

Worst Case Analysis of DRAM Latency in Multi-Requestor Systems

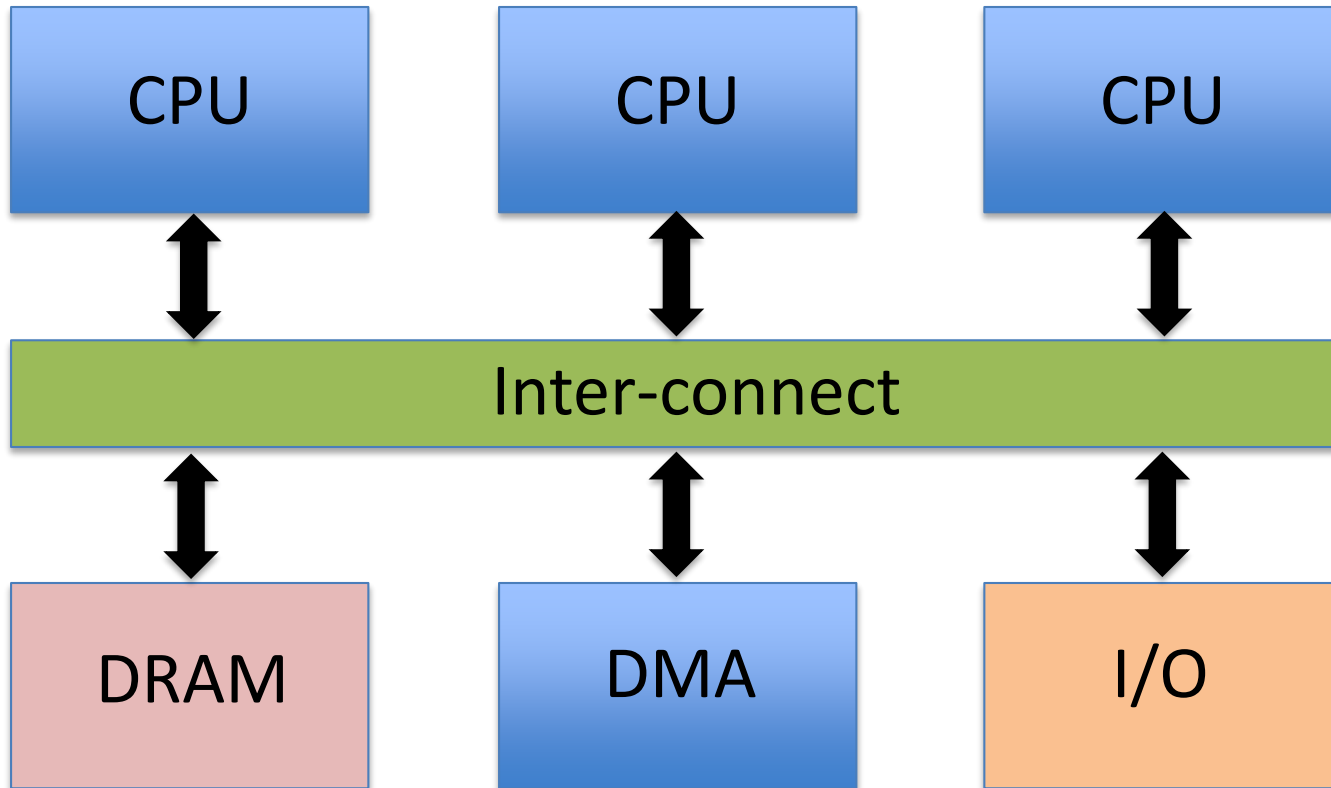
Zheng Pei Wu

Yogen Krish

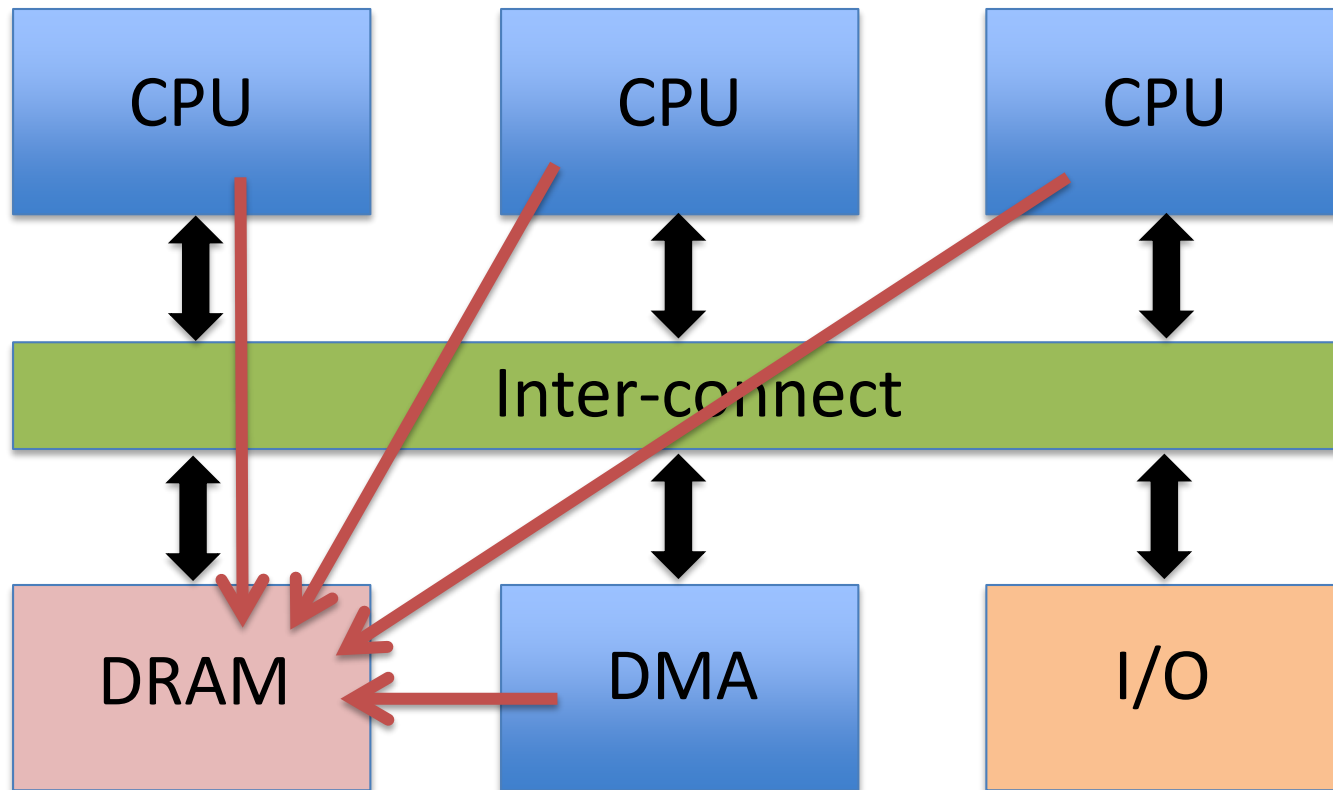
Rodolfo Pellizzoni



Multi-Requestor Systems

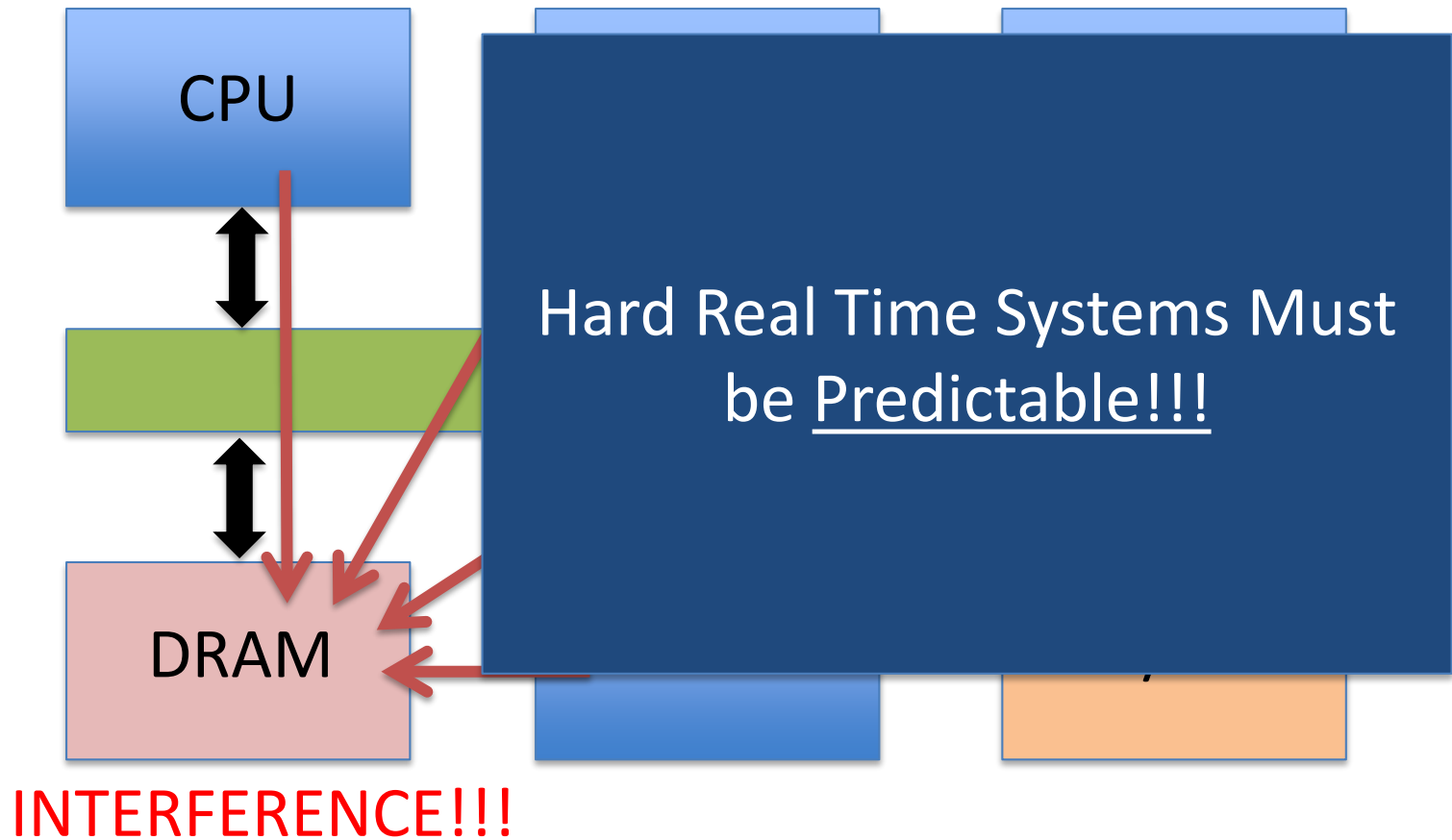


Multi-Requestor Systems



INTERFERENCE!!!

Multi-Requestor Systems



Multi-Requestor Systems

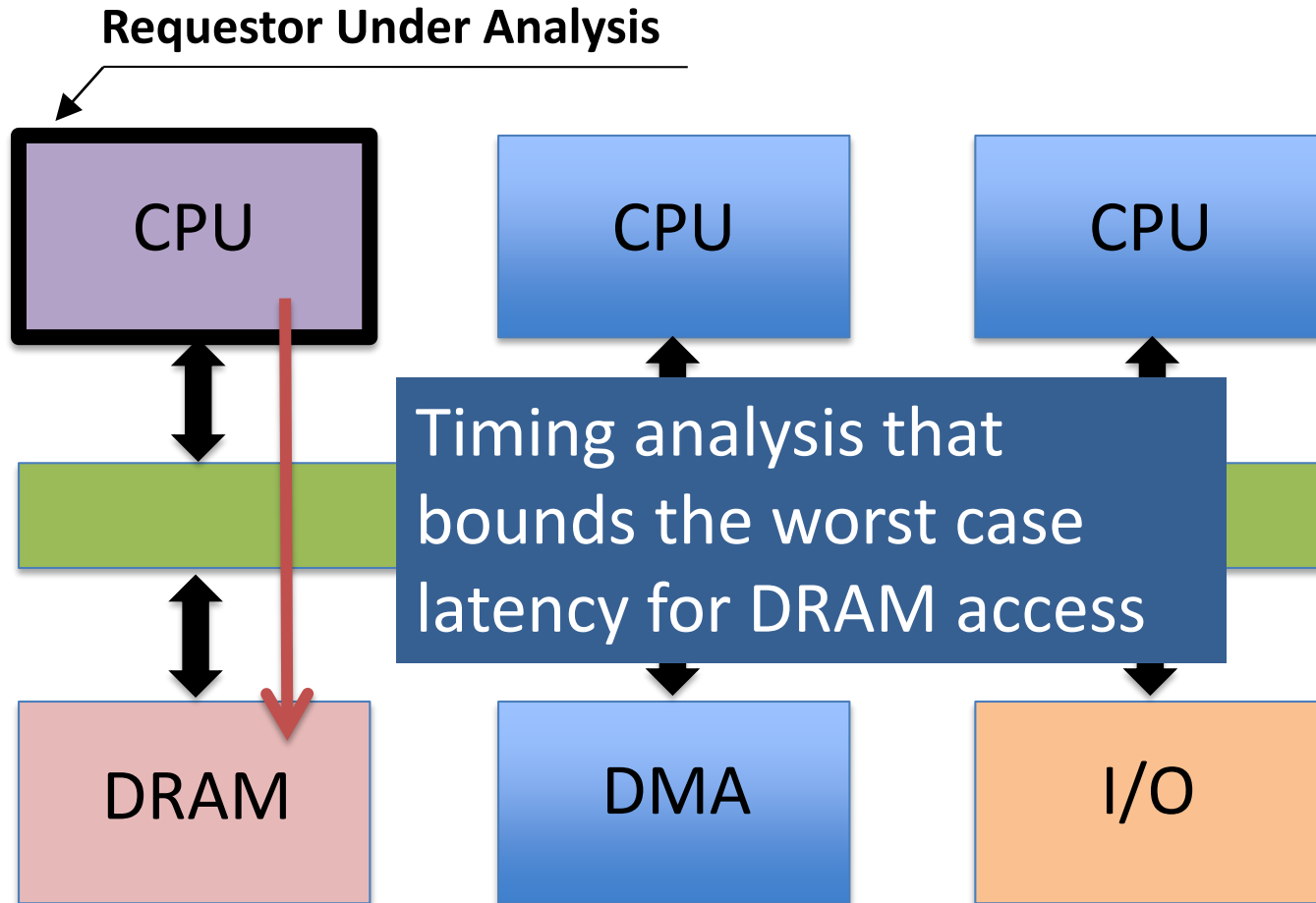
- Schedulability Analysis: needs WCET as input
- WCET depends on hardware platform
- WCET: needs Latency to access shared resource (e.g. cache, DRAM)
- Existing approaches can bound the interference but they **assume the latency for DRAM access is constant**

