

IoT 2025  
Vienna, Austria

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

15th International Conference on the Internet of Things (IoT 2025)

Elena Pretel<sup>1</sup>, Boris Sedlak<sup>2</sup>, Víctor Casamayor-Pujol<sup>2</sup>, Elena Navarro<sup>1</sup>, Víctor López-Jaquero<sup>1</sup>,  
Pascual González<sup>1</sup>, and Schahram Dustdar<sup>2,3</sup>

<sup>1</sup> University of Castilla-La Mancha, Albacete, Spain

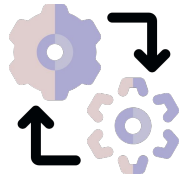
<sup>2</sup> Universitat Pompeu Fabra, Barcelona, Spain

<sup>3</sup> TU Wien, Vienna, Austria





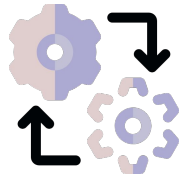
# Theoretical Foundations: Digital Twins (DTs)



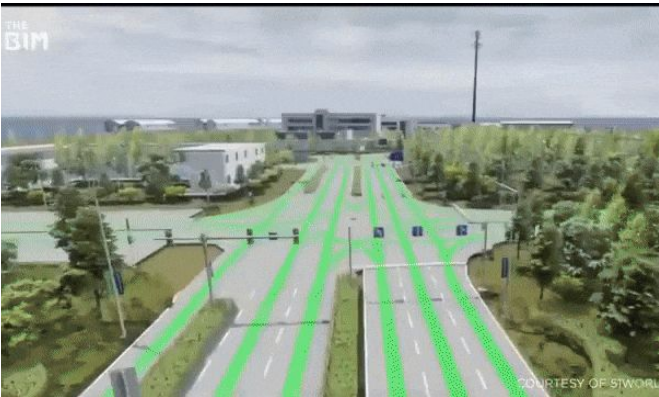
A DT is a virtual replica that allows monitoring and controlling a physical (processing) system.



# Theoretical Foundations: Digital Twins (DTs)



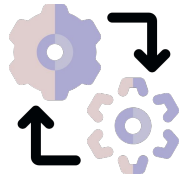
A DT is a virtual replica that allows monitoring and controlling a physical (processing) system.



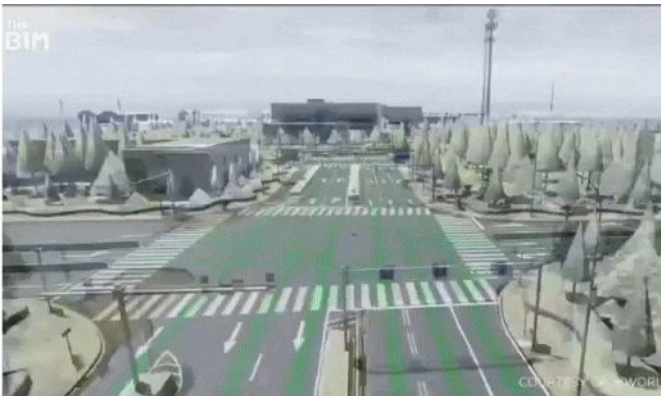
Model Shanghai in Unreal Engine



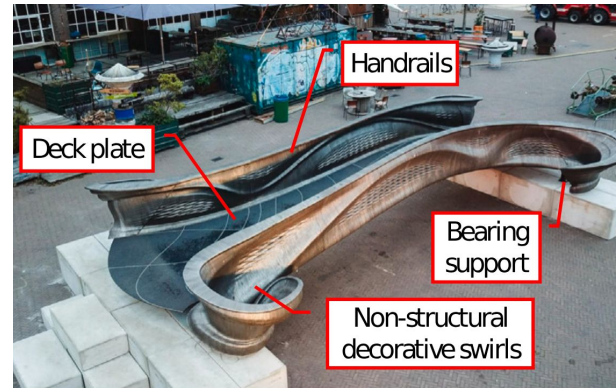
# Theoretical Foundations: Digital Twins (DTs)



A DT is a virtual replica that allows monitoring and controlling a physical (processing) system.



