

Accelerating Generalized Linear Models with MLWeaving: A One-Size-Fits-All System for Any-Precision Learning

Zeke Wang, Kaan Kara, Hantian Zhang, Gustavo Alonso, Onur Mutlu, Ce Zhang

Systems Group, ETH Zurich



Outline

• • • •

Quick Background (5mins)

Stochastic Gradient Descent (SGD)

Synchronous vs. Asynchronous

Low Precision

MLWeaving (10mins)

Arbitrary-precision Training

MLWeaving Memory Layout

MLWeaving Hardware Design

Efficient Synchronous Design

OK, how does SGD work?



Stochastic Gradient Descent (SGD)

P1: Model can be staled, especially when running on multiple cores.

Linear Regression

$$\min_{x} \frac{1}{2} \sum_{r} (A_r x^T - b_r)^2$$

$$A_r = get_data()$$

$$x = get_model()$$

$$g = comp_grad(x, A_r)$$

x = x - g

set_model(x)

Two Interesting Properties

SGD on the CPU: synchronous or asynchronous?

Sync. Single-Core SGD: Low Throughput



Causes Problem When Using Multiple Cores.

Async. Multi-Core SGD: High Throughput

Multi-core SGD relies on asynchrony.

HogWild! [1]

ModelAverage [2]

[1] Hogwild: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent. In NIPS. 2011.

[2] Parallelized Stochastic Gradient Descent. In NIPS. 2010.

Synchrony vs. Asynchrony on CPUs

• • • • •

	Hardware Efficiency (Throughput)	Statistical Efficiency (Convergence Rate)
Single-core SGD (Synchrony)	Low 🙁	High 😃
Multi-core SGD (Asynchrony)	High 🕛	Low 🙁

Why Low Precision?

Why Low Precision?

• • • •



Full precision

1.310245

X 0.602069

0.788857897

Low precision

about 1.3

X about 0.6

about 0.78

Relax, It is only Machine Learning.

Current Hardware Supports Limited Precisions

• • • • •



GPU









Char (8-bit), Short (16-bit)

FP8 (8-bit), FP16 (16-bit)

INT8 (8-bit)

Goal of This Work

• • • •



For Generalized Linear Model training, can we enable things that cannot be well done on CPUs?



Outline

• • • •

Quick Background (5mins)

Stochastic Gradient Descent (SGD)

Synchronous vs. Asynchronous

Low Precision

MLWeaving (10mins)

Arbitrary-precision Training

MLWeaving Memory Layout

MLWeaving Hardware Design

Efficient Synchronous Design

Two Goals of Arbitrary-precision Training

1, One hardware design and one copy of dataset support any-precision training.

2, Our design achieves linear speedup with lower precision.

MLWeaving Memory Layout



 $1^{st} \text{ feature} \qquad 2^{nd} \text{ feature}$ $1^{st} \text{ row A} \quad A_1^{[1]} \quad A_1^{[2]} \quad A_1^{[3]} \quad A_1^{[4]} \rightarrow A_2^{[1]} \quad A_2^{[2]} \quad A_2^{[3]} \quad A_2^{[4]} \rightarrow A_2^{[4]}$ $2^{nd} \text{ row B} \quad B_1^{[4]} \quad B_1^{[2]} \quad B_2^{[3]} \quad B_3^{[3]} \quad B_3^{[4]} \rightarrow B_2^{[4]} \quad B_2^{[2]} \quad B_2^{[3]} \quad B_2^{[4]} \qquad B_2^{[4]} \rightarrow B_2^{[4]$

MLWeaving:

Observation 1: Often memory bandwidth bound

1st row A

Observation 2: Low precision (e.g., 8 bit fixed point) often provides reasonable quality

Observation 3: Different training task might need different precision level even on the same dataset

Can we store the data in a new data structure that efficiently supports arbitrary precision data movement?

MLWeaving Memory Layout



1st feature 2nd feature 1st row A $A_1^{[1]}$ $A_1^{[2]}$ $A_1^{[3]}$ $A_1^{[4]} \rightarrow A_2^{[1]}$ $A_2^{[2]}$ $A_2^{[3]}$ $A_2^{[4]}$ 2nd row B $B_1^{[3]}$ $B_1^{[2]}$ $B_1^{[3]}$ $B_1^{[4]}$ $B_2^{[4]}$ $B_2^{[4]}$ $B_2^{[2]}$ $B_2^{[2]}$ $B_2^{[4]}$

MLWeaving:

Observation 1: Often memory bandwidth bound

1st row A

Observation 2: Low precision (e.g., 8 bit fixed point) often provides reasonable quality

Observation 3: Different training task might need different precision level even on the same dataset

Can we store the data in a new data structure that efficiently supports arbitrary precision data movement?

MLWeaving Memory Layout





Observation 1: Often memory bandwidth bound

Observation 2: Low precision (e.g., 8 bit fixed point) often provides reasonable quality

Observation 3: Different training task might need different precision level even on the same dataset

Can we store the data in a new data structure that supports arbitrary precision data movement?



