

# Tensor Relational Algebra for Distributed Machine Learning System Design

Binhang Yuan, Dimitrije Jankov, Jia Zou,  
Yuxin Tang, Daniel Bourgeois, Chris Jermaine



# Current ML System



Systems like TensorFlow and PyTorch include:

- A front-end:
  - An API for specifying the model;
  - An automatic differentiation engine.
- A back-end:
  - A compute engine responsibly for (parallel/distributed) execution of SGD optimization;
  - Common parallel strategies include: *data parallelism*; *pipeline parallelism* and *operation partitioning*.

# What is the Problem?

- Current ML runtime is optimized for particular device(s) and computation paradigms, without careful consideration for distribution.
- Database systems are very successful for declarative parallel/distributed processing.
- However, database is built on set theory which is not appropriate for ML computation.
- Our goal is to bridge this gap!

# Run it in a Relational Style

- Encode each computation to a group of relational operations, e.g., join, selection, and aggregation;
- Encode the whole computation flow graph into a query plan;
- Let the system automate the optimization in a distributed runtime.

# Matrix Multiplication Example

Let us take matrix multiply in FC as an example:

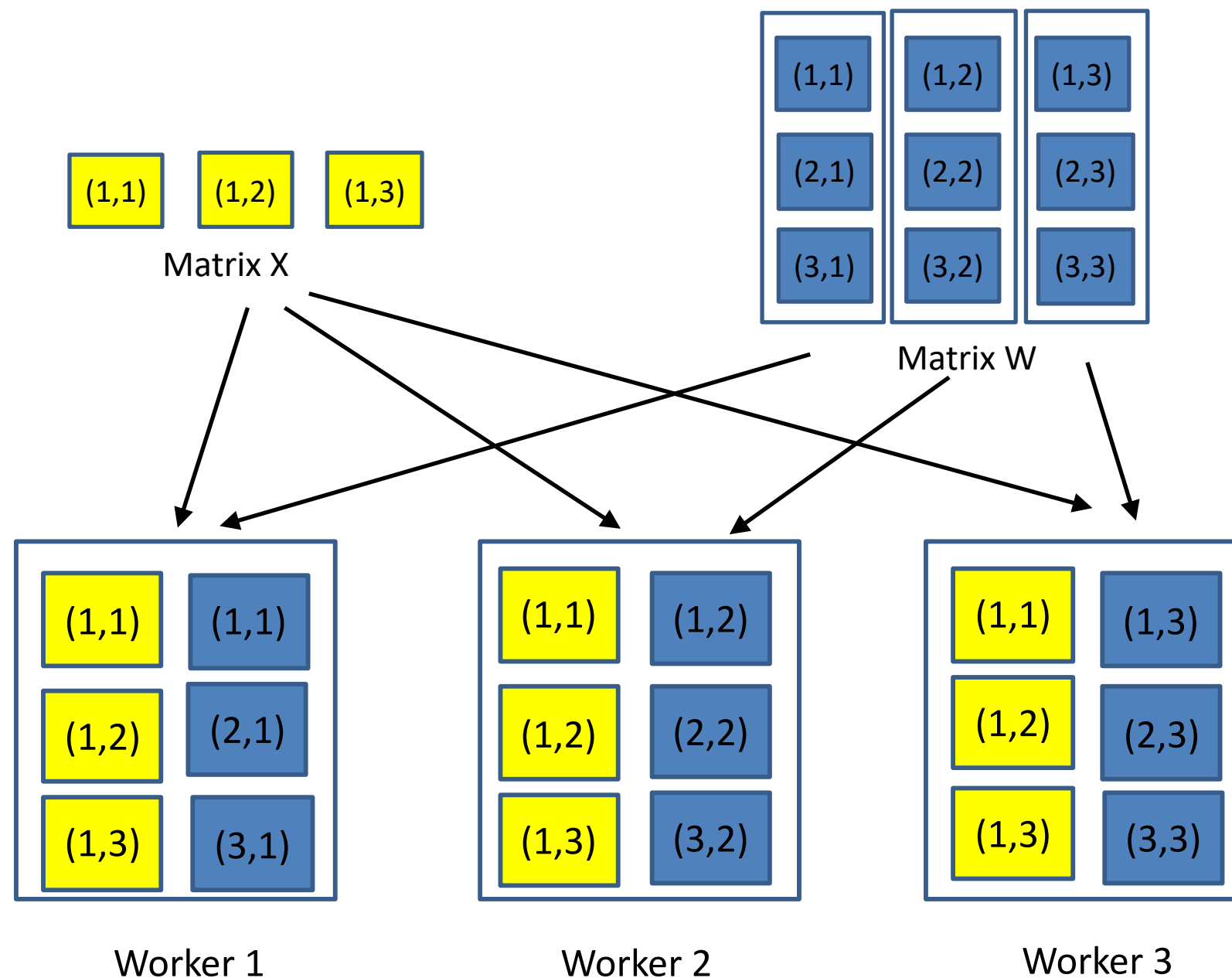
Assume two matrices are stored as blocks in two relations;



