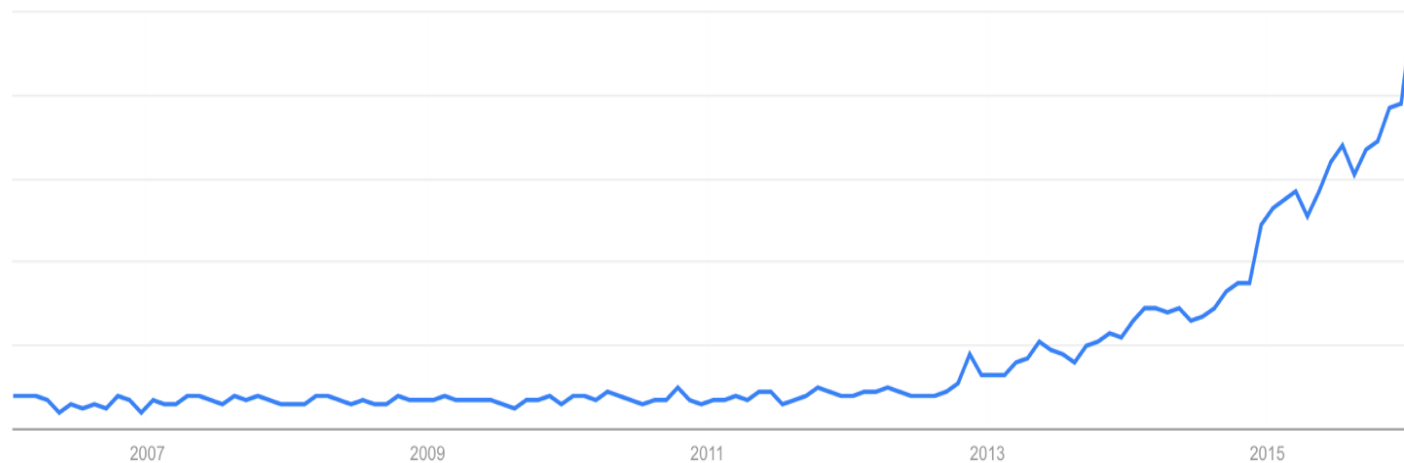# Tofu: Parallelizing Deep Learning Systems with Automatic Tiling

## Minjie Wang

# Deep Learning



"Deep Learning" trend in the past 10 years
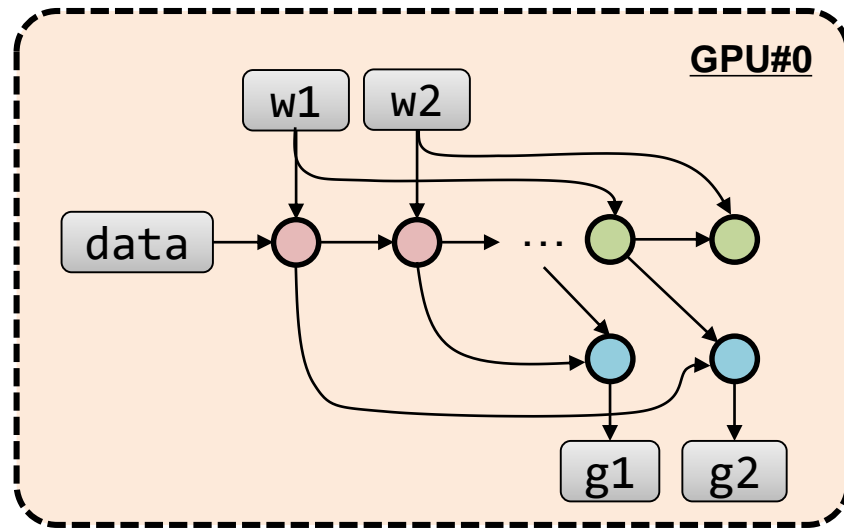
# State-of-art DL system is based on dataflow

```python
import tensorflow as tf
...  # generate data and weight
act1 = tf.matmult(data, w1)
act2 = tf.matmult(act1, w2)
...
grad_act2 = tf.matmult(w3.T, grad_act3)
grad_act1 = tf.matmult(w2.T, grad_act2)
...
grad_w2 = tf.matmult(act1.T, grad_act2)
grad_w1 = tf.matmult(data.T, grad_act1)
...  # update weights using gradients
```
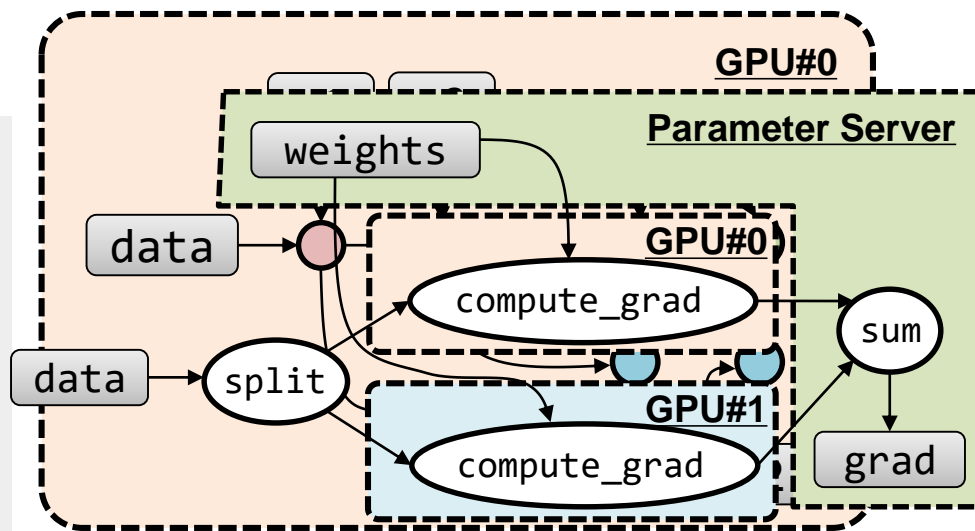
⬤ *Forward propagation*

