

# Network-Scale Maintenance Planning for Infrastructures Using Reinforcement Learning (Formulation, Part 2)

Zachary Hamida Postdoc  
James-A. Goulet Professeur

 Polytechnique Montréal, Canada  
Département des génies civil, géologique et des mines

August 6, 2021

Funding:  
Transportation Ministry of Quebec (MTQ)

# 🗺️ Outline

---

## Recap & Definitions

## Reinforcement Learning - Context

## Problem Formulation

---

# Recap & Definitions

Visual Inspections (VI): Network-scale monitoring technique

Source: Google images

# Recap & Definitions

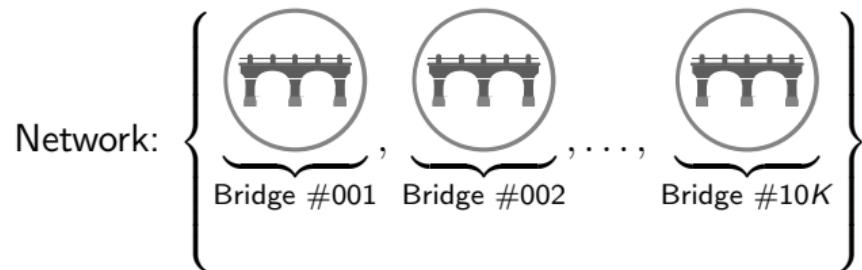
Visual Inspections (VI): Network-scale monitoring technique

Network:

Source: Google images

# Recap & Definitions

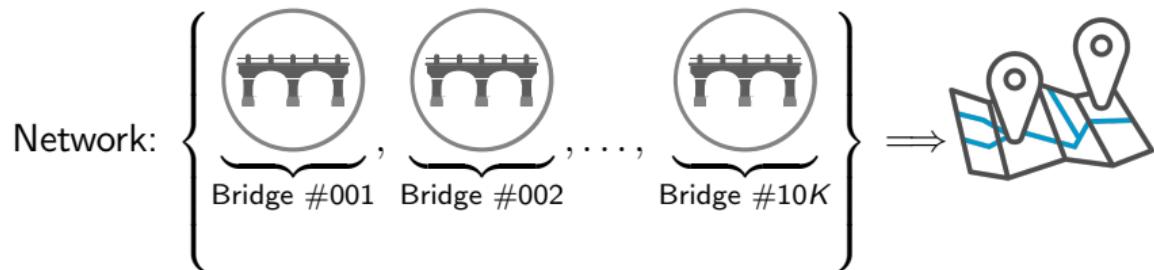
Visual Inspections (VI): Network-scale monitoring technique



Source: Google images

# Recap & Definitions

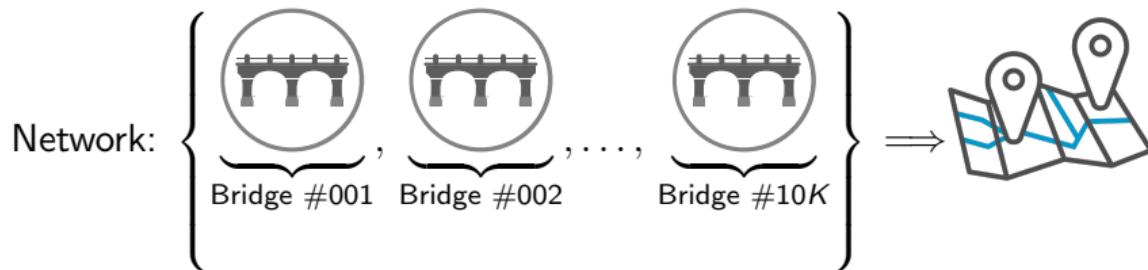
Visual Inspections (VI): Network-scale monitoring technique



Source: Google images

# Recap & Definitions

Visual Inspections (VI): Network-scale monitoring technique

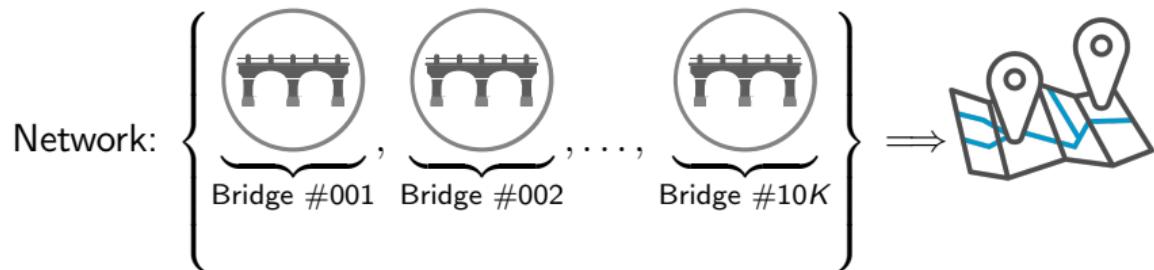


Bridge: {type, material, traffic, location, ... }

Source: Google images

# Recap & Definitions

Visual Inspections (VI): Network-scale monitoring technique



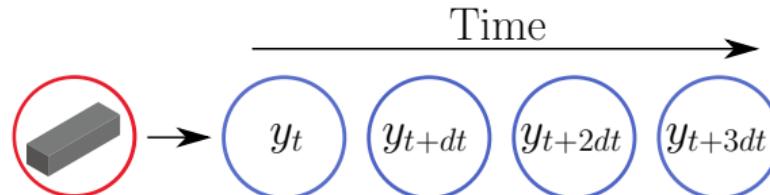
Bridge: {type, material, traffic, location, ... }

