

IBM Semiconductors

# Cross-Process Defect Attribution using Potential Loss Analysis

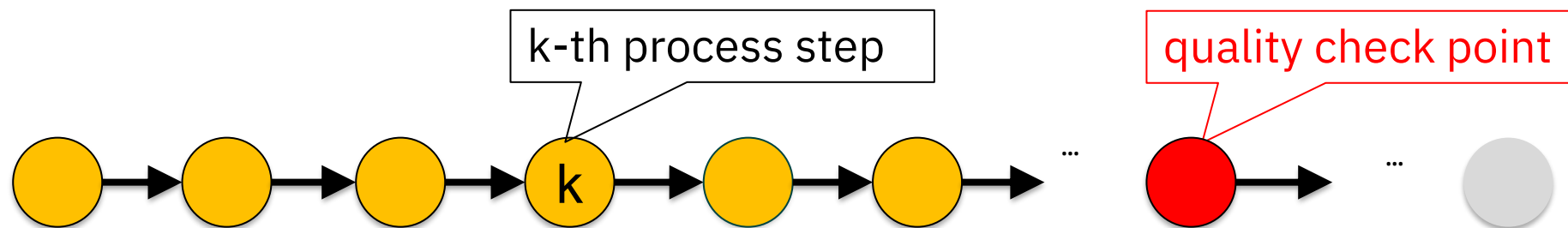
**Tsuyoshi (“Ide-san”) Ide** and Kohei Miyaguchi\*  
IBM Semiconductors, T.J. Watson Research Center  
\*IBM Research – Tokyo (currently with LY Corporation)

T. Ide and K. Miyaguchi, “Cross-process defect attribution suing potential loss analysis,” Proceedings of the 2025 Winter Simulation Conference (WSC 25), to appear; <https://arxiv.org/abs/2508.00895>

# Agenda

- Problem description and background
- Interventional causal attribution and its challenges
- Cross-process attribution with potential loss analysis (PLA)

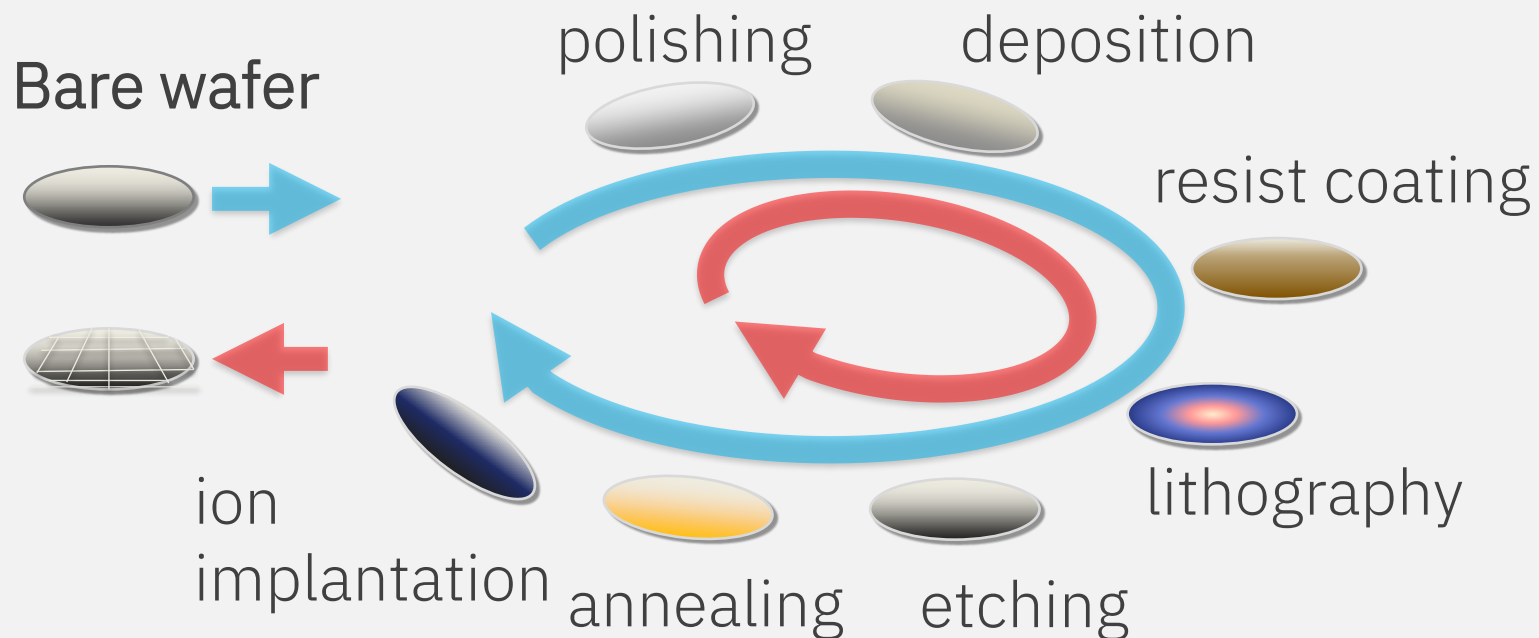
# Target task: cross-process defect attribution



