

Implementation and Utilization of Instruction Level Parallelism in Modern CPU Architectures

Justin Reed

Advisor:
David Wonnacott

Background: How can we increase the speed of the CPU?
Fundamental ideas behind hardware parallelism.

Background
Control Unit

Already familiar with the
basics of pipelining? Skip to
the next section.

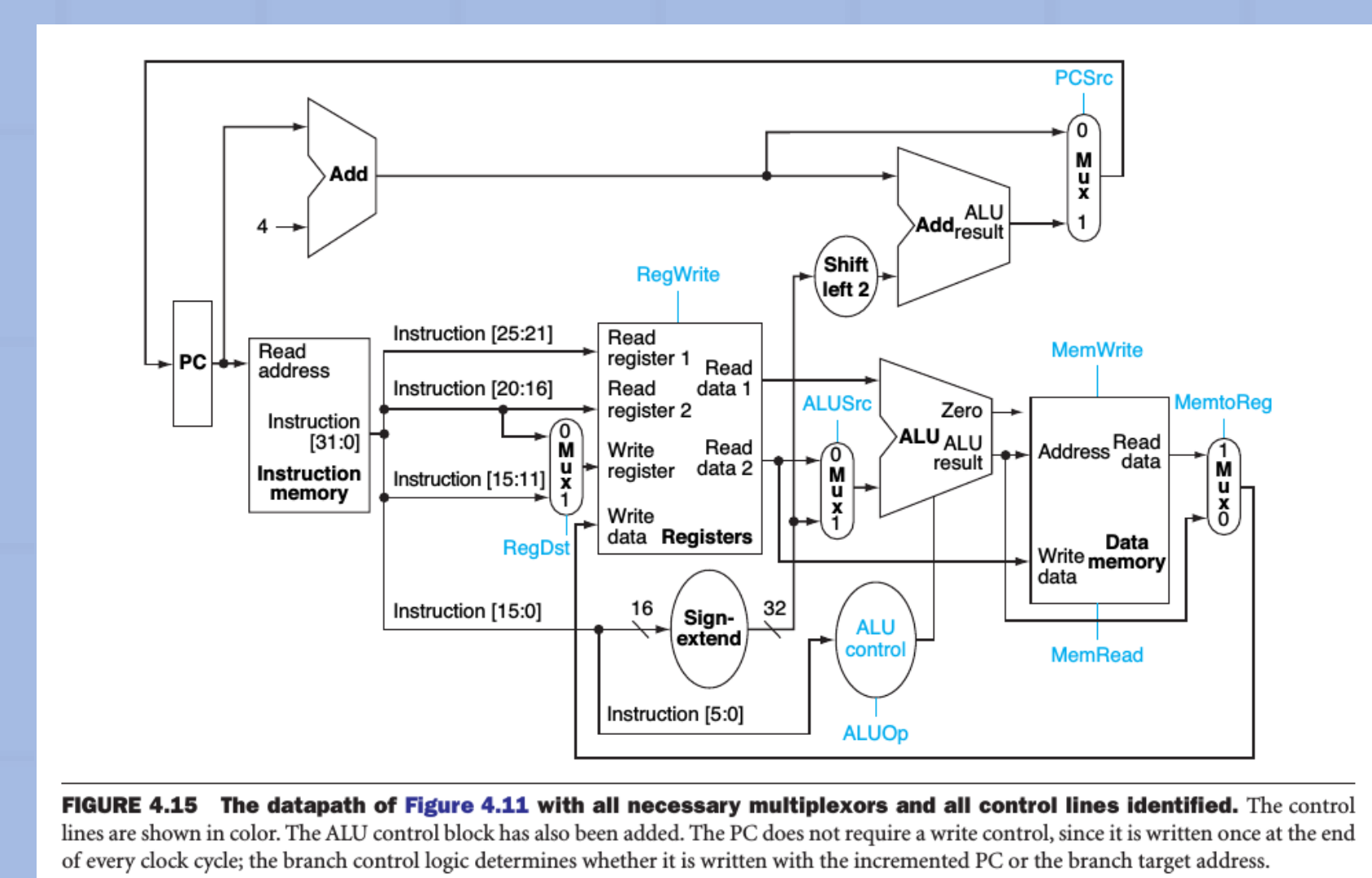


Fig 1. Decrease clock cycle length using multicycling [PH17a]

