



Analyzing the Performance of Graph Neural Networks with Pipe Parallelism

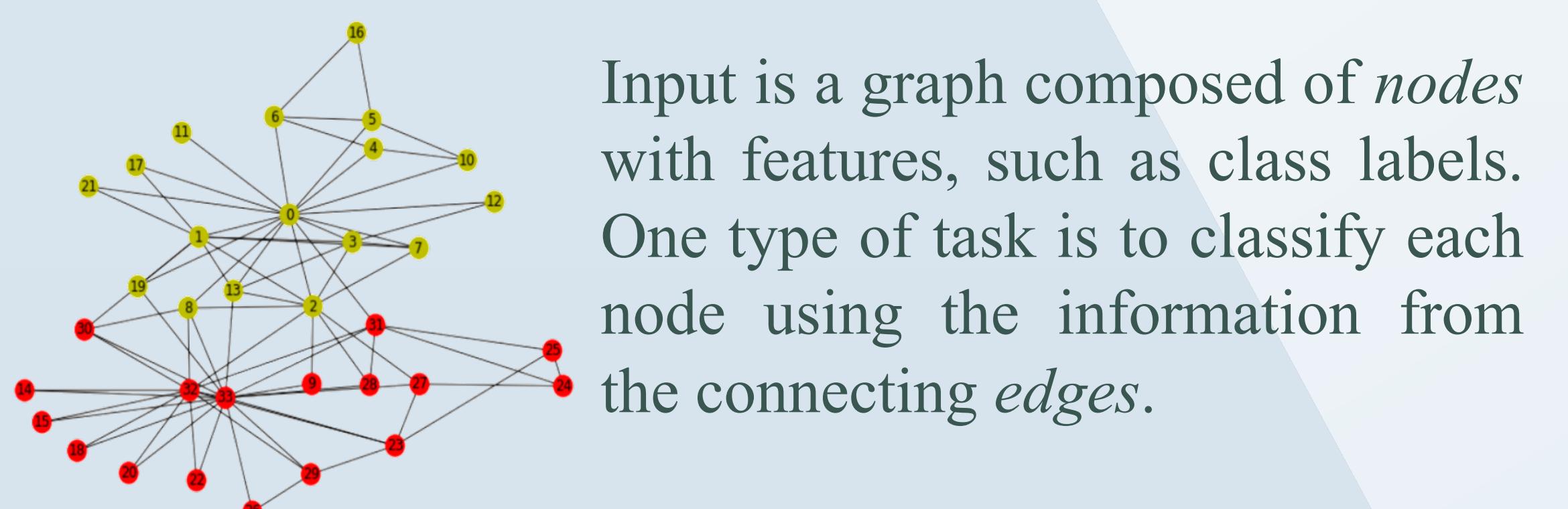
Matthew T. Dearing & Xiaoyan (Angela) Wang

MLSys

Overview & Motivation

- Many datasets in the real world can be structured as *graphs* or *networks*, such as social groups, chemical interactions, protein models, and phase transitions.
- Graph neural networks (GNN) generalize neural networks for deep learning on arbitrary irregular graphs.
- Because current GNN models do not scale well and datasets of interest continue to increase in complexity and size, improving computational efficiency is important to further research applicability.
- We study parallelizing GNNs using *existing* tools and frameworks known to be successful for traditional deep learning.

