



# STAMP: Accelerating Second-order DNN Training Via ReRAM-based Processing-inMemory Architecture

Yilong Zhao (Speaker), Fangxin Liu\*, Mingyu Gao, Xiaoyao Liang, Qidong Tang,

Chengyang Gu, Tao Yang, Naifeng Jing, and Li Jiang\*

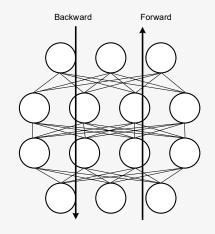
Shanghai Jiao Tong University, Shanghai Qi Zhi Institute, Tsinghua University

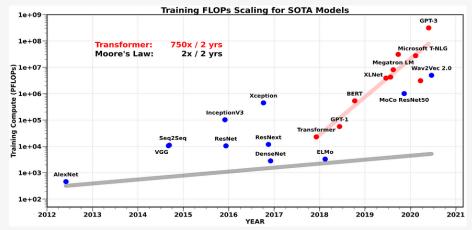
June 4, 2025

## **Background**



Training DNN models requires significant computation.

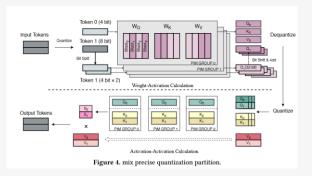




Gholami, et. al. Al and Memory Wall

© Current approaches to alleviate training burden.

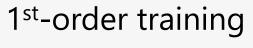
e.g.



Quantization Methods. (Introduce quantize/dequantize overhead.)



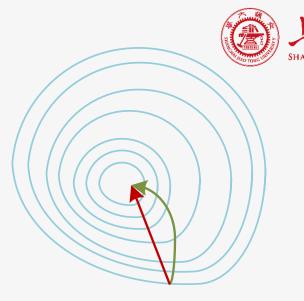
## **Background – Second-order Training**

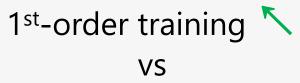


$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$$

2<sup>nd</sup>-order training 
$$\theta = \theta - \eta \cdot A^{-1} \nabla_{\theta} J(\theta) G^{-1}$$

Utilize the **Second-order Information (SOI) Matrix's Inversion** for a more accurate step and direction
-> faster convergence





2<sup>nd</sup>-order training \(^\) **FEWER STEPS!** 



## **Background – Second-order Training**

1<sup>st</sup>-order training

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$$

2<sup>nd</sup>-order training  $\theta = \theta - \eta \cdot A^{-1} \nabla_{\theta} J(\theta) G^{-1}$ 

Utilize the Second-order Information (SOI) Matrix's Inversion for a more accurate step and direction -> faster convergence

#### Why it is not widely used?

SOI matrix brings too much computation/storage Loss  $O(n^2)$  overhead on **GPUs**,  $O(n^2)$  storage,  $O(n^3)$  computing complexity, -> longer step time



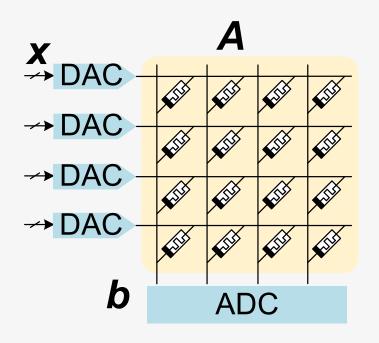
1st-order training \( \scalen\) vs

2<sup>nd</sup>-order training \(^\) **FEWER STEPS!** 



## **Background – ReRAM Accelerator**





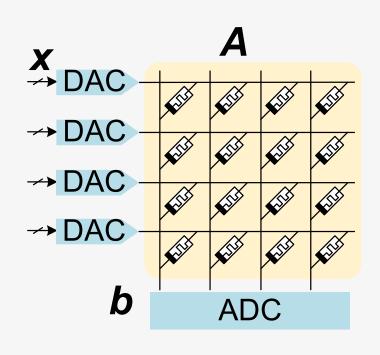
$$b = x \cdot A$$

ReRAM-based Vector-Matrix Multiplication (**VMM**)



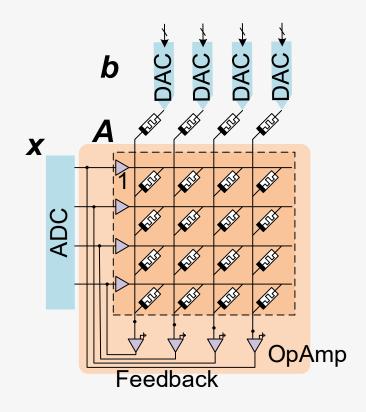
## **Background – ReRAM Accelerator**





$$b = x \cdot A$$

ReRAM-based Vector-Matrix Multiplication (**VMM**)



$$x = b \cdot A^{-1}$$

ReRAM-based Matrix Inversion (**INV**)

