

# Speeding up Query Execution

Andres Freund

PostgreSQL Developer & Committer  
Citus Data – citusdata.com - @citusdata

<http://anarazel.de/talks/pgday-nordic-2017-03-21/jit.pdf>

# Motivation

Scalability  
Efficiency  
Performance  
Concurrency

OLTP  
VS  
OLAP  
VS  
HTAP

Online Transactional Processing  
VS  
Online Analytical Processing  
VS  
Hybrid Transactional / Analytical Processing

# Amdahl's Law

Two independent parts **A** **B**

Original process



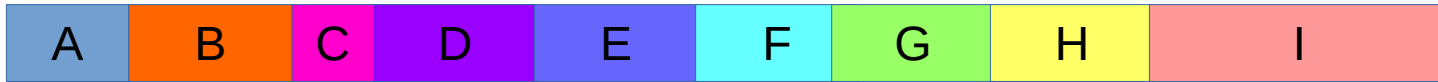
Make **B** 5x faster



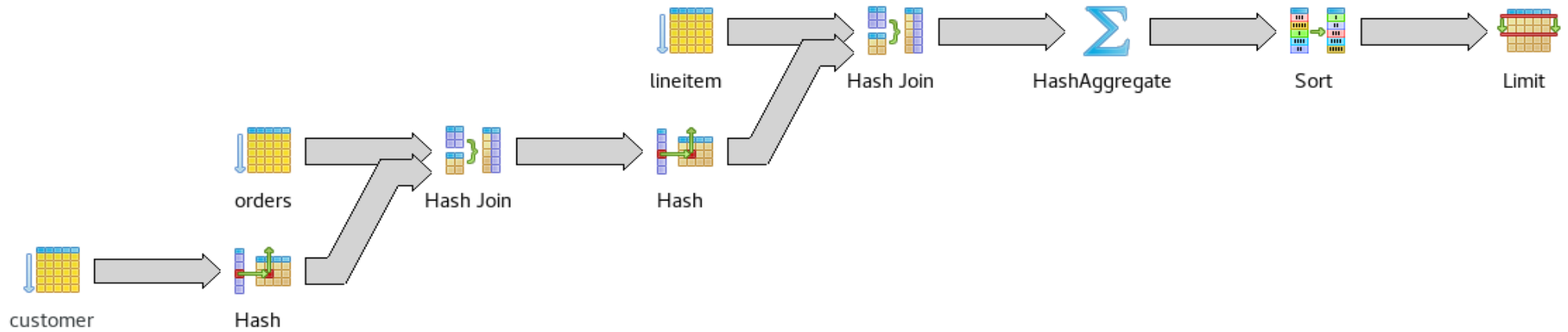
Make **A** 2x faster



<https://commons.wikimedia.org/wiki/File:Optimizing-different-parts.svg>



# Batch Tuple Processing

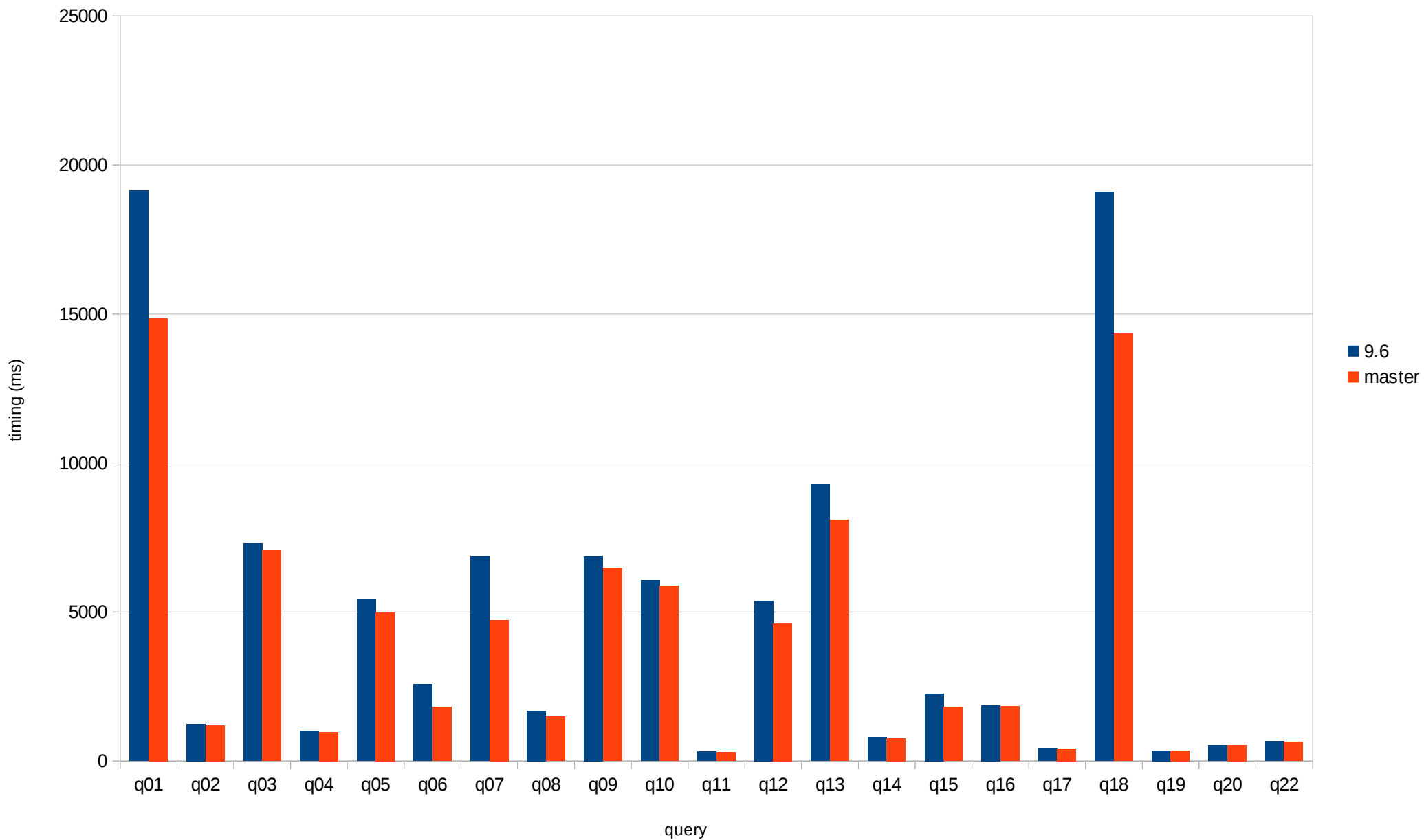


```
SELECT
    l_returnflag,
    l_linestatus,
    sum(l_quantity) as sum_qty,
    sum(l_extendedprice) as sum_base_price,
    sum(l_extendedprice * (1 - l_discount)) as sum_disc_price,
    sum(l_extendedprice * (1 - l_discount) * (1 + l_tax)) as sum_charge,
    avg(l_quantity) as avg_qty,
    avg(l_extendedprice) as avg_price,
    avg(l_discount) as avg_disc,
    count(*) as count_order
FROM
    lineitem
WHERE
    l_shipdate <= date '1998-12-01' - interval '74 days'
GROUP BY
    l_returnflag,
    l_linestatus
ORDER BY
    l_returnflag,
    l_linestatus;
```

```
Sort (cost=2153027.37..2153027.39 rows=6 width=68)
    (actual time=19742.591..19742.591 rows=4 loops=1)
    Sort Key: l_returnflag, l_linestatus
    Sort Method: quicksort Memory: 25kB
-> HashAggregate (cost=2153027.19..2153027.29 rows=6 width=68)
    (actual time=19742.557..19742.559 rows=4 loops=1)
    Group Key: l_returnflag, l_linestatus
-> Seq Scan on lineitem (cost= rows=29727516 width=36)
    (actual time=0.012..4514.186 rows=29713242 loops=1)
    Filter: (l_shipdate <= '1998-09-18 00:00:00'::timestamp without time zone)
    Rows Removed by Filter: 286553
Planning time: 0.326 ms
Execution time: 19742.705 ms
```

