# At the Edge of AI at the Edge
## Ultra efficient AI on low-power devices
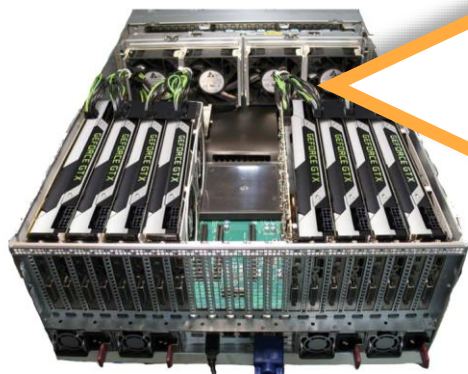
Mohammad Rastegari

mohammad@xnor.ai

May 2018

XNOR.AI

# AI is confined to the cloud

far from the users at the edge

XNOR.AI bridges the growing divide between AI models dependent on the cloud and devices running at the edge
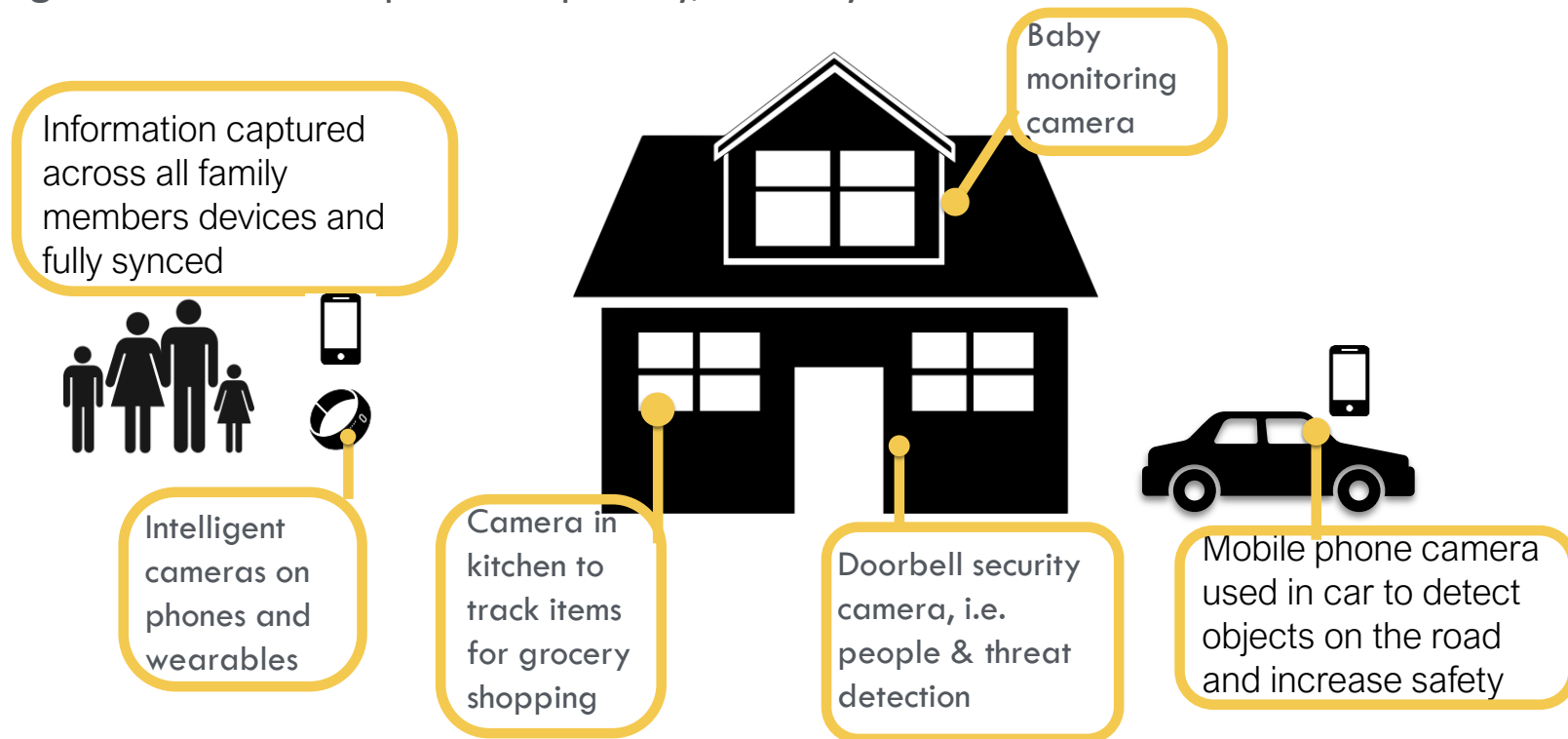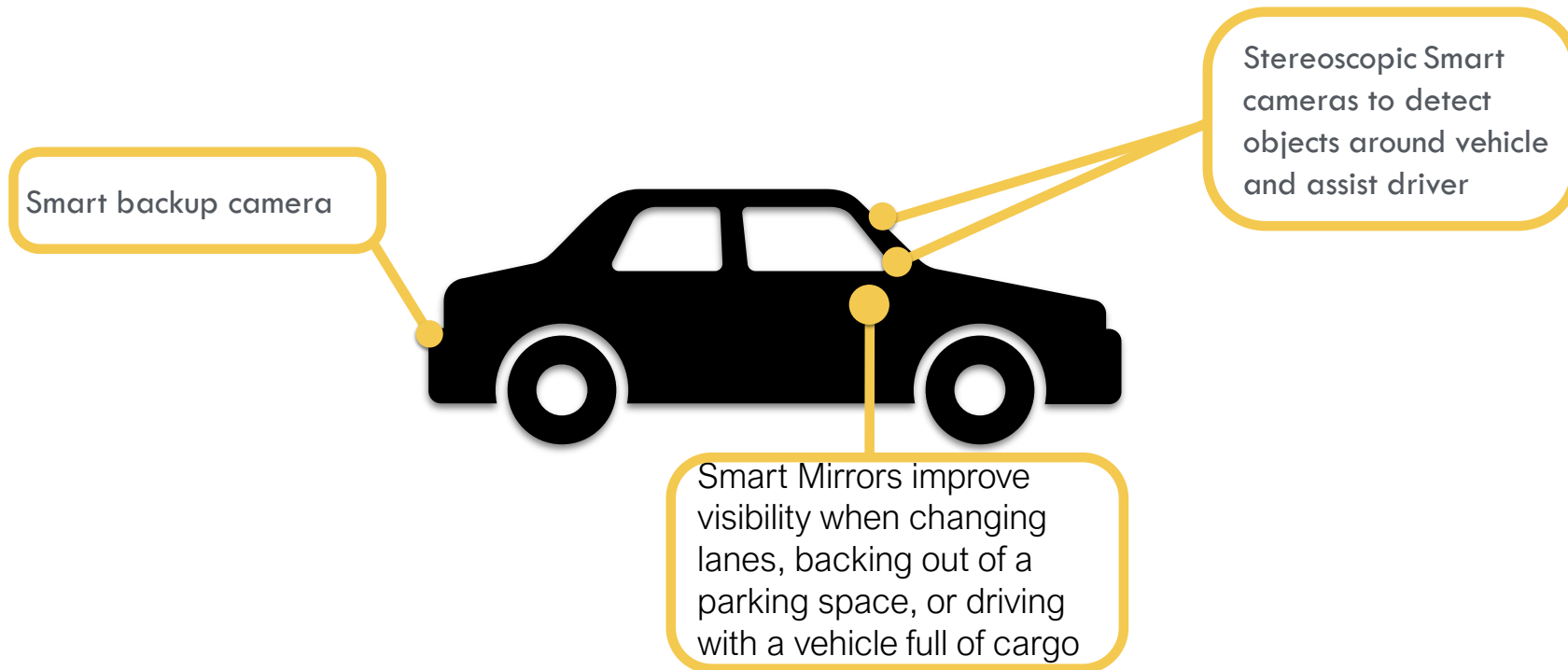
Deep learning models reliant on the cloud
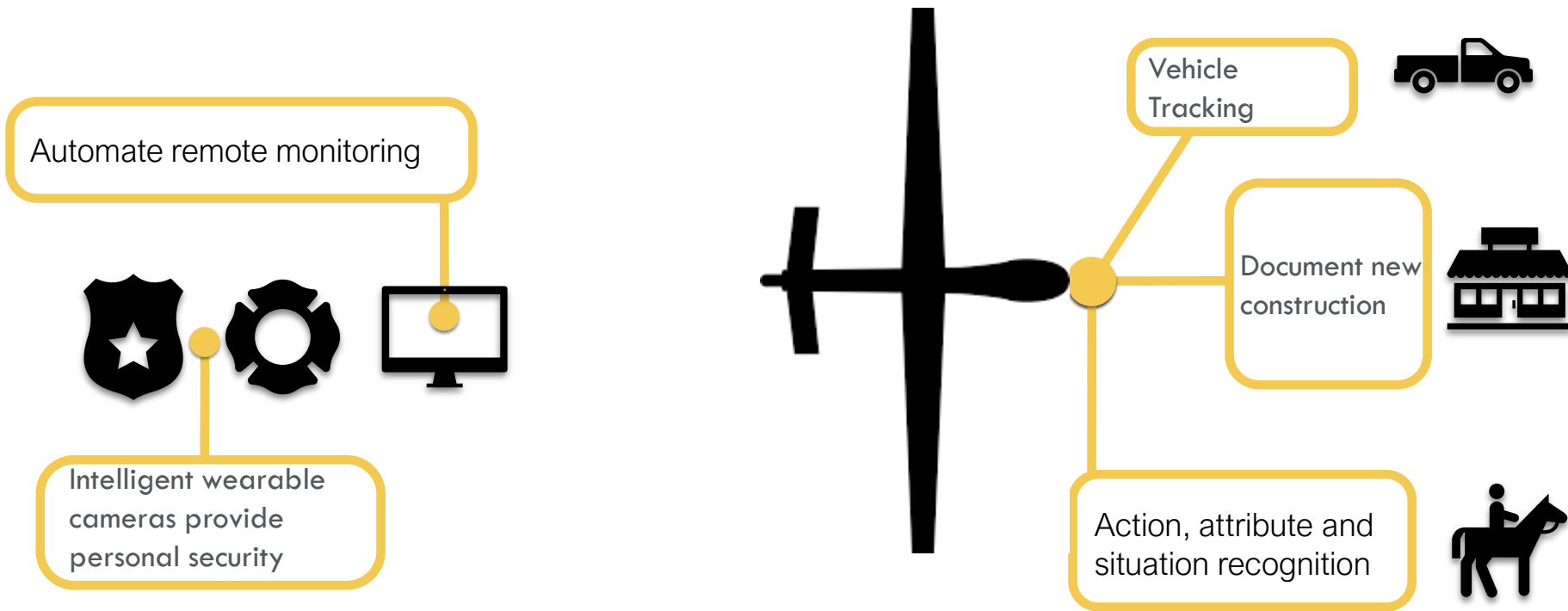
Growing demand for edge devices

