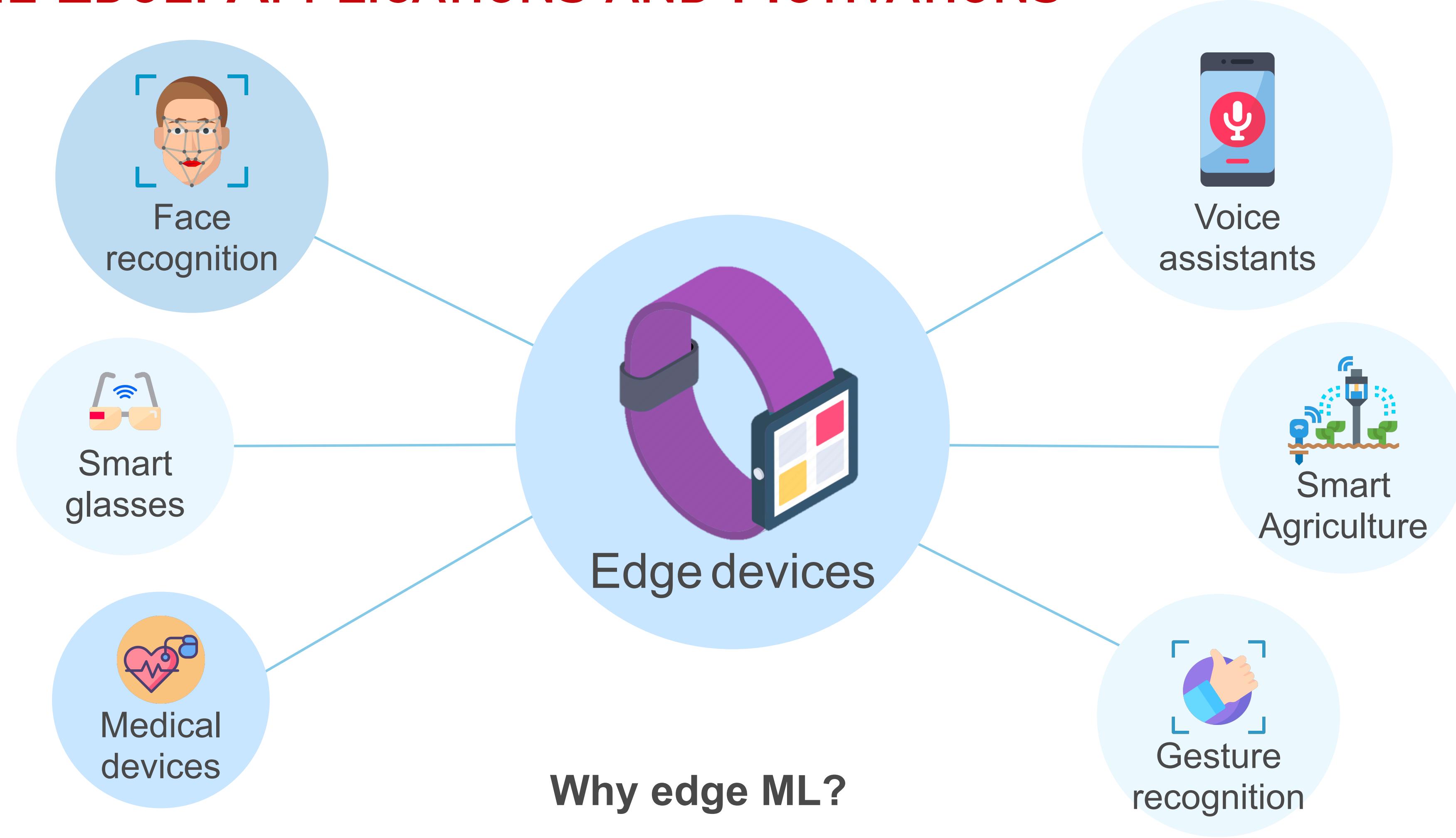


TIANEN CHEN, TAYLOR KEMP, AND YOUNGHYUN KIM  
UNIVERSITY OF WISCONSIN–MADISON

# SYNTHNET: A HIGH-THROUGHPUT YET ENERGY-EFFICIENT COMBINATIONAL LOGIC NEURAL NETWORK

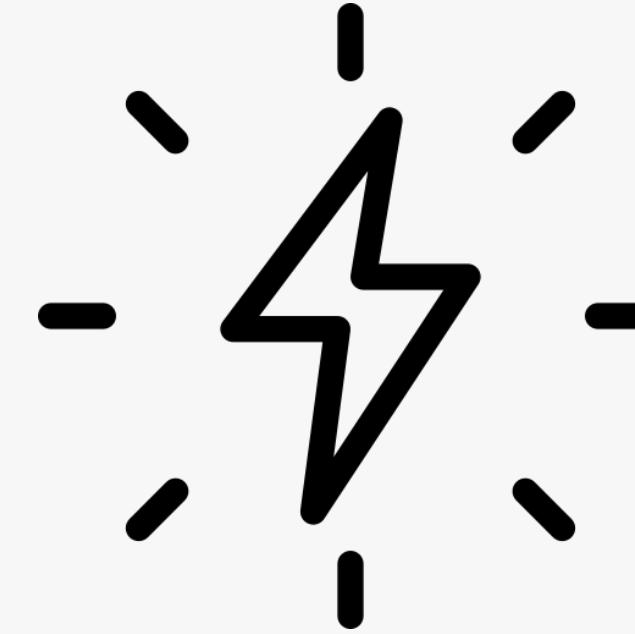
# ML AT THE EDGE: APPLICATIONS AND MOTIVATIONS



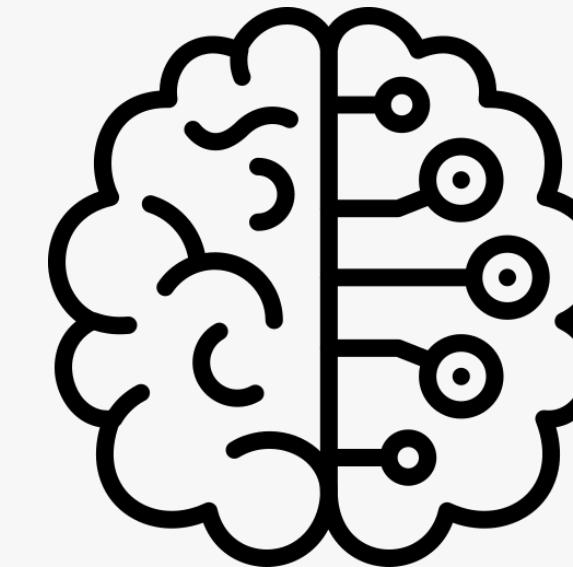
## Why edge ML?

