

Exploiting 3D Memory for Accelerated In-Network Processing of Hash Joins in Distributed Databases



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Johannes Wirth¹, Jaco A. Hofmann¹, Lasse Thostrup², Andreas Koch¹, Carsten Binnig²

1) Embedded Systems and Applications Group, TU Darmstadt

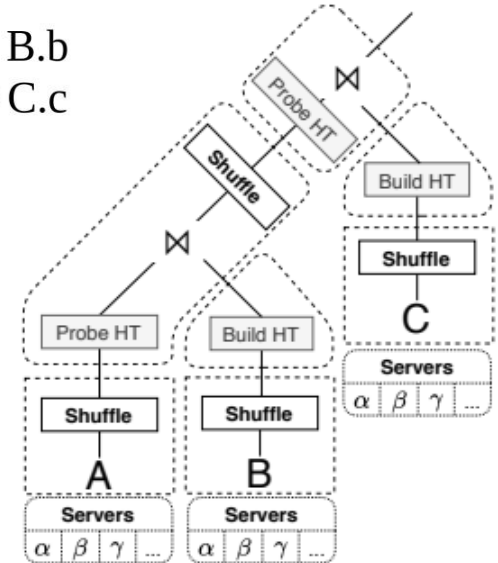
2) Data Management Lab, TU Darmstadt

BACKGROUND

SQL Join Operation

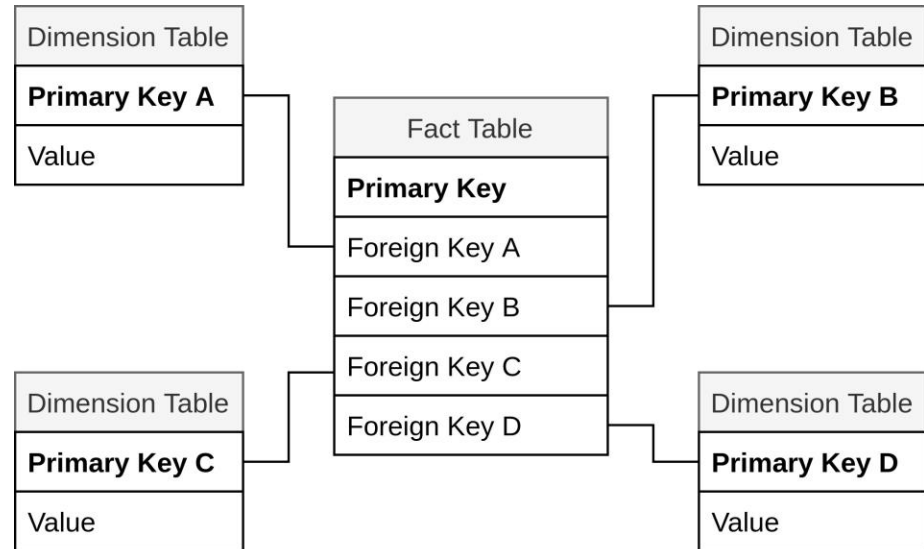
- Common database operation
- Merges two tables/relations
 - Combine entries with matching values
- Implementation: hash join algorithm
- Problematic in distributed databases
→ Shuffling

```
SELECT A.a,B.d,C.e  
FROM A,B,C  
WHERE A.b = B.b  
      AND A.c = C.c
```



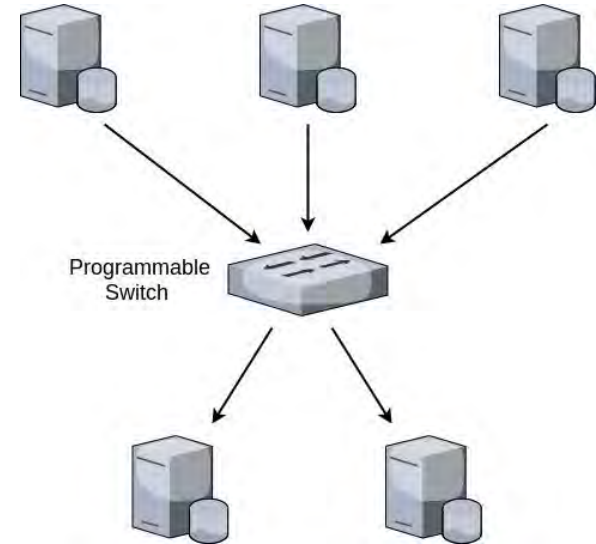
Star Schema

- Large fact table
- Smaller dimension tables
- Foreign key in fact table points to entry in corresponding dimension table
- SQL Join to combine them



