

The Future of Machine Learning is Sparse

Future is Sparse @ SC'23

11/17/23

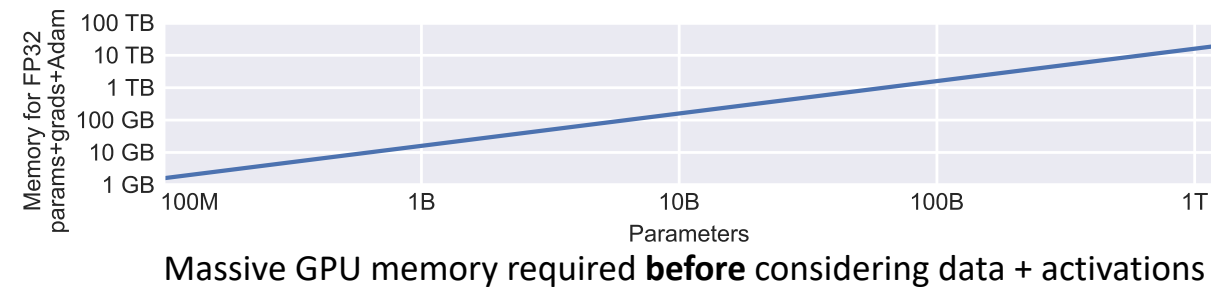
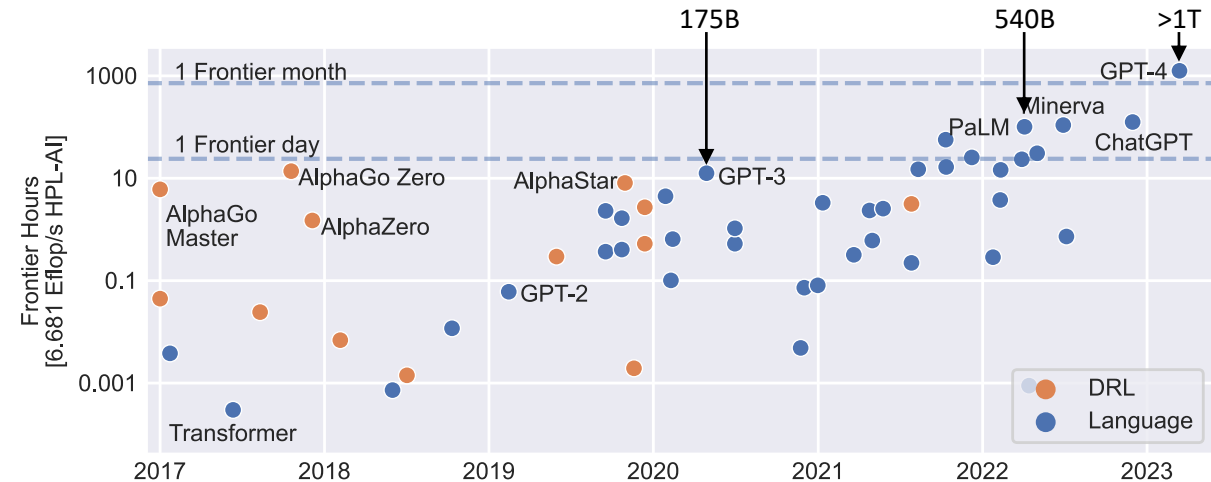


Tal Ben-Nun
Computer Scientist



The large-scale era of Machine Learning

- Training time measured in Exaflop days
- Models (reportedly) exceed 1T parameters
 - and are actually better as they grow...
- Without reducing memory consumption, we will not have the capacity to expand
- Are all the parameters and features *necessary*?

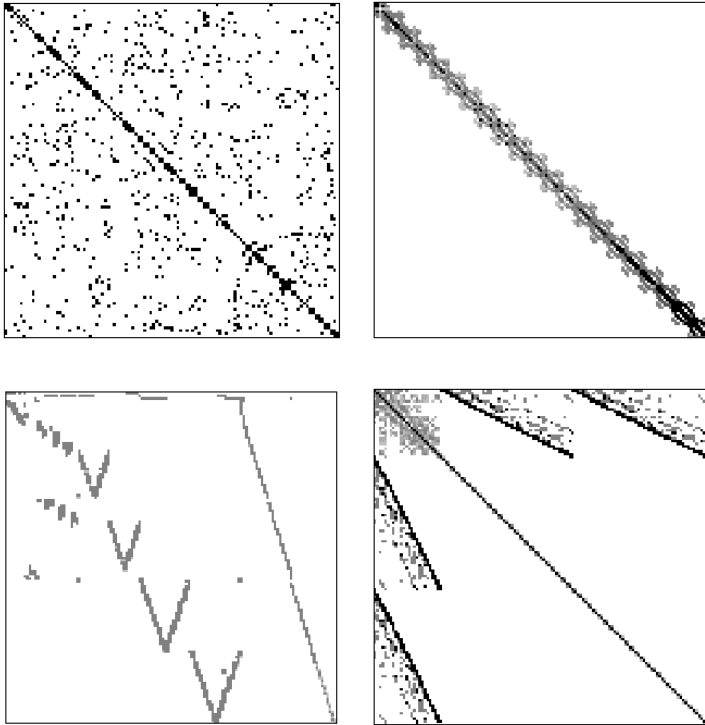


Sevilla et al., "Parameter, Compute, and Data Trends in Machine Learning", 2021

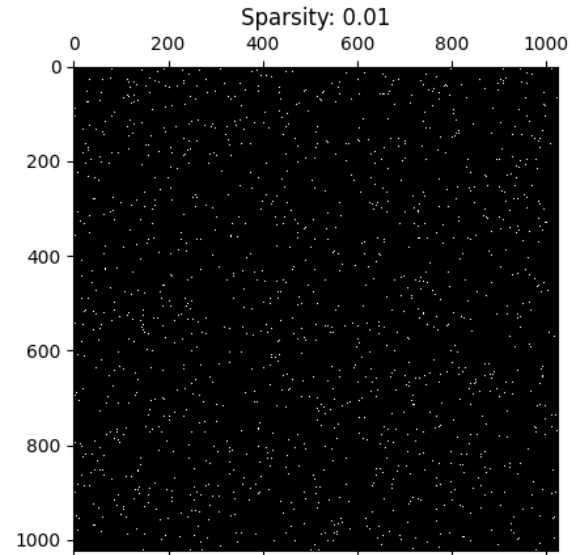
Sparsity in HPC

vs.

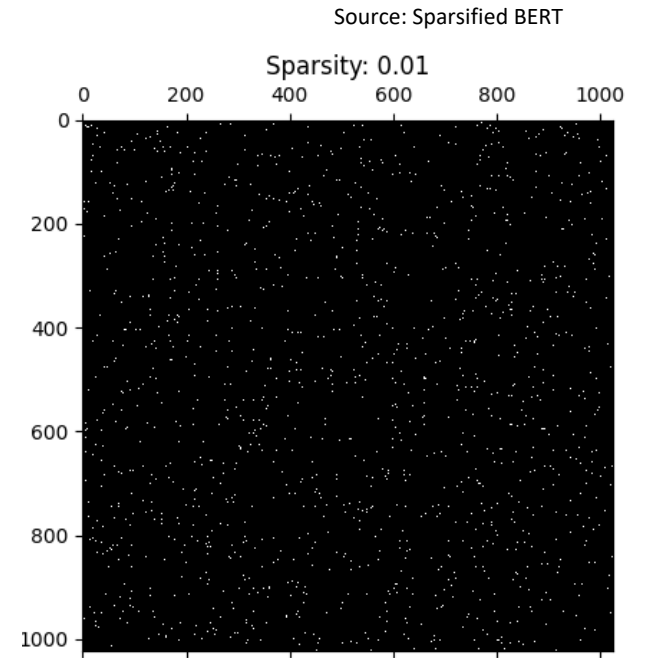
Machine Learning



Source: SuiteSparse Matrix Collection



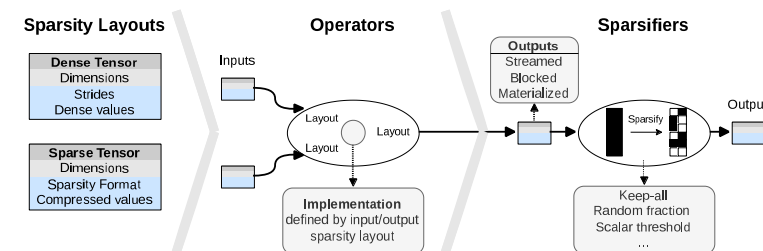
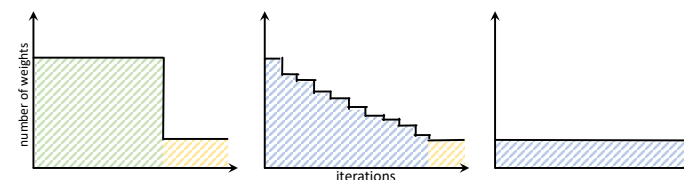
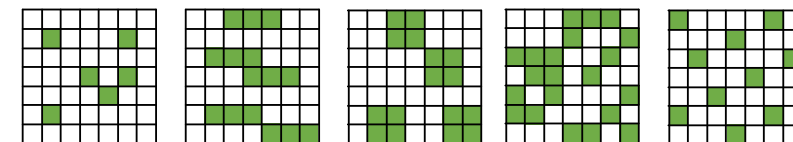
W_K



W_Q

Overview

- Sparse elements in deep learning
- Representations
- Scheduling strategies
- Hardware/software co-design and research tools



Primer on deep learning

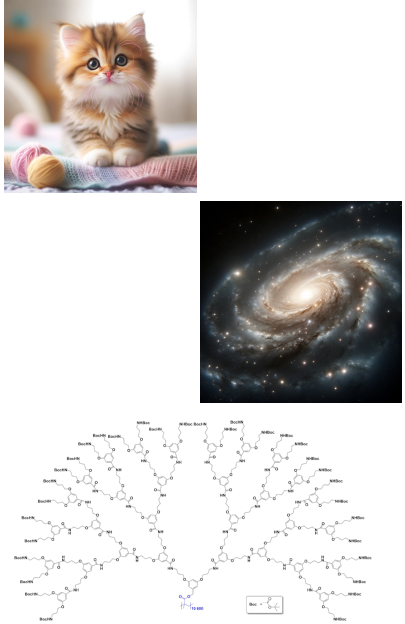
$$f(\mathbf{w}; x)$$

Input distribution \mathcal{X}

Output distribution \mathcal{Y}

Primer on deep learning

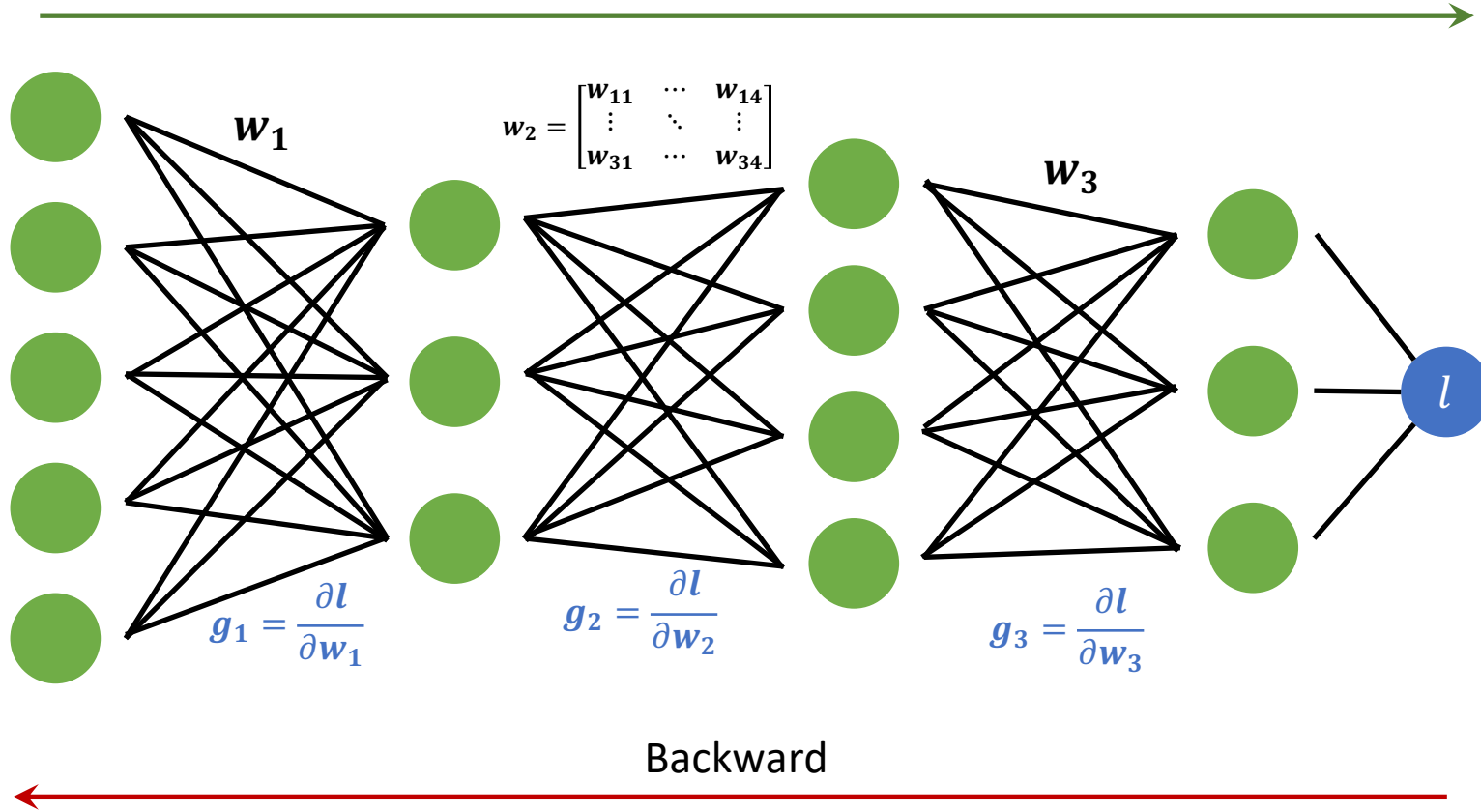
“This is a...”



