# Skitter: A Distributed Stream Processing Framework with Pluggable Distribution Strategies
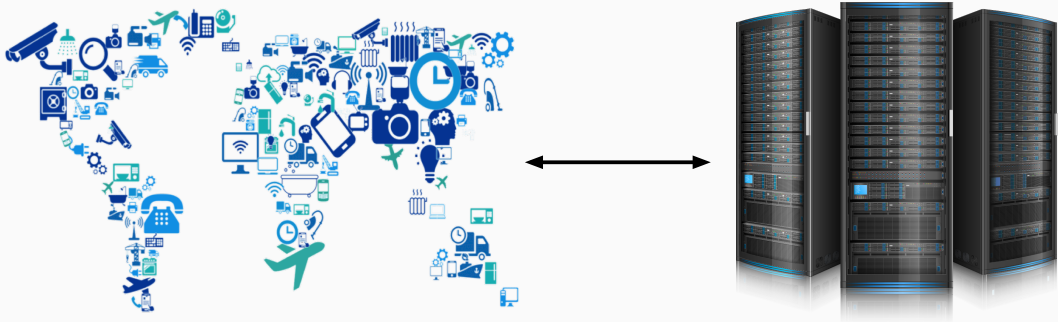
**Mathijs Saey**, Joeri De Koster, Wolfgang De Meuter

mathijs.saey@vub.be

# Reactive Big Data Applications



- Respond to real-time data streams
- Volume of incoming data requires execution on a cluster

$(ad, product, session)$

$(product, session)$

$(ad, conversion_\%)$

Build application by combining operations into a DAG.

# Distribution Strategies

The distribution strategy of an operation determines how it is distributed over the cluster.



- Spawning workers.
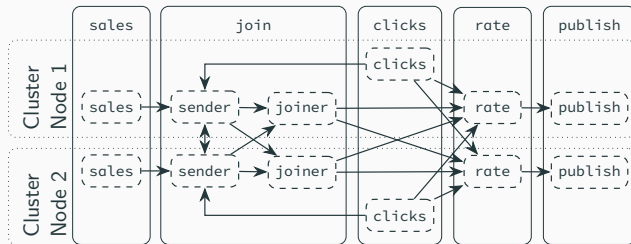- Communication between workers.
- Role performed by each worker.
- Partitioning of state between workers.

# Distribution Strategies



The distribution strategy of an operation determines how it is distributed over the cluster.

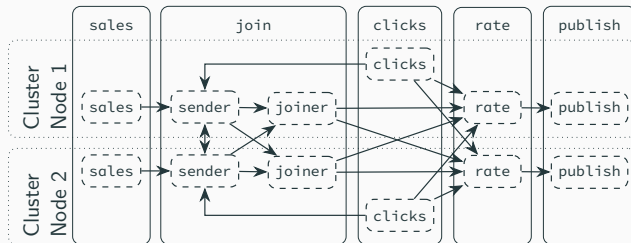- Spawning workers.
- Communication between workers.
- Role performed by each worker.
- Partitioning of state between workers.

**Importance**

Distribution strategies are key to the performance of a distributed stream processing application.

## Distribution Strategies

The distribution strategy of an operation determines how it is distributed over the cluster.



- Spawning workers.
- Communication between workers.
- Role performed by each worker.
- Partitioning of state between workers.

**Importance**

Distribution strategies are key to the performance of a distributed stream processing application.

**Goal**

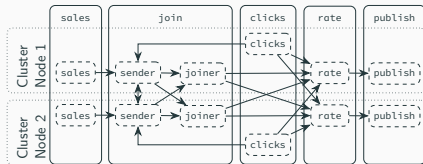We need a DSPF which makes it easy to select the appropriate distribution strategy.

```
clicks = source()
sales = source()

sales.join(clicks)
     .where(… -> …)
     .equalTo(… -> …)
     .apply(… -> …)
     .union(clicks.map(… -> …))
     .keyBy(… -> …)
     .reduce(…, … -> …)
     .map(… -> …)
     .publish()
```



