



# Addressing Deployment Challenges in Data Stream Processing

**Corso di Sistemi e Architetture per Big Data**  
A.A. 2021/22  
Valeria Cardellini

Laurea Magistrale in Ingegneria Informatica

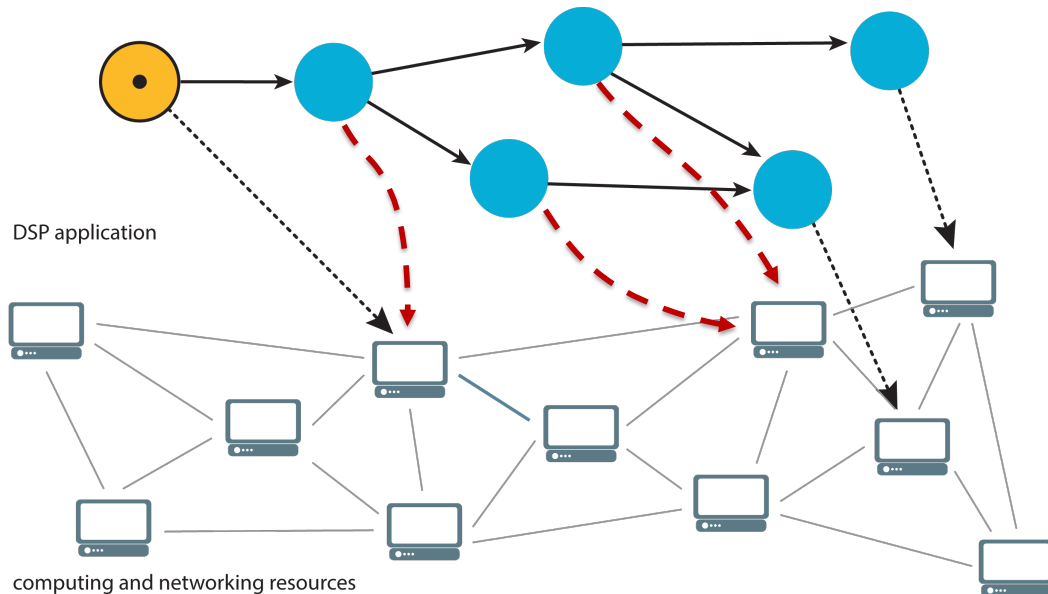
## DSP deployment challenges

---

- Let's consider two challenges when deploying DSP applications
  - a) How to place DSP operators on underlying computing infrastructure (i.e., **operator placement**)
  - b) How to determine and adapt at run-time the number of replicas per operator (i.e., **operator parallelism**)

# DSP operator placement

- Goal: determine which distributed computing nodes should **host** and **execute** each application operator, with the goal of optimizing **application QoS**



## Placement: Edge-Cloud continuum

- Edge/Fog + Cloud computing**: allows to increase scalability and availability, reduce latency, network traffic, and power consumption
- But placement becomes more challenging



## Placement: challenges

---

- Significant **network latencies**
  - E.g., geo-distributed resources
- **Heterogeneous** computing and networking resources
  - E.g., capacity limits , business constraints
- Computing/network resources can be **unavailable**
- **Data movement** around the network
- Plus peculiarities of DSP applications:
  - Computational requirements may be **unknown** a-priori and **change** continuously
  - Long-running applications

→ Need to adapt to internal and external changes

## Placement: frameworks

---

- Most frameworks use simple placement policies
- Apache Storm
  - Round Robin as default strategy
  - Resource Aware Scheduler as alternative
    - [https://storm.apache.org/releases/2.0.0/Resource\\_Aware\\_Scheduler\\_overview.html](https://storm.apache.org/releases/2.0.0/Resource_Aware_Scheduler_overview.html)
    - Takes into account resource availability on machines and resource requirements of workloads
    - But requires user to specify memory and CPU requirements for individual topology components

# Placement: different approaches

---

- Several operator placement policies in literature that address the problem but:
  - Different **assumptions** (system model, application topology, QoS attributes and metrics, ...)
  - Different **objectives**
  - Not easily comparable
- Main methodologies:
  - **Mathematical programming**
    - Optimal operator placement problem: **NP-hard**
    - Does not scale well, but provides useful insights
  - **Heuristics**
    - Majority of policies

# Placement: different approaches

---

- Who is the decision maker?
  - **Centralized** placement strategies
    - Require global view (full resource and network state, application state, workload information)
    - ✓ Capable of determining optimal global solution
    - ✗ Scalability
  - **Decentralized** placement strategies
    - Take decision based only on local information
    - ✓ Scalability, better suited for run-time adaptation
    - ✗ Optimality is not guaranteed



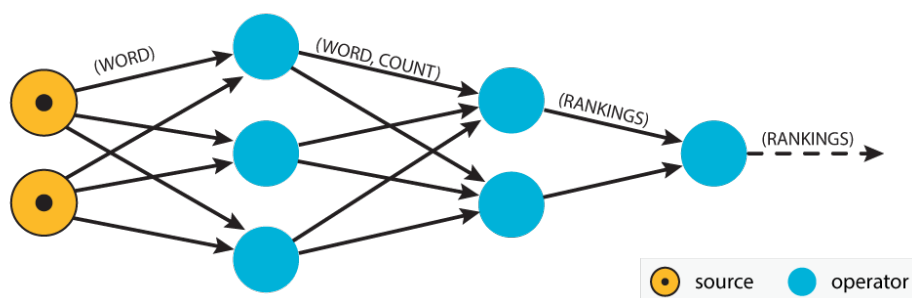
# ODP: Optimal DSP Placement

- We proposed ODP
  - Centralized policy for optimal placement of DSP applications
  - Formulated as Integer Linear Programming (ILP) problem
- Our goals:
  - To compute the **optimal placement** (of course!)
  - To provide a **unified general formulation** of the placement problem for DSP applications (but not only!)
  - To consider multiple **QoS attributes** of applications and resources
  - To provide a **benchmark** for heuristics

