



Universitat  
Pompeu Fabra  
*Barcelona*

# Robust Service Orchestration for Computing Continuum Systems

---

Boris Sedlak



Co-funded by  
the European Union



Funded by  
the European Union  
NextGenerationEU



Logo of the Plan de Recuperación, Transformación y Resiliencia, featuring a stylized "R" and "T".  
Plan de Recuperación,  
Transformación y Resiliencia

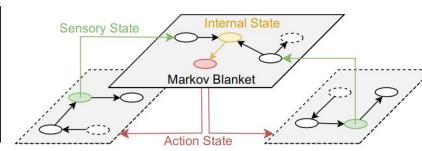
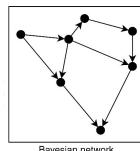
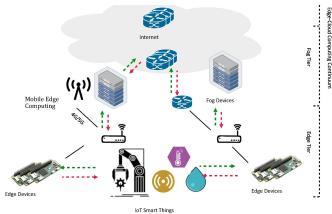


# Abstract

**Hidden in main presentation**

While the emergence of Edge computing promises major improvements in IoT processing, the heterogeneity and volatility of Edge infrastructures make service orchestration increasingly complex. Yet, to maintain robust system operation, we must be certain—or at least certain enough—about the expected outcome of our actions, such as shifting resources or workloads. This talk highlights Active Inference (AIF), an agent-based framework inspired by neuroscience, which supports operators in developing a causal understanding of the underlying generative processes. By continuously exploring and interacting with their environment, AIF agents develop models of the systems they govern and their interdependencies. This, in turn, enables them to predict the outcomes of actions when composing systems into larger architectures. The talk first outlines how AIF, as a concept, is inherently designed for acting under uncertainty, and then presents examples of how AIF can establish robust understanding of the environment to ensure continuous and adaptive system operation.

# Structure of the Talk



Motivation &  
Common Problems

Methodologies &  
Latest Results

Future Work



I – Who am I ??

# Boris Sedlak

**Postdoc @ UPF Barcelona**

Distributed Intelligence & Systems-Engineering Lab

**PhD @ TU Wien, Vienna**

Distributed Systems Group

**Software Engineer @ Lotterien**

Agile Development and Testing

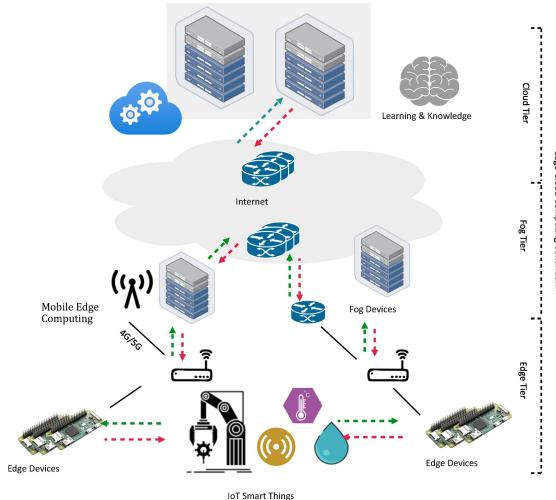


# I – Common Concepts

**Computing Continuum (CC)** as a composition of multiple processing tiers that stretch from IoT and edge computing, over Fog resources, to distant Cloud centers

Combines the benefits of all its tiers, i.e., low-latency and privacy-protecting computation from Edge, high availability and virtually unlimited processing resources from Cloud

Smart Cities are a common instance of distributed systems, where interconnected services (e.g., traffic surveillance or road surveillance) collaborate based on collected sensor data



Example of a Computing Continuum architecture [1]

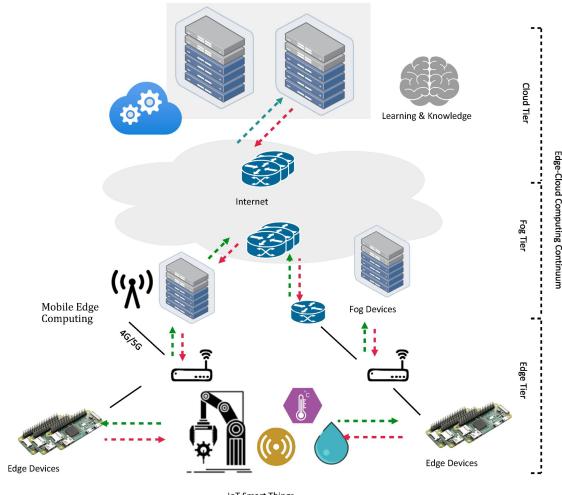
[1] P. Donta, I. Murturi, V. Casamayor, B. Sedlak, and S. Dustdar; **Exploring the Potential of Distributed Computing Continuum Systems** (2023)

# I – Common Concepts

**Computing Continuum (CC)** as a composition of multiple processing tiers that stretch from IoT and edge computing, over Fog resources, to distant Cloud centers

Combines the benefits of all its tiers, i.e., low-latency and privacy-protecting computation from Edge, high availability and virtually unlimited processing resources from Cloud

**Smart Cities** are a common instance of distributed systems, where interconnected services (e.g., traffic surveillance or road surveillance) collaborate based on collected sensor data



Example of a Computing Continuum architecture [1]

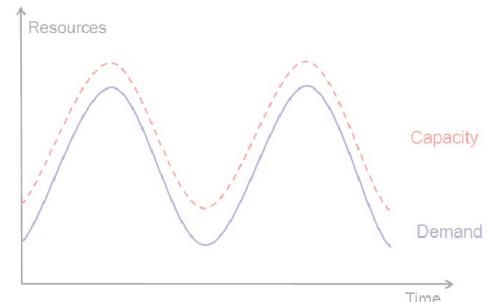
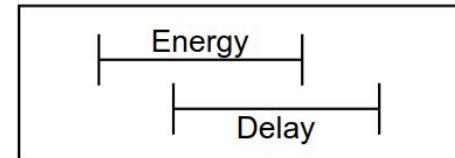
[1] P. Donta, I. Murturi, V. Casamayor, B. Sedlak, and S. Dustdar; **Exploring the Potential of Distributed Computing Continuum Systems** (2023)

# I – Common Concepts (cont.)

**Service Level Objectives (SLOs)** specify requirements that must be ensured throughout operation (e.g., latency  $< t$ ).

**Elasticity Strategies** scale a system according to current demand; e.g., if performance is insufficient, allocate more resources. However, what if this does not fulfill SLOs?

**Service Level Agreements (SLAs)** as binding agreement between service provider and consumer. However, very limited support in resource-restricted environments



Elasticity allocates the right amount of resources [2]

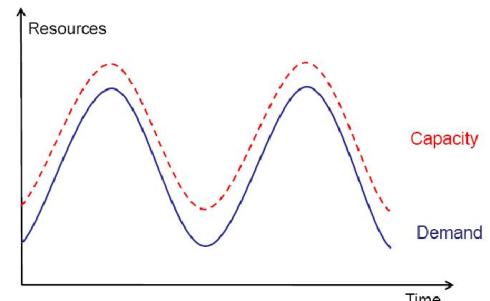
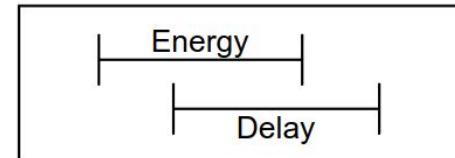
[2] Ricciardi et al., **Saving Energy in Data Center Infrastructures** (2011)

# I – Common Concepts (cont.)

**Service Level Objectives (SLOs)** specify requirements that must be ensured throughout operation (e.g., latency  $< t$ ).

**Elasticity Strategies** scale a system according to current demand; e.g., if performance is insufficient, allocate more resources. However, what if this does not fulfill SLOs?

**Service Level Agreements (SLAs)** as binding agreement between service provider and consumer. However, very limited support in resource-restricted environments

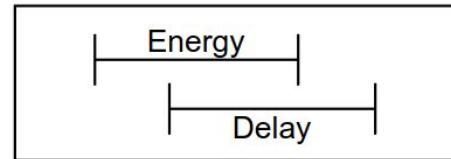


Elasticity allocates the right amount of resources [2]

[2] Ricciardi et al., Saving Energy in Data Center Infrastructures (2011)

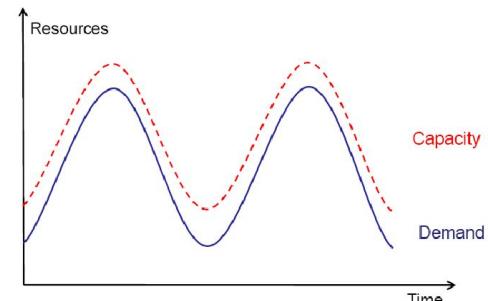
# I – Common Concepts (cont.)

**Service Level Objectives (SLOs)** specify requirements that must be ensured throughout operation (e.g., latency  $< t$ ).



