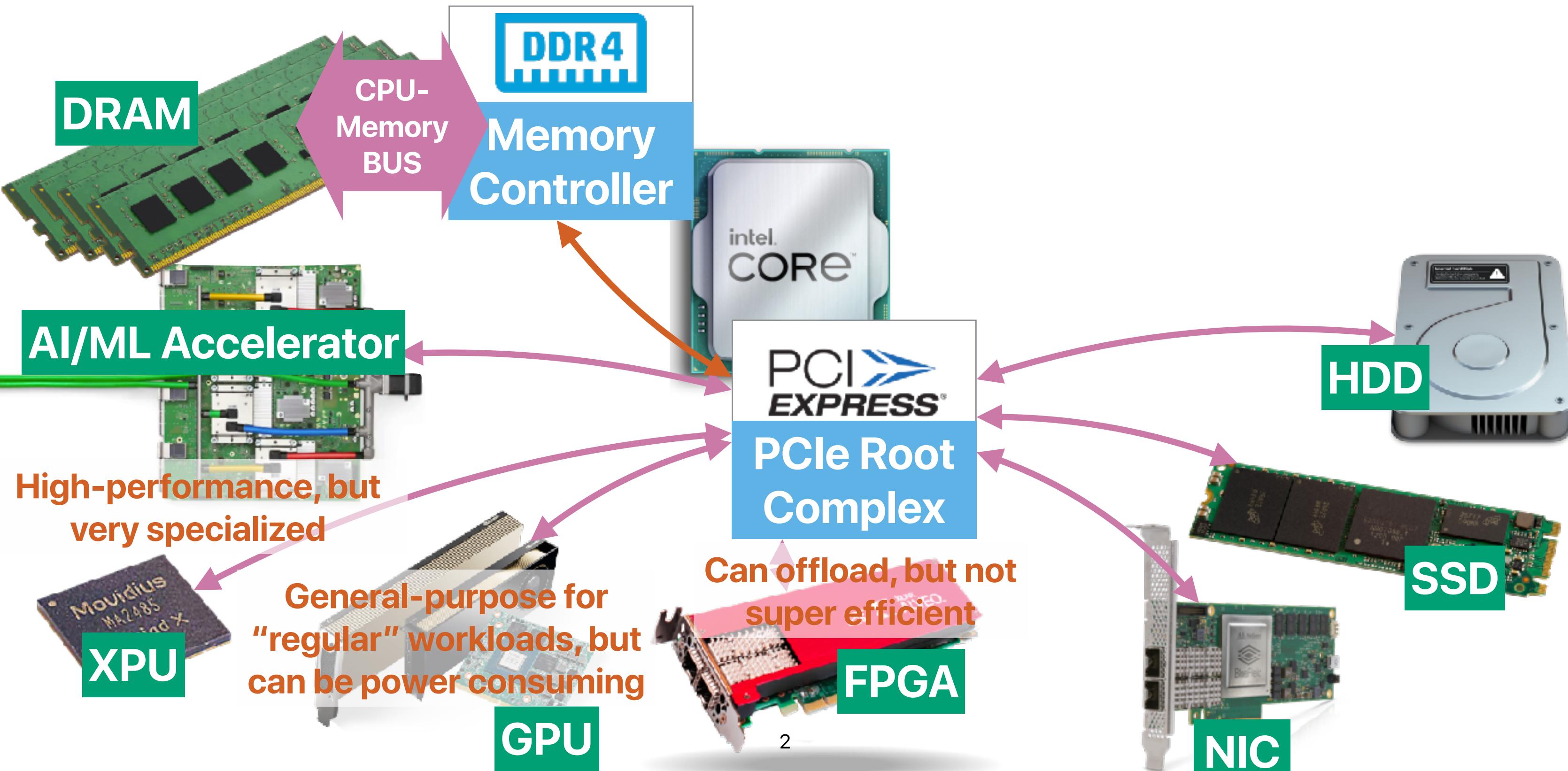


General purpose computing on hardware accelerators (2)

Hung-Wei Tseng

Recap: The landscape of modern computers



Recap: The concept of approximate computing on NPUs

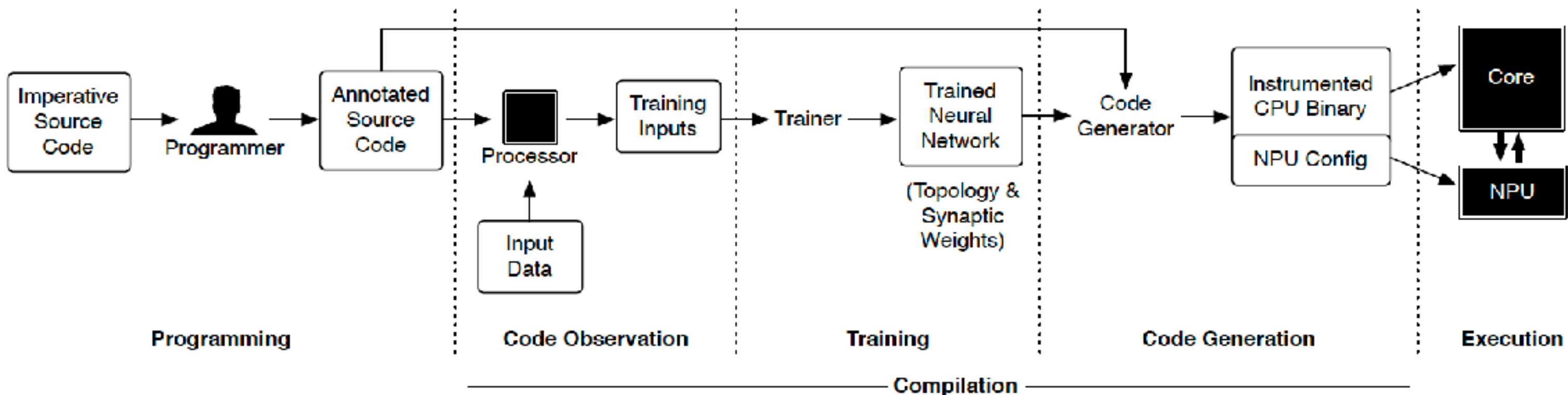


Figure 1: The Parrot transformation at a glance: from annotated code to accelerated execution on an NPU-augmented core.

Recap: GPTPU framework

```
#include <stdio.h>
#include <stdlib.h>
#include <gptpu.h>

// The TPU kernel
void *kernel(openctpu_buffer *matrix_a,
             openctpu_buffer *matrix_b,
             openctpu_buffer *matrix_c)
{
    // invoke the TPU operator
    openctpu_invoke_operator(conv2D, SCALE, matrix_a, \
                             matrix_b, matrix_c);
    return 0;
}

int main(int argc, char **argv)
{
    float *a, *b, *c; // pointers for raw data
    openctpu_dimension *matrix_a_d, *matrix_b_d, *matrix_c_d;
    openctpu_buffer * tensor_a, * tensor_b, * tensor_c;
    int size; // size of each dimension

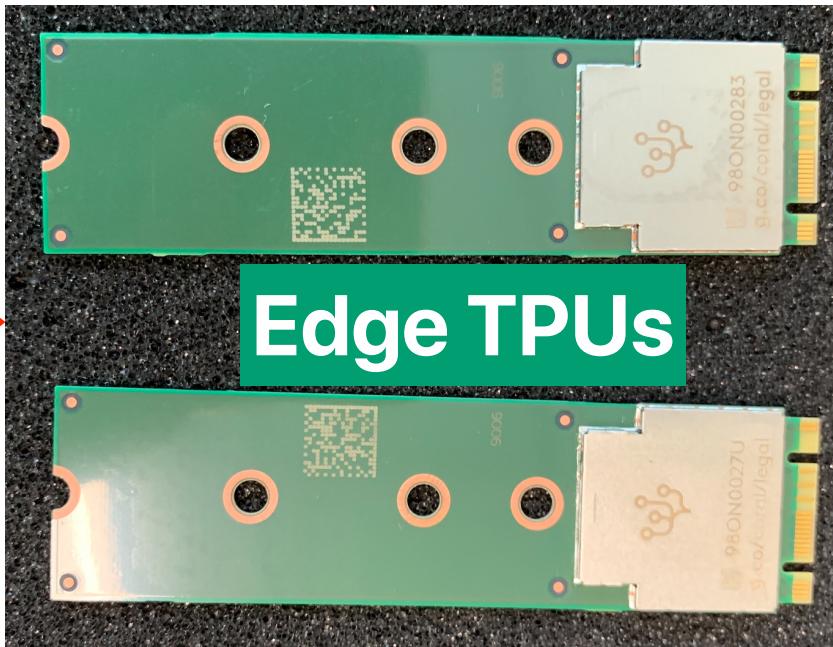
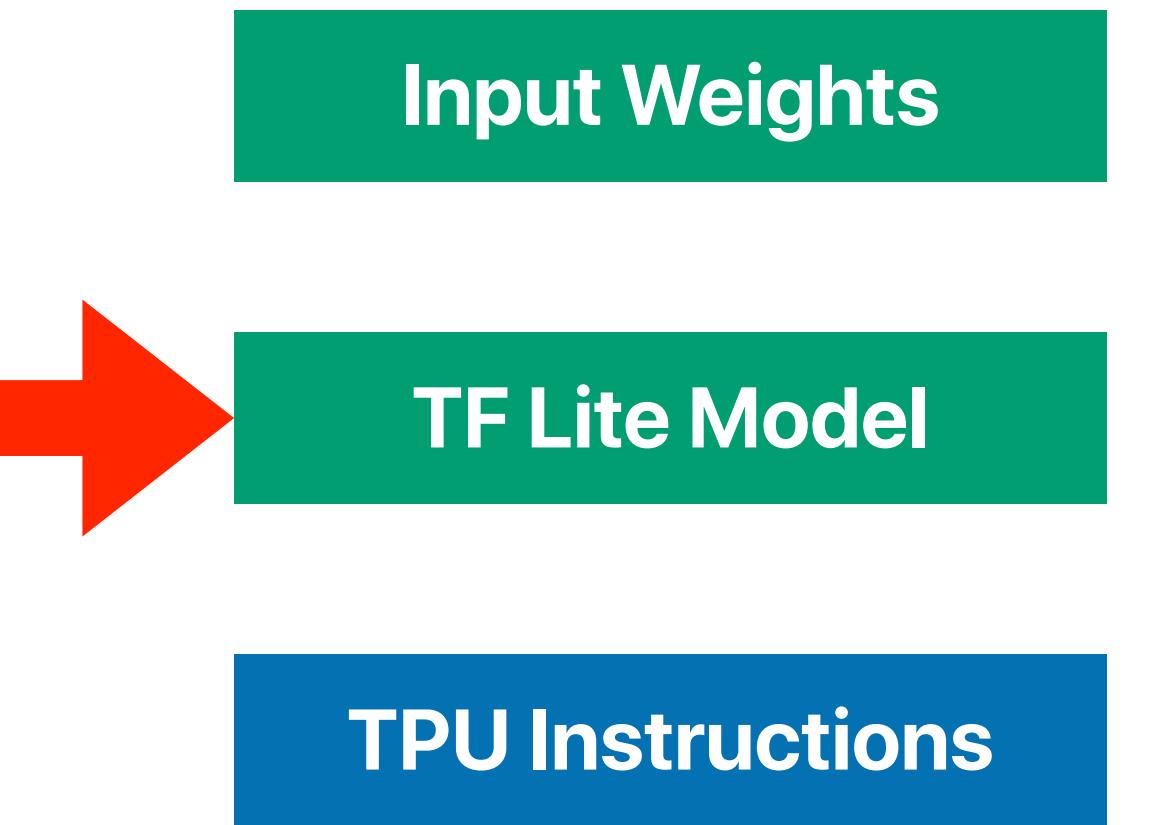
    // skip: data I/O and memory allocation/initialization

    // describe a 2-D tensor (matrix) object for a
    matrix_a_d = openctpu_alloc_dimension(2, size, size);
    // describe a 2-D tensor (matrix) object for b
    matrix_b_d = openctpu_alloc_dimension(2, size, size);
    // describe a 2-D tensor (matrix) object for c
    matrix_c_d = openctpu_alloc_dimension(2, size, size);

    // create/fill the tensor a from the raw data
    tensor_a = openctpu_create_buffer(matrix_a_d, a);
    // create/fill the tensor b from the raw data
    tensor_b = openctpu_create_buffer(matrix_b_d, b);
    // create/fill the tensor c from the raw data
    tensor_c = openctpu_create_buffer(matrix_c_d, c);

    // enqueue the matrix_mul TPU kernel
    openctpu_enqueue(kernel, tensor_a, tensor_b, tensor_c);
    // synchronize/wait for all TPU kernels to complete
    openctpu_sync();

    // skip: the rest of the program
    return 0;
}
```



Can you compare Approx. on NPU & GPTPU? Under what situations would you prefer each?