# How the System Executes the plan?

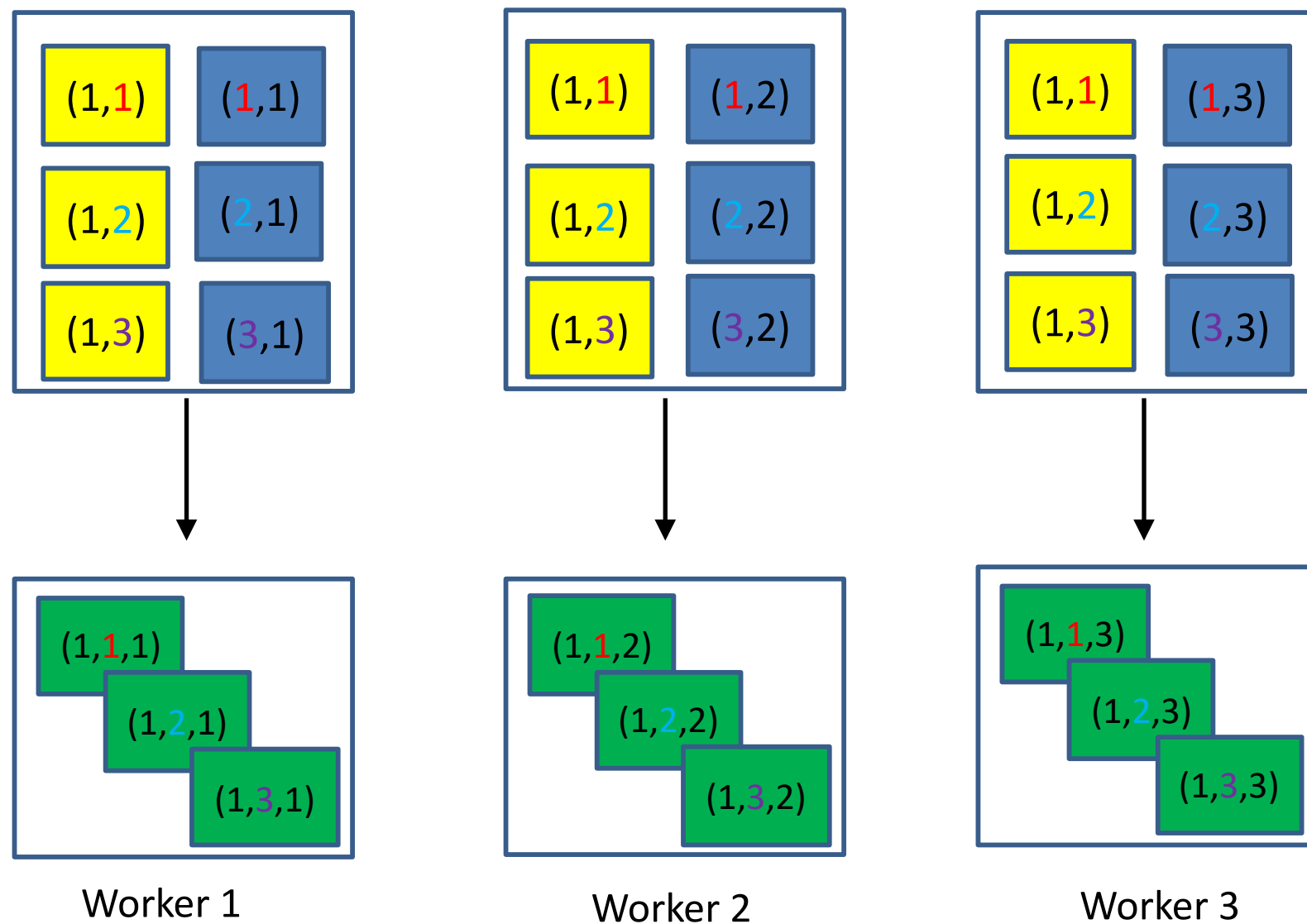
Given 3 worker nodes;

