



Hewlett Packard
Enterprise



Causal Inference for Digital Twins in GPU Data Centers and Smart Grids

Presented by,
Pavana Prakash, Hewlett Packard Labs

Rolando P. Hong Enriquez[‡], **Pavana Prakash**[‡], Ebad Taheri[‡], Aditya Dhakal[‡], Matthias Maiterth^{*}, Wesley Brewer^{*}, Dejan Milojicic[‡]

[‡]Hewlett Packard Labs, ^{*}Oak Ridge National Laboratory (ORNL)

May 7, 2025



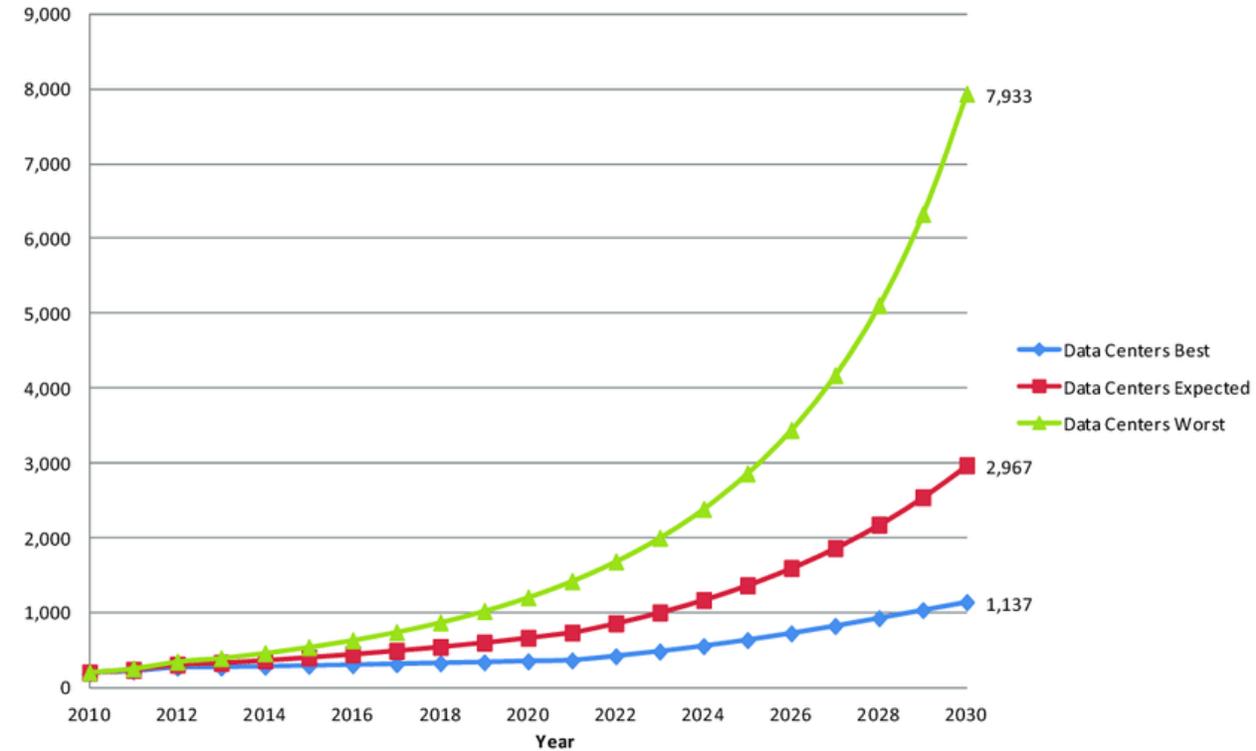
Hewlett Packard
Labs

Motivation



Data Centers – Global distribution

Electricity usage (TWh) of Data Centers 2010-2030

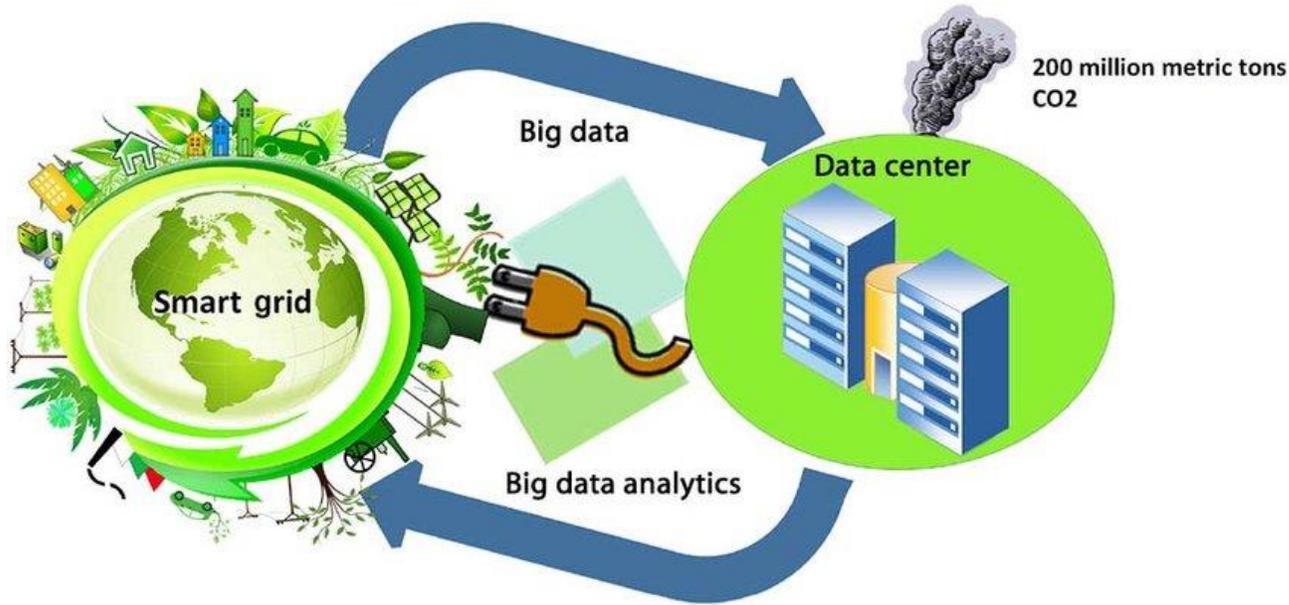


Data Centers - Global electricity usage

<https://www.datacentermap.com/datacenters/>

Andrae, Anders & Edler, Tomas. (2015). On Global Electricity Usage of Communication Technology: Trends to 2030. Challenges. 6. 117-157. 10.3390/challe6010117.