#### **Motivation**



#### Use ReRAM+2<sup>nd</sup>-order training to accelerate DNN training!

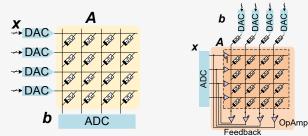
#### 2<sup>nd</sup>-order training

$$\theta = \theta - \eta \cdot A^{-1} \nabla_{\theta} J(\theta) G^{-1}$$

#### **SOI Matrix Overhead**

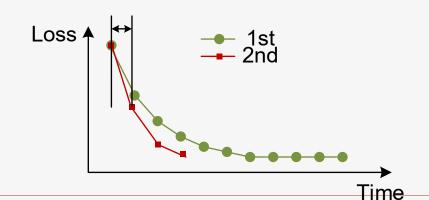
 $O(n^2)$  storage overhead  $O(n^3)$  computing complexity

#### **ReRAM-based Accelerator**



#### **ReRAM's Feature**

- → Process in Memory
- $\rightarrow O(1)$  Matrix Inversion Time

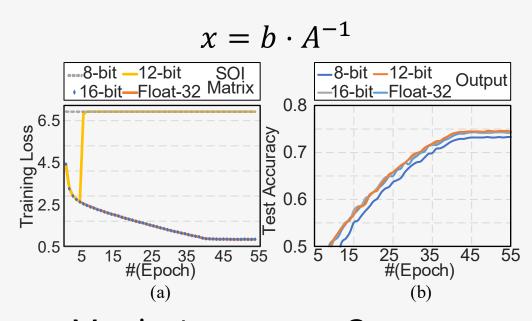


#### Aim:

Fewer Steps without apparently enlarging step time

-> Faster training





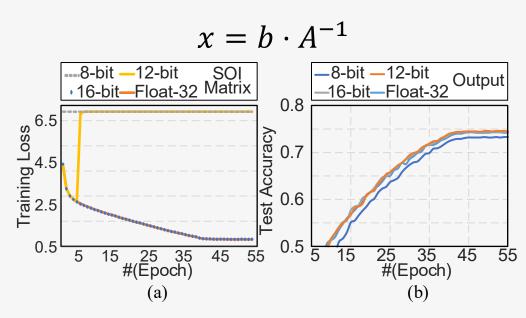
Matrix *A*: At least **16** bits

Output *x*: At least **12** bits

2<sup>nd</sup>-order training requires **HIGH**-precision Matrix Inversion



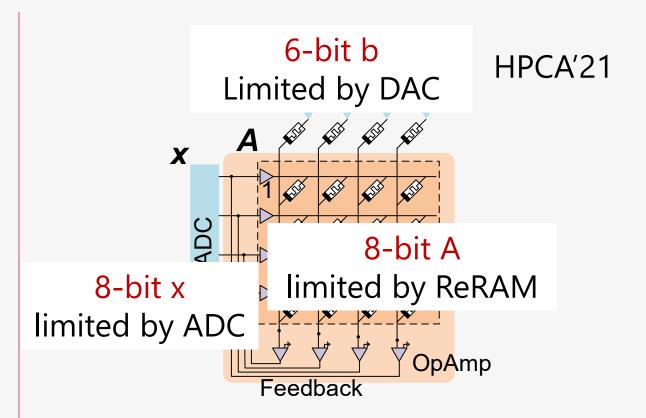




Matrix *A*: At least **16** bits

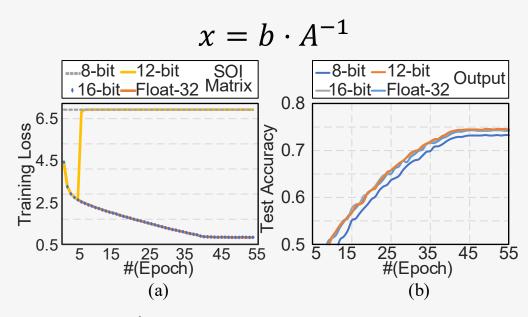
Output *x*: At least **12** bits

