

Low Precision Arithmetic and Quantization

CS6787 Lecture 9 — Spring 2026

Final Project Proposal Discussion

Split up into groups of 4–5

Do not be in a group with your project partners

Each of you presents a 2-minute pitch

Then discuss after everyone has pitched

Reminder: Final Project Requirements

- **Implement a machine learning system** to solve a problem
- Use one or more of the **techniques we discussed in class**
 - The mere use of a LLM in the project does not constitute a technique
- To achieve an **improvement over some baseline method**
 - Measuring both statistical performance and hardware performance
 - Or at least evaluate and attempt to achieve such a speedup
- Otherwise, very **open-ended**
 - **Groups of up to three**

Project proposals due **NEXT MONDAY**

- The main body should be about one page in length.
- It should describe the project you intend to do.
- It should contain at least one citation of a relevant paper that we did not cover in class.
- It should include some preliminary or exploratory work you've already done, that helps to support the idea that your project is feasible.
 - Don't need a lot of work, just a nonzero amount of work supporting feasibility.
- In addition to the one-page text proposal, one short **experiment plan** per person

Experiment plan

- The hypothesis
- The proxy
- The protocol
- Expected results

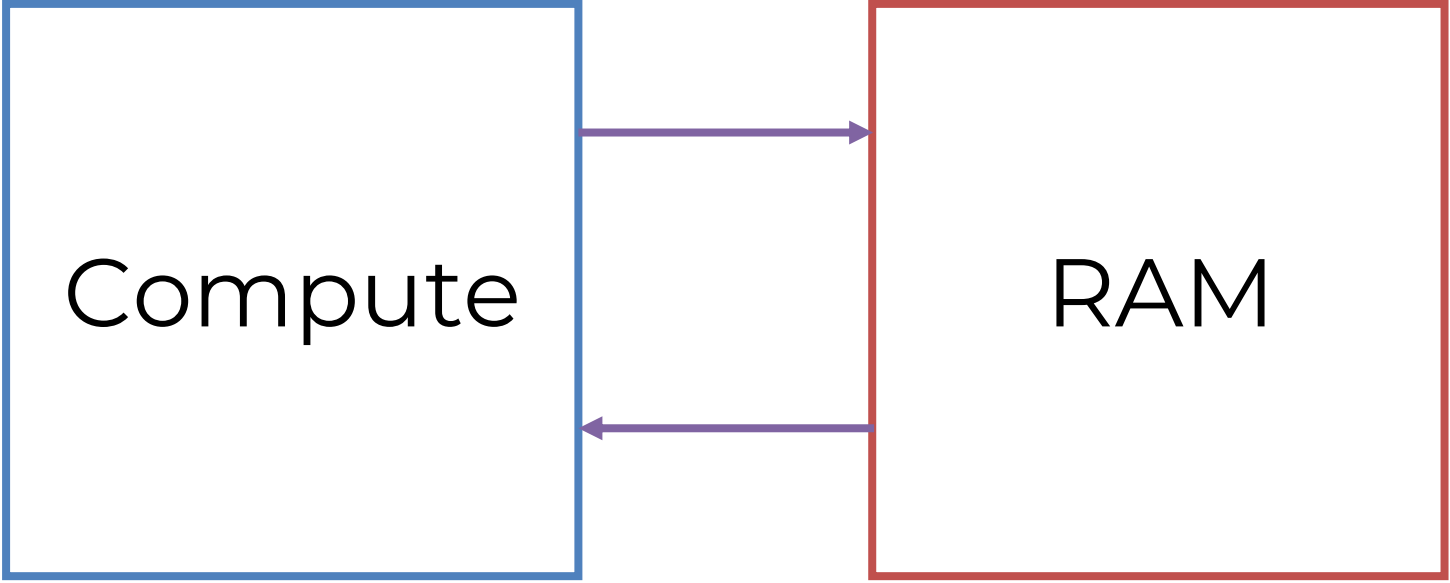
Low Precision Arithmetic and Quantization

CS6787 Lecture 9 — Spring 2026

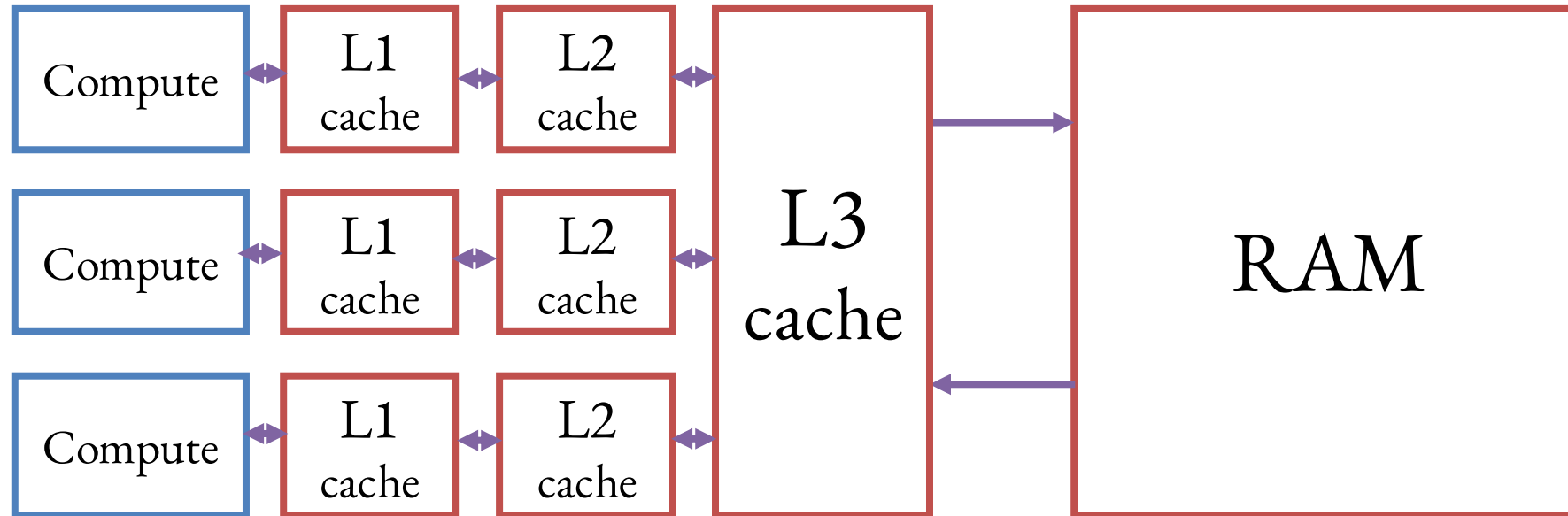
Memory as a Bottleneck

- So far, we've just been talking about **compute**
 - e.g. techniques to decrease the amount of compute by decreasing iterations
- But machine learning systems need to process **huge amounts of data**
- Need to **store, update, and transmit** this data
- As a result: **memory** is of critical importance
 - Many applications are memory-bound