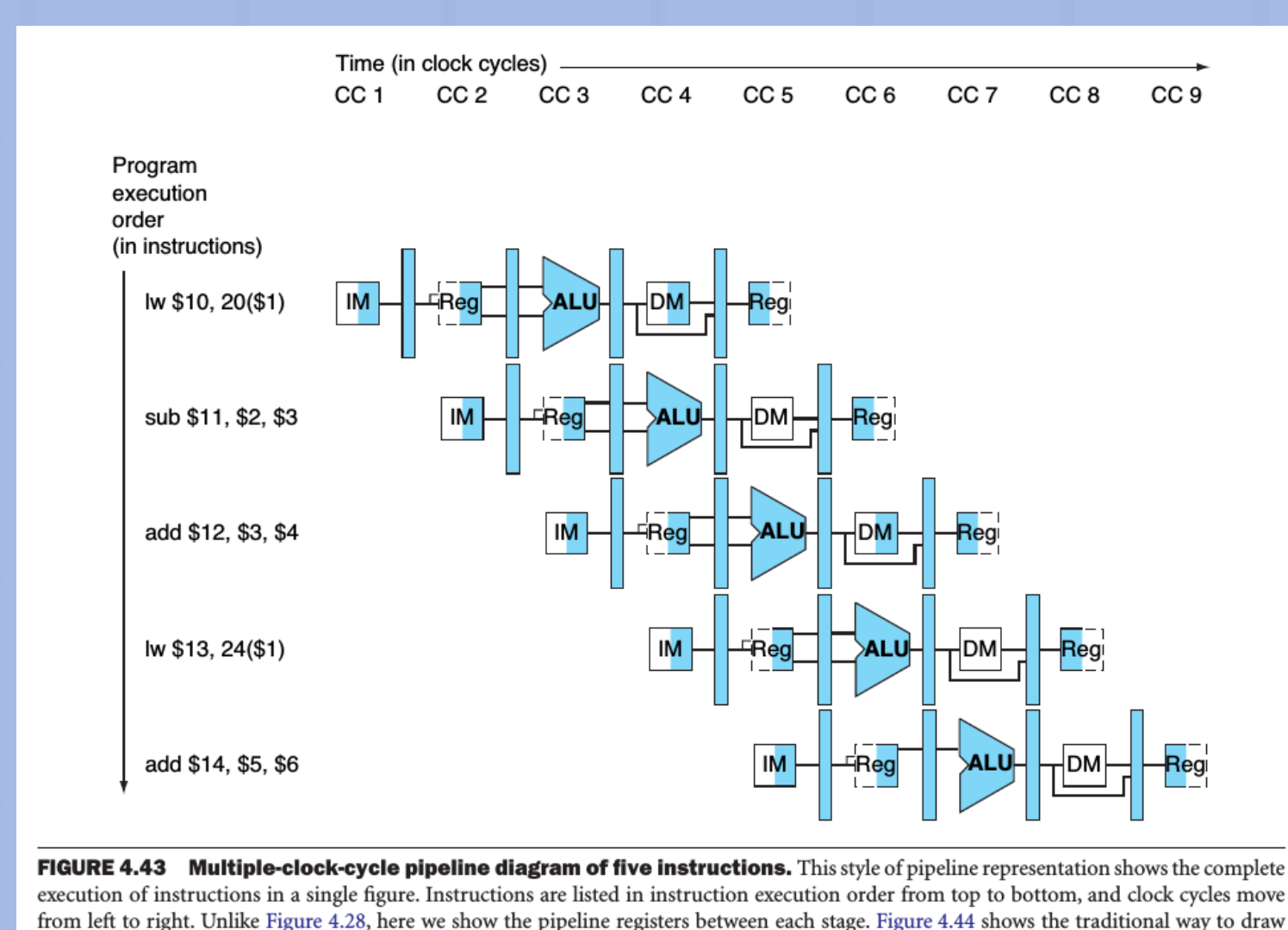


Fig 2. Maximize potential of multicycled CPU with pipelining [PH17a]