- Limited set of operators.
- Fixed strategy for each operator.
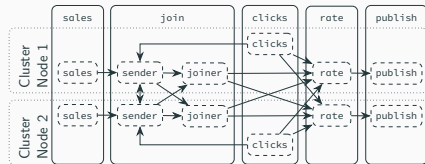
## 2 – Low-level DSPFs: Wiring DAGs in Storm

```
b = TopologyBuilder();

b.setSpout("sales", SalesSpout(), 2)
b.setSpout("clicks", ClicksSpout(), 2)

b.setBolt("join-sender", JoinSendBolt(), 2)
 .localGrouping("clicks")
 .localGrouping("sales")
b.setBolt("join-joiner", JoinBolt(), 8)
 .customGrouping("join-sender", JoinBGrouping())

b.setBolt("rate", RateBolt(), 2)
 .fieldsGrouping("clicks", "ad-id")
 .fieldsGrouping("join-joiner", "ad-id")
b.setBolt("publish", PublishBolt(), 2)
 .localGrouping("rate")
```



- Flexible, low-level model.
- Difficult to express strategies.

  - Scattered distribution logic.
  - Tangled distribution and application logic.
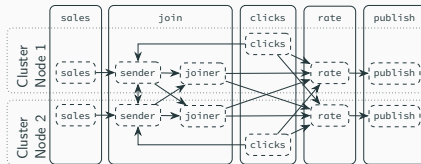  - No support for different worker types.

```
b = TopologyBuilder();

b.setSpout("sales", SalesSpout(), 2)
b.setSpout("clicks", ClicksSpout(), 2)

b.setBolt("join-sender", JoinSendBolt(), 2)
 .localGrouping("clicks")
 .localGrouping("sales")
b.setBolt("join-joiner", JoinBolt(), 8)
 .customGrouping("join-sender", JoinBGrouping())

b.setBolt("rate", RateBolt(), 2)
 .fieldsGrouping("clicks", "ad-id")
 .fieldsGrouping("join-joiner", "ad-id")
b.setBolt("publish", PublishBolt(), 2)
 .localGrouping("rate")
```



- Flexible, low-level model.
- Difficult to express strategies.

  - Scattered distribution logic.
  - Tangled distribution and application logic.
  - No support for different worker types.
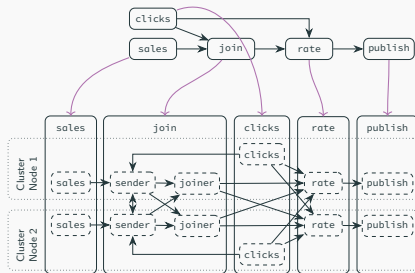
## 2 – Low-level DSPFs: Wiring DAGs in Storm

```
b = TopologyBuilder();

b.setSpout("sales", SalesSpout(), 2)
b.setSpout("clicks", ClicksSpout(), 2)

b.setBolt("join-sender", JoinSendBolt(), 2)
  .localGrouping("clicks")
  .localGrouping("sales")
b.setBolt("join-joiner", JoinBolt(), 8)
  .customGrouping("join-sender", JoinBGrouping())

b.setBolt("rate", RateBolt(), 2)
  .fieldsGrouping("clicks", "ad-id")
  .fieldsGrouping("join-joiner", "ad-id")
b.setBolt("publish", PublishBolt(), 2)
  .localGrouping("rate")
```



- Flexible, low-level model.
- Difficult to express strategies.

  ○ Scattered distribution logic.
  ○ Tangled distribution and
    application logic.
  ○ No support for different
    worker types.

## Problem Statement



DSPFs: Distribution Over a Cluster

- High-level model to express applications.
- Flexible model to express distribution strategies.
- In a modular fashion.

Novel DSPF with Pluggable Distribution Strategies

- Programming model

  **Dual** Separate abstractions for data processing and distribution logic.

  **Open** Strategies and operations can be implemented as needed.

- Implementation in Elixir

```
workflow do
  …
end
```

```
defoperation Rate, … do
  defcb key(data) do
    …
  end

  defcb react(data) do
    …
  end
end
```

```
defstrategy KeyedState do
  defhook deploy(args) do
    …
  end

  defhook deliver(data) do
    …
  end

  defhook process(data, state, role) do
    …
  end
end
```

```
workflow do
  node(ClicksSource, as: clicks)
  clicks.out ~> join.right
  clicks.out ~> rate.clicks

  node(SalesSource, as: sales)
  ~> node(Join, with: FastJoin, as: join)
  ~> node(Rate, with: KeyedState, as: rate)
  ~> node(Publish)
end
```