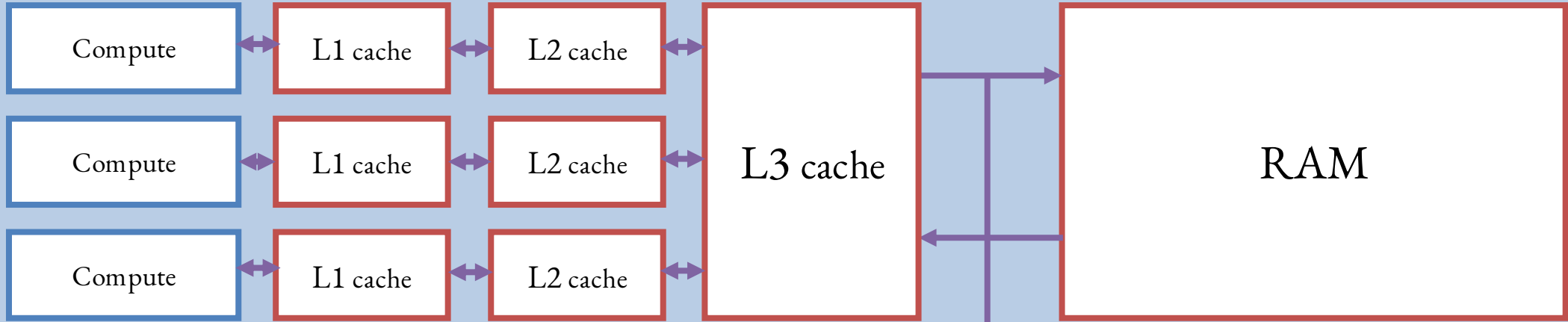
Memory: The Simplified Picture



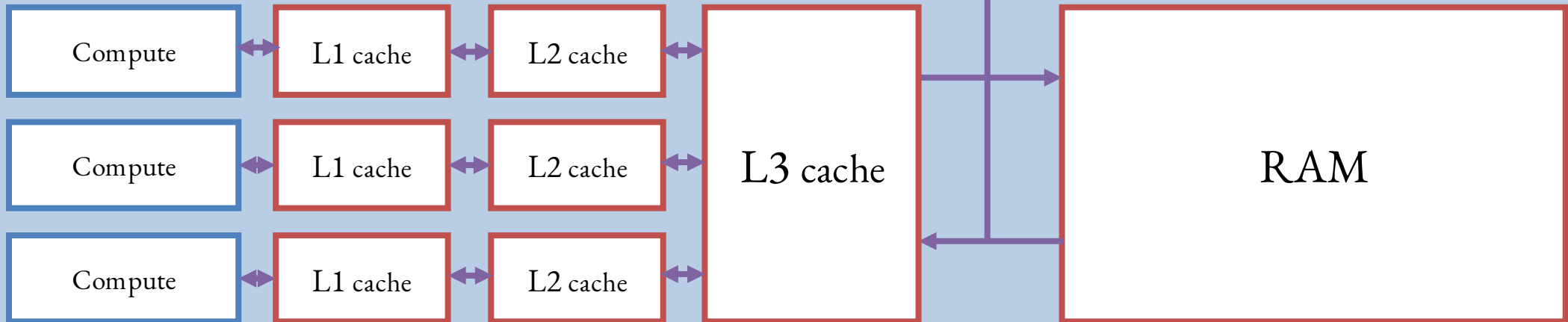
Memory: The Multicore Picture



Socket 1



Socket 2



What can we learn from these pictures?

- Many more **memory** boxes than **compute** boxes
 - And even more as we zoom out
- Memory has a **hierarchical structure**
- **Locality matters**
 - Some memory is closer and easier to access than others
 - Also have standard concerns for CPU cache locality

What limits us?

- **Memory capacity**

- How much data can we store locally in RAM and/or in cache?

- **Memory bandwidth**

- How much data can we load from some source in a fixed amount of time?

- **Memory locality**

- Roughly, how often is the data that we need stored nearby?

- **Power**

- How much energy is required to operate all of this memory?

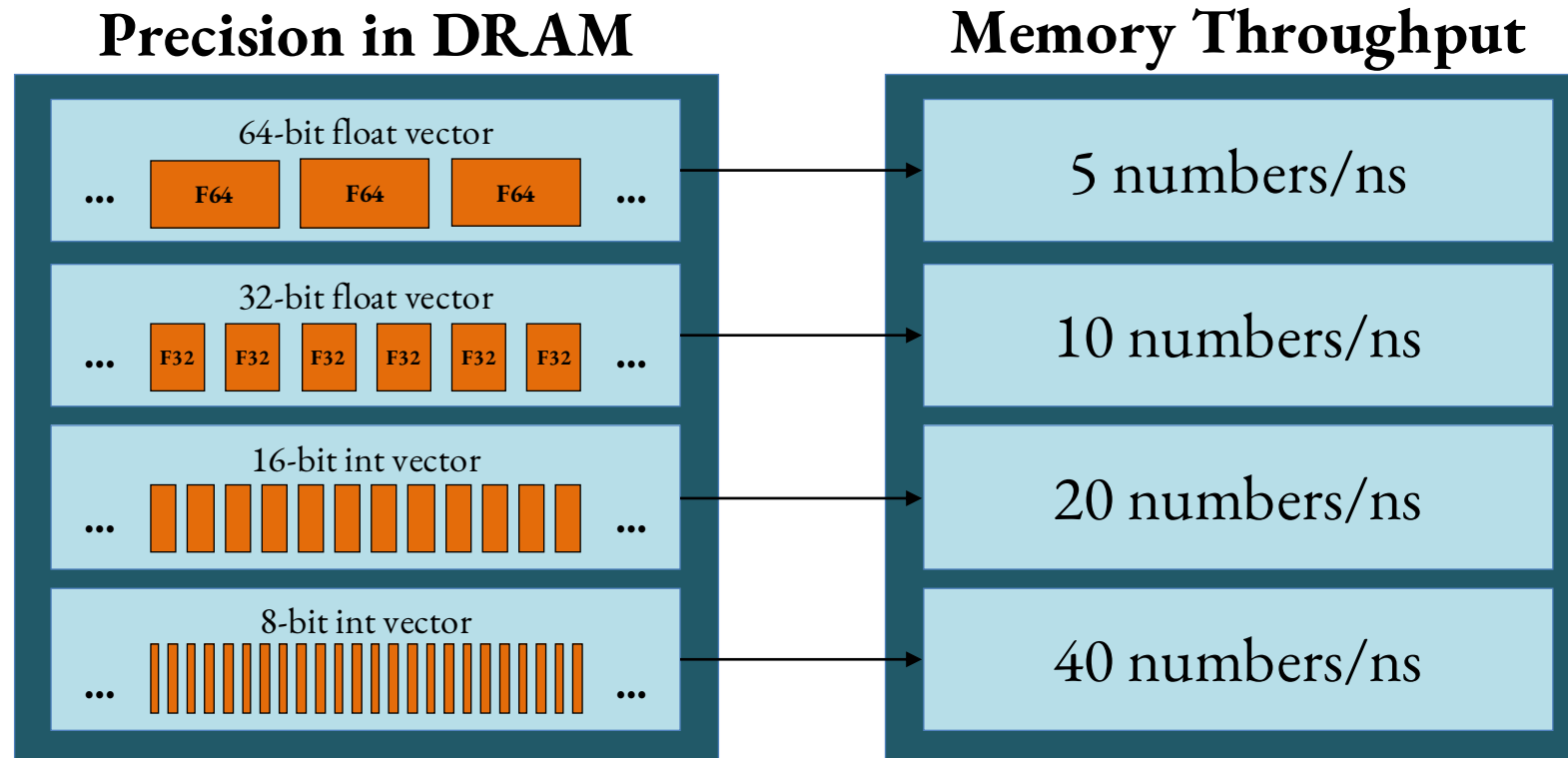
One way to help:
Low-Precision Arithmetic

Low-Precision Arithmetic

- Traditional ML systems use 32-bit or 64-bit **floating point numbers**
- **But do we actually need this much precision?**
 - Especially when we have inputs that come from noisy measurements
- Idea: instead use **8-bit or 16-bit numbers** to compute
 - Can be either floating point or fixed point
 - On an FPGA or ASIC can use arbitrary bit-widths