2<sup>nd</sup>-order training requires **HIGH**-precision Matrix Inversion





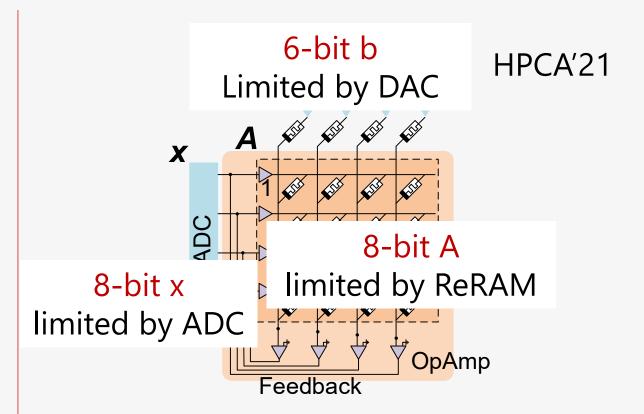




Matrix *A*: At least **16** bits

Output *x*: At least **12** bits

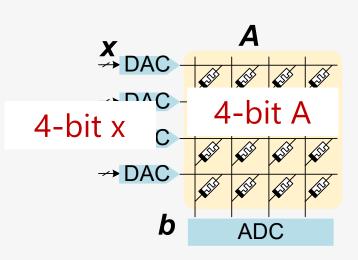
2<sup>nd</sup>-order training requires **HIGH**-precision Matrix Inversion



ReRAM-based Matrix Inversion's precision is **not enough** 



• How this problem solved in ReRAM-based Vector-Matrix Multiplication

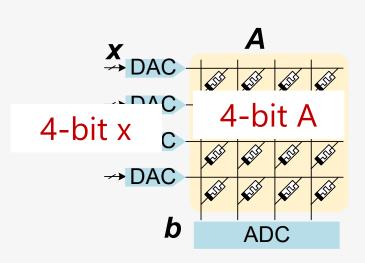


We need 8-bit x, 8-bit A? HOW?

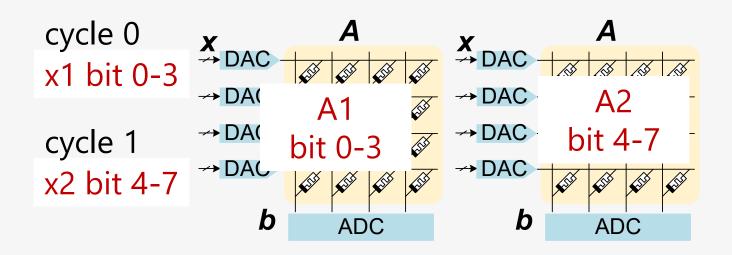




• How this problem solved in ReRAM-based Vector-Matrix Multiplication



We need 8-bit x, 8-bit A? HOW?

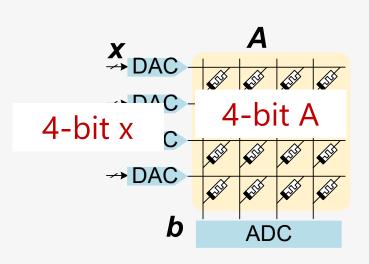


$$b = x1 \cdot A1 \cdot 2^{-8} + x1 \cdot A2 \cdot 2^{-4} + x2 \cdot A1 \cdot 2^{-4} + x2 \cdot A2$$

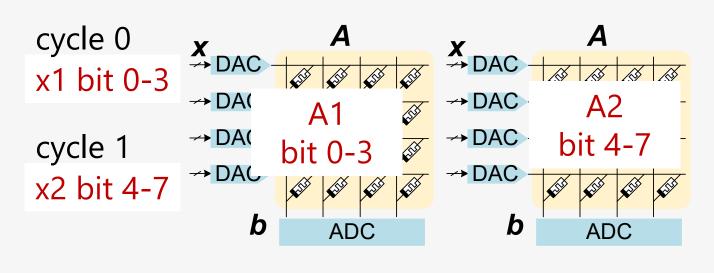




#### • How this problem solved in ReRAM-based Vector-Matrix Multiplication



We need 8-bit x, 8-bit A? HOW?