Example: Semiconductor Manufacturing

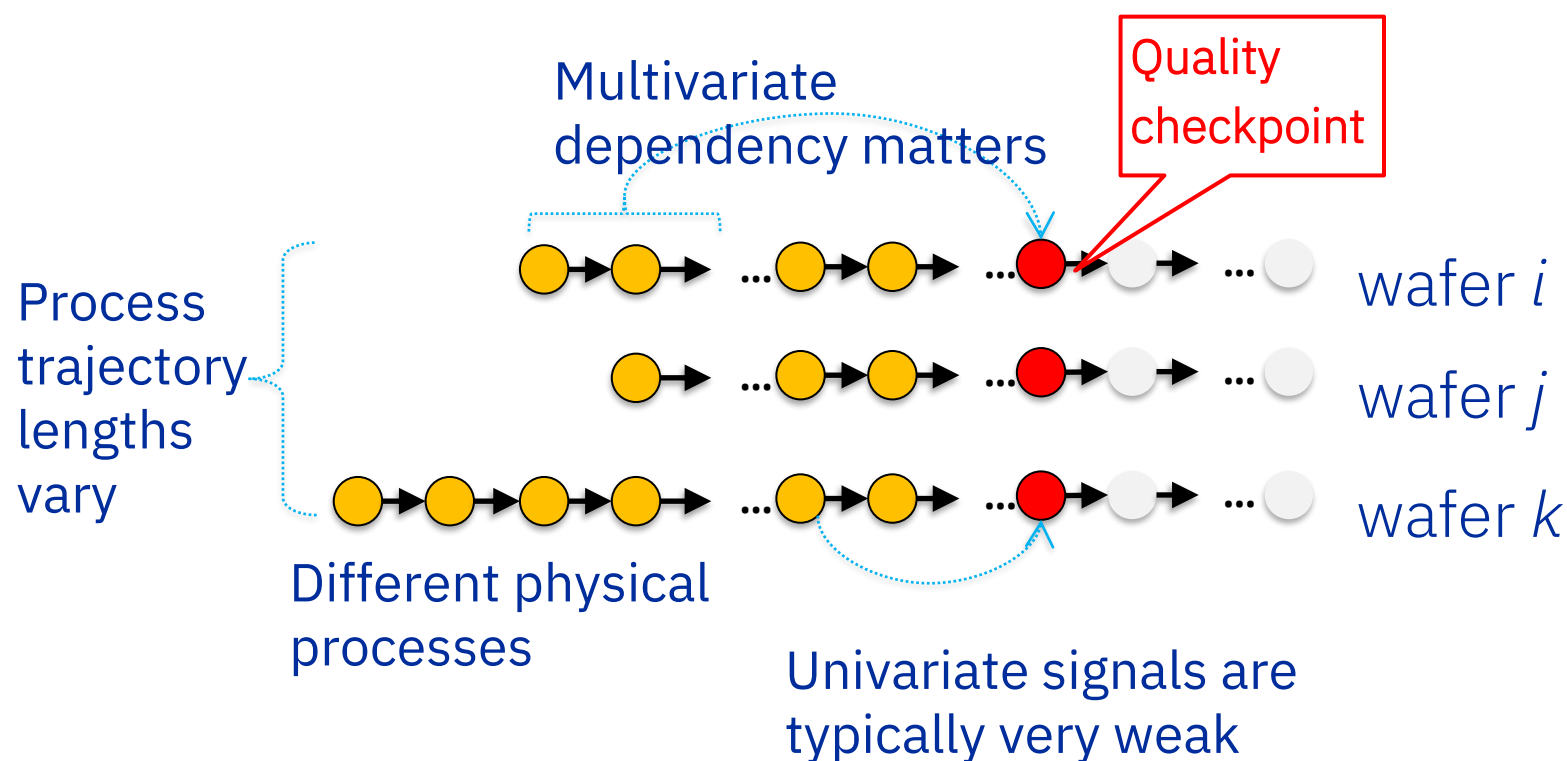
FEOL: device  
fabrication

BEOL: wiring  
formation



# Target task: cross-process defect attribution

Problem: Given a wafer quality metric value, compute the **attribution score** for each of the upstream process steps.

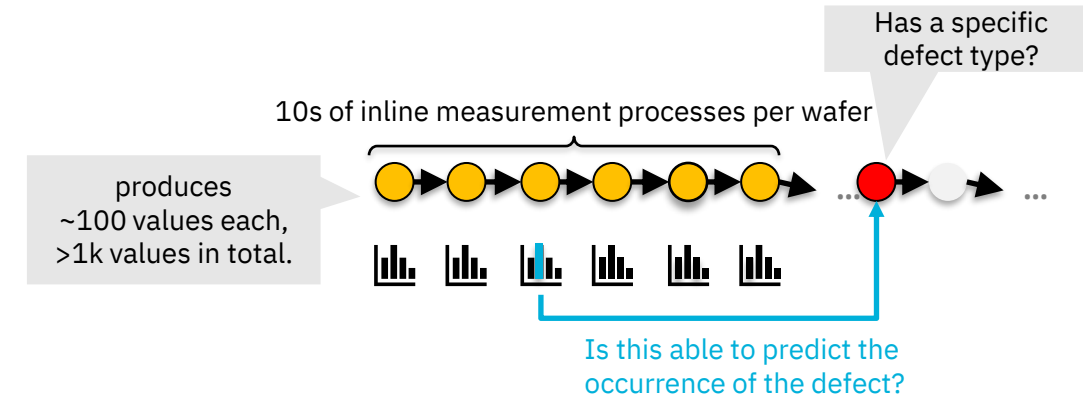