Intelligent cameras that preserve privacy, security and bandwidth at home

Baby monitoring camera

Information captured across all family members devices and fully synced

Intelligent cameras on phones and wearables

Camera in kitchen to track items for grocery shopping

Doorbell security camera, i.e. people & threat detection

Mobile phone camera used in car to detect objects on the road and increase safety

Intelligent cameras for Advanced Driver Assistance Systems (ADAS)



Smart backup camera

Stereoscopic Smart cameras to detect objects around vehicle and assist driver

Smart Mirrors improve visibility when changing lanes, backing out of a parking space, or driving with a vehicle full of cargo

Intelligent cameras provide real-time tracking, recognition & detection on device

Automate remote monitoring

Intelligent wearable cameras provide personal security

Vehicle Tracking

Document new construction

Action, attribute and situation recognition

# What we do

## Image Tagging



Man
Dog
Tree
Grass
Park

## Image Enhancement



## Tracking



## Action Recognition



Jogging    Running

## Object Detection



## Scene Recognition



## Segmentation





NATURAL
LANGUAGE
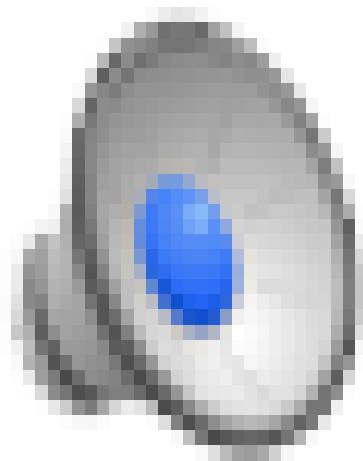
# Object Detection: An Expensive Task in AI

# Fine Grain Categories

XNOR.AI provides fast and efficient AI at the edge



Server GPU $1,000 · Desktop CPU $500 · Mobile CPU $600 · Raspberry Pi-3 $35 · NonoPi-A64 $25 · Pine-64 $15 · Raspberry Pi-0 $5

State-of-the-Art AI: all the way to Pi Zero

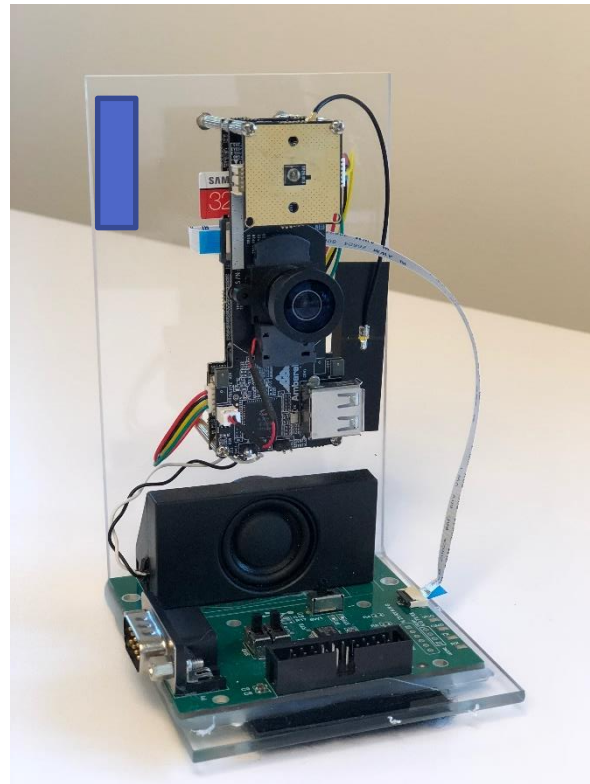XNOR $5 deep learning machine…
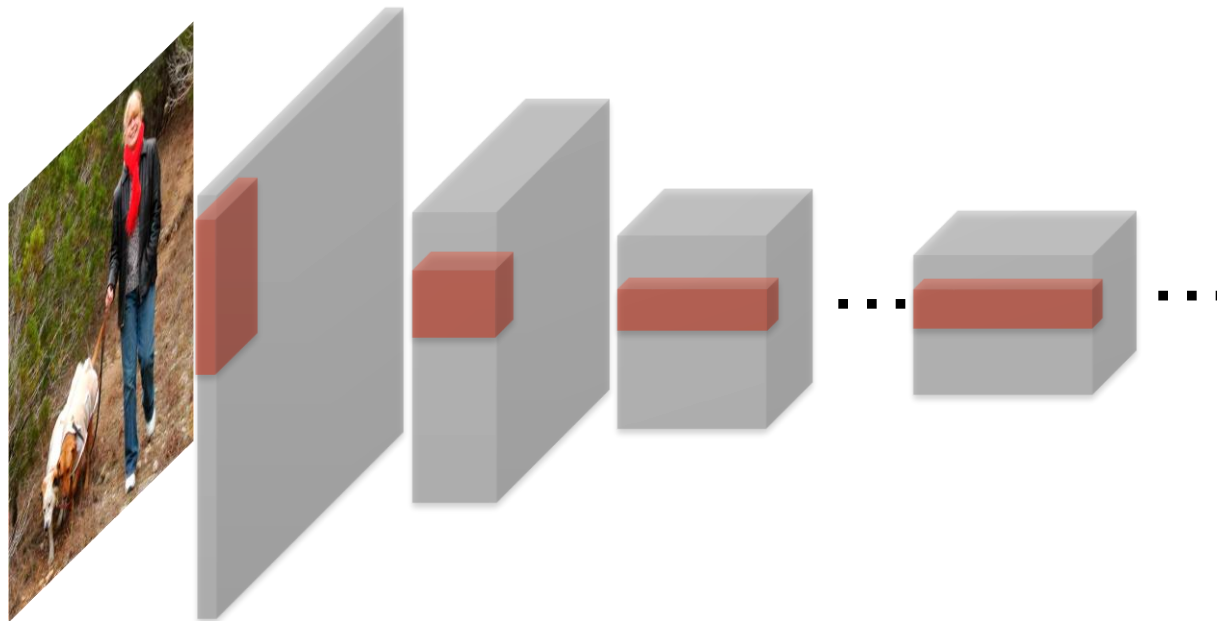… on Raspberry Pi Zero

