



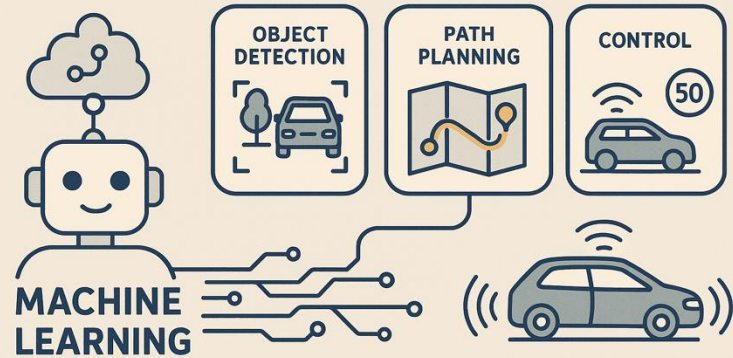
# SAVE: Software-Implemented Fault Tolerance for Model Inference against GPU Memory Bit Flips

Wenxin Zheng, Bin Xu, Jinyu Gu, Haibo Chen

*IPADS Lab, Shanghai Jiao Tong University*

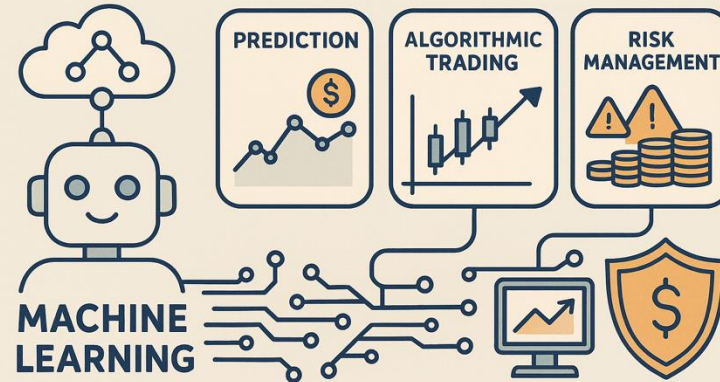
# Characteristics of Machine Learning in the Edge: Result Critical

## MACHINE LEARNING IN AUTONOMOUS DRIVING



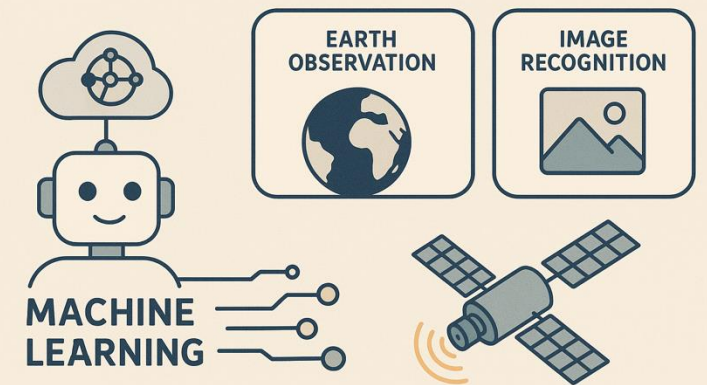
Autonomous Driving

## MACHINE LEARNING IN FINANCE



Financial System

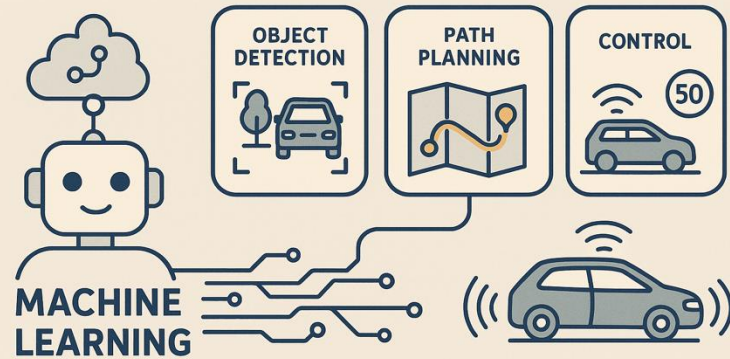
## MACHINE LEARNING IN SATELLITES



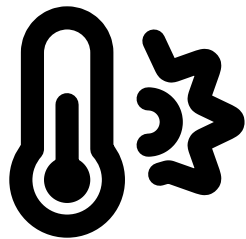
Satellite in Orbit

# Characteristics of Machine Learning in the Edge: Result Critical

## MACHINE LEARNING IN AUTONOMOUS DRIVING

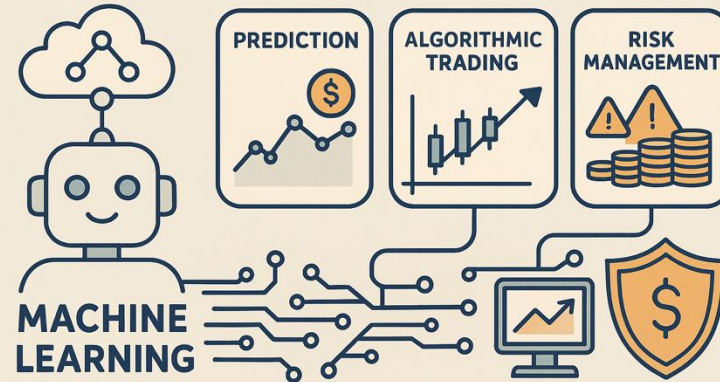


Autonomous Driving

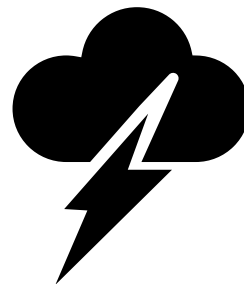


Changing Temperature

## MACHINE LEARNING IN FINANCE

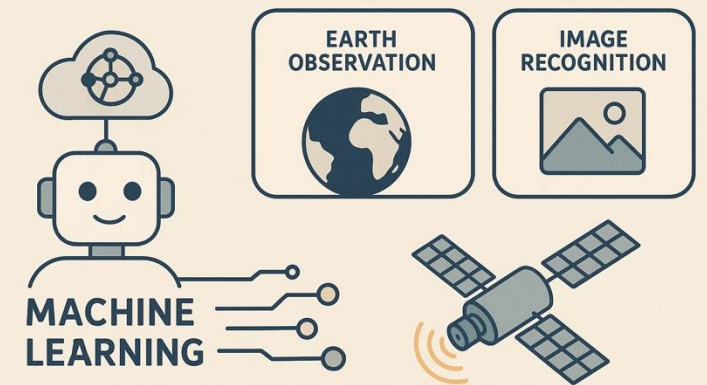


Financial System

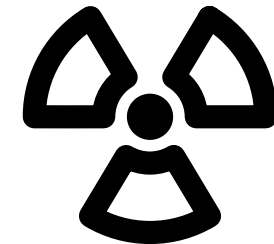


Unstable Voltage

## MACHINE LEARNING IN SATELLITES

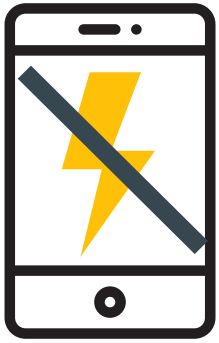


Satellite in Orbit

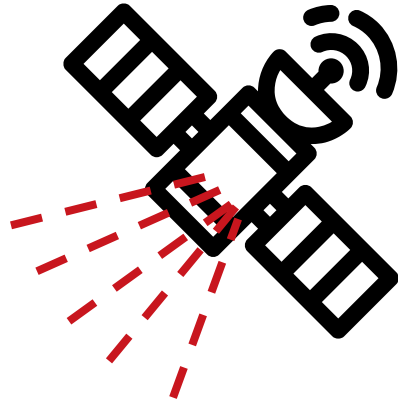


Radiation

# Characteristics of Machine Learning in the Edge: Bit Sensitive



**3,900** bit flips/day



**16 million** bit flips/day

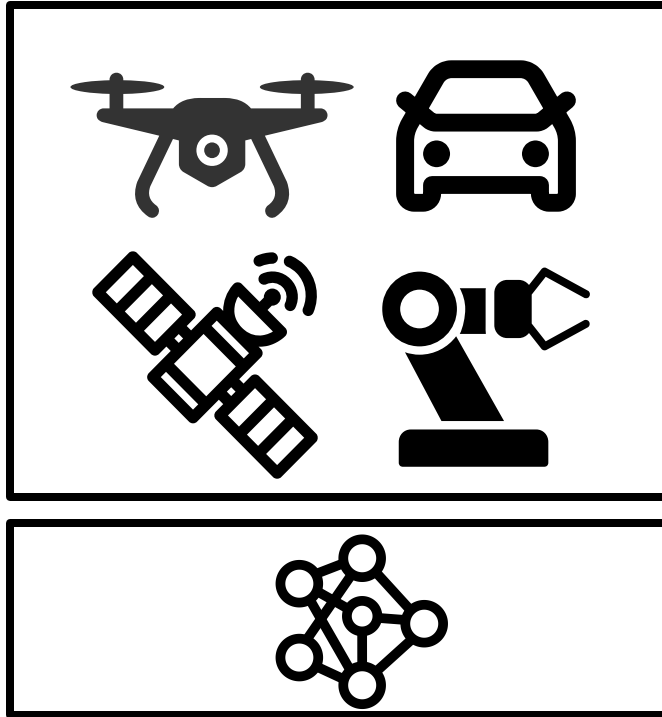
Memory is **not** always reliable

## Model degradation due to bit flips

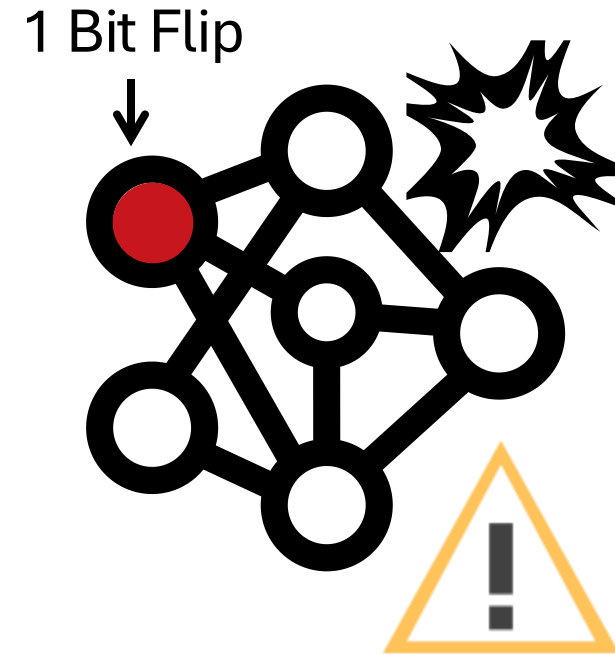
	Model	Flipped Bits	Accuracy Degradation
NaN [29]	EfficientNet	3	82.00% → 0%
FIASM [46]	Softmax <sup>(1)</sup>	1	97.4% → 62%
TBD [40]	CV <sup>(2)</sup>	1	up to 99% drop
BFA [73]	ResNet18 <sup>(3)</sup>	13	69.8% → 0.1%
DHB [37]	ResNet20 <sup>(3)</sup>	28	81.39% drop
TBT [74]	ResNet18 <sup>(3)</sup>	84	92% ASR <sup>(4)</sup>
TA-LBF [9]	ResNet18 <sup>(3)</sup>	1507	100% ASR <sup>(4)</sup>
T-BFA [75]	ResNet18 <sup>(3)</sup>	27	100% ASR <sup>(4)</sup>
Our Experiments	ResNet	1	99% Drop
	ViT		
	CogACT		

