



#### A Data-Driven Approach to Dataflow-Aware Online Scheduling for Graph Neural Network Inferences

Pol Puidgemont<sup>1</sup>, Enrico Russo<sup>2</sup>, Axel Wassington<sup>1</sup>, Abhijit Das<sup>1</sup>, Sergi Abadal<sup>1</sup>, Maurizio Palesi<sup>2</sup>

<sup>1</sup> Universitat Politècnica de Catalunya, Spain <sup>2</sup> University of Catania, Italy





#### Overview



**ASP-DAC 2025** 

Efficiently running Graph Neural Networks inferences in online, real-time, or streaming contexts requires optimal mapping and dataflow configurations, which must adapt dynamically to the changing input graph.



#### Overview



Efficiently running Graph Neural Networks inferences in online, real-time, or streaming contexts requires optimal mapping and dataflow configurations, which must adapt dynamically to the changing input graph.



and edges (relationships between points).



From left to right. Citronella molecule. Adiacency matrix representation. Graph representation. From https://distill.pub/2021/gnn-intro/

**ASP-DAC 2025** 

#### Graph Neural Networks (GNNs) are a specialized type of neural network designed to process and analyze graph-structured data, which consists of nodes (data points)



logistics, smart energy grids, cybersecurity, drug discovery, etc.





**ASP-DAC 2025** 

#### GNNs are increasingly employed in a wide range of applications: social network analysis, recommender systems, fraud detection, transport networks and





A Graph Neural Network (GNN) is a parametrized transformation of graph attributes designed to preserve permutation invariance and graph symmetries. In our case, the input graph is represented by an adjacency matrix and real-valued feature arrays for the nodes, as illustrated below:



Node feature array

1	0	1	0	1
0	1	1	0	1
1	1	1	0	0
0	0	0	1	1
1	1	0	1	1

Α **Adjacency Matrix** 



# Background

#### **Graph Neural Networks**

Networks) consists in a two-phase procedure applied node-wise.



# A commonly adopted transformation in each layer (e.g. in Graph Convolutional

# Background

#### **Graph Neural Networks**

Networks) consists in a two-phase procedure applied node-wise.



# A commonly adopted transformation in each layer (e.g. in Graph Convolutional



multiplication:



#### The computation for many GNN models can be expressed as a dense-sparse matrix



# Background **Spatial Accelerators for GNNs**



We focus on large on programmable spatial accelerators with high parallelism Opportunities. These spatial accelerators can efficiently execute both the SpMM and Dense GEMM kernels of the two GNN layers phases.





# Background **Spatial Accelerators for GNNs**



Each accelerator feature a global buffer and a set of PEs interconnected with a Network-on-Chip. **Different unrolling dimensions and dataflow capabilities can be supported.** The PEs could be splitted between the two phases or be capable of runing computations for both phases.





#### Background Heterogenous Dataflow Multi-Accelerator Systems



Flexibility in supported dataflows and interconnects can come with additional overhead. A multi-dataflow accelerator can also be obtianed combining multiple accelerators with limited flexibility in a single system, i.e. a heterogenous dataflow multi-accelerator system, similarly, to what has been proposed for traditional DNNs by Kwon et al.

Kwon et al. "Heterogeneous dataflow accelerators for multi-DNN workloads." HPCA 2021



#### Background **GNN Mapping**

#### a) Sequential Loop Nest

# Aggregation for v in range(0, V, T\_Va): for f in range(0, F, T\_Fa): parallel-for v1 in range(0, T\_Va, 1): for n in range(A\_node[v+v1], A\_node[v+v1+1]): parallel-for f1 in range(0, T\_Fa, 1): parallel-for n1 in range(0, T\_N, 1): M[v+v1, f+f1] += H[ A\_edge[n+n1], f+f1] # Combination for v in range(0, V, T\_Fc): for g in range(0, G, T\_Gc): for f in range(0, F, T\_Fc): parallel-for v1 in range(0, T\_Vc): parallel-for g1 in range(0, T\_G): parallel-for f1 in range(0, T\_Fc): H\_new[v+v1, g+g1] += M[v+v1, f+f1] \* W[f+f1, g+g1] # Legend H: input node features  $\mathbf{H}^{(k-1)}$  H\_new: output node features  $\mathbf{H}^{(k)}$ M: intermediate matrix A\_node: adj. matrix CSR node array A\_edge: adj. matrix CSR edge array W: weights W

#### Similarly to traditional DNNs, the algorithm for the two phases can be represented as two loop nests executed sequentially.