Source: Google images

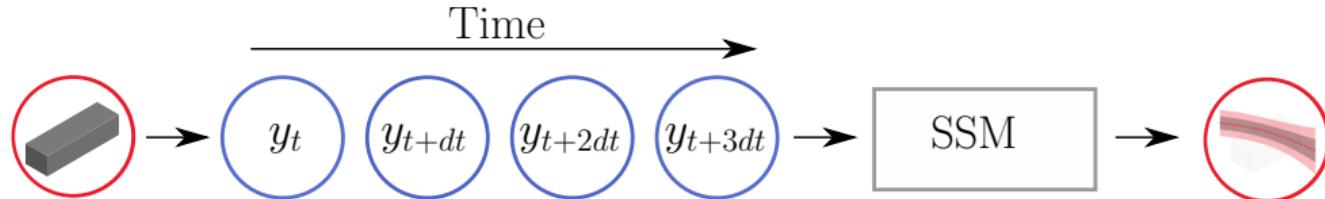
# Deterioration Framework



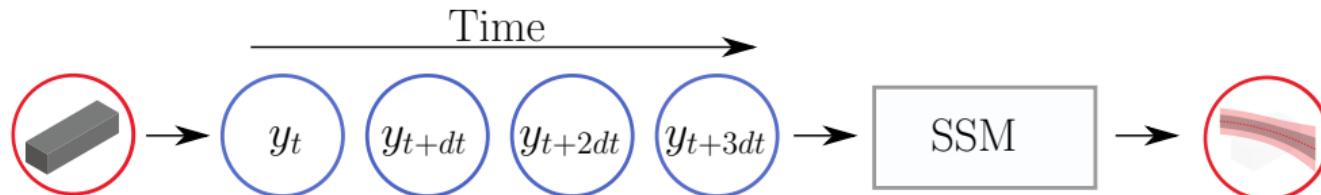
# Deterioration Framework



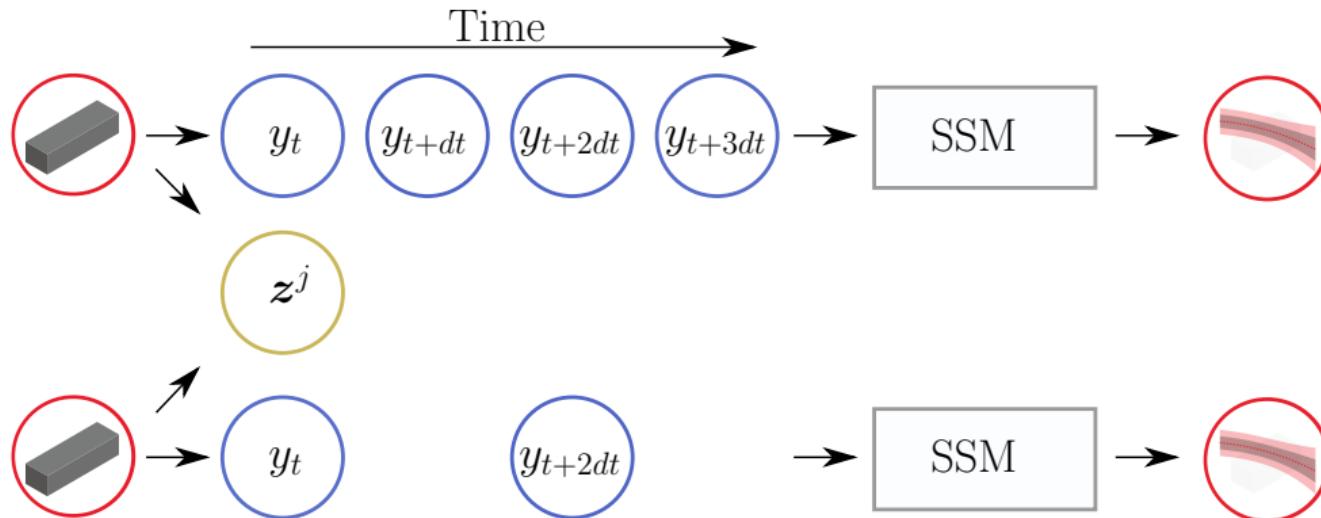
# Deterioration Framework



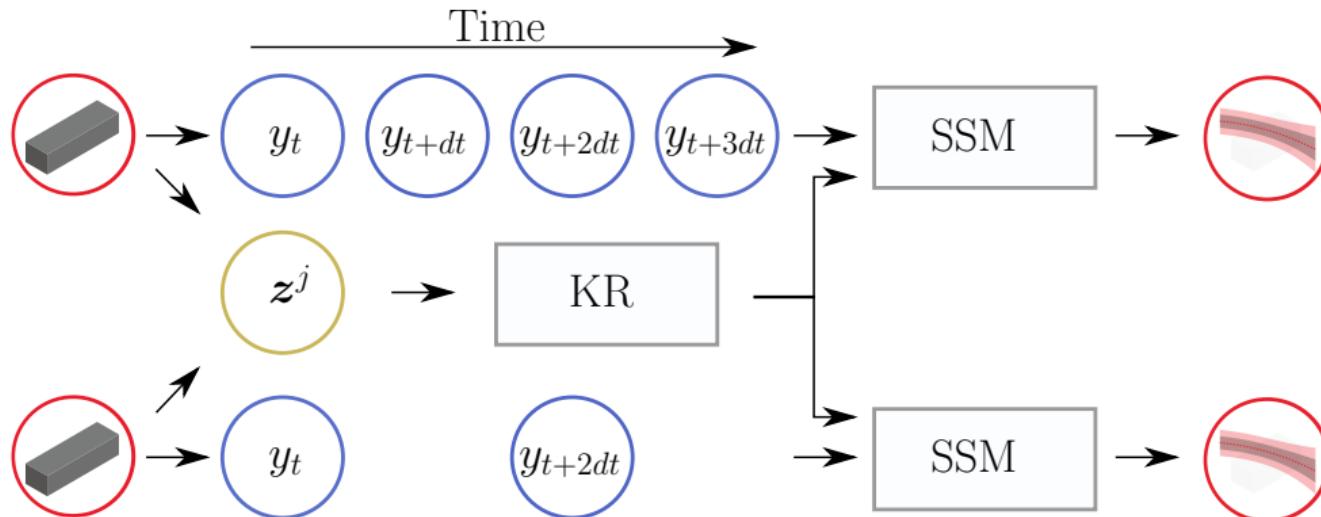
# Deterioration Framework



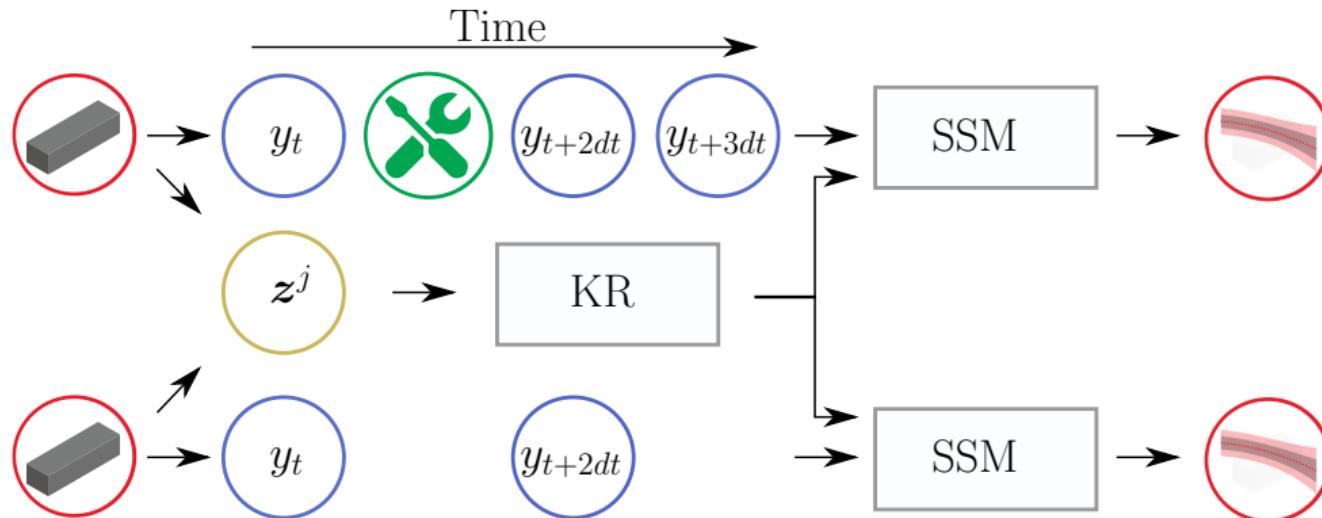
# Deterioration Framework



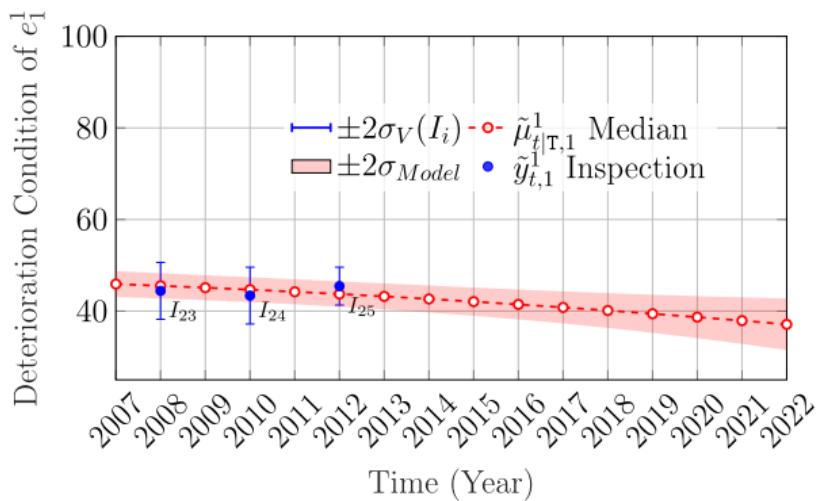
# Deterioration Framework



# Deterioration Framework

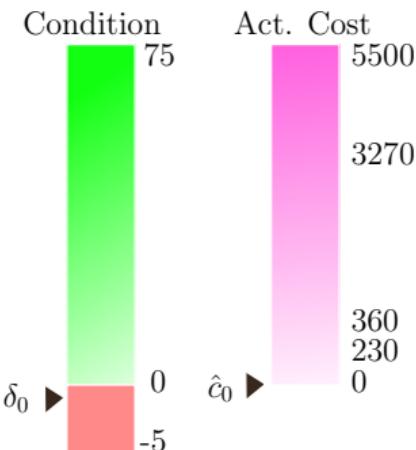
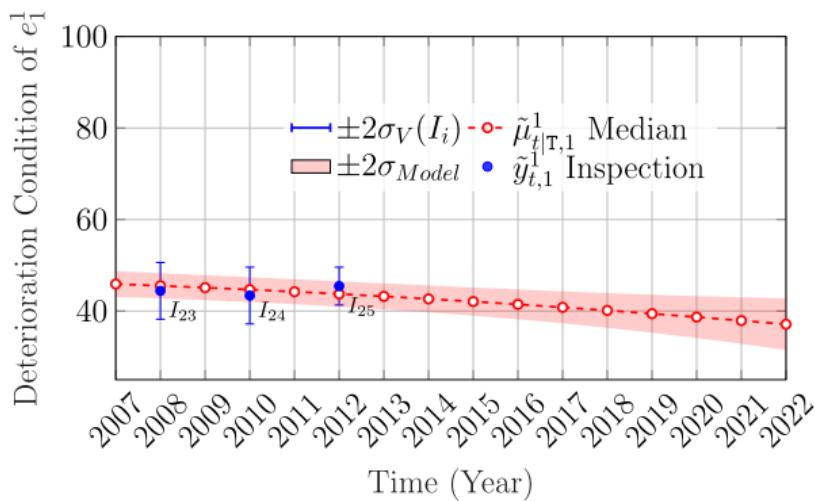


# Element-Level Analyses



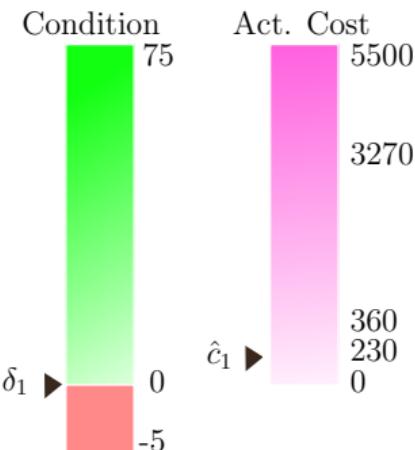
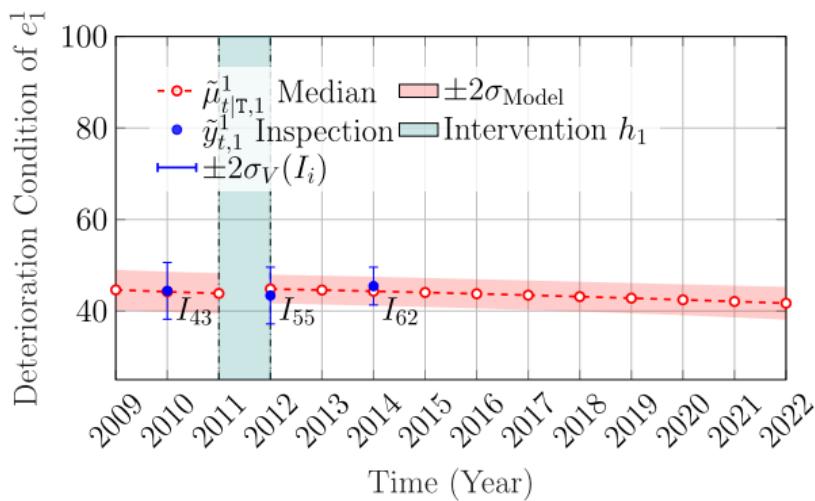
Source: MTQ, Manual of Inspection

