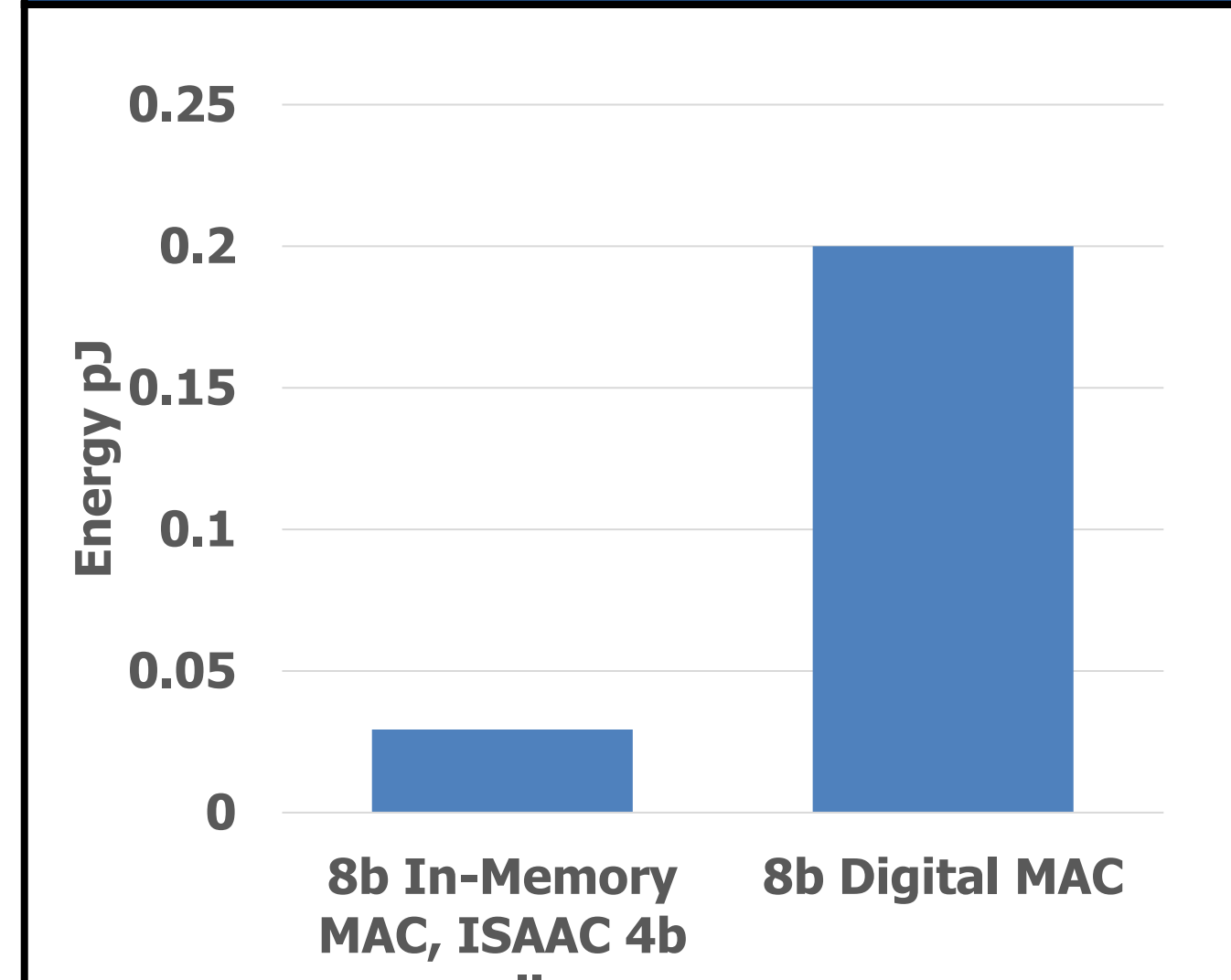


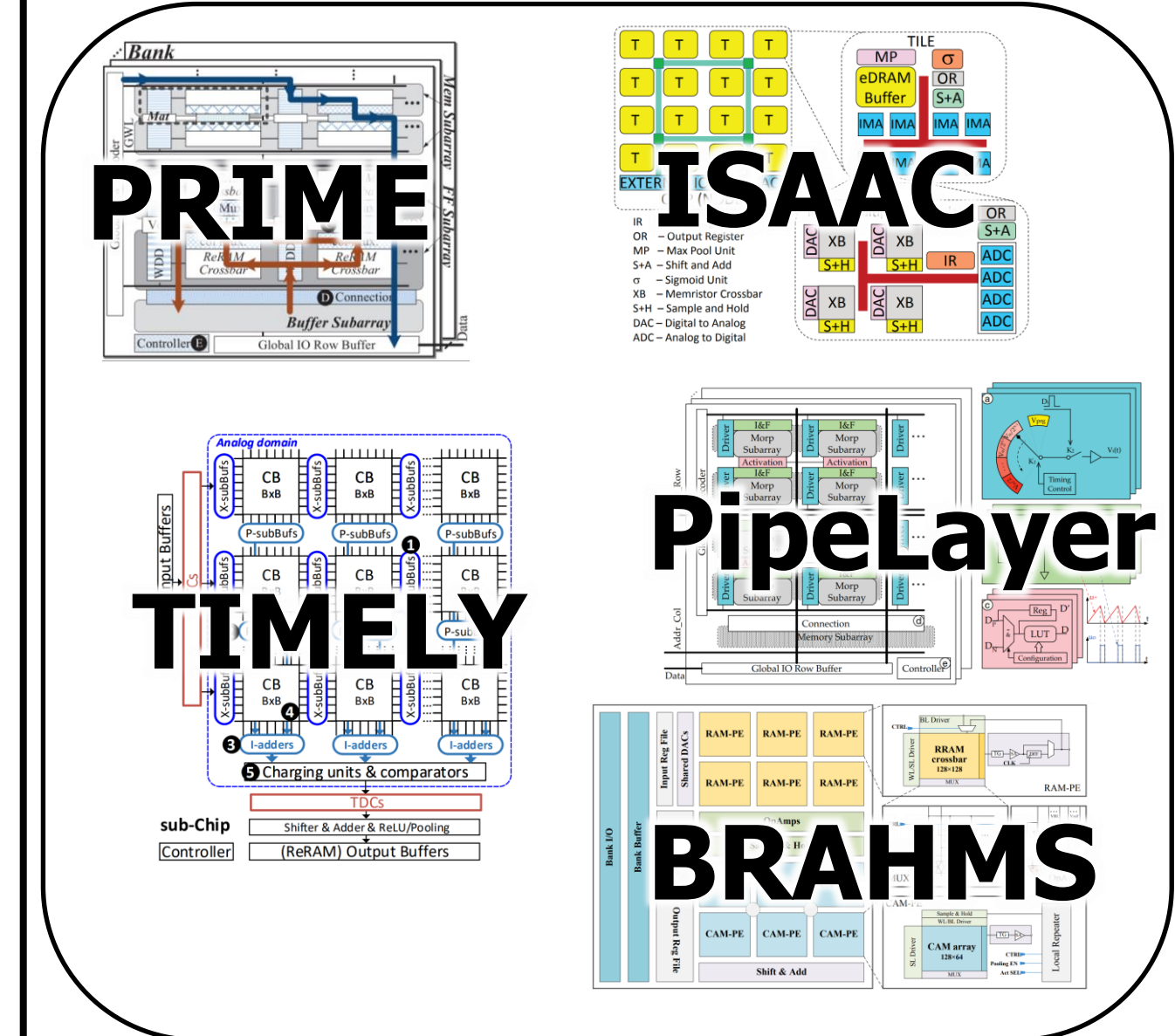
Architectural Evaluation of Processing-In-Memory Systems

Tanner Andrulis, Joel Emer, Vivienne Sze, Massachusetts Institute of Technology

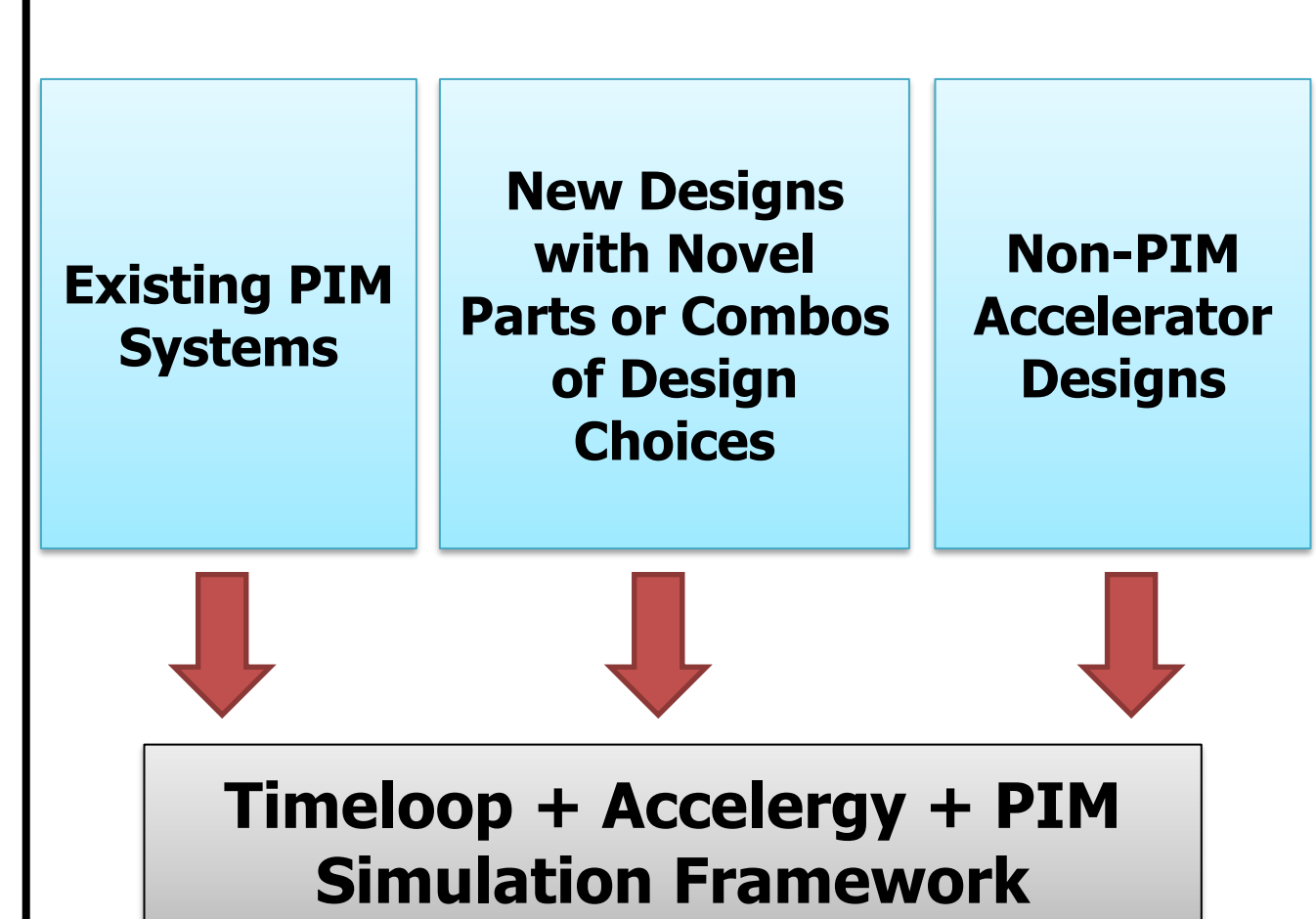
Motivation



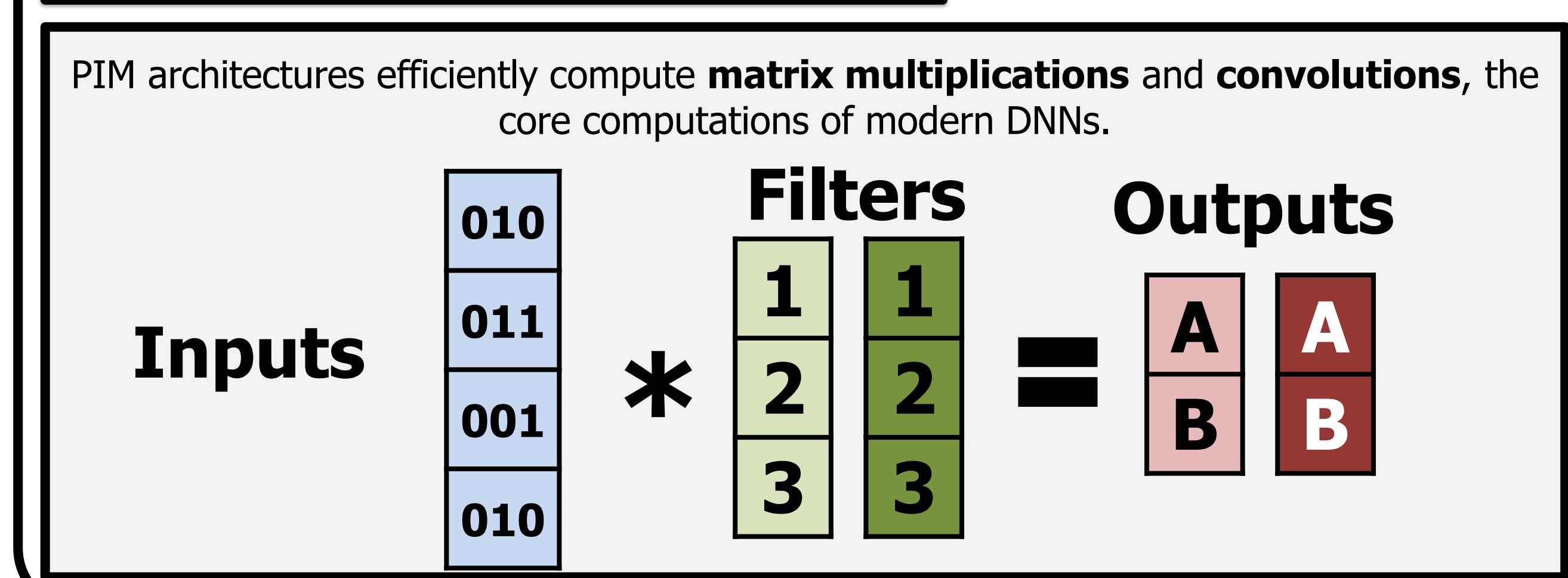
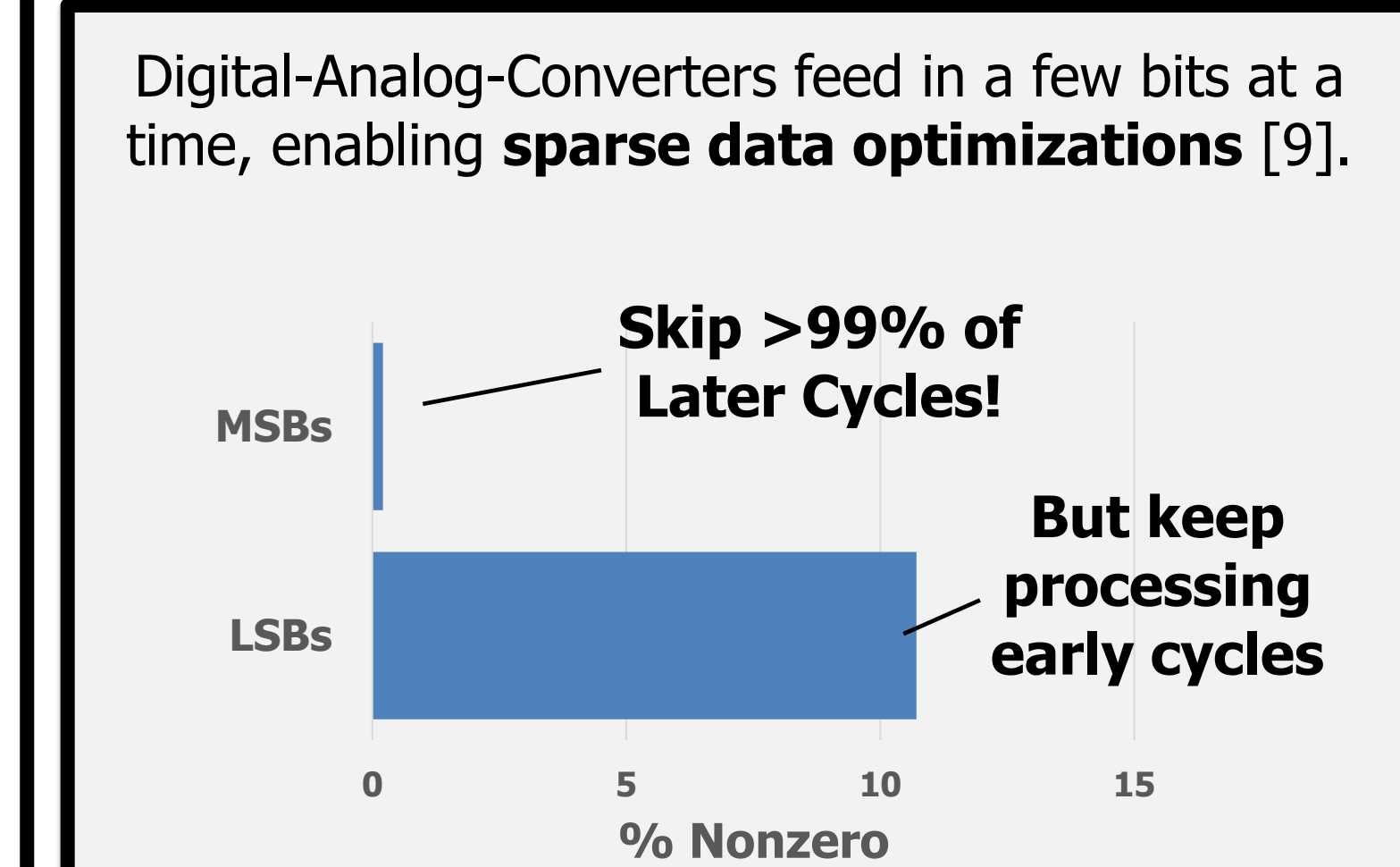
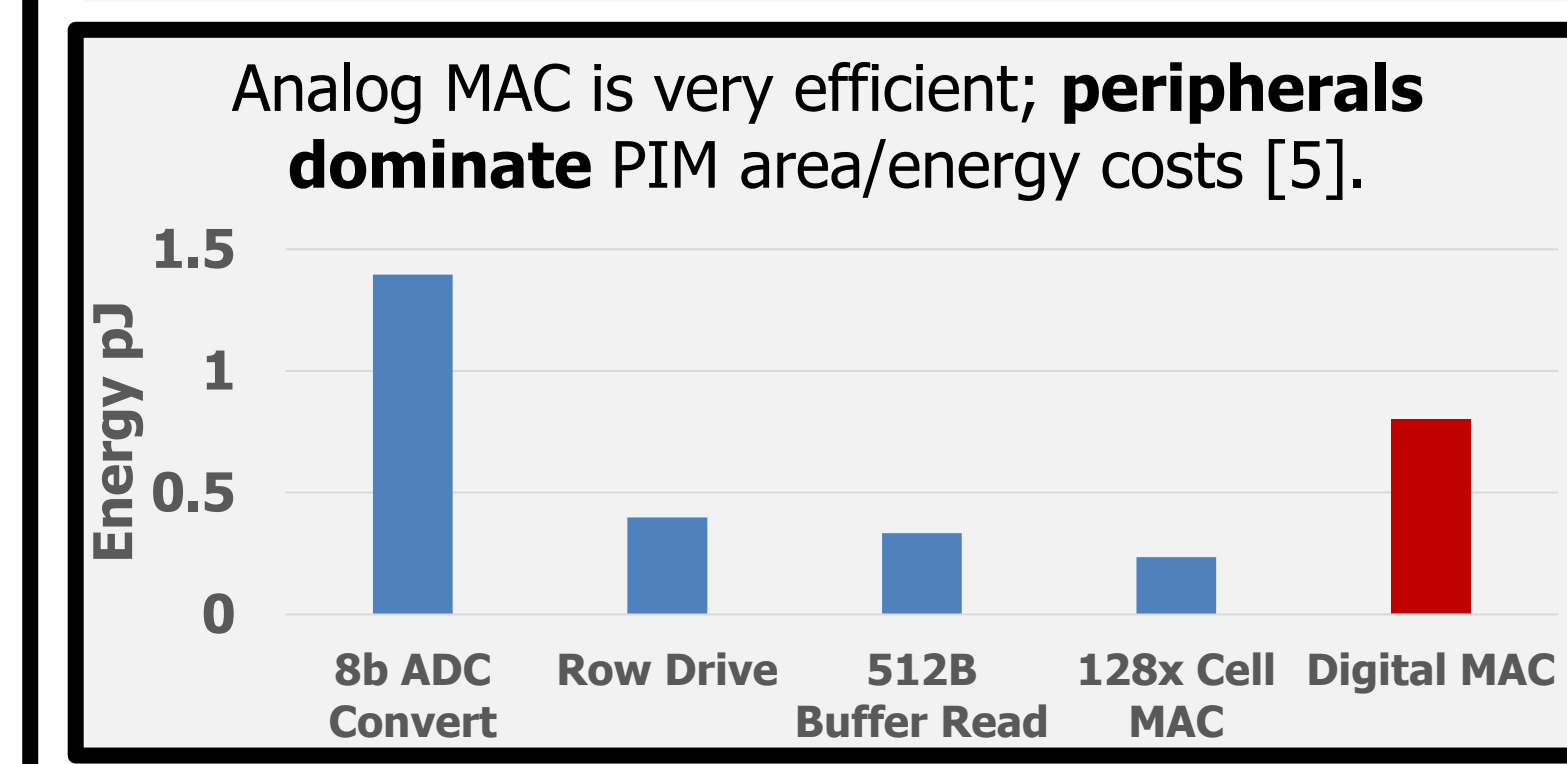
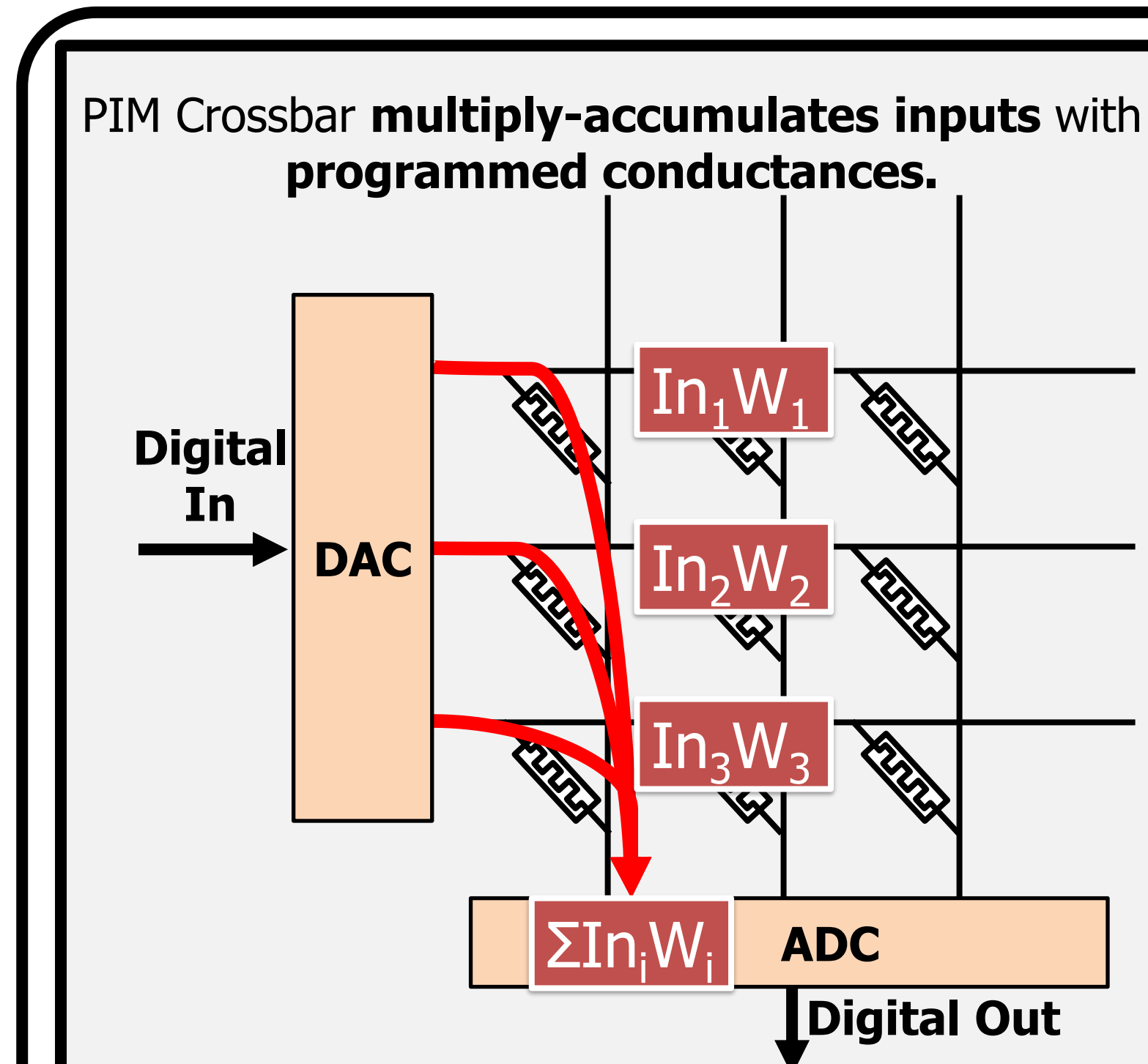
PIM efficiency has inspired many designs



