

tinyML[®] Talks

Enabling Ultra-low Power Machine Learning at the Edge

“Universal CNN Accelerator Intended for Edge-Based AI Inference”

Rastislav Struharik - University of Novi Sad

Germany Area Group – April 7, 2021



www.tinyML.org



tinyML Talks Sponsors



tinyML Strategic Partner



EDGE IMPULSE



maxim
integrated™



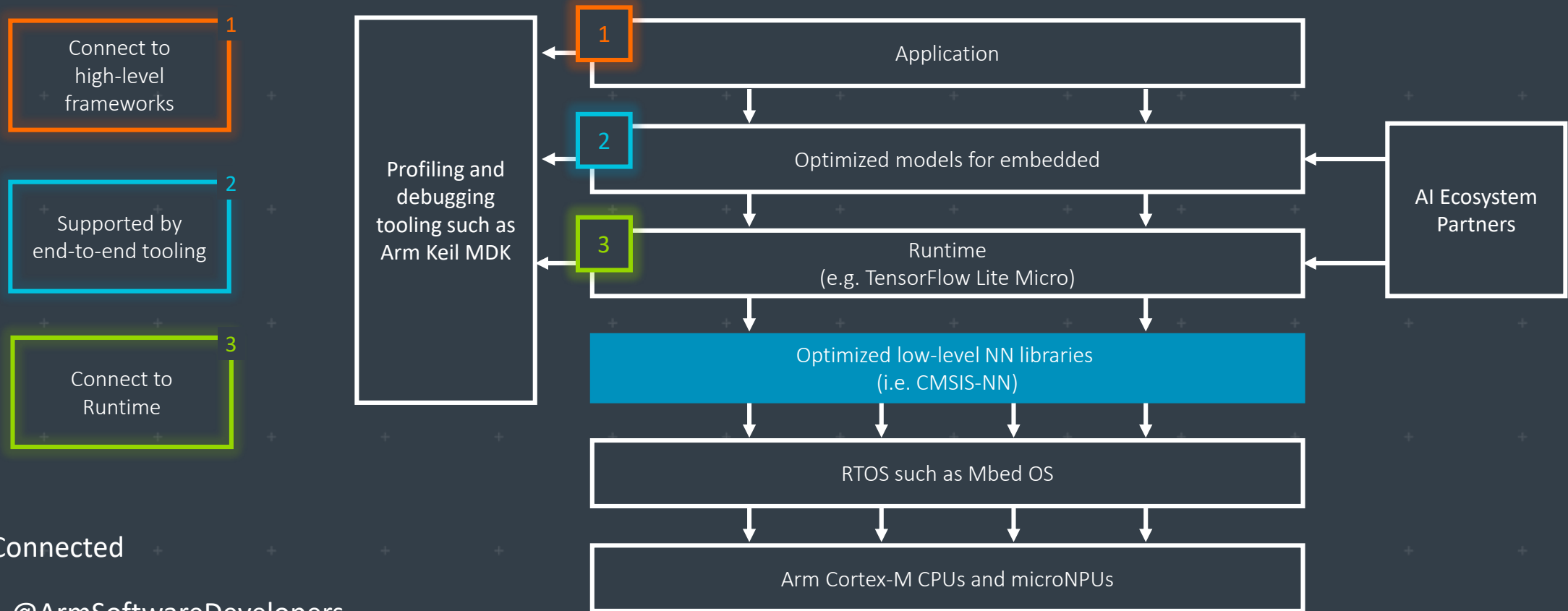
Reality AI®



SynSense

Additional Sponsorships available – contact Olga@tinyML.org for info

Arm: The Software and Hardware Foundation for tinyML



Stay Connected

 @ArmSoftwareDevelopers

 @ArmSoftwareDev

Resources: developer.arm.com/solutions/machine-learning-on-arm



WE USE AI TO MAKE OTHER AI FASTER, SMALLER AND MORE POWER EFFICIENT



Automatically compress SOTA models like MobileNet to <200KB with **little to no drop in accuracy** for inference on resource-limited MCUs



Reduce model optimization trial & error from weeks to days using Deeplite's **design space exploration**



Deploy more models to your device without sacrificing performance or battery life with our **easy-to-use software**

BECOME BETA USER bit.ly/testdeeplite

mobilityXlab

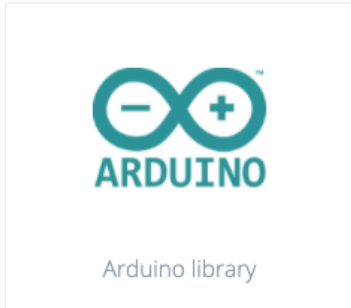
arm



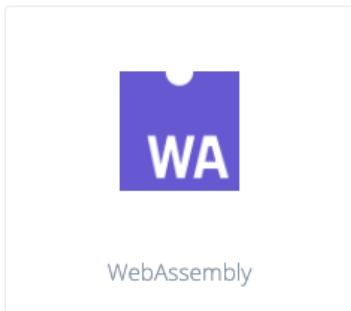
TinyML for all developers



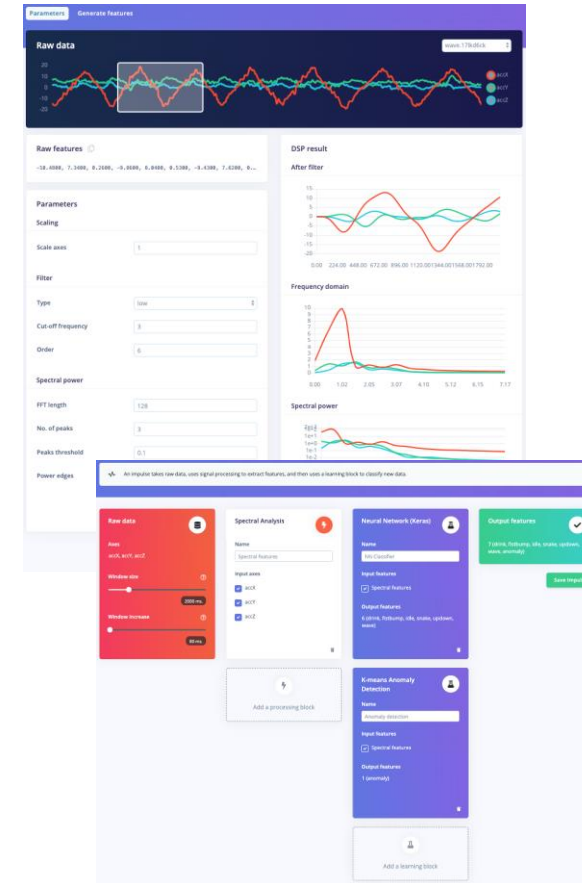
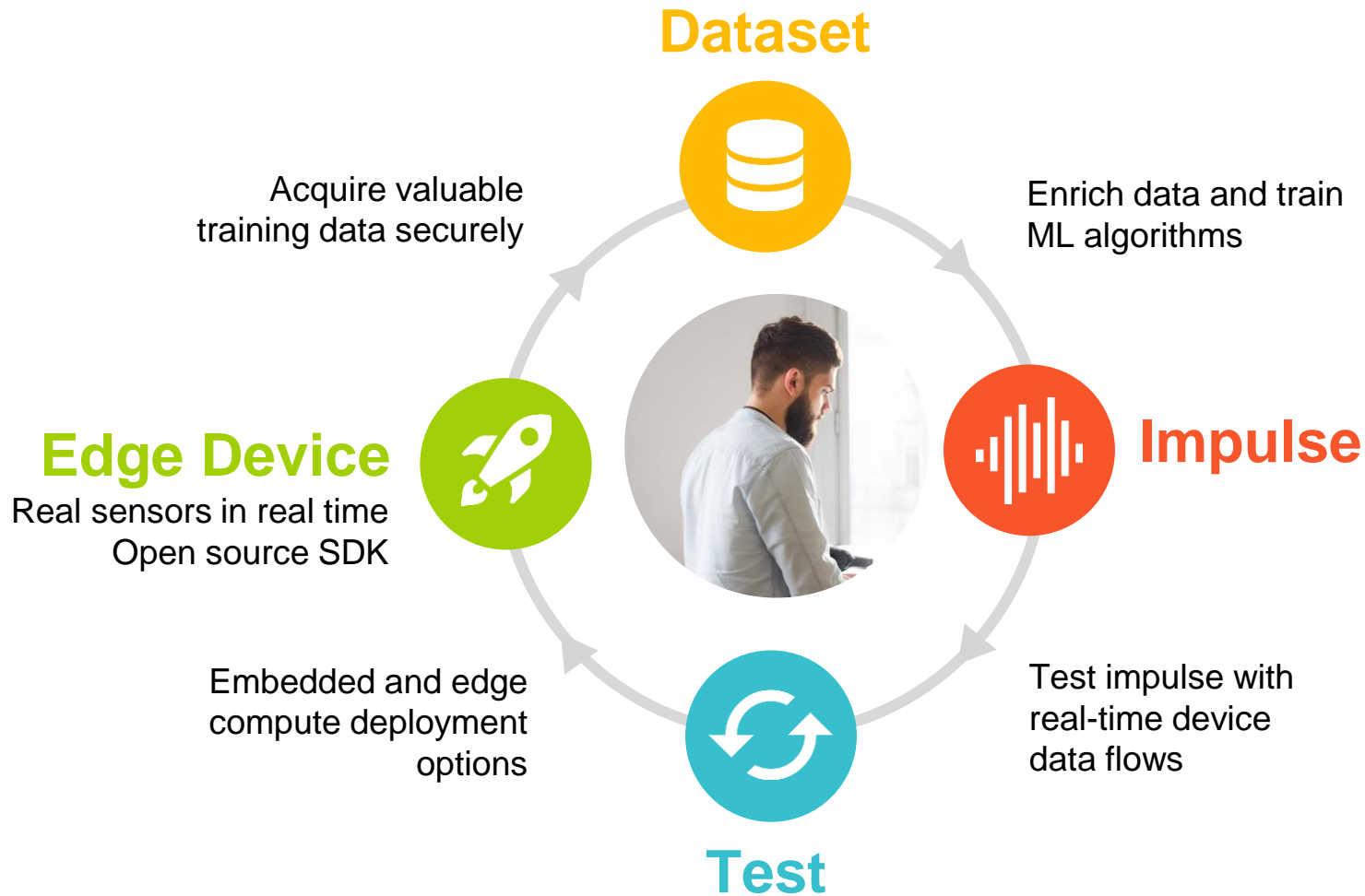
C++ library



