# AI Application Processing Requirements

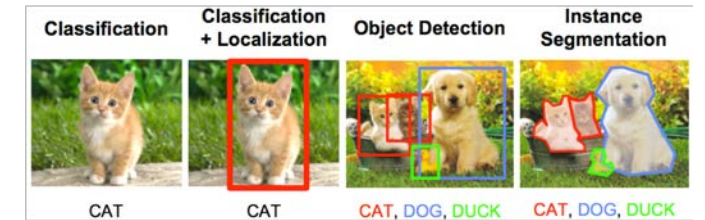| Low | Medium | High |
|---|---|---|



**Low**
- Sensor analysis
- Activity Recognition (motion sensors)
- Stress Analysis or Attention Analysis

**Medium**
- Audio & sound
- Speech Recognition
- Object detection

**High**
- Computer Vision
- Multiple Objects Detection/Classification/Tracking
- Speech Synthesis

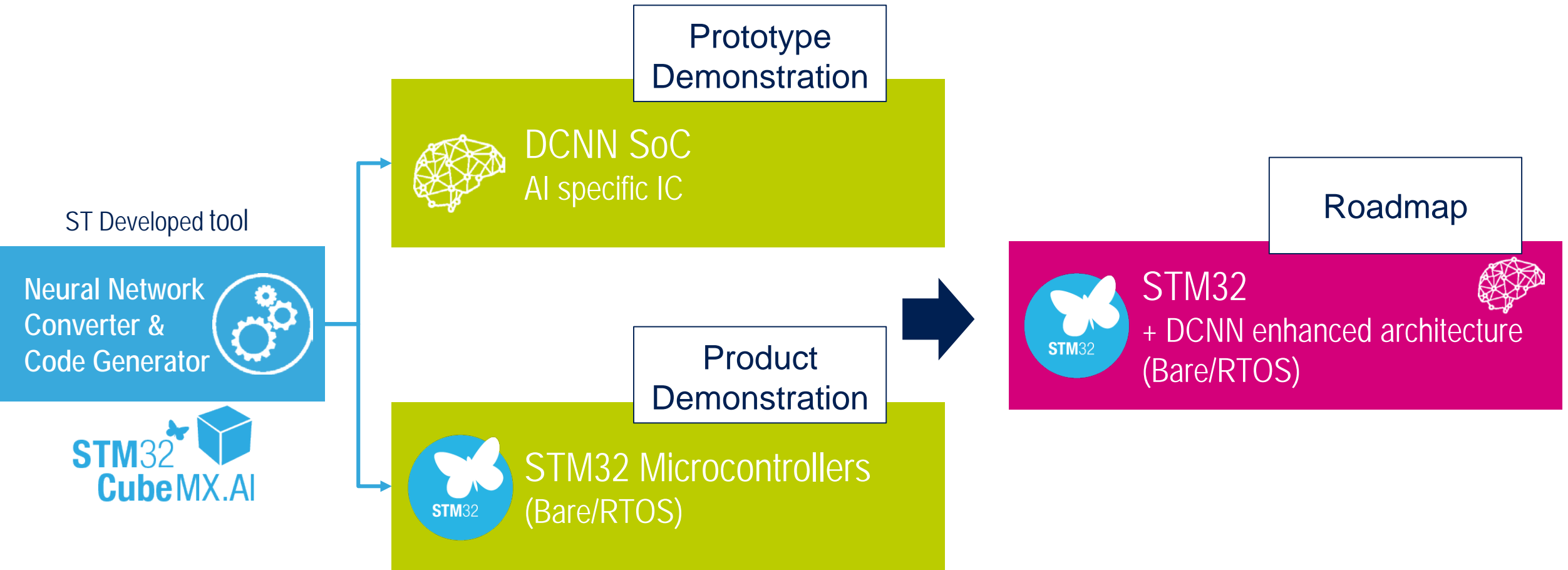STM32 → From IP embedded in MCU/MPU to dedicated SOC

STM32 CubeMX.AI

- Audio use cases with individual commands
- Classic motion sensor use cases

- Mandatory to support advanced Audio and Video complex use cases.