The system can determine to broadcast **X** and partition **W**:



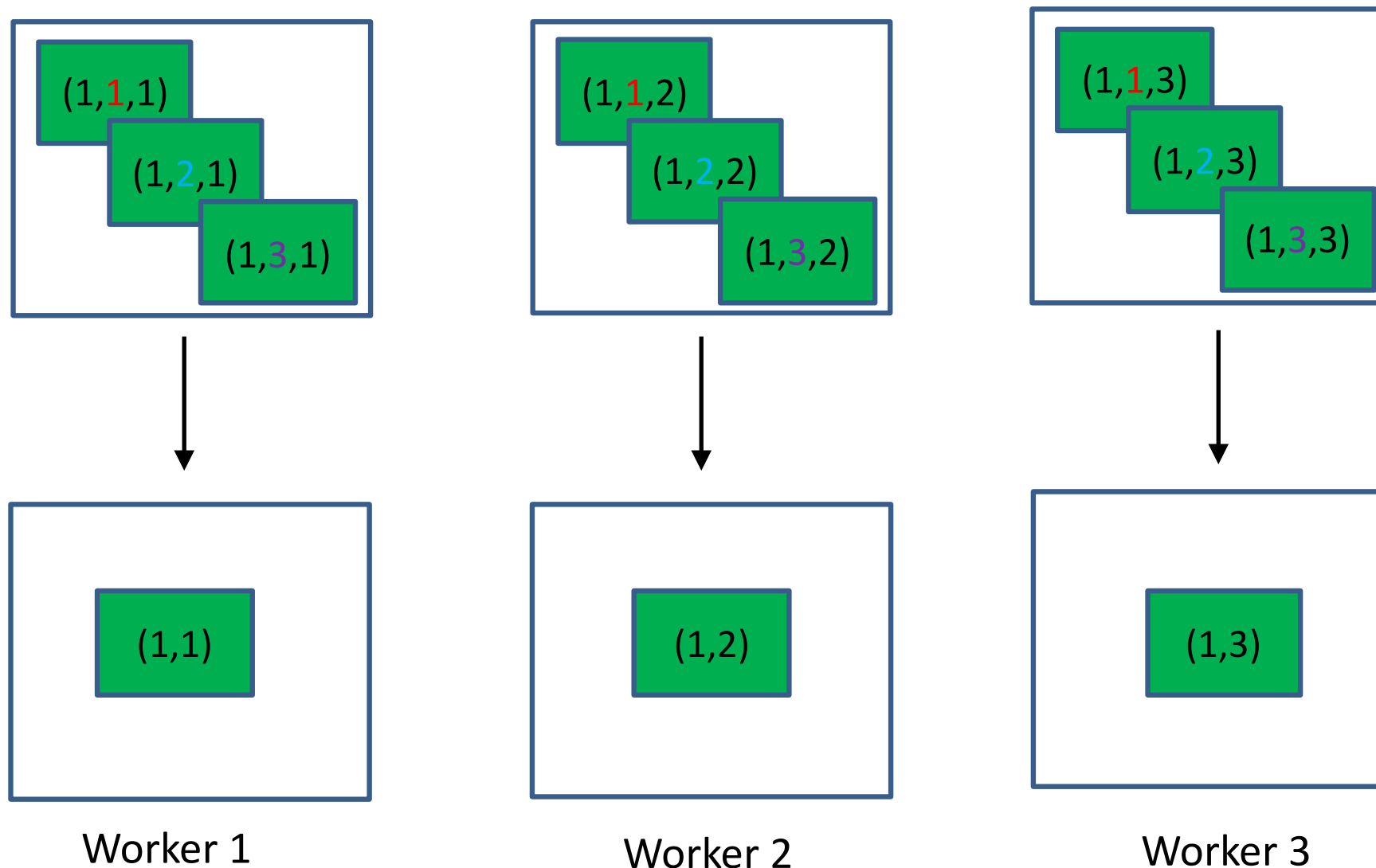
# How the System Executes the plan?

