EVA: Energy-Efficient Virtual Node Accelerator for Graph Neural Networks

Dhruva Digesh Barfiwala

Parima Devanshu Mehta



Motivation

Graph Neural Networks

- Handle non-Euclidean data
- Existing Models: GIN, GCN, GAT, GraphSAGE, etc.





Social Networks



Existing GNNs fail to distinguish simple graph structures



Existing GNN Inference Accelerators (e.g., GRIP, BoostGCN) are not suitable for real-time inference

- Critical tasks in chemistry, medicine, etc. demand high accuracy and low <u>inference</u> latency
- Real-time inference of GNNs augmented with Virtual Nodes remains relatively unexplored





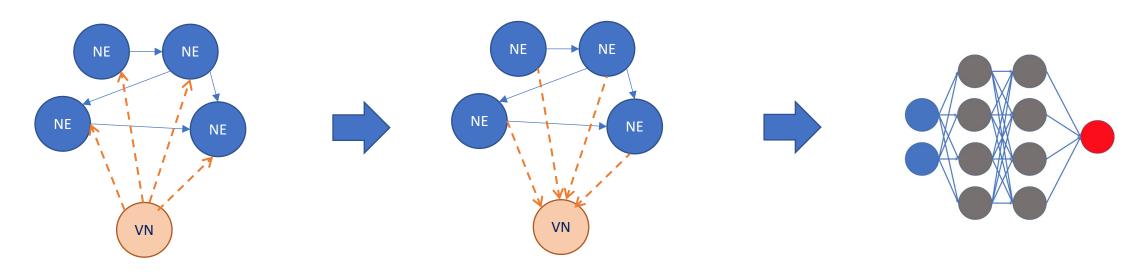
Molecules

Virtual Node Augmentation

- Virtual node is a dummy node connected to all graph nodes
- Enhances the representation power of GNNs by providing a shortcut during message passing

Table 15: Results for ogbg-molhiv.

Method	Additional	Virtual	ROC-AUC (%)		
	Features	Node	Training	Validation	Test
GCN	X	1	88.65±1.01	83.73±0.78	74.18 ±1.22
		×	$88.65 {\scriptstyle\pm2.19}$	82.04 ± 1.41	$76.06 {\pm} 0.97$
			$90.07{\scriptstyle\pm4.69}$	$83.84 {\pm} 0.91$	$75.99 {\scriptstyle\pm1.19}$
GIN	X	/	93.89±2.96	84.1 ±1.05	75.2 ±1.30
		×	$88.64 {\pm} 2.54$	$82.32{\scriptstyle\pm0.90}$	$75.58 {\pm} 1.40$
	✓	/	92.73 ±3.80	84.79 ±0.68	77.07 ±1.49



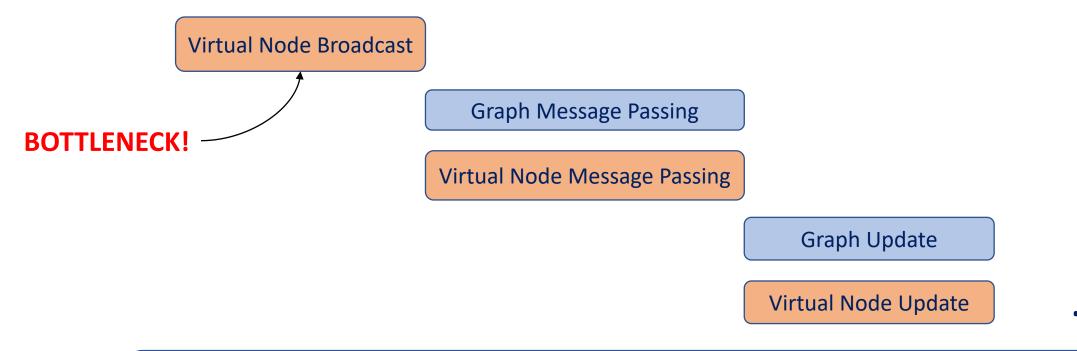
Virtual Node Broadcast

Virtual Node Message Passing

Virtual Node Update

Challenges with Virtual Node

- Virtual node broadcast falls on the critical path for graph message passing
- Latency and resource intensive involves all nodes of the graph
- One Layer of GNN augmented with Virtual Node:



GOAL

Energy efficient virtual node augmentation to Graph Isomorphism Network (GIN) with minimal latency and area overhead

Setup

Model Training

- Trained <u>GIN-Virtual Model</u> implemented in PyTorch
- Model specifications:
 - # Layers: 5
 - Emb Dim: 100
- Dataset: <u>ogbg-molhiv</u>
- Saved pretrained weights

