

**LLM.int8()
+
OPTQ**

Martin Jaggi

- int8():
Run a transformer with **8bit** ops instead of 32
- some models can now be used on *Google Colab* that previously couldn't

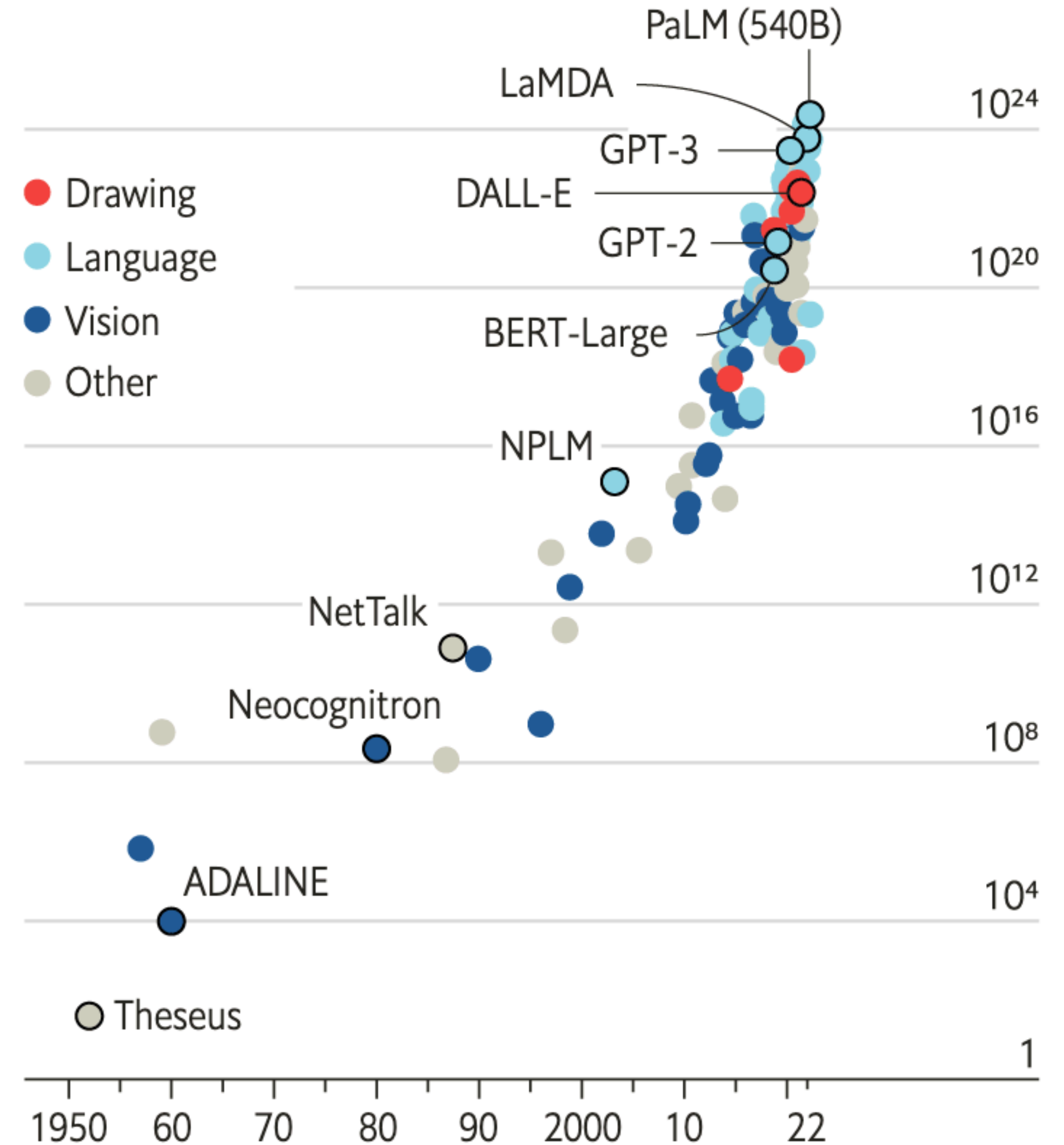
- int8():
Run a transformer with **8bit** ops instead of 16
- reduce the **memory footprint** of a large model by 2x (compared to FP16)
- some models can now be used on *Google Colab* that previously couldn't

1) Problem importance

- LLMs are large
- Inference cost is crucial (memory & compute)

