# Contributions to
# Large-Scale Data Processing Systems
## PhD Defense

**Matthieu Caneill**

February 5, 2018

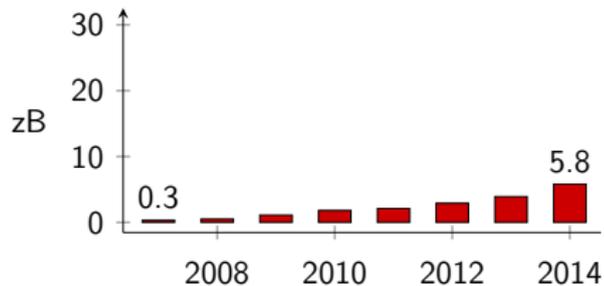**Daniel Hagimont**
INP Toulouse
ENSEEIHT

**Jean-Marc Menaud**
IMT Atlantique
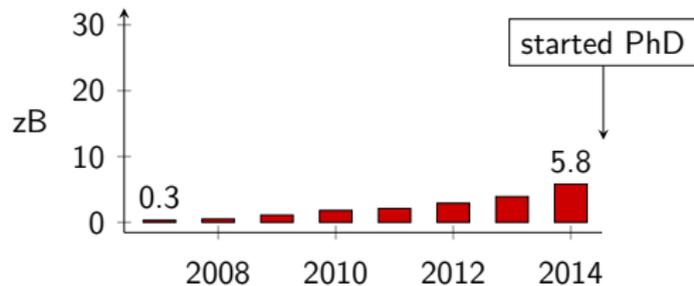
**Sihem Amer-Yahia**
CNRS / Université
Grenoble Alpes

**Noël De Palma**
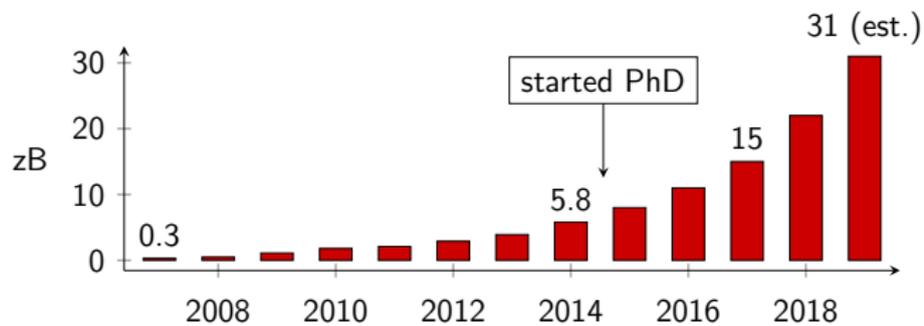Université
Grenoble Alpes
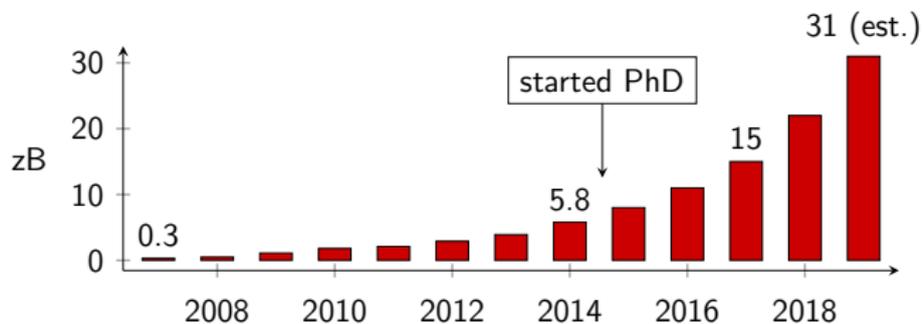
## Worldwide data production

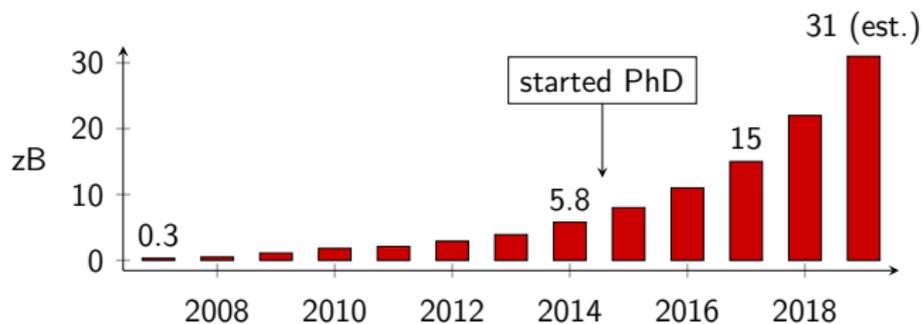Worldwide data production

## Worldwide data production

## Worldwide data production



$1 zetabyte = 1000 exabytes = 10^6 petabytes = 10^9 terabytes$

## Worldwide data production



$1 zetabyte = 1000 exabytes = 10^6 petabytes = 10^9 terabytes$

(1 zetabyte is 2 billion times my hard drive)

Applications

▶ Genome sequencing and querying (human: 3 B base pairs)

Applications

- ▶ Genome sequencing and querying (human: 3 B base pairs)
- ▶ Web and social networks (Facebook: 600 TB/day in 2014)

Applications

- Genome sequencing and querying (human: 3 B base pairs)
- Web and social networks (Facebook: 600 TB/day in 2014)
- Particle physics (CERN: 1 PB/s of collision data)
- etc.

### Applications

- Genome sequencing and querying (human: 3 B base pairs)
- Web and social networks (Facebook: 600 TB/day in 2014)
- Particle physics (CERN: 1 PB/s of collision data)
- etc.

### Problems

- Data management at scale

## Applications

- ▶ Genome sequencing and querying (human: 3 B base pairs)
- ▶ Web and social networks (Facebook: 600 TB/day in 2014)
- ▶ Particle physics (CERN: 1 PB/s of collision data)
- ▶ etc.

## Problems

- ▶ Data management at scale
- ▶ Data processing in reasonable time

## Applications

- ▶ Genome sequencing and querying (human: 3 B base pairs)
- ▶ Web and social networks (Facebook: 600 TB/day in 2014)
- ▶ Particle physics (CERN: 1 PB/s of collision data)
- ▶ etc.

## Problems

- ▶ Data management at scale
- ▶ Data processing in reasonable time
- ▶ . . . and reasonable price

How to design. . .