Motivation

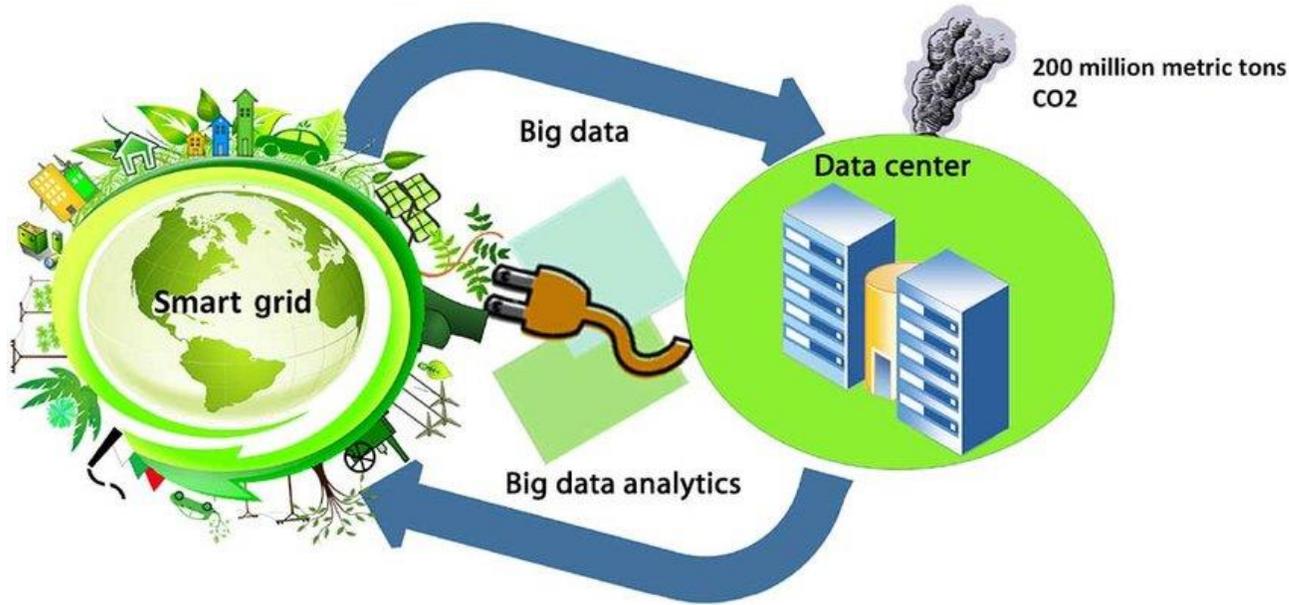


Integration (Data Centres + Smart Grids)

Challenges

- These integrations would require massive amounts of data feeding increasingly larger models
- Need combination of
 - Advanced data analytics
 - Digital Twins for system insights
 - AI models with improved accuracy
- AI models → larger models
 - Higher power consumption
 - Diminished explainability.

Motivation



Integration (Data Centres + Smart Grids)

Solutions

1. Build less complex (i.e., energy efficient) AI models with higher causal explanatory power
2. Explore the efficacy of current causal inference methods.

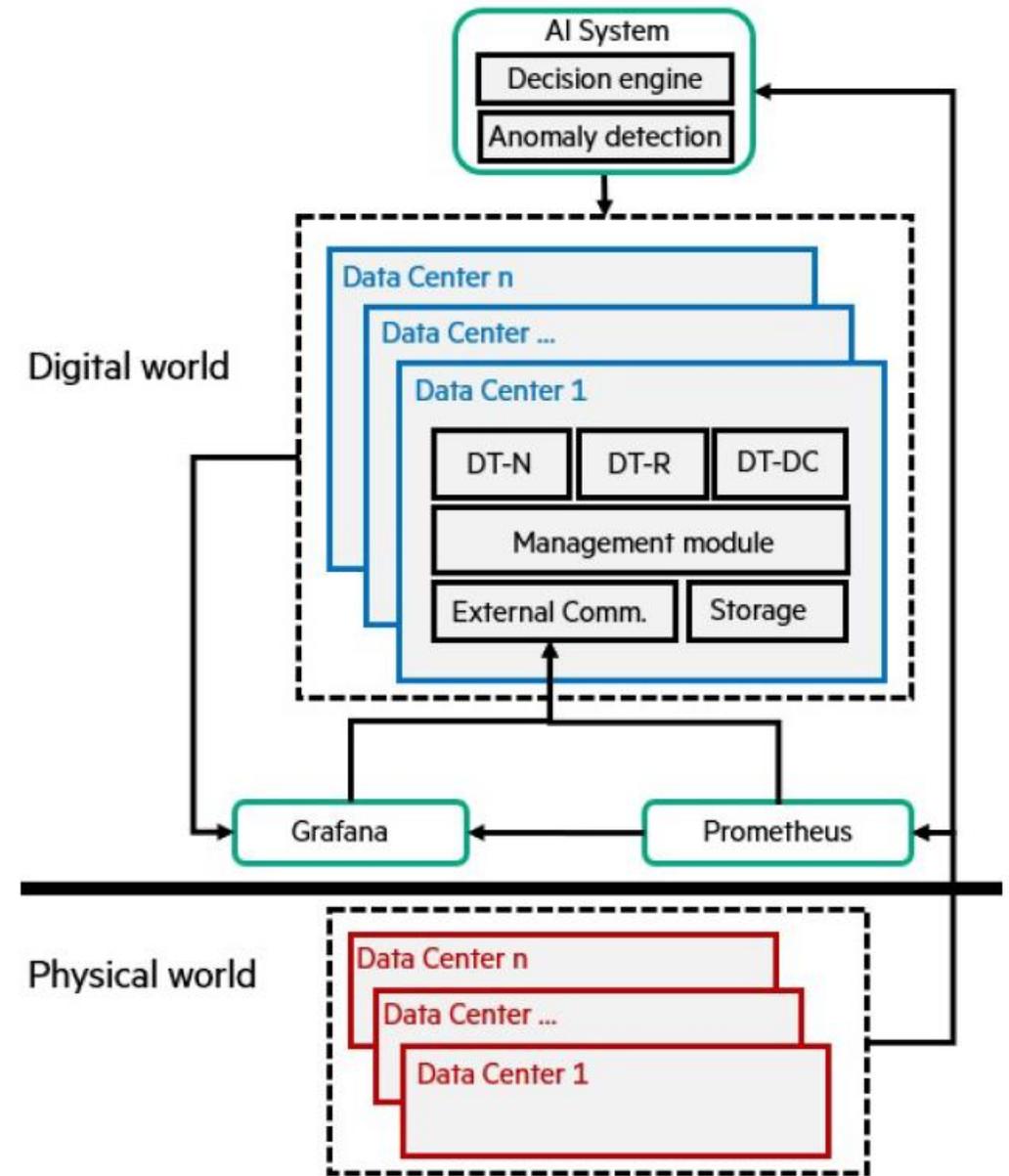
Need for Causality

- **Correlation** is a statistical measure of association between two variables that expresses the direction and extent to which they are related
 - Example: GPU temperature increase -- system failure
 - May lead to false alarms or missed root causes
- **Causality** refers to a **cause-and-effect relationship**, where one event or condition directly influences another
 - Example: Cooling failure (e.g., malfunctioning fans) -- GPU temperature rise -- thermal throttling -- system failure
 - Helps identify root cause analysis, improves prediction accuracy, leading to informed decision-making



Our Solution

- Digital Twins: a virtual representation of the data center infrastructure
 - Incorporated causal insights into its monitoring and simulation capabilities
 - Enhanced real-time monitoring and prediction of GPU performance and failures
 - Perform “what if?” analyses to improve and simplify optimization



Digital Twin with Causal Analyzer

- Causal metrics:

1. Average Treatment Effect (ATE)

$$\text{ATE} = \mathbb{E}[Y(1) - Y(0)]$$

$$\text{ATE} = \frac{1}{N} \sum_{i=1}^N Y_i(1) - Y_i(0)$$

- $Y(1)$ - outcome when the treatment is applied
- $Y(0)$ - outcome without treatment
- $\mathbb{E}[\cdot]$ - expectation operator, which averages the difference over the entire population
- N - total number of individuals in the population.