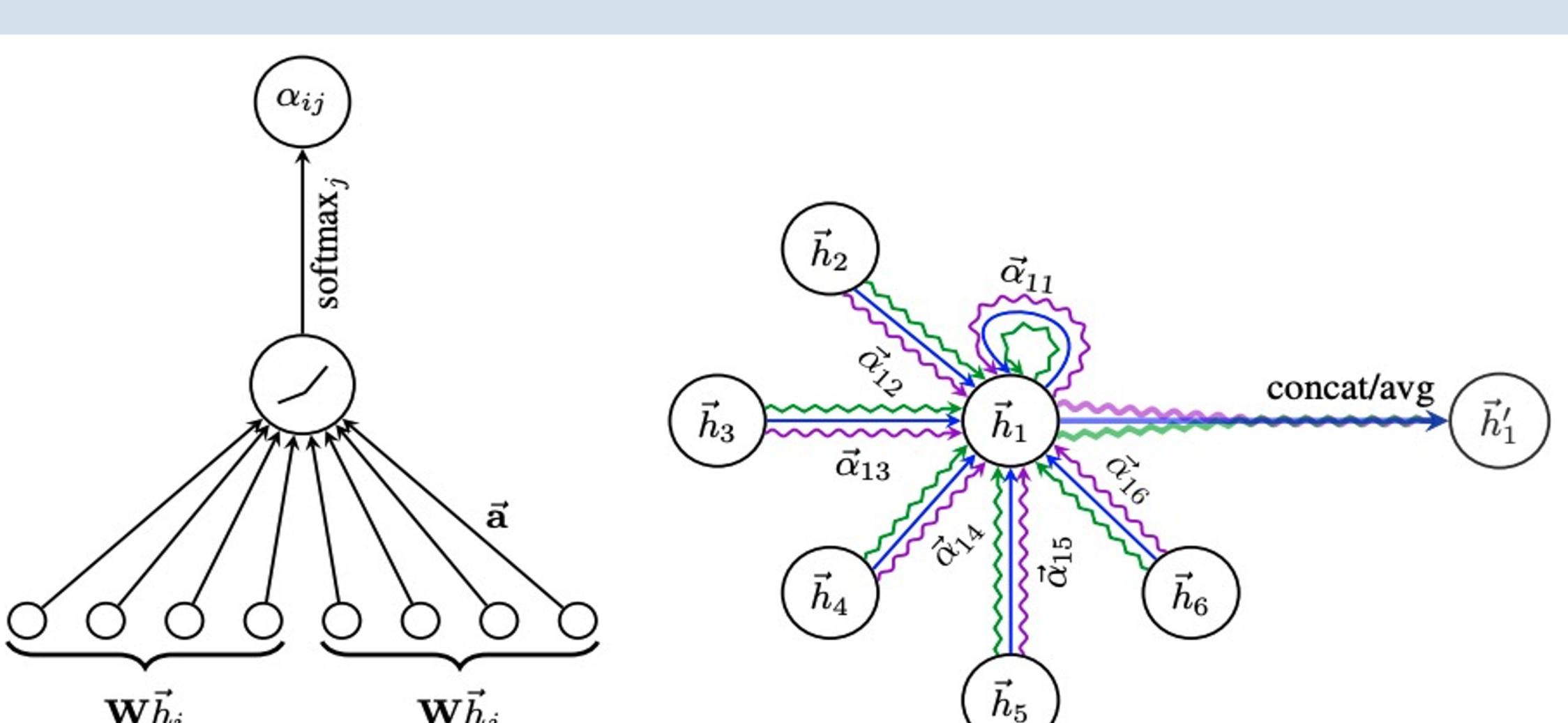
Graph Neural Networks



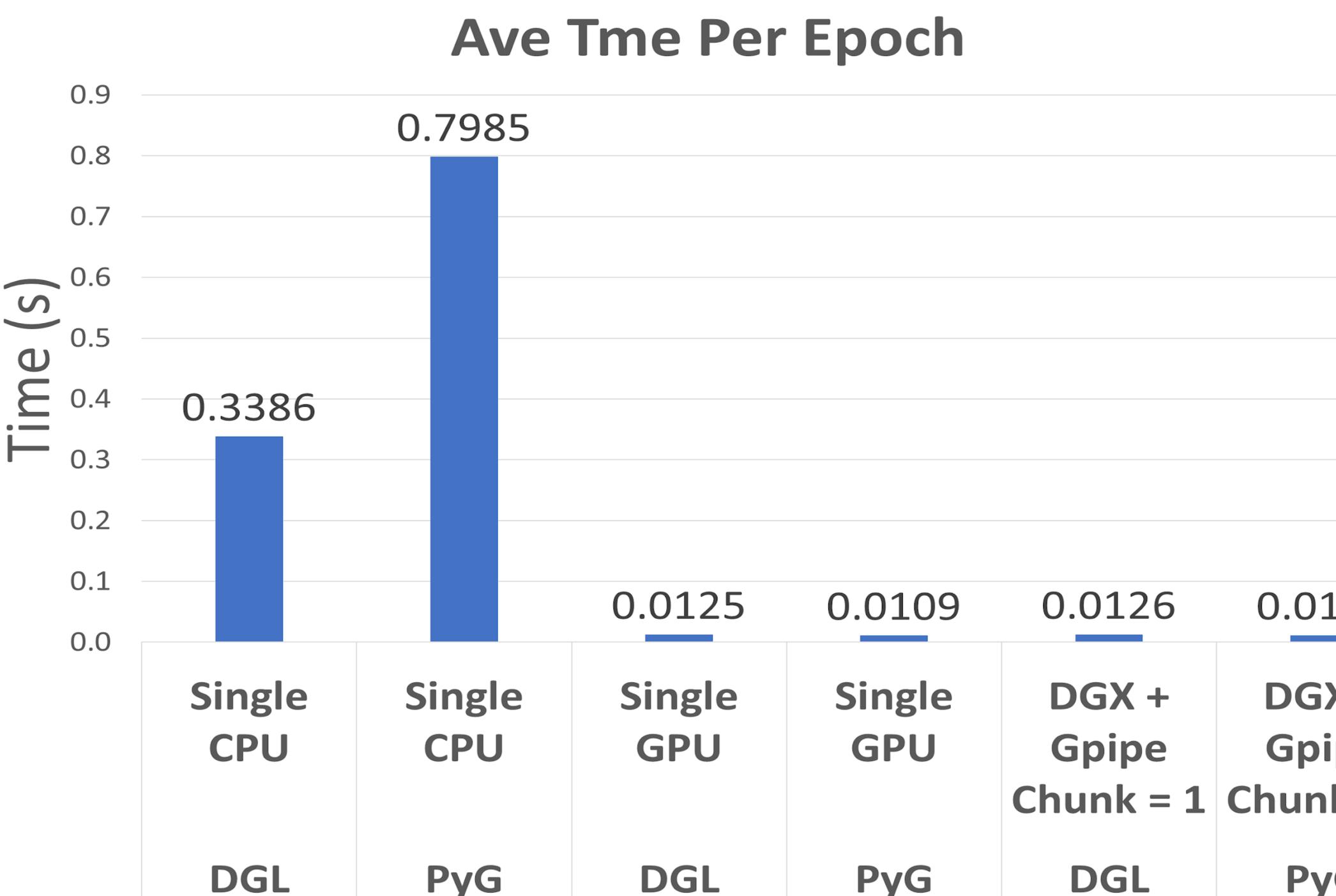
- GNNs build a neural network atop a simultaneous message passing paradigm that aggregates neighborhood information.
- Each NN layer is a function of the previous layer *and* the input graph.

Graph Attention Network (GAT)

- GAT (Veličković et al. 2017) is a GNN with *attention layers* on top of graph convolutions.
- The entire network consists of two GAT layers with a multi-head attention layer to compute attention coefficients for the neighboring nodes, which is used to propagate information through the network.