- Less network communication load
- Reduced latency
- Enhanced privacy

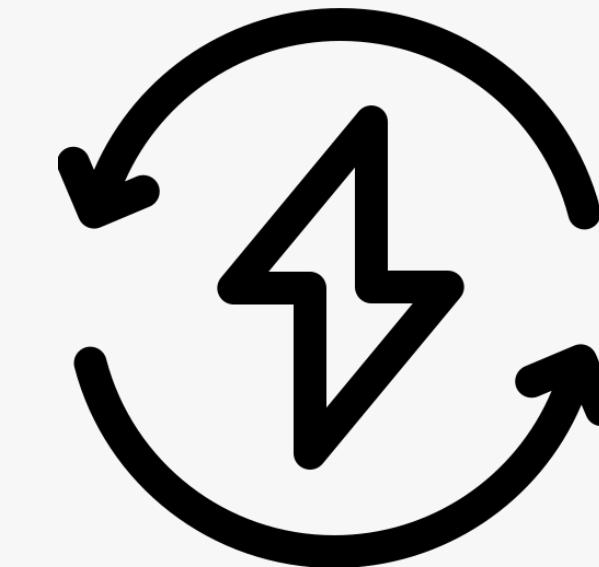
# ML AT THE EDGE: CHALLENGES



Energy constraints  
on edge devices



Energy-efficient ML  
systems



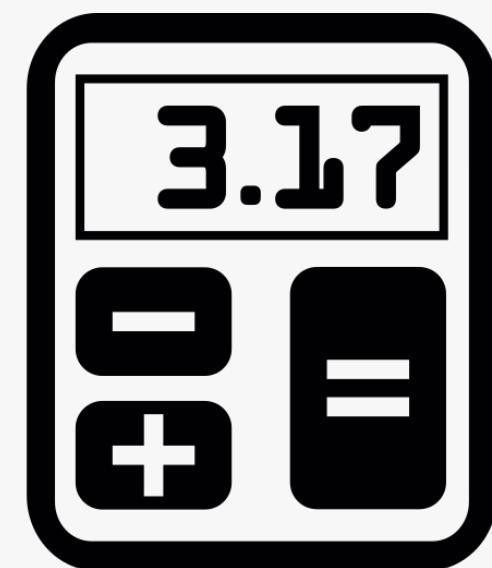
Energy-aware network  
design

## **X Heavy memory access overheads**

- Major source of energy consumption in conventional load-and-compute **neural processing element (NPE)** architectures
- **Multiply-and-accumulate (MAC)** operations dominate memory accesses

# TECHNIQUES FOR REDUCING MEMORY ACCESSES

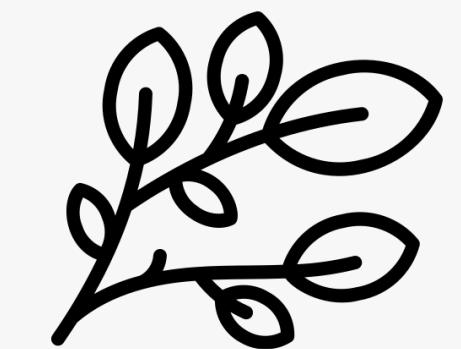
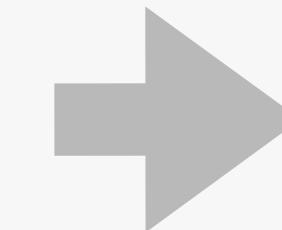
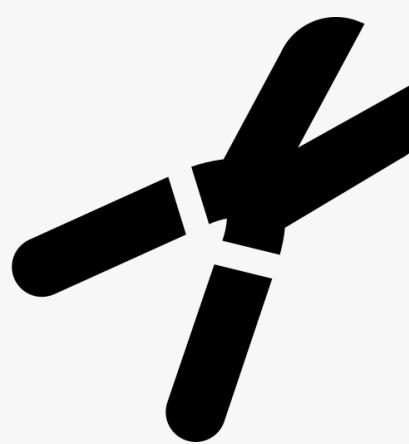
## *Quantization*



**0110  
1001  
1010**

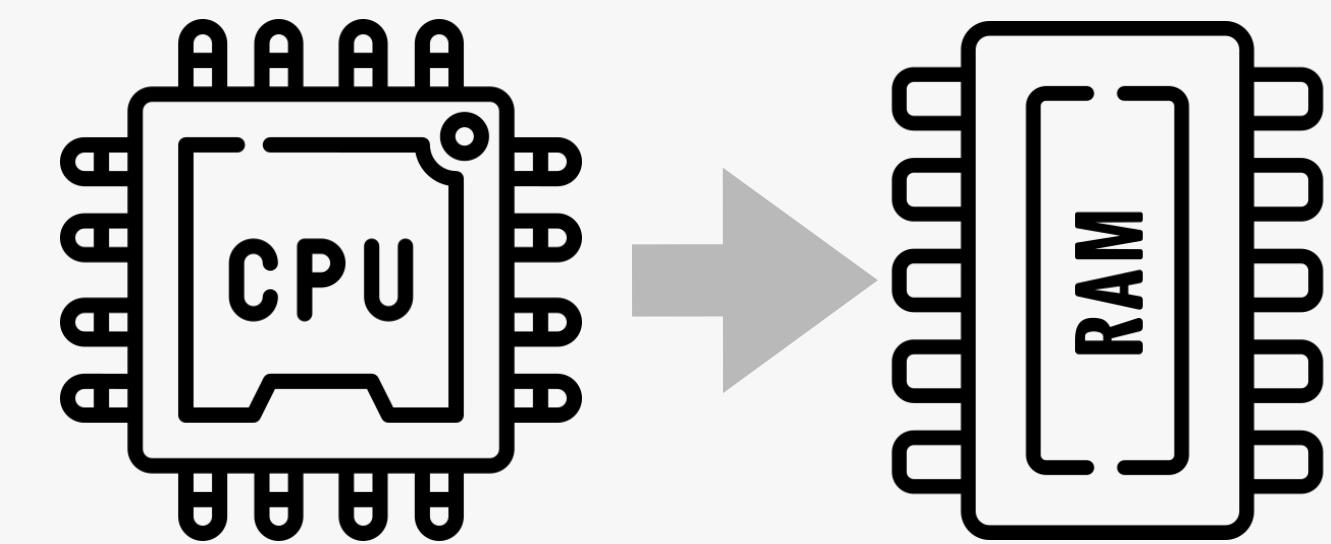
Replace floating-points with numbers with simpler representations

