# Towards a Hybrid Cloud & Edge Orchestrator for Mining Exascale Distributed Streams

**Herodotos Herodotou**

Cyprus University of Technology

herodotos.herodotou@cut.ac.cy

# Acknowledgments

## Nicolas Kourtellis

## Maciej Grzenda

## Piotr Wawrzyniak

## Albert Bifet

MSc/PhD Computer Science

Assistant Professor

2012 · 2014 · 2022

Microsoft

aster data
big data. fast insights.

YAHOO! RESEARCH

Jim Gray Doctoral Dissertation Award Honorable Mention

# Data Intensive Computing Research Lab

## Research Areas

- Large-scale data processing systems (e.g., MapReduce, Spark)
- Centralized and distributed database systems
- Cloud computing (compute, storage, and networking)
- Data-driven applications (maritime, tourism, social computing)

## Research Team

- Supervise: 1 postdoc & 3 PhD students
- Co-supervise: 2 postdocs & 1 PhD student

## Equipment

- 15-node local private cluster
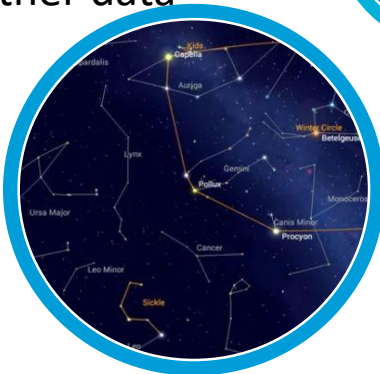
# Our Data-driven World

## Massive Data

**Analysis** → **Insights**

### Medicine
- MRI & CT scans
- Patient records

### Entertainment
- Internet images
- Movie & music clips

### Science
- Genomics
- Astronomy
- Weather data

### Business
- Product sales
- Stock market
- Customer data

### Humanities
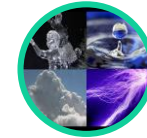- Social interactions
- Historical documents

### Insights
→ Ad placement

→ Scientific breakthroughs

→ Business process efficiencies

→ Personalized recommendations

→ Improved healthcare

→ Fraud detection

# Evolution of Big Data Analytics Systems

Stream Data Analytics

Dataflow Analytics

MapReduce Analytics

Database Analytics

Data Warehouses

2000s

2005s

2010s

2015s

2020s

IaaS

PaaS

**However**

➢ Cloud heterogeneity & management overhead

➢ Vendor lock-in

➢ 4Vs already stress current, non-scalable infra

**S2CE: Stream AI + Cloud + Edge**

➢ 10s of billions of devices

➢ 4/5 Vs: Volume, Velocity, Variety, Veracity

➢ Need for data modeling to extract Value (5$^{th}$ V)

**S2CE**
**Telco Cloud**

Upstream ————

Downstream ————

**S2CE** edge node  O

- Big data stream processing systems
- Cloud resource management and tuning
- Distributed stream processing at the edge
- Machine and deep learning over data streams
- Data transformation techniques

Google Cloud Dataflow

Azure

amazon Kinesis

# Cloud Resource Management and Tuning



**Decisions:**

- Task parallelism
- Micro-batch size
- ...

- Number of cores
- Memory settings
- ...

- Number of nodes
- VM/Container specs
- ...

# Cloud Resource Management and Tuning – Approaches

| Cost Modeling | Use cost models & statistics to find optimal settings |
|---|---|
| Simulation-based | Use simulator to estimate application performance |
| Experiment-driven | Execute application with different settings iteratively |
| Machine Learning | Use machine learning to model application performance |
| Adaptive | Change configurations while application is running |

- ➢ Benefits of edge:
  - ❖ Reduce end-to-end latency and communication costs
  - ❖ Enable services to react to events locally
  - ❖ Offload processing from the cloud
- ➢ Challenges:
  - ❖ Computing, storage, and network resources are constrained
  - ❖ Deployment of data stream processing applications onto heterogeneous infrastructure has been proven to be NP-hard

# Distributed Stream Processing at the Edge – Projects

| Aspect | EdgeX Foundry | Azure IoT Edge | Apache Edgent | CORD | Akraino Edge Stack |
|---|---|---|---|---|---|
| Interface | Restful API | Web service | API | API or XOS-GUI | N/A |
| OS support | Various | Various | Various | Ubuntu | Linux |
| Programming framework | Not provided | Java, .NET, C, Python, etc. | Java | Shell, Python | N/A |
| Applications | IoT | Unrestricted | IoT | Unrestricted | Unrestricted |
| Deployment | Dynamic | Dynamic | Static | Dynamic | Dynamic |
| Target user | General users | General users | General users | Network ops | Network ops |
| Virtualization | Container | Container | JVM | VM & Container | VM & Container |
| Limitation | Lack of program-ming interface | Azure services chargeable | Limited to data analysis | Unable to be offline | Unable to be offline |
| Scalability | Scalable | Scalable | Not scalable | Scalable | Scalable |
| Mobility | Not supported | Not supported | Not supported | Support | Support |

