# Schrödinger's FP:
# Dynamic Adaptation of Floating-Point Containers for Deep Learning Training
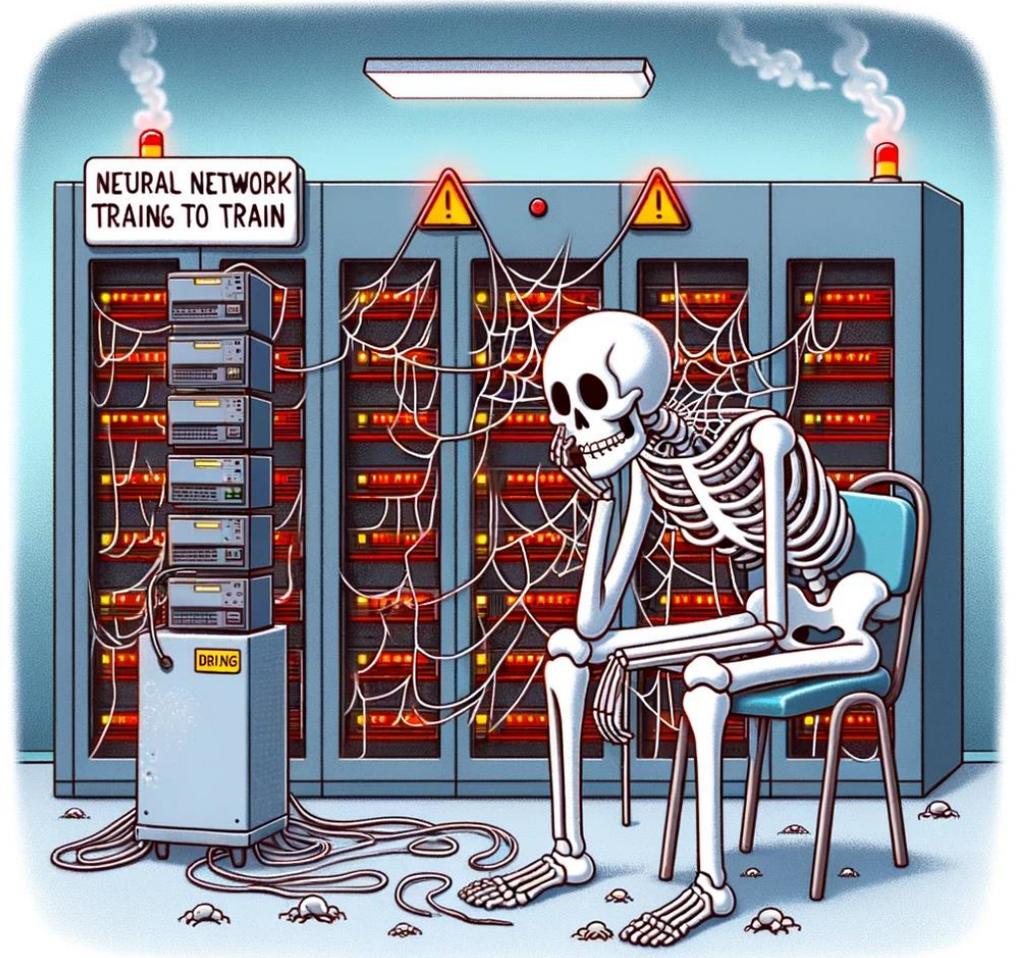
**Miloš Nikolić**, Enrique Torres Sanchez, Jiahui Wang, Ali Hadi Zadeh,
Mostafa Mahmoud, Ameer Abdelhadi, Kareem Ibrahim, and Andreas Moshovos

The Edward S. Rogers Sr. Department
of Electrical & Computer Engineering
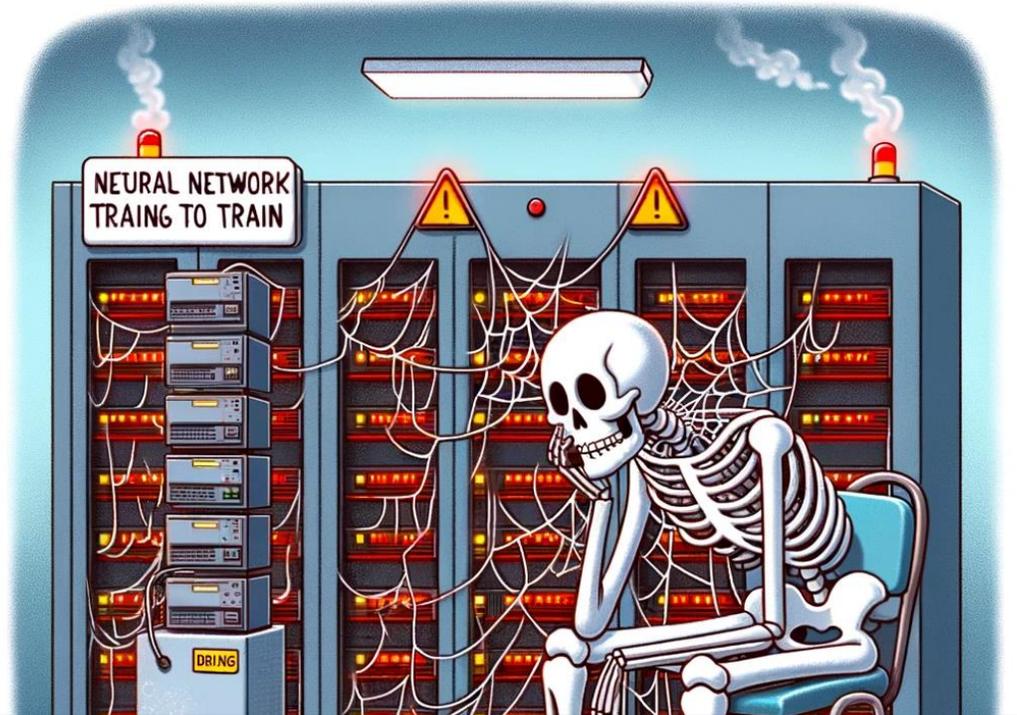UNIVERSITY OF TORONTO

# Training is Expensive!

- Neural Networks are becoming ubiquitous
- But they are expensive
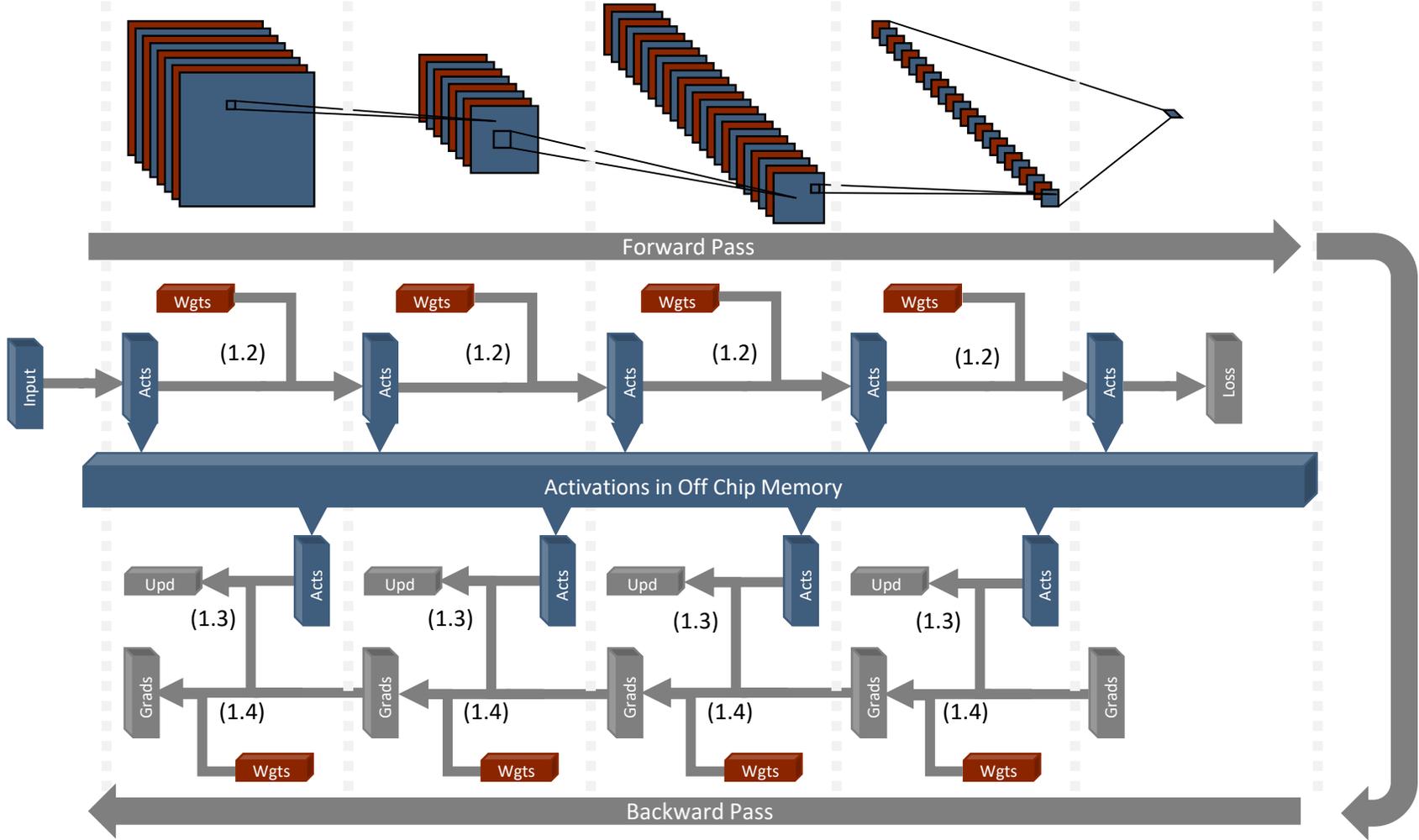  - Energy
  - Time