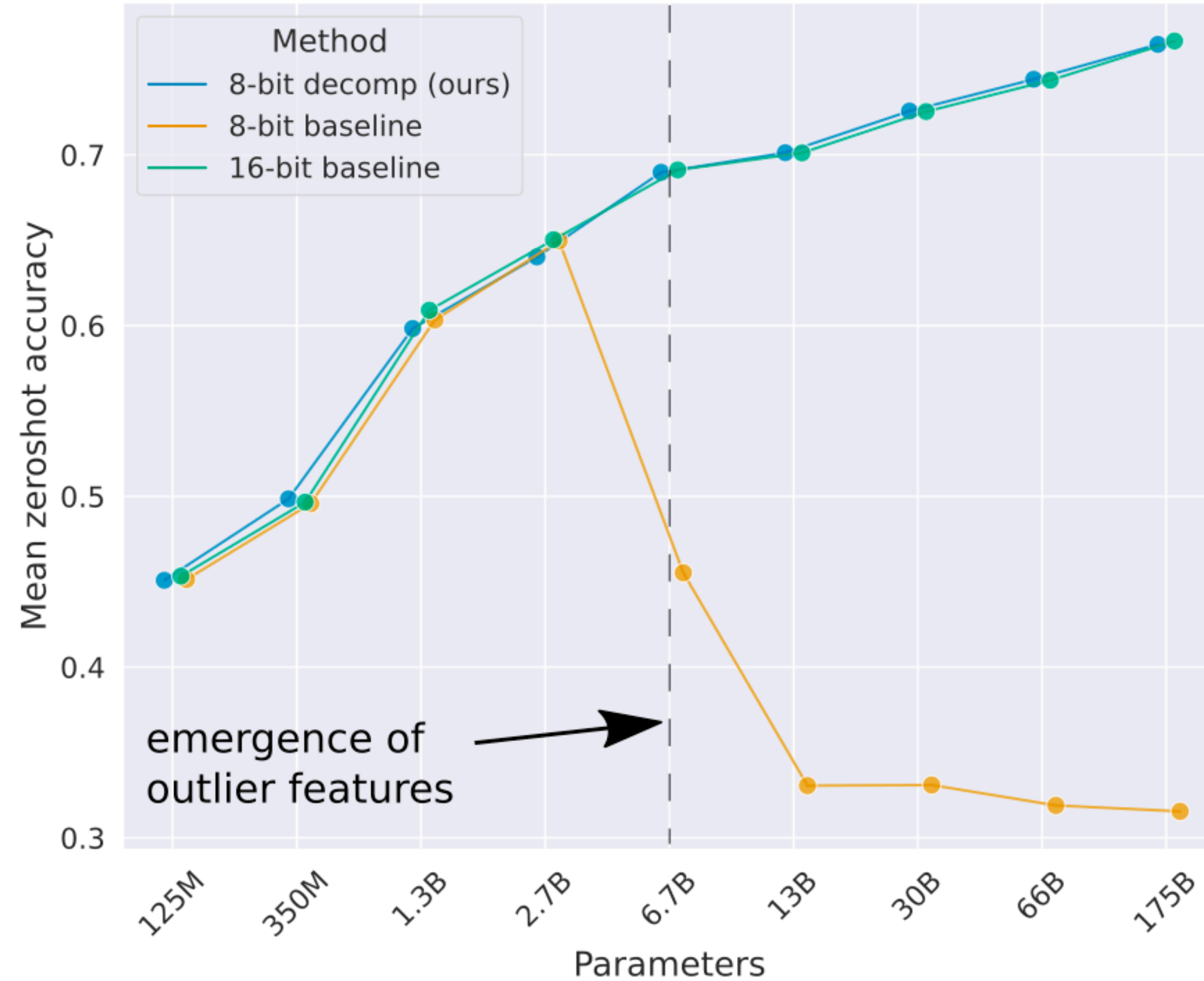
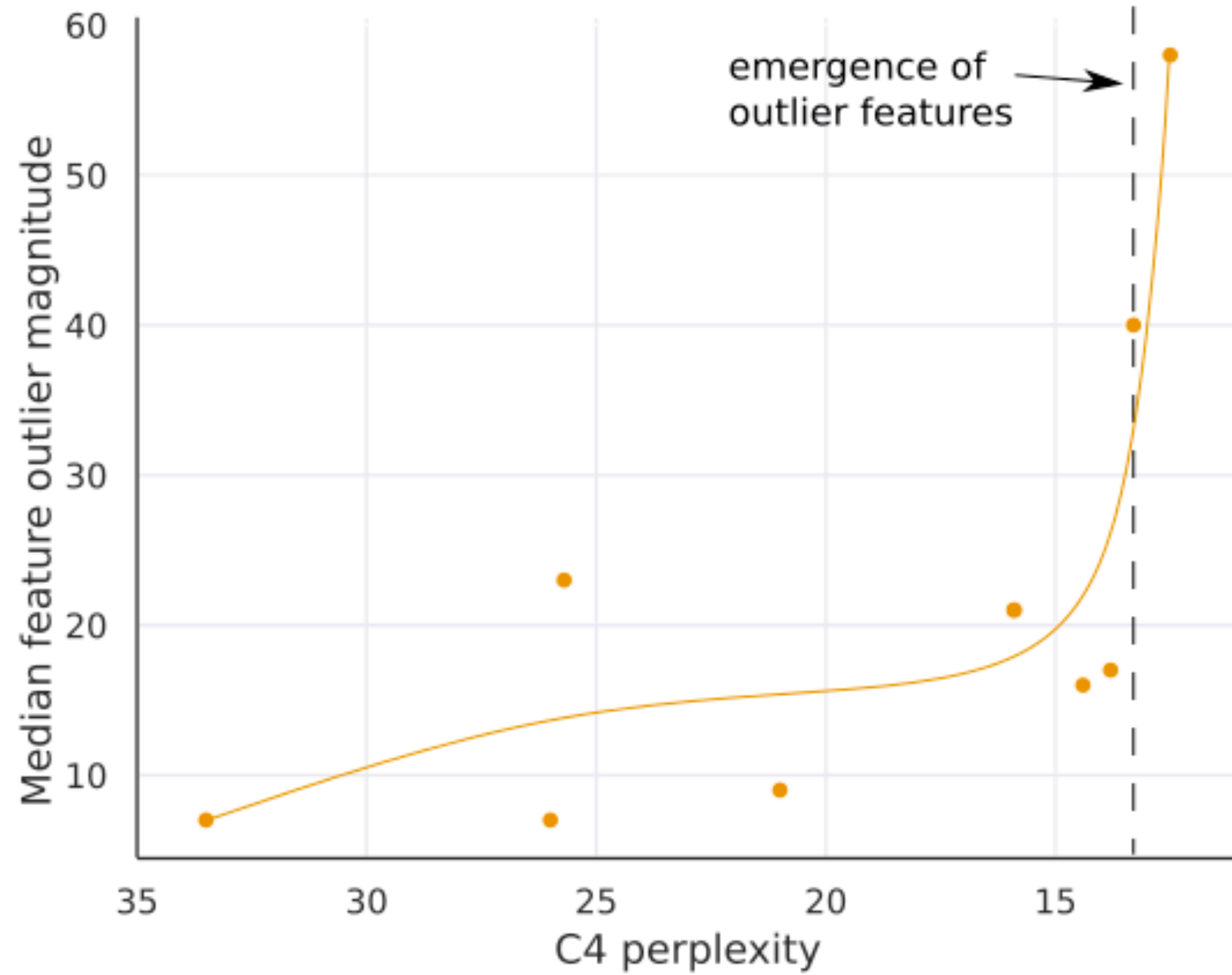
