# Service Orchestration in the Computing Continuum: Structural Challenges and Vision

Boris Sedlak ⓘ, Víctor Casamayor Pujol ⓘ, Ildefons Magrans de Abril ⓘ, *Universitat Pompeu Fabra, Barcelona*

Praveen Kumar Donta ⓘ, *Stockholm University, Stockholm, Sweden*

Adel N. Toosi ⓘ, *University of Melbourne, Melbourne, Australia*

Schahram Dustdar ⓘ, *ICREA, Barcelona, Spain; and TU Wien, Vienna, Austria*

*Abstract*—*The Computing Continuum (CC) integrates different layers of processing infrastructure—from Edge to Cloud—to optimize service quality through ubiquitous and reliable computation. Compared to central architectures, however, heterogeneous and dynamic infrastructure increases the complexity for service orchestration. To guide research, this article first summarizes structural problems of the CC, and then, envisions an ideal solution for autonomous service orchestration across the CC. As one instantiation, we show how Active Inference—a concept from neuroscience—can support self-organizing services in continuously interpreting their environment to optimize service quality. Still, we conclude that no existing solution achieves our vision, but that research on service orchestration faces several structural challenges. Most notably: provide standardized simulation and evaluation environments for comparing the performance of orchestration mechanisms. Together, the challenges outline a research roadmap toward resilient and scalable service orchestration in the CC.*

## INTRODUCTION

Sensory data from Internet of Things (IoT) devices form the backbone of pervasive applications. As running example, consider a digital twin of an entirely smart city that allows clients to interact with content through Augmented Reality (AR) across surfaces, like handhelds or windows. To enter and navigate such environments in real-time, computation is shifted from Cloud centers towards Edge devices; thus, it is possible to process and render content on nearby devices. While this reduces latency, Edge devices offer limited and less predictable resources. During environmental changes, e.g., when facing higher load at rush hours, Edge devices require fallback mechanisms that ensure service quality. To combine the strengths of both worlds, Cloud and Edge layers are integrated in one composite architecture—the Computing Continuum (CC). Thus, applications can place latency-aware services on the Edge, and use the abundance of Cloud resources for hosting the remaining services.

While the CC aims to improve Quality of Experience (QoE), it introduces numerous orchestration challenges [1]. Most notably, by distributing services from one application across different physical devices with heterogeneous characteristics, it gets complex to predict application behavior (e.g., after scaling up one service). This is further complicated by deploying applications across multiple vendors and jurisdictions, as providers may be unwilling to share their full system state. This demands solutions that allow parties with partial observability to collaborate towards common optima, while accounting for their individual behavior.

Driven by the numbers of IoT devices and the need to process data nearby, the last decade has produced a large amount of research to cope with the heterogeneity and distribution of processing systems. Common topics revolve around latency-aware service placement [2], or forming agreements between service and device providers [3]. These mechanisms are often based on static architectures and require a priori understanding of services and devices, e.g., benchmarking the service quality across device types. However,

---

individual devices in the CC will (dis)appear dynamically, which can include new device types that are not benchmarked yet. Also, clients can always redefine the desired service operation—usually specified through *Service Level Objectives* (SLOs)—Thus putting the service in an unknown context. To summarize these problems, this article provides a structured overview of hypotheses that describe the CC, its applications, and challenges for service orchestration.

To guide research on service orchestration—also inspired by the well-known *vision of autonomic computing* [4]—this article outlines how continuous service interpretation and adaptation can help optimize SLO fulfillment. To instantiate this design, we present a preliminary implementation using Active Inference (AIF) [5]—a concept from neuroscience that aims to create self-organizing components that maintain internal requirements fulfilled: First, we model the interactions between a service and its environment through a behavioral Markov blanket (MB)—a description of *how* a service interprets its current state and *which* corrective action to take [6]. Thus, services can decide how to react depending on the context, e.g., by shifting computation or services accordingly. Second, whenever the context changes, e.g., after dynamically reconfiguring SLOs, these behavioral models must be updated. Therefore, we wrap each component in a continuous action-perception cycle and adjust its MBs according to environmental feedback [7]. Third, we analyze interactions between services, and additionally, their hosting devices, by composing their MBs [8]. This enhances collaboration within the CC because services can estimate how local actions impact dependent services and the corresponding resource demand.

This implementation, using AIF, allows collaborative agents to continuously model and understand the environment—a first step to address fundamental problems of dynamic CC systems. However, when mapping the implementation to the structural hypotheses and problems, we still identify multiple challenges that impede research on service orchestration. Most notably, we see a clear gap for: (1) large-scale, standardized simulation environments that simplify testing hypotheses and comparing solutions, (2) continuous and context-aware mechanisms for accurate ML inference, and (3) seamless infrastructure composition that allows clients to use infrastructure from multiple vendors or individuals. To achieve structural improvement for service computing, we summarize these challenges in more detail under three key challenge areas.

## STRUCTURAL DEFINITION

Research on CC systems is in an early stage, resulting in numerous beliefs and views about their inherent challenges, like in [3], [7], [9]. In this section, we provide a structural perspective to orchestration problems in CC systems and design a general guideline for solutions to them. Our core objective: create a holistic solution that optimizes processing requirements throughout CC applications; in practice, this means capturing the desired state through SLOs (e.g., expected latency or quality) and making them first-class citizens during the entire system operation. This abstract solution invites and permits different instantiations—the one provided by us in the next section, being one of them.