## ■ In current practice...

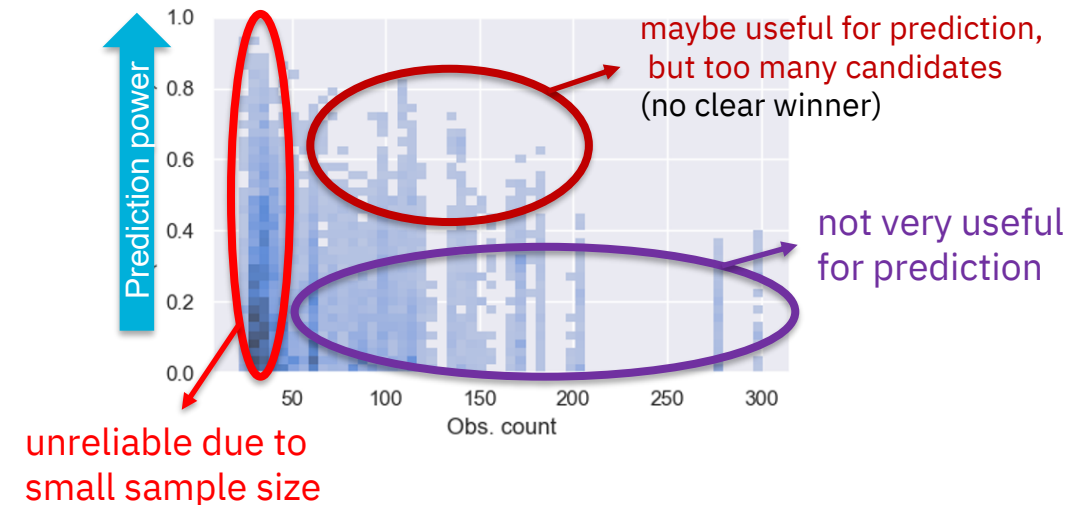
- The only viable approach is to run as many wafers as possible under varying conditions.
- Then, relatively simple statistical analysis is applied.
- This approach requires significant domain expertise to decide on parameter choices.
- This semi-manual approach is reaching its limits as technology nodes advance.

