

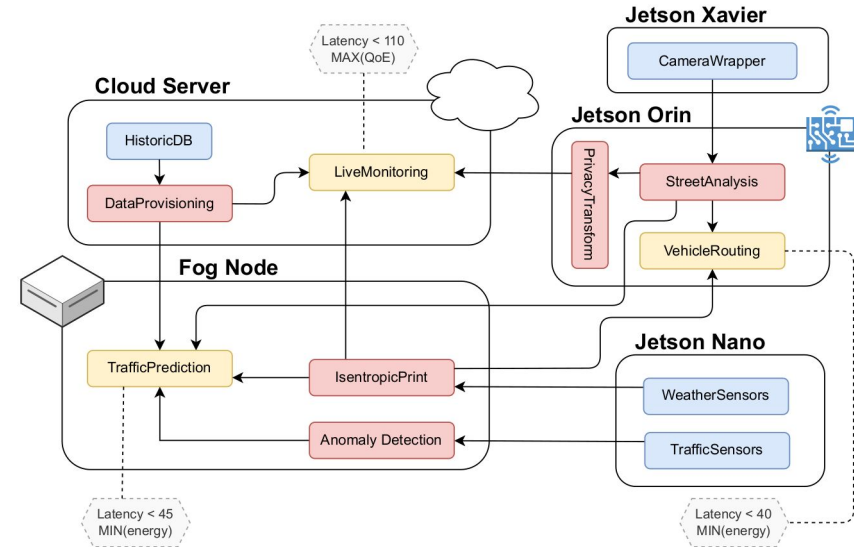
Diffusing High-level SLOs in Microservice Pipelines

TU Wien: Boris Sedlak, Victor Casamayor Pujol, Praveen Kumar Donta, Schahram Dustdar



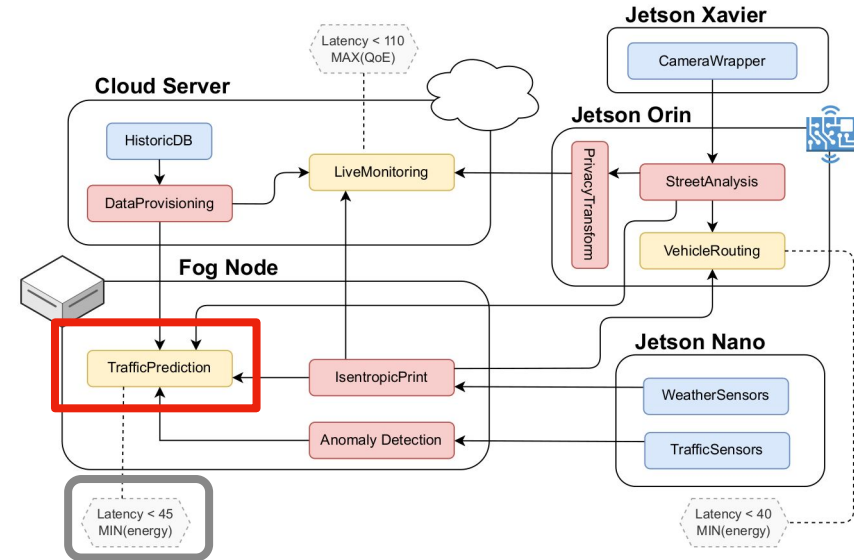
Problem Statement

- ❑ Microservice pipelines **distributed** over computing infrastructure, i.e., Edge to Cloud, unclear implications of individual services to high-level requirements, i.e., **SLOs**
 - ❑ Multiple tenant and vendors (= stakeholders) specify **opposing** SLOs, e.g., minimize energy and response time, which results in **contradicting** service configurations
 - ❑ Stakeholders not aware what their SLO implies for lower level components and services, e.g., **energy** → cpu
-
- ❑ Requires multiple layers of SLOs that specify application requirements through **fine-grained** control mechanism, i.e., diffusing high-level SLOs into lower level SLOs



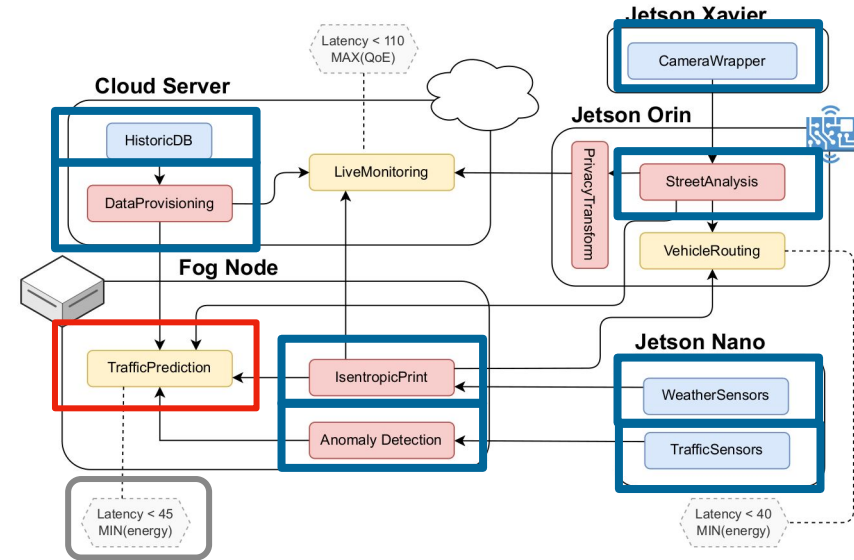
Problem Statement

- ❑ Microservice pipelines **distributed** over computing infrastructure, i.e., Edge to Cloud, unclear implications of individual services to high-level requirements, i.e., **SLOs**
 - ❑ Multiple tenant and vendors (= stakeholders) specify **opposing** SLOs, e.g., minimize energy and response time, which results in **contradicting** service configurations
 - ❑ Stakeholders not aware what their SLO implies for lower level components and services, e.g., **energy** → cpu
-
- ❑ Requires multiple layers of SLOs that specify application requirements through **fine-grained** control mechanism, i.e., diffusing high-level SLOs into lower level SLOs



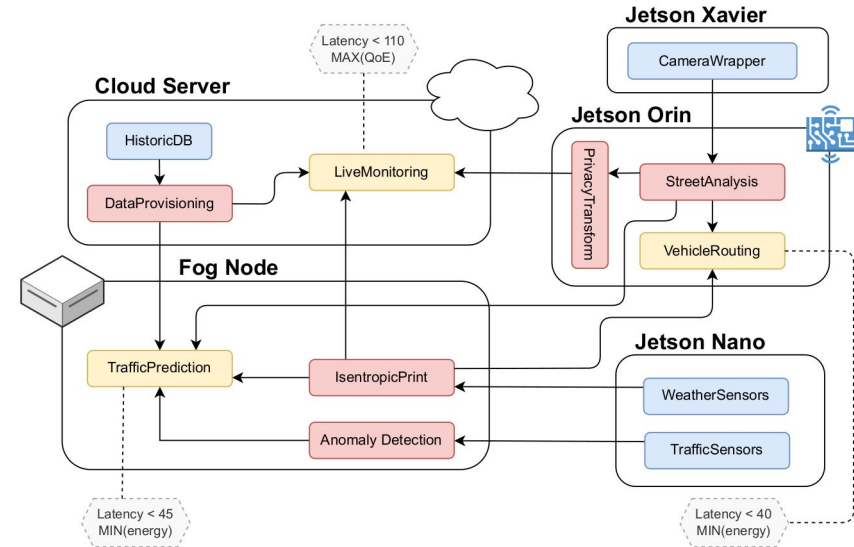
Problem Statement

- ❑ Microservice pipelines **distributed** over computing infrastructure, i.e., Edge to Cloud, unclear implications of individual services to high-level requirements, i.e., **SLOs**
 - ❑ Multiple tenant and vendors (= stakeholders) specify **opposing** SLOs, e.g., minimize energy and response time, which results in **contradicting** service configurations
 - ❑ Stakeholders not aware what their SLO implies for lower level components and services, e.g., **energy** → cpu
-
- ❑ Requires multiple layers of SLOs that specify application requirements through **fine-grained** control mechanism, i.e., diffusing high-level SLOs into lower level SLOs



