# Clarinet: WAN-Aware Optimization for Analytics Queries

Raajay Viswanathan,  Ganesh Ananthanarayanan, Aditya Akella
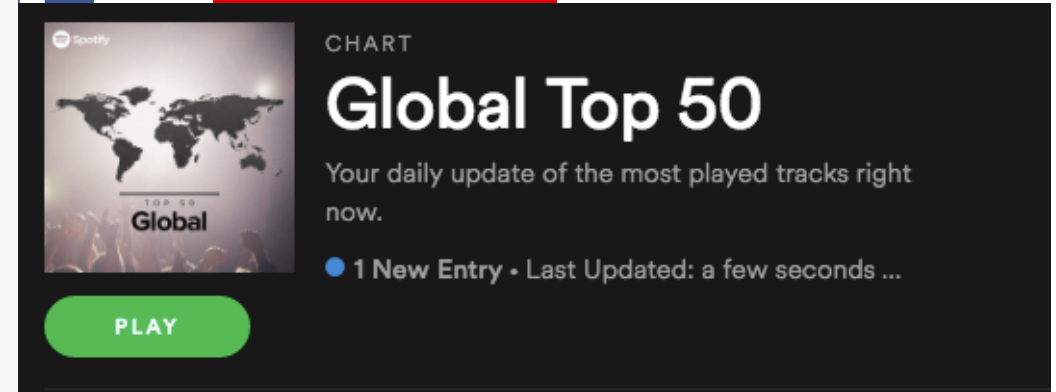
# Overview

- Web apps hosted on multiple DCs → Low latency access to end-user

# Overview



- Web apps hosted on multiple DCs → Low latency access to end-user

# Overview



- Web apps hosted on multiple DCs → Low latency access to end-user
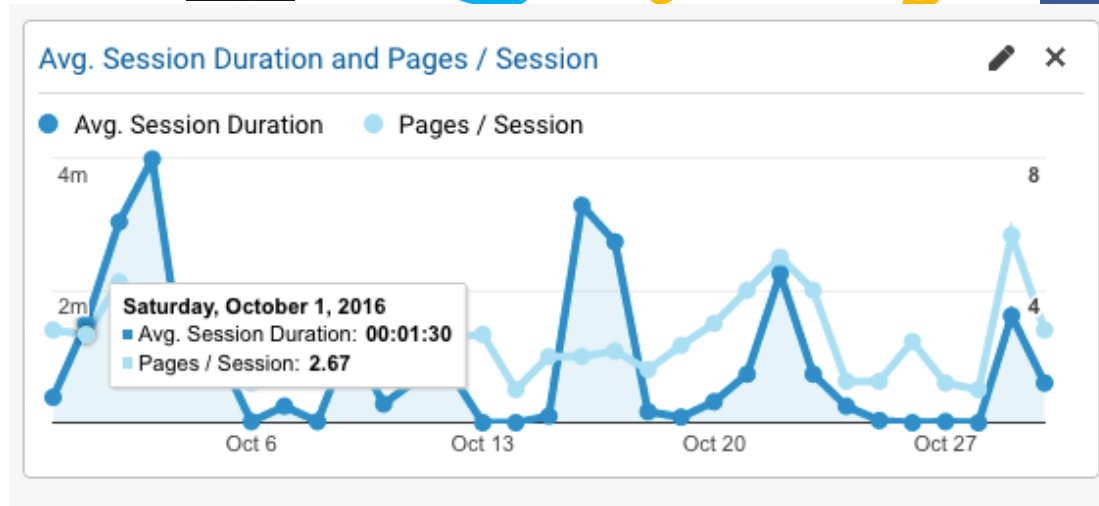
# Overview

- Web apps hosted on multiple DCs → Low latency access to end-user

# Overview
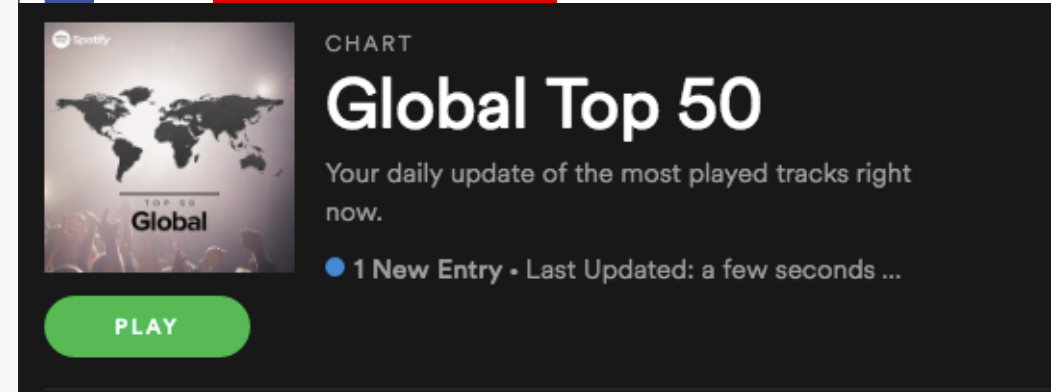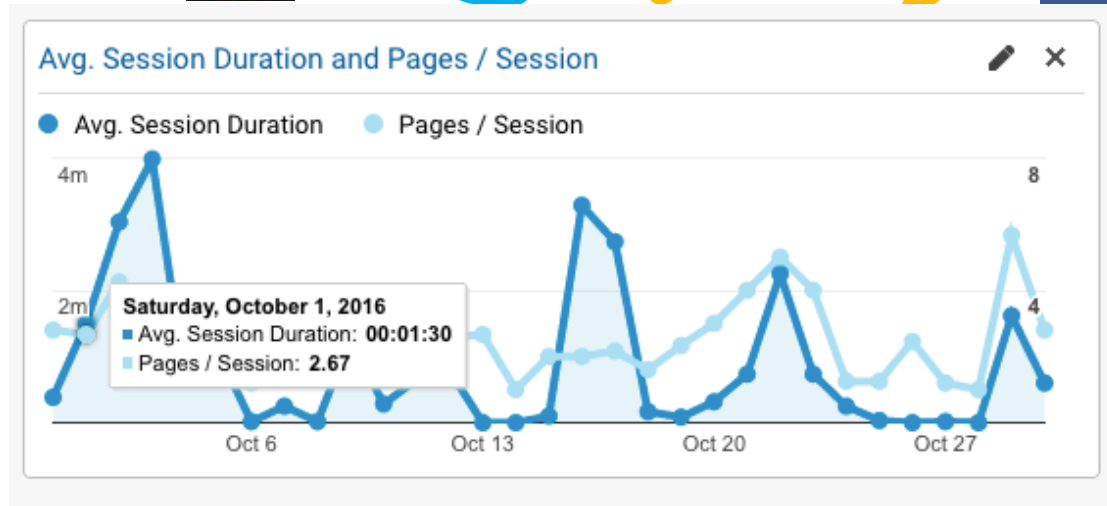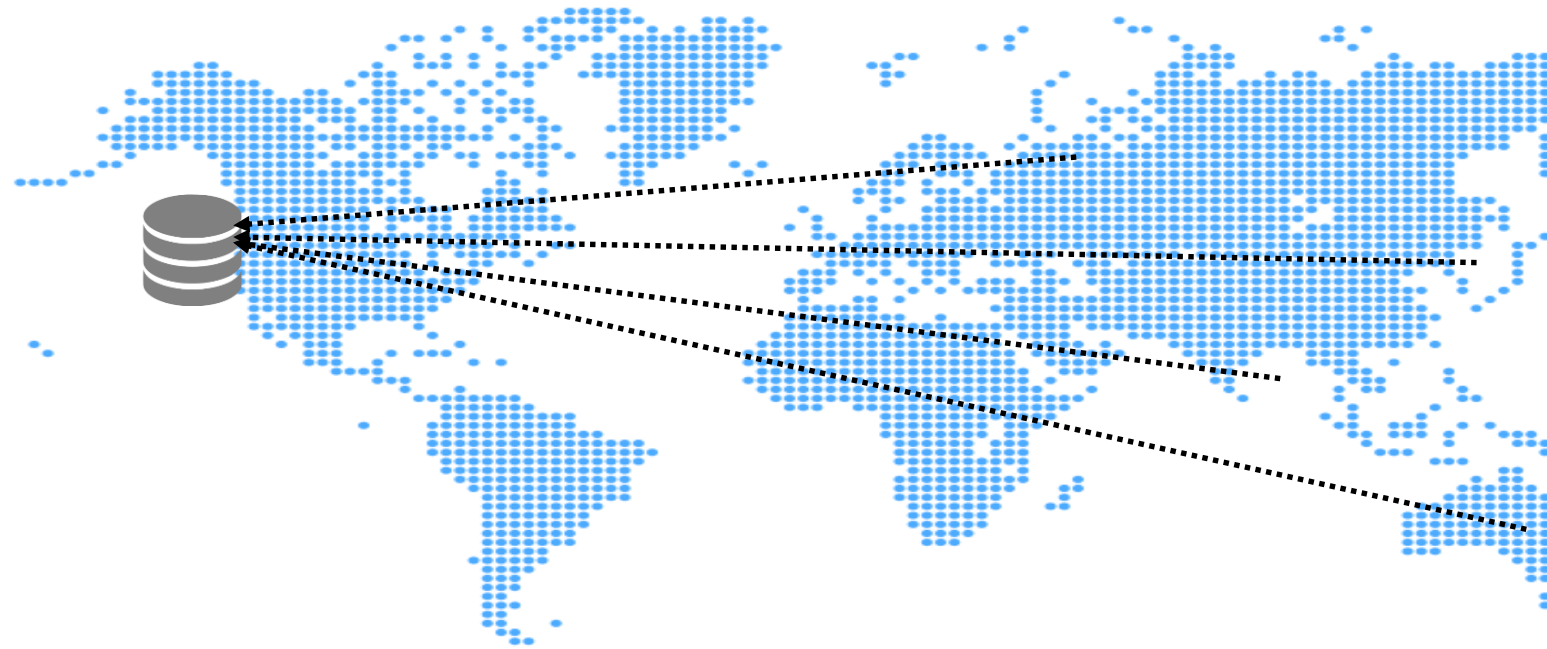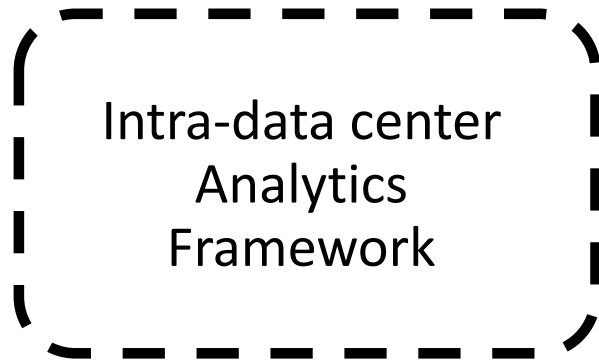
- Web apps hosted on multiple DCs → Low latency access to end-user
- Need efficient methods to analyze data located in **multiple** data centers

