# AI Acceleration
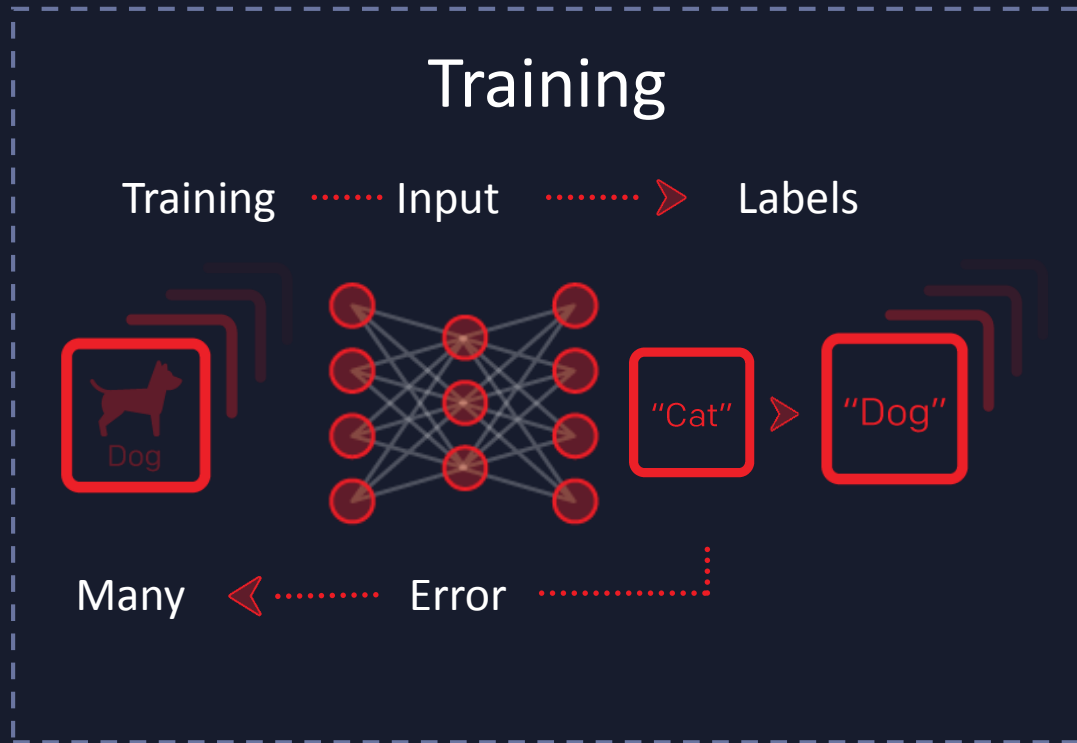
Shaun Purvis

# Training vs. Inference



**Training**

Training ⋯⋯ Input ⋯⋯▷ Labels

"Cat" ▷ "Dog"