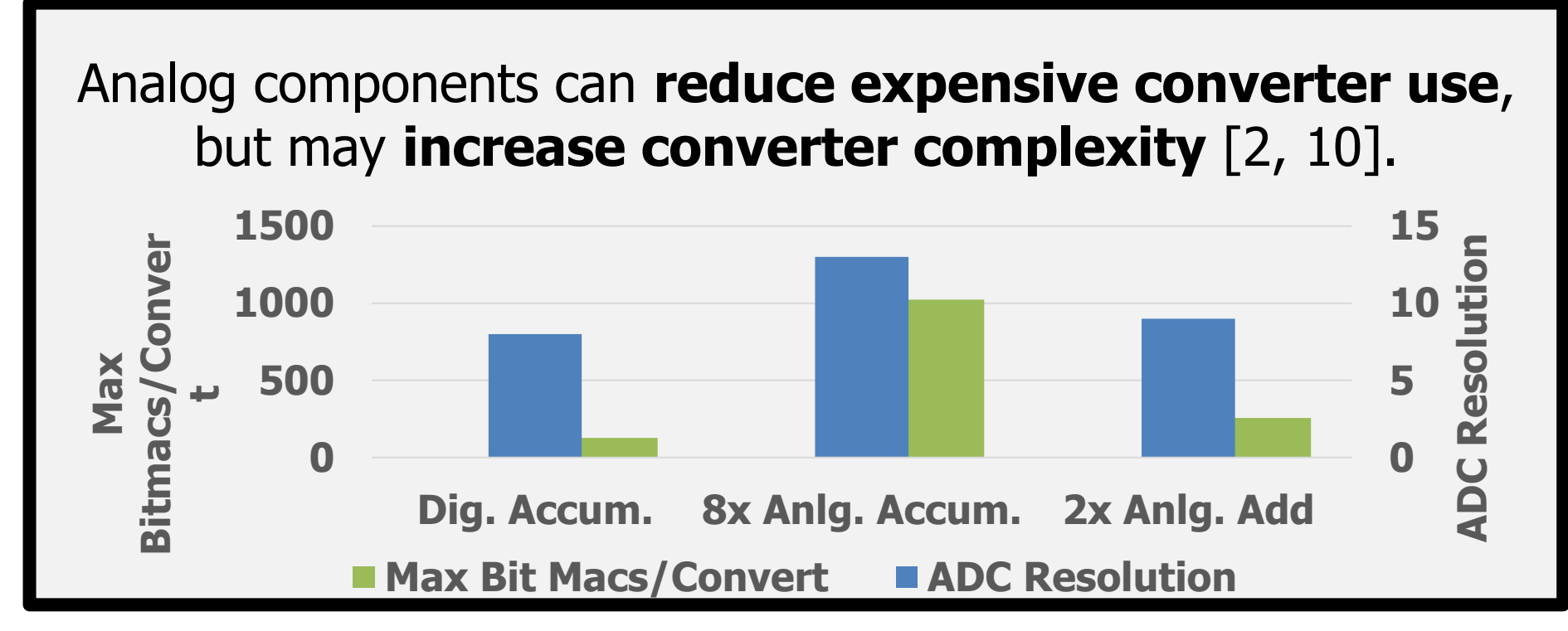
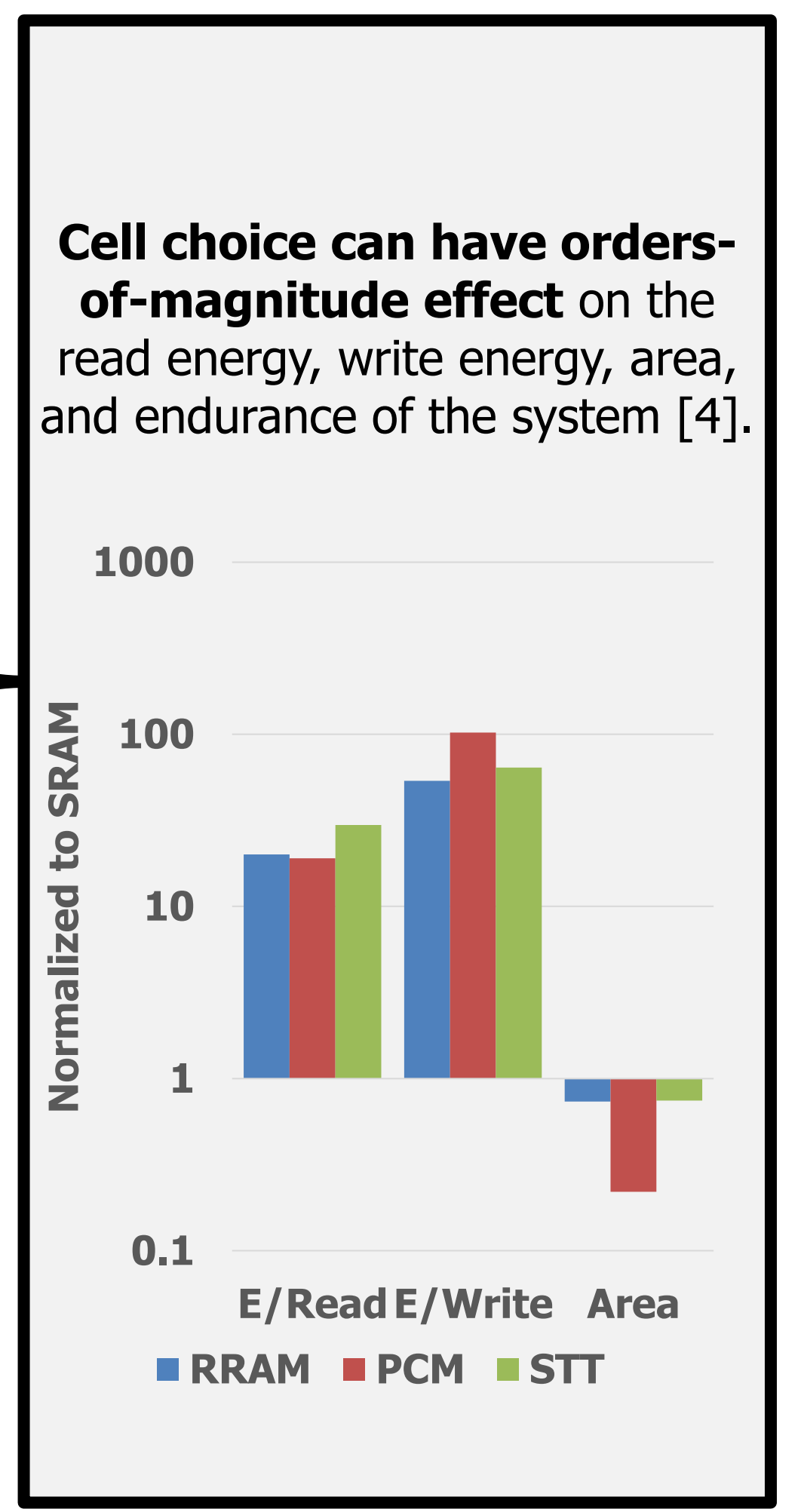
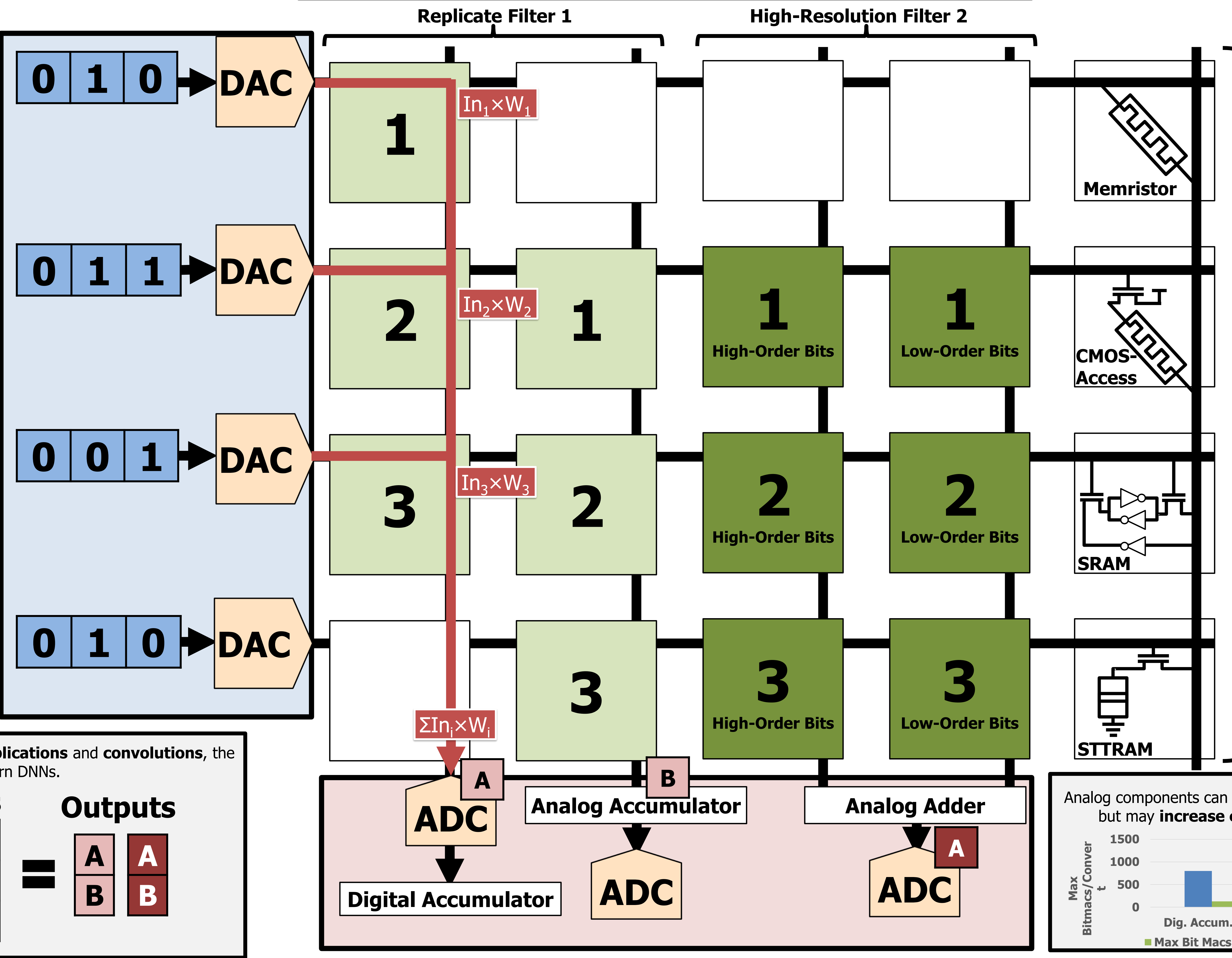
Understand the space with PIM Modeling Framework



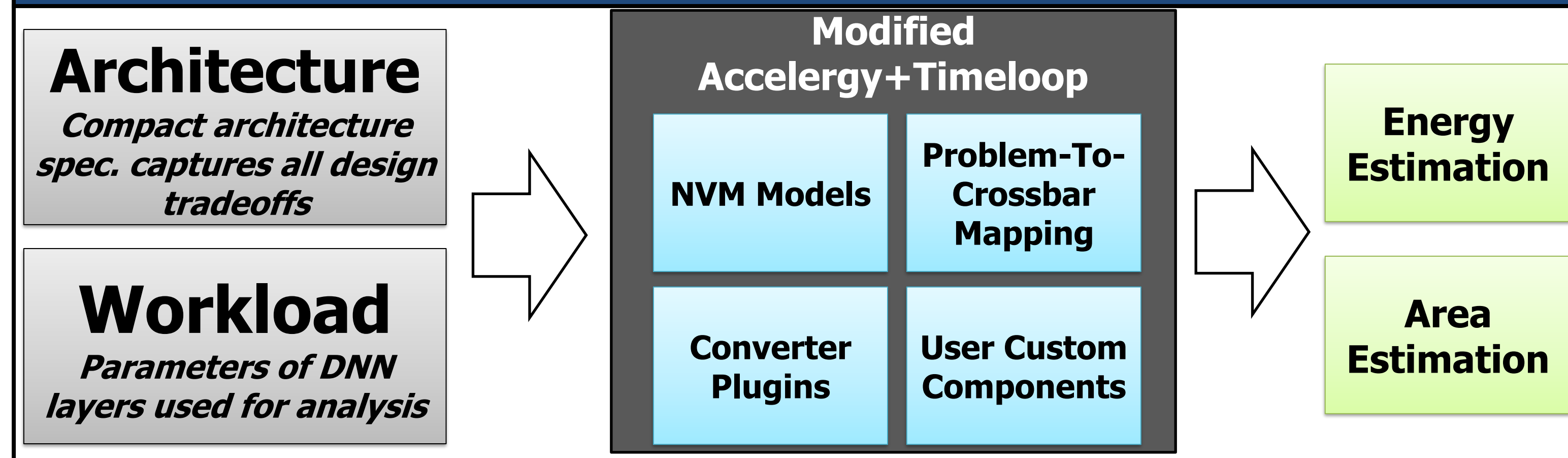
- ✓ Full Understanding of Design Space
- ✓ Fair Comparison of Various Systems and Design Choices
- ✓ Fair Comparison of PIM and Non-PIM Architectures
- ✓ Discovery of New Architectures



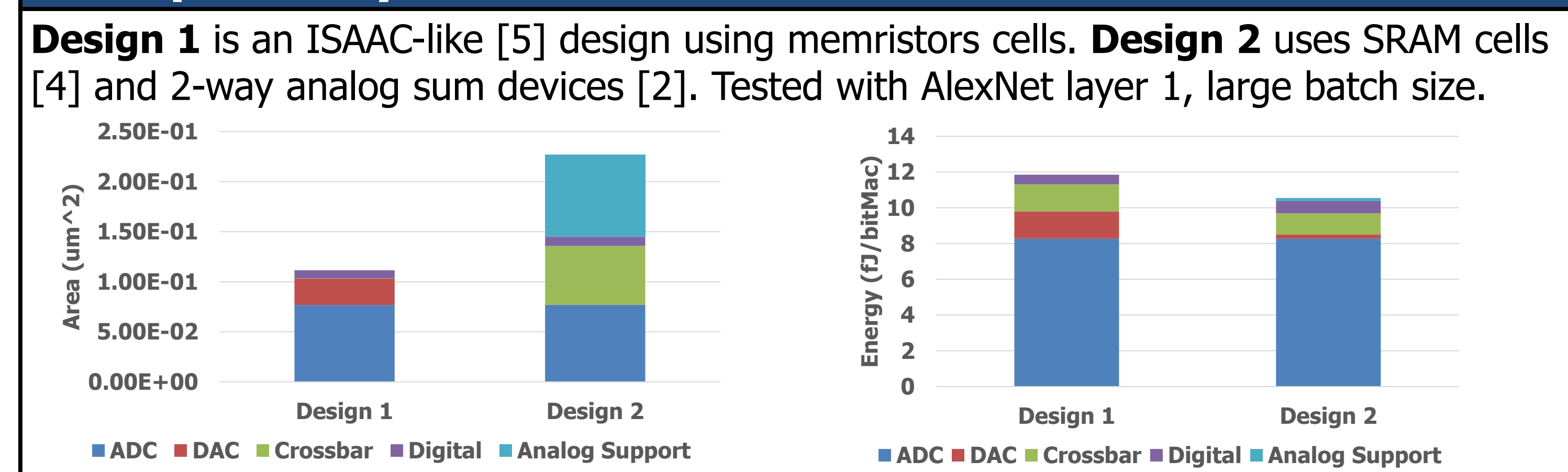