V. Cardellini, V. Grassi, F. Lo Presti, M. Nardelli, [Optimal Operator Placement for Distributed Stream Processing Applications](#), DEBS '16

## ODP: model

### DSP application



### Operators

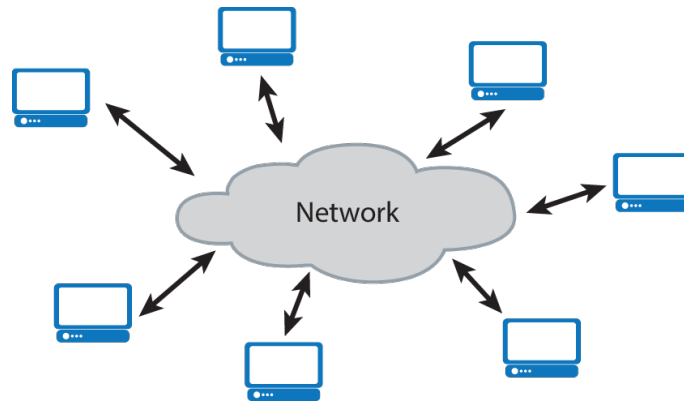
- $C_i$ : required computing resources
- $R_i$ : execution time per data unit

### Data streams

- $\lambda_{i,j}$ : data rate from operator  $i$  to  $j$

# ODP: model

## Computing and network resources



### Computing resources

- $C_u$ : amount of resources
- $S_u$ : processing speed
- $A_u$ : resource availability

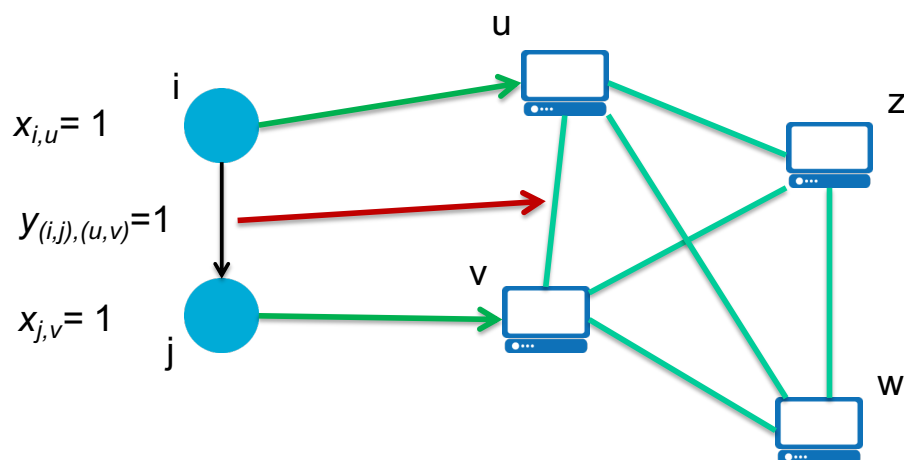
### (Logical) Network links

- $d_{u,v}$ : network delay from  $u$  to  $v$
- $B_{u,v}$ : bandwidth from  $u$  to  $v$
- $A_{u,v}$ : link availability

# ODP: model

## Decision variables

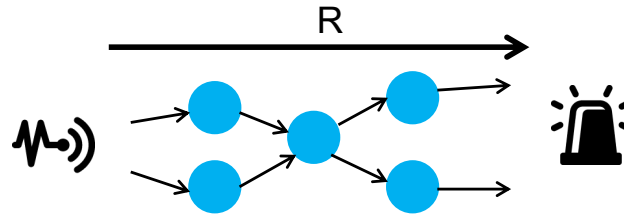
- Determine where to map DSP operators and data streams



# ODP: some QoS metrics

- Response time**

max end-to-end delay between sources and destination



- Application availability**

probability that all components/links are up and running

- Inter-node traffic**

overall network data rate

- Network usage**

in-flight bytes

$$\sum_{\text{links} \in l} \text{rate}(l) \text{Lat}(l)$$

## ODP: optimal problem formulation

Tunable knobs to set the optimal placement goals

Latency

Availability

Network bandwidth and node capacity constraints

Assignment and integer constraints

$$\max_{\mathbf{x}, \mathbf{y}, r} F(\mathbf{x}, \mathbf{y}, r)$$

subject to:

$$r \geq \sum_i \sum_u \frac{R_i}{S_u} x_{i,u} + \sum_{(i,j)} \sum_{(u,v)} d_{(u,v)} y_{(i,j),(u,v)} \quad \forall p \in \pi_G$$

$$a(\mathbf{x}, \mathbf{y}) = \sum_i \sum_u a_u x_{i,u} + \sum_{(i,j)} \sum_{(u,v)} a_{(u,v)} y_{(i,j),(u,v)}$$

$$B_{(u,v)} \geq \sum_{(i,j)} \lambda_{(i,j)} y_{(i,j),(u,v)} \quad \forall u \in V_{res}, v \in V_{res}$$

$$\sum_i C_i x_{i,u} \leq C_u \quad \forall u \in V_{res}$$

$$\sum_u x_{i,u} = 1 \quad \forall i \in V_{dsp}$$

$$x_{i,u} = \sum_v y_{(i,j),(u,v)} \quad \forall (i,j) \in E_{dsp}, u \in V_{res}$$

