

# tinyML<sup>®</sup> Talks

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

## “Physics-Aware Auto Tiny Machine Learning”

Swapnil Sayan Saha – Algorithm Development Engineer, STMicroelectronics Inc.

March 5, 2024



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



Thank you, **tinyML Strategic Partners**,  
for committing to take tinyML to the next Level, together



T I N Y



TALKS  
*webcast*

# Executive Strategic Partners

**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



Accelerate Your Edge Compute

**SYNTIANT**

Making Edge AI A Reality

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

# Platinum Strategic Partners

T I N Y



TALKS  
*webcast*

# *embed* UR





**DEPLOY VISION AI  
AT THE EDGE AT SCALE**

**SONY**



# Gold Strategic Partners

Build the  
Future of tinyML

on **arm**



T I N Y



TALKS  
*webcast*



**EDGE IMPULSE**

# The Leading Development Platform for Edge ML

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

Decarbonization

Digitalization



Driving decarbonization and digitalization. Together.

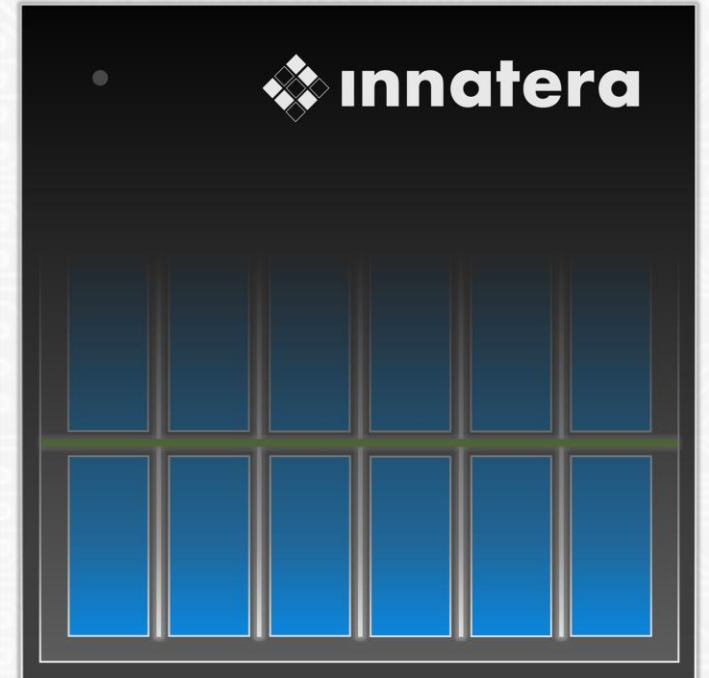
**Infineon serving all target markets as**  
**Leader in Power Systems and IoT**

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



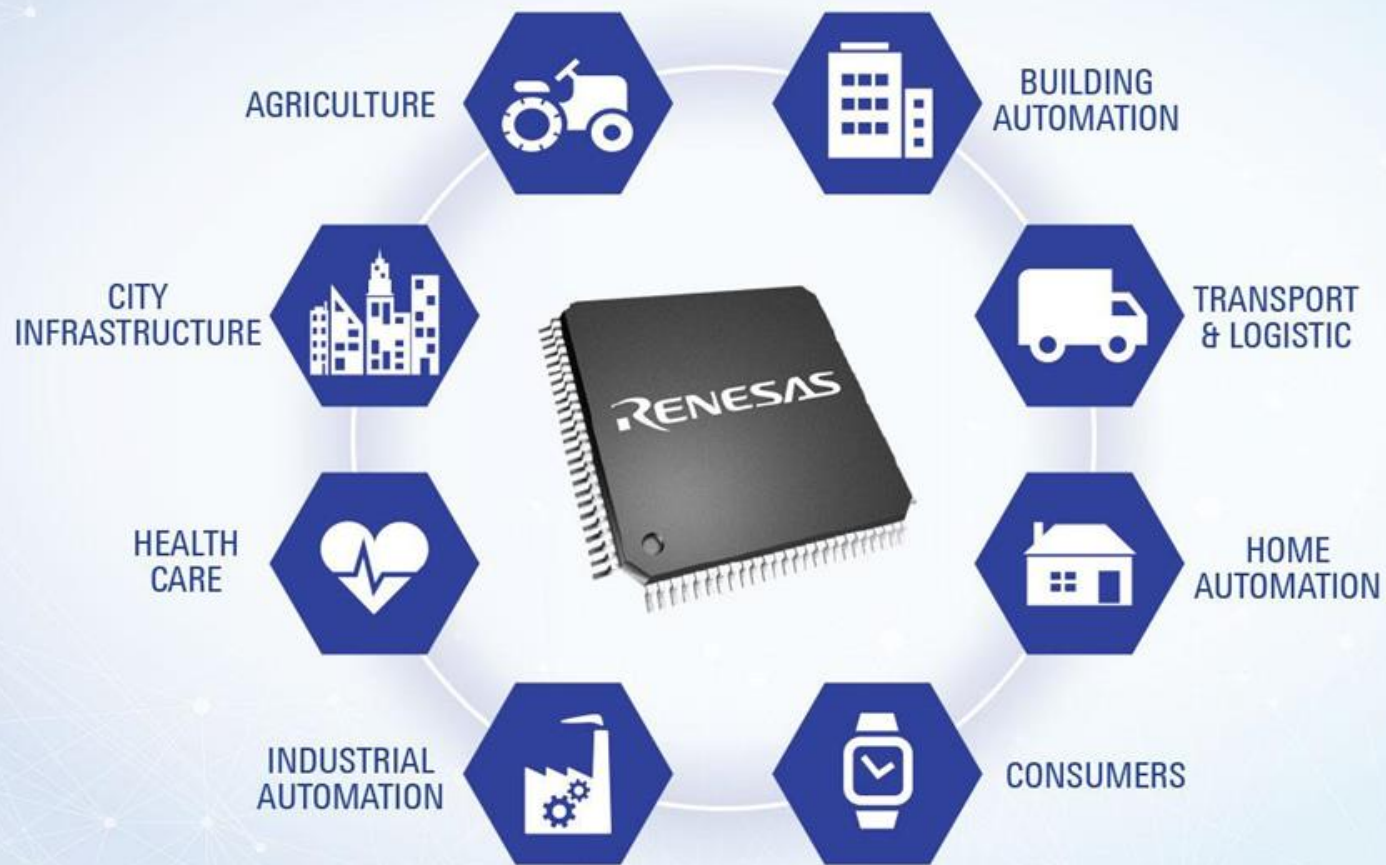


# NEUROMORPHIC INTELLIGENCE FOR THE SENSOR-EDGE



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

**Renesas is enabling the next generation of AI-powered solutions that will revolutionize every industry sector.**



[renesas.com](https://www.renesas.com)



life.augmented

**STMicroelectronics provides extensive solutions to make tiny Machine Learning easy**



# ENGINEERING EXCEPTIONAL EXPERIENCES

We engineer exceptional experiences for consumers in the home, at work, in the car, or on the go.

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





T I N Y



# Silver Strategic Partners



brainchip



GREENWAVES  
TECHNOLOGIES



£Grovety Inc.



Nota AI



QORVO





# Join Growing tinyML Communities:



19.6k members in  
49 Groups in 41 Countries

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

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



4.2k members  
&  
14.5k 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](http://www.youtube.com/tinyML)



**tinyML**  
4.33K subscribers

**11.8k subscribers, 662 videos with 432k views**

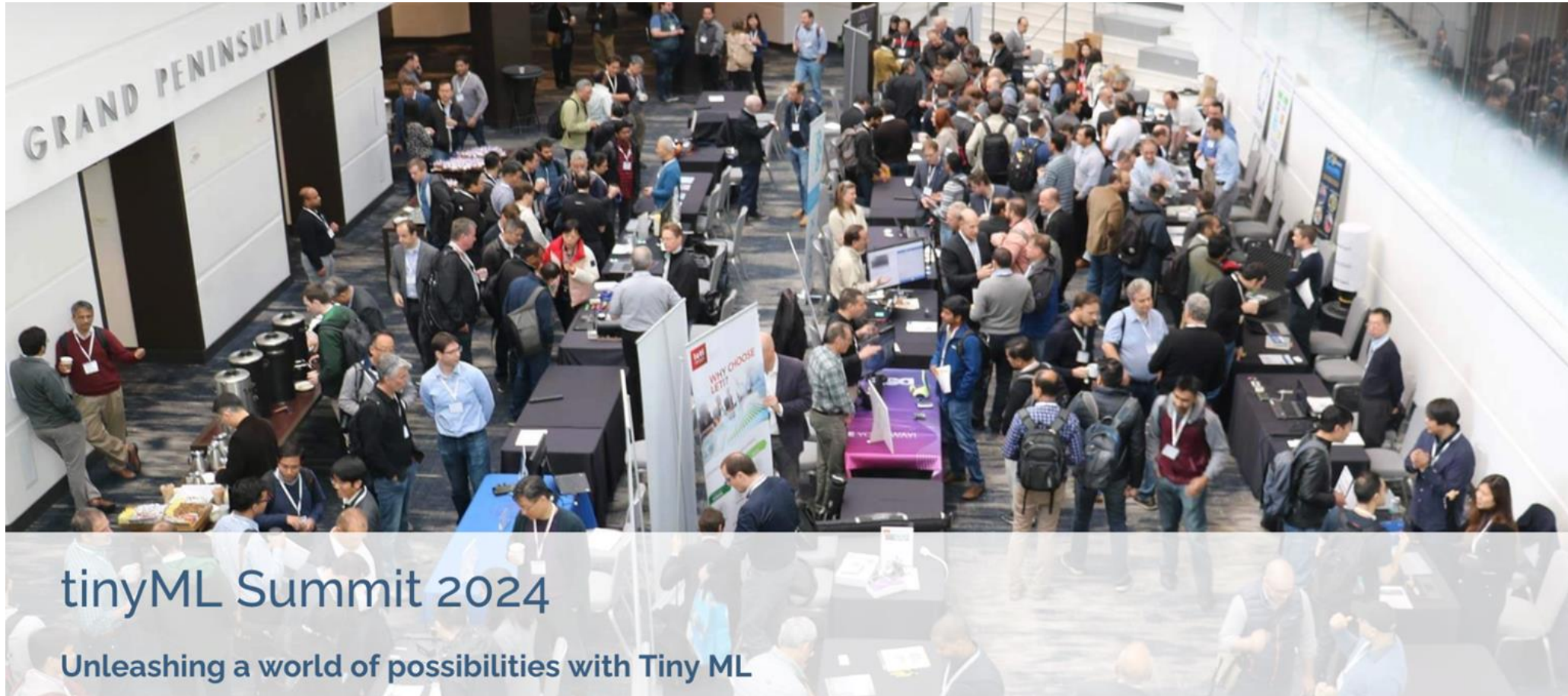
HOME VIDEOS PLAYLISTS COMMUNITY CHANNELS ABOUT

