

tinyML[®] Talks

Enabling Ultra-low Power Machine Learning at the Edge

“On-Device Domain Adaptation for Noise-Robust Keyword Spotting”

Cristian Cioflan – Doctoral Student, ETH Zürich

November 16, 2023



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

T I N Y



TALKS
webcast

Platinum Strategic Partner



**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

Decarbonization

Digitalization



Driving decarbonization and digitalization. Together.

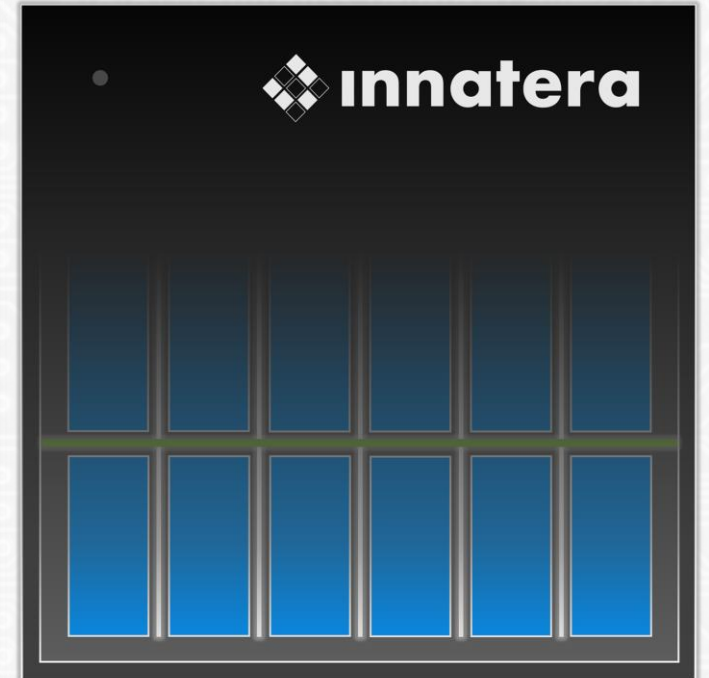
Infineon serving all target markets as
Leader in Power Systems and IoT

www.infineon.com



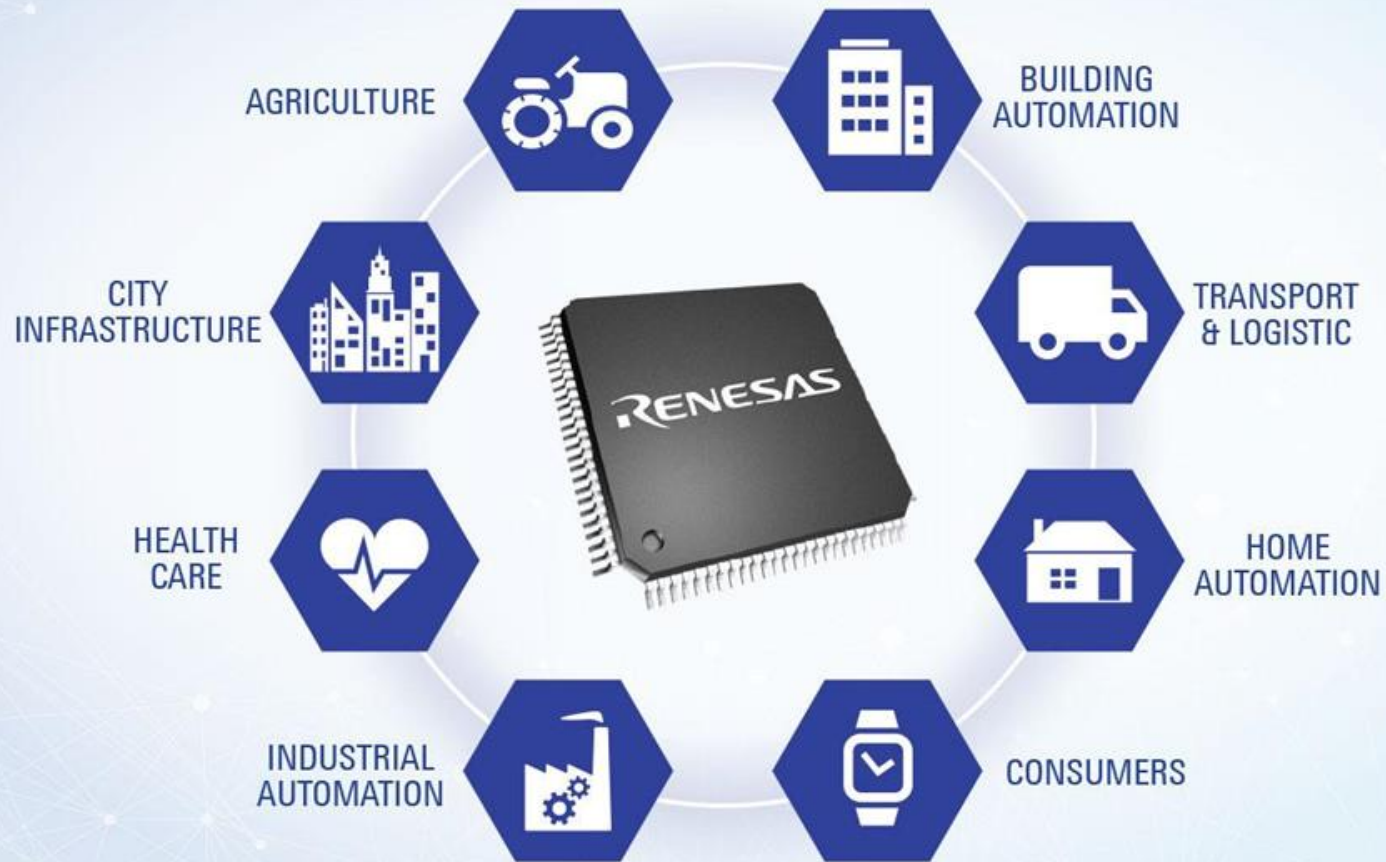


NEUROMORPHIC INTELLIGENCE FOR THE SENSOR-EDGE



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



T I N Y



Silver Strategic Partners



brainchip



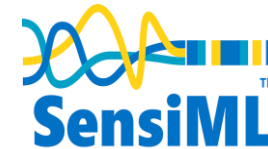
GREENWAVES
TECHNOLOGIES



£Grovety Inc.



NotaAI





Join Growing tinyML Communities:



17.6k 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/>



4k members
&
13k 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



tinyML
4.33K subscribers

10.9k subscribers, 633 videos with 391k views

HOME VIDEOS PLAYLISTS COMMUNITY CHANNELS ABOUT

13:24	33:27	32:39	36:41	34:03	34:58
On Device Learning Forum - Professors...	On Device Learning - Manuel Roveri: Is on-...	On Device Learning Forum - Warren Gros...	On Device Learning Forum - Yiran Chen...	On Device Learning Forum - Hiroku...	On Device Learning Forum - Song Han: O...
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
1:13	1:07:43	53:41	45:46	51:01	1:03:24
tinyML Smart Weather Station Challenge - ...	tinyML Talks Singapore...	tinyML Talks Shenzhen: Data...	tinyML Talks Singapore...	tinyML Smart Weather Station with Syntiant...	tinyML Trailblazers August with Vijay...
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
58:50	34:36	55:01	59:51	59:48	58:09
tinyML Auto ML Tutorial with SensiML	tinyML Auto ML Tutorial with Qeexo	tinyML Talks Germany: Neural network...	tinyML Trailblazers with Yoram Zylberberg	tinyML Auto ML Tutorial with Nota AI	tinyML Auto ML Tutorial with Neuton
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
1:02:30	34:31	1:00:30	1:06:44	1:53:07	42:13
tinyML Challenge 2022: Smart weather...	tinyML Talks South Africa - What is...	tinyML Talks: The new Neuromorphic Anal...	tinyML Talks Shenzhen: 分享主题...	tinyML Auto ML Forum - Paneldiscussion	tinyML Auto ML Forum - Demos
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

tinyAI Forum on PdM & Anomaly Detection 2023



Interactive live webinar December 5, 2023 at 8AM Pacific Time
Registration is free of charge

tinyML Research Symposium

April 22, 2023

Call for Papers



Research Symposium - April 22, 2024

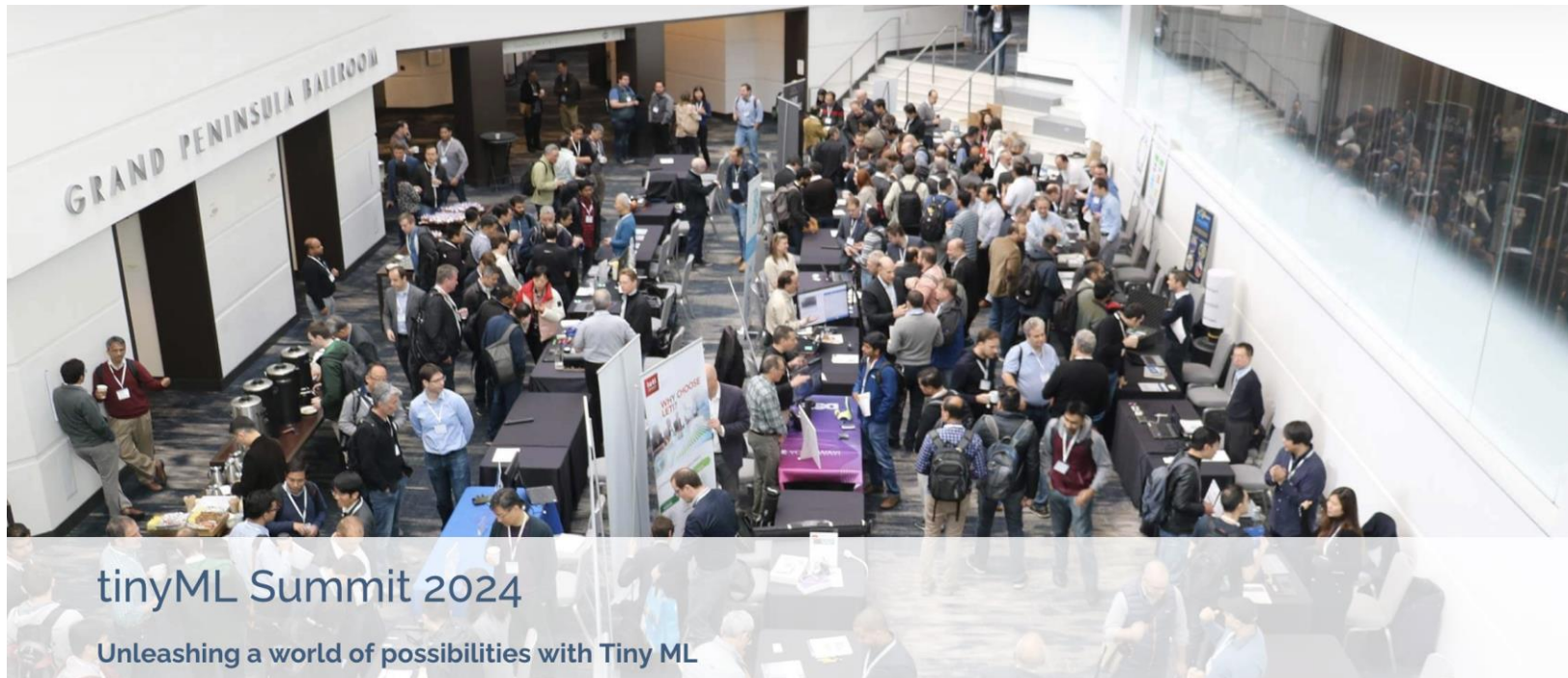
The tinyML research symposium serves as a flagship venue for related research at the intersection of machine learning applications, algorithms, software, and hardware in deeply embedded machine learning systems.

[Call for Papers](#)



tinyML Summit April 23-24, 2024

Call for Presentations and Posters



2023 Edge AI Technology Report

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



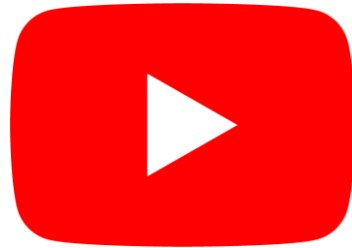


Reminders

Slides & Videos will be posted tomorrow



tinyml.org/forums



youtube.com/tinyml



Please use the Q&A window for your questions





Cristian Cioflan



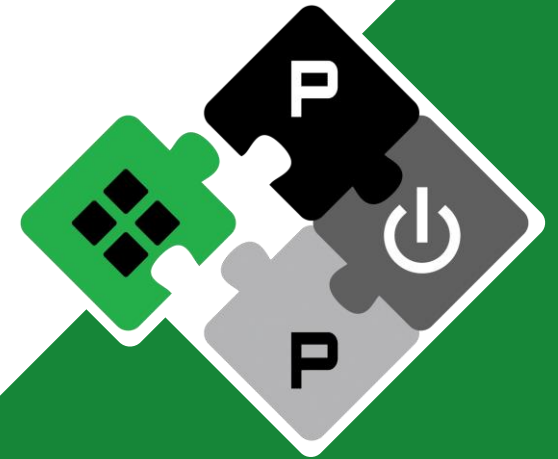
Cristian Cioflan received his B.Sc. degree in Electrical Engineering and Information Technology from University Politehnica of Bucharest in 2018 and his M.Sc. degree from the Swiss Federal Institute of Technology Zürich (ETHZ) in 2020. He is currently pursuing a Ph.D. in the Digital Circuits and Systems group of Prof. Luca Benini. His research interests include audio processing in low-power embedded systems, on-device continual learning, and neural architecture search for energy-efficient learning.

On-Device Domain Adaptation for Noise-Robust Keyword Spotting

Integrated Systems Laboratory (ETH Zürich)

PULP Platform

Open Source Hardware, the way it should be!