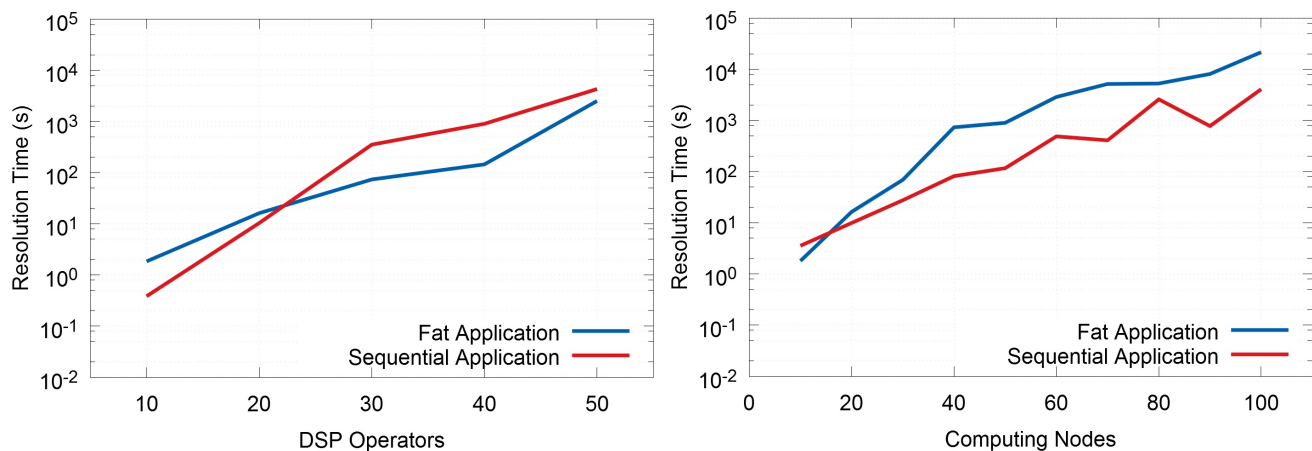
$$x_{j,v} = \sum_u y_{(i,j),(u,v)} \quad \forall (i,j) \in E_{dsp}, v \in V_{res}$$

$$x_{i,u} \in \{0, 1\} \quad \forall i \in V_{dsp}, u \in V_{res}$$

$$y_{(i,j),(u,v)} \in \{0, 1\} \quad \forall (i,j) \in E_{dsp}, (u,v) \in E_{res}$$

# ODP: scalability issue

Placement problem is **NP-hard**: does not scale well!



We need **heuristics** to compute the placement in a feasible amount of time

## Centralized placement heuristics

- Example of centralized heuristic that aims to **reduce inter-node traffic**
- **Aniello et al.:** **co-locate pairs** of communicating tasks on same computing node as to **minimize inter-node communication** and balance CPU demand

**Greedy heuristic** – Key idea:

- Rank task pairs according to exchanged traffic
- For each pair:
  - » If task pairs have not been yet assigned, assign them to same node
  - » If either is assigned, consider least loaded node and those where they have been assigned. Work out the configuration which minimizes the inter-process traffic

# Decentralized placement heuristic

- Heuristics goal: **reduce network usage**
  - Network usage metric combines link latencies and exchanged data rates among DSP operators:

$$\sum_{\text{links} \in I} \text{rate}(I) \text{Lat}(I)$$

- **Pietzuch et al.** exploit **spring relaxation** idea:
  - DSP application regarded as a system of springs, whose **minimum energy configuration** corresponds to minimizing network usage
- **Features**
  - Decentralized policy to minimize network impact
  - Adaptive to change in network conditions

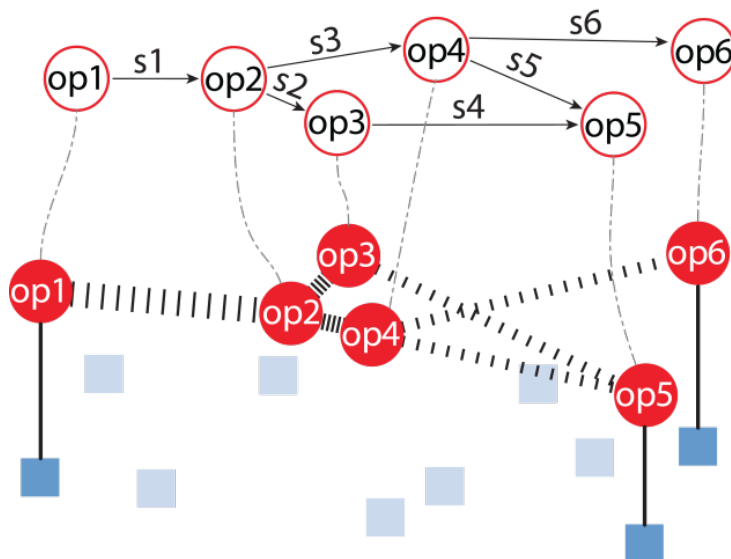
P. Pietzuch et al., [Network-aware operator placement for stream-processing systems](#), ICDE '06

