

# BAQET: BRAM-aware Quantization for Efficient Transformer Inference via Stream-based Architecture on an FPGA

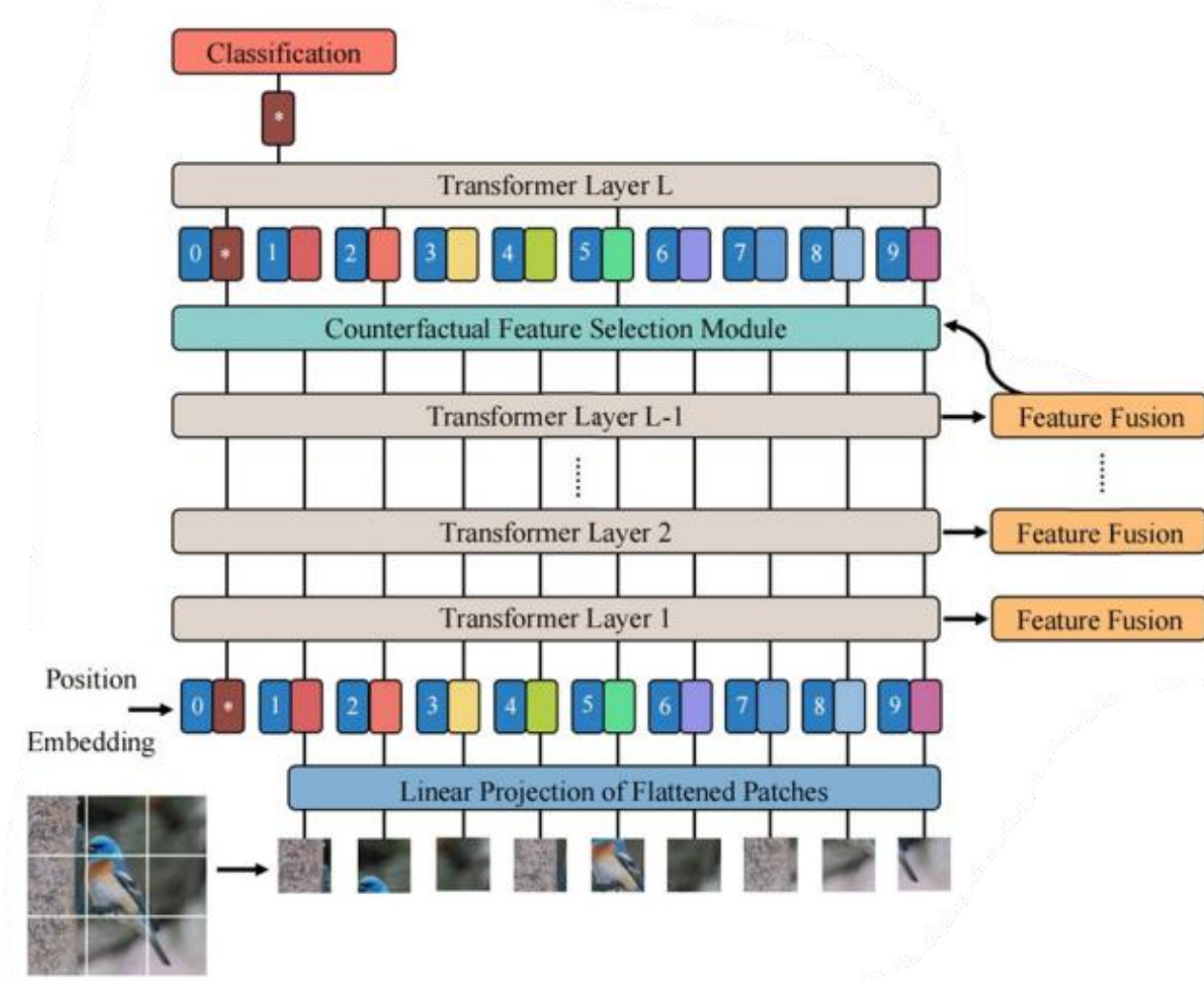