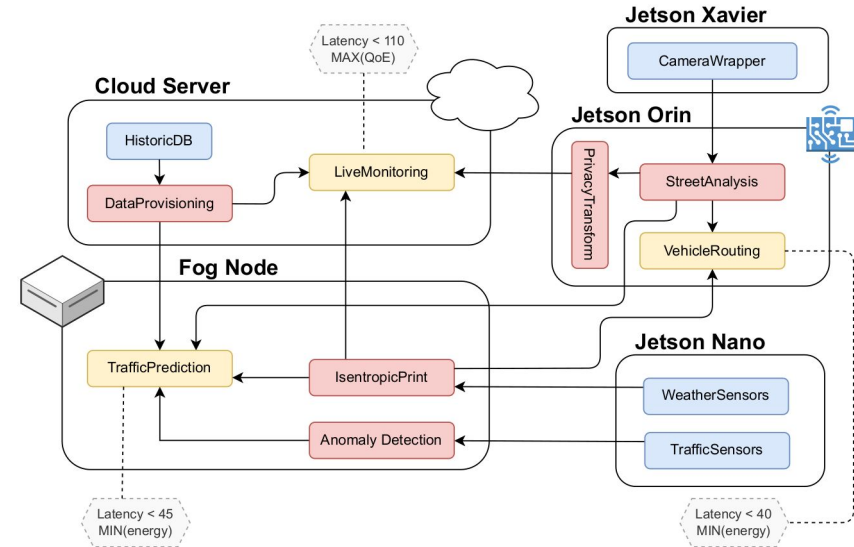
Problem Statement

- ❑ Microservice pipelines **distributed** over computing infrastructure, i.e., Edge to Cloud, unclear implications of individual services to high-level requirements, i.e., **SLOs**
 - ❑ Multiple tenant and vendors (= stakeholders) specify **opposing** SLOs, e.g., minimize energy and response time, which results in **contradicting** service configurations
 - ❑ Stakeholders not aware what their SLO implies for lower level components and services, e.g., **energy** → cpu
-
- ❑ Requires multiple layers of SLOs that specify application requirements through **fine-grained** control mechanism, i.e., diffusing high-level SLOs into lower level SLOs



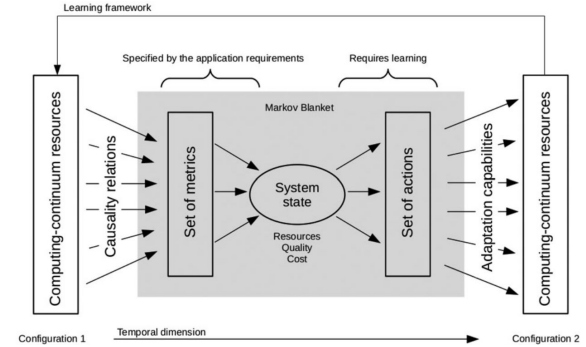
Problem Statement

- ❑ Microservice pipelines **distributed** over computing infrastructure, i.e., Edge to Cloud, unclear implications of individual services to high-level requirements, i.e., **SLOs**
 - ❑ Multiple tenant and vendors (= stakeholders) specify **opposing** SLOs, e.g., minimize energy and response time, which results in **contradicting** service configurations
 - ❑ Stakeholders not aware what their SLO implies for lower level components and services, e.g., **energy** → cpu
-
- ❑ Requires multiple layers of SLOs that specify application requirements through **fine-grained** control mechanism, i.e., diffusing high-level SLOs into lower level SLOs

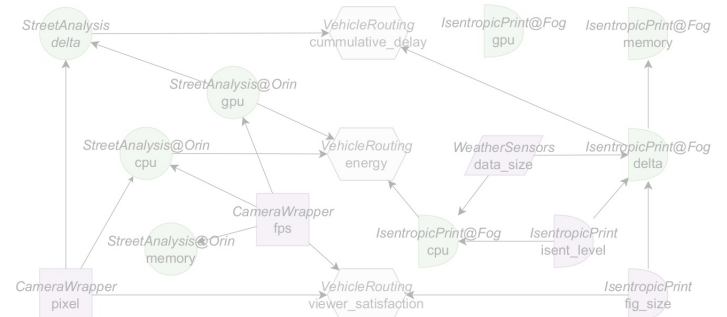


MBs express how systems interact with their environment through **perception** (e.g., metrics) and **actions** (e.g., scaling)

Model the **behavior** of a system in changing environments; focuses on variables that have an impact on taken actions



Behavioral Markov blanket for a system [1]

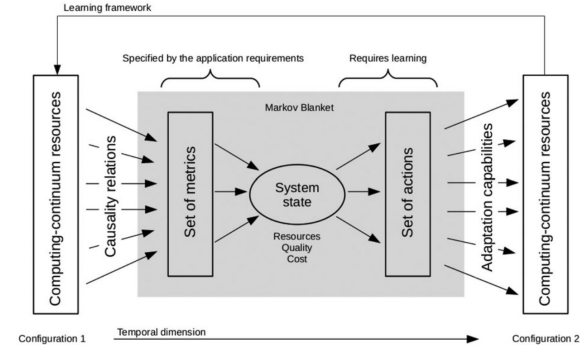


Causal dependencies between hierarchical microservices

[1] Dustdar, Casamayor Pujol, and Donta; On Distributed Computing Continuum Systems (2023)

MBs express how systems interact with their environment through **perception** (e.g., metrics) and **actions** (e.g., scaling)

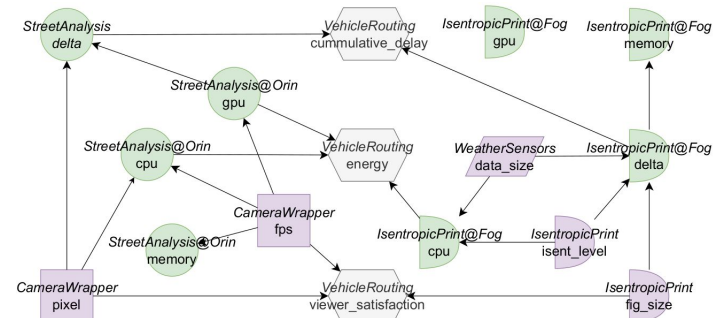
Model the **behavior** of a system in changing environments; focuses on variables that have an impact on taken actions



Behavioral Markov blanket for a system [1]

Note two fundamental properties of Bayesian Network (BN)

