



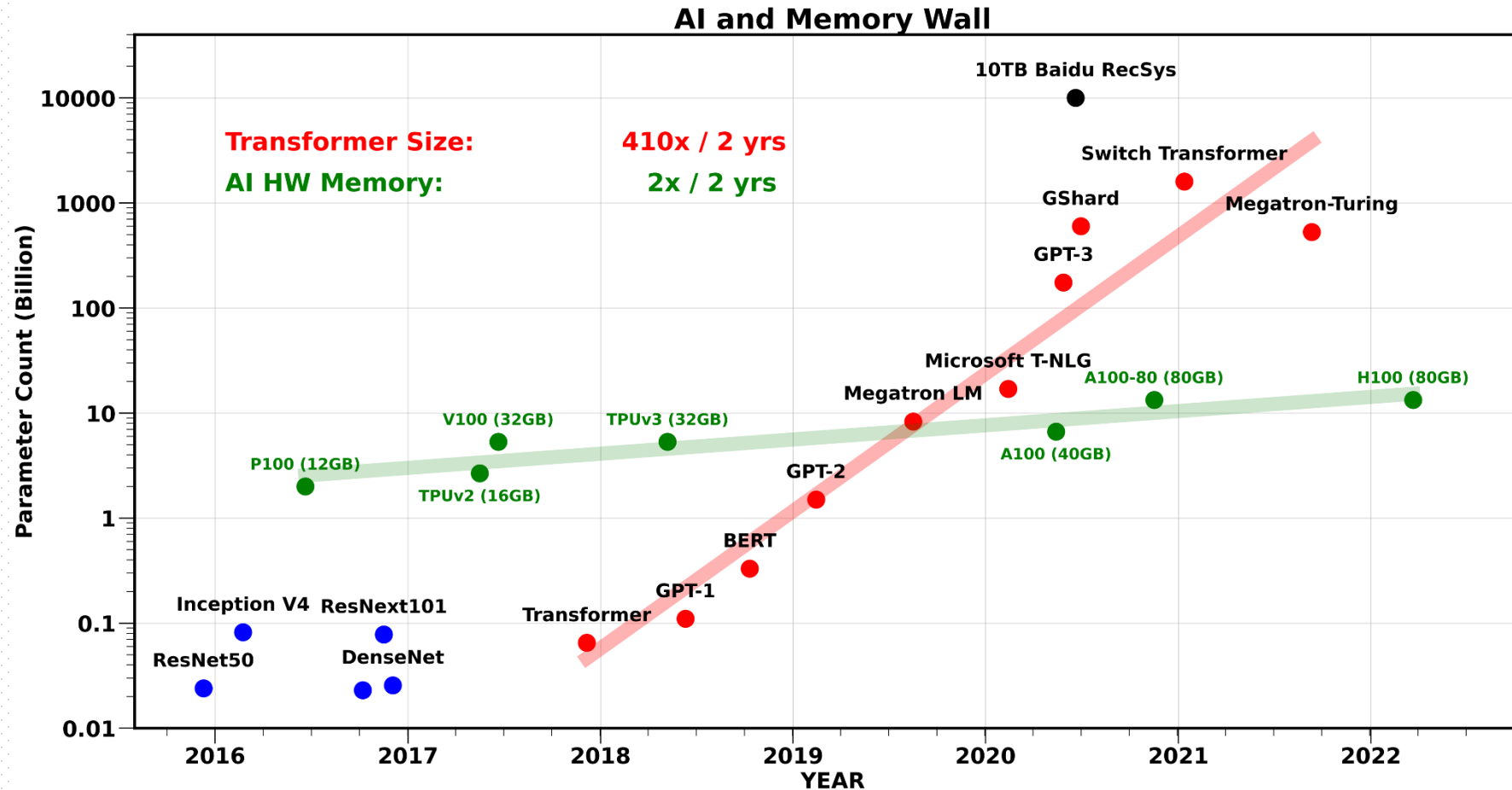
nnScaler: Constraint-Guided Parallelization Plan Generation for Deep Learning Training

Zhiqi Lin[†], Youshan Miao[‡], Quanlu Zhang[‡], Fan Yang[‡], Yi Zhu[‡], Cheng Li[†],
Saeed Maleki[◇], Xu Cao[‡], Ning Shang[‡], Yilei Yang[‡], Weijiang Xu[‡],
Mao Yang[‡], Lintao Zhang[△], Lidong Zhou[‡]