Low Precision and Memory

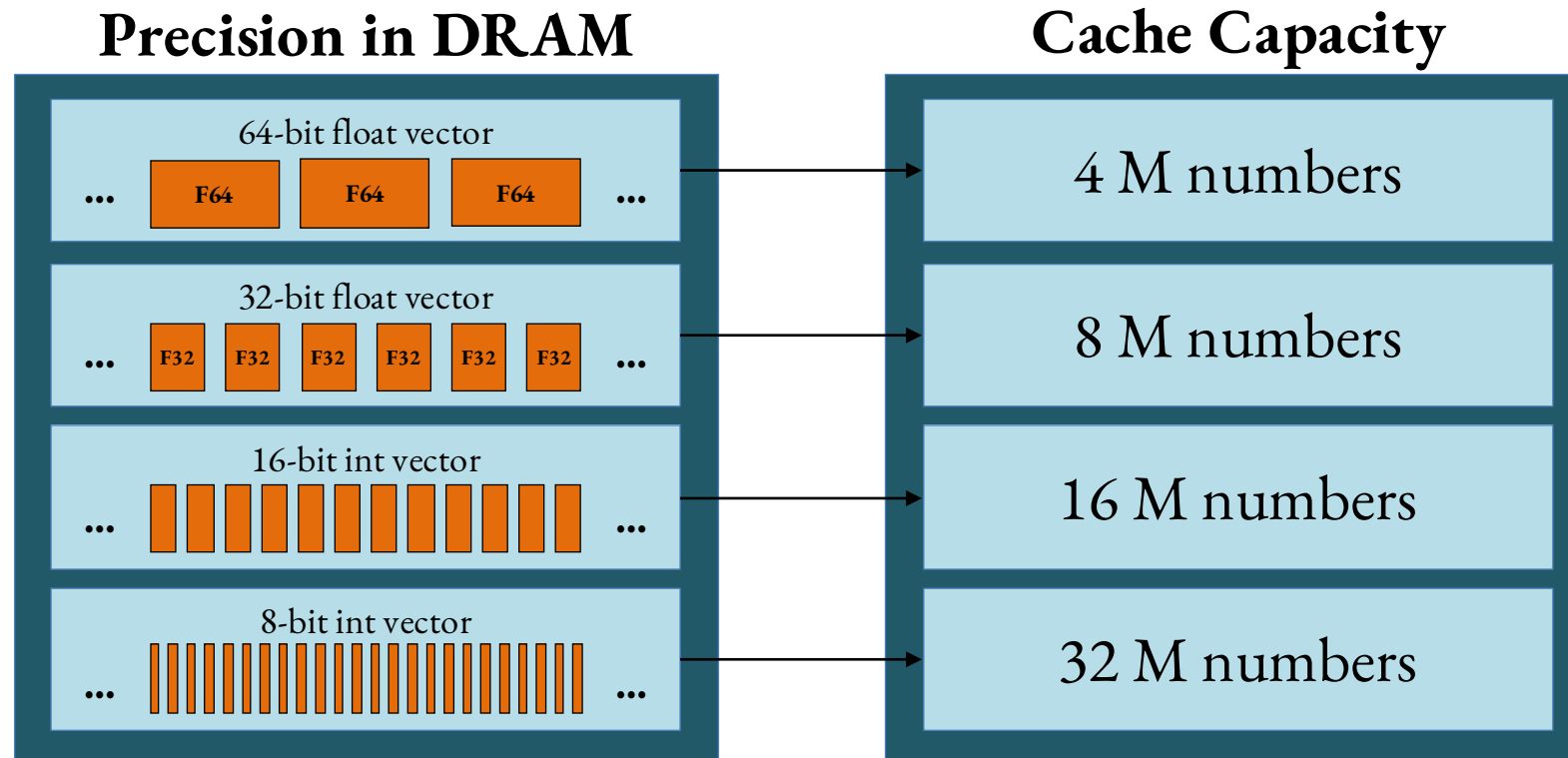
- Major benefit of low-precision: **uses less memory bandwidth**



(assuming ~40 GB/sec memory bandwidth)

Low Precision and Memory

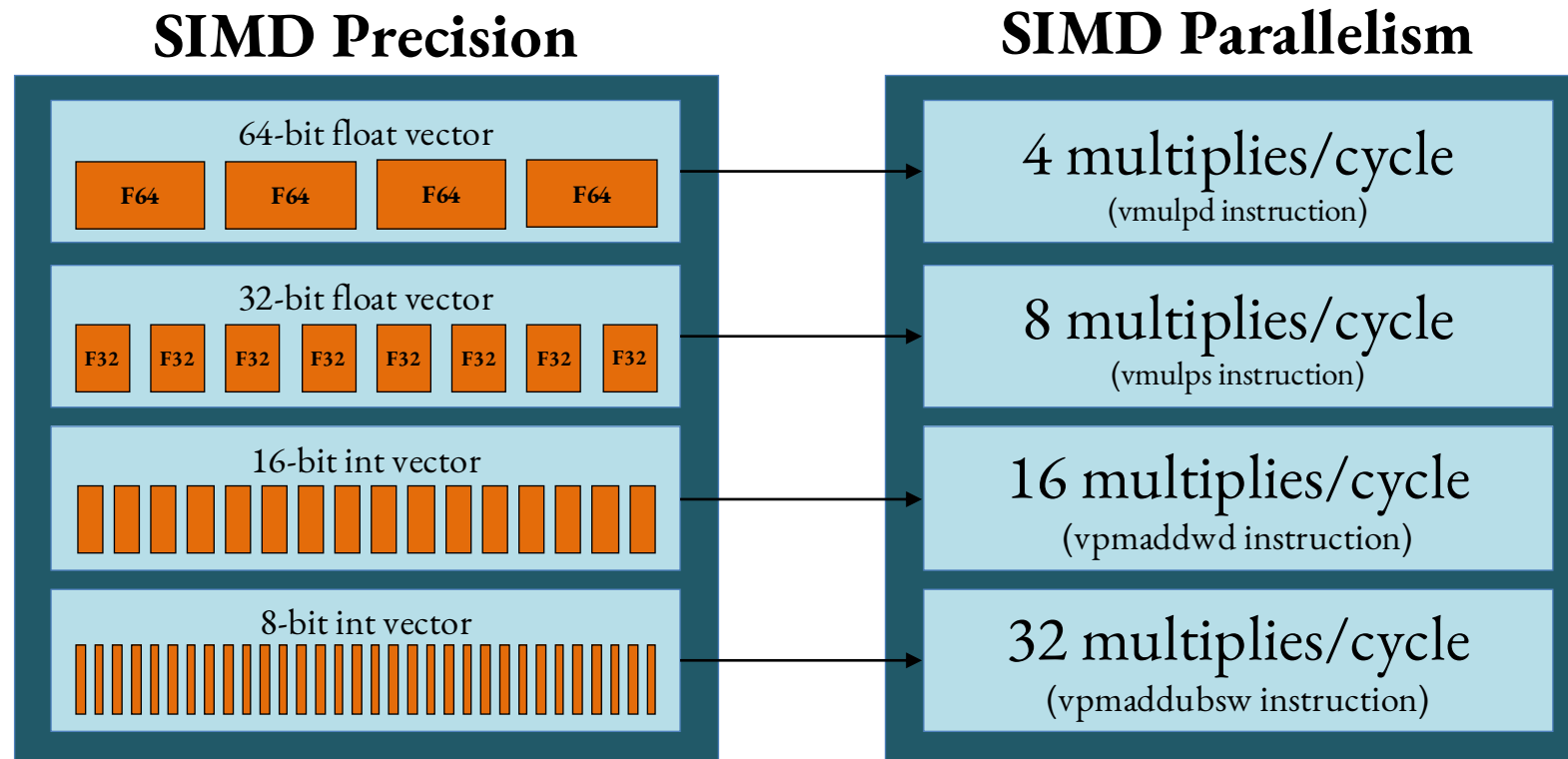
- Major benefit of low-precision: **takes up less space**



(assuming ~32 MB cache)