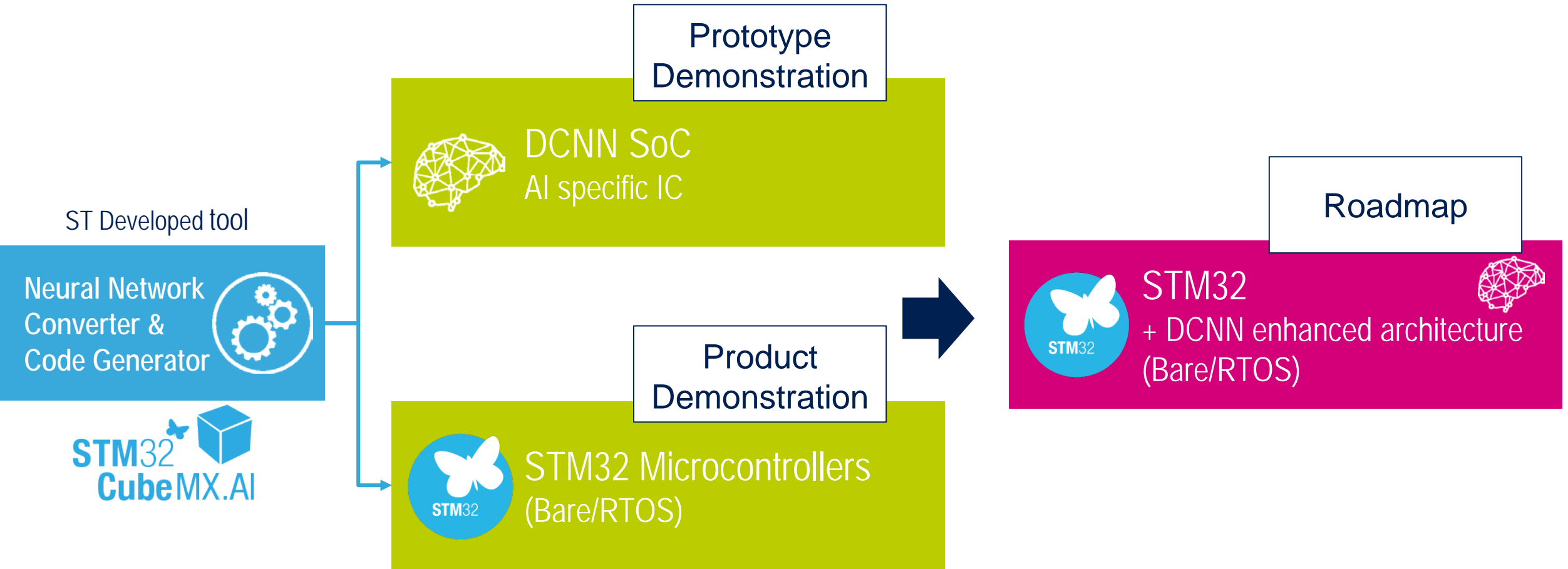life.augmented

# ST Solutions for Artificial Intelligence

ST Developed tool

**Neural Network Converter & Code Generator**

STM32 CubeMX.AI

Prototype Demonstration

**DCNN SoC**
AI specific IC

Product Demonstration

**STM32 Microcontrollers (Bare/RTOS)**

Roadmap

**STM32** + DCNN enhanced architecture (Bare/RTOS)

DCNN = Deep Convolutional Neural Network

life.augmented

# ST Solutions for Artificial Intelligence

ST Developed tool

**Neural Network Converter & Code Generator**

STM32 **Cube**MX.AI

Prototype Demonstration

DCNN SoC
AI specific IC

Product Demonstration

STM32 Microcontrollers
(Bare/RTOS)

Roadmap

STM32
+ DCNN enhanced architecture
(Bare/RTOS)

DCNN = Deep Convolutional Neural Network

*life.augmented*

# ST Enables AI on STM32

## Benefits of Deep Learning now available across all STM32 portfolio

Input your Framework dependent, Pre-Trained neural network into **STM**32**Cube**MX.AI