# Training is Expensive!

- Neural Networks are becoming ubiquitous
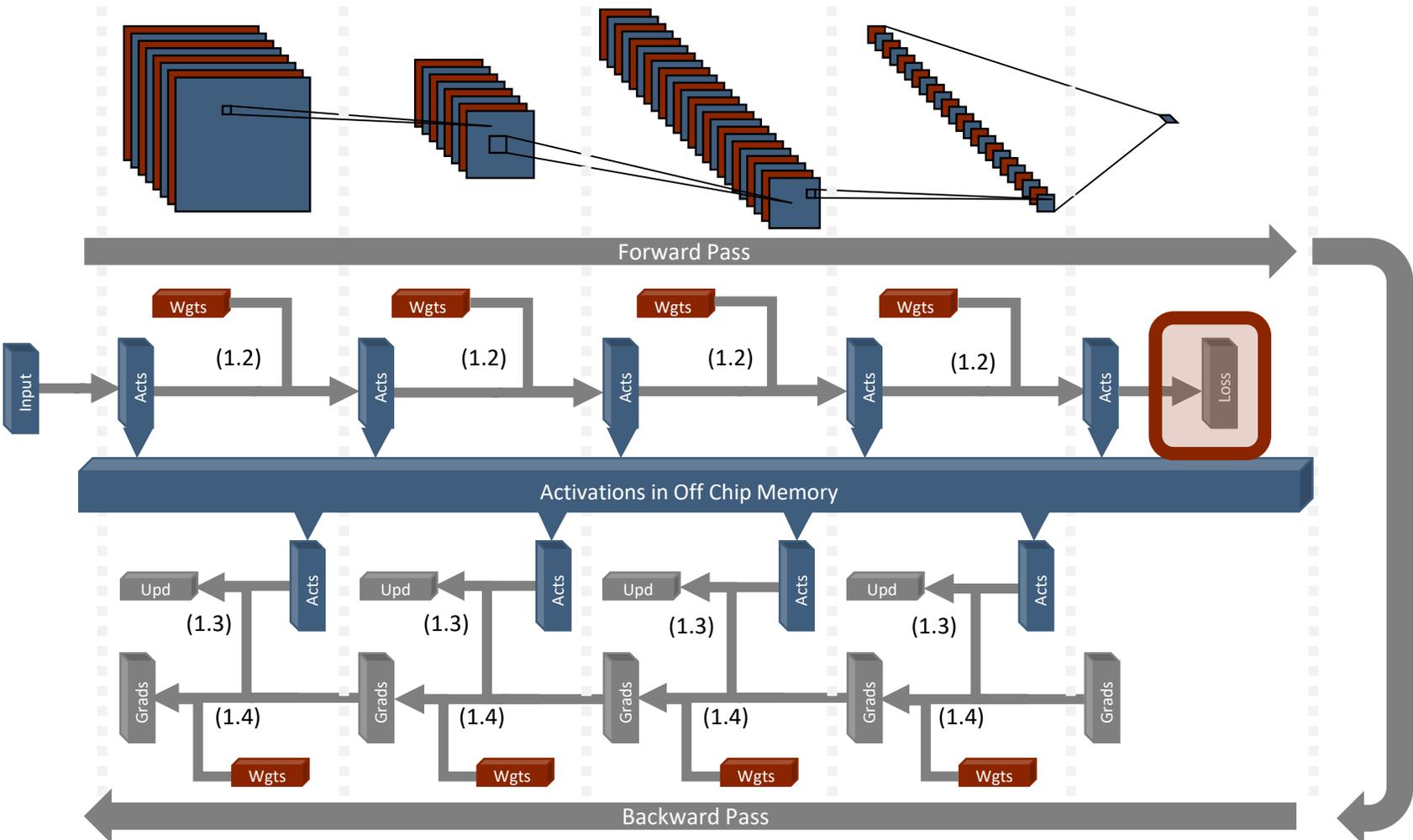- But they are expensive
  - Energy
  - Time