Approx. on NPU v.s. GPTPU

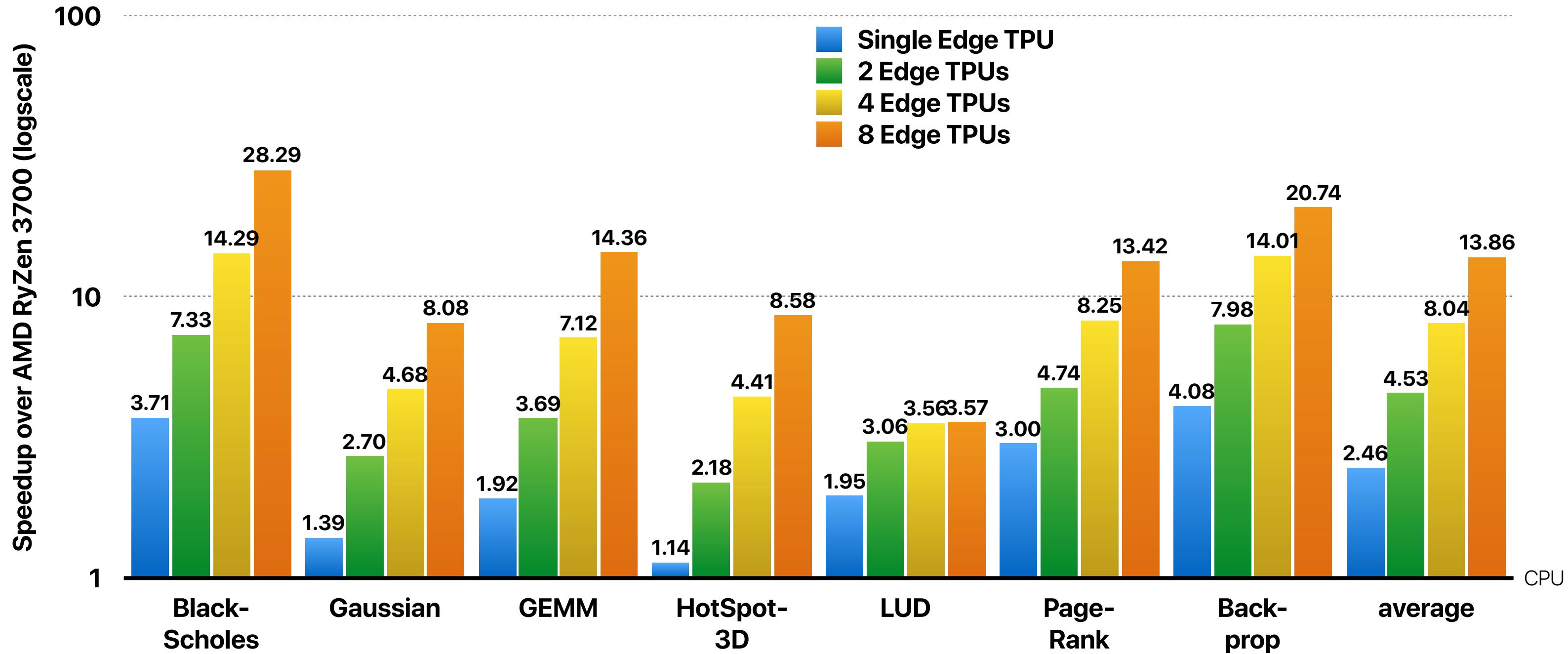
Recap: Approx. on NPU v.s. GPTPU

	Approx. on NPU	GPTPU
Performance	Better if the algorithm is more complex, but worse otherwise	Depending on the complexity of the algorithm
Result quality	Approximate, not accurate	Exact — if the hardware supports the desired precision
Programming efforts	Easier — the programmer does not need to take care of programming on AI/ML accelerators	Difficult — the programmer has to map the computation into operators that the hardware supports
Applicability	Any algorithm on an application tolerating inexactness	Only algorithms that map well to the hardware architecture

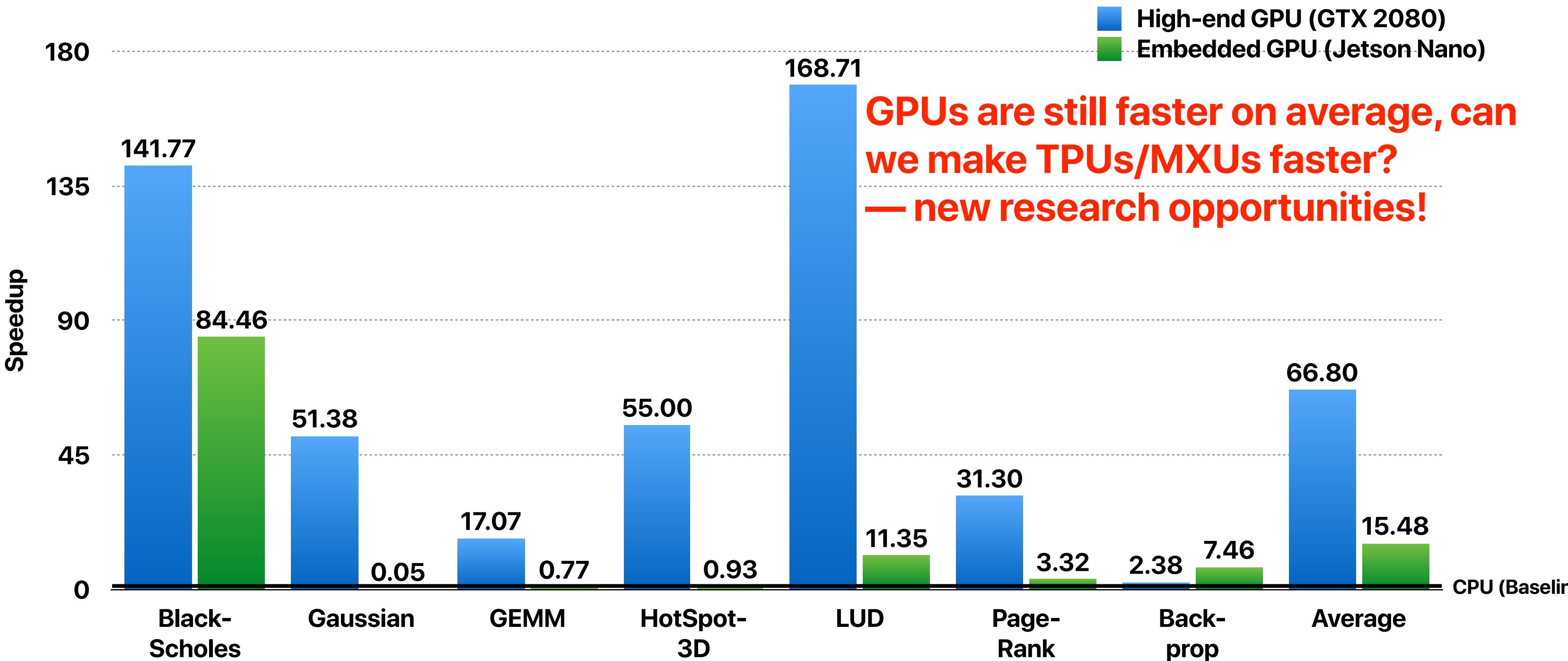
Outline

- TCUDB: Accelerating Database with Tensor Processors
- RTNN: accelerating neighbor search using hardware ray tracing

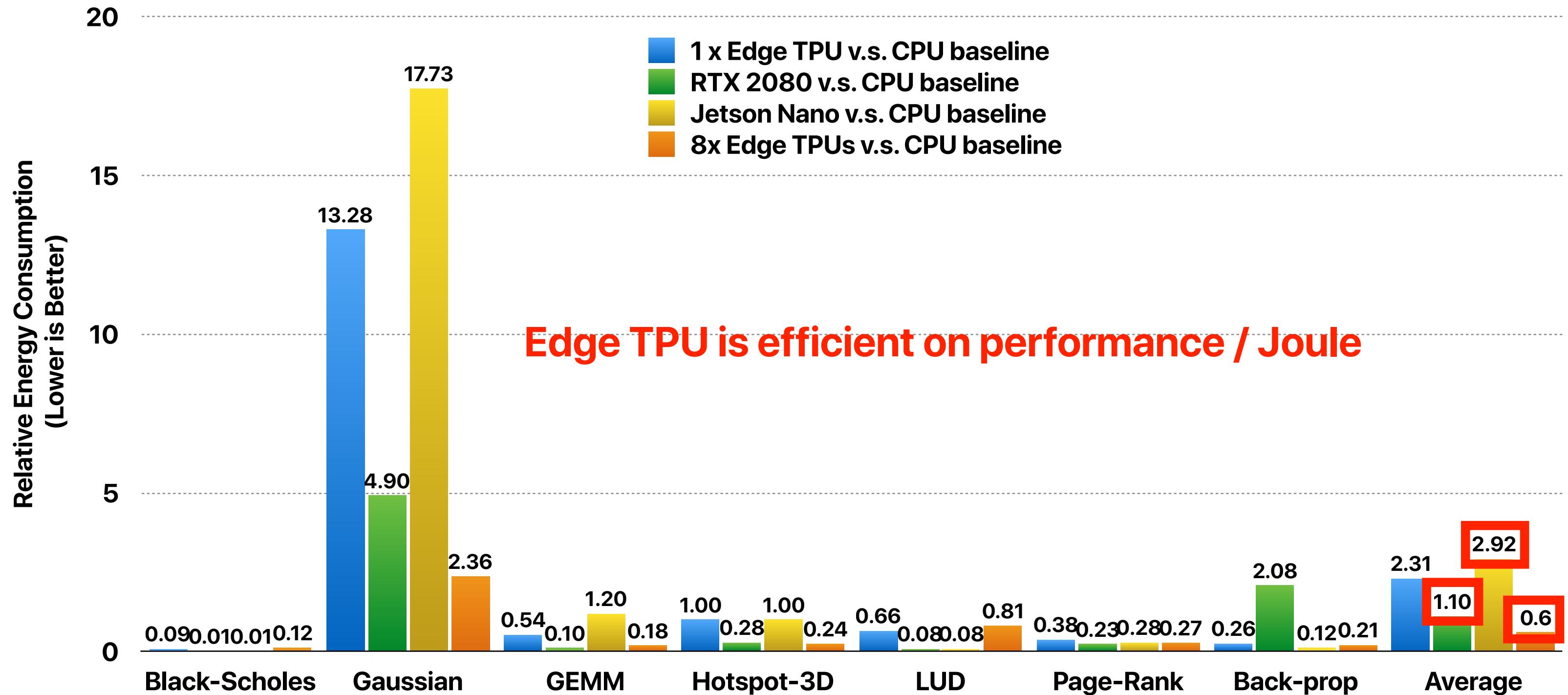
Performance for general-purpose Workloads



