

**M. BESTA, P. RENC, R. GERSTENBERGER, P. S. LABINI, A. ZIOGAS, T. CHEN, L. GIANINAZZI, F. SCHEIDL, K. SZENES,
A. CARIGET, P. IFF, G. KWASNIEWSKI, R. KANAKAGIRI, C. GE, S. JAEGER, J. WAS, F. VELLA, T. HOEFLER**



High-Performance and Programmable Attentional Graph Neural Networks with Global Tensor Formulations



Graphs are Powerful and Ubiquitous!

Social sciences

Biology

Chemistry

Engineering

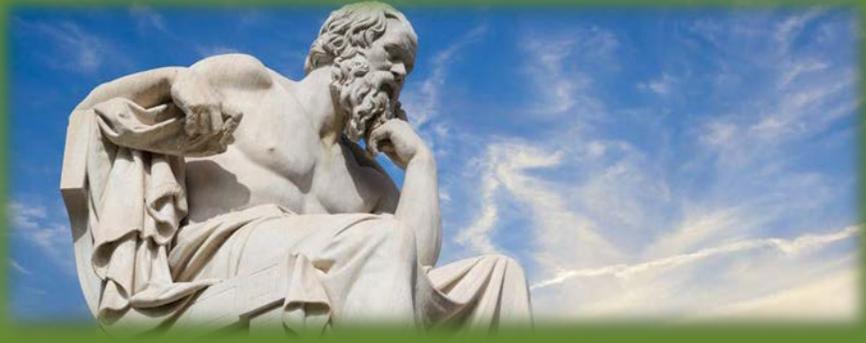
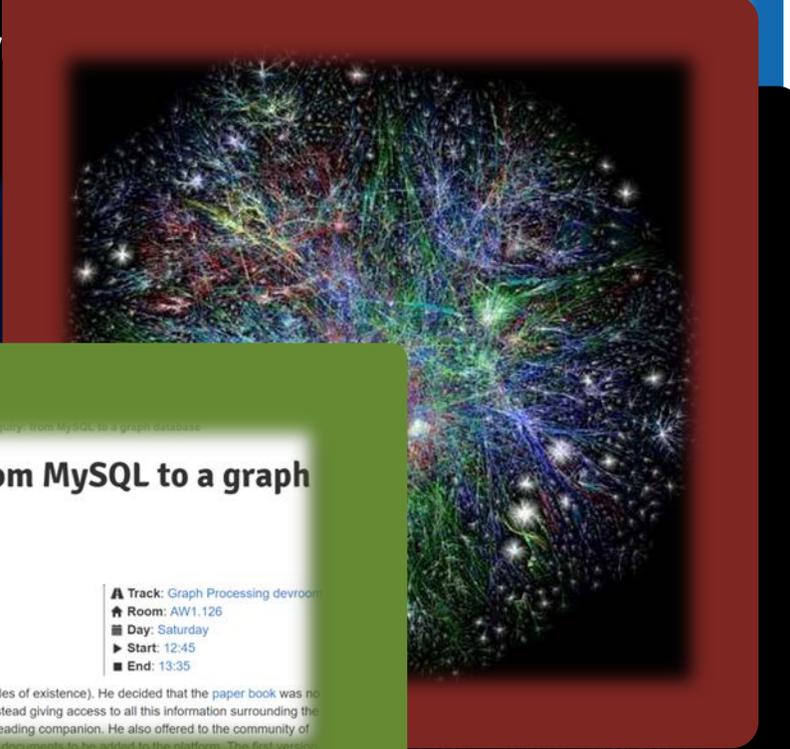
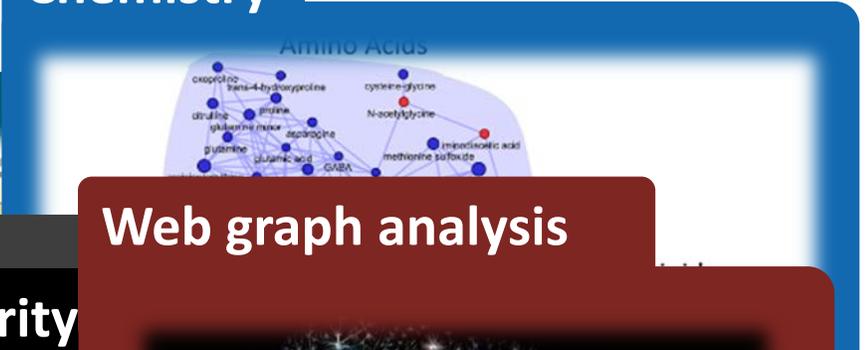
Communication

Web graph analysis

Medicine

Cybersecurity

...even philosophy 😊



FOSDEM 2016 / Schedule / Events / Developer rooms / Graph Processing / Modeling a Philosophical Inquiry: from MySQL to a graph database

Modeling a Philosophical Inquiry: from MySQL to a graph database

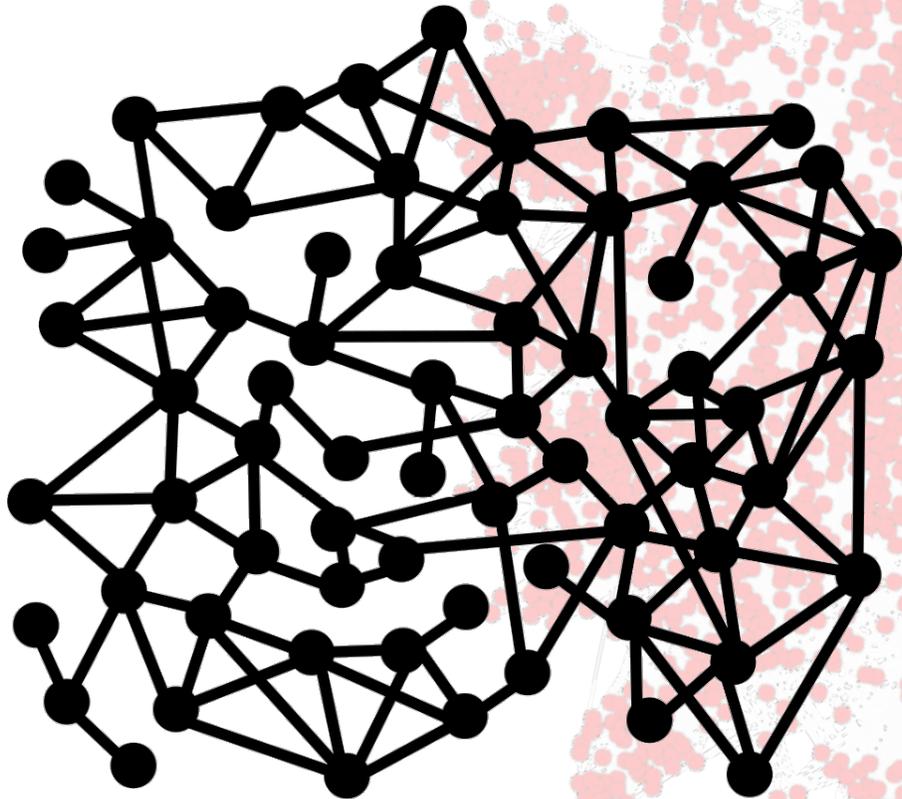
The short story of a long refactoring process

- Track: Graph Processing devroom
- Room: AW1.126
- Day: Saturday
- Start: 12:45
- End: 13:35

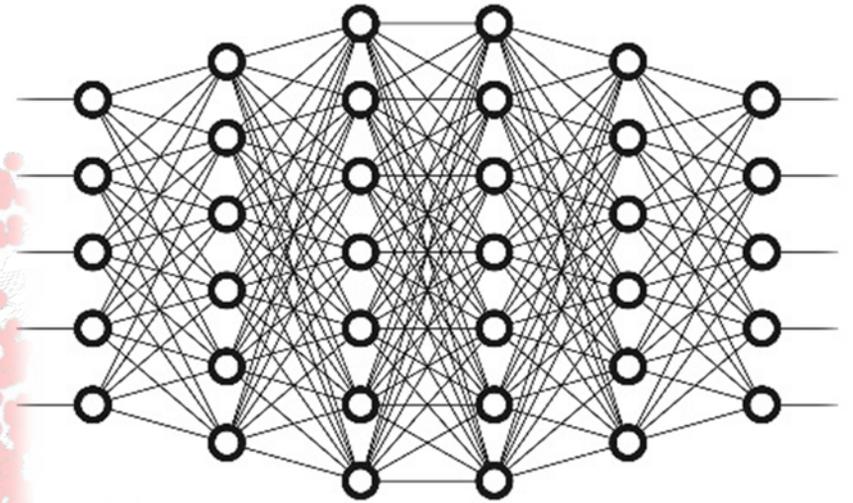
Bruno Latour wrote a book about philosophy (an inquiry into modes of existence). He decided that the paper book was no place for the numerous footnotes, documentation or glossary, instead giving access to all this information surrounding the book through a web application which would present itself as a reading companion. He also offered to the community of readers to submit their contributions to his inquiry by writing new documents to be added to the platform. The first version...

Graphs + Deep Learning = Graphs Neural Networks (GNNs)

In the last 5 years, learning on graphs exploded



+



Let's See Some Recent Success Stories of GNNs

Article

A graph placement methodology for fast chip design

<https://doi.org/10.1038/s41586-021-03...>

Received: 3 November 2020

Accepted: 13 April 2021

Published online: 9 June 2021

 Check for updates

Article

Advancing mathematics by guiding human intuition with AI

<https://doi.org/10.1038/s41586-021-04086-x>

Received: 10 July 2021

Accepted: 30 September 2021

Published online: 1 December 2021

Open access

Alex Da
Nenad
Marc La

The pra
formul

Article

Highly accurate protein structure prediction with AlphaFold

<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

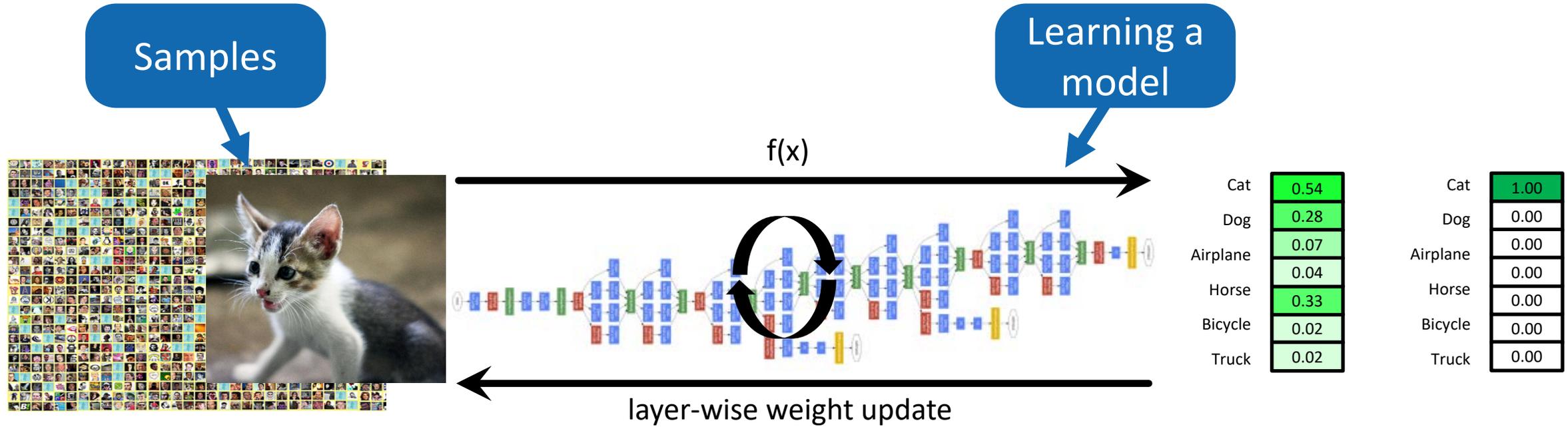
Open access

 Check for updates

John Jumper^{1,4}, Richard Evans^{1,4}, Alexander Pritzel^{1,4}, Tim Green^{1,4}, Michael Figurnov^{1,4}, Olaf Ronneberger^{1,4}, Kathryn Tunyasuvunakool^{1,4}, Russ Bates^{1,4}, Augustin Židek^{1,4}, Anna Potapenko^{1,4}, Alex Bridgland^{1,4}, Clemens Meyer^{1,4}, Simon A. A. Kohl^{1,4}, Andrew J. Ballard^{1,4}, Andrew Cowie^{1,4}, Bernardino Romera-Paredes^{1,4}, Stanislav Nikolov^{1,4}, Rishub Jain^{1,4}, Jonas Adler¹, Trevor Back¹, Stig Petersen¹, David Reiman¹, Ellen Clancy¹, Michal Zielinski¹, Martin Steinegger^{2,3}, Michalina Pacholska¹, Tamas Berghammer¹, Sebastian Bodenstein¹, David Silver¹, Oriol Vinyals¹, Andrew W. Senior¹, Koray Kavukcuoglu¹, Pushmeet Kohli¹ & Demis Hassabis^{1,4}

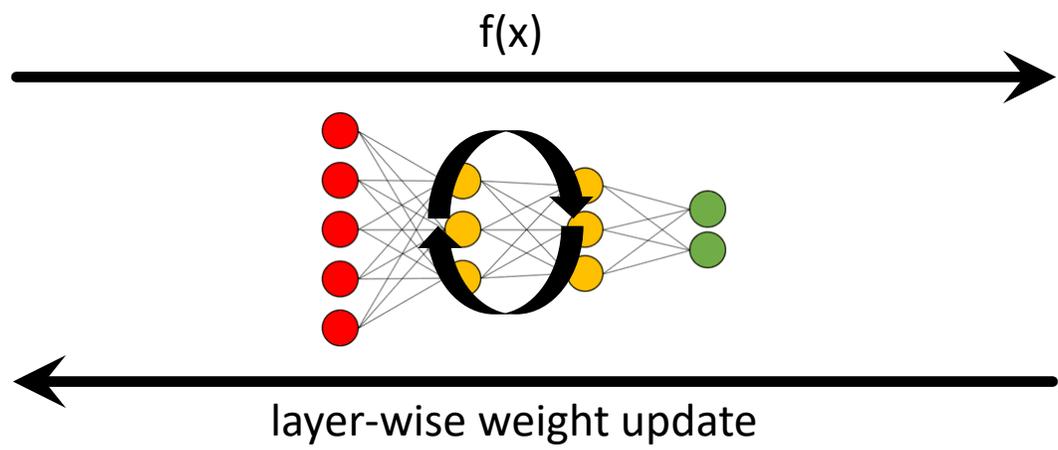
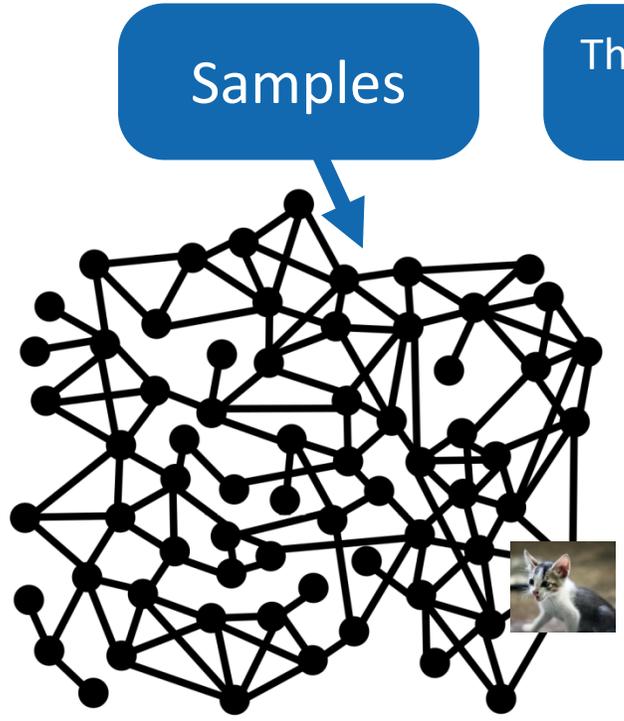
Proteins are essential to life, and understanding their structure can facilitate a

Deep Learning (DL) in a Nutshell

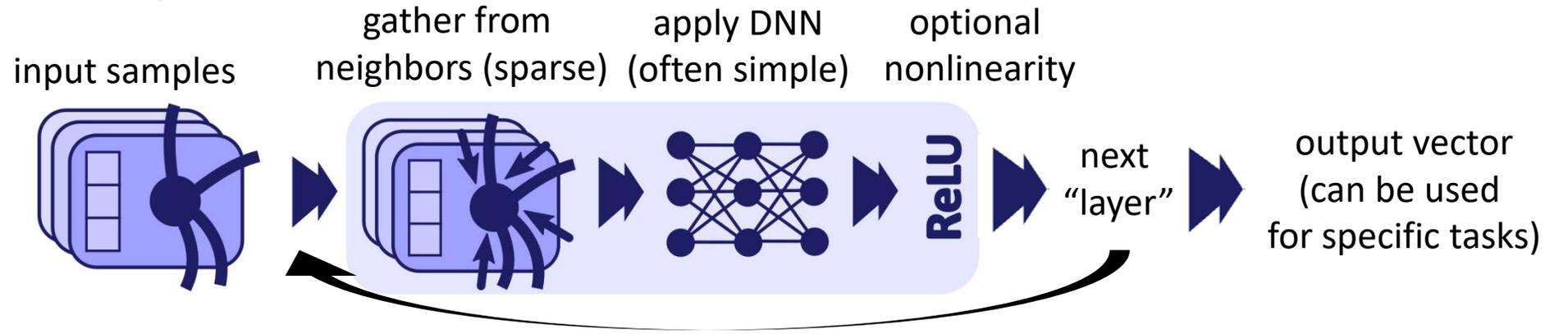


A Primer on Graph Neural Networks (GNNs)

These could still be photos, but now forming **explicit relations**, e.g., two photos are related if they were taken at the same place.