Even **1 Bit flip** can disable whole model

# Protecting Machine Learning in the Edge

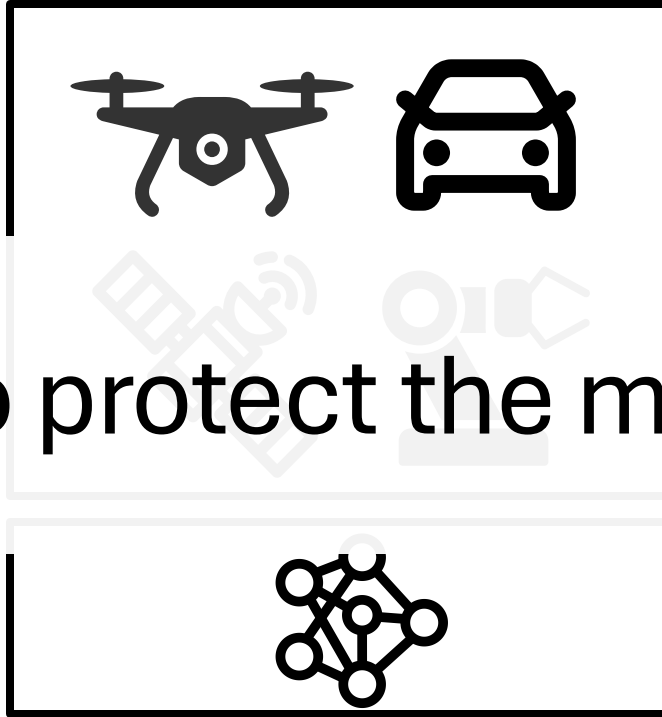


Result Critical

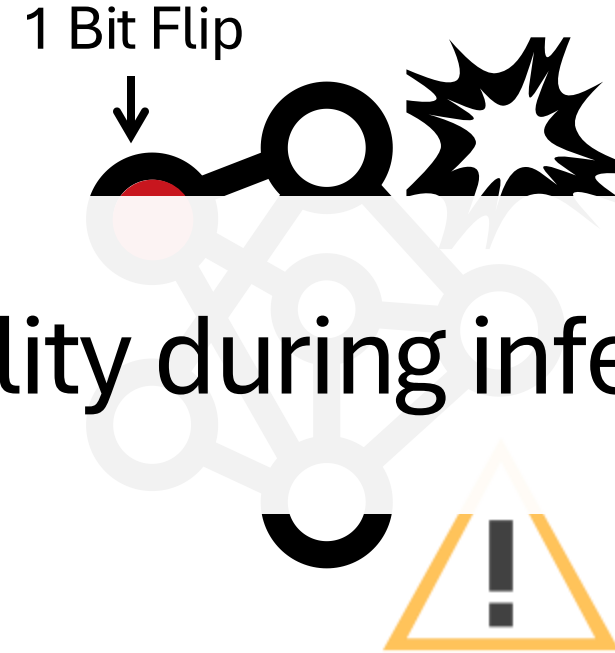


Bit Sensitive

# Protecting Machine Learning in the Edge



Result Critical



Bit Sensitive

How to protect the model reliability during inference?

