# [FKeras](#): A Sensitivity Analysis Tool for Edge Neural Networks

Olivia Weng, **Andres Meza**, Quinlan Bock, Benjamin Hawks, Javier Campos, Nhan Tran, Javier Duarte, Ryan Kastner

Presented on November 2nd, 2023 at Fast ML for Science @ ICCAD 2023

**‡ Fermilab**    **UC San Diego**    **hls 4 ml**

# Why analyze the **sensitivity** of edge NNs?

# Why analyze the sensitivity of edge NNs?

- Edge NNs can be sensitive creatures

# Why analyze the sensitivity of edge NNs?

- Edge NNs can be sensitive creatures

- They are put through a lot:

# Why analyze the sensitivity of edge NNs?

- Edge NNs can be sensitive creatures

- They are put through a lot:
    - Pruning
    - Quantization
    - Hardware faults

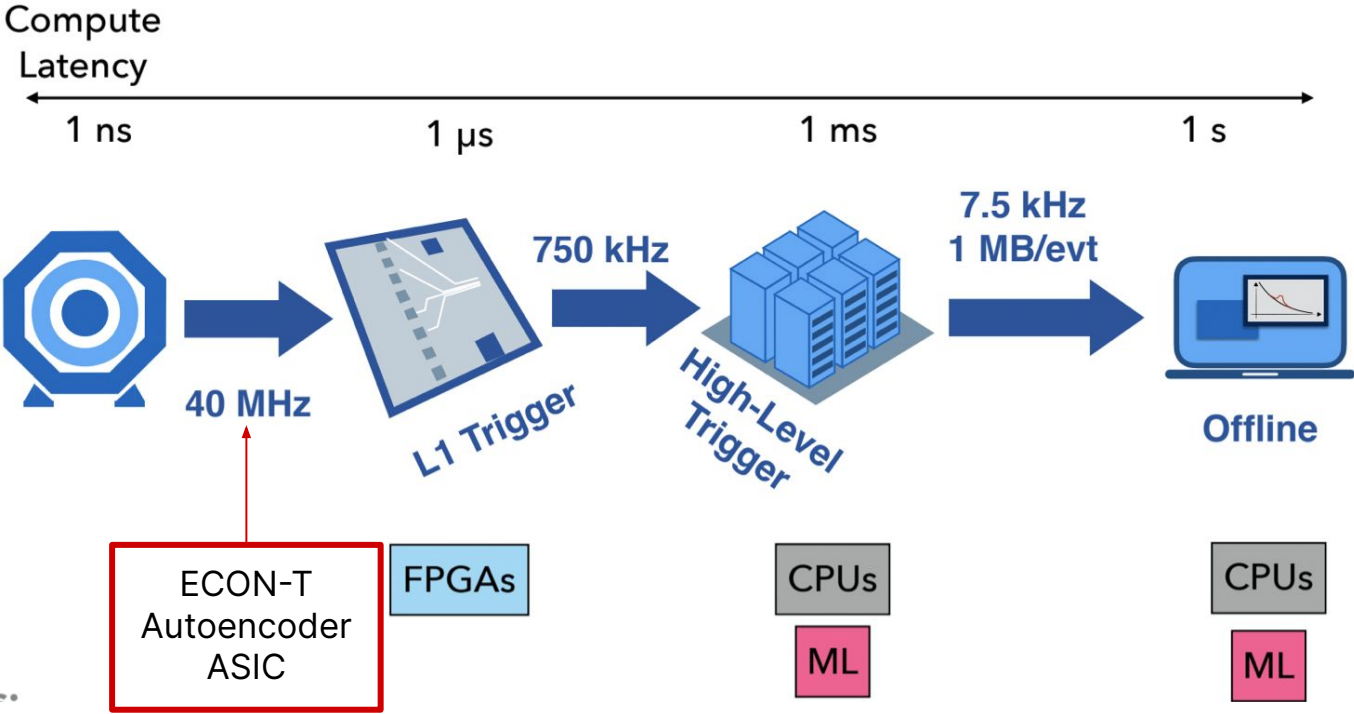# Why analyze the sensitivity of edge NNs?

- Edge NNs can be sensitive creatures

- They are put through a lot:
  - Pruning
  - Quantization
  - **Hardware faults**

# Why analyze the sensitivity of edge NNs?

- Edge NNs can be sensitive creatures

- They are put through a lot:
  - Pruning
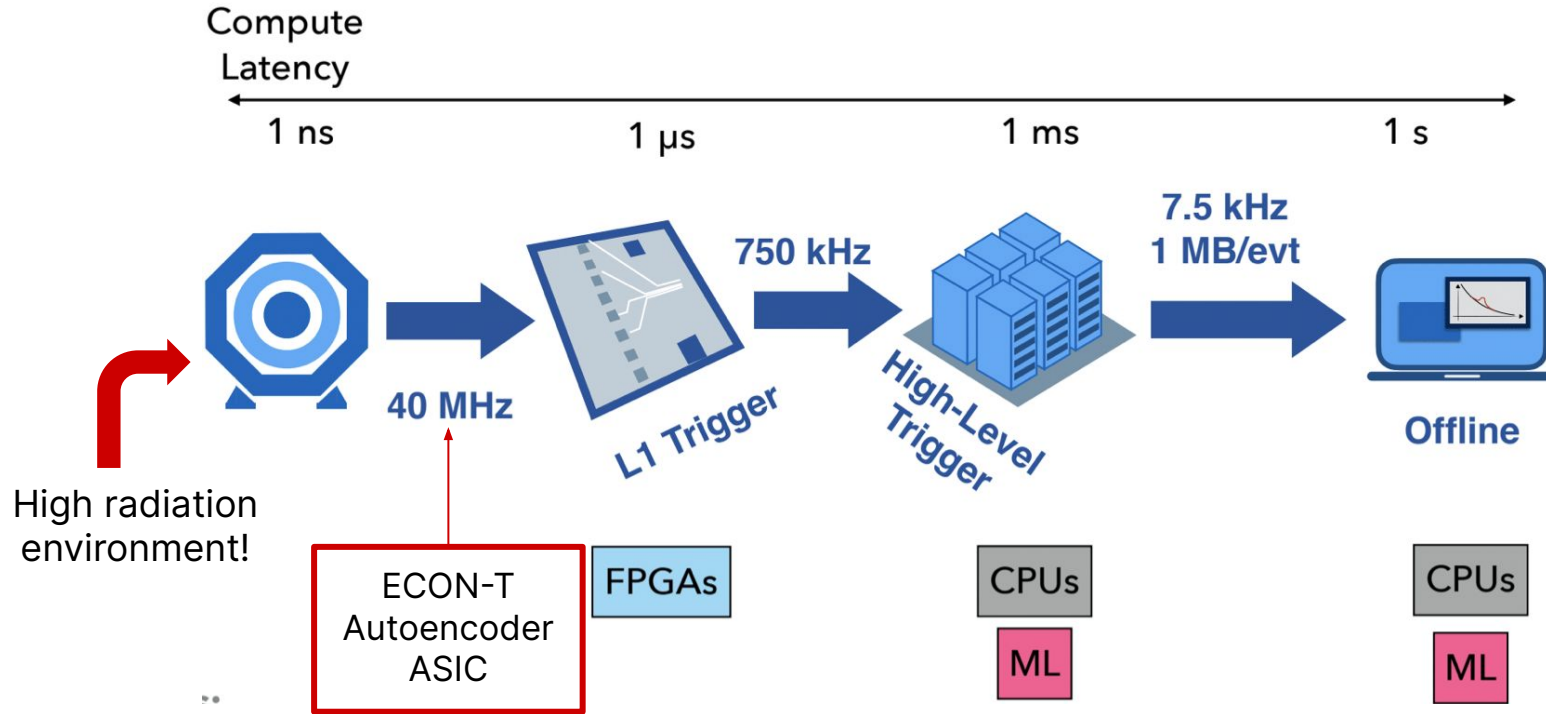  - Quantization
  - **Hardware faults**

When do **hardware faults** occur?

# Example: LHC's CMS Data Processing Pipeline

Ref: https://indico.fnal.gov/event/46746/contributions/210450/attachments/141293/177902/hirschauer_AE_CPAD_19mar2020.pdf

# Example: LHC's CMS Data Processing Pipeline



Compute Latency

1 ns — 1 μs — 1 ms — 1 s

40 MHz → L1 Trigger → 750 kHz → High-Level Trigger → 7.5 kHz 1 MB/evt → Offline

High radiation environment!

ECON-T Autoencoder ASIC
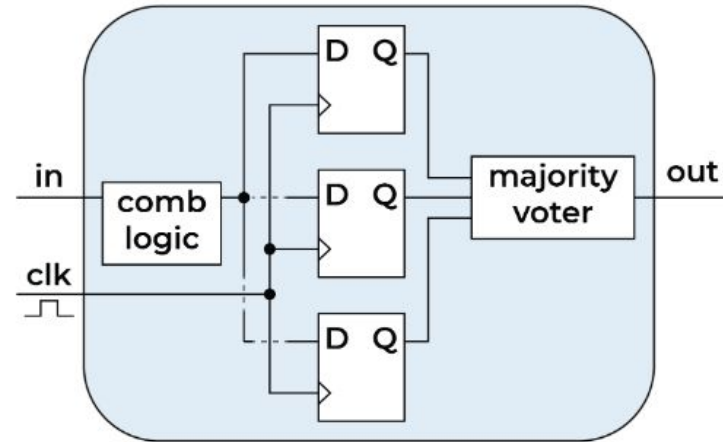
FPGAs

CPUs

ML

CPUs

ML

# How does the ECON-T Autoencoder **tolerate** radiation?

# ECON-T Radiation Tolerance

- ECON-T employs **triple modular redundancy (TMR)** to its registers

Ref: Di Guglielmo et al. A reconfigurable neural network ASIC for detector front-end data compression at the HL-LHC. IEEE Trans. Nucl. Sci.'21.