V. Cardellini - SABD 2021/22

16

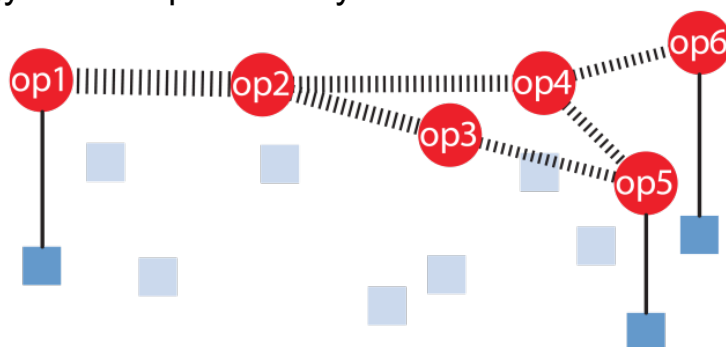
# Decentralized placement heuristic

1. Represents DSP application as an equivalent system of springs



## Decentralized placement heuristic

2. Determines operator placement in the cost space by minimizing the elastic energy of the equivalent system

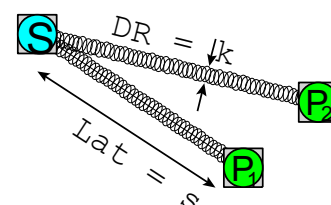


Network of springs tries to minimize potential energy  $E$

$$E = \sum_{l \in L} DR(l) Lat(l)^2$$

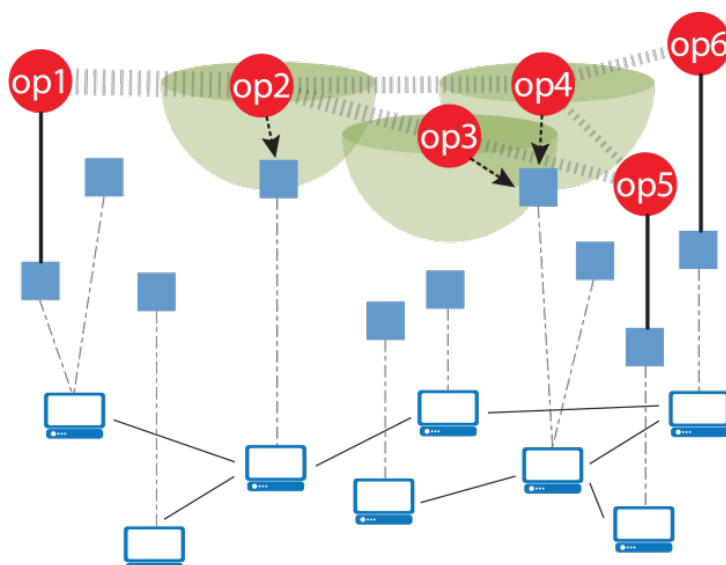
Streams as springs, that restore a force  $F = \frac{1}{2} \cdot k \cdot s$ :

- $k$  (spring constant): exchanged data rate on link
- $s$  (spring extension): latency on link



## Decentralized placement heuristic

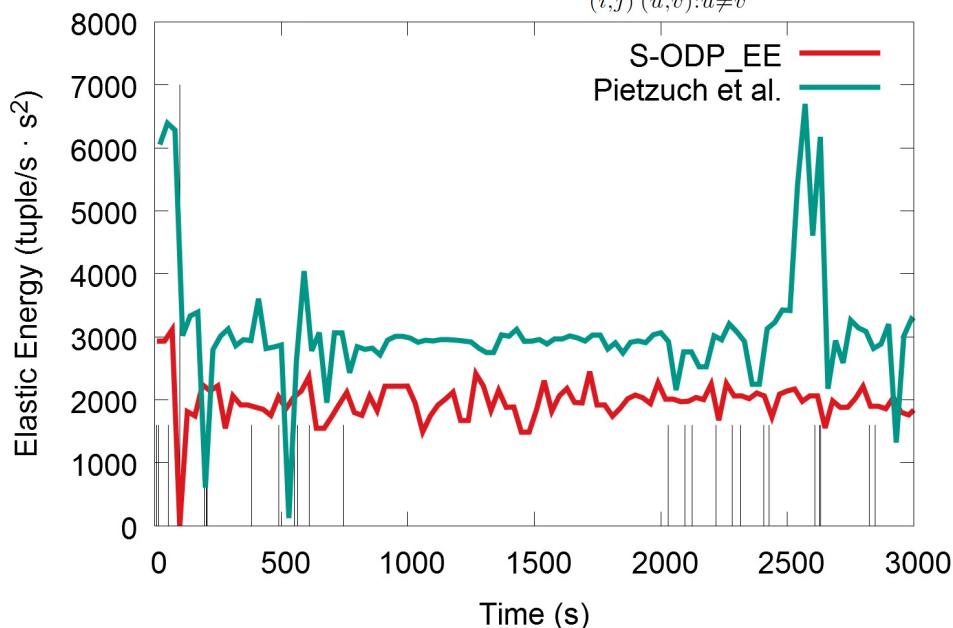
3. Maps its decision back to physical nodes