⬤ *Backward propagation (input gradients)*

⬤ *Backward propagation (weight gradients)*

# What if I have many GPUs?

# Data parallelism with manual distribution

```python
import tensorflow as tf
...    # generate data and weight
data1, data2 = tf.split(data, axis=0)
with tf.device('/gpu:0'):
    grad1 = compute_grad(data1, weights)
with tf.device('/gpu:1'):
    grad2 = compute_grad(data2, weights)
with tf.device('/ps'):
    grad = aggregate(grad1, grad2)
    ...    # update weights using gradients
...    # update weights using gradients
```
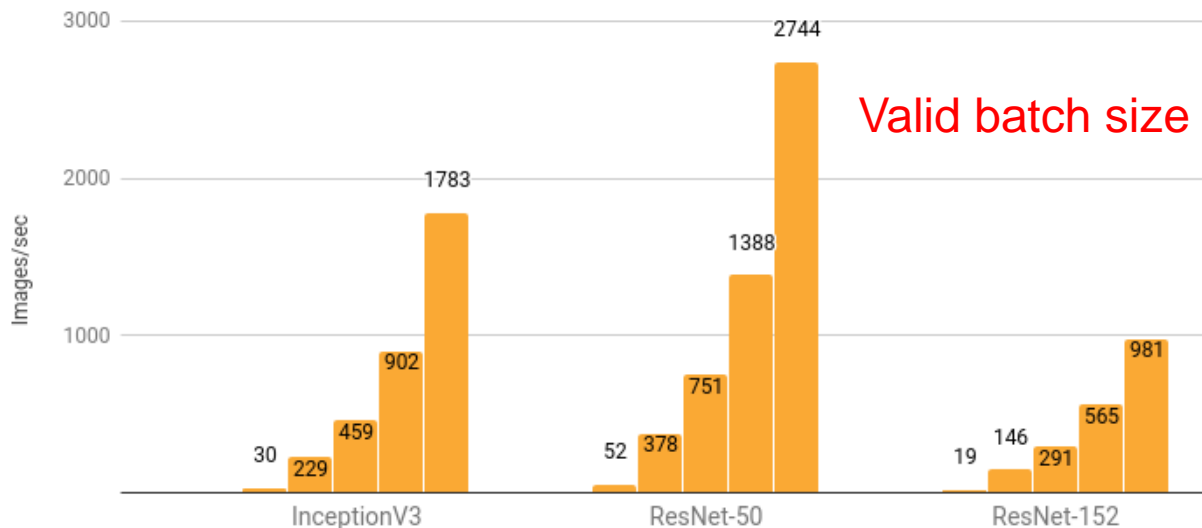


**Manual Distribution &
Device assignment**

# Scalability secret of data parallelism

Training: NVIDIA® Tesla® K80 synthetic data (1,8,16,32, and 64)

Valid batch size = 64 * 64 = 4096



| Options | InceptionV3 | ResNet-50 | ResNet-152 | Alexnet | VGG16 |
|---|---|---|---|---|---|
| Batch size per GPU | 64 | 64 | 64 | 512 | 64 |
| Optimizer | sgd | sgd | sgd | sgd | sgd |

* Numbers from https://www.tensorflow.org/performance/benchmarks
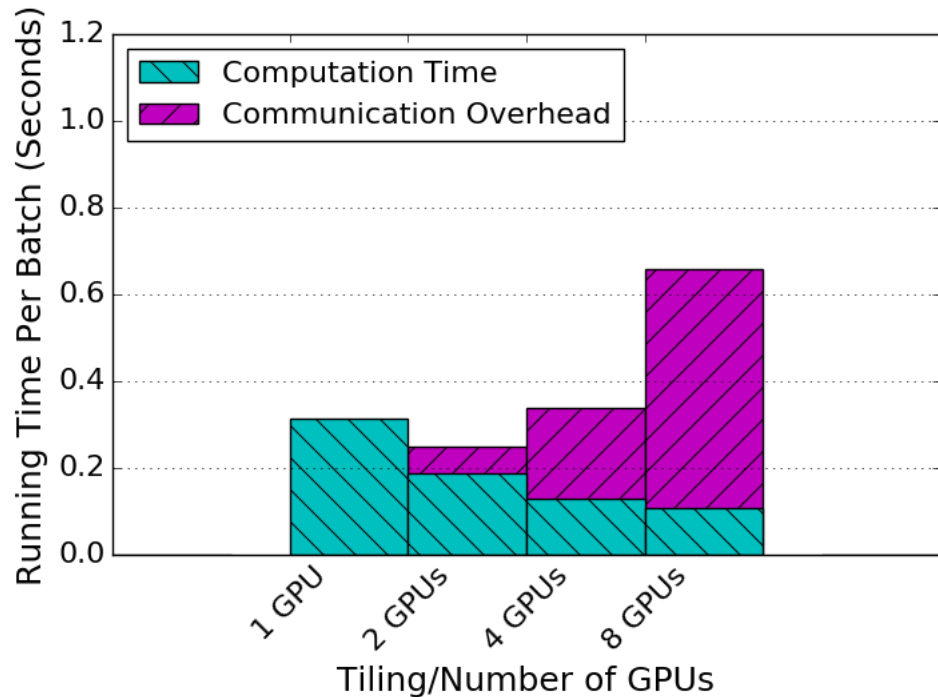
# Large batch size harms model accuracy



Inception Network on Cifar-10 dataset

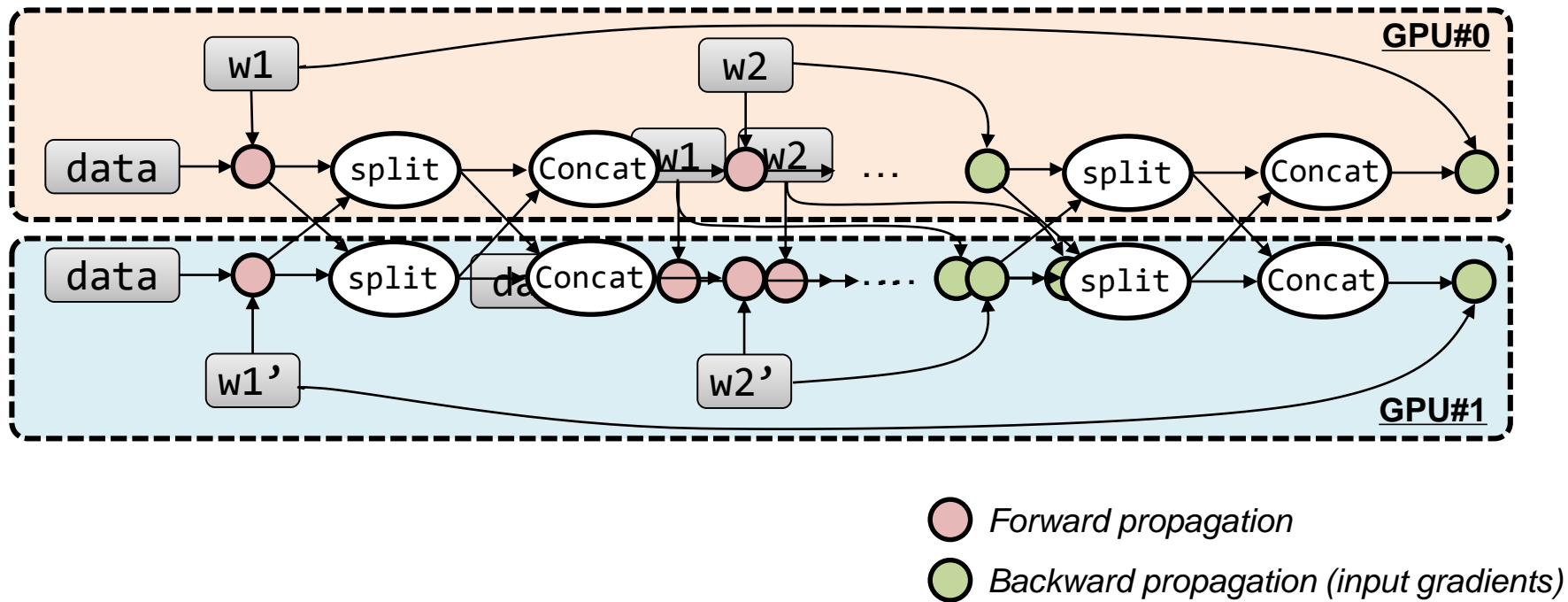# Data parallelism bottlenecked by communication