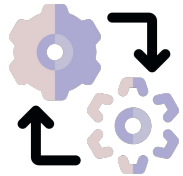
Model Shanghai in Unreal Engine



Create a bridge from a 3D printer



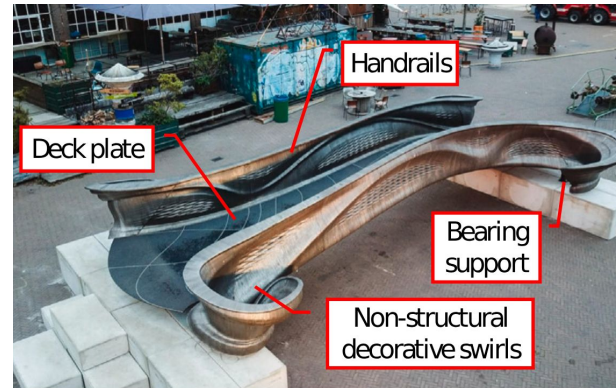
# Theoretical Foundations: Digital Twins (DTs)



A DT is a virtual replica that allows monitoring and controlling a physical (processing) system.



Model Shanghai in Unreal Engine



Create a bridge from a 3D printer

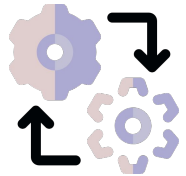


Run object detection with YOLOv8





# Theoretical Foundations: Digital Twins (DTs)



A DT is a virtual replica that allows monitoring and controlling a physical (processing) system.



Model Shanghai in Unreal Engine



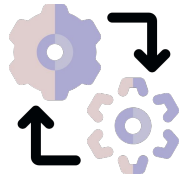
deploy  
optimise



Run object detection with YOLOv8



# Theoretical Foundations: Digital Twins (DTs)



A DT is a virtual replica that allows monitoring and controlling a physical (processing) system.



Model Shanghai in Unreal Engine



deploy  
optimise



Run object detection with YOLOv8

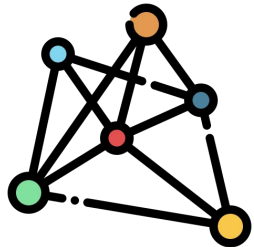
❑ **Predictability** – ensure accuracy of predictions despite **changing environment**

So what actually is **Active Inference**?  
What can it **contribute** to Digital Twins?





# Theoretical Foundations: Active Inference (AIF)



**Origin:** From neuroscience—systems (e.g., humans) act to minimize surprise or uncertainty. Agents create internal representations of their (processing) environment that allow them interpreting and acting efficiently.

**Core idea:** AIF agents constantly predict sensory inputs, compare to reality, and update their internal model. Actions not merely reactive, but reduce uncertainty.

## AIF cycle:

- **Perception** → infers hidden states
- **Action** → changes the environment to meet expectations.
- **Learning** → updates beliefs from experience



# Theoretical Foundations: Active Inference (AIF)



**Internalize bus schedule**

<https://www.pexels.com/photo/men-waiting-on-bus-stop-in-london-18111366/>

<https://www.japan-experience.com/plan-your-trip/to-know/understanding-japan/stepping-stones>



# Theoretical Foundations: **Active Inference (AIF)**



**Internalize bus schedule**



**Building hypothesis of actions**

<https://www.pexels.com/photo/men-waiting-on-bus-stop-in-london-18111366/>

<https://www.japan-experience.com/plan-your-trip/to-know/understanding-japan/stepping-stones>





# Theoretical Foundations: Active Inference (AIF)



**Internalize bus schedule**



**Building hypothesis of actions**



**Predict IoT Performance**

<https://www.pexels.com/photo/men-waiting-on-bus-stop-in-london-18111366/>

<https://www.japan-experience.com/plan-your-trip/to-know/understanding-japan/stepping-stones>