- ▶ An industrial system to handle monitoring data and make predictions about future failures?

How to design...

- An industrial system to handle monitoring data and make predictions about future failures?
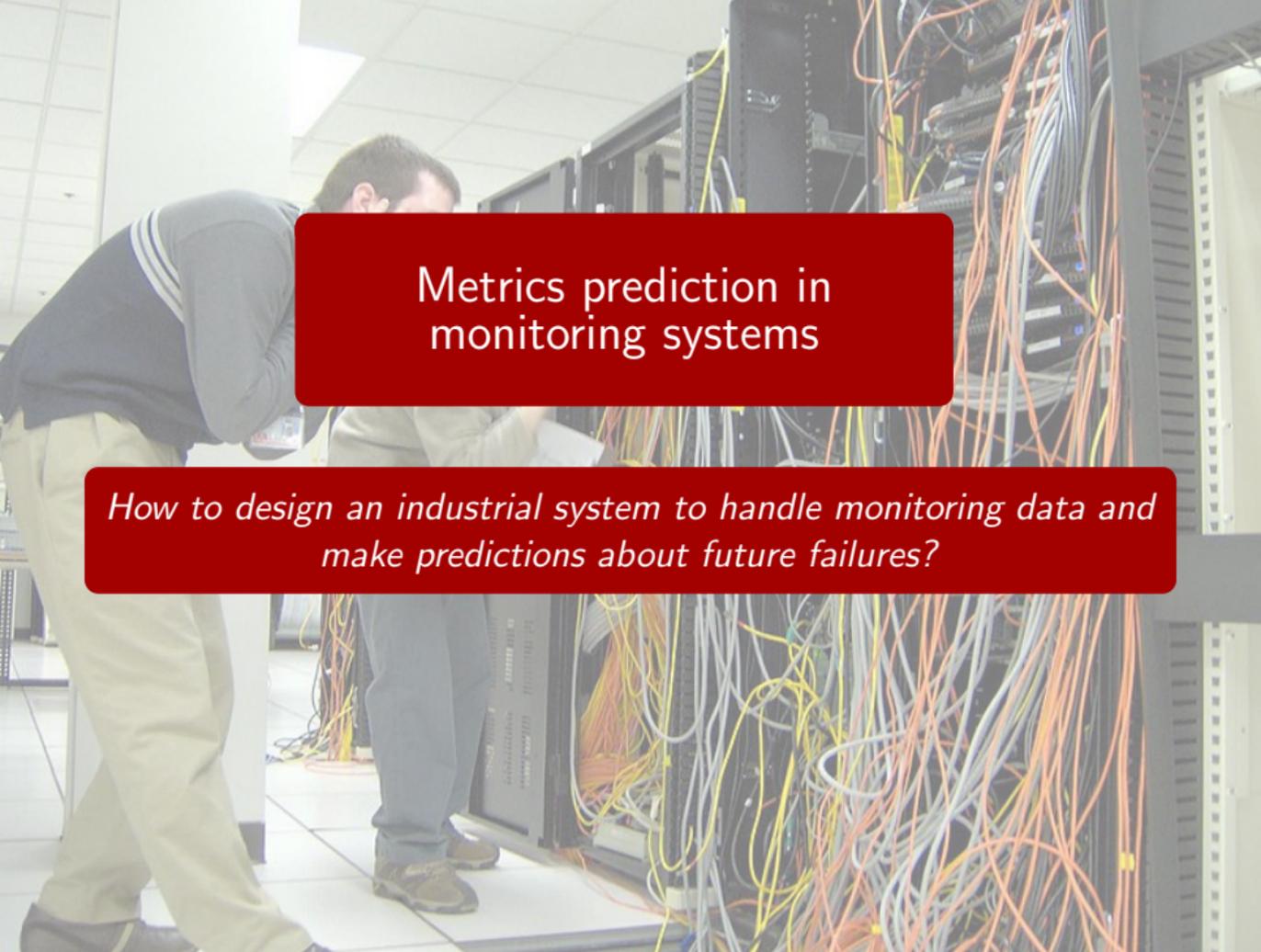- An algorithm to improve locality in distributed streaming engines?

How to design...

- ▶ An industrial system to handle monitoring data and make predictions about future failures?
- ▶ An algorithm to improve locality in distributed streaming engines?
- ▶ A framework to compose data processing algorithms in a descriptive fashion, while reasoning on high level abstractions?

Structure of this presentation

1. Online metrics prediction in monitoring systems
2. Locality data routing
3. $\lambda$-blocks
4. Conclusion

Metrics prediction in monitoring systems

How to design an industrial system to handle monitoring data and make predictions about future failures?

Actors and roles of Smart Support Center

- **Coservit**: Monitoring services
- **HP**: Cloud computing, hardware
- **LIG – AMA**: Machine learning
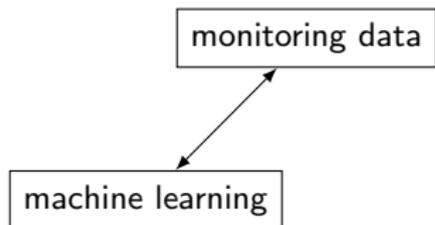- **LIG – ERODS**: Cloud computing, systems

## Scope of Smart Support Center

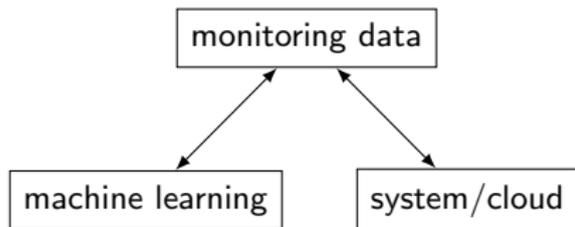monitoring data

► Monitoring insights

## Scope of Smart Support Center



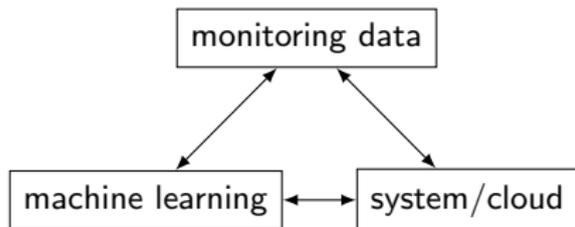- ▶ Monitoring insights
- ▶ Failure prediction