Many ◁⋯⋯ Error
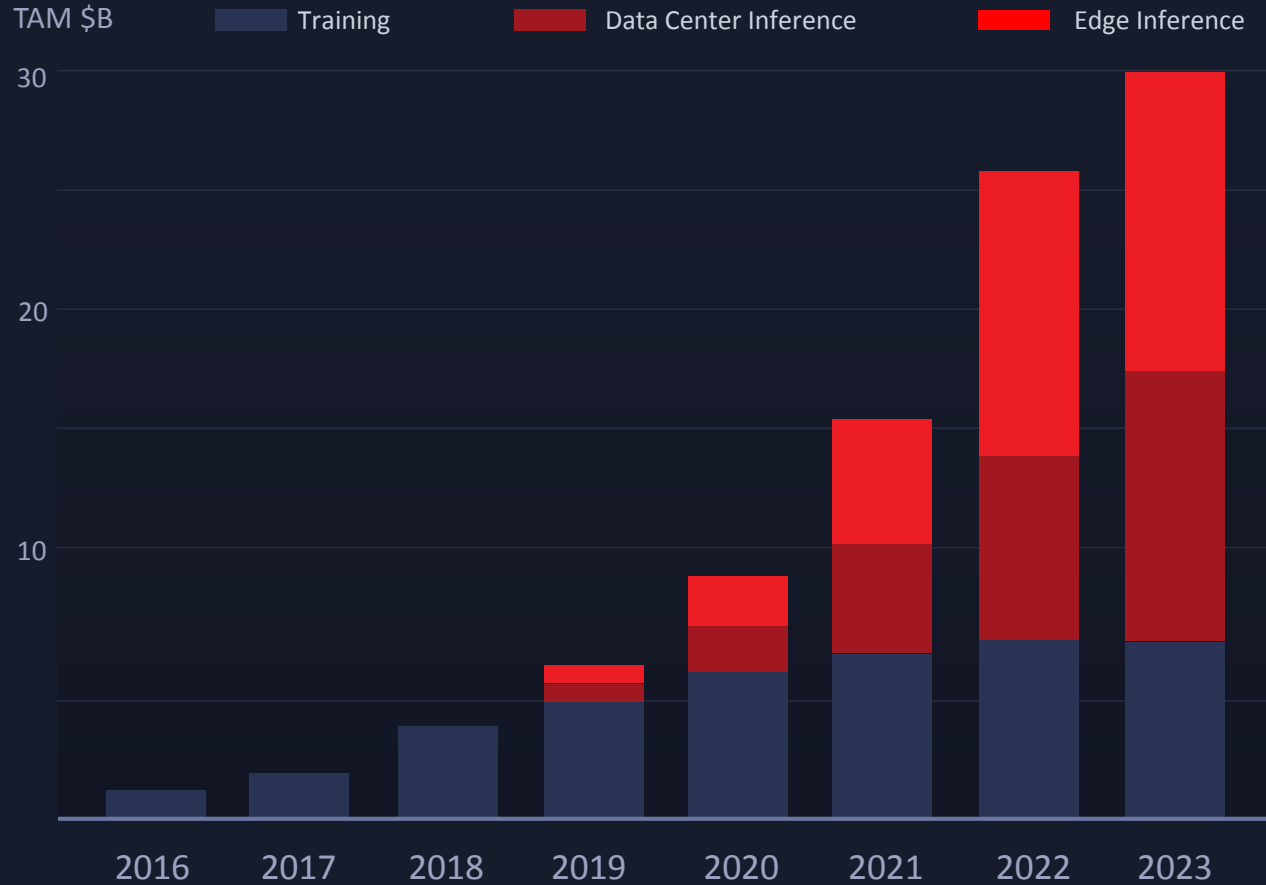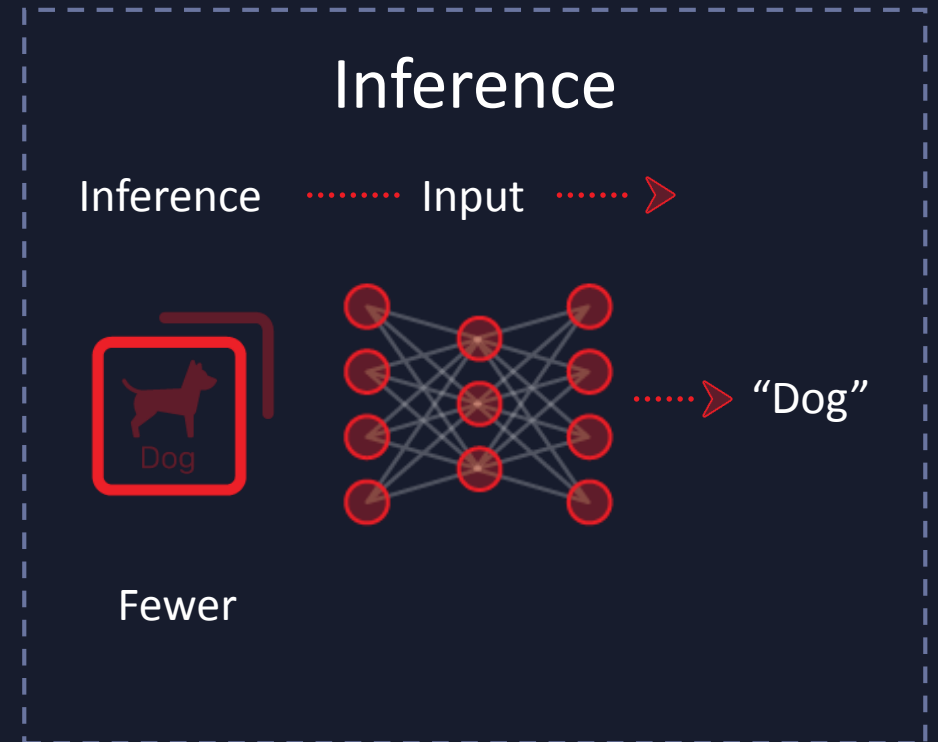
Migrate trained model to inference hardware

**Inference**

Inference ⋯⋯ Input ⋯⋯▷

"Dog"

# Inference Projected Growth

**TAM $B** ▪ Training ▪ Data Center Inference ▪ Edge Inference

Barclays Research, Company Reports May 2018

## Inference

Inference ┄┄ Input ┄┄ ▷

Fewer ┄┄ ▷ "Dog"

# Inference Space

# Inference Challenges

The rate of AI innovation

Performance at low latency

Low power consumption

Whole app acceleration

## Inference

Inference ·········· Input ·········· ▷

Dog

·········· ▷ "Dog"

Fewer

# The Rate of AI Model Innovation

APPLICATIONS

Classification

Object Detection

Segmentation

Speech Recognition

Recommendation Engine

Anomaly Detection

CNN

RNN, LSTM

MLP

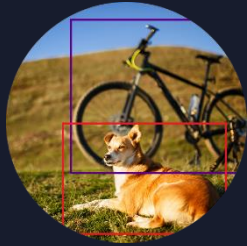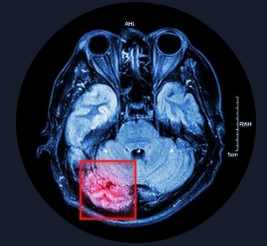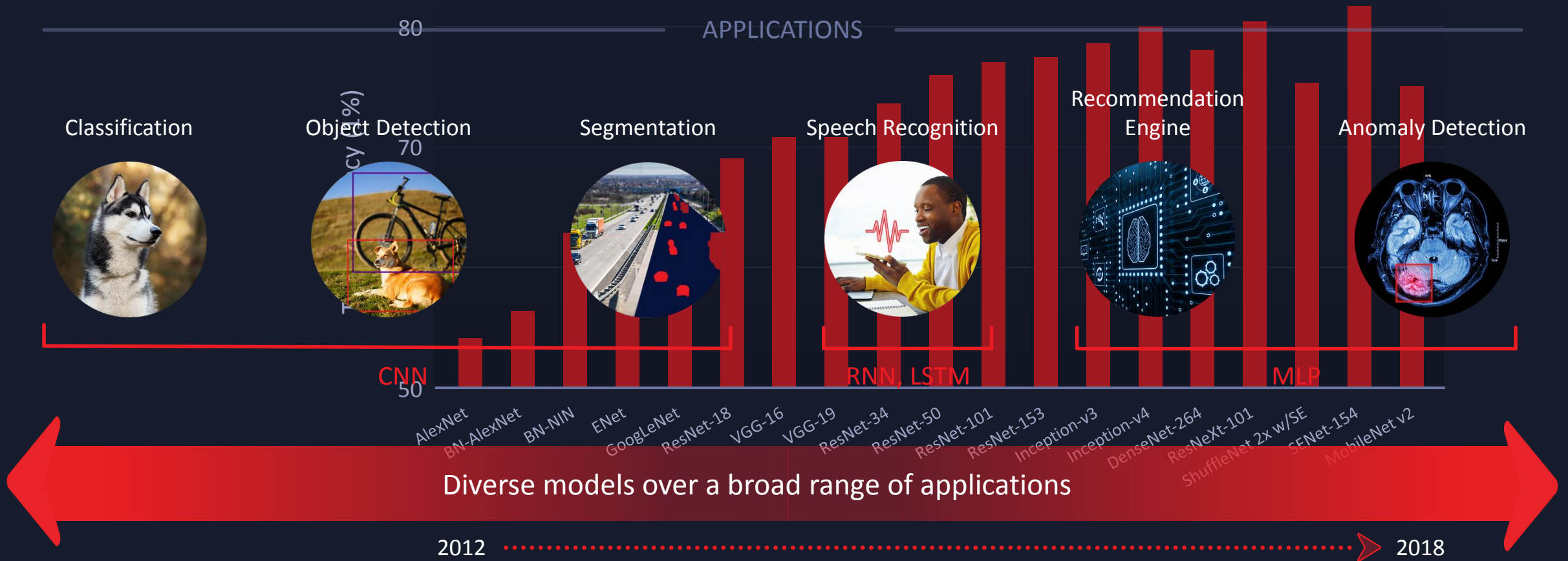Diverse models over a broad range of applications

