

# Analog eFlash Drives AI/ML Acceleration

AIML-102-1

---

Presenter: Dave  
Eggleston Microchip/SST

# What's So Challenging About **Digital** Edge Inference?



Flash Memory Summit

## Want:

- Fast
- Power efficient
- Low cost
- Small



## Get (1 or more):

- Slow
- Power Hog
- Expensive
- Big



# What's So Challenging About **Digital** Edge Inference?



Flash Memory Summit

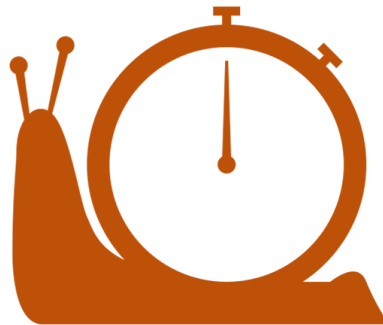
## Want:

- Fast
- Power efficient
- Low cost
- Small



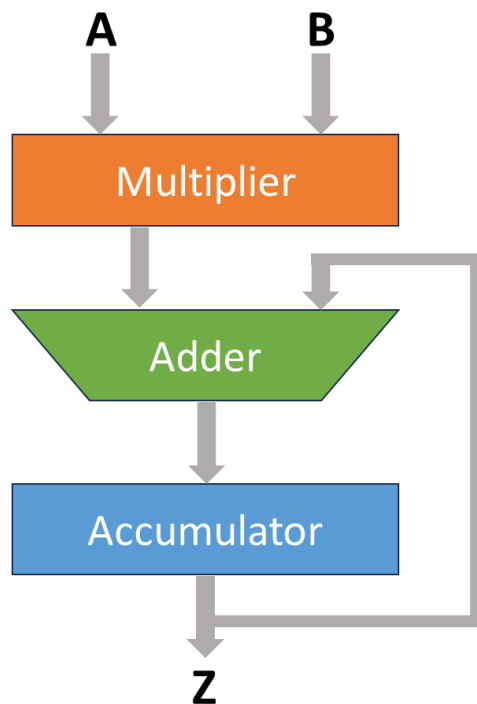
## Get (one or more):

- Slow
- Power hog
- Expensive
- Big





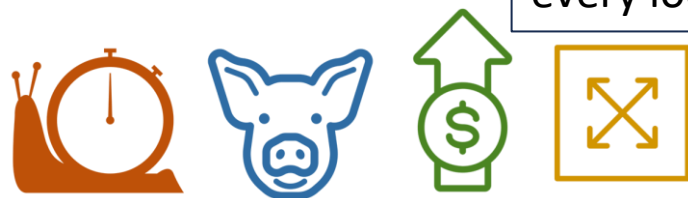
# Why? Multiply and Accumulate (MAC) Operation Analysis



Read A  
Read B  
Multiply  
Read Z  
Add Z  
Update (**write**) Z  
Repeat...  
Repeat...  
Repeat...  
...

Three memory **reads** and one memory **write** every loop

AlexNet → **3B** DRAM accesses!



Digital MAC



$$\text{Output } I_j = \sum_{t=1}^3 V_t \times G_{tj}$$

AnyNet → **Zero** DRAM accesses!

Write each array element with conductance **G**

Input V  
Output I  
DONE!  
Next layer...

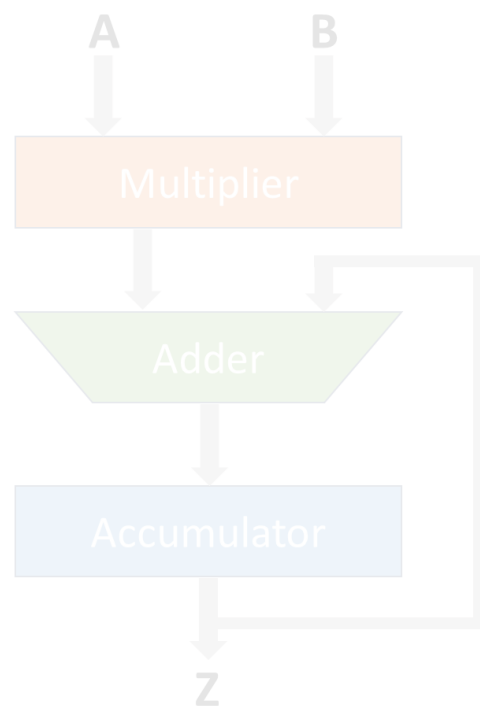
**Single cycle** MAC operations!