# Element-Level Analyses



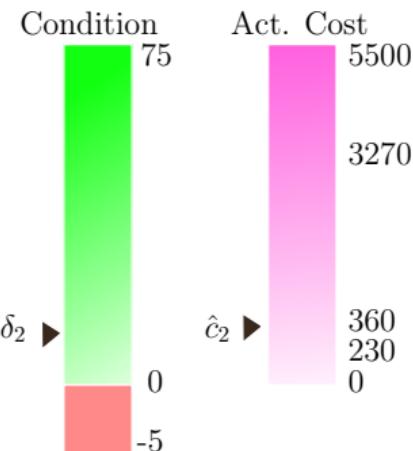
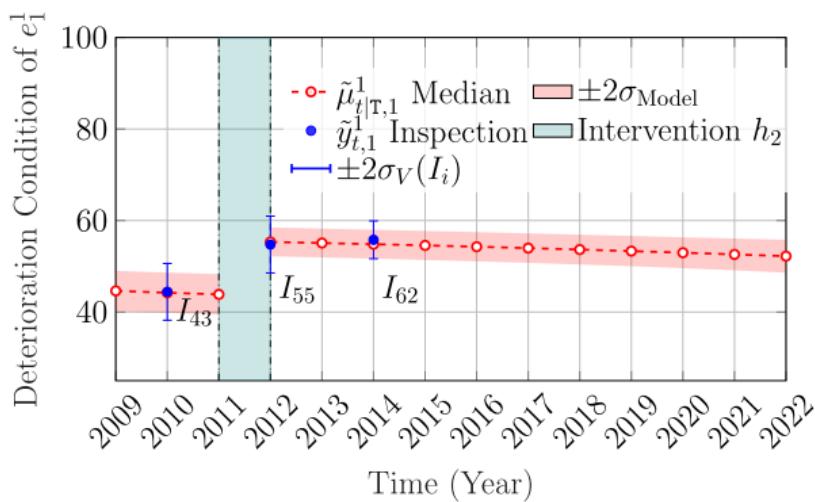
Source: MTQ, Manual of Inspection

# Element-Level Analyses



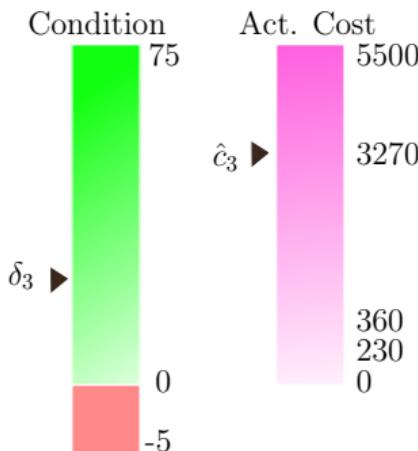
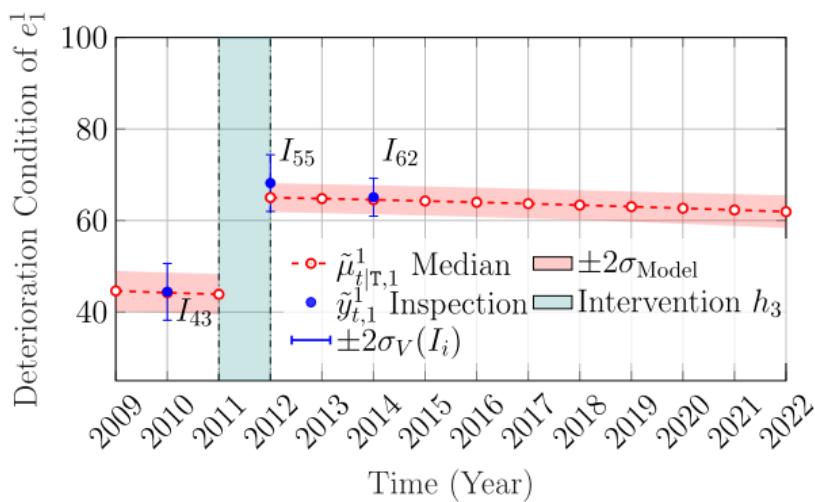
Source: MTQ, Manual of Inspection

# Element-Level Analyses



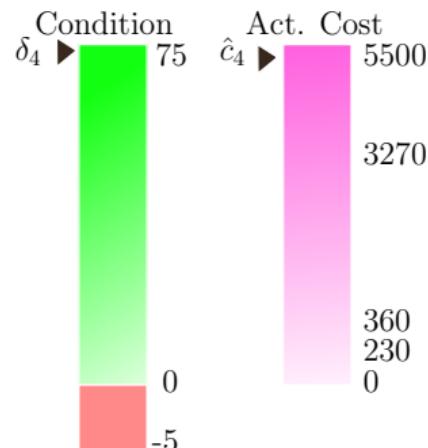
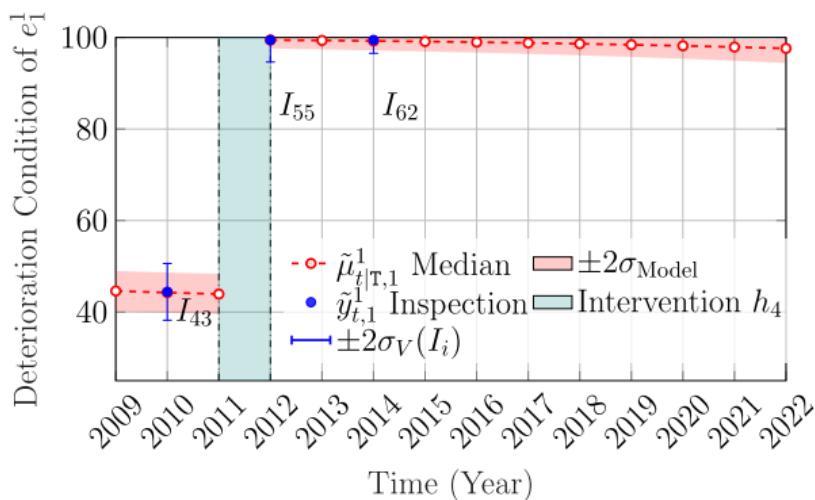
Source: MTQ, Manual of Inspection

# Element-Level Analyses



Source: MTQ, Manual of Inspection

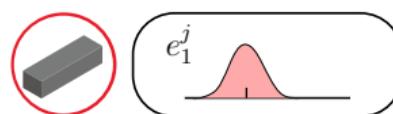
# Element-Level Analyses



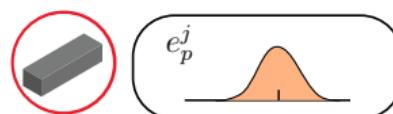
Source: MTQ, Manual of Inspection

# System-Level Analyses

At any time  $t$ :

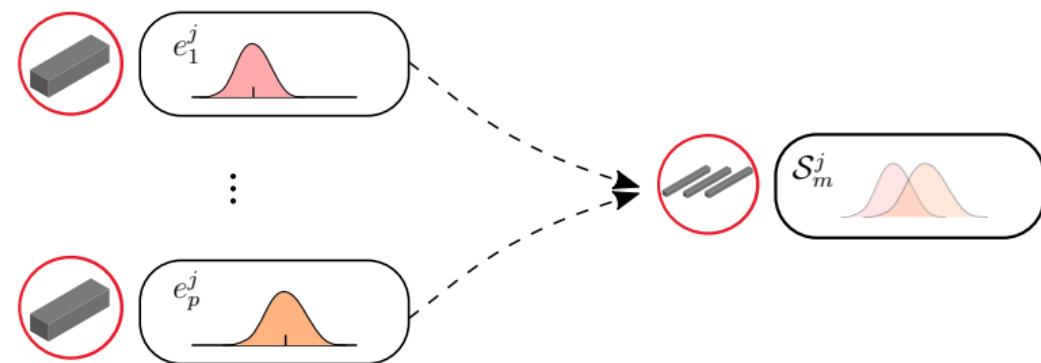


⋮



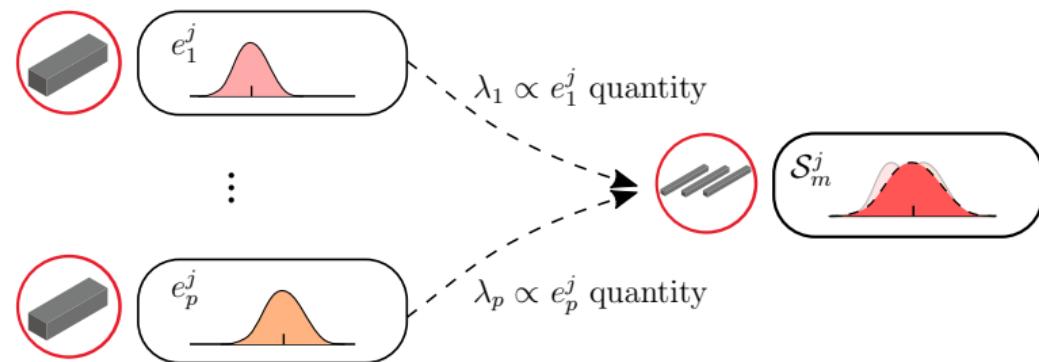
# System-Level Analyses

At any time  $t$ :

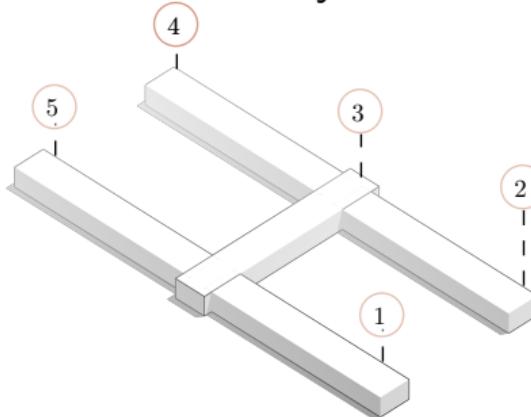


# System-Level Analyses

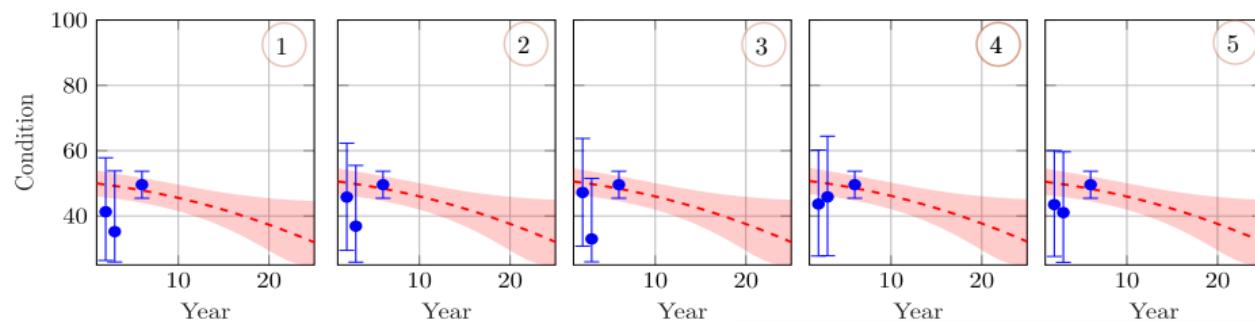
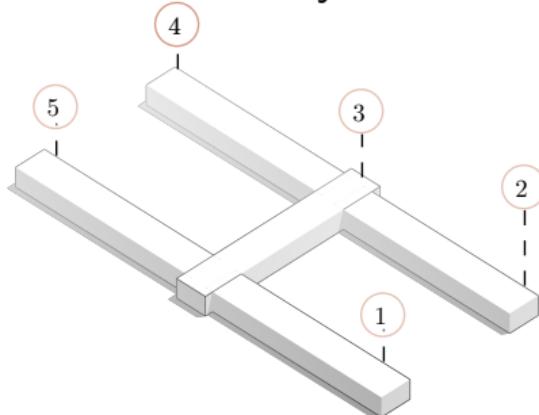
At any time  $t$ :



