

Performance Modeling and System Design Insights for Scientific AI Foundation Models

Shashank Subramanian
NERSC, Berkeley Lab



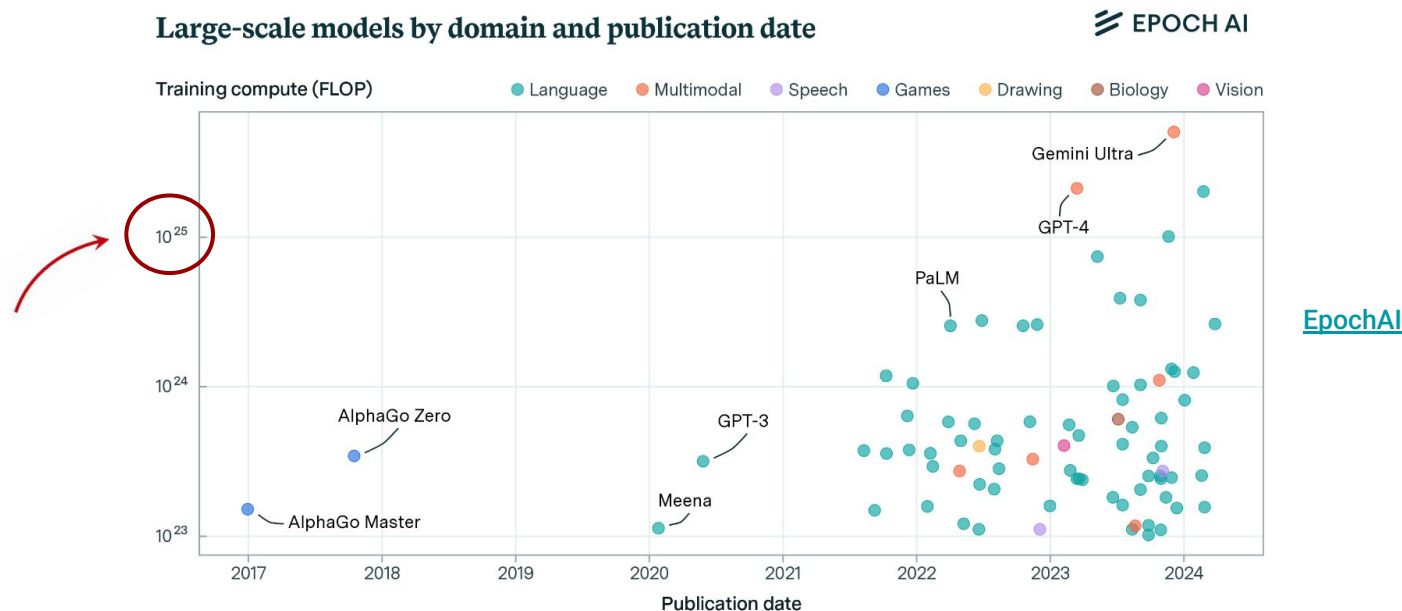
U.S. DEPARTMENT OF
ENERGY

Office of
Science



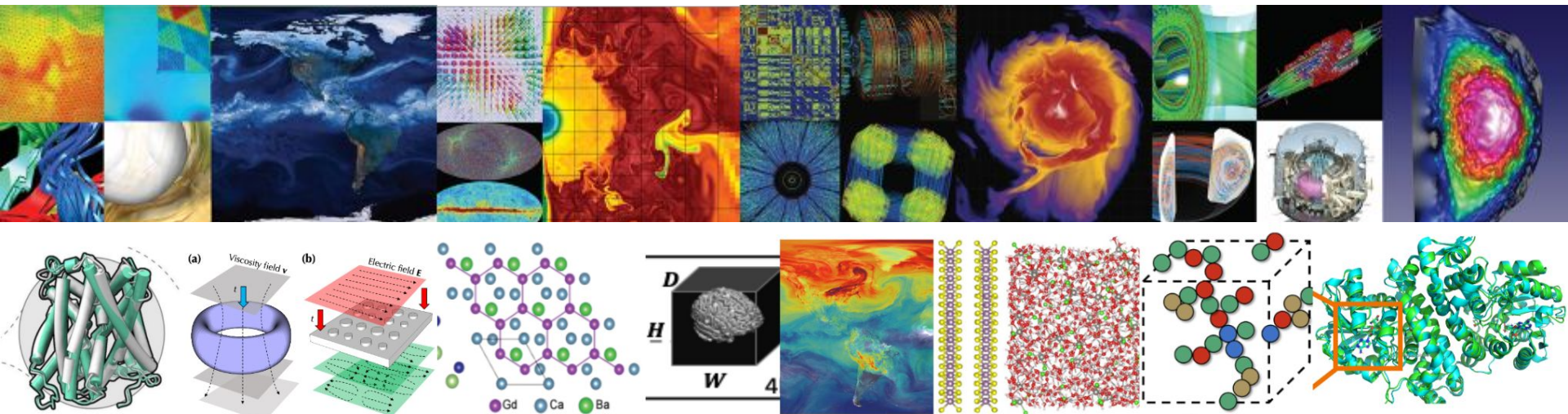
Comprehensive Performance Modeling and System Design Insights for Foundation Models, PMBS,
[SC24](#), [Github](#)

AI Foundation Models are **Expensive**



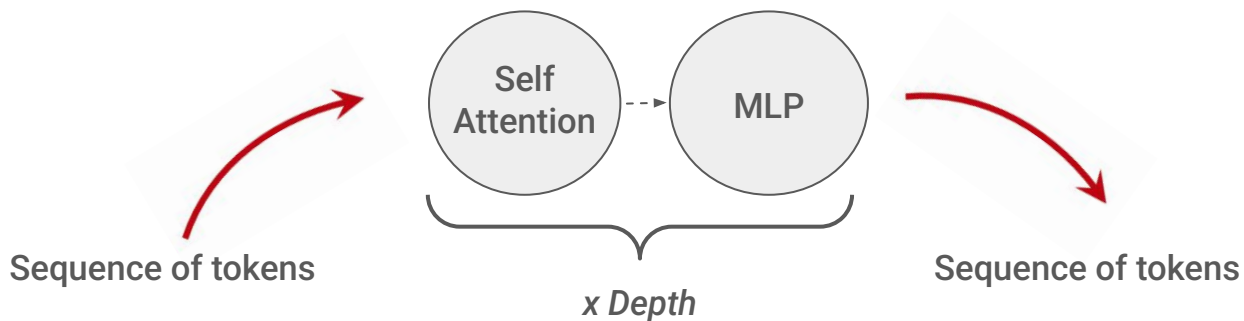
- Transformers are the workhorse: Scaling properties, flexible, SOTA results

Large-scale AI Models are Growing in Science

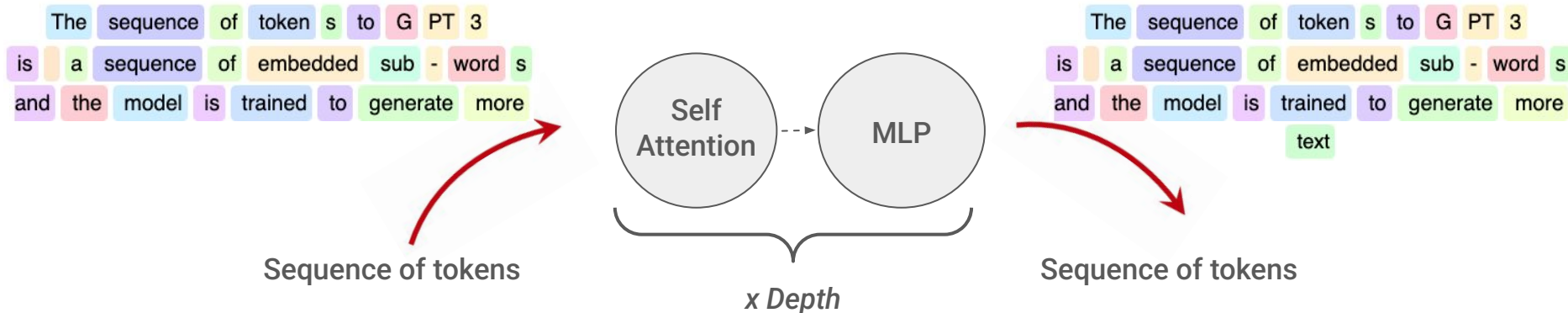


- Range of scientific simulation tasks is enormous
 - Weather/climate, fusion, seismic, fluids, proteins, material sciences, high-energy physics
- Surge of transformer models as possible *foundations* for downstream tasks

Transformers in Science may **Operate in Different Regimes**

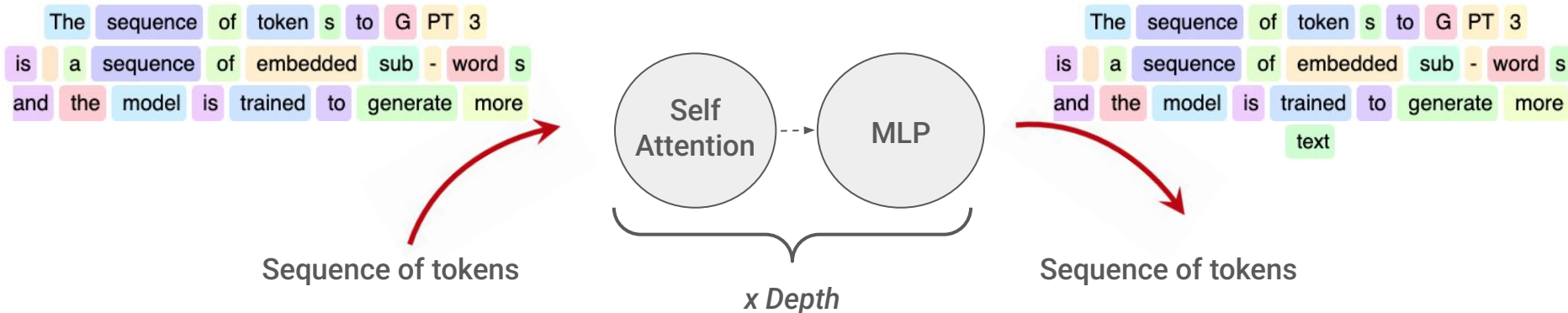


Transformers in Science may Operate in Different Regimes



- A Large Language Model (LLM) example: GPT3

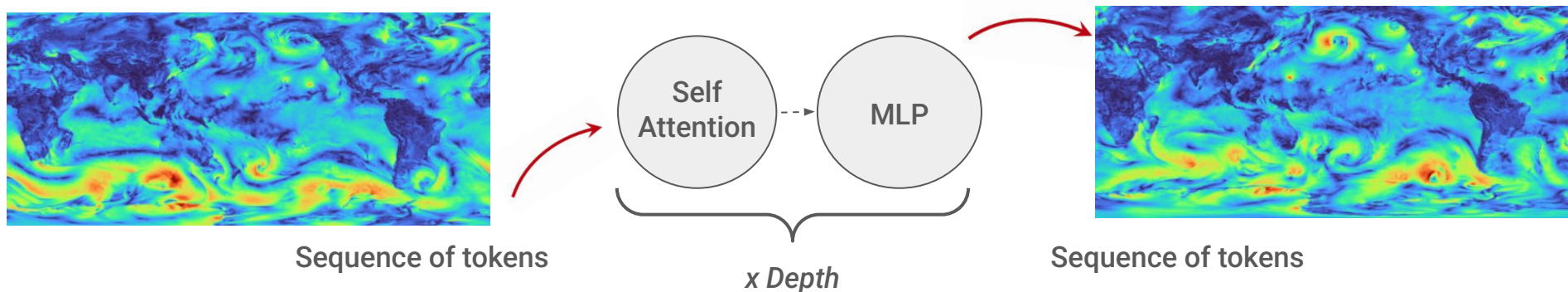
Transformers in Science may Operate in Different Regimes



- A Large Language Model (LLM) example: GPT3

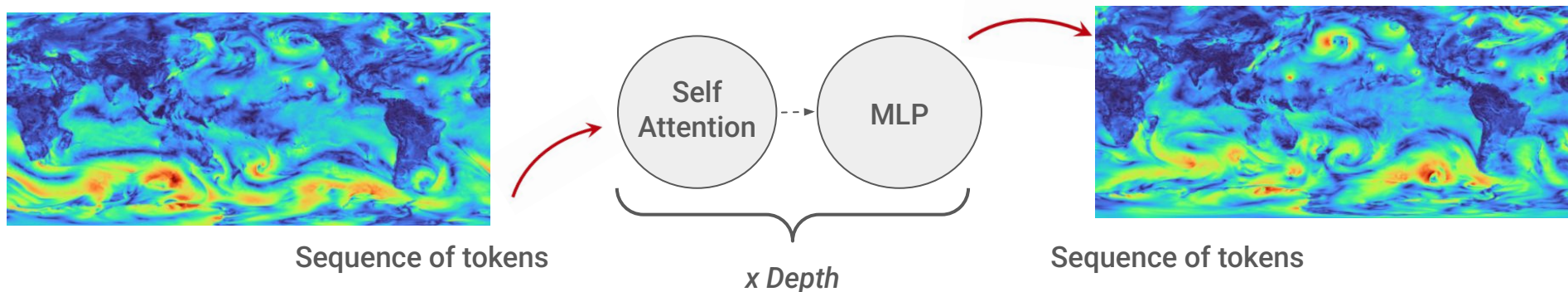
- #Parameters can be huge ~ **billions to trillions** of parameters
- Process a sequence of $O(1K)$ tokens (usually **2K, 4K, 8K** tokens in pre-training)
- MLP FLOPs are large (compared to S/A)
- GPT3-1T on **3072 A100 GPUs** takes **84 days** to train on 450B tokens
- Understood reasonably well

