

BinaryCoP: Binary Neural Network-based COVID-19 Face-Mask Wear and Positioning Predictor

Team Number: xohw21-142

Team Members: Nael Fafous¹, Manoj Rohit Vemparala², Alexander Frickenstein²

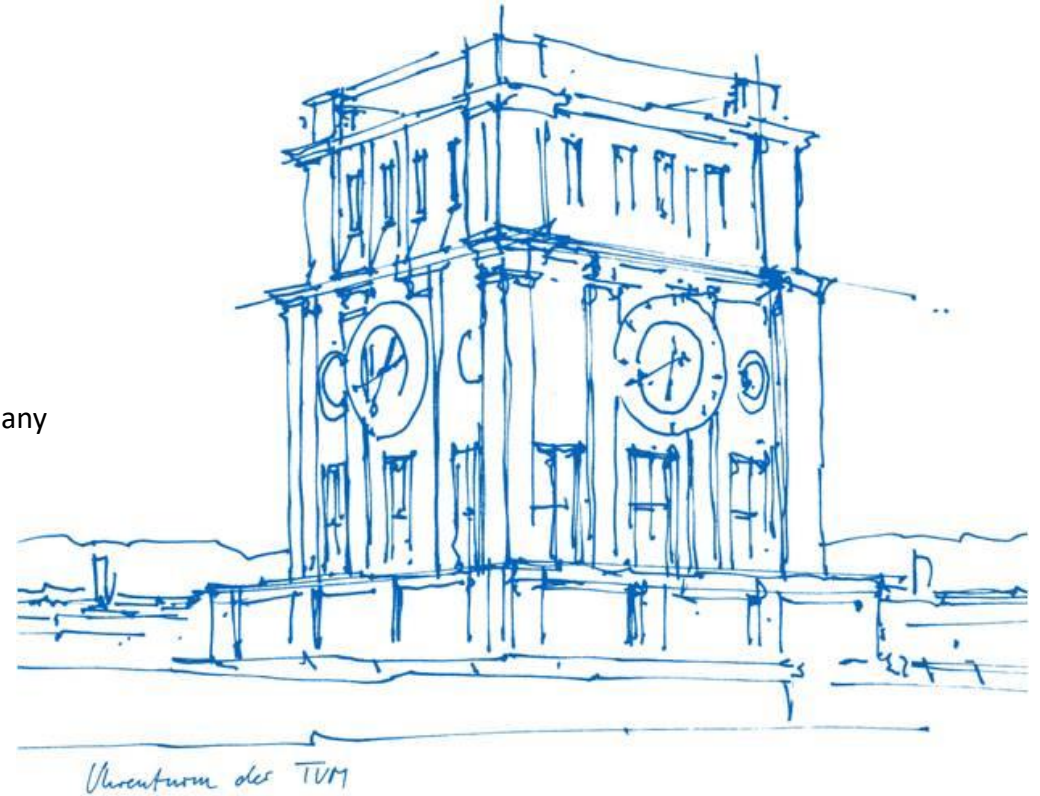
Supervisor: Walter Stechele¹

¹: Department of Electrical and Computer Engineering, Technical University of Munich, Munich, Germany

²: Autonomous Driving, BMW Group, Munich, Germany

¹{nael.fafous, walter.stechele}@tum.de

²{manoj-rohit.vemparala, alexander.frickenstein}@bmw.de



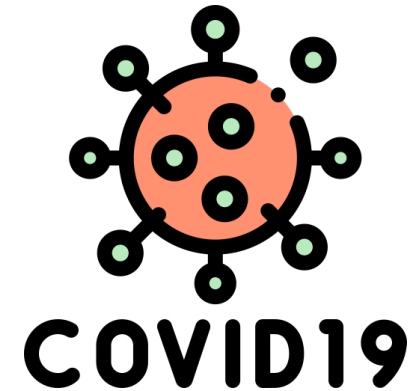
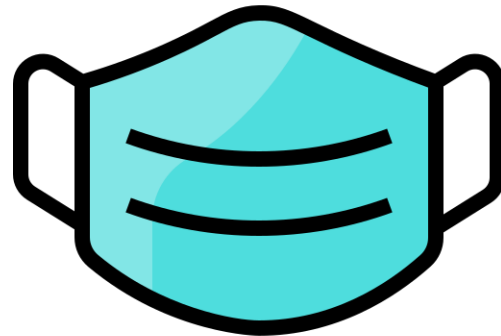
Outline



- **Problem Statement**
- **Common Challenges for AI Deployment**
 - Synthetic Data
 - Binary Neural Networks
 - Algorithm Interpretability
- **Hardware Dimensioning**
- **Comparison with Existing Work**

Problem Formulation

- Face-masks are a simple yet effective solution to mitigate the spread of COVID-19 [1]
- Most public indoor spaces have a mandatory rule on wearing masks
- **Compliance** to rules is hard to guarantee
- **Correct wearing position** of masks hard to assert