Automatic and fast generation of an STM32-optimized library

**STM**32**Cube**MX.AI guarantees interoperability with state-of-the-art Deep Learning Design Frameworks
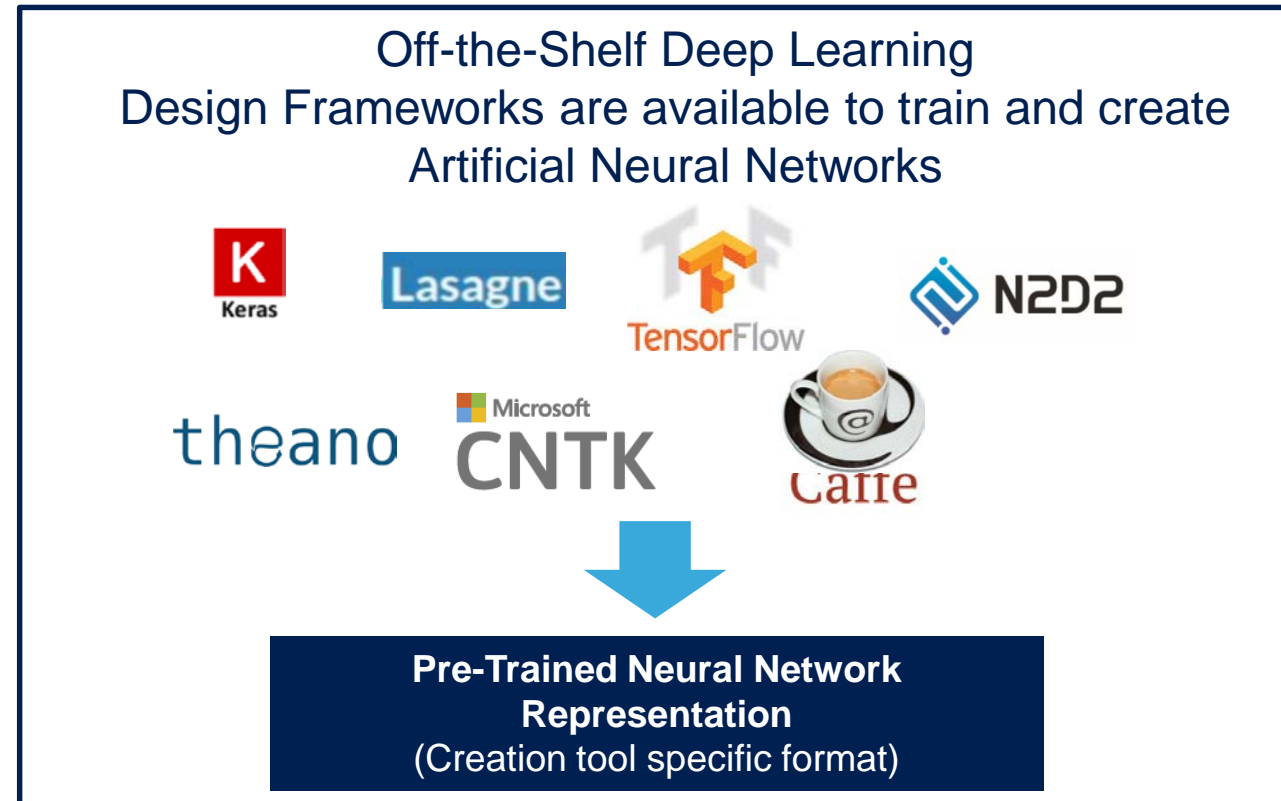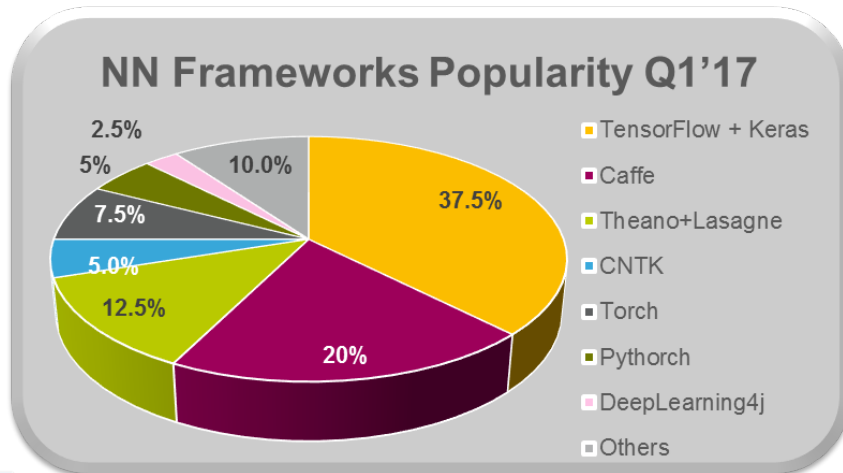
life.augmented

# Neural Networks Ecosystem

## Benefits of Deep Learning now available across all STM32 portfolio

How to create a Neural Network?