### Problem Hypotheses

In the following, we provide three sets of hypotheses that define the problem domain, its management, and ways to model this. First, the **domain hypotheses**, are assumptions on inherent problems of the CC.

D.1 **Heterogeneous continuum**. The computing continuum spans all computational tiers, implying a high level of heterogeneity in hardware or network capacity, and energy consumption.

D.2 **Dynamic conditions**. Resource availability, workload intensity, network conditions, and SLOs from clients will vary over time. This can change at multiple time scales, including fast-paced periodic fluctuations, or slow-paced behavioral drifts.

D.3 **Partial and decentralized observability**. Due to distribution, communication delays, and administrative boundaries, no management component can observe the full system state instantaneously.

D.4 **Large-scale**. The CC infrastructure creates a highly-distributed, dense network of computing nodes that spans vast geographical areas.

Second, we define **application hypotheses** that define CC services and their management scope.

A.1 **Interdependent services**. Applications in the CC are composed of multiple interacting services or microservices, making their performance dependent on upstream services and their resources.

A.2 **Multi-tenant and multi-vendor applications**. Applications are composed of services and computing nodes that belong to different parties, so that no entity has control over the entire system.

A.3 **SLO-based objectives**. Each application (or service) is associated with Service Level Objectives (SLOs), expressed in relation to measurable metrics, that define the desired application state.

A.4 **Service instrumentation**. Each service exposes

monitoring data (e.g., metrics, logs, traces) and control interfaces (e.g., scale, migrate, change configuration) to observe and act on services.

Third, we describe the **modeling hypotheses**, so which abstraction are valid to describe the problem.

M.1 **State-space representation**. The CC can be described by high-dimension state vectors—$\mathcal{S}$—that comprise resource, service, and SLO states. Still, agents can observe only a subset of the true state.

M.2 **Action-space representation**. The control interfaces exposed by services and devices mark the boundaries of the CC's action space—$\mathcal{A}$. Technically, agents could invoke all types of action remotely; however, we limit their scope and permissions to a hierarchical subset of actions.

M.3 **Stochastic dynamics**: State transitions are stochastic—identical actions can lead to different resulting states due to environmental variability.

Under these hypotheses, we can design solutions that autonomously manage applications within the CC, e.g., by scaling or reconfiguring services to fulfill SLOs. In the following, we envision an abstract solution design for service orchestration in the CC that builds upon the well-known *vision of autonomic computing* [4].

## Abstract Solution Design

To optimize service orchestration in the CC, we provide an abstract design guideline split into two parts: a high-level agentic lifecycle for monitoring and optimizing processes, and an instantiation of this cycle using custom tools and mechanisms. By following this agentic cycle, intelligence can be created at a desired scale, thus interlacing the CC with self-adaptive components. First, regardless of the infrastructure or application, smart components must implement a closed-loop architecture [4]; this forms the **invariant** solution part:

- Observation $\rightarrow$ Knowledge: Monitoring data is collected and interpreted to transform observations $o_t$ to internal state representations $\hat{s}_t$; for example, monitor CPU load to be aware of device utilization.
- Knowledge $\rightarrow$ Decision-making: an action $a_t \in \mathcal{A}$ is inferred based on the state $\hat{s}_t$, an internal knowledge model $\mathcal{M}$, and SLO preferences $\mathcal{O}$; for example, note excessive device utilization, and infer that offloading computation will improve latency.
- Decision-making $\rightarrow$ Actuation: the action $a_t$ is executed in the CC's infrastructure and applications. For example, scaling an application horizontally.
- Actuation $\rightarrow$ Observation: actions modify the system state, generating new observations $o_{t+1}$.
- Coordination $\leftrightarrow$ all components: exchange summaries $\xi_t$ of states, intents, and constraints to im-

prove coordination between decision makers; for example, using knowledge transfer or distillation.

Any control mechanism for the CC must implement these steps in a recurrent loop of the form

$$o_t \rightarrow \hat{s}_t \rightarrow a_t \rightarrow o_{t+1}.$$

While the structure of autonomous control loops is invariant, the individual steps can be instantiated differently to create customized implementations of agents; this forms the **variant** part of the solution:

V.1 **Architecture**. Centralized or distributed agents with different hierarchical depth or granularities, e.g., agents per service, application, region etc.

V.2 **Configuration**. Which actions and states are available at each control layer, how often are control loops (action-perception) executed, and what operational overhead is tolerable per agent.

V.3 **Modeling**. Design the state representation purely metric-based, or graph-based with dependencies. Keep probabilistic beliefs about hidden states or directly observe states. Model state and action spaces as discrete or continuous variables.

V.4 **Algorithm**. Decision models can be based on heuristic rules, numerical optimizers, active inference, (multi-agent) RL, or other techniques.

Based on the presented hypotheses and the abstract solution, the next section shows a custom implementation for parts of this abstract solution.

## PRELIMINARY IMPLEMENTATION

This section outlines a three-level implementation for autonomous service orchestration in CC systems (see Figure 1) that follows the abstract solution design and its agentic lifecycles. First, we model dependencies between services, the processing environment, and SLO fulfillment, e.g., quantify the impact of GPU resources on *energy efficiency* and *latency*. Next, we update these models through a continuous action-perception cycle to keep them up-to-date with changing context; lastly, we compose models into larger structures that quantify dependencies between services, and their resource demand according to services' configuration.