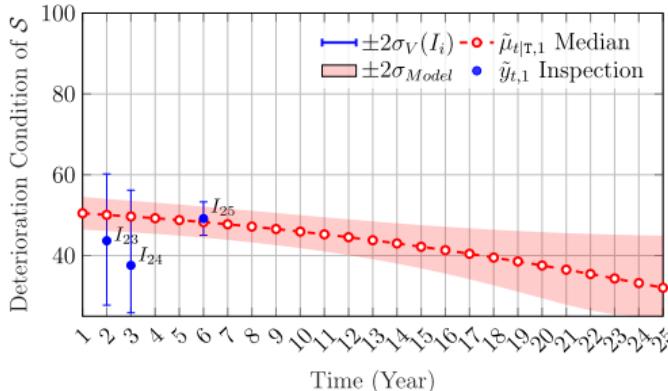
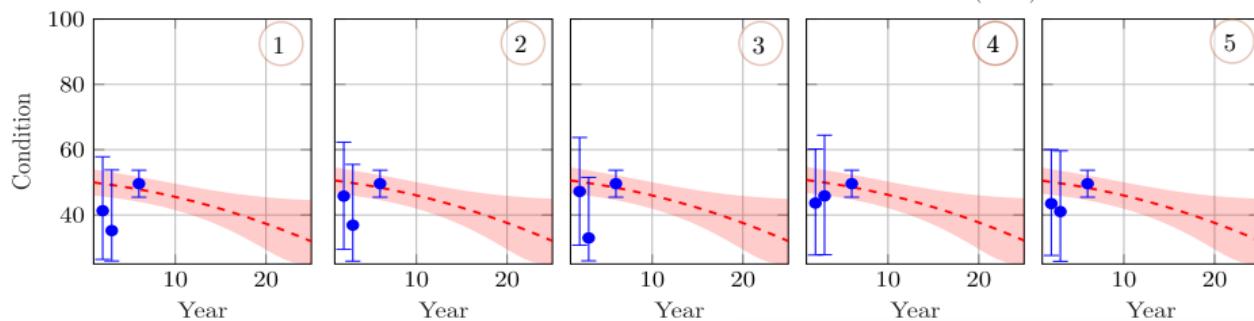
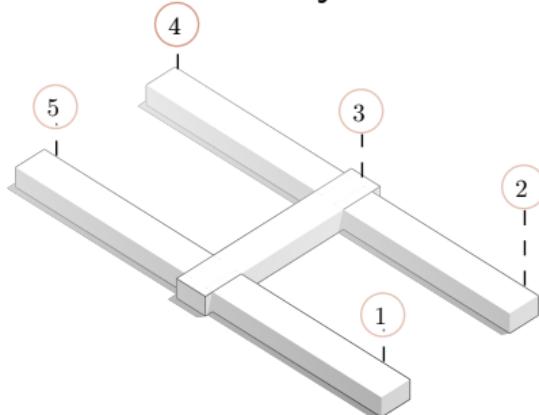
# System-Level Analyses



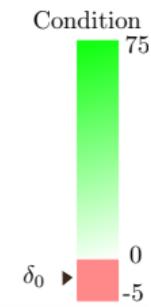
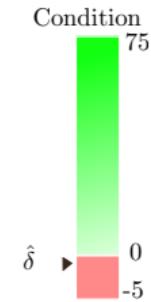
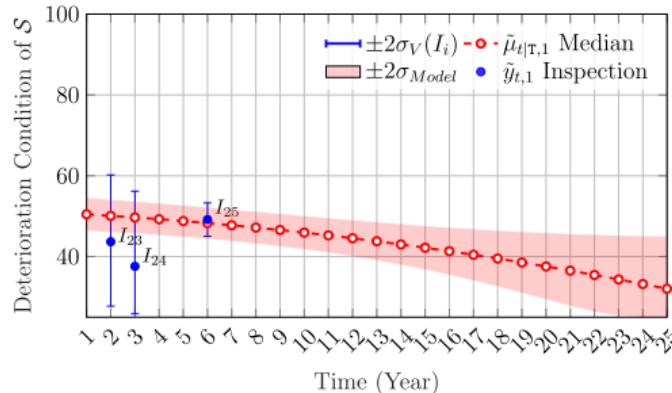
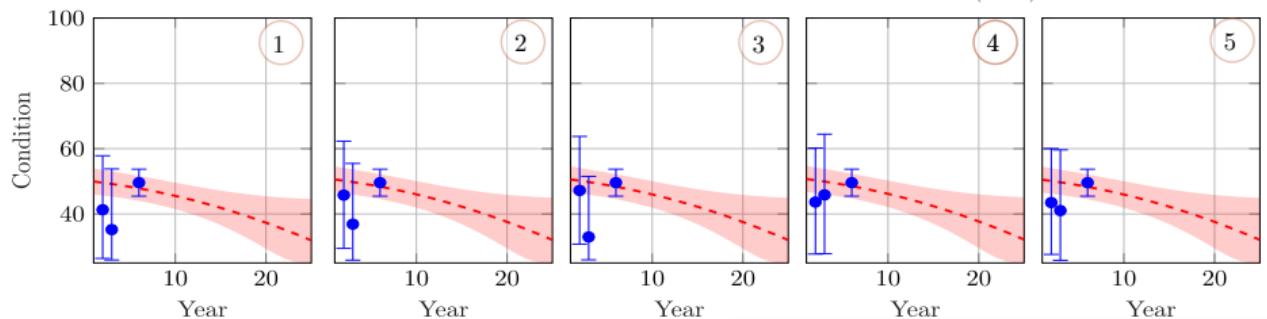
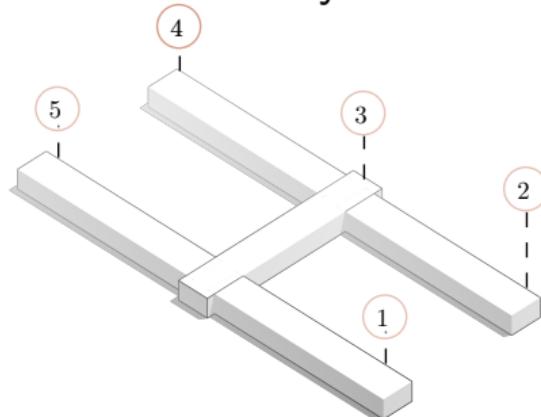
# System-Level Analyses



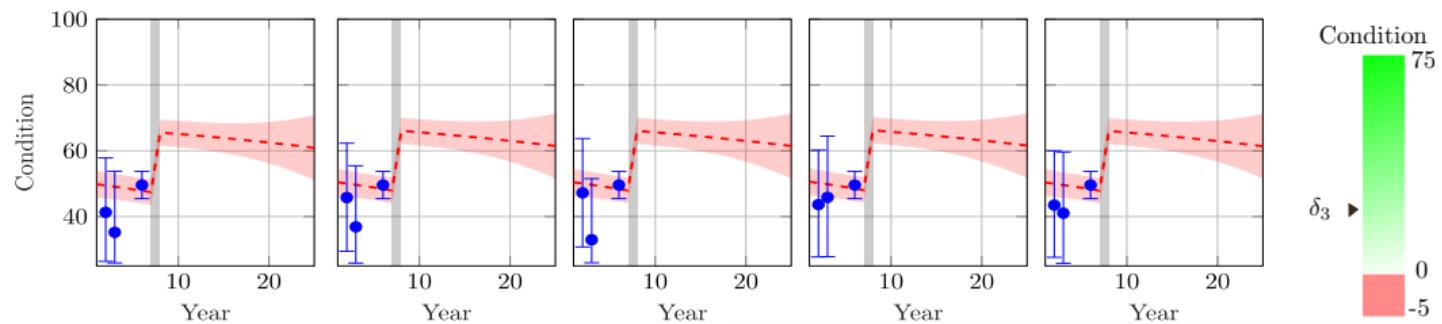
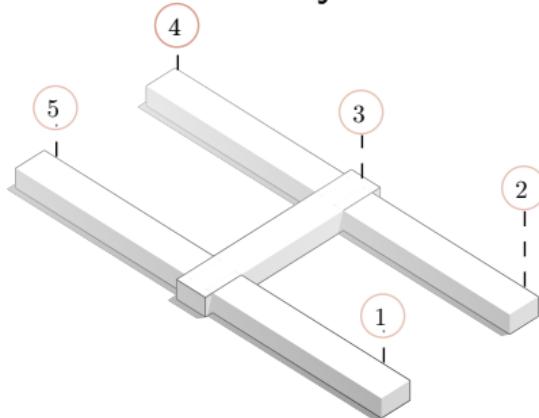
# System-Level Analyses