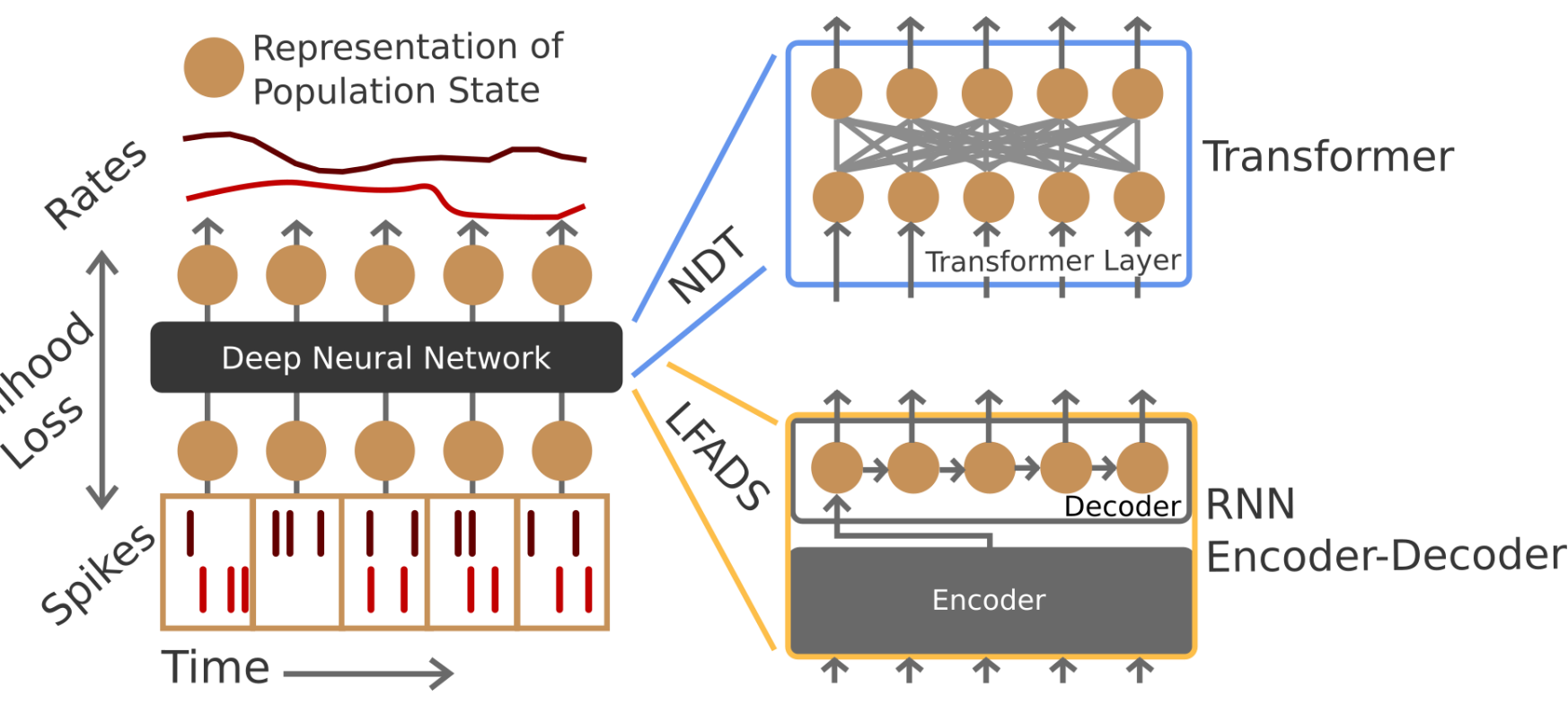
Ling-Chi Yang<sup>1</sup>, Chi-Rui Chen<sup>1</sup>, Trung Le<sup>2</sup>, Bo-Cheng Lai<sup>1</sup>, Scott Hauck<sup>2</sup>, Shih-Chieh Hsu<sup>2</sup>