106 views · 4 days ago	138 views · 4 days ago	54 views · 4 days ago	47 views · 4 days ago	132 views · 4 days ago	137 views · 4 days ago
122 views · 4 days ago	262 views · 2 weeks ago	511 views · 3 weeks ago	229 views · 3 weeks ago	265 views · 3 weeks ago	286 views · 1 month ago
351 views · 1 month ago	462 views · 2 months ago	374 views · 2 months ago	133 views · 2 months ago	287 views · 2 months ago	336 views · 2 months ago
378 views · 2 months ago	214 views · 2 months ago	448 views · 2 months ago	159 views · 2 months ago	190 views · 2 months ago	545 views · 2 months ago



# tinyML Summit

## April 22 - 24, 2024 - Register now!





# tinyML Awards 2024



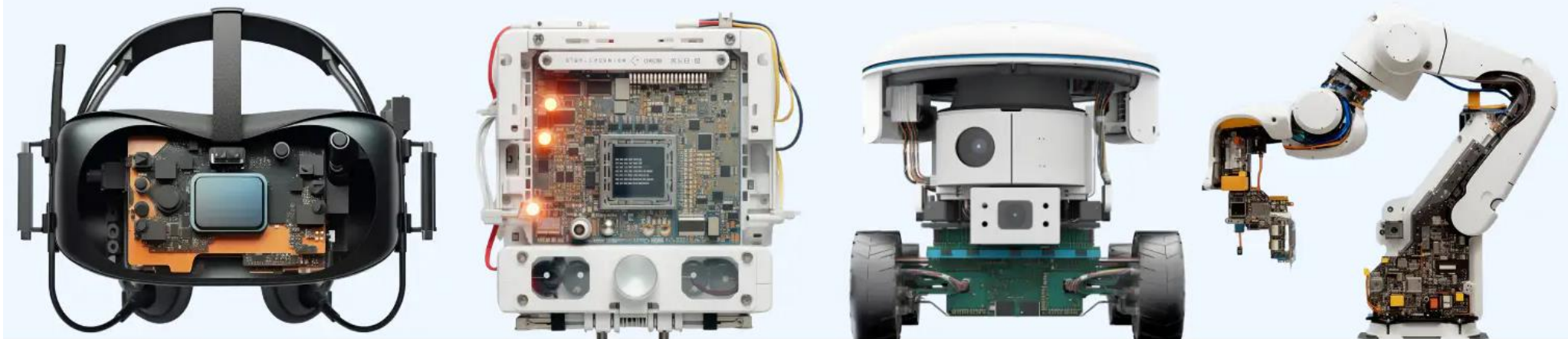
## tinyML Awards 2024

- Best Product (*Tiny ML chip, Audio or Vision Application, Sensor Application Product*)
- Best Prototype
- Sustainable Future Pioneer Award



# 2023 Edge AI Technology Report

The guide to understanding the state of the art in hardware & software in Edge AI.





# Reminders

Slides & Videos will be posted tomorrow



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



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



Please use the Q&A window for your questions





## Swapnil Sayan Saha



Swapnil Sayan Saha is an algorithm development engineer at STMicroelectronics Inc. He received his Ph.D. and M.S. in Electrical and Computer Engineering from the University of California, Los Angeles in 2023 and 2021 respectively, and B.Sc. in Electrical and Electronics Engineering from the University of Dhaka in 2019. His research explores how rich, robust, and complex inferences can be made from sensors onboard low-end embedded systems within tight resource budgets in a platform-aware fashion. To date, he has published more than 25 peer-reviewed articles/patents and received more than 30 awards in robotics, technical, and business-case forums worldwide.





life.augmented

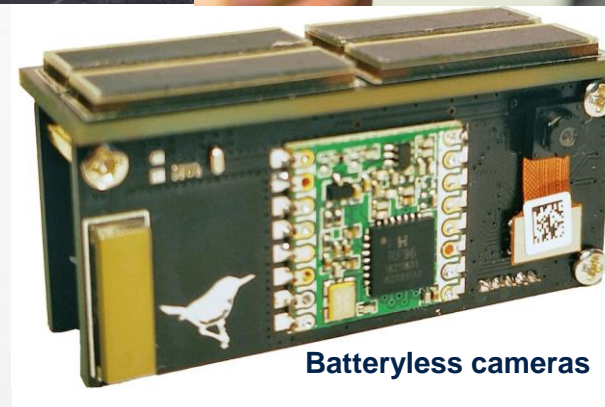
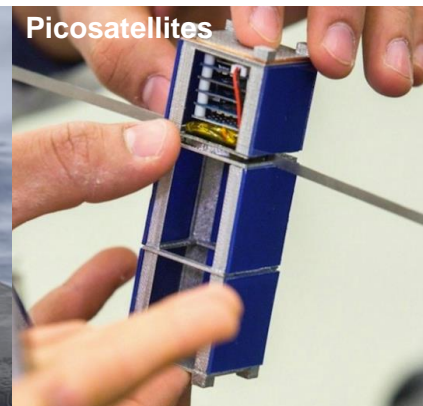
# Physics-aware auto tiny machine learning

**Swapnil Sayan Saha, Ph.D.**  
Algorithm Development Engineer  
STMicroelectronics



# Tiny machine learning

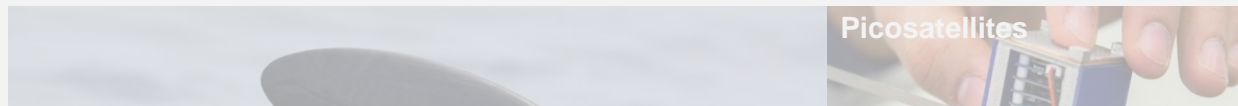
Hardware and software suites that enable **always-on**, **ultralow power**, and **on-device data analytics**



Enables applications that need to make **“complex inferences”** for **“time-critical”** and **“remote”** applications from **“unstructured data”** independent of large systems.

# Tiny machine learning

Hardware and software suites that enable always-on, ultra-low power, and on-device data analytics.

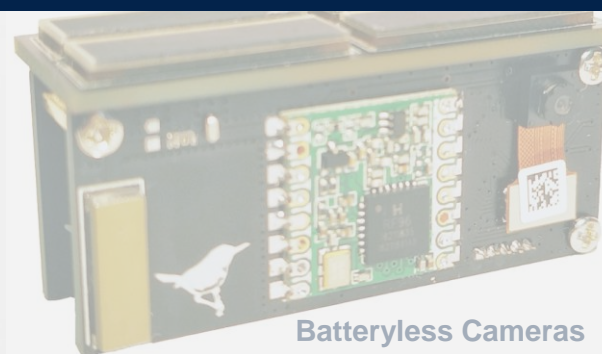


Picosatellites

**First generation efforts focused on squeezing standalone neural networks within the resource bounds of tinyML platforms**



Micro-Aerial Vehicles

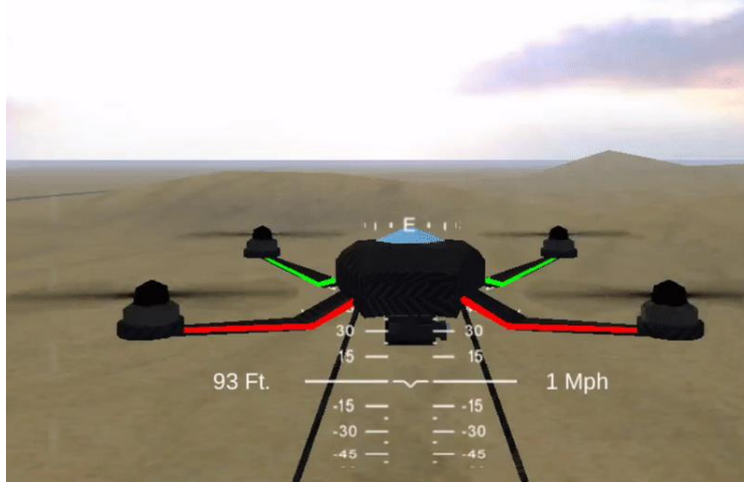


Batteryless Cameras

large systems.

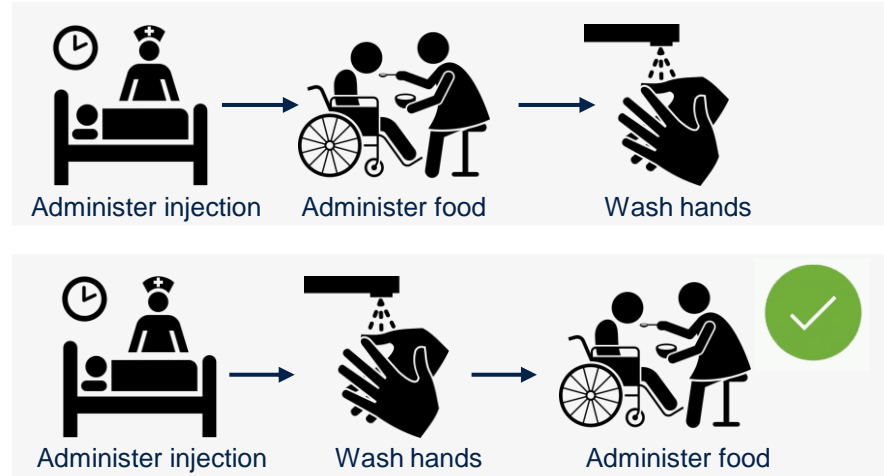
# Challenge 1: obeying physics, rules, and constraints

**Goal: land the quadrotor as fast as possible**



Quadrotor crash due to unsafe bank angle caused by a neural flight controller.