# Theoretical Foundations: Active Inference (AIF)

Create Digital Twin by **using** X

$X = \{\text{Physics Engine, Kafka, Active Inference, etc}\}$

**Using** Active Inference → creating a Digital Twin



# Theoretical Foundations: Active Inference (AIF)

~~Create Digital Twin by using X~~

$X = \{\text{Physics Engine, Kafka, Active Inference, etc}\}$

**Using** Active Inference → creating a Digital Twin





# Proposal and Architecture

## Case Study: AIF-Enabled DT in Edge Computing

**Physical twin:** Edge device with limited CPU power

**Two concurrent processing services:**

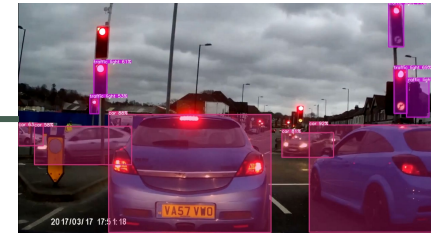
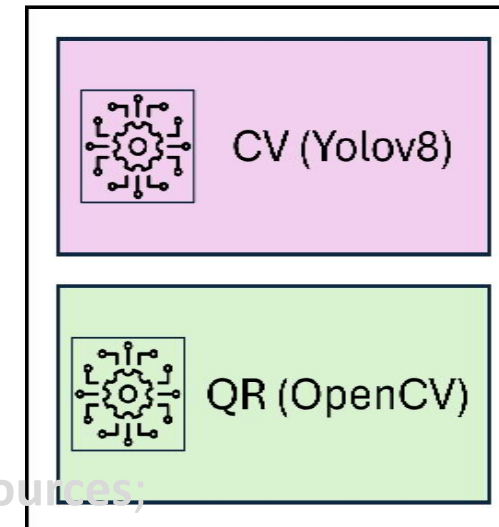
- YOLOv8 video inference
- QR code reading service

**Observe and configure services:**

- CV: throughput, data quality, cores, model size
- QR: throughput, data quality, cores

Both services compete for limited CPU resources;  
thus, creating a dynamic optimization scenario

**Physical Twin**  
Edge Device





# Proposal and Architecture

## Case Study: AIF-Enabled DT in Edge Computing

**Physical twin:** Edge device with limited CPU power

**Two concurrent processing services:**

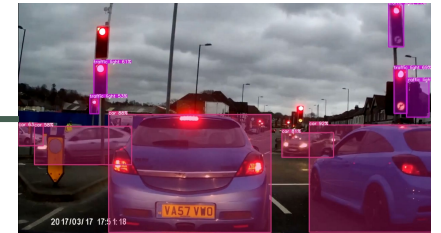
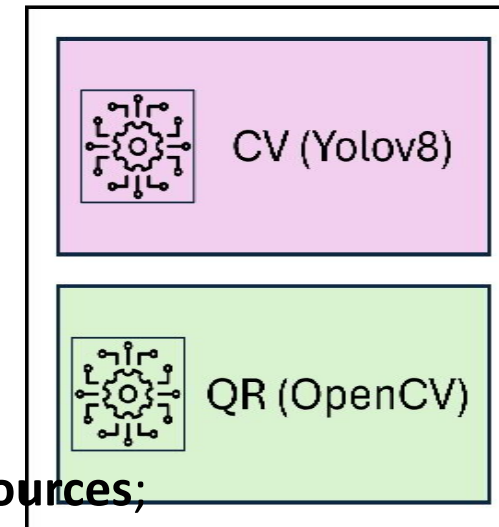
- YOLOv8 video inference
- QR code reading service

**Observe and configure services:**

- CV: throughput, data quality, cores, model size
- QR: throughput, data quality, cores

Both services **compete** for **limited CPU resources**;  
thus, creating a dynamic optimization scenario

