

# RedMULE: A Compact FP16 Matrix-Multiplication Accelerator for Adaptive Deep Learning on RISC-V-Based Ultra-Low-Power SoCs

Yvan Tortorella, Luca Bertaccini, Davide Rossi, Luca Benini, Francesco Conti



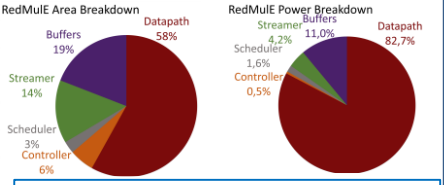
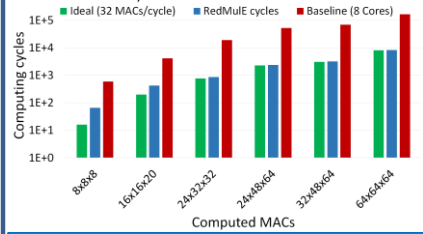
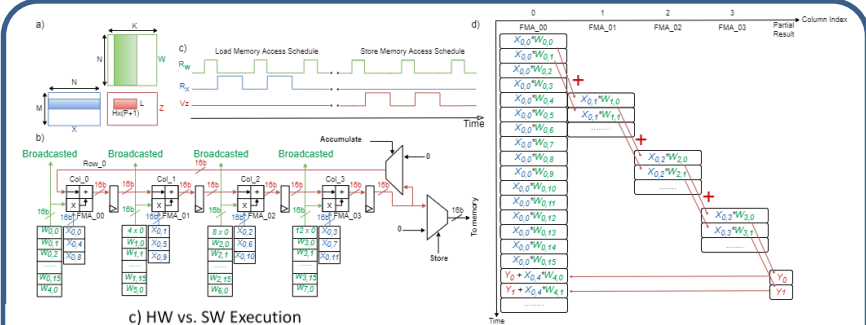
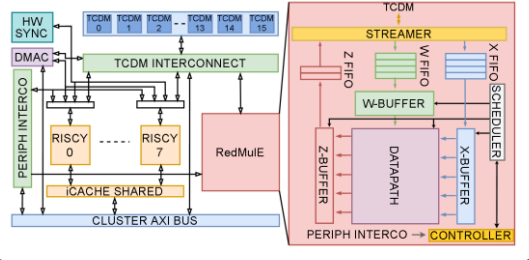
The request for **Energy-efficient** hardware enabling **extreme-edge** Deep Learning (DL) is increasing

**Narrow integer** operations suffice for extreme-edge **inference**, helping reducing average power consumption and increase energy efficiency

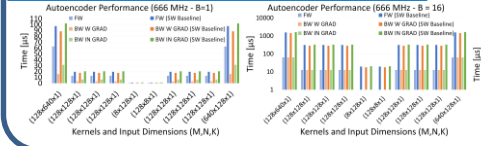
**Training on-the-edge** is **challenging** and hard to achieve in **sub-100 mW** domain due to **FP operations** requirement for accuracy and precision

**RedMULE accelerates FP16 matrix multiplications to fasten online training of generalized DL models**

**First PULP-Based reduced precision matrix multiplication accelerator enabling extreme-edge training at very low HW cost**



Up to **22x** HW speedup and **31.6 MAC/cycle** (98.8% utilization), **42 GFLOPS** at **666 MHZ**



Occupies **0.07 mm<sup>2</sup>** (14% of the PULP Cluster)  
Up to **688 GFLOPS/W**, **43.5 mW** (average Cluster power consumption) at **475 MHz**, with **30 GFLOP/sec** (from post-layout power simulations)

**2.6x** Speedup over SW (B = 1)  
**24.4x** speedup over SW (B = 16)

Category	Design	Tech	Area mm <sup>2</sup>	Freq MHz	Vol V	Power mW	Perf GOPS	Energy Eff GOPS/W	Mac Units	Precision
GPU	NVIDIA A100	7	4410	300000	-	300000	-	-	256	FP16
	EIE	45	40.8	800	1.0	278	46	166	168	INT16
Inference Chips	Zeng et al.	65	2.14	250	-	478	1152	2410	256	INT8
	Simba	16	6	168	0.82	-	9100	-	1024	INT8
Training Chips	IBM	7	19.6	1000	0.55	4400	8000	1800	4096	FP16
	Combrion-Q	45	888	1000	0.75	13000	12800	980	1024	INT8
HPC	Manticore	22	888	500	0.6	200	25	188	24	FP64
	Mat-Mul Acc.	Anders et al.	14	0.024	1000	0.9	900	54	50	FP16
Our Work	PULP (w/ RedMULE)	22	0.5	475	0.65	43.5	30	688	32	FP16
	DarkSide	65	3.85	200	1.1	88.1	12.6	152	32	FP16

We presented **RedMULE**, an example of accelerator to enable **on-chip leaning** efficiently at **low HW cost**

**Future work:**

- explore HW solution for efficient **inference and training** on lower precision (e.g. **Hybrid-FP8**)
- Introduction of **reconfiguration** features (e. g. **zero skipping**)