

Mocktails: Capturing the Memory Behaviour of Proprietary Mobile Architectures

Mario Badr¹, Carlo Delconte², Isak Edo¹, Radhika Jagtap², Matteo Andreozzi², and Natalie Enright Jerger¹

¹ University of Toronto ² Arm



What is Mocktails?

- An open-source tool for synthesizing memory requests
 - Focuses on proprietary mobile workloads running on specialized architectures
- Industry: Generate and distribute models
 - Models hide industry secrets and proprietary details of the workload and architecture
 - Models are **84%** smaller than trace files
- Academia: Generate memory requests from the models to evaluate memory
 - Maximum **7.8%** error on read row hits and **2.8%** error on write row hits

Presentation Outline



Background

Heterogeneity Memory hierarchy Statistical simulation



Mocktails

Modeling Requests Uncovering patterns Synthesizing

Evaluation

Memory controller



Conclusion

Summary of results

A System-on-Chip

- General-purpose cores (CPU)
- Graphics-Processing Unit (GPU)
- Display-Processing Unit (DPU)
- Video Processing Unit (VPU)

Significant real estate allocated to accelerators.



Heterogeneity is Increasing



Apple SoCs – The Cache Hierarchy

□ L1 □ L2 ■ L3



Heterogeneous Systems-on-Chip

- Specialized hardware for commonly used workloads
- Spend area to b Memory hierarchy?
- More IP blocks = varying demands on memory

Academic research in SoC memory hierarchies.

Proprietary IP blocks increasingly used in SoCs.

Statistical simulation can bridge the gap.

Statistical Simulation



Up Next: Mocktails



Background

Heterogeneity Memory hierarchy Statistical simulation



Mocktails

Modeling Requests Uncovering patterns Synthesizing

Evaluation

Conclusion

Memory controller



Summary of results

Modeling a Memory Request



An Example Workload

- Lots of variability
- Hard to find a pattern in the memory accesses
- We can zoom in



Modeling Addresses		Time 1	Address D	Stride
		Z	A	-12
		3	C	8
 Given a starting 		4	В	-4
address, model the High	variability in st	trides is c	lifficult to	4
strides				
	moderac	curately.		16896
 Stride models used in 	moderac	eurately. 8	X	-8
 Stride models used in prior art 	moderac	8 9	X Y	16896 -8 4
 Stride models used in prior art 	moderac	8 9 10	X Y A	16896 -8 4 -16884
 Stride models used in prior art 	moderac	8 9 10 11	X Y A Y	16896 -8 4 -16884 16884
 Stride models used in prior art 	moderac	8 9 10 11 12	X Y A Y B	16896 -8 4 -16884 16884 -16880

Temporal Partitioning			Time 1	Address D	Stride
			2	A	-12
		Time Interval 1	3	С	8
 Divide the sequence of 			4	В	-4
requests in two	no ir	a interval 2 has high variability in			
Two starting addTwo stride mode		stride v	0		
		Time Interval 2	8	Х	-8
Each partition has different behaviour			9	Y	4
different benaviour			10	А	-16884
			11	Y	16884
 Temporal partitioning used in prior art 			12	В	-16880

Spatial Partitioning				Time	Address	Stride
				1	D	
				2	А	-12
				3	С	8
LI			Spatial Partition 1	4	В	-4
 Divide requests into separate address ranges In the 				5	С	4
		ial n	ortitioning r	aducas th		0
		lai partitioning reduces the variance				-8
		ne st	ride feature	for both	partitions.	4
 Each partition has different behaviour 						
			Time	Address	Stride	
			7	W		
 Spatial partitioning used in prior art But tuned for CPUs 		Spatial Partition 2	8	Х	-8	
			9	Y	4	
			11	Y	0	

Partitioning in Two Dimensions

- Temporal: reduces variability in delta times
- Spatial: reduces variability in strides



Write Request

Carefully Dividing Requests

- Dynamic Spatial Partitioning
 - Find reque
 Dynamic spatial partitioning adapts to the memory access behaviour of the workload and device.
- Spatial partitioning uncovers variable-sized time intervals
 - Phases with different start times and durations

Modeling Each Partition

- Each partition consists of a sequence of memory requests
- Model each partition independently
- Save each partition's:
 - Start time
 - Start address







Modeling Each Feature

- Model each feature independently
- Features that do not change:
 - Constant value
- Features that do change:
 - Markov chain



Synthesizing Requests

- Each model is used to generate requests
 - Need initial time and address
- Requests pushed into a priority queue
 - Ordered by timestamp



Up Next: Evaluation



Background

Heterogeneity Memory hierarchy Statistical simulation

Modeling Requests Uncovering patterns Synthesizing

Evaluation

Conclusion

Mocktails

Memory controller



Summary of results

Methodology

- Proprietary memory access traces from Arm
 - CPU, DPU, GPU, VPU devices
- Trace-based simulation with gem5
 - Baseline: Arm traces
 - Requests sent to main memory over a crossbar
- Memory Controller Configuration
 - 4 channels: Each channel has a read and write queue
 - Dynamic scheduling: First-ready, first-come first-serve

Model Comparison

- Perform hierarchical partitioning
 - Model each partition with two different approaches
 - Configuration: Temporal then Spatial (i.e., 2L-TS)
 - 500,000 cycle time intervals
- Mocktails Approach
 - Markov chain or Constant value for each feature (i.e., the McC model)
- STM Approach
 - A statistical simulation technique for the CPU
 - Weighted coin flip for operation feature
 - Markov chain with history for stride feature
 - Other features use Mocktails approach

Absolute Accuracy of Row Hits





Write Row Hits (DPU)





Write Queue Length (Manhattan GPU)



- - Baseline --- McC ······ STM



Up Next: Conclusion



Background

Heterogeneity Memory hierarchy Statistical simulation

Modeling Requests Uncovering patterns Synthesizing

Evaluation

Mocktails

Memory controller



Conclusion

Summary of results

Summary of Results

- Mocktails is accurate for proprietary mobile applications that use specialized hardware
 - Maximum 7.8% error on read row hits and 2.8% error on write row hits
- Mocktails is accurate for traditional CPU benchmarks
 - SPEC CPU2006: **5.6%** average error on L1 cache miss rates
 - More details and exploration in the paper
- Mocktails profiles are distributable
 - 84% smaller than trace files
 - Do not include proprietary details

Mocktails is Open Source

- <u>https://github.com/mariobadr/statistical-simulation</u>
- Mocktails Models of proprietary IP blocks
- Source code for...
 - Mocktails model/trace generation
 - Prior work comparisons (i.e., STM, HRD)
 - gem5 integration
 - Miscellaneous utilities and libraries