2. Prediction effectiveness

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

- AI model: Binary classification model to predict failures
- Dataset: Summit*
 - Node-level data on power, energy, job scheduling and failures
 - Observed over three years with 27,648 V100 GPUs.

*Woong Shin, Vladyslav Oles, Anna Schmedding, George Ostrouchov, Evgenia Smirni, Christian Engelmann, and Feiyi Wang. 2023. *OLCF Summit Supercomputer GPU Snapshots During Double-Bit Errors and Normal Operations. Technical Report.* Oak Ridge National Lab.(ORNL), Oak Ridge, TN (United States).

Causality Evaluation and Results

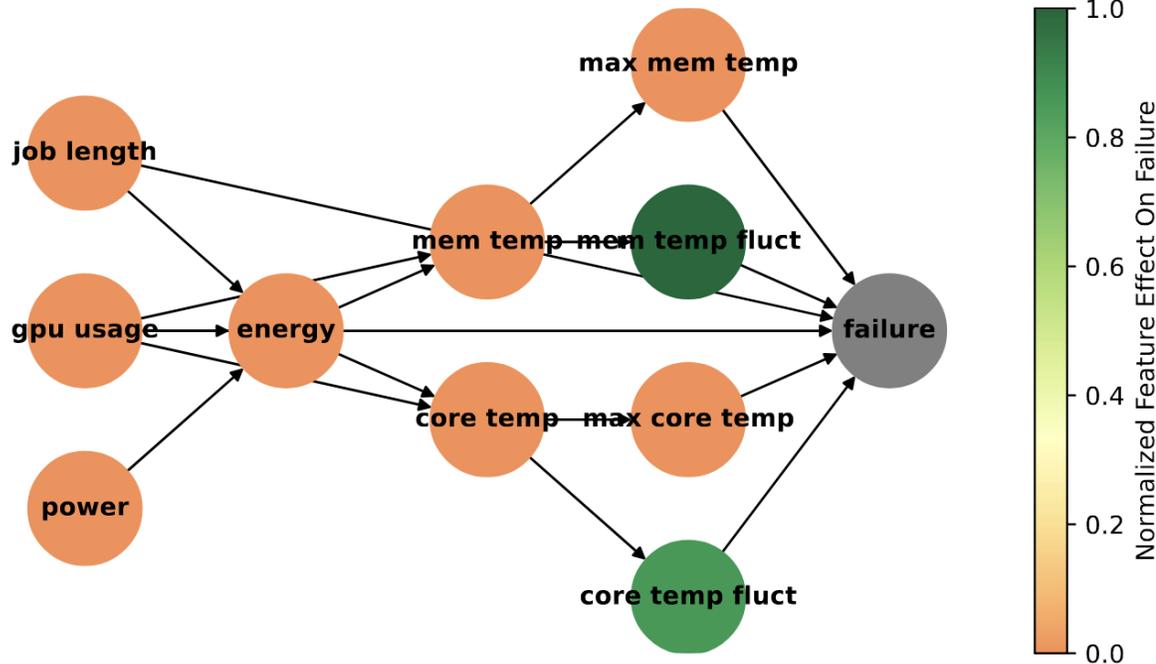


Figure 1:

- Graph shows feature relationships and their impact on failures, with color-coded nodes representing feature influence
- Core and memory temperature fluctuations are identified as the most significant causes of failures.

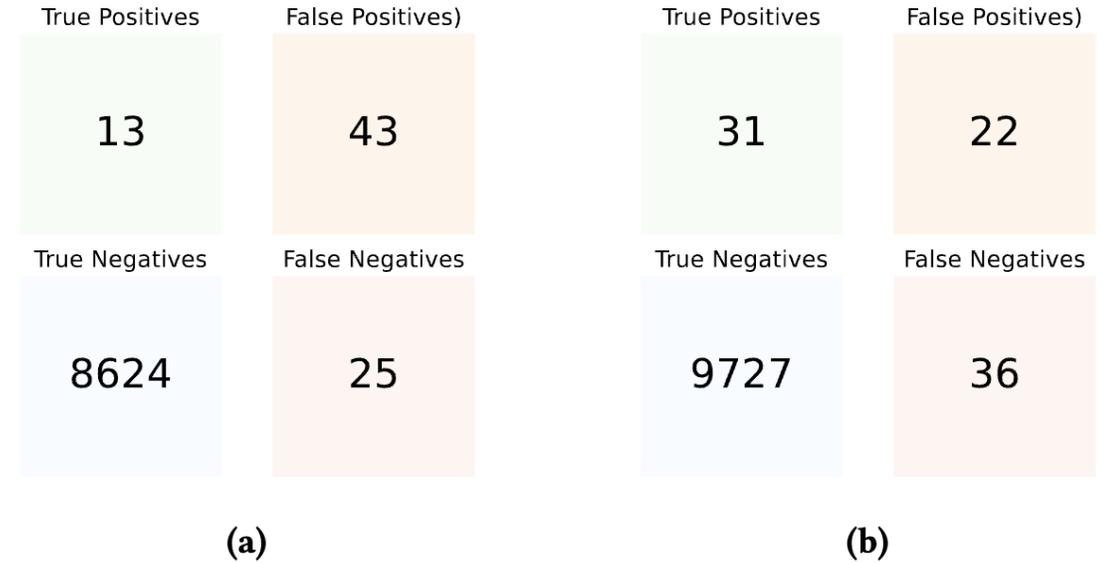
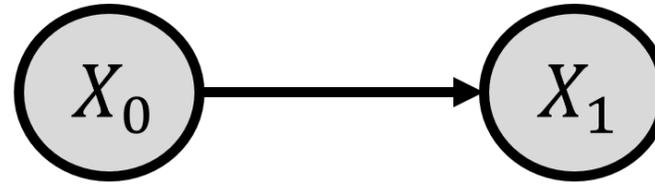


Figure 2: Failure prediction with
 (a) All features and
 (b) Two features suggested by causal analysis.

- Our causal analysis approach improves prediction efficiency: Precision from 0.271 to 0.577 and Recall from 0.316 to 0.493.
- Reduces model complexity: Utilizes only 2 features instead of 10.

Data for Causality (Generated – Structural Models)

- Causal analysis in digital twins uses different data types, affecting model accuracy
- This study focuses on **GPU failure prediction** in data centers
- We use **ground-truth validation data:** experimental (real) and mathematically generated (synthetic) to test our methods.

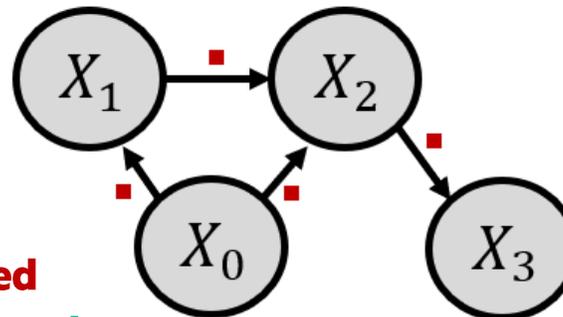


$$X_0 \sim N(0,1)$$

$$X_1 = 1/2 X_0 + \sqrt{3}/2 Z_1$$

with $Z_1 \sim N(0,1)$ and independent from X_0

Structural equations for a causal pair formed by a two-node linear gaussian system.