## *Pruning*



Remove connections in order to discriminate against low-magnitude weights

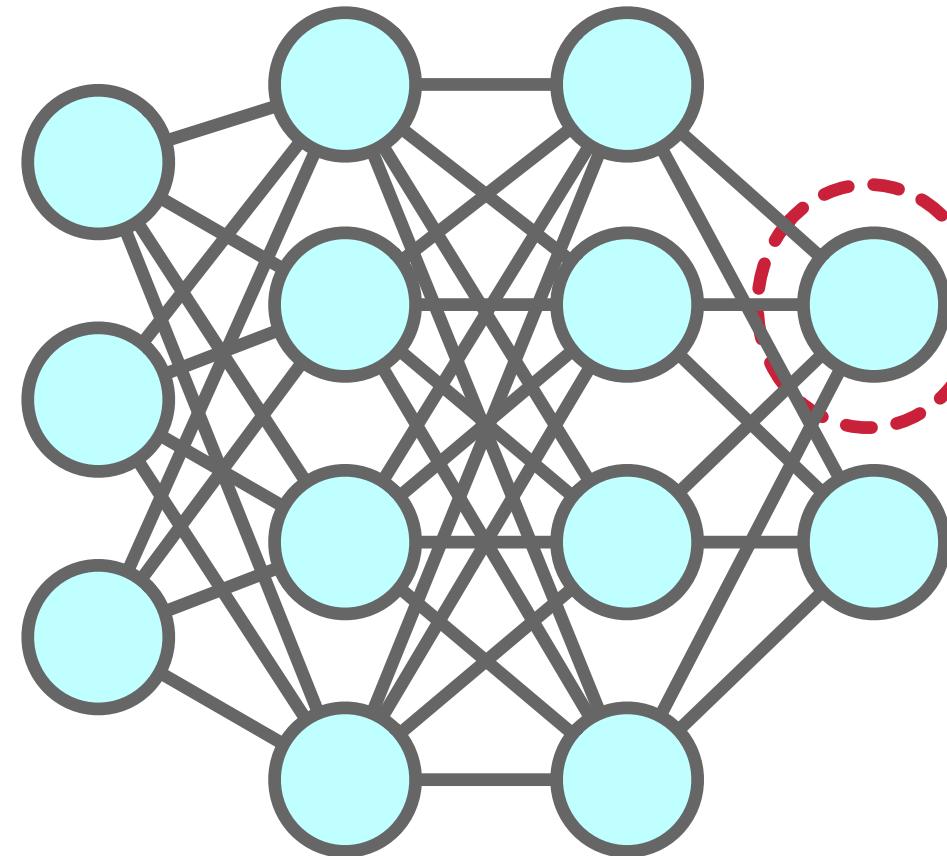
## *Processing in memory*



Perform computation **in the memory** without fetching data to the processor

# A NEW APPROACH: LOGIC-BASED NEURAL NETWORKS

*Neurons are implemented as a Boolean logic circuit*

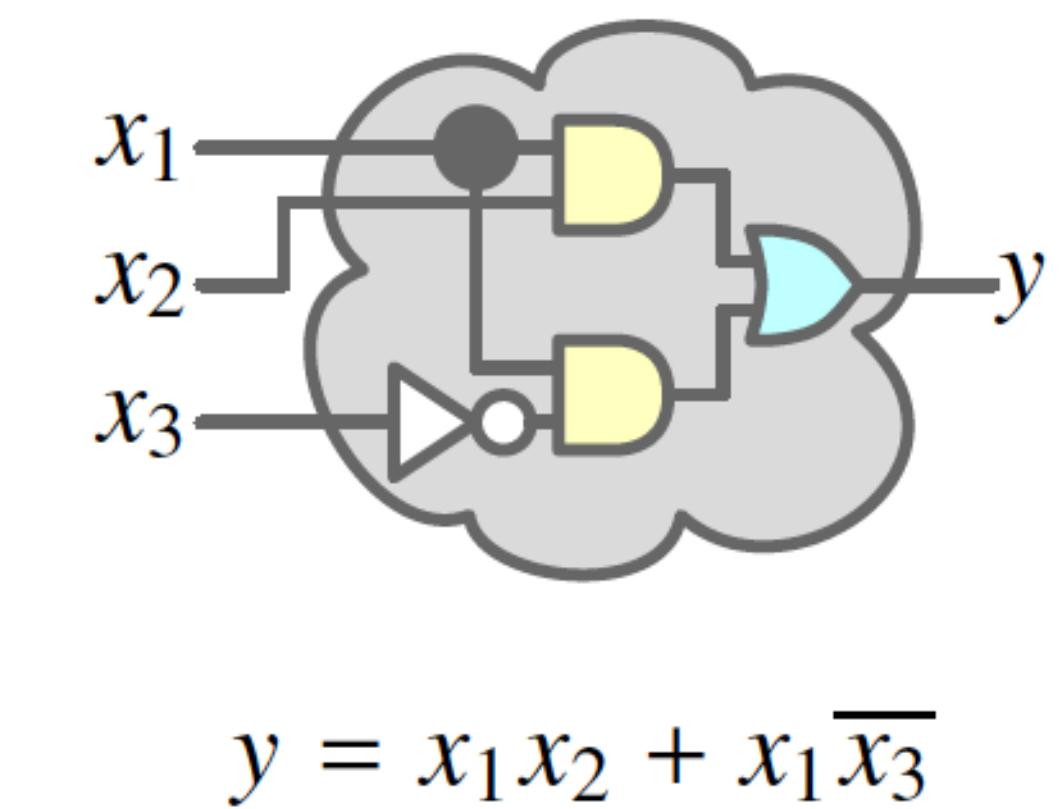


McCulloch-Pitts neuron

$$S = \sum_j a^j \times w^j$$
$$x_1 \xrightarrow{w_1 = 10}$$
$$x_2 \xrightarrow{w_2 = 1}$$
$$x_3 \xrightarrow{w_3 = -5}$$
$$\Sigma$$
$$b = 6$$
$$S \geq b \rightarrow y$$

$x_1$	$x_2$	$x_3$	$\Sigma$	$y$
0	0	0	-6	0
0	0	1	-16	0
0	1	0	-4	0
0	1	1	-14	0
1	0	0	14	1
1	0	1	4	0
1	1	0	16	1
1	1	1	6	1