Arduino library

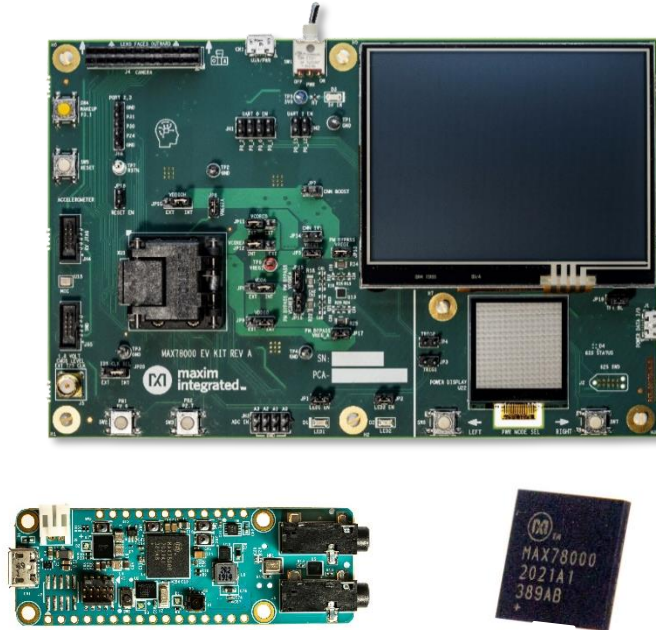


WebAssembly



Maxim Integrated: Enabling Edge Intelligence

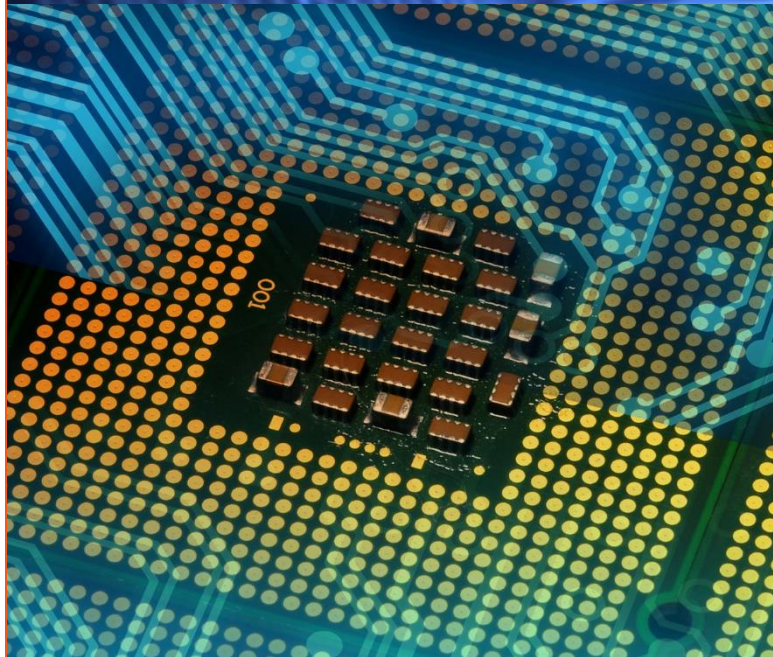
Advanced AI Acceleration IC



The new MAX78000 implements AI inferences at low energy levels, enabling complex audio and video inferencing to run on small batteries. Now the edge can see and hear like never before.

www.maximintegrated.com/MAX78000

Low Power Cortex M4 Micros



Large (3MB flash + 1MB SRAM) and small (256KB flash + 96KB SRAM, 1.6mm x 1.6mm) Cortex M4 microcontrollers enable algorithms and neural networks to run at wearable power levels.

www.maximintegrated.com/microcontrollers

Sensors and Signal Conditioning



Health sensors measure PPG and ECG signals critical to understanding vital signs. Signal chain products enable measuring even the most sensitive signals.

www.maximintegrated.com/sensors

Qeexo AutoML

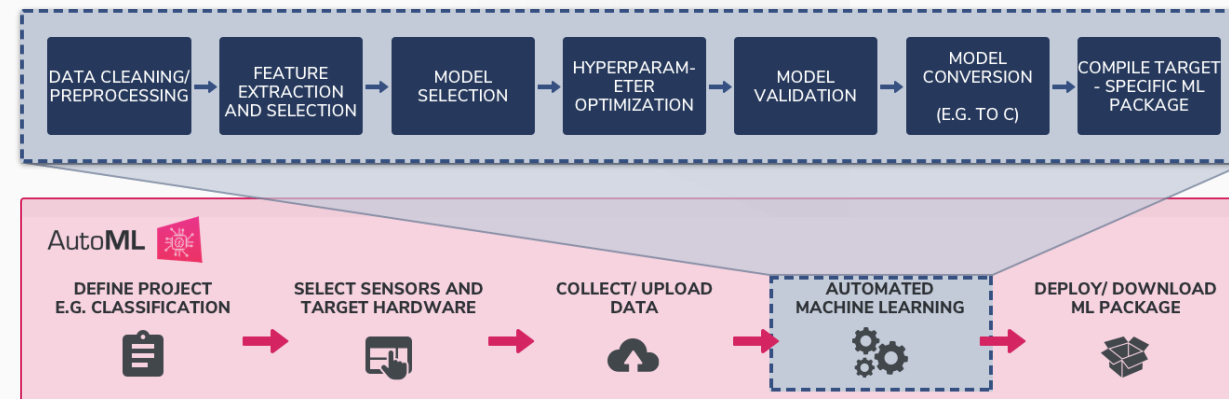


Automated Machine Learning Platform that builds tinyML solutions for the Edge using sensor data

