

Role of Tensors in Machine Learning

Anima Anandkumar

Amazon AI & Caltech

Tensors: Beyond 2D world

Scalar



Vector



Matrix

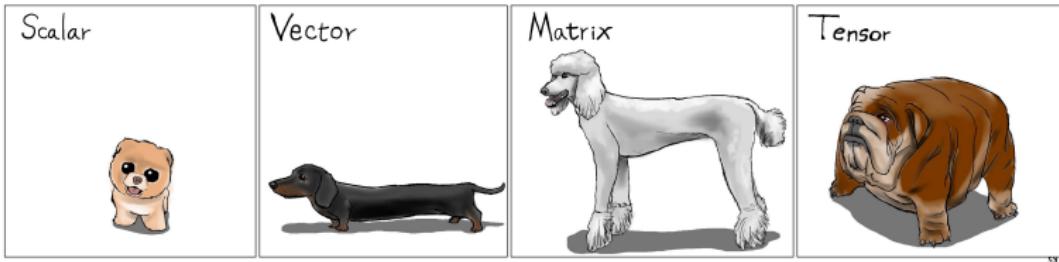


Tensor

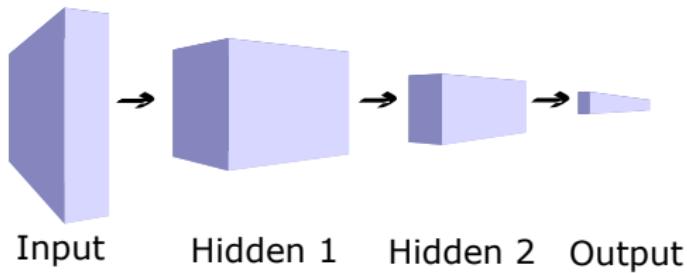


Modern data is inherently multi-dimensional

Tensors: Beyond 2D world



Modern data is inherently multi-dimensional

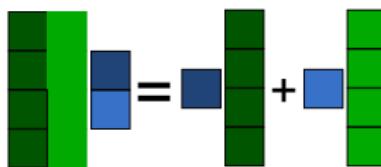


Tensor Contraction

Extends the notion of matrix product

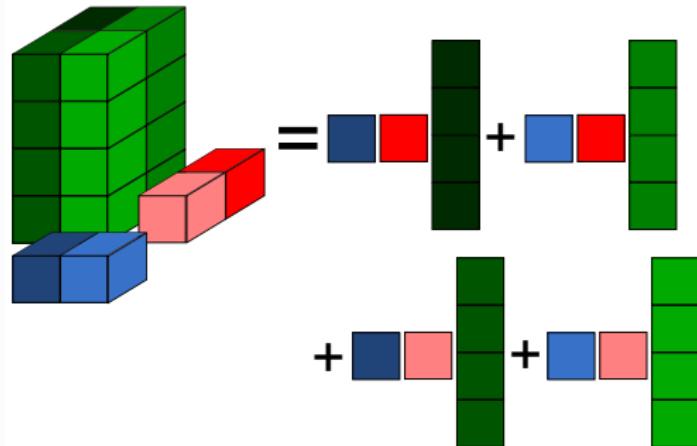
Matrix product

$$Mv = \sum_j v_j M_j$$



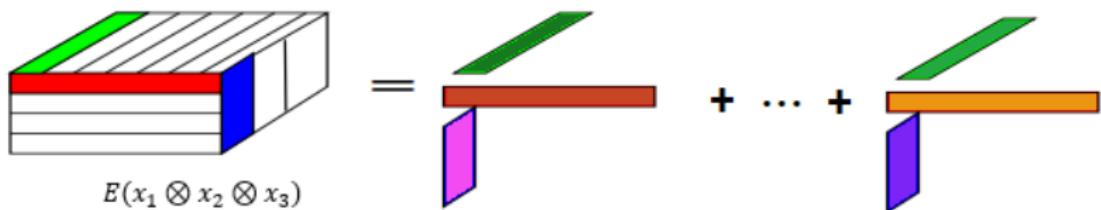
Tensor Contraction

$$T(u, v, \cdot) = \sum_{i,j} u_i v_j T_{i,j,:}$$



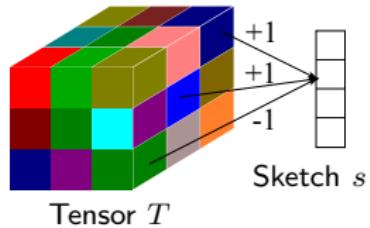
Tensor Decompositions

$$E(x_1 \otimes x_2) = \text{[red bar]} + \dots + \text{[orange bar]}$$


$$E(x_1 \otimes x_2 \otimes x_3) = \text{[red block]} + \dots + \text{[orange block]}$$


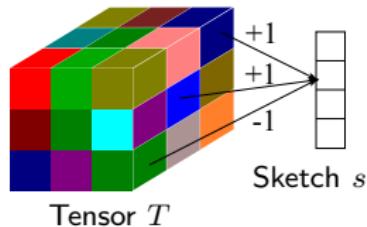
Tensor Sketches

- Dimensionality reduction through sketching.
 - Complexity independent of tensor order:
exponential gain!



Tensor Sketches

- Randomized dimensionality reduction through sketching.
 - ▶ Complexity independent of tensor order: **exponential gain!**

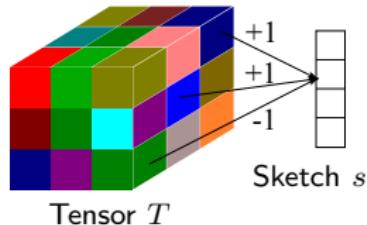


Applications

- Tensor Decomposition via Sketching by Wang, Tung, Smola, **A**, NIPS'15.
- Compact Tensor Pooling for Visual Question and Answering by Shi, Anubhai, Furlanello, **A**, CVPR 2017 VQA workshop.

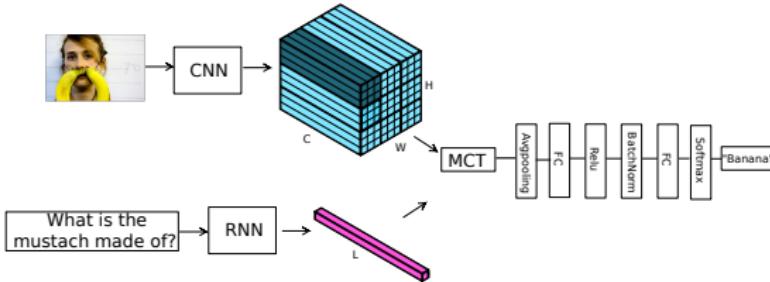
Tensor Sketches

- Randomized dimensionality reduction through **sketching**.
 - ▶ Complexity independent of tensor order: **exponential gain!**



Applications

- Tensor Decomposition via Sketching by Wang, Tung, Smola, **A**, NIPS'15.
- Compact Tensor Pooling for Visual Question and Answering by Shi, Anubhai, Furlanello, **A**, CVPR 2017 VQA workshop.



Outline

1 Introduction

2 Tensor Contractions

3 Speeding up Tensor Contractions

4 Tensor Sketches

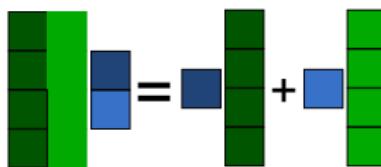
5 Conclusion

Tensor Contraction

Extends the notion of matrix product

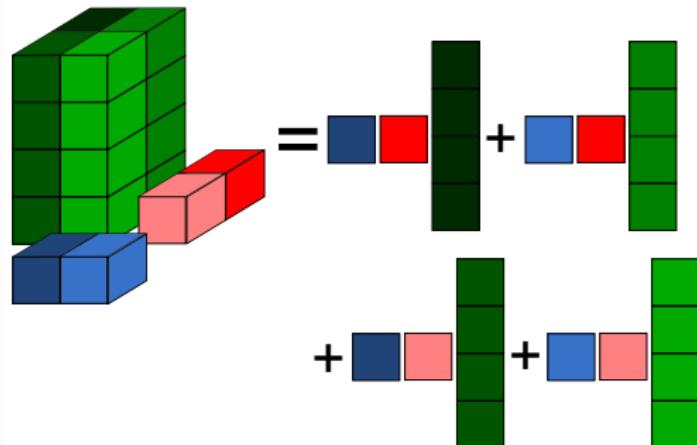
Matrix product

$$Mv = \sum_j v_j M_j$$

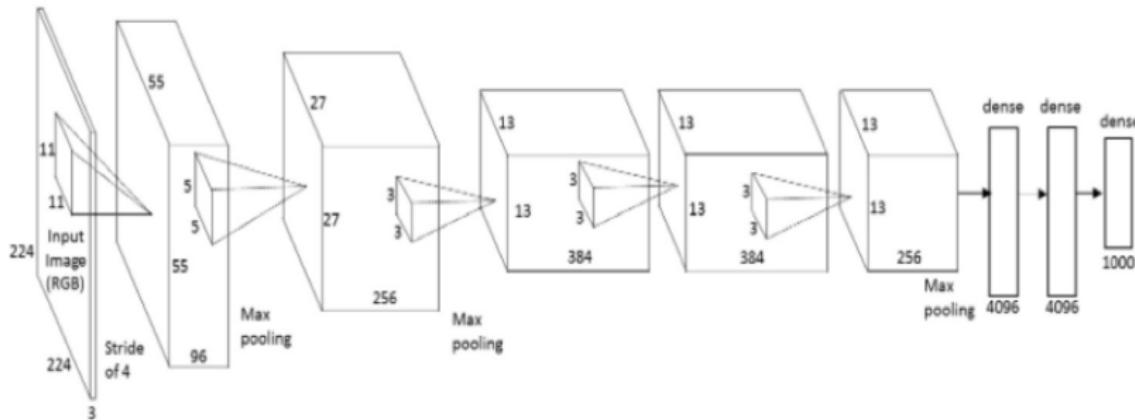


Tensor Contraction

$$T(u, v, \cdot) = \sum_{i,j} u_i v_j T_{i,j,:}$$

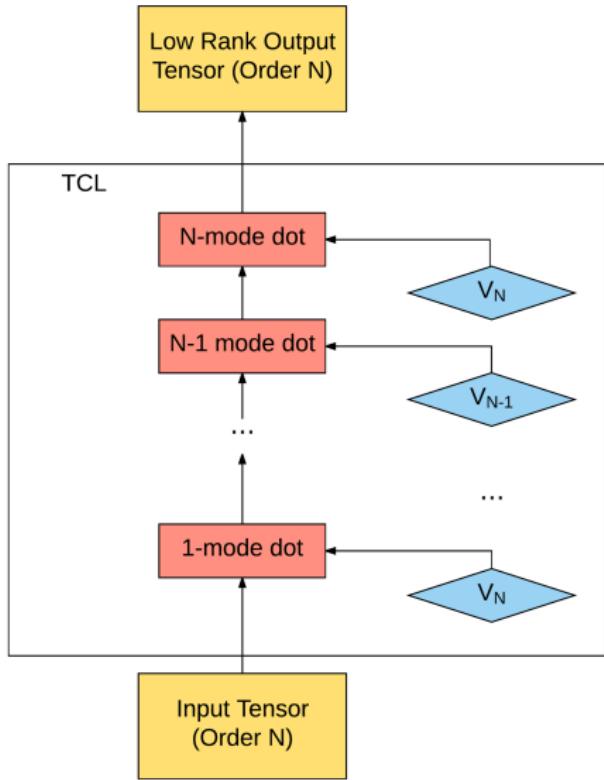
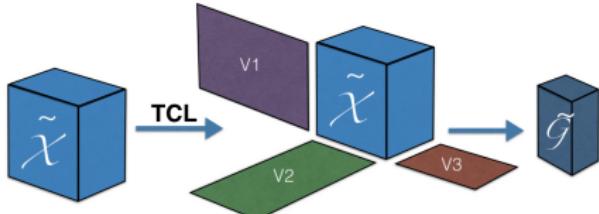


Employing Tensor Contractions in Alexnet



Replace fully connected layer with tensor contraction layer

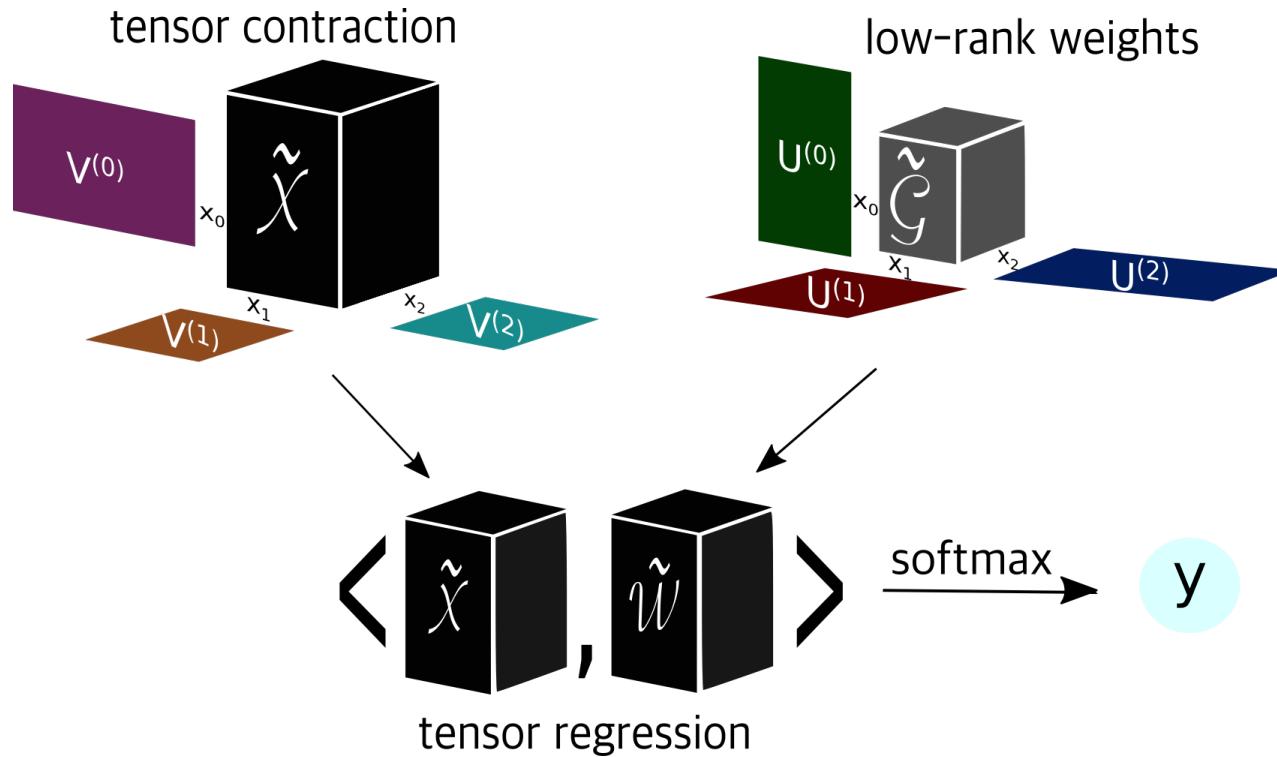
Enabling Tensor Contraction Layer in Mxnet