In-Network Processing

- Process in-flight data
→ Reduce network traffic
- Programmable switches (P4, Tofino, ...)
 - Limited flexibility
 - Only small memory
 - Mostly stateless operations



OUR CONTRIBUTION

- Network accelerator for SQL Join operation
- Multiple joins at once
- Tailored to Star Schema
- Builds on prior work [9]
→ Make use of HBM to improve performance

High-Performance In-Network Data Processing

Jaco Hofmann
Embedded Systems and
Applications Group
TU Darmstadt
Hochschulstr. 10
64289 Darmstadt
hofmann@esa.tu-
darmstadt.de

Lasse Thostrup
Data Management
TU Darmstadt
Hochschulstr. 10
64289 Darmstadt
lasse.thostrup@cs.tu-
darmstadt.de

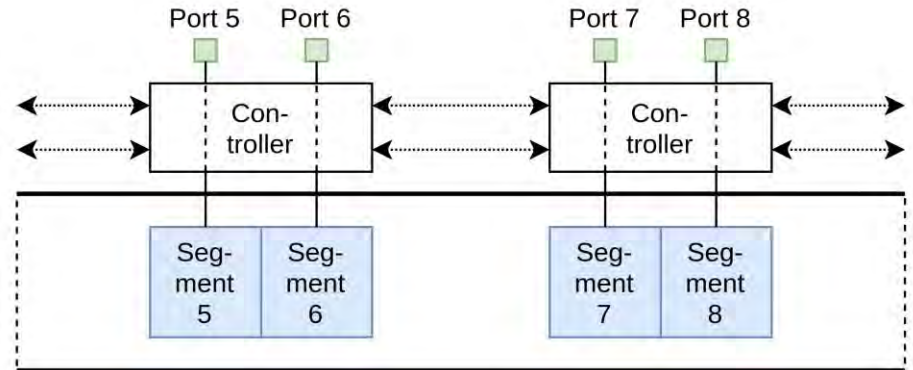
Tobias Ziegler
Data Management
TU Darmstadt
Hochschulstr. 10
64289 Darmstadt
tobias.ziegler@cs.tu-
darmstadt.de

Carsten Binnig
Data Management
TU Darmstadt
Hochschulstr. 10
64289 Darmstadt
carsten.binnig@cs.tu-
darmstadt.de

Andreas Koch
Embedded Systems and
Applications Group
TU Darmstadt
Hochschulstr. 10
64289 Darmstadt
koch@esa.tu-
darmstadt.de

High-Bandwidth Memory

- Fast on-chip memory
- Xilinx FPGAs:
 - 4 – 16 GB
 - 460 GB/s
- Highly parallel (32 channels)
- Optional switch: not used here



Operation

Operation

```
SELECT A.a,B.e,C.f,D.g  
FROM A,B,C,D  
WHERE A.b = B.b  
      AND A.c = C.c  
      AND A.d = D.d
```