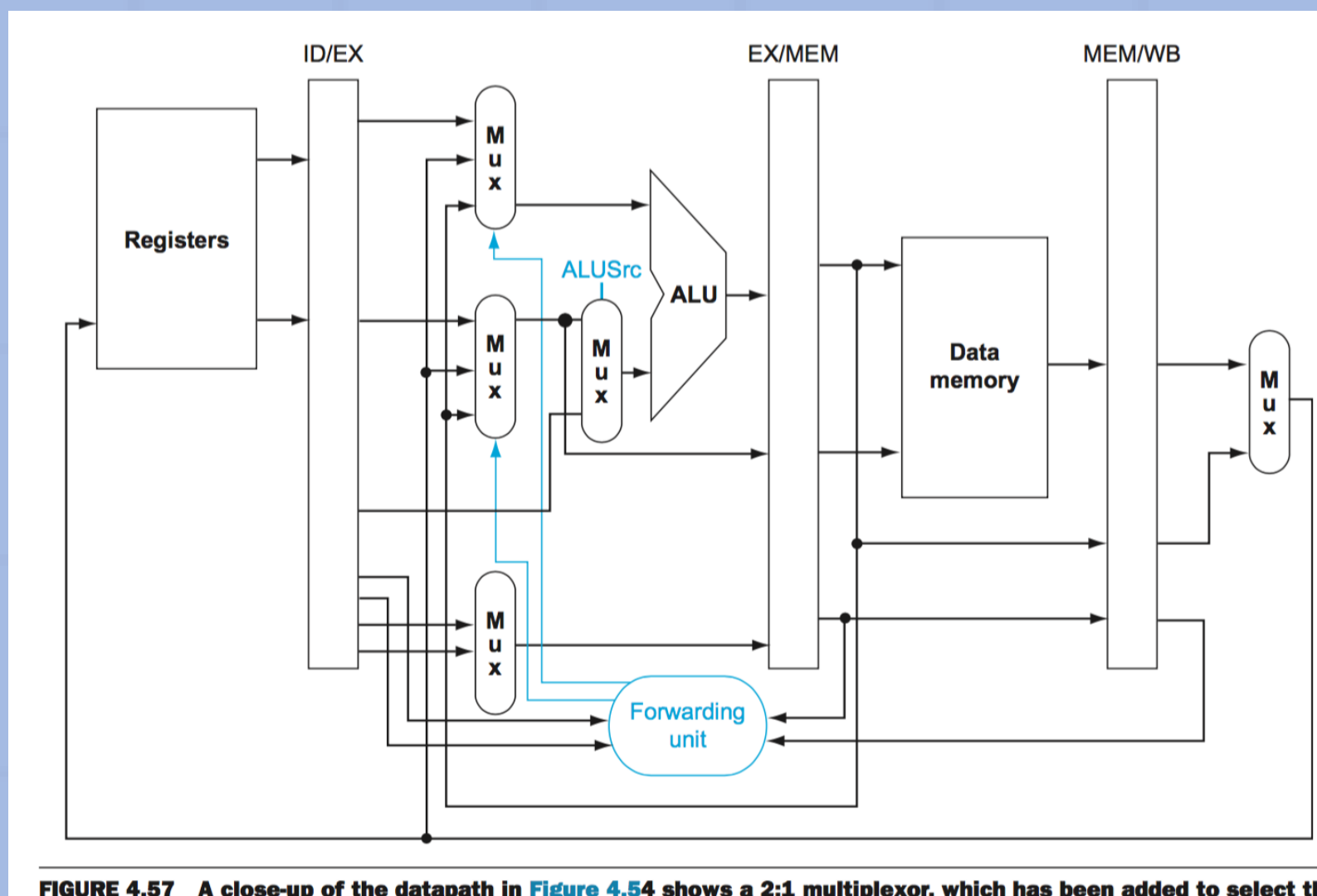


Fig 3. Technical implementation of basic pipelined CPU [PH17a]

Summary: CPU pipelining
and its implications

Though the multi-cycle instruction is slightly slower than its single-cycle equivalent, it has the key advantage of making a pipelined architecture possible. Such an architecture begins executing new instructions before the previous ones have completed, taking advantage of the fact that each sub-instruction only needs a portion of the chip to execute. Multiple instructions are therefore able to execute simultaneously, improving throughput by a factor close to the number of sub-instructions per multi-cycle instruction.

Simultaneous Multithreading: Making use of instruction and thread level parallelism at the same time.

