

# Temporal and SFQ Pulse-Streams Encoding for Area-Efficient Superconducting Accelerators

Lead author: Patricia Gonzalez-Guerrero

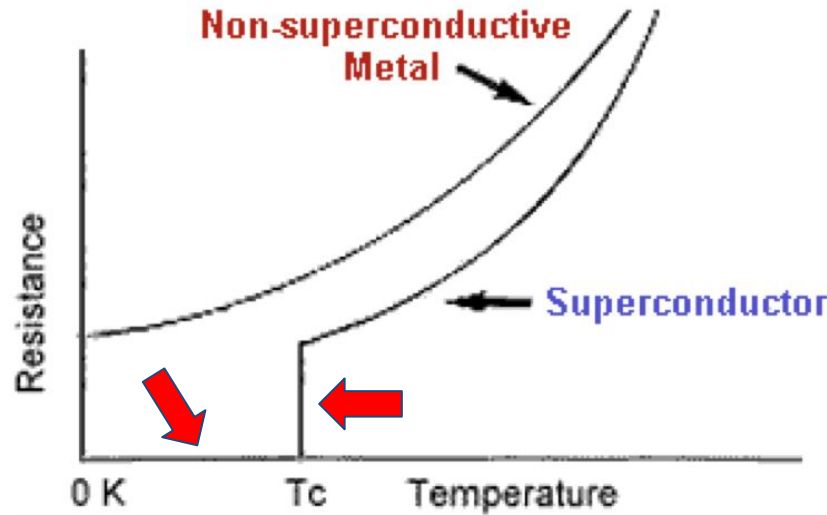
Speaker: George Micheliogiannakis

Computer Architecture Group  
Computational Research Division

ASPLOS 2022

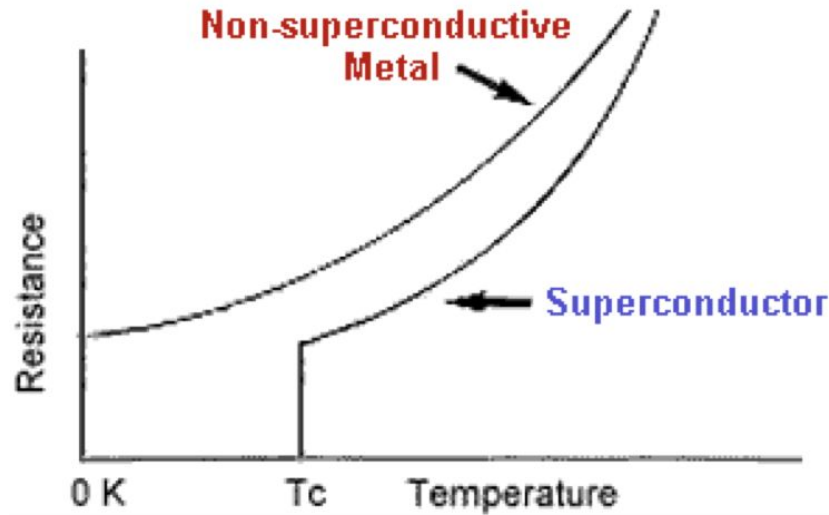


# Superconductivity is a good candidate for the future of computing

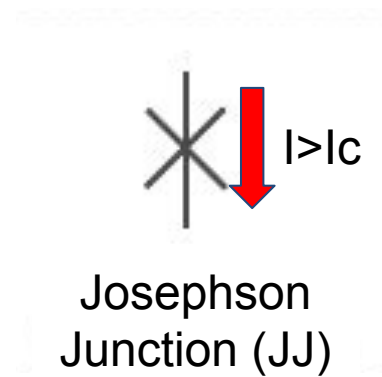


Gallardo et al, "Superconductivity observation in a  $(\text{CuInTe}_2)_{1-x}(\text{NbTe})_x$  alloy with  $x=0.5$ "

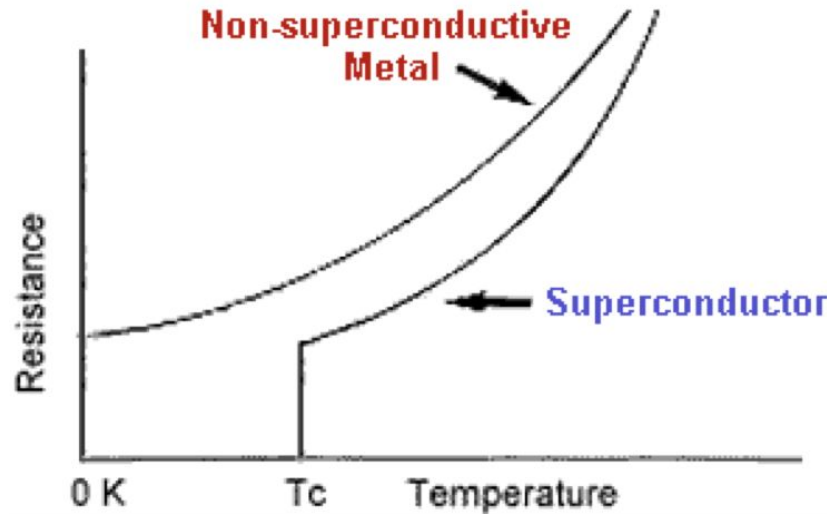
# Superconductivity is a good candidate for the future of computing



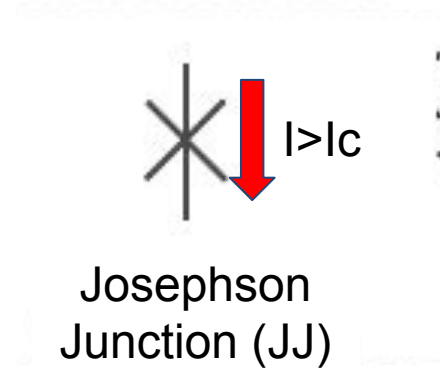
Gallardo et al, "Superconductivity observation in a  $(\text{CuInTe}_2)_{1-x}(\text{NbTe})_x$  alloy with  $x=0.5$ "



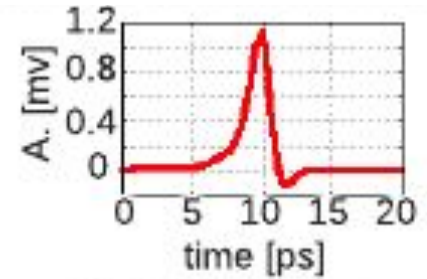
# Superconductivity is a good candidate for the future of computing



Gallardo et al, "Superconductivity observation in a  $(\text{CuInTe}_2)_{1-x}(\text{NbTe})_x$  alloy with  $x=0.5$ "

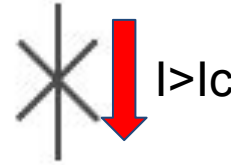


Single Flux Quantum (SFQ) pulse



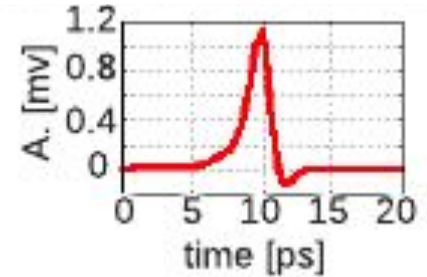
# Superconductivity is a good candidate for the future of computing

- Logic Families:
  - Rapid Single Flux Quantum (RSFQ)
  - Energy Efficient Rapid Flux Quantum (ERSFQ)
- 10X-100X faster operation frequency than CMOS.
  - 20GHz - 80GHz
- ~10X lower active energy consumption than CMOS.



