# Adaptive query optimization in PostgreSQL

Oleg Ivanov

Postgres Professional

What is query optimization?

How does PostgreSQL optimize queries?

What is adaptive query optimization?
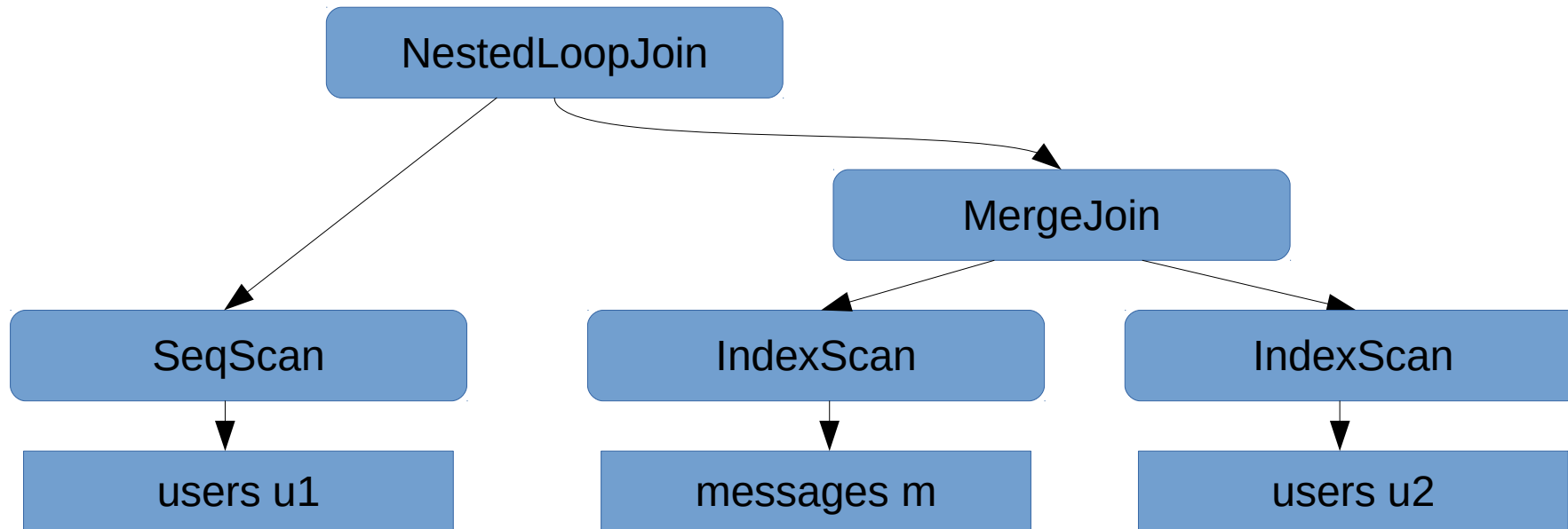
Machine learning and kNN method.

How to use machine learning for adaptive query optimization?

How much can it improve PostgreSQL performance?
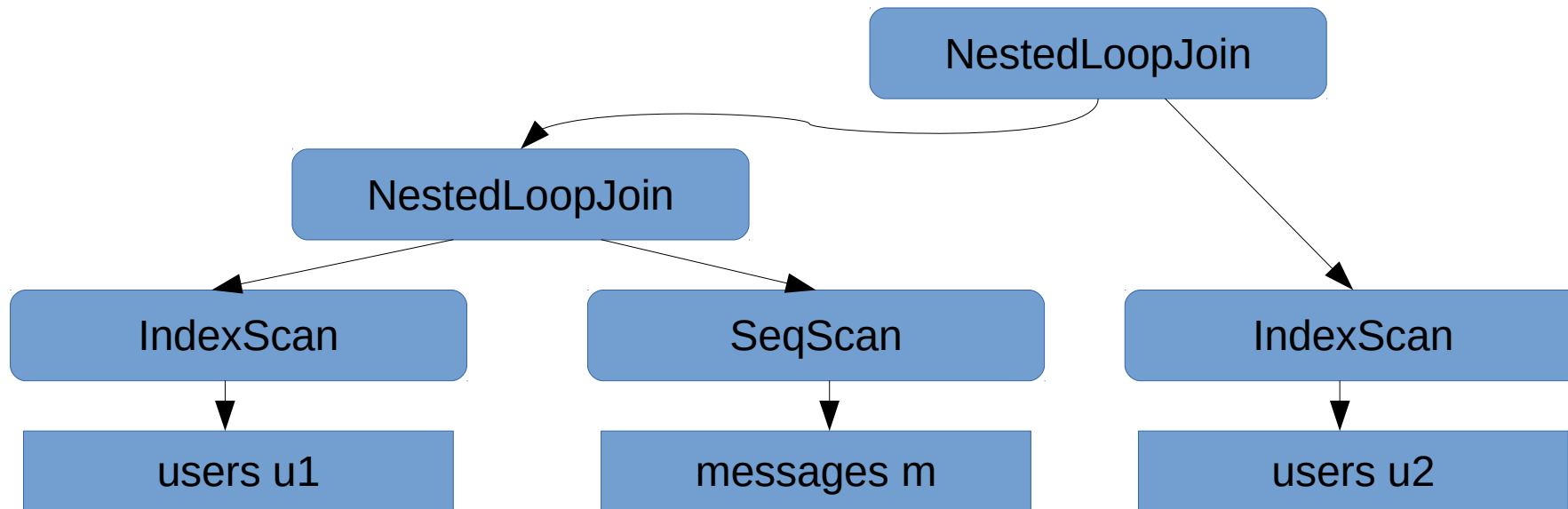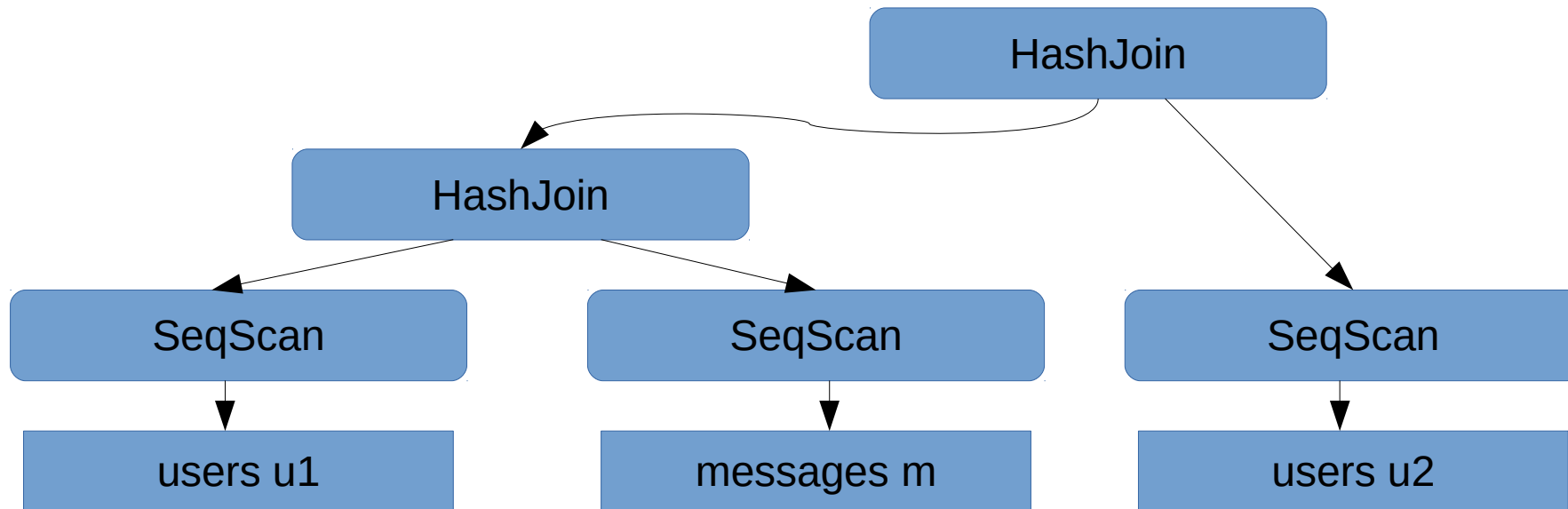
Implementaion details: AQO.

# What is query optimization?

```
SELECT *
FROM users AS u1, messages AS m, users AS u2
WHERE u1.id = m.sender_id AND m.receiver_id = u2.id;
```

# What is query optimization?

```
SELECT *
FROM users AS u1, messages AS m, users AS u2
WHERE u1.id = m.sender_id AND m.receiver_id = u2.id;
```

# What is query optimization?

```
SELECT *
FROM users AS u1, messages AS m, users AS u2
WHERE u1.id = m.sender_id AND m.receiver_id = u2.id;
```

# What is query optimization?

```
EXPLAIN SELECT *
FROM users AS u1, messages AS m, users AS u2
WHERE u1.id = m.sender_id AND m.receiver_id = u2.id;
                                QUERY PLAN
-----------------------------------------------------------------------------
 Hash Join  (cost=540.00..439429.44 rows=10003825 width=27)
   Hash Cond: (m.receiver_id = u2.id)
   ->  Hash Join  (cost=270.00..301606.84 rows=10003825 width=23)
         Hash Cond: (m.sender_id = u1.id)
         ->  Seq Scan on messages m  (cost=0.00..163784.25 rows=10003825 width=19)
         ->  Hash  (cost=145.00..145.00 rows=10000 width=4)
               ->  Seq Scan on users u1  (cost=0.00..145.00 rows=10000 width=4)
   ->  Hash  (cost=145.00..145.00 rows=10000 width=4)
         ->  Seq Scan on users u2  (cost=0.00..145.00 rows=10000 width=4)
(9 rows)
```

# What is query optimization?

Plan execution cost

Plan node execution cost

Plan node cardinality

```
EXPLAIN SELECT *
FROM users AS u1, messages AS m, users AS u2
WHERE u1.id = m.sender_id AND m.receiver_id = u2.id;
                              QUERY PLAN
--------------------------------------------------------------------------------
 Hash Join  (cost=540.00..439429.44 rows=10003825 width=27)
   Hash Cond: (m.receiver_id = u2.id)
   ->  Hash Join  (cost=270.00..301606.84 rows=10003825 width=23)
         Hash Cond: (m.sender_id = u1.id)
         ->  Seq Scan on messages m  (cost=0.00..163784.25 rows=10003825 width=19)
         ->  Hash  (cost=145.00..145.00 rows=10000 width=4)
               ->  Seq Scan on users u1  (cost=0.00..145.00 rows=10000 width=4)
   ->  Hash  (cost=145.00..145.00 rows=10000 width=4)
         ->  Seq Scan on users u2  (cost=0.00..145.00 rows=10000 width=4)
(9 rows)
```

# How does PostgreSQL optimize queries?

## Cost-based query optimization

System R (1974)

Choose the cheapest plan
among all the possible plans

# How does PostgreSQL optimize queries?

$$Cost = n_s c_s + n_r c_r + n_t c_t + n_i c_i + n_o c_o$$

| | | |
|---|---|---|
| $c_s$ | seq_page_cost | 1.0 |
| $c_r$ | random_page_cost | 4.0 |
| $c_t$ | cpu_Tuple_cost | 0.01 |
| $c_i$ | cpu_Index_tuple_cost | 0.005 |
| $c_o$ | cpu_Operator_cost | 0.0025 |

# How does PostgreSQL optimize queries?

```
SELECT * FROM users
WHERE age < 25;
```

SeqScan

Data

IndexScan

$$Cost = n_s c_s + n_o \cdot c_o$$

$$n_s = N_{pages}$$

$$n_o = N_{tuples}$$

# How does PostgreSQL optimize queries?

```
SELECT * FROM users
WHERE age < 25;
```

SeqScan

Data

IndexScan

Index

Data

$$Cost = n_s c_s + n_o \cdot c_o$$
$$n_s = N_{pages}$$
$$n_o = N_{tuples}$$

$$Cost = n_r \cdot c_r$$
$$n_r = Cardinality$$

# How does PostgreSQL optimize queries?

```
SELECT * FROM users
WHERE age < 25;
```



*Cardinality*

# How does PostgreSQL optimize queries?

```
SELECT *
FROM users AS u1, messages AS m, users AS u2
WHERE u1.id = m.sender_id AND m.receiver_id = u2.id;
```

# How does PostgreSQL optimize queries?

```
SELECT *
FROM users AS u1, messages AS m, users AS u2
WHERE u1.id = m.sender_id AND m.receiver_id = u2.id;
```

# How does PostgreSQL optimize queries?

```
SELECT *
FROM users AS u1, messages AS m, users AS u2
WHERE u1.id = m.sender_id AND m.receiver_id = u2.id;
```
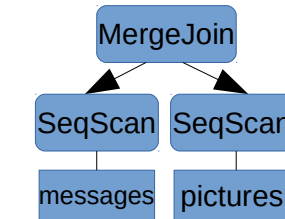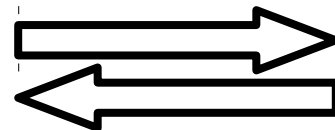
# How does PostgreSQL optimize queries?