# Univariate correlation analysis is often not informative

- Univariate analysis typically yields weak predictive signals.
  - Example: measurement-based univariate defect prediction
    - ✓ Trained a univariate binary classifier to distinguish between good and bad wafers.
    - ✓ Accuracy tends to decrease as the observation count increases.
- Traditional correlation metrics apply to populations, not to individual wafers.
  - We need a wafer-wise diagnostic algorithm.



Defect occurrence prediction with single measurement values

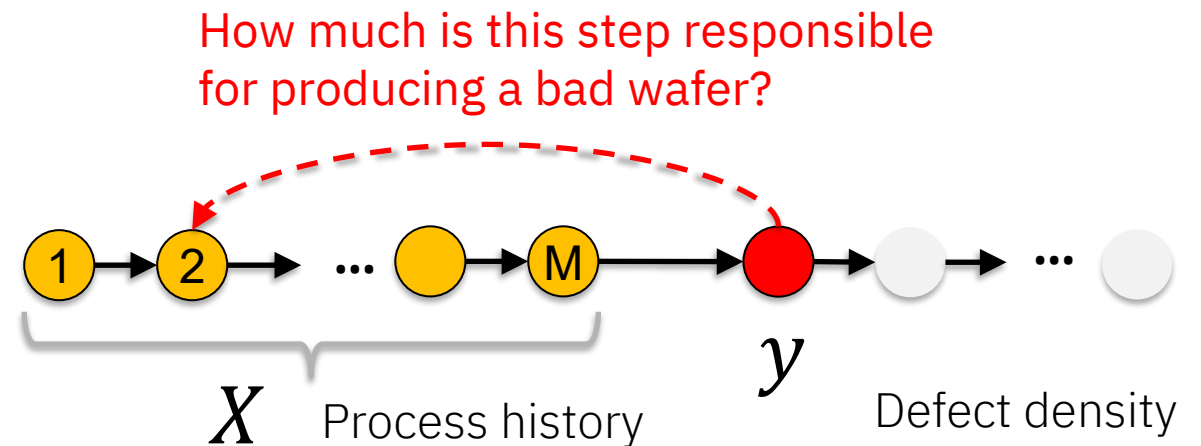


# Agenda

- Problem description and background
- Interventional causal attribution and its challenges
- Cross-process attribution with potential loss analysis (PLA)

# Defect attribution: Data assumptions and problem setting

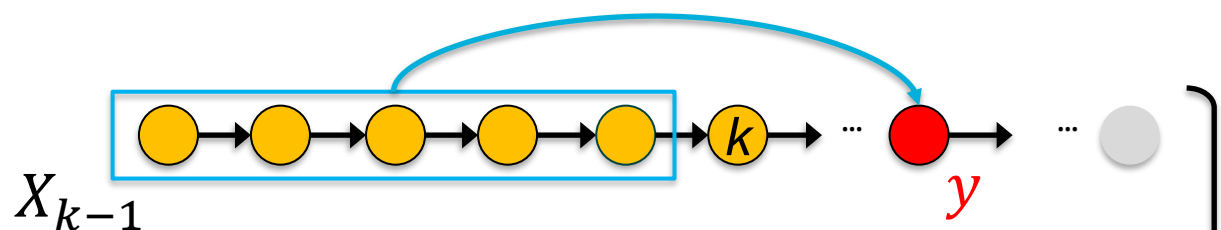
- Data  $D = \{(X^{(n)}, y^{(n)}) \mid n = 1, \dots, N\}$ 
  - $X^{(n)}$  : process trajectory  $(x_1^{(n)}, \dots, x_{L^{(n)}}^{(n)})$ . Each process is assumed to have a vector representation ( $\rightarrow$  discussed later).
  - $y^{(n)}$  : Product badness such as defect density (a real number).
- Task: wafer-wise defect attribution
  - Given a wafer quality metric value, compute the attribution score for each of the process steps.