Performance of the TCL

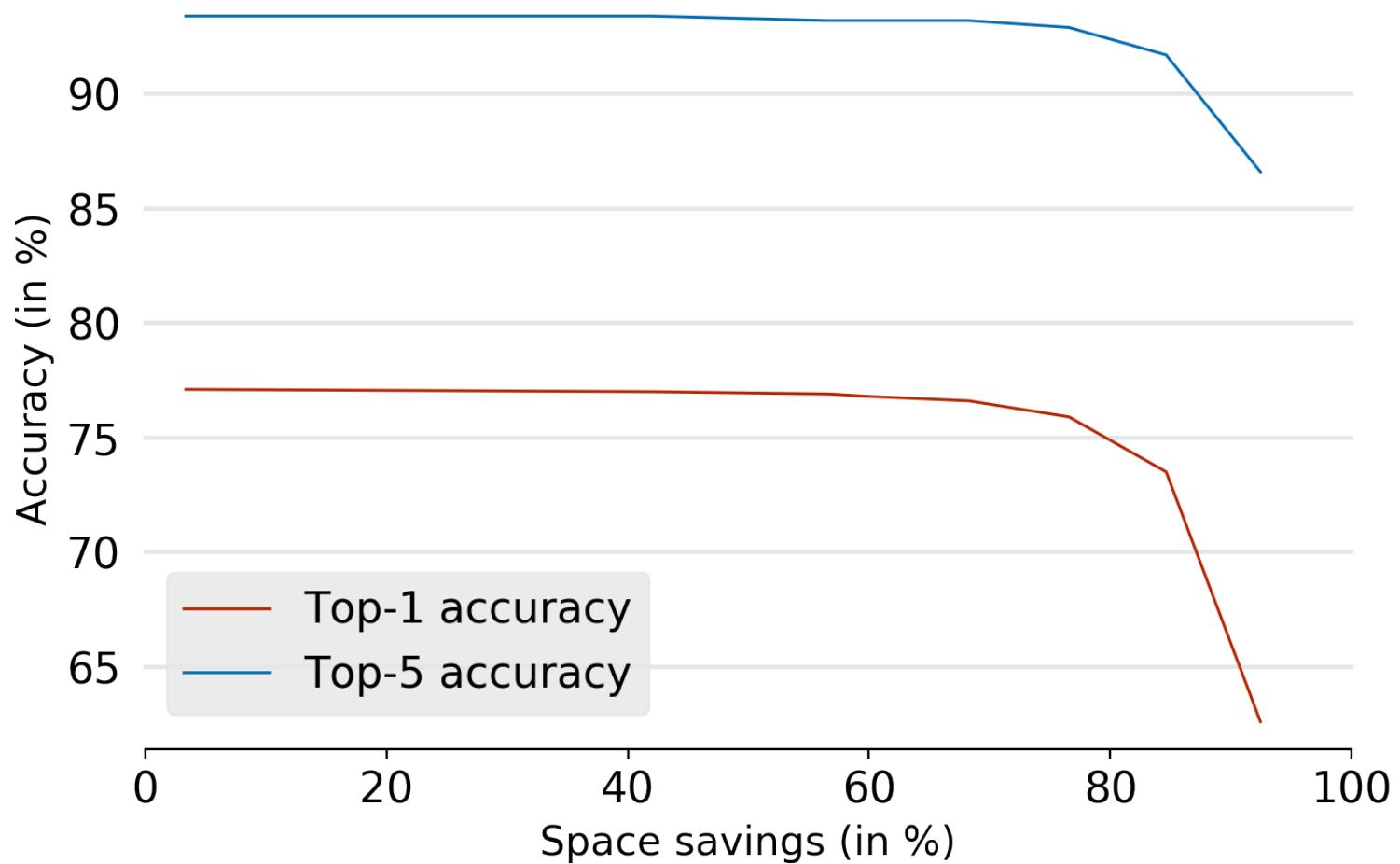
- Trained end-to-end
- On ImageNet with VGG:
 - 65.9% space savings
 - performance drop of 0.6% only
- On ImageNet with AlexNet:
 - 56.6% space savings
 - Performance improvement of 0.5%

Low-rank tensor regression



Tensor Regression Networks, J. Kossaifi, Z.C.Lipton, A.Khanna,
T.Furlanello and A.Anandkumar, ArXiv pre-publication

Performance and rank



Outline

1 Introduction

2 Tensor Contractions

3 Speeding up Tensor Contractions

4 Tensor Sketches

5 Conclusion

Speeding up Tensor Contractions

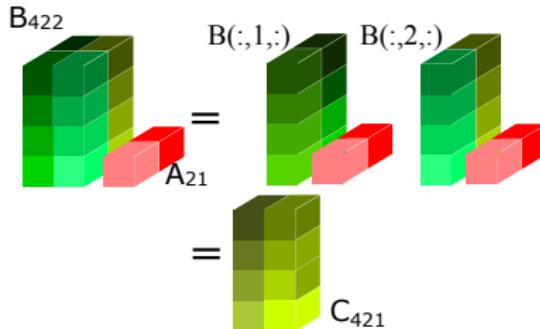
- ① Tensor contractions are a core primitive of multilinear algebra.
- ② BLAS 3: Unbounded compute intensity (no. of ops per I/O)

Consider single-index contractions: $C_{\mathcal{C}} = A_{\mathcal{A}} B_{\mathcal{B}}$

Speeding up Tensor Contractions

- ① Tensor contractions are a core primitive of multilinear algebra.
- ② BLAS 3: Unbounded compute intensity (no. of ops per I/O)

Consider single-index contractions: $C_{\mathcal{C}} = A_{\mathcal{A}} B_{\mathcal{B}}$



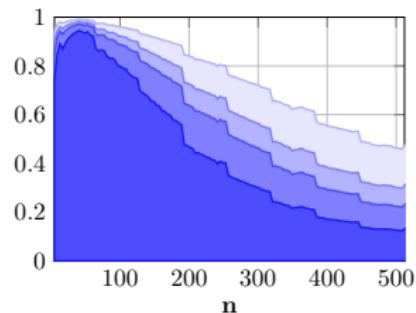
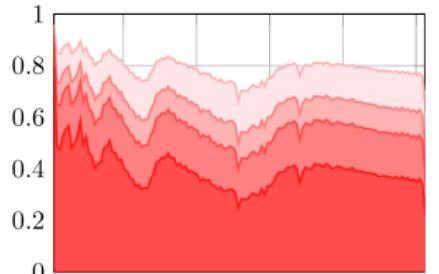
e.g. $C_{mnp} = B_{mnk} A_{kp}$

Speeding up Tensor Contraction

Explicit permutation dominates,
especially for **small tensors**.

Consider $C_{mnp} = A_{km} B_{pkn}$.

- ① $A_{km} \rightarrow A_{mk}$
- ② $B_{pkn} \rightarrow B_{kpn}$
- ③ $C_{mnp} \rightarrow C_{mpn}$
- ④ $C_{m(pn)} = A_{mk} B_{k(pn)}$
- ⑤ $C_{mpn} \rightarrow C_{mnp}$



(Top) CPU. (Bottom) GPU. The fraction of time spent in copies/transpositions. Lines are shown with 1, 2, 3, and 6 transpositions.

Existing Primitives

GEMM

- Suboptimal for many small matrices.

Pointer-to-Pointer BatchedGEMM

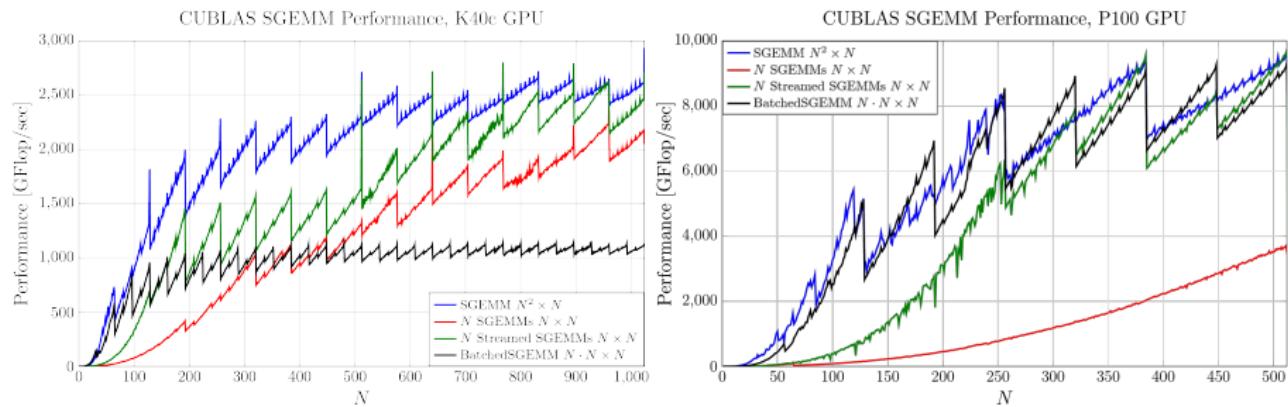
- Available in MKL 11.3 β and cuBLAS 4.1

$$C[p] = \alpha \text{op}(A[p]) \text{op}(B[p]) + \beta C[p]$$

```
cublas<T>gemmBatched(cublasHandle_t handle,
                        cublasOperation_t transA, cublasOperation_t transB,
                        int M, int N, int K,
                        const T* alpha,
                        const T** A, int ldA,
                        const T** B, int ldB,
                        const T* beta,
                        T** C, int ldC,
                        int batchCount)
```

Existing Primitives

Pointer-to-Pointer BatchedGEMM



A new primitive: StridedBatchedGEMM

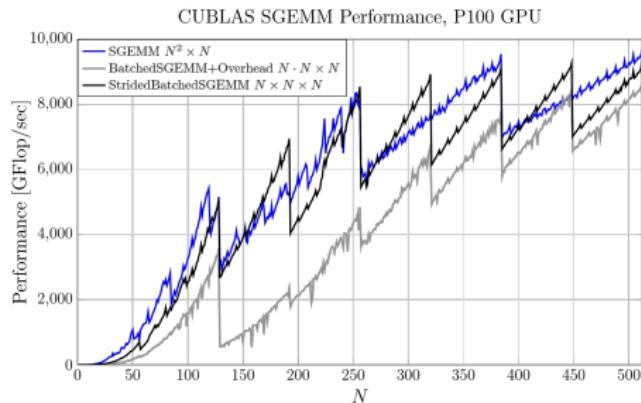
$$C[p] = \alpha \operatorname{op}(A[p]) \operatorname{op}(B[p]) + \beta C[p]$$

- Pointer-to-Pointer BatchedGEMM requires memory allocation and pre-computation.
- **Solution:** StridedBatchedGEMM with fixed strides.
 - ▶ Special case of Pointer-to-pointer BatchedGEMM.
 - ▶ No Pointer-to-pointer data structure or overhead.

```
cublas<T>gemmStridedBatched(cublasHandle_t handle,
                               cublasOperation_t transA, cublasOperation_t transB,
                               int M, int N, int K,
                               const T* alpha,
                               const T* A, int ldA1, int strideA,
                               const T* B, int ldB1, int strideB,
                               const T* beta,
                               T* C, int ldc1, int strideC,
                               int batchCount)
```

A new primitive: StridedBatchedGEMM

- Performance on par with pure GEMM (P100 and beyond).



StridedBatchedGEMM

Documentation in cuBLAS 8.0:

```
## grep StridedBatched -A 17 /usr/local/cuda/include/cublas_api.h
2320:CUBLASAPI cublasStatus_t cublasSgemmStridedBatched (cublasHandle_t handle,
2321-                                         cublasOperation_t transa,
2322-                                         cublasOperation_t transb,
2323-                                         int m,
2324-                                         int n,
2325-                                         int k,
2326-                                         const float *alpha, // host or device pointer
2327-                                         const float *A,
2328-                                         int lda,
2329-                                         long long int strideA, // purposely signed
2330-                                         const float *B,
2331-                                         int ldb,
2332-                                         long long int strideB,
2333-                                         const float *beta, // host or device pointer
2334-                                         float *C,
2335-                                         int ldc,
2336-                                         long long int strideC,
2337-                                         int batchCount);
...
...
```