## Scope of Smart Support Center



- ► Monitoring insights
- ► Failure prediction
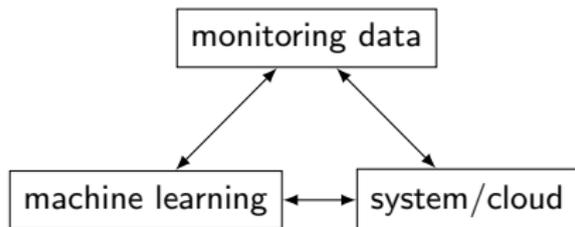- ► Infrastructure scaling

## Scope of Smart Support Center



- Monitoring insights
- Failure prediction
- Infrastructure scaling

## Scope of Smart Support Center



- ▶ Monitoring insights
- ▶ Failure prediction
- ▶ Infrastructure scaling
- ▶ More server uptime

Challenges

- Scale monitoring infrastructure (from 1 to $N$ nodes)

## Challenges

- Scale monitoring infrastructure (from 1 to $N$ nodes)
- System design for low latency analytics

Challenges

- Scale monitoring infrastructure (from 1 to $N$ nodes)
- System design for low latency analytics
- Fault tolerance

## Metrics

▶ Monitoring metric: observation point on a server in a datacenter

## Metrics

▶ Monitoring metric: observation point on a server in a datacenter

▶ CPU load, memory, service status

## Metrics

- Monitoring metric: observation point on a server in a datacenter
- CPU load, memory, service status
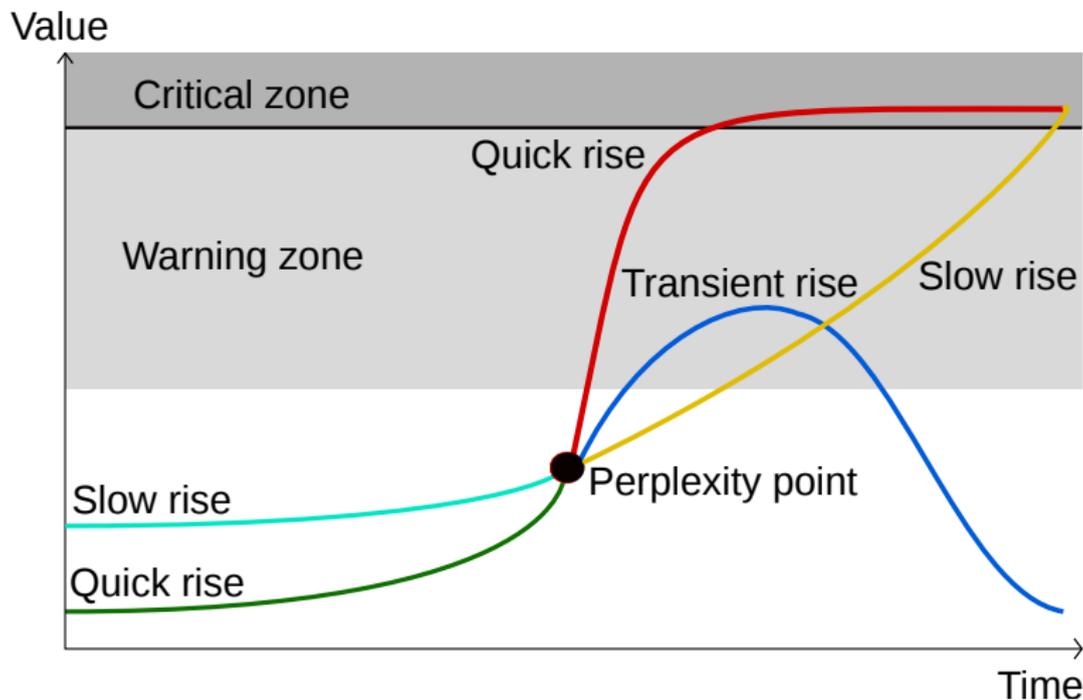- Reported by agents, processed, and stored

## Metrics

- Monitoring metric: observation point on a server in a datacenter
- CPU load, memory, service status
- Reported by agents, processed, and stored
- Computed as time-series

## Metrics
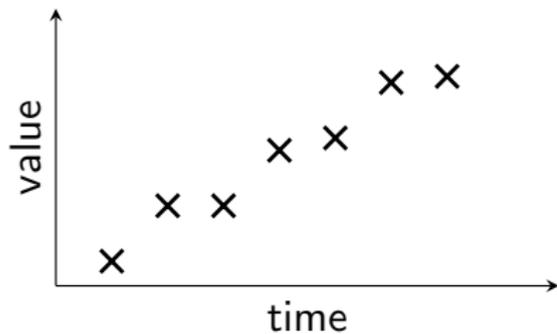
- Monitoring metric: observation point on a server in a datacenter
- CPU load, memory, service status
- Reported by agents, processed, and stored
- Computed as time-series
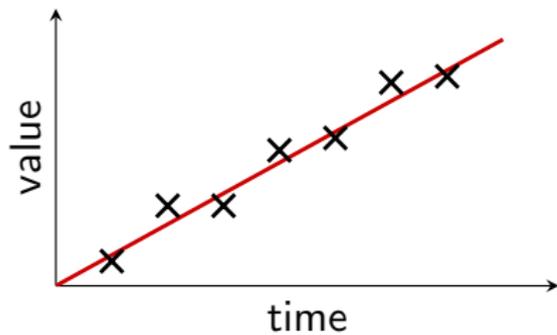- Associated to thresholds: warning and critical
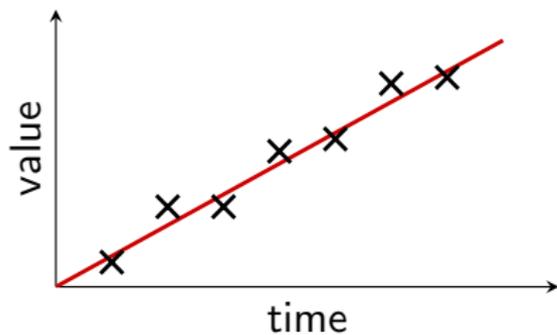
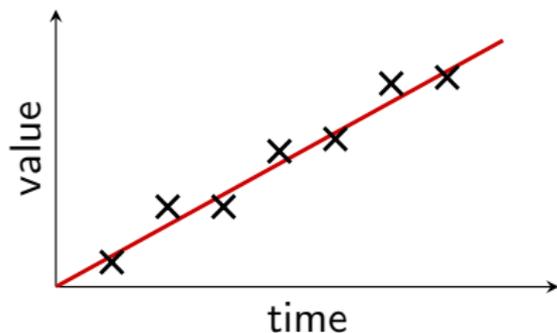Metrics behaviour: 6 scenarios

