





T. HOEFLER

arXiv:2506.06941v1

From Big Large Language Models to Fast R

Three Eras in The Age of Computation.

with contributions by the whole SPCL deep learning team (M. Besta, J. Barth, E. Schr Belk, S. Scott, D. Goel, M. Castro) and collaborators (D. Alistarh and others) DP2E-AI, Paris, June, 2025



Hybrid reasoning model that pushes the frontier for coding and AI agents, featuring a 200K context window

dels

rosoft Azure (M. Heddes, J.

The Illusion of Thinking:

Understanding the Strengths and Limitations of Reasoning Models

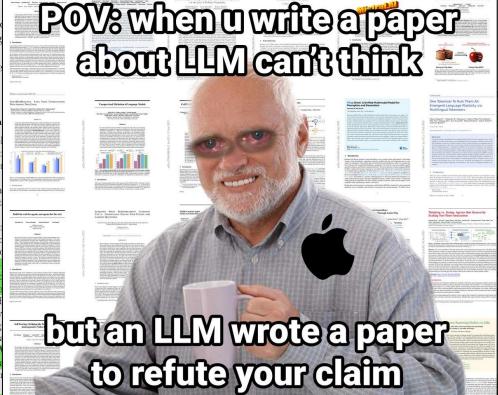
via the Lens of Prol

Parshin Shojaee*† Iman Mirza Maxwell Horton Samy Beng

Abstra

(LRMs) that generate detailed thinking processes demonstrate improved performance on reasoning be ing properties, and limitations remain insufficiently cus on established mathematical and coding benchn ever, this evaluation paradigm often suffers from da into the reasoning traces' structure and quality. In gaps with the help of controllable puzzle environme tional complexity while maintaining consistent logic of not only final answers but also the internal reas "think". Through extensive experimentation acros face a complete accuracy collapse beyond certain of intuitive scaling limit: their reasoning effort increas declines despite having an adequate token budget. counterparts under equivalent inference compute, complexity tasks where standard models surprising tasks where additional thinking in LRMs demonstrated where both models experience complete collapse. computation: they fail to use explicit algorithms also investigate the reasoning traces in more dept and analyzing the models' computational behavior and ultimately raising crucial questions about th

Recent generations of frontier language mode



The Illusion e Illusion of Thinking

> A Comn on Shojaee et al. (2025)

A. Lawsen

June 10, 2025

Abstract

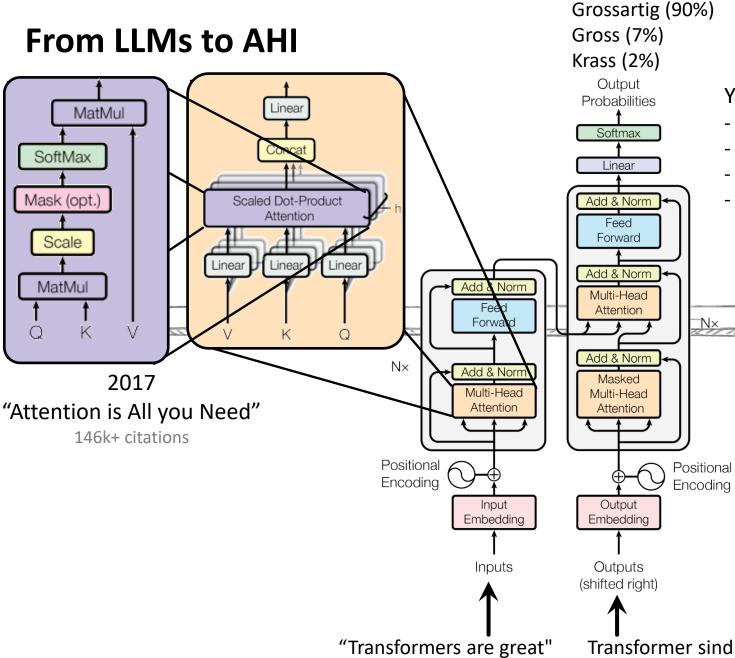
25) report that Large Reasoning Models (LRMs) exhibit "accuracy colizzles beyond certain complexity thresholds. We demonstrate that their ect experimental design limitations rather than fundamental reasoning failveals three critical issues: (1) Tower of Hanoi experiments systematically token limits at reported failure points, with models explicitly acknowledgn their outputs; (2) The authors' automated evaluation framework fails to easoning failures and practical constraints, leading to misclassification o Most concerningly, their River Crossing benchmarks include mathematances for $N \geq 6$ due to insufficient boat capacity, yet models are scored ving these unsolvable problems. When we control for these experimental g generating functions instead of exhaustive move lists, preliminary exiple models indicate high accuracy on Tower of Hanoi instances previously failures. These findings highlight the importance of careful experimental AI reasoning capabilities.











You can explain the computation to your grandmother!

- Three simple kernels: MMM, Softmax, Layernorm
- >95%+ matrix multiplication
- Great fit for HPC GPUs
- Easy to parallelize

Text is encoded as tokens (very important!)

- Tokens are offsets into learned vector tables
- Often learned based on statistics
- Most common sub-strings (e.g., Byte Pair Encoding)
- Think of them as vectors
- Word2vec: "Efficient Estimation of Word Representations in Vector Space" (45k+ citations)







From LLMs to AHI

Poor English input: I eated the purple berries. Good English output: I ate the purple berries.

Poor English input: Thank you for picking me as your designer. I'd appreciate it. Good English output: Thank you for choosing me as your designer. I appreciate it. Poor English input: The mentioned changes have done. or I did the alteration that you

requested. or I changed things you wanted and did the modifications.

Good English output: The requested changes have been made. or I made the alteration that you

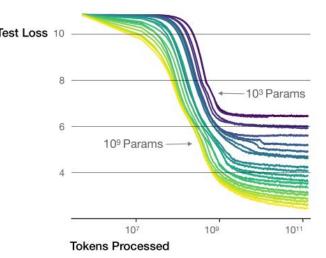
requested. or I changed things you wanted and made the modifications.

Poor English input: I'd be more than happy to work with you in another project.

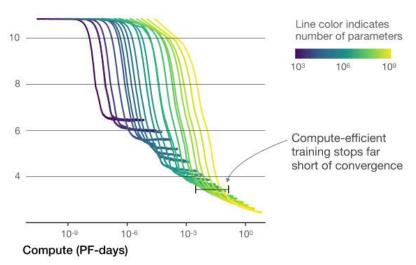
Good English output: I'd be more than happy to work with you on another project.

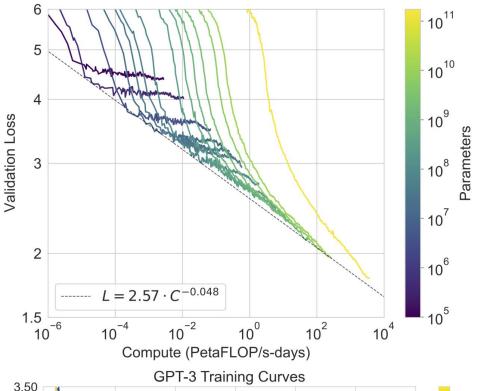
Scaling Laws for Neural Language Models

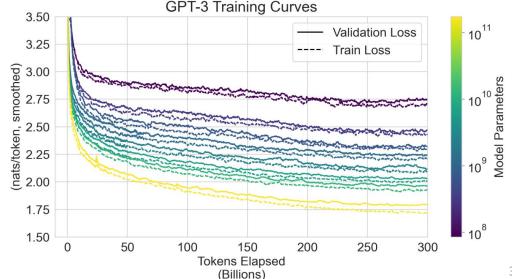
Larger models require fewer samples to reach the same performance



The optimal model size grows smoothly with the loss target and compute budget













From LLMs to AHI

Microsoft invests \$1 billion in OpenAI to pursue holy grail of artificial intelligence

Building artificial general intelligence is OpenAl's ambitious goal
By James Vincent | Jul 22, 2019, 10:08am EDT

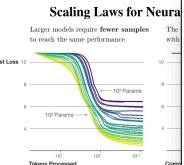
"BERT: Pre-training of Deep Bidirectional

Transformers for Language Understanding"

Toslg upveils Doio supercomputer: world's

Tesla unveils Dojo supercomputer: world's new most powerful AI training machine

Fred Lambert - Aug. 20th 2021 3:08 am PT 💆 @FredericLambert



2020 - **GPT-3** (202 "Language Mo Few-Shot Lea

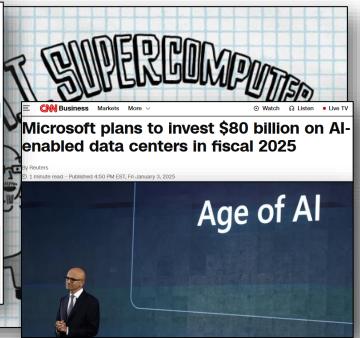
37k+ citati

Trump's AI Push: Understanding
The \$500 Billion Stargate
Initiative

Garth Friesen Contributor ©
Specialist in global markets, economics and alternative investments.

Deplated Jan 24, 2025, 07:25em EST







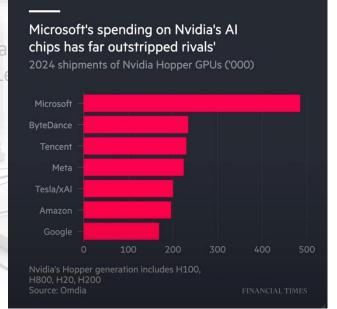
BABY STEPS Google artificial intelligence supercomputer creates its own 'Al child' that can outperform its human-made rivals

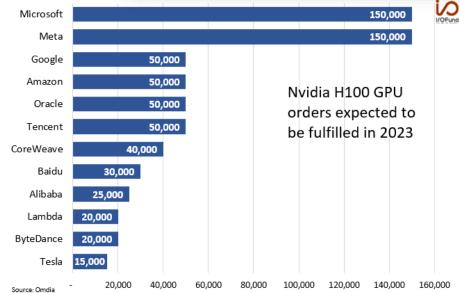
The NASNet system was created by a neural network called AutoML earlier this year

The NASNet system was created by a neural network called AutoML earlier this year

oding 15:22, 5 Dec 2017 | **Updated**: 11:27, 6 Dec 2017

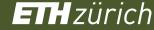
5.22, 5 bec 2017 | **Opunces**. 11.27, 6 bec 201











Supercomputers fuel Modern Al



Tesla unveils Dojo supercomputer: world's new most powerful AI training machine

Fred Lambert - Aug. 20th 2021 3:08 am PT 🂆 @FredericLambert



Microsoft invests \$1 billion in OpenAl to pursue holy grail of artificial intelligence

Building artificial general intelligence is OpenAl's ambitious goal

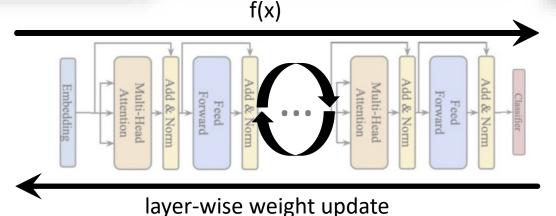
By James Vincent | Jul 22, 2019, 10:08am EDT







A robot may not injure a human being or, through inaction, allow a human being to come to ____



PaLM-540B: 1.4 trillion tokens

ImageNet (22k): A few TB

• Actually: the whole internet!