# The Rate of AI Model Innovation: Classification

APPLICATIONS

80

Classification    Object Detection    Segmentation    Speech Recognition    Recommendation Engine    Anomaly Detection

70

CNN    RNN, LSTM    MLP

50

AlexNet, BN-AlexNet, BN-NIN, ENet, GoogLeNet, ResNet-18, VGG-16, VGG-19, ResNet-34, ResNet-50, ResNet-101, ResNet-153, Inception-v3, Inception-v4, DenseNet-264, ResNeXt-101, ShuffleNet 2x w/SE, SENet-154, MobileNet v2

Diverse models over a broad range of applications

2012 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . 2018

# Network Complexity is Growing

AlexNet

GoogLeNet

DenseNet

| | FP32 | INT32 | | FP/INT16 | | INT8 | INT6 | INT4 | INT2 | INT1 |
|---|---|---|---|---|---|---|---|---|---|---|
| 2012 | 2013 | | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | |

## Inference is Moving to Lower Precision

### RELATIVE ENERGY COST

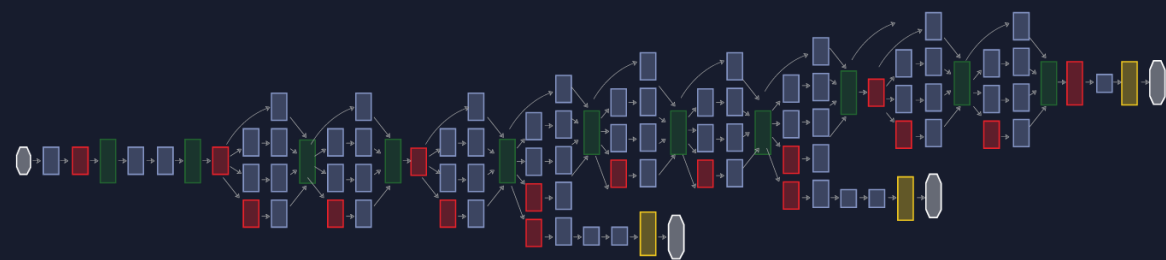| Operation: | Energy (pJ) | |
|---|---|---|
| 8b Add | 0.03 | |
| 16b Add | 0.05 | |
| 32b Add | 0.1 | |
| 16b FP Add | 0.4 | |
| 32b FP Add | 0.9 | |

# Rate of Innovation Outpaces Silicon Cycles

AlexNet

GoogLeNet

DenseNet

**Silicon lifecycle**

# Only Adaptable Hardware Addresses Inference Challenges

Custom data flow

Custom memory hierarchy

Custom precision

Domain Specific Architectures
(DSAs)
on Adaptable Platforms
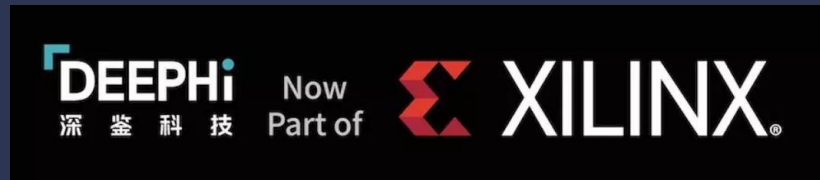
# DeePhi Joins Xilinx