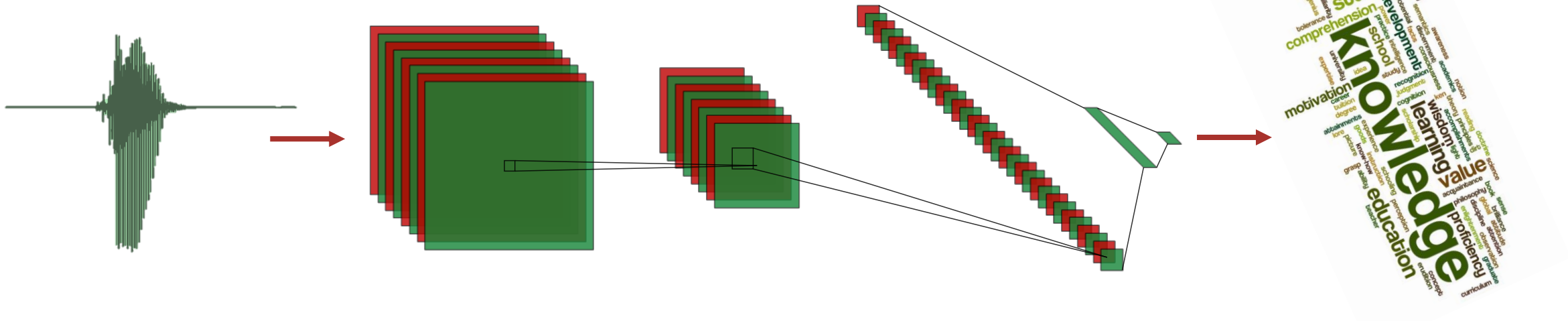
@pulp_platform 

pulp-platform.org 

youtube.com/pulp_platform 

Keyword Spotting – a definition

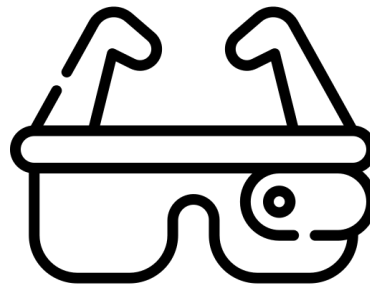
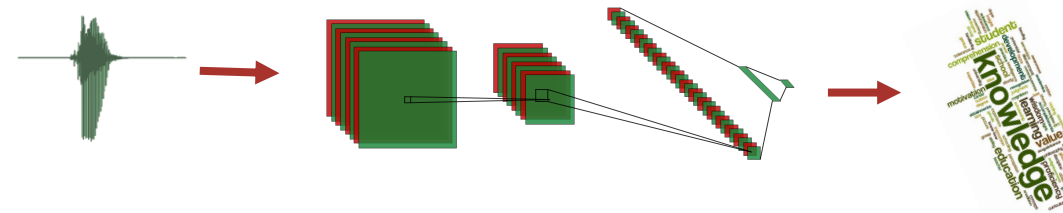
- **Keyword spotting – AI/ML branch**
 - Process an audio signal
 - Recognize a **target** word from a **predefined set**
 - (to **control** embedded devices)



Keyword Spotting – a definition

- **Keyword spotting – AI/ML branch**

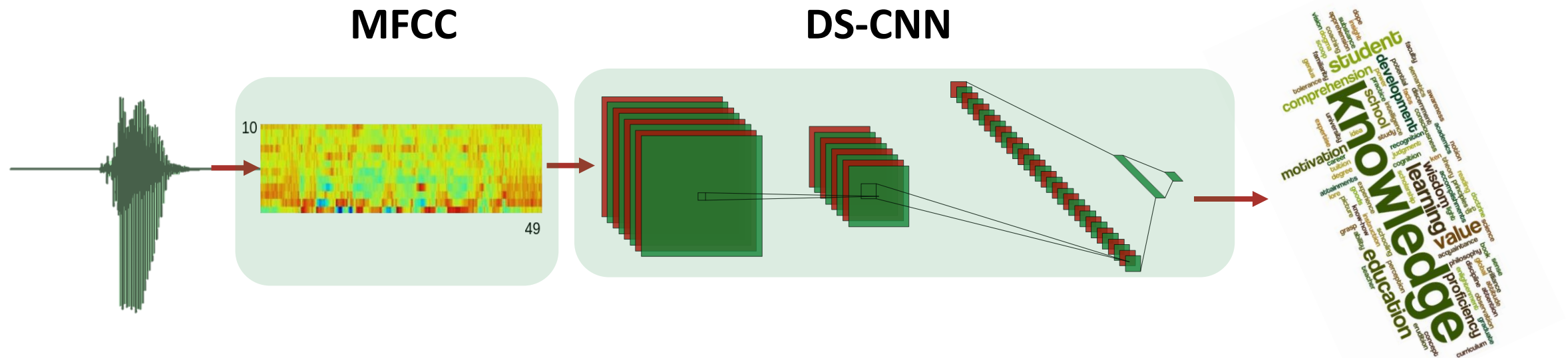
- **Process** an audio signal
- **Recognize** a **target** word from a **predefined set**
- (to **control** embedded devices)



Keyword Spotting – a definition



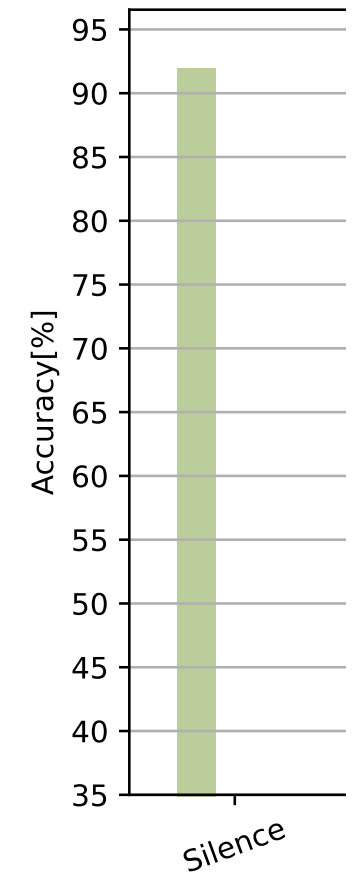
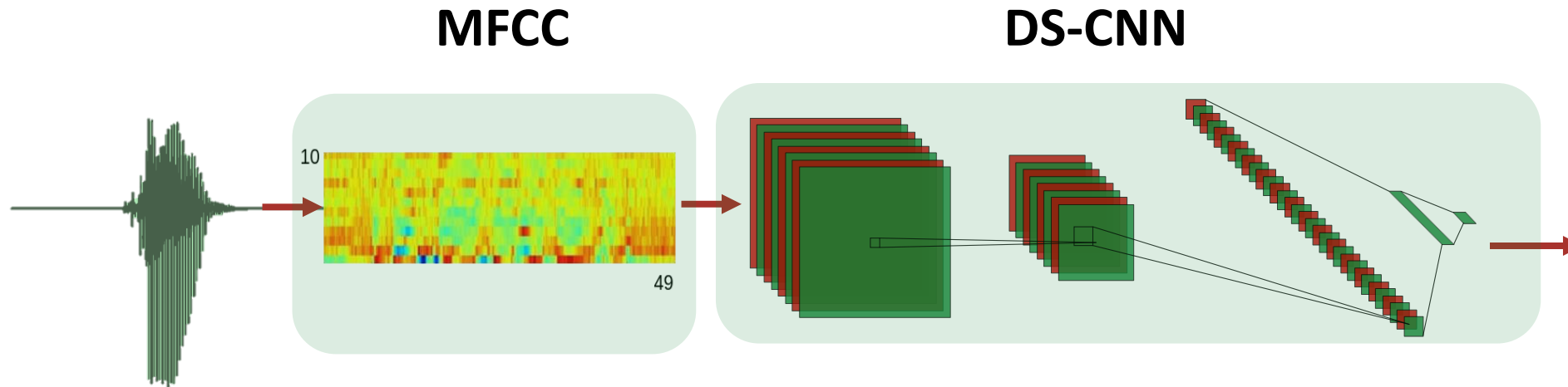
- **Keyword spotting – AI/ML branch**
 - **Process** an audio signal
 - **Recognize** a **target** word from a **predefined** set
 - (to **control** embedded devices)



92% accuracy in clean conditions



- **Keyword spotting – AI/ML branch**
 - **Process** an audio signal
 - **Recognize** a **target** word from a **predefined set**
 - (to **control** embedded devices)



The problem – noisy environments



The problem – noisy environments

- Additive background noise pollutes our signal
- Negative Signal-to-Noise ratios are common
 - -10 dB

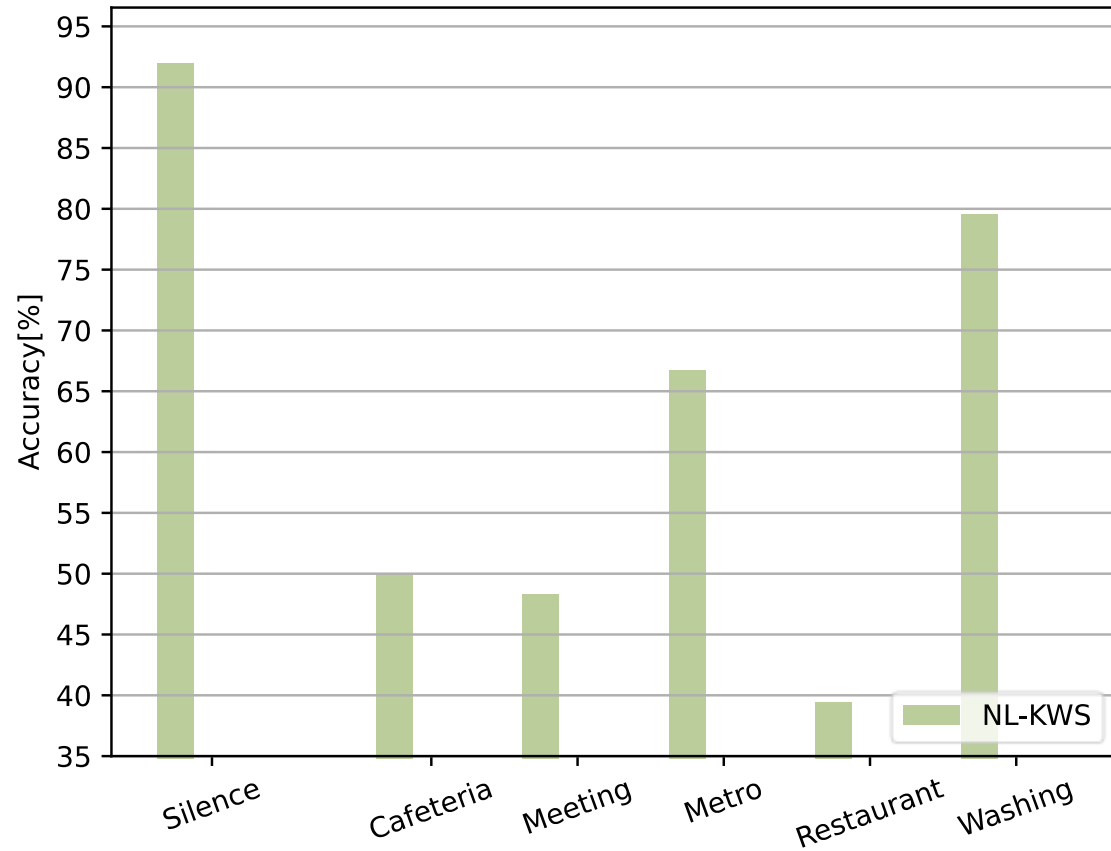
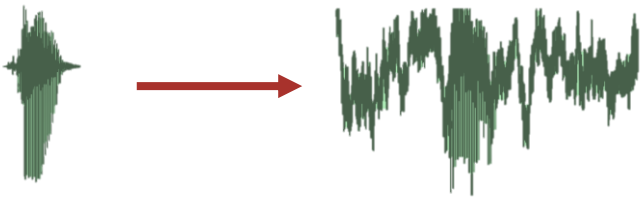
- Let us assume a Noiseless Keyword Spotting model
- Can a NL-KWS system still recognize the word?



Can a KWS system still recognize the words?