The system will do local join and apply multiplication on blocks:



# How the System Executes the plan?

The system will do aggregation: local aggregation and apply addition on blocks:





# What Should Be the Right Abstractions?

- Answer three fundamental questions:
  - *What should be the right formalization for express distributed ML computations?*
    - ➔ **Tensor Relational Algebra (TRA).**
    - What is the right abstraction to implement such formalization, especially for a distributed runtime?
      - ➔ **Implementation Abstraction (IA).**
      - What are the unique optimizations we can adopt for this implementation abstraction?
        - ➔ **A rule based relational optimizer.**

# Tensor Relational Algebra (TRA)

- A **tensor relation** includes tuples of  $\langle \text{key}, \text{array} \rangle$ .
- Support UDF-style operations:  
*Aggregate  $\Sigma$ ; Join  $\bowtie$ ; ReKey; Filter  $\sigma$ ; Transform  $\lambda$ ;  
Tile; Concat.*
- This abstraction is expressive enough to represent any tensor manipulation allowed by Einstein Notation.

# TRA — Operations

- **Join:**  $\bowtie: \left( (\mathbb{Z}^*)^g \times (\mathbb{Z}^*)^g \times \left( T^{(r_l, b_l)} \times T^{(r_r, b_r)} \rightarrow T^{(r_o, b_o)} \right) \right) \rightarrow \left( R^{(k, r, b)} \rightarrow R^{(g, r, b)} \right)$
- **Aggregation:**  $\Sigma: \left( (\mathbb{Z}^*)^g \times \left( T^{(r, b)} \times T^{(r, b)} \rightarrow T^{(r, b)} \right) \right) \rightarrow \left( R^{(k, r, b)} \rightarrow R^{(g, r, b)} \right)$
- **Rekey:**  $\text{ReKey}: \left( (\mathbb{Z}^*)^{k_i} \rightarrow (\mathbb{Z}^*)^{k_o} \right) \rightarrow \left( R^{(k_i, r, b)} \rightarrow R^{(k_o, r, b)} \right)$
- **Transform:**  $\lambda: \left( T^{(r_i, b_i)} \rightarrow T^{(r_o, b_o)} \right) \rightarrow \left( R^{(k, r_i, b_i)} \rightarrow R^{(k, r_o, b_o)} \right)$
- **Filter:**  $\sigma: \left( (\mathbb{Z}^*)^k \rightarrow \{\text{true}, \text{false}\} \right) \rightarrow \left( R^{(k, r, b)} \rightarrow R^{(k, r, b)} \right)$
- **Tile:**  $\text{Tile}: (\mathbb{Z}^* \times \mathbb{Z}^*) \rightarrow \left( R^{(k, r, b)} \rightarrow R^{(k+1, r, b')} \right)$
- **Concat:**  $\text{Concat}: (\mathbb{Z}^* \times \mathbb{Z}^*) \rightarrow \left( R^{(k, r, b)} \rightarrow R^{(k-1, r, b')} \right)$

# Implementation Abstraction (IA)

- TRA can formalize ML operations, but is not sufficient to specify different distributed implementations.
- IA is to bridge this gap.
- A physical tensor relation includes tuples of  $\langle \text{key}, \text{array}, \text{site} \rangle$ :
- Extend the TRA to operations preferred in Distributed runtimes.