**Goal: guide a nurse perform correct action sequences**

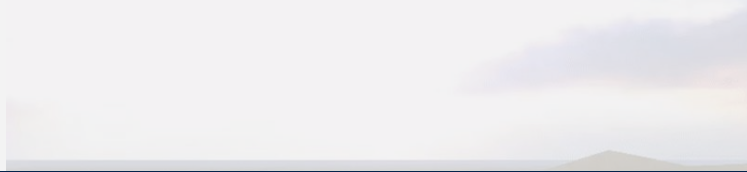


Invalid and valid sanitary protocol sequence in a nurse care setting for complex event processing.

**Physics matter more at the edge; standalone neural networks cannot assure that the learned distributions obey the rules and physics of the underlying system**

# Challenge 1: obeying physics, rules, and constraints

Goal: land the quadrotor as fast as possible



Goal: guide a nurse perform correct action sequences



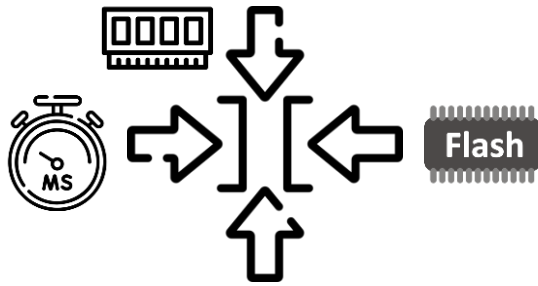
**Additional routines (called symbolic programs) must be jointly deployed with the neural network**

Quadrotor crash due to unsafe bank angle caused by a neural flight controller.

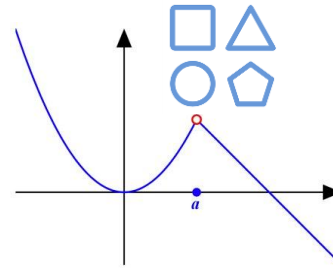
Invalid and valid sanitary protocol sequence in a nurse care setting for complex event processing.

**Physics matter more at the edge; standalone neural networks cannot assure that the learned distributions obey the rules and physics of the underlying system**

# Challenge 2: synthesizing platform-aware neurosymbolic programs



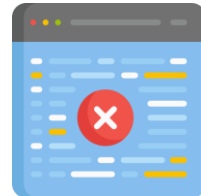
Fitting neural and symbolic routines within resource limits



Dealing with mixed and discontinuous parameter spaces



Arbitrary cost function formulation

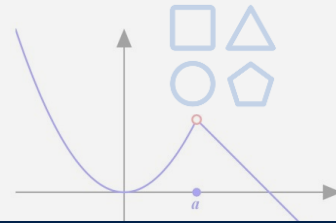
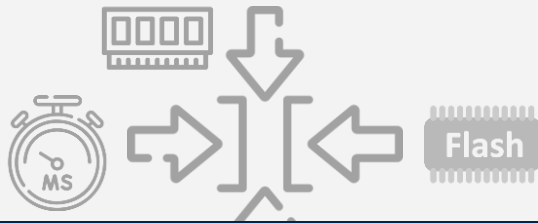


Dealing with runtime faults due to more moving parts

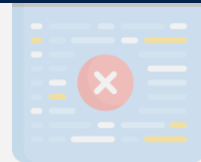


Finding proxies for so many program combinations

# Challenge 2: synthesizing platform-aware neurosymbolic programs



**Finding the optimal synergy between neural and symbolic components within the tight resource constraints is challenging**

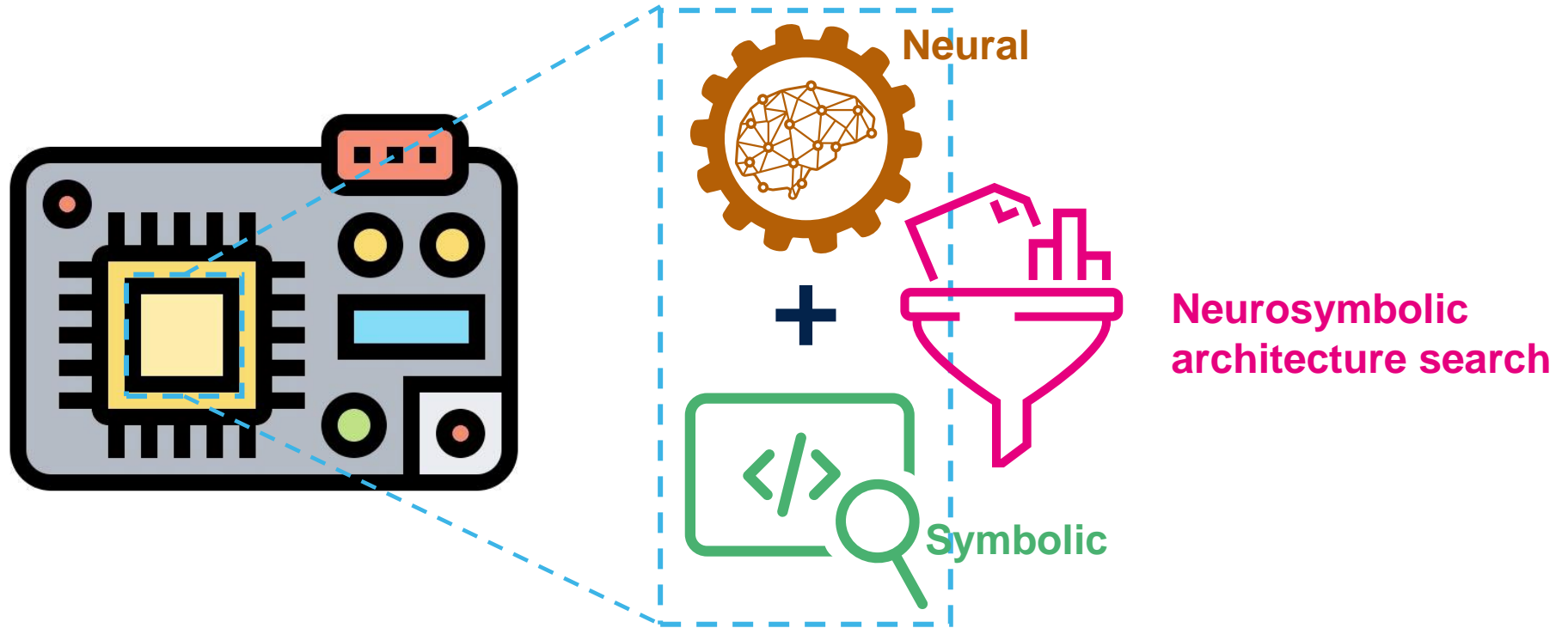


Dealing with runtime faults due to more moving parts



Finding proxies for so many program combinations

# Neurosymbolic auto tiny machine learning





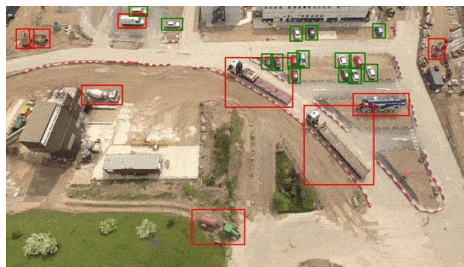
# What is neurosymbolic artificial intelligence?

A program containing neural and human-readable (symbolic) code

**Performant**

**Non-parametric  
(flexible and  
scalable)**

**Robust**



Aerial object detection



**Neurosymbolic AI**

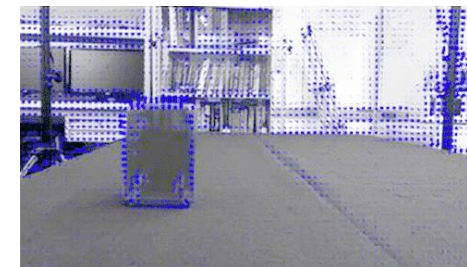
Self-navigating  
surveillance quadrotor



