# Training Graph Neural Networks with 1000 Layers

3/4/2025

Arun Cheriakara Joseph,
Faculty Of Math

UNIVERSITY OF
WATERLOO

# Motivation

- Graphs Are Everywhere

- Real-World Examples

    - Netflix, Facebook, Amazon, Google Scholar

    - Internet = Vast Web Graph

UNIVERSITY OF
WATERLOO

# Why Graph Neural Networks (GNNs)?

- Designed for graph-structured data

- Message passing between nodes

- Impressive results on **small** graphs

UNIVERSITY OF
**WATERLOO**

# The Scaling Problem

- Memory bottleneck for large graphs

- Deep GNNs → More parameters → More memory

- Sampling/partitioning exist but not always ideal

UNIVERSITY OF
WATERLOO

# Inspiration from NLP & Efficient Architectures

- NLP success: GPT, BERT, GPT-3 (scale improves performance)

- Efficient methods in  NLP:

    - Reversible networks (RevNets)

    - Weight tying & equilibrium models

UNIVERSITY OF
WATERLOO

# Proposed Methods

- Reversible Connections

- Weight-Tied Networks

- Graph Equilibrium (Implicit) GNN

UNIVERSITY OF
WATERLOO

# Setup

- Input: A graph **G = (V, E)**

- N=|V|, M=|E|

- Adjacency matrix $A \in \mathbb{R}^{N \times N}$

- Node feature matrix $X \in \mathbb{R}^{N \times D}$

- (Optional) Edge feature matrix $U \in \mathbb{R}^{M \times F}$

- GNN operator: $f_w(X, A, U) \rightarrow X'$

UNIVERSITY OF
**WATERLOO**

# Over-parameterized GNNs

- Deep GNNs often use **residual connections**

  - $X' = f_w(X, A, U) + X$

- Memory cost for L-layer GNN: $O(LND)$

  - L is the number of GNN layers,

  - N is the number of vertices,

  - D is the size of vertex features.

- Activation storage > parameter storage

UNIVERSITY OF
**WATERLOO**

# Grouped Reversible GNN

- Inspired by **RevNets** + grouped convolutions

- Split features into C groups

- Reversible updates → No storing intermediate states

- Memory complexity → $O(ND)$

UNIVERSITY OF
**WATERLOO**

# Forward Pass in Grouped Reversible GNNs

- Divide input vertex feature matrix X into C groups.

- Compute an exchange term

- Process each group iteratively

$$X_0' = \sum_{i=2}^{C} X_i \tag{2}$$

$$X_i' = f_{w_i}(X_{i-1}', A, U) + X_i, \; i \in \{1, \cdots, C\}, \tag{3}$$

UNIVERSITY OF
WATERLOO

# Efficient Backpropagation

- Recompute $X_i$

- Recompute $X_0'$

- Reconstruct $X_1$

$$X_i = X_i' - f_{w_i}(X_{i-1}', A, U), \ i \in \{2, \cdots, C\} \quad (4)$$

$$X_0' = \sum_{i=2}^{C} X_i \quad (5)$$

$$X_1 = X_1' - f_{w_1}(X_0', A, U). \quad (6)$$

UNIVERSITY OF
WATERLOO

# Avoiding Dropout Issues

- Dropout adds randomness, making reconstruction difficult.

- Solution: Use **shared dropout patterns** across all layers.

- This prevents memory from scaling with depth.

$$Memory\ Complexity: O(ND)$$

UNIVERSITY OF
**WATERLOO**

# Memory Complexity & Practical Benefits

- Activations: $O(ND)$, not $O(LND)$

- Works with GCN, GAT, etc.

- Enables 1000+ layers on one GPU

UNIVERSITY OF
**WATERLOO**

# Weight-Tied GNNs: The Idea

- Instead of learning different weights for each layer, we **reuse the same function** across layers.

- Reduces parameter count, making the model more **memory-efficient** and **generalizable**.

$$f_w^{(1)} := f_w^{(2)} \ldots := f_w^{(L)}, \qquad (9)$$

UNIVERSITY OF
**WATERLOO**

# Two Types of Weight-Tied GNNs

- **Weight-Tied Residual GNNs**:

  1. Shares weights across layers.

  **2. Still requires storing activations**, so memory complexity remains **O(LND)**.

- **Weight-Tied Reversible GNNs**:

  1. Combines weight-tying with **reversibility**.

  2. Memory is reduced to **O(ND)** (independent of depth).

UNIVERSITY OF
**WATERLOO**

# Extending Weight-Tying to Groups

- Instead of weight-tying for all features, we divide features into **C groups**.