“Pruning a neural network”

Input distribution \mathcal{X}

Forward



Backward

“...presentation”

Cat

Ω_M, σ_8, n_s

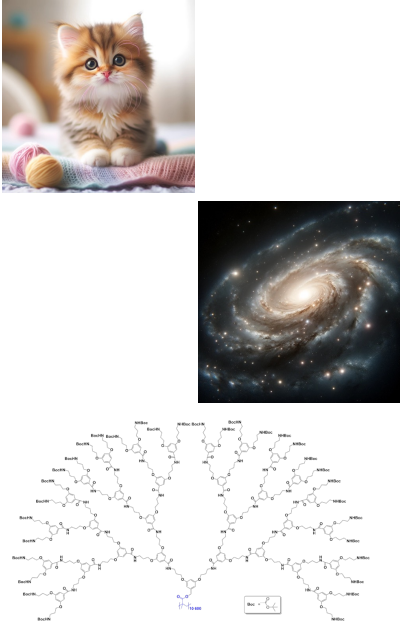
1.3 g/cc



Output distribution \mathcal{Y}

Where do we encounter sparsity?

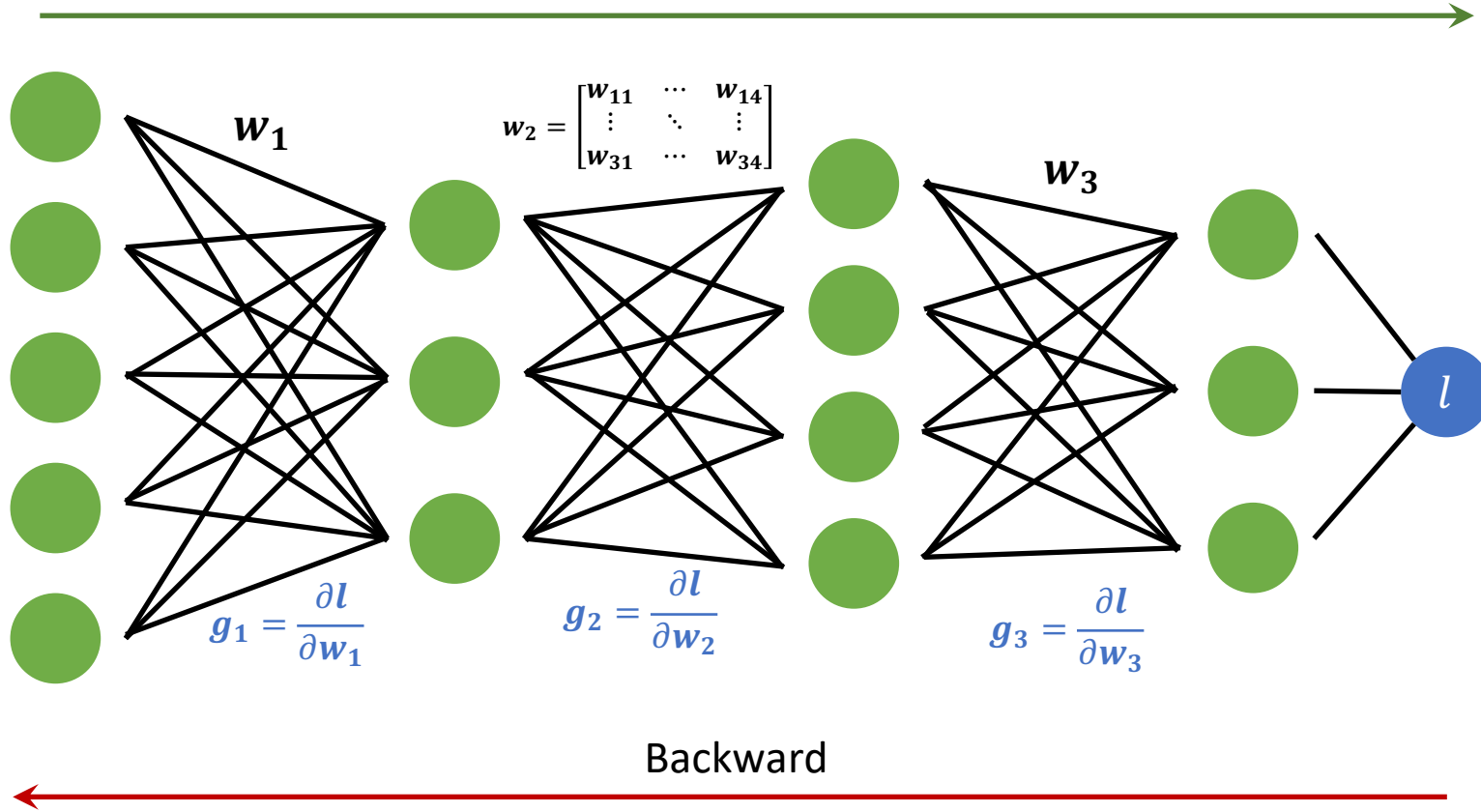
“This is a...”



“Pruning a neural network”

Input distribution \mathcal{X}

Forward



“...presentation”