Multi-Requestor Systems

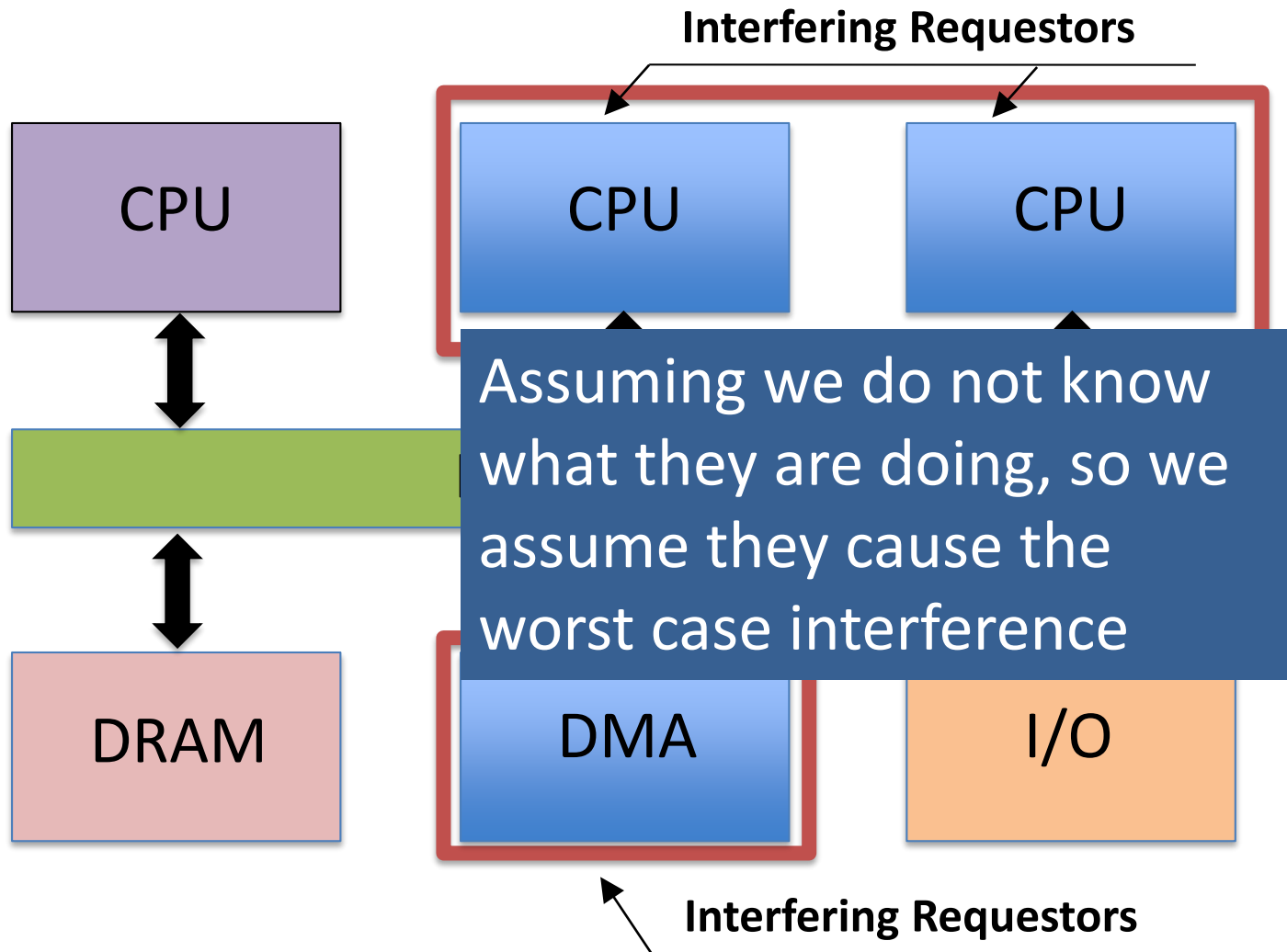
- Scheduling
- WCET d
- WCET: n
cache, D
- Existing approaches can bound the interference but they **assume the latency for DRAM access is constant**

Problem:
DRAM latency is variable and
changes depending on its state source (e.g.

Contribution



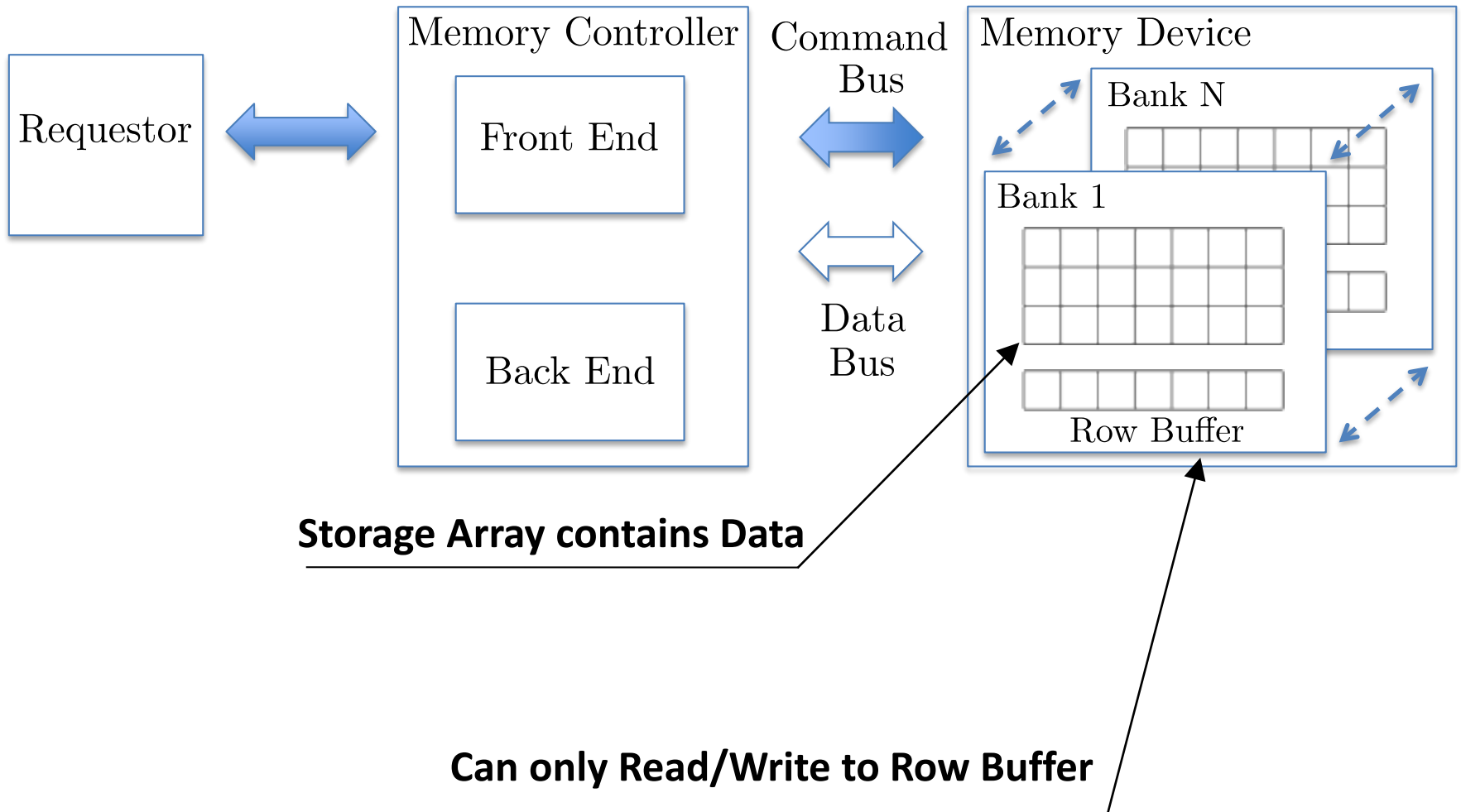
Contribution



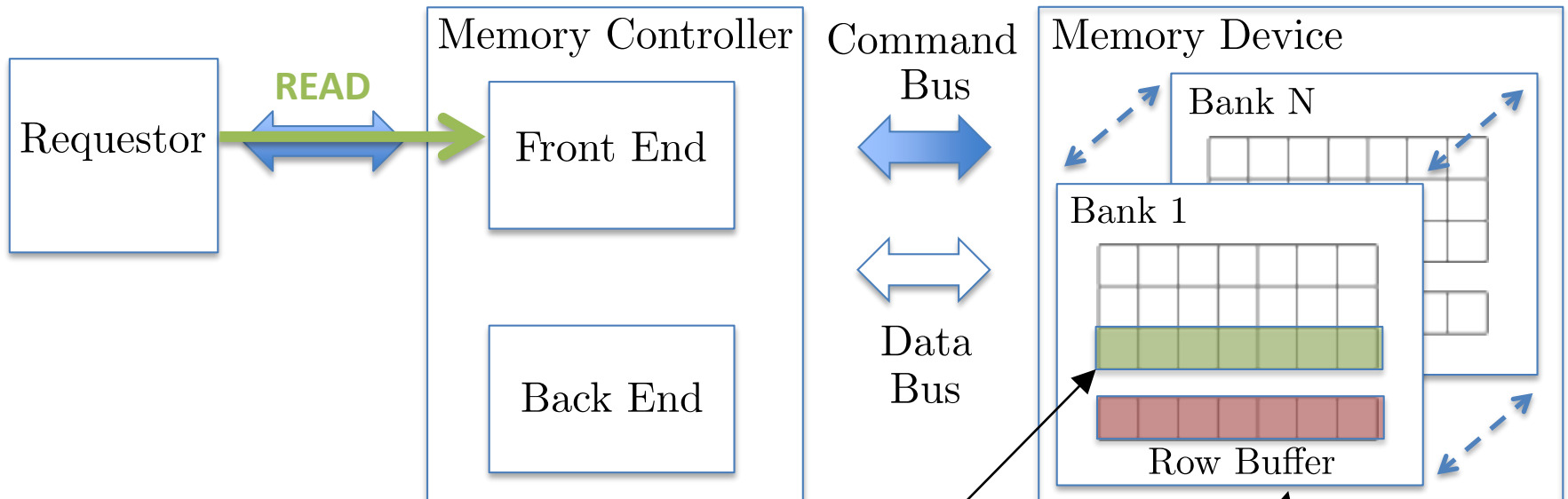
Outline

1. Background & Related Work
2. Memory Controller Model
3. Worst Case Latency Analysis
4. Results & Conclusion

Background



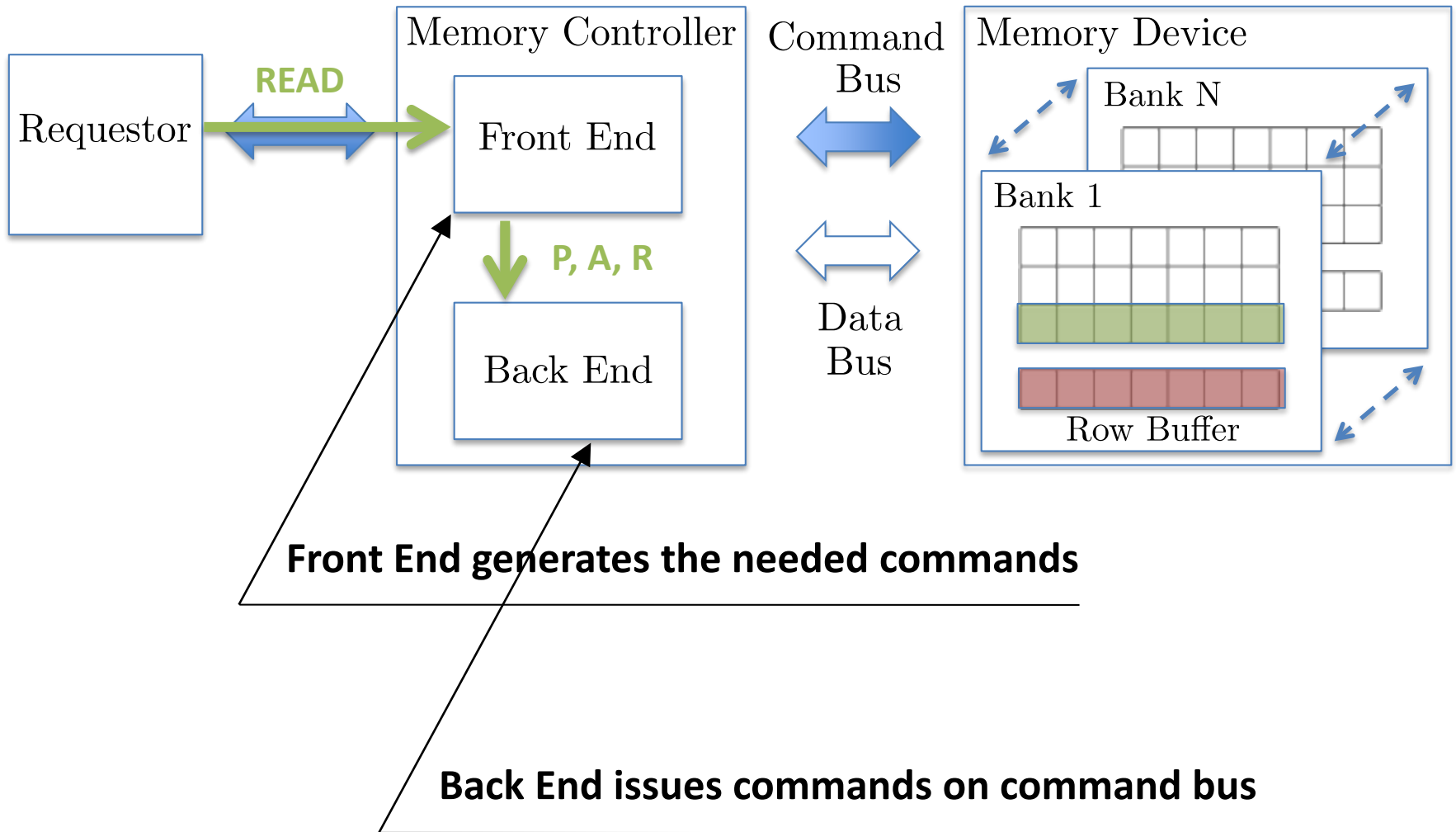
Background



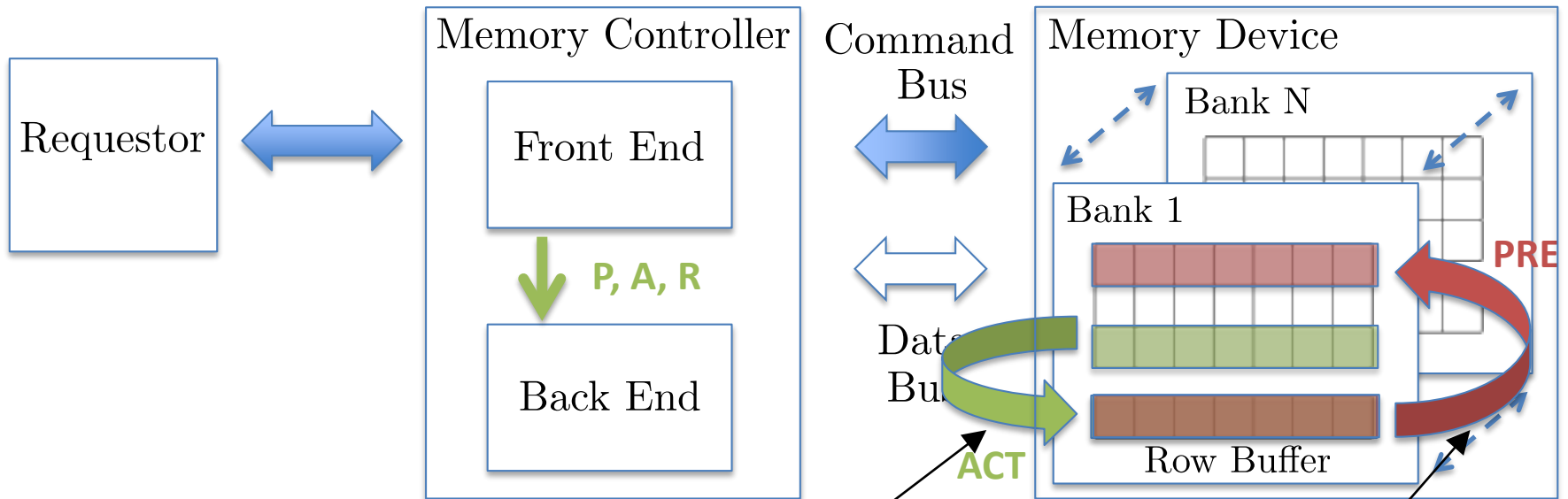
Targeting Data in this Row

Row Buffer contain data from a different row

Background



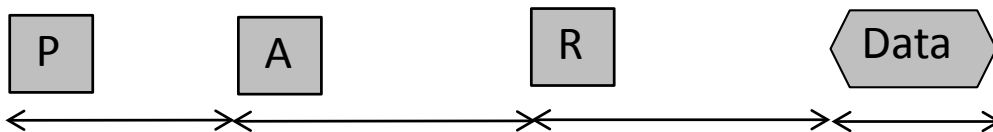
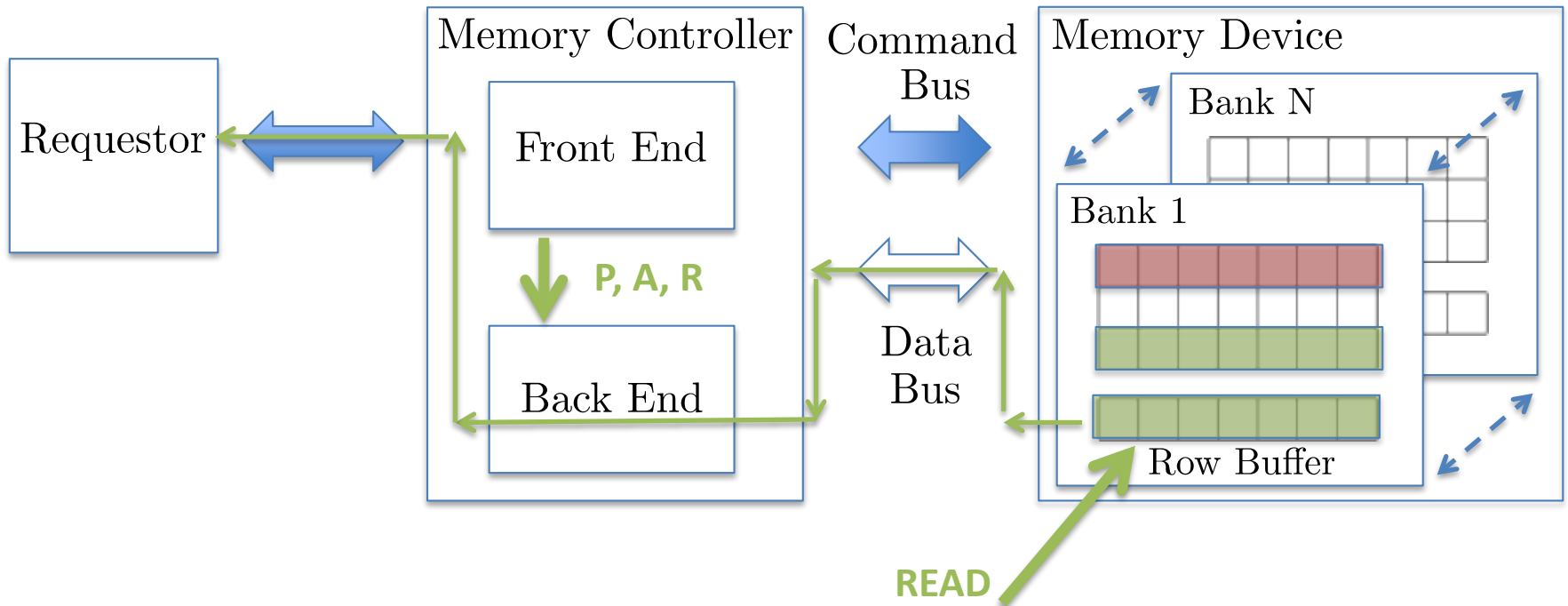
Background