Key Features

- Supports 17 ML methods:
 - Multi-class algorithms: GBM, XGBoost, Random Forest, Logistic Regression, Gaussian Naive Bayes, Decision Tree, Polynomial SVM, RBF SVM, SVM, CNN, RNN, CRNN, ANN
 - Single-class algorithms: Local Outlier Factor, One Class SVM, One Class Random Forest, Isolation Forest
- Labels, records, validates, and visualizes time-series sensor data
- On-device inference optimized for low latency, low power consumption, and small memory footprint applications
- Supports Arm® Cortex™- M0 to M4 class MCUs

End-to-End Machine Learning Platform



For more information, visit: www.qeexo.com

Target Markets/Applications

- Industrial Predictive Maintenance
- Smart Home
- Wearables
- Automotive
- Mobile
- IoT



Reality AI[®]

Add Advanced Sensing to your Product with Edge AI / TinyML

<https://reality.ai>



info@reality.ai



[@SensorAI](https://twitter.com/SensorAI)



[Reality AI](https://www.linkedin.com/company/reality-ai)

**Pre-built Edge AI sensing modules,
plus tools to build your own**

Reality AI solutions

Prebuilt sound recognition models for
indoor and outdoor use cases

Solution for industrial anomaly detection

Pre-built automotive solution that lets cars
“see with sound”

Reality AI Tools[®] software

Build prototypes, then turn them into
real products

Explain ML models and relate the function
to the physics

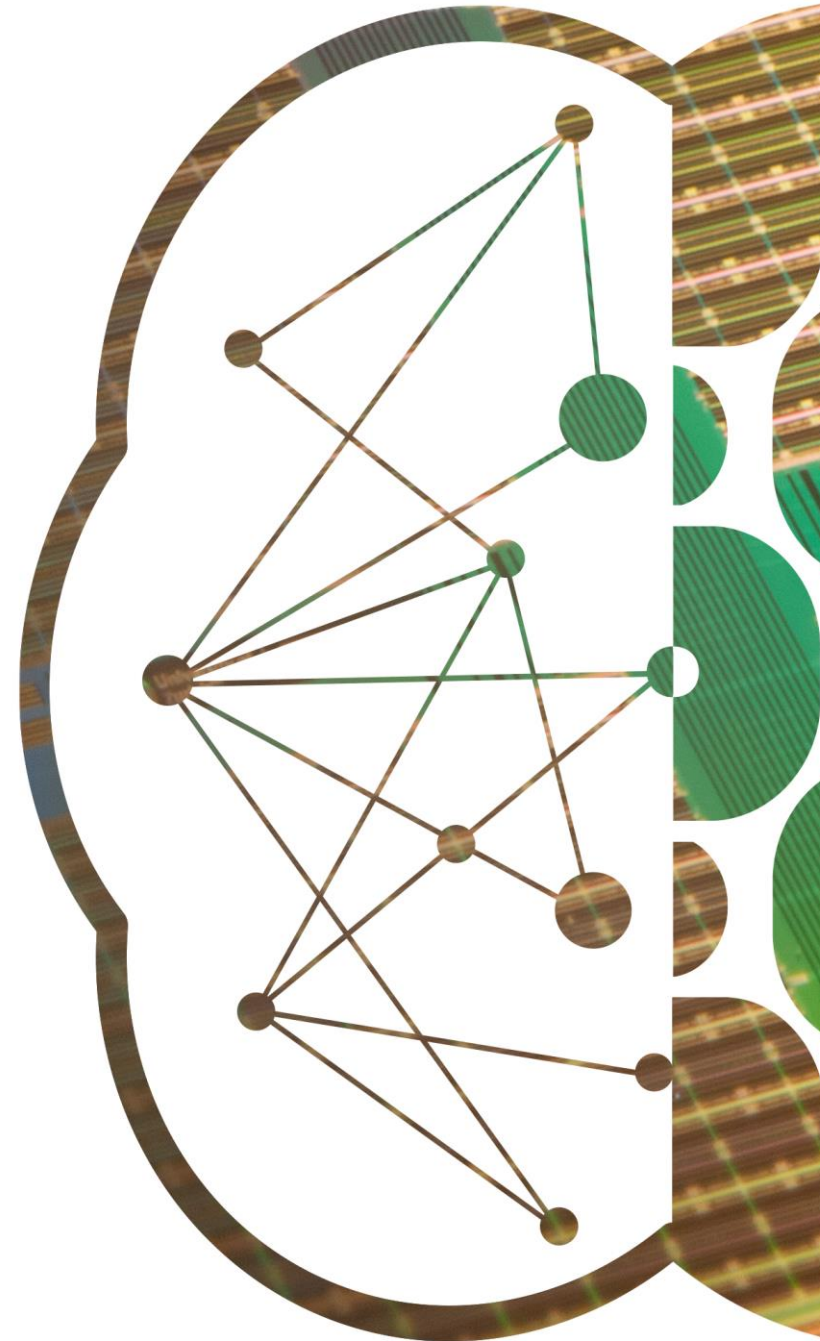
Optimize the hardware, including
sensor selection and placement



SynSense

SynSense builds **sensing and inference** hardware for **ultra-low-power** (sub-mW) **embedded, mobile and edge** devices. We design systems for **real-time always-on smart sensing**, for audio, vision, IMUs, bio-signals and more.

<https://SynSense.ai>



KIM-Labs

Our offer for small and medium-sized enterprises (SMEs):

- AI trainings
- AI events
- AI development projects for your application
- AI Working Groups & Expert Tables
- AI Co-Working Spaces & Prototyping Trainings





Next tinyML Talks

Date	Presenter	Topic / Title
Tuesday, April 13	Bernhard Suhm Machine Learning Product Manager, MathWorks	Deploying AI to Embedded Systems

Webcast start time is 8 am Pacific time

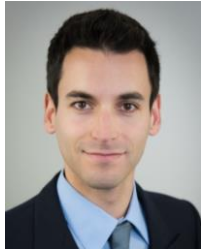
Please contact talks@tinymml.org if you are interested in presenting



Local Committee in Germany



Alexis Veynachter,
Master Degree in Control Engineering, Senior Field Application
Engineer
Infineon 32bits MCUs for Sensors, Fusion & Control



Carlos Hernandez-Vaquero
Software Project Manager, IoT devices
Robert Bosch



Prof. Dr. Daniel Mueller-Gritschneider
Interim Head - Chair of Real-time Computer Systems
Group Leader ESL - Chair of Electronic Design Automation
Technical University of Munich



Marcus Rüb
Researcher in the field of TinyML
Hahn-Schickard

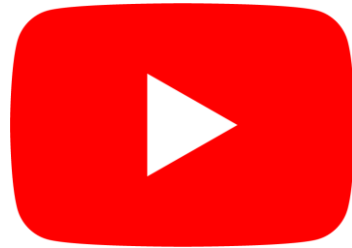


Reminders

Slides & Videos will be posted tomorrow

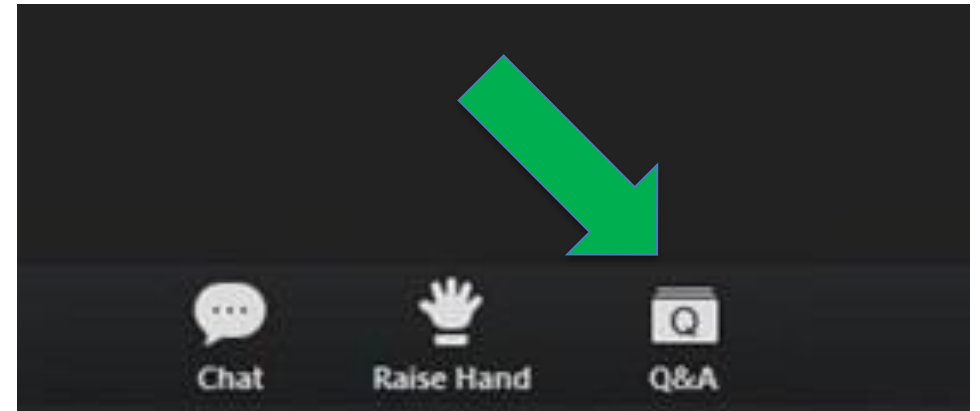


tinyml.org/forums



youtube.com/tinyml

Please use the Q&A window for your questions



Rastislav Struharik



Rastislav Struharik is a full professor at the Department of Power, Electronics and Telecommunications, Faculty of Technical Sciences, University of Novi Sad, Serbia. He received his PhD in Electronics in 2009, in the area of hardware acceleration of machine learning algorithms. During his academic career he has published more than 35 papers in international journals and conferences, mainly focusing on the hardware acceleration of machine learning algorithms, such as Decision Trees, Support Vector Machines, Artificial Neural Networks, Convolutional Neural Networks, and Ensemble Classifiers, targeting both learning and inference algorithms. For the past three years he has also been working as the chief architect for the IDS own FPGA IP core technology for hardware acceleration of Convolutional Neural Networks, intended for edge AI applications.