# Centralized Aggregation is Wasteful

# Centralized Aggregation is Wasteful

```
SELECT * … FROM .. WHERE .. ;
```

Intra-data center Analytics Framework

# Centralized Aggregation is Wasteful



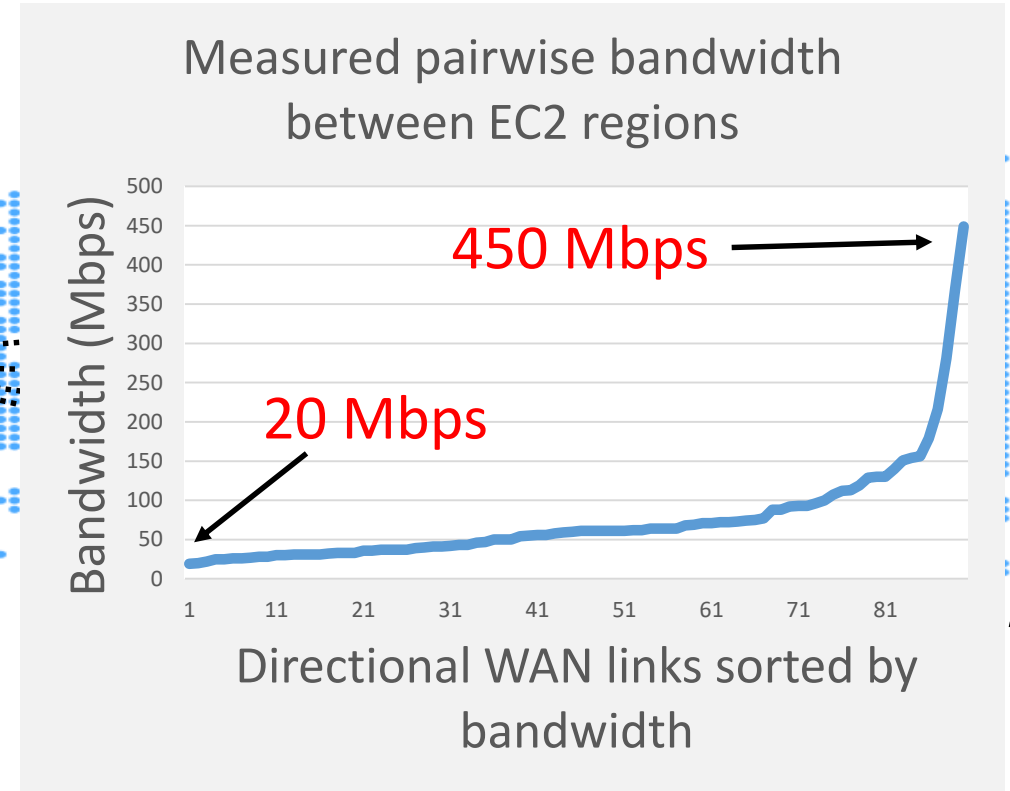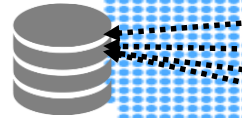- Available WAN bandwidth is limited → Aggregation latency overhead

# Centralized Aggregation is Wasteful

```
SELECT * … FROM .. WHERE .. ;
```

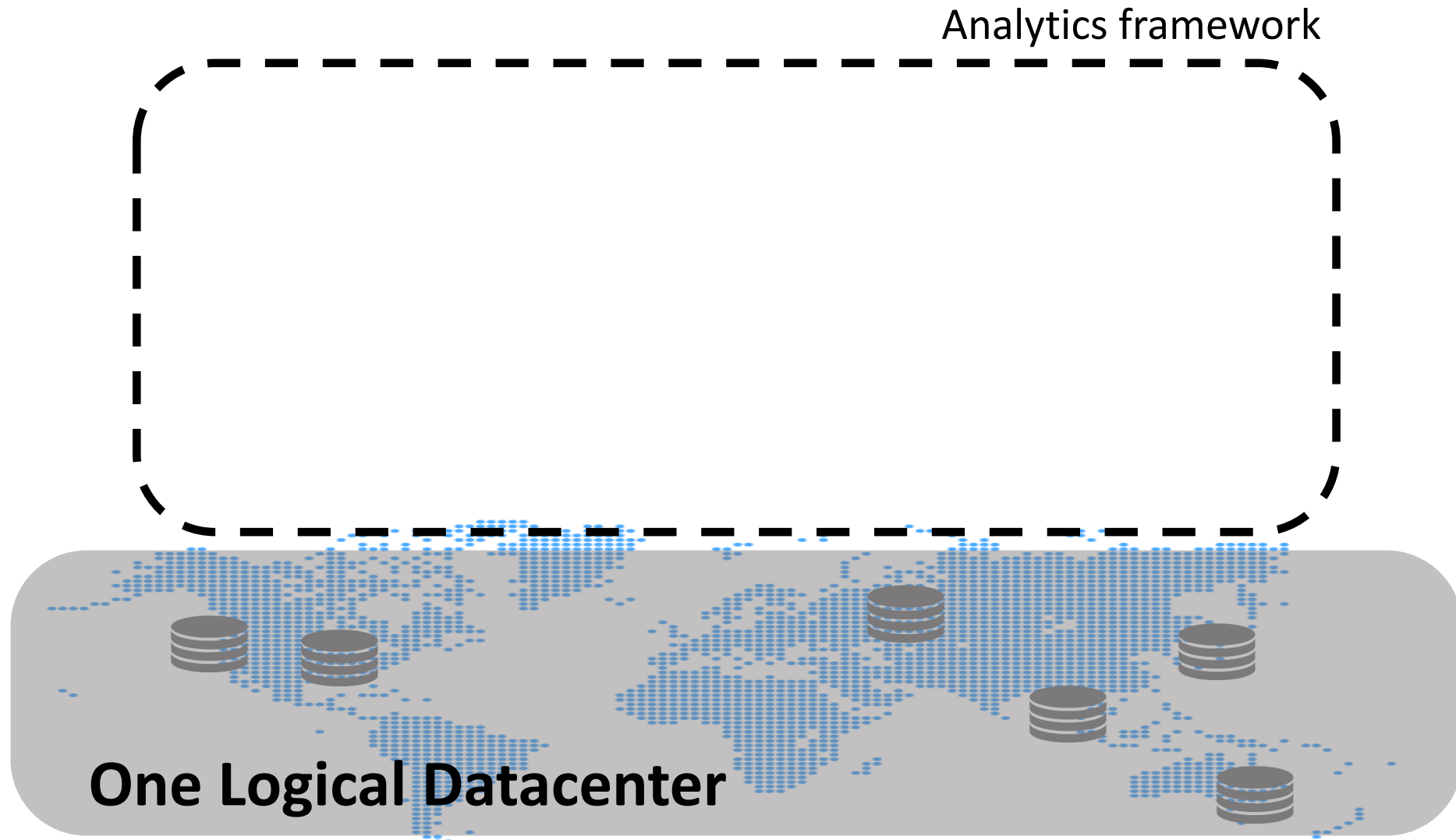Intra-data center Analytics Framework

$$$$

$$$

$$$$

- Available WAN bandwidth is limited → Aggregation latency overhead

- WAN links are expensive → High data transfer cost

# Geo-distributed Analytics

# Geo-distributed Analytics

Analytics framework

**One Logical Datacenter**

# Geo-distributed Analytics



Analytics framework

**Distributed Storage Layer**

**One Logical Datacenter**

4

# Geo-distributed Analytics

```
SELECT * … FROM .. WHERE .. ;
```

Analytics framework

**Query Optimizer**

*Multi-stage parallelizable jobs*

**Distributed Execution Layer**

**Distributed Storage Layer**

**One Logical Datacenter**

# Geo-distributed Analytics

`SELECT * … FROM .. WHERE .. ;`

Geo-distributed Analytics framework

**Query Optimizer**

*Multi-stage parallelizable jobs*

Requires WAN-aware optimization

**Distributed Execution Layer**

**Distributed Storage Layer**

**One Logical Datacenter**

4

# Geo-distributed Analytics



SELECT * … FROM .. WHERE .. ;