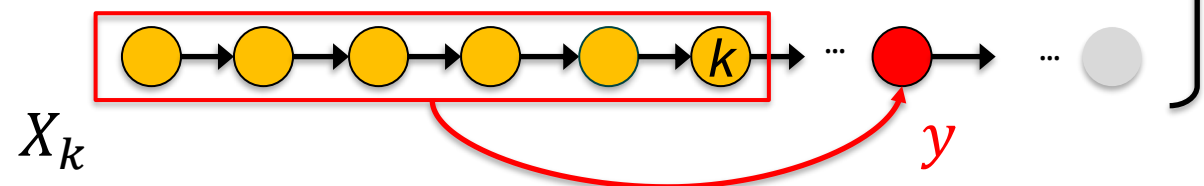
# Interventional Causal Attribution with Partial trajectory regression (PTR)

- PTR computes the attribution score of the  $k$ -th process by evaluating the impact of  $k$ 's “participation” in the process trajectory.

A: Prediction using a partial trajectory **not** including  $k$



B: Prediction using a partial trajectory including  $k$



} Any significant difference?

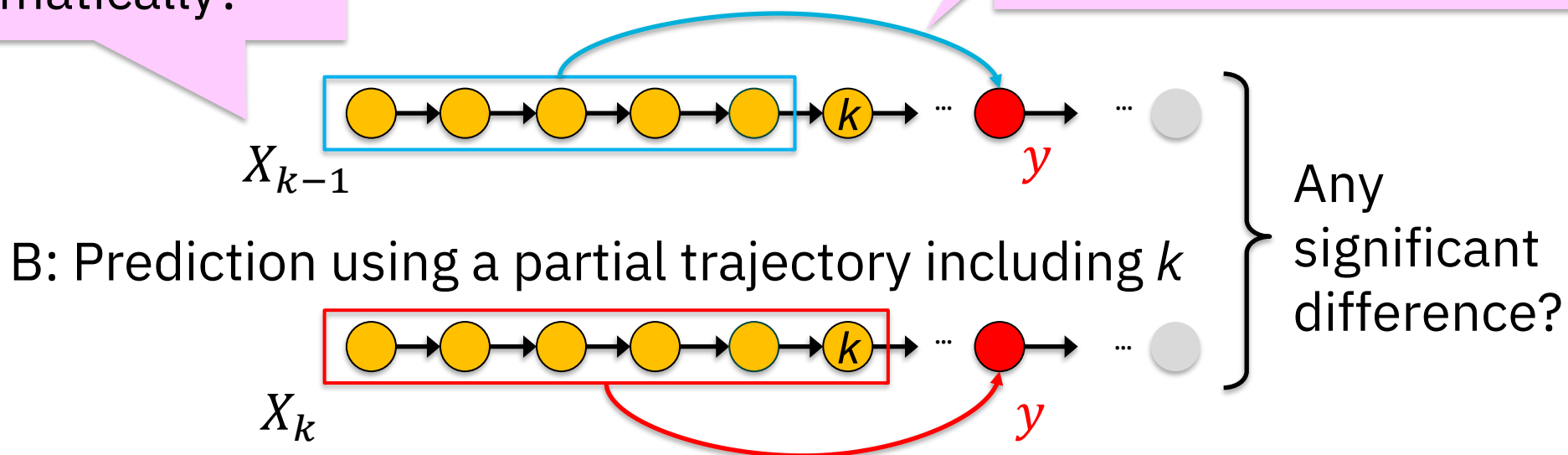


# Interventional Causal Attribution with Partial trajectory regression (PTR)

- PTR computes the attribution score of the  $k$ -th process by evaluating the impact of its “participation” in the process

How do we represent the process trajectories mathematically?

