#### 



Frontend implementation of Al-machine learning neural networks for in-pixel radiation-hard edge compute

Farah Fahim on behalf of Smart Pixels collaboration April 2023

## Case Study I: 28 nm implementation for CMS Phase III R&D – reconfigurable NN for data filtering in-pixel



#### LHC Grand challenge



Extreme Environment

Total ionizing dose : ~ 1 Grad

Single event effects : Flux of > 20 MeV hadrons = 10 MHz/cm<sup>2</sup>. Target 1 error every ~10s

Data Volume

1 billion sensor channels and 1 billion collisions/sec producing raw data at >Pb/s

More than 99.9% of data generated does not provide useful physics information

Occupancy ~ 3 GHz/cm<sup>2</sup> in vertex detectors

Requires more granular detectors  $\rightarrow$  increasing data volume, power consumption, more complex on-detector electronics. **Per chip ~ 1Mpixels; > 1 Tbps of zero-suppressed data.** 

## **Edge computing**

What is the best approach?

- Data compression
- Data filtering
- Data featurization
- Where should the algorithm be implemented?
- Near sensor: Al/ML mixed-foundry chiplet:

Standalone general purpose

Does not overcome IO bottleneck (1 Tbps off-sensor)

Edge of sensor/ Readout integrated circuit:

Easier implementation (larger area and power), optimally utilizes chip bandwidth On-chip transfer

• In-pixel: At the data source:

Minimizes data movement

Power and area constraint







## Design Methodology: Physics driven hardware co-design

ALGORITHM

DEVELOPMENT

ML Model

- Algorithm development based on Physics data
- hIs4mI simplifies the design of on-chip ML accelerators

   hIs4mI directives | << | HLS directives |</li>
   C++ library of ML functionalities optimized for HLS
- TMR4sv\_hls: Triple Modular Redundancy tool for System Verilog & HLS



## **LHC Pixel detector**



- 4 layers closest to the beam pipe
- Vertex position for high momentum tracks (smaller pixels => higher resolution)
- ~5B pixels across ~30,000 chips
- R&D for Phase II started in 2012; chips in production ~ 2022
- 20B pixels; starting R&D 2022 (expected upgrade in 2034)

Technology	65nm CMOS	28 nm CMOS
Pixel size	50x50 μm <sup>2</sup>	25x25 μm <sup>2</sup>
Pixels	394x400 = 157.6k	788x800 = 0.63M
Detection threshold	~1000e-	~500e-
Hit rate	< 3GHz/cm <sup>2</sup>	< 3GHz/cm <sup>2</sup>
Trigger rate	1MHz	40MHz (?)
Digital buffer	12.5 μs	(?)
Readout data rate	1-4 links @ 1.28Gbps	1 photonic link @ 30 Gbps
Radiation tolerance	500Mrad at -15°C	1Grad at -15°C
Power	1 W /cm <sup>2</sup>	1 W /cm <sup>2</sup>
		52 Fermi



#### **Concept: Real-time tracking**



Cluster shapes and Pulse information for filtering out low  $p_T$  particles

- NN classifier identifies and saves clusters from tracks with  $p_T > 2 \text{ GeV}$
- $\geq$  95% data reduction by saving only high p<sub>T</sub>
- Low power implementation

#### Compact algorithms for data reduction through featurization

- Predict physics information (x,y,θ,φ) and meaningful error (UQ) on particle position, angle
- Potential for reduction of track seeds → saves time & computing resources down the line

#### Technology development to enable on-sensor computing

- Ultra low power in-memory compute chips
- 3D integration for optimized data processing
- Leverage emerging technologies such as novel CMOS compatible memristors

S Fermilab

## Algorithm development for momentum classification



## Smart pixel dataset





- State-of-the-art dataset for developing algorithms for implementation on-ASIC
- Simulated MIP interactions in a **futue pixel detector** Located at radius of 30 mm
  - 3.8 T magnetic field
  - $50 \times 12.5 \mu m$  pixels
  - 100 µm thick sensor
  - Time steps of 250 picoseconds
- Use track parameters taken from CMS data in Run 2

Simulation with time-sliced <u>PixelAV</u>, including E field and weighting field