Transformers in Science may **Operate in Different Regimes**



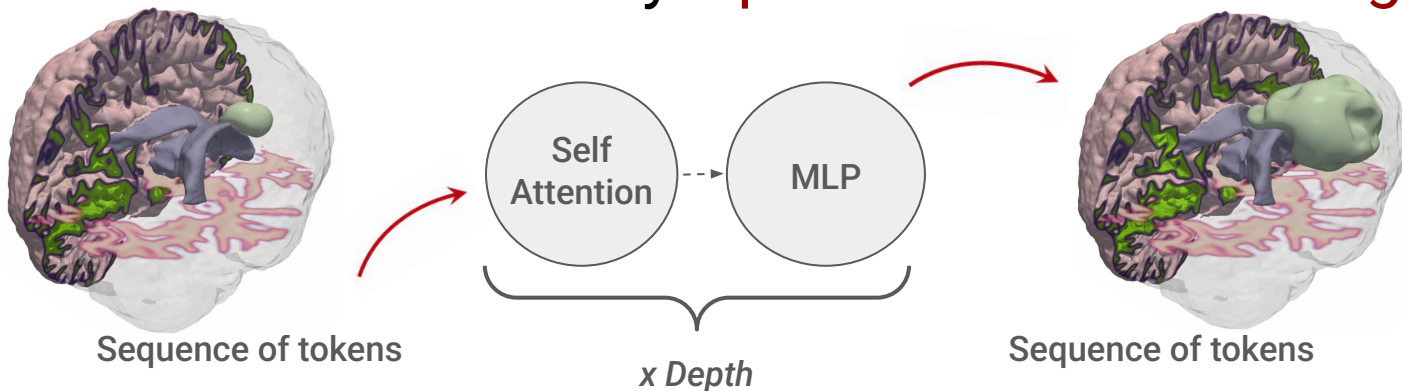
- A Scientific Surrogate example: Transformer for global weather forecasting

Transformers in Science may **Operate in Different Regimes**



- A Scientific Surrogate example: Transformer for global weather forecasting
 - #Parameters are moderate ~ **million to billion** parameters
 - Process a sequence of **O(1M) tokens** (can be compressed to O(100K) tokens)
 - S/A FLOPs are large (compared to MLP)
 - **A small model could be more expensive than a trillion parameter LLM!**
 - [?] Days on [?] GPUs on [?] tokens. Less understood

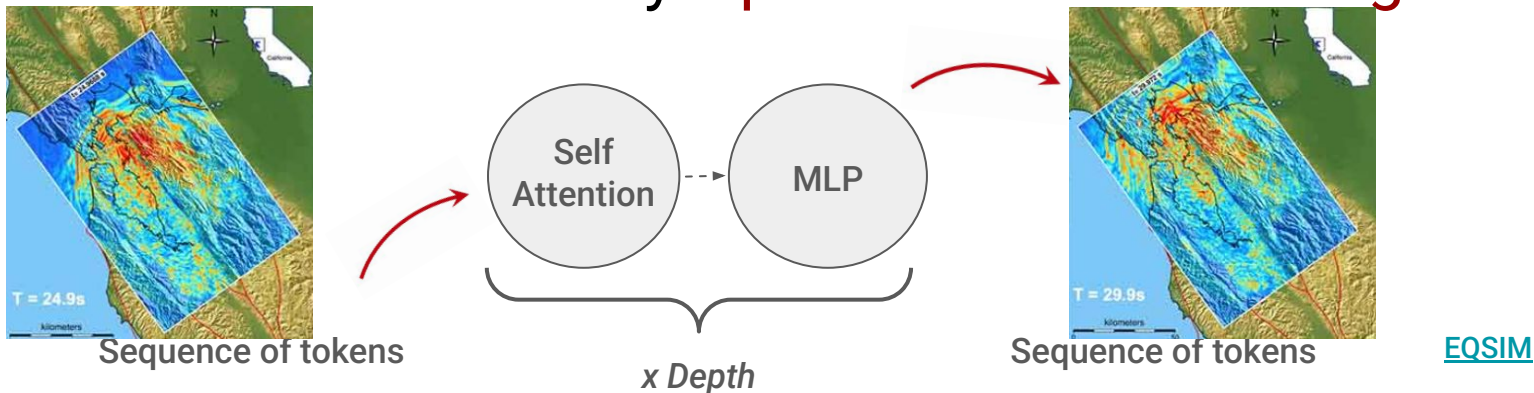
Transformers in Science may **Operate in Different Regimes**



- **A Scientific Surrogate example:**

- #Parameters are moderate ~ **million to billion** parameters
- Process a sequence of **$O(1M)$ tokens**
- S/A FLOPs are large (compared to MLP)
- **A small model could be more expensive than a trillion parameter LLM!**
- [?] Days on [?] GPUs on [?] tokens. Less understood

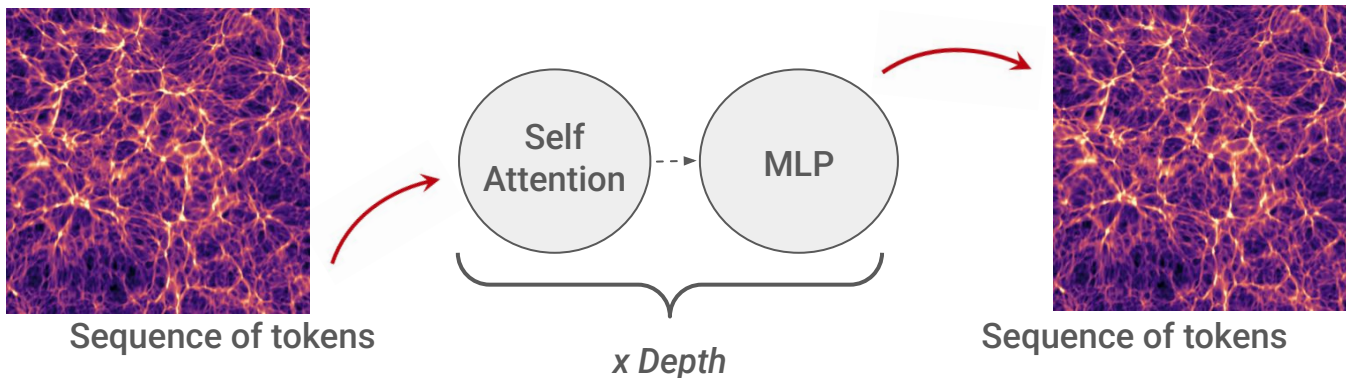
Transformers in Science may **Operate in Different Regimes**



- A Scientific Surrogate example:

- #Parameters are moderate ~ **million to billion** parameters
- Process a sequence of **O(billions) tokens**
- S/A FLOPs are large (compared to MLP)
- **A small model could be more expensive than a trillion parameter LLM!**
- [?] Days on [?] GPUs on [?] tokens. Less understood

Transformers in Science may **Operate in Different Regimes**



[ACCEL2](#)

- **A Scientific Surrogate example:**

- #Parameters are moderate ~ **million to billion** parameters
- Process a sequence of **O(billions) tokens**
- S/A FLOPs are large (compared to MLP)
- **A small model could be more expensive than a trillion parameter LLM!**
- [?] Days on [?] GPUs on [?] tokens. Less understood

Performance Modeling can be Valuable

- Understand **Costs/Bottlenecks** and analyze **Sensitivity of Performance**
 - What bottlenecks w.r.t parallelization strategies?
 - Different Transformer regimes (LLMs vs Science)?
 - Different system hardware (specifically network/NVLINK effects)?
 - Different system scales (10s vs 1000s of accelerators)?

Performance Modeling can be Valuable

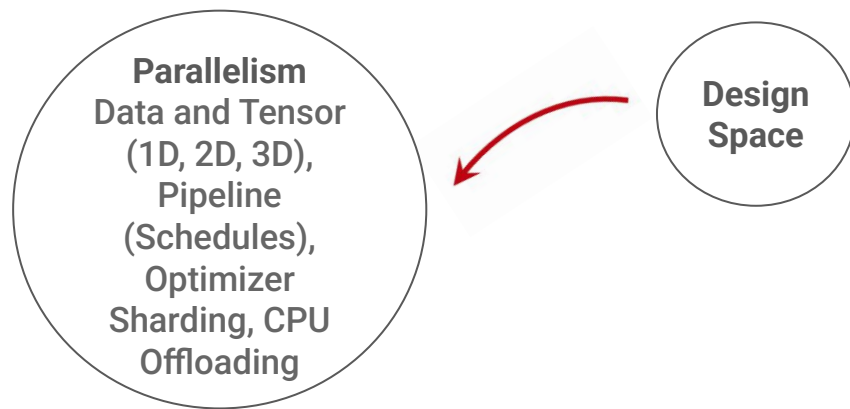
- Understand **Costs/Bottlenecks** and analyze **Sensitivity of Performance**
 - What bottlenecks w.r.t parallelization strategies?
 - Different Transformer regimes (LLMs vs Science)?
 - Different system hardware (specifically network/NVLINK effects)?
 - Different system scales (10s vs 1000s of accelerators)?
- **Value-add** for:
 - Users (researchers, engineers)
 - Optimal ways to parallelize AI models? Architecture search with performance in mind?
 - Systems design
 - Which aspects of the HPC system are crucial? Alternate design choices?

AI Performance Modeling is Challenging

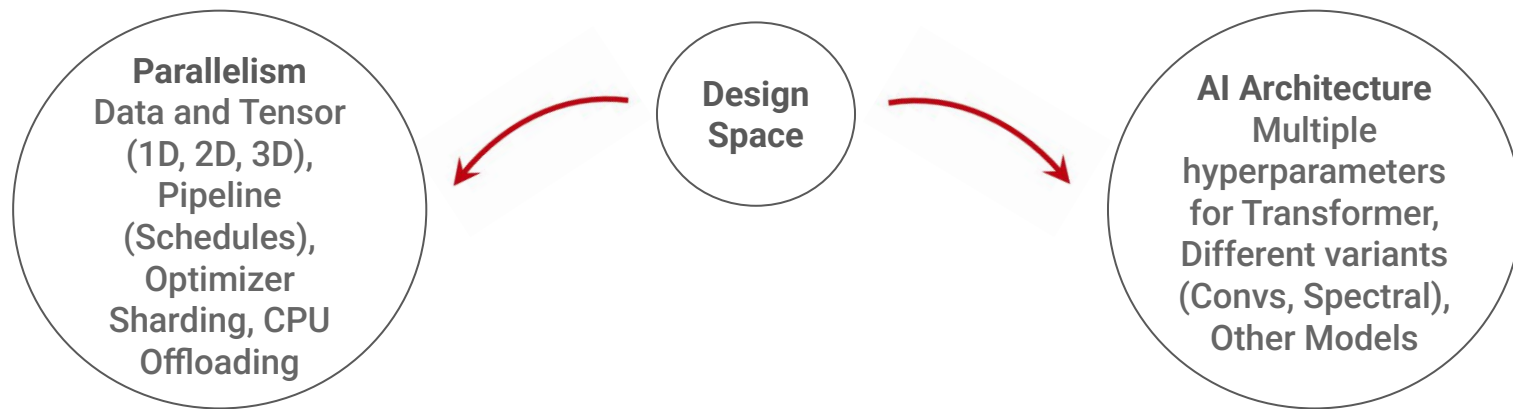


Design
Space

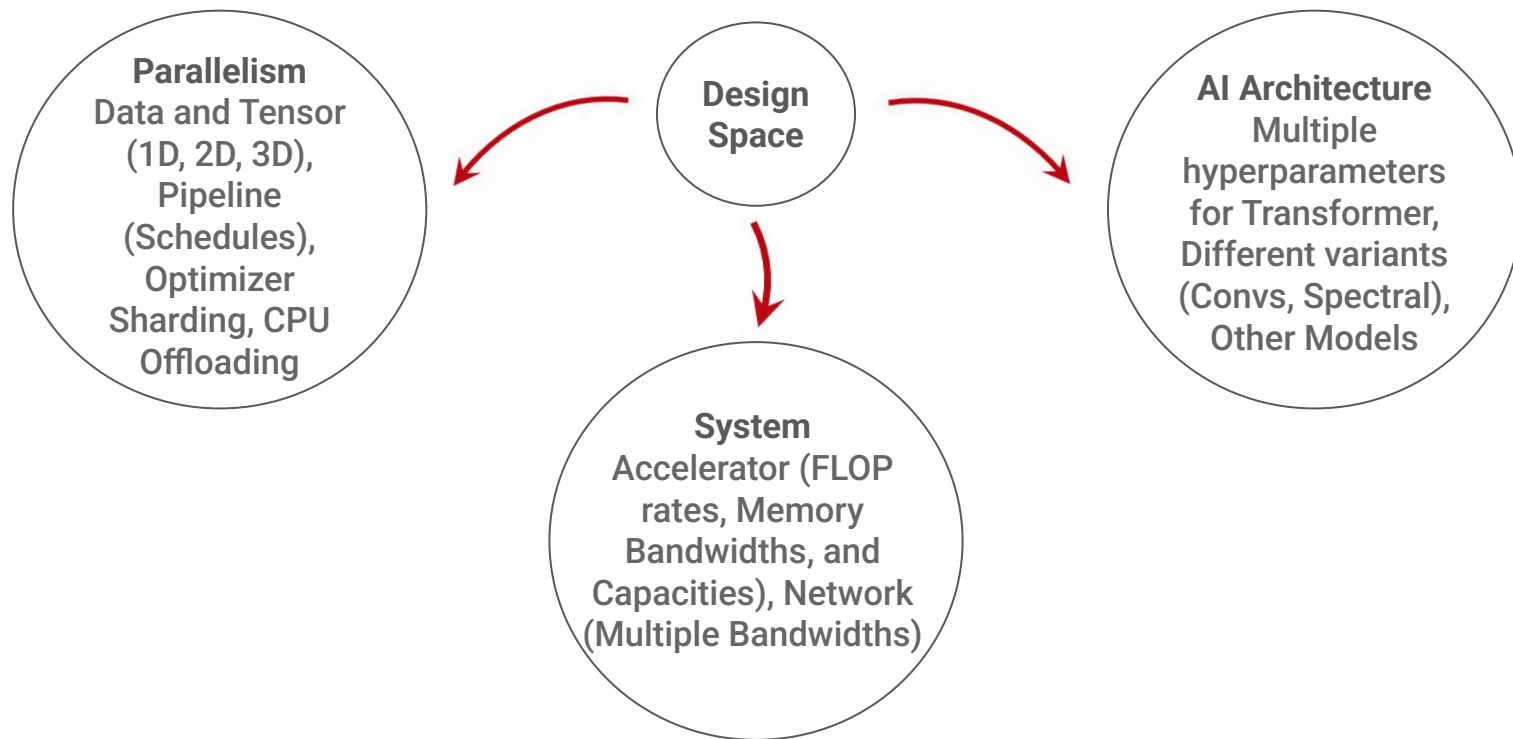
AI Performance Modeling is **Challenging**