- (1) Variables reflecting high-level SLO are **leaf node**; otherwise constraining childs; always has leaves
- (2) Parameter variables always located at **root nodes**; conditional independent of all other nodes

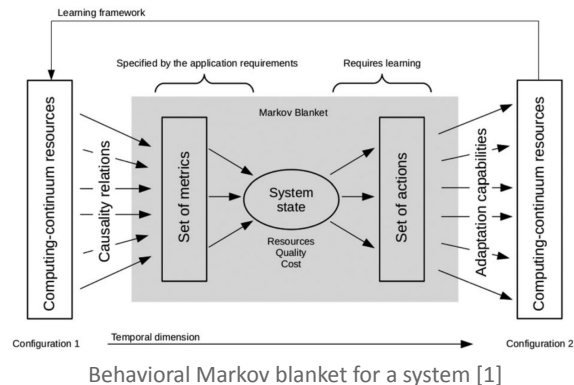


Causal dependencies between hierarchical microservices

[1] Dustdar, Casamayor Pujol, and Donta; On Distributed Computing Continuum Systems (2023)

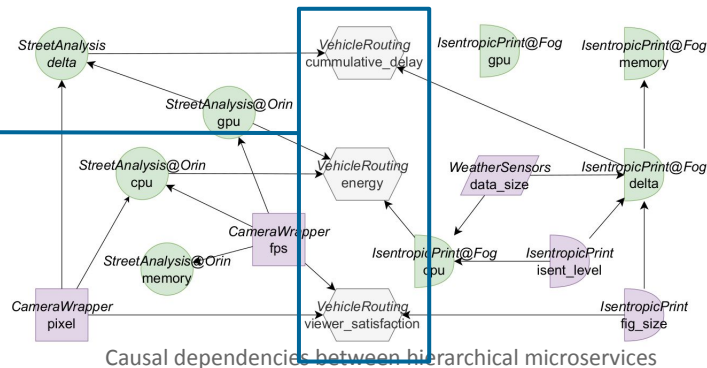
MBs express how systems interact with their environment through **perception** (e.g., metrics) and **actions** (e.g., scaling)

Model the **behavior** of a system in changing environments; focuses on variables that have an impact on taken actions



Note two fundamental properties of Bayesian Network (BN)

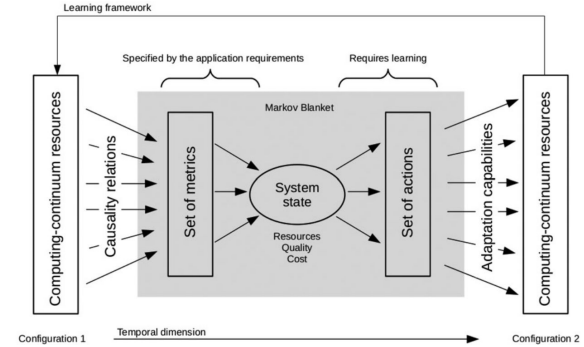
- (1) Variables reflecting high-level SLO are **leaf node**; otherwise constraining childs; always has leaves
- (2) Parameter variables always located at **root nodes**; conditional independent of all other nodes



[1] Dustdar, Casamayor Pujol, and Donta; On Distributed Computing Continuum Systems (2023)

MBs express how systems interact with their environment through **perception** (e.g., metrics) and **actions** (e.g., scaling)

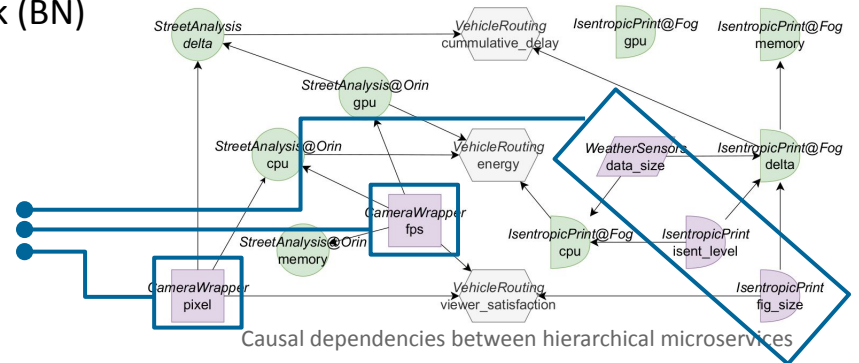
Model the **behavior** of a system in changing environments; focuses on variables that have an impact on taken actions



Behavioral Markov blanket for a system [1]

Note two fundamental properties of Bayesian Network (BN)

- (1) Variables reflecting high-level SLO are **leaf node**; otherwise constraining childs; always has leaves
- (2) Parameter variables always located at **root nodes**; conditional independent of all other nodes



[1] Dustdar, Casamayor Pujol, and Donta; On Distributed Computing Continuum Systems (2023)

RQ-1) *How can high-level SLOs be translated to lower-level objectives?*

Fulfilling high-level SLOs requires equilibrium among **all** components, hence they which require **clear configurations** to achieve this

RQ-2) *How restrictive should low-level SLOs be?*

Predicting SLO behavior is not black-white, but **continuous**. What are desirable states for lower-level SLO and how to achieve them

RQ-3) *Where do SLO conflicts occur and how can they be resolved?*

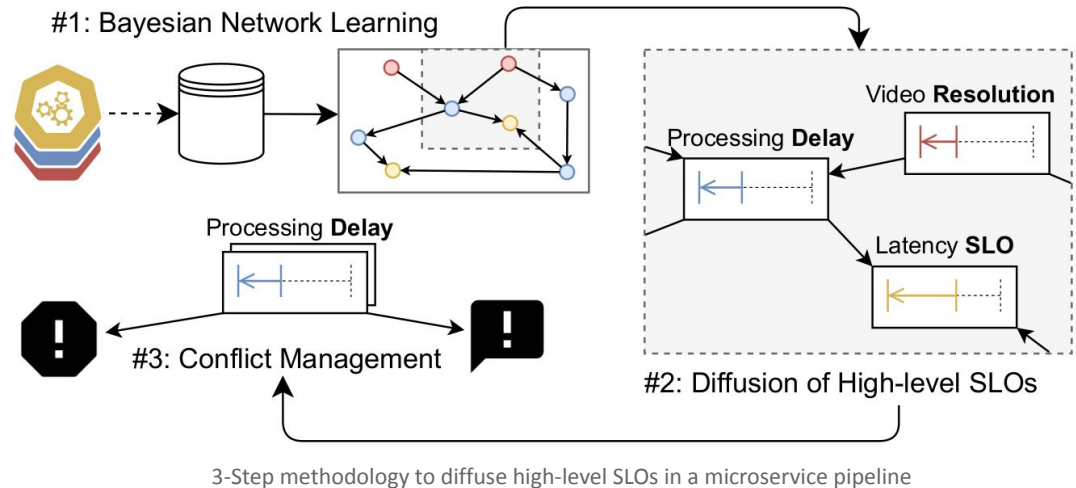
Stakeholders cannot maintain a **reasonable** overview over multiple competing SLOs, e.g., energy vs. performance, how to resolve this

3-Step approach

#1 Extract BN as a probabilistic view into the service execution

#2 diffuse high-level SLOs into lower level ones and param. assignments

#3 identify conflicting variables and resolve them as far as possible

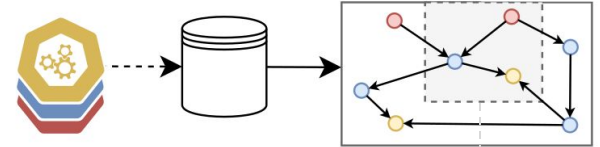


#1 Bayesian Network Learning (BNL)

