Tensor graph superoptimization for graph neural networks

Project Proposal by Antonia Boca

Motivation

• Tensat [1] uses *equality saturation* to optimize computation graphs



- It is built on top of egg [2] and TASO [3]
- It performs equality saturation for a bunch of operators...
- ... but not for GNN operators!

Graph Neural Networks

- Most popular GNNs use the `message-passing' framework
- Each architecture has an associated aggregation operator



The basic GNN operator

• A "layer" in a GNN is usually an iteration of the aggregation & update operators

$$\begin{aligned} \mathbf{h}_{u}^{(k+1)} &= \mathrm{UPDATE}^{(k)}(\mathbf{h}_{u}^{(k+1)}, \mathrm{AGGREGATE}^{(k)}(\{\mathbf{h}_{v}^{(k+1)}, \forall v \in \mathcal{N}(u)\})) \ (1) \\ &= \mathrm{UPDATE}^{(k)}(\mathbf{h}_{u}^{(k+1)}, \mathbf{m}_{\mathcal{N}(u)}^{(k)}) \end{aligned}$$
(2)

• GCN Convolution [4]:

$$\mathbf{h}_{u}^{(k)} = \sigma(\mathbf{W}^{(k)} \sum_{v \in \mathcal{N}(u) \cup \{u\}} rac{\mathbf{h}_{v}}{\sqrt{|\mathcal{N}(u)||\mathcal{N}(v)|}})$$

Operator	Description Inputs		Inputs	Type signature			
ewadd ewmul matmul	Element-wise addition Element-wise multiplication Matrix multiplication	in in activatio	$put_1, input_2$ $put_1, input_2$ on, input_1, input_2	$\begin{array}{c} (T,T) \rightarrow T \\ (T,T) \rightarrow T \\ (N,T,T) \rightarrow T \end{array}$			
conv ^a relu tanh	Grouped convolution Relu activation Tanh activation	stride _h , stride _w	, pad., act., input, weight input input	$\begin{array}{c} (\mathrm{N},\mathrm{N},\mathrm{N},\mathrm{N},\mathrm{T},\mathrm{T}) \to \mathrm{T} \\ \mathrm{T} \to \mathrm{T} \\ \mathrm{T} \to \mathrm{T} \end{array}$			
sigmoid poolmax poolavg transpose ^b	Sigmoid activation Max pooling Average pooling Transpose	input, kernel _{{h} , input, kernel _{{h} , input, input	input w_{j} , stride $\{h,w\}$, pad., act. w_{j} , stride $\{h,w\}$, pad., act. t permutation	$T \rightarrow T$ $(T, N, N, N, N, N, N) \rightarrow T$ $(T, N, N, N, N, N, N) \rightarrow T$ $(T, S) \rightarrow T$			
enlarge c concat _n d split e	Pad a convolution kernel with zeros Concatenate Split a tensor into two Get the first output from split	ing axis, in	point, ref-input $put_1, \ldots, input_n$ axis, input input	Graph Neural Network Operators			
split ₁	Get the second output from split		input	Name	SparseTensor	edge_weight	edge_attr
merge ^f reshape ^g input	Update weight to merge grouped conv Reshape tensor Input tensor	w in i	weight, count input, shape identifier ^h	GCNConv (Paper)	\checkmark	\checkmark	
weight noop ⁱ	Weight tensor Combine the outputs of the graph	identifier ^h input ₁ , input ₂		ChebConv (Paper)		\checkmark	
				SAGEConv (Paper)	~		
				GraphConv (Paper)	\checkmark	\checkmark	
				GatedGraphConv (Paper)	√	\checkmark	
Decide on a subs of GNN operators			rs to	ResGatedGraphConv (Paper)	\checkmark		
	optimize			GATConv (Paper)	1		\checkmark
				GATv2Conv (Paper)	~		~
				TransformerConv (Paper)	\checkmark		\checkmark

5

bipartite

 \checkmark

 \checkmark

 \checkmark

 \checkmark

 \checkmark

 \checkmark

static

 \checkmark

 \checkmark

 \checkmark

 \checkmark

 \checkmark

 \checkmark

lazy

 \checkmark

 \checkmark

 \checkmark

 \checkmark

 \checkmark

Project idea

- Implement equality saturation for a subset of GNN operators using the egg library
 - Extend Tensat with these new operators
 - Evaluate extended Tensat
- Experiment settings:
 - Run extended Tensat on GNN architectures (e.g. GCN, GAT, GraphSAGE) and benchmark datasets (e.g. Reddit and Citation networks)
- Evaluate architectures optimized by extended Tensat:
 - Against architectures optimized by normal Tensat
 - Against architectures optimized by FlexFlow [5] (built on top of Unity)

Workplan

- Preparation
 - Understand how GCNConv, GATConv, SAGEConv are implemented in pytorch-geometric
 - Learn (basic) Rust
 - Understand the implementation of Tensat operator optimisations
- Come up with optimisations
 - Implement optimisations using egg
 - Extend the Tensat framework
- Run experiments on GNN architectures and datasets
 - NVIDIA core; plan to run this on AWS

Questions and suggestions

References

- [1] Max Willsey, Chandrakana Nandi, Yisu Remy Wang, Oliver Flatt, Zachary Tatlock, and Pavel Panchekha. 2021. Egg: Fast and extensible equality saturation. Proc. ACM Program. Lang. 5, POPL, Article 23 (January 2021), 29 pages.
- [2] Yang, Yichen, et al. Equality saturation for tensor graph superoptimization. Proceedings of Machine Learning and Systems 3 (2021): 255-268.
- [3] Z. Jia, O. Padon, J. Thomas, T. Warszawski, M. Zaharia, A. Aiken. TASO: Optimizing Deep Learning Computation with Automated Generation of Graph Substitutions. *SOSP*, 2019.
- [4] Kipf, Thomas N., and Max Welling. "Semi-supervised classification with graph convolutional networks." *arXiv preprint arXiv:1609.02907 (2016).*
- [5] Colin Unger, Zhihao Jia, Wei Wu, Sina Lin, Mandeep Baines, Carlos Efrain Quintero Narvaez, Vinay Ramakrishnaiah, Nirmal Prajapati, Pat McCormick, Jamaludin Mohd-Yusof, Xi Luo, Dheevatsa Mudigere, Jongsoo Park, Misha Smelvanskiy, and Alex Aiken. Unity: Accelerating DNN Training Through Joint Optimization of Algebraic Transformations and Parallelization. In Proceedings of the Symposium on Operating Systems Design and Implementation (OSDI), July 2022.