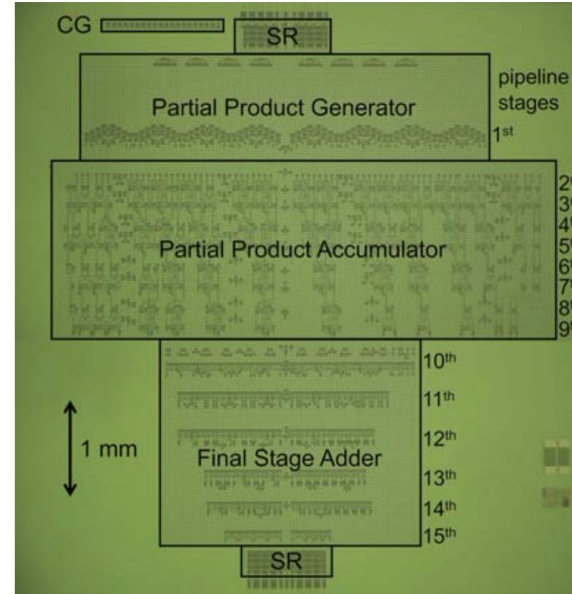
Josephson Junction (JJ)

Single Flux Quantum (SFQ) pulse



# Is superconductivity a good candidate for the future of computing, though?

- Extremely constrained area density.
- Manufacturing limitations restrict the number of JJs per chip to ~30000 JJs.
- Current hardware accelerators such as Google's TPU require up to ~64000 multipliers.

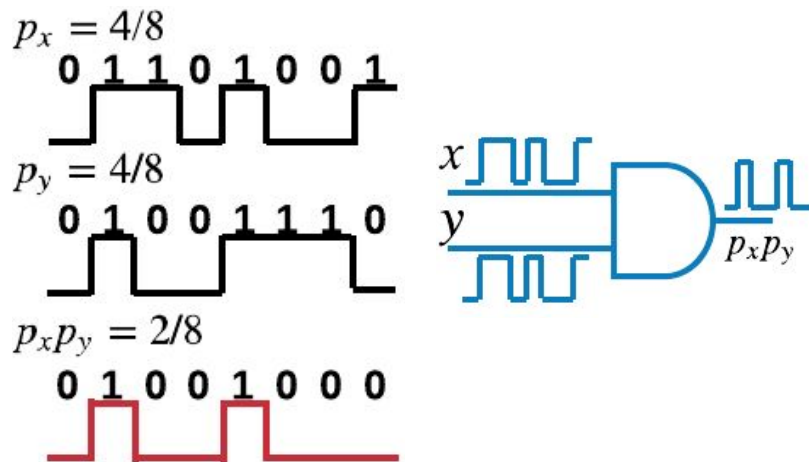


State of  
the art  
RSFQ  
Binary  
Multiplier  
~17000  
JJs

ISSCC'2019.

# Unary encoding may address the area density challenge

## CMOS stochastic computing

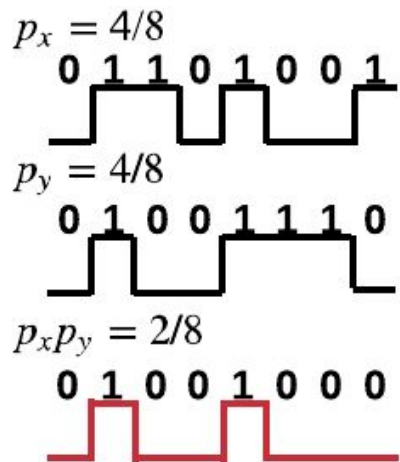


Good for arithmetic!

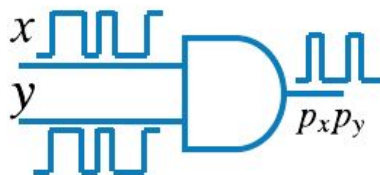
# Unary encoding may address the area density challenge

Advait Madhavan, et al. Race Logic: A Hardware Acceleration for Dynamic Programming Algorithms. In ISCA '14

## CMOS stochastic computing

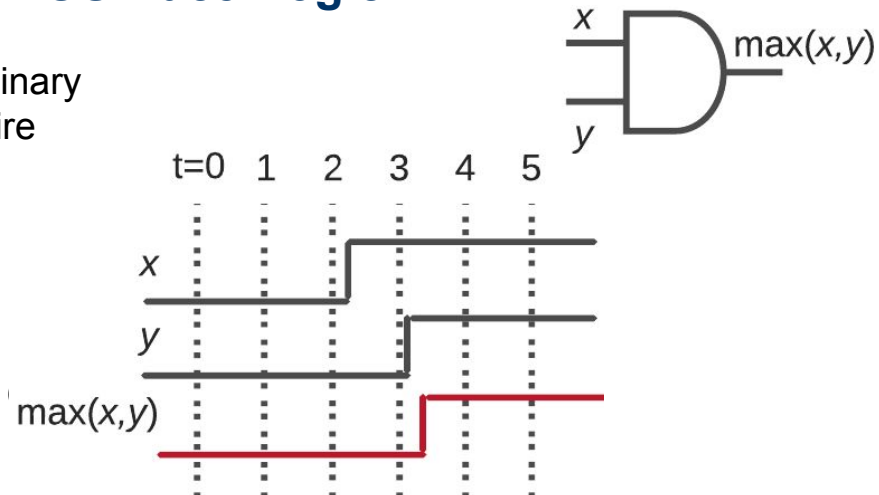


Many wires per variable from binary get replaced with a single wire



Good for arithmetic!

## CMOS Race-Logic



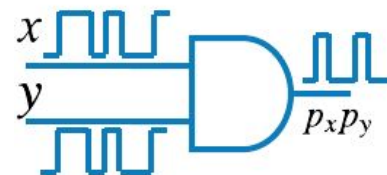
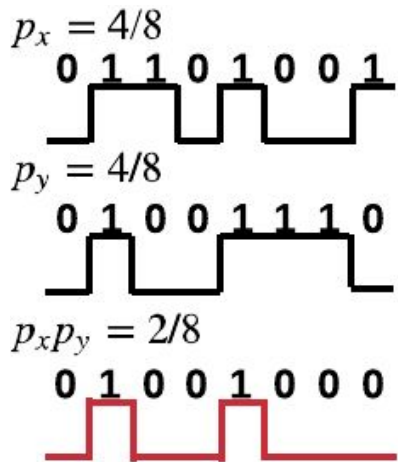
Good for dynamic programming algorithms!