# **Background**GNN Mapping



However, to optimise data movement, the computation of the two phases can be interleaved or overlapped. Hence, three possible **inter-phase dataflows** can be selected: sequential, sequential pipeline and parallel pipeline. Each requiring specific tiling strategies.





# **Background** GNN Mapping



Furthermore, for both phases the **tiling sizes** and **unrolling dimensions** have to be chosen as in traditional DNNs mappings, but the N dimension, i.e. the number of neighbours, varies for each node and for each graph.





# Problem **DNN vs GNN Mapping**

#### DNN

A DNN layer can be **mapped offline** considering its shape and architecture. The optimal mapping does not depend on the specific layer input, e.g. image.

DNN mapping usually involves only dense tensors.

**Inter-layer** optimization can be applied.

#### GNN

A GNN layer computation also depends on contingent input graph topology, hence, the optimal mapping must be found for each input instance.

GNN mapping has to consider **sparsity** of the adjacency matrix.

**Inter-phase** dataflow can be optimized.







#### Problem

# The best mapping of a GNN has to be found for each input graph instance and there is no *one-fits-all* solution

	GNNBenchmark.Pattern -	0.00	0.00	0.05	0.00	0.62	0.00	0.32	0.02	0.00	0.00	0.00	0.00	0.00	0.00
	GNNBenchmark.TSP -	0.00	0.00	0.62	0.00	0.33	0.00	0.00	0.03	0.00	0.00	0.00	0.00	0.00	0.00
	Graphlaxy.Large –	0.28	0.00	0.26	0.02	0.19	0.00	0.20	0.00	0.00	0.00	0.00	0.00	0.00	0.00
ISEL	Graphlaxy.Medium –	0.20	0.02	0.13	0.00	0.19	0.15	0.20	0.04	0.02	0.00	0.00	0.00	0.00	0.00
	Graphlaxy.Small –	0.25	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.02	0.00	0.52	0.00	0.05
Date	QOPTLib –	0.44	0.00	0.02	0.00	0.15	0.00	0.22	0.07	0.00	0.00	0.00	0.05	0.00	0.02
	SuiteSparse.Mix –	0.37	0.00	0.16	0.00	0.23	0.00	0.00	0.21	0.00	0.02	0.00	0.00	0.00	0.00
	TUDataset.ENZYMES –	0.02	0.00	0.00	0.00	0.00	0.00	0.06	0.13	0.00	0.21	0.12	0.08	0.07	0.07
	TUDataset.PROTEINS -	0.05	0.00	0.02	0.00	0.10	0.00	0.00	0.14	0.00	0.33	0.20	0.06	0.00	0.00
Т	UDataset.REDDIT-BINARY –	0.00	0.00	0.18	0.00	0.60	0.00	0.00	0.04	0.00	0.12	0.02	0.00	0.00	0.00
		(49) 1	(89') 1	(199) (2)	(99) 1	(PP) 	(PP, 0)	(89'N)	(Sed)	sed.	(sed)	(sed)	ા હિ <sup>ર</sup> ે	۔ جھر چھر	ા હર્જ

Optimal (Inter-phase Dataflow, Tiling Scheme) configuration





#### Problem

Many GNNs application scenarios involve real-time or streaming GNN inferences such as fraud detection, malicious user detection in social networks, load balancing in energy grids, smart cities, etc.



Mapping and scheduling of GNN inferences requires online optimization.

#### Problem

#### How to find the best GNN mapping/dataflow on-the-fly?



Thus, each GNN request can translate to multiple **GNN layer inference requests** that are stored in a ready queue of jobs ready to be run on the serving system.

The relevant information to choose the optimal dataflow is the graph structure (num. of nodes, num. of edges, density, degree distribution, etc.) and the layer shape (number of output and input features).



#### Each GNN request is associated to a GNN model, an input graph and embeddings.

- GNN Layer Weights
- **GNN Layer Shape**
- Input Graph
- Input Embeddings



of latency for an incoming inference request.



# We investigated using a trained ML model to predict the best mapping in terms



dataflow and tiling strategy).

Trained ML Models for each mapping config.



**ASP-DAC 2025** 

#### We considered a latency prediction model each mapping configuration (inter-phase



latency for better online scheduling algorithms.