Analog MAC



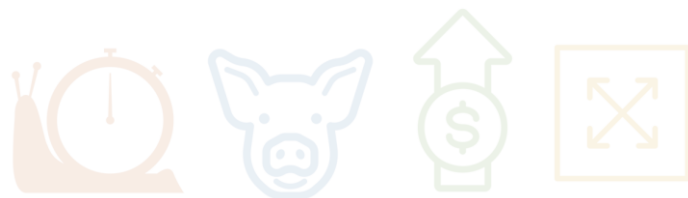
# Why? Multiply and Accumulate (MAC) Operation Comparison



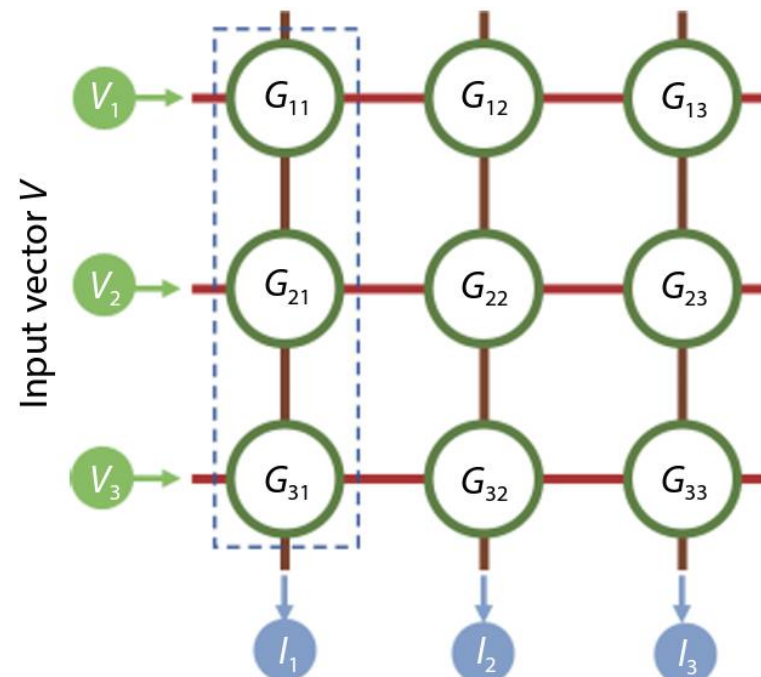
Read A  
Read B  
Multiply  
Read Z  
Add Z  
Update (Write) Z  
Repeat...  
Repeat...  
Repeat...  
...

3 memory Reads  
& 1 memory Write every loop

AlexNet → 3B DRAM accesses!



Digital MAC



$$\text{Output } I_j = \sum_{t=1}^3 V_t \times G_{tj}$$

AnyNet → **Zero** DRAM accesses!

**Write** each array element with conductance **G**

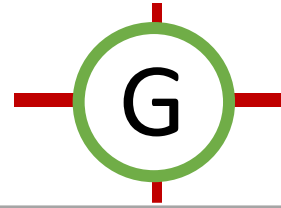
**Input V**  
**Output I**  
**DONE!**  
Next layer...

**Single cycle MAC** operations!



Analog MAC

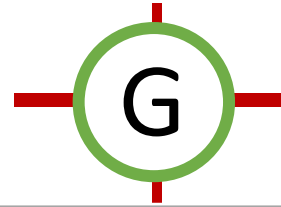
# eFlash is the *BEST* Choice for Analog MAC Array Element



	eFlash	eRRAM	eMRAM	ePCM	FeFET
Bits per cell capability	5 bpc+	2bpc	1bpc	2bpc	TBD
Foundry nodes available	>100	1	5	0	0
Production devices shipped	100B+	0+	1M?	N/A	N/A
Analog capable	YES	Maybe	No		
Analog IP available	YES	No	N/A		
Analog IP demos	YES	N/A			
Analog IP in production	YES				



# eFlash is the *BEST* Choice for Analog MAC Array Element



	eFlash	eRRAM	eMRAM	ePCM	FeFET
Bits per cell capability	5 bpc+	2 bpc	1 bpc	2 bpc	TBD
Foundry nodes available	>100	1	5	0	0
Production devices shipped	100B+	0+	1M?	N/A	N/A
Analog capable	YES	Maybe	No		
Analog IP available	YES	No	N/A		
Analog IP demos	YES	N/A			
Analog IP in production	YES				