Morris Swartz

Farah Fahim | In-Pixel Al 4/26/2023

on <u>zenodo</u>

9

#### **Co-relation between cluster shape and momentum**



- 1/3<sup>rd</sup> of all clusters are associated with tracks
- Untracked clusters are mainly due to loopers, broken clusters & other backgrounds
- Low p<sub>T</sub> positively charged tracks increase cluster size
- Low p<sub>T</sub> negatively charged tracks decrease cluster size
- Strong correlation between y-size and y-local on the module
- Multi-level classifier for selecting high momentum tracks

Morris Swartz

🛟 Fermilab

## Data filtering to reject 95% of data



Conservatively reject:

< 0.2 GeV	≥ 6%
< 0.5 GeV	≥ 36%
< 1 GeV	≥ 70%
< 2 GeV	≥ 94%

- Only hits that are reconstructed into tracks are simulated
- Rejecting lower p<sub>T</sub> hits (often single pixel) will increase the fraction rejected
- Generating data sets associated with untracked clusters, noise hits and other background



11 4/26/2023 Farah Fahim | In-Pixel Al

Morris Swartz, Jennet Dickinson

## **Algorithm performance**



• Model 1: y-size: count no. of pixels with hit

 Model 2: y-profile: add total charge accumulated across all pixels in the ydirection

 Model 3: y-profile with timing: add total charge accumulated across all pixels in the y-direction every 250ps – 8 data points over 4ns of charge integration

$\mathbf{Model}$	Data Reduction	False Positives	False Negatives
Model 1	77.9~%	13.4%	20.7~%
Model 2	65.8~%	10.9%	13.4~%
Model $3$	57.4~%	8.9~%	8.1%



Jieun Yoo

## Algorithm co-design with frontend architecture



- Lowering noise threshold increases slightly improves efficiency
- Increasing no. of bits does not improve efficiency

Threshold $[e^-]$	$0.2~{ m GeV}$	$0.5~{ m GeV}$	$1 { m ~GeV}$	$2  {\rm GeV}$
0	78.2~%	88.2~%	91.0~%	92.6~%
400	77.5~%	85.6~%	87.5~%	88.3~%
600	74.4~%	82.2~%	84.3~%	85.9~%
800	75.3~%	82.2~%	84.3~%	86.0~%



## **Retraining local algorithms**



- Sensor pixels can be divided into various zones to exploit the correlation with y\_local
- The algorithm is reprogrammable and weights are trained based on position dependent data
- 2×2 cm sensor module: 3200 rows & 100 cols
- Every region needs minimum 16 rows and maximum ~ 64 rows
- 16×16 sensor pixel = 8×32 ROIC pixels



## SuperPixel concept: Algorithm co-design



Total Digital area for 256 pixels

> V1: 64(625×4 - 625) = 120,000 μm<sup>2</sup> V2: 64(625×4 - 825) = 104,000 μm<sup>2</sup>





## **Design space exploration for hardware co-design**



Target digital area V1= 120,000 μm<sup>2</sup> Target digital area V2 = 104,000 μm<sup>2</sup>

- Original model is too large for the given area
- Exploring area reduction without significantly decreasing the model accuracy
- Reducing the number of neurons in the first hidden layer and the bit-width of the weights and activation function



## Design space exploration for hardware co-design



- · Compact: Reduce size by reducing the number of parameters
- Energy Efficient: Eliminate clocking and make the design fully combinatorial
- No. of operations: 75 GOPS (6b INT) per super pixel: 0.5 TOPS/mm<sup>2</sup>: 200 TOPS/chip

17 4/26/2023 Farah Fahim | In-Pixel Al

Giuseppe Di Guglielmo

🛟 Fermilab

## **Readout Integrated Circuit: Traditional Implementation**



## Front-end pixel architecture: Synchronous ADC



- Signal processing within 1 BXClk cycle (25 ns)
- Insensitive to pileup



#### **Front-end pixel architecture**





- 25µm × 25µm pixel AC coupled ADC
- 2-bit ADC:

 $\rightarrow$  version 1: Single ended

 $\rightarrow$  version 2: Differential

• Auto zero in every pixel for threshold correction



## **Charge Sensitive Preamplifier**



- Regulated cascode for gain stage
- Detector Capacitance 30fF
- Leakage current compensation (up to 50 nA)

Result	Mode 1	Mode 2	unit
Power	3	4.2	μW
ENC	61	52	e-
Charge sensitivity	40.85	42.3	μV/e-
Phase Margin	65	78	degree
Dynamic Range	13	13	ke-