PaLM-540B: 118 (complex) layers
 540 bn parameters (1 TiB in fp16)
 2048-token "sentences"

harm 0.74 injury 0.28 now 0.07 never 0.04 pain 0.33 boat 0.02 house 0.02

now 0.00
never 0.00
pain 0.00
boat 0.00
house 0.00

1.00

harm

PaLM-540B: 256k token dict

takes weeks to train

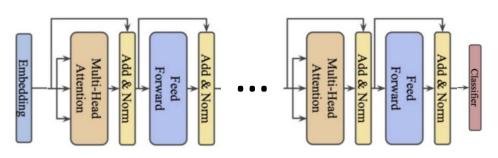






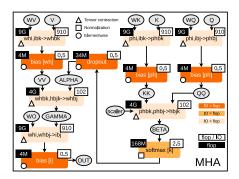
Data Movement Is All You Need: A Case Study on Optimizing Transformers (arXiv:2007.00072)

GPT/BERT encoder



Express as Dataflow





OpenAl booth at NeurIPS 2019 in Vancouver, Canada Image Credit: Khari Johnson / VentureBeat

Last week, OpenAI published a paper detailing GPT-3, a machine learning model that achieves strong results on a number of natural language benchmarks. At 175 billion parameters, where a parameter affects data's prominence in an overall prediction, it's the largest of its kind. And with a memory size exceeding 350GB, it's one of the priciest, costing an estimated \$12 million to train.

	highly		
Operator class	optimized	% flop	% Runtime
Tensor contract	ion	99.80	61.0
Statistical norm	alization	0.17	25.5
Element-wise		0.03	13.5
		0.2%	39%

Our performance improvement for BERT-large

- 30% over PyTorch
- 20% over Tensorflow + XLA
- 8% over DeepSpeed

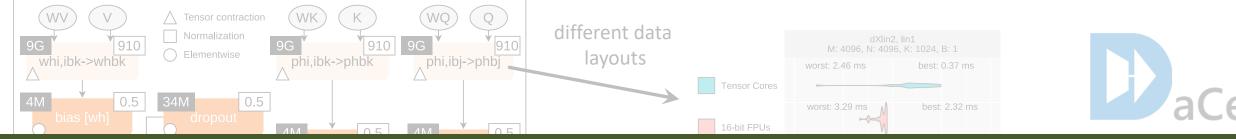
est. savings on AWS over PyTorch: \$85k for BERT, \$3.6M GPT-3







Data Movement Is All You Need: A Case Study on Optimizing Transformers (arXiv:2007.00072)





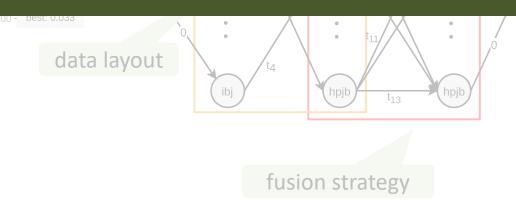
These Ideas are in Production Today



Reducing Cost Remains Imperative to Continue Scaling

Full BERT encoder layer performance (ms)

	TF+XLA			
Forward	3.2	3.45	2.8	2.63
Backward	5.2	5.69	4.8	4.38









From LLMs to AHI

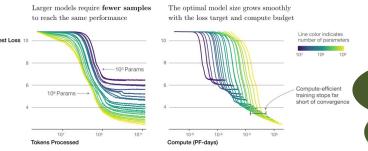


2018 - **BERT**

"BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding"

122k+ citations

Scaling Laws for Neural Language Models



2020 - **GPT-3** (2020, scaling laws)

"Language Models are Few-Shot Learners"

37k+ citations

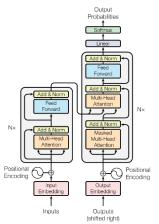
How to turn this into a business serving millions of customers?



2017 - Transformers

"Attention is All you Need"

146k+ citations



2019 - **GPT-2**

"Language Models are Unsupervised Multitask Learners"

14k+ citations



2022 – **ChatGPT** (RLHF, 2023, DPO)

"Training language models to follow instructions with human feedback"

14k+ citations









From LLMs to AHI

How to turn this into a business serving millions of customers?

Scaling Laws for Neural Language Models
models require fewer samples
in the same performance

The optimal model size grows smoothly with the loss target and compute budget

2020, scaling laws)

iguage Models are



Compete through Openness

2023 – **Llama** (Qwen, Grok, etc.) "LLaMA: Open and Efficient

Foundation Language Models"

11k+ citatio

Needs even more (pre)training compute!

Reduce cost

more better data,
more training compute

optimize models computationally

reduce hardware cost and increase efficiency

2022 – **ChatGPT** (RLHF, 2023, DPO)
"Training language models to follow instructions with human feedback"

14k+ citations

era of data scaling.



Optimization
Determines the
Future of Al





Optimization Determines the Future of Al

We need a Scientific Approach to it

Next, let's see how to improve cost by 1,000x



Moving Data is Most Expensive!

Techniques to Shrink ML Data



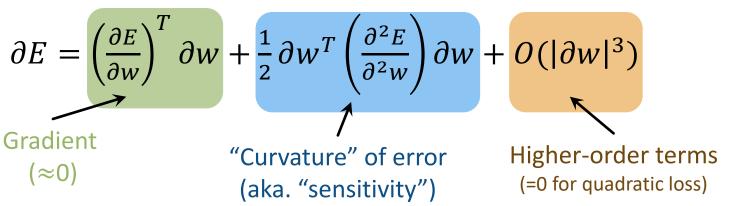


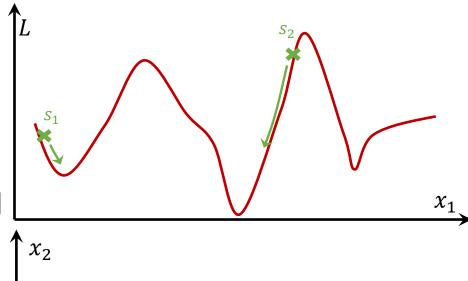


Quantization – Running Gigantic LLMs on Reasonable Systems (arXiv:2210.17323



- Brains have limited precision! Why are we computing with FP32?
 - For technical reasons (SGD, optimization, how we quantize)
 - Neurons in Hippocampus can "reliably distinguish 24 strengths" [1] 4.6 bits of information!
- PaLM-540B has up to 540 billion parameters
 - 1.08 TiB in FP16/BF16, 540 GiB in FP8 🕾
 - Rounding to <5 bits is not so simple
 - Requires some foundation and many tricks
- Consider "error landscape" of a trained model with weights w [2]











Quantization - Running Gigantic LLMs on Reasonable Systems (arXiv:2210.17323



• Quantization objective for low precision rounded weights \widehat{w} argmin $_{\widehat{w}} \|wx - \widehat{w}x\|^2$

Solve PTQ optimization problem row by row of w

- Round row and push the error forward using the inverse Hessian
- Update Hessian for each column

Tricks

- Block updates for better locality (10x speedup)
- Use Cholesky to invert Hessian (higher stability)
- Work one transformer block at a time (6 operators fit in memory)
- Use quantized input from previous blocks for block i

Results

- Generative inference 2-4x faster
- 3 bits → 66 GiB, fits in a single (high-end) A100 GPU!

Model	FP16	1024	512	256	128	64	32	3-bit
OPT-175B	8.34	11.84	10.85	10.00	9.58	9.18	8.94	8.68
BLOOM	8.11	11.80	10.84	10.13	9.55	9.17	8.83	8.64

GPTQ: ACCURATE POST-TRAINING QUANTIZATION FOR GENERATIVE PRE-TRAINED TRANSFORMERS

A PREPRINT

Elias Frantar*

Klosterneuburg, Austria elias.frantar@ist.ac.at

Torsten Hoefler

ETH Zurich Switzerland htor@inf.ethz.ch

Saleh Ashkboos

ETH Zurich Switzerland saleh.ashkboos@inf.ethz.ch

Dan Alistarh

IST Austria & Neural Magic, Inc. Klosterneuburg, Austria dan.alistarh@ist.ac.at

ARSTRACT

Generative Pre-trained Transformer (GPT) models set themselves apart through breakthrough performance across complex language modelling tasks, but also by their extremely high computational and storage costs. Specifically, due to their massive size, even inference for large, highly-accurate GPT models may require multiple performant GPUs to execute, which limits the usability of such models. While there is emerging work on relieving this pressure via model compression, the applicability and performance of existing compression techniques is limited by the scale and complexity of GPT models. In this paper, we address this challenge, and propose GPTQ, a new one-shot weight quantization method based on approximate second-order information, that is both highly-accurate and highly-efficient. Specifically, GPTQ can quantize GPT models with 175 billion parameters in approximately four GPIL hours, exclusion the highly-limit the day to the parameters in approximately four GPIL hours, exclusion the highly-limit they are previous the previous the highly-limit proposed or the previous the highly-limits are previous the highly-limits and the proposed or the previous the highly-limits are previous the highly-limits and the proposed of the previous the highly-limits are previous the highly-limits and the proposed of the proposed or the proposed or the proposed of the proposed or the

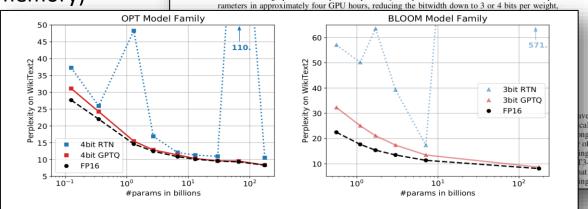


Figure 1: Quantizing OPT models to 4 and BLOOM models to 3 bit precision, comparing GPTQ with the FP16 baseline and round-to-nearest (RTN) [34, 5].

Table 6: 2-bit GPTQ quantization results with varying group-sizes; perplexity on WikiText2.





Quantization Reduces Data by an Order of Magnitude

10x

How to Go Further?



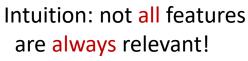




Model Sparsification ... (arXiv:2102.00554)



- For technical reasons (training, implementation etc.)
- We may want to shift towards sparse!

