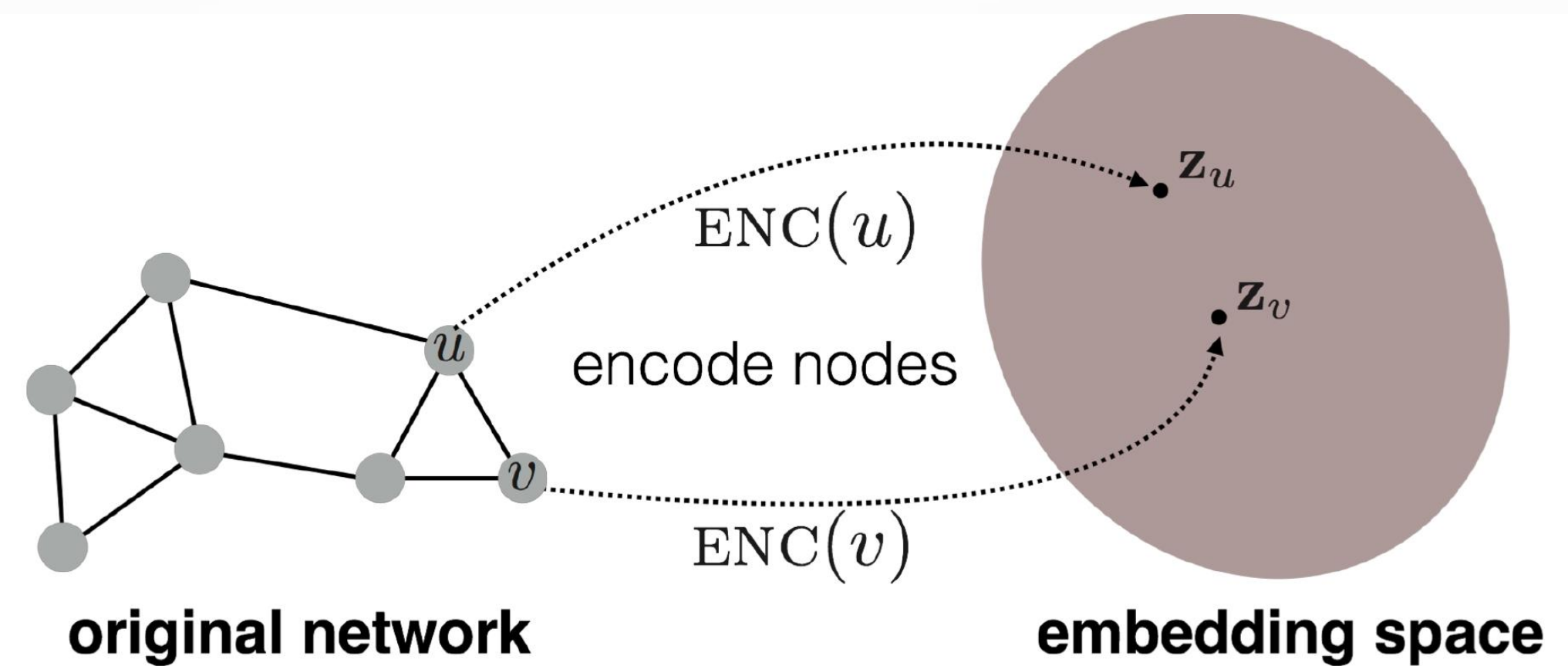
Pregel, PowerGraph, Spark GraphX

图+深度学习 = 图神经网络



非规则化的数据

Graph Embedding



图神经网络的应用领域



社交领域



推荐领域