**Goal: Improve energy efficiency and performance of Training**

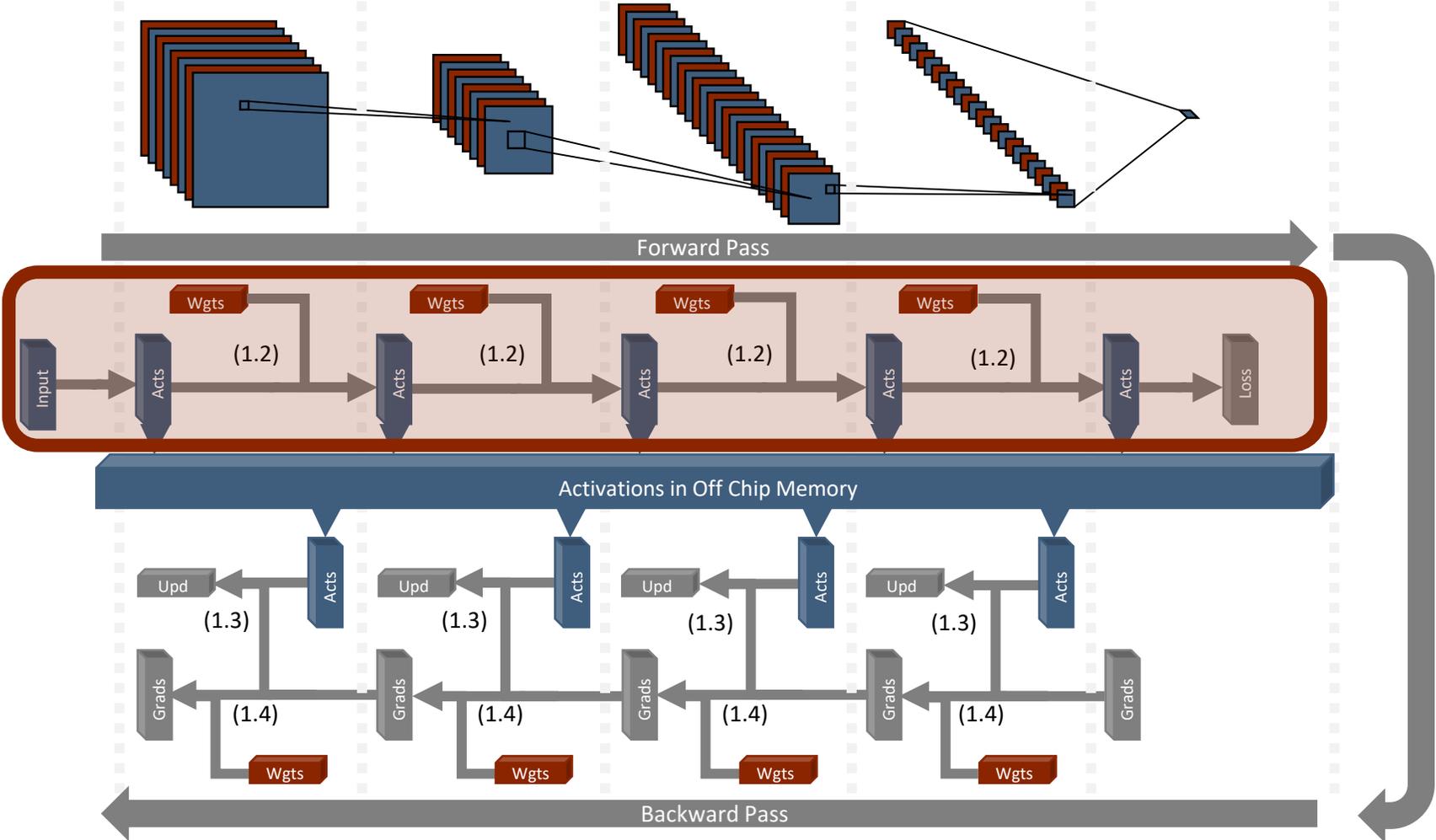# Gradient Descent – Overview

# Gradient Descent – Overview
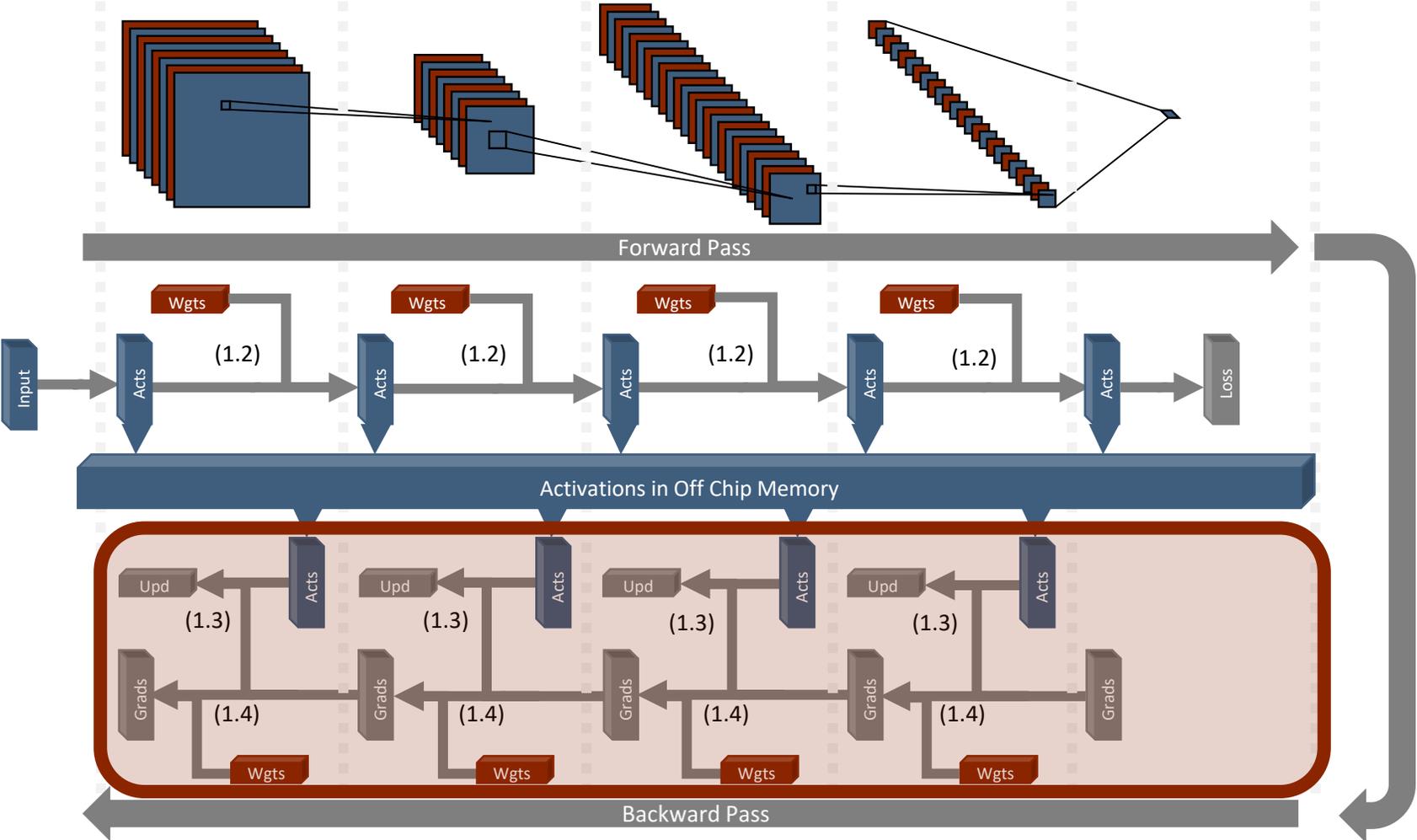
- Loss function

# Gradient Descent – Overview

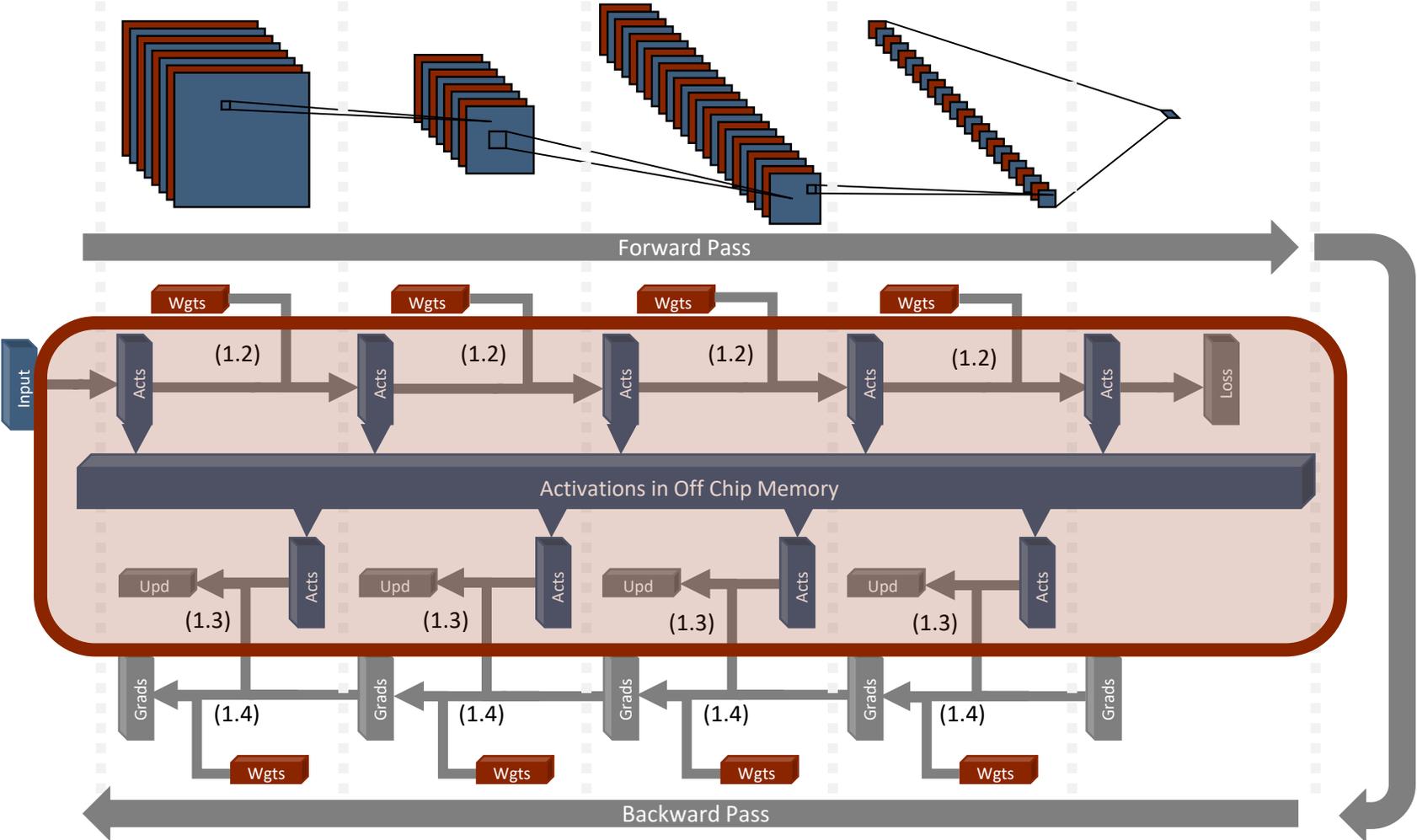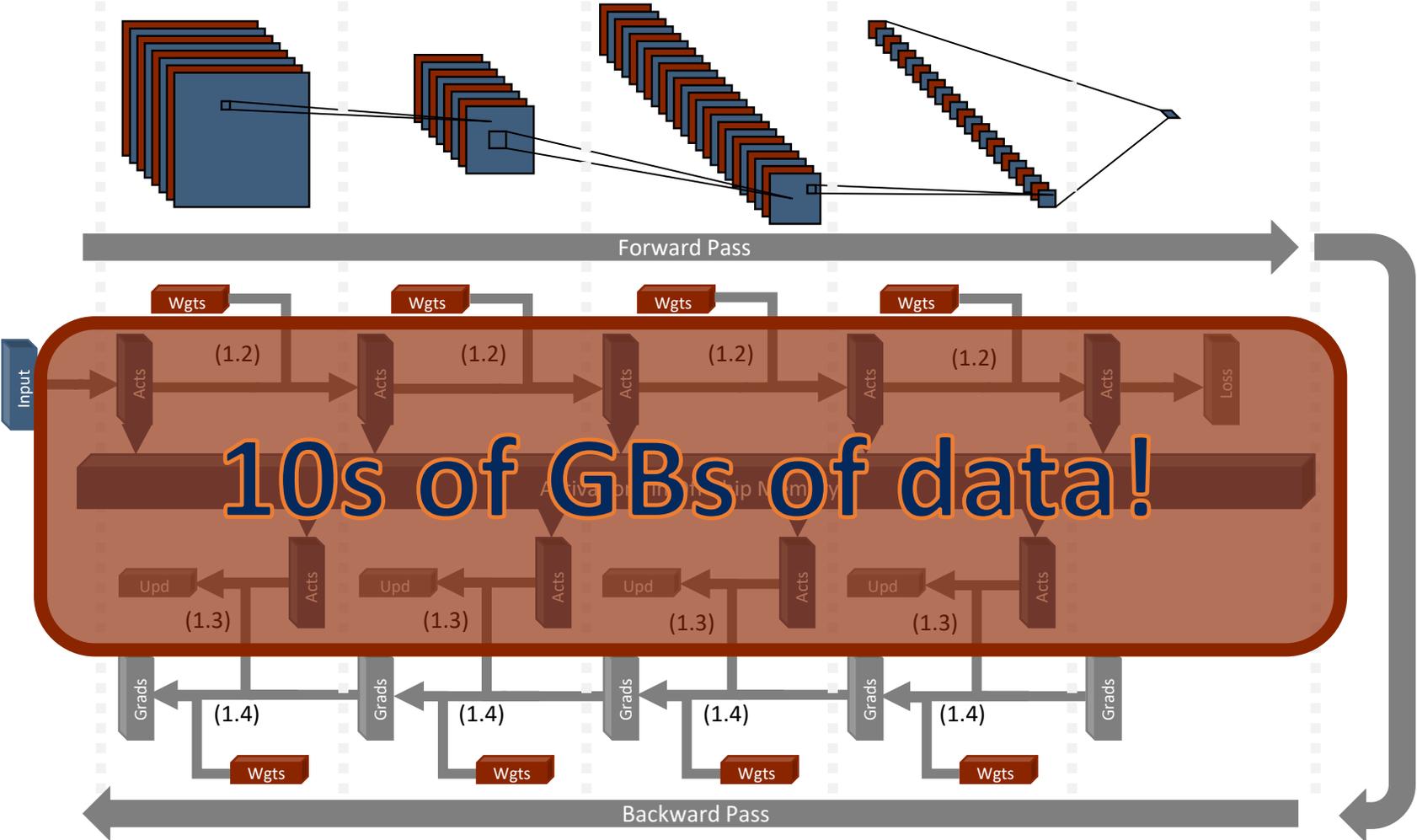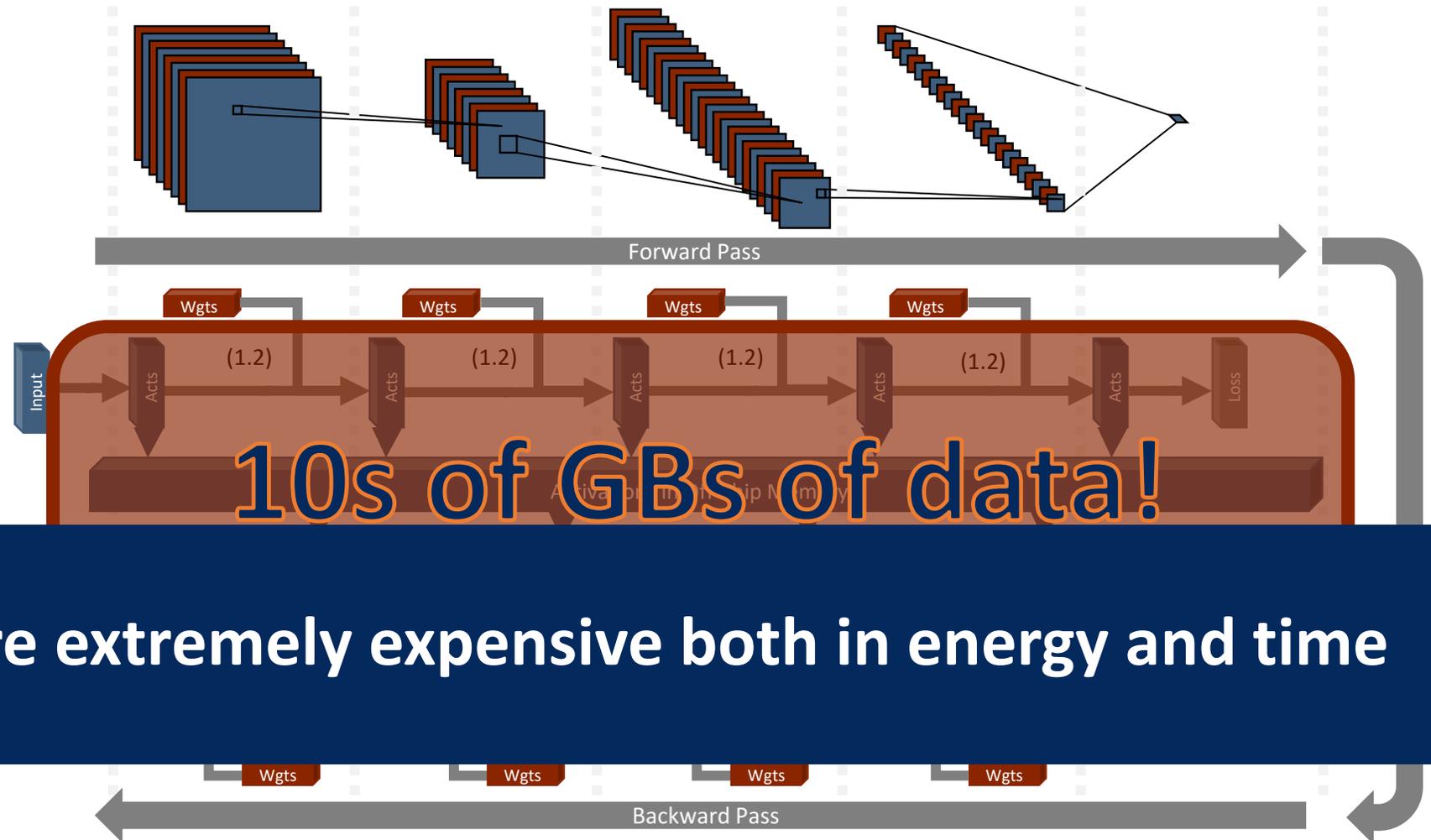- Loss function
- Forward pass

