

Research context and motivation

- At the basis of Deep Learning (DL) algorithms are convolutions and matrix multiplications, which require the computation of many dot products and simple scalar multiplications between features and weights.
- These operations are typically executed by multiply-and-accumulate (MAC) units.
- Edge devices are often resource constrained devices that need low execution latency.

Addressed research questions/problems

- When running DL applications on edge devices, energy and latency of these MAC units have to be minimized.
- Mixed-precision quantization is a smart technique to squeeze DNNs without losing accuracy, but it requires precision-scalable hardware support.
- Data precision may vary across different applications, but also within the same application, therefore hardware need to be runtime reconfigurable in precision.
- To speedup the execution of DNN models, faster MAC units and hardware accelerators are needed.

Novel contributions

- We propose a new reconfigurable modified Radix-4 Booth signed multiplier with 4:2 compressor tree and final adder with Sum-Together (ST) mode (Fig. 1-(a)).
 - Other than 16-bit full precision multiplication, it can be reconfigured to support dot-product computations at reduced precision (8 and 4 bits).
 - It exploits normal alignment of partial products in a standard multiplier (Fig. 3), enabling the computation of dot products when two or four scalar inputs are packed in each operand without any additional costs.
 - When used in low-precision configurations, the multiplier reduces the cycles of MACs.
- We show how the flexibility of ST multipliers can be exploited in layer-specific DNN accelerators (Fig. 1-(b)): 2D-Convolution, Depth-wise Convolution and Fully- Connected.
 - We show how these accelerators are obtained with High-Level Synthesis (HLS) (Fig. 3).

Results & Future work

- On a 28-nm technology, at the cost of limited overhead in area and power compared to a non-reconfigurable design, our multiplier/dot-product unit is superior to other reconfigurable units proposed in the literature (Tab. 1).
- We made a design-space exploration (DSE) via Catapult HLS varying accelerators memory sizes, MAC units parallelism, clock frequency (and even ST multipliers type).
- The results of the DSE show many reconfigurable Pareto points, especially for low-precision configurations, which dominate the non-reconfigurable counterparts (Fig. 4).
- Our findings allow the designers to select the best ST-based accelerator depending on the target, either high performance, low area (or low power).
- In future, we plan to integrate these ST-based accelerators into a System-on-Chip with a RISC-V processor to compute heterogeneously quantized DNNs.

[2] M. Gautschi et al., "Near-threshold RISC-V core with DSP extensions for scalable IoT endpoint devices," IEEE TVLSIS, vol. 25, no. 10, pp. 2700–2713, 2017.
[3] Zhang, Z. Li, and Q. Zheng, "Design of a configurable fixed-point multiplier for digital signal processor," in Proc. PrimeAsia, pp. 217–220, 2009.
[4] R. Lin, "Reconfigurable parallel inner product processor architectures," IEEE TVLSIS, vol. 9, no. 2, pp. 261–272, 2001.

Submitted and published works

- Published: Urbinati, L., Ricci, M., Turvani, G., Tobon Vasquez, J. A., Vipiana, F., and Casu, M. R., "A Machine-Learning Based Microwave Sensing Approach to Food Contaminant Detection", IEEE International Symposium on Circuits and Systems (ISCAS), Seville, Spain, 2020, pp. 1–5.
- Published: Ricci, M., Stitic, B., Urbinati, L., Di Guglielmo, G., Tobon Vasquez, J. A., Carloni, L. P., Vipiana, F., Casu, M. R., "Machine-Learning-Based Microwave Sensing: A Case Study for the Food Industry", IEEE Journal of Emerging and Selected Topics on Circuits and Systems, vol. 11, no. 3, 2021, pp. 503–514.
- Published: Urbinati, L., and Casu, M. R., "A Reconfigurable Depth-Wise Convolution Module for Heterogeneously Quantized DNNs", IEEE International Symposium on Circuits and Systems (ISCAS), Austin, Texas, 2022, pp. 1–5.
- Accepted: Urbinati, L., and Casu, M. R., "A Reconfigurable Depth-Wise Convolution Module for Heterogeneously Quantized DNNs", IEEE International Symposium on Circuits and Systems (ISCAS), Austin, Texas, 2022.
- Accepted: Urbinati, L., and Casu, M. R., "A Reconfigurable 2D-Convolution Accelerator for DNNs Quantized with Mixed-Precision", Springer on Lecture Notes in Electrical Engineering (LNEE), ApplePies, Genova, Italy, 2022.
- Accepted: Urbinati, L., and Casu, M. R., "A Reconfigurable Multiplier/Dot-Product Unit for Precision-Scalable Deep Learning Applications", Springer on Lecture Notes in Electrical Engineering (LNEE), Società Italiana di Elettronica, Pizzo, Italy, 2022.
- Submitted: Urbinati, L., and Casu, M. R., "High-Level Design of Precision-Scalable DNN Accelerators Based on Sum-Together Multipliers", IEEE Design, Automation and Test in Europe Conference (DATE), Antwerp, Belgium, 2023.

