



The NebulaStream Platform

Data and Application Management for the Internet of Things



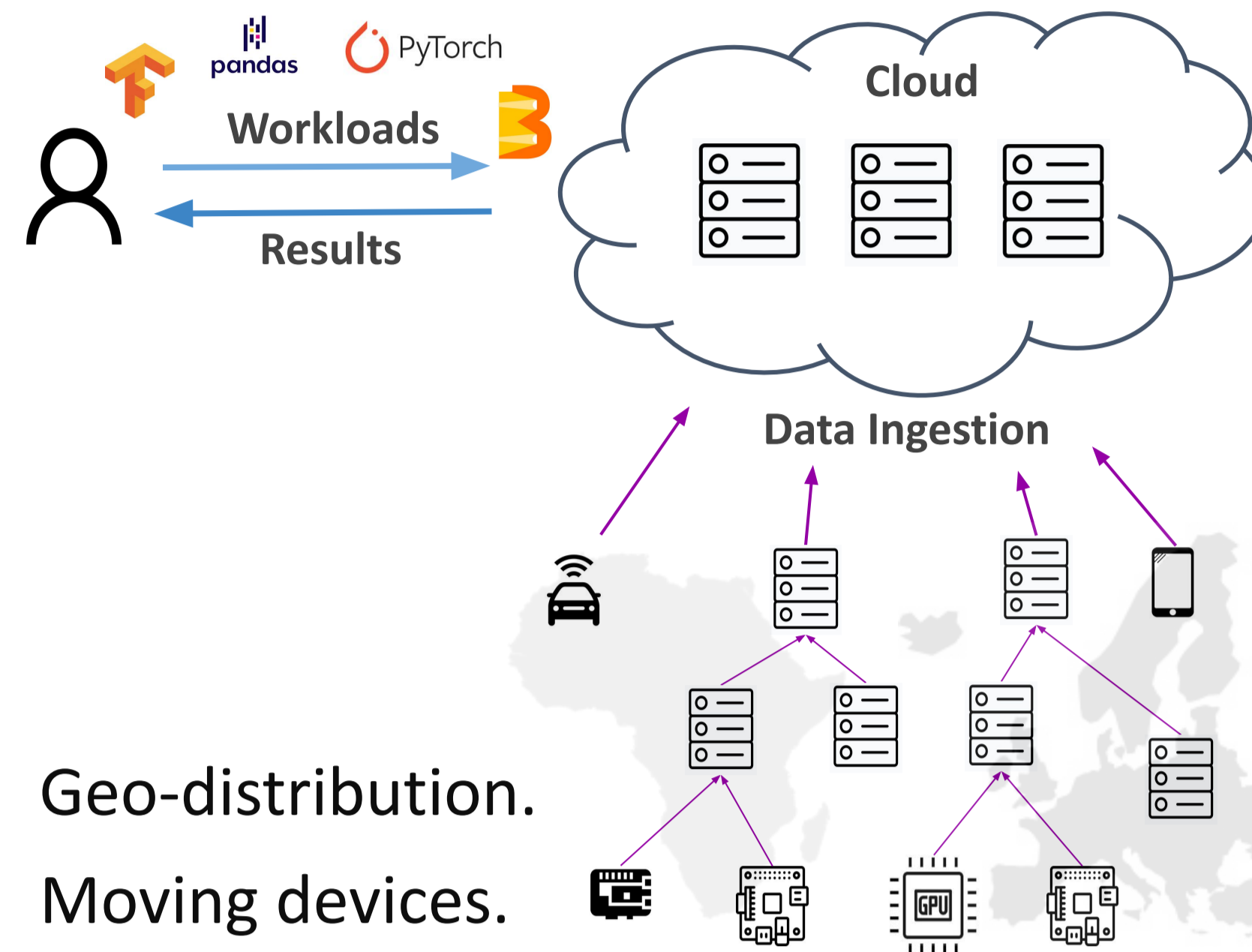
Overview

NebulaStream is a **scalable, adaptive, and efficient** data management platform for the Internet-of-Things

Research Goals:

- Efficient execution for thousands of concurrently running queries.
- Fast deployment to massive and dynamic topologies combining edge and cloud.
- Full utilization of heterogeneous hardware resources and accelerators.
- Support for complex analytical workloads involving stateful operators and UDFs.