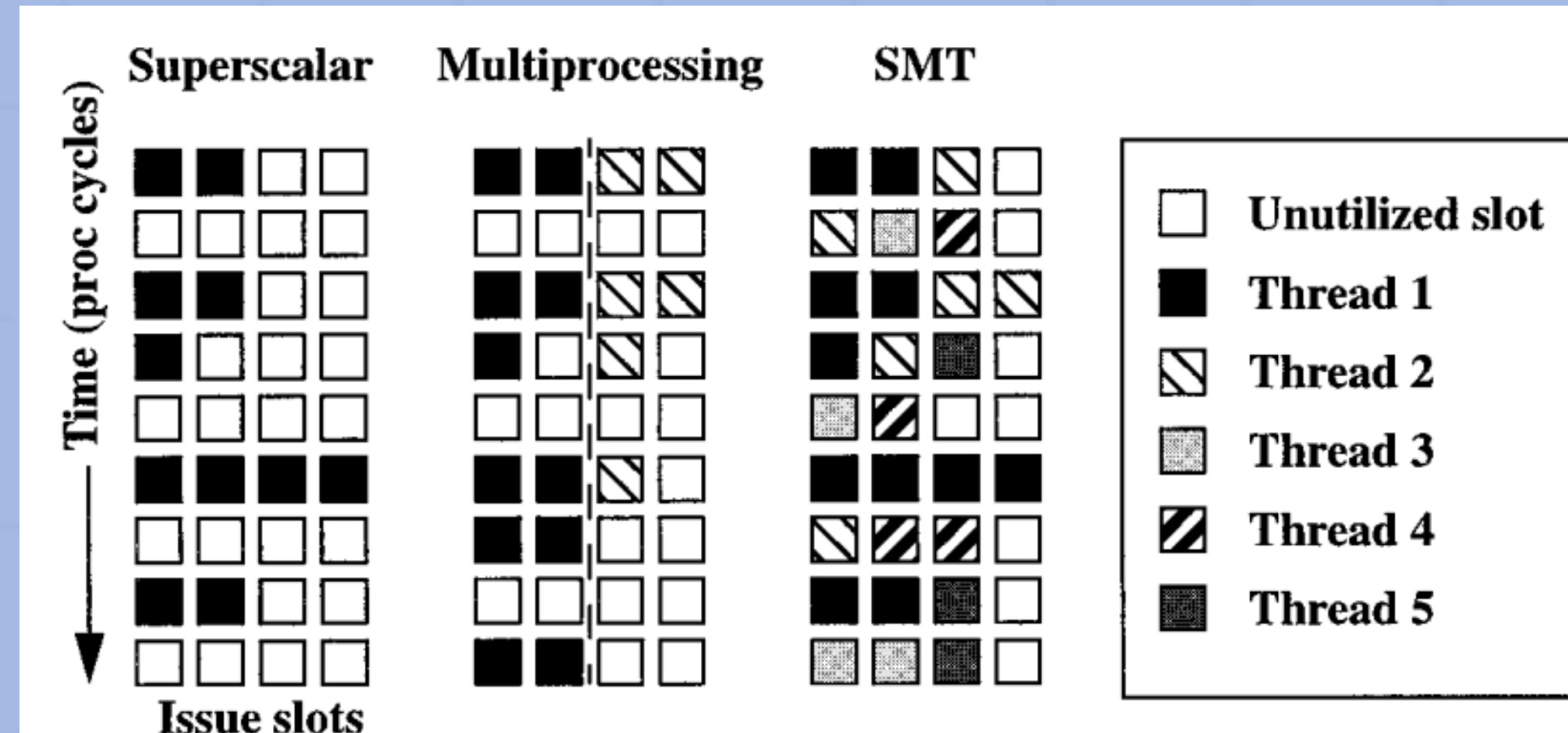


Fig 4. How TLP and ILP are viewed as equivalent by the SMT core [LEL+ 97]

Summary: why SMT matters

The advantage of a SMT CPU is greatest for single-thread applications, as these programs are implicitly converted to multithreaded applications at execution time, allowing for thread-level parallelism where there previously wasn't any. But even for multithreaded applications, SMT CPU's maintained a clear advantage. Since SMT can consider ILP and TLP interchangeably, the SMT CPU achieves high performance for both single and multi-threaded applications. On average, the performance increase from the SMT architecture is 52-percent greater. This performance increase is so decisive that nearly all modern general-purpose processors now make use of SMT.