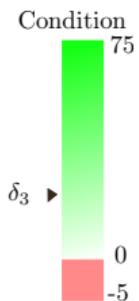
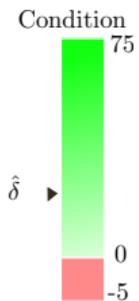
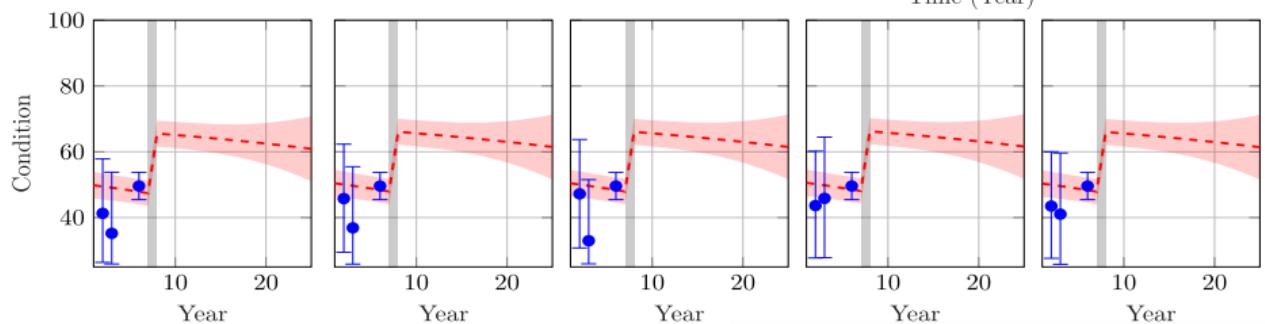
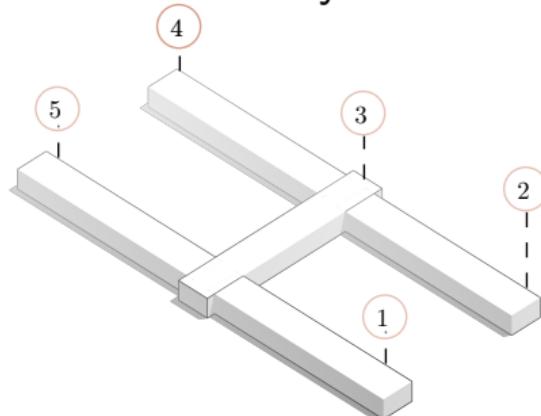
# System-Level Analyses



# System-Level Analyses

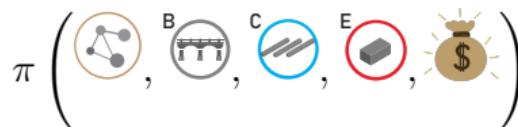


# System-Level Analyses



# Research Objective

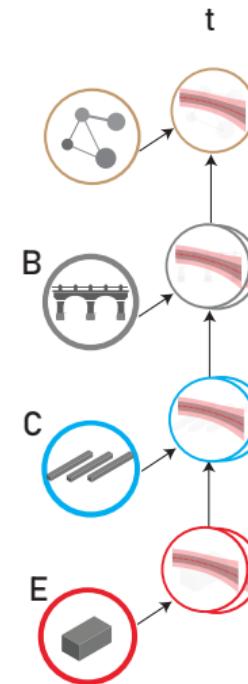
Find a policy for maintenance:

$$\pi \left( \text{A} \left( \text{B}, \text{C}, \text{D}, \text{E} \right) \right)$$


# Research Objective

Find a policy for maintenance:

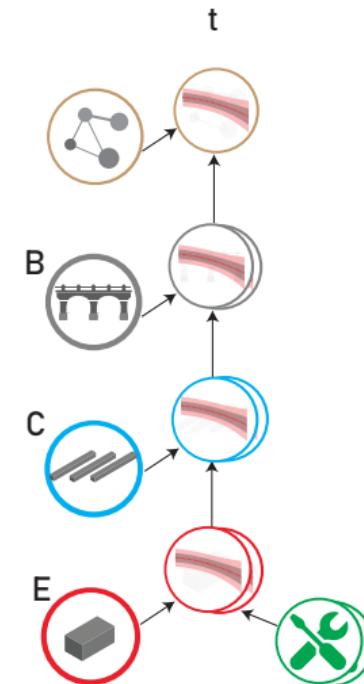
$$\pi \left( \text{A} \left( \text{graph icon} \right), \text{B} \left( \text{bridge icon} \right), \text{C} \left( \text{wires icon} \right), \text{E} \left( \text{brick icon} \right), \text{D} \left( \text{money bag icon} \right) \right)$$



# Research Objective

Find a policy for maintenance:

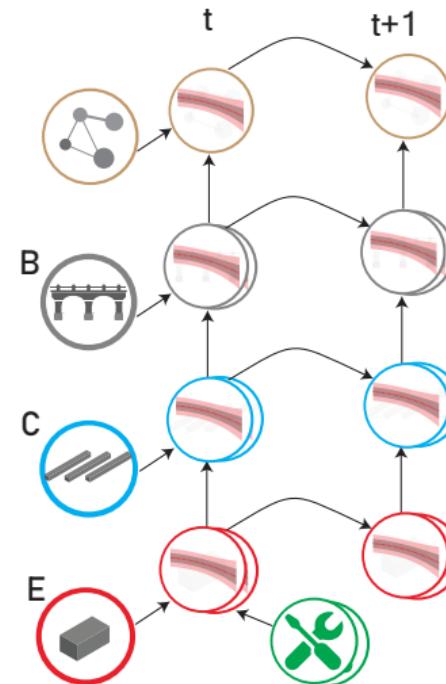
$$\pi \left( \text{A} \left( \text{graph icon} \right), \text{B} \left( \text{bridge icon} \right), \text{C} \left( \text{wires icon} \right), \text{E} \left( \text{brick icon} \right), \text{F} \left( \text{money bag icon} \right) \right)$$



# Research Objective

Find a policy for maintenance:

$$\pi \left( \text{A} \left( \text{graph icon} \right), \text{B} \left( \text{bridge icon} \right), \text{C} \left( \text{wires icon} \right), \text{E} \left( \text{brick icon} \right), \text{D} \left( \text{money bag icon} \right) \right)$$



# Reinforcement Learning

Source: Google Images

# Reinforcement Learning

▷ Agent: decision maker.

Source: Google Images

# Reinforcement Learning

- ▷ Agent: decision maker.
- ▷ Environment: e.g., video game.

Source: Google Images

# Reinforcement Learning

- ▷ Agent: decision maker.
- ▷ Environment: e.g., video game.
- ▷ Actions: e.g., key press or movement.

Source: Google Images

# Reinforcement Learning

- ▷ Agent: decision maker.
- ▷ Environment: e.g., video game.
- ▷ Actions: e.g., key press or movement.
- ▷ Rewards: e.g., achieving the goal.

Source: Google Images

# Reinforcement Learning

- ▷ Agent: decision maker.
- ▷ Environment: e.g., video game.
- ▷ Actions: e.g., key press or movement.
- ▷ Rewards: e.g., achieving the goal.

**Learn a policy that maximizes the total expected discounted rewards.**

Source: Google Images

# Environment

Employing the deterioration model as an environment requires:

# Environment

Employing the deterioration model as an environment requires:

- ▷ Adjustments for the actions' exploration:

# Environment

Employing the deterioration model as an environment requires:

- ▷ Adjustments for the actions' exploration:
  - ✓ Account for the boundaries of the feasible condition.

# Environment

Employing the deterioration model as an environment requires:

- ▷ Adjustments for the actions' exploration:
  - ✓ Account for the boundaries of the feasible condition.
  - ✓ Account for repeating actions.

# Environment

Employing the deterioration model as an environment requires:

- ▷ Adjustments for the actions' exploration:
  - ✓ Account for the boundaries of the feasible condition.
  - ✓ Account for repeating actions.
- ▷ Defining measures of success:

# Environment

Employing the deterioration model as an environment requires:

- ▷ Adjustments for the actions' exploration:
  - ✓ Account for the boundaries of the feasible condition.
  - ✓ Account for repeating actions.
- ▷ Defining measures of success:
  - ✓ Compatibility (Elements and Category).

# Environment

Employing the deterioration model as an environment requires:

- ▷ Adjustments for the actions' exploration:
  - ✓ Account for the boundaries of the feasible condition.
  - ✓ Account for repeating actions.
- ▷ Defining measures of success:
  - ✓ Compatibility (Elements and Category).
  - ✓ Criticality (Elements and Category).

# Environment

Employing the deterioration model as an environment requires:

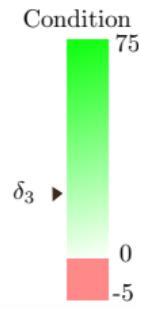
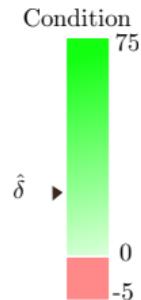
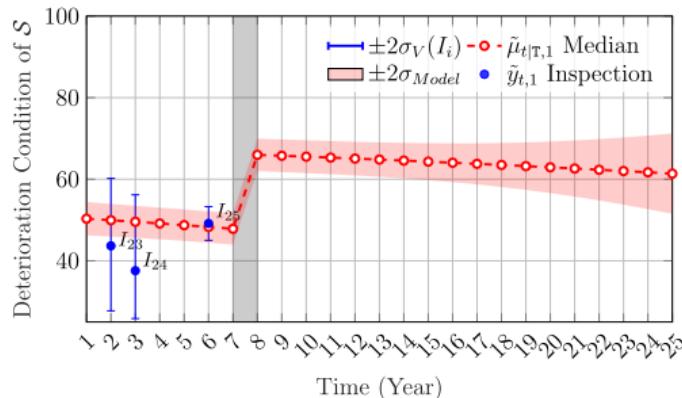
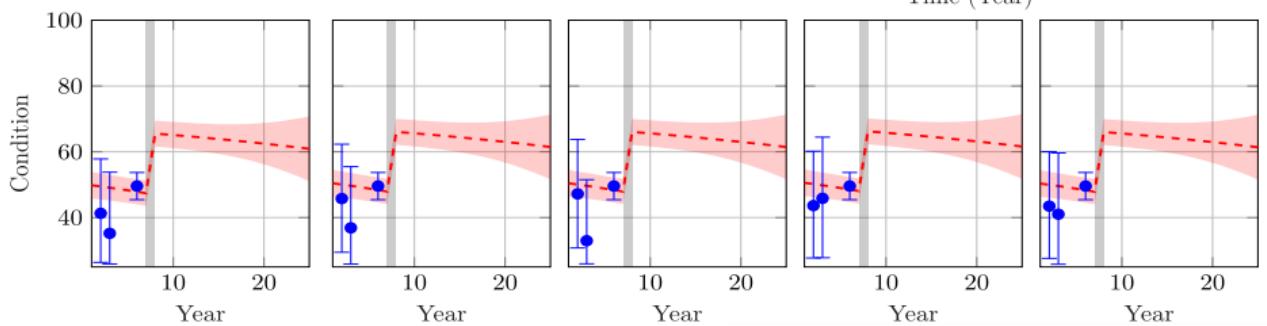
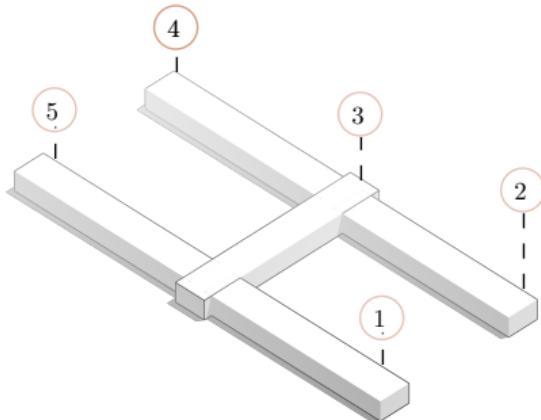
- ▷ Adjustments for the actions' exploration:
  - ✓ Account for the boundaries of the feasible condition.
  - ✓ Account for repeating actions.
- ▷ Defining measures of success:
  - ✓ Compatibility (Elements and Category).
  - ✓ Criticality (Elements and Category).
  - ✓ Variability (Elements).