Modular AI at Edge

**Ambarella S5L**

- **Very low power (~2x lower than Pi Zero)**
- **Standard AI model for object detection**
  - **1 fps**
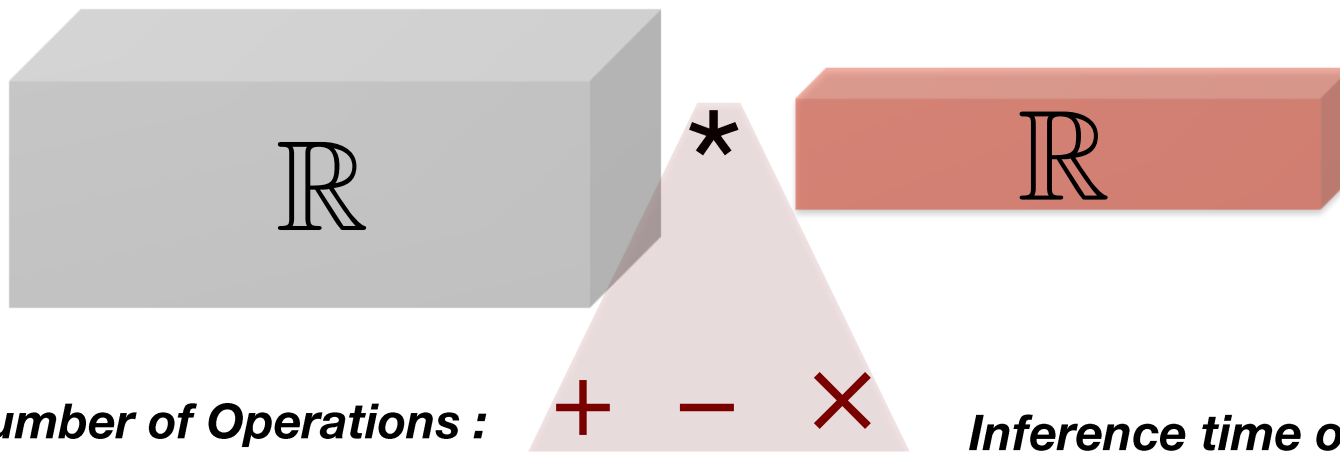- **XNOR AI Model for object detection**
  - **17 fps**

How our technology works

# Convolutional Neural Network

# GPU !



**Number of Operations :**

- AlexNet →1.5B FLOPs
- VGG → 19.6B FLOPs

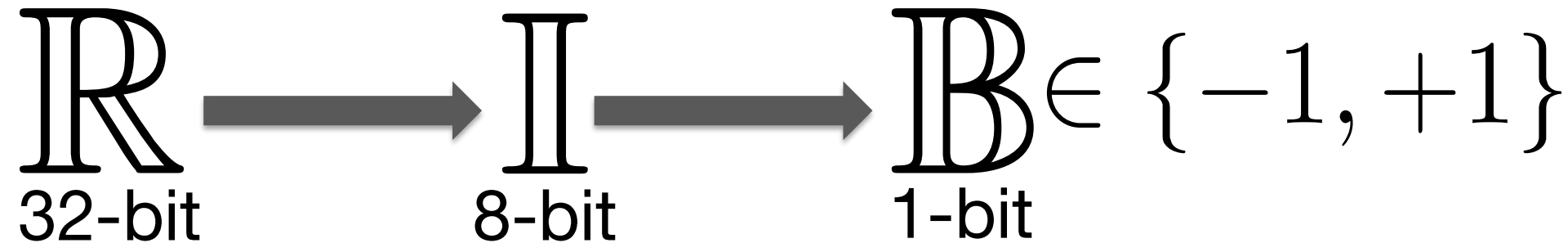**Inference time on CPU :**

- AlexNet →~3 fps
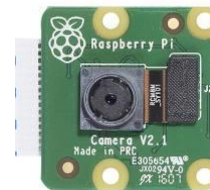- VGG → ~0.25 fps

# Lower Precision

Reducing Precision
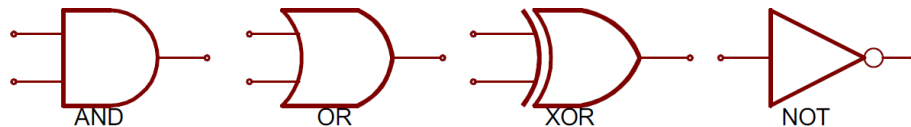- Saving Memory
- Saving Computation

$$\mathbb{R} \longrightarrow \mathbb{I} \longrightarrow \mathbb{B} \in \{-1, +1\}$$