Custom data flow

Custom memory hierarchy

Custom precision



Pruning          Quantization

Patented Compression Technology
- Reduces DL accelerator footprint
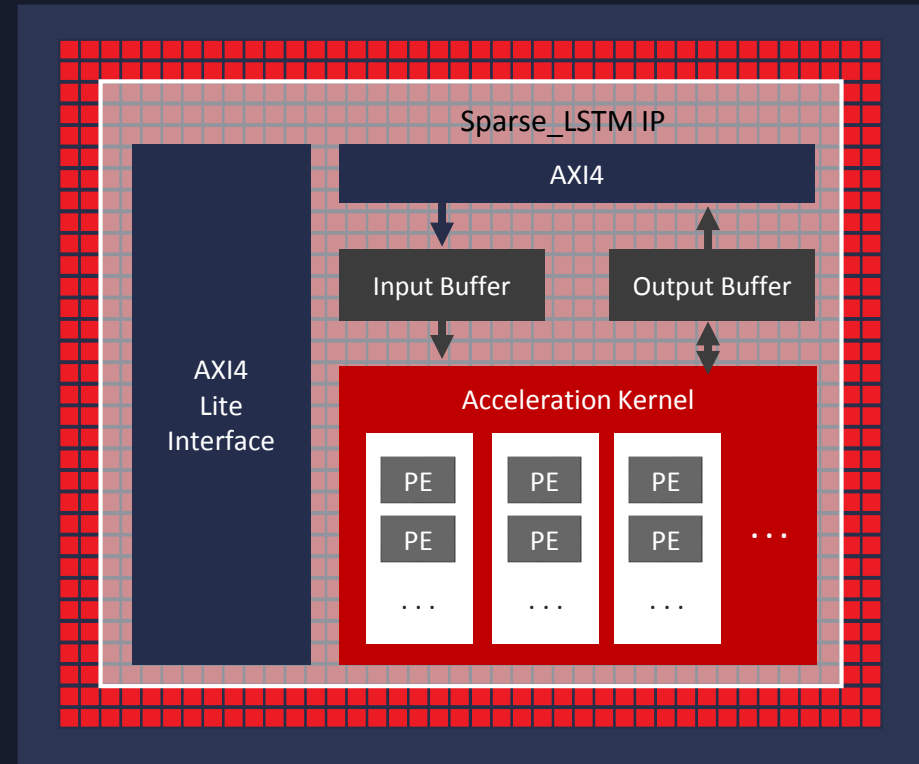- Increases performance per watt

# Example: DeePhi LSTM

**Custom data flow**

LSTM for speech recognition

**Custom memory hierarchy**

Sparse matrix implementation in memory
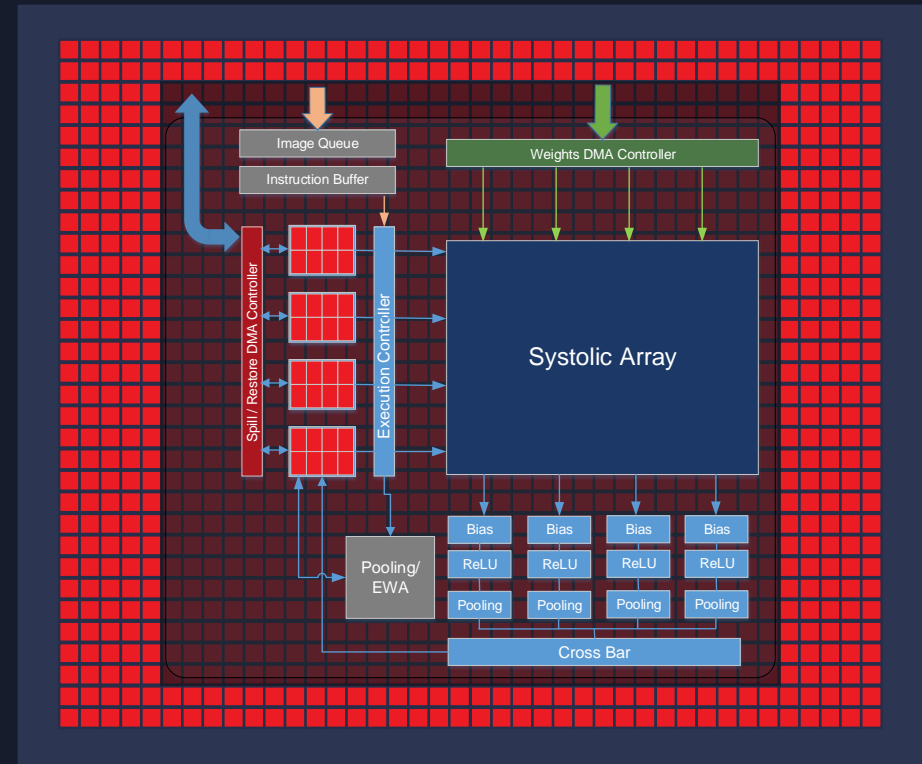
**Custom precision**