$$b = x1 \cdot A1 \cdot 2^{-8} +$$

$$x1 \cdot A2 \cdot 2^{-4} +$$

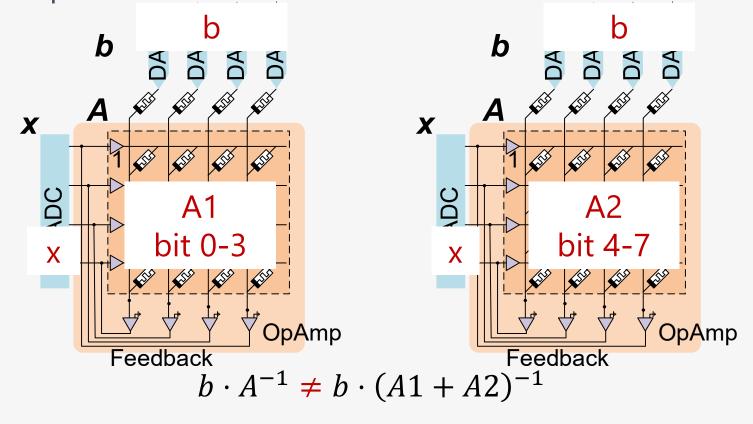
$$x2 \cdot A1 \cdot 2^{-4} +$$

$$x2 \cdot A2$$

Traditional **BIT SLICE SCHEME** 



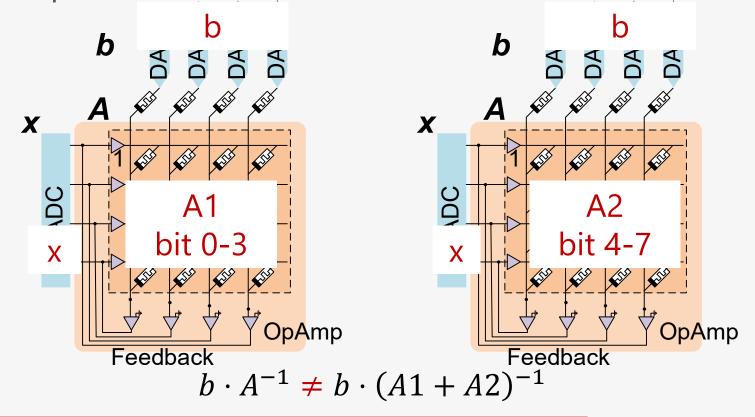
It is more complicated in ReRAM-based Matrix Inversion







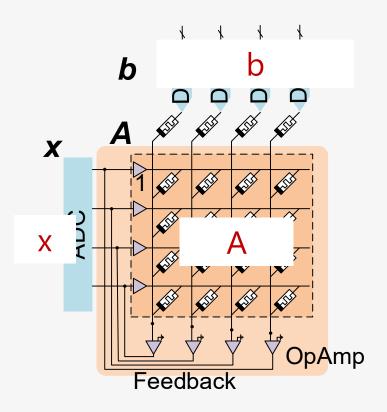
It is more complicated in ReRAM-based Matrix Inversion



For Matrix Inversion, Traditional Bit Slice Scheme Doesn't Work! Because Matrix Inversion doesn't have distributive law



