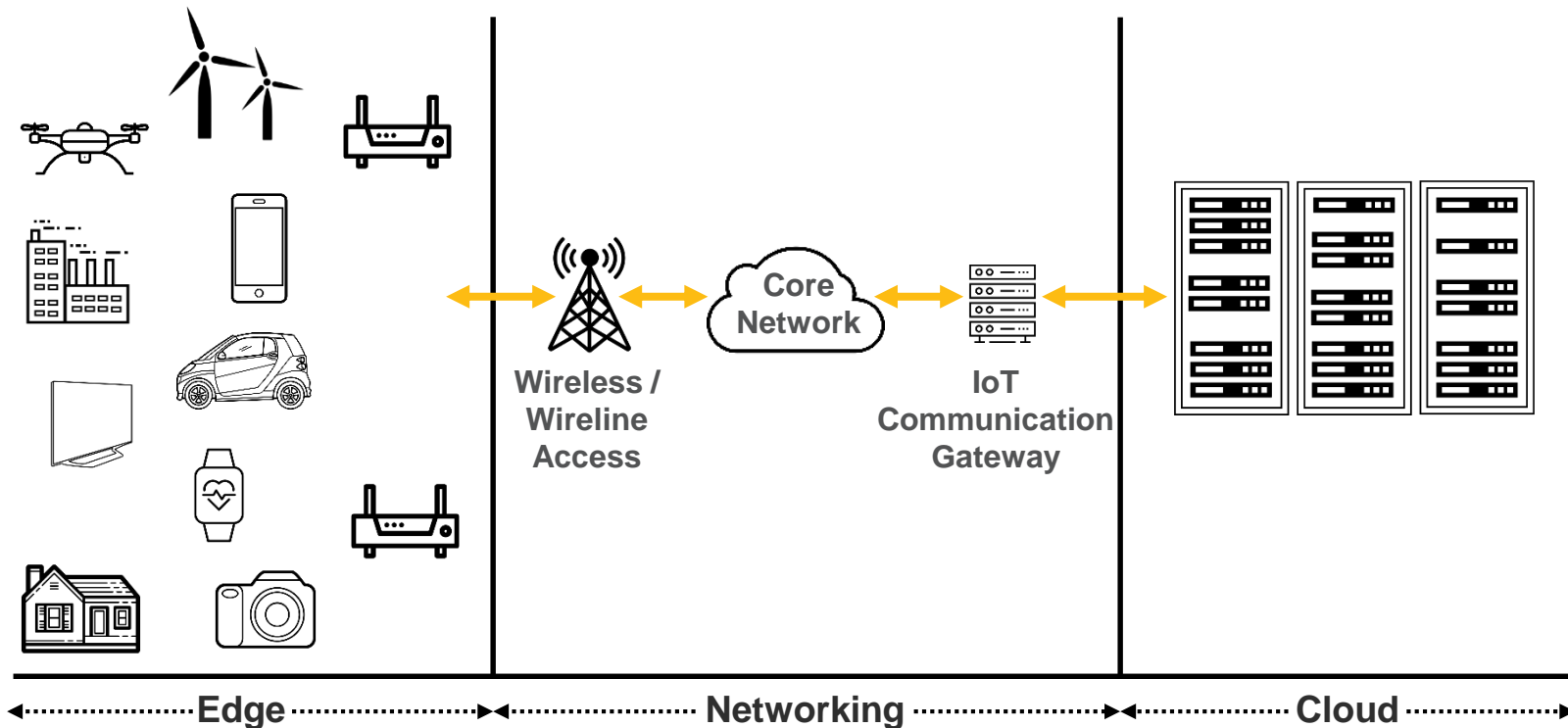# Architecting Always-On, Context-Aware, On-Device AI Using Flexible Low-power FPGAs

Deepak Boppana – Senior Director Product & Segment Marketing
Gordon Hands – Director Solutions Marketing