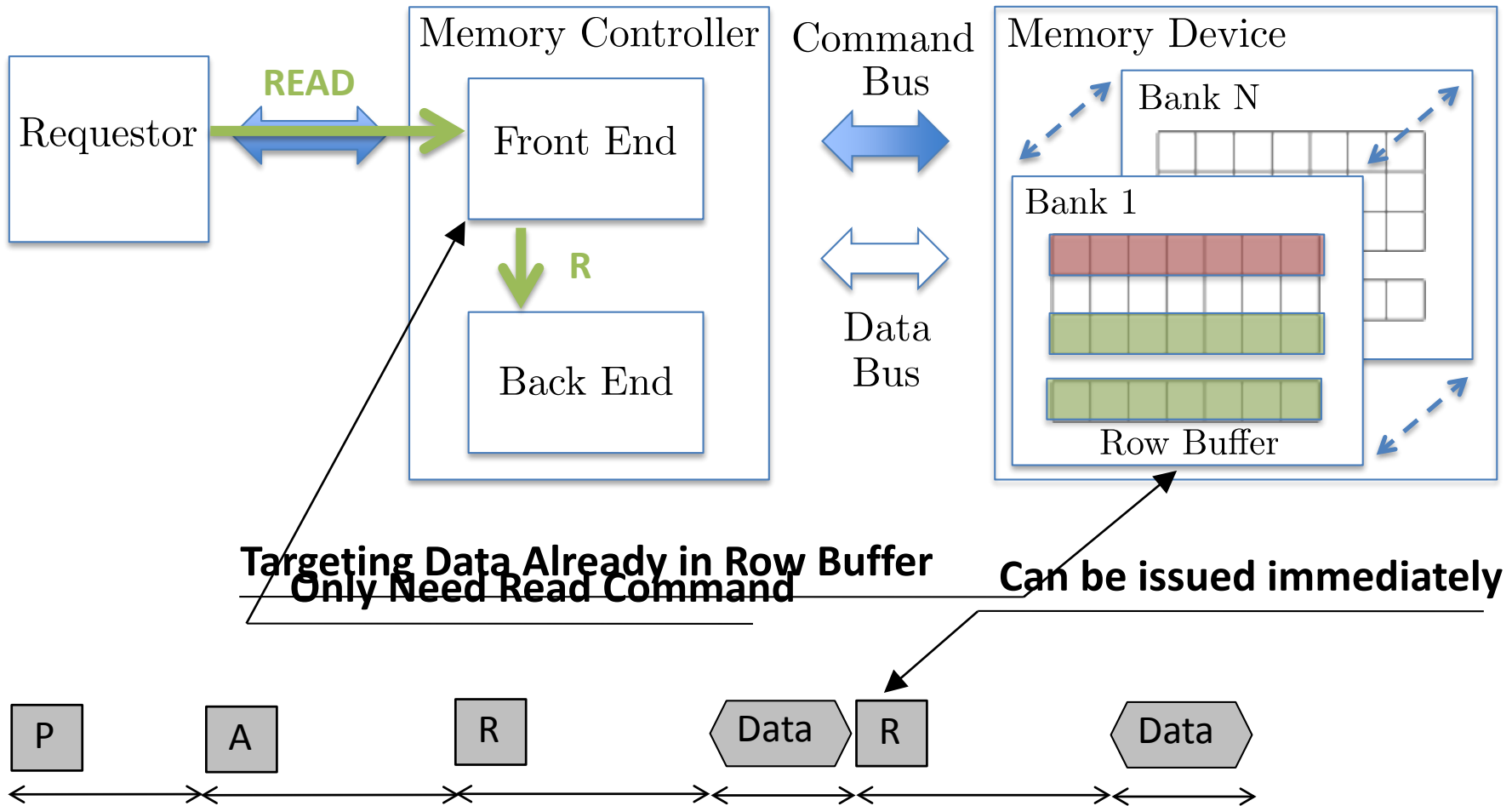
ACT: Load the data from array into buffer
PRE: Charge buffer, store the data back into array



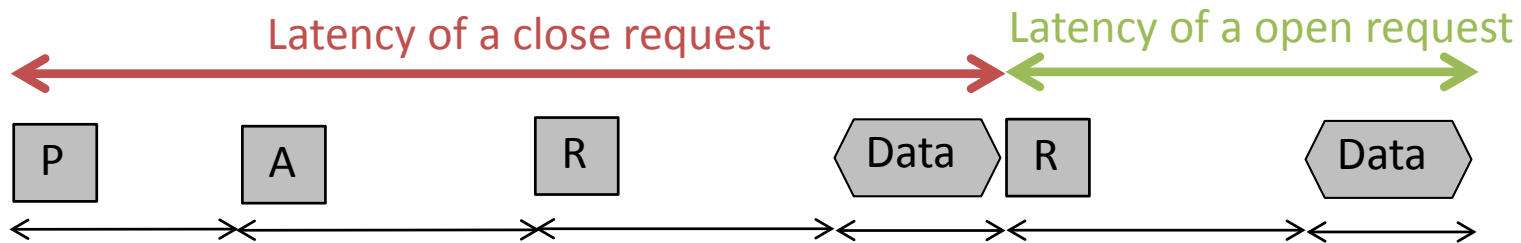
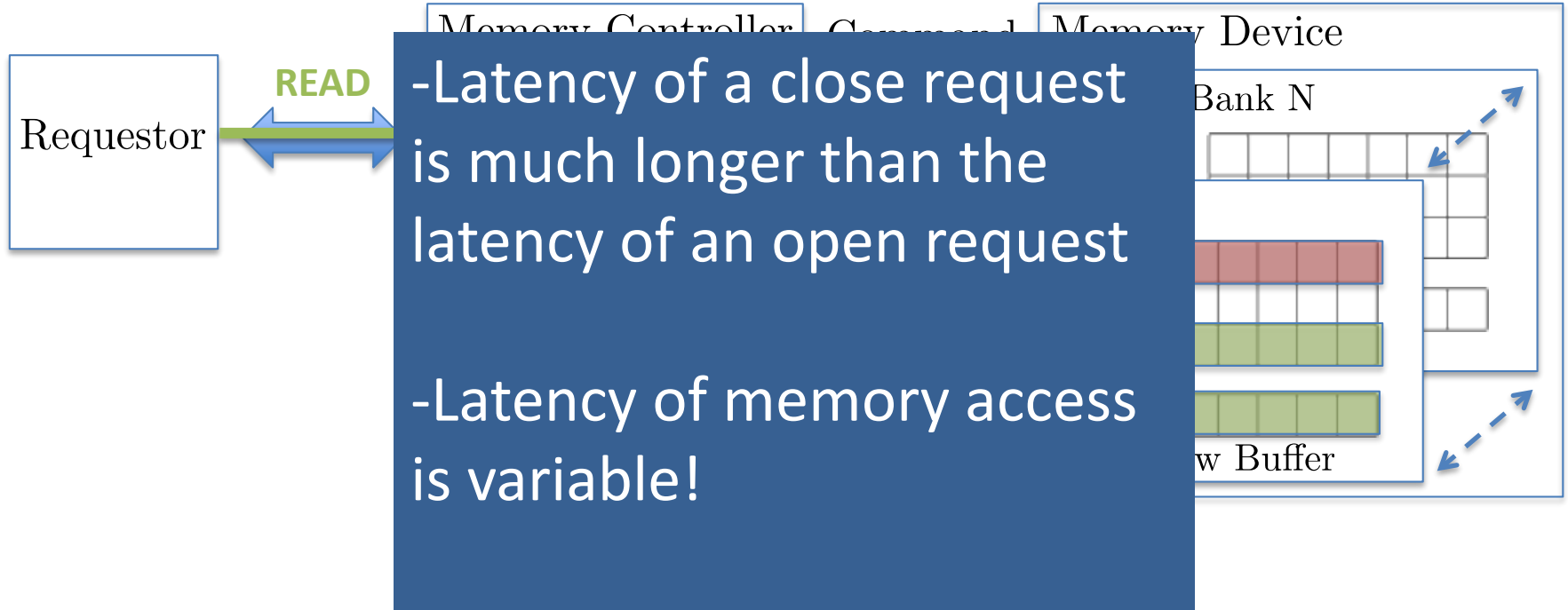
Background



Background



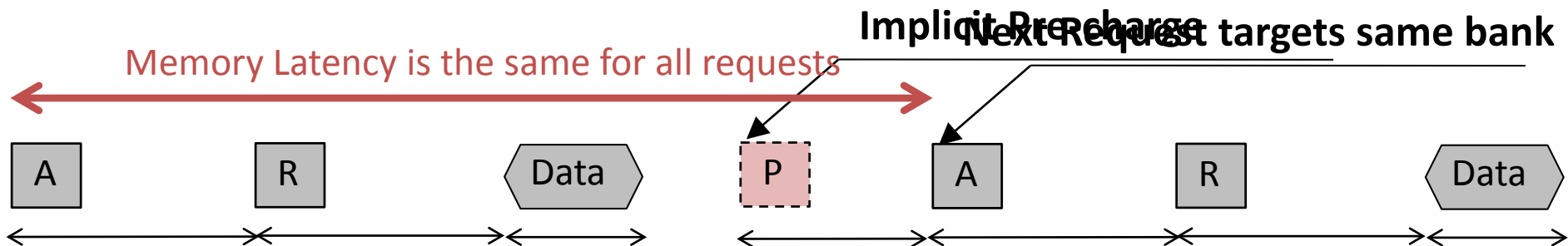
Background



Predictable Memory Controllers

- Close Row Policy:

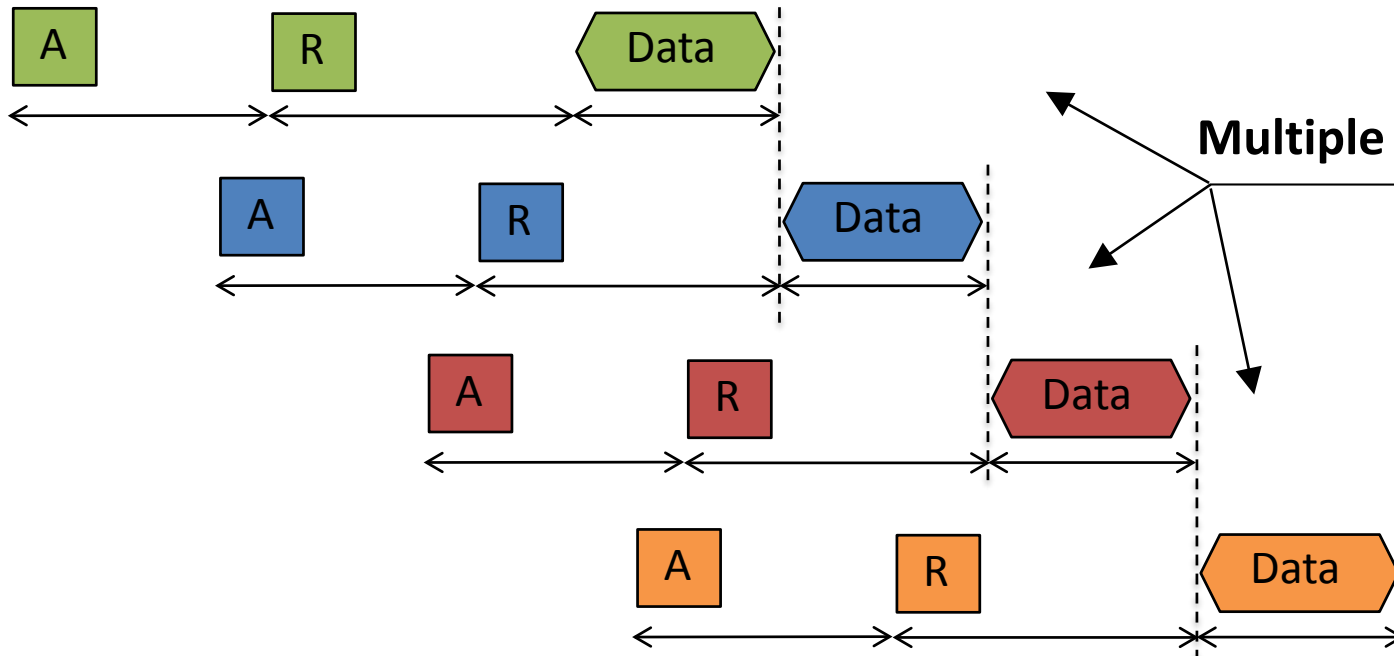
- After each request, the controller must explicitly re-charge the row. -Can not take advantage of automatic locality (row hits)
- Latency is much longer than open request