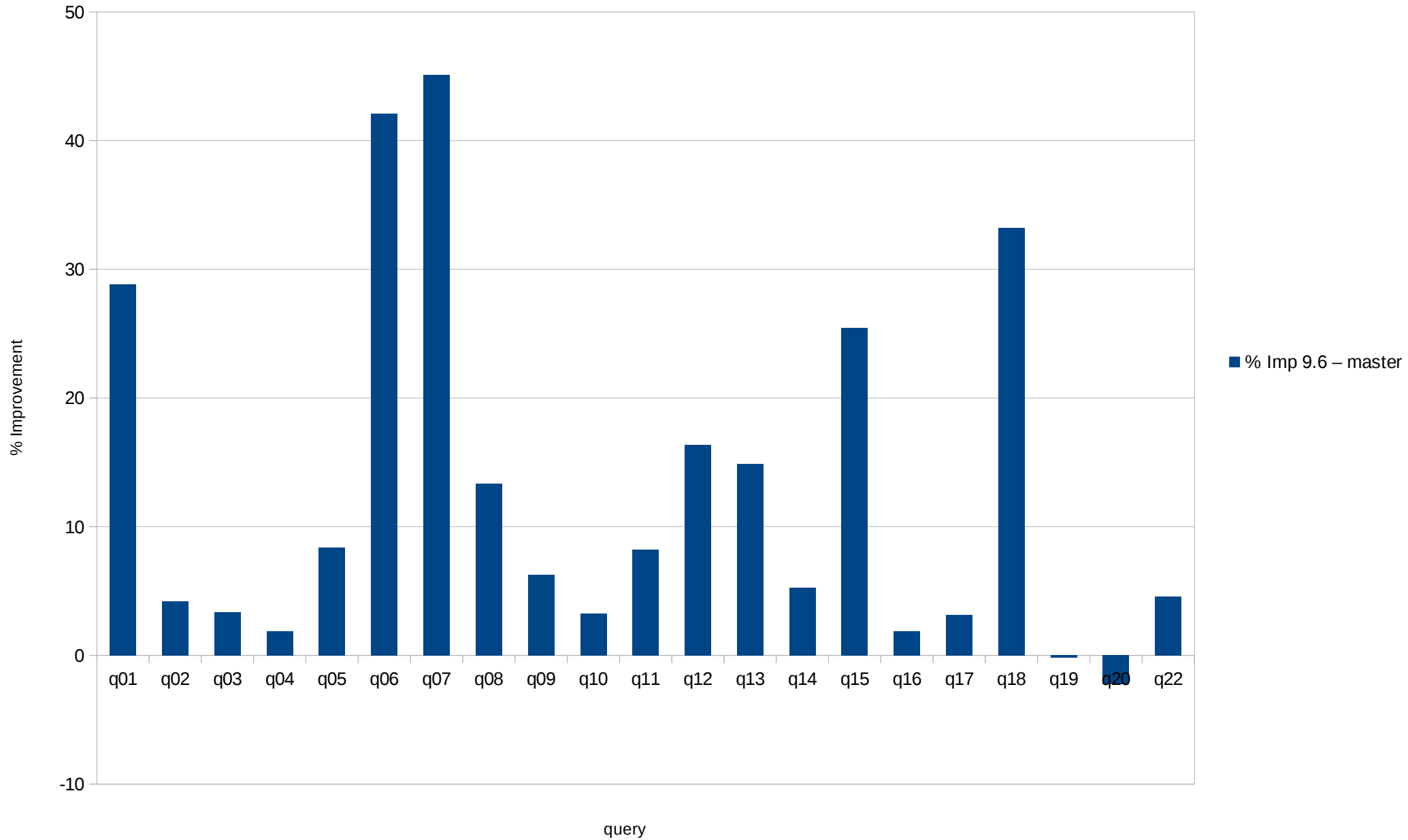
# TPCH Timing

Non-Parallized, Scale 5



# TPCH Improvements

Non-Parallelized, Scale 5



# Expression Evaluation

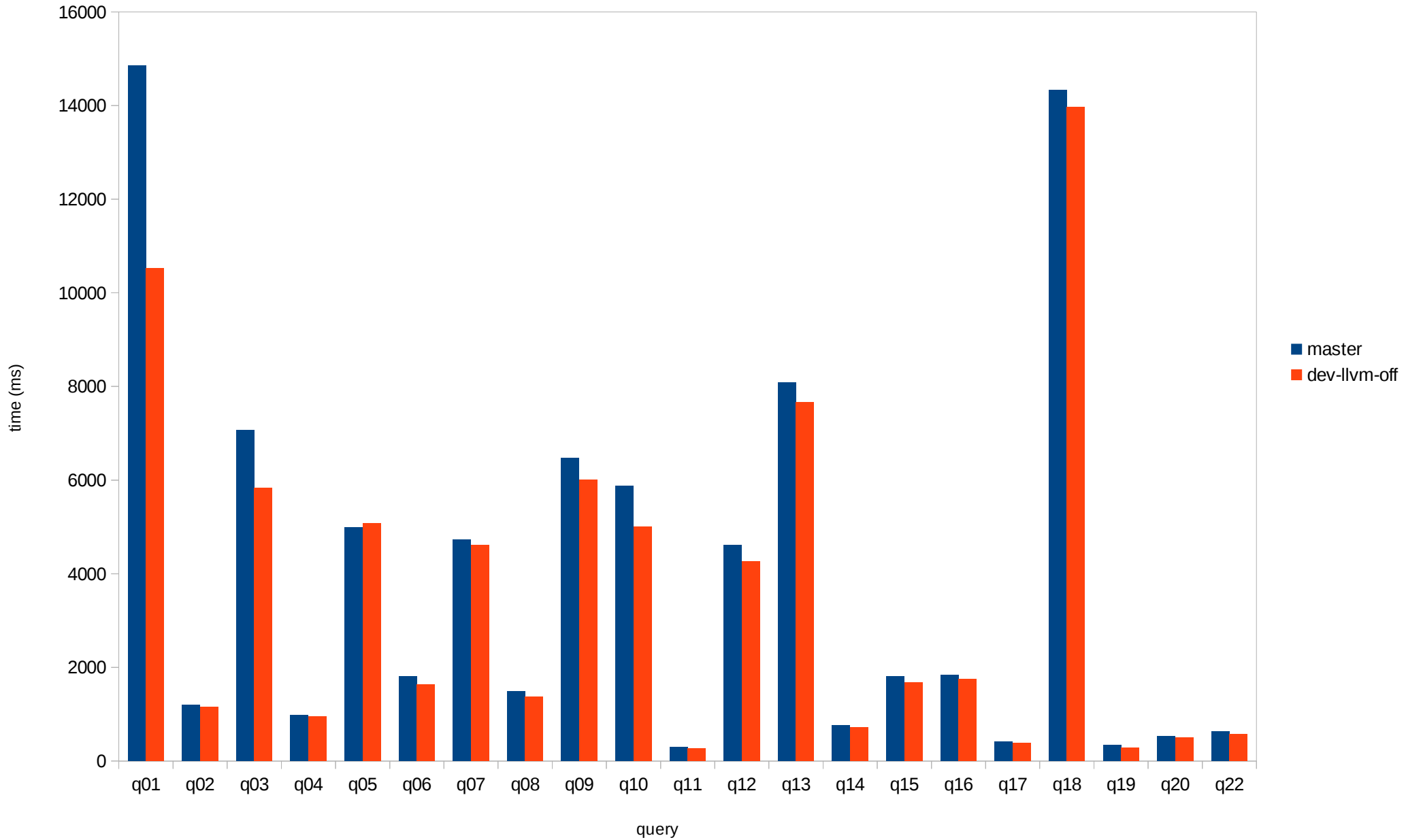
- SELECT list, of, columns,  $a + b$  as expr
- WHERE clauses
- GROUP BY clauses
- ...

# Expression Evaluation Methodology

- Old:
  - Tree Walk
  - Callback Based
  - Recursive
- New
  - “Byte code” generated internally
  - No / Fewer Callbacks
  - “Direct Threaded”
  - Non-Recursive