# Current Solution

## Hardware Solution

- TMR
- ECC
- Hardware Reliable Only

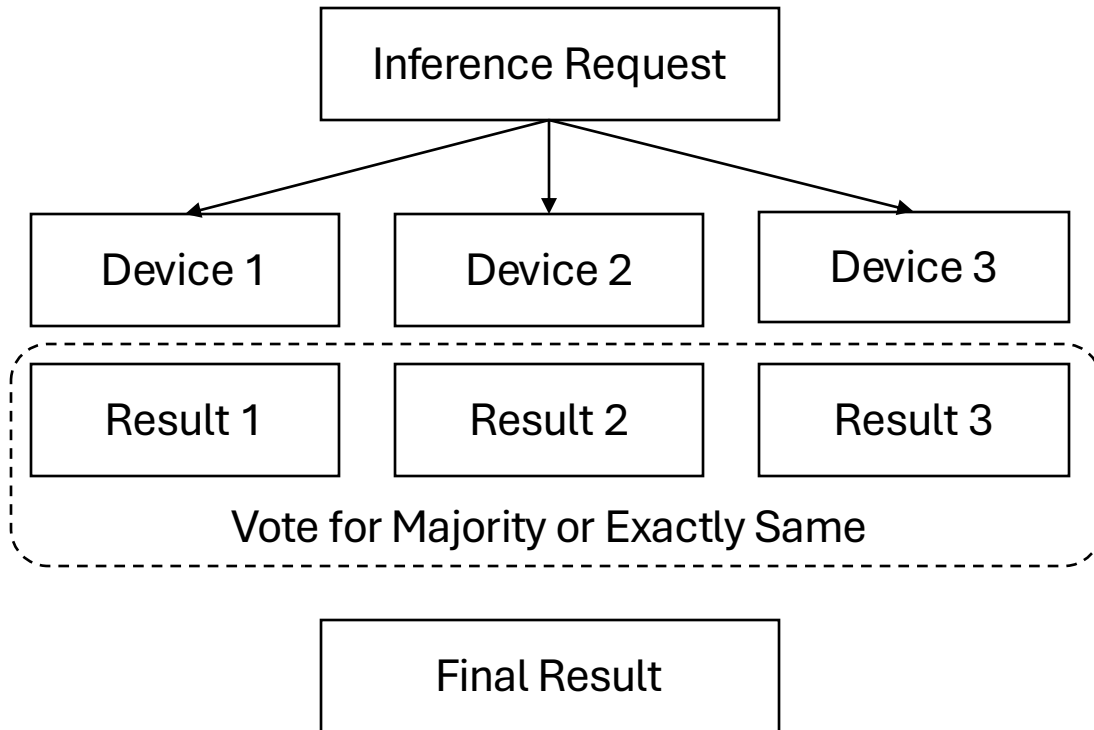
## Software Solution

- Check Distribution of Activations
- Robust Activation Function

# Hardware Solution

## TMR Solution

Triple Modular Redundancy

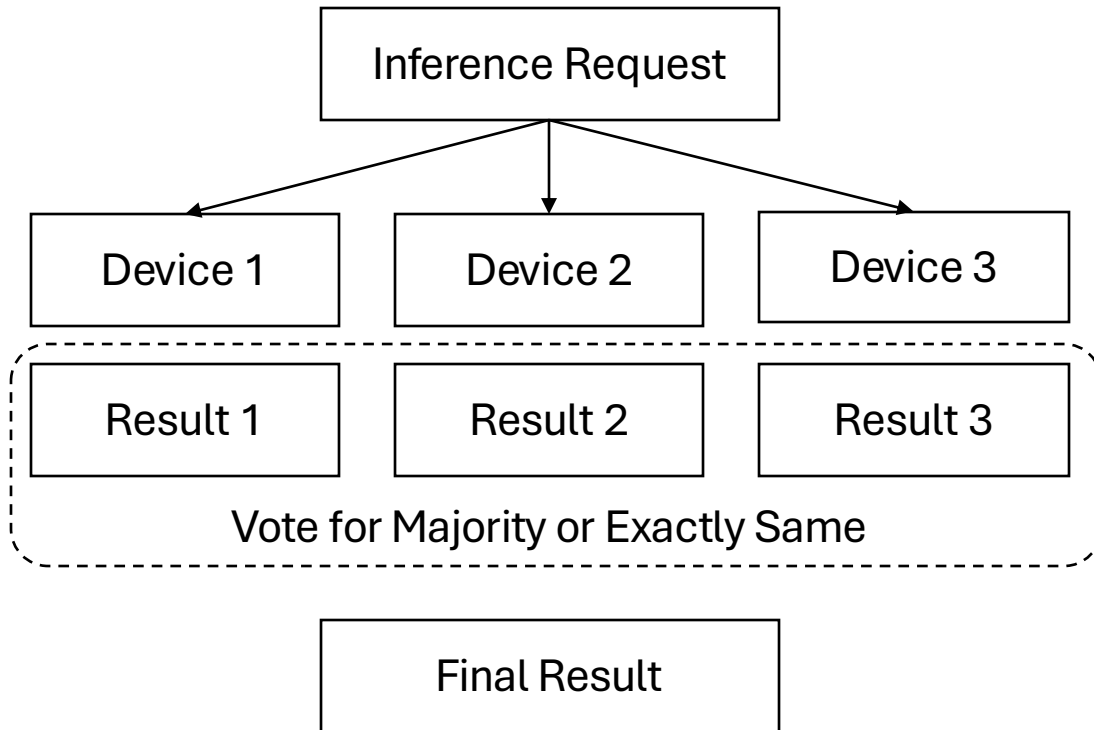


Require **3x** of computation resources

# Hardware Solution

## TMR Solution

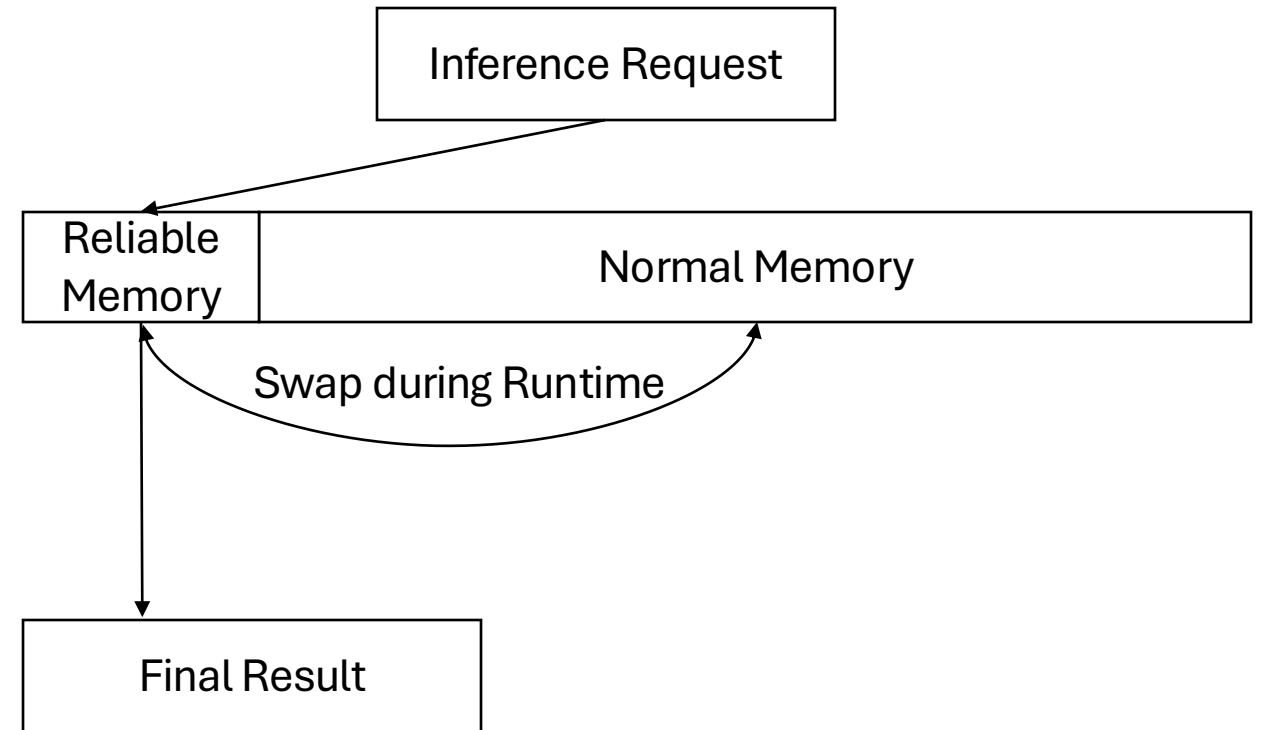
Triple Modular Redundancy



Require **3x** of computation resources

## Hardware Reliable Memory Solution

Use reliable memory for all calculation

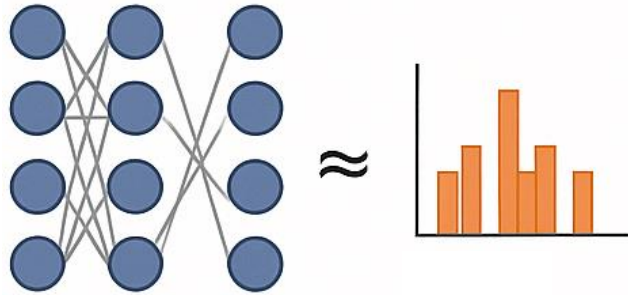


Additional **1000x** runtime overhead or infeasible

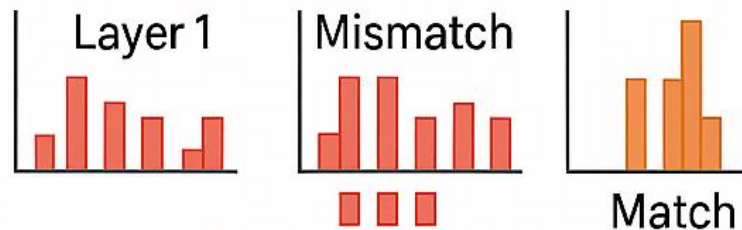
# Software Solution

Dr.DNA Solution [ASPLOS'24]

Patch the output with profiled distribution



Step 1: Profile the Distribution



Step 2: Runtime Identification and patching

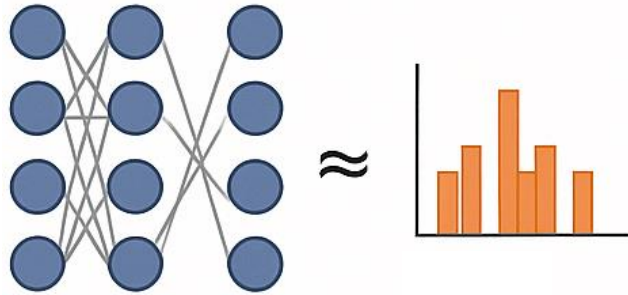
**Accuracy Affected**

Cannot handle Out-of-Distribution Data

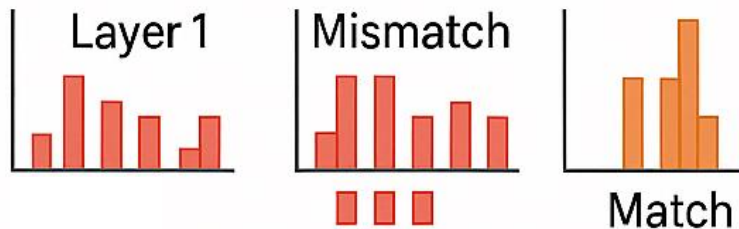
# Software Solution

Dr.DNA Solution [ASPLOS'24]

Patch the output with profiled distribution



Step 1: Profile the Distribution



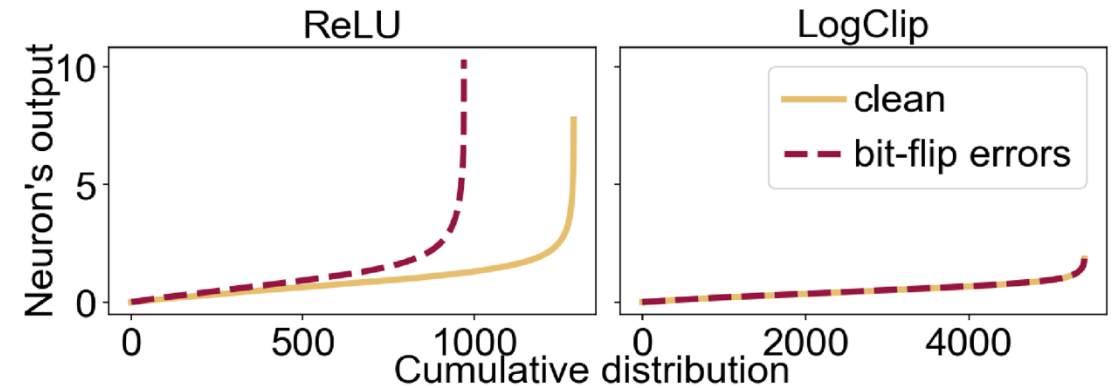
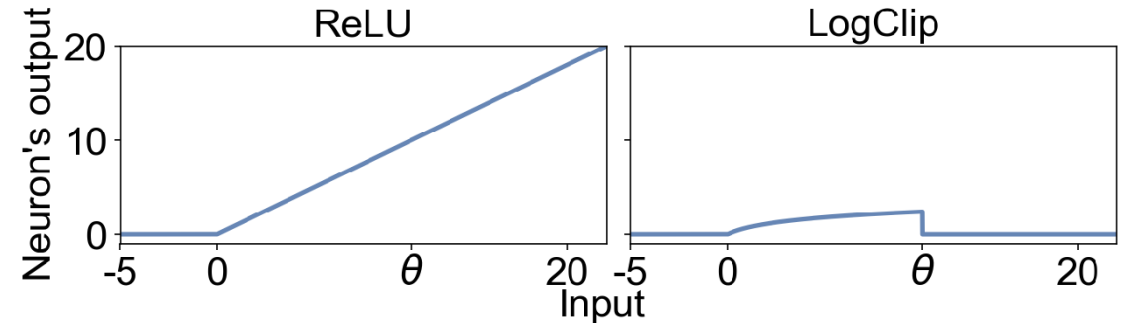
Step 2: Runtime Identification and patching

**Accuracy Affected**

Cannot handle Out-of-Distribution Data

RedNet Solution [2]

Change the Activation function into Robust one



**Retrain Required**

Cannot handle activation function other than ReLU

[1] Ma, Dongning, et al. "Dr. DNA: Combating silent data corruptions in deep learning using distribution of neuron activations." *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*. 2024

[2] Wang, Meiqi, et al. "A case for application-aware space radiation tolerance in orbital computing." arXiv preprint arXiv:2407.11853 (2024)



# Our Key Insights



# Our Key Insights

Insight-1: Not all hardware bits are created equal.

# Our Key Insights

Insight-1: Not all hardware bits are created equal.

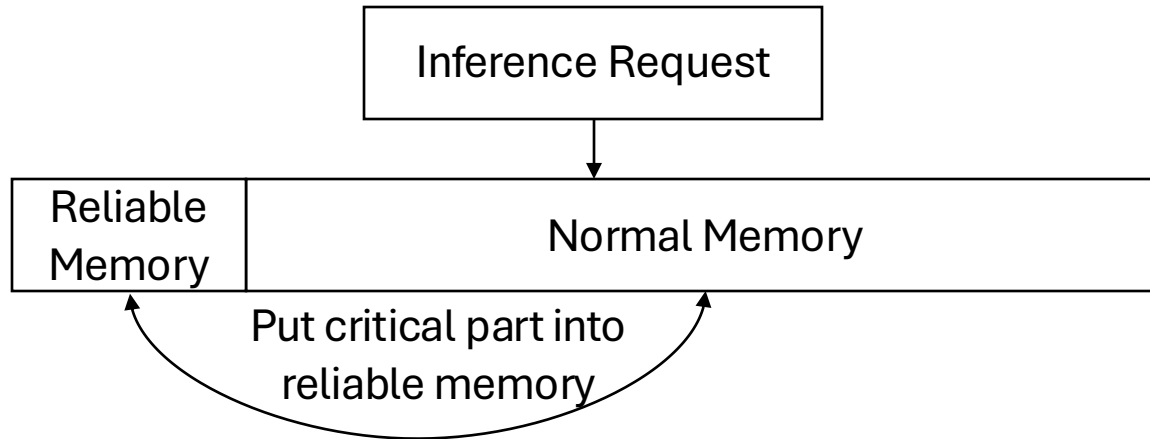
Reliable Memory	Normal Memory
--------------------	---------------

Reliable Memory:

- Soft Error Detection in SRAM
- A small portion of DRAM

# Our Key Insights

Insight-1: Not all hardware bits are created equal.

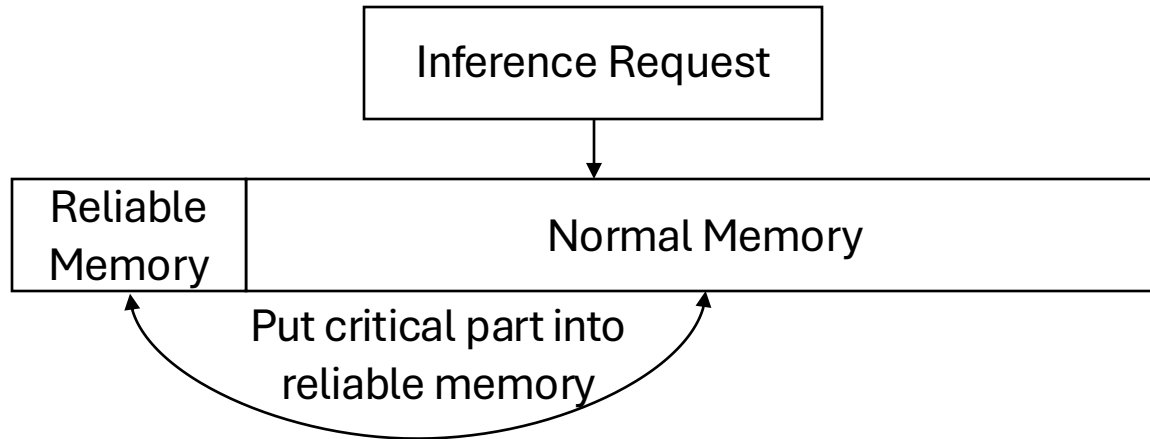


Reliable Memory:

- Soft Error Detection in SRAM
- A small portion of DRAM

# Our Key Insights

Insight-1: Not all hardware bits are created equal.