# ECON-T Radiation Tolerance

- ECON-T employs **triple modular redundancy (TMR)** to its registers

Ref: Di Guglielmo et al. A reconfigurable neural network ASIC for detector front-end data compression at the HL-LHC. IEEE Trans. Nucl. Sci.'21.

# ECON-T Radiation Tolerance

- ECON-T employs **triple modular redundancy (TMR)** to its registers



Ref: Di Guglielmo et al. A reconfigurable neural network ASIC for detector front-end data compression at the HL-LHC. IEEE Trans. Nucl. Sci.'21.
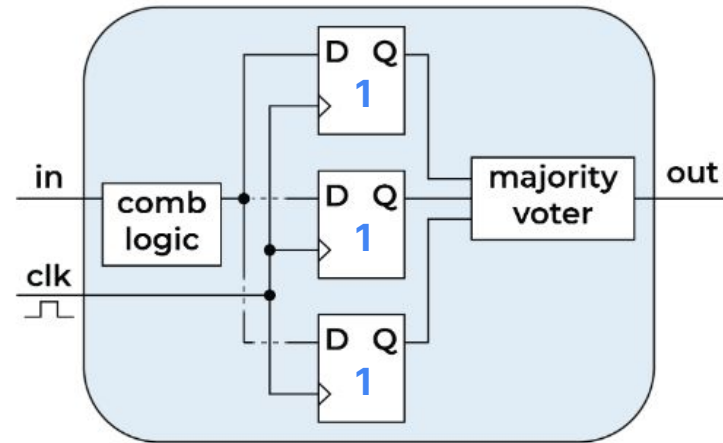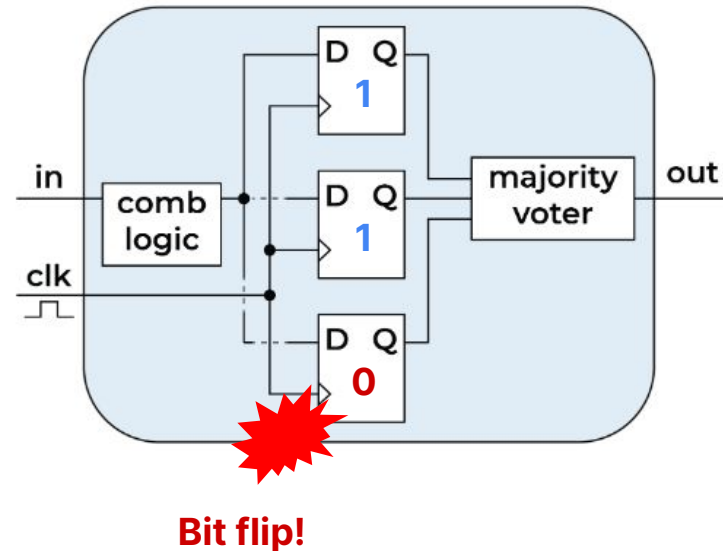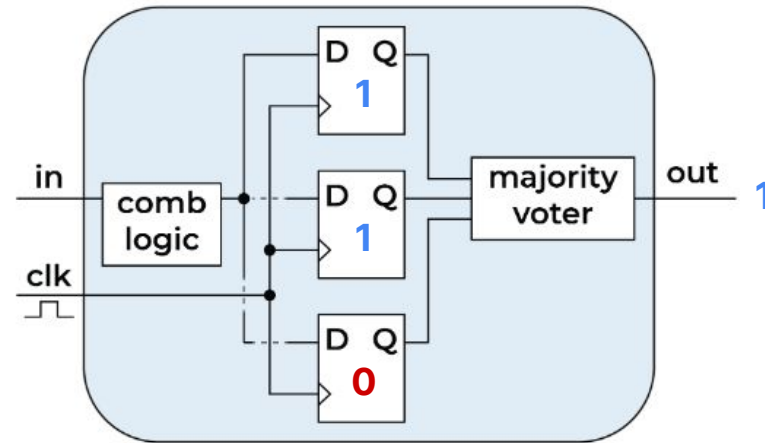
# ECON-T Radiation Tolerance

- ECON-T employs **triple modular redundancy (TMR)** to its registers

# ECON-T Radiation Tolerance

- ECON-T employs **triple modular redundancy (TMR)** to its registers

# ECON-T Radiation Tolerance

- ECON-T employs **triple modular redundancy (TMR)** to its registers



**TMR incurs ≥200% area overhead!**

Ref: Di Guglielmo et al. A reconfigurable neural network ASIC for detector front-end data compression at the HL-LHC. IEEE Trans. Nucl. Sci.'21.

How can we **reduce** radiation tolerance **costs**?

Observation: Tolerance only applied to **hardware**

<u>Observation</u>: Tolerance only applied to **hardware**

What about **software**?

How should we assess
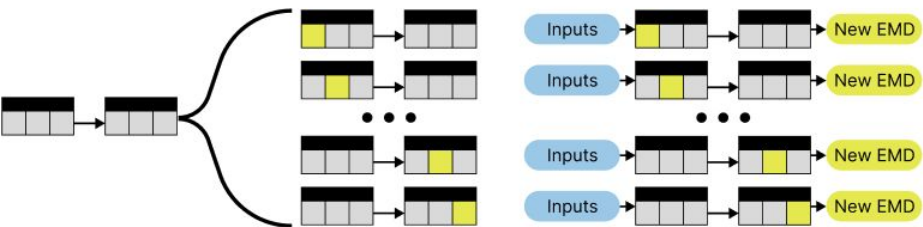the **fault sensitivity** of NN **software**?

# FKeras

- A library that assesses the fault sensitivity of (Q)Keras models

# FKeras

- A library that assesses the fault sensitivity of (Q)Keras models

- Current features include:

  - Bit-level fault injection with fine-grained control

  - Bit-level sensitivity metrics for ranking weight bits

# FKeras

- A library that assesses the fault sensitivity of (Q)Keras models

- Current features include:

  - Bit-level fault injection with fine-grained control

  - Bit-level sensitivity metrics for ranking weight bits

Bit-level fault injection with fine-grained control

Bit-level sensitivity metrics for ranking weight bits (without fault injection)
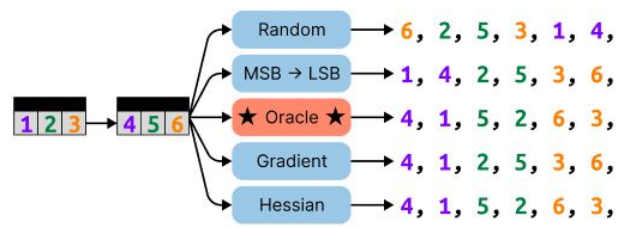
# FKeras

- A library that assesses the fault sensitivity of (Q)Keras models

- Current features include:

  - Bit-level fault injection with fine-grained control

  - Bit-level sensitivity metrics for ranking weight bits

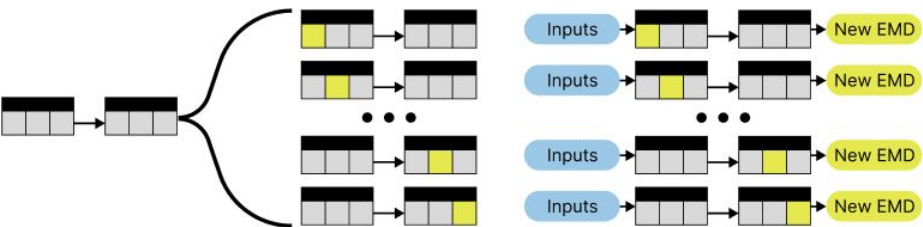

Bit-level fault injection with fine-grained control

Bit-level sensitivity metrics for ranking weight bits (without fault injection)

# ECON-T: Fault Injection Campaign (Setup)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models

# ECON-T: Fault Injection Campaign (Setup)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models

**Small Pareto** (Total Weight Bits: 10,240)

| DENSE |
|---|
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

→

| DENSE |
|---|
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

# ECON-T: Fault Injection Campaign (Setup)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models

**Small Pareto** (Total Weight Bits: 10,240)

| DENSE |
|---|
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

→

| DENSE |
|---|
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720)

| CONV-2D |
|---|
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

→

| DENSE |
|---|
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

27

# ECON-T: Fault Injection Campaign (Setup)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models

**Small Pareto** (Total Weight Bits: 10,240)

| DENSE |
|---|
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

| DENSE |
|---|
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720)

| CONV-2D |
|---|
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

| DENSE |
|---|
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

**Large Pareto** (Total Weight Bits: 61,344)

| CONV-2D |
|---|
| Kernel: 5×5×1×32 |
| Bias: 32 |
| Total Weights: 800 |
| Quantization: |
| • Total: 5 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 4,000 |

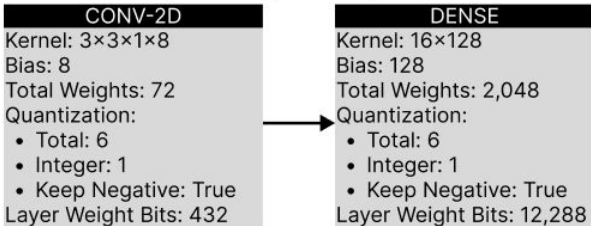| DENSE |
|---|
| Kernel: 16×512 |
| Bias: 512 |
| Total Weights: 8,192 |
| Quantization: |
| • Total: 7 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 57,344 |