- Each group has its own tied function:

$$f_{w_i}^{(1)} := f_{w_i}^{(2)} \ldots := f_{w_i}^{(L)}, \; i \in \{1, \cdots, C\} \qquad (10)$$

- Helps retain expressiveness while reducing parameters.

UNIVERSITY OF
**WATERLOO**

# Weight-Tied Memory Complexity

| Model Type | Memory Complexity |
|---|---|
| Regular Deep GNN | O(LND) |
| Weight-Tied Residual GNN | O(LND) |
| Weight-Tied Reversible GNN | O(ND) |

UNIVERSITY OF
**WATERLOO**

# Deep Equilibrium GNN: Overview

- Instead of explicitly stacking layers, the model **converges to a fixed point**.

- The equilibrium equation:

$$Z^* = f_w^{\text{DEQ}}(Z^*, X, A, U), \qquad (11)$$

- The network **iterates until it reaches this stable state**.

- No need to define a fixed depth for the model.

UNIVERSITY OF
**WATERLOO**

# DEQ-GNN Equations

- The forward pass of DEQ-GNN is implemented with a root-finding algorithm (e.g. Broyden's method)

- The gradients are obtained by implicitly differentiating through the equilibrium node state for the backward pass.

$$Z' = \text{GraphConv}(Z_{\text{in}}, A, U) \tag{12}$$

$$Z'' = \text{Norm}(Z' + X) \tag{13}$$

$$Z''' = \text{GraphConv}(\text{Dropout}(\text{ReLU}(Z'')), A, U) \tag{14}$$

$$Z_{\text{o}} = \text{Norm}(\text{ReLU}(Z''' + Z')), \tag{15}$$

UNIVERSITY OF
WATERLOO

# Memory Efficiency in DEQ-GNNs

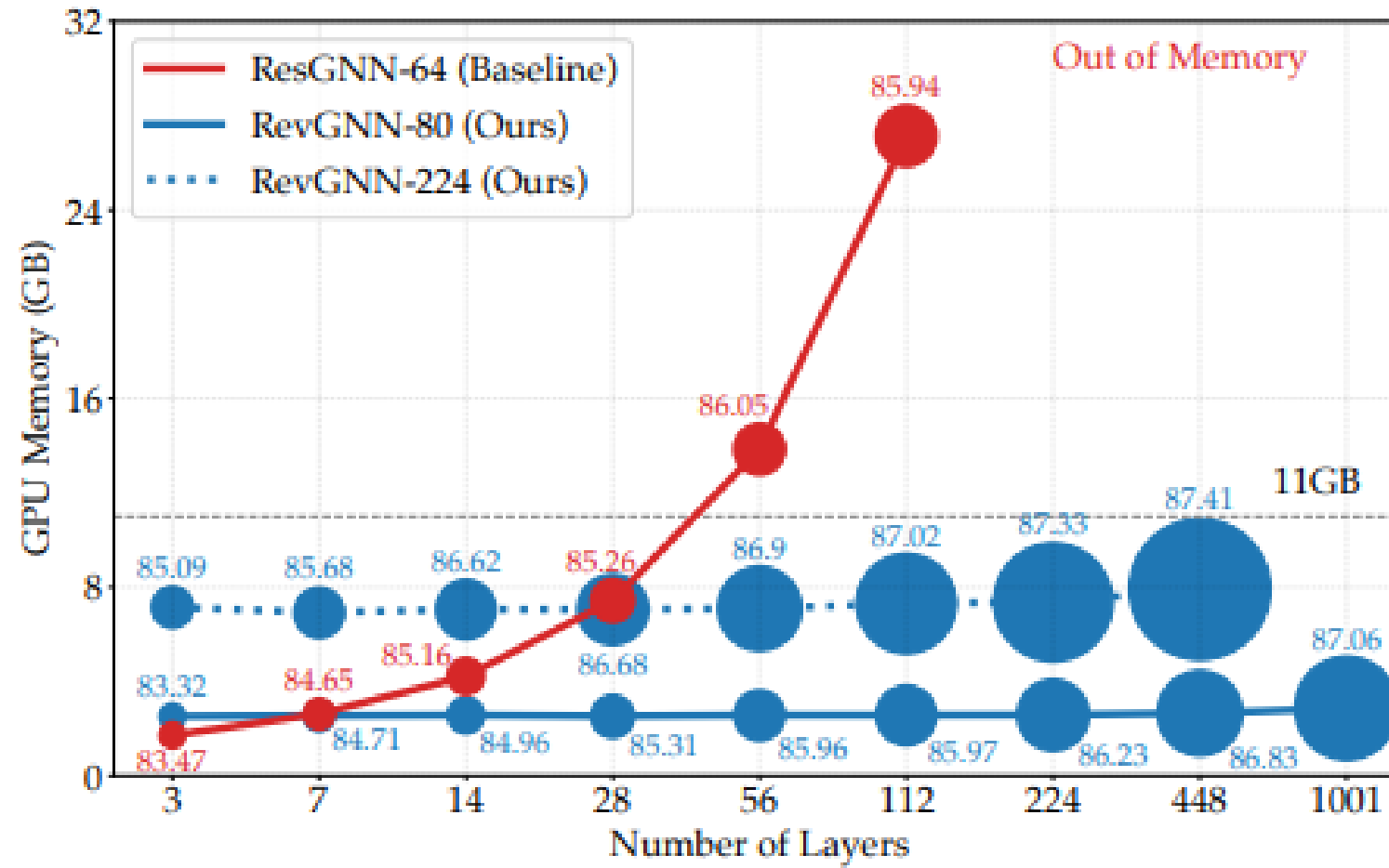| Model Type | Memory Complexity |
|---|---|
| Regular Deep GNN | O(LND) |
| DEQ-GNN | O(ND) |

