

# tinyML<sup>®</sup> Talks

*Enabling Ultra-low Power Machine Learning at the Edge*

## **“Energy-Efficiency and Security for TinyML and EdgeAI: A Cross-Layer Approach”**

Dr. Muhammad Shafique – Professor, New York University Abu Dhabi

February 1, 2022



[www.tinyML.org](http://www.tinyML.org)



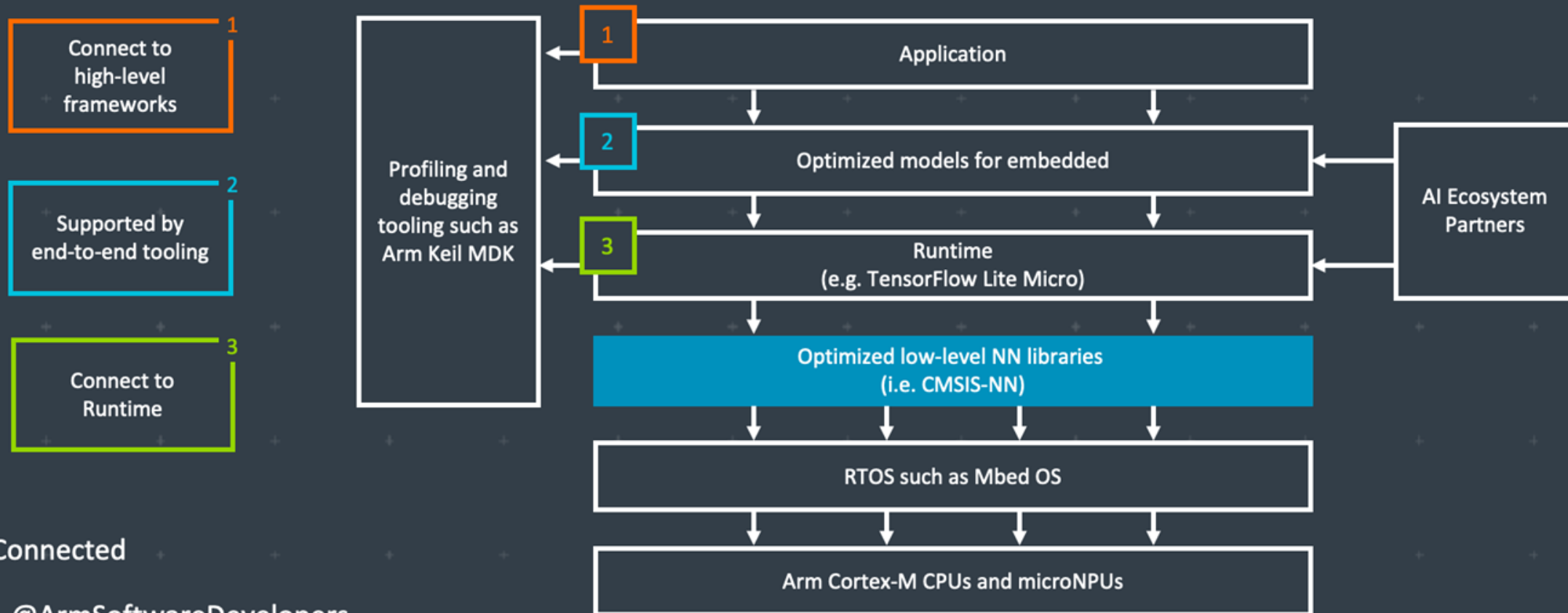
# tinyML Talks Strategic Partners



Additional Sponsorships available – contact [Olga@tinyML.org](mailto:Olga@tinyML.org) for info

# Executive Strategic Partners

# Arm: The Software and Hardware Foundation for tinyML



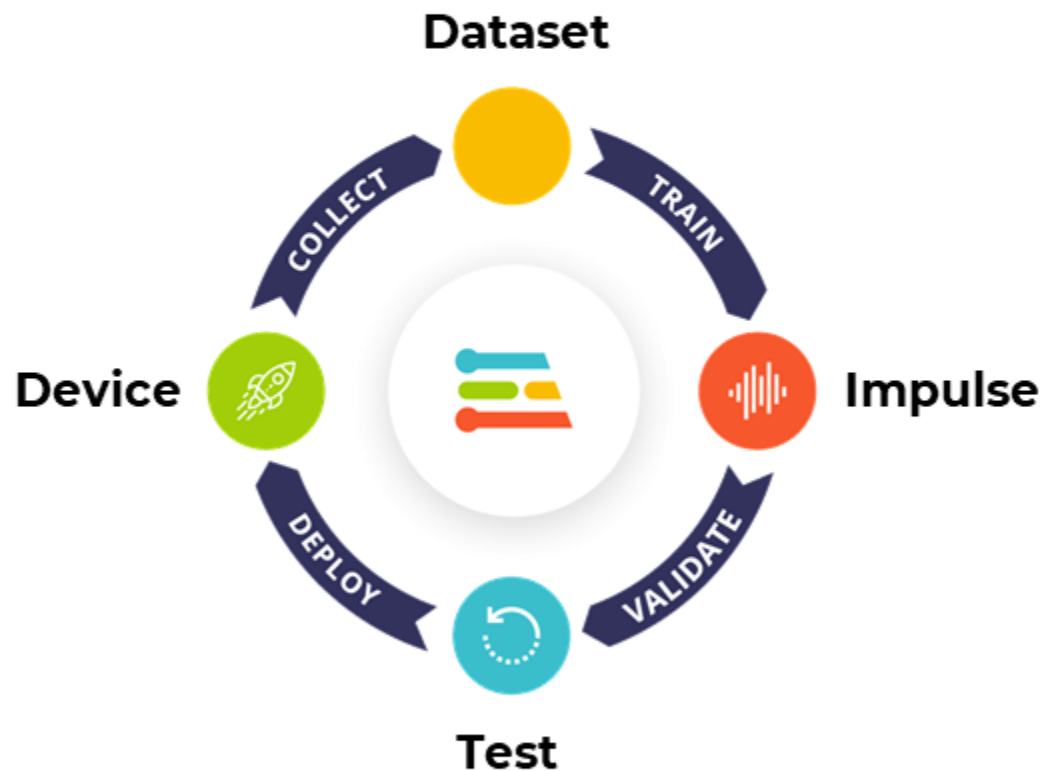
Stay Connected

 @ArmSoftwareDevelopers

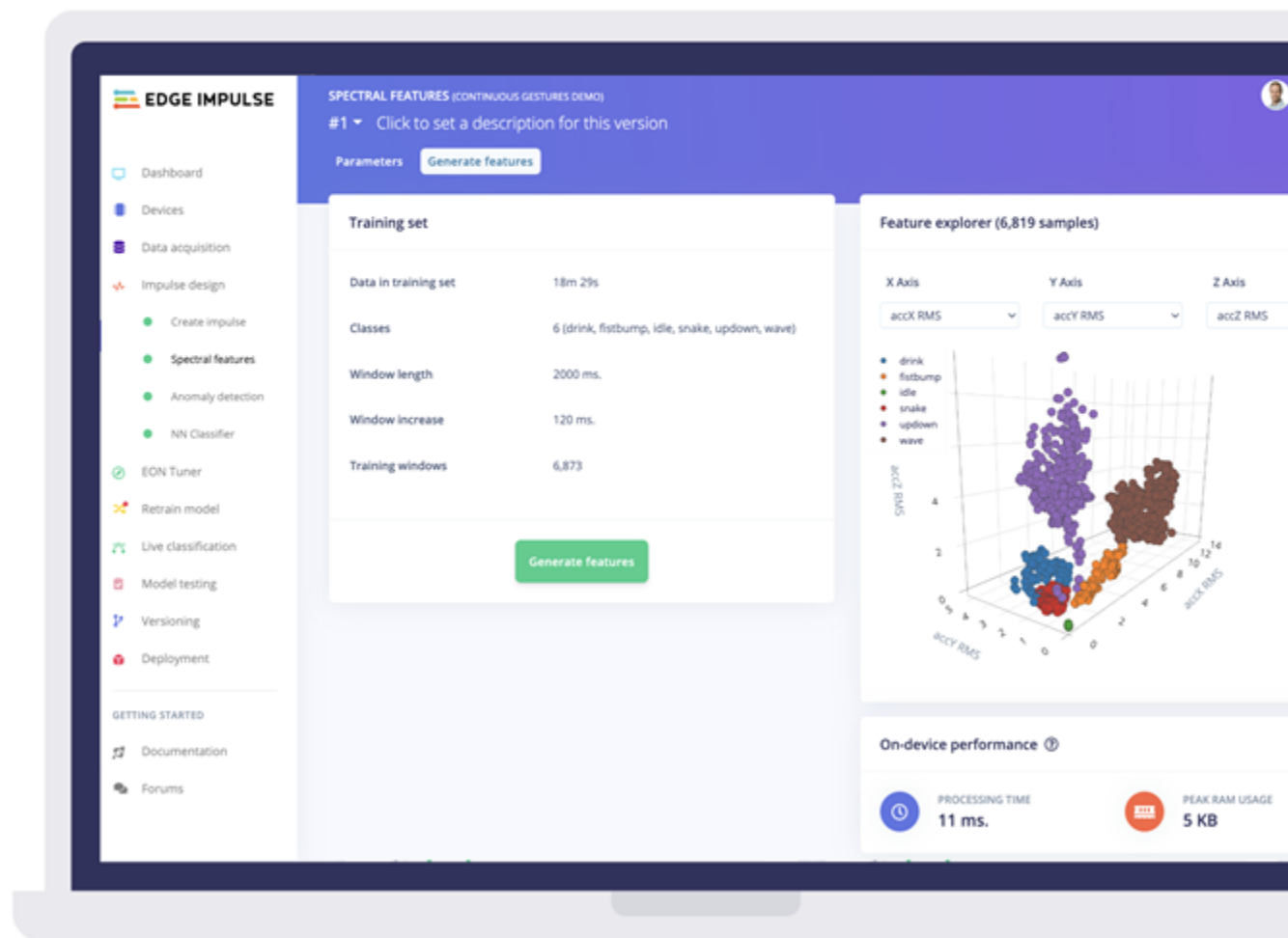
 @ArmSoftwareDev

Resources: [developer.arm.com/solutions/machine-learning-on-arm](https://developer.arm.com/solutions/machine-learning-on-arm)

# EDGE IMPULSE The leading edge ML platform



[www.edgeimpulse.com](http://www.edgeimpulse.com)



**Qualcomm**  
AI research

# Advancing AI research to make efficient AI ubiquitous

## Power efficiency

Model design, compression, quantization, algorithms, efficient hardware, software tool

## Personalization

Continuous learning, contextual, always-on, privacy-preserved, distributed learning

## Efficient learning

Robust learning through minimal data, unsupervised learning, on-device learning

## A platform to scale AI across the industry



### Perception

Object detection, speech recognition, contextual fusion



### Reasoning

Scene understanding, language understanding, behavior prediction



### Action

Reinforcement learning for decision making



Edge cloud



Cloud



IoT/IIoT



Automotive



Mobile

# SYNTIANT

End-to-End  
Deep Learning  
Solutions  
for  
TinyML & Edge AI



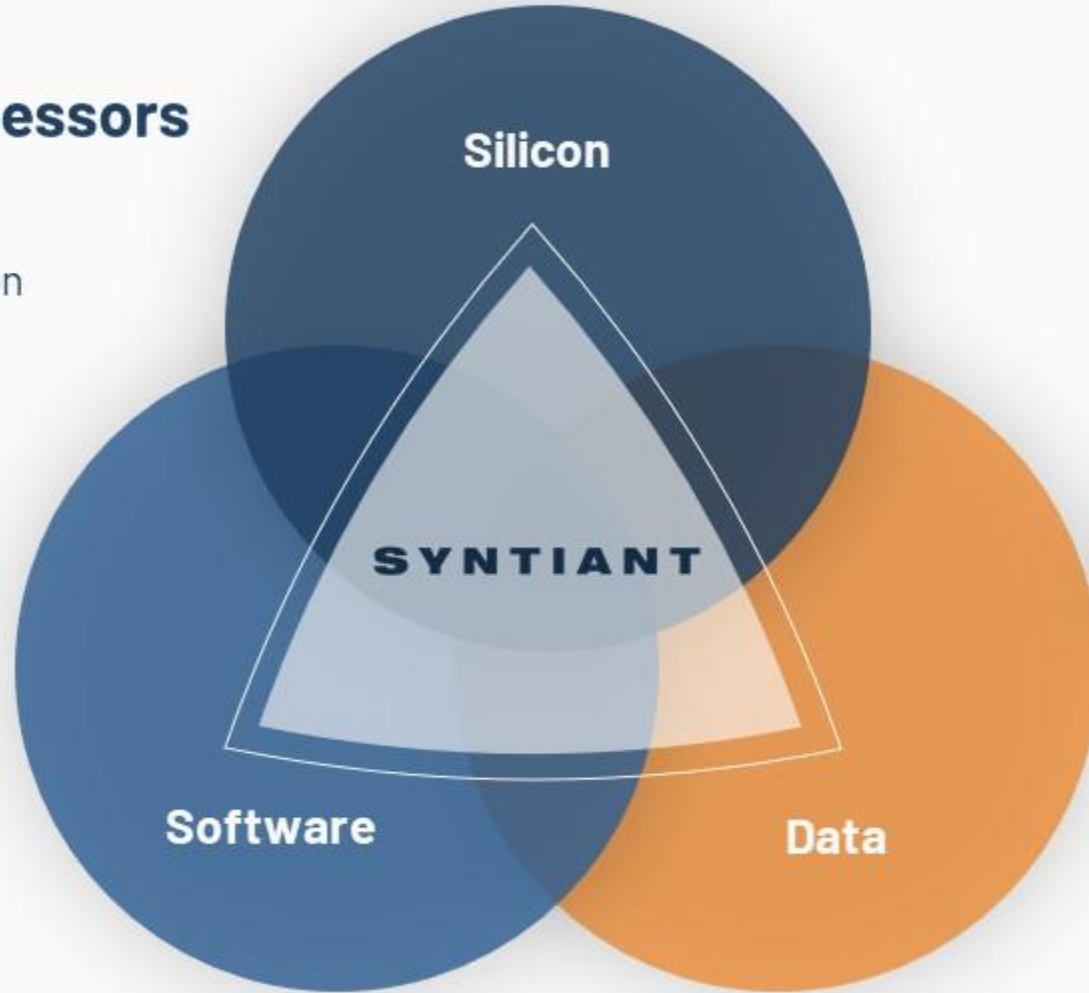
## Neural Decision Processors

- At-Memory Compute
- Sustained High MAC Utilization
- Native Neural Network Processing