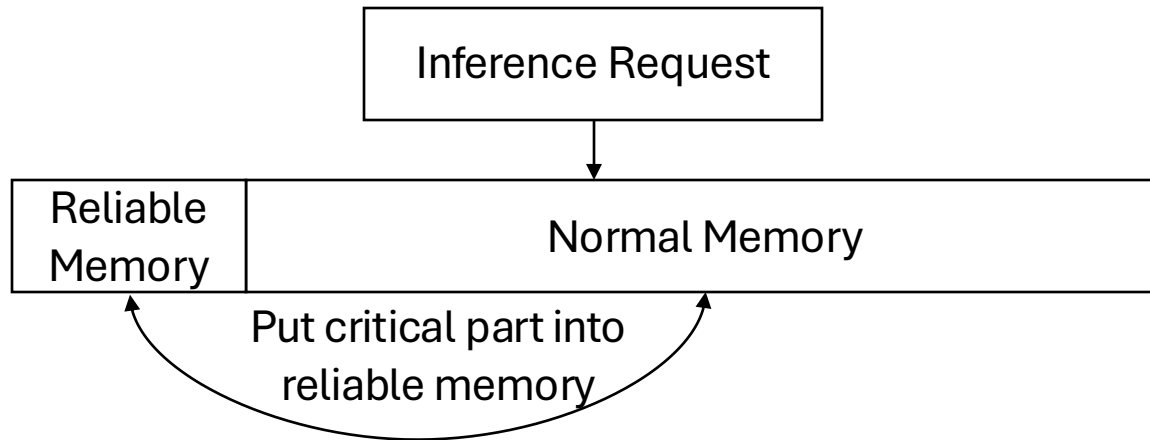
Reliable Memory:

- Soft Error Detection in SRAM
- A small portion of DRAM

Insight-2: Not all bit flips in software are fatal or silent.

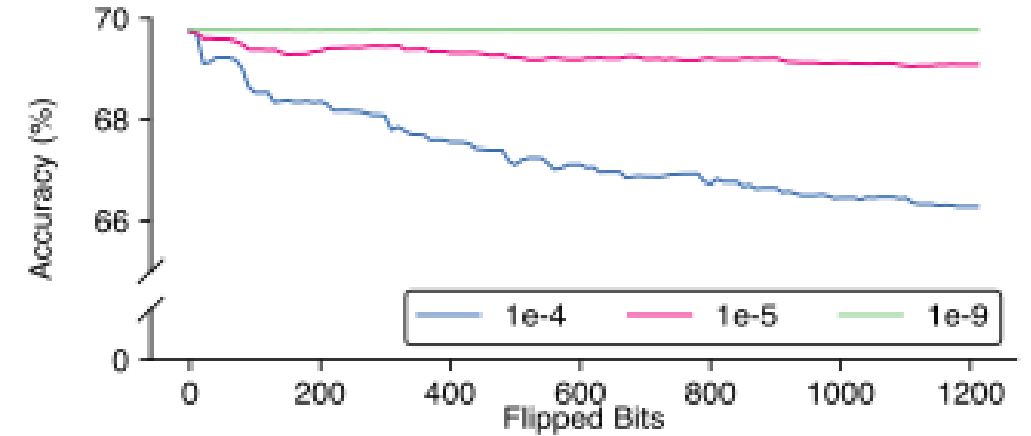
# Our Key Insights

Insight-1: Not all hardware bits are created equal.



Reliable Memory:

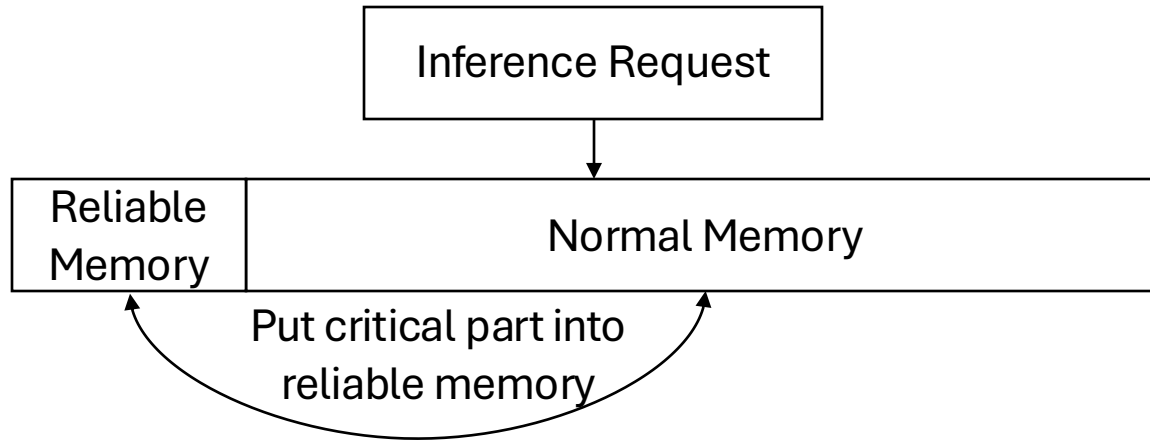
- Soft Error Detection in SRAM
- A small portion of DRAM



Insight-2: Not all bit flips in software are fatal or silent.

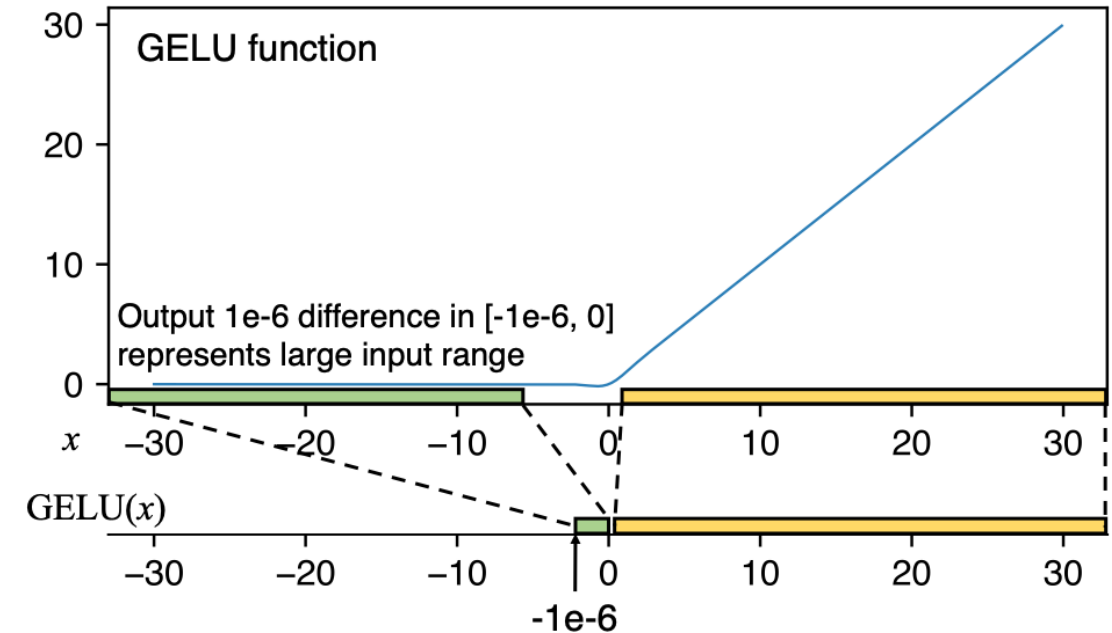
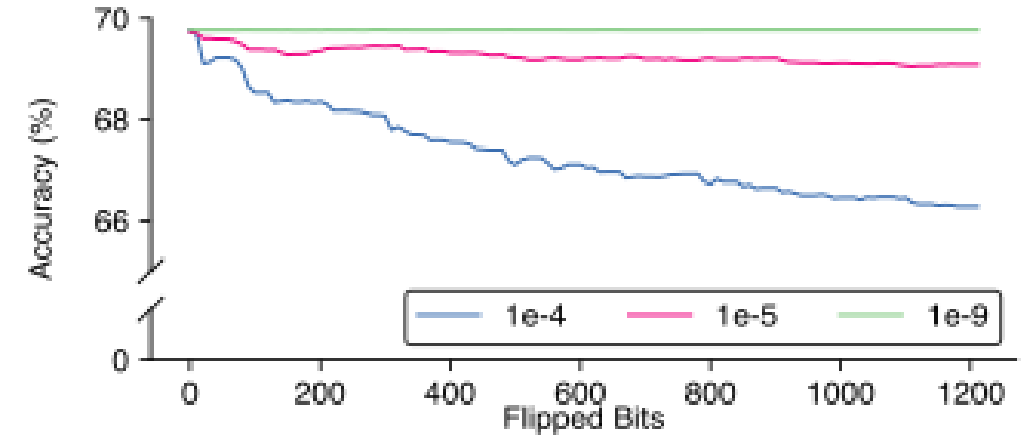
# Our Key Insights

Insight-1: Not all hardware bits are created equal.



Reliable Memory:

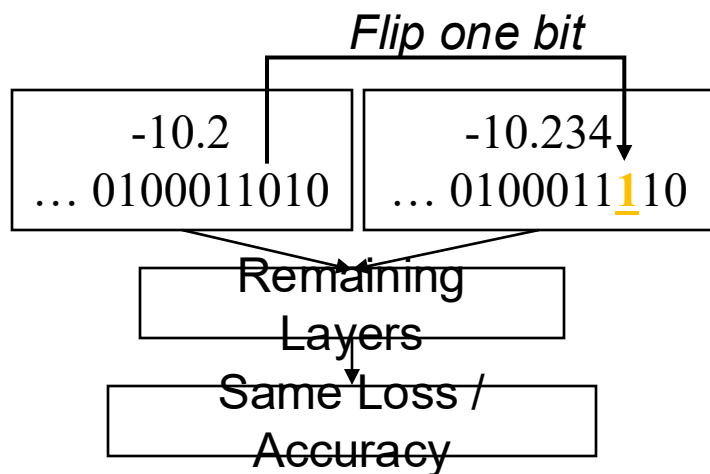
- Soft Error Detection in SRAM
- A small portion of DRAM



Insight-2: Not all bit flips in software are fatal or silent.