Extracting causal pairs (red squares) from complex causal models (Simpson paradox).

$$X_0 \sim N(0,1)$$

$$X_1 = s(1 - X_0) + \sqrt{3/20} Z_1$$

$$X_2 = \tanh(2X_1) + 3/2 X_0 - 1 + \tanh(Z_2)$$

$$X_3 = 5 \tanh\left(\frac{X_2 - 4}{5}\right) + 3 + \frac{1}{\sqrt{10}} Z_3$$

where Z_1, Z_2, Z_3 are mutually independent and independent from X_0 and $s(x) = \log(1 + \exp(x))$ is the softplus function

$(X_0, X_1), (X_0, X_2), (X_1, X_2),$ and (X_2, X_3)

Data for Causality (Generated – Dynamical and Chaotic Systems)

- **Dynamical systems** like coupled logistic map, model interdependent relationships between variables to simulate causal effects by adjusting coupling terms:

$\mathbf{x} \rightarrow \mathbf{y}$: X causes Y

$\mathbf{x} \leftarrow \mathbf{y}$: Y causes X

$\mathbf{x} \leftrightarrow \mathbf{y}$: X and Y cause each other

$\mathbf{x} \perp \mathbf{y}$: No causal relationship between X and Y

- **Chaotic systems**, such as the Rössler and Lorenz attractors, can drive the behavior of other systems, generating complex, synthetic causal data for testing models.

X - Rössler attractor

$$\dot{x} = -(y + z)$$

$$\dot{y} = x + \alpha y$$

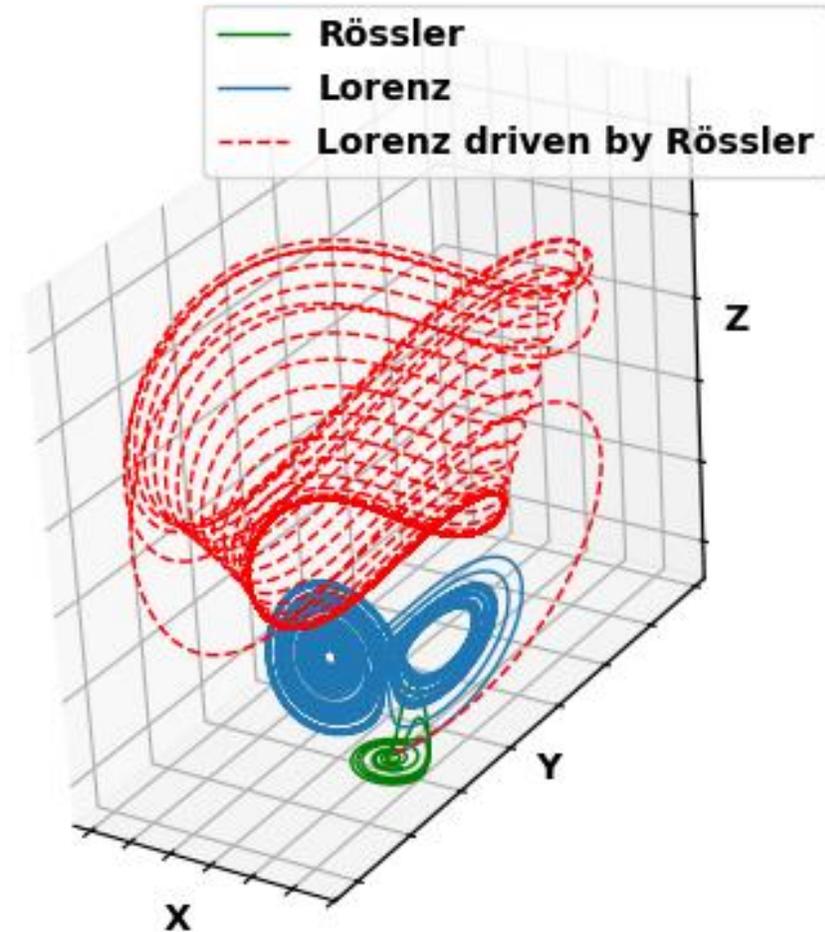
$$\dot{z} = \beta + z(x - \gamma)$$

Y - Lorenz attractor

$$\dot{x} = \delta(y - x)$$

$$\dot{y} = \epsilon x - y - xz + C y$$

$$\dot{z} = xy - \zeta z$$



Evolution of free and coupled (i.e., causally linked) chaotic dynamical systems

Data for Causality (resumed)

Datasets used in our work

Dataset	Source	Type
Summit ¹	Experimental	Mixed
Smart Grid ¹	Experimental	Mixed
inHouse	Experimental	Causal Pairs
Csuite	Generated	Mixed
Couple Logistic Maps	Generated	Causal Pairs
Driven attractors	Generated	Mixed
BioGeoScience	Generated	DAG
Health BMI	Generated	Mixed

¹ Target dataset.



Causal Inference – Working Definition

“The process of inferring whether or not an event A is caused by another event B”. (...from philosophy/psychology)

“From a set of N measured variables describing a certain process, find the subset of M causally-related variables describing the same process with maximal accuracy and minimal model capacity. Extract the causal graph between the M variables.”

There are 3 major theories on which causal inference methods have been developed,

- (a) causal calculus
- (b) information theory, and
- (c) dynamical system theory.



Causality Inference – Causal Calculus

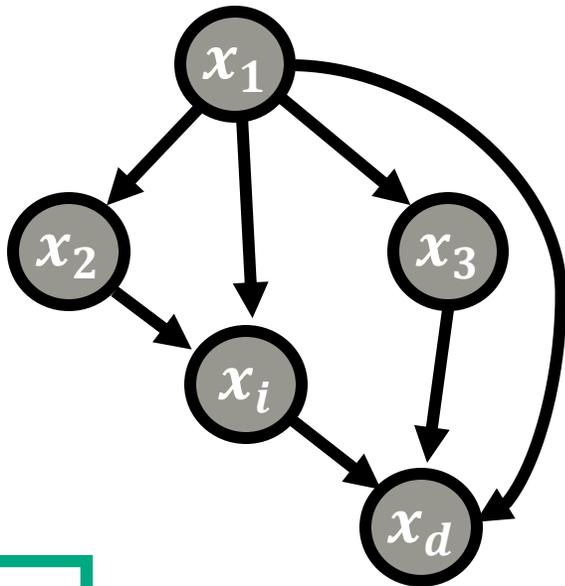


Judea Pearl

Turing Award - 2011

“For fundamental contributions to artificial intelligence through the development of a calculus for probabilistic and causal reasoning.”

$$p(x_1, \dots, x_d) = \prod_{j=1}^d p(x_j | x_{pa(j)})$$



This formulation assumes that a variable or node is independent of previous non-parent nodes and therefore captures dependency between variables or nodes using Bayesian networks.

Causality Inference – Causal Calculus

Fetal ECG dataset

Multichannel ECG recording

- **Channel D** comes directly from the baby's head at labor
- **Channels (Ab1, ...Ab4)** are indirect measures taken from the abdomen of the mother.
- **Channel E**, annotations (causally unrelated).