## ML Training Pipeline

- Enables Production Quality Deep Learning Deployments



## Data Platform

- Reduces Data Collection Time and Cost
- Increases Model Performance

T I N Y



TALKS  
*webcast*

# Platinum Strategic Partners





# 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](https://bit.ly/testdeeplite)

mobilityXlab

arm





**KLIKA · TECH**

GLOBAL IOT SOLUTIONS



# Reality AI<sup>®</sup>

## Add Advanced Sensing to your Product with Edge AI / TinyML

<https://reality.ai>



[info@reality.ai](mailto: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<sup>®</sup> 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

# BROAD AND SCALABLE EDGE COMPUTING PORTFOLIO

## Microcontrollers & Microprocessors

### Arm® Core



Arm® Cortex®-M 32-bit MCUs  
Arm ecosystem, Advanced security, Intelligent IoT



Arm®-based High-end 32 & 64-bit MPUs  
High-resolution HMI, Industrial network & real-time control



Arm® Cortex®-M0+ Ultra-low Power 32-bit MCUs  
Innovative process tech (SOTB), Energy harvesting

**Renesas Synergy™** Arm®-based 32-bit MCUs for Qualified Platform  
Qualified software and tools

### Renesas Core



Ultra-low Energy 8 & 16-bit MCUs  
Bluetooth® Low Energy, SubGHz, LoRa®-based Solutions



High Power Efficiently 32-bit MCUs  
Motor control, Capacitive touch, Functional safety, GUI



40nm/28nm process Automotive 32-bit MCUs  
Rich functional safety and embedded security features

## Core technologies

### AI

A broad set of high-power and energy-efficient embedded processors

### Security & Safety

Comprehensive technology and support that meet the industry's stringent standards



### Digital & Analog & Power Solution

Winning Combinations that combine our complementary product portfolios

### Cloud Native

Cross-platforms working with partners in different verticals and organizations

T I N Y



TALKS  
*webcast*

# Gold Strategic Partners

T I N Y

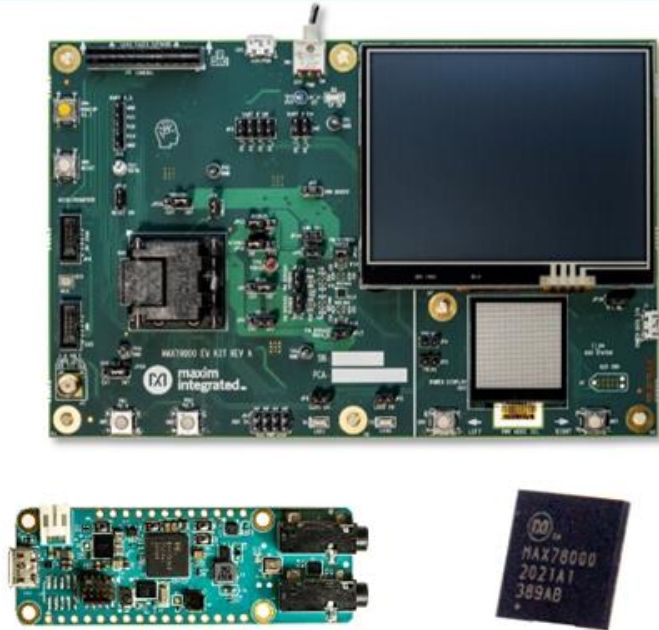


TALKS  
*webcast*



## Maxim Integrated: Enabling Edge Intelligence

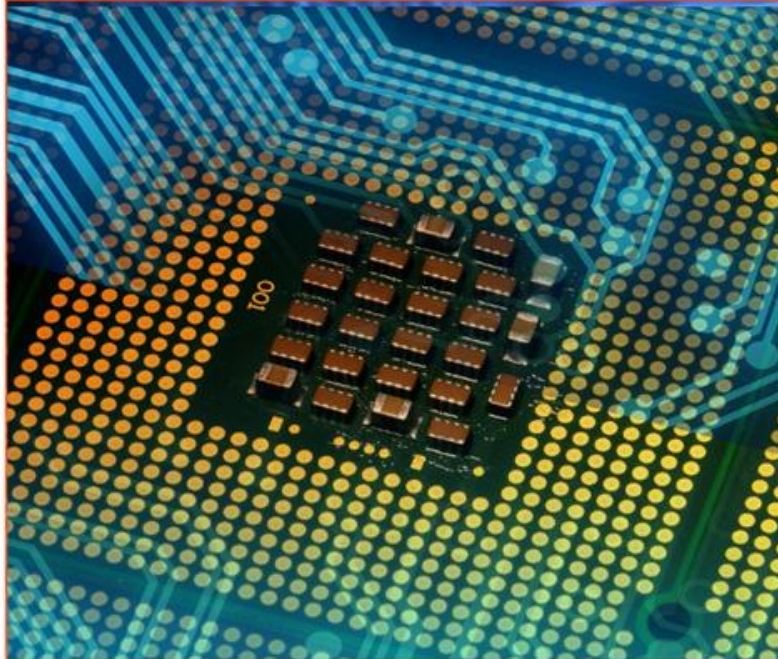
### Advanced AI Acceleration IC



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](http://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](http://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](http://www.maximintegrated.com/sensors)



# Latent AI

Adaptive AI for the Intelligent Edge

[Latentai.com](https://latent.ai)



T I N Y



TALKS  
*webcast*

# Micr .ai

T I N Y



TALKS  
*webcast*

NXP



**seeed** studio

**The IoT Hardware Enabler**



# Build Smart IoT Sensor Devices From Data

SensiML pioneered TinyML software tools that auto generate AI code for the intelligent edge.

- End-to-end AI workflow
- Multi-user auto-labeling of time-series data
- Code transparency and customization at each step in the pipeline

We enable the creation of production-grade smart sensor devices.



[sensiml.com](https://sensiml.com)

T I N Y



TALKS  
*webcast*



life.augmented



# 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>



# Silver Strategic Partners

AONdevices



The logo for Grovety Inc. features a green lightning bolt icon followed by the text "Grovety Inc." in a bold, uppercase sans-serif font.





# tinyML Summit 2022

Miniature dreams can come true...

March 28-30, 2022

Hyatt Regency San Francisco Airport

<https://www.tinymml.org/event/summit-2022/>

*The Best Product of the Year and the Best Innovation of the Year awards are open for nominations between **November 15** and **February 28**.*

# tinyML Research Symposium 2022

March 28, 2022

<https://www.tinymml.org/event/research-symposium-2022>

More sponsorships are available: [sponsorships@tinyML.org](mailto:sponsorships@tinyML.org)





# tinyML Trailblazers Series

Success Stories with Joel Rubino  
(CEO & Co-founder of Cartesium)

**LIVE ONLINE February 2nd, 2022 at 8 am PST**



Register now!





# Join Growing tinyML Communities:



8k members in  
42 Groups in 33 Countries

**tinyML - Enabling ultra-low Power ML at the Edge**

<https://www.meetup.com/tinyML-Enabling-ultra-low-Power-ML-at-the-Edge/>



2.6k members  
&  
4.6k followers

**The tinyML Community**

<https://www.linkedin.com/groups/13694488/>





Subscribe to  
**tinyML YouTube Channel**  
for updates and notifications  
*(including this video)*  
[www.youtube.com/tinyML](https://www.youtube.com/tinyML)



**tinyML**  
4.33K subscribers

**5.9k subscribers, 347 videos with 174k views**

HOME VIDEOS PLAYLISTS COMMUNITY CHANNELS ABOUT

tinyML Summit 2021 Partner Session: Productio... 146 views • 1 month ago

tinyML Summit 2021 Partner Session: System... 131 views • 1 month ago

tinyML Talks Massimo Banzi: tiny machine learnin... 1.4K views • 1 month ago

tinyML Summit 2021 Partner Session: It's an SN... 420 views • 1 month ago

tinyML Summit 2021 Partner Session: The... 104 views • 1 month ago

tinyML Summit 2021 Partner Session: Innovativ... 110 views • 1 month ago

Deploying AI to Embedded Systems 430 views • 1 month ago

tinyML Summit 2021 Breaking News on... 472 views • 1 month ago

tinyML Summit 2021 Partner Session: TinyML... 160 views • 1 month ago

tinyML Summit 2021 Breaking News on... 257 views • 1 month ago

tinyML Summit 2021 Partner Session on... 299 views • 1 month ago

tinyML Summit 2021 Breaking News on... 173 views • 1 month ago

tinyML Summit 2021 Partner Session: Machine... 212 views • 1 month ago

tinyML Summit 2021 Partner Session: Tiny and... 109 views • 1 month ago

tinyML Summit 2021 tiny Talks: Real-World... 55 views • 1 month ago

tinyML Summit 2021 tiny Talks: Environmental Nois... 173 views • 1 month ago

tinyML Summit 2021 Breaking News on... 151 views • 1 month ago

tinyML Summit 2021 tiny Talks: Insights from a Mult... 222 views • 1 month ago

tinyML Summit 2021 tiny Talks: TinyML Software... 121 views • 1 month ago

tinyML Summit 2021 Keynote: Data-Free Model... 240 views • 1 month ago

tinyML Summit Partner Session: Pushing the AI... 82 views • 1 month ago

tinyML Summit 2021 Partner Session: Low pow... 129 views • 1 month ago

tinyML Summit 2021 Partner Session: How... 95 views • 1 month ago

tinyML Summit 2021 Partner Session: Low-pow... 168 views • 1 month ago

tinyML Summit 2021 Partner Session: TinyML is... 115 views • 1 month ago

tinyML Summit 2021 Keynote: Efficient Audio... 1.4K views • 1 month ago

tinyML Summit 2021 tiny Talks: An Introduction to a... 180 views • 1 month ago

tinyML Summit 2021 tiny Talks: Hardware Aware... 215 views • 1 month ago

tinyML Summit 2021 tiny Talks: Hardware Aware... 75 views • 1 month ago

tinyML Summit 2021 tiny Talks: Neutrino: A BlackBo... 168 views • 1 month ago

tinyML Summit 2021 Panel Discussion: tinyML... 1:01:15

tinyML Summit 2021 tiny Talks: Person Detection... 18:26

tinyML Summit 2021 tiny Talks: Using Neural... 19:03

tinyML Summit 2021 Keynote: Adaptive Neural... 55:15

tinyML Summit 2021 Keynote: milliJoules for... 99:43

tinyML Summit 2021 Market Opportunities for Edge AI 51:28



# Next tinyML Talks

Date	Presenter	Topic / Title
Tuesday, February 8	Stefano Cadario, Director Product Management, IoT Group, Arm	Get Ahead of the Curve: Develop Software in the Cloud for the Ethos-U55 and Cortex-M55 Processors

Webcast start time is 8:00 am Pacific time

Please contact [talks@tinymml.org](mailto:talks@tinymml.org) if you are interested in presenting



# Reminders

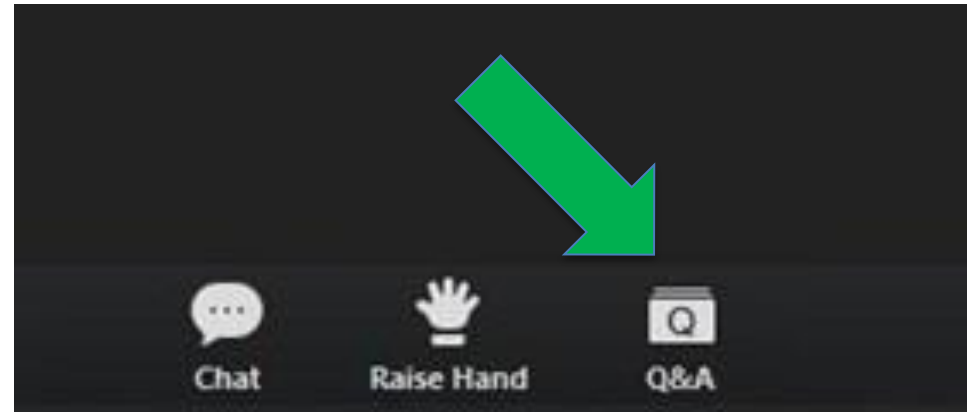
Slides & Videos will be posted tomorrow

Please use the Q&A window for your questions



[tinyml.org/forums](https://tinyml.org/forums)

[youtube.com/tinyml](https://youtube.com/tinyml)





# Dr. Muhammad Shafique



M. Shafique is an Associate Professor in the Division of Engineering, New York University (NYU) Abu Dhabi, UAE, and Global Network Associate Professor in the Tandon School of Engineering, NYU-NY, USA. He is also a CoPI / Investigator in multiple Centers, i.e., Center of AI and Robotics, Center of Quantum Computing, Center of Cyber Security, and Center for InTeraCTIng urban nEtworkS.