Geo-distributed Analytics framework

**Query Optimizer**

*Multi-stage parallelizable jobs*

Requires WAN-aware optimization

**Distributed Execution Layer**

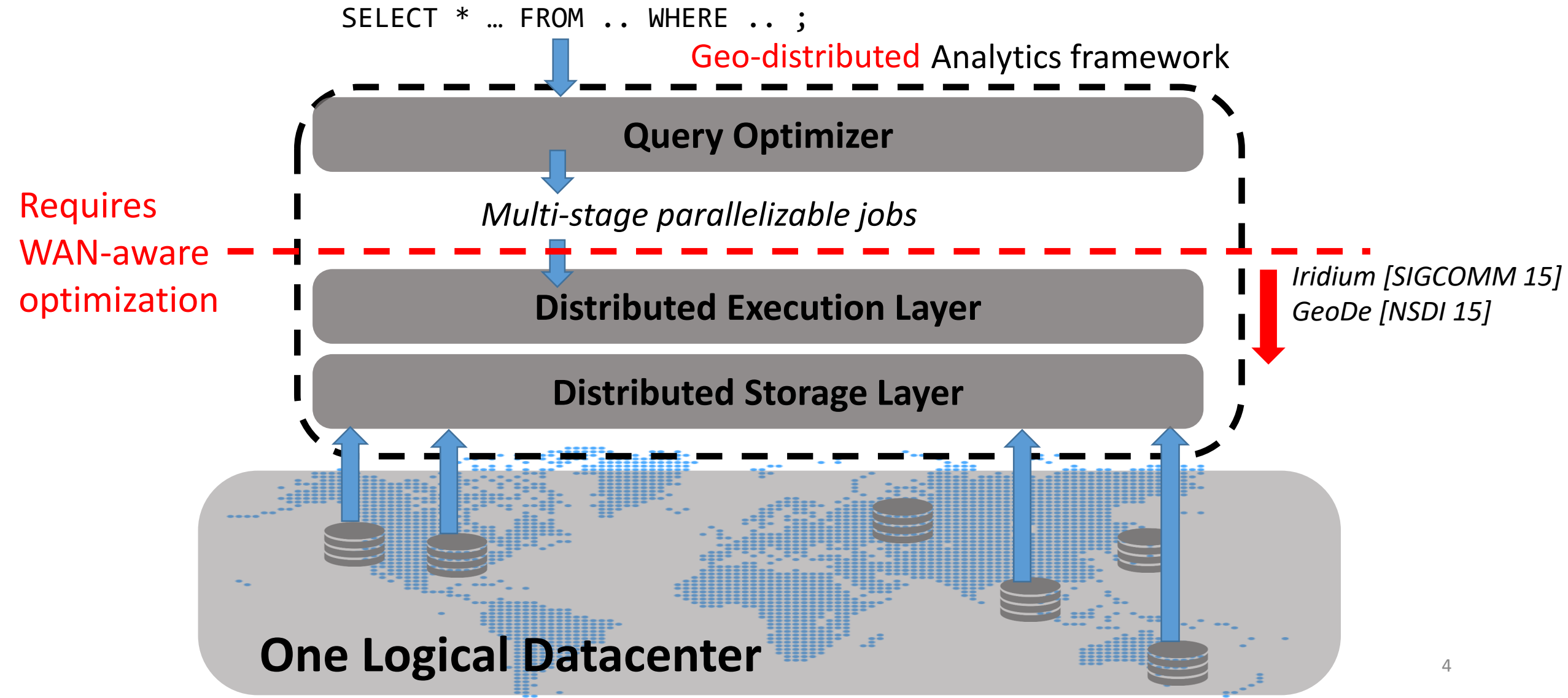*Iridium [SIGCOMM 15]*
*GeoDe [NSDI 15]*

**Distributed Storage Layer**

**One Logical Datacenter**

4

# Geo-distributed Analytics

# Geo-distributed Analytics

```
SELECT * … FROM .. WHERE .. ;
```

Geo-distributed Analytics framework

Requires WAN-aware optimization

**Query Optimizer**

*Multi-stage parallelizable jobs*

**Distributed Execution Layer**

**Distributed Storage Layer**

*2.7x* reduction in query runtime
***Clarinet***

**One Logical Datacenter**

4

# WAN Aware Query Optimization

DC₂

T2

DC₁

T1

DC₃

T3

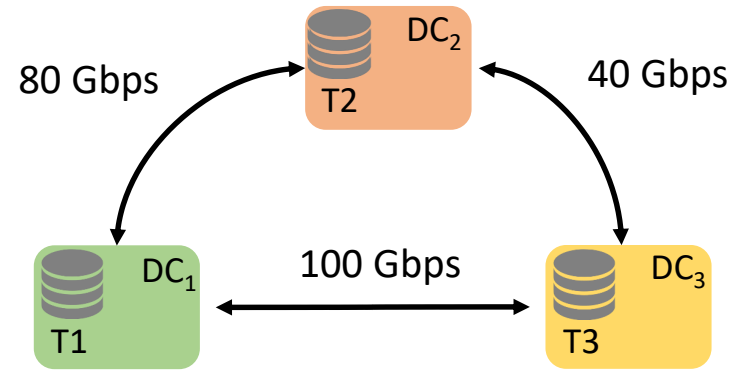*T1, T2, T3: Tables storing click logs*

# WAN Aware Query Optimization



T1, T2, T3: Tables storing click logs

# WAN Aware Query Optimization
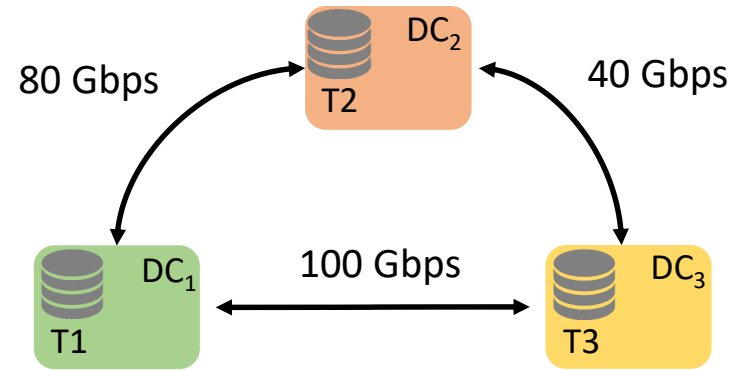


80 Gbps

40 Gbps

100 Gbps

*T1, T2, T3: Tables storing click logs*

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM  T1, T2, T3
WHERE T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```

# WAN Aware Query Optimization



80 Gbps

40 Gbps

100 Gbps

*T1, T2, T3: Tables storing click logs*

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM  T1, T2, T3
WHERE T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```
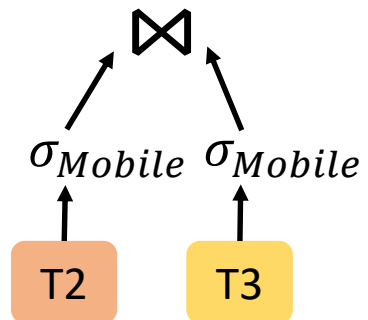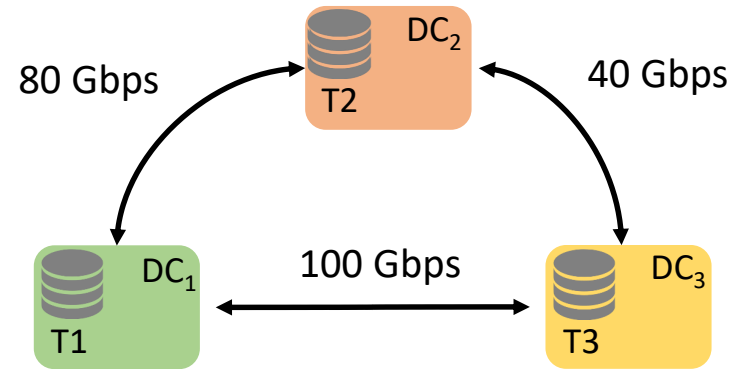
# WAN Aware Query Optimization



T1, T2, T3: Tables storing click logs

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM  T1, T2, T3
WHERE T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```
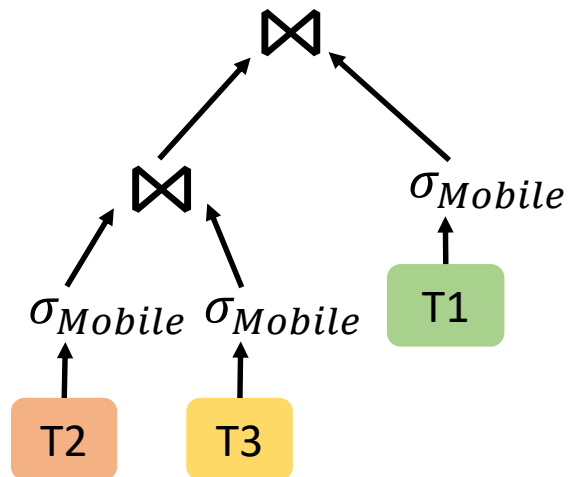
# WAN Aware Query Optimization



T1, T2, T3: Tables storing click logs

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM  T1, T2, T3
WHERE T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```
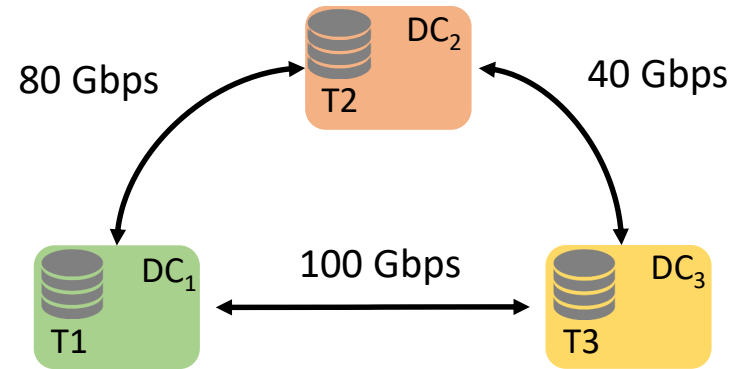
# WAN Aware Query Optimization