<sup>1</sup> National Yang Ming Chiao Tung University, Taiwan | <sup>2</sup> University of Washington, USA

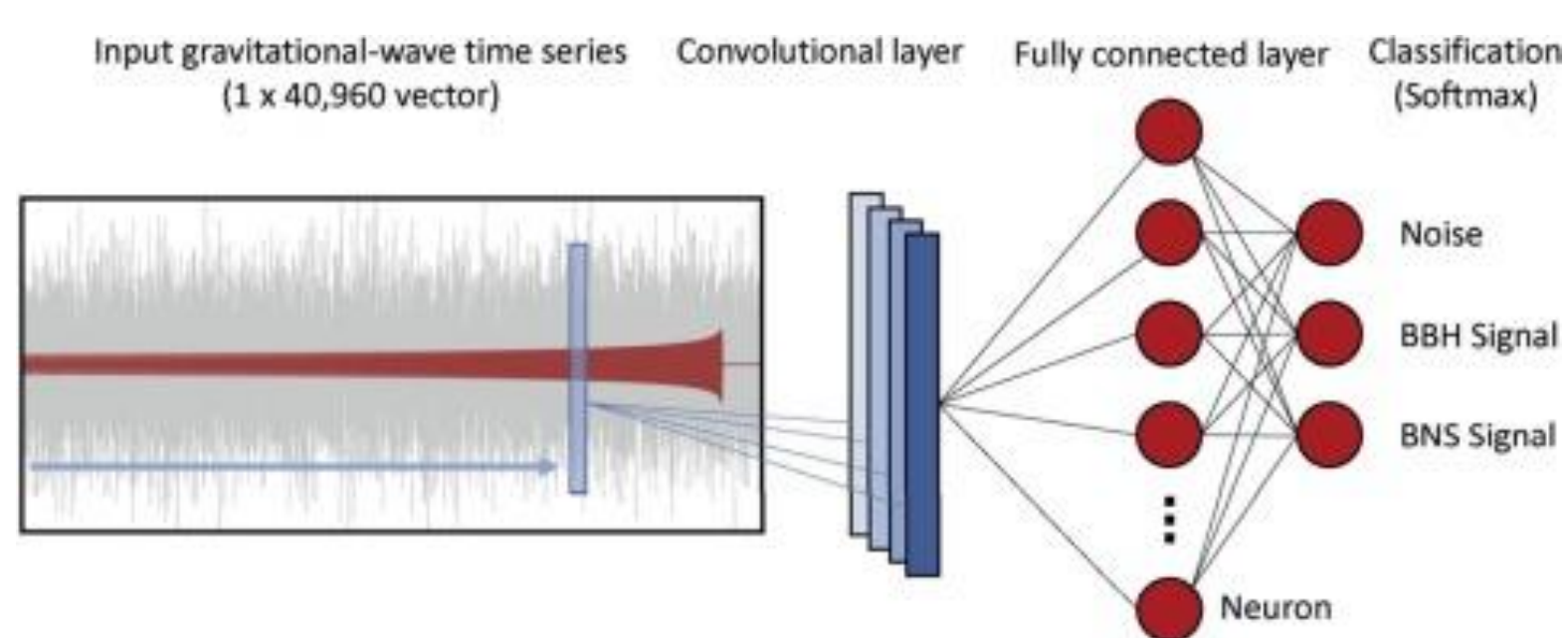
## Background

Transformer on GPU can't handle low-latency applications

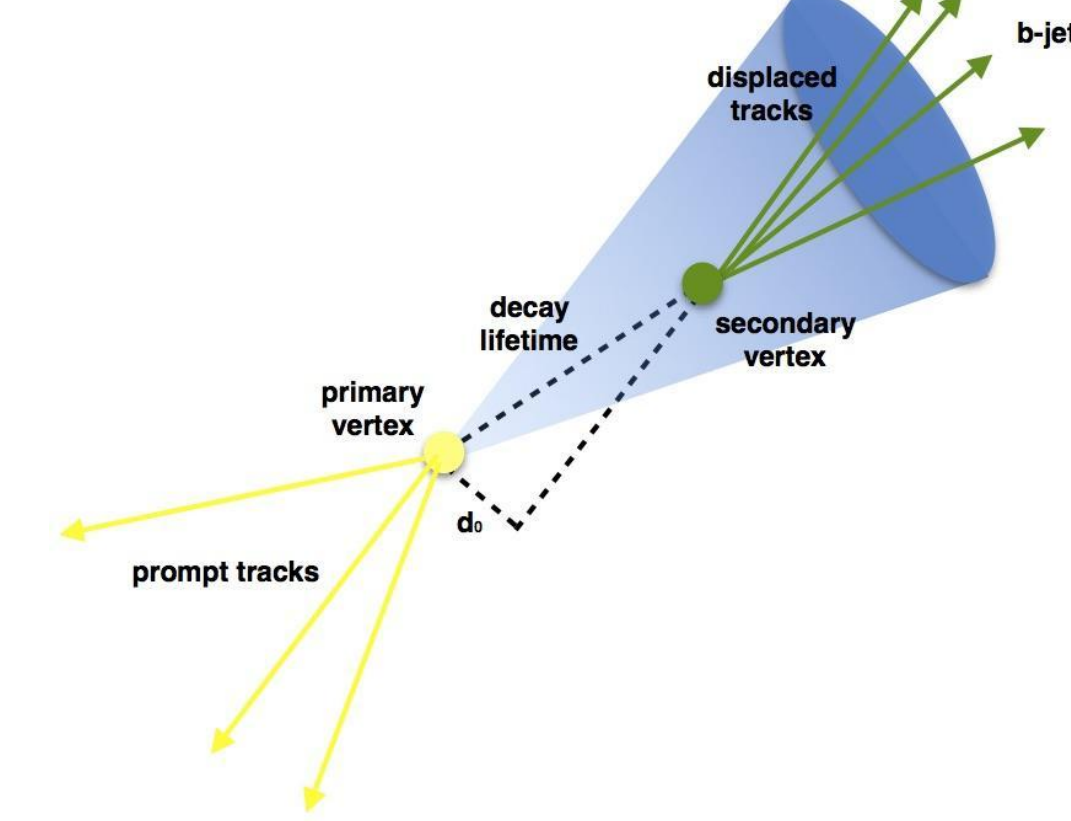
- Representation Learning for Neural Population Activity[1]
- Real-time Image Classification