Cat

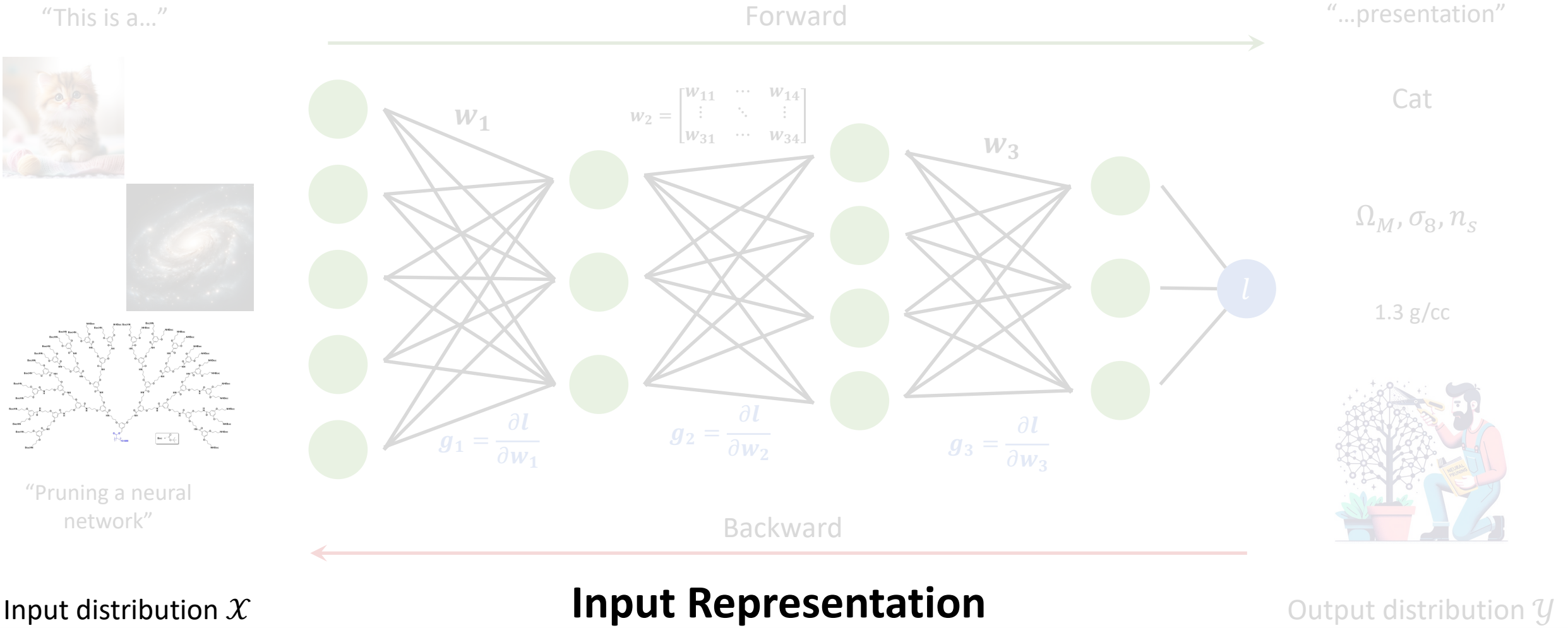
Ω_M, σ_8, n_s

1.3 g/cc

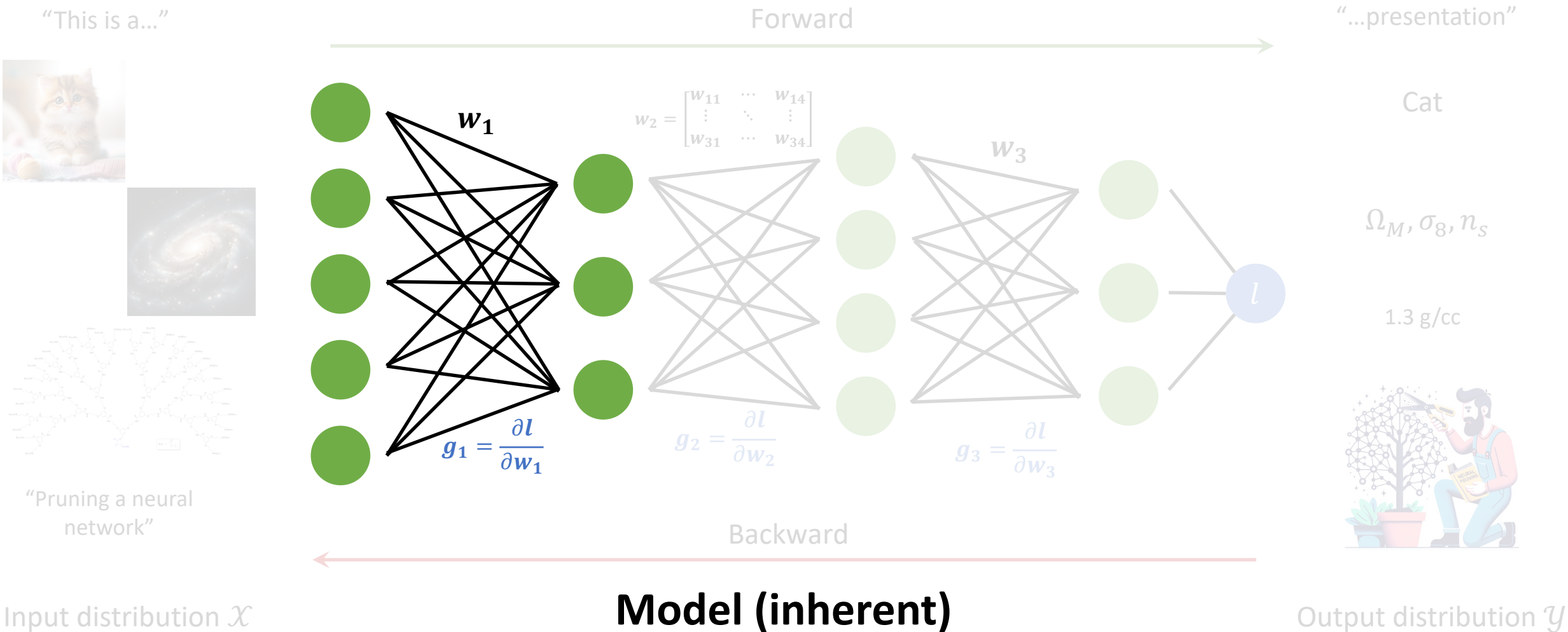


Output distribution \mathcal{Y}

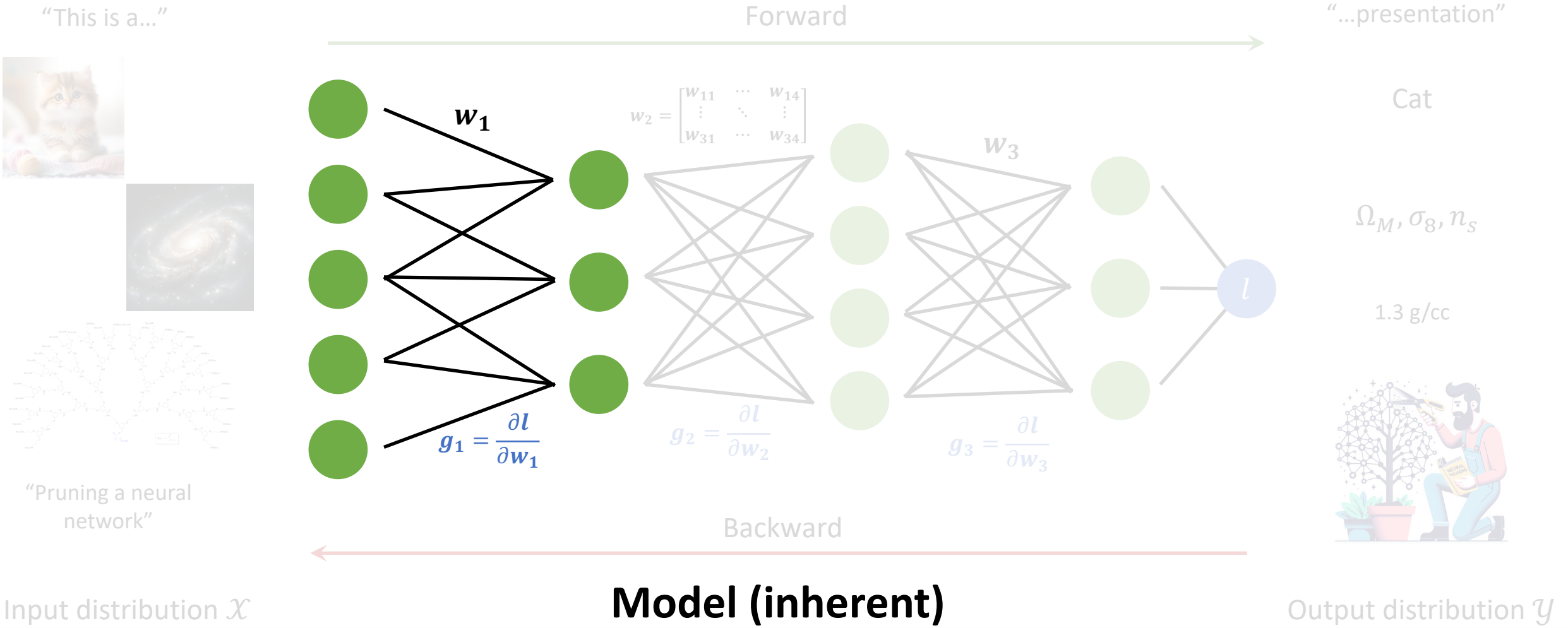
Where do we encounter sparsity?



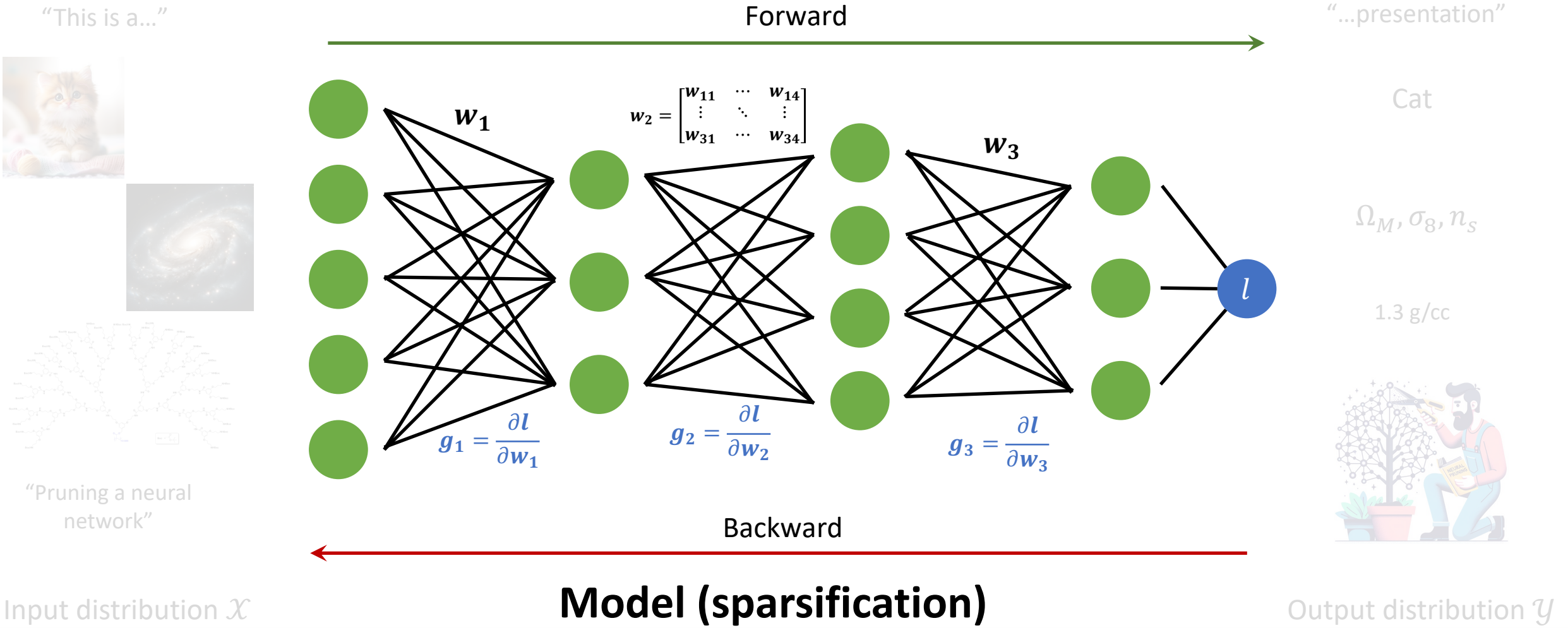
Where do we encounter sparsity?



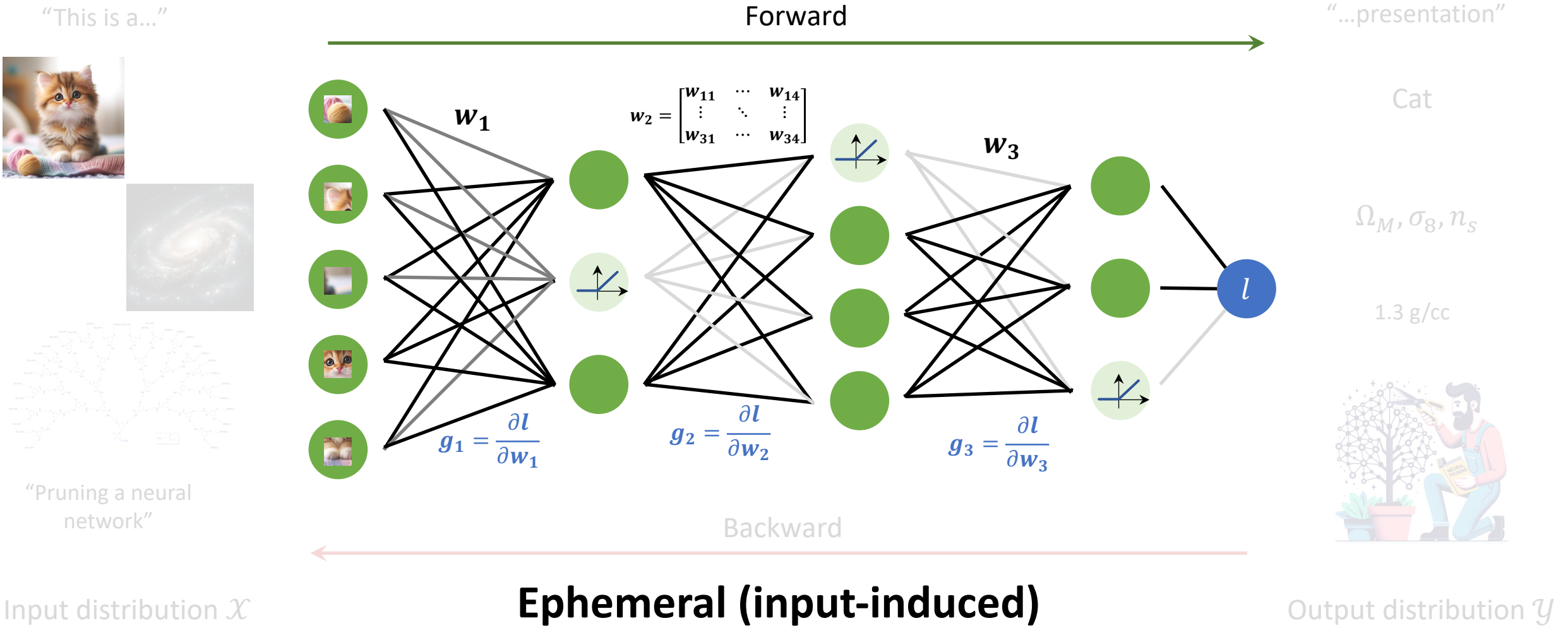
Where do we encounter sparsity?



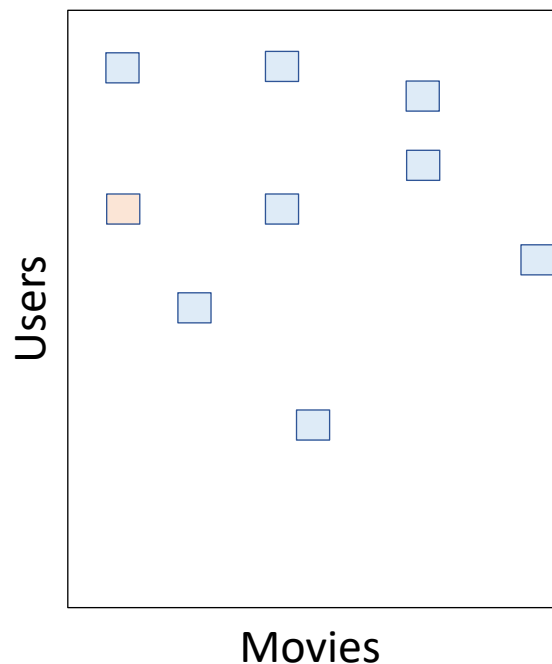
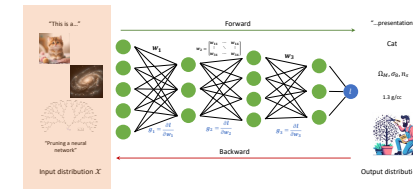
Where do we encounter sparsity?



Where do we encounter sparsity?