# Unary encoding may address the area density challenge

Advait Madhavan, et al. Race Logic: A Hardware Acceleration for Dynamic Programming Algorithms. In ISCA '14

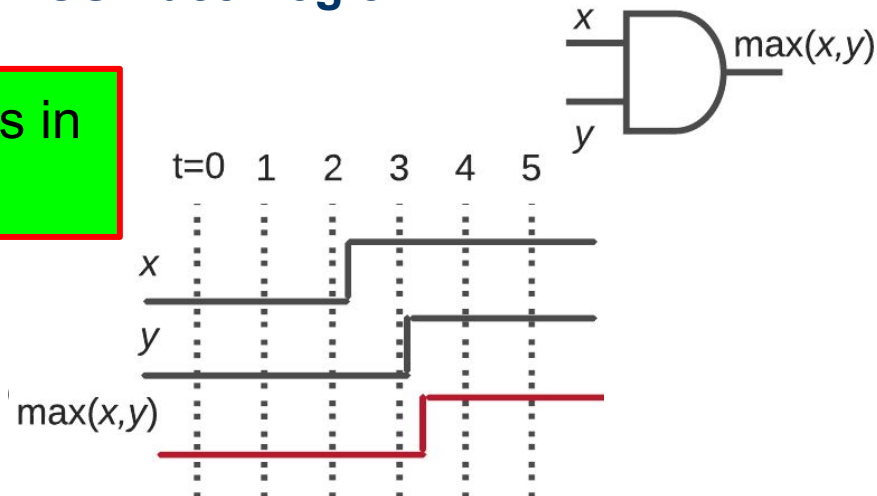
## CMOS stochastic computing



Good for arithmetic!

+90% savings in area

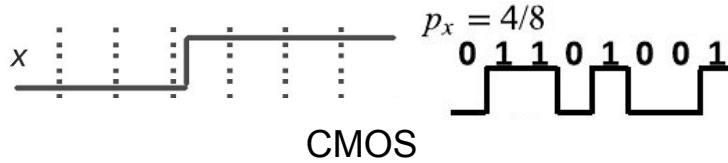
## CMOS Race-Logic



Good for dynamic programming algorithms!

Can we leverage the best of pulse-streams and race logic to build an area-efficient superconducting unary computing architecture?

## 1. Superconducting Unary SFQ encoding



## 2. Building Blocks



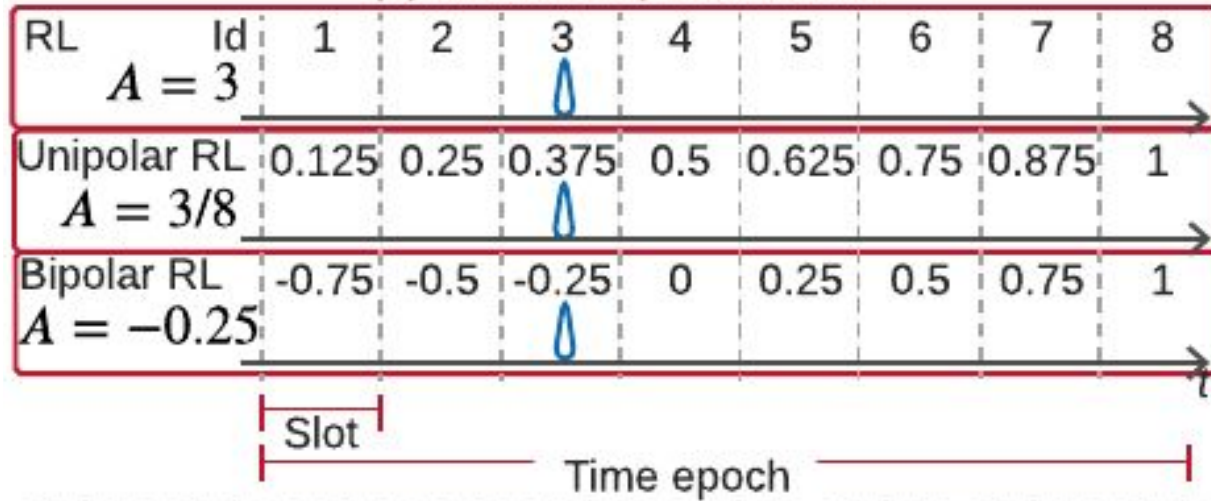
BINARY MULTIPLIER AND ADDER

## 3. Superconducting Accelerators

- Processing Element (PE) array for CGRAs or ANN
- Dot Product Unit
- Finite Impulse Response Filter

# In superconducting race logic data is mapped to the time the SFQ pulse arrives

Georgios Tzimpragos et al. 2020. A Computational Temporal Logic for Superconducting Accelerators. In ASPLOS '20

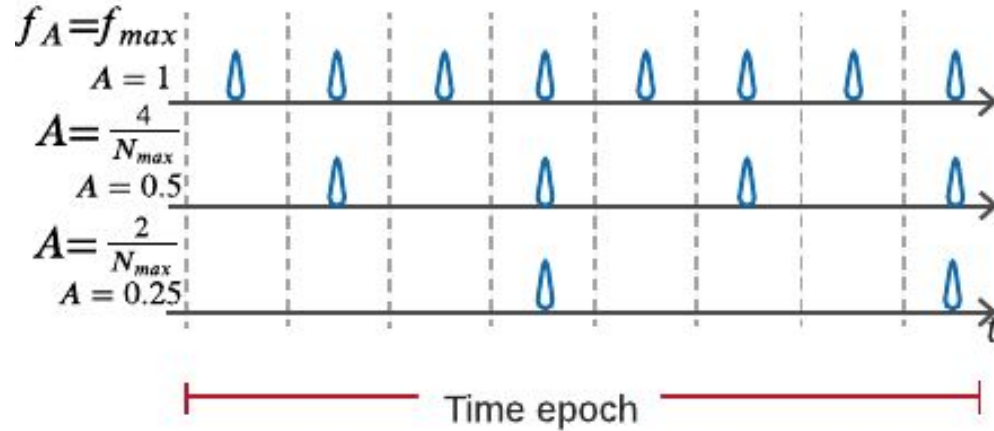


$$N_{max} = 8$$

$$A_b = 2A_u - 1$$

To obtain bipolar representation

# In a superconducting pulse-stream data is mapped to the frequency of the SFQ pulses



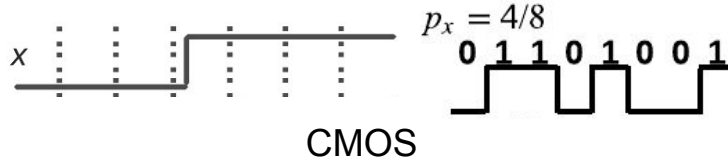
$$f_{max} \longrightarrow N_{max} = 8$$

$$A = n/N_{max}$$

$$A_b = 2A_u - 1$$

To obtain bipolar representation

## 1. Superconducting Unary encoding



RL	Id	1	2	3	4	5	6	7	8
$A = 3$									
Unipolar RL		0.125	0.25	0.375	0.5	0.625	0.75	0.875	1
$A = 3/8$									
Bipolar RL		-0.75	-0.5	-0.25	0	0.25	0.5	0.75	1
$A = -0.25$									