Compare 8 Edge TPUs with GPUs



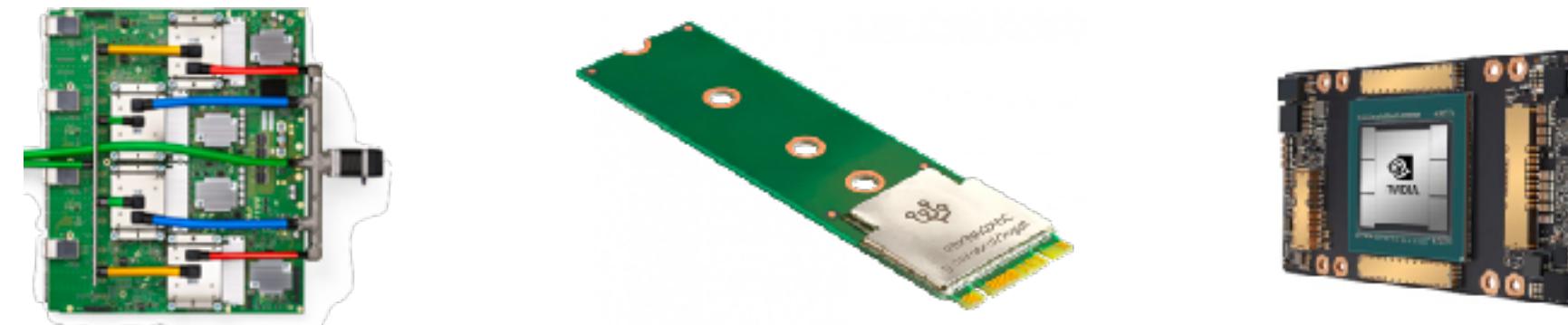
Compare to GPUs



Lessons learned from developing “GPTPU”

- Reverse engineering is interesting but not worth the time
- Developing software frameworks on hardware with steady support
- The accessibility to low-level and backward compatible APIs
- Estimating the performance gain of design using the right/appropriate metrics
- Algorithm implementations should fit the most efficient operations
- The precision support is still below the demand of applications beyond AI/ML

TPU v.s. Edge TPU v.s. Tensor Cores



	TPU v4	Edge TPU	RTX 3090 Tensor Cores
Year	2021	2018	2020
Precision	BF16	INT8	FP16/BF16/TF32
Accessibility	Google Cloud Only	Public Available	Public Available
Programming Framework	Tensorflow/PyTorch	Tensorflow/TFLite/C++	CUDA/Tensorflow/Pytorch
TOPS	275	4	285
Power (Watt)	170	2	350
TOPS/W	1.62	2	0.814285714285714
Device Memory	32 GB	8MB	24 GB

Lessons learned from developing “GPTPU”

- Reverse engineering is interesting but not worth the time
- Work with the company to get important insights/hints/feedbacks
- Developing software frameworks on hardware with steady support
- The accessibility to low-level and backward compatible APIs
- Estimating the performance gain of design using the right/appropriate metrics
- Algorithm implementations should fit the most efficient operations
- The precision support is still below the demand of applications beyond AI/ML

TCUDB: Accelerating Database with Tensor Processors

Yu-Ching Hu, Yuliang Li* and Hung-Wei Tseng

University of California, Riverside and Megagon Labs*

In The ACM SIGMOD/PODS 2022 International Conference on Management of Data.

<https://github.com/escalab/TCUDB>

Who is my most loyal customer?

$O(m \times n)$

```
SELECT A.customer, B.brand, SUM(A.Quantity * B.Price) AS value FROM A INNER JOIN B WHERE ON
A.ProductID = B.ProductID GROUP BY A.customer, B.brand;
```

A m records

Customer	ProductID	Brand	Quant.
Abe	i9-12900	Intel	1
Abe	i7-12700	Intel	1
Abe	RTX3080	NVIDIA	1
Abe	660p	Intel	2
Bob	RyZen 5800G	AMD	1
Bob	980Pro	Samsung	1
Cindy	RTX3080	NVIDIA	1
Cindy	i7-12700	Intel	2
Cindy	660p	Intel	2
Diana	RyZen 5800G	AMD	1
Diana	RX5000	AMD	1
Diana	980Pro	Samsung	1

B n records

Brand	ProductID	Price
Intel	i9-12900	600
Intel	i7-12700	400
Intel	660p	100
AMD	RX5000	500
AMD	RyZen 5800G	200
NVIDIA	RTX3080	1200
Samsung	980Pro	100

Joined A*B

Customer	ProductID	Brand	Quant.	Price	Value
Abe	i9-12900	Intel	1	600	600
Abe	i7-12700	Intel	1	400	400
Abe	RTX3080	NVIDIA	1	1200	1200
Abe	660p	Intel	2	100	200
Bob	RyZen 5800G	AMD	1	400	400
Bob	980Pro	Samsung	1	100	100
Cindy	RTX3080	NVIDIA	1	1200	1200
Cindy	i7-12700	Intel	2	400	800
Cindy	660p	Intel	2	100	200
Diana	RyZen 5800G	AMD	1	400	400
Diana	RX5000	AMD	1	500	500
Diana	980Pro	Samsung	1	100	100

Compare and sum

Customer	Brand	Value
Abe	Intel	1200
Abe	NVIDIA	1200
Bob	AMD	400
Bob	Samsung	100
Cindy	Intel	1000
Cindy	NVIDIA	1200
Diana	AMD	900
Diana	Samsung	100

Memory copy

Compare & Element-wise Memory copy

Vector multiplications

$O(2 \times m \times n)$ $O(m \times n)$

Who is my most loyal customer?

```
SELECT A.customer, B.brand, SUM(A.Quantity * B.Price) AS value FROM A INNER JOIN B WHERE ON  
A.ProductID = B.ProductID GROUP BY A.customer, B.brand;
```

Customer	ProductID	Brand	Quant.
Abe	i9-12900	Intel	1
Abe	i7-12700	Intel	1
Abe	RTX3080	NVIDIA	1
Abe	660p	Intel	2
Bob	RyZen 5800G	AMD	1
Bob	980Pro	Samsung	1
Cindy	RTX3080	NVIDIA	1
Cindy	i7-12700	Intel	2
Cindy	660p	Intel	2
Diana	RyZen 5800G	AMD	1
Diana	RX5000	AMD	1
Diana	980Pro	Samsung	1

B

Brand	ProductID	Price
Intel	i9-12900	600
Intel	i7-12700	400
Intel	660p	100
AMD	RX5000	500
AMD	RyZen 5800G	200
NVIDIA	RTX3080	1200
Samsung	980Pro	100

O(n)

Customer	Intel	AMD	NVIDIA	Samsung
Abe	1200	0	1200	0
Bob	0	400	0	100
Cindy	1000	0	1200	0
Diana	0	900	0	100

Left

Customer/ ProductID	i9-1290	i7-1270	660	RX500	RyZen 5800G	RTX308	980Pr o
	0	0	p	0	0	0	0
Abe	1	1	2	0	0	1	0
Bob	0	0	0	0	1	0	1
Cindy	0	2	2	0	0	1	0
Diana	0	0	0	1	1	0	1

O(m × n × k)

Matrix multiplications

This could be done by 1 Tensor Core

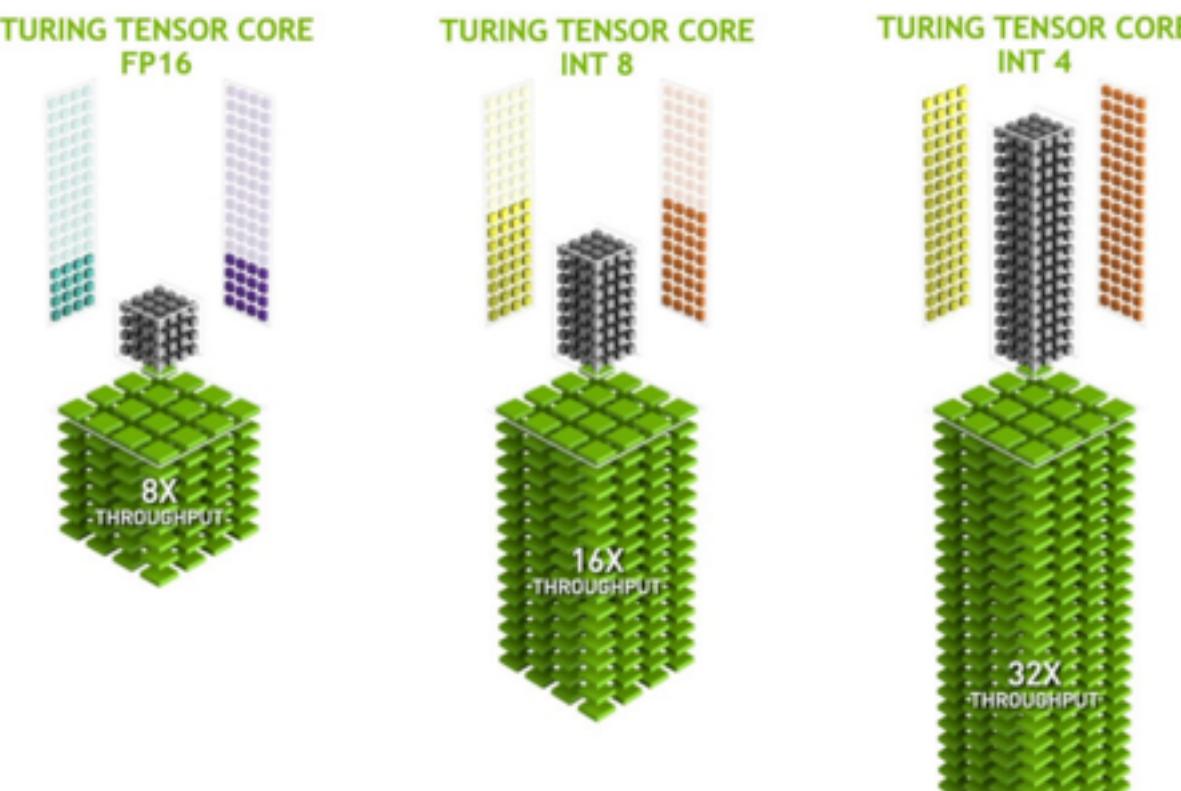
ProductID/ Brand	Intel	AMD	NVIDIA	Samsung
i9-12900	600	0	0	0
i7-12700	400	0	0	0
660p	100	0	0	0
RX5000	0	500	0	0
RyZen 5800G	0	400	0	0
RTX3080	0	0	1200	0
980Pro	0	0	0	100

O(m)

Memory copy

Recap: Efficiency of Tensor Cores

The NVIDIA Tesla T4 GPU includes 2,560 CUDA Cores and 320 Tensor Cores, delivering up to 130 TOPs (Tera Operations per second) of INT8 and up to 260 TOPS of INT4 inferencing performance (see Appendix A, *Turing TU104 GPU* for more Tesla T4 specifications). Compared to CPU-based inferencing, the Tesla T4, powered by the new Turing Tensor Cores, delivers up to 40X higher inference performance⁶ (see Figure 9).



