#### why I love memory latency

high throughput through latency masking hardware multithreading

.

.

.

.

۰

WALID NAJJAR UC RIVERSIDE

#### the team











#### **the cheerleaders:** Walid A. Najjar, Vassilis J. Tsotras, Vagelis Papalexakis, Daniel Wong

#### the players:

- Robert J. Halstead (Xilinx)
- Ildar Absalyamov (Google)
- Prerna Budhkar (Intel)
- Skyler Windh (Micron)
- Vassileos Zois (UCR)
- Bashar Roumanous (UCR)
- Mohammadreza Rezvani (UCR)



- + background
- filament execution model
- \* application SpMV
- + application Selection
- application Group by aggregation
- + conclusion

#### current DRAM trends



K. K. Chang. Understanding and Improving Latency of DRAM-Based Memory Systems. PhD thesis, Carnegie Mellon University, 2017.

UCR







		The second second	South With	
	ellellh	core I/	Min III	
				a go
			แม้ไหม่อนได้	36.
				-
<b>Core</b>				
Coto				
	Sharz		<b>SARA</b>	
	-Lo Cal			
(Coltan-			ore	
<b>Sole</b>				
				88888
	emory	ontol	er	
				arara
<b>建筑成成实际情况表出的</b> 统计			的的自然性的自然的	State.

Intel Core-i7-5960x		
Launch Date	Q3'14	
Clock Freq.	3.0 GHz	
Cores / Threads	8 / 16	
L3 Cache	20 MB	
Memory Channels	4	
Memory Bandwidth	68 GB/s	



#### • caches > 80% of CPU area,

- big energy sink!
- application must have temporal and/ or spatial locality
- up to 7 levels storage hierarchy (register file, L1, L2, L3, memory, disk cache, disk)



Intel Core-i7-5960x		
Launch Date	Q3'14	
Clock Freq.	3.0 GHz	
Cores / Threads	8 / 16	
L3 Cache	20 MB	
Memory Channels	4	
Memory Bandwidth	68 GB/s	



- caches > 80% of CPU area,
- big energy sink!
- application must have temporal and/ or spatial locality

up to 7 levels
 storage hierarchy
 (register file, L1, L2,
 L3, memory, disk
 cache, disk)



Intel Core-i7	-5960x
Launch Date	Q3'14
Clock Freq.	3.0 GHz
Cores / Threads	8 / 16
L3 Cache	20 MB
Memory Channels	4
Memory Bandwidth	68 GB/s

# caches everywhere! reduces <u>average</u> memory latency assumes <u>locality</u>



#### where cycles go - cache misses



[DaMoN'13] Tözün, P. et al. OLTP in wonderland: where do cache misses come from in major OLTP components? UCR

#### where cycles go - cache misses



[DaMoN'13] Tözün, P. et al. OLTP in wonderland: where do cache misses come from in major OLTP components?

[TODS'07] Chen, S. et al. Improving Hash Join Performance through Prefetching





BIG data (x petaB)



































# You have to know the past to understand the present

- Carl Sagan



- + latency <u>mitigations</u>, aka "<u>caching</u>"
  - \* effective only with spatial/temporal localities
- streaming: latency incorporated into pipeline
  requires spatial locality
- + latency masking, aka hardware multithreading
  - \* fast context switch after memory access
  - \* small thread context
  - \* concurrency == outstanding memory requests









- \* main architect of
  - \* Denelcor HEP (1982)
  - \* Horizon (1985-88) @ IDA SRC
  - \* Tera MTA (1999) later Cray XMT





- \* main architect of
  - \* Denelcor HEP (1982)
  - \* Horizon (1985-88) @ IDA SRC
  - \* Tera MTA (1999) later Cray XMT





#### \* main architect of

- \* Denelcor HEP (1982)
- \* Horizon (1985-88) @ IDA SRC
- \* Tera MTA (1999) later Cray XMT