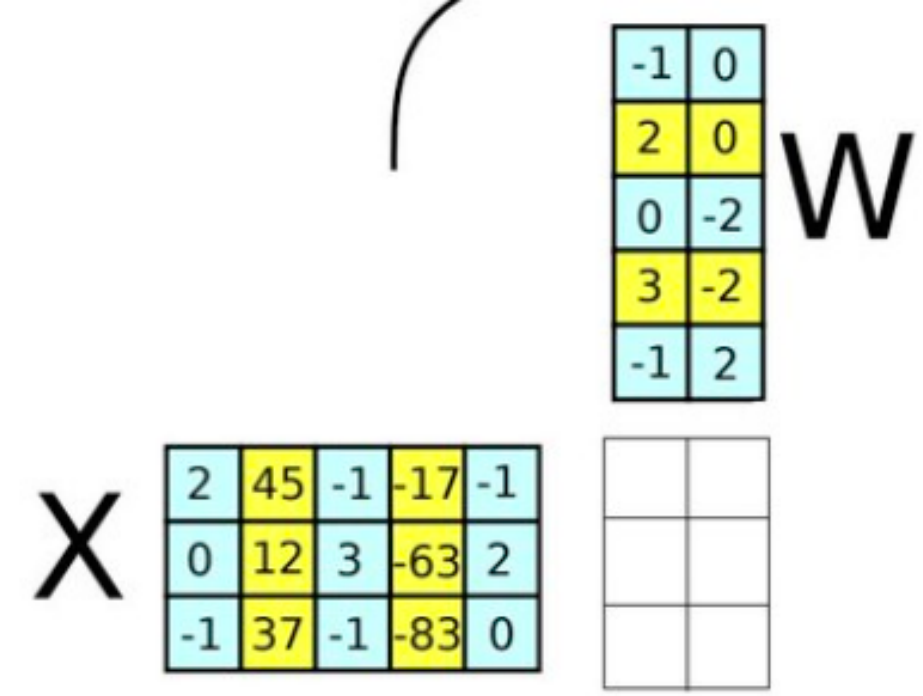
The blessings of scale