**Elasticity Strategies** scale a system according to current demand; e.g., if performance is insufficient, allocate more resources. However, what if this does not fulfill SLOs?

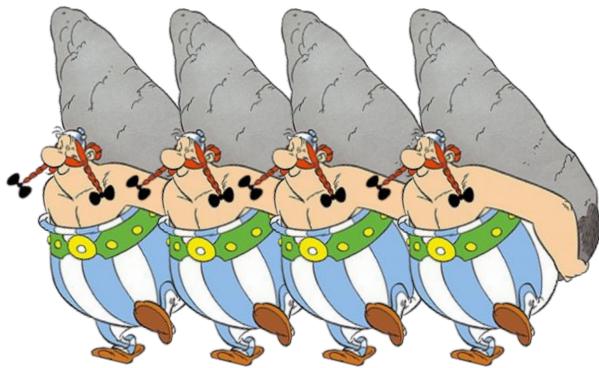
**Service Level Agreements (SLAs)** as binding agreement between service provider and consumer. However, very limited support in resource-restricted environments



Elasticity allocates the right amount of resources [2]

[2] Ricciardi et al., Saving Energy in Data Center Infrastructures (2011)

# I – Fundamental Problems



Homogeneous Resources (Cloud)

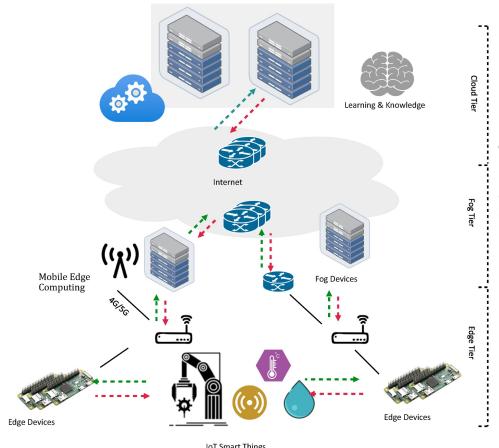


Heterogeneous Resources (CC)

# I – Fundamental Problems



Homogeneous Resources (Cloud)

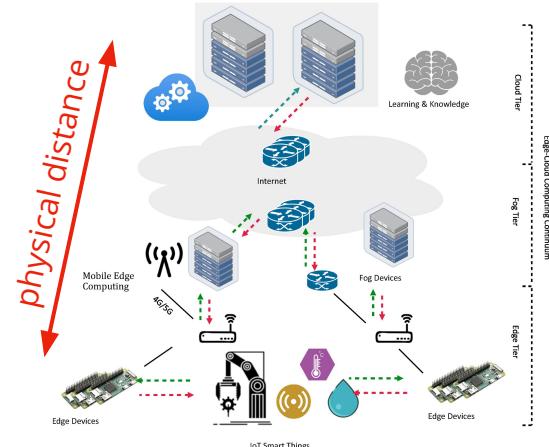


Heterogeneous Resources (CC)

# I – Fundamental Problems



Homogeneous Resources (Cloud)

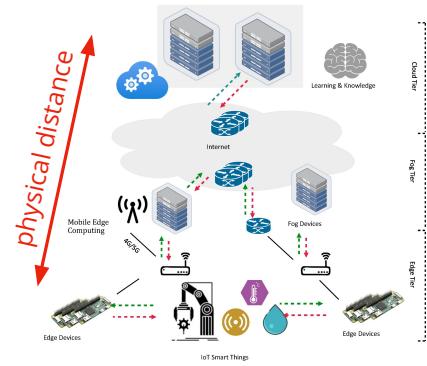


Heterogeneous Resources (CC)

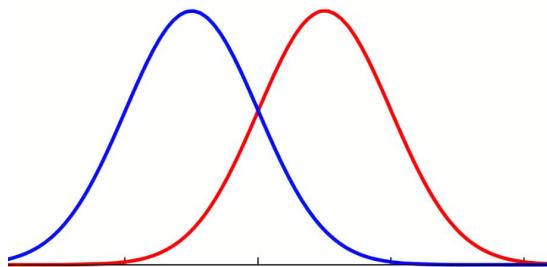
# I – Fundamental Problems

**Heterogeneity** prevents simple task distribution between available devices. Don't know *a priori* what will be the *behavior* of service X on device Y

**Physical distribution** creates latency between different distributed devices; transferring state information (e.g., CPU load) creates lots of traffic on the network.



# I – Fundamental Problems (cont.)



System dynamics change over time



Low trust into black-box models

# I – Fundamental Problems (cont.)

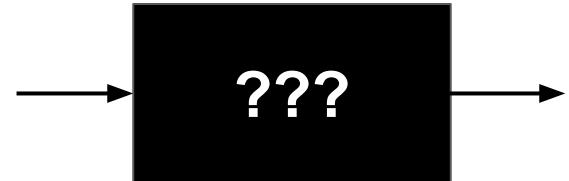
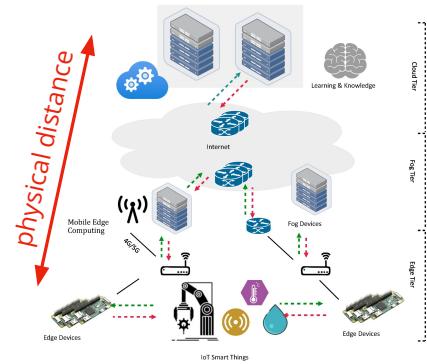
**Heterogeneity** prevents simple task distribution between available devices. Don't know a-priori what will be the *behavior* of service X on device Y

**Physical distribution** creates latency between different distributed devices; transferring state information (e.g., CPU load) creates lots of traffic on the network.

---

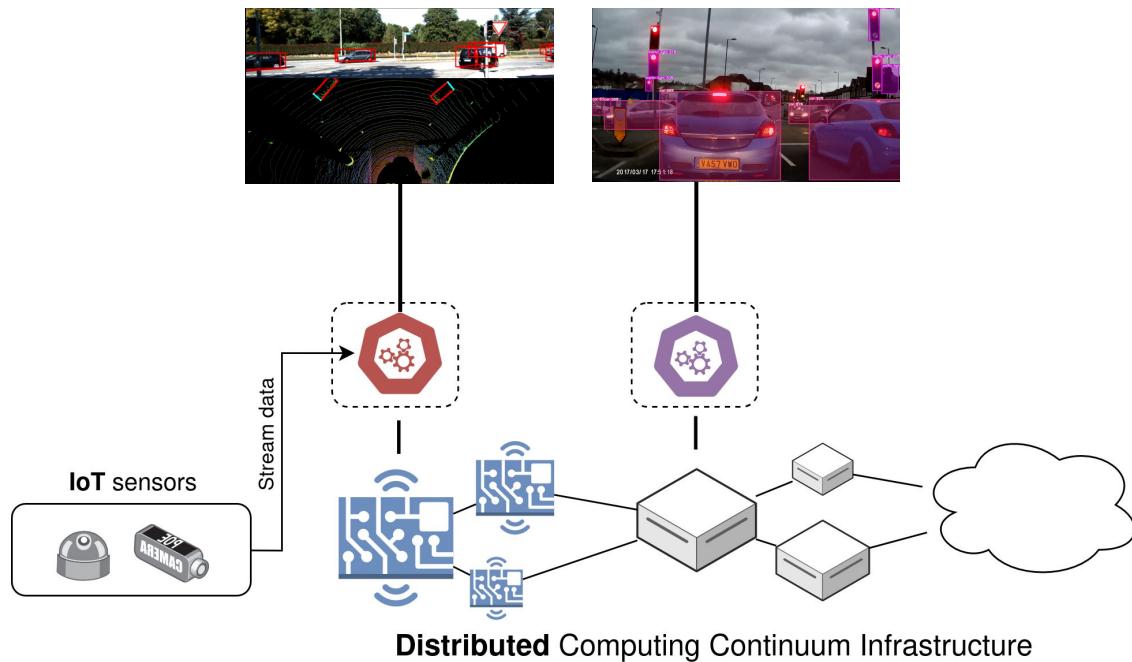
**Variable drifts** cause any ML model to lose accuracy over time. Cannot train models in one-shot operations but must incorporate changes continuously.