# Gradient Descent – Overview

- Loss function
- Forward pass
- Backward pass
  - Update
  - Gradient

# Gradient Descent – Overview

- Loss function
- Forward pass
- Backward pass
  - Update
  - Gradient

# Gradient Descent – Overview

- Loss function
- Forward pass
- Backward pass
  - Update
  - Gradient

# Gradient Descent – Overview

- Loss function
- Forward pass
- Backward pass
  - Update
  - Gradient



**10s of GBs of data!**

**Off-chip accesses are extremely expensive both in energy and time**

# Narrow Datatypes can help

- Efficient Datatype reduce costs!
  - Memory
  - Compute

# Narrow Datatypes can help

- Efficient Datatype reduce costs!
  - Memory
  - Compute

FP32

# Narrow Datatypes can help

- Efficient Datatype reduce costs!
  - Memory
  - Compute

# Narrow Datatypes can help

- Efficient Datatype reduce costs!
  - Memory
  - Compute

| FP8 | FP8 | FP8 | FP8 |
|-----|-----|-----|-----|

# Narrow Datatypes can help

- Efficient Datatype reduce costs!
  - Memory
  - Compute

| FP6 | FP6 | FP6 | FP6 | FP6 | ... |

# Which Datatype? Where? When?

- There are many great options

# Which Datatype? Where? When?

- There are many great options

- But it's a trade-off
  - Accuracy vs. Efficiency
  - All have edge cases



4

# Which Datatype? Where? When?

- There are many great options

- But it's a trade-off
  - Accuracy vs. Efficiency
  - All have edge cases



4

# Which Datatype? Where? When?

- There are many great options

- But it's a trade-off
  - Accuracy vs. Efficiency
  - All have edge cases

# Which Datatype? Where? When?

- There are many great options

- But it's a trade-off
  - Accuracy vs. Efficiency
  - All have edge cases

- Guess
  - Confirm
  - Repeat



4

# Which Datatype? Where? When?

- There are many great options

- But it's a trade-off
  - Accuracy vs. Efficiency
  - All have edge cases

- Guess
  - Confirm
  - Repeat

- Lost Potential

# Which Datatype? Where? When?

- There are many great options

- But it's a trade-off
  - Accuracy vs. Efficiency
  - All have edge cases

- Guess
  - Confirm
  - Repeat

**Automate the datatype selection!** 😎

# Our Goal

- Oracle?

# Our Goal

- Oracle?

- Learn the Datatype
  - Use gradient descent

# Our Goal

- Oracle?

- Learn the Datatype
  - Use gradient descent

- Monitor Loss
  - Adapt accordingly

# Our Goal

- Oracle?
- Learn the Datatype
  - Use gradient descent
- Monitor Loss
  - Adapt accordingly
- Result
  - Less bits
  - Less energy
  - Less time
  - Bigger models

# Datatype – Floating Point

- $Value = (-1)^S \times M \times 2^E$

| S | M | M | M | M | M | M | M | E | E | E | E | E | E | E | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Datatype – Floating Point

- $Value = (-1)^S \times M \times 2^E$
- Sign is trivial

| S | M | M | M | M | M | M | E | E | E | E | E | E | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Datatype – Floating Point

- $Value = (-1)^S \times M \times 2^E$
- Sign is trivial
- Mantissa
  - Precision

# Datatype – Floating Point

- $Value = (-1)^S \times M \times 2^E$
- Sign is trivial
- Mantissa
  - Precision
- Exponent
  - Range

| S | M | M | M | M | M | M | M | E | E | E | E | E | E | E | E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Schrödinger's FP

- Machine Learning
  - *Quantum Mantissa and Exponent* (**4.7x** reduction)

- Black Box Sampling
  - *BitWave* (**3.2x** reduction)

- Exponent compression
  - *Gecko* (boost to **5.6x** and **4.6x**)

- Hardware IP blocks



7

# Schrödinger's FP

- Machine Learning
  - *Quantum Mantissa and Exponent*
- Black Box Sampling
  - *BitWave*
- Exponent compression
  - *Gecko*
- Hardware IP blocks

| BitWave | Quantum Mantissa and Exponent |
|---------|-------------------------------|

Datatype Selection

| Gecko |
|-------|

Lossless Compression

| IP Blocks |
|-----------|

Implementation

# Per Tensor Datatype

- ResNet18
  - Per Tensor?
  - Single datatype?

# Per Tensor Datatype

- ResNet18
  - Per Tensor?
  - Single datatype?

# Per Tensor Datatype

- ResNet18
  - Per Tensor?
  - Single datatype?

# Per Tensor Datatype

- ResNet18
  - Per Tensor?
  - Single datatype?

# Per Tensor Datatype

- ResNet18
  - Per Tensor?
  - Single datatype?

# Per Tensor Datatype

- ResNet18
  - Per Tensor?
  - Single datatype?
- Less than 6b on average!
  - Much better than network wide!

# Per Tensor Datatype

- ResNet18
  - Per Tensor?
  - Single datatype?
- Less than 6b on average!
  - Much better than network wide!



**No chance to guess effectively and expect to train well!**

# *Quantum Mantissa Loss*

- Our Loss
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according to priority
    - Tensor footprint
    - # of operations

# Quantum Mantissa Loss

- Our Loss    **Original Loss**
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according to priority
    - Tensor footprint
    - # of operations

# *Quantum Mantissa Loss*

- Our Loss
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$  **Mantissa Lengths**
  - Weighted according to priority
    - Tensor footprint
    - # of operations

# Quantum Mantissa Loss

- Our Loss

  **Mantissa Weights**

  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according to priority
    - Tensor footprint
    - # of operations

# Quantum Mantissa Loss

- Our Loss

  **Regularizer Strength**

  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according to priority
    - Tensor footprint
    - # of operations

# Quantum Mantissa Bitlength

- Int Datatype

# Quantum Mantissa Bitlength

- Int Datatype
- Non-Int Datatype

# Quantum Mantissa Bitlength

- Int Datatype
- Non-Int Datatype

# Quantum Mantissa Bitlength

- Int Datatype
- Non-Int Datatype

# Quantum Mantissa Bitlength

- Int Datatype
- Non-Int Datatype

# Quantum Exponent Range

- Distribution



Exponents distribution biased - ImageNet - ResNet18

# Quantum Exponent Range

- Distribution



Exponents distribution biased - ImageNet - ResNet18

# Quantum Exponent Range

- Distribution
- Compress range



Exponents distribution biased - ImageNet - ResNet18

# Quantum Exponent Range

- Distribution
- Compress range
- $Value = (-1)^S \times M \times 2^E$



Exponents distribution biased - ImageNet - ResNet18

# Quantum Exponent Range

- Distribution
- Compress range
- $Value = (-1)^S \times M \times 2^E$

# Quantum Exponent Range

- Distribution
- Compress range
- $Value = (-1)^S \times M \times 2^E$
- Range(x) overlay



Exponents distribution biased - ImageNet - ResNet18

# *Quantum Exponent Range*

- Distribution
- Compress range
- $Value = (-1)^S \times M \times 2^E$
- Range(x) overlay

# Quantum Exponent

- Our Loss
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according priority

# Quantum Exponent

- Our Loss **Original Loss**
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according priority

# Quantum Exponent

- Our Loss

  **Exponent Lengths**

  - $L = L_O + \gamma \times \left( \sum \lambda_i \times \alpha_i \right)$
  - Weighted according priority

# Quantum Exponent

- Our Loss
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according priority

**Exponent Weights**

# Quantum Exponent

- Our Loss

  **Regularizer Strength**

  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according priority

# Quantum Exponent

- Our Loss
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according priority

- Int Datatype

# Quantum Exponent

- Our Loss
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according priority

- Int Datatype

- Non-Int Datatype

# Quantum Exponent

- Our Loss
  - $L = L_O + \gamma \times (\sum \lambda_i \times \alpha_i)$
  - Weighted according priority
- Int Datatype
- Non-Int Datatype
- Determine boundaries

# Schrödinger's FP

- Machine Learning
  - *Quantum Mantissa and Exponent*

- Black Box Sampling
  - *BitWave*

- Exponent compression
  - *Gecko*

- Hardware IP blocks



**BitWave** | **Quantum Mantissa and Exponent**

Datatype Selection

**Gecko**

Lossless Compression

**IP Blocks**

Implementation

# *BitWave*

- Observe Loss
  - Adjust bitlengths
    - Mantissa
    - Exponent
- Training is forgiving
- Network wide