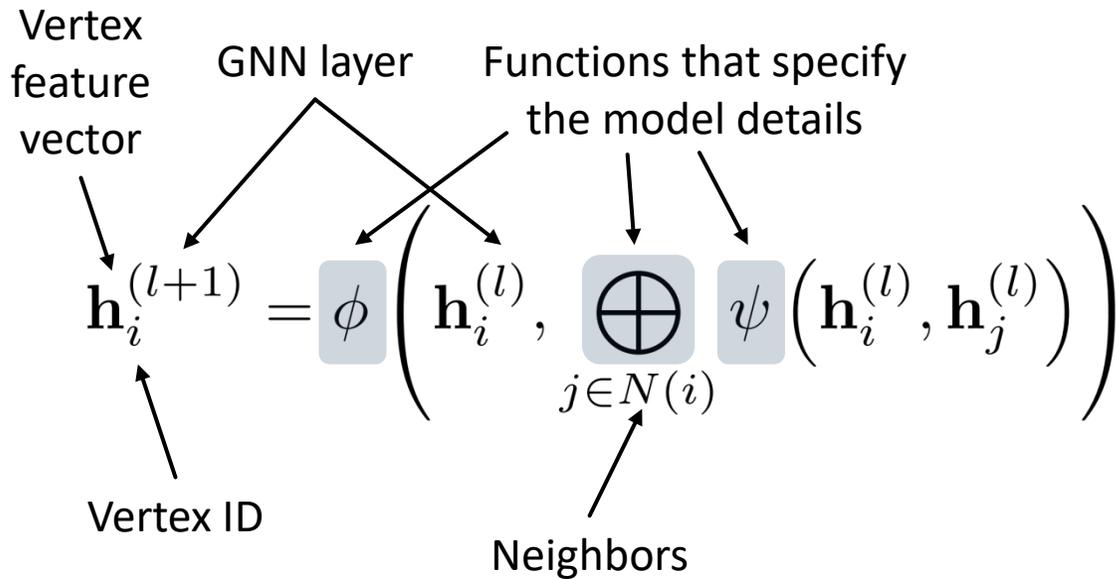
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.04	Horse	0.00
Bicycle	0.33	Bicycle	0.00
Truck	0.02	Bicycle	0.00
		Truck	0.00



Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges

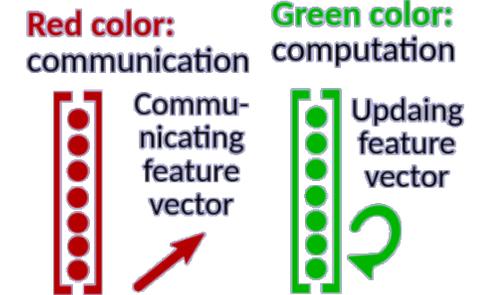


Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$



Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

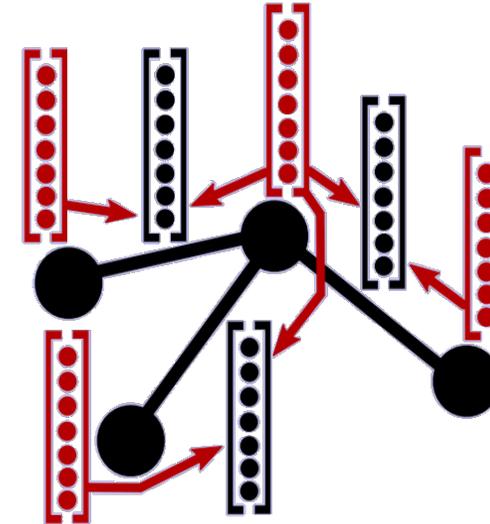
Formulations based on functions operating on single vertices & edges

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

Red color: communication
Green color: computation

Communicating feature vector

Updating feature vector

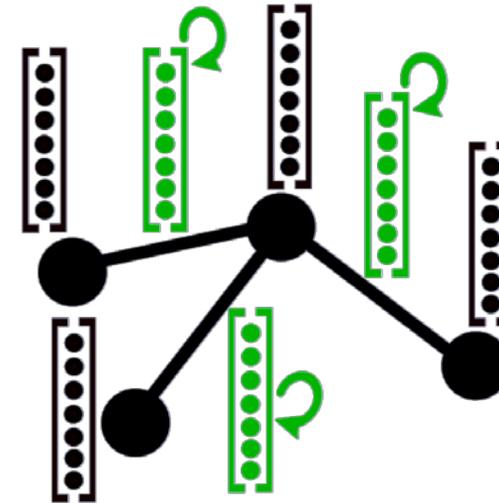
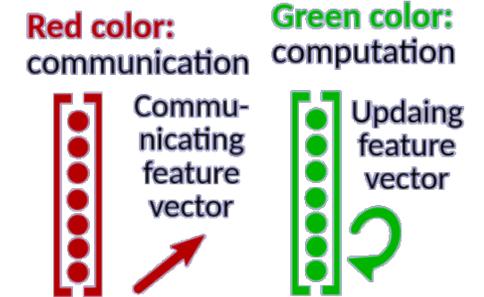


Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

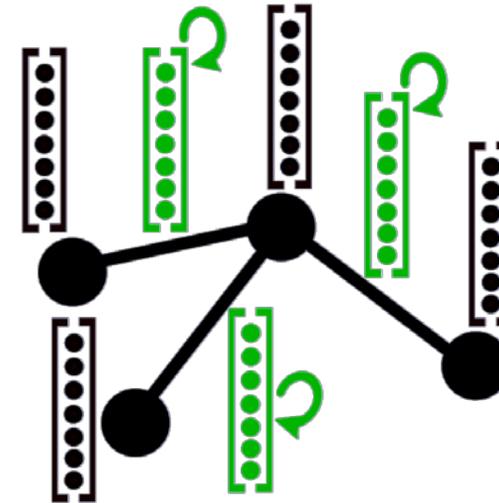
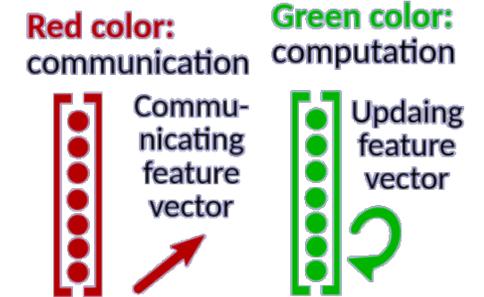


Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in \mathcal{N}(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

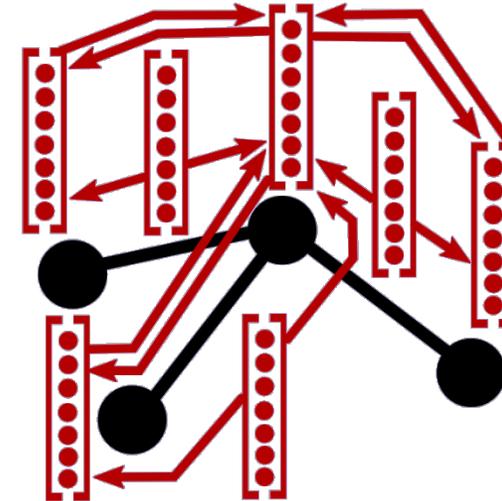
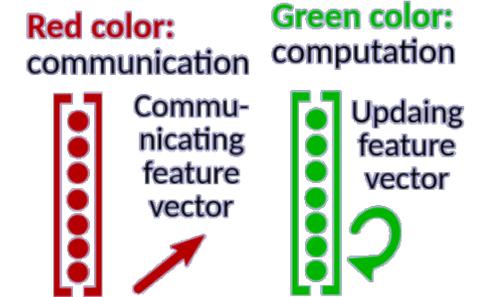


Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in \mathcal{N}(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

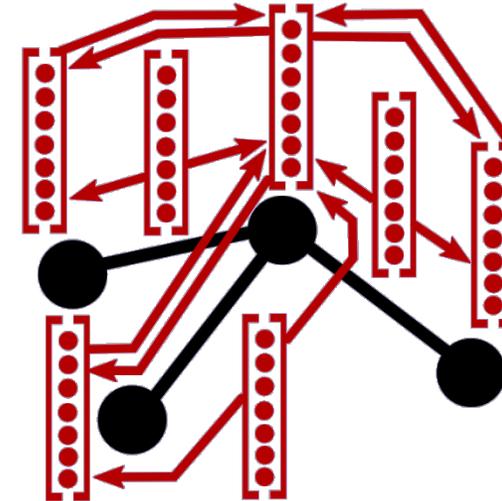
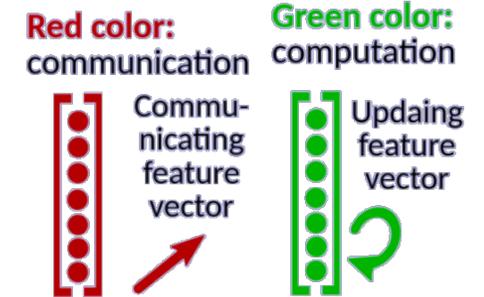


Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

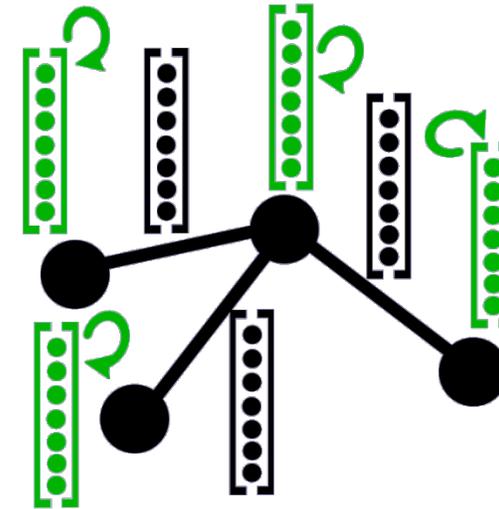
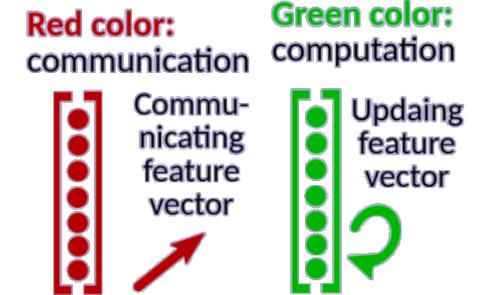


Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges

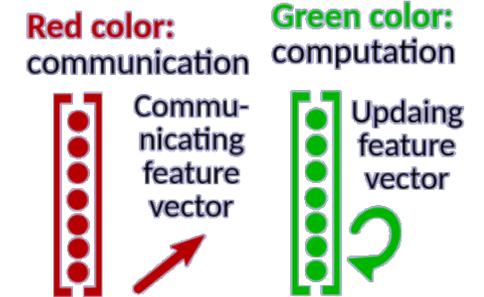
$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$



Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges



$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

Taxonomy of Mathematical Formulations of GNNs

Local GNN formulations

Formulations based on functions operating on single vertices & edges

$$\begin{matrix} \bullet & \bullet & \bullet \\ \vdots & & \\ \bullet & & \end{matrix} \mathbf{h}_{i-1}^{(l+1)} = \phi \left(\begin{matrix} \mathbf{h}_{i-1}^{(l+1)} \oplus \psi \left(\mathbf{h}_{i-1}^{(l)}, \mathbf{h}_j^{(l)} \right) \\ j \in N(i-1) \end{matrix} \right)$$

$$\begin{matrix} \bullet & & \\ \vdots & & \\ \bullet & & \end{matrix} \mathbf{h}_i^{(l+1)} = \phi \left(\begin{matrix} \mathbf{h}_i^{(l)}, \oplus \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \\ j \in N(i) \end{matrix} \right)$$

$$\begin{matrix} \bullet & & \\ \vdots & & \\ \bullet & & \end{matrix} \mathbf{h}_{i+1}^{(l+1)} = \phi \left(\begin{matrix} \mathbf{h}_{i+1}^{(l+1)} \oplus \psi \left(\mathbf{h}_{i+1}^{(l)}, \mathbf{h}_j^{(l)} \right) \\ j \in N(i+1) \\ \bullet & \bullet & \bullet \end{matrix} \right)$$

Taxonomy of Mathematical Formulations of GNNs

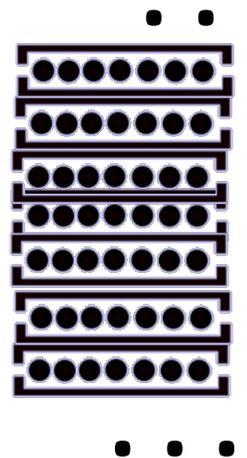
Local GNN formulations

Formulations based on functions operating on single vertices & edges

$$\mathbf{h}_{i-1}^{(l+1)} = \phi \left(\mathbf{h}_{i-1}^{(l+1)} \oplus \bigoplus_{j \in N(i-1)} \psi \left(\mathbf{h}_{i-1}^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

$$\mathbf{h}_{i+1}^{(l+1)} = \phi \left(\mathbf{h}_{i+1}^{(l+1)} \oplus \bigoplus_{j \in N(i+1)} \psi \left(\mathbf{h}_{i+1}^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$



Global GNN formulations

Formulations based on operations on matrices grouping all vertex & edge related vectors

All vertex feature vectors grouped together

$$\mathbf{H}^{l+1} = \sigma(\mathbf{Z}), \quad \mathbf{Z} = \Psi \mathbf{H} \mathbf{W}$$

Model parameters

Model details (a transformation of, among others, the adjacency matrix)

Taxonomy of Mathematical Formulations of GNNs

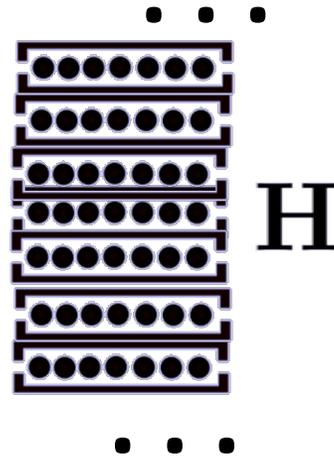
Local GNN formulations

Formulations based on functions operating on single vertices & edges

“Per-vertex” formulations can’t expose data reuse!

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

$$\mathbf{h}_{i+1}^{(l+1)} = \phi \left(\mathbf{h}_{i+1}^{(l)}, \bigoplus_{j \in N(i+1)} \psi \left(\mathbf{h}_{i+1}^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$



Global GNN formulations

Formulations based on operations on matrices grouping all vertex & edge related vectors

Global formulations can utilize optimal linear algebra algorithms!

- Communication-avoiding 2.5D MMM
- Tiling
- Kernel fusion
- ...

Taxonomy of Mathematical Formulations of GNNs

Local GNN

Problem:

Finding global formulations may be challenging

Formulations based on
single vertices & edges

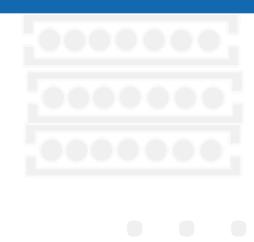
Grouping all vertex & edge related vectors
matrices

“Per-vertex” formulations can’t

All Global formulations can utilize
linear algebra algorithms!

Global formulations are known for simple
models such as Convolutional GNNs

$\mathbf{h}_i^{(l+1)}$

$$\mathbf{h}_{i+1}^{(l+1)} = \phi \left(\mathbf{h}_{i+1}^{(l+1)} \oplus_{j \in N(i+1)} \psi \left(\mathbf{h}_{i+1}^{(l)} \mathbf{h}_j^{(l)} \right) \right)$$


- Tiling
- Kernel fusion
- ...

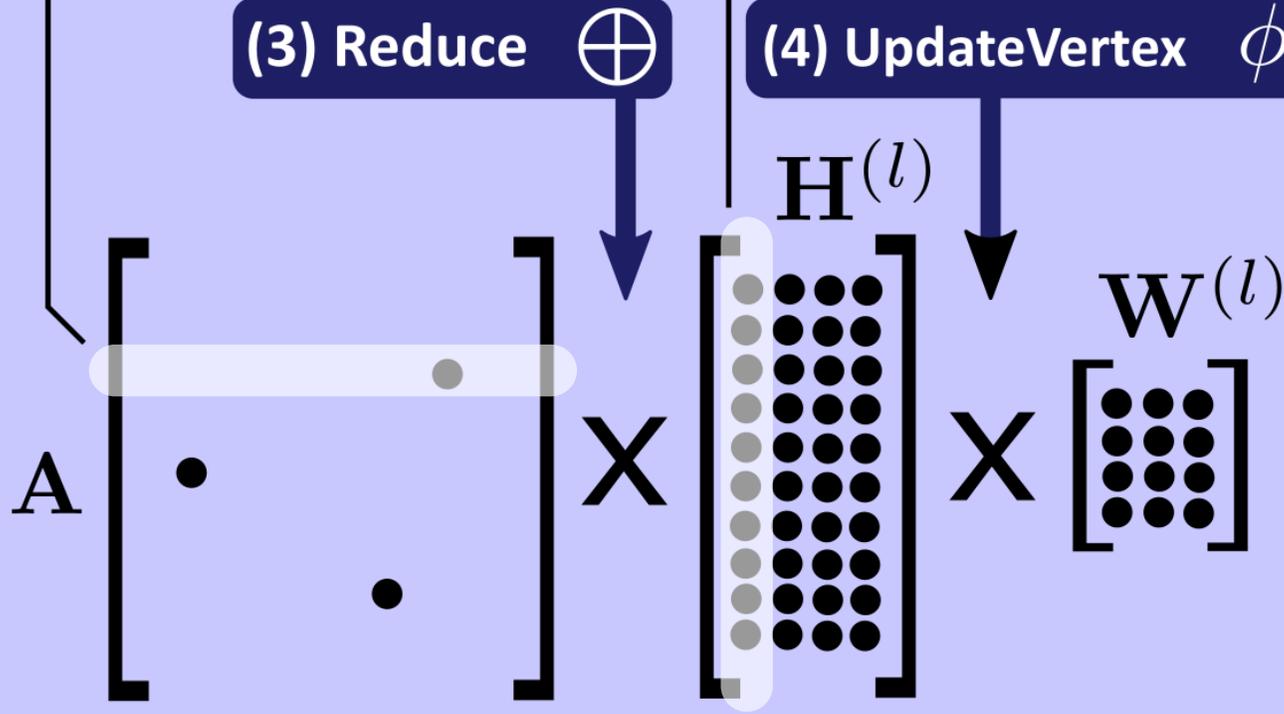
Communication-avoiding 2.5D

Global Formulations of GNN Models

The simplest model: Graph Convolution Network [1]

Highlighted row corresponds to the neighbors of a specific vertex v , whose feature vector is being computed

Highlighted column corresponds to the specific feature f that is being computed for vertex v



$$H^{(l+1)} = A \times H^{(l)} \times W^{(l)}$$

What are the global formulations of more complex models, such as attentional GNNs?

Also, why do we care?

Article

Highly accurate protein structure prediction with AlphaFold

Key technique?

Graph Attention Networks

<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

Check for updates

John Jumper^{1,4,5}, Richard Evans^{1,4}, Alexander Pritzel^{1,4}, Tim Green^{1,4}, Michael Figurnov^{1,4}, Anna Potapenko^{1,4}, Andrew J. Ballard^{1,4}, Andrew Senior^{1,4}, Rishub Jain^{1,4}, Jonas Adler^{1,4}, Michal Zielinski¹, Martin Tobias¹, Sebastian Bodenstein¹, David Pfaff¹, Pushmeet Kohli¹ & Demis

Proteins are essential to life, and understanding their structure can facilitate a

[1] T. Kipf et al. Semi-Supervised Classification with Graph Convolutional Networks. ICLR. 2017.

Attention in GNN Models

Convolutional GNN

Attentional GNN

The contribution of neighbors is **fixed**
e.g., sum

The contribution of neighbors is **learnable**

$$\mathbf{h}_i^{(l+1)} = \phi \left(\mathbf{h}_i^{(l)}, \bigoplus_{j \in N(i)} \psi \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} \right) \right)$$

Static, binary matrix
adjacency matrix of the graph

$$\mathbf{H}^{l+1} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$$

Matrix Ψ with dynamic attention scores

Attention in GNN Models

Convolutional GNN

Attentional GNN

We provide *generic* global formulations for any attentional GNNs, both for the **forward** and the **backward** propagation pass

Static, binary matrix adjacency matrix of the graph

$$\mathbf{H}^{l+1} = \left[\begin{array}{ccc} \bullet & & \\ \bullet & \bullet & \\ \bullet & & \bullet \end{array} \right] \left[\begin{array}{c} \bullet \\ \bullet \end{array} \right] \left[\begin{array}{cc} \bullet & \bullet \\ \bullet & \bullet \end{array} \right]$$

Matrix Ψ with dynamic attention scores

Attention in GNN Models – Forward Pass

$$\mathbf{H}^{l+1} = \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right]$$

$$\mathbf{H}^{l+1} = \sigma(\mathbf{Z}), \quad \mathbf{Z} = \Psi \mathbf{H} \mathbf{W}$$

Non-linearity \uparrow σ
 A sparse $n \times n$ tensor with attention scores, model specific \uparrow Ψ
 Features from previous layer \leftarrow \mathbf{H}
 weights \leftarrow \mathbf{W}

Formulating ψ is the „crux“ of devising a concrete formulation for a specific model

Vanilla Attention

Ψ $\Psi = \mathcal{A} \odot \mathbf{H}_x$

$$\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \odot \left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \right)$$

SDDMM

Graph Attention Network (GAT)

Ψ \star $\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \odot \left(\underbrace{\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right]}_{\text{rep}} \odot \underbrace{\left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \right)}_{\text{SDDMM}} + \underbrace{\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right]}_{\text{MM}} \odot \underbrace{\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right]}_{\text{rep}} \right)$

$$\Psi = \text{sm}(\mathcal{T}) \quad \mathcal{T} = \mathcal{A} \odot \exp(\sigma(\mathbf{C}))$$

$$\mathbf{C} = \text{rep}_n^T((\mathbf{H}'\underline{\mathbf{a}})^T) + \text{rep}_n(\mathbf{H}'\bar{\mathbf{a}})$$

Attention-based GNN (AGNN)

Ψ \star $\Psi = \mathcal{A} \odot \mathbf{H}_x \oslash \mathbf{n}_x$

$$\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \odot \left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \right) \oslash \left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \right)$$

SDDMM

Attention in GNN Models – Forward Pass

$$\mathbf{H}^{l+1} = \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right]$$

$$\mathbf{H}^{l+1} = \sigma(\mathbf{Z}), \quad \mathbf{Z} = \Psi \mathbf{H} \mathbf{W}$$

Non-linearity \uparrow σ
 A sparse $n \times n$ tensor with attention scores, **model specific** \uparrow Ψ
 Features from previous layer \leftarrow \mathbf{H}