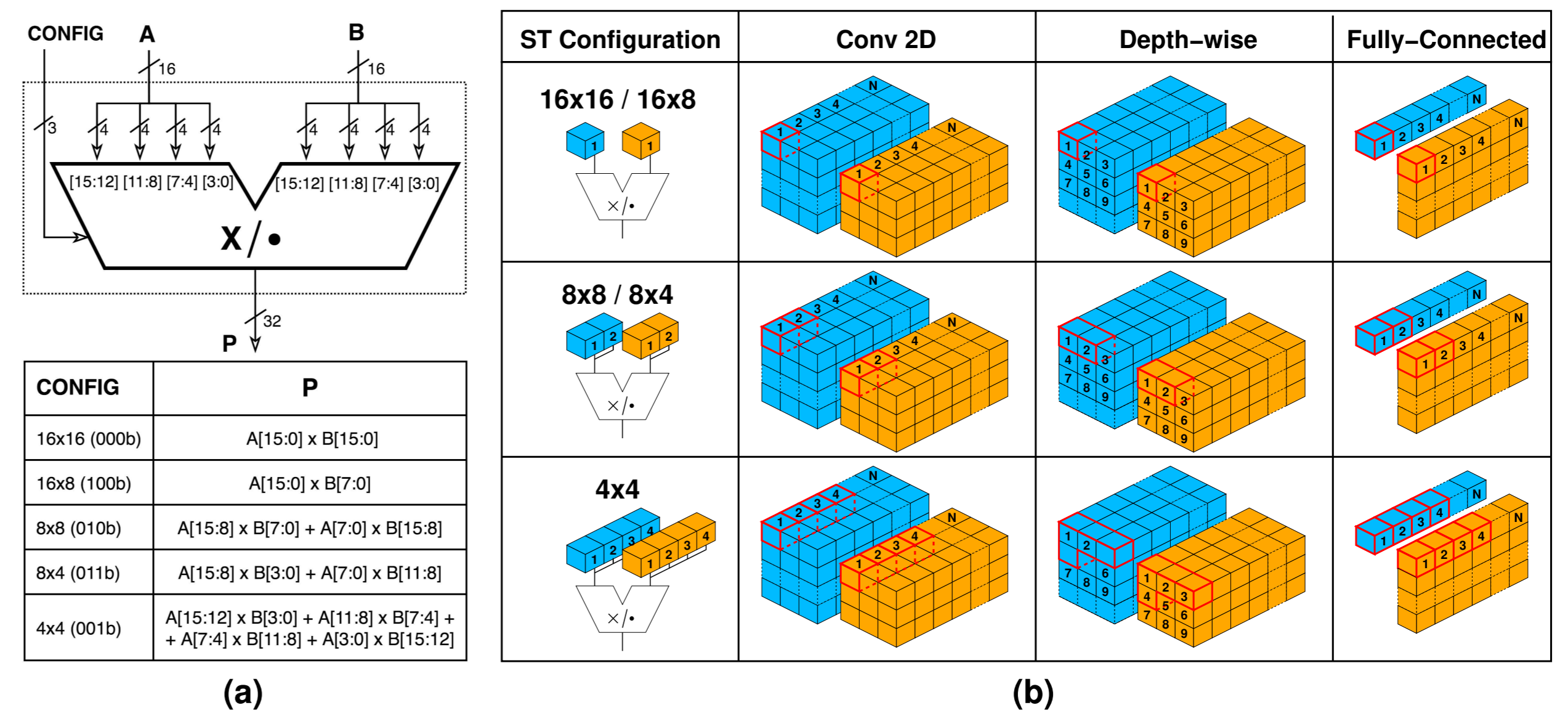


Fig. 1: ST multiplier and its five configurations (a), and working principles of our ST-based DNN accelerators (b).