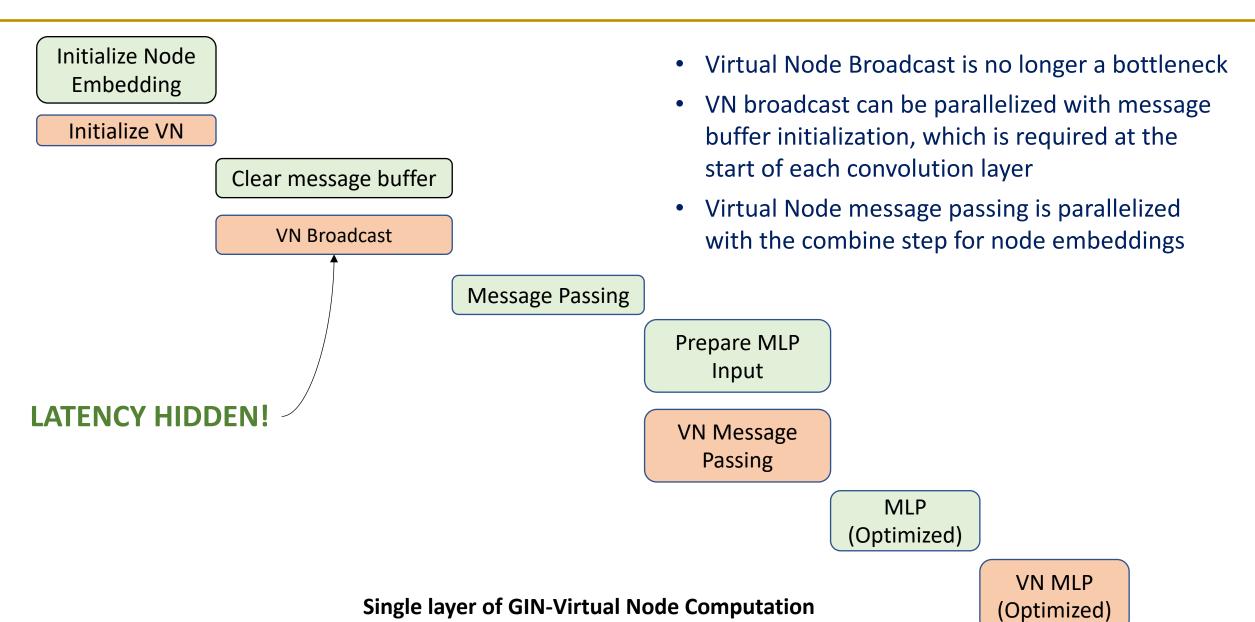
Collect Model Weights and Golden Output

- Employed weight fusion to reduce HW complexity
- Organized the fused weights into binaries
- Collected output from PyTorch implementation for Golden C verification

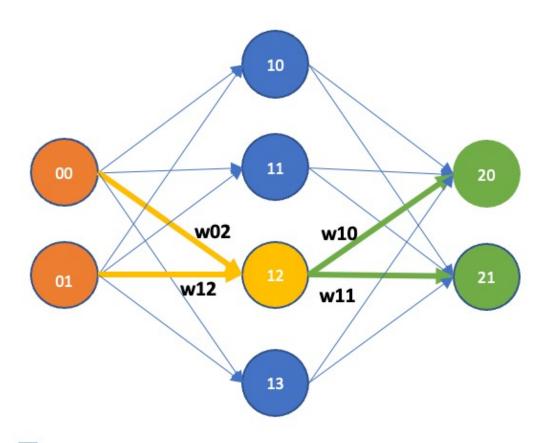
Golden C Implementation

- Translated PyTorch GIN-Virtual model to Golden C
- Automated C header generation
- Verified functional correctness against PyTorch implementation
- MSE of **10**⁻⁴ was observed

Base HLS with Virtual Node



Base HLS: VN MLP Optimization



- MLP is <u>pipelined</u> across the hidden MLP layer dimensions
- MLP input is <u>partitioned</u> to ensure each hidden layer dimension can be computed in a single cycle
- MLP output is <u>partitioned</u> to allow broadcast of hidden layer in one cycle

- One dimension of hidden layer is computed in two cycles (one for ReLU)
- One dimension of hidden layer is broadcast in one cycle

Further Optimization Opportunities



Base HLS Shortcomings:

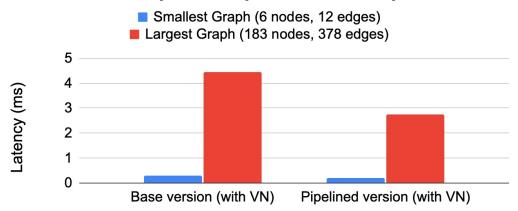
- No pipelining
- Not scalable for large graphs
- Needs aggressive partitioning to further reduce latency



Pipelining to the rescue:

- Inter-layer pipelining
- Suitable for graphs of all sizes
- Requirements:
 - Degree and neighbor table
 - Double buffering

Base vs. Pipelined Optimization Comparison



Optimization type

^{*}Note: Pipelined optimization latency is obtained by averaging min-max latencies.