*Liu et al. A survey on edge computing systems and tools. IEEE*

**Existing libraries/systems:**

- Massive Online Analysis (MOA)
  - algorithms for streaming classification, clustering, and change detection

- Vowpal Wabbit
  - based on the perceptron algorithm with a focus on reinforcement learning

- Jubatus
  - tight coupling between the ML library and the underlying custom-built DSPS

- Apache SAMOA
  - distributed computation of several ML algorithms over four DSPSs

**Challenges:**

➤ Data Availability

  ❖ a train, test and predict approach is not applicable for stream data; models are susceptible to changes

➤ Real-Time Streaming

  ❖ need to reduce time to train models dramatically

➤ Concept Drift

  ❖ models must adapt to patterns evolving over time by detecting changes quickly

➤ IID Random Variables

  ❖ statistically independent variables cannot be guaranteed for the overall population

➢ **Data Preprocessing**
 ❖ filtering, format conversion, and multiplexing/demultiplexing

➢ **Dimensionality reduction**
 ❖ statistical inference methods using hashing projections
 ❖ different subspace tracking methods

➢ **Stream sampling**
 ❖ allows one-pass algorithms for analysing big data streams
 ❖ uniform or biased sampling along with reduction of the problem space

➢ **Synthetic data stream generator**
 ❖ infer underlying statistical distributions of the real data
 ❖ do not work for streams with concept drifts
 ❖ protecting privacy and confidentiality is hard

✘ means no support
! means partial support
✓ means good/full support

| Features/Capabilities | Apache Storm | Apache Samza | Apache Spark | Apache Flink | Apache Apex | Apache Beam | Google CD | MS Azure ML | AWS Kinesis | MOA | Vowpal Wabbit | Jubatus | Apache SAMOA | Desired Platform |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Stream integration components | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ! | ! | ! | ✘ | ✘ | ! | ✓ | ✓ |
| Data preprocessing and fusion | ! | ✘ | ! | ! | ✘ | ✘ | ! | ✘ | ! | ! | ! | ✓ | ! | ✓ |
| Built-in synthetic data generator | ✘ | ✘ | ! | ! | ! | ! | ✘ | ✘ | ✓ | ✓ | ✘ | ✘ | ! | ✓ |
| Stream-based machine learning | ! | ! | ✓ | ✓ | ! | ! | ✓ | ! | ! | ✓ | ✓ | ✓ | ✓ | ✓ |
| Stream-based deep learning | ✘ | ✘ | ! | ✘ | ! | ✘ | ✘ | ✘ | ✘ | ✘ | ✘ | ✘ | ✘ | ✓ |
| Resource management | ✓ | ! | ✓ | ✓ | ✓ | ✓ | ✓ | ! | ✓ | ✘ | ✘ | ! | ! | ✓ |
| Cloud-Edge orchestration | ✘ | ✘ | ✘ | ✘ | ✘ | ✘ | ! | ! | ! | ✘ | ✘ | ✘ | ✘ | ✓ |
| Distributed Platform | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✘ | ✘ | ✓ | ✓ | ✓ |
| Open license (Apache preferred) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✘ | ✘ | ✘ | ! | ! | ! | ✓ | ✓ |

| Expected Industrial Challenge |
|---|
| Heterogeneity<br>Scalability<br>Data-in-motion and data-at-rest |
| Hybrid (central+edge) big data architectures<br>Decentralization & edge |
| Data/AI/predictive/prescriptive analytics<br>Stream analytics frameworks & processing<br>Advanced business analytics |
| Heterogeneity<br>Semantic interoperability<br>Data quality<br>Distributed trust infrastructures |

Scalable Data processing

Edge vs. Cloud infrastructure

ML/DL-based analytics

Data fusion and input / output

*BDVA: https://www.bdva.eu

# S2CE Design Objectives

| Expected Industrial Challenge | (Objective) How does S2CE address the challenge? |
|---|---|
| Heterogeneity<br>Scalability<br>Data-in-motion and data-at-rest | (O1) Handling diverse types of cloud computing resources<br>(O1) Distributed and parallelized dynamic analytics for real-time learning<br>(O1) Processing data seamlessly at the same time without extra system overhead |
| Hybrid (central+edge) big data architectures<br>Decentralization & edge | (O2) Optimizing an efficient mixture of central and edge resources<br>(O2) Computing at edge for faster, more scalable, energy efficient processing |
| Data/AI/predictive/prescriptive analytics<br>Stream analytics frameworks & processing<br>Advanced business analytics | (O3) Using distributed deep and machine learning<br>(O3) Minimal development effort, scalability, processing speed<br>(O3) Intelligence to empower companies for accurate, instant, data-driven decisions |
| Heterogeneity<br>Semantic interoperability<br>Data quality<br>Distributed trust infrastructures | (O4) Handling diverse data, modeling, and input/output interfaces<br>(O4) Facilitating data and model exchange between vertical data silos<br>(O4) Providing curation methods for data filtering, quality assessment, improvement<br>(O4) Managing data in anonymized and decentralized fashion |