**LATTICE**
SEMICONDUCTOR

# Rapidly Emerging Edge Computing Trend
## Driven by Latency, Privacy, and Bandwidth Limitations



**Wireless / Wireline Access**

**Core Network**

**IoT Communication Gateway**

◄┅┅┅┅ **Edge** ┅┅┅┅►◄┅┅┅┅ **Networking** ┅┅┅┅►◄┅┅┅┅ **Cloud** ┅┅┅┅►

Unit growth for edge devices with AI will explode increasing over 110% CAGR over the next five years – *Semico Research*

# Always-on, On-device AI Applications

Human Presence Detection Example

**Smart Home Appliance**

LCD turns on when needed

**Consumer Electronics**

TV turns off when no one is present

**Smart DoorBell**

Rings automatically when needed

**Vending Machine**

LCD turns on when needed

**Security Camera**

Alerts when intruder present, not a cat

**Smart Doors**

Opens when person is present

# Always-on, On-device AI Applications
## Other Examples

**Smart speakers**
Key phrase detection

**Retail store cameras**
Face tracking

**Selfie drones**
Face tracking

**Toll gate camera**
Vehicle classification

**Machine vision**
Object counting

**After market automotive cameras**
Speed sign detection

SPEED LIMIT 25

LATTICE
SEMICONDUCTOR
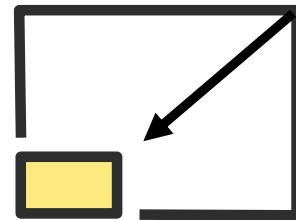
# Always-on, On-device AI Requirements
Unmet Need for Ultra-Low Power, Scalable, and Flexible Inferencing

**Few mWs of Power Consumption**
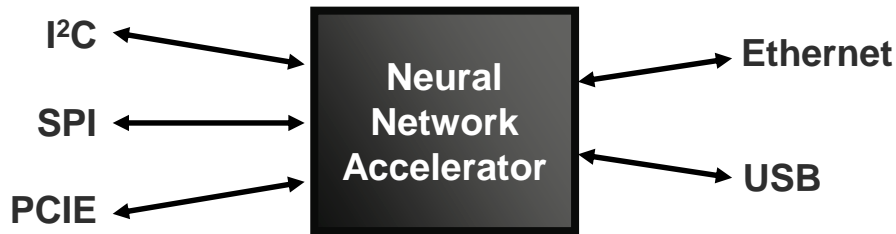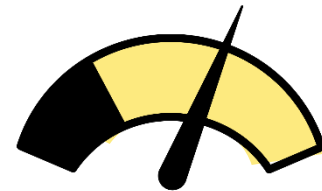
**Few $s of BOM Cost Adder**