AI training runs, estimated computing resources used
Floating-point operations, selected systems, by type, log scale



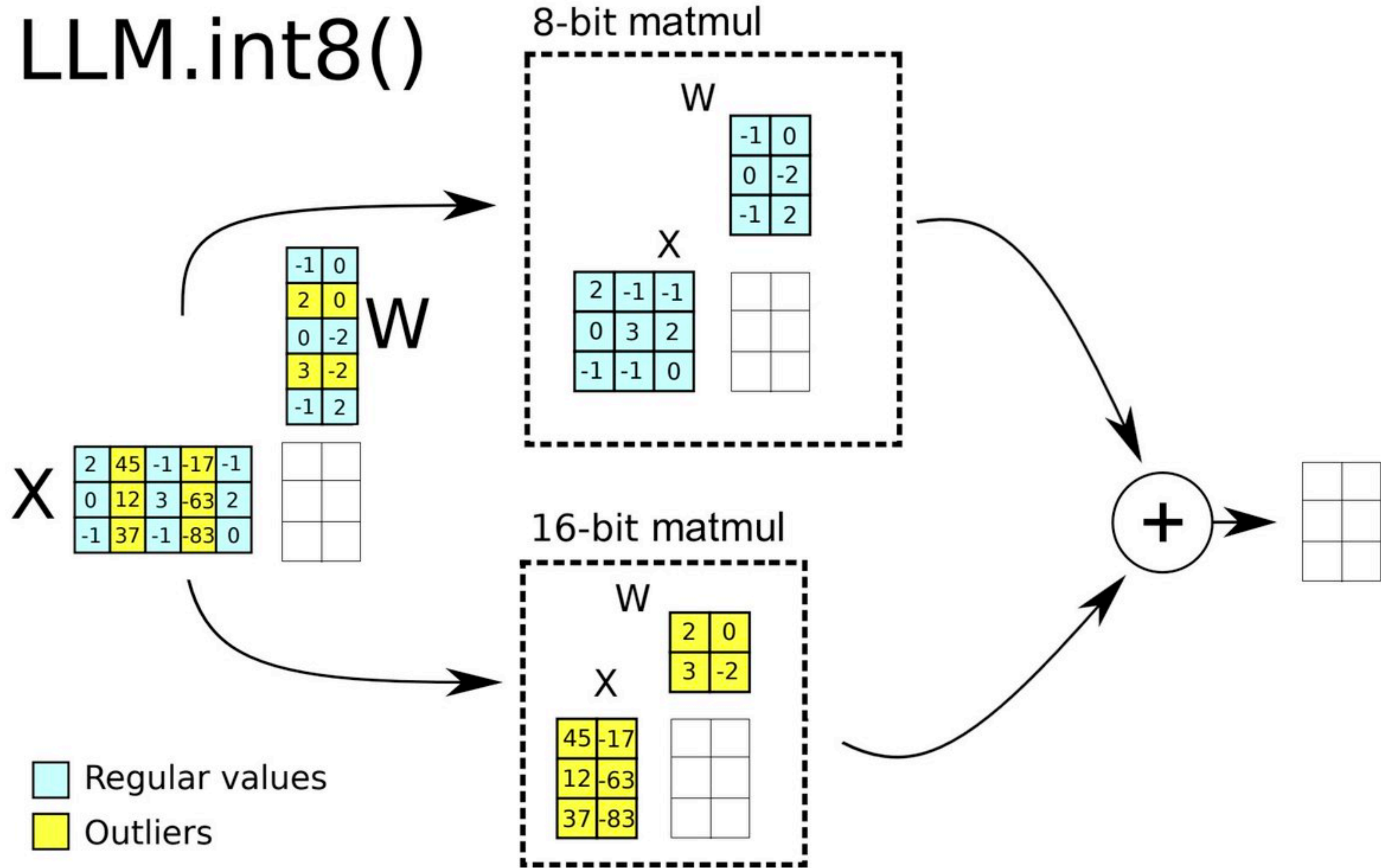
2) Insights

- Emergence of outlier features



3) Solution

- Treat regular & outlier values separately
- Integration into hugging face



EPFL 4) Takeaway

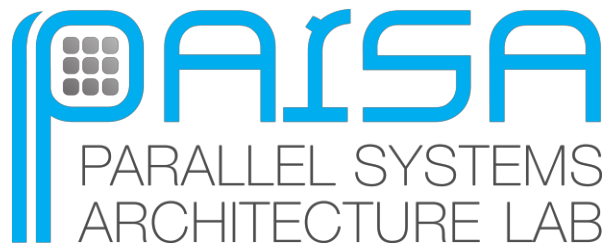
- reduce the **memory footprint** of a large model by *(only)* 2x.
- some models can now be used on *(insert particular hardware)* that previously couldn't
- any huggingface model can directly be converted
- no obvious speedup (only memory saving)

- Two papers addressing same problem:
 - LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale
NeurIPS 2022 (*arXiv Aug 2022*)
 - OPTQ: Accurate Quantization for Generative Pre-trained Transformers
ICLR 2023 (*arXiv Oct 2022*)

FAST: DNN training under variable precision Block Floating Point with Stochastic Rounding

