



SpokeN-100

A Cross-Lingual Benchmarking Dataset for The Classification of Spoken Numbers in Different Languages

René Groh, Nina Goes and Andreas M. Kist
Friedrich-Alexander-Universität Erlangen-Nürnberg
{rene.groh, nina.goes, andreas.kist}@fau.de

SpokeN



Problem

Speech to numerical input



Automatic recognition of a pronounced phone number

What system is capable of:

→ + 4 9 7 2 1 1 7 8 3 3 2 4 4 1

Reality:

→ + 49 72 11 78 33 24 41

**We want a local system that is able to perform this task
– is there a dataset for this?**

- **Speech Commands/AudioMNIST: numbers 0 to 9**
- **Most datasets are monolingual**
- **Combining different data sets leads to inconsistent data quality**

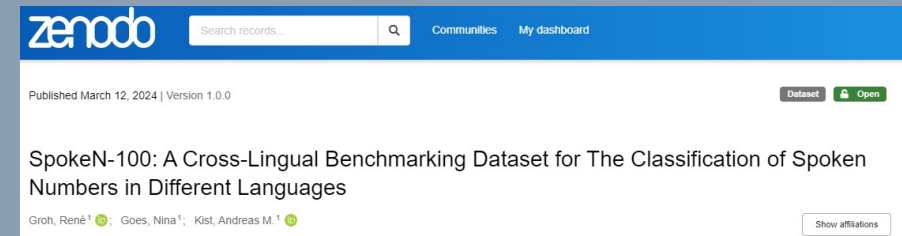
SpokeN-100

Classification of spoken numbers



Spoken numbers from 0 to 99

32 different speaker



Four different languages: English, Mandarin, German and French

In total: 12,800 audio samples

Transparent, open-source, reproducible

→ All data is **artificially generated**

SpokeN-100 Dataset

Completely artificially generated data

LLM

	English	German	French	Mandarin
0	"Zero"	"Null"	"Zéro"	"零"
1	"One"	"Eins"	"Un"	"一"
2	"Two"	"Zwei"	"Deux"	"二"
...