Slot

Time epoch

## 2. Building Blocks

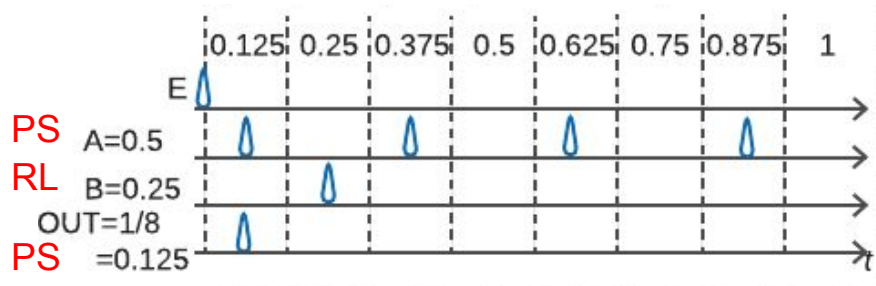


BINARY MULTIPLIER AND ADDER

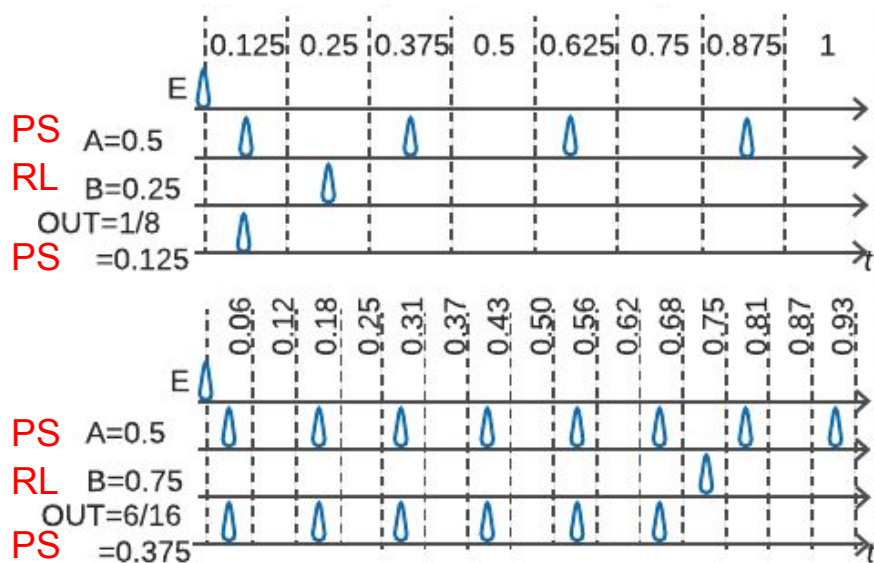
## 3. Superconducting Accelerators

- Processing Element (PE) array for CGRAs or ANN
- Dot Product Unit
- Finite Impulse Response Filter

# Building Blocks: Multipliers

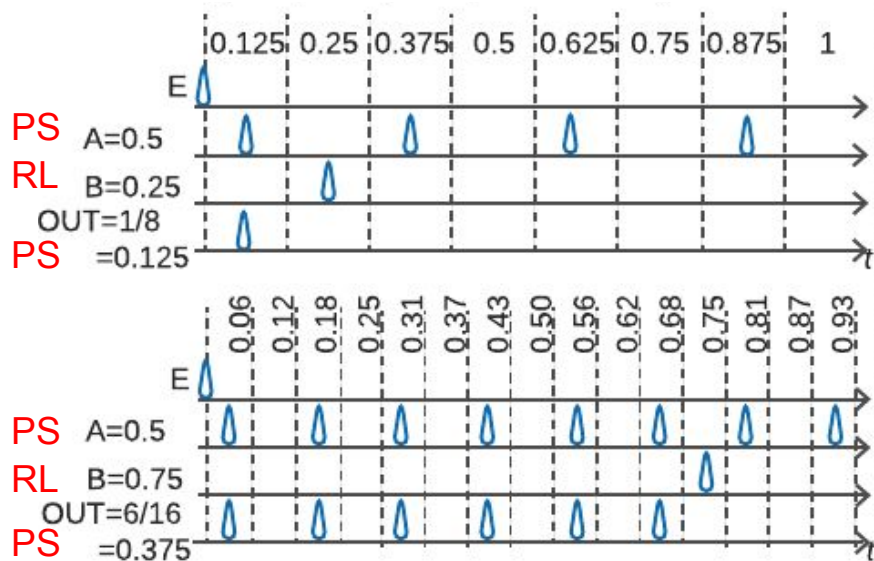


# Building Blocks: Multipliers

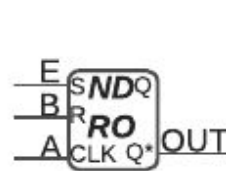




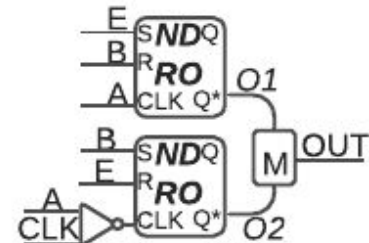
# Building Blocks: Multipliers



Essentially a CMOS XOR

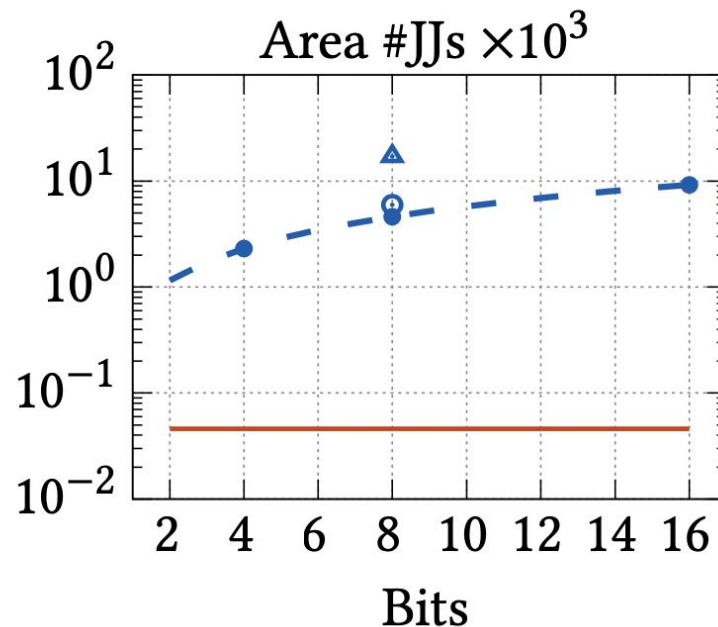
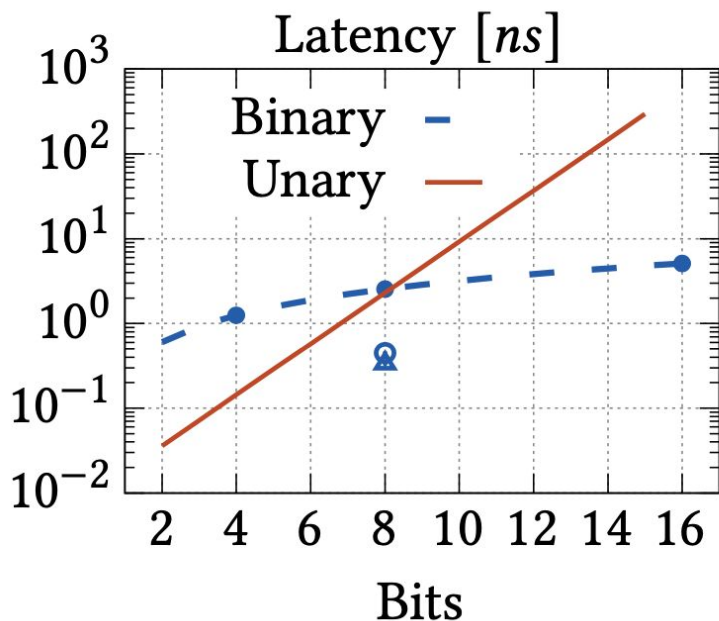