# ODP as benchmark

Distributed placement heuristic that minimizes network usage

$$\text{Pietzuch et al. : } \min EE(y) = \min \sum_{(i,j)} \sum_{(u,v): u \neq v} \lambda_{(i,j)} d_{(u,v)}^2 y_{(i,j),(u,v)}$$



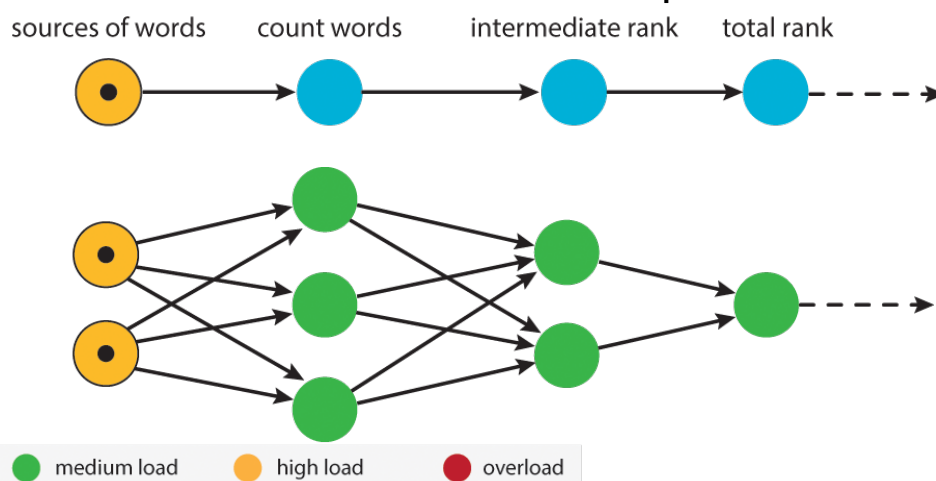
V. Cardellini, V. Grassi, F. Lo Presti, M. Nardelli, [Optimal Operator Placement for Distributed Stream Processing Applications](#), DEBS '16

V. Cardellini - SABD 2021/22

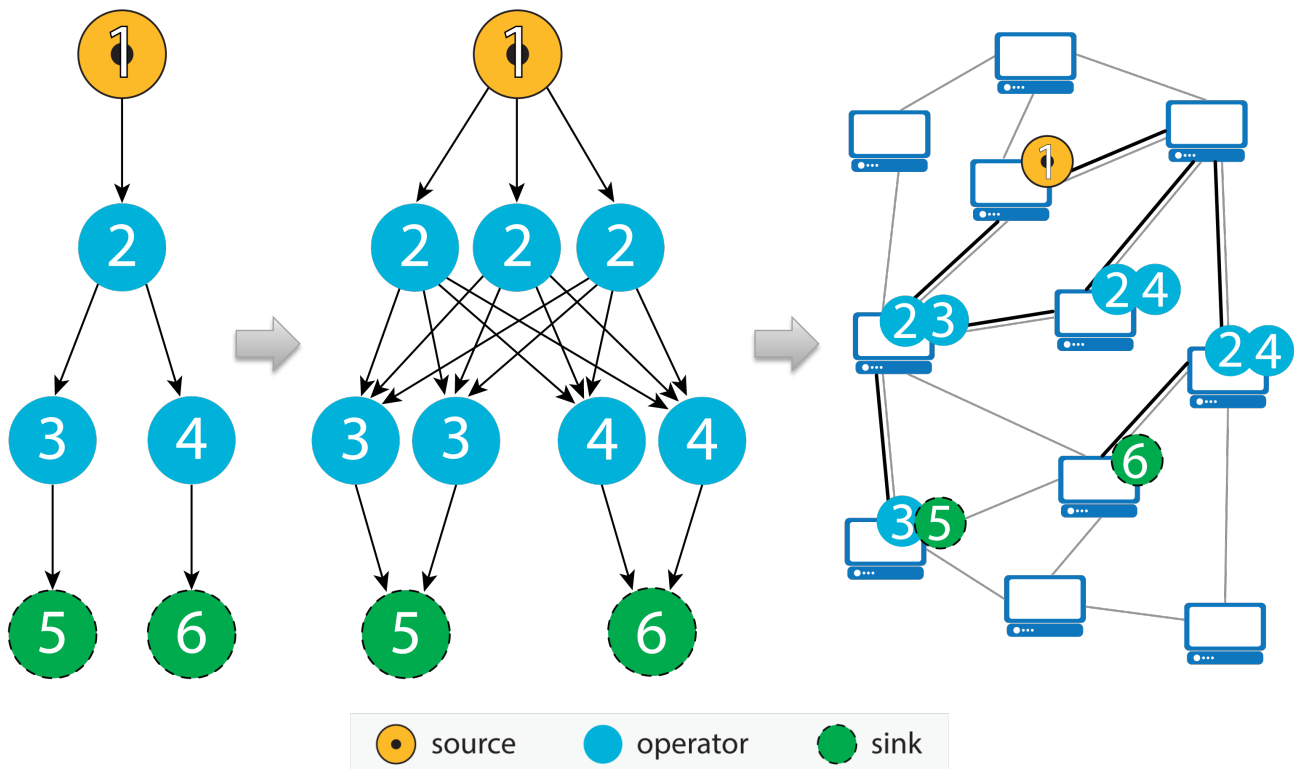
20

## Not only placement

- Stream processing workloads are characterized by:
  - High volume and production rate
- Exploit replication (i.e., [operator parallelism](#)): concurrent execution of multiple operator replicas on different data portions
- How to determine the number of replicas?