32-bit          8-bit          1-bit

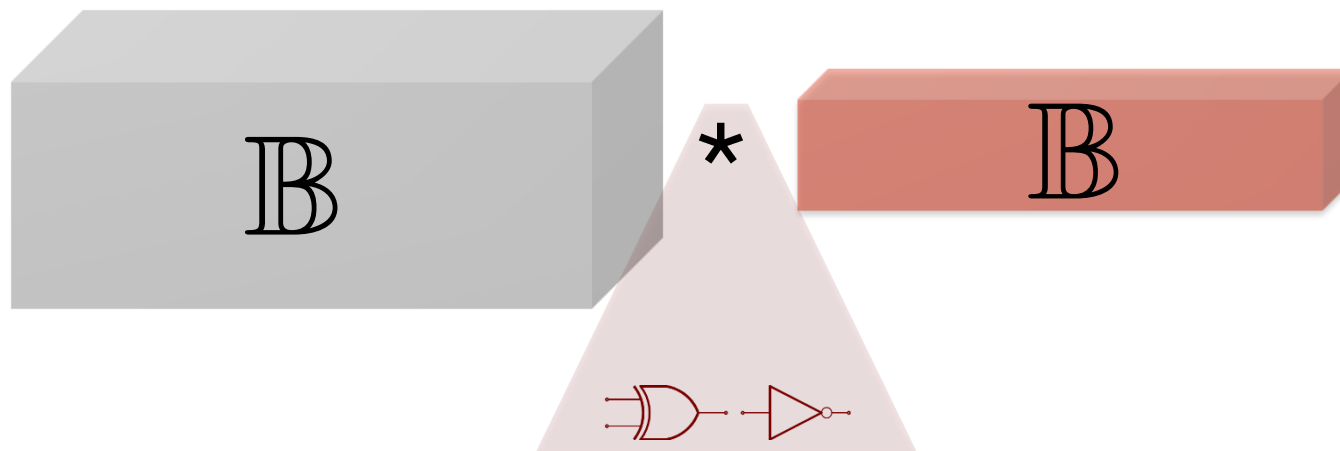| {-1,+1} | {0,1} |
|---------|-------|
| MUL | XNOR |
| ADD, SUB | Bit-Count (popcount) |

# Why Binary?

- Binary Instructions
  - AND, OR, XOR, XNOR, PopCount (Bit-Count)

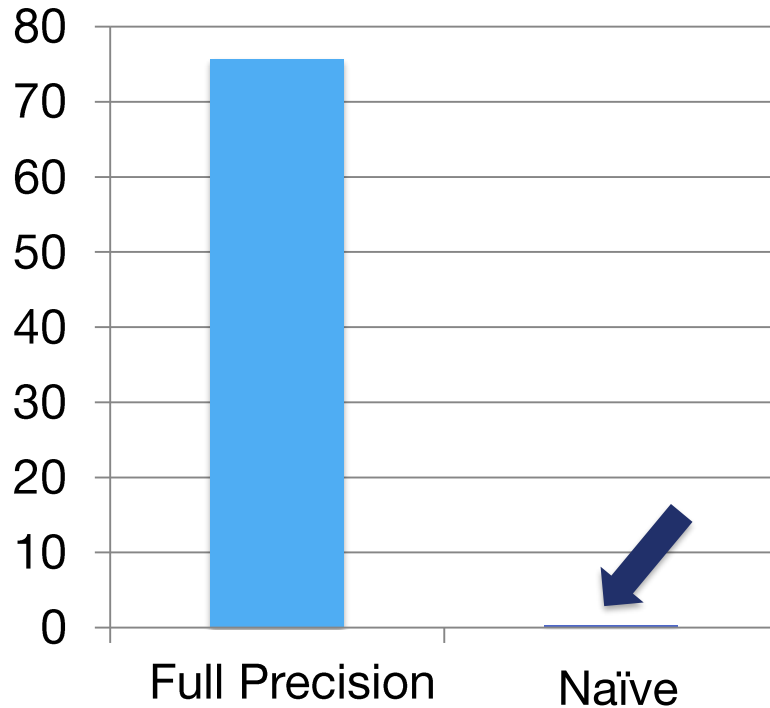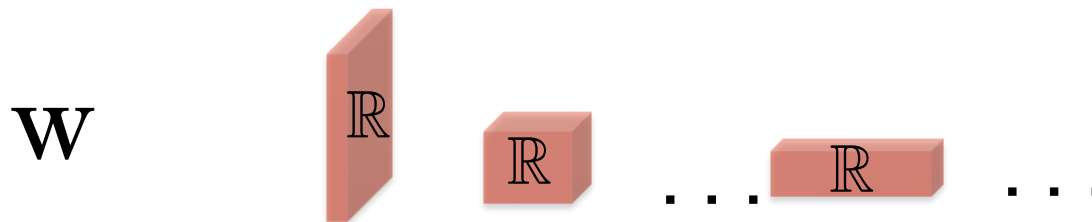- Low Power Device

- Easy to Implement in hardware
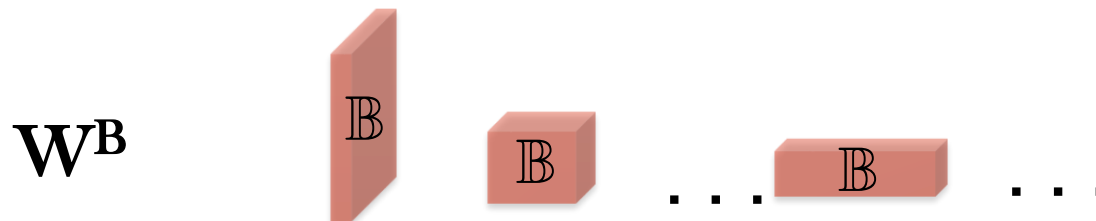
**Training Binary Weight Networks**