## Optimization method

**Full search**

HashJoin

HashJoin

SeqScan    SeqScan    SeqScan

users    messages    pictures

Cost = 439429

MergeJoin

SeqScan    SeqScan

messages    pictures

Cost = 304528

## Plan's cost estimation

# How does PostgreSQL optimize queries?

**Optimization method**

Full search

HashJoin

HashJoin

SeqScan  SeqScan  SeqScan

users  messages  pictures

Cost = 439429

MergeJoin

SeqScan  SeqScan

messages  pictures

Cost = 304528
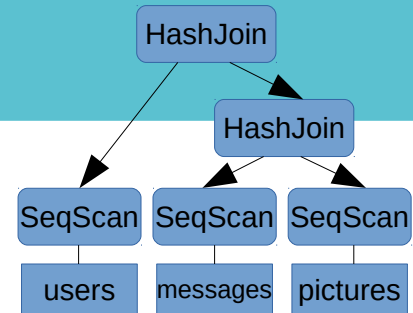
**Plan's cost estimation**

# How does PostgreSQL optimize queries?

**Optimization method**
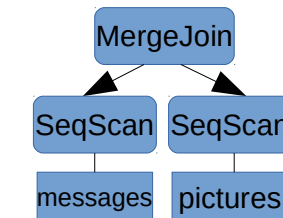
**Dynamic programming**

or

**Genetic algorithm**

HashJoin

HashJoin

SeqScan  SeqScan  SeqScan

users  messages  pictures

Cost = 439429

MergeJoin

SeqScan  SeqScan

messages  pictures

Cost = 304528

**Plan's cost estimation**

# Dynamic programming over subsets

- System R
- Time complexity: $3^n$
- Memory consumption: $2^n$
- Always finds the cheapest plan

# Genetic algorithm

- PostgreSQL
- Common and flexible method
- Can be stopped on every iteration
- No guarantees

# How does PostgreSQL optimize queries?

**Optimization method**

**Dynamic programming**

or

**Genetic algorithm**

HashJoin
HashJoin
SeqScan  SeqScan  SeqScan
users  messages  pictures

Cost = 439429

MergeJoin
SeqScan  SeqScan
messages  pictures

Cost = 304528

**Plan's cost estimation**

TPC-DS qualification benchmark

Dataset:
The TPC Benchmark™DS (TPC-DS)
http://www.tpc.org/tpcds/

Error: 300 times

Error: 4 times

TPC-DS qualification benchmark

TPC-DS qualification benchmark

Actual cardinality

Predicted cardinality

Execution time, ms

Total cost

Dataset:
The TPC Benchmark™DS (TPC-DS)
http://www.tpc.org/tpcds/

```
SELECT * FROM users
WHERE age < 25;
```

$$Selectivity \simeq 0.3$$

$$Cardinality = N_{tuples} \cdot Selectivity$$

```
SELECT * FROM users
WHERE age < 25 AND city = 'Moscow';
```

Only selectivities of individual clauses
(i.e. *marginal* selectivities)
are known

$$Selectivity_{age} = \frac{1}{3}$$

$$Selectivity_{city} = \frac{1}{7}$$

$$Selectivity_{age,city} = ?$$
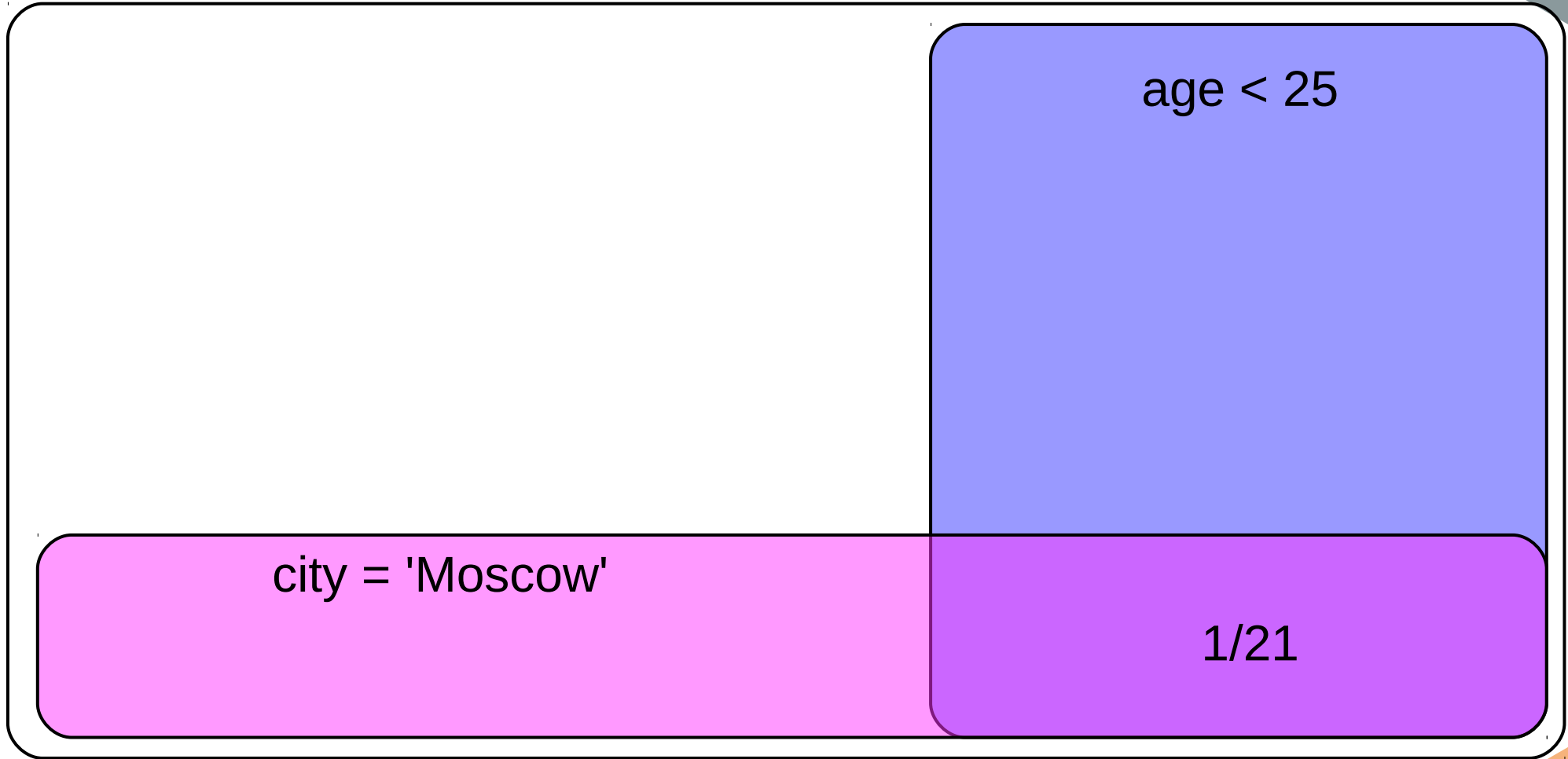
1/3
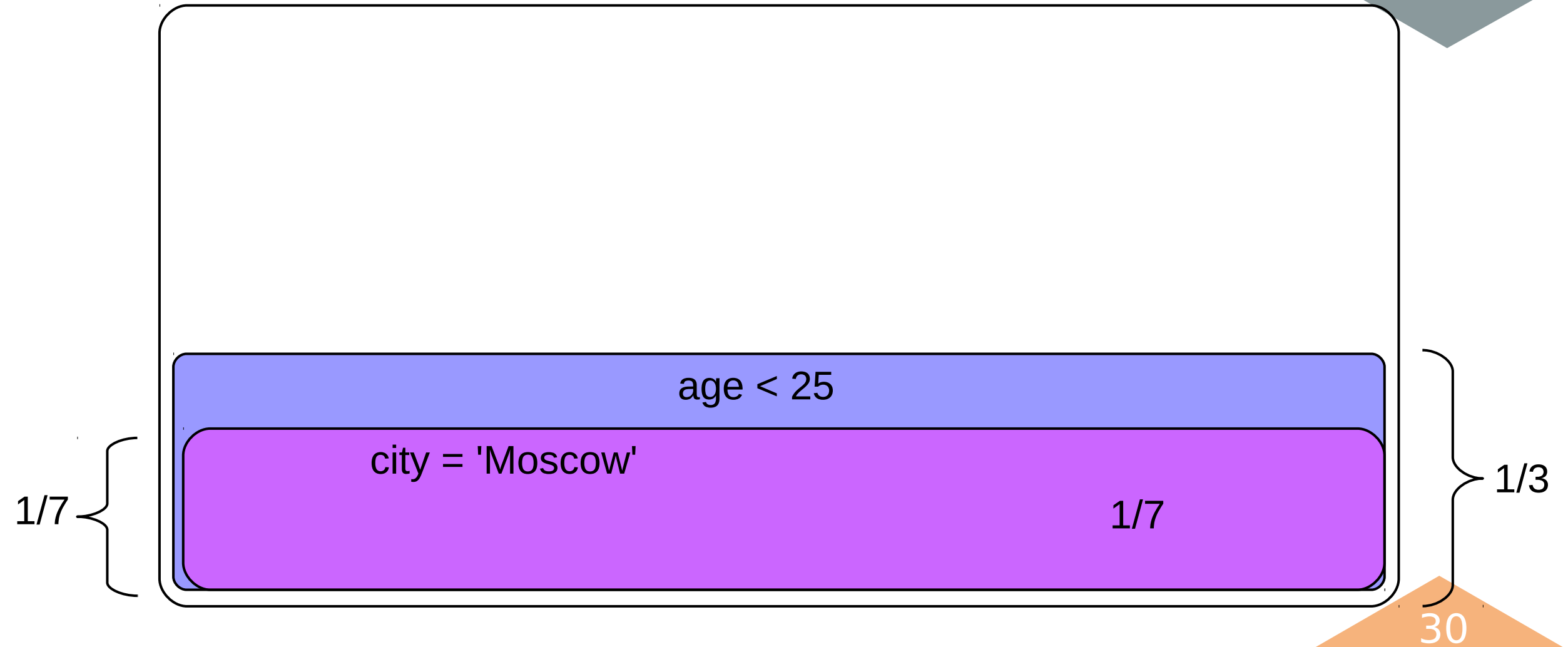
age < 25

city = 'Moscow'

1/21

1/7

```
SELECT * FROM users
WHERE age < 25 AND city = 'Moscow';
```

## Only selectivities of individual clauses are known

### The clauses are considered to be independent:

$$Selectivity_{age,city} = Selectivity_{age} \cdot Selectivity_{city}$$

With the exception of $Selectivity_{25 < age \; AND \; age < 57} = Selectivity_{25 < age < 57}$

age < 25

city = 'Moscow'

1/7

1/7

1/3

age < 25

1/3

city = 'Moscow'