Benchmark Training: Single Devices

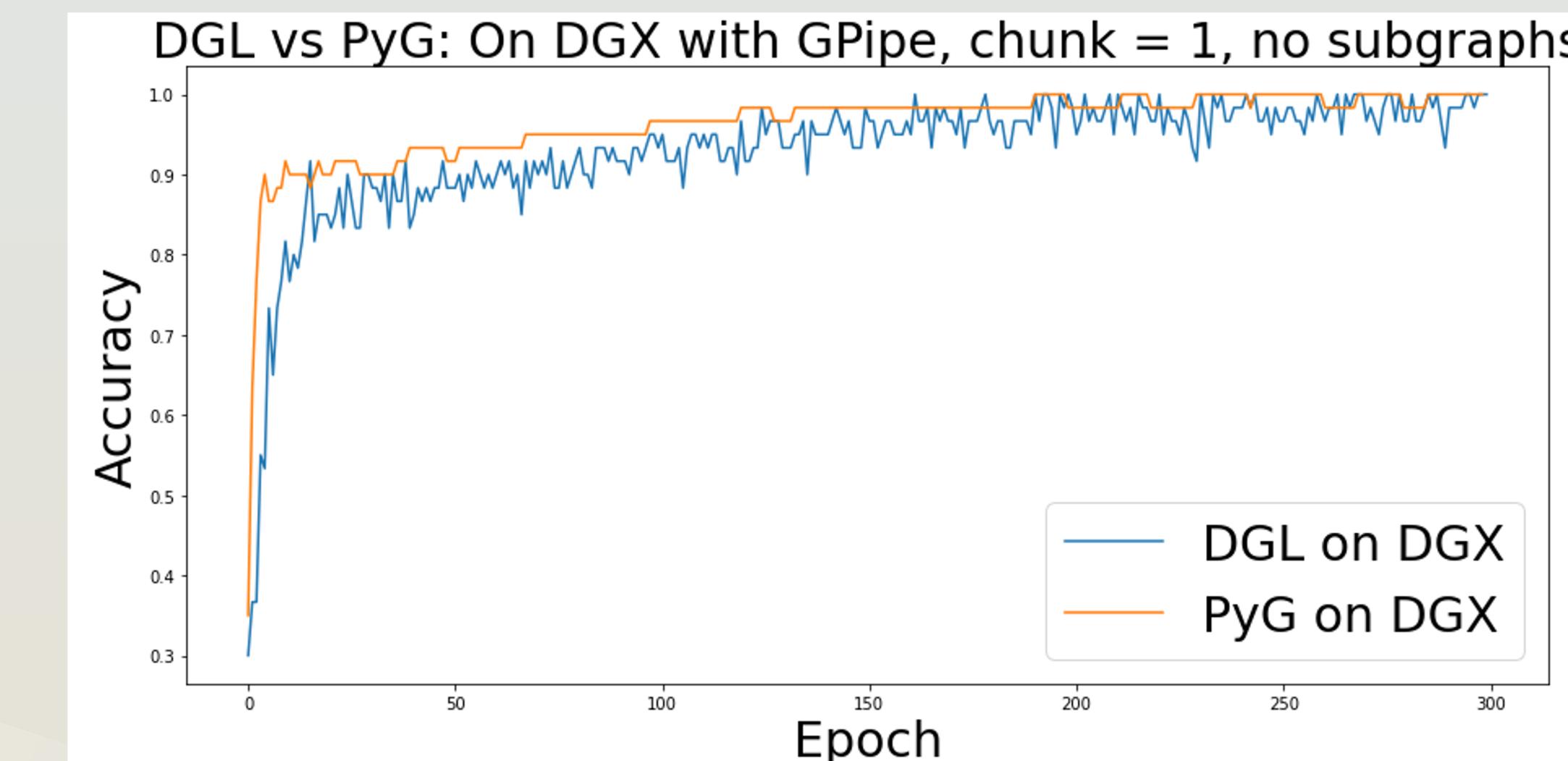


As expected, single GPU performs much faster than single CPU for both PyG and DGL frameworks.