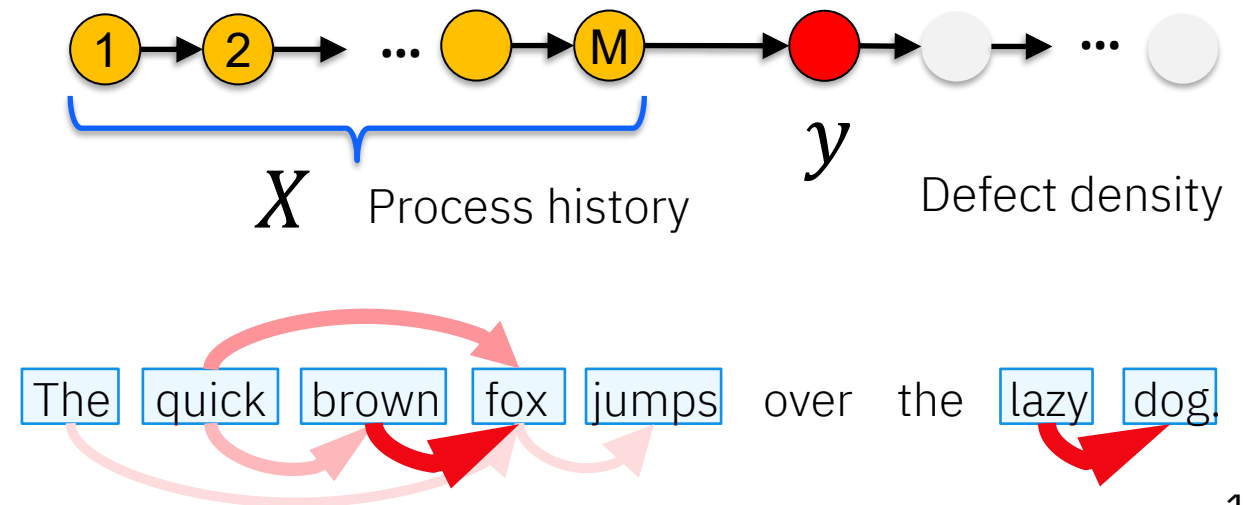
How do we predict  $y$  from a **partial** process trajectory?



# Representing process trajectories as vector sequences:

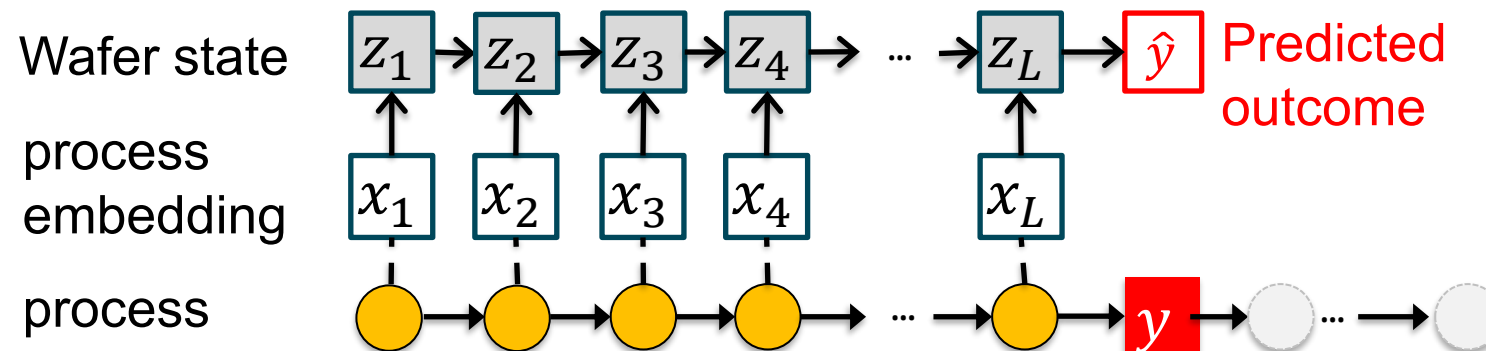
## Process embedding (“proc2vec”)