>80% of the total running time is for communication on 8 cards

5-layer MLP; Hidden Size = 8192; Batch Size = 512

# An alternative way: Model Parallelism

# MP is hard to program

```python
# Original MLP code.
def mlp(data, weights):
    # Forward Propagation.
    fwd = [data]
    for i in xrange(FLAGS.num_layers):
        fwd.append(tf.matmul(fwd[-1], weights[i])) # forward matmult
    # Backward Propagation.
    targets = []
    last = fwd[-1]
    for i in reversed(xrange(FLAGS.num_layers)):
        dw = tf.matmul(fwd[i], last, transpose_a=True) # matmult: grad
        last = tf.matmul(last, w[i], transpose_b=True) # matmult: bp
        # update
        targets.append(dw)
    return targets
```
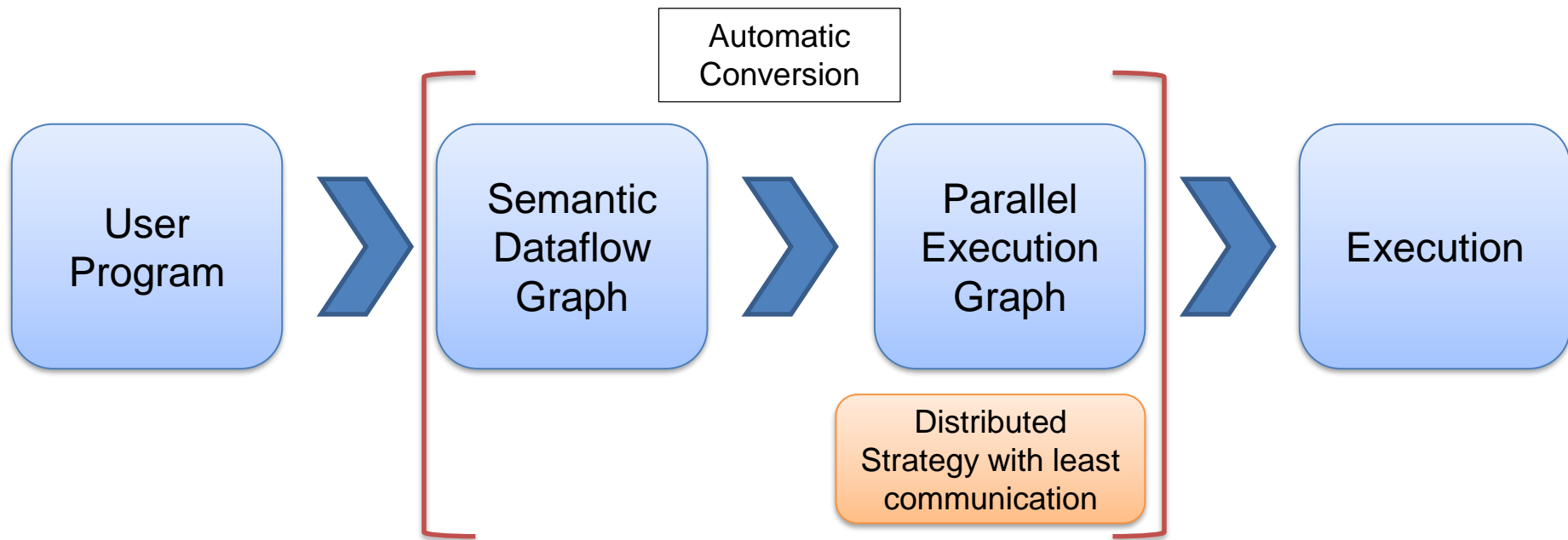
```python
# Manual Model Parallelism implementation for a MLP network.
def model_par_mlp(data, weights):
    # Partition weights on row.
    w = []
    for i in xrange(FLAGS.num_layers):
        w.append([])
        for j in xrange(FLAGS.num_workers):
            with tf.device('/job:worker/task:%d' % j):
                w[i].append(tf.get_variable(
                        name='w%d' % j,
                        shape=[slice_size,feature_size],
                        trainable=True))
    # Forward Propagation.
    fwd = []
    last = data
    for i in xrange(FLAGS.num_layers):
        with tf.name_scope('fc_ff%d' % i):
            fwd.append(last)
            tmp = []
            for j in xrange(FLAGS.num_workers):
                with tf.device('/job:worker/task:%d' % j):
                    y = tf.matmul(last[j], w[i][j])   # forward matmult
                    # split the result so we can do balanced reduction.
                    tmp.append(tf.split(split_dim=1, num_split=FLAGS.num_workers, value=y))
            # Reduce the result.
            red = []
            for j in xrange(FLAGS.num_workers):
                with tf.device('/job:worker/task:%d' % j):
                    red.append(tf.accumulate_n([s[j] for s in tmp]))
            last = red
    # Backward Propagation.
    targets = []
    for i in reversed(xrange(FLAGS.num_layers)):
        with tf.name_scope('fc_bp%d' % i):
            # Concatenate input tensors.
            tmp = []
            for j in xrange(FLAGS.num_workers):
                with tf.device('/job:worker/task:%d' % j):
                    tmp.append(tf.concat(concat_dim=1, values=last))
            last = []
            for j in xrange(FLAGS.num_workers):
                with tf.device('/job:worker/task:%d' % j):
                    dy = tf.matmul(tmp[j], w[i][j], transpose_b=True) # matmult: bp
                    last.append(dy)
                    dw = tf.matmul(fwd[i][j], tmp[j], transpose_a=True) # matmult: grad
                    targets.append(dw) # update
    return targets
```

# What is the best strategy for distribution?

- No one-size-fits-all
  - DP and MP suit different situations (parameter shapes, batch sizes).
  - Different layers might be suited for different strategies **(hybrid parallelism)**.
    - Use data parallelism for convolution layers; use model parallelism for fully-connected layers.
- DP and MP can be combined in a single layer
  - DistBelief (Dean, 2012)
  - Impossible to program with manual distributed strategy!