# Using Bit Robustness

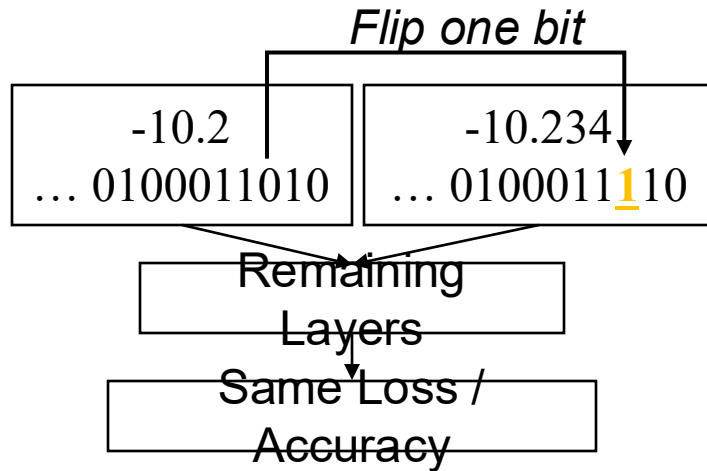
## Robust Bits



**Ignore**

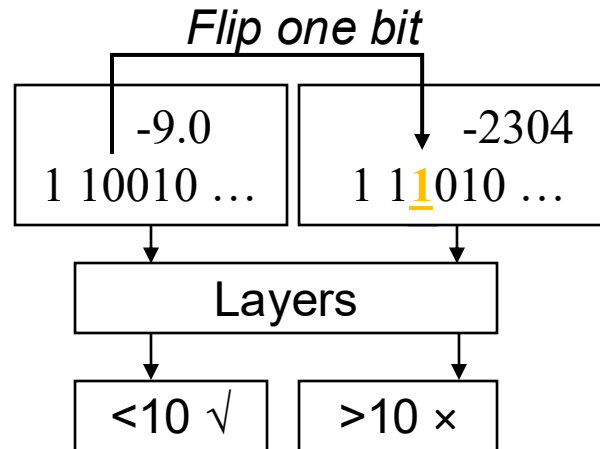
# Using Bit Robustness

## Robust Bits



**Ignore**

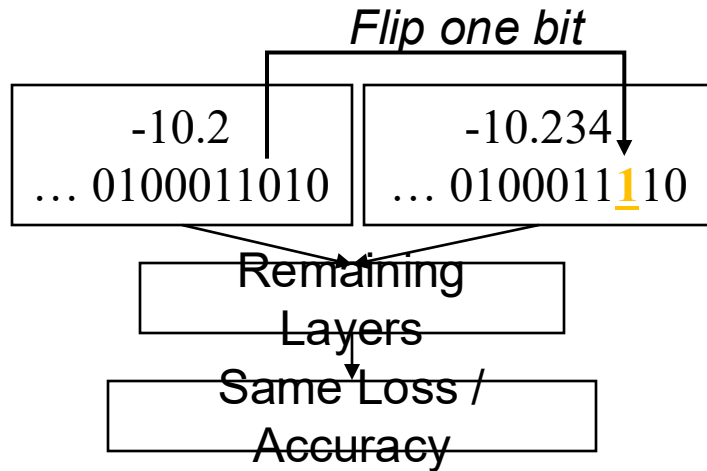
## Ranging Bits



**Range Check**

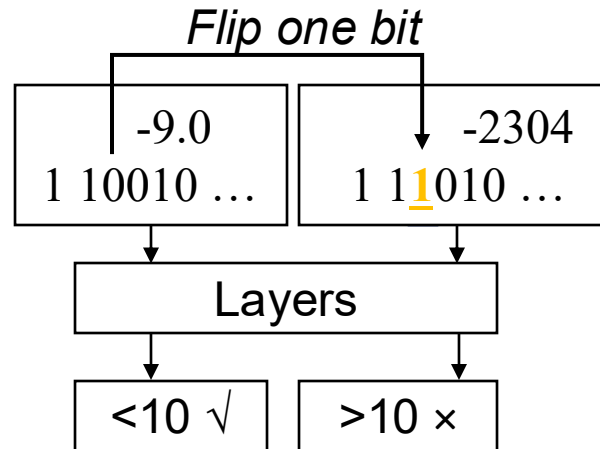
# Using Bit Robustness

## Robust Bits



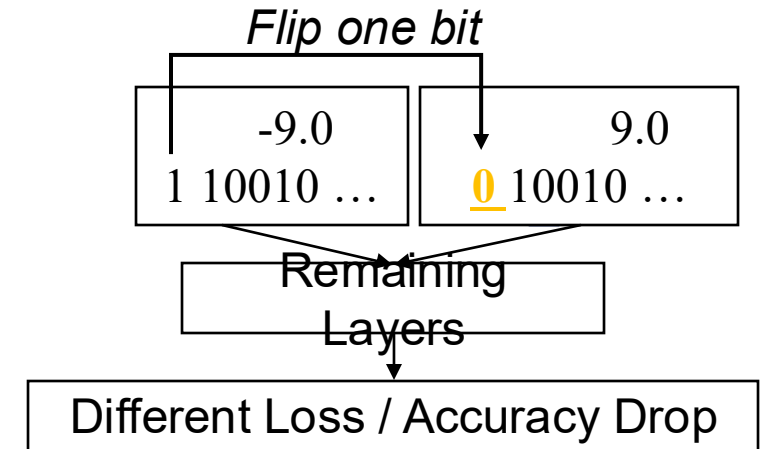
**Ignore**

## Ranging Bits



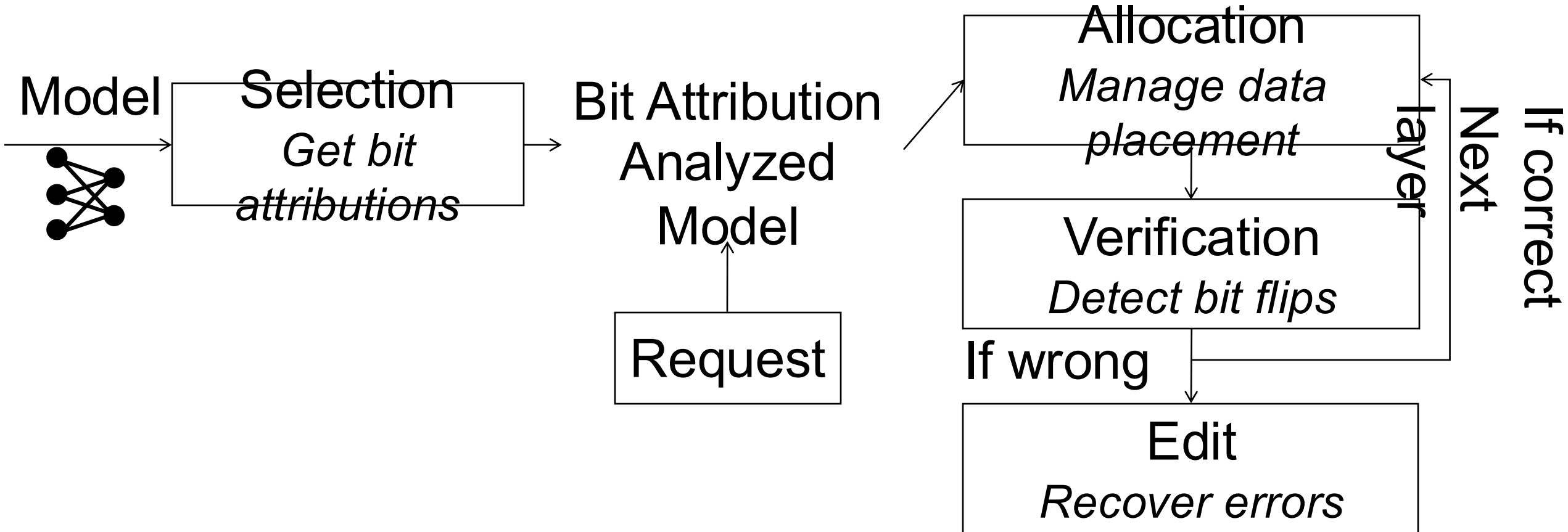
**Range Check**

## Vulnerable Bits

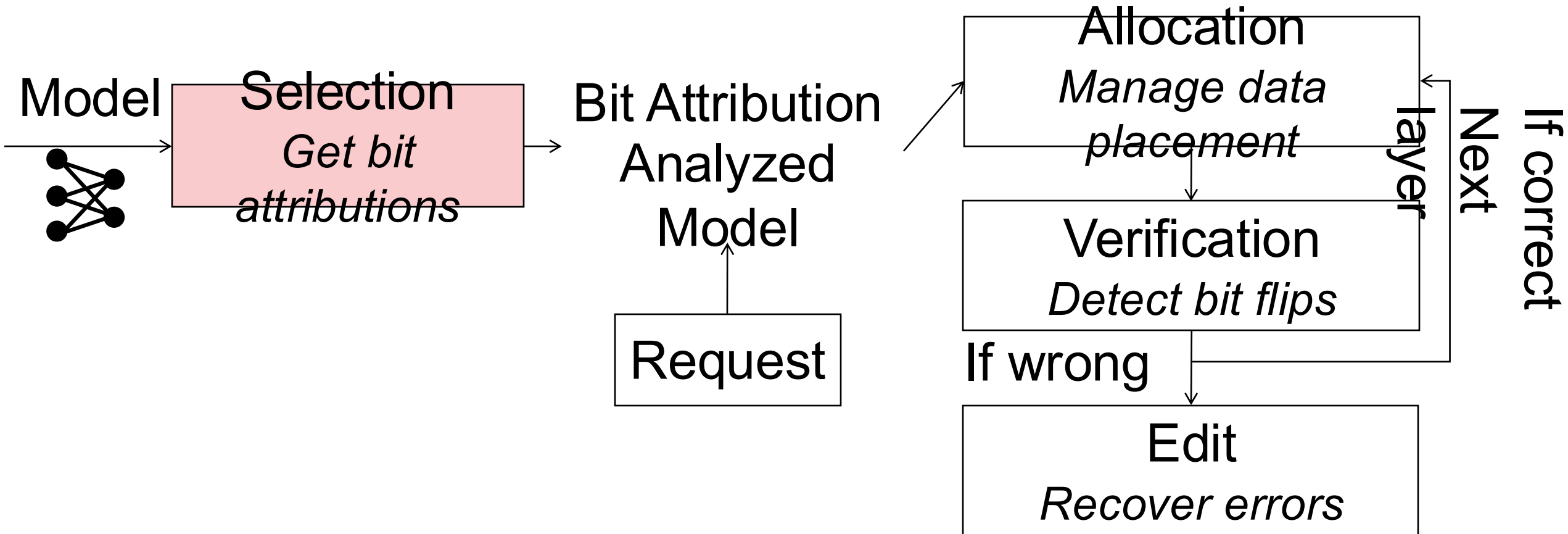


**DMR**

# SAVE Design



# SAVE Design



- **One-time** 2-Phase Analysis:

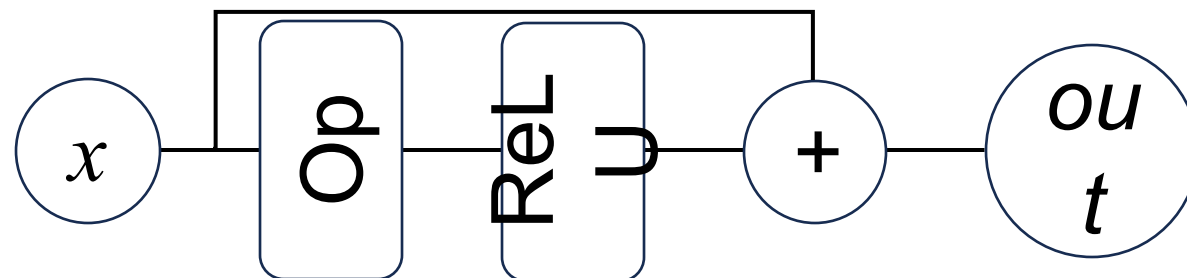
- Range Analysis: Use DAG to propagate the expected range
- Bit Attributions: Map ranges into per value bit attributions

# Selection Stage-Range Analysis

```
def NonResidualExample(x):
```

```
    x = relu(Op(x)) + x
```

```
    return x
```

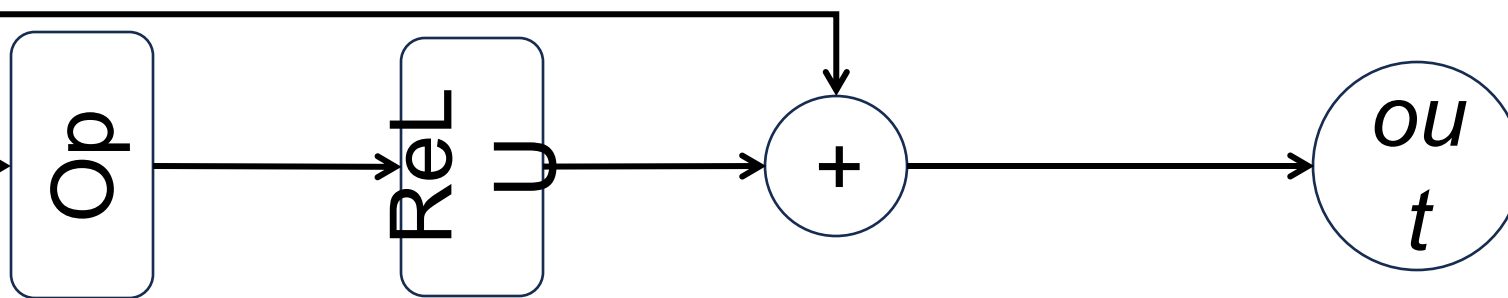


*Rangin*

$g$

$x$

0 to 1



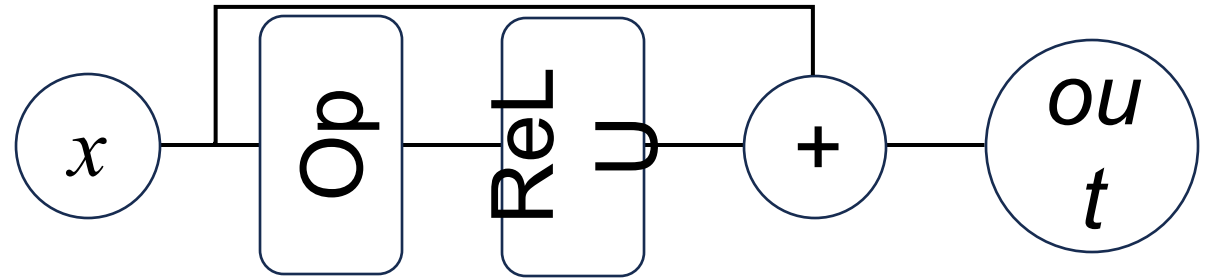
① *Input Constraint*

# Selection Stage-Range Analysis

```
def NonResidualExample(x):
```

```
    x = relu(Op(x)) + x
```

```
    return x
```



*Rangin*

$g$

$x$

0 to 1

② Propagate

$Op$

$ReLU$

$+$

$out$

① Input Constraint

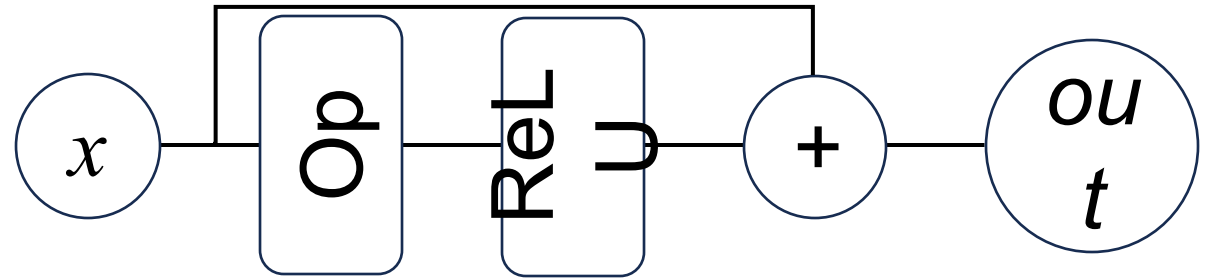
-0.75 to -0.1

# Selection Stage-Range Analysis

```
def NonResidualExample(x):
```

```
    x = relu(Op(x)) + x
```

```
    return x
```



*Rangin*

*g*

*x*

0 to 1

② Propagate

*Op*

*ReLU*

*+*

*out*

① Input Constraint

-0.75 to -0.1

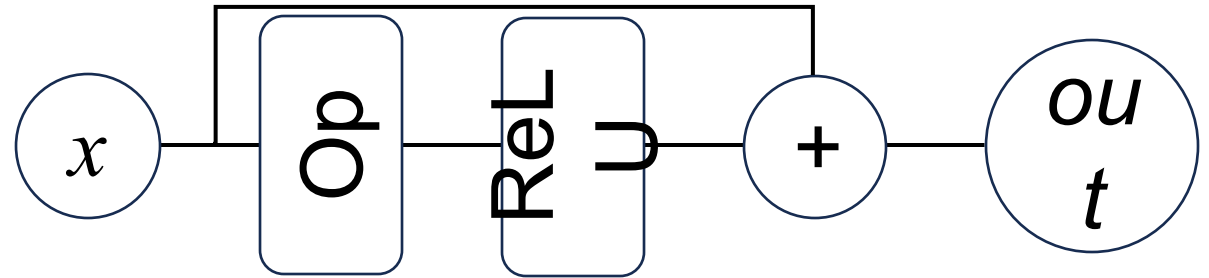
0

# Selection Stage-Range Analysis

```
def NonResidualExample(x):
```

```
    x = relu(Op(x)) + x
```

```
    return x
```



*Rangin*

*g*

*x*

0 to 1

② Propagate

*Op*

*ReLU*

*+*

*out*

① Input Constraint

-0.75 to -0.1

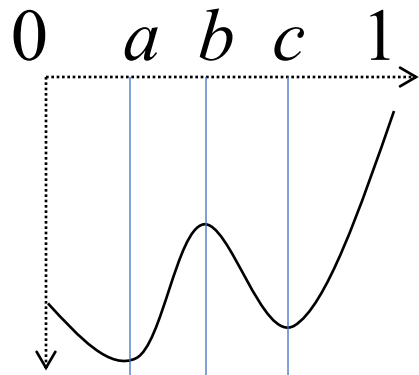
0

0 to 1

(0 + 0 to 1 + 0)

0 to 1

# Selection Stage-Bit Attribution Analysis

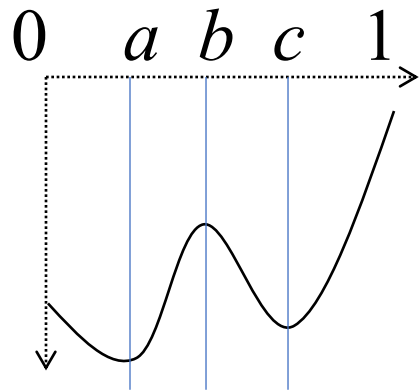


Input	Output	Ranging
0 to a	-0.65 to -0.75	9 Bits
...	...	...
c to 1	-0.7 to -0.1	6 Bits

Value	Binary
-0.65	<u>1 01111110</u> 0...
...	...
-0.75	<u>1 01111110</u> 1...

**9 Ranging B** *Map to bits*

# Selection Stage-Bit Attribution Analysis



Input	Output	Ranging
0 to a	-0.65 to -0.75	9 Bits
...	...	...
c to 1	-0.7 to -0.1	6 Bits

Value	Binary
-0.65	<u>1 01111110</u> 0...
...	...
-0.75	<u>1 01111110</u> 1...

**Ranging Bits**

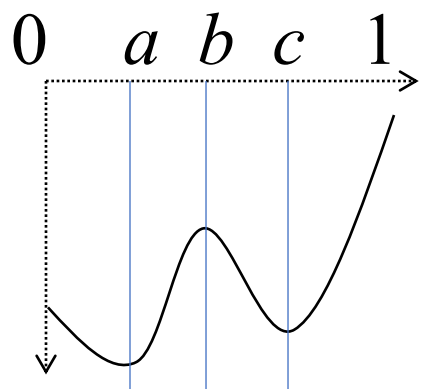
**9 Ranging B** *Map to bits*

Range	Binary(Significand)	Value
-0.65	...10011 <u>001100110</u>	-0.65
to		
-0.75	...10011 <u>100001110</u>	-0.65001

**9 Robust Robust Bits**

*Robust by Numerical*

# Selection Stage-Bit Attribution Analysis



Input	Output	Ranging
0 to a	-0.65 to -0.75	9 Bits
...	...	...
c to 1	-0.7 to -0.1	6 Bits

Value	Binary
-0.65	<u>1 01111110</u> 0...
...	...
-0.75	<u>1 01111110</u> 1...

**Ranging Bits**

**9 Ranging B** *Map to bits*

Range	Binary(Significand)	Value
-0.65	...10011 <u>001100110</u>	-0.65
to		
-0.75	...10011 <u>100001110</u>	-0.65001

**9 Robust Robust Bits**

*Robust by Numerical*

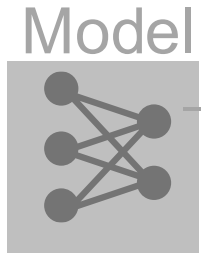


Any negative values give same results.

**31 Robust Bits**

*Robust by Op*

# SAVE Design



Bit Attribution  
Analyzed Model



layer

If correct  
Next



If wrong

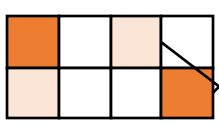


- Place the value with less robust bits into reliable memory

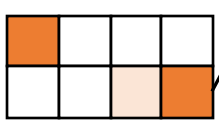
Input



Layer 1



Layer 2



...