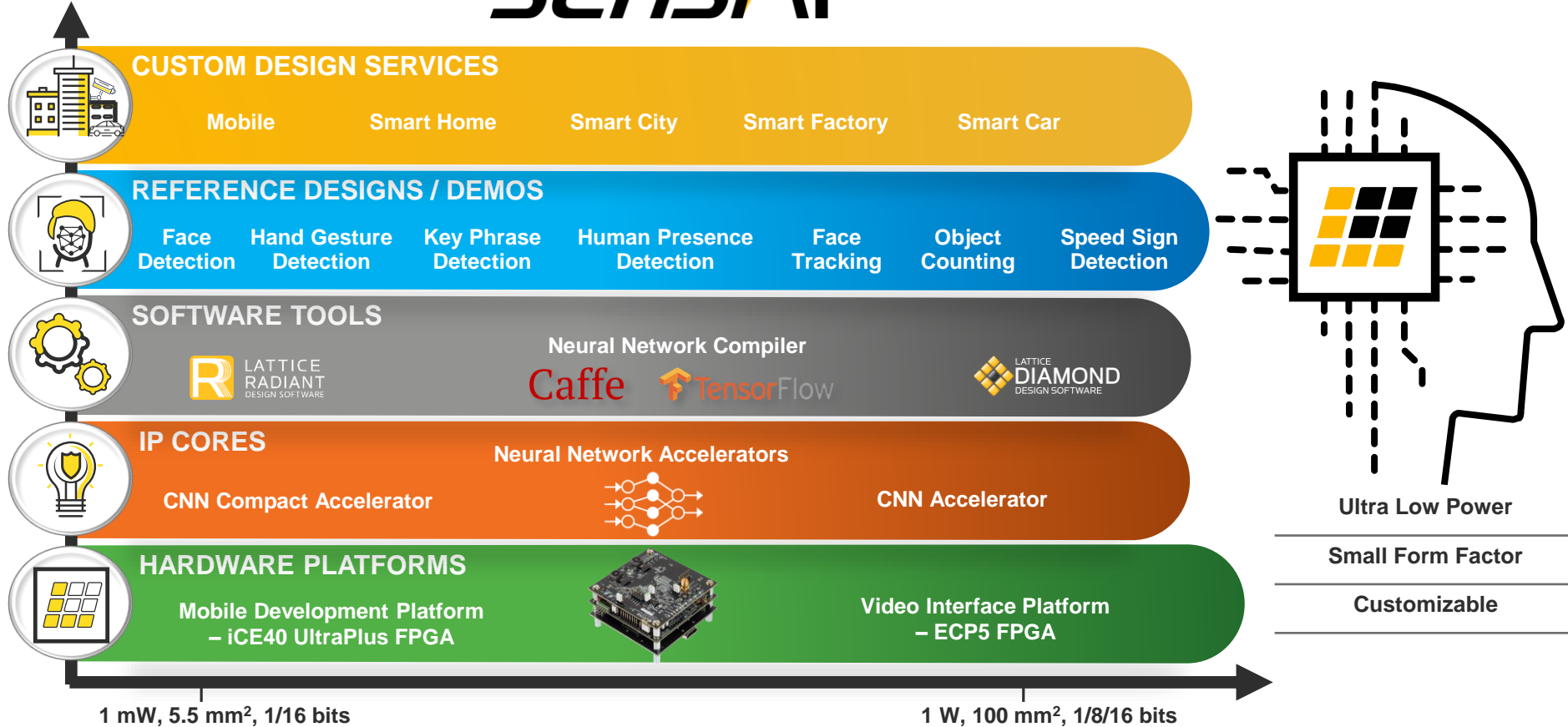**Few mm$^2$ of Board Area**

I$^2$C

SPI

PCIE

**Neural Network Accelerator**

Ethernet

USB

**Flexible Legacy Interface Support**

**Customized Performance/Accuracy**

# LATTICE sensAI

**CUSTOM DESIGN SERVICES**

Mobile          Smart Home          Smart City          Smart Factory          Smart Car

**REFERENCE DESIGNS / DEMOS**

Face Detection | Hand Gesture Detection | Key Phrase Detection | Human Presence Detection | Face Tracking | Object Counting | Speed Sign Detection