Formulating ψ is the „crux“ of devising a concrete formulation for a specific model

Vanilla Attention

Ψ $\Psi = \mathcal{A} \odot \mathbf{H}_x$

$$\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \odot \left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \right)$$

SDDMM

Graph Attention Network (GAT)

Ψ \star

$$\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \odot \left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \odot \left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] + \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \right) \right)$$

$\Psi = \text{sm}(\mathcal{T}) \quad \mathcal{T} = \mathcal{A} \odot \exp(\sigma(\mathbf{C}))$
 $\mathbf{C} = \text{rep}_n^T((\mathbf{H}'\underline{\mathbf{a}})^T) + \text{rep}_n(\mathbf{H}'\bar{\mathbf{a}})$

Attention-based GNN (AGNN)

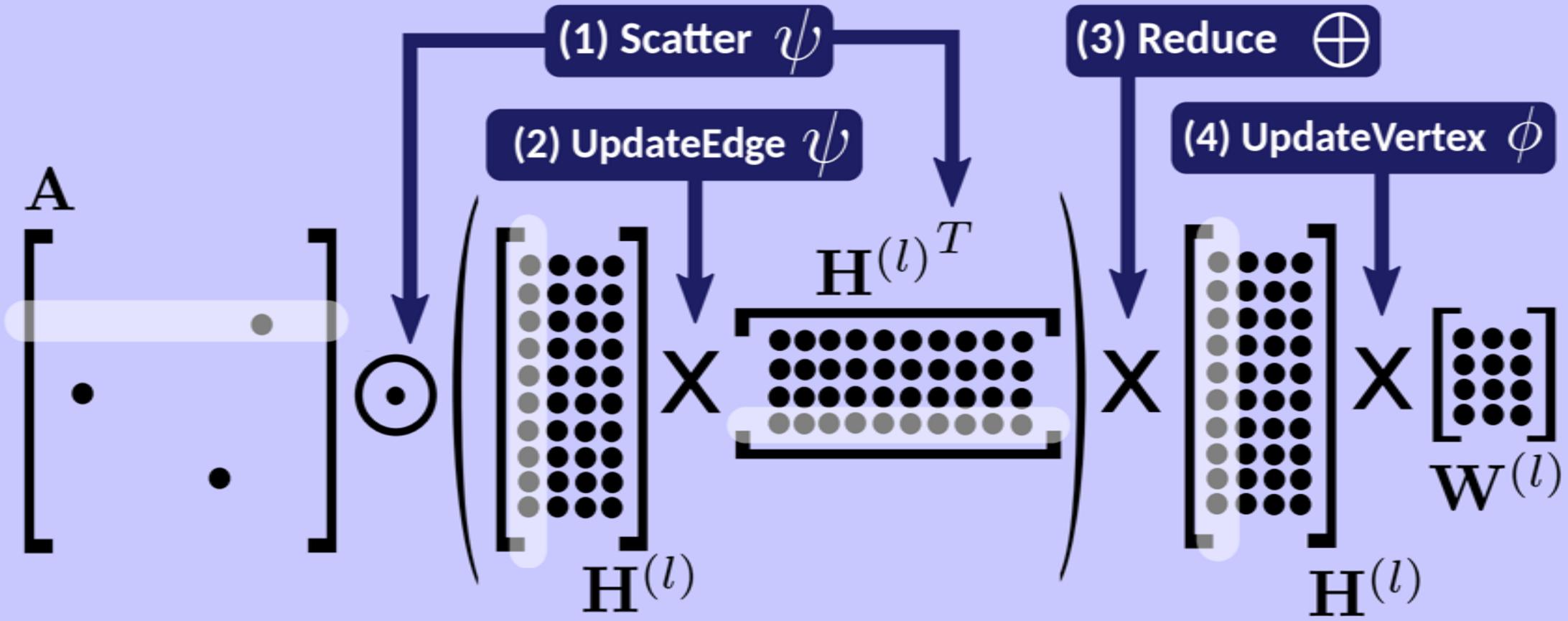
Ψ \star $\Psi = \mathcal{A} \odot \mathbf{H}_x \odot \mathbf{n}_x$

$$\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \odot \left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \odot \left(\left[\begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right] \right) \right)$$

SDDMM

Global Formulations of GNN Models

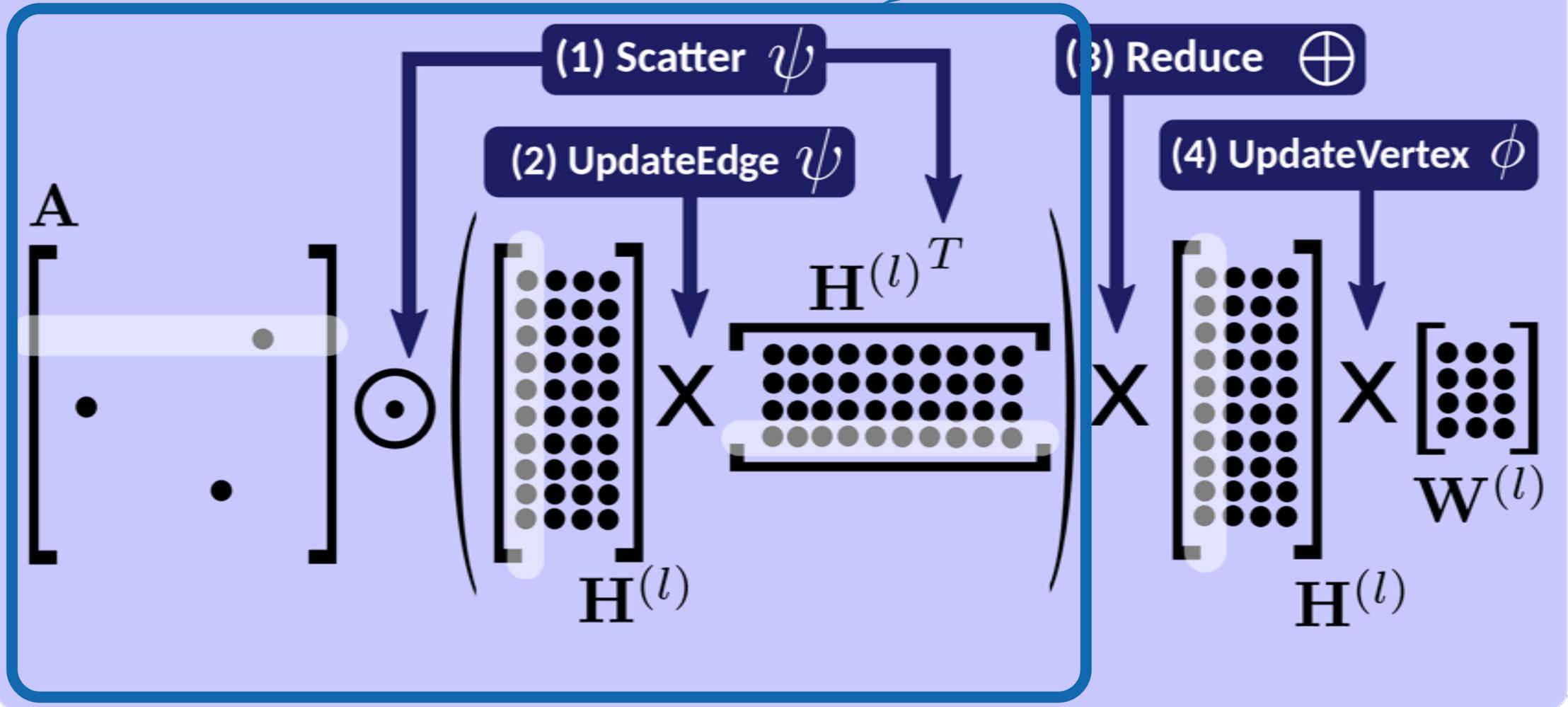
Example model: Graph Attention Network based on Dot Product (Vanilla Attention)



Global Formulations of GNN Models

This is our ψ for the Vanilla Attention Model

Example model: Graph Attention Network based on Dot Product (Vanilla Attention)



Attention in GNN Models – Forward Pass

Formulating ψ is the „crux“ of devising a concrete formulation for a specific model

$$\mathbf{H}^{l+1} = \sigma(\mathbf{Z}), \quad \mathbf{Z} = \Psi \mathbf{H} \mathbf{W}$$

σ : Non-linearity
 Ψ : A sparse $n \times n$ tensor with attention scores, model specific
 \mathbf{H} : Features from previous layer
 \mathbf{W} : weights

Vanilla Attention

$$\Psi = \mathcal{A} \odot \mathbf{H}_x$$

SDDMM

Graph Attention Network (GAT)

$$\Psi = \text{sm}(\mathcal{T}) \quad \mathcal{T} = \mathcal{A} \odot \exp(\sigma(\mathbf{C}))$$

$$\mathbf{C} = \text{rep}_n^T((\mathbf{H}'\underline{\mathbf{a}})^T) + \text{rep}_n(\mathbf{H}'\bar{\mathbf{a}})$$

Diagrams showing matrix operations: \odot (SDDMM), $+$ (MM), and rep (repetition).

Attention-based GNN (AGNN)

$$\Psi = \mathcal{A} \odot \mathbf{H}_x \odot \mathbf{n}_x$$

SDDMM

Graph Attention Network (GAT)

Local ψ formulation is very involving – how to obtain the global formulation?

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$


 Vector concatenation

Graph Attention Network (GAT)

Local ψ formulation is very involving – how to obtain the global formulation?

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{\mathcal{N}}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$

This is the softmax normalization, we'll get to it later 😊

Graph Attention Network (GAT)

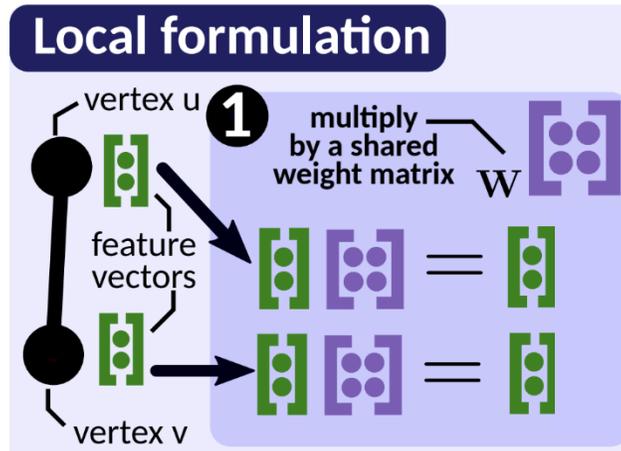
Local ψ formulation is very involving – how to obtain the global formulation?

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u \right] \right)\right)}{\sum_{y \in \hat{\mathcal{N}}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y \right] \right)\right)} \mathbf{h}_u$$

Graph Attention Network (GAT)

Local ψ formulation is very involving – how to obtain the global formulation?

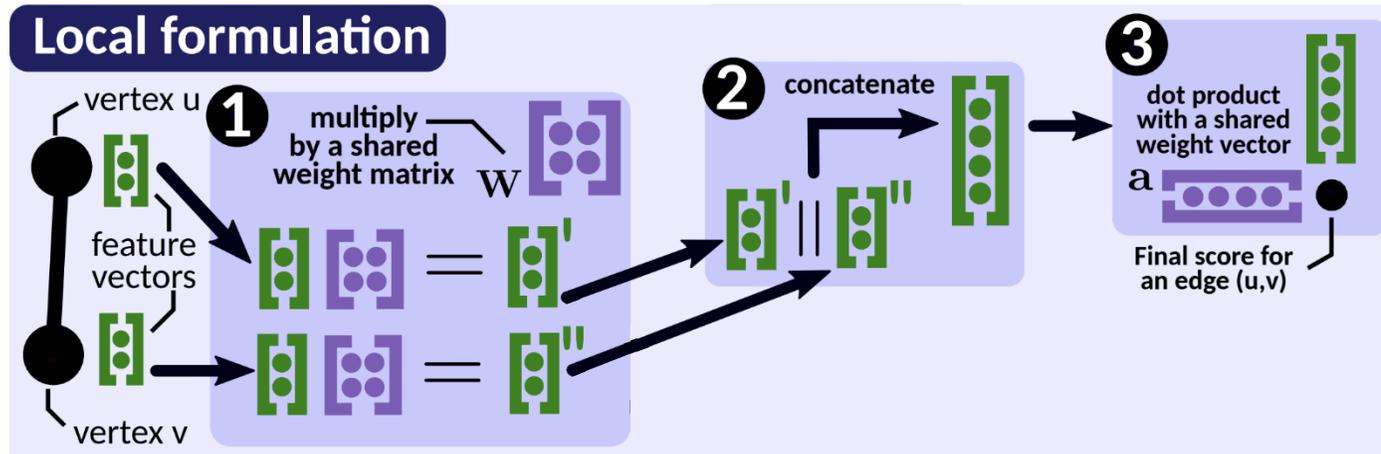
$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u \right] \right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y \right] \right)\right)} \mathbf{h}_u$$



Graph Attention Network (GAT)

Local ψ formulation is very involving – how to obtain the global formulation?

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$



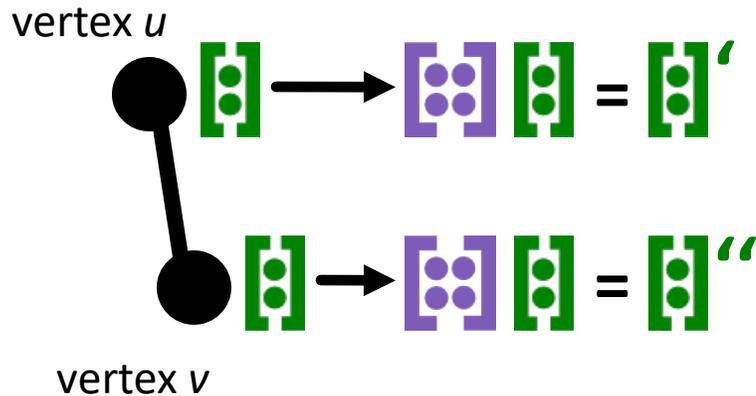
Graph Attention Network (GAT)

Local ψ formulation is very involving – how to obtain the global formulation?



$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$

Local formulation



multiply by shared weight matrix

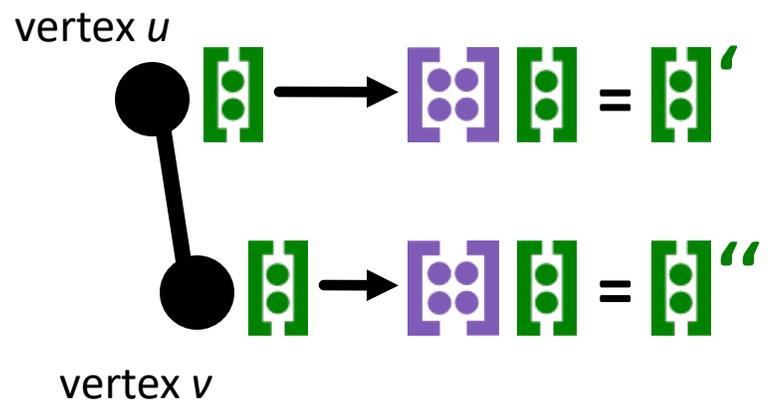
Graph Attention Network (GAT)

Local ψ formulation is very involving – how to obtain the global formulation?

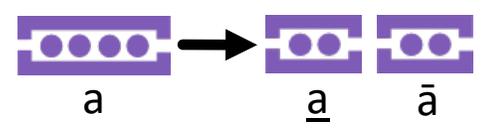


$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$

Local formulation



multiply by shared weight matrix



vector split

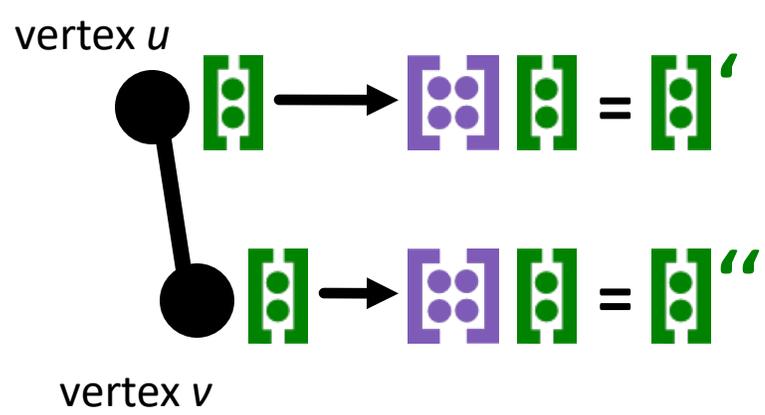
Graph Attention Network (GAT)

Local ψ formulation is very involving – how to obtain the global formulation?

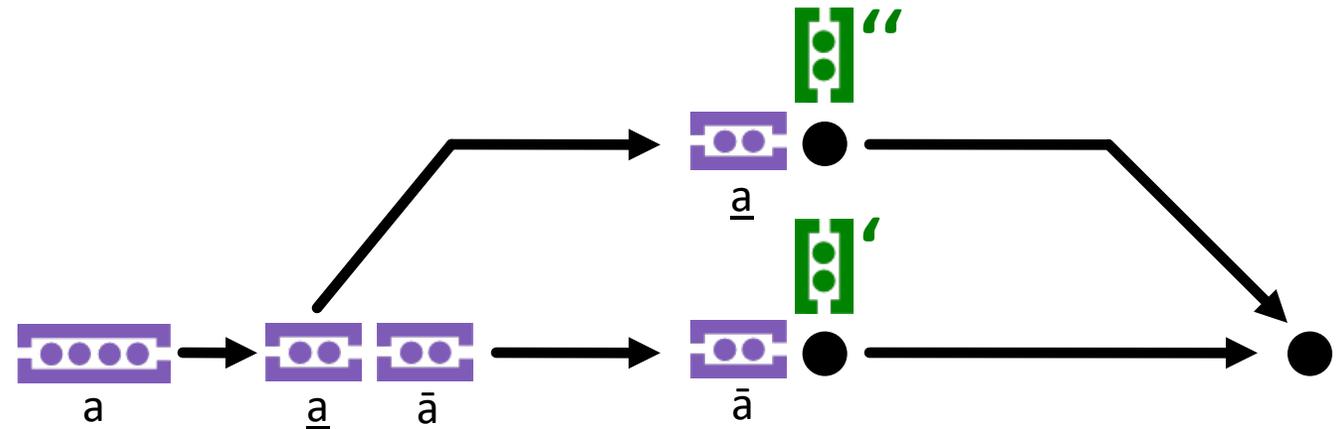


$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{\mathcal{N}}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$

Local formulation



multiply by shared weight matrix



vector split

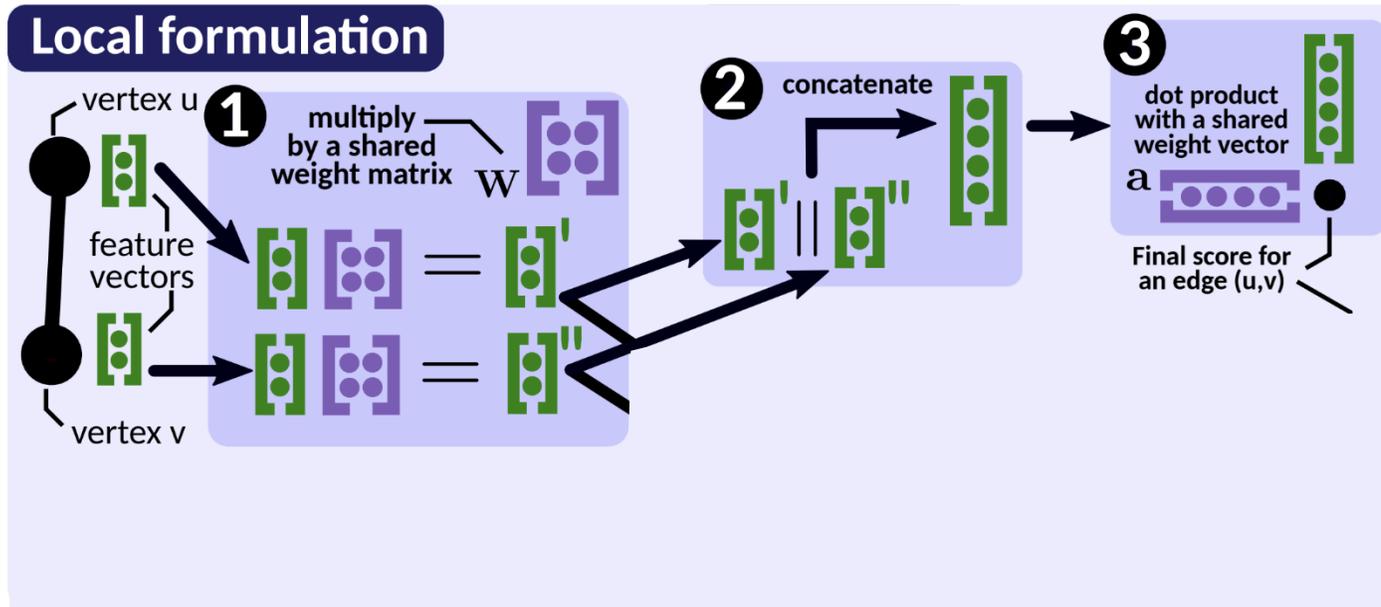
two dot products

sum partial sums

Graph Attention Network (GAT)

Local ψ formulation is very involving – how to obtain the global formulation?