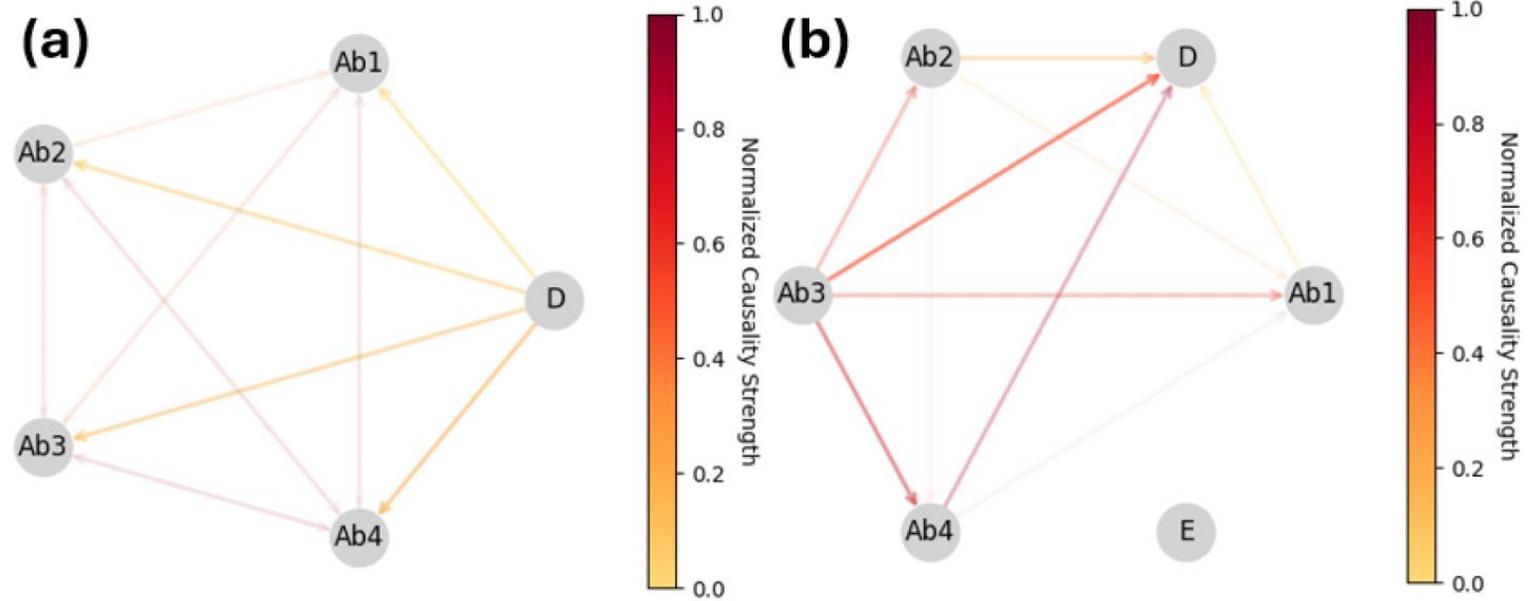


Fig. 6. Performance of causal calculus on the fetal ECG dataset [35].(a) Ground truth - polynomial graph ($p(x) = x^6$) (b) Generated data -polynomial graph ($p(x) = x^5$). This method seems to recognize valid causal relationships and correctly excludes E as a causal variable.



Causality Inference – Information Theory (I)



Norbert Wiener

Theory of Predictions - 1956

The statement “X causes Y”, implies that we have better chances of predicting the correct value of variable Y by including the past and present of values of X, than by using variable Y alone.

Transfer entropy – Schreiber, 2000

$$TE_{X,Y} = \sum_{y_{t+1}, y_t^n, x_t^m} p(y_{t+1}, y_t^n, x_t^m) \log \left(\frac{p(y_{t+1} | y_t^n, x_t^m)}{p(y_{t+1} | y_t^n)} \right)$$

Causality Inference – Information Theory (II)

Transfer entropy – Bivariate vs Multivariate

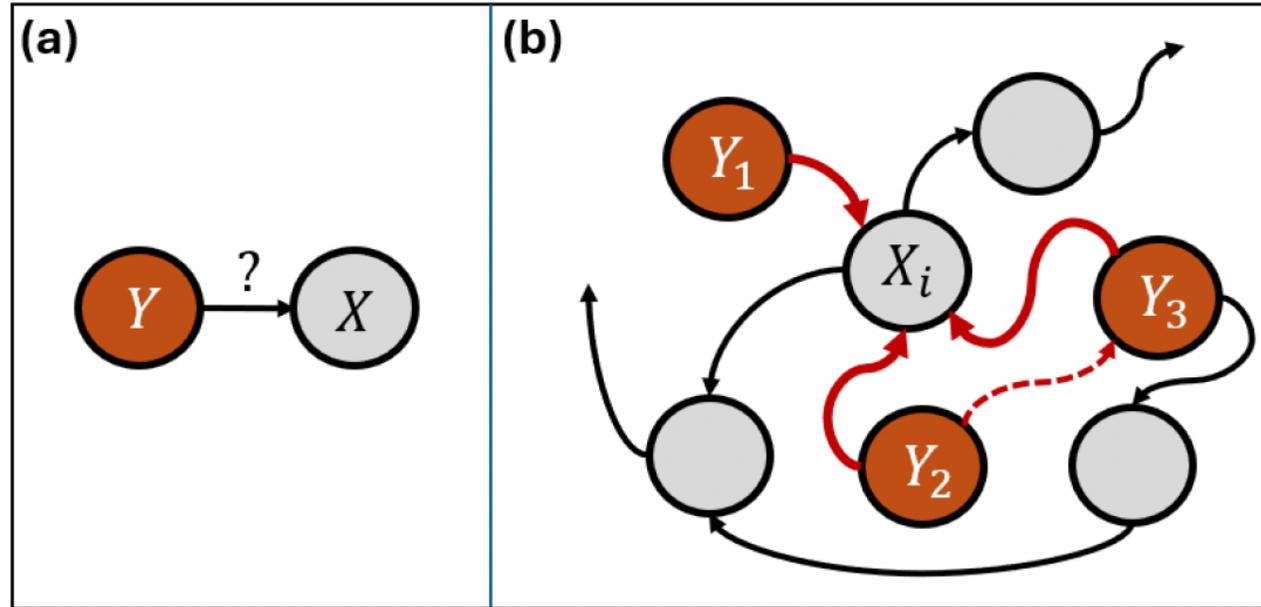


Fig. 7. Transfer entropy. (a) Bivariate case. Aims to determine if there is a causal link between selected nodes. (b) Multivariate case. Aims to find the set of nodes $V_{X_i} = \{Y_1, Y_2, Y_3\}$ that can be used to predict the next state of node X with statistical significance. Repeat and optimized the procedure for all the nodes X_i in the network of potential causal variables

Causality Inference – Information Theory

Fetal ECG dataset

Multichannel ECG recording

- **Channel D** comes directly from the baby's head at labor
- **Channels (Ab1, ...Ab4)** are indirect measures taken from the abdomen of the mother.
- **Channel E**, annotations (causally unrelated).