**SOFTWARE TOOLS**

LATTICE RADIANT DESIGN SOFTWARE

Neural Network Compiler

Caffe          TensorFlow

LATTICE DIAMOND DESIGN SOFTWARE

**IP CORES**

Neural Network Accelerators

CNN Compact Accelerator          CNN Accelerator

**HARDWARE PLATFORMS**

Mobile Development Platform – iCE40 UltraPlus FPGA

Video Interface Platform – ECP5 FPGA

1 mW, 5.5 mm$^2$, 1/16 bits          1 W, 100 mm$^2$, 1/8/16 bits

Ultra Low Power

Small Form Factor

Customizable

LATTICE SEMICONDUCTOR

# Flexible and Scalable Inferencing at the Edge

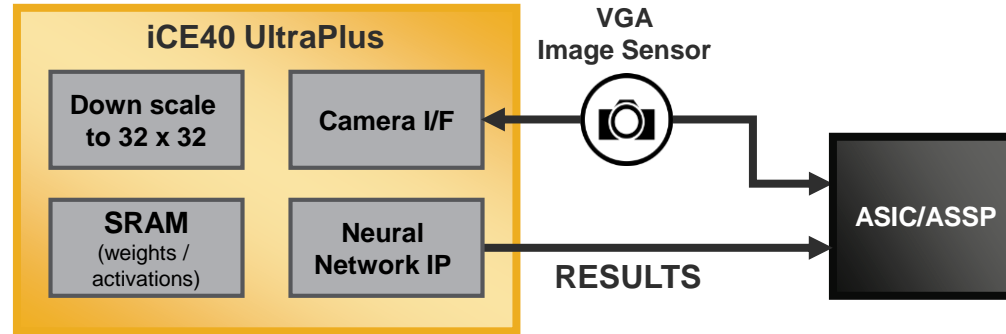From under 1 mW to 1 W with Lattice sensAI

# Stand-alone, Integrated FPGA Solution



- Always-on, integrated solutions on ECP5 or iCE40 UltraPlus FPGA
- Low latency and secure implementation
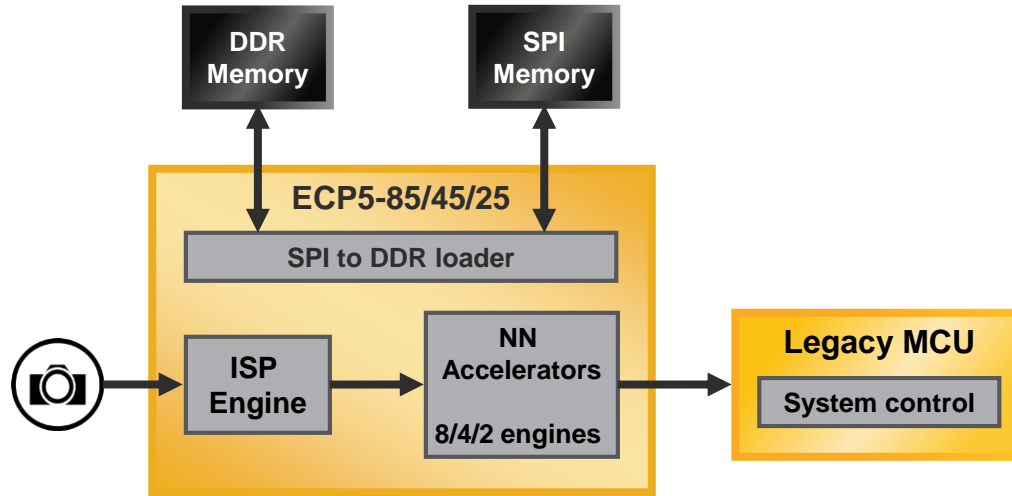- Small form factor packages from 5.5 mm$^2$ to 100 mm$^2$
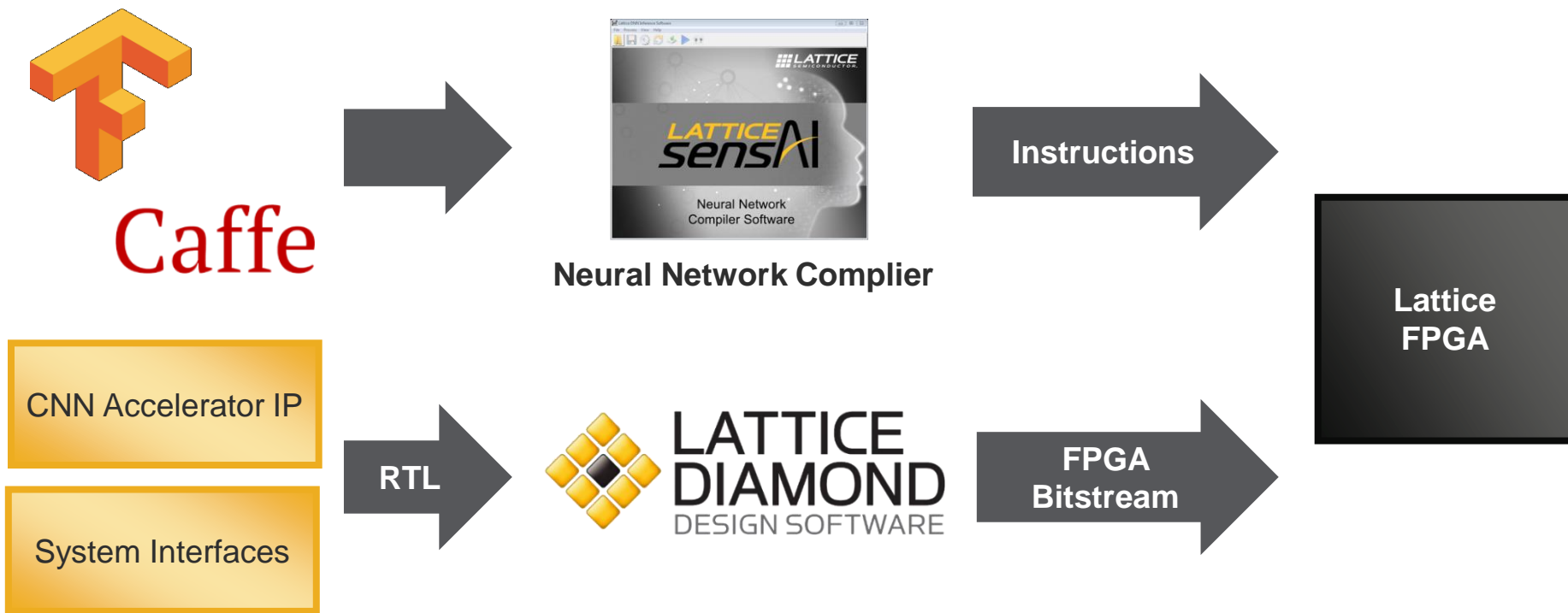
# FPGA as Activity Gate to ASIC/ASSP