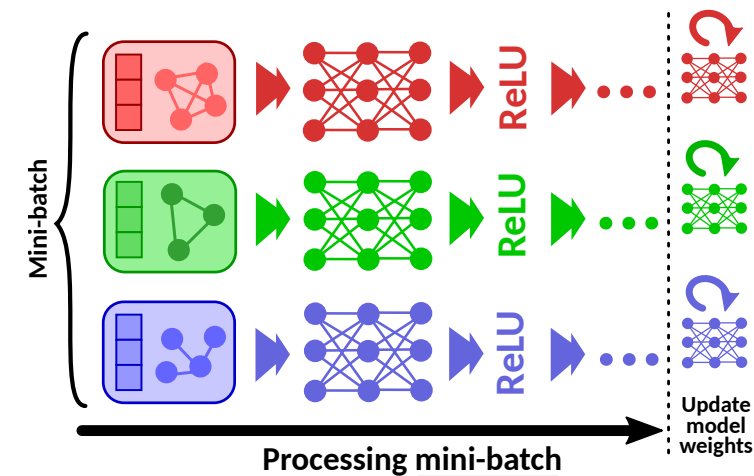


Input representations



“X likes The Barbie Movie, what else might they like?”
 Model output: “Oppenheimer”

Data parallelism



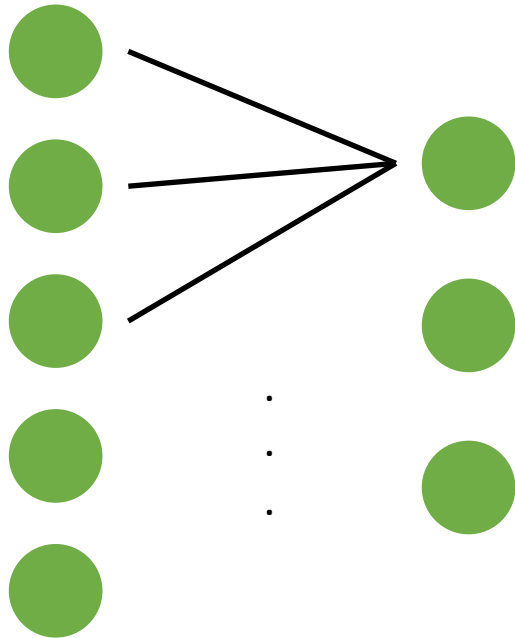
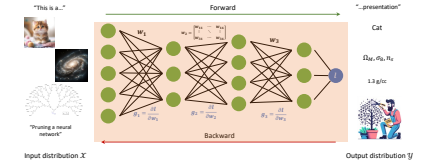
Pipeline parallelism



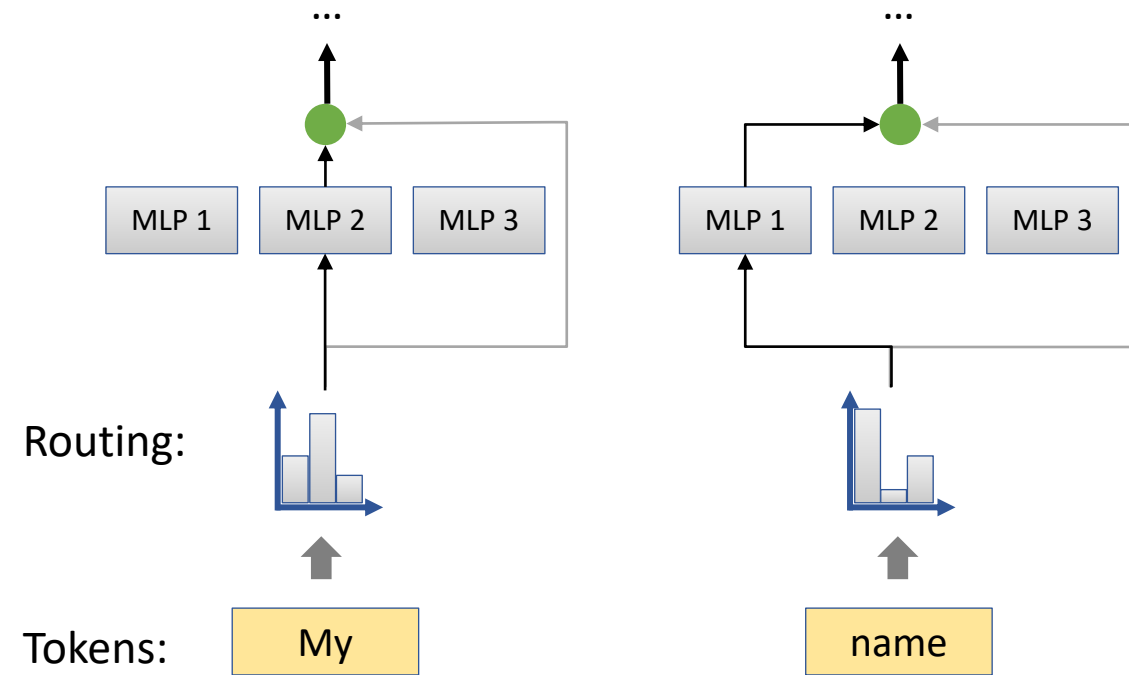
Recommendation Systems

Graph Neural Networks

Inherent sparsity in models



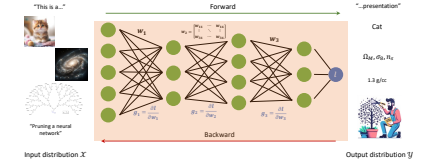
Convolution Operator



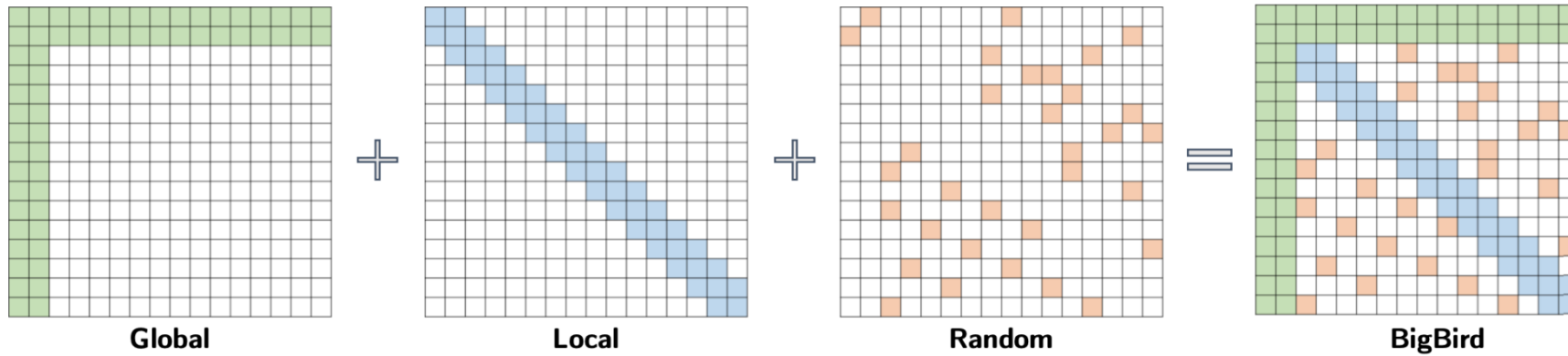
Mixture of Experts

Fedus, Zoph, Shazeer. "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity", JMLR'22

Inherent sparsity in models – sparse attention



- Attention is $O(n^2)$. n is now 128,000



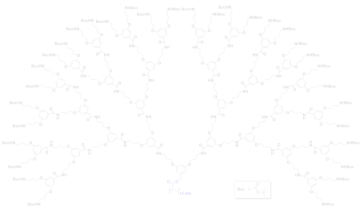
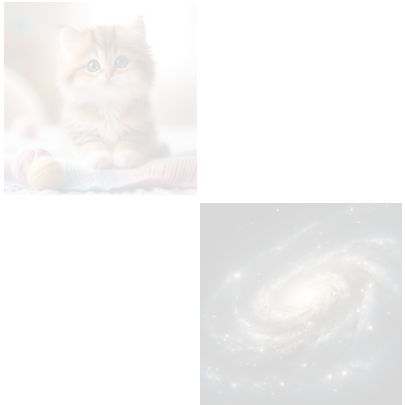
Source: <https://blog.research.google/2021/03/constructing-transformers-for-longer.html>

Model Architecture	Complexity per Layer	Sequential Operation
Recurrent	$O(n)$	$O(n)$
Transformer, (Vaswani et al., 2017)	$O(n^2)$	$O(1)$
Sparse Transformer, (Child et al., 2019)	$O(n\sqrt{n})$	$O(1)$
Reformer, (Kitaev et al., 2020)	$O(n \log(n))$	$O(\log(n))$
Linformer	$O(n)$	$O(1)$

Zaheer et al. "Big Bird: Transformers for Longer Sequences". NeurIPS 2020
 Wang et al. "Linformer: Self-Attention with Linear Complexity". arXiv:2006.04768

What about tried and true models?

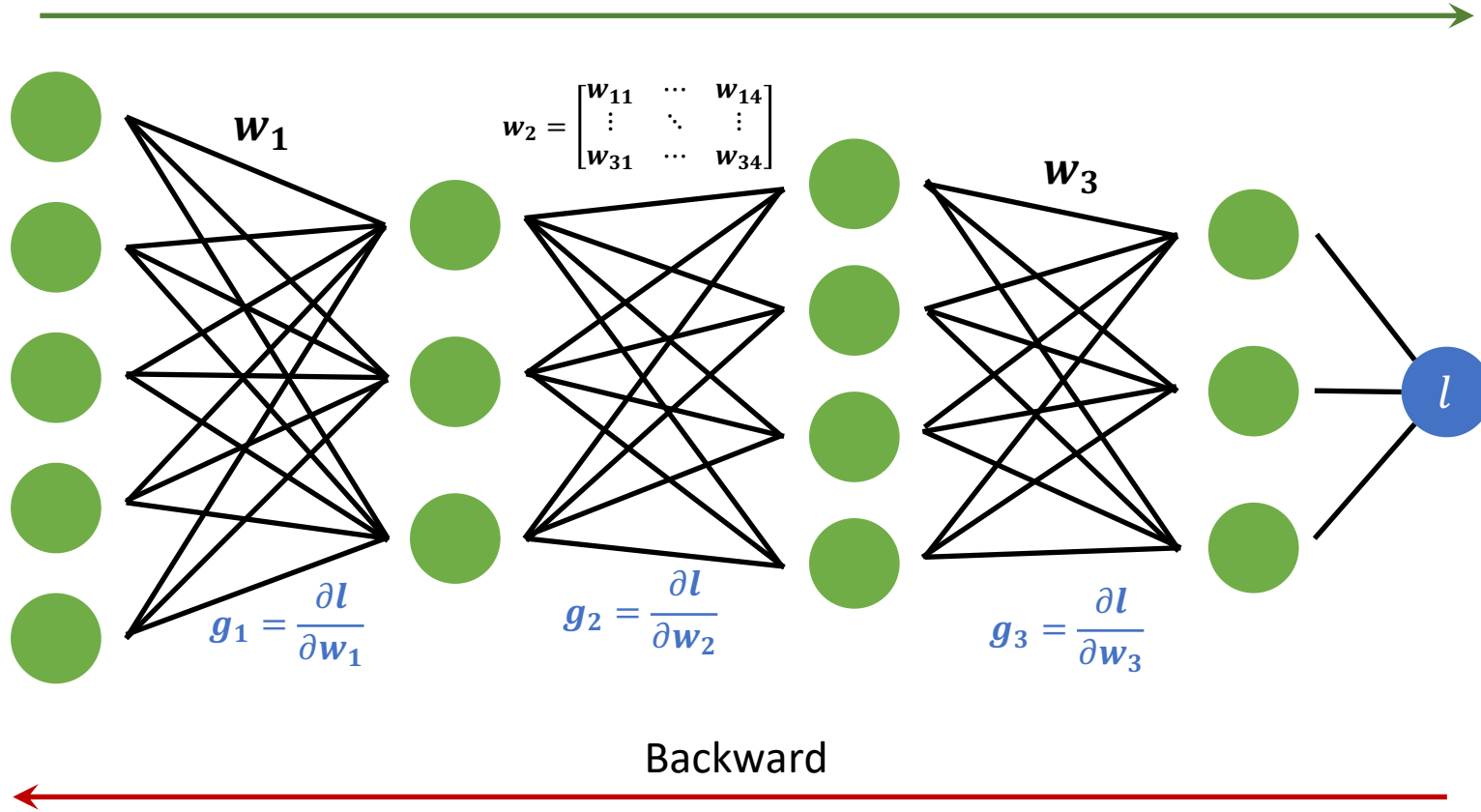
“This is a...”



“Pruning a neural network”

Input distribution \mathcal{X}

Forward



Backward

Model (sparsification)

“...presentation”

Cat

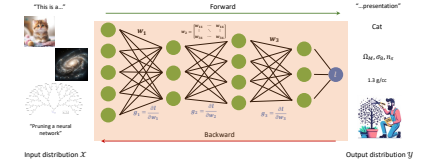
Ω_M, σ_8, n_s

1.3 g/cc



Output distribution \mathcal{Y}

Why should we sparsify?