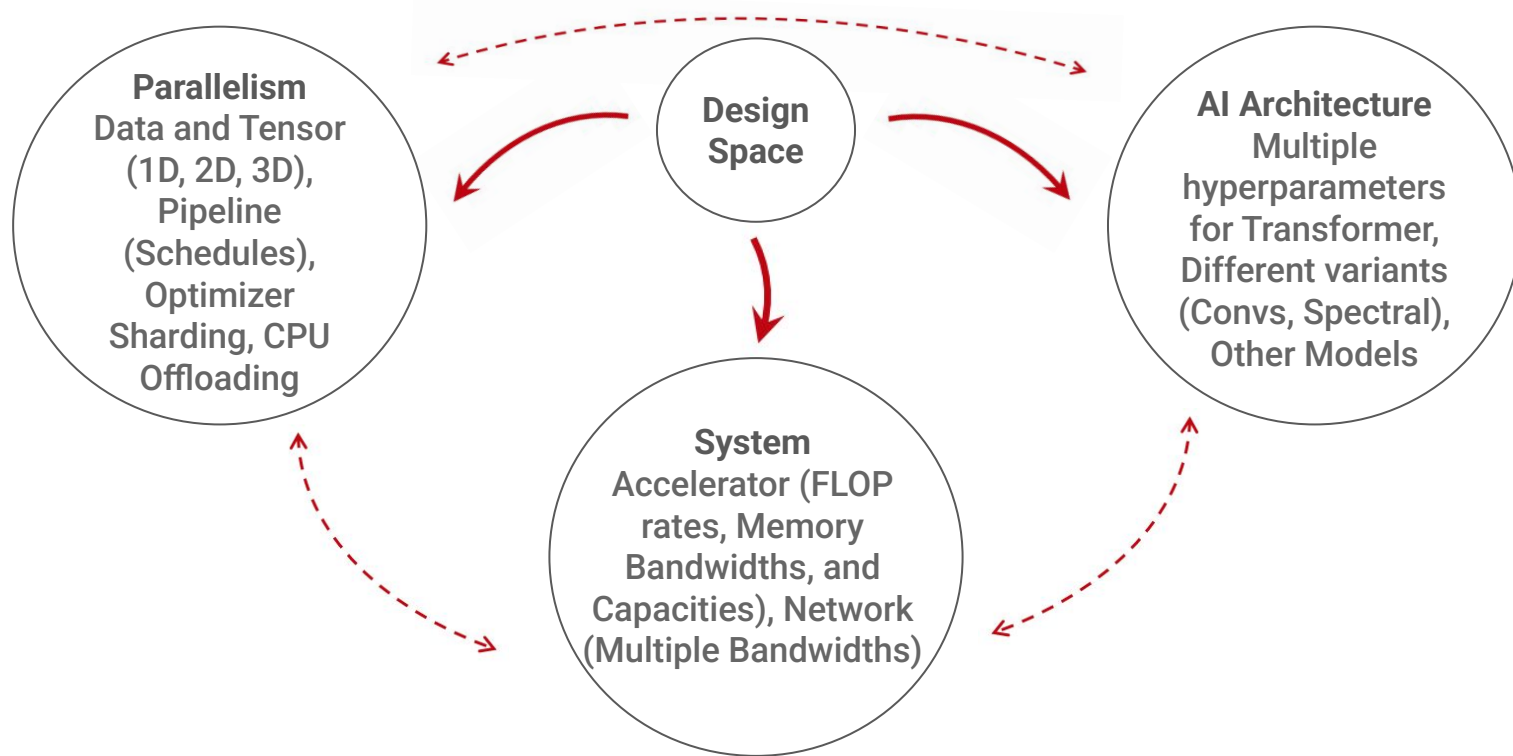
AI Performance Modeling is **Challenging**



AI Performance Modeling is **Challenging**



AI Performance Modeling is **Challenging**



Analytical and Parameterized Models can be Valuable

AMPeD: An Analytical Model for Performance in Distributed Training of Transformers, [ISPASS23](#)
Calculon. A Methodology and Tool for High-Level Co-Design of Systems and Large Language Models. [SC23](#)
Comprehensive Performance Modeling and System Design Insights for Foundation Models, PMBS, [SC24](#), [Github](#)



U.S. DEPARTMENT OF
ENERGY

Office of
Science



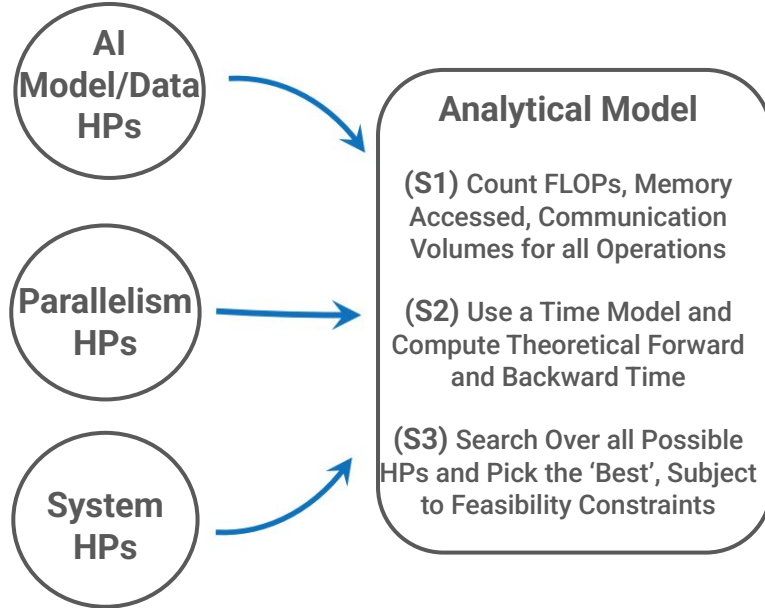
Analytical and Parameterized Models can be Valuable

AI
Model/Data
HPs

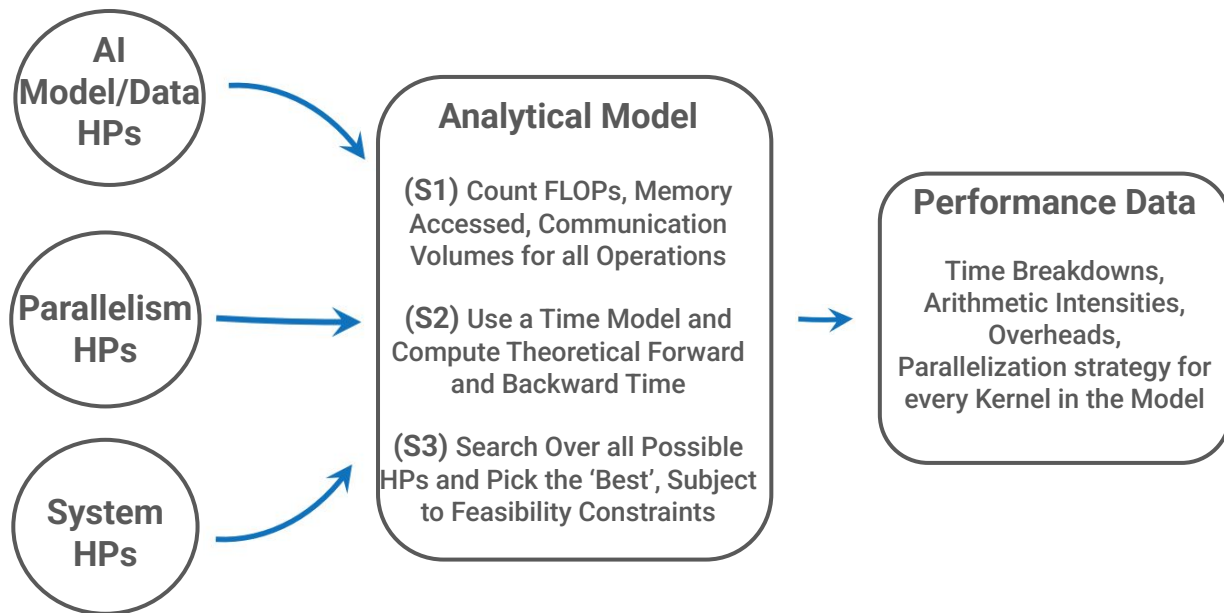
Parallelism
HPs

System
HPs

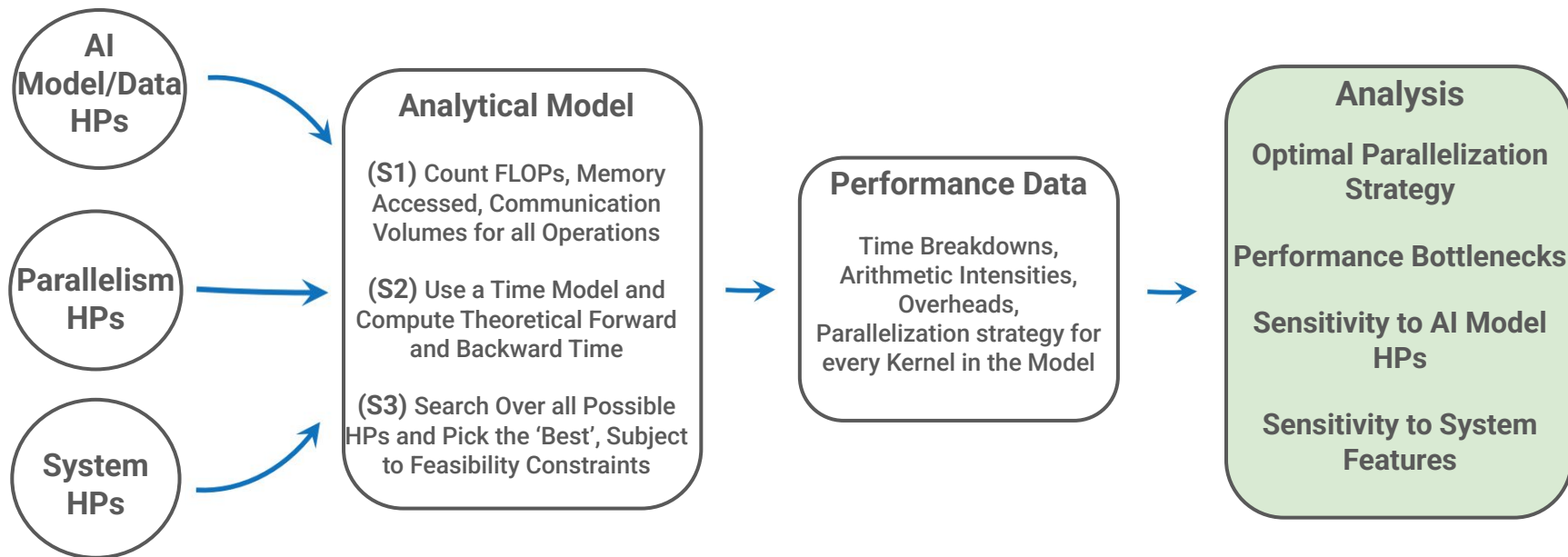
Analytical and Parameterized Models can be Valuable



Analytical and Parameterized Models can be Valuable



Analytical and Parameterized Models can be Valuable



AMPeD: An Analytical Model for Performance in Distributed Training of Transformers, [ISPASS23](#)
Calculon. A Methodology and Tool for High-Level Co-Design of Systems and Large Language Models. [SC23](#)
Comprehensive Performance Modeling and System Design Insights for Foundation Models, PMBS, [SC24](#), [Github](#)

Analyze **Varying Needs** for Transformers in Science

- Counting FLOPs, communication volume is dependent on the parallelism
- Long sequence lengths may necessitate N-D parallelism

Operation	Partitioned Tensor Shapes	Type	Vol
1D TP over n_t GPUs			
<i>SA</i>			
$\tilde{\mathbf{X}} = \text{LN}(\mathbf{X})$	$\tilde{\mathbf{X}} : (b, l, e), \mathbf{X} : (b, \frac{l}{n_t}, e),$	\mathcal{AG}	ble
$\mathbf{Q} = \tilde{\mathbf{X}}\mathbf{W}_{\mathbf{Q}}$	$\mathbf{Q} : (b, \frac{h}{n_t}, l, e_h), \mathbf{W}_{\mathbf{Q}} : (e, \frac{e}{n_t}),$	-	0
$\mathbf{A} = \mathbf{Q}\mathbf{K}^T$	$\mathbf{A} : (b, \frac{h}{n_t}, l, l), \mathbf{K} : (b, \frac{h}{n_t}, l, e_h)$	-	0
$\mathbf{S} = \mathbf{A}\mathbf{V}$	$\mathbf{S} : (b, \frac{h}{n_t}, l, e_h), \mathbf{V} : (b, \frac{h}{n_t}, l, e_h)$	-	0
$\mathbf{Y} = \mathbf{S}\mathbf{W}_{\mathbf{P}}$	$\mathbf{Y} : (b, \frac{l}{n_t}, e), \mathbf{W}_{\mathbf{P}} : (\frac{e}{n_t}, e)$	\mathcal{RS}	ble
<i>MLP</i>			
$\tilde{\mathbf{Y}} = \text{LN}(\mathbf{Y})$	$\tilde{\mathbf{Y}} : (b, l, e), \mathbf{Y} : (b, \frac{l}{n_t}, e),$	\mathcal{AG}	ble
$\mathbf{Z} = \tilde{\mathbf{Y}}\mathbf{W}_1$	$\mathbf{Z} : (b, l, f/n_t), \mathbf{W}_1 : (e, \frac{f}{n_t})$	-	0
$\mathbf{X} = \mathbf{Z}\mathbf{W}_2$	$\mathbf{X} : (b, \frac{l}{n_t}, e), \mathbf{W}_2 : (\frac{f}{n_t}, e)$	\mathcal{RS}	ble

Analyze **Varying Needs** for Transformers in Science

- Counting FLOPs, communication volume is dependent on the parallelism
- Long sequence lengths may necessitate N-D parallelism