Table A

a	b	c	d	Server
a1	1	1	1	α
a2	2	1	3	β
a3	3	2	2	γ
a4	1	2	1	α
a5	2	3	3	β
a6	3	3	2	γ

Table B

b	e	Server
1	e1	α
2	e2	β
3	e3	γ

Table C

c	f	Server
1	f1	α
2	f2	β
3	f3	γ

Table D

d	g	Server
1	g1	α
2	g2	β
3	g3	γ

Hash Table: 2 Buckets, 2 Slots per Bucket

Collision Handling: Use next free Slot in Bucket

Hashing: $\text{bucket}(x) = x \% 2$

Hashing 1

Hashing: $\text{bucket}(x) = x \% 2$

Collision Handling: Use next free Slot in Bucket

Table B

b	e	Server
1	e1	α
2	e2	β
3	e3	γ

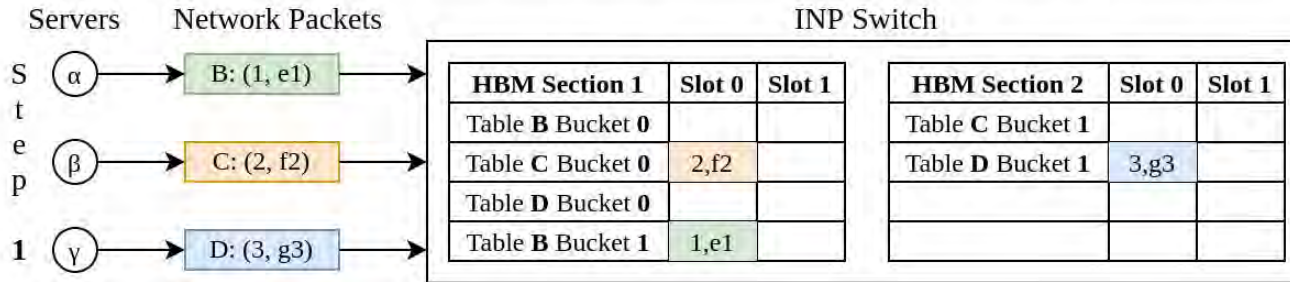
Table C

c	f	Server
1	f1	α
2	f2	β
3	f3	γ

Table D

d	g	Server
1	g1	α
2	g2	β
3	g3	γ

Hashing Phase: Transfer dimension tables B,C,D to INP switch and store in Hash Tables



Hashing 2

Hashing: $\text{bucket}(x) = x \% 2$

Collision Handling: Use next free Slot in Bucket

Table B

b	e	Server
1	e1	α
2	e2	β
3	e3	γ

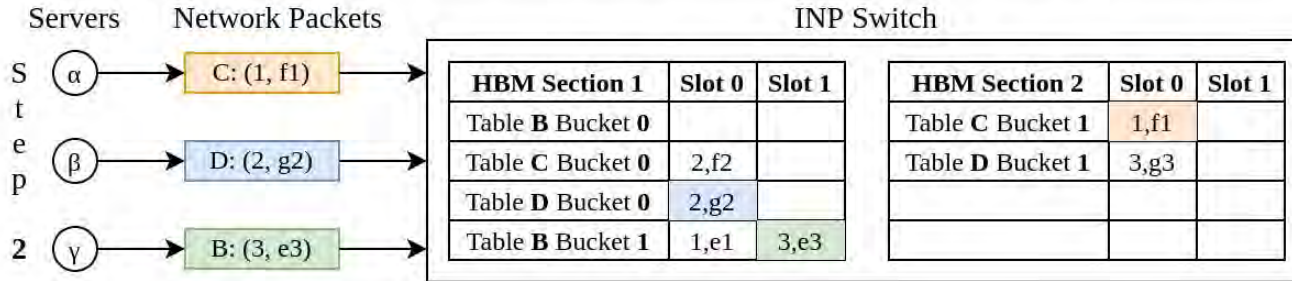
Table C

c	f	Server
1	f1	α
2	f2	β
3	f3	γ

Table D

d	g	Server
1	g1	α
2	g2	β
3	g3	γ

Hashing Phase: Transfer dimension tables B,C,D to INP switch and store in Hash Tables