## BEHAVIORAL MARKOV BLANKETS

A MB marks the statistical boundary between an entity and its environment—this defines what it can directly observe and control, and separates it from whatever lies beyond these limits. We use this concept to model *how* a service (or rather its SLOs) is affected by its direct observations and the actions it could take. For
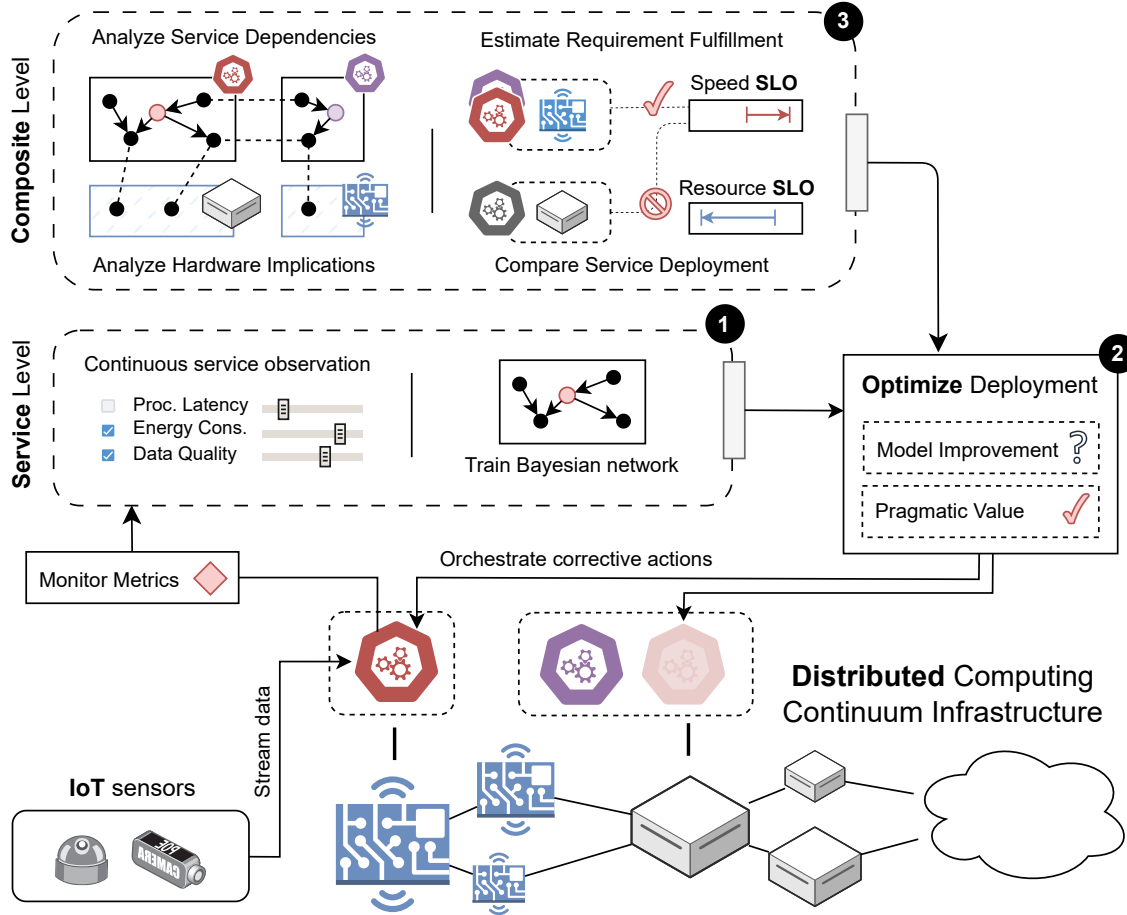
**FIGURE 1**: Three-level implementation for autonomous service orchestration in the Computing Continuum: ❶ we interpret current service states (e.g., *why* was an SLO violated) by training behavioral Markov blankets (MBs) from processing metrics; ❷ agents continuously optimize the service operation according to their internal understanding of the environment (i.e., their MB) and current SLOs. Lastly, ❸ to optimize SLO fulfillment throughout the CC, we compose MBs to quantify dependencies between services and hosting devices.

example, how much *load* an AR application can take before dropping its *performance,* and how lowering the *rendering quality* can recover *performance*. The precise relations and conditional dependencies between these variables are modeled through a Bayesian Network (BN); for instance, the BN could express a function of how the rendering *quality* impacts service *performance*, and given that we change the *quality* to a certain level, what *performance* to expect.

The MB contains all actions that an application or service can initiate itself; when choosing between actions, the BN can answer which action promises the highest SLO improvement. For instance, given that application *performance* is declining, the choice might be between scaling up *resources* or decreasing *quality*. Conditioned on the current state—which tells that additional *resources* are either expensive or depleted—it is possible to make an informed decision [6]. This greatly improves applications' flexibility because individual ser-

vices can act according to changing context.

However, relations and dependencies that make up the MB might either be unknown at design time or change dynamically; to that extent, we continuously update the MBs at runtime using observational data. In our AR overlay example, we monitor the application and collect operational metrics (e.g., rendering *latency*) to evaluate SLOs fulfillment. Whenever SLOs change (e.g., AR application forced to save *energy*), this merely changes the desired metrics distribution. These steps correspond to ❶ in Figure 1. Thus, we analyze the precise impact of environmental dynamics and heterogeneous factors during runtime.