**Low-trust** into black-box ML models that cannot give human-verifiable chains of thought. E.g., debugging



## Distributed processing

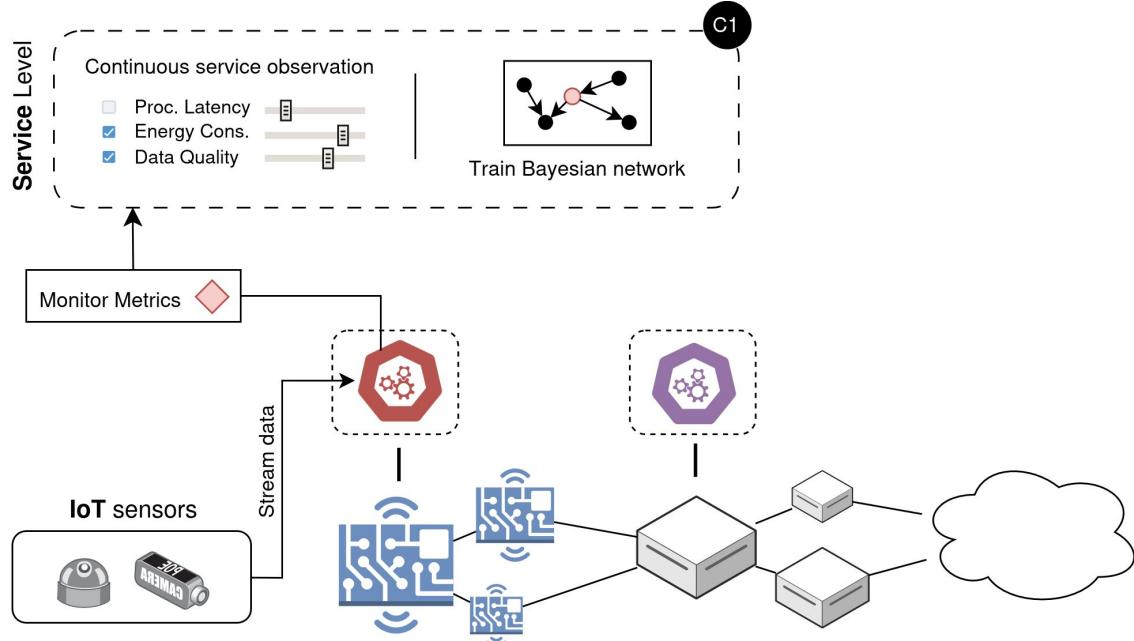
- Distributed CC infrastructure composed of various **device types** at different **locations**
- Sensor data continuously streamed from IoT devices to different processing services



# I – Overview of our Approach (C1) **Hidden in main presentation**

## Service interpretation

- Monitor processing across CC by collecting metrics; eval. real-time SLO fulfillment
- Train a model (i.e., BN) for predicting SLO fulfillment of **individual** services in different deployment configurations [7]
- Retrain model according to SLO prediction accuracy; slows down as model improves [9]



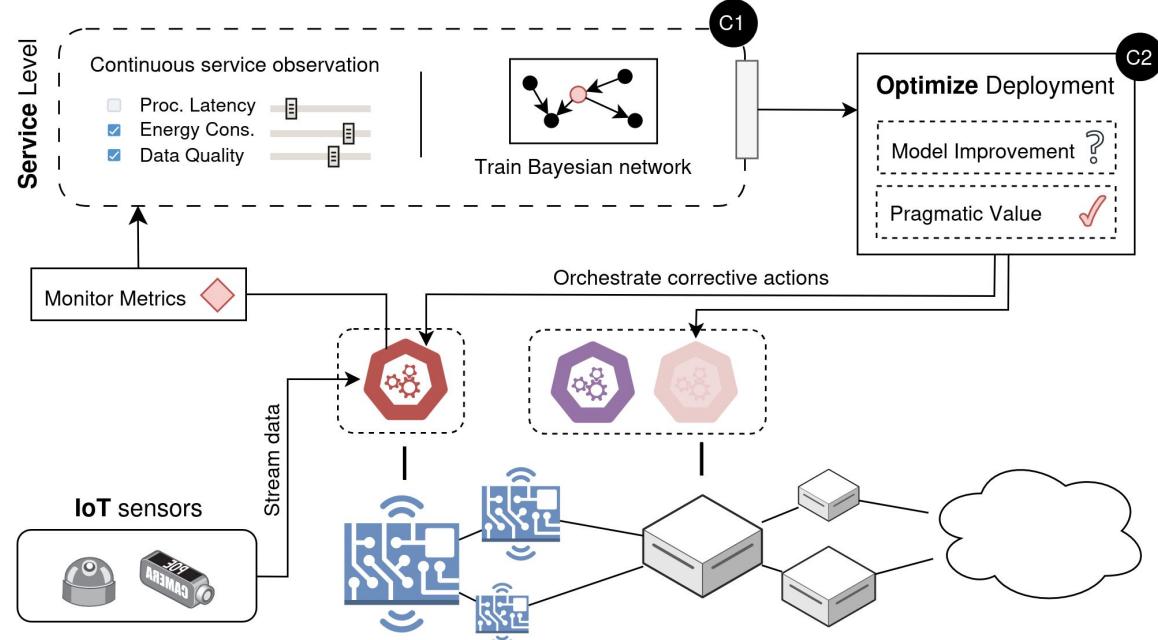
[7] Sedlak et al., Designing Reconfigurable Intelligent Systems with Markov Blankets, at ICSOC 2023

[9] Sedlak et al., SLO-Aware Task Offloading Within Collaborative Vehicle Platoons, at ICSOC 2024

# I – Overview of our Approach (C2) **Hidden in main presentation**

## Service adaptation

- Optimize **local SLO fulfillment** by reconfiguring processing services; choose between avail. elasticity strategies [5]
- Escape local optima through continuous exploration; check poss. **model improvement** [12]
- Decisions can be empirically verified and interpreted; useful for non-technical explanation



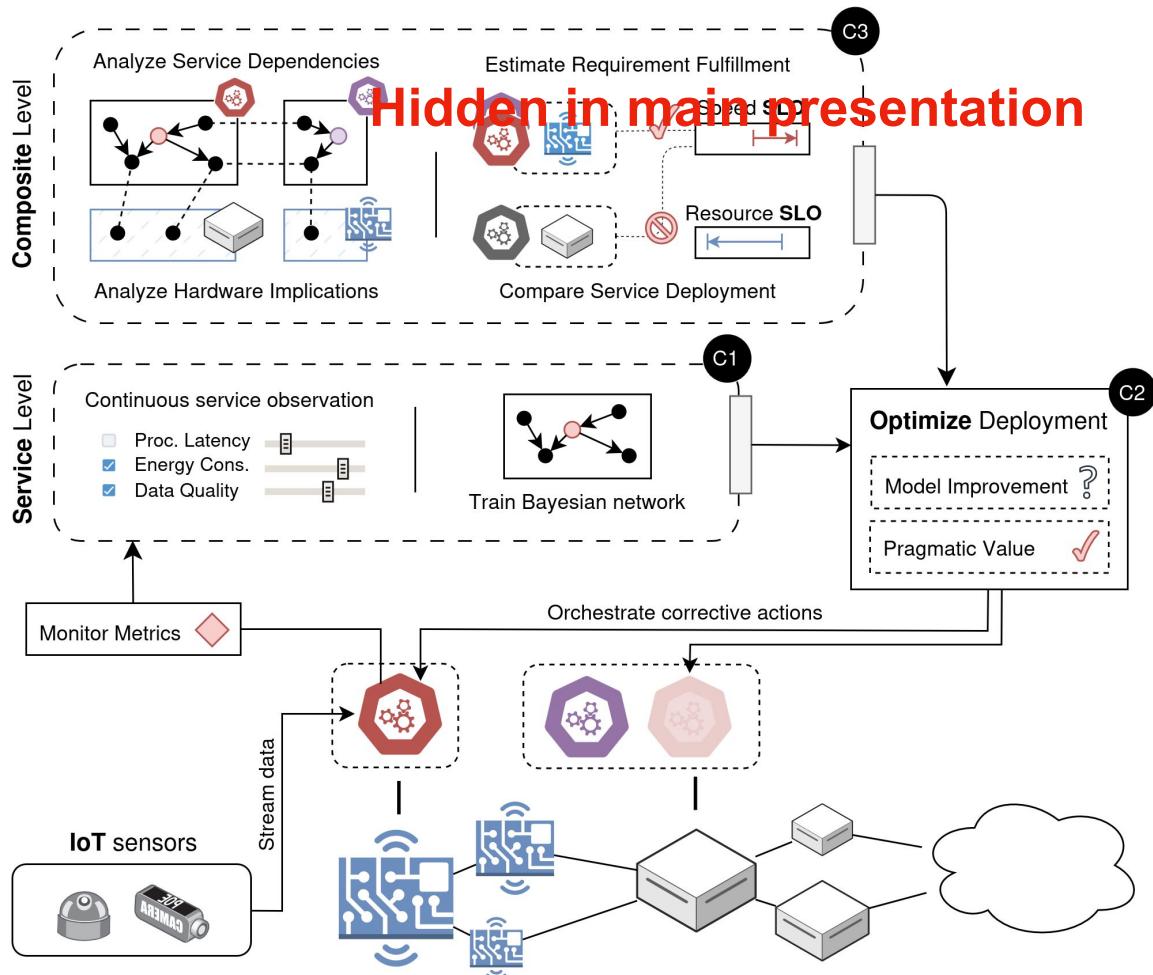
[5] Sedlak et al., From Metrics to Multidimensional Elasticity Strategies, at IEEE Services EDGE 2023