知识图谱



生命科学

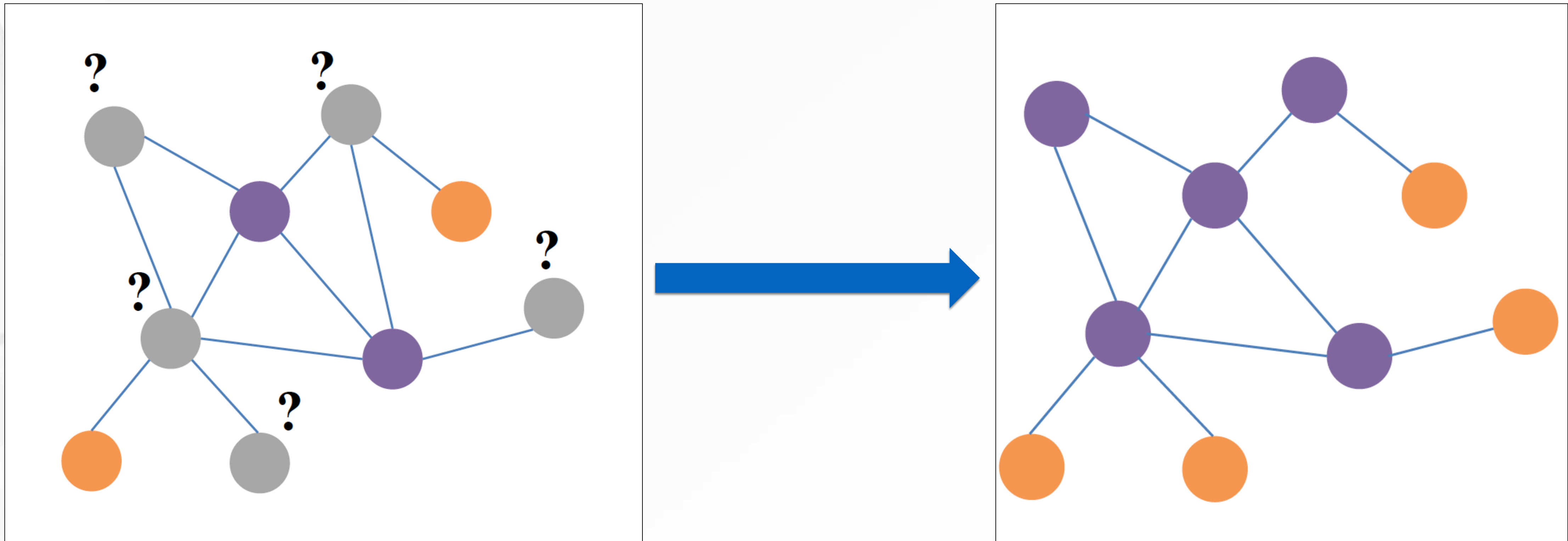


反作弊

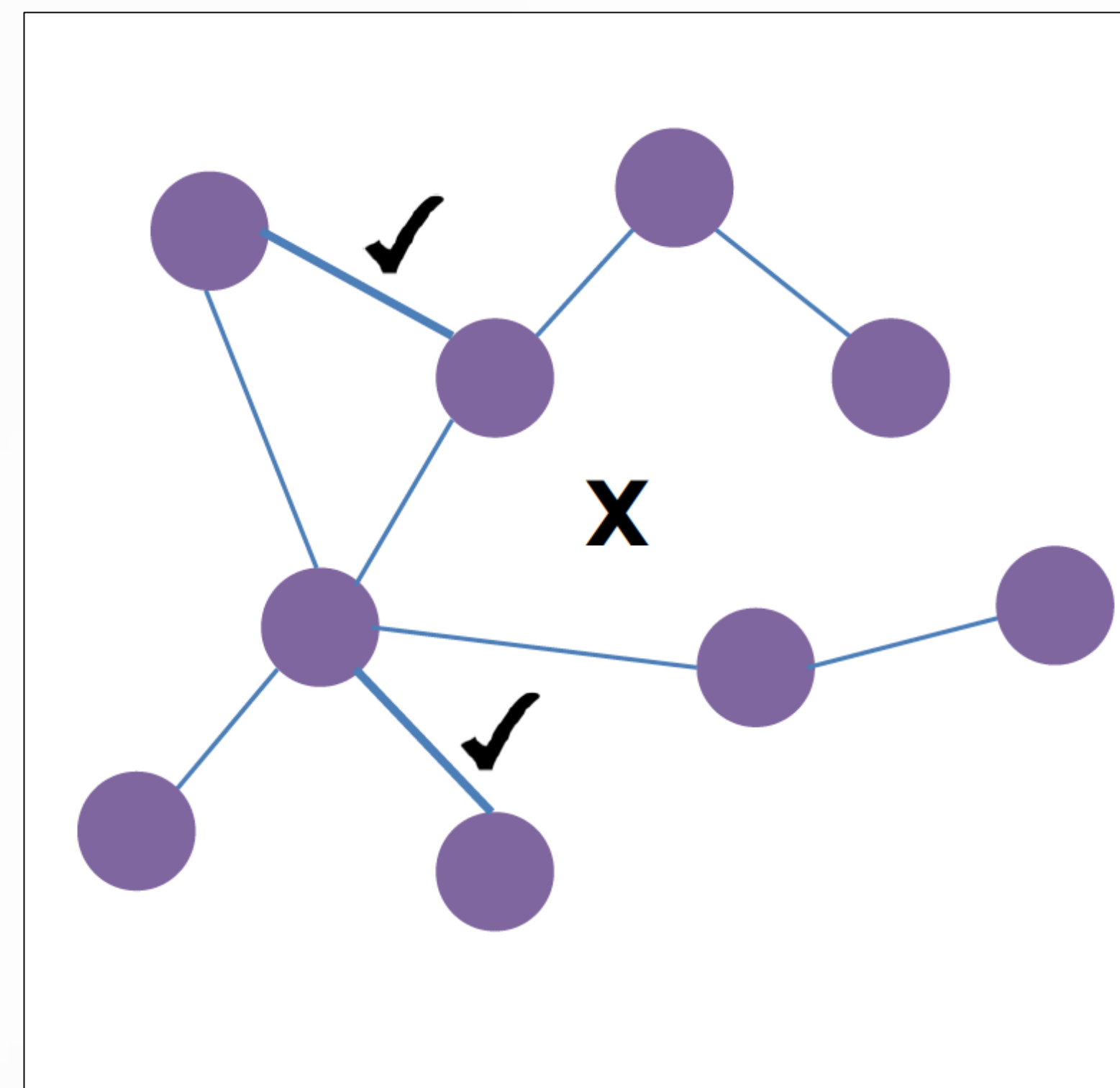
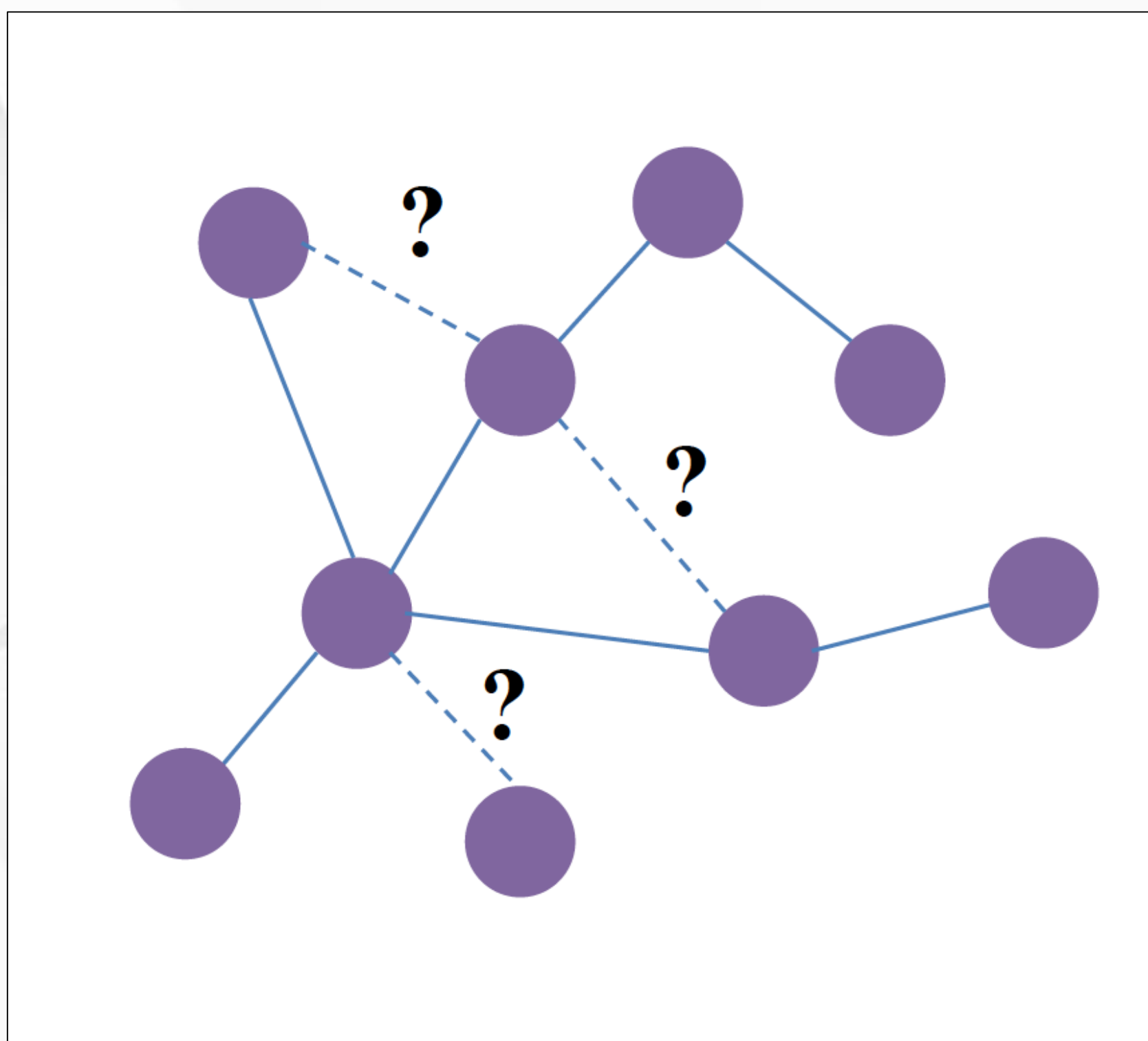