Trained ML Models for each mapping config.



**ASP-DAC 2025** 

# The ensemble model allows us to find both ranking of suboptimal mappings and the



#### In particular, we considered all the 3 inter-phase dataflows and the following 8 tiling configurations totaling 24 possible mapping configurations.

Tiling		Aggregation	Combination			
ω	T_Va	T_Fa	T_N	T_Vc	T_Gc	T_Fc
a	*	$\min(\text{Num. PEs}, F)$	1	*	1	min(Num. PEs,F)
b	*	$\min(2,F)$	$\lfloor F/2 \rfloor$	*	1	$\min(2,F)$
с	*	$\min(8,F)$	$\lfloor F/2 \rfloor$	*	1	$\min(8,F)$
d	*	1	1	*	1	1
e	*	$\min(18,F)$	$\lfloor F/2 \rfloor$	*	1	$\min(18,F)$
f	*	1	$\min(18,V)$	*	1	1
g	*	$\min(18,F)$	1	*	1	$\min(85,F)$
h	*	1	$\min(18,V)$	*	1	$\min(85,F)$



# **Proposed solution Feature Extraction**



We considered input features of GNN inference requests for latency prediction models.

Simple features represent basic GNN layer characteristics, such as input/output features and nodes, while composite features are latency estimation formulas based on these attributes.

**ASP-DAC 2025** 

		Description	S	ymbol and/or Equation
		Number of nodes and edges		V and E
1		Spatial tiling factors	T_Va	a,T_Fa,T_N,T_Vc,T_Gc,T_Fc
		Graph density		$E \times V^{-2}$
	base	Clustering coefficient, measuring the proba- bility that the adjacent nodes of a node are connected		[24]
		Seven node degrees quantiles normalized with respect to maximum degree including minimum degree		$Q_i,  \forall i \in \{1,, 7\}$
		Number of operations, assuming a dense matrix mul. for the aggregation phase	<i>S</i> <sub>1</sub>	$V \times F \times (G + \text{Mean Degree})$
		Estimation of the number of cycles for the combination phase	$S_2$	$\frac{V \times F \times G}{T_V c \times T_F c \times T_G}$
	eatures	Estimation of the number of cycles for the aggregation phase (dense mat. mul.)	<i>S</i> <sub>3</sub>	$\frac{V \times F \times \text{Mean Degree}}{T_N \times T_Fa}$
	base+f	Estimation of the latency for sequential inter-phase dataflow	$S_4$	$S_1 + S_3$
		Cycles estimation for aggregation phase assuming CSR encoding	$S_5$	$\sum_{v}^{\mathcal{W}} \frac{N_{v} \times F}{\min(N_{v}, T_N) \times T_F}$
	·	Cycles estimation for sequential inter-phase dataflow assuming CSR encoding	<i>S</i> <sub>6</sub>	$S_3 + \frac{S_5}{T_Va}$

25



# **Proposed solution ML Models**



#### The latency regressors were implemented as gradient boosting trees using LightGBM framework



# **Proposed solution ML Models**



**ASP-DAC 2025** 

#### For better generalization capability we train the models to predict the logarithm of the layency.



# **Proposed solution** Training

#### The latency targets were obtained using **STONNE-Omega Simulator**

Garg, Raveesh, et al. "Understanding the design-space of sparse/dense multiphase GNN dataflows on spatial accelerators." 2022 IEEE International Parallel and Distributed Processing Symposium

#### The training was performed on datasets of synthetic graphs generated using Graphlaxy

A. Wassington, S. Abadal, "Bias reduction via cooperative bargaining in synthetic graph dataset generation.", 2022





#### **Experiments** Prediction Accuracy in Offline Single Accelerator Setup

Configurable Inter-Phase Dataflow Configurable Unrolling Dimensions and Tiling Sizes



We evaluated the total execution time for different graph datasets on highly configurable single accelerator system. The optimal total execution time is achieved when each request is executed with the best dataflow configuration.