## ACTIVE INFERENCE CYCLE

Equipped with their own behavioral MBs, smart services can perceive and act on their environment without external instructions. For instance, a service agent

can approximate how much load a service can endure before dropping its quality, and take a corrective action, e.g., decreasing sensor *resolution*, to restore service *quality*. When processing sensor data on an Edge device—like movement data of habitants in a smart city—the demand spikes at rush hour and video quality drops at night; to ensure service quality, the agent must act timely. However, running decision-making in the Cloud requires sharing vast amounts of service and device states, which congests the network, slows down reaction, and poses legal problems for private data. Additionally, dynamic IoT environments are prone to variable drifts—leading to inaccurate models.

To ensure that corrective actions show the expected and desired outcome, we wrap a service and its decision model—the MB—into a continuous action-perception cycle, guided by AIF; notice how this implements the abstract control loop presented earlier. In each iteration, the service agent evaluates how accurately it predicted the effects of its last action and updates its MB. In more detail, this means adjusting the conditional variable dependencies according to new observations, for instance, because it overestimated the effect of adjusting the sensor *resolution*. Afterwards, the agent chooses between actions that promise high information gain for improving the MB accuracy, or such that fulfill its preferences, e.g., low rendering latency. This is shown by ❷ in Figure 1.

To allow agents persist over time, AIF aims to minimize *free energy* [5]—an upper bound for uncertainty in a system. Originating from neuroscience, AIF can be used to guide the behavior of humans and service agents alike. Within the control loop, our service agent aims to develop an accurate (i.e., unsurprising) model of the environment, and use it to ensure each component's SLOs. This sets AIF apart from pure RL, because AIF prioritizes model accuracy over solely maximizing cumulative reward [10]. By integrating AIF into service orchestration, we ensure that software components have an accurate understanding of how their local actions can optimize SLO fulfillment.

## COLLABORATIVE ORCHESTRATION

Every computing service has specific demands: video processing services might favor infrastructure with integrated GPUs, or force other services to change their configuration, e.g., downstream services might depend on an upstream service to improve its quality. Other services might be less computationally intense, but require low-latency placement close to an IoT data source. If we deploy services without considering these constraints, it inevitably leads to significant SLO vi-

olations. Often an optimal solution is non-trivial: for example, if we have three services that require GPU support and only resources to accommodate one, how to allocate the resources? This shows that service orchestration in resource-restricted environments cannot be dictated by a single actor, but requires their collaboration to find a global and fair optimum.

To support decentralized and collaborative actuation, we analyze the dependencies between services and hosting devices, more precisely, between the variables located in their MBs. Thus, we explore whether actions have an impact on nearby entities, e.g., if it affects the performance and resource consumption of downstream services, or any service co-located on the same device. Consequently, individual components become aware how local actions (e.g., reconfiguring a service) impact the entire system; this is shown by ❸ in Figure 1. Whenever SLOs change throughout the application, e.g., the desired quality of one service within a pipeline, the upstream services will explore their action space to improve downstream SLO fulfillment. Again, this is a powerful upgrade for the autonomy of computing services because they can jointly find and execute an optimal corrective action.

## PRELIMINARY RESULTS

To ensure the soundness of our three-layer methodology, we evaluated the individual parts through a series of experiments [6], [7], [8]. We summarize the key findings below. In particular, we showed that: (1) training BN and extracting the MB for autonomous processing services did not introduce a critical overhead for common edge devices (e.g., Nvidia Jetson), (2) MBs allowed services to ensure SLO fulfillment by taking optimal actions; at the same time, the BN structure allowed empirically interpreting and verifying these actions, and (3) for composite microservice applications, the MBs of individual services could be equally composed to analyze dependencies between services and optimize global SLO fulfillment.

## SUITABILITY OF THIS IMPLEMENTATION

Considering the extension of structural problems in the CC, no single solution can wrap all—this also shows shortcomings of existing works [2], [3], [11]. Rather, solutions must be modular and evolve organically. In this section, we discuss aspects covered by our preliminary solution, and critical open research gaps.

Starting from the **domain hypotheses**, our approach considers heterogeneity (hardware capacity, energy consumption, network constraints) in the processing environment (D.1) and models dynamic envi-

ronmental conditions as factors inside the MB (D.2); however, it focuses on fast-paced fluctuations, leaving a gap for hierarchical solutions operating at different time scales. The largest gap stems from the difficulty and operational overhead of evaluating orchestration mechanisms in large-scale environments (D.4).

Our approach uses local observations to quantify their implication to global objectives (D.3); however, connecting to the **application hypotheses**, there is a large gap for multi-vendor solutions (A.2) to exchange observations and agree on joint policies. Still, our approach makes use of clear monitoring and action interfaces between services (A.4), which it uses to quantify dependencies between services, their SLOs (A.3), and the processing hardware. Thus, the presented AIF agents follow the **invariant** action-perception cycle and the stochastic dynamics of the processing environment (M.3). While we evaluated different **variants** for agents' algorithms and configurations, it remains to analyze different hierarchical and evolving architectures.