$(ad_1, click)\ (ad_2, click)\ (ad_1, click)$

$(ad_1, sale)$

rate

$(ad_1, 0.5)\ (ad_2, 0)\ (ad_1, 1)\ (ad_1, 0)$

- ◎◯◯ skitter. calls strategy hooks (meta level) in response to events.
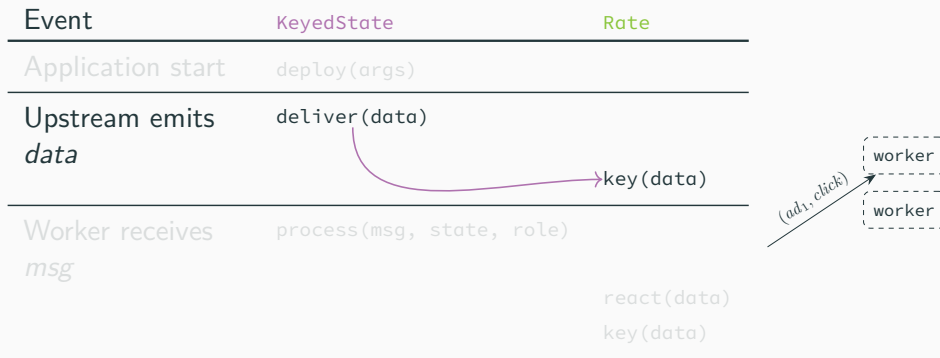- Strategy calls operation callbacks (base level) to handle data processing logic.

| Event | KeyedState | Rate |
|---|---|---|
| Application start | deploy(args) | |
| Upstream emits *data* | deliver(data) | |
| | | key(data) |
| Worker receives *msg* | process(msg, state, role) | |
| | | react(data) |
| | | key(data) |

- skitter. calls strategy hooks (meta level) in response to events.
- Strategy calls operation callbacks (base level) to handle data processing logic.

| Event | KeyedState | Rate |
|---|---|---|
| Application start | deploy(args) | |
| Upstream emits *data* | deliver(data) | |
| | | key(data) |
| Worker receives *msg* | process(msg, state, role) | |
| | | react(data) |
| | | key(data) |

worker

worker

- skitter. calls strategy hooks (meta level) in response to events.
- Strategy calls operation callbacks (base level) to handle data processing logic.

| Event | KeyedState | Rate |
|---|---|---|
| Application start | deploy(args) | |
| Upstream emits *data* | deliver(data) | key(data) |
| Worker receives *msg* | process(msg, state, role) | |
| | | react(data) |
| | | key(data) |

worker

worker

*(ad₁, click)*

- skitter. calls strategy hooks (meta level) in response to events.
- Strategy calls operation callbacks (base level) to handle data processing logic.

| Event | KeyedState | Rate |
|---|---|---|
| Application start | deploy(args) | |
| Upstream emits *data* | deliver(data) | |
| | | key(data) |
| Worker receives *msg* | process(msg, state, role) | |
| | | →react(data) |
| | | →key(data) |



$(ad_1, 0)$

worker

worker

- ◎⟳⟩ skitter. calls strategy hooks (meta level) in response to events.
  - Hooks are fixed and defined by Skitter.
- Strategy calls operation callbacks (base level) to handle data processing logic.
  - Callbacks to be implemented are defined by the strategy.

```
defstrategy KeyedState do
  defhook deploy(args)

  defhook deliver(data) do
    …
    call(:key, args: [data])
    …
  end

  defhook process(data, state, role) do
    …
    call(:key, args: [data])
    call(:react, state: state, args: [data])
    …
  end
end
```

```
defoperation Rate, … do
  defcb key(data) do
    …
  end

  defcb react(data) do
    …
  end
end
```

# Evaluation

## Research Questions

**Qualitative**     *Modularity*     Does Skitter enable the expression of distribution strategies in a modular fashion?

**Quantitative**     *Performance*     Does Skitter influence the performance characteristics of distribution strategies?

                             *Overhead*     Do the Skitter language abstractions introduce a significant amount of overhead?

                             *Impact*     Can application performance be improved by selecting an alternative strategy?