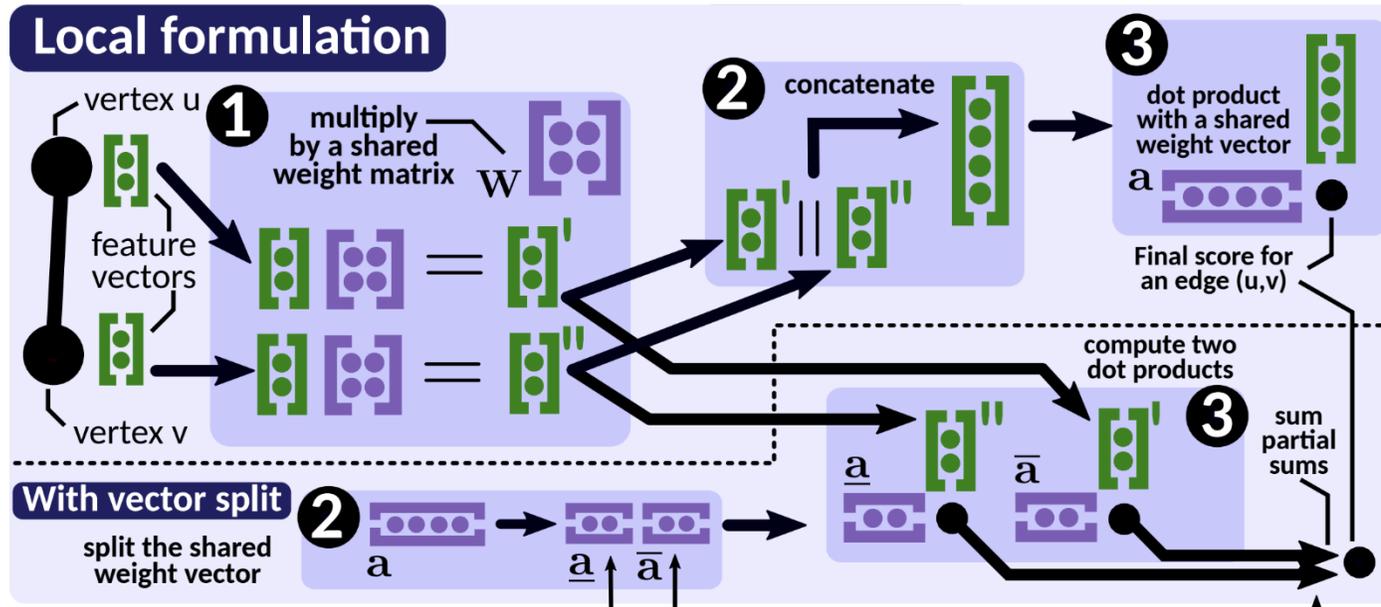
$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$



Graph Attention Network (GAT)

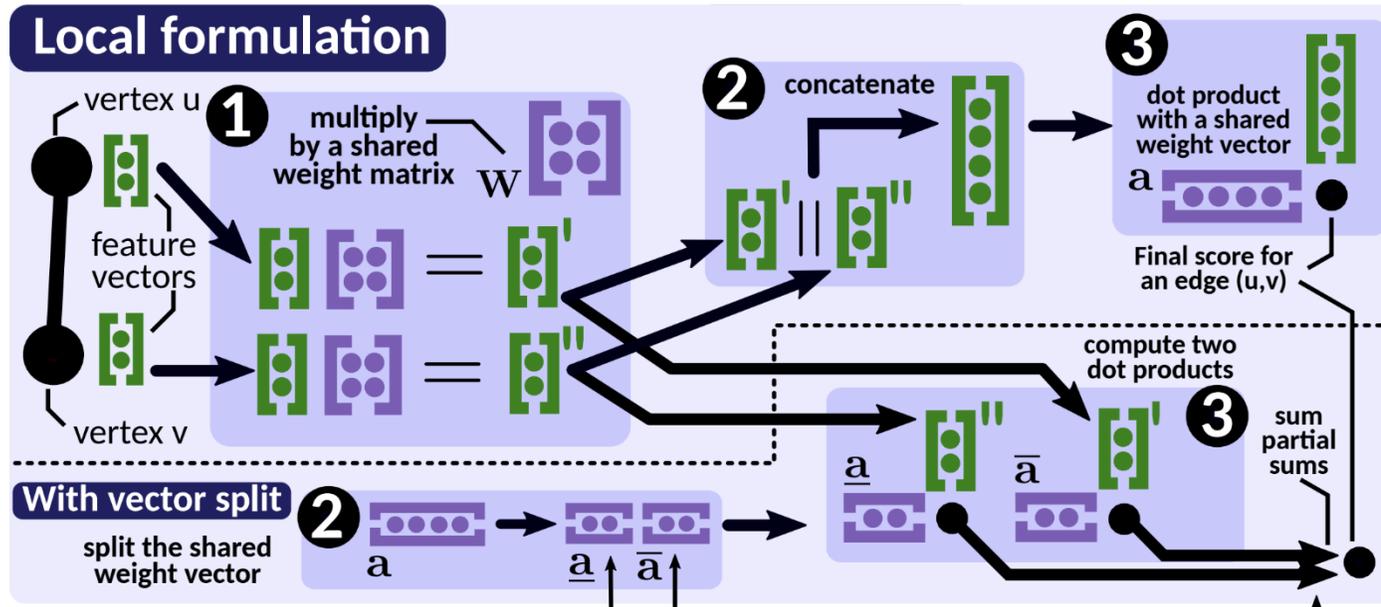
Local ψ formulation is very involving – how to obtain the global formulation?

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$



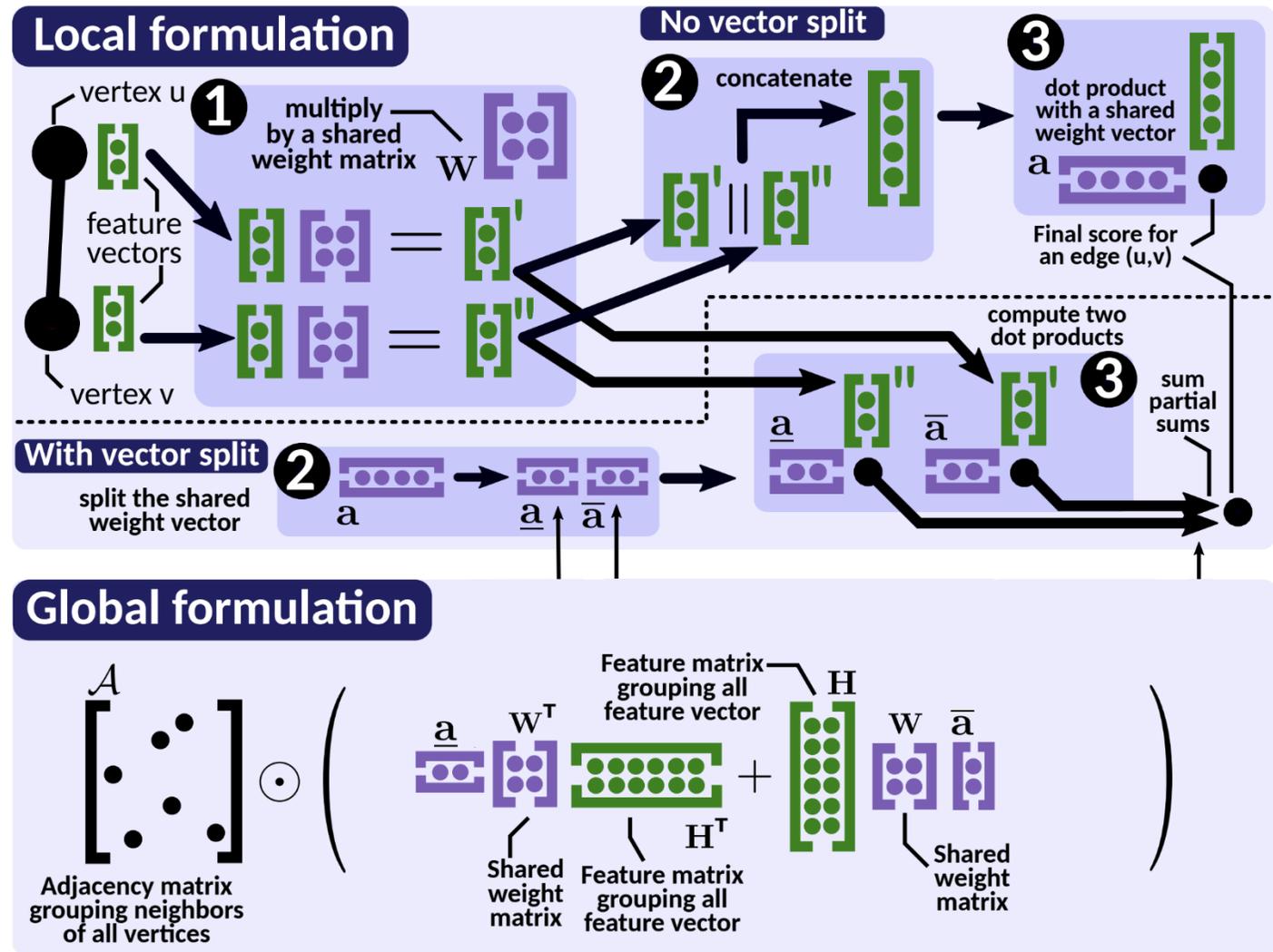
Graph Attention Network (GAT)

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$



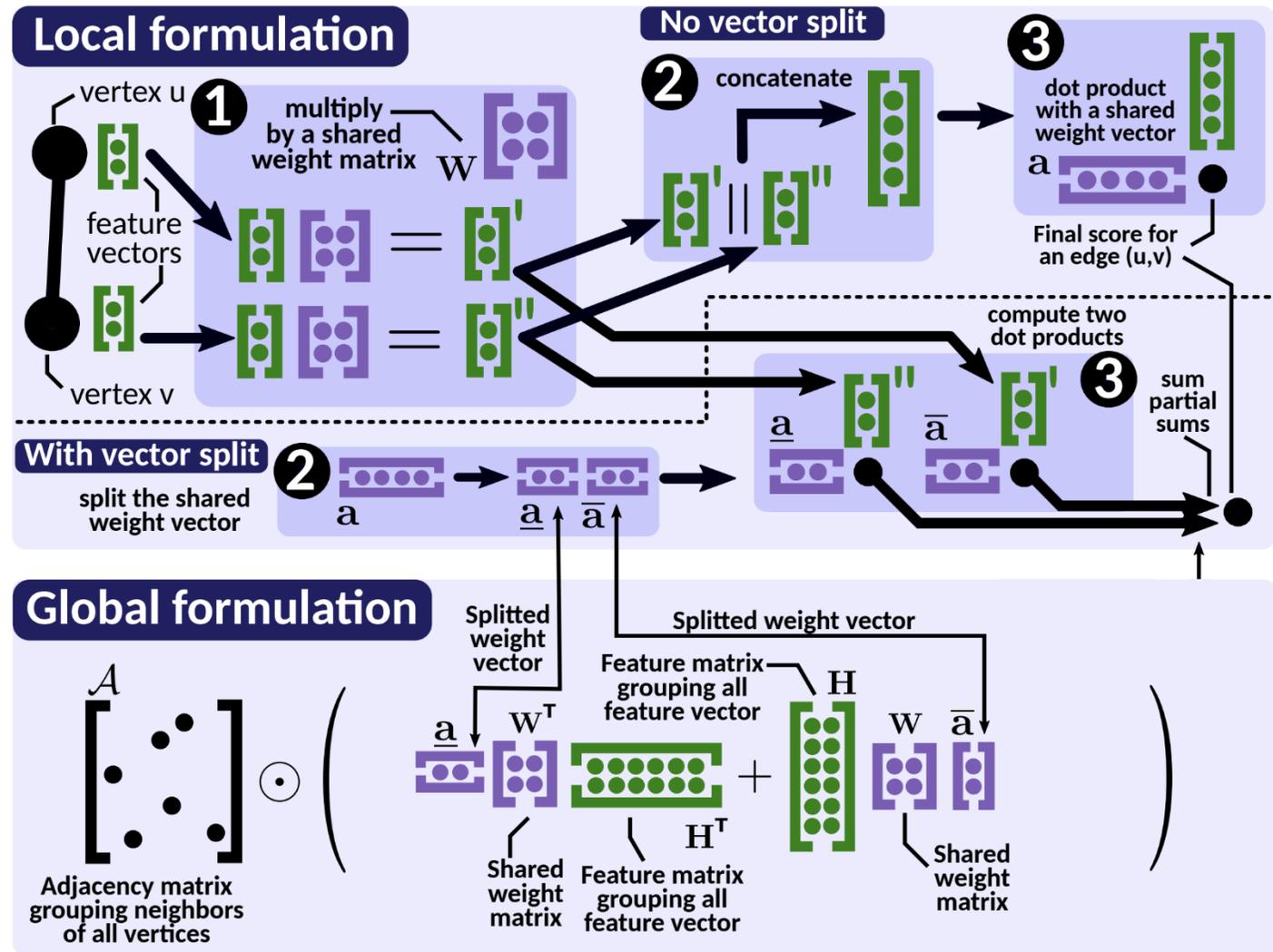
Graph Attention Network (GAT)

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(a^T \cdot \left[Wh_v || Wh_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(a^T \cdot \left[Wh_v || Wh_y\right]\right)\right)} h_u$$



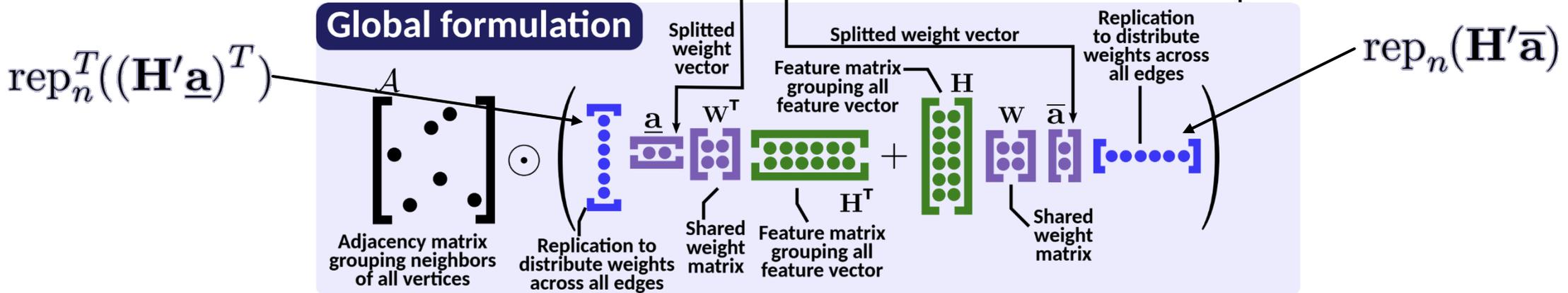
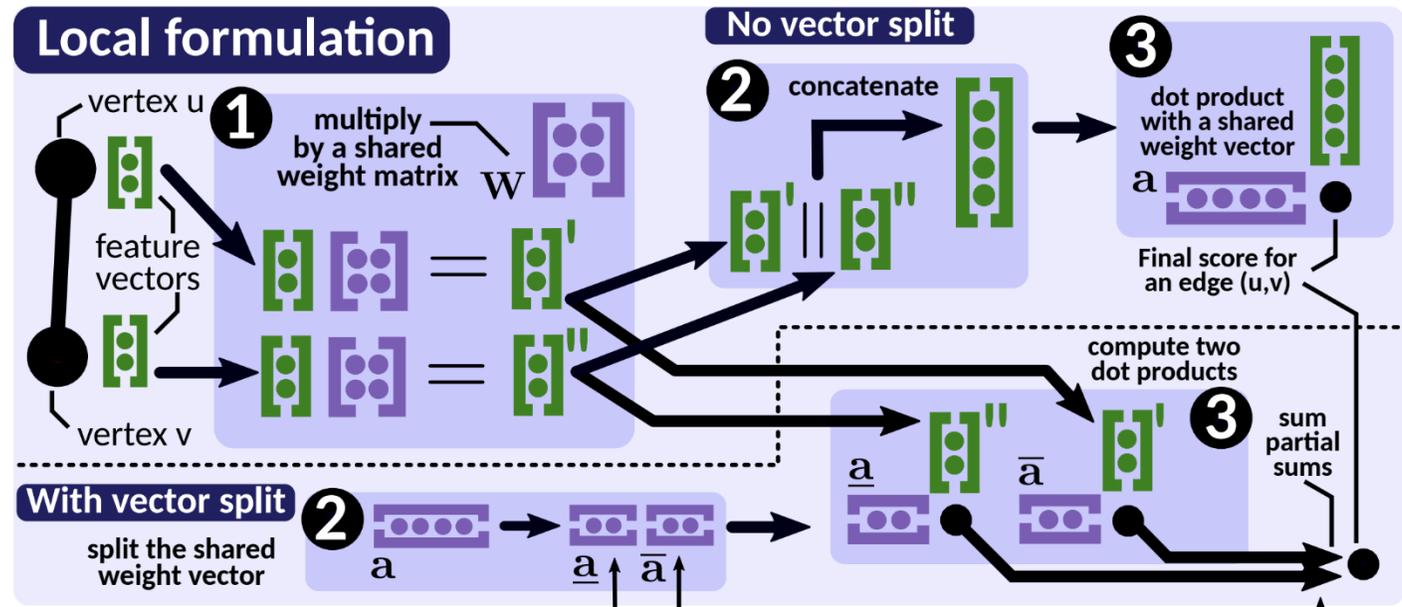
Graph Attention Network (GAT)

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(a^T \cdot \left[Wh_v || Wh_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(a^T \cdot \left[Wh_v || Wh_y\right]\right)\right)} h_u$$