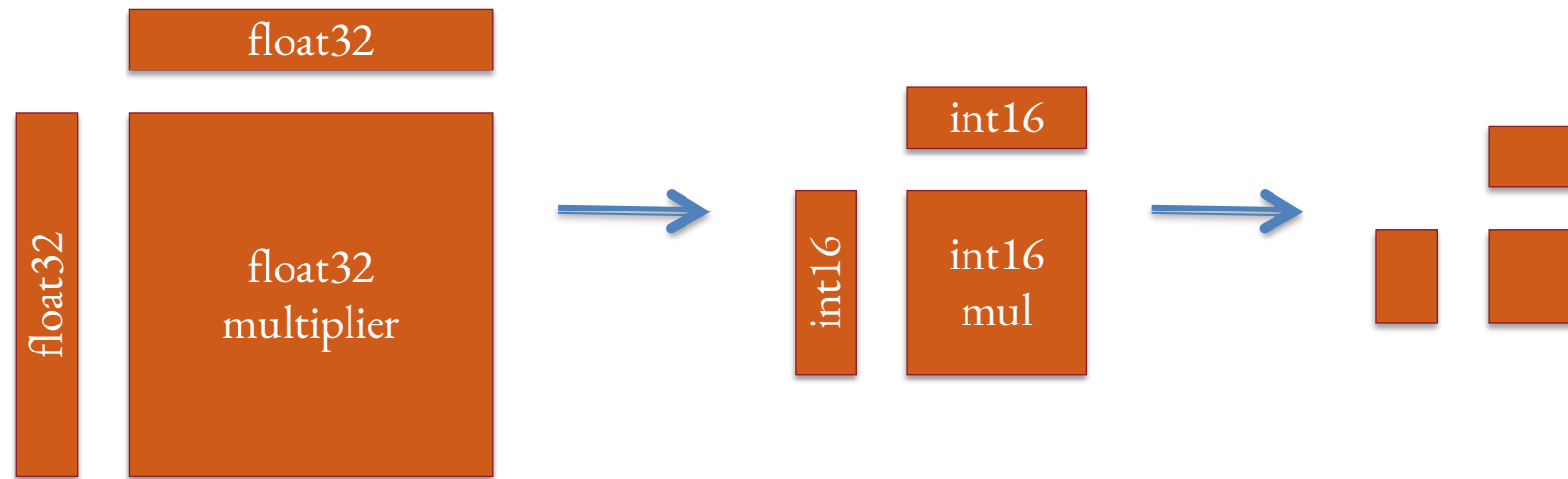
Low Precision and Parallelism

- Another benefit of low-precision: use **SIMD instructions** to get more parallelism on CPU

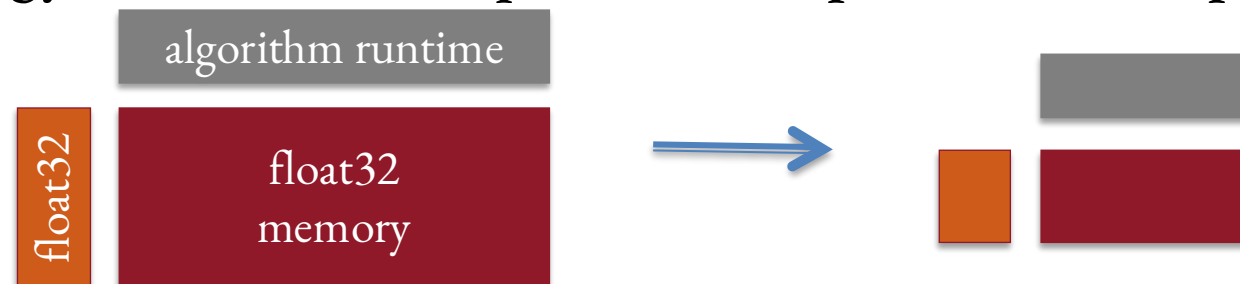


Low Precision and Power

- Low-precision computation can even have a super-linear effect on energy



- Memory energy can also have quadratic dependence on precision



Effects of Low-Precision Computation

- **Pros**

- Fit more numbers (and therefore more training examples) in memory
- Store more numbers (and therefore larger models) in the cache
- Transmit more numbers per second
- Compute faster by extracting more parallelism
- Use less energy

- **Cons**

- Limits the numbers we can represent
- Introduces **quantization error** when we store a full-precision number in a low-precision representation

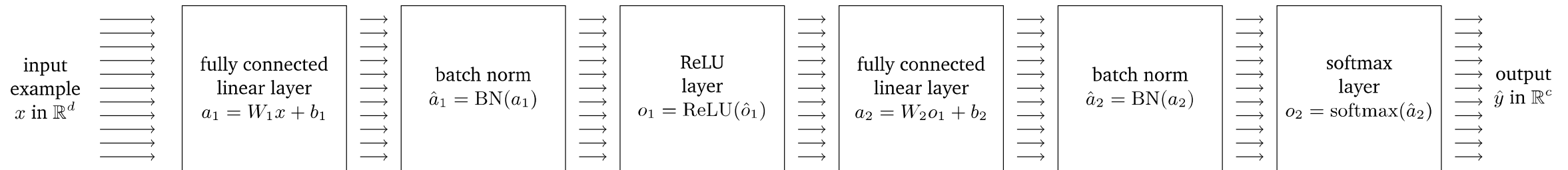
Numeric Formats in Machine Learning

How do we represent numbers as bit patterns on a computer?

A representative setup: DNN training

Many of the large-scale learning tasks we want to accelerate are deep learning tasks.

A **deep neural network (DNN)** looks like this:



Many layers connected to each other in series.

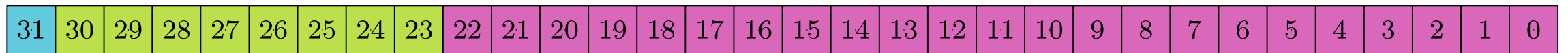
To train, we compute the loss gradient and run stochastic gradient descent:

$$w_{t+1} = w_t - \alpha_t \nabla f(w_t; x_t)$$

The standard approach

Single-precision floating point (FP32)

- 32-bit floating point numbers



sign

8-bit exponent

23-bit mantissa

- Usually, the represented value is

$$\text{represented number} = (-1)^{\text{sign}} \cdot 2^{\text{exponent}-127} \cdot 1.b_{22}b_{21}b_{20} \dots b_0$$

- Has a machine epsilon (measures relative error) of $\epsilon_{\text{machine}} \approx 6.0 \times 10^{-8}$

An example

- Let's convert the number **-6.5** to floating point.

$$6.5 = 13 \times 2^{-1} = (8 + 4 + 1) \times 2^{-1}$$

$$= 1101_b \times 2^{-1} = 1.101_b \times 2^2$$

$$= 1.101_b \times 2^{(129-127)}$$

$$= 1.101_b \times 2^{(10000001_b-127)}$$

1 10000001 10100000000000000000000000000000

What is the machine epsilon?

Or, confusingly,
twice this.

- Represents the relative error of the floating-point format
 - One half the distance between 1 and the next-largest floating point number
 - If there are m mantissa bits, $\varepsilon_{\text{machine}} \approx 2^{-m-1}$
 - Because the smallest representable number > 1 is $1 + 2^{-m}$

Relative error bound. If $x \in \mathbb{R}$ is any number in range of the format, and \hat{x} is the nearest number representable in the format, then

$$|\hat{x} - x| \leq \varepsilon_{\text{machine}} \cdot |x|.$$

Similarly, if $x, y \in \mathbb{R}$ are two floating-point numbers, \star is any primitive numerical operation (e.g. $+$, \times , etc.), and \otimes is the floating-point “version” of that op, then

$$|(x \otimes y) - (x \star y)| \leq \varepsilon_{\text{machine}} \cdot |x \star y|.$$

A low-precision alternative FP16/Half-precision floating point

- 16-bit floating point numbers



1-bit
sign

5-bit
exponent

10-bit
significand

- Usually, the represented value is

$$x = (-1)^{\text{sign bit}} \cdot 2^{\text{exponent} - 15} \cdot 1.\text{significand}_2$$

Numeric properties of 16-bit floats

- A larger machine epsilon (**larger rounding errors**) of $\varepsilon_{\text{machine}} = 4.9 \times 10^{-4}$
 - Compare 32-bit floats which had $\varepsilon_{\text{machine}} \approx 6.0 \times 10^{-8}$
- A smaller overflow threshold (**easier to overflow**) at about 6.5×10^4
 - Compare 32-bit floats where it's 3.4×10^{38}
- A larger underflow threshold (**easier to underflow**) at about 6.0×10^{-8} .
 - Compare 32-bit floats where it's 1.4×10^{-45}

With all these drawbacks, does anyone use this?