Extract **causal** dependencies between dependent services; training data collected centrally and used for BNL, must be captured in parallel or joined through **interface variables**

Renders a **composite graph** for a services tree starting from the consumer, i.e., restricts the number of variables per case

#1: Bayesian Network Learning

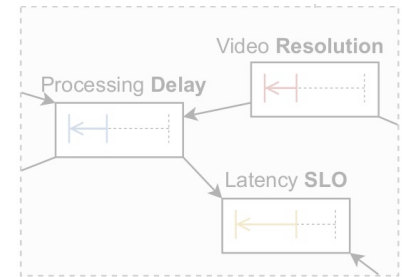


#2 Diffusion of High-level SLOs

Recursively traverse the **parents** of high-level SLOs and constraint states of lower-level variables and parameters; visiting variables multiple times constrains them further

Evaluates each low-level state's probability to **satisfying** high-level constraints; only considers a state iff $p > \max(X) * \lambda$

X = list of all state's probabilities to fulfill SLOs; λ = acceptance range



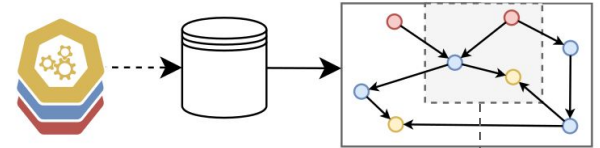
#2: Diffusion of High-level SLOs

#1 Bayesian Network Learning (BNL)

Extract **causal** dependencies between dependent services; training data collected centrally and used for BNL, must be captured in parallel or joined through **interface variables**

Renders a **composite graph** for a services tree starting from the consumer, i.e., restricts the number of variables per case

#1: Bayesian Network Learning

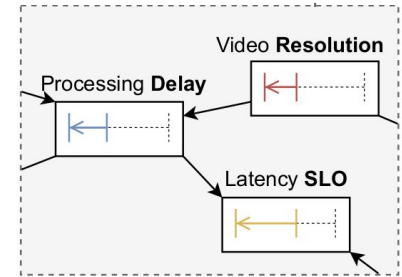


#2 Diffusion of High-level SLOs

Recursively traverse the **parents** of high-level SLOs and constraint states of lower-level variables and parameters; visiting variables multiple times constrains them further

Evaluates each low-level state's probability to **satisfying** high-level constraints; only considers a state iff $p > \max(\mathbf{X}) * \lambda$

\mathbf{X} = list of all state's probabilities to fulfill SLOs; λ = acceptance range



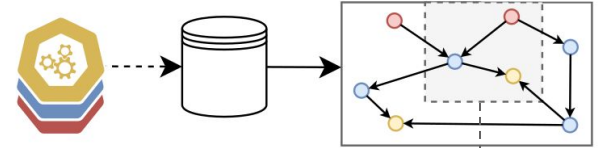
#2: Diffusion of High-level SLOs

#1 Bayesian Network Learning (BNL)

Extract **causal** dependencies between dependent services; training data collected centrally and used for BNL, must be captured in parallel or joined through **interface variables**

Renders a **composite graph** for a services tree starting from the consumer, i.e., restricts the number of variables per case

#1: Bayesian Network Learning



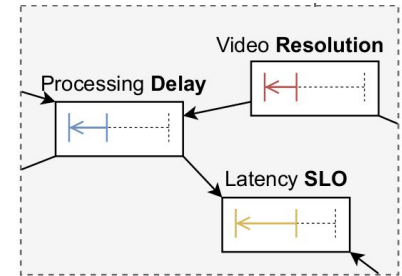
#2 Diffusion of High-level SLOs

Recursively traverse the **parents** of high-level SLOs and constraint s... roles and parameters; visiting vari... trains them further

Evaluates e... ability to **satisfying** high-level constraints, only considers a state iff $p > \max(\mathbf{X}) * \lambda$

Complexity grows with the number of states

\mathbf{X} = list of all state's probabilities to fulfill SLOs; λ = acceptance range



#2: Diffusion of High-level SLOs

#3 Conflict Management

Identify **conflicting variables** and **resolve** them as far as possible; merge the inferred low-level SLOs and try to find an **intersection** among the assignments

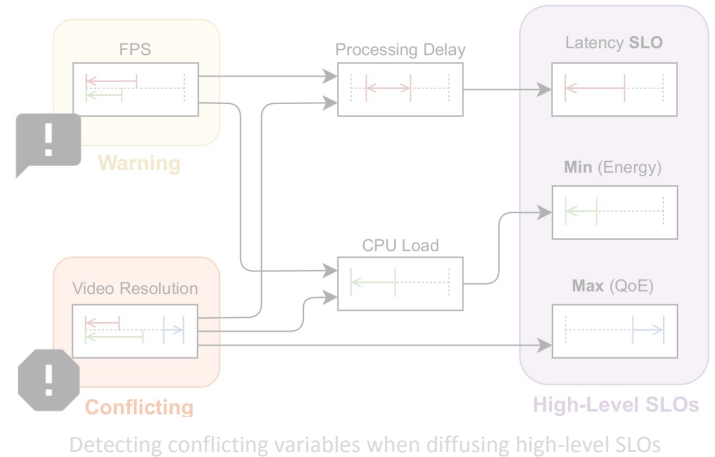
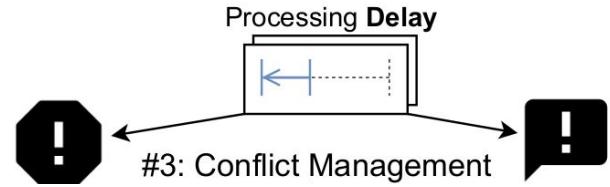
$$\text{INTER}(L_v) = \bigcap_{i,j=1;i \neq j}^n L_i \cap L_j \neq \emptyset$$

Warning

issued in case that variables were visited and constrained multiple times, but it was possible to merge inferred SLOs

Conflict

highlight to stakeholders which variables are conflicting and which high-level SLO are responsible, i.e., **Max** (QoE)



#3 Conflict Management

Identify **conflicting variables** and **resolve** them as far as possible; merge the inferred low-level SLOs and try to find an **intersection** among the assignments

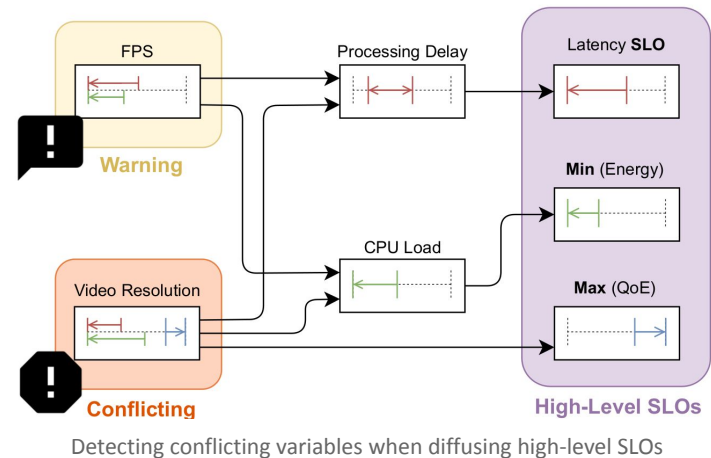
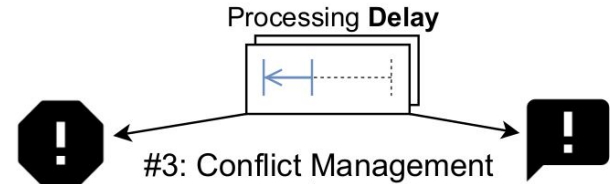
$$\text{INTER}(L_v) = \bigcap_{i,j=1;i \neq j}^n L_i \cap L_j \neq \emptyset$$