Tensor Contraction with Extended BLAS Primitives

$$C_{mnp} = A_{**} \times B_{***}$$

$$C_{mnp} \equiv C[m + n \cdot \text{ldC1} + p \cdot \text{ldC2}]$$

Case	Contraction	Kernel1	Kernel2	Case	Contraction	Kernel1	Kernel2
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	4.1	$A_{kn}B_{kmp}$	$C_{mn[p]} = B_{km(p)}^\top A_{kn}$	
1.2	$A_{mk}B_{kpn}$	$C_{mn[p]} = A_{mk}B_{k[p]n}$	$C_{m[n]p} = A_{mk}B_{kp[n]}$	4.2	$A_{kn}B_{kpm}$	$C_{mn[p]} = B_{k[p]m}^\top A_{kn}$	
1.3	$A_{mk}B_{nkp}$	$C_{mn[p]} = A_{mk}B_{nk[p]}^\top$		4.3	$A_{kn}B_{mkp}$	$C_{mn[p]} = B_{mk[p]}A_{kn}$	
1.4	$A_{mk}B_{pkn}$	$C_{m[n]p} = A_{mk}B_{pk[n]}^\top$		4.4	$A_{kn}B_{pkm}$		
1.5	$A_{mk}B_{npk}$	$C_{m(np)} = A_{mk}B_{(np)k}^\top$	$C_{mn[p]} = A_{mk}B_{n[p]k}^\top$	4.5	$A_{kn}B_{mpk}$	$C_{mn[p]} = B_{m[p]k}A_{kn}$	
1.6	$A_{mk}B_{pnk}$	$C_{m[n]p} = A_{mk}B_{p[n]k}^\top$		4.6	$A_{kn}B_{pmk}$		
2.1	$A_{km}B_{knp}$	$C_{m(np)} = A_{km}^\top B_{k(np)}$	$C_{mn[p]} = A_{km}^\top B_{kn[p]}$	5.1	$A_{pk}B_{kmn}$	$C_{(mn)p} = B_{k(mn)}^\top A_{pk}^\top$	$C_{m[n]p} = B_{km[n]}^\top A_{pk}^\top$
2.2	$A_{km}B_{kpn}$	$C_{mn[p]} = A_{km}^\top B_{k[p]n}$	$C_{m[n]p} = A_{km}^\top B_{kp[n]}$	5.2	$A_{pk}B_{knm}$	$C_{m[n]p} = B_{k[n]m}^\top A_{pk}^\top$	
2.3	$A_{km}B_{nkp}$	$C_{mn[p]} = A_{km}^\top B_{nk[p]}^\top$		5.3	$A_{pk}B_{mkn}$	$C_{m[n]p} = B_{mk[n]}A_{pk}^\top$	
2.4	$A_{km}B_{pkn}$	$C_{m[n]p} = A_{km}^\top B_{pk[n]}^\top$		5.4	$A_{pk}B_{nkm}$		
2.5	$A_{km}B_{npk}$	$C_{m(np)} = A_{km}^\top B_{(np)k}^\top$	$C_{mn[p]} = A_{km}^\top B_{n[p]k}^\top$	5.5	$A_{pk}B_{mnk}$	$C_{(mn)p} = B_{(mn)k}A_{pk}^\top$	$C_{m[n]p} = B_{m[n]k}A_{pk}^\top$
2.6	$A_{km}B_{pnk}$	$C_{m[n]p} = A_{km}^\top B_{p[n]k}^\top$		5.6	$A_{pk}B_{nmk}$		
3.1	$A_{nk}B_{kmp}$	$C_{mn[p]} = B_{km[p]}^\top A_{nk}^\top$		6.1	$A_{kp}B_{kmn}$	$C_{(mn)p} = B_{k(mn)}^\top A_{kp}$	$C_{m[n]p} = B_{km[n]}^\top A_{kp}$
3.2	$A_{nk}B_{kpm}$	$C_{mn[p]} = B_{k[p]m}^\top A_{nk}^\top$		6.2	$A_{kp}B_{knm}$	$C_{m[n]p} = B_{k[n]m}^\top A_{kp}$	
3.3	$A_{nk}B_{mkp}$	$C_{mn[p]} = B_{mk[p]}A_{nk}^\top$		6.3	$A_{kp}B_{mkn}$	$C_{m[n]p} = B_{mk[n]}A_{kp}$	
3.4	$A_{nk}B_{pkm}$			6.4	$A_{kp}B_{nkm}$		
3.5	$A_{nk}B_{mpk}$	$C_{mn[p]} = B_{m[p]k}A_{nk}^\top$		6.5	$A_{kp}B_{mnk}$	$C_{(mn)p} = B_{(mn)k}A_{kp}$	$C_{m[n]p} = B_{m[n]k}A_{kp}$
3.6	$A_{nk}B_{pmk}$			6.6	$A_{kp}B_{nmk}$		

Tensor Contraction with Extended BLAS Primitives

Case	Contraction	Kernel1	Kernel2	Kernel3
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	$C_{m[n]p} = A_{mk}B_{k[n]p}$
6.1	$A_{kp}B_{kmn}$	$C_{(mn)p} = B_{k(mn)}^\top A_{kp}$	$C_{m[n]p} = B_{km[n]}^\top A_{kp}$	

Example: Mappings to Level 3 BLAS routines

- Case 1.1, Kernel2: $C_{mn[p]} = A_{mk}B_{kn[p]}$

```
cublasDgemmStridedBatched(handle,
                            CUBLAS_OP_N, CUBLAS_OP_N,
                            M, N, K,
                            &alpha,
                            A, ldA1, 0,
                            B, ldB1, ldB2,
                            &beta,
                            C, ldc1, ldc2,
                            P)
```

Tensor Contraction with Extended BLAS Primitives

Case	Contraction	Kernel1	Kernel2	Kernel3
1.1	$A_{mk}B_{knp}$	$C_{m(np)} = A_{mk}B_{k(np)}$	$C_{mn[p]} = A_{mk}B_{kn[p]}$	$C_{m[n]p} = A_{mk}B_{k[n]p}$
6.1	$A_{kp}B_{kmn}$	$C_{(mn)p} = B_{k(mn)}^\top A_{kp}$	$C_{m[n]p} = B_{km[n]}^\top A_{kp}$	

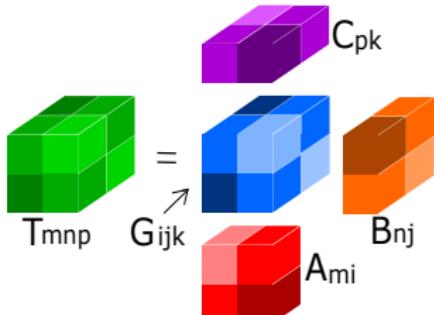
Example: Mappings to Level 3 BLAS routines

- Case 6.1, Kernel2: $C_{m[n]p} = B_{km[n]}^\top A_{kp}$

```
cublasDgemmStridedBatched(handle,
                            CUBLAS_OP_T, CUBLAS_OP_N,
                            M, P, K,
                            &alpha,
                            B, ldB1, ldB2,
                            A, ldA1, 0,
                            &beta,
                            C, ldC2, ldC1,
                            N)
```

Applications: Tucker Decomposition

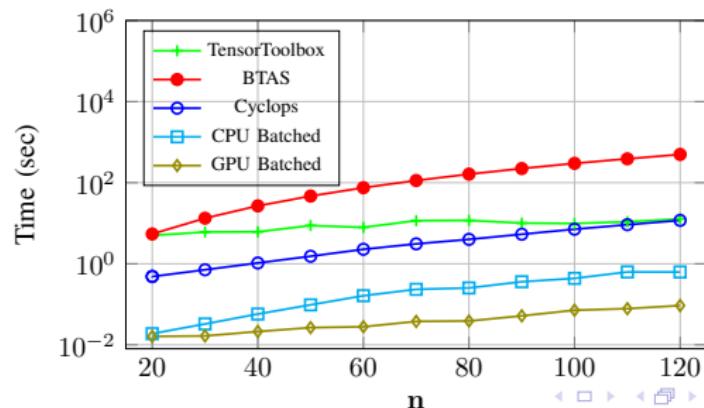
$$T_{mnp} = G_{ijk} A_{mi} B_{nj} C_{pk}$$



Main steps in the algorithm

- $Y_{mjk} = T_{mnp} B_{nj}^t C_{pk}^t$
- $Y_{ink} = T_{mnp} A_{mi}^{t+1} C_{pk}^t$
- $Y_{ijp} = T_{mnp} B_{nj}^{t+1} A_{mi}^{t+1}$

Performance on Tucker decomposition:



Applications: FFT

Low-Communication FFT for multiple GPUs involves tensor contractions.

$$\begin{array}{ll} T_{pib} = S2T_{ij\bar{s}}^{(p)} S_{pj(b+s)} & \implies T_{pib} = S2T_{i(j\bar{s})}^{(p)} S_{p(j\bar{s})b} \\ M_{pq\bar{b}} = S2M_{qi} S_{pib} & \implies M_{pq[b]} = S_{pi[b]} S2M_{qi}^T \\ M_{pq\bar{b}'} = M2M_{qm}^- M_{pm\bar{b}-} + M2M_{qm}^+ M_{pm\bar{b}+} & \implies M_{pq[b']} = M_{pM[b]} M2M_{qM}^T \\ r_p = 1_{ib} S_{pib} = 1_{qb} M_{pq\bar{b}} & \implies r_p = 1_{(qb)} M_{p(qb)} \\ L_{pn\bar{b}} = M2L_{nms}^{(p)} M_{pm(b+s)} & \implies L_{pn\bar{b}} = M2L_{n(ms)}^{(p)} M_{p(ms)\bar{b}} \\ L_{pq\bar{b}\pm} = L2L_{qm}^\pm L_{pm\bar{b}'} & \implies L_{pq[b]} = L_{pM[b']} M2M_{qM} \\ T_{pib} = L2T_{iq} L_{pq\bar{b}} & \implies T_{pi[b]} = L_{pq[b]} S2M_{qi} \end{array}$$

- StridedBatchedGEMM composes 75%+ of the runtime.
 - ▶ Essential to the performance.
 - ▶ Two custom kernels are precisely interleaved GEMMs.
- 2 P100 GPUs: **1.3x** over cuFFXT.
- 8 P100 GPUs: **2.1x** over cuFFXT.

Summary of this Section

- Tensor contractions provide significant space savings in deep learning.
- Fast tensor contractions using extended BLAS primitive:
StridedBatchedGEMM
- Avoids explicit transpositions or permutations.
- **10x**(GPU) and **2x**(CPU) speedup on small/moderate sized tensors.
- **Available in cuBLAS 8.0**

Summary of this Section

- Tensor contractions provide significant space savings in deep learning.
- Fast tensor contractions using extended BLAS primitive:
StridedBatchedGEMM
- Avoids explicit transpositions or permutations.
- **10x**(GPU) and **2x**(CPU) speedup on small/moderate sized tensors.
- **Available in cuBLAS 8.0**
- Future work:
 - ▶ Exceptional case kernels/performance/interface??
 - ▶ Library Optimizations
 - ★ Matrix stride zero – Persistent Matrix Strided Batched GEMM

Outline

1 Introduction

2 Tensor Contractions

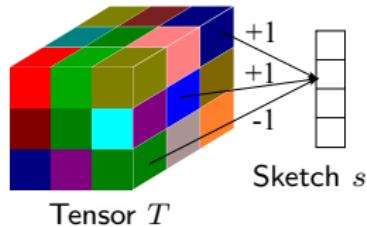
3 Speeding up Tensor Contractions

4 Tensor Sketches

5 Conclusion

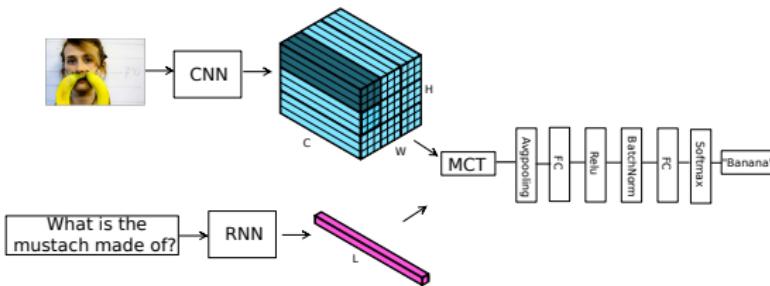
Tensor Sketches

- Randomized dimensionality reduction through sketching.
 - ▶ Complexity independent of tensor order: exponential gain!



Applications

- Tensor Decomposition via Sketching
- Visual Question and Answering

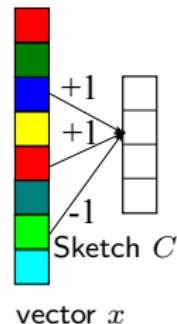


Tensor Sketching

- Dimensionality reduction through sketching.

Count Sketch for vector x

- $C[h[i]]+ = s[i]x[i]$, for $s[i] \in \{-1, +1\}^n$

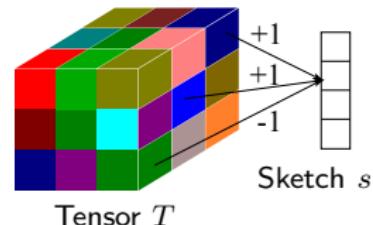


vector x

Count Sketch for outer products $x \otimes y$

- Convolution of count sketches

$$\begin{aligned}C(x \otimes y, h, s) &= C(x, h, s) * C(y, h, s) \\&= FFT^{-1}(FFT(C(x, h, s))FFT(C(y, h, s)))\end{aligned}$$



Tensor T

Accelerated tensor low-rank decomposition

Symmetric tensor CP decomposition

For symmetric tensor $T \in \mathbb{R}^{n \times n \times n}$, find $\{(\lambda_i, u_i)\}_{i=1}^k$ to minimize

$$\|T - \sum_{i=1}^k \lambda_i u_i^{\otimes 3}\|_F^2.$$

- Wide application in data mining and latent variable models.
- **Tensor power iteration:** $u^{(t+1)} = T(I, u^{(t)}, u^{(t)}) / \|T(I, u^{(t)}, u^{(t)})\|_2$.
- Accelerated tensor power iteration via sketching:
 - TENSORSKETCH: $s(T) \in \mathbb{R}^b$, for $n < b \ll n^3$.

$$\begin{aligned}[T(I, u, u)]_i &\approx \langle s(T), s(u \otimes u \otimes e_i) \rangle \\&= \langle \mathcal{F}(s(T)), \mathcal{F}(s(u)) \cdot * \mathcal{F}(s(u)) \cdot * \mathcal{F}(s(e_i)) \rangle \\&= \langle \mathcal{F}^{-1}(\mathcal{F}(s(T)) \cdot * \overline{\mathcal{F}(s(u))} \cdot * \overline{\mathcal{F}(s(u))}), s(e_i) \rangle.\end{aligned}$$

- Time complexity: $O(n^3) \rightarrow O(n + b \log b)$.