**Physical Twin**  
Edge Device





# Proposal and Architecture

## Case Study: AIF-Enabled DT in Edge Computing

**Physical twin:** Edge device with limited CPU power

**Two concurrent processing services:**

- YOLOv8 video inference
- QR code reading service

**Observe and configure services:**

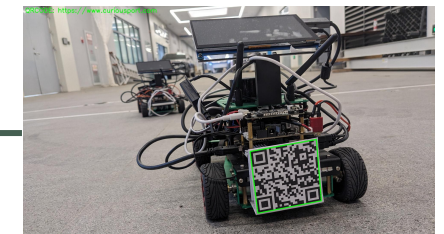
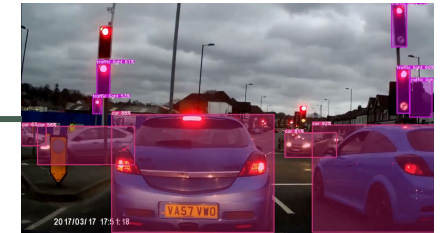
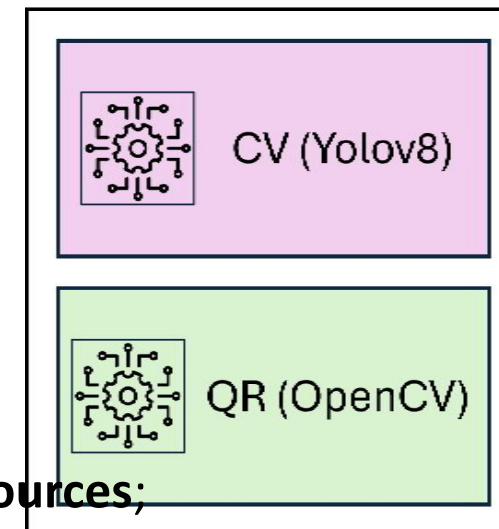
- CV: throughput, data quality, cores, model size
- QR: throughput, data quality, cores

**Predict**

**Actions**

Both services **compete** for **limited CPU resources**;  
thus, creating a dynamic optimization scenario