1/7

```
SELECT * FROM users
WHERE age < 12 AND married = true;
```

age < 12

married = true
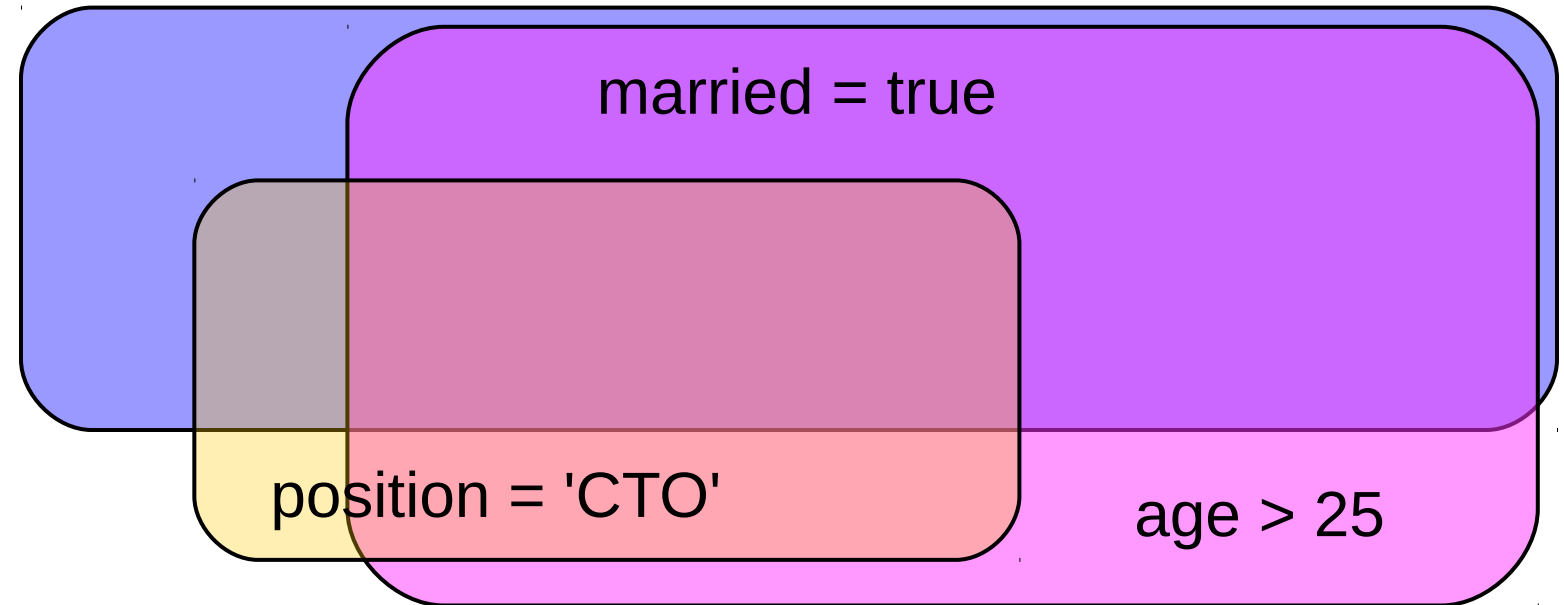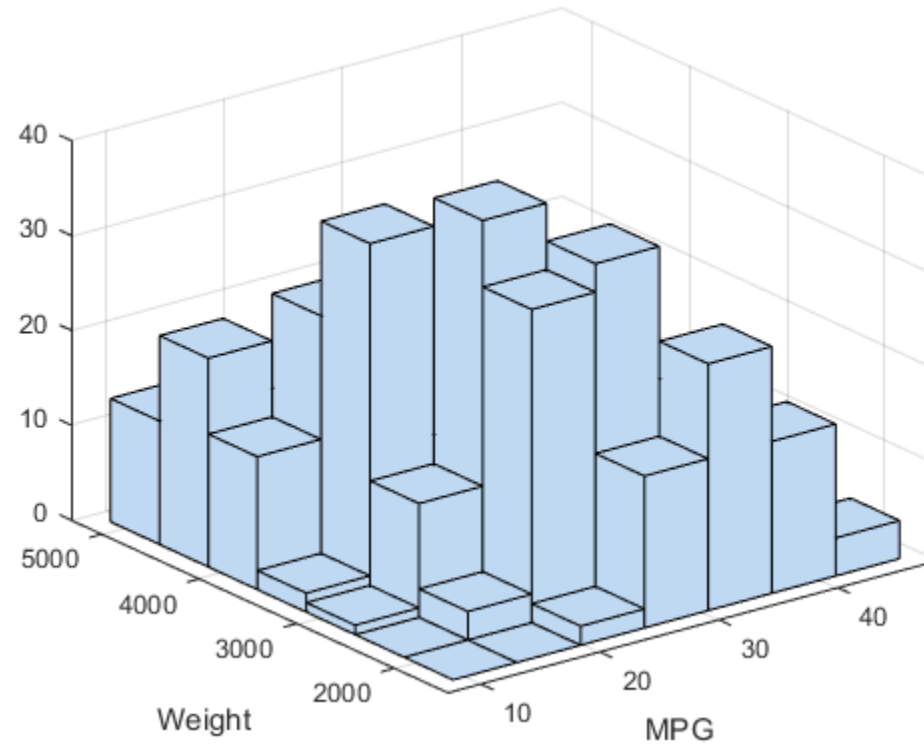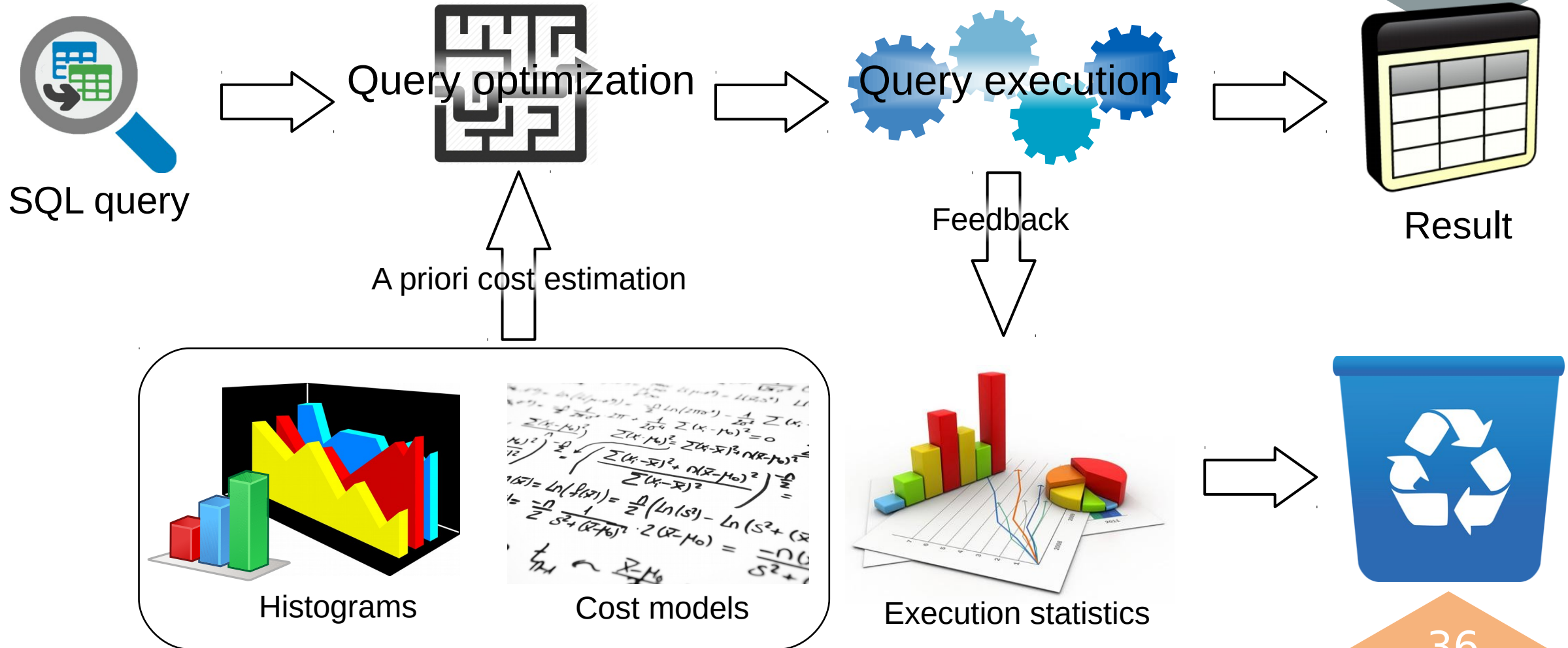
```
SELECT * FROM users
WHERE age > 25 AND married = true
AND position = 'CTO';
```



married = true

position = 'CTO'

age > 25
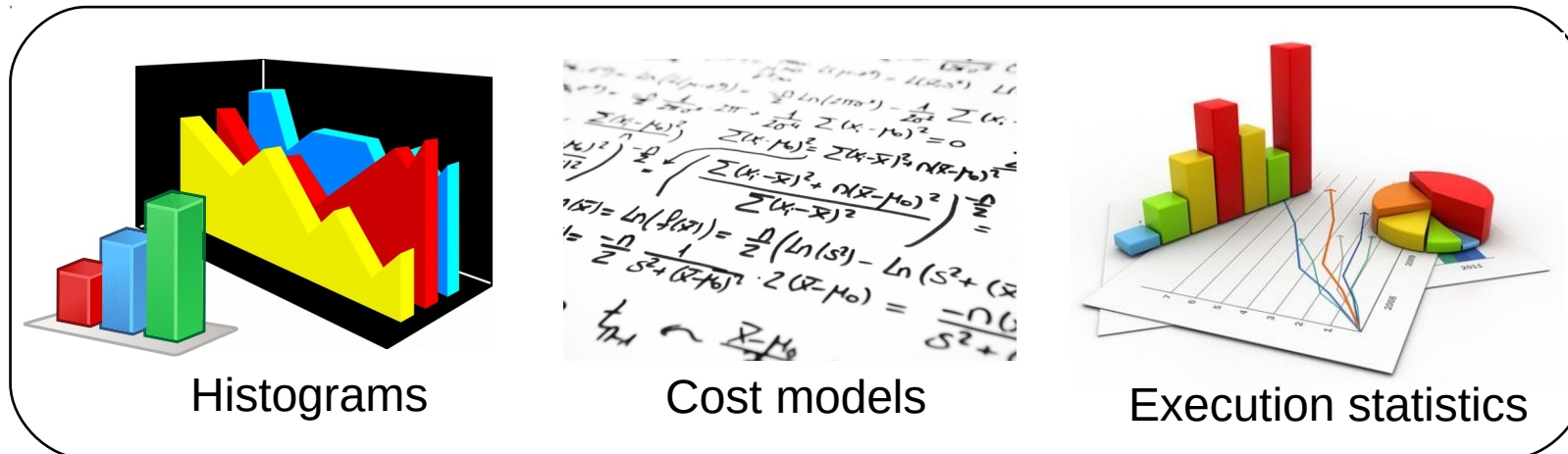
# Multidimensional histograms

# What is adaptive query optimization?

SQL query → Query optimization → Query execution → Result

A priori cost estimation

Feedback

Histograms    Cost models    Execution statistics

# What is adaptive query optimization?



SQL query

Query optimization

Query execution

Result

A priori cost estimation

Feedback

Histograms

Cost models

Execution statistics

# Machine learning

# K Nearest Neighbours method

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Age** | 25 | 47 | 55 | 32 | 22 | 45 | 28 |
| **Salary** | 50 | 120 | 100 | 80 | 30 | 90 | ? |

# K Nearest Neighbours method

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Age | 25 | 47 | 55 | 32 | 22 | 45 | 28 |
| Salary | 50 | 120 | 100 | 80 | 30 | 90 | ? |