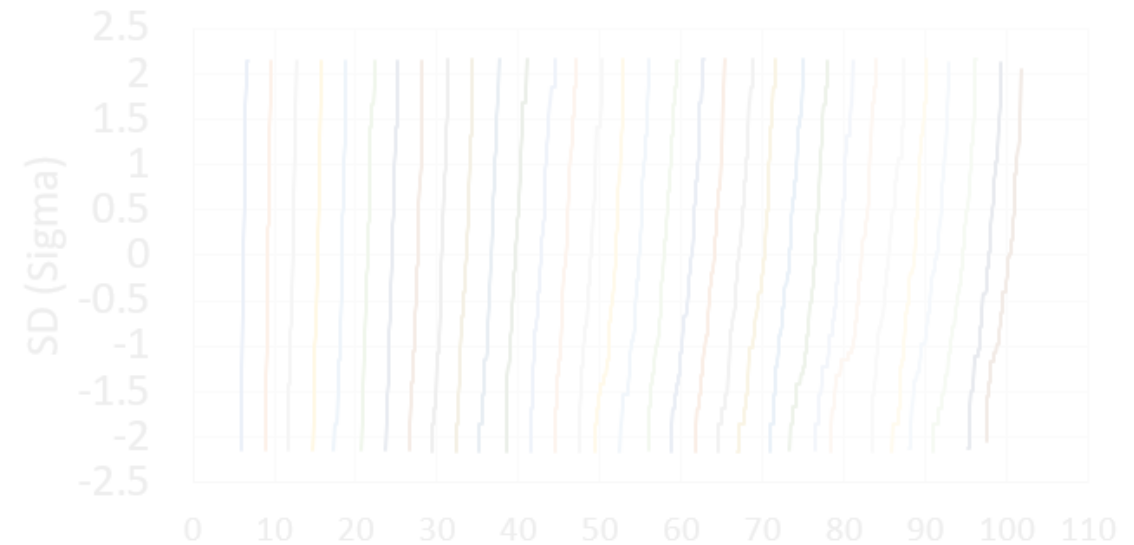
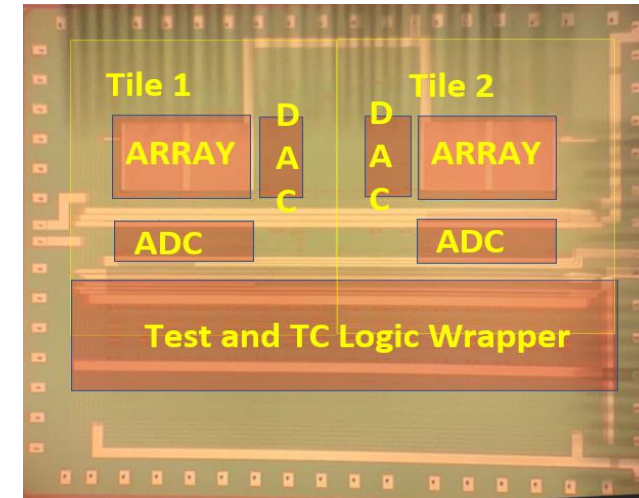


# eFlash 5 bpc (PLC) Demonstrated in 2X nm Silicon



Flash Memory Summit

- Silicon fully functional (2021)
  - All macro functions validated
  - 5-bits per cell = 32 levels
  - Read neural multiply operation realized
  - Neural net functionality realized over temp – MNIST ~95% (MLP), ~97% (CNN)
- 5 bpc MW and bit placement
  - 0 to 100 nA range
  - Sigma <1.3%
    - Single electron resolution



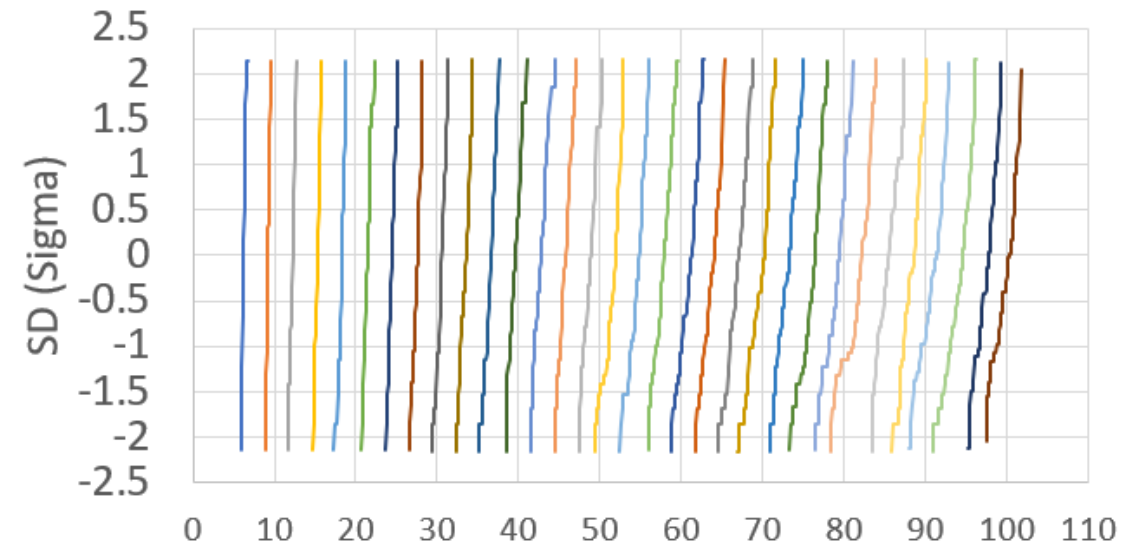
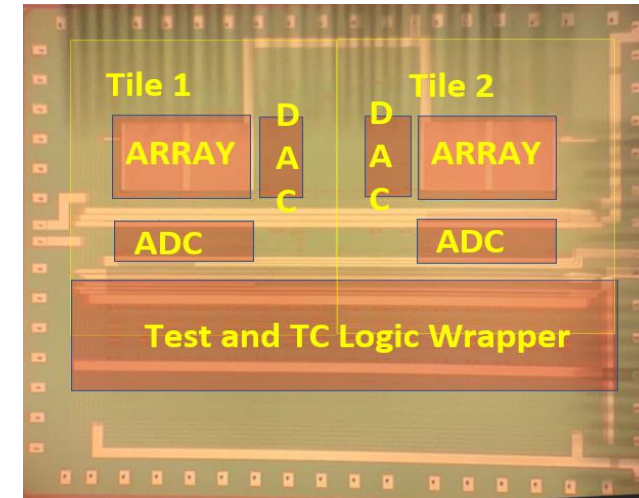