T1, T2, T3: Tables storing click logs

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM   T1, T2, T3
WHERE T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```
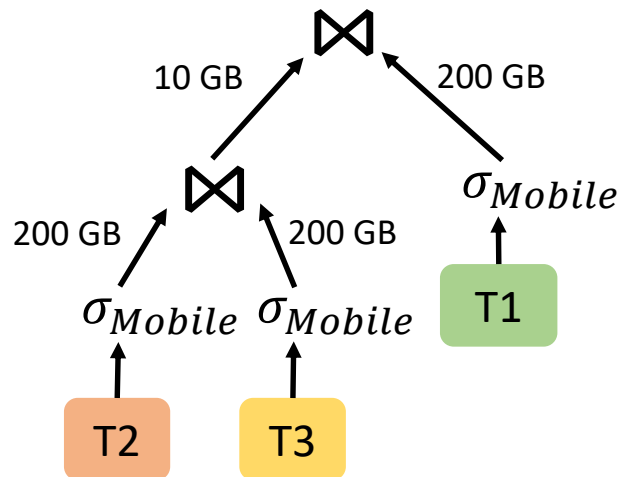
Plan running time: 41 s

# WAN Aware Query Optimization



T1, T2, T3: Tables storing click logs

```
QUERY
SELECT  T1.user, T1.latency, T2.latency, T3.latency
FROM    T1, T2, T3
WHERE  T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```
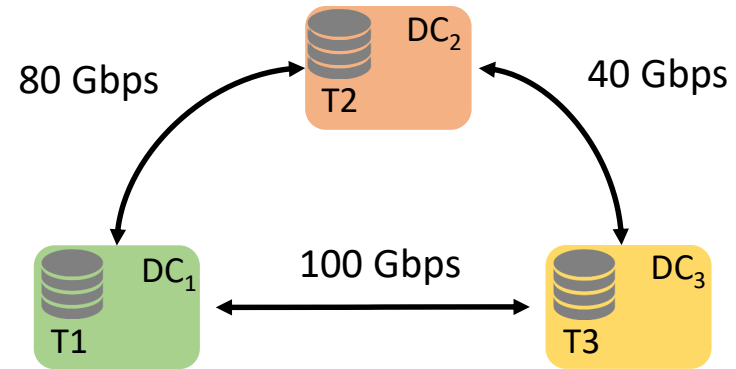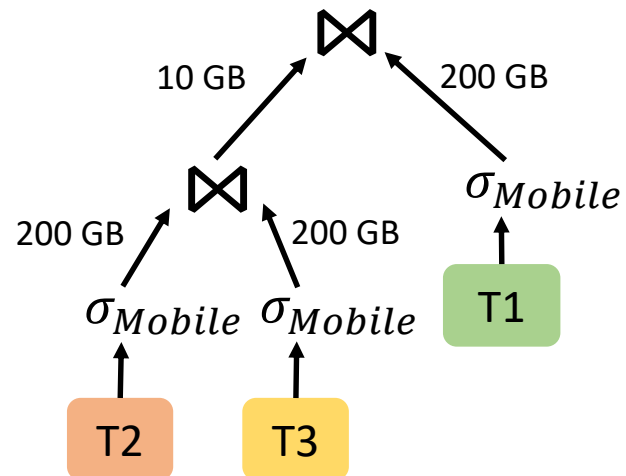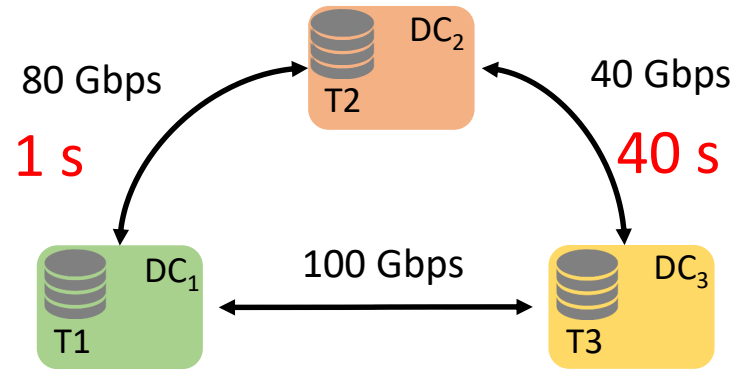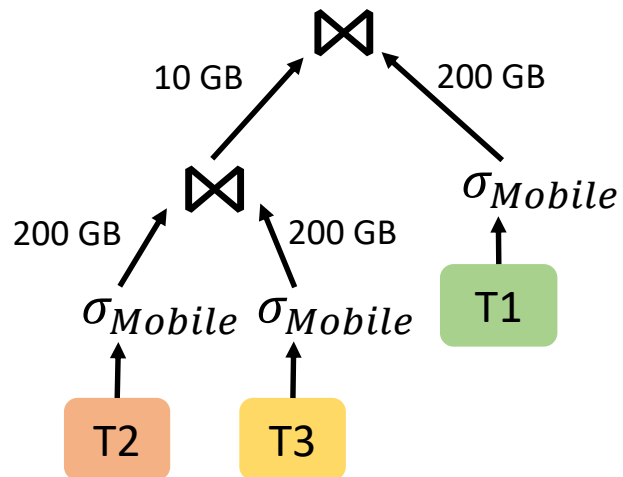
# WAN Aware Query Optimization



80 Gbps

**1 s**

DC₂ — T2

**WAN-only bottleneck**

40 Gbps

**40 s**

DC₁ — T1

100 Gbps

DC₃ — T3

*T1, T2, T3: Tables storing click logs*

Plan running time: **41 s**

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM   T1, T2, T3
WHERE  T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```

5

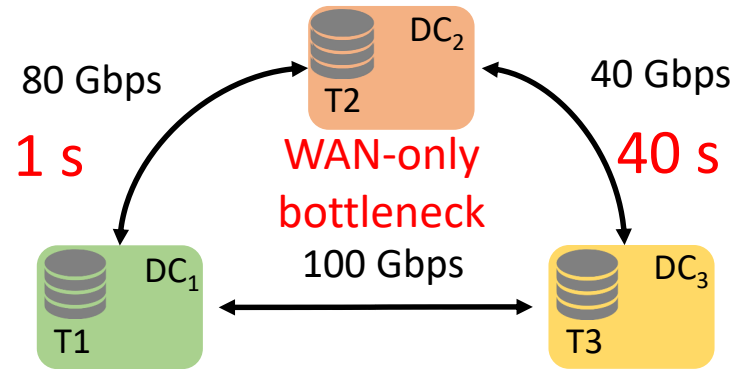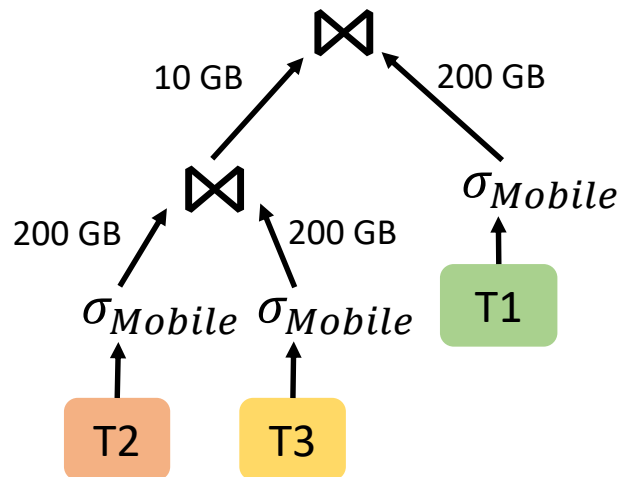# WAN Aware Query Optimization



T1, T2, T3: Tables storing click logs

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM  T1, T2, T3
WHERE T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```

# WAN Aware Query Optimization



80 Gbps
1 s
40 Gbps
40 s
WAN-only bottleneck
100 Gbps

DC₂ — T2
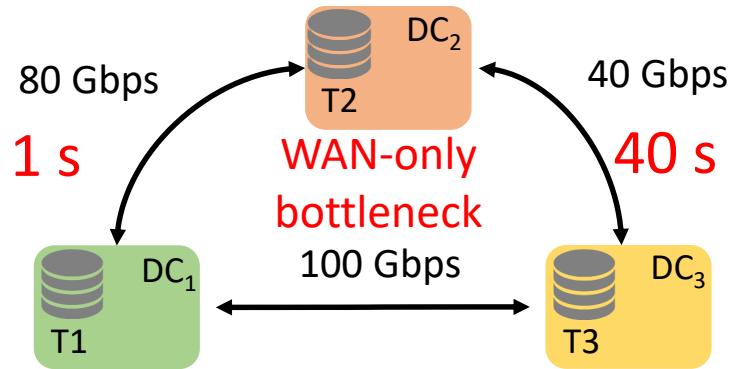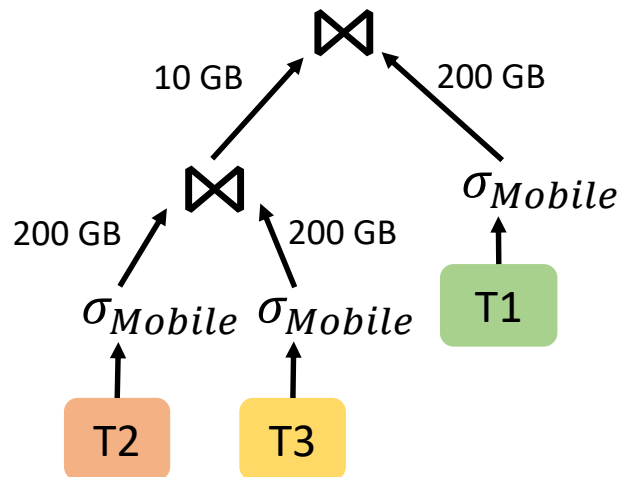DC₁ — T1
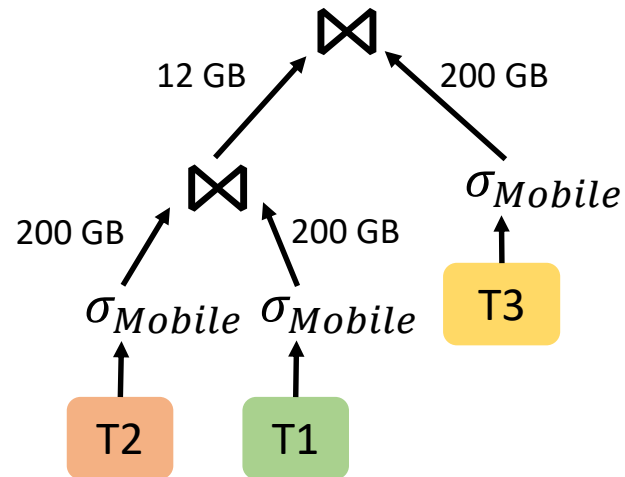DC₃ — T3

T1, T2, T3: Tables storing click logs

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM   T1, T2, T3
WHERE  T1.user == T2.user AND T1.user == T3.user
AND T1.device == T2.device == T3.device == "mobile";
```
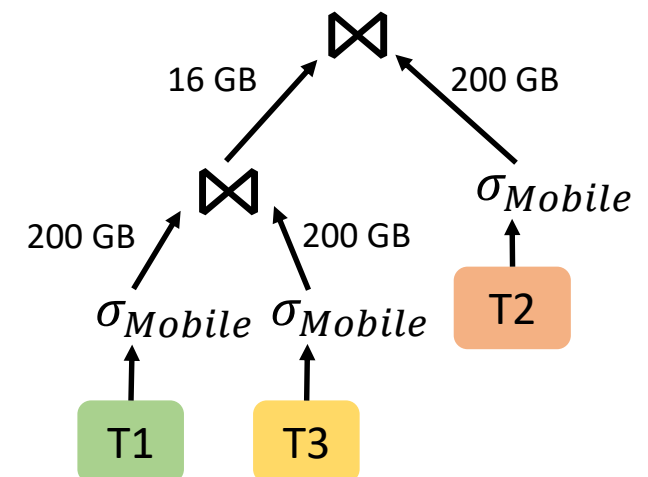
Plan running time: 41 s