[12] Sedlak et al., Active Inference on the Edge: A Design Study, at PerconAI 2024

# I – Overview of our Approach (C3)

## Service collaboration

- Compose individual models to overarching representation; quantify service dependencies and **impact on hardware** [4]
- Collaborative orchestration considers service relations to optimize global SLO fulfillment
- Exchanging and merging BNs between services to speed up the **onboarding** of new service types or devices types [13]



[4] Sedlak et al., Markov Blanket Composition of SLOs, at IEEE Service EDGE 2024

[13] Sedlak et al., Equilibrium in the Computing Continuum through Active Inference, at Elsevier FGCS, 2024

# I – Fundamental Problems (cont.)

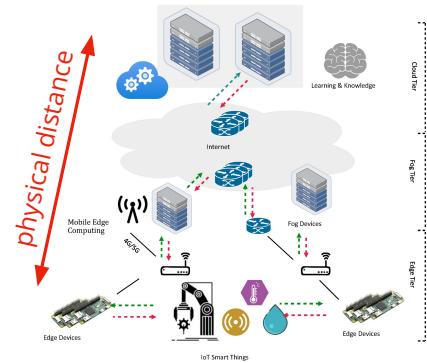
Heterogeneity prevents simple task distribution between available devices. Don't know *a priori* what will be

**Decentralized** decision-making based on local observations and understanding

Physical distribution  
CPU load) creates lots of traffic on the network.

**Variable drifts** cause any ML model to lose accuracy over time. Cannot train models in one-shot operations but must incorporate changes continuously.

**Low-trust** into black-box ML models that cannot give human-verifiable chains of thought. E.g., debugging



# I – Fundamental Problems (cont.)

Heterogeneity prevents simple task distribution between available devices. Don't know *a priori* what will be

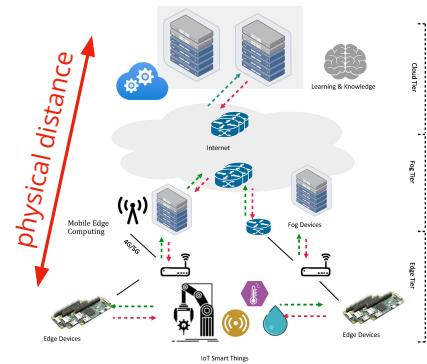
**Decentralized** decision-making based on local observations and understanding

Physical distribution  
CPU load) creates lots of traffic on the network.

Variability over time  
but more

**Continuous** and lifelong learning of environment, verifiable through structural causal models

Low-latency  
human-verifiable chains of thought. E.g., debugging



## II – Markov Blanket (MB)

Interactions between **systems** (e.g., human in world) can be expressed through MBs; fulfill Markov property – so decisions can be taken based on system's current state

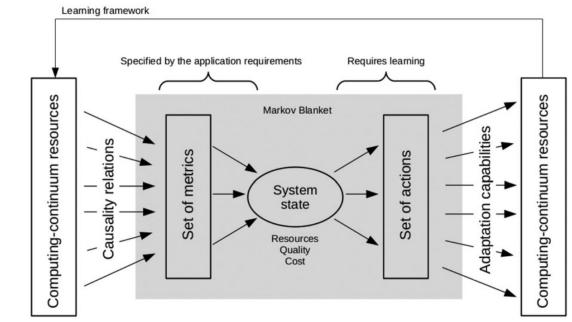
Creates **formal boundary** between a system and external states; within MB, discard all variables that show no impact on internal states and requirements (i.e., SLO fulfillment)

---

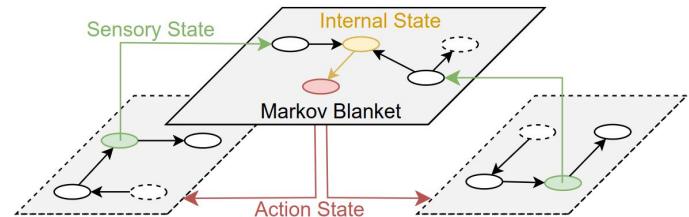
Provides clear interfaces for **sensory** and **action states**; policy (e.g., scaling) as a mapping between these states

[3] Dustdar et al., **On Distributed Computing Continuum Systems** (2023)

[4] Sedlak et al., **Markov Blanket Composition of SLOs**, at IEEE Services Edge 2024



Behavioral Markov blanket of a system [3]

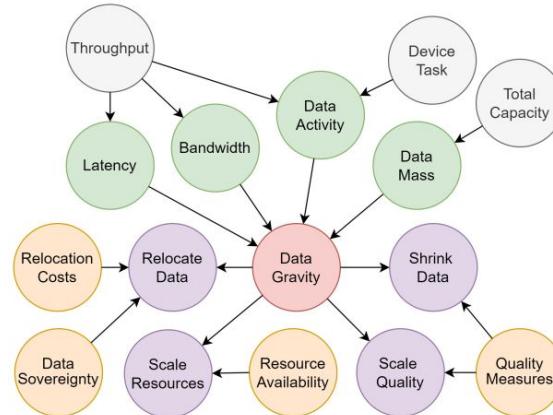


Action-perception cycle between multiple entities [4]

MB: Expresses how to evaluate a composite SLO and how to react according to the current device context

## Behavioral model

Internal state (●) evaluates objectives and how these relate to external sensory inputs (●); can interact with the world through action, i.e., elasticity strategies (●), which are influenced by contextual factors (●)



Example of a behavioral model for data gravity [5]

[5] Sedlak et al., **Controlling Data Gravity and Data Friction: From Metrics to Multidimensional Elasticity Strategies**, at IEEE Service EDGE 2023

## II – Stream Processing Scenarios

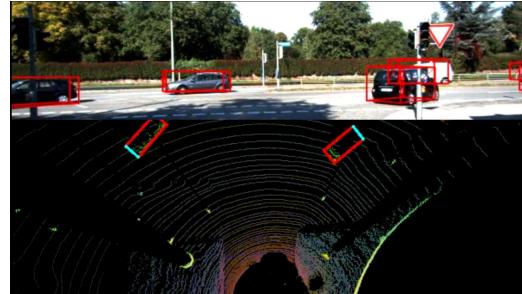
Commonly addressed use cases revolve around continuous **stream processing**; in case **time-critical** adaptations are required, this poses a higher need for sophisticated adaptation mechanisms.

Video Processing (Yolo V8)



Object detection in a video stream using Yolo [6]

Mobile Mapping (Lidar)



Creating a mobile map from binaries using Lidar [6]

QR Scanner (OpenCV)

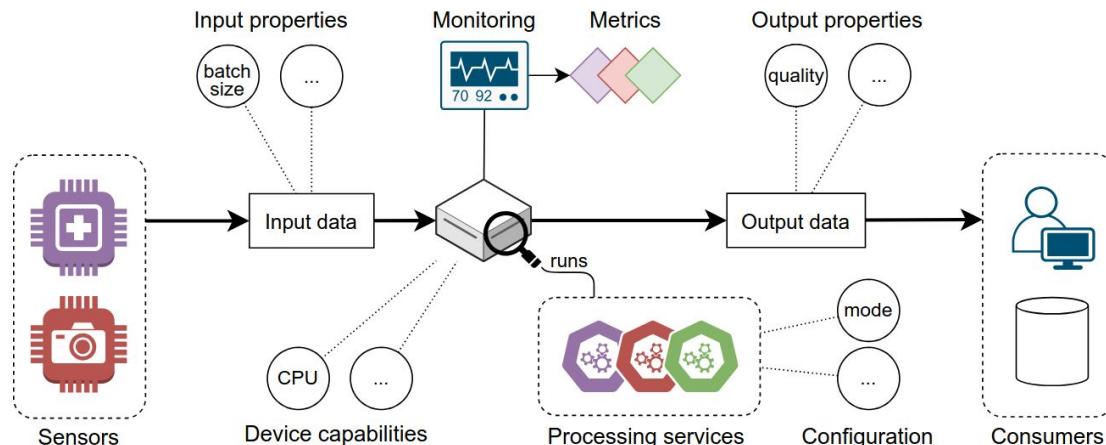


QR code scanning in a video using OpenCV [6]

[6] Sedlak et al., **Adaptive Stream Processing on Edge Devices through Active Inference** (Scheduled for 2025 at Springer ES)

## II – Stream Processing Scenarios

Commonly addressed use cases revolve around continuous **stream processing**; in case **time-critical** adaptations are required, this poses a higher need for sophisticated adaptation mechanisms.