Operation	Partitioned Tensor Shapes	Type	Vol
1D TP over n_t GPUs			
<i>SA</i>			
$\tilde{\mathbf{X}} = \text{LN}(\mathbf{X})$	$\tilde{\mathbf{X}} : (b, l, e), \mathbf{X} : (b, \frac{l}{n_t}, e),$	\mathcal{AG}	ble
$\mathbf{Q} = \tilde{\mathbf{X}}\mathbf{W}_{\mathbf{Q}}$	$\mathbf{Q} : (b, \frac{h}{n_t}, l, e_h), \mathbf{W}_{\mathbf{Q}} : (e, \frac{e}{n_t}),$	-	0
$\mathbf{A} = \mathbf{Q}\mathbf{K}^T$	$\mathbf{A} : (b, \frac{h}{n_t}, l, l), \mathbf{K} : (b, \frac{h}{n_t}, l, e_h)$	-	0
$\mathbf{S} = \mathbf{A}\mathbf{V}$	$\mathbf{S} : (b, \frac{h}{n_t}, l, e_h), \mathbf{V} : (b, \frac{h}{n_t}, l, e_h)$	-	0
$\mathbf{Y} = \mathbf{S}\mathbf{W}_{\mathbf{P}}$	$\mathbf{Y} : (b, \frac{l}{n_t}, e), \mathbf{W}_{\mathbf{P}} : (\frac{e}{n_t}, e)$	\mathcal{RS}	ble
<i>MLP</i>			
$\tilde{\mathbf{Y}} = \text{LN}(\mathbf{Y})$	$\tilde{\mathbf{Y}} : (b, l, e), \mathbf{Y} : (b, \frac{l}{n_t}, e),$	\mathcal{AG}	ble
$\mathbf{Z} = \tilde{\mathbf{Y}}\mathbf{W}_1$	$\mathbf{Z} : (b, l, f/n_t), \mathbf{W}_1 : (e, \frac{f}{n_t})$	-	0
$\mathbf{X} = \mathbf{Z}\mathbf{W}_2$	$\mathbf{X} : (b, \frac{l}{n_t}, e), \mathbf{W}_2 : (\frac{f}{n_t}, e)$	\mathcal{RS}	ble

Analyze **Varying Needs** for Transformers in Science

- Counting FLOPs, communication volume is dependent on the parallelism
- Long sequence lengths may necessitate N-D (4D) parallelism

Operation	Partitioned Tensor Shapes	Type	Vol
2D TP over $n_1 \times n_2$ grid of GPUs			
<i>SA</i>			
$\tilde{\mathbf{X}} = \text{LN}(\mathbf{X})$	$\tilde{\mathbf{X}} : (b, \frac{l}{n_2}, e), \mathbf{X} : (b, \frac{l}{n_1 n_2}, e),$	\mathcal{AG}	$b \frac{l}{n_2} e$
$\mathbf{Q} = \tilde{\mathbf{X}} \mathbf{W}_{\mathbf{Q}}$	$\mathbf{Q} : (b, \frac{h}{n_1}, \frac{l}{n_2}, e_h), \mathbf{W}_{\mathbf{Q}} : (e, \frac{e}{n_1}),$	-	0
$\mathbf{A} = \mathbf{Q} \mathbf{K}^T$	$\mathbf{A} : (b, \frac{h}{n_1}, \frac{l}{n_2}, l), \mathbf{K} : (b, \frac{h}{n_1}, l, e_h)$	\mathcal{AG}	$bl \frac{e}{n_1}$
$\mathbf{S} = \mathbf{A} \mathbf{V}$	$\mathbf{S} : (b, \frac{h}{n_1}, \frac{l}{n_2}, e_h), \mathbf{V} : (b, \frac{h}{n_1}, l, e_h)$	\mathcal{AG}	$bl \frac{e}{n_1}$
$\mathbf{Y} = \mathbf{S} \mathbf{W}_{\mathbf{P}}$	$\mathbf{Y} : (b, \frac{l}{n_1 n_2}, e), \mathbf{W}_{\mathbf{P}} : (\frac{e}{n_1}, e)$	\mathcal{RS}	$b \frac{l}{n_2} e$
<i>MLP</i>			
$\tilde{\mathbf{Y}} = \text{LN}(\mathbf{Y})$	$\tilde{\mathbf{Y}} : (b, \frac{l}{n_2}, e), \mathbf{Y} : (b, \frac{l}{n_1 n_2}, e),$	\mathcal{AG}	$b \frac{l}{n_2} e$
$\mathbf{Z} = \tilde{\mathbf{Y}} \mathbf{W}_1$	$\mathbf{Z} : (b, \frac{l}{n_2}, \frac{f}{n_1}), \mathbf{W}_1 : (e, \frac{f}{n_1})$	-	0
$\mathbf{X} = \mathbf{Z} \mathbf{W}_2$	$\mathbf{X} : (b, \frac{l}{n_1 n_2}, e), \mathbf{W}_2 : (\frac{f}{n_1}, e)$	\mathcal{RS}	$b \frac{l}{n_2} e$

Analyze **Varying Needs** for Transformers in Science

- Counting FLOPs, communication volume is dependent on the parallelism
- Long sequence lengths may necessitate N-D (4D) parallelism

Operation	Partitioned Tensor Shapes	Type	Vol
2D TP over $n_1 \times n_2$ grid of GPUs			
<i>SA</i>			
$\tilde{\mathbf{X}} = \text{LN}(\mathbf{X})$	$\tilde{\mathbf{X}} : (b, \frac{l}{n_2}, e), \mathbf{X} : (b, \frac{l}{n_1 n_2}, e),$	<i>AG</i>	$b \frac{l}{n_2} e$
$\mathbf{Q} = \tilde{\mathbf{X}} \mathbf{W}_{\mathbf{Q}}$	$\mathbf{Q} : (b, \frac{h}{n_1}, \frac{l}{n_2}, e_h), \mathbf{W}_{\mathbf{Q}} : (e, \frac{e}{n_1}),$	-	0
$\mathbf{A} = \mathbf{Q} \mathbf{K}^T$	$\mathbf{A} : (b, \frac{h}{n_1}, \frac{l}{n_2}, l), \mathbf{K} : (b, \frac{h}{n_1}, l, e_h)$	<i>AG</i>	$bl \frac{e}{n_1}$
$\mathbf{S} = \mathbf{A} \mathbf{V}$	$\mathbf{S} : (b, \frac{h}{n_1}, \frac{l}{n_2}, e_h), \mathbf{V} : (b, \frac{h}{n_1}, l, e_h)$	<i>AG</i>	$bl \frac{e}{n_1}$
$\mathbf{Y} = \mathbf{S} \mathbf{W}_{\mathbf{P}}$	$\mathbf{Y} : (b, \frac{l}{n_1 n_2}, e), \mathbf{W}_{\mathbf{P}} : (\frac{e}{n_1}, e)$	<i>RS</i>	$b \frac{l}{n_2} e$
<i>MLP</i>			
$\tilde{\mathbf{Y}} = \text{LN}(\mathbf{Y})$	$\tilde{\mathbf{Y}} : (b, \frac{l}{n_2}, e), \mathbf{Y} : (b, \frac{l}{n_1 n_2}, e),$	<i>AG</i>	$b \frac{l}{n_2} e$
$\mathbf{Z} = \tilde{\mathbf{Y}} \mathbf{W}_1$	$\mathbf{Z} : (b, \frac{l}{n_2}, \frac{f}{n_1}), \mathbf{W}_1 : (e, \frac{f}{n_1})$	-	0
$\mathbf{X} = \mathbf{Z} \mathbf{W}_2$	$\mathbf{X} : (b, \frac{l}{n_1 n_2}, e), \mathbf{W}_2 : (\frac{f}{n_1}, e)$	<i>RS</i>	$b \frac{l}{n_2} e$

Analyze **Varying Needs** for Transformers in Science

- Counting FLOPs, communication volume is dependent on the parallelism
- Long sequence lengths may necessitate N-D (4D) parallelism

Operation	Partitioned Tensor Shapes	Type	Vol
2D TP over $n_1 \times n_2$ grid of GPUs			
<i>SA</i>			
$\tilde{\mathbf{X}} = \text{LN}(\mathbf{X})$	$\tilde{\mathbf{X}} : (b, \frac{l}{n_2}, e), \mathbf{X} : (b, \frac{l}{n_1 n_2}, e),$	<i>AG</i>	$b \frac{l}{n_2} e$
$\mathbf{Q} = \tilde{\mathbf{X}} \mathbf{W}_{\mathbf{Q}}$	$\mathbf{Q} : (b, \frac{h}{n_1}, \frac{l}{n_2}, e_h), \mathbf{W}_{\mathbf{Q}} : (e, \frac{e}{n_1}),$	-	0
$\mathbf{A} = \mathbf{Q} \mathbf{K}^T$	$\mathbf{A} : (b, \frac{h}{n_1}, \frac{l}{n_2}, l), \mathbf{K} : (b, \frac{h}{n_1}, l, e_h)$	<i>AG</i>	$bl \frac{e}{n_1}$
$\mathbf{S} = \mathbf{A} \mathbf{V}$	$\mathbf{S} : (b, \frac{h}{n_1}, \frac{l}{n_2}, e_h), \mathbf{V} : (b, \frac{h}{n_1}, l, e_h)$	<i>AG</i>	$bl \frac{e}{n_1}$
$\mathbf{Y} = \mathbf{S} \mathbf{W}_{\mathbf{P}}$	$\mathbf{Y} : (b, \frac{l}{n_1 n_2}, e), \mathbf{W}_{\mathbf{P}} : (\frac{e}{n_1}, e)$	<i>RS</i>	$b \frac{l}{n_2} e$
<i>MLP</i>			
$\tilde{\mathbf{Y}} = \text{LN}(\mathbf{Y})$	$\tilde{\mathbf{Y}} : (b, \frac{l}{n_2}, e), \mathbf{Y} : (b, \frac{l}{n_1 n_2}, e),$	<i>AG</i>	$b \frac{l}{n_2} e$
$\mathbf{Z} = \tilde{\mathbf{Y}} \mathbf{W}_1$	$\mathbf{Z} : (b, \frac{l}{n_2}, \frac{f}{n_1}), \mathbf{W}_1 : (e, \frac{f}{n_1})$	-	0
$\mathbf{X} = \mathbf{Z} \mathbf{W}_2$	$\mathbf{X} : (b, \frac{l}{n_1 n_2}, e), \mathbf{W}_2 : (\frac{f}{n_1}, e)$	<i>RS</i>	$b \frac{l}{n_2} e$

Analyze **Varying Needs** for Transformers in Science

- Long sequence lengths may necessitate 4D parallelism
- Different choices for Matrix Multiplies: SUMMA also possible

$$C = \sum_{\kappa}^{n_b-1} A^{\kappa} B^{\kappa}$$

Algorithm 1 $C = AB$ using SUMMA

```
1: Input:  $A_{ij}, B_{ij}$ 
2: Output:  $C_{ij}$ 
3:  $C = 0$ 
4: for  $\kappa = 0 \rightarrow n_b - 1$  do
5:   for  $i = 0, \dots, n_1 - 1$  Broadcast  $A_i^{\kappa}$  to  $i^{th}$  process row
6:   for  $j = 0, \dots, n_2 - 1$  Broadcast  $B_j^{\kappa}$  to  $j^{th}$  process col
7:    $C_{ij} = C_{ij} + A_i^{\kappa} B_j^{\kappa}$ 
8: end for
9: return  $C_{ij}$ 
```

SUMMA: Scalable Universal Matrix Multiplication Algorithm, [Link](#)

Analyze **Varying Needs** for Transformers in Science

- Long sequence lengths may necessitate 4D parallelism
- Different choices for Matrix Multiplies: SUMMA also possible

Operation	Partitioned Tensor Shapes	Type	Vol
2D TP with SUMMA over $n_1 \times n_2$ grid of GPUs			
SA			
$\tilde{\mathbf{X}} = \text{LN}(\mathbf{X})$	$\tilde{\mathbf{X}} : (b, \frac{l}{n_2}, \frac{e}{n_1}), \mathbf{X} : (b, \frac{l}{n_2}, \frac{e}{n_1}),$	\mathcal{AR}	$b \frac{l}{n_2} e$
$\mathbf{Q} = \tilde{\mathbf{X}} \mathbf{W}_{\mathbf{Q}}$	$\mathbf{Q} : (b, \frac{h}{n_1}, \frac{l}{n_2}, e_h), \mathbf{W}_{\mathbf{Q}} : (\frac{e}{n_2}, \frac{e}{n_1}),$	\mathcal{B}	V_1
$\mathbf{A} = \mathbf{Q} \mathbf{K}^T$	$\mathbf{A} : (b, \frac{h}{n_1}, \frac{l}{n_2}, l), \mathbf{K} : (b, \frac{h}{n_1}, l, e_h)$	\mathcal{AG}	$bl \frac{e}{n_1}$
$\mathbf{S} = \mathbf{A} \mathbf{V}$	$\mathbf{S} : (b, \frac{h}{n_1}, \frac{l}{n_2}, e_h), \mathbf{V} : (b, \frac{h}{n_1}, l, e_h)$	\mathcal{AG}	$bl \frac{e}{n_1}$
$\mathbf{Y} = \mathbf{S} \mathbf{W}_{\mathbf{P}}$	$\mathbf{Y} : (b, \frac{l}{n_2}, e), \mathbf{W}_{\mathbf{P}} : (\frac{e}{n_1}, e)$	\mathcal{RS}	$b \frac{l}{n_2} e$
MLP			
$\tilde{\mathbf{Y}} = \text{LN}(\mathbf{Y})$	$\tilde{\mathbf{Y}} : (b, \frac{l}{n_2}, \frac{e}{n_1}), \mathbf{Y} : (b, \frac{l}{n_2}, \frac{e}{n_1}),$	\mathcal{AR}	$b \frac{l}{n_2} e$
$\mathbf{Z} = \tilde{\mathbf{Y}} \mathbf{W}_1$	$\mathbf{Z} : (b, \frac{l}{n_2}, \frac{f}{n_1}), \mathbf{W}_1 : (\frac{e}{n_2}, \frac{f}{n_1})$	\mathcal{B}	V_2
$\mathbf{X} = \mathbf{Z} \mathbf{W}_2$	$\mathbf{X} : (b, \frac{l}{n_2}, \frac{e}{n_1}), \mathbf{W}_2 : (\frac{f}{n_2}, \frac{e}{n_1})$	\mathcal{B}	V_3

$$V_1 = ble/n_2 + e^2/n_1$$

SUMMA: Scalable Universal Matrix Multiplication Algorithm, [Link](#)