He received his Ph.D. in Computer Science from Karlsruhe Institute of Technology (KIT), Germany in 2011. From Sep.2016 to Aug.2020, he was a Full Professor of Computer Architecture and Robust Energy-Efficient Technologies (CARE-Tech.) at the Institute of Computer Engineering, Vienna University of Technology (TU Wien).

Dr. Shafique has received ACM SIGDA Outstanding New Faculty Award, AI-2000 Most Influential Scholar Award in 2020, ASPIRE Award for Research Excellence, and multiple best paper awards and nominations at flagship conferences.

# Who Ruled the World!

## Age of Power

**Man-Power (#), Skills, Strength, Courage, etc.**



## Age of Resources and Industry

**Fuel, Industrial Tech., Economic Politics, etc.**



## Age of Data and AI

*Data is the New Fuel*

**Innovation in Technology is the New Politics  
Nation-wide Race for Dominance in AI**



# Outline

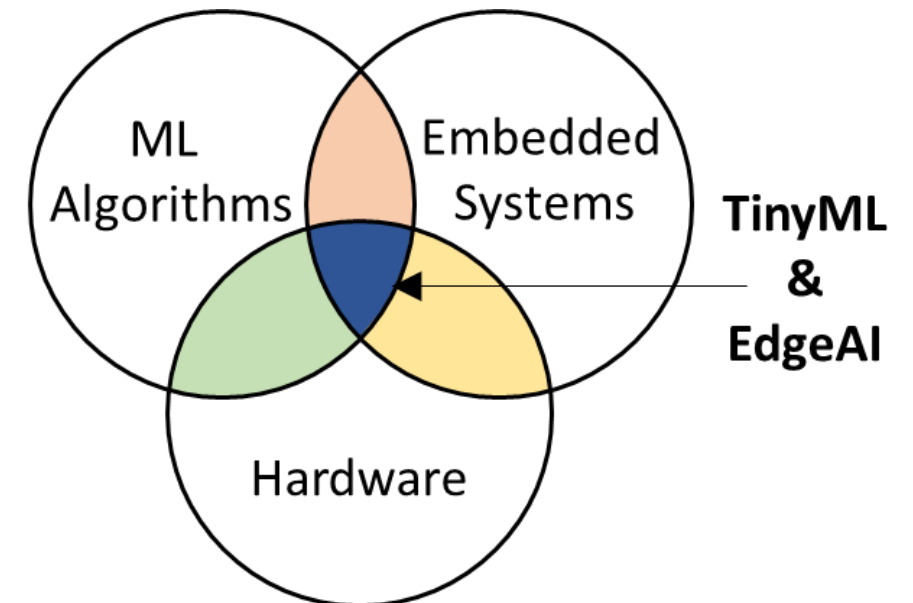
- ❑ What are TinyML and EdgeAI?
- ❑ Applications
- ❑ Cross-Layer Design Flow
- ❑ Future Research Directions



# TinyML and EdgeAI: Unique Features?



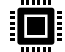
*Enabling on-device data analytics, predictions, & intelligence at extremely low power*

- Fastest-growing field of machine learning
- Combination of embedded systems, algorithms and hardware
- On-device ML under limited resources
- Stringent design constraints
- Always-on use-cases
- Battery-operated devices
- Scalable to trillions of sensors

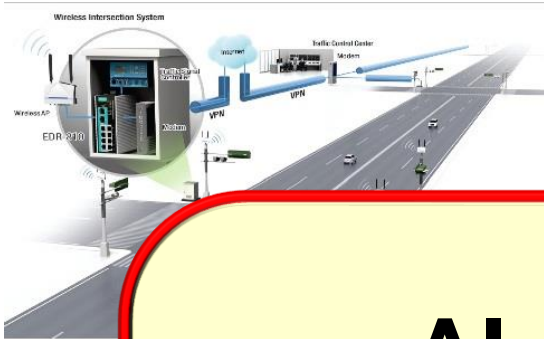


# TinyML and EdgeAI

□ Fundamentally different from machine learning in the cloud

	Cloud AI 	➔ Edge/Mobile AI 	➔ TinyML 
Hardware	NVIDIA DGX A100	S21, iPhone 13, NVIDIA Jetson	STM32F769 Microcontroller
Memory	1 TB System Memory + 320 GB GPU Memory	2 - 12 GB	~512 KB
Storage	>15 TB	16 - 512 GB	~2 MB
Applications	Model Training, Big Data Analytics	Embedded processing Continual Learning	In-/near-sensor processing
		<b>Tight Constraints</b>	<b>Extreme Constraints</b>

# Smart Cyber Physical Systems & Internet-of-Things



AI / ML is inevitable, we have to efficiently **infer knowledge** from the big data, and **derive predictions**

<https://>

Systems  
[analysis/automotive-driving-revolution/](https://analysis/automotive-driving-revolution/)



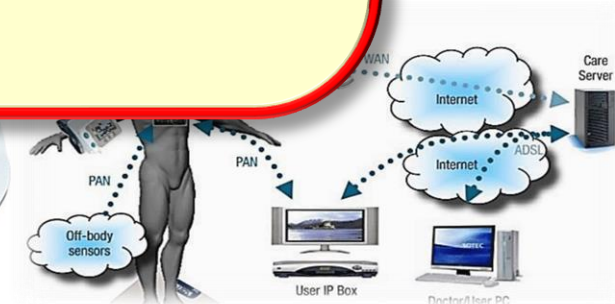
CP Factory  
Wireless communication via RFID, NFC and WLAN  
**Industry 4.0: Smart Industrial Automation**  
<https://vimeo.com/145877805>



**Smart Houses**  
<https://www.linkedin.com/pulse/smart-homes-private-secure-future-intelligent-home-tripti-jha>



**Smart Grids**  
[http://solutions.3m.com/wps/portal/3M/en\\_EU/SmartGrid/EU-Smart-Grid/](http://solutions.3m.com/wps/portal/3M/en_EU/SmartGrid/EU-Smart-Grid/)



**Smart Health Care**

# Smart CPS & IoT => The Robustness Challenge!

... should consider

**Robustness**

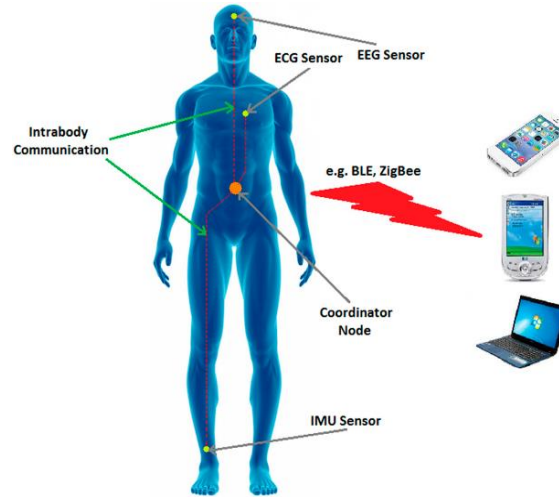
- Reliability
- Security

**Performance**

- Throughput
- Latency

**Others**

- Adaptability
- Safety
- Privacy
- Interoperability



**Smart Healthcare**  
(Energy and time constraints)



**Norwegian C-130 crash (2012)**

[https://en.wikipedia.org/wiki/2012\\_Norwegian\\_C-130\\_crash](https://en.wikipedia.org/wiki/2012_Norwegian_C-130_crash)



**Failure of F-22 Raptor (2007)**

<http://www.dailytech.com/Lockheeds+F22+Raptor+Gets+Zapped+by+International+Date+Line/article6225.htm>



Satellite imagery of the Northeastern United States taken before and during the blackout



Toronto, on the evening of August 14, 2003

**Northeast blackout of 2003**

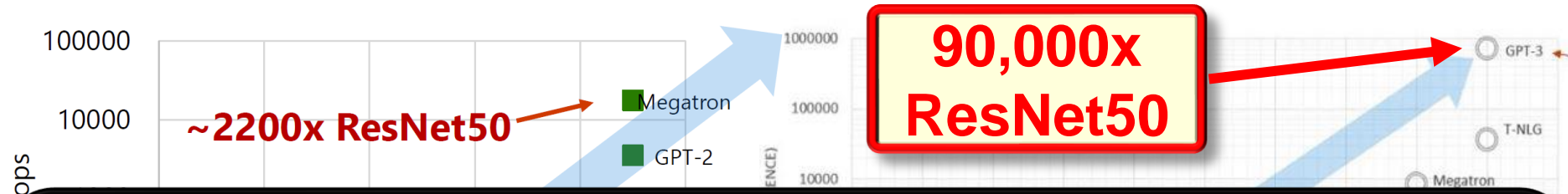
[https://en.wikipedia.org/wiki/Northeast\\_blackout\\_of\\_2003](https://en.wikipedia.org/wiki/Northeast_blackout_of_2003)

**Hacking Jeep Cherokee 4x4 (2015)**

Sent the instructions through Entertainment systems

- Control the steering <https://www.ophtek.com/4-real-life-examples-iot-hacked/>
- Control the braking system

# Complexity: Exponential Growth in Model Sizes!



**Human Brain => 20W**  
**Efficiency Gap => 1,000x → 100,000x!!!**

Source: Eric Chung, "Accelerating Microsoft's AI Ambitions", Microsoft, Azure AI and Advanced Architectures Group, 2019.

Source: <https://www.microsoft.com/en-us/research/blog/a-microsoft-custom-data-type-for-efficient-inference/>.

## Challenging Question

How to process **huge amount of data** in **robust & energy-efficient** way, while considering tinyML / EdgeAI constraints?

# Robustness for Machine Learning: News Feed



## Beware: Galaxy S10's Facial Recognition Easily Fooled with a Photo

Jesus Diaz · Freelance Writer  
 Updated Mar 11, 2019

## Self-driving Uber kills Arizona woman in first fatal crash involving pedestrian

Tempe police said car was in autonomous mode at the time of the crash and that the vehicle hit a woman who later died at a hospital



### Hackers trick a Tesla into veering into the wrong lane

<https://www.youtube.com/watch?v=a7L51u23YoM>



GOOGLE SELF DRIVING CAR CRASHES INTO A BUS

Tesla Model 3: Autopilot engaged during fatal crash

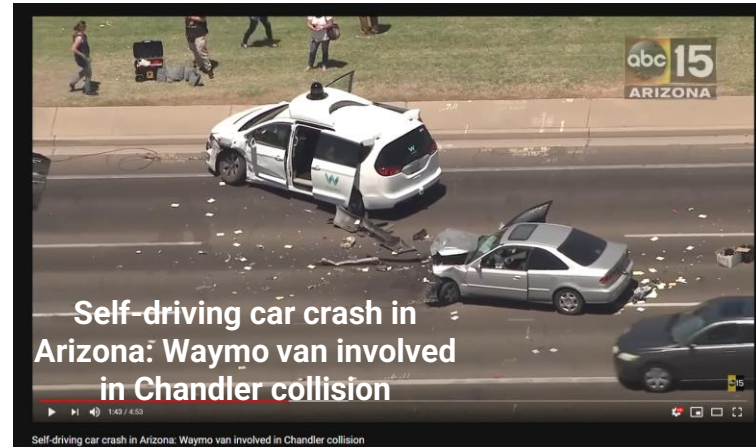
© 17 May 2019

BBC



## Tesla driver dies in first fatal crash while using autopilot mode

The autopilot sensors on the Model S failed to distinguish a white tractor-trailer crossing the highway against a bright sky



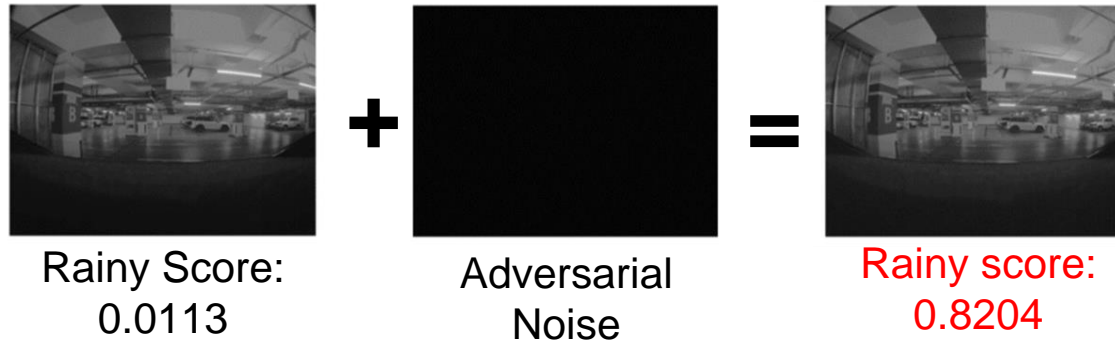
## Self-driving car crash in Arizona: Waymo van involved in Chandler collision

<https://www.technologyreview.com/f/613254/hackers-trick-teslas-autopilot-into-veering-towards-oncoming-traffic/>

