# RaiLog RCA : Railway Log-based Root Cause Analysis

Nadia Chouchani[1], Alexandre Trilla, Ossee Yiboe, Rajesh Rajendran, and
Ankit Bhoge

Alstom Transport, D&IS France,
`nadia.chouchani@alstomgroup.com`

**Abstract.** Improving safety and operational efficiency in the rail transport industry relies on a precise understanding of the root causes behind system failures. In this research, we propose RaiLog RCA, a comprehensive root cause analysis approach leveraging log data from railway operating systems. This approach is designed to detect the Point of Incipient Failure through the analysis of real-world time-series data. By constructing a structural causal model and applying probabilistic counterfactual analysis, RaiLog RCA provides actionable insights that enhance the root cause discovery from identification of the time of appearance of anomalies and the associated causal graph. The paper presents the experimental results that demonstrate the accuracy performance of the proposed approach.

**Keywords:** Root Cause Analysis, Probabilistic inference, Failure prediction, Railway

## 1 Introduction

Root Cause Analysis (RCA) is a troubleshooting method of problem solving used for identifying the root causes of faults or failures [22]. RCA is a form of deductive inference since it requires an understanding of the underlying causal mechanisms for the potential root causes and the problem, i.e., what is typically found in the context of Predictive Maintenance through the Failure Mode, Mechanism, and Effect Analysis documentation. RCA can be decomposed into four steps:

1. *Identify* and *describe* the problem clearly.
2. Establish a *timeline* from the normal situation until the failure finally occurs, through the *Point of Incipient Failure* [24].
3. Distinguish between the *root cause* and other causal factors.
4. Establish a *causal graph* between the root cause and the observed problem.

The trigger signal of the RCA is given by the failure timestamp (i.e., the point in time when the failure variable is observed). Then, RCA yields a list of potential root cause variables along with their probabilities, which aligns with the way complex systems fail [6]. The variables that comprise the data are required

to be representative enough to help the developers and engineers pinpoint the source of the observed problems through the root causes and their effects [21].

This work presents RaiLog RCA, a multi-level diagnosis method that signals railway systems anomalies and pinpoints their root causes through statistical inference. Initially, the proposed approach tackles the detection of anomalies by conducting a *coarse-grained* diagnosis to determine if the asset under test shows a normal or an abnormal working condition. Then, it explains the potential reasons why the strange behavior occurred. Finally, it exploits a probabilistic model and conducts a *fine-grained* diagnosis by determining the root cause of the observed anomaly.

## 2   RaiLog RCA: Causal Inference

### 2.1   Causal Graph Generation

The causal links among the variables $X$ that build the model of a system are assumed to be most effectively represented using the tools from the field of Causality. In this sense, the Structural Causal Model (SCM) is the framework that can most generally capture such directed associations [15]. The SCM defines a set of assignments governing their specific functional associations $f$, along with some independent noise $U$ that accounts for everything that is not explicitly included in the model:

$$X_j := f_j(PA_j, U_j) \ , \tag{1}$$

where $PA_j$ represents the direct causes of the $X_j$ variable.

If enough knowledge and experience from the field is available from the subject matter experts, i.e., strictly complying with the RCA requirements, then a complete SCM may be developed right from the start. However, this is not the typical use-case scenario in complex industrial settings, and data generally needs to be carefully leveraged to drive the development of the causal model.

Whenever the structure of the model is to be inferred from the observed variables, assumptions need to be made about the data generating process, constraints need to be applied, and usually the statistical methods of the algorithms yield different graphs that explain the same factual data [7].

In a multivariate environment, the most straightforward approach is led by the so-called "constraint-based" discovery methods. These traditional approaches iteratively build the causal graph by utilizing a score such as the $p$-value of conditional independence tests. As a general technique, the Peter-Clark (PC) algorithm is described [19]. PC is a causal network structure learning algorithm that copes well with high dimensionality and can often also identify the direction of contemporaneous links [17]. It is consistent under i.i.d. sampling assuming no latent confounders, i.e., common causes, so that all relevant variables need to be observed in the data. Its outcome is a Markov Equivalence Class, and thus it is likely to have different graphical representations that explain the same observed data. The PC algorithm is especially suited to discover causality in combination with the Fisher-Z independence test because it requires less constraints for the input data [10].