N=99

→ ChatGPT 3.5

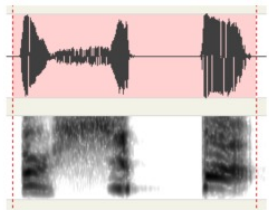


Generative Audio AI

ElevenLabs  32 speaker

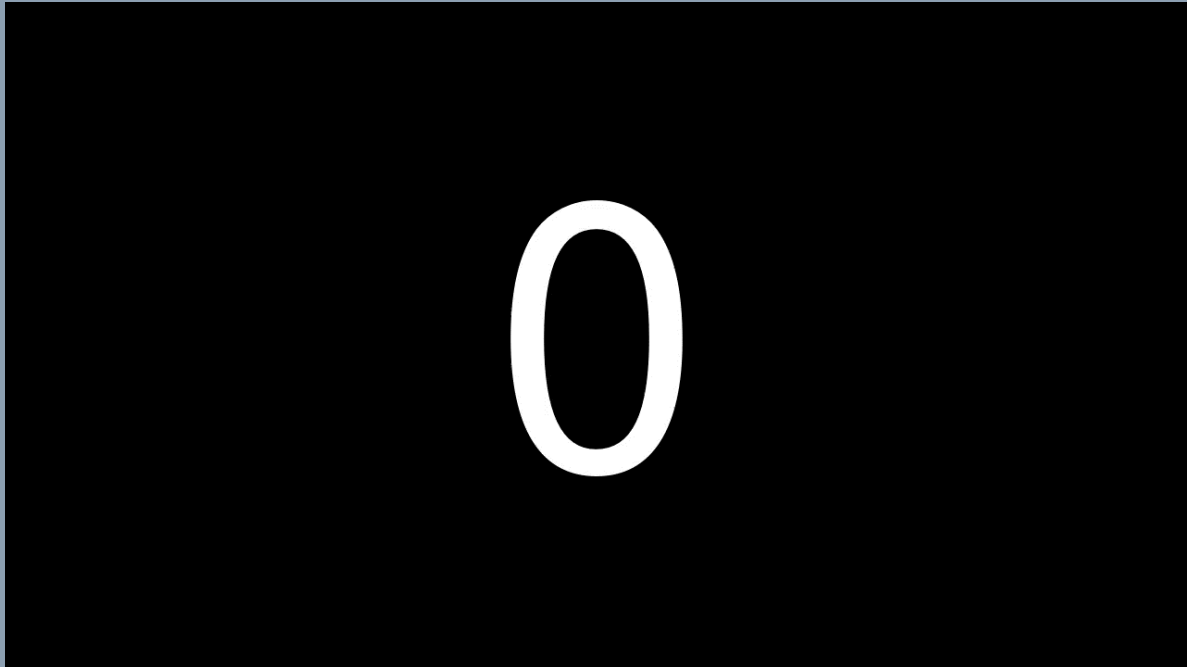


Audio Segmentation

WebMAUS 

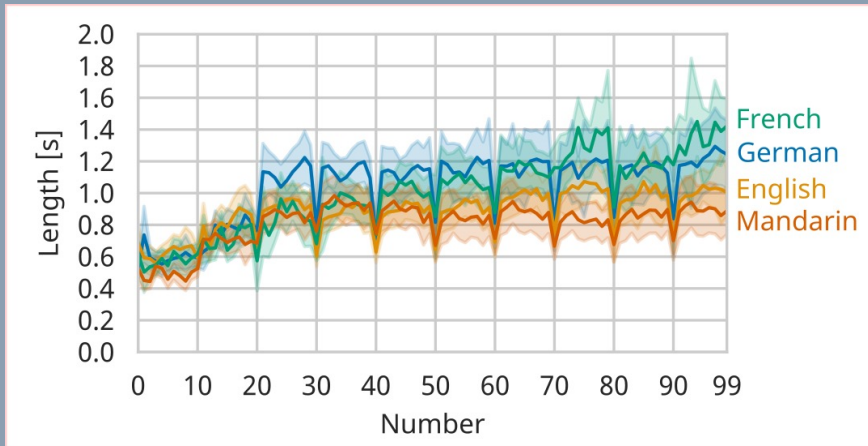
SpokeN-100 Dataset

Completely artificially generated data

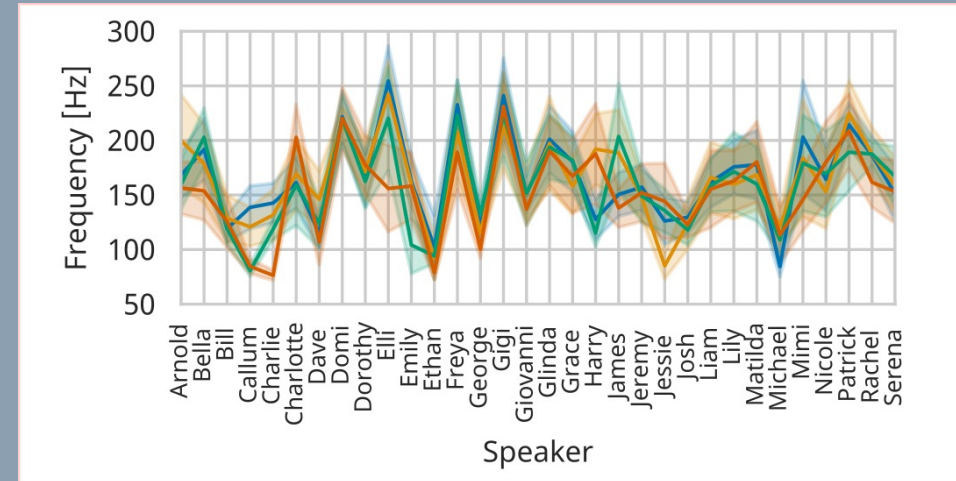
A large, white, sans-serif digit '0' is centered within a solid black rectangular frame. The digit is clean and minimalist, with a consistent stroke width.