Comprehensive Performance Modeling and System Design Insights for Foundation Models, PMBS,

[SC24](#), [Github](#)

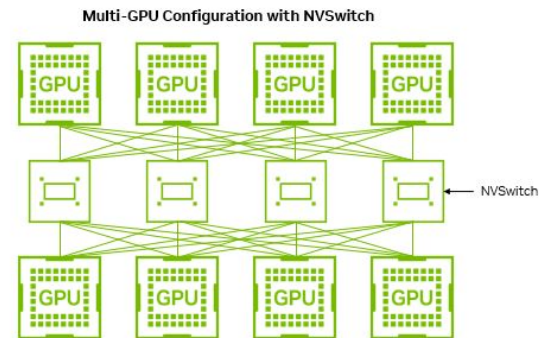
Two Transformer Variants on Different Systems

- Large GPT3 (1T, 2K) on a few trillion tokens
- Large ViT (80B, 250K) on decades of weather data

Two Transformer Variants on Different Systems

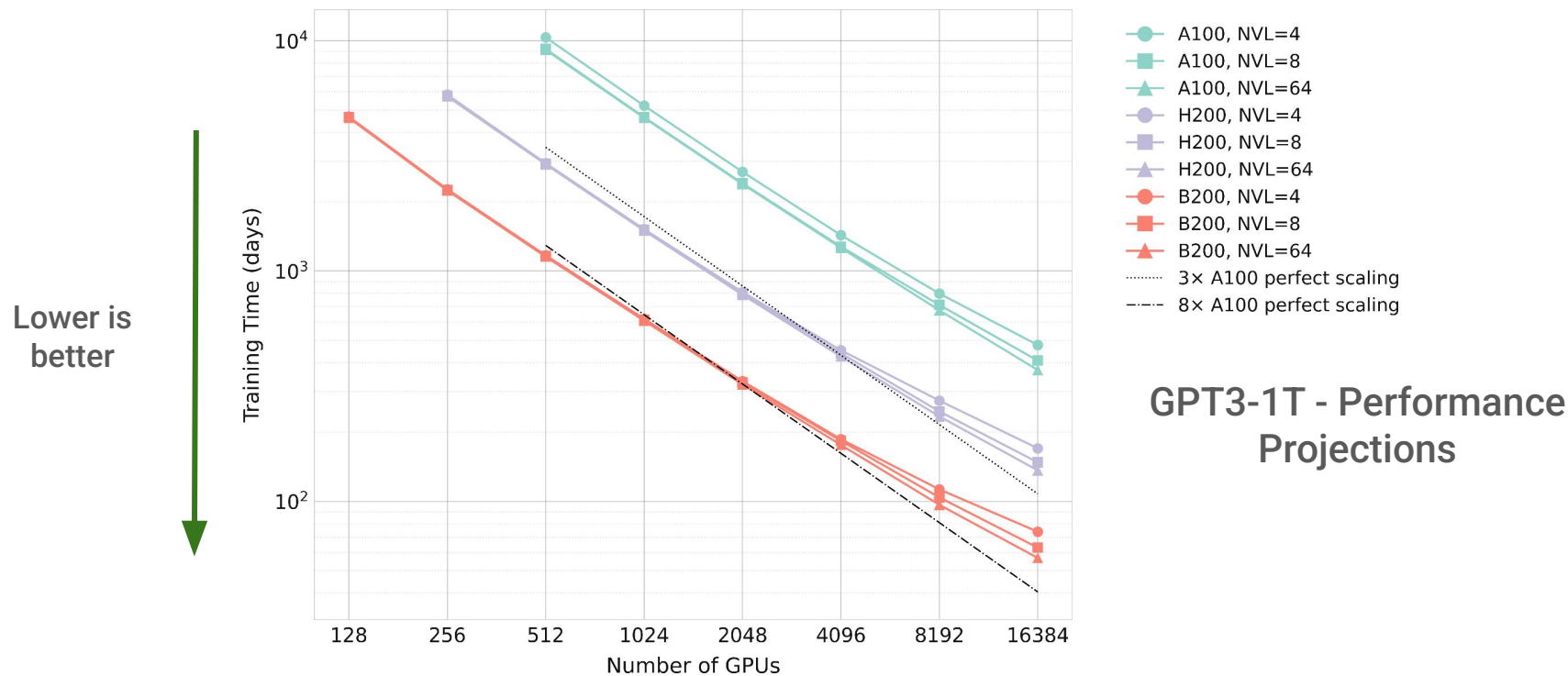
- Large GPT3 (1T, 2K) on a few trillion tokens
- Large ViT (80B, 250K) on decades of weather data

- Three NVIDIA GPU generations: A100, H200, B200
- Three NVLINK/NVL through NVSWITCH domain sizes: 4, 8, 64

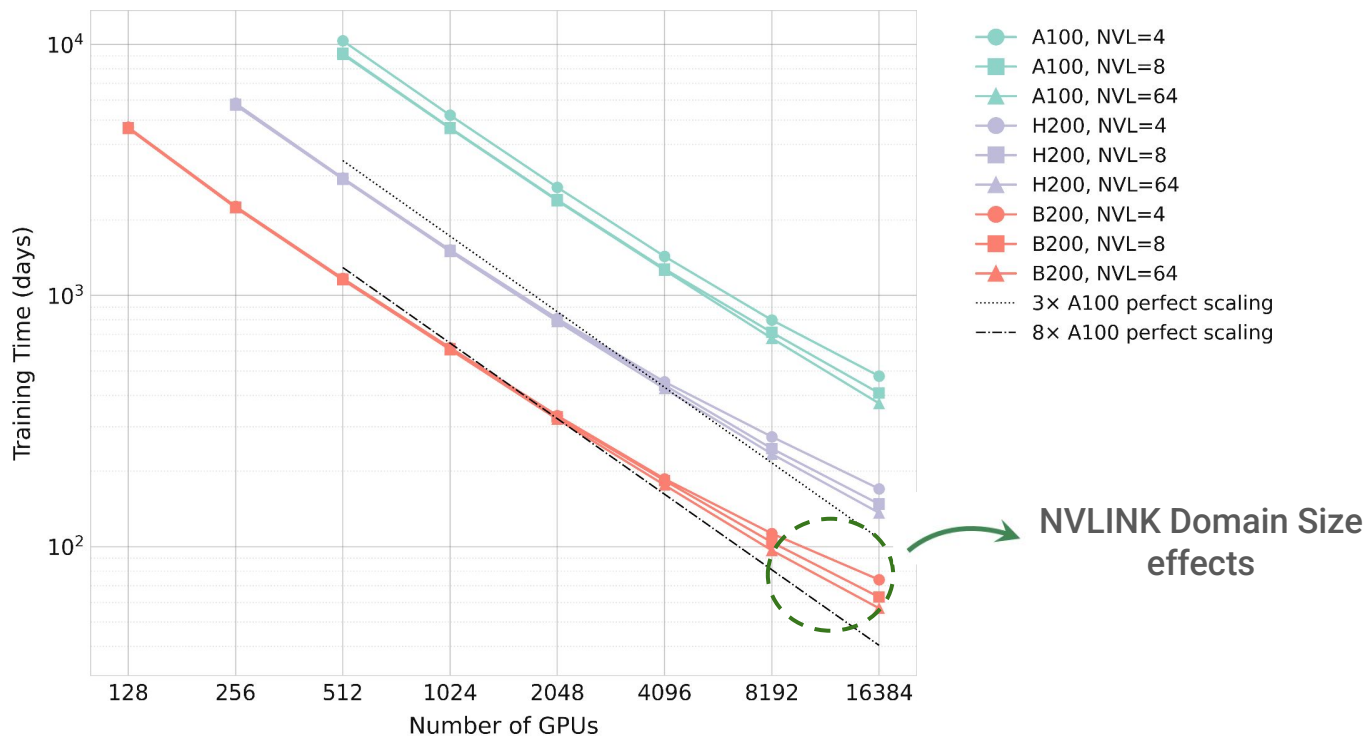


NVSWITCH

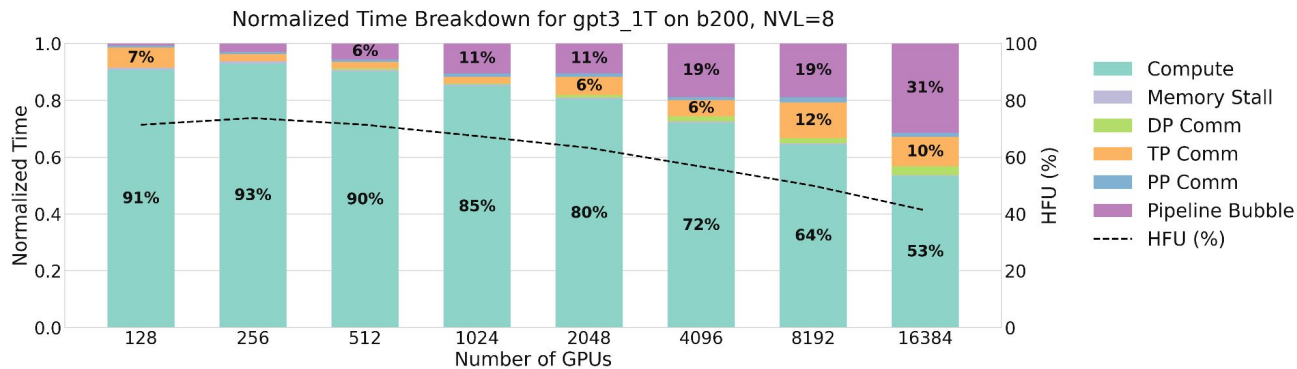
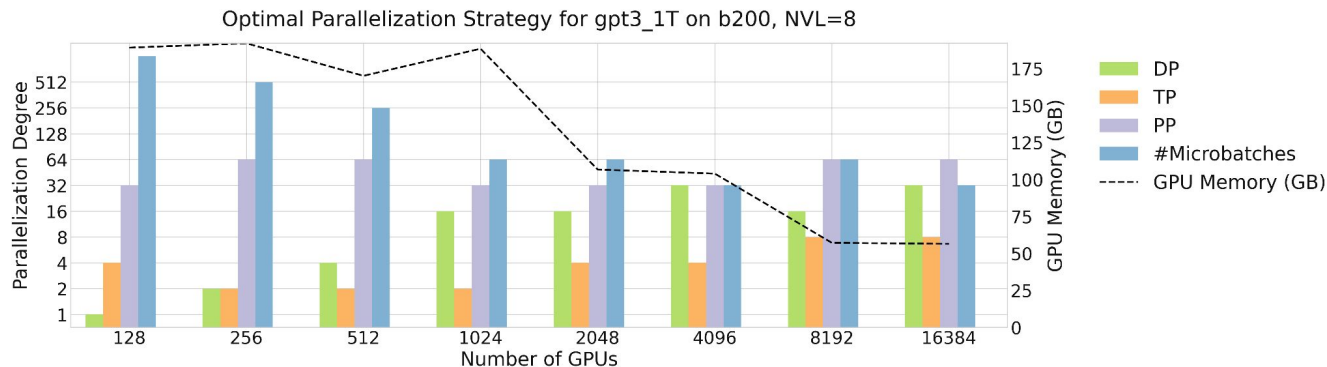
Provides a High-level View of Scaling Behavior



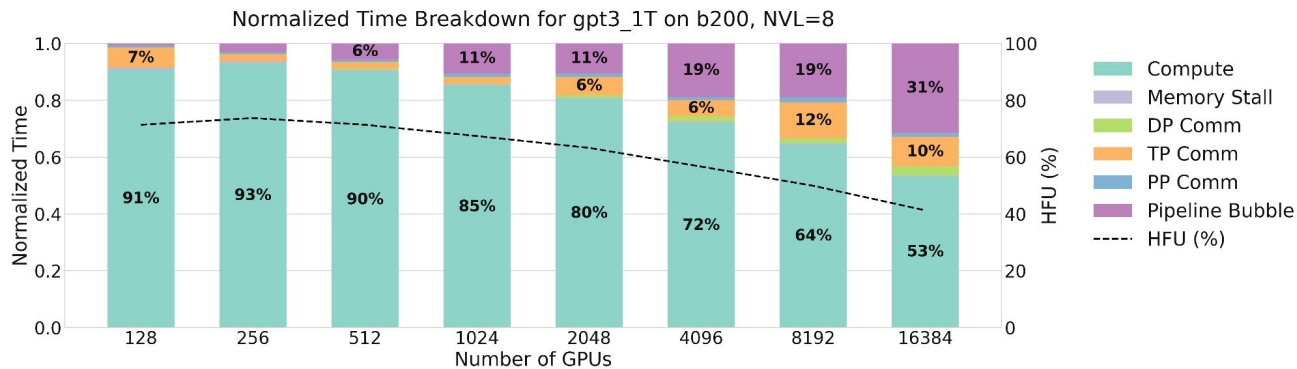
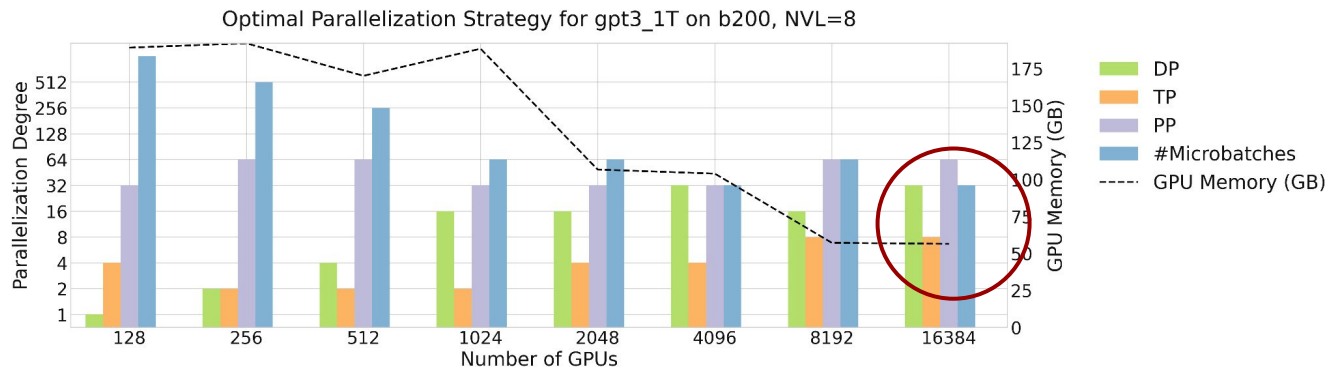
Provides a High-level View of Scaling Behavior