Unifying Edge and Cloud



- Geo-distribution.
- Moving devices.
- Hierarchical topologies.
- Diverse use cases and workloads.
- Heterogeneous compute resources.
- Control over data acquisition.

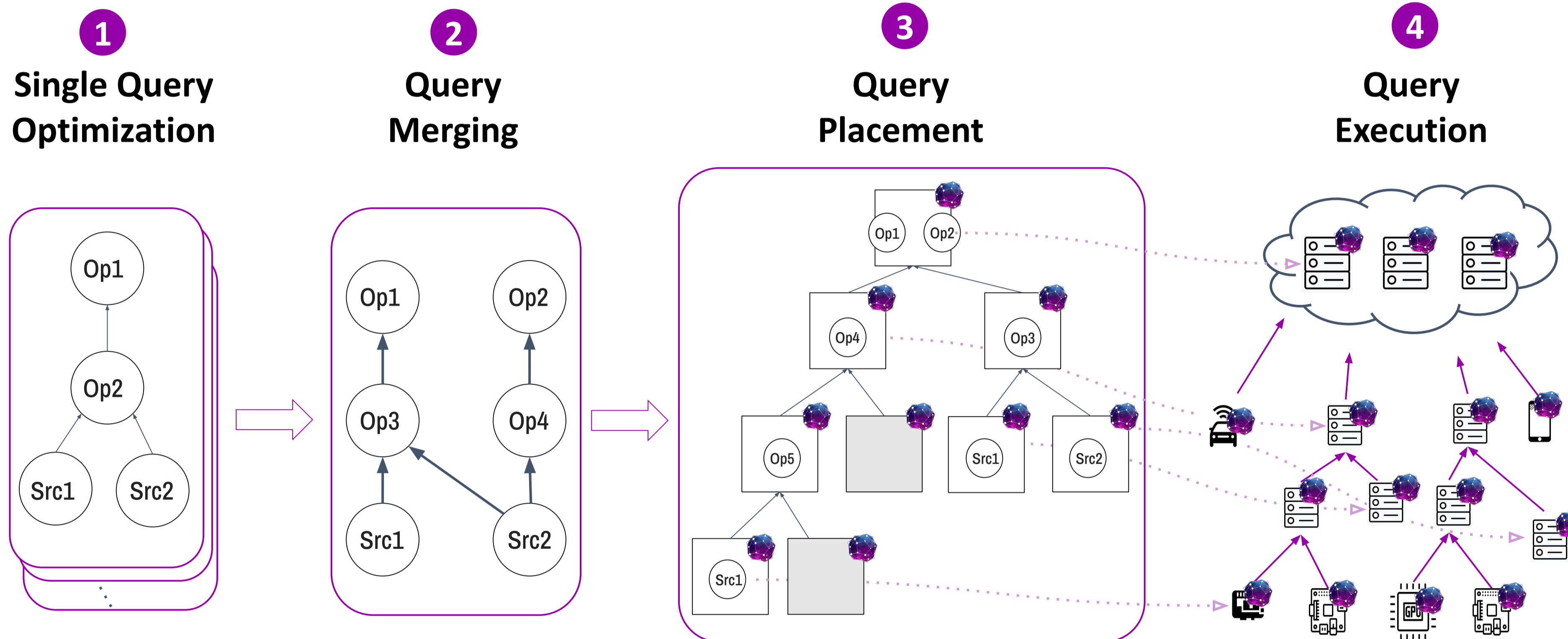
An Example Application

How much energy do our clients consume per day?

```
Query::from("smart_meter")
  .filter(Attribute("type")==="electricity")
  .window(TumblingWindow(days(1)))
  .byKey(Attribute("owner_id"))
  .apply(Sum("value"))
```