- Approach 1: Use measurement data (e.g., CDs) as a surrogate for process steps.
  - Straightforward, but data may have many missing entries due to wafer sampling.
- Approach 2: Use process data and apply (deep) embedding.
  - $\mathbf{x} = \text{ReLU}(\mathbf{W}\mathbf{u} + \mathbf{b})$ ,
    - ✓  $\mathbf{u}$ : some process data;  $\mathbf{W}$ ,  $\mathbf{b}$ : trainable parameters
  - Data hungry, prone to overfit.
- Approach 3: Symbolic embedding
  - Create a synthetic process token
    - ✓ “Process token” =  $\text{eqp\_id} \oplus \text{recipe\_id} \oplus \dots \oplus \text{tool\_trace}$
  - Employ multidimensional scaling based on a similarity matrix among the tokens



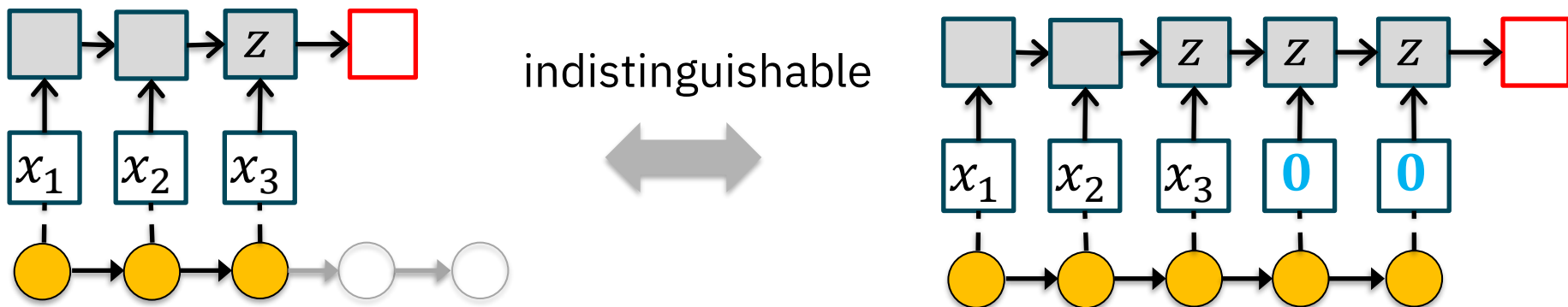
# Learning partial trajectory regression model

- Typically, a prediction function has a fixed number of input slots.
  - For 3 processes, it would be like  $f(x_1, x_2, x_3)$ , a 3-slot function.
  - Hence, cannot handle process trajectories with different lengths.
- State-space model (or RNN) eliminates this limitation
  - Partial prediction by a length-k trajectory:  $\hat{y} = f(z_k = \text{RNN}(x_1, \dots, x_k))$ 
    - ✓  $f$ : A parametric function with trainable parameters
    - ✓  $z_k$ : Latent state vector after observing  $x_k$
- Off-the-shelf RNN (e.g., LSTM, GRU) models need adjustments.
  - → next page



# Catch: You can't simply zero off downstream processes

- RNNs are generally data hungry. Careful model customization is needed.
  - Process timestamp needs special treatment (→ paper)
- RNNs may introduce a significant bias.
  - Example:  $k=3$  partial prediction is indistinguishable with  $k=5$  partial prediction with **zeroed-off input**.
    - ✓ i.e., RNN's partial trajectory prediction = full trajectory prediction with a zeroed-off process sequence



# Agenda

- Problem description and background
- Interventional causal attribution and its challenges
- Cross-process attribution with potential loss analysis (PLA)