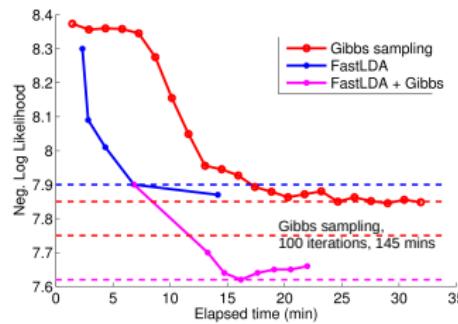
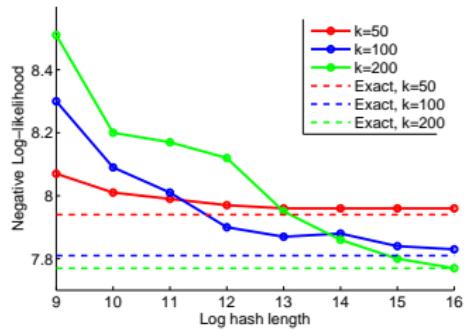
Efficient spectral method for topic modeling

Topic modeling

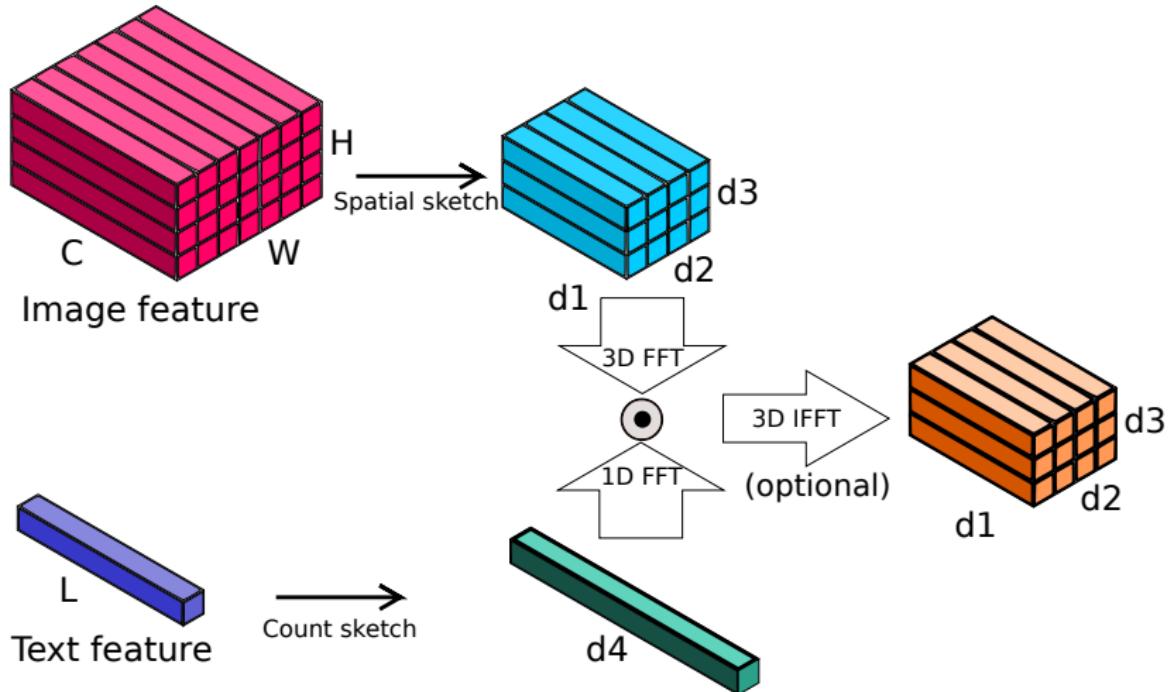
V : vocabulary size; k : number of topics. Recover topic distributions $\mu_1, \dots, \mu_k \in \mathbb{R}^V$ from N unlabeled documents.

Figure 1: Negative log-likelihood and running time (min) on Wikipedia dataset.

k		like.	time	$\log_2 b$	iters	k	like.	time	$\log_2 b$	iters
200	Spectral	7.49	34	12	-	300	7.39	56	13	-
	Gibbs	6.85	561	-	30		6.38	818	-	30
	Hybrid	6.77	144	12	5		6.31	352	13	10

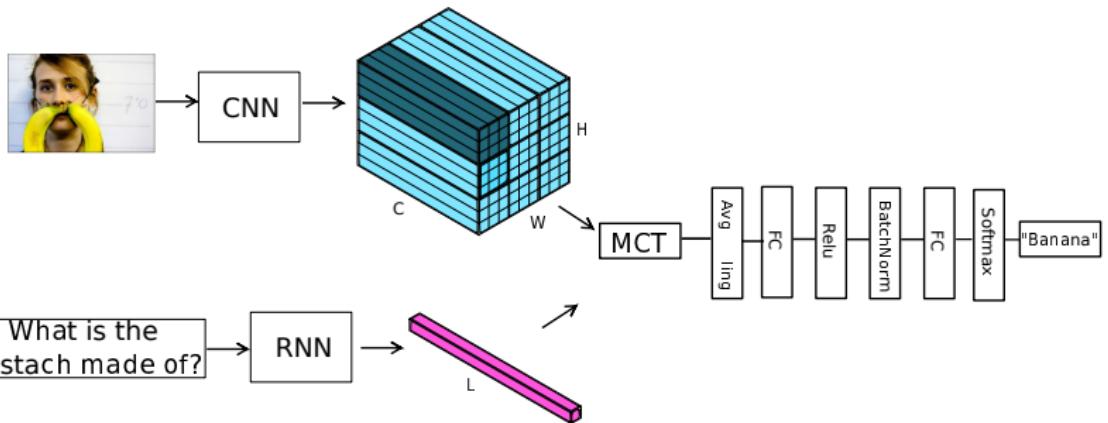


Multimodal Tensor Pooling



MCT in Visual Question & Answering

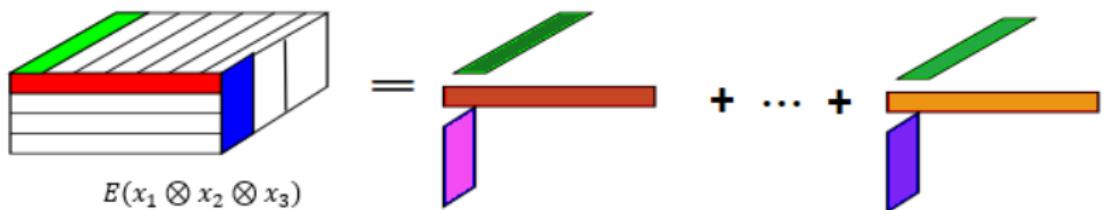
•



Tensor Decompositions

$$E(x_1 \otimes x_2) = \text{[red bar]} + \dots + \text{[orange bar]}$$


$E(x_1 \otimes x_2)$

$$E(x_1 \otimes x_2 \otimes x_3) = \text{[red block]} + \dots + \text{[orange block]}$$


$E(x_1 \otimes x_2 \otimes x_3)$

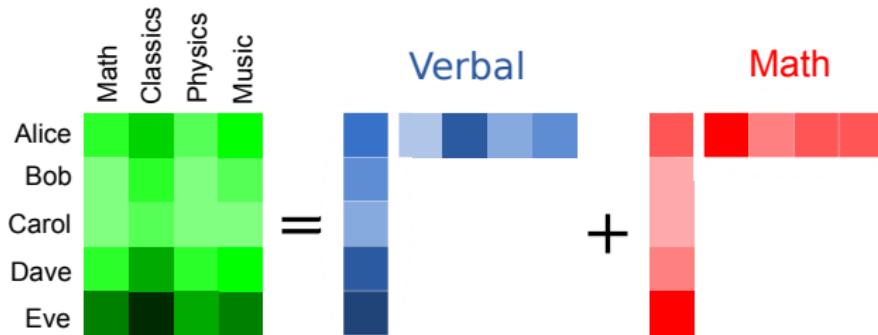
Example: Discovering Latent Factors

	Math	Classics	Physics	Music
Alice	High	High	High	High
Bob	Medium-High	Medium-High	Medium-High	Medium-High
Carol	Medium-Low	Medium-Low	Medium-Low	Medium-Low
Dave	Low	Low	Low	Low
Eve	Very Low	Very Low	Medium-Low	Medium-Low

- List of scores for students in different tests
- Learn **hidden factors** for **Verbal** and **Mathematical** Intelligence [C. Spearman 1904]

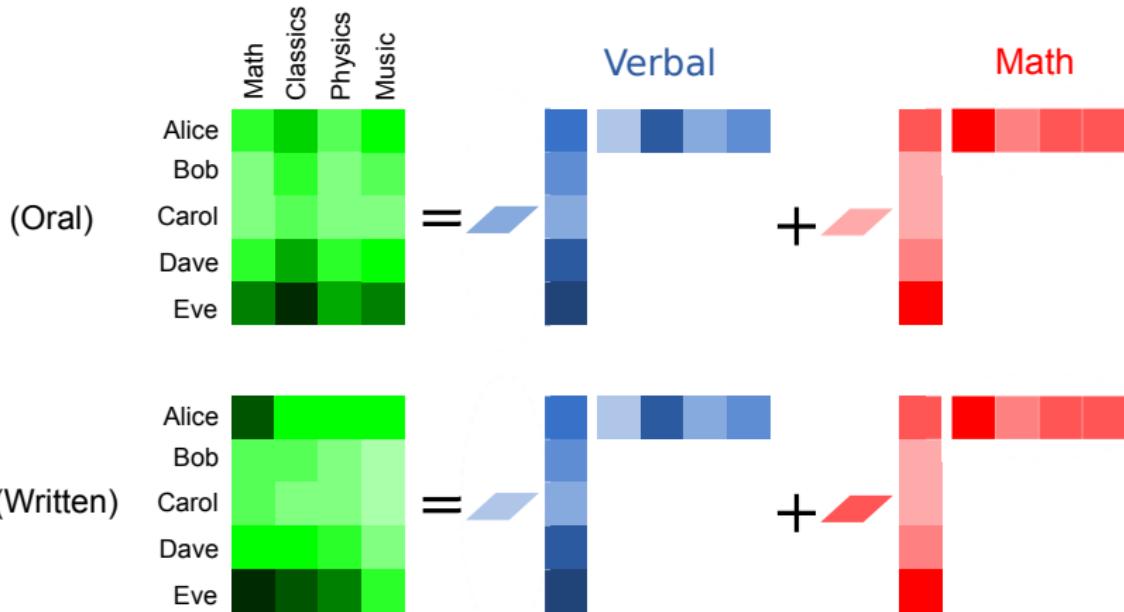
$$\begin{aligned}\text{Score}(\text{student}, \text{test}) &= \text{student}_{\text{verbal-intlg}} \times \text{test}_{\text{verbal}} \\ &\quad + \text{student}_{\text{math-intlg}} \times \text{test}_{\text{math}}\end{aligned}$$

Matrix Decomposition: Discovering Latent Factors



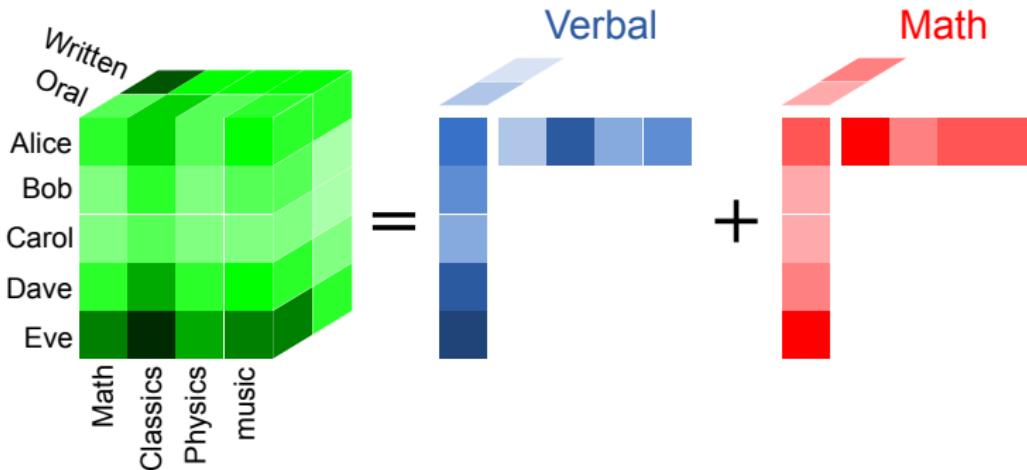
- Identifying **hidden factors** influencing the observations
- Characterized as **matrix decomposition**

Tensor: Shared Matrix Decomposition



- Shared decomposition with different scaling factors
- Combine matrix slices as a tensor

Tensor Decomposition



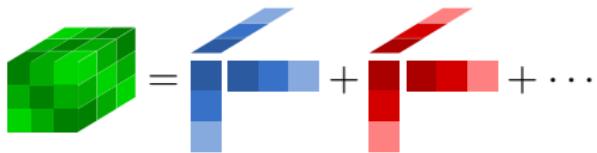
- Outer product notation:

$$T = u \otimes v \otimes w + \tilde{u} \otimes \tilde{v} \otimes \tilde{w}$$

⇓

$$T_{i_1, i_2, i_3} = u_{i_1} \cdot v_{i_2} \cdot w_{i_3} + \tilde{u}_{i_1} \cdot \tilde{v}_{i_2} \cdot \tilde{w}_{i_3}$$

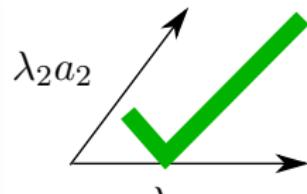
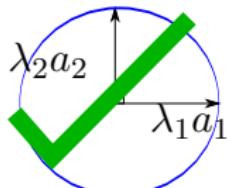
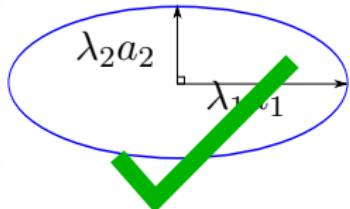
Identifiability under Tensor Decomposition



$$T = \lambda_1 a_1^{\otimes 3} + \lambda_2 a_2^{\otimes 3} + \dots,$$

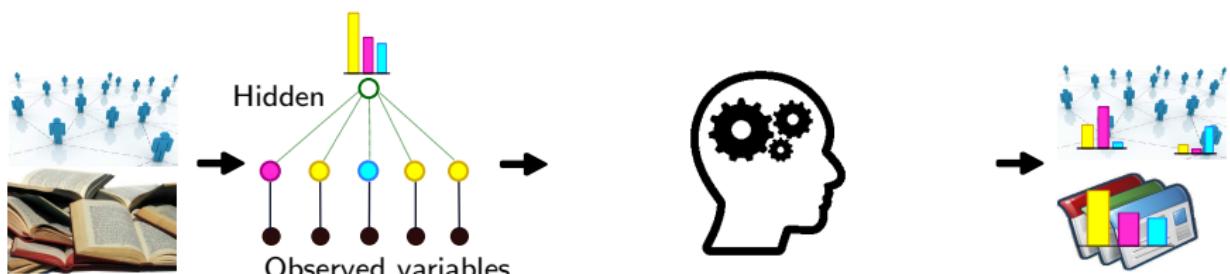
Uniqueness of Tensor Decomposition [J. Kruskal 1977]

- Above tensor decomposition: unique when rank one pairs are linearly independent
- Matrix case: when rank one pairs are orthogonal



Unsupervised Learning via Probabilistic Models

Data → Model → **Learning Algorithm** → Predictions

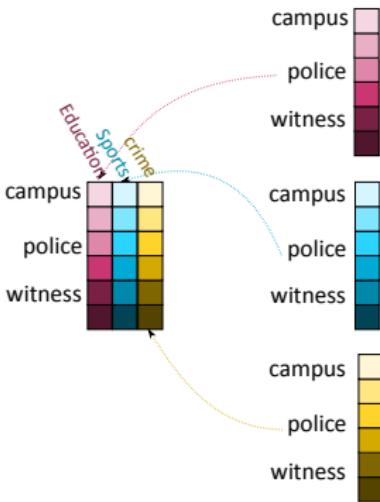


Challenges in High dimensional Learning

- Dimension of $x \gg$ dim. of latent variable h .
- Learning is like finding needle in a haystack.
- Computationally & statistically challenging.