Unipolar SFQ multiplier



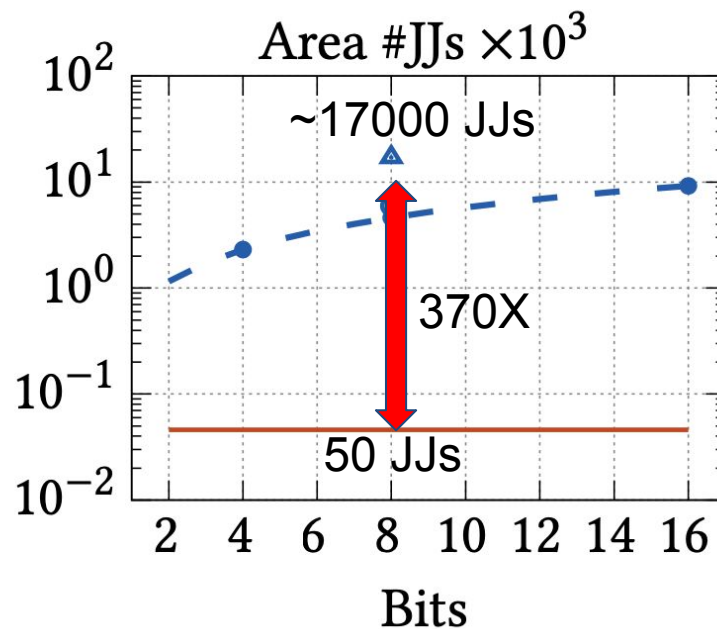
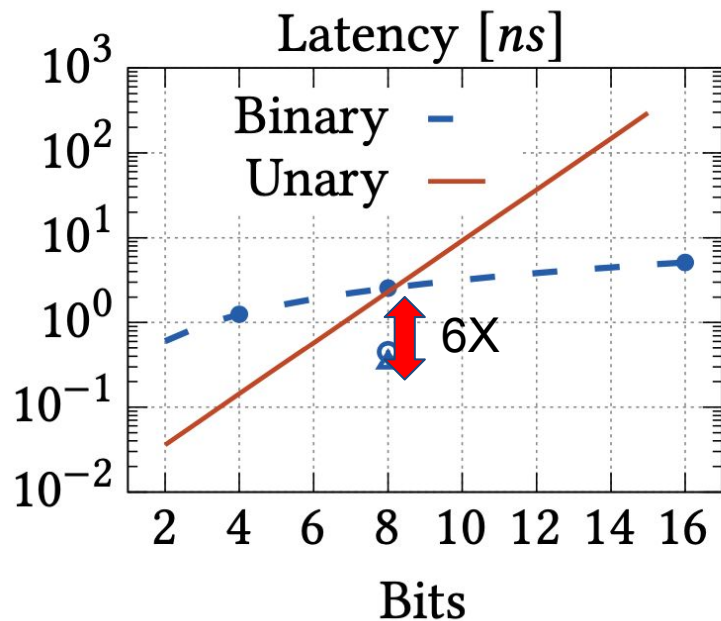
Bipolar SFQ multiplier

# The superconducting unary multiplier exposes an area-latency tradeoff

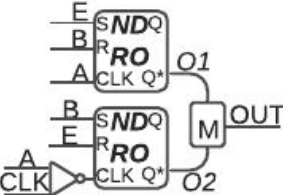
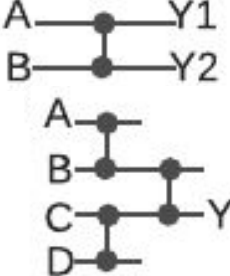
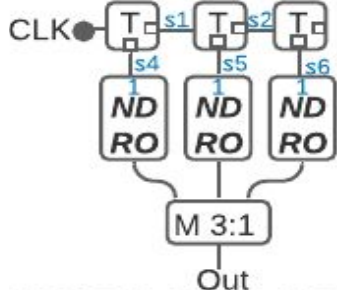
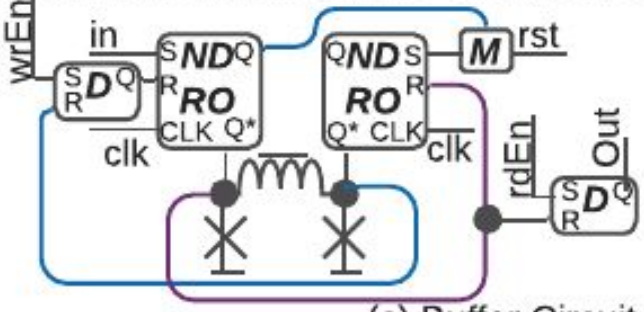


Bipolar unary multiplier

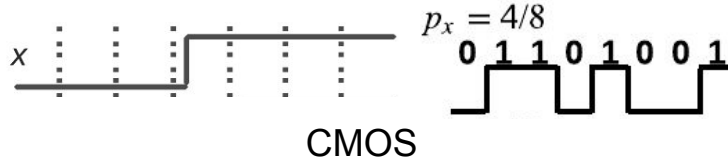
# The superconducting unary multiplier exposes an area-latency tradeoff



# The building blocks for the superconducting unary architecture

Multiplier	Adder	Pulse-stream generator	Memory cell for RL
			
<p><b>25X-370X</b> <b>Area Savings</b></p>	<p><b>11X-200X</b> <b>Area Savings</b></p>	<p><b>10% Area penalty</b></p>	<p><b>2.5X-1.3X</b> <b>Area penalty</b></p>

## 1. Superconducting Unary encoding



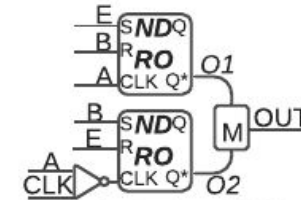
RL	Id	1	2	3	4	5	6	7	8
A = 3				1					
Unipolar RL		0.125	0.25	0.375	0.5	0.625	0.75	0.875	1
A = 3/8				1					
Bipolar RL		-0.75	-0.5	-0.25	0	0.25	0.5	0.75	1
A = -0.25				1					

Slot | Time epoch

## 2. Building Blocks



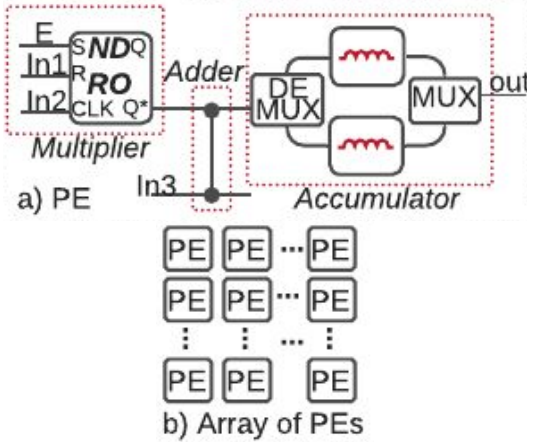
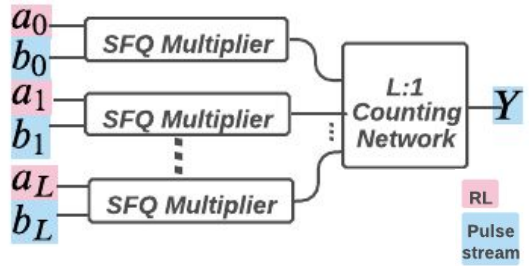
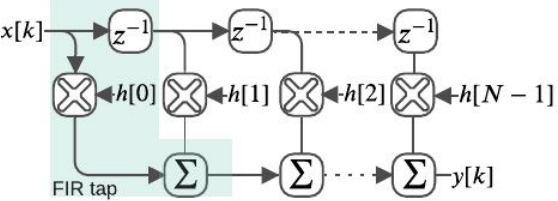
BINARY MULTIPLIER AND ADDER