# IA Operations

- Assign tuple to site:
  - **Broadcast**:  $\text{BROADCAST} : \mathcal{R}^{(k,r,b,s)} \rightarrow \mathcal{R}^{(k,r,b,s)}$ , after this we will have  $\text{ALL}(R) = \text{true}$
  - **Shuffle**:  $\text{SHUFFLE} : 2^{\{1\dots k\}} \rightarrow \left( \mathcal{R}^{(k,r,b,s)} \rightarrow \mathcal{R}^{(k,r,b,s)} \right)$ , after this we will have  $\text{PART}_D(R) = \text{true}$
- local operations:
  - **Local Join**:  $\bowtie^L : \left( (\mathbb{Z}^*)^g \times (\mathbb{Z}^*)^g \times \left( T^{(r_l,b_l)} \times T^{(r_r,b_r)} \rightarrow T^{(r_o,b_o)} \right) \right) \rightarrow \left( \mathcal{R}^{(k_l,r_l,b_l,s)} \times \mathcal{R}^{(k_r,r_r,b_r,s)} \rightarrow \mathcal{R}^{(k_l+k_r-g,r_o,b_o,s)} \right)$
  - **Local Aggregate**:  $\Sigma^L : \left( (\mathbb{Z}^*)^k \times \left( T^{(r,b)} \times T^{(r,b)} \rightarrow T^{(r,b)} \right) \right) \rightarrow \left( \mathcal{R}^{(k_l,r_l,b_l,s)} \times \mathcal{R}^{(k_r,r_r,b_r,s)} \rightarrow \mathcal{R}^{(k_l+k_r-g,r_o,b_o,s)} \right)$
  - **Local filter**:  $\sigma^L : \left( (\mathbb{Z}^*)^g \rightarrow \{\text{true}, \text{false}\} \right) \rightarrow \left( \mathcal{R}^{(k,r,b,s)} \rightarrow \mathcal{R}^{(k,r,b,s)} \right)$
  - **Local map**:  $\lambda^L : \left( \left( (\mathbb{Z}^*)^{k_i} \rightarrow ((\mathbb{Z}^*)^{k_o})^m \right) \times \left( T^{(r_i,b_i)} \rightarrow (T^{(r_o,b_o)})^m \right) \right) \rightarrow \left( \mathcal{R}^{(k_i,r_i,b_i,s)} \rightarrow \mathcal{R}^{(k_o,r_o,b_o,s)} \right)$

# Optimization — Equivalence

- Equivalence of IA expressions: Equivalent input physical tensor relations leads to equivalent output physical tensor relations.
- Simple equivalence rules:
  - Kernel function composition;
  - Equivalent partitioning.
- Domain specific equivalence rules:
  - E.g., different IA representations for distributed matrix multiplication.

# Optimization — Equivalence

- Domain specific equivalence rules for matrix multiplication.
- Target TRA:  $\Sigma_{(\langle 0,2 \rangle, \text{matAdd})} \left( \bowtie_{(\langle 1 \rangle, \langle 0 \rangle, \text{matMul})} (R_X, R_Y) \right)$
- IA can provide implementations for:
  - Broadcast based matrix multiplication (BMM).
  - Cross product based matrix multiplication (CPMM).
  - Replication based matrix multiplication (RMM).

# Optimization — Cost Model

- A cost model to consider network traffic.
- In contrast to the classic query optimization (where estimating selectivity is usually difficult) — the continuity constraints make this free.
- Give input tensor relation(s), the output tensor relation's frontier (for the key) and bound (for the array) can be estimated.
- For a tensor relation  $R$ , suppose  $f$  is the estimated floating numbers in  $R$ ,  $s$  is the number of sites, then:
  - Broadcast has a cost of  $f \times s$ ;
  - Shuffle has a cost of  $s$ .