Recent bottlenecks:
Reducing memory latency

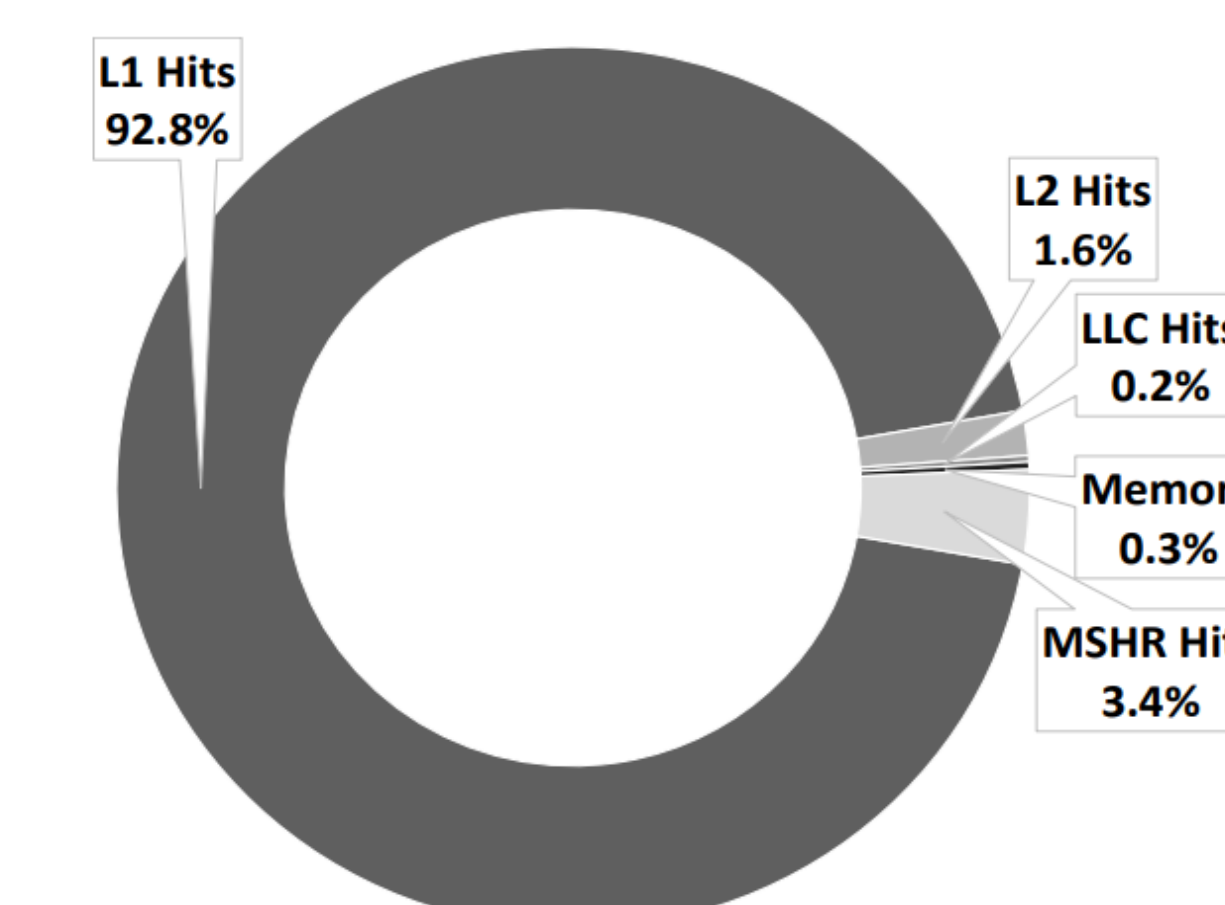


Fig 8. Frequency of cache-misses by memory type [SBG522]