10 GB
200 GB
200 GB
200 GB
σMobile
σMobile σMobile
T2   T3   T1

Chosen by *network agnostic* query optimizer

Plan A

Plan running time: 20.96 s

12 GB
200 GB
200 GB
200 GB
σMobile
σMobile σMobile
T2   T1   T3

Plan B

Plan running time: 17.6 s

16 GB
200 GB
200 GB
200 GB
σMobile
σMobile σMobile
T1   T3   T2

Plan C

5

# WAN Aware Query Optimization

# WAN Aware Query Optimization



DC₂ — T2 — WAN-only bottleneck

80 Gbps — **1 s**
40 Gbps — **40 s**

DC₁ — T1
100 Gbps
DC₃ — T3

```
QUERY
SELECT T1.user, T1.latency, T2.latency, T3.latency
FROM   T1, T2, T3
WHERE  T1.user == T2.user AND T1.user == T3.user
AND    T1.device == T2.device == T3.device == "mobile";
```
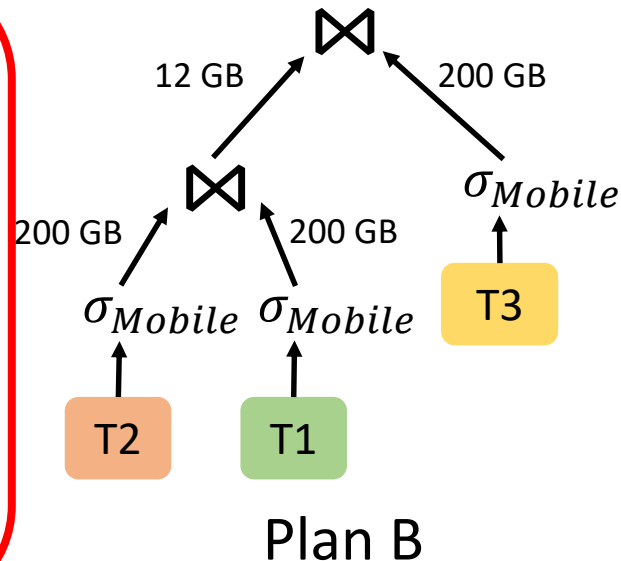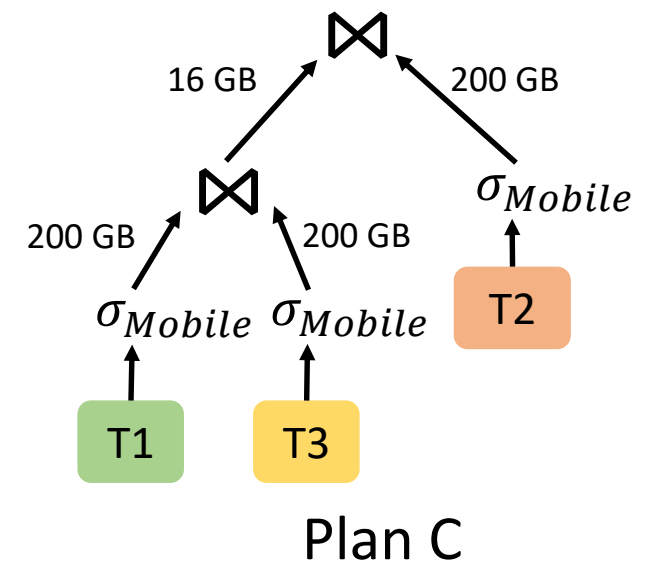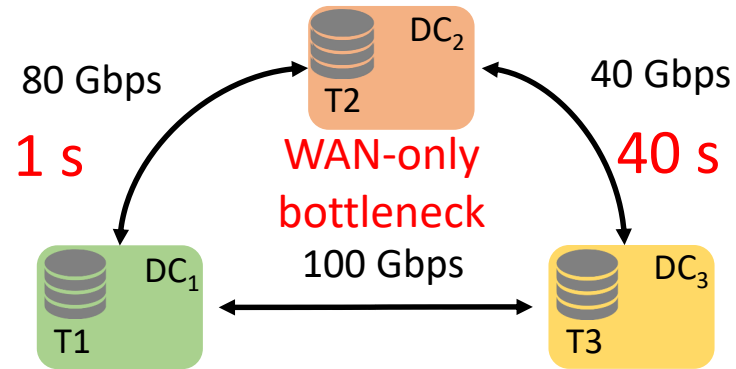
**WAN-aware query optimizer** that uses network **transfer duration** to choose query plans

**Plan A**
⋈
200 GB   200 GB
$\sigma_{Mobile}$   $\sigma_{Mobile}$   $\sigma_{Mobile}$
T2   T3   T1
*Chosen by network agnostic query optimizer*