# ECON-T: Fault Injection Campaign (Setup)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models
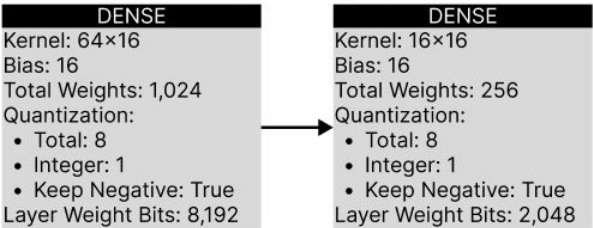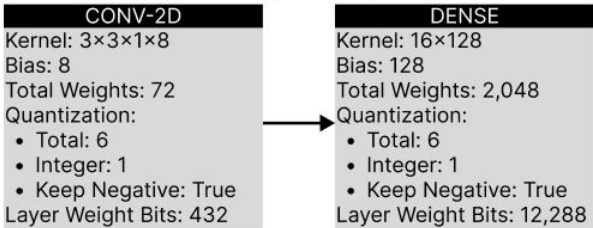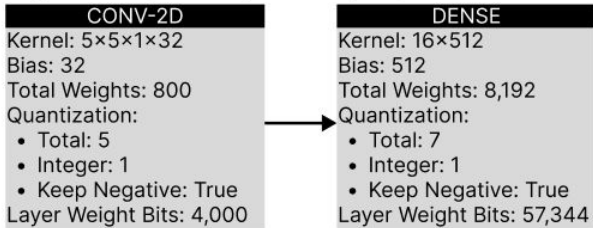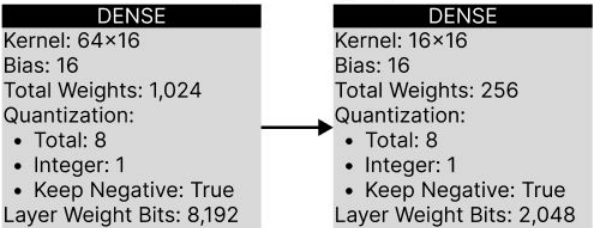
**Small Pareto** (Total Weight Bits: 10,240)

| DENSE |
|---|
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

| DENSE |
|---|
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720)

| CONV-2D |
|---|
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

| DENSE |
|---|
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

**Large Pareto** (Total Weight Bits: 61,344)

| CONV-2D |
|---|
| Kernel: 5×5×1×32 |
| Bias: 32 |
| Total Weights: 800 |
| Quantization: |
| • Total: 5 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 4,000 |

| DENSE |
|---|
| Kernel: 16×512 |
| Bias: 512 |
| Total Weights: 8,192 |
| Quantization: |
| • Total: 7 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 57,344 |

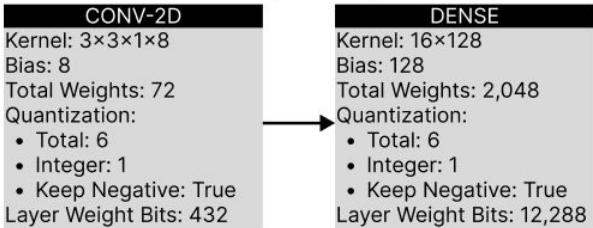- Conceptually, each fault injection campaign:

# ECON-T: Fault Injection Campaign (Setup)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models

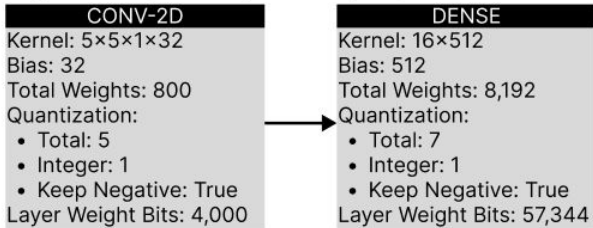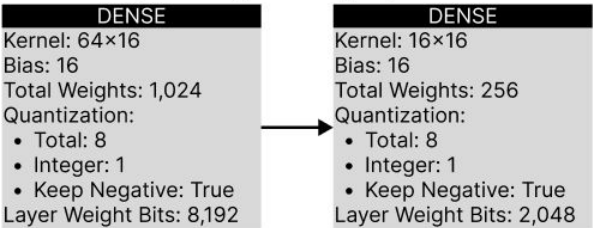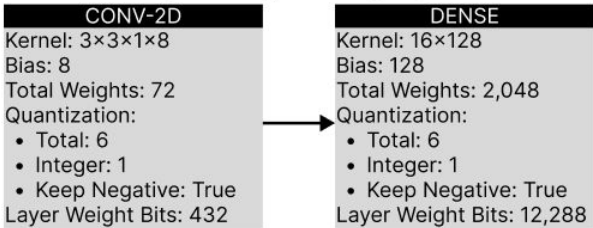**Small Pareto** (Total Weight Bits: 10,240)

| DENSE |
| --- |
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

→

| DENSE |
| --- |
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720)

| CONV-2D |
| --- |
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

→

| DENSE |
| --- |
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

**Large Pareto** (Total Weight Bits: 61,344)

| CONV-2D |
| --- |
| Kernel: 5×5×1×32 |
| Bias: 32 |
| Total Weights: 800 |
| Quantization: |
| • Total: 5 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 4,000 |

→

| DENSE |
| --- |
| Kernel: 16×512 |
| Bias: 512 |
| Total Weights: 8,192 |
| Quantization: |
| • Total: 7 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 57,344 |

- Conceptually, each fault injection campaign:



1. Generates X "faulty" variants by flipping a single weight bit

# ECON-T: Fault Injection Campaign (Setup)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models

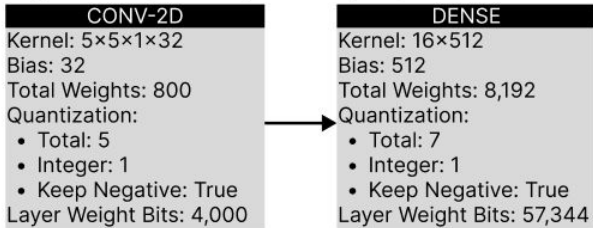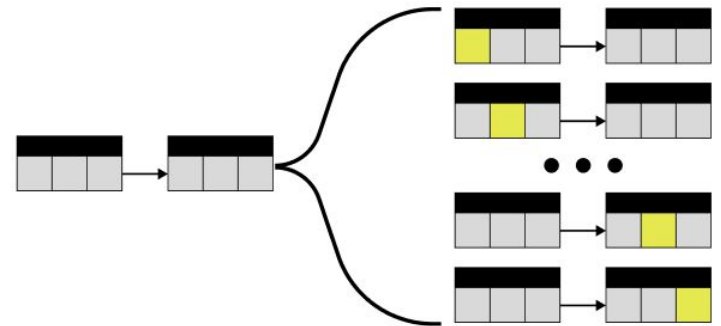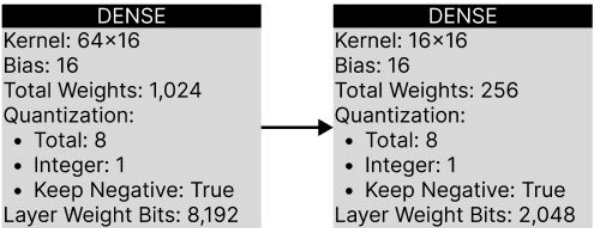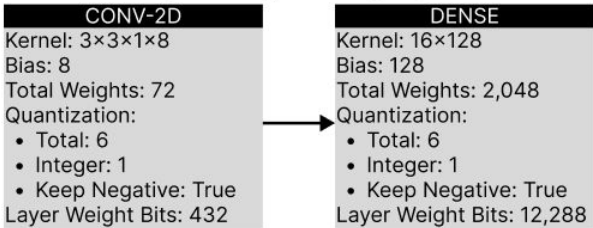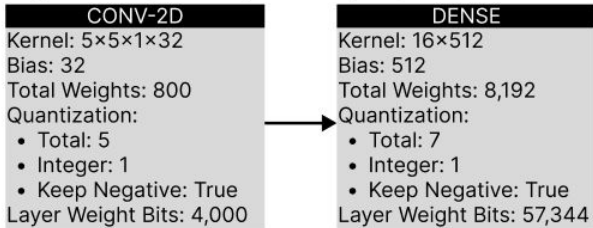**Small Pareto** (Total Weight Bits: 10,240)

| DENSE |
| --- |
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

| DENSE |
| --- |
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720)

| CONV-2D |
| --- |
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

| DENSE |
| --- |
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

**Large Pareto** (Total Weight Bits: 61,344)

| CONV-2D |
| --- |
| Kernel: 5×5×1×32 |
| Bias: 32 |
| Total Weights: 800 |
| Quantization: |
| • Total: 5 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 4,000 |

| DENSE |
| --- |
| Kernel: 16×512 |
| Bias: 512 |
| Total Weights: 8,192 |
| Quantization: |
| • Total: 7 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 57,344 |

- Conceptually, each fault injection campaign:



1. Generates X "faulty" variants by flipping a single weight bit

2. Measures the new EMD on a set of test inputs

# ECON-T: Fault Injection Campaign (Setup)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models
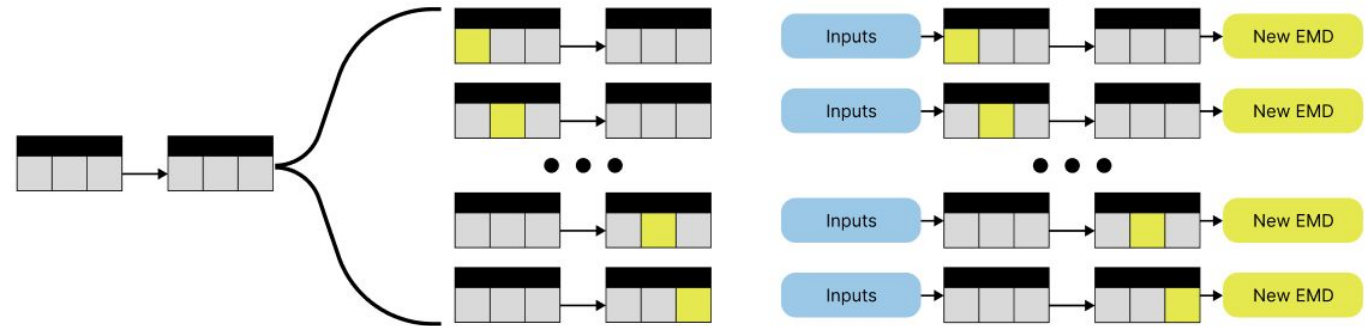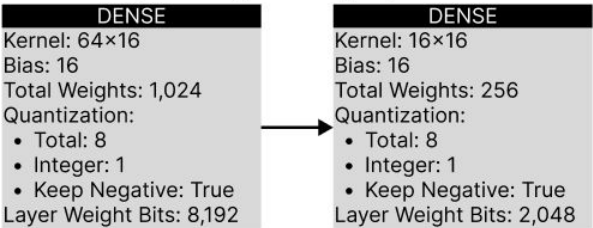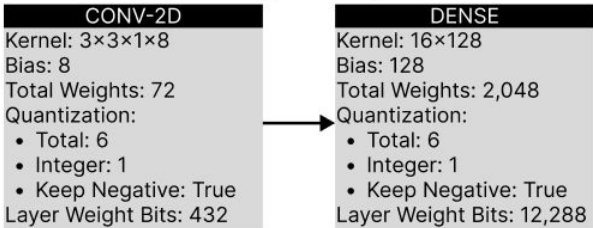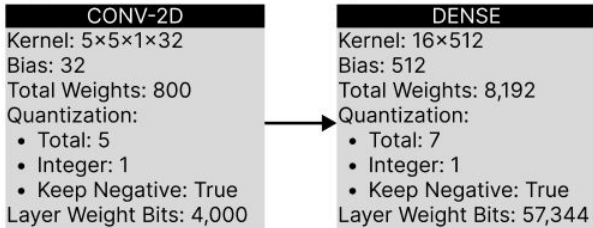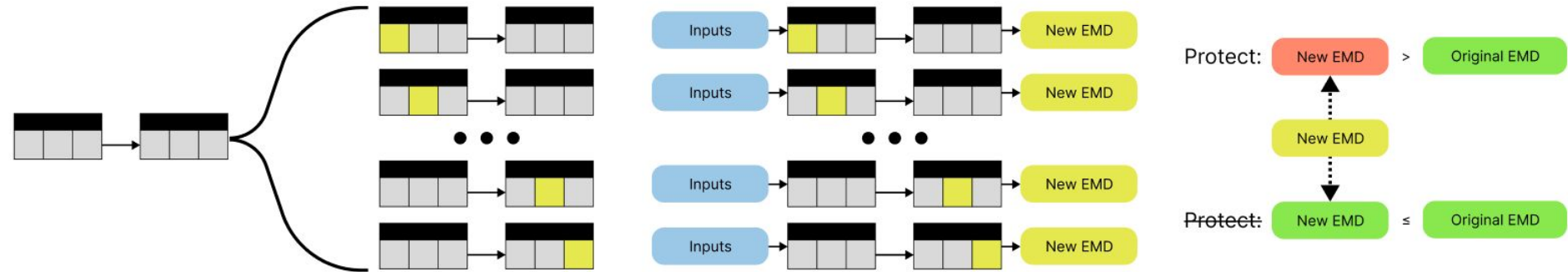
**Small Pareto** (Total Weight Bits: 10,240)

| DENSE |
|---|
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

→

| DENSE |
|---|
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720)

| CONV-2D |
|---|
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

→

| DENSE |
|---|
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

**Large Pareto** (Total Weight Bits: 61,344)

| CONV-2D |
|---|
| Kernel: 5×5×1×32 |
| Bias: 32 |
| Total Weights: 800 |
| Quantization: |
| • Total: 5 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 4,000 |

→

| DENSE |
|---|
| Kernel: 16×512 |
| Bias: 512 |
| Total Weights: 8,192 |
| Quantization: |
| • Total: 7 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 57,344 |

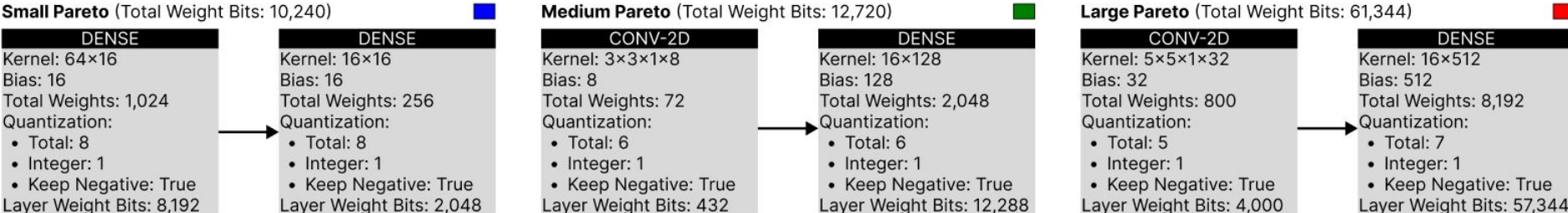- Conceptually, each fault injection campaign:



1. Generates X "faulty" variants by flipping a single weight bit

2. Measures the new EMD on a set of test inputs

3. Determines the weight bits to protect

# ECON-T: Fault Injection Campaign (Results)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models
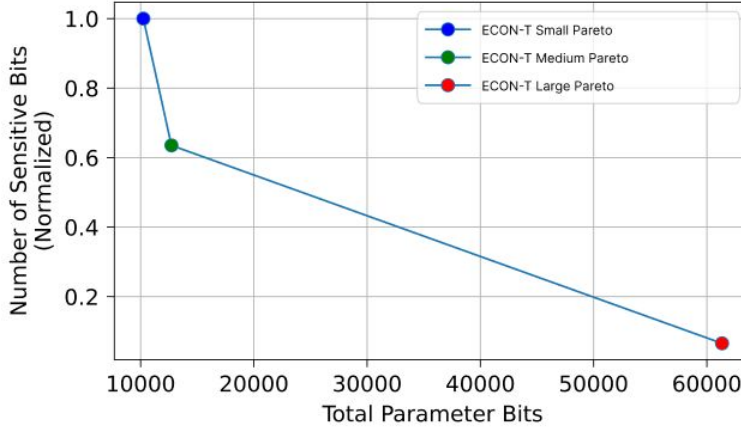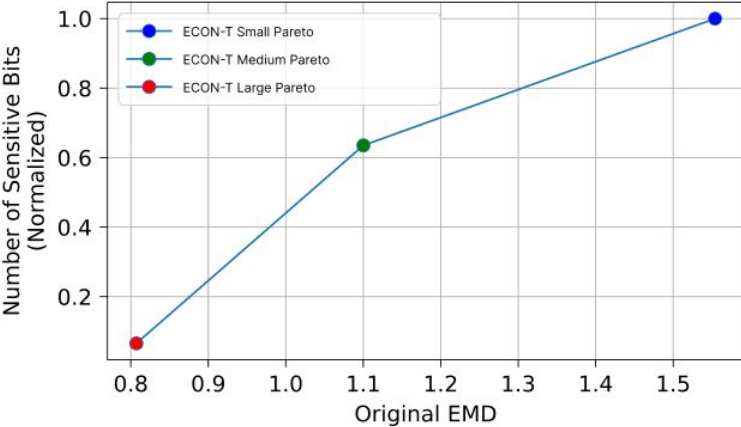
**Small Pareto** (Total Weight Bits: 10,240) ▪

| DENSE |
|---|
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

→

| DENSE |
|---|
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720) ▪

| CONV-2D |
|---|
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

→

| DENSE |
|---|
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

**Large Pareto** (Total Weight Bits: 61,344) ▪

| CONV-2D |
|---|
| Kernel: 5×5×1×32 |
| Bias: 32 |
| Total Weights: 800 |
| Quantization: |
| • Total: 5 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 4,000 |

→

| DENSE |
|---|
| Kernel: 16×512 |
| Bias: 512 |
| Total Weights: 8,192 |
| Quantization: |
| • Total: 7 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 57,344 |

# ECON-T: Fault Injection Campaign (Results)

- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models

**Small Pareto** (Total Weight Bits: 10,240) 🟦

| DENSE |
|---|
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

→

| DENSE |
|---|
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720) 🟩

| CONV-2D |
|---|
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

→

| DENSE |
|---|
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

**Large Pareto** (Total Weight Bits: 61,344) 🟥

| CONV-2D |
|---|
| Kernel: 5×5×1×32 |
| Bias: 32 |
| Total Weights: 800 |
| Quantization: |
| • Total: 5 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 4,000 |

→

| DENSE |
|---|
| Kernel: 16×512 |
| Bias: 512 |
| Total Weights: 8,192 |
| Quantization: |
| • Total: 7 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 57,344 |