## RESEARCH CHALLENGES

The presented implementation covers large parts of the inherent problems in CC systems, yet it leaves numerous open research challenges—we argue that many of them are structural and generally obstruct research on service orchestration. Below, we recapitulate the most critical gaps and formulate three challenges for accelerating research on service computing. Most importantly, we lack: (1) large-scale, standardized simulation/emulation environments for quickly testing hypotheses and comparing solutions, (2) context-aware and resilient mechanisms for updating and invoking orchestration mechanisms, and (3) seamless infrastructure composition that allows clients to combine devices from multiple vendors and their individual repertoire.

### SYSTEM & ACTION SIMULATION

Empirical evaluations are a standard way to demonstrate technical soundness of hypotheses. However, scaling them to larger environments requires substantial on-demand resources, which are owned by few large-scale companies. For researchers outside this league, the cost of physical experiments (e.g., setup time or hardware expenses) can quickly outweigh benefits [12]. At the same time, extracting and implementing baselines from reference literature is often a tedious and time-consuming task. To lower the resource barrier and enable continuous integration with empirical or emulated environments, we emphasize the following two-stage challenge alongside Figure 2.
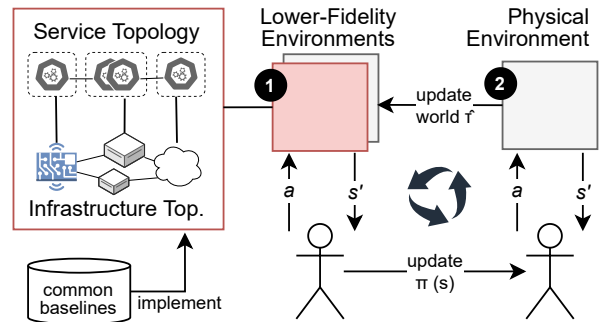


**FIGURE 2**: Missing an ecosystem for ❶ pretraining agents in lower-fidelity environments (e.g., simulation) implementing common baselines and infrastructure; ❷ agents are deployed in physical devices and pretraining environments improve through real-world feedback.

*Baseline Standardization*  The rise of RL has been supported by an entire ecosystem of tools that allow developers to quickly test their ideas on standardized problems (e.g., *OpenAI Gymnasium*[1]) or compare them with contemporary algorithms (e.g., *Stable Baselines*[2]–SB3). While service computing also knows standardized problem instances, such as the *Death-StarBench* suite,[3] running its manifold of microservices requires sophisticated hardware and slows down ML training, since results can only be obtained in real time. Moreover, although baselines such as SB3 are easily accessible, they typically contain generic versions of algorithms that can easily be outperformed by more optimized implementations [13], thus limiting the rigor and comparability of experimental results.

To simplify the training, testing, and deployment of ML solutions in service computing, we identify a substantial gap in lower-fidelity environments (e.g., simulation) that allow researchers and developers to (1) quickly obtain results for training and testing, and (2) compare their solutions with state-of-the-art algorithms. Addressing this challenge would simplify preparation and execution of experiments—lowering the entry barrier for R&D and SMEs—and improve results rigor through standardized baseline comparison.

*Digital Twinning*  Simple simulation environments, such as *OpenAI Gym*, provide hands-on ways to engage with ML technologies, but offer limited practical value for pretraining agents on real-world problems.
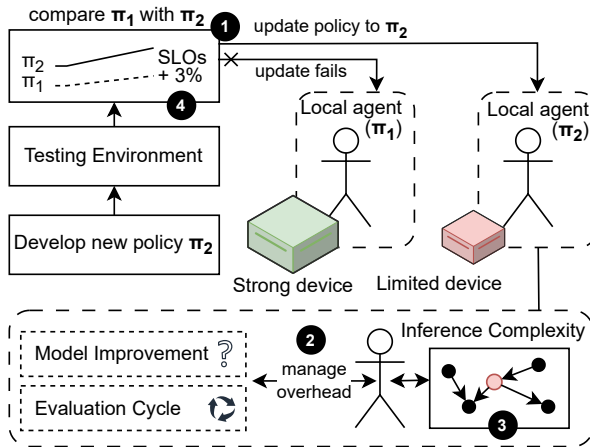
---

**FIGURE 3**: Updates to orchestration policies (e.g., from $\pi_1$ to $\pi_2$) must **1** balance the impact to performance, and tolerate inconsistent model rollout; **2** agents must minimize their execution overhead according to the context, and keep the user in the loop to **3** develop new metrics and **4** align with user intents.

More sophisticated simulators, such as NVIDIA Isaac[4], combine physics engines and real-world benchmarks to create more realistic environments. This mechanism—integrating benchmarks and real-world feedback—is a core feature of digital twins and may form the backbone of self-evolving simulators [14].

A complete ecosystem for accurate benchmarking of CC applications could combine different *fidelity levels*, where simulated results are gradually validated in emulated platforms (e.g., iContinuum [9]) or physical testbeds. This ecosystem could also help replicate realistic system operations from physical setups down to simulation environments; if simulations face realistic conditions (e.g., node failures), it also improves their rigor. Not only would it be costly to train the agent entirely on real node failures, but these negative situations also occur rarely, so the agents can use the lower-fidelity levels to explore these situations thoroughly, and use higher-fidelity environments for fine-tuning.