# Tofu automatically distributes DL training

| User Program | → | Semantic Dataflow Graph | → | Parallel Execution Graph | → | Execution |

Automatic Conversion

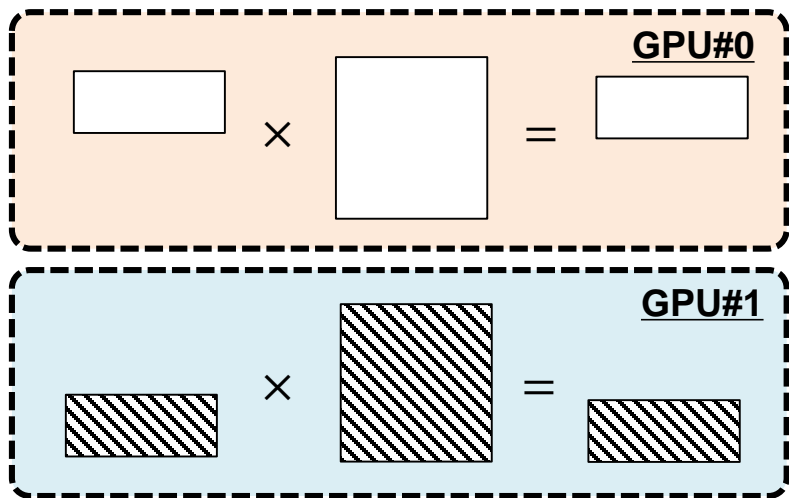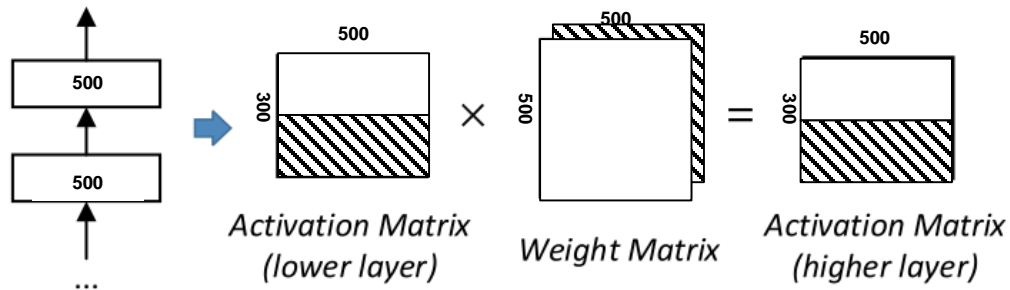Distributed Strategy with least communication

**Tofu**

# Challenges

- What are the different ways to distribute each tensor operator?
- What is the globally optimal way of distribution
  - that minimizes communication?
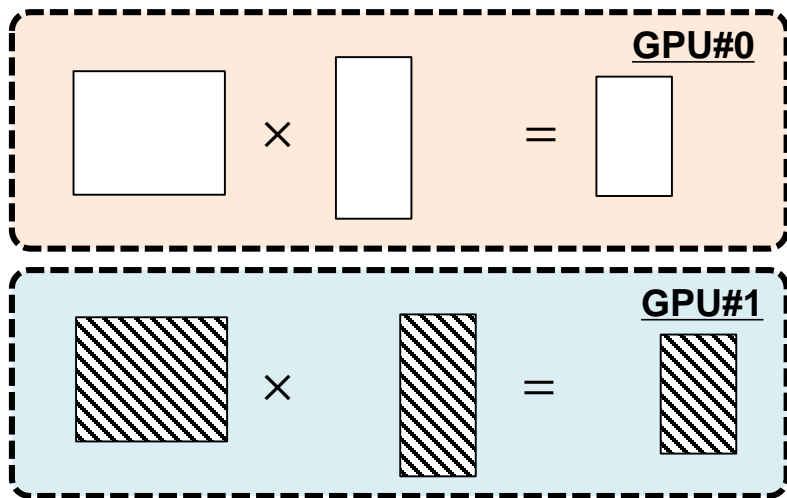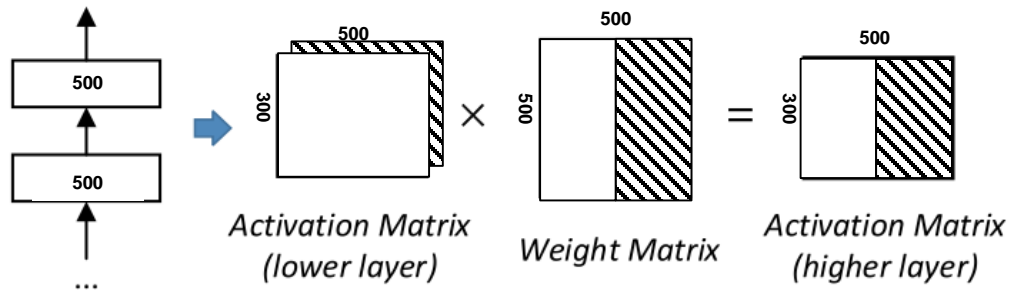
# Different ways of distributing matrix multiplication

Batch size: 300



*Activation Matrix (lower layer)* × *Weight Matrix* = *Activation Matrix (higher layer)*

GPU#0

GPU#1

➢ Activation Matrix (lower layer) is **row-partitioned**

➢ Weight Matrix is **replicated**

➢ Acitvation Matrix (higher layer) is **row-partitioned**

➢ Data parallelism

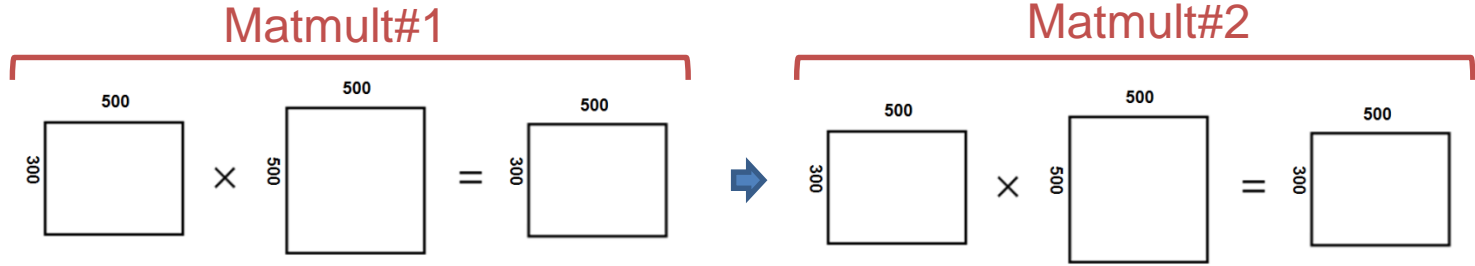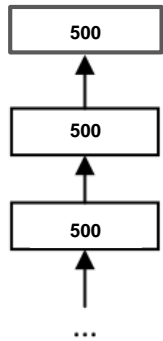# Different ways of distributing matrix multiplication

