

# 【动手做】 用DGL库和 GraphSAGE完成链接预测

# 讨论内容

## 1. 工具准备

- Python和Anaconda
- DGL-Deep Graph Library

## 2. 知识准备

- 链接预测
- GraphSAGE的原理

## 3. 案例演示

- 代码运行流程
- 代码演示

## 4. 参考资料

AchieveFun

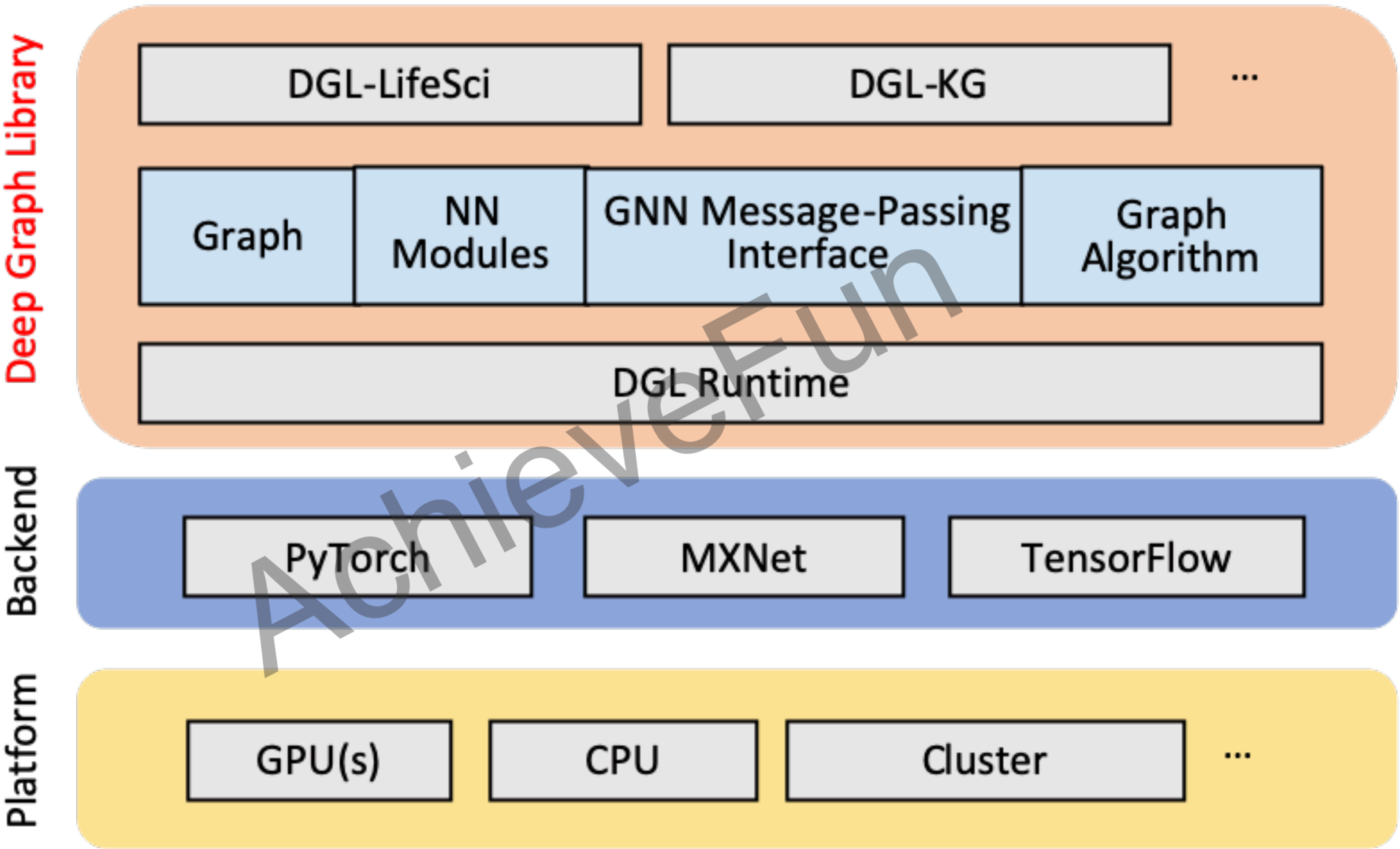
# 1 工具准备：Python和Anaconda

- [illegible]