Half-precision floating point support

- Supported on most **modern machine-learning-targeted GPUs**
 - E.g. efficient implementation as far back as NVIDIA Pascal GPUs

Pascal Hardware Numerical Throughput

GPU	DFMA (FP64 TFLOP/s)	FFMA (FP32 TFLOP/s)	HFMA2 (FP16 TFLOP/s)	DP4A (INT8 TIOP/s)	DP2A (INT16/8 TIOP/s)
GP100 (Tesla P100 NVLink)	5.3	10.6	21.2	NA	NA
GP102 (Tesla P40)	0.37	11.8	0.19	43.9	23.5
GP104 (Tesla P4)	0.17	8.9	0.09	21.8	10.9

Table 1: Pascal-based Tesla GPU peak arithmetic throughput for half-, single-, and double-precision fused multiply-add instructions, and for 8- and 16-bit vector dot product instructions. (Boost clock rates are used in calculating peak throughputs. TFLOP/s: Tera Floating-point Operations per Second. TIOP/s: Tera Integer Operations per Second. <https://devblogs.nvidia.com/parallelforall/mixed-precision-programming-cuda-8/>)

- Good empirical results for **deep learning**

Another common option

Bfloat16 — “brain floating point”

- Another 16-bit floating point number



1-bit
sign

8-bit
exponent

7-bit
significand

$$\epsilon_{\text{machine}} = 3.9 \times 10^{-3}$$

- Main benefit: numeric range is **now the same as single-precision float**
 - Since it looks like a truncated 32-bit float
 - This is useful because ML applications are **more tolerant to quantization error** than they are to overflow

A more recent option fp8

- An 8-bit floating point number: e5m2



1-bit
sign

5-bit
exponent

2-bit
significand

- Now supported in TensorCores on NVIDIA GPUs

A more recent option fp8

- An 8-bit floating point number: e4m3



1-bit
sign

4-bit
exponent

3-bit
significand

- Now supported in TensorCores on NVIDIA GPUs

A more recent option

fp4

- A 4-bit floating point number: e2m1



1-bit
sign

2-bit
exponent

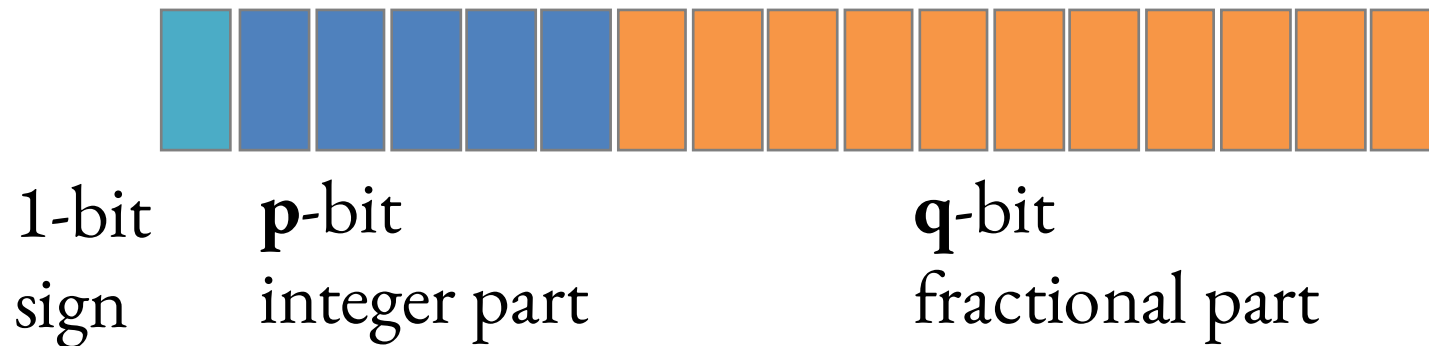
1-bit
significand

- Now supported in TensorCores on NVIDIA GPUs

An alternative to low-precision floating point

Fixed point numbers

- $p + q + 1$ -bit fixed point number



- The represented number is

$$x = (-1)^{\text{sign bit}} (\text{integer part} + 2^{-q} \cdot \text{fractional part})$$
$$= 2^{-q} \cdot \text{whole thing as signed integer}$$

Arithmetic on fixed point numbers

- **Simple and efficient**

- Can just use preexisting integer processing units
- **Lower power** than floating point operations with the same number of bits

- **Mostly exact**

- Can always convert to a higher-precision representation to avoid overflow

- Can represent a **much narrower range of numbers than float**

- **Has an absolute error bound, not relative error bound**

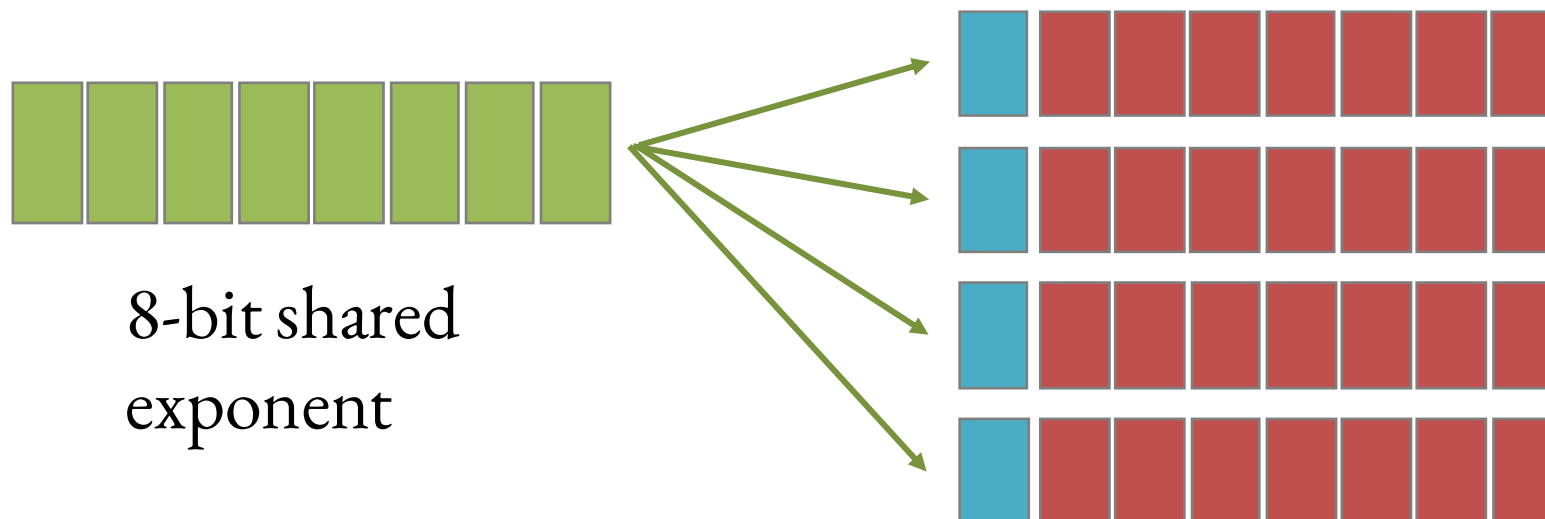
Support for fixed-point arithmetic

- **Anywhere integer arithmetic is supported**
 - CPUs, GPUs
 - Although not all GPUs support 8-bit integer arithmetic
 - And AVX2 does not have all the 8-bit arithmetic instructions we'd like
- Particularly effective on **FPGAs and ASICs**
 - Where floating point units are costly
- Some support for **4-bit int on GPUs**