- AC coupled
- Auto-zero for offset cancellation
- Exploring an alternate clockless architecture

Result	v1	v2	unit	comments
Power	3.7	5.2	$\mu W$	Per pixel
Total Equivalent Dispersion Charge	90	100	e-	Includes: • ENC • Qth • Baseline fluctuation • kickback
Min Threshold	430	475	e-	4.75 σ
Analog Area	678	845	μm²	

**Benjamin Parpillon** 

#### Grouping 4 pixels to create an analog island



23 4/26/2023 Farah Fahim | In-Pixel Al

**Benjamin Parpillon** 

#### HLS implementation – multi-classifier algorithm





24 4/26/2023 Farah Fahim | In-Pixel Al

## **SuperPixel implementation**



- Floorplan with analog pixels with power and bias grid
- Red highlight corresponds to the data classifier algorithm
- White highlight corresponds to the registers for programming the neural network



 Full triple modular redundancy for SEE mitigation



#### **Current status and next steps**



- Analog prototype ROIC with 16 x 32 pixels just been manufactured. In the process of starting tests
- Prototype with embedded algorithm with two super pixels expected to be submitted May 2023
- Next version of chip to include

Multiple super-pixels and peripheral data transfer of selected clusters

Explore further data reduction through featurization in the periphery



#### Algorithm development for featurization



## Predicting x,y, angle of incidence to sensor ( $\alpha$ , $\beta$ )



- Where does the particle cross the sensor mid-plane
- Angle of incidence to sensor (perpendicular and in-bending plane of the magnet)
- 3D tracking algorithms are computationally expensive! Grows with number of possible combinations of hits
- With angle and reliable uncertainty, area to search for next hit is defined. Save resources by not looking at bad combinations



Jennet Dickinson, Gauri Pradhan, Lindsey Gray

## **Prediction with uncertainty (x,y)**



20 40 Full precision • Full precision + time Quantized ..... 10 15 20

30

.

Full precision Full precision + time

**Fermilab** 

Jennet Dickinson, Gauri Pradhan, Lindsey Gray

## Prediction with uncertainty ( $\alpha$ )





30





Jennet Dickinson, Gauri Pradhan, Lindsey Gray

**Fermilab** 

## **Prediction with uncertainty (β)**

#### $\beta$ from y-profile



Jennet Dickinson, Gauri Pradhan, Lindsey Gray

## **Implementation ideas**

#### Multi-foundry chiplet

- Serialize and transfer selected clusters off-chip on a high-speed link
- AI Chiplet: 7/5/3nm node
- Highest overhead and power consumption



#### **ROIC edge**

- Transfer selected clusters to chip periphery
- Data transfer across chip ~ 2cm is about ~1pF
- High speed inference engine



#### 3D-IC

- Highly parallel processing enabled by vertical integration
- Low power data transfer
- Reduced clock speed



## Neuromorphic Implementation: How to exploit the time information?



#### Neuromorphic frontend concept based on Dynamic Vision Sensors



- Pulse shape encoded in spike train
- Ultra-fast signal processing ~ 100ps
- Digital processing power constraints
- Analog processing device constraints (speed)





## **Neuromorphic Algorithm development**

# Abisko Codesign Overview



## **Neuromorphic Processing – Basic Model of Operation**



- Computation with sparse events
- Binary valued events



- Inherent temporal dynamics
- Neuron fires with charge exceeding threshold
- Multiple tunable parameters
  - Synaptic weights, delays
  - Neuronal thresholds, leaks
  - Network architecture
- Promising solutions in designing energyefficient and area-constrained systems



Shruti Kulkarni, Aaron Young

## **Evolutionary Optimization for Neuromorphic Systems (EONS)**

- SNN Parameters:
  - Synaptic Weight and Delays
  - Neuronal Thresholds and leak
  - Network topology
- Often produces networks with very sparse connectivity
- Useful for co-design
- Can be slow to train



Schuman, Catherine D., et al. "An evolutionary optimization framework for neural networks and neuromorphic architectures." 2016 International Joint Conference on Neural Networks (IJCNN). IEEE, 2016.