# 1 工具准备： DGL-Deep Graph Library

- 正式网站 <https://www.dgl.ai/>
- GitHub <https://github.com/dmlc/dgl/>
- [Overview of Deep Graph Library \(DGL\)](#)
- 开始 <https://www.dgl.ai/pages/start.html>
- Deep Graph Library (DGL) 是一个Python软件包，用于在现有DL框架（当前支持PyTorch，MXNet和TensorFlow）之上轻松实现图神经网络模型系列。它提供消息传递的通用控制，通过自动分批处理和高度可调的稀疏矩阵内核进行速度优化以及多GPU / CPU训练，以缩放到数亿个节点和边缘的图形。

# DGL架构

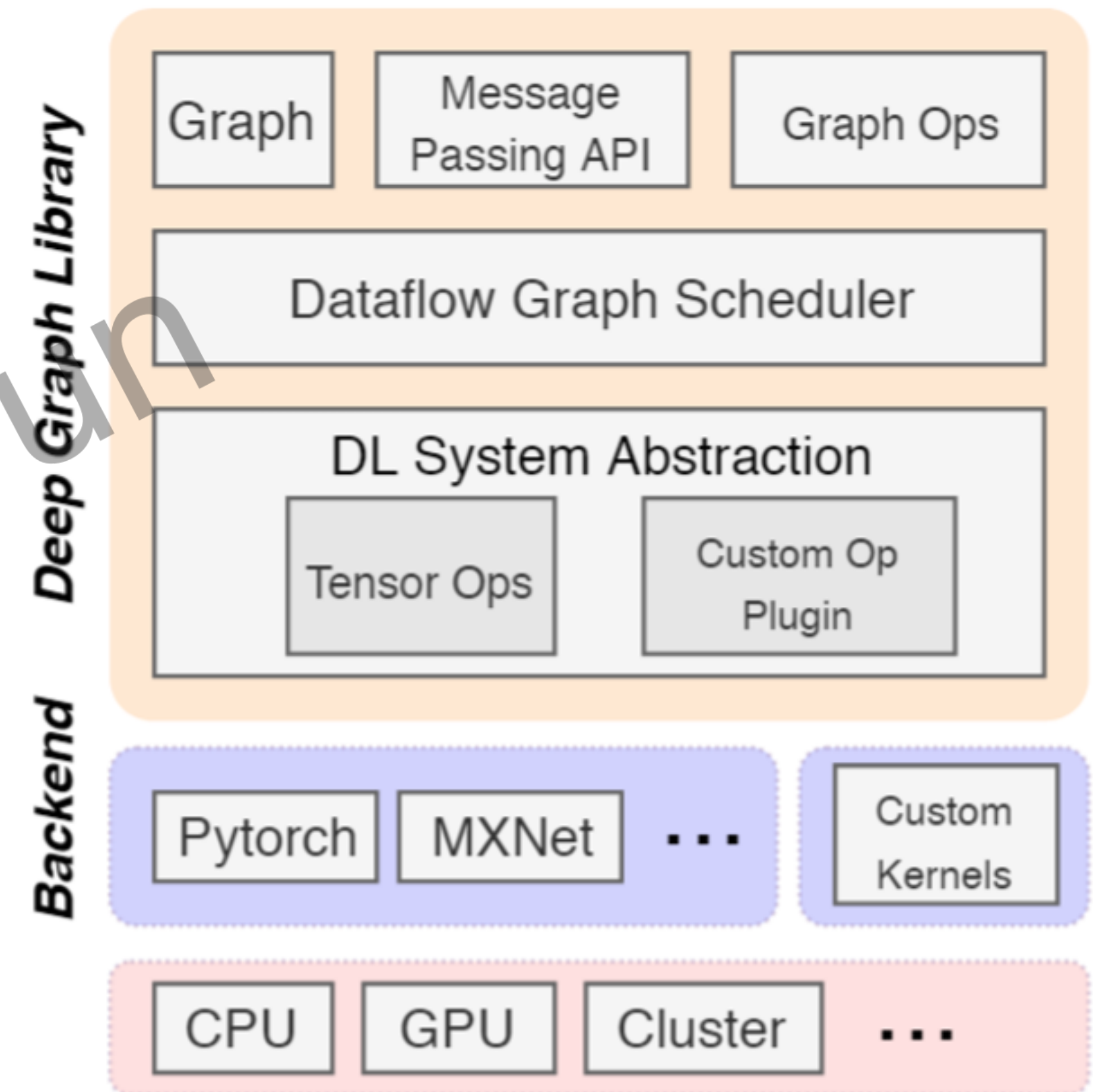


图片来源: [Overview of Deep Graph Library \(DGL\)](#)  Watermarkly



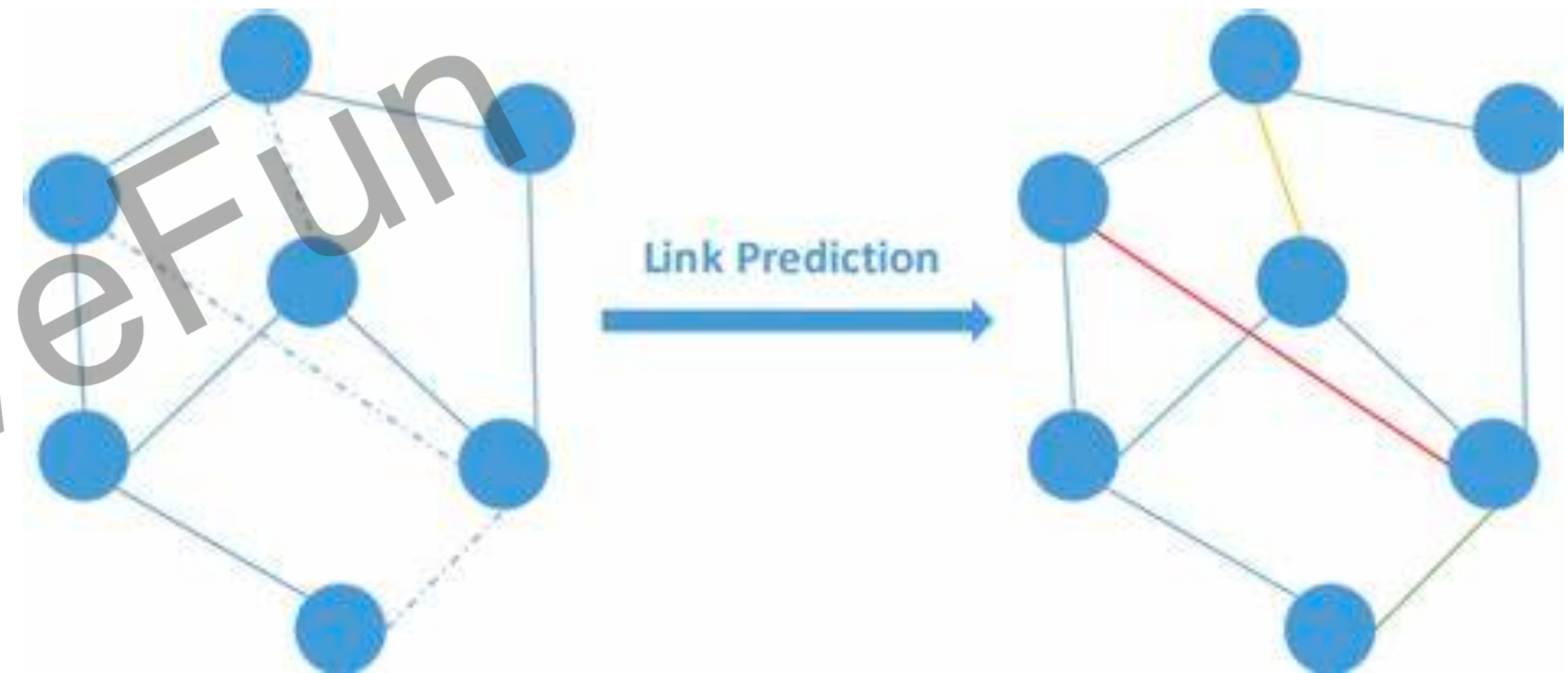
# DGL meta-objectives & architecture

- Forward and backward compatible
  - **Forward**: easy to develop new models
  - **Backward**: seamless integration with existing frameworks (MXNet/Pytorch/Tensorflow)
- **Fast and Scalable**



## 2 知识准备：链接预测

- 基于已有的知识图谱，预测和推断出实体之间可能存在的关系，或者预测可能存在的实体和关系。预测两个特定节点之间是否存在边。
- 有助于丰富和完善知识图谱的内容，使其更加全面和精确。
- 应用于：推荐系统、社交网络分析、生物信息学、知识图谱补全



图片来源: <https://www.researchgate.net/publication/350110330/figure/fig1/AS:1002329977393155@1615985497381/An-example-of-link-prediction.opm>



## 2 知识准备：GraphSAGE的原理

- 对节点的周围邻居节点进行随机采样，使得图网络模型可以应用到大规模图谱上。
- 基于采样 (**S**ample) 和聚集 (**AggreGatE**) 操作
- 对红色节点采样到其周围的蓝色节点，并进行特征聚合。使用更通用的聚合函数 AGG (Aggregation)：任意地可将一组向量聚合成一个向量的可微分方程。