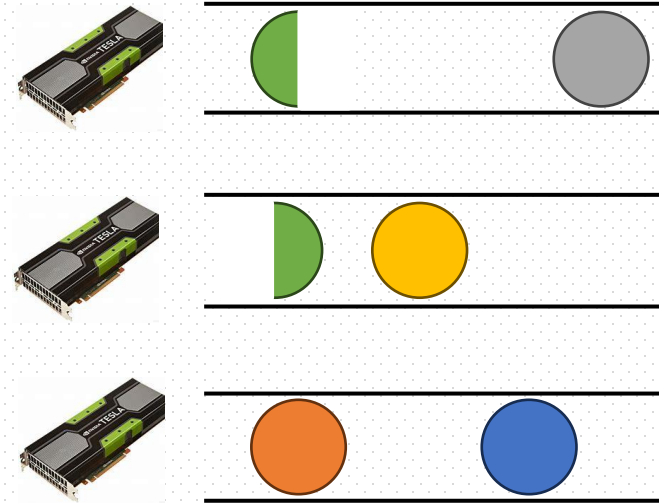
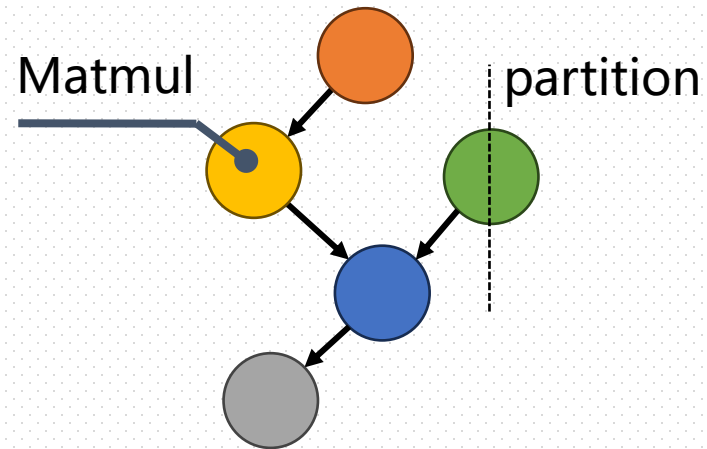
[†]University of Science and Technology of China,
[‡]Microsoft Research, [◇]xAI, [△]BaseBit Technologies



Large Models Need Parallelization



Parallelization Plans Matter



①

Graph partition

- Nodes: operators
- Edges: tensors

②

Spatial-temporal schedule

- Operator placement
- Determine execution order

Parallelization plans affect training performance greatly

Hard to Find Efficient Plans

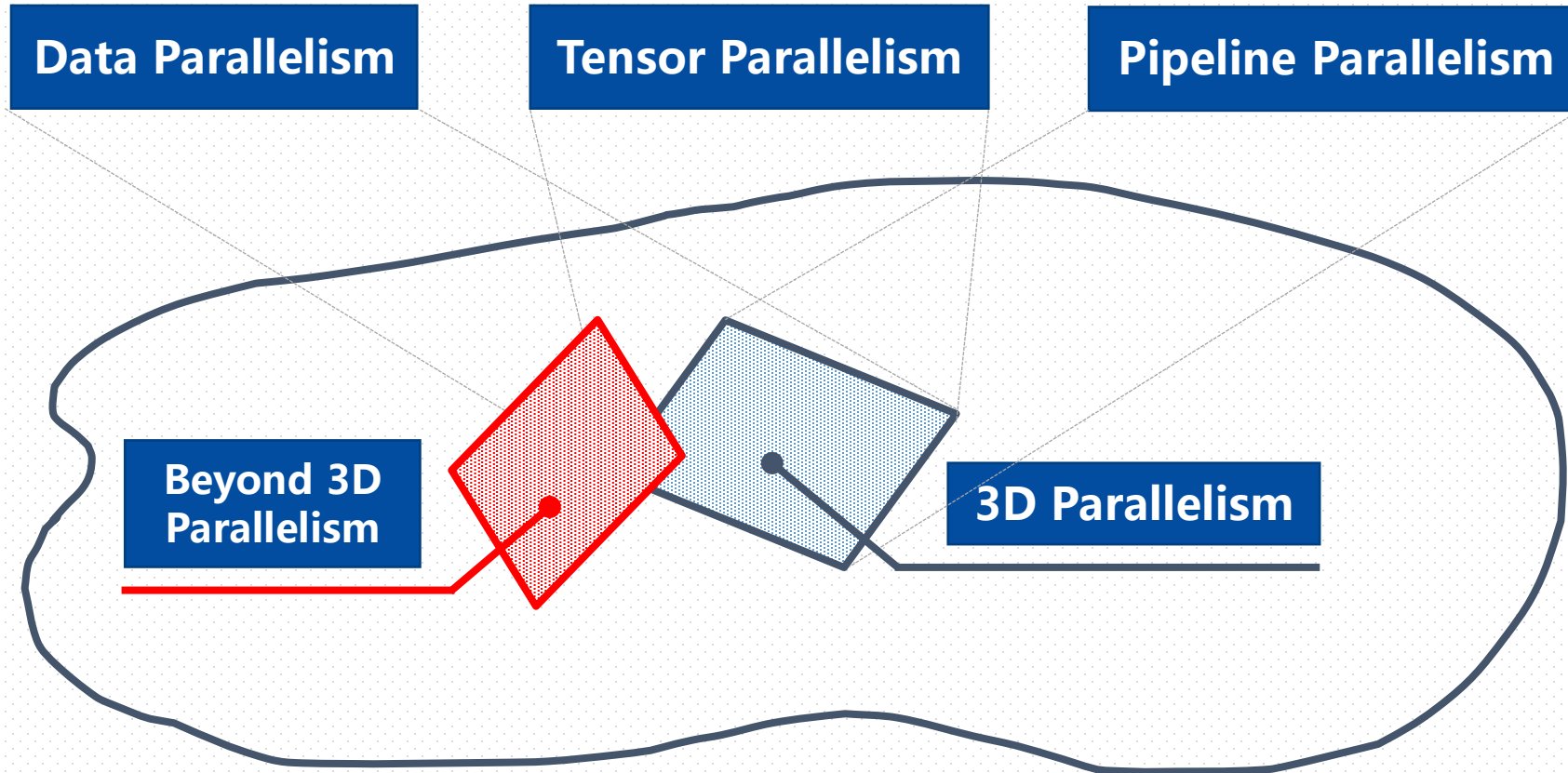
Given a model and GPUs

- ❑ Many operators
- ❑ Many partition choices
- ❑ Many placement choices
- ❑ Different execution orders for independent operators in one GPU