MLWeaving:

1st row A

MLWeaving Memory Layout





Observation 1: Often memory bandwidth bound

Observation 2: Low precision (e.g., 8 bit fixed point) often provides reasonable quality

Observation 3: Different training task might need different precision level even on the same dataset

Can we store the data in a new data structure that supports arbitrary precision data movement?



MLWeaving:



MLWeaving Memory Layout





Observation 1: Often memory bandwidth bound

Observation 2: Low precision (e.g., 8 bit fixed point) often provides reasonable quality

Observation 3: Different training task might need different precision level even on the same dataset

Can we store the data in a new data structure that supports arbitrary precision data movement?



MLWeaving Memory Layout



Observation 1: Often memory bandwidth bound

Observation 2: Low precision (e.g., 8 bit fixed point) often provides reasonable quality

Observation 3: Different training task might need different precision level even on the same dataset

Can we store the data in a new data structure that supports arbitrary precision data movement?



<u>MLWeaving does not work out on CPUs</u>. CPU does not have custom instruction for MLWeaving memory layout and then we have to group bits from different memory locations before the further computing.

Outline

• • • •

Quick Background (5mins)

Stochastic Gradient Descent (SGD)

Synchronous vs. Asynchronous

Low Precision

MLWeaving (10mins)

Arbitrary-precision Training

MLWeaving Memory Layout

MLWeaving Hardware Design

Efficient Synchronous Design

MLWeaving Hardware Design: Key Idea

• • • •

MLWeaving memory layout:

Key idea of MLWeaving hardware design:



To use bit-serial multiplier to enable efficient data processing from the MLWeaving memory layout.

The modern CPU does not support bit-serial multiplier.

Please read our paper about how bit-serial multiplier works.

Custom Computation for MLWeaving

MLWeaving hardware design:





Bit-serial multiplier + MLWeaving memory layout enable any-precision ML training.

MLWeaving's Performance: Almost Linear Speedup with Lower Precision



Memory traffic vs. Precision 3231.8

1615.9

16

32



• • • •

Quick Background (5mins)

Stochastic Gradient Descent (SGD)

Synchronous vs. Asynchronous

Low Precision

MLWeaving (10mins)

Arbitrary-precision Training

MLWeaving Memory Layout

MLWeaving Hardware Design

Efficient Synchronous Design

Synchronous SGD or asynchronous SGD on custom hardware?

SGD on Custom Hardware: The Best of Two Worlds

• • • • •

	Hardware Efficiency (Throughput)	Statistical Efficiency (Convergence Rate)
Single-core (Synchrony)	Low 🙁	High 😃
Multi-core (Asynchrony)	High 🕛	Low 🙁
Custom hardware (Synchrony)	High 🕛	High 🕛

Original Synchronous Implementation: Compute-Bound

• • • • •



Original Synchronous Implementation: Compute-Bound

• • • •



Optimal Synchronous Implementation: Memory-Bound Cvcles 1st batch 2rd batch Model Read Original: Compute-bound 1st batch Dot 2rd batch Product 1st gradient 2rd gradient Model **Observation:** Custom hardware can Write update the model (thousands of weights) at the granularity level: 64 Cycles 1st batch 2rd batch Model weights, not the whole model. Read 1st batch Dot 2rd batch With Chaining: Memory-bound Product 2rd gradient Model Write

Optimal Synchronous Implementation: Memory-Bound Cvcles 1st batch 2rd batch Model Read Original: Compute-bound 1st batch Dot 2rd batch Product 1st gradient 2rd gradient Model **Observation:** Custom hardware can Write update the model (thousands of weights) at the granularity level: 64 Cycles 1st batch 2rd batch Model weights, not the whole model. Read 1st batch Dot 2rd batch With Chaining: Memory-bound Product 1st gradient 2rd gradient Model Write *High throughput: "sync" is as fast as "async".* **Gap**: gradient from 64 weights

Effect of Sync. Design

• • • •



Training loss vs. Number of Epochs

ModelAverage and *Hogwild* on the multi-core CPU: Async. *MLWeaving* on the *custom hardware* : Sync.

Outline

• • • •

Quick Background (5mins)

Stochastic Gradient Descent (SGD)

Synchronous vs. Asynchronous

Low Precision

MLWeaving (10mins)

Arbitrary-precision Training

MLWeaving Memory Layout

MLWeaving Hardware Design

Efficient Synchronous Design

System Integration

• • • •

Intel HARP2 Platform:

Intel Broadwell 14-core CPU Intel Arria 10 FPGA

We integrate MLWeaving into DoppioDB that is an open source solution for FPGA-enhanced databases based on MonetDB. (1) CREATE INDEX mlweaving_on_t1 ON create_mlweaving('t1');
(2) CREATE INDEX model_on_t1 ON train_mlweaving('t1', numEpochs, ...);
(3) SELECT * FROM infer_mlweaving(model_on_t1, 't1', labelIndex);



If you are interested to play with MLWeaving, please visit our demo DoppioDB 2.0 in the Hollywood ballroom, 11:00-12:30, 27-29 Aug, 2019. End-to-End Performance: MLWeaving

• • • •



ModelAverage and Hogwild on an Intel CPU: 14 cores, AVX2-enhanced, 8-bit dataset.