Truth table



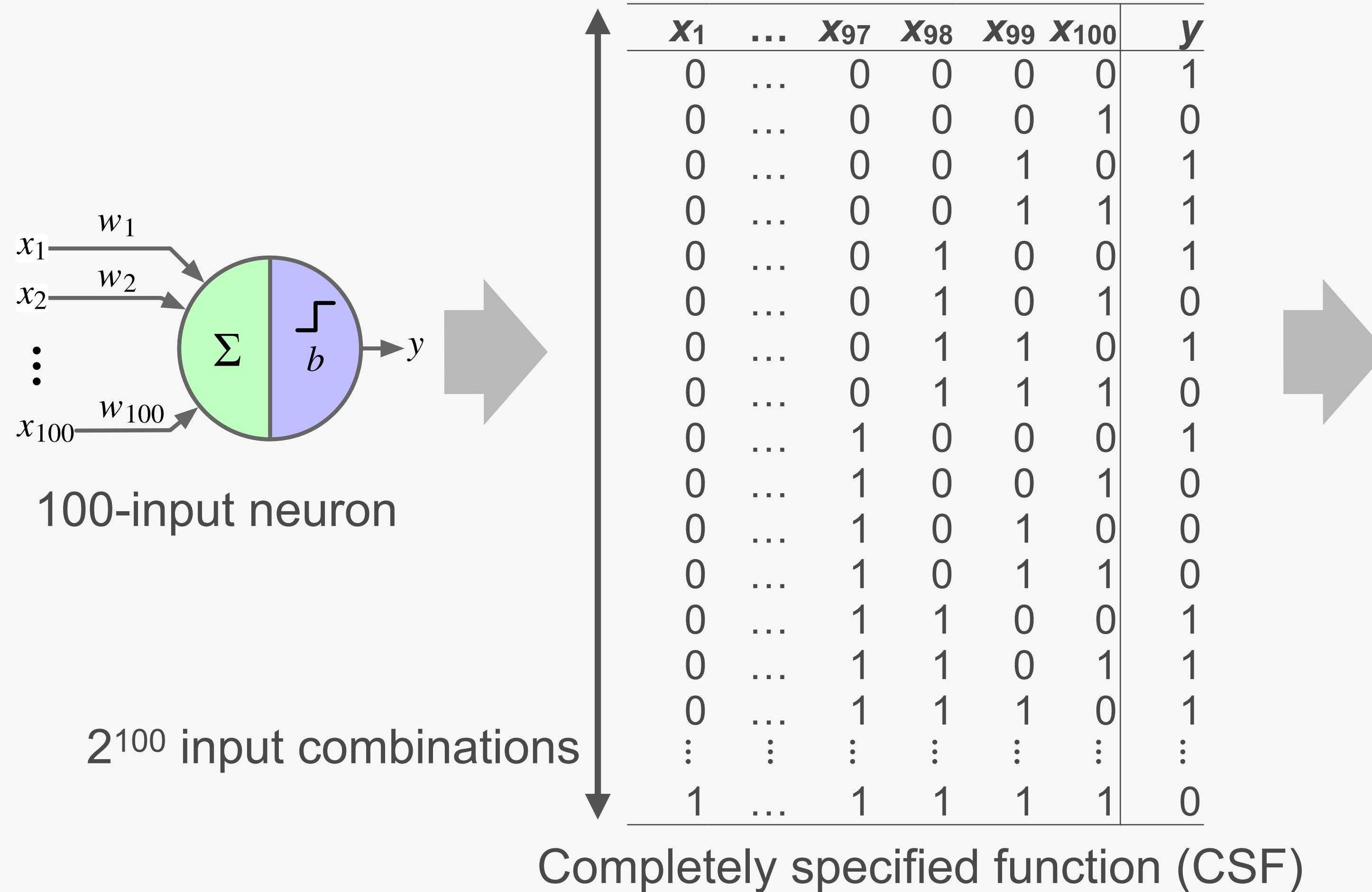
Logic implementation

Pre-evaluated Boolean mappings functions  
of the neurons

Weights are hardwired into the logic, not  
stored in the memory

✓ No memory access overheads

# LOGIC SIZE REDUCTION WITH INCOMPLETELY SPECIFIED FUNCTIONS



*Not feasible*

The diagram shows a table for an Incompletely Specified Function (ISF) with 100 inputs and 1 output. The table has columns for  $x_1, \dots, x_{97}, x_{98}, x_{99}, x_{100}$  and  $y$ . The first row shows an input pattern  $0|6|1 \dots |6$  with  $y=1$ . The second row shows an input pattern  $7$  with  $y=0$ . Subsequent rows show various input patterns with  $y$  values of 0 or 1. Red 'X' marks are placed in the  $y$  column for inputs where the function value is unknown or unspecified. A large arrow points from the previous CSF table to this ISF table.

Input	$x_1$	$\dots$	$x_{97}$	$x_{98}$	$x_{99}$	$x_{100}$	$y$
0 6 1 ... 6	0	...	0	1	1	0	1
7	0	...	1	0	1	0	0
0 6 1 ... 6	0	...	0	0	0	1	X
0 6 1 ... 6	0	...	0	0	1	0	X
0 6 1 ... 6	0	...	0	1	0	0	X
0 6 1 ... 6	0	...	0	1	0	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1 ... 6	0	...	1	1	1	1	X
0 6 1 ... 6	0	...	1	0	0	0	X
0 6 1 ... 6	0	...	1	0	0	1	X
0 6 1 ... 6	0	...	1	1	0	0	X
0 6 1 ... 6	0	...	1	1	0	1	X
0 6 1 ... 6	0	...	1	1	1	0	X
0 6 1							

# MOTIVATION: UNEXPLOITED ERROR RESILIENCE

Observed for  
many inputs

→ **0|6|1|...|6**

Observed only  
for 1 input

Input	$x_1$	...	$x_{97}$	$x_{98}$	$x_{99}$	$x_{100}$	$y$
	0	...	0	0	0	0	X
	0	...	0	0	0	1	X
	0	...	0	0	1	0	X
	0	...	0	0	1	1	X
	0	...	0	1	0	0	X
	0	...	0	1	0	1	X
	0	...	0	1	1	0	1
	0	...	0	1	1	1	X
	0	...	1	0	0	0	X
	0	...	1	0	0	1	X
	0	...	1	0	1	0	0
	0	...	1	0	1	1	X
	0	...	1	1	0	0	X
	0	...	1	1	0	1	X
	0	...	1	1	1	0	X
	⋮	⋮	⋮	⋮	⋮	⋮	⋮
	1	...	1	1	1	1	X

*Not all input combinations are  
equally important*

← **Important for accuracy**

← **Not important for accuracy**

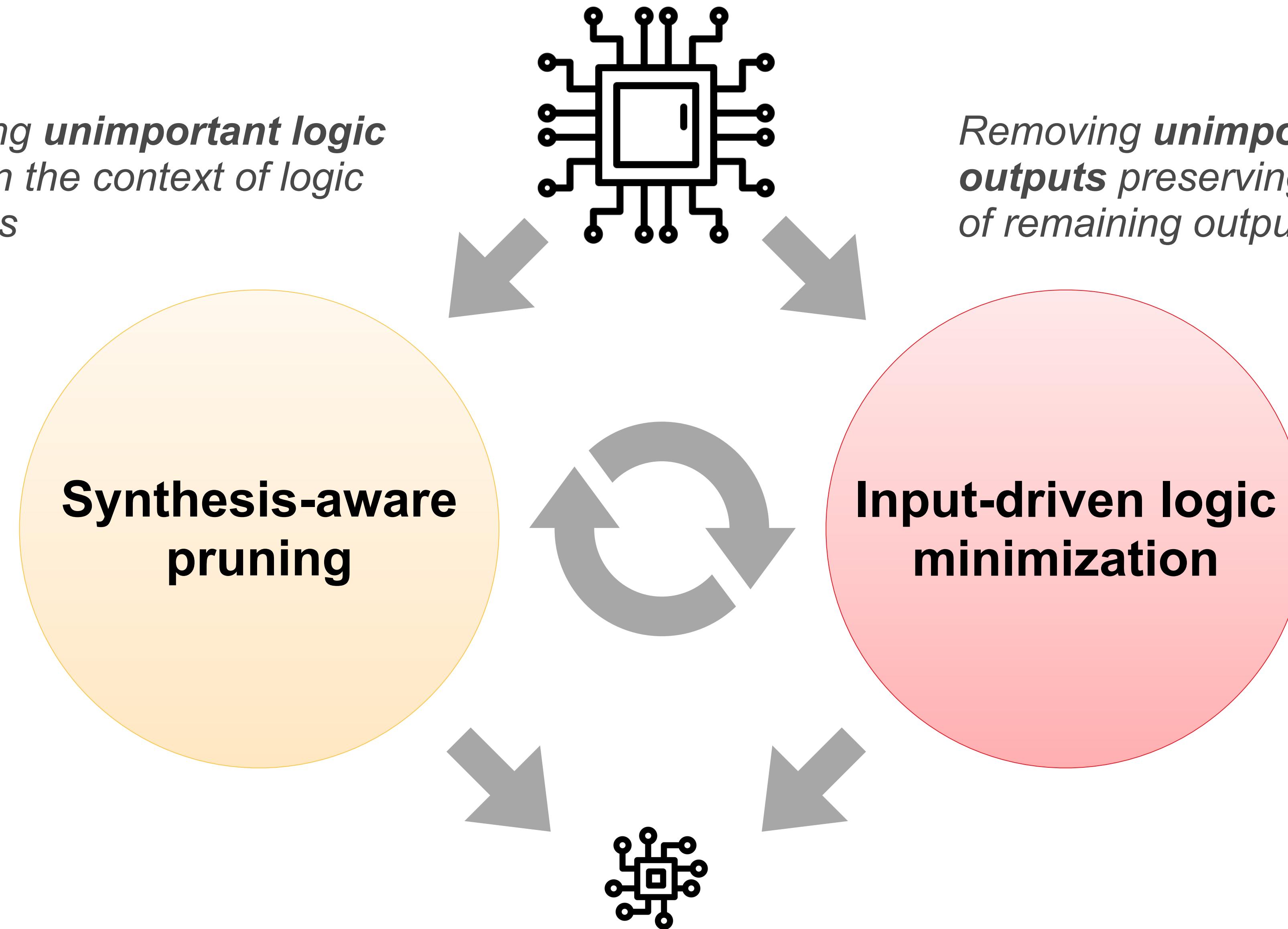
= Opportunity for further logic  
reduction without accuracy  
degradation

*Question: How do we identify and remove unimportant outputs?*

# PROPOSED METHOD: TWO-COORDINATE LOGIC MINIMIZATION

*Removing **unimportant logic inputs** in the context of logic synthesis*

*Removing **unimportant logic outputs** preserving importance of remaining outputs*