SpokeN-100 Dataset

Descriptive statistics of the generated data



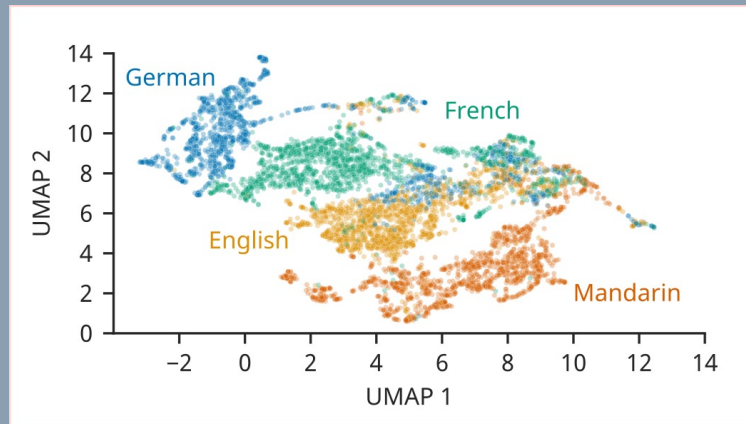
- Audio length increases as the numerical value increases
- Mandarin most efficient
- German numbers are the longest, except for higher numbers where French takes over



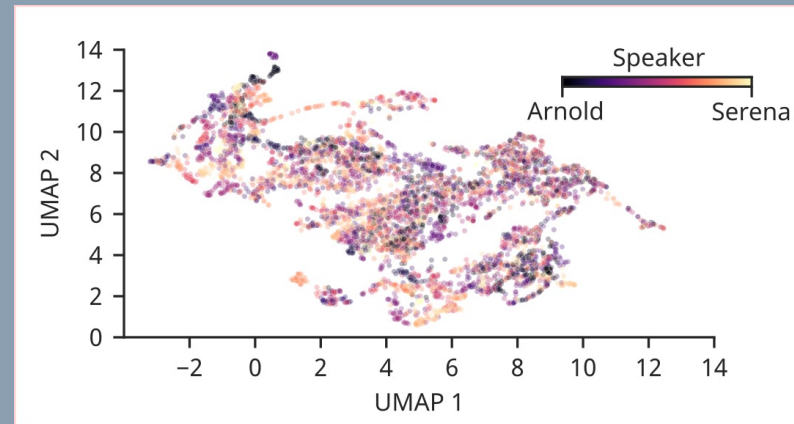
- Fundamental frequency similar for every speaker and every language
- Wide variety of fundamental frequencies

SpokeN-100

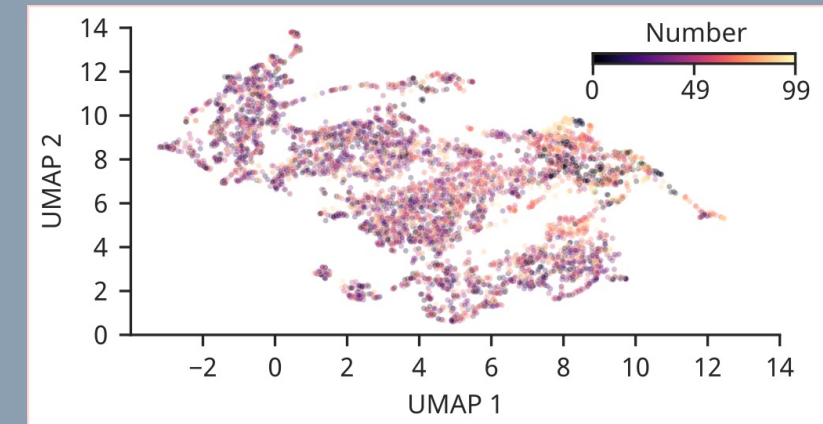
Dimensionality reduction of Mel-Spectrograms using UMAP