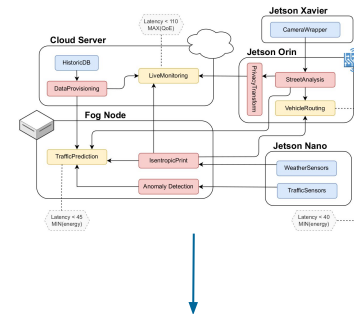
Warning

issued in case that variables were visited and constrained multiple times, but it was possible to merge inferred SLOs

Conflict

highlight to stakeholders which variables are conflicting and which high-level SLO are responsible, i.e., **Max** (QoE)





Combination of 12 microservices **plugged together**; service implementations all collected in [GitHub](#) repository

Microservice pipelines consist of 4 producers, 5 workers, and 3 consumer services; deployed on **different hosts** and feature distinct numbers of **configuration parameters**

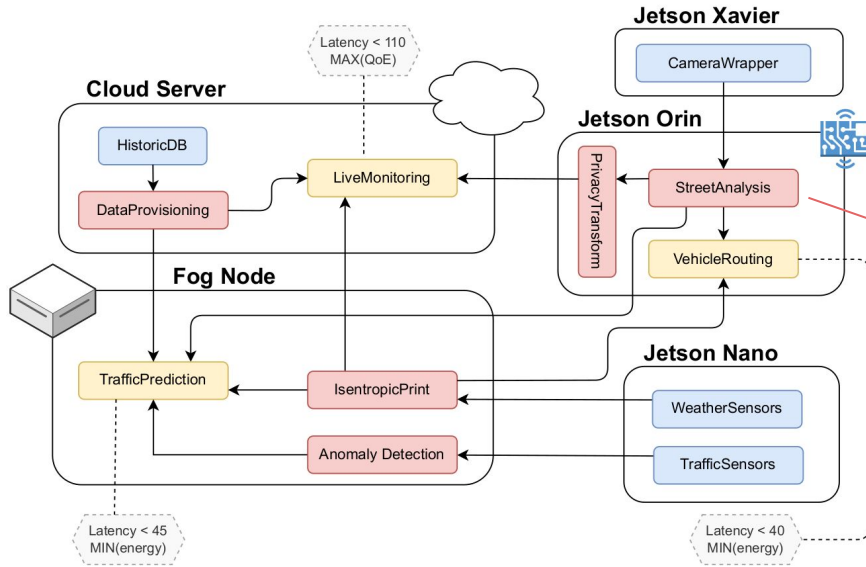
Services executed on the physical setup, provide all the metrics required for the 3-step methodology

TABLE II: Microservices available for evaluation

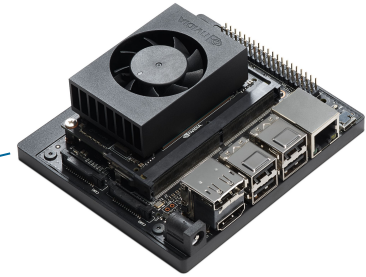
	ID	type	param / var	host
<i>TrafficSensors</i>	28	Producer	1 / 1	<i>Xavier</i>
<i>HistoricDB</i>		Producer	2 / 2	<i>Server</i>
<i>CameraWrapper</i>	29	Producer	2 / 2	<i>Nano</i>
<i>WeatherSensors</i>	30	Producer	1 / 1	<i>Xavier</i>
<i>AnomalyDetection</i>	28	Worker	0 / 5	<i>Fog</i>
<i>HistoricProvision</i>		Worker	2 / 7	<i>Server</i>
<i>StreetAnalysis</i>	31	Worker	0 / 4	<i>Orin</i>
<i>PrivacyTransform</i>	29	Worker	0 / 6	<i>Orin</i>
<i>IsentropicPrint</i>	30	Worker	2 / 6	<i>Fog</i>
<i>TrafficPrediction</i>		Consumer	0 / 2	<i>Fog</i>
<i>VehicleRouting</i>		Consumer	0 / 3	<i>Orin</i>
<i>LiveMonitoring</i>		Consumer	0 / 3	<i>Server</i>

<https://github.com/borissedlak/deploymentOptimizer/tree/main/SOSE>

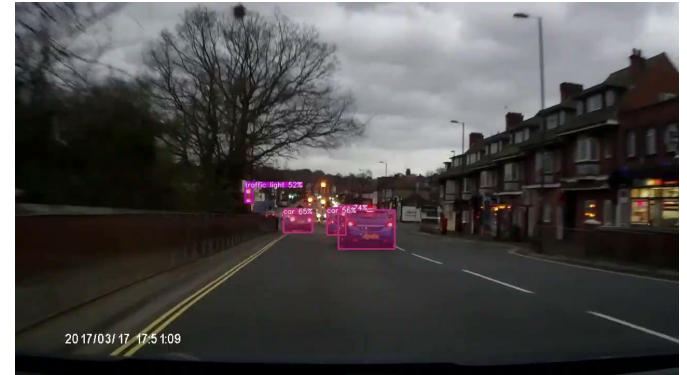
Experimental Setup (2)



Network of microservice pipelines deployed over the computing continuum



Nvidia Jetson Orin board (image from [here](#))



CV Service with Yolov8 running on the produced videos

Results: RQ-1 SLO Diffusion

Given a set of **high-level SLOs** and a **overarching BN**
 Diffuse the SLOs to **lower-level variables** and **parameters**

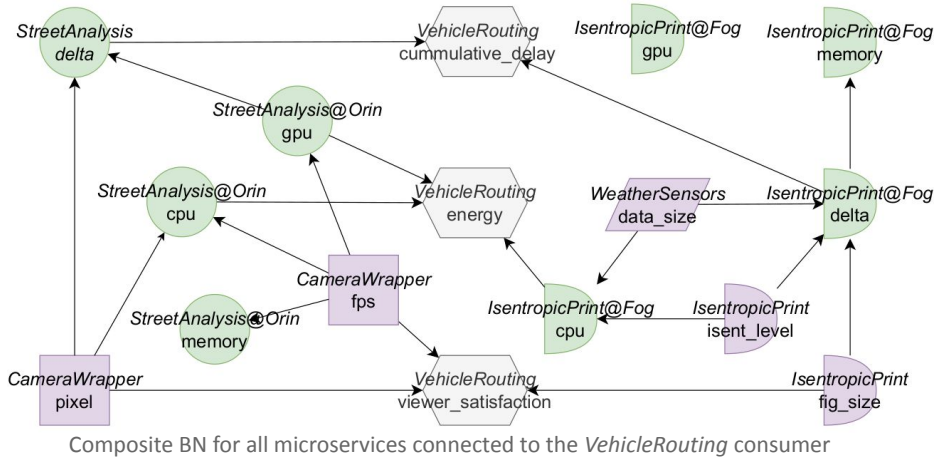


TABLE III: SLOs and parameters inferred for *VehicleRouting*

Microservice	Variable	States	SLO / Param
VehicleRouting	cumm_delay	≤ 45 ms	High-level
	energy	≤ 19 W	
	viewer_sat	—	
<i>StreetAnalysis</i>	delta	≤ 35 ms	Low-level
	cpu (Orin)	≤ 21 %	
	gpu (Orin)	≤ 40 %	
	delta	≤ 37 ms	
<i>IsentropicPrint</i>	cpu (Fog)	≤ 17 %	Parameter
	memory	—	
<i>CameraWrapper</i>	pixel	= 480 p	Parameter
	fps	= 15 f	
<i>IsentropicPrint</i>	fig_size	≤ 50 p	Parameter
	isent_level	≤ 200 k	
<i>WeatherSensors</i>	data_size	≤ 30 pi	Parameter