➤ Gravitational Wave Anomaly Detection



➤ B-tagging Detection



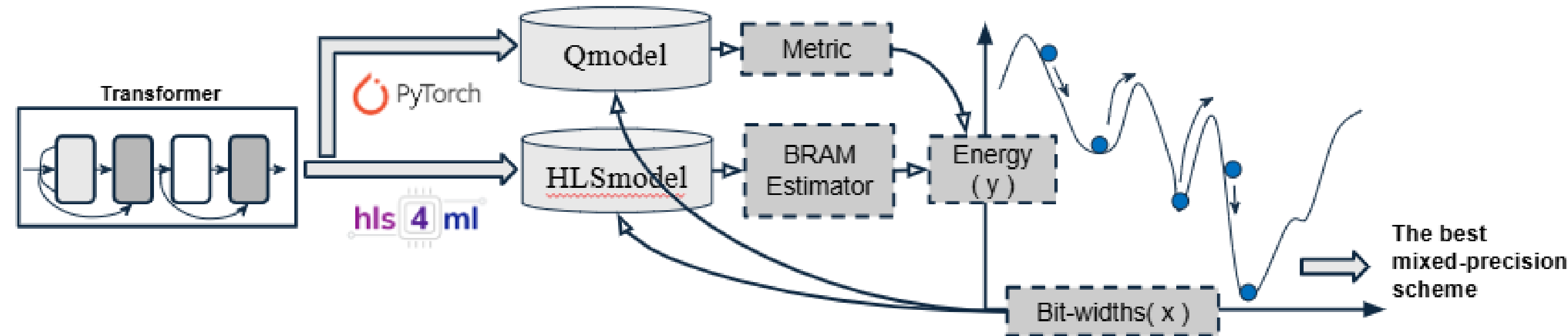
## Challenge & Contribution

### Challenges

- Large memory footprint of fully on-chip design
- Hard to determine mixed-precision scheme due to large search space