#### Experiments **Prediction Accuracy**

Dataset Name	MAPE↓	Top-1   acc. (%) ↑	Top-3 acc. (%) ↑	Improve   random (%) ↑	ment over   best fixed (%) ↑	Degradation over optimal (%)↓
Graphlaxy.Medium	3.78	91.28	99.62	93.63	58.42	0.56
Graphlaxy.Large	4.83	83.63	98.20	96.27	8.99	0.86
Graphlaxy.Small	17.36	46.94	75.51	97.93	0.49	5.12
SparseSuite.Mix	13.32	62.79	81.40	89.79	8.04	5.24
QOPTLib	24.99	56.41	74.36	82.45	20.98	14.94
PATTERN	40.17	33.13	68.56	84.04	-6.30	9.40
TSP	29.13	62.18	92.85	87.38	2.18	4.22
ENZYMES	20.23	19.52	40.41	91.95	-3.29	12.10
PROTEINS	15.82	30.45	56.76	91.49	1.24	11.56
REDDIT-BINARY	14.16	35.08	76.45	85.96	0.74	8.68

Less than 15% degradation in execution time over optimal with the predicted dataflows.



#### Experiments

#### Ablation study on composite features and log prediction

Dataset Name	Model Name	MAPE ↓	Degradation over optimal (%) $\downarrow$
	base	35.77	48.12
Crophlaw Madium	base+features	10.76	4.06
Grapmaxy.Medium	base+log	4.95	0.84
	base+features+log	3.78	0.56
	base	97.04	92.68
Sporco Suito Mix	base+features	31.68	12.52
SparseSuite.Mix	base+log	19.09	5.87
	base+features+log	13.32	5.24
	base	515.94	92.94
OODTI ;h	base+features	152.67	86.60
QUEILID	base+log	38.22	21.18
	base+features+log	24.99	14.94

**ASP-DAC 2025** 

Considering **composite** features and predicting **logarithm of the latency** achieves the highest accuracy.



## Experiments **Online Scheduling Setup: Multi-Accelerator System**



#### Augmented Ready Queue



# Experiments



**ASP-DAC 2025** 

#### **Online Scheduling Setup: Multi-Accelerator System**

#### Augmented Ready Queue



### Experiments **Online Scheduling Setup: Multi-Accelerator System**



**ASP-DAC 2025** 

#### Augmented Ready Queue



# Experiments

#### Heterogenous Dataflow Multi-Accelerator System

For the online scheduling experiments, we considered a system consisting in **3 sub-accelerators**.

Each sub-accelerator (SA) support a different inter-phase dataflow (Sequential, Seq. Pipeline, Parallel Pipeline), but it is flexible in terms of loop unrolling dimensions.

Furthermore each SA features **512 processing elements** with 64B local buffer each.





#### Experiments **Considered Online Scheduling Algorithm**

We considered a shortest-job-next algorithm which minimizes job waiting time. We assumed a work-conserving scenario in which once a SA is free, the GNN layer in the ready queue with shortest predicted execution time is chosen to be executed next.



#### Experiments

#### **Online scheduling performance**



baselines using random tiling selection.

The proposed approach, i.e. a shortest-job-first based on the latency predictions, allows 83.88% and 99.95% reduction in execution time and turnaround time respectively, with respect to the best performing feasible scheduling algorithm, namely shortest-job-first based on number of nodes

**ASP-DAC 2025** 



✦ := Not feasible in practice



#### Experiments

#### **Online scheduling performance**



Algorithm performance for online scheduling policies with baselines using random tiling selection.

While, with respect to the best considered non-feasible scheduling algorithm, time, a shortest-job-first based on actual execution times, the proposed solution achieves 1.07x and 1.60x higher mean execution time and mean turnaround time, respectively.



✦ := Not feasible in practice



#### Experiments **Online scheduling overhead**



**ASP-DAC 2025** 

In our setup, compiled gradient boosting tree predictors run on the CPU parallel to inference executions. The mean latency prediction phase duration for each graph was only **12.3% of the mean job waiting time**, avoiding any turnaround time increase.

Notably, the latency prediction phase is independent of GNN request sizes, suggesting that overhead could be further reduced for larger graphs.



# Conclusion

- We presented a data-driven approach to dataflow latency evaluation of GNN workloads, based on gradient boosting trees.
- We showed the usefulness of such predictors in an online scheduling scenario featuring multi-dataflow GNN accelerators.
- Several limitations remain to be addressed in future works. For instance, only a subset of possible mappings has been considered, and the proposed methodology requires to train a model for each tiling/dataflow configuration.
- Future work could also focus on more complex online scheduling algorithms informed by the predictions, customized for the specific hardware system.





enrico.russo@phd.unict.it

**ASP-DAC 2025** 

# Thanks