#### \* latency masking

- \* = do useful work while waiting for memory = <u>multithreading</u>
- switch to a ready process on long-latency I/O operations
- \* => higher throughput
- 50 years ago it was I/O latency => multi-programming



#### \* main architect of

- \* Denelcor HEP (1982)
- \* Horizon (1985-88) @ IDA SRC
- \* Tera MTA (1999) later Cray XMT

#### \* latency masking

- \* = do useful work while waiting for memory = <u>multithreading</u>
- switch to a ready process on long-latency I/O operations
- \* => higher throughput
- 50 years ago it was I/O latency => multi-programming



#### \* main architect of

- \* Denelcor HEP (1982)
- \* Horizon (1985-88) @ IDA SRC
- \* Tera MTA (1999) later Cray XMT

#### Intervention \* latency masking

- \* = do useful work while waiting for memory = <u>multithreading</u>
- switch to a ready process on long-latency I/O operations
- \* => higher throughput
- 50 years ago it was I/O latency => multi-programming

\* https://www.microsoft.com/en-us/research/people/burtons/



## **Tera MTA/Cray XMT**

UCR

- + supercomputer, barrel processor
- + 128 independent threads/processor
- + 2 inst/cycle issued: 1 ALU, 1 memory
- 1 full/empty bit per word insures correct ordering of memory accesses
- + longest memory latency: 128 cycles
- \* randomized memory: reduces bank conflicts
- processor & memory nodes on a 3D torus interconnection network
- + 4096 nodes: sparsely populated



#### 128 hardware threads, in round robin



#### 2 instructions issue per cycle





- CPU that switches between threads of execution on <u>every cycle</u>.
  - also known as "interleaved" or "fine-grained" temporal multithreading.
  - each thread is assigned its own program counter and other hardware registers (architectural state).
  - + guarantee each thread will execute one instruction every n cycles.
  - n-way barrel processor acts like n separate processors,
    - each running at ~ 1/n the original speed.





- + background
- + filament execution model
- Application SpMV
- + application Selection
- + application Group by aggregation
- + conclusion



- + trade bandwidth for latency
- + custom or semi custom processor/accelerator
- + no need for large register files and thread state
  - smaller data path
  - more processors
  - more threads per processor (engine or core)
  - much higher throughput!
- suited for large scale irregular applications:
  - data analytics, databases, sparse linear algebra: SpMV, SpMM, graph algorithms, bioinformatics

## hardware multithreading

UCR














- Call them <u>filament</u>, to distinguish from "threads"
- Preempting executing filaments once it access memory
- Ready & Waiting Filament Queues fit 1000s of thread states
- Massive parallelism: ≥ 4,000 waiting filaments

UĽŔ







#### 8 filament termination









generate i	DATA PATH
	MEMORY
wait queue	



generate i	DATA PATH
	MEMORY
wait queue	

















## a very important



### a very important and often forgotten Little law

## a very important and often forgotten Little law

a queueing system



## a very important and often forgotten Little law



a queueing system

$$s \rightarrow \lambda$$

- N = customers in system, w = wait time
- ★  $\forall$  distribution of **λ** and **S** (service time)



 $s \rightarrow \lambda$ 

- N = customers in system, w = wait time
- ♦  $\forall$  distribution of **λ** and **S** (service time)

 $N = \lambda.w$ 



**S** - Λ

- N = customers in system, w = wait time
- ★  $\forall$  distribution of **λ** and **S** (service time)





S - Λ

- N = customers in system, w = wait time
- ★  $\forall$  distribution of **λ** and **S** (service time)

- L = pipeline latency
- B = memory bandwidth
- C = parallelism





- N = customers in system, w = wait time
- ♦  $\forall$  distribution of **λ** and **S** (service time)



high bandwidth + hyper-pipelining -> massive parallelism

## regular vs. irregular applications

+Regular Applications

- \* Have good temporal and spatial locality
- \* Caches are well suited for these types of applications

#### +Regular Applications

- \* Have good temporal and spatial locality
- \* Caches are well suited for these types of applications