Batch size: 300



> Activation Matrix (lower layer) is **replicated**

> Weight Matrix is **column-partitioned**

> Acitvation Matrix (higher layer) is **column-partitioned**

> Model Parallelism

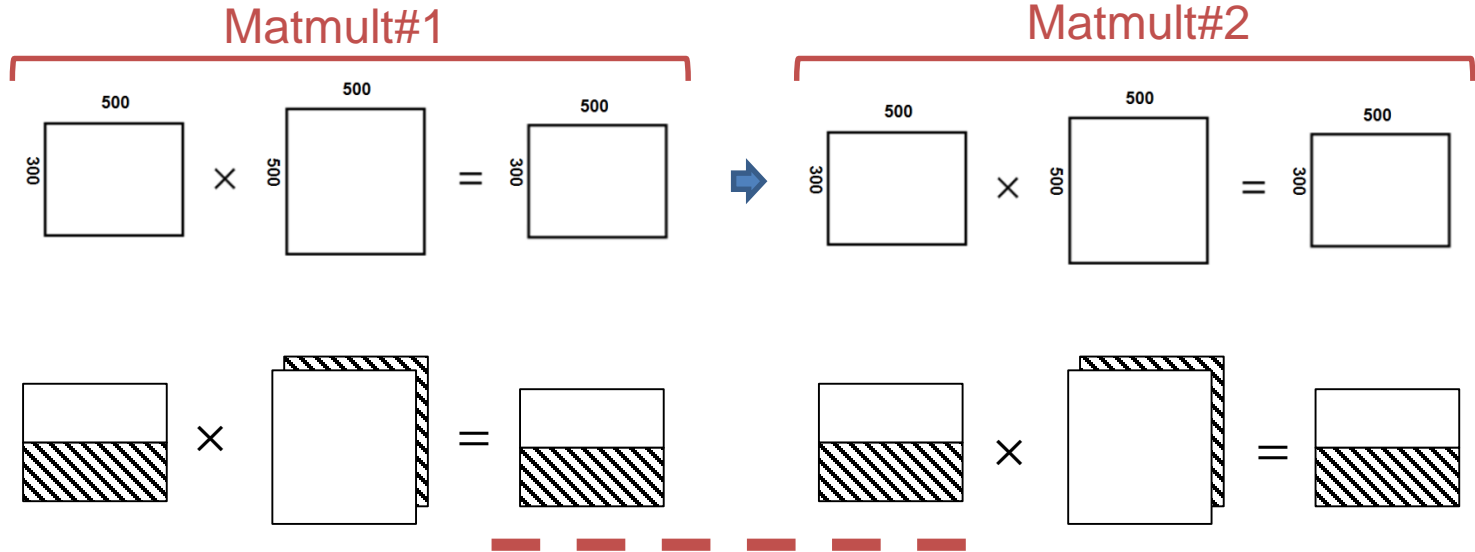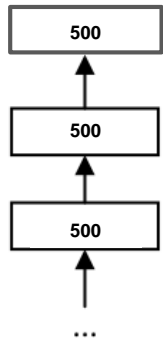# Operators can have different strategies

- Different matrix multiplications may choose different strategies.
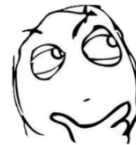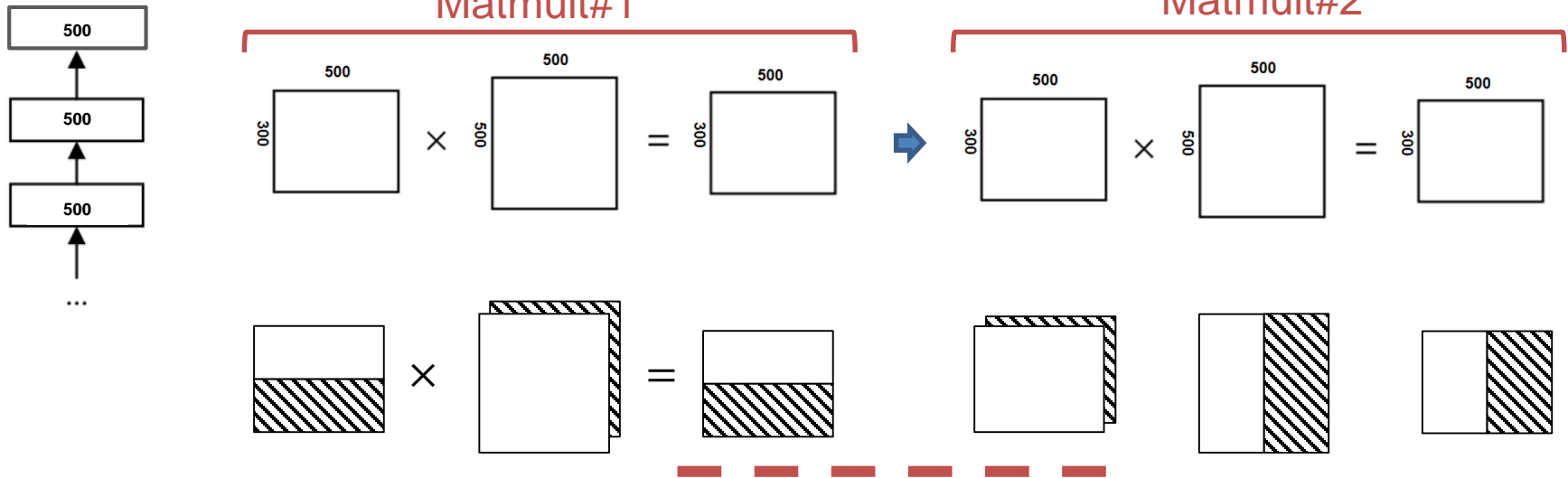
# Operators can have different strategies

- No communication if the output matrix satisfies the input partition.
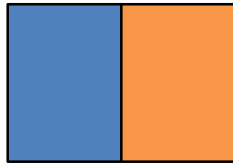


Matmult#1

Matmult#2

No Communication!