Utilize maximal ILP with out-of-order and in-order execution

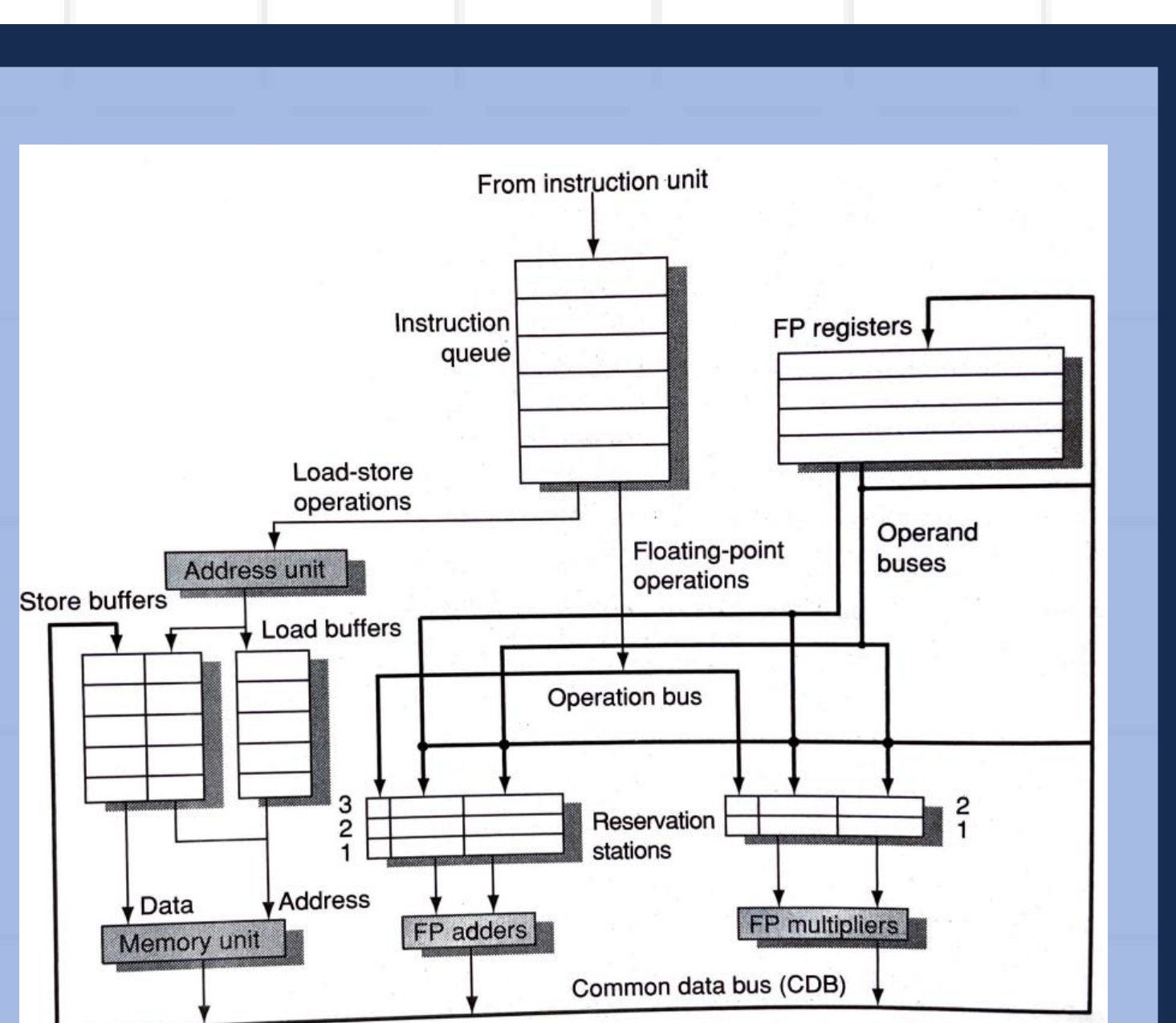


Fig 5. An out-of-order execution processor diagram [PH17b]