线上支付

GNN的应用 – 分类



GNN的应用 – 关系预测



GNN大规模应用的四大挑战

规模庞大

数百亿甚至数千亿点，
数千亿甚至数万亿的
边

点边异构

同构、异构、
多边、多属性

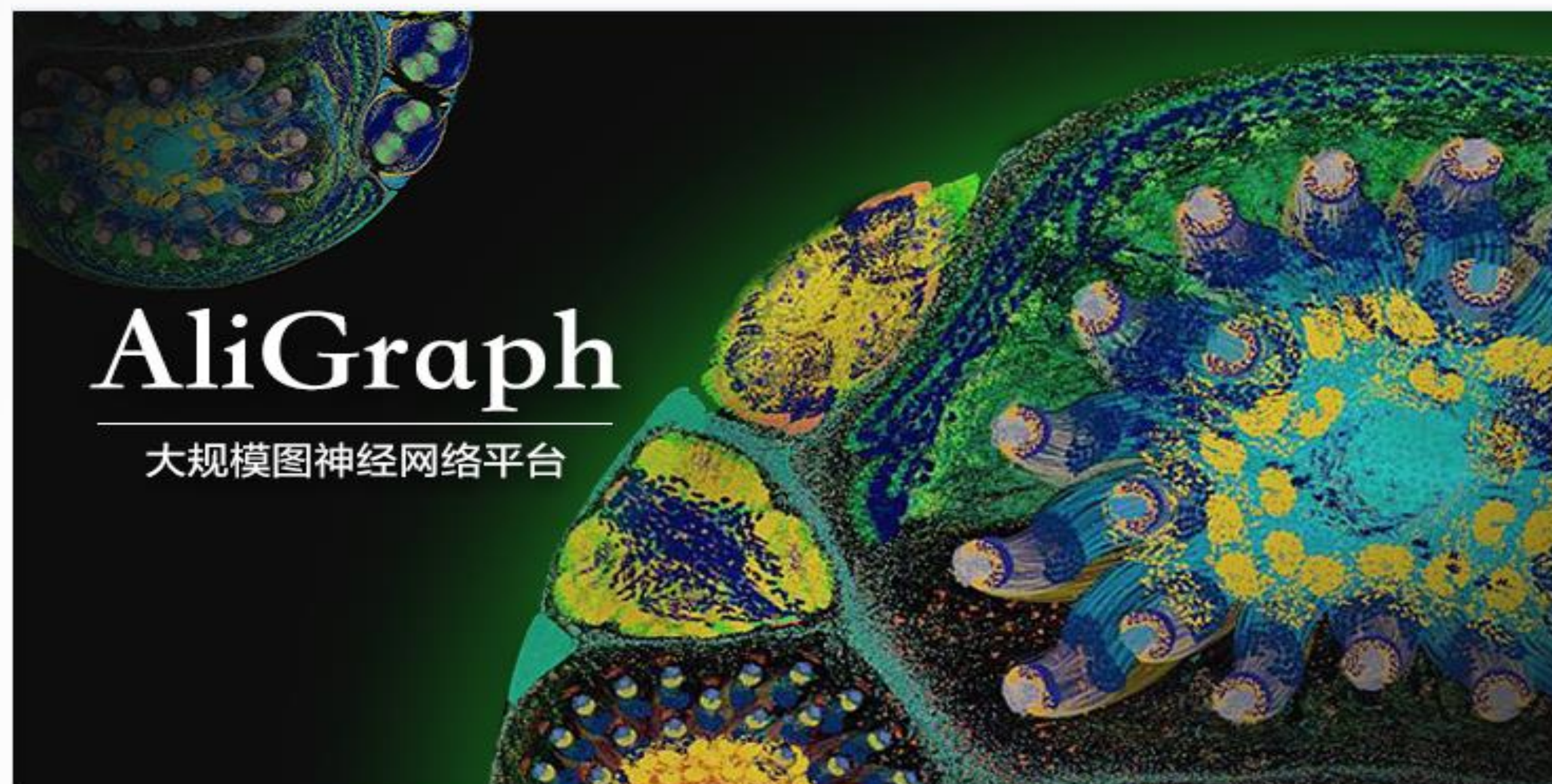
属性丰富

点属性，边属性

动态变化

节点、边的增删更新

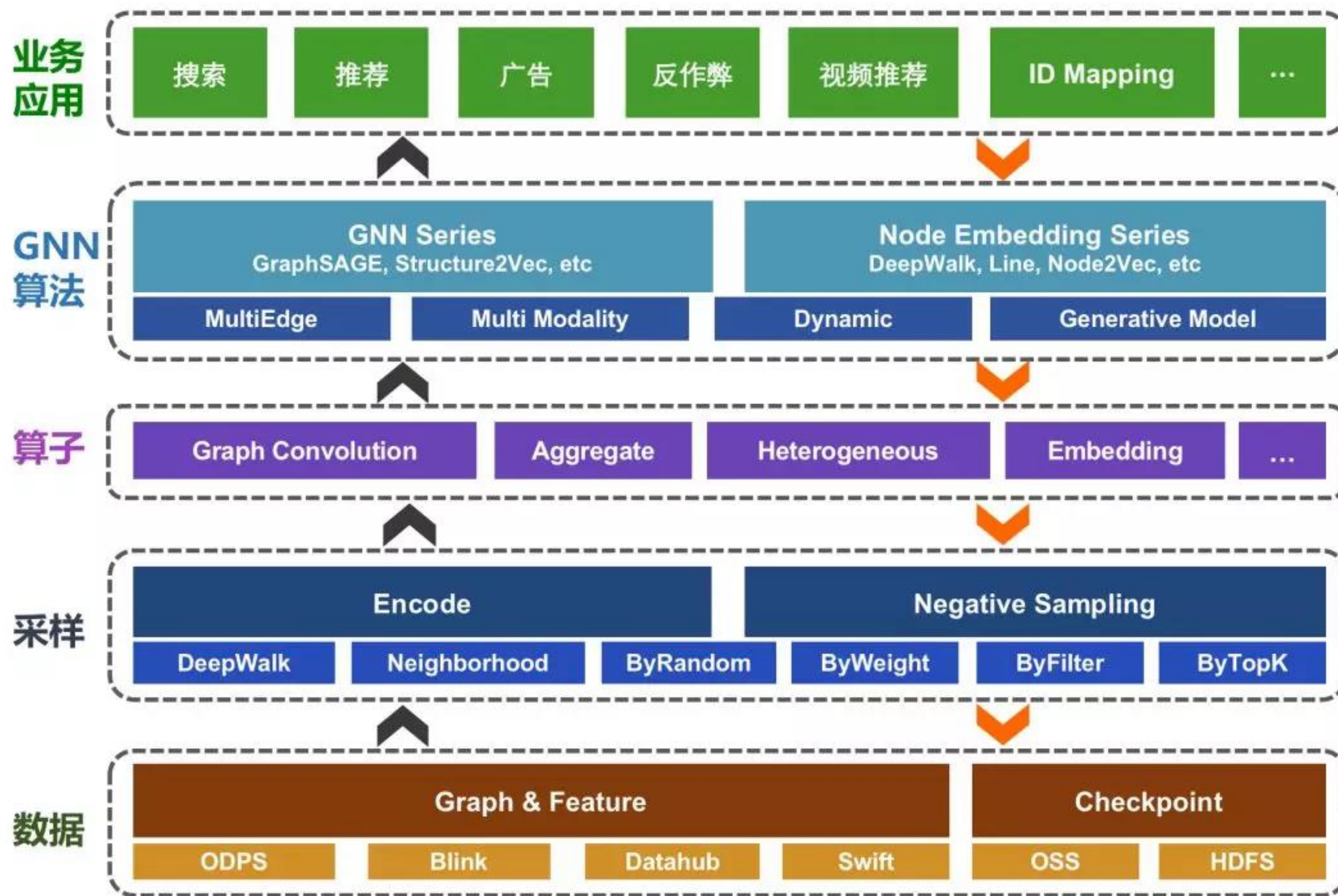
采样&建模&训练一体化 的GNN平台

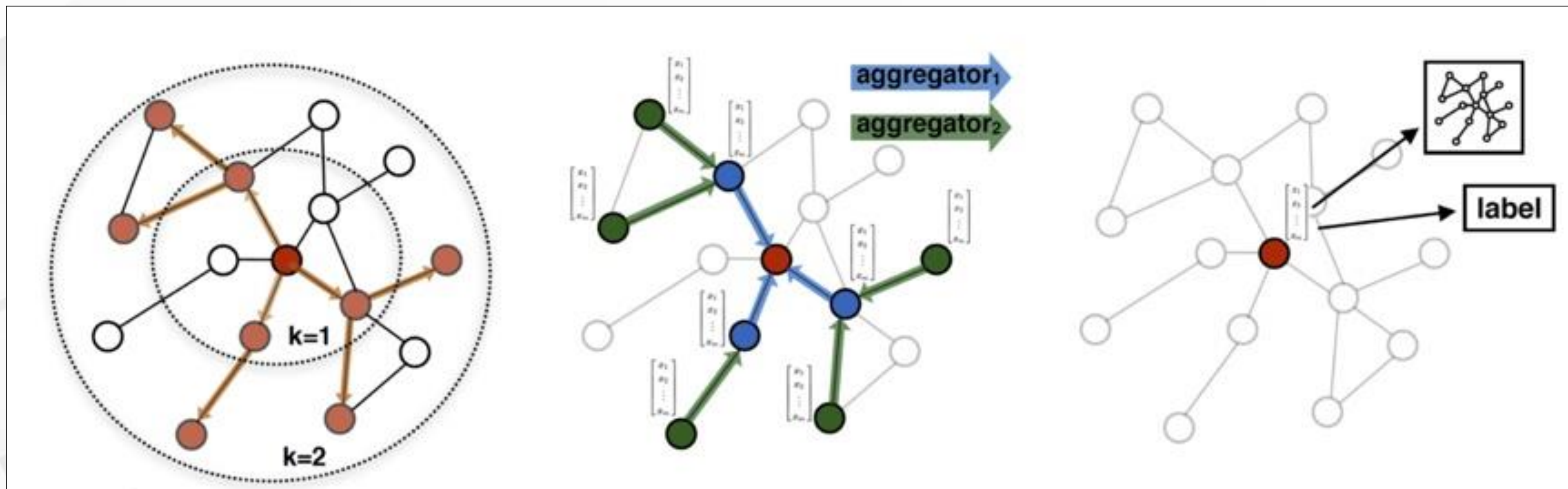


AliGraph：大规模图神经网络平台

AliGraph是新一代大规模图神经网络平台，其认知计算模型比起现有的深度学习技术有了突破性进展，被誉为人工智能2.0模型。

AliGraph系统架构





sample

aggregate

combine

Algorithm 1: GNN Framework

Input: network \mathcal{G} , embedding dimension $d \in \mathbb{N}$, a vertex feature \mathbf{x}_v for each vertex $v \in \mathcal{V}$ and the maximum hops of neighbors $k_{max} \in \mathbb{N}$.