Extracting Topics from Documents



The New York Times

COLLEGE FOOTBALL

At Florida State, Football Clouds Justice

By MARK MANTIKI and WALTER BOGDANICH OCT. 18, 2010

Now, an examination by The New York Times of police and court records, along with interviews with victims, attorneys and law enforcement officials, has raised doubt that, for years, the Florida State football program's treatment of the Winston complaint was in keeping with the way the *Florida* newspaper's editors believe have soft-pedaled allegations of wrongdoing by Seminoles football players. From criminal mischief and motor vehicle theft to domestic violence, arrests have been avoided, investigations have stalled and players have escaped serious consequences.

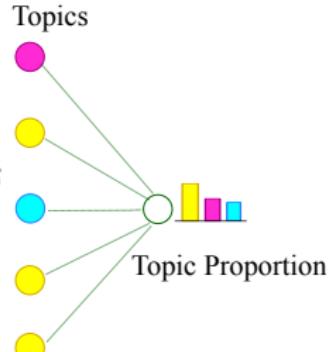
In a community where well-being and economic well-being are so tightly bound to the fortunes of the nation's top-ranked college football team, law enforcement officers are freely attached to a suspect's football connection. Those ties are cited repeatedly in *police* reports examined by The Times. What's more, dozens of officers work second jobs directing traffic and providing security at home football games and many express their devotion to the Seminoles on social media.

On Jan. 30, 2009, a female student at Florida State spotted the man she believed had raped her the previous month. After *posting* his name, Jairis Winston, the reporter told him to the Tallahassee *police*.

In the 2½ months since, Florida State officials have said little about how they handled the case, which is now *pending*. As The Times reported last April, the Tallahassee *police* also failed to respond to a formal Department of Justice request for the rape accusation. It did not become public until November, when a Tampa reporter, Matt Baker, acting on a tip, sought records of the *police investigation*.

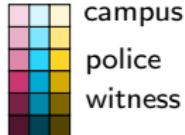
Upon learning of Mr. Baker's inquiry, Florida State, having shewn little curiosity about the rape accusation, suddenly took a keen interest in the journalist seeking to report it, according to emails obtained by The Times.

"Can you share any details on the requesting source?" David Perry, the university's *police* chief, asked the Tallahassee *police*. Several hours later, Mr.



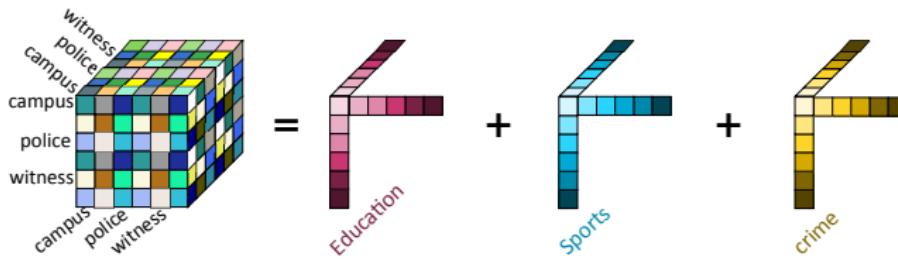
A., D. P. Foster, D. Hsu, S.M. Kakade, Y.K. Liu. "Two SVDs Suffice: Spectral decompositions for probabilistic topic modeling and latent Dirichlet allocation," NIPS 2012.

Tensor Methods for Topic Modeling



- Topic-word matrix $\mathbb{P}[\text{word} = i | \text{topic} = j]$
- Linearly independent columns

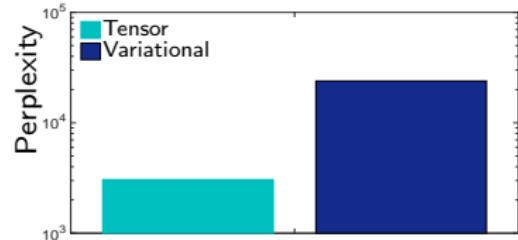
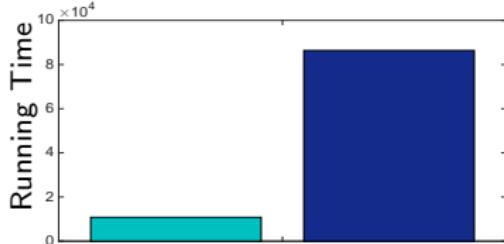
Moment Tensor: Co-occurrence of Word Triplets



Tensors vs. Variational Inference

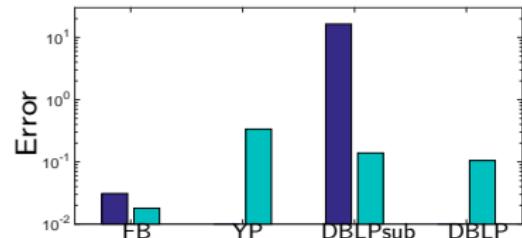
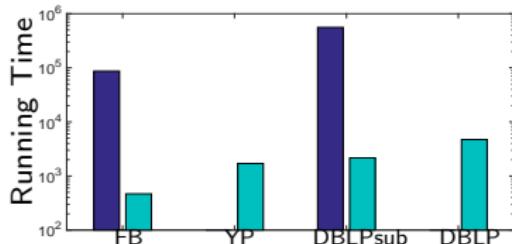
Criterion: Perplexity = $\exp[-\text{likelihood}]$.

Learning Topics from PubMed on Spark, 8mil articles



Learning network communities from social network data

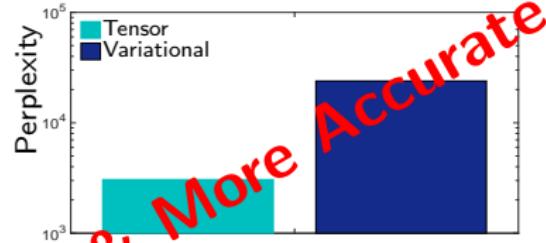
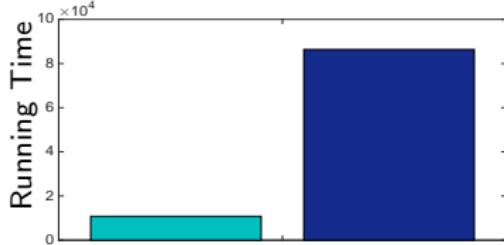
Facebook $n \sim 20k$, Yelp $n \sim 40k$, DBLP-sub $n \sim 1e5$, DBLP $n \sim 1e6$.



Tensors vs. Variational Inference

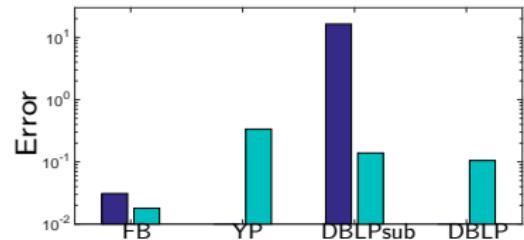
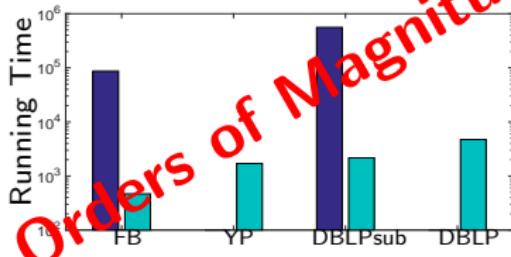
Criterion: Perplexity = $\exp[-\text{likelihood}]$.

Learning Topics from PubMed on Spark, 8mil articles



Learning network communities from social network data

Facebook $n \sim 20k$, Yelp $n \sim 40k$, DBLP-sub $n \sim 1e5$, DBLP $n \sim 1e6$.

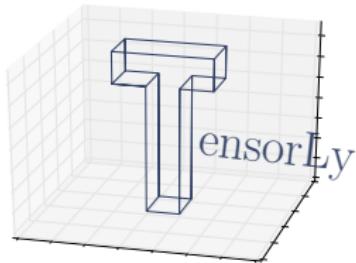


TensorLy: Tensor Learning in Python

- ▶ Pure Python
- ▶ Integrated in the Python ecosystem
- ▶ Minimal dependencies (NumPy, SciPy and Matplotlib)
- ▶ Easy to use and extend
- ▶ Fast
- ▶ Extensively tested (unit-tests)
- ▶ Exhaustive documentation
- ▶ Open source, BSD licensed

Tensorly yours,

Try it: pip install tensorly
<https://tensorly.github.io>



Contributions welcome!



3. MXNet

- Imperative and Declarative Programming
- Language Support
- Backend and Automatic Parallelization