IT'S SO EASY.

WITH IDS IMAGING SOLUTIONS

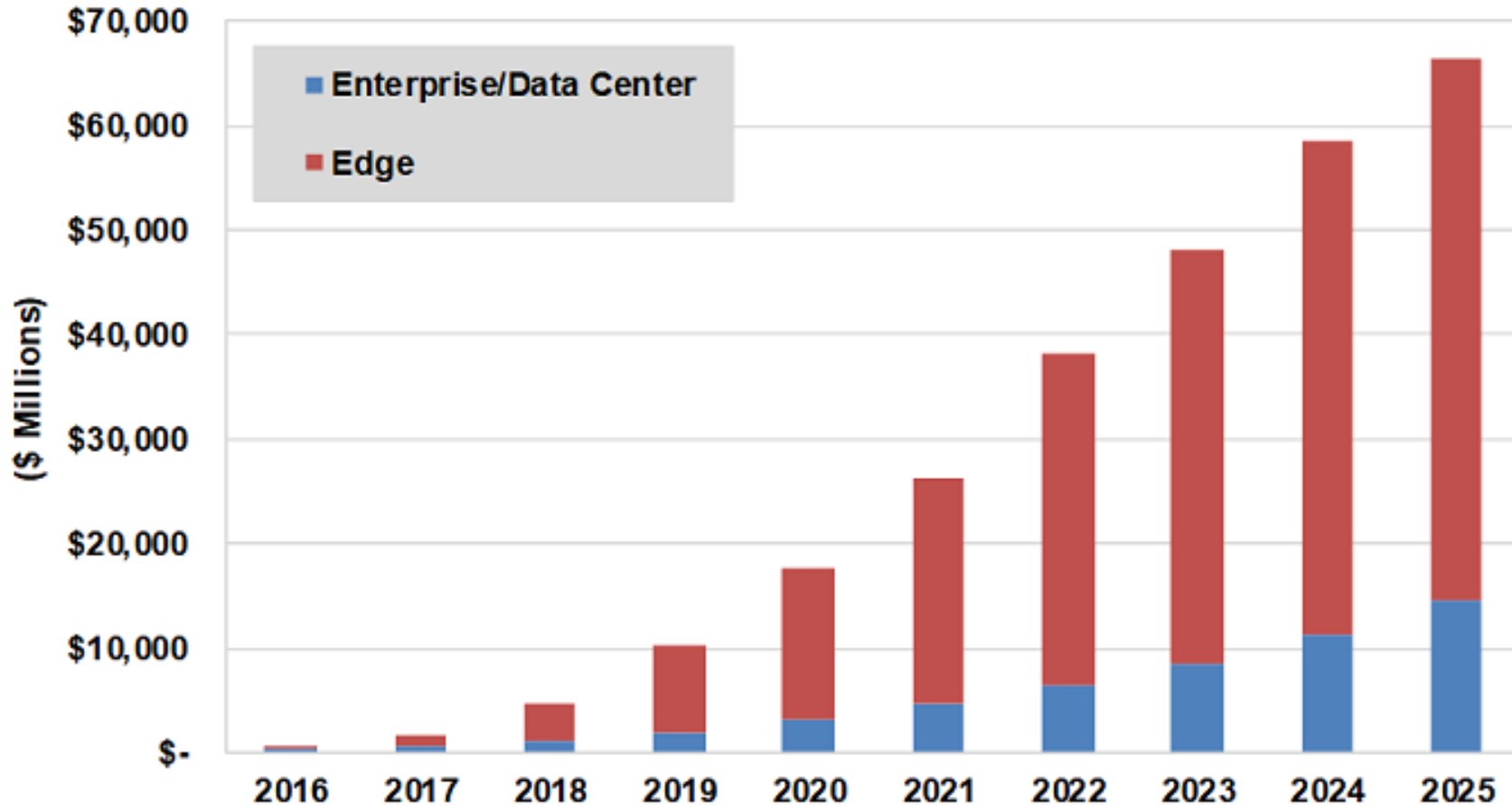


Universal CNN Accelerator for Edge-based AI Inference

Rastislav Struharik



Edge AI – The next big thing?