## Linear regression

Linear regression

## Linear regression
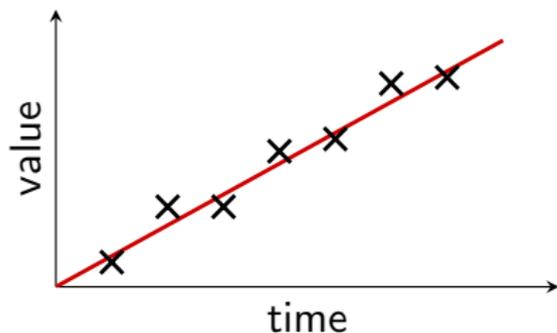


- Ability to identify local trends (few hours)

## Linear regression



- Ability to identify local trends (few hours)
- Fast to compute

## Linear regression



- ► Ability to identify local trends (few hours)
- ► Fast to compute
- ► Good candidate to avoid false positives (peaks)

## Linear regression
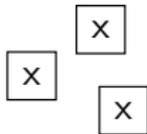


- Ability to identify local trends (few hours)
- Fast to compute
- Good candidate to avoid false positives (peaks)
- Library: MLlib (part of Apache Spark)

## System architecture



Monitoring agents

## System architecture



Monitoring agents

## System architecture

## System architecture

Desired properties

- Scalable: up to a few servers (150 CPU cores) to handle Coservit's load

System architecture

Desired properties

- Scalable: up to a few servers (150 CPU cores) to handle Coservit's load
- End-to-end fault tolerance: metrics can never be lost

Desired properties

- Scalable: up to a few servers (150 CPU cores) to handle Coservit's load
- End-to-end fault tolerance: metrics can never be lost
- Performances: "fast" to compute metrics predictions

## Setup

- Hardware: 4 servers (16–28 cores, 128–256 GB RAM)
- Dataset: Replay on production data recorded at Coservit
- 424 206 metrics, 1.5 billion data points monitored on 25 070 servers

## Evaluation



Figure: swap memory

Figure: swap memory

## Evaluation



Figure: swap memory

## Evaluation



Figure: physical memory

## Evaluation



Figure: physical memory

## Evaluation



Figure: physical memory

Evaluation



Figure: disk partition

● ● ○ ○ ○

## Evaluation



Figure: disk partition

Evaluation



Figure: disk partition

Metric blacklisting

- Some metrics are too volatile and hard to predict

## Evaluation

### Metric blacklisting

- Some metrics are too volatile and hard to predict
- To avoid false positives/negatives, and save resources, they are blacklisted

Evaluation

## Metric blacklisting

- Some metrics are too volatile and hard to predict
- To avoid false positives/negatives, and save resources, they are blacklisted
- Root Mean Square Error evaluated weekly

Evaluation

## Metric blacklisting

- ▶ Some metrics are too volatile and hard to predict
- ▶ To avoid false positives/negatives, and save resources, they are blacklisted
- ▶ Root Mean Square Error evaluated weekly
- ▶ Metrics (temporarily) blacklisted if their RMSE $>$ threshold

## Metric blacklisting

- Some metrics are too volatile and hard to predict
- To avoid false positives/negatives, and save resources, they are blacklisted
- Root Mean Square Error evaluated weekly
- Metrics (temporarily) blacklisted if their RMSE > threshold
- 58.5% of the metrics have a low RMSE $\rightarrow$ good predictions

## Evaluation

### CPU load and memory consumption



(a) master

(b) slave-1

Figure: Running on 4 machines and 100 cores for 15 minutes.

### Time repartition



Figure: Time repartition for predicting a metric.

Load handling

- ▶ End-to-end process for the prediction of 1 metric: 1 second.

Evaluation

## Load handling

- ► End-to-end process for the prediction of 1 metric: 1 second.
- ► One monitoring server (with 24 cores) can handle the load of 1440 metrics (at worst), which is 85 servers on average.

## Load handling: linear scaling



Figure: Amount of metrics handled in 15 minutes.

## Evaluation

### Load handling: linear scaling



Figure: Amount of metrics handled in 15 minutes.

Evaluation

## Load handling: linear scaling



Figure: Amount of metrics handled in 15 minutes.

## Related work

### Positioning

No published work exhibits the same system (end-to-end system for monitoring metrics prediction, storage and blacklisting).

### Prediction models

- Hardware failures [CAS12]

### Positioning

No published work exhibits the same system (end-to-end system for monitoring metrics prediction, storage and blacklisting).

### Prediction models

- Hardware failures [CAS12]
- Capacity planning (e.g. Microsoft Azure [mic])

### Positioning

No published work exhibits the same system (end-to-end system for monitoring metrics prediction, storage and blacklisting).

### Prediction models

- Hardware failures [CAS12]
- Capacity planning (e.g. Microsoft Azure [mic])
- Datacenter temperature (e.g. Thermocast [LLL$^+$11])

## Related work

### Positioning

No published work exhibits the same system (end-to-end system for monitoring metrics prediction, storage and blacklisting).

### Prediction models

- Hardware failures [CAS12]
- Capacity planning (e.g. Microsoft Azure [mic])
- Datacenter temperature (e.g. Thermocast [LLL$^+$11])
- Monitoring metrics (e.g. Zabbix [zab] with manual tuning)

# Locality data routing

*How to design an algorithm to improve locality in distributed streaming engines?*

**Actors**
Collaboration with **Vincent Leroy** (SLIDE)
and **Ahmed El-Rheddane** (ERODS).

Goals

- Real-time message handling

Goals

- Real-time message handling
- Real-time metric calculations

Goals

- Real-time message handling
- Real-time metric calculations
- Parallelization

Goals

- ▶ Real-time message handling
- ▶ Real-time metric calculations
- ▶ Parallelization
- ▶ Fault-tolerance

Apache Storm $\rightarrow$ topologies.

Apache Storm $\rightarrow$ topologies.



Figure: Trending hashtags topology.

$S$ sends tweets, operator $A$ extract hashtags, $B$ converts them to lowercase, and $C$ counts the frequency of each hashtag.

Apache Storm $\rightarrow$ topologies.