# Operator placement and replication



V. Cardellini - SABD 2021/22

22

## ODRP: Opt. DSP Replication and Placement

- We proposed **ODRP**
  - Centralized policy for **optimal replication and placement** of DSP applications
  - Formulated as Integer Linear Programming (ILP) problem that extends ODP
- Our goals:
  - Jointly determine optimal **number of replicas** and their **placement**
  - Consider multiple **QoS attributes** of applications and resources
  - Provide a **unified general formulation**
  - Provide a **benchmark** for heuristics
- Limitation: scalability, in practice we need heuristics

V. Cardellini, V. Grassi, F. Lo Presti, M. Nardelli, [Optimal operator replication and placement for distributed stream processing systems](#), *ACM Perf. Eval. Rew.*, 2017.

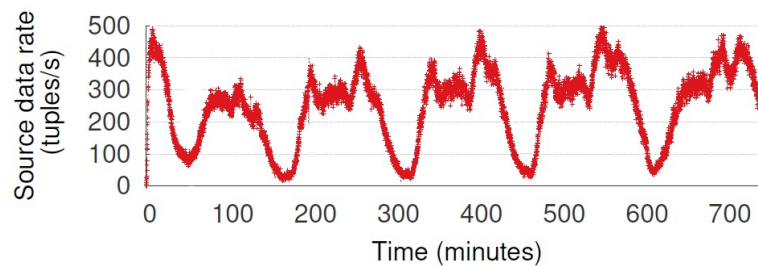
V. Cardellini - SABD 2021/22

23



# DSP deployment challenges

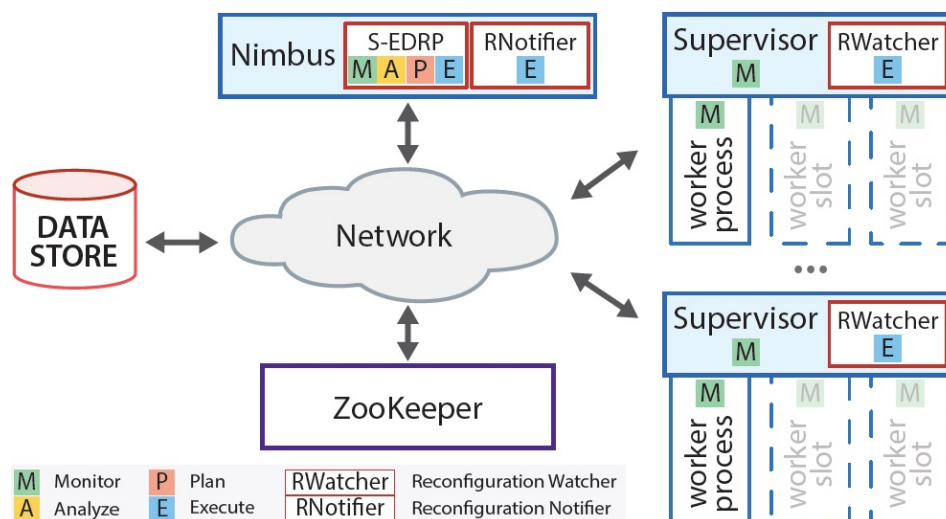
- How to self-adapt at run-time the deployment?
- DSP applications are:
  - long-running
  - subject to varying workloads
  - with computational requirements unknown a-priori



- Which main mechanisms do we need for run-time adaptation?
  - **Migration**: move operators from one node to another
  - **Elastic scaling**: change parallelism at application and/or infrastructure level

## EDRP: Elastic DSP in Storm

- Elastic DSP Replication and Placement (EDRP)
  - We augmented Distributed Storm with MAPE capabilities and optimal centralized placement and reconfiguration policy that keeps into account reconfiguration costs



# EDRP: still some limitations

- Centralized optimization algorithms do not scale for large problem instances
- Centralized MAPE architecture does not scale in geo-distributed environments
  - Distributed components but logic is still centralized
  - But fully distributed solutions have limitations
- Which solution for Edge-Cloud continuum?  
**Decentralize MAPE**

