

TAP: Efficient Derivation of Tensor Parallel Plans for Large Neural Networks



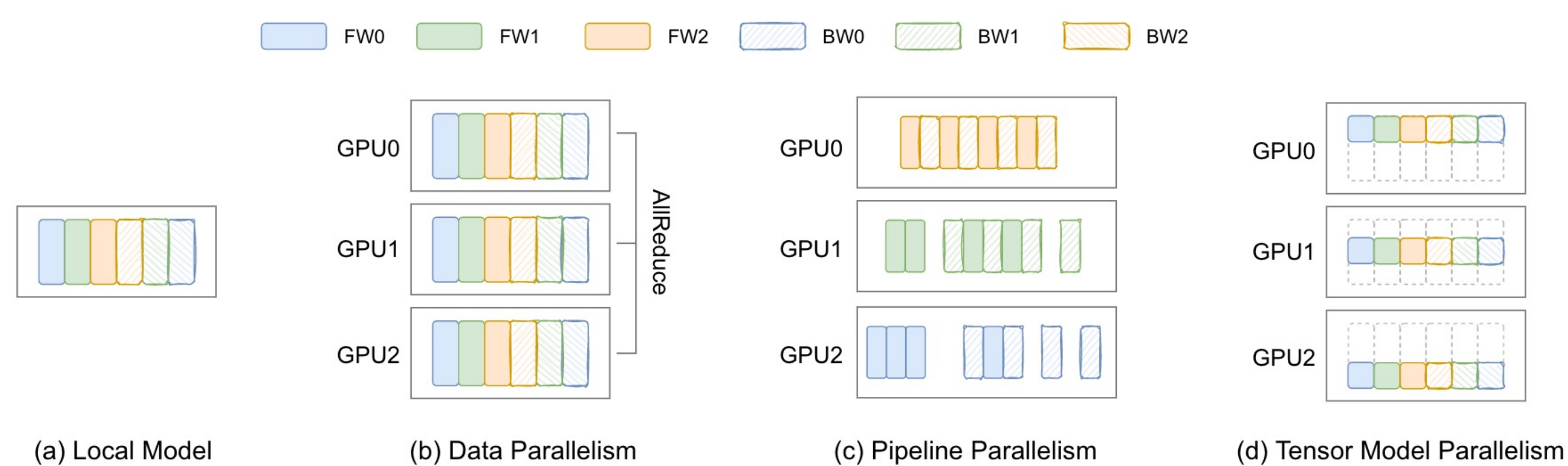
Ziji Shi^{1,2}, Le Jiang², Jie Zhang², Xianyan Jia², Yong Li²,
Chencan Wu², Jialin Li¹, Wei Lin²
¹National University of Singapore, ²Alibaba Group



Challenge: Memory Wall and Tensor Parallelism

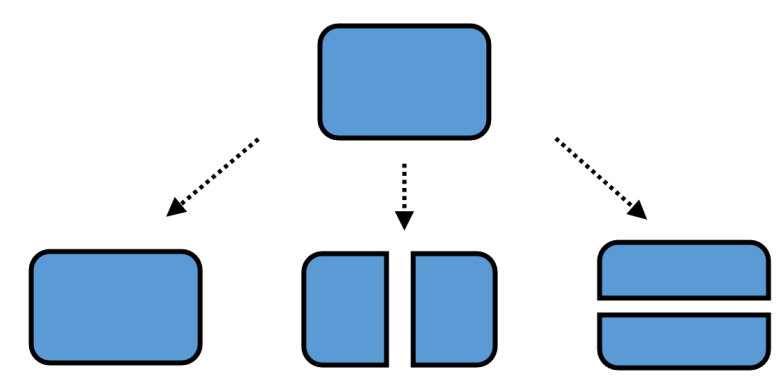
a. Memory wall for large neural networks training

- In the last decade, the model size has increased ~240X every two years, while the GPU memory has only doubled within the same amount of time.
- The memory wall problem has drawn much attention on model parallelism, where the model weights are sharded.



b. Tensor parallelism

- Split the tensor and distribute on different devices
- A more general approach for training large models
 - When a **gigantic layer** cannot be fitted into accelerator memory
 - When pipeline parallelism cannot work well due to **imbalanced pipeline**
- Problem
 - The size of decision space is too large**
 - Each 2D tensor has 3 possible choices, thus brute force method is $O(3^N)$