# Gradient approach to kNN

2/3

1/3

1/3    5/3

| | | | |
|---|---|---|---|
| Age | 27 | 47 | 28 |
| Salary | 53 | 103 | ? |

# How to use machine learning
# for adaptive query optimization?

# The object is a node with its subtree



NestedLoopJoin

u1.id = messages.sender_id

MergeJoin

u2.id = messages.receiver_id

SeqScan

IndexScan

IndexScan

u2.married = true
**AND**
u2.age < 25

users u1

messages m

users u2

# Machine learning problem statement

Object is a plan node

Features
- users.id = messages.receiver_id        0.0001
- users.married = **const**        0.73
- users.age < **const**        0.23
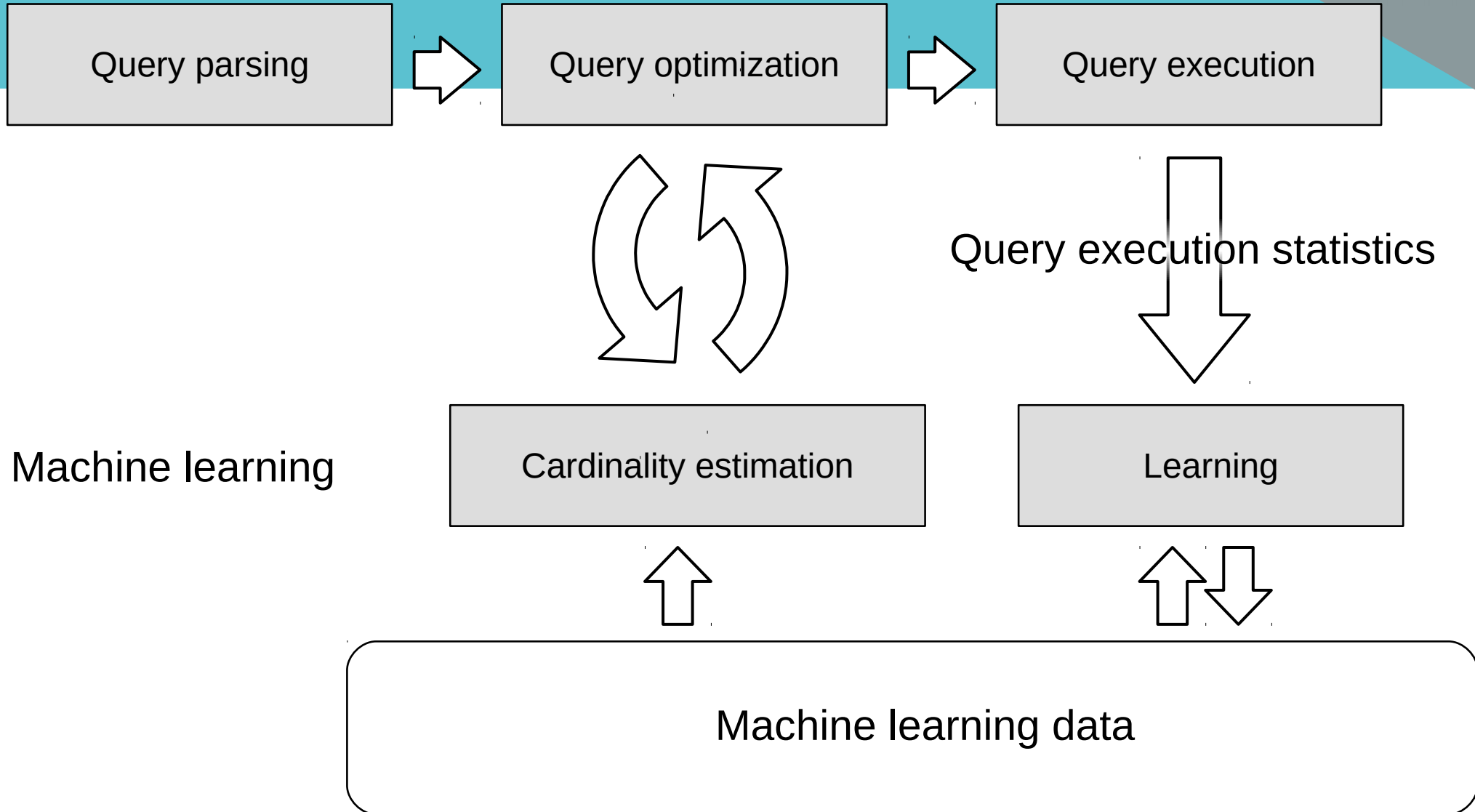
Hidden value     Node cardinality        ?

# Workflow

| Query parsing | ➡ | Query optimization | ➡ | Query execution |
|---|---|---|---|---|

Query execution statistics

Machine learning

| Cardinality estimation | | Learning |
|---|---|---|

Machine learning data
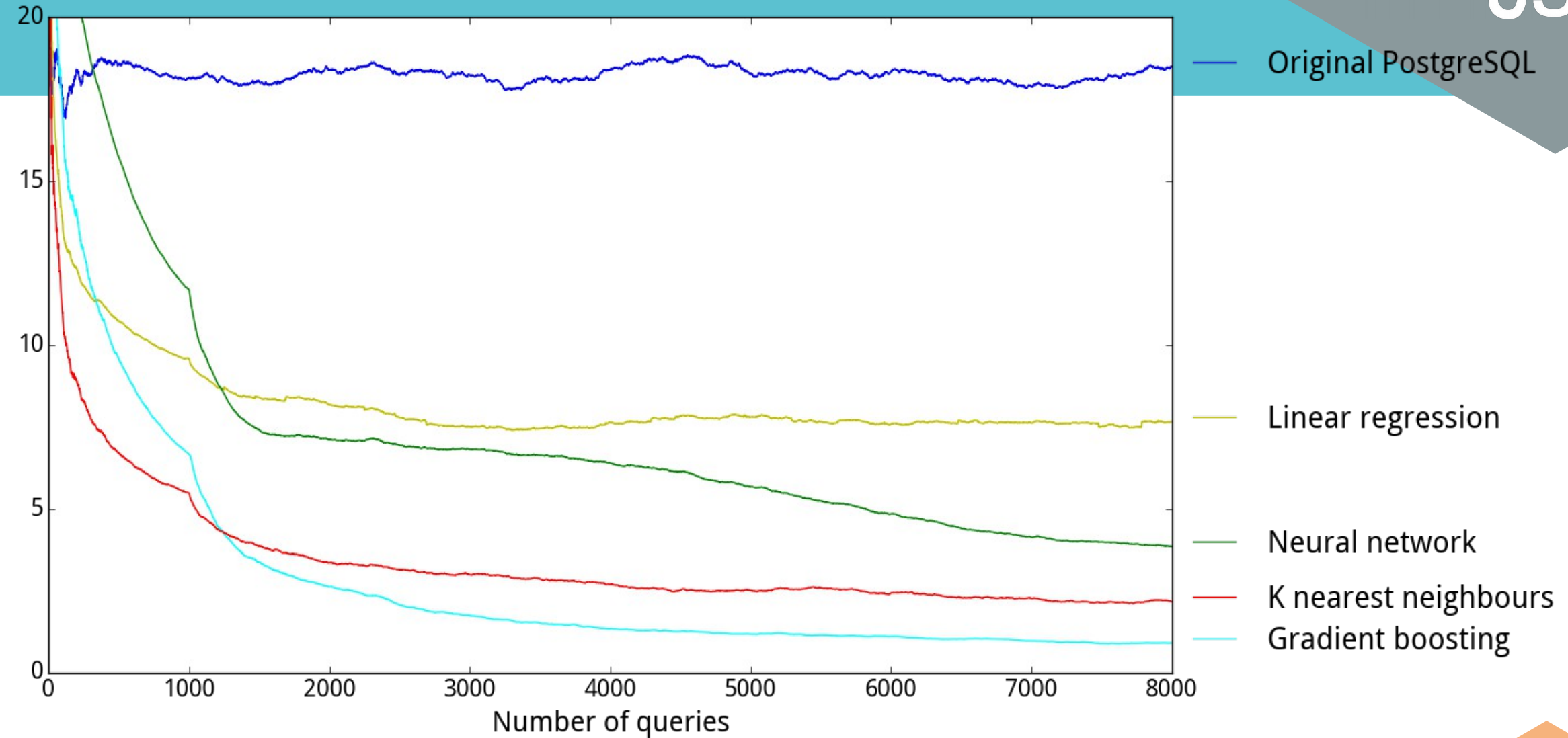
# Theoretical properties

- Will it converge?
Yes, in the finite number of steps

- How fast will it converge?
Don't know (in practice in a few steps)

- What guarantees on obtained plans or regressor do we have?
Predictions are correct for all executed paths
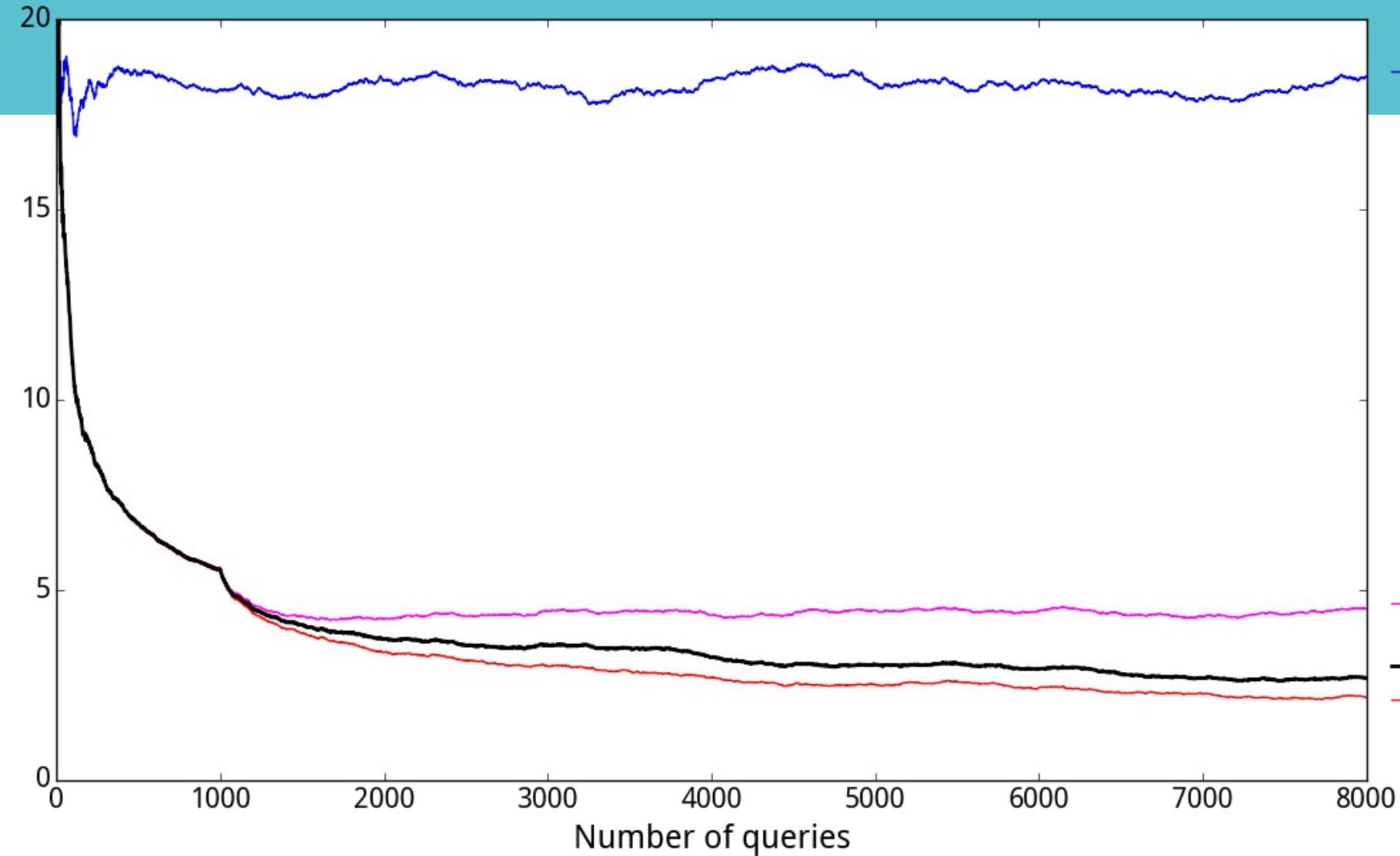With perfect cost model obtained plans are not worse

# How much can it improve PostgreSQL performance?

Experimental evaluation

# Estimation error



**PROFESSIONAL**

**Postgres**

Legend:
- Original PostgreSQL
- Linear regression
- Neural network
- K nearest neighbours
- Gradient boosting