**Dynamic Causal Bayesian Network Model** The proposed approach initially infers the causal relations from observational time series event data. Nevertheless, no family or method for causal discovery in time series stands out in all situations with different characteristics [2].

An initial baseline is obtained with the PC algorithm (using the Fisher-Z independence test) on data augmented with time lags. The event count-based transformation naturally lends itself to the application of this technique as long as the counts approximate a Gaussian distribution, which can be asserted using the Lilliefors normality test [11].

However, the direct application of PC discovery may not be advised for certain time series cases, and other more involved methods using more powerful statistical tests (including explicit time lags) should be explored on top of it. In consequence, the Momentary Conditional Independence (PC-MCI) test is considered [18], which has a stronger causal detection ability based on partial correlation tests.

Once the structural graph that binds the variables is determined, the functional associations of the SCM may be learned, and this work adopts a stochastic interpretation of the world. Therefore, it treats all $X(t)$ as random variables, and the resulting SCM statistically describes their (conditional) probability distributions. In this sense, Dynamic Causal Bayesian Networks (DCBN) are generative model that yield a factorized representation of a stochastic process. They represent a probability distribution over the possible histories of a time-invariant process; their advantage with respect to classical probabilistic temporal models like a Markov chain is that a DCBN is a stochastic transition model factored over a number of random variables, over which a set of conditional dependency assumptions is defined [4].

Considering $n$ time-dependent discrete random variables $X_1^t, X_2^t, ..., X_n^t$ and a directed acyclic graph that relates them causally, a DCBN is essentially their replication over time slices $t-\Delta$ (creating the so-called discretization steps), with the addition of a set of arcs in the graph representing the transition model, which is defined through the distribution $P(X_i^t|X_j^{t-\Delta})$, for all time-related variables $i$ and $j$. Arcs connecting nodes at different time-slices ($\Delta > 0$) are called interslice edges, while arcs connecting nodes at the same slice ($\Delta = 0$) are called intraslice edges. The joint probability distribution of the DCBN is shown as follows:

$$P(\mathbf{X}) = \prod_{\forall j} P\left(X_j^{t-\Delta}|PA_j^{t-\Delta}\right) \ . \tag{2}$$

The graphical nature of such Bayesian networks allows seeing relationships among different variables, and their conditional dependencies enable performing probabilistic inference [1]. Specifically, DCBN are powerful tools for knowledge representation and inference under *uncertainty* [16].

**Incipient Failure Prediction** The learned DCBN shall be used to estimate the probability of the Failure variable $X_F$ in time $P_F(t)$, which is the sink node

in the model that represents the eventual system crash, given the observed data (i.e., the root causes and their effects):

$$P_F(t) = P(X_F^t | PA_F^t) \ . \tag{3}$$

Ideally, the probability of observing a high count of failure events should be a monotonically increasing function (in time) until the moment of system failure.

To detect if an anomaly is present, the Point of Incipient Failure $T$ should be determined. This is the moment in time when the system *starts* developing an abnormal behavior that will eventually lead to the crash. Also, this is where the root cause of the observed anomaly is reasonably expected to be found. A possible strategy to determine this instant can be defined by the minimum-time significant-second-derivative of the probability of Failure, as this is the *first* inflection point with a *minimally relevant* increase $\Theta$ of risk (it may not be the greatest absolute increase, but it shall be one with the precedence in time):

$$T = \min_t \left( \frac{\partial^2}{\partial t^2} P_F(t) > \Theta \right) \ . \tag{4}$$

While there may be many different ways to express this criterion, by using a $\Theta$ threshold parameter the subject matter experts can be easily involved in the design of the solution.

## 2.2   Path Likelihood Estimation

The hypothesis of isolation is a methodological requirement of the sciences for research; hence, the useful fiction of the isolated "causal chain" or "singled-out path" in the structure will work to the extent to which such an isolation takes place, and this is often the case in definite respects during limited intervals of time. Moreover, since every isolable process is causal, anomalies can emerge solely as a result of external perturbations [5].