- iCE40 UltraPlus FPGA for always-on detection of key-phrases or objects
- Wakes-up a high performance ASIC/ASSP for further analytics only when required
- Reduces overall system power consumption

# FPGA as a Co-Processor to MCU



- Scalable performance/power with ECP5 based neural network acceleration
- ECP5 based IO flexibility to seamlessly interface to on-board legacy devices including sensors
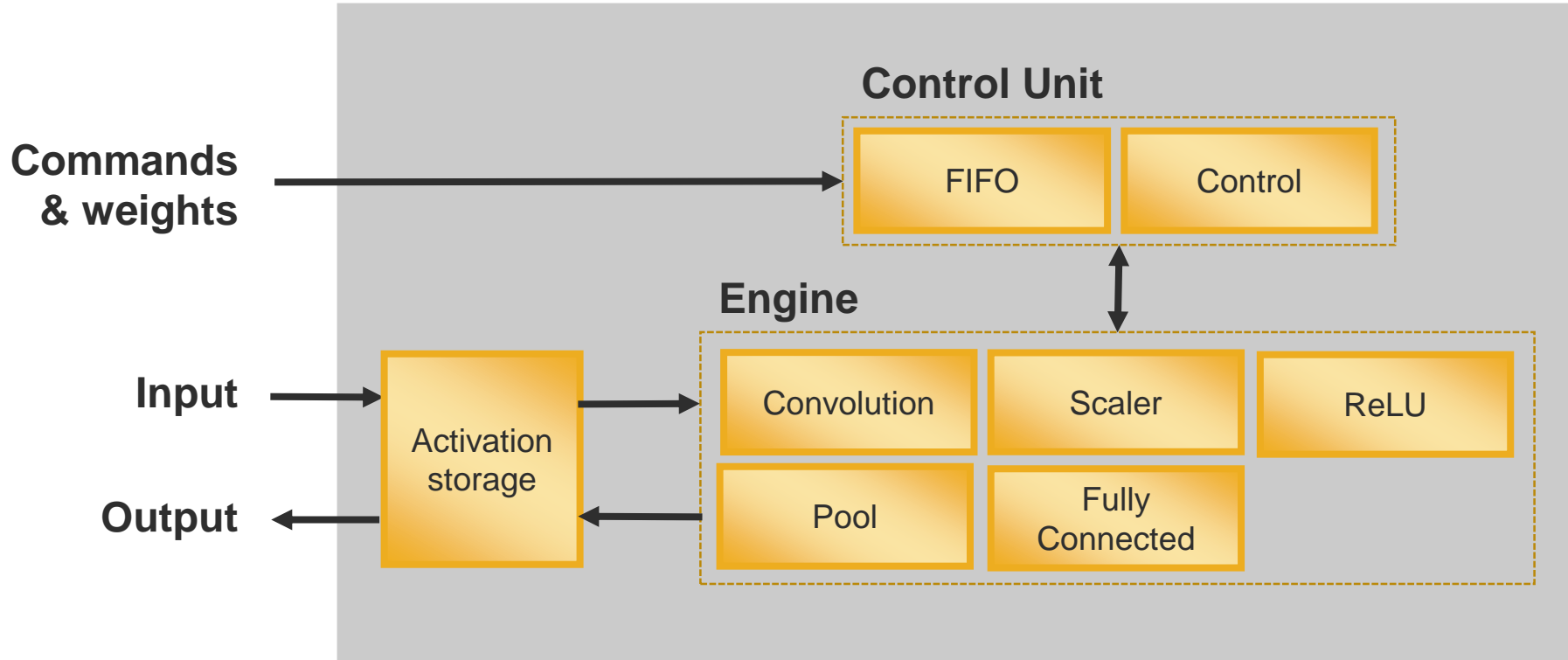- Low-end MCU for flexible system control

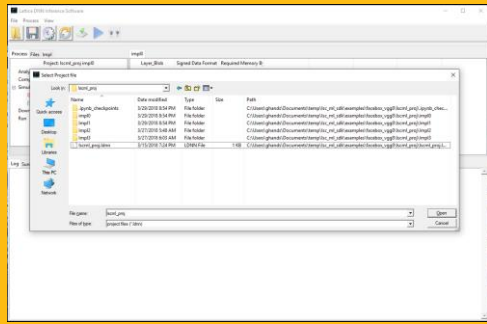# Delivering Edge CNN Acceleration in Lattice FPGA



**Neural Network Complier**

**Instructions**

**Lattice FPGA**

CNN Accelerator IP

System Interfaces

**RTL**

**FPGA Bitstream**

# CNN Accelerator IP Architecture

# CNN Compact Accelerator IP Architecture



**Control Unit**

FIFO | Control

**Commands & weights**

**Engine**

Convolution | Scaler | ReLU

Pool | Fully Connected

**Input**

Activation storage

**Output**

# Translating Trained Neural Network Into Lattice CNN Accelerator Instructions

## 1. Load



## 2. Review



## 3. Analyze



## 4. Compile



## 5. Simulate

# On-device AI – Complex Optimization

| Attributes / Design Factors | Device | | Network | | |
|---|---|---|---|---|---|
| | # of Engines | Local Memory | Input Size | Number of Multipliers | Bit Widths |
| Power (W) | 🟥 | 🟩 | | | |
| Device Size | 🟥 | 🟥 | | | |
| Performance (fps) | 🟩 | 🟩 | 🟥 | 🟥 | 🟥 |
| Accuracy (%) | | | 🟩 | 🟩 | 🟩 |
| Small Object (% fov) | | | 🟩 | | |

**Correlation Between Design Factors and Product Attributes**

LATTICE
SEMICONDUCTOR.