When Two Choices Are not Enough:
Balancing at Scale in Distributed Stream Processing

Muhammad Anis Uddin Nasir, Gianmarco De Francisci Morales, Nicolas Kourtellis, Marco Serafini

ICDE'16
115 citations



Scalable Distributed Stream Join Processing

Qian Lin, Beng Chin Ooi, Zhengkui Wang, Cui Yu

SIGMOD/PODS'15
113 citations

Comparison of multiple distribution strategies
Performance evaluation in Storm

# Experimental Setup

| Benchmark | Strategy | Label |
| --- | --- | --- |
| WordCount | D-Choices | D-C |
| | W-Choices | W-C |
| | Partial Key Grouping | PKG |
| | Key Grouping | KG |
| | Shuffle Grouping | SG |
| Join | Join-Matrix | JM |
| | Join-Biclique | JB |
| | Join-Biclique ContRand | JB-CR |

- 3 implementations: Storm, Skitter, ad-hoc (Elixir)
- Used to compare modularity and performance (average throughput)

# Q1: Modularity

## Question

How modular are distribution strategies in Skitter compared to the state of the art (Storm)?

| Benchmark | Strategy | Label |
|-----------|----------|-------|
| WordCount | D-Choices | D-C |
| | W-Choices | W-C |
| | Partial Key Grouping | PKG |
| | Key Grouping | KG |
| | Shuffle Grouping | SG |
| Join | Join-Matrix | JM |
| | Join-Biclique | JB |
| | Join-Biclique ContRand | JB-CR |

- Measure LOC added or modified to change distribution strategy.
- Categorize LOC based on abstractions offered by framework.

## Q1: Modularity (Join)

### Question

How modular are distribution strategies in Skitter compared to the state of the art (Storm)?

|  | Strategy | Storm | | | Skitter | | |
|---|---|---|---|---|---|---|---|
|  |  | Topology | Component | Grouping | Workflow | Operation | Strategy |
| Q5 | JB | 29 | 162 | 46 | 3 | 0 | 119 |
| | JB-CR | 29 | 162 | 61 | 3 | 0 | 134 |
| Q7 | JB | 22 | 162 | 46 | 2 | 0 | 119 |
| | JB-CR | 22 | 162 | 61 | 2 | 0 | 134 |

# Q1: Modularity (Join)

## Question

How modular are distribution strategies in Skitter compared to the state of the art (Storm)?

| | Strategy | Storm | | | Skitter | | |
|---|---|---|---|---|---|---|---|
| | | Topology | Component | Grouping | Workflow | Operation | Strategy |
| Q5 | JB | 29 | 162 | 46 | 3 | 0 | 119 |
| | JB-CR | 29 | 162 | 61 | 3 | 0 | 134 |
| Q7 | JB | 22 | 162 | 46 | 2 | 0 | 119 |
| | JB-CR | 22 | 162 | 61 | 2 | 0 | 134 |

# Q1: Modularity (Join)

## Question

How modular are distribution strategies in Skitter compared to the state of the art (Storm)?

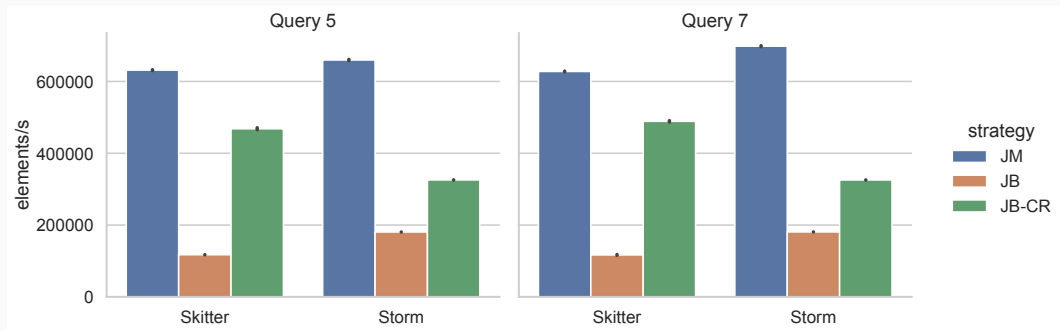| | Strategy | Storm | | | Skitter | | |
|---|---|---|---|---|---|---|---|
| | | Topology | Component | Grouping | Workflow | Operation | Strategy |
| Q5 | JB | 29 | 162 | 46 | 3 | 0 | 119 |
| | JB-CR | 29 | 162 | 61 | 3 | 0 | 134 |
| Q7 | JB | 22 | 162 | 46 | 2 | 0 | 119 |
| | JB-CR | 22 | 162 | 61 | 2 | 0 | 134 |

## Q2: Performance

**Question**

Do strategies implemented in Skitter maintain their performance characteristics?