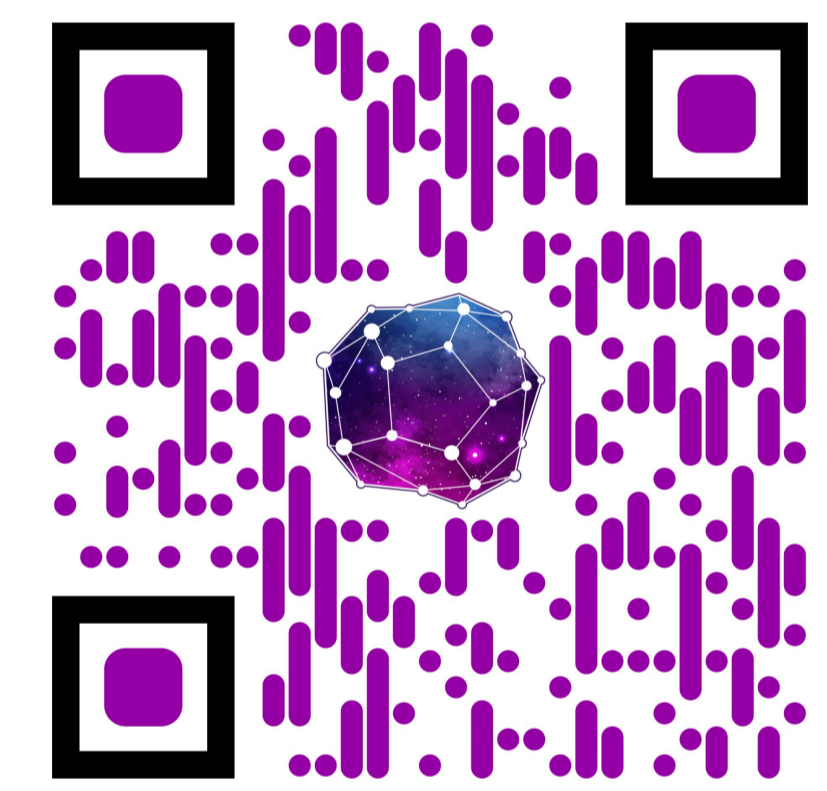


Distributed Query Deployment Flow in NebulaStream



Visit our website

Learn about NebulaStream. Try out our examples. Join our project and collaborate with us!



NebulaStream

<https://nebula.stream>
nebulastream@dima.tu-berlin.de

How to optimize query plans?

Concurrently-running streaming queries often consume the same sources or perform similar tasks.

We use Semantic Stream Query Merging:

- To derive semantic operator signatures.
- To identify sharing opportunities via constraint solving even for syntactically different queries.

```
Query 1
Query::from("car")
  .map(Attribute("speed") = Attribute("speed")*1.6)
  .filter(Attribute("speed") > 100)
  .map(Attribute("over_speed") = true)
  .sink(...)
```

```
Query 2
Query::from("car")
  .filter(Attribute("speed")*1.6 > 100)
  .map(Attribute("speed") = Attribute("speed")*1.6)
  .map(Attribute("over_speed") = true)
  .sink(...)
```

Common Signature

Conditions: car.source == "car" and car.speed * 1.6 > 100
Columns: car.speed = car.speed * 1.6
 car.id = car.id
 car.loc = car.loc
 car.over_speeding = true

Enable resource-efficient execution of semantically-equivalent queries

Where to place operators?

Unified edge/cloud environments consist of heterogeneous nodes with very different resources.