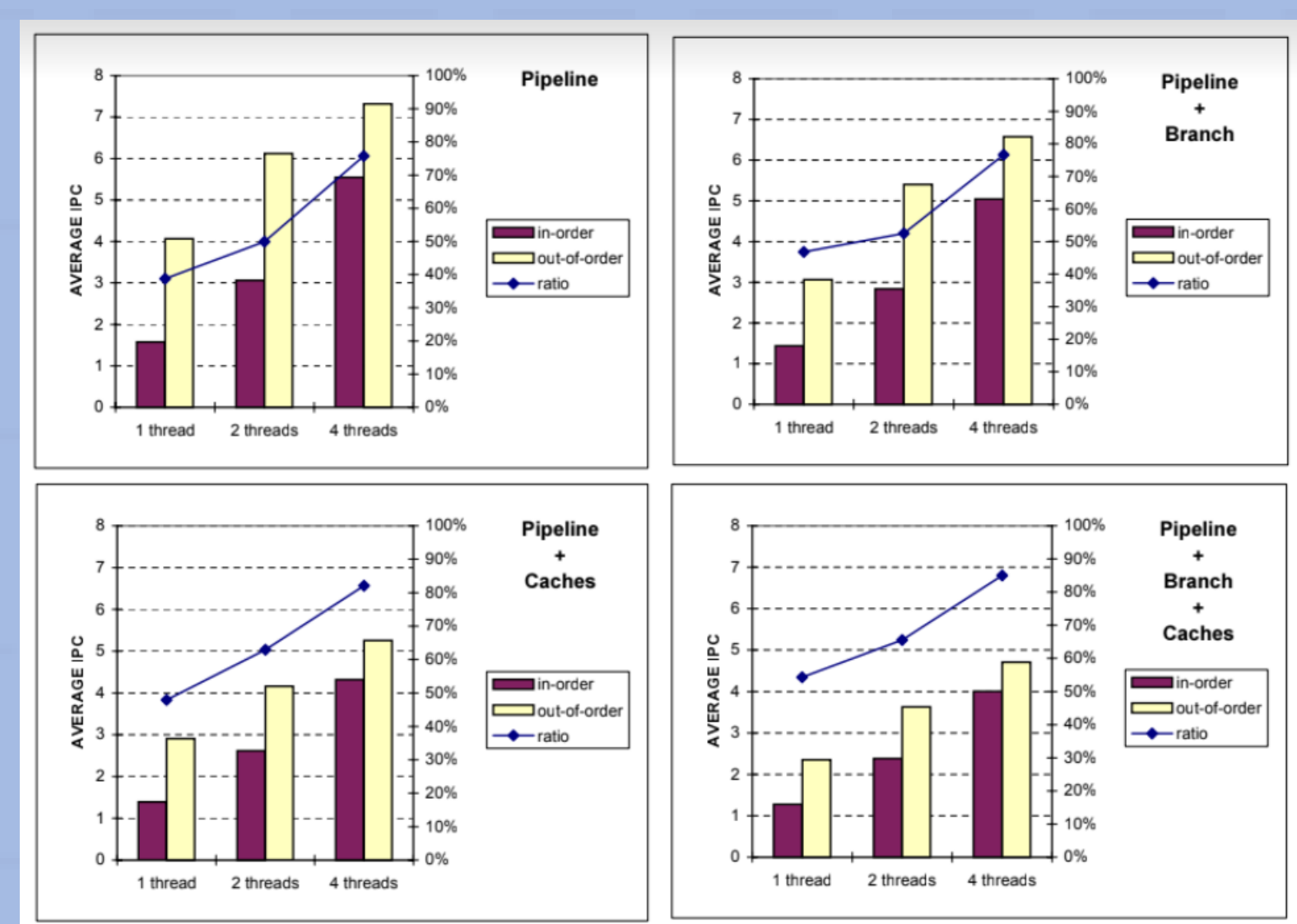


Fig 6. OOO loses its edge in the context of SMT [HS99]