Figure: Trending hashtags topology.

$S$ sends tweets, operator $A$ extract hashtags, $B$ converts them to lowercase, and $C$ counts the frequency of each hashtag.

Division into tasks $\rightarrow$ distribution and parallelization made easy.

### States are associated to keys

For example, the operator C can keep the list of trending hashtags (values) per location (keys).

## Parallelization

To keep a consistent state, same keys must be routed to the same instance.



Figure: Tasks A and B are stateless, C is stateful.

## Situation

Let's have two stateful operators, each with two instances.

## Situation
Let's have two stateful operators, each with two instances.



Server 1

### Goal
Minimize the traffic between the machines: $A1 \rightarrow B2$ and $A2 \rightarrow B1$. By default, $locality = 1/parallelism$

## Situation

Let's have two stateful operators, each with two instances.



Server 1

$S$

$A_1$ → $B_1$

$A_2$ → $B_2$

Server 2

### Goal

Minimize the traffic between the machines: $A1 \to B2$ and $A2 \to B1$. By default, $locality = 1/parallelism$

### Constraint

Keep a good load balance between the machines.

## Keys correlation

Dynamically instrument the keys couples and represent them with a bipartite graph.

## Keys correlation

Dynamically instrument the keys couples and represent them with a bipartite graph.

## Keys correlation

Dynamically instrument the keys couples and represent them with a bipartite graph.



Routing tables

- $S$: $Asia \rightarrow A_1$
      $Oceania \rightarrow A_2$

- $A_1$: **#java $\rightarrow$ B$_1$**
        **#ruby $\rightarrow$ B$_1$**
        $\#python \rightarrow B_2$

- $A_2$: **#python $\rightarrow$ B$_2$**
        $\#java \rightarrow B_1$
        $\#ruby \rightarrow B_1$

## Keys correlation

Dynamically instrument the keys couples and represent them with a bipartite graph.



### Routing tables

- $S$: $Asia \rightarrow A_1$
  $Oceania \rightarrow A_2$

- $A_1$: **#java $\rightarrow$ B$_1$**
  **#ruby $\rightarrow$ B$_1$**
  $\#python \rightarrow B_2$

- $A_2$: **#python $\rightarrow$ B$_2$**
  $\#java \rightarrow B_1$
  $\#ruby \rightarrow B_1$

Graph partitioning $\rightarrow$ optimized routing, favorizing local links.

**Message**:
**Posted from**:

| Server 1 |
| $A_1$ → $B_1$ |
| $S$ |
| $A_2$ → $B_2$ |
| Server 2 |

**S**

| key | route |
|-----|-------|
|     |       |
|     |       |

**A**

| key | route |
|-----|-------|
|     |       |
|     |       |

**Message**: #python doesn't have braces
**Posted from**: Oceania

| S | |
|-----|-------|
| **key** | **route** |
| | |
| | |

| A | |
|-----|-------|
| **key** | **route** |
| | |
| | |

**Message**: #python doesn't have braces
**Posted from**: Oceania

Server 1

| S | |
|---|---|
| **key** | **route** |
| Oceania | A1 |
| | |

| A | |
|---|---|
| **key** | **route** |
| python | B2 |
| | |

Server 2

**Message**: #java is a verbose language
**Posted from**: Asia



Server 1

Server 2

| S | |
|---------|-------|
| **key** | **route** |
| Oceania | A1 |
| | |

| A | |
|---------|-------|
| **key** | **route** |
| python | B2 |
| | |

**Message**: #java is a verbose language
**Posted from**: Asia

Server 1

S

$A_1$ ⟶ $B_1$

$A_2$ ⟶ $B_2$

Server 2

**Message**:
**Posted from**:

| **S** | |
|-------|-------|
| **key** | **route** |
| Oceania | A1 |
| Asia | A2 |

| **A** | |
|--------|-------|
| **key** | **route** |
| python | B2 |
| java | B1 |

Reconfiguration is computed and applied

**Message**:
**Posted from**:

**S**

| key | route |
|---------|-------|
| Oceania | A1 |
| Asia | A2 |

**A**

| key | route |
|--------|-------|
| python | B1 |
| java | B2 |

Reconfiguration is computed and applied

Correlation between **Oceania**/**python** and **Asia**/**java**

**Message**: #python is pretty cool!
**Posted from**: Oceania

| S | |
|---|---|
| **key** | **route** |
| Oceania | A1 |
| Asia | A2 |

| A | |
|---|---|
| **key** | **route** |
| python | B1 |
| java | B2 |

**Message**: #python is pretty cool!
**Posted from**: Oceania

Server 1

| S | |
|---|---|
| **key** | **route** |
| Oceania | A1 |
| Asia | A2 |

| A | |
|---|---|
| **key** | **route** |
| python | B1 |
| java | B2 |

Server 2

## Trends evolve with time

Correlations between keys change frequently.

## Trends evolve with time

Correlations between keys change frequently.



Figure: #nevertrump, in March 2016

Locality decay

- Keys correlations evolve with time.

Locality decay

- ▶ Keys correlations evolve with time.
- ▶ Routing tables optimized by examining old data lead to decreased locality.

## Locality decay

- Keys correlations evolve with time.
- Routing tables optimized by examining old data lead to decreased locality.

## Reconfiguration

- We re-compute the tables every $N$ minutes.

## Locality decay

- Keys correlations evolve with time.
- Routing tables optimized by examining old data lead to decreased locality.

## Reconfiguration

- We re-compute the tables every $N$ minutes.
- Difficulty: keep the state consistent.

Solution: online reconfiguration protocol
- update the routing tables in a live system
- without losing any message and state

① Get statistics

## Reconfiguration protocol



① Get statistics
② Send statistics

① Get statistics
② Send statistics
*Partition graph, compute routing tables*

① Get statistics
② Send statistics
*Partition graph, compute routing tables*
③ Send reconfiguration

① Get statistics
② Send statistics
*Partition graph, compute routing tables*
③ Send reconfiguration
④ Send ACK

## Reconfiguration protocol



① Get statistics
② Send statistics

*Partition graph, compute routing tables*

③ Send reconfiguration
④ Send ACK
⑤ Propagate

① Get statistics
② Send statistics
*Partition graph, compute routing tables*
③ Send reconfiguration
④ Send ACK
⑤ Propagate
⑥ Transfer key states

① Get statistics
② Send statistics

*Partition graph, compute routing tables*

③ Send reconfiguration
④ Send ACK
⑤ Propagate
⑥ Transfer key states

*Propagate to next operator*

## Datasets

- From Flickr and Twitter
- Fields: location (country or place), hashtag
- Size: 173M records (Flickr), 100M (Twitter)

## Datasets

- From Flickr and Twitter
- Fields: location (country or place), hashtag
- Size: 173M records (Flickr), 100M (Twitter)

## Setup

- $8\times$ 128 GB RAM, 20 cores.
- Computation of aggregated statistics (stateful workers).
- Parallelism (2..6), network speed (1Gb/s | 10Gb/s), message size (0..20kB).