Hashing 3

Hashing: $\text{bucket}(x) = x \% 2$

Collision Handling: Use next free Slot in Bucket

Table B

b	e	Server
1	e1	α
2	e2	β
3	e3	γ

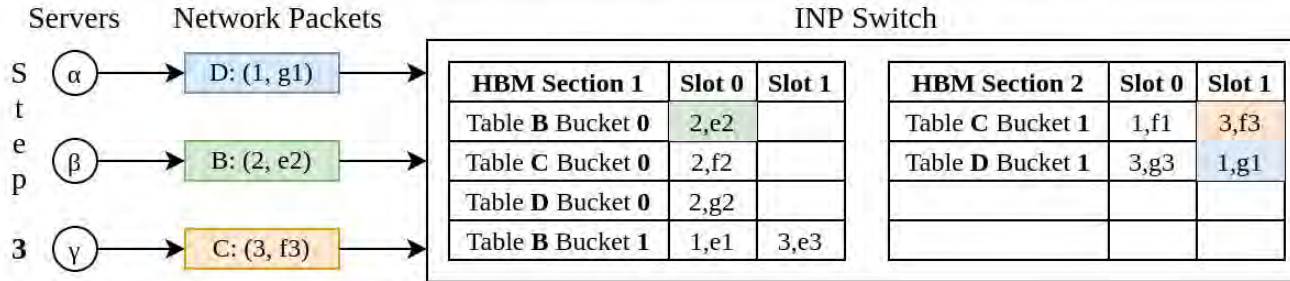
Table C

c	f	Server
1	f1	α
2	f2	β
3	f3	γ

Table D

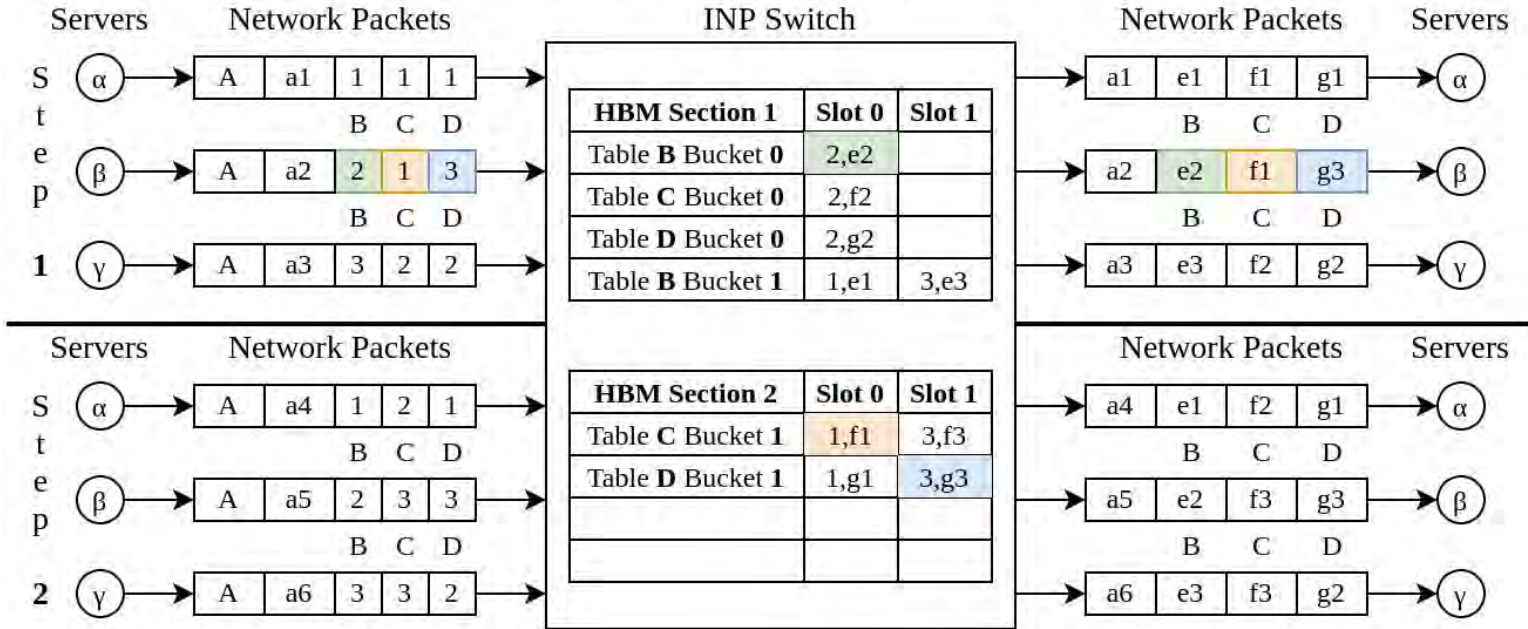
d	g	Server
1	g1	α
2	g2	β
3	g3	γ

Hashing Phase: Transfer dimension tables B,C,D to INP switch and store in Hash Tables