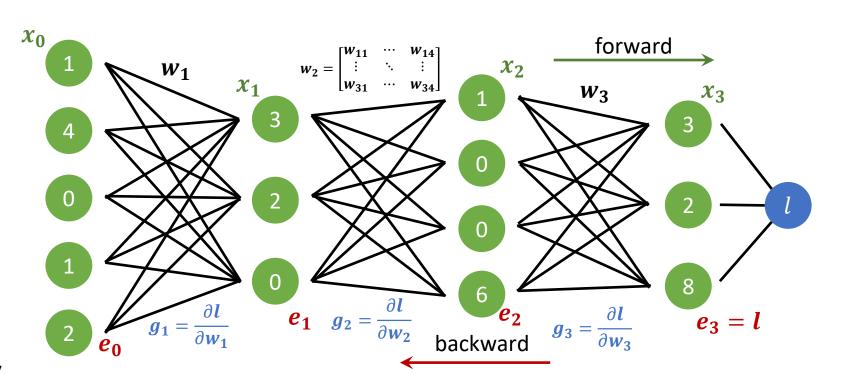


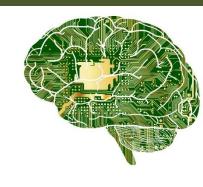
- Represent as (sparse) vector space
- ✓ Less overfitting
- ✓ Interpretability
- ✓ Parsimony

the f_t_re wi_l b_ sp_rs_

Key results:

- 95% sparse ResNet-52,
 BERT, or GPT models
- Essentially same quality
- Up to 20x cheaper!











a 4 V 005 2102

Much more (at least two hours more)

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

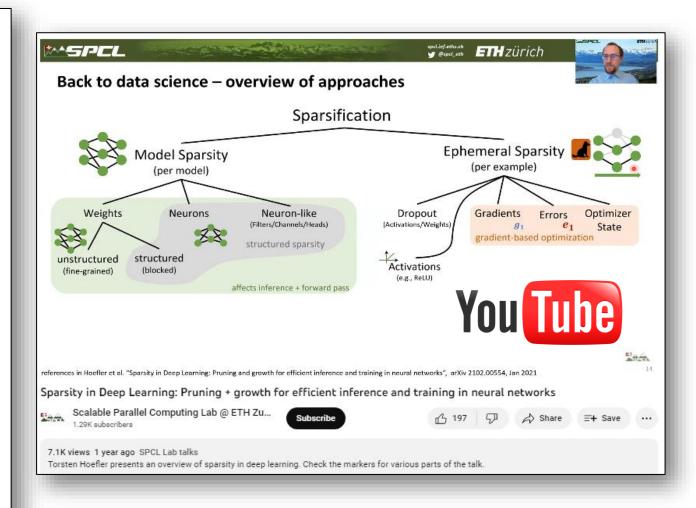
The growing energy and performance costs of deep learning have driven the community to reduce the size of neural networks by selectively pruning components. Similarly to their biological counterparts, sparse networks generalize just as well, if not better than, the original dense networks. Sparsity can 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 datum of experience -

Albert Einstein, 1933

1 INTRODUCTION

Deep learning shows unparalleled promise for solving very complex real-world problems in areas such as computer vision, natural language processing, knowledge representation, recommendation systems, drug discovery, and many more. With this development, the field of machine learning is moving from traditional feature engineering to neural architecture engineering. However, still





arXiv:2306.03078v1





17

The next step: Sparse-Quantized Representations - SpQR

SpQR: A Sparse-Quantized Representation for Near-Lossless LLM Weight Compression

Tim Dettmers*†
University of Washington

Ruslan Svirschevski* HSE University & Yandex Vage Egiazarian* HSE University & Yandex

Denis Kuznedelev* Yandex & Skoltech Elias Frantar IST Austria Saleh Ashkboos ETH Zurich Alexander Borzunov HSE University & Yandex

Torsten Hoefler ETH Zurich Dan Alistarh
IST Austria & NeuralMagic

Abstract

published at ICLR'24

Recent advances in large language model (LLM) pretraining have led to highquality LLMs with impressive abilities. By compressing such LLMs via quantization to 3-4 bits per parameter, they can fit into memory-limited devices such as laptops and mobile phones, enabling personalized use. However, quantization down to 3-4 bits per parameter usually leads to moderate-to-high accuracy losses, especially for smaller models in the 1-10B parameter range, which are well-suited for edge deployments. To address this accuracy issue, we introduce the Sparse-Quantized Representation (SpQR), a new compressed format and quantization technique which enables for the first time near-lossless compression of LLMs across model scales, while reaching similar compression levels to previous methods. SpQR works by identifying and isolating outlier weights, which cause particularlylarge quantization errors, and storing them in higher precision, while compressing all other weights to 3-4 bits, and achieves relative accuracy losses of less than 1% in perplexity for highly-accurate LLaMA and Falcon LLMs. This makes it possible to run 33B parameter LLM on a single 24 GB consumer GPU without any performance degradation at 15% speedup thus making powerful LLMs available to consumer without any downsides. SpQR comes with efficient algorithms for both encoding weights into its format, as well as decoding them efficiently at runtime³. Specifically, we provide an efficient GPU inference algorithm for SpQR which yields faster inference than 16-bit baselines at similar accuracy, while enabling memory compression gains of more than 4x.







Model Compression Enables

100x

More Efficient Processing

Which Makes Data Movement Even More Important!

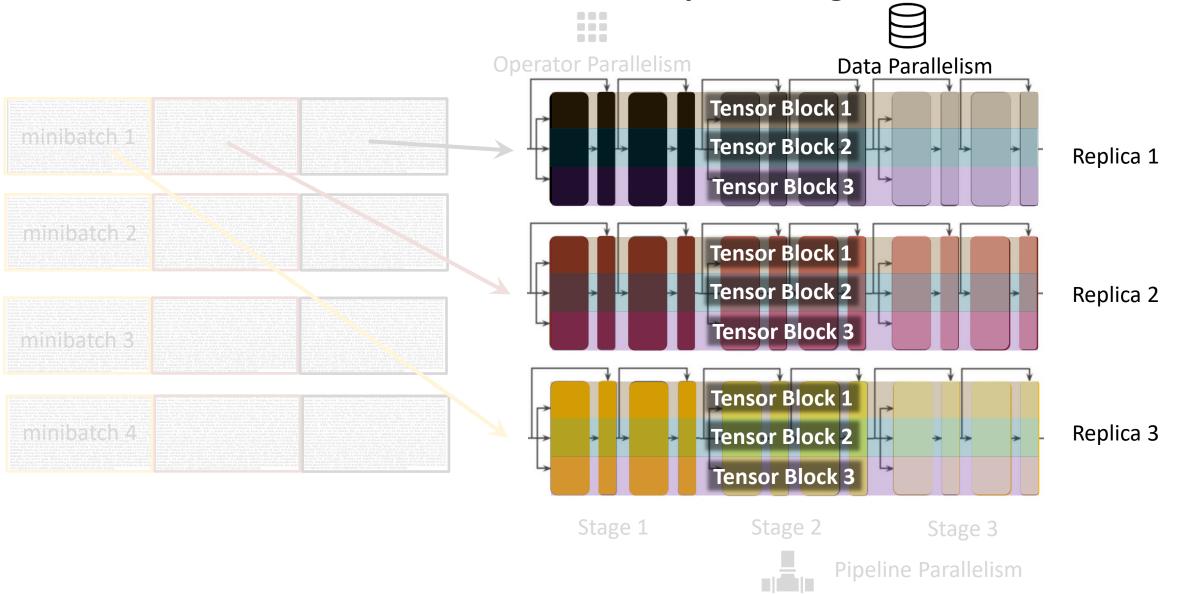
Especially in the Network!







The Three Dimensions of Parallelism in Deep Learning (arXiv:1802.09941)





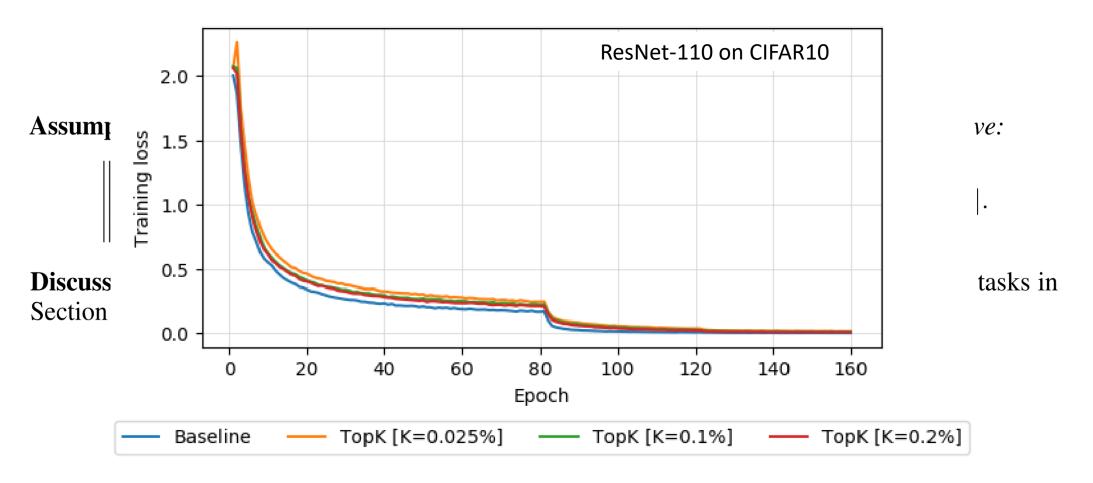




Data-parallel Gradient Sparsification — Top-k SGD (arXiv:1809.10505)



- Turns out 90-99.9% of the smallest gradient values can be skipped in the summation at similar accuracy
 - Accumulate the skipped values locally (convergence proof, similar to async. SGD with implicit staleness bounds [1])

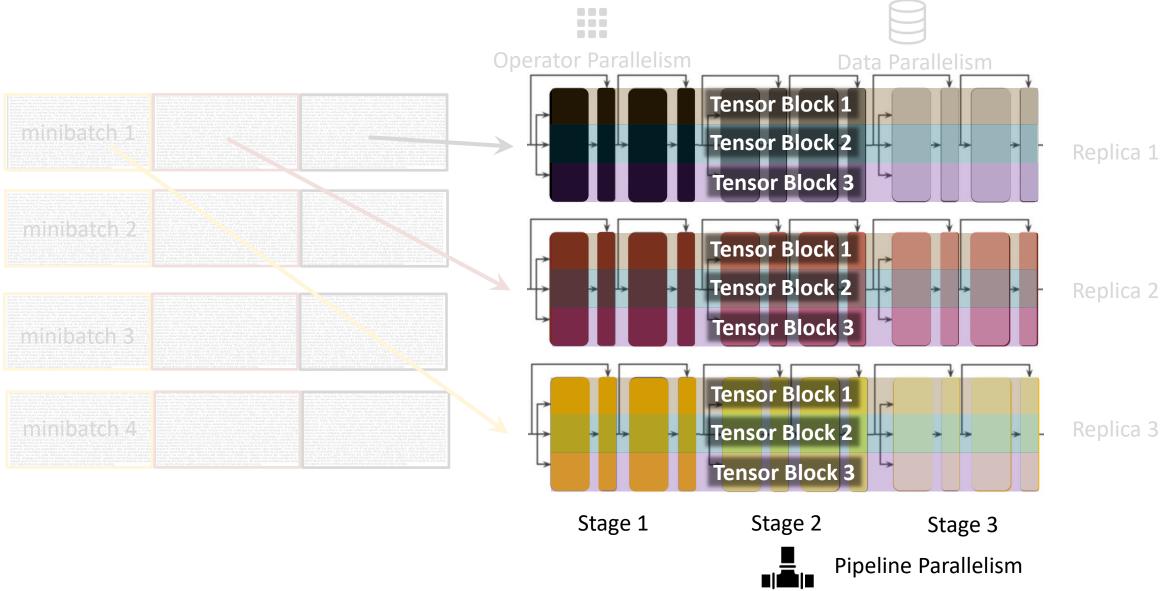








The Three Dimensions of Parallelism in Deep Learning (arXiv:1802.09941)