# Eliminating the bias of partial trajectory analysis: Potential Loss Analysis (PLA)

- PLA: The attribution score for  $k$  is evaluated, given an optimal downstream trajectory:
  - $\min_{x_{k+1}, \dots, x_L} F(\mathbf{z}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_L)$ 
    - ✓  $\mathbf{z}_k$  : latent wafer state at process  $k$
    - ✓ We don't use  $F(\mathbf{z}_k, \mathbf{0}, \dots, \mathbf{0})$ .
- This trajectory optimization problem can be solved as a Bellman equation.
  - → Next page

attribution score

$$\alpha_k = F(\text{“optimal” downstream trajectory}) - F(\text{“optimal” downstream trajectory})$$

“optimal” downstream trajectory

“optimal” downstream trajectory

# Formalizing PLA as a reinforcement learning problem (1/2)

- $\min_{\mathbf{x}_{k+1}, \dots} F(\mathbf{z}_k, \mathbf{x}_{k+1}, \dots) = \min_{\mathbf{x}_{k+1}, \dots} \mathbb{E}[\sum_{t=1}^{\infty} C(\mathbf{z}_{k+t}) \mid \mathbf{z}_k]$ 
  - terminal reward model:  $C(\mathbf{z}) = \begin{cases} y(\mathbf{z}), & \mathbf{z} \in (\text{terminal state}) \\ 0, & \text{otherwise} \end{cases}$
- Bellman equation
  - $F^*(\mathbf{z}) \equiv \min_{\mathbf{x}_1, \dots} F(\mathbf{z}, \mathbf{x}_1, \dots) = \min_{\mathbf{x}_1} \{ C(\mathbf{z}) + \underbrace{\sum_{\mathbf{z}_2} p(\mathbf{z}_2 \mid \mathbf{x}_1, \mathbf{z}_1)}_{\text{transition model (assumed deterministic)}} F^*(\mathbf{z}_2) \}$
- The optimization problem we solve (with  $F^\theta(\mathbf{z})$  approximating  $F^*(\mathbf{z})$ ):
  - $\max_{\theta} \underbrace{\sum_{\mathbf{z}} \rho(\mathbf{z})}_{\text{empirical density of } \mathbf{z} \text{ (under deterministic assumption)}} F^\theta(\mathbf{z}) \quad \text{s.t.} \quad F^\theta(\mathbf{z}) \leq C(\mathbf{z}) + F^*(\mathbf{z}'), \quad \forall (\mathbf{z} \rightarrow \mathbf{z}'),$

# Formalizing PLA as a reinforcement learning problem (2/2)

- Final objective function to be maximized

- $$R(\theta \mid \mu) = \frac{1}{N} \sum_{n=1}^N \left[ \frac{\mu}{L^{(n)}} \sum_{t=1}^{L^{(n)}} F^{\theta}(z_t^{(n)}) - \underbrace{\frac{1}{2} \{y^{(n)} - F^{\theta}(z_{L^{(n)}}^{(n)})\}^2}_{\text{squared loss}} - \frac{1}{2} \sum_{t=1}^{L^{(n)}-1} \{F^{\theta}(z_{t+1}^{(n)}) - F^{\theta}(z_t^{(n)})\}^2 \right]$$

- This provides the partial prediction function and attribution model simultaneously.

- The time difference of the F function can be directly parameterized:

- $G^{\theta}(\mathbf{z}_t, \mathbf{z}_{t-1}) \equiv F^{\theta}(\mathbf{z}_t) - F^{\theta}(\mathbf{z}_{t-1}) = \text{ReLU}_{\theta}(\mathbf{z}_t \oplus \mathbf{z}_{t-1})$
  - This function provides the attribution score for the k-th process.
    - ✓ Yields a positive number



# 2nm process defect diagnosis example

- The graphs plot the cumulative attribution score.
  - Big jump = likely root cause
  - The model was trained with 727 wafers.
- PTR approach does not provide meaningful signals.
- PLA successfully detects likely root cause
  - In this case, they correspond to too long waiting hours at a certain piece of equipment.