# Schrödinger's FP

- Machine Learning
  - *Quantum Mantissa and Exponent*

- Black Box Sampling
  - *BitWave*

- Exponent compression
  - *Gecko*

- Hardware IP blocks



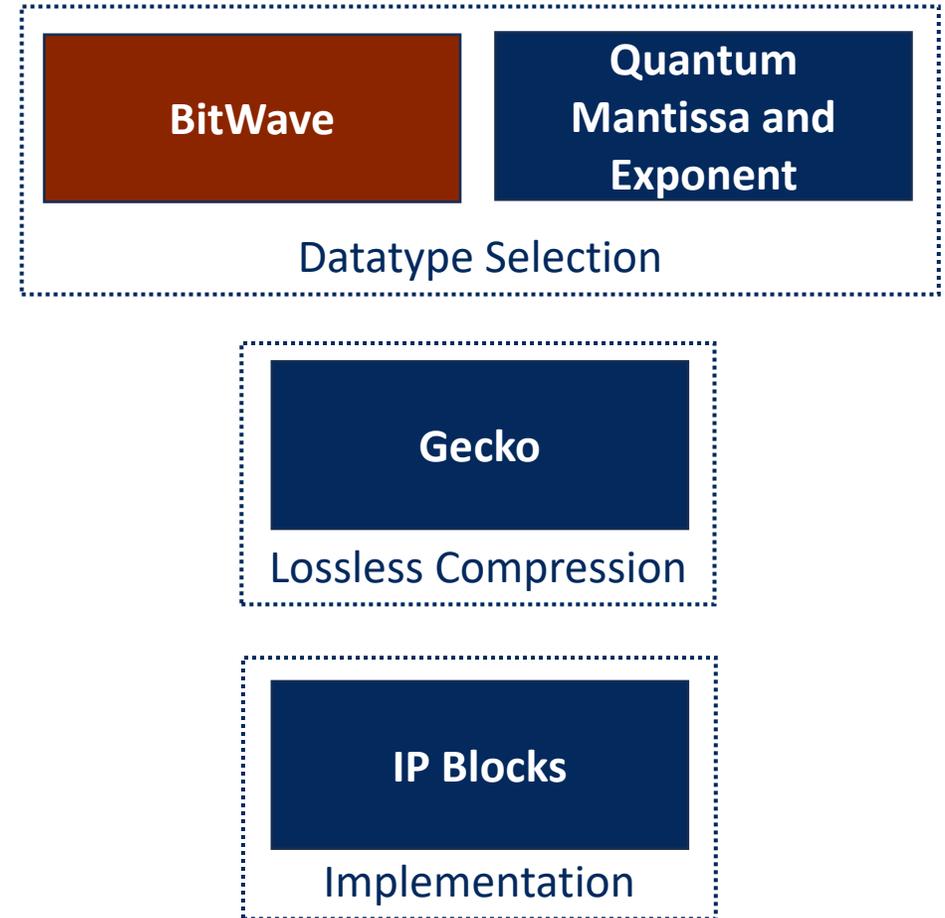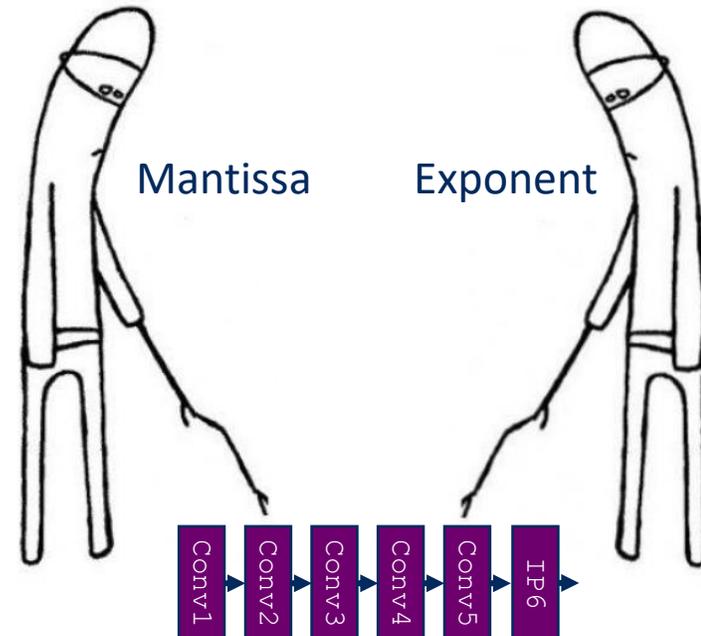| BitWave | Quantum Mantissa and Exponent |
|---------|-------------------------------|

Datatype Selection

| Gecko |
|-------|

Lossless Compression

| IP Blocks |
|-----------|

Implementation

# Exponent – *Gecko*

- Distribution



Exponents distribution biased - ImageNet - ResNet18

# Exponent – *Gecko*

- Distribution
- Group



Exponents distribution biased - ImageNet - ResNet18

# Exponent – *Gecko*

- Distribution
- Group

# Exponent – *Gecko*

- Distribution
- Group

# Massive footprint reduction on ResNet18!

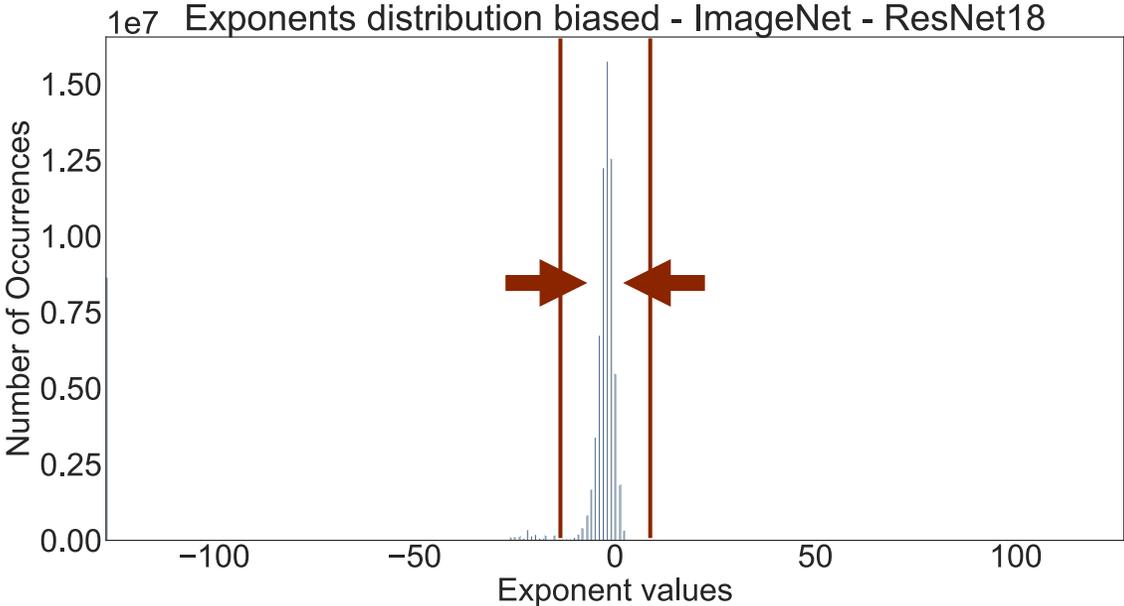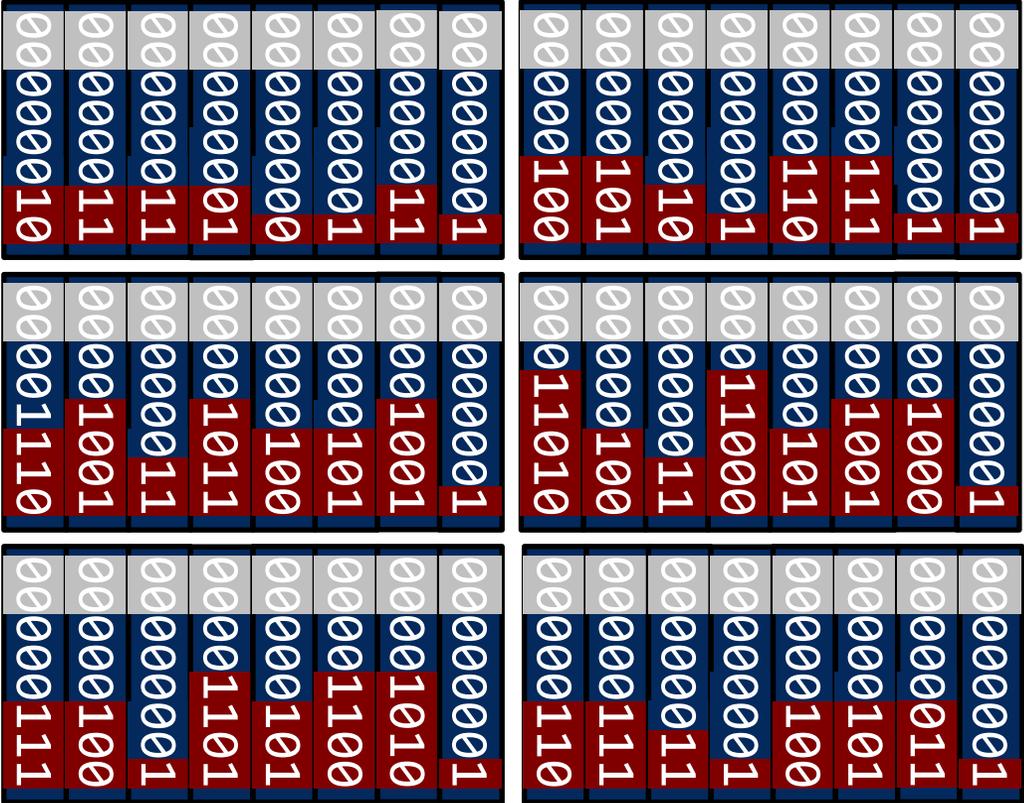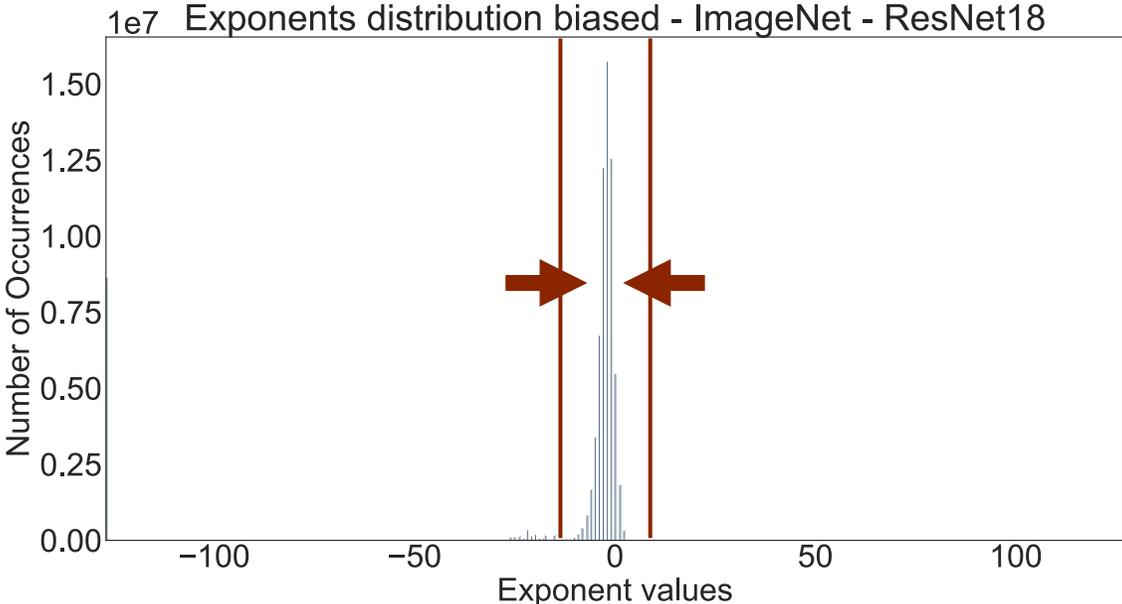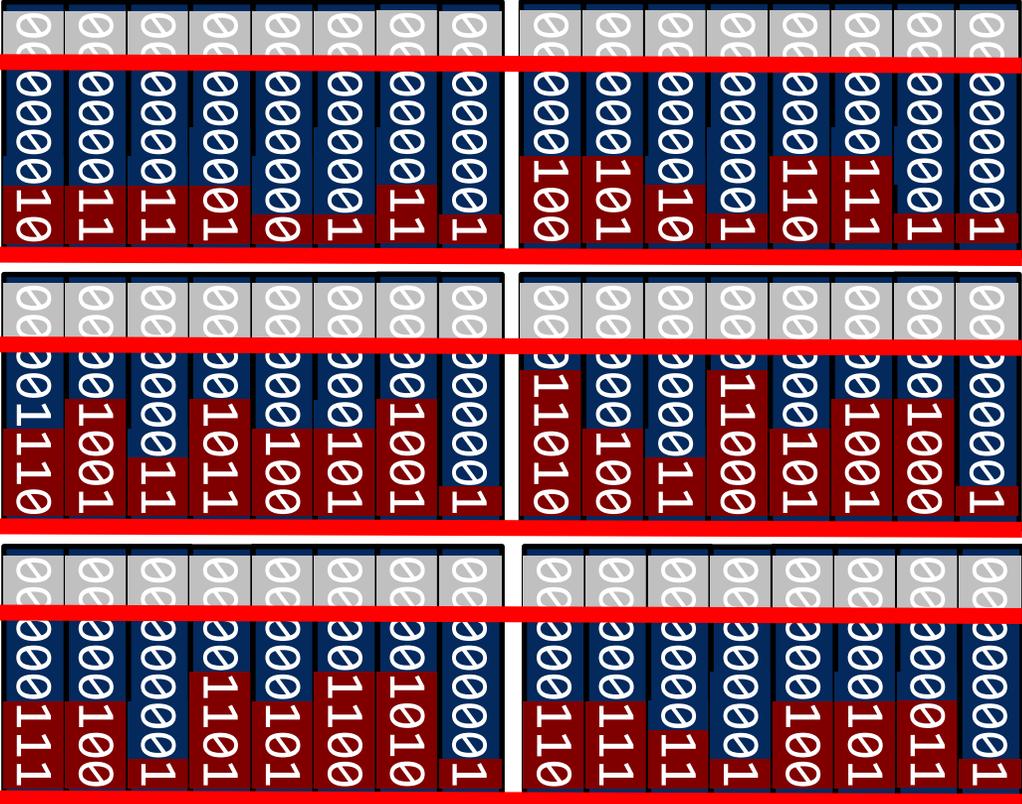# Massive footprint reduction on ResNet18!