# Experiments

- **Aims:**

- *Can the proposed optimization generate efficient plan(s) for a particular task over a particular dataset?*
- *Can the end-to-end performance be competitive to the STOA HPC or ML systems?*

- **Benchmarks:**

- Distributed matrix multiplication (MM).
- Nearest Neighbor Search (NNS) in Riemannian metric space.
- **Feed Forward Neural Network SGD iteration (FFNN).**

# Experimental Task — FFNN

- 2-layer FFNN for classification:
  - MiniBatch ( $\mathbf{X} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{Y} \in \mathbb{R}^{N \times L}$ ),
  - Weight matrix of layer 1  $\mathbf{W}_1 \in \mathbb{R}^{D \times H}$ ,
  - Weight matrix of layer 2  $\mathbf{W}_2 \in \mathbb{R}^{H \times L}$ ;
- Settings:
  - Google speech recognition:  
 $N = 10^4$ ,  $D = 1600$ ,  $L = 10$ ,  $H = 10^5, 1.5 \times 10^5, 2 \times 10^5$ ;
  - Amazon 14k XML:  
 $N = 10^3$ ,  $D = 597540$ ,  $L = 14588$ ,  $H = 10^3, 3 \times 10^3, 5 \times 10^3$ ;
- IA implementation:
  - TRA-DP, TRA-MP;
- Compared with:
  - TensorFlow(TF)/PyTorch data parallel implementation;
  - A careful HPC implementation based on ScaLapack;
  - Dask — a popular distributed python analytic package.

# Experimental Results — FFNN1

## 2-layer FFNN for Google Speech

Cluster	CPU			GPU		
Nodes	2	5	10	2	5	10
100k Neurons						
PyTorch-DP	11.16	6.15	4.75	0.99	1.19	1.27
TF-DP	11.93	7.32	5.51	0.87	1.13	1.17
ScaLAPACK	8.32	4.97	2.79	NA	NA	NA
Dask	62.57	56.57	49.63	NA	NA	NA
TRA-DP	11.62	6.51	5.20	1.49	1.59	1.63
TRA-MP	26.56	28.71	29.09	7.01	11.56	Fail
200k Neurons						
PyTorch-DP	17.25	11.94	9.30	Fail	2.09	2.42
TF-DP	21.36	13.21	11.21	1.52	2.12	2.46
ScaLAPACK	17.18	10.05	5.06	NA	NA	NA
Dask	136.66	112.72	104.01	NA	NA	NA
TRA-DP	17.89	12.51	9.69	1.49	1.59	1.63
TRA-MP	37.82	54.23	59.84	Fail	Fail	Fail

## Predicted Cost

	TRA-DP	TRA-MP
100k Neurons	<b>9.7x10<sup>8</sup></b>	1.0x10 <sup>10</sup>
200k Neurons	<b>1.9x10<sup>9</sup></b>	2.0x10 <sup>10</sup>

# Experimental Results — FFNN2

## 2-layer FFNN for Amazon-14k XML (seconds)

Cluster	CPU			GPU		
Nodes	2	5	10	2	5	10
1k Neurons						
PyTorch-DP	9.74	10.29	10.34	2.67	3.76	4.20
TF-DP	Fail	Fail	Fail	Fail	Fail	Fail
ScaLAPCK	8.16	6.65	2.47	NA	NA	NA
Dask	45.40	42.15	29.34	NA	NA	NA
TRA-DP	12.50	14.29	15.68	4.67	4.69	4.73
TRA-MP	<b>3.86</b>	<b>2.79</b>	<b>1.7</b>	<b>0.4</b>	<b>0.37</b>	<b>0.35</b>
5k Neurons						
PyTorch-DP	34.05	46.53	50.17	Fail	Fail	Fail
TF-DP	Fail	Fail	Fail	Fail	Fail	Fail
ScaLAPACK	23.21	11.65	8.33	NA	NA	NA
Dask	246.56	143.86	127.26	NA	NA	NA
TRA-DP	44.12	68.54	75.15	Fail	Fail	Fail
TRA-MP	<b>18.59</b>	<b>8.07</b>	<b>5.57</b>	Fail	<b>0.59</b>	<b>0.48</b>

## Predicted Cost

	TRA-DP	TRA-MP
1k Neurons	$3.7 \times 10^9$	<b><math>1.0 \times 10^7</math></b>
5k Neurons	$1.8 \times 10^{10}$	<b><math>5.0 \times 10^7</math></b>

**Thank you!**