## MODEL UPDATE & INFERENCE

Service orchestration requires accurate policies, e.g., for fine-grained scaling of resources or replicas. However, as the environment changes, policies quickly become outdated or fail to capture the new context, e.g., as clients change SLOs. Thus, scaling agents might fail to predict the next state and reward accurately, or exceed the available timeframe for inferring a policy. This

is particularly critical for increasingly smaller devices and autonomous components that might not possess sufficient resources to run sophisticated mechanisms for finding and updating policies. We summarize four problems alongside Figure 3: first, continuously update policies and rollout changes despite dynamic disruptions, and second, manage agents' execution overhead on devices with limited resources. Third, we stress new new performance metrics to identify bottlenecks and disruptions within the environment; and fourth, require a human-in-the-loop to ensure correctness.

*Continuous Learning at Scale* To maintain their perception, predictions, and actuation accurate, autonomous agents must adapt to new data—a process known as continuous learning. Given the extension and highly-distributed nature of CC applications, this requires communicating and coordinating policy updates at unprecedented scale. Consider alone the impact of inconsistent or brittle network conditions during model rollout: the outcome will be agents with poorly synchronized models, causing undesired effects, like unfair resource allocations among tenants. From a networking perspective, we argue that dynamic disruptions cannot be completely prevented, and hence, agents require robust orchestration mechanisms to (co-)operate gracefully despite diverging policy versions. Simulating failures for improving the training stability can be supported by common tools, like Chaos Mesh[5].

At the same time, the overhead of continuously integrating model updates must be weighed against the accuracy gained in return, and the cold-start introduced by any model update. This creates an optimization problem in which agents' hyperparameters (e.g., update rate or accuracy thresholds) need to be adjusted at numerous points in the architecture. However, current ML solutions often rely on static or central hyperparameter optimization—as common in Federated Learning (FL) [15]—without the options to dynamically adapt and reconfigure the learning process. This solicits mechanisms that align the update rate with hardware limitations [11], however, the challenge is integrating such context-aware approaches at larger scale, and performing dynamic hyperparameter optimizations at increasingly smaller components.

*Context-aware Orchestration Quality* Time-critical domains, like robotic control [13], frequently require complex operations on sensory data. While the CC—and in particular Edge computing—can deploy services in

---

[4]https://developer.nvidia.com/isaac

[5]https://github.com/chaos-mesh/chaos-mesh

close vicinity, this exposes services to dynamic environmental conditions, like fluctuating load and limited resources. Under these circumstances, scaling agents must equally be able to orchestrate services under strict latency requirements. For example, an AR-supported robot navigating a smart city might offload and scale nearby services in real-time. To not impede other services running on these devices, scaling agents must consider their own computation overhead, the effects on co-located applications, and deduce hardware constraints for inference; for example, prioritize fast autoscaling cycles over their high accuracy. Depending on this context, scaling agents could adjust their internal policy calculation, e.g., by switching between ML models [14], or early-exiting NNs [16].

*Novel Performance Metrics*   Understanding a system often depends on interpreting metrics that provide insight into hidden states [7]. As CC environments become increasingly dynamic and complex—featuring heterogeneous devices, fluctuating resource availability, and frequent mobility—there is a growing need for context-aware metrics that reflect operating conditions. For example, *next-level metrics* could account for environmental volatility, multi-layer coordination costs, or a system's ability to stabilize performance under uncertainty. Creating such metrics may require capturing new states or combining existing ones. Incorporating these metrics into simulation and testing frameworks would enable more realistic assessments of orchestration, scheduling, and adaptation mechanisms.

*Intent Control & Human-in-the-Loop*   Agents can be trained to ensure SLO fulfillment, e.g., through RL and control-loop systems [2]. This, however, does not guarantee that system behavior actually aligns with strategic goals of human operators. To ensure a match, orchestration must shift from rigid SLO definitions to intent-based control: operators define high-level intents, which are continuously diffused to lower-level SLOs and configurations. In this loop, Large Language Models (LLMs) could serve as critical intermediaries, translating high-level human intents (e.g., "prioritize low latency for emergency streams") into scaling actions (e.g., offload services to make space for an emergency service) [17]. This would democratize access to complex infrastructure, as non technical clients (e.g., city planners) can deploy applications and rely on automatic orchestration according to their intents.

## INTERPLAY OF INFRASTRUCTURE

A vision of CC systems is to combine infrastructure layers into a single platform, where services are op-
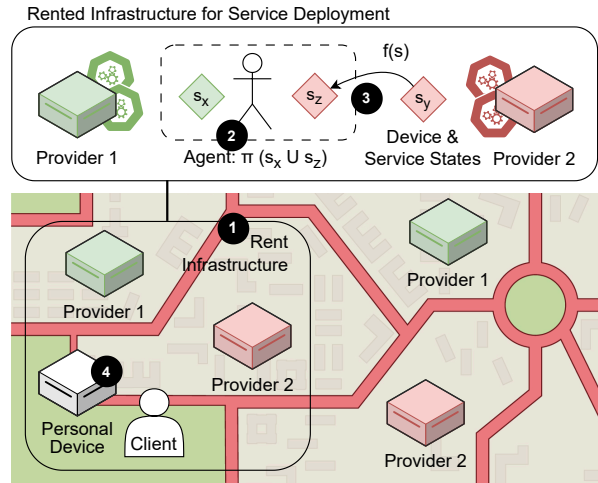


Rented Infrastructure for Service Deployment

FIGURE 4: Client wants to ❶ rent and combine the nearest infrastructure from multiple vendors; agents must ❷ orchestrate and cooperate towards global objective, ❸ despite limited visibility; ideally, this allows ❹ seamless interplay with client-owned devices.

timized autonomously, encapsulated from the application provider. However, looking at today's provider landscape, e.g., for Cloud infrastructure, there is no interplay between providers [18]: clients must weigh between different providers and their subset of the overall infrastructure. Another promise of the CC—or Edge computing—is processing confidential data close to the data source. While Edge devices are integrated into providers' portfolios, like AWS Outposts,[6] data is still processed outside the customer's sovereignty. To support the collaboration of infrastructure providers and empower citizens to take control of personal data, we formulate four challenges alongside Figure 4.