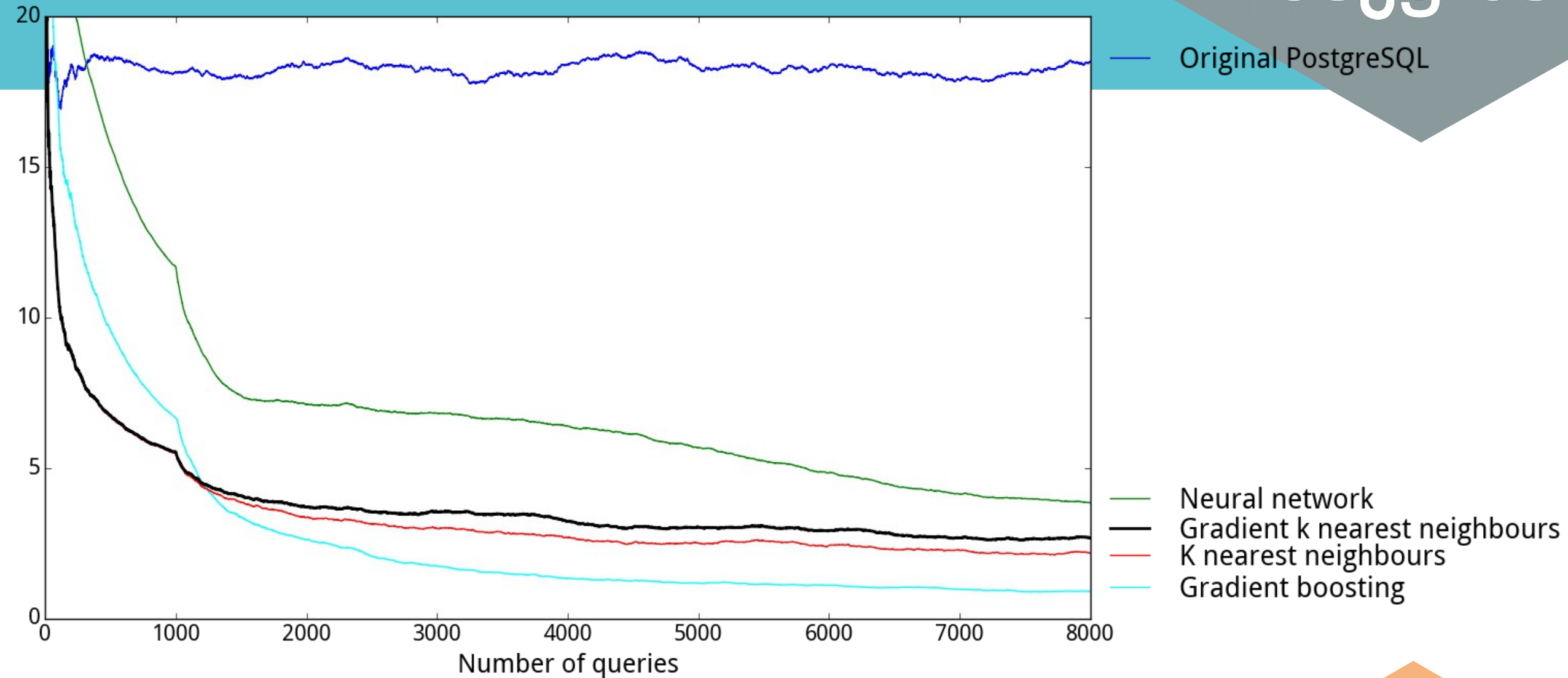X-axis: Number of queries

# Estimation error


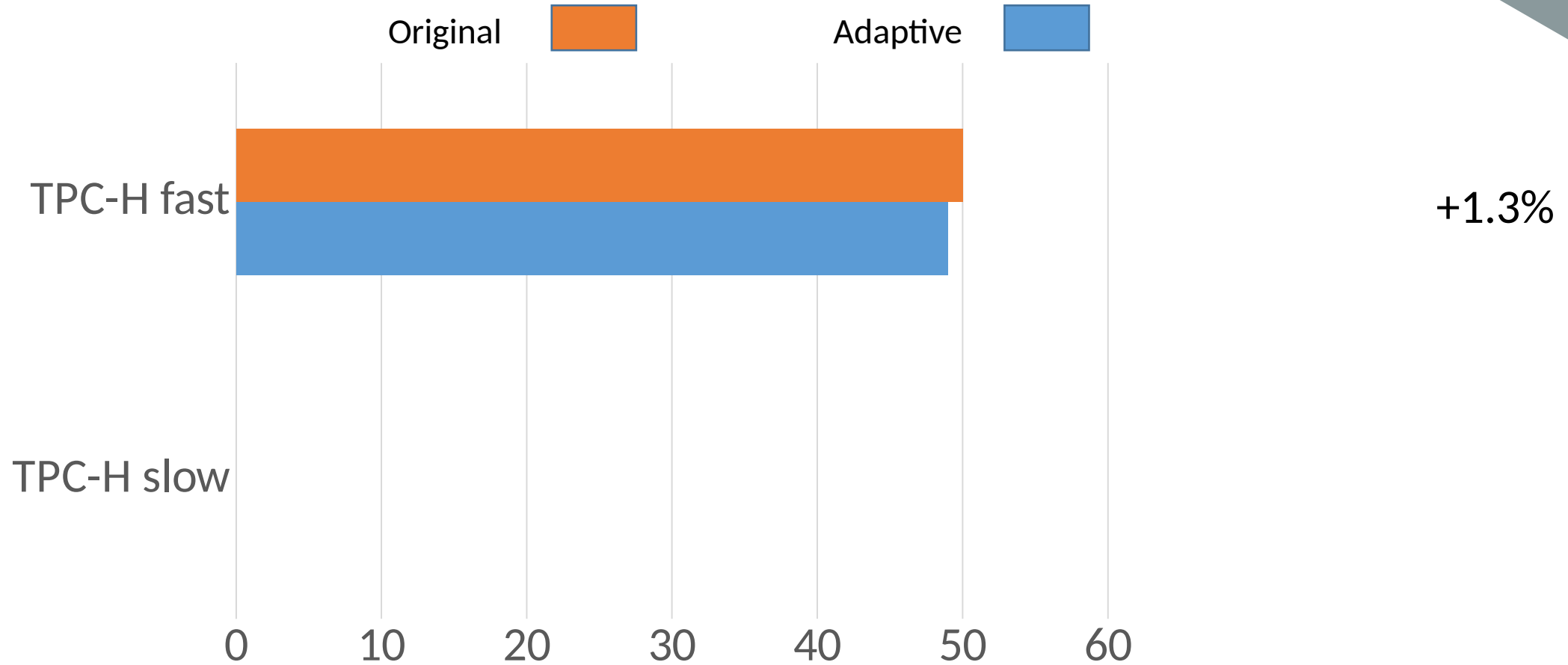
Original PostgreSQL

k nearest neighbours limited

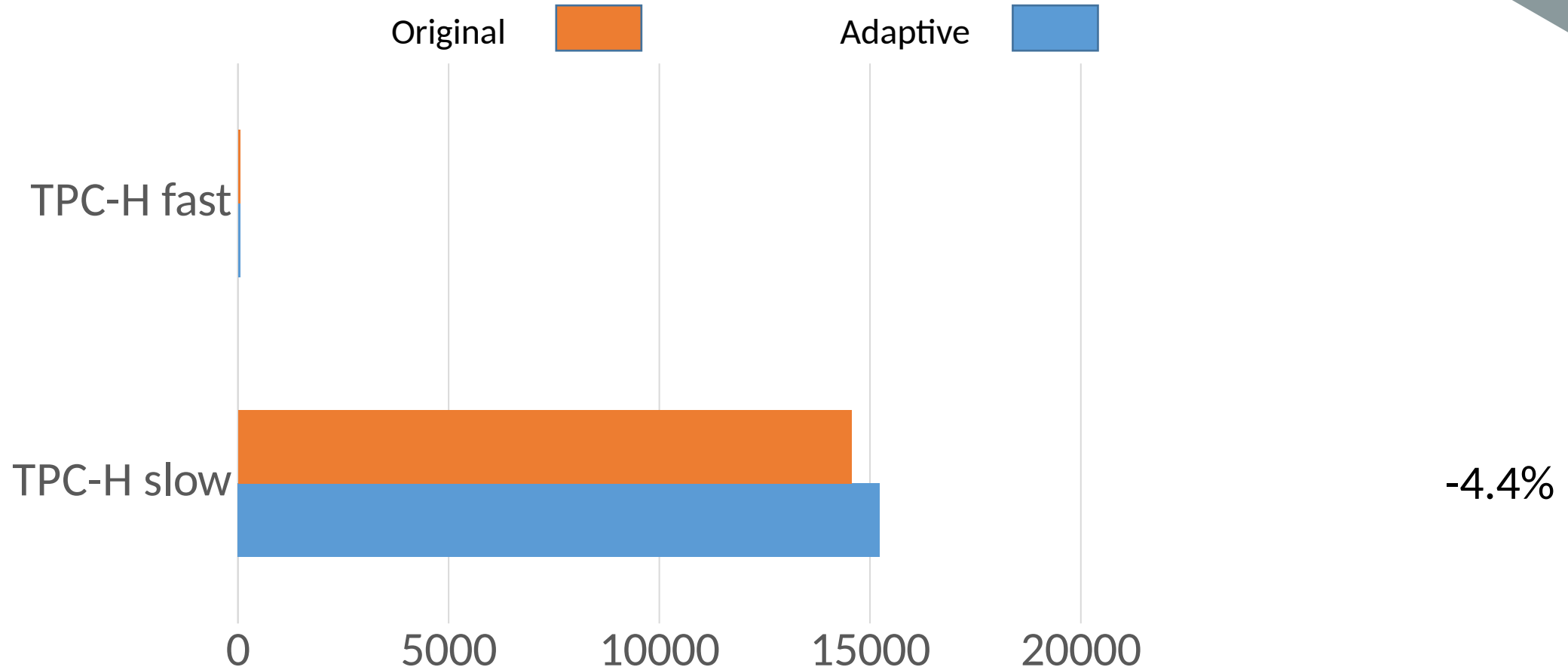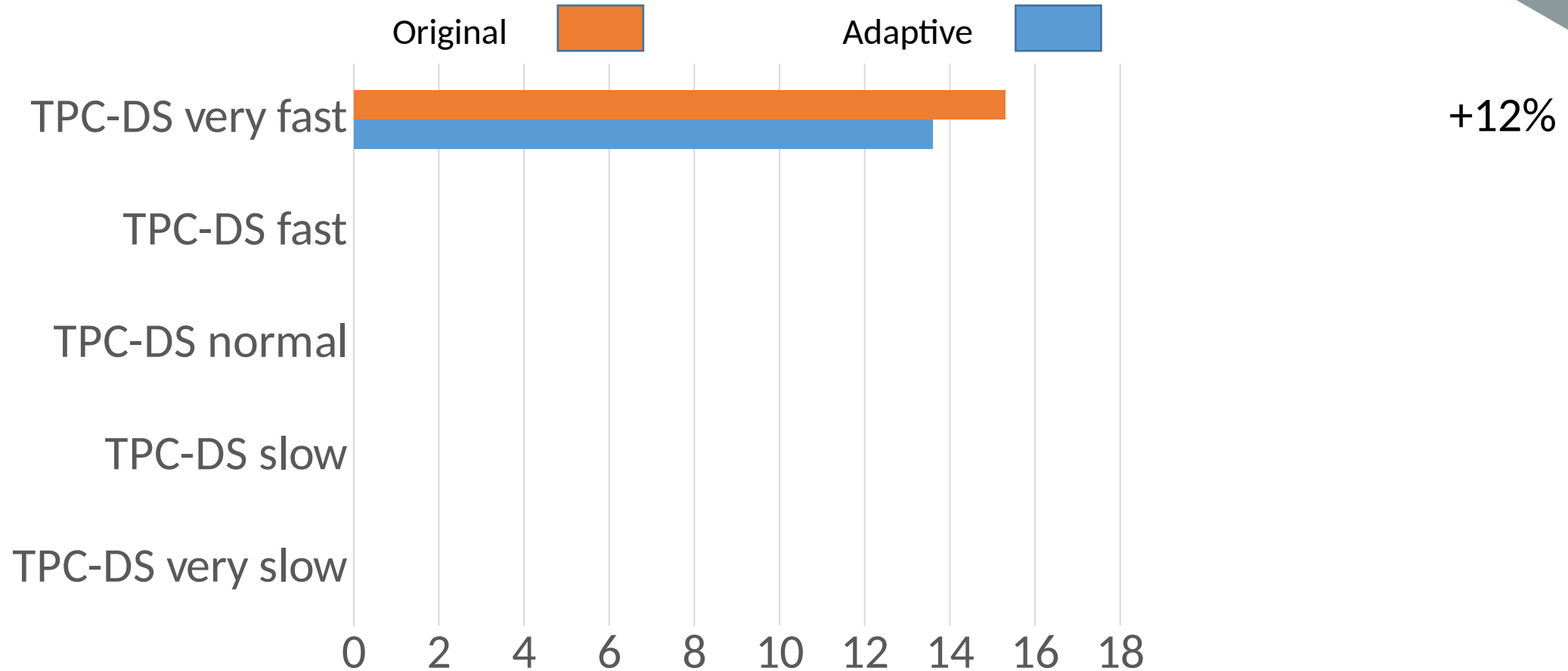Gradient k nearest neighbours