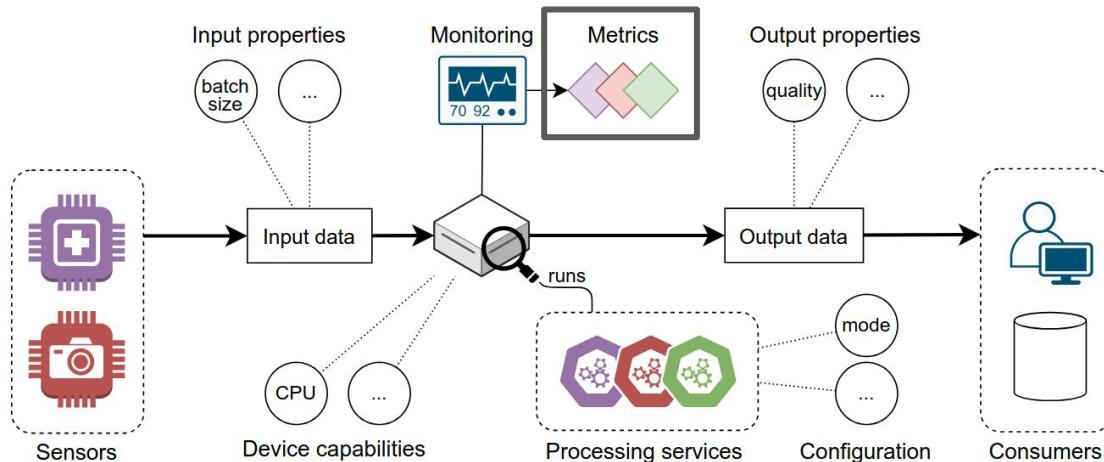


Abstract representation of a monitored stream processing service [6]

[6] Sedlak et al., **Adaptive Stream Processing on Edge Devices through Active Inference** (Scheduled for 2025 at Springer ES)

## II – Stream Processing Scenarios

Commonly addressed use cases revolve around continuous **stream processing**; in case **time-critical** adaptations are required, this poses a higher need for sophisticated adaptation mechanisms.



Abstract representation of a monitored stream processing service [6]

fps	pixel	cores	change_flag	timestamp
31	800	3	False	2024-11-30
32	800	3	False	2024-11-30
32	800	3	False	2024-11-30
32	800	3	False	2024-11-30
32	800	3	False	2024-11-30
32	800	3	False	2024-11-30
32	800	3	False	2024-11-30
32	800	3	False	2024-11-30
32	800	3	False	2024-11-30

Example of processing metrics as tabular data [6]

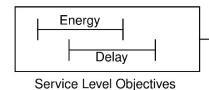
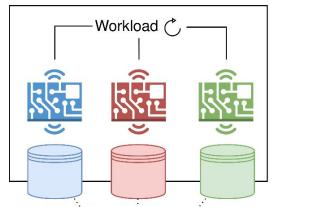
[6] Sedlak et al., *Adaptive Stream Processing on Edge Devices through Active Inference* (Scheduled for 2025 at Springer ES)

## II – Basic Methodology

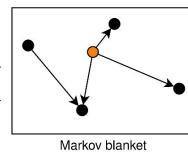
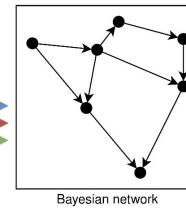
**Optimize** the system (i.e., its SLO fulfillment) according to model

**Model:** use processing metrics and properties to train a *generative model* that describes the variable relations within the environment.

### 1. Bayesian Network Learning



### 2. Markov Blanket Selection



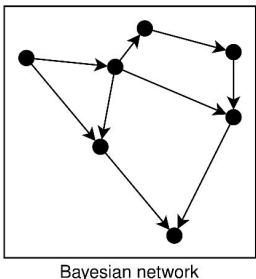
### 3. Explainable Inference

[7]

[7] Sedlak et al., Designing Reconfigurable Intelligent Systems with Markov Blankets, at ICSOC 2023

## II – Basic Methodology

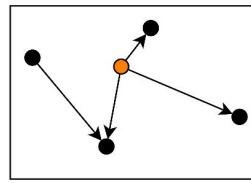
### Bayesian Network Learning



Bayesian network

- Structure Learning**  
Various algorithms (e.g., HCS)  
Directed Acyclic Graph (DAG)
- Parameter Learning**  
Various algorithms (e.g., MLE)  
Conditional Prob. Table (CPT)

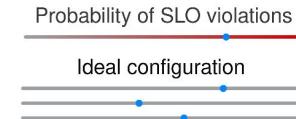
### Markov Blanket Selection



Markov blanket

- Causality filter**  
Extract subset of variables  
that impact SLO fulfillment
- Behavioral MB**  
MB now contains contextual  
factors & elasticity strategies

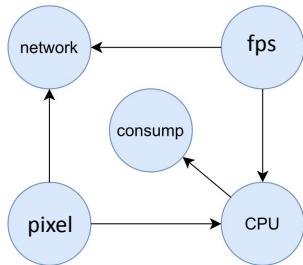
### Knowledge Extraction



- Conditional Inference**  
Estimate impact of different  
deployment configuration
- Optimize SLOs**  
Adjust processing services  
according to inferred policy

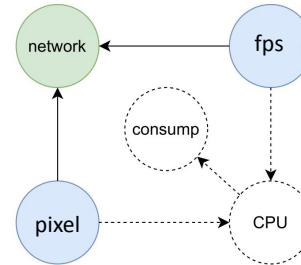
## II – Basic Methodology

### Bayesian Network Learning



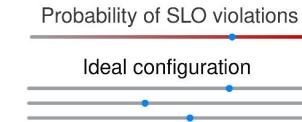
- Structure Learning**  
Various algorithms (e.g., HCS)  
Directed Acyclic Graph (DAG)
- Parameter Learning**  
Various algorithms (e.g., MLE)  
Conditional Prob. Table (CPT)

### Markov Blanket Selection



- Causality filter**  
Extract subset of variables  
that impact SLO fulfillment
- Behavioral MB**  
MB now contains contextual  
factors & elasticity strategies

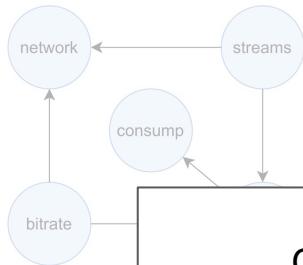
### Knowledge Extraction



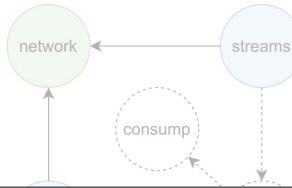
- Conditional Inference**  
Estimate impact of different  
deployment configuration
- Optimize SLOs**  
Adjust processing services  
according to inferred policy

## II – Basic Methodology

### Bayesian Network Learning



### Markov Blanket Selection



### Knowledge Extraction



Core methodology applied in multiple application.

#### Structure Learning

Various algorithms (e.g., HCS)  
Directed Acyclic Graph (DAG)

#### Causality infer

Extract subset of variables  
that impact SLO fulfillment

#### Conditional Inference

Estimate impact of different  
deployment configuration

#### Parameter Learning

Various algorithms (e.g., MLE)  
Conditional Prob. Table (CPT)

#### Behavioral MB

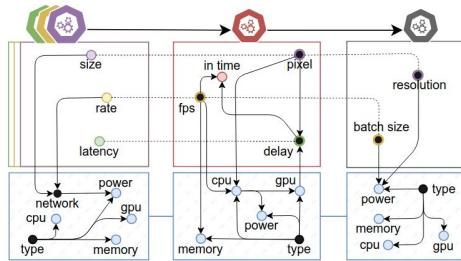
MB now contains contextual  
factors & elasticity strategies

#### Optimize SLOs

Adjust processing services  
according to inferred policy

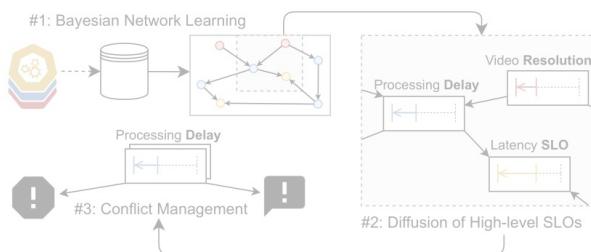
## II – Multiple Applications

### Transitive Requirements [4]



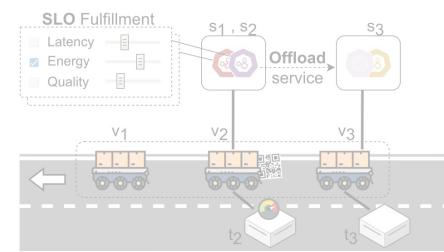
Optimize SLO fulfillment by deploying microservices over **heterogeneous hardware**; use BNs to analyze & optimize service/device dependencies

### Diffusing High-Level SLOs [8]



Find configurations for subsystems according to **high-level SLO** values; create hierarchy of dependencies and infer lower-level configurations

### SLO-Aware Offloading [9]



Offload microservices over **heterogeneous hardware**; estimate effects of service swapping on SLO fulfillment