- Accuracy drops between 12% and 52%
- Stationary vs non-stationary noises
- Speech noise
 - @ 0 dB – 43% accuracy loss
 - Hard to separate target from noise

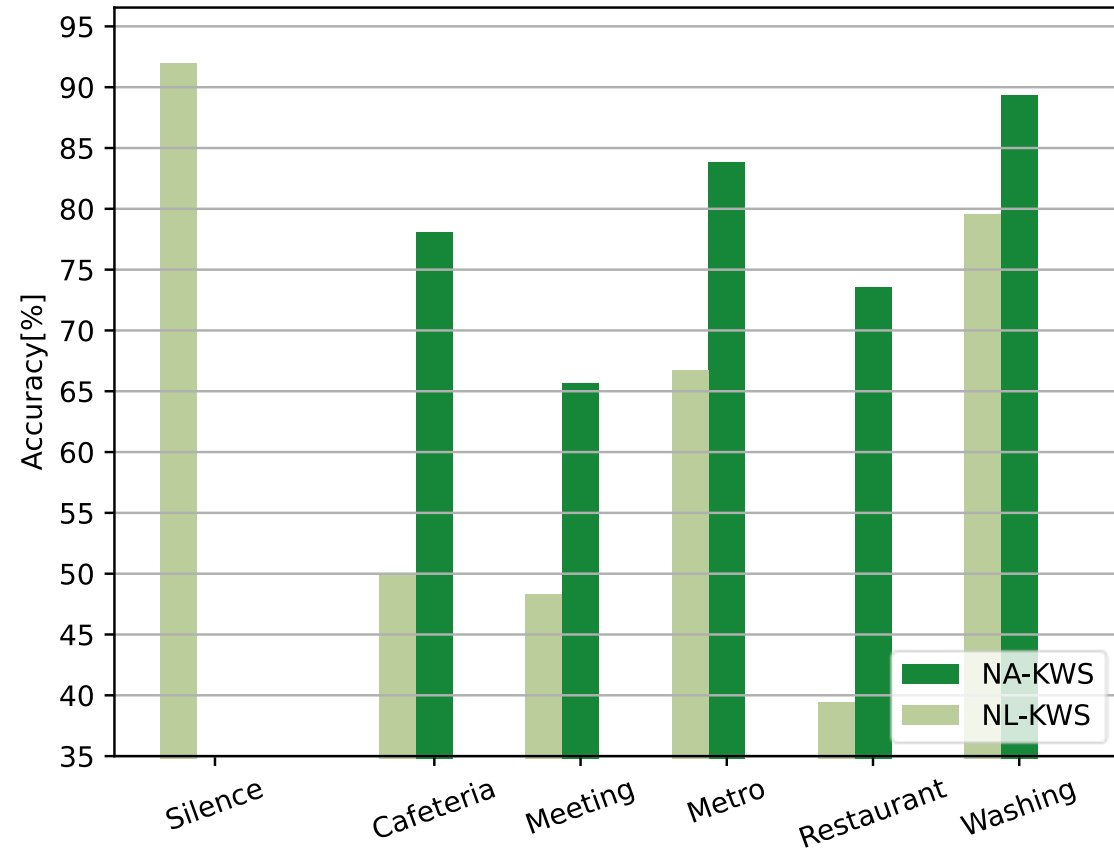
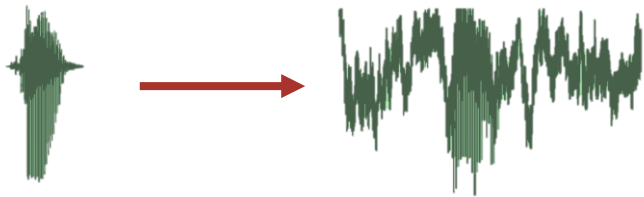


Noise-Aware Keyword Spotting



- **NL-KWS vs NA-KWS**

- Augment training samples with diverse noise types (and SNRs...)
- Improved results
 - Between **10%** and **34%** over NL-KWS

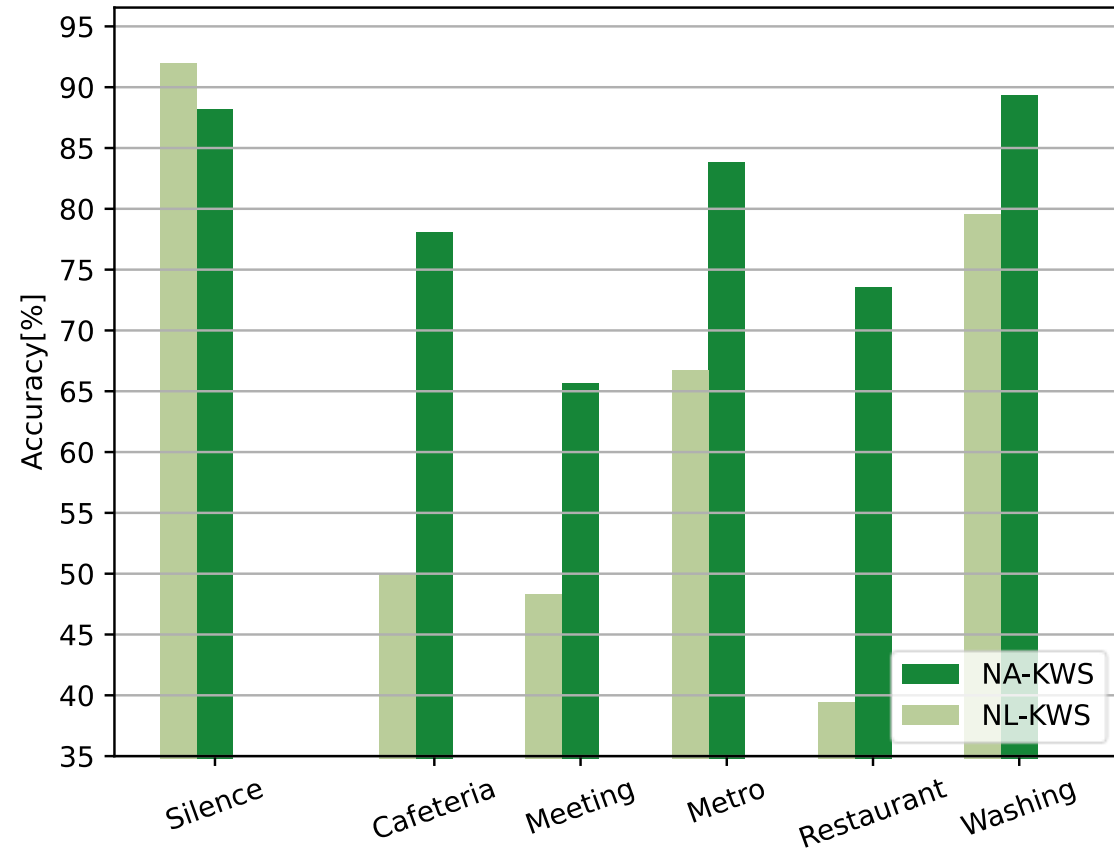
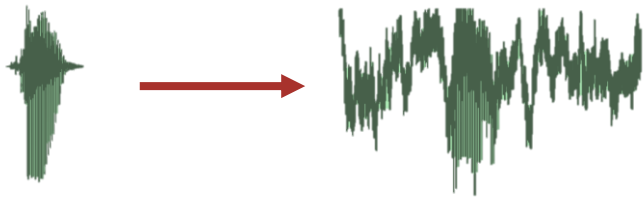


Noise-Aware Keyword Spotting



- **NL-KWS vs NA-KWS**

- Augment training samples with diverse noise types (and SNRs...)
- Improved results
 - Between **10%** and **34%** over NL-KWS
 - **-4%** on *Silence*



Can we improve KWS accuracy
through direct noise exposure?

Can we improve KWS accuracy
through direct noise exposure?

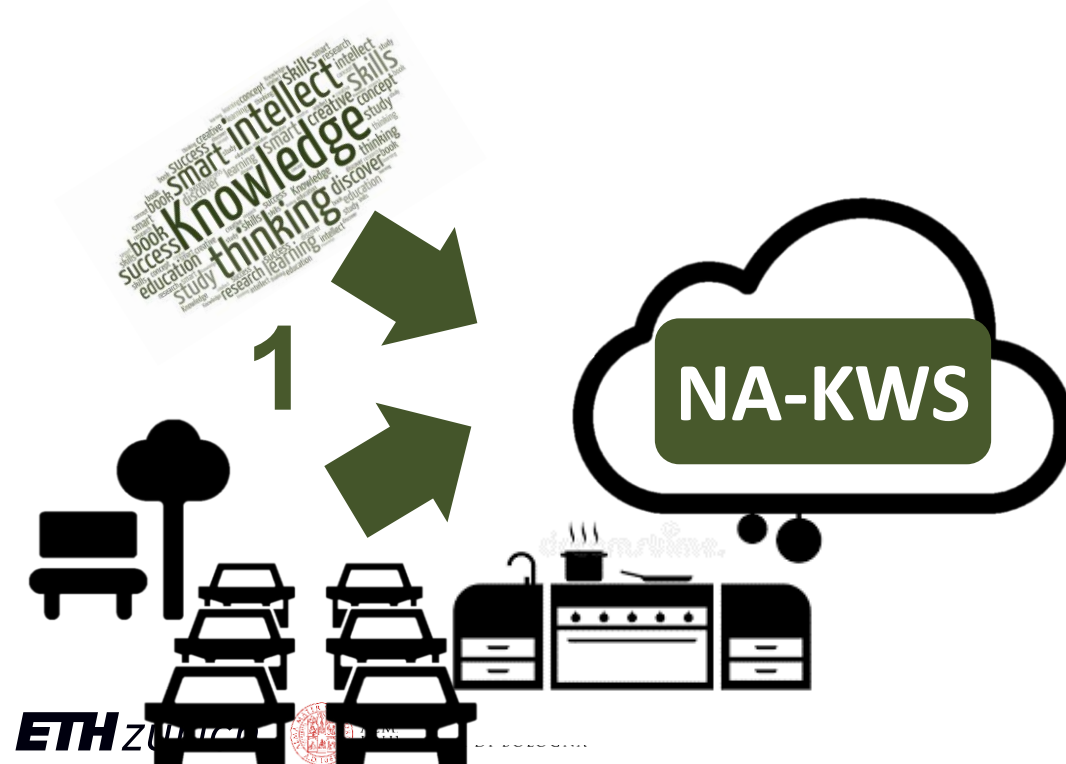
On-Device Domain Adaptation

On-Device Domain Adaptation – the methodology



1. Noise-Aware Keyword Spotting

- Data augmentation for improved generalization
- Google Speech Commands v2 x DEMAND



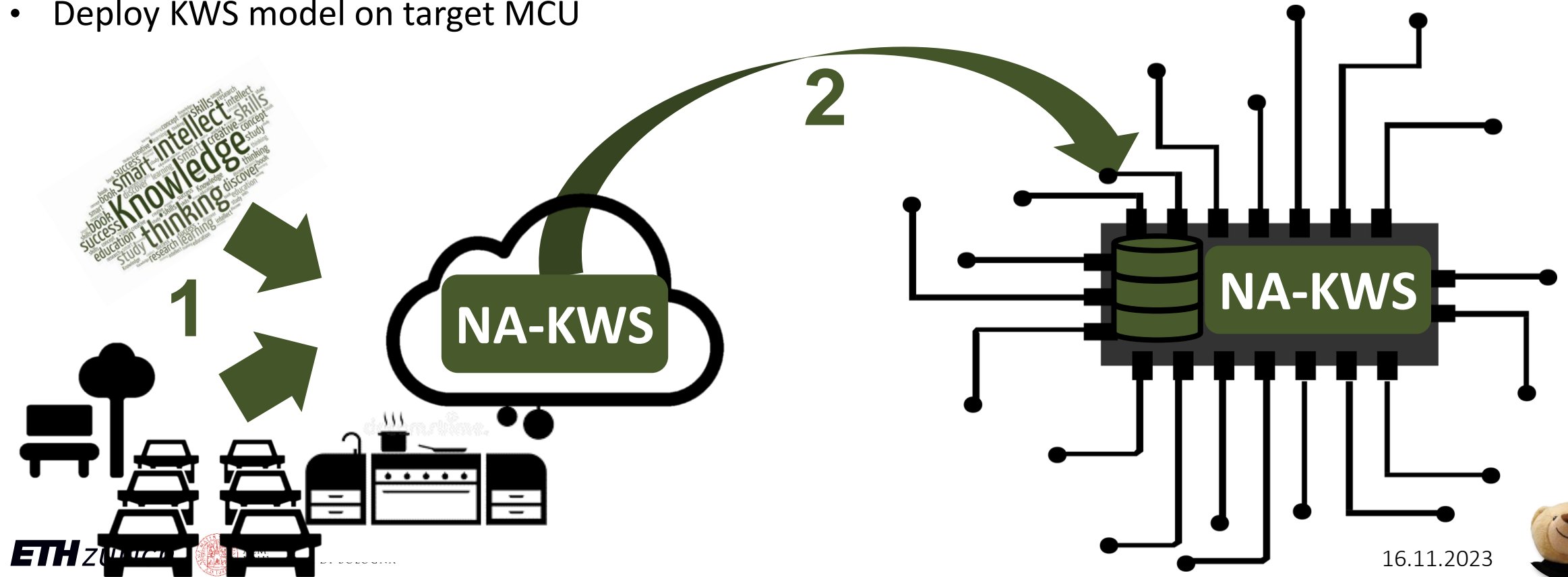
On-Device Domain Adaptation – the methodology



1. Noise-Aware Keyword Spotting

2. Deployment

- Store pre-recorded utterances
- Deploy KWS model on target MCU



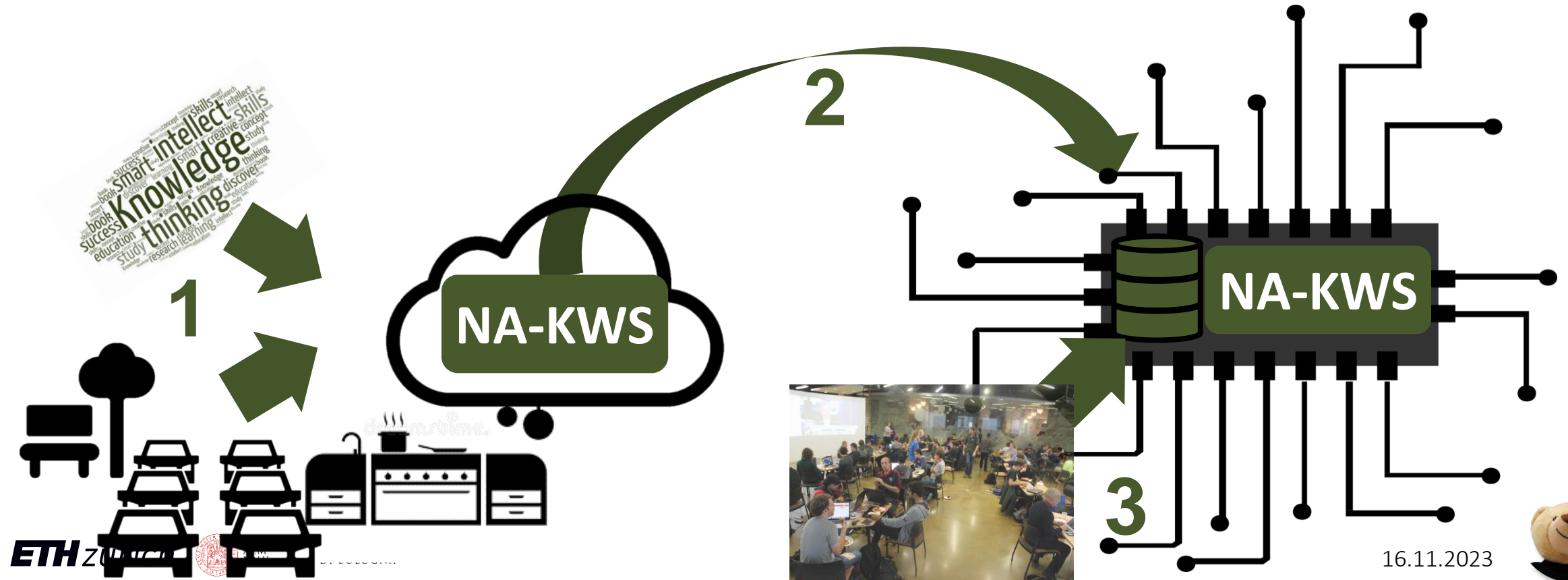
On-Device Domain Adaptation – the methodology