1. Define a problem
2. Collect/find/label data
3. Define topology
4. Design with

**NN Frameworks Popularity Q1'17**

- TensorFlow + Keras — 37.5%
- Caffe — 20%
- Theano+Lasagne — 12.5%
- CNTK — 5.0%
- Torch — 7.5%
- Pythorch — 5%
- DeepLearning4j — 2.5%
- Others — 10.0%

Off-the-Shelf Deep Learning
Design Frameworks are available to train and create
Artificial Neural Networks

Keras    Lasagne    TensorFlow    N2D2

theano    Microsoft CNTK    Caffe

**Pre-Trained Neural Network Representation**
(Creation tool specific format)

# Neural Networks
## Available Now for STM32

Benefits of Deep Learning now available across all STM32 portfolio

**Off-the-shelf** :
Pre-trained Neural
Network Model

Deep Learning
Framework dependent

**STM32 CubeMX.AI**

**Embedded Solution**
Optimized Neural
Network Code
generated for STM32

**Deep Learning SW Solution**

... brings your AI-based innovation to the existing STM32 Portfolio

# STM32 CubeMX.AI : Architecture

## Benefits of Deep Learning now available across all STM32 portfolio

**Off-the-shelf** :
Pre-trained Neural Network Model

Deep Learning Framework dependent
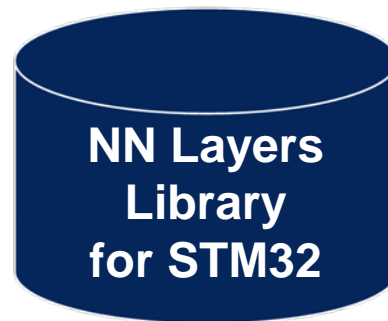
**Neural Network Exporter**

DL Framework Independent Neural Network Representation

**Code Generator**

**NN Layers Library for STM32**

Neural Networks API's

**Embedded Solution**
Optimized Neural Network Code generated for STM32

This optimized STM32 neural network model can be included into the user project (using KEIL, IAR, OpenSTM32) and can be compiled and ported onto the final device for field trials

# Artificial Intelligence is Everywhere



Gaming

Drone

Companion Robot

Mobile

Domestic Robot

Surveillance

Security/Eye tracking

Virtual/augmented Reality

Smart home

Relationship robot

Retail

# Artificial Intelligence for Everything

## Analysis

**Where am I ?**

- Scene classification (audio, video, environmental sensors)

**Which objects are in the scene, where are they?**

- Video object detection/classification

**What am I doing?**

- Activity recognition (audio, video, inertial sensors)

**What's happening?**

- Event recognition (audio, video, inertial sensors, environmental sensors).

## User Interaction

- Command detection (audio)
- Speech Recognition (audio)
- Gesture Recognition (inertial sensors, video)
- User identification and mood detection (audio, video)

## Continuous Learning

- How can I detect unpredictable, unclassified events in dynamic environments?
- Recurrent networks (audio, video, inertial sensors, environmental sensors)

life.augmented

# Distributed Artificial Intelligence is a Must
## to Increase Systems Efficiency

NODES

| | | |
|---|---|---|
| **Processing** **Connectivity** **Security** **Sensing & Actuating** | **Processing** **Connectivity** **Security** | **Processing** **Connectivity** **Security** |
| Sensor Data / Actions | Sensor Data / Actions | |
| 1 Sensor | 100 Sensors | 10,000 Sensors |

Processing Requirements

Connectivity/Bandwidth Requirements

# Neural Networks are Key
## for Intelligent Nodes

**What is an Artificial Neural Network?**

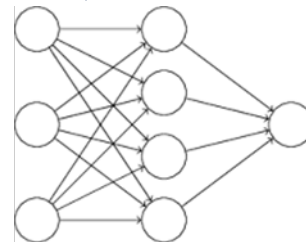**How to design ANN's?**

### Artificial Intelligence

#### Machine Learning

**Deep Learning**

**Subset of Machine Learning algorithms based on Artificial Neural Networks (DNN, CNN, RNN, SOANN, ect)**

- Models for input-output transfer function approximation inspired by biological neural networks.
- Capable of modeling and processing highly time varying and non linear relationships between inputs and outputs.
- Exponentially faster and more efficient than traditional computer processing models for typical AI uses like detection, classification, prediction.