**Interpretable**

**Data efficient**

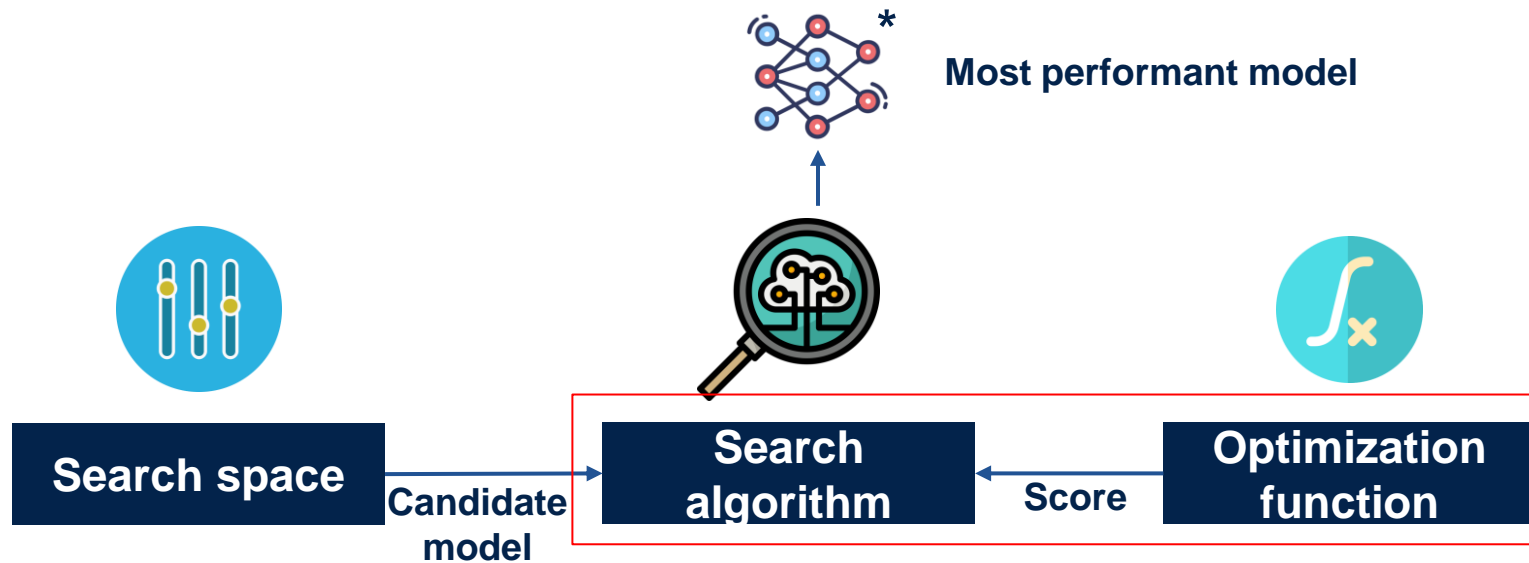
**Physics and context-  
aware**



Optical flow

# What is neural architecture search?

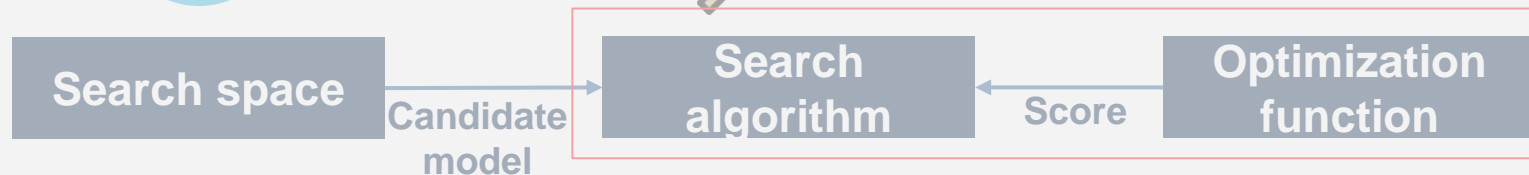
Automatically find the most performant neural network architecture from a hyperparameter space within some constraints



# What is neural architecture search?

Automatically find the most performant neural network architecture from a hyperparameter space within some constraints

**We adopt Mango, a black-box Bayesian optimizer, which can efficiently handle mixed and discontinuous search spaces**



# Gradient-free Bayesian optimization: surrogate function

Two components: surrogate function and acquisition function

$$\hat{f}(\Omega) \sim \mathcal{GP}(\mu(\Omega), k(\Omega, \Omega'))$$

A surrogate function approximates an optimization function, e.g., gaussian process

Gaussian process provides tractable assessment of uncertainty under data scarcity

Non-parametric model using mean  $\mu$  and Matern kernel function  $k$  over the search space  $\Omega$

Sandha, Sandeep Singh, Mohit Aggarwal, Igor Fedorov, and Mani Srivastava. "Mango: A python library for parallel hyperparameter tuning." in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3987-3991. IEEE, 2020.

Sandha, Sandeep Singh, Mohit Aggarwal, Swapnil Sayan Saha, and Mani Srivastava. "Enabling hyperparameter tuning of machine learning classifiers in production." In *2021 IEEE Third International Conference on Cognitive Machine Intelligence (CogMI)*, pp. 262-271. IEEE, 2021.

# Gradient-free Bayesian optimization: acquisition function

Two components: surrogate function and acquisition function

$$\Omega_t = \arg \max_{\Omega} \underbrace{(\mu_{t-1}(\Omega))}_{\text{gear icon}} + \beta^{0.5} \underbrace{\sigma_{t-1}(\Omega)}_{\text{magnifying glass icon}}$$

An acquisition function selects the next promising set of points to sample

Mango adopts Monte Carlo sampling with upper confidence bound, using adaptive exploration factor  $\beta$

First term: goodness of sampled point (exploitation);  
second term: uncertainty of sampled point (exploration);  
does not get stuck in local optima

# Gradient-free Bayesian optimization: acquisition function

Two components: surrogate function and acquisition function

$$\Omega_t = \arg \max_{\Omega} (\underbrace{\mu_{t-1}(\Omega)}_{\text{exploitation}} + \beta^{0.5} \underbrace{\sigma_{t-1}(\Omega)}_{\text{exploration}})$$

**Adaptive exploration factor finds near-optimal values at the boundary of violating deployability constraints with 90% theoretical guarantees**

bound, using adaptive exploration factor  $\beta$

First term: goodness of sampled point (exploitation);  
second term: uncertainty of sampled point (exploration);  
does not get stuck in local optima

# Formulating the neurosymbolic optimization function

$$\min \mathbf{f}(\Omega), \text{ s.t. } \mathbf{f}(\Omega) \leq \mathbf{b}$$

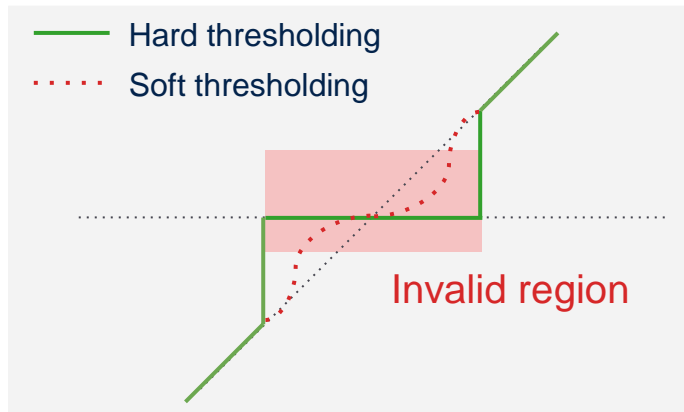
$$\min f_{\text{opt}}, \quad f_{\text{opt}} = \lambda_1 f_{\text{error}}(\Omega) + \lambda_2 f_{\text{flash}}(\Omega) + \lambda_3 f_{\text{SRAM}}(\Omega) + \lambda_4 f_{\text{latency}}(\Omega)$$

**GP-UCB solves a non-linear program with constraints**

**Goal: construct a fault-free neurosymbolic program such that latency and error are minimized, while the memory usage is maximized within device memory limits**

**Search space  $\Omega$  contains both neural and symbolic hyperparameters, trainable weights, neural operators, and symbolic program atoms**

# Fast and guaranteed deployment: hard thresholding



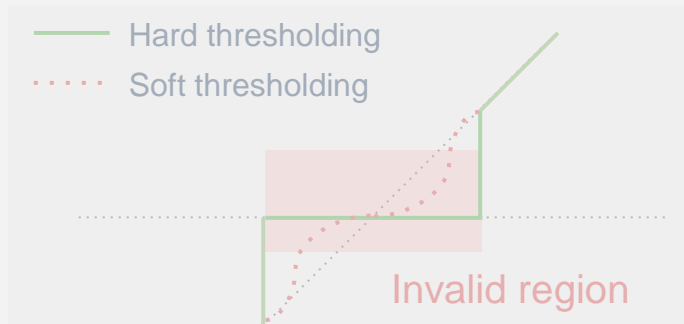
To guarantee deployability and optimize at the execution level, we perform platform-in-the-loop search

If a model induces faults, we do not train the model; the search algorithm is penalized by a constant large number (hard thresholding)

Thanks to GP-UCB, the search algorithm is able to observe the small valid linear region where memory usage and accuracy are proportional



# Fast and guaranteed deployment: hard thresholding

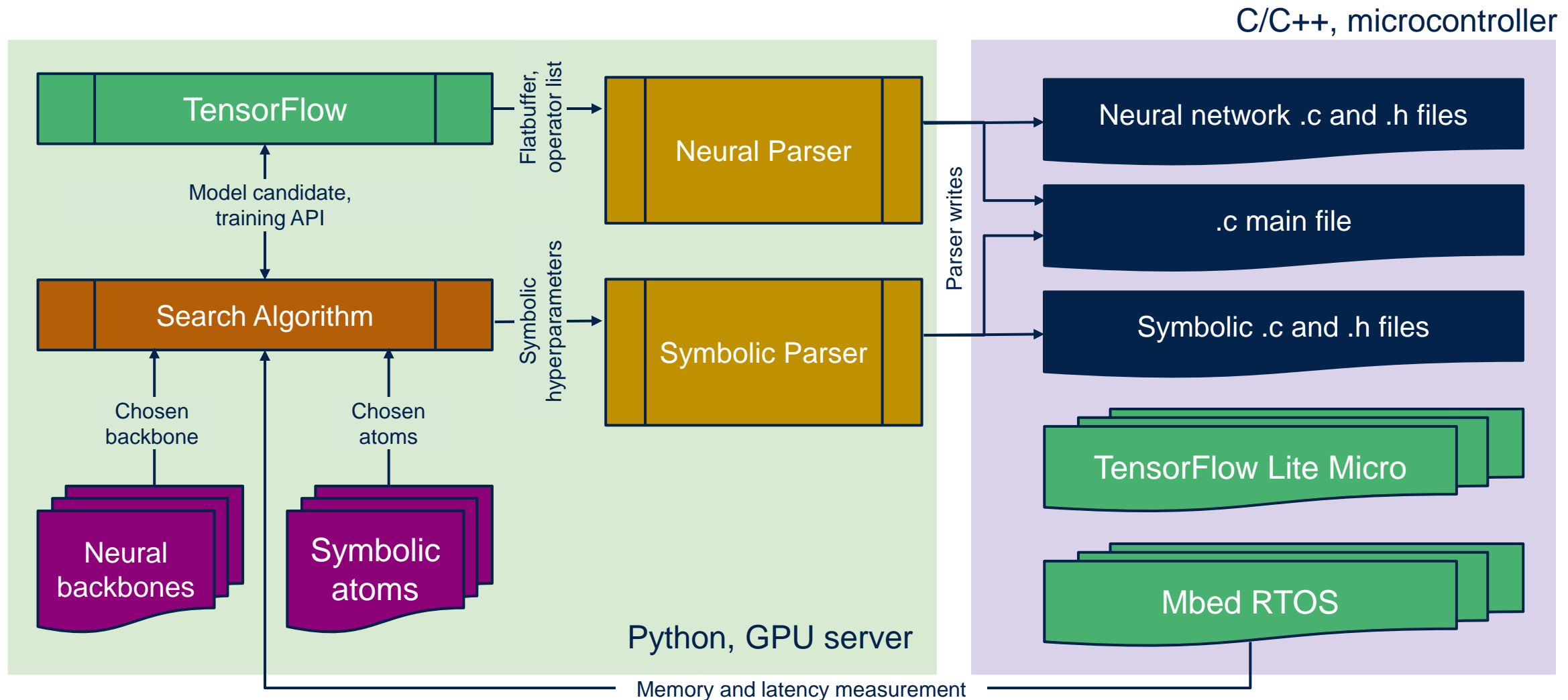


**Platform-in-the-loop + hard thresholding =  
50% faster than proxy + soft thresholding**

If a model induces faults, we do not train the model; the search algorithm is penalized by a constant large number (hard thresholding)