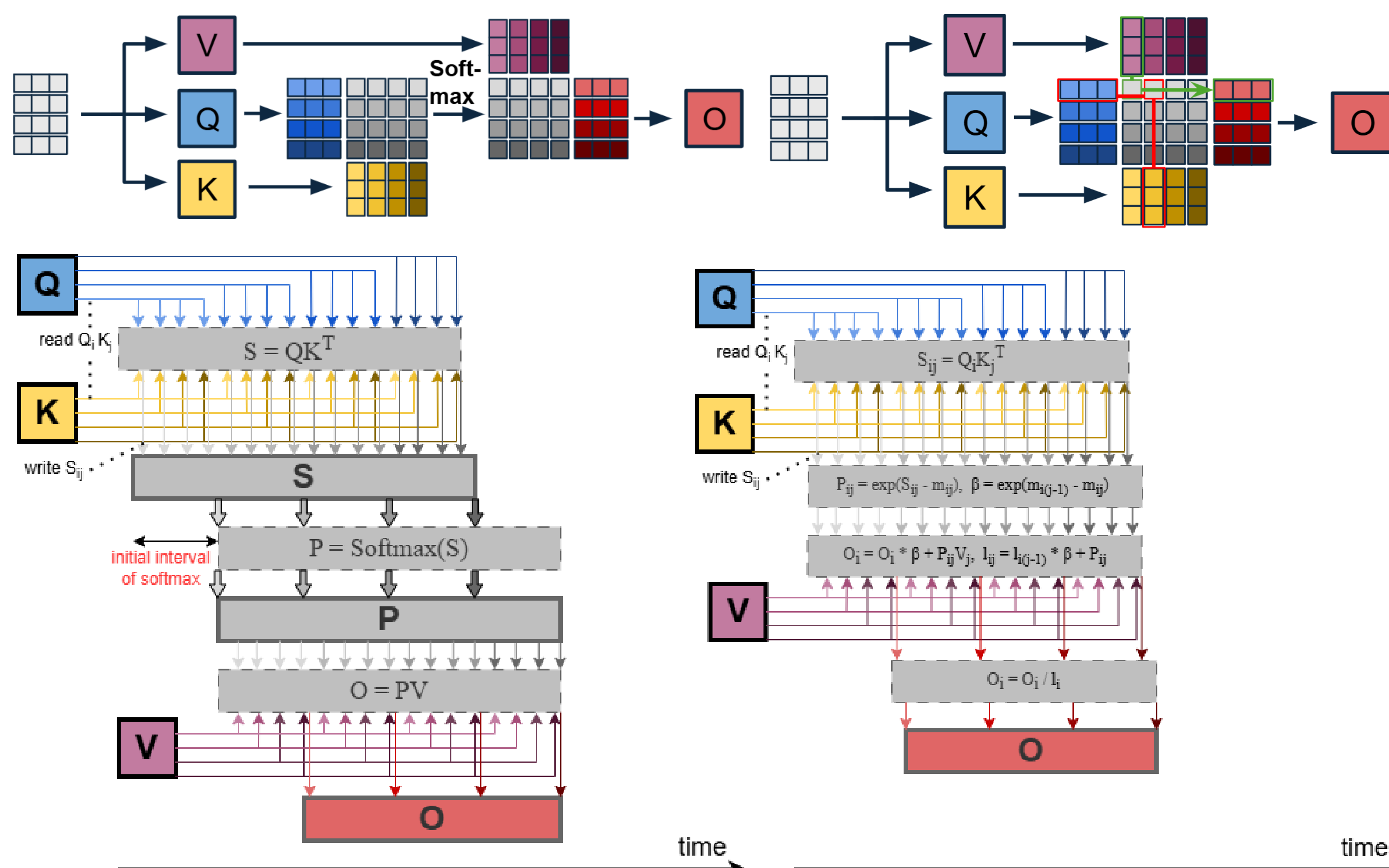
### Contributions

- Improved **Flash Attention[2]** algorithm in **stream-based architecture** on FPGA
- Estimating the number of BRAMs in hls4ml and formulating the problem, adopting **simulated annealing** to find better solution.