# Environment

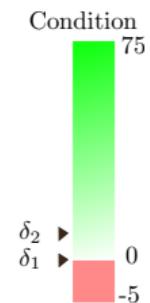
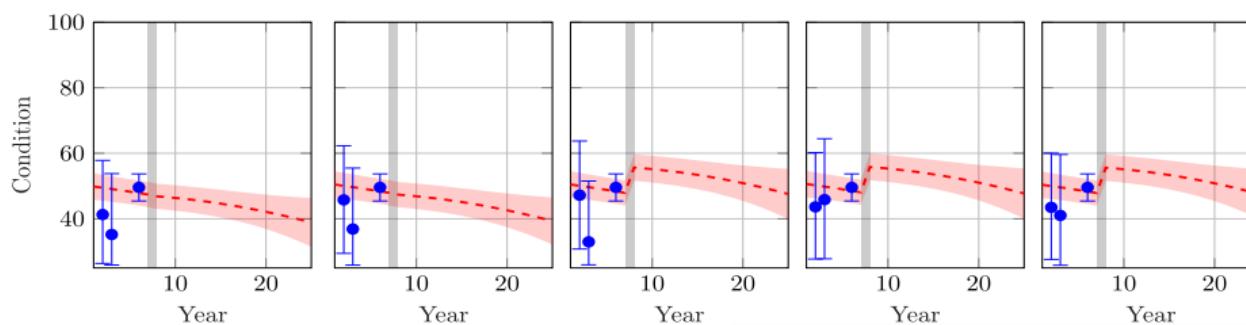
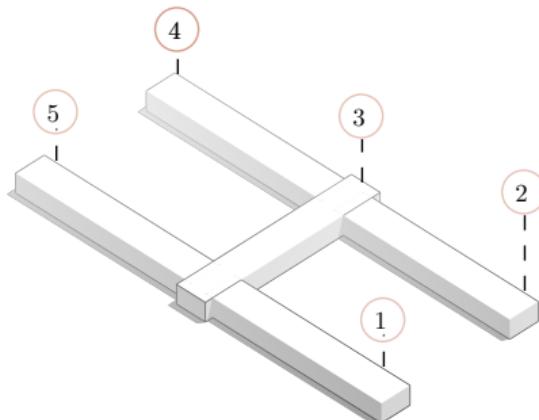
Employing the deterioration model as an environment requires:

- ▷ Adjustments for the actions' exploration:
  - ✓ Account for the boundaries of the feasible condition.
  - ✓ Account for repeating actions.
- ▷ Defining measures of success:
  - ✓ Compatibility (Elements and Category).
  - ✓ Criticality (Elements and Category).
  - ✓ Variability (Elements).
  - ✓ Action frequency (Category).

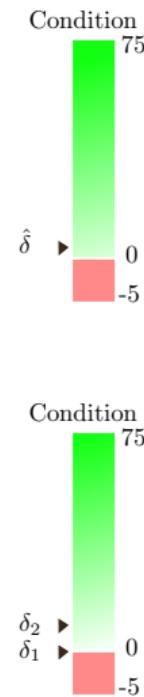
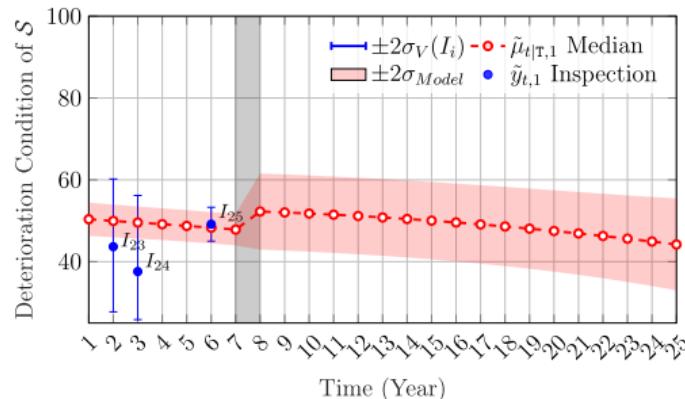
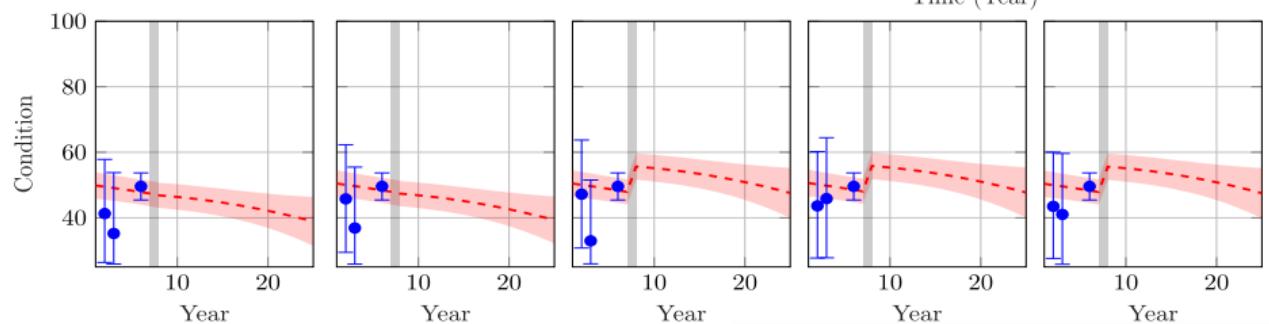
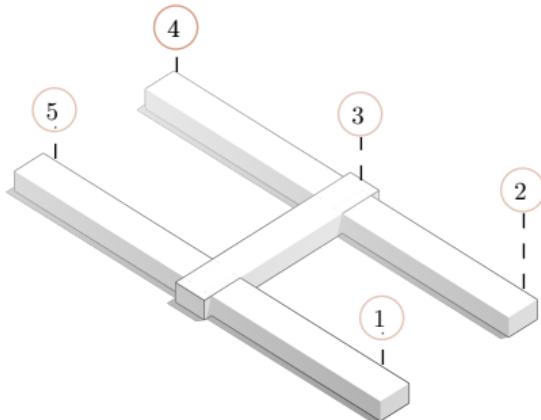
## Hierarchical Reinforcement Learning (HRL)



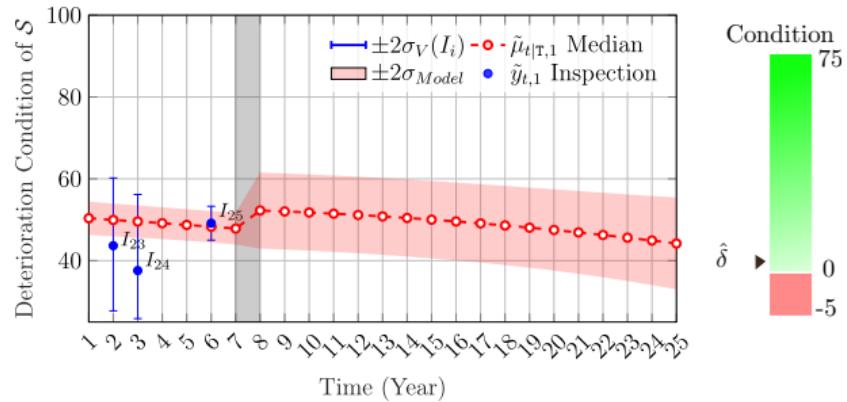
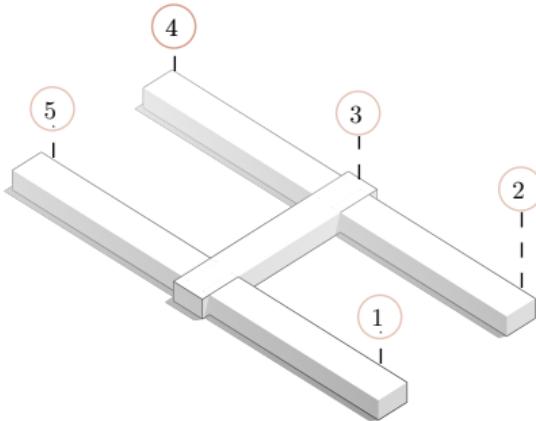
## Hierarchical Reinforcement Learning (HRL)