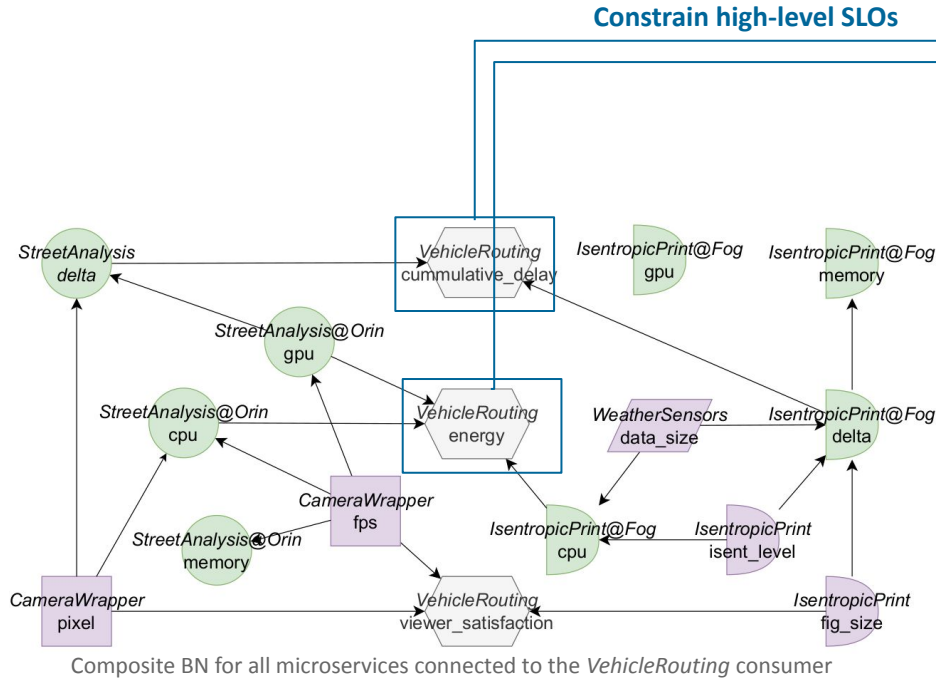
Inferred low-level SLOs and parameter assignments for high-level SLOs

Results: RQ-1 SLO Diffusion

TABLE III: SLOs and parameters inferred for *VehicleRouting*

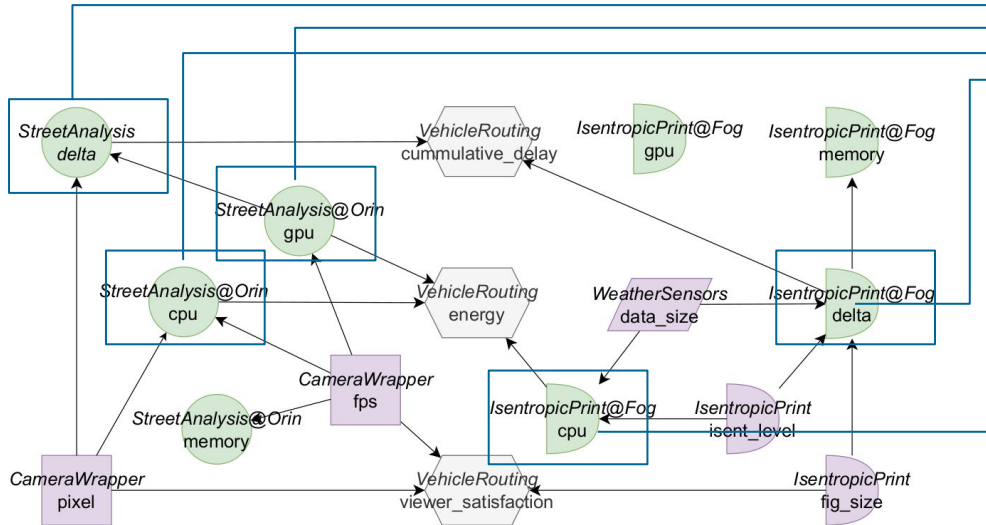
Microservice	Variable	States	SLO / Param
VehicleRouting	cumm_delay	≤ 45 ms	High-level
	energy	≤ 19 W	
	viewer_sat	—	
<i>StreetAnalysis</i>	delta	≤ 35 ms	Low-level
<i>StreetAnalysis</i>	cpu (Orin)	≤ 21 %	
<i>StreetAnalysis</i>	gpu (Orin)	≤ 40 %	
<i>IseotropicPrint</i>	delta	≤ 37 ms	Low-level
	cpu (Fog)	≤ 17 %	
<i>CameraWrapper</i>	pixel	= 480 p	Parameter
<i>CameraWrapper</i>	fps	= 15 f	
<i>IseotropicPrint</i>	fig_size	≤ 50 p	
<i>IseotropicPrint</i>	isent_level	≤ 200 k	
<i>WeatherSensors</i>	data_size	≤ 30 pi	
<i>WeatherSensors</i>	data_size	≤ 30 pi	

Inferred low-level SLOs and parameter assignments for high-level SLOs



Results: RQ-1 SLO Diffusion

Constrain low-level SLOs



Composite BN for all microservices connected to the *VehicleRouting* consumer

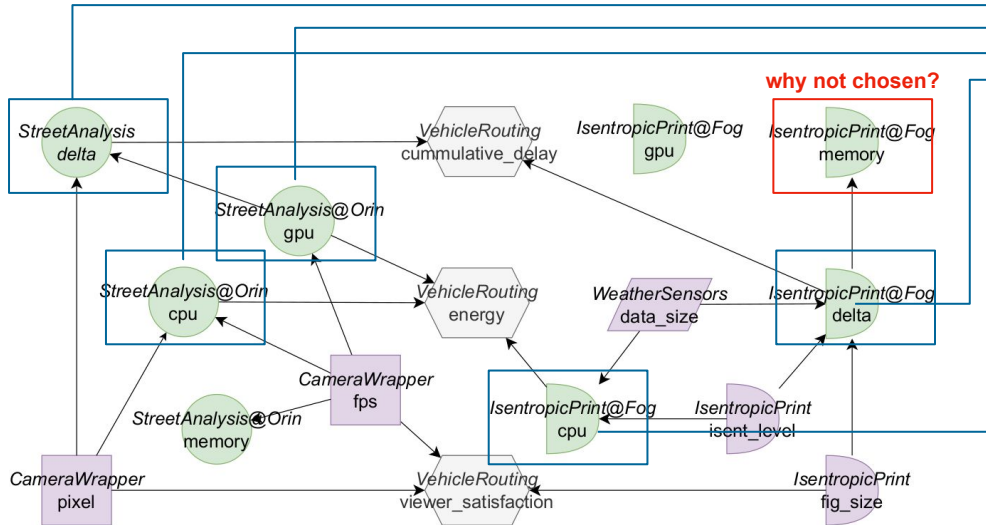
TABLE III: SLOs and parameters inferred for *VehicleRouting*