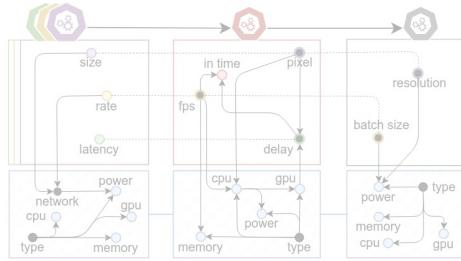
[4] Sedlak et al., **Markov Blanket Composition of SLOs**, at IEEE Services EDGE 2024

[8] Sedlak et al., **Diffusing High-level SLO in Microservice Pipelines**, at IEEE SOSE 2024

[9] Sedlak et al., **SLO-Aware Task Offloading Within Collaborative Vehicle Platoons**, at ICSOC 2024

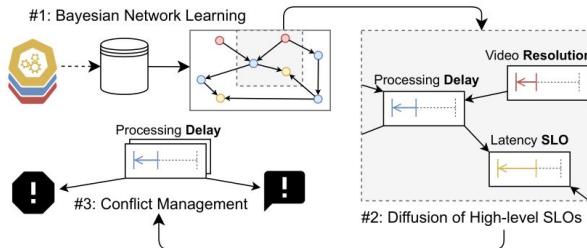
## II – Multiple Applications

### Transitive Requirements [4]



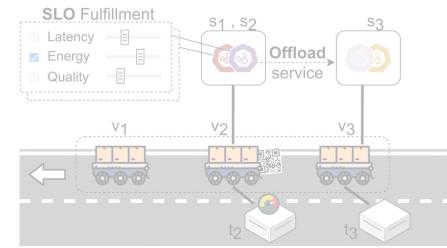
Optimize SLO fulfillment by deploying microservices over heterogeneous hardware;  
use BNs to analyze & optimize service/device dependencies

### Diffusing High-Level SLOs [8]



Find configurations for subsystems according to **high-level SLO** values;  
create hierarchy of dependencies and infer lower-level configurations

### SLO-Aware Offloading [9]



Offload microservices over heterogeneous hardware;  
estimate effects of service swapping on SLO fulfillment

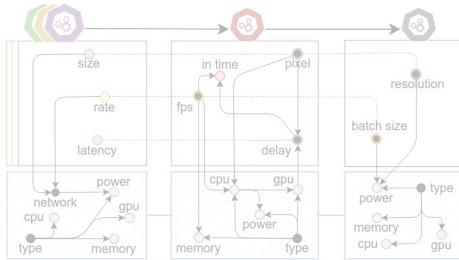
[4] Sedlak et al., **Markov Blanket Composition of SLOs**, at IEEE Services EDGE 2024

[8] Sedlak et al., **Diffusing High-level SLO in Microservice Pipelines**, at IEEE SOSE 2024

[9] Sedlak et al., **SLO-Aware Task Offloading Within Collaborative Vehicle Platoons**, at ICSOC 2024

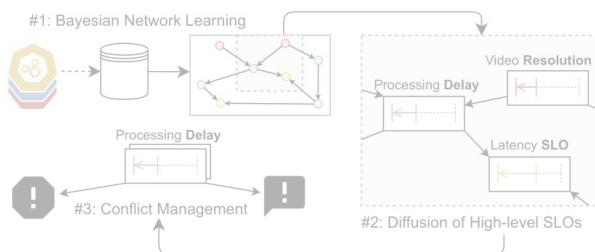
## II – Multiple Applications

### Transitive Requirements [4]



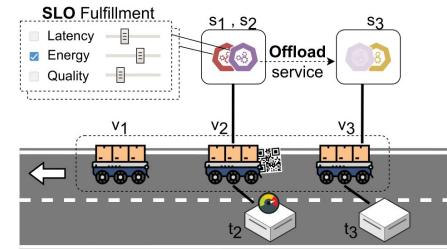
Optimize SLO fulfillment by deploying microservices over **heterogeneous hardware**;  
use BNs to analyze & optimize service/device dependencies

### Diffusing High-Level SLOs [8]



Find configurations for subsystems according to **high-level SLO** values;  
create hierarchy of dependencies and infer lower-level configurations

### SLO-Aware Offloading [9]



Offload microservices over **heterogeneous hardware**;  
estimate effects of service swapping on SLO fulfillment

[4] Sedlak et al., **Markov Blanket Composition of SLOs**, at IEEE Services EDGE 2024

[8] Sedlak et al., **Diffusing High-level SLO in Microservice Pipelines**, at IEEE SOSE 2024

[9] Sedlak et al., **SLO-Aware Task Offloading Within Collaborative Vehicle Platoons**, at ICSOC 2024

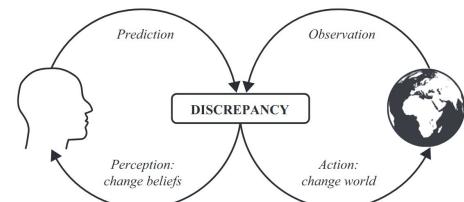
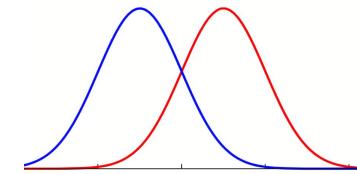
## II – Refining our Approach

### Known Shortcomings

- (1) BNL requires large amounts of training data upfront
- (2) can't visit all possible states, so where to start exploring?
- (3) over time, models get distorted due to variable drifts

### Active Inference

Concept from **neuroscience** developed by Friston et al. [10,11]; allows agents to interact with their environment by learning the underlying **generative models** to persist over time



Action-perception cycle in Active Inference [11]

[10] Parr et al., **Active Inference: The Free Energy Principle in Mind, Brain, and Behavior** (2022)

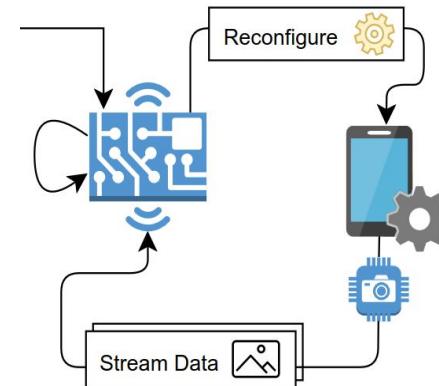
[11] Friston et al., **Designing ecosystems of intelligence from first principles** (2024)

## II – Active Inference in CC Systems

Mapping between neuroscience and distributed computing systems [6,12,13]; understanding processing requirements (i.e., SLOs) as a form of **homeostasis**, e.g., cell temperature

Create autonomous components that identify how to ensure requirements and resolve them independently, clear modelling between higher-level and low-level components

Simplify service orchestration in large-scale distributed systems; decentralized decision-making of individual components **avoids** transferring service states to the Cloud



Ensure internal requirements [13]

[6] Sedlak et al., **Adaptive Stream Processing on Edge Devices through Active Inference** (Scheduled for 2025 at Springer ES)

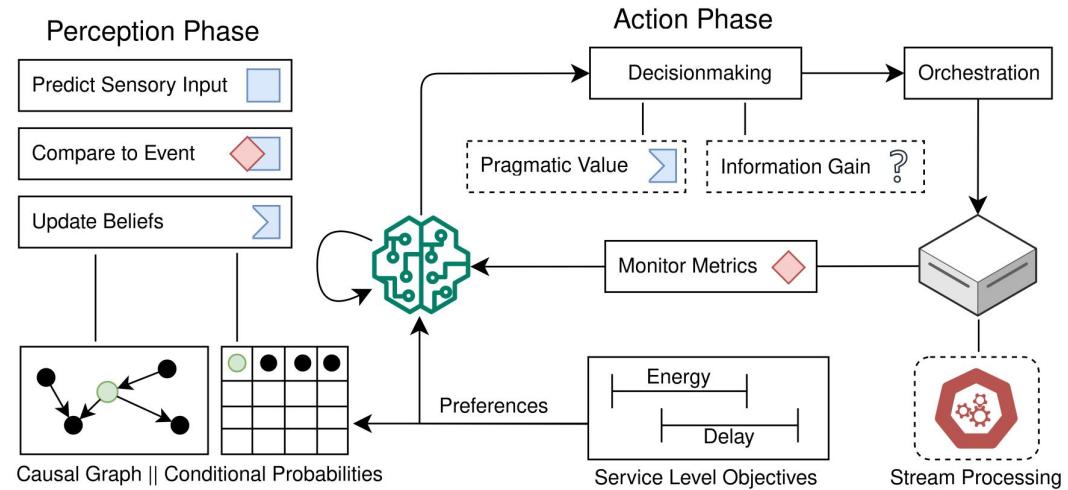
[12] Sedlak et al., **Active Inference on the Edge: A Design Study**, at PerconAI 2024

[13] Sedlak et al., **Equilibrium in the Computing Continuum through Active Inference**, Elsevier FGCS (2024)

## II – Active Inference Architecture

### Approach

- (1) **Specify** ideal runtime behavior through SLOs
- (2) **AIF agents** monitor their environment & collect metrics
- (3) **Perception phase** predicts expected SLO fulfillment and adjusts the generative model
- (4) **Action phase** orchestrates the processing environment to optimize both SLOs and model