- Languages cluster together in the 2D visualization



- No clear clusters emerge when color-coded based on speakers



- No clear clusters emerge when color-coded based on spoken number

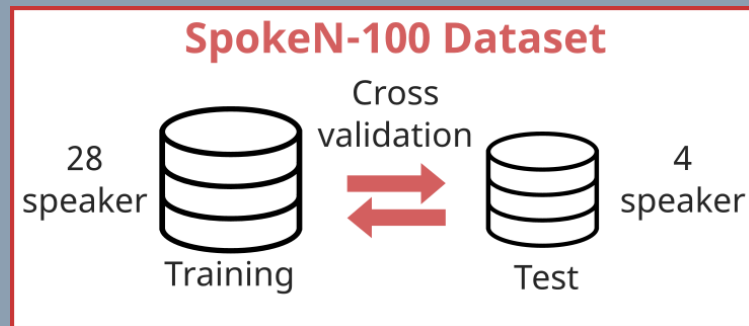
Benchmark task design

We introduce two classification tasks

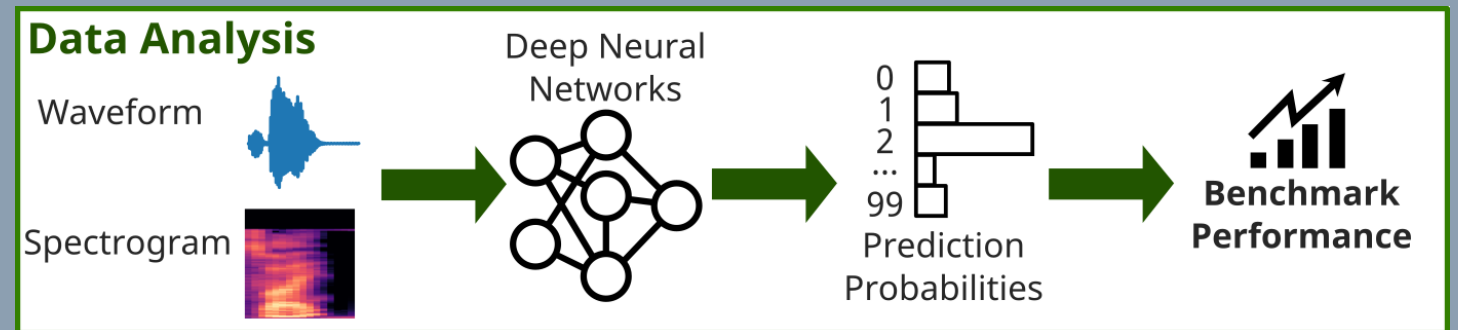
Classification of

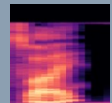
a) Languages (Four classes)

b) Spoken numbers for any language (100 classes)