1. Noise-Aware Keyword Spotting
2. Deployment

3. Noise exposure

- Record on-site noise
- Augment pre-recorded samples



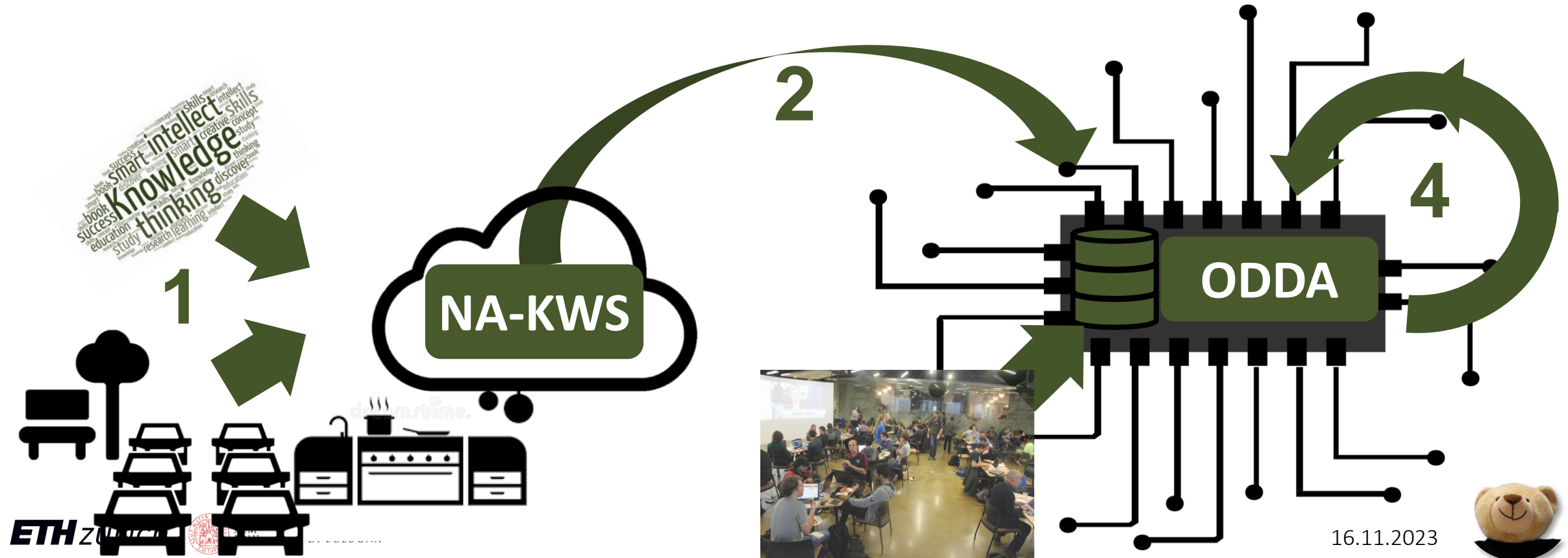
On-Device Domain Adaptation – the methodology



1. Noise-Aware Keyword Spotting
2. Deployment

3. Noise exposure
4. ODDA [1]

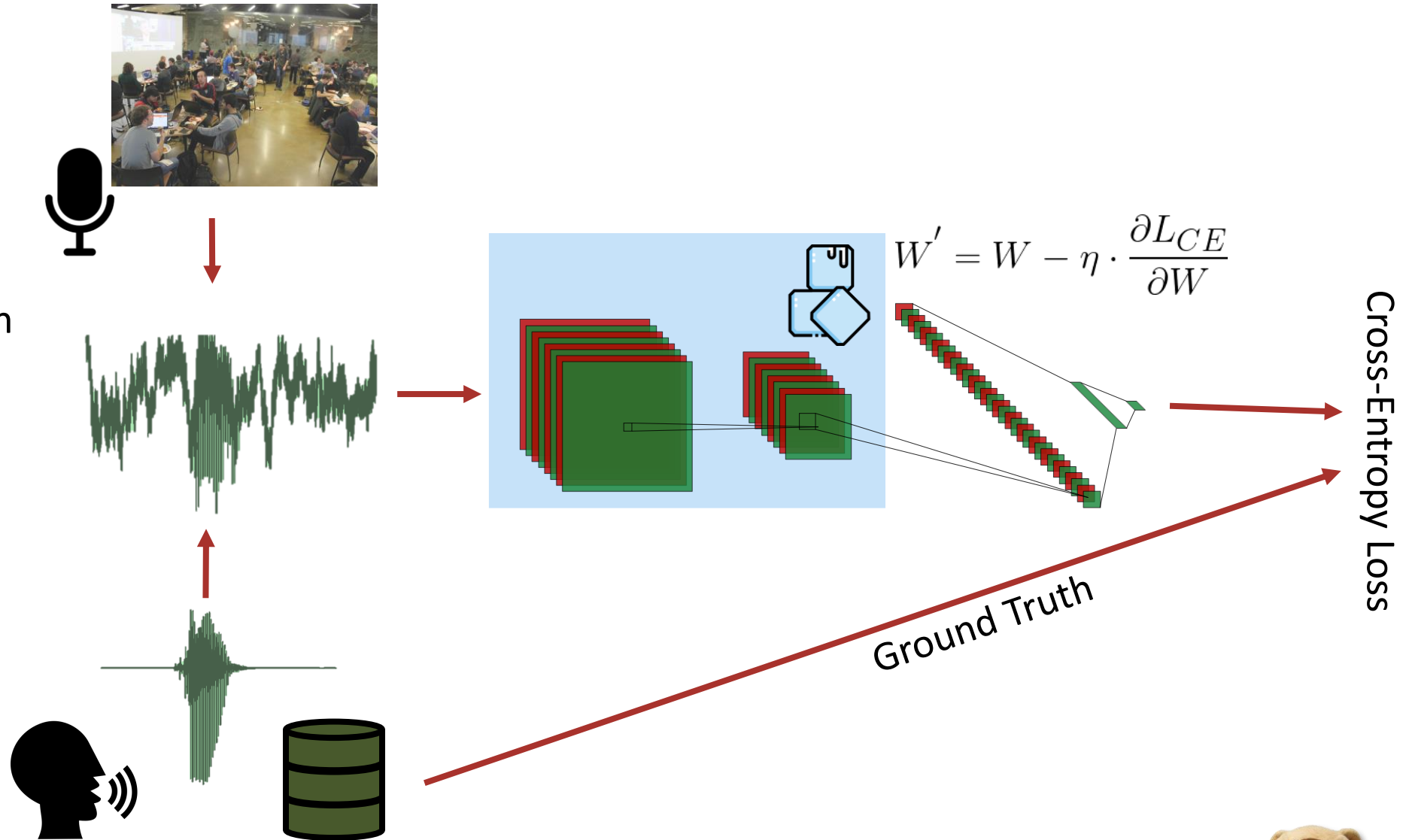
- Refine NA-KWS using augmented utterances



4. ODDA – Refine NA-KWS using augmented utterances



- a. Combine additive background noise with pre-recorded utterances
- b. Compute CE loss using partially frozen model prediction
- c. Selectively update model parameters



Exposure to environmental noise for specialization



On-Device Domain Adaptation

1. **Pretrain** a Noise-Aware Keyword Spotting network
2. **Deploy** NA-KWS (and clean utter.) on embedded device
3. **Augment** utterances with recorded environmental noise
4. **Refine** the network's parameters on the target platform
5. **Recognize** utterances with domain-adapted KWS system



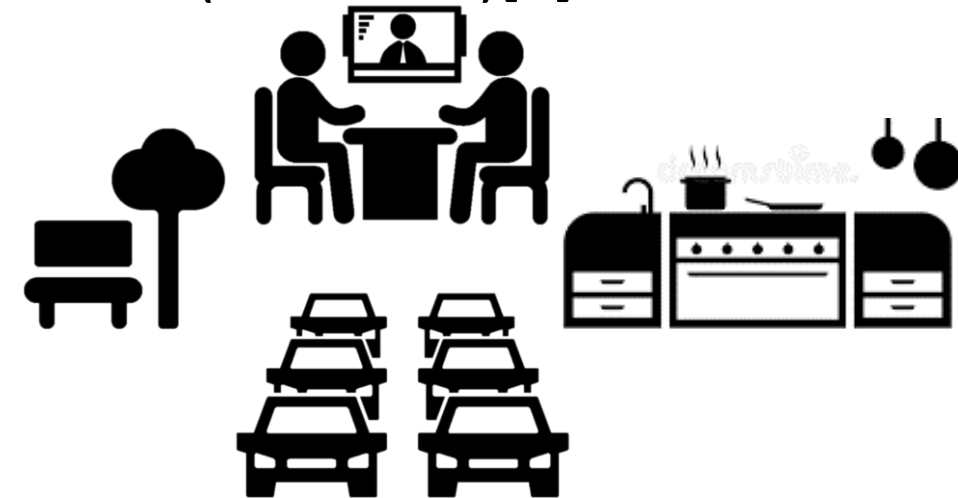


Google Speech Commands v2 (GSCv2)[2]



- 1-second audio @ 16 kHz sample rate
- silence & unknown – 10%
- NA-KWS/ODDA:validation:test – 80:10:10

Diverse Environments Multichannel Acoustic Noise Database (DEMAND)[3]



- 18 noises in real-world conditions
- 5 LOO settings
 - cafeteria, meeting, metro, restaurant, washing
- SNR: {-10, 10} dB



How well does ODDA perform?