Each tensor core operation performs 4x4x4 MMA in one cycle
~ 128 FP operations in conventional scalar models

FLOPS of 8K by 8K matrix multiplications =

$$8192 \times 8192 \times 8192 \times 2$$

$$\frac{8192 \times 8192 \times 8192 \times 2}{2560} = 429496730 \text{ CUDA core cycles}$$

$$\frac{8192 \times 8192 \times 8192 \times 2}{320 \times 128} = 26843546 \text{ Tensor core cycles}$$



Tensor cores make matrix processing 16x faster

TCU friendly queries

Two-way natural join

```
SELECT A.Val, B.Val FROM A, B  
WHERE A.ID = B.ID;
```

Multi-way joins

```
SELECT A.Val, B.Val, C.Val FROM A,B,C  
WHERE  
A.ID_1 = B.ID_1 AND B.ID_2 = C.ID_2;
```

Group-by aggregates
over joins

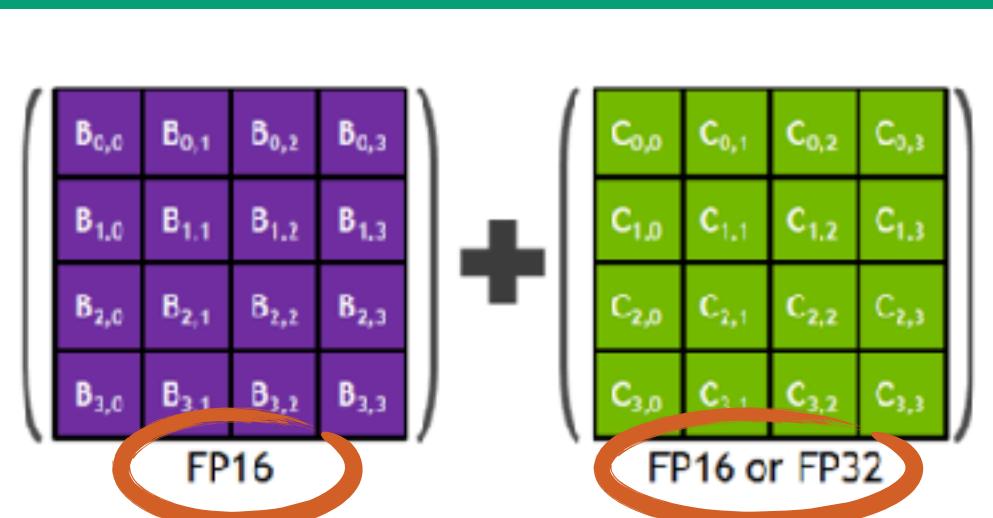
```
SELECT SUM(A.Val), B.Val FROM A, B  
WHERE A.ID = B.ID  
GROUP BY B.Val;
```

Aggregate join
without group-by

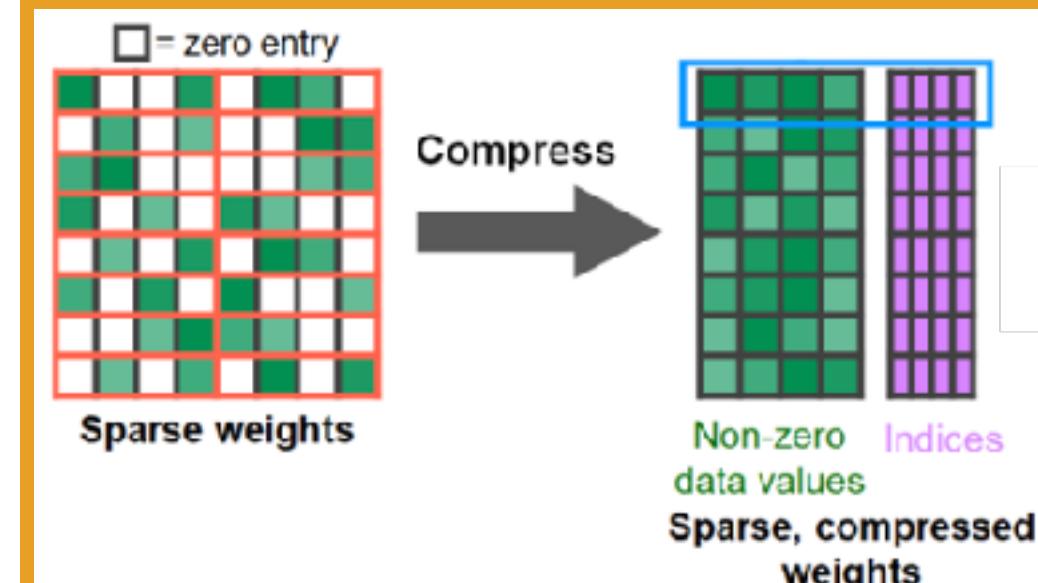
```
SELECT SUM(A.Val * B.Val)  
FROM A, B  
WHERE A.ID = B.ID;
```

Challenges

Limited precision



Data-dependent code optimization

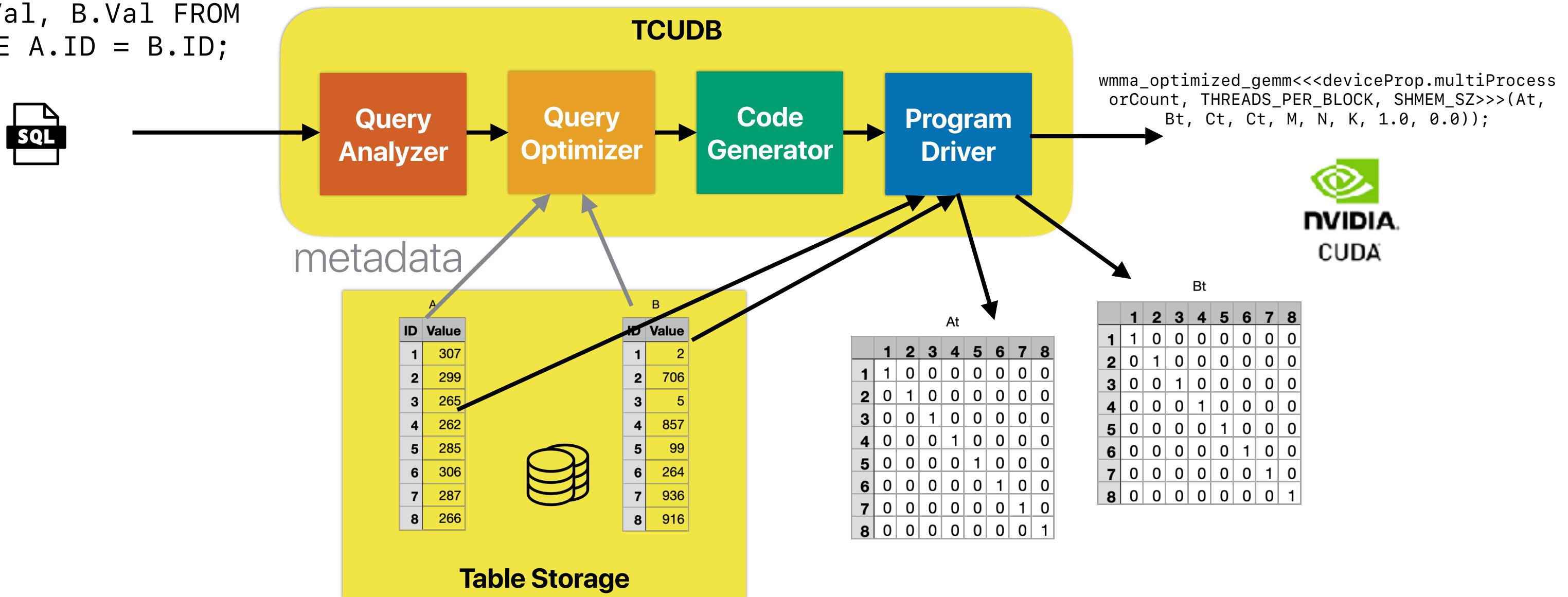


Limited device memory

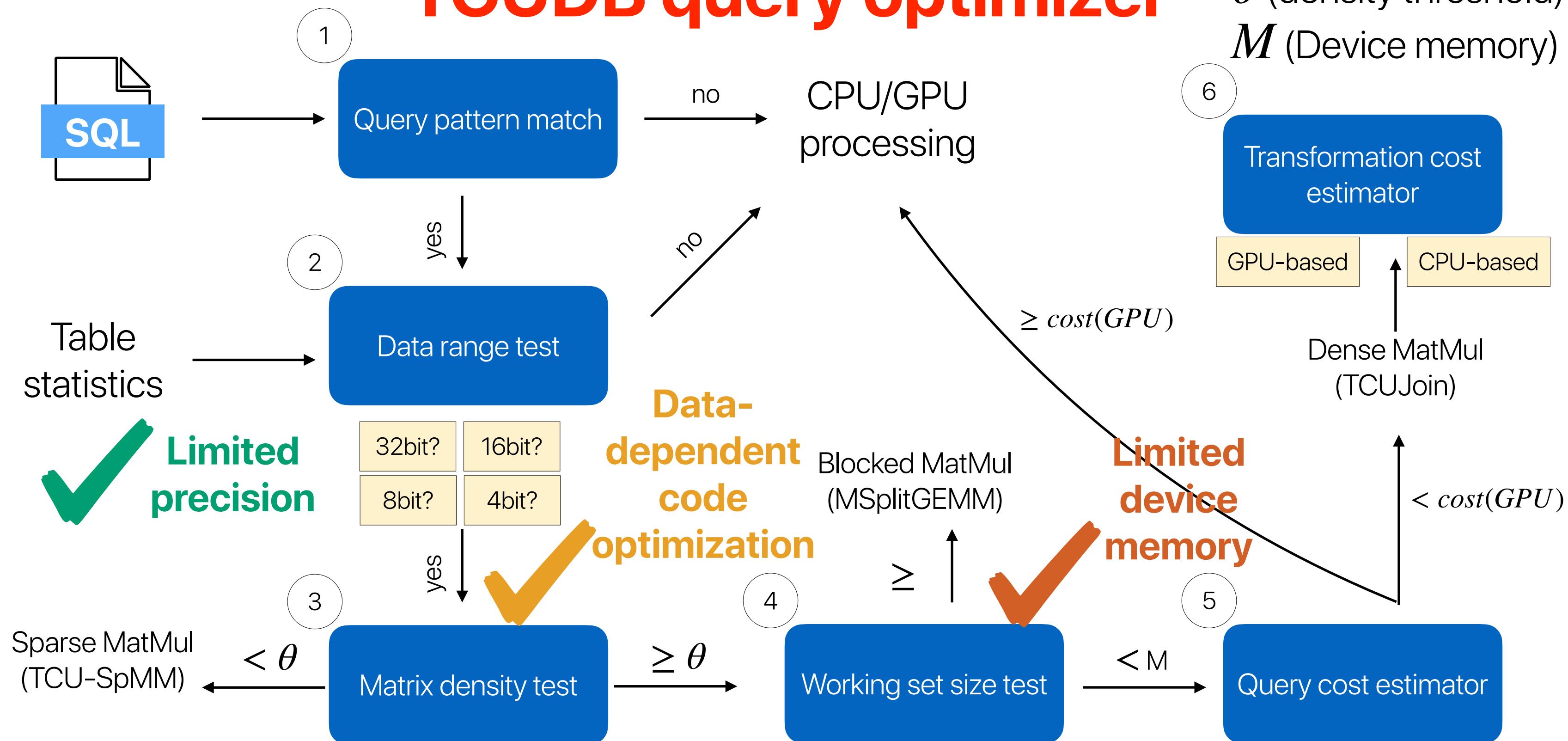
Peak TF32 Tensor TFLOPS ¹	NA
Peak INT8 Tensor TOPS ¹	261
Peak INT4 Tensor TOPS ¹	522
Frame Buffer Memory Size and Type	24576MB GDDR6
Memory Interface	384-bit
Memory Clock (Data Rate)	14 Gbps
Memory Bandwidth (GB/sec)	672 GB/sec
ROPs	95
Pixel Fill-rate (Gigapixels/sec)	169.9

The TCUDB system architecture

```
SELECT A.Val, B.Val FROM  
A, B WHERE A.ID = B.ID;
```