Thanks to GP-UCB, the search algorithm is able to observe the small valid linear region where memory usage and accuracy are proportional

# Implementing automatic platform-in-the-loop



# Comparing tinyML NAS strategies

Method	Search strategy	Profiler	Search space	Optimization terms
SpArSe	Gradient-driven Bayesian	Analytical	Conv2D	Error, SRAM, flash
MCUNet	Evolutionary (weight sharing)	Lookup tables	Conv2D	Error, SRAM, flash, latency
MicroNets	One-shot DNAS	Analytical	Conv2D	Error, SRAM, flash, latency
$\mu$ NAS	Evolutionary (no weight sharing)	Analytical	Conv2D	Error, SRAM, flash, latency
iNAS	Reinforcement learning	Lookup tables	Conv2D (execution level)	Error, flash, latency, volatile buffer, power cycle energy
UDC	DNAS with exploration/exploitation	Analytical	Conv2D, pruning, quantization	Error, flash
<b>TinyNS</b>	<b>Gradient-free Bayesian with exploration/exploitation</b>	<b>Platform-in-the-loop</b>	<b>Any supported ML operator and symbolic atoms</b>	<b>Any scalar term</b>

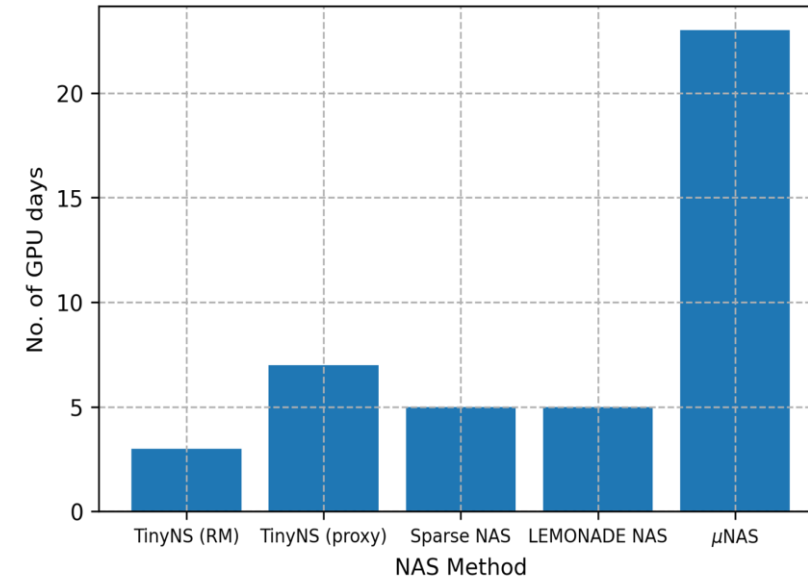
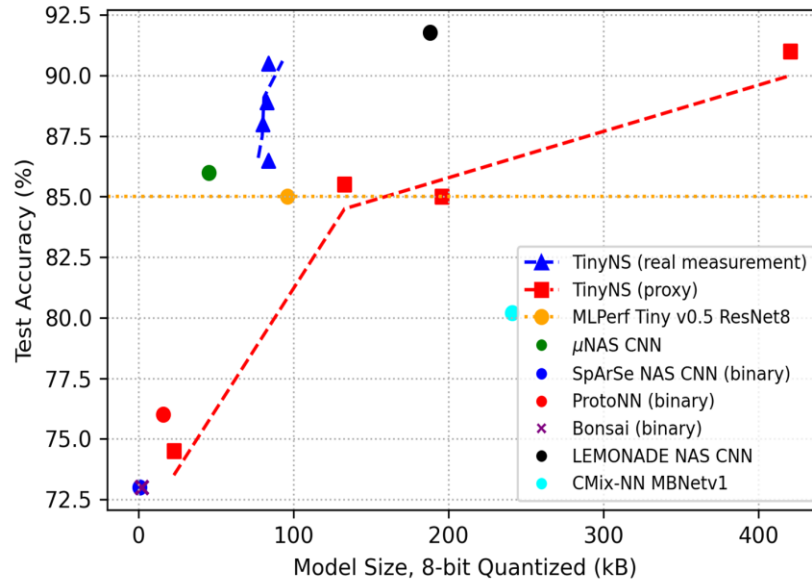
**Operates on mixed and arbitrary parameter spaces and optimization function terms**

**Platform-in-the-loop guarantees deployability by taking execution level dynamics into account**

**Efficient over RL, DNAS, and evolutionary search algorithms; combines the best features of other search algorithms in one package**

# MLPerf Tiny v0.5 inference benchmark (CIFAR10)

Backbone: ResNet8

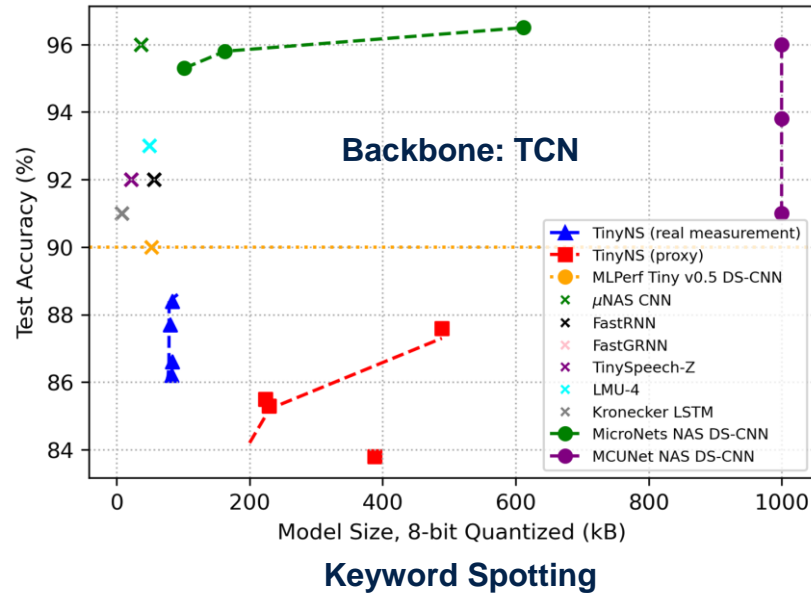
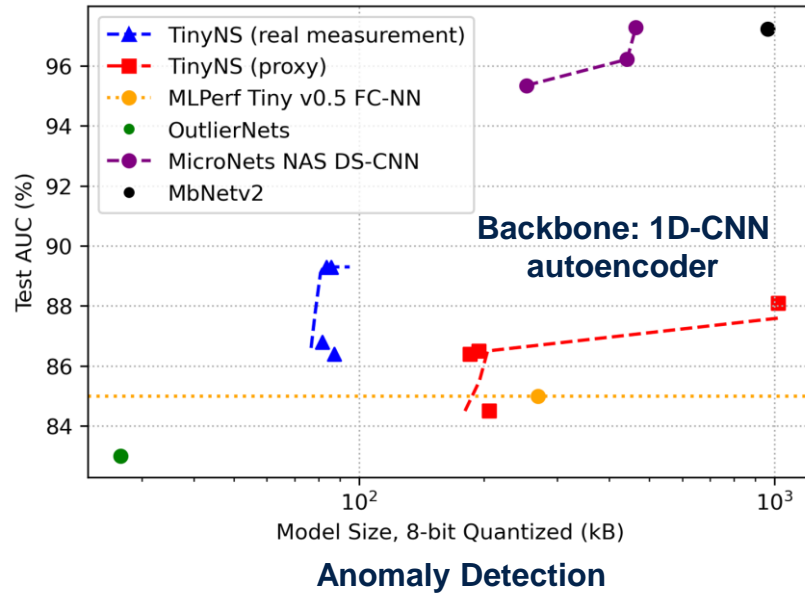


Exceeds benchmark accuracy by 4.3%, outperforms baselines by 4.5%-17.5%

1.7x-7.7x lower convergence time than baselines

Platform-in-the-loop provides models that have 1.6%-5.5% higher accuracy and consumes 4.2x lower flash, while converging 2.3x faster