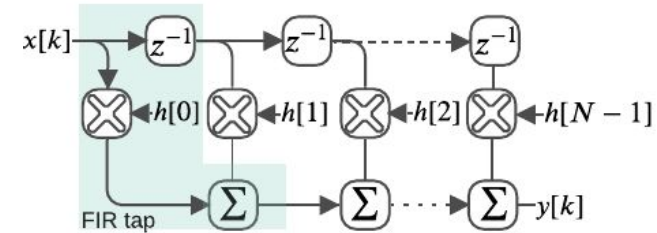
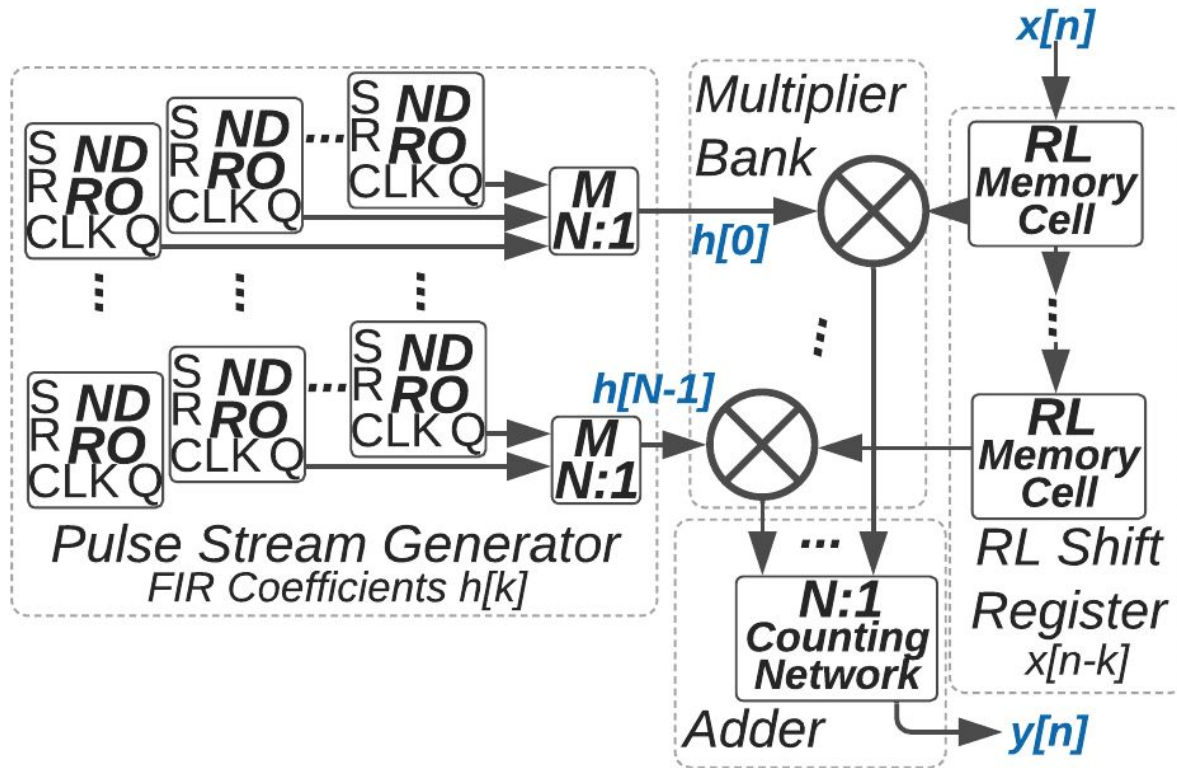
## 3. Superconducting Accelerators

- Processing Element (PE) array for CGRAs or ANN
- Dot Product Unit
- Finite Impulse Response Filter

# Evaluate the proposed superconducting architecture with three applications

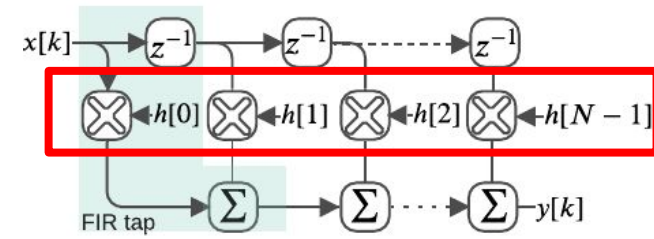
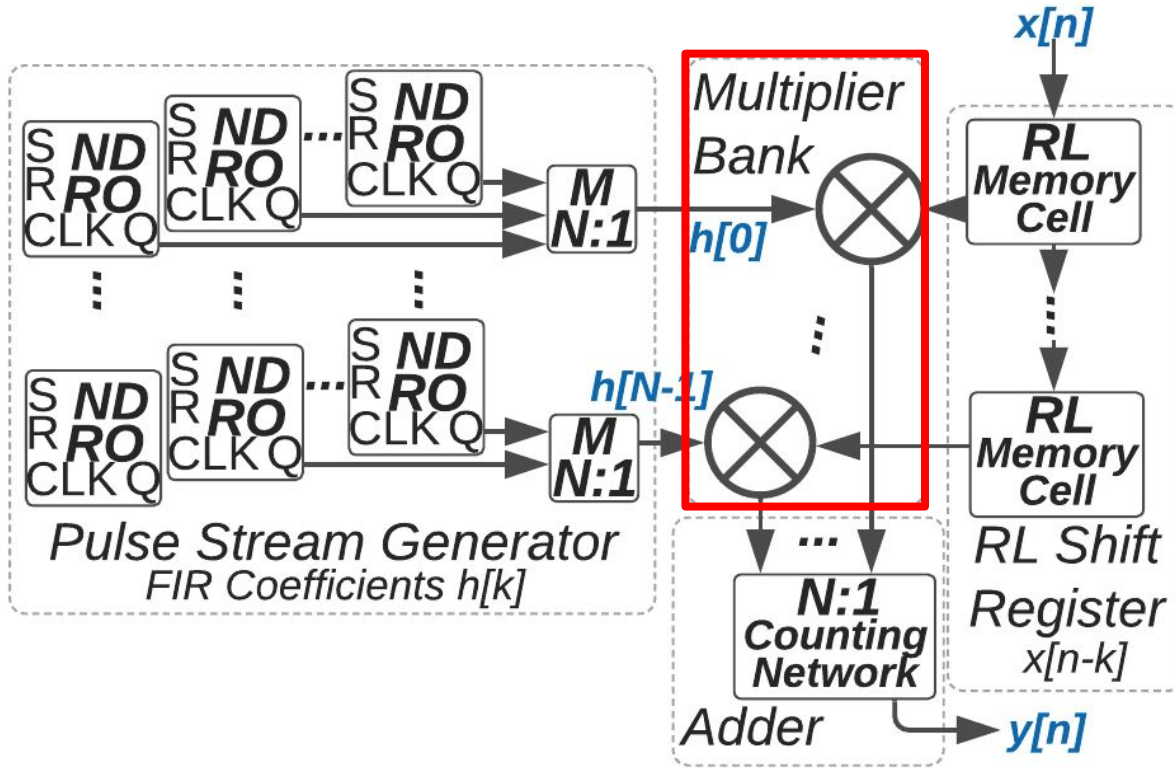
Processing Element for CGRA and Neural Networks	Dot Product Unit for DSP and Neural Networks	Programmable Finite Impulse Response Filter
 <p>a) PE</p> <p>b) Array of PEs</p>	 $y = a \cdot b = \sum_{i=0}^{L-1} a[i]b[i]$	 $y[n] = \sum_{k=0}^{N-1} h[k]x[n-k]$