Data in PIM-Crossbars can be replicated to complete multiple convolution steps / vector multiplications at once or stored across multiple devices to increase resolution [2].



Infrastructure



Example Analysis



References

- [1] P. Chi et al., "PRIME: A Novel Processing-In-Memory Architecture for Neural Network Computation in RoRAM-Based Main Memory," in 2016 ACM/IEEE 33rd Annual International Symposium on Computer Architecture (ISCA), Jun. 2016, pp. 27–39. doi: 10.1109/ISCA.2016.13.
- [2] W. Li, P. Xu, Y. Zhao, H. Li, Y. Xie, and Y. Lin, "TIMELY: Pushing Data Movements and Interfaces in PIM Accelerators Towards Local and in Time Domain," arXiv:2005.01206 [cs, stat], May 2020. Accessed: Nov. 26, 2021. [Online]. Available: <https://arxiv.org/abs/2005.01206>
- [3] A. Parashar et al., "Timeloop: A Systematic Approach to DNN Accelerator Evaluation," in 2019 IEEE International Symposium on Performance Analysis of Systems and Software (SPASS), Mar. 2019, pp. 304–315. doi: 10.1109/SPASS.2019.00042.
- [4] A. Pentecost, A. Hankins, M. Donato, M. Homjadi, G. Y. Wei, and D. Brooks, "NVMExplains: A Framework for Cross-Stack Comparisons of Embedded Non-Volatile Memories," arXiv:2109.01188 [cs, stat], Jan. 2022. Accessed: Feb. 24, 2022. [Online]. Available: <https://arxiv.org/abs/2109.01188>
- [5] A. Stadler et al., "ISAAC: A Convolutional Neural Network Accelerator with In-Situ Analog Arithmetic in Crossbars," in 