# eFlash 5 bpc (PLC) Demonstrated in 2X nm Silicon



Flash Memory Summit

- Silicon fully functional (2021)
  - All macro functions validated
  - 5-bits per cell = 32 levels
  - Read neural multiply operation realized
  - Neural net functionality realized over temp – MNIST ~95% (MLP), ~97% (CNN)
- 5 bpc MW and bit placement
  - 0 to 100 nA range
  - Sigma <1.3%
    - Single electron resolution



# Analog MACs Handle Large Weights with Ease!

- Why is using analog superior for edge inference acceleration?
  - Parallel in-memory compute the matrix vector product is very power efficient (no data movement)
  - Analog matrix vector product is computed in constant time
  - Large weight matrix takes the same time to compute the matrix vector product as a small weight

- Strategy divide and conquer!
  - Small tensor input; use Digital
  - Large tensor input; use Analog



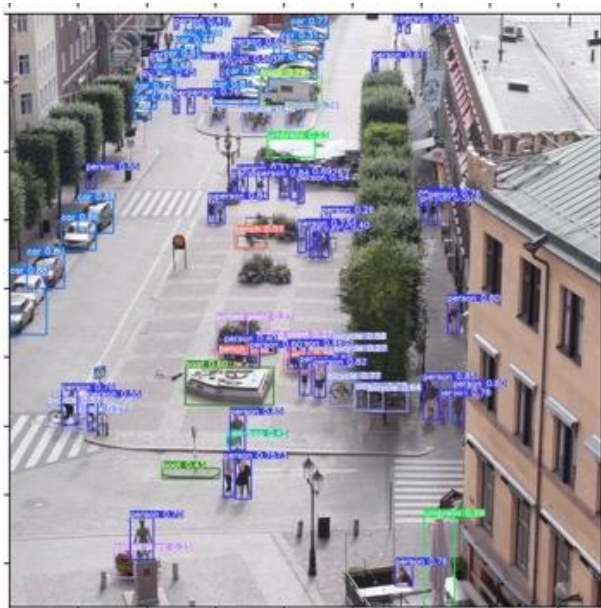
# Analog MACs Handle Large Weights with Ease!

- Why is using analog superior for edge inference acceleration?
  - Parallel in-memory compute the matrix vector product is very power efficient (no data movement)
  - Analog matrix vector product is computed in constant time
  - Large weight matrix takes the same time to compute the matrix vector product as a small weight
- Strategy divide and conquer!
  - Small tensor input; use **Digital**
  - Large tensor input; use **Analog**