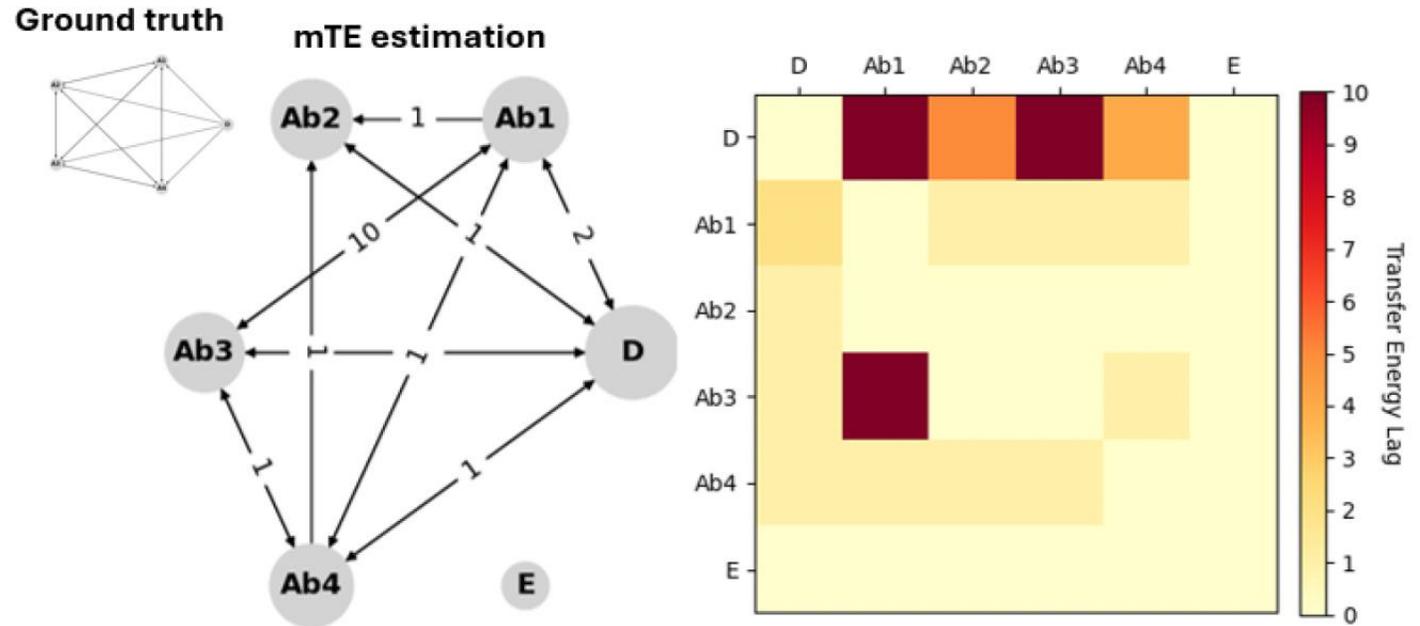


Fig. 9. Performance of Multivariate Transfer Energy (mTE) on the fetal ECG [35]. The ground truth for this model is shown in the inset. The method recognizes valid causal relationships, particularly it leaves out the only node that does not have causal relationships with the others. The graph polynomial extracted from the estimated model ($p(x) = x^6 + 7x^4 - 10x^3 - 4x^2$) is different from the one extracted from the ground truth graph ($p(x) = x^5$). The estimated graph is also not isomorphic with the ground truth; its edit distance from it is 7.

Causality Inference – Dynamical Systems

Fetal ECG dataset

Multichannel ECG recording

- **Channel D** comes directly from the baby's head at labor
- **Channels (Ab1, ...Ab4)** are indirect measures taken from the abdomen of the mother.
- **Channel E**, annotations (causally unrelated).

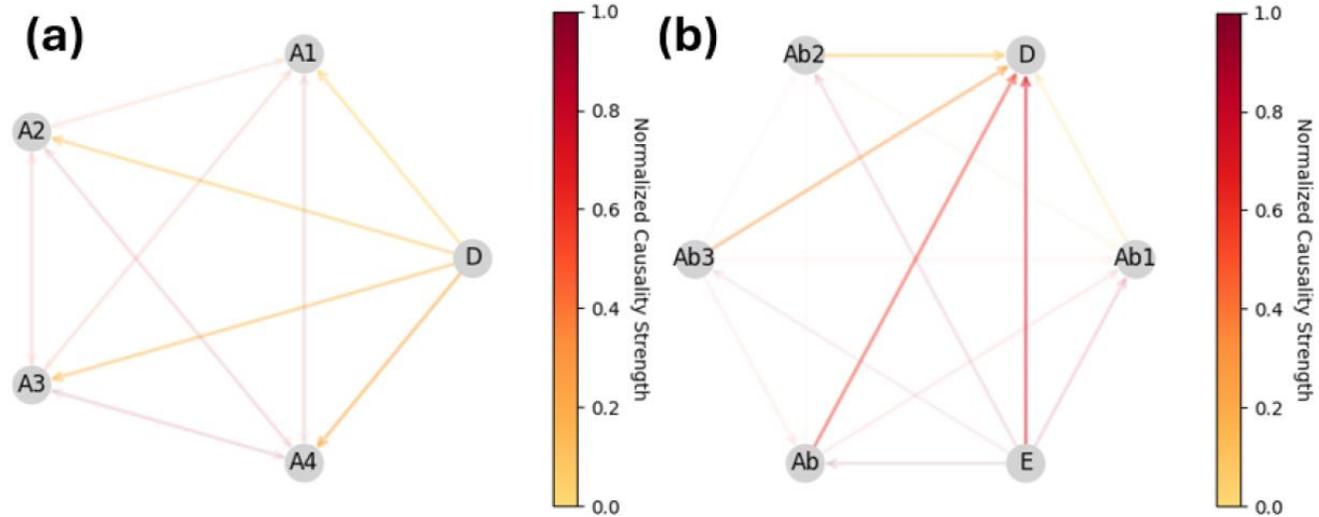
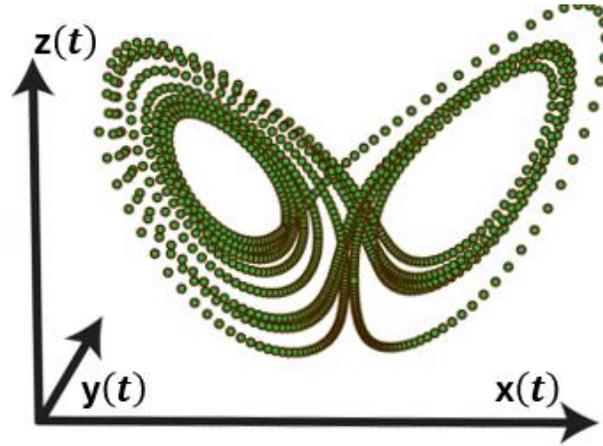


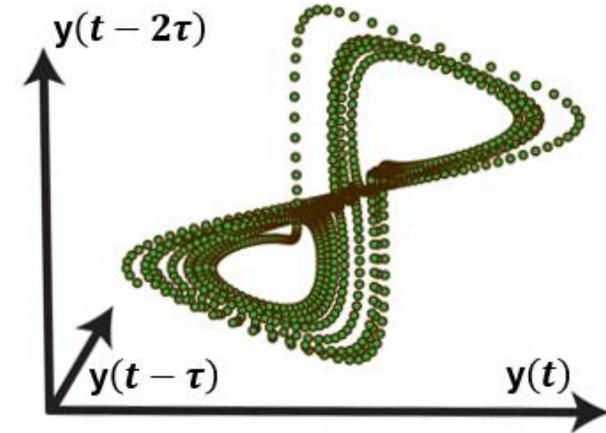
Fig. 13. Performance of Convergent Cross Mapping (CCM) on the fetal ECG [35].(a) Ground truth bas build based on several observations. Node E are non-causally related annotations. Node D is measured on the baby's head directly. The other measures are performed on the mother's abdomen (b) CCM estimation based on recorded data. Although CCM fails to exclude E from the graph, it does seem to recognize valid causal relationships. The graph polynomial extracted from the estimated model ($p(x) = x^6$) is different from the one extracted from the ground truth graph ($p(x) = x^5$). The estimated graph is also not isomorphic with the ground truth; its edit distance from it is 12.