- Compare the *relative* performance of Storm and Skitter implementations of the same experiments.

**Question**

Do strategies implemented in Skitter maintain their performance characteristics?

# Conclusion



https://soft.vub.ac.be/~mathsaey/skitter/

# Skitter: A Distributed Stream Processing Framework with Pluggable Distribution Strategies

**Mathijs Saey**, Joeri De Koster, Wolfgang De Meuter

mathijs.saey@vub.be

```elixir
workflow do
  source()
  ~> flatmap(&String.split/1, with: RepartitionedOutput)
  ~> keyed_reduce(fn word -> word end, fn count -> count + 1 end, 0)
  ~> print()
end

workflow do
  node(SomeSource)
  ~> node(FlatMap, args: [&String.split/1], with: RepartitionedOutput)
  ~> node(KeyedReduce, args: [fn word -> word end, fn count -> count + 1 end, 0])
  ~> node(Print)
end
```

```
defoperation Rate, in: [sales, clicks], out: conversion_rate, strategy: KeyedState do
  initial_state {0, 0}

  defcb key(data), do: data.ad_id

  defcb react(data) do
    {clicks, sales} = state()
    {new_clicks, new_sales} = case port_of(data) do
      :sales -> {clicks, sales + 1}
      :clicks -> {clicks + 1, sales}
    end
    state <~ {new_clicks, new_sales}
    {data.ad_id, new_sales / new_clicks} ~> conversion_rate
  end
end
```

```
defstrategy KeyedState do
  defhook deploy(args) do
    Remote.on_all_workers(fn -> local_worker(Map.new(), :aggregator) end)
    |> Enum.map(fn {remote, worker} -> worker end)
  end

  defhook deliver(data) do
    key = call(:key, args: [data]).result
    aggregators = deployment()
    idx = rem(Murmur.hash_x86_32(key), length(aggregators))
    worker = Enum.at(aggregators, idx)
    send(worker, data)
  end

  defhook process(data, state_map, :aggregator) do
    key = call(:key, args: [data]).result
    state = Map.get(state_map, key, initial_state())
    res = call(:react, state: state, args: [data])
    emit(res.emit)
    Map.put(state_map, key, res.state)
  end
end
```
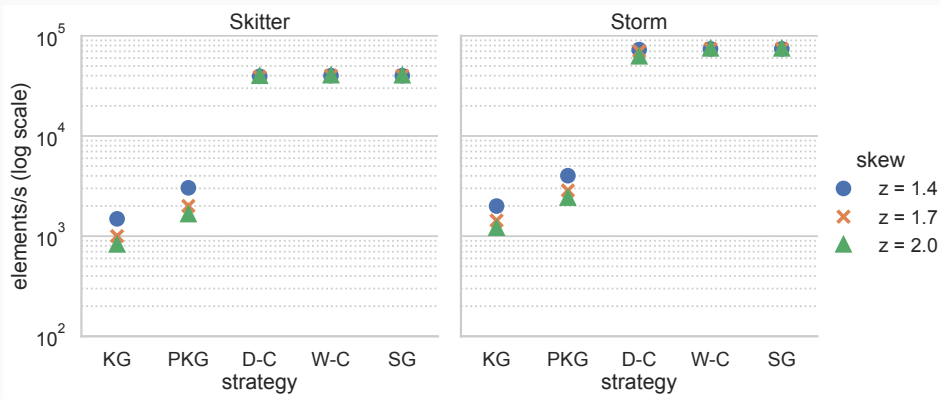
## Q1: Modularity (WordCount)

| Strategy | Storm | | | Skitter | | |
|----------|----------|-----------|----------|----------|-----------|----------|
| | Topology | Component | Grouping | Workflow | Operation | Strategy |
| SG | 1 | 0 | 0 | 1 | 0 | 8 |
| PKG | 1 | 0 | 0 | 1 | 0 | 46 |
| W-C | 1 | 0 | 29 | 1 | 0 | 71 |
| D-C | 1 | 0 | 59 | 1 | 0 | 107 |
| PKG† | 4 | 38 | 0 | 1 | 4 | 65 |
| W-C† | 4 | 38 | 29 | 1 | 4 | 90 |
| D-C† | 4 | 38 | 59 | 1 | 4 | 126 |