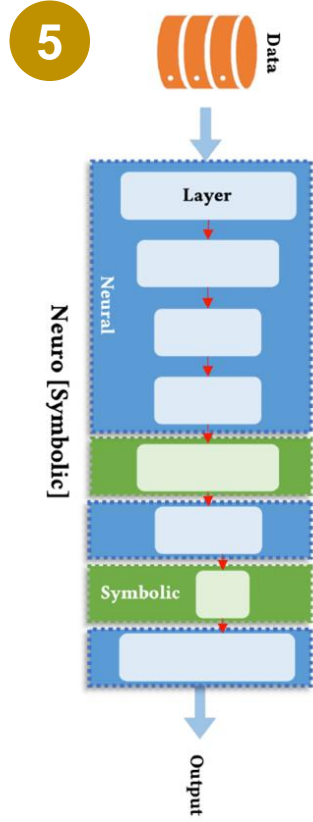
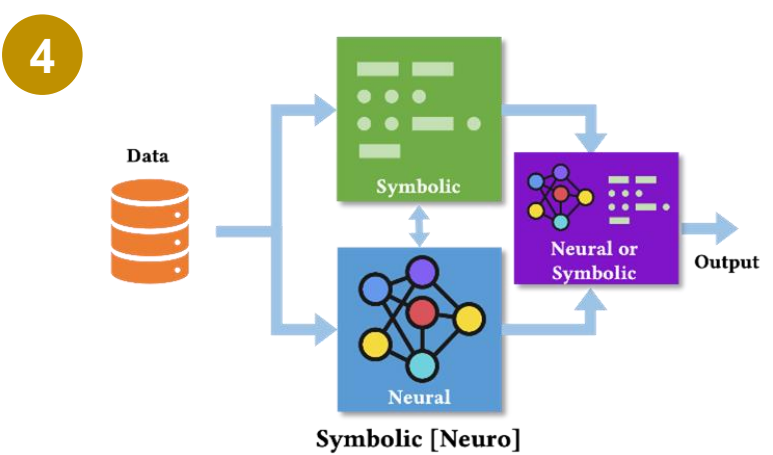
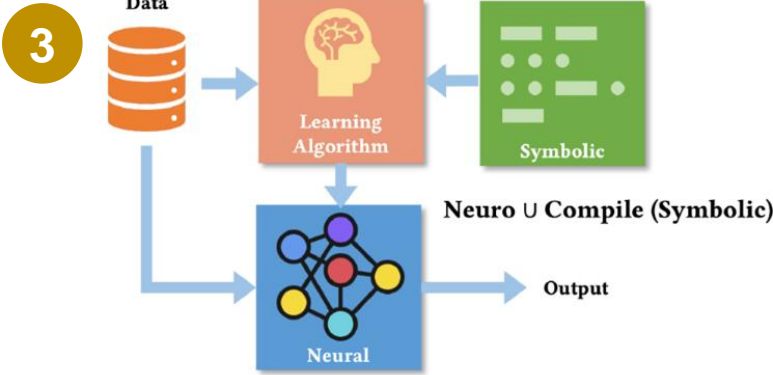
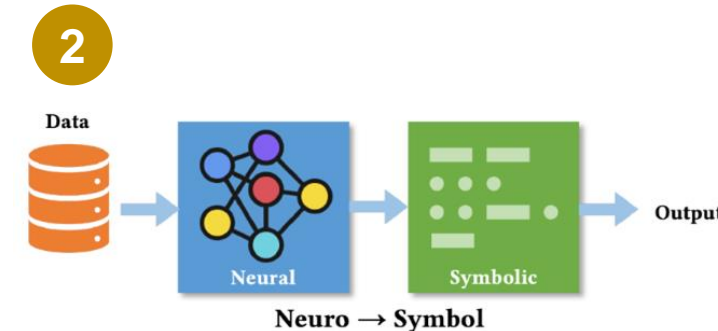
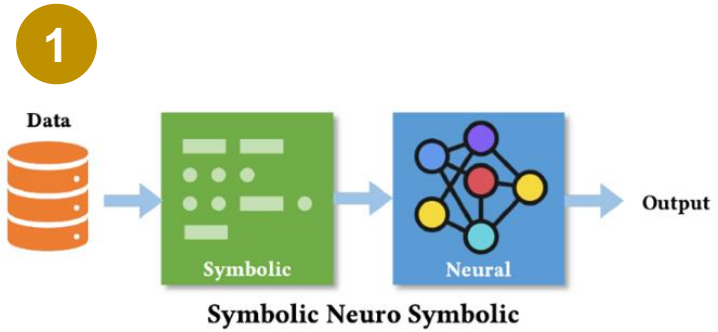
# MLPerf Tiny v0.5 inference benchmark – search space matters



**Anomaly Detection (ToyADMOS): exceeds benchmark accuracy by 4%, outperforms OutlierNets by 6.3%, guarantees deployability over MbNetv2**

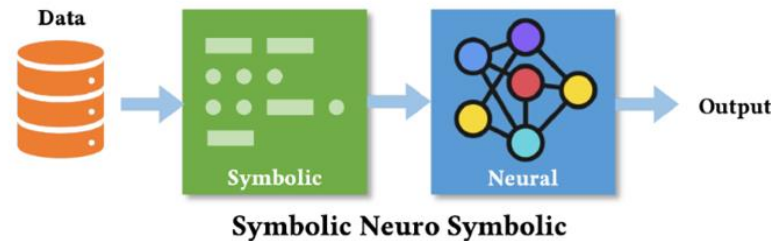
**Keyword spotting (Speech): incorrect backbone leads to suboptimal pareto-frontier, stressing importance of a search space containing several models**

# Neurosymbolic AI taxonomy



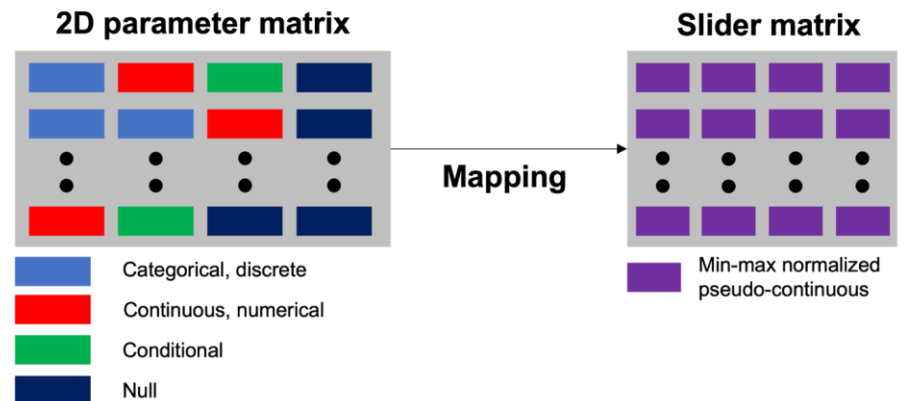
# Symbolic neurosymbolic – problem formulation

A series of independent domain-engineered functions is applied on the input dataset  $X$ , followed by a single ML model



Symbolic search space: 2D hyperparameter matrix, with each row corresponding to the arguments of each function (binary mask)

Neural search space: multiple model backbones, each considered for use at each step using an ordinal mask



# Example: co-optimizing features and a single backbone

Application: gesture recognition using a temporal CNN; operates on handcrafted features

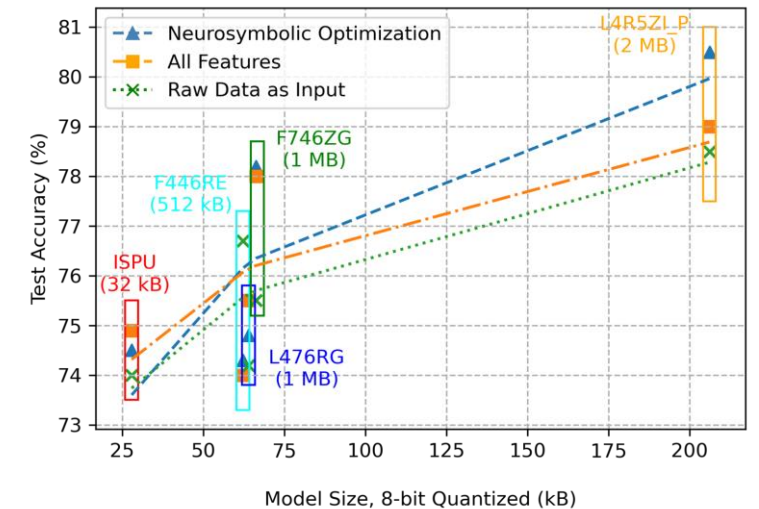
Features											
Mean	IQR	Maximum	Median	Variance	MAD	Abs. Energy	Entropy	Peak-to-Peak	FFT mean coeff.	Fundamental Frequency	Abs. Energy



# Example: co-optimizing features and a single backbone

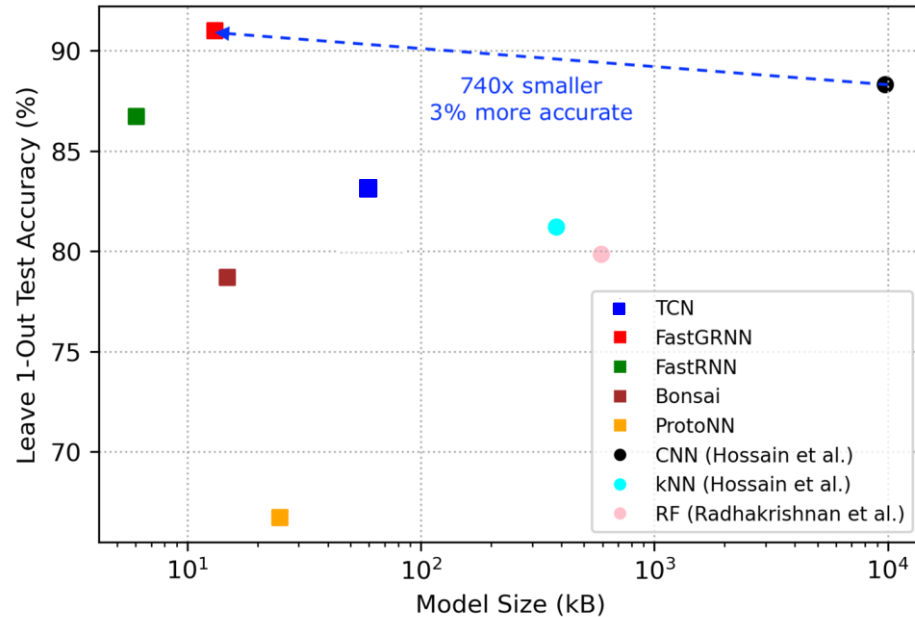
Application: gesture recognition using a temporal CNN; operates on handcrafted features