Source: Tractica

AI processing

→ on the Edge, where things are going to get interesting

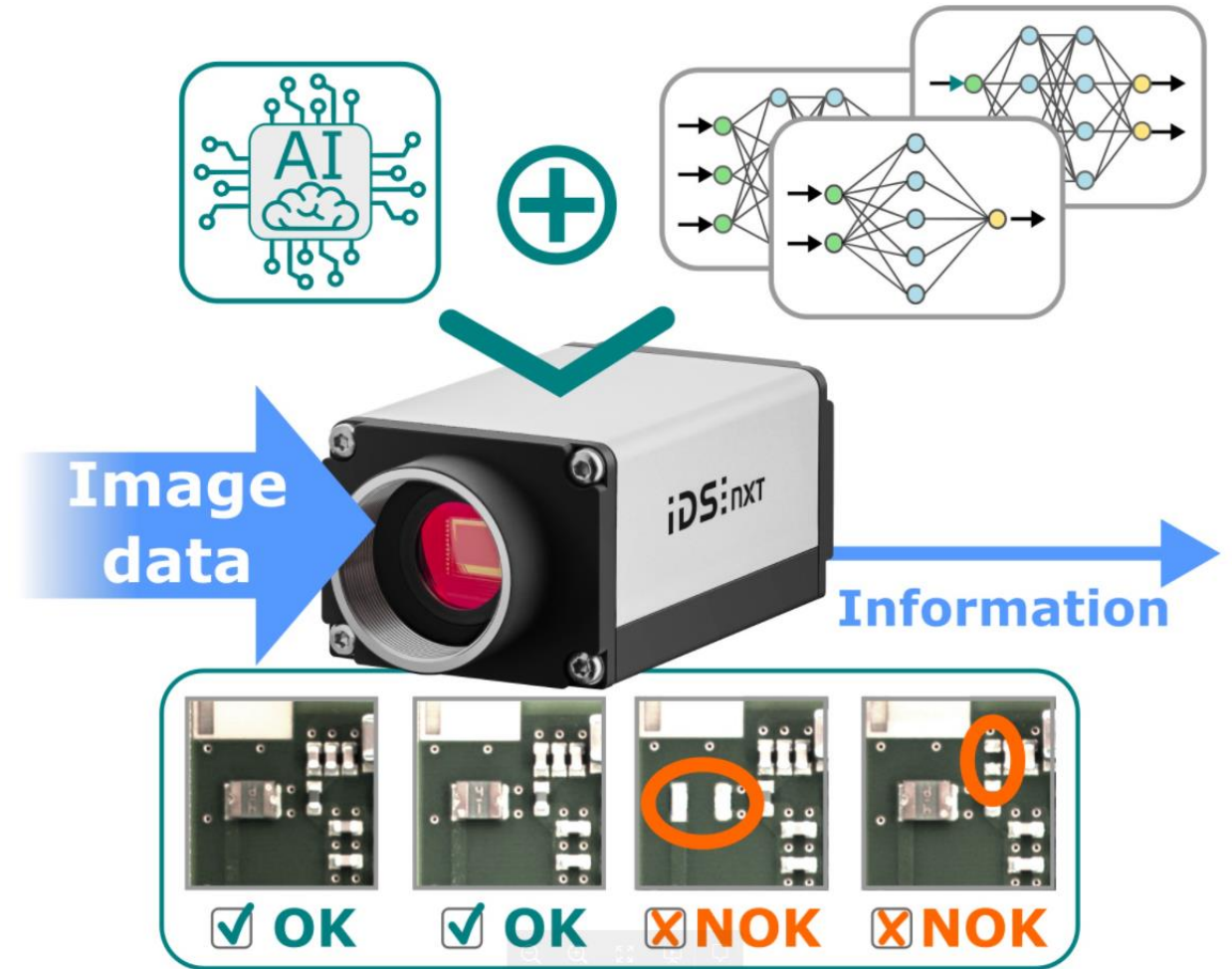
AI training

→ will stay in the cloud, where Nvidia's GPUs are most suited for the job

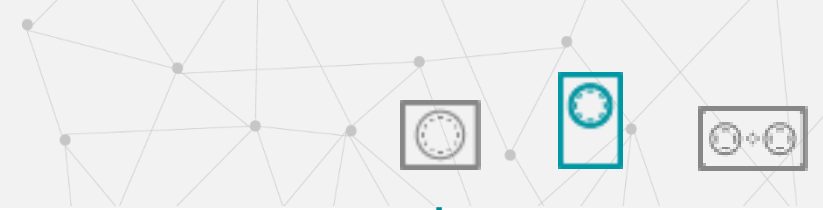
It's so easy!

Why AI on the edge?

- Bandwidth and Latency
- Security and Privacy
- Reliability and Availability
- Customization



It's so easy!



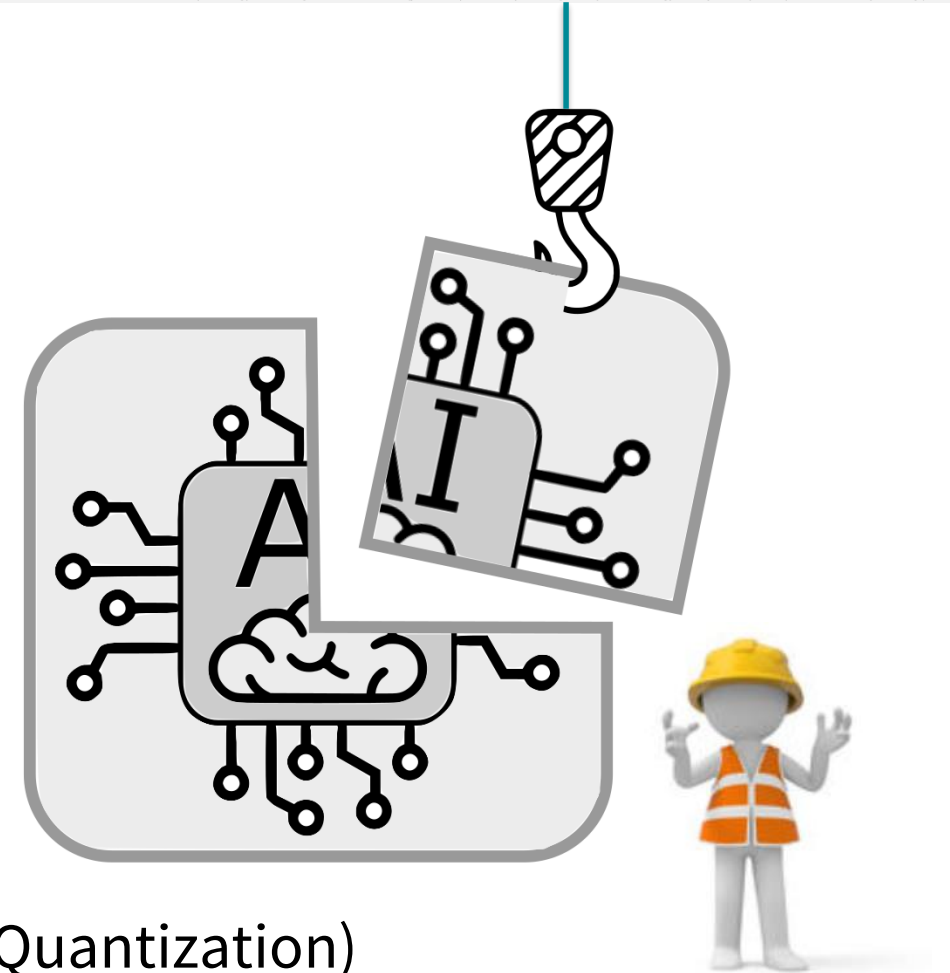
Challenges in deploying AI on the edge

- **Edge Computing**

- is all about efficiency
- requires low latency

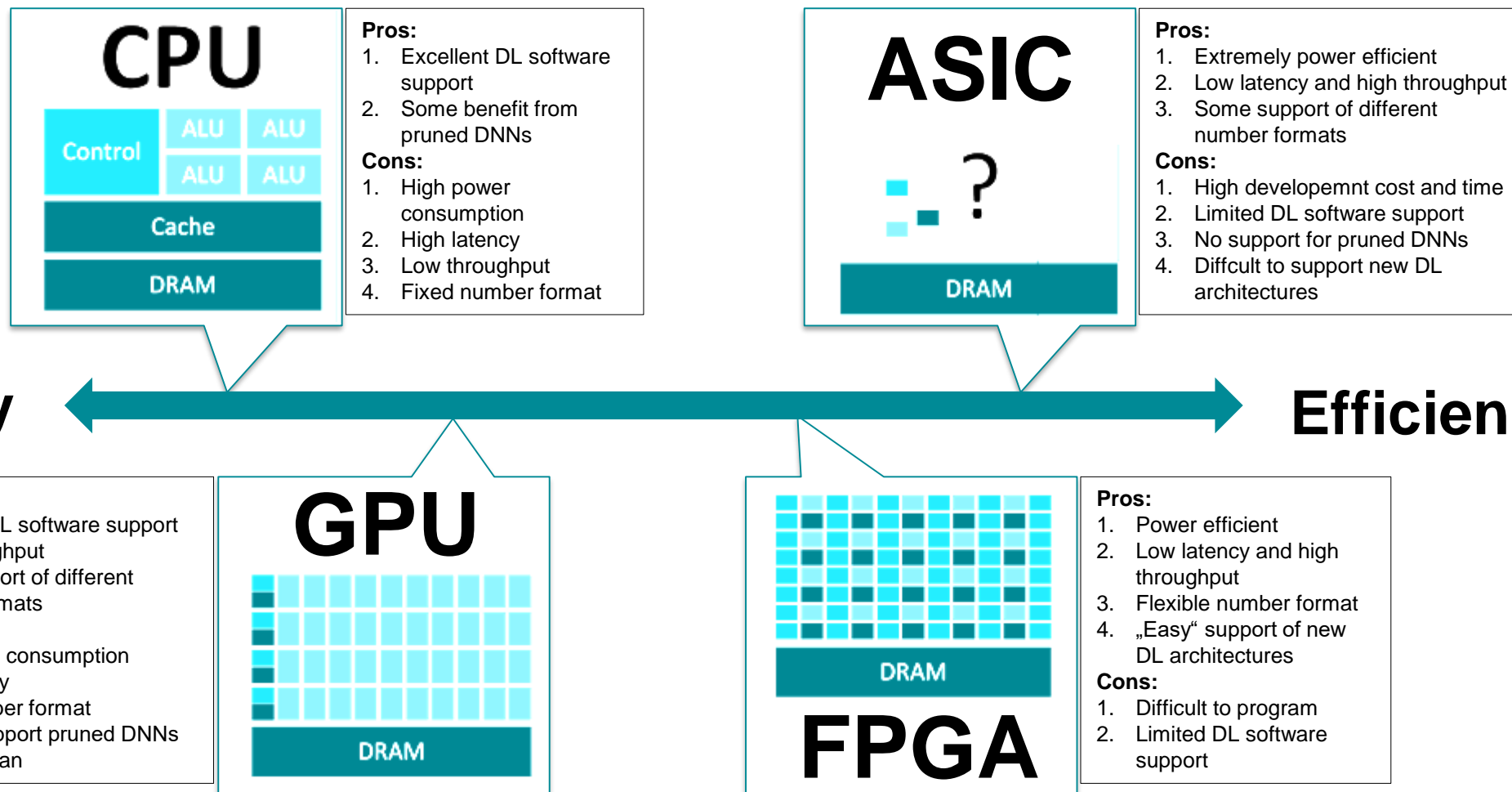
- **Edge AI requires:**

- Parameter Efficient Networks (MobileNet, EfficientNet, SqueezeNet, ...)
- Network Compression (by using Pruning and Quantization)
- Customized Computing Systems