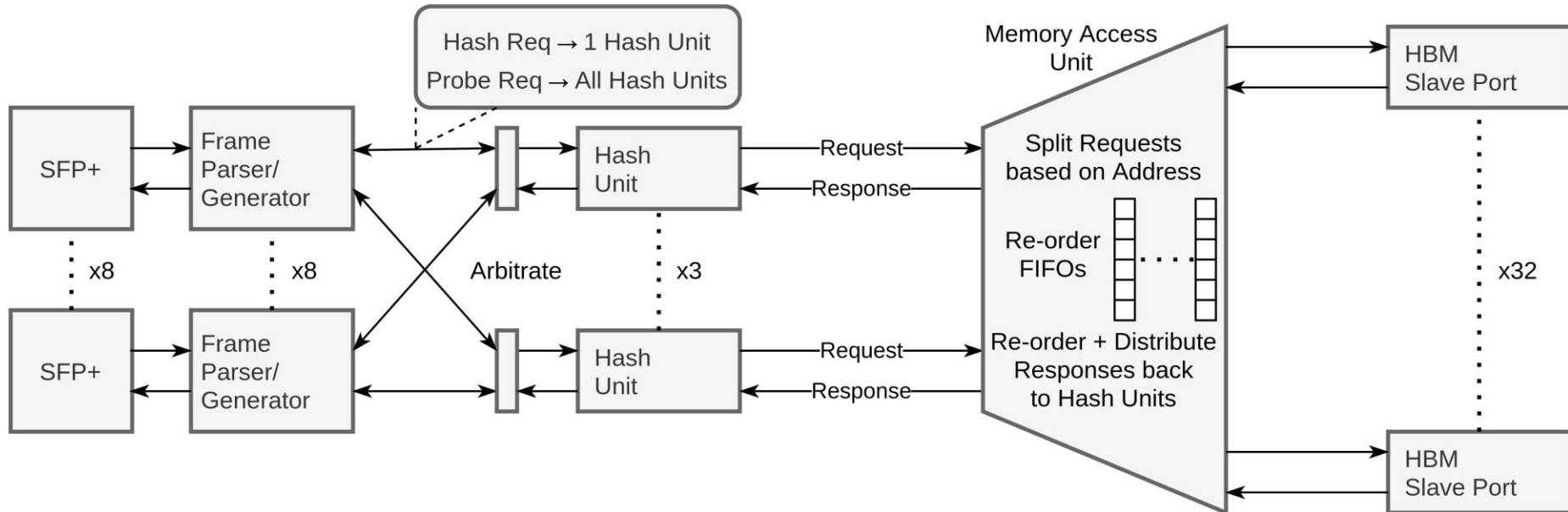
Probing

Probing Phase: Query INP switch with fact table A to retrieve join results via Hash Tables