Microservice	Variable	States	SLO / Param
VehicleRouting	cumm_delay	≤ 45 ms	High-level
	energy	≤ 19 W	
	viewer_sat	—	
StreetAnalysis	delta	≤ 35 ms	Low-level
	cpu (Orin)	≤ 21 %	
	gpu (Orin)	≤ 40 %	
	delta	≤ 37 ms	
IsentropicPrint	cpu (Fog)	≤ 17 %	Parameter
	memory	—	
CameraWrapper	pixel	= 480 p	Parameter
CameraWrapper	fps	= 15 f	
IsentropicPrint	fig_size	≤ 50 p	
IsentropicPrint	isent_level	≤ 200 k	
WeatherSensors	data_size	≤ 30 pi	

Inferred low-level SLOs and parameter assignments for high-level SLOs

Results: RQ-1 SLO Diffusion

Constrain low-level SLOs



Composite BN for all microservices connected to the *VehicleRouting* consumer

TABLE III: SLOs and parameters inferred for *VehicleRouting*

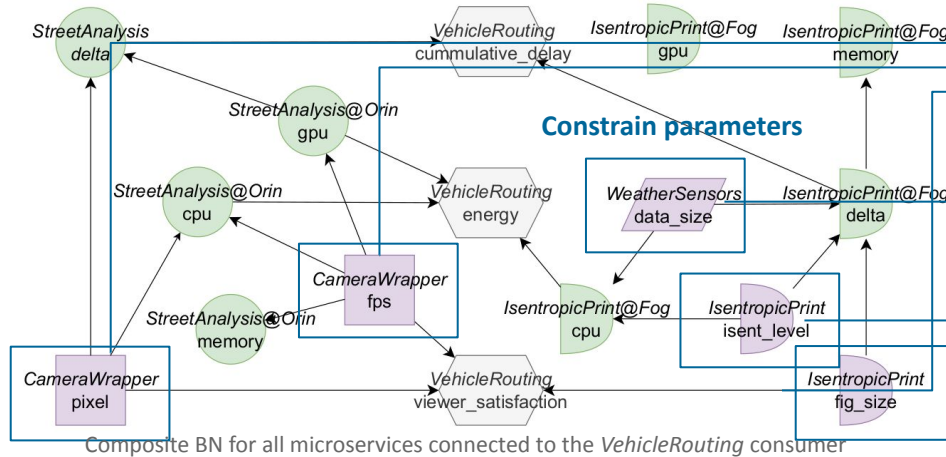
Microservice	Variable	States	SLO / Param
VehicleRouting	cumm_delay	≤ 45 ms	High-level
	energy	≤ 19 W	
	viewer_sat	—	
StreetAnalysis	delta	≤ 35 ms	Low-level
	cpu (Orin)	≤ 21 %	
	gpu (Orin)	≤ 40 %	
	delta	≤ 37 ms	
IsentropicPrint	cpu (Fog)	≤ 17 %	Parameter
	pixel	= 480 p	
CameraWrapper	fps	= 15 f	Parameter
IsentropicPrint	fig_size	≤ 50 p	
IsentropicPrint	isent_level	≤ 200 k	
WeatherSensors	data_size	≤ 30 pi	

Inferred low-level SLOs and parameter assignments for high-level SLOs

Results: RQ-1 SLO Diffusion

TABLE III: SLOs and parameters inferred for *VehicleRouting*

Microservice	Variable	States	SLO / Param
VehicleRouting	cumm_delay	≤ 45 ms	High-level
	energy	≤ 19 W	
	viewer_sat	—	
<i>StreetAnalysis</i>	delta	≤ 35 ms	Low-level
<i>StreetAnalysis</i>	cpu (Orin)	≤ 21 %	
<i>StreetAnalysis</i>	gpu (Orin)	≤ 40 %	
<i>IseotropicPrint</i>	delta	≤ 37 ms	
<i>IseotropicPrint</i>	cpu (Fog)	≤ 17 %	Parameter
<i>CameraWrapper</i>	pixel	= 480 p	
<i>CameraWrapper</i>	fps	= 15 f	
<i>IseotropicPrint</i>	fig_size	≤ 50 p	
<i>IseotropicPrint</i>	isent_level	≤ 200 k	
<i>WeatherSensors</i>	data_size	≤ 30 pi	



Inferred low-level SLOs and parameter assignments for high-level SLOs

Results: RQ-1 SLO Diffusion (2)

Setup the evaluation environment

- 1) **Deploy** microservices over the architecture
- 2) Set **parameter** according inferred thresholds
- 3) Measure the **actual** SLO fulfillment
- 4) Evaluate **alternative** param. configurations
- 5) **Compare** results and find min / max

What can we report?

- Inferred configuration close to optimal
- Discrepancy occur either due to **flexible boundaries**, or conflicts between SLOs

TABLE III: SLOs and parameters inferred for *VehicleRouting*

Microservice	Variable	States	SLO / Param
VehicleRouting	cumm_delay	≤ 45 ms	High-level
	energy	≤ 19 W	
	viewer_sat	—	
<i>StreetAnalysis</i>	delta	≤ 35 ms	Low-level
<i>StreetAnalysis</i>	cpu (Orin)	≤ 21 %	
<i>StreetAnalysis</i>	gpu (Orin)	≤ 40 %	
<i>IsentropicPrint</i>	delta	≤ 37 ms	
<i>IsentropicPrint</i>	cpu (Fog)	≤ 17 %	
<i>CameraWrapper</i>	pixel	= 480 p	Parameter
<i>CameraWrapper</i>	fps	= 15 f	
<i>IsentropicPrint</i>	fig_size	≤ 50 p	
<i>IsentropicPrint</i>	isent_level	≤ 200 k	
<i>WeatherSensors</i>	data_size	≤ 30 pi	

TABLE IV: High-level SLO fulfillment of inferred and alternative assignment for all three evaluated applications

Microservice	High-level SLO	% Min	% Fulfill	% Max
VehicleRouting	cumm_delay ≤ 45	0.00	0.94	1.00
	min(energy)	0.53	0.99	1.00
TrafficPrediction	cumm_delay ≤ 40	0.00	0.83	0.90
LiveMonitoring	cumm_delay ≤ 110	0.13	0.93	1.00
	max(viewer_sat.)	0.00	1.00	1.00

Results: RQ-1 SLO Diffusion (2)

Setup the evaluation environment

- 1) **Deploy** microservices over the architecture
- 2) Set **parameter** according inferred thresholds
- 3) Measure the **actual** SLO fulfillment
- 4) Evaluate **alternative** param. configurations
- 5) **Compare** results and find min / max

What can we report?