Output: embedding result \mathbf{h}_v of each vertex $v \in \mathcal{V}$

```
1  $\mathbf{h}_v^{(0)} \leftarrow \mathbf{x}_v$ 
2 for  $k \leftarrow 1$  to  $k_{max}$  do
3   for each vertex  $v \in \mathcal{V}$  do
4      $S_v \leftarrow \text{SAMPLE}(\text{Nb}(v))$ 
5      $\mathbf{h}'_v \leftarrow \text{AGGREGATE}(\mathbf{h}_u^{(k-1)}, \forall u \in S)$ 
6      $\mathbf{h}_v^{(k)} \leftarrow \text{COMBINE}(\mathbf{h}_v^{(k-1)}, \mathbf{h}'_v)$ 
7   normalize all embedding vectors  $\mathbf{h}_v^{(k)}$  for all  $v \in \mathcal{V}$ 
8  $\mathbf{h}_v \leftarrow \mathbf{h}_v^{(k_{max})}$  for all  $v \in \mathcal{V}$  return  $\mathbf{h}_v$  as the embedding result for all  $v \in \mathcal{V}$ 
```

AliGraph五大特点

大规模
图存储

分布式
采样

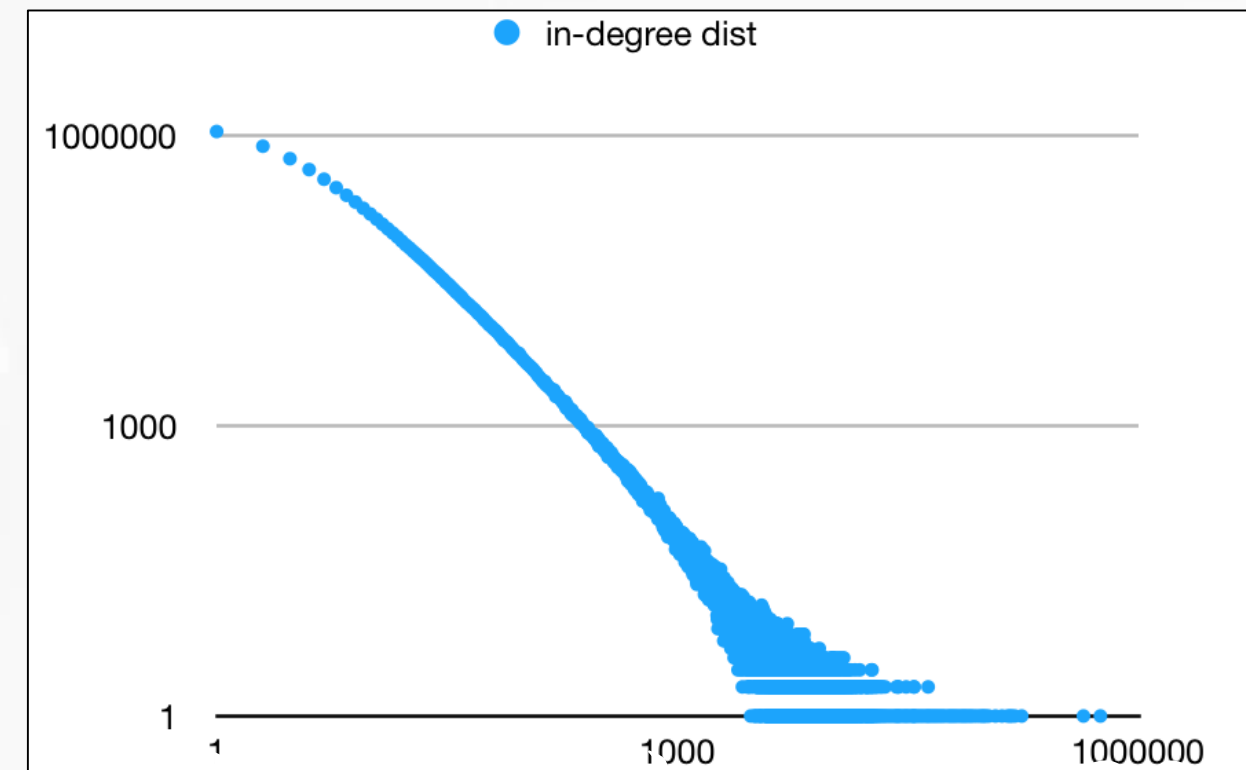