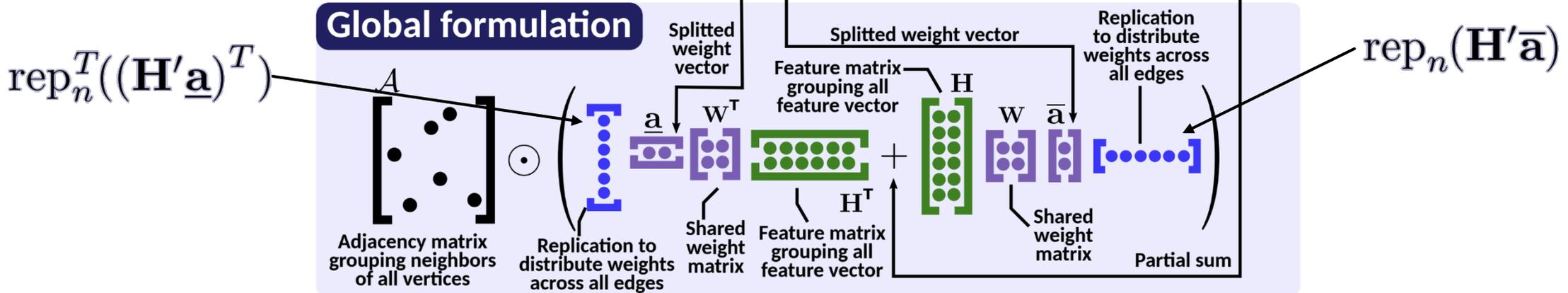
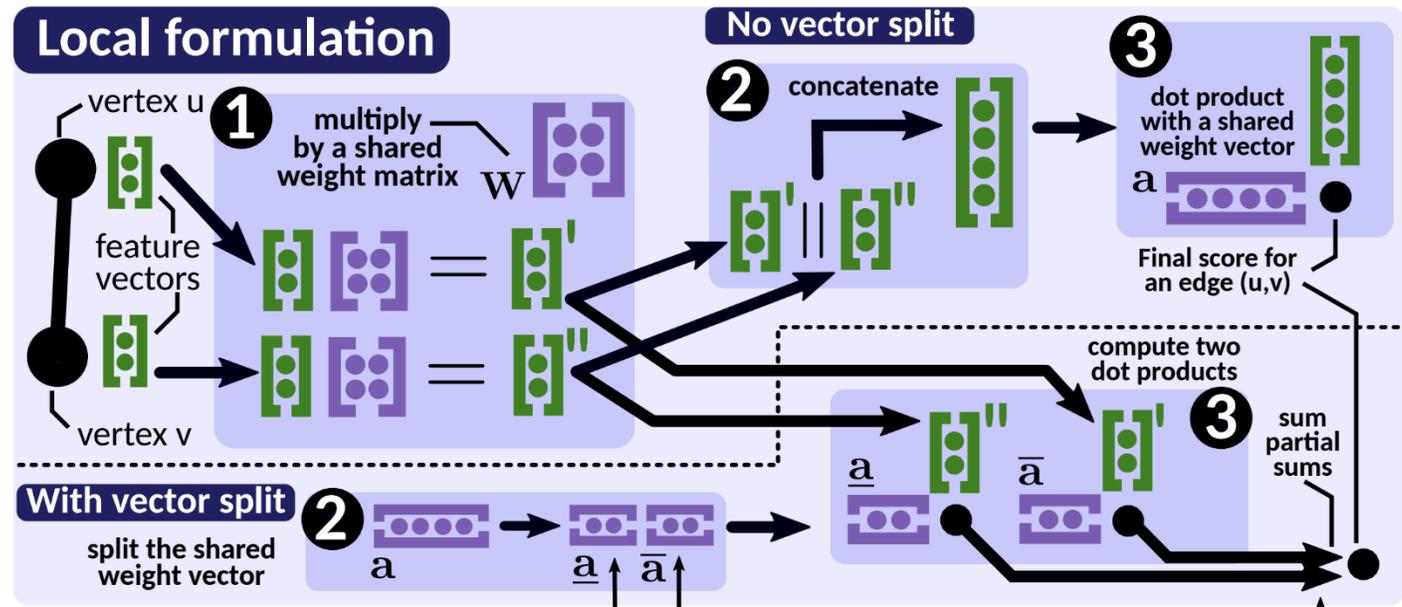
Graph Attention Network (GAT)

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(a^T \cdot \left[Wh_v || Wh_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(a^T \cdot \left[Wh_v || Wh_y\right]\right)\right)} h_u$$



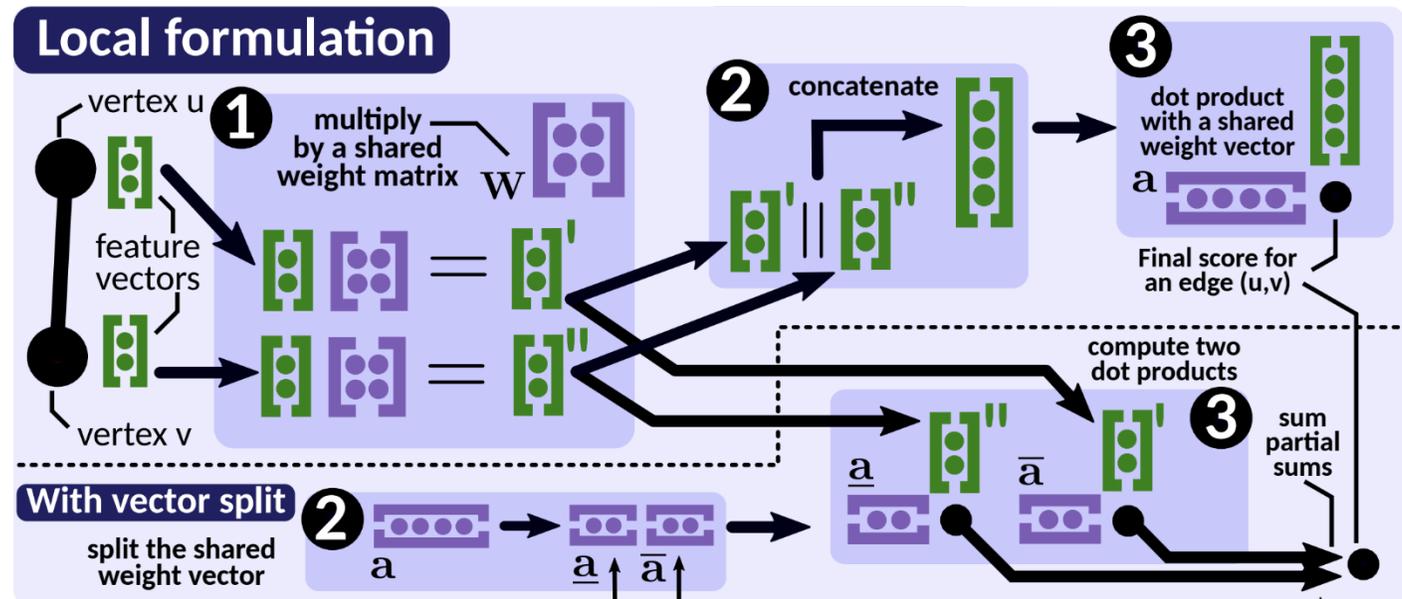
Graph Attention Network (GAT)

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(a^T \cdot \left[Wh_v || Wh_u\right]\right)\right)}{\sum_{y \in \hat{N}(v)} \exp\left(\sigma\left(a^T \cdot \left[Wh_v || Wh_y\right]\right)\right)} h_u$$

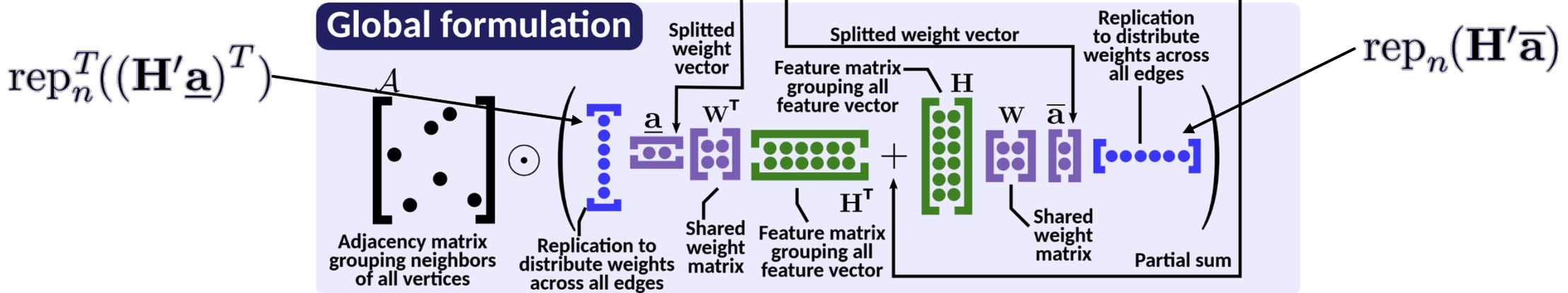


Graph Attention Network (GAT)

$$\psi_{v,u} = \frac{\exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_u\right]\right)\right)}{\sum_{y \in \hat{\mathcal{N}}(v)} \exp\left(\sigma\left(\mathbf{a}^T \cdot \left[\mathbf{W}\mathbf{h}_v \parallel \mathbf{W}\mathbf{h}_y\right]\right)\right)} \mathbf{h}_u$$



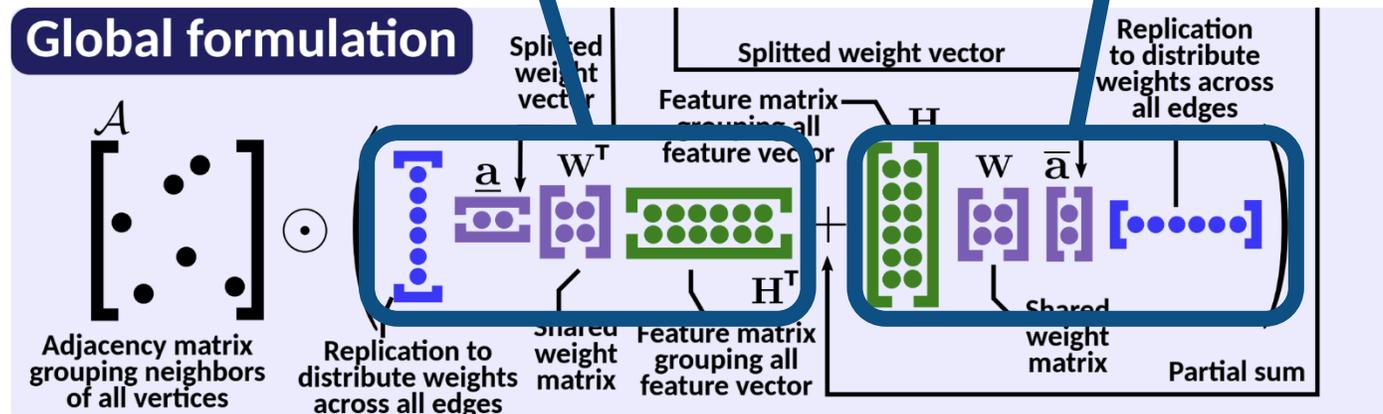
Element-wise operations are not shown (negligible from performance perspective)



Graph Attention Network (GAT)

$$\Psi = \text{sm}(\mathcal{X}) \quad \mathcal{X} = \mathcal{A} \odot \exp(\sigma(\mathbf{C}))$$

$$\mathbf{C} = \text{rep}_n^T((\mathbf{H}'\underline{\mathbf{a}})^T) + \text{rep}_n(\mathbf{H}'\bar{\mathbf{a}})$$

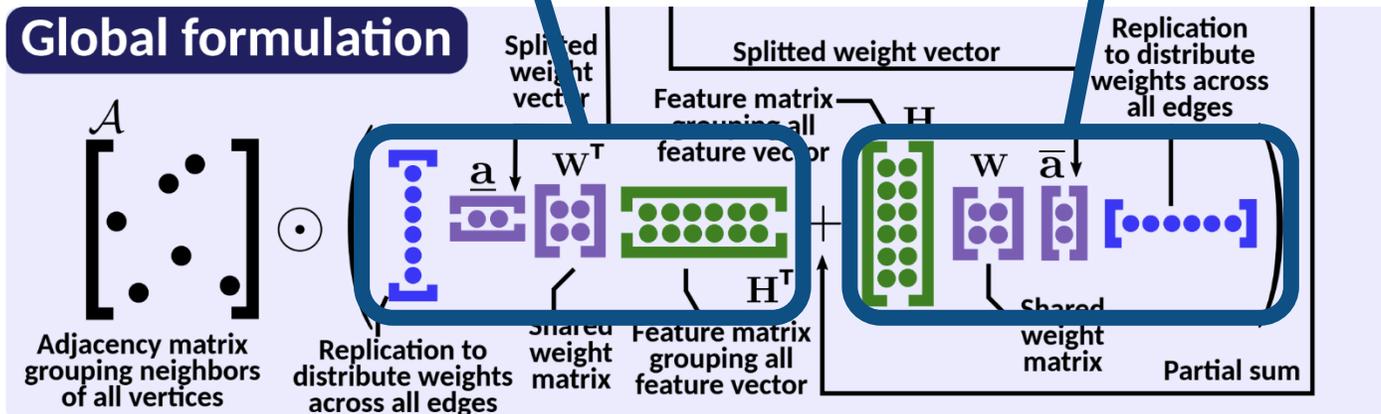


Graph Attention Network (GAT)

Let's see how softmax (sm) looks like in the global formulation

$$\Psi = \text{sm}(\mathcal{X}) \quad \mathcal{X} = \mathcal{A} \odot \exp(\sigma(\mathbf{C}))$$

$$\mathbf{C} = \text{rep}_n^T((\mathbf{H}'\underline{\mathbf{a}})^T) + \text{rep}_n(\mathbf{H}'\bar{\mathbf{a}})$$



Global Formulations of GNN Kernels – Softmax

Tensor algebra expression

$$\begin{aligned} \text{sm}(\mathcal{X}) &= \exp(\mathcal{X}) \otimes \text{rs}_n(\exp(\mathcal{X})) \\ &= \exp(\mathcal{X}) \otimes (\exp(\mathcal{X}) \mathbf{1}\mathbf{1}^T) \end{aligned}$$

Example

$$\mathcal{X} = \begin{bmatrix} 0 & 0 & 0 & 2 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 4 \\ 2 & 0 & 0 & 0 & 3 \\ 1 & 0 & 4 & 3 & 0 \end{bmatrix} \xrightarrow[\text{(1) element-wise exponentiation of non-zero elements}]{\mathcal{X}_1 = \exp(\mathcal{X})} \begin{bmatrix} 0 & 0 & 0 & e^2 & e \\ 0 & 0 & e & 0 & 0 \\ 0 & e & 0 & 0 & e^4 \\ e^2 & 0 & 0 & 0 & e^3 \\ e & 0 & e^4 & e^3 & 0 \end{bmatrix} \xrightarrow[\text{(2) multiplication by a column vector of ones}]{\mathcal{X}_2 = \text{sum}(\mathcal{X}_1)} \begin{bmatrix} 0 & 0 & 0 & e^2 & e \\ 0 & 0 & e & 0 & 0 \\ 0 & e & 0 & 0 & e^4 \\ e^2 & 0 & 0 & 0 & e^3 \\ e & 0 & e^4 & e^3 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} e^2 + e \\ e \\ e + e^4 \\ e^2 + e^3 \\ e + e^4 + e^3 \end{bmatrix}$$

$$\mathcal{X}_3 = \text{rep}_n(\mathcal{X}_2) \xrightarrow[\text{(3) multiplication by a row vector of ones}]{\text{(3) multiplication by a row vector of ones}} \begin{bmatrix} e^2 + e \\ e \\ e + e^4 \\ e^2 + e^3 \\ e + e^4 + e^3 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} e^2 + e \\ e \\ e + e^4 \\ e^2 + e^3 \\ e + e^4 + e^3 \end{bmatrix} \dots$$

part 1 of rs()

part 2 of rs()

Global Formulations of GNN Kernels – Softmax

Tensor algebra expression

$$\begin{aligned} \text{sm}(\mathcal{X}) &= \exp(\mathcal{X}) \oslash \text{rs}_n(\exp(\mathcal{X})) \\ &= \exp(\mathcal{X}) \oslash (\exp(\mathcal{X}) \mathbf{1}\mathbf{1}^T) \end{aligned}$$

Example

$$\mathcal{X} = \begin{bmatrix} 0 & 0 & 0 & 2 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 4 \\ 2 & 0 & 0 & 0 & 3 \\ 1 & 0 & 4 & 3 & 0 \end{bmatrix} \xrightarrow[\substack{\text{(1) element-wise} \\ \text{exponentiation of} \\ \text{non-zero elements}}]{\mathcal{X}_1 = \exp(\mathcal{X})} \begin{bmatrix} 0 & 0 & 0 & e^2 & e \\ 0 & 0 & e & 0 & 0 \\ 0 & e & 0 & 0 & e^4 \\ e^2 & 0 & 0 & 0 & e^3 \\ e & 0 & e^4 & e^3 & 0 \end{bmatrix} \xrightarrow[\substack{\text{(2) multiplication} \\ \text{by a column} \\ \text{vector of ones}}]{\mathcal{X}_2 = \text{sum}(\mathcal{X}_1)} \begin{bmatrix} 0 & 0 & 0 & e^2 & e \\ 0 & 0 & e & 0 & 0 \\ 0 & e & 0 & 0 & e^4 \\ e^2 & 0 & 0 & 0 & e^3 \\ e & 0 & e^4 & e^3 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} e^2 + e \\ e \\ e + e^4 \\ e^2 + e^3 \\ e + e^4 + e^3 \end{bmatrix}$$

$$\xrightarrow[\substack{\text{(3) multiplication} \\ \text{by a row vector} \\ \text{of ones}}]{\mathcal{X}_3 = \text{rep}_n(\mathcal{X}_2)} \begin{bmatrix} e^2 + e \\ e \\ e + e^4 \\ e^2 + e^3 \\ e + e^4 + e^3 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} e^2 + e \\ e \\ e + e^4 \\ e^2 + e^3 \\ e + e^4 + e^3 \end{bmatrix} \cdots \xrightarrow[\substack{\text{(4) element-} \\ \text{-wise division} \\ \text{of edge scores} \\ \text{(normalization)}}]{\mathcal{X}_1 \oslash \mathcal{X}_3} \begin{bmatrix} 0 & 0 & 0 & e^2 & e \\ 0 & 0 & e & 0 & 0 \\ 0 & e & 0 & 0 & e^4 \\ e^2 & 0 & 0 & 0 & e^3 \\ e & 0 & e^4 & e^3 & 0 \end{bmatrix} \oslash \begin{bmatrix} e^2 + e \\ e \\ e + e^4 \\ e^2 + e^3 \\ e + e^4 + e^3 \end{bmatrix} \cdots$$

Global Formulations of GNN Kernels – Backward Pass

Generic formulation

$$\mathbf{G}^{l-1} = \sigma' \left(\mathbf{Z}^{l-1} \right) \odot \Gamma^l$$

$$\mathbf{Y}^l = \mathbf{H}^{lT} \Psi \left(\mathcal{A}^T, \mathbf{H}^l \right) \mathbf{G}^l + \mathbf{G}^l \mathbf{W}^{lT} \mathbf{H}^{lT} \frac{\partial \Psi}{\partial \mathbf{W}^l}$$

Matrix view

$$\mathbf{G} = \begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix} \odot \begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix}$$

$$\mathbf{Y} = \begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix} \begin{bmatrix} \bullet & \bullet \\ \bullet & \bullet \\ \bullet & \bullet \\ \bullet & \bullet \end{bmatrix} \begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix}$$

MSpMM
 (or: MM + SpMM)

Global Formulations of GNN Kernels – Backward Pass

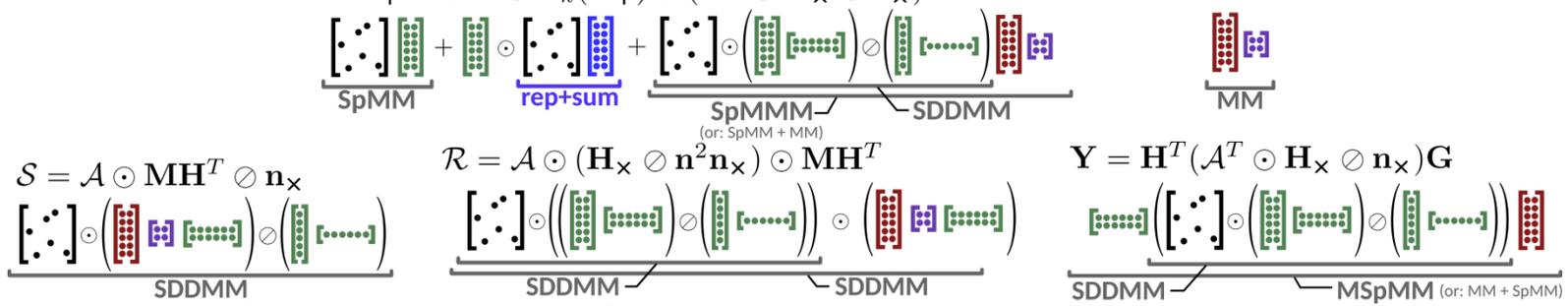
Vanilla Attention (VA)

$$\Gamma = \mathbf{N}_+ \mathbf{H} + (\mathcal{A}^T \odot \mathbf{H}_x) \mathbf{M} \quad \mathbf{N} = \mathcal{A} \odot (\mathbf{M} \mathbf{H}^T) \quad \mathbf{M} = \mathbf{G} \mathbf{W}^T \quad \mathbf{Y} = \mathbf{H}^T (\mathcal{A}^T \odot \mathbf{H}_x) \mathbf{G}$$