MLWeaving on an FPGA: 3-bit dataset.

If you are interested to play with MLWeaving, please visit our demo DoppioDB 2.0 in the Hollywood ballroom, 11:00-12:30, 27-29 Aug, 2019.



Any Questions?



For GLM training, can we enable things that cannot be well done on CPUs/GPUs?

Any-precision Training

High-throughput Sync. Design

Problem? Cache-coherence is expensive, especially for dense data!

How Bit-serial Multiplier Deals with Low Precision?

• • • • •	4321
<u>4-bit:</u>	X 0 0 2 0
	<mark>8642</mark> 0
	4 3 2 0
<u>3-bit:</u>	X 0 0 2 0
	86400
	4300
<u>2-bit:</u>	X 0 0 2 0
	86000
	4000
<u>1-bit:</u>	X 0 0 2 0
	80000
No	rmal Multiplier

Each bit should be binary, but we use decimal for ease of understanding.

How Bit-serial Multiplier Deals with Low Precision?

• • • •	• 4 3 2 1	Initialization
<u>4-bit:</u>	X 0 0 2 0	<u>mittanzation.</u>
	86420	4 3 2 1
	4 3 2 0	
<u>3-bit:</u>	X 0 0 2 0	
	86400	
	4300	
<u>2-bit:</u>	X 0 0 2 0	
	86000	Sum = 0.0000
	4000	
<u>1-bit:</u>	X 0 0 2 0	
	80000	
	Normal Multiplier	Bit-serial Multiplier (BSM)

How Bit-serial Multiplier Deals with Low Precision?

Bit-serial Multiplier: 1-Bit Precision

Bit-serial Multiplier: 1-Bit Precision

Bit-serial Multiplier: 1-Bit Precision

Bit-serial Multiplier: 2-Bit Precision

Normal Multiplier

Bit-serial Multiplier: 2-Bit Precision

Normal Multiplier

Bit-serial Multiplier: 2-Bit Precision

Normal Multiplier

Bit-serial Multiplier: 3-Bit Precision

• • • • •	4321	3th Cycle.	
<u>4-bit:</u>	X 0 0 2 0	<u>J Cycle.</u>	
_	86420	4 3 2 1	
	4 3 2 0		
<u>3-bit:</u>	X 0 0 2 0		Momory
-	86400		Memory
	4300		Hardware
<u>2-bit:</u>	X 0 0 2 0		
-	86000	$S_{\rm U}m = \frac{8}{6} \frac{6}{0} \frac{0}{0} \frac{0}{0}$	
	4000		
<u>1-bit:</u>	X 0 0 2 0		
	80000		

Normal Multiplier

Bit-serial Multiplier: 3-Bit Precision

Normal Multiplier

Bit-serial Multiplier: 3-Bit Precision

Bit-serial Multiplier: 4-Bit Precision

• • • • •	4 3 2 1	Ath Cycle.
<u>4-bit:</u>	X 0 0 2 0	
_	<mark>8642</mark> 0	4 3 2 1
	4320	
<u>3-bit:</u>	X 0 0 2 0	Momory
	86400	MEIIIOIY
	4300	Hardware
<u>2-bit:</u>	X 0 0 2 0	
-	86000	$Sum = \frac{86400}{100}$
	4000	
<u>1-bit:</u>	X 0 0 2 0	
	80000	

Normal Multiplier

Bit-serial Multiplier: 4-Bit Precision

• • • • •	4 3 2 1	Ath Cycle.
<u>4-bit:</u>	X 0 0 2 0	<u></u>
-	<mark>8642</mark> 0	4 3 2 1
	4320	
<u>3-bit:</u>	X 0 0 2 0	1 means 1. Momory
	864 00	Internuty
	4300	Hardware Hardware
<u>2-bit:</u>	X 0 0 2 0	X 0020 DSIVI Sum += 20 * 1
	86000	Sum = 86400
	4000	
<u>1-bit:</u>	X 0 0 2 0	
	80000	

Normal Multiplier

Bit-serial Multiplier: 4-Bit Precision

• • • • •	4 3 2 1	Ath Cycle.
<u>4-bit:</u>	X 0 0 2 0	
_	<mark>8642</mark> 0	4 3 2 1
	4320	
<u>3-bit:</u>	X 0 0 2 0	Momory
	86400	MEMORY
	4300	Hardware Hardware
<u>2-bit:</u>	X 0 0 2 0	
-	86000	Sum = 86420
	4000	Done with 4-bit precision
<u>1-bit:</u>	X 0 0 2 0	
-	80000	

Normal Multiplier