#### Translating dataset to spikes



Shruti Kulkarni, Aaron Young

## **Initial training results**

#### • Dataset:

95% high pT samples, 5% low pT samples

#### • Training:

With 70k samples

Balanced dataset

SNN size: 688 neurons, 746 synapses

#### • Test results:

With ~650k samples (unbalanced)

F1 score: 0.7

Accuracy: 0.60





## Y-profile equivalent: "OR"



#### **Unbalanced dataset**

- Training samples: 2.6 million; Test samples: ~1.3 million
- Test accuracy: 0.62; Test f1 score: 0.73

Trair	ning set	
True labels/ Predicted	Class 0	Class 1
True low pT (0)	70,229	68,480
True High pT(1)	929,830	1,613,667

Test	set	
True labels/ Predicted	Class 0	Class 1
True low pT (0)	35,997	39,051
True High pT(1)	463,962	814,966





## Fitness function with penalty term

## Introduce penalty term that is high for misclassification of very high pT samples, and small for low pT (class 1) samples

Prediction error:

E = |pred - target|

EONS fitness:



#### Classification is heavily biased towards high pT samples



- Trained for 1000 epochs
- Network: 548 nodes, 202 edges
- Total Test samples: 1,353,976
- Testing f1 score: 0.92
- Testing Accuracy: 0.94

🛾 🛟 Fermilab

Shruti Kulkarni, Aaron Young

## Y-profile with spatial encoding

- SNN trained on samples with:
  - pT threshold at 0.5
  - y-local  $\in [-1.0, 1.0]$
- # inputs: 27 (26 + bias), #o/ps: 2
- Training Samples: 345,383

- Trained network:
  - # Neurons: 46, # synapses: 89
- Performance evaluated with pT threshold at 2GeV
- **Test Samples**: 173510
- Data Classification (>2GeV): 91.43%
- Data Rejection: 19.05%





#### **Neuromorphic Materials for edge computation**

Analog processing potentially enables lower energy, faster, and/or more complex on-pixel processing – or enables larger, more sophisticated neural networks



🛠 Fermilab

Patrick Xiao, Alec Talin, Sapan Agarwal

Co-

design

## Thank you

• Fermilab

Nhan Tran, Doug Berry, Lindsey Gray, Jennet Dickinson, Gauri Pradhan, Benjamin Parpillon, Guiseppe Di Guglielmo

• Northwestern University

Manuel Blanco Valentin, Seda Memik

• UIC

Corrinne Mills, Jieun Yoo, Amit Trivedi

• Johns Hopkins

Morris Swartz, Petar Maksimovic

• ORNL

Shruti Kulkarni, Aaron Young

Sandia National Lab

Patrick Xiao, Alec Talin, Sapan Agarwal



#### **Case Study II - 65 nm implementation**

- to understand challenges and strategies to overcome them
- non-reconfigurable, non-programmable NN for data compression



#### In-pixel AI/ML to enable Mfps camera operation

 $\Box$  Charge integration  $\rightarrow$  Digital Conversion  $\rightarrow$  Data Compression  $\rightarrow$  Off-Chip Data transfer

□ While ADC is converting the signal for this cycle, readout IC is transferring data from previous cycle

□ Single chip: (400 x 400 pixels x 10b x 1Mfps) 1.6Tbps data → 32Gbps enabled by 50x lossy data compression

□ 1x 32 Gbps photonic link or 3 x 10.24 Gbps link per chip (~5cm<sup>2</sup>)

	Charge Integration	Analog to Digital Conversion	Data Compression	Off-chip Data Transfer
Subchip = 1024 pixels for 100 Kfps operation	~1µs	~9µs	Full-readout	1024*10b @ 1GSPS = 10µs
Subchip = 1024 pixels for <b>100kfps operation</b>	~1µs	~9µs	Data compression PCA: 50x 1 cycle: 160n AE: 70x 30 cycles: 4.8µs	220b@ 1GSPS = 220n 150b@ 1GSPS = 150n

#### **Physics Background: Ptychography**



Huang, P., Du, M., Hammer, M., Miceli, A., & Jacobsen, C. (2021). Fast digital lossy compression for X-ray ptychographic data. Journal of synchrotron radiation, 28 (Pt 1), 292–300.