K nearest neighbours
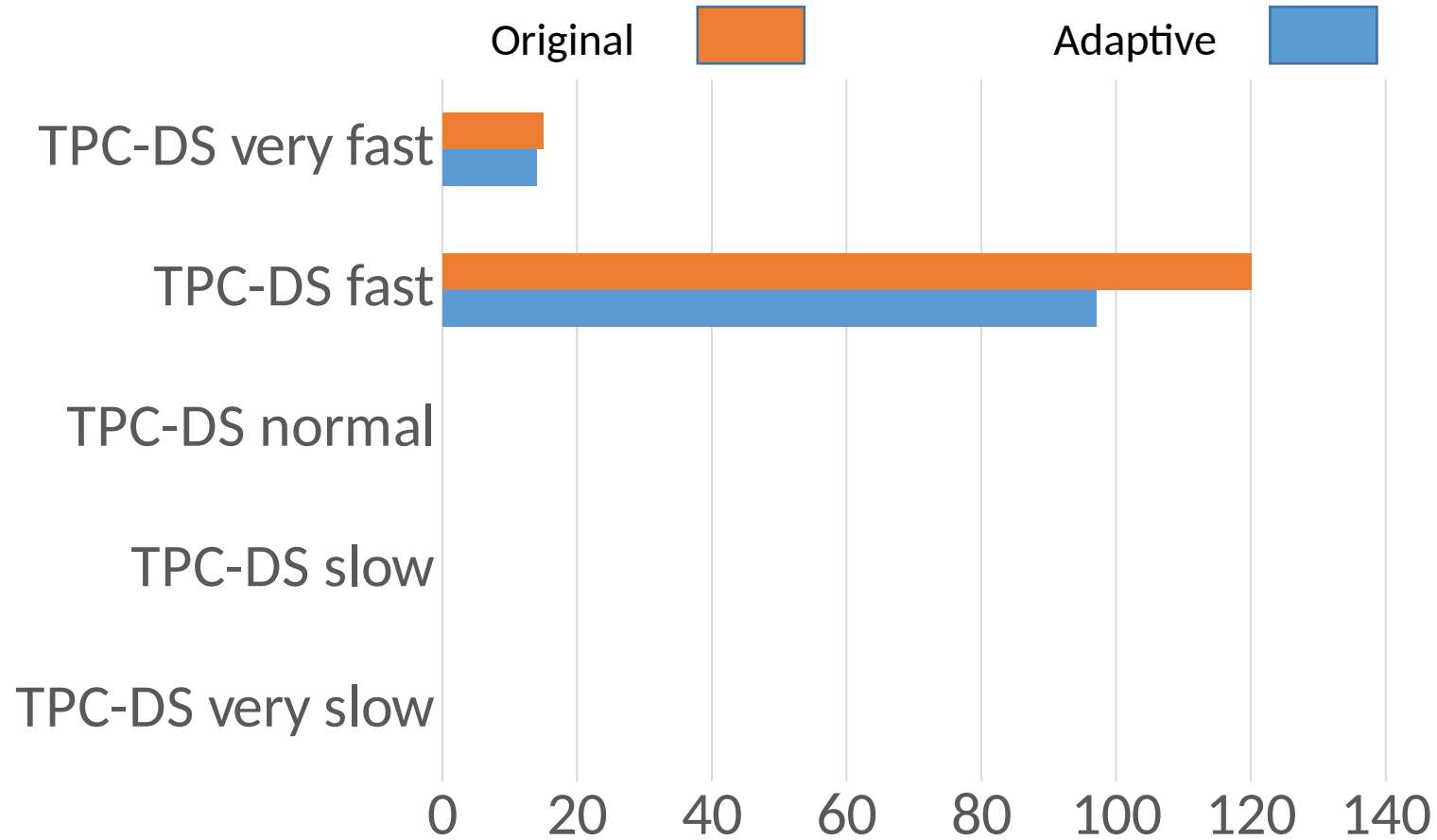
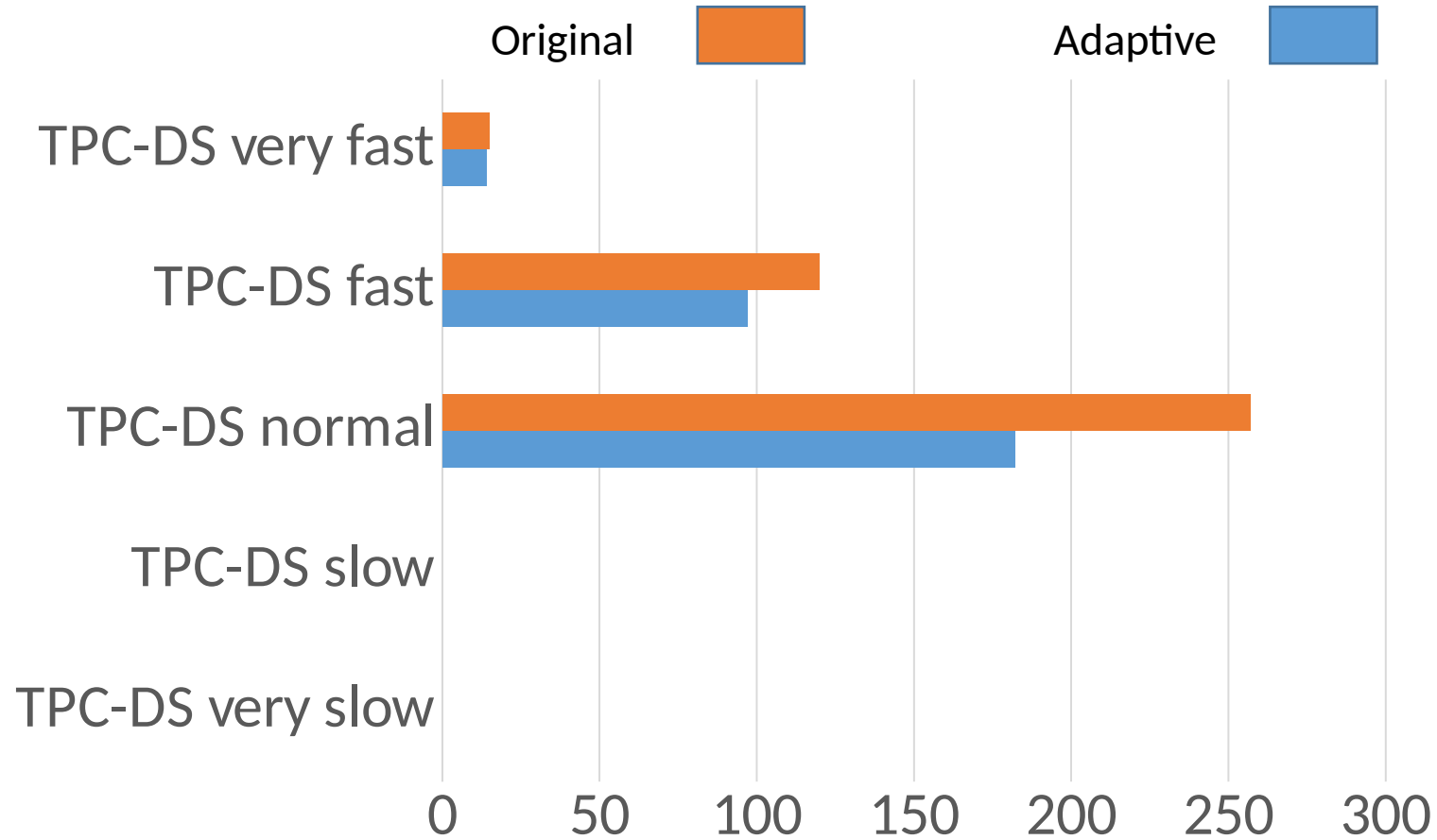# Estimation error

# Performance improvement



+1.3%

53

# Performance improvement

# Performance improvement



+12%

# Performance improvement



Original Adaptive

- TPC-DS very fast
- TPC-DS fast
- TPC-DS normal
- TPC-DS slow
- TPC-DS very slow