- Reduces model parameters
- Improves generalization (Occam's Razor)
- Necessity: input representation or infeasibility
- State-of-the-art: 95% sparse ResNet-50, 50% sparse GPT models run at essentially same quality, up to 20x cheaper!

<https://jmlr.org/papers/volume22/21-0366/21-0366.pdf>

Sparsity in Deep Learning: Pruning and growth for efficient inference and training in neural networks

TORSTEN HOEFLER, ETH Zürich, Switzerland
DAN ALISTARH, IST Austria, Austria
TAL BEN-NUN, ETH Zürich, Switzerland
NIKOLI DRYDEN, ETH Zürich, Switzerland
ALEXANDRA PESTE, IST Austria, Austria

The growing energy and performance costs of deep learning have driven the community to reduce the size of neural networks by selectively pruning components. Similarly to their biological counterparts, sparse networks generalize just as well, if not better than, the original dense networks. Sparsity can reduce the memory footprint of regular networks to fit mobile devices, as well as shorten training time for ever growing networks. In this paper, we survey prior work on sparsity in deep learning and provide an extensive tutorial of sparsification for both inference and training. We describe approaches to remove and add elements of neural networks, different training strategies to achieve model sparsity, and mechanisms to exploit sparsity in practice. Our work distills ideas from more than 300 research papers and provides guidance to practitioners who wish to utilize sparsity today, as well as to researchers whose goal is to push the frontier forward. We include the necessary background on mathematical methods in sparsification, describe phenomena such as early structure adaptation, the intricate relations between sparsity and the training process, and show techniques for achieving acceleration on real hardware. We also define a metric of pruned parameter efficiency that could serve as a baseline for comparison of different sparse networks. We close by speculating on how sparsity can improve future workloads and outline major open problems in the field.

The supreme goal of all theory is to make the irreducible basic elements as simple and as few as possible without having to surrender the adequate representation of a single

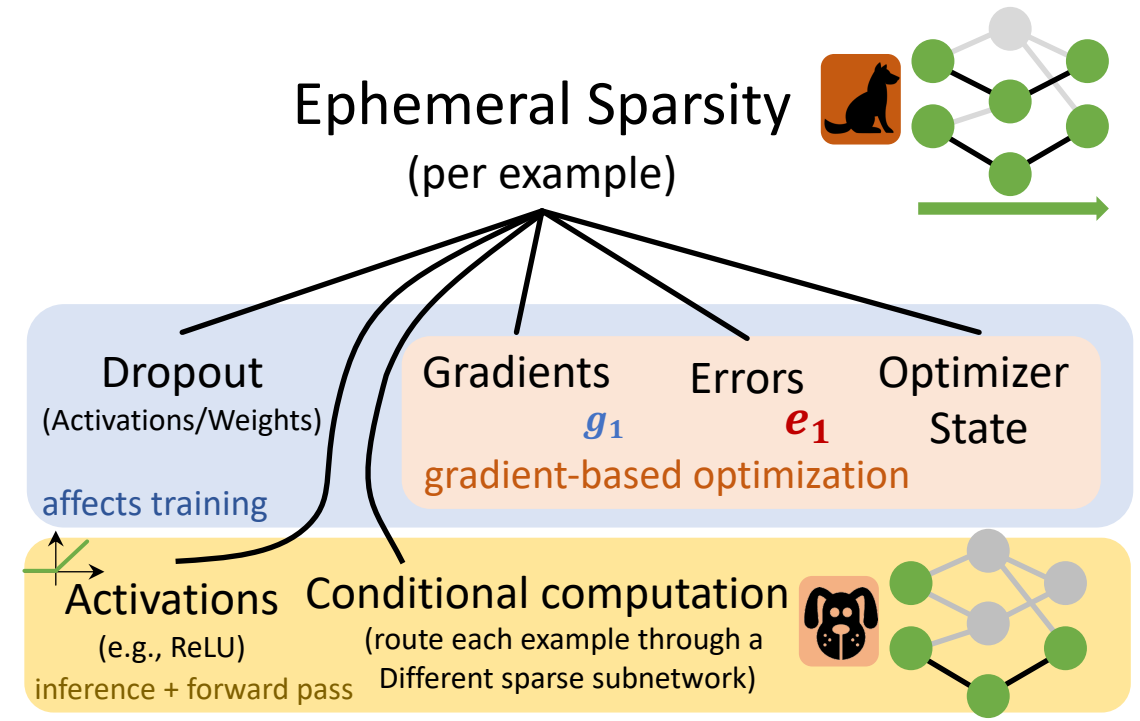
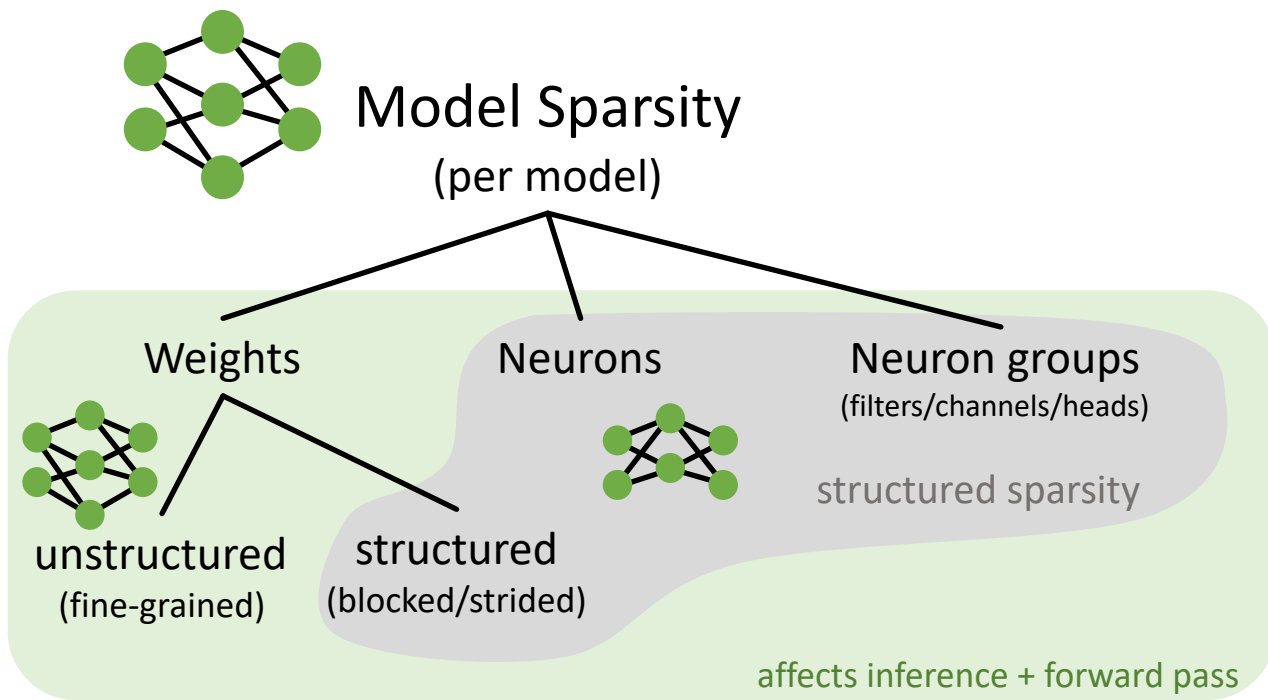
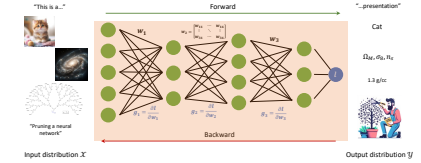
1 INTRODUCTION

Deep learning shows unparalleled promise for solving very complex real-world systems, drug discovery, and many more. With this development, the field is moving from traditional feature engineering to neural architecture engineering



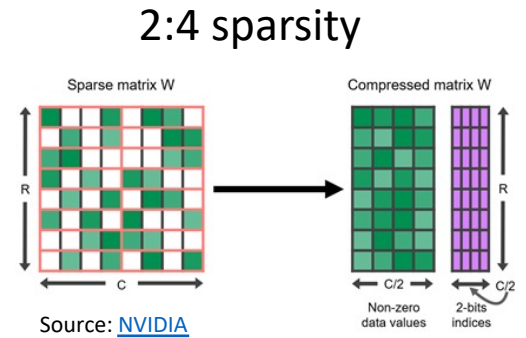
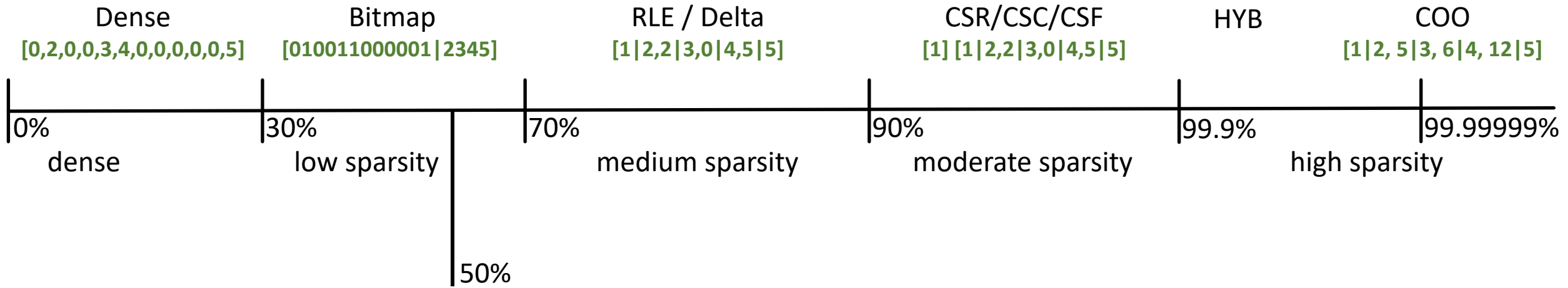
arXiv:2102.00554v1 [cs.LG] 31 Jan 2021

A taxonomy of model sparsification



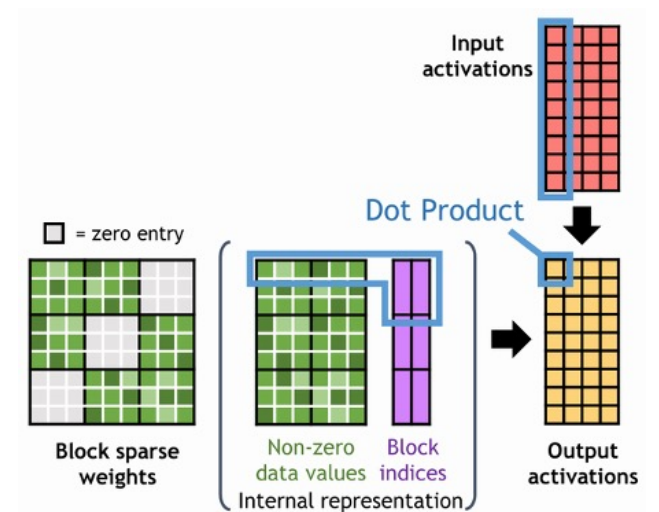
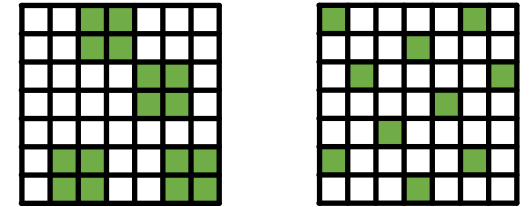
Hoefler, Alistarh, Ben-Nun, Dryden, Peste. "Sparsity in Deep Learning: Pruning and growth for efficient inference and training in neural networks", JMLR'21

Sparse representations in ML



Considerations in picking a representation

- Sparsity can be conditioned
 - The more constraints applied, the worse the end result
- Training and inference differ
 - Transposed representation necessary for backpropagation
 - Sparse representation may change during training!
- Hardware support
 - NVIDIA Sparse Tensor Cores
 - CSR/CSC can be used effectively for inference



Source: [NVIDIA](#)

Operation order in GNNs matters!

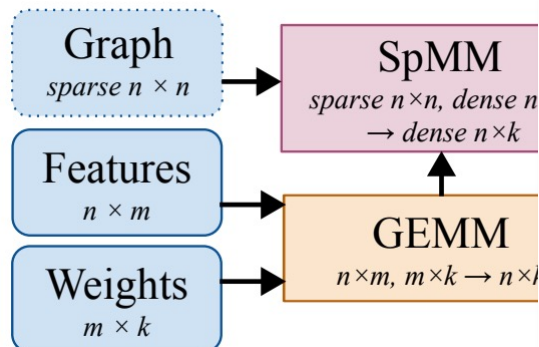


Table 1: Summary of differences between alternative GCN implementations. Entries in the “Usage condition” column were derived to minimize the total number of operations and memory transfers in the computation. In each pair of alternative schemes, the GEMMs are executed on the same shapes, while the SpMMs are executed on different shapes.