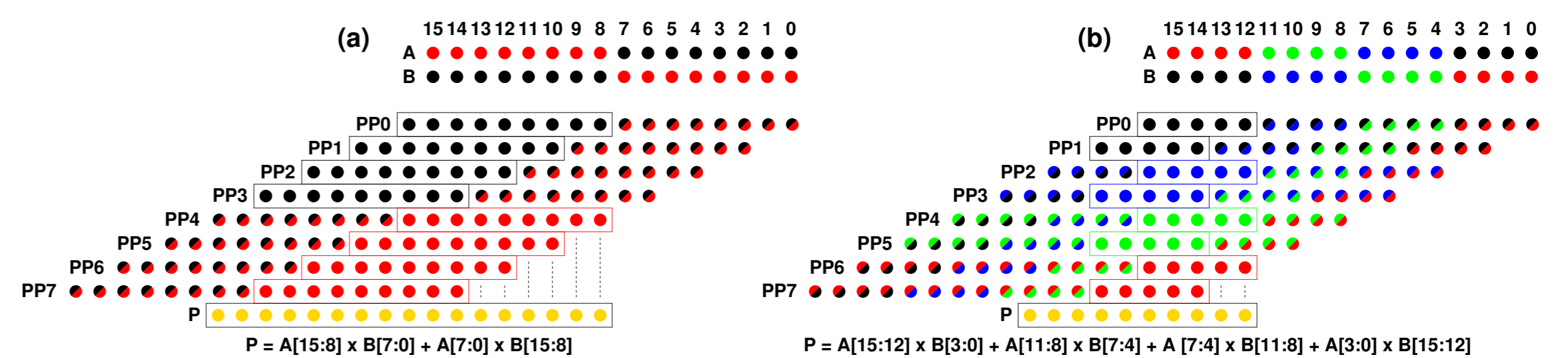


Fig. 2: Alignment of PPI partial products for CONFIG 8x8 (a) and 4x4 (b).