Available options when implementing CNNs on the edge

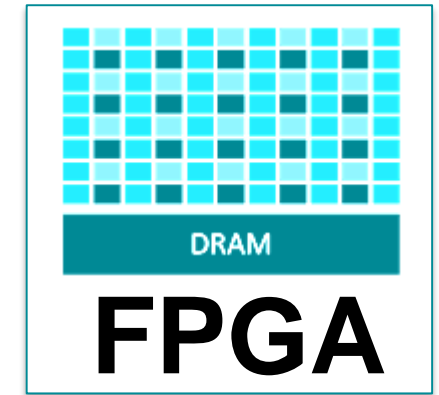


It's so easy!



Benefits of using FPGA technology for DL acceleration

- **Customizable**
- **Optimizable** for specific types of architectures, CNNs
 - • reduces power requirements
 - higher performance
- **Longer lifespan** → 2-5 times that of GPUs
- **More resistant** to rugged settings and environmental factors
- **Reconfigurable**
 - • ideal when algorithms change frequently,
 - clear advantage over ASIC solutions

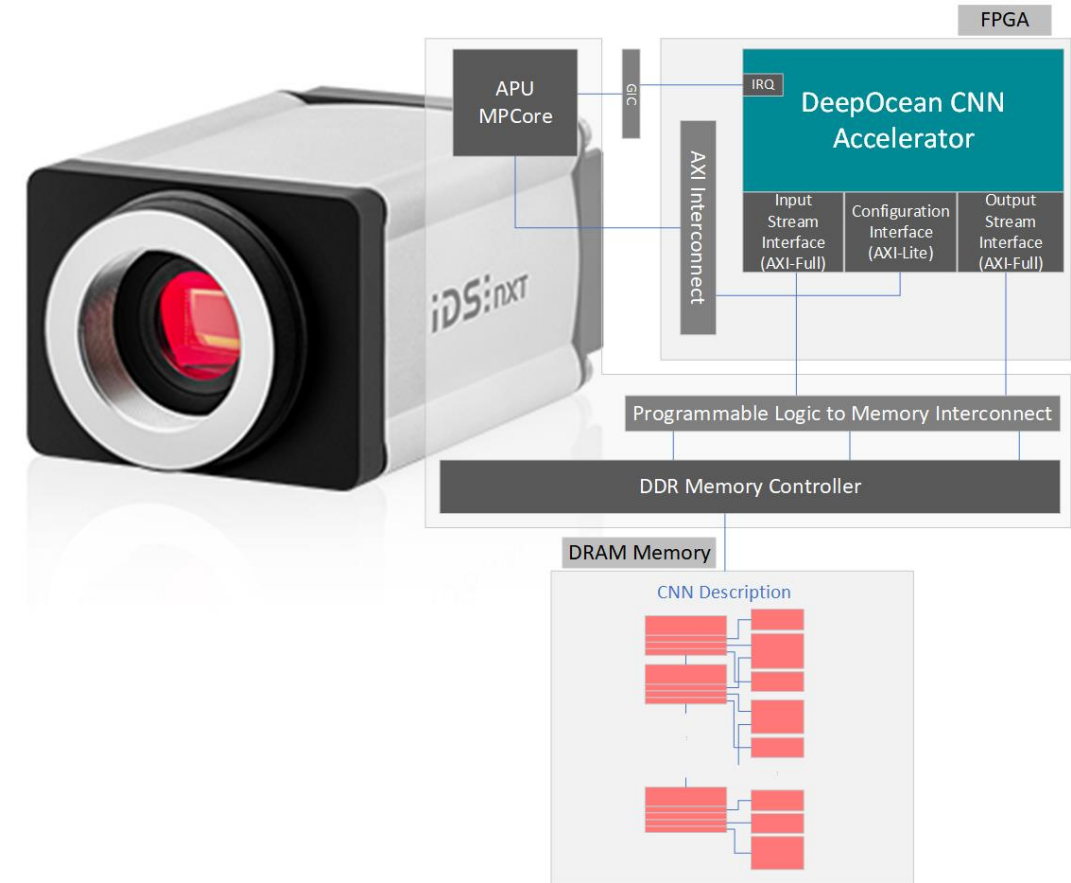


It's so easy!



IDS FPGA implementation strategy

- **ONE, UNIVERSAL** architecture
- **Any CNN** can run on developed CNN accelerator
- No “Difficult to Program” issue
- Great **flexibility**
- Support of **newest DL algorithms**
- Good DL software **support**



It's so easy!



CNN implementation flow using „Deep Ocean” accelerator

STEP 1:

Prune selected CNN

(DeepCompressor Tool)



STEP 2:

Translate pruned CNN

1. Select number representation
2. Create binary description

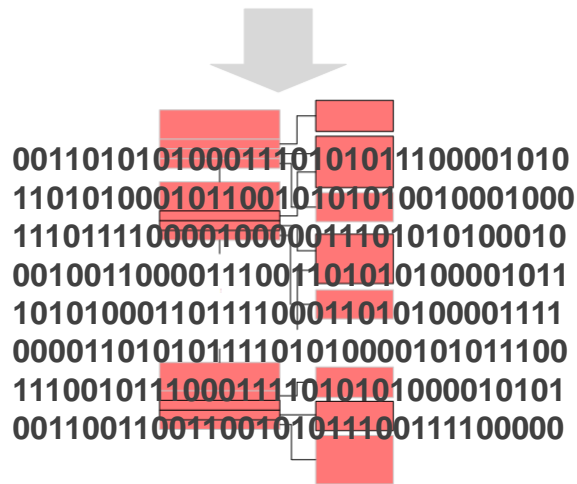
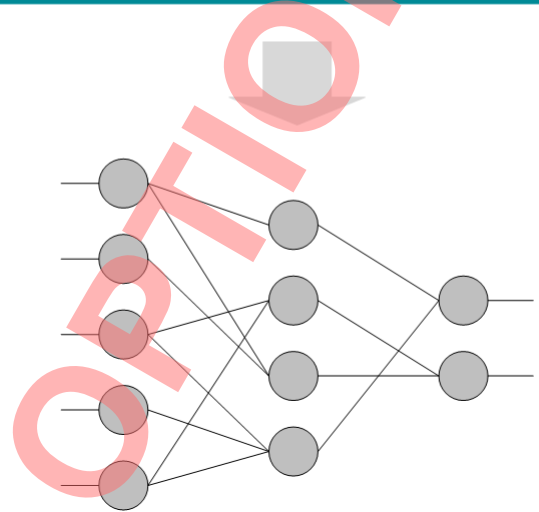
(DeepTranslator Tool)



STEP 3:

Run the CNN on Camera

1. Download CNN binary description to camera
2. Configure and start DeepOcean CNN accelerator



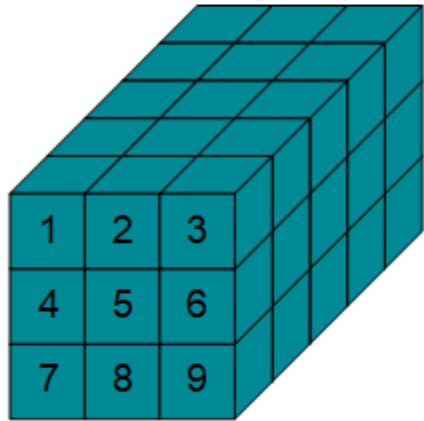
It's so easy!