Operation	Scheme (Figure)	GEMM input sizes	SpMM input sizes	# Transient values	Usage condition
GCN forward	Transform-first (1a)	nm, mk	nn, nk	nk	$k < m$
	Propagate-first (1b)		nn, nm	nm	$k \geq m$
GCN backward with feature gradient	Fused-propagate (12a)	nm, mk	nn, nk	nk	$k < 2m$
	Split-propagate (12b)	nm, nk	$2 \times: nn, nm$	$2nm$	$k \geq 2m$
GCN backward no feature gradient	Fused-propagate (12a)	nm, mk	nn, nk	nk	$k < m$
	Split-propagate (12b)	nm, nk	nn, nm	nm	$k \geq m$

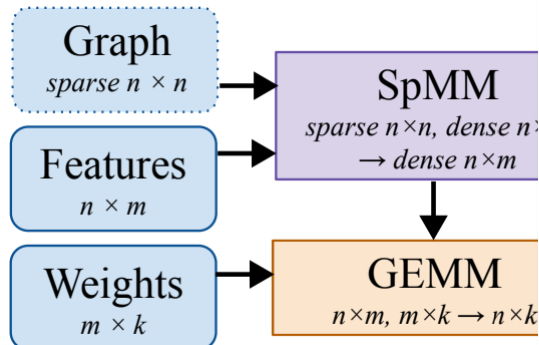
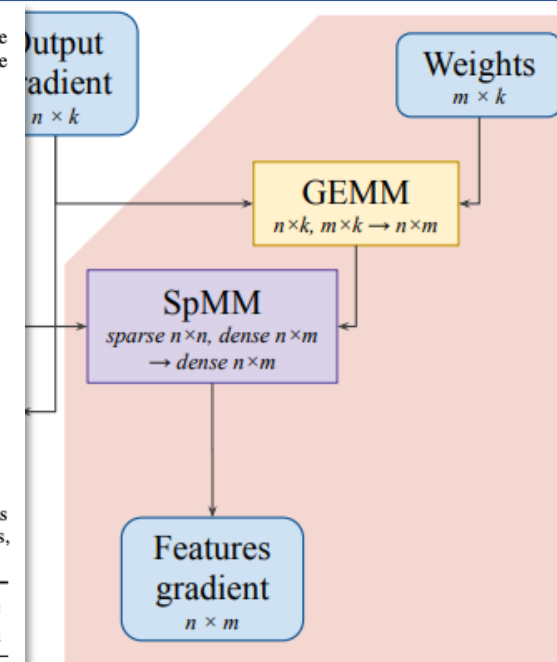


Table 2: Comparison of GCN computation schemes with caching. The case presented in Figures 1b and 12b is omitted because it always executes more operations than the scheme using caching. In each pair of alternative schemes, the GEMMs are executed on the same shapes, while the SpMMs are executed on different shapes.

Operation	Forward scheme (Figure)	Backward scheme (Figure)	GEMM input sizes	SpMM input sizes	# Transients (fwd; bwd)	Usage condition
GCN, with feature gradient	Transform-first (1a)	Fused-propagate (12a)	nm, mk	$2 \times: nn, nk$	$nk; nk$	$k < m$
	Propagate-first with caching (13a)	Split-propagate with caching (13b)	nm, mk	$2 \times: nn, nm$	$nm; nm$	$k \geq m$
GCN, no feature gradient	Transform-first (1a)	Fused-propagate (12a)	nm, mk	$2 \times: nn, nk$	$nk; nk$	$k < \frac{1}{2}m$
	Propagate-first with caching (13a)	Split-propagate with caching (13b)	nm, nk	nn, nm	$nm; 0$	$k \geq \frac{1}{2}m$

Forward



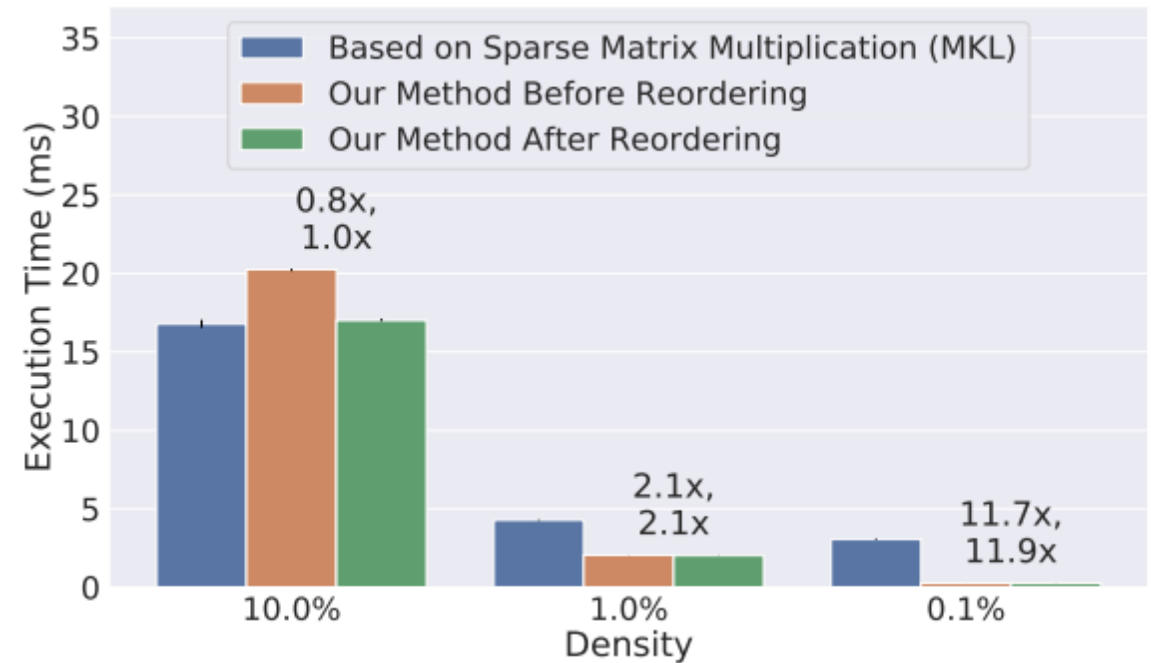
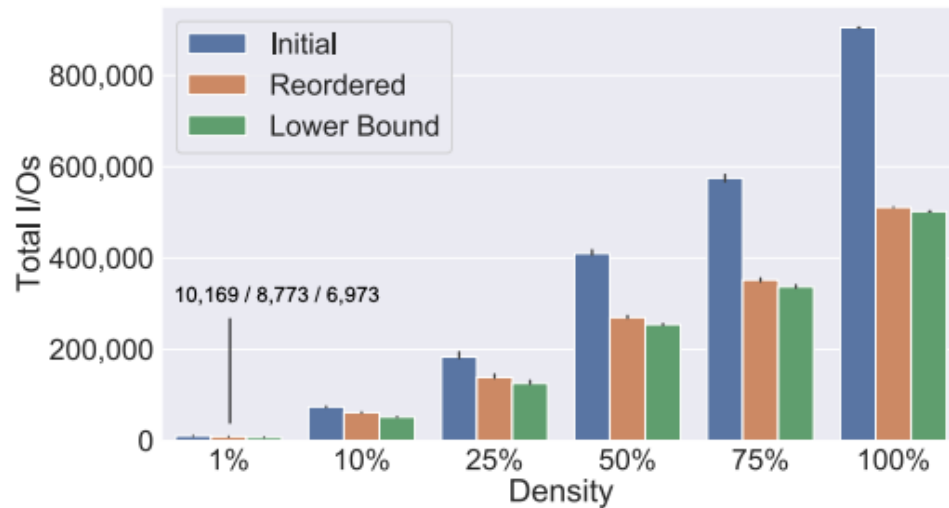
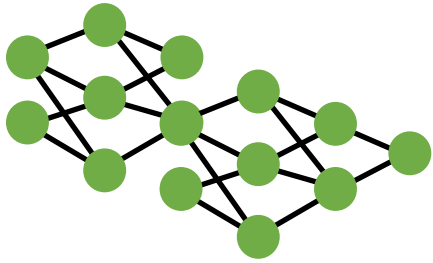
required only if there is a preceding layer

Backward

Up to 1.94x speedup over PyTorch Geometric!

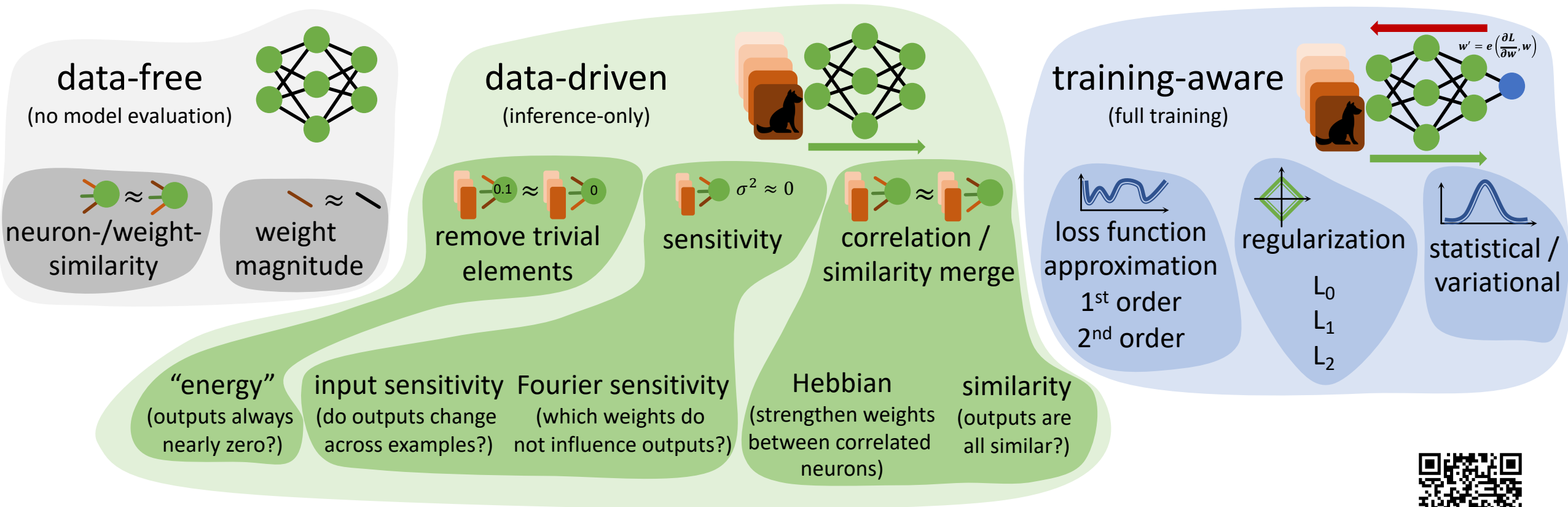
Bazinska et al. “Cached Operator Reordering: A Unified View for Fast GNN Training”. arXiv:2308.12093