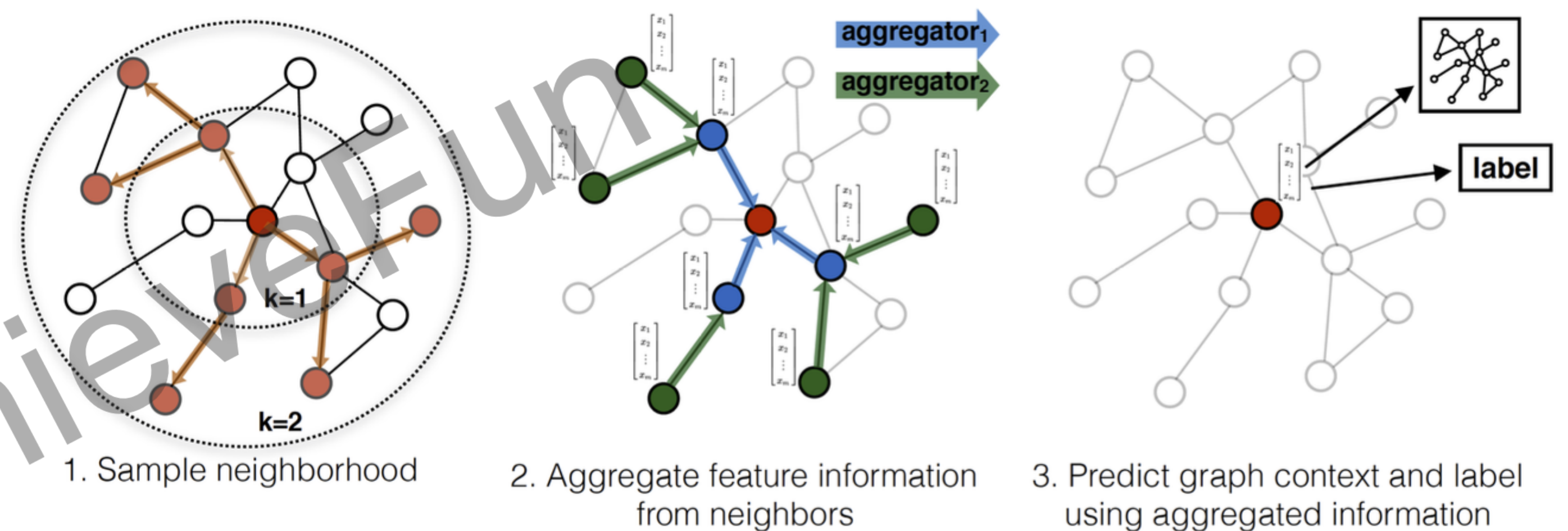


Figure 1: Visual illustration of the GraphSAGE sample and aggregate approach.

- paperswithcode: <https://paperswithcode.com/method/graphsage>
- <https://snap.stanford.edu/graphsage/>