*Naive Solution:*

1. Train a network with real value parameters
2. Binarize the weight filters

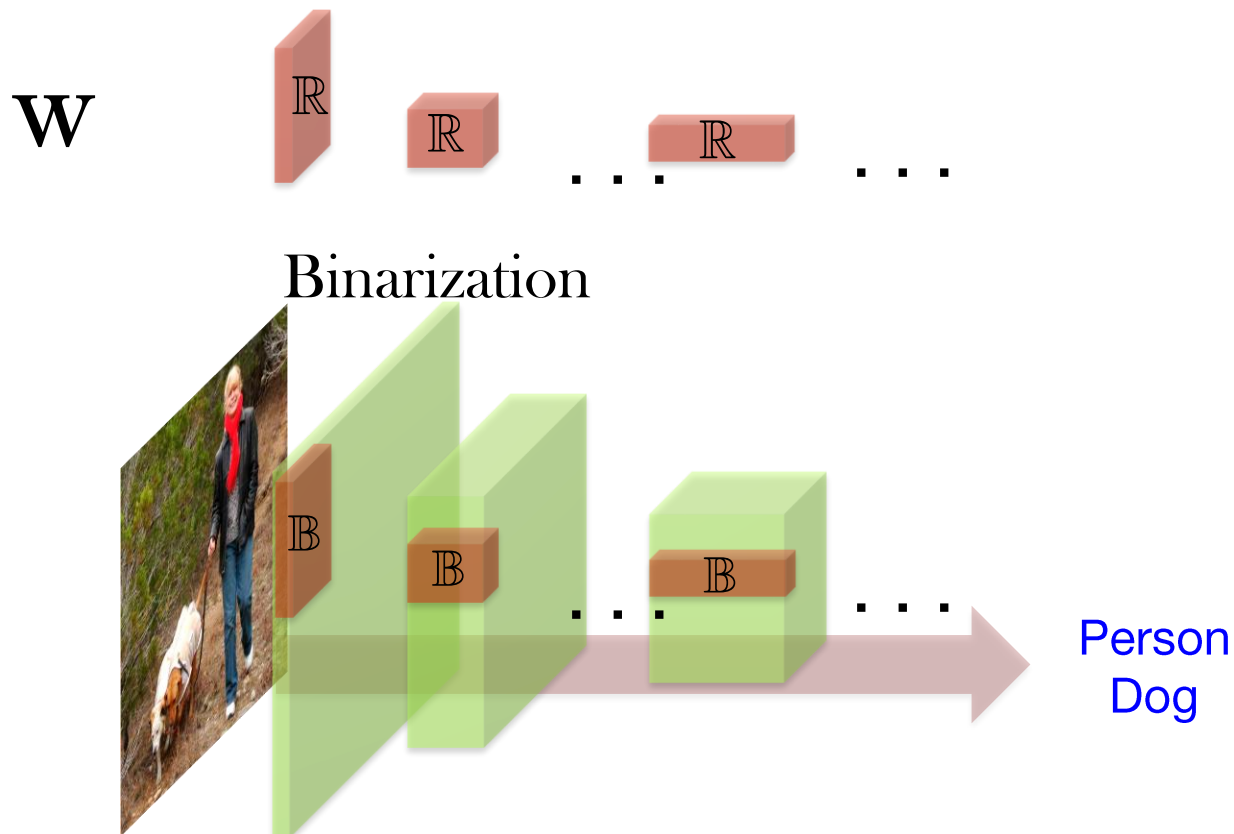ResNet-50 Top-1 (%) ILSVRC2012

$$\mathbf{W}$$
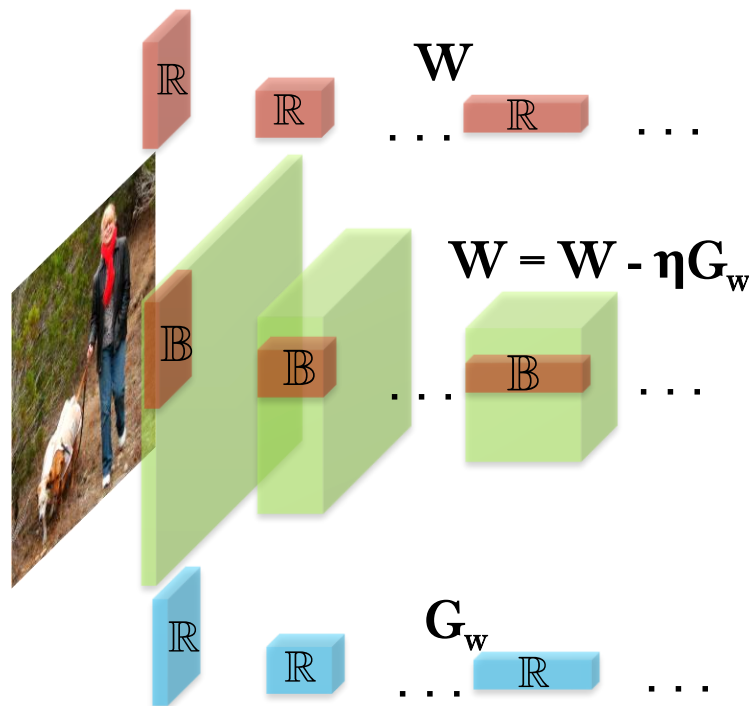
$$\mathbb{R} \qquad \mathbb{R} \qquad \ldots \mathbb{R} \ldots$$