We explore different operator placement strategies:

- | | |
|---|--|
| Constructive Approach: | Cost-based Approach: |
| <ul style="list-style-type: none"> • Bottom Up • Top Down | <ul style="list-style-type: none"> • Random Search • Integer Linear Programming • Genetic Algorithm |

Challenges:

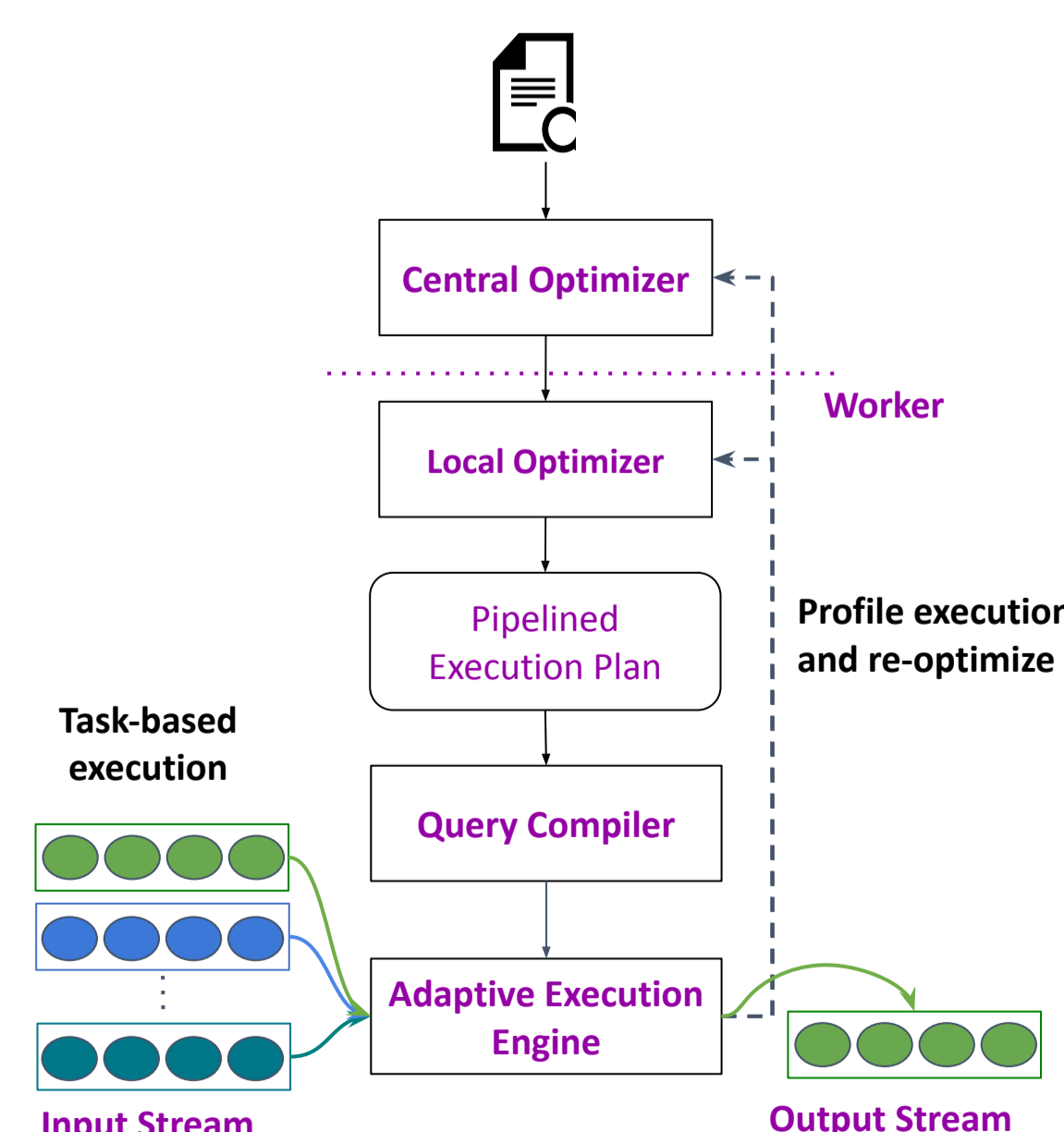
- Millions of devices
- Resource heterogeneity
- Hierarchical infrastructure
- Geo-distributed data sources
- Changing data-characteristics and statistics

Enable user to make trade-offs to speed-up query deployment

How to leverage heterogeneous resources?