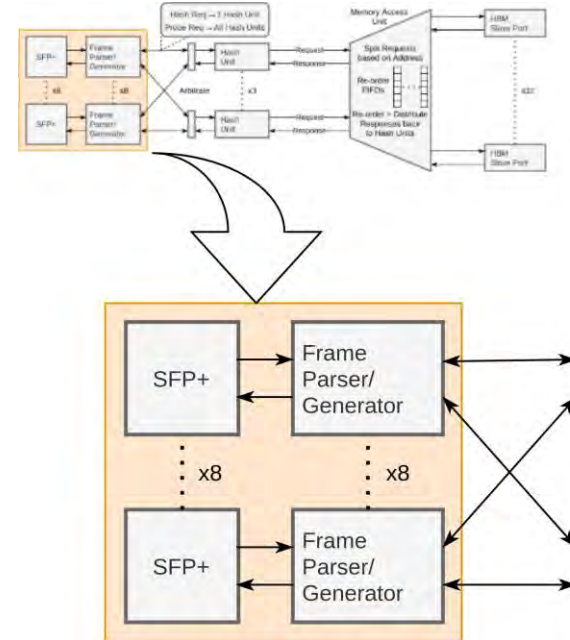


ARCHITECTURE

Architecture Overview

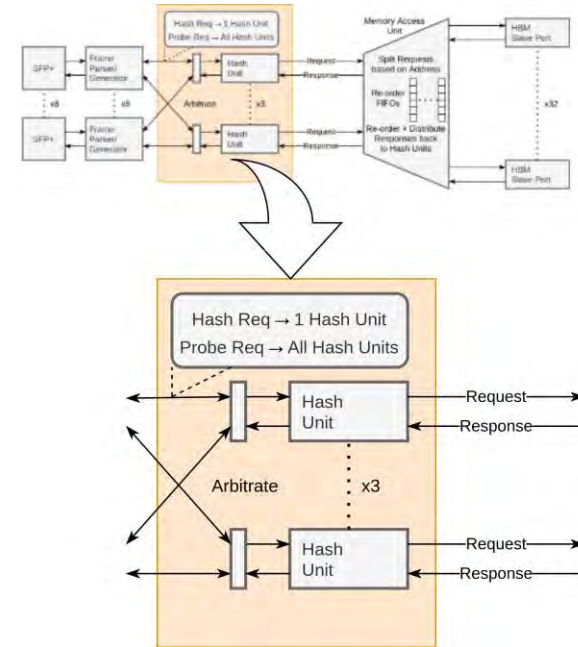


- Parse/create ethernet frames
- 10G Ethernet
- Separate clock domain (156 MHz)
- Must not block
 - Drop frame if queue full
 - Re-request data from sender



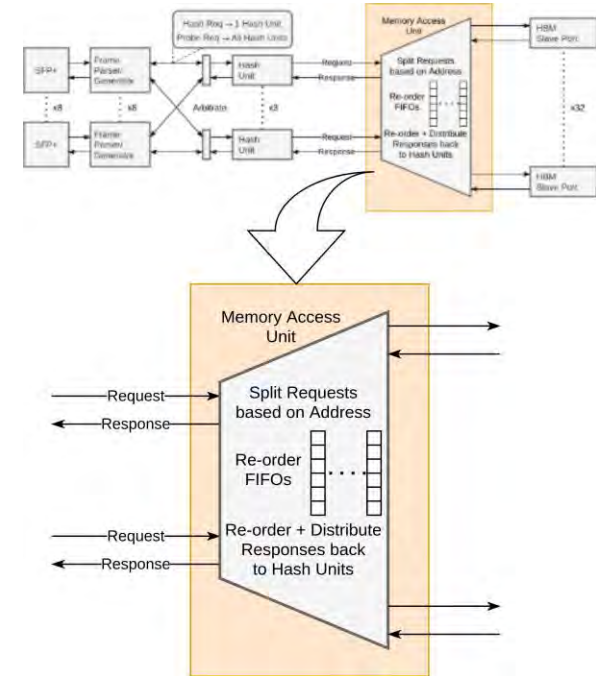
Hash Unit

- One hash unit per join
- Throughput optimized
- Four slots per bucket
 - Bucket width matches HBM width
- Hashing: apply integer hash function to key
 - $\text{Bucket} = \text{hash} \% \#\text{buckets}$



Memory Access

- Connect hash units to all HBM ports
 - Challenge: Routability
 - Standard interconnect not feasible
- Intermediate stage (4 HBMs per group)
- Re-order responses to hash units
 - Keep track of original request order