## How to decentralize control?

- Many patterns for decentralized control
  - Each one having pros and cons

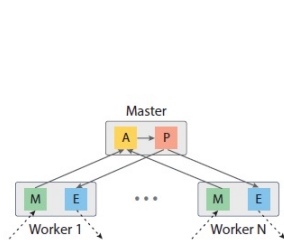


Figure 1: Hierarchical MAPE: master-worker pattern

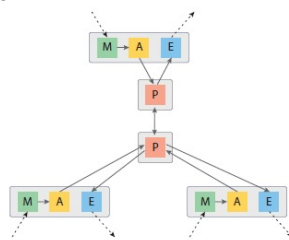


Figure 2: Hierarchical MAPE: regional pattern

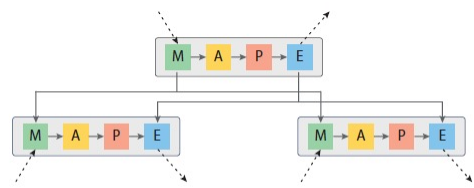


Figure 3: Hierarchical MAPE: hierarchical control pattern

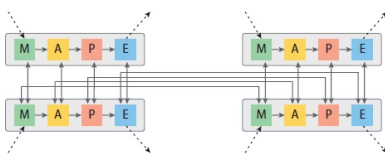


Figure 4: Flat MAPEs: coordinated control pattern

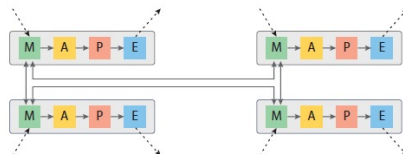
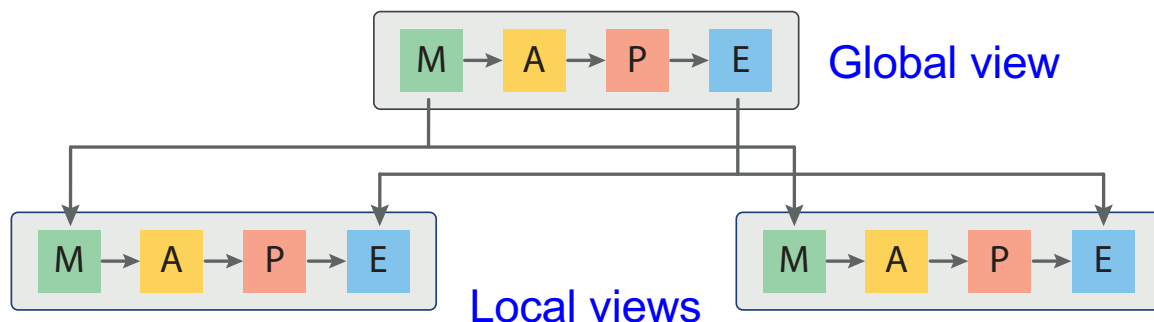


Figure 5: Flat MAPEs: information sharing pattern

D. Weyns et al., [On patterns for decentralized control in self-adaptive systems](#). In *Software Engineering for Self-Adaptive Systems II*, 2013

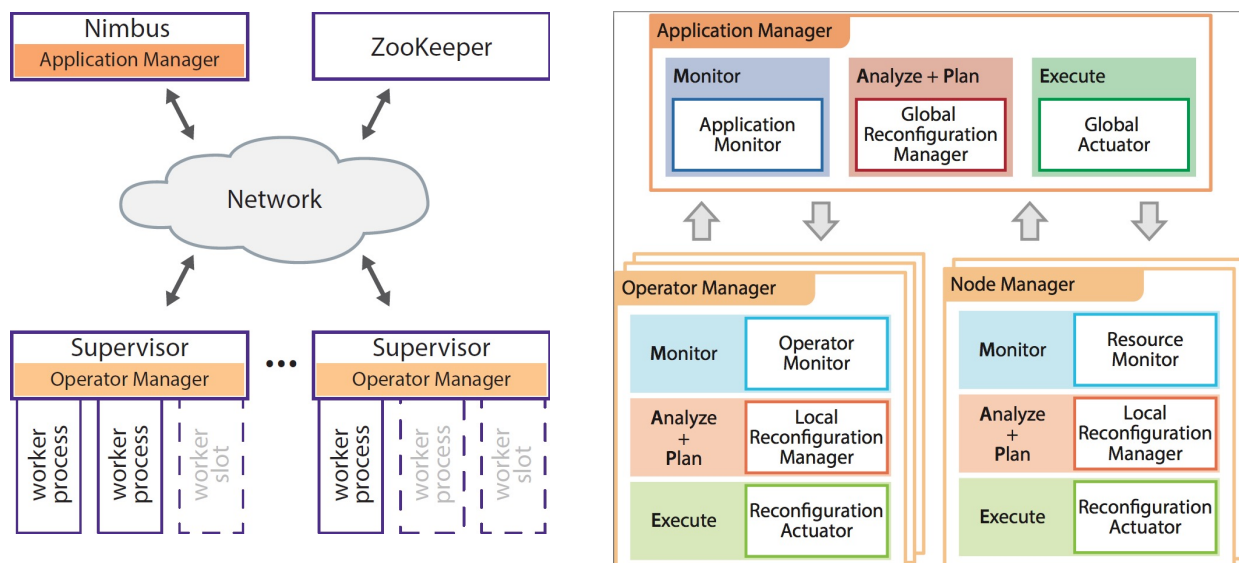
# How to decentralize control?

- Our approach:
  - **Hierarchical distributed** architecture to support run-time adaptation
  - Based on efficient distribution of **MAPE** control loops



## EDF: Elastic and Distributed DSP Framework

- Augmented Distributed Storm with MAPE capabilities and elasticity control



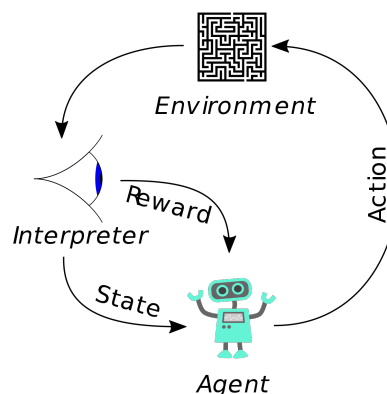
## EDF: Local elasticity policy

---

- Limited local view of the system (e.g., utilization level and input data rate of its operator)
- Two classes of elasticity policies
  - Classic threshold-based policy
    - Cons: empirical experience to choose thresholds
  - Based on **Reinforcement Learning**
    - Collection of techniques revolving around the basic idea of learning to make optimal decisions through interaction with controlled system
    - Goal: learn to select good actions *online*, based on paid costs (or gained reward)
    - *Pros*: *what* the user aims to obtain, instead of *how* it should be obtained