#### +Irregular Applications

- \* Have poor temporal and spatial locality
- \* Caches are NOT well suited for these types of applications
- \* Multithreading is an alternative
  - Can mask long memory latencies
  - Needs high bandwidth to support a large number of concurrent threads

#### +Regular Applications

- \* Have good temporal and spatial locality
- \* Caches are well suited for these types of applications

#### +Irregular Applications

- \* Have poor temporal and spatial locality
- \* Caches are NOT well suited for these types of applications
- \* Multithreading is an alternative
  - Can mask long memory latencies
  - Needs high bandwidth to support a large number of concurrent threads

+Won't increasing the cache size mitigate latency and improve performance?





















## **Convey HC-2ex architecture**





Convey HC-2ex		
Clock Freq.	150 MHz	
Memory Controllers	8	
Memory Channels	16 / FPGA	
Memory Bandwidth	76.8 GB/s	
Outstanding Requests	~ 500 Requests / Channel	
Data bus width	8 Byte	





- multi bank memory architecture
- randomized memory allocation


- + background
- filament execution model
- application SpMV
- + application Selection
- application Group by aggregation
- conclusion

# selection query



- selection is a database operator
  - selects all rows (tuples) in a relation (table)
  - that satisfy a set of conditions
  - expressed as a logical expression on the attributes
  - attributes (columns) are elements of a tuple
  - e.g. ((a1 > 10) ^ (a2>0)) < ((a10<100) ^ (a13<0))</p>
- + a very common database operator
  - poor locality
  - relation can be store row or column wise

# selection example



row/ attribute	a1	a2	a3	a4	а5	a6
0	5	-70	CA	1992	0	12
1	10	70	CO	1990	1	4
2	15	-10	AZ	2001	1	24
3	20	3	NV	1989	0	18

row/ attribute	a1	a2	a3	a4	а5	a6
0	5	-70	CA	1992	0	12
1	10	70	CO	1990	1	4
2	15	-10	AZ	2001	1	24
3	20	3	NV	1989	0	18



row/ attribute	a1	a2	a3	a4	а5	a6
0	5	-70	CA	1992	0	12
1	10	70	CO	1990	1	4
2	15	-10	AZ	2001	1	24
3	20	3	NV	1989	0	18



#### query example:

query A:  $(a5=1) \lor (a3 = NV) \lor [(a1<20)\land(a2>0)]$ query B:  $(a2<0) \land (a3=CA) \land (a5=0)$ 

row/ attribute	a1	a2	a3	a4	а5	a6
0	5	-70	CA	1992	0	12
1	10	70	CO	1990	1	4
2	15	-10	AZ	2001	1	24
3	20	3	NV	1989	0	18



#### query example:

query A:  $(a5=1) \lor (a3 = NV) \lor [(a1<20)\land(a2>0)]$ query B:  $(a2<0) \land (a3=CA) \land (a5=0)$ 



row/ attribute	a1	a2	a3	a4	а5	a6
0	5	-70	CA	1992	0	12
1	10	70	CO	1990	1	4
2	15	-10	AZ	2001	1	24
3	20	3	NV	1989	0	18

#### query example:

query A:  $(a5=1) \lor (a3 = NV) \lor [(a1<20)\land(a2>0)]$ query B:  $(a2<0) \land (a3=CA) \land (a5=0)$ 

#### qualifying rows

row/ query	qA	qB
0	Ν	Y
1	Y	Ν
2	Y	Ν
3	Ν	Ν

# selection example (2)



#### qA predicate control block

#### query example:

query A: (a5=1) ∨ (a3 = NV) ∨ [(a1<20)∧(a2>0)] query B: (a2<0) ∧ (a3=CA) ∧ (a5=0)

predicate control block (PCB) encodes the selection logic:

- early termination (T or F)
- only the required attributes are fetched (no caching)

	cond	lition	action			
attri bute	ор	const	if T	if F		
a5	=	1	Terminate T	a3		
a3	=	NV	Terminate T	a5		
a1	<	20	a2	Terminate F		
a2	>	0	Terminate T	Terminate F		