92 Levels: x, A

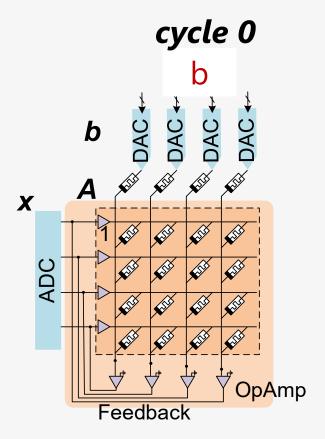






2 Levels: x, A

Level *x* 16-bit ->2\*8 bits

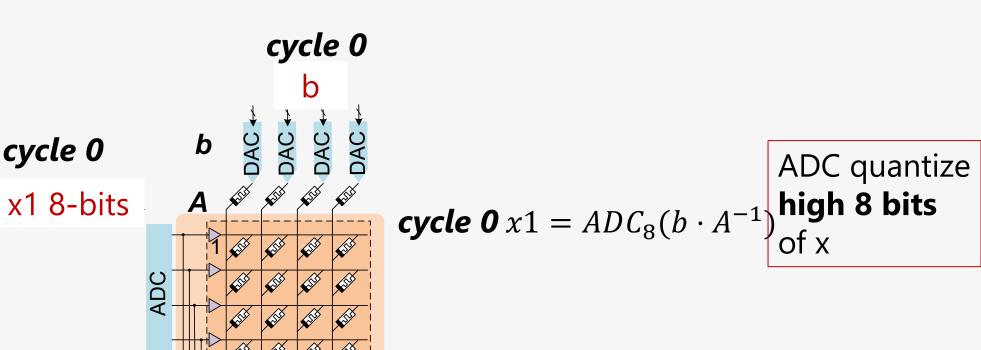


cycle 0



@2 Levels: x, A

Level *x* 16-bit ->2\*8 bits

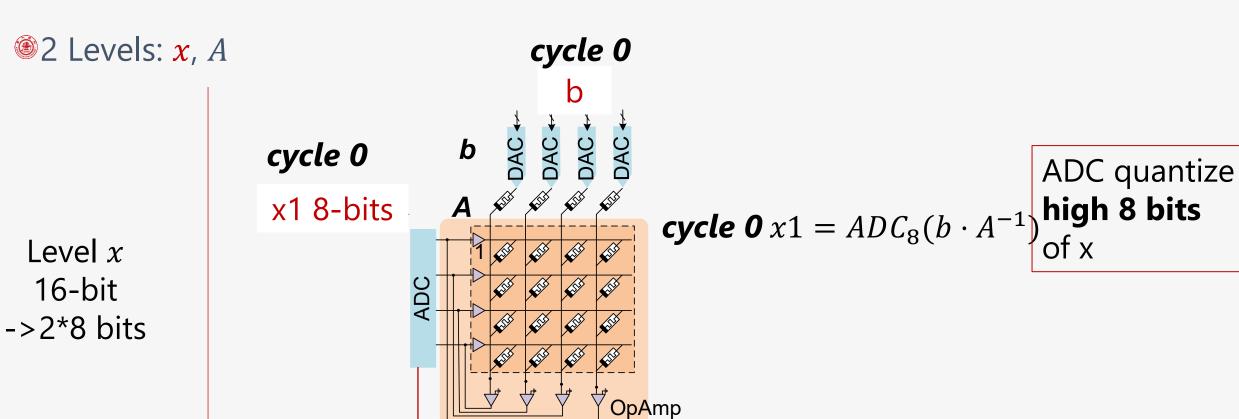


**OpA**mp

Feedback





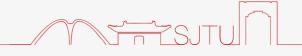


In Analog Field:

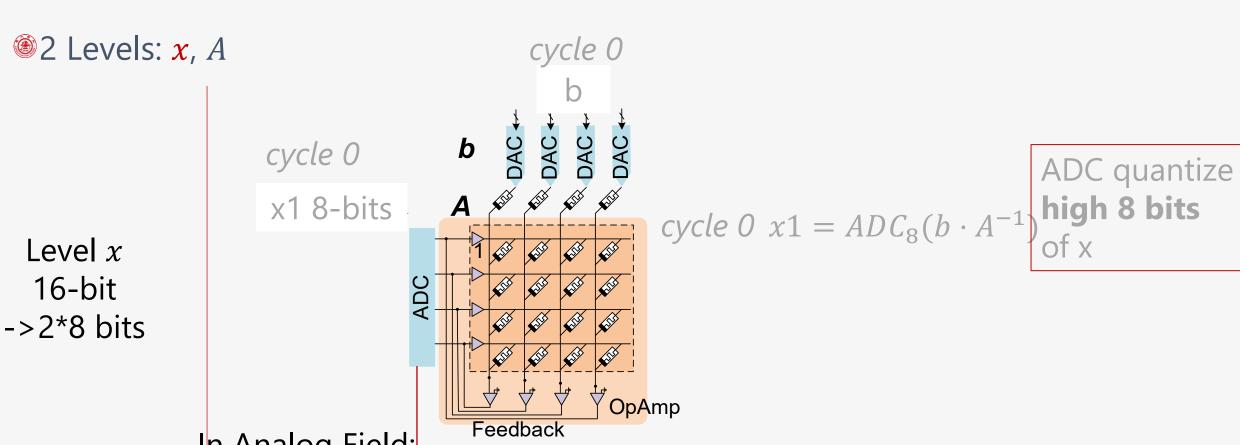
$$x1 + x2 = b \cdot A^{-1}$$

Feedback

quantized omitted (lower part)







In Analog Field:

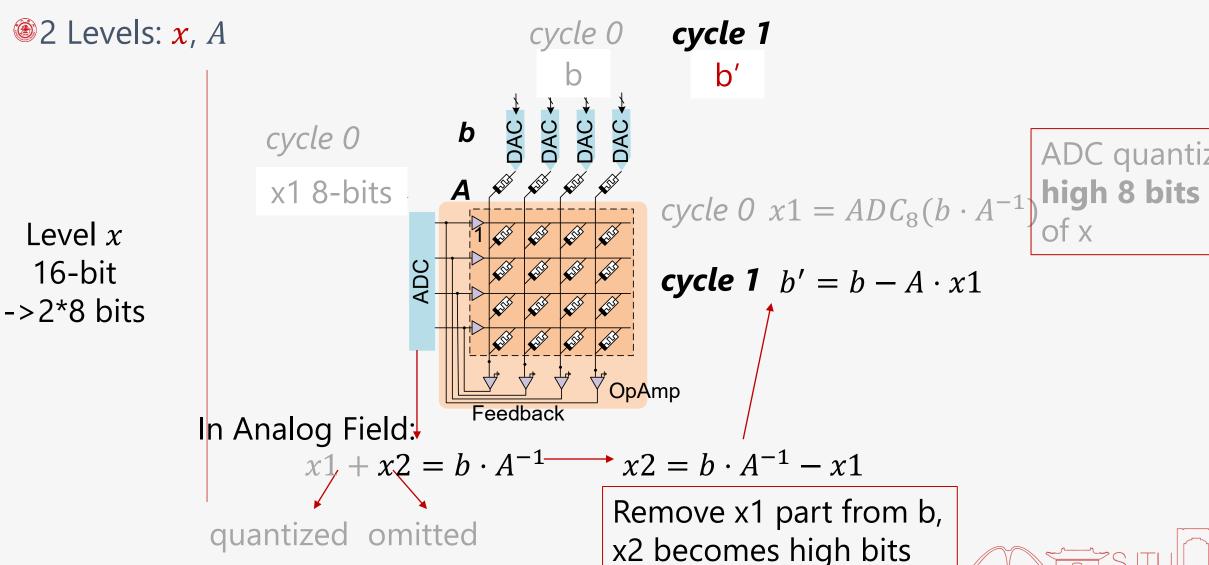
$$x1 + x2 = b \cdot A^{-1} \qquad x2 = b \cdot A^{-1} - x1$$