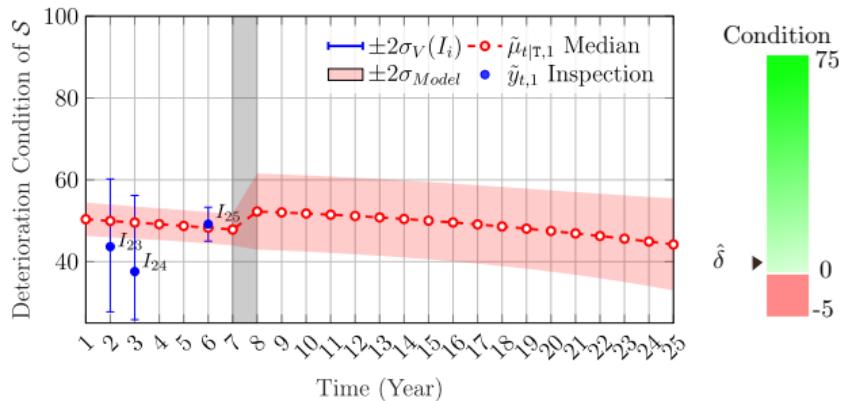
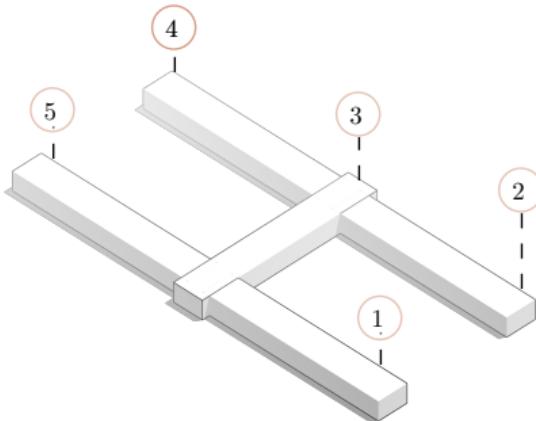
## Hierarchical Reinforcement Learning (HRL)



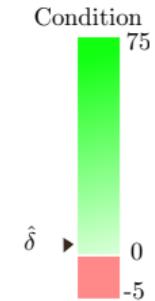
## Hierarchical Reinforcement Learning (HRL)



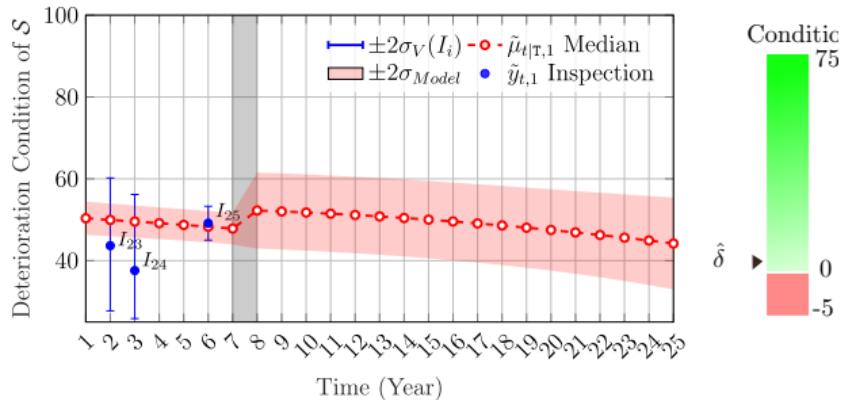
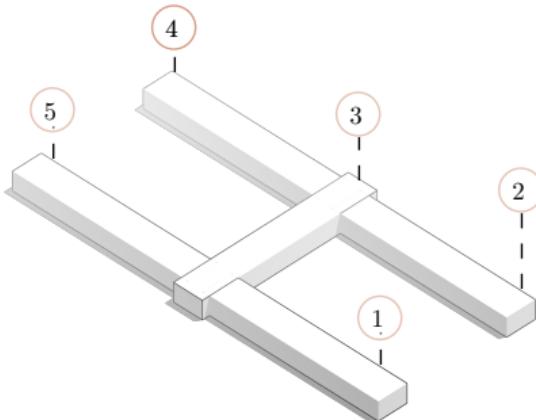
## Hierarchical Reinforcement Learning (HRL)



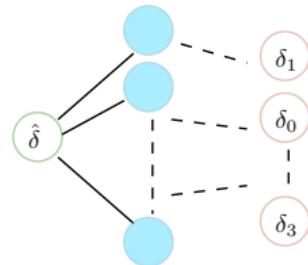
$$L = \|\mathbb{E}[\boldsymbol{\delta}] - \hat{\boldsymbol{\delta}}\|$$



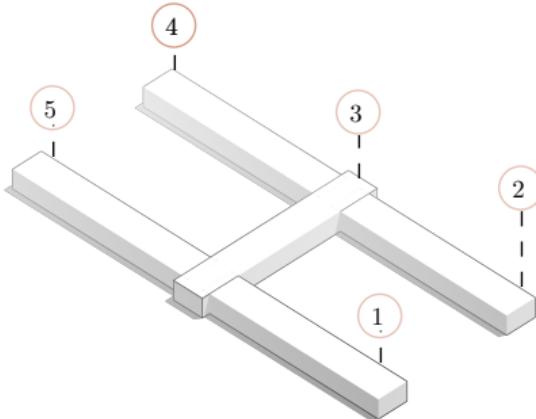
## Hierarchical Reinforcement Learning (HRL)



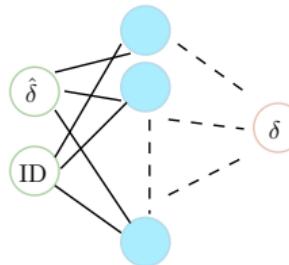
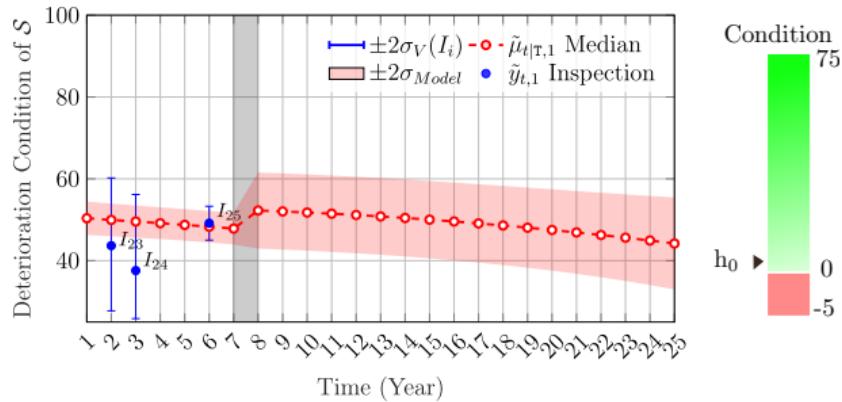
$$L = \|\mathbb{E}[\boldsymbol{\delta}] - \hat{\boldsymbol{\delta}}\|$$



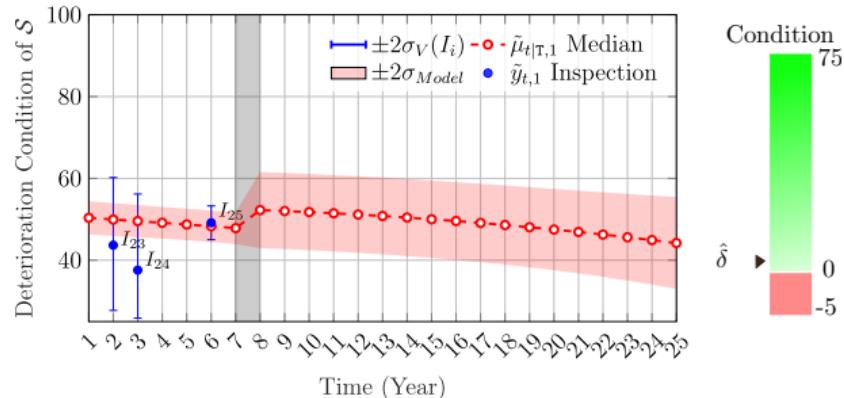
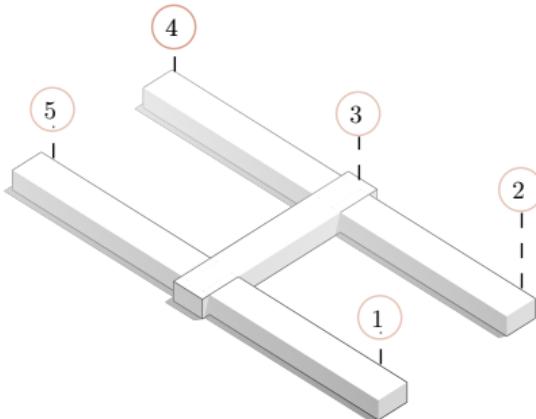
## Hierarchical Reinforcement Learning (HRL)



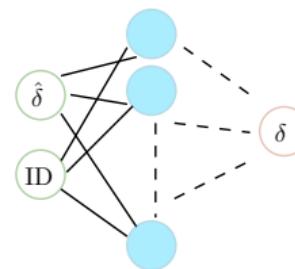
$$L = \left\| \frac{\delta_i}{E} - \hat{\delta} \right\|$$



## Hierarchical Reinforcement Learning (HRL)



$$L = \|\frac{\delta_i}{E} - \hat{\delta}\| \times \hat{c}_i$$



## RL Sub-task (Local Policy $\pi^e$ ):

## RL Sub-task (Local Policy $\pi^e$ ):

▷ state:  $\hat{\delta}_i$ , ID

## RL Sub-task (Local Policy $\pi^e$ ):

- ▷ state:  $\hat{\delta}_i$ , ID
- ▷ actions:  $a_0, a_1, a_2, a_3, a_4$  corresponding to the effects  $\delta_0, \delta_1, \delta_2, \delta_3, \delta_4$

## RL Sub-task (Local Policy $\pi^e$ ):

- ▷ state:  $\hat{\delta}_i$ , ID
- ▷ actions:  $a_0, a_1, a_2, a_3, a_4$  corresponding to the effects  $\delta_0, \delta_1, \delta_2, \delta_3, \delta_4$
- ▷ rewards:  $\hat{r}_i = -\|\hat{\delta} - \frac{\delta_i}{E}\| \times \hat{c}_i$ .

## RL Sub-task (Local Policy $\pi^e$ ):

- ▷ state:  $\hat{\delta}_i$ , ID
- ▷ actions:  $a_0, a_1, a_2, a_3, a_4$  corresponding to the effects  $\delta_0, \delta_1, \delta_2, \delta_3, \delta_4$
- ▷ rewards:  $\hat{r}_i = -\|\hat{\delta} - \frac{\delta_i}{E}\| \times \hat{c}_i$ .
- ▷ transition model:  $\hat{\delta}_{i+1} = \hat{\delta}_i - \frac{\delta_i}{E}$ .