UNIVERSITY OF
**WATERLOO**

# Analysis of Different Deep GNN Architectures

- **Dataset:** ogbn-proteins (OGB benchmark)

- **Training method:** Mini-batch training with random partitioning

- **Baseline Model:** ResGNN (Li et al., 2020)

- **Metric:** ROC-AUC score (higher is better)

UNIVERSITY OF
**WATERLOO**

# Baseline GNN – ResGNN / Reversible GNN - RevGNN

- **Baseline** ResGNN

  - It's a **pre-activation residual** GNN

- Reversible GNN - RevGNN

  - No need to store all activations

  - Constant memory usage

  - Scales beyond 1000 layers
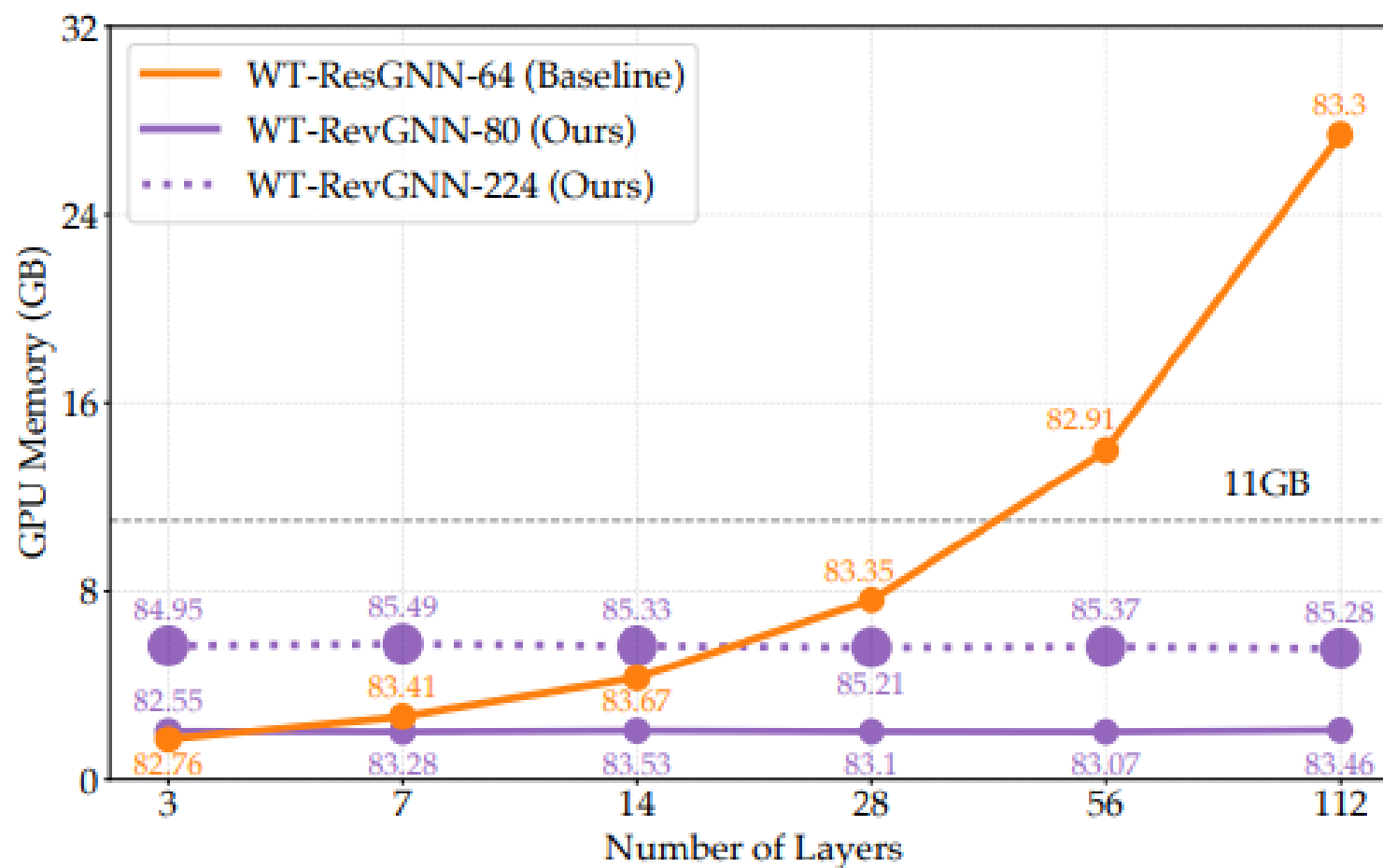
UNIVERSITY OF
**WATERLOO**