## Reinforcement Learning

---



- We considered:
  - Baseline **model-free** learning algorithm (Q-learning)
  - **Model-based** learning algorithm that exploits what is known or can be estimated about the system dynamics

Sutton and Barto, [Reinforcement Learning: An Introduction](#), 2020

## EDF: Local elasticity policy based on RL

- At each step RL agent performs an **action**, looking at current **state**
- Chosen action causes payment of **immediate cost** and transition to a new state
- To minimize expected **long-term (discounted) cost**, RL agent keeps estimates  $Q(s, a)$ 
  - Q-function: expected long-run cost that follows the execution of action  $a$  in state  $s$ :

---

**Algorithm 1** RL-based Operator Elastic Control Algorithm

---

```
1: Initialize the Q functions
2: loop
3:   choose a scaling action  $a_i$  (based on current estimates of Q)
4:   observe the next state  $s_{i+1}$  and the incurred cost  $c_i$ 
5:   update the  $Q(s_i, a_i)$  functions based on the experience
6: end loop
```

---

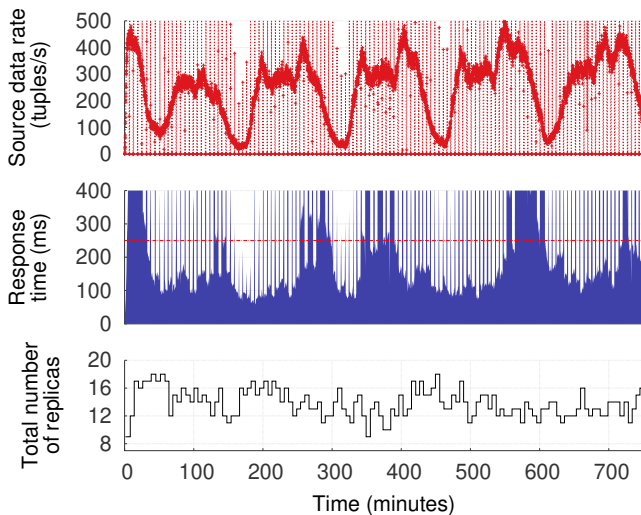
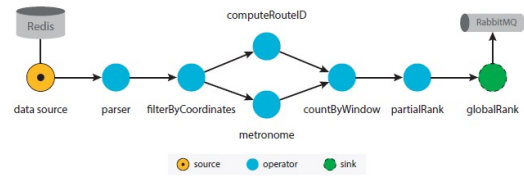
## EDF: Local elasticity policy based on RL

- **Q-learning**: classic model-free RL algorithm
- Q-learning: **choose next action**
  1. Either **exploits** agent knowledge about system, i.e., the current estimates  $Q$ , by greedily selecting the action that minimizes the estimated future costs
  2. Or **explores** by selecting a random action to improve its system knowledge
    - We consider  $\epsilon$ -greedy action selection method
- Q-learning: **update step**

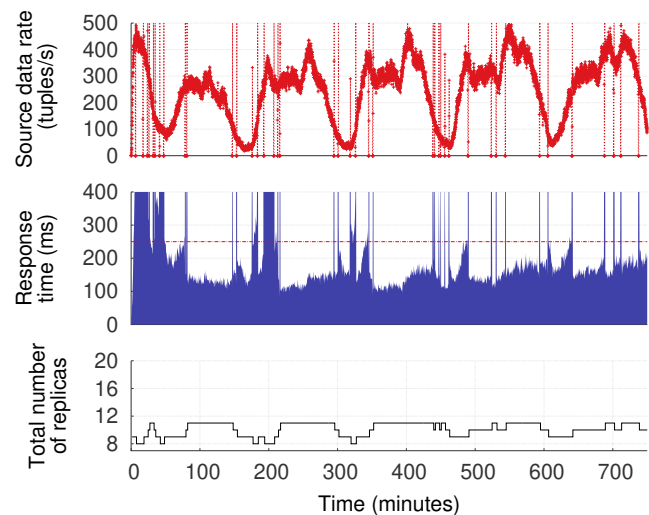
$$Q(s_i, a_i) \leftarrow (1 - \alpha)Q(s_i, a_i) + \alpha \left[ c_i + \gamma \min_{a' \in \mathcal{A}(s_{i+1})} Q(s_{i+1}, a') \right]$$

## EDF: Some results

- DSP application: DEBS'15 GC
- Q-learning vs. model-based RL



Avg. response time = 475 ms  
Downtime = 11%



Avg. response time = 176 ms  
Downtime = 3.2%

V. Cardellini - SABD 2021/22

34

## Exploit advanced RL techniques

- We have exploited more advanced RL techniques to tackle with heterogeneous resources
  - To deal with large state spaces, Function Approximation and Deep Learning techniques can be integrated into RL algorithms
  - Goal: to build approximate representations of state space and achieve near-optimal solutions with reduced memory demand
- See our tutorial at Performance 2021: [Reinforcement Learning for Run Time Performance Management in the Cloud/Edge](#)

## Other DSP deployment challenges

---

- How to manage DSP applications in Edge/Fog and Mobile Computing platforms?
- What about serverless DSP in the Edge-Cloud continuum?
- How to provide security guarantees?

### **Thesis opportunities**