quantized omitted

Remove x1 part from b, x2 becomes high bits



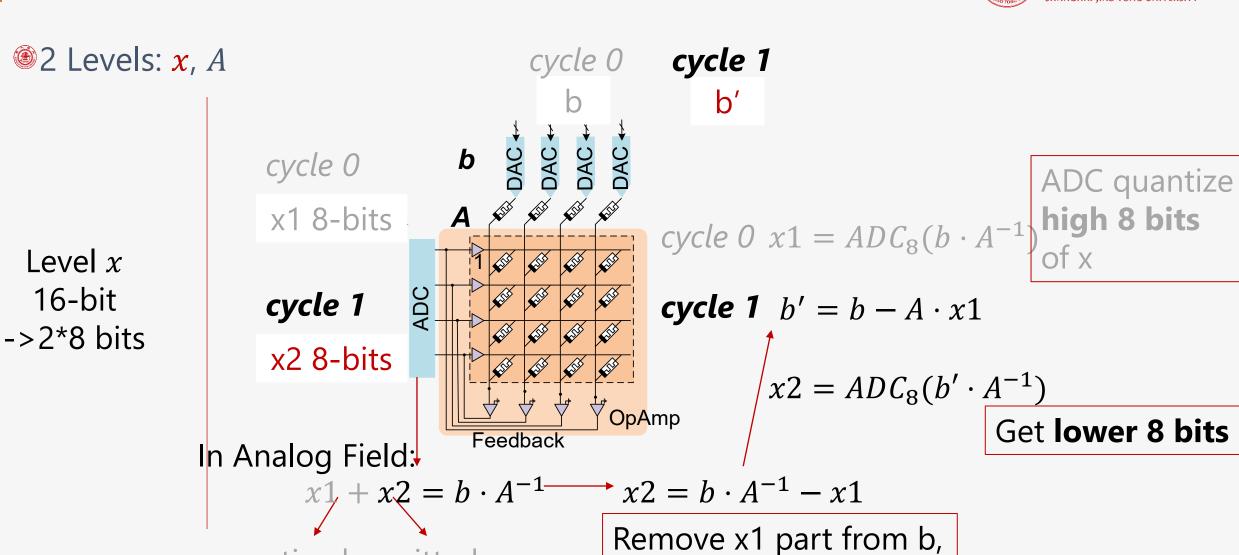




ADC quantize

quantized omitted





x2 becomes high bits



2 Levels: x, A

Level *A*16-bit
-> 2\*8 bits

8 bit  $A_H$  & 8 bit  $A_L$ 

**Taylor Expansion:** 





2 Levels: x, A

Level *A* 16-bit ->2\*8 bits

8 bit *A<sub>H</sub>* & 8 bit *A<sub>L</sub>* 

**Taylor Expansion:** 

$$A^{-1} \cdot b = (A_H + A_L)^{-1} \cdot b$$
  
=  $A_H^{-1} \cdot (I - P + P^2 - P^3 + \cdots) \cdot b$   
 $(P = A_H^{-1} \cdot A_L)$ 