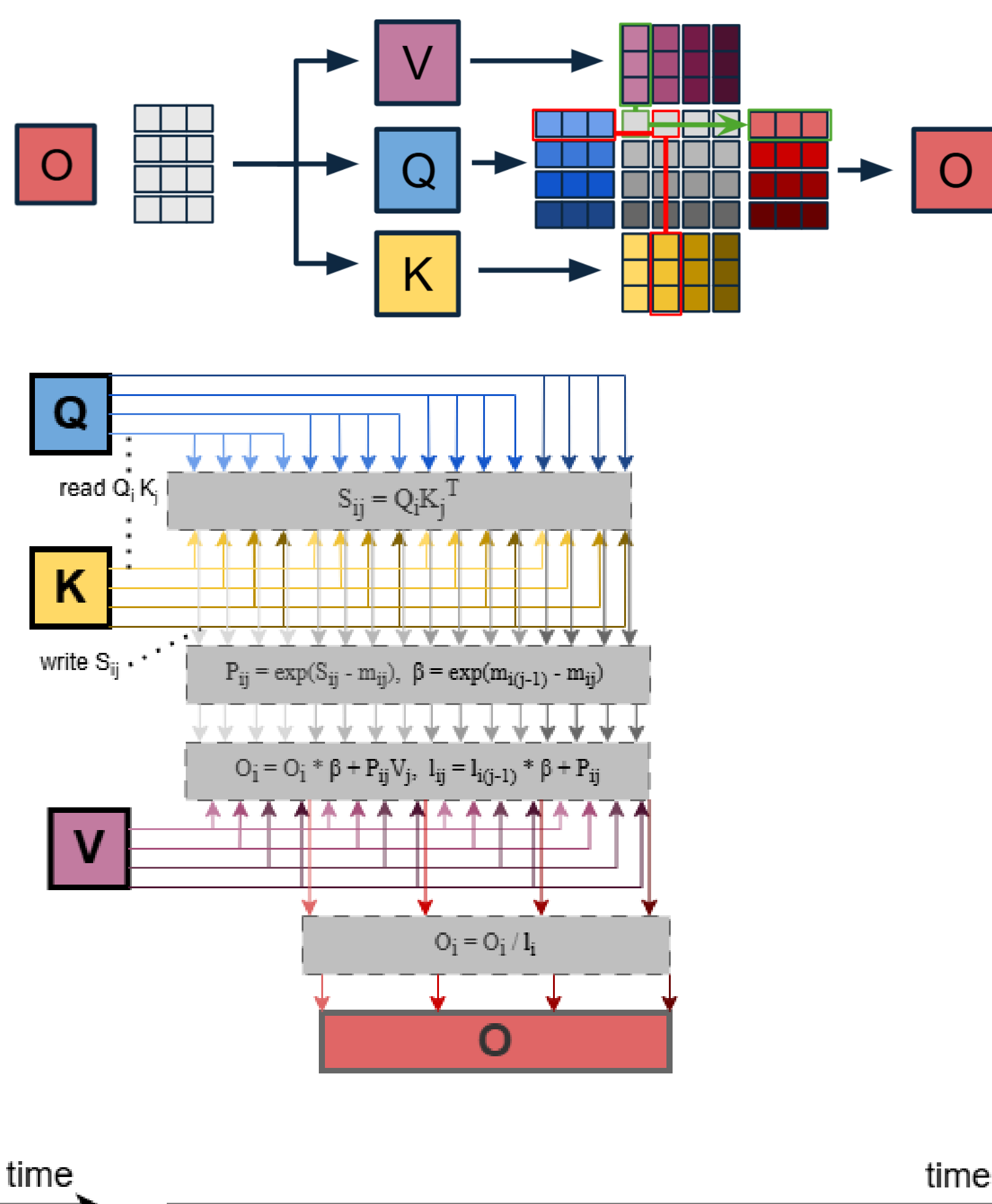


## Algorithm Optimizations

### Standard attention



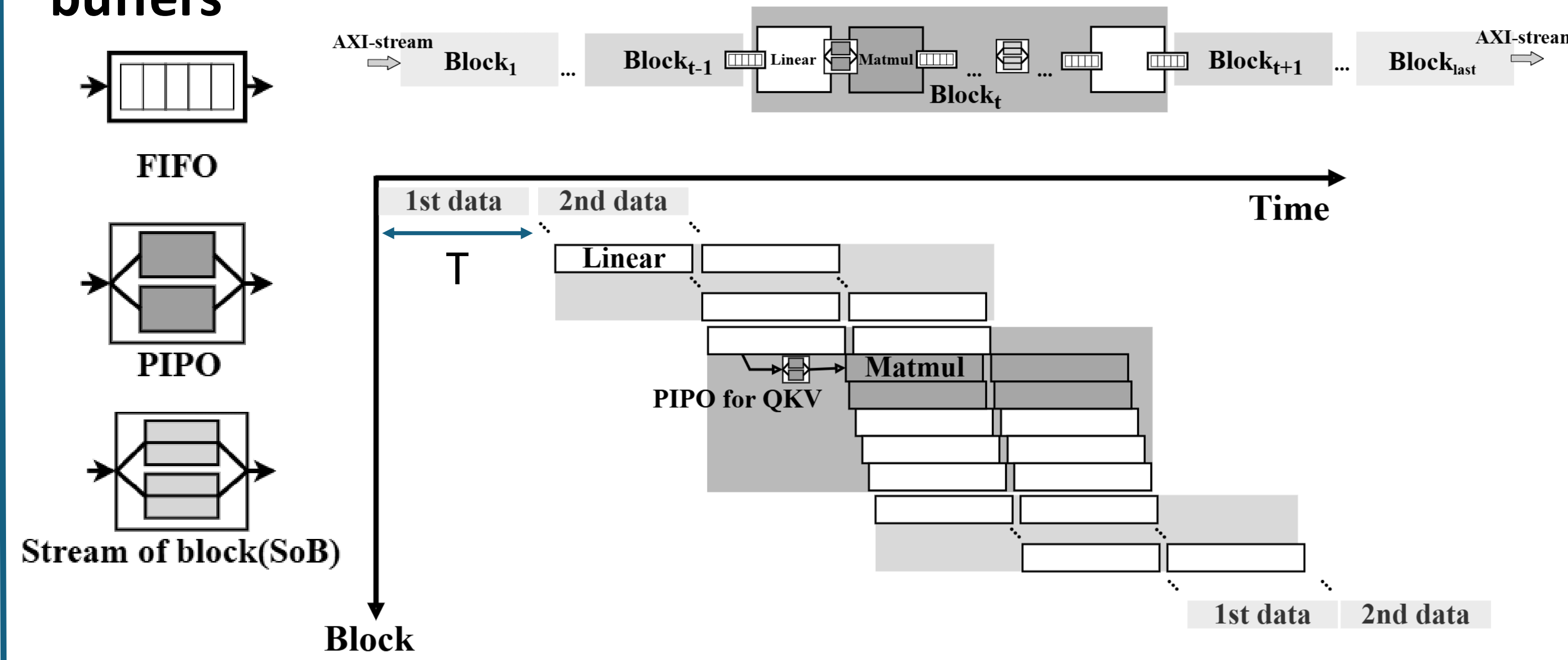
### Flash attention



- Saved memory footprint of attention and improved latency of Softmax

## Stream-based Architecture

Boost performance and save memory footprint with various buffers



- Overall latency for single batch =  $T \cdot (\text{number of transformer layer} + 1)$