*Multi-Vendor Governance*   Collaboration between different infrastructure providers is a necessity for optimizing the customer experience; ideally, service consumers could process data on the nearest device—despite provider boundaries. The fragmented infrastructure landscape is primarily caused by business interests, but regardless, brings numerous technical challenges, such as distributing the governance between multiple vendors. Consider the following example: we scale an application horizontally and place replicas on devices of operator *A*, and also operator *B*; how would these operators divide the responsibility to ensure the application goals, or more critically, recompense the client in case of SLA violations?

---

[6] https://aws.amazon.com/edge/services/

To tackle this challenge, we see a design-time and runtime part to it: during **design time**, agreements must be formed between infrastructure providers, and between providers and clients. This process can be supported through smart contracts [3] and might lead to a multitude of SLAs. However, as SLOs are violated, there is a gap for accounting mechanisms that identify the partial faults of providers and automatically recompense clients accordingly. During **runtime**, the service deployment must be optimized across vendors—an upgrade to existing single-vendor mechanisms [2].

*Partial observability*   CC systems use infrastructure from multiple vendors that operate across different data jurisdictions. The constraints imposed by these jurisdictions, together with vendors' privacy regulations, restrict data monitoring and remote intervention capabilities. These limitations motivate research into multi-agent algorithms capable of performing effectively even when information sharing is only possible between some jurisdictions, vendors, or agents, and when data sharing must comply with specific privacy constraints.

*Partial cooperation*   Even when vendors act honestly and comply with their individual SLAs, a multi-vendor setup can still drift into non-cooperative behavior. This resembles multiple musicians playing their parts perfectly yet producing a dissonant ensemble. In practice, vendors may prioritize actions that safeguard their own SLA compliance rather than collaborating with others. Under tight conditions, this can lead to *handoff timing optimizations* (accelerating local processing in ways that shift load to downstream vendors), or *selective throttling* (limiting request or resource usage for certain clients). Hence, we must detect non-cooperative vendor behavior without breaking their autonomy.

Altogether, the CC demands decentralized orchestration mechanisms that support: partial visibility, partial cooperation, and continuous updates. Combining these features, however, easily results in oscillating actions and cascading failures. Therefore, risk management must equally be decentralized so that individual services are always accountable for failures they cause. This incentivizes service providers to develop a clear model of the quality and reliability they can provide; otherwise, SLA violations will prove costly and service consumers will avoid using the service.

*Democratizing Infrastructure*  The CC poses the unique opportunity to allow service users to integrate their own infrastructure into a larger computing platform. First, this would allow clients to constrain confidential operations to local or neighboring device, which they fully trust and over which they exhibit a certain sovereignty. Second, by integrating their devices into the CC, clients would benefit from higher availability and fault tolerance—like offloading computations to remote nodes if they exceed local capabilities.

To empower service consumers in controlling the flow and processing of data, there is active research on: access and governance of personal data [19], and creation of trusted execution environments [20]. A logical next step would be to integrate these mechanisms into the CC, and allow individuals to add personal infrastructure. This faces numerous challenges, most critically, ensuring the security of consumer devices, but also service orchestration. Namely, we lack mechanisms to scale services between trusted environments, and business models that reward device owners for providing their own personal infrastructure.

## CONCLUSION

This article provides a structural perspective into the Computing Continuum (CC) and the challenges that its heterogeneous and dynamic infrastructure raises for service orchestration. To guide research in that domain, we first envision an autonomous solution for service adaptation, and second, show how Active Inference (AIF)—a concept from neuroscience—could fill this gap. Following an agentic lifecycle, services can interpret their context, maintain accurate internal models, and reason about interdependencies with other services. While our implementation addresses fundamental problems—like heterogeneous and dynamic operating conditions—numerous aspects remain unaddressed by contemporary research. To bridge this gap, we formulate three research areas that aim to simplify the execution and interpretation of experiments, the continuous and efficient updates of AI models, and the interplay of infrastructure from multiple providers and clients. We invite other researchers to break down these structural barriers together and accelerate research on service orchestration in the CC.

## ACKNOWLEDGMENTS

## REFERENCES

1. P. Plebani, S. Schulte, D. A. Tamburri, and S. Dustdar, "Service-Oriented Computing: A Trajectory for

Research to 2030," *IEEE Internet Computing*, vol. 28, no. 3, pp. 59–63, May 2024.

2. S. Lan, Z. Duan, S. Lu, B. Tan, S. Chen, Y. Liang, and S. Chen, "SLA-ORECS: an SLA-oriented framework for reallocating resources in edge-cloud systems," *Journal of Cloud Computing*, Jan. 2024.