# Adversarial Attacks on **Tesla Autopilot** by Tencent Keen Security Lab

## Digital Adversarial Examples

- ❑ Insert the noise into the DNN input

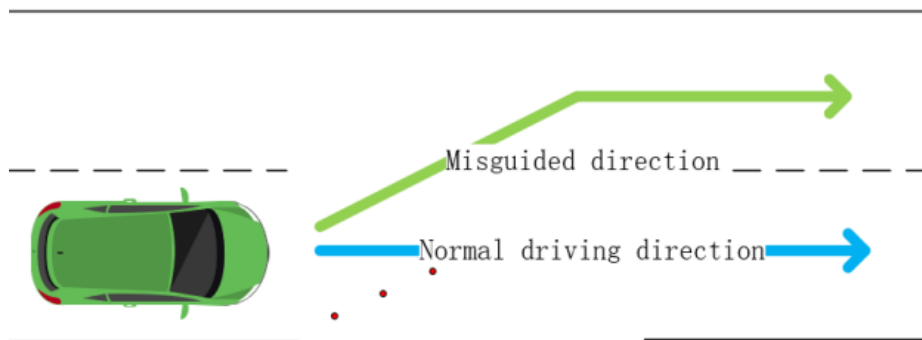


## Black-Box Attack

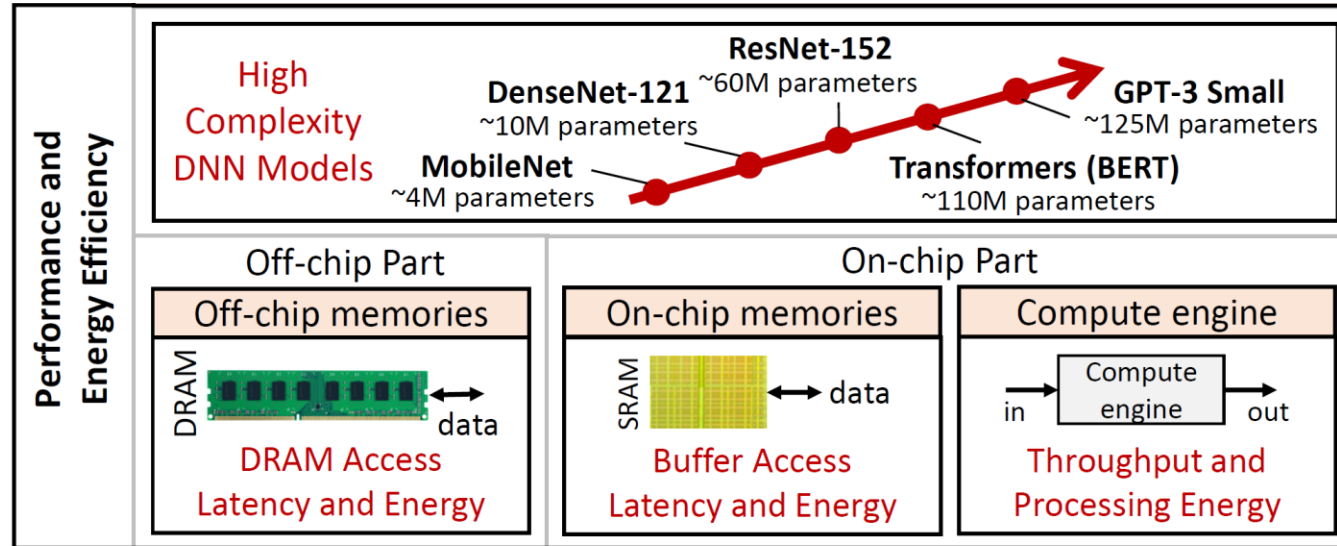


## Physical World Adversarial Examples

- ❑ Place the small stickers on the ground



# Overview of Challenges for EdgeAI & tinyML





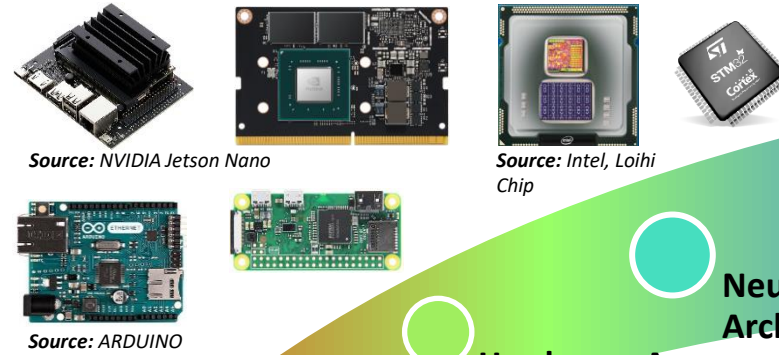
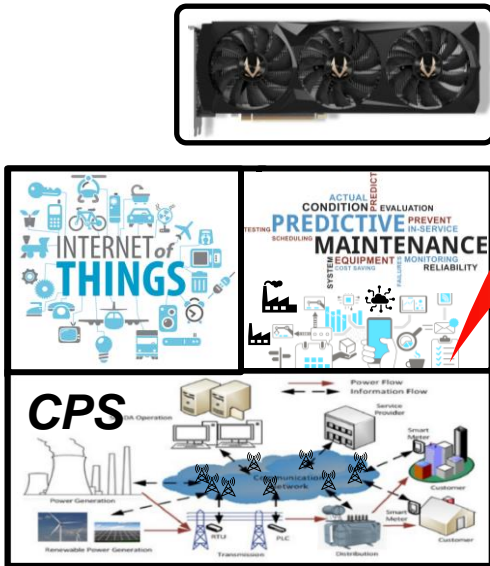
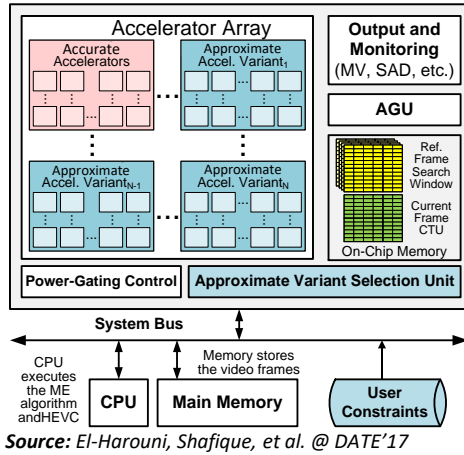
# Cross-Layer Design Flow

- ❑ Frameworks enable seamless integration of algorithms and optimizations at all layers, developed by the community.
  - ❑ Design and optimize ML models for ultra-low power devices



- ❑ Hardware accelerators
  - ❑ Specialized hardware for accelerating vector/matrix multiplication
- ❑ DNN Optimization
  - ❑ Neural Architecture Search (NAS), Pruning and Quantization

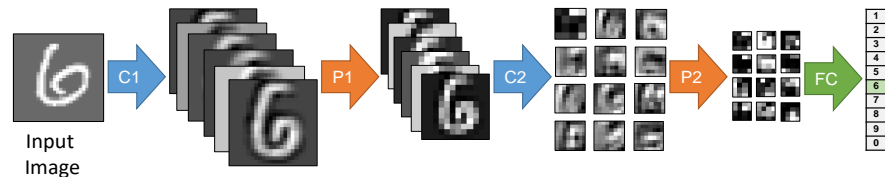
# Embedded AI @ eBrain Lab: A Multi-Dimensional Research Challenge



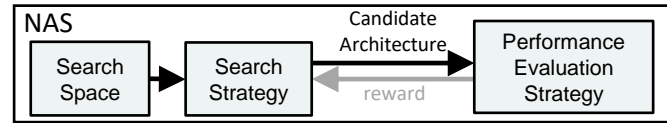
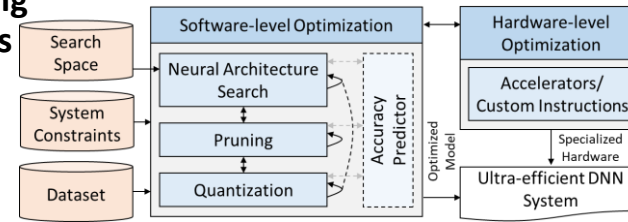
**Software (GPUs)**

**Accelerators + Approximate Computing**

**Deep Learning Architectures**



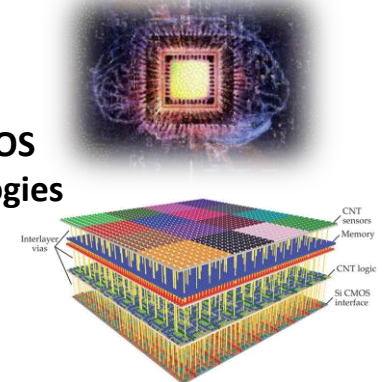
**Hardware-Aware Neural Architecture Search (NAS) + Optimization**



**Neuromorphic Architectures**

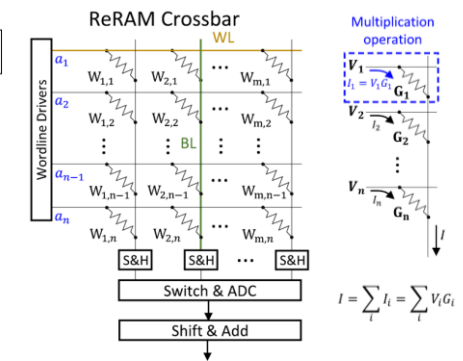
**Post-CMOS Technologies**

**TinyML & EdgeAI**



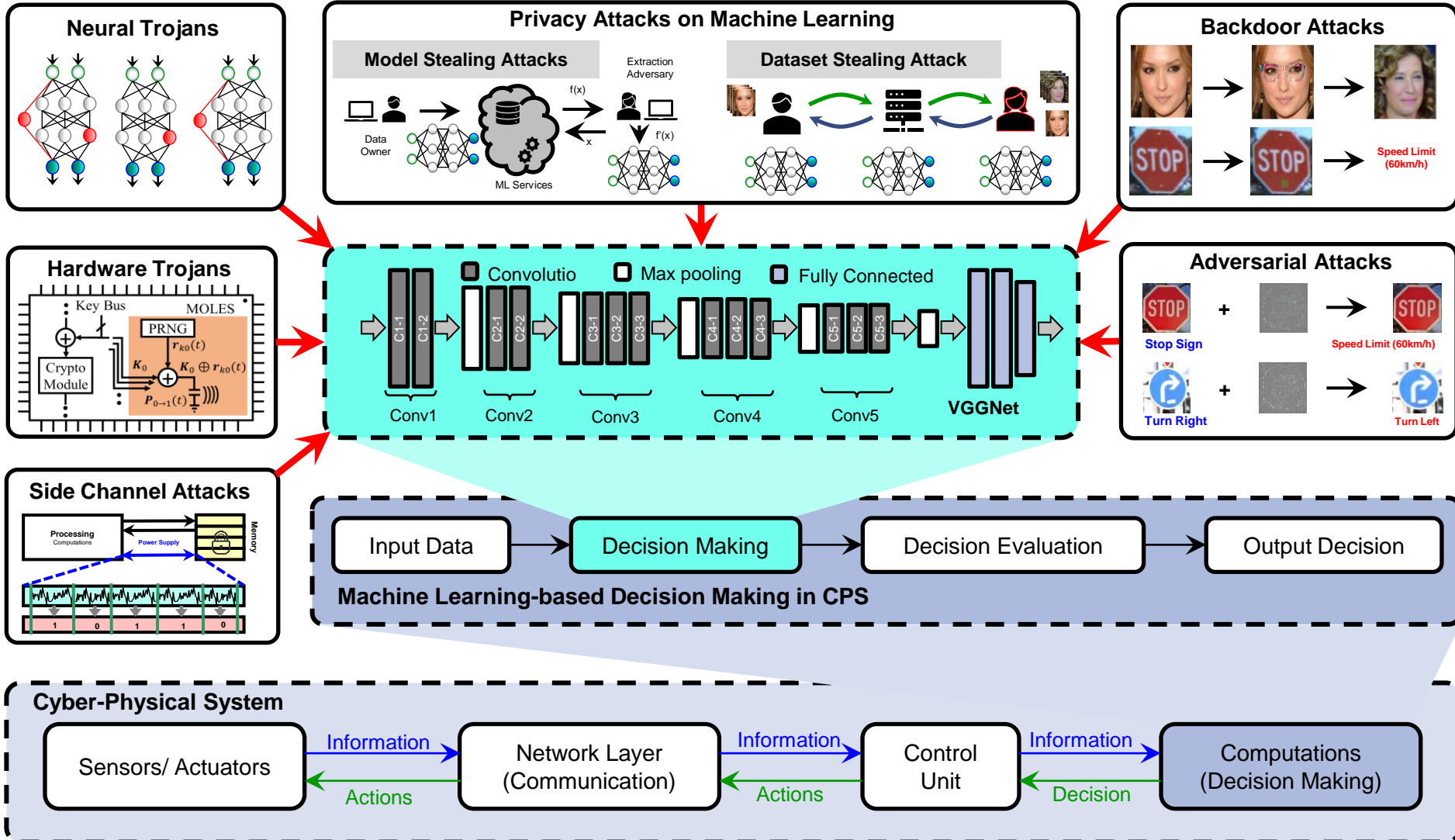
Source: M. M. Shulaker et al., Nature 547, 74 (2017) and R. Mark Wilson; Physics Today 70, 14-16 (2017)

**In-Memory Computing**



Source: Hanif, Manglik, Shafique @ IEEE Access'20

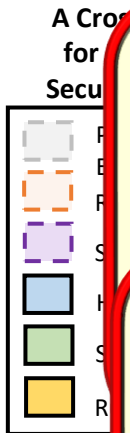
# ML Security Research @ eBrain Lab



- M. A. Hanif, F. Khalid, R. V. W. Putra, S. Rehman, M. Shafique, "Robust Machine Learning Systems: Reliability and Security for Deep Neural Networks", in IOLTS-2018, Platja d'Aro, Spain, pp. 257 - 260.
- F. Kriebel, S. Rehman, M. A. Hanif, F. Khalid, M. Shafique, "Robustness for Smart Cyber-Physical Systems and Internet-of-Things: From Adaptive Robustness Methods to Reliability and Security for Machine Learning", ISVLSI-2018, Hong Kong, China, pp. 581-586.