A powerful hybrid approach

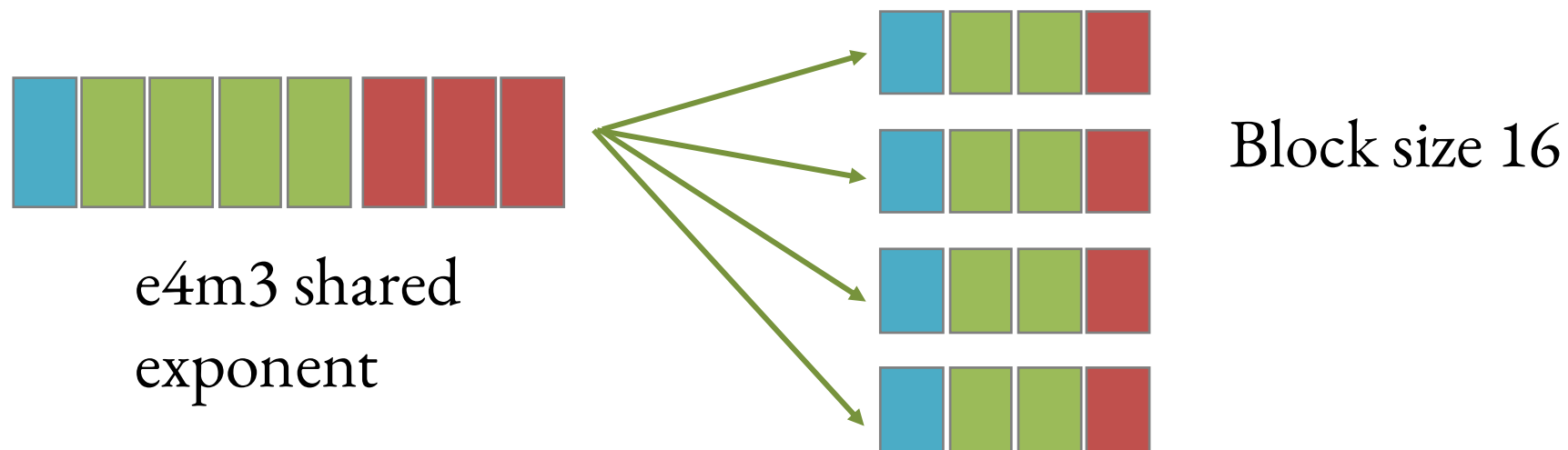
Block Floating Point

- Motivation: when storing a vector of numbers, often these numbers all lie in the same range.
 - So they will have the same or similar exponent, if stored as floating point.
- **Block floating point** shares a single exponent among multiple numbers.



A powerful hybrid approach nvfp4

- Motivation: when storing a vector of numbers, often these numbers all lie in the same range.
 - So they will have the same or similar exponent, if stored as floating point.
- **Block floating point** shares a single exponent among multiple numbers.



A more specialized approach

Custom Quantization Points

- Even more generally, we can just have a list of 2^b numbers and say that these are the numbers a particular low-precision string represents
 - We can think of the bit string as indexing a number in a dictionary
- Gives us total freedom as to range and scaling
 - But **computation can be tricky**
- Some **research into using this with hardware support**
 - *“The ZipML Framework for Training Models with End-to-End Low Precision: The Cans, the Cannots, and a Little Bit of Deep Learning”* (Zhang et al 2017)

Vector quantization

- Group **K** numbers together into a **K**-dimensional vector and quantize this to a **K**-dimensional codebook.
- Ways to make the codebook:
 - Choose via k-Means
 - Choose it to be a lattice with nice packing properties
 - E.g. **E8 lattice**
 - Learn it via SGD
- **What is this most useful for? Weights? Activations?**

Low-precision formats in general

- These are some of the **most common formats used in ML**
 - ...but we're not limited to using only these formats!
- There are **many other things we could try**
 - For example, floating point numbers with different exponent/mantissa sizes
 - Fixed point numbers with nonstandard widths
- Problem: there's **no hardware support** for these other things yet, so it's hard to get a sense of how they would perform.
 - Need to **simulate**

Other Numerical Formats Used Rarely

- **BigFloats**

- Higher-precision floating-point numbers that are implemented in software
- Are sometimes necessary when you need very high precision, such as for very poorly conditioned problems

- Exact arithmetic with **rational numbers**

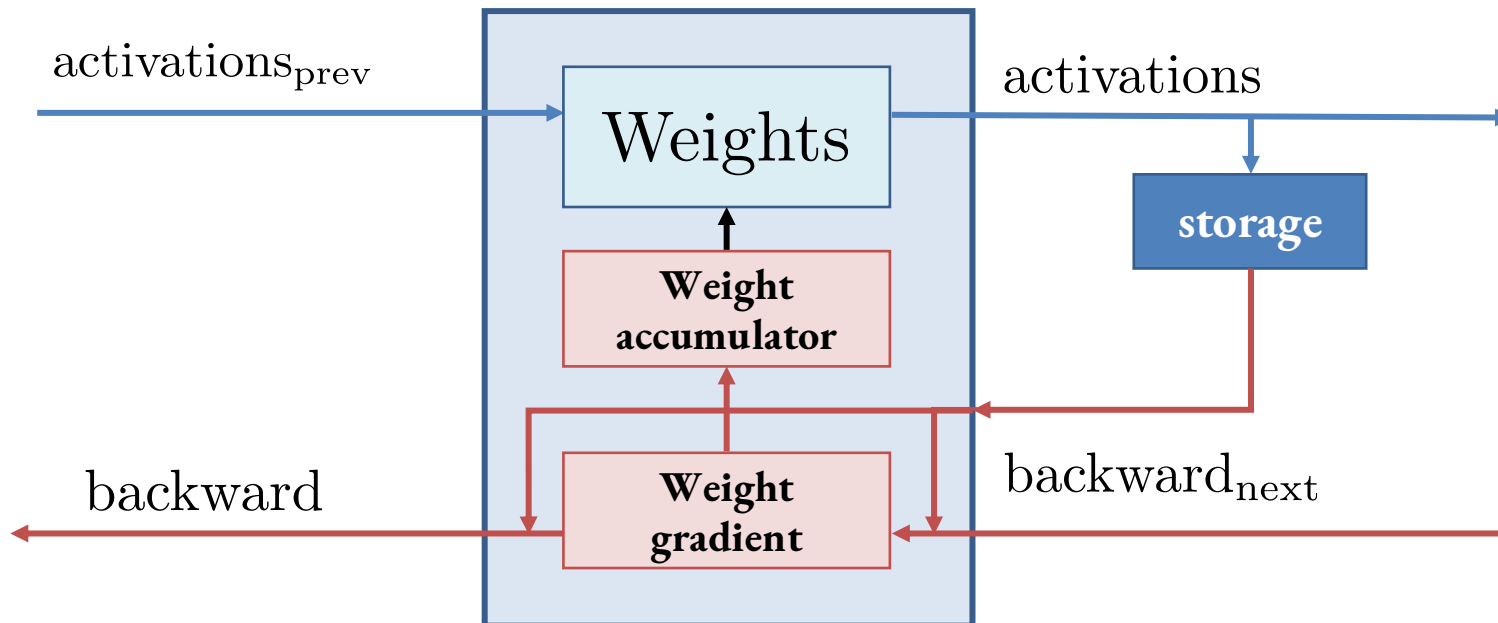
- Lets you do arithmetic with no error
- Numbers have variable length, because they require arbitrarily large integers
- Can also support countable **field extensions of the rational numbers**
- But these are very rarely used because of performance implications

Low-Precision SGD

Using low-precision arithmetic for training

How is precision used for training

- Recall our training diagram
 - Each of these signals forms a class of numbers
 - Generally, we assign a precision to each of the classes, and different classes can have different precisions