# Examples for Illustration

| | Architecture | Number of Multiplications | Input Size | Quantization |
|---|---|---|---|---|
| **Face Detection** | VGG style | 290,816 | 32*32*3 | 16-bit fixed point |
| | VGG style | 14,353,920 | 90*90*3 | 16-bit fixed point |
| **Human Presence Detection** | VGG style | 8,570,880 | 64*64*3 | 16-bit fixed point |
| | VGG style | 338,558,976 | 128*128*3 | 16-bit fixed point |

LATTICE
SEMICONDUCTOR.

# Image Based Neural Networks on Lattice FPGAs

# Image Based Neural Networks Lattice Hardware



Himax HM01B0 UPduino Shield



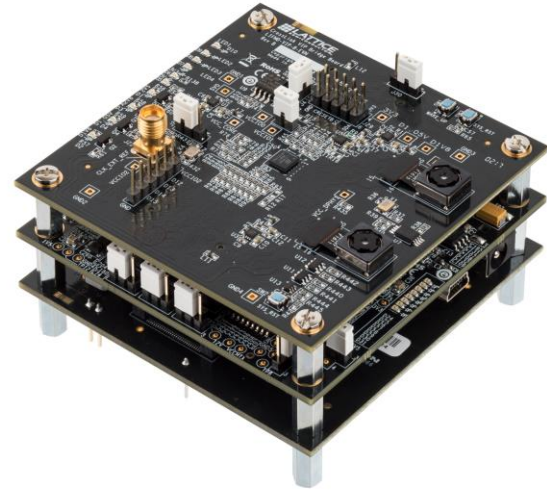Embedded Vision Development Kit

LATTICE
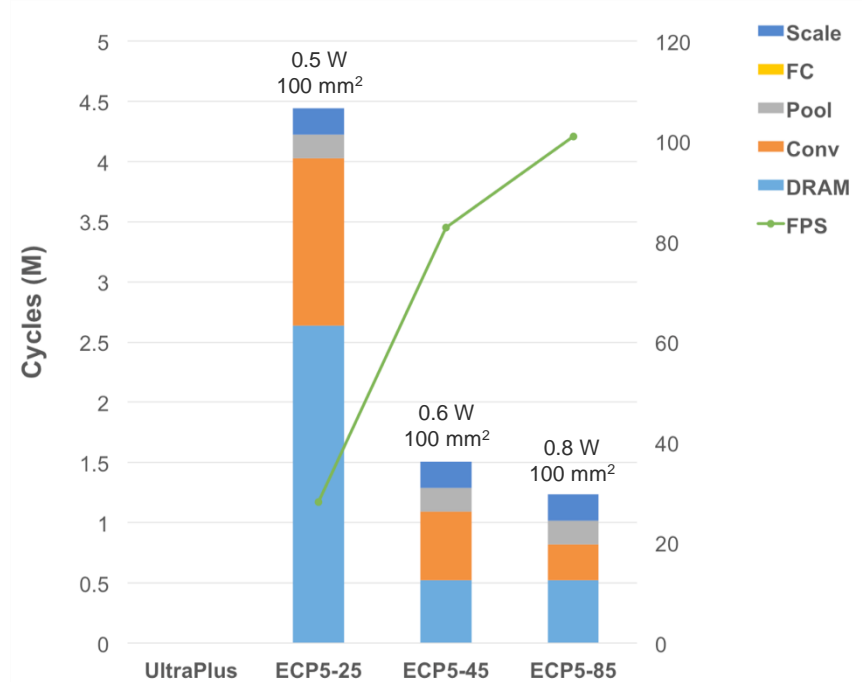SEMICONDUCTOR.