Predictable Memory Controllers

- Interleaving Banks

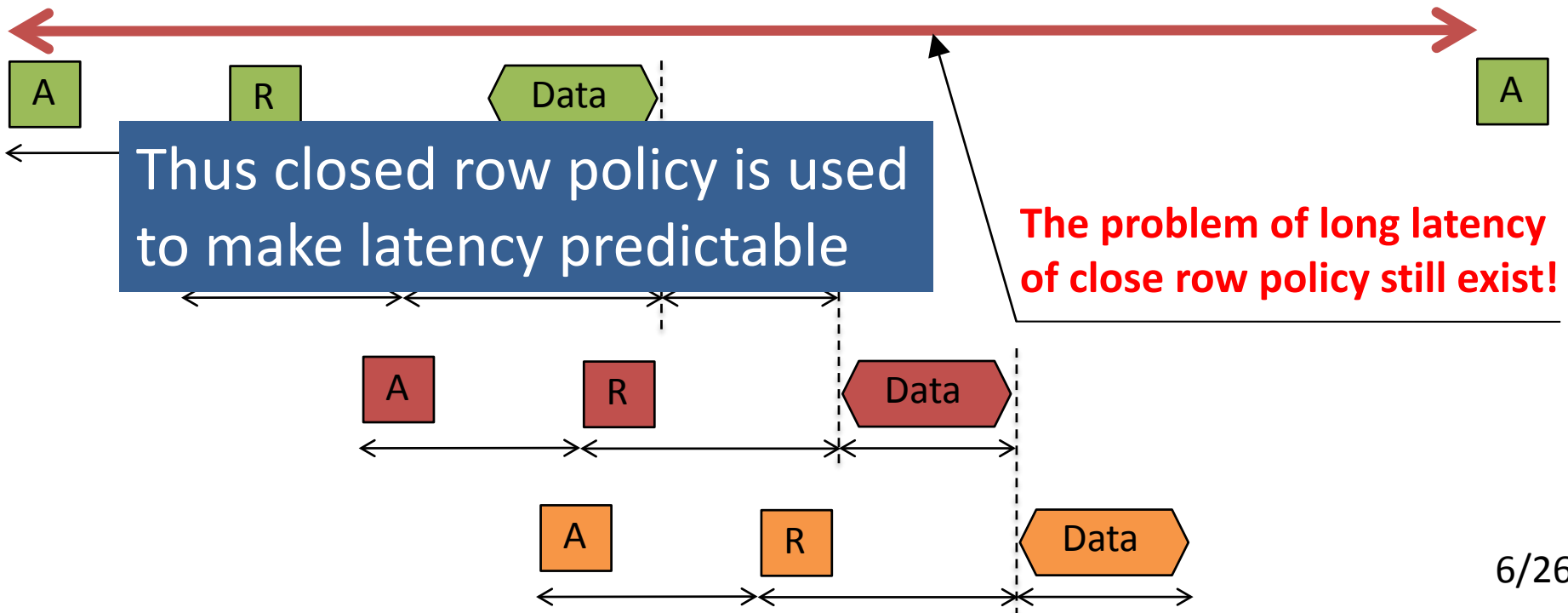
Accessing data in multiple banks



Predictable Memory Controllers

- Interleaving Banks

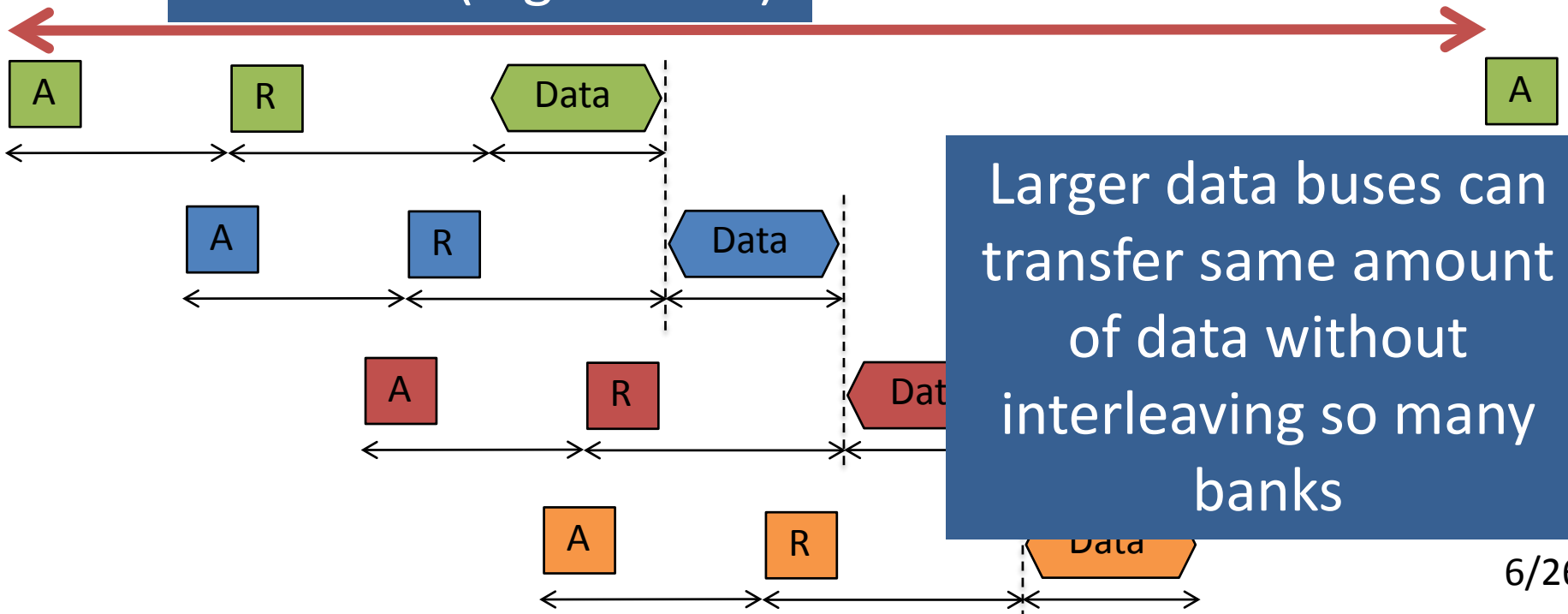
Problem: requestors can close each other's row buffer since they can access all banks



Predictable Memory Controllers

- Interleaving Banks

This is good for system with small DRAM data bus width (e.g. 16 bits)



Predictable Memory Controllers

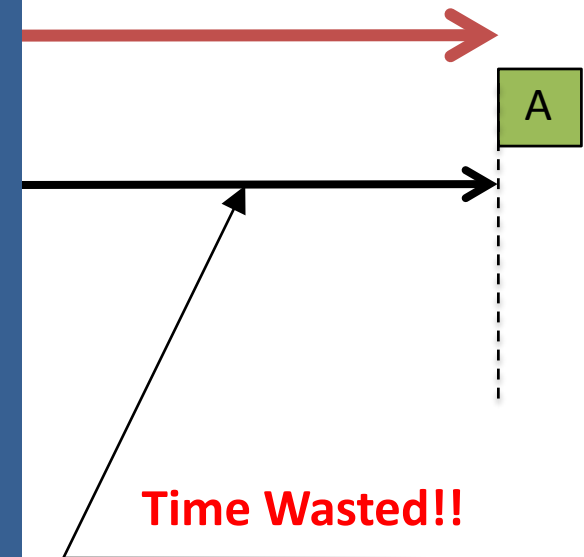
- Interleaving Banks

Interleaving two banks for wider data bus (e.g. 32 bits)



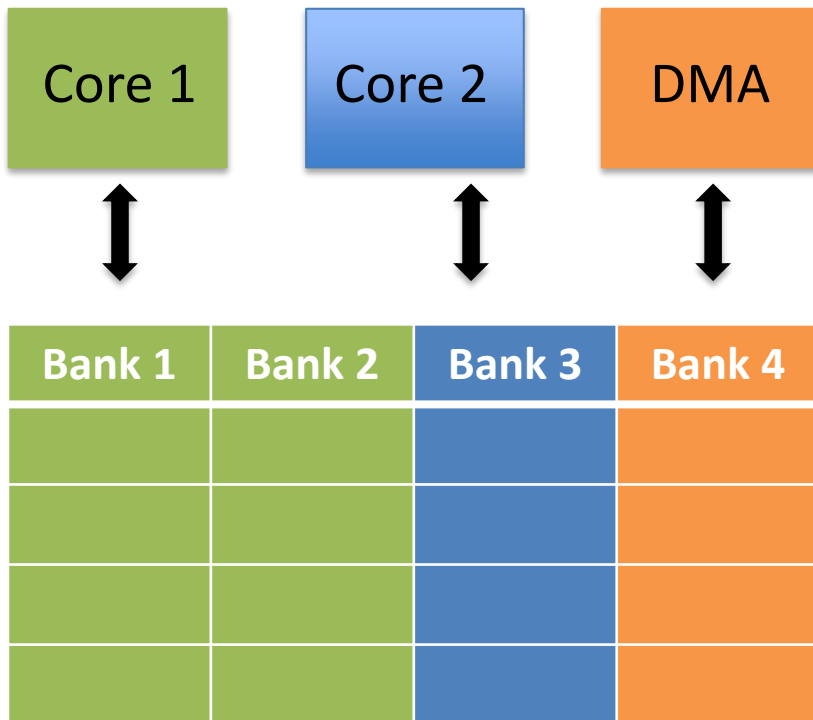
Interleaving Problems:

1. Requestors can close each other's rows (interference)
2. Must be used with close row policy to make latency predictable
3. For wider data bus, effectiveness of interleaving is diminished



Predictable Memory Controllers

- Private Banks



- Can partition banks to either requestors or tasks
- This can be done by:
 - Hardware if Memory controller supports
 - By compiler
 - In OS, using virtual memory

Related Work

- AMC[1] and Predator [2]:
 - Close Row Policy
 - Interleaved Bank
- Conservative Open-Page [3]:
 - Interleaved Bank
 - Leave row open for a small window of time
- PRET DRAM Controller [4]:
 - Close Row Policy
 - Private Bank

Our Approach

- Private Bank

- eliminate
requirements

Challenge:

1. Analysis is more complex
2. More than 20 timing constraints
3. Latency depends on the dynamic state of DRAM

other

- Open Row Policy

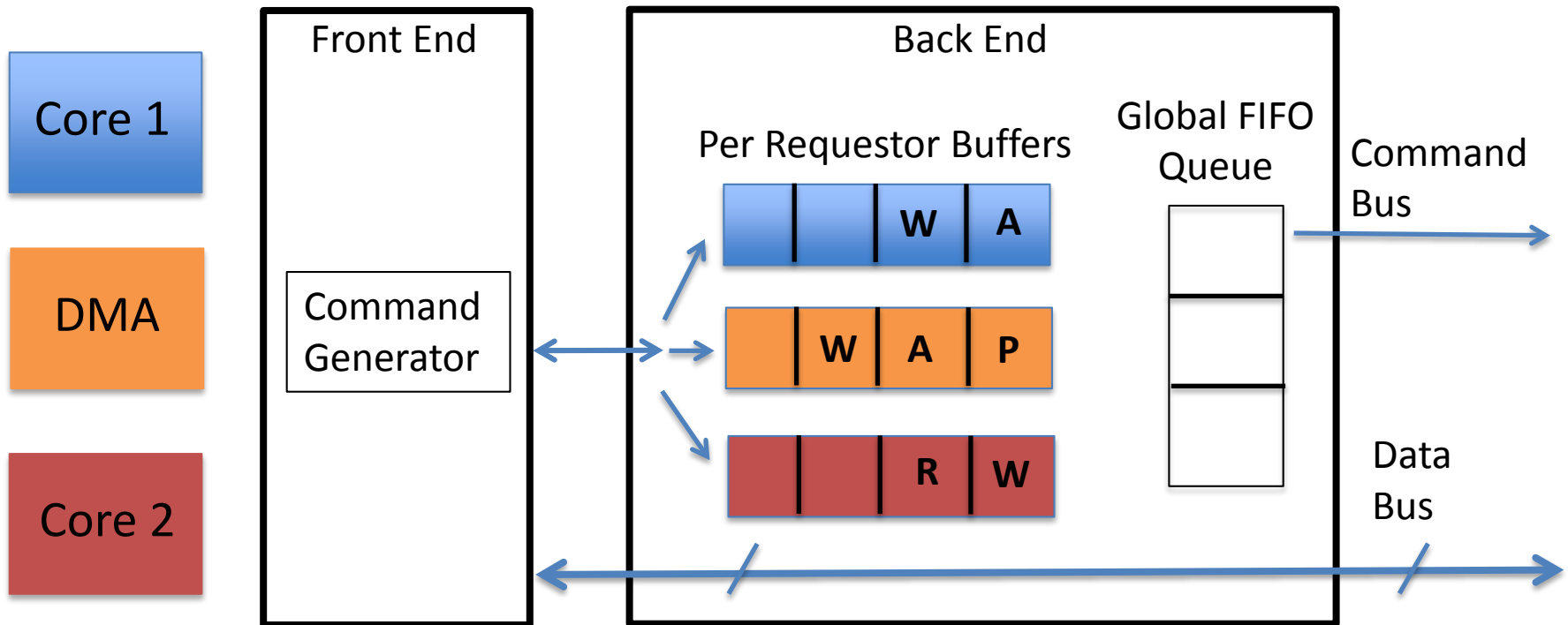
- reduce latency and take advantage of row hit ratio (locality)

Outline

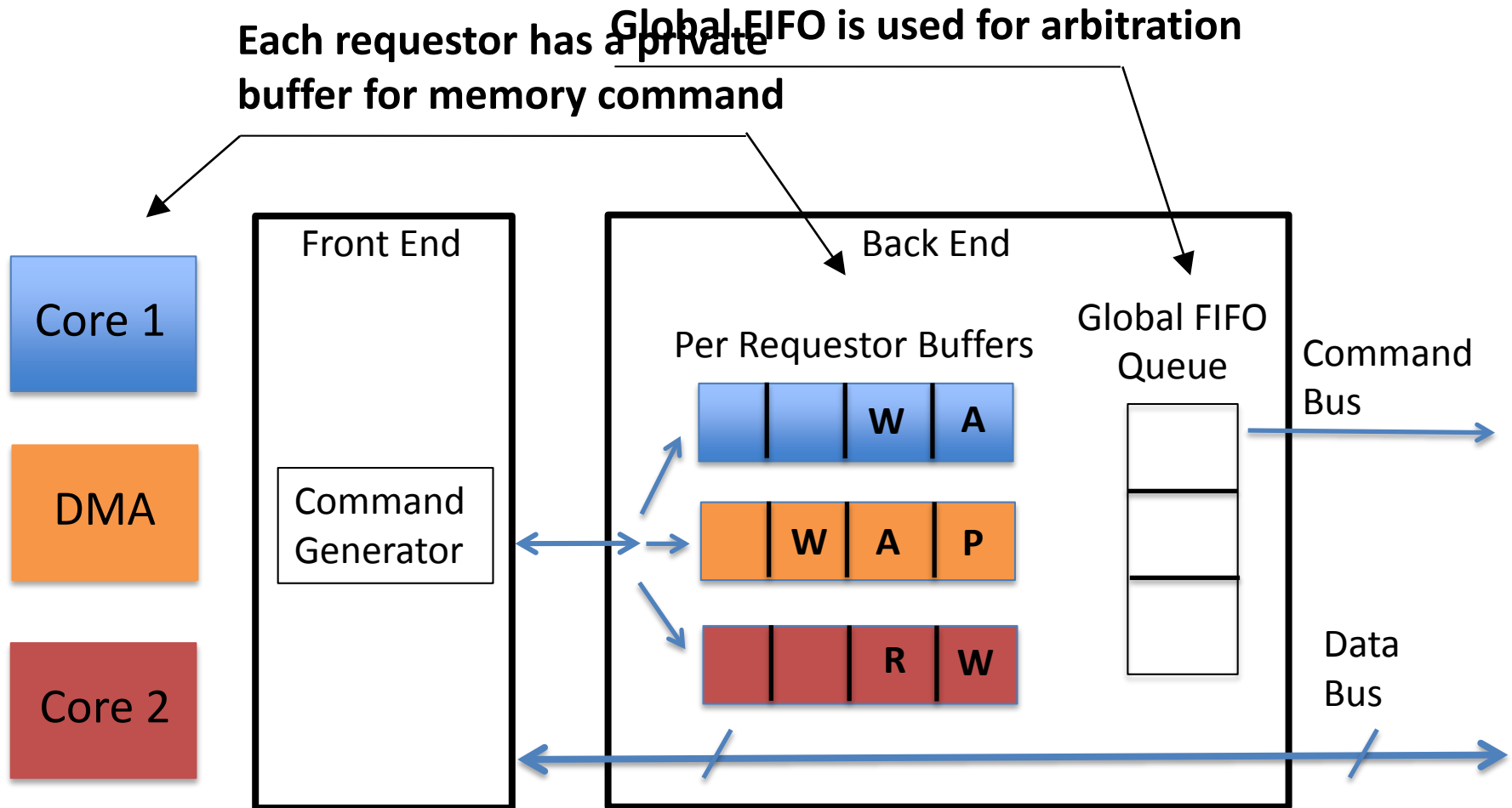
1. Background & Related Work
2. Memory Controller Model
3. Worst Case Latency Analysis
4. Results & Conclusion