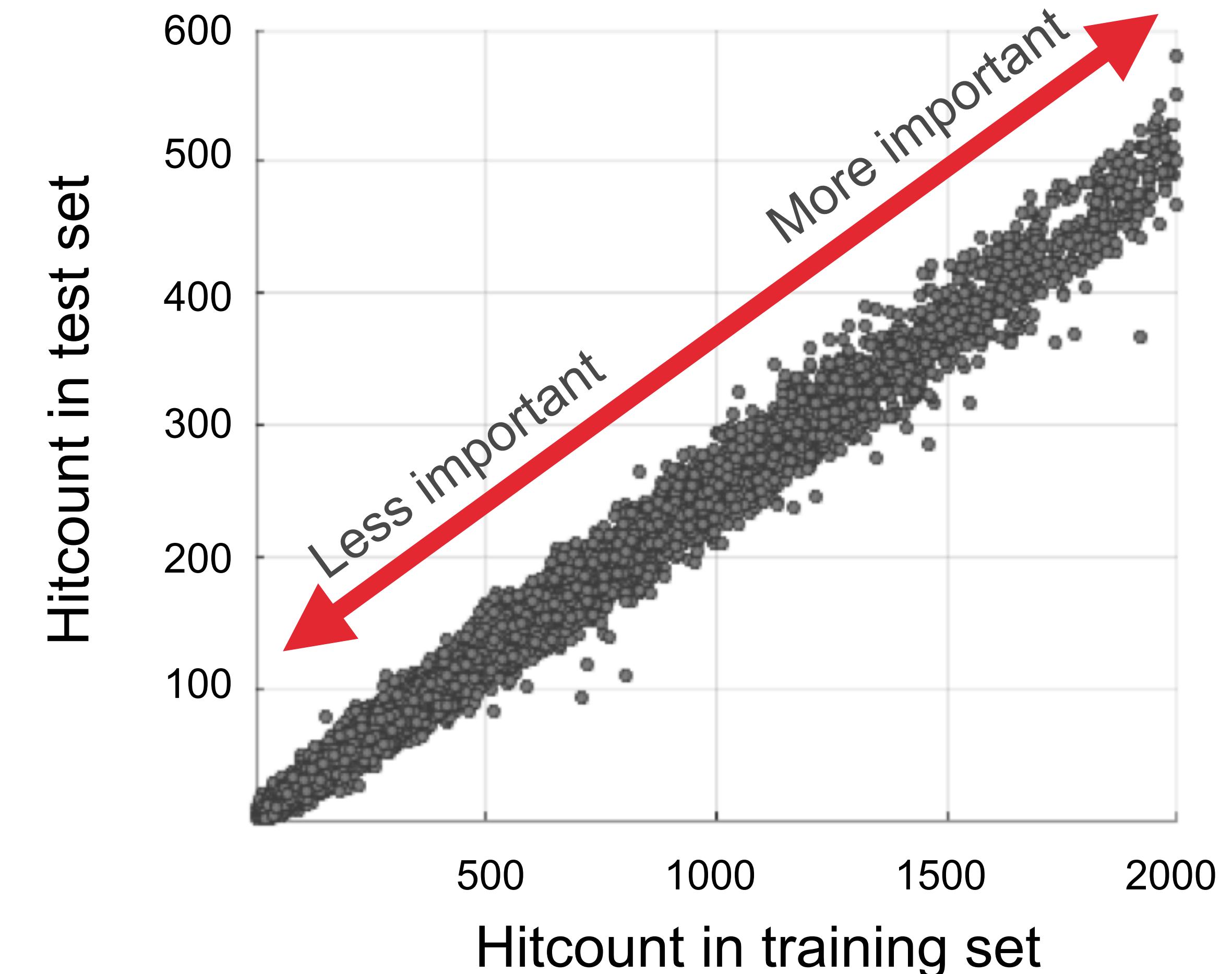
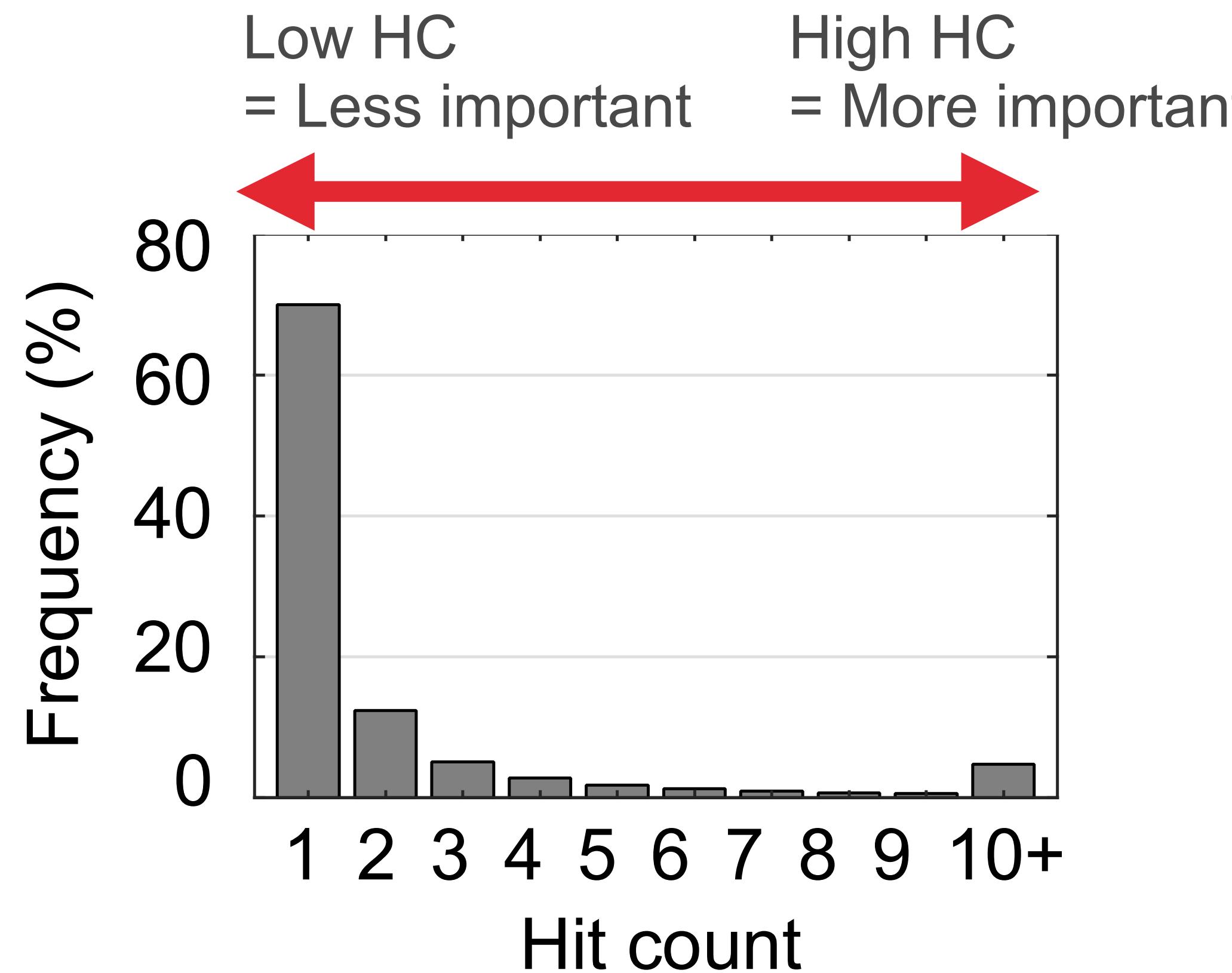
# HITCOUNT: AN IMPORTANCE METRIC OF INPUTS