# Operators can have different strategies

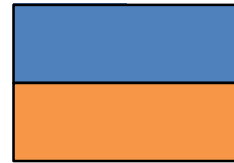- Communication happens when matrices need to be re-partitioned.

# Communication Cost

- Communication happens when matrices need to be re-partitioned.
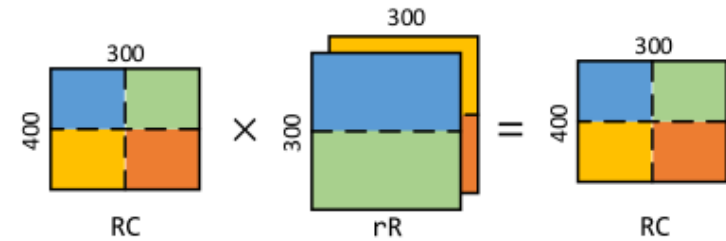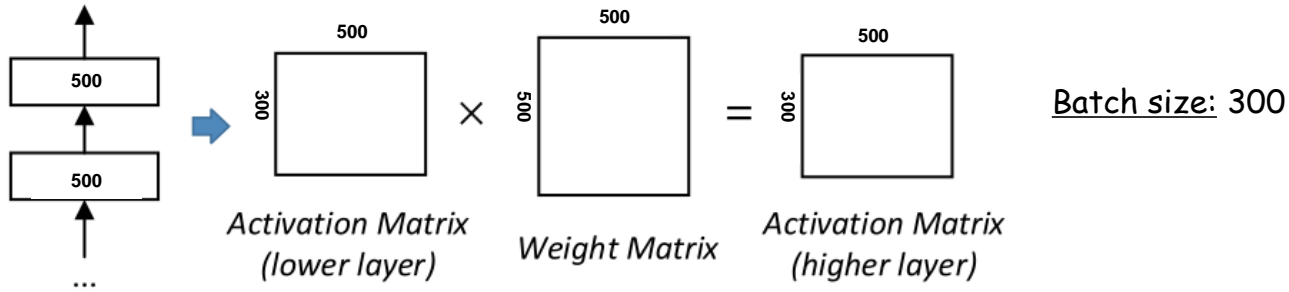- Communication cost == partition conversion cost.

C

R

# Finding optimal strategy with minimal communication

- Each operator has several distribution decisions.
  - DP and MP are one of them.
- Looking at one operator at a time is **not** optimal.
- Finding strategy with minimal communication cost for a general graph is NP-Complete.
- Tofu finds optimal strategy for deep learning in polynomial time:
  - "Layer-by-layer" propagations → graph with long diameter.
  - Use dynamic programming algorithm to find optimal strategy.

# Combined strategies for one operator



500

500
300 Activation Matrix (lower layer)

×

500
500 Weight Matrix

=

500
300 Activation Matrix (higher layer)

Batch size: 300

300
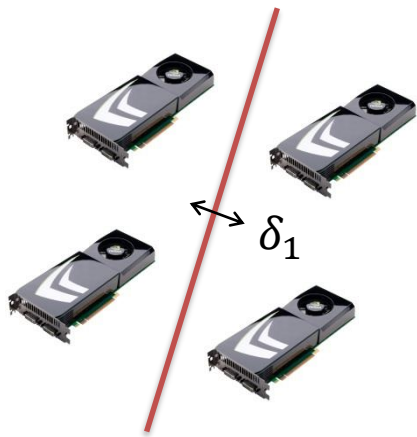400 RC

×

300
300 rR

=
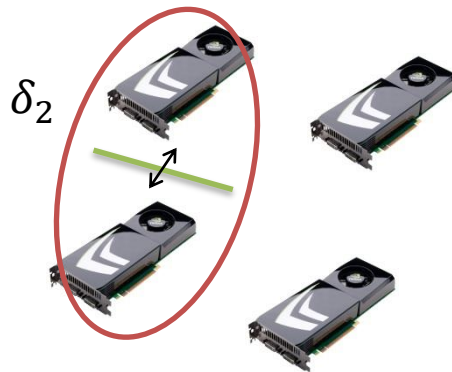
300
400 RC

# Combined strategy is sometimes better

- Fully-connected layer of 500 neurons with batch size 300.
- One combined strategy on 16 GPUs:
  - Model parallelism into 4 groups of GPUs (each group has 4 GPUs).
  - Data parallelism within each group.
  - Saves >33.3% communications than DP and MP.

# Find combined strategies

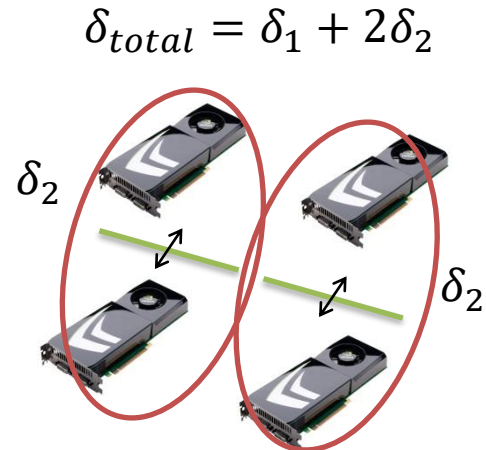- Solve the problem recursively.
- Proved to be optimal.



$$\delta_{total} = \delta_1 + 2\delta_2$$

Step 1: Partition to two groups

Step 2: Apply the algorithm again on one of the group

Step 3: Apply the same strategy to the other group due to symmetry.

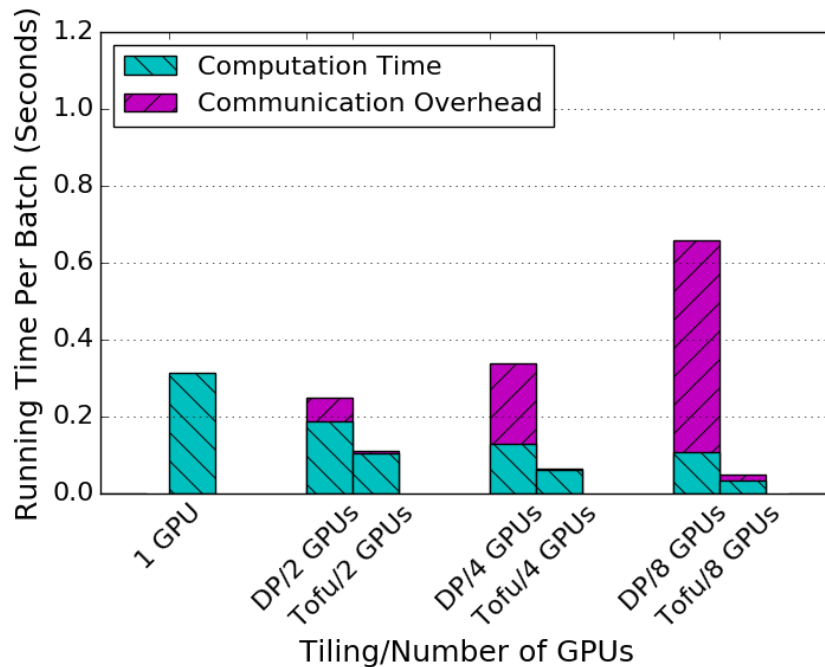# Tofu Evaluation Setup

- Implemented in MXNet's NNVM dataflow optimization library.