# ECON-T: Fault Injection Campaign (Results)

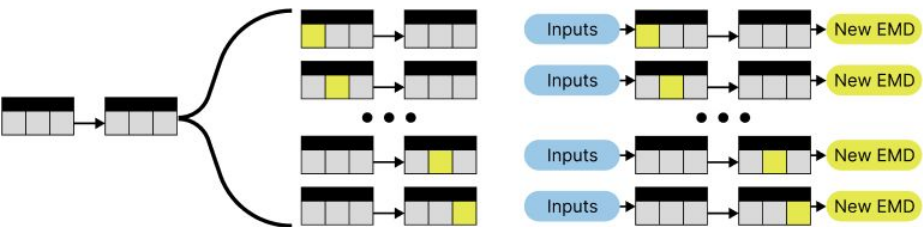- We perform an exhaustive, bit-level fault injection campaign for 3 Pareto-optimal ECON-T models

**Small Pareto** (Total Weight Bits: 10,240) ▮ (blue)

| DENSE |
|---|
| Kernel: 64×16 |
| Bias: 16 |
| Total Weights: 1,024 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 8,192 |

→

| DENSE |
|---|
| Kernel: 16×16 |
| Bias: 16 |
| Total Weights: 256 |
| Quantization: |
| • Total: 8 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 2,048 |

**Medium Pareto** (Total Weight Bits: 12,720) ▮ (green)

| CONV-2D |
|---|
| Kernel: 3×3×1×8 |
| Bias: 8 |
| Total Weights: 72 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 432 |

→

| DENSE |
|---|
| Kernel: 16×128 |
| Bias: 128 |
| Total Weights: 2,048 |
| Quantization: |
| • Total: 6 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 12,288 |

**Large Pareto** (Total Weight Bits: 61,344) ▮ (red)

| CONV-2D |
|---|
| Kernel: 5×5×1×32 |
| Bias: 32 |
| Total Weights: 800 |
| Quantization: |
| • Total: 5 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 4,000 |

→

| DENSE |
|---|
| Kernel: 16×512 |
| Bias: 512 |
| Total Weights: 8,192 |
| Quantization: |
| • Total: 7 |
| • Integer: 1 |
| • Keep Negative: True |
| Layer Weight Bits: 57,344 |



35

Fault injection campaigns are expensive...

# Fault injection campaigns are expensive...

Can we **quantify** fault sensitivity a priori?

# FKeras

- A library that assesses the fault sensitivity of (Q)Keras models

- Current features include:

  - Bit-level fault injection with fine-grained control
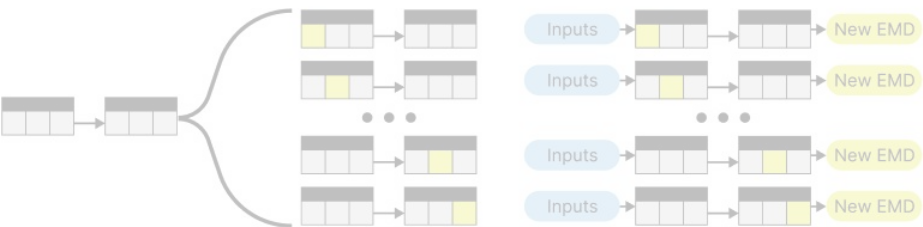
  - Bit-level sensitivity metrics for ranking weight bits

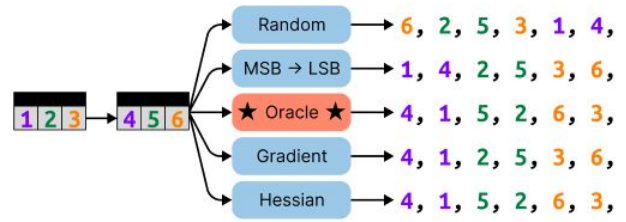

Bit-level fault injection with fine-grained control

Bit-level sensitivity metrics for ranking weight bits (without fault injection)

# FKeras

- A library that assesses the fault sensitivity of (Q)Keras models

- Current features include:
  - Bit-level fault injection with fine-grained control
  - Bit-level sensitivity metrics for ranking weight bits



Bit-level fault injection with fine-grained control



Bit-level sensitivity metrics for ranking weight bits (without fault injection)
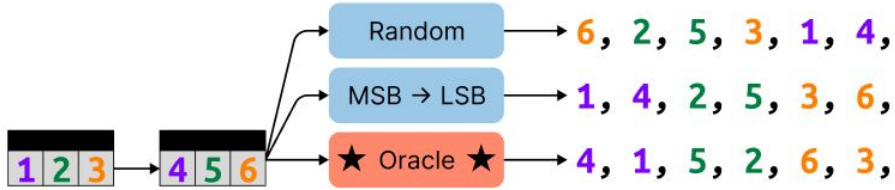
# Bit-level Sensitivity Metrics

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

# Bit-level Sensitivity Metrics

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity

  - High sensitivity: New EMD >> Original EMD

  - Low sensitivity: New EMD ≤ Original EMD

- <u>Example:</u> Assume we have a model with two layers which each have one 3-bit weight



Original Model

# Bit-level Sensitivity Metrics

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity

  - High sensitivity: New EMD >> Original EMD

  - Low sensitivity: New EMD ≤ Original EMD

- Example: Assume we have a model with two layers which each have one 3-bit weight

- Ranking Metric:

  - Oracle (requires fault injection)



Original Model          Ranking Metric          High → Low Ranking

42

# Bit-level Sensitivity Metrics

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- Example: Assume we have a model with two layers which each have one 3-bit weight

- Ranking Metric:
  - Oracle (requires fault injection)
  - Random



Original Model          Ranking Metric          High → Low Ranking

43

# Bit-level Sensitivity Metrics

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- <u>Example:</u> Assume we have a model with two layers which each have one 3-bit weight

- Ranking Metric:
  - Oracle (requires fault injection)
  - Random
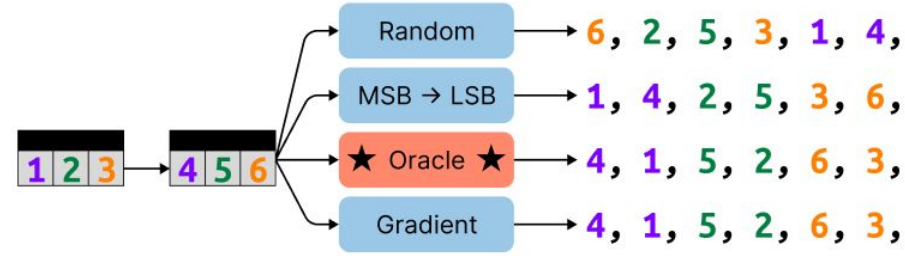  - Most significant bit → least significant bit



Original Model        Ranking Metric        High → Low Ranking

44

# Bit-level Sensitivity Metrics

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- Example: Assume we have a model with two layers which each have one 3-bit weight

- Ranking Metric:
  - Oracle (requires fault injection)
  - Random
  - Most significant bit → least significant bit
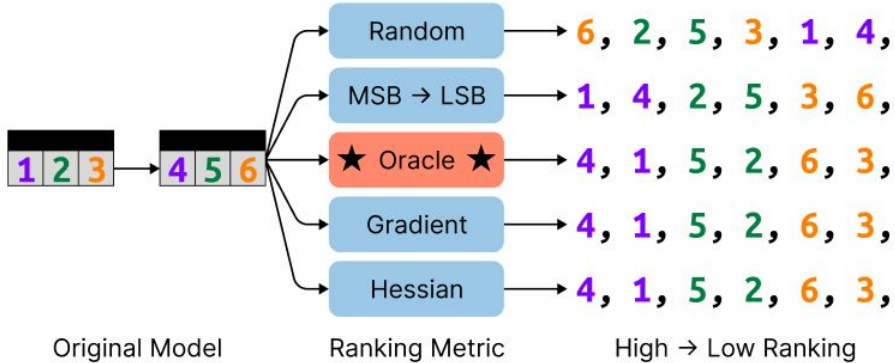  - Gradient (computed at weight level)



| | |
|---|---|
| Random | **6**, **2**, **5**, **3**, **1**, **4**, |
| MSB → LSB | **1**, **4**, **2**, **5**, **3**, **6**, |
| ★ Oracle ★ | **4**, **1**, **5**, **2**, **6**, **3**, |
| Gradient | **4**, **1**, **5**, **2**, **6**, **3**, |

Original Model     Ranking Metric     High → Low Ranking

# Bit-level Sensitivity Metrics

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- Example: Assume we have a model with two layers which each have one 3-bit weight

- Ranking Metric:
  - Oracle (requires fault injection)
  - Random
  - Most significant bit → least significant bit
  - Gradient (computed at weight level)
  - Hessian  (computed at weight level)