稀疏内
核优化

自创
SGCN

线性
扩展

- 分布式的图存储
- 支持百亿点的规模，可伸缩
- 基于Vertex Cut的分片
- Worker基于出入度的缓存优化

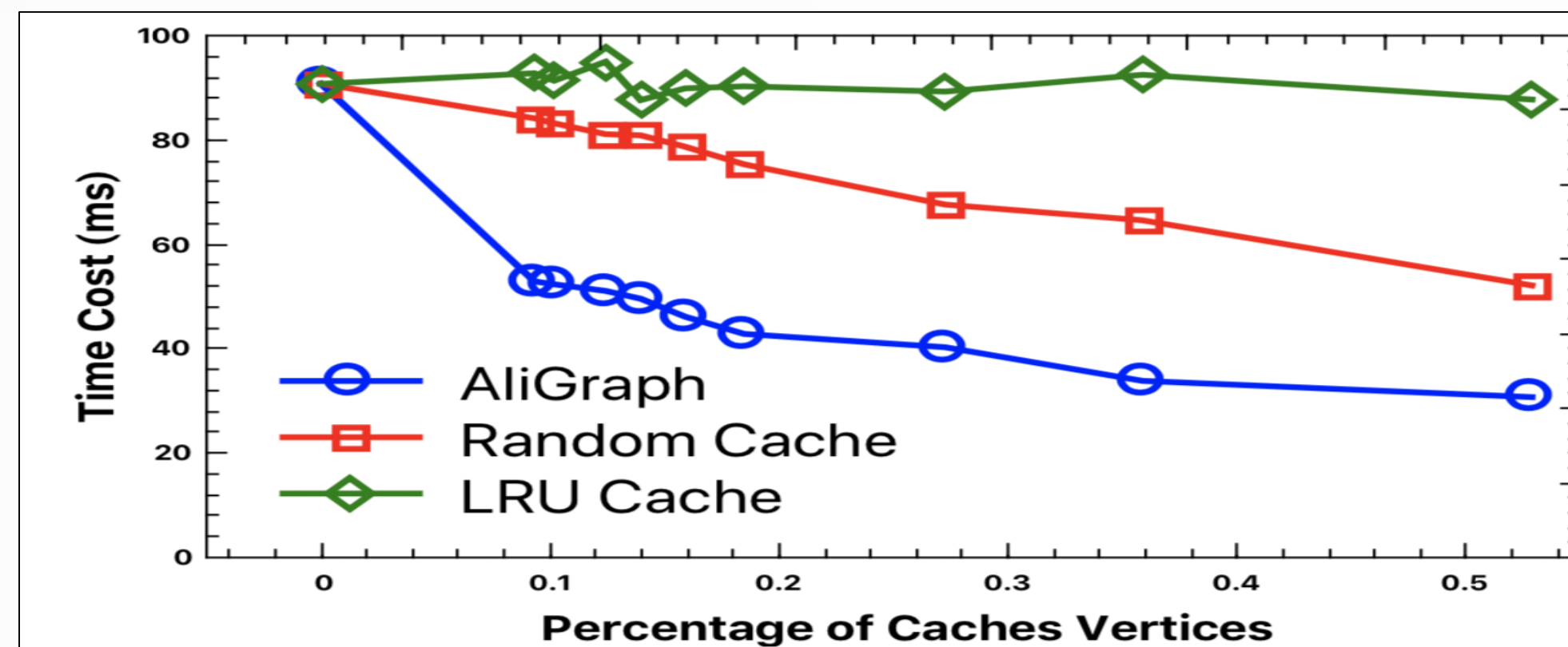
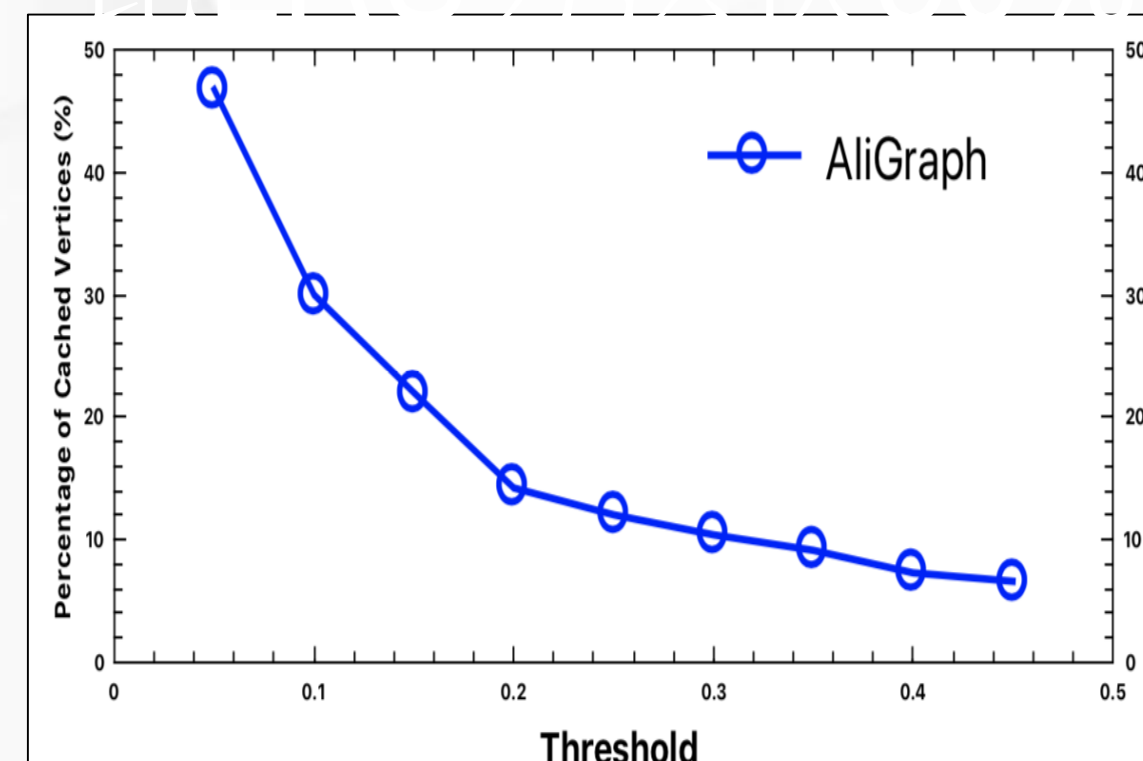
缓存效果



Input: graph \mathcal{G} , partition number p , cache depth h , threshold $\tau_1, \tau_2, \dots, \tau_h$
Output: p subgraphs

```
1 Initialize  $p$  graph servers
2 for each edge  $e = (u, v) \in \mathcal{E}$  do
3    $j = \text{ASSIGN}(u)$ 
4   Send edge  $e$  to the  $j$ -th partition
5 for each vertex  $v \in V$  do
6   for  $k \leftarrow 1$  to  $h$  do
7     Compute  $D_i^{(k)}(v)$  and  $D_o^{(k)}(v)$ 
8     if  $\frac{D_i^{(k)}(v)}{D_o^{(k)}(v)} \geq \tau_k$  then
9       Cache the 1 to  $k$ -hop out-neighbors of  $v$  on each partition where  $v$  exists
```