Pruning CNNs

Before pruning

Every weight value is different from zero



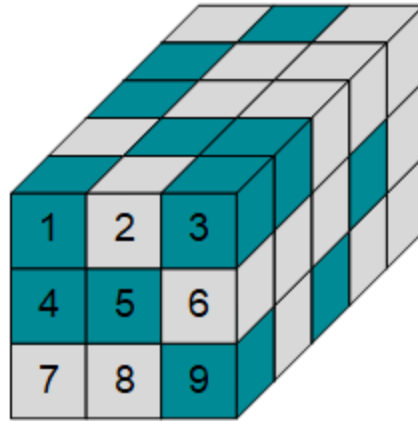
Original 3x3x5 dense kernel

Pruning Algorithm
Removes specified
% of "least
important" weights



After pruning

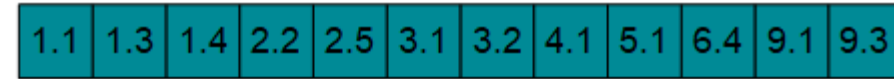
Some weight values equal to zero
(shown in light grey)



Sparse 3x3x5 kernel



A list of non-zero valued weights

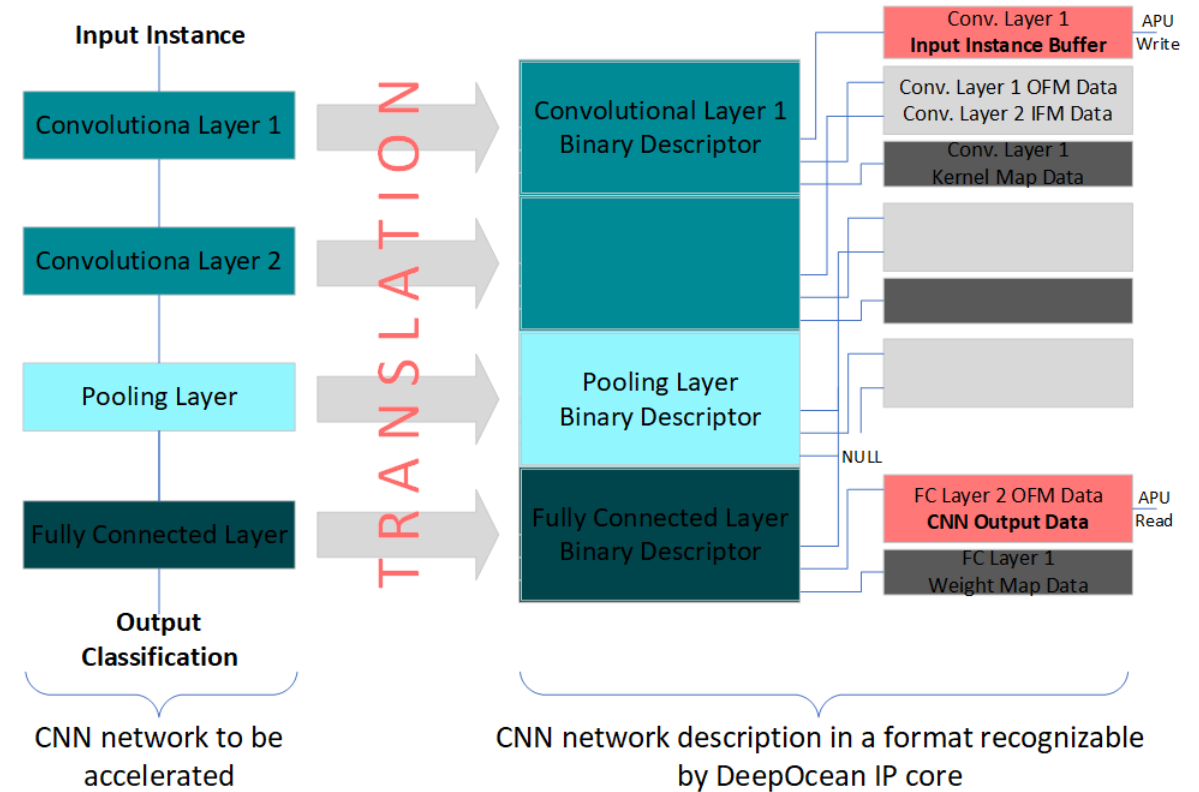


It's so easy!



Preparing CNN to run on „Deep Ocean” accelerator

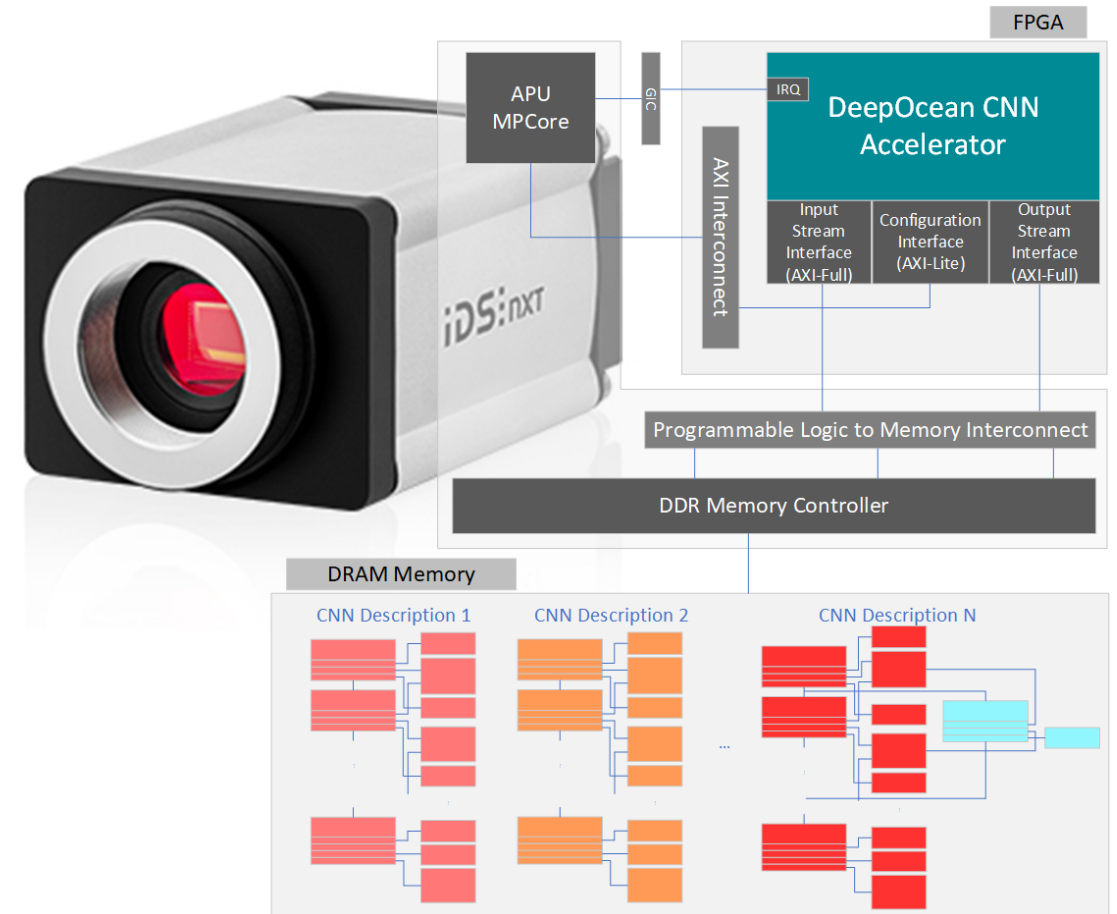
- **Translate** CNN network model into a binary format
- Determine **optimal** fixed-point number representation
- Floating-point to fixed-point **conversion**
- **Compress** sparse kernels





Benefits of using linked list CNN representation

- Extremely **easy to change** CNN network
- **Dynamic Switching** between different CNN networks is **extremely fast**
- Enables changing CNN network topology „**on-the-fly**“



It's so easy!