12 bit weights, 16 bit activations

Sparse_LSTM IP

AXI4

AXI4
Lite
Interface

Input Buffer

Output Buffer

Acceleration Kernel

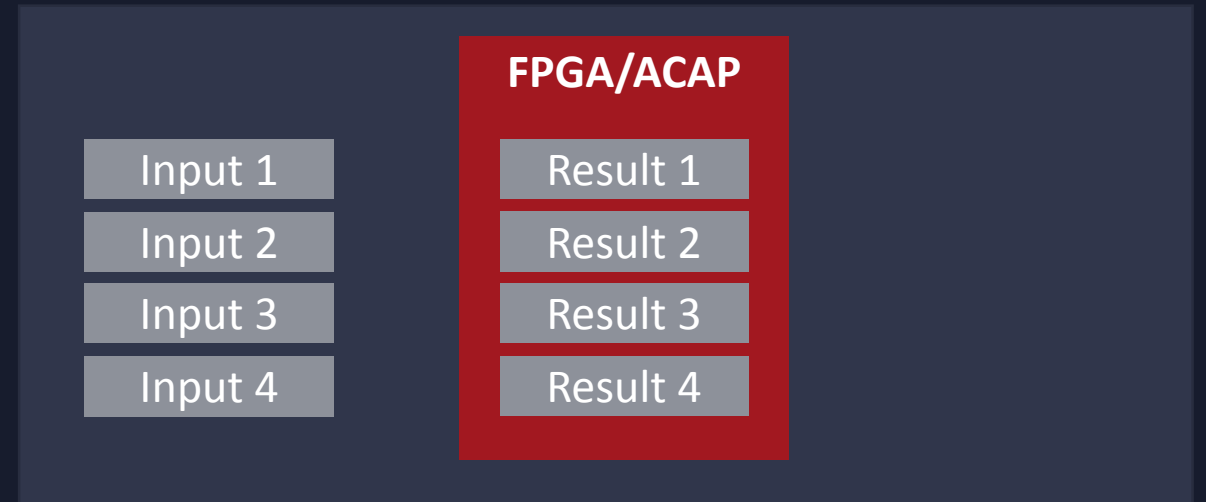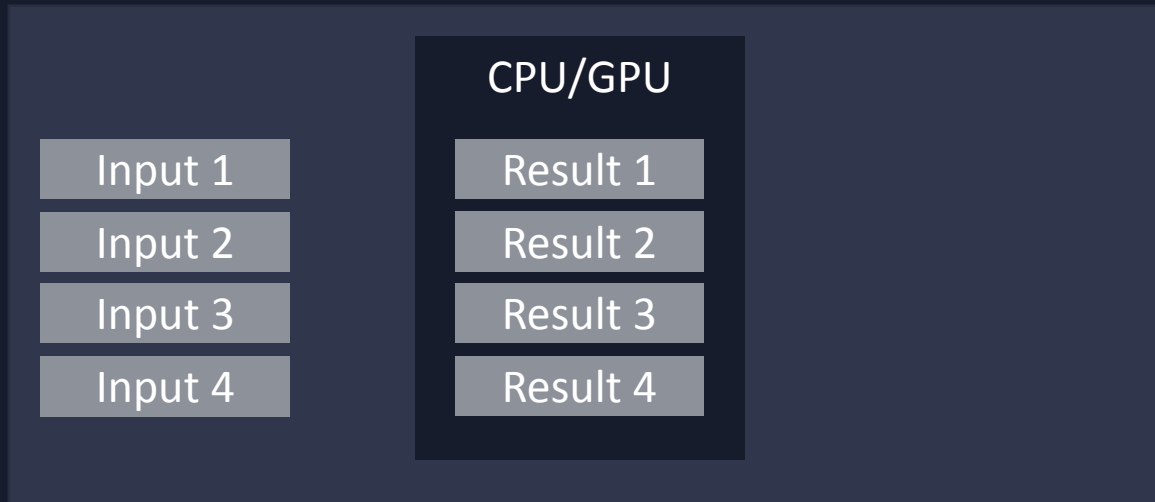| PE | PE | PE |
| PE | PE | PE |
| . . . | . . . | . . . |

. . .

# Example: xDNN

**Custom data flow**

Optimized for latest CNN

**Custom memory hierarchy**

Optimized on-chip memory

**Custom precision**
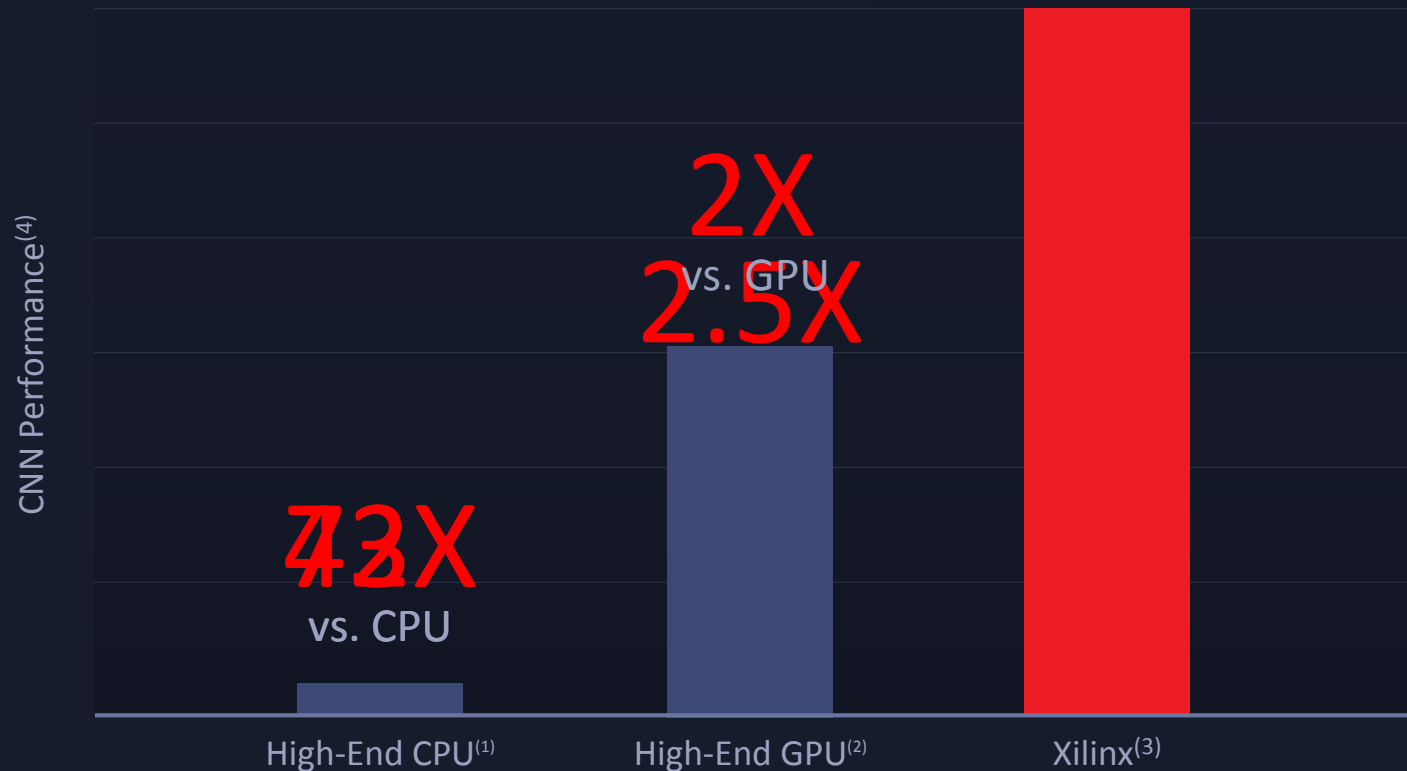
Int8

# Low Latency is Critical for Inference

50ms latency response

3ms latency response

**CPU/GPU**

| Input 1 | Result 1 |
| Input 2 | Result 2 |
| Input 3 | Result 3 |
| Input 4 | Result 4 |