*Graphs from ogbg-molhiv dataset

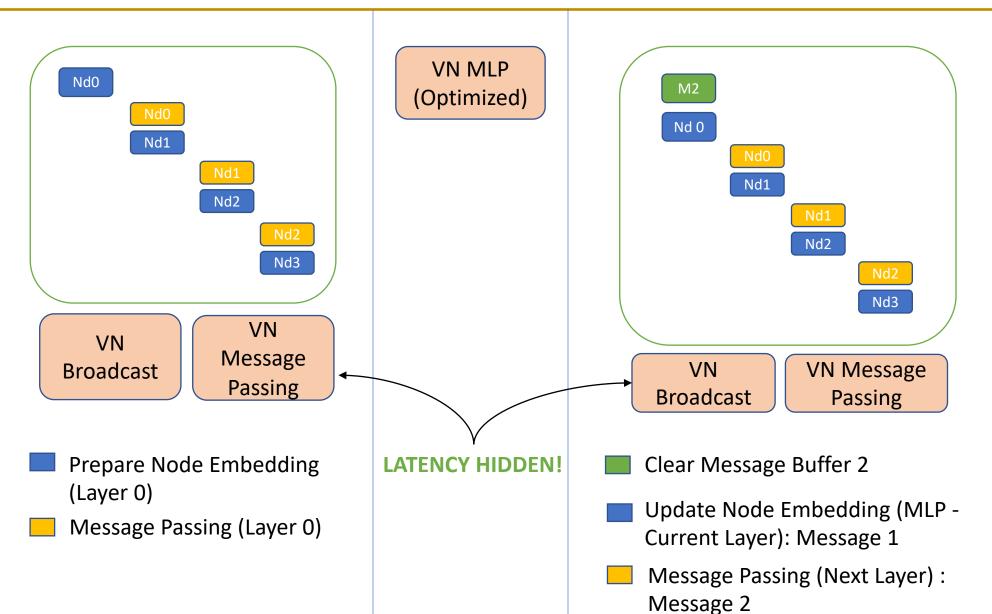
Pipelined HLS with Virtual Node

Clear Message 1

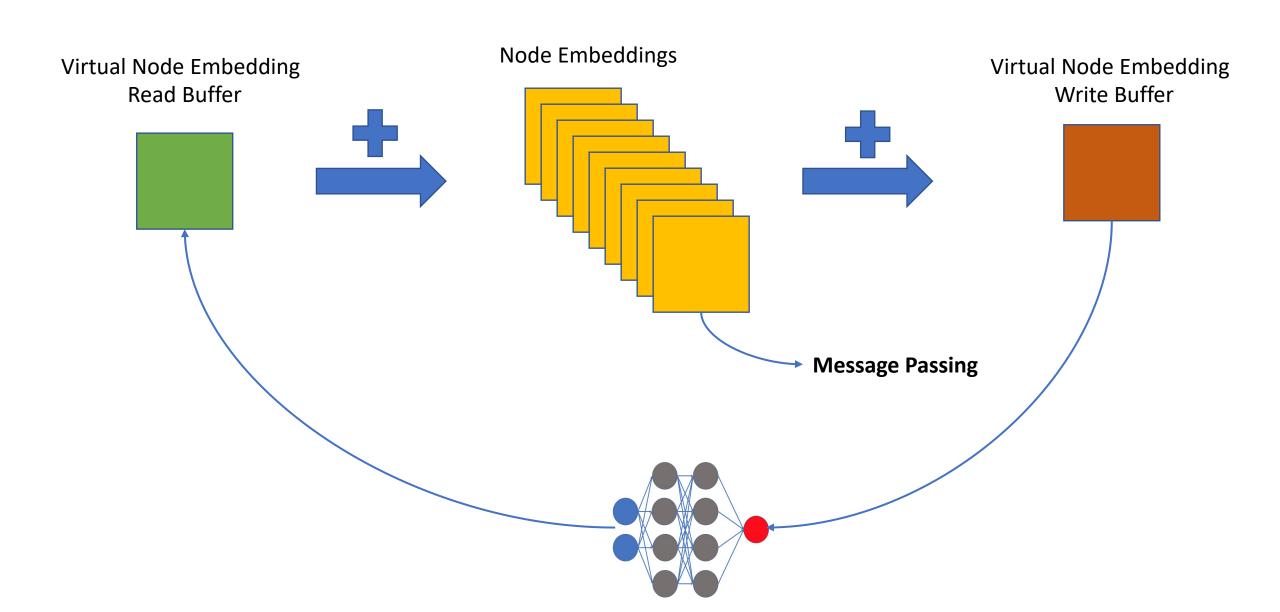
Clear Message 2

Prepare
Degree/Neighbor
Table

Initialize VN



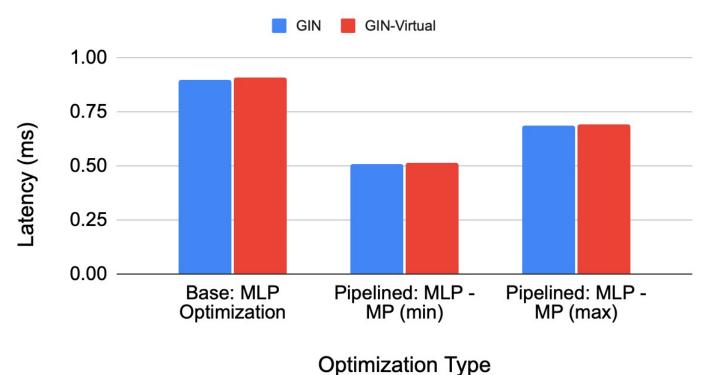
Read-Write Buffers for Virtual Node



Latency Overhead Analysis

Board	Dataset	Graph	
Xilinx Alveo U280	ogbg-molhiv	19 nodes, 40 edges	

Performance Overhead of GIN + Virtual Node



Latency Overhead of ~10 us due to Virtual Node MLP

Test rocauc (1 epoch):

GIN: 66.4

GIN-Virtual: 68.05

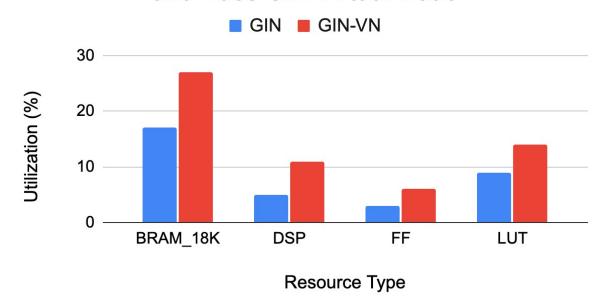
Area Overhead Analysis

Graph from ogbg-molhiv dataset: 19 nodes, 40 edges

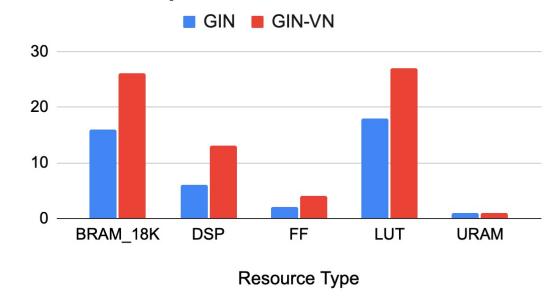
Total Available Resource on Xilinx Alveo U280

BRAM_18K	DSP	FF	LUT	URAM
4032 9024		2607360	1303680	960

Resource Utilization Comparison of Base GIN and Base GIN-Virtual Node



Resource Utilization Comparison of Pipelined GIN and Pipelined GIN-Virtual Node



Utilization (%)

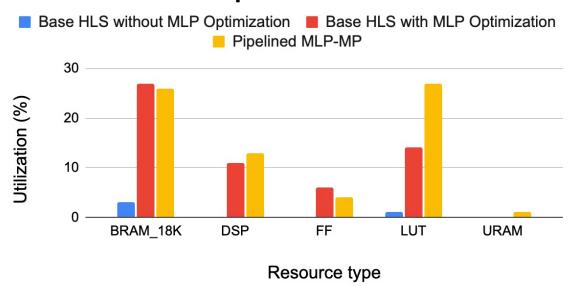
VN: Speed up and Area Analysis

- Comparison against Base HLS (without MLP Optimization) augmented with Virtual Node
- Graph from ogbg-molhiv dataset: 19 nodes, 40 edges

Optimization	None	MLP	Pipelined MLP- MP
Latency (ms)	50.116	0.907	0.6043
Speedup	1	55.25x	82.93x

*Note: Pipelined optimization latency is obtained by averaging min-max latencies.

Area analysis w.r.t Base HLS without MLP Optimization



Note: Net Compute Time and Resource Utilization on U280 estimates obtained from Vitis HLS 2021.1

Future Work

Multiple VN

Increase the number of virtual nodes and analyze the impact on accuracy, performance, and resource utilization

Intra-layer Pipelining Explore pipelining between message passing and MLP within the same layer by processing the message for a given destination

Large Graphs

Pipeline loading graphs from DRAM with message passing and explore data streaming

Thank You!

Acknowledgement: We thank Prof. Callie Hao and the Sharc Lab at Georgia Tech for providing us with a GIN model to augment Virtual Node for our project.