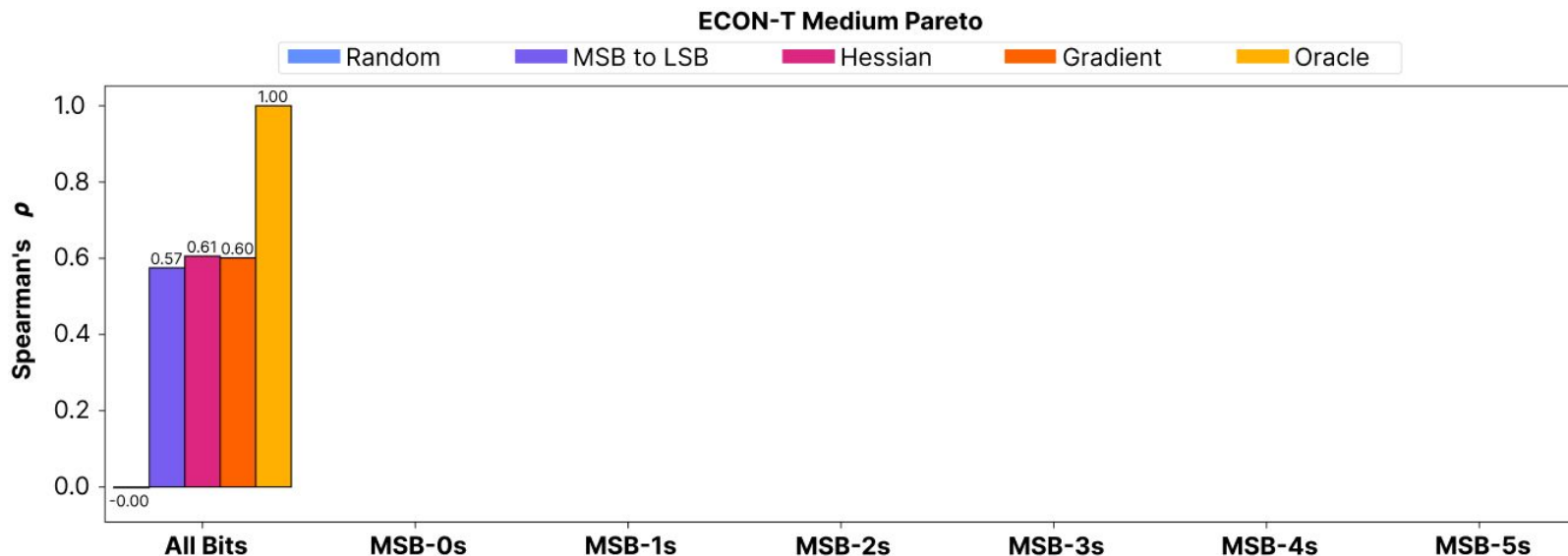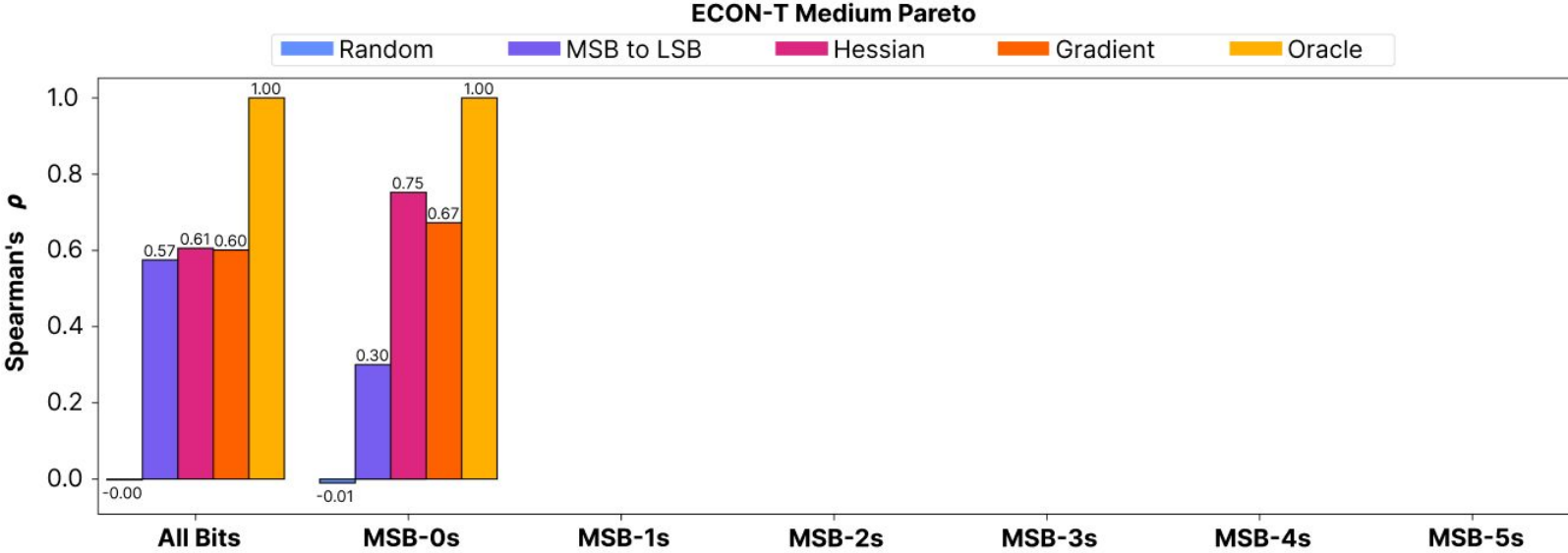| Original Model | Ranking Metric | High → Low Ranking |
|---|---|---|
| 1 2 3 → 4 5 6 | Random | 6, 2, 5, 3, 1, 4, |
|  | MSB → LSB | 1, 4, 2, 5, 3, 6, |
|  | ★ Oracle ★ | 4, 1, 5, 2, 6, 3, |
|  | Gradient | 4, 1, 5, 2, 6, 3, |
|  | Hessian | 4, 1, 5, 2, 6, 3, |

# ECON-T: Bit-level Sensitivity Metrics (Results)

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- How do our metrics compare with a perfect but costly oracle ranking?

# ECON-T: Bit-level Sensitivity Metrics (Results)

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- How do our metrics compare with a perfect but costly oracle ranking?
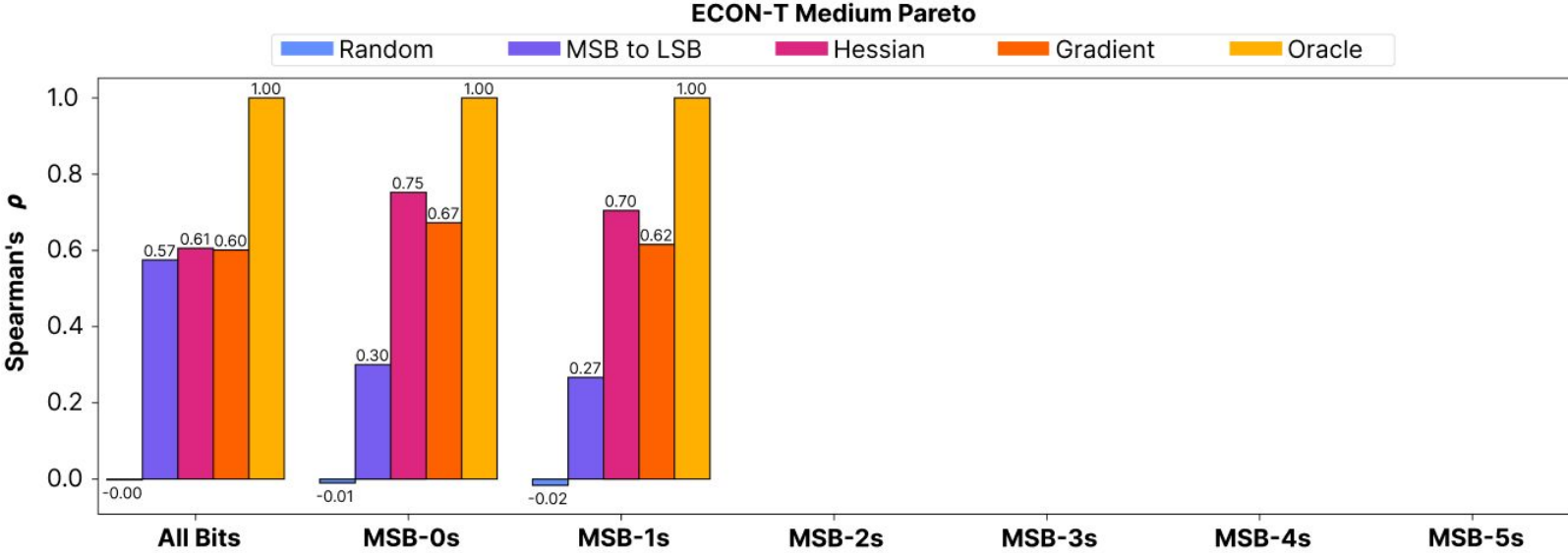


ECON-T Medium Pareto

# ECON-T: Bit-level Sensitivity Metrics (Results)

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

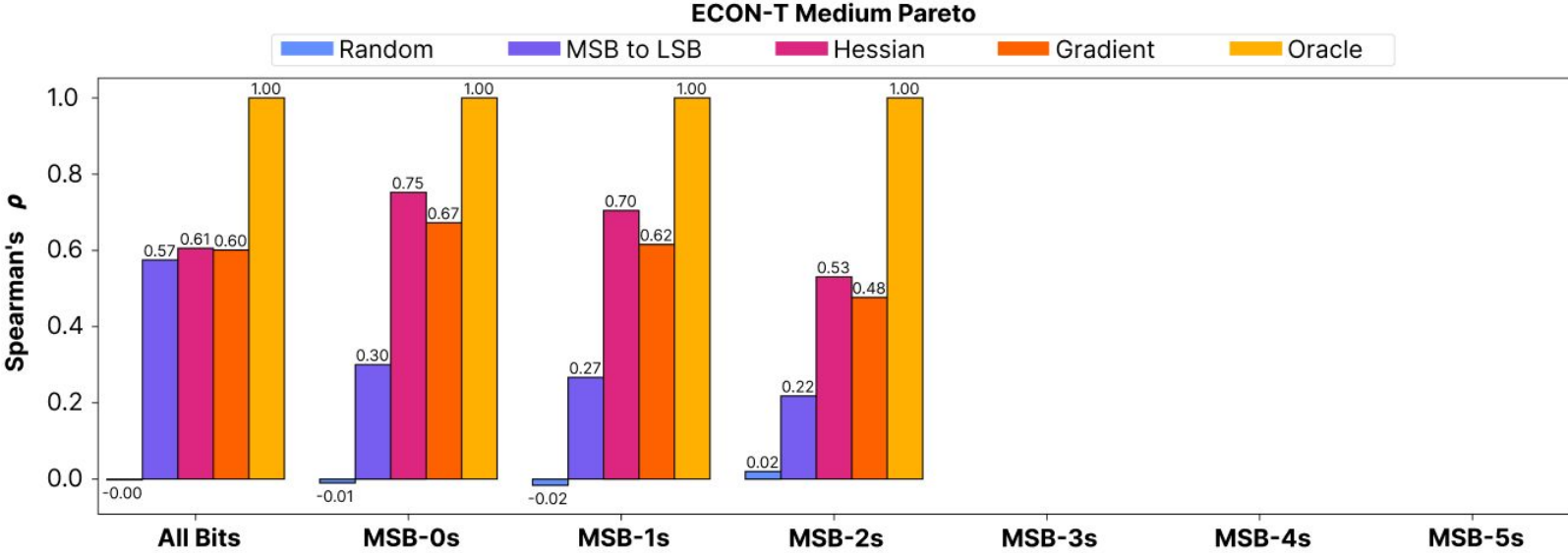- How do our metrics compare with a perfect but costly oracle ranking?