#### qB predicate control block

condition			action			
attri bute	ор	const	if T	if F		
a2	<	0	a3	Terminate F		
a3	=	CA	a5	Terminate F		
a5	=	0	Terminate T	Terminate F		

# selection CDFG

# UCR

#### Explanations:

- PCB: process control block (stored locally on the FPGA)
- The state of each filament is the current index into the PCB
- Multiple queries can be executed simultaneously on the same relation, the state becomes the (PCB #, index)
- data is returned from memory in the order fetched, on one channel





- Each row consists of 8 fixed size 64-bit columns.
- Query selectivity was varied from 0% to 100% (0%: no row qualifies, 100% all rows qualify)
- Dataset varied from 8M 128M.
- Results were obtained on the following platforms:

Device	Make & Model	Clock (MHz)	Cores	Memory Size (GB)	Memory Bandwidth, GB/s
CPU	Intel Xeon E5-2643	3,300	8	128	51.2
GPU	Nvidia Titan X	1,500	3,584	12	480.0
FPGA	Virtex6 – 760	150	64	64	76.8

#### bandwidth utilization - selection



filament has the highest bandwidth utilization

#### predicate evaluation/sec



#### peak = theoretical max throughput

### runtime (msec)

UCR



at only 150 MHz

FPL 2019 - © W. Najjar

- Lineitem table has 16 columns Measured selectivity is 1.91%
- CPU and GPU access common attributes only once.
- MTP treats them as two independent attributes.

# **TPC-H** benchmark

- Query Q6 for our evaluation because all select conditions apply on the columns of one table and has 5 predicates
- Query Q6:
  - SELECT \* FROM lineitem
  - WHERE I\_shipdate>= date '1995-01-01' AND
  - I\_shipdate<date '1996-01-01' AND
  - I discount between 0.04 AND 0.06 AND
  - l\_quantity<24;

0 **CPU Scalar CPU SIMD** 

Throughput

8

6

2

Tuples/s)

Billion <sup>-</sup>





GPU

MTP



- Memory Fetches: number of columns fetched for evaluation
- Avg Evaluations per Row: predicate comparisons per row
- Effective Bandwidth: total size of the "Lineitem" table processed per unit time / theoretical peak bandwidth, prescribed in [1].
- Peak Bandwidth Utilization: actual amount of data fetched per unit time / theoretical peak bandwidth

	CPU	CPU	CDII	Filamont
	Scalar	SIMD	GFU	
Memory Fetches	100%	18.75%	18.75%	11.4%
Avg Evaluations (per Row)	3	3	3	1.83
Effective Bandwidth Speedup	0.47	3.44	2.01	6.13
Peak Bandwidth Utilization	47.6%	61.2%	36.6%	70.3%

[1] B. Sukhwani, et al. Database Analytics Acceleration Using FPGAs. Parallel Architectures and Compilation Techniques, 2012.



- + background
- filament execution model
- application SpMV
- application Selection
- + application Group by aggregation
- conclusion

#### synchronization in irregular applications

### synchronization in irregular applications





FPL 2019 - © W. Najjar



# synchronization in irregular applications



# synchronization in irregular applications



#### FPL 2019 - © W. Najjar

### synchronizing between threads



- \*Non-deterministic # threads
  - requires inter-thread synchronization
    - only FPGA platforms that supports in-memory synchronization: Convey MX
    - Heavy overhead



- \*Non-deterministic # threads
  - requires inter-thread synchronization
    - only FPGA platforms that supports in-memory synchronization: Convey MX
    - Heavy overhead
- Move the synchronization on-chip
  - use a CAM to cache recently seen nodes: reduces number of memory accesses
  - \* map threads (nodes) to engines deterministically
  - deal with load balancing



- + CAM caches address of all pending memory accesses
  - \* all accesses are checked in CAM first
    - on miss, address is inserted into CAM
    - on hit, access is deferred
    - read hits get the returned value
    - write hits are deferred => strict serialization