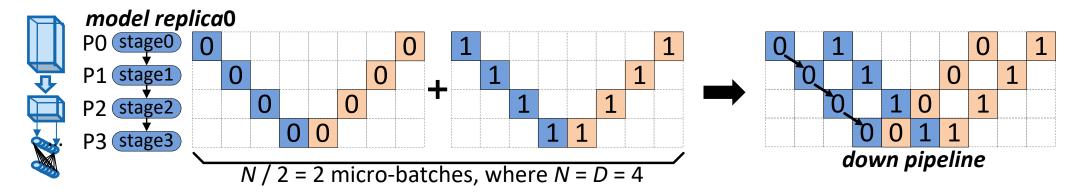


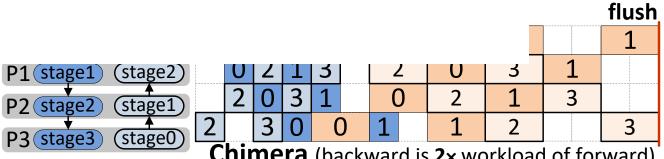




Bidirectional Pipelines – Meet Chimera (arXiv: 2107.06925v3)







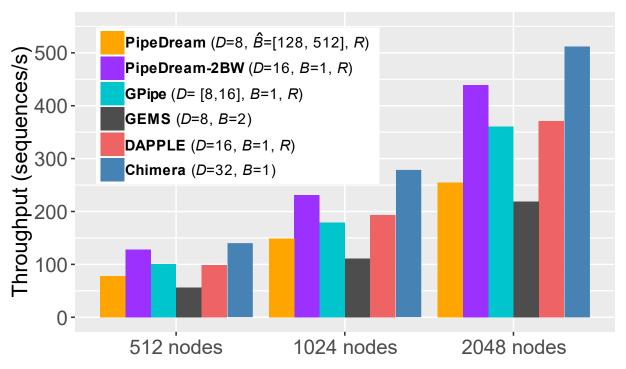






Chimera Weak Scaling (arXiv: 2107.06925v3)





Weak scaling for GPT-2 on Piz Daint (512 to 2048 GPU nodes)

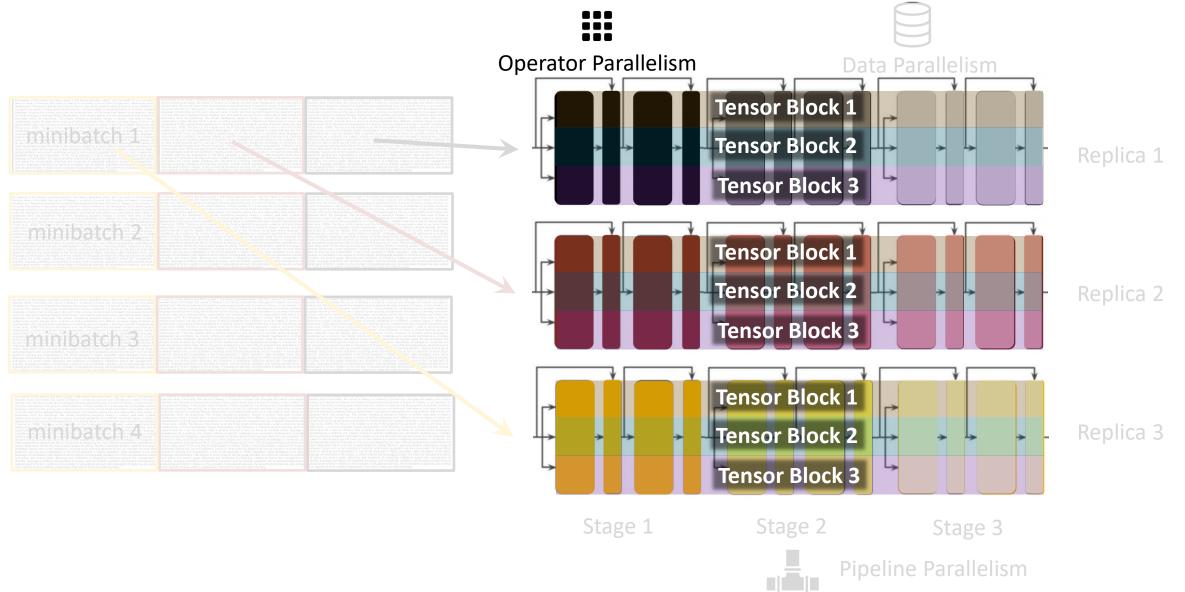
- 1.38x 2.34x speedup over synchronous approaches (GPipe, GEMS, DAPPLE)
 - Less bubbles
 - More balanced memory thus no recomputation
- 1.16x 2.01x speedup over asynchronous approaches (PipeDream-2BW, PipeDream)
 - More balanced memory thus no recomputation
 - Gradient accumulation thus low synch frequency







The Three Dimensions of Parallelism in Deep Learning (arXiv:1802.09941)





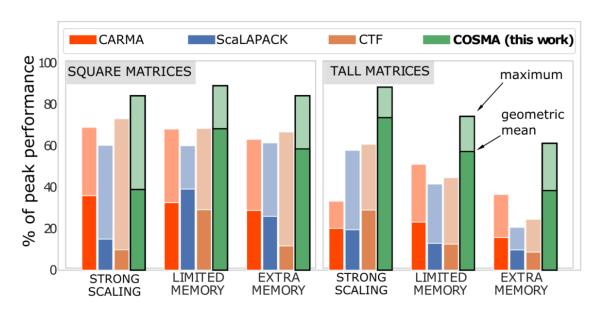
All MMM!



Operator Parallelism, i.e., Parallel Matrix Matrix Multiplication Remember those?

- Large MMMs dominate large language models!
 - e.g., GPT-3 multiples 12,288x12,288 matrices600 MiB in fp32 and 1.9 Tflop
 - generative inference multiplies tall & skinny matrices
- Distribute as operator parallelism
 - Heaviest communication dimension!
 Requires most optimization!
- COSMA [1] communication-optimal distributed MMM
 - Achieves tight I/O lower bound of $Q \ge \min \left\{ \frac{2mnk}{p\sqrt{S}} + S, 3\left(\frac{mnk}{p}\right)^{\frac{2}{3}} \right\}$
 - Uses partial replication with an outer-product schedule See paper for details and proofs!
- AutoDDL [2] combines operator-parallel models into communication-avoiding data distribution

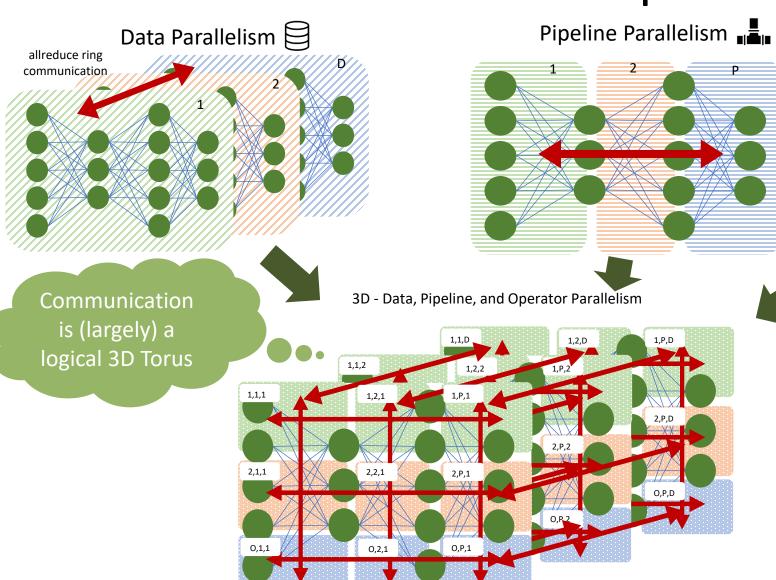
Operator class	% flop	% Runtime	
Tensor contraction	99.80	61.0	
Statistical normalization	0.17	25.5	
Element-wise	0.03	13.5	



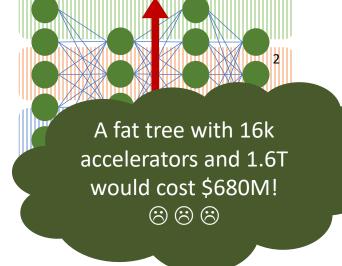




Communications in 3D Parallelism in Deep Learning (arXiv:2209.01346)



Operator Parallelism



AI bandwidth today / yesterday (and growing!)

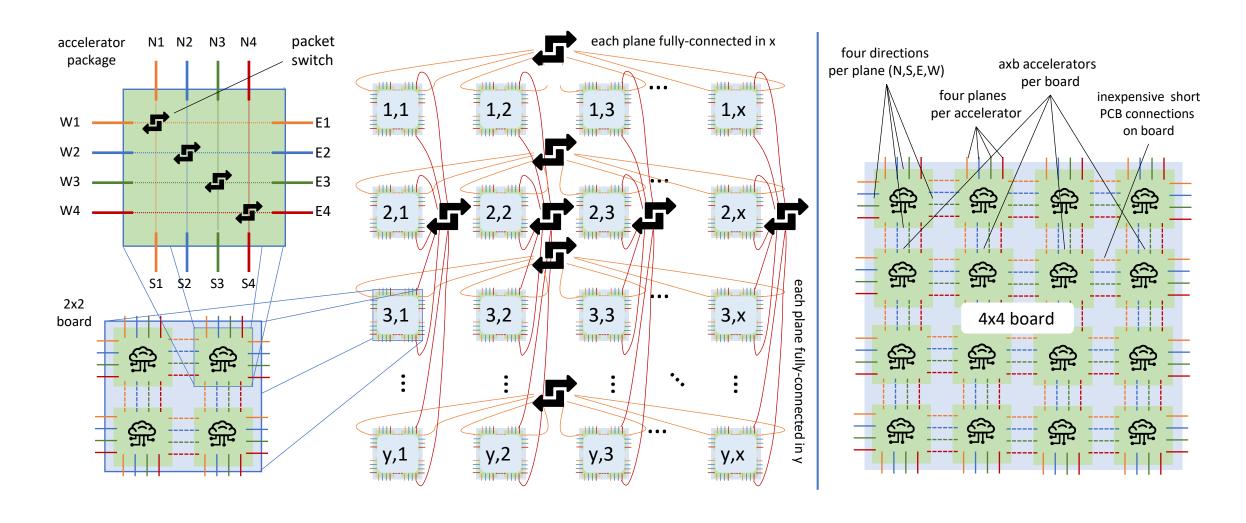
- Google TPUv2 ('21): 1T
- AWS Trainium ('21): 1.6T
- DGX-2 (A100, '21): 4.8T (islands of NVLINK)
- Tesla Dojo ('22): 128T
 - → Broadcom TH5 / NVIDIA Spectrum 4: 51.2T