Memory Controller Model

We focus on the back end latency
ignore **CONSTANT** front end delay

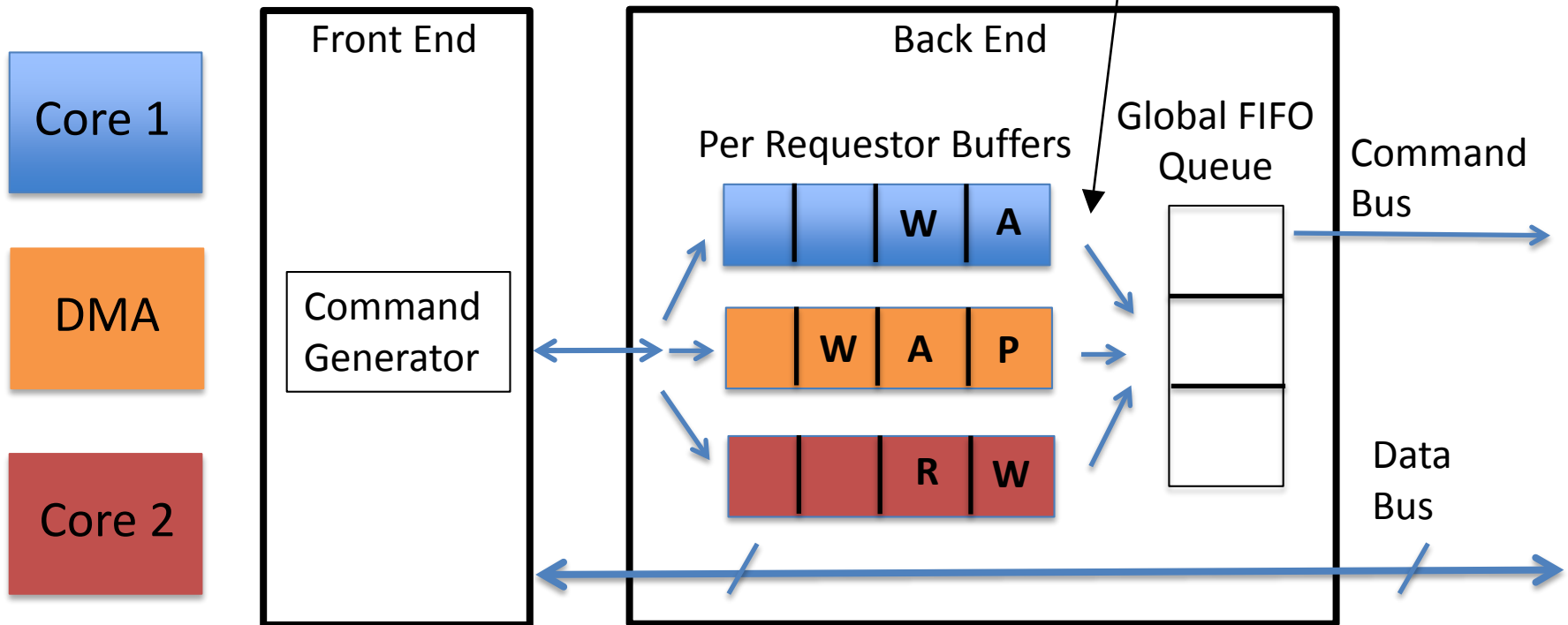


Memory Controller Model



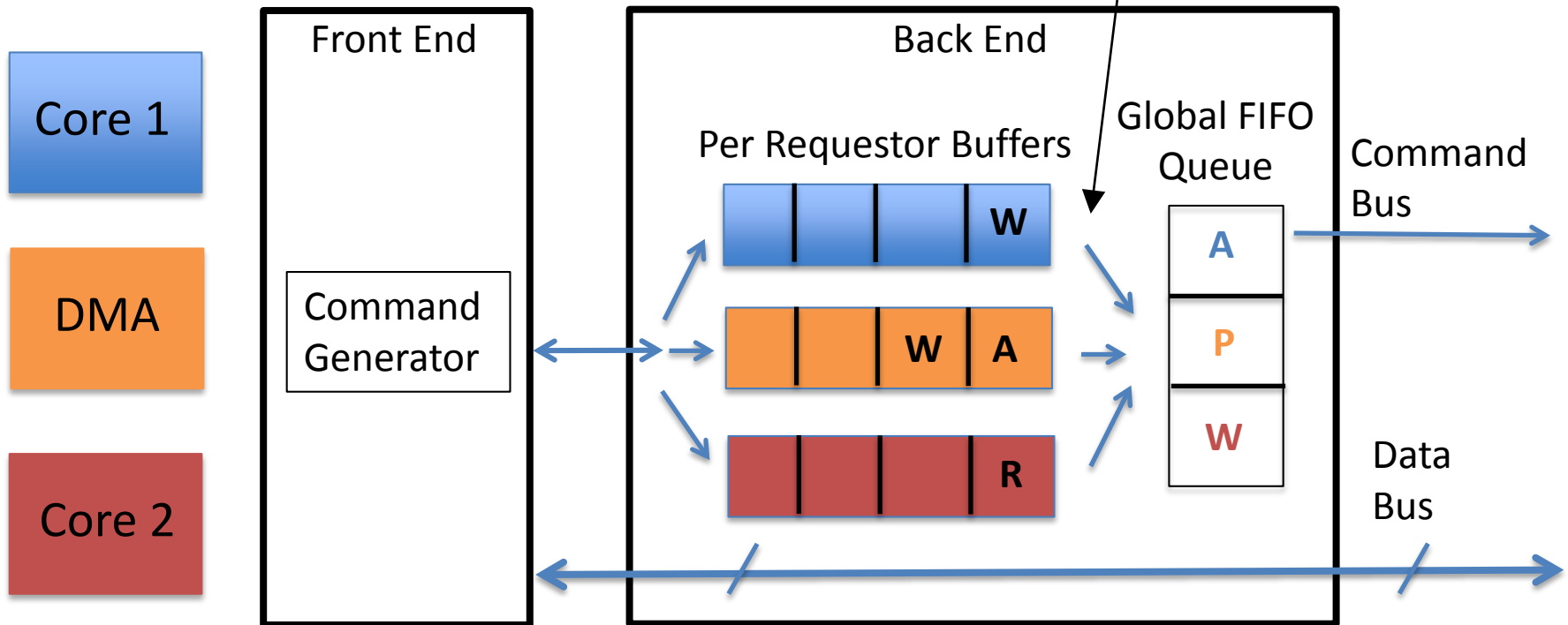
Memory Controller Model

Command at head of each private buffer are inserted into the FIFO



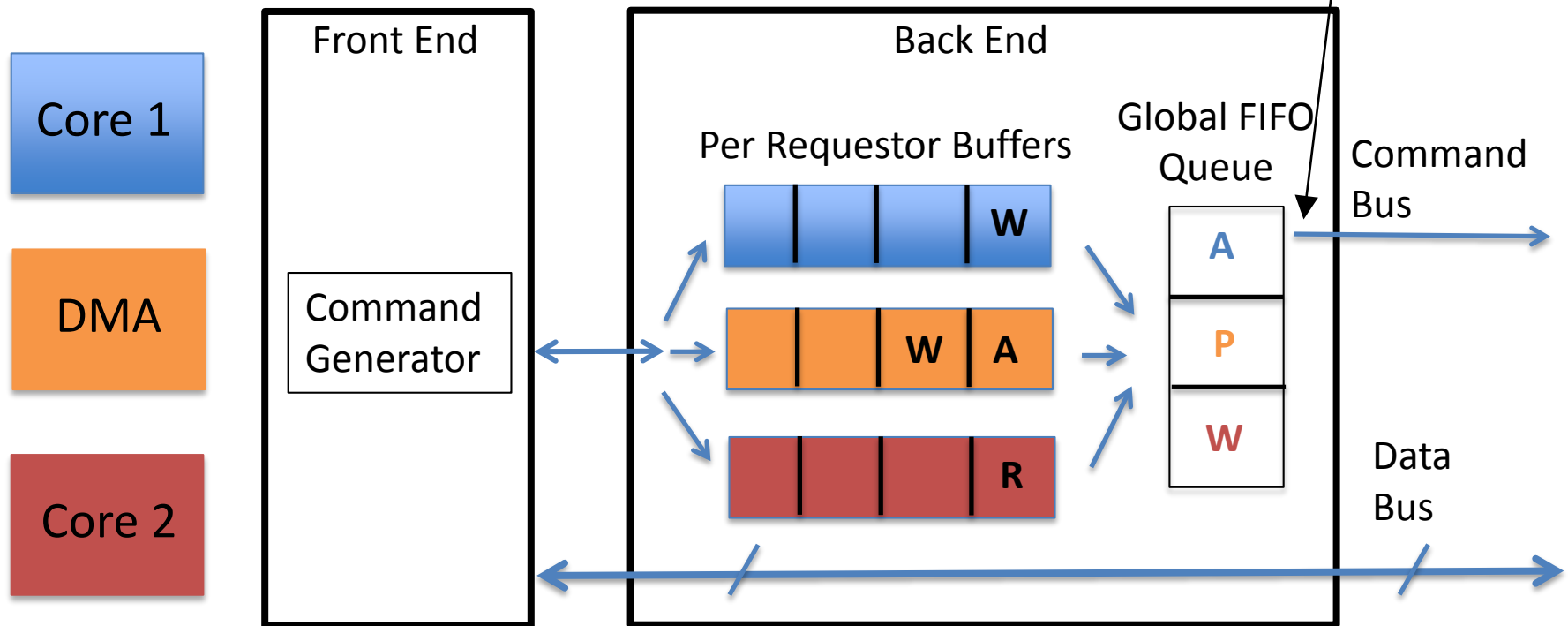
Memory Controller Model

Command at head of each private buffer are inserted into the FIFO



Memory Controller Model

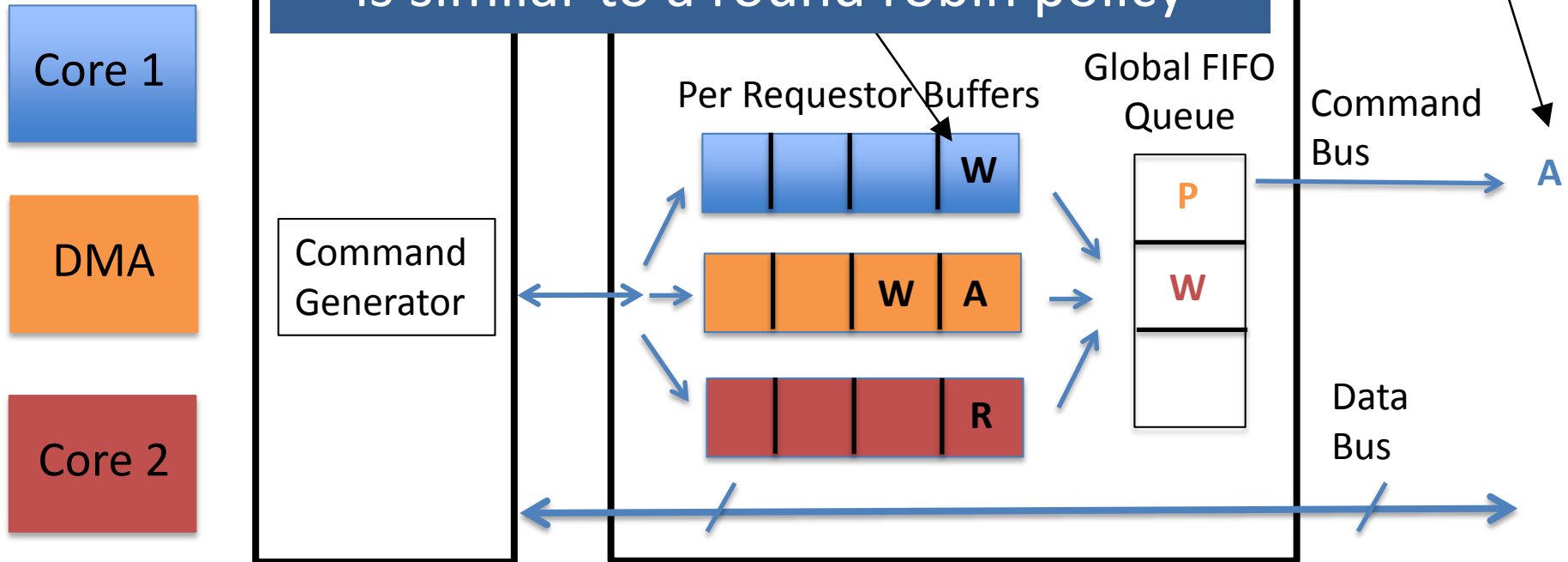
Controller scan the global FIFO from front to end for a command that can be issued



Memory Controller Model

Next command must wait until timing constraints are satisfied before it can be inserted into FIFO

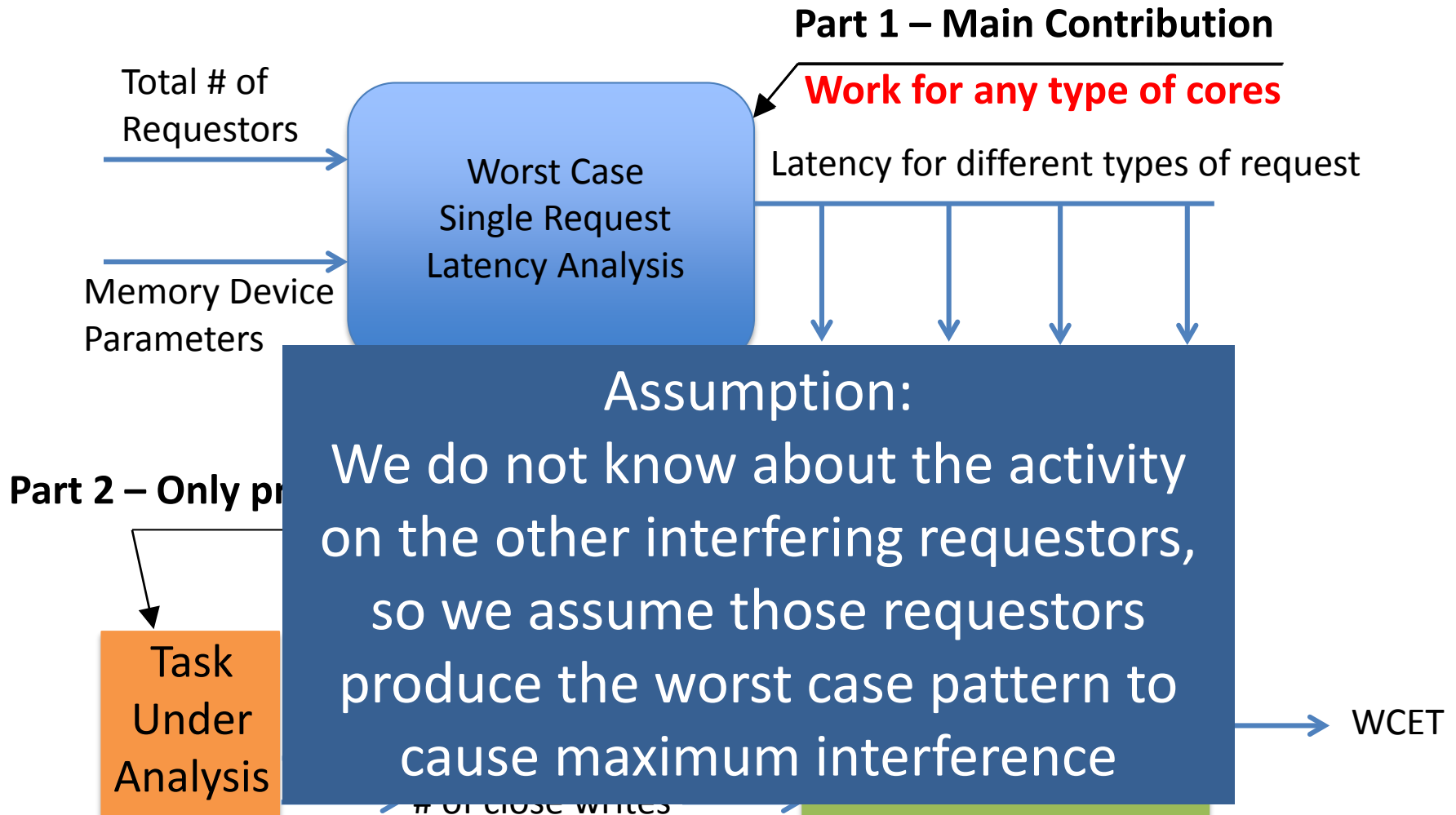
Intuitively, the arbitration is fair and is similar to a round robin policy