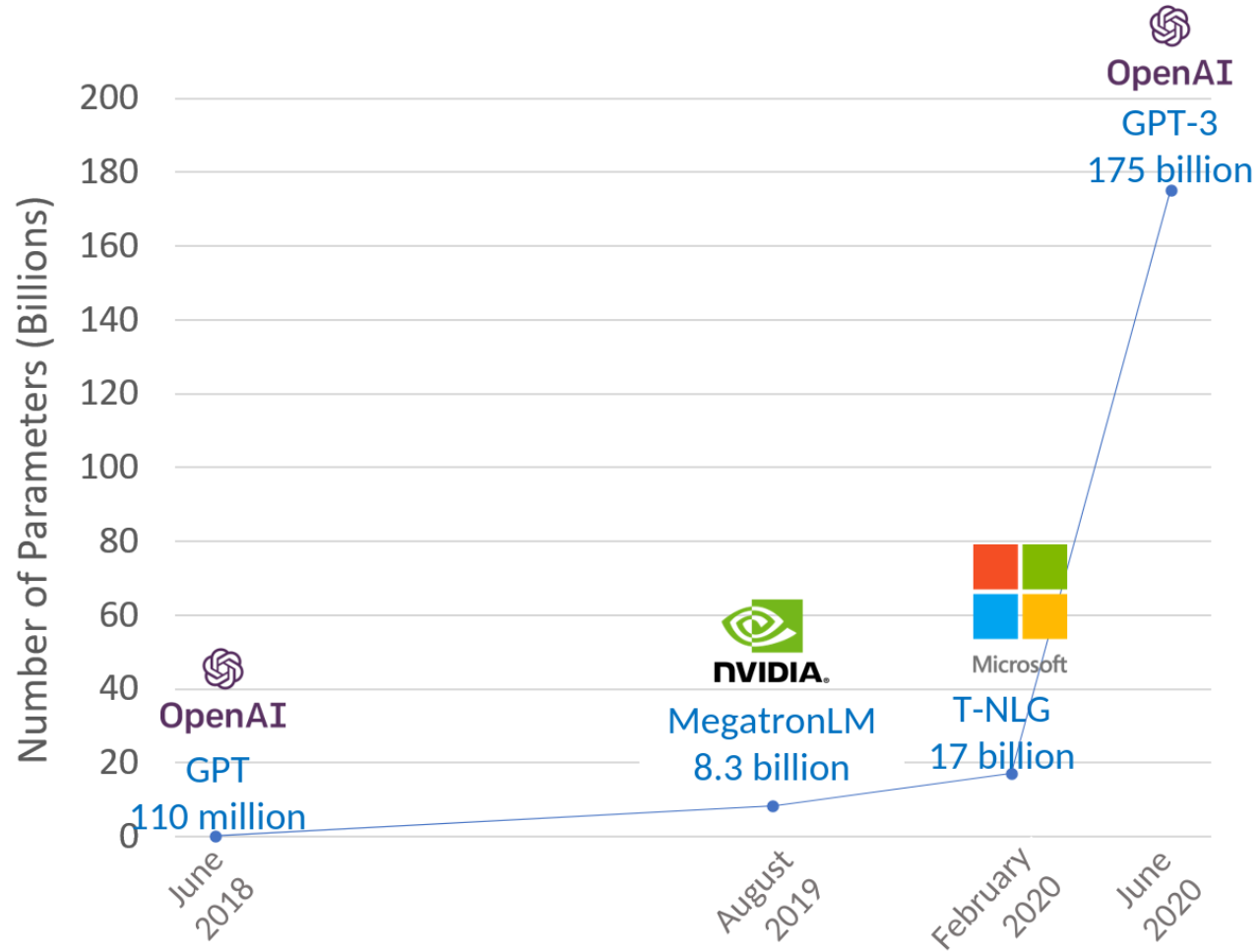
CS-723 Presentation 4th April, 2023

By Ayan Chakraborty



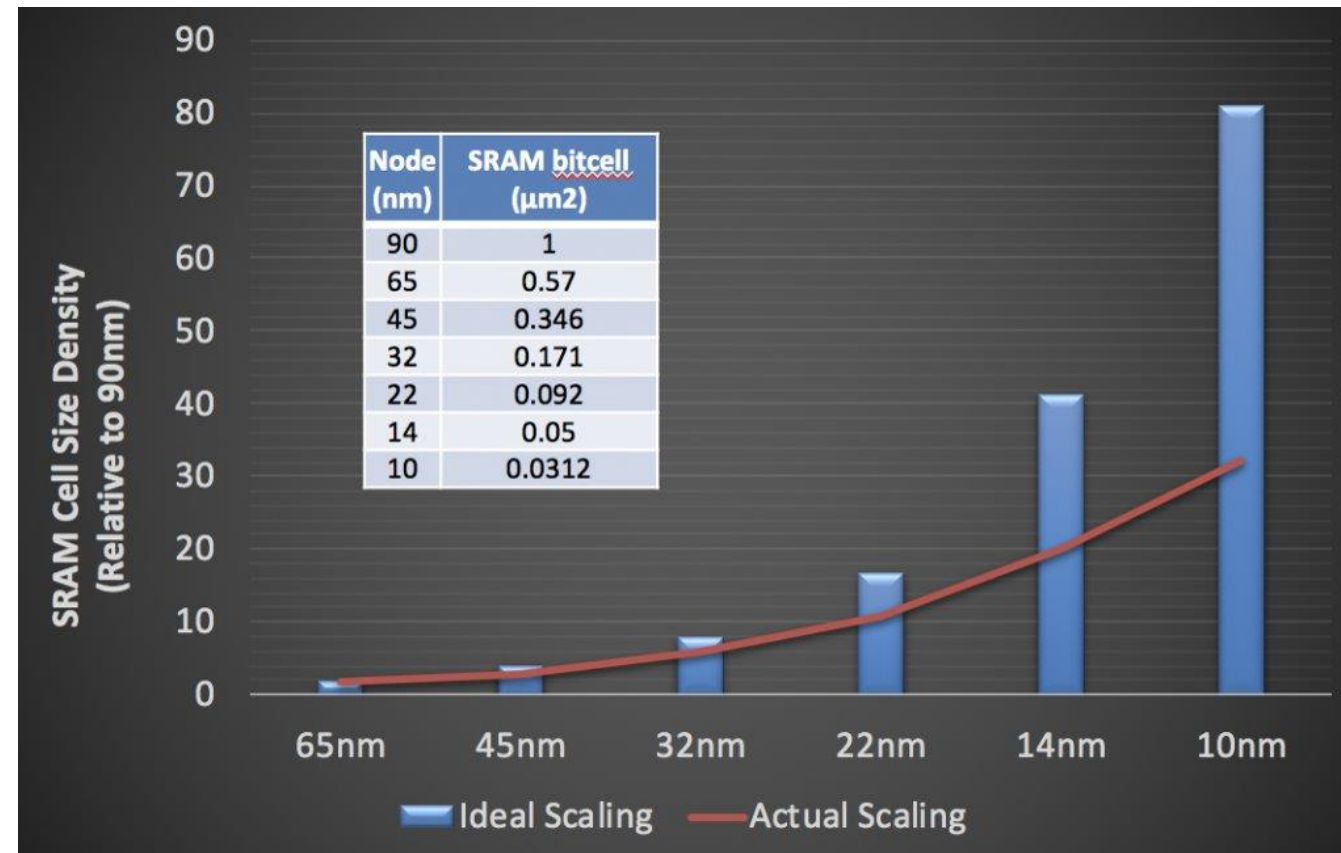
Elevator Pitch

Exponential growth of DNNs



Commensurate increase in computational cost and carbon emissions!

Slowdown of Moore's Law



Need for alternate ways to increase Arithmetic Density

Use new numerical encodings for DNNs

- Block Floating Point promising candidate for DNNs
- Combines the advantages of Fixed point and Floating point
- Previous work does not manage to reduce precision below 8 bits for training

- The paper shows how you can use variable precision to use even lower precision for most operations