- **Hitcount (HC)**: Number of encounters of input combination
- High HC  
= High input encounter frequency  
= Important input combination
- Output is specified only when the input is encountered at least once during training (i.e.,  $HC \geq 1$ )

Input	$x_1$	...	$x_{97}$	$x_{98}$	$x_{99}$	$x_{100}$	$y$	$HC$
	0	...	0	0	0	0	x	0
	0	...	0	0	0	1	x	0
	0	...	0	0	1	0	x	0
	0	...	0	0	1	1	x	0
	0	...	0	1	0	0	x	0
	0	...	0	1	0	1	x	0
<b>0 6 7 ... 6</b>	0	...	0	1	1	0	1	<b>10000</b>
	0	...	0	1	1	1	x	0
	0	...	1	0	0	0	x	0
	0	...	1	0	0	1	x	0
<b>7</b>	0	...	1	0	1	0	0	<b>1</b>
	0	...	1	0	1	1	x	0
	0	...	1	1	0	0	x	0
	0	...	1	1	0	1	x	0
	0	...	1	1	1	0	x	0
	:	:	:	:	:	:	:	0
	1	...	1	1	1	1	x	0

# HITCOUNT: UNEQUAL IMPORTANCE OF INPUTS

Dataset: MNIST



Only a very small subset of input combinations are important

HC in the training set well represents the importance of input combinations in the test set

# METHOD 1: SYNTHESIS-AWARE PRUNING

- Magnitude-based pruning – Eliminate inputs with low magnitudes
- Update HC:** Add HCs of merged rows  $HC' = \sum_{i \in I} HC_i$

Remove low-magnitude input  $x_2$

$x_1$	$x_2$	$x_3$	$y$	$HC$
10	1	-5		
0	0	0	0	2000
0	0	1	X	0
0	1	0	X	0
0	1	1	X	0
1	0	0	1	5
1	0	1	0	10
1	1	0	1	150
1	1	1	1	1000

$x_1$	$x_2$	$x_3$	$y$	$HC$
10	1	-5		
0	X	0	0	2000
0	X	1	X	0
0	X	0	X	0
0	X	1	X	0
1	X	0	1	5
1	X	1	0	10
1	X	0	1	150
1	X	1	1	1000

$x_1$	$x_3$	$y'$	$HC'$
10	-5	?	2000+0=2000
0	0	?	0+0=0
0	1	?	5+150=155
1	0	?	10+1000=1010

Merge rows and update HC

## METHOD 1: SYNTHESIS-AWARE PRUNING

- **Update output:** Weighted average of merged rows

$$y' = \text{round} \left( \frac{\sum_{i \in I} y_i \times HC_i}{\sum_{i \in I} HC_i} \right)$$

The diagram illustrates the synthesis-aware pruning process. It shows a transition from an initial dataset to an updated output table.

**Initial Dataset:**

$x_1$	$x_2$	$x_3$	$y$	$HC$
10	1	-5		
0	X	0	0	2000
0	X	1	X	0
0	X	0	X	0
0	X	1	X	0
1	X	0	1	5
1	X	1	0	10
1	X	0	1	150
1	X	1	1	1000

**Output Table:**

$x_1$	$x_3$	$y'$	$HC'$
10	-5		
0	0	round((0\times2000+0\times0)/2000)= 0	2000
0	1	round((X\times0+X\times0)/0)= X	0
1	0	round((1\times5+1\times150)/155)= 1	155
1	1	round((0\times10+1\times1000)/1010)= 1	1010

**Update output:** This label points to the final value in the  $y'$  column of the output table, which is 1.

## METHOD 2: INPUT-DRIVEN LOGIC MINIMIZATION

- Unspecify output (switch to Don't-Care) for unimportant inputs
- Makes the truth table even more sparse, resulting in even smaller logic
- HC remains unchanged

Original Truth Table:

$x_1$	$x_2$	$x_3$	$y$	HC
10	1	-5		
0	0	0	0	2000
0	0	1	X	0
0	1	0	X	0
0	1	1	X	0
1	0	0	1	5
1	0	1	0	10
1	1	0	1	150
1	1	1	1	1000

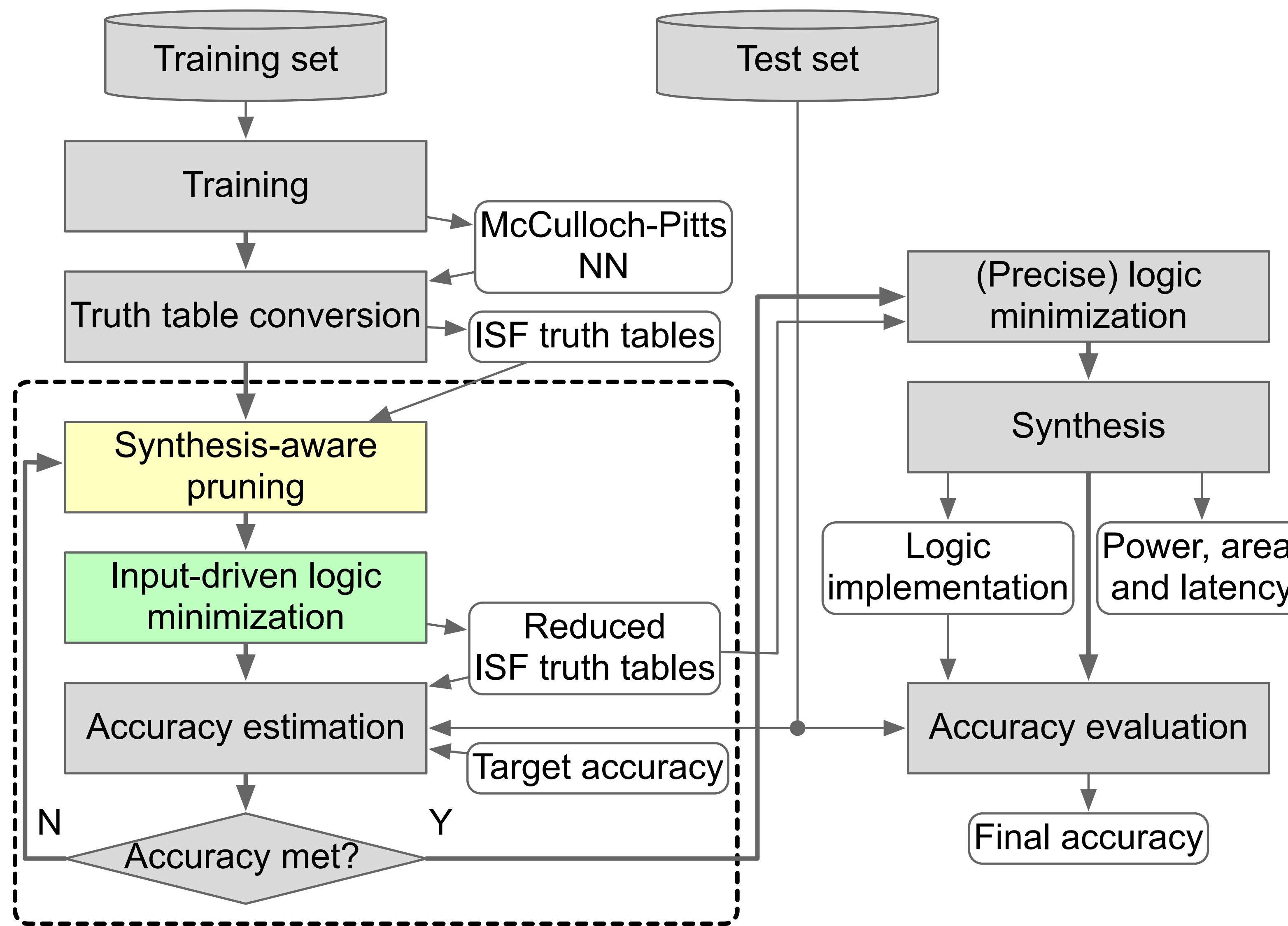
Modified Truth Table (Output unspecified for rows 3, 4, 5, 6, 7, 8):