# Massive footprint reduction on ResNet18!

# Massive footprint reduction on ResNet18!

# Massive footprint reduction on ResNet18!

# Massive footprint reduction on ResNet18!

# Massive footprint reduction on ResNet18!

# Massive footprint reduction on ResNet18!

# Similar Task Performance!

# Similar Task Performance!

# Similar Task Performance!

# Similar Task Performance!
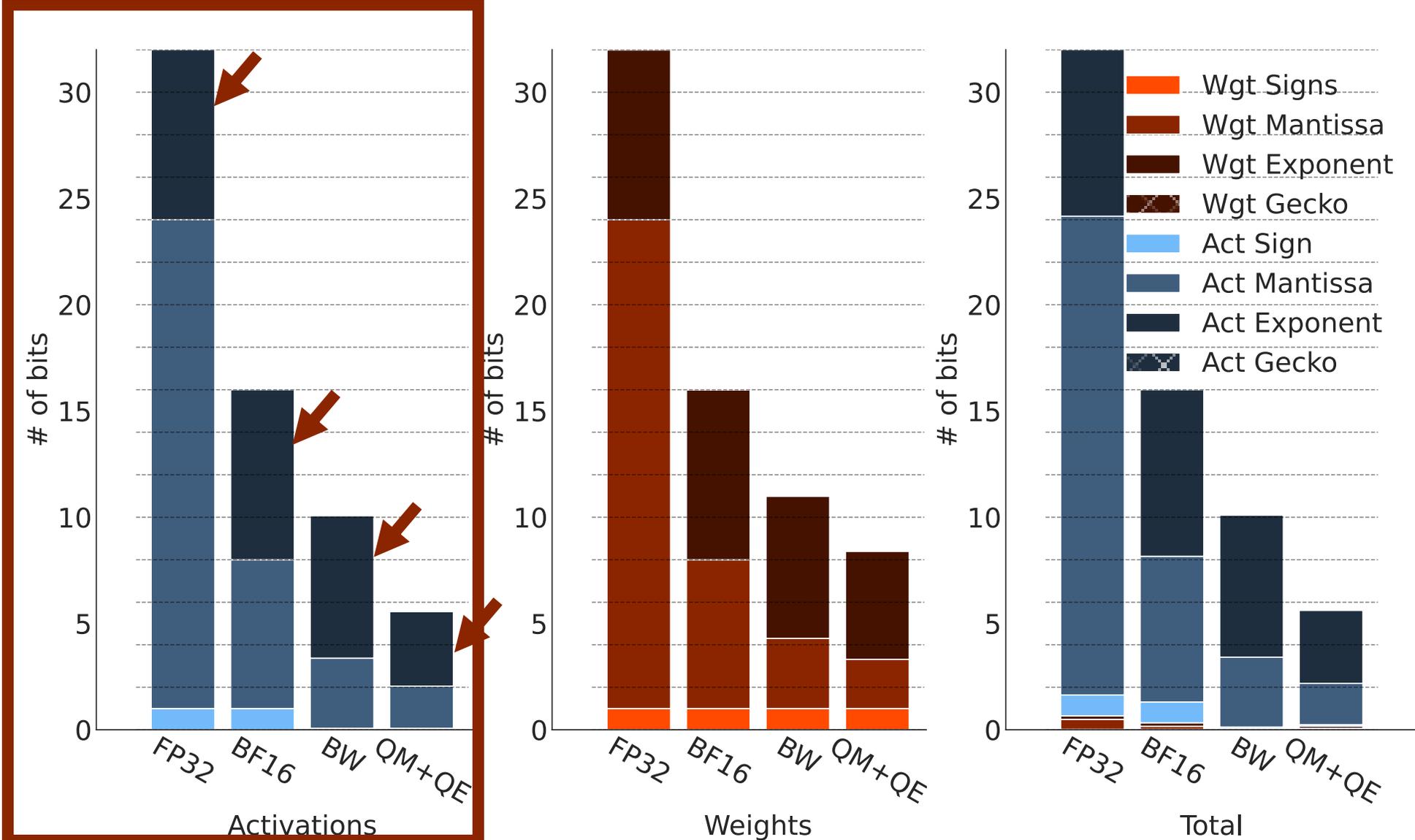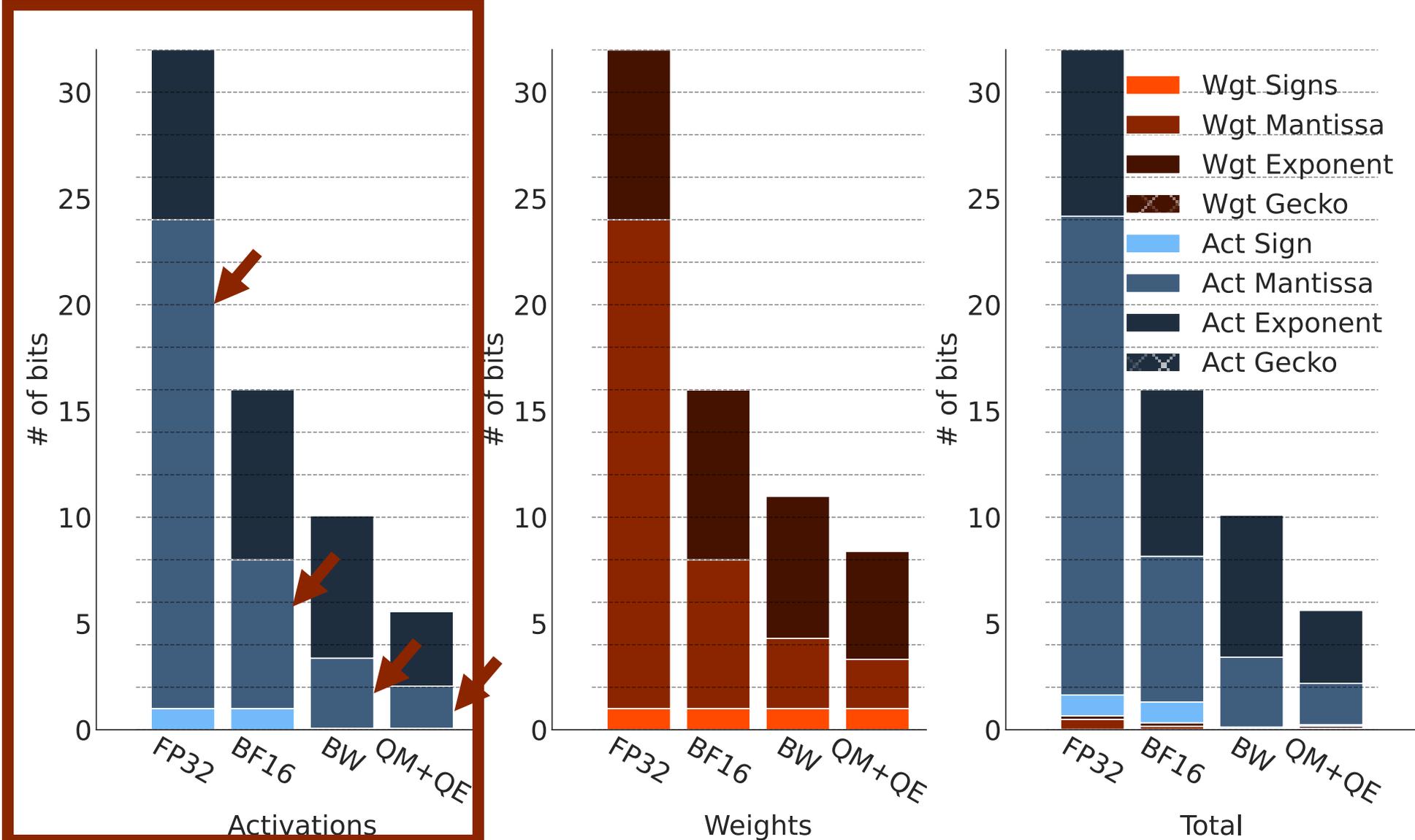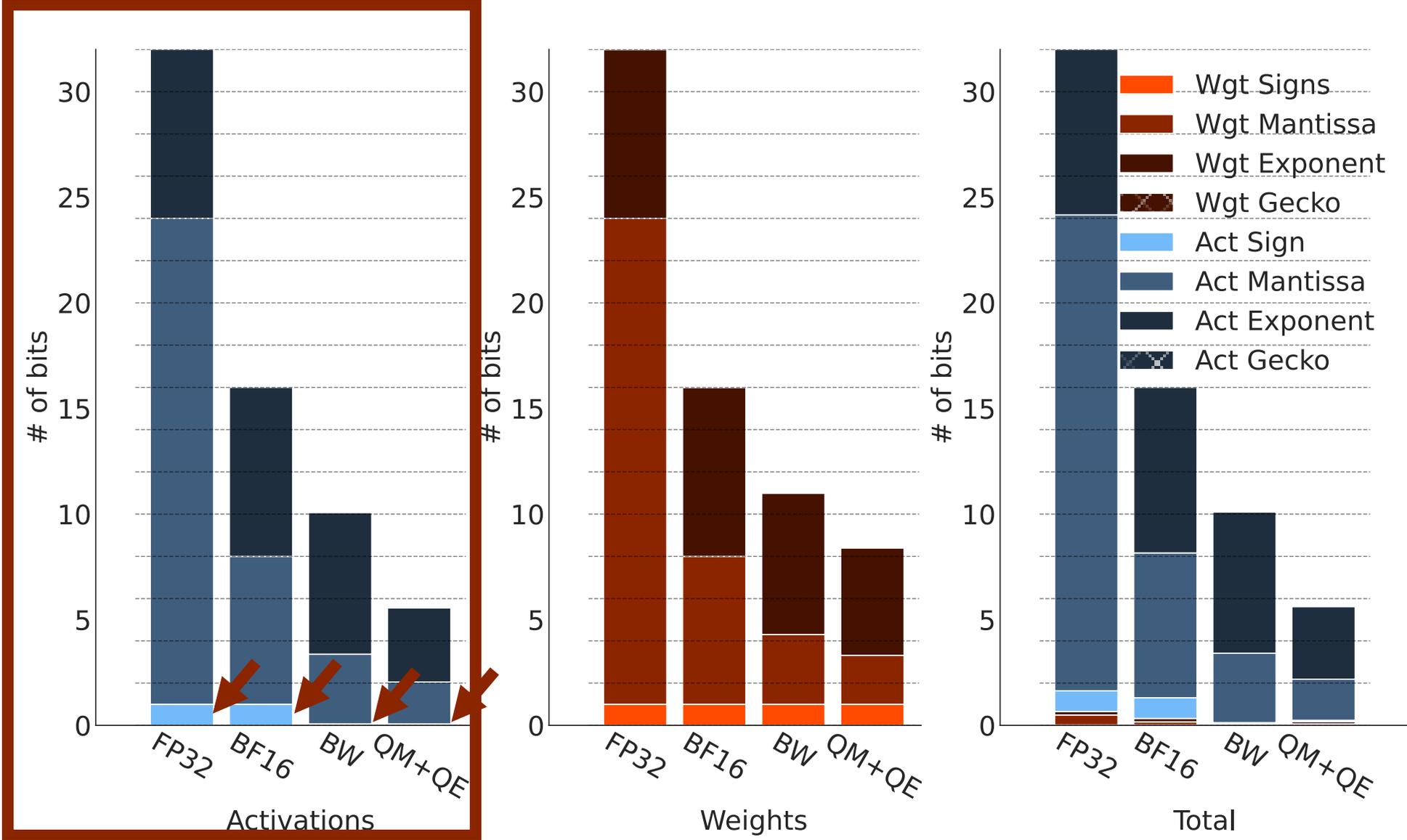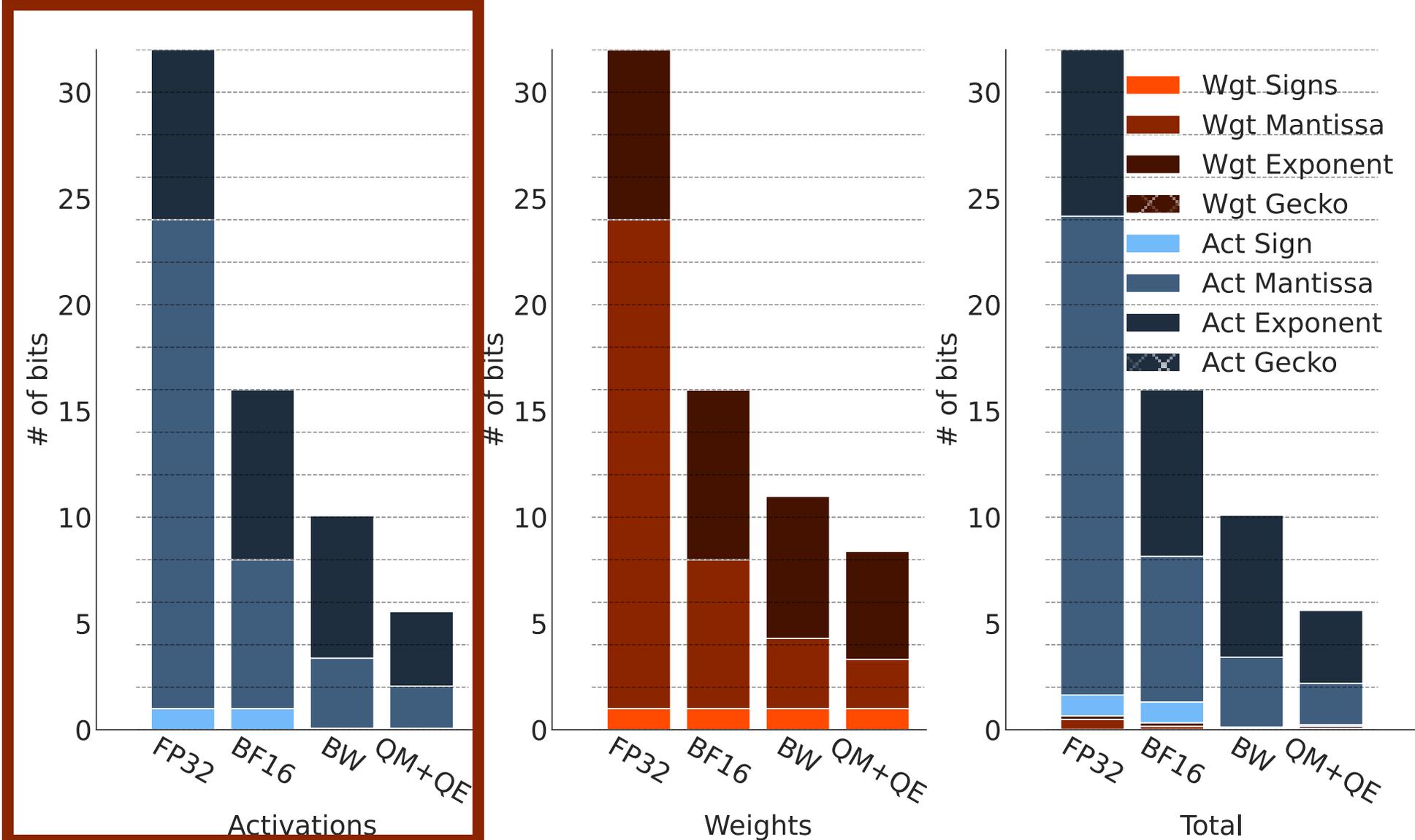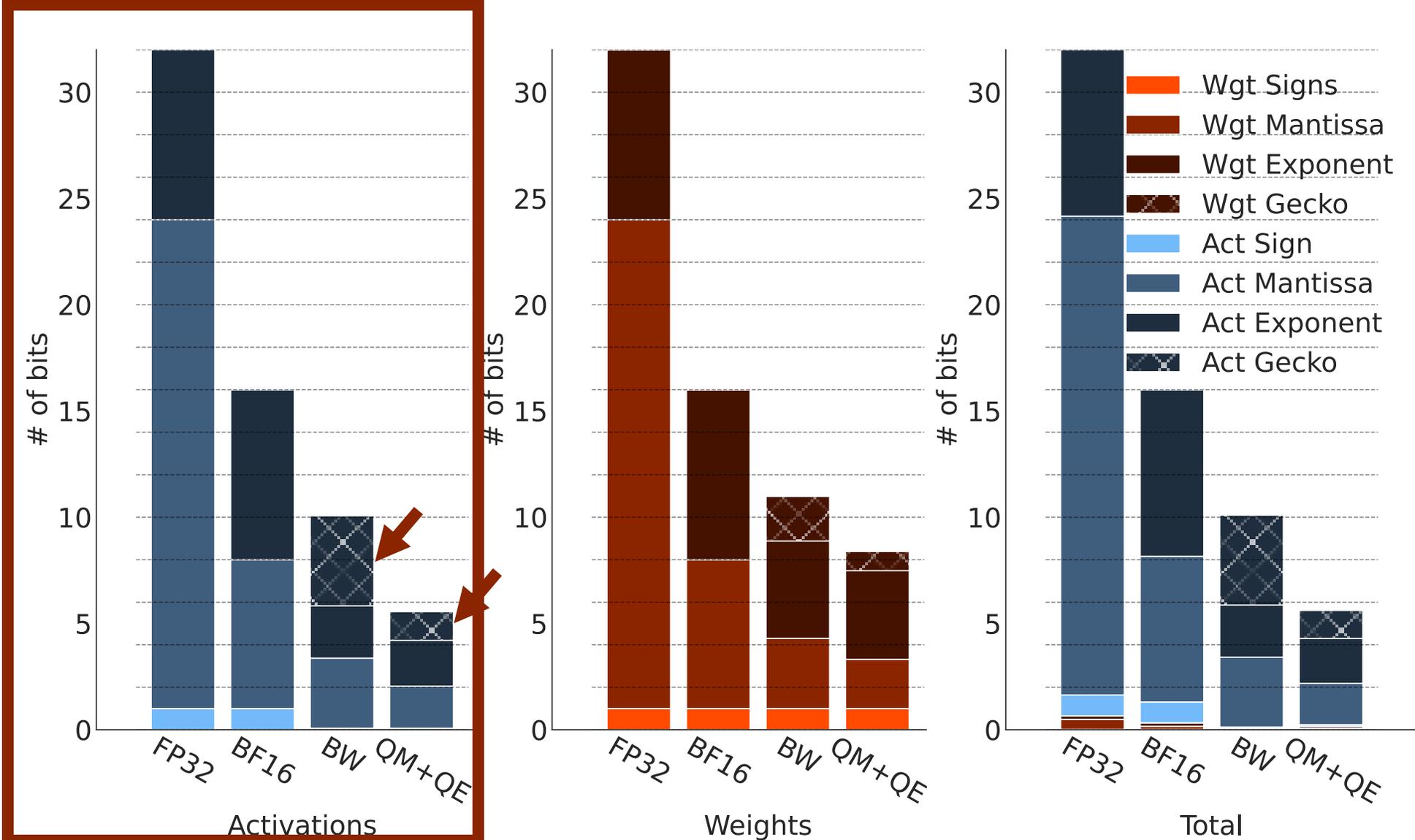
# Similar Task Performance!