Robust Bit Count

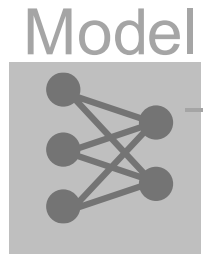
More

Less

Hardware  
Reliable

Other Memory

# SAVE Design



Bit Attribution  
Analyzed Model



- Using different verification method
  - Robust bits: Ignore
  - Ranging bits: Range Check
  - Vulnerable bits: Double Modular Redundancy

GPU

*Compute Op1*

GPU  
Async

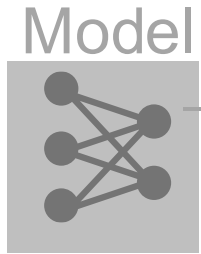
CPU

layer

If correct  
Next

If wrong

# SAVE Design



Bit Attribution  
Analyzed Model



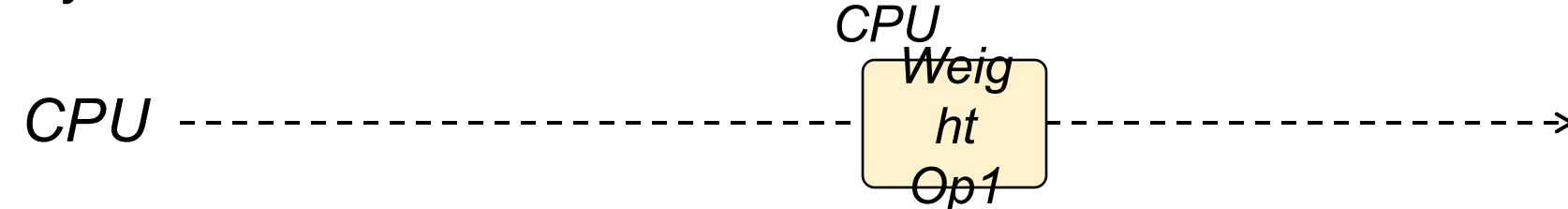
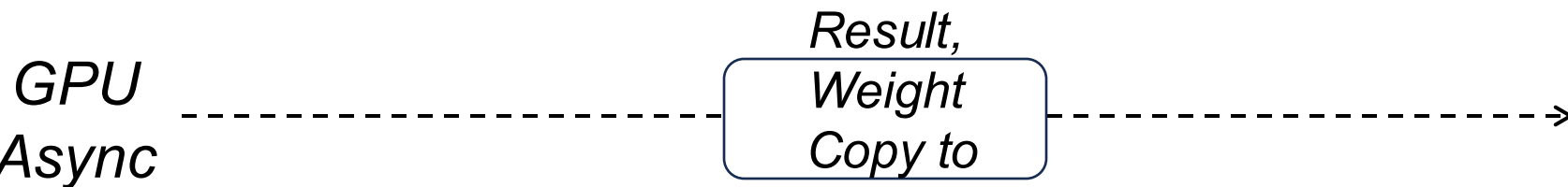
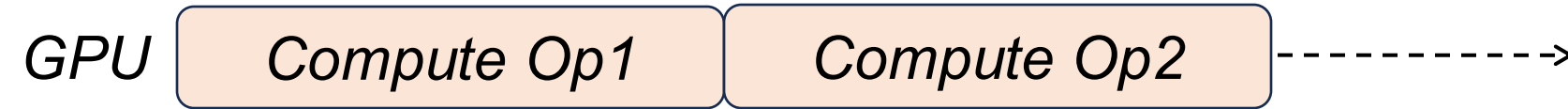
layer

If correct  
Next

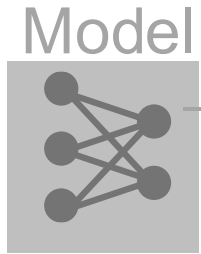
- Using different verification method
  - Robust bits: Ignore
  - Ranging bits: Range Check
  - Vulnerable bits: Double Modular Redundancy



If wrong



# SAVE Design



Bit Attribution  
Analyzed Model



layer

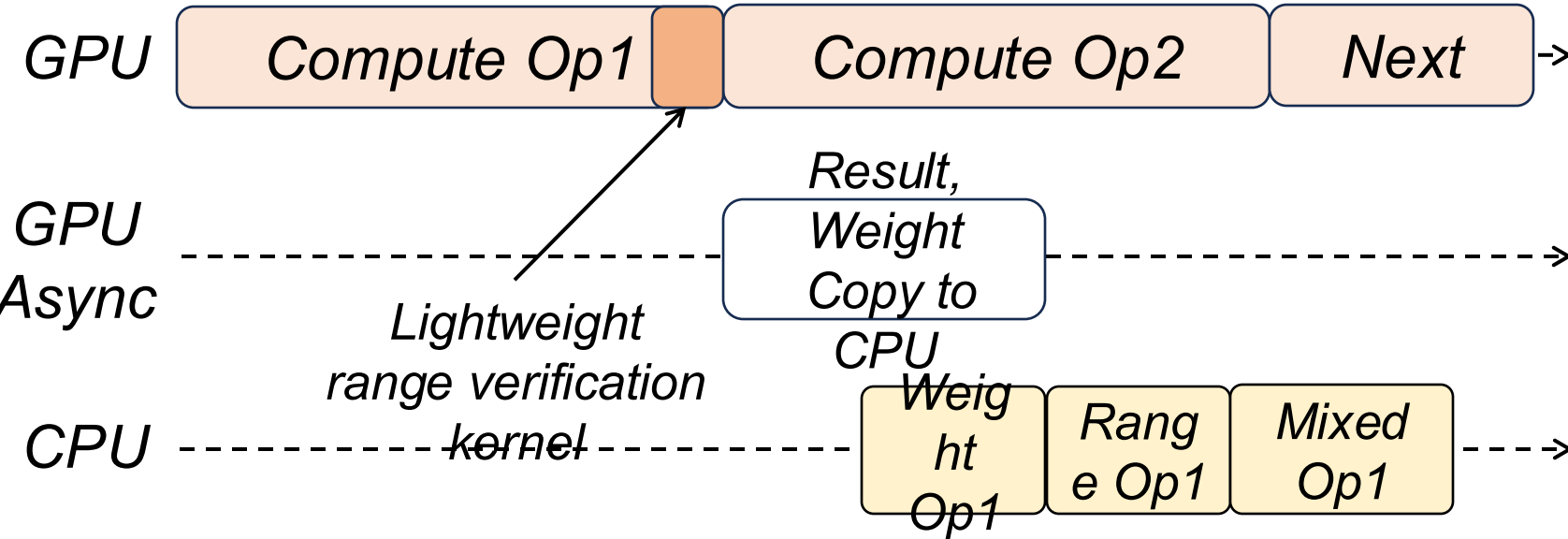
If correct  
Next



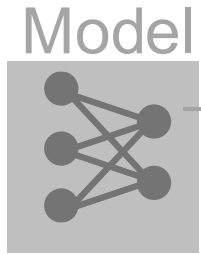
If wrong



- Using different verification method
  - Robust bits: Ignore
  - Ranging bits: Range Check
  - Vulnerable bits: Double Modular Redundancy