Causality Inference – Dynamical Systems (I)

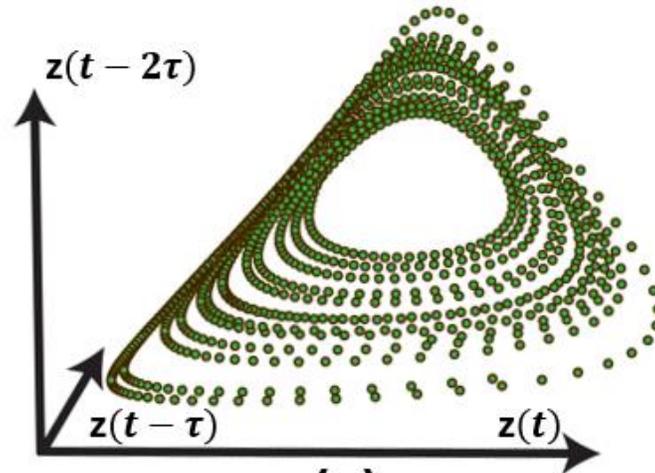
1. A dynamical system is the mathematical representation of the time evolution of a set of variables
2. We can seldom measure all the variables of high dimensional dynamical systems and therefore study the system in its original manifold M .
3. We can reconstruct a shadow manifold (e.g., M_x) that is topologically equivalent to M .
4. M_x is built using past (i.e., time-delayed) values of x as variables.



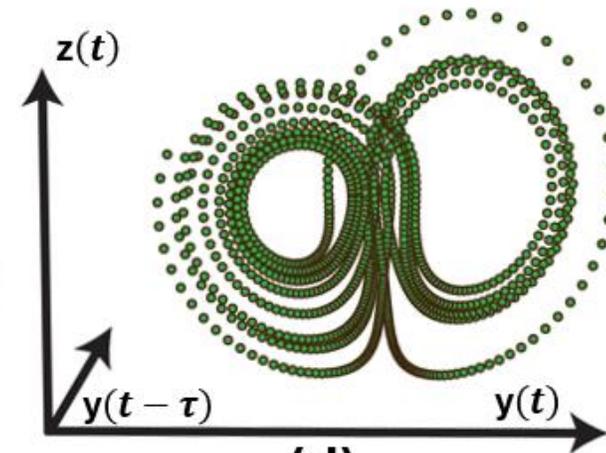
(a)



(b)



(c)

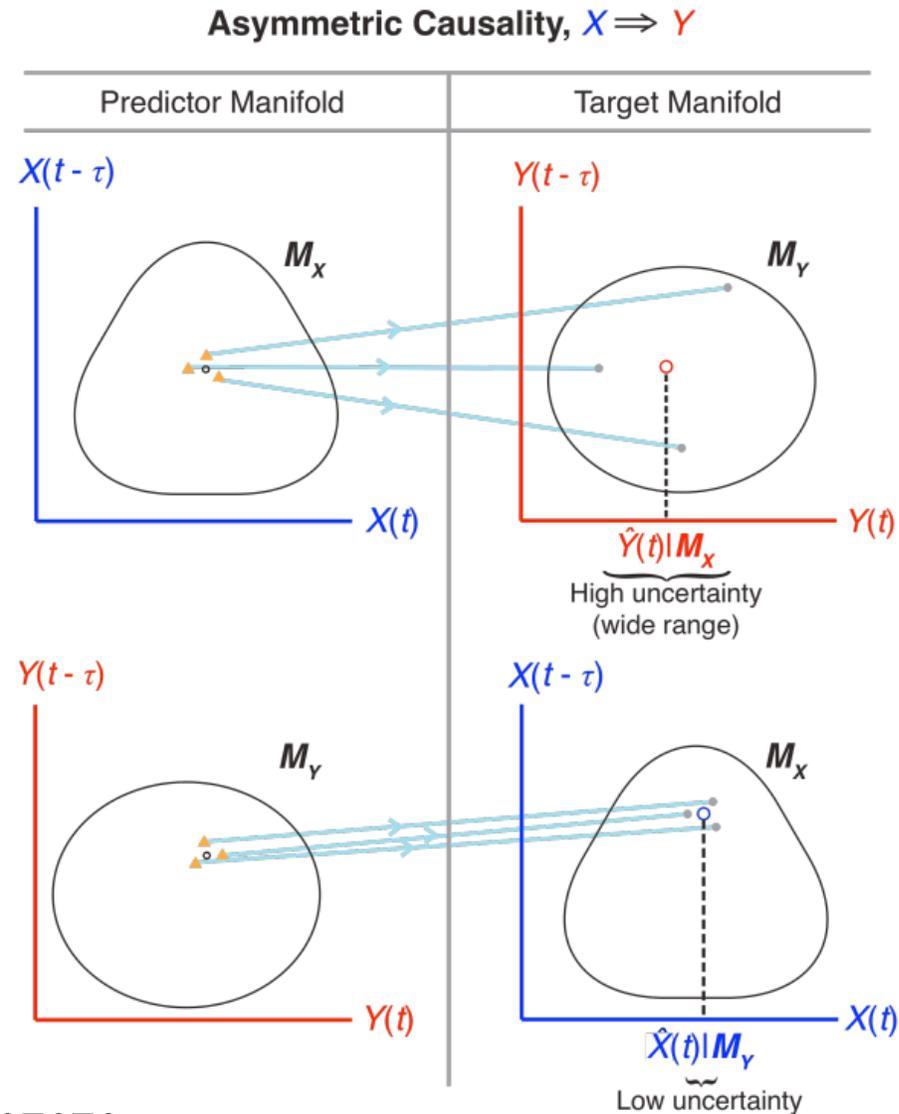


(d)

Causality Inference – Dynamical Systems (II)

Convergent Cross Mappings

- If variable X and Y are causally related, they belong to the same dynamical system
- To measure causality here we use the shadow manifolds (M_X, M_Y)
- By measuring correlations between M_X, M_Y , and their shared manifold M , we indirectly assess the causality between X and Y .