Great speed-up when network is the bottleneck.

Great speed-up when network is the bottleneck.

Highly dependent on message size.

Throughput (Ktuples/s) on 10Gb/s network, parallelism 6



(a) message size=4kB          (b) message size=8kB

Throughput (Ktuples/s) on 1Gb/s network, parallelism 6



(a) message size=4kB          (b) message size=8kB

## Evaluation – Flickr

Average throughput with 1Gb/s network, 4kB message size



Figure: Average throughput, measured after the first reconfiguration.

Locality, with parallelism 6

Locality, with parallelism 6

Locality, with parallelism 6

Locality when changing the number of collected key correlations

Scheduling: placement of operators on servers

- Using the topology [ABQ13]
- Using observed communication patterns [ABQ13]
- Using observed and/or estimated CPU and memory patterns [FB15, PHH$^+$15]

Scheduling: placement of operators on servers

- Using the topology [ABQ13]
- Using observed communication patterns [ABQ13]
- Using observed and/or estimated CPU and memory patterns [FB15, PHH+15]

Load balancing: limit impact of data skew

- Partial key grouping [NMG+15]
- Special routing for frequent keys [RQA+15]

## Related work

### Scheduling: placement of operators on servers

- ▶ Using the topology [ABQ13]
- ▶ Using observed communication patterns [ABQ13]
- ▶ Using observed and/or estimated CPU and memory patterns [FB15, PHH$^+$15]

### Load balancing: limit impact of data skew

- ▶ Partial key grouping [NMG$^+$15]
- ▶ Special routing for frequent keys [RQA$^+$15]

### Co-location of correlated keys

- ▶ Databases partitions [CJZM10], social networks [BJJL13]

GOTO E SPAGHETTI CODE

λ-blocks

*How to design a framework to compose data processing algorithms in a descriptive fashion, while reasoning on high level abstractions?*

Design goals

► A data processing abstraction

Design goals

- A data processing abstraction
- A graph of code blocks to represent an end-to-end processing system

Design goals

- A data processing abstraction
- A graph of code blocks to represent an end-to-end processing system
- Separation of concerns: low-level data operations, high-level data processing programs

Design goals

- A data processing abstraction
- A graph of code blocks to represent an end-to-end processing system
- Separation of concerns: low-level data operations, high-level data processing programs
- Maximize reuse of code

Design goals

- A data processing abstraction
- A graph of code blocks to represent an end-to-end processing system
- Separation of concerns: low-level data operations, high-level data processing programs
- Maximize reuse of code
- Compatible with existing (specialized) frameworks and possibility to mix them

Design goals

- A data processing abstraction
- A graph of code blocks to represent an end-to-end processing system
- Separation of concerns: low-level data operations, high-level data processing programs
- Maximize reuse of code
- Compatible with existing (specialized) frameworks and possibility to mix them
- Graph manipulation toolkit

Design goals

- A data processing abstraction
- A graph of code blocks to represent an end-to-end processing system
- Separation of concerns: low-level data operations, high-level data processing programs
- Maximize reuse of code
- Compatible with existing (specialized) frameworks and possibility to mix them
- Graph manipulation toolkit
- Bring simplicity to large-scale data processing

Topologies

```python
"""Counts system users.
"""

def main():
    with open('/etc/passwd') as f:
        return len(f.readlines())

if __name__ == '__main__':
    print(main())
```

```python
"""Counts system users.
"""


def main():
    with open('/etc/passwd') as f:
        return len(f.readlines())


if __name__ == '__main__':
    print(main())



$ wc -l /etc/passwd
```

```python
"""Counts system users.
"""

def main():
    with open('/etc/passwd') as f:
        return len(f.readlines())

if __name__ == '__main__':
    print(main())
```

```
$ wc -l /etc/passwd
```

## Topologies

```
---
name: count_users
description: Count number of system users
modules: [lb.blocks.foo]
---
- block: readfile
  name: my_readfile
  args :
    filename: /etc/passwd

- block: count
  name: my_count
  inputs :
    data: my_readfile.result
```

## Blocks

- read_http
- plot_bars
- show_console
- write_line
- write_lines
- split
- concatenate
- map_list
- flatMap
- flatten_list
- group_by_count
- sort
- get_spark_context
- spark_readfile
- spark_text_to_words
- spark_map
- spark_filter

- spark_flatMap
- spark_mapPartitions
- spark_sample
- spark_union
- spark_intersection
- spark_distinct
- spark_groupByKey
- spark_reduceByKey
- spark_aggregateByKey
- spark_sortByKey
- spark_join
- spark_cogroup
- spark_cartesian
- spark_pipe
- spark_coalesce
- spark_repartition
- spark_reduce

- spark_collect
- spark_count
- spark_first
- spark_take
- spark_takeSample
- spark_takeOrdered
- spark_saveAsTextFile
- spark_countByKey
- spark_foreach
- spark_add
- spark_swap
- twitter_search
- cat
- grep
- cut
- head
- tail

```python
@block(engine='localpython')
def take(n: int=0):
    """Truncates a list of integers.

    :param int n: The length of the desired result.
    :input List[int] data: The list of items to truncate.
    :output List[int] result: The truncated result.
    """
    def inner(data: List[int])->ReturnType[List[int]]:
        assert n <= len(data)
        return ReturnEntry(result=data[:n])
    return inner
```

**count_pb**

## Sub-topologies

```
---
name: count_pb
---
- block: filter
  name: filter
  args:
    contains: error
  inputs:
    data: $inputs.data
- block: count
  name: count
  inputs:
    data: filter.result
```