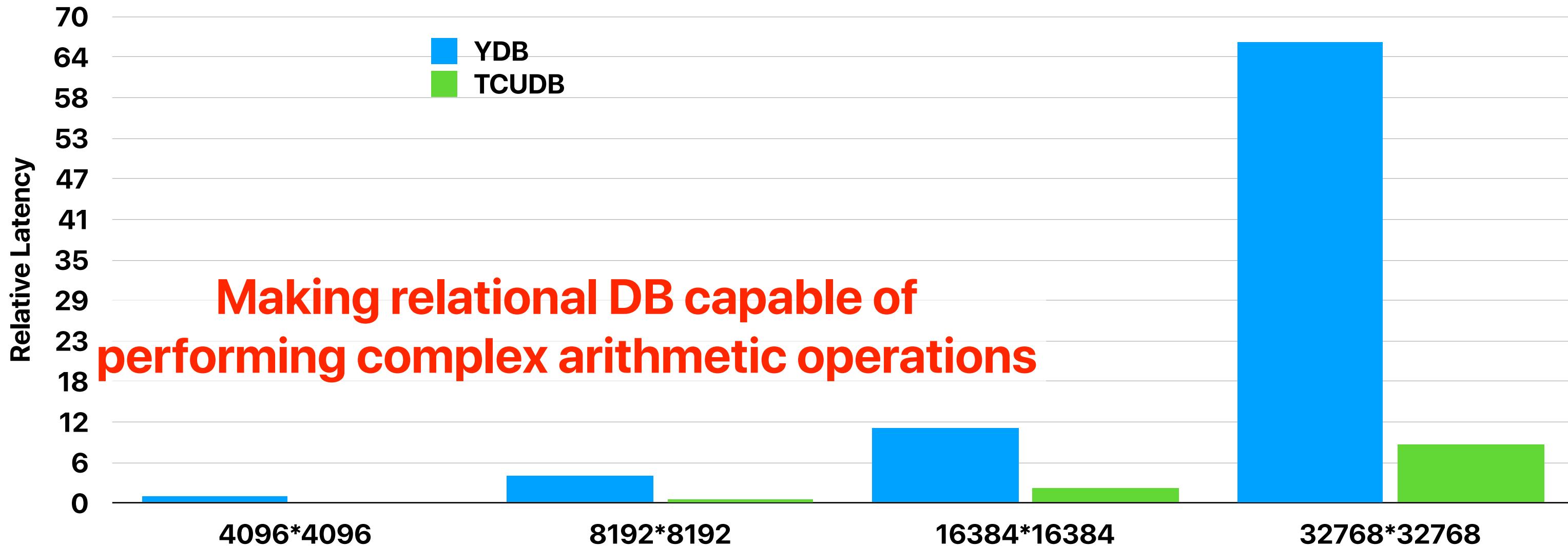


TCUDB query optimizer



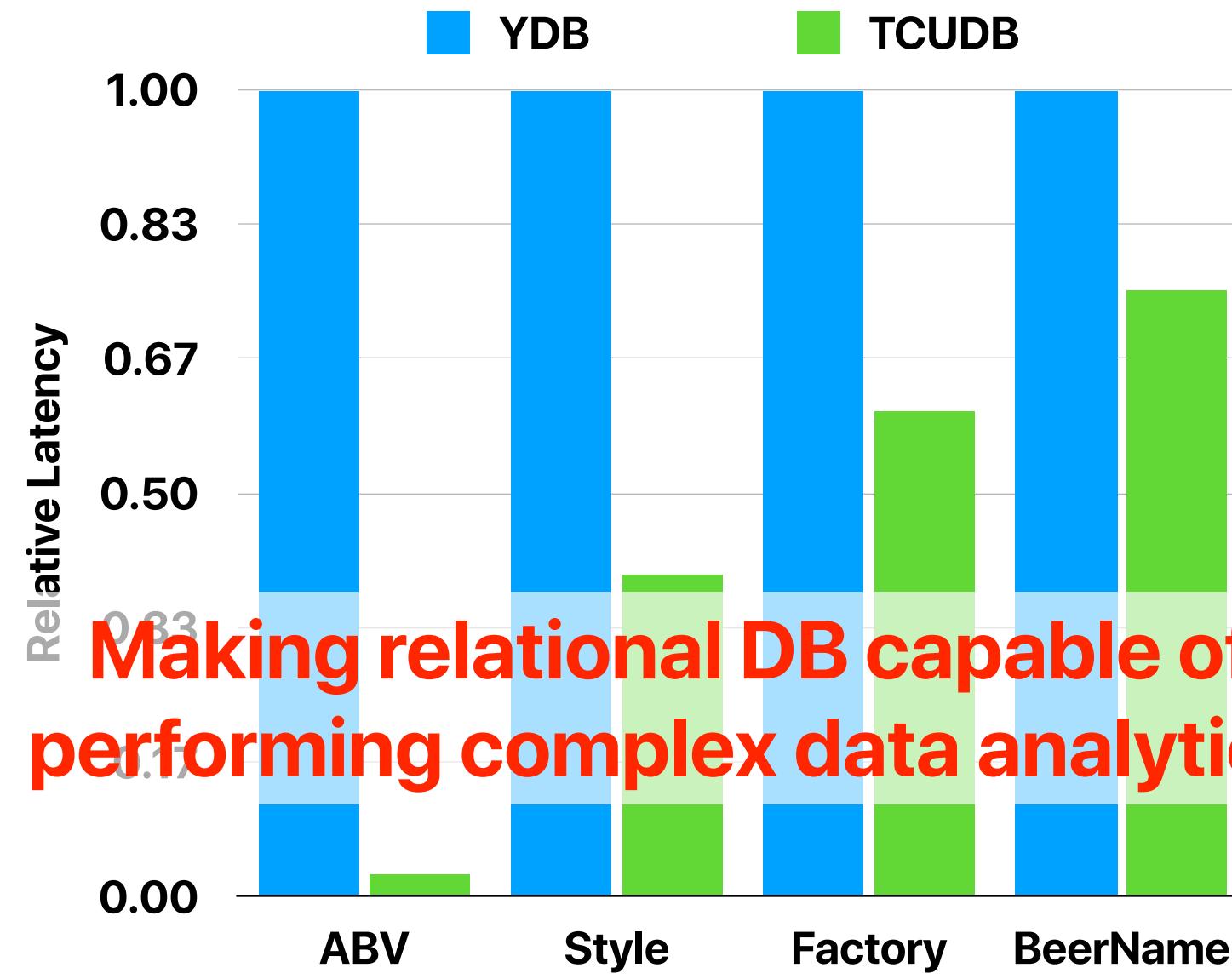
Matrix Multiplication Query

- `SELECT A.col_num, B.row_num, SUM(A.val * B.val) as res FROM A, B WHERE A.row_num = B.col_num GROUP BY A.col_num, B.row_num;`

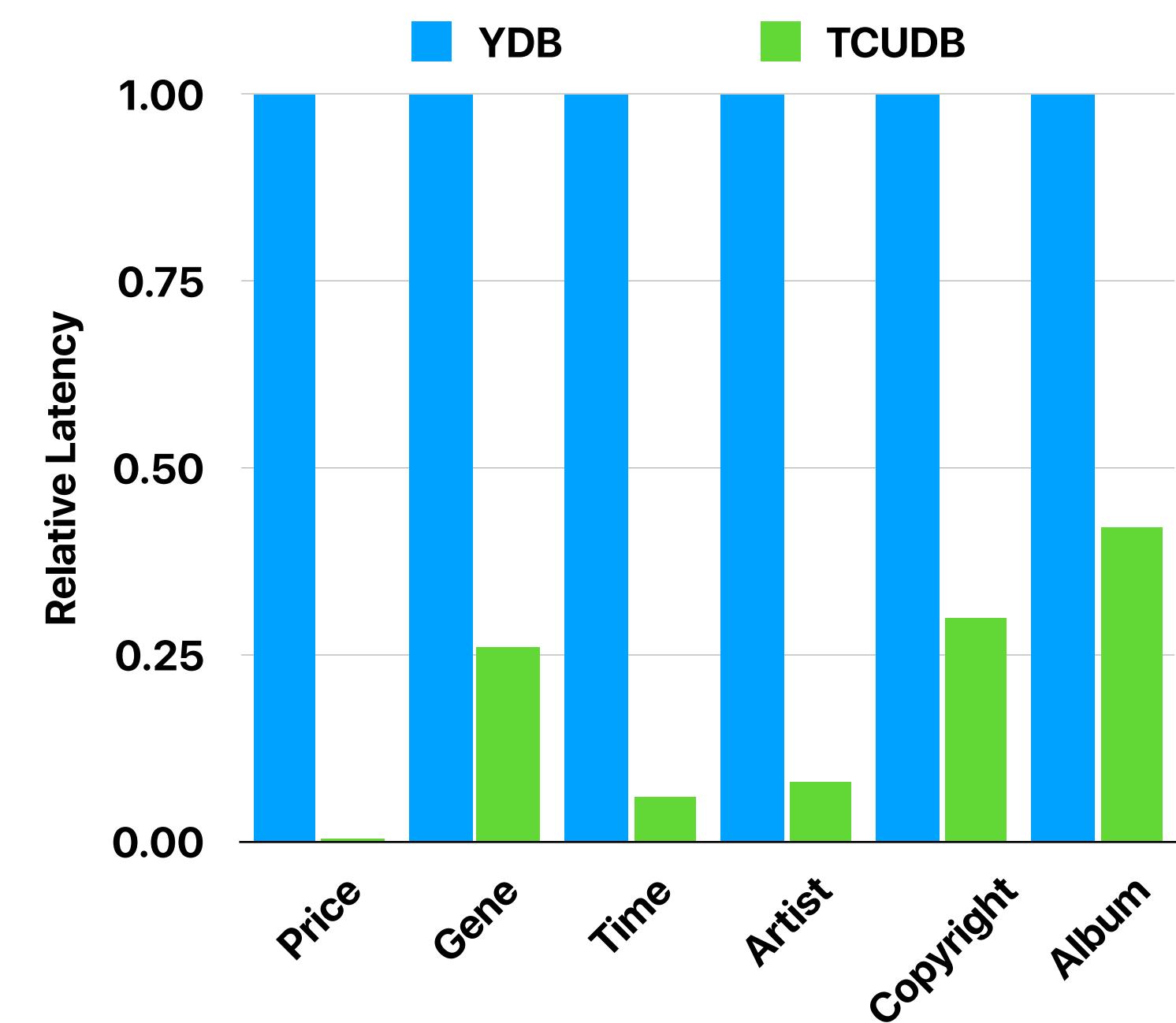


Entity matching

- `SELECT TABLE_A.ID, TABLE_A.BEER_NAME, TABLE_B.ID, TABLE_B.BEER_NAME FROM TABLE_A, TABLE_B WHERE TABLE_A.ABV = TABLE_B.ABV;`
- `SELECT TABLE_A.ID, TABLE_A.SONG, TABLE_B.ID, TABLE_B.SONG FROM TABLE_A, TABLE_B WHERE TABLE_A.ARTIST = TABLE_B.ARTIST;`