2 Levels: x, A

Level *A* 16-bit ->2\*8 bits

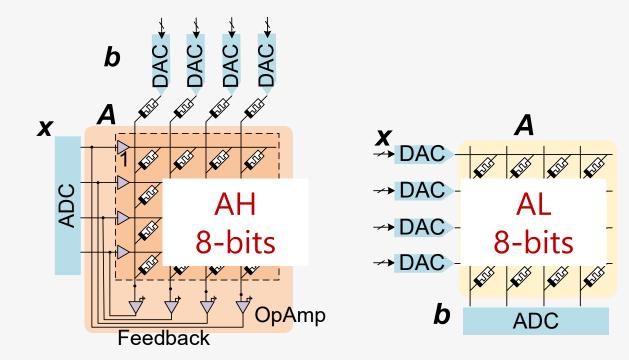
8 bit *A<sub>H</sub>* & 8 bit *A<sub>L</sub>* 

**Taylor Expansion:** 

$$A^{-1} \cdot b = (A_H + A_L)^{-1} \cdot b$$

$$= A_H^{-1} \cdot (I - P + P^2 - P^3 + \cdots) \cdot b$$

$$(P = A_H^{-1} \cdot A_L)$$



init: 
$$r = A_H^{-1} \cdot b$$
,  $x = r$ 





Level *A* 16-bit ->2\*8 bits

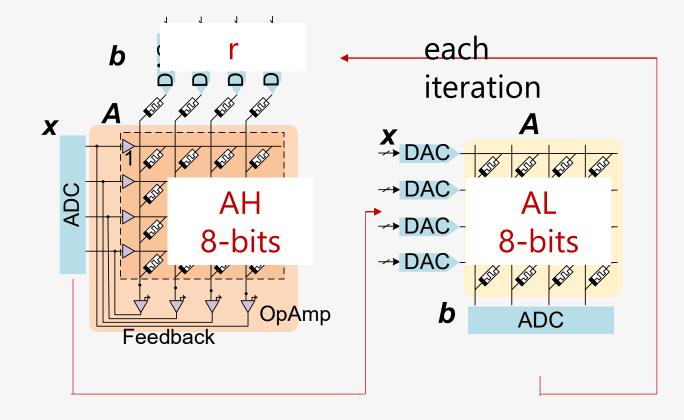
8 bit *A<sub>H</sub>* & 8 bit *A<sub>L</sub>* 

**Taylor Expansion:** 

$$A^{-1} \cdot b = (A_H + A_L)^{-1} \cdot b$$

$$= A_H^{-1} \cdot (I - P + P^2 - P^3 + \cdots) \cdot b$$

$$(P = A_H^{-1} \cdot A_L)$$





# calculate one term in each iteration

init: 
$$r = A_H^{-1} \cdot b$$
,  $x = r$   
Iteration 0:

$$\begin{array}{cc} \boldsymbol{P} & r = A_H^{-1} A_L r \\ x = r \end{array}$$





Level *A* 16-bit ->2\*8 bits

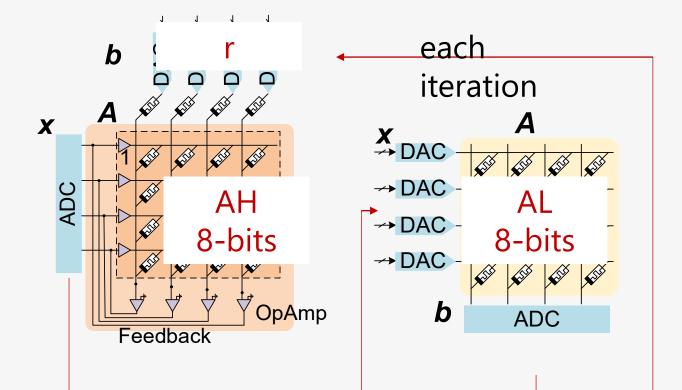
8 bit *A<sub>H</sub>* & 8 bit *A<sub>L</sub>* 

**Taylor Expansion:** 

$$A^{-1} \cdot b = (A_H + A_L)^{-1} \cdot b$$

$$= A_H^{-1} \cdot (I - P + P^2 - P^3 + \cdots) \cdot b$$

$$(P = A_H^{-1} \cdot A_L)$$





## calculate one term in each iteration

init: 
$$r = A_H^{-1} \cdot b$$
,  $x = r$   
Iteration 0:

$$\begin{array}{ccc}
\mathbf{P} & r = A_H^{-1} A_L r \\
x & -= r
\end{array}$$

Iteration 1:

$$\begin{array}{cc} \boldsymbol{P^2} & r = A_H^{-1} A_L r \\ & x \mathrel{+}= r \end{array}$$