2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA), Jun. 2016, pp. 14–26. doi: 10.1109/ISCA.2016.12.
- [6] Y. Song, X. Qian, H. Li, and Y. Chen, "PipeLayer: A Pipelined RoRAM-Based Accelerator for Deep Learning," in 2017 IEEE International Symposium on High Performance Computer Architecture (HPCA), Feb. 2017, pp. 541–552. doi: 10.1109/HPCA.2017.55.
- [7] T. Song, X. Chen, X. Zhang, and Y. Han, "BRAHMS: Beyond Conventional RoRAM-based Neural Network Accelerators Using Hybrid Analog Memory System," in 2021 56th ACM/IEEE Design Automation Conference (DAC), Dec. 2021, pp. 1033–1038. doi: 10.1109/DAC18074.2021.9588247.
- [8] Y. N. Wu, J. S. Emer, and V. Sze, "Accelry: An Architecture-Level Energy Estimation Methodology for Accelerator Designs," in 2019 IEEE/ACM International Conference on Computer-Aided Design (ICCAD), Nov. 2019, pp. 1–8. doi: 10.1109/ICCAD45719.2019.8942149.
- [9] J. Zhang, H. Tang, F. Chen, Y. Wang, and H. Li, "Exploring Bit-Slice Sparsity in Deep Neural Networks for Efficient RoRAM-Based Deployment," arXiv:1909.08496 [cs, stat], Nov. 2019. Accessed: Dec. 14, 2021. [Online]. Available: <https://arxiv.org/abs/1909.08496>
- [10] T. Chou, W. Tang, J. Botmer, and Z. Zhang, "CASCADE: Connecting RoRAMs to Extend Analog Dataflow in An End-To-End In-Memory Processing Paradigm," in Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture, Columbus OH USA, Oct. 2019, pp. 114–125. doi: 10.1145/3352460.3358328.