图片来源: <https://arxiv.org/pdf/1706.02216.pdf>



# 2 知识准备： GraphSAGE的聚合函数

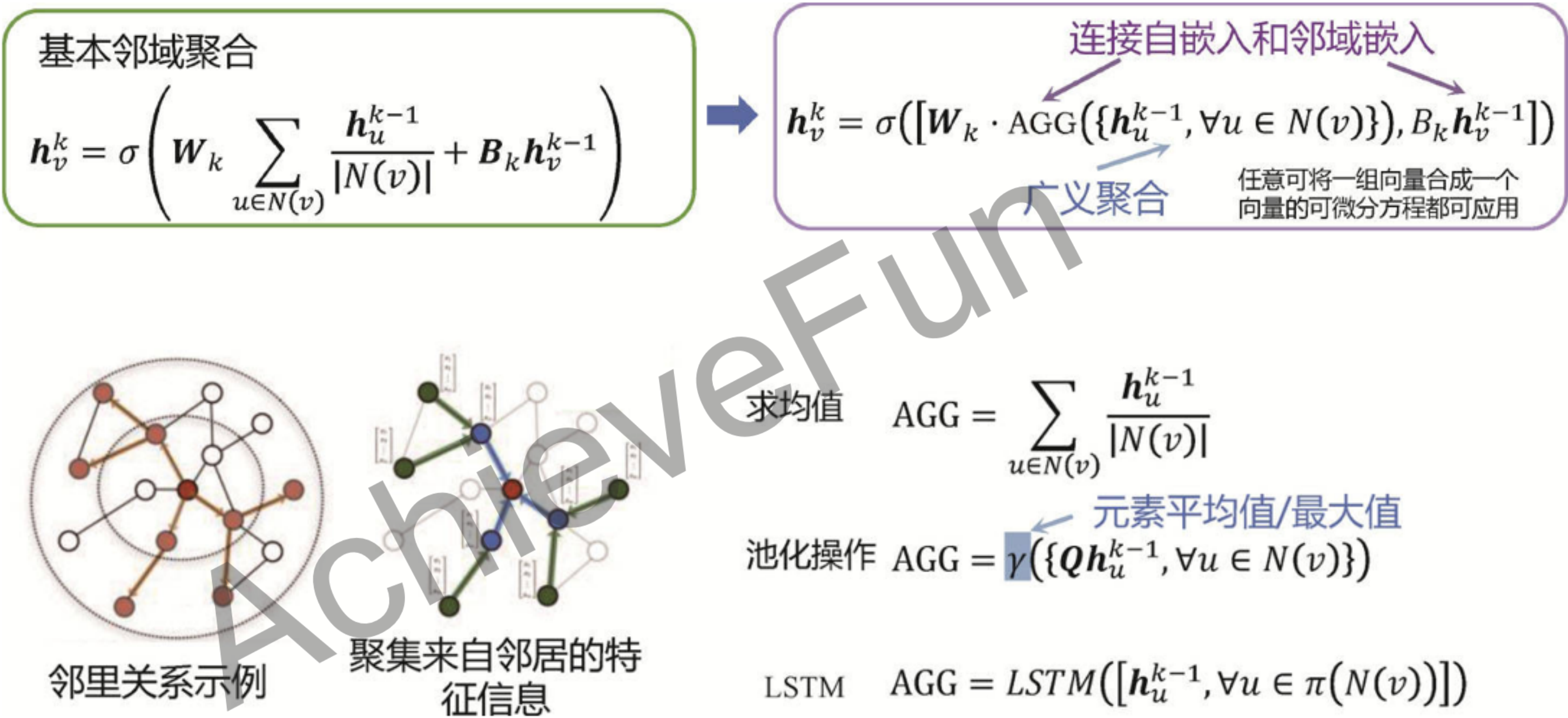


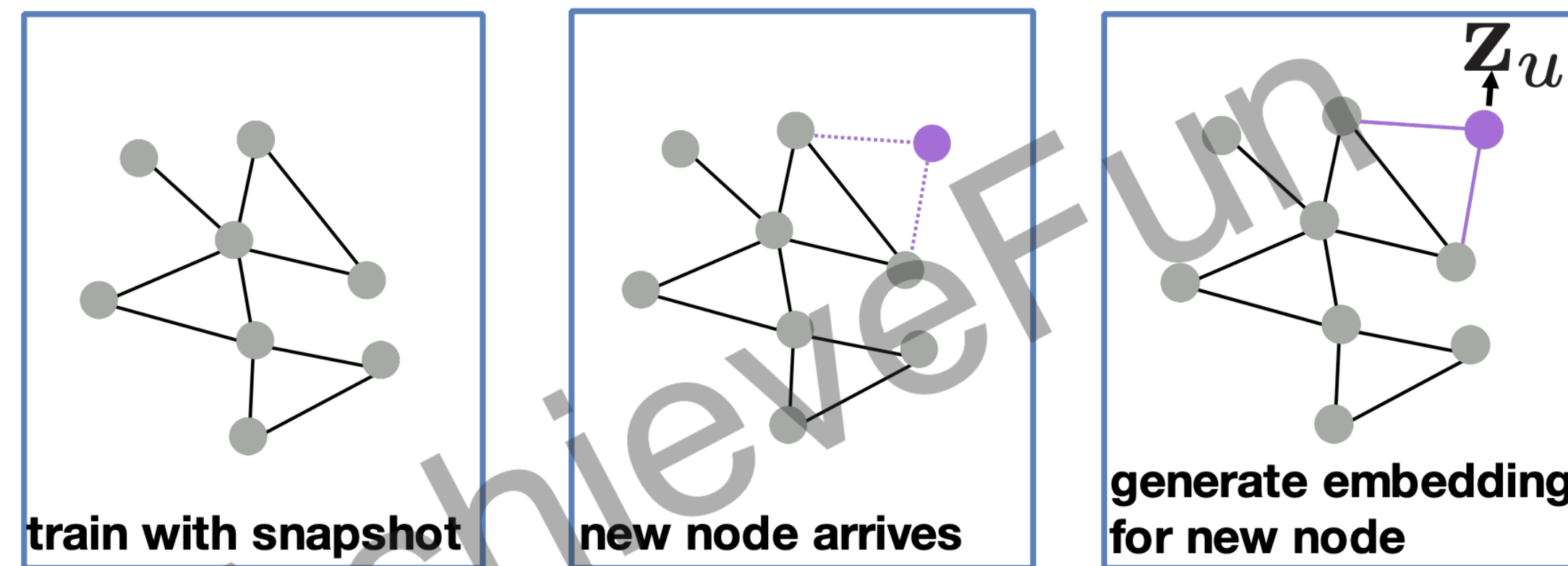
图8-30 GraphSAGE：对邻居节点随机采样

图片来源：陈华钧. 知识图谱导论 (Chinese Edition) (p. 303). Kindle Edition

- 使用三种类型的聚合函数，即：求均值（Mean）、池化操作以及使用LSTM网络。

## 2 知识准备： GraphSAGE的归纳能力

### Inductive Capability



Many application settings constantly encounter previously unseen nodes.  
e.g., Reddit, YouTube, GoogleScholar, ....