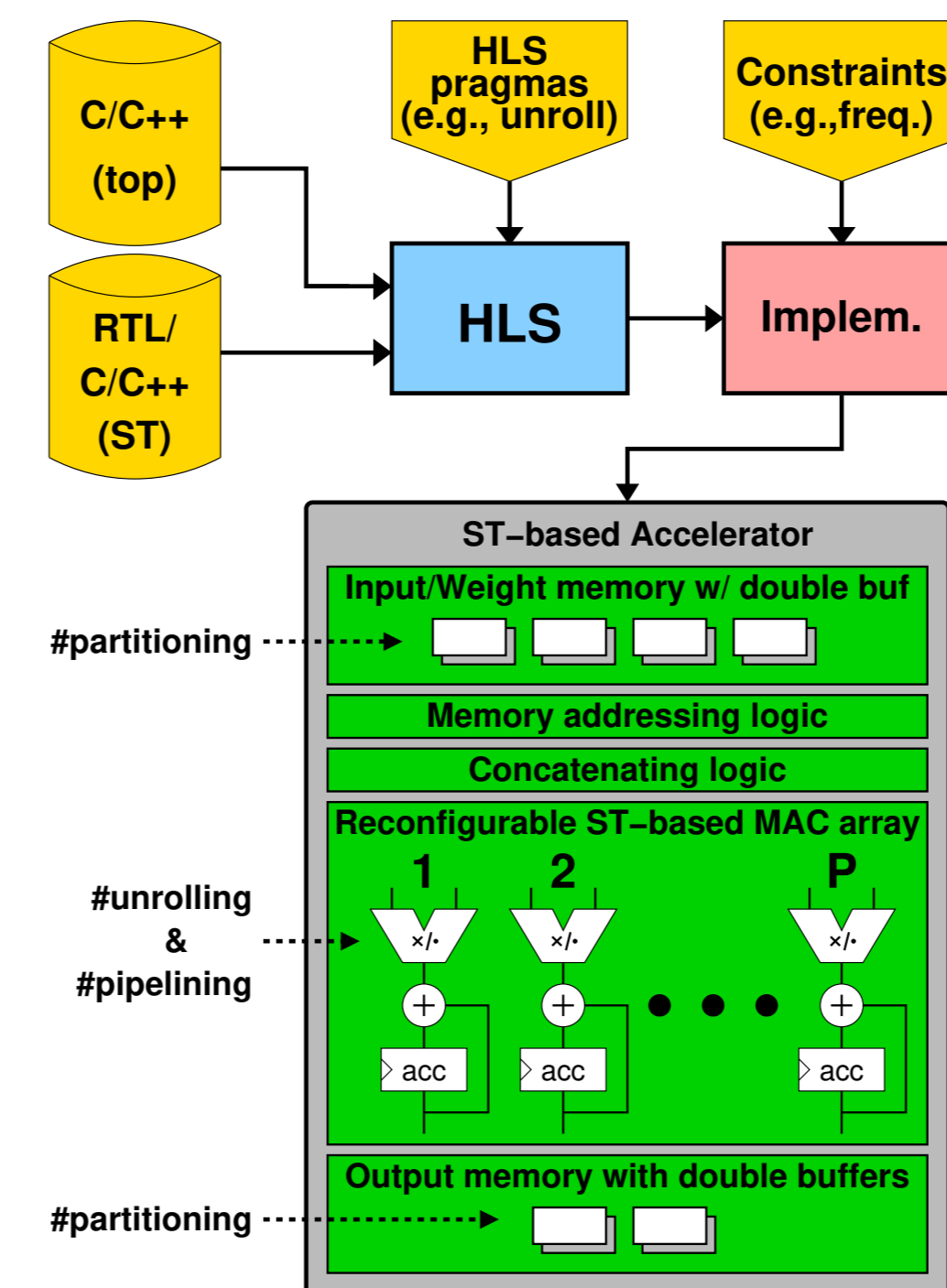


Fig. 3: HLS flow and architecture of ST-based accelerators.

Multiplier	Area [μm^2] @ 1GHz	Avg. Power [mW] @ 1GHz
Non-reconfig.	1133	0.791
Reconfig. [2]	1747 (+54%)	0.903 (+14%)
Reconfig. [3]	1718 (+52%)	0.896 (+13%)
Reconfig. [4]	1629 (+44%)	0.885 (+12%)
Reconfig. Ours	1248 (+10%)	0.893 (+13%)