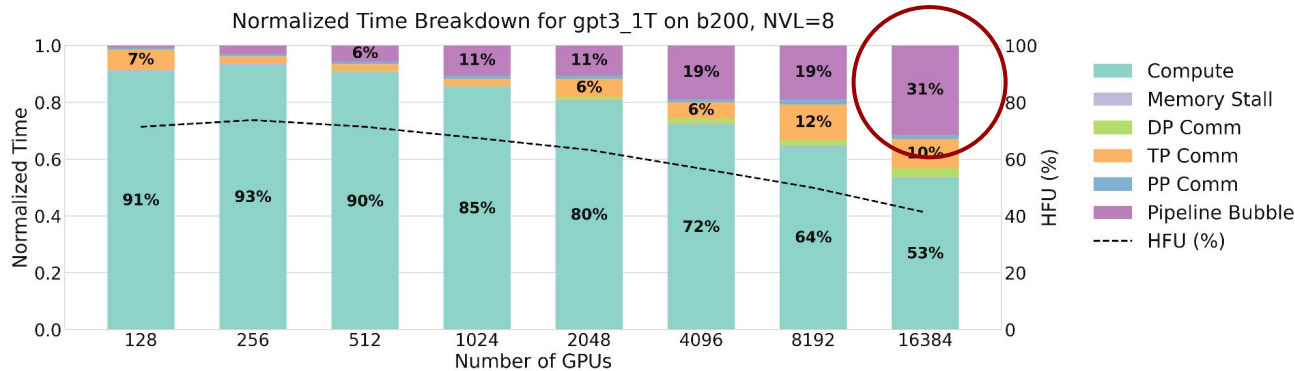
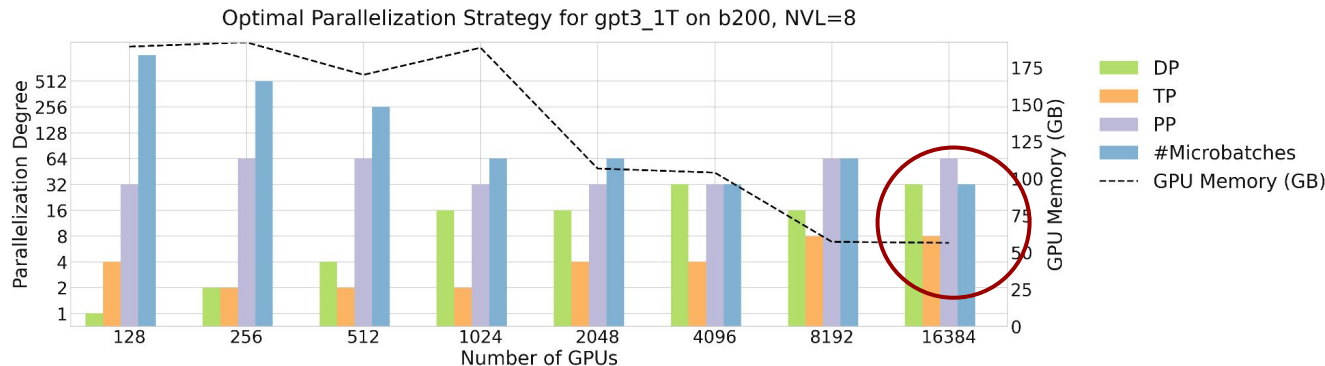
Exposes **Bottlenecks** and **Optimal Parallelism**



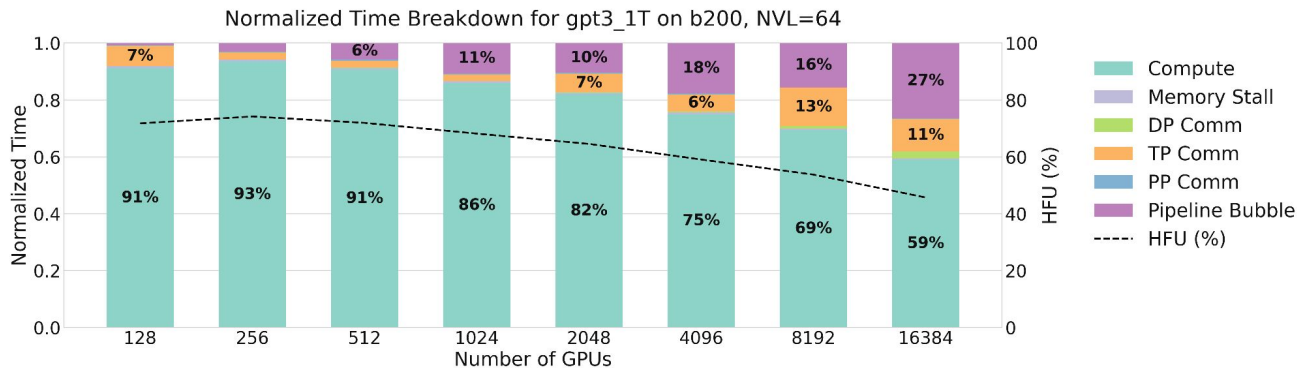
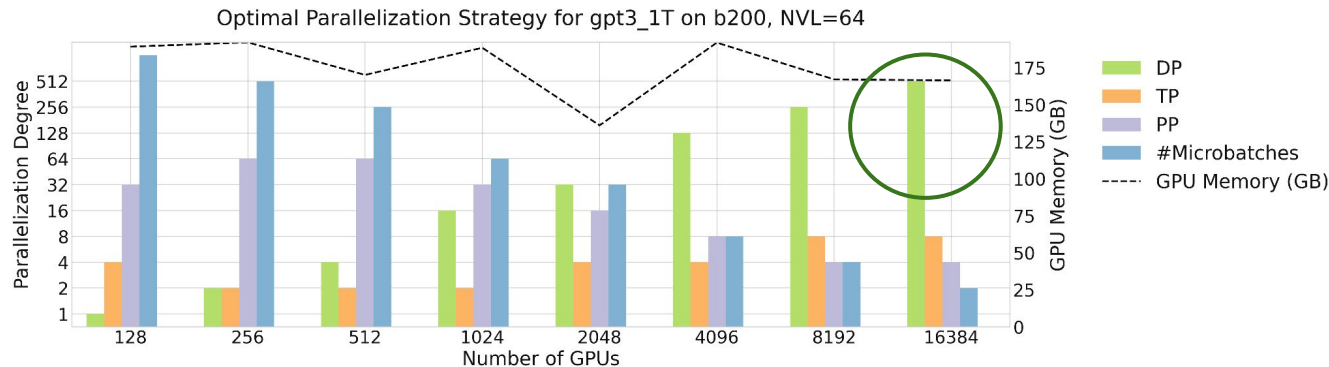
Exposes **Bottlenecks** and **Optimal Parallelism**



Exposes **Bottlenecks** and **Optimal Parallelism**

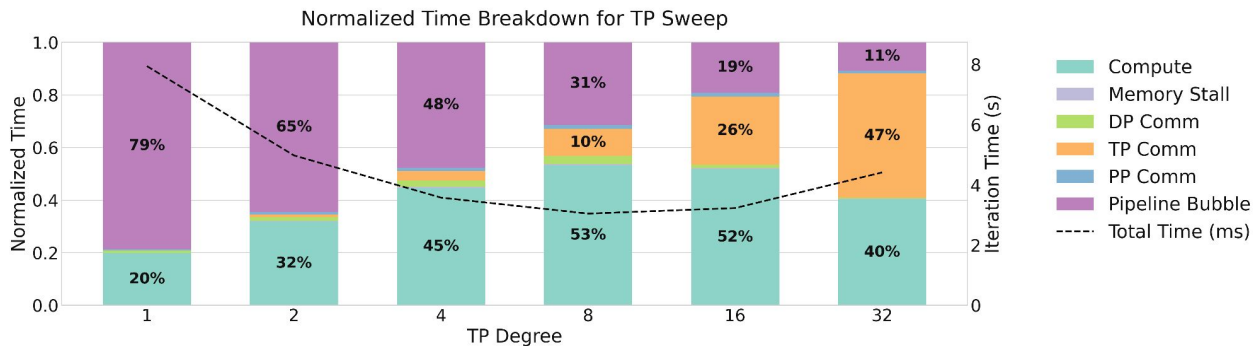
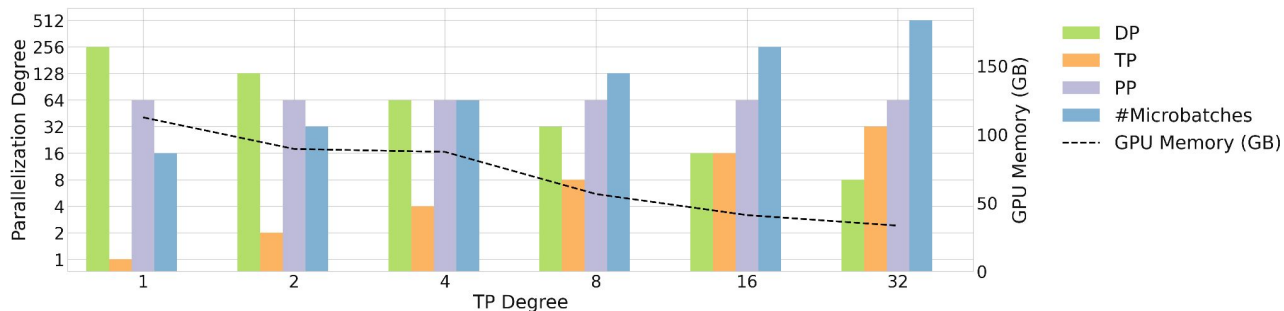


Larger NVLINKs Favor High Data Parallelism



Probe the Model to Get Deeper Insights

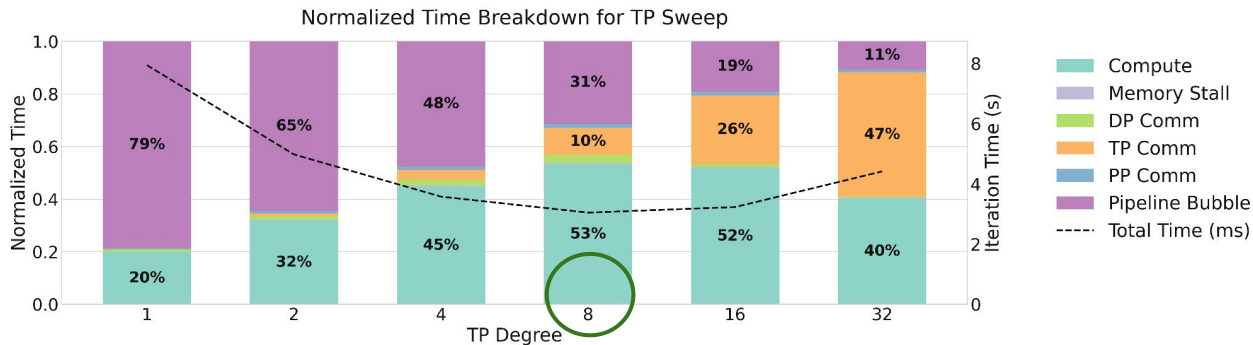
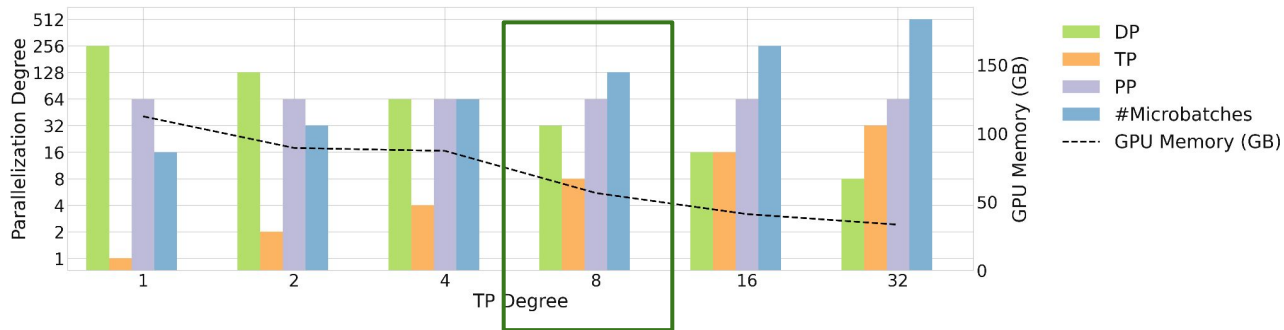
TP Sweep for gpt3_1T on b200, NVL=8
Total GPUs=16384, Fixed PP=64, Micro Batch=1