**Plan B**
⋈
200 GB   200 GB
$\sigma_{Mobile}$   $\sigma_{Mobile}$   $\sigma_{Mobile}$
T2   T1   T3

**Plan C**
⋈
200 GB   200 GB
$\sigma_{Mobile}$   $\sigma_{Mobile}$   $\sigma_{Mobile}$
T1   T3   T2

5

# Outline

1. Motivation

2. Challenges in choosing query plan based on WAN transfer durations

3. Solution
   - Single query
   - Multiple simultaneous queries

4. Experimental Evaluation

# Other factors also affect query plan run time

# Other factors also affect query plan run time

# Other factors also affect query plan run time



Map Reduce Job

# Other factors also affect query plan run time



Map Reduce Job

# Other factors also affect query plan run time

# Other factors also affect query plan run time



Map Reduce Job

While evaluating different query plans

1. Plan A:      41 s
2. Plan B:      20.96
3. Plan C:      17.6 s

# Other factors also affect query plan run time



Tasks are placed uniformly across $DC_1$ and $DC_2$

Map Reduce Job

While evaluating different query plans

1. Plan A:    ~~41 s~~    20.5 s
2. Plan B:    ~~20.96~~    11.2 s
3. Plan C:    17.6 s

# Other factors also affect query plan run time



Tasks are placed uniformly across $DC_1$ and $DC_2$

Map Reduce Job

While evaluating different query plans

1. Plan A: ~~41 s~~ 20.5 s
2. Plan B: ~~20.96~~ 11.2 s
3. Plan C: 17.6 s

# Other factors also affect query plan run time

While evaluating different query plans

1. Plan A:    ~~41 s~~    20.5 s
2. Plan B:    ~~20.96~~   11.2 s
3. Plan C:    17.6 s

# Other factors also affect query plan run time

REDUCE: JOIN

Tasks are placed uniformly across $DC_1$ and $DC_2$

DC$_2$

T2

40 Gbps

80 Gbps

200 GB    200 GB

$\sigma_C$    $\sigma_C$

200 GB    200 GB

MAP: SELECT    MAP: SELECT

## Choose query plan based on:
## 1.    Best available task placements

W

1.    Plan A:    ~~41 s~~    20.5 s
2.    Plan B:    ~~20.96~~    11.2 s
3.    Plan C:    17.6 s

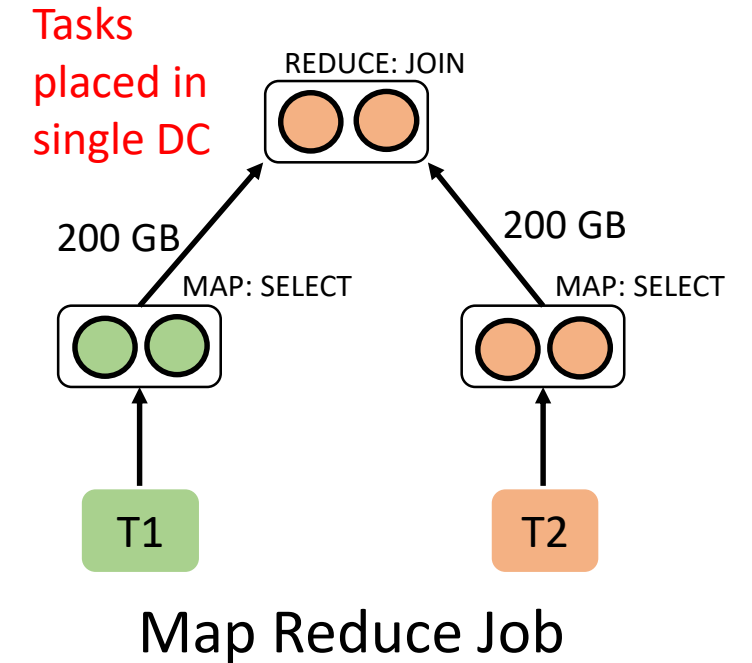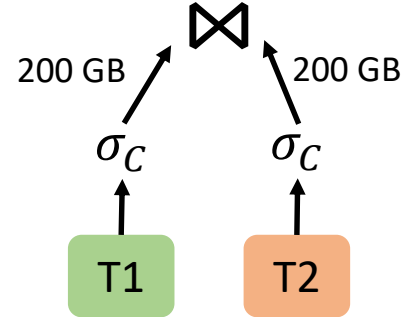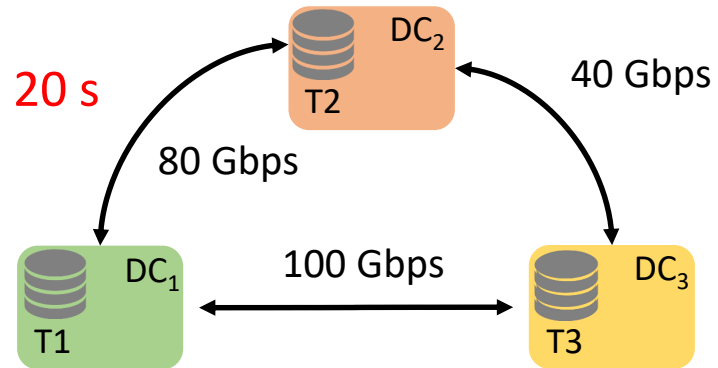# Other factors also affect query plan run time



Tasks are placed uniformly across $DC_1$ and $DC_2$

REDUCE: JOIN

40 Gbps

80 Gbps

200 GB    200 GB

$\sigma_C$    $\sigma_C$

200 GB    200 GB

MAP: SELECT    MAP: SELECT

$DC_2$

T2

**Choose query plan based on:**
1. **Best available task placements**
2. **Schedule of network transfers**

W

1.  Plan A:    ~~41 s~~    20.5 s
2.  Plan B:    ~~20.96~~    11.2 s
3.  Plan C:    17.6 s

# Joint plan selection, placement and scheduling

# Joint plan selection, placement and scheduling

```
SELECT * FROM … WHERE.. ;
```

Query Optimizer

**Multiple** query plans (join orders) per query

# Joint plan selection, placement and scheduling

SELECT * FROM … WHERE.. ;

Query Optimizer

Logical plan to physical plan

**Multiple** query plans (join orders) per query

Assign parallelism for each stage

# Joint plan selection, placement and scheduling

```
SELECT * FROM … WHERE.. ;
```



**Multiple** query plans (join orders) per query

Assign parallelism for each stage

Network aware task placement and scheduling for each query plan

Choose plan with smallest run time for execution

# Joint plan selection, placement and scheduling

```
SELECT * FROM … WHERE.. ;
```

Query Optimizer

**Multiple** query plans (join orders) per query

Clarinet binds query to plan lower in the stack

Clarinet

Network aware task placement and scheduling for each query plan

Choose plan with smallest run time for execution

8

# Network aware placement and scheduling

# Network aware placement and scheduling

- Task placement decided greedily one stage at a time
  - Minimize per stage run time

# Network aware placement and scheduling

- Task placement decided greedily one stage at a time
  - Minimize per stage run time

- Scheduling of network transfers
  - Determines start times of inter-DC network transfers

# Network aware placement and scheduling

- Task placement decided greedily one stage at a time
  - Minimize per stage run time

- Scheduling of network transfers
  - Determines start times of inter-DC network transfers
  - Formulate a Binary Integer Linear Program to solve scheduling
  - Factors transfer dependencies

# How to extend the late-binding strategy to multiple queries?

# Queries affect each others' run time

# Queries affect each others' run time



**QUERY 1**
```
SELECT …
device == "mobile"
…;
```

**QUERY 2**
```
SELECT …
genre == "pc"
…;
```

# Queries affect each others' run time



Same query plan (Plan C) for Query 1 and Query 2

**QUERY 1**
SELECT …
device == "mobile"
…;

**QUERY 2**
SELECT …
genre == "pc"
…;

# Queries affect each others' run time



80 Gbps

40 Gbps

100 Gbps

Contention increases query run time

16 GB          200 GB

$\sigma_C$

200 GB          200 GB

$\sigma_{Mobile}$   $\sigma_{Mobile}$

T1          T3

T2

16 GB          200 GB

$\sigma_R$

200 GB          200 GB

$\sigma_{PC}$   $\sigma_{PC}$

T1          T3

T2

Same query plan (Plan C) for Query 1 and Query 2

**QUERY 1**
SELECT …
device == "mobile"
…;

**QUERY 2**
SELECT …
genre == "pc"
…;

11

# Queries affect each others' run time

# Queries affect each others' run time



80 Gbps

40 Gbps

100 Gbps

No contention of network links

**QUERY 1**
```
SELECT …
device == "mobile"
…;
```

**QUERY 2**
```
SELECT …
genre == "pc"
…;
```

Different query plans for Query 1 (Plan C) and Query 2 (Plan B)

Plan C:

⋈

16 GB    200 GB

⋈ (200 GB, 200 GB)     $\sigma_C$

$\sigma_{Mobile}$  $\sigma_{Mobile}$    T2

T1    T3

Plan B:

⋈

12 GB    200 GB

⋈ (200 GB, 200 GB)     $\sigma_C$

$\sigma_{PC}$  $\sigma_{PC}$    T3

T2    T1

# Queries affect each others' run time

# Iterative Shortest Job First

QUERY A

QUERY B

QUERY C



- Best combination → minimize average completion
  - Computationally intractable

# Iterative Shortest Job First



- Best combination → minimize average completion
  - Computationally intractable

- Iterative Shortest Job First (SJF) scheduling heuristic
  1. Pick shortest physical query plan in each iteration

# Iterative Shortest Job First



QUERY A

QUERY B

QUERY C

QO

QO

QO

Iter 1: 10    18    12     5      8     20    30

Clarinet

- Best combination → minimize average completion
  - Computationally intractable

- Iterative Shortest Job First (SJF) scheduling heuristic
  1. Pick shortest physical query plan in each iteration

# Iterative Shortest Job First



QUERY A    QUERY B    QUERY C

QO    QO    QO

Iter 1:   10   18   12    5   8    20   30

Clarinet

- Best combination → minimize average completion
  - Computationally intractable

- Iterative Shortest Job First (SJF) scheduling heuristic
  1. Pick shortest physical query plan in each iteration

# Iterative Shortest Job First

QUERY A        QUERY B        QUERY C

QO             QO             QO

Iter 1:  10    18    12      5      8      20      30

Clarinet

Link 1 [ B1 ]

Link 2

0       5       t

- Best combination → minimize average completion
  - Computationally intractable

- Iterative Shortest Job First (SJF) scheduling heuristic
  1. Pick shortest physical query plan in each iteration

- Reserve bandwidth to guarantee completion time

# Iterative Shortest Job First



- Best combination → minimize average completion
  - Computationally intractable

- Iterative Shortest Job First (SJF) scheduling heuristic
  1. Pick shortest physical query plan in each iteration

- Reserve bandwidth to guarantee completion time

# Iterative Shortest Job First



QUERY A      QUERY B      QUERY C

QO    QO    QO

Iter 2:   15    18    17      25    30

Clarinet

Link 1   B1

Link 2

0      5      t

- Best combination → minimize average completion
  - Computationally intractable

- Iterative Shortest Job First (SJF) scheduling heuristic
  1. Pick shortest physical query plan in each iteration

- Reserve bandwidth to guarantee completion time

12

# Iterative Shortest Job First



QUERY A  QUERY B  QUERY C

QO   QO   QO

Iter 2:  15    18    17        25    30

Clarinet

Link 1   B1

Link 2

0        5              t

- Best combination → minimize average completion
  - Computationally intractable

- Iterative Shortest Job First (SJF) scheduling heuristic
  1. Pick shortest physical query plan in each iteration