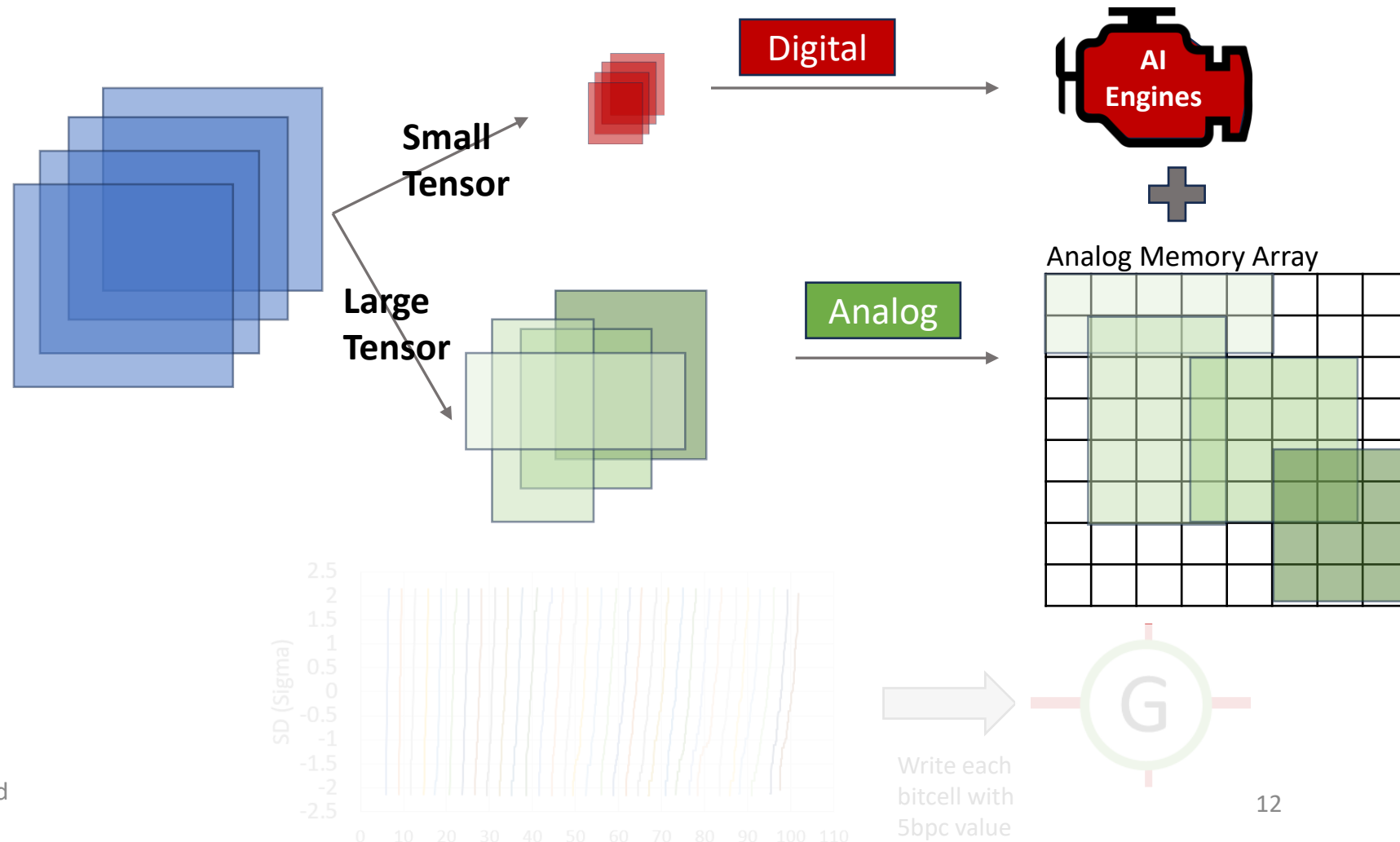


# Divide the Tensor Inputs Between **Analog** and **Digital**

- Example: Video object detection YOLOv5

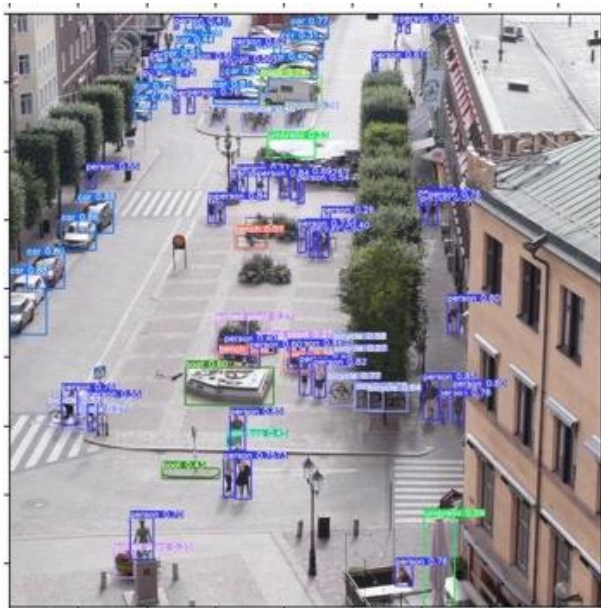


Security video outside a train station

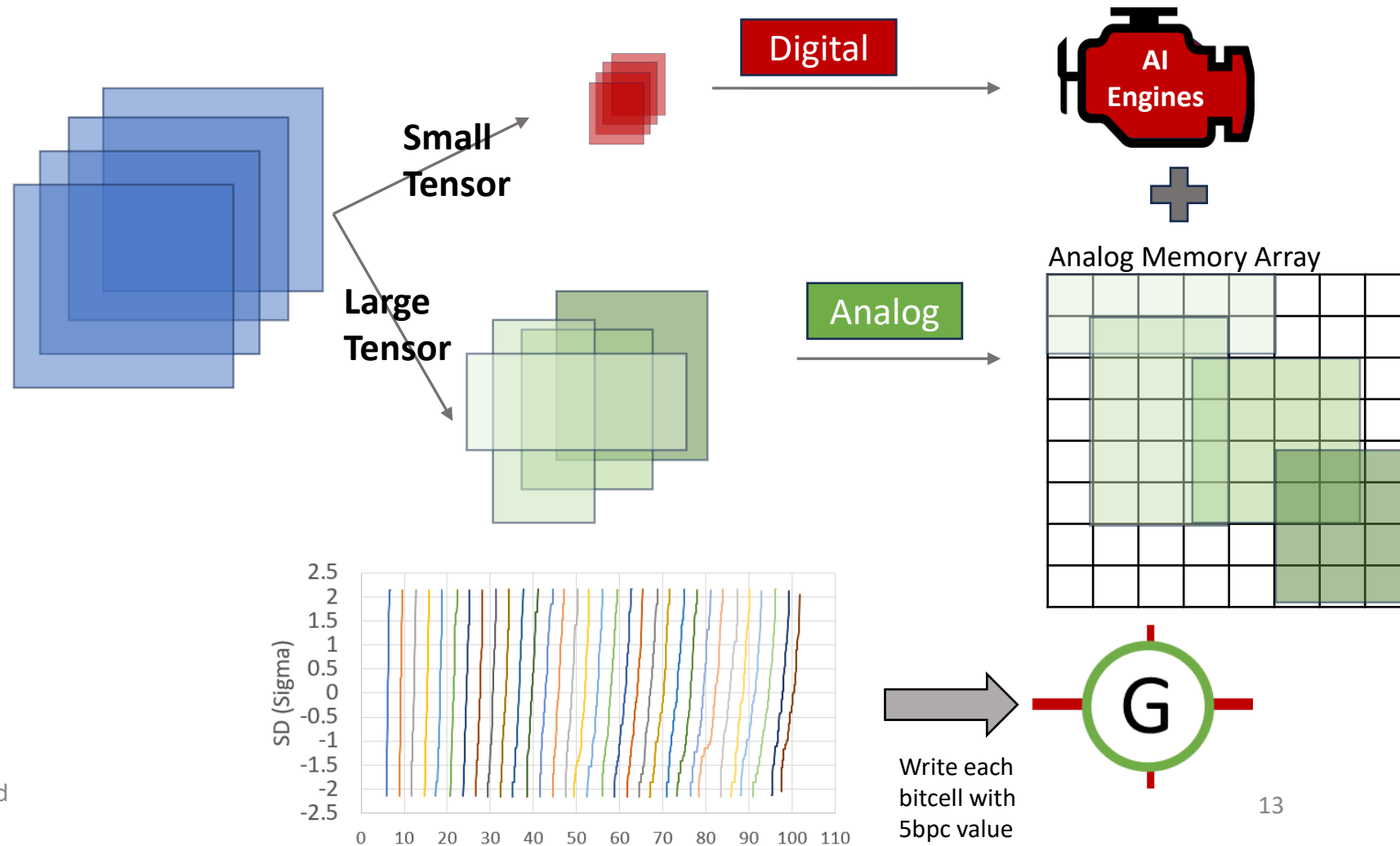


# Divide the Tensor Inputs Between **Analog** and **Digital**

- Example: Video object detection YOLOv5

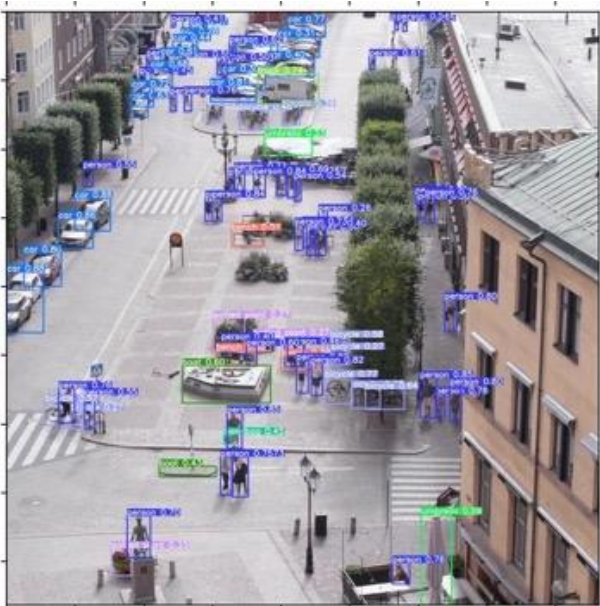