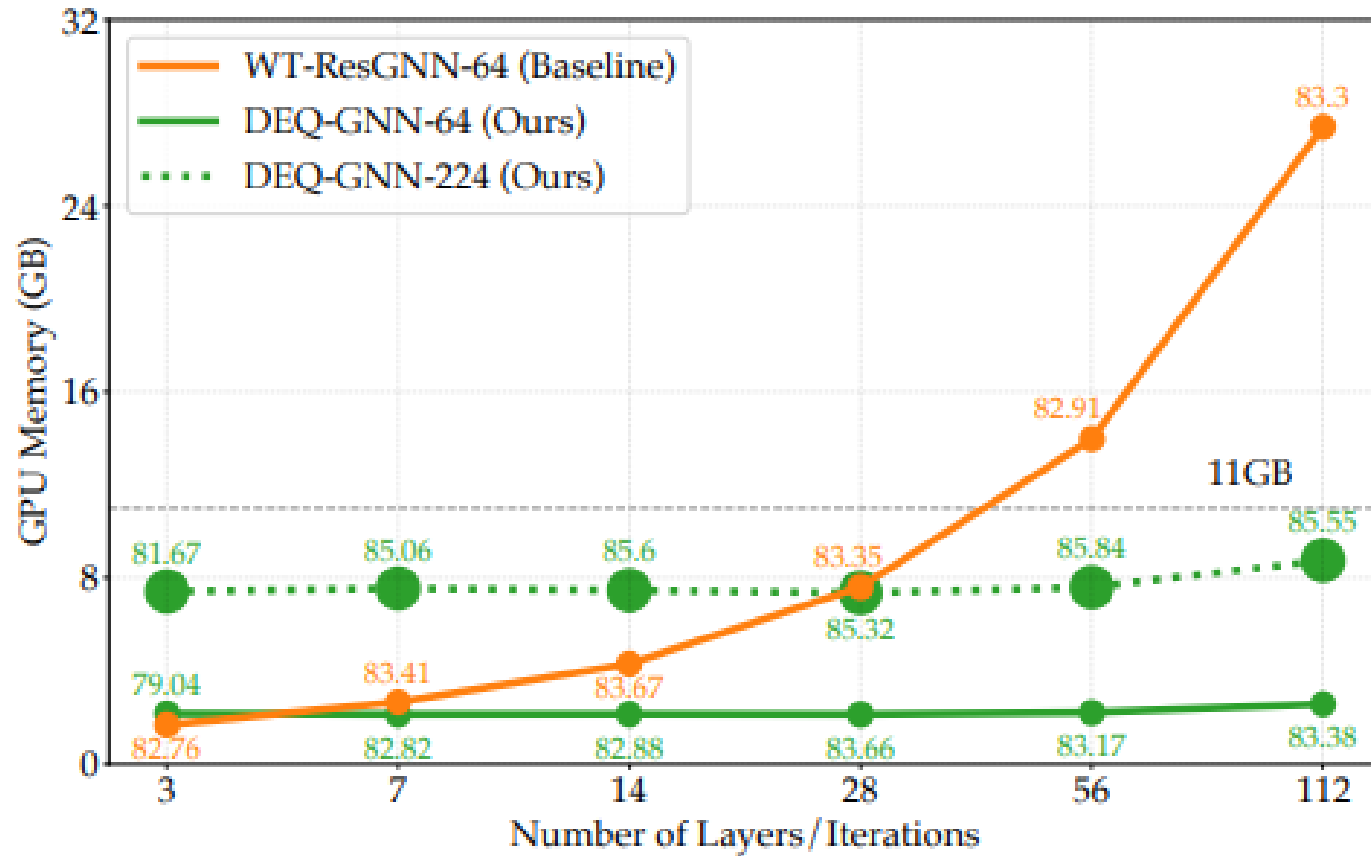# ResGNN VS RevGNN

UNIVERSITY OF
**WATERLOO**

# WT-ResGNN vs. WT-RevGNN

- **Sharing the same parameters** across all layers

- Parameter count **does not grow** with depth

- **WT-ResGNN**: Residual, weight-tied

- **WT-RevGNN**: Reversible, weight-tied

- Same number of parameters, **different memory usage**

UNIVERSITY OF
**WATERLOO**

# WT-ResGNN vs. WT-RevGNN

# Equilibrium GNN

UNIVERSITY OF
WATERLOO

# Over-parameterized Deep GNNs

- Experiments on **OGB** datasets (ogbn-proteins, ogbn-arxiv)

- Show **state-of-the-art** results with reversible GNN

UNIVERSITY OF
**WATERLOO**

# ogbn-Protein Results

Table 1. **Results on the *ogbn-proteins* dataset compared to SOTA.** RevGNN-Deep has 1001 layers with 80 channels each. It achieves SOTA performance with minimal GPU memory for training. RevGNN-Wide has 448 layers with 224 channels each. It achieves the best accuracy while consuming a moderate amount of GPU memory.

| Model | ROC-AUC ↑ | Mem ↓ | Params |
|---|---|---|---|
| GCN (Kipf & Welling) | 72.51 ± 0.35 | 4.68 | 96.9k |
| GraphSAGE (Hamilton et al.) | 77.68 ± 0.20 | 3.12 | 193k |
| DeeperGCN (Li et al.) | 86.16 ± 0.16 | 27.1 | 2.37M |
| UniMP (Shi et al.) | 86.42 ± 0.08 | 27.2 | 1.91M |
| GAT (Veličković et al.) | 86.82 ± 0.21 | 6.74 | 2.48M |
| UniMP+CEF (Shi et al.) | 86.91 ± 0.18 | 27.2 | 1.96M |
| Ours (RevGNN-Deep) | 87.74 ± 0.13 | **2.86** | 20.03M |
| Ours (RevGNN-Wide) | **88.24** ± 0.15 | 7.91 | 68.47M |

UNIVERSITY OF
WATERLOO

# ogbn-arxiv Results

Table 2. **Results on the *ogbn-arxiv* dataset compared to SOTA.** RevGCN-Deep has 28 layers with 128 channels each. It achieves SOTA performance with minimal GPU memory. RevGAT-Wide has 5 layers with 1068 channels each. RevGAT-SelfKD denotes the student models with 5 layers and 768 channels. It achieves the best accuracy while consuming a moderate amount of GPU memory.