Binarization

$$\mathbf{W^B}$$

$$\mathbb{B} \qquad \mathbb{B} \qquad \ldots \mathbb{B} \ldots$$

*Train for binary weights:*

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.     Load a random input image $\mathbf{X}$
4.     $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.     $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6.     Forward pass with $\alpha, \mathbf{W^B}$
7.     Compute loss function $\mathbf{C}$
8.     $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} = $ Backward pass with $\alpha, \mathbf{W^B}$
9.     Update $\mathbf{W}$ $(\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}})$



$$\mathbf{W}$$

$$\mathbf{W} = \mathbf{W} - \eta\mathbf{G_w}$$

$$\mathbf{G_w}$$

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.     Load a random input image $\mathbf{X}$
4.     $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.     $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6.     Forward pass with $\alpha, \mathbf{W^B}$
7.     Compute loss function $\mathbf{C}$
8.     $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} = $ Backward pass with $\alpha, \mathbf{W^B}$
9.     Update $\mathbf{W}$ ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)

- 15x Smaller
- 10x Faster
- 200% power efficiency

XNOR.AI

Raspber r y Pi Zer o

$5

Xnor.ai IP

R * R ≈ [ B * B ] ⊙ β ⊙ α

sign(x)   sign(**W**)
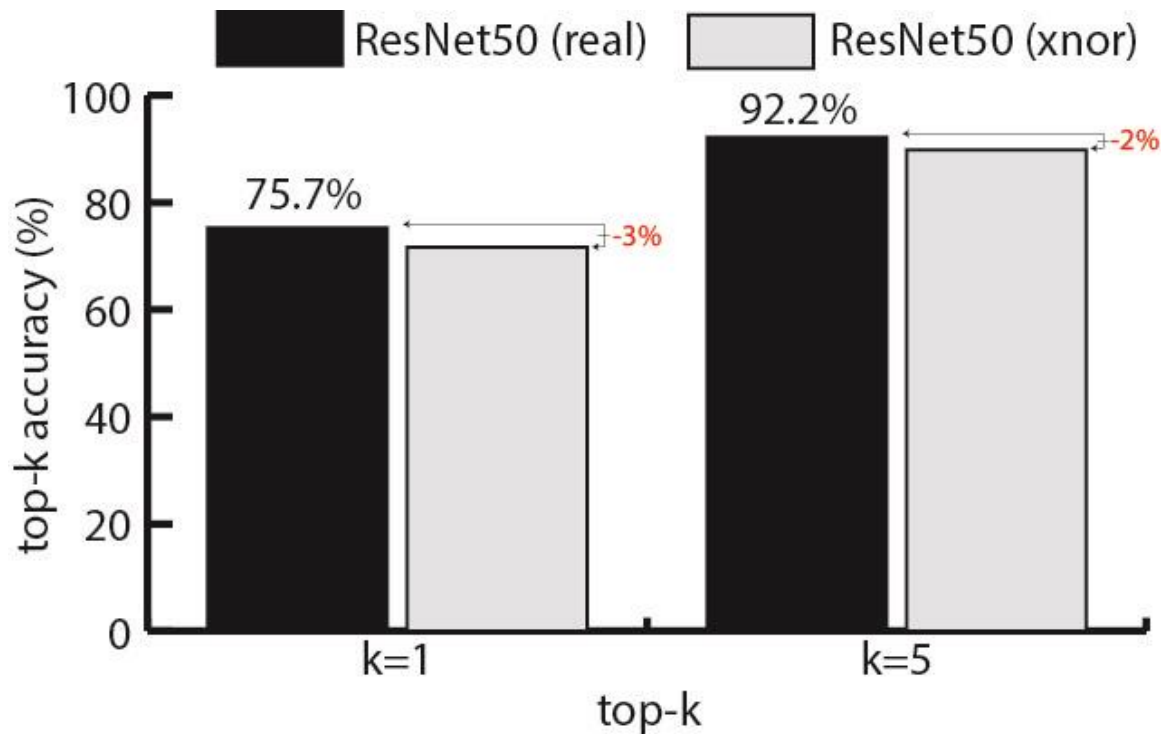
1. Randomly initialize W
2. For *iter* = 1 to *N*
3.    Load a random input image X
4.    $W^B$ = sign(W)
5.    $\partial = \frac{kWk_{-1}}{n}$
6.    Forward pass with $\partial$, $W^B$
7.    Compute loss function C
8.    $\frac{\partial C}{\partial W}$ = Backward pass with $\partial$, $W^B$
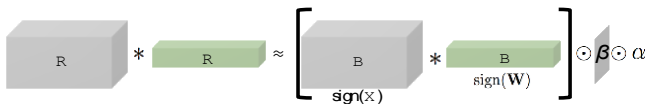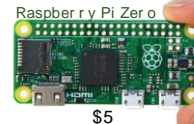9.    Update W (W = W − $\frac{\partial C}{\partial W}$)

BNorm | Activ | Conv W | Pool

Machine Learning

Code Optimization    Computer Architecture

How to integrate with XNOR.AI?