State-of-the-art SPSs do not fully utilize available hardware resources.

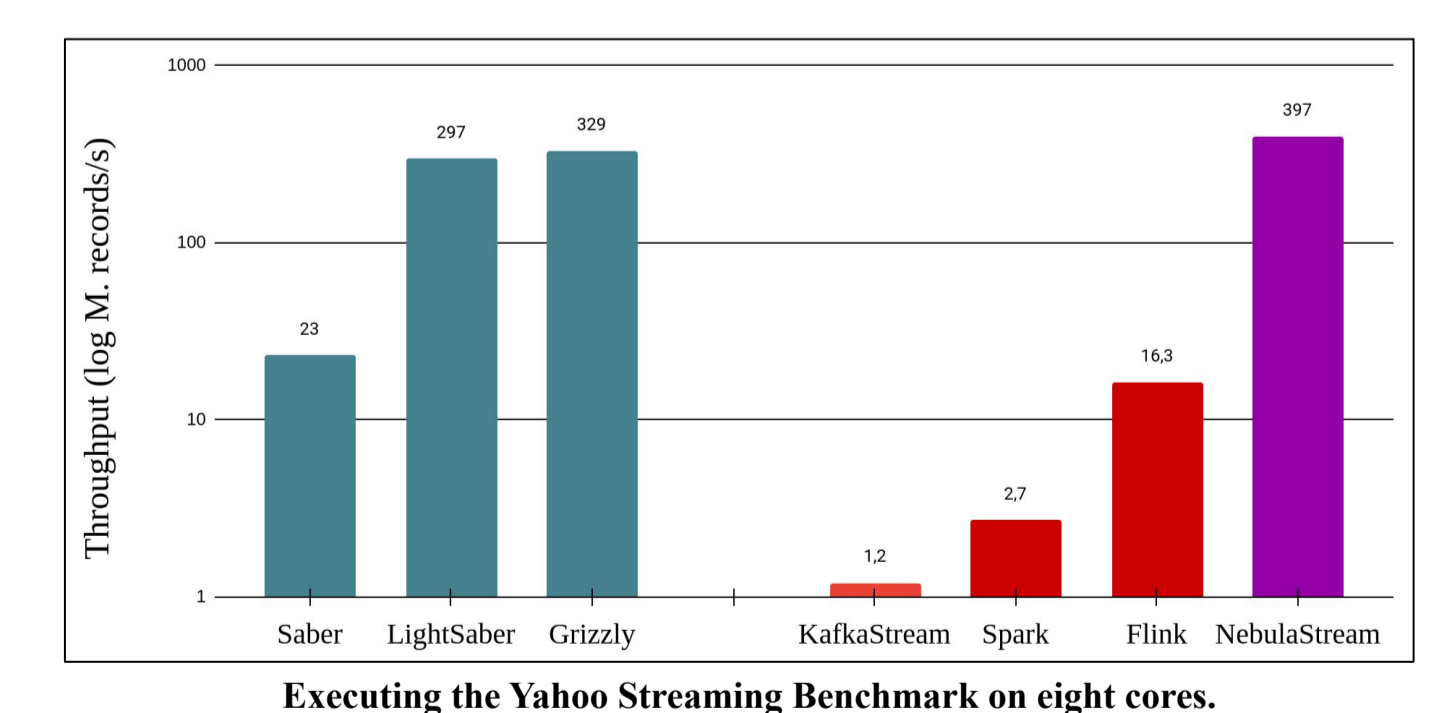
We rely on adaptive query compilation to specialize the execution to data and hardware characteristics.



Enable hardware-conscious and data-conscious query execution

Evaluation

NebulaStream combines the performance of research prototypes with the generality of mature SPSs.



Query merging is crucial to support a high number of concurrent queries.