# The FIR Filter accelerator shows the strengths and challenges of the proposed unary architecture



$$y[n] = \sum_{k=0}^{N-1} h[k]x[n-k]$$

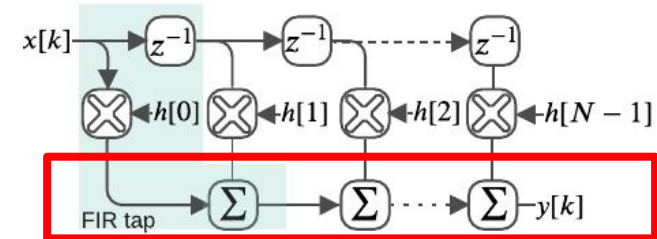
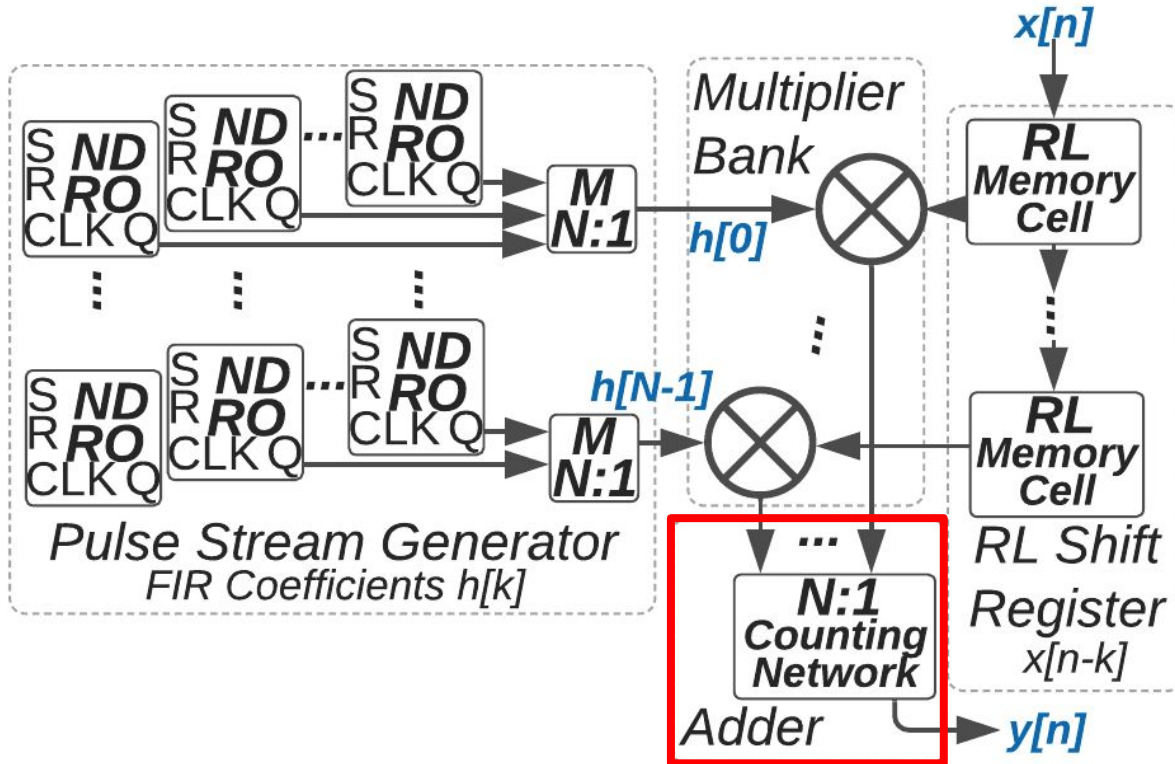
# The FIR Filter accelerator shows the strengths and challenges of the proposed unary architecture