**FPGA/ACAP**

| Input 1 | Result 1 |
| Input 2 | Result 2 |
| Input 3 | Result 3 |
| Input 4 | Result 4 |

High throughput OR low latency

High throughput AND low latency

# Low Latency: Xilinx's Unique Advantage

## Latency Insensitive Inference



**2X**
vs. GPU

**2.5X**

**73X**
**42X**
vs. CPU

CNN Performance[4]

High-End CPU[1]    High-End GPU[2]    Xilinx[3]

## AI Inference Acceleration

Leveraging AI Engines

Majority of Adaptable & Scalar Engines available for Whole App Acceleration

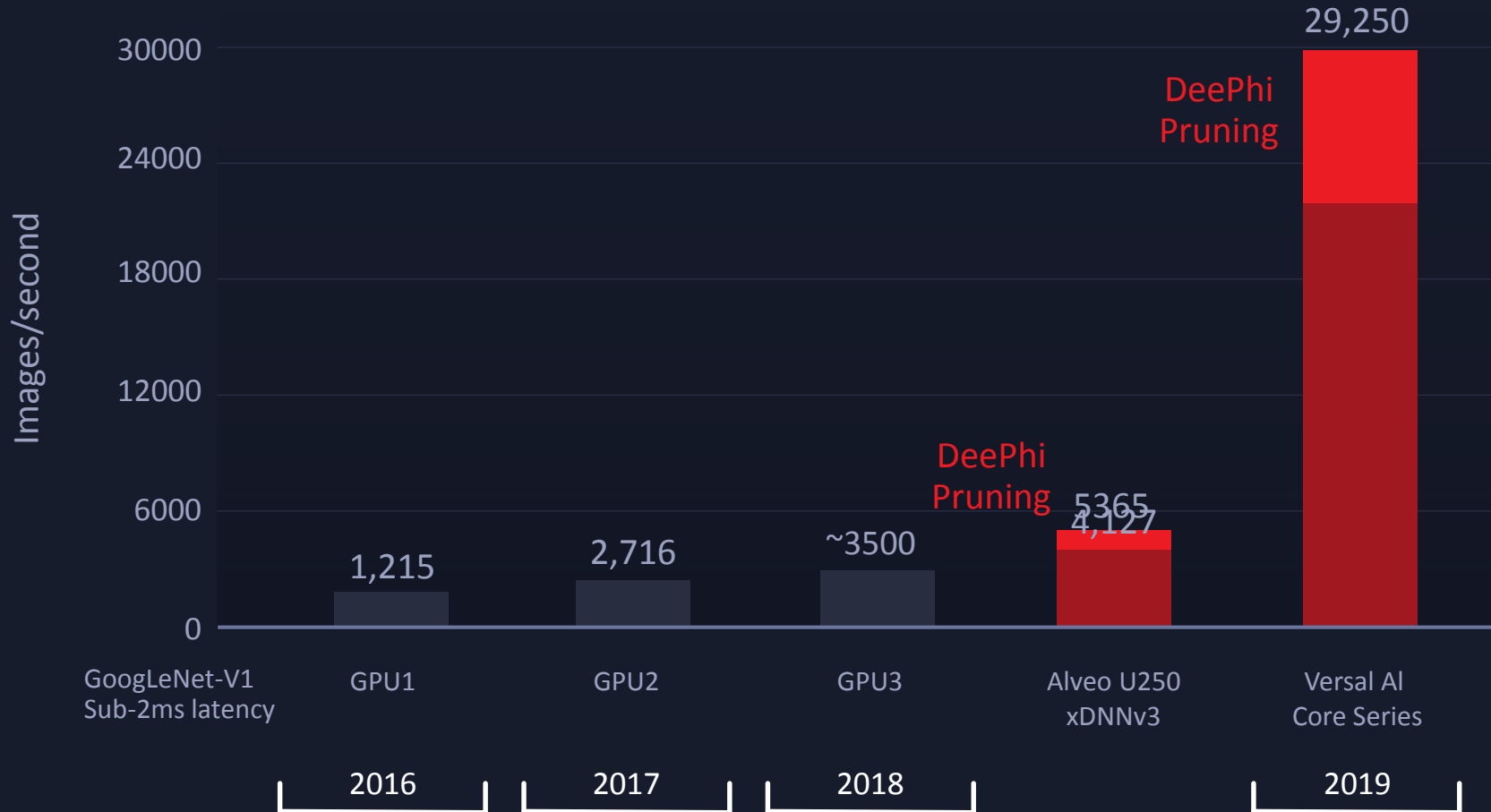(1)  Measured on EC2 Xeon Platinum 8124 Skylake, c5.18xlarge AWS instance, Intel Caffe: https://github.com/intel/caffe
(2)  V100 numbers taken from Nvidia Technical Overview, "Deep Learning Platform, Giant Leaps in Performance and Efficiency for AI Services"
(3)  Versal Core Series
(4)  GoogLeNet V1 throughput (Img/sec)