**Physical Twin**  
Edge Device





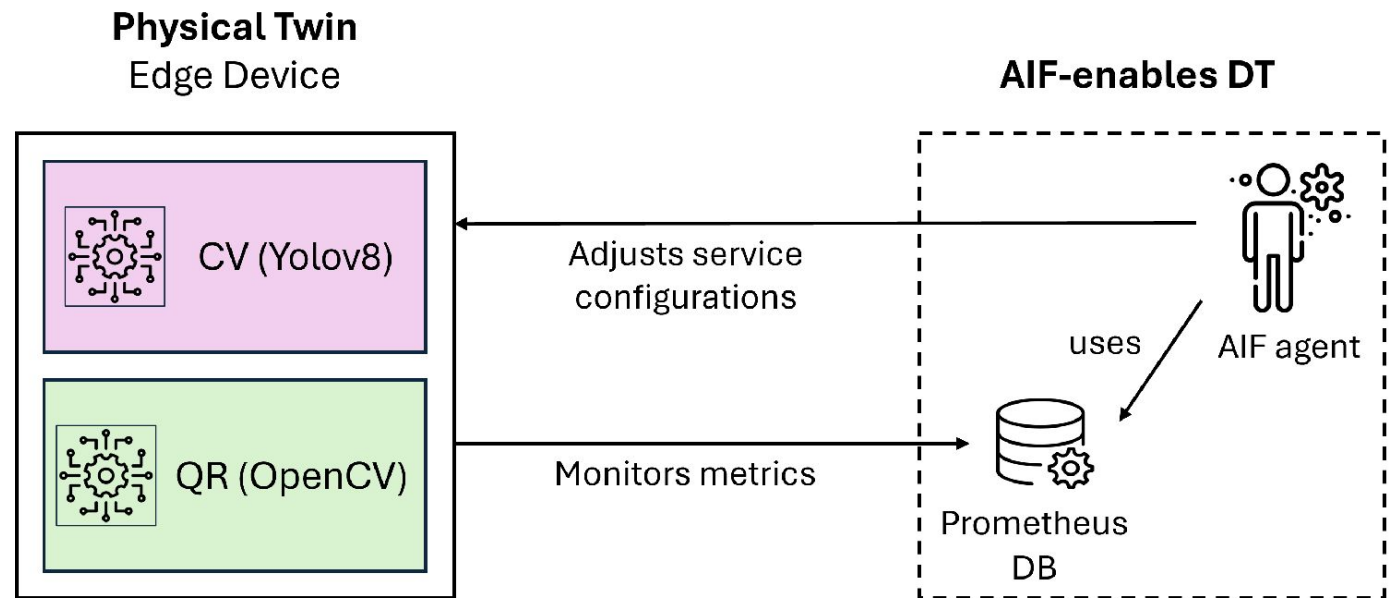
# Proposal and Architecture

## Architecture of the AIF-Enabled DT

**Prometheus Database:** collects and stores real-time service metrics (CPU usage, throughput, data quality, etc.)

**AIF Agent:** use service observations, update beliefs, and act to optimize model accuracy and service performance

**Closed-loop control:** DT observes → infers → acts → updates continuously





# Proposal and Architecture

AIF agent is modeled as a POMDP

Component	Details / Variables in this Case Study
<b>States (S)</b>	CV throughput, CV data quality, CV cores, CV model size, QR throughput, QR data quality, QR cores. <b>Constraint:</b> $CV\_cores + QR\_cores \leq 8$ .
<b>Actions (U)</b>	- <b>CV actions (U<sub>CV</sub>)</b> : increase/decrease <i>quality</i> , <i>CPU cores</i> , <i>model size</i> , or <i>no-op</i> . - <b>QR actions (U<sub>QR</sub>)</b> : increase/decrease <i>quality</i> , <i>CPU cores</i> , or <i>no-op</i> . Combined as a <b>joint action</b> ( $U = [U_{CV}, U_{QR}]$ ).
<b>Transition Model (B)</b>	Encodes probability of moving from one state to another after an action, give small prior about variable dependencies (e.g., cores $\rightarrow$ throughput)
<b>Preferences (C)</b>	Define the agents “goals” or Service Level Objectives (SLOs). - <b>CV</b> : throughput $\geq 3$ FPS, high quality $\geq 288$ px, large model $\geq$ YOLOv8m. - <b>QR</b> : throughput $\geq 40$ FPS, high quality $\geq 800$ px.



# Proposal and Architecture

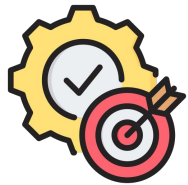
AIF agent is modeled as a POMDP

Component	Details / Variables in this Case Study
<b>States (S)</b>	CV throughput, CV data quality, CV cores, CV model size, QR throughput, QR data quality, QR cores. <b>Constraint:</b> $CV\_cores + QR\_cores \leq 8$ .
<b>Actions (U)</b>	<ul style="list-style-type: none"><li>- <b>CV actions (U<sub>CV</sub>):</b> increase/decrease <i>quality</i>, <i>CPU cores</i>, <i>model size</i>, or <i>no-op</i>.</li><li>- <b>QR actions (U<sub>QR</sub>):</b> increase/decrease <i>quality</i>, <i>CPU cores</i>, or <i>no-op</i>.</li></ul> Combined as a <b>joint action</b> ( $U = [U_{CV}, U_{QR}]$ ).
<b>Transition Model (B)</b>	Encodes probability of moving from one state to another after an action, give small prior about variable dependencies (e.g., cores $\rightarrow$ throughput)
<b>Preferences (C)</b>	Define the agents “goals” or Service Level Objectives (SLOs). <ul style="list-style-type: none"><li>- <b>CV:</b> throughput <math>\geq 3</math> FPS, high quality <math>\geq 288</math>px, large model <math>\geq</math> YOLOv8m.</li><li>- <b>QR:</b> throughput <math>\geq 40</math> FPS, high quality <math>\geq 800</math>px.</li></ul>





# Experiments and Results



**Goal:** Validate whether the **AIF-enabled DT** can

- **Predict** the physical twin's performance (RQ1)
- **Take** actions to meet **Service Level Objectives (SLOs)** (RQ2)

## Experimental Setup



- **Platform:** Edge device (8 CPU cores) running two services
- **Environment:** Docker containers + Prometheus monitoring
- **AIF implementation:** pymdp library
- **Duration:** 50 AIF cycles per experiment, 30 repetitions with different seeds
- **Metrics:** RMSE (prediction accuracy), SLO fulfillment rate (%)