## Sub-topologies

```
---
name: count_pb
---
- block: filter
  name: filter
  args:
    contains: error
  inputs:
    data: $inputs.data
- block: count
  name: count
  inputs:
    data: filter.result
```

```
---
name: foo_errors
---
- block: readfile
  name: readfile
  args:
    filename: foo.log

- topology : count_pb
  name: count_pb
  bind_in :
    data: readfile.result
  bind_out :
    result: count.result

- block: print
  name: print
  inputs:
    data: count_pb.result
```

## Architecture



Block libraries

## Architecture

- Verification (e.g. type checking)

- Verification (e.g. type checking)
- Instrumentation

- Verification (e.g. type checking)
- Instrumentation
- Caching

- Verification (e.g. type checking)
- Instrumentation
- Caching
- Debugging tools

- Verification (e.g. type checking)
- Instrumentation
- Caching
- Debugging tools
- Optimizations

- Verification (e.g. type checking)
- Instrumentation
- Caching
- Debugging tools
- Optimizations
- Monitoring

- Verification (e.g. type checking)
- Instrumentation
- Caching
- Debugging tools
- Optimizations
- Monitoring
- Program reasoning and semantics

- Reasoning on the computation graph as a high-level object

- Reasoning on the computation graph as a high-level object
- Plugin system

- Reasoning on the computation graph as a high-level object
- Plugin system
- Hooks:
  - `before_graph_execution`
    pre-processing, optimizations, verifications

## Graph manipulations

- Reasoning on the computation graph as a high-level object
- Plugin system
- Hooks:
  - `before_graph_execution`
    pre-processing, optimizations, verifications
  - `after_graph_execution`
    post-processing

- Reasoning on the computation graph as a high-level object
- Plugin system
- Hooks:
  - `before_graph_execution`
    pre-processing, optimizations, verifications
  - `after_graph_execution`
    post-processing
  - `before_block_execution`
    observation, optimizations

## Graph manipulations

- Reasoning on the computation graph as a high-level object
- Plugin system
- Hooks:
    - `before_graph_execution`
      pre-processing, optimizations, verifications
    - `after_graph_execution`
      post-processing
    - `before_block_execution`
      observation, optimizations
    - `after_block_execution`
      observation

```python
by_block = {} # timing by block: begin, duration

@before_block_execution
def store_begin_time(block):
    name = block.fields['name']
    by_block[name]['begin'] = time.time()
```

```python
by_block = {} # timing by block: begin, duration

@before_block_execution
def store_begin_time(block):
    name = block.fields['name']
    by_block[name]['begin'] = time.time()


@after_block_execution
def store_end_time(block, results):
    name = block.fields['name']
    by_block[name]['duration'] = \
      time.time() - by_block[name]['begin']
```

```python
@after_graph_execution
def show_times(results):
    longest_first = sorted(by_block, reverse=True)
    for blockname in longest_first:
        print('{}\t{}'.format(
            blockname,
            by_block[blockname]['duration'])
```

Graph manipulation example: instrumentation

| block | duration (ms) |
|---|---|
| read http | 818 |
| write lines | 54 |
| grep | 49 |
| split | 20 |

Setup

- Wordcount over https: local machine, 8 cores, 16 GB RAM
- Wordcount over disk: local machine, 8 cores, 16 GB RAM
- PageRank on Spark: Spark on 1 server (24 cores, 128 GB RAM)

## Performances



Figure: Wordcount over https: Twitter feed.

## Evaluation

### Performances



Figure: Wordcount over disk: Wikipedia dataset.

## Performances



Figure: PageRank on Wikipedia hyperlinks with Spark.

Maximum overhead measured per topology: 50 ms

Dataflow programming

- ▶ ML pipelines: scikit-learn [PVG+11], Spark [The17a], Orange framework [DCE+13]
- ▶ Real-time: Apache Beam [apa], StreamPipes [RKHS15]

## Related work

### Dataflow programming

- ML pipelines: scikit-learn [PVG$^+$11], Spark [The17a], Orange framework [DCE$^+$13]
- Real-time: Apache Beam [apa], StreamPipes [RKHS15]

### Blocks programming

- Recognition over recall, immediate feedback [BGK$^+$17]

## Dataflow programming

- ML pipelines: scikit-learn [PVG$^+$11], Spark [The17a], Orange framework [DCE$^+$13]
- Real-time: Apache Beam [apa], StreamPipes [RKHS15]

## Blocks programming

- Recognition over recall, immediate feedback [BGK$^+$17]

## Graphs from configuration

- Pyleus [Yel16], Storm Flux [The17b]

## Dataflow programming

- ML pipelines: scikit-learn [PVG$^+$11], Spark [The17a], Orange framework [DCE$^+$13]
- Real-time: Apache Beam [apa], StreamPipes [RKHS15]

## Blocks programming

- Recognition over recall, immediate feedback [BGK$^+$17]

## Graphs from configuration

- Pyleus [Yel16], Storm Flux [The17b]

## Other

- "Serverless" architectures and stateless functions [JVSR17]

## Context

Computer systems to process large quantities of data.

### Context

Computer systems to process large quantities of data.

### Problems: how to design...

- An industrial system to handle monitoring data and make predictions about future failures?
- An algorithm to improve locality in distributed streaming engines?
- A framework to compose data processing algorithms in a descriptive fashion, while reasoning on high level abstractions?

**Metrics prediction**     **Locality routing**     **$\lambda$-blocks**

## Contributions

| | **Metrics prediction** | **Locality routing** | **$\lambda$-blocks** |
|---|---|---|---|
| **What it is** | Industrial system | Online routing library | Data processing abstraction |

## Contributions

|  | **Metrics prediction** | **Locality routing** | **$\lambda$-blocks** |
|---|---|---|---|
| **What it is** | Industrial system | Online routing library | Data processing abstraction |
| **Layer** | End-to-end | Low | High |

|  | **Metrics prediction** | **Locality routing** | **$\lambda$-blocks** |
|---|---|---|---|
| **What it is** | Industrial system | Online routing library | Data processing abstraction |
| **Layer** | End-to-end | Low | High |
| **Improves** | **Uptimes** | **Throughput** | **Programmability** |

Metrics prediction in monitoring systems

- ▶ Predictions on long-term global trends
- ▶ Ticketing mechanism

## Metrics prediction in monitoring systems

- ▶ Predictions on long-term global trends
- ▶ Ticketing mechanism

## Locality data routing

- ▶ Replace binary locality/non-locality with distance
- ▶ Smarter way to determine when to reschedule
- ▶ Extend to more complex topologies

$\lambda$-blocks

- ▶ Explore more graph manipulation abstractions (complexity analysis, serialization, verification. . . )
- ▶ Streaming and online operations
- ▶ Tight integration with clusters (data storage, caches, etc)

Thanks! Questions?