| Model | ACC ↑ | Mem ↓ | Params |
|---|---|---|---|
| GraphSAGE (Hamilton et al.) | $71.49 \pm 0.27$ | 1.99 | 219k |
| GCN (Kipf & Welling) | $71.74 \pm 0.29$ | 1.90 | 143k |
| DAGNN (Liu et al.) | $72.09 \pm 0.25$ | 2.40 | 43.9k |
| DeeperGCN (Li et al.) | $72.32 \pm 0.27$ | 21.6 | 491k |
| GCNII (Chen et al.) | $72.74 \pm 0.16$ | 17.0 | 2.15M |
| GAT (Veličković et al.) | $73.91 \pm 0.12$ | 5.52 | 1.44M |
| UniMP_v2 (Shi et al.) | $73.97 \pm 0.15$ | 25.0 | 687k |
| Ours (RevGCN-Deep) | $73.01 \pm 0.31$ | **1.84** | 262k |
| Ours (RevGAT-Wide) | $74.05 \pm 0.11$ | 8.49 | 3.88M |
| Ours (RevGAT-SelfKD) | $\mathbf{74.26} \pm 0.17$ | 6.60 | 2.10M |

UNIVERSITY OF
**WATERLOO**

# ogbn-arxiv Results on different GNN operators

Table 3. **Results with different GNN operators on the ogbn-arxiv.** All GAT models use label propagation. #L and #Ch denote the number of layers and channels respectively. *Baselines are in italic.*