- + CAM caches address of all pending memory accesses
  - \* all accesses are checked in CAM first
    - on miss, address is inserted into CAM
    - on hit, access is deferred
    - read hits get the returned value
    - write hits are deferred => strict serialization
- + some writes need not be done in memory
  - \* e.g. aggregation: increment count locally ==> reduced traffic to memory
  - very effective on some workloads



- + CAM caches address of all pending memory accesses
  - \* all accesses are checked in CAM first
    - on miss, address is inserted into CAM
    - on hit, access is deferred
    - read hits get the returned value
    - write hits are deferred => strict serialization
- + some writes need not be done in memory
  - \* e.g. aggregation: increment count locally ==> reduced traffic to memory
  - very effective on some workloads
- + each CAM guards a slice of the address space
  - \* partitioning of the address space
  - \* deal with load balancing



- A crucial operation for any OLAP workload
- Aggregation is challenging to implement in hardware logic
  - \* Not all tuples will create a new node (e.g. update existing node)
  - \* Every tuple reads and writes to the table
- Content Addressable Memories (CAMs)
  - \* Used to enforce memory locks and merge jobs together
  - \* Splits nodes among multiple hash tables, need to be merged



Halstead et al. *FPGA-based Multithreading for In-Memory Hash Joins*, in 7th Biennial Conf. on Innovative Data Systems Research (CIDR'15), Jan. 4-7, 2015, Asilomar, CA.

# thread control flow graph



# thread control flow graph





# aggregation engine



- Main Memory **FPGA** Logic Registers: Base addresses, sizes, etc. **CAM:** Filter Keys Tuple 0 Tuple Request Aggregation Relation Tuple 1 \* Jobs with the same CAM (filter keys) N alguT Hash Function key are merged HT request OxFFFF FFFF CAM: HT Lock Hash FIFO CAM Table 0x0000 003F (HT Lock) Wait for Lock \* Locks individual hash FIFO OxFFFF FFFF table locations HT Lookup list data Linked Lists Linked List 0x0000 0003 Search New LL list data Linked List Job FIFO 0x0000 002A \* Read through to find Rec. LL Analvze Node Job FIFO list data match **OxFFFF FFFF** Update/Insert
- Update/Insert node
  Design is limited by the
- 128 CAM location



- Mapping of nodes to engines done statically
  - \* bandwidth between FPGAs is not sufficient
  - \* may becomes the design bottleneck
- Relation keys can be in multiple nodes in separate hash tables
  - \* keys are not duplicated within a hash table
- After the job is completed
  - \* multiple hash table must be merged: a streaming operation
  - \* can be done on FPGA:
    - stream list chains for each hash table bucket
    - merge matching keys together

# **concurrent FPGA engines**



- + target platform: Convey HC-2ex
  - \* 4 Xilinx Virtex 6 FPGAs x 16 Memory channels/FPGA

# two designs

- replicated cores (4 channels/core, 4 cores/FPGA)
- multiplexed cores (2.5 channels/core, 6 cores/FPGA) write channel shared by 2 cores



**Replicated Engine Design** 



Multiplexed Engine Design

### fine v/s coarse-grain locking








- \* more parallel filaments
- \* higher pressure on CAM



higher pressure on CAM

#### who wins?

## experimental evaluation FGL v/s CGL



- fine-grain ~4x higher throughput
- note: constant throughput v/s relation cardinality

## comparison to software

- 5 aggregation algorithms
  - Independent Tables Hardware-oblivious [VLDB'07] Shared Table

  - Partition & Aggregate Hardware-conscious [VLDB'07]
  - Hybrid Table \*
  - Hybrid Table Partition with Local Aggregation Table

[ICDE'13] Balkesen, C. et al. Main-memory Hash Joins on Multi-core CPUs: Tuning to the underlying hardware. [VLDB'07] J. Cieslewicz et al. Adaptive Aggregation on Chip Multiprocessors. [DAMON'11] Y. Ye et al. Scalable aggregation on multicore processors