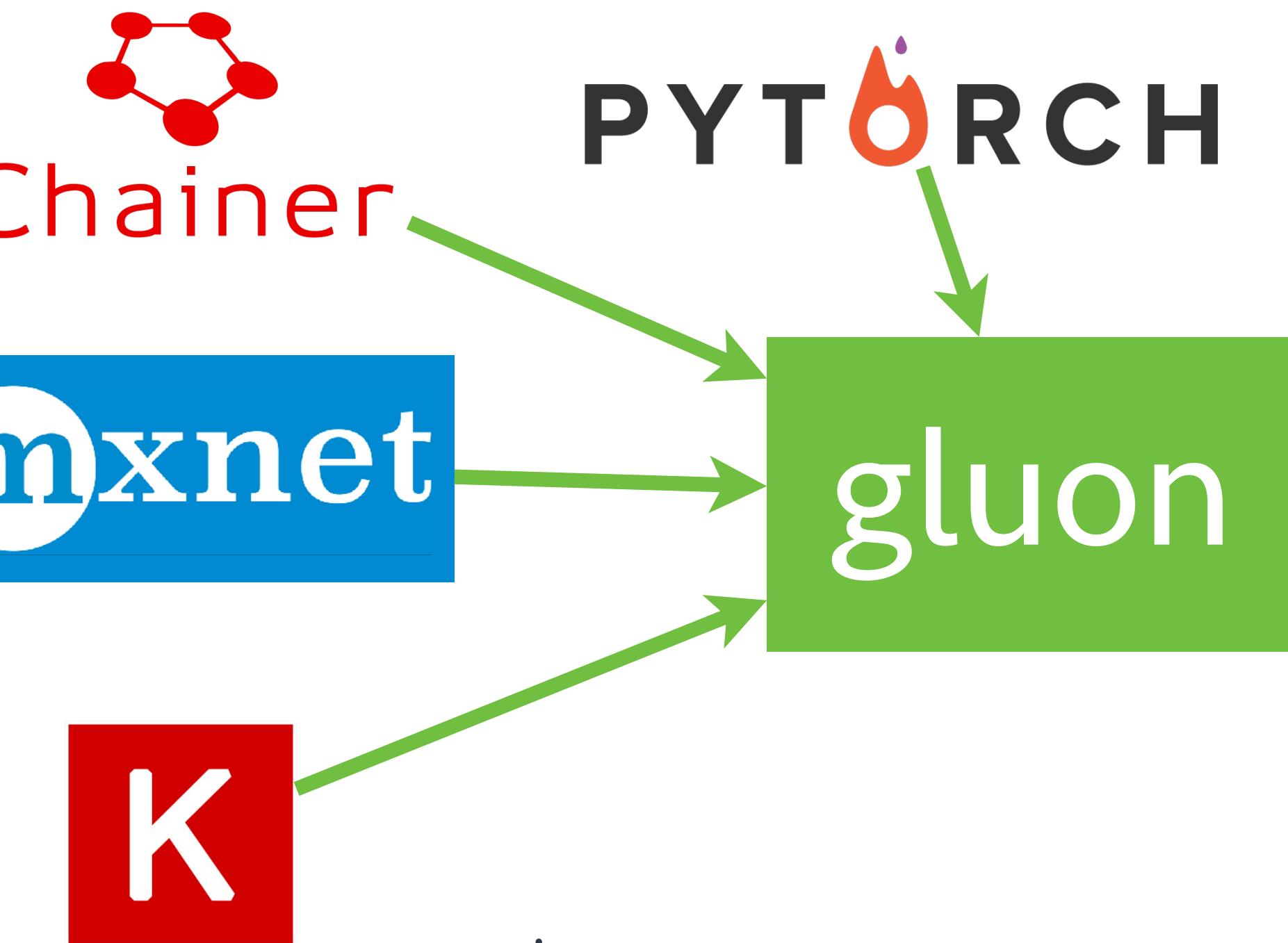
imperative

symbolic

theano

Caffe

Microsoft
CNTK



before

2012

2013

2014

2015

2016

2017

Back-end System



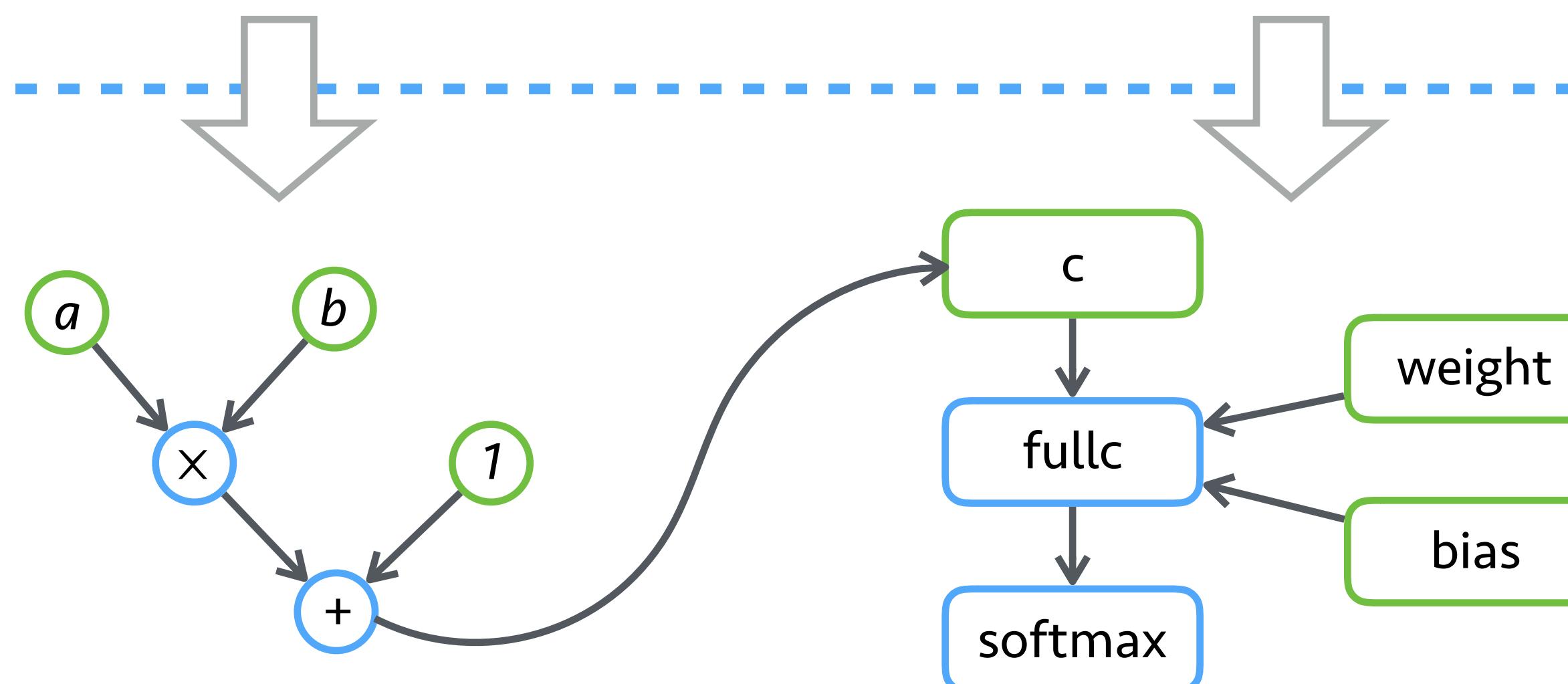
Front-end

Back-end



```
import mxnet as mx  
a = mx.nd.zeros((100, 50))  
b = mx.nd.ones((100, 50))  
c = a * b  
c += 1
```

```
import mxnet as mx  
net = mx.symbol.Variable('data')  
net = mx.symbol.FullyConnected(  
    data=net, num_hidden=128)  
net = mx.symbol.SoftmaxOutput(data=net)  
texec = mx.module.Module(net)  
texec.forward(data=c)  
texec.backward()
```



- ❖ Optimization
 - ✓ Memory optimization
 - ✓ Operator fusion
- ❖ Scheduling
 - ✓ Auto-parallelization

In summary

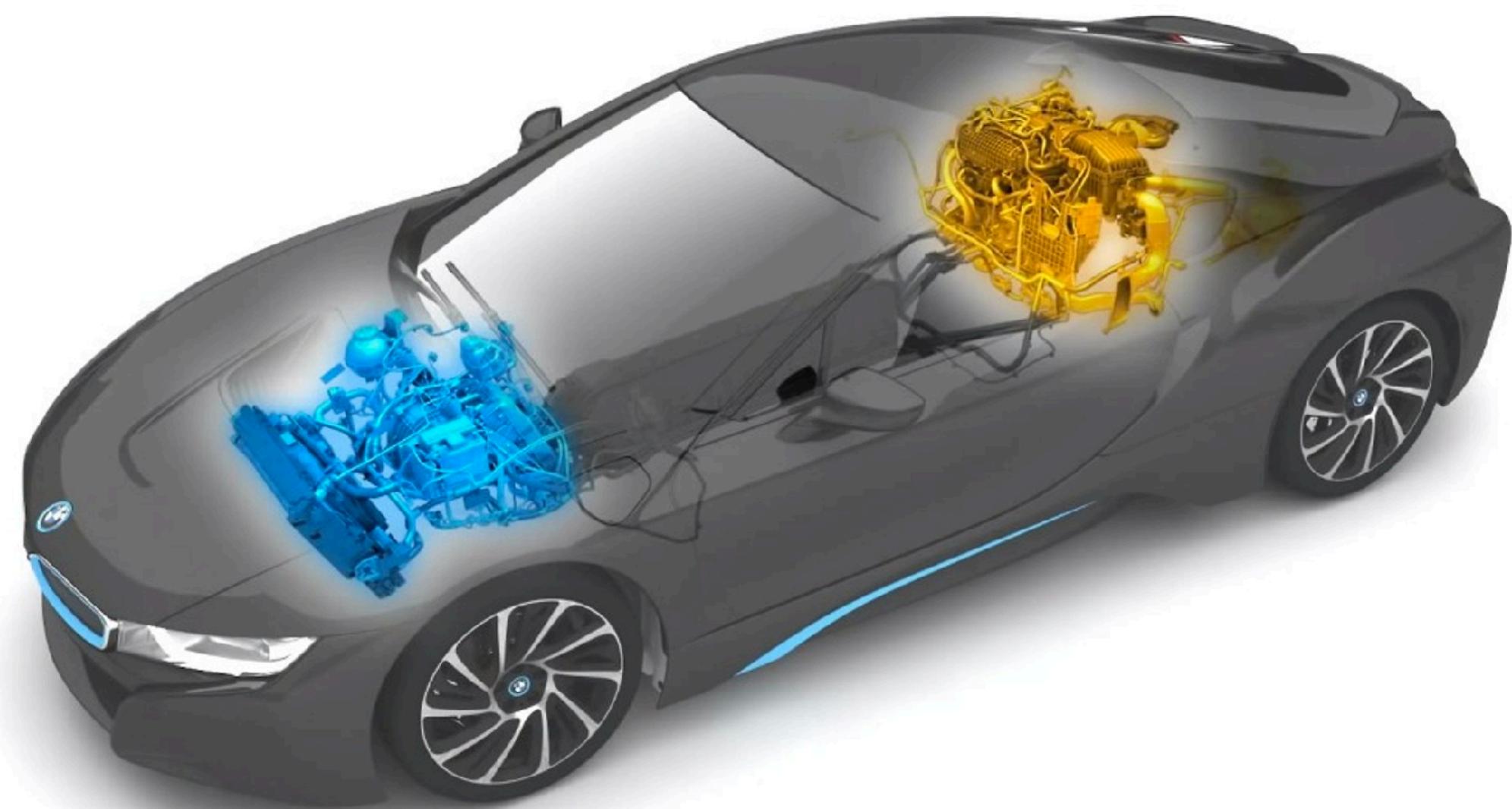
♦ Symbolic

- ❖ efficient & portable
- ❖ but hard to use



♦ Gluon

- ❖ imperative for developing
- ❖ symbolic for deploying



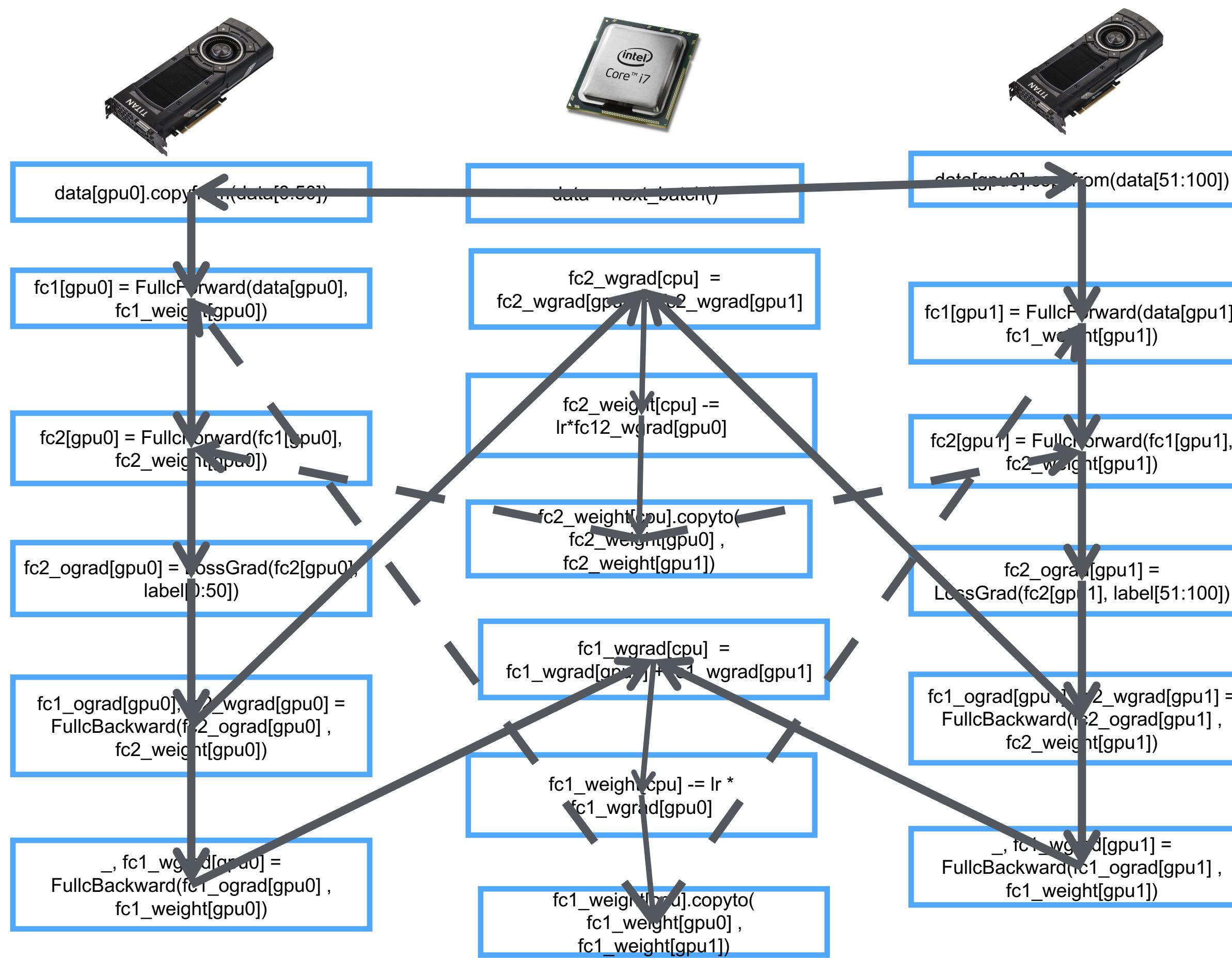
♦ Imperative

- ❖ flexible
- ❖ may be slow



Writing Parallel Programs is Painful

Dependency graph for 2-layer neural networks with 2 GPUs



Each forward-backward-update involves $O(\text{num_layer})$, which is often 100—1,000, tensor computations and communications