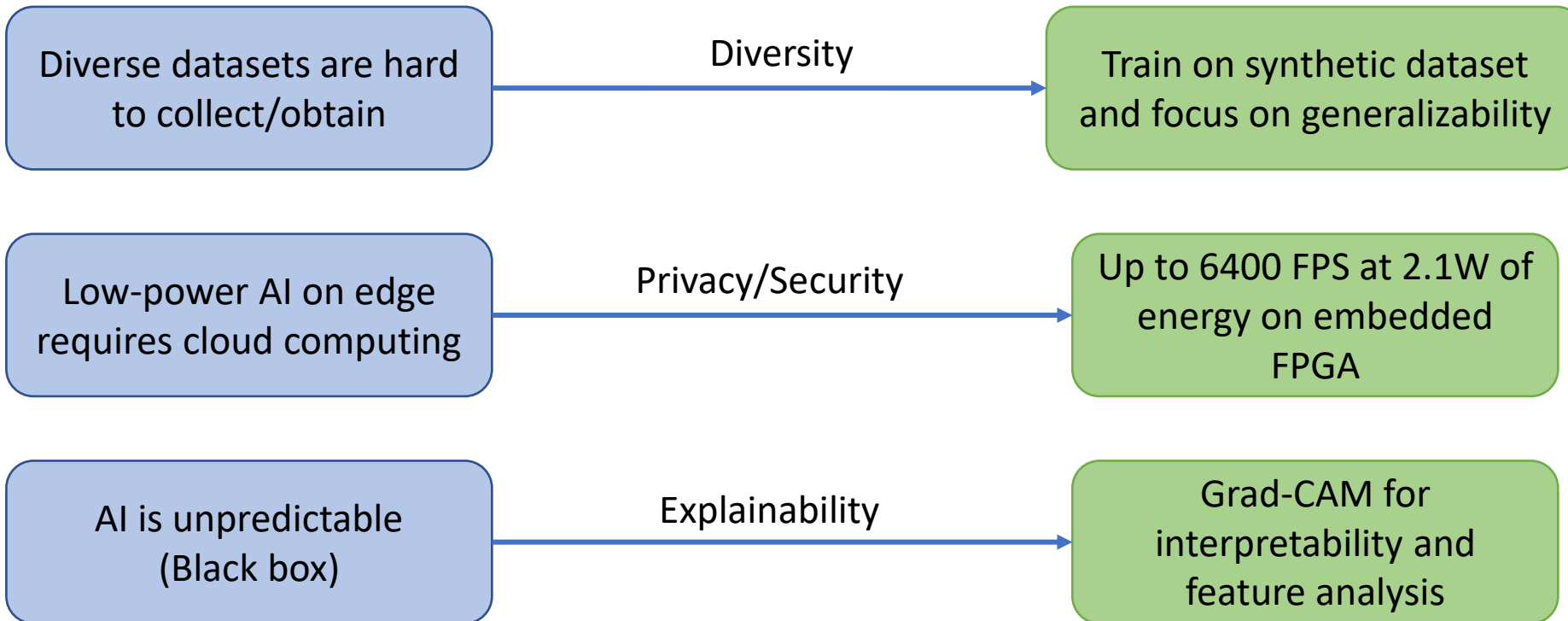


Problem Statement

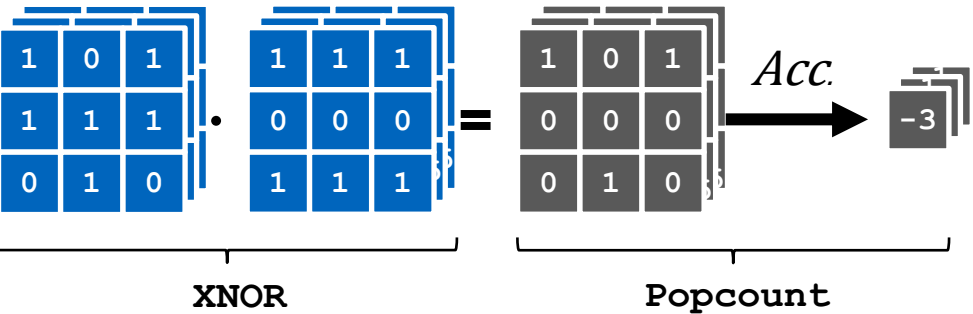


Deploy **accurate, unbiased** image classification algorithms, which can be used at entrances or speed-gates to check **correct mask wear and positioning**, with all processing on **low-power, cheap, edge** hardware to **preserve privacy** of passing users.

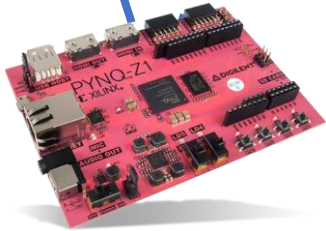
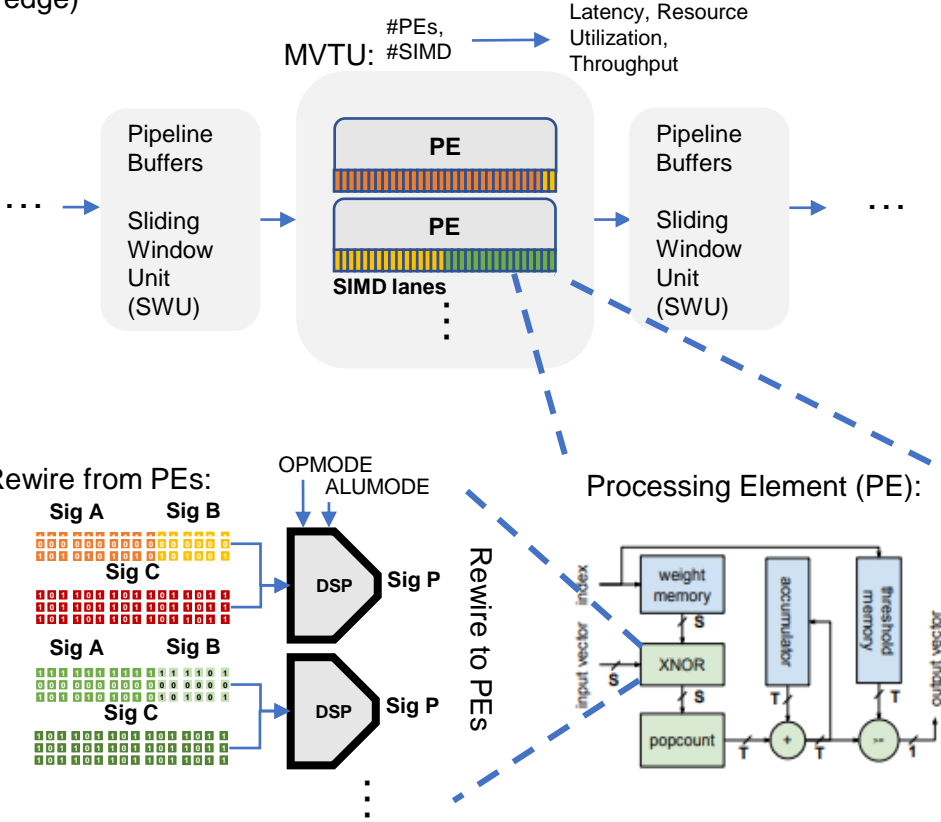
Common Challenges for AI Deployment



Binary Neural Network (low memory, good generalization on synthetic data)



FINN-based Accelerator (privacy preserving, high-performance on edge)



PYNQ board (or PYNQ -based)

Synthetic Data Training and Interpretability

Label	Raw	BCoP CNV	Label	Raw	BCoP CNV
Correctly Masked			Nose Exposed		
Correctly Masked			Nose Exposed		
Correctly Masked			Nose Exposed		
Chin Exposed			Chin Exposed		
Chin Exposed			Chin Exposed		
Chin Exposed			Chin Exposed		

- Avoid Region-Local Data
- Maintain Subject Diversity
- Assert correct features being learnt

Dataset Diversity

- **Diversity**
 - Large-scale **natural face** datasets exist
 - But the number of real-world **masked face images** is limited

Large dataset predominantly collected in East Asia [2]



Maintain
Diversity



Synthetically generated on top of natural images [3]



Dataset Diversity

Should we sacrifice real-world accuracy for diversity?