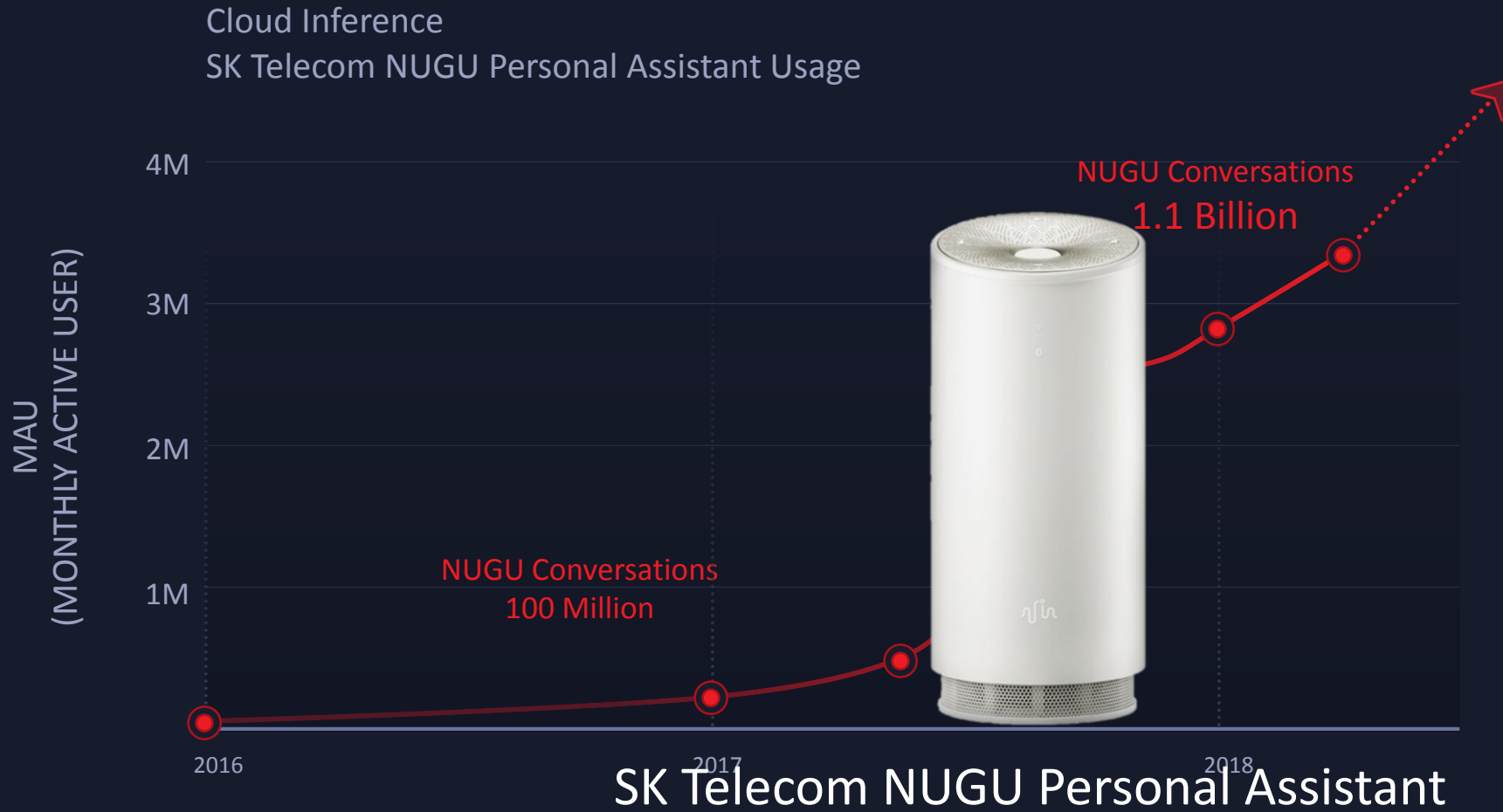# Low-Latency CNN Inference Performance



30000

24000

DeePhi
Pruning

29,250

18000

Images/second

12000

DeePhi
Pruning

6000

1,215

2,716

~3500

5365
4,127

0

GoogLeNet-V1
Sub-2ms latency

GPU1

GPU2

GPU3

Alveo U250
xDNNv3

Versal AI
Core Series

2016
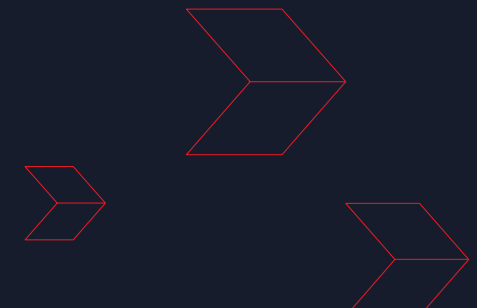
2017

2018

2019

DeePhi Pruning
Technology

# 1.3x-8x

Performance improvement
based on the
network

# Power Is Critical for Inference Applications

Cloud Inference
SK Telecom NUGU Personal Assistant Usage

NUGU Conversations
1.1 Billion

NUGU Conversations
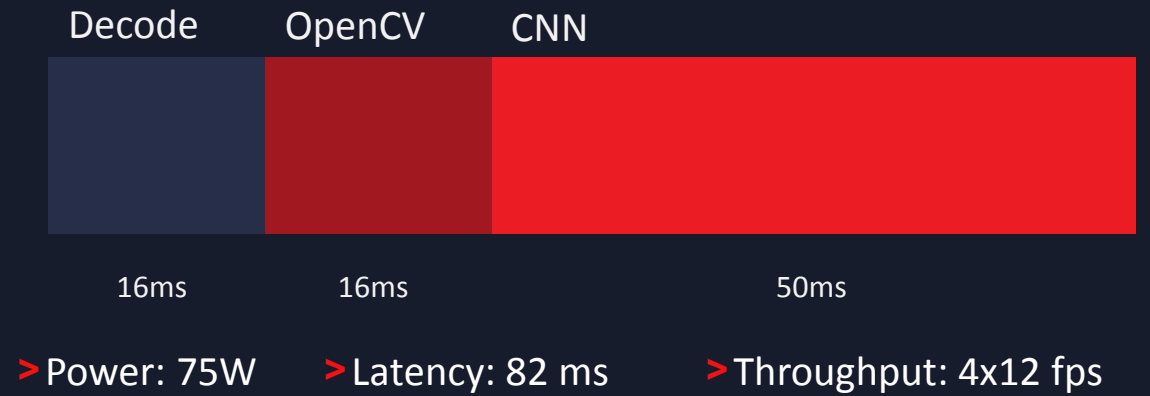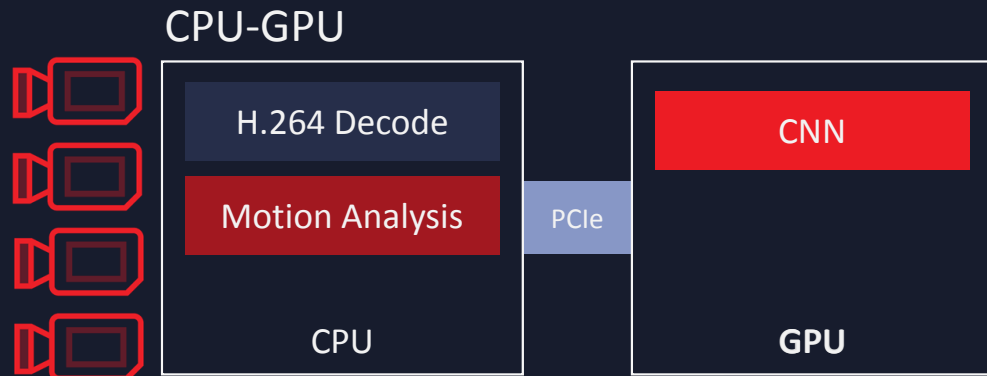100 Million

MAU
(MONTHLY ACTIVE USER)

4M

3M

2M

1M

2016

2017

2018

SK Telecom NUGU Personal Assistant

# 16x
## Perf/watt
## vs. GPU

# Whole Application Acceleration:
# Smart City / Security

## CPU-GPU

| CPU | | GPU |