$$y[n] = \sum_{k=0}^{N-1} h[k]x[n-k]$$

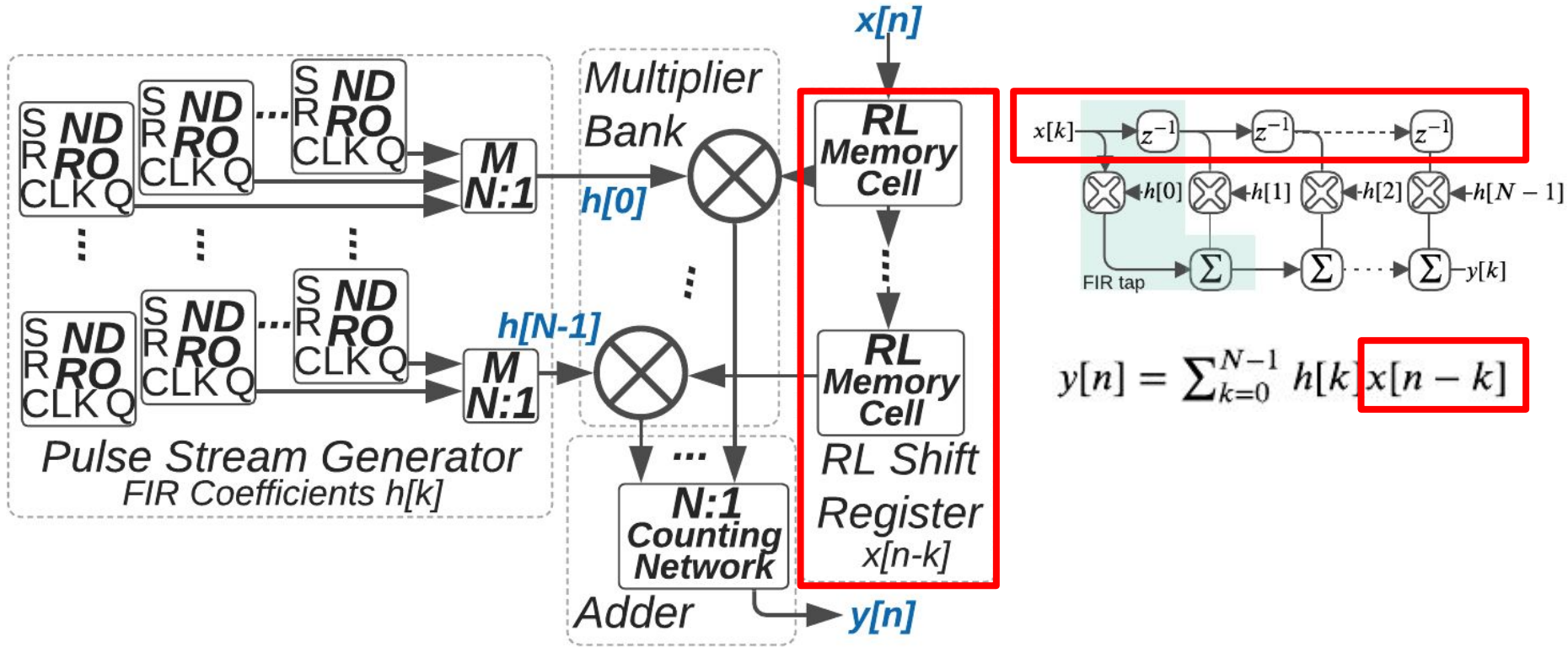


# The FIR Filter accelerator shows the strengths and challenges of the proposed unary architecture

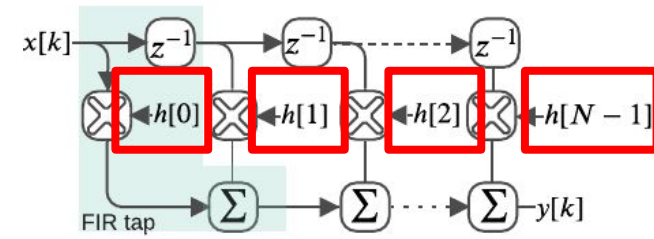
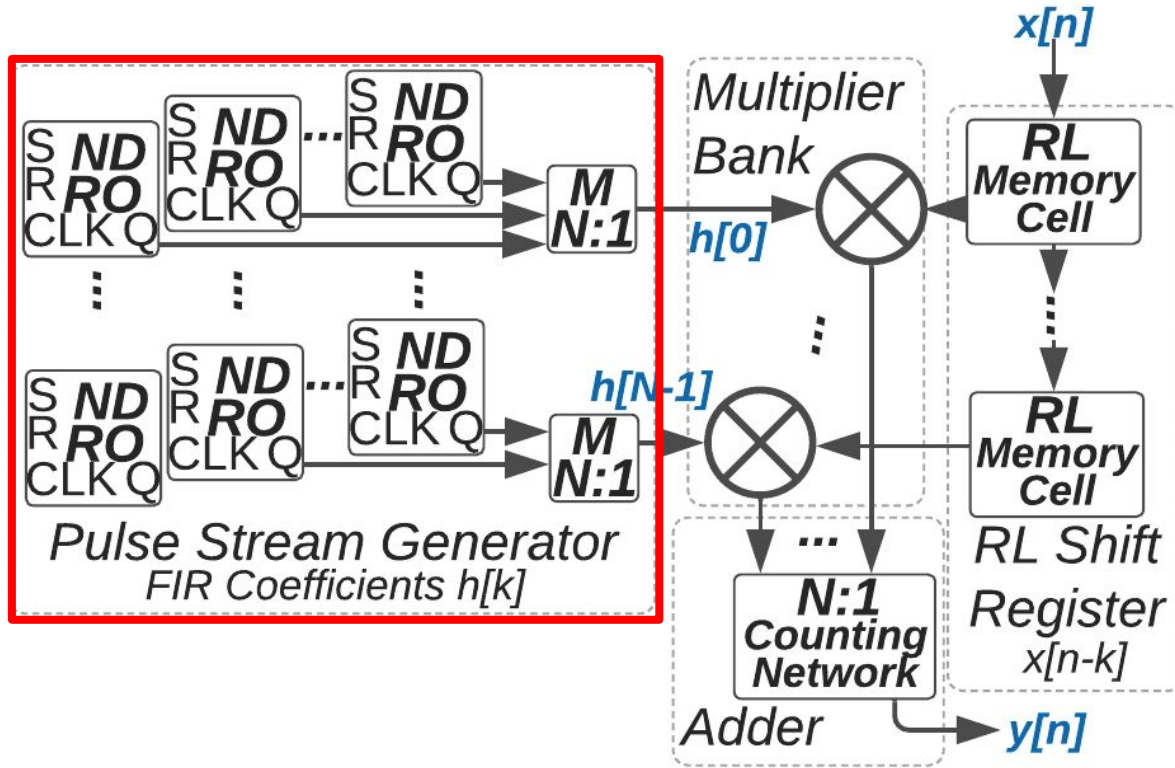


$$y[n] = \sum_{k=0}^{N-1} h[k]x[n-k]$$

# The FIR Filter accelerator shows the strengths and challenges of the proposed unary architecture

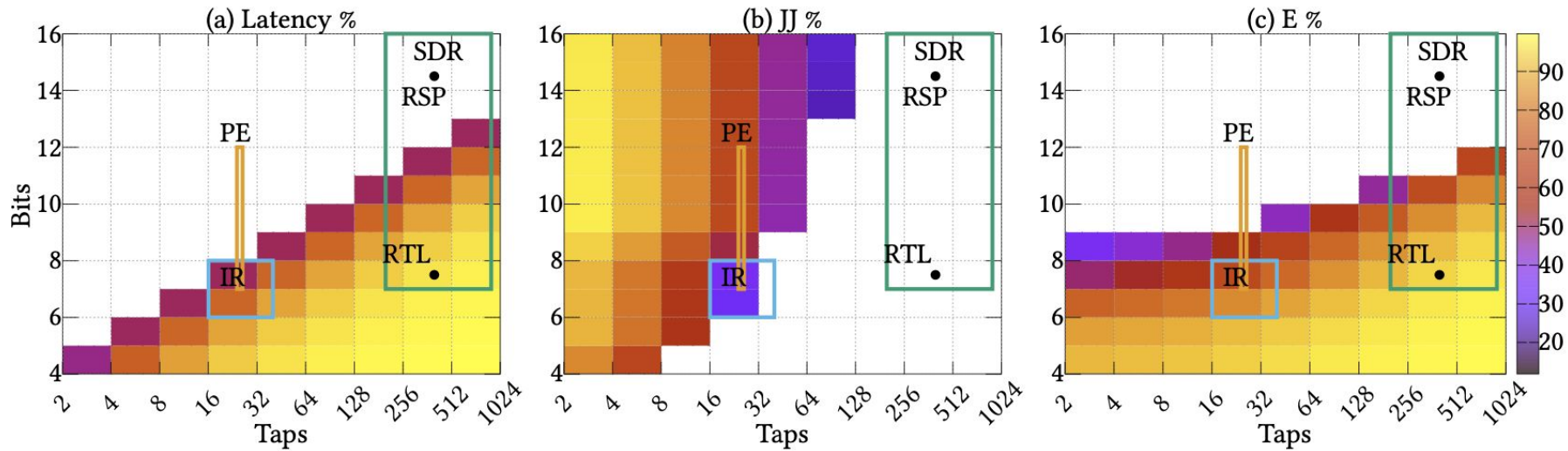


# The FIR Filter accelerator shows the strengths and challenges of the proposed unary architecture



$$y[n] = \sum_{k=0}^{N-1} h[k] x[n-k]$$

# FIR shows a design space where superconducting Unary yields significant advantages over superconducting binary



# Thank you for your attention

## Contact information:

[Ig4er@lbl.gov](mailto:Ig4er@lbl.gov), [mihelog@lbl.gov](mailto:mihelog@lbl.gov)