# Our Cross-Layer TinyML and Edge AI Framework: *An Overview*



## Class-Blind Pruning (IJCNN'19)

**190x – 15x memory savings**

## DRAM Access Energy Savings (TVLSI'21)

**~45% for AlexNet, VGG-16, MobileNet, and SqueezeNet**

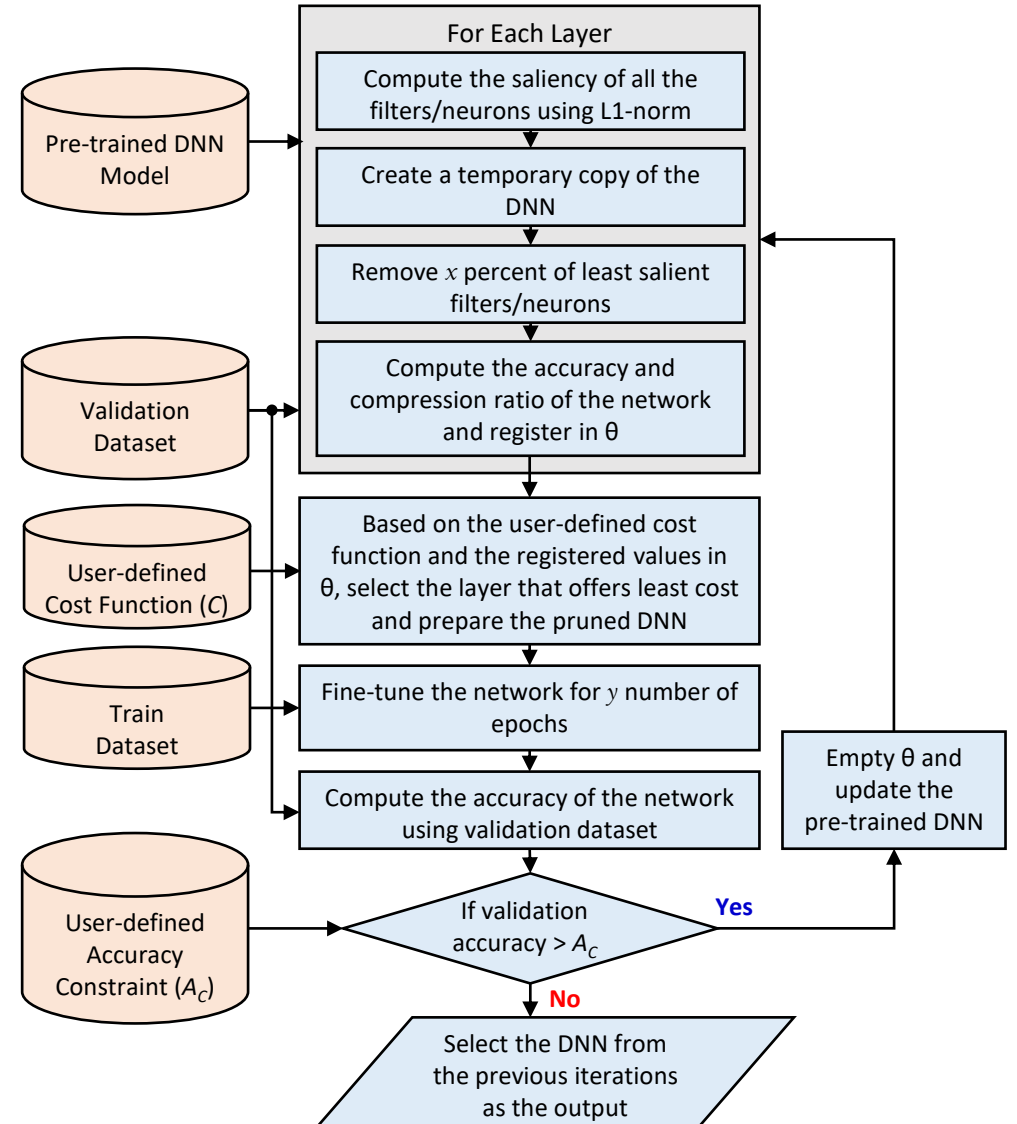
## able imations

## (DAC'19)

**1.5x Energy Efficiency  
@ *NO Accuracy Loss***

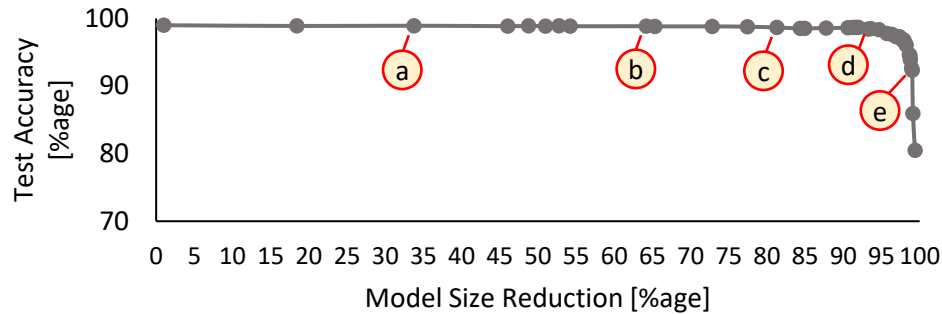
# Structured Pruning Methodology

- ❑ Step 1: **Compute the sensitivity** of the layers of the given DNN to pruning using a user-defined cost function
- ❑ Step 2: **Remove  $x$  percent filters/neurons** from the least sensitive layer
- ❑ Step 3: **Fine-tune the network** for  $y$  number of epochs
- ❑ Step 4: **Compare the accuracy** with the defined accuracy constraint
- ❑ Step 5: Continue pruning if the accuracy is greater than the defined constraint