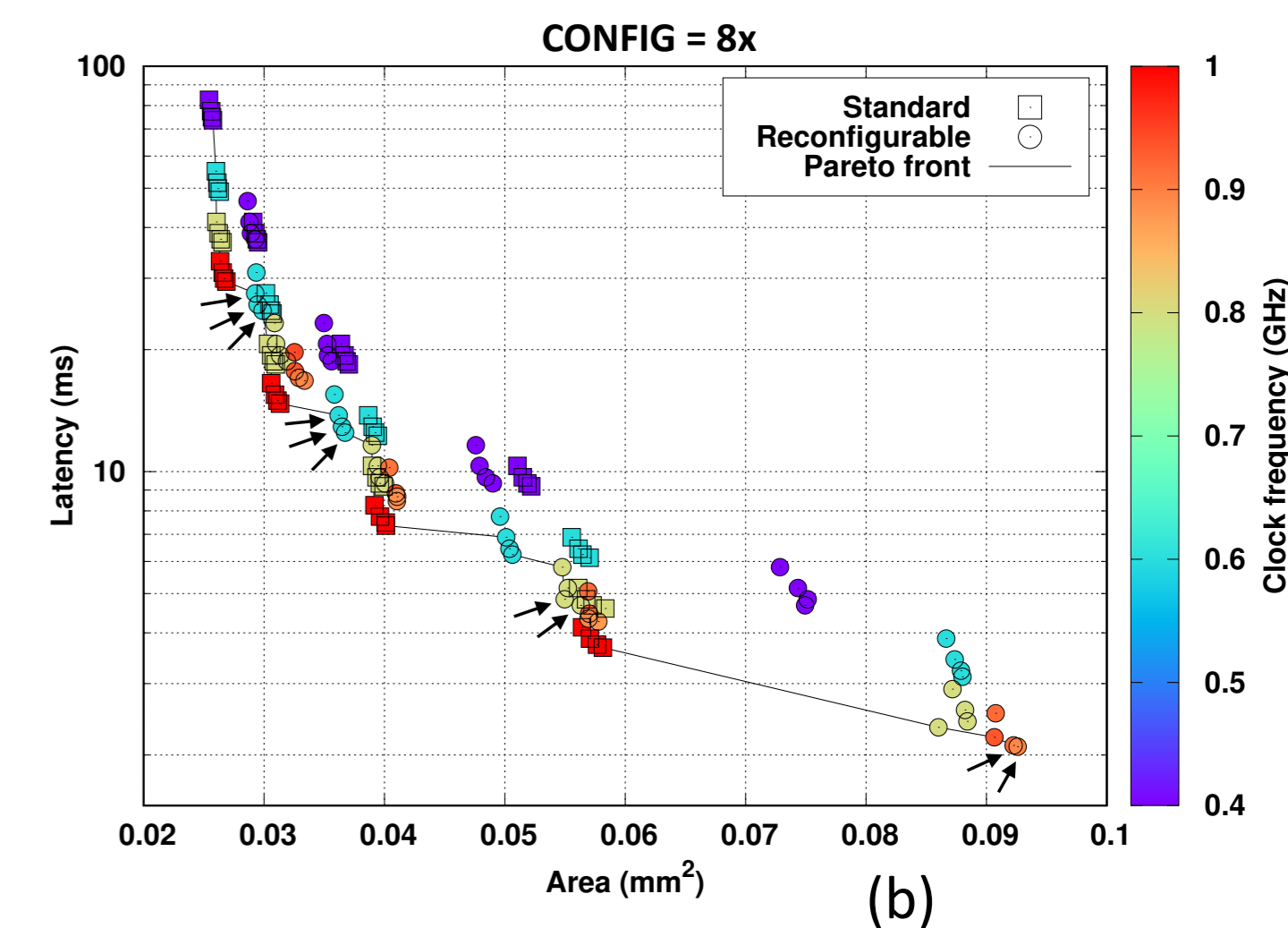


Fig. 4: Latency vs area design-space of the 2D-Conv ST-based accelerator for CONFIG=8x (b) and 4x (c).

Tab. 1. Reconfigurable multipliers vs the baseline non-reconfigurable 16-bit Booth multiplier (at 1-ns target Tck).

List of attended classes

- 01UJBRV – Adversarial training of neural networks (03/06/2021, 25.00)
- 01TRARV – Big data processing and programming (23/03/2021, 33.33)
- 02LWHRV – Communication (08/02/2021, 6.67)
- 01RRPRV – Lean startup e lean business for l'innovation management (15/7/2021, 33.33)
- 01MNFUI – Parallel and distributed computing (26/7/2021, 41.67)
- 01UNYRV – Personal branding (10/1/2022, 1.33)
- 01RISRV – Public speaking (07/4/2021, 6.67)
- 01TAHIU – Quantum computing (18/6/2021, 33.33)
- 01MMRRV – Tecniche numeriche avanzate per l'analisi ed il progetto di antenne (09/6/2021, 33.33)
- 01DOCRV – The Hitchhiker's Guide to the Academic Galaxy [...] (16/6/2022, 26.67)
- 01UNXRV – Thinking out of the box (9/12/2020, 1.33)
- 01SWPRV – Time management (19/1/2021, 2.67)
- External – 2021 IEEE Summer School of Information Engineering, University of Padova (16/7/2021, 19.0)
- External – 2021 Summer School Società Italiana di Elettronica, University of Trieste (9/7/2021, 17.0)
- External – Structuring Machine Learning Projects, Coursera (Andrew Ng) (26/3/2021, 5.0)