Graph Attention Network (GAT)

$$\Gamma = \mathbf{H} \odot (\text{rs}_k(\mathbf{F}) \odot \text{rep}_n^T(\mathbf{a}') + \text{rs}_k(\mathbf{F}^T) \odot \text{rep}_n^T(\mathbf{a}'')) \quad \mathbf{Y} = \mathbf{H}^T \text{sm}(\mathcal{T})^T \mathbf{G} \quad \mathbf{F} = \sigma'(\mathbf{C}) \odot (\text{rs}_n(\mathcal{T}) \odot (\mathbf{G} \mathbf{H}^T) + \text{rs}_n^2(\mathcal{T}) \odot \text{rs}_n(\mathcal{T} \mathbf{H}' \odot \mathbf{G}))$$

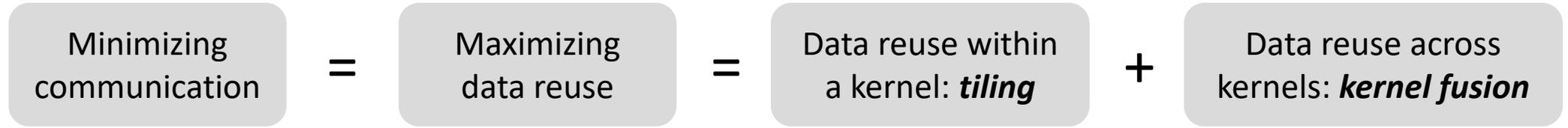
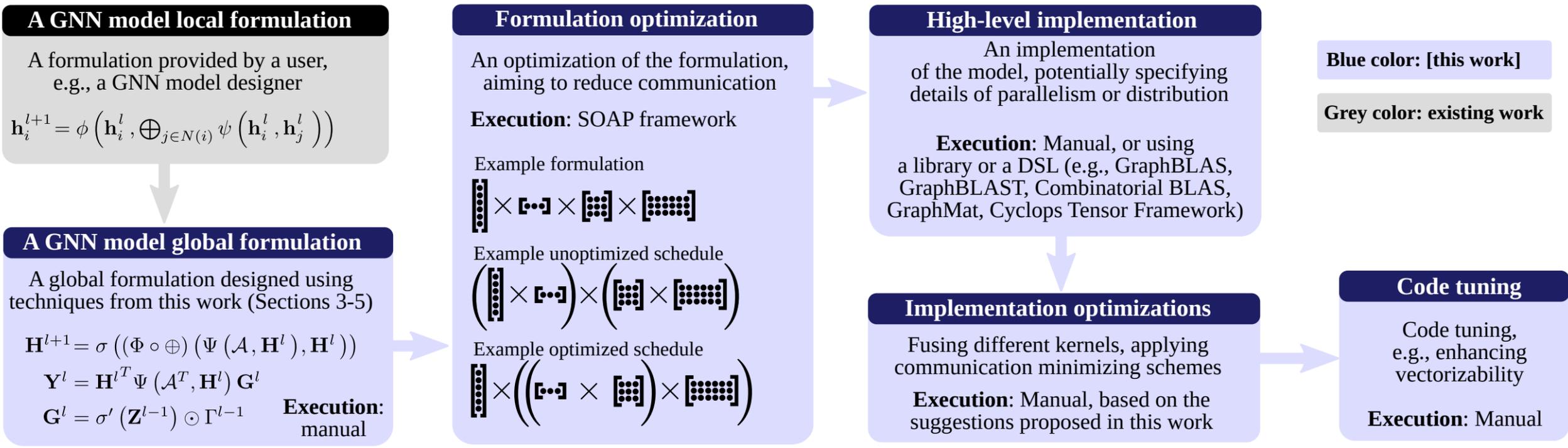

Attention-based GNN (AGNN) [1]

$$\Gamma = \mathbf{S}_+ \mathbf{H} + \mathbf{H} \odot \text{rs}_k(\mathcal{R}_+) + (\mathcal{A}^T \odot \mathbf{H}_x \odot \mathbf{n}_x) \mathbf{M} \quad \mathbf{M} = \mathbf{G} \mathbf{W}^T$$


$$\mathbf{S} = \mathcal{A} \odot \mathbf{M} \mathbf{H}^T \odot \mathbf{n}_x \quad \mathcal{R} = \mathcal{A} \odot (\mathbf{H}_x \odot \mathbf{n}^2 \mathbf{n}_x) \odot \mathbf{M} \mathbf{H}^T \quad \mathbf{Y} = \mathbf{H}^T (\mathcal{A}^T \odot \mathbf{H}_x \odot \mathbf{n}_x) \mathbf{G}$$

[1] K. Thekumparampil et al. Attention-based Graph Neural Network for Semi-supervised Learning. arXiv:2018.

The Entire Optimization Toolchain



Evaluation

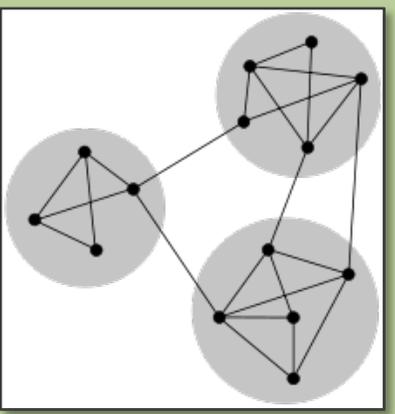


CSCS Cray Piz Daint supercomputer

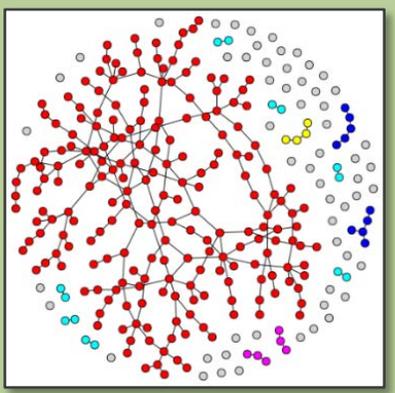
- Cray XC50 nodes
- Intel Xeon E5-2690 v3, 12 cores
- Single NVIDIA Tesla P100 per node
- 64 GB RAM per node

Considered Graph Datasets

Synthetic graphs



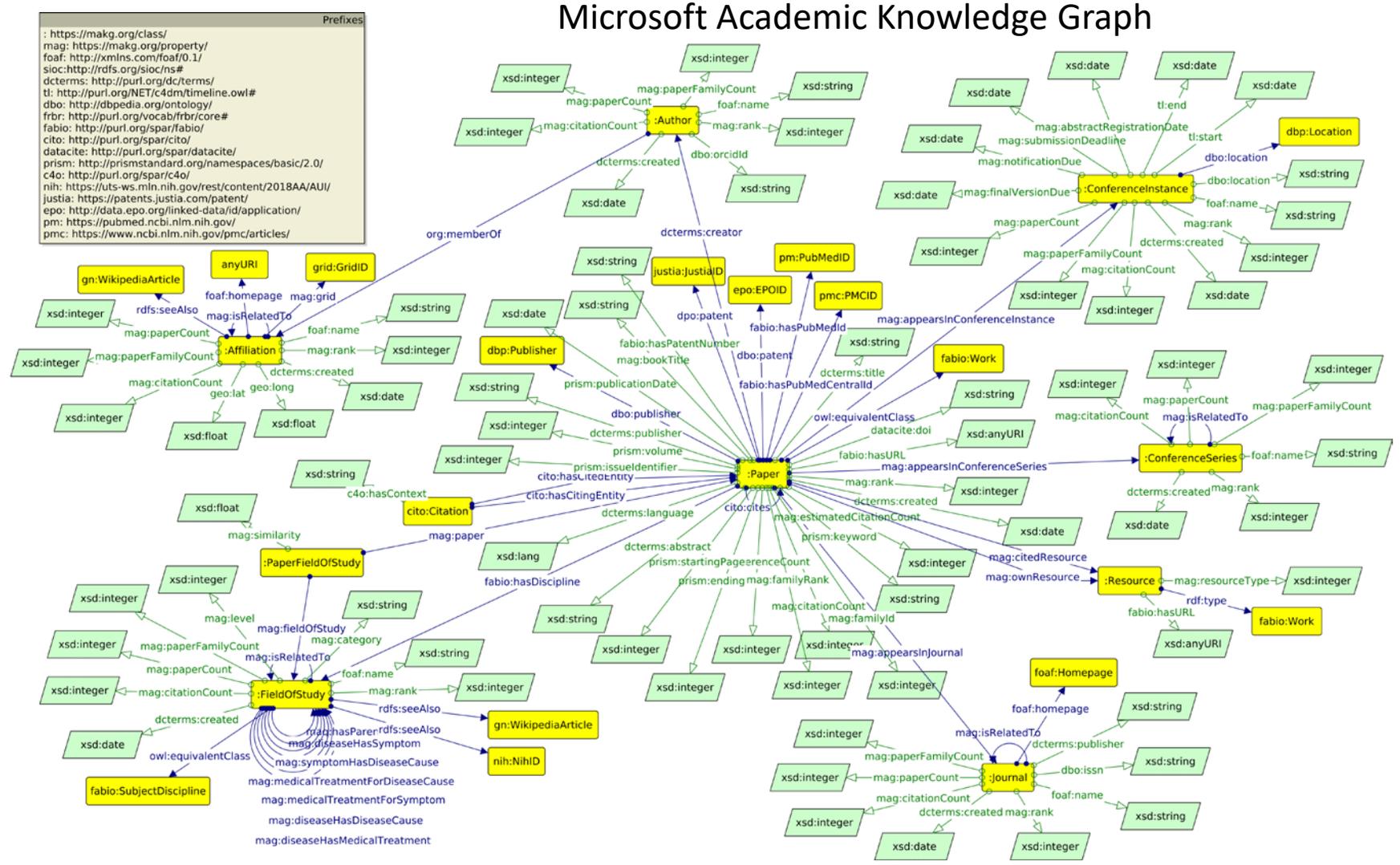
Kronecker [1]



Erdős-Rényi [2]

```

Prefixes
: https://makg.org/class/
mag: https://makg.org/property/
foaf: http://xmlns.com/foaf/0.1/
sioc: http://rdfs.org/sioc/ns#
dcterms: http://purl.org/dc/terms/
ti: http://purl.org/NET/c4dm/timeline.owl#
dbo: http://dbpedia.org/ontology/
frbr: http://purl.org/vocab/frbr/core#
fabio: http://purl.org/spar/fabio/
cito: http://purl.org/spar/cito/
datacite: http://purl.org/spar/datacite/
prism: http://prismstandard.org/namespaces/basic/2.0/
c4o: http://purl.org/spar/c4o/
nih: https://uts-ws.nlm.nih.gov/rest/content/2018AA/AU/
justia: https://patents.justia.com/patent/
epo: http://data.epo.org/linked-data/rd/application/
pm: https://pubmed.ncbi.nlm.nih.gov/
pmc: https://www.ncbi.nlm.nih.gov/pmc/articles/
    
```

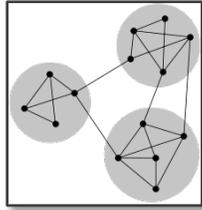


Microsoft Academic Knowledge Graph

[1] J. Leskovec et al. Kronecker Graphs: An Approach to Modeling Networks. J. Mach. Learn. Research. 2010.
 [2] P. Erdos and A. Renyi. On the evolution of random graphs. Pub. Math. Inst. Hun. A. Science. 1960.

Strong Scaling

Kronecker [1]



- | | | |
|------------------|----------------|------------|
| GAT - This Work | GAT - DistDGL | GAT - DGL |
| VA - This Work | VA - DistDGL | VA - DGL |
| AGNN - This Work | AGNN - DistDGL | AGNN - DGL |

$p = 1\%$

$p = 0.01\%$

Sparsity: $p = m / (n * n)$

k=16

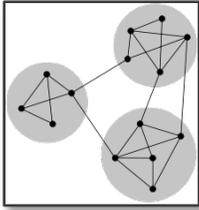
k=128

#features

[1] J. Leskovec et al. Kronecker Graphs: An Approach to Modeling Networks. J. Mach. Learn. Research. 2010.

Strong Scaling

Kronecker [1]



$\rho = 1\%$

$n = 131k, m = 172M$

$n = 262k, m = 687M$

#vertices

#edges

k=16

k=128

—○— GAT - This Work

—◇— VA - This Work

—□— AGNN - This Work

-○- GAT - DistDGL

-◇- VA - DistDGL

-□- AGNN - DistDGL

● GAT - DGL

◆ VA - DGL

■ AGNN - DGL

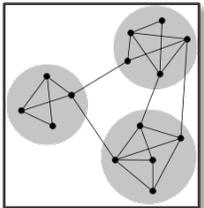
$\rho = 0.01\%$

$n = 1M, m = 110M$

$n = 2.1M, m = 440M$

Strong Scaling

Kronecker [1]



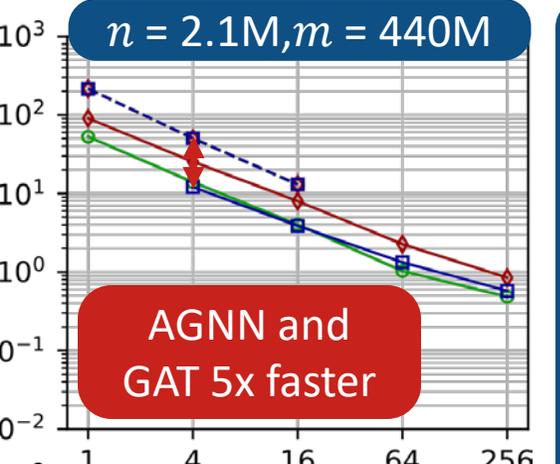
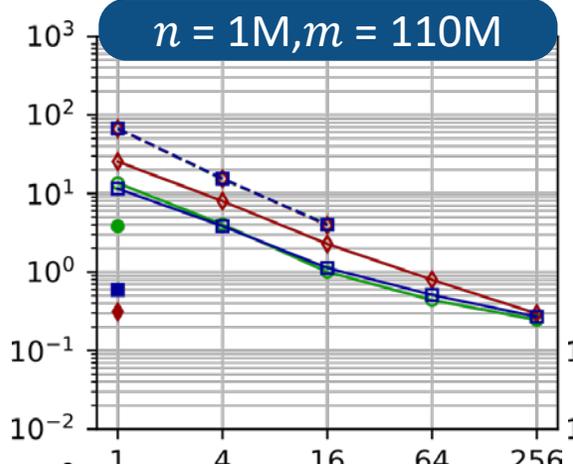
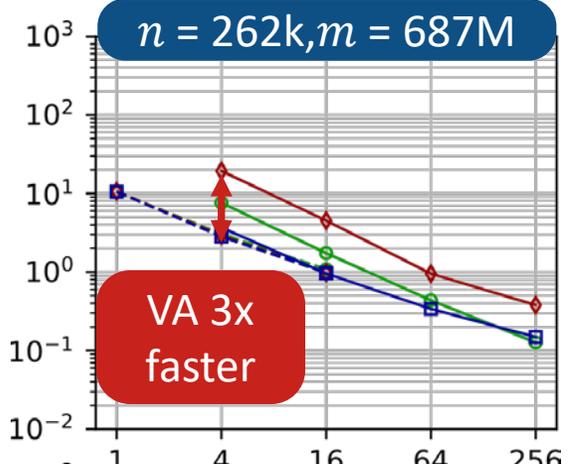
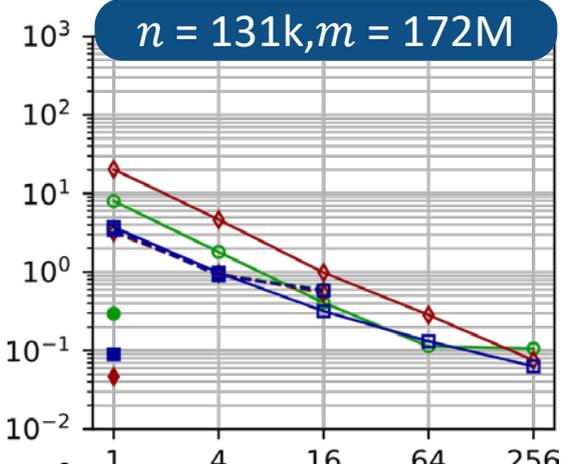
- GAT - This Work
- GAT - DistDGL
- GAT - DGL
- ◇— VA - This Work
- ◇- VA - DistDGL
- ◆ VA - DGL
- AGNN - This Work
- AGNN - DistDGL
- AGNN - DGL

$\rho = 1\%$

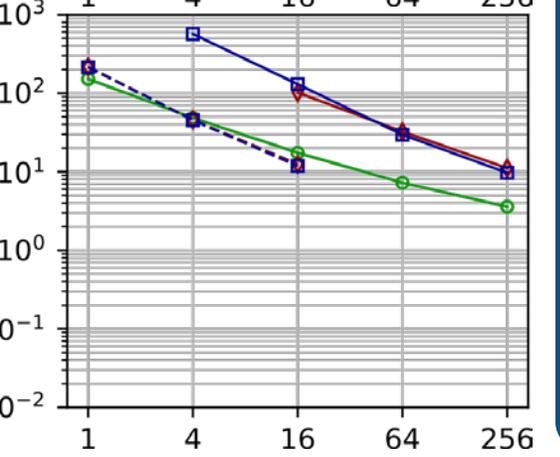
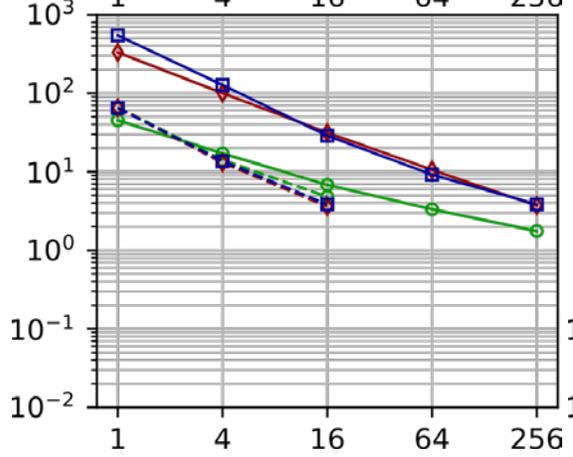
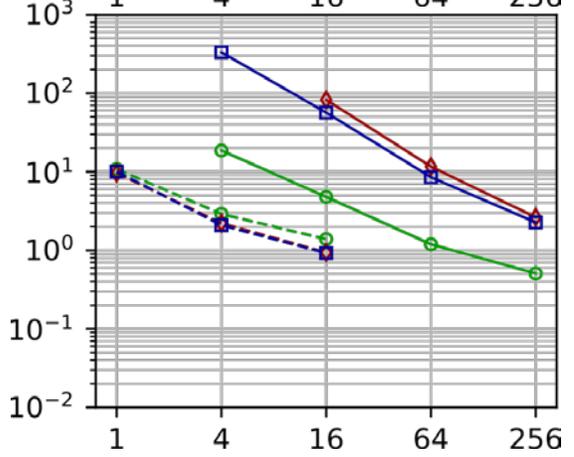
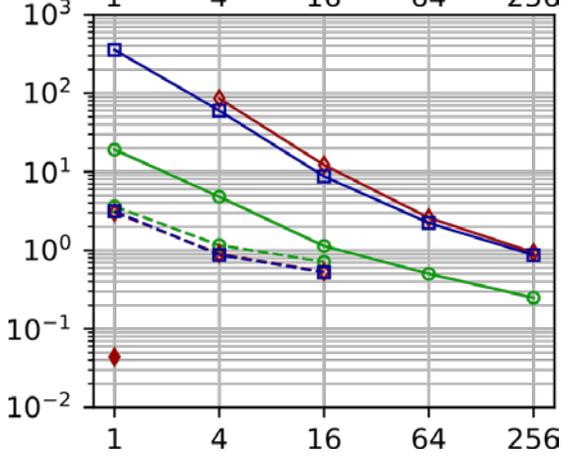
$\rho = 0.01\%$

k=16

runtime [s]