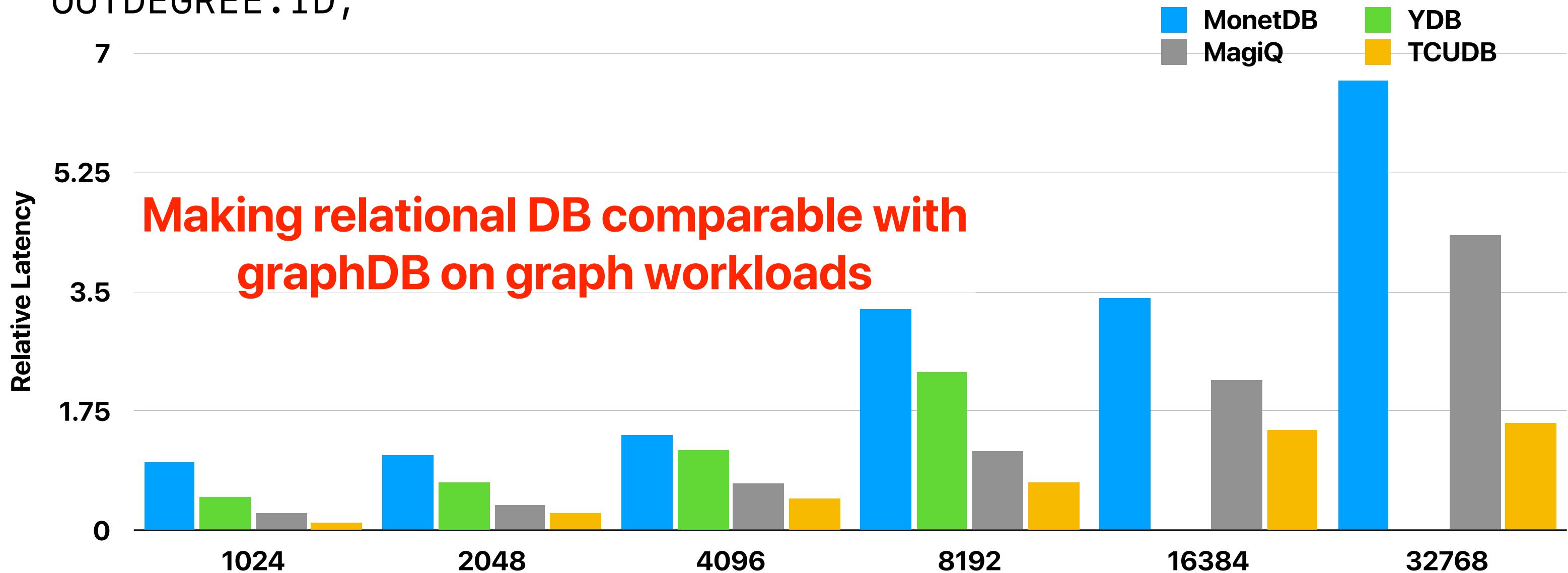


Making relational DB capable of performing complex data analytics



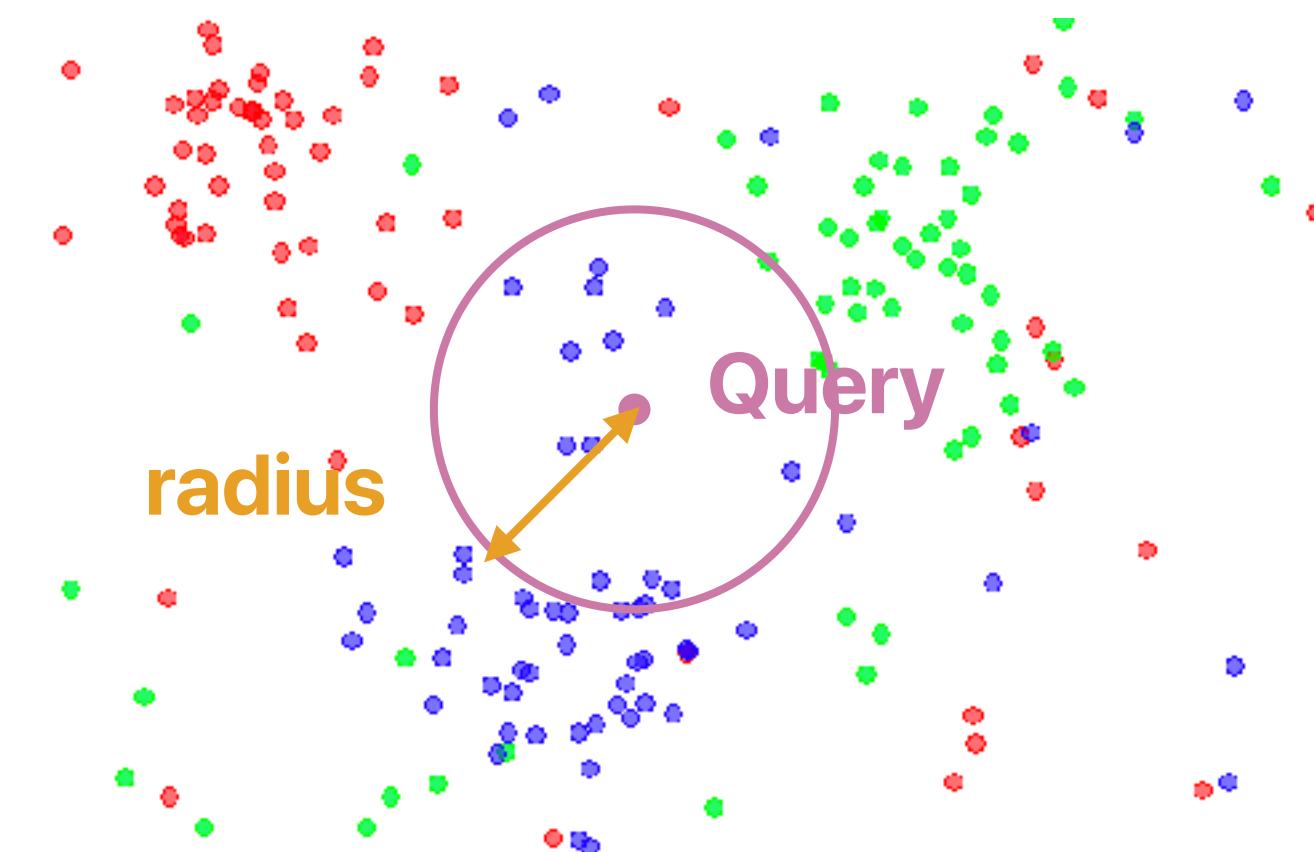
PageRank

- `SELECT SUM(@alpha * PAGERANK.rank / OUTDEGREE.DEGREE) + (1-@alpha)/@num_node FROM PAGERANK, OUTDEGREE WHERE PAGERANK.ID = OUTDEGREE.ID;`



Nearest neighbor search

- Fixed-radius search (Range search) — returns all the neighbors within a fixed radius.
- K nearest neighbor search (kNN) — returns the nearest K neighbors of a query



Can you map the NN problem as a
ray tracing problem?

NN Search as Ray Tracing?

RTNN: accelerating neighbor search using hardware ray tracing

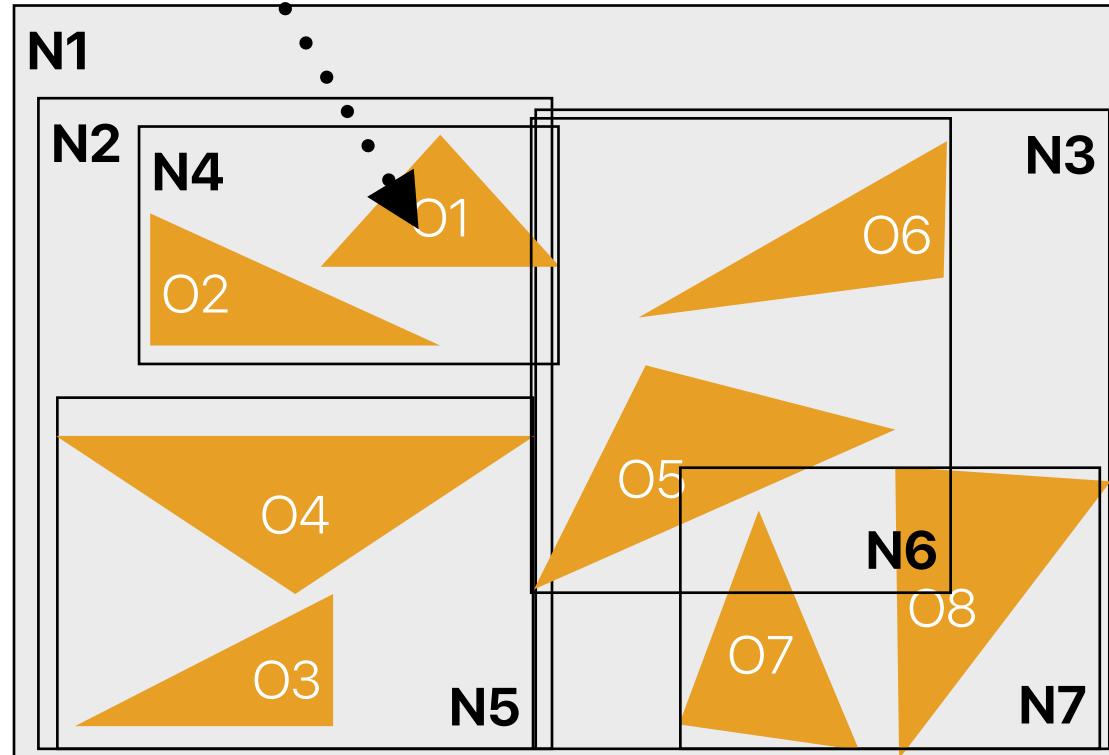
Yuhao Zhu

University of Rochester

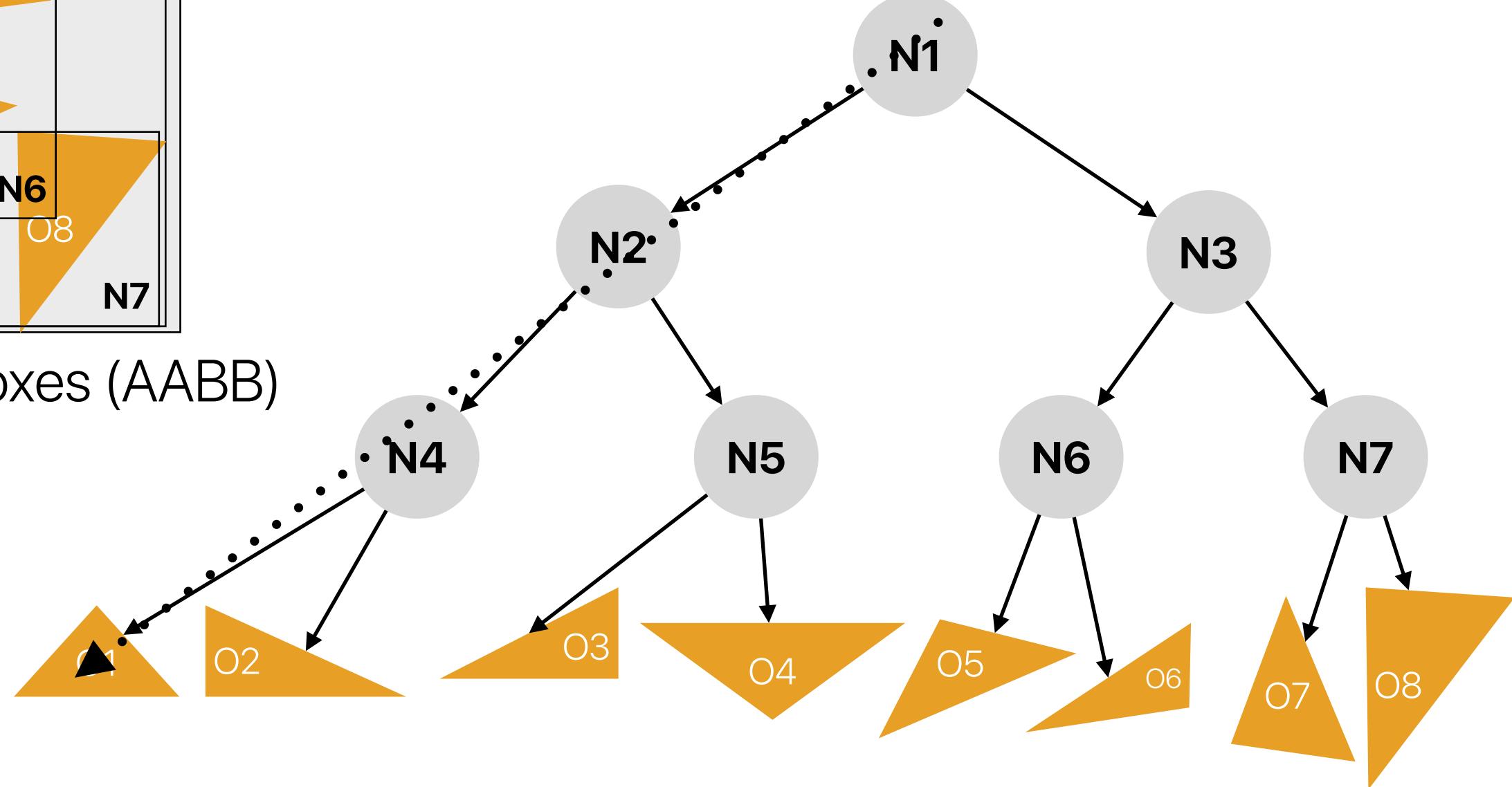
In Proceedings of the 27th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '22)

<https://github.com/horizon-research/rtnn>

Bounding volume hierarchy



axis-aligned bounding boxes (AABB)