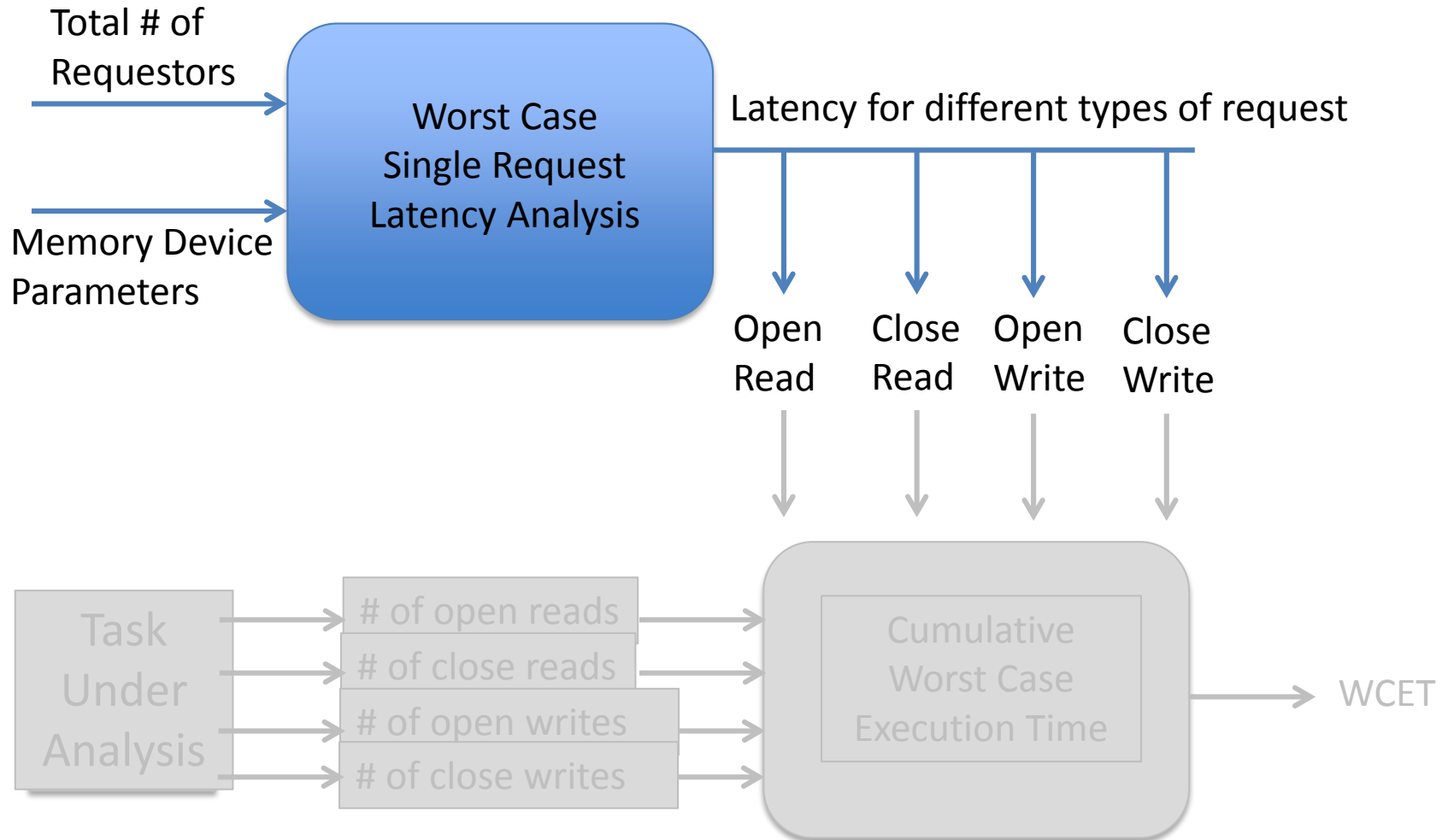
Outline

1. Background & Related Work
2. Memory Controller Model
- 3. Worst Case Latency Analysis**
4. Results & Conclusion

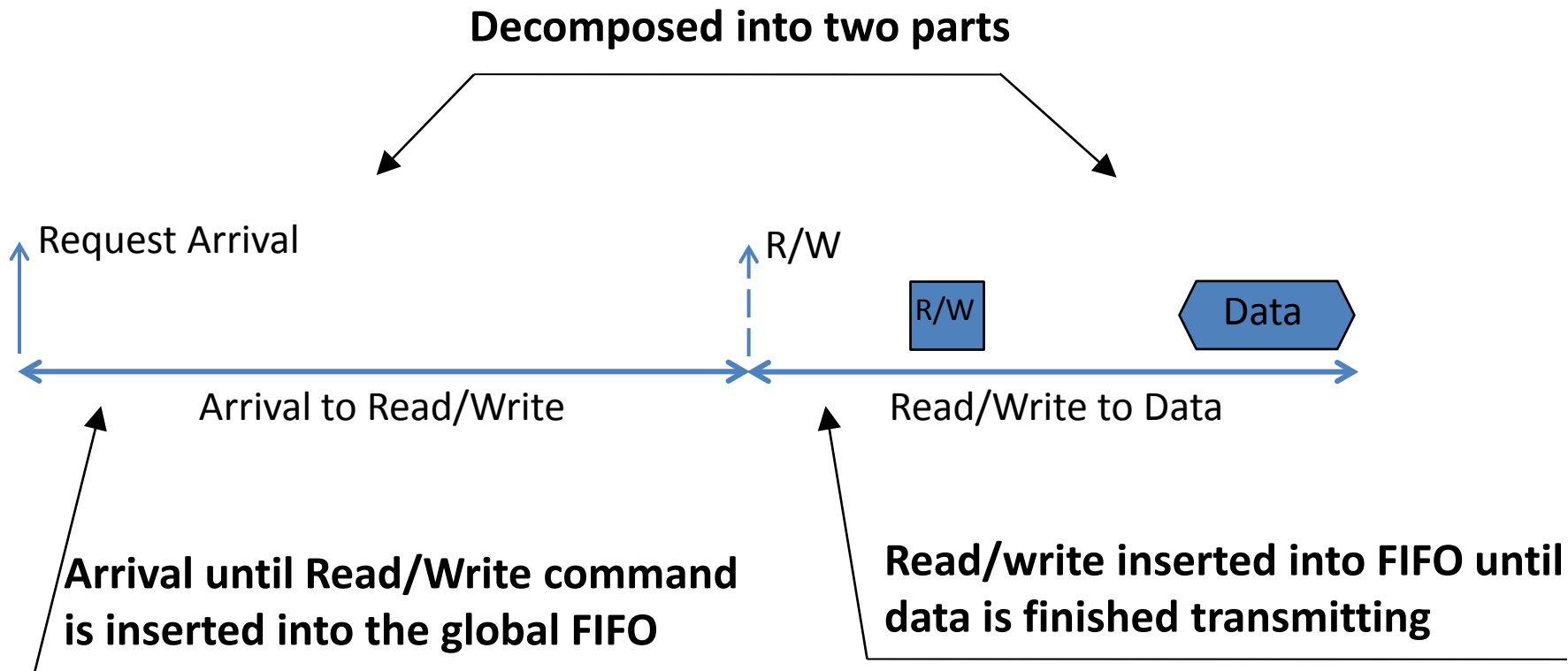
Worst Case Analysis



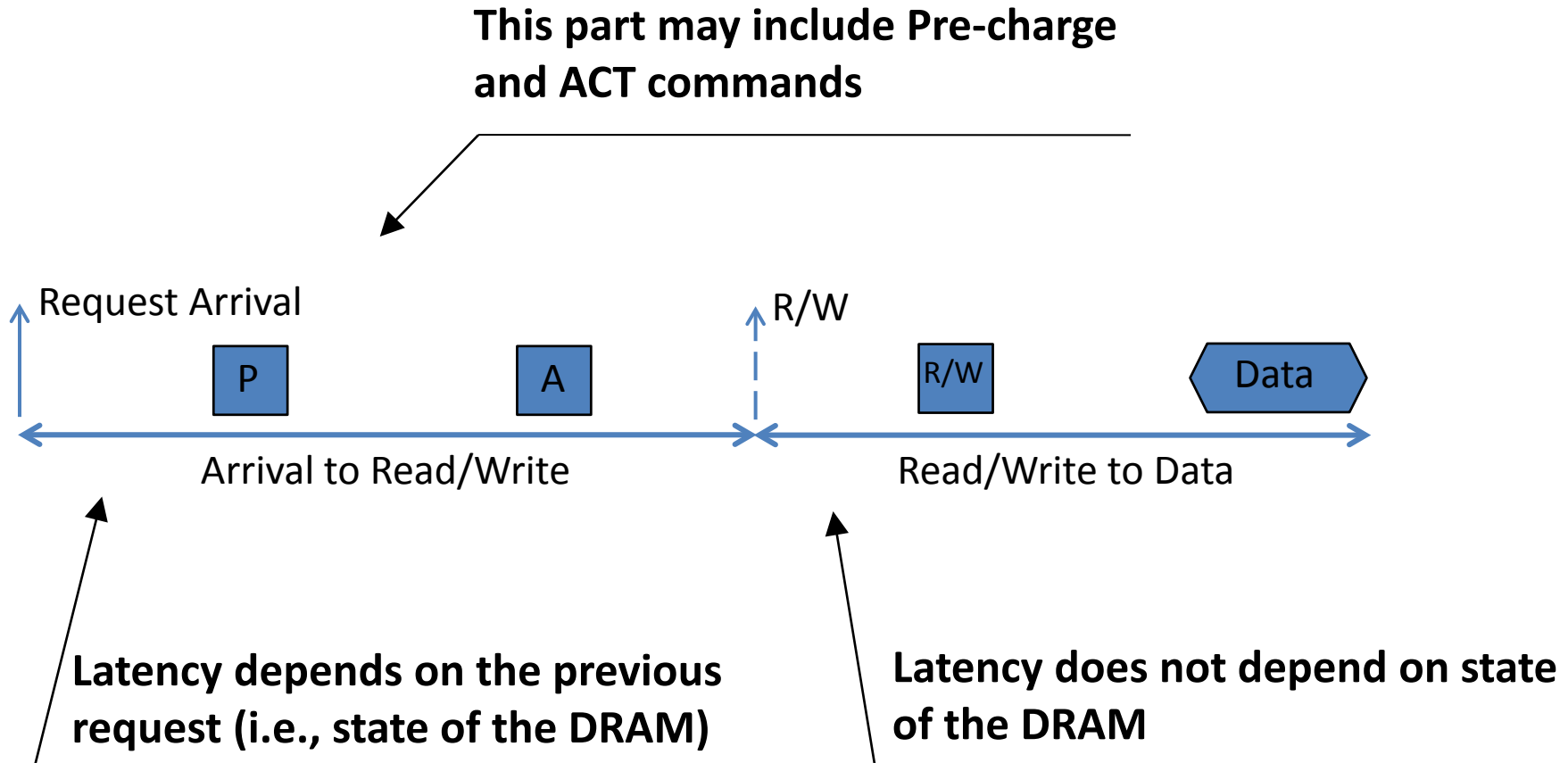
Worst Case Analysis



Single Request Latency

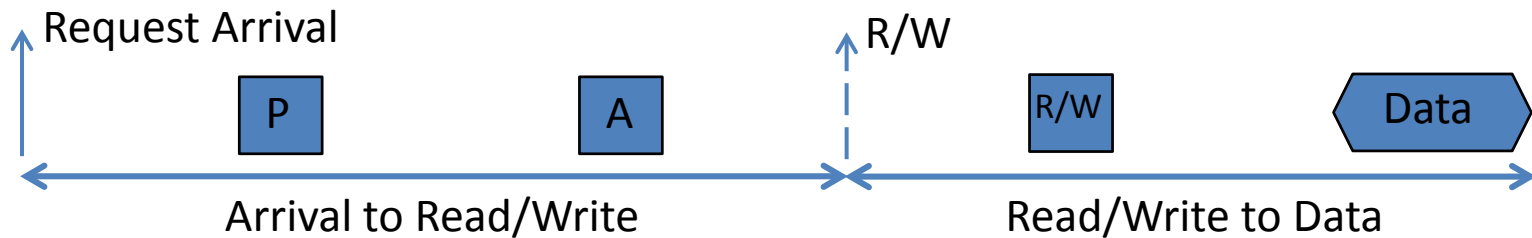


Single Request Latency

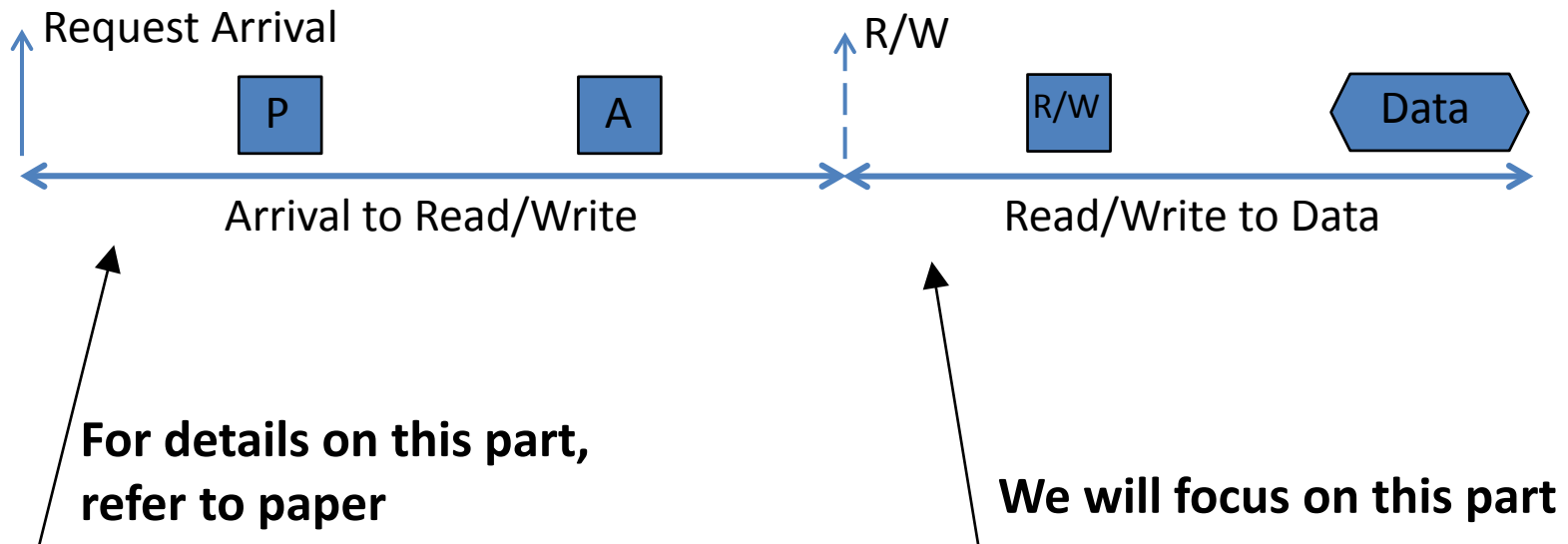


Single Request Latency

Both parts depends on the # of interfering requestors as well as DRAM timing constraints