Synthetic data generation can be improved to close the gap [4]

BNNs have proven regularization characteristics [5]

Initialize on synthetic data,
then fine-tune on real-world data

Synthetically added mask type: [4]



Mask pattern variation:



Mask color and intensity:



Binary Neural Networks

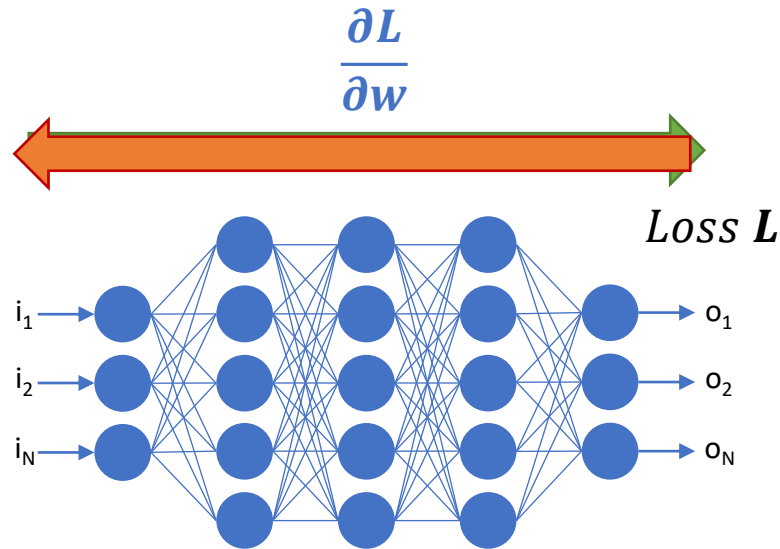


- Restrict all weights w and activations a to $\{-1, 1\}$ the *sign* function:

$$b = \text{sign}(w) = \begin{cases} -1, & w < 0 \\ 1, & w \geq 0 \end{cases}$$

- Where w is the **latent weight** representation (backward-pass) and b is the **binarized weight** (forward-pass)
- **Latent weights** w are represented in **full-precision** to allow fine adjustments during training
- In final deployment, only b **binarized weights** are kept.

Binary Neural Networks: Gradient Flow Problem



- Inputs $i \in I$ result in outputs $o \in O$
- Outputs $o \in O$ result in a loss L when compared to true class of I
- Weights $w \in W$ are updated in backpropagation to minimize L
- Stochastic gradient decent (SGD) applies the chain rule to update each weight for a mini-batch input:

Update w to minimize L
i.e. find minimization slope:

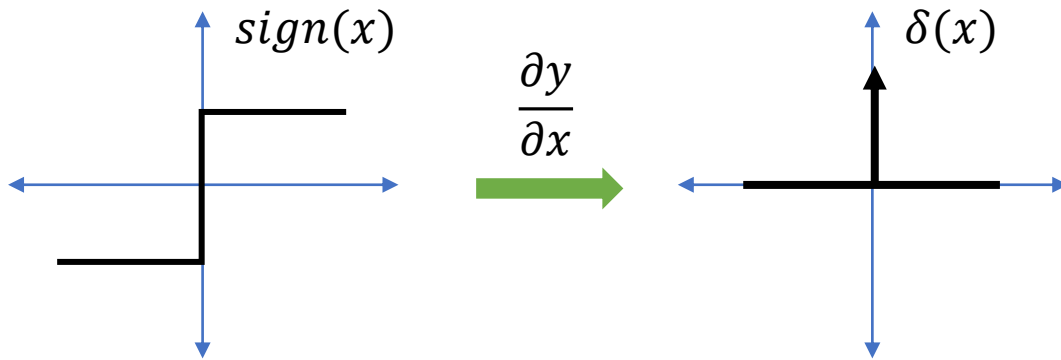
$$\frac{\partial L}{\partial w} = \frac{\partial v}{\partial w} \boxed{\frac{\partial z}{\partial v} \frac{\partial k}{\partial z}} \frac{\partial L}{\partial k}$$

Intermediate operations between L and w

Binary Neural Networks

- The *sign* function is an intermediate operation in BNNs:

$$y = \text{sign}(x) = \begin{cases} -1, & w < 0 \\ 1, & w \geq 0 \end{cases}$$



$$\frac{\partial L}{\partial w} = \frac{\partial v}{\partial w} \boxed{\frac{\partial z}{\partial v} \frac{\partial k}{\partial z}} \frac{\partial L}{\partial k}$$

