

Benchmarking Small Language Models on an Industry-grade, High-end Microcontroller

 **EDGEAI** FOUNDATION Research Symposium 2026

Pietro Firpo¹, Riccardo Berta¹, Francesco Bellotti¹, Luca Lazzaroni¹, Danilo Pietro Pau²

¹ University of Genoa

² STMicroelectronics

Introduction

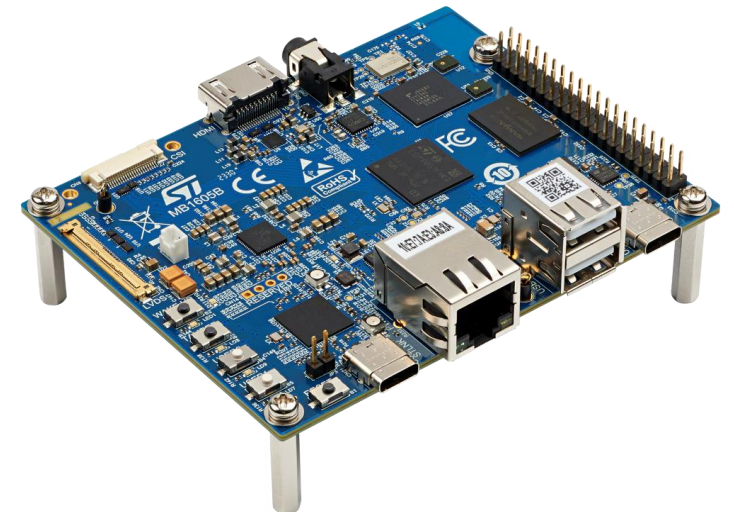
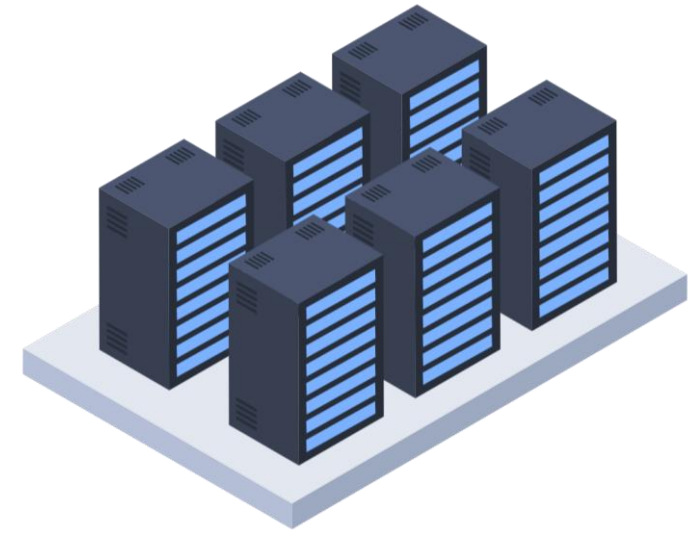
From the Cloud to the Edge

From LLMs to SLMs

- Size and HW requirements
- Benchmarking: accuracy-oriented vs latency oriented

Chosen a platform:

- RQ1: What inference throughput can SLMs achieve on an industry-grade MCU-MPU hybrid?
- RQ2: How does throughput scale with model size, and how predictive is parameter count of deployability?
- RQ3: Which throughput regimes correspond to viable embedded language-based applications?



State of the Art

SLMs:

- DistilBERT, ALBERT: smaller size, similar accuracy
- SLMquant, TensorSLM: compression techniques for LMs

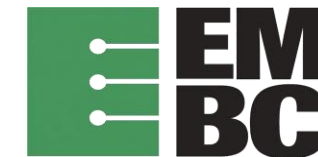
LM Benchmarking:

- GLUE: accuracy
- MLPerf: latency and throughput. High-end HW
- SLM-bench: throughput on Nvidia Jetson AGX Orin

MCU benchmarking:

- MLPerf Tiny: TinyML on MCUs
- EEMBC: embedded AI for signal processing or vision