#### Ptychography: Imaging using coherent interference patterns in diffraction from object

- Large quantities of data with inherent redundancy

#### Solution: Image Compression for X-Ray Imaging

**Detector goal: 1 Mframe/s** 

This prototype: Frontend at 100 Kframes/s and 50 - 70x data compression in the pixel



#### Why data compression instead of zero-suppression?

Data rate:

## 400 × 400 (pixels per chip) × 10b (ADC) × 1 MHz $\sim$ 1.6 Tbps

#### **Data Sparsification:**

- Zero suppression (overhead address ~ 18-20b per 12b data)
- Simulated data (~97% zeros) vs. Noisy data (~60% zeros)











compression - Explored two approaches

#### **XROCKET Pixel frontend**

- Analog frontend consists of preamplifier + CDS + serial SAR
- Digital pixel backend consists of serial SAR logic
- Synchronous binary tree priority encoder







## **XROCKET**

Data rates and design choices

- Have to compress data on chip
- Spend ~ 1pJ/b for data transfer offchip (takes lower power to process the data on-chip)
- Should this be done in-pixel or at the periphery?





	1 pixel	32 x 32 pixels	400 x 400 pixels	Camera with 4 x 4 chips
Size	55 × 55 μm²	1.76 × 1.76 mm <sup>2</sup>	2.2 × 2.5 cm <sup>2</sup>	
Power	50µW	51.2mW	8W	128W
Raw data	10b/10μs = 1Mbps	1Gbps	1.6Tbps	25.6Tbps
With 70x compression		14Mbps	23 Gbps	360Gbps
				🕹 🕹 Fermi

#### 



## **Algorithm Co-Design**

#### Algorithm 1: Principal Component Analysis (PCA)



Calculated Weight Distribution, PCA

#### Goal: Reasonable Reconstruction of original diffraction image

#### 

- Analogous to principal components
- Weights are thus mathematically calculated: for inverse of eigenimages of diffraction input



#### PCA Algorithm: C++ Realization



Material: Matrix Factorization approaches to lost compression of ptychographic data. Journal of synchrotron radiation, Preprint, 43-50

#### Output Size: $30 \rightarrow$ based on factor of compression, here 50x

Going from  $1024 \times 10$ -bit to  $30 \times 7$ -bit

#### **PCA Loop Limits**

- Number of rows : 32 0
- Number of columns : 32 0
- Output Size Ο

- **Eigenimage dimensions**
- : 30  $\leftrightarrow$  Number of Eigenvalues
- Intermediate Data
  - 16-bit: PCA sensitive to loss of precision during accumulation
    - Inverse Eigenimage Matrix

with 4 x 4 pixels pixels chips TOPS-INT12 (data ~0.06 ~0.9 ~14 rate ~ 100 Kfps) TOPS-INT12 (data ~140 ~0.6 ~9 rate ~ 1 Mfps)

🛠 Fermilab

Google edge TPU: 4TOPS (Int 8)











Eigenimage 2

40

-10

-15



Eigenimage 3



PCA Reconstruction, N=20



PCA Reconstruction, N=100



PCA Reconstruction, N=5



PCA Reconstruction, N=30



PCA Reconstruction, N=300



PCA Reconstruction, N=10



PCA Reconstruction, N=50

PCA Reconstruction, N=1000



**‡** Fermilab



-20

-40

- 10

- 5

- 0

-5

L 10<sup>0</sup>

- 10<sup>1</sup>

- 10<sup>3</sup>

- 10<sup>2</sup>

#### Can you use eigen images from one training set for another data set



- The eigenimages are the basis of the diffraction patterns, not the initial images, and the characteristics of the diffraction pattern depend more on the experimental setup than the actual object.
- For the same experimental setup, the eigenimages from two objects are roughly the same and can be interchangeable
- Eigen values generated from the mandrill image as the training data
- PCA performed on Algae data set
- Acceptable image reconstruction performance
- OK to hardcode weights !!



mandrill: training data



#### Algorithm 2: AutoEncoder (AE)



#### Fully-Connected Dense Layer "Encoder" Matrix

- mimic PCA implementation structure
- Outputs remain as a 30-value latent space

#### **Weight Computation**

- Determined from Quantization-Aware model training

Output Size:  $30 \rightarrow$  based on factor of compression, here 70x

- Going from 1024 x 10-bit to 30 x 5-bit











#### AutoEncoder (AE) Algorithm

#### **Pre-Processing Inputs**

- reduce input precision (10'b ADC output) to 5'b: re-scaling
- Goal: Improve training of weights
  - Via reducing full-scale and thus increasing occupancy
- 5'b "Pseudo-square root" achieved by creating 3 gain regions: x32; x2; ÷2