# Face Detect Implementations

## 32 x 32 Input



## 90 x 90 Input



* Running at 5 frames per second

# Human Presence Detect Implementations



64 x 64 Input

128 x 128 Input

* Running at 5 frames per second

LATTICE
SEMICONDUCTOR.

# Bringing It Together

| Network | Smallest Object | Device Size / Power / Performance | | | |
|---|---|---|---|---|---|
| | | UltraPlus 1 – 7 mW* 5.5 mm² | ECP5-25 0.5 W 100 mm² | ECP5-45 0.6 W 100 mm² | ECP5-85 0.8 W 100 mm² |
| Face Detection 32 x 32 Input | 50% | 465 | 3360 | 4511 | 5251 |
| Face Face Detection 90 x 90 Input | 20% | -- | 28 | 82 | 101 |
| Human Presence Detect 64 x 64 Input | 20% | 18 | 115 | 161 | 338 |
| Human Presence Detect 128 x 128 Input | 10% | -- | 2.3 | 3.5 | 5.4 |

* Running at 5 frames per second

LATTICE SEMICONDUCTOR.

# Summary

- AI at the edge solves real world problems

- FPGAs can implement AI standalone or in conjunction with other components

- sensAI stack components provide edge AI building blocks

  - Silicon, soft IP, tools, development boards & reference designs

- Configurable engine size and bit widths coupled with multiple target devices allows system optimization

  - 1 mW – 1 W

  - 5.5 mm$^2$ – 100 mm$^2$

LATTICE
SEMICONDUCTOR.

# Resources

Please visit [latticesemi.com/sensAI](latticesemi.com/sensAI) for more information and downloads

- 4 ECP5 Based Reference Designs / Demonstrations – Free
- 4 iCE40 Based Reference Designs / Demonstrations – Free
- CNN Accelerator IP – Free Evaluation
- CNN Compact Accelerator IP – Free
- Neural Network Compiler – Free
- Embedded Vision Development Kit – $199 Promotional Price
- Himax HM01B0 UPduino Shield – Available November ~$49

LATTICE
SEMICONDUCTOR.

# Empowering Product Creators to Harness Embedded Vision

The Embedded Vision Alliance (www.Embedded-Vision.com) is a partnership of 90+ leading embedded vision technology and services suppliers, and solutions providers

Mission: Inspire and empower product creators to incorporate visual intelligence into their products

The Alliance provides low-cost, high-quality technical educational resources for product developers

**Register for updates at www.Embedded-Vision.com**

The Alliance enables vision technology providers to grow their businesses through leads, ecosystem partnerships, and insights

**For membership, email us: membership@Embedded-Vision.com**

# Join us at the Embedded Vision Summit
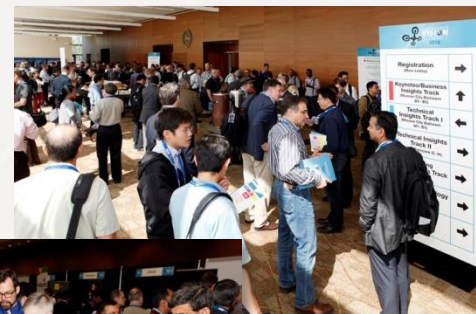## May 20-23, 2019—Santa Clara, California

*The only industry event focused on enabling product creators to create "machines that see"*

- *"Awesome!  I was very inspired!"*

- *"Fantastic. Learned a lot and met great people."*

- *"Wonderful speakers and informative exhibits!"*

**Embedded Vision Summit 2019 highlights:**

- **Inspiring keynotes** by leading innovators

- High-quality, practical **technical, business and product talks**

- Exciting **demos** of the latest apps and technologies

**Visit [www.EmbeddedVisionSummit.com](http://www.EmbeddedVisionSummit.com) to sign up for updates**

# Q & A

---

**LATTICE**
SEMICONDUCTOR.