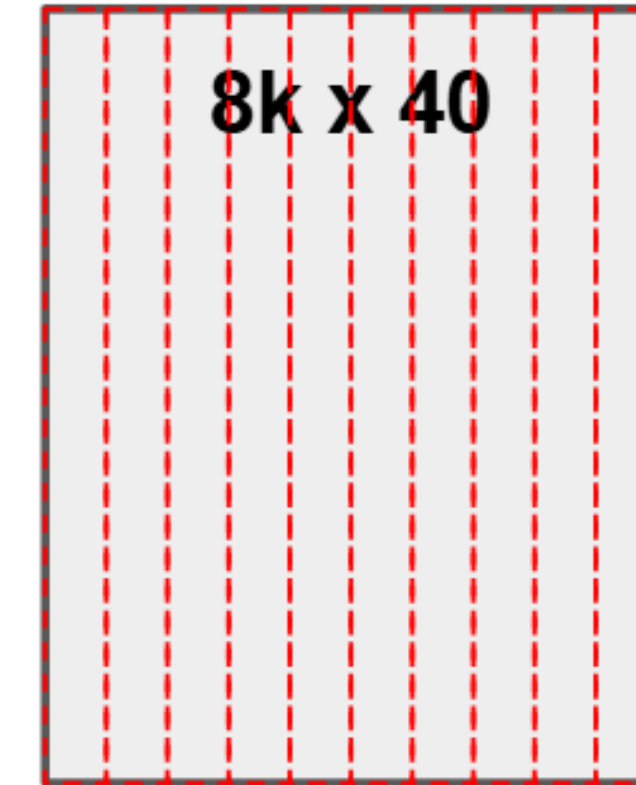
## BRAM-aware Quantization

- Formula :  $Q(x) = \frac{\text{bitwidth}_{int}}{\text{bitwidth}_{total}} \cdot \text{Clip}(\text{Round}(\frac{\text{bitwidth}_{total}}{\text{bitwidth}_{int}} \cdot x))$

- BRAM estimation

### Fixed Primitive

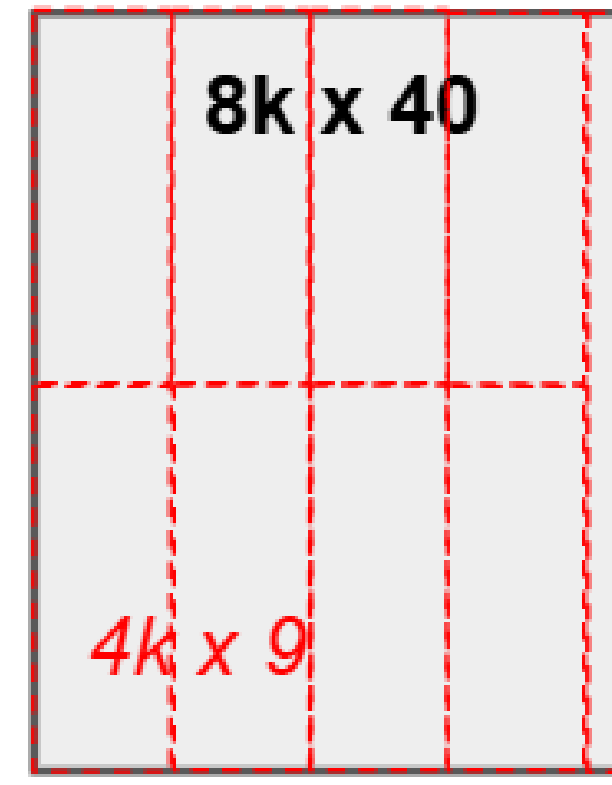
The core builds the memory by concatenating this single primitive type in width and depth.



For ROM

### Minimum Area

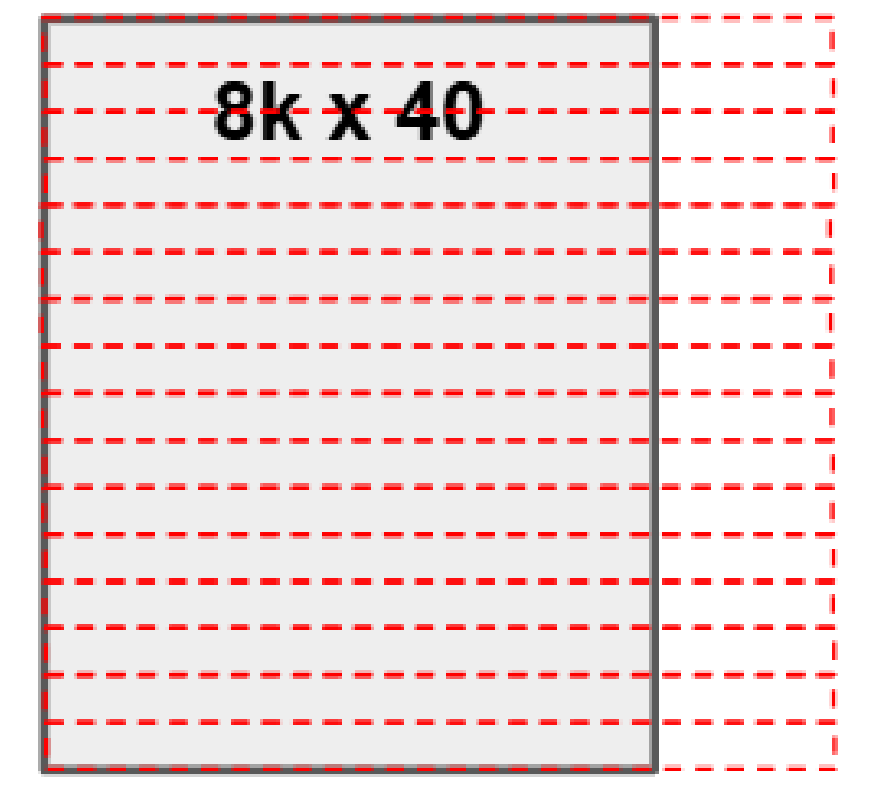
use the fewest resources while maximizing performance by reducing output multiplexing



For RAM

### Low Power

the minimum number of block RAM primitives are enabled during a Read or Write operation.