#### Note INT6 vs. INT 12 for PCA – area improvement

#### Weight Comparison: AE vs. PCA

- Percentage of Zero-Value Weights: PCA is 77.98 % vs. AE is 8.69% -
- Weight Value Patterns per Output: none for PCA, regionalized for AE -





## Key takeaways

- Lossy data compression is acceptable for Ptychography due to the repetitive nature of the data
- Weights can be hardcoded since these don't need to be recomputed for new images
- Heat maps of the weights from PCA and AE indicate the different implementation methods might be needed
- We wanted to avoid dead areas in a large area sensor and avoid moving data across the large ROIC and chose to investigate in-pixel implementation instead of chip periphery



#### 



## High-Level Synthesis Implementation

#### **Overcoming Congestion**

- Fully-Dense Architecture for for both AE and PCA
- Challenge: Physical Interconnectivity of inputs and weights

#### Post synthesis results for DNN.syn\_opt12.v3 (beast1) Instance Cell count Cell area Net Area Comb. Buffer Flop Total area Inverter 165,211 41,214,600 852,903,339 DNN 553,451,040 299,452,299 485 828 640 1.871.280 24,536,520 Latch, Clock-Gate, Mem and Macro areas are omitted as they are all 0 -Pixel size to Ratio (H/W): 370073815 Core Utilization: 1.043757 50μm x 50 μm Cell Utilization: 104% -1598. occupancy 1589.0

#### Closeup of Congestion Map: zoomed around the chip's central regions

62 4/26/2023 Farah Fahim | In-Pixel AI

#### Increased pixel size to $75 \mu$ m x 75 $\mu$ m

**Fermilab** 

#### **HLS Directives - Experiments**

<pre>void top (a,b,c,d) {     func_A(a,b,11);     func_B(c,11,12);     func_C(12,d)     return d; }</pre>	Nuc.A Inc.B Nuc.C
< 8 cycles	(B) With Dataflaw Pinalining

Xilinx: Pipelining

	veid xep()-( for_methdar()=0(2)=0(1-) ( = e(0)= H(0)= 4(0) ) ] ]	i.		
Iterations				
	lafied Loop	Partially Un	rolled Loop	Unvolled Los
Read b(3) Read b(3	Read b(1) Read b(1)	Read b(3)	Read NO1	Read b(1
Read all Read all	Read of 1 Read of 0	Read c[2]	Read (C)	Read of t
		Read b(2)	Read 1075	Read b(2
		Read (2)	Read c(D)	Read of 2
Write a[2] Write a[2	Write a[1] Write a[2]			Read b(1
Write a[2] Write a[2	( washi washi			Read bit
Write a[3] Write a[3	( Write a[1] Write a[2]	•	•	Read b(1 Read c(1) Read b(1)
Write a[3] Write a[3	( Wolfer all) Wolfer all	ente alti		Read (c)) Read (c)) Read (c))
Write a(2) Write a(2	( Wolla a[1] Wolla a[2]	a a and a col and a col	- 	Read (2) Read (2) Read (2) Read (2)
With a(2) With a(2	( webs a()) webs a(d)	e Notice a(2) Notice a(2)	* * **********************************	Read 60 Read 60 Read 60
withe a(2) withe a(2	() Weiten aft) Weiten aft)	vinte a(2) vinte a(2)	e Normaliji Normaliji	Read b(1) Read c(1) Read c(1)
winte a(2) winte a(2	(	We file a (2) We file a (2)	nota ajti nota ajti	Read b(1) Read b(2) Read b(2)
With a(2) With a(2	(	a Withou (2) Withou (2)	nota ajti nota ajti	Read b(1) Read d(1) Read d(1)
withe a(3) within a(3	− − − − − − − − − − − − − − − − − − −	a Without (C) Without (C)	- 19-5a a(*) 19-5a a(*)	Read bit
with a [3] with a [3	() Werke a(t) Merke a(t)	e Teo teo a Ci Teo teo a Ci	antin afti antin afti	Read bit Read bit Rea
wine (3) wine (3	() Werka a(t) Kanka a(t)	e Victorați Victorați	- 19:30-4(-) 19:30-4(-)	Read bit Read bit Rea