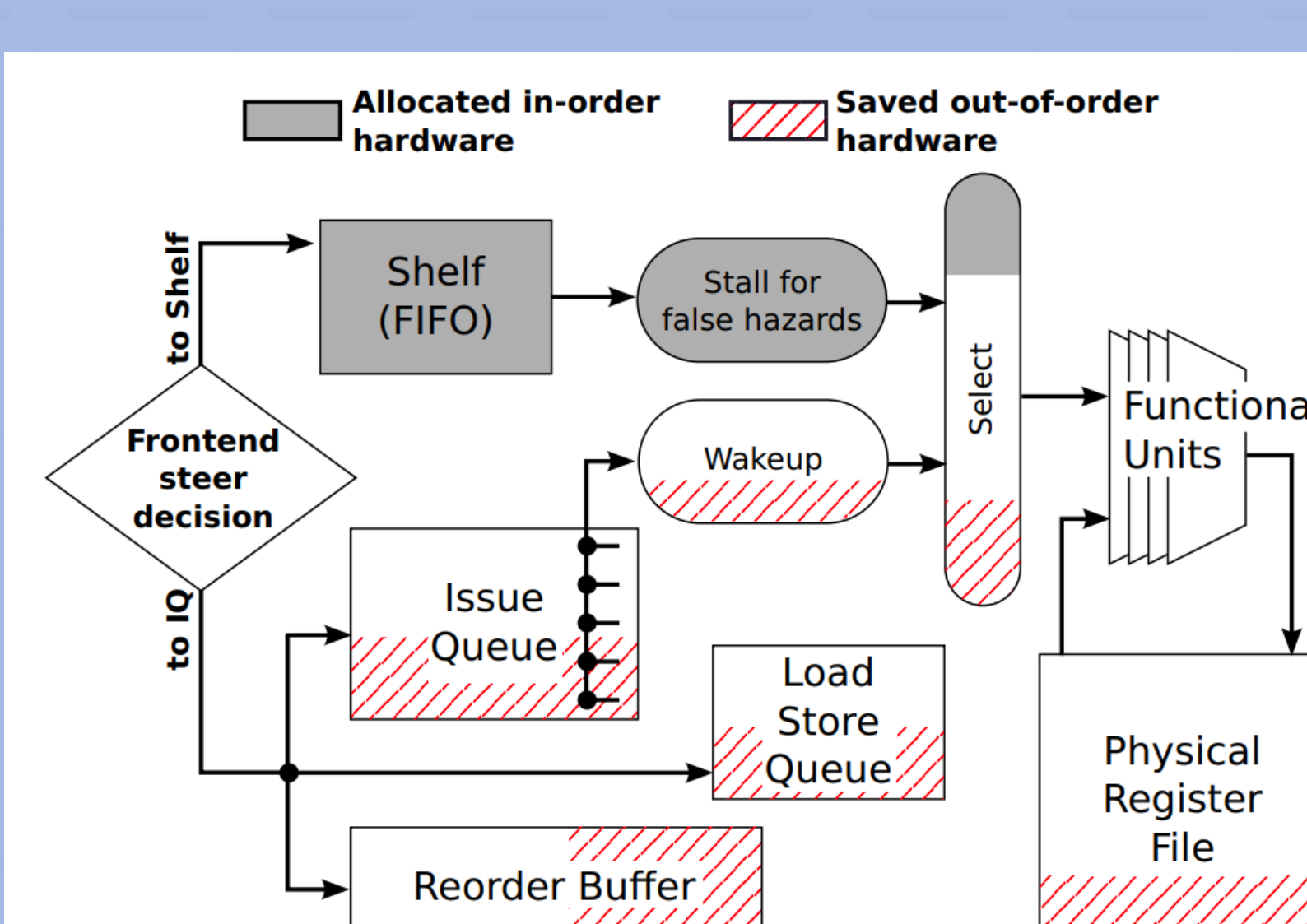


Fig 7. Solution: use a hybrid architecture [SW16]

Sources and
future work

While the results of [HS99] indicate an important result regarding in-order execution, the experiment was done on a 1999 processor. One possible area of future work would be to replicate this paper on a modern processor.

[SW16] Faissal M. Sleiman and Thomas F. Wenisch, "Efficiently Scaling Out-of-Order Cores for Simultaneous Multithreading," in *Proceedings of the 44th International Symposium on Computer Architecture, ISCA '16* (Seoul, Republic of Korea: IEEE Press, 2016), 431–443. <https://doi.org/10.1109/ISCA.2016.65>

[LEL+ 97] Jack L. Lu et al., "Converting Thread-Level Parallelism to Instruction-Level Parallelism via Simultaneous Multithreading," *ACM Trans. Comput. Syst.* 15, no. 3 (August 1997): 322–54. <https://doi.org/10.1145/263105.263108>

[PH17a] Patterson, David A., Author, John L. Hennessy, and Perry Alexander, *Computer Organization and Design: The Hardware/Software Interface*. (Amsterdam : Boston: Elsevier/Morgan Kaufmann, 0, 2017)

[PH17b] J.L. Hennessy and D.A. Patterson, *Computer Architecture: A Quantitative Approach*. (Elsevier Science, 2017). <https://books.google.com/books?id=c8bhuDrsAAQ&pg=PA40>

[HS99] S. Hry and A. Sene, "Out-of-Order Execution May Not Be Cost-Effective on Processors Featuring Simultaneous Multithreading," in *Proceedings of the 19th International Symposium on High-Performance Computer Architecture*, 1999, 64–67. <https://doi.org/10.1109/HPCA.1999.744331>

[SBG522] Subraman Shukla et al., "Register File Prefetching," in *Proceedings of the 40th Annual International Symposium on Computer Architecture, ISCA '22* (New York, NY, USA: Association for Computing Machinery, 2022), 410–23. <https://doi.org/10.1145/3530596.3532288>