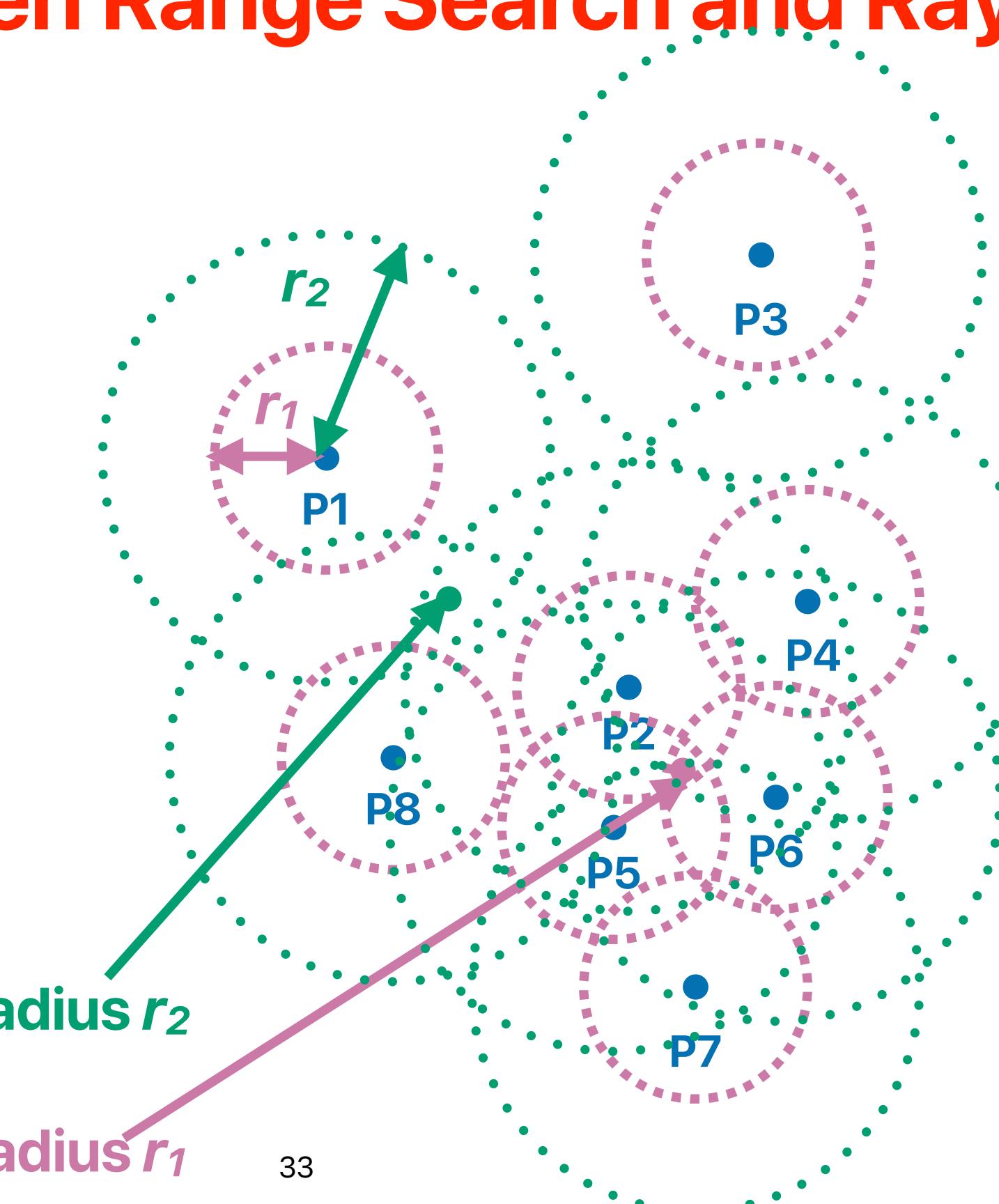
Similarities between Range Search and Ray Tracing

Answer: P1, P2, P8

Q_2 with radius r_2

Answer: P2, P5, P6

Q_1 with radius r_1



Similarities between kNN and Ray Tracing

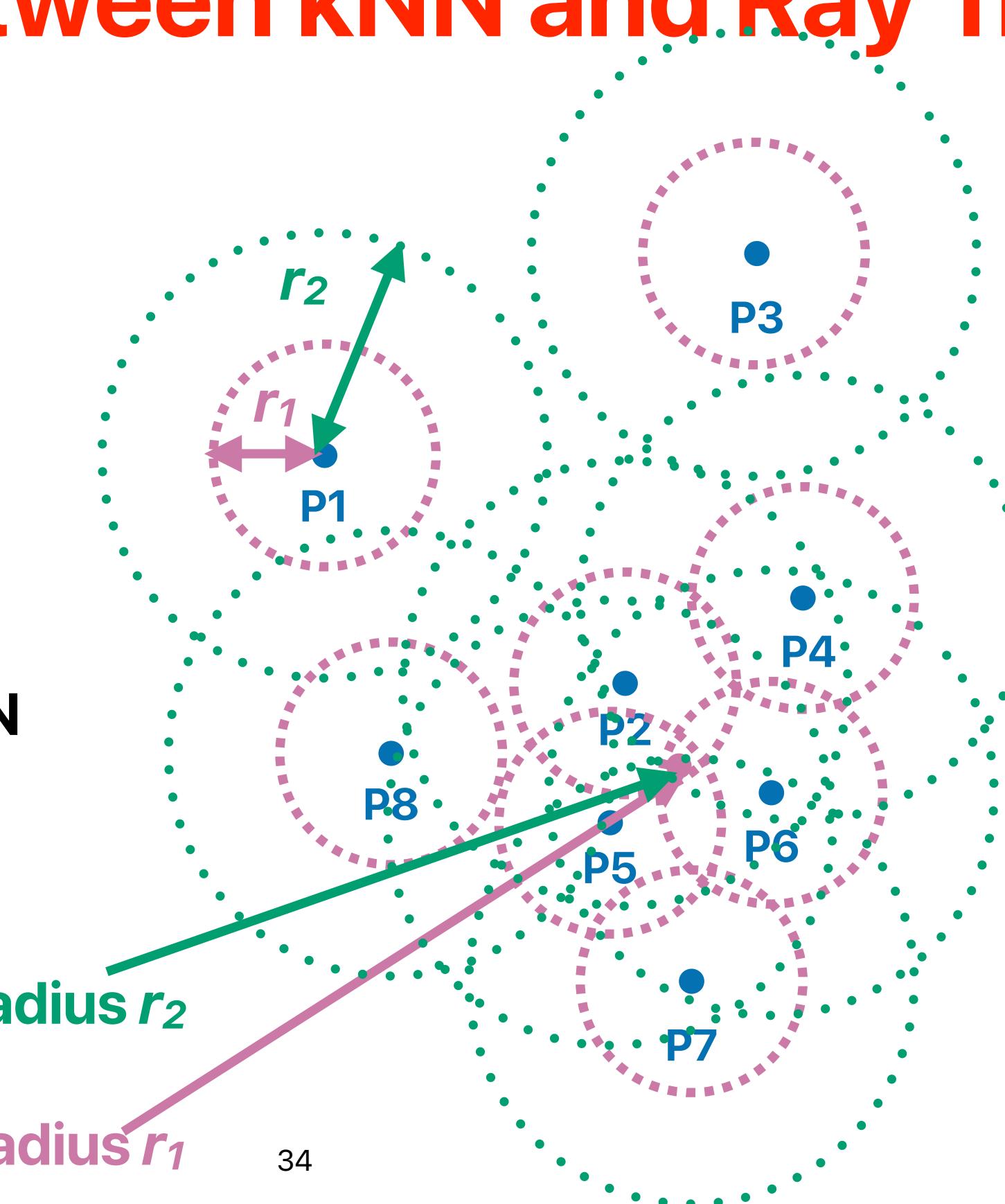
Assume we want to find "5" NN

Answer: P2, P4, P5, P6, P7

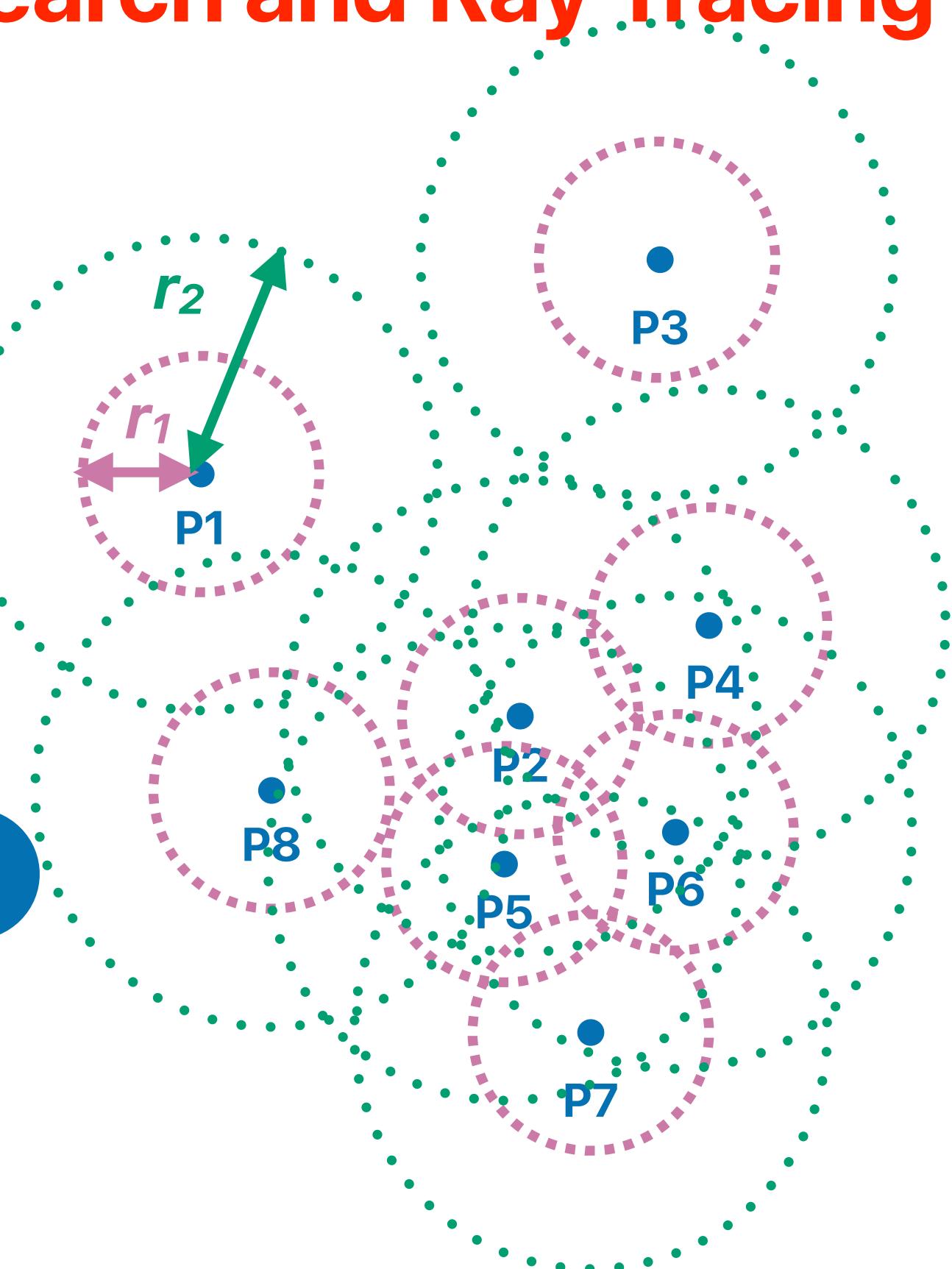
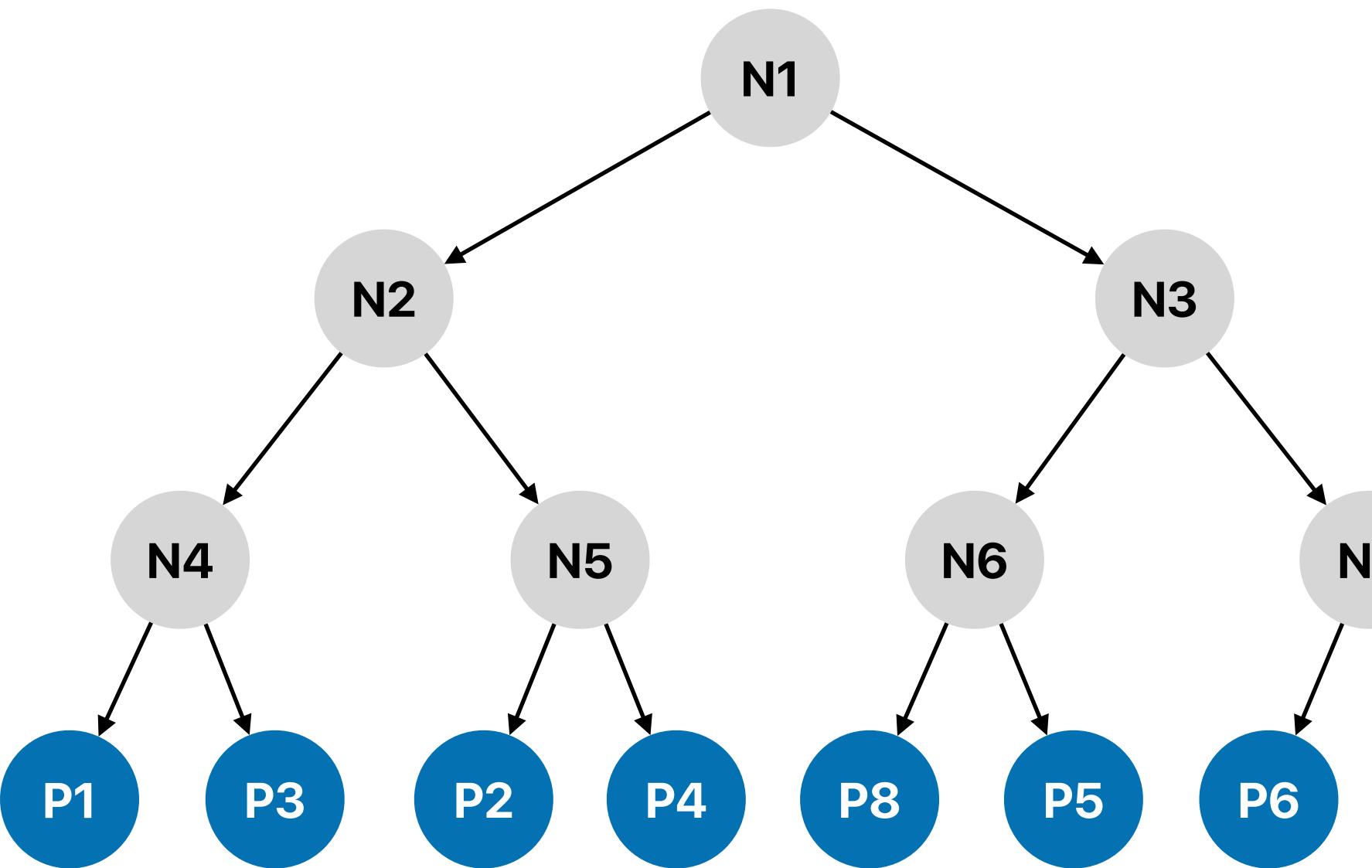
Q₁ with radius r₂