Microcontroller (SRAM, Flash)	Features											
	Mean	IQR	Maximum	Median	Variance	MAD	Abs. Energy	Entropy	Peak-to-Peak	FFT mean coeff.	Fundamental Frequency	Abs. Energy
ISPU (8,32)												
F446RE (128, 512)												
L476RG (128, 1024)												
F746ZG (320, 1024)												
L4R5ZI_P (640, 2048)												

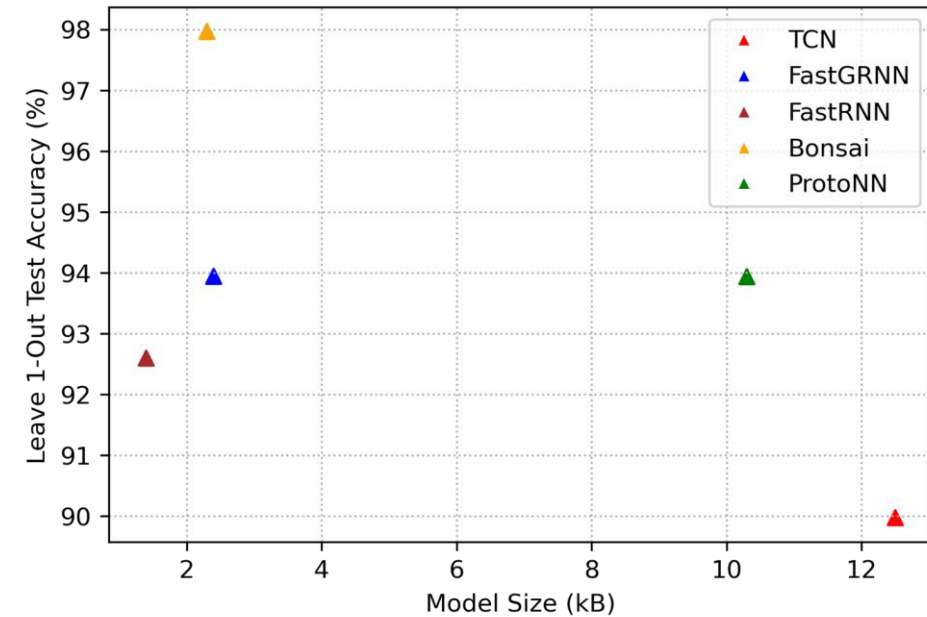


Extracting all features is computationally intensive. TinyNS picks the most important features when resources are scarce to maximize accuracy

# Example: optimizing over multiple backbones



Activity detection using earables



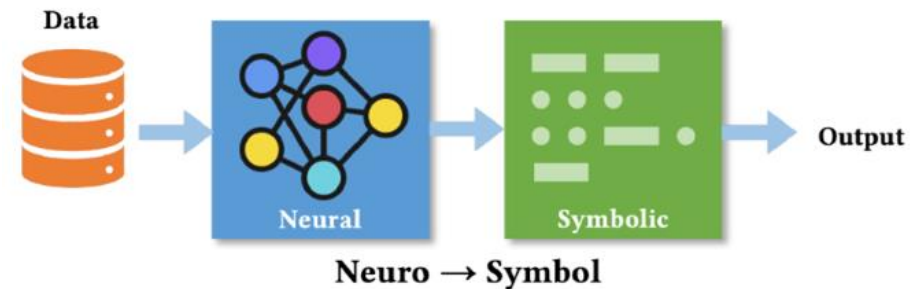
Fall detection using earables

**98-740x smaller, 3-6% more accurate models for human activity detection using earables over baselines**

**Activity detection under 6-12 KB of memory, fall detection under 2 KB of memory**

# Neuro → symbol – problem formulation

A single ML model operates on the input data, followed by either a single domain-engineered function or a program graph



Given a collection of logical, relational, arithmetic, and conditional operators, program decision trees can be synthesized conditioned upon a finite tree count and depth

Tree enumeration algorithm to generate all possible paths to Decision A and B; or optimize parameters of a pre-defined tree

# Example: co-optimizing neural detector and symbolic tracker

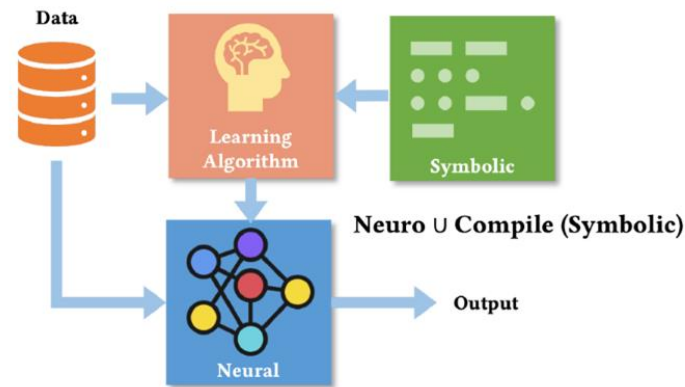


Constraint	Flash Usage (MB)	Performance		Neural hyperparameters				Symbolic hyperparameters	
		MOTA	IDF	Kernel Size	Stacks	Head Convolution	Activations	Rendering	Confidence
None (handcrafted)	238	36.5	55.0	1	1	128	True	0.4	0.5
250 MB limit	238	36.1	54.6	1	1	150	True	0.3	0.4
500 MB limit	270	38.0	57.2	9	1	100	False	0.7	0.5

**Achieves human-level performance ( $\pm 1\%$ ) of program hand-tuned using hundreds of human hours in 3 GPU days**

# Neuro U compile (symbolic) – problem formulation

Single ML model operates on the data, while the symbolic rules are expressed in two ways

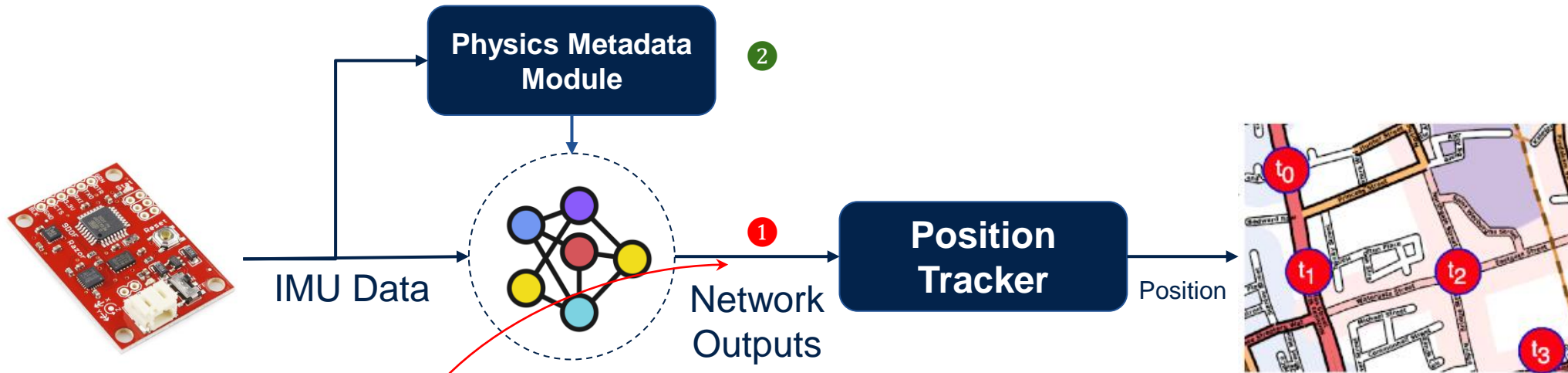


1. Add more regularizer terms in the NAS optimization function

2. Add physics metadata channel as additional inputs to the model

$$\min f_{\text{opt}}, \quad f_{\text{opt}} = \lambda_1 f_{\text{error}}(\Omega') + \lambda_2 f_{\text{flash}}(\Omega') + \lambda_3 f_{\text{SRAM}}(\Omega') + \lambda_4 f_{\text{latency}}(\Omega') + \lambda_5 f_{\text{rule 1}}(\Omega')$$

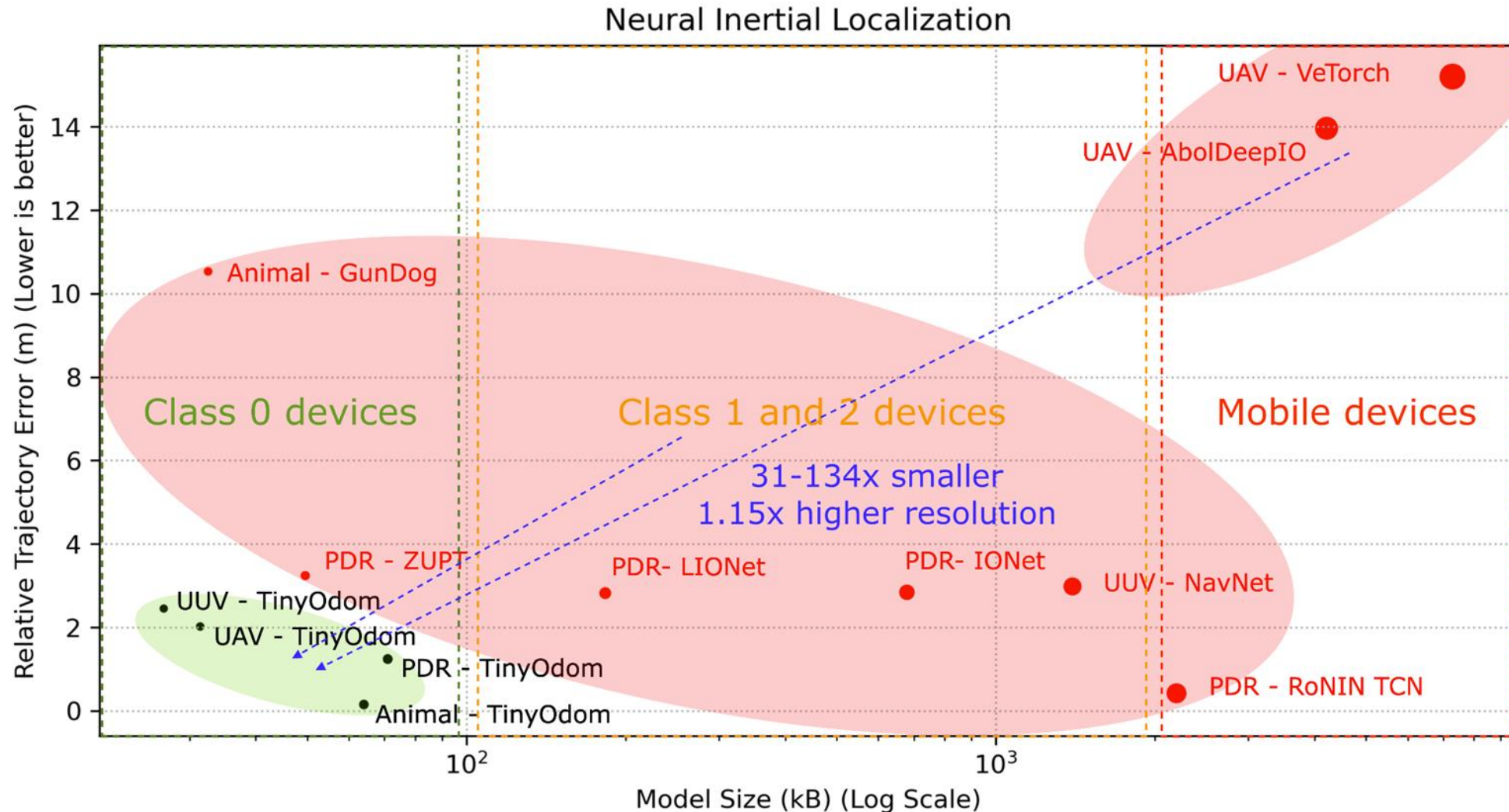
# Example: physics-aware neural inertial navigation