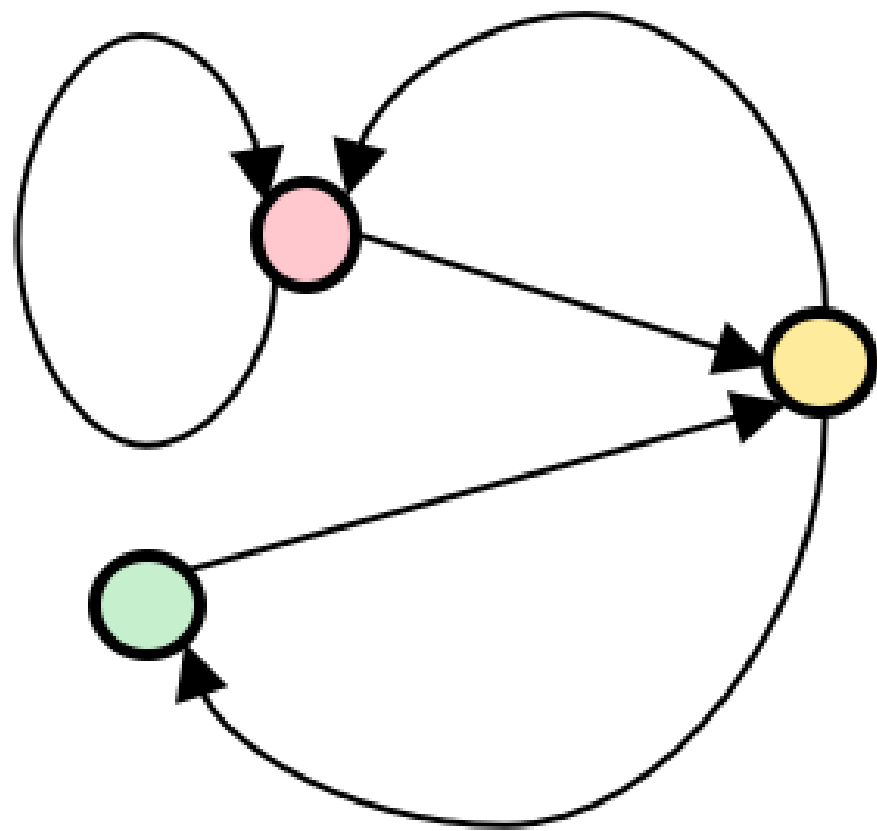
缓存加速：比随机方法快40%-50%，
比LRU方法快50%-60%



缓存加速：比随机方法快40%-50%，比LRU方法快50%-60%

- 支持丰富的采样策略
- 支持多跳的采样功能
- 毫秒级的采样性能
- 模块化的采样设计
- 本地/全局的负采样

图连接关系:

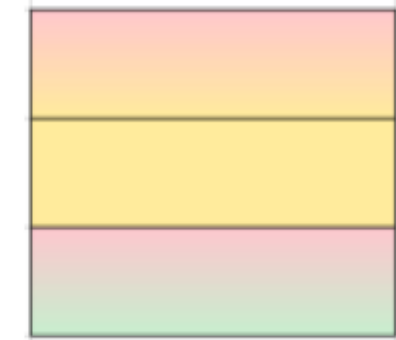
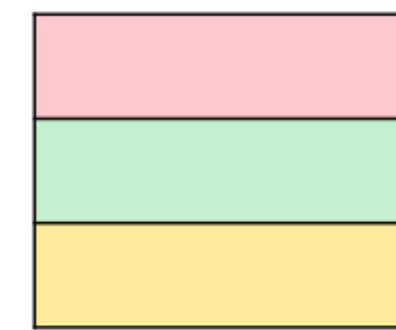


图矩阵表示:

V

1	0	1
0	0	1
1	1	0

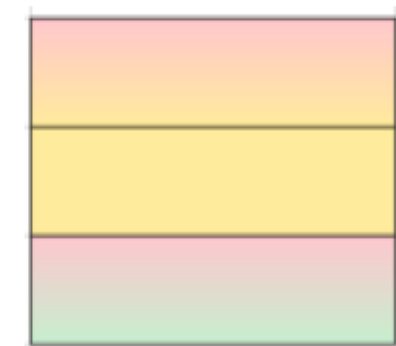
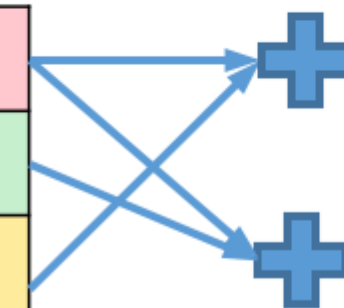
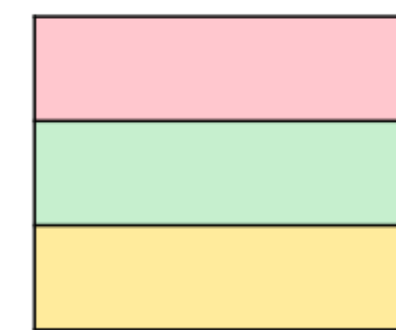
V