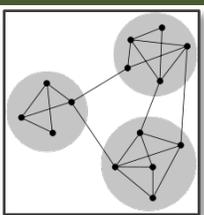
k=128



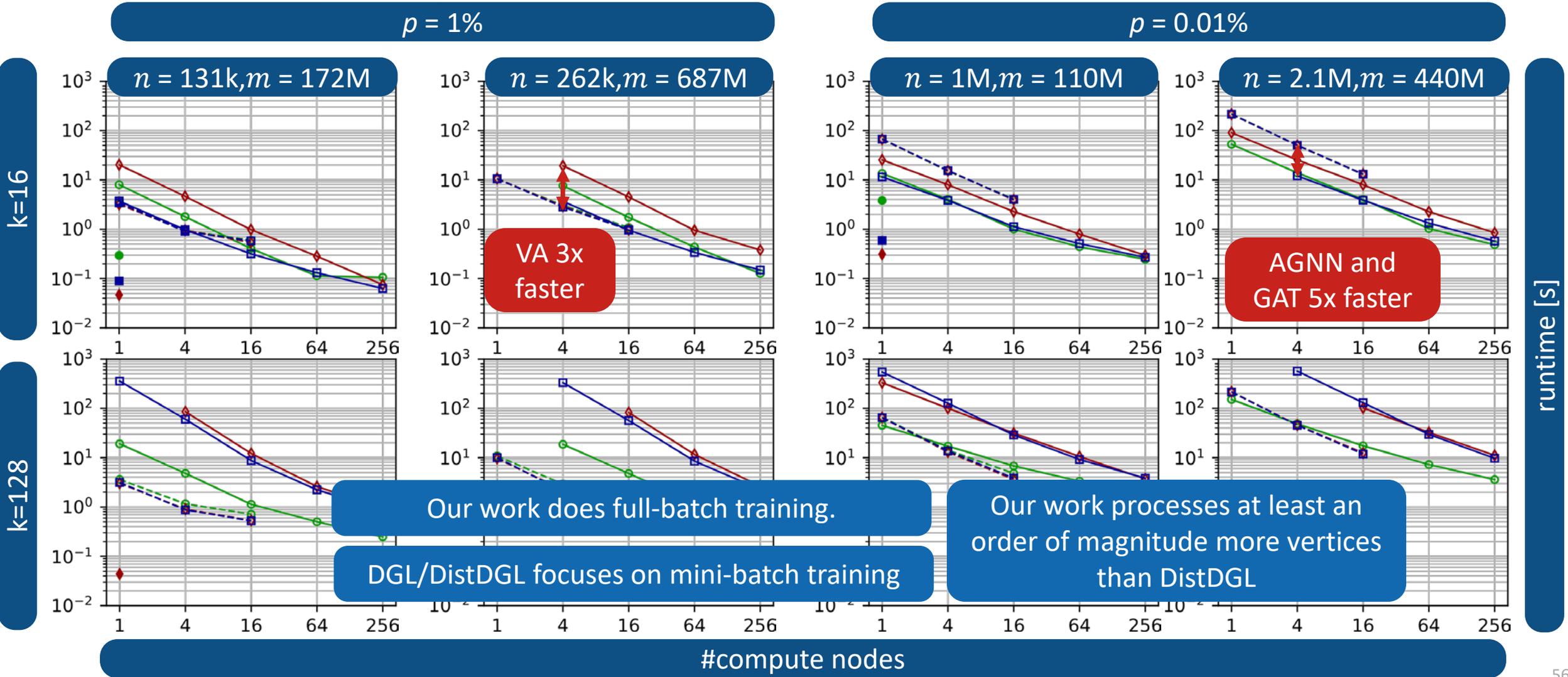
#compute nodes

Strong Scaling

Kronecker [1]



- GAT - This Work
- GAT - DistDGL
- GAT - DGL
- ◇— VA - This Work
- ◇- VA - DistDGL
- ◆ VA - DGL
- AGNN - This Work
- AGNN - DistDGL
- AGNN - DGL



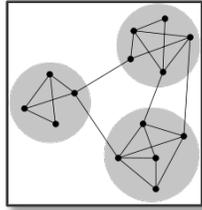
Our work does full-batch training.
DGL/DistDGL focuses on mini-batch training

Our work processes at least an order of magnitude more vertices than DistDGL

runtime [s]

Weak Scaling

Kronecker [1]



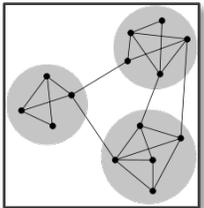
- GAT - This Work
- ◇— VA - This Work
- AGNN - This Work
- GAT - DistDGL
- ◇- VA - DistDGL
- AGNN - DistDGL

$p = 0,1\%$

$p = 0,01\%$

[1] J. Leskovec et al. Kronecker Graphs: An Approach to Modeling Networks. J. Mach. Learn. Research. 2010.

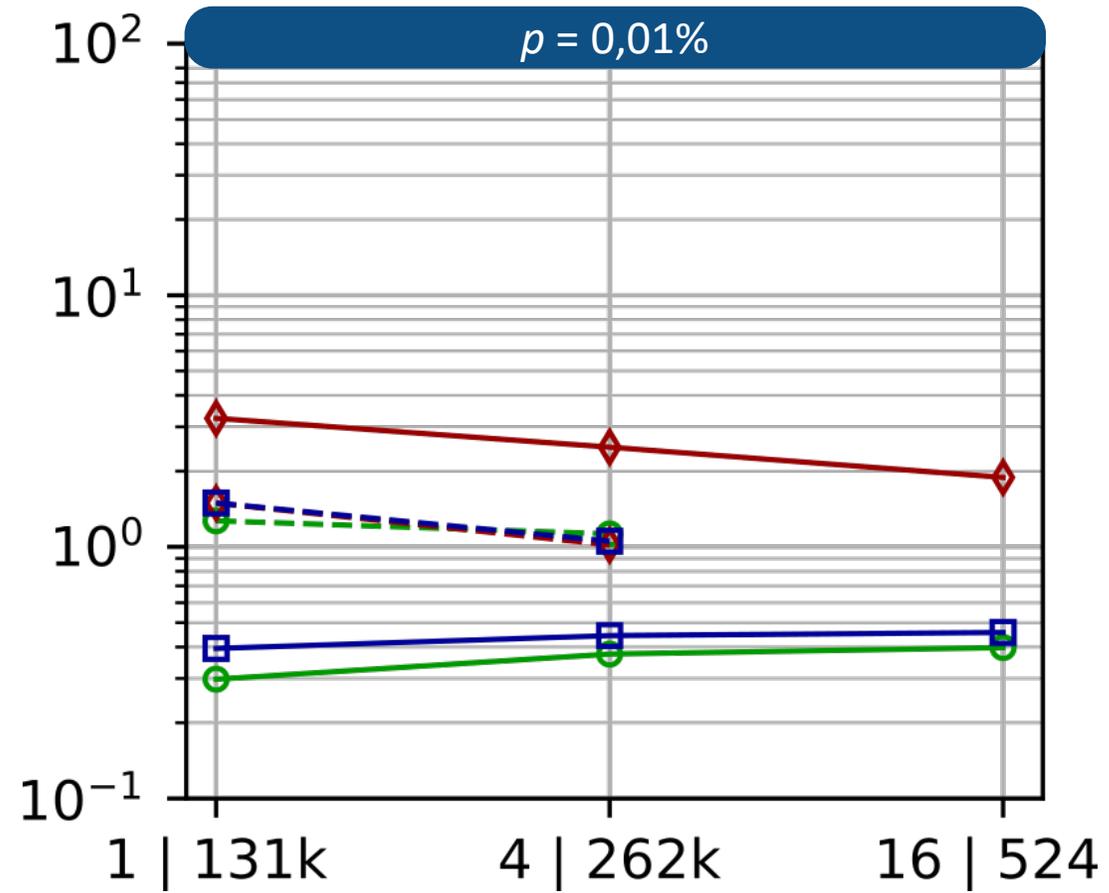
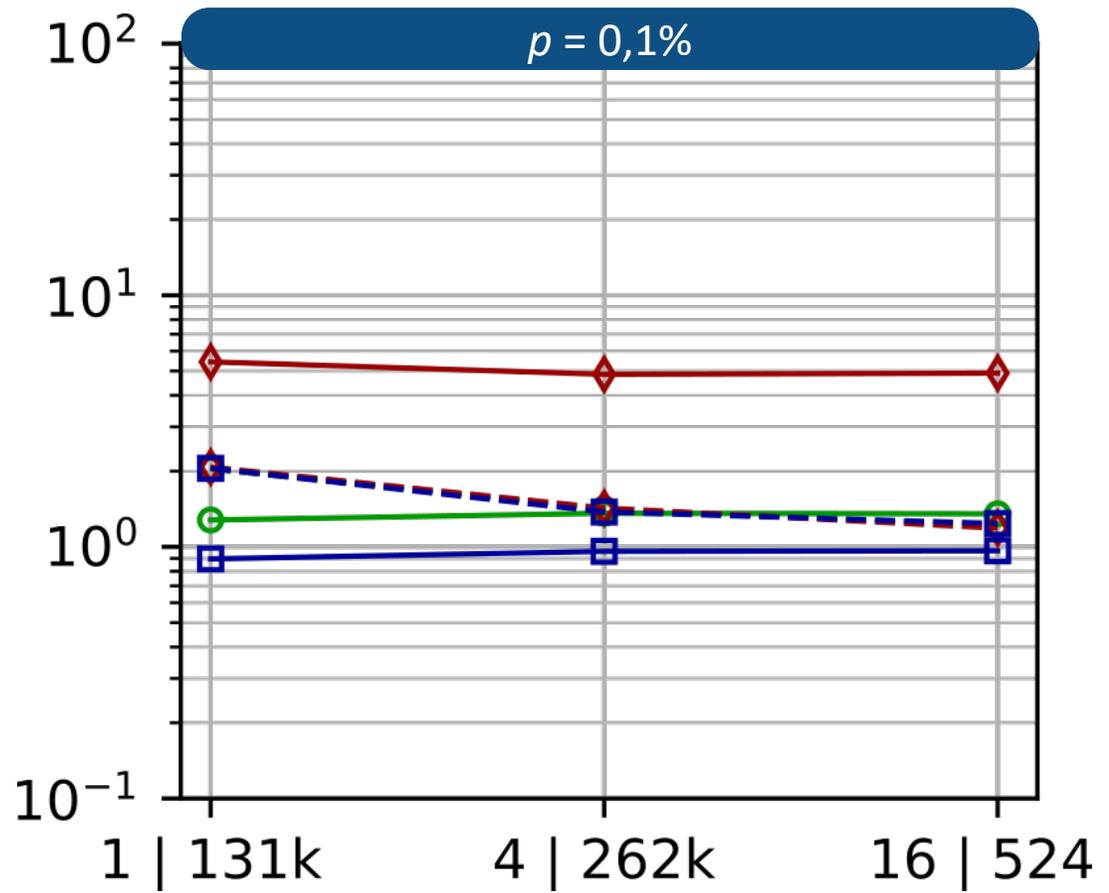
Weak Scaling



Kronecker [1]

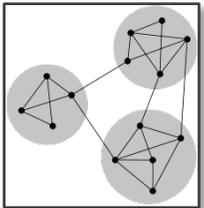
- GAT - This Work
- ◇— VA - This Work
- AGNN - This Work
- -○- - GAT - DistDGL
- -◇- - VA - DistDGL
- -□- - AGNN - DistDGL

runtime [s]



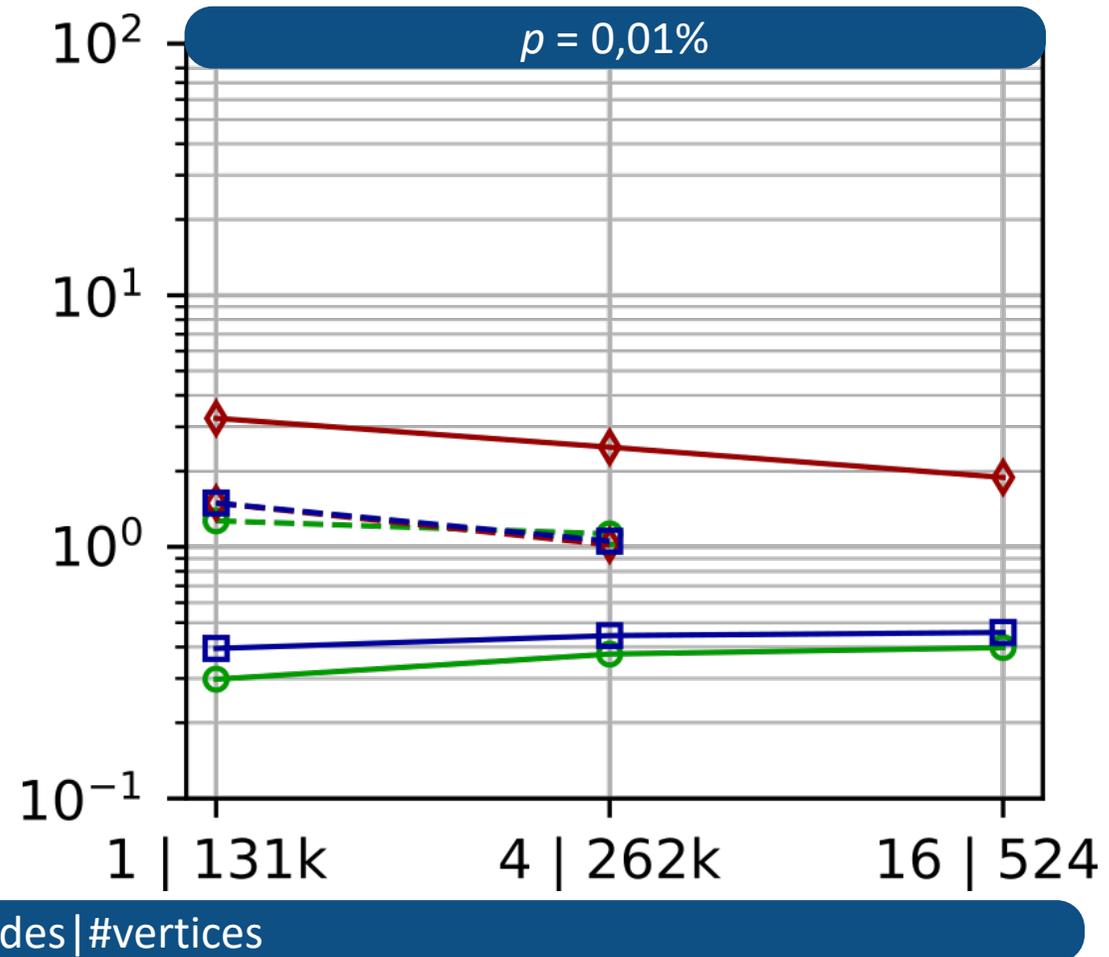
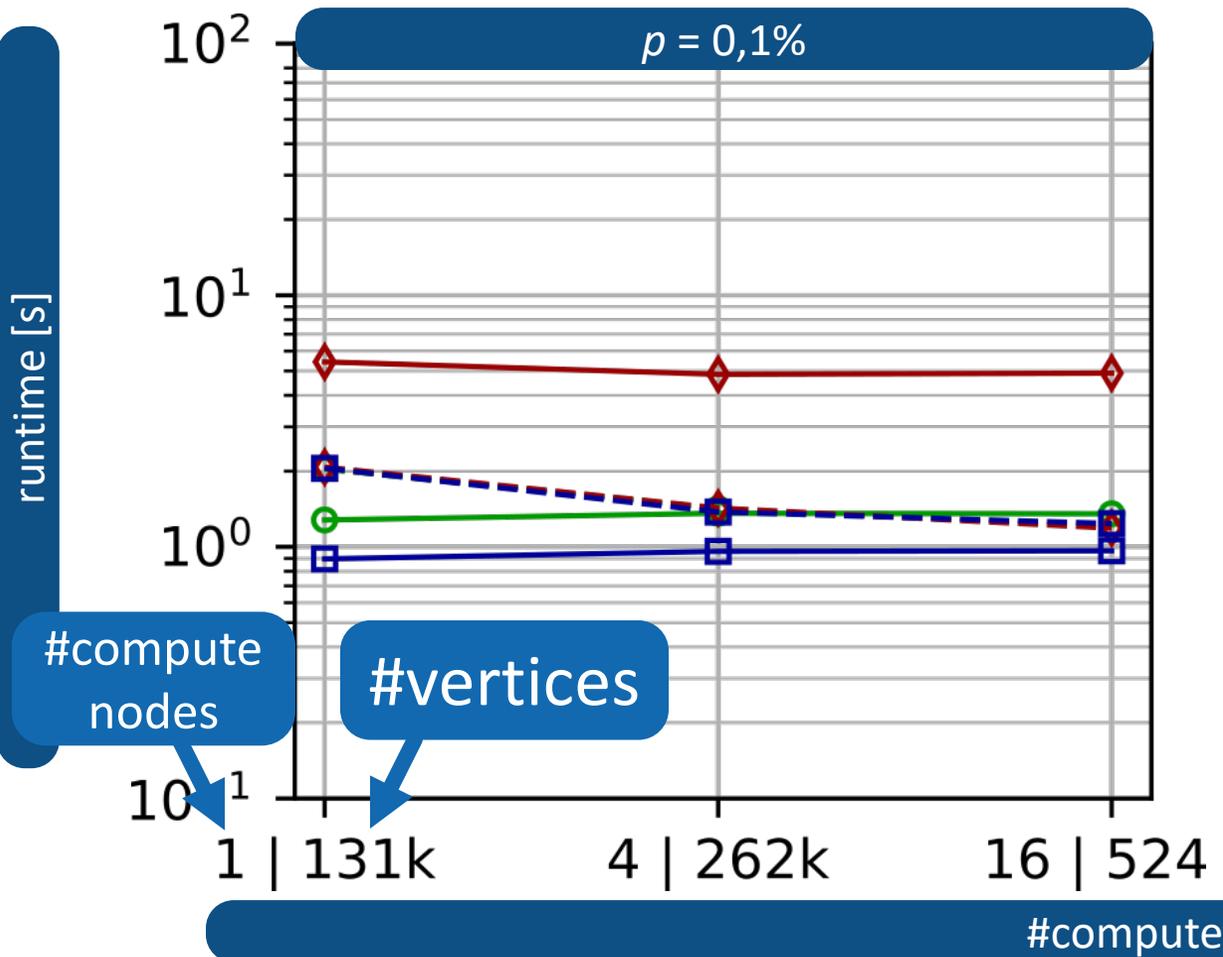
#compute nodes | #vertices

Weak Scaling



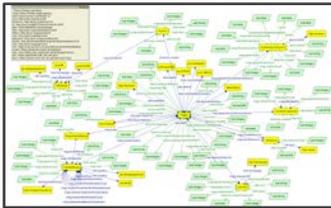
Kronecker [1]

- GAT - This Work
- ◇— VA - This Work
- AGNN - This Work
- -○- - GAT - DistDGL
- -◇- - VA - DistDGL
- -□- - AGNN - DistDGL