NN evaluation can be performed as graph traversal



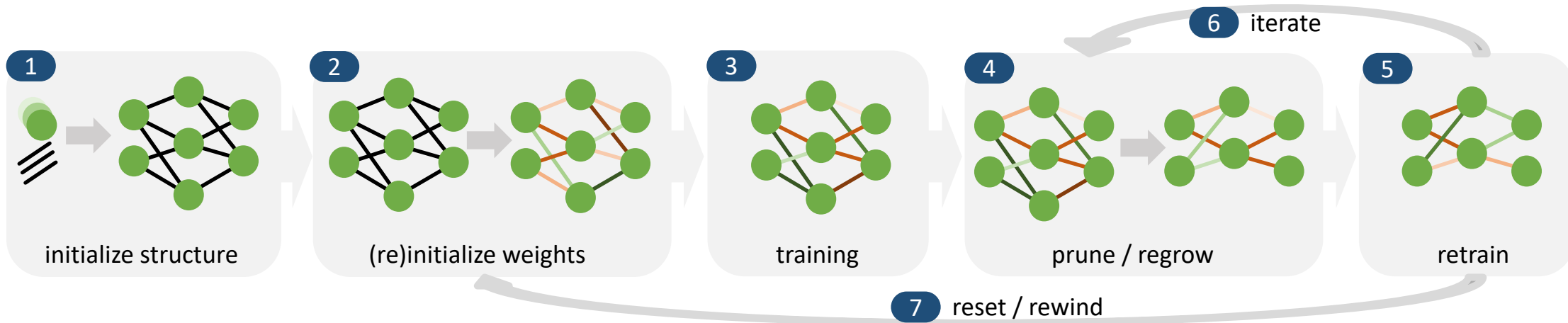
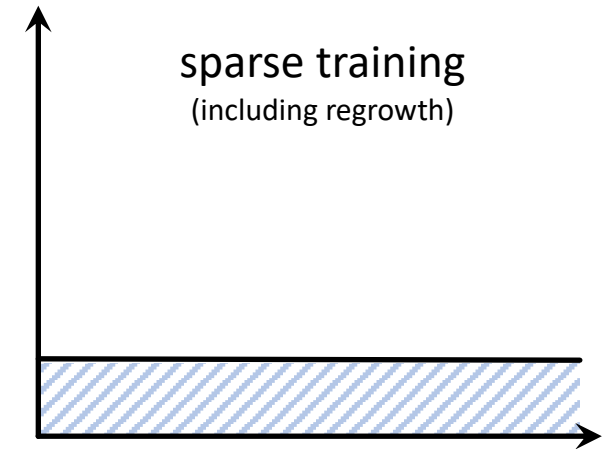
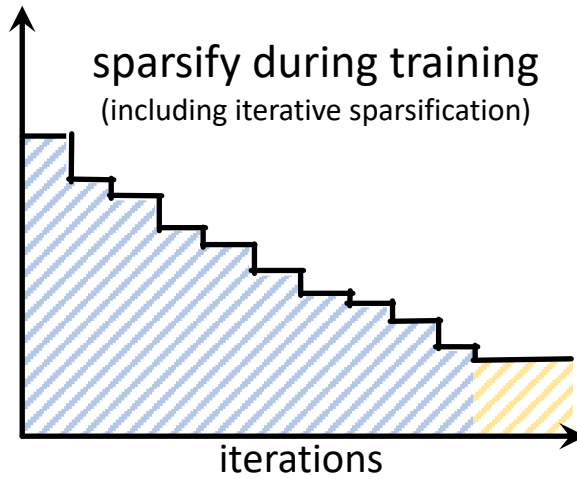
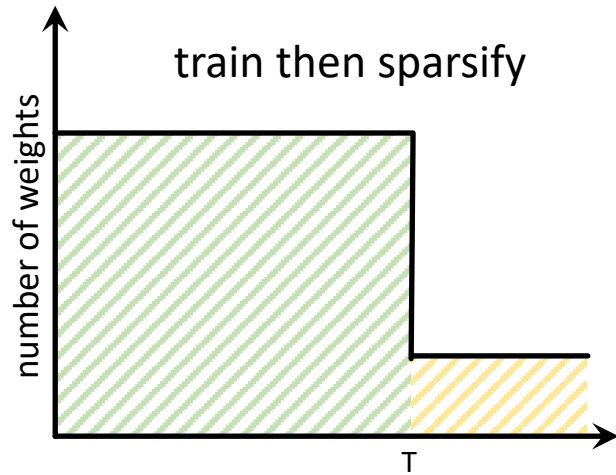
Performance on BERT_{LARGE} vs. MKL

Model sparsification techniques



Hoefler, Alistarh, Ben-Nun, Dryden, Peste. “Sparsity in Deep Learning: Pruning and growth for efficient inference and training in neural networks”, JMLR’21

When to sparsify?

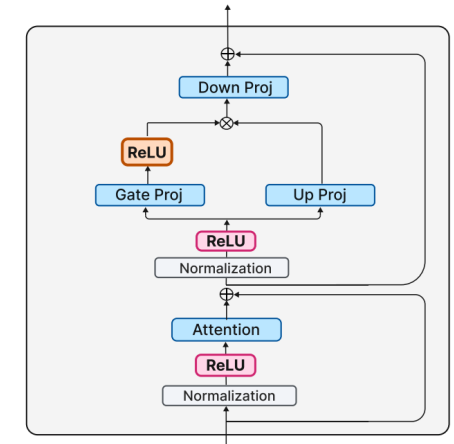


When to sparsify?

- Early structure adaptation takes place
- In large models: we may not have the budget to retrain
 - SparseGPT uses second-order information on a set of examples
- Model can be adapted: Inject ReLU to promote activation sparsity
- Repositories such as the Sparse Zoo contain recipes for many models

“Best” practices:

- Gradual magnitude pruning will get you most of the way to 90%
 - Higher sparsity or less drop will require more advanced techniques
- Be mindful of which layers you sparsify and their position in the model



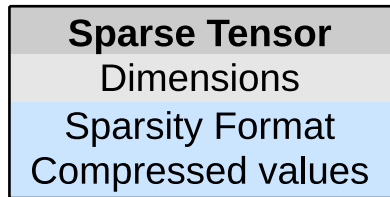
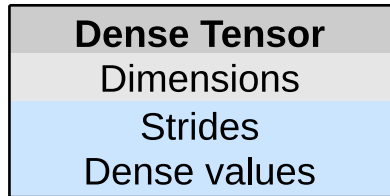
Model Name	Text Generation	Items/second throughput	ms latency	MB file size	Perplexity
opt-1.3b-opt_pretrain-pruned50_quantW8A8	0.22	4.55K	1.1	17.4	
codegen_mono-350m-bigpython_bigquery_thepile-pruned50	13.3	182	774.7	3.9	
mpt-7b-dolly_mpt_pretrain-pruned50	8.2	122	16.1	16.0	
mpt-7b-mpt_chat_mpt_pretrain-pruned50_quantized	114	42.7	3.2	24.1	
opt-2.7b-opt_pretrain-pruned50_quantW8A8	0.136	7.34K	1.9	15.9	

Frantar & Alistarh, “SparseGPT: Massive Language Models Can Be Accurately Pruned in One-Shot”. arXiv:2301.00774
Mirzadeh et al. “ReLU Strikes Back: Exploiting Activation Sparsity in Large Language Models”. arXiv:2310.04564
Sparse Zoo, <https://sparsezoo.neuralmagic.com/>

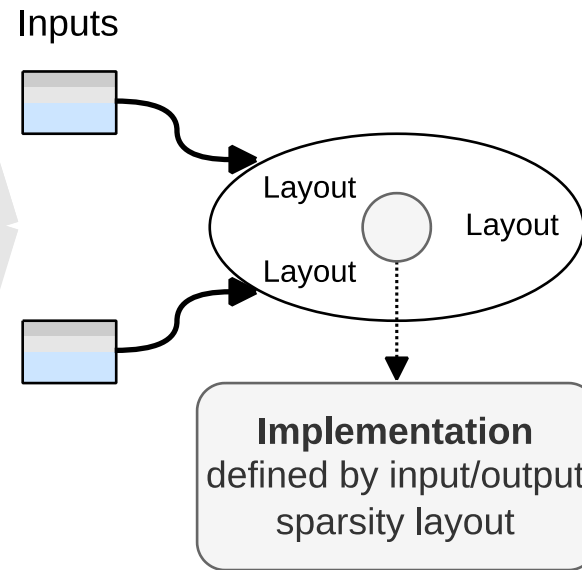
Programming Sparse Models – Meet PyTorch STen



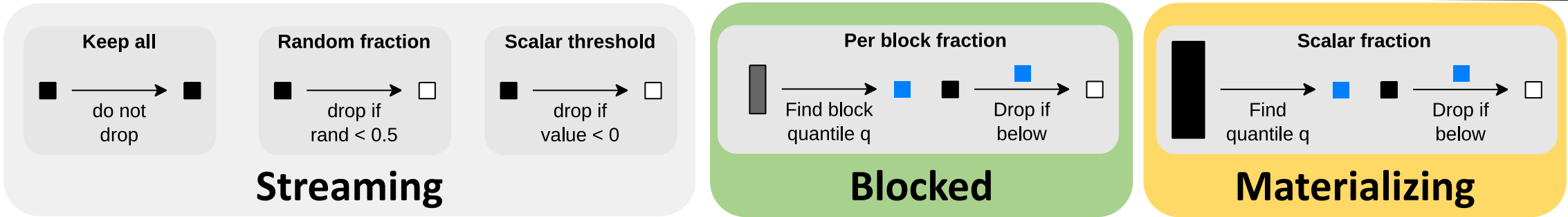
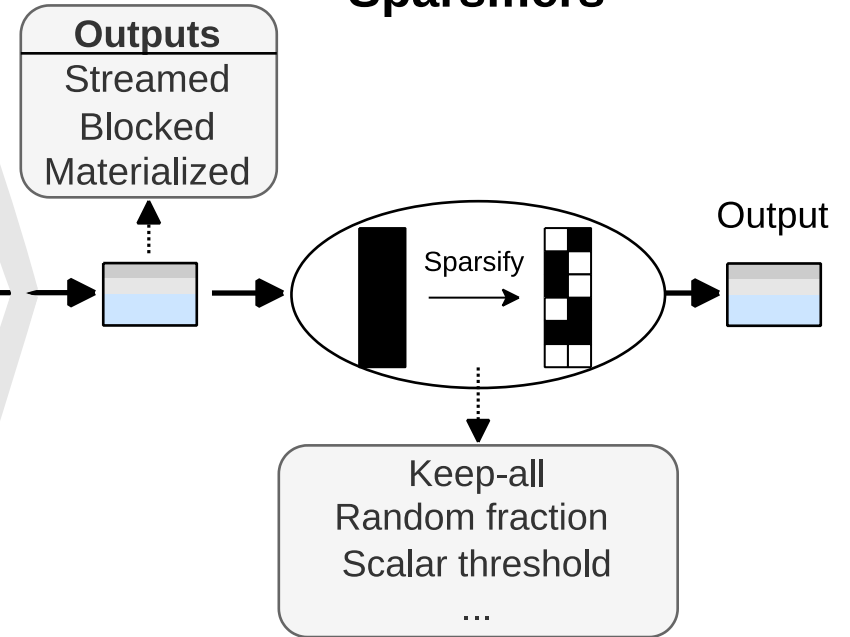
Sparsity Layouts



Operators



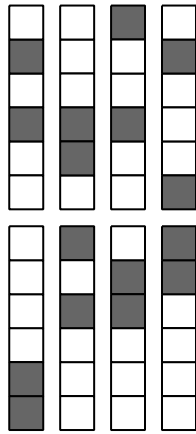
Sparsifiers



Ivanov et al. "STen: Productive and Efficient Sparsity in PyTorch". arXiv:2304.07613.