# ECON-T: Bit-level Sensitivity Metrics (Results)

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- How do our metrics compare with a perfect but costly oracle ranking?



ECON-T Medium Pareto

# ECON-T: Bit-level Sensitivity Metrics (Results)

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- How do our metrics compare with a perfect but costly oracle ranking?
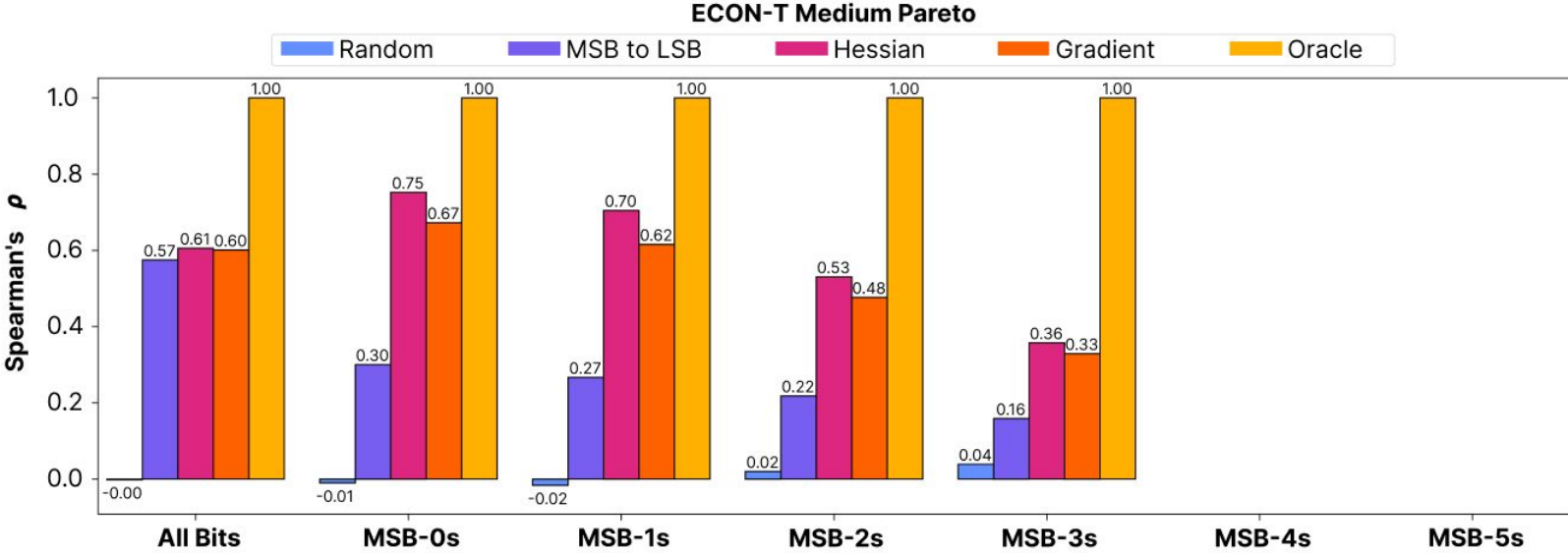


ECON-T Medium Pareto

# ECON-T: Bit-level Sensitivity Metrics (Results)

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- How do our metrics compare with a perfect but costly oracle ranking?



**ECON-T Medium Pareto**

Legend: Random, MSB to LSB, Hessian, Gradient, Oracle

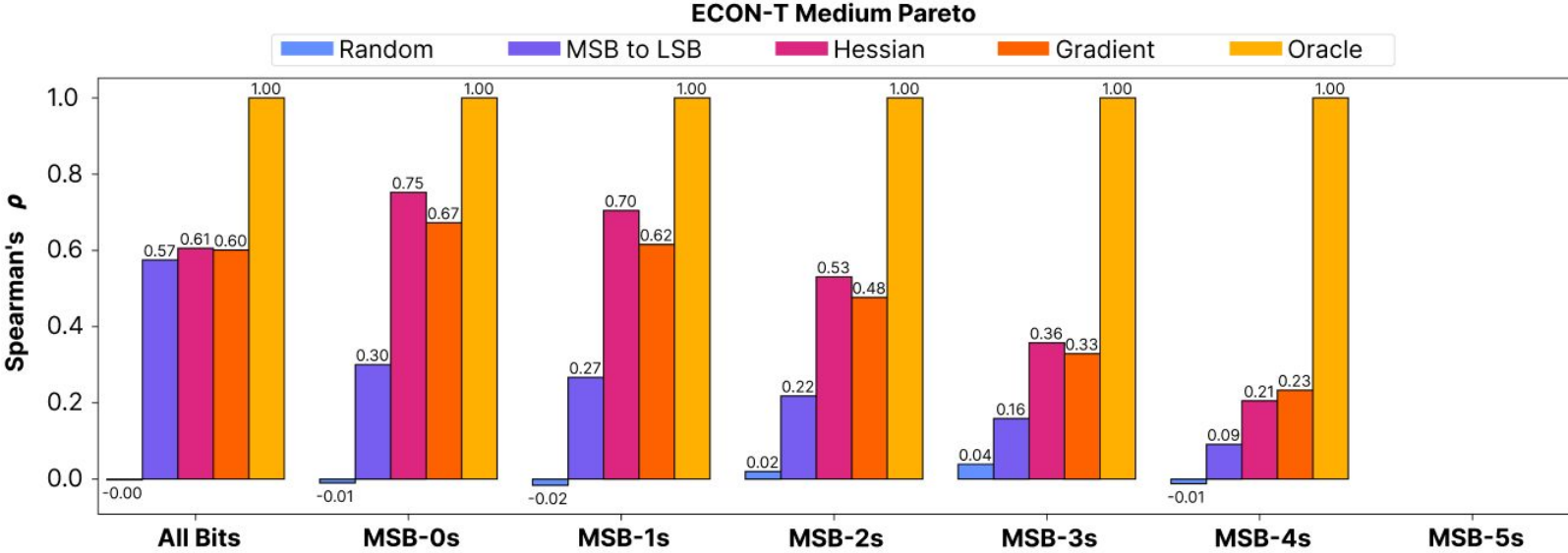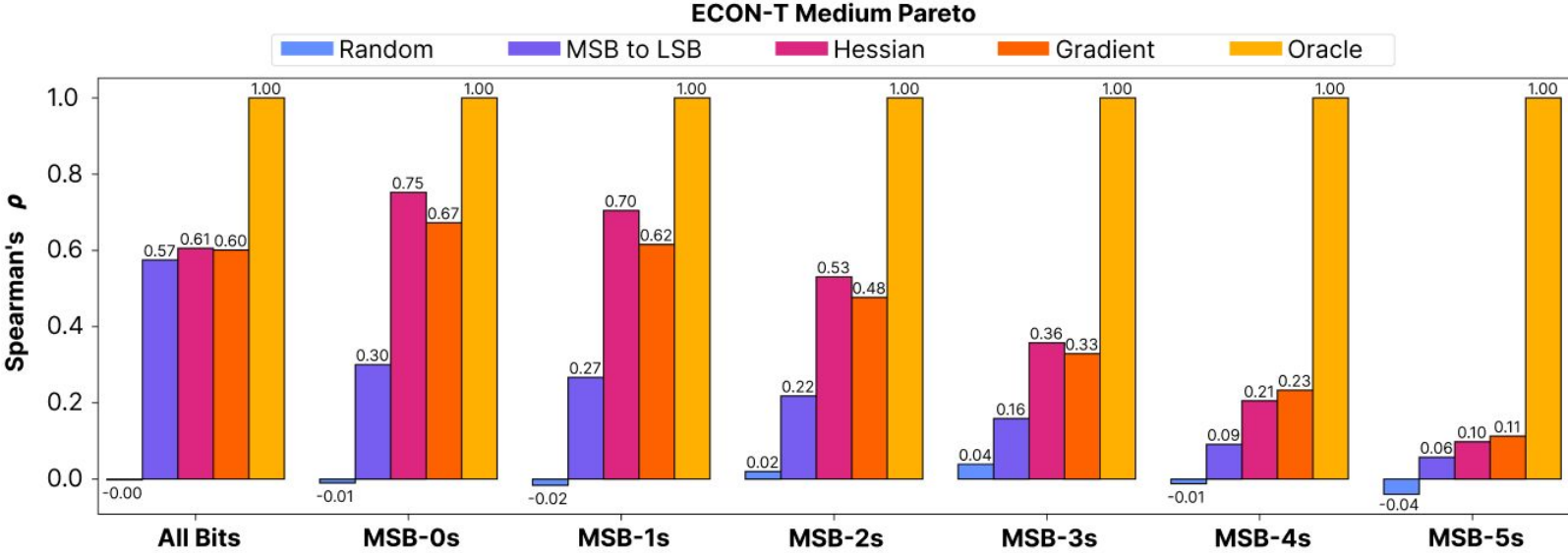| Category | Random | MSB to LSB | Hessian | Gradient | Oracle |
|---|---|---|---|---|---|
| All Bits | -0.00 | 0.57 | 0.61 | 0.60 | 1.00 |
| MSB-0s | -0.01 | 0.30 | 0.75 | 0.67 | 1.00 |
| MSB-1s | -0.02 | 0.27 | 0.70 | 0.62 | 1.00 |
| MSB-2s | 0.02 | 0.22 | 0.53 | 0.48 | 1.00 |
| MSB-3s | 0.04 | 0.16 | 0.36 | 0.33 | 1.00 |
| MSB-4s | | | | | |
| MSB-5s | | | | | |

Y-axis: Spearman's $\rho$

# ECON-T: Bit-level Sensitivity Metrics (Results)

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- How do our metrics compare with a perfect but costly oracle ranking?