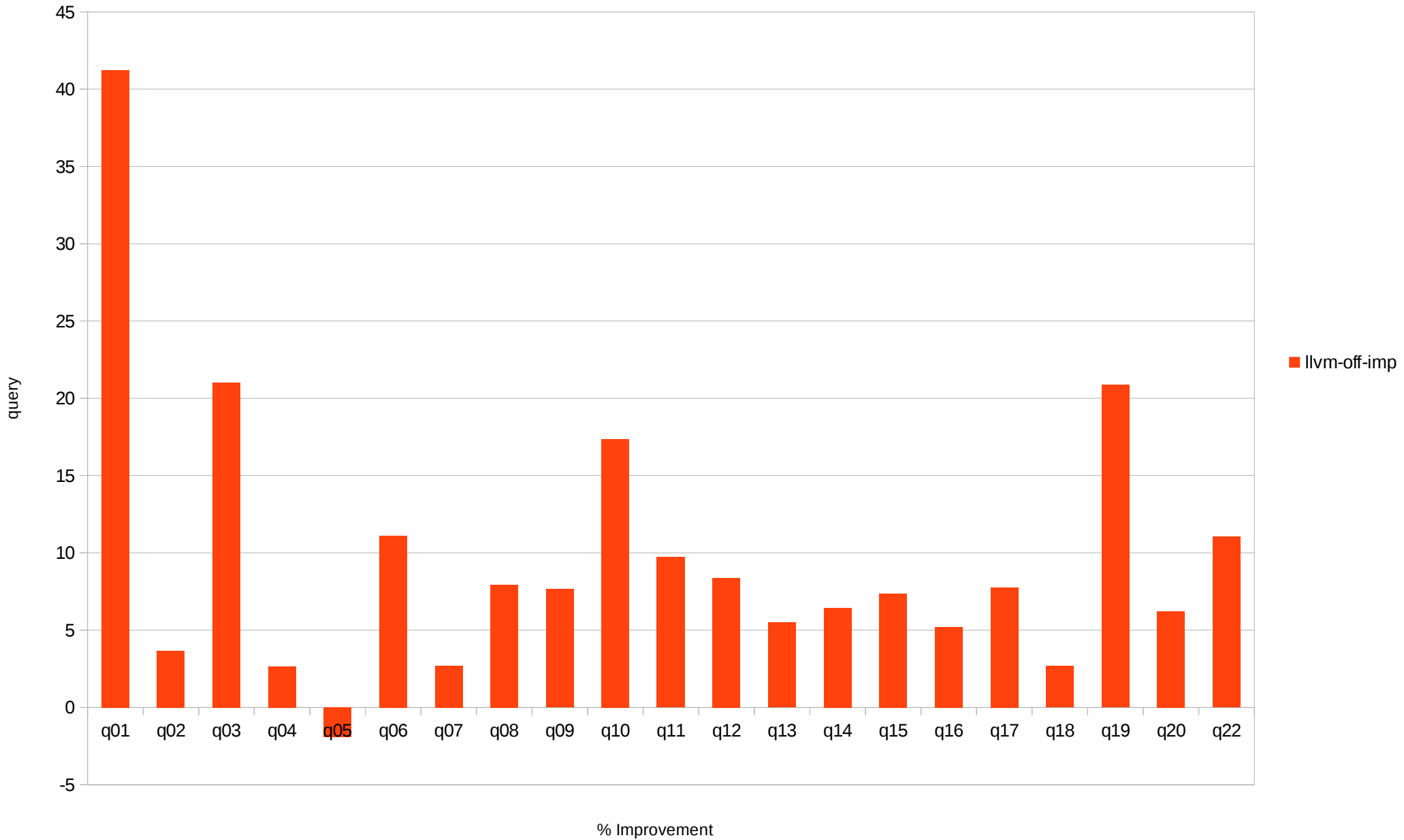
# TPCH Timings

Not Parallelized, Scale 5



# TPCH Improvement

Not Parallelized, Scale 5

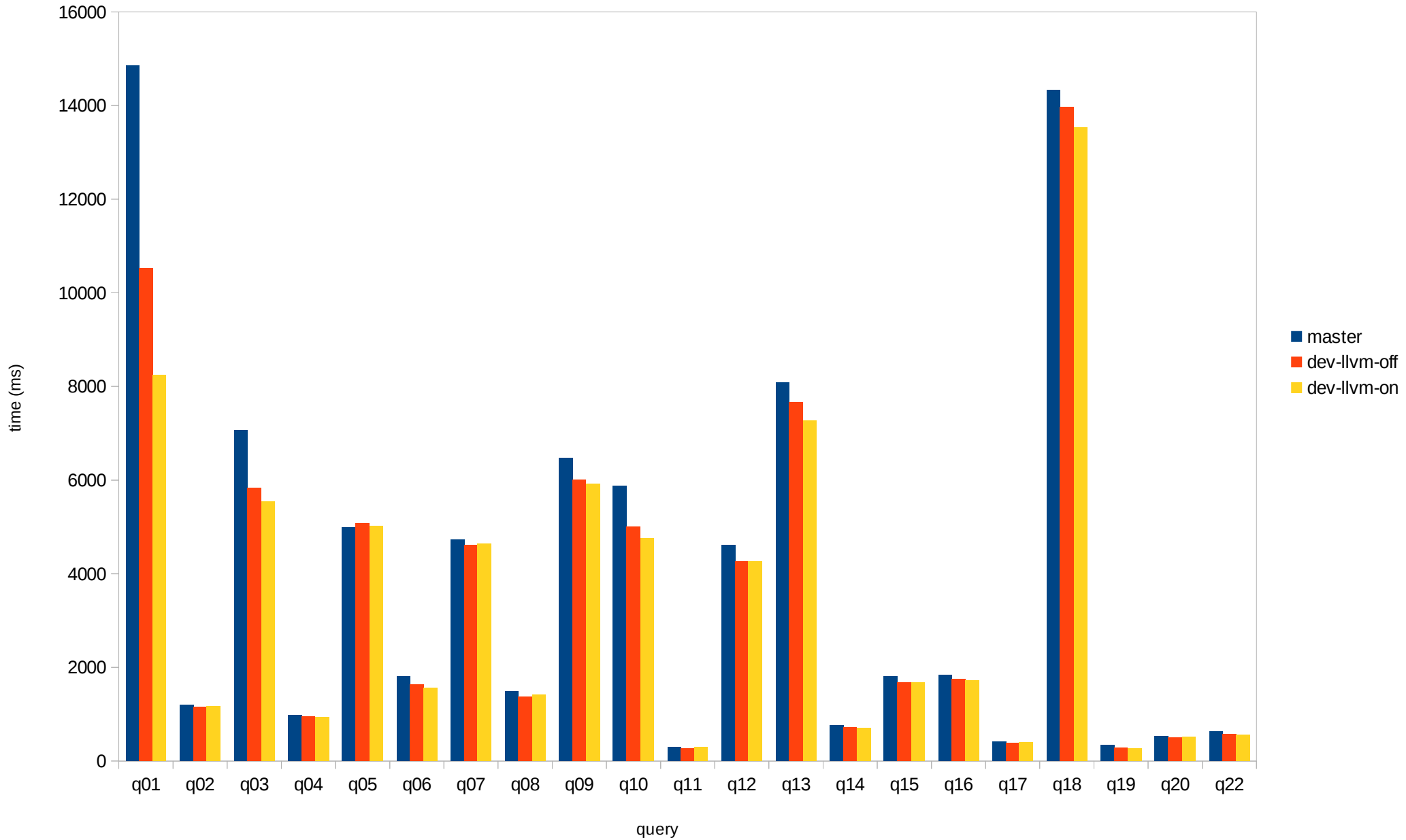


# Just In Time Compilation

- Interpretation has a lot of “jumps”
- Interpretation calls a lot of “unknown functions”
- Native Code doesn't have those Issues
- Only do so when beneficial
  - Generating a native function is expensive (~0.5-5ms)
- Can be used in a lot of places