Concerning the analysis of a DCBN for RCA, estimating the most likely time-sequence chain of variables for the observed anomaly event adds explanatory value in an industrial environment. In the DCBN, each node represents an event count or state change of a variable, and the arcs represent causal–temporal relationships between the nodes. In this setting, probabilistic temporal logic determines that causes and effects are steady state formulas, the properties of which hold for the system at a certain point in time [20], and this allows for each formula to be a path formula too where multiple variables are involved. Therefore, the causal paths shall be given by the structure of the graph: a search algorithm shall be used to traverse it and find all the routes $S$ from the different root nodes to the sink Failure node.

For the Point of Incipient Failure $T$, the most likely causal path $S^*$ that explains the anomaly data can be determined after the exhaustive search among all the potential paths $S$ and their respective probabilities:

$$S^* = \max_{s \in S} P(s|\hat{s}); \ t = T \ , \tag{5}$$

where $s$ represents a structural path from a source node to the sink node (i.e., the failure event variable).

Conditioning on the variables not in the path under analysis ($\hat{s}$) is important to block spurious associations. This is especially relevant in the case of descendants, because in the event of an anomaly, the parent/ancestor variables are preferred as precedents [12].

Finally, in addition to putting the focus on the most expected behavior, one could argue that the root cause may also have occurred in the most unexpected/irregular setting [23], assuming that the most commonly experienced issues will have already been solved. This alternative perspective may also be covered in the proposed approach by minimizing the path likelihood probability.

## 2.3   Causal Inference

Beyond probabilistic inference, Causal Inference provides the tools that allow estimating causal conclusions even in the absence of a true experiment, given that certain assumptions are fulfilled. These assumptions increase in strength as is defined in Pearl's Causal Hierarchy (PCH) abstraction [3], i.e., Associational, Interventional, and Counterfactual, which is summarized as follows.

At the bottom of the hierarchy there's the Associational Layer, which describes the observational distribution of the factual data through their joint probability function $P(X)$. From this point forward, interesting quantities, i.e., the queries $X_Q$, can be directly computed given some evidence $X_E$, through their conditional probability:

$$P(X_Q|X_E) = \frac{P(X_Q, X_E)}{P(X_E)} \ . \tag{6}$$

This level of analysis displays a degree sophistication akin to classical (un)supervised Machine Learning techniques. As such, it is subject to confounding bias.

**Interventional Analysis** The Interventional Layer describes an actionable distribution, which endows causal information at the population level, i.e., to better understand the general behavior of the system. This level of analysis can be achieved through actual experimentation via Randomized Control Trials, or through statistical adjustments that smartly combine observed conditional probabilities to reduce spurious associations in the estimation. Pearl's *do*-calculus is likely to be the most effective approach to determine the identifiability of causal effects by applying the following three rules: 1) insertion/deletion of observations, 2) action/observation exchange, and 3) insertion/deletion of actions [13].

**Counterfactual Analysis** Finally, the Counterfactual Layer at the top of the hierarchy describes a potential distribution (possibly at the individual failure level) driven by hypothetical speculations over data that may contradict the facts. Conducting this estimation requires the following three steps [14]:

1. **Abduction**: Beliefs about the world are initially updated by taking into account all the evidence $E$ given in the context. Formally, the exogenous noise probability distributions $P(U)$ in the SCM are updated to $P(U|E)$.
2. **Action**: Interventions are then conducted to reflect the counterfactual assumptions, and a new causal model is thus created.
3. **Prediction**: Finally, counterfactual reasoning occurs over the new model using the updated knowledge.

To systematically explore these counterfactual worlds, Algorithmic Recourse applies a specific goal-driven rationale[9], where such environments are simulated via inference through (atomic) interventions $\alpha$ in time on a specific abnormal instance in order to revert the anomaly [8], i.e., to lower the risk of failure $X_F$. This is expected to help in the recognition and understanding of the general root causes that lead to the system failure [12].

Formally, the specific retrospective reasoning that these counterfactuals explore on the anomaly, i.e., the Point of Incipient Failure at $t = T$, can be stated as:

$$P(X_F