Action and perception cycles performed by the AIF agent to create an accurate model and shape the world [6]

[6] Sedlak et al., **Adaptive Stream Processing on Edge Devices through Active Inference** (Scheduled for 2025 at Springer ES)

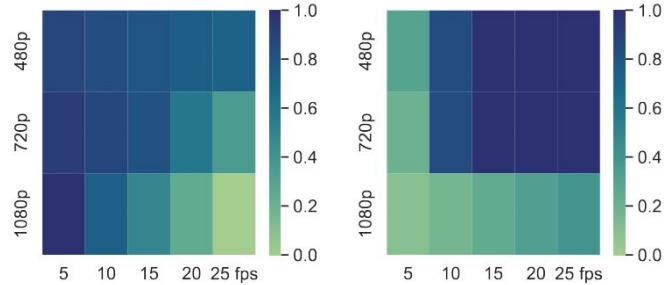
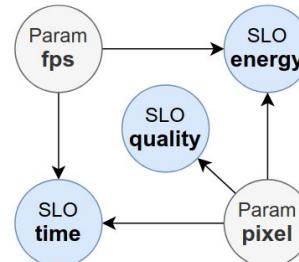
## II – Active Inference Architecture (cont.)

### Interpretable behavior [6]

- Empirically verify variable relations in the BNs, e.g., increasing quality (pixel) leads to high energy usage; adjust **parameters** (i.e., pixel & fps) according to SLOs
- Quantified preferences of the agent: (1) expected SLO fulfillment or (2) potential model improvement; determine the behavior of the scaling agent

### Continuous composition [4]

- Gradually create increasingly accurate models for individual processing services; continuously compose to estimate the **impact** they have on each other



[4] Sedlak et al., **Markov Blanket Composition of SLOs**, at IEEE Services EDGE 2024

[6] Sedlak et al., **Adaptive Stream Processing on Edge Devices through Active Inference** (Scheduled for 2025 at Springer ES)

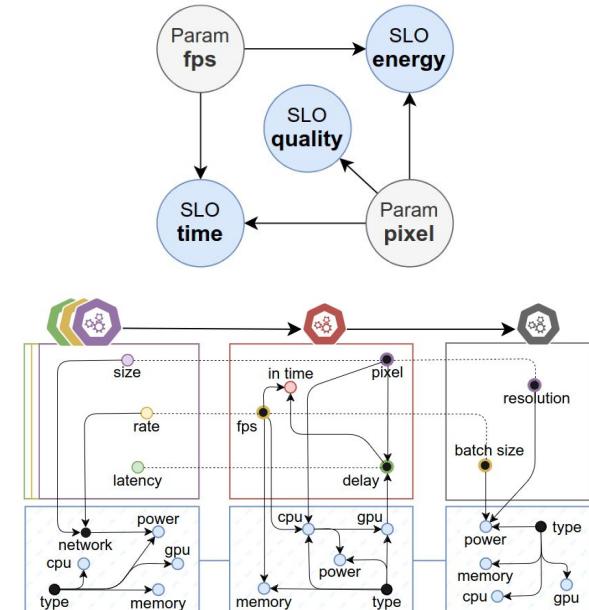
## II – Active Inference Architecture (cont.)

### Interpretable behavior [6]

- Empirically verify variable relations in the BNs, e.g., increasing quality (pixel) leads to high energy usage; adjust **parameters** (i.e., pixel & fps) according to SLOs
- Quantified preferences of the agent: (1) expected SLO fulfillment or (2) potential model improvement; determine the behavior of the scaling agent

### Continuous composition [4]

- Gradually create increasingly accurate models for individual processing services; continuously compose to estimate the **impact** they have on each other



[4] Sedlak et al., **Markov Blanket Composition of SLOs**, at IEEE Services EDGE 2024

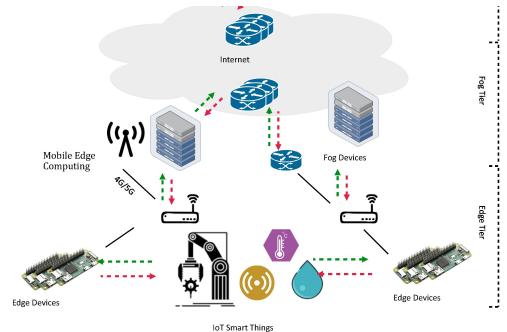
[6] Sedlak et al., **Adaptive Stream Processing on Edge Devices through Active Inference** (Scheduled for 2025 at Springer ES)

## III – Summary Contemporary Challenges

**IoT & Edge** enable large-scale distributed services that optimize our daily routines; CC as underlying infrastructure for supporting these services

**Resource** limitations and device heterogeneity complicate service orchestration; device & service behavior not guaranteed, leads to violated SLOs

**Missing** explainability for black box ML models; leads to low trust and non-interpretable behavior



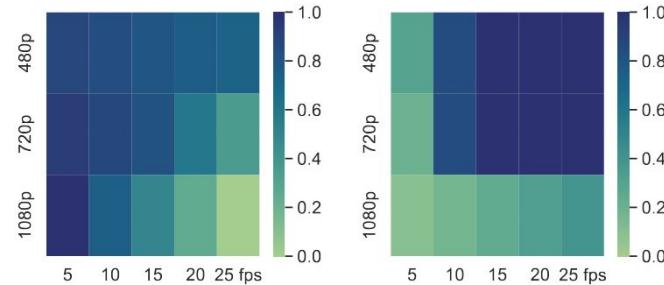
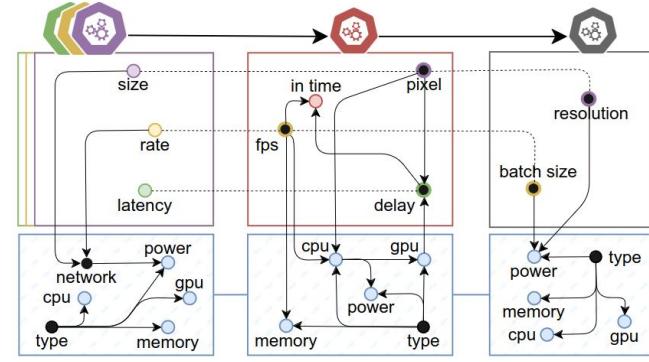
## III – Summary

### Contributions & Results

**Monitor** processing services closely, analyze their behavior, and infer optimal elasticity strategy; train a MB that combines all these factors in one model

**Dynamically** optimize SLOs fulfillment for a network of processing services and heterogeneous devices; orchestrate services (i.e., placement, configuration, replication, etc) according to current context

**Active Inference** as a natural fit to train behavioral MBs and keep them accurate; balance continuously between ensuring SLOs and improving the model



## III – Summary

### Limitations & Challenges (1/2)

**High cost** of running empirical experiments; slow training progress, needs training environment.

→ Train agents in simulated environments; creating accurate environment, like for **digital twinning**

Overhead of extracting baselines from publication; low rigor from generic algorithms, e.g., SB3

→ Support baseline comparison and standardized problem instances in evaluation environment

Governance limited to individual vendor; SLAs only supported by one provider and its subset of nodes  
 → Create overarching models for orchestration and client compensation; technical & economical part

#### Active Inference for Digital Twins: Predicting and Optimising IoT Processing Service Performance

Elena Pretel  
 LoIIS Research Group,  
 DIA, University of Castilla-La Mancha  
 Albacete, Spain  
 emanuela.pretel@uclm.es

Boris Sedlak  
 Vienna University of Technology  
 Vienna, Austria  
 b.sedlak@stg.tuw.ac.at

Victor Cammarer-Pujol  
 Universitat Pompeu Fabra  
 Barcelona, Spain  
 victor.cammarer@upf.edu

Pascual Gonzalez  
 LoIIS Research Group,  
 DIA, University of Castilla-La Mancha  
 Albacete, Spain  
 pascual.gonzalez@uclm.es

Elena Navarro  
 LoIIS Research Group,  
 DIA, University of Castilla-La Mancha  
 Albacete, Spain  
 elena.navarro@uclm.es

Sahabram Dastidar  
 Vienna University of Technology  
 and Universitat Pompeu Fabra  
 Wien and Barcelona, Austria and Spain  
 dastidar@dtv.tuw.ac.at

**Abstract**  
 Digital Twins (DTs) are emerging as enablers for real-time monitoring and control in IoT systems. In this work, we propose a DT enabled with Active Inference (AI), a framework rooted in the Free Energy Principle (FEP) that provides a way to predict and self-adjust decision-making capabilities. Our approach is evaluated on a realistic edge computing scenario where two co-located teleprocessing services compete for resources. The proposed AI-based DT outperforms a competitor for limited CPU resources. The DT continuously infers physical entities' states and optimizes their QoS to meet their Service Level Objectives (SLO). In a series of experiments, we validate the predictive accuracy and control effectiveness of the AI-enabled DT. The DT achieves an average SLO fulfillment above 90% for both services, and predicted throughput improvements are consistent with real observations confirmed through statistical testing.