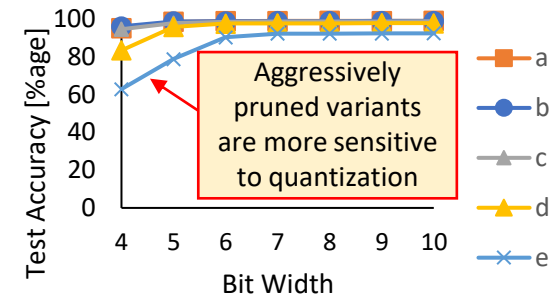


# Results using LeNet-5 trained with MNIST Dataset

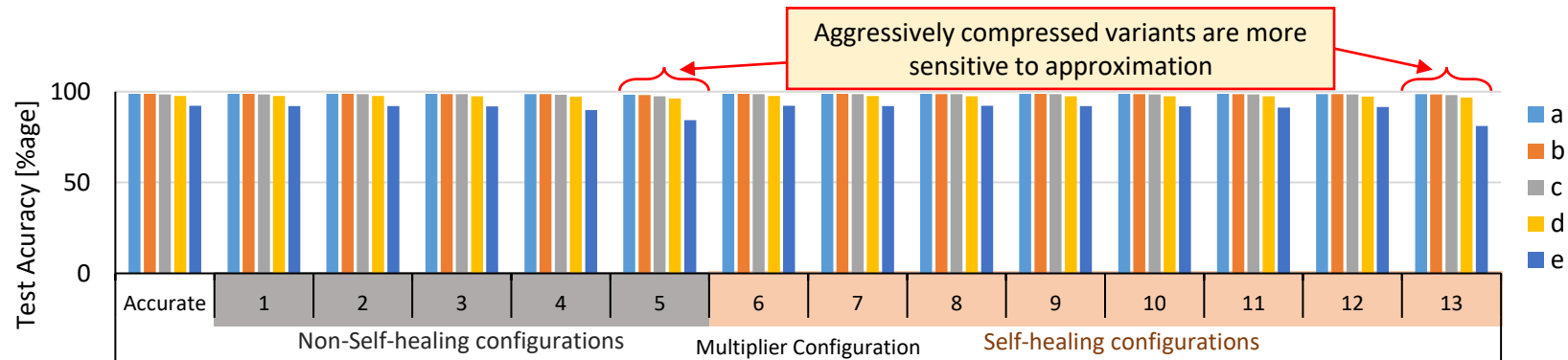
## 1. Structured Pruning



## 2. Quantization

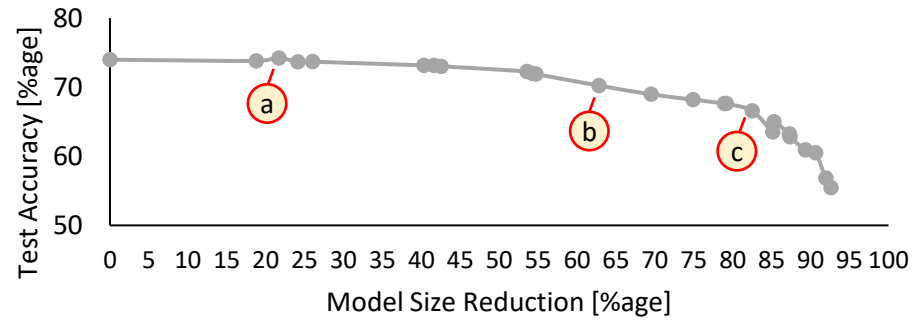


## 3. Hardware Approximation

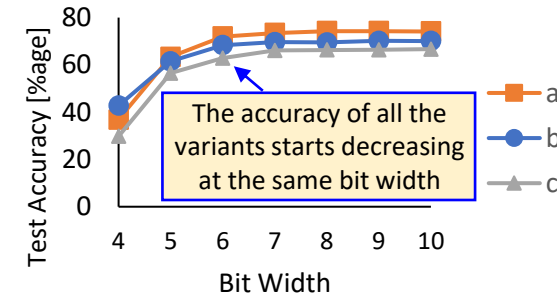


# Results using LeNet-5 trained with Cifar10 Dataset

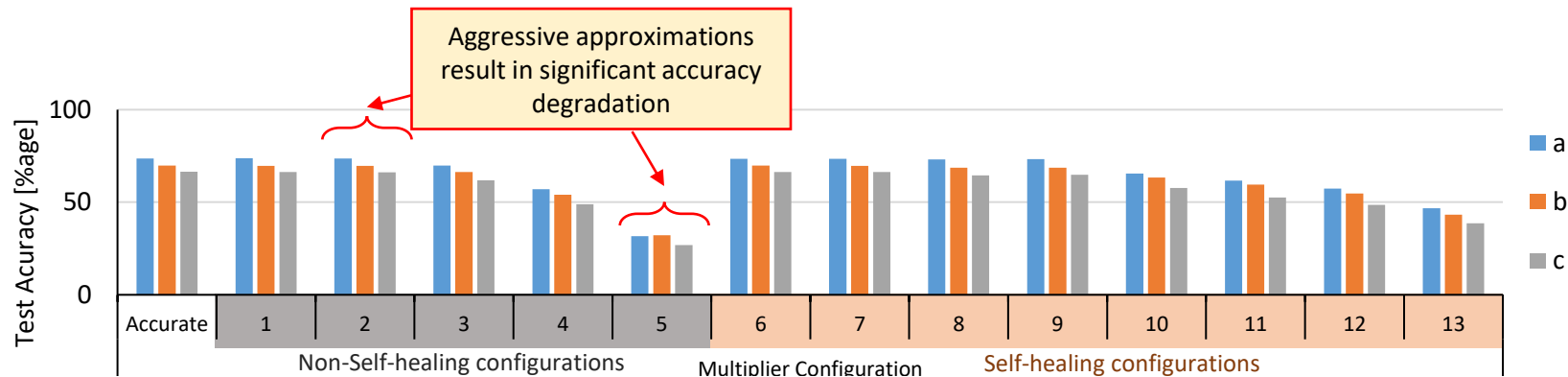
## 1. Structured Pruning



## 2. Quantization

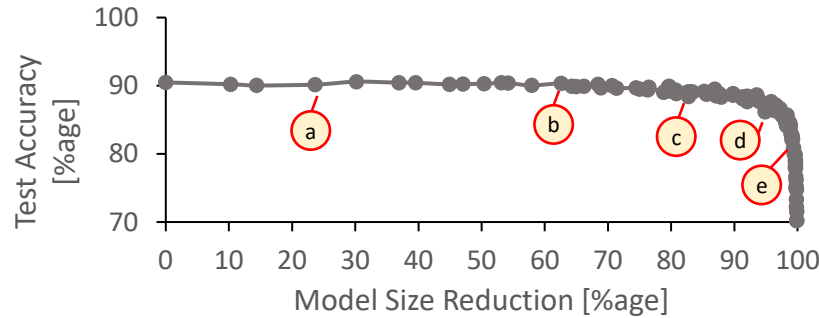


## 3. Hardware Approximation

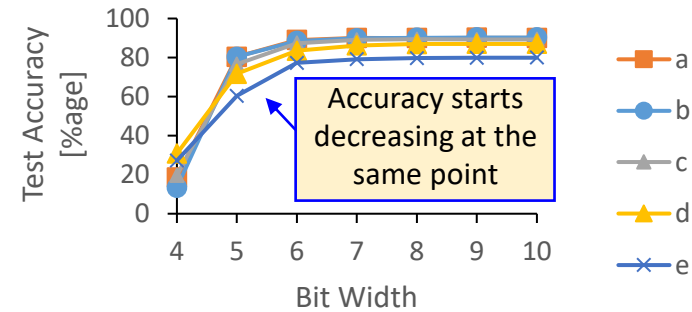


# Results using VGG11 trained with Cifar10 Dataset

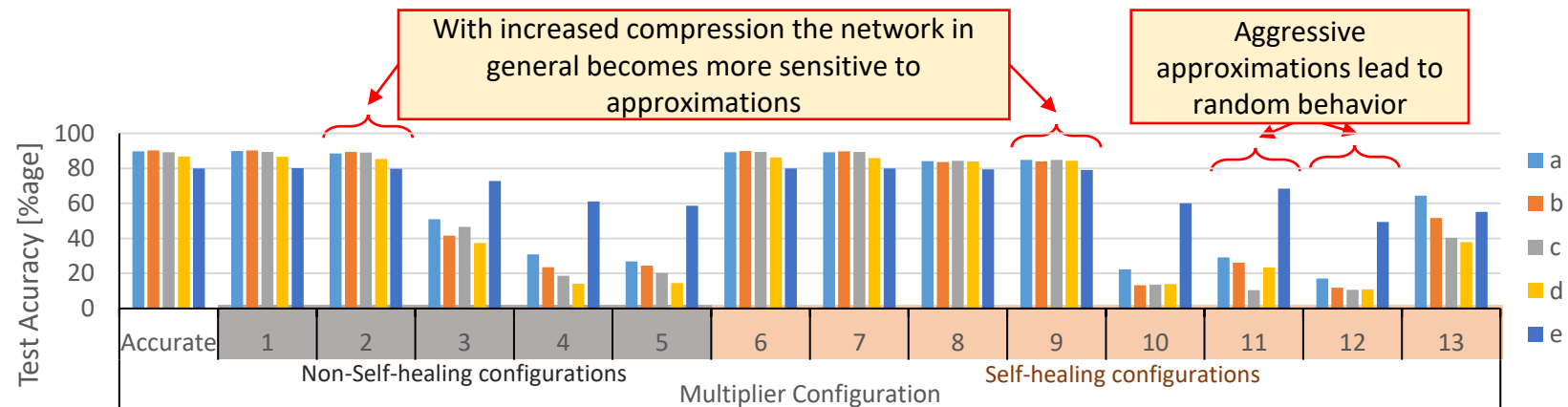
## 1. Structured Pruning



## 2. Quantization



## 3. Hardware Approximation

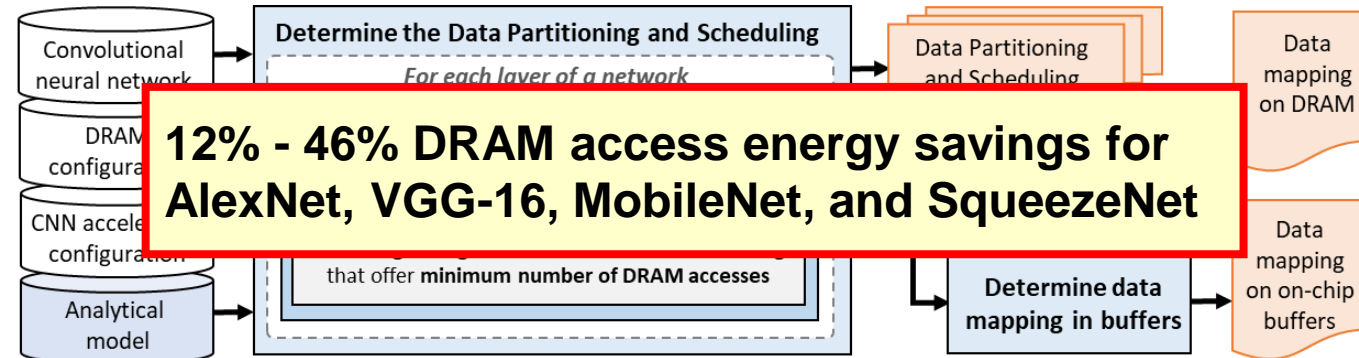




# Memory Optimizations

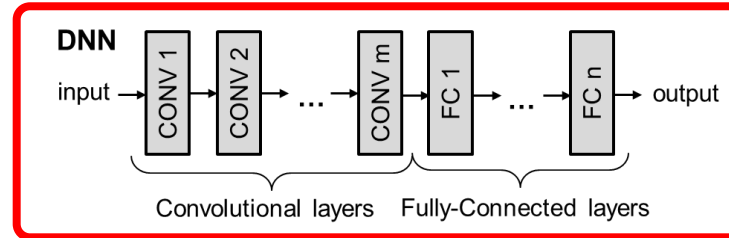
## Energy-Efficient Memory Accesses for DNN Accelerators (IEEE TVLSI'21)

1



**12% - 46% DRAM access energy savings for AlexNet, VGG-16, MobileNet, and SqueezeNet**

2



## Generic DRAM Mapping for Energy-Efficient DNNs (DAC'20)

*Our Novel Contributions*

- DRMap: A Generic DRAM Mapping
- Design Space
- Analytical Model of the EDPs of DRAM Mapping Policies

**Pseudo-code of DRMap**

```
for (ch = 0; ch < # channels; ch++) {
```

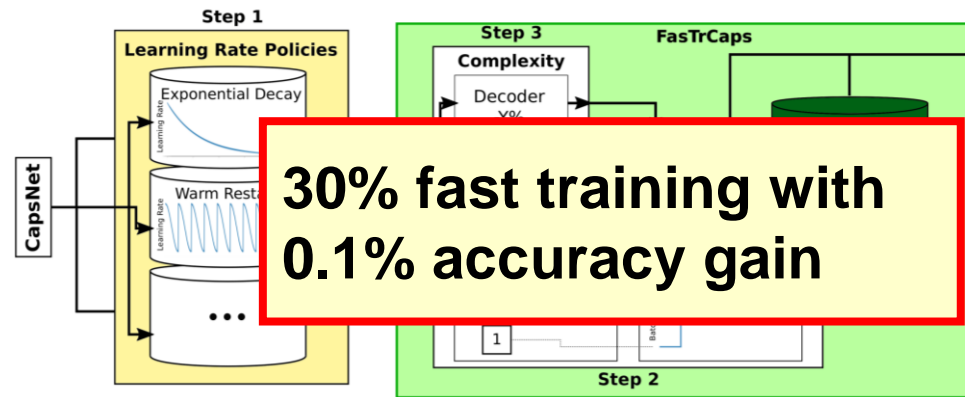
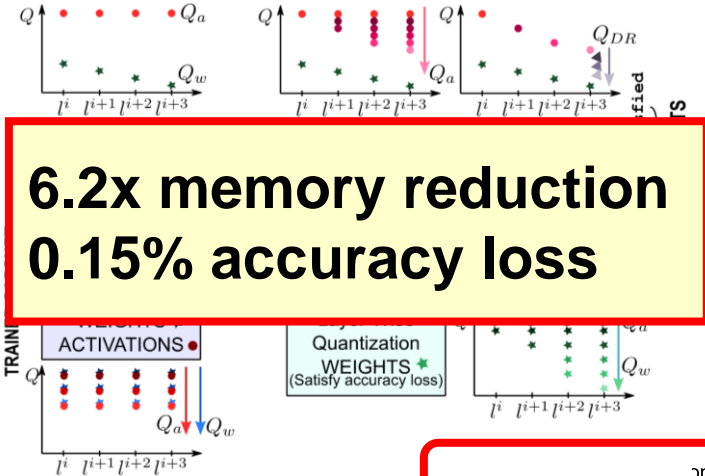
Partitioning each data type

**Compared to other mapping policies and reuse schedules,**

- up to 96% EDP improvements in DDR3
- up to 94% EDP improvements in SALP architectures

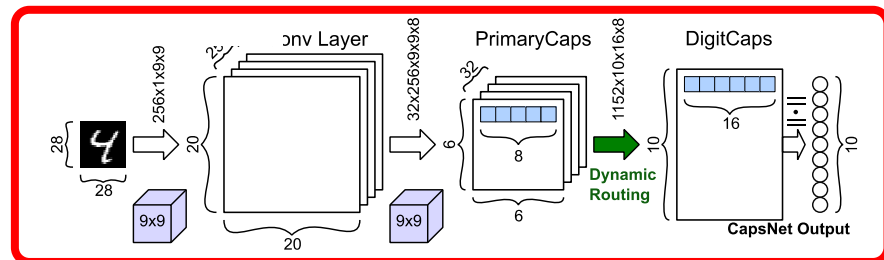
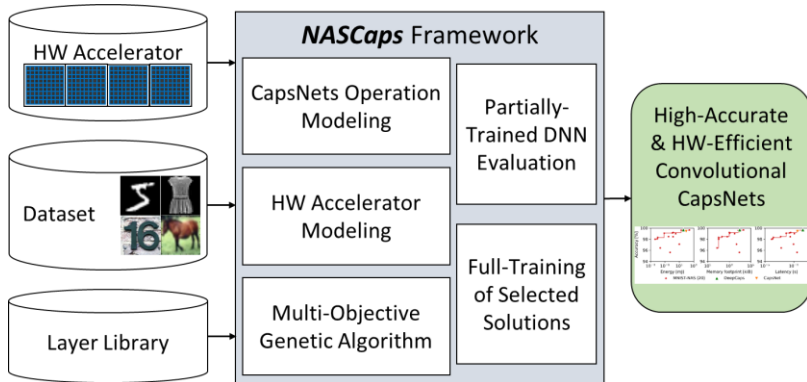
# Capsule Networks Research

**1**  
**QCaps: Quantization Framework (DAC'20)**

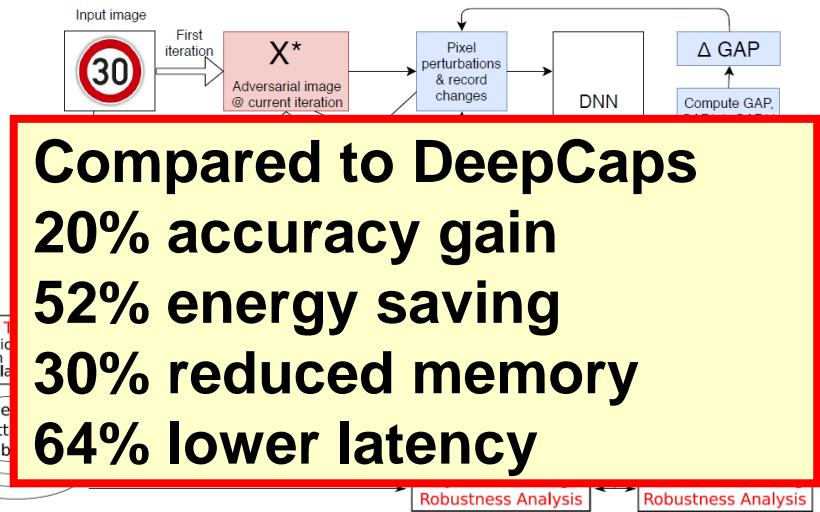


**2**  
**FasTrCaps: Fast Training (IJCNN'20)**

**3**  
**NASCaps: NAS Framework for CapsNet (ICCAD'20)**



**4**  
**RobCaps: Security & Robustness (under Review)**

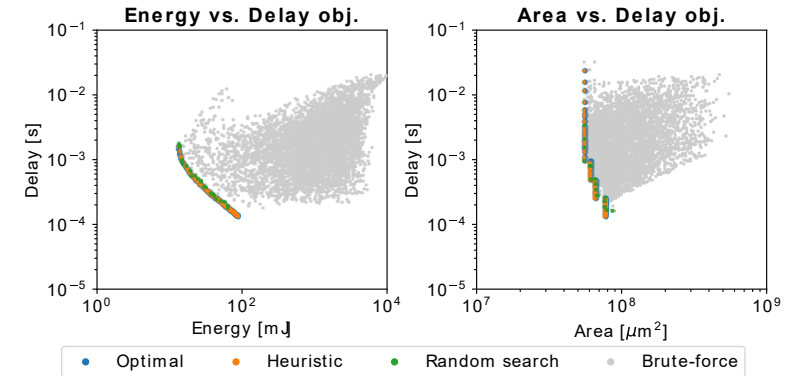
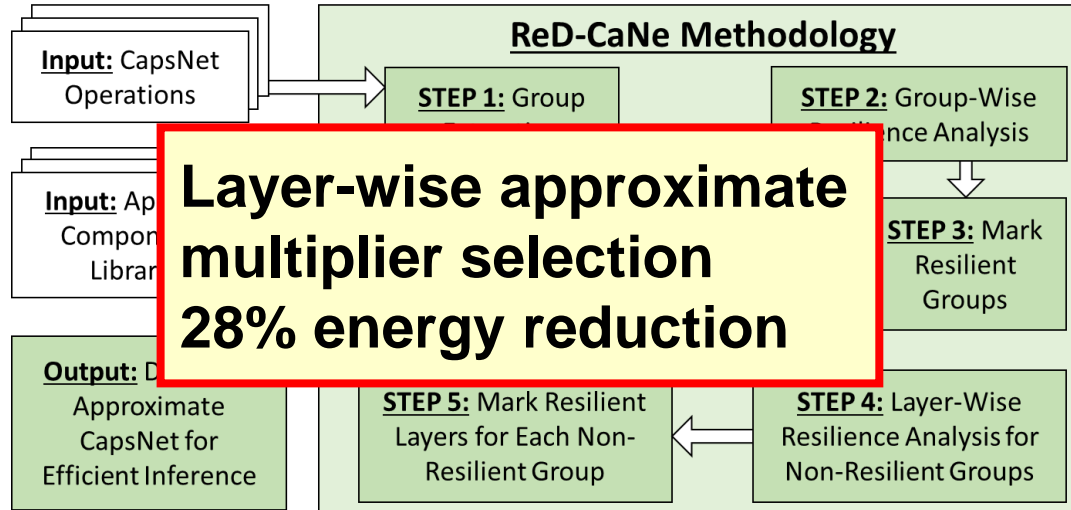