# S2CE Architecture Overview

**S2CE**

**Upstream S2CE**

Data-at-rest

Data-in-motion

**Upstream DSPE Infrastructure**

**Input Interface**

**Transformations**
- Data Fusion & Preprocessing
- Model Interpretations, Evolution & Explainability
- True data stream generators

**Machine Learning**
- Machine learning algorithms for streaming data
- Adaptive Deep Learning for streaming data

**Cloud Management**
- Cloud Resource Allocation
- Cloud/Engine Optimization
- App Optimization & Selftuning

**Edge Management**
- Energy-efficient Edge
- Multi-Access Edge
- Cloud-Edge Movement

**Cloud/Edge Computing Infrastructure Interface**

**Output Interface**

**Downstream S2CE**

Interpretable Models

Performance Indicators

Predictions & Rules

**Downstream DSPE Infrastructure**

**Cloud & Edge**
| Local | Flink | Storm | Samza | Spark | ... |

**S2CE**

**Upstream S2CE**

Data-at-rest

Data-in-motion

**Input Interface**

**Transformations**

Data Fusion & Preprocessing

Model Interpretations, Evolution & Explainability

True data stream generators

**Machine Learning**

Machine learning algorithms for streaming data

Adaptive Deep Learning for streaming data

**Cloud Management**

Cloud Resource Allocation

Cloud/Engine Optimization

App Optimization & Selftuning

**Edge Management**

Energy-efficient Edge

Multi-Access Edge

Cloud-Edge Movement

**Cloud/Edge Computing Infrastructure Interface**

**Output Interface**

**Downstream S2CE**

Interpretable Models

Performance Indicators

Predictions & Rules

**Upstream DSPE Infrastructure**

**Downstream DSPE Infrastructure**

**Cloud & Edge**

Local | Flink | Storm | Samza | Spark | ...

**dicl** — Data Intensive Computing Research Lab

**S2CE**

**Upstream S2CE**

Data-at-rest

Data-in-motion

**Input Interface**

**Transformations**
- Data Fusion & Preprocessing
- Model Interpretations, Evolution & Explainability
- True data stream generators

**Machine Learning**
- Machine learning algorithms for streaming data
- Adaptive Deep Learning for streaming data

**Cloud Management**
- Cloud Resource Allocation
- Cloud/Engine Optimization
- App Optimization & Selftuning

**Edge Management**
- Energy-efficient Edge
- Multi-Access Edge
- Cloud-Edge Movement

**Upstream DSPE Infrastructure**

**Cloud/Edge Computing Infrastructure Interface**

**Output Interface**

**Downstream S2CE**

Interpretable Models

Performance Indicators

Predictions & Rules

**Downstream DSPE Infrastructure**

**Cloud & Edge**
Local | Flink | Storm | Samza | Spark | …

S2CE

**Upstream S2CE**

**Input Interface**

**Transformations**
- Data Fusion & Preprocessing
- Model Interpretations, Evolution & Explainability
- True data stream generators

**Machine Learning**
- Machine learning algorithms for streaming data
- Adaptive Deep Learning for streaming data

**Cloud Management**
- Cloud Resource Allocation
- Cloud/Engine Optimization
- App Optimization & Selftuning

**Edge Management**
- Energy-efficient Edge
- Multi-Access Edge
- Cloud-Edge Movement

**Output Interface**

**Cloud/Edge Computing Infrastructure Interface**

Data-at-rest

Data-in-motion

**Upstream DSPE Infrastructure**

**Downstream S2CE**

Interpretable Models

Performance Indicators

Predictions & Rules

**Downstream DSPE Infrastructure**

**Cloud & Edge**

| Local | Flink | Storm | Samza | Spark | ... |

**S2CE**

**Upstream S2CE**

Data-at-rest

Data-in-motion

**Upstream DSPE Infrastructure**

**Input Interface**

**Transformations**

Data Fusion & Preprocessing

Model Interpretations, Evolution & Explainability

True data stream generators

**Machine Learning**

Machine learning algorithms for streaming data

Adaptive Deep Learning for streaming data

**Cloud Management**

Cloud Resource Allocation

Cloud/Engine Optimization

App Optimization & Selftuning

**Edge Management**

Energy-efficient Edge

Multi-Access Edge

Cloud-Edge Movement

**Output Interface**

**Cloud/Edge Computing Infrastructure Interface**

**Downstream S2CE**

Interpretable Models

Performance Indicators

Predictions & Rules

**Downstream DSPE Infrastructure**

**Cloud & Edge**

| Local | Flink | Storm | Samza | Spark | ... |

- APIs:
  - High-performant & secured
  - Scalable to support voluminous, fast & rich data

Number of implemented ML algorithms

Big data stream processing support

Cloud Provisioning and Orchestration

Edge Preprocessing and Movement

Data-driven Decisions based on Streams

Innovation Exchange with Apache Ecosystem's Open Source

Unifying Efforts for a Stream ML Library

# Thank you!

Herodotos Herodotou

herodotos.herodotou@cut.ac.cy

https://dicl.cut.ac.cy