## Question
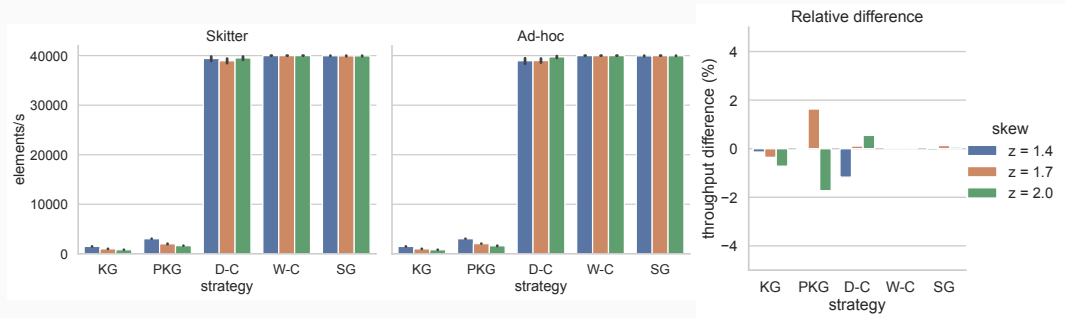
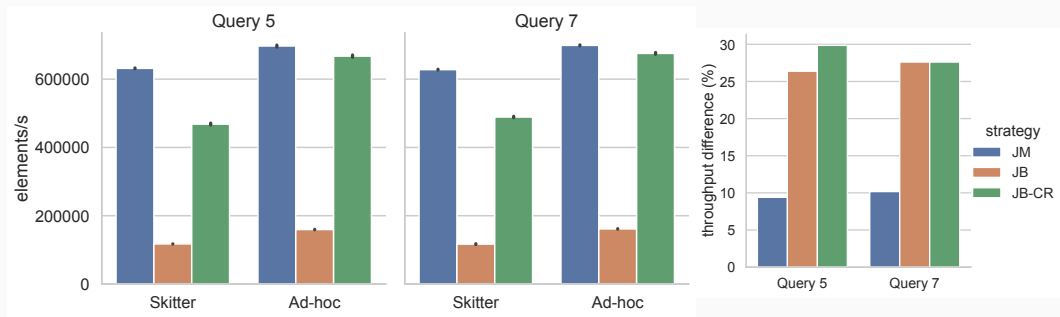Do strategies implemented in Skitter maintain their performance characteristics?

## Question

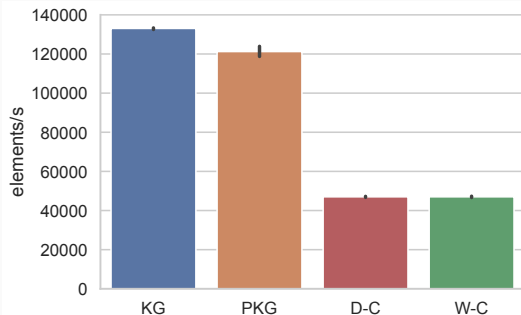Do the abstractions introduced by Skitter introduce additional overhead?

## Question

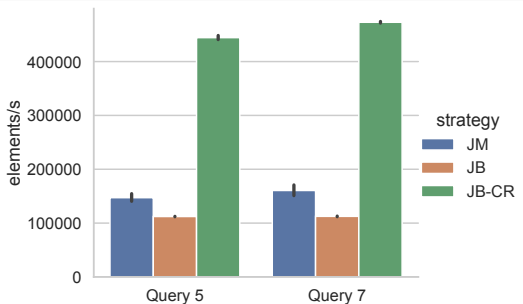Do the abstractions introduced by Skitter introduce additional overhead?

# Q4: Impact

## Question

Can we improve performance by changing distribution strategy?



WordCount benchmark with key merging and no skew ($z = 0$).



Join benchmark handling 80GB of data.