Tile size \* data bit width

- Cost function of multi-objective optimization

$$E = \text{penalty} \times \left( \frac{\# \text{ of BRAM}}{\text{available } \# \text{ of BRAM}} \right) + \alpha \times \left( \frac{\pm \text{ accuracy}}{\text{best accuracy}} \right)$$

- The more BRAM used, the worse the timing closure. A larger penalty is needed to reduce BRAM usage.
- Users can adjust the alpha value based on different applications to give more weight to accuracy.

## Experiments

- FPGA board : AMD Xilinx Virtex Ultrascale+ U55C
- Development platform : Vitis HLS 2023.2

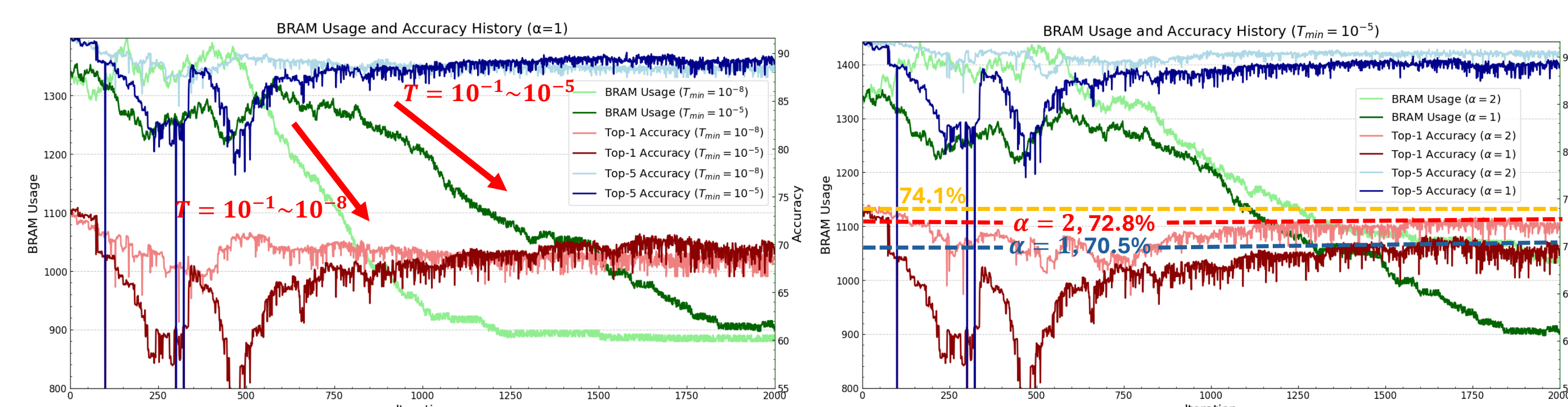
### Impact of tile size

- Larger tile size need higher bandwidth of RAM
- Higher tile size can improve power efficiency

	$B_{dk}, B_{ffn}$				
	1, 12	2, 24	4, 48	8, 96	16, 192
BRAM	73.96%	76.26%	69.44%	53.08%	99.80%
DSP	10.28%	15.69%	26.48%	47.57%	89.27%
LUT	12.19%	17.09%	23.57%	38.32%	72.66%
FF	3.30%	4.15%	6.34%	9.10%	16.41%
URAM	30%	30%	30%	60%	60%
Interval (ms)	12.54	6.27	3.136	1.56	0.78
Power (W)	12.03	14.332	18.23	28.32	51.203
Latency (ms)	36.87	18.441	9.22	4.61	2.3
WNS	0.189	0.041	0.031	-0.038	-1.516
Efficiency (GOPS/W)	16.2406	27.2641	42.8552	55.4560	61.3446

Work	Model		Board	Frequency (Mhz)	Power (W)	Throughput (GOPS)	Efficiency (GOPS/W)	
FTRANS	RoBERTa		VCU118	N/A	25	170	6.8	
ViA	Swin-T		U50	300	39	309.6	7.94	
SwiftTron	DeiT-S		ASIC	143	33.64	407	12.09	
SWAT	Swin-T		U50	200	14.35	301.9	21.04	
Ours	DeiT-T	NDT	U55C	200	7.161	166.27	23.21	
					13.32	9.086	308.77	33.98
		4,48			18.23	781.15	42.85	
		8,96			28.32	1563.67	55.21	

### Impact of alpha and temperature range



[1] Trung Le and Eli Shlizerman. 2022. Stndt: Modeling neural population activity with spatiotemporal transformers. Advances in Neural Information Processing Systems 35 (2022), 17926–17939.

[2] Tri Dao. 2023. FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning. <https://doi.org/10.48550/arXiv.2307.08691>.