Co-designing an AI Supercomputer with Unprecedented and Cheap Bandwidth





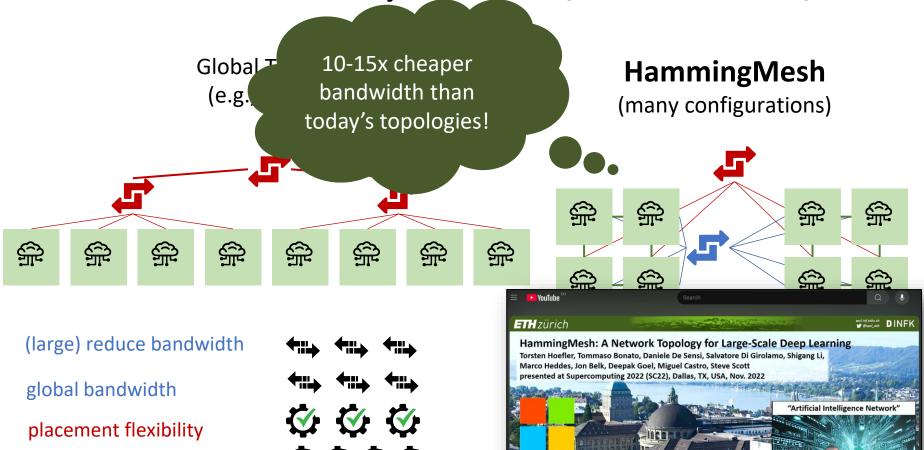
injection bandwidth

The Age of Computation

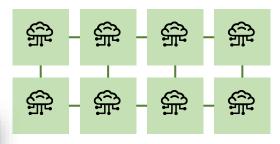


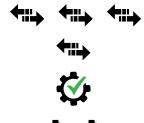


Bandwidth-cost-flexibility Tradeoffs (arXiv:2209.01346)



Local Topology (e.g., 2D Torus)



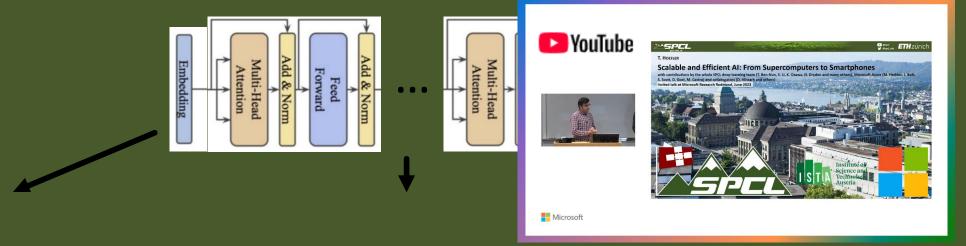








Three Systems Dimensions in Large-scale Super-learning ...





Altogether, we discussed a cost / performance improvement of

>1,000x

What now?

















Pre-training as we know it will unquestionably end...because we have but one internet

Let's teach them to reason!



[Lei et al., August'23]

Ilya Sutskever

"Let's proceed step by step" ⓒ

Basic Input-Output (IO)

Input

[Wang et al., March'22] https://github.com/princeton-nlp/tree-of-thought-llm

[Long, May'23]

https://github.com/jieyilong/tree-of-thought-puzzle-solver

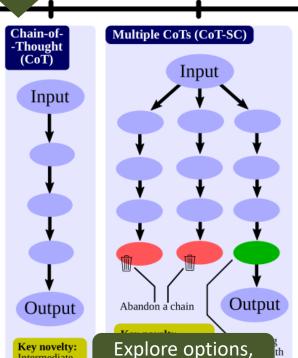
[Yao et al., May'23]

Sort the numbers "3, 2, 4, 5, 7, 12, 5, 6"

2017 - Transfc

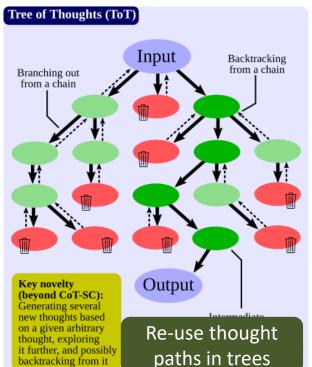
To sort "4, 6, 1, 8", I first split them into sets "4, 6" and "1, 8". Then I sort the sets and then I merge them sorted.

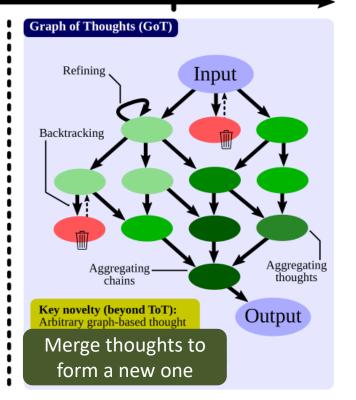
Sort the numbers "3, 2, 4, 5, 7, 12, 5, 6"



majority vote.

LLM thoughts











From LLMs to AHI

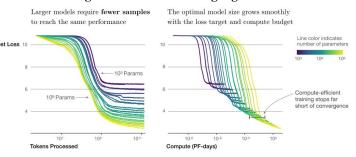


2018 - **BERT**

"BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding"

122k+ citations

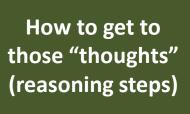
Scaling Laws for Neural Language Models



2020 - **GPT-3** (2020, scaling laws) "Language Models are Few-Shot Learners" 37k+ citations

LLaMA by Meta

2023 – **Llama** (Qwen, Grok, etc.) "LLaMA: Open and Efficient Foundation Language Models" 11k+ citations





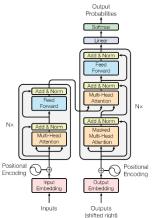
era of model size scaling

era of data scaling

2017 - Transformers

"Attention is All you Need"

146k+ citations



2019 - **GPT-2**

"Language Models are Unsupervised Multitask Learners"

14k+ citations



2022 – **ChatGPT** (RLHF, 2023, DPO)

"Training language models to follow instructions with human feedback"

14k+ citations

ChatGPT

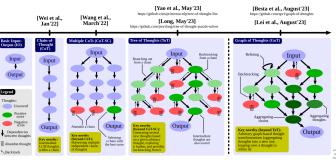
2023 – Chain of Thought

Reasoning (SC-CoT, ToT, GoT, etc.)

"Chain-of-Thought Prompting

Elicits Reasoning in Large Language Models"

8k+ citations

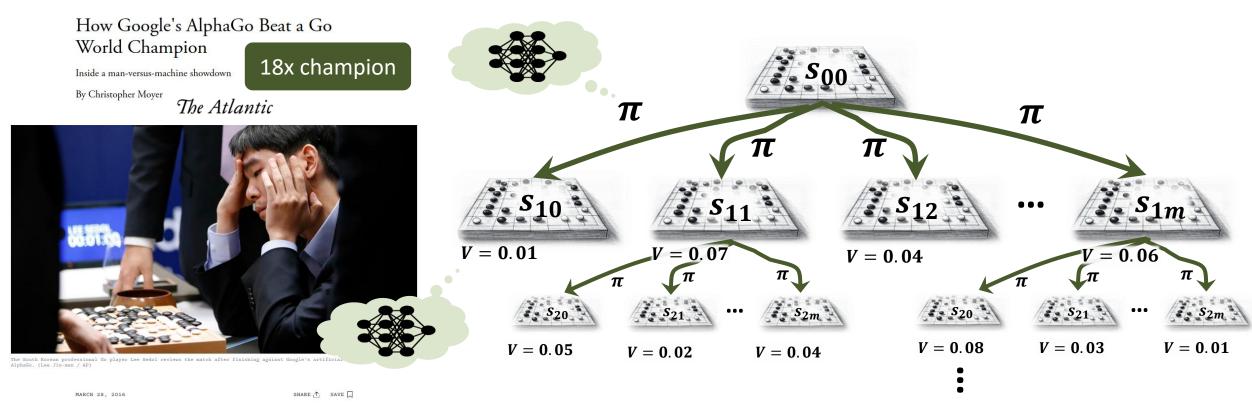








A Detour to Go Playing – AlphaGo vs. Lee Sedol (considered best Go player at the time)



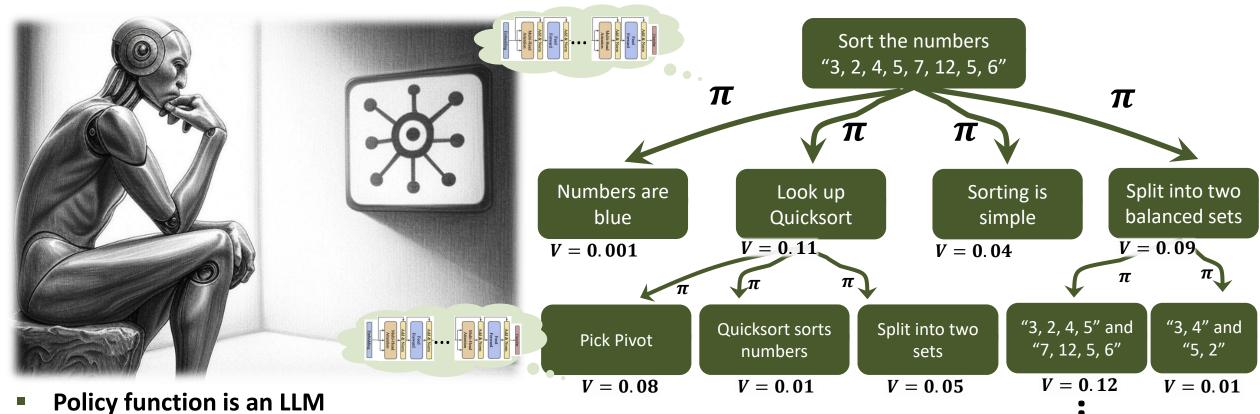
- (Monte Carlo) Tree Search (MCTS) samples multiple tree searches to some depth and propagates final values up the path, which keeps statistics for each state, action pair (edge)
 - Up to 1,600 expansions per move for AlphaGo Zero
 - Depth is decided by the value network (no fixed depth rollout)
- At the end, choose most promising action from root and prepare next move







Unifying LLMs and Reinforcement Learning into Large Reasoning Models (LRMs)



- - Fine-tuned with a special loss function to generate next best reasoning step (a bit tricky, needs multiple evals)
- Value function is another LLM
 - Replace final token output layer with a regression to a value (train on known examples, e.g., math tasks)
- During inference, still do MCTS search to cover reasoning paths
 - Extremely expensive! Up to thousands of inferences per reasoning step!