8-fold cross-validation

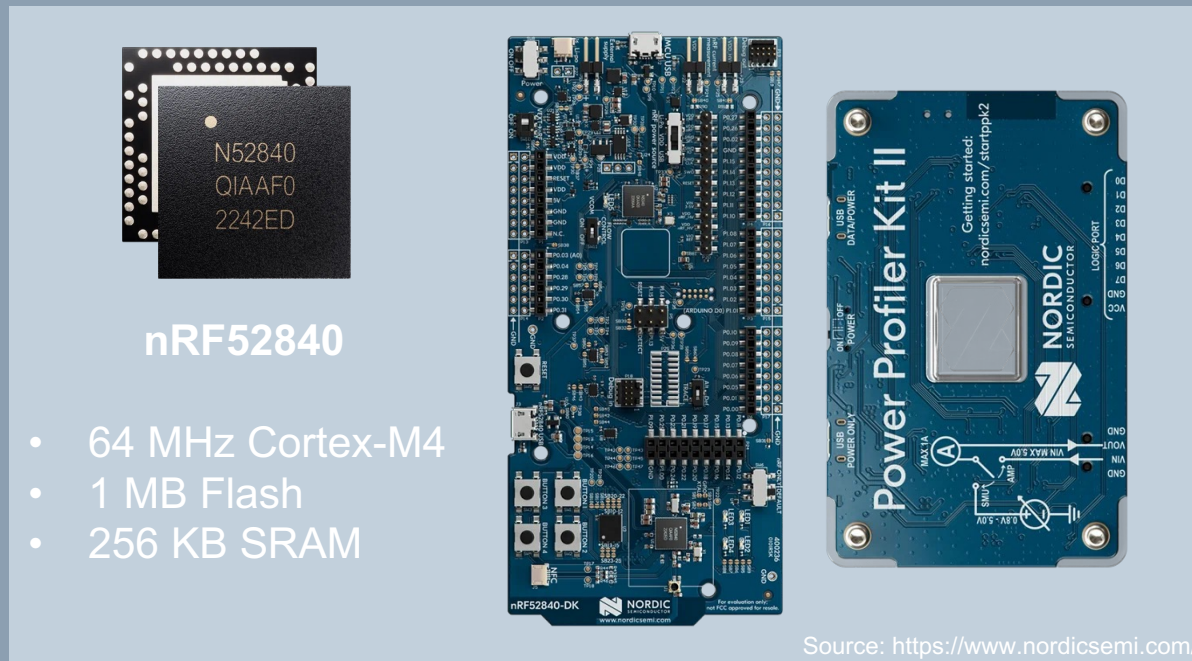




Architecture
1D CNN
2D CNN
RNN
EfficientNet-B0 CNN
Transformer

→ **Not** deployable to a microcontroller

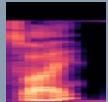
Groh and Kist – „End-to-end evolutionary neural architecture search for microcontroller units”, 2023



→ **Goal:** Maximize fitness which is a metric that combines multiple objectives (accuracy, inference time, energy consumption)