EVALUATION

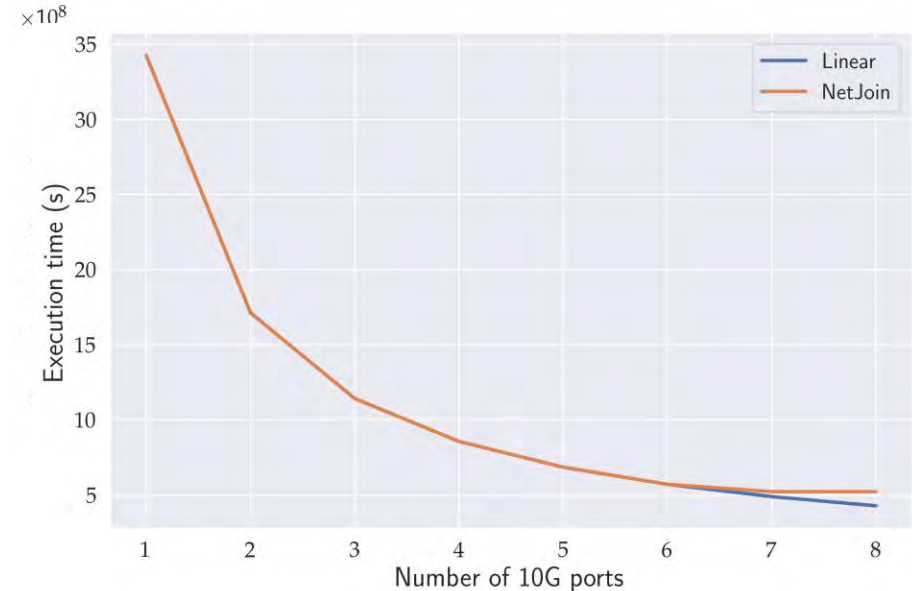
Evaluation Setup

- Scenario: fact table joined with three dimension tables
- Evaluated on Bittware XUP-VVH (8GB HBM; 8 x 10G Ethernet)
- 8 database servers

- Performance scaling
- Comparison to software baseline & prior work

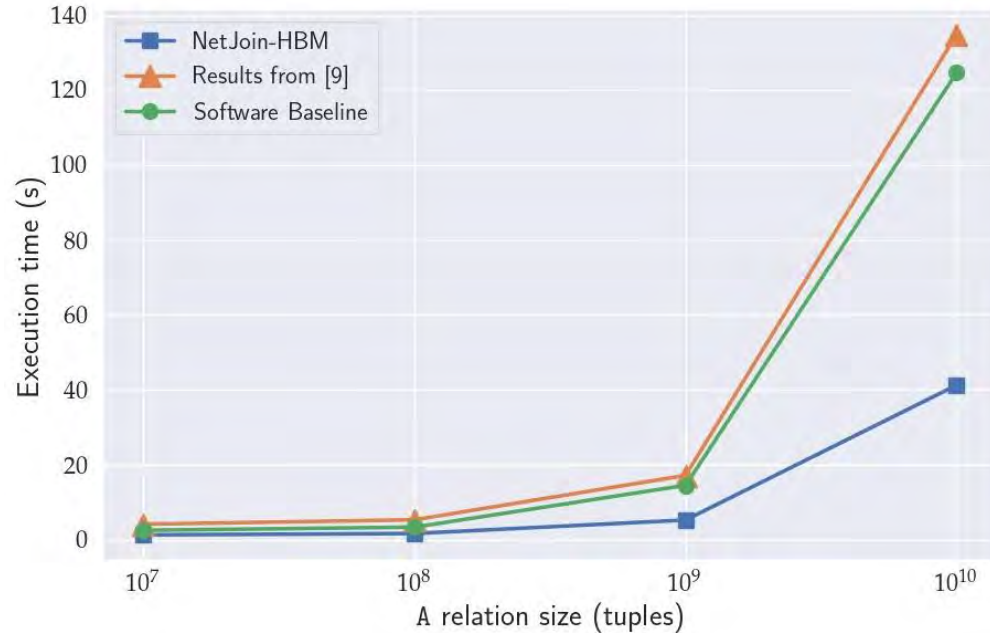
System Performance

- Linear scaling until 6 x 10G
- Similar for both phases
- Probing faster than hashing
- Probing caps at 250 MIOPS
→ bottleneck hash units (250 MHz)