O3-preview

Gemini 1.5 Pro (002)

Claude 3.5 Sonnet

(2024-10-22)

o1-preview

o1-mini

GPT-40

(2024-08-06)

Grok 2 Beta





With RLMs to AHI

GPQA: A Graduate-Level Google-Proof **Q&A Benchmark**

Human PhDs:

34% outside their field 81% inside their field

03:

87% in all fields

We present GPQA, a challenging dataset of 448 multiple-choice questions written by domain experts in biology, physics, and chemistry. We ensure that the questions are high-quality and extremely difficult: experts who have or are pursuing PhDs in the corresponding domains reach 65% accuracy (74% when discounting clear mistakes the experts identified in retrospect), while highly skilled non-expert validators only

FRONTIERMATH: A BENCHMARK FOR EVALUATING ADVANCED MATHEMATICAL REASONING IN AI

25.2% Dec.'24



IMPRESSIONS OF OUR RESEARCH-LEVEL PROBLEMS



2024 – Strawberry **RL** (o1, o3, etc.) "Learning to Reason with LLMs"

Codeforce Elo rating

			O SERIES PERF	ORMANCE / ARC-AGI SEM	I-PRIVATE EVAL	
	100%		,	STEM GRAD	03 1	88% HIGH (TUNED) •
	75%		● AVG. MTURKER	76%		
S	50%	• KAGGLE	31% 32%	A har requ	ke"	
	25% · 7.	.80%	25% • 01 MED 01 LOW 3.33% • 01 PREVIEW		alization capab very few exam	
	0%	• 01-MINI \$1	.0 \$1	0.0 \$1 Cost Per Task	0.00 \$1,0	hir 000.0

			EIO-IVIIVIK
9/	100	ARC-AGI Semi-Private v1 Scores Over Time	3000+
ı		o3 tuned high (unreleased) o3 tuhed low	2700-2999
	80	(unreleased)	2400-2699
()	60		2200-2399
Score (%)		ol Pro	2000-2199
Sco	40	o1 high	1800-1999
		o1-pyeview	1600-1799
	20		1400-1599
	0	GPT-2 GPT-3 GPT-40	1200-1399
	9.01.0	paparat artain apparat apparat apparat	1000-1199
201		က် ကို	Up to 999

Nov.'24

	Elo-MMR	Title	Division	Number	Percentile	CF at same rank (spread)		
	3000+	Legendary Grandmaster	1	8	99.99	3382+		
	2700-2999	International Grandmaster	1	37	99.95	3010-3329 (372)		
	2400-2699	Grandmaster	1	255	99.7	2565-3010 (445)		
	2200-2399	International Master	1	560	99.1	2317-2565 (248)		
	2000-2199	Master	1	2089	97	2088-2317 (229)		
	1800-1999	Candidate Master	o3 achieves 2727 \rightarrow					
1600-1799 Expert								
99.95 th percentile of					le ot			

competitive

programmers!

Up to 818









Chollet: Calling something like o1 "an LLM" is about as accurate as calling AlphaGo "a convnet"

We are NOT done yet!



202

an

N

S

C

V





If you want to know more how this works or want to build one yourself!

Reasoning Language Models: A Blueprint

Maciej Besta^{1†}, Julia Barth¹, Eric Schreiber¹, Ales Kubicek¹, Afonso Catarino¹, Robert Gerstenberger¹, Piotr Nyczyk², Patrick Iff¹, Yueling Li³, Sam Houliston¹, Tomasz Sternal¹, Marcin Copik¹, Grzegorz Kwaśniewski¹, Jürgen Müller³, Łukasz Flis⁴, Hannes Eberhard¹, Hubert Niewiadomski², Torsten Hoefler¹

[†] Corresponding author ¹ETH Zurich ²Cledar ³BASF SE ⁴Cyfronet AGH

Abstract—Reasoning language models (RLMs), also known as Large Reasoning Models (LRMs), such as OpenAl's o1 and o3, DeepSeek-V3, and Alibaba's QwQ, have redefined Al's problem-solving capabilities by extending large language models (LLMs) with advanced reasoning mechanisms. Yet, their high costs, proprietary nature, and complex architectures—uniquely combining Reinforcement Learning (RL), search heuristics, and LLMs—present accessibility and scalability challenges. To address these, we propose a comprehensive blueprint that organizes RLM components into a modular framework, based on a survey and analysis of all RLM works. This blueprint incorporates diverse reasoning structures (chains, trees, graphs, and nested forms), reasoning strategies (e.g., Monte Carlo Tree Search, Beam Search), RL concepts (policy, value models and others), supervision schemes (Outcome-Based and Process-Based Supervision), and other related concepts (e.g., Test-Time Compute, Retrieval-Augmented Generation, agent tools). We also provide detailed mathematical formulations and algorithmic specifications to simplify RLM implementation. By showing how schemes like LLaMA-Berry, QwQ, Journey Learning, and Graph of Thoughts fit as special cases, we demonstrate the blueprint's versatility and unifying potential. To illustrate its utility, we introduce x1, a modular implementation for rapid RLM prototyping and experimentation. Using x1 and a literature review, we provide key insights, such as multi-phase training for policy and value models, and the importance of familiar training distributions. Finally, we discuss scalable RLM cloud deployments and we outline how RLMs can integrate with a broader LLM ecosystem. Our work demystifies RLM of Three Pillars of Reasoning Language Models (RLMs)

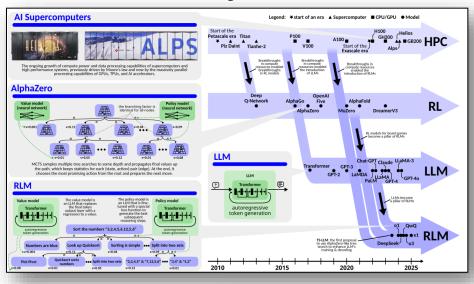
Index Terms—Reasoning Language Model, Large Reasoning Model, LRM, Reasoning LLMs, Reinforcement Learning for LLMs, MCTS for

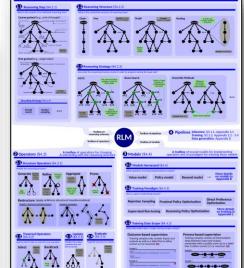
Introduction

Reasoning Language Models (RLMs), such as OpenAI's o1 [116], o3 [76], and Alibaba's OwO [148], also referred to as Large Reasoning Models (LRMs)¹, represent a transformative breakthrough in AI, on par with the advent of Chat-GPT [114]. These advanced systems have fundamentally redefined AI's problem-solving capabilities, enabling nuanced reasoning, improved contextual understanding, and robust decision-making across a wide array of domains, reshaping science [45], industries [21], governance [52], and numerous other aspects of human life [46], [75], [80], [143], [144]. By extending the capabilities of standard large language

fosters innovation, aiming to mitigate the gap between "rich AI" and "po Reasoning Language Models (RLMs) See §2.1.3 See §2.1.1 See §2.1.2 Pillar 1: Pillar 2: Reinforcement High-Performance Large Language Models (LLMs) Learning (RL) Computing (HPC) Examples: GPT-4o, LLaMA, Qwen. mples: AlphaZero, AlphaGo, eviously driven by Moore's law and now by the massively el processing capabilities of GPUs, TPUs, and AI accelerators, HPC is the foundation of LLMs, RL, and RLMs.

Hierarchy of Language Models Language Models (LMs) Large Language Models (LLMs) Reasoning Language Models (RLMs) See §2.1.1 Capable of System 1 Thinking: Capable of System 2 Thinking: Examples: GPT-40, LLaMA, Qwer xamples: o1, o3, DeepSeek-V3, QwQ Explicit RLMs (see §2.4.2) Implicit RLMs (see §2.4.1) Example: QwQ





growing disadvantage, threatening to stifle innovation and reinforce systemic inequities. As RLMs become integral to



The Age of Computation





With RLMs to AHI

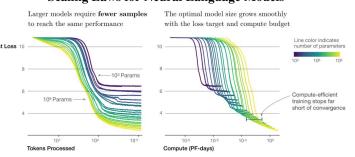


2018 - **BERT**

"BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding"

122k+ citations

Scaling Laws for Neural Language Models



2020 - **GPT-3** (2020, scaling laws)

"Language Models are
Few-Shot Learners"

37k+ citations

LLaMA by Meta

2023 – **Llama** (Qwen, Grok, etc.) "LLaMA: Open and Efficient Foundation Language Models" 11k+ citations



2024 – **Strawberry RL** (o1, o3, etc.) "Learning to Reason

_ with LLMs"

era of model size scaling

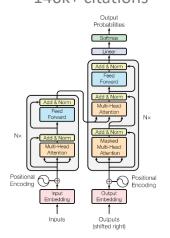
era of data scaling

era of reasoning scaling

2017 - Transformers

"Attention is All you Need"

146k+ citations



2019 - **GPT-2**

"Language Models are Unsupervised Multitask Learners"

14k+ citations



2022 – **ChatGPT** (RLHF, 2023, DPO)

"Training language models to follow instructions with human feedback"

14k+ citations

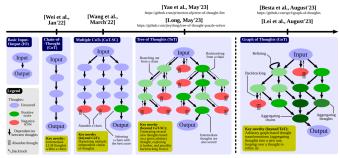
S ChatGPT

2023 – Chain of Thought

Reasoning (SC-CoT, ToT, GoT, etc.)

"Chain-of-Thought Prompting
Elicits Reasoning in Large Language Models"

8k+ citations



The Age of Computation







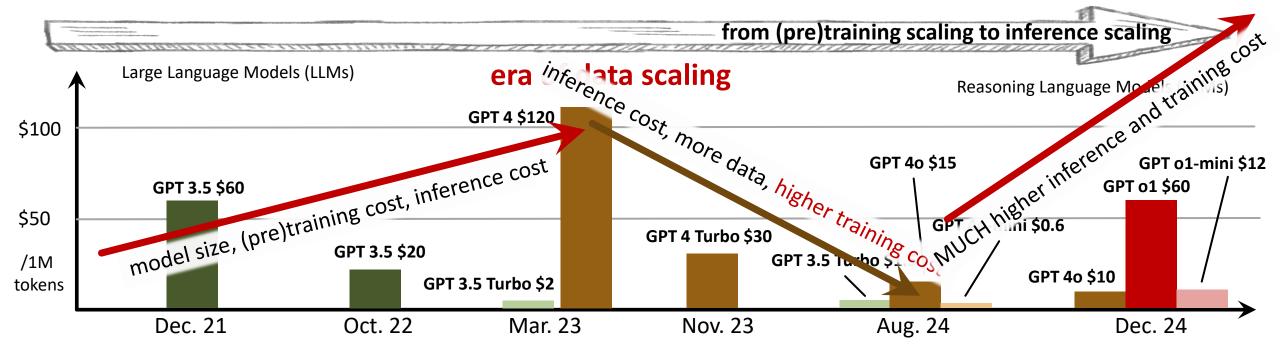
Development of Computation Requirement with RLMs

o3 > \$5000 / task

era of model size scaling

Efficient Language Models (ELMs)

era of reasoning scaling



We need cheaper compute



We need cheaper systems (networking!)



Principles: high local bandwidth, reliability, cost





Networks Converge

The Datacenter will be a Supercomputer

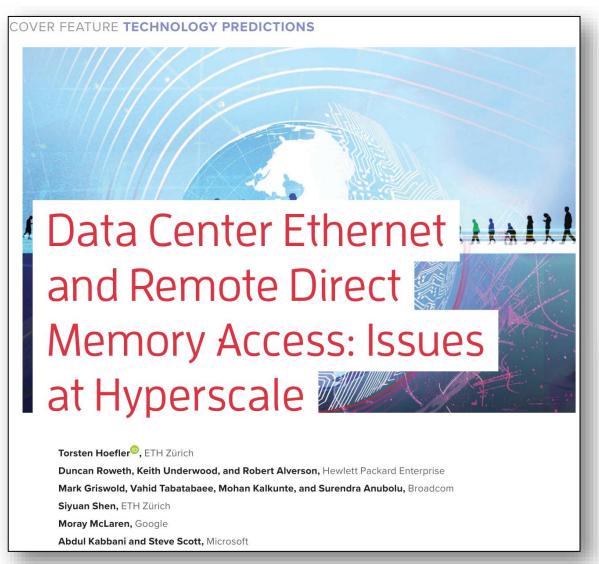








Ultra Ethernet Set Out to Create the Best AI/ML and HPC Interconnect!