**Compared to DeepCaps**  
**20% accuracy gain**  
**52% energy saving**  
**30% reduced memory**  
**64% lower latency**

# Capsule Networks Research

5

Approximate CapsNet Design (DATE'20)

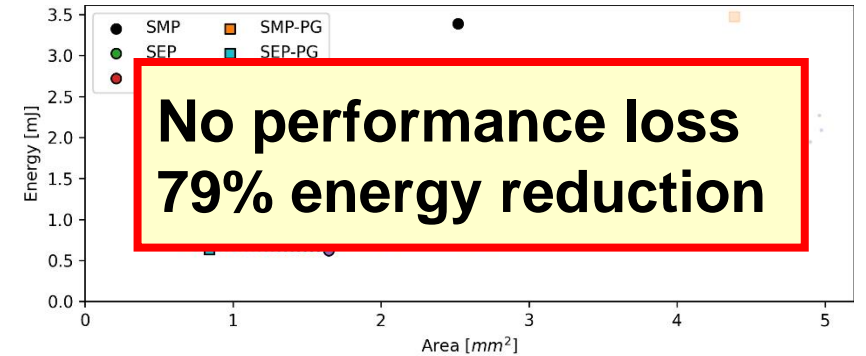
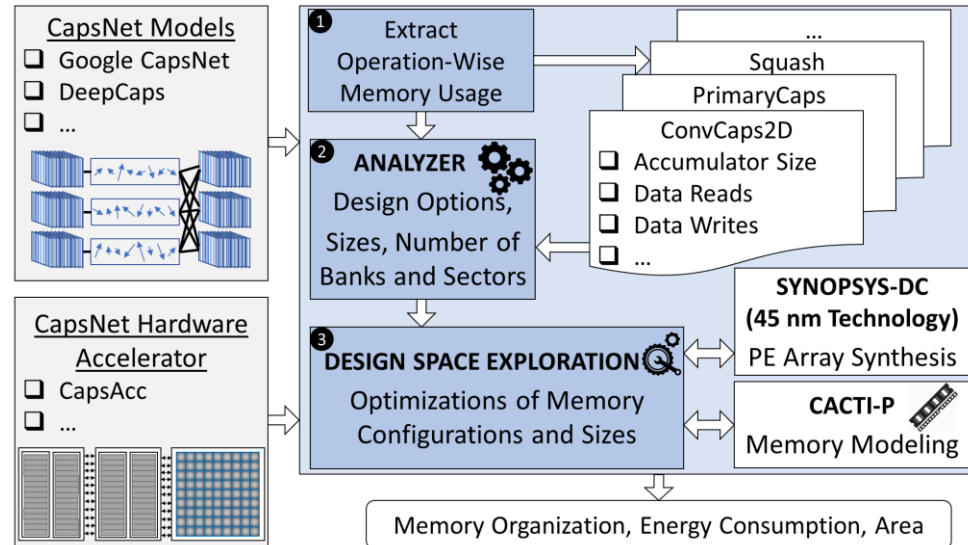


6

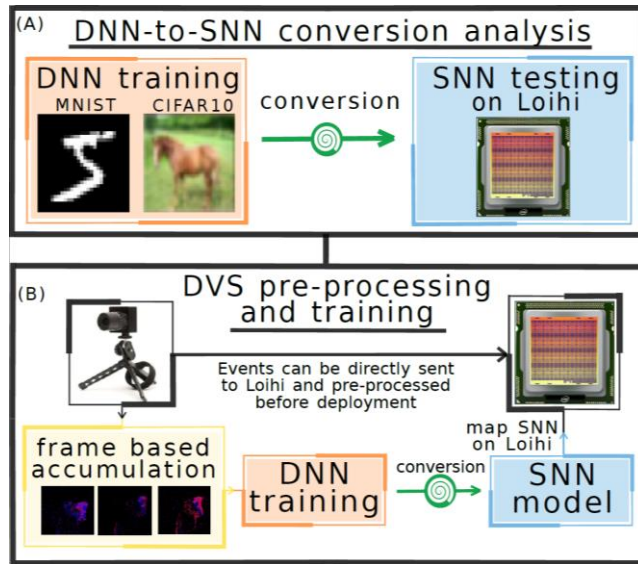
DSE of the PE Array for CapsNet Accelerators (IEEE TVLSI'21)

DESCNet: Scratchpad Memory Design for CapsNet Hardware (IEEE TCAD'20)

7

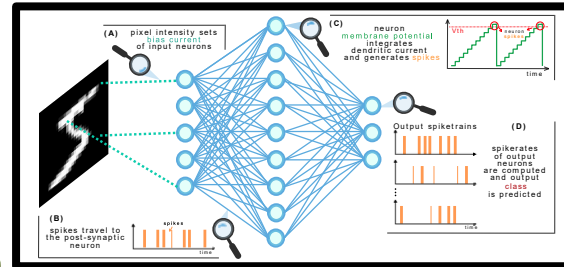
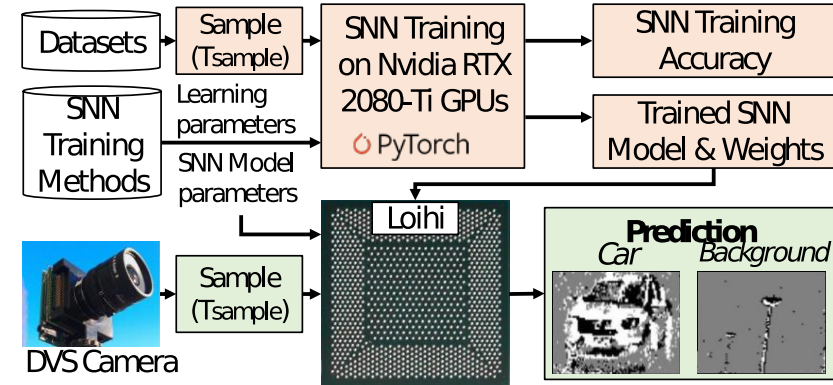


# Neuromorphic Computing using Intel's Loihi

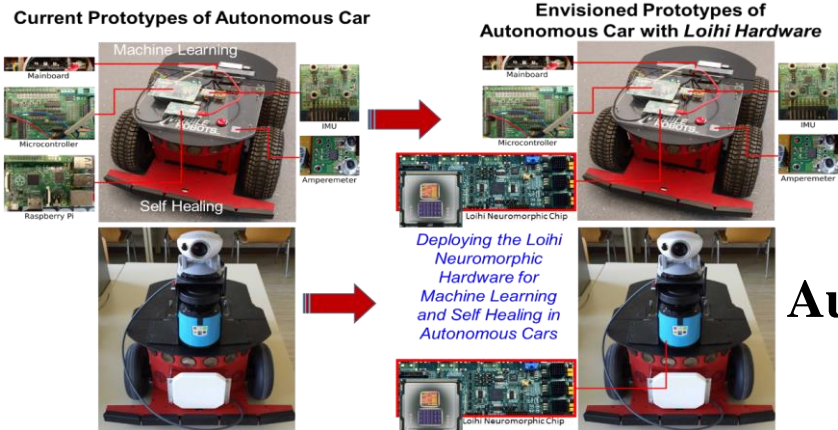


## SNN Mapping over Intel's Loihi Processor (IJCNN'20)

1



## 2 DVS-Based Car vs. Background Classification on Intel's Loihi (IJCNN'21)



## Autonomous Driving

## Smart Farming

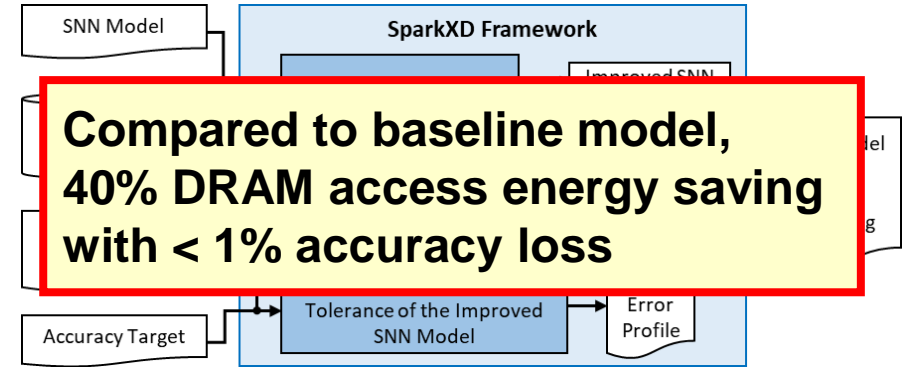


# Spiking Neural Networks Research

## Energy-Aware Optimizations and Learning Methods (IEEE TCAD'20)

1

- Our Novel Contributions
- Compared to state-of-the-art model,
    - 7.5x memory saving
    - 3.5x energy improvement in training
    - 1.8x energy improvement in inference

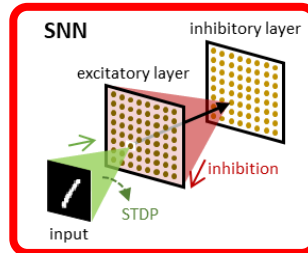


- Compared to baseline model, 40% DRAM access energy saving with < 1% accuracy loss

## SNN with Unsupervised Continual Learning (DAC'21)

3

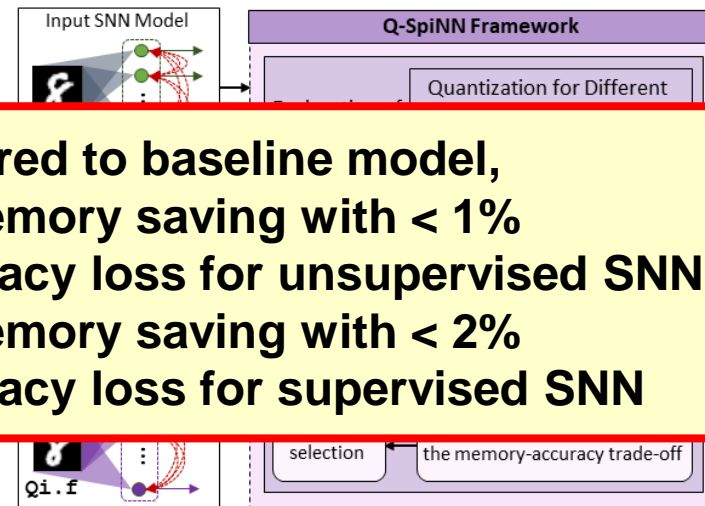
- Compared to state-of-the-art model,
  - 51% energy saving in training
  - 37% energy saving in inference
  - 21% accuracy gain for the most recently learned task
  - 8% accuracy gain for the previously learned tasks



## Resilient and Energy-Efficient SNN Inference (DAC'21)

2

- Compared to baseline model,
  - 4x memory saving with < 1% accuracy loss for unsupervised SNN
  - 2x memory saving with < 2% accuracy loss for supervised SNN

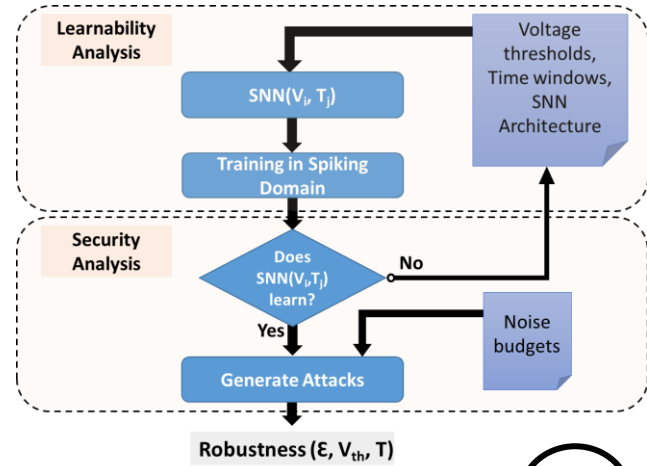


## Quantization for SNNs (IJCNN'21)

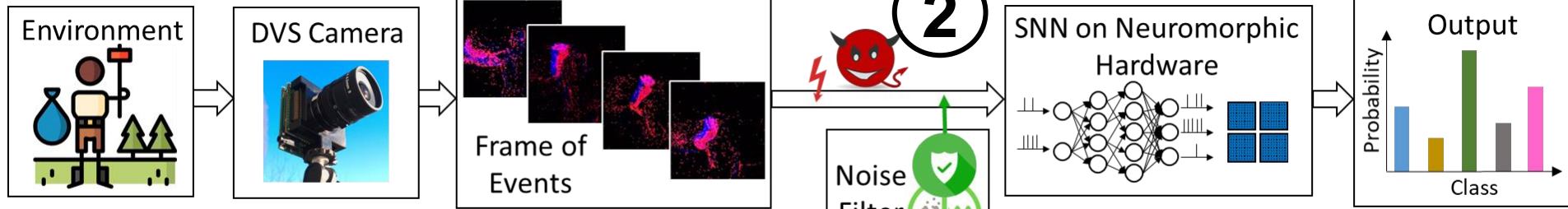
4

# Security for SNNs & Neuromorphic Computing

## Robust SNN Design against Adversarial Attacks (DATE'21) ①

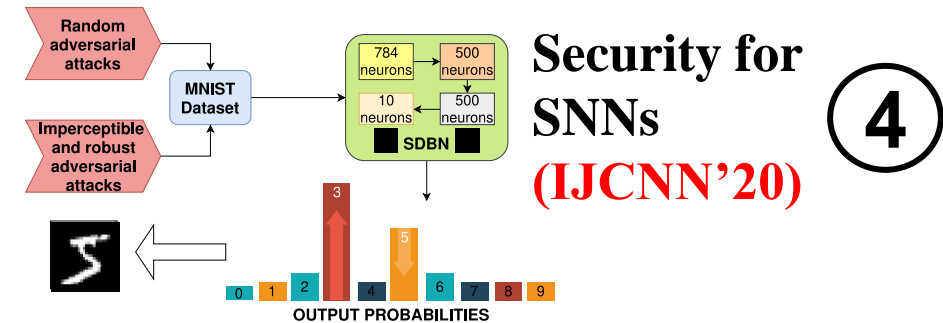
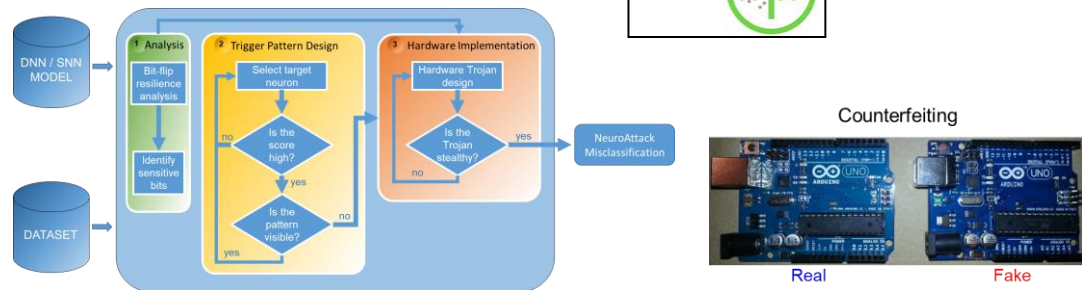