Security video outside a train station

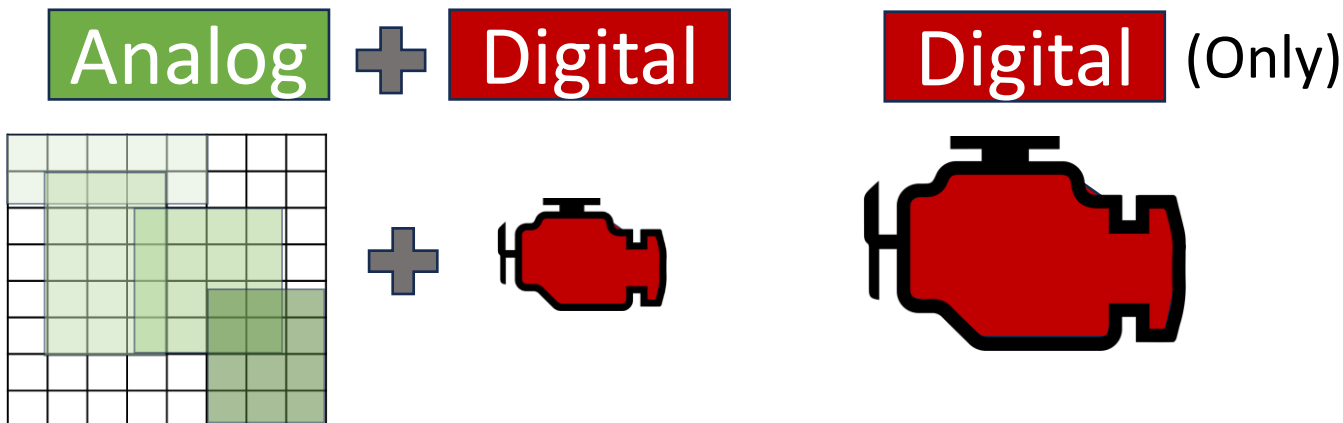


# Performance Difference: Analog vs. Digital

- Example: Video object detection YOLOv5



Security video outside a train station



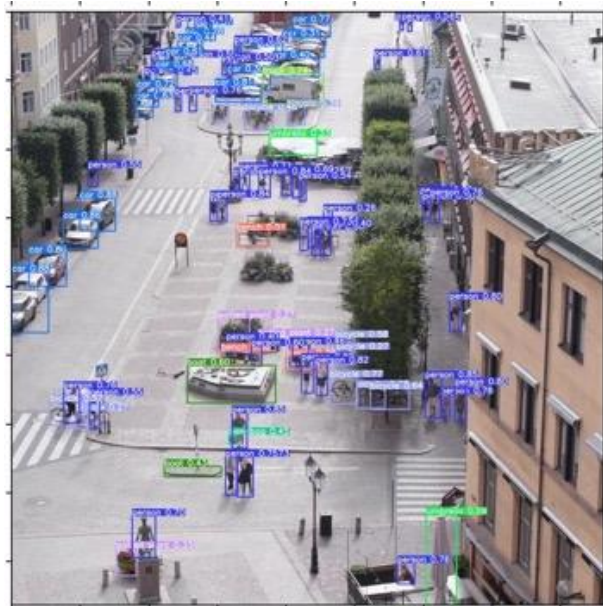
Difference

Die Cost	1/32	1	97%↓
Frames per sec	150 fps	200 fps	75%↓
Power	3W	60W	20x↓
FPS/W	50	3.33	15x↑
FPS/\$	3	0.125	24x↑