Focus Dataset – Smart Grids

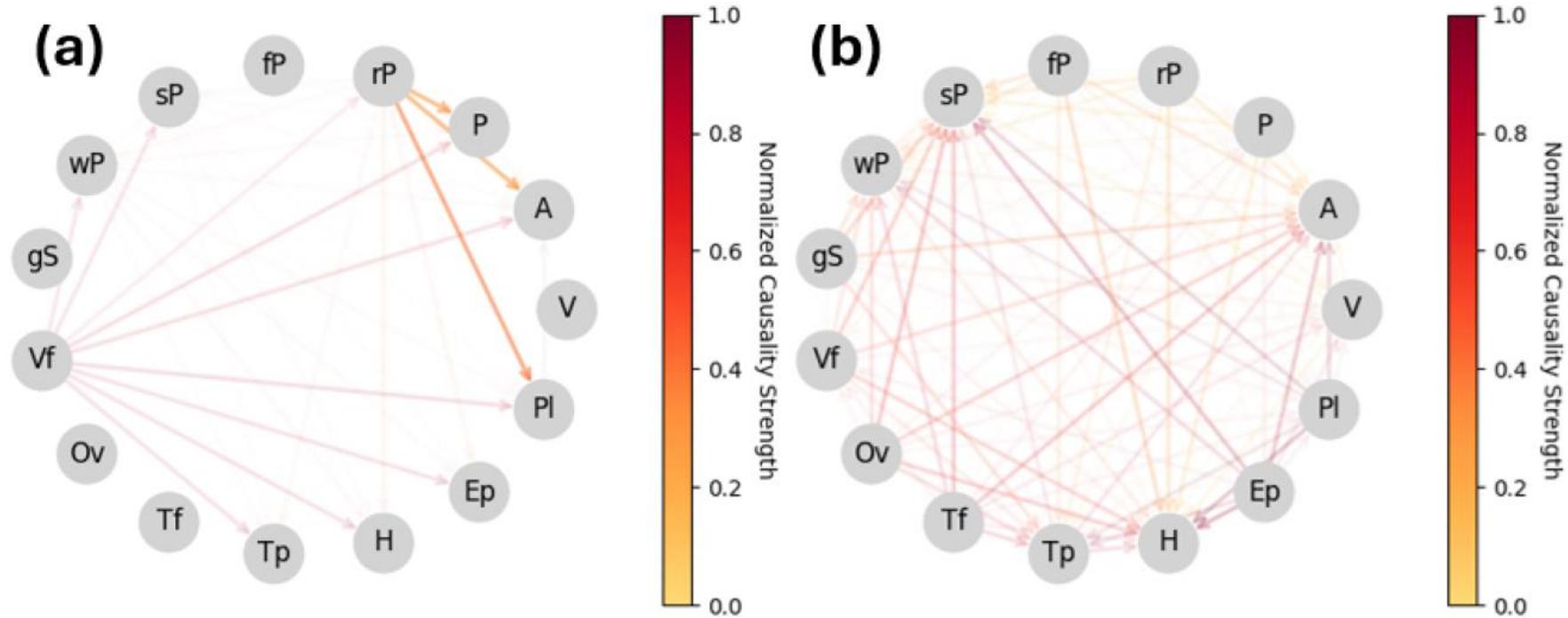


Fig. 15. Causality inference of the smart grid target dataset. [29] using:(a) Causal calculus, and (b) Convergent Cross Mapping (CCM).



Focus Dataset – Data Centers

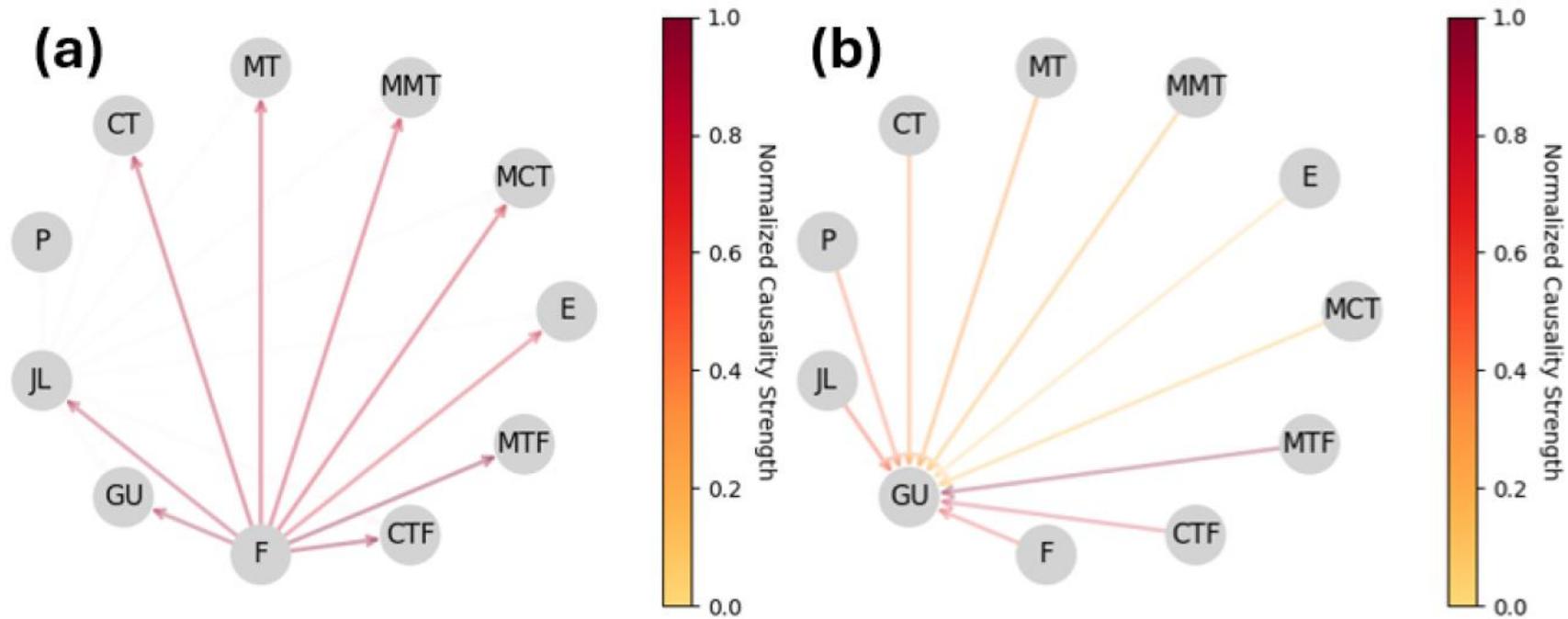


Fig. 14. Causality inference of the Summit target dataset. [28] using:(a) Causal calculus, and (b) Convergent Cross Mapping (CCM).



Causality Inference – Resumed

TABLE II
CAUSAL GRAPH ESTIMATION ON SELECTED DATASETS.¹

	Causal Calc.	mTE	CCM
Fetal ECG	■	■	■
BioGeoScience	■	■	■
Driven Attractors	■	□	■
Smart Grid	■	□	■
Summit	■	□	■

¹ Other results are shown in the Appendix.

Success (■) or failure (□) of using the causal inference methods to generate causal graphs or DAGs for selected datasets



Conclusions

1. We curated and tested datasets for causality inference, which is an important requirement given the scarcity of those datasets.
2. Using both generated and experimental datasets with known ground truth, we validated three types of causality inference methods respectively based on causal calculus, information theory and dynamical system theory
3. We applied the validated methods to two target datasets: (a) Summit - GPU data centres, and (b) Smart Grids.

The resulting new causal insights suggest new ways to performs meaningful and hopefully energy-efficient predictions for future digital twin solutions on complex infrastructures.



Future work

1. Continue studies with enriched target datasets. We are currently experimenting with data center cooling data, but data from other systems are still lacking.
2. Several other causality methods need to be tested. Research on the interaction between causality metrics and AI algorithms is, therefore, a promising long-term area of research for us.



Thank you

Pavana Prakash
Research Scientist, Hewlett Packard Labs
prakash@hpe.com

Rolando P. Hong Enriquez: rhong@hpe.com

RP Hong Enriquez ‡, P Prakash ‡, E Taheri ‡, A Dhakal‡, M Maiterth*, W Brewer*, D Milojcic‡
‡Hewlett Packard Labs, *Oak Ridge National Laboratory (ORNL)