Hybrid [DAMON'11]

## experimental evaluation



#### Hardware Region

	FPGA board	Virtex-6 760	
	# FPGAs	2	
	Clock Freq.	150 MHz	
	Engines per FPGA	4/3	
	Memory Channels	32	
	Memory Bandwidth (total)	38.4 GB/s	
	Software	Pogion	
	Soltwale	negion	
C	PU	Intel Xeon E5-2643	
C #	CPUs	Intel Xeon E5-2643 1	
C # C	CPUs Cores / Threads	Intel Xeon E5-2643 1 4 / 8	
C # C	CPUs CPUs Cores / Threads	Intel Xeon E5-2643 1 4 / 8 3.3 GHz	
C # C L	CPUs CPUs Cores / Threads Clock Freq. 3 Cache	Intel Xeon E5-2643 1 4 / 8 3.3 GHz 10 MB	

#### Software Approaches

- Individual Table
- Shared Table (atomic)
- \* Hybrid Table
- \* Partitioning
- \* PLAT

### Datasets

- Uniform, Zipf 0.5, Heavy Hitter, Moving Cluster, Self Similar
- 8M to 64M tuples/ benchmarks
- 1K to 4M cardinality

## uniform workload



Uniform - 256 M tuples



Heavy Hitter - 256 M tuples

UCR



## bandwidth utilization



### + FPGA implementation achieves much higher memory bandwidth utilization





- + background
- filament execution model
- application SpMV
- + application Selection
- application Group by aggregation
- + conclusion



- + Big data is: **BIG** and multi-dimensional
  - \* access patterns to data are extremely irregular
  - \* no ideal storage scheme constant poor locality

### +Other applications

\* string matching, hash join, group by aggregation, selection, outer joins, sorting

## +Challenges

- \* inter-filament synchronization
- \* dynamic load re-balancing (at run time)

## publications



#### http://www.cs.ucr.edu/~najjar/publications.html

- P. Budhkar, I. Absalyamov, V. Zois, S. Windh, W. A. Najjar and V. J. Tsotras. <u>Accelerating In-</u> <u>Memory Database Selections Using Latency Masking Hardware Threads</u>, in ACM TACO, 2019.
- I. Absalyamov, R. J. Halstead, P. Budhkar, W. A. Najjar, S. Windh, V. J. Tsotras. <u>FPGA-Accelerated Group-by Aggregation Using Synchronizing Caches</u>. 12th Int. Workshop on Data Management on New Hardware (DaMoN 2016), Co-located with ACM SIGMOD/PODS, San Francisco, USA, June 27, 2016.
- S. Windh, P. Budhkar and W. A. Najjar. <u>CAMs as Synchronizing Caches for Multithreaded</u> <u>Irregular Applications on FPGAs</u>. ICCAD 2015: 331-336
- R. J. Halstead, I. Absalyamov, W. A. Najjar and V. J. Tsotras. <u>FPGA-based Multithreading for In-Memory Hash Joins</u>, in 7th Biennial Conference on Innovative Data Systems Research (CIDR 15), January 4-7, 2015, Asilomar, California.
- E. B. Fernandez, J. Villarreal, S. Lonardi, and W. A. Najjar. <u>FHAST: FPGA-based acceleration of</u> <u>Bowtie in hardware</u>, in IEEE/ACM Trans. on Computational Biology and Bioinformatics, # 99, Feb. 2015.
- R. J. Halstead, W. A. Najjar and O. Huseini. <u>SpVM Acceleration with Latency Masking Threads</u> on <u>FPGAs</u>. Technical Report UCR-CSE-2014-04001

## conclusion



## advantages of filament execution

- very small state allows
  - fast context switching between filaments
  - small storage for waiting queues
- masking the latency of memory and non volatile storage
  - no reliance on locality, no caches
  - Iower energy consumption
- limitations: custom or semi-custom accelerators



# thank you!

# questions?