ODDA – Qualitative analysis

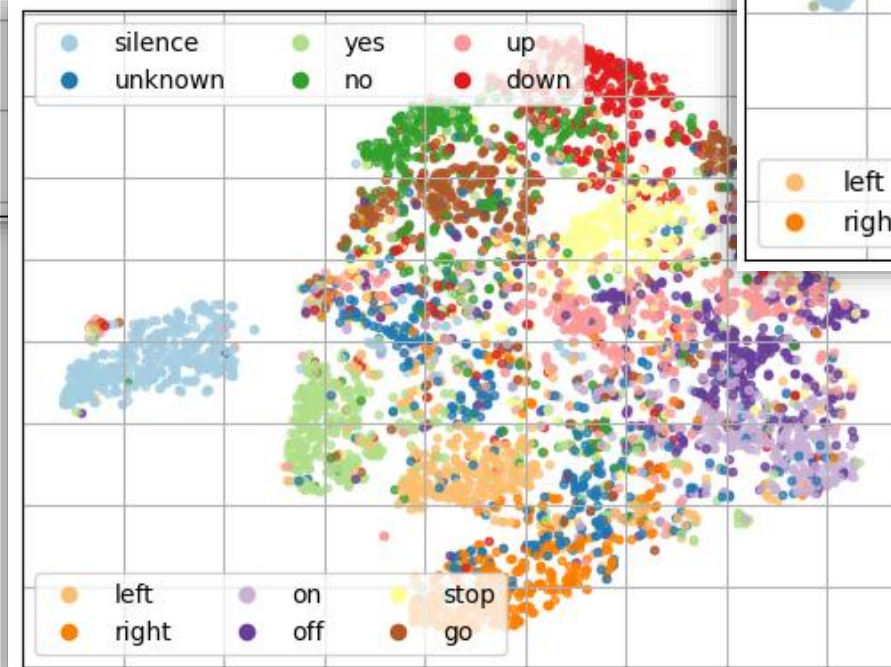
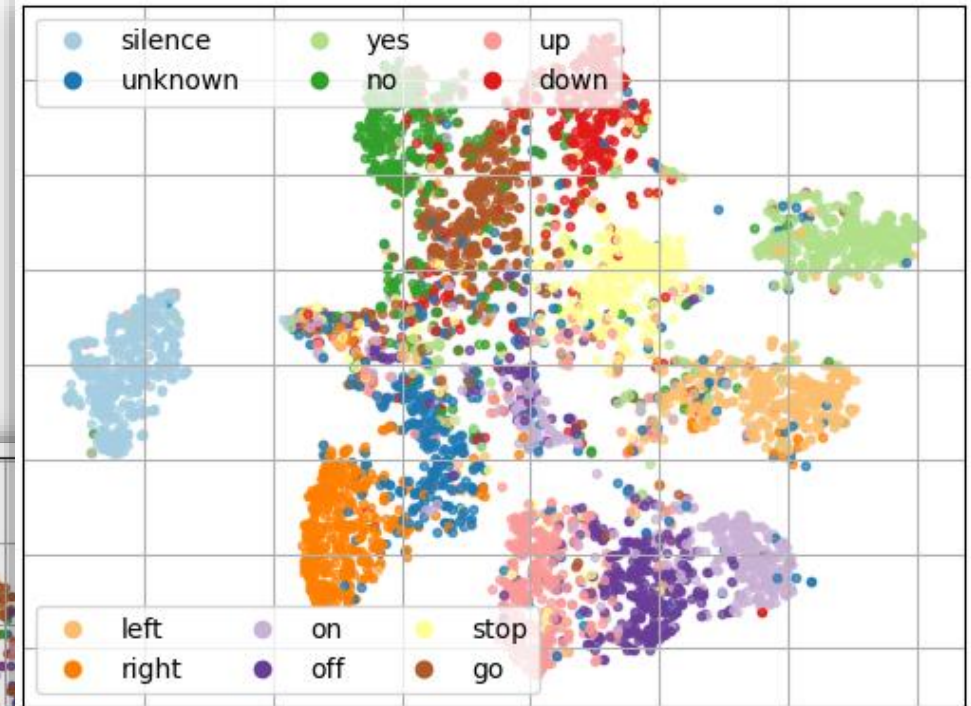
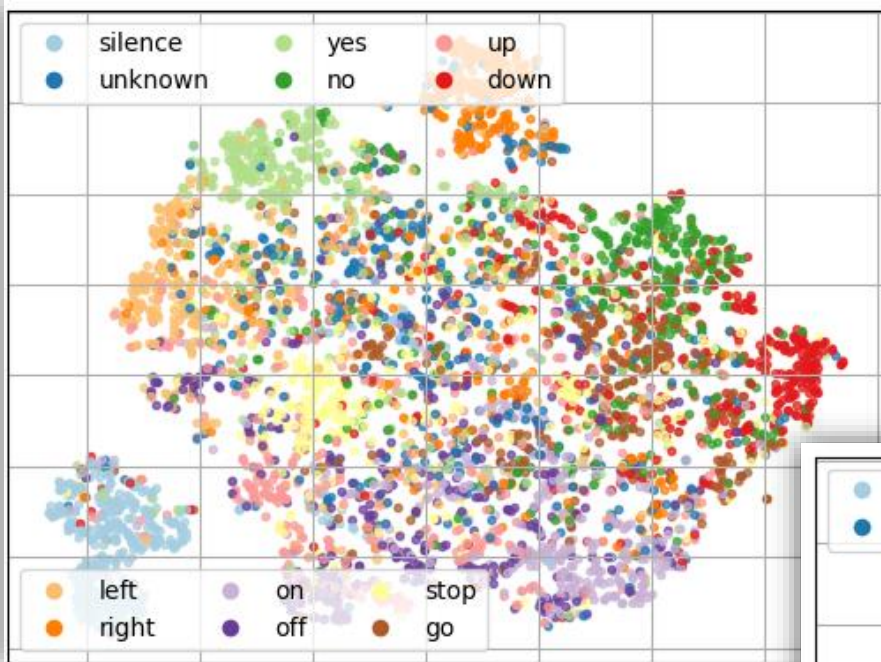


NL-KWS

ODDA

- t-SNE visualization of pre-classifier features
- *meeting* noise

NA-KWS



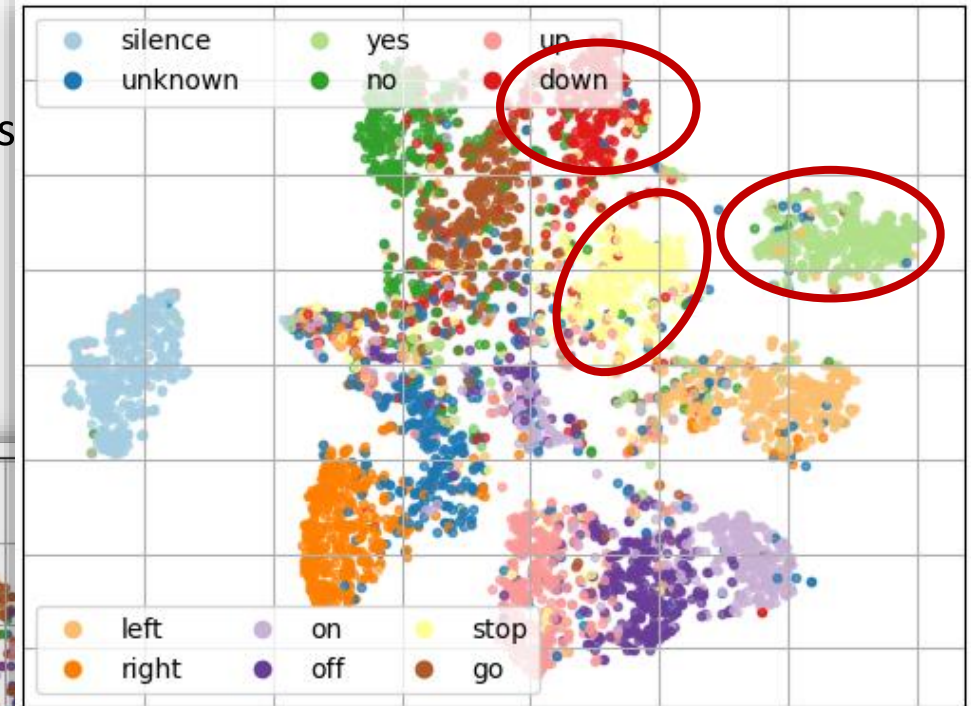
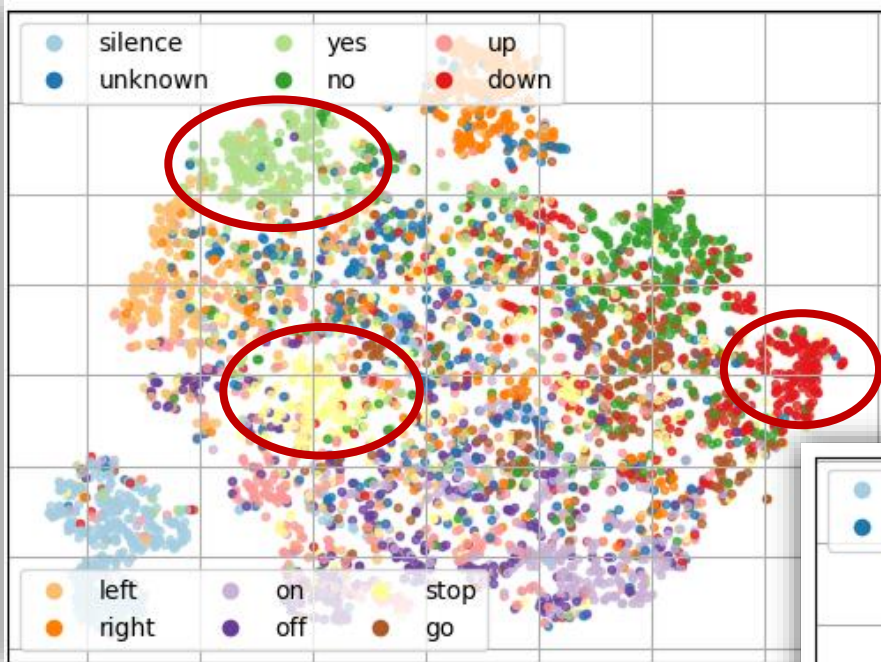
Exposure to env. noise → feature separation



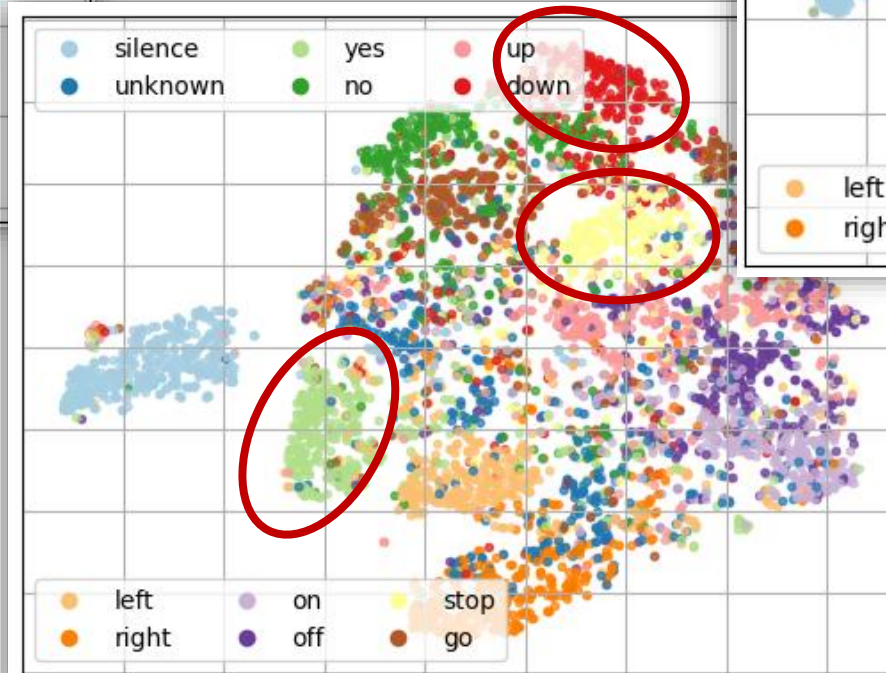
NL-KWS

ODDA

- *Yes, stop, down*
- **Marginal** improvements



NA-KWS



Exposure to env. noise → feature separation

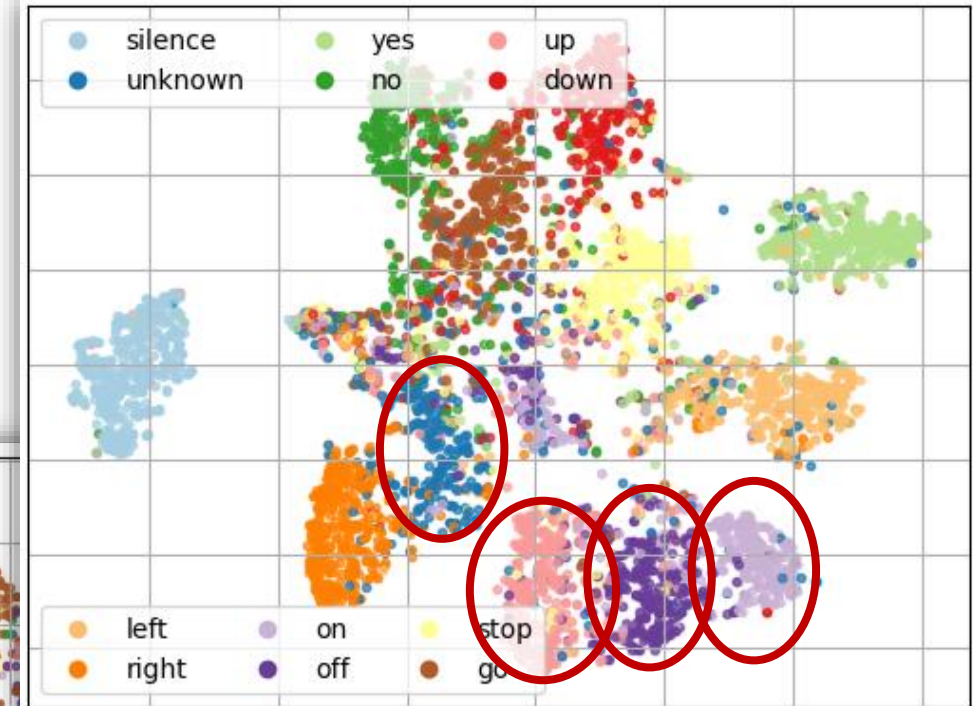
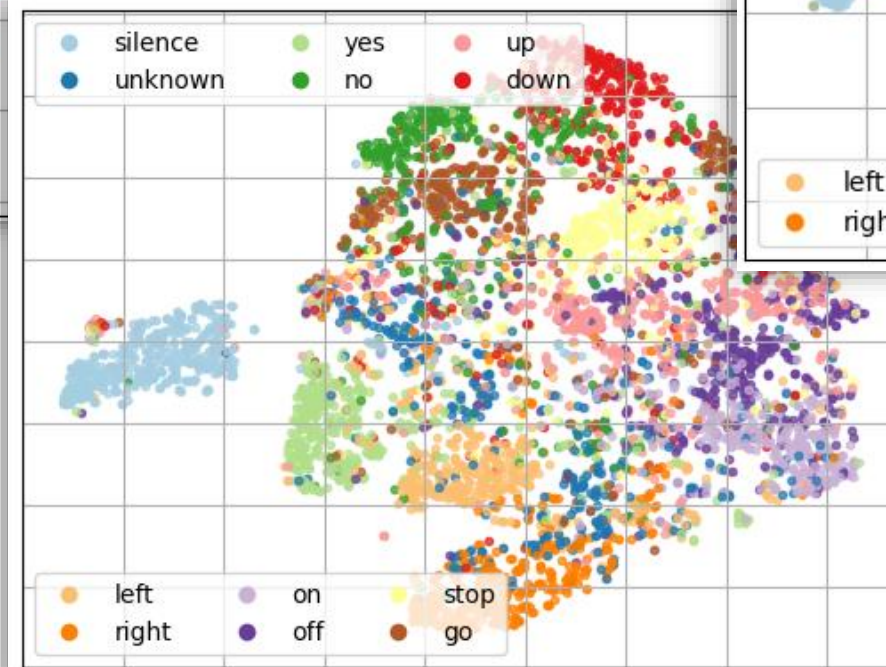
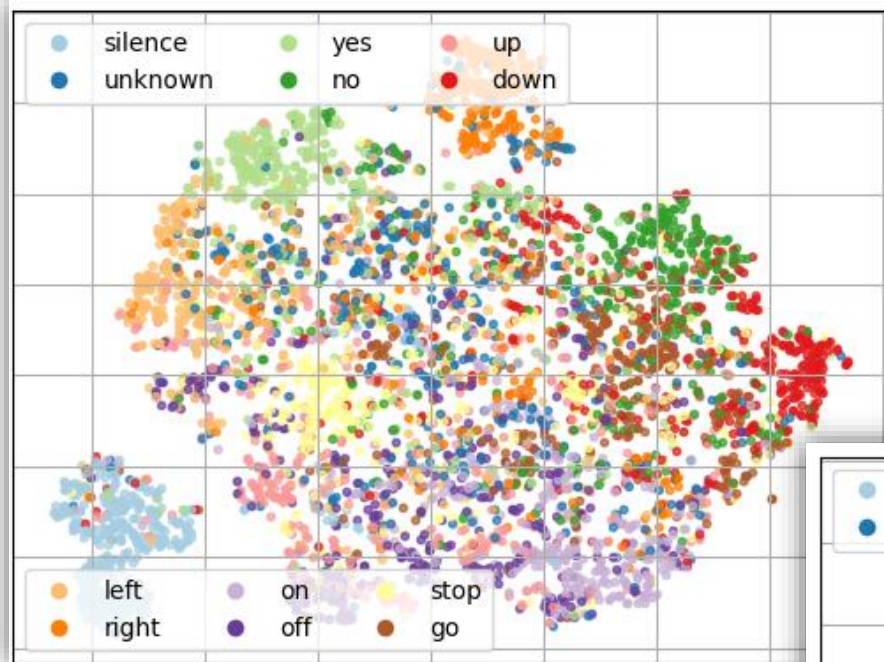