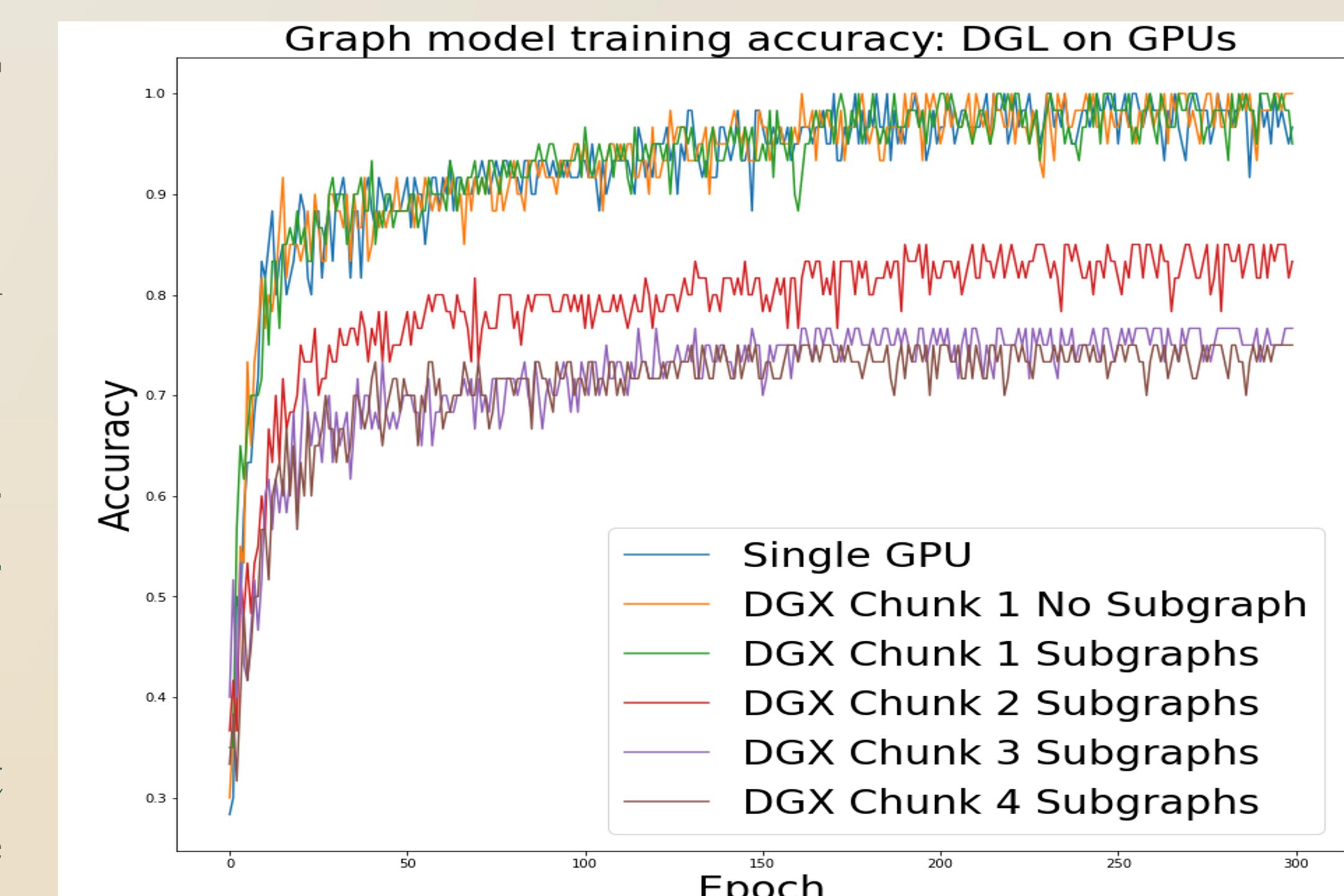
No significant performance improvement using DGX (4 GPUs) and GPipe with “chunk size” = 1 (i.e., no micro-batching) compared to single GPU.

Refactoring the GNN into a Sequential container that only passes a single argument required extracting the graph nodes and embedding features into a tuple of tensors so that the graph could be re-built in the convolutional layers.

Results: Accuracy degraded



DGL and PyG implemented on DGX with GPipe without micro-batching offer similar performance.



GPipe micro-batches the graph nodes *sequentially*, from which the subgraphs are generated. Significant information loss resulted from ignored edge relationships across the batches.

DGL and PyG provide sophisticated single graph micro-batching solutions, but GPipe would not use these batches.

Future work

- Implement other GNN layer models, such as GraphConv, with these frameworks.
- PubMed is too small. Experiment with massive datasets: *Can performance over single GPU increase on extremely large graphs?*

Acknowledgements

Thank you to Prof. Ian Foster, Prof. Rick Stevens, and Peng Ding for the guidance and feedback on this paper. We also thank the DGX team, Daniel Murphy-Olson and Ryan Aydelott, and the Computing, Environment, and Life Sciences directorate at Argonne National Laboratory.

Argonne
NATIONAL LABORATORY