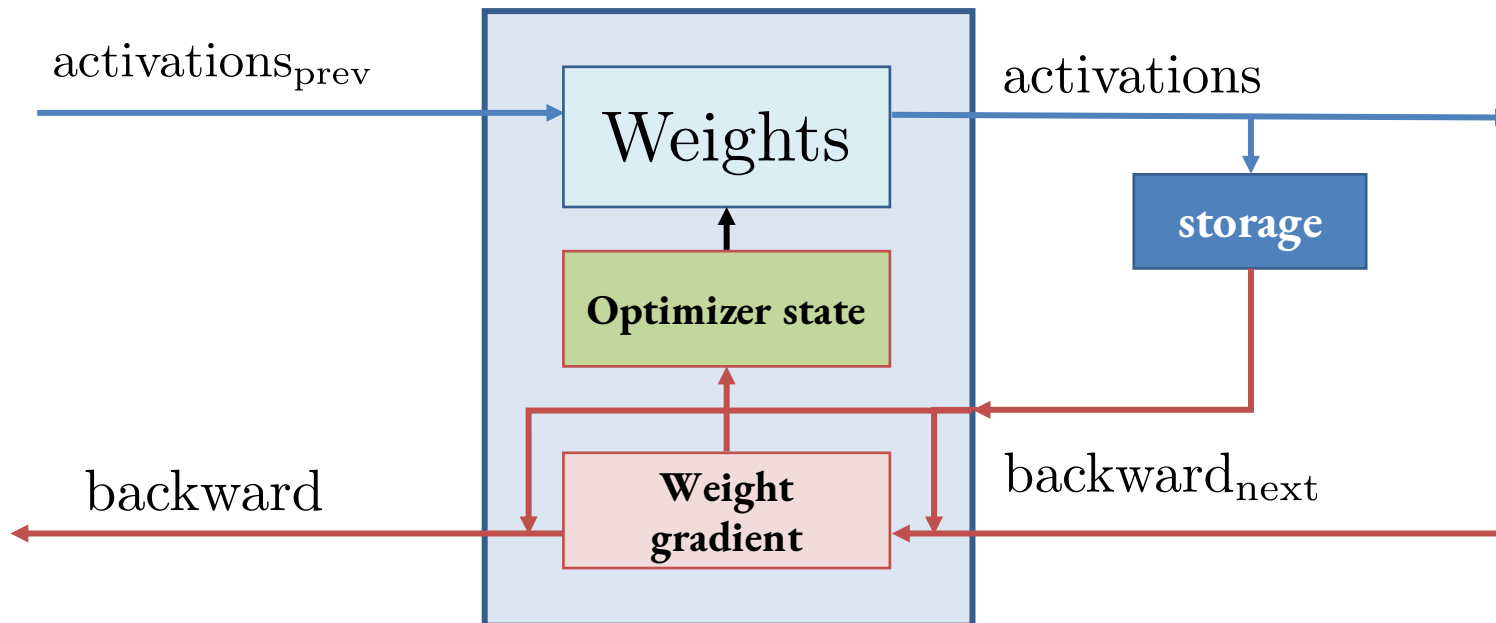


Number classes extended from
“Understanding and Optimizing
Asynchronous Low-Precision Stochastic
Gradient Descent,” ISCA 2017:

- **D**ataset numbers
- **M**odel/weight numbers
- **G**radient numbers
- **C**ommunication numbers
- **A**ctivation numbers
- **B**ackward pass numbers
- **W**eight accumulator
- **L**inear layer accumulator

How is precision used for training

- Recall our training diagram
 - Each of these signals forms a class of numbers
 - Generally, we assign a precision to each of the classes, and different classes can have different precisions



Nobody uses that format 😊

We almost always say something like **W8A16** to describe a format's precision, because (besides the weights and activations) the other numbers are set via a standard recipe.

Optimizer state = float32
Backward = same as activations

Quantize classes independently

- Using low-precision for different number classes has **different effects on throughput**.
 - Quantizing the **dataset numbers** is something we generally don't care about because datasets are small relative to other tensors in the system!
 - Quantizing the **model numbers** improves cache capacity and saves on compute
 - Quantizing the **activation numbers** saves on memory for backprop
 - Quantizing the **communication numbers** saves on expensive inter-worker memory bandwidth

Quantize classes independently

- Using low-precision for different number classes has **different effects on statistical efficiency and accuracy**.
 - Quantizing the **dataset numbers** means you're solving a (slightly) different problem
 - Quantizing the **model numbers** adds noise to each gradient step, and often means you can't exactly represent the solution
 - Quantizing the **gradient numbers** can add errors to each gradient step
 - Quantizing the **activation numbers** changes what the model computes, making it somehow less continuous

Quantization-Aware Training & The Straight-Through Estimator

- A “function” where

$$Q_{\text{straight-thru}}(x) = \text{round}(x)$$

$$Q'_{\text{straight-thru}}(x) = 1$$

- Just round in the forward pass, and pretend the round is not there in the backward pass

The basic recipe for training in low-precision: **Mixed Precision Training**

- Use fp16 to store model and activations wherever this is possible without significant loss of precision
- Use **loss scaling** to stop small gradient values & backward signals from underflowing
- Keep optimizer state in fp32

Published as a conference paper at ICLR 2018

MIXED PRECISION TRAINING

Sharan Narang*, **Gregory Diamos**, **Erich Elsen**[†]

Baidu Research

{sharan, gdiamos}@baidu.com

Paulius Micikevicius*, **Jonah Alben**, **David Garcia**, **Boris Ginsburg**, **Michael Houston**,
Oleksii Kuchaiev, **Ganesh Venkatesh**, **Hao Wu**

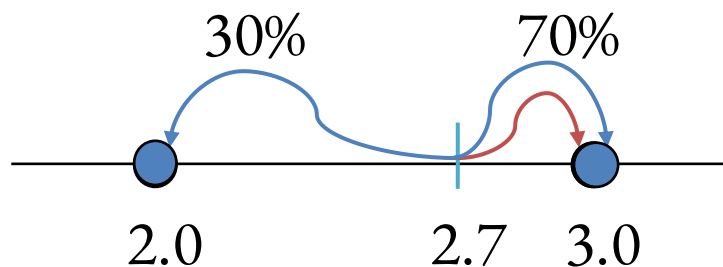
NVIDIA

{paulium, alben, dagarcia, bginsburg, mhouston,
okuchaiev, gavenkatesh, skyw}@nvidia.com

Theoretical Guarantees for Low Precision

- Reducing precision adds noise in the form
- Two approaches to rounding:
 - **nearest rounding** – round to nearest number
 - **stochastic rounding** – round randomly: $E[Q(x)] = x$

Using this, we can prove **guarantees** that SGD converges with a low precision model.



Why stochastic rounding?