0  20  40  60  80  100  120  140

+24%

# Performance improvement



TPC-DS very fast, TPC-DS fast, TPC-DS normal, TPC-DS slow, TPC-DS very slow

Original — Adaptive

+41%

0  50  100  150  200  250  300
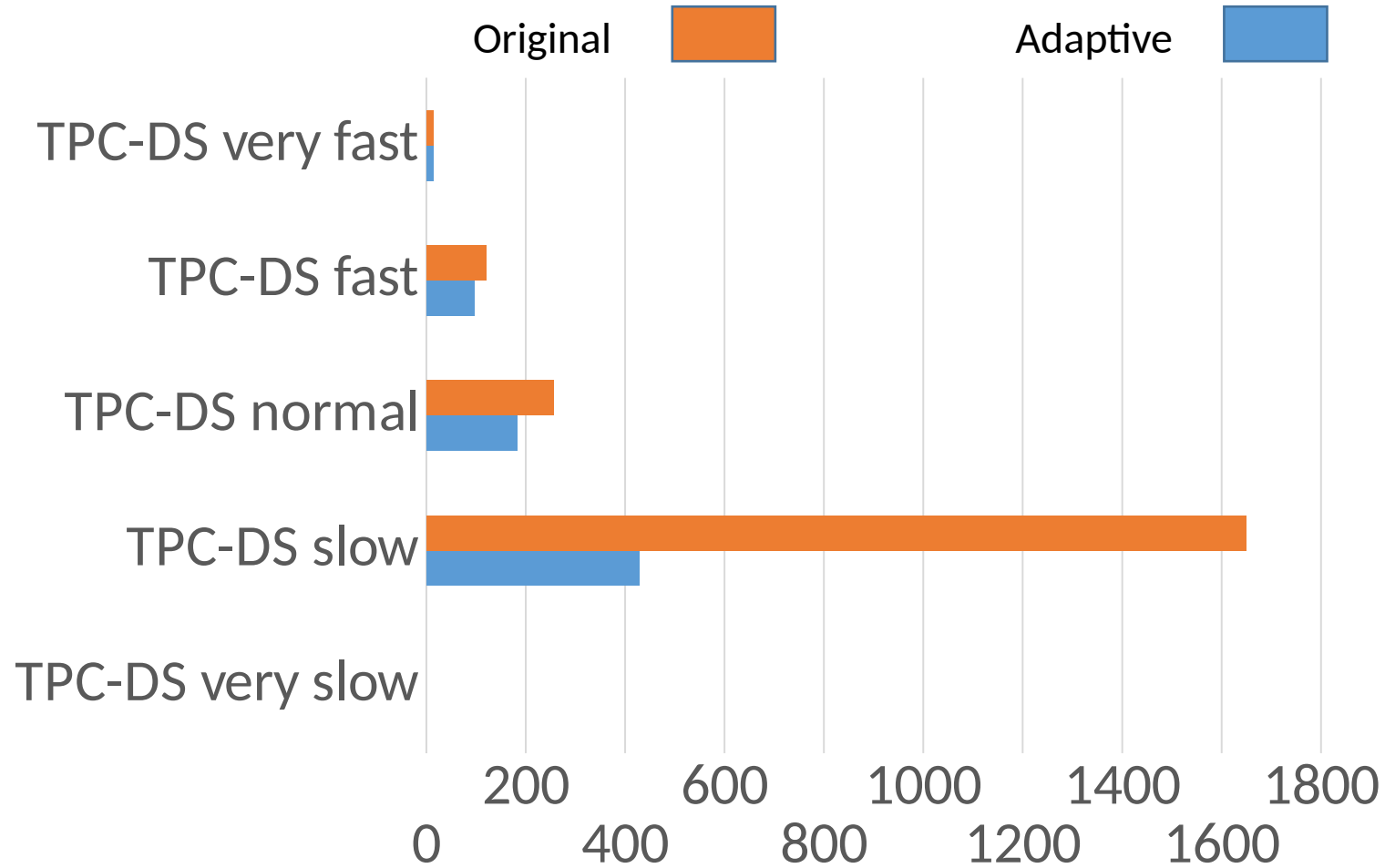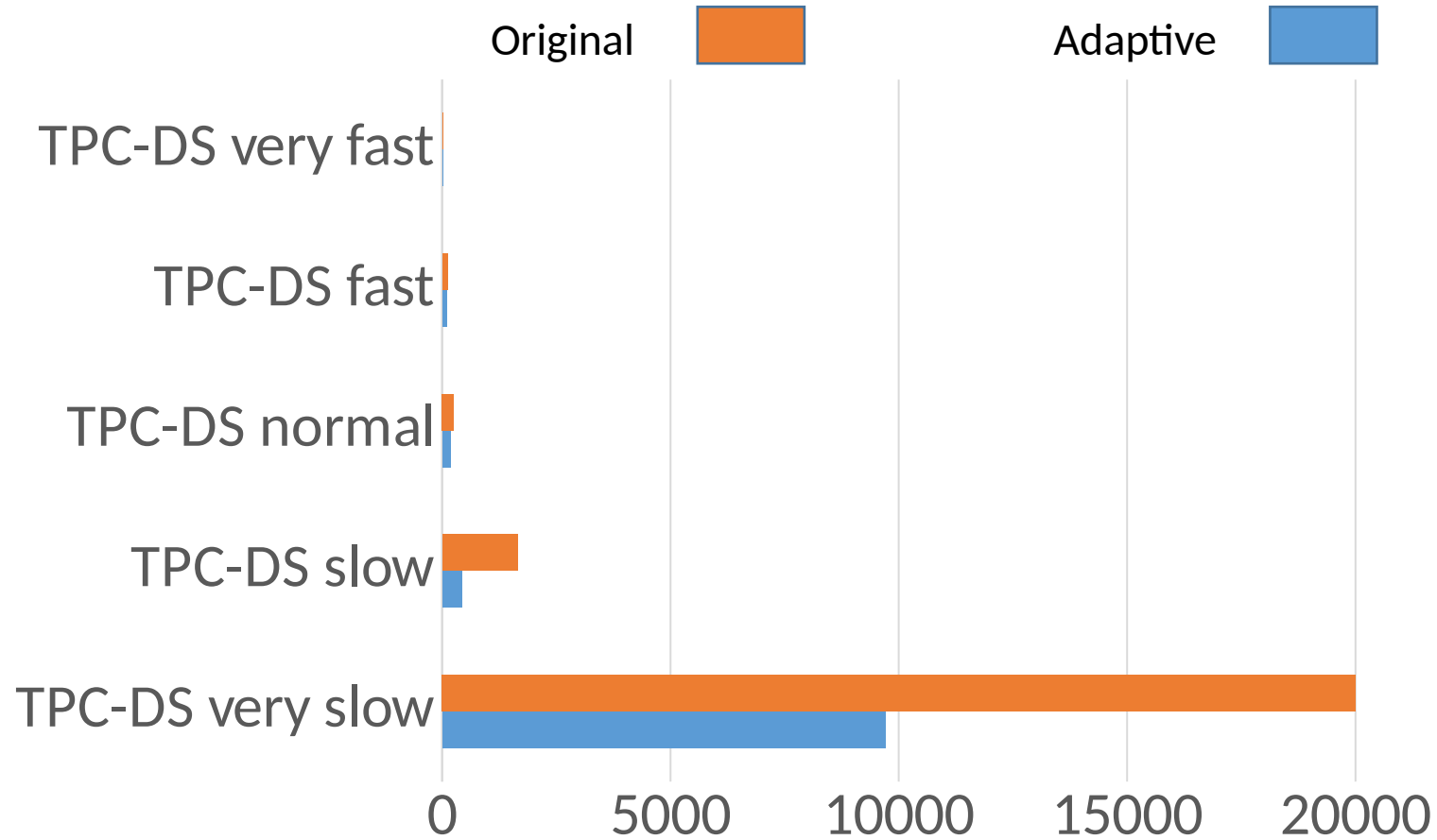
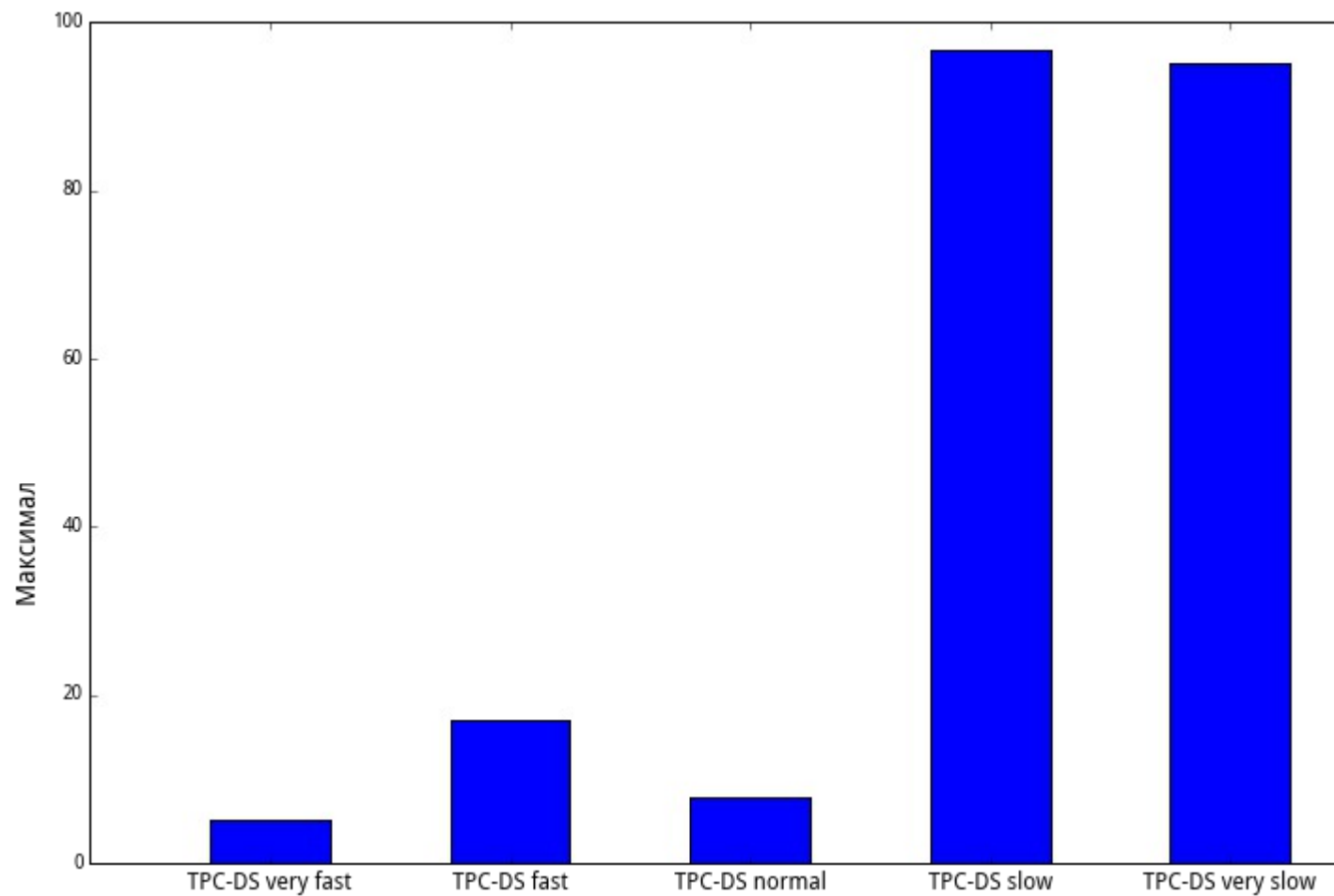# Performance improvement