$x_1$	$x_2$	$x_3$	$y'$	HC
10	1	-5		
0	0	0	0	2000
0	0	1	X	0
0	1	0	X	0
0	1	1	X	0
1	0	0	X	5
1	0	1	X	10
1	1	0	X	150
1	1	1	X	1000

Low HC {

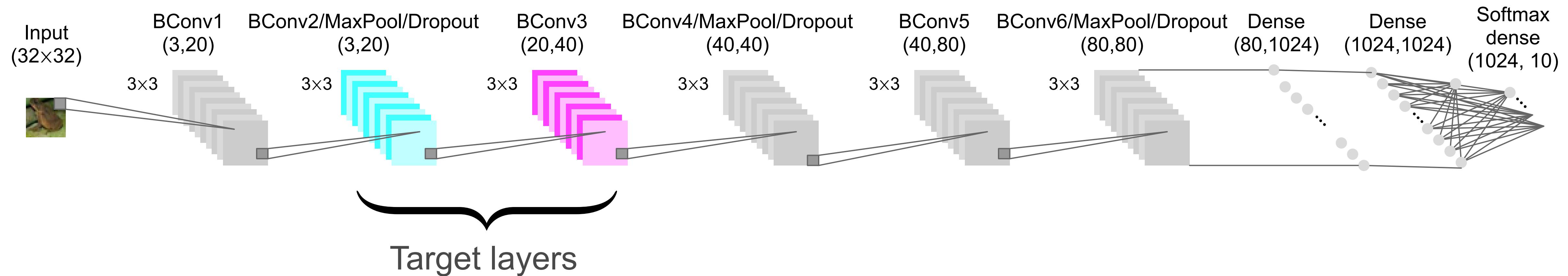
Output unspecified ↑

# TWO METHODS APPLIED IN DESIGN FLOW



- Synthesis-aware pruning and input-driven logic minimization iteratively applied
- Accuracy estimated from reduced truth tables without synthesis
- Repeated until target accuracy is met

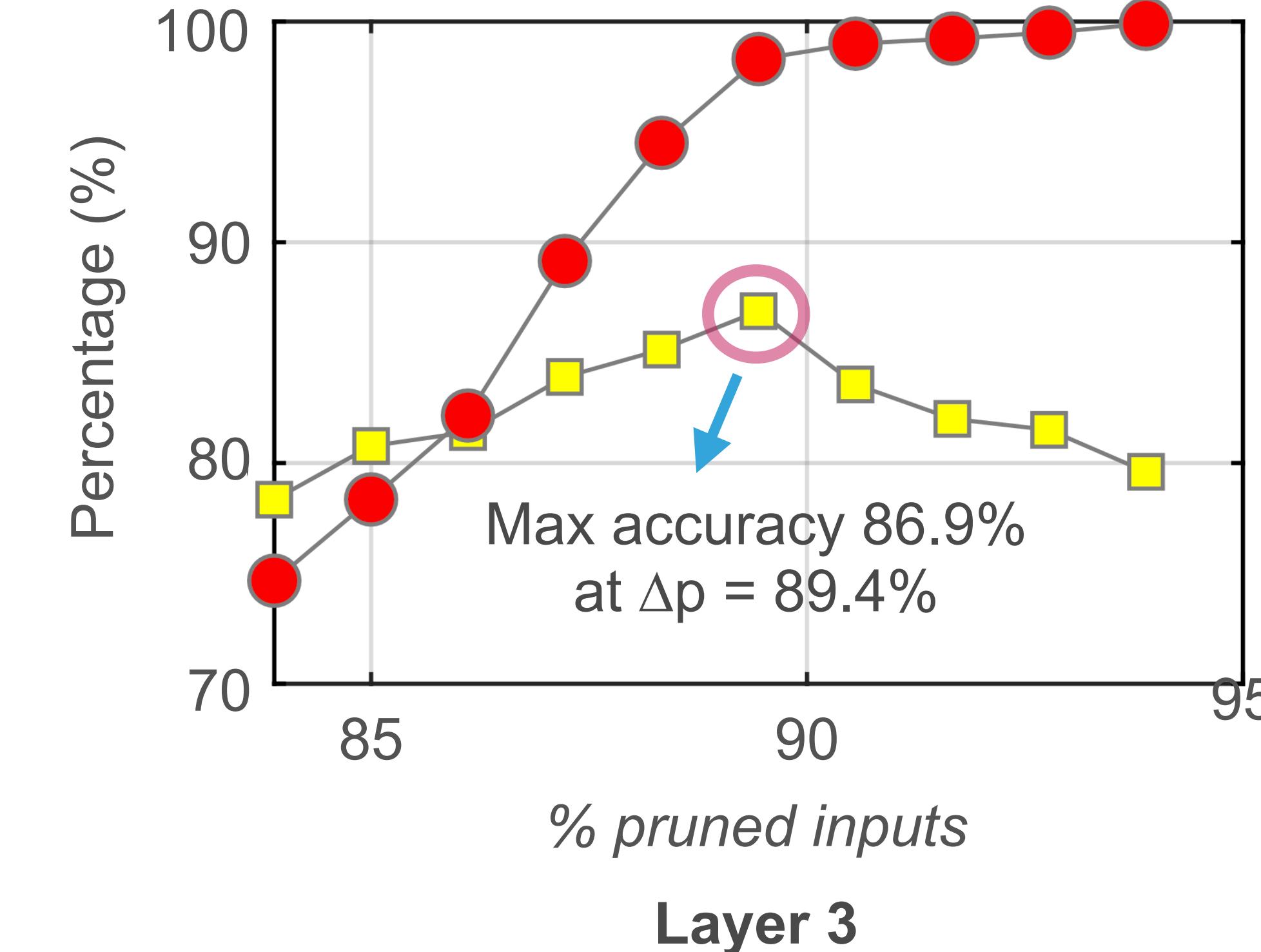
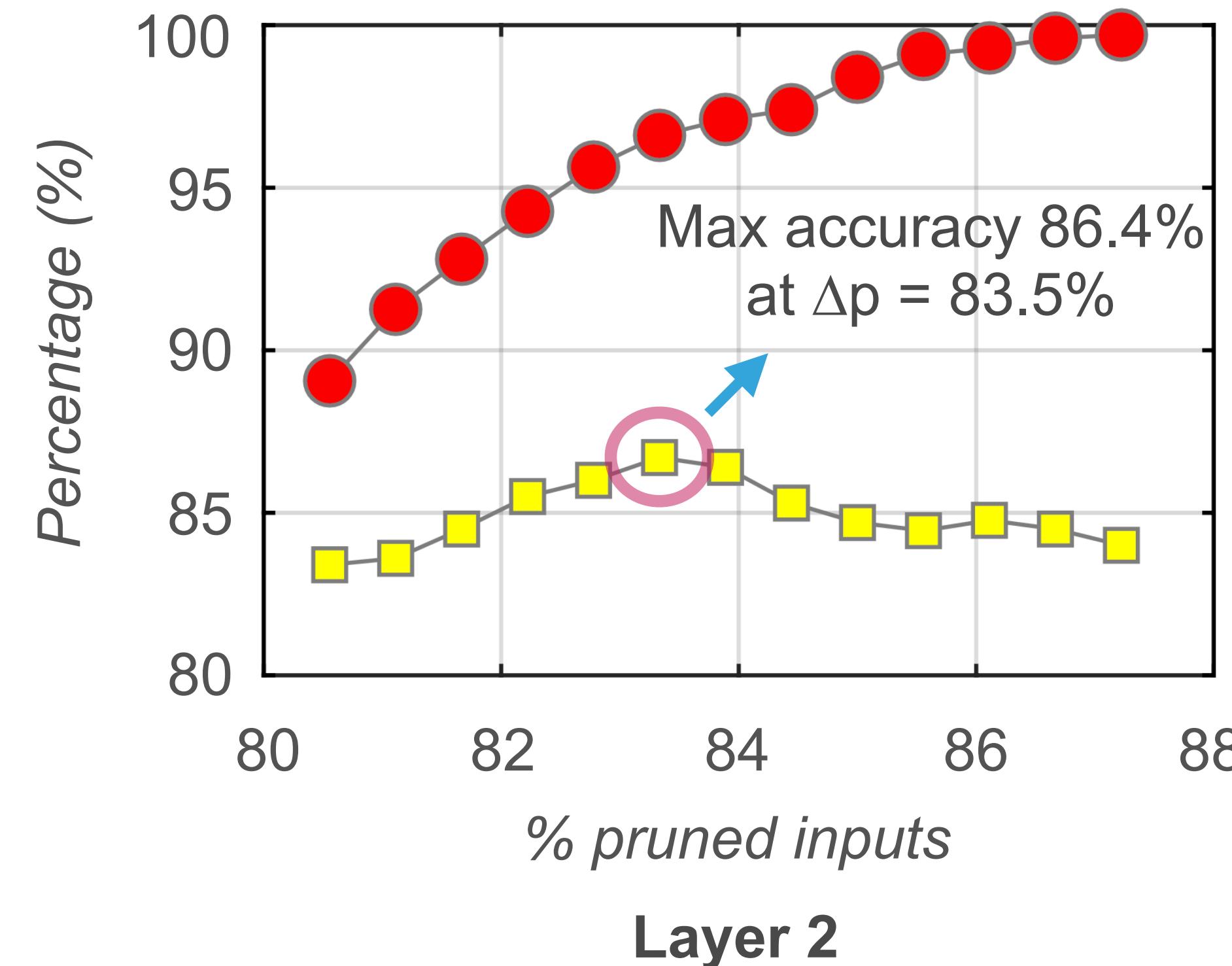
# EXPERIMENTAL SETUP