| H.264 Decode | PCIe | CNN |
| Motion Analysis | | |

Decode · OpenCV · CNN

16ms · 16ms · 50ms

> Power: 75W   > Latency: 82 ms   > Throughput: 4x12 fps

## CPU-Xilinx FPGA

| CPU | | FPGA |
| H.264 Decode | PCIe | CNN |
| | | Motion Analysis |

16ms · 0.9ms · 9.2ms

> Power: 50W   > Latency: 26.1 ms   > Throughput: 4x38 fps

# Whole Application Acceleration: Online Video Streaming



**1 Aup2603**

48 ZU7EV

| Decode H.264 | ML / Analytics | Process | Encode HEVC |
|---|---|---|---|

Video transcoding + AI analytics

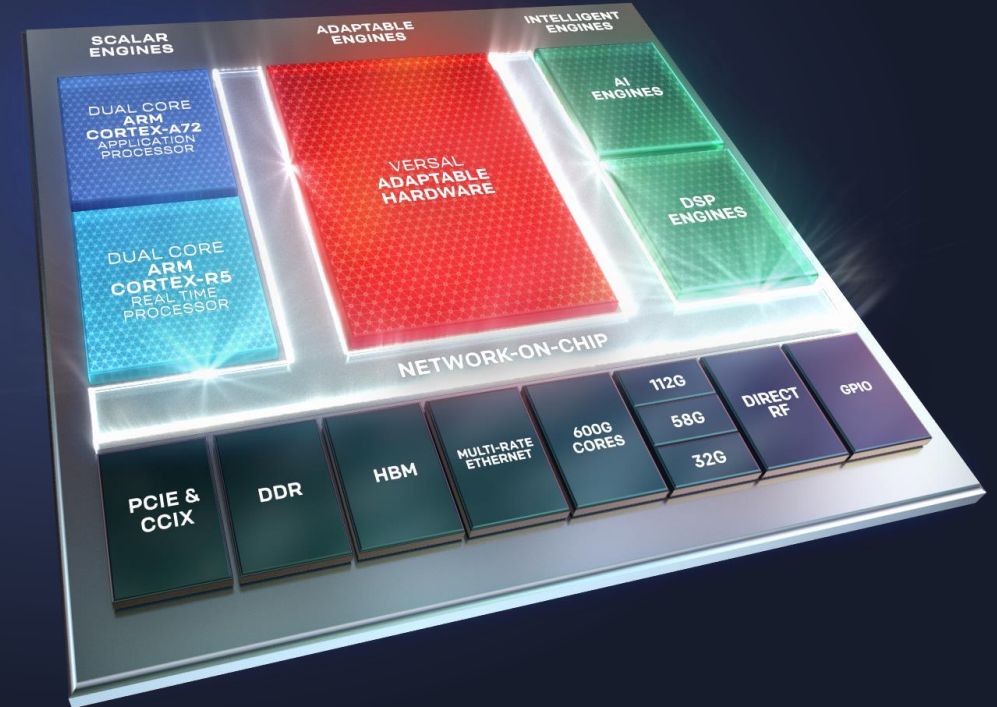Energy ────── 10x

Performance ────── 3.3x

**30 E5 Servers**

# Only Xilinx Adaptable Devices Can:

Match the speed of AI innovation

Give the best performance at low latency

Give the best power results

Accelerate the whole application

# Xilinx

> Building
> the Adaptable,
> Intelligent World