- Imagine running SGD with a low-precision model with update rule

$$w_{t+1} = \tilde{Q} (w_t - \alpha_t \nabla f(w_t; x_t, y_t))$$

- Here, \mathbf{Q} is an unbiased quantization function
- In expectation, this is **just gradient descent**

$$\begin{aligned} \mathbf{E}[w_{t+1} | w_t] &= \mathbf{E} \left[\tilde{Q} (w_t - \alpha_t \nabla f(w_t; x_t, y_t)) \middle| w_t \right] \\ &= \mathbf{E} [w_t - \alpha_t \nabla f(w_t; x_t, y_t) | w_t] \\ &= w_t - \alpha_t \nabla f(w_t) \end{aligned}$$

Implementing stochastic rounding

- To implement an unbiased to-integer quantizer:

sample $u \sim \text{Unif}[0, 1]$, then set $Q(x) = \lfloor x + u \rfloor$

- Why is this unbiased?

$$\begin{aligned}\mathbf{E}[Q(x)] &= \lfloor x \rfloor \cdot \mathbf{P}(Q(x) = \lfloor x \rfloor) + (\lfloor x \rfloor + 1) \cdot \mathbf{P}(Q(x) = \lfloor x \rfloor + 1) \\ &= \lfloor x \rfloor + \mathbf{P}(Q(x) = \lfloor x \rfloor + 1) = \lfloor x \rfloor + \mathbf{P}(\lfloor x + u \rfloor = \lfloor x \rfloor + 1) \\ &= \lfloor x \rfloor + \mathbf{P}(x + u \geq \lfloor x \rfloor + 1) = \lfloor x \rfloor + \mathbf{P}(u \geq \lfloor x \rfloor + 1 - x) \\ &= \lfloor x \rfloor + 1 + (\lfloor x \rfloor + 1 - x) = x.\end{aligned}$$

Doing stochastic rounding efficiently

- We still need an efficient way to do unbiased rounding
- **Pseudorandom number generation can be expensive**
 - E.G. doing C++ rand or using Mersenne twister takes many clock cycles
- Empirically, we can use **very cheap** pseudorandom number generators
 - And still get good statistical results
 - For example, we can use XORSHIFT which is just a cyclic permutation

Limitations of stochastic rounding

- Technique only makes sense when we're summing up a bunch of independently rounded values
- Works best for the accumulators in the optimizer!
- But in the Mixed Precision recipe, we store those accumulators in full-precision anyway
 - ...so there's not much point in the stochastic rounding
- Also it introduces **a lot of noise** to the training process.



Benefits of low-precision On a real device...

NVIDIA H100 Tensor Core GPU

Technical Specifications	
	H100 SXM
FP64	34 teraFLOPS
FP64 Tensor Core	67 teraFLOPS
FP32	67 teraFLOPS
TF32 Tensor Core	989 teraFLOPS ²
BFLOAT16 Tensor Core	1,979 teraFLOPS ²
FP16 Tensor Core	1,979 teraFLOPS ²
FP8 Tensor Core	3,958 teraFLOPS ²
INT8 Tensor Core	3,958 TOPS ²

<https://resources.nvidia.com/en-us-tensor-core/nvidia-tensor-core-gpu-datasheet>

Drawbacks of low-precision

- The draw back of low-precision arithmetic is the **low precision!**
- Low-precision computation means we accumulate **more rounding error** in our computations
- These rounding errors can add up throughout the learning process, resulting in **less accurate learned systems**
- The trade-off of low-precision: **throughput/memory vs. accuracy**

A modern LLM example

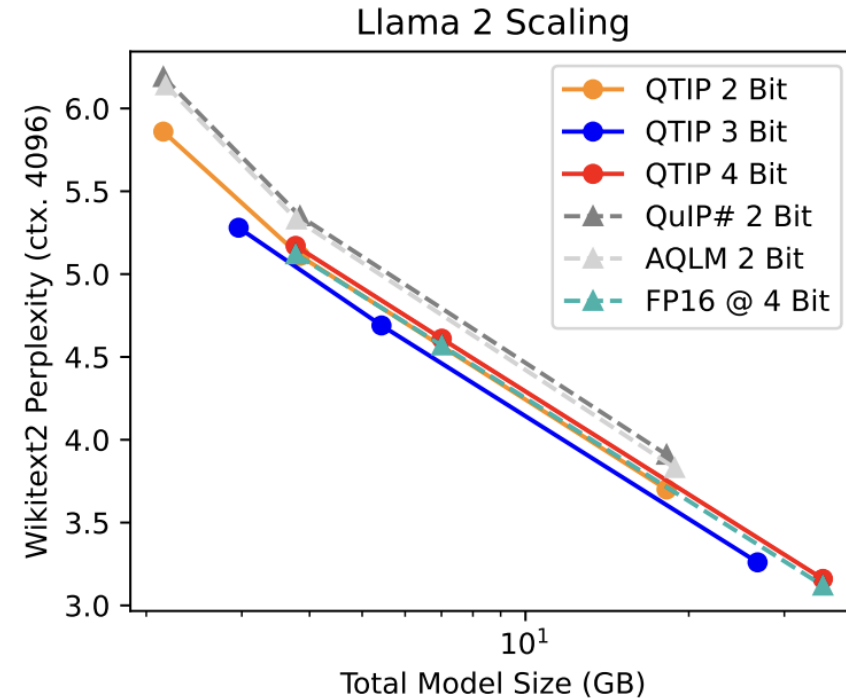
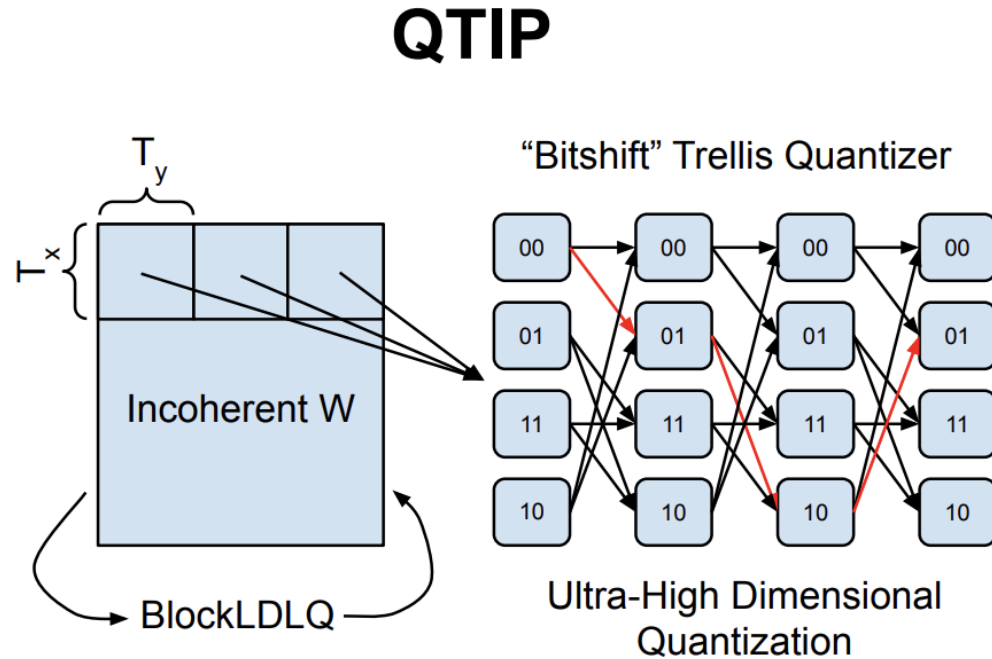


Figure 1: QTIP performs ultra-high dimensional (> 100) quantization by using Trellis Coded Quantization, which has linear cost in dimension. This enables QTIP to outperform Vector Quantization-based approaches (QuIP#, AQLM) that are limited to low dimensions. With QTIP, 2 bit models scale better than theoretically optimal 4 bit models.

A modern LLM example

openai/gpt-oss-120b like 4.59k Follow OpenAI 32.7k

Text Generation Transformers Safetensors gpt_oss vllm conversational Eval Results 8-bit precision mxfp4 arxiv:2508.10925 License: apache-2.0

Model card Files and versions xet Community 198

Deploy Use this model

Edit model card

Downloads last month
4,727,242



Safetensors

Model size 120B params Tensor type BF16 · U8 Chat template Files info

Inference Providers NEW

Groq +6

Text Generation

Examples

Input a message to start chatting with openai/gpt-oss-120b.



[Try gpt-oss](#) · [Guides](#) · [Model card](#) · [OpenAI blog](#)

Highlights

- MXFP4 quantization:** The models were post-trained with MXFP4 quantization of the MoE weights, making `gpt-oss-120b` run on a single 80GB GPU (like NVIDIA H100 or AMD MI300X) and the `gpt-oss-20b` model run within 16GB of memory. All evals were performed with the same MXFP4 quantization.

Next time...

- **Post-training quantization & compression!**
- **How can we leverage low-precision arithmetic to make large models work on small devices?**