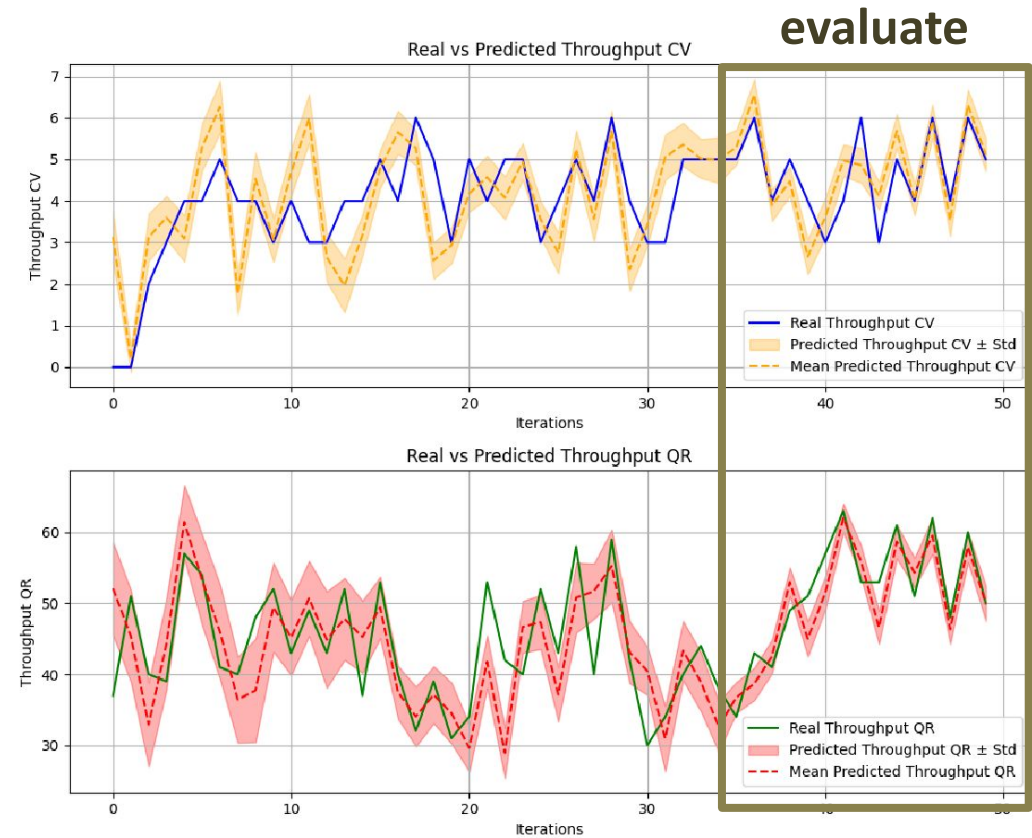
# Experiments and Results

## Experiment 1 (RQ1):

Predict the physical twin's performance (predicted vs. real throughput) in **last 15** iterations.

First 35 iterations learning transition model

Test	Metric	p-value	Statistical Conclusion
Mann–Whitney U	Distributions with $\alpha=0.05$	CV: 0.491 QR: 0.087	No significant difference <b>predictions <math>\approx</math> real values</b>





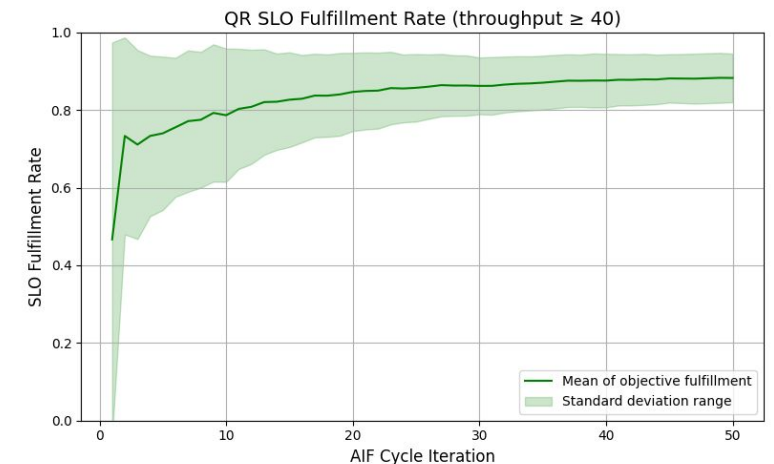
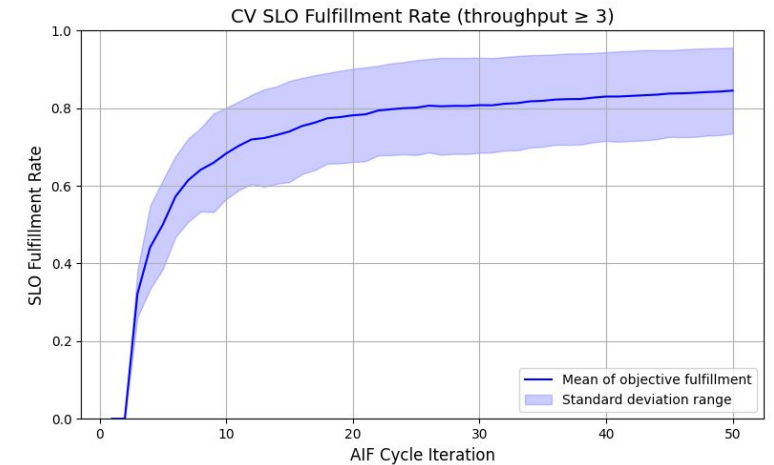
# Experiments and Results

## Experiment 2 (RQ2):

Use DT to optimize physical twin's performance (i.e., SLOs).