+285%

# Performance improvement



+115%

59

# Maximum acceleration

# Overheads
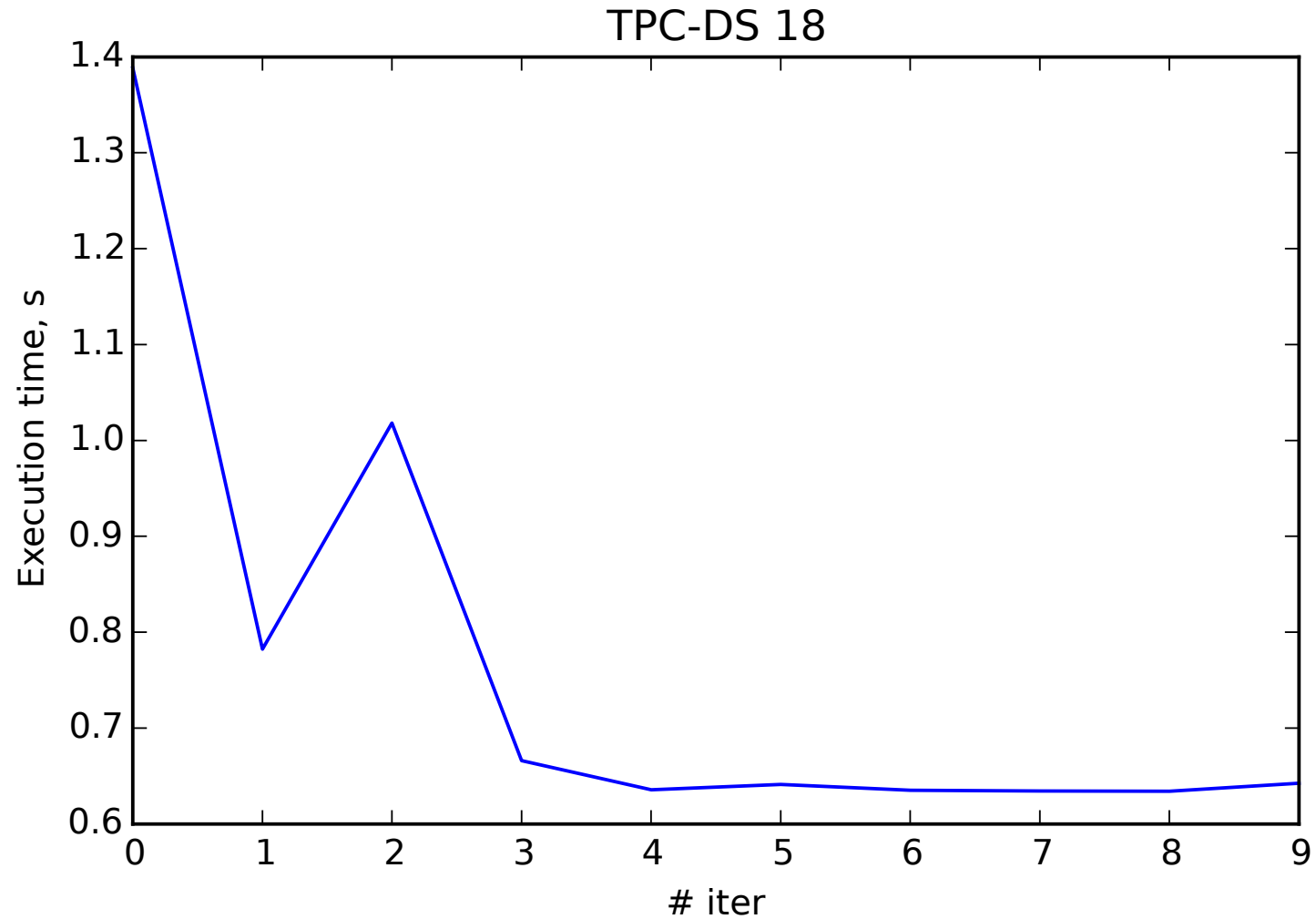
Experimental evaluation

Slowdown for genetic algorithm is not more than 2 seconds

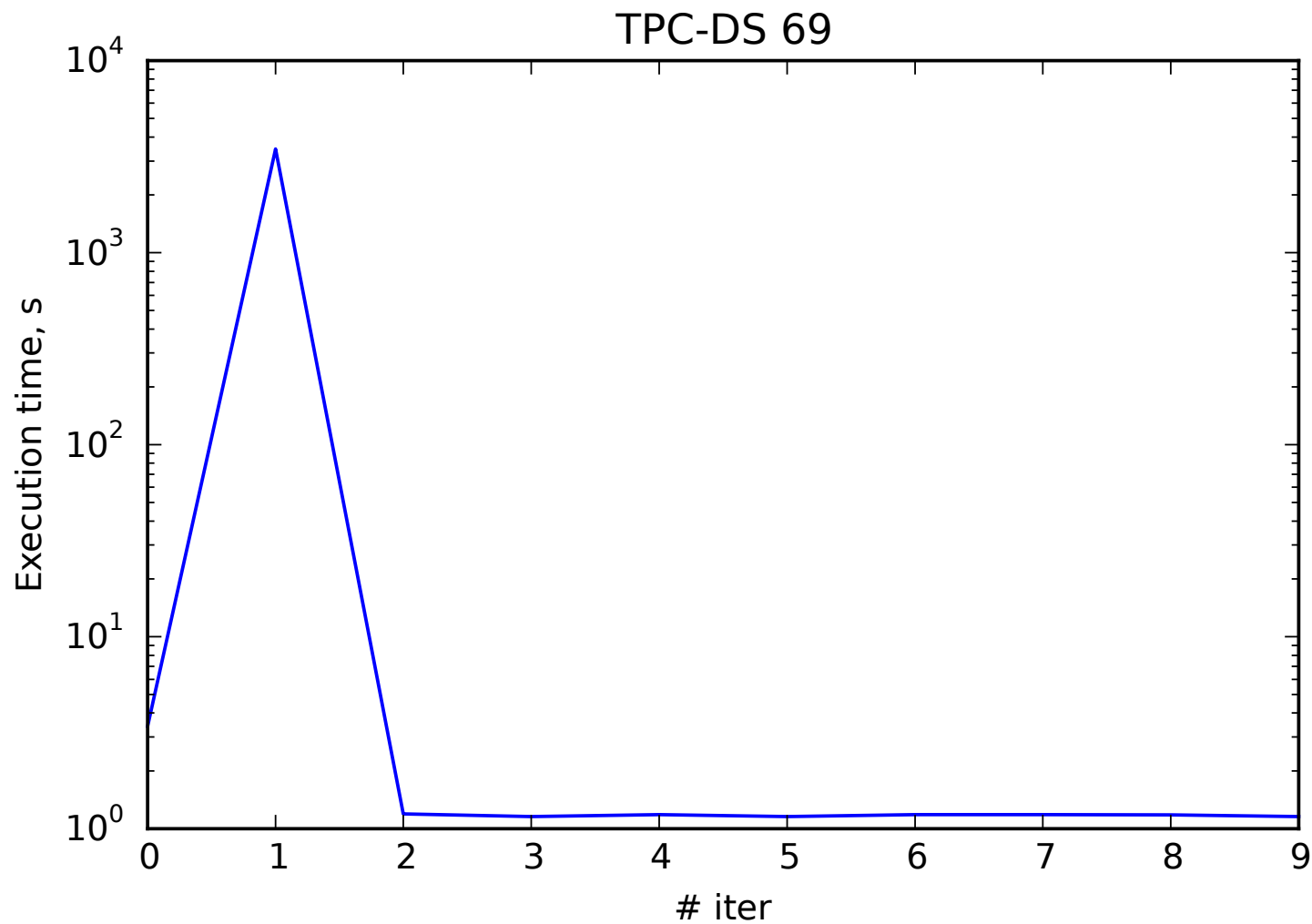Slowdown for dynamic programming is not more than 30 ms

# Applicability

Complex analytical queries
with a repeating pattern.

# Learning progress
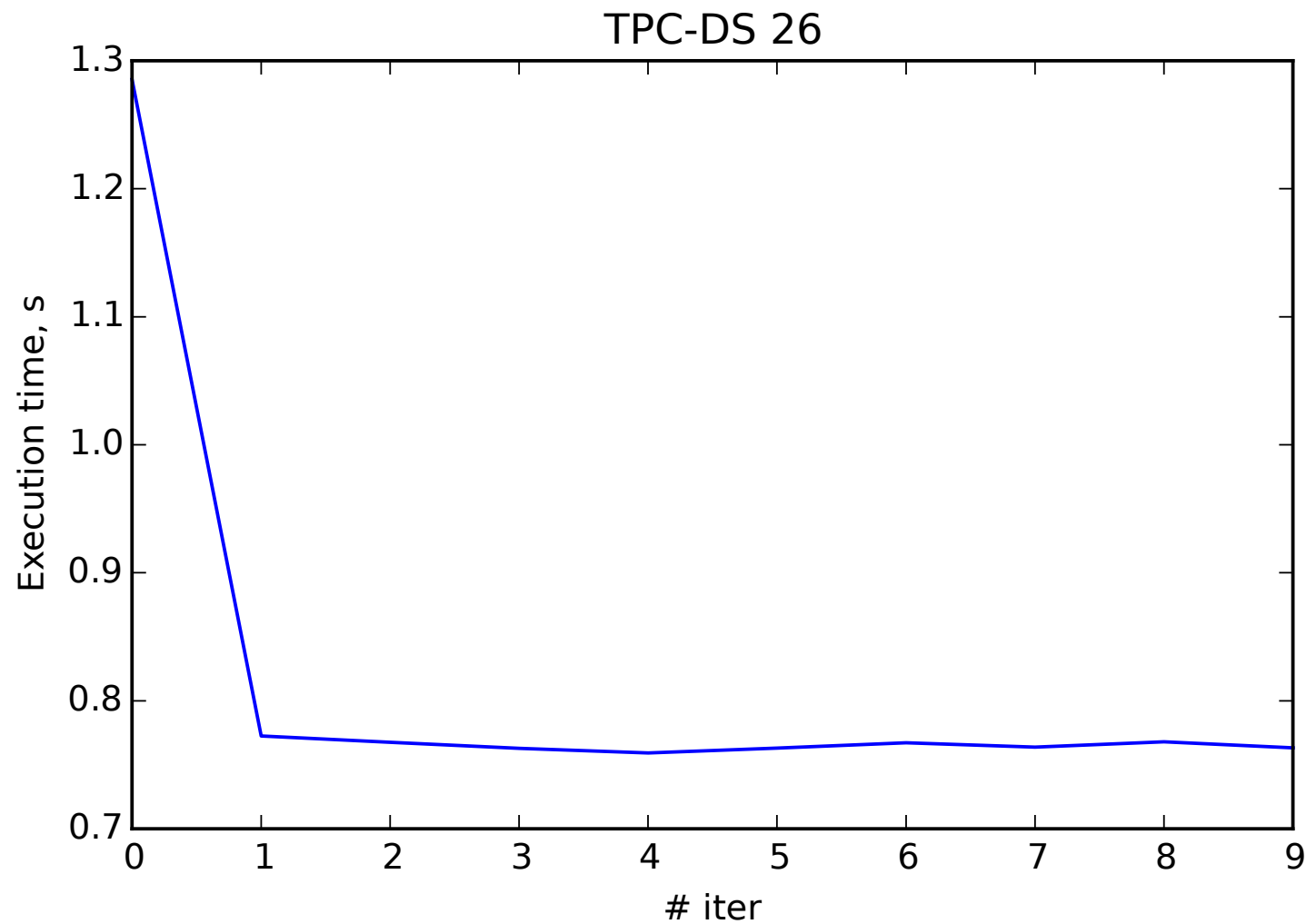


TPC-DS 18

# Learning progress

TPC-DS 69

# Learning progress

TPC-DS 26

# AQO: adaptive query optimization

Current code for vanilla PostgreSQL (extension + patch):
https://github.com/tigvarts/aqo

Available in Postgres Pro Enterprise

# AQO: adaptive query optimization

For some queries we don't need AQO.

So we need a mechanism to determine whether the query needs AQO.

# AQO: adaptive query optimization

*Query type* is the set of queries, which differ only in their constants.

Query type:
**SELECT** * **FROM** users **WHERE** age > **const AND** city = **const**;

Queries:
**SELECT** * **FROM** users **WHERE** age > 18 **AND** city = 'Moscow';
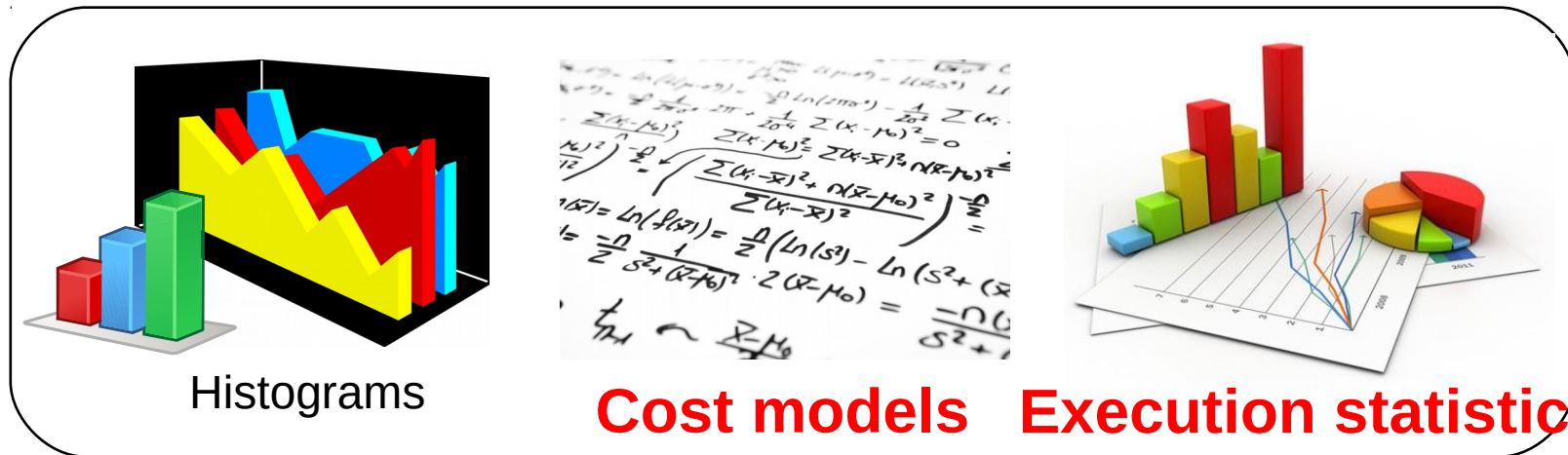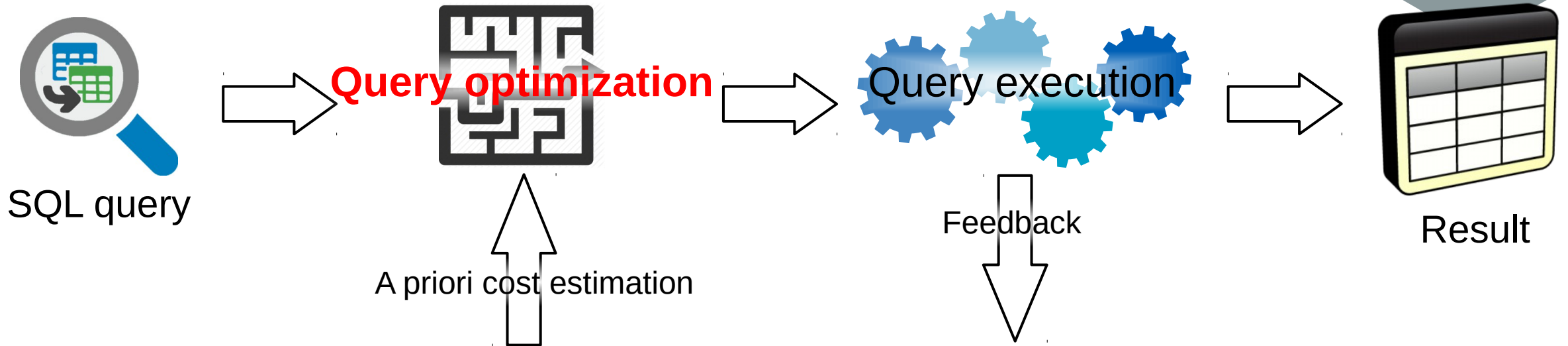**SELECT** * **FROM** users **WHERE** age > 65 **AND** city = 'Kostroma';
…

# AQO: adaptive query optimization

aqo.mode:

- Disabled for all query types
- Enabled for all query types
- Use manual settings for known query types, ignore others
- Use manual settings for known query types, tries to tune others automatically

# What is next?



SQL query → **Query optimization** → Query execution → Result

A priori cost estimation

Feedback

Histograms     **Cost models   Execution statistics**

# Questions

Contacts:
- o.ivanov@postgrespro.ru
- +7 (916) 377-55-63

# Postgres Professional

**http://postgrespro.ru/**

**+7 495 150 06 91**

**info@postgrespro.ru**

postgrespro.ru