Fix #GPUs and look around the optimal configuration

Probe the Model to Get Deeper Insights

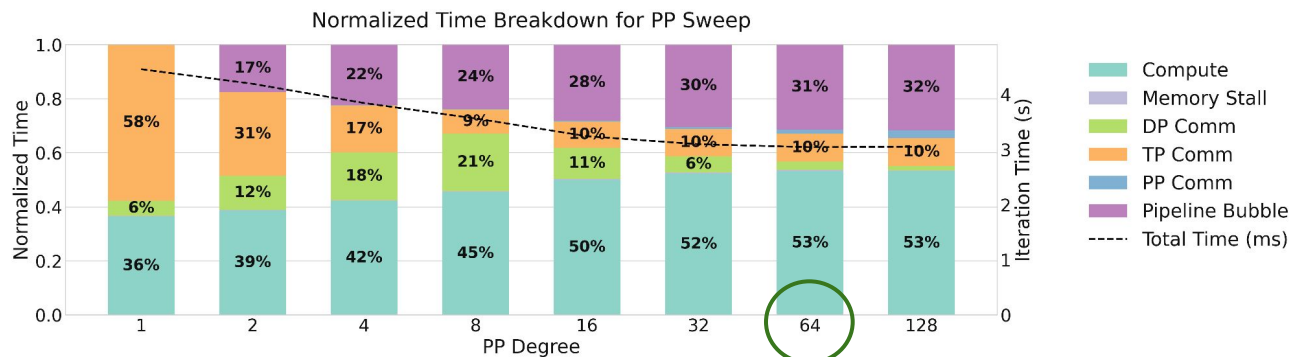
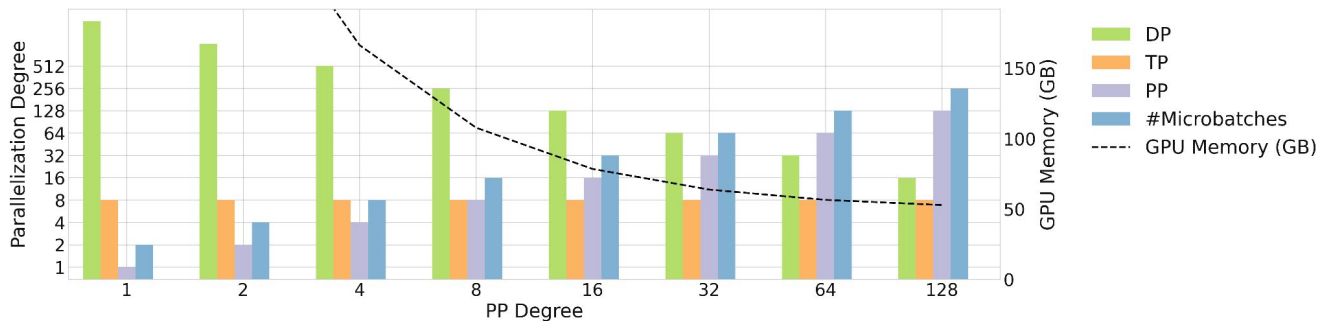
TP Sweep for gpt3_1T on b200, NVL=8
Total GPUs=16384, Fixed PP=64, Micro Batch=1



Fix #GPUs and look around the optimal configuration

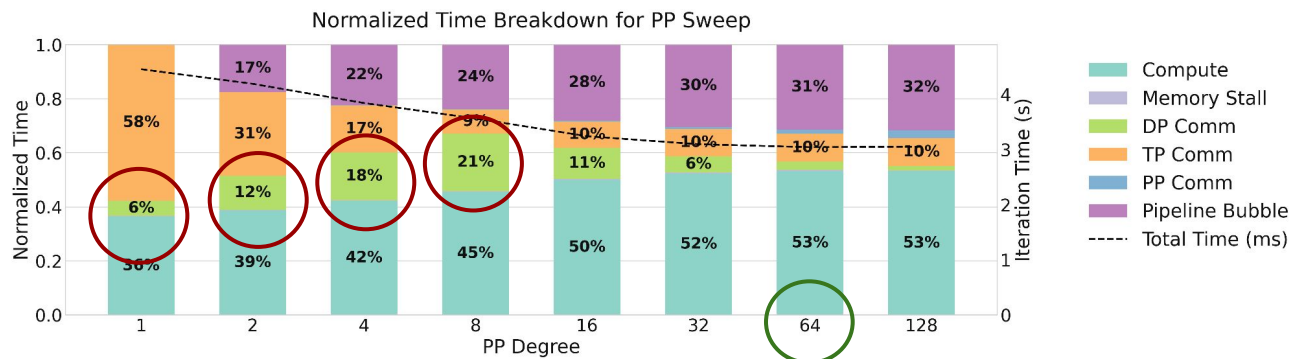
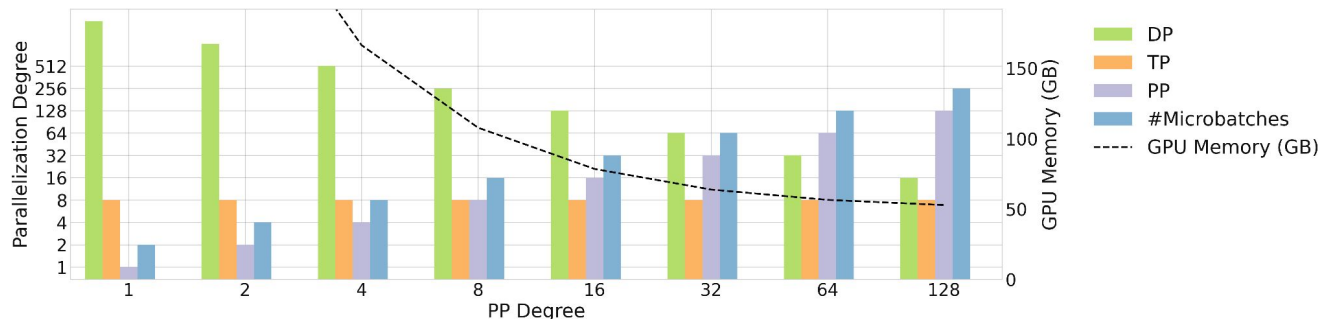
Placement of GPUs Matters

PP Sweep for gpt3_1T on b200, NVL=8
Total GPUs=16384, Fixed TP=8, Micro Batch=1



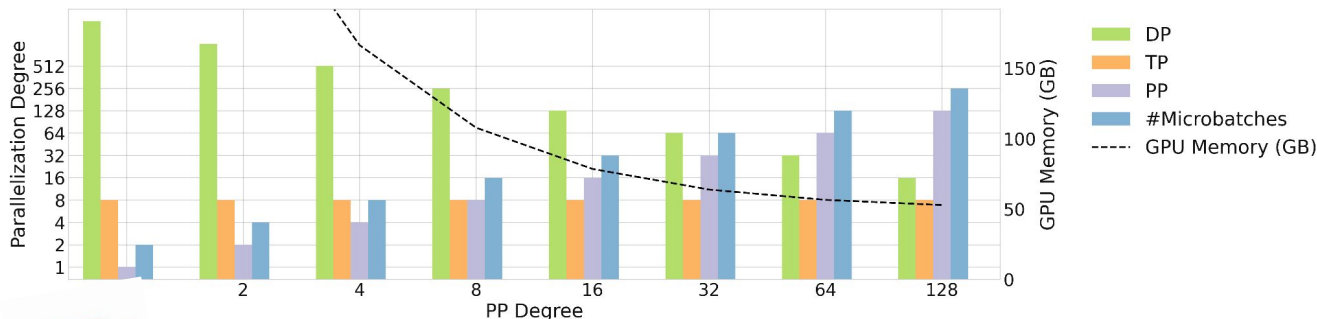
Placement of GPUs Matters

PP Sweep for gpt3_1T on b200, NVL=8
Total GPUs=16384, Fixed TP=8, Micro Batch=1

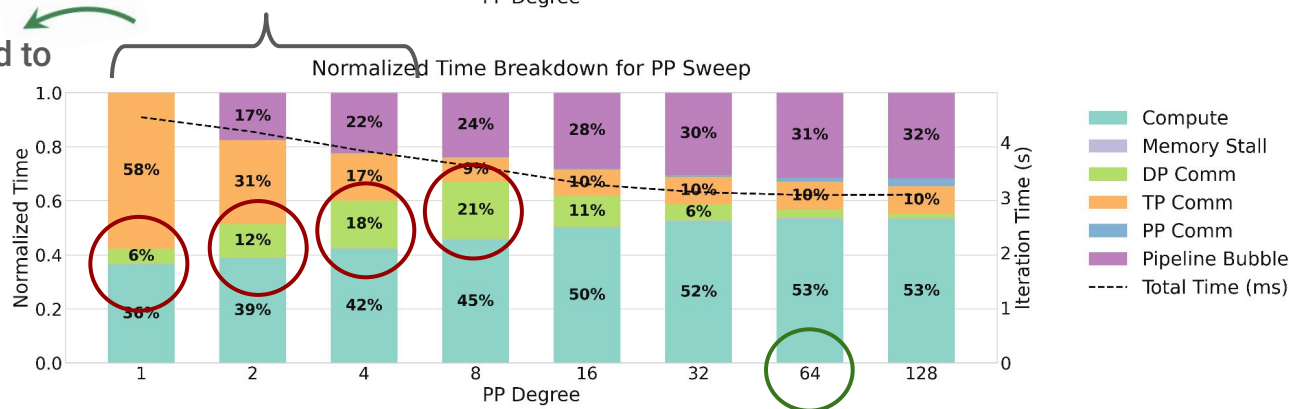


Placement of GPUs Matters

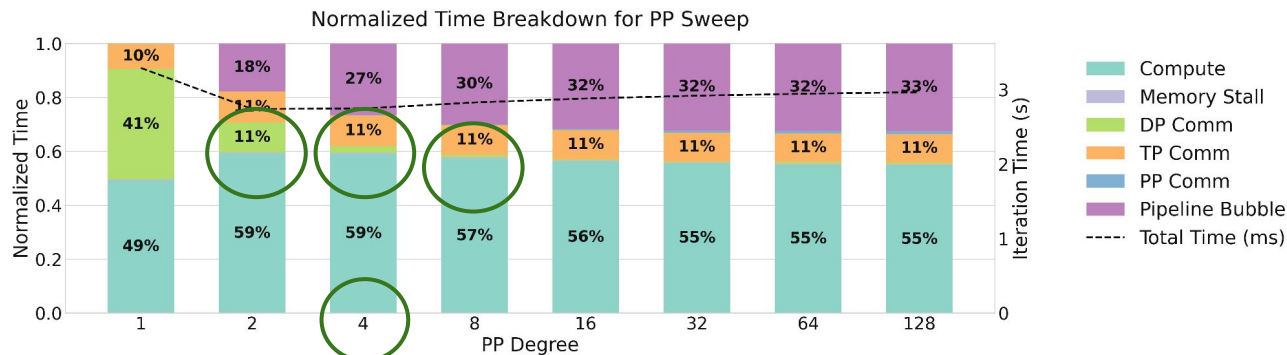
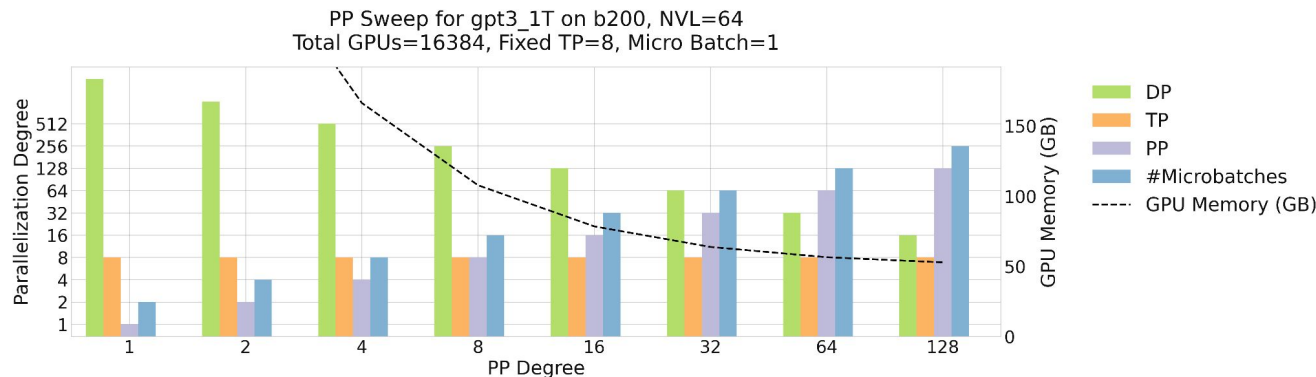
PP Sweep for gpt3_1T on b200, NVL=8
Total GPUs=16384, Fixed TP=8, Micro Batch=1