- Multi-GPU evaluation
  - Amazon p2.8xlarge instance
  - 8 NVIDIA GK210 GPUs (4 K80)
  - 12GB memory per card
  - Connected by PCI-e (160Gbps bandwidth)

Under submission. Contact wmjlyjemaine@gmail.com for more details.

# Communication Overhead Evaluation

- Per batch running time of a 4-layer MLP for DP and Tofu.
- Hidden layer size: 8192; Batch size: 512

# Real Deep Neural Networks Evaluation
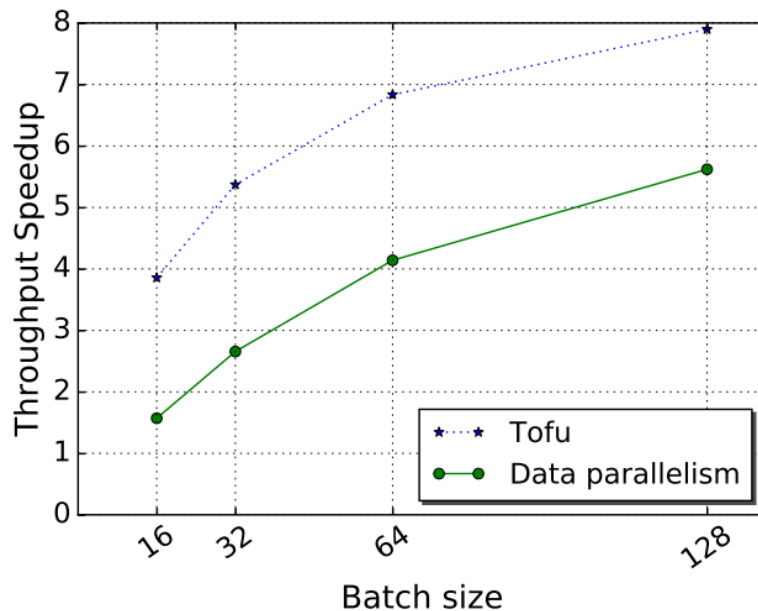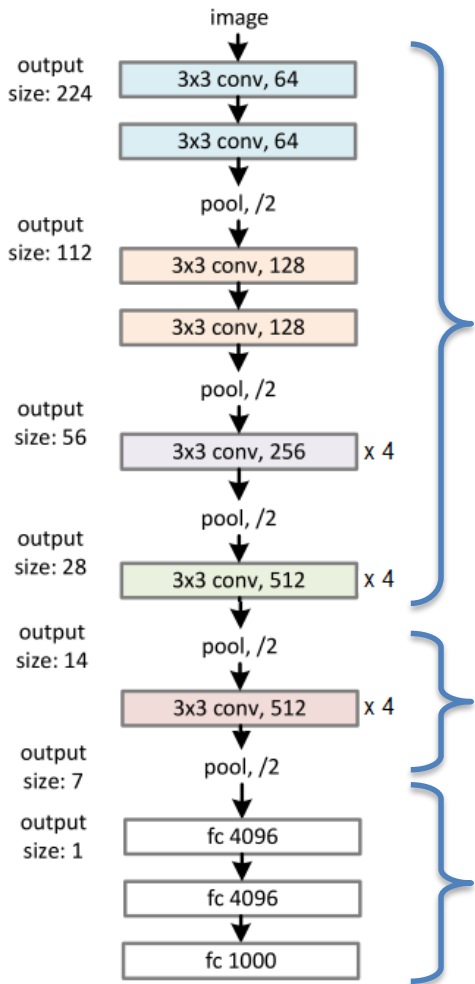
- Experimental setup: 1 machine, 8 cards.



Figure 13: VGG throughput speedup.

Batch Size: 64    VGG-19

Tofu's tiling for VGG-19 on 8 GPUs

**Data Parallelism**

**Hybrid Parallelism**
- 8 GPUs into 4 groups
- Data parallelism among groups
- Model parallelism within each group (tile on channel)

**Model Parallelism**
- Tile on both row and column for weight matrices

# Recap

- Data parallelism suffers from *batch-size-dilemma.*
- Other parallelisms exist but are hard to program.
  - Model parallelism, hybrid parallelism, combined parallelism, etc.
- Tofu automatically parallelizes deep learning training
  - Figure out distributed strategies for each operator.
  - Combine strategies recursively.
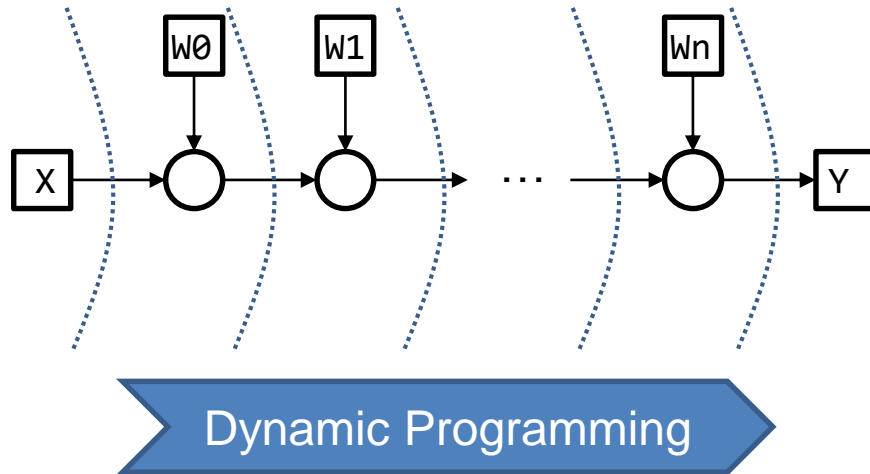  - Proved to have least communication cost.