- The paper designs a custom HW accelerator to get the maximum benefits

Increase hardware efficiency while maintaining accuracy

What is the Problem?

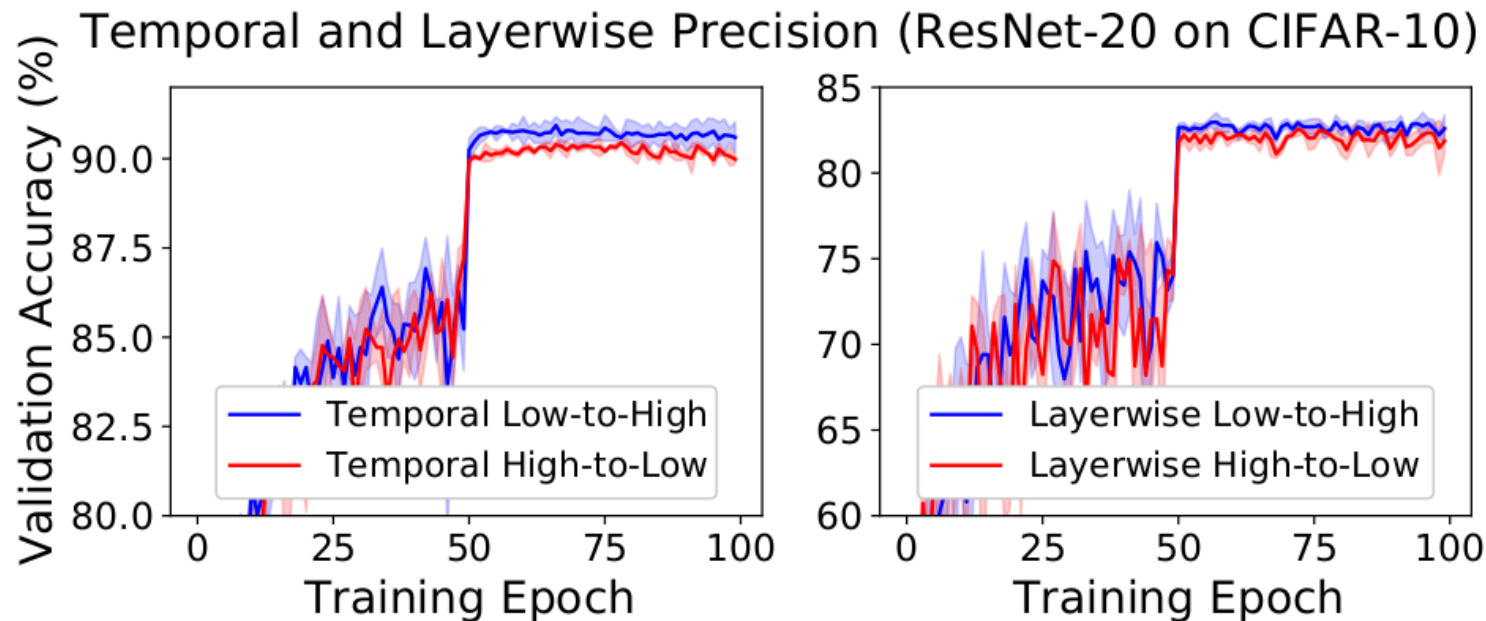
Unexplored design space of BFP

- Prior work does not manage to use BFP with precision below 8 bits to get high accuracy.
- Variable precision and Mixed precision training used for other formats before but not studied for BFP.
- Big unexplored design space of BFP configurations below 8 bits.
- HW support for BFP in general, and more specifically for variable precision BFP not studied well in prior work.

