

Mach-RT: A Many Chip Architecture for Ray Tracing

Elena Vasiou, Konstantin Shkurko, Erik Brunvand, Cem Yuksel

University of Utah

More Chips! Handle it!



- Higher framerates
- Lower energy

Ray Tracing Performance



- No predictable memory access pattern
 - Is unlike Raster which is able to leverage streaming
- Large energy and time overheads due to irregular DRAM accesses





- Coalescing memory access
 - [Aila & Laine 2009], [Gribble & Ramani 2008], [Spjut et al. 2009]



- Coalescing memory access
 - [Aila & Laine 2009], [Gribble & Ramani 2008], [Spjut et al. 2009]
- Batching and packets for scheduling
 - [Myungbae 2017], [Boulos et al. 2007]



- Coalescing memory access
 - [Aila & Laine 2009], [Gribble & Ramani 2008], [Spjut et al. 2009]
- Batching and packets for scheduling
 - [Myungbae 2017], [Boulos et al. 2007]
- Ordered ray generation and sorting
 - [Eisenacher et al. 2013], [Purcell et al. 2002]



- Better data layout
 - [Hapala et al. 2013], [Viitanen et al. 2017], [Meister & Bittner 2018]



- Better data layout
 - [Hapala et al. 2013], [Viitanen et al. 2017], [Meister & Bittner 2018]
- Compression
 - [Keely 2014], [Ylitie et al. 2017], [Benthin et al. 2018]

Dual Streaming







Traversal Order





Drawbacks





XRay Duplication Early Termination

Memory Traffic





Single Large Chip





Board of Chips





Board of Chips











----- Generic Connection

Board Final





Pixel Distribution



Interleaved pixel assignment

1	2	3	4	1	2
3	1	2	З	4	1
2	3	4	1	2	3
4	1	2	3	4	1
2	3	4	1	2	3
4	1	2	З	4	1

Improvements



- Because of on chip rays:
 - Eliminate ray traffic to DRAM

Improvements



- Because of on chip rays:
 - Eliminate ray traffic to DRAM
 - Ray duplication effect reduced

Improvements



- Because of on chip rays:
 - Eliminate ray traffic to DRAM
 - Ray duplication effect reduced
 - Early Ray Termination

Results



Fairy Forest 174K





Vegetation 1.1M Crytek Sponza 262K





Dragon Sponza 6.6M Dragon Box 870K





San Miguel 10.5M

Hardware Simulation



- Cycle accurate simulator
 - Render to frame completion
- Large simulation times
 - Milliseconds of simulated behavior

System Configuration



	MACH_RT
Frequency	2.0GHz
Threads	512/chip
DRAM bandwidth	512GB/s
L3 cache	64MB
On chip Memory	114MB

Single Large Chip





Single Large Chip: Render Time



Single Large Chip: Render Time



Single Large Chip: Energy



Multiple Chips





Ours



Dual Streaming [Skhurko et al. 2017] STRaTA [Kopta et al. 2013]





- Energy and Bandwidth efficient
- Limited on chip ray queues
- Reconfigurable pipelines

Multiple Chips: Render Time



Number of Threads

Multiple Chips: Render Time



Number of Threads

Multiple Chips: Render Time



Multiple Chips: Energy



U

Cache Lines Transferred



Dual Streaming STRaTA

Cache Lines Transferred



Comparisons to DXR







- Reduced gains
 - Work starvation
 - Large impact of scene fetches

Conclusion



- Multiple chip architecture for ray tracing
 - Is doable by using Dual Streaming
 - No ray traffic to DRAM





Thank you!