**ACM Reference Format:** Elena Pretel, Victor Cammarer-Pujol, Boris Sedlak, Navarro, Victor, Elena, Pretel, Victor, Cammarer-Pujol, Boris, Sedlak, Navarro, Victor, Elena, Navarro, Sahabram Dastidar. 2023. Active Inference for Digital Twins: Predicting and Optimising IoT Processing Service Performance. In *Proceedings of the 2023 ACM SIGMOD/PODS Joint Conference on Data and Cloud Computing (SIGMOD/PODS '23)*, November 20–24, 2023, Vienna, Austria. ACM, New York, NY, USA, Article 10 (2023), 10 pages. <https://doi.org/10.1145/3558808.3589776>

**1 Introduction**  
 The fast advancement of technology over the last few decades has transformed the way systems are monitored, controlled, and optimized. The Internet of Things (IoT) is one of the main drivers of this evolution, enabling the creation of smart cities and healthcare. The rise of mobile devices, Internet of Things (IoT) sensors, and cloud computing have paved the way for the emergence of Digital Twins (DTs).

In large-scale and distributed IoT environments, such as smart manufacturing, smart cities, and smart healthcare, it is challenging to manage the complexity of the system being monitored and controlled. DTs provide a way to monitor and manage these IoT systems by creating virtual replicas of physical entities operating under real-time data flows [1].

The DT concept was originally coined by Gómez and Vivero in 2000 [2] and is defined as “a system that consists of (1) a physical space, (2) a virtual space, and (3) bidirectional connection between both spaces allowing the management” [4]. Since then, the concept has been refined and extended. DTs are now considered as virtual representations of physical entities that allow real-time monitoring, control, and optimization of the system to improve its behavior [5]. One of the main properties of a DT is the predictability [6]. This property is crucial for a DT to be able to forecast the future behaviour of its physical twin [5]. This property is crucial to enable intelligent decision making within the DT, in particular when leveraging AI based techniques.

Upcoming paper at ACM IoT [14]

## III – Summary

### Limitations & Challenges (1/2)

**High cost** of running empirical experiments; slow training progress, needs training environment.

→ Train agents in simulated environments; creating accurate environment, like for **digital twinning**

**Overhead** of extracting baselines from publication; low rigor from generic algorithms, e.g., SB3

→ Support baseline comparison and standardized problem instances in evaluation environment

Governance limited to individual vendor; SLAs only supported by one provider and its subset of nodes  
 → Create overarching models for orchestration and client compensation; technical & economical part

#### Active Inference for Digital Twins: Predicting and Optimising IoT Processing Service Performance

Elena Pretel  
LoIIS Research Group,  
DIA, University of Castilla-La Mancha  
Albacete, Spain  
mariamaria.pretel@uclm.es

Boris Sedlak  
Vienna University of Technology  
Vienna, Austria  
b.sedlak@stg.tuw.ac.at

Victor Cammarero-Pujol  
Universitat Pompeu Fabra  
Barcelona, Spain  
victor.cammarero@upf.edu

Elena Navarro  
LoIIS Research Group,  
DIA, University of Castilla-La Mancha  
Albacete, Spain  
e.lnavarro@uclm.es

Sahabram Dastidar  
Vienna University of Technology  
and Universitat Pompeu Fabra  
Vienna, Austria  
s.dastidar@tuw.ac.at

Victor López-Jaquero  
LoIIS Research Group,  
DIA, University of Castilla-La Mancha  
Albacete, Spain  
victor.lopez@uclm.es

Pascual González  
LoIIS Research Group,  
DIA, University of Castilla-La Mancha  
Albacete, Spain  
pascual.gonzalez@uclm.es

ACM Reference Format:

Digital Twin (DT) are emerging as enablers for real-time monitoring and control in IoT systems. In this work, we propose a DT enabled with Active Inference (AI), a framework rooted in the Free Energy Principle (FEP). The AI module provides decision-making capabilities. Our approach is evaluated on a realistic edge computing scenario where two co-located teleprocessing services are connected via a 5G network. The DT is compared to a competitor for limited CPU resources. The DT continuously infers physical entities' states and optimizes their performance to meet Service Level Objectives (SLO). In a series of experiments, we validate the predictive accuracy and control effectiveness of the AI-enabled DT. The DT achieves an average SLO fulfillment above 90% for both services, and predicted throughput improvements are consistent with real observations confirmed through statistical testing.

**CCS Concepts**  
• Computer systems organization → Real-time system architecture

**Keywords**  
Digital Twins, Predictability, Active Inference, Internet of Things, IoT



Upcoming paper at ACM IoT [14]

## III – Summary

### Limitations & Challenges (1/2)

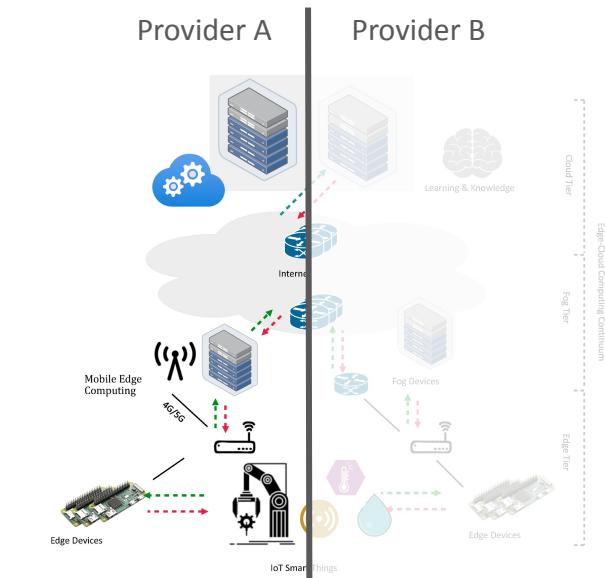
**High cost** of running empirical experiments; slow training progress, needs training environment.

→ Train agents in simulated environments; creating accurate environment, like for **digital twinning**

**Overhead** of extracting baselines from publication; low rigor from generic algorithms, e.g., SB3

→ Support baseline comparison and standardized problem instances in evaluation environment

**Governance** limited to individual vendor; SLAs only supported by one provider and its subset of nodes  
→ Create overarching models for orchestration and client compensation; technical & economical part



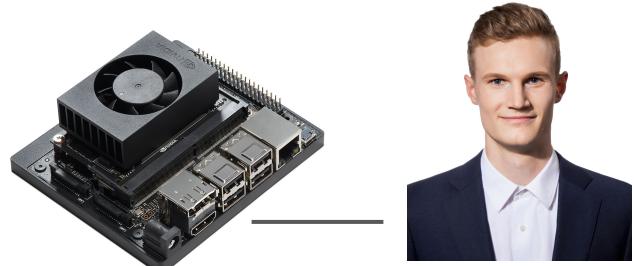
CC infrastructure split between providers

## III – Summary

### Limitations & Challenges (2/2)

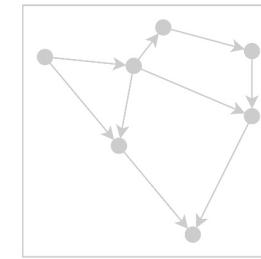
**Missing** sovereignty over data processing; devices rarely in possession of service consumers

→ Integrate personal devices into infrastructure; combine with personal wallets and trusted env.



**Rigid** inference quality for recommending actions; prone to violate operational boundaries, e.g., time

→ Consider resources and context when inferring actions, e.g., smaller graph; create **elastic certainty**



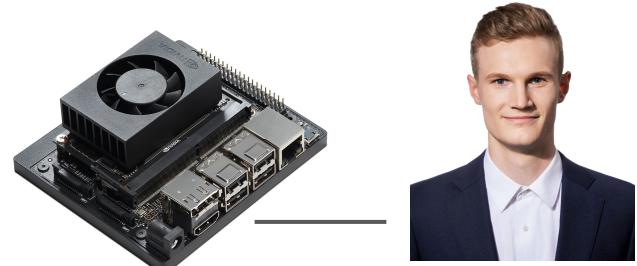
Bayesian network

## III – Summary

### Limitations & Challenges (2/2)

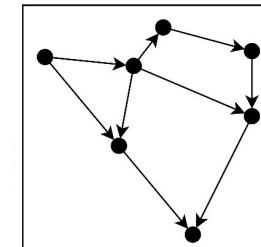
**Missing** sovereignty over data processing; devices rarely in possession of service consumers

→ Integrate personal devices into infrastructure; combine with personal wallets and trusted env.



**Rigid** inference quality for recommending actions; prone to violate operational boundaries, e.g., time

→ Consider resources and context when inferring actions, e.g., smaller graph; create **elastic certainty**



Bayesian network