CUSTOMER INPUT

Training data + network infrastructure

XNOR-NET

Retrain deep learning model in binary format

**XNOR.AI deep learning** Deep learning platform for model creation and optimization

XNOR CORE

Binarized model input

Reference algorithm

**XNOR.AI CORE**

A proprietary code base that automatically translates binarized deep learning models into highly optimized programs that run on the leading hardware platforms

HARDWARE PLATFORMS

Intel          ARM          NVIDIA

CUSTOM INTEGRATIONS

**Product integrations for devices** Custom work done as part of our services agreement
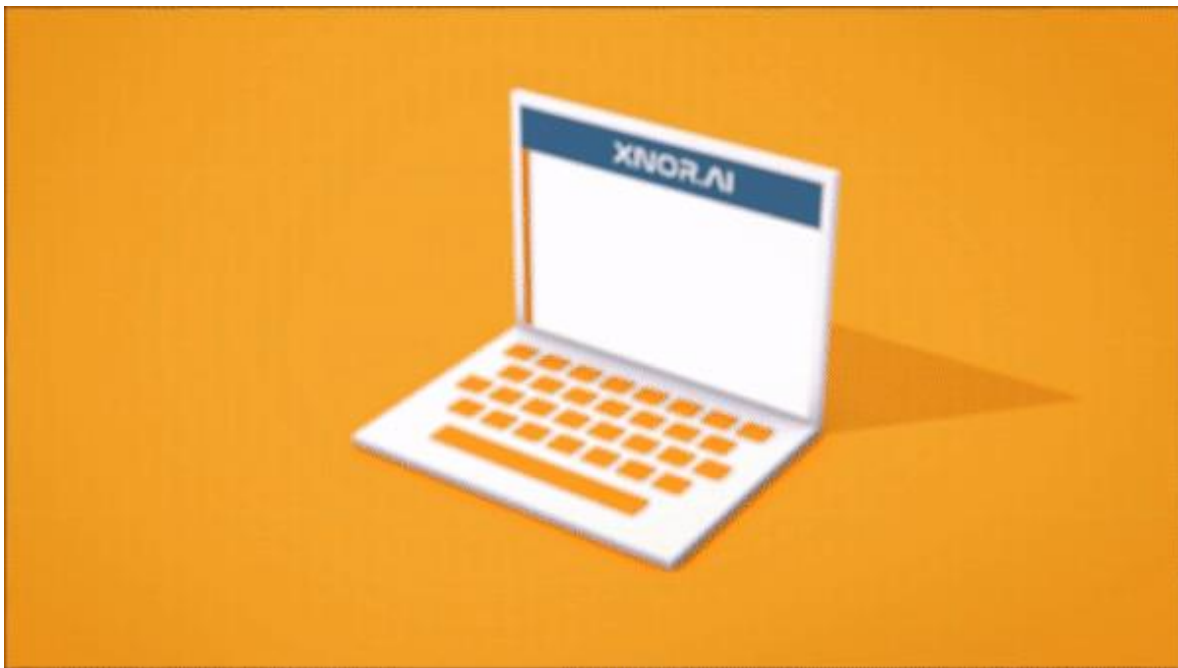
# Developers:

AI Everywhere
For Everyone

# Developers Platform



**MODEL SELECTION INTERFACE**

| TASK | DEVICE TYPE | MEMORY CONSTRAINT | LATENCY CONSTRAINT | POWER CONSTRAINT |
|---|---|---|---|---|
| ○ Scene Recognition | ◉ iPhone | | | |
| ○ Image Tagging | ☐ iPhone 6, 6s 6 Plus, 6s Plus | 1 MB | 100 FPS | 0.001 JPI |
| ◉ Object Detection | ☑ iPhone 7, 7 Plus | 10 MB | | |
| ☑ Person | ☑ iPhone 8, 8 Plus | | 10 FPS | 0.01 JPI |
| ☐ Tree | ☑ iPhone X | | | |
| ☐ Building | ○ Android Phones | 100 MB | | |
| ☐ Car | ○ Smartwatch | | 1 FPS | 0.1 JPI |
| ☐ Traffic Light | ○ ARM-based Platforms | 1 GB | | |
| ⋮ | ○ Intel-based Platforms | | 0.1 FPS | 1 JPI |
| ○ Image Enhancement | ○ NVIDIA Platforms | | | |
| ○ Object Segmentation | ○ Drones | | | |
| ○ Action Recognition | ○ X86 Families | | | |
| ○ Speech Recognition | ⋮ | 10 GB | 0.01 FPS | 10 JPI |
| ○ Text Recognition | | | | |
| ⋮ | | | | |

# XNOR.AI technology powers multiple domains

✓ Aerospace & Surveillance

✓ Driver Assisted Systems

✓ Retail

✓ Consumer Mobile

# AI Everywhere

Founded:
2017 by Professor Ali Farhadi
and Dr. Mohammad Rastegari

Intellectual Property:
Highly strategic patent portfolio covering
efficient AI at the edge

Press:
New York Times
Tech Crunch

Board members:
Ali Farhadi (CEO), Oren Etzioni (CEO of
Allen Institute for AI), Matt McIlwain
(Madrona Venture Group)

XNOR Innovations in AI:
XNOR-Net, Yolo, Yolo9000, LCNN, Neural Speed
Reading, understanding actions (imSitu), question
answering (BiDAF)

Awards:
Best paper at CVPR 2017
CVPR 2016 People's Choice Award

# Thank you !!!

## Learn more
## Visit our table #809
## www.xnor.ai

Mohammad Rastegari  |  Chief Technology Officer  |  mohammad@xnor.ai