What are the insights?

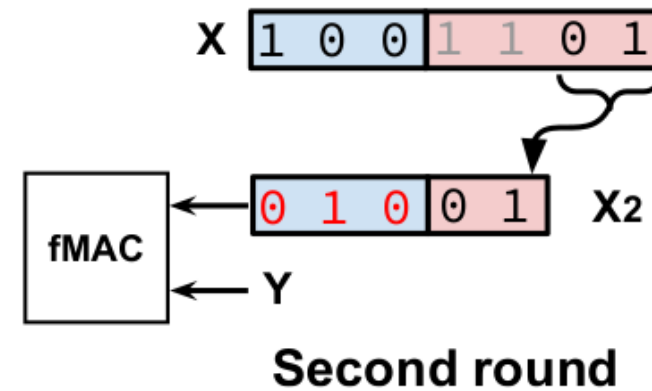
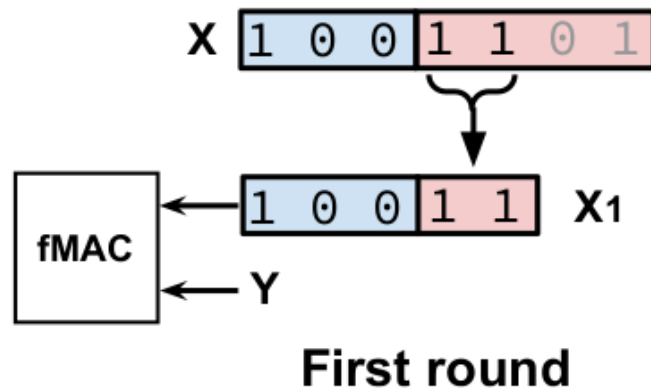
Variable precision

- The precision of the model parameters does not need to stay fixed during training and even across model layers.
- Early layers of a model and the initial epochs of a training session are more tolerant to lower precision.



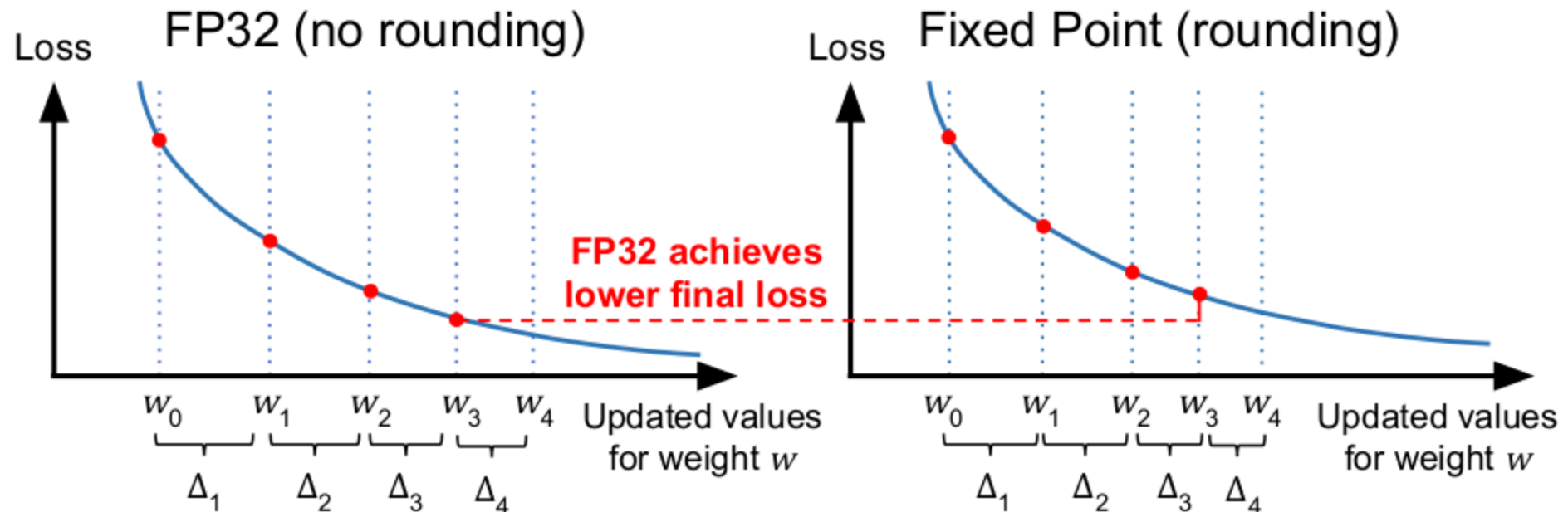
HW support for variable precision

- Sub-divide the computation into 2-bit chunks
- Allows the same custom HW unit to implement arithmetic operations involving higher precision mantissa by simply running multiple passes of the unit.