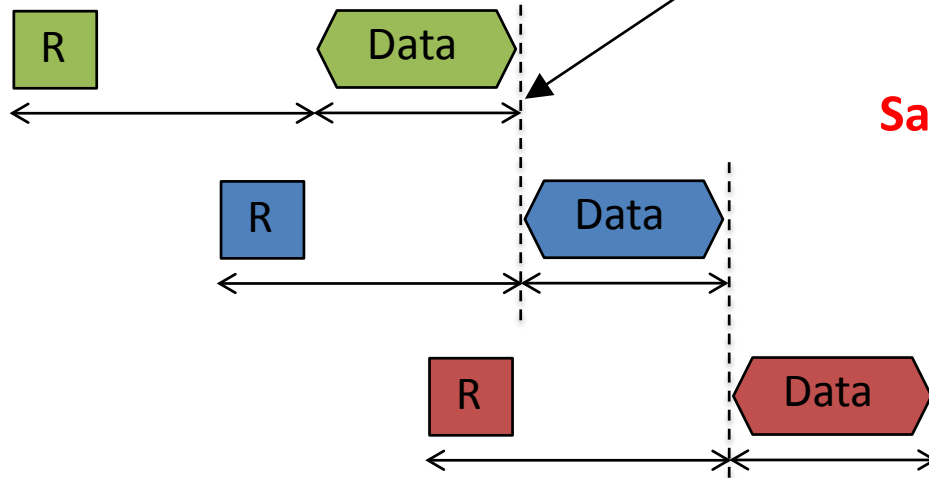


Single Request Latency



Read/Write to Data Latency

Read to Read has no timing constraints,
only contention on the data bus

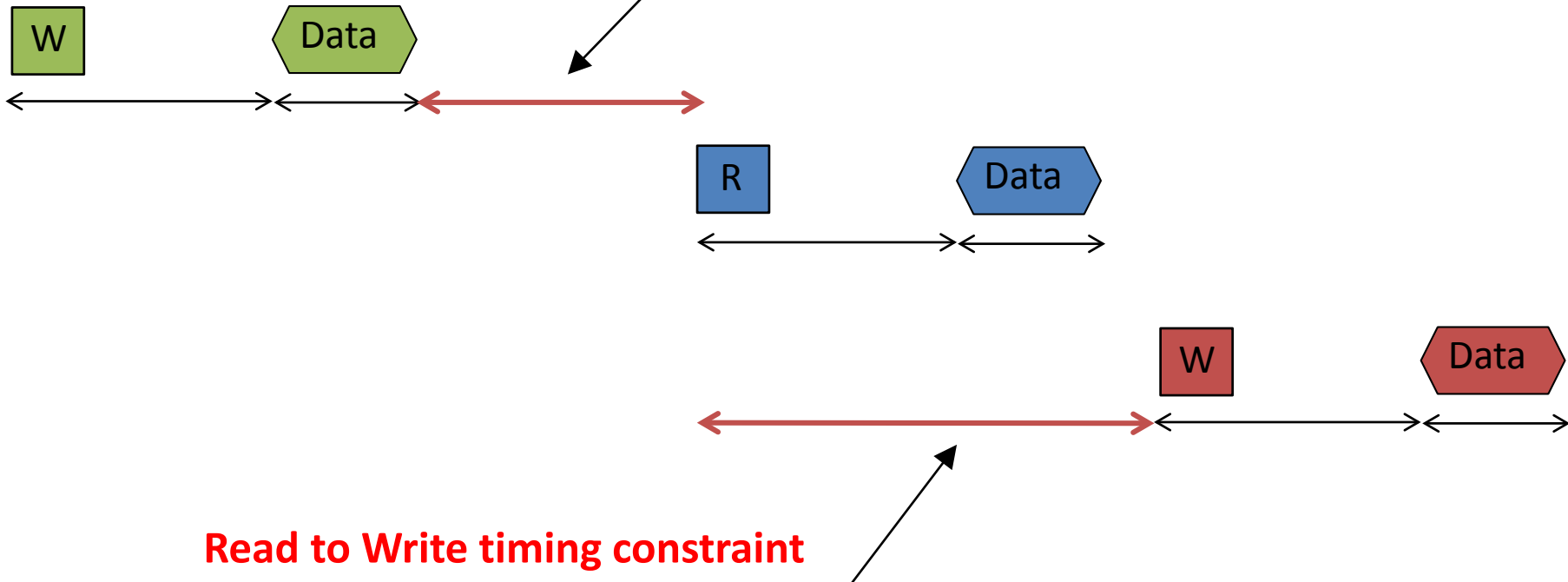


Same for Write to Write

Read/Write to Data Latency

Therefore, an alternation of read and write commands produce longer latency

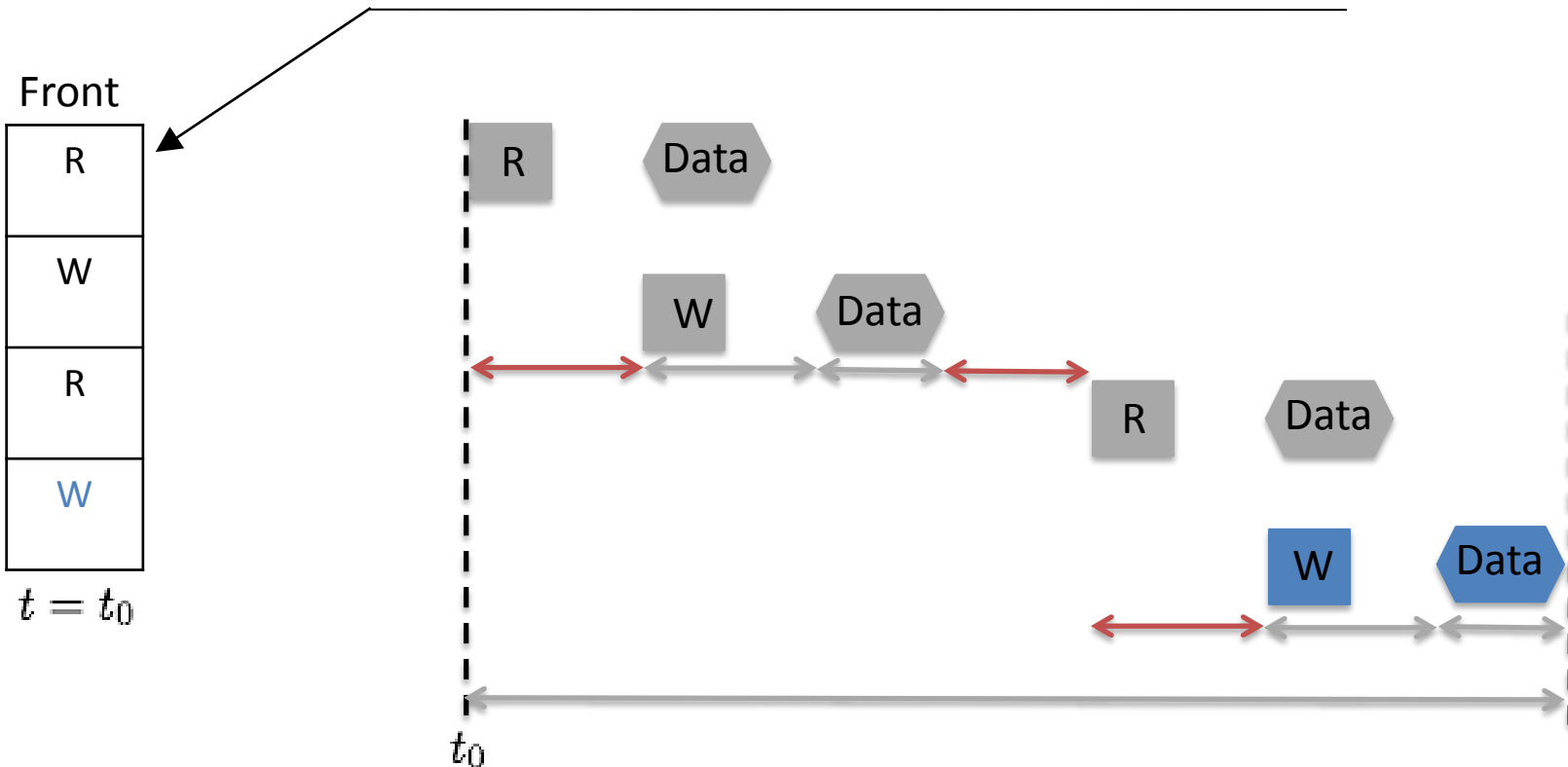
Write to Read timing constraint



Read/Write to Data Latency

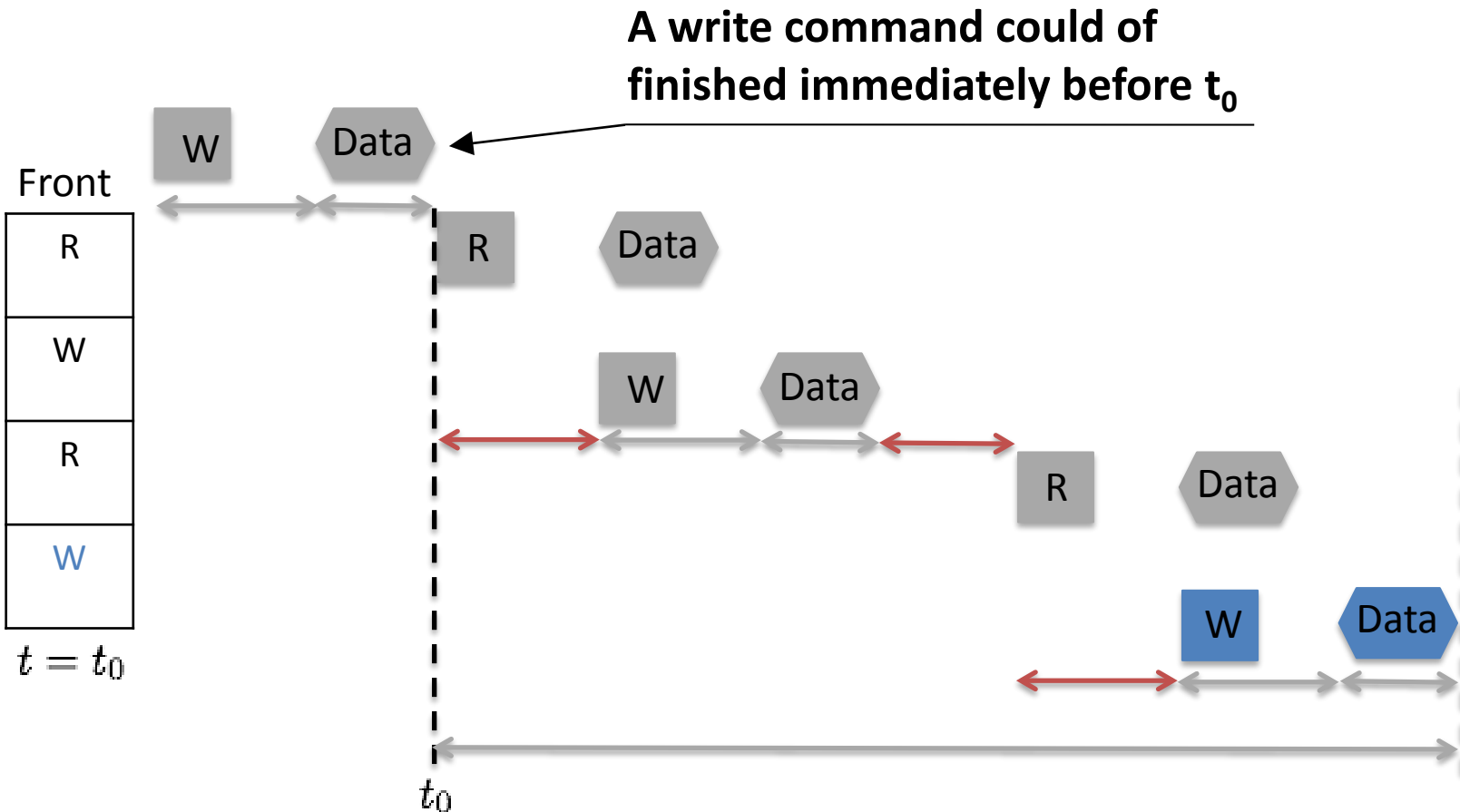
- Interference on Write command

All other requestors inserts R/W commands to create maximum interference



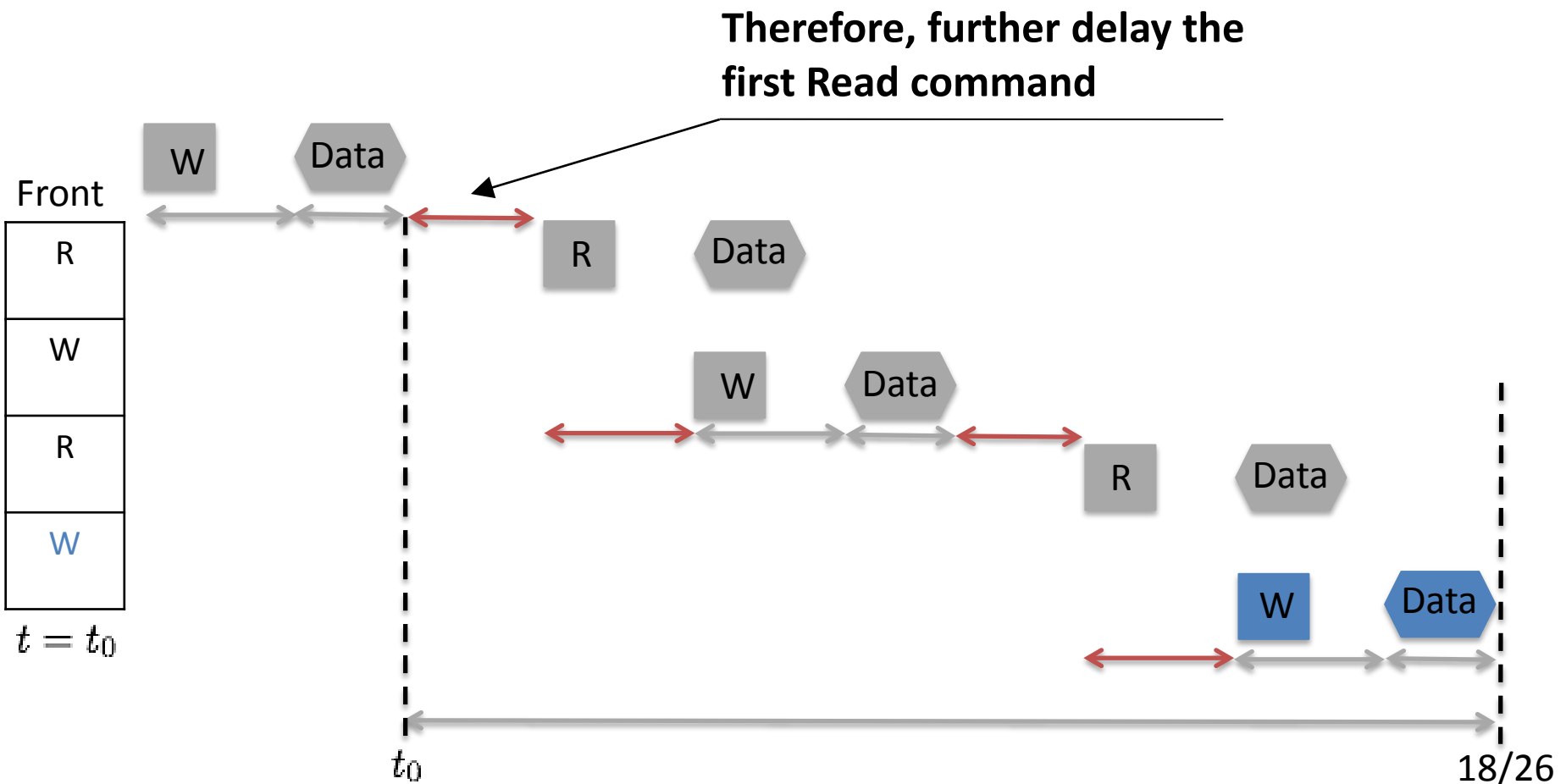
Read/Write to Data Latency

- Interference on Write command

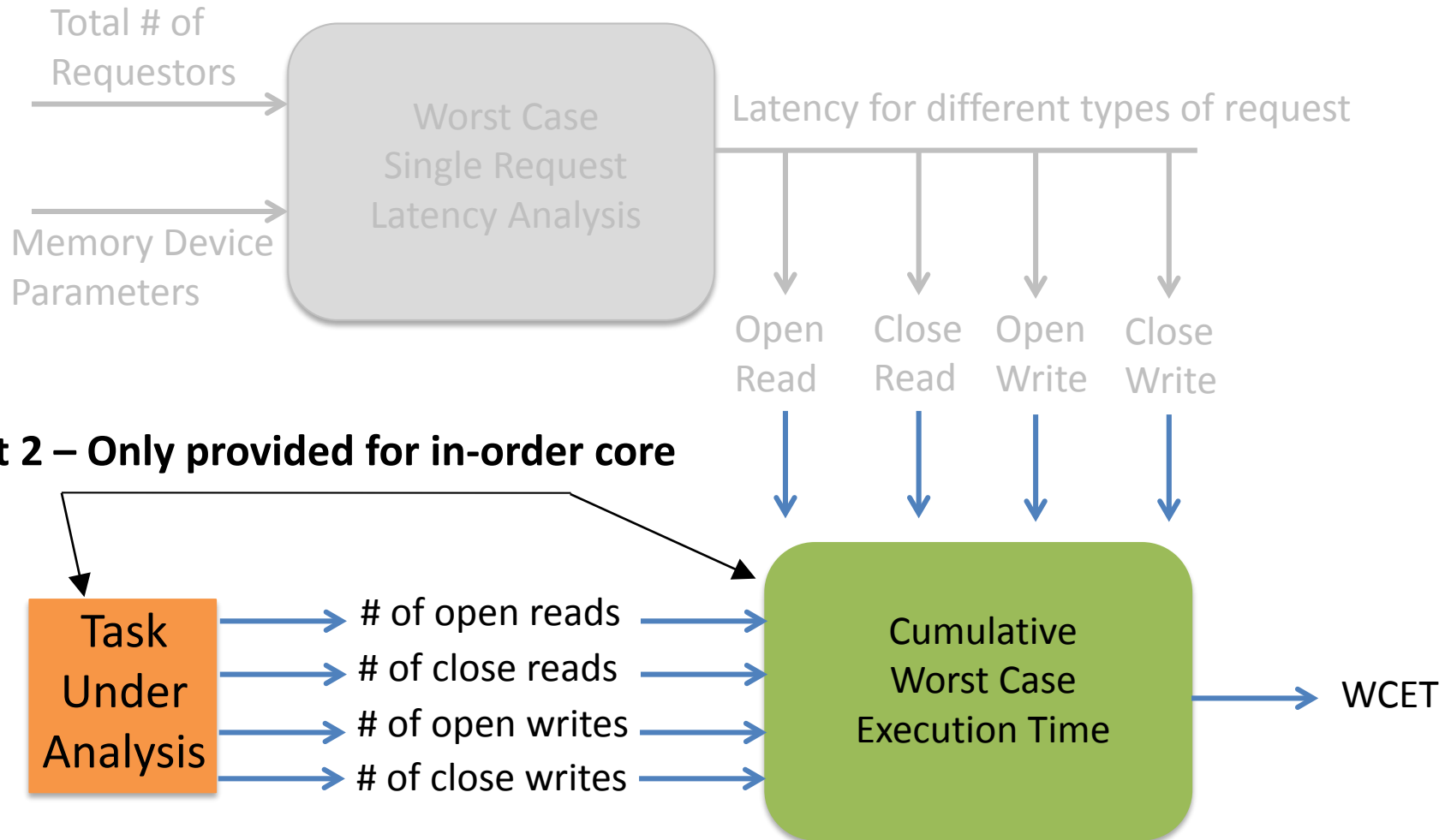


Read/Write to Data Latency

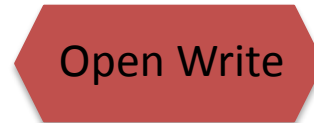
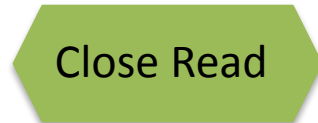
- Interference on Write command



Worst Case Analysis



Cumulative Latency



Task Under Analysis:



Cumulative Latency

Open F

Worst case request order depends on input value, code path, cache state, etc.