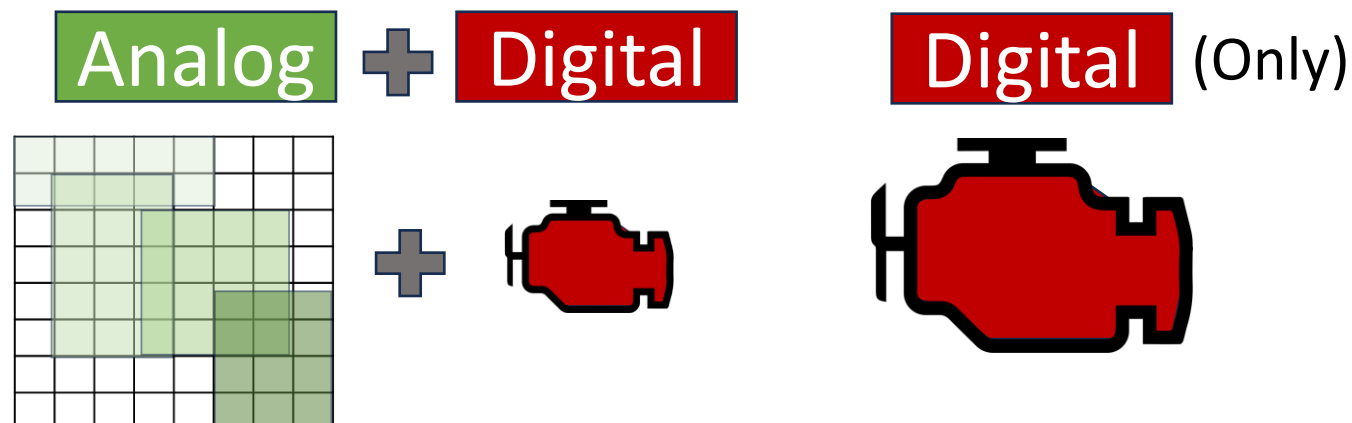


# Performance Difference: Analog vs. Digital

- Example: Video object detection YOLOv5



Security video outside a train station

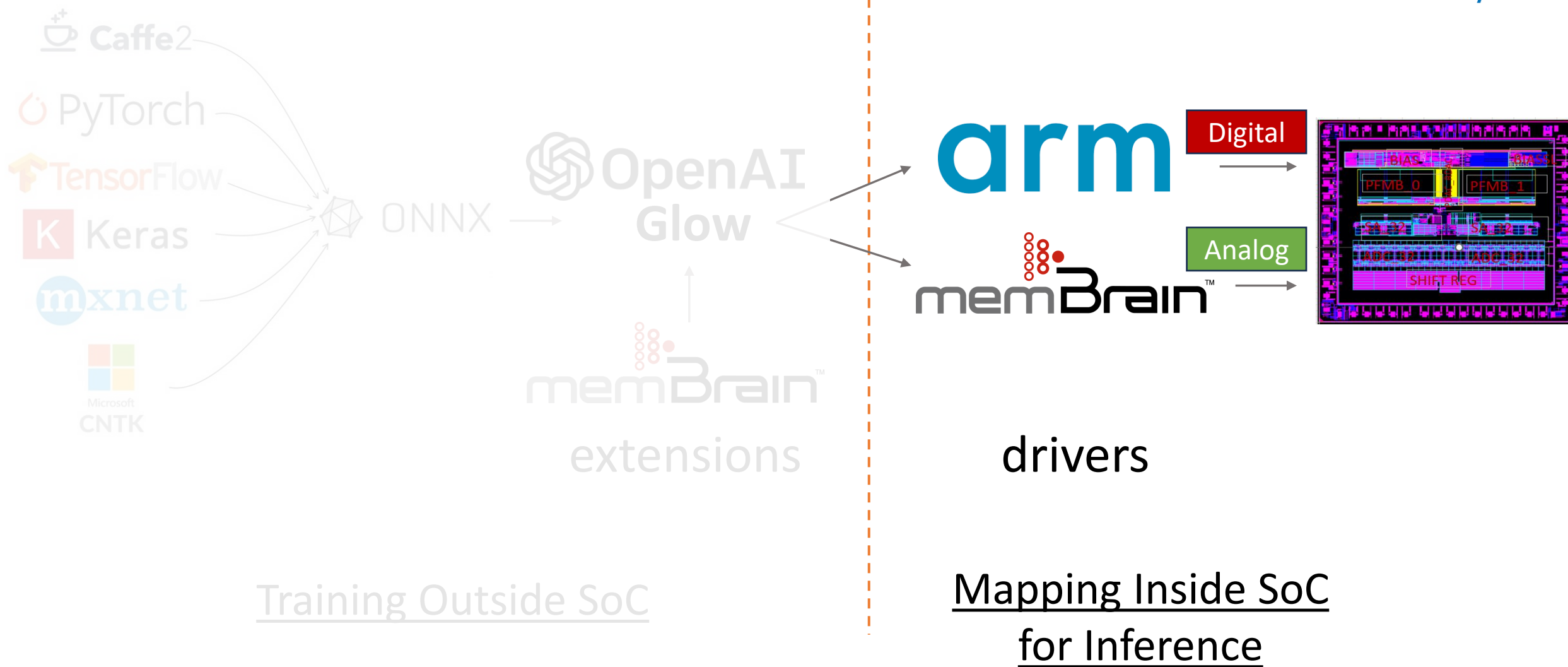


	Difference		
Die Cost	1/32	1	97%↓
Frames per sec	150 fps	200 fps	75%↓
Power	3W	60W	20x↓
FPS/W	50	3.33	15x↑
FPS/\$	3	0.125	24x↑





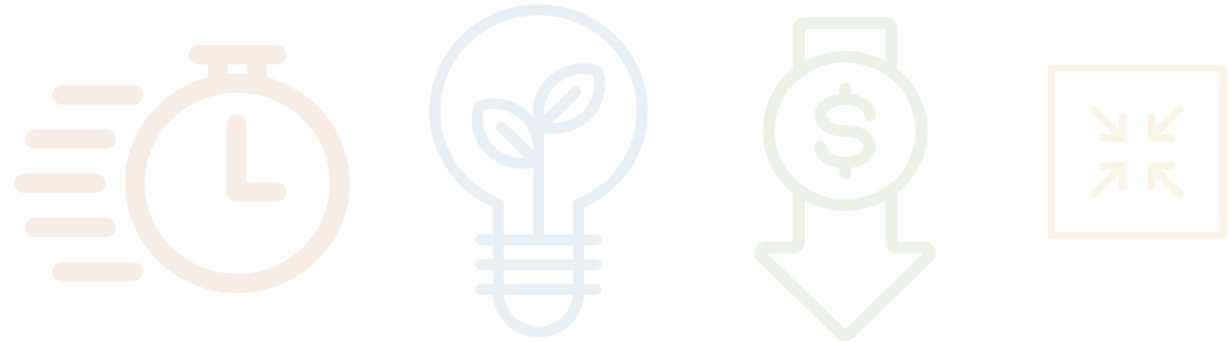
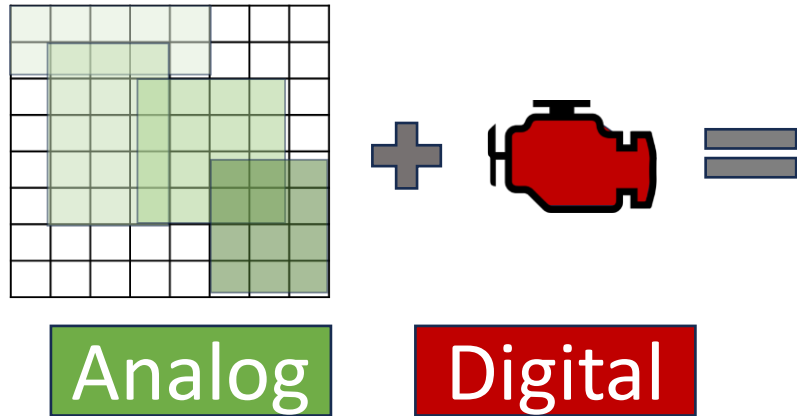
# Software Flow Example: eFlash Analog Hardware Aware



# Best of Analog + Best of Digital = Best Edge Inference Solution!



Flash Memory Summit



## Get:

- Fast
- Power efficient
- Low cost
- Small



# Best of **Analog** + Best of **Digital** = Best Edge Inference Solution!



Flash Memory Summit



Get:

- Fast
- Power efficient
- Low cost
- Small



Want to see and learn more?

- Visit our SST demo booth at the AI/ML Edge Hardware Summit
- September 12 – 14 in Santa Clara
- [aihwedgesummit.com](https://aihwedgesummit.com)



# AI HARDWARE & EDGE AI SUMMIT

# Thank You!



Flash Memory Summit

**190+ Employees**  
**More than 170 Engineers**

SST EU/Korea Sales:  
CBrown@SST.com

SST US/Japan Sales:  
DEggleston@SST.com

San Jose

Hermitage

UK, France, Italy

Ho Chi Minh City

Shanghai

Hsinchu

SST China/TW Sales:  
GChen@SST.com