NL-KWS

ODDA

- *Unknown, up, off, on*
- **Significant improvements**

NA-KWS



Exposure to env. noise → feature separation

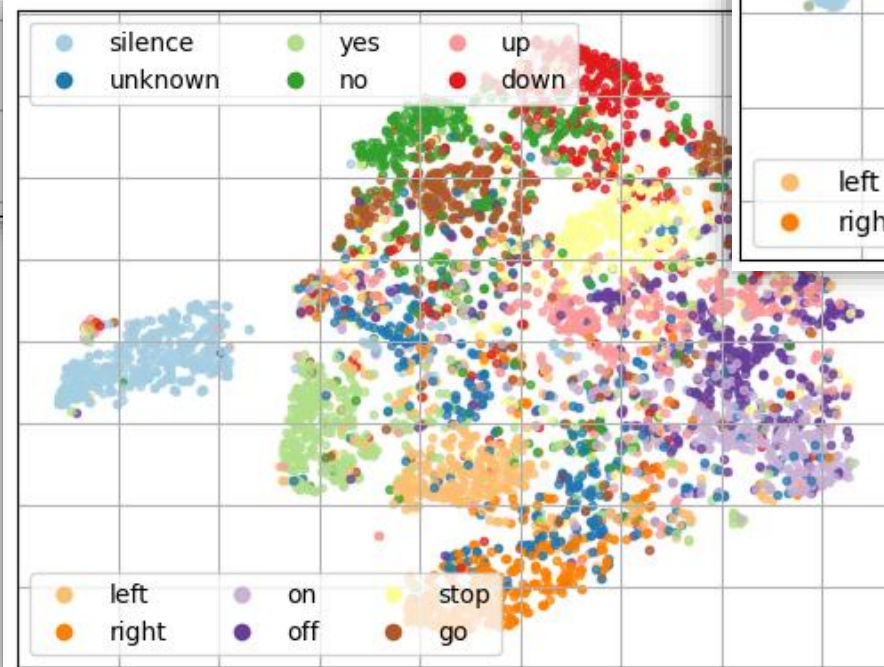
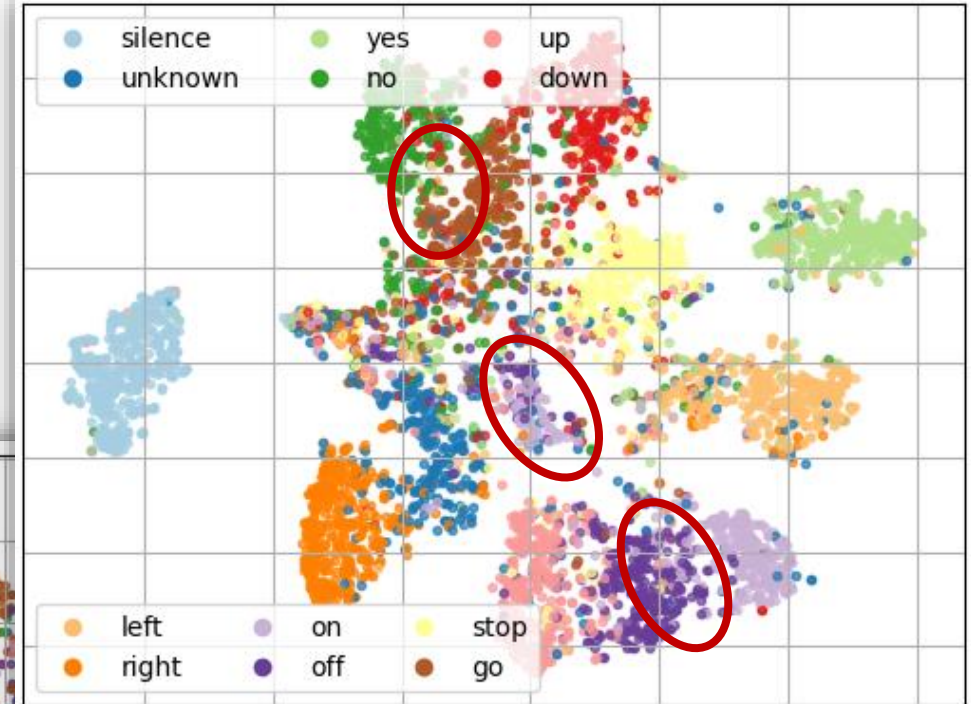
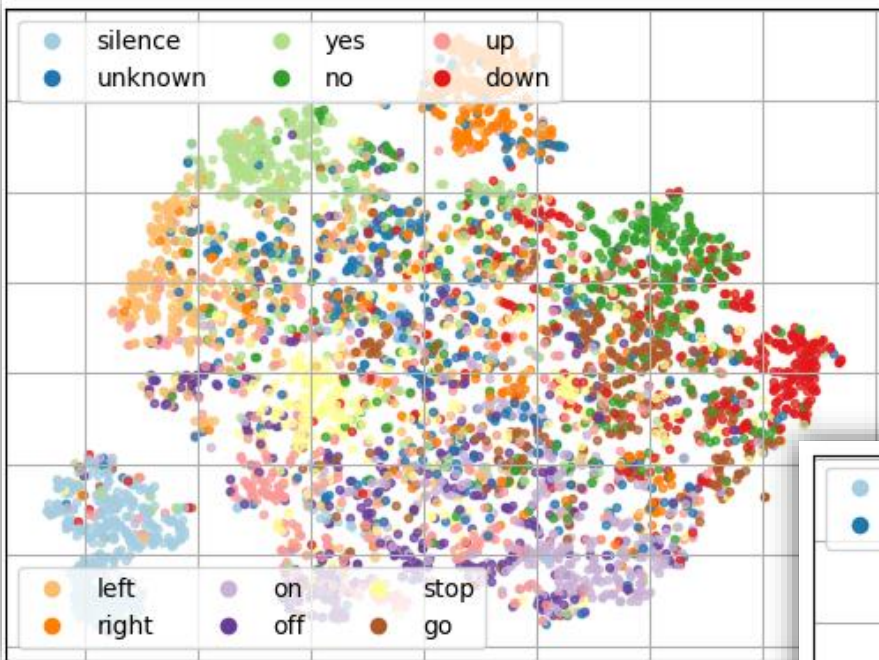


NL-KWS

ODDA

- *On, off*
- *No, go*
 - Short, phonetically similar

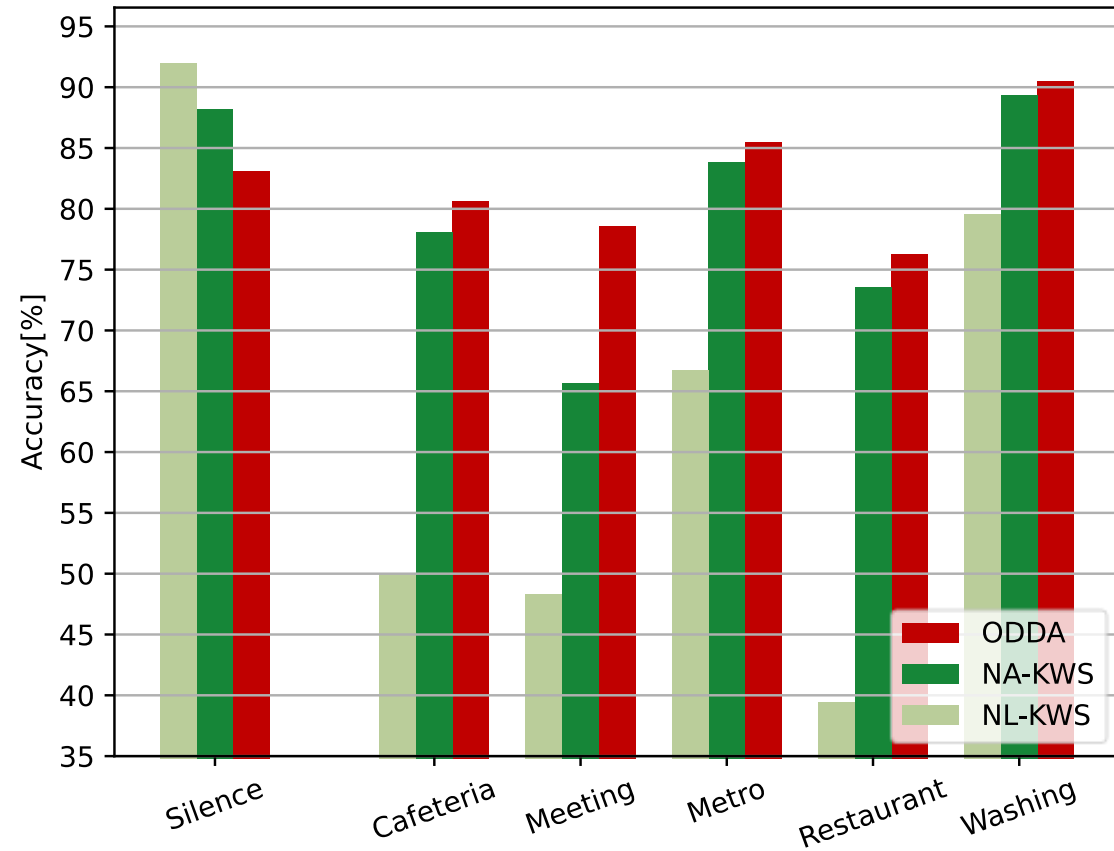
NA-KWS



ODDA offers accuracy improvements on any noise...



- **For 0 dB SNRs**
 - **12% improvement** on *meeting* over NA-KWS
 - **+30%** over NL-KWS
 - **37% improvement** on *restaurant* over NL-KWS
 - **+3%** over NA-KWS



... and any SNR

- For 0 dB SNRs
 - 12% improvement on *meeting* over NA-KWS
 - +30% over NL-KWS
 - 37% improvement on *restaurant* over NL-KWS
 - +3% over NA-KWS
- ODDA's effect scales with SNR
 - 5-6% to <1%

		SNR _{rest}			
		-10	0	30	AVG
SNR _{ODDA}	[-10,10]	59.68	80.96	85.86	79.23
	0	54.34	81.87	87.71	79.29
	[0,10]	51.52	81.06	88.93	76.49
	NA-KWS	49.10	75.85	87.61	74.49



... and any SNR



- **For 0 dB SNRs**
 - **12% improvement** on *meeting* over NA-KWS
 - **+30%** over NL-KWS
 - **37% improvement** on *restaurant* over NL-KWS
 - **+3%** over NA-KWS
- **ODDA's effect scales with SNR**
 - 5-6% to <1%
 - Up to **10%** for known env. SNR

		SNR _{rest}			
		-10	0	30	AVG
SNR _{ODDA}	[-10,10]	59.68	80.96	85.86	79.23
	0	54.34	81.87	87.71	79.29
	[0,10]	51.52	81.06	88.93	76.49
	NA-KWS	49.10	75.85	87.61	74.49





Institut für Integrierte Systeme – ETH Zürich

Gloriastrasse 35
Zürich, Switzerland

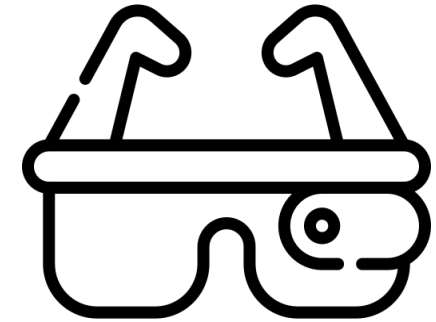
DEI – Università di Bologna

Viale del Risorgimento 2
Bologna, Italy

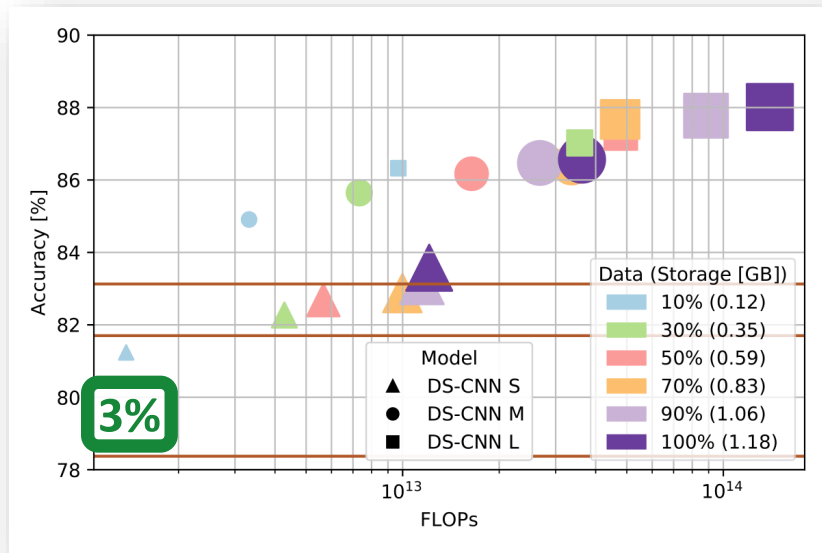
We are improving accuracy,
but at what cost?