Advantages of dynamic CNN switching

STEP 1:

Object Localization/Tracking
+ Classification
Run Object Localization CNN

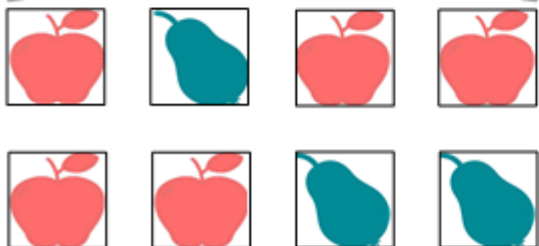
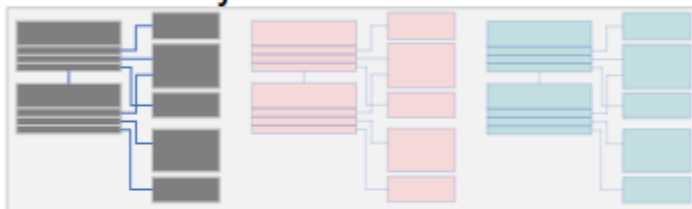
To Difficult to Train!



STEP 1:

Object Localization/Tracking

Run Object Localization CNN

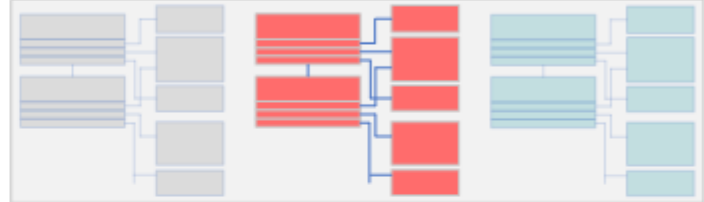


... so easy.

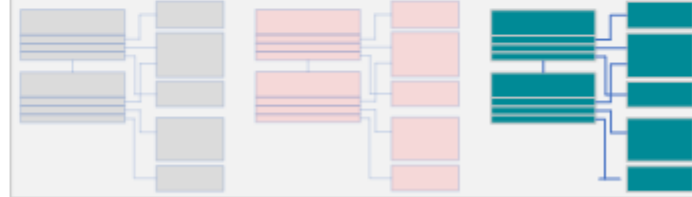
STEP 2:

Conditional Object Classification

Run Classification Network CNN 1



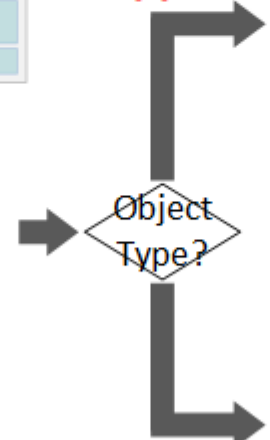
Run Classification Network CNN 2



Apple

Pear

Object Type?





„Deep ocean” inference latency and throughput

- Target FPGA device Xilinx ZCU3
- Operating at 245 MHz
- Using 64 Compute Cores

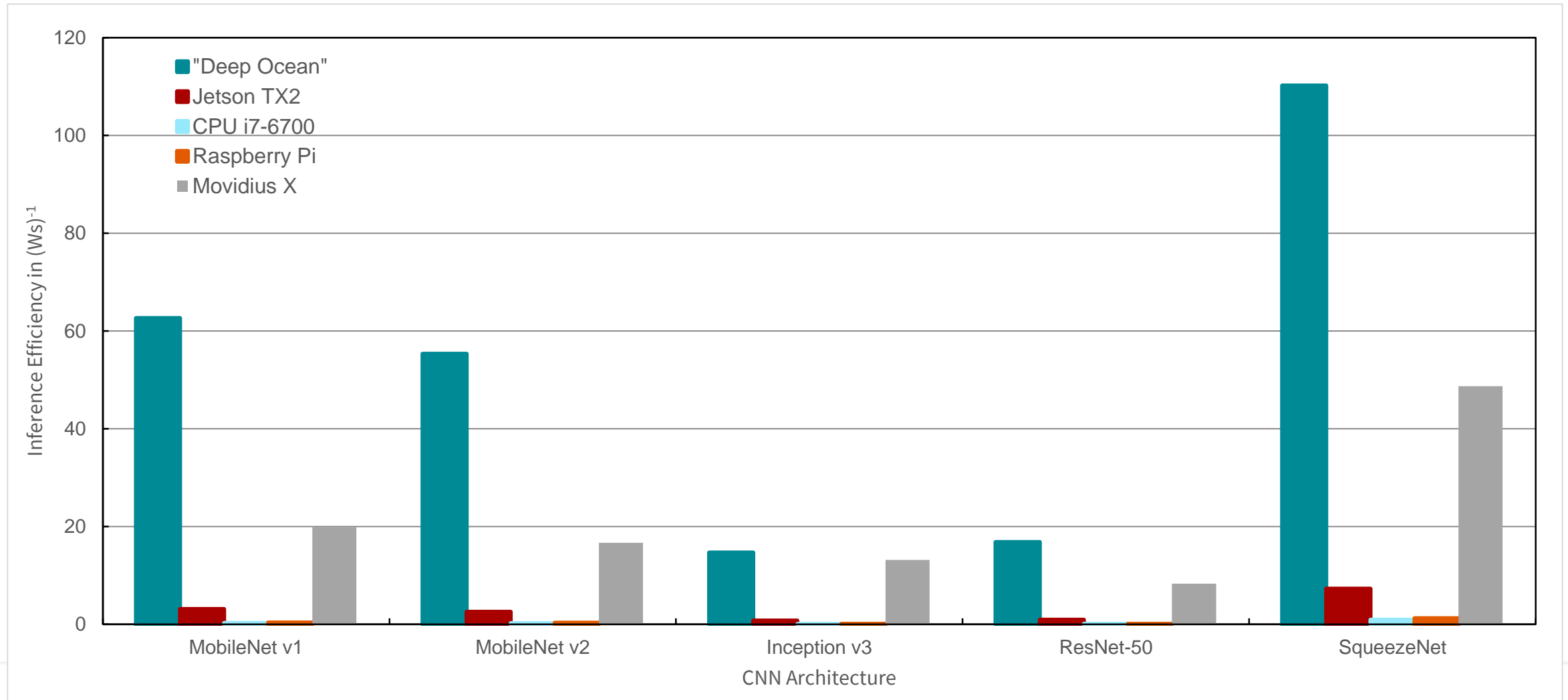
CNN Network	Input Image Size	Latency [ms]	Frame Rate [fps]
MobileNet v1	224x224	26.16	38.23
MobileNet v2	224x224	29.63	33.75
MobileNet v3 large	224x224	34.97	28.60
MobileNet v1 SE	224x224	40.1	24.94
SqueezeNet	224x224	14.87	67.25
EfficientNet B0	224x224	51.79	19.31
MnasNet	224x224	33.43	29.91
Inception v3	299x299	111.62	8.96
ResNet-18	224x224	48.62	20.57
ResNet-34	224x224	61.02	16.39
ResNet-50	224x224	97.39	10.27
MobileNet v1 SSD	300x300	48.20	20.75
MobileNet v2 SSD	300x300	50.65	19.74

Actual performance of IDS NXT camera can differ from numbers presented above

It's so easy!



„Deep ocean” power efficiency





THANK YOU!

IDS Imaging Development Systems GmbH

Dimbacher Str. 6-8, 74182 Obersulm | Tel.: +49 7134 96196-0 | www.ids-imaging.com

It's so easy!



Copyright Notice

This presentation in this publication was presented as a tinyML® Talks webcast. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

www.tinyML.org