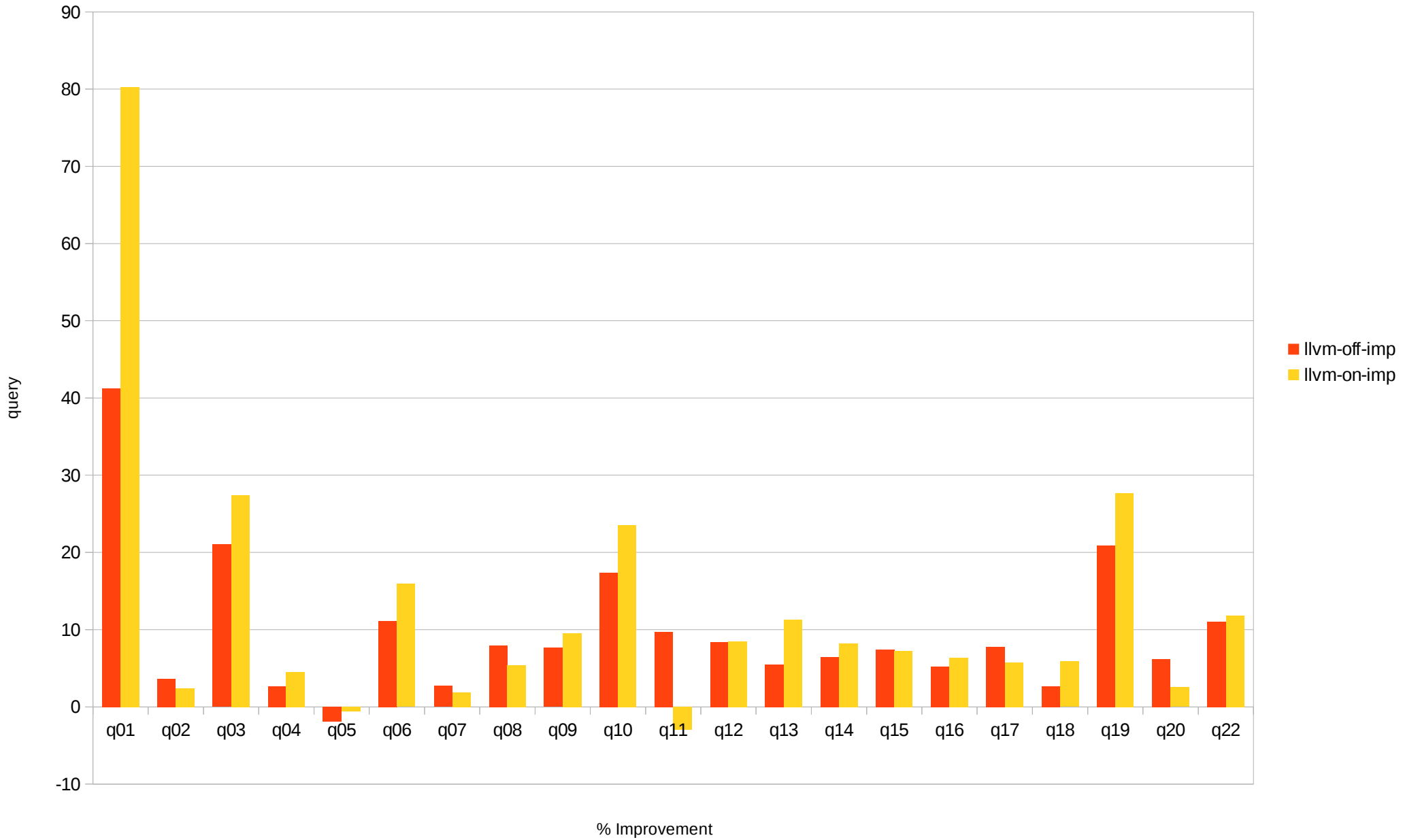
# TPCH Timings

Not Parallelized, Scale 5



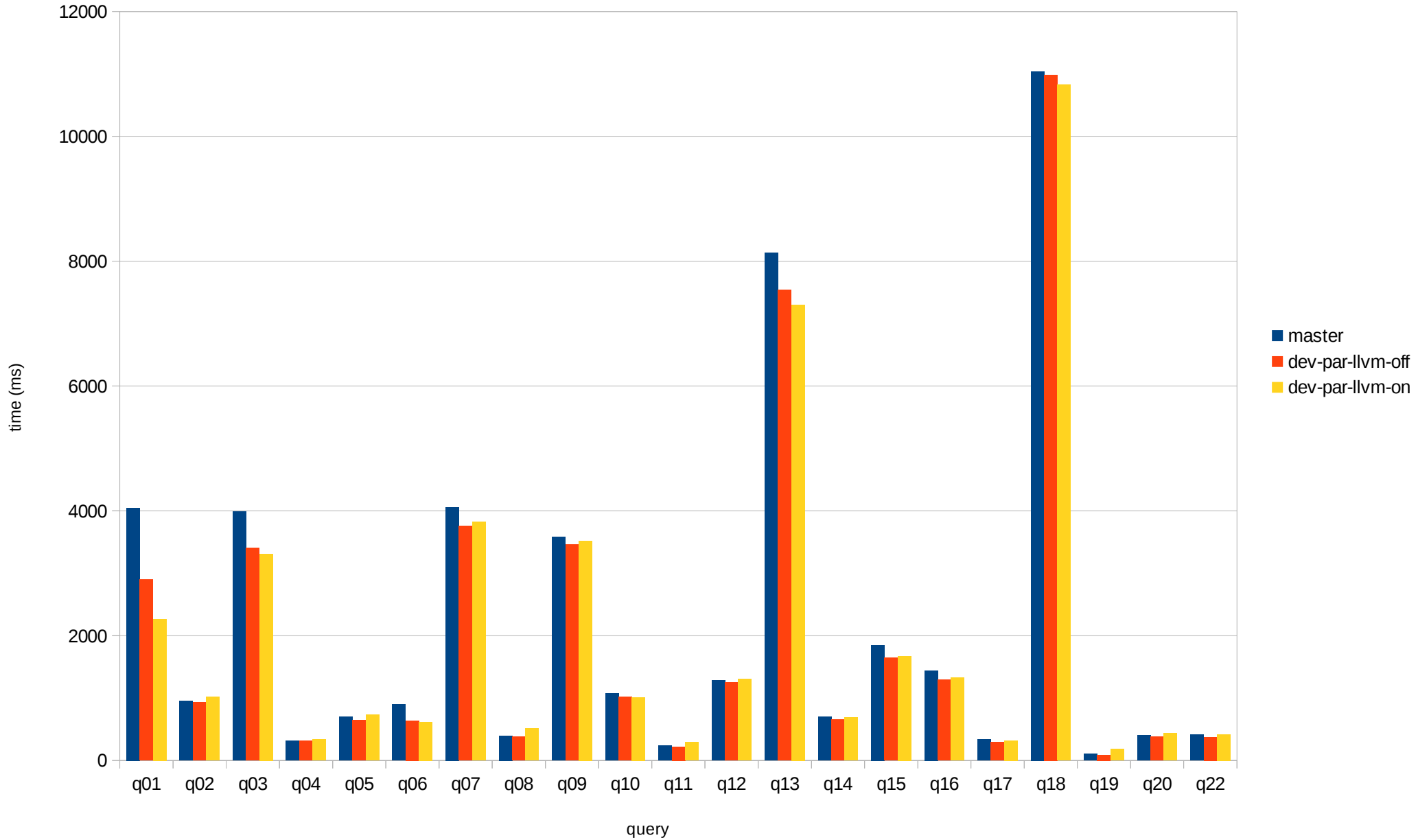
# TPCH Improvement

Not Parallelized, Scale 5



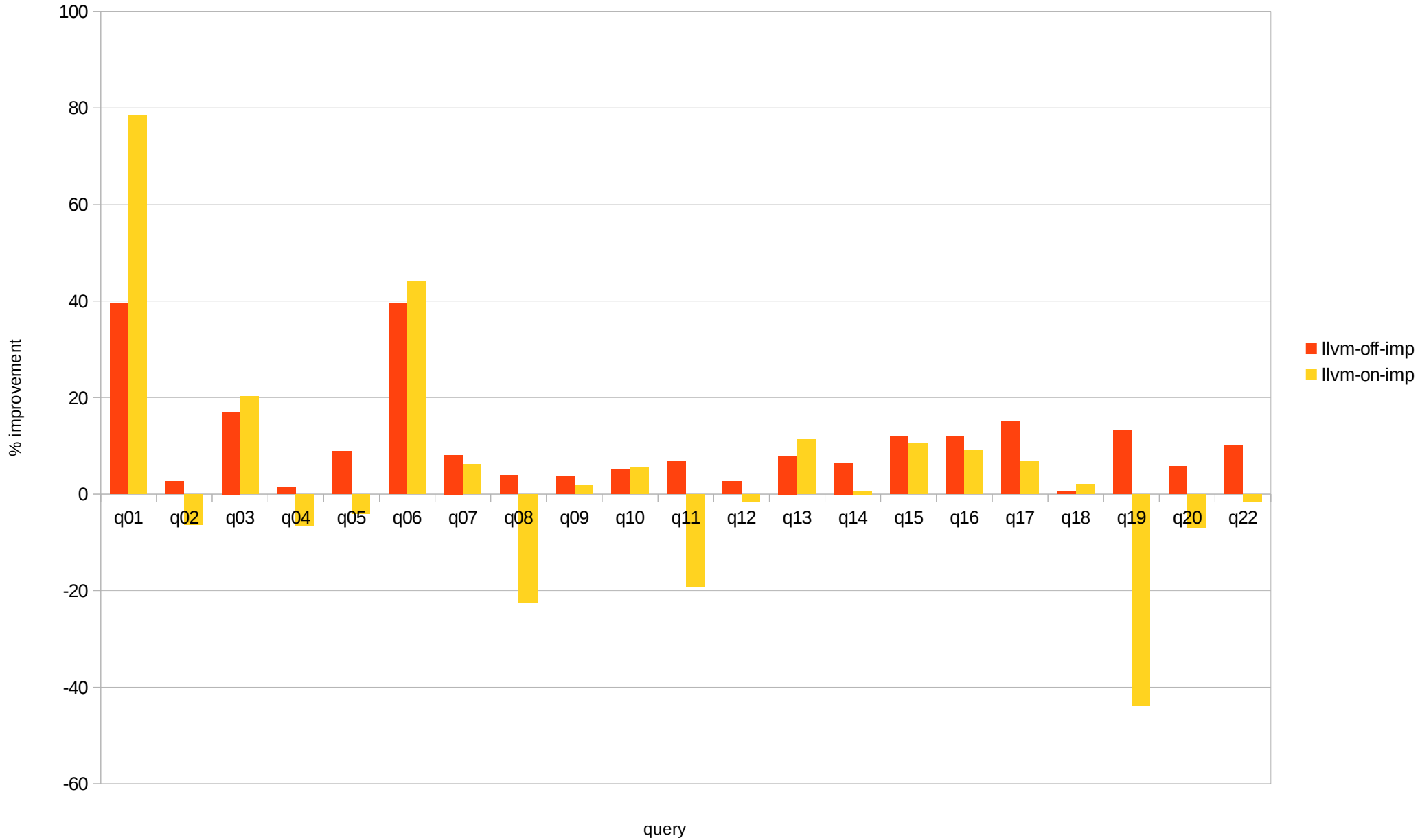
# TPCH Timings

Parallelized (6), Scale 5



# TPCH Timings

Parallelized (6), Scale 5



# Tuple Deforming

- On-Disk → Memory Representation
- Expensive, don't know columns ahead
- Improving Code: 1 Month, up to 3%
- JIT: ~1 Month, up to 80 %

# Better JIT

- Inline functions
  - WHERE  $(a + b) < c$ 
    - float8pl / int4pl / \*pl
    - float8lt / int4lt / \*lt
  - AVG(a + b)
    - int8pl
    - int8\_avg\_accum
    - numeric\_poly\_avg
  - Function calls are expensive, content cheap to execute
- Better Code

# Batch Tuple Processing #2

- Other Bottlenecks Are Gone
- Experiments show 0 - 40% in TPC-H queries

# Speeding up Query Execution

Andres Freund

PostgreSQL Developer & Committer  
Citus Data – citusdata.com - @citusdata

<http://anarazel.de/talks/pgday-nordic-2017-03-21/jit.pdf>