Strong Scaling

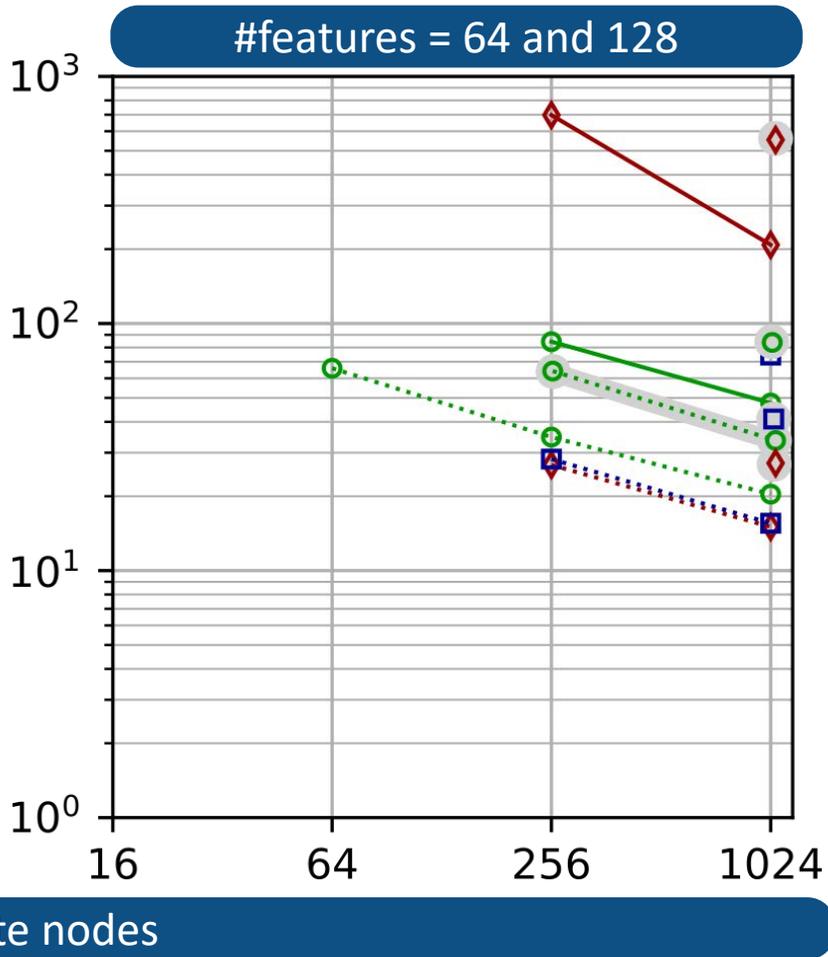
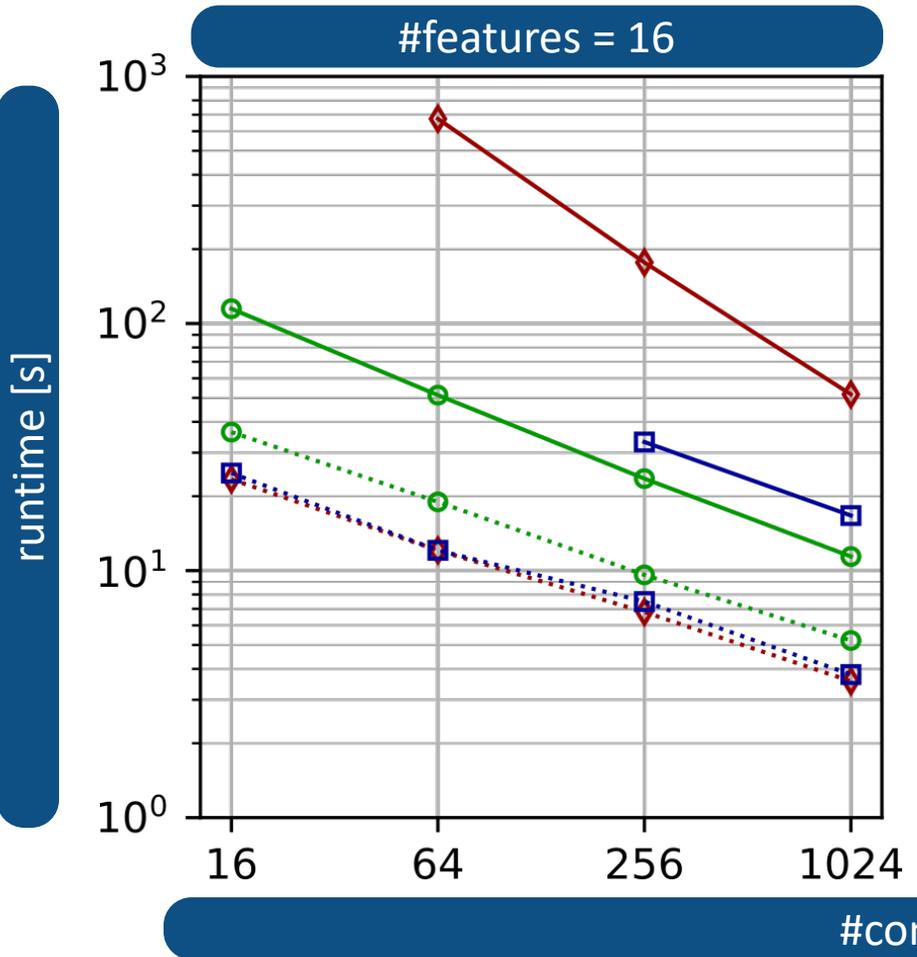
Microsoft Academic Knowledge Graph



Standard real world graph dataset

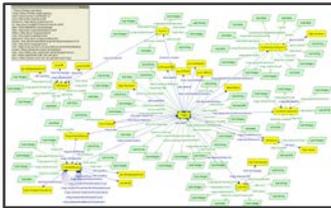
111 million vertices and 3.2 billion edges

- GAT - Training
- ◇— VA - Training
- AGNN - Training
- -○- - GAT - Inference
- -◇- - VA - Inference
- -□- - AGNN - Inference



Strong Scaling

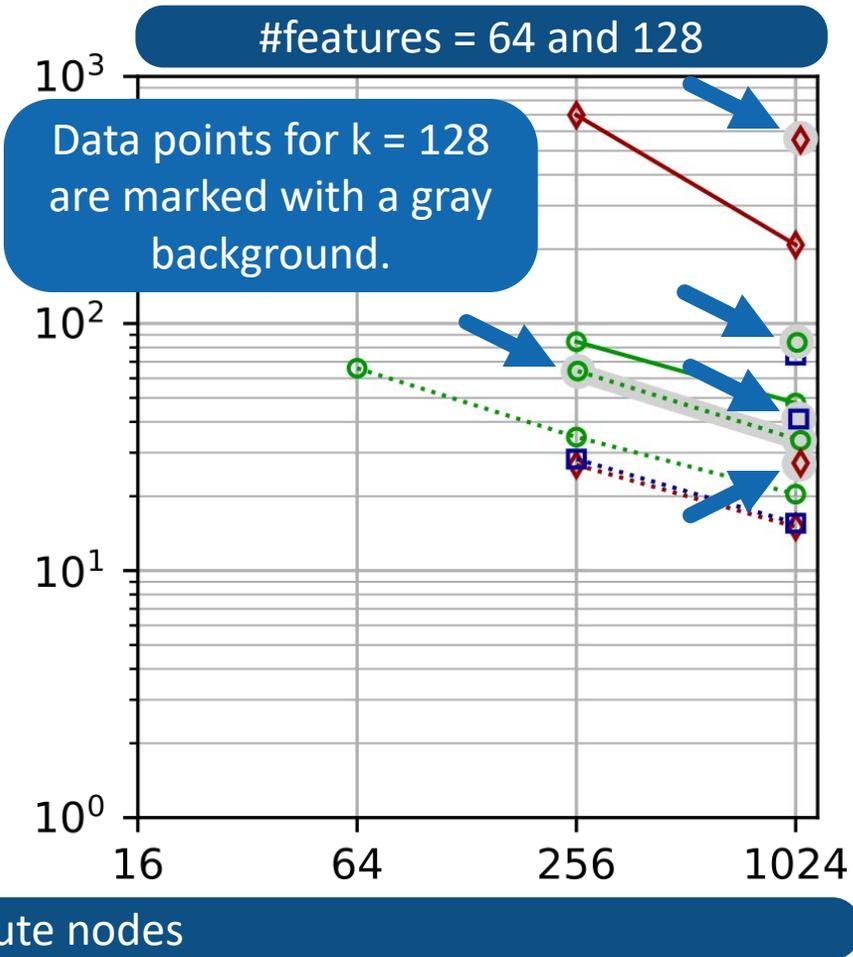
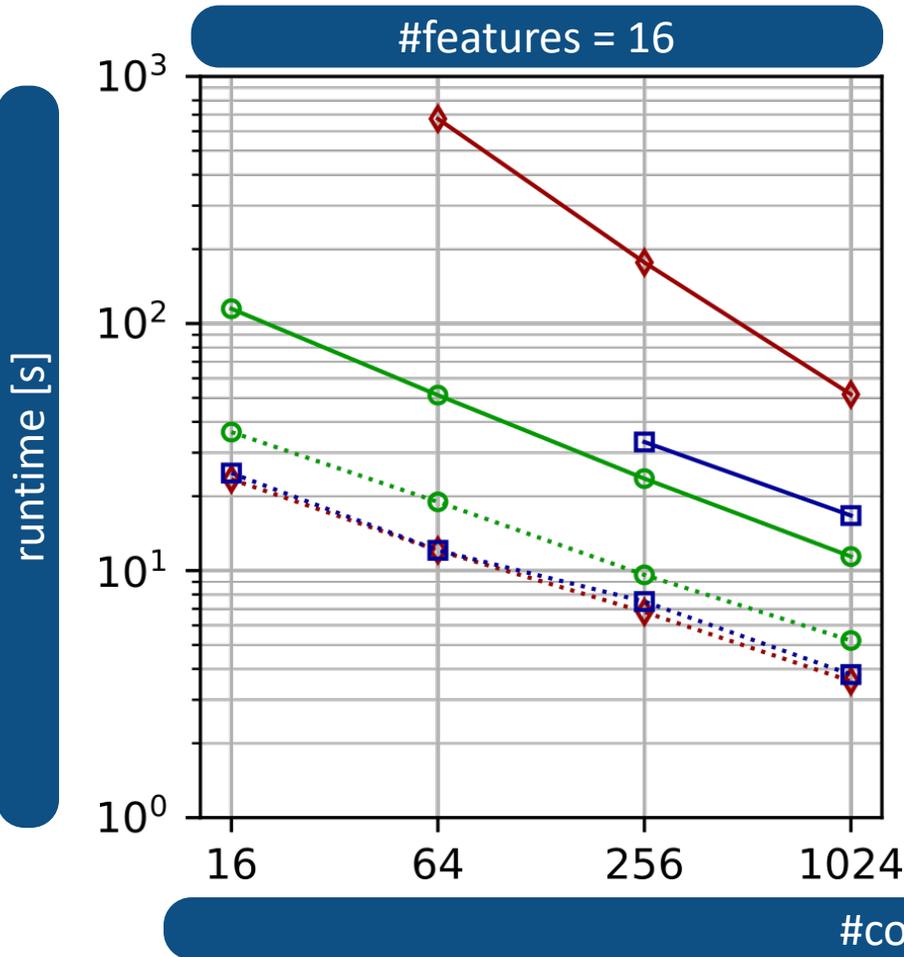
Microsoft Academic Knowledge Graph



Standard real world graph dataset

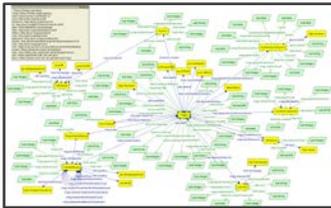
111 million vertices and 3.2 billion edges

- GAT - Training
- ◇— VA - Training
- AGNN - Training
- -○- - GAT - Inference
- -◇- - VA - Inference
- -□- - AGNN - Inference



Strong Scaling

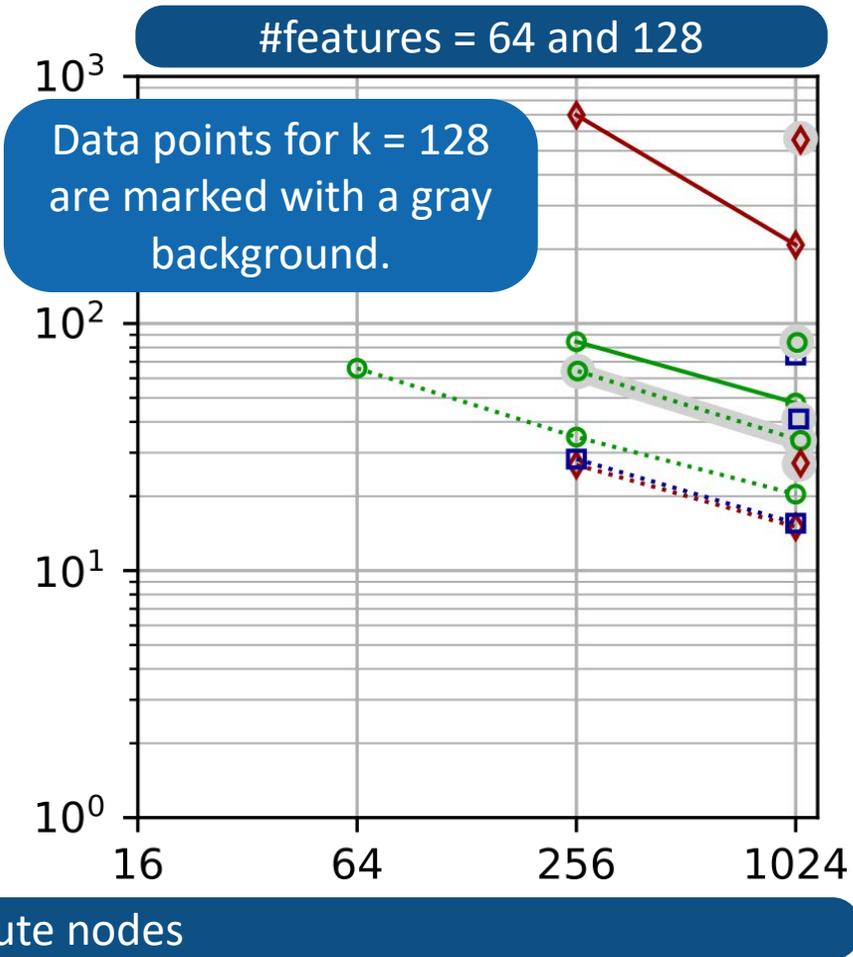
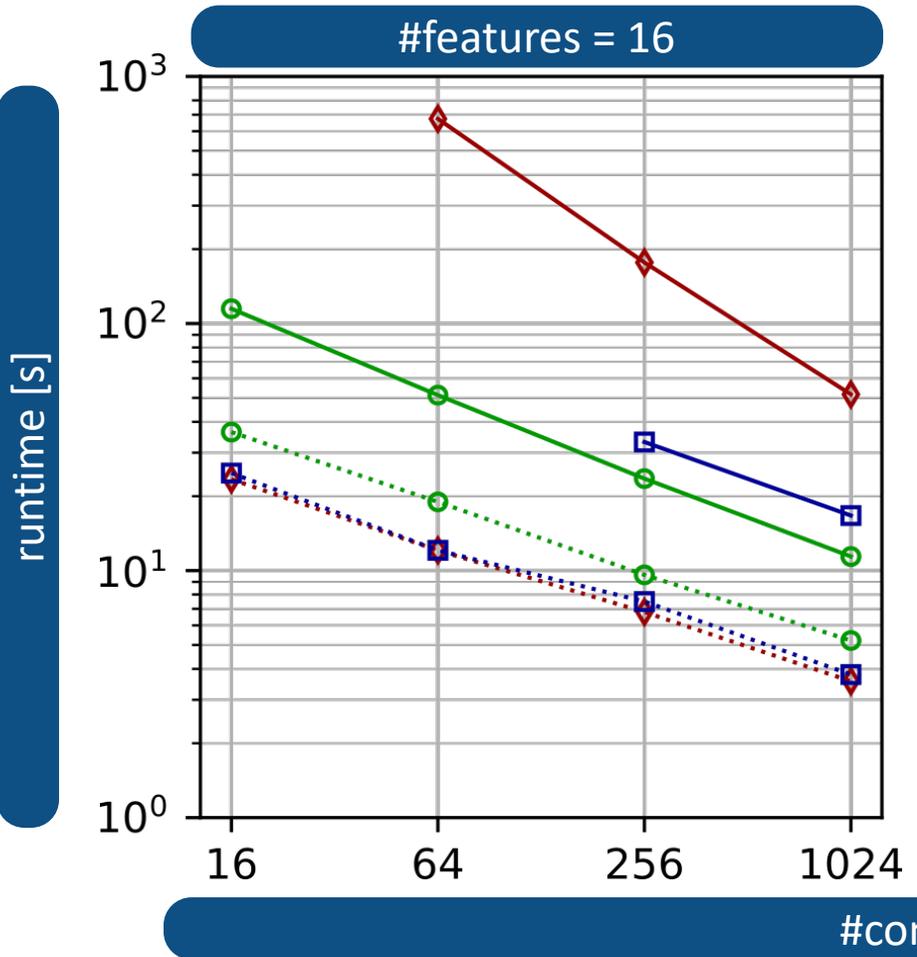
Microsoft Academic Knowledge Graph



Standard real world graph dataset

111 million vertices and 3.2 billion edges

- GAT - Training
- ◇— VA - Training
- AGNN - Training
- -○- - GAT - Inference
- -◇- - VA - Inference
- -□- - AGNN - Inference



For k = 64, all models except GAT require 256 compute nodes.

For k = 128, all models except GAT require 1,024 compute nodes.

Conclusions

More of SPCL's research:

Attention in GNN Models – Forward Pass

Formulating ψ is the „crux“ of devising a concrete formulation for a specific model

$$\mathbf{H}^{l+1} = \sigma(\mathbf{Z}), \quad \mathbf{Z} = \Psi \mathbf{H} \mathbf{W}$$

Non-linearity: σ
A sparse $n \times n$ tensor with attention scores, model specific: Ψ
Features from previous layer: \mathbf{H}
weights: \mathbf{W}

Vanilla Attention

 $\Psi = \mathcal{A} \odot \mathbf{H}_x$

SDDMM

Graph Attention Network (GAT)

 $\Psi = \text{sm}(\mathcal{T}) \quad \mathcal{T} = \mathcal{A} \odot \exp(\sigma(\mathbf{C}))$
 $\mathbf{C} = \text{rep}_n^T(\mathbf{H}'\mathbf{a})^T + \text{rep}_n(\mathbf{H}'\mathbf{a})$

SDDMM

Attention-based GNN (AGNN)

 $\Psi = \mathcal{A} \odot \mathbf{H}_x \otimes \mathbf{n}_x$

SDDMM

Global Formulations of GNN Kernels – Backward Pass

Generic formulation

$$\mathbf{G}^{l-1} = \sigma'(\mathbf{Z}^{l-1}) \odot \Gamma^l$$

$$\mathbf{Y}^l = \mathbf{H}^{lT} \Psi (\mathcal{A}^T, \mathbf{H}^l) \mathbf{G}^l + \mathbf{G}^l \mathbf{W}^{lT} \mathbf{H}^{lT} \frac{\partial \Psi}{\partial \mathbf{W}^l}$$

Matrix view

$$\mathbf{G} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \odot \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

MSPMM (or: MM + SpMM)

Graph Attention Network (GAT)

$$\psi_{v,u} = \frac{\exp(\sigma(\mathbf{a}^T \cdot \|\mathbf{W}\mathbf{h}_v\| \|\mathbf{W}\mathbf{h}_u\|))}{\sum_{i \in \mathcal{N}(v)} \exp(\sigma(\mathbf{a}^T \cdot \|\mathbf{W}\mathbf{h}_v\| \|\mathbf{W}\mathbf{h}_i\|))}$$

Local formulation

- multiply by a shared weight matrix
- concatenate
- dot product with a shared weight vector \mathbf{a}

Final score for an edge (v,u)

sum partial sums

Global formulation

Adjacency matrix grouping neighbors of all vertices

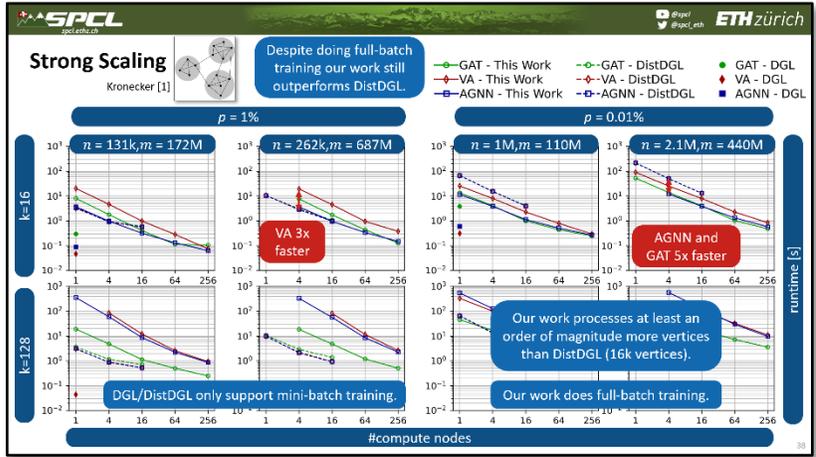
Replication to distribute weights across all edges

Shared weight matrix

Feature matrix grouping all feature vectors

Shared weight matrix

Partial sum



youtube.com/@spcl **175+ Talks**

twitter.com/spcl_eth **1.4K+ Followers**

github.com/spcl **2K+ Stars**

... or spcl.ethz.ch