## Signature algorithm



$$H(B) = h(B.name, \quad \text{block name (not instance name)}$$
$$B.args, \quad \text{list of (name, value) tuples}$$
$$B.inputs) \quad \text{list of (name, H(block), connector) tuples}$$

# λ-blocks

## Evaluation

### Engine instrumentation



Figure: Wordcount program running under different setups.
(1) Startup (modules import, etc); (2) Blocks registry creation, block modules import; (3) Plugin import; (4) YAML parsing and graph creation; (5) Graph checks; (6) Graph execution.

# Metrics prediction in monitoring systems

## Database schema

**metrics**

| | |
|---|---|
| metric_id | uuid |
| metric_name | text |
| group_id | uuid |

**predictions**

| | |
|---|---|
| metric_id | uuid |
| timestamp | int |
| predicted_values | list |

**measurements**

| | |
|---|---|
| metric_id | uuid |
| timestamp | int |
| warn | text |
| crit | text |
| max | double |
| min | double |
| value | double |
| metric_name | text |
| metric_unit | text |

- *Data Center operators verifying network cable integrity*, CC-BY-SA, https://commons.wikimedia.org/wiki/File: Dc_cabling_50.jpg
- *Tokyo metro map*, http://bento.com/subtop5.html
- *Goto e spaghetti code*, http://blogbv2.altervista.org/ HD/il-goto-e-la-buona-programmazione-parte-ii/

# Bibliography I

📄 Leonardo Aniello, Roberto Baldoni, and Leonardo Querzoni.
Adaptive online scheduling in storm.
In *Proceedings of the 7th ACM International Conference on Distributed Event-based Systems*, DEBS '13, pages 207–218. ACM, 2013.

📄 Apache Beam.
https://beam.apache.org/.

📄 David Bau, Jeff Gray, Caitlin Kelleher, Josh Sheldon, and Franklyn Turbak.
Learnable programming: Blocks and beyond.
*Commun. ACM*, 60(6):72–80, May 2017.

📄 Xiao Bai, Arnaud Jégou, Flavio Junqueira, and Vincent Leroy.
Dynasore: Efficient in-memory store for social applications.
In *Middleware 2013 - ACM/IFIP/USENIX 14th International Middleware Conference, Beijing, China, December 9-13, 2013, Proceedings*, pages 425–444, 2013.

📄 T. Chalermarrewong, T. Achalakul, and S. C. W. See.
Failure prediction of data centers using time series and fault tree analysis.
In *2012 IEEE 18th International Conference on Parallel and Distributed Systems*, pages 794–799, Dec 2012.

📄 Carlo Curino, Evan Jones, Yang Zhang, and Sam Madden.
Schism: A workload-driven approach to database replication and partitioning.
*Proc. VLDB Endow.*, 3(1-2):48–57, September 2010.

📄 Janez Demšar, Tomaž Curk, Aleš Erjavec, Črt Gorup, Tomaž Hočevar, Mitar Milutinovič, Martin Možina, Matija Polajnar, Marko Toplak, Anže Starič, Miha Štajdohar, Lan Umek, Lan Žagar, Jure Žbontar, Marinka Žitnik, and Blaž Zupan.
Orange: Data mining toolbox in python.
*Journal of Machine Learning Research*, 14:2349–2353, 2013.

📄 Lorenz Fischer and Abraham Bernstein.
Workload scheduling in distributed stream processors using graph partitioning.
In *2015 IEEE International Conference on Big Data, Big Data 2015, Santa Clara, CA, USA, October 29 - November 1, 2015*, pages 124–133, 2015.

📄 Eric Jonas, Shivaram Venkataraman, Ion Stoica, and Benjamin Recht.
Occupy the cloud: Distributed computing for the 99%.
*arXiv preprint arXiv:1702.04024*, 2017.

📄 Lei Li, Chieh-Jan Mike Liang, Jie Liu, Suman Nath, Andreas Terzis, and Christos Faloutsos.
Thermocast: A cyber-physical forecasting model for datacenters.
In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '11, pages 1370–1378, New York, NY, USA, 2011. ACM.

📄 Microsoft cloud azure.
https://docs.microsoft.com/en-us/azure/
machine-learning/
machine-learning-algorithm-choice.

📄 Muhammad Anis Uddin Nasir, Gianmarco De Francisci
Morales, David García-Soriano, Nicolas Kourtellis, and Marco
Serafini.
The power of both choices: Practical load balancing for
distributed stream processing engines.
In *31st IEEE International Conference on Data Engineering,
ICDE*, pages 137–148, 2015.

Boyang Peng, Mohammad Hosseini, Zhihao Hong, Reza Farivar, and Roy Campbell.
R-storm: Resource-aware scheduling in storm.
In *Proceedings of the 16th Annual Middleware Conference*, Middleware '15, pages 149–161, New York, NY, USA, 2015. ACM.

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay.
Scikit-learn: Machine learning in Python.
*Journal of Machine Learning Research*, 12:2825–2830, 2011.

📄 Dominik Riemer, Florian Kaulfersch, Robin Hutmacher, and Ljiljana Stojanovic.
Streampipes: solving the challenge with semantic stream processing pipelines.
In *Proceedings of the 9th ACM International Conference on Distributed Event-Based Systems*, pages 330–331. ACM, 2015.

📄 Nicoló Rivetti, Leonardo Querzoni, Emmanuelle Anceaume, Yann Busnel, and Bruno Sericola.
Efficient key grouping for near-optimal load balancing in stream processing systems.
In *Proceedings of the 9th ACM International Conference on Distributed Event-Based Systems*, DEBS '15, pages 80–91, New York, NY, USA, 2015. ACM.

📄 The Apache Spark developers.
ML Pipelines.
https:
//spark.apache.org/docs/latest/ml-pipeline.html,
2017.

📄 The Apache Storm developers.
Flux.
http://storm.apache.org/releases/2.0.0-SNAPSHOT/
flux.html, 2017.

📄 YelpArchive.
Pyleus.
https://github.com/YelpArchive/pyleus, 2016.

# Bibliography IX

📄 Zabbix prediction triggers.
https://www.zabbix.com/documentation/3.0/manual/
config/triggers/prediction.