A combinatorial search space

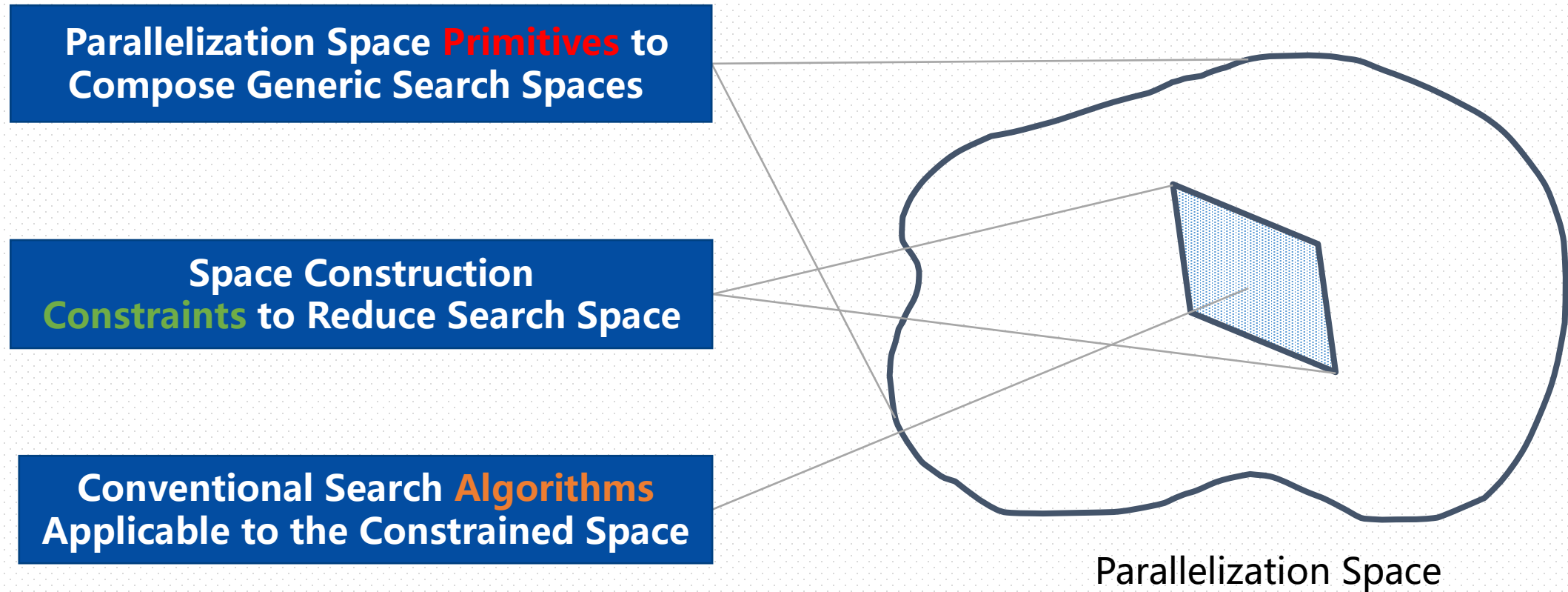
Worse when model and the cluster become larger

Large Space vs. Efficient Plan?



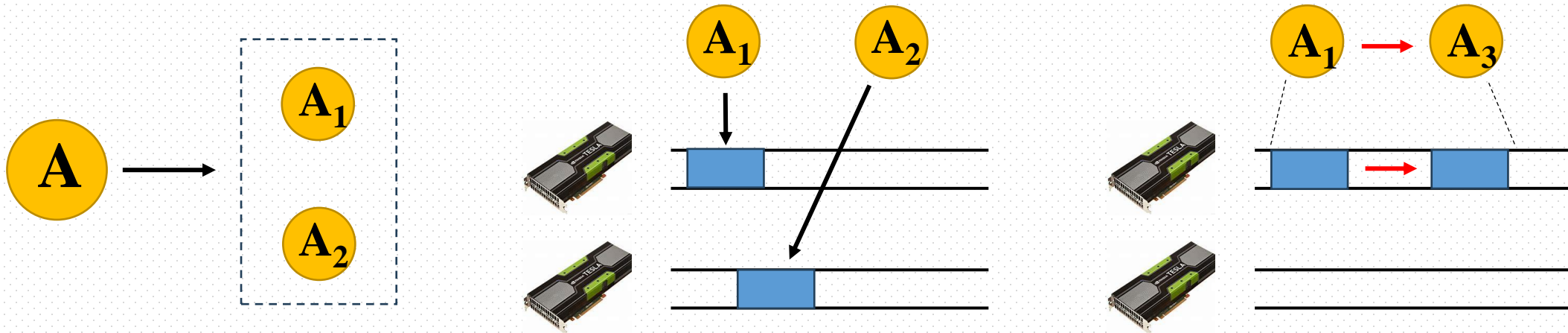
Current practice: limiting plan search within a well-studied search space