Results

Benchmark results for MCU-optimized architectures

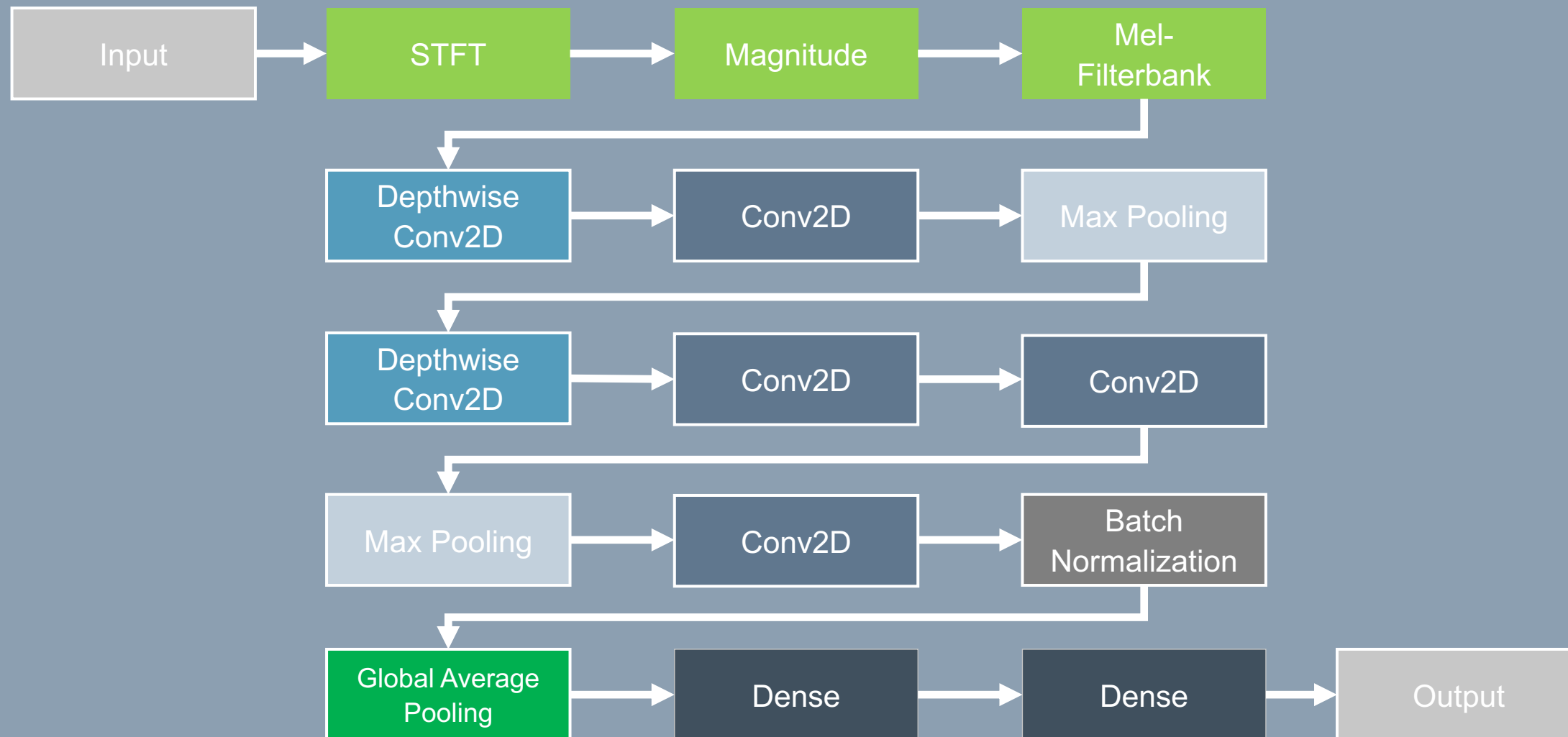


Architecture
EfficientNet-B0 CNN
Most accurate
Fastest
Fittest

→ Proof-of-concept: huge **optimization potential**

Results

Fittest model architecture



- SpokeN-100: Artificially generated speech recognition dataset with spoken numbers from 0 to 99 in four different languages
- It can serve as a benchmark for tinyDL datasets
- Data is not noisy, as with most other data sets → perfect for deep learning training, as you can control the amount of noise yourself
- High practical relevance for tinyDL applications: robot navigation, interaction with wearable devices, ...
- Possibility to expand the database as required (voice cloning)



<https://zenodo.org/records/10810044>

Acknowledgements



Martin Fernholz
Prof. Dr. Tobias Bonhoeffer



Prof. Dr. Dave Berry
Dr. Dinesh Chhetri



Prof. Dr. Rebecca Leonard



Prof. Dr. Nicole Li-Jessen

Universitätsklinikum
Erlangen



Prof. Dr. Michael Döllinger
PD Dr. Anne Schützenberger



Prof. Dr. Youri Maryn
Dr. Monique Verguts



Prof. Dr. Melda Kunduk



Prof. Dr. Matthias Echernach



Prof. Dr. Cara Stepp



Prof. Dr. Aaron Johnson



Prof. Dr. Cate Madill

Bayerisches Staatsministerium für
Wissenschaft und Kunst



Fonds de recherche
Santé



Bundesministerium
für Bildung
und Forschung



@rene__gr
@anki_xyz



rgroh1996
ankilab

Copyright Notice

This presentation in this publication was presented at the tinyML[®] Research Symposium 2024. 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