Founding Members





















white Paper on <u>ultraethernet.org</u>

Overview of and Motivation for the Forthcoming Ultra Ethernet Consortium Specification

Networking Demands of Modern Al Jobs

Networking is increasingly important for efficient and cost-effective training of AI models. Large Language Models (LLMs) such as GPT-3, Chinchilla, and PALM, as well as recommendation systems like DLRM and DHEN, are trained on clusters of thousands of GPUs.







Ecosystem is quicky growing



Today 10 steering companies, 26 general member companies, 54 contributor members



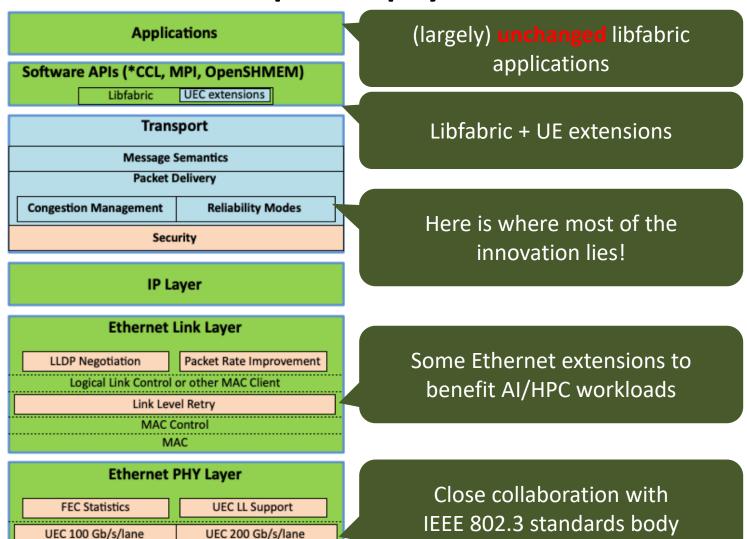
Chair's view of the Transport WG Meeting in March'24 (60+ members on site, 800+ total now)







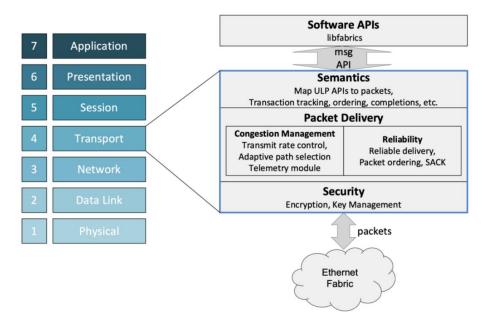
Ultra Ethernet's philosophy





UE enables cheap high-performance hardware implementations of an optimized transport over (legacy)

Ethernet networks while enabling vendor innovation



PMA





Congestion Management



Reliability Modes

Ultra Ethernet's key features compared RoCE and others

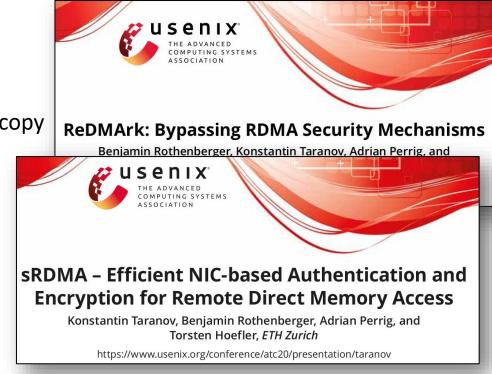


Transport

Message Semantics
Packet Delivery

Security

- Lossy (best effort) & lossless operation
 - Solves all PFC/blocking issues!
- Flexible ordering in packet and message delivery
 - Reliable Ordered Delivery (ROD), Reliable Unordered Delivery (RUD/RUDI), Unreliable Unordered Delivery (UUD)
- "State of the art" (2024), easily configured congestion control mechanisms
 - Sender and Receiver-based mechanisms over lossy networks
 - Supports trimming extensions
- Multi-path packet spraying
 - Adaptive routing with ordering using existing switches (ECMP), zero copy
- Switch offload (i.e., In-Network Collectives)
 - Cheap & effective
- Security as a first-class citizen co-designed with the transport
 - Addressing issues in RoCE
- Ethernet Link and Physical layer enhancements (optional)
 - See previous slide





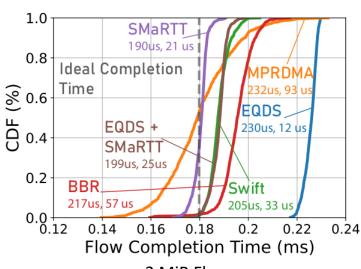


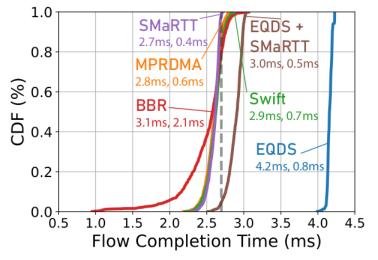


SMaRTT-REPS enables Modern Packet Spraying

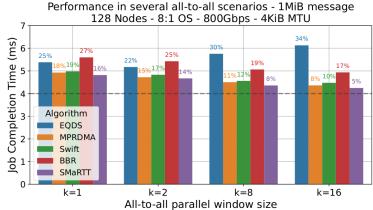


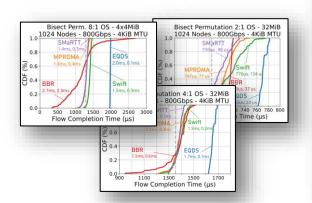
"State of the art" (2024), easily configured congestion control mechanisms





2 MiB Flows
Permutation traffic on 8:1 oversubscribed fat tree









37 lines simple pseudo-code

SMaRTT-REPS: Sender-based Marked Rapidly-adapting Trimmed & Timed Transport with Recycled Entropies

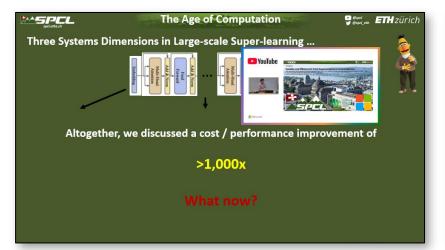
Tommaso Bonato ETH Zürich Microsoft	Abdul Kabl Microsoft		:
Rong Pan ^{AMD}	Yanfang AMD	Le Costin Raiciu Broadcom Inc.	
Mark Handley Broadcom Inc.	Timo Schne ETH Züric		
Ahmad Ghalayini ^{Microsoft}	Daniel Al		
Adrian C Micro		Torsten Hoefler ETH Zürich Microsoft	

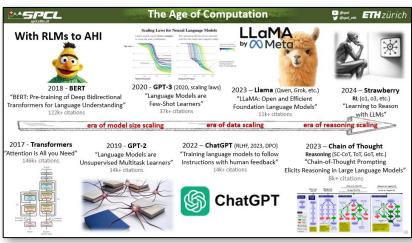


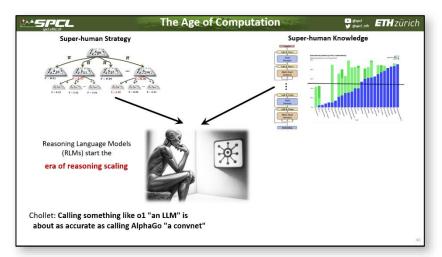


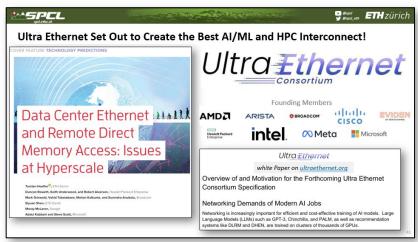


Key Points and Conclusions









More of SPCL's research:



... or <u>spcl.ethz.ch</u>



Want to join our efforts?
We're looking for excellent
Postdocs, PhD students, and Visitors.
Talk to me!







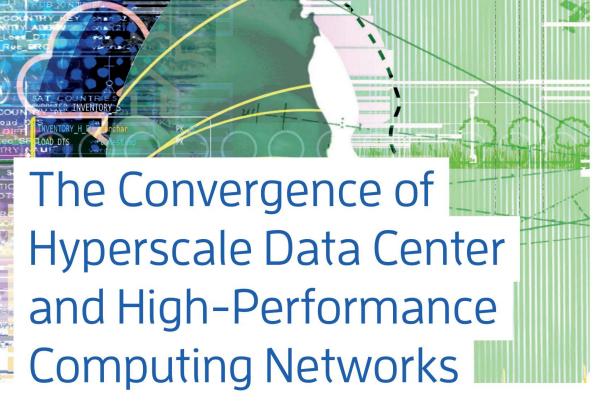








COVER FEATURE TECHNOLOGY PREDICTIONS

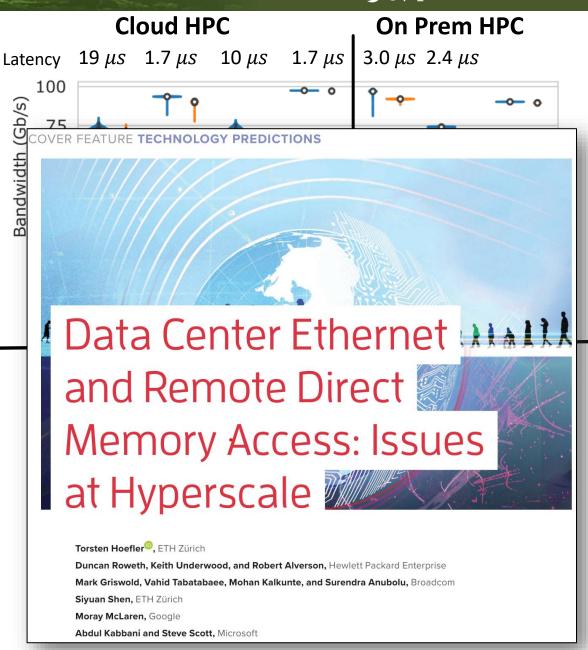


Torsten Hoefler, ETH Zurich

Ariel Hendel, Scala Computing

Duncan Roweth, Hewlett Packard Enterprise

We discuss the differences and commonalities between network technologies used in supercomputers and data centers and outline a path to convergence at multiple layers. We predict that emerging smart networking solutions will accelerate that convergence.