Stochastic Rounding

- Randomly round up or round down when converting tensors from FP32 to BFP.
- Critical to maintaining training stability when using very low precision.



What is the solution?

Schedule for variable precision

- Use 4-bit mantissas in combination with 2-bit mantissas
- Check if relative improvement of using the higher precision is smaller than a threshold for a given layer and iteration

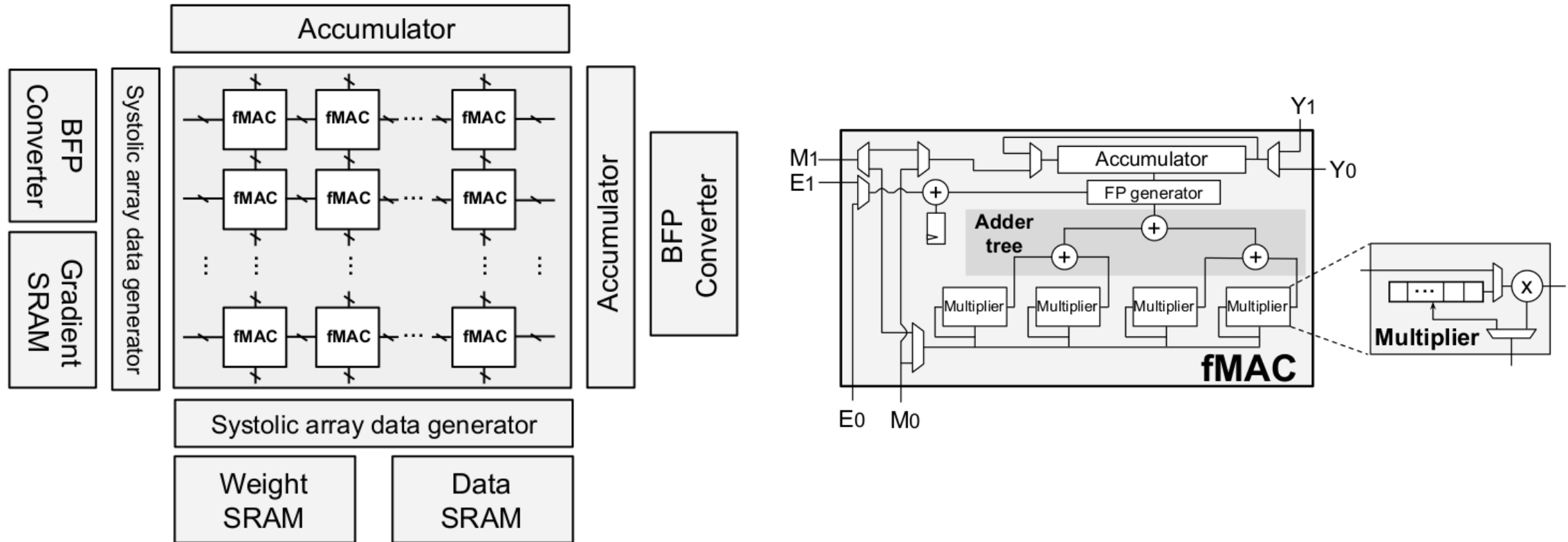
$$r(X) = \frac{\sum_n |BFP(X_n, 4) - BFP(X_n, 2)|}{\sum_n |BFP(X_n, 2)|}$$

- Scale the threshold linearly across the layers and iterations

$$\varepsilon(l, i) = \alpha - \beta \frac{i}{I} - \beta \frac{l}{L}$$

HW accelerator for variable precision BFP

- Systolic array based accelerator with custom MAC unit.



- Novel memory layout to store variable precision BFP values.

What is the take-away message?

Dynamic precision is the way to go

- Dynamic precision enables the usage of BFP precisions below 8 bits to get comparable high performance.
- Dynamic precision allows most operations to be done at a much reduced precision.
- Building a custom HW accelerator allows you to extract the maximum benefits of dynamic precision BFP operations.

Will this paper win the Test of Time award?

No

- ML highly volatile field.
- Unlikely only using BFP by itself will turn out to be a good choice.
- Does not have any new insights which can be applied to any number format.
- Unclear how HW would generalize well for any mantissa width.
- Unlikely results shown in the paper will generalize to current LLM models.

Why should this paper not have appeared at HPCA?

Highly empirical with no new insights

- The idea of mixed precision training is not new, and the paper simply adapts it for BFP without providing any insight as to why it works.
- The design choices in the paper are mostly based on empirical results with no theoretical backing.
- The training accuracies are not too high, and it is hard to believe SOTA performance for current models can be gained using only 4-bits and 2-bits
- Design not open sourced, makes it hard to reproduce results

Thank you