- Inferred configuration close to optimal
- Discrepancy occur either due to **flexible boundaries**, or conflicts between SLOs

TABLE III: SLOs and parameters inferred for *VehicleRouting*

Microservice	Variable	States	SLO / Param
VehicleRouting	cumm_delay	≤ 45 ms	High-level
	energy	≤ 19 W	
	viewer_sat		
<i>StreetAnalysis</i>	delta	≤ 35 ms	Low-level
<i>StreetAnalysis</i>	cpu (Orin)	≤ 21 %	
<i>StreetAnalysis</i>	gpu (Orin)	≤ 40 %	
<i>IsentropicPrint</i>	delta	≤ 37 ms	
<i>IsentropicPrint</i>	cpu (Fog)	≤ 17 %	
<i>CameraWrapper</i>	pixel	= 480 p	Parameter
<i>CameraWrapper</i>	fps	= 15 f	
<i>IsentropicPrint</i>	fig_size	≤ 50 p	
<i>IsentropicPrint</i>	isent_level	≤ 200 k	
<i>WeatherSensors</i>	data_size	≤ 30 pi	

TABLE IV: High-level SLO fulfillment of inferred and alternative assignment for all three evaluated applications

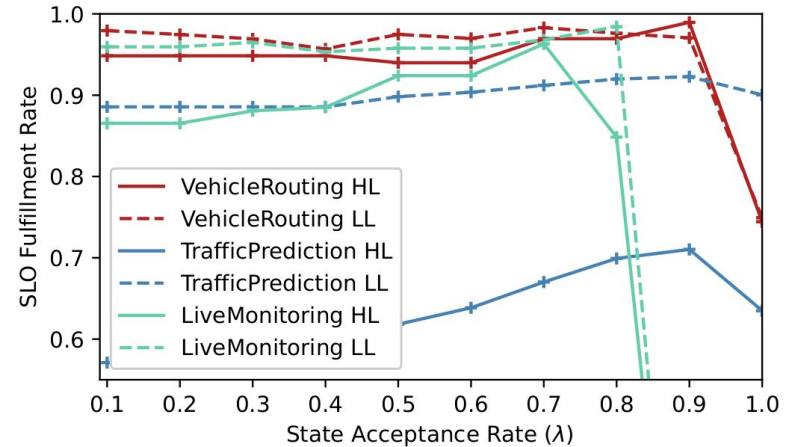
Microservice	High-level SLO	% Min	% Fulfill	% Max
VehicleRouting	cumm_delay ≤ 45	0.00	0.94	1.00
	min(energy)	0.53	0.99	1.00
TrafficPrediction	cumm_delay ≤ 40	0.00	0.83	0.90
LiveMonitoring	cumm_delay ≤ 110	0.13	0.93	1.00
	max(viewer_sat.)	0.00	1.00	1.00

Results: RQ-2 Acceptance Range

When **constraining** lower-level states of variables, how **restrictive** should lower-level boundaries be?

Vary acceptance rate (λ) from [0.1,1.0]; means very **loose** or **restrictive** for lower-level SLOs and params; for 3 consumers evaluate SLOs for microservices

Lower acceptance rate **improves** SLO fulfillment; when acceptance rate **too narrow**, not possible anymore to find satisfying parameter assignments



Results: RQ-3 Conflicts and Resolution

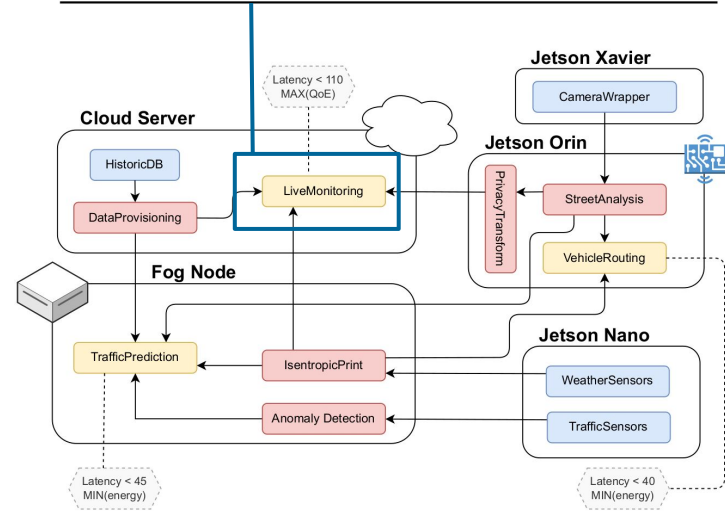
Conflicts between high-level SLOs appear at the lower-level SLOs and parameter assignments

Choose microservices around *LiveMonitoring* and specify increasingly tight performance boundary, extended with SLOs for opposing targets, i.e., **min**(energy) and/or **max**(satisfact.)

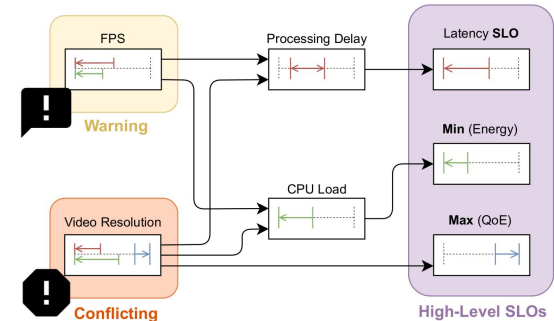
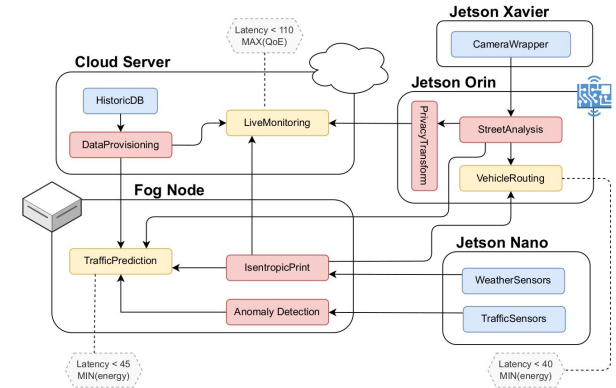
Wide boundaries (first row) allow to infer parameter assignments for most cases, but tighter constraints results in most variables finding no intersection

TABLE V: Conflicts among high-level SLOs for *LiveMonitor*

cumm_delay	min(energy)	max(customer_sat)	both
≤ 120 ms	✓	✓	$\{fps\}$
≤ 100 ms	$\{pixel\}$	✓	$\{ \wedge \cup pixel \}$
≤ 50 ms	$\{ \wedge \cup fps \}$	$\{gpu, pixel, fps\}$	$\{ \wedge \cup gpu \}$
≤ 40 ms	$\{ \wedge \cup batch \}$	$\{ \wedge \}$	$\{ \wedge \cup fig_size \}$
≤ 25 ms	$\{ \wedge \}$	$\{ \wedge \cup cache_db \}$	$\{ \wedge \cup cache_db \}$



- Stakeholders **unaware** of implications of high-level SLOs to lower-level components and parameters and whether specifies SLOs are **conflicting**
- Requires mechanisms to diffuse high-level SLOs to composite microservices pipelines
- Diffusion traverses BN to constrain lower-level variables to states that fulfill higher-level goals
- Evaluation for microservice networks; high-level goals were diffused to service parameters; could highlight conflicting SLOs to stakeholders



Let's **discuss!**

Please come forward with any **question** you have



@

boris.sedlak@dsg.tuwien.ac.at