TinyML constraints

- **(Extreme) Edge AI/ML**
- **Embedded, miniaturized devices**
 - Limited storage (e.g., model parameters)
 - Limited memory (e.g., activations, gradients)
- **Real-time operation**
 - Limited computational resources
 - Maximize throughput, minimize latency
- **Always-on, battery operated devices**
 - (Ultra-) Low-Power systems



Resource-constrained ODDA

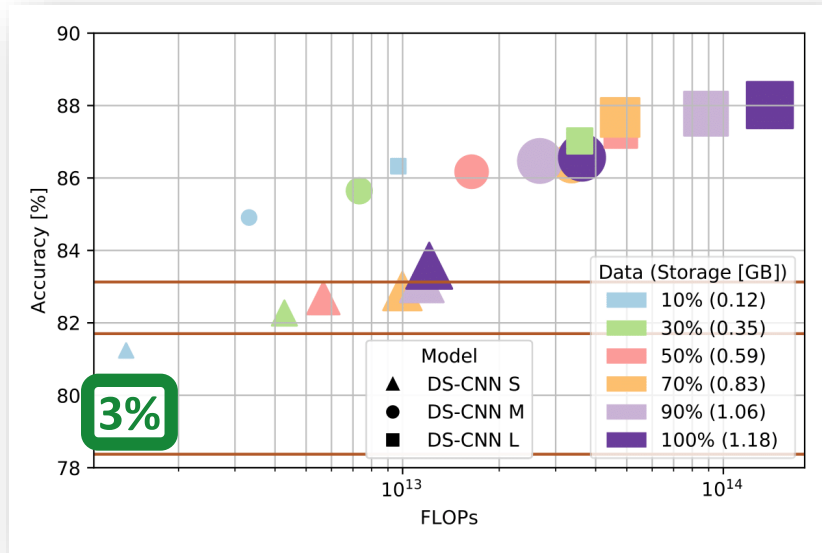


Storage

- 120 MB – 310 samples/class
 - +3% over NA-KWS
 - +5% on *meeting* with 3 MB
- Increasing model size
 - +3.7% for M model
 - +5% for L model

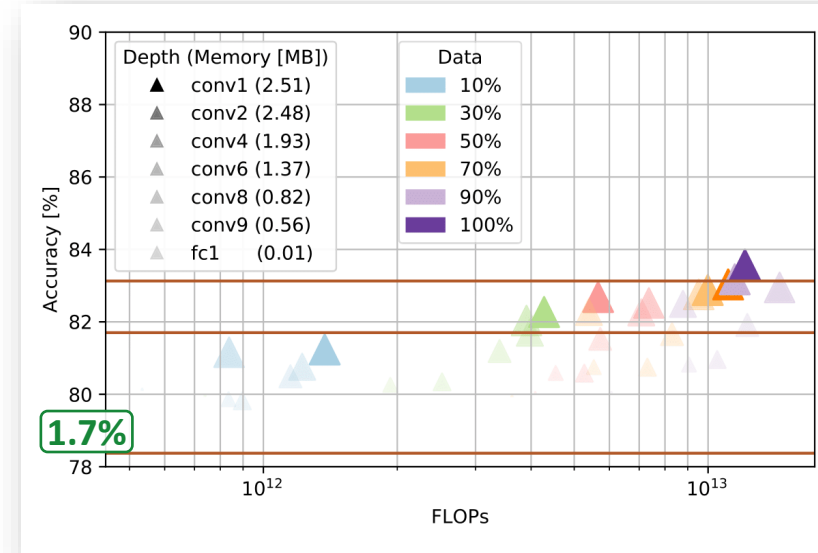


Resource-constrained ODDA



Storage

- 120 MB – 310 samples/class
 - +3% over NA-KWS
 - +5% on *meeting* with 3 MB
- Increasing model size
 - +3.7% for M model
 - +5% for L model

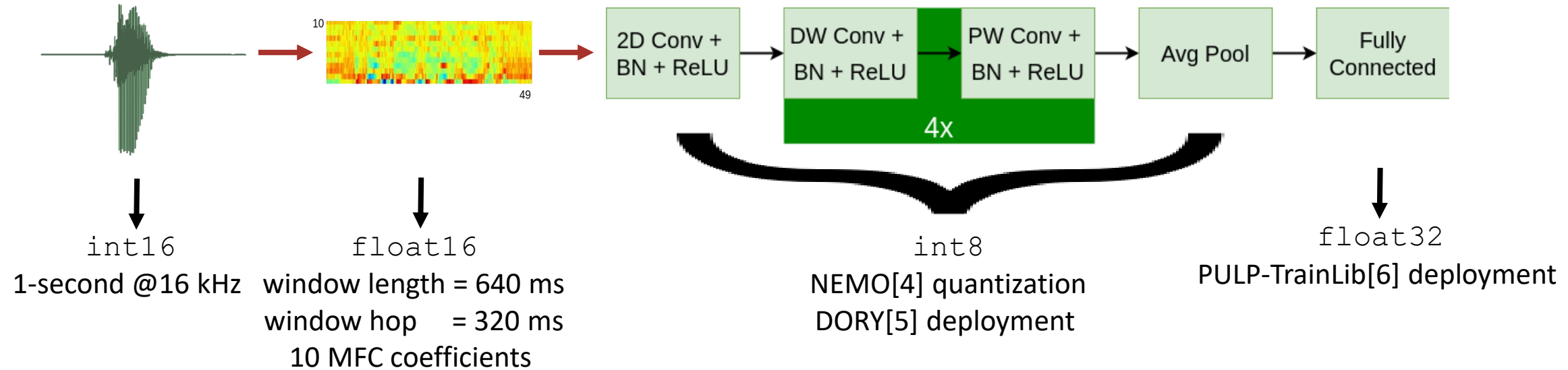


Memory

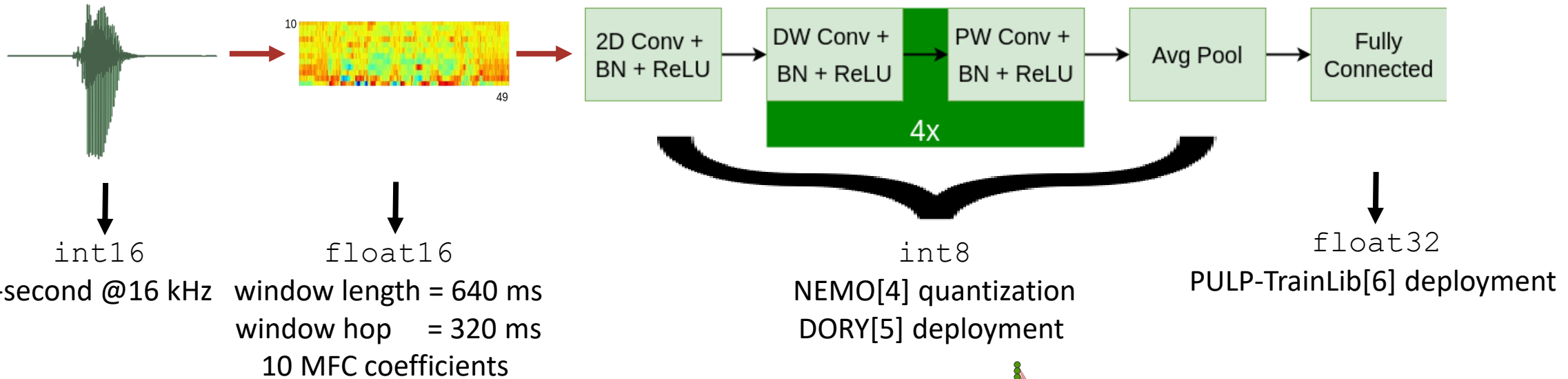
- +1.7% by refining the only classifier
 - 10 kB (i.e., parameters, activations)



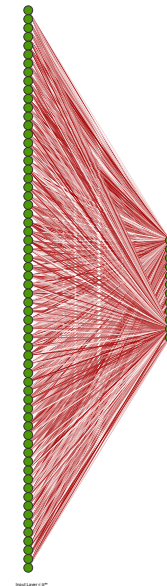
Model deployment for On-Device Domain Adaptation



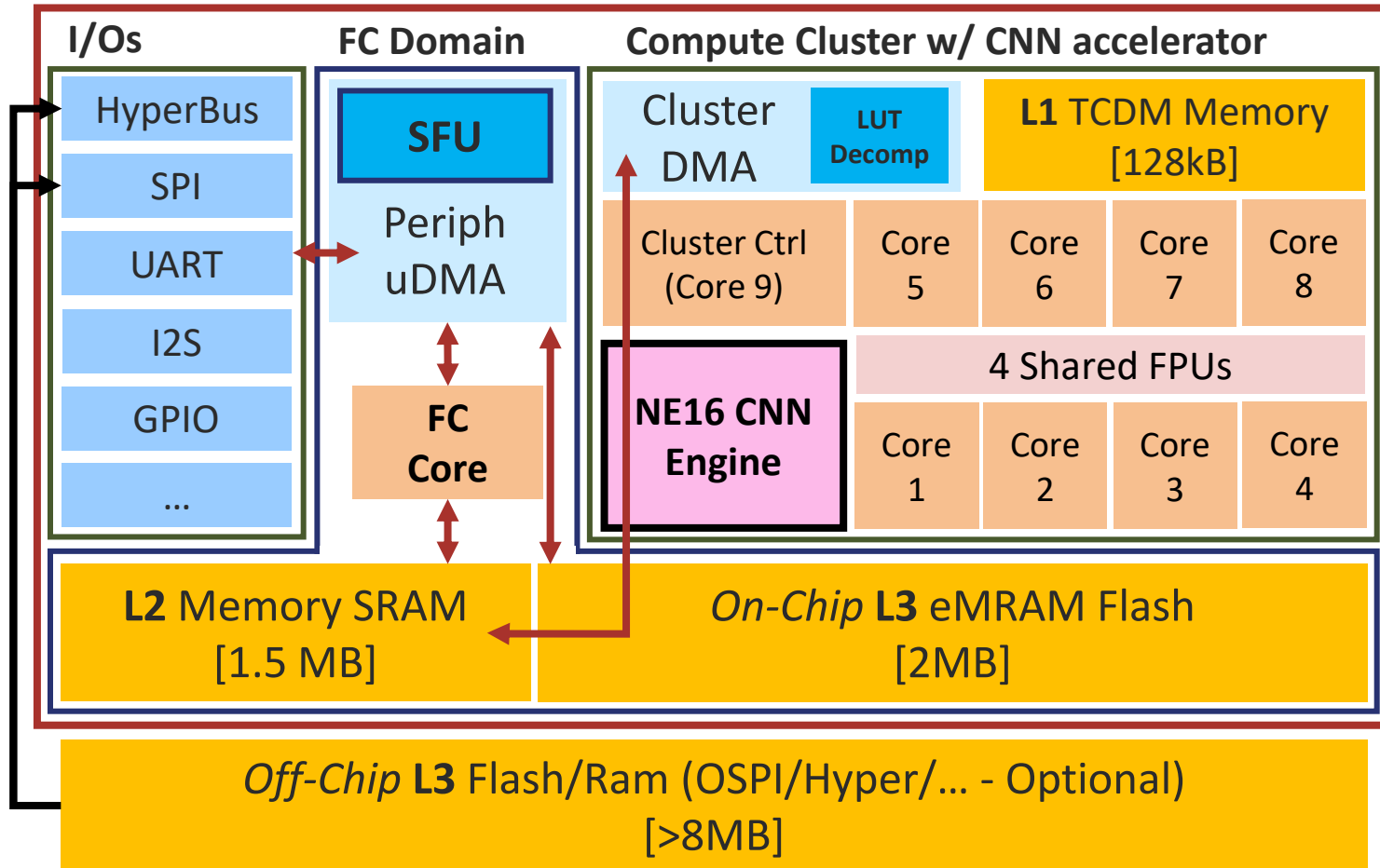
Model deployment for On-Device Domain Adaptation



- Fully Connected structure
 - 64 input neurons
 - 12 output neurons
 - **768 weights**
- PULP-Trainlib
 - Softmax
 - Cross-Entropy loss