Same clean accuracy than CNN  
75% higher accuracy for large perturbations



## Adversarial Perturbations for Dynamic Vision Sensors (IROS'21, IJCNN'21)

## Fault-Injection Attacks (IJCNN'20) ③

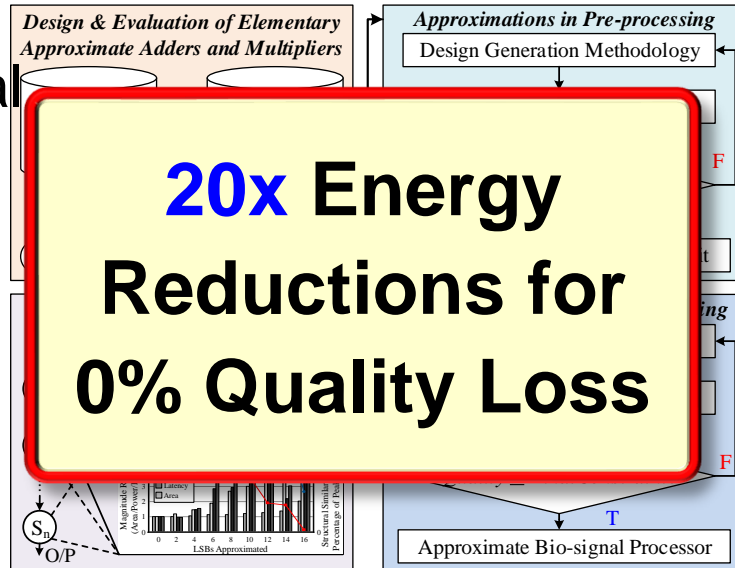


## Security for SNNs (IJCNN'20) ④

# Energy-Efficient IoT-Healthcare and AI

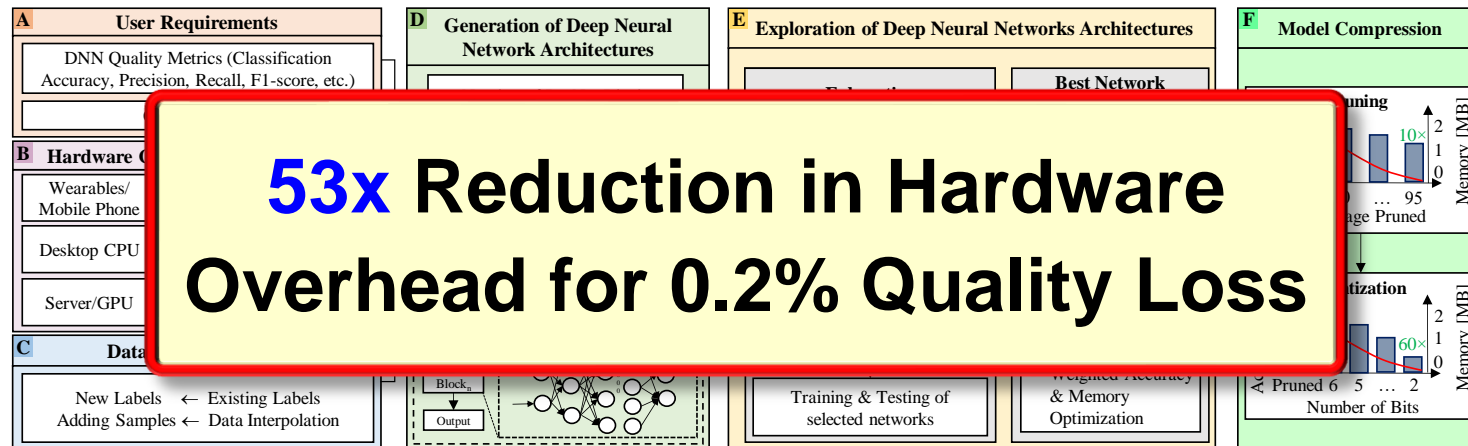
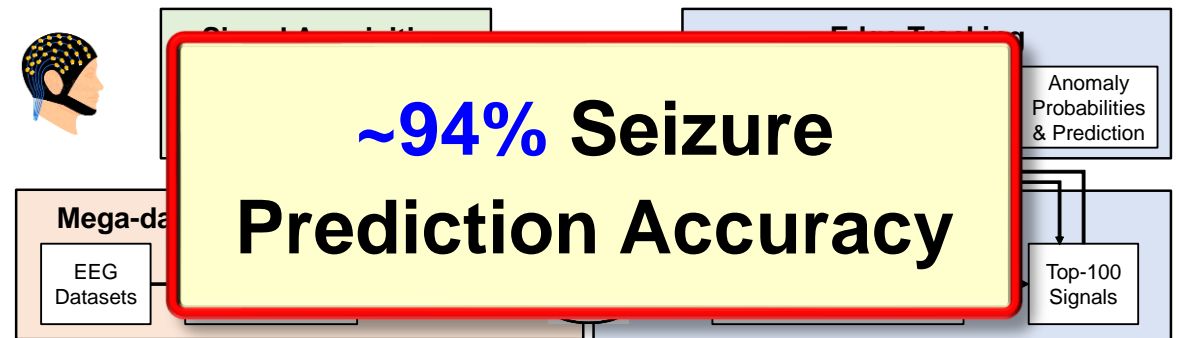
Methodology for Approx. Bio-signal Processing: **DAC'19**

①



Cloud-Edge Framework for EEG Monitoring and Real-time Anomaly Prediction: **DAC'20**

②



NAS for HW-Constrained Healthcare DNNs: **(IEEE IoT'21)**

③

# EdgeAI for Healthcare: Moore4Medical EU Project



Src: Google Images

## Next Generation Ultrasound



- Data Acquisition
- 3D Reconstruction
- Edge Processing
- AI algorithms for detecting fetus' anatomical features
- Hardware accelerator for high throughput feature extraction
- Closed-loop system for real-time user feedback

- Investigating **DL architectures** and **statistical ML techniques** for classification, segmentation, and anatomical feature extraction
- Evaluating requirements of proposed algorithms to develop **energy-efficient hardware accelerators for edge processing**
- Develop **FPGA prototype** to demonstrate the efficacy of the accelerator and deployability of the HW-SW system



# Future Research Directions

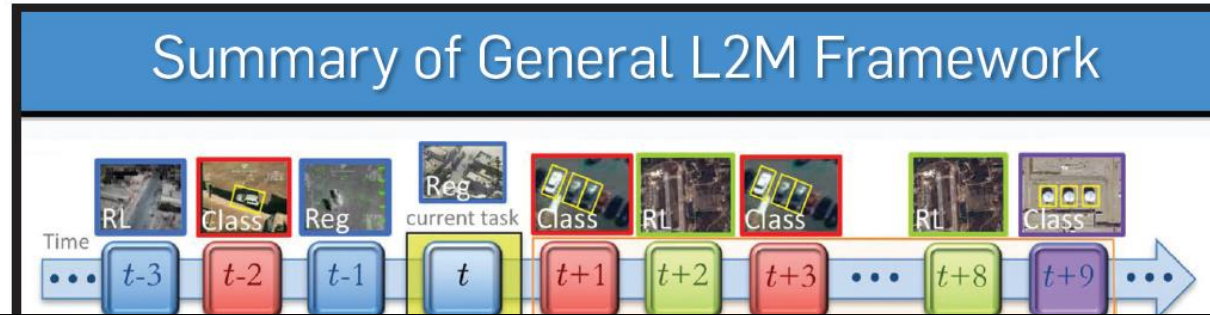
- New computing paradigms such as near-/in-memory computing and approximate computing
- It is not all about deep learning. Conventional machine learning models can offer better performance in some scenarios.
- Optimization frameworks for all types of systems, as the selection is limited in some scenarios due to other constraints, e.g., cost.
- Novel techniques for training and optimizing machine learning models
- Interpretability of models to ensure robustness

# Summary

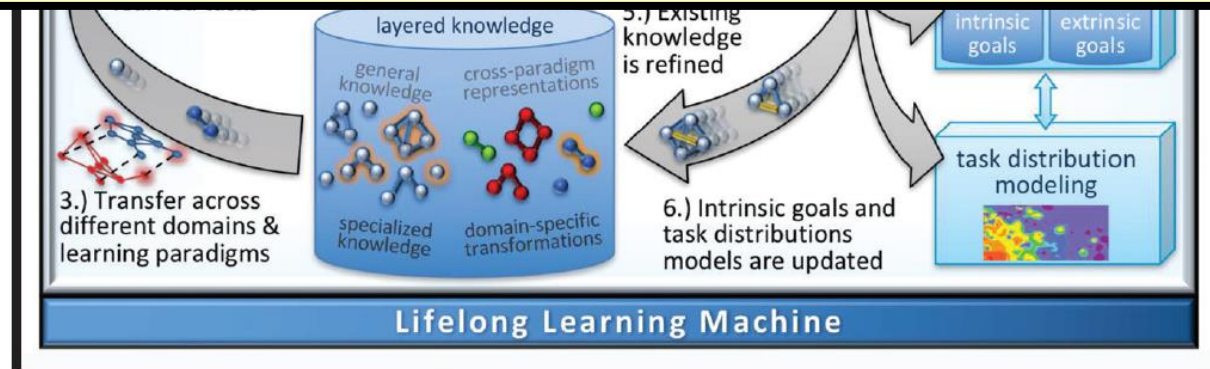
- ❑ **Artificial Intelligence** has proliferated almost everywhere, **that's for a good reason!** => *the big data challenge!*
  - ❑ Cloud, Fog, Edge, ..., In-Sensor / In-Situ
- ❑ **Required: High-Throughput, Energy-Efficient, & Robust Designs**
- ❑ **Our System-Level Framework**
  - ❑ Optimizations across the Software & Hardware stacks
  - ❑ Specialized hardware accelerators, dataflows, memory, self-healing approximations, hardware-aware NAS, ...
  - ❑ Selective Tile Processing for energy-efficient object detection
  - ❑ Robustness: Analyzing security attacks and hardware-level faults.

**A system level approach requires bridging the gap between the AI/ML community & System designers (HW + SW)**

# Lifelong Learning in Artificial Neural Networks



**“In a few years, much of what we consider AI today won’t be considered AI without lifelong learning”**

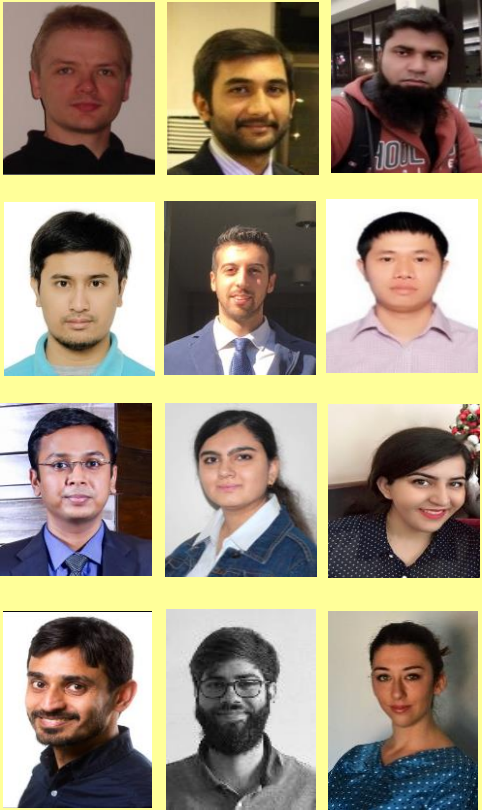


Data and image source:  
“Lifelong Learning in Artificial  
Neural Networks” in  
Communications of the ACM



# My Research Team and Collaborators

## Post-Docs and PhDs



## MS/BS Students



## Key Collaborators



## Previous Students





# 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](http://www.tinyml.org)**