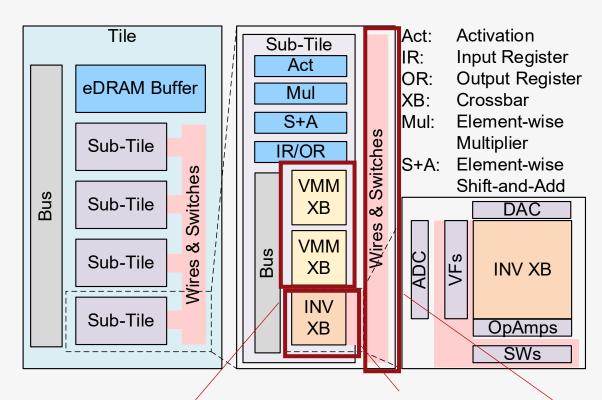
Iteration 2:

 $P^3$ 

x becomes more and more accurate

#### **STAMP Architecture**





For Vector-matrix multiplication

For Matrix Inversion

Connect INV crossbars for scaling

## **Evaluation**

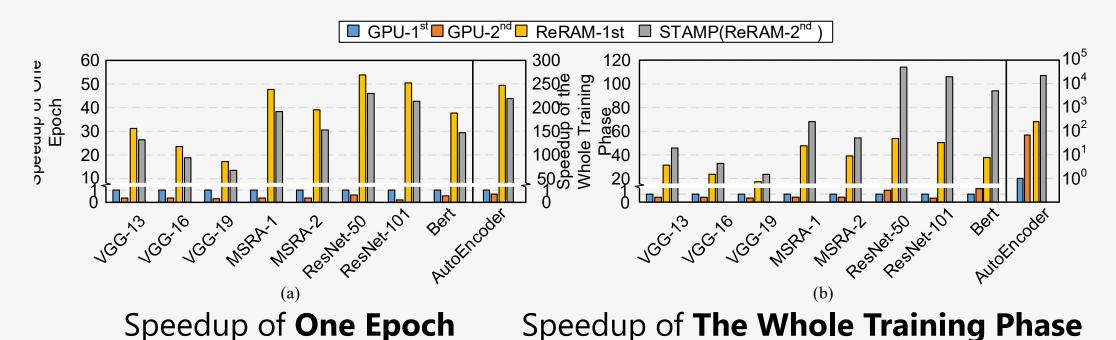


Baselines	Tesla V100 GPU, ReRAM-based 1st-order training	
Benchmark	VGG-13/16/19, MSRA, ResNet, Bert, autoencoder	
Simulation Model	DAC ADC Hyper	ISCA'16
	ОрАМР	HPCA'21
	DRAM buffer	CACTI 7
Architecture Setup	VMM:INV	16:1
	Cycle	100ns
	Crossbar size	256



#### **Performance**



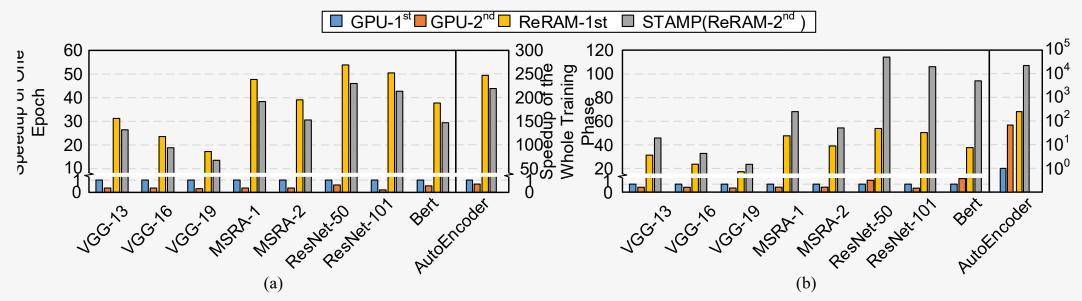


STAMP is 68x better than GPU-1<sup>st</sup> Compared to ReRAM-1<sup>st</sup>, Epoch +21.5%, Overall Training 11.4x Faster



#### **Performance**





Speedup of **One Epoch** 

Speedup of The Whole Training Phase

STAMP is 68x better than GPU-1<sup>st</sup> Compared to ReRAM-1<sup>st</sup>, Epoch +21.5%, Overall Training 11.4x Faster

- On ReRAM, 2<sup>nd</sup>-order training always better than 1<sup>st</sup>-order training
- On GPU, 2<sup>nd</sup>—order training is worse than 1<sup>st</sup>-order training due to SOI overhead

## **Reduced ReRAM writing**



STAMP can **reduce 55.7% write** number compared to first-order training on ReRAM-based accelerator, as there are fewer epochs.

**Enhance endurance of ReRAM Accelerator** 



## **Summary**



We use 2<sup>nd</sup>-order training & ReRAM-based architecture to accelerate DNN training.

We propose a high-precision matrix inversion algorithm based on low-precision ReRAM circuits.