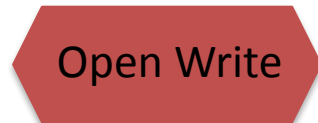
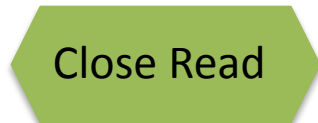
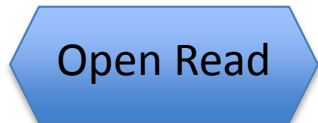
Write

If worst case request order is known, we can sum the latency of each request

Task Under Analysis:



Cumulative Latency



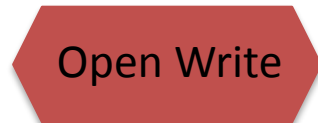
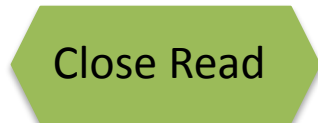
Static Analysis tools can be used to obtain safe bound for # of each type of request

own, we
est

Task Under Analysis:



Cumulative Latency



This problem can be solved in constant time; see paper for detail

Task Under Analysis:



Which pattern leads to worst case latency?



Outline

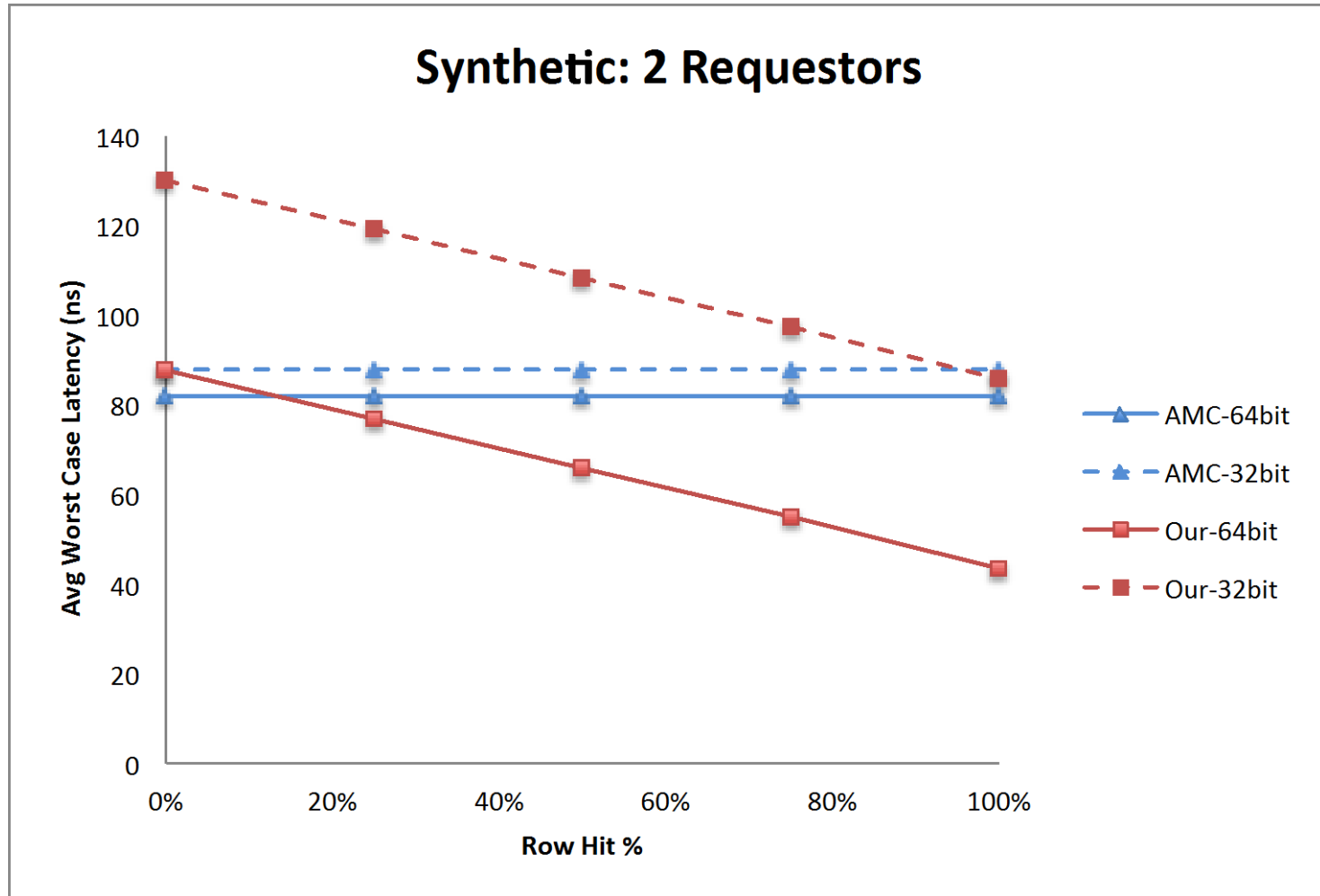
1. Background & Related Work
2. Memory Controller Model
3. Worst Case Latency Analysis
 - Single Request Latency
 - Cumulative Latency
4. Results & Conclusion

Results

- Comparison against Analyzable Memory Controller [1]
 - Since they use fair arbitration (Round Robin) which is similar to our approach
- Synthetic Benchmarks
 - Used to show how worst case latency varies as parameters are changed
- CHStone Benchmarks
 - Memory traces are obtained from gem5 simulator
 - Memory traces are used as input the worst case analysis

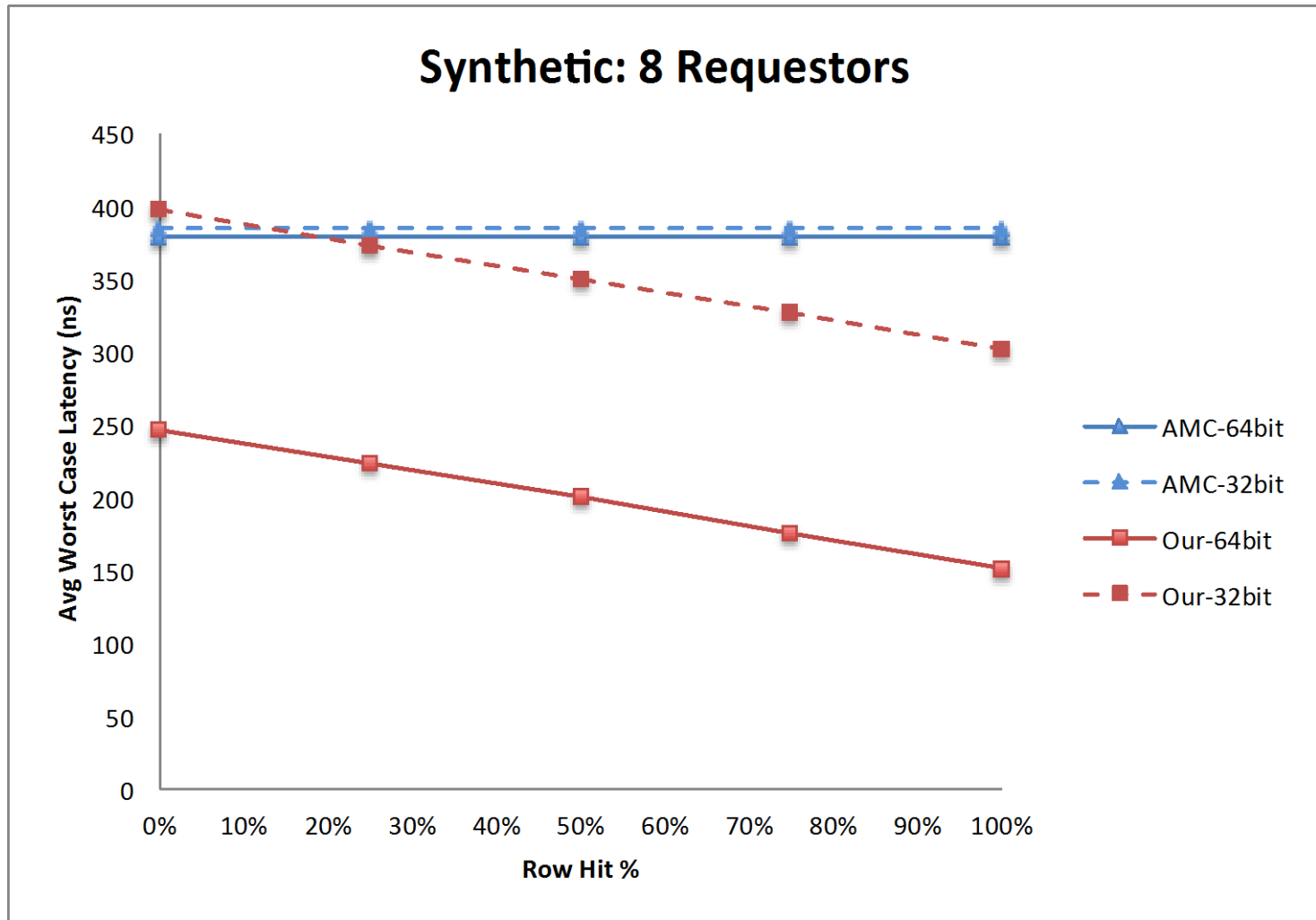
Results

- Synthetic Benchmarks



Results

- Synthetic Benchmarks



Results

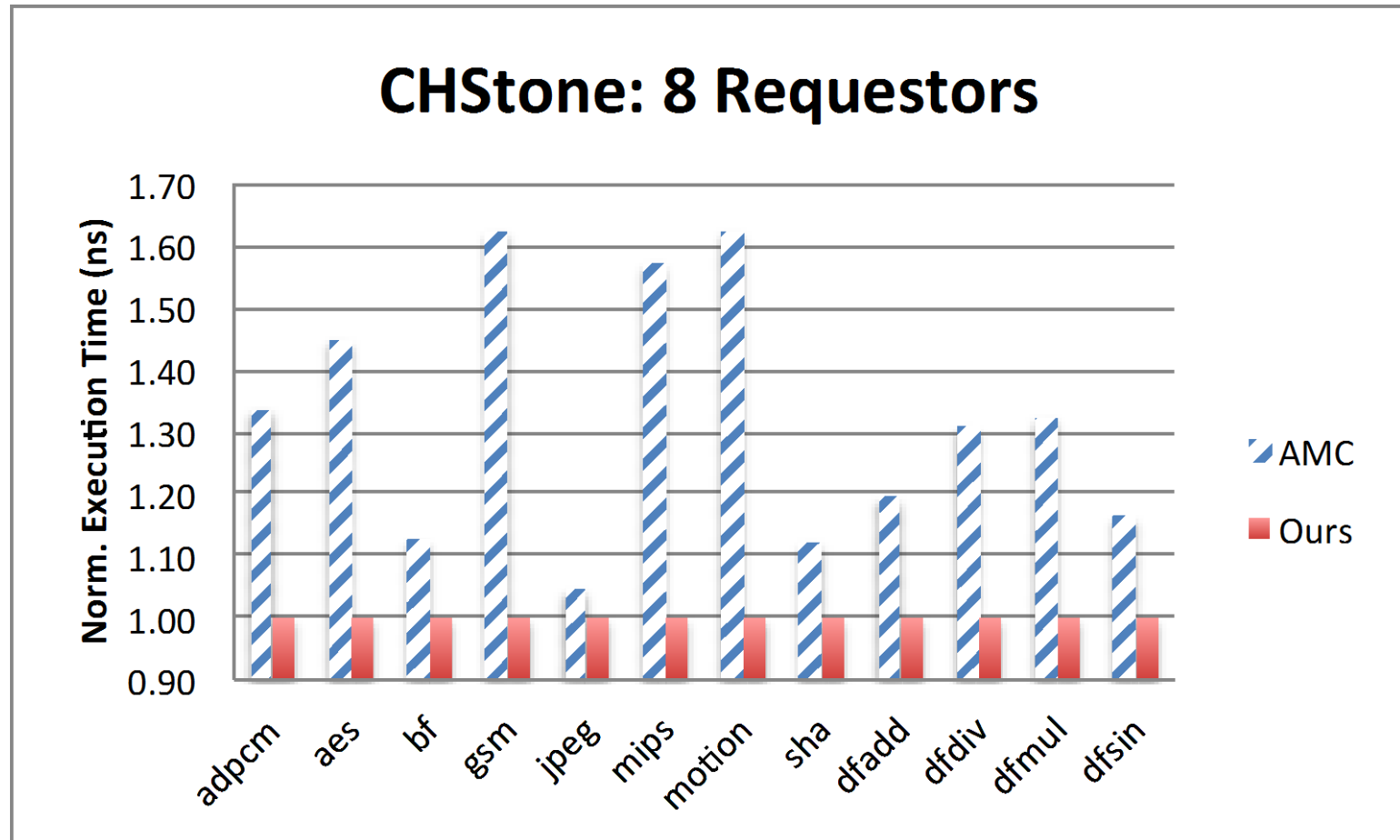
- As memory devices becomes faster, the difference between open and close access is getting larger and therefore close row is becoming too pessimistic

50% Row Hit Ratio, 4 Requestors, 20% Writes

Devices	800D (ns)	1066F (ns)	1333H (ns)	1600K (ns)	1866L (ns)	2133N (ns)	% better
AMC (64 bits)	185	185.27	180.9	178	169.84	163	11.89%
Our (64 bits)	125.2	112.47	104.85	102.18	96.97	92.85	25.84%

Results

- CHStone Benchmarks for 64bits bus



Conclusion

- A novel worst case analysis that takes dynamic state into account
- Open row policy can reduce memory latency as devices are becoming faster
- Private bank scheme is used to eliminate row buffer interference from other requestors

Future Work

- Discussion of shared data
- Bus utilization is still poor due to read/write switching
- Read/Write optimization to reduce latency bound
- Handle Multiple Ranks
- Implementation in hardware

References

- [1] M. Paolieri, E. Quiñones, F. Cazorla, and M. Valero, “An Analyzable Memory Controller for Hard Real-Time CMPs,” *Embedded Systems Letters, IEEE*, vol. 1, no. 4, pp. 86–90, 2009.
- [2] B. Akesson, K. Goossens, and M. Ringhofer, “Predator: a predictable SDRAM memory controller,” in CODES+ISSS, 2007, pp. 251–256.
- [3] S. Goossens, B. Akesson, and K. Goossens, “Conservative Open- page Policy for Mixed Time-Criticality Memory Controllers,” in *DATE*, 2013.
- [4] J. Reineke, I. Liu, H. D. Patel, S. Kim, and E. A. Lee, “Pret dram controller: Bank privatization for predictability and temporal isolation,” in CODES+ISSS, 2011, pp. 99–108.