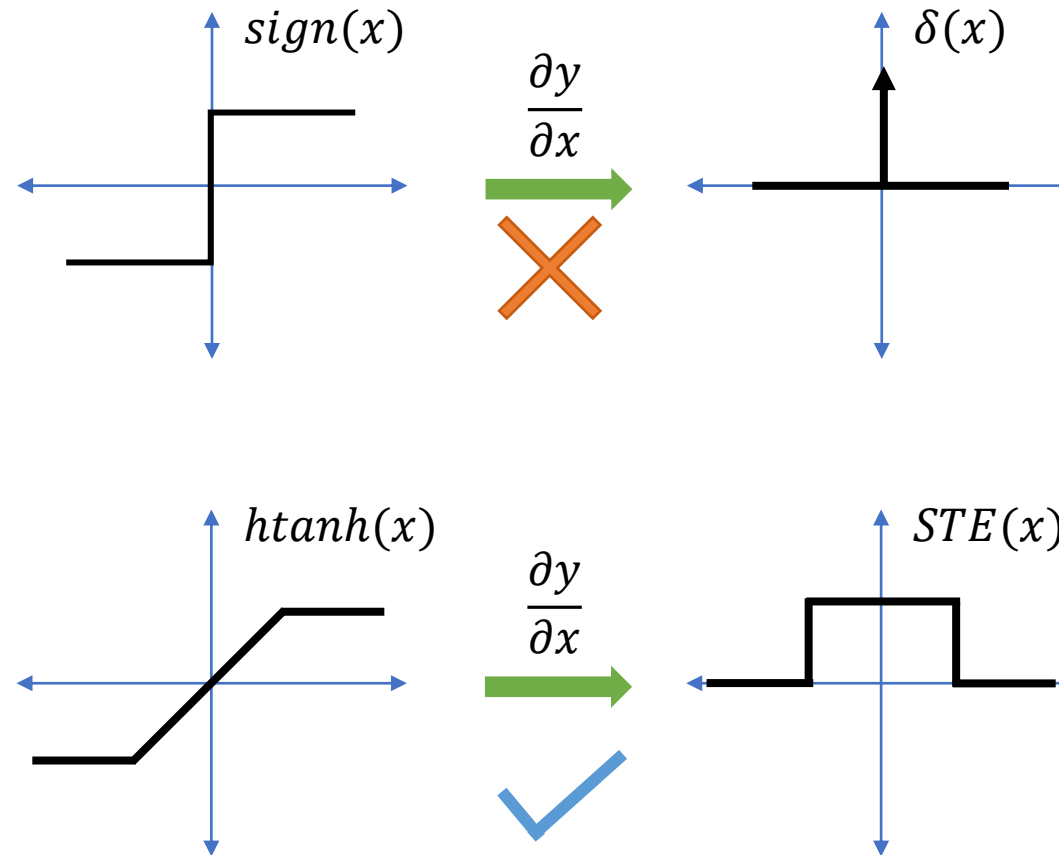
Intermediate operations between L and w

Derivative of *sign* blocks gradient flow
because it is **zero almost everywhere!**

Discreteness of binarization is incompatible with **Gradient Decent optimization**

Binary Neural Networks: Gradient Flow Problem

- Solution: **approximate** derivative of $sign$ with $htanh$ (or other):



Gradient flow
possible again!

Binary Neural Networks

- So what do we get for all this?

High-Precision

Binary

Multiplications



XNOR

Accumulation



PopCount

BatchNorm



Comparator

32 or 16 Bits/Operand



1-Bit Operands

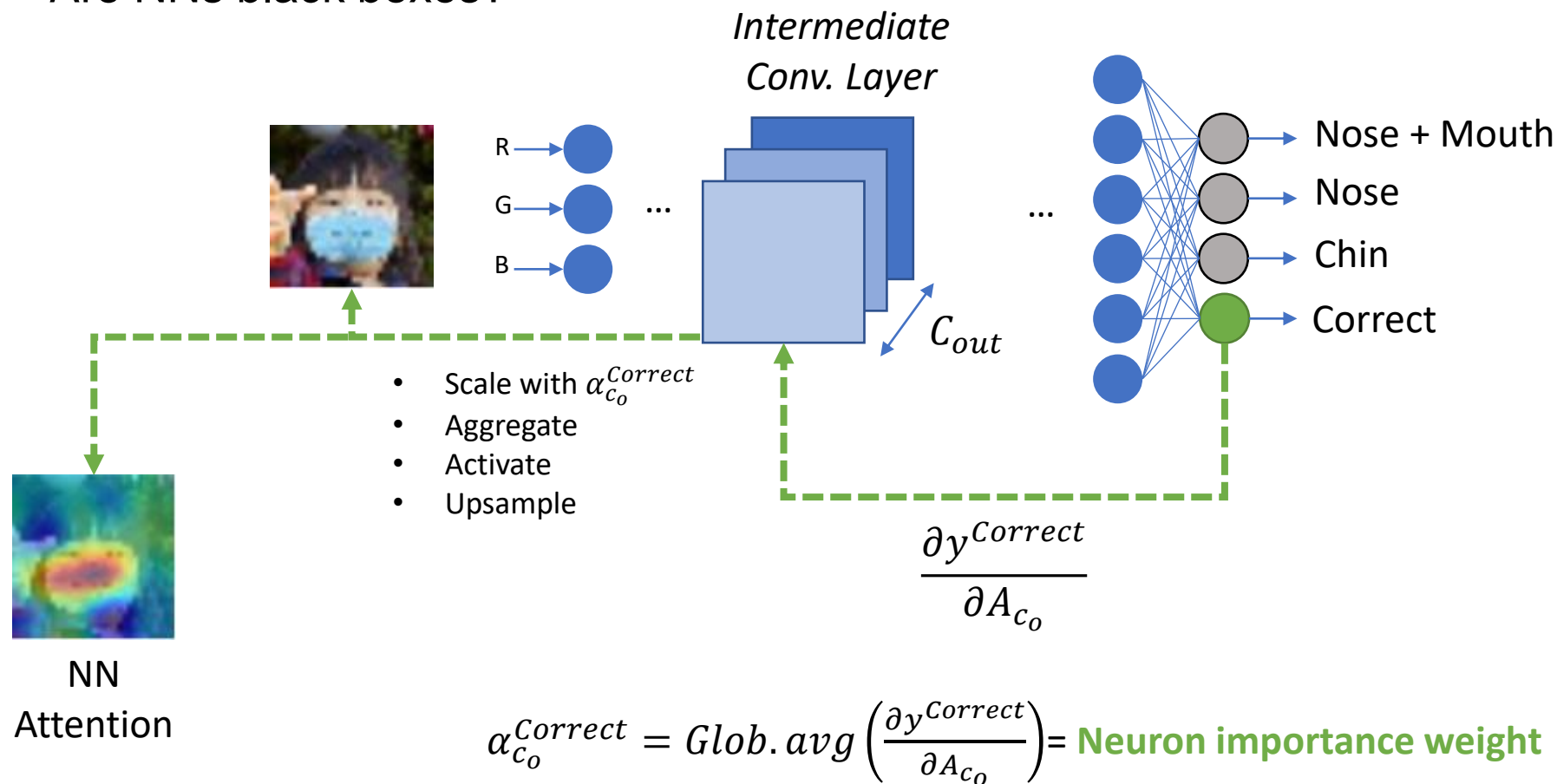
- **Bonus:**

BNNs' approximate training makes it harder to overfit on training data

Good candidate for training on Synthetic Data

Algorithm Explainability: Grad-CAM [6]