On-Device Domain Adaptation target: GAP9

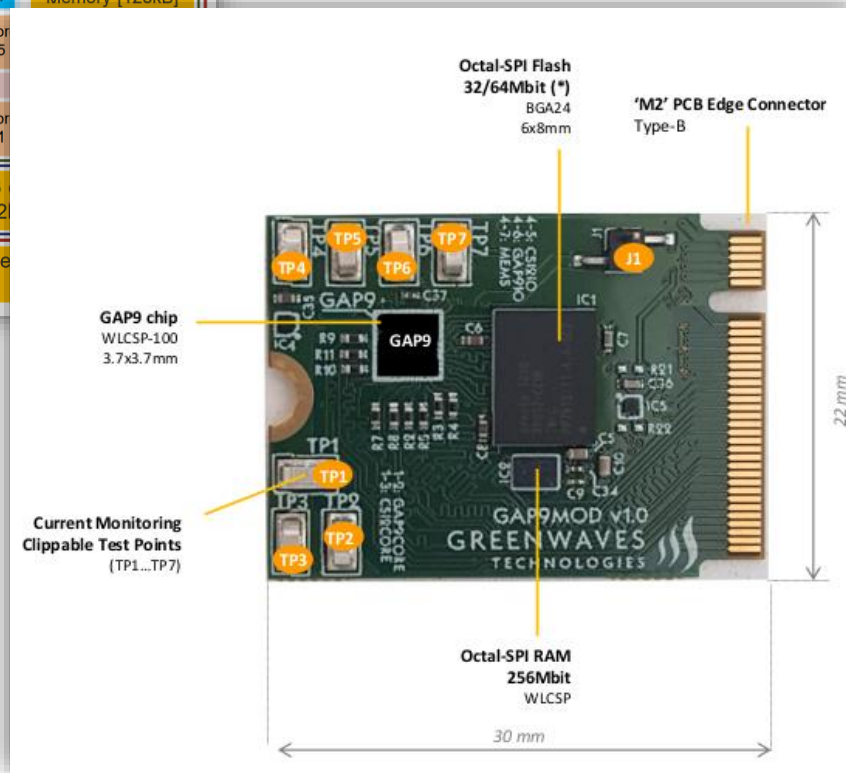
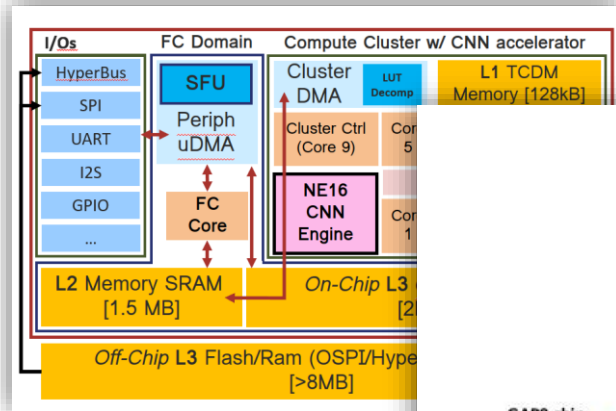


GAP9 Architecture

- 4 frequency domains
- Hierarchical memory architecture
- Heterogeneous compute units
 - 10 general purpose RISC-V cores
 - 4 shared FPUs



On-Device Domain Adaptation target: GAP9

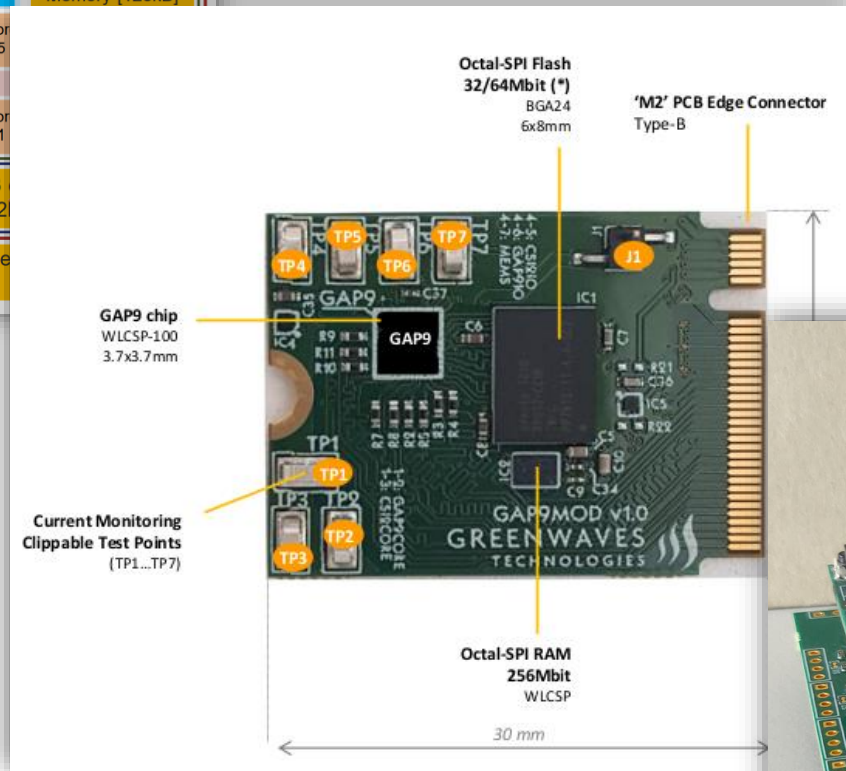
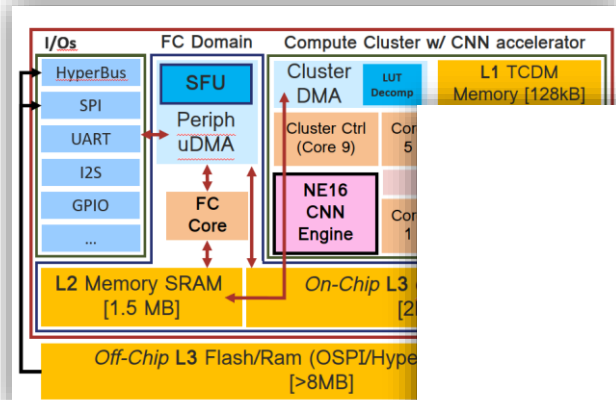


GAP9Mod

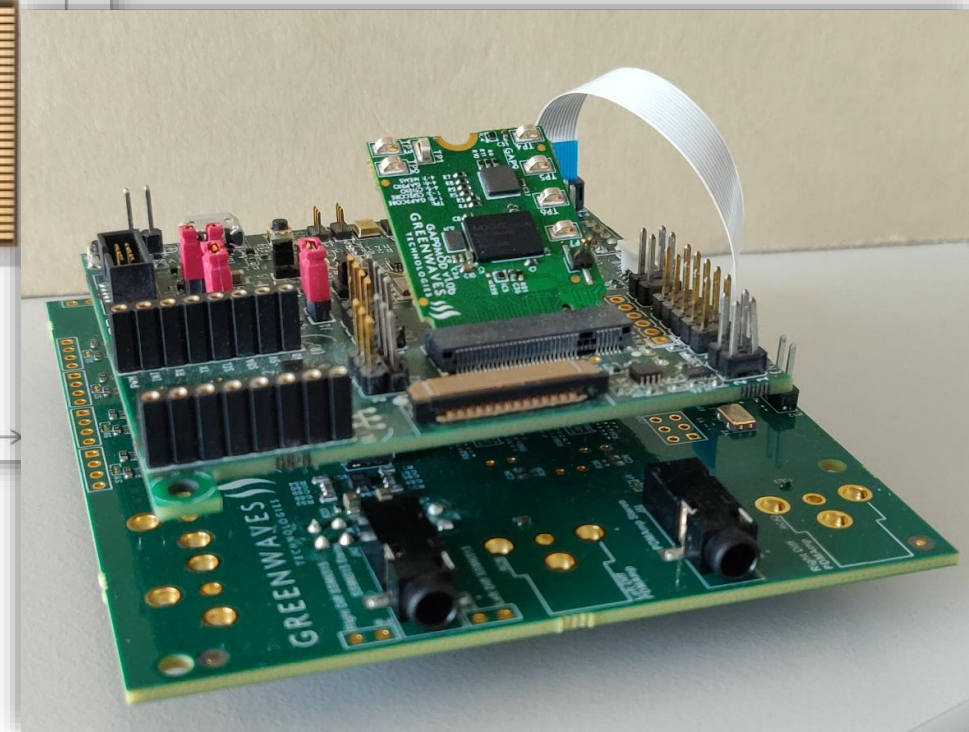
- GAP9 chip
- Octal-SPI RAM (32 Mbit)
- Octal-SPI Flash (256 Mbit)



On-Device Domain Adaptation target: GAP9



- Evaluation Kit
- Audio Add-On
 - 4 PDM microphones
 - 2 DAC+amplifier



ODDA pipeline on GAP 9 – a practical example



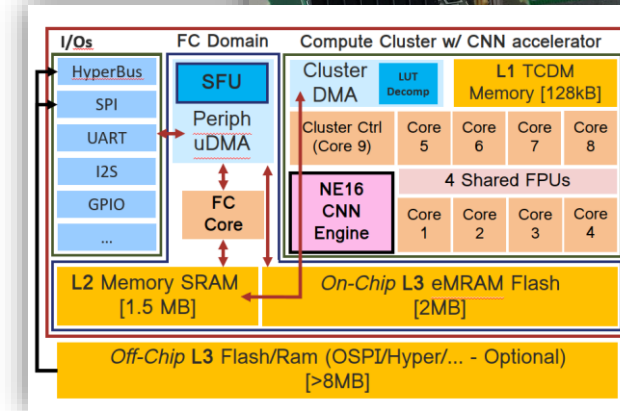
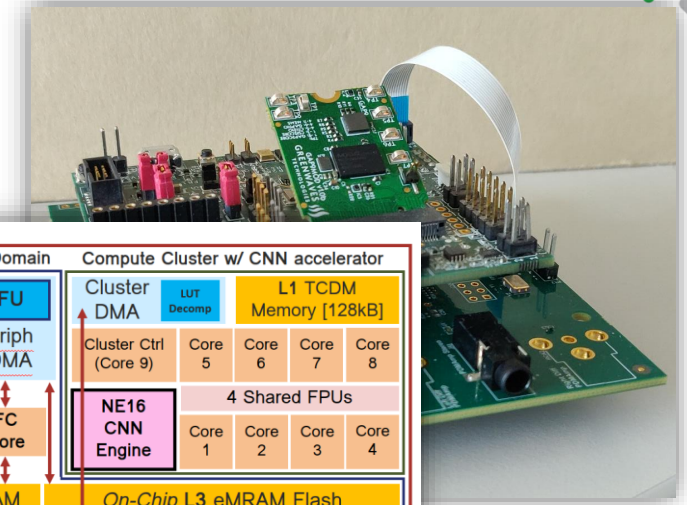
1. Flash pre-recorded .WAVs, model parameters

2. Streaming inference

1. Always-on recording using Audio Add-on mic(s)
2. MFCC computation using FPU
3. Backbone inference on 9-core cluster
4. Classification using 9-core cluster

3. On-Device Domain Adaptation

1. Record on-site noise, augment pre-recorded samples
2. MFCC computation using FPU
3. Backbone inference on 9-core cluster
4. Classification using 9-core cluster
5. Compute loss, compute gradients, update weights



ODDA pipeline on GAP 9 – a practical example



GAP9 in ultra-low-power mode ($f_{FC}=50$ MHz, $f_{CL}=50$ MHz, $V_{dd}=650$ mV)

Step	Data acq.	MFCC	Backbone	Classifier	Inference
Latency [ms]	50	20.9	26.4	1.6	50
Energy [mJ]	0.54	0.34	0.46	0.02	1.36

- Inference triggered every 50 ms; energy consumption as low as 1.36 mJ

Step	Data prep.	MFCC	Backbone	Classifier	Update	Training
Latency [ms]	2.5	20.9	26.4	1.6	2.3	53.7
Energy [mJ]	0.03	0.34	0.46	0.02	0.03	0.88

- For 1 sample/class for evaluation and 10 samples/class for training:
 - Training process completes in 7.4 s
 - ODDA (incl. evaluation and data acquisition) – energy consumption of 117.5 mJ



Conclusions



- **ODDA improves over NA-KWS by specializing on the target noise**
 - Accuracy **gains up to 12%** over NA-KWS, **37%** over NL-KWS at **0 dB** for **DS-CNN S**
 - Average gains of **10% over NA-KWS for – 10 dB**
- **Enables inexpensive on-device learning**
 - **100 samples (3 MB)** for **+5% on *meeting***
 - As little as **10 kB of memory**
- **Demonstrates practical value on GAP9**
 - **0.88 mJ** per training epoch
 - **ODDA completes in 7.4 s**, requiring as low as **118 mJ**



Cristian Cioflan

cioflanc@iis.ee.ethz.ch

github.com/pulp-platform/odda-for-kws



Institut für Integrierte Systeme – ETH Zürich

Gloriastrasse 35
Zürich, Switzerland

DEI – Università di Bologna

Viale del Risorgimento 2
Bologna, Italy



Bibliography



- [1] C. Cioflan et al., “Towards On-device Domain Adaptation for Noise-Robust Keyword Spotting”, AICAS, 2022
- [2] Pete Warden, “Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition”, 2018
- [3] Joachim Thiemann, Nobutaka Ito & Emmanuel Vincent , “DEMAND: a collection of multi-channel recordings of acoustic noise in diverse environments”, ICA, 2013
- [4] Francesco Conti, "Technical Report: NEMO DNN Quantization for Deployment Mode", 2020;
<https://github.com/pulp-platform/nemo>
- [5] A. Burrello et al., "DORY: Automatic End-to-End Deployment of Real-World DNNs on Low-Cost IoT MCUs", IEEE Transactions on Computers, 2021; <https://github.com/pulp-platform/dory>
- [6] D. Nadalini et al., "PULP-TrainLib: Enabling On-Device Training for RISC-V Multi-core MCUs Through Performance-Driven Autotuning", Embedded Computer Systems: Architectures, Modeling, and Simulation, 2022; <https://github.com/pulp-platform/pulp-trainlib>





Copyright Notice

This multimedia file is copyright © 2023 by tinyML Foundation. All rights reserved. It may not be duplicated or distributed in any form without prior written approval.

tinyML[®] is a registered trademark of the tinyML Foundation.

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