Need to generate new embeddings “on the fly”

Representation Learning on Networks: [snap.stanford.edu/proj/embeddings-www](https://snap.stanford.edu/proj/embeddings-www), WWW 2018

30

图片来源: <https://snap.stanford.edu/proj/embeddings-www/files/nrltutorial-part2-gnns.pdf>

模型可以泛化到一些新的、未出现在训练集中的节点。

### 3 案例演示

- 案例来源: [https://docs.dgl.ai/tutorials/blitz/4\\_link\\_predict.html](https://docs.dgl.ai/tutorials/blitz/4_link_predict.html)
- 代码:
  - [https://docs.dgl.ai/downloads/ad3a2962e53be21909260c39376a994a/4\\_link\\_predict.py](https://docs.dgl.ai/downloads/ad3a2962e53be21909260c39376a994a/4_link_predict.py)
  - [https://docs.dgl.ai/downloads/514d96075aeaa53fc7a3d9873b7b963f/4\\_link\\_predict.ipynb](https://docs.dgl.ai/downloads/514d96075aeaa53fc7a3d9873b7b963f/4_link_predict.ipynb)

# 使用GNN进行链接预测的方法

- 例子：预测引用网络中两篇论文之间是否存在引用关系（引用或被引用）。
- 将链接预测问题表述为如下二元分类问题：
  - a. 将图中的边视为正例。
  - b. 抽取一些不存在的边（即节点对之间没有边）作为负例。
  - c. 将正例和负例分为训练集和测试集。
  - d. 用任何二元分类指标（如曲线下面积 (AUC)）评估模型。
- 这种做法来自 `SEAL` <https://papers.nips.cc/paper/2018/file/53f0d7c537d99b3824f0f99d62ea2428-Paper.pdf>，尽管这里的模型没有使用他们的节点标记思想。
- 在某些领域，如大规模推荐系统或信息检索，偏爱强调 top-K 预测的良好性能的指标。在这种情况下，可能需要考虑平均精度等其他指标，并使用其他负抽样方法。



# 环境准备

1. 安装Anaconda <https://www.anaconda.com/download/>
2. 在Anaconda建立环境dgl, 根据dgl文档(<https://www.dgl.ai/pages/start.html>)确定**Python版本**
3. 根据代码, 使用pip install -r requirements.txt 安装依赖库并验证库的存在pip list
  - matplotlib
  - torch
  - scikit-learn
4. 安装dgl <https://www.dgl.ai/pages/start.html>

例如: conda install -c dglteam dgl

# 代码运行流程

1. 加载图形和特征 Loading graph and features
2. 准备训练集和测试集 Prepare training and testing sets <https://graphsandnetworks.com/the-cora-dataset>, <https://github.com/topics/cora-dataset>
3. 选择GraphSAGE 模型 Define a GraphSAGE model
4. 正图、负图和应用网格 Positive graph, negative graph, and apply\_edges
5. 训练模型

## 4 参考资料

- [https://docs.dgl.ai/guide\\_cn/training-link.html](https://docs.dgl.ai/guide_cn/training-link.html)
- [https://docs.dgl.ai/tutorials/blitz/4\\_link\\_predict.html](https://docs.dgl.ai/tutorials/blitz/4_link_predict.html)
- [https://docs.dgl.ai/downloads/ad3a2962e53be21909260c39376a994a/4\\_link\\_predict.py](https://docs.dgl.ai/downloads/ad3a2962e53be21909260c39376a994a/4_link_predict.py)
- [https://docs.dgl.ai/downloads/514d96075aeaa53fc7a3d9873b7b963f/4\\_link\\_predict.ipynb](https://docs.dgl.ai/downloads/514d96075aeaa53fc7a3d9873b7b963f/4_link_predict.ipynb)
- <https://arxiv.org/pdf/1706.02216.pdf>