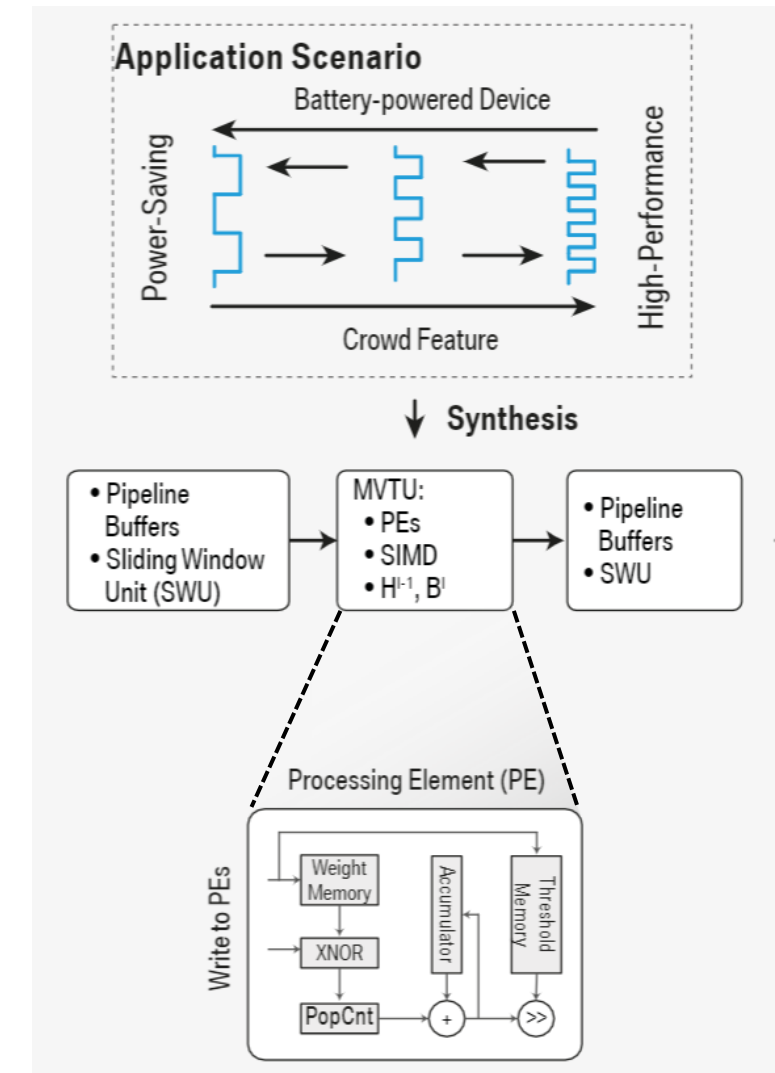
- Are NNs black boxes?



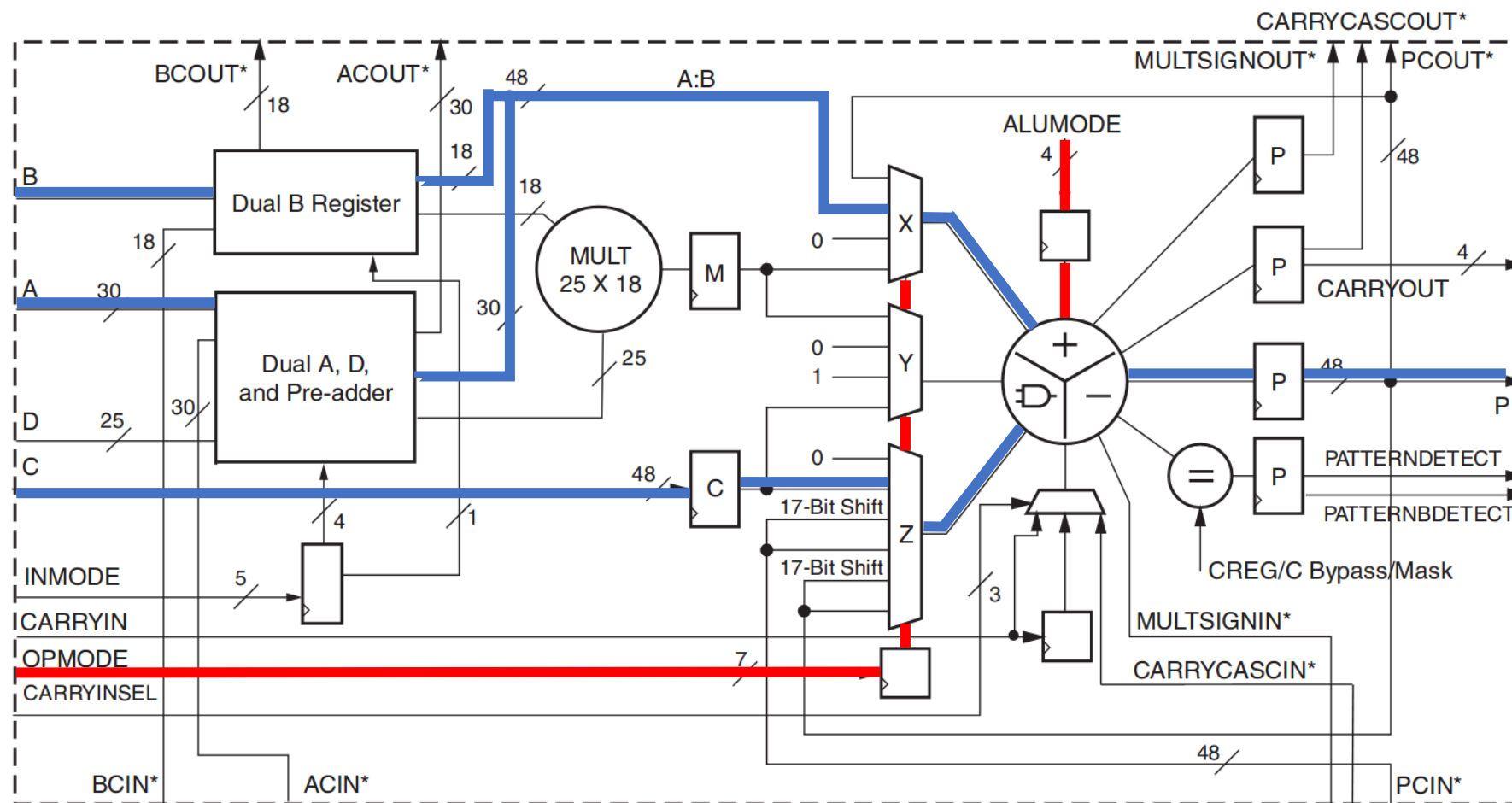
NN's learning can be observed to understand its decisions and learning patterns