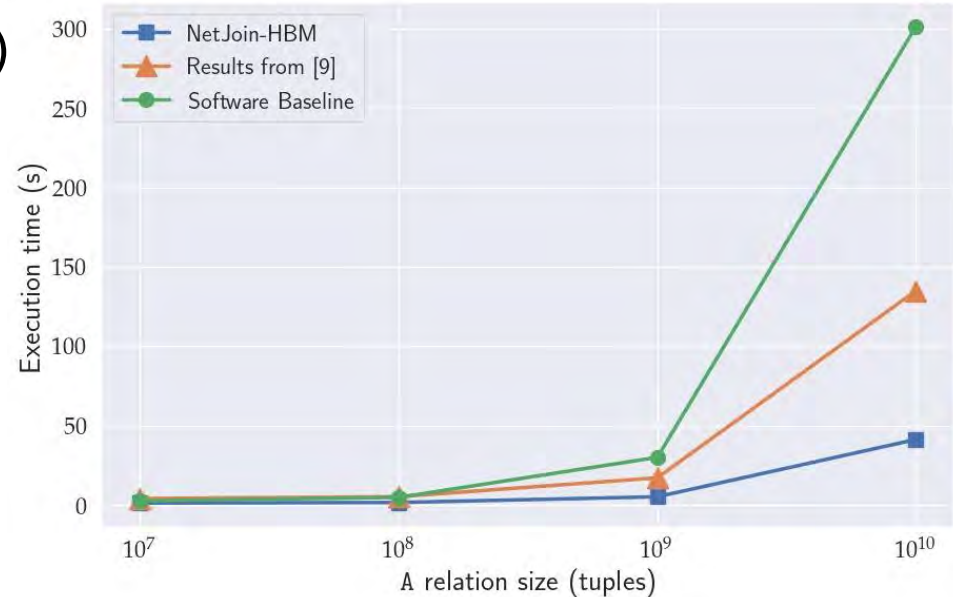
Baseline Comparison

- Dimension table fixed (100×10^6)
- Fact table size varied
- Equally distributed
→ Optimal for baseline
- Up to three times faster

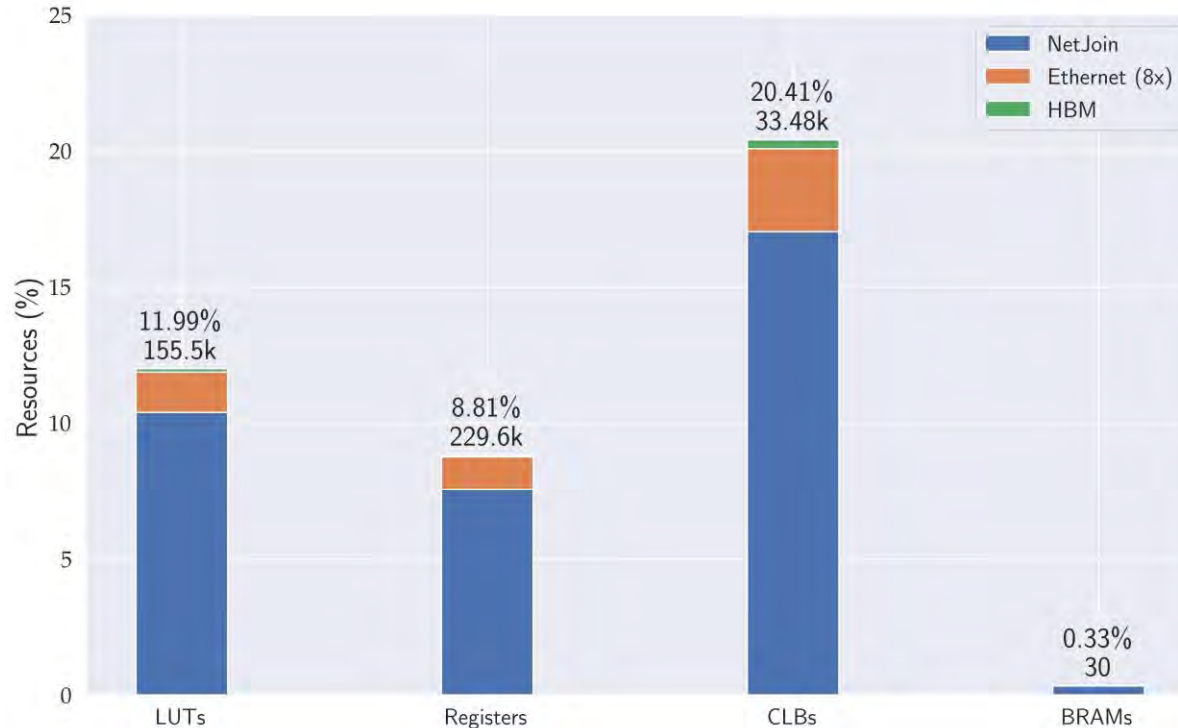


Baseline Comparison (Skewed)

- Dimension table fixed (100×10^6)
- Fact table size varied
- Skewed distribution (34% - 2%)
- Up to 7 times faster vs baseline
- Up to 3 times faster vs prior work



Resource Usage



Conclusion

- Handles at least 60 Gbit/s
- Up to 7x faster than software baseline
- Limitations
 - Limited HBM memory
 - Bottleneck: hash units
 - No FPGA-based switch

THANKS FOR YOUR ATTENTION!

E-Mail: wirth@esa.tu-darmstadt.de