Filling the gap: no data on LM latency/throughput on MCUs/MPUs



Experimental Settings

Metric: Inference Throughput

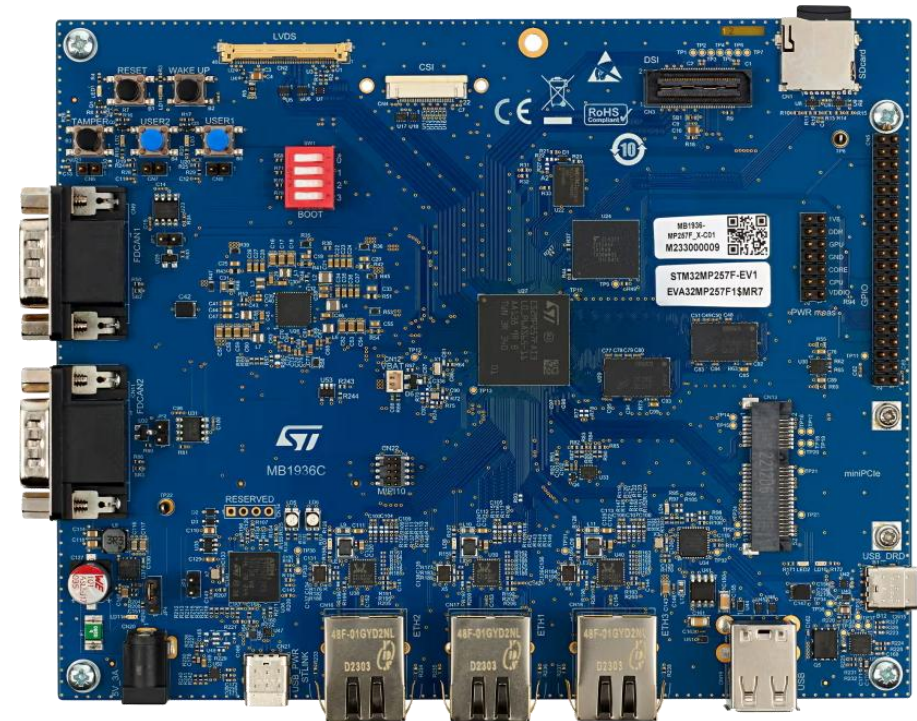
Platform: STM32MP257F-EV1

- Compromise between ultra low power and GPU powered

Component	Description
Application cores	2 × Arm Cortex-A35
Real-time core	Arm Cortex-M33
AI accelerator	VeriSilicon GC8000UL NPU
RAM	4 GB
Clock frequency	1.5 GHz

Inference frameworks:

- Llama2.c: weights and activations in FP32
- Llama.cpp: Q8 weights, Q16 activations
- Ollama: Q4 weights, Q8 activations