- Reserve bandwidth to guarantee completion time

# Iterative Shortest Job First



QUERY A

QUERY B

QUERY C

QO

QO

QO

Iter 2:   15      18      17                    25      30

Clarinet

Link 1

Link 2

B1   A1

A2

0        5    7              15   t

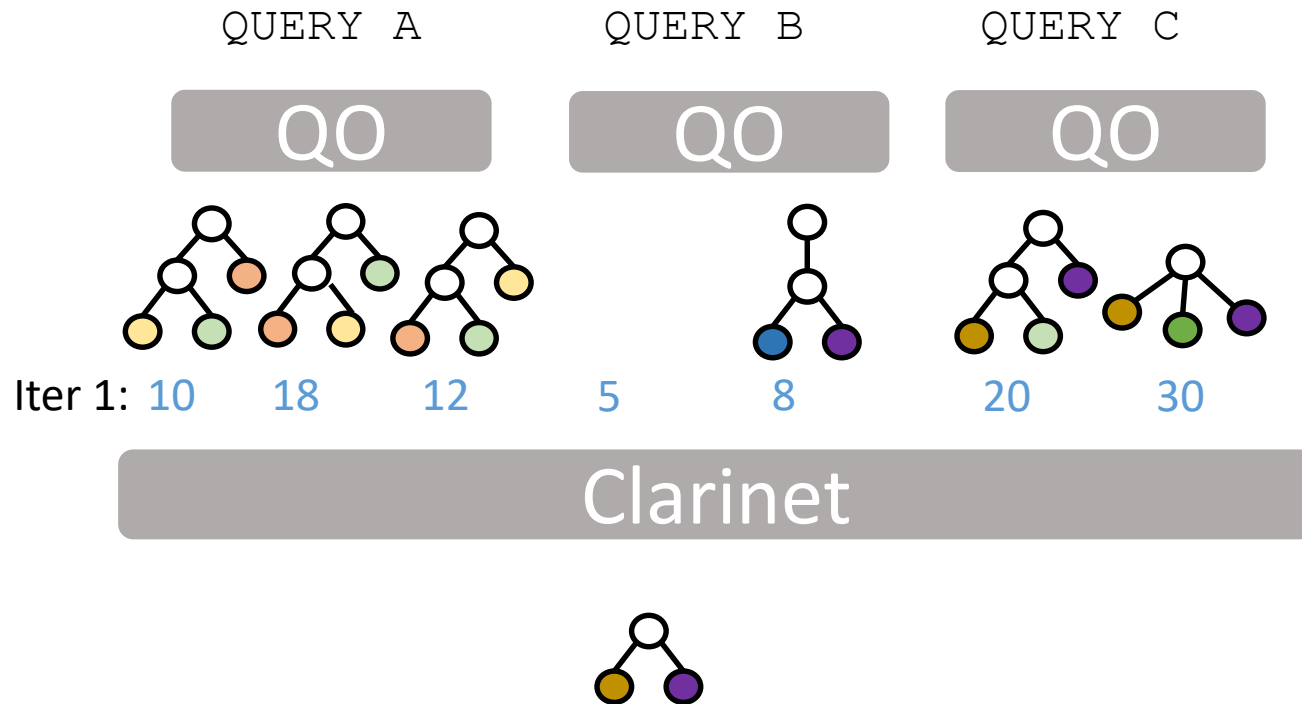- Best combination → minimize average completion
  - Computationally intractable

- Iterative Shortest Job First (SJF) scheduling heuristic
  1. Pick shortest physical query plan in each iteration

- Reserve bandwidth to guarantee completion time

# Avoid fragmentation and improve completion time

# Avoid fragmentation and improve completion time

- SJF & reservation leads to bandwidth fragmentation

# Avoid fragmentation and improve completion time

- SJF & reservation leads to bandwidth fragmentation



Scheduled in SJF order

# Avoid fragmentation and improve completion time

- SJF & reservation leads to bandwidth fragmentation

Scheduled in SJF order

Dominant transfers execute sequentially

Link 1

Link 2

B1   A1

A2

0          10  12          22   t

# Avoid fragmentation and improve completion time

- SJF & reservation leads to bandwidth fragmentation

Scheduled in SJF order

Dominant transfers execute sequentially

# Avoid fragmentation and improve completion time

- SJF & reservation leads to bandwidth fragmentation



Scheduled in SJF order

Alternate schedule with same query plans

# Avoid fragmentation and improve completion time

- SJF & reservation leads to bandwidth fragmentation

Scheduled in SJF order          Alternate schedule with same query plans

Li

Link 2    Extended Idling

0        10   12        22   t        0   2        12    t

**Re-arranging transfers resulting in deviation from SJF schedule can help**

# k-Shortest Jobs First Heuristic

Offline schedule

# k-Shortest Jobs First Heuristic



Offline schedule

Link *n*

Link 2

Link 1

t

- Identify transfers of k-shortest yet incomplete jobs

# k-Shortest Jobs First Heuristic



Offline schedule

- Identify transfers of k-shortest yet incomplete jobs
- Relax transfer schedule → Start as soon as link is free and task is available

14

# k-Shortest Jobs First Heuristic



Offline schedule

Link *n*

Link 2

Link 1

t

- Identify transfers of k-shortest yet incomplete jobs
- Relax transfer schedule → Start as soon as link is free and task is available
-  Best 'k' ←Prior observations (or) through offline simulations

# Clarinet Implementation



QUERY 1     QUERY 2     QUERY 3

QO     QO     QO

Batch of queries

Existing Query Optimizers

# Clarinet Implementation

QUERY 1          QUERY 2          QUERY 3

QO               QO               QO

Batch of queries

Existing Query Optimizers
- Modified Hive to generate multiple plans

# Clarinet Implementation

QUERY 1          QUERY 2          QUERY 3

Batch of queries

Existing Query Optimizers
- Modified Hive to generate multiple plans
- QOs control set of generated plans
- Existing optimizations are applied
  - Push down Select
  - Partition pruning

# Clarinet Implementation



QUERY 1    QUERY 2    QUERY 3

QO    QO    QO

Clarinet

Execution framework

Batch of queries

Existing Query Optimizers
- Modified Hive to generate multiple plans
- QOs control set of generated plans
- Existing optimizations are applied
    - Push down Select
    - Partition pruning

Enforces Clarinet's schedule

# Clarinet Implementation



Batch of queries

Existing Query Optimizers
- Modified Hive to generate multiple plans
- QOs control set of generated plans
- Existing optimizations are applied
  - Push down Select
  - Partition pruning

Enforces Clarinet's schedule
- Modified Tez's DAGScheduler

# Clarinet Implementation

QUERY 1  QUERY 2  QUERY 3

QO    QO    QO

Clarinet

Execution framework

Batch of queries
Online query arrivals

Existing Query Optimizers
- Modified Hive to generate multiple plans
- QOs control set of generated plans
- Existing optimizations are applied
  - Push down Select
  - Partition pruning

Enforces Clarinet's schedule
- Modified Tez's DAGScheduler
- Fairness guarantees

# Evaluation

Compare Clarinet with following GDA approaches:

# Evaluation

Compare Clarinet with following GDA approaches:

1. Hive
2. Hive + Iridium
3. Hive + Reducers in single DC

# Evaluation

Compare Clarinet with following GDA approaches:

1. Hive                                    : <span style="color:red">WAN agnostic</span> task placement + scheduling
2. Hive + Iridium
3. Hive + Reducers in single DC

# Evaluation

Compare Clarinet with following GDA approaches:

1. Hive : WAN agnostic task placement + scheduling
2. Hive + Iridium : WAN aware task placement across DCs
3. Hive + Reducers in single DC

# Evaluation

Compare Clarinet with following GDA approaches:

1. Hive                                    : WAN agnostic task placement + scheduling
2. Hive + Iridium                     : WAN aware task placement across DCs
3. Hive + Reducers in single DC   : Distributed filtering + central aggregation

# Evaluation

Compare Clarinet with following GDA approaches:

1. Hive                          : WAN agnostic task placement + scheduling
2. Hive + Iridium              : WAN aware task placement across DCs
3. Hive + Reducers in single DC    : Distributed filtering + central aggregation

- Geo-Distributed Analytics stack across 10 EC2 regions

# Evaluation

Compare Clarinet with following GDA approaches:

1. Hive : WAN agnostic task placement + scheduling
2. Hive + Iridium : WAN aware task placement across DCs
3. Hive + Reducers in single DC : Distributed filtering + central aggregation

- Geo-Distributed Analytics stack across 10 EC2 regions

- Workload:
  - 30 batches of 12 randomly chosen TPC-DS queries

# Evaluation: Reduction in average completion time

| GDA Approach Vs. Hive | Average Gains |
|---|---|
| Clarinet | 2.7x |
| Hive + Iridium | 1.5x |

# Evaluation: Reduction in average completion time

| GDA Approach Vs. Hive | Average Gains |
|---|---|
| Clarinet | 2.7x |
| Hive + Iridium | 1.5x |

Clarinet chooses a different plan for **75%** of queries

# Evaluation: Reduction in average completion time

| GDA Approach Vs. Hive | Average Gains |
|---|---|
| Clarinet | 2.7x |
| Hive + Iridium | 1.5x |



**Data from a single batch 12 queries**

Clarinet chooses a different plan for **75%** of queries

# Evaluation: Reduction in average completion time

| GDA Approach Vs. Hive | Average Gains |
|---|---|
| Clarinet | 2.7x |
| Hive + Iridium | 1.5x |