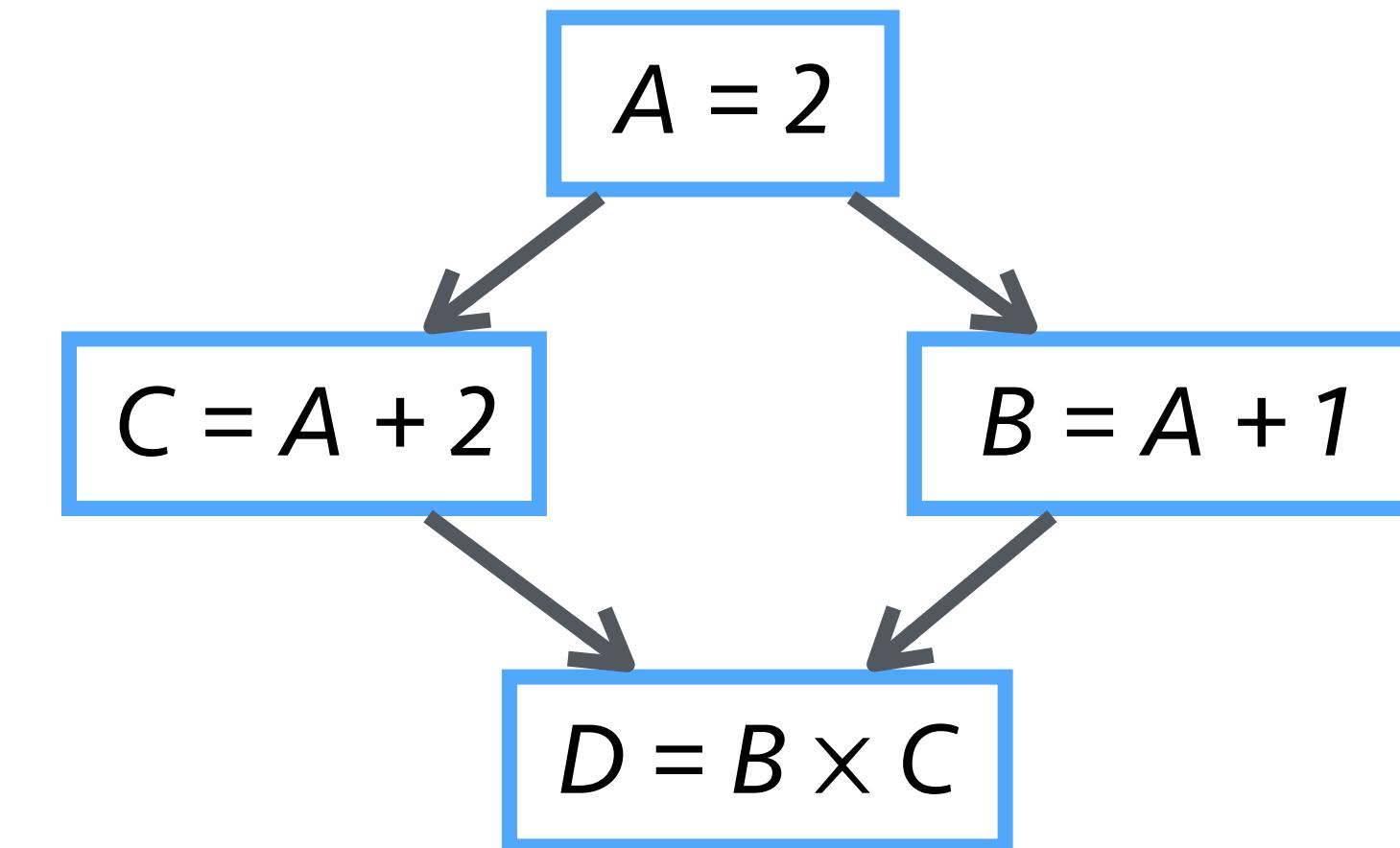


Auto Parallelization

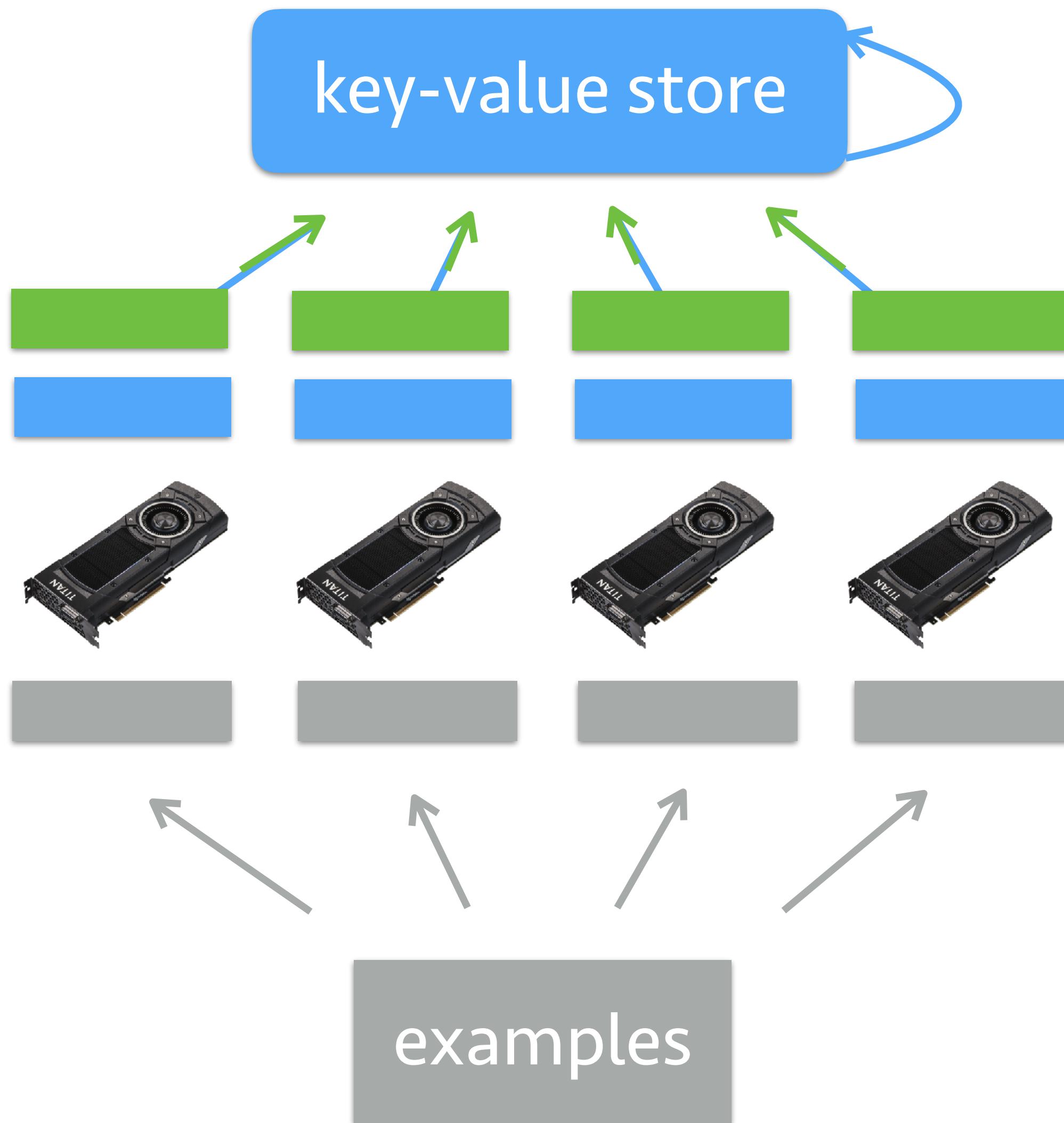
Write **serial** programs

```
>>> import mxnet as mx  
>>> A = mx.nd.ones((2,2)) *2  
>>> C = A + 2  
>>> B = A + 1  
>>> D = B * C  
>>> D.wait_to_read()
```

Run in **parallel**



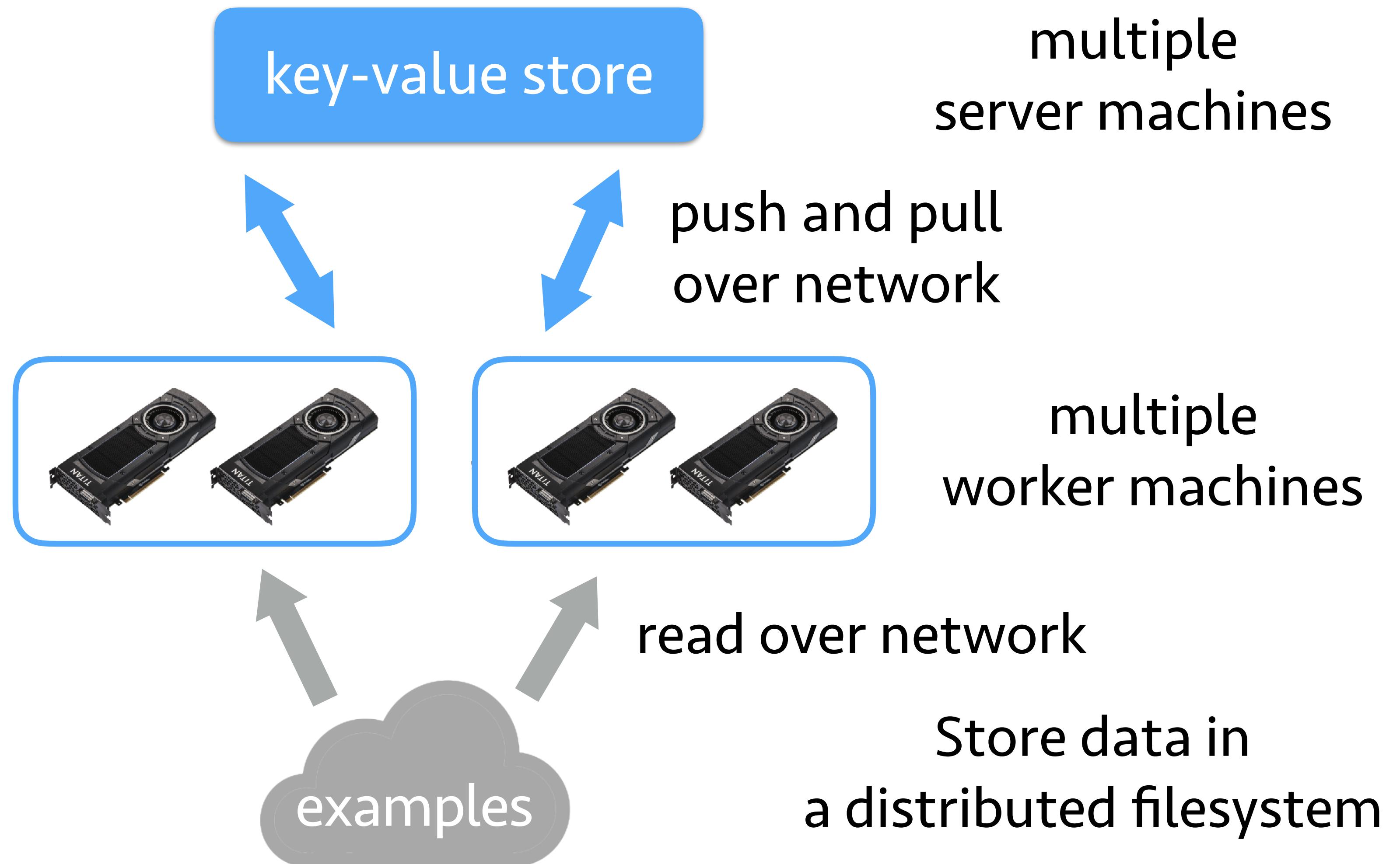
Data Parallelism



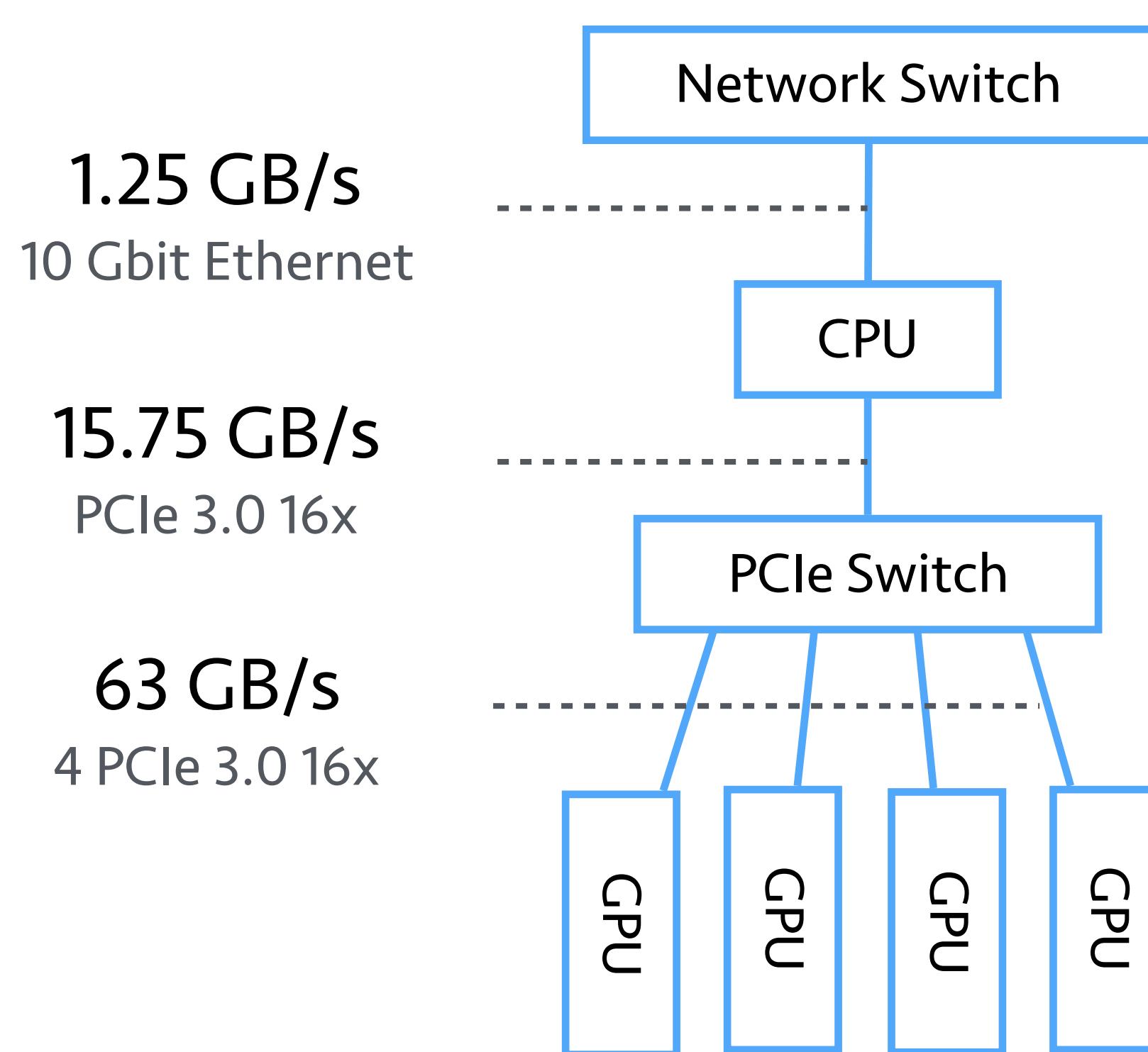
1. Read a data partition
2. Pull the parameters
3. Compute the gradient
4. Push the gradient
5. Update the parameters

Distributed Computing

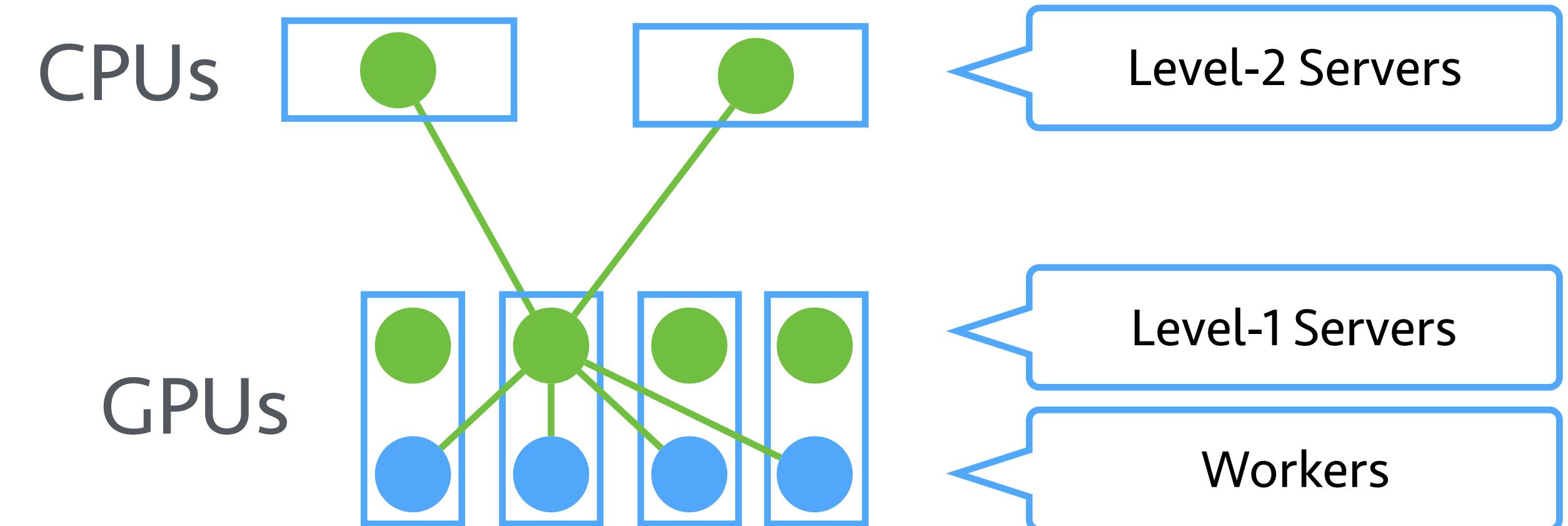
A user does not need
to change the codes
when using multiple
machines