Hardware Acceleration: FINN framework [7]

- BNNs can have a very small memory footprint
 - Accelerator can pre-load all weights
 - Computation units (MVTUs) can be parameterized (PEs, SIMD)
 - Match-throughput of MVTUs



Hardware Acceleration: Bit-Packing XNORs into DSPs



*These signals are dedicated routing paths internal to the DSP48E1 column. They are not accessible via fabric routing resources.

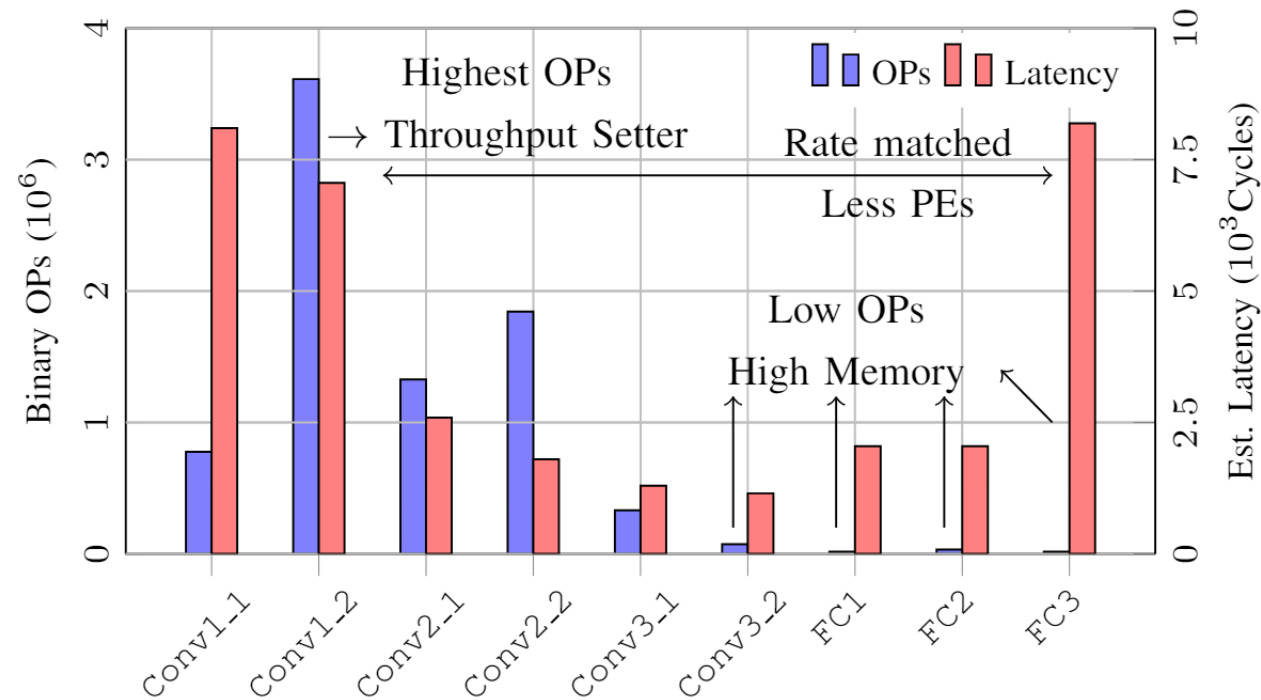
BinaryCoP Prototypes

- 3 BinaryCoP variants

	High Accuracy	Low Memory	Low Logic
Network	CNV	n -CNV	μ -CNV
Arch. $L \mid [C_i, C_o]$ $K = 3 \forall \text{ Conv}$	Conv1_1 \mid [3, 64] Conv1_2 \mid [64, 64] Conv2_1 \mid [64, 128] Conv2_2 \mid [128, 128] Conv3_1 \mid [128, 256] Conv3_2 \mid [256, 256] FC1 \mid [512] FC2 \mid [512] FC3 \mid [4]	Conv1_1 \mid [3, 16] Conv1_2 \mid [16, 16] Conv2_1 \mid [16, 32] Conv2_2 \mid [32, 32] Conv3_1 \mid [32, 64] Conv3_2 \mid [64, 64] FC1 \mid [128] FC2 \mid [128] FC3 \mid [4]	Conv1_1 \mid [3,16] Conv1_2 \mid [16, 16] Conv2_1 \mid [16, 32] Conv2_2 \mid [32, 32] Conv3_1 \mid [32, 64] FC1 \mid [128] FC2 \mid [4]
PE Count SIMD lanes	16, 32, 16, 16, 4, 1, 1, 1, 4 3, 32, 32, 32, 32, 32, 4, 8, 1	16, 16, 16, 16, 4, 1, 1, 1, 1 3, 16, 16, 32, 32, 32, 4, 8, 1	4, 4, 4, 4, 1, 1, 1 3, 16, 16, 32, 32, 16, 1

BinaryCoP Prototypes

- Throughput Matching – Weight Memory Fragmentation
 - Example: BinaryCoP- n -CNV



Prototype Hardware Results

- Masked-FaceNet Dataset
- 32x32 input resolution

Prototype	LUT	BRAM	DSP	Power [W]		Thr.put [FPS]	Latency [ms]	Acc. [%]
				Idle	Inf.			
CNV	26060	124	24		2.212	3049	1.58	98.10
<i>n</i> -CNV	20425	10.5	14	1.65*	2.122	6460	0.31	93.94
μ -CNV	11738	14	27		2.028	1646	0.81	93.78

*Required by the board and ARM-Cortex A9 processor. Accelerator is idle.