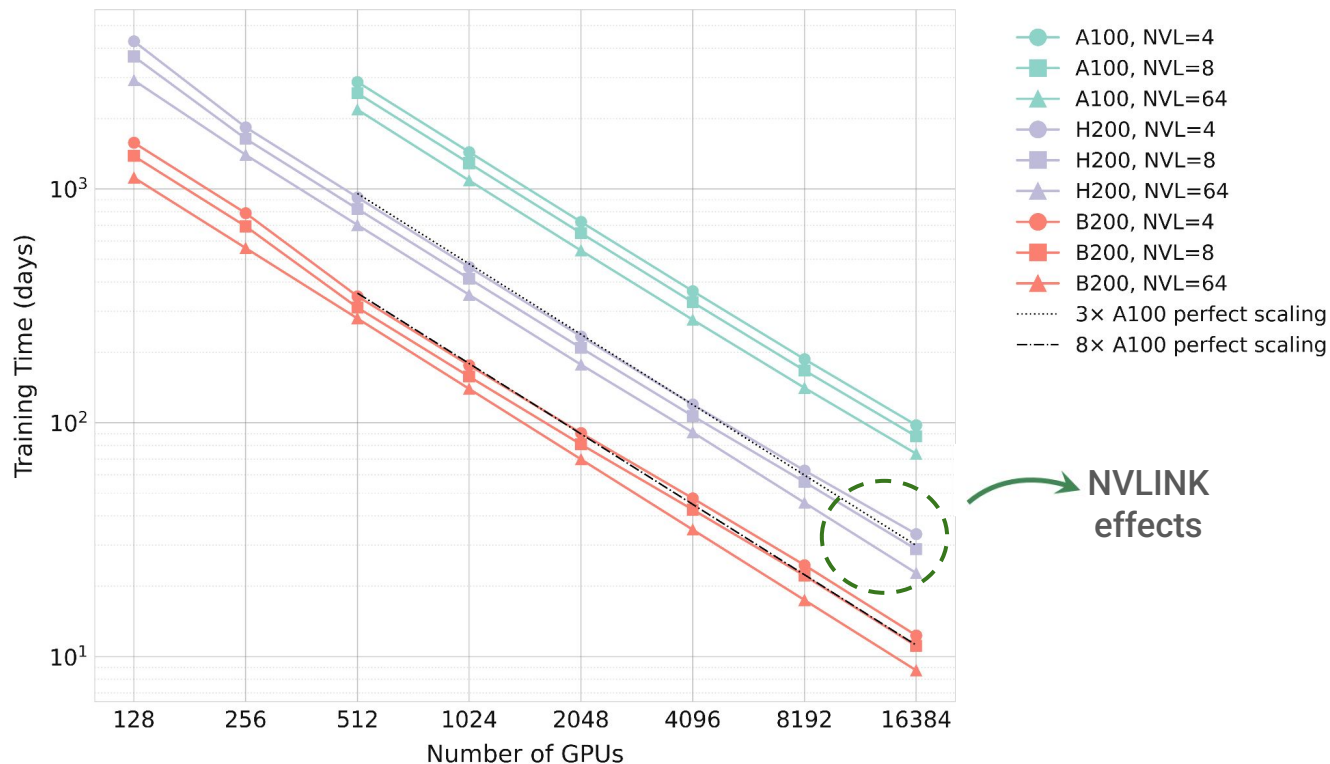
DP GPUs allocated to
NVLINK



Placement of GPUs Favor Data Parallelism for Large NVL

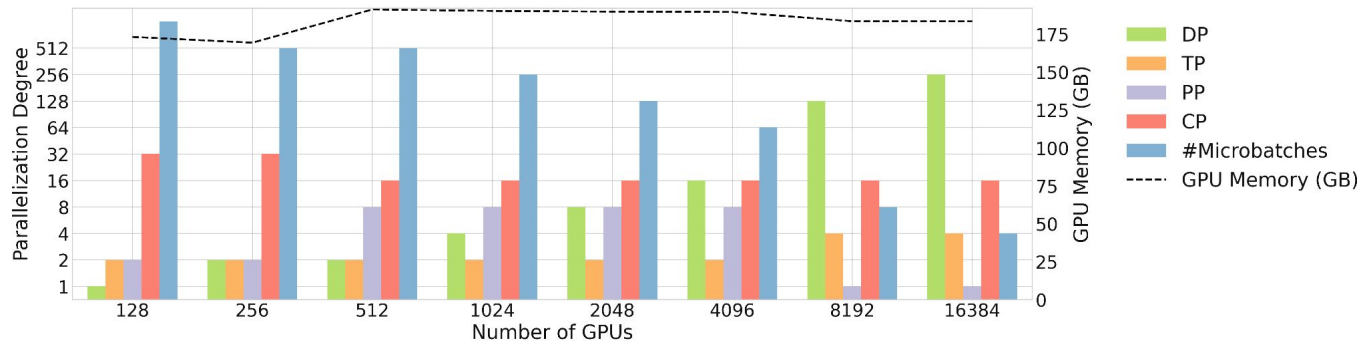


Transformer in Science is **More Sensitive** to the Network

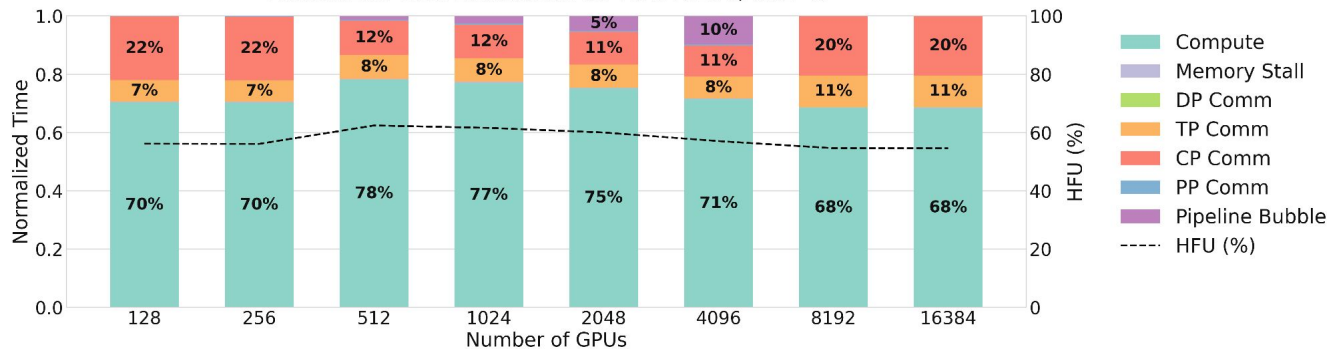


Long Sequences Need 4D Parallelism

Optimal Parallelization Strategy for vit on b200, NVL=8

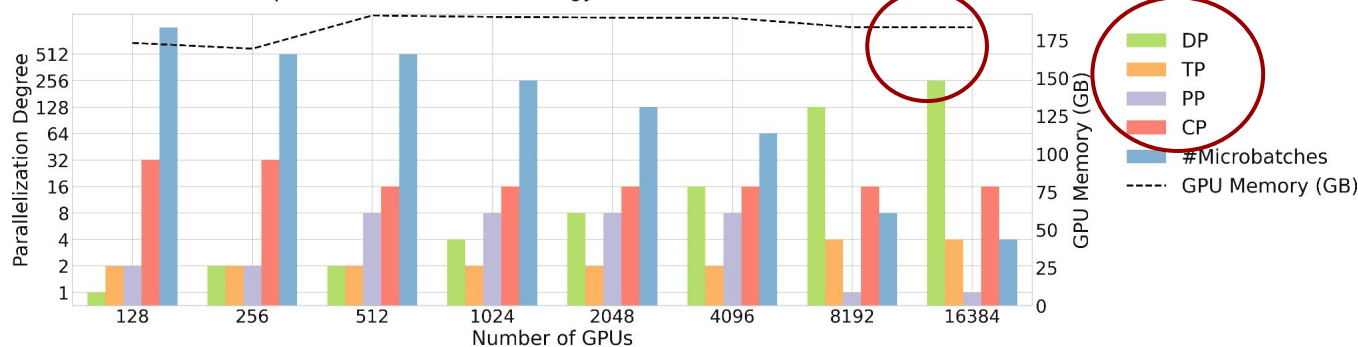


Normalized Time Breakdown for vit on b200, NVL=8

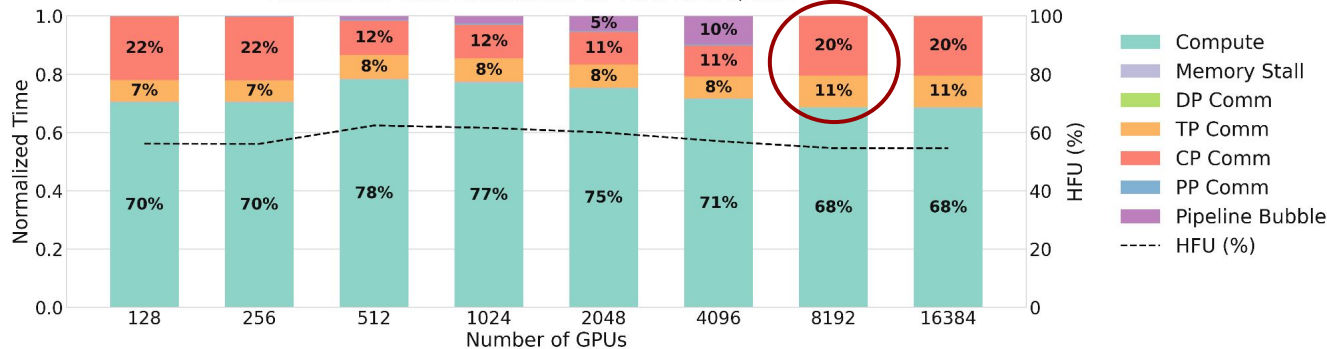


Long Sequences Need 4D Parallelism

Optimal Parallelization Strategy for vit on b200, NVL=8

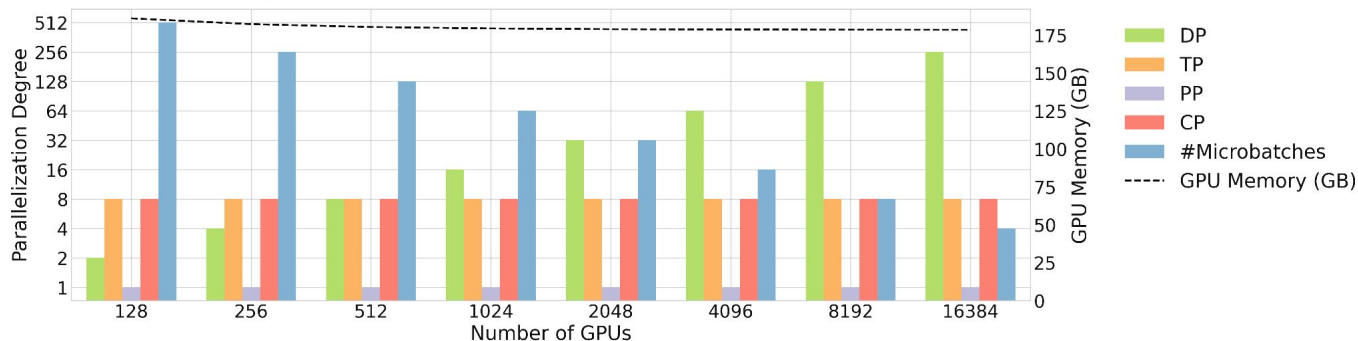


Normalized Time Breakdown for vit on b200, NVL=8

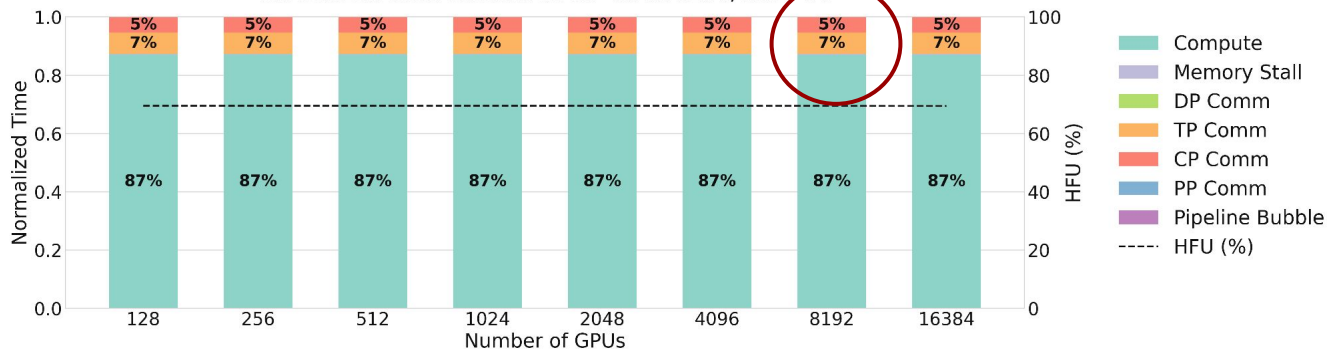


Larger NVLINK Drops **Communication Costs**

Optimal Parallelization Strategy for vit on b200, NVL=64

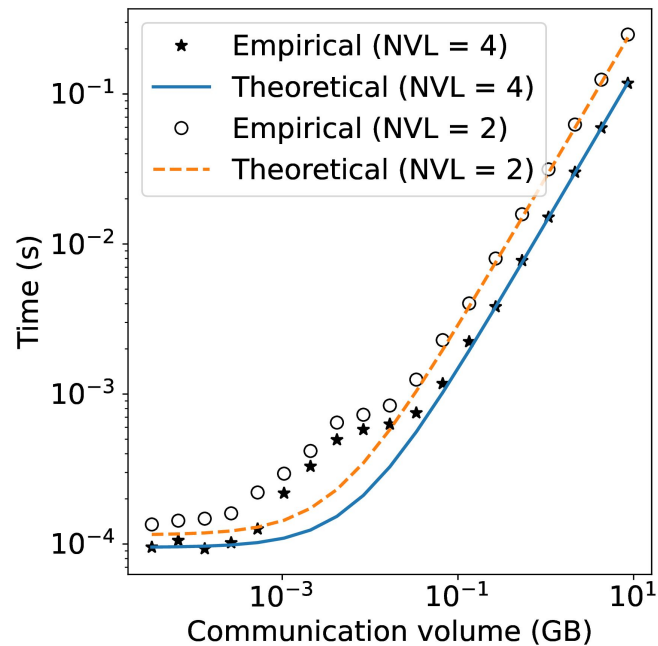


Normalized Time Breakdown for vit on b200, NVL=64



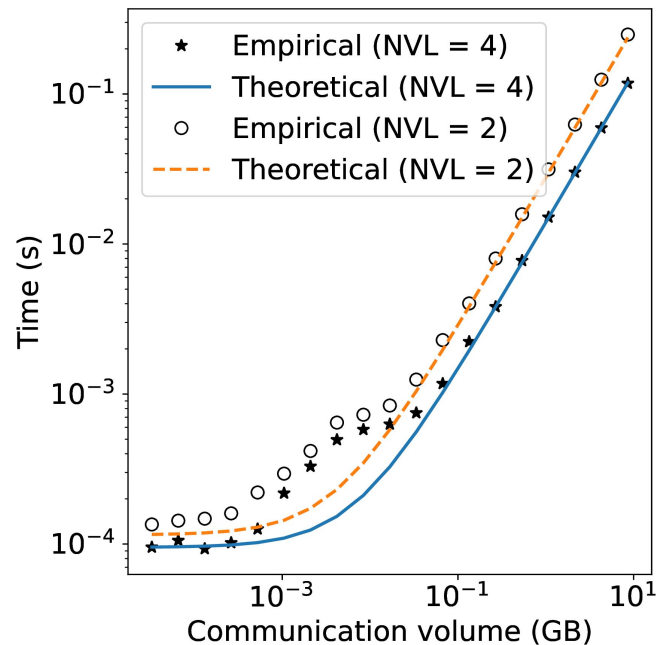
Validation with Megatron-LM

- Validated time models on the Perlmutter supercomputer
 - 4-way NVLINK domain



Validation with Megatron-LM

- Validated time models on the Perlmutter supercomputer
 - 4-way NVLINK domain
- Validated throughput numbers on 512 GPUs
 - GPT3 (175B) and ViT (32K)
- ~10% errors in iteration time
 - Controlled GPU placement with Megatron flags
 - Overlap flags, *FlashAttention*, other optimizations in sync with model
 - Validated sub-optimal configurations as well
- SUMMA validation challenging
 - [ColossalAI](#) in future work



Some Key Takeaways

- Placement of GPUs on high-bandwidth domain affects the optimal parallelism
 - Software codebases to be flexible to this
 - NVLINK domains help expose “easier” parallelisms from the software POV
- LLMs benefit from large NVLINKs at pre-training scales
 - Fine-tuning scales can leverage other parallelization strategies to be less sensitive
 - HBM capacity is underutilized for the largest scales
- Science Transformers benefit uniformly from NVLINK due to memory pressure
 - Demand 4D parallelism (data + pipeline + 2D tensor + optimizer sharding)
 - Capacity is more critical (High capacity, low bandwidth alternatives?)
- 4D/ND (SUMMA/context) parallelism can give you good speedups

Thank You!