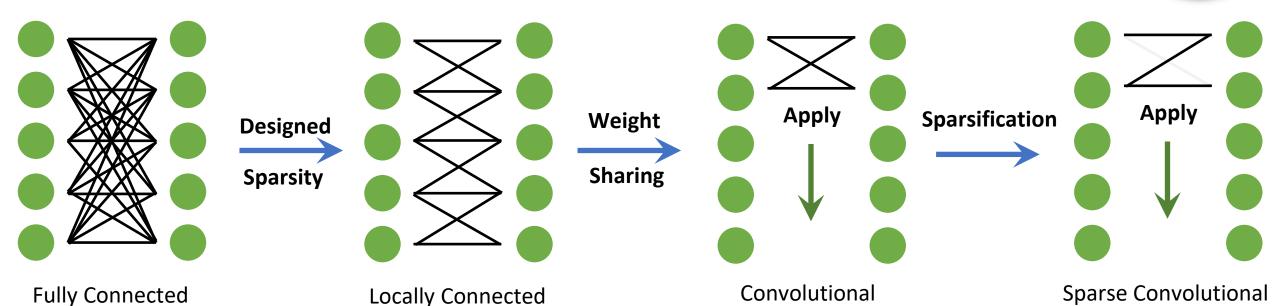


An example to relate to – CNNs from a sparsity viewpoint

"With all things being equal, the simplest explanation tends to be the right one"

- William of Ockham, ~1300





universal approximation hard to train

inductive locality bias reduce training complexity

inductive translational equivariance bias image recognition/object detection

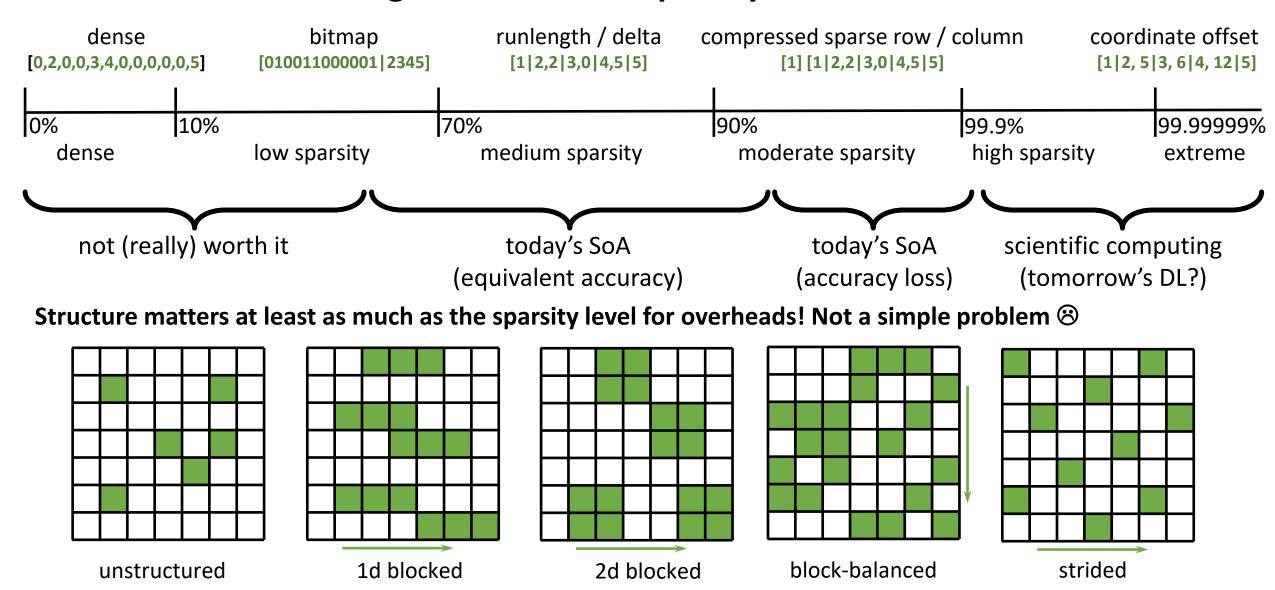
sparsification reduce representational complexity (MDL, Occam)







Performance and storage overheads of sparsity

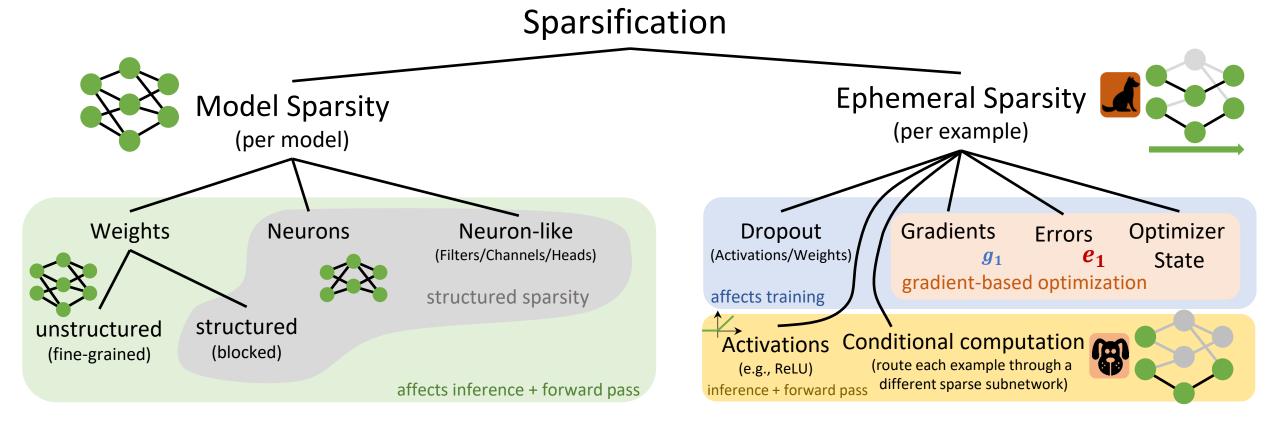








Back to data science – overview of approaches



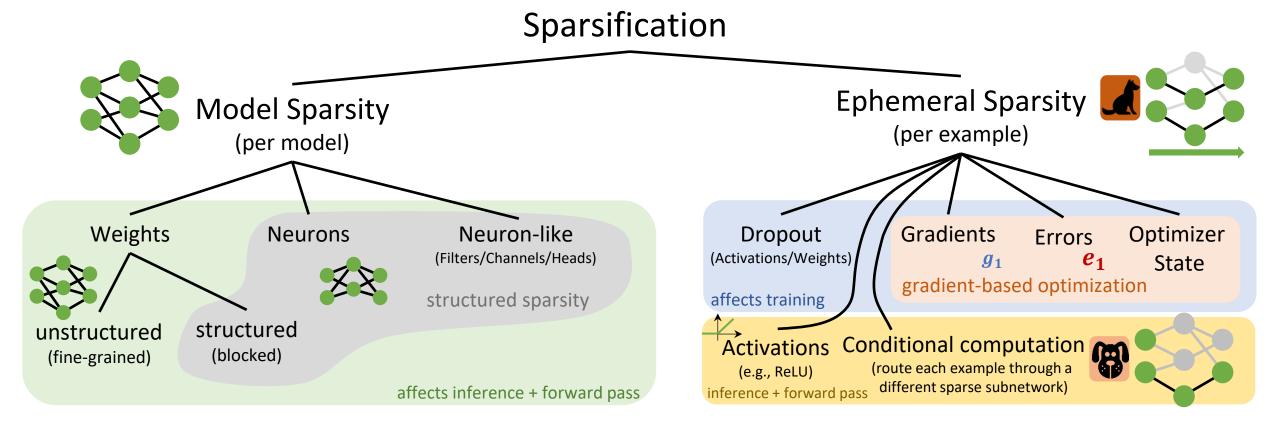
Quite complex, isn't it? It'll get better ©







Back to data science – overview of approaches



Quite complex, isn't it? It'll get better ©

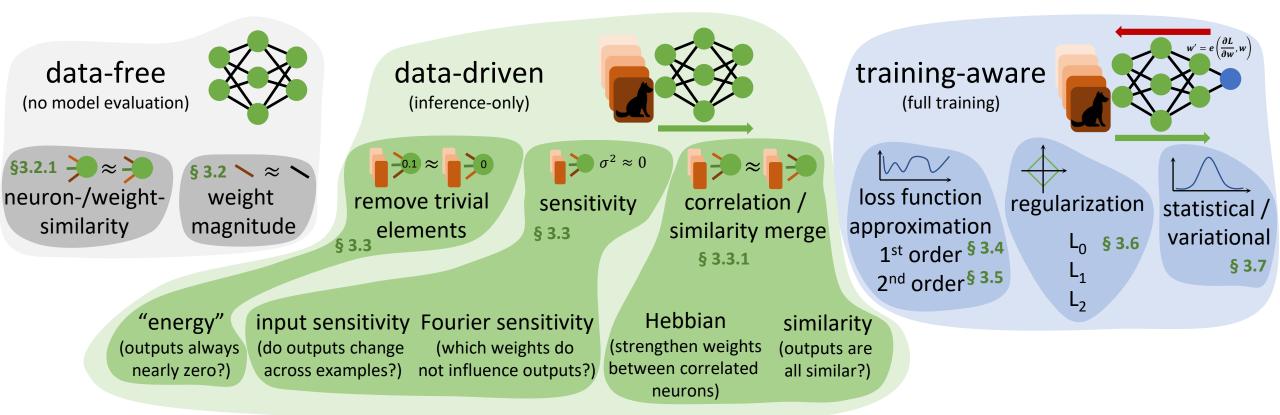






What, when, now the how to sparsify / remove elements!

- Simplest scheme: leave k out train $\binom{n}{k}$ models to convergence $\mathfrak S$
 - Various selection schemes by some importance metric
 Whirlwind overview of various metrics and selection techniques then focus on some





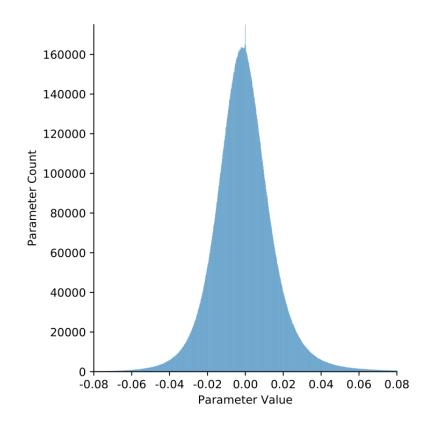


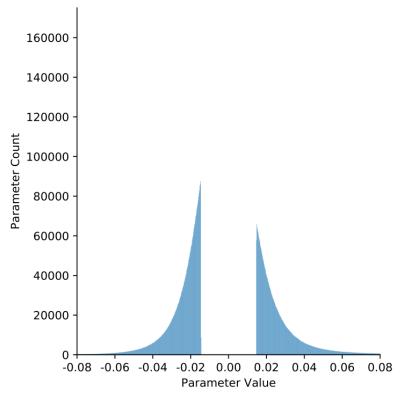


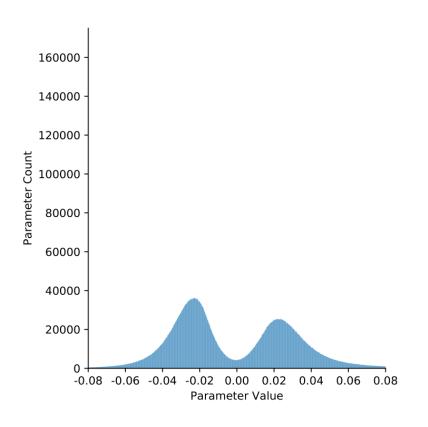
[Data free] Magnitude-based pruning

- Remove weights with smallest absolute magnitude |w|
 - Most popular and simplest either by absolute value or select top-k e.g., ResNet-50









(a) Dense Network (76.0%)

(b) 70% Pruned (36.1%)

(c) After 3-epoch Retraining (71.4%)