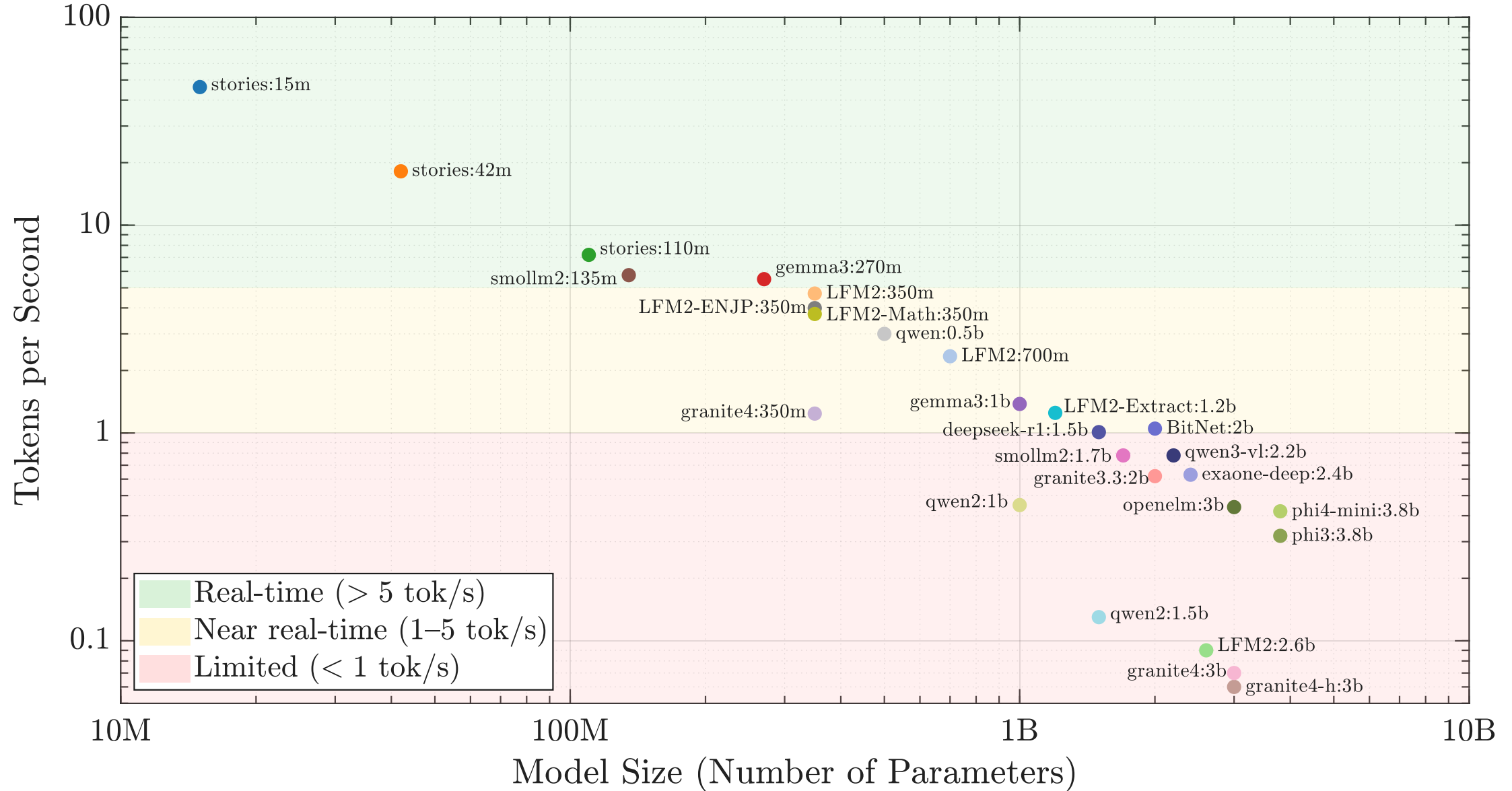
Experimental Settings

Models tested

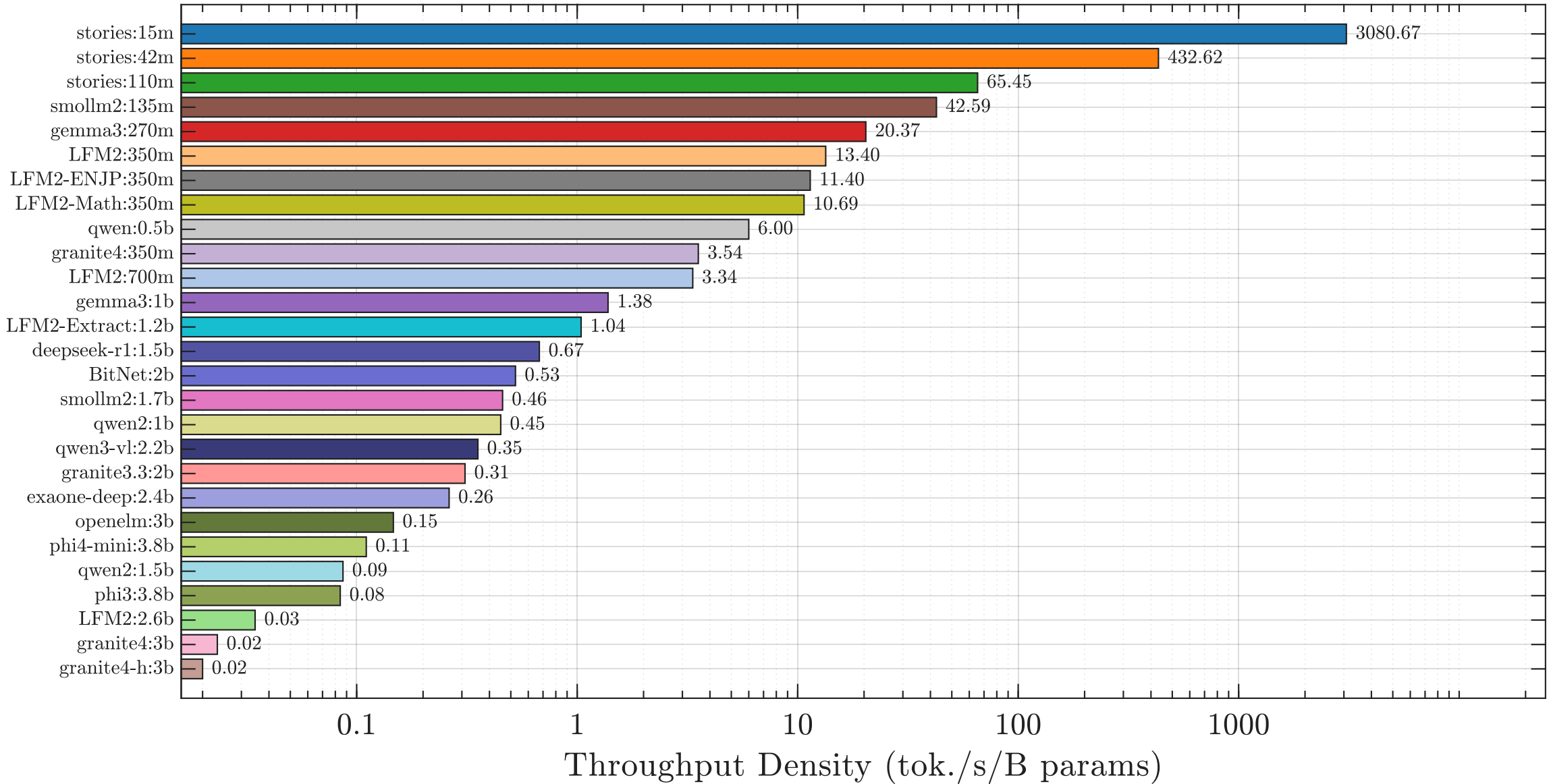
Model Name	General purpose	Params	Encoder layers	Inference framework
stories:15m	No (story generation)	15M	6	llama2.c
stories:42m	No (story generation)	42M	8	llama2.c
stories:110m	No (story generation)	110M	12	llama2.c
smollm2:135m	Yes	135M	30	ollama
smollm2:1.7b	Yes	1.7B	24	ollama
LFM2-ENJP:350m	No (EN-JP translation)	350M	16	llama.cpp
LFM2-Math:350m	No (math reasoning)	350M	16	llama.cpp
LFM2-Extract:1.2b	No (information extraction)	1.2B	16	llama.cpp
LFM2:350m	Yes	350M	16	ollama
LFM2:700m	Yes	700M	16	ollama
LFM2:2.6b	Yes	2.6B	30	ollama
deepseek-r1:1.5b	Yes	1.5B	28	ollama

Model Name	General purpose	Params	Encoder layers	Inference framework
gemma3:270m	Yes	270M	18	ollama
gemma3:1b	Yes	1B	26	ollama
granite3.3:2b	Yes	2B	40	ollama
granite4:350m	Yes	350M	28	ollama
granite4:3b	Yes	3B	40	ollama
granite4-h:3b	Yes	3B	40	ollama
qwen:0.5b	Yes	500M	24	ollama
qwen2:1b	Yes	1B	16	llama.cpp
qwen2:1.5b	Yes	1.5B	28	llama.cpp
qwen3-vl:2.2b	Yes	2.2B	28	ollama
BitNet:2b	Yes	2B	30	proprietary
exaone-deep:2.4b	Yes	2.4B	30	ollama
openelm:3b	Yes	3B	36	ollama
phi3:3.8b	Yes	3.8B	32	ollama
phi4-mini:3.8b	Yes	3.8B	32	ollama

Results



Results



Conclusions

RQ1: inference throughput on hybrid MCU-MPU

Broad range, 0.1 tok./s < throughput < 10 tok./s

RQ2: throughput scaling wrt model size and relationship between param. count and deployability

Throughput influenced by more than size. Difficult scaling

RQ3: throughput regimes for applications

> 5 tok./s: RT and latency-sensitive applications

1-5 tok./s: near-RT and agentic

< 1 tok./s: severely limited, highly restrictive use cases



Thanks for your attention