AIF agent starts with minimal prior understanding  
(i.e., hint on variable dependency, but no **distributions**)

Test	Metric	p-value	Statistical Conclusion
Wilcoxon Signed-Rank	Median vs. 0.8 threshold	CV: 0.001 QR: 0.010	Median > 0.8 <b>Objectives met significantly</b>





# Experiments and Results

## Experiment 2 (RQ2):

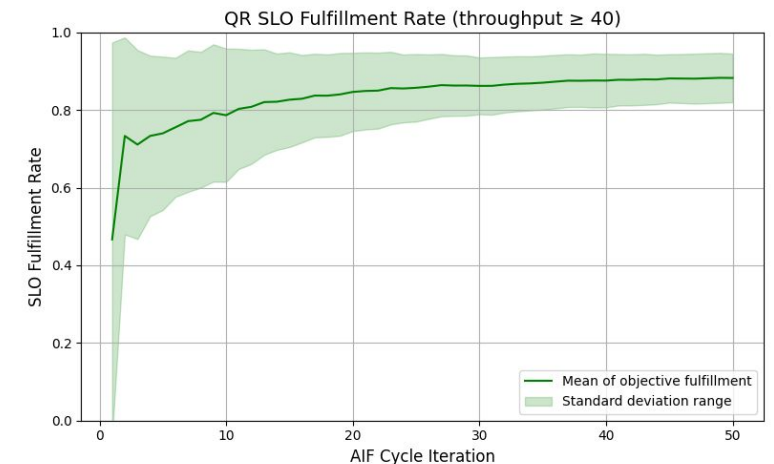
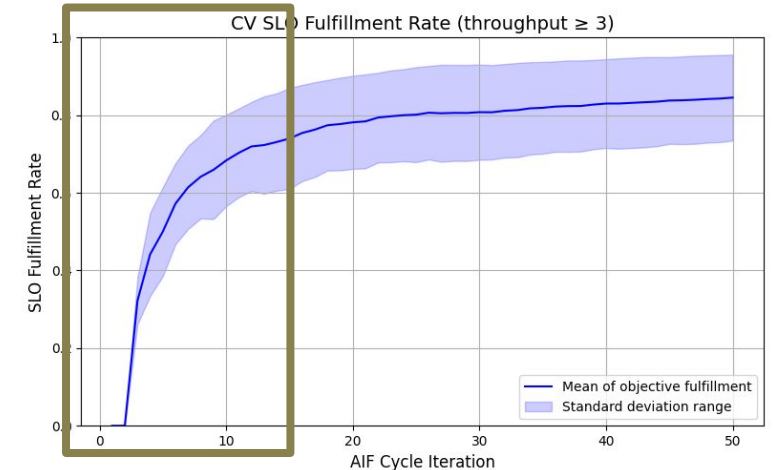
Use DT to optimize physical twin's performance (i.e., SLOs).

AIF agent starts with minimal prior understanding  
(i.e., hint on variable dependency, but no **distributions**)

Test	Metric	p-value	Statistical Conclusion
Wilcoxon Signed-Rank	Median vs. 0.8 threshold	CV: 0.001 QR: 0.010	Median > 0.8 <b>Objectives met significantly</b>



steep!





# Conclusion and Future Work

Complexity of IoT systems requires real-time monitoring and control.

Introduced a **novel approach** to enhance **predictability** of DT using **Active Inference (AIF)**.

Proposed an **AIF-enabled DT architecture** capable of **perception, prediction, learning and control** in dynamic IoT environments.

Validated that the **AIF-enabled DT**:

- Accurately **predicts** physical twin performance (RQ1).
- Effectively **selects actions** to meet Service Level Objectives (RQ2).
- Achieves >80% **SLO fulfillment** across experiments (RQ2).

**Future Work:** Compare **AIF-enabled DTs** with contemporary ML solutions (e.g., RL).