Low Power

High Accuracy

Low Logic
(Synthesizable on Z7010)

Low Memory
High Throughput


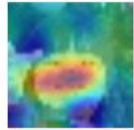

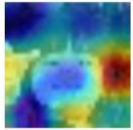

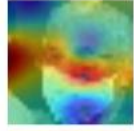
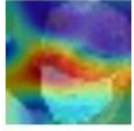
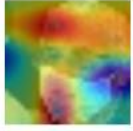




Explainability Results

- Confusion matrix of BinaryCoP-CNV



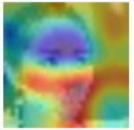




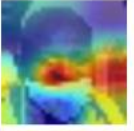

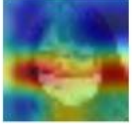
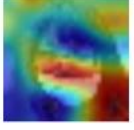
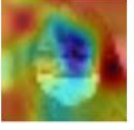
		Predicted Class			
		Correct	Nose	N+M	Chin
True Class	Correct	7125 98%	41 1%	1 0%	90 1%
	Nose	26 0%	7042 98%	94 2%	26 0%
	N+M	4 0%	79 1%	5651 98%	9 0%
	Chin	107 1%	41 1%	7 0%	7363 98%

Chin area too small,
can be overlooked by BNN







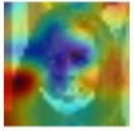
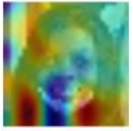


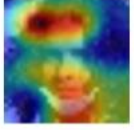
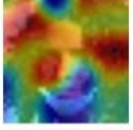
Explainability Results: Grad-CAM Interpretation

Label	Raw	BCoP CNV	BCoP n -CNV	FP32
Correctly Masked				
Correctly Masked				
Correctly Masked				


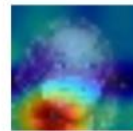
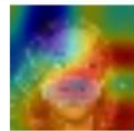
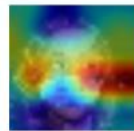



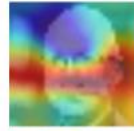

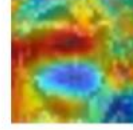
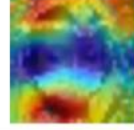
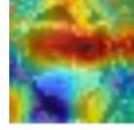
(a) Correctly masked Grad-CAM

Label	Raw	BCoP CNV	BCoP n -CNV	FP32
Nose Exposed				
Nose Exposed				
Nose Exposed				

(b) Nose exposed Grad-CAM

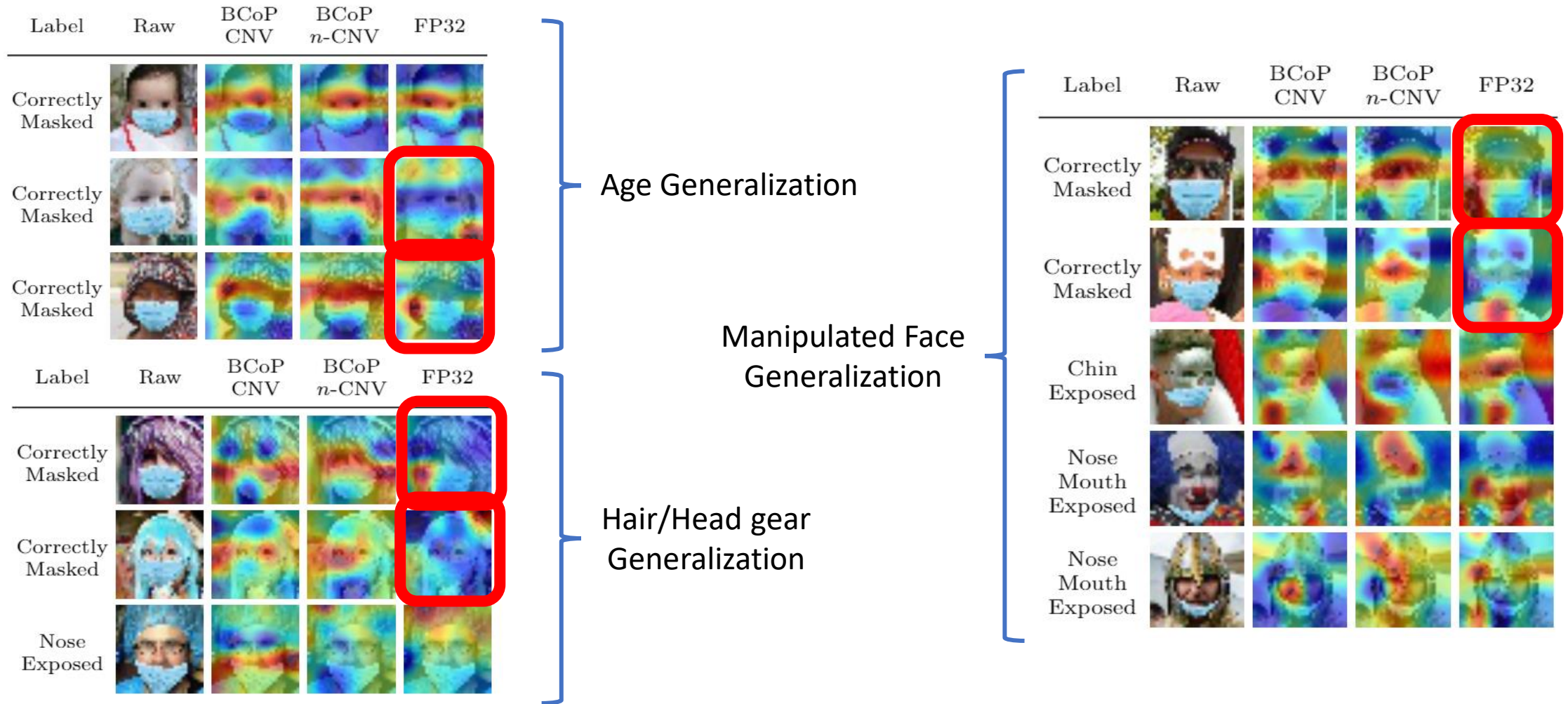
Label	Raw	BCoP CNV	BCoP n -CNV	FP32
Nose Mouth Exposed				
Nose Mouth Exposed				
Nose Mouth Exposed				