# SAVE Design



Selection  
*Get bit attributions*

Bit Attribution  
Analyzed Model

Allocation  
*Manage data placement*

layer

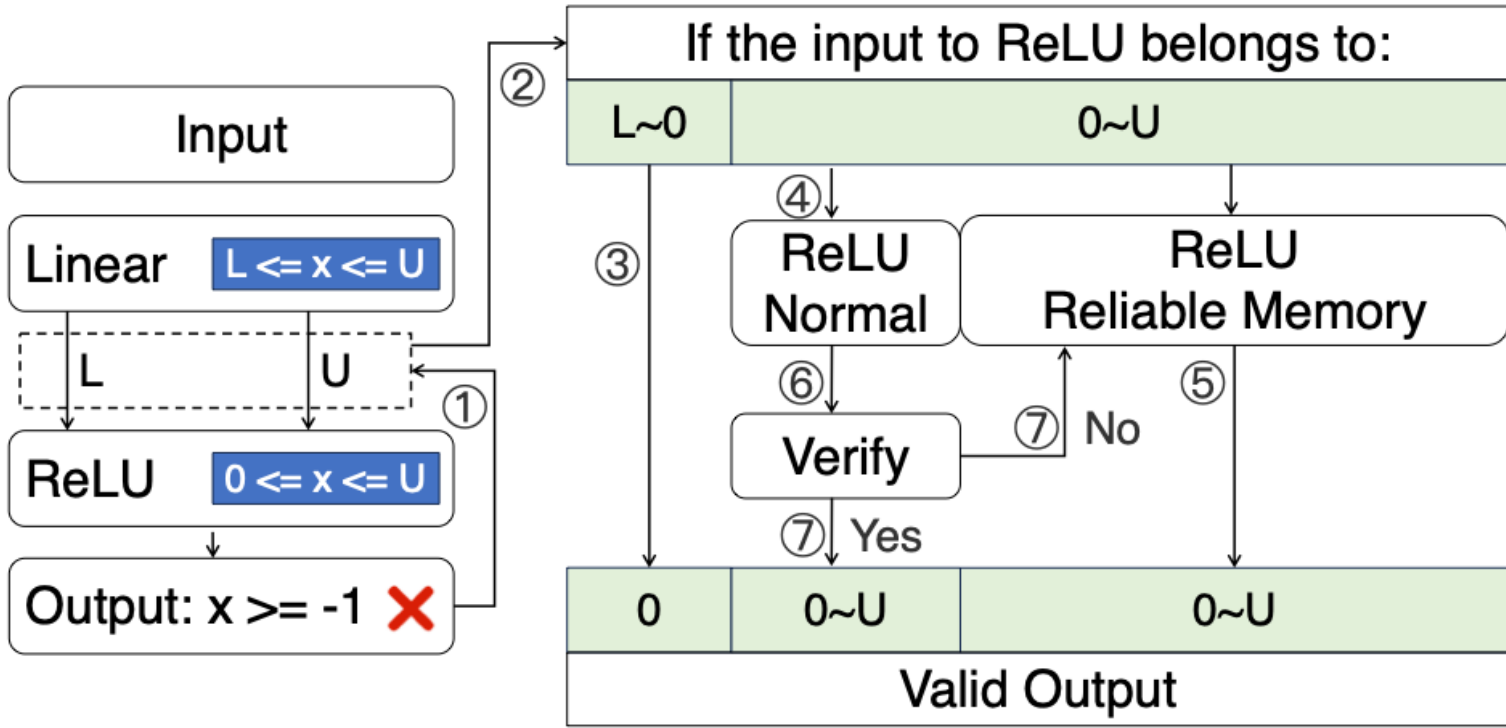
If correct  
Next

Verification  
*Detect bit flips*

If wrong

Edit  
*Recover errors*

- Restart the inference from the fault model layer



## New Metric: AccurateLatency

$$\textit{AccurateLatency} = \textit{Latency} \times (1 + \textit{ErrorRate})$$

Assume the second inference must be right

## New Metric: AccurateLatency

$$\textit{AccurateLatency} = \textit{Latency} \times (1 + \textit{ErrorRate})$$

Assume the second inference must be right

Correct Inference  
1 - ErrorRate

Latency

$$\textit{Latency} \times (1 - \textit{ErrorRate})$$

# New Metric: AccurateLatency

$$\text{AccurateLatency} = \text{Latency} \times (1 + \text{ErrorRate})$$

Assume the second inference must be right

Correct Inference  
1 - ErrorRate

Latency

Wrong Inference  
ErrorRate

Latency

$$\text{Latency} \times (1 - \text{ErrorRate})$$

# New Metric: AccurateLatency

$$\text{AccurateLatency} = \text{Latency} \times (1 + \text{ErrorRate})$$

Assume the second inference must be right

Correct Inference  
1 - ErrorRate

Latency

$$\text{Latency} \times (1 - \text{ErrorRate})$$

Wrong Inference  
ErrorRate

Latency

Latency

$$2 \times \text{Latency} \times \text{ErrorRate}$$

# New Metric: AccurateLatency

$$\text{AccurateLatency} = \text{Latency} \times (1 + \text{ErrorRate})$$

Assume the second inference must be right

Correct Inference  
1 - ErrorRate

Latency

$$\text{Latency} \times (1 - \text{ErrorRate})$$

Wrong Inference  
ErrorRate

Latency

Latency

$$2 \times \text{Latency} \times \text{ErrorRate}$$

$$\text{Total: } \text{Latency} \times (1 + \text{ErrorRate})$$

# New Metric: AccurateLatency

$$\text{AccurateLatency} = \text{Latency} \times (1 + \text{ErrorRate})$$

Assume the second inference must be right

Correct Inference  
1 - ErrorRate

Latency

$$\text{Latency} \times (1 - \text{ErrorRate})$$

Wrong Inference  
ErrorRate

Latency

Latency

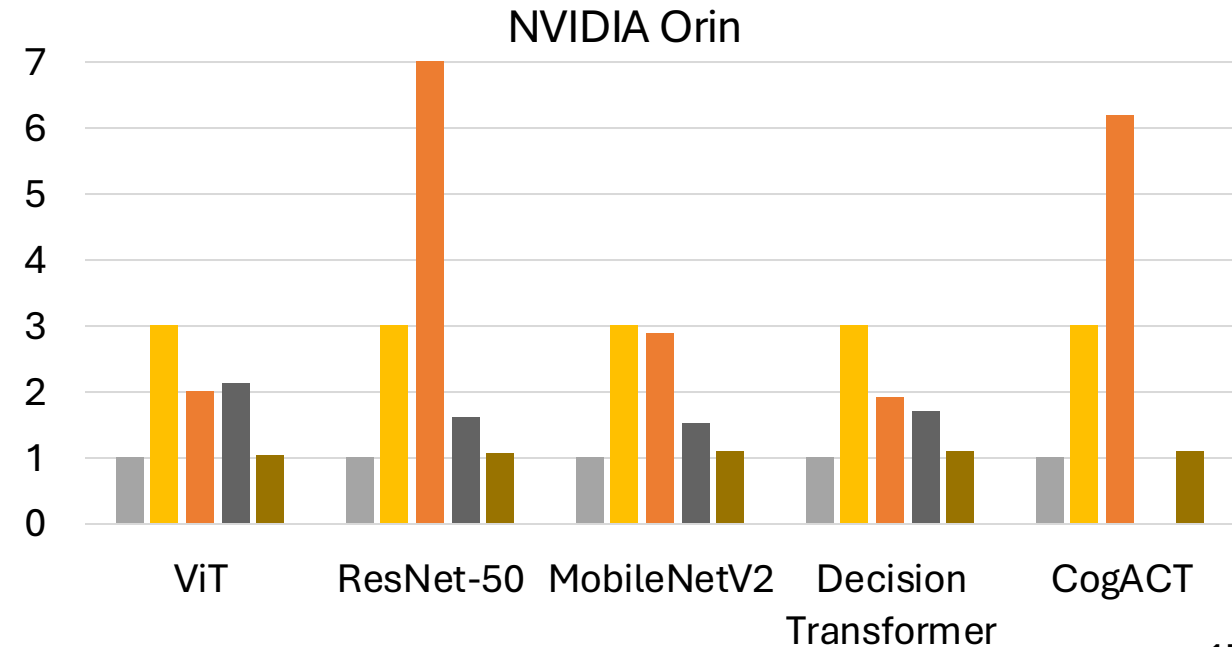
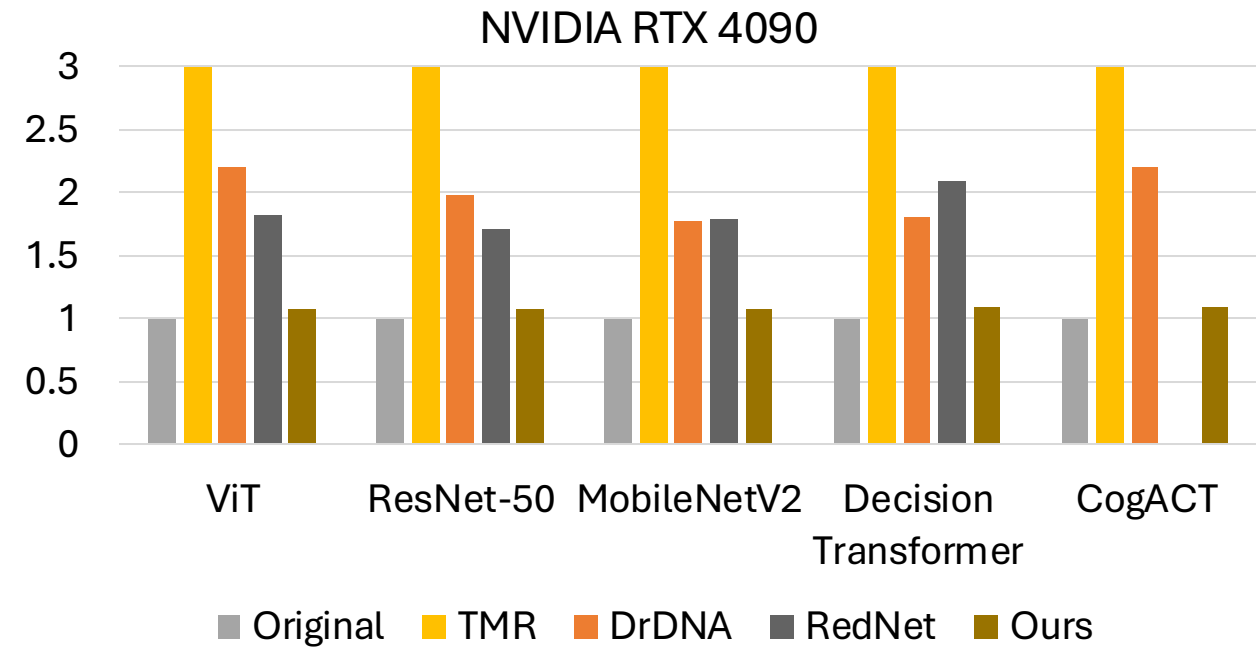
$$2 \times \text{Latency} \times \text{ErrorRate}$$

$$\text{Total: } \text{Latency} \times (1 + \text{ErrorRate})$$

**Add correctness into the latency**

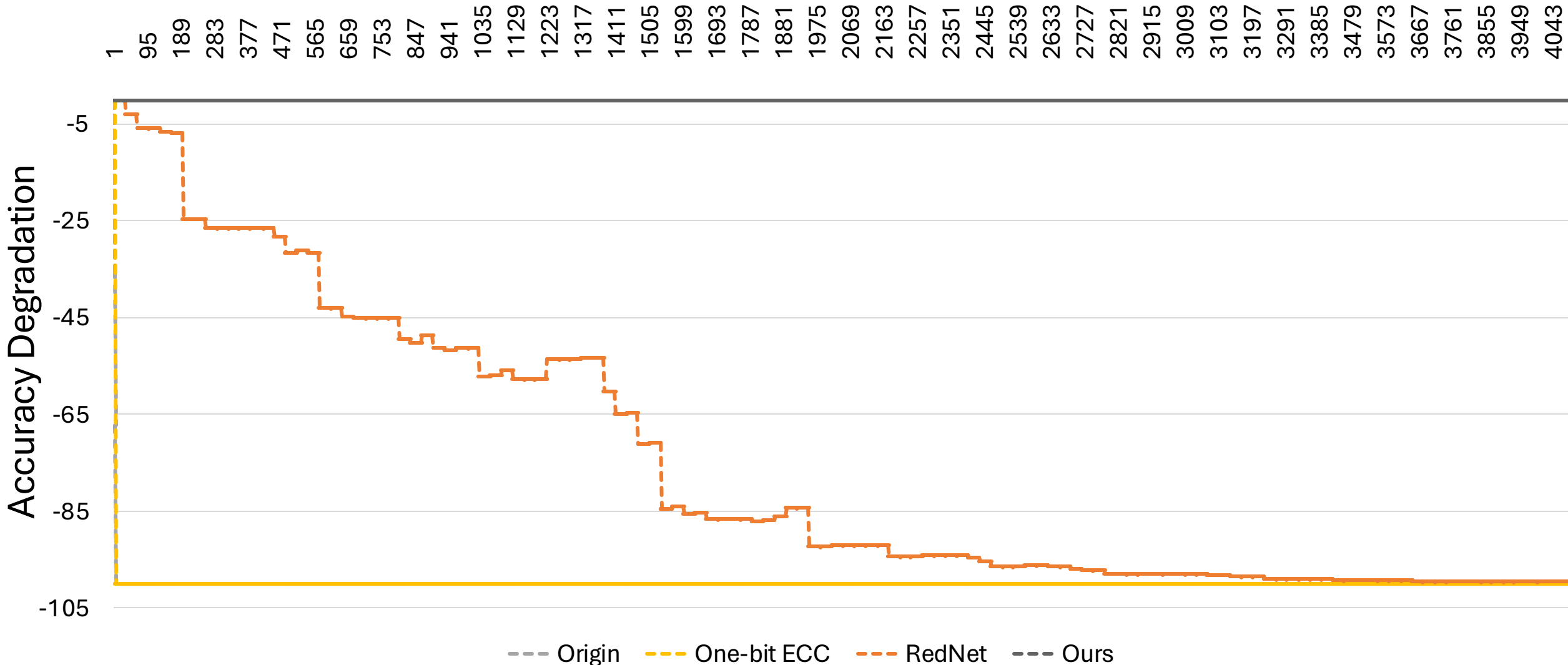
# End-to-end Performance

- Performance
  - Overall Accurate Latency **<9%**
  - Accuracy Degradation **<1%**
- Human Efforts
  - Modification to user code: **0 LoC**
  - Retraining Time: **0 Second**



# Bit Flip Resistance

Bit Flips Count



SAVE can keep the model accuracy **under 4000 bit flips** across different models

# Conclusion



- Software-implemented fault tolerance for machine learning models
- Key reason for Improvement: The differences in the bit properties inherent in computations within machine learning
- Evaluation shows that the overhead is as small as 9%
  - General to all models
  - No need for additional retraining or fine-tuning

# Thanks!