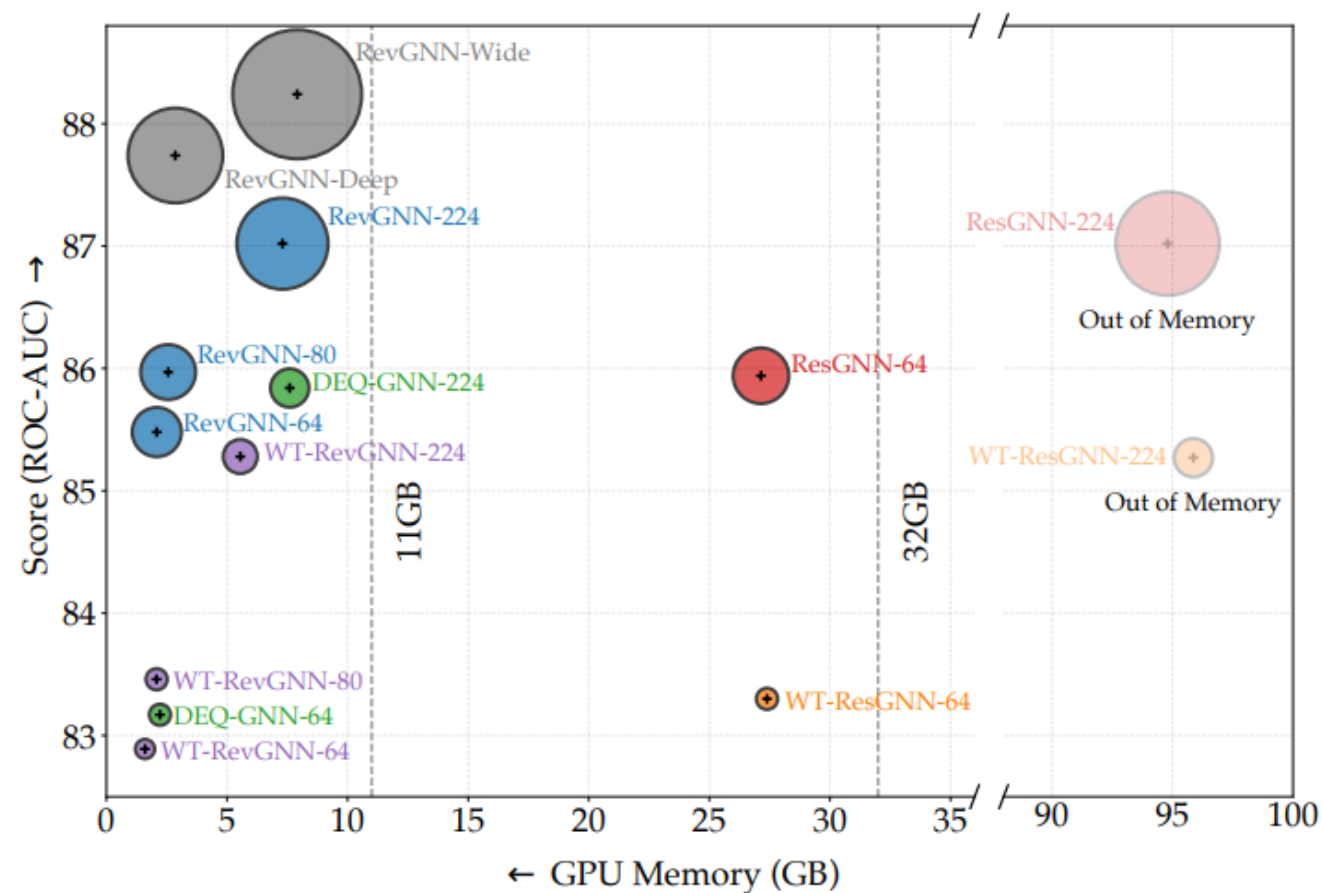
| Model | #L | #Ch | ACC ↑ | Mem ↓ | Params |
|---|---|---|---|---|---|
| *ResGCN* | 28 | 128 | 72.46 ± 0.29 | 11.15 | 491k |
| RevGCN | 28 | 128 | 73.01 ± 0.31 | **1.84** | 262k |
| RevGCN | 28 | 180 | **73.22** ± 0.19 | 2.73 | 500k |
| *ResSAGE* | 28 | 128 | 72.46 ± 0.29 | 8.93 | 950k |
| RevSAGE | 28 | 128 | 72.69 ± 0.23 | **1.17** | 491k |
| RevSAGE | 28 | 180 | **72.73** ± 0.10 | 1.57 | 953k |
| *ResGEN* | 28 | 128 | 72.32 ± 0.27 | 21.63 | 491k |
| RevGEN | 28 | 128 | 72.34 ± 0.18 | **4.08** | 262k |
| RevGEN | 28 | 180 | **72.93** ± 0.10 | 5.67 | 500k |
| *ResGAT* | 5 | 768 | 73.76 ± 0.13 | 9.96 | 3.87M |
| RevGAT | 5 | 768 | 74.02 ± 0.18 | **6.30** | 2.10M |
| RevGAT | 5 | 1068 | **74.05** ± 0.11 | 8.49 | 3.88M |