3. P. Kochovski, V. Stankovski, S. Gec, F. Faticanti, M. Savi, D. Siracusa, and S. Kum, "Smart Contracts for Service-Level Agreements in Edge-to-Cloud Computing," *Journal of Grid Computing*, vol. 18, no. 4, pp. 673–690, Dec. 2020.

4. J. Kephart and D. Chess, "The vision of autonomic computing," *Computer*, vol. 36, no. 1, Jan. 2003.

5. T. Parr, G. Pezzulo, and K. J. Friston, *Active Inference: The Free Energy Principle in Mind, Brain, and Behavior*. The MIT Press, 2022.

6. B. Sedlak, V. C. Pujol, P. K. Donta, and S. Dustdar, "Designing Reconfigurable Intelligent Systems with Markov Blankets," in *ICSOC*, 2023, pp. 42–50.

7. ——, "Equilibrium in the Computing Continuum through Active Inference," *Future Generation Computer System*, vol. 160, pp. 92–108, 2024.

8. ——, "Markov Blanket Composition of SLOs," in *2024 IEEE International Conference on Edge Computing and Communications (EDGE)*, 2024.

9. N. Akbari, A. N. Toosi, J. Grundy, H. Khalajzadeh, M. S. Aslanpour, and S. Ilager, "iContinuum: An Emulation Toolkit for Intent-Based Computing Across the Edge-to-Cloud Continuum," in *2024 IEEE International Conference on Cloud Computing*, Jul. 2024.

10. A. Tschantz, B. Millidge, A. K. Seth, and C. L. Buckley, "Reinforcement Learning through Active Inference," in *Bridging AI and Cognitive Science*, 2020.

11. A. Piaseczny, K. C. Shisher, S. Wang, and C. G. Brinton, "RCCDA: Adaptive Model Updates in the Presence of Concept Drift under a Constrained Resource Budget," in *NeurIPS*, 2025.

12. A. B. Bowden, "Designing Field Experiments to Integrate Research on Costs," *AERA Open*, Jan. 2023.

13. L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, and A. Madry, "Implementation Matters in Deep RL: A Case Study on PPO and TRPO," in *ICLR*, Apr. 2020.

14. E. Pretel, B. Sedlak, V. Casamayor-Pujol, E. Navarro, V. López-Jaquero, P. González, and S. Dustdar, "Active Inference for Digital Twins: Predicting and Optimising IoT Processing Service Performance," in *Proceedings of the 15th International Conference on the Internet of Things (IoT 2025)*, 2025.

15. A. Danilenka, A. Furutanpey, V. C. Pujol, B. Sedlak, A. Lackinger, M. Ganzha, M. Paprzycki, and S. Dustdar, "Adaptive Active Inference Agents for Heterogeneous and Lifelong Federated Learning," Oct. 2024.

16. Y. Chen, Z. Niu, M. Roveri, and G. Casale, "CEED: Collaborative Early Exit Neural Network Inference at the Edge," in *IEEE INFOCOM*, May 2025.

17. Z. Su, M. Islam, M. Goudarzi, and A. Toosi, "LLM-Driven Intent-Based Privacy-Aware Orchestration Across the Cloud–Edge Continuum," in *ACSW*, 2026.

18. D. Saxena, R. Gupta, and A. K. Singh, "A Survey and Comparative Study on Multi-Cloud Architectures: Emerging Issues And Challenges For Cloud Federation," Aug. 2021.

19. R. Verborgh, "Re-decentralizing the Web, For Good This Time," in *Linking the World's Information: Essays on Tim Berners-Lee's Invention of the World Wide Web*, 1st ed., Sep. 2023, vol. 52, pp. 215–230.

20. A. Muñoz, R. Ríos, R. Román, and J. López, "A survey on the (in)security of trusted execution environments," *Computers & Security*, 2023.

**Boris Sedlak** is Postdoctoral researcher at the Engineering department of UPF, Barcelona, in the Distributed Intelligence and Systems-Engineering Lab (DISL). Contact him at `boris.sedlak@upf.edu`.

**Víctor Casamayor Pujol** is a Tenure track professor at the Engineering department of UPF, Barcelona, in the Distributed Intelligence and Systems-Engineering Lab (DISL). Contact him at `victor.casamayor@upf.edu`.

**Ildefons Magrans de Abril** is currently working as a Postdoctoral researcher at the Engineering department of UPF, Barcelona, in the Distributed Intelligence and Systems-Engineering Lab (DISL). Contact him at `ildefons.magrans@upf.edu`.

**Praveen Kumar Donta** (SM'22) is Associate Professor (Docent) at the Department of Computer and Systems Sciences, Stockholm University, Sweden. His research includes learning in Distributed Continuum Systems. Contact him at `praveen@dsv.su.se`.

**Adel N. Toosi** is an Associate Professor of Computer Systems and Director of the *Dis*tributed Systems and *Net*work Applications Lab (*DisNet* Lab) at the University of Melbourne, Australia. Contact him at `adel.toosi@unimelb.edu.au`.

**Schahram Dustdar** is Full Professor of Computer science heading the Research Division of Distributed Systems at the TU Wien, Austria. He is also part-time research professor at ICREA | UPF Barcelona. Contact him at `dustdar@dsg.tuwien.ac.at`.