# Q & A

YOU SHALL NOT PASS
zipmeme

# One-cut Tiling Algorithm

- Given a dataflow graph $G$, find $\mathcal{T}_{min}: M_G \mapsto \{\mathrm{R,C,r}\}$ such that the communication cost of *all* matrix multiplications are minimized.
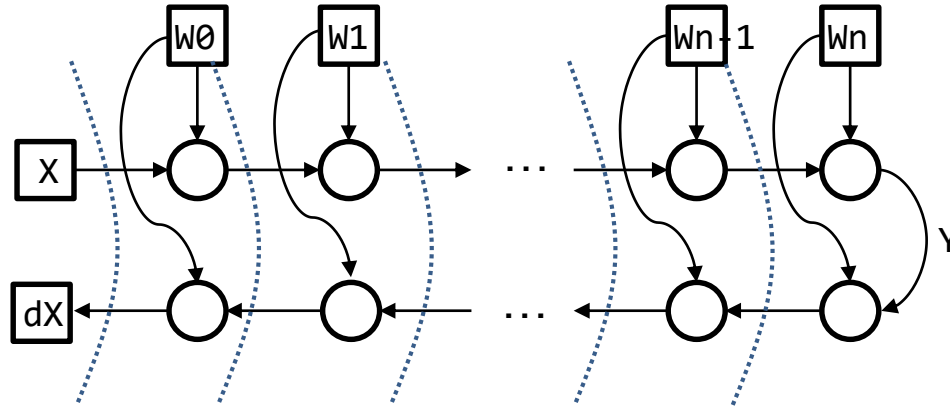
- Case #1:

$$XW_0 W_1 \dots W_n = Y$$



Dynamic Programming

# One-cut Tiling Algorithm

- Case #2:

$$XW_0 W_1 \dots W_n = Y$$
$$dX = YW_n^T W_{n-1}^T \dots W_0^T$$



Dynamic Programming

# One-cut Tiling Algorithm

- Organize nodes in the dataflow graph into levels, such that for any node, **all** its neighbors are contained in the adjacent levels.
- BFS is one way to produce such levels.
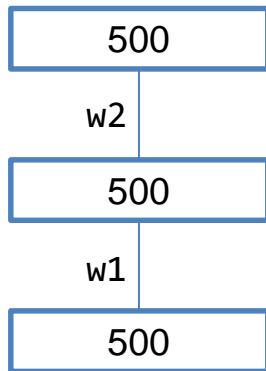- Dynamic Programming:

*Initial condition:*

$$g_0(\tau_0) = \text{level\_cost}_0(\phi, \tau_0)$$

*DP equation* $(l \geqslant 1)$:

$$g_l(\tau_l) = \min_{\tau_{l-1}} \left\{ \text{level\_cost}_l(\tau_{l-1}, \tau_l) + g_{l-1}(\tau_{l-1}) \right\}$$

# Which One is Better?

**ToyNet Configuration**

```
        500
         |
        w2
         |
        500
         |
        w1
         |
        500
```

nGPUs: 16
Batch size: 300

Parameter (gradients) size:
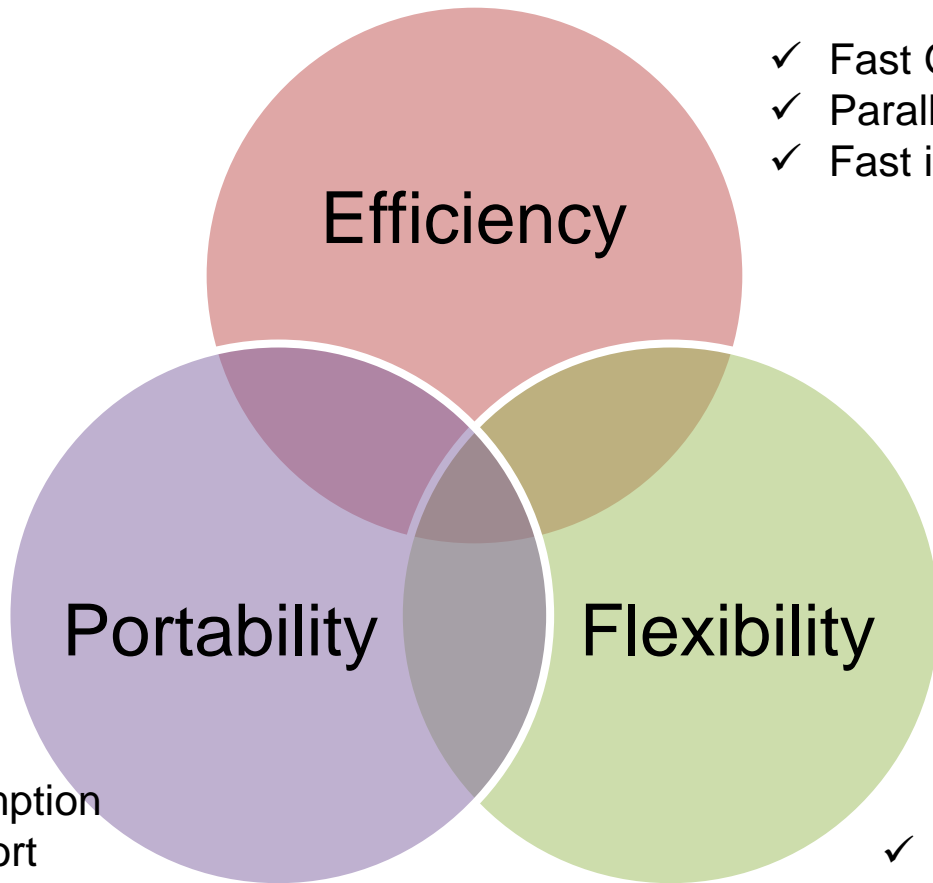  500 * 500 * 2 = 500K
Activation (gradients) size:
  500 * 300 * 2 = 300K

✓ Data Parallelism
  • 500K * 2 * 4B * 16 = 64MB

✓ Model Parallelism
  • 300K * 2 * 4B * 16 = 38.4MB

✓ Hybrid Parallelism
  • 4 groups of GPUs, each group has 4 GPUs
  • Model Parallelism among groups
    • 300K * 2 * 4B * 4 = 9.6MB
  • Data Parallelism within each group
    • 500K / 4 * 2 * 4B * 4 = 4MB
  • 9.6MB + 4 * 4MB = 25.6MB
  • Save 33.3% communications!

# Single Card Different Tilings

- Per batch running time for a 4-layers MLP network.
- Hidden layer size: 8192
- Partition dataflow to 8 workers but put them on the same GPU.

| Batch Size | Single GPU | Single GPU w/ Tofu partitions |
|---|---|---|
| 512 | 0.31s | 0.19s |
| 1024 | 0.56s | 0.39s |
| 2048 | 1.13s | 0.73s |

# Construct Parallel Execution Graph

- Three-phase computation



*Inputs Conversion Phase*   *Computation Phase*   *Outputs Conversion Phase*

Semantic dataflow

Execution dataflow

# Construct Parallel Execution Graph

- Dataflow graph for tiling conversion.



R          Split          Shuffle          Concat          C