#### Niansong Zhang: Unrolling

#### Trade latency for resources

#### **Inherent Data Patterns**

- Diffraction "rings": results from regional resources may cancel
- Reduce amount of latching/storage



#### 



# Overcoming Congestion

**RTL+Logic Synthesis Level** 

**Physical Design Level** 

#### **HLS Experiments: Design Conclusions**

HLS Solution	AE		PCA	
	Latency	Area (mm <sup>2</sup> )	Latency	Area (mm <sup>2</sup> )
modular	30	0.549	30	1.516
in-line	1	1.700	1	0.652

**AE** Restructured the weight buffers to perform 1,024 multiplication in parallel and efficiently pipeline them across each of the 30 output values

#### vs.

**PCA** Inlined all of the HLS-code functions and unrolled all the loops to take advantage of the sparsity of weights



## **Different approaches for AE & PCA**



AE one multiplier per pixel pipeline over 30 clock cycles

PCA: Constant propagation. Multipliers are replaced by fixed point adders



#### **Physical Design**

#### **Three Branches:**

- Optimizing Pitch between pixels
- Logic Flow: Accumulation
  - Initial RTL S tree (System Verilog)  $\rightarrow$  1D tree  $\rightarrow$  changed to H tree
- Pixel Arrangement to decrease routing lengths







2D hierarchy



Initial S-Tree Logic Flow

67 4/26/2023 Farah Fahim | In-Pixel AI

#### **Physical Design**

#### **Three Branches:**

- Optimizing Pitch between pixels
- Logic Flow: Accumulation
  - Initial RTL S tree  $\rightarrow$  changed to H tree
- Pixel Arrangement to decrease routing lengths







## Changes to HLS for optimization and better layout routing

```
1 #define H 32 // height
2 #define W 32 // width
3 #define C 30 // channels
5 void ae top(
     input_t inputs [H*W],
     weights_t weights [H*W][C], // Channel last
     biases_t biases[C],
8
     outputs_t outputs[C]) {
9
10
     accum_t accum[C];
11
12
     for (u32 \ j = 0; \ j < C; \ j++)
13
       accum[j] = biases[j];
14
     }
15
16
     for (u32 \ i = 0; i < H*W; i++)
17
       for (u32 \ i = 0; \ i < C; \ i++)
18
         accum[j] += inputs[i] * weights[i][j];
19
20
21
22
     for (u32 \ j = 0; \ j < C; \ j++) {
23
24
       outputs[j] = accum[j];
25
26 }
```

Multiplication last Creates a 1D hierarchy

```
1 #define H 32 // height
 2 #define W 32 // width
 3 #define C 30 // channels
 5 #define B 4 // block size
7 void ae_top(
     input_t inputs [H*W],
 8
     weights_t weights[C][H*W], // Channel first
9
     biases_t biases[C],
10
     outputs_t outputs[C]) {
11
12
     for (u32 \ i = 0; \ i < C; \ i++) \{ // Pipeline \}
13
       accum_t accum = biases[j];
14
       for (u32 \ i = 0; \ i < (H*W)/B; \ i++) \{ // Unroll \}
15
           accum t sub accum = 0.0;
16
           for (k = 0; k < B; k++) { // Submodule, unroll
17
                sub accum += inputs[i*B+k] *
18
                              weights[i][i*B+k];
19
20
           accum += sub_accum;
21
22
23
       outputs[j] = acc;
24
25 }
```

Multiplication first Create a 2D hierarchy



#### **Physical Design**

#### **Three Branches:**

- Optimizing Pitch between pixels
- Logic Flow: Accumulation
  - Initial RTL S tree  $\rightarrow$  changed to H tree
- Pixel Arrangement to decrease routing lengths



<sup>80%</sup> occupancy cutoff @pitch=112  $\rightarrow$  proceed with 110



#### With pitch=110, post-Place&Route:

- No congestion issues
- Logic distributed throughout quadrant
- DRC clean



#### What about a reprogrammable version

PCA –

 Massive multiplier size (12b multipliers much larger than 6b multiplier) – Maybe Not the best path forward

Autoencoder -

- only: 1024 multipliers (6b) already on chip
- Pixel area for encoding weights: 30x6b weights adding registers in the pixel area; (maybe SRAMs would be smaller) if pixel size grows to 75μm x 75μm