## RL Sub-task (Local Policy $\pi^e$ ):

- ▷ state:  $\hat{\delta}_i$ , ID
- ▷ actions:  $a_0, a_1, a_2, a_3, a_4$  corresponding to the effects  $\delta_0, \delta_1, \delta_2, \delta_3, \delta_4$
- ▷ rewards:  $\hat{r}_i = -\|\hat{\delta} - \frac{\delta_i}{E}\| \times \hat{c}_i$ .
- ▷ transition model:  $\hat{\delta}_{i+1} = \hat{\delta}_i - \frac{\delta_i}{E}$ .

**Use Q-learning to find the policy  $\pi^e$  for the RL problem.**

## RL Sub-task (Local Policy $\pi^e$ ):

- ▷ state:  $\hat{\delta}_i$ , ID
- ▷ actions:  $a_0, a_1, a_2, a_3, a_4$  corresponding to the effects  $\delta_0, \delta_1, \delta_2, \delta_3, \delta_4$
- ▷ rewards:  $\hat{r}_i = -\|\hat{\delta} - \frac{\delta_i}{E}\| \times \hat{c}_i$ .
- ▷ transition model:  $\hat{\delta}_{i+1} = \hat{\delta}_i - \frac{\delta_i}{E}$ .

Use Q-learning to find the policy  $\pi^e$  for the RL problem.

But, what about the maintenance policy  $\pi$  ?

# RL Main Task (Global Policy $\pi$ ):

# RL Main Task (Global Policy $\pi$ ):

- ▷ state:  $\mu_{t|T}, \dot{\mu}_{t|T}$

# RL Main Task (Global Policy $\pi$ ):

- ▷ state:  $\mu_{t|T}, \dot{\mu}_{t|T}$
- ▷ actions:  $\hat{\delta}_t$  (continuous action)

## RL Main Task (Global Policy $\pi$ ):

- ▷ state:  $\mu_{t|\mathcal{T}}, \dot{\mu}_{t|\mathcal{T}}$
- ▷ actions:  $\hat{\delta}_t$  (continuous action)
- ▷ rewards:  $r_t = \sum_i^E \hat{c}_i.$

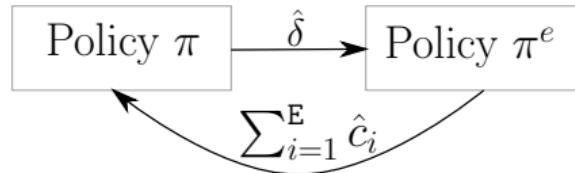
## RL Main Task (Global Policy $\pi$ ):

- ▷ state:  $\mu_{t|\mathcal{T}}, \dot{\mu}_{t|\mathcal{T}}$
- ▷ actions:  $\hat{\delta}_t$  (continuous action)
- ▷ rewards:  $r_t = \sum_i^E \hat{c}_i.$

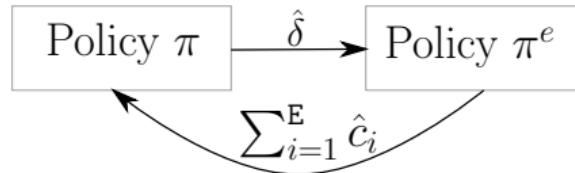
**Use Q-learning (with policy gradient) to solve the RL problem.**

## Simultaneous Training for $\pi^e$ and $\pi$ (HRL):

# Simultaneous Training for $\pi^e$ and $\pi$ (HRL):

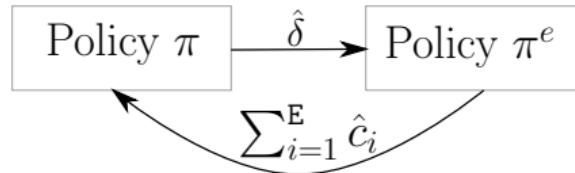


## Simultaneous Training for $\pi^e$ and $\pi$ (HRL):



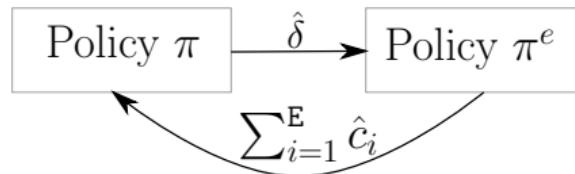
**Problem:**  $\|\mathbb{E}[\delta] - \hat{\delta}\| \gg 0$

## Simultaneous Training for $\pi^e$ and $\pi$ (HRL):



**Problem:**  $\|\mathbb{E}[\delta] - \hat{\delta}\| \gg 0 \rightarrow \sum_{i=1}^E \hat{c}_i$  and  $\hat{\delta}$  **do not** correspond to each other.

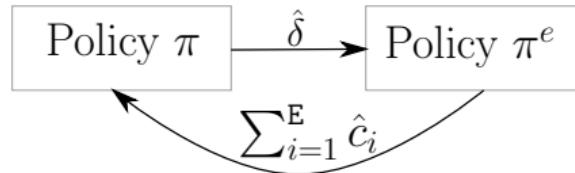
## Simultaneous Training for $\pi^e$ and $\pi$ (HRL):



**Problem:**  $\|\mathbb{E}[\delta] - \hat{\delta}\| \gg 0 \rightarrow \sum_{i=1}^E \hat{c}_i$  and  $\hat{\delta}$  **do not** correspond to each other.

**Solution:** replace  $\hat{\delta}$  with  $\bar{\delta}$  that maximizes the log probability  $\log \pi^e(a_i|\bar{\delta})$  as in:

## Simultaneous Training for $\pi^e$ and $\pi$ (HRL):

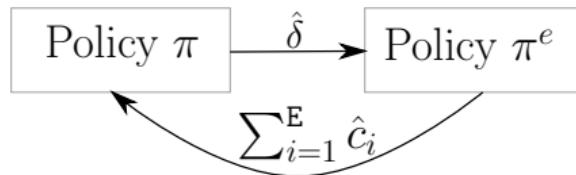


**Problem:**  $\|\mathbb{E}[\delta] - \hat{\delta}\| \gg 0 \rightarrow \sum_{i=1}^E \hat{c}_i$  and  $\hat{\delta}$  **do not** correspond to each other.

**Solution:** replace  $\hat{\delta}$  with  $\bar{\delta}$  that maximizes the log probability  $\log \pi^e(a_i|\bar{\delta})$  as in:

$$\log \pi^e(a_i|\bar{\delta}) \propto -\frac{1}{2} \sum_{i=1}^E \|a_i - \pi^e(\bar{\delta}_i, \text{ID})\|^2$$

## Simultaneous Training for $\pi^e$ and $\pi$ (HRL):



**Problem:**  $\|\mathbb{E}[\delta] - \hat{\delta}\| \gg 0 \rightarrow \sum_{i=1}^E \hat{c}_i$  and  $\hat{\delta}$  **do not** correspond to each other.

**Solution:** replace  $\hat{\delta}$  with  $\bar{\delta}$  that maximizes the log probability  $\log \pi^e(a_i|\bar{\delta})$  as in:

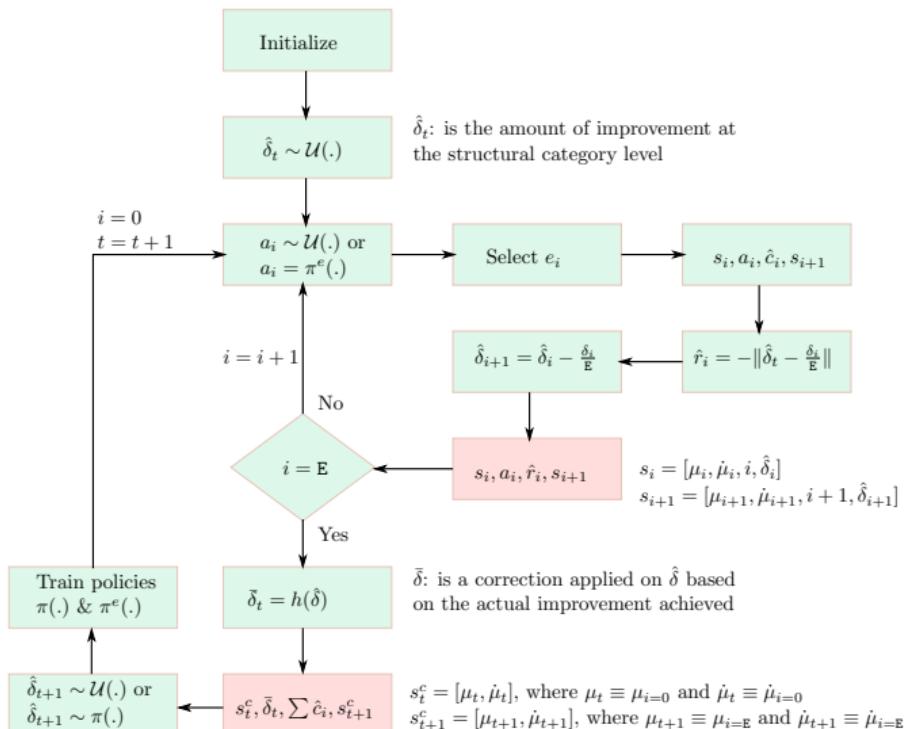
$$\log \pi^e(a_i|\bar{\delta}) \propto -\frac{1}{2} \sum_{i=1}^E \|a_i - \pi^e(\bar{\delta}_i, \text{ID})\|^2$$

To maximize this quantity, candidate goals  $\bar{\delta}$  are sampled based on the original goal  $\hat{\delta}$ .

# Full HRL Framework

## Hierarchical Reinforcement Learning (HRL)

## Full HRL Framework



# Next Steps

- ▷ Improve environment design (e.g., repeating actions).

# Next Steps

- ▷ Improve environment design (e.g., repeating actions).
- ▷ Design the reward functions.

# Next Steps

- ▷ Improve environment design (e.g., repeating actions).
- ▷ Design the reward functions.
- ▷ Account for the budget in the framework.