

Distributed Inference with Minimal Off-Chip Traffic for Transformers on Low-Power MCUs

Severin Bochem^{*}, Victor J.B. Jung[†], Arpan Suravi Prasad[†], Francesco Conti[‡], Luca Benini[†][‡]

*D-ITET, ETH Zurich; †Integrated Systems Laboratory, ETH Zurich; ‡DEI, University of Bologna

PULP Platform Open Source Hardware, the way it should be!



Transformer Model Architecture





ETH zürich

Idea: Distribute Workload across mutliple chips

- One chip is too small to hold full model weights
- Distribute model across multi-chip system to run from on-chip memory





Partitioning Scheme



Previous work **duplicated weights** or uses **model pipelining**, both of which is infeasible for autoregressive Transformer inference!



ETH zürich

Speedup of Distributed Inference



Using more than two chips achieves **above linear speedup** by keeping **layer weights in on-chip** memory



ETH zürich

Scalability Study for TinyLlama



Partitioning becomes **less efficient** as tensor sizes per chip scale down





Conclusion



- Proposed and benchmarked a Transformer partitioning scheme
- No weight duplication and small chip-to-chip traffic
- Achieves super-linear speedup for various models

Want to hear more? See more results? Stop by at the poster session!



Paper on Arxiv