Answer: P2, P5, P6

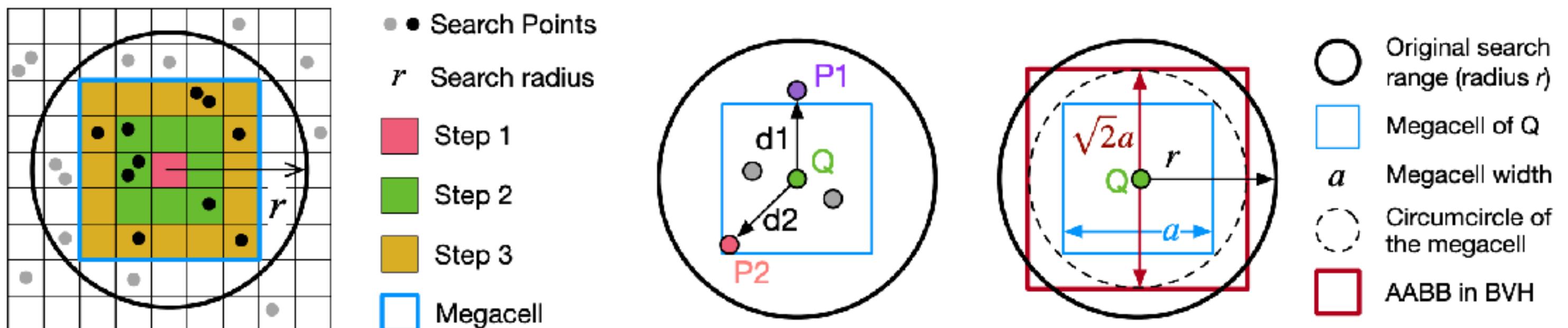
Q₁ with radius r₁



Similarities between Range Search and Ray Tracing



Ideas behind RTNN



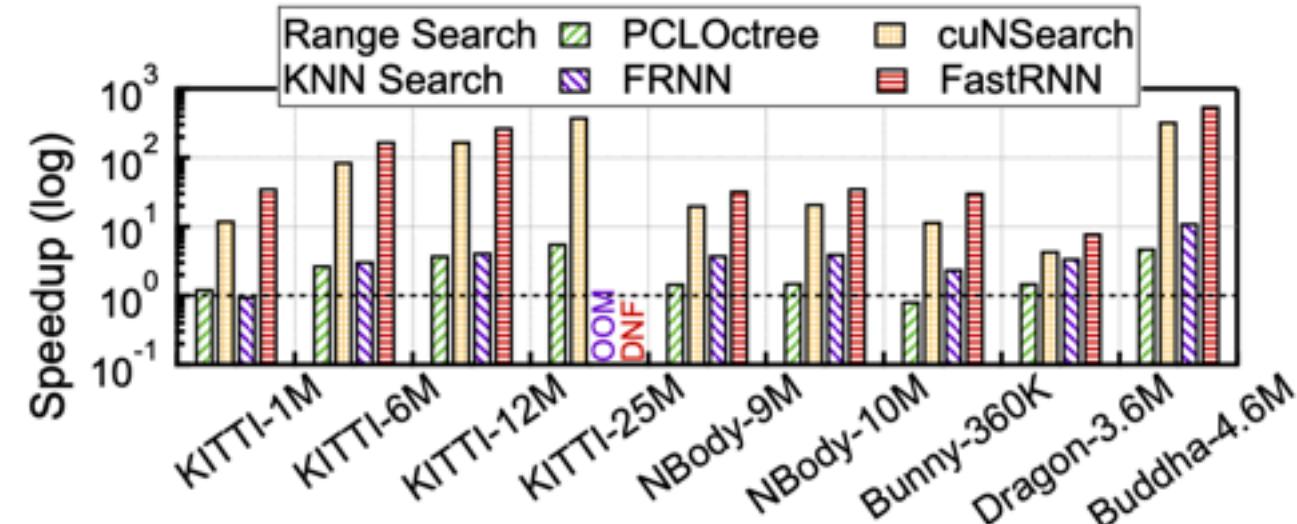
(a) Calculating the megacell of a query by incrementally growing from the cell that contains the query. The growth stops when either the sphere boundary is reached or at least K neighbors are found.

(b) P_1 is outside of Q 's megacell that contains 3 points, but is among the 3 nearest neighbors of Q (as $d_1 < d_2$).

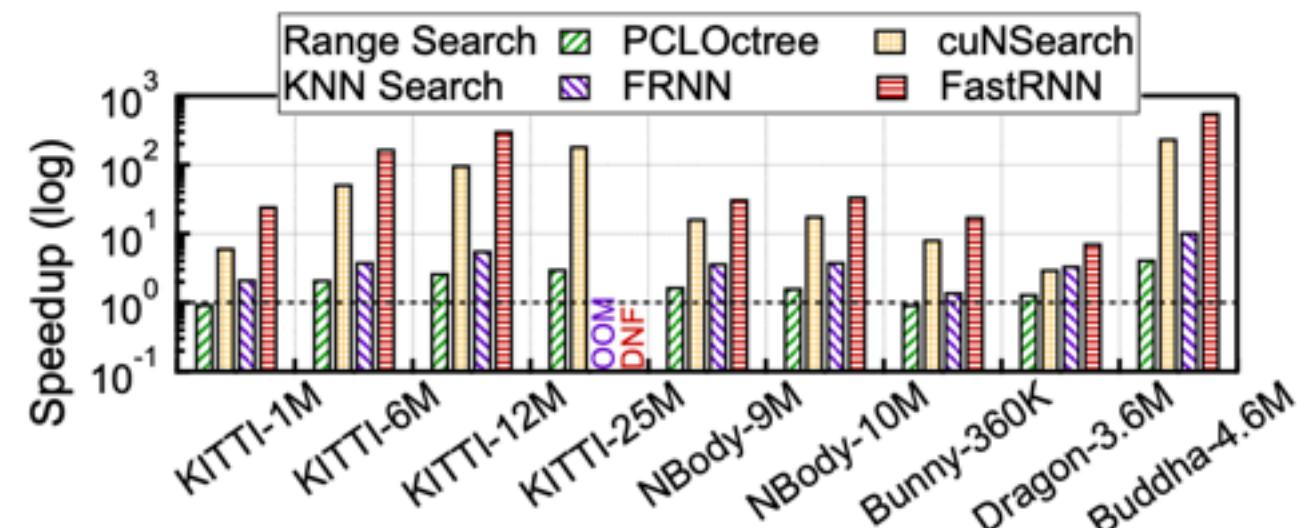
(c) For KNN search, the AABB must circumscribe the sphere that circumscribes the megacell. The AABB width is $\sqrt{2}a$ for 2D search (illustrated here) and $\sqrt{3}a$ for 3D search (as the AABB is 3D and the circumcircle becomes the circumsphere), where a is the megacell width.

Fig. 10. Determining the megacell and AABB size used for neighbor search.

Performance



(a) Speedups on RTX 2080.



(b) Speedups on RTX 2080Ti.

Fig. 11. Speedup of RTNN over the baselines (log-scale). OOM denotes that the baseline ran out of memory; DNF denotes that the baseline did not finish within the time that would have given RTNN a 1,000 speedup.

Summary of General-purpose computing on “something”

- Map your problem as a problem that the target hardware can solve
 - Machine learning models (NPU)
 - Machine learning mathematical operations (GP Edge TPU)
 - Matrix multiplications (TCUDB)
 - Ray Tracing (RTNN)
- Due to the domain specific nature of accelerators, we have to program in a more domain specific way, “currently”.
- Not necessarily the most performant, but typically more energy-efficient

Roogle Claimed ideas

- Claimed ideas
 - Setup and benchmarking SmartSSDs
 - Parallel processing using matrix parallelism
 - Benchmarking fault-tolerant mechanisms in machine learning
 - Benchmarking and optimizing multimedia machine learning applications on heterogeneous computing resources
 - LLM on Edge TPUs/Embedded Systems
- Contact me if you've confirmed your project ideas but not seeing it in the slide — also send me your team member names
- Make an appointment through Google Calendar if you need consultant on projects

Roogle Other interesting ideas

- Can we use RT cores to accelerate other problems beyond identified ones?
- Can you improve the efficiency of existing problems on Tensor Cores (matrix decompositions, security algorithms)?
- Can you use other accelerated libraries (e.g., cuDNN, cuFFT) for applications beyond their domains?
- Can you think about some data representations that can make important compute kernels (e.g., machine learning algorithms, matrix multiplications) more efficient or secure?

Roogle Project Presentations

- Make an appointment on 2/1 and 2/6 through the Google Calendar
- 12 minute presentation with 3 minute Q & A
- Why & what & how!!! — considering you're giving a presentation at Apple's keynote
 - 6-minute why — why should everyone care about this problem? Why is this still a problem?
 - 4-minute what — what are you proposing in this project to address the problem?
 - 2-minute how — expected platforms/engineering efforts, milestones and workload distribution among members
 - Please reference this article to make a good presentation <https://cseweb.ucsd.edu/~swanson/GivingTalks.html>

Electrical Computer Science Engineering

277

つづく