UNIVERSITY OF
**WATERLOO**

# Analysis of Complexities

Table 4. **Comparison of complexities.** $L$ is the number of layers, $D$ is the number of hidden channels, $N$ is of the number of nodes, $B$ is the batch size of nodes and $R$ is the number of sampled neighbors of each node. $K$ is the maximum Broyden iterations.

| Method | Memory | Params | Time |
|---|---|---|---|
| Full-batch GNN | $\mathcal{O}(LND)$ | $\mathcal{O}(LD^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2)$ |
| GraphSAGE | $\mathcal{O}(R^L BD)$ | $\mathcal{O}(LD^2)$ | $\mathcal{O}(R^L ND^2)$ |
| VR-GCN | $\mathcal{O}(LND)$ | $\mathcal{O}(LD^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2 + R^L ND^2)$ |
| FastGCN | $\mathcal{O}(LRBD)$ | $\mathcal{O}(LD^2)$ | $\mathcal{O}(RLND^2)$ |
| Cluster-GCN | $\mathcal{O}(LBD)$ | $\mathcal{O}(LD^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2)$ |
| GraphSAINT | $\mathcal{O}(LBD)$ | $\mathcal{O}(LD^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2)$ |
| Weight-tied GNN | $\mathcal{O}(LND)$ | $\mathcal{O}(D^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2)$ |
| RevGNN | $\mathcal{O}(ND)$ | $\mathcal{O}(LD^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2)$ |
| WT-RevGNN | $\mathcal{O}(ND)$ | $\mathcal{O}(D^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2)$ |
| DEQ-GNN | $\mathcal{O}(ND)$ | $\mathcal{O}(D^2)$ | $\mathcal{O}(K\|A\|_0 D + KND^2)$ |
| RevGNN + Subgraph Sampling | $\mathcal{O}(BD)$ | $\mathcal{O}(LD^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2)$ |
| WT-RevGNN + Subgraph Sampling | $\mathcal{O}(BD)$ | $\mathcal{O}(D^2)$ | $\mathcal{O}(L\|A\|_0 D + LND^2)$ |
| DEQ-GNN + Subgraph Sampling | $\mathcal{O}(BD)$ | $\mathcal{O}(D^2)$ | $\mathcal{O}(K\|A\|_0 D + KND^2)$ |

UNIVERSITY OF
**WATERLOO**

# Results

UNIVERSITY OF
**WATERLOO**

Thank You!