- Model: A CNN with binarized activation convolutional layers
- Target layers: 2nd and 3rd conv layers
- Dataset: CIFAR-10
- Validation: Last 5000 samples of training set as validation set after 200 epochs

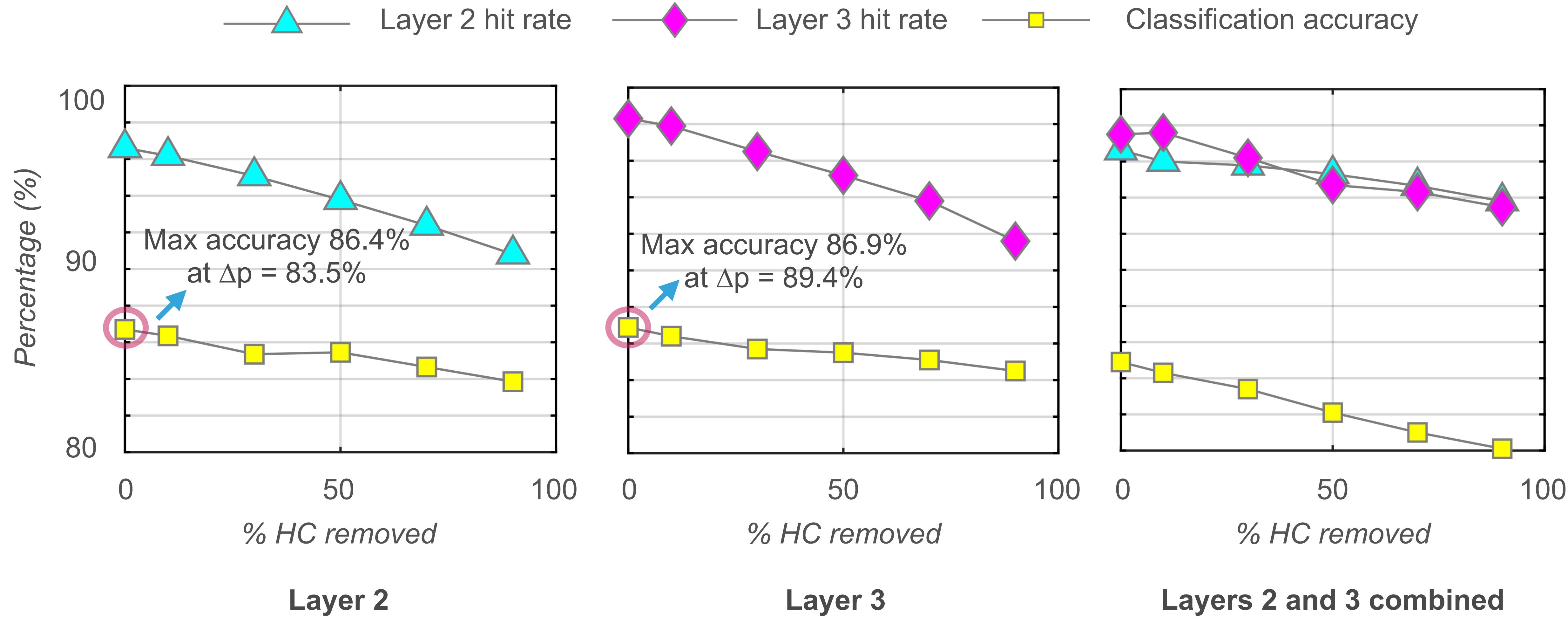
# RESULTS: SYNTHESIS-AWARE PRUNING

—●— Hit rate: Probability of non-X output during the inference  
—■— Classification accuracy



80%–90% inputs can be pruned, resulting in increased hit rate and accuracy

# RESULTS: INPUT-DRIVEN LOGIC MINIMIZATION



High accuracy maintained even with ~90% HC removed

# HARDWARE EVALUATION

Layer	CLNN (Proposed)			SCALE-Sim energy per image (nJ)	Energy reduction per image
	Energy per cycle (pJ)	Cycles per image	Energy per image (nJ)		
2	541	1,024	553	5,682	90.3%
3	249	256	64	11,360	99.4%

- Synthesis: Synopsys Design Compiler using TSMC 45 nm library
- Baseline: Systolic array simulated using SCALE-Sim

**90+% energy reduction per image**

[1] A. Samajdaret al., “SCALE-Sim: Systolic CNN accelerator simulator,” arXiv preprint arXiv:1811.02883, 2018

# CONCLUSION

- Identified **unexploited energy-saving opportunities** in combinational logic neural network implementation
  - Introduced **hitcount** as a unique importance metric of inputs in combinational logic neural network
- Proposed **two methods for reducing model size** in the context of logic synthesis
  - **Synthesis-aware pruning** and **input-driven logic minimization**
  - Exploits **error resilience** of neural networks to further reduce logic size and enhance scalability
- **90+%** power reduction while maintaining a competitive accuracy of **82%** on the CIFAR-10 dataset