Solution Outline of nnScaler



Empowering developers to find their own plans

Space Construction Primitives



Transformation

$\text{op-trans}(\text{op}, \text{algo}, \text{num})$

Choices of operator partitioning schemes

Placement

$\text{op-assign}(\text{op}, \text{device})$

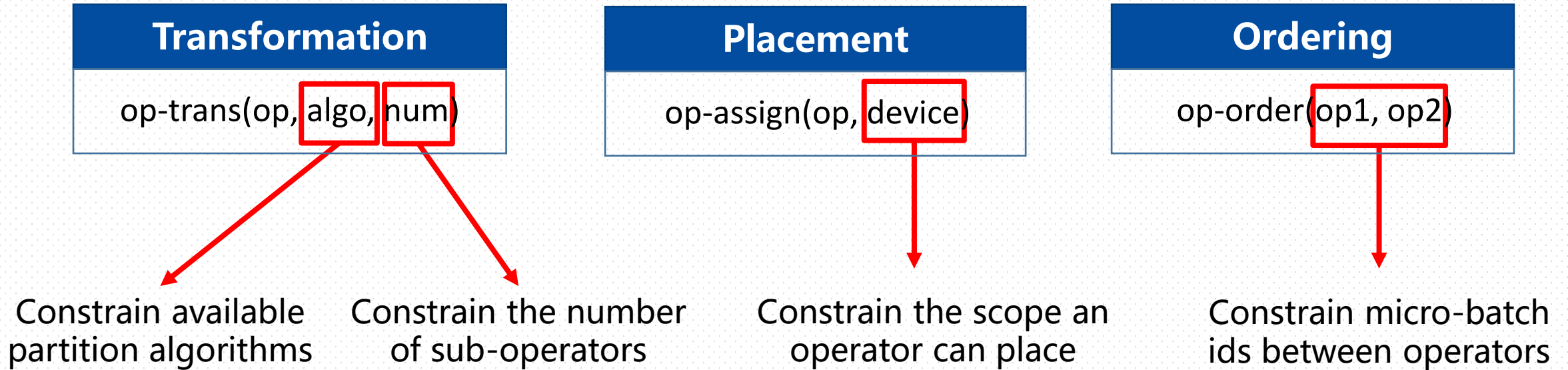
Assign an operator to a device

Ordering

$\text{op-order}(\text{op1}, \text{op2})$

Order independent operators in a device

Constraints in the Primitives



Constraints => reduce the search space

Constraint as a Powerful Abstraction

Constraints: the insights of domain experts

Primitives	Constraints
① $\text{sub_ops} = \text{op_trans}(\text{op}, \text{algo}, n)$	$n = \mathbf{D} $
② $\text{op_assign}(\text{sub_op}_i, d_i)$	$d_i, d_j \in \mathbf{D}$
③ $\text{op_assign}(\text{sub_op}_j, d_j)$	$d_i \neq d_j$

Constraints of data/tensor parallelism

Primitives	Constraints
① $\text{op_assign}(G_i, d_i)$	$d_i, d_j \in \mathbf{D}$
② $\text{op_assign}(G_j, d_j)$	$d_i \neq d_j$

Primitives	Constraints
① $\text{op_order}((fG_i, m), (fG_i, n))$	$m < n$
② $\text{op_order}((bG_i, m), (bG_i, n))$	
③ $\text{op_order}((fG_i, m+ofst), (bG_i, m))$	$ofst = \mathbf{D} - i$
④ $\text{op_order}((bG_i, m), (fG_i, m+ofst+1))$	$m \geq 0$

Constraints of pipeline parallelism (1F1B)

Existing Parallelism as Constraints

Operators	Primitives	Constraints
$\text{op} \in \{\text{Attn} \cup \text{FF}\}$	$\text{sub_ops} = \text{op_trans}(\text{op}, \text{algo}, n)$	$n = C \cdot \mathbf{D}_i $
	$\text{op_assign}(\text{sub_op}_i^j, d_i)$	$0 \leq j < C $ $d_i \in \mathbf{D}_i$

CoShare

Primitives	Constraints
① $\text{op_order}((f_1G_i, m+2), (f_2G_i, m+1))$	$m \geq 0$
② $\text{op_order}((f_2G_i, m+1), (f_3G_i, m))$	$m > 0$
③ $\text{op_order}((f_3G_i, m), (bG_i, m-ofst))$	$m > ofst$

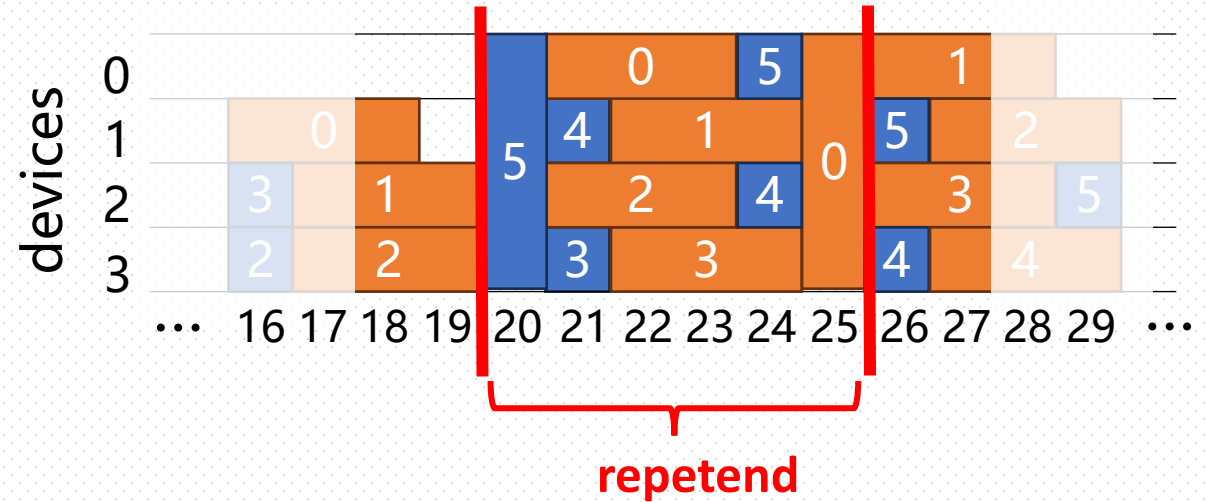
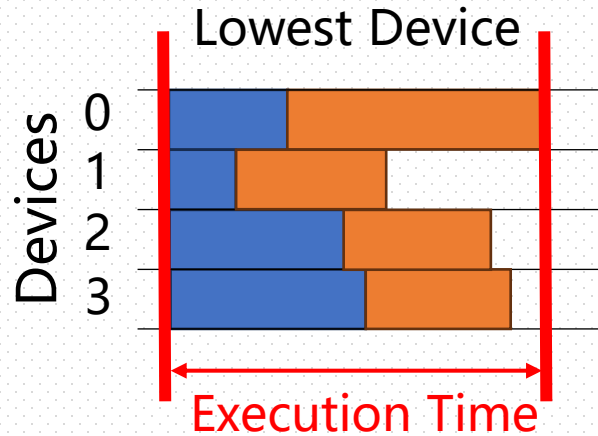
3F1B

Operators	Primitives	Constraints
$\text{op} \in \mathbf{E}$	$\text{sub_ops} = \text{op_trans}(\text{op}, \text{algo}, n)$	$n = \mathbf{D} $
	$\text{op_assign}(\text{sub_op}_i, d_i)$	$d_i \in \mathbf{D}$
$\text{ops} \notin \mathbf{E}$	$\text{staged_spmd}(\text{ops}, \mathbf{D})$	$0 \leq i < D $

Interlaced Pipeline

New Constraints for Emerging Models

Plan Search Policy



$$\text{minimize } \max_{d \in D} \sum_{op \in d_{op}} Comp_{op} + Comm_{op}$$

Transformation & Placement

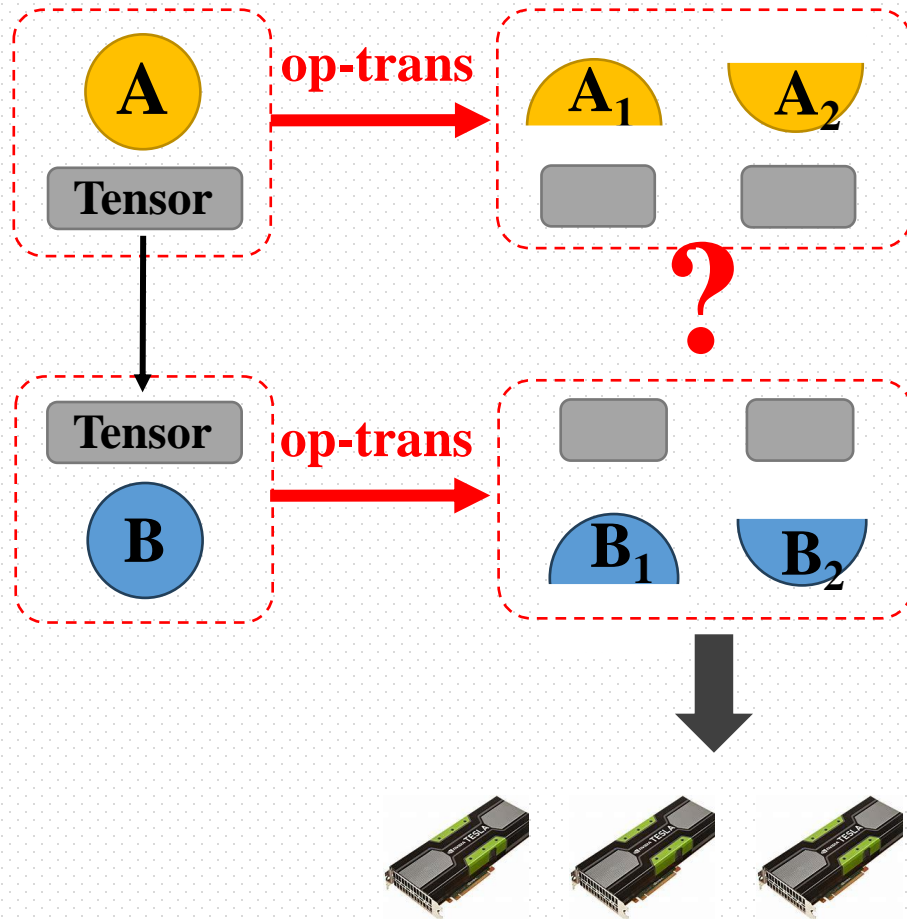
single micro-batch
minimize device execution time
 <DP, ILP> search



Temporal Ordering

multiple micro-batches
maximize device utilization
 <Tessel> search

Plan Materialization for Execution

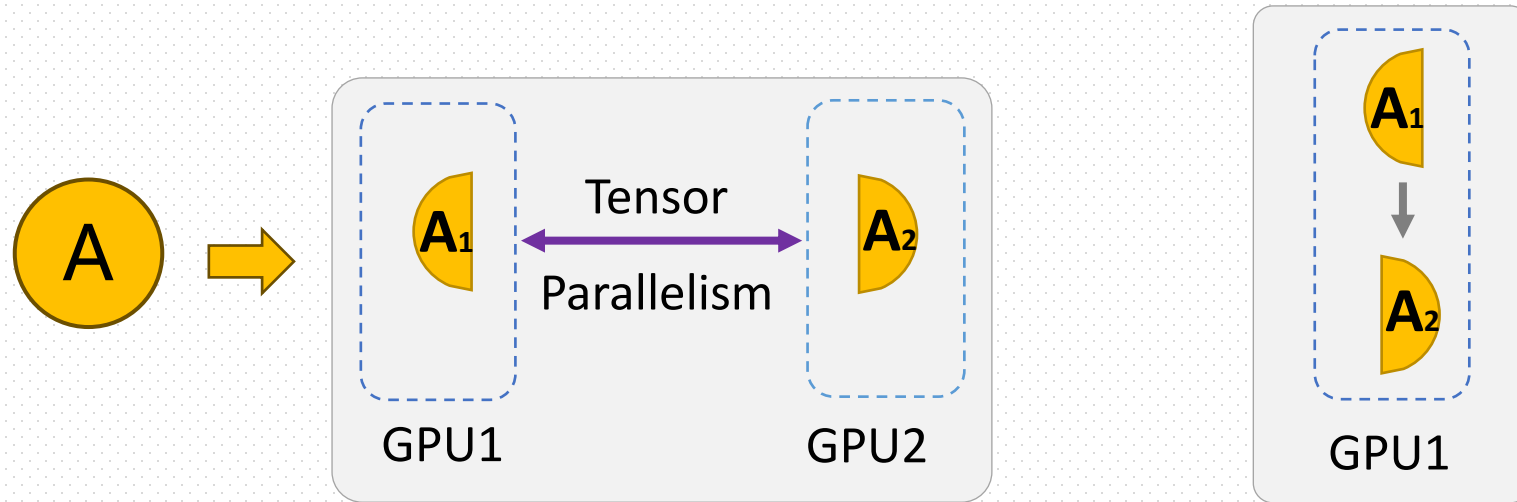


Additional Considerations

- **Tensor lineage** during transformation
- Efficient **communications** for equivalency
- Overall plan **correctness**
- Executable: PyTorch code

vTensor-pTensor abstraction (Section 6)

New Plan: CoShard



- **Reduced** peak memory usage
- **Lower** communication cost
- Beyond tensor parallelism

Tensor Parallelism

partitioned operators must be placed on **different** devices

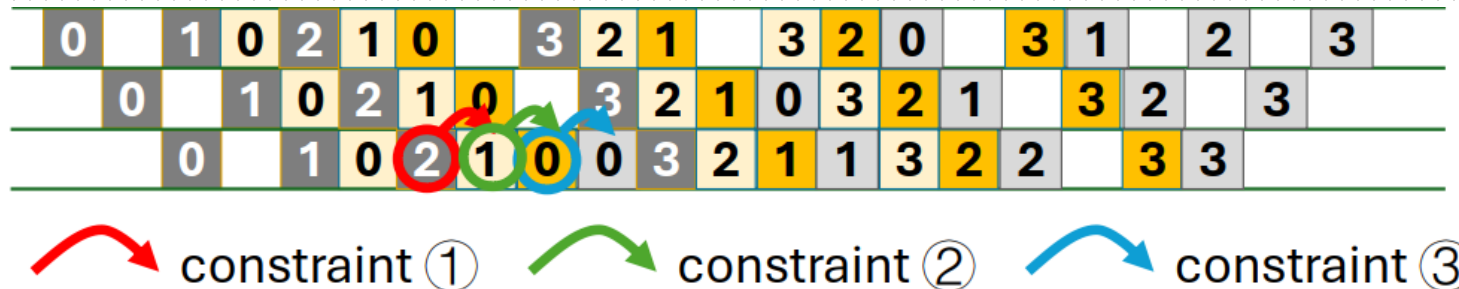
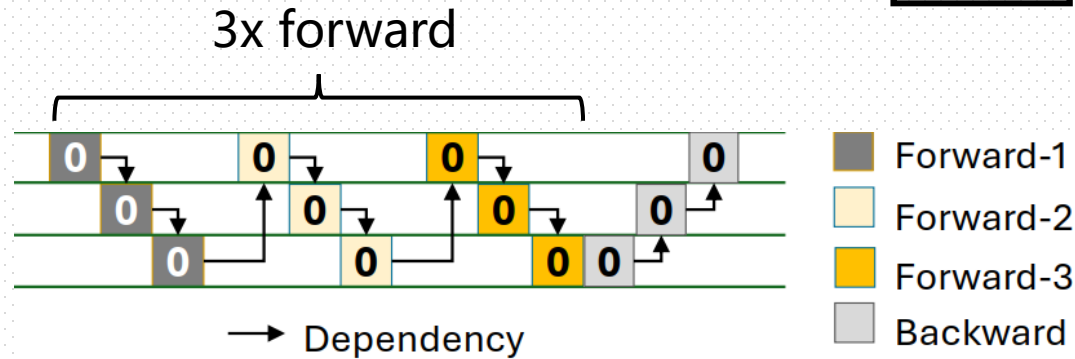
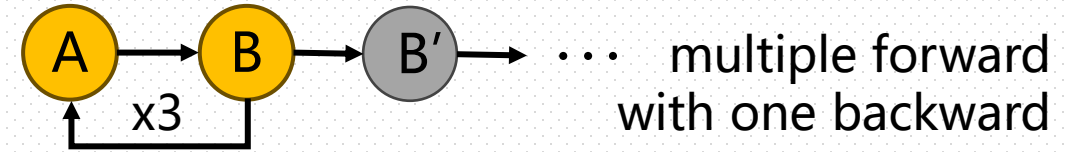
Coshard

partitioned operators can be **co-located** on a same device

Operators	Primitives	Constraints
$op \in \{\text{Attn} \cup \text{FF}\}$	$\text{sub_ops} = \text{op-trans}(op, \text{algo}, n)$	$n = C \cdot \mathbf{D}_i $
	$\text{op-assign}(\text{sub_op}_i^j, d_i)$	$0 \leq j < C $ $d_i \in \mathbf{D}_i$

New Plan: 3F1B

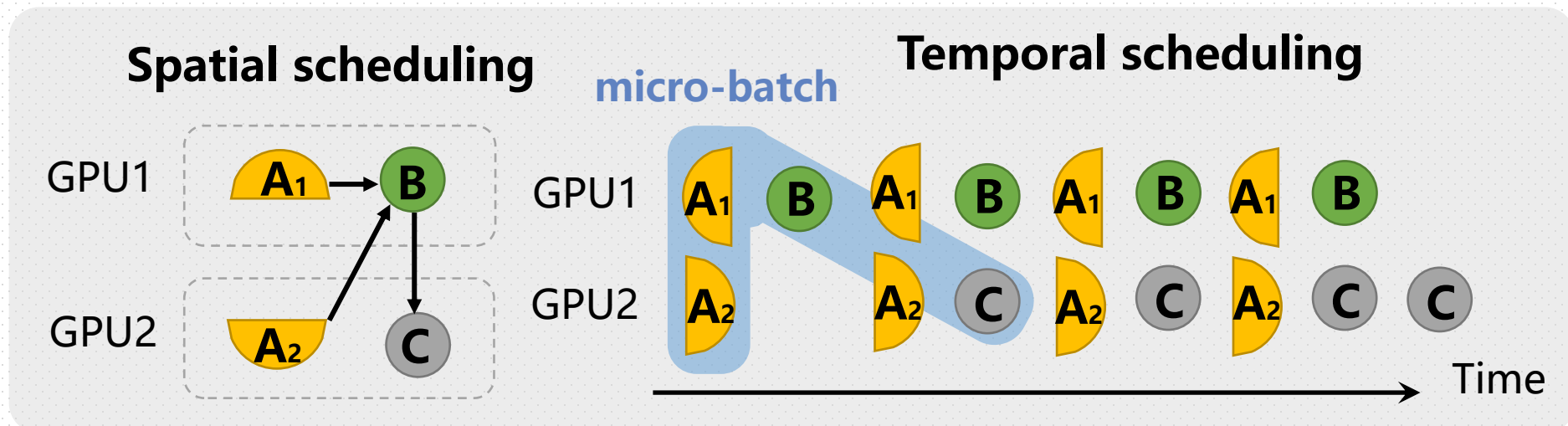
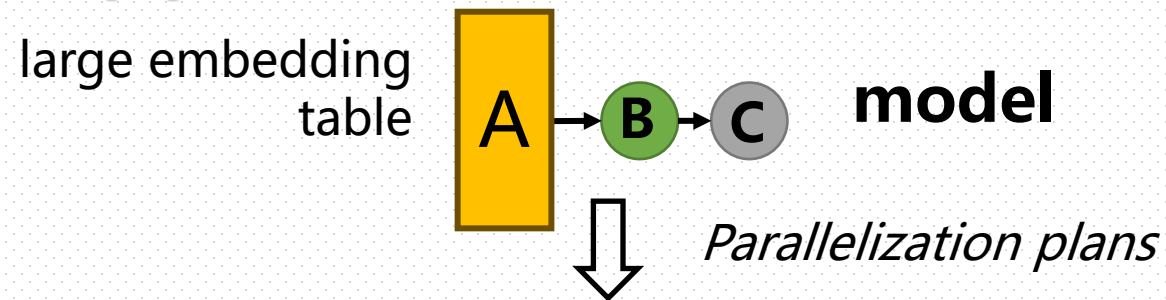
- AlphaFold: 3F1B Parallelization Plan**



Primitives	Constraints
① $op_order((f_1G_i, m+2), (f_2G_i, m+1))$	$m \geq 0$
② $op_order((f_2G_i, m+1), (f_3G_i, m))$	$m > 0$
③ $op_order((f_3G_i, m), (bG_i, m-ofst))$	$m > ofst$

New Plan: Interlaced Pipeline

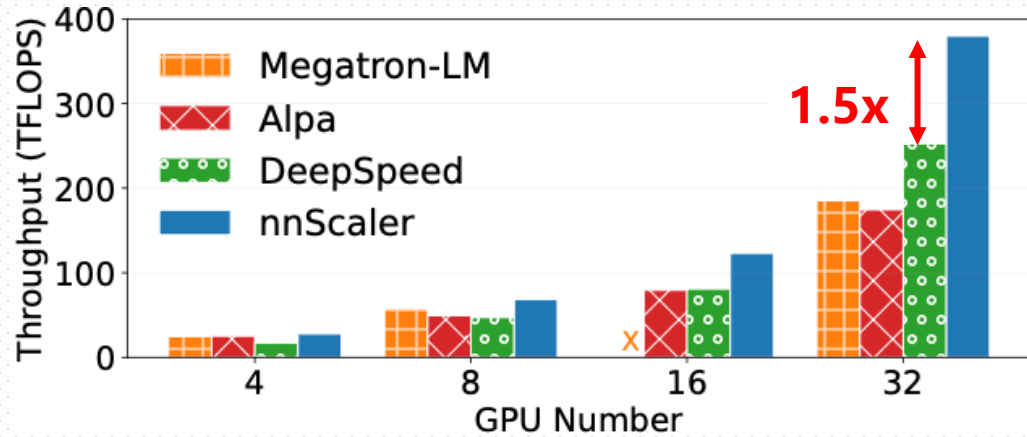
- T5: Interlaced pipeline**



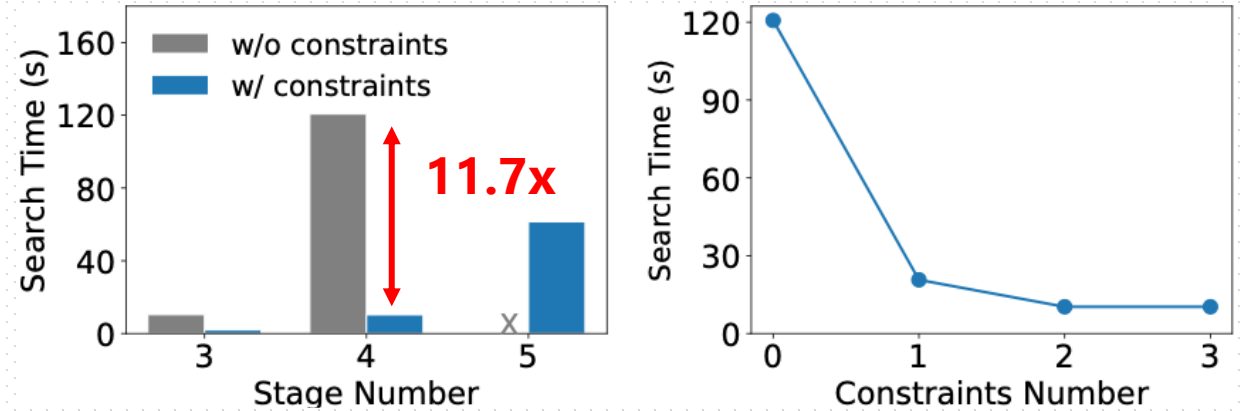
**full-device tensor parallelism
for embedding**

Operators	Primitives	Constraints
$op \in \mathbf{E}$	$sub_ops =$ $op_trans(op, algo, n)$ $op_assign(sub_op_i, d_i)$	$n = \mathbf{D} $ $d_i \in \mathbf{D}$ $0 \leq i < \mathbf{D} $
$ops \notin \mathbf{E}$	$staged_spmd(ops, \mathbf{D})$	

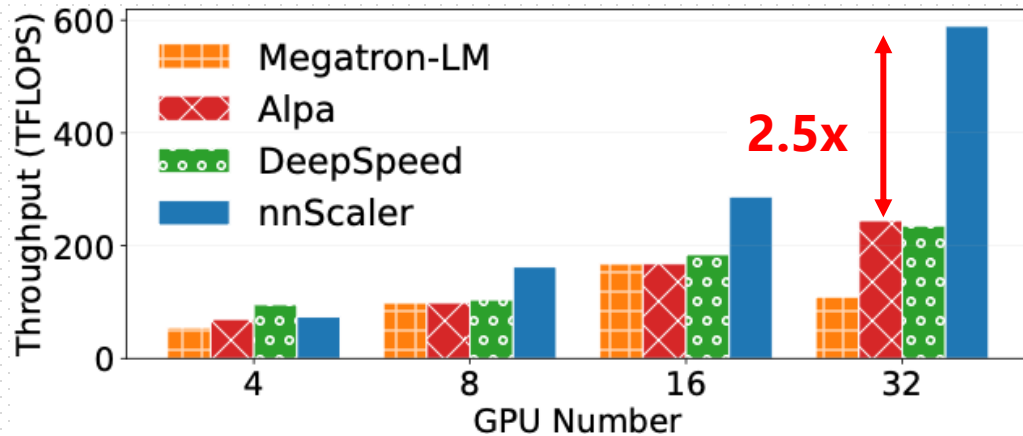
Evaluation: Performance



Swin-Transformer Model



AlphaFold 2 Search Time

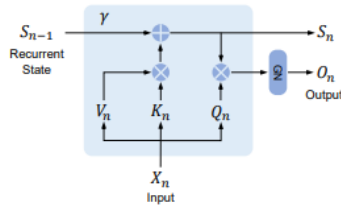


T5 Model

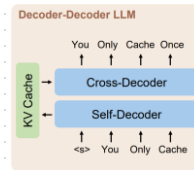
- DGX-2 32x V100-32GB
- Training Throughput Improvement:
 - **1.5-2.5x**
- Search Speed Improvement with Constraints:
 - **11.7x**

nnScaler in Practice

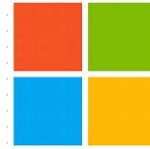
PreTrain and PostTrain (Finetune)



RetNet



YOCO



Phi-3



Graphormer



Llama2

...



DGX-2, A100, H100



AMD

MI-200

Deployment

512x GPUs (NV, AMD)
92B large model

Summary

The primitives and constraints:

abstractions for training flexibility and efficiency

nnScaler:

a powerful tool to facilitate model training



nnScaler: Constraint-Guided Parallelization Plan Generation for Deep Learning Training

Thank you for listening
Q & A