**Data from a single batch 12 queries**

Clarinet chooses a different plan for **75%** of queries

# Evaluation: Reduction in average completion time

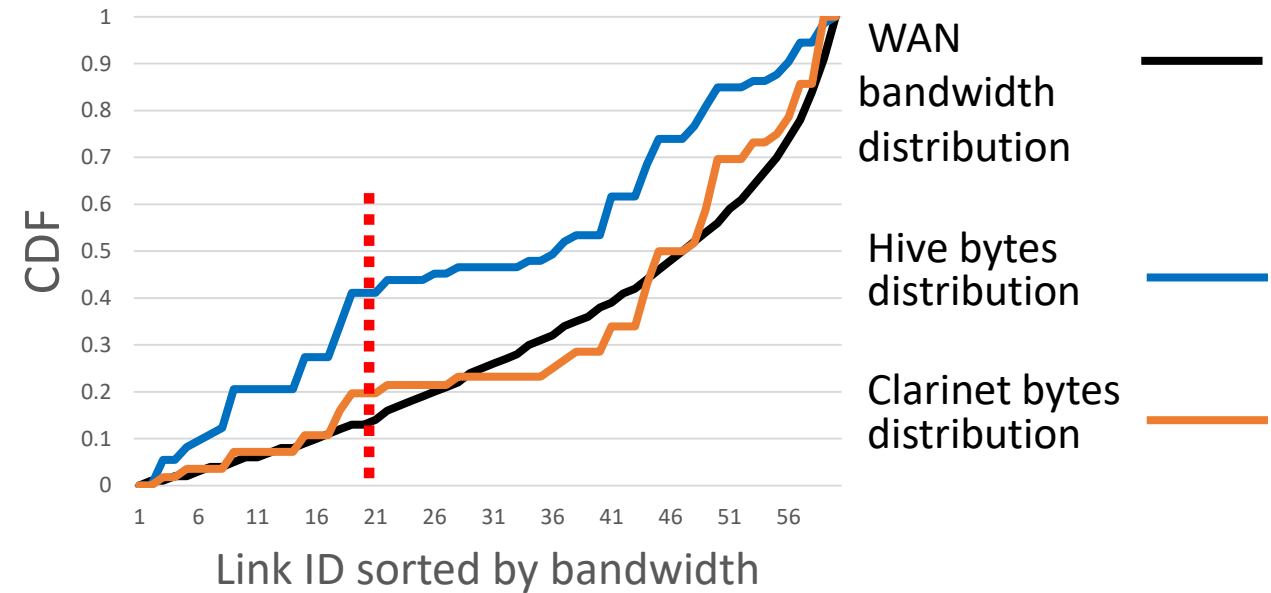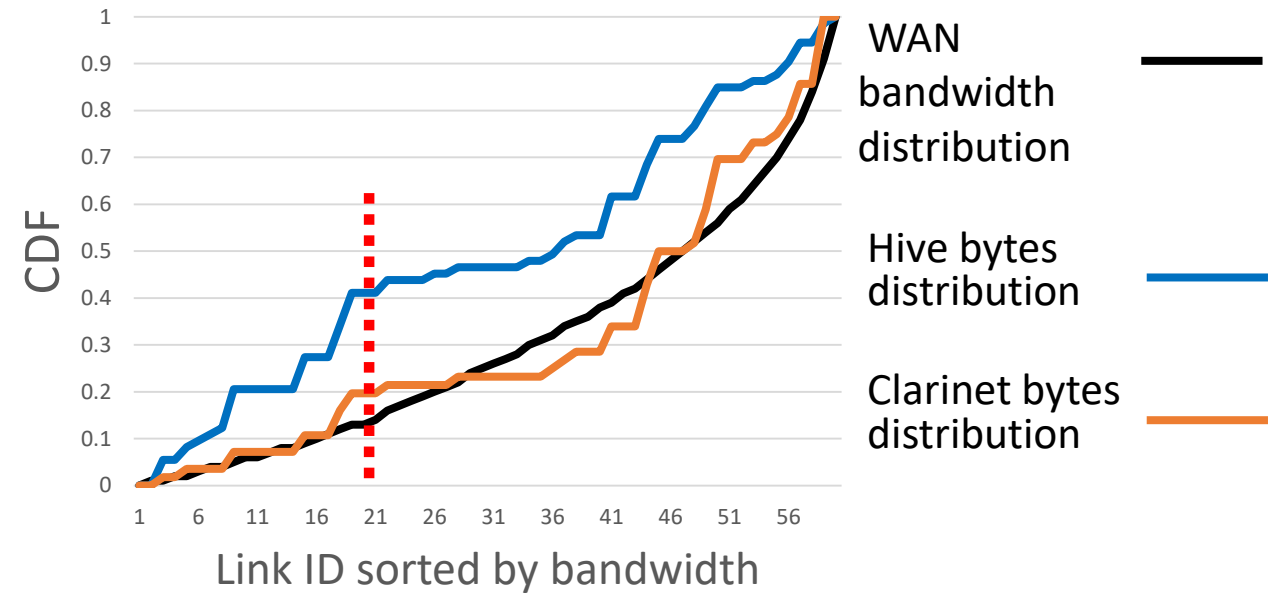| GDA Approach Vs. Hive | Average Gains |
| --- | --- |
| Clarinet | 2.7x |
| Hive + Iridium | 1.5x |
| Hive + Reducers in single DC | 0.6x |



**Data from a single batch 12 queries**

Clarinet chooses a different plan for **75%** of queries

# Evaluation: Optimization overhead

# Evaluation: Optimization overhead

1. Generate multiple query plans

2. Iterative multi-query plan selection

# Evaluation: Optimization overhead

1. Generate multiple query plans
   - Up to 64 plans in less than 5 s

2. Iterative multi-query plan selection

# Evaluation: Optimization overhead

1. Generate multiple query plans
   - Up to 64 plans in less than 5 s

2. Iterative multi-query plan selection
   - Max. 15 s for batches with 12 queries
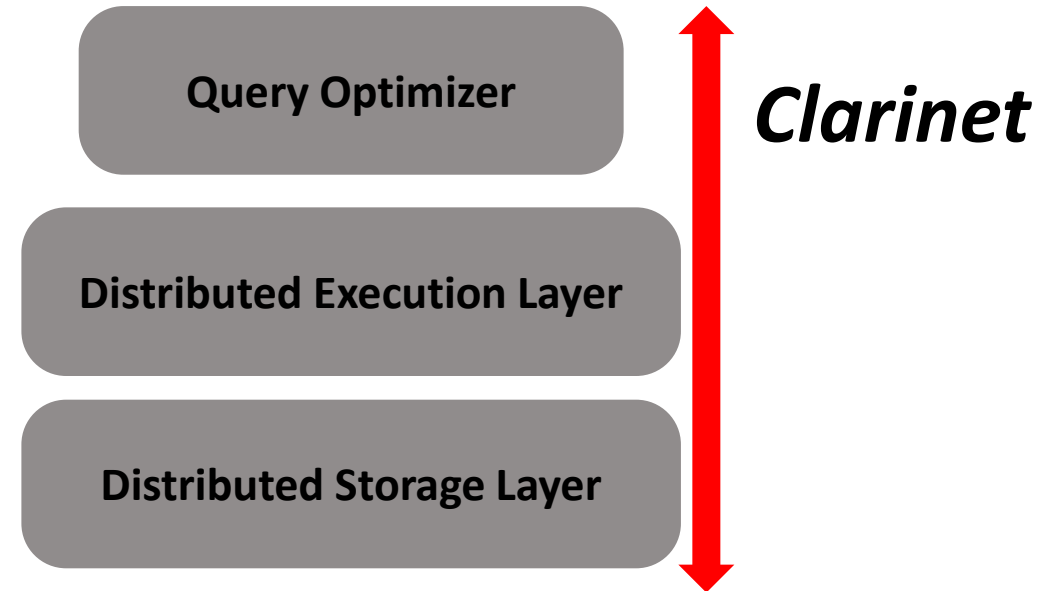
# Evaluation: Optimization overhead

1. Generate multiple query plans
   - Up to 64 plans in less than 5 s

2. Iterative multi-query plan selection
   - Max. 15 s for batches with 12 queries

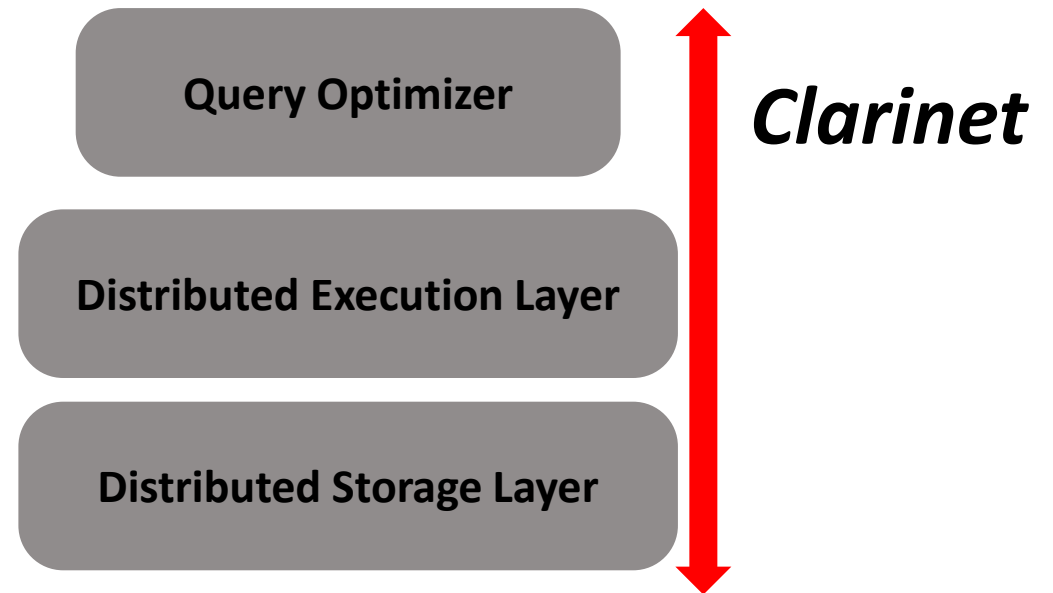Insignificant w.r.t. query running times (order of 10's of minutes)

# Summary

- WAN-awareness in QO + cross-layer optimization

# Summary

- WAN-awareness in QO + cross-layer optimization

- Presented a **scalable** way to implement **multi-query optimization** with minimal overhead

# Summary

- WAN-awareness in QO + cross-layer optimization

- Presented a **scalable** way to implement **multi-query optimization** with minimal overhead



**Query Optimizer**

**Distributed Execution Layer**

**Distributed Storage Layer**

***Clarinet***
**2.7x** Reduction in average completion time