Neural Network

Example: A neural network trained to classify an object in a picture.

Off-the-Shelf Deep Learning Design Frameworks are available to design and train Artificial Neural Networks



**Pre-Trained Neural Network Representation**
(Creation tool specific format)

life.augmented

# ST Enables AI @ the Edge

**Pre-Trained Neural Network Representation** (Creation tool specific format)

**Automatic ST Tool**
**Neural Network Optimizer & Code Generator**

**Code for ARM Cortex-M Microcontroller**

**Code for ST AI-Specific IC**

## STM32 Microcontrollers
(Bare/RTOS)

**Benefits**
- Available today
- Runs on any STM32

**Typical use cases**
- Smart Industry: preventive maintenance
- Wearable: Human activity recognition
- Consumer : Entry level IoT node
- Audio : Intelligent microphones

## ULTRA HIGH POWER/AREA EFFICIENCY NEURAL NETWORK SOC

**Benefits**
- Performance/power optimized
- Intensive data & processing applications

**Typical use cases**
- Smart Driving: In-vehicle driver monitoring
- Smart surveillance cameras
- Consumer : AR, Drones, Robots

life.augmented

# HW Accelerated Deep Learning SoC



A configurable, scalable and design time parametric **Convolutional Neural Network Processing Engine**
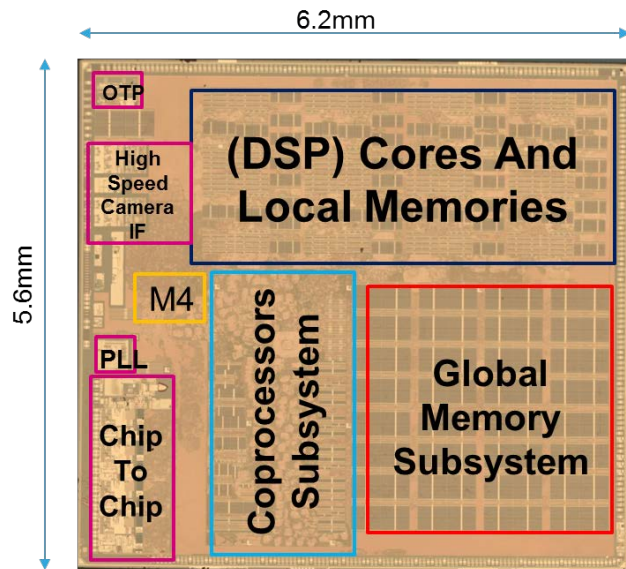
DCNN Convolutional Layers accounts for more than 90% DCNN operations, hence **8 Convolution HW Accelerators** allow high efficiency in area vs GOPS vs power consumption

In addition to **ARM Cortex M4, 8 DSP Clusters** allow both programmability and flexible mapping of diversified, custom DCNN's

**Embedded Memory** enables further reduction of power consumption required by IOT applications.

# DCNN SoC Test Chip Main Features



6.2mm

5.6mm

OTP

High Speed Camera IF

(DSP) Cores And Local Memories

M4

PLL

Chip To Chip

Coprocessors Subsystem

Global Memory Subsystem

(*) 1 MAC defined as 2 OPS (ADD + MUL)

| Technology | FD-SOI 28nm |
|---|---|
| Package | FBGA 15x15x1.83 |
| Clock freq | 200MHz – 1.175GHz |
| Supply voltages | 0.575V – 1.1V digital – 1.8V I/O |
| On-chip RAM | 4x1 MB (Global), 8x192 KB (DSP), 128 KB (Host) |
| Host | ARM® Cortex®-M4 |
| DSPs Nr | 16 |
| Peak DSP performance (1.175GHz, 1.1V) | 75 GOPS (dual 16b MAC loop) |
| Convolutional Accelerators Nr | 8 |
| CA size (including local memory) | 0.27 sqmm |
| Max Tot CAs performance (1.175GHz, 1.1V) | 676 GOPS |
| Power (**) @200MHz, 0.6V, 2 CAs | 37.5mW @ 10 FPS |
| Power (**) @1GHz, 1V, 8 CAs | 600mW @ 60FPS (not optimized yet) |
| CAs Peak Efficiency @200Mhz, 0.575V (Alexnet) | 2.9 TOPS/W |

life.augmented