ECON-T Medium Pareto

Legend: Random, MSB to LSB, Hessian, Gradient, Oracle

| | Random | MSB to LSB | Hessian | Gradient | Oracle |
|---|---|---|---|---|---|
| All Bits | -0.00 | 0.57 | 0.61 | 0.60 | 1.00 |
| MSB-0s | -0.01 | 0.30 | 0.75 | 0.67 | 1.00 |
| MSB-1s | -0.02 | 0.27 | 0.70 | 0.62 | 1.00 |
| MSB-2s | 0.02 | 0.22 | 0.53 | 0.48 | 1.00 |
| MSB-3s | 0.04 | 0.16 | 0.36 | 0.33 | 1.00 |
| MSB-4s | -0.01 | 0.09 | 0.21 | 0.23 | 1.00 |
| MSB-5s | | | | | |

Y-axis: Spearman's $\rho$

# ECON-T: Bit-level Sensitivity Metrics (Results)

- We provide bit-level sensitivity metrics for ranking weight bits from high to low sensitivity
  - High sensitivity: New EMD >> Original EMD
  - Low sensitivity: New EMD ≤ Original EMD

- How do our metrics compare with a perfect but costly oracle ranking?



ECON-T Medium Pareto

Legend: Random, MSB to LSB, Hessian, Gradient, Oracle

| | All Bits | MSB-0s | MSB-1s | MSB-2s | MSB-3s | MSB-4s | MSB-5s |
|---|---|---|---|---|---|---|---|
| Random | -0.00 | -0.01 | -0.02 | 0.02 | 0.04 | -0.01 | -0.04 |
| MSB to LSB | 0.57 | 0.30 | 0.27 | 0.22 | 0.16 | 0.09 | 0.06 |
| Hessian | 0.61 | 0.75 | 0.70 | 0.53 | 0.36 | 0.21 | 0.10 |
| Gradient | 0.60 | 0.67 | 0.62 | 0.48 | 0.33 | 0.23 | 0.11 |
| Oracle | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |

Spearman's $\rho$

# FKeras: Future Work

- We want to use FKeras to:

  - Analyze more edge NNs and datasets

    - How much can our metrics speed up fault injection campaigns?

  - Perform NN design space exploration that considers fault sensitivity using our metrics

    - How does fault sensitivity interact with performance, area, etc?

# Thank you! Questions?
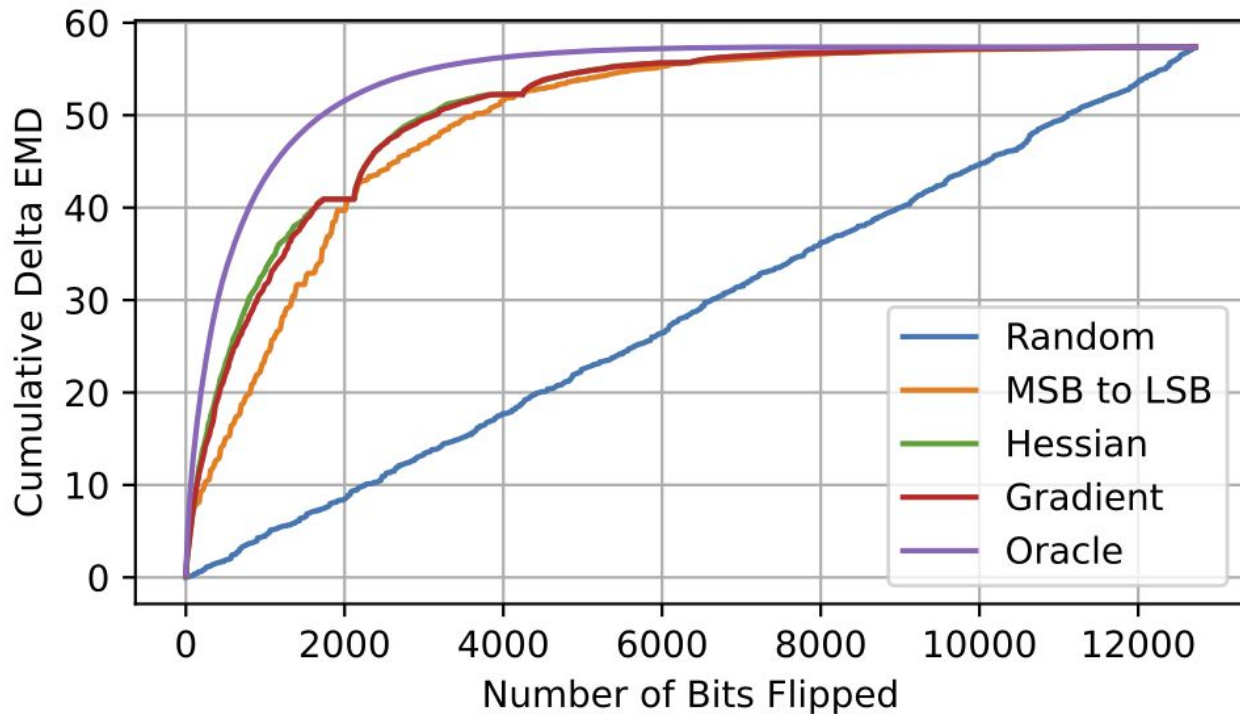
# Thank you! Questions?

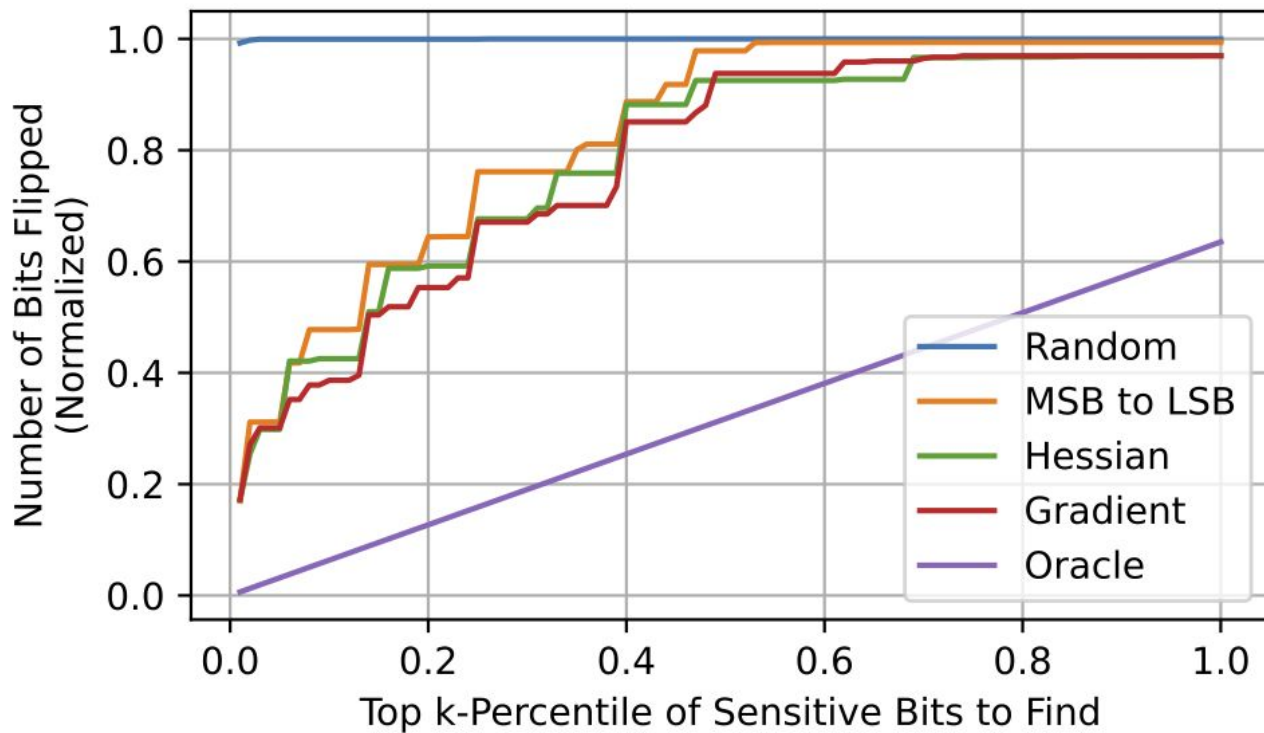**FKeras Repo: https://github.com/KastnerRG/fkeras**

# Backup

# How do our metrics compare with a perfect ranking?

# How do our metrics compare with a perfect ranking?

# How to use bit-level sensitivity rankings?

# How to use bit-level sensitivity rankings?

- Speed up fault injection campaigns

# How to use bit-level sensitivity rankings?

- Speed up fault injection campaigns

- Design space exploration