Scale to Multiple GPU Machines

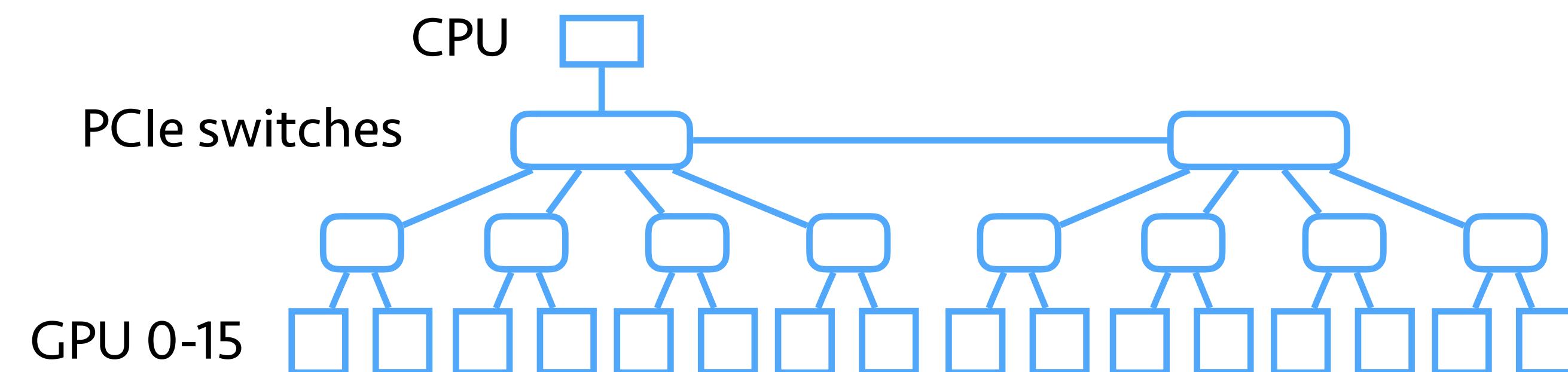


Hierarchical parameter server

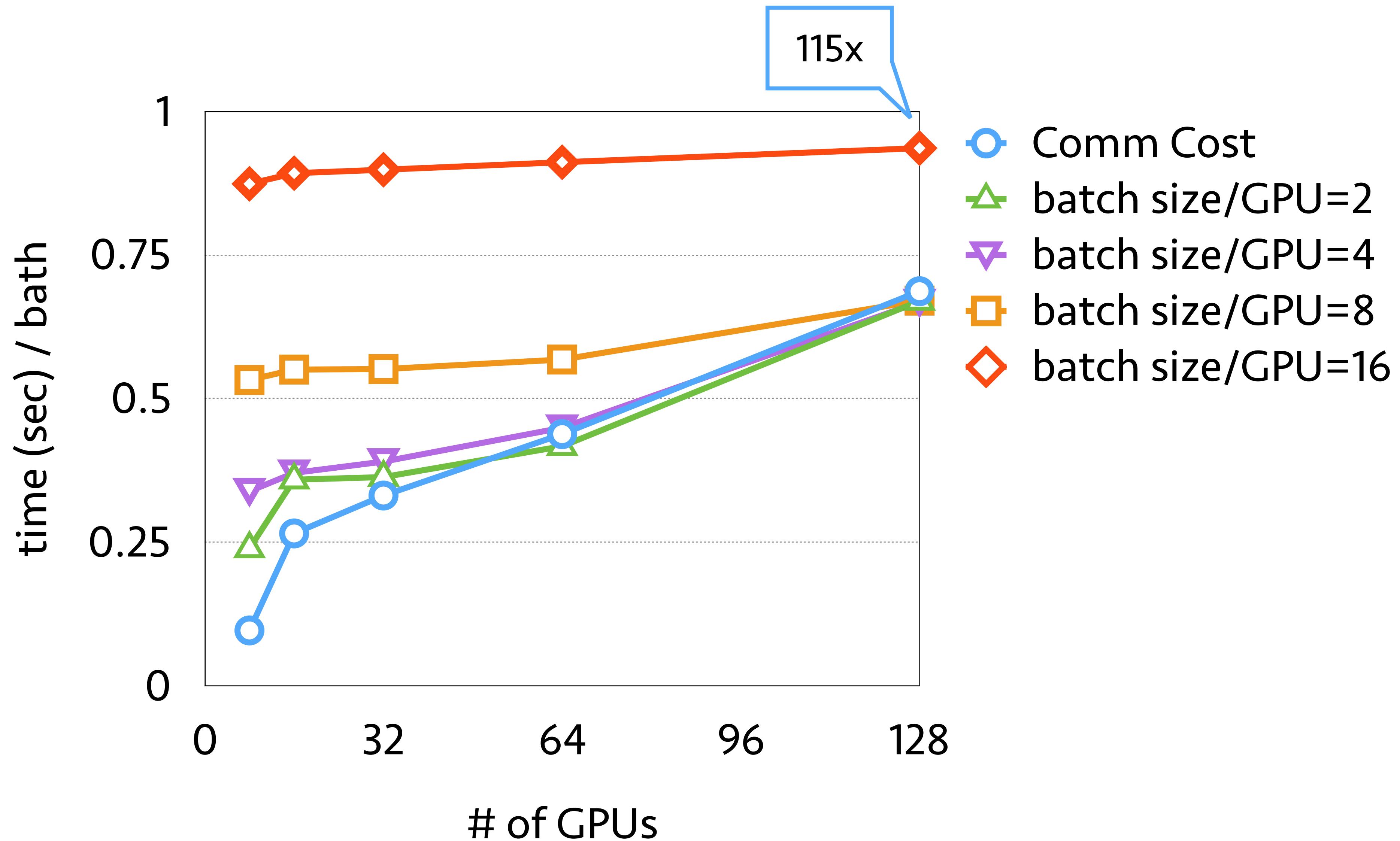


Experiment Setup

- ❖ IMAGENET
 - ✓ 1.2 million images with 1000 classes
- ❖ Resnet 152-layer model
- ❖ EC2 P2.16xlarge
- ❖ Minibatch SGD
- ❖ Synchronized Updating

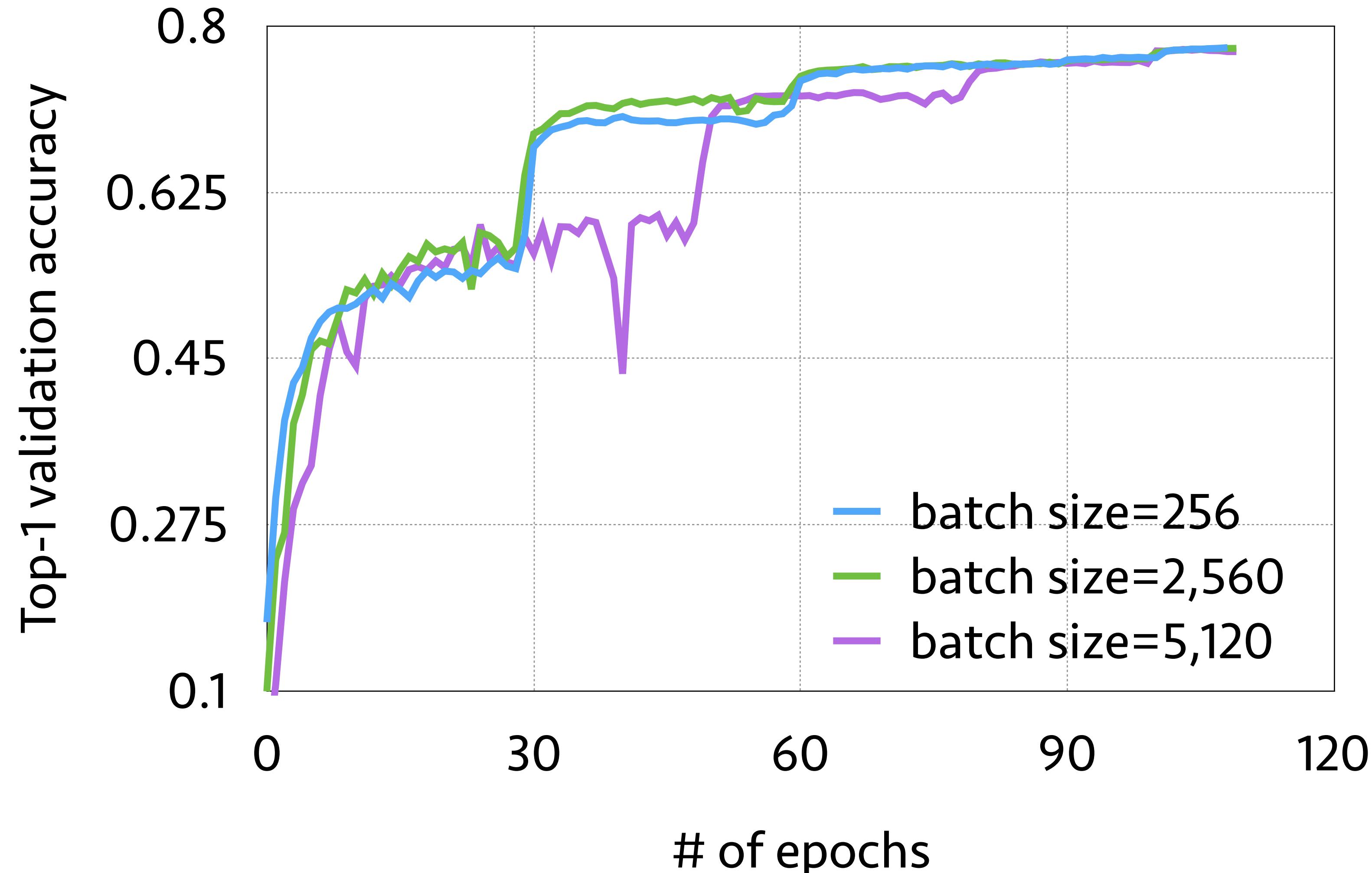


Scalability over Multiple Machines

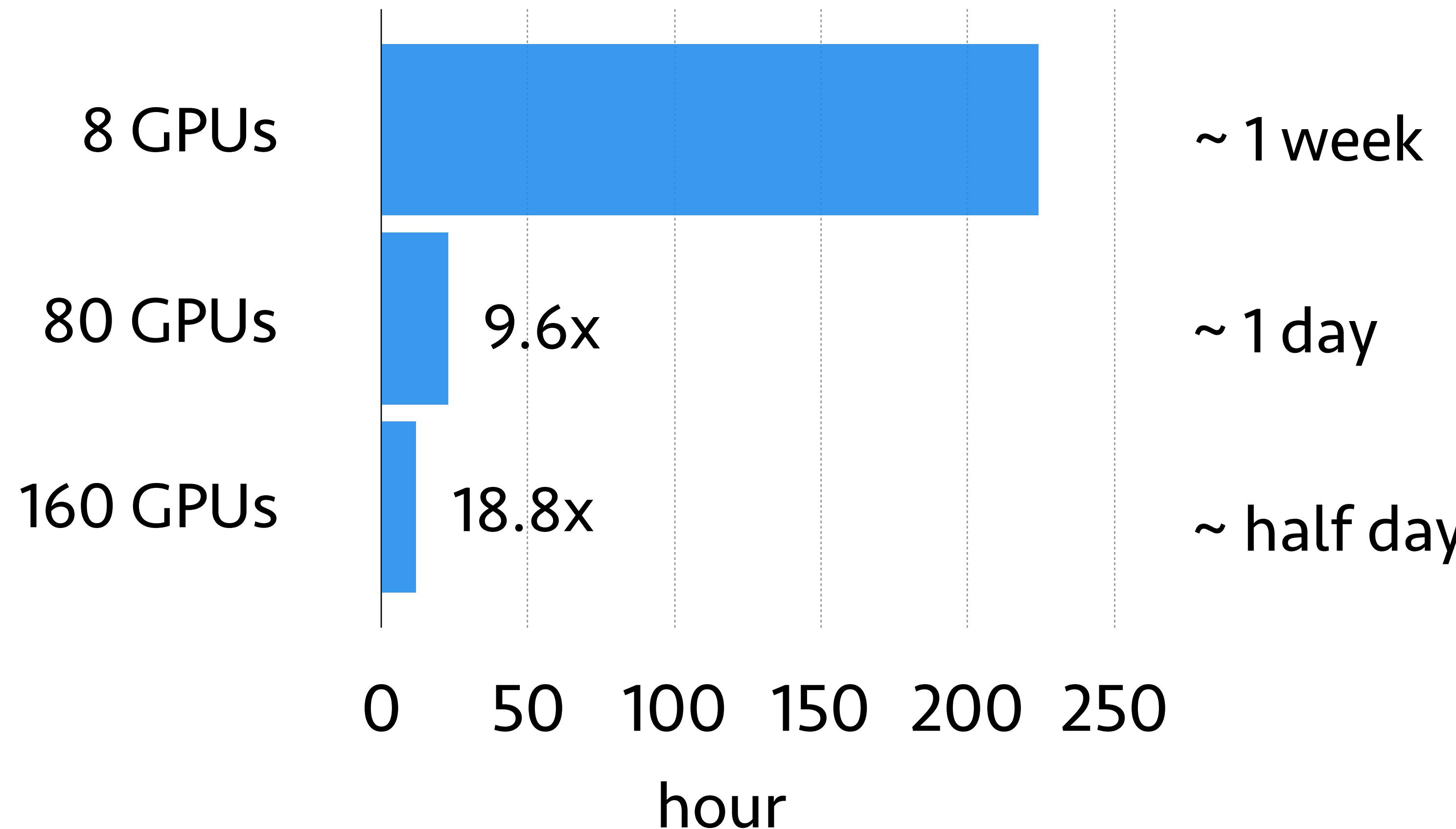


Convergence

- ❖ Increase learning rate by 5x
- ❖ Increase learning rate by 10x, decrease it at epoch 50, 80



Time to achieve 22.5% top-1 accuracy



Outline

1 Introduction

2 Tensor Contractions

3 Speeding up Tensor Contractions

4 Tensor Sketches

5 Conclusion

Conclusion

Tensors are the future of ML

- Tensor contractions: space savings in deep architectures.
- New primitives speed up tensor contractions: extended BLAS
- Tensor sketches compress efficiently while preserving information
- Tensor sketches enable multi-modal tasks such as visual Q&A