$$\underbrace{(v_{x,k}, v_{y,k})}_{\text{1}} = \gamma_{\theta^*}(\underbrace{\hat{\mathbf{a}}_{q:q+n}^I}_{\text{2}}, \hat{\mathbf{w}}_{q:q+n}^I, \underbrace{\hat{\mathbf{m}}_{q:q+n}^I}_{\text{2}}, \underbrace{c_k(I\hat{\mathbf{a}})}_{\text{2}})$$

- 1 Velocity and magneto-centric DNN regresses velocities and uses magnetic North as an additional anchor point.
- 2 A physics metadata module supplies latent information about whether valid translational movements have occurred.

# Example: physics-aware neural inertial navigation

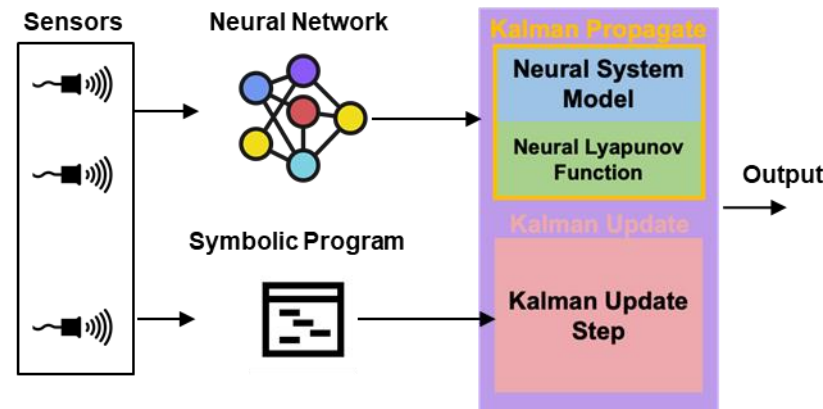


Saha, Swapnil Sayan, Sandeep Singh Sandha, Luis Antonio Garcia, and Mani Srivastava. "Tinyodom: Hardware-aware efficient neural inertial navigation." *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 6, no. 2 (2022): 1-32.

Saha, Swapnil Sayan, Yayun Du, Sandeep Singh Sandha, Luis Antonio Garcia, Mohammad Khalid Jawed, and Mani Srivastava. "Inertial Navigation on Extremely Resource-Constrained Platforms: Methods, Opportunities and Challenges." In *2023 IEEE/ION Position, Location and Navigation Symposium (PLANS)*, pp. 708-723. IEEE, 2023.

# Symbolic [neuro] – problem formulation

Use Kalman filter theory to combine a noisy neural system model with noisy symbolic measurement updates



Separate neural and non-neural parts in Kalman propagate. Neural network provides a black box mapping

Use the linearized Jacobian of the neural network w.r.t the past state and inputs in the Lyapunov function

$$\hat{\mathbf{x}}_{k+1|k} = \mathbf{A}\hat{\mathbf{x}}_k + g(\mathbf{u}_{k+1}), \quad g(\cdot) = f(y_\theta(\cdot))$$

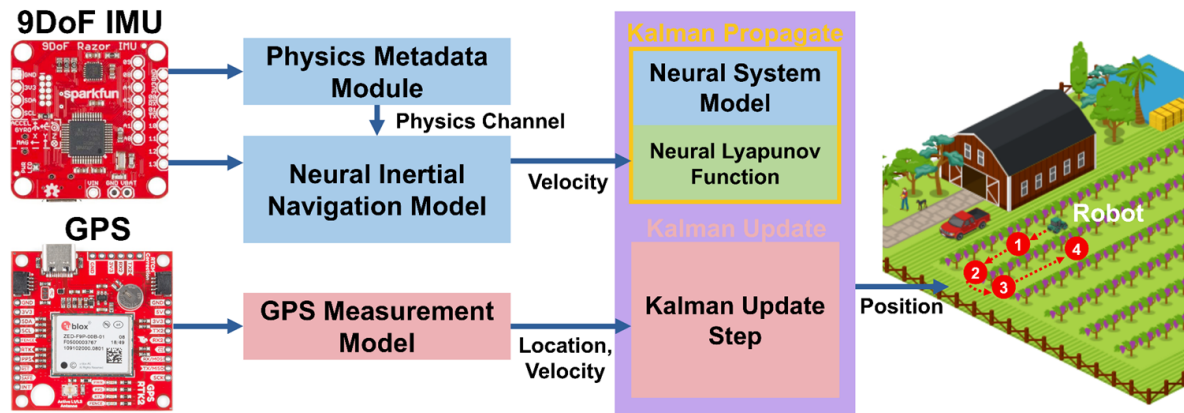
Linear state evolution (known)      Non-linear state evolution (known)

$$\mathbf{P}_{k+1|k} = \mathbf{A}\mathbf{P}_k\mathbf{A}^T + \mathbf{B}_{k+1}\mathbf{U}_k\mathbf{B}_{k+1}^T, \quad \mathbf{B}_{k+1} = \left. \frac{\partial g}{\partial \mathbf{u}} \right|_{\hat{\mathbf{x}}_k, \mathbf{u}_{k+1}}$$

Sensor Allan parameters      Jacobian term



# Example: neural-Kalman filtering



Method (1 Hz GPS)	Median Absolute Trajectory Error (m)	Median Relative Trajectory Error (m)
UKF-M GPS/INS	4.35	0.21
EKF GPS/INS	2.24	0.35
GPS only	1.89	0.40
<b>Neural-Kalman (ours)</b>	<b>1.36</b>	<b>0.35</b>

**Application: Tracking agricultural robots using neural inertial navigation and GPS; neural network provides a model-free evolution of the robot dynamics**

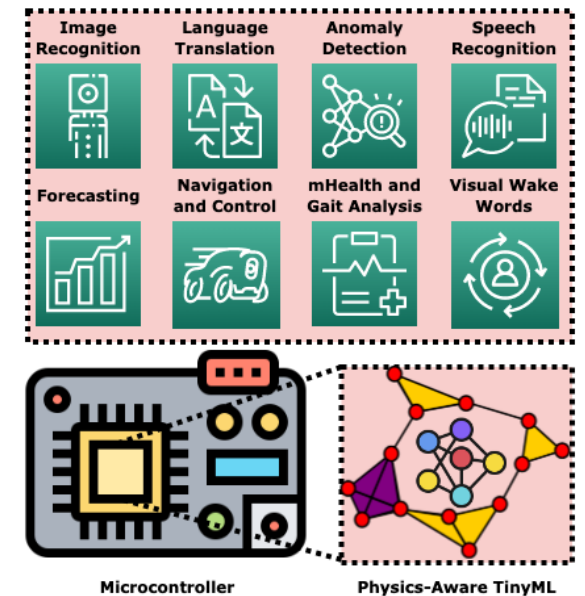
**Neural-Kalman filter combines smoothness and short-term accuracy of neural networks with long-term precision of noisy GPS/GNSS updates under 1 MB of memory**

# Conclusion

Neurosymbolic tiny machine learning enables context-aware, physics-aware, robust, interpretable, and performant edgeAI systems

TinyNS automates the process of generating neurosymbolic programs for TinyML platforms

Enables a broad spectrum of new applications for wearables, robots, automobiles, and environmental sensors

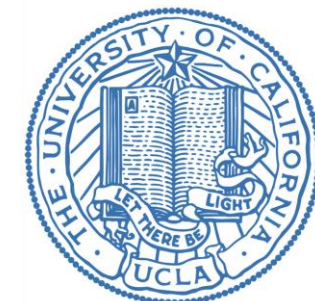




life.augmented

Try TinyNS:

<https://github.com/nesl/neurosymbolic-tinyml>



# Our technology starts with You

© STMicroelectronics - All rights reserved.

ST logo is a trademark or a registered trademark of STMicroelectronics International NV or its affiliates in the EU and/or other countries.

For additional information about ST trademarks, please refer to [www.st.com/trademarks](http://www.st.com/trademarks).

All other product or service names are the property of their respective owners.



life.augmented



# Copyright Notice

This multimedia file is copyright © 2024 by tinyML Foundation. All rights reserved. It may not be duplicated or distributed in any form without prior written approval.

tinyML<sup>®</sup> is a registered trademark of the tinyML Foundation.

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



# 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)**