# Massive footprint reduction across the board!



Bar chart titled "Massive footprint reduction across the board!" showing Footprint (y-axis, 0.0 to 1.0) for various models with legend "SFP Q + G". Models on x-axis: ResNet18, ResNet50, MobileNet V2, DLRM, ViT, GPT-2, BERT - CoLA, BERT - SST-2, BERT - MRPC, BERT - STS-B, BERT - QQP, BERT - MNLI, BERT - QNLI, Geo Mean.

# Massive footprint reduction across the board!

# Massive footprint reduction across the board!

# Massive footprint reduction across the board!

# Massive footprint reduction across the board!

# Massive footprint reduction across the board!

# Massive footprint reduction across the board!

# Massive footprint reduction across the board!

# How to exploit?

- Between on- and off-chip memory

# How to exploit?

- Between on- and off-chip memory
- Seamless

# How to exploit?

- Between on- and off-chip memory
- Seamless
- Evaluate in hardware
  - Can be done in software

# How to exploit?

- Between on- and off-chip memory
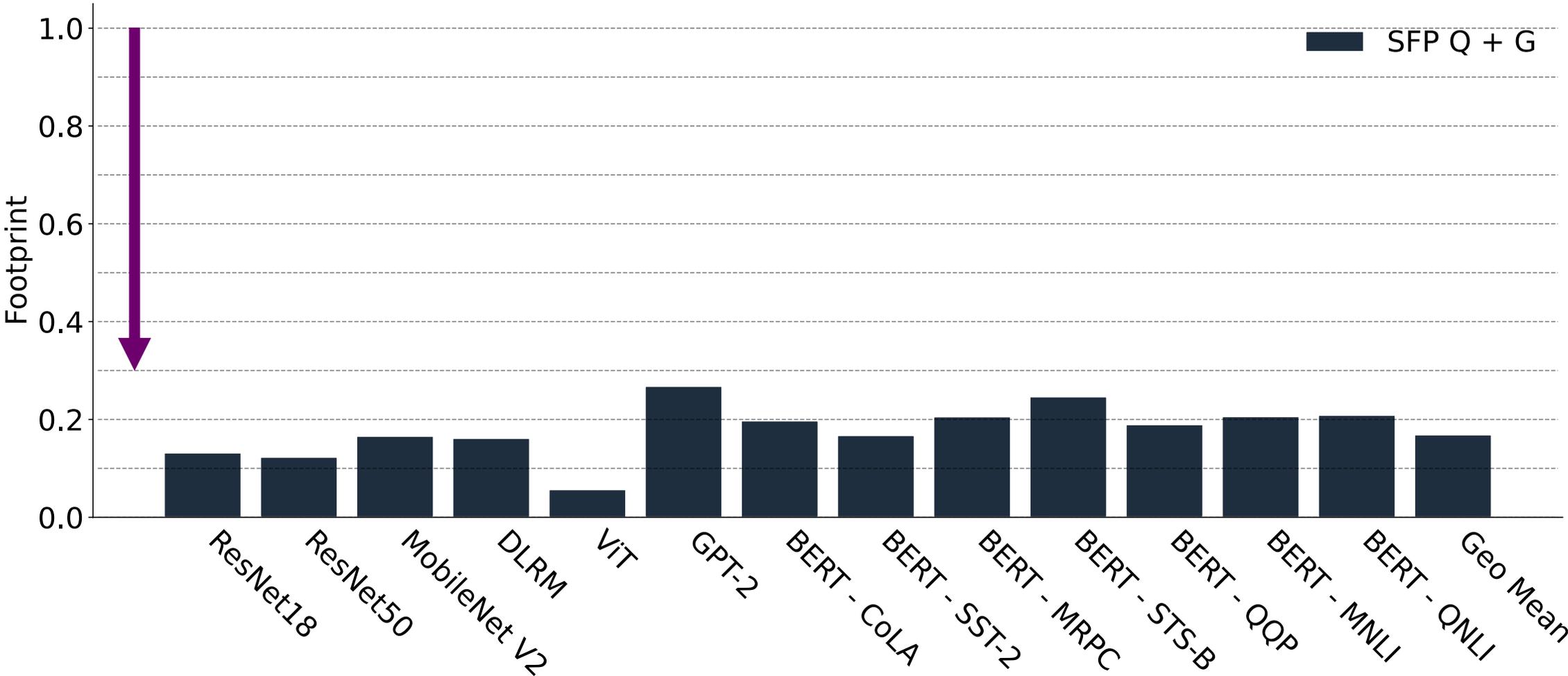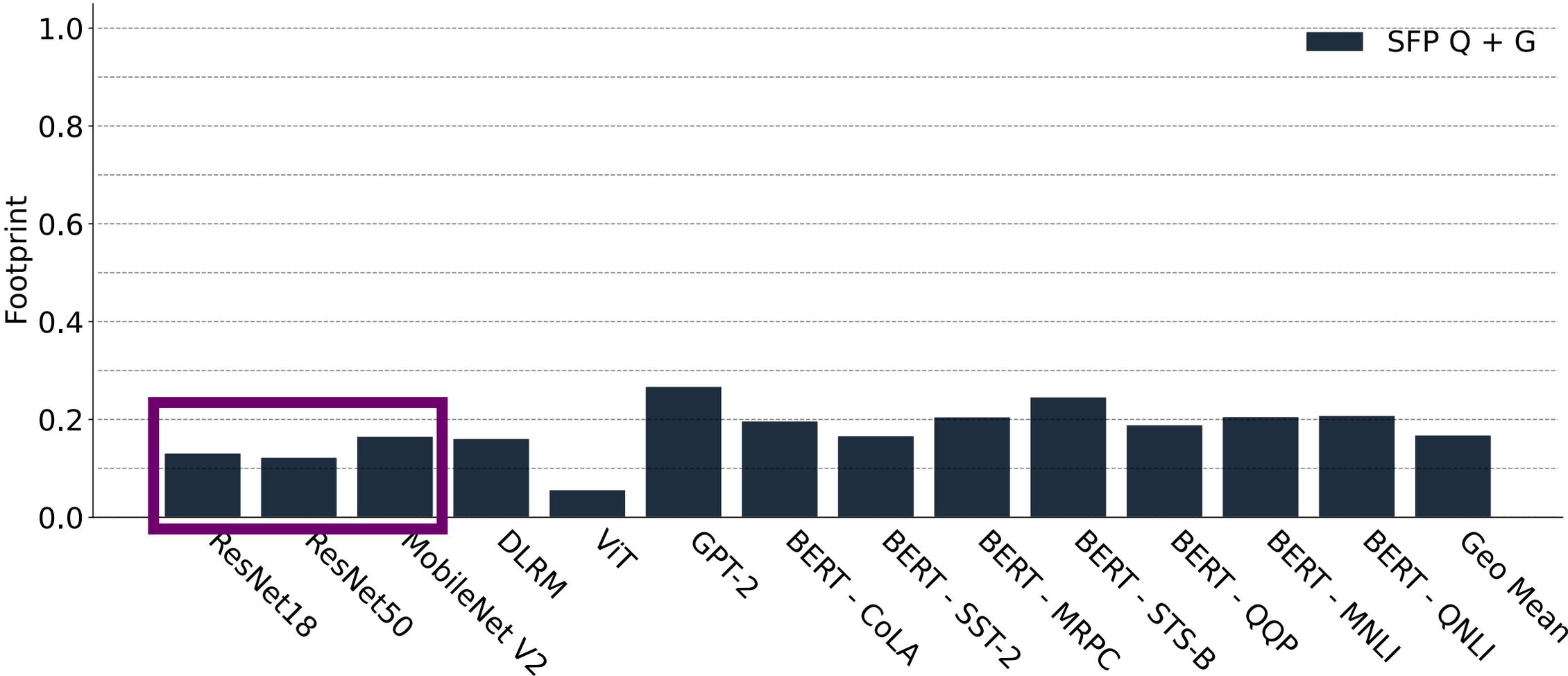- Seamless
- Evaluate in hardware
  - Can be done in software
- vs. FP8 baseline
  - **2.6x** performance
  - **2.3x** energy efficiency



21

# More in the paper!

- *BitWave* description & results
- Behavior through time
- Results for every network
- Hardware IP blocks
- Simulation results
- Simulation details

# Conclusion

- **Schrödinger's FP**: Automatic datatype selection and compression methods

# Conclusion

- **Schrödinger's FP**: Automatic datatype selection and compression methods
  - Machine Learning
    - **4.7x** & **5.6x** (*+G*)

# Conclusion

- **Schrödinger's FP**: Automatic datatype selection and compression methods
  - Machine Learning
    - **4.7x** & **5.6x** (*+G*)
  - Sampling
    - **3.2x** & **4.6x** (*+G*)

# Conclusion

- **Schrödinger's FP**: Automatic datatype selection and compression methods
  - Machine Learning
    - **4.7x** & **5.6x** (*+G*)
  - Sampling
    - **3.2x** & **4.6x** (*+G*)
  - Optionally *+Gecko*

# Conclusion

- **Schrödinger's FP**: Automatic datatype selection and compression methods
  - Machine Learning
    - **4.7x** & **5.6x** (*+G*)
  - Sampling
    - **3.2x** & **4.6x** (*+G*)
  - Optionally *+Gecko*

## Questions?

**Miloš Nikolić** <milos.nikolic@mail.utoronto.ca>