STen Performance

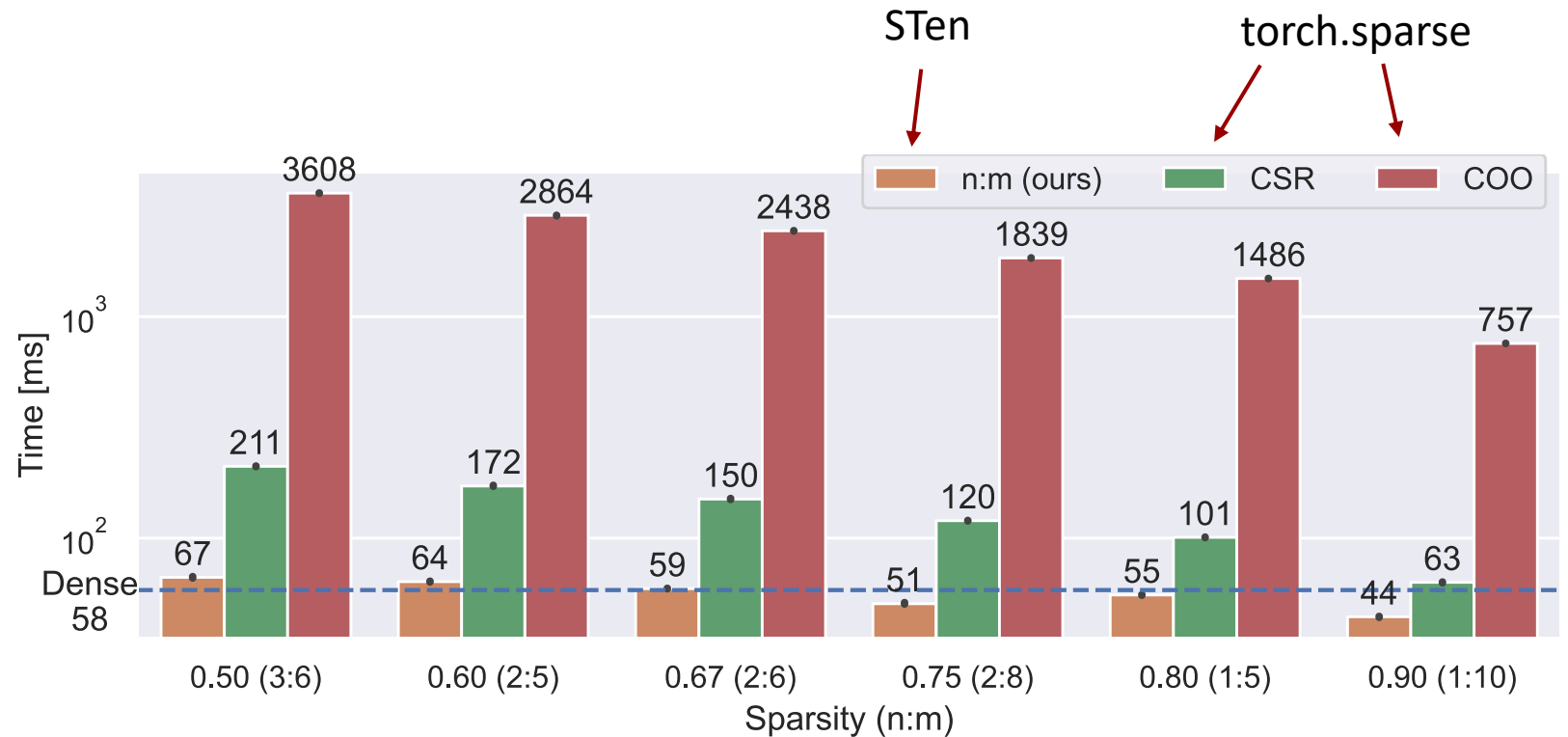
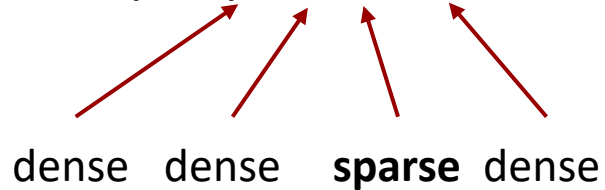
2:6 sparse format



Custom implementation of matrix multiplication:

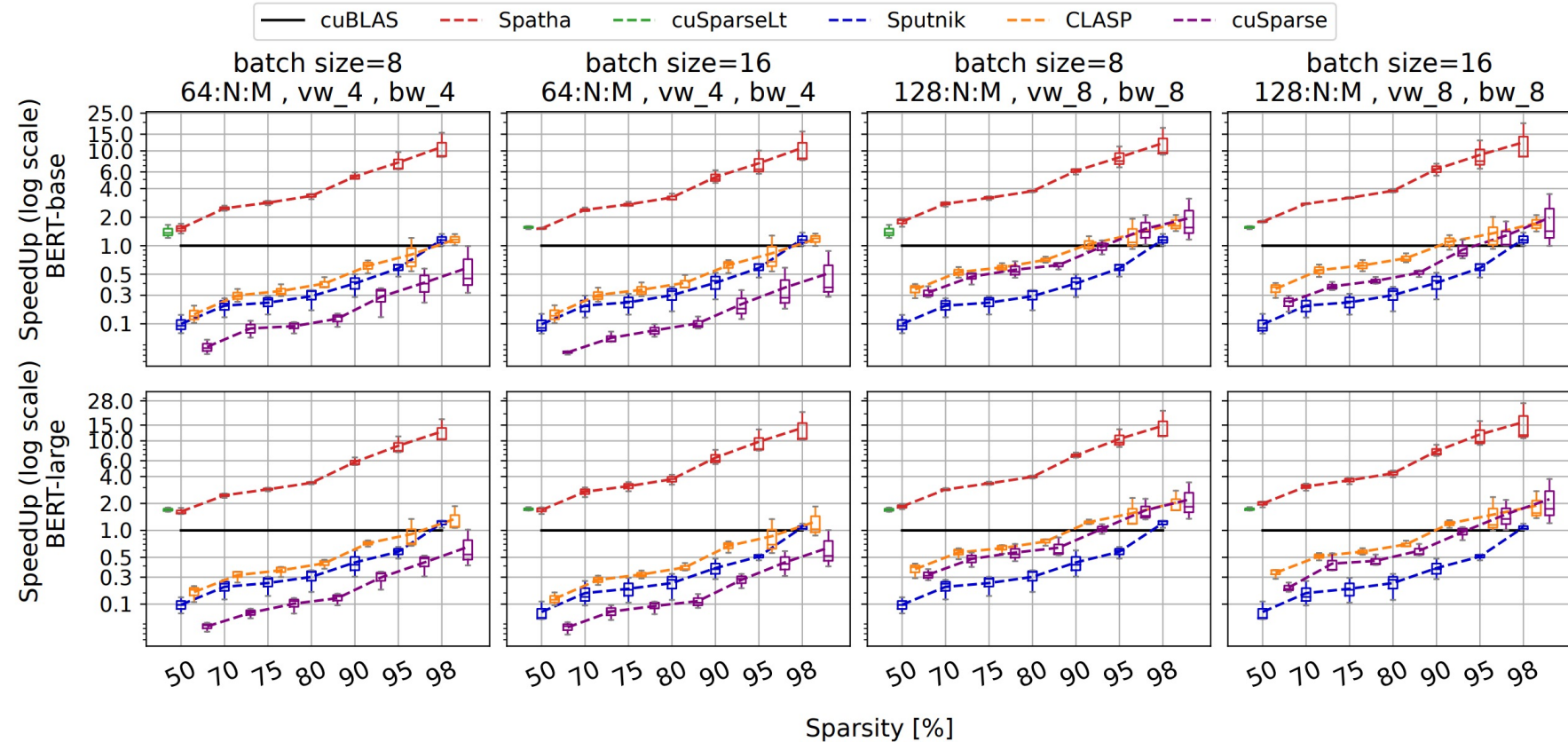
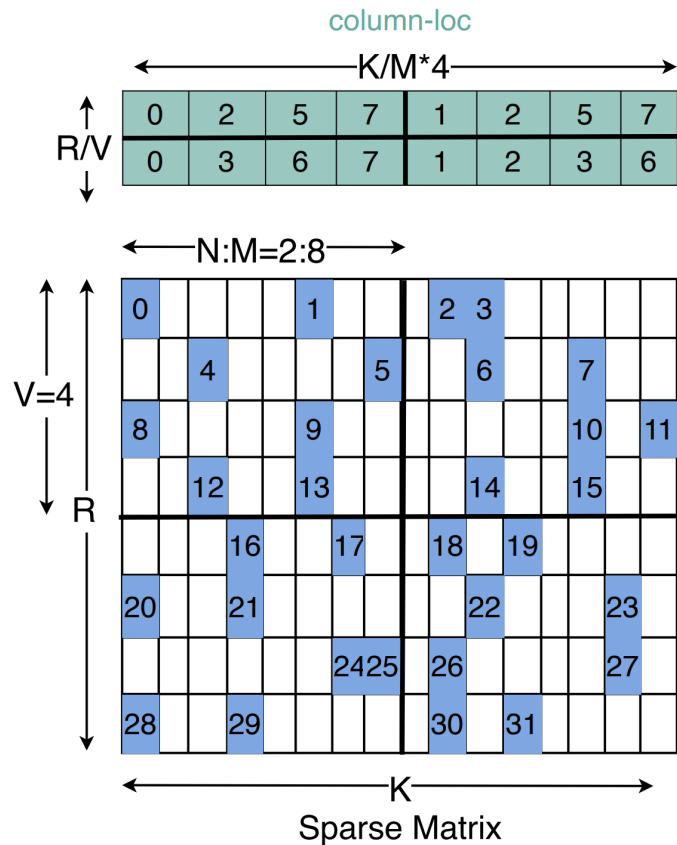
sparse @ dense -> dense

Linear layer: $y = x \mathbf{W} + b$

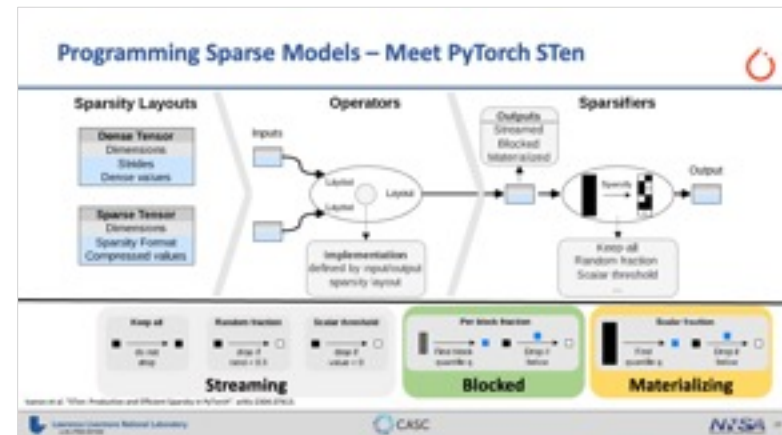
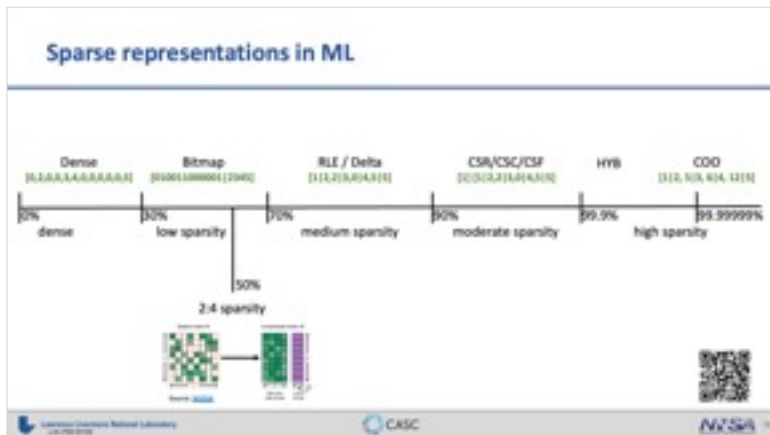
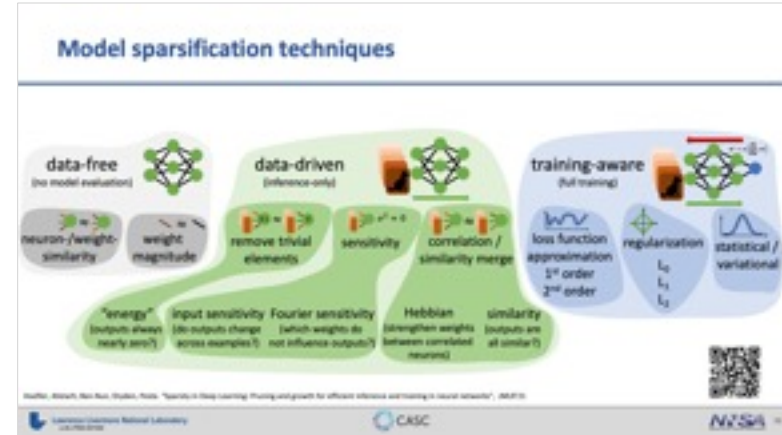
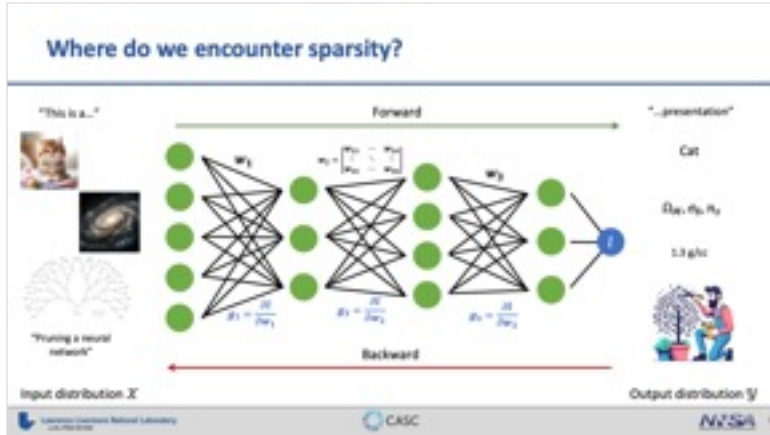


32% speedup

VENOM Performance



Conclusion





CASC

Center for Applied
Scientific Computing



Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.