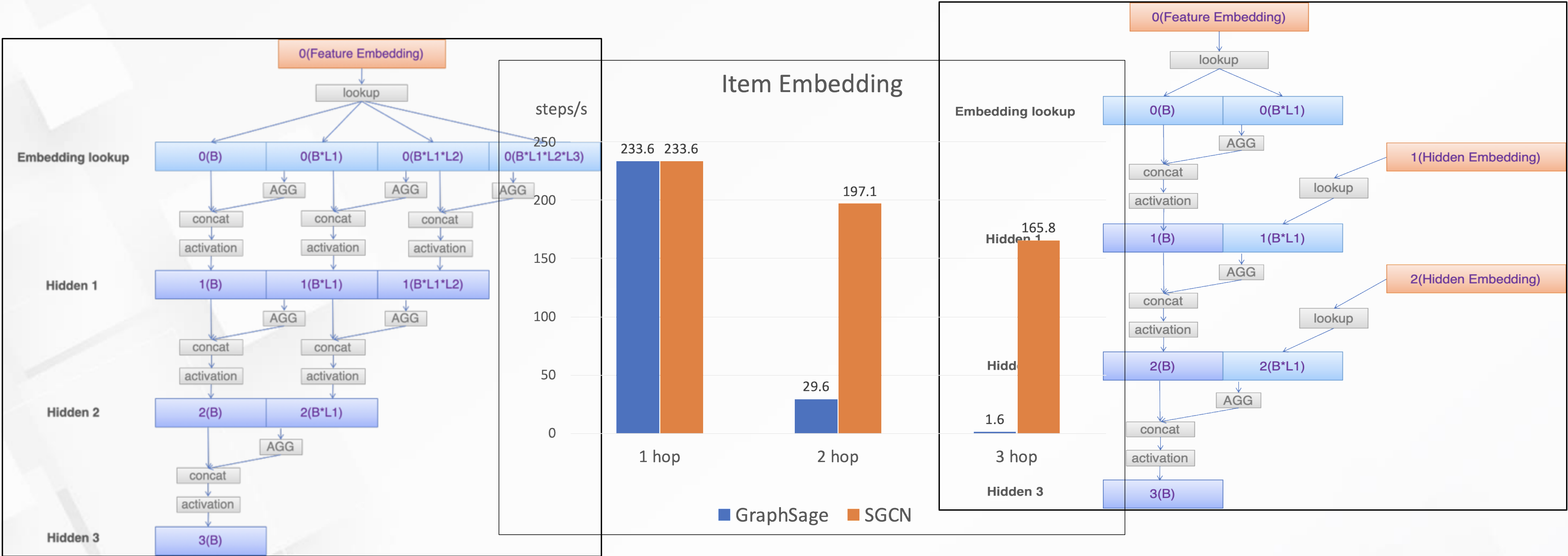
V

1	0	1
0	0	1
1	1	0

V



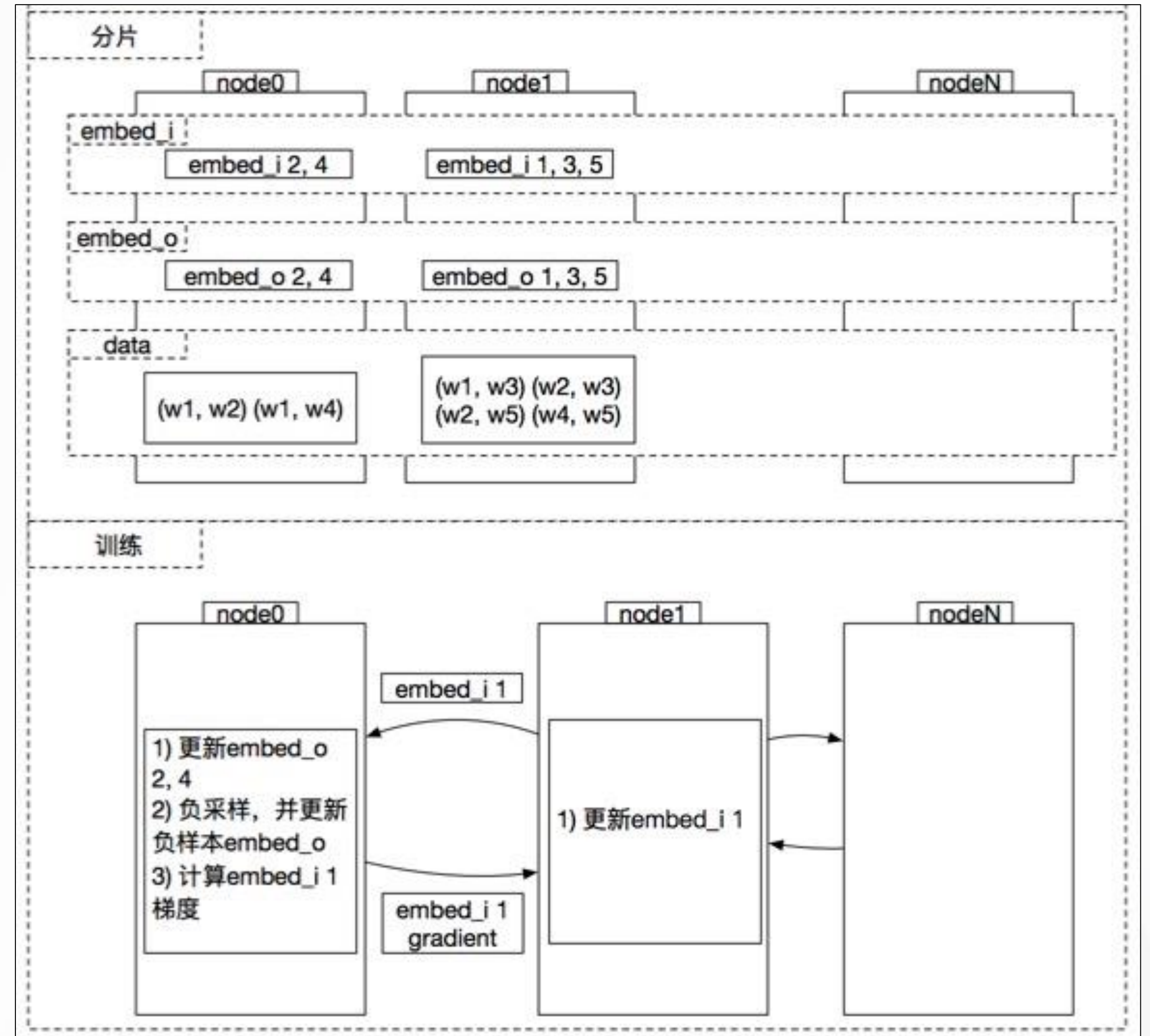
加速原始版本GCN达3X!



去中心化

Distribute Negative Sampling

多机异步训练



以点为中心

Field	Type	Label	Comment
node_id	bigint		
feature	string		

边表混合

Field	Type	Label	Comment
src_id	bigint		
src_type	string		
dst_id	bigint		
dst_type	string		
edge_type	string		
val	double		

采样&建模&训练一体

```
g = ss.Graph()  
    .add_nodes("item")  
    .sample("topk")  
    .batch(128)  
    .model(GraphSage())  
  
g.train()
```

如何使用

机器学习PAI交流5群

71人



扫一扫群二维码，立刻加入该群。

搜索

常用组件

保存的分组

源 / 目标

数据预处理

特征工程

统计分析

机器学习

深度学习

二部图GraphSage...

MXNet

TensorFlow

格式转换组件

Caffe

强化学习

时间序列

文本分析

网络分析

工具

金融板块

废弃栏(15天后会下线)

gnn

运行 部署 Auto ML

二部图GraphSage...

字段设置

参数设置

执行调优 体验新版

evaluation轮数

200

encoding schema

u-i-u

中间层embedding dimension

256

最后层embedding dimension

64

user侧邻居编码

20

item侧邻居编码

20

user特征个数

1

item特征个数

1

user总数

二部图GraphSage嵌入算法



阿里云开发者社区

扫码加入社群
与志同道合的码友一起
Code Up

大数据计算开发者...



该群属于“阿里云ACE”部门群，仅组织内部成员可以加入，如果组织外部人员收到此分享，需要先申请加入该组织。

谢谢！