- A neural network can have thousands to millions of tensors, taking several months to find the best plan
- Expert annotation can be difficult, stiff, and error-prone**
 - Requires deep understanding on both system and neural network
 - Annotations need to be updated when system/model changes
 - Incorrect annotations may result in training halt or even failure

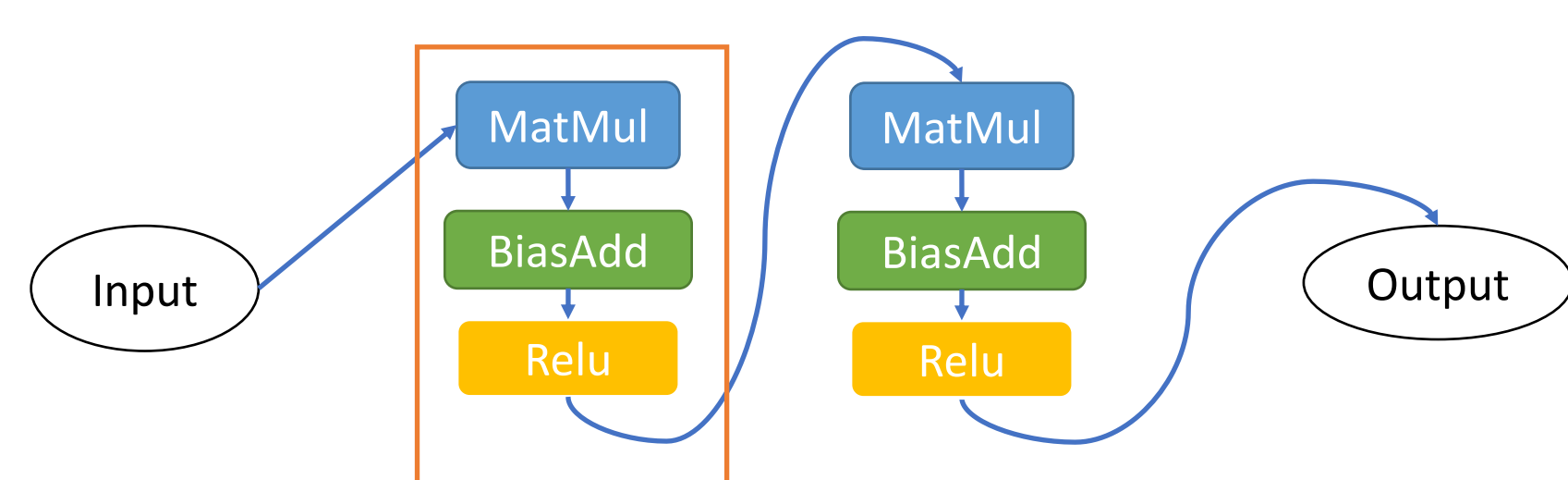
Approach: Tensor Auto Parallelism (TAP)

a. Related work

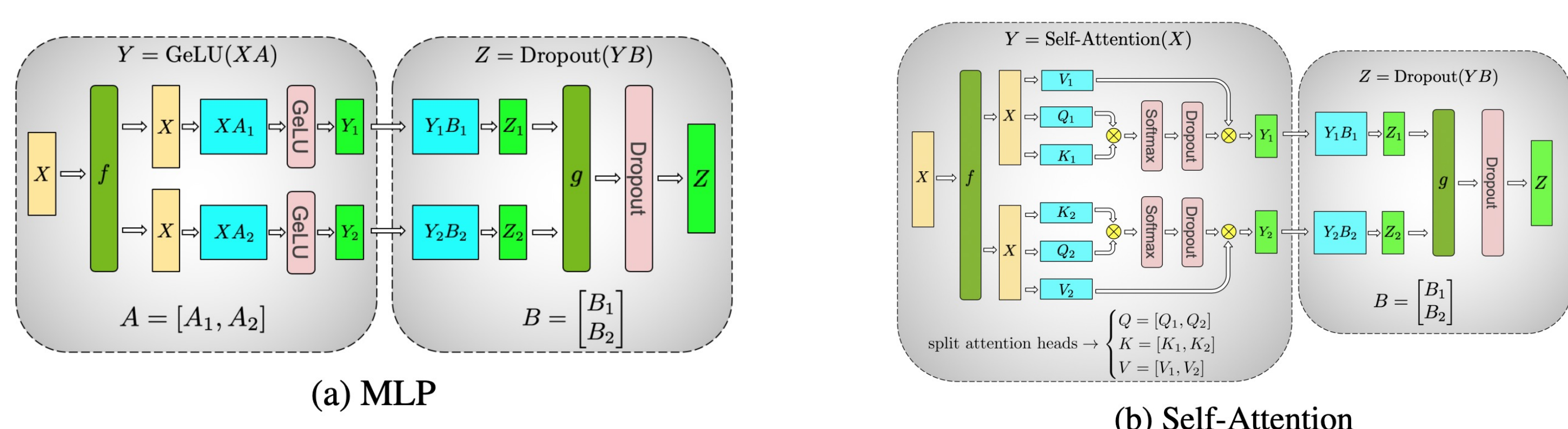
- Expert-annotation driven
 - MeshTensorFlow(NeurlPS'18), GSPMD (2021), Whale (ATC'22)...
- Automatic Parallelism
 - FlexFlow(ICML'18), Tofu(EuroSys'18), Unity(OSDI'22), Alpa(OSDI'22)...

b. Key observations

- A neural network can be represented as a directed acyclic graph, within which only contains a limited set of unique subgraphs.
 - Layers (eg. dense, self-attention)
 - Composite operators (eg. softmax)

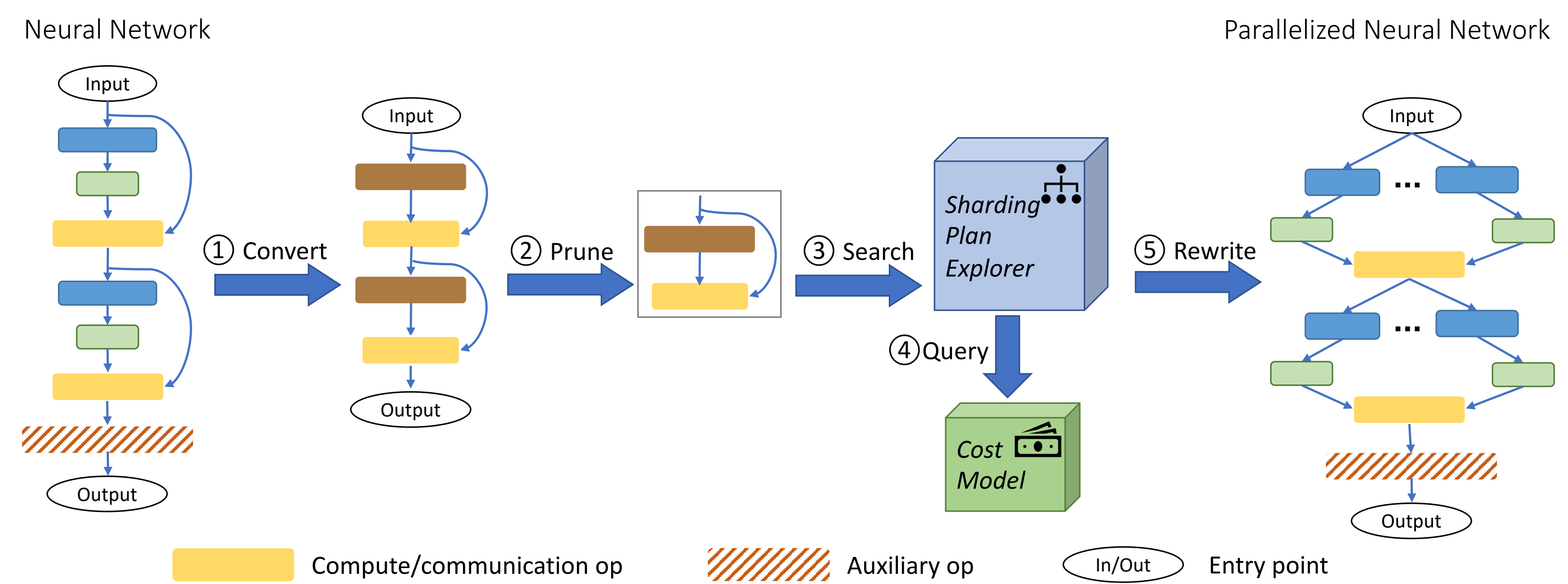


- The parallel strategy for similar layers are similar
 - For instance, Megatron-LM adopts a similar strategy for the Transformer layers



Credit: Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism (2019)

Design Overview



a. Workflow

- Convert** the original graph to a coarser graph representation
 - Group operators serving a similar functionality
 - Blue and yellow nodes were fused into a brown node
- Prune**: Discover unique subgraphs by gradually expanding the smaller subgraphs

```

Algorithm 1 Graph Pruning
1: procedure PRUNEGRAPH(modelDef, minDuplicate)
2:   nodeTree ← ∅
3:   maxDepth ← modelDef.depth
4:   for all depth ∈ maxDepth...1 do
5:     nodeTree[depth] ← longestCommonPrefix(modelDef.nodes.name)
6:     opCount = findSimilarBlk(nodeTree[depth])
7:     if opCount ≥ minDuplicate then
8:       subgraphs.append(nodeTree[depth])
9:     else
10:      break
11:   end if
12: end for
13: return subgraphs
14: end procedure
    
```

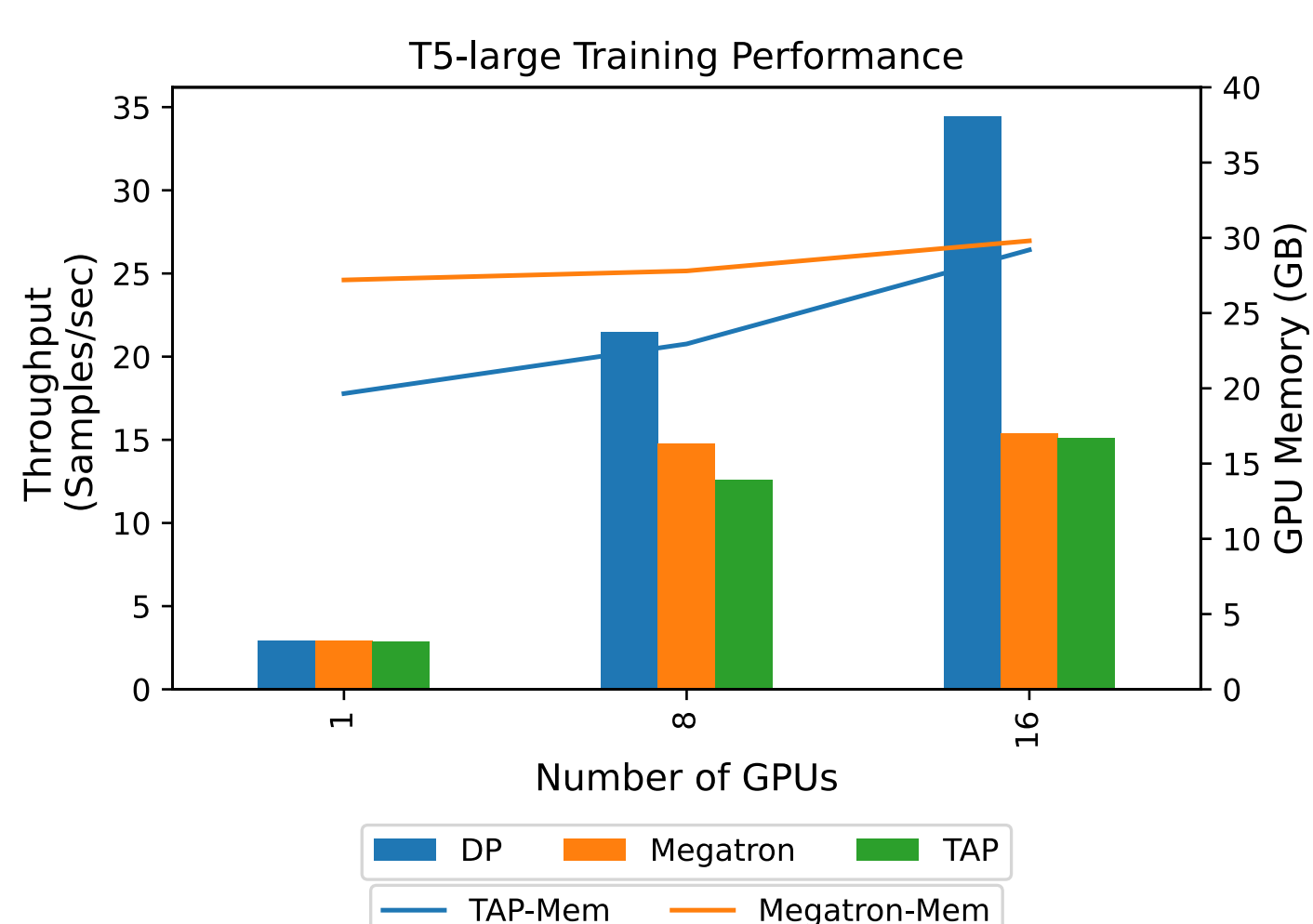
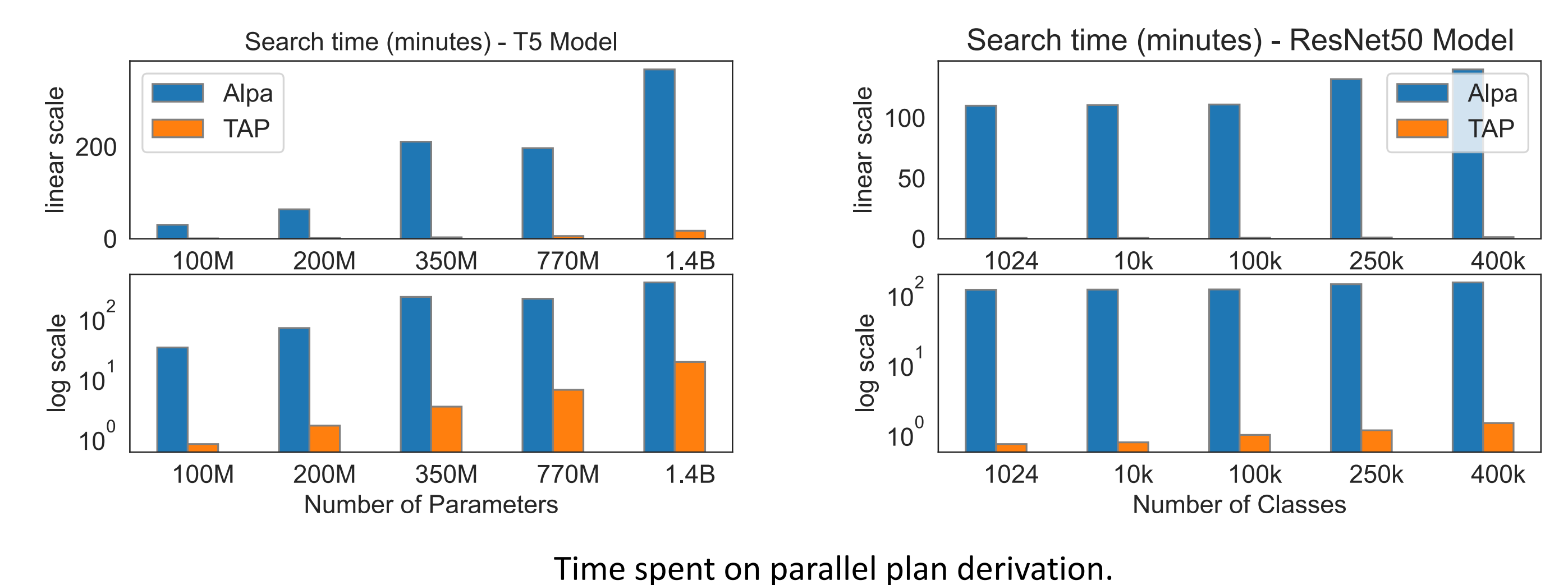
- Search** within the subgraphs by enumerating all possible tensor parallel plans
- Query** the analytical cost model for all plans and find the one that minimize communication
- Validate and **rewrite** the plan into original graph format

b. Usage

```

1. Example with TAP on 2 workers each with 8 GPUs
import tensor_auto_parallel as tap
mesh = [2, 8]
tap.auto_parallel(tap.split(mesh))
model_def()
    
```

Evaluations



Performance comparison with expert-designed Megatron-LM.

- 20-160X faster** in plan derivation compared to SoTA Auto Parallel Framework (Alpa)
- Discovered strategy has **comparable performance to Megatron**
- TAP also discovers new partially sharded plan that perform well when the memory is not so constrained.

Conclusion

We present TAP, an automatic parallelism framework that efficiently discovers tensor parallel plans for large models. Leveraging the observation that subgraphs widely exist in neural networks, we design a system that runs at sub-linear end-to-end complexity w.r.t. to model size.