(c) Mouth + nose exposed Grad-CAM

Label	Raw	BCoP CNV	BCoP n -CNV	FP32
Chin Exposed				
Chin Exposed				
Chin Exposed				

(d) Chin exposed Grad-CAM

Explainability Results: Diversity and Generalizability



Demo

Input:

BinaryCoP Output:

Mode: ⚡ Performance

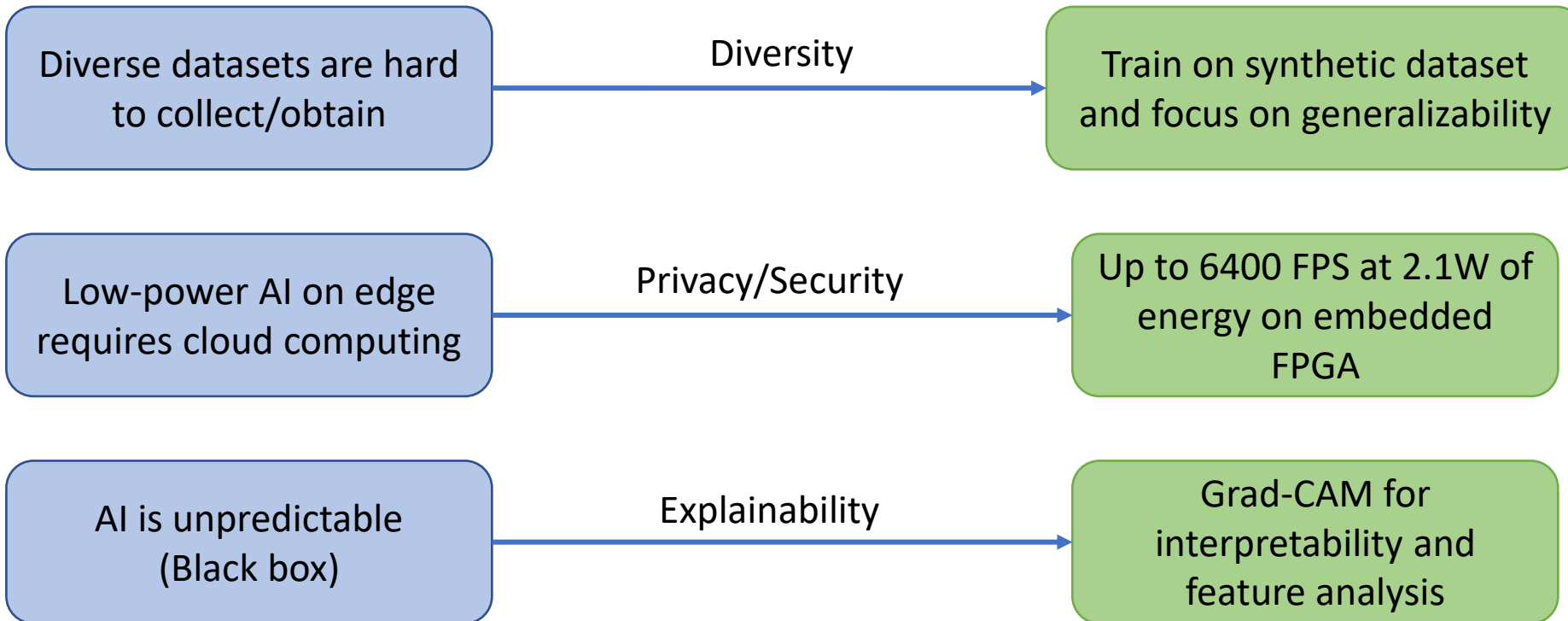
```
root@pynq:/home/xilinx/jupyter_notebooks/MASKEDFACE-BNN-PYNQ# python3 BinaryCoP-Video-DEMO.py
```

Comparison with Other Works

- Dataset for Mask Wear problem are not standardized
- Tasks range from classification to object detection
- Resolutions range from 960×544 to 32x32

Work	Pro	Con
NVIDIA [8]	<ul style="list-style-type: none">• High Res• Object Localization	<ul style="list-style-type: none">• 500€ Hardware• 20-30 Watts
Amazon [9]	<ul style="list-style-type: none">• High Res• Other PPE supported	<ul style="list-style-type: none">• Cloud Processing
Wang et al. [10]	<ul style="list-style-type: none">• Serverless, In-browser Processing	<ul style="list-style-type: none">• Requires Web-assembly Browser (tested on iPhone, iPad, MacBook)• \$\$\$ Expensive Devices
"CheckYourMask" Hammoudi et al. [11]	<ul style="list-style-type: none">• Edge Processing through Android app	<ul style="list-style-type: none">• Designed for self-check, not surveillance• Not optimized for power, low-battery application

Conclusion



Resources and References



- [1] When and how to use masks. [Online]. Available: <https://www.who.int/emergencies/diseases/novel-coronavirus-2019/advice-for-public/when-and-how-to-use-masks>
- [2] <https://github.com/X-zhangyang/Real-World-Masked-Face-Dataset>
- [3] <https://github.com/cabani/MaskedFace-Net>
- [4] <https://github.com/aqeelanwar/MaskTheFace>
- [5] M. Courbariaux et al. “Binaryconnect: Training deep neural networks with binary weights during propagations”
- [6] R. Selvaraju et al., “Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization”
- [7] Y. Umuroglu et al., “FINN: A Framework for Fast, Scalable Binarized Neural Network Inference”
- [8] A. Kulkarni et al., “Implementing a real-time, ai-based, face mask detector application for covid-19”
- [9] T. Agrawal et al., ““Automatically detecting personal protective equipment on persons in images using amazon rekognition”
- [10] Z. Wang et al, “Wearmask:Fast in-browser face mask detection with serverless edge computing for covid-19”
- [11] K. Hammoudi et al., “Validating the correct wearing of protection mask by taking a selfie: Design of a mobile application “checkyourmask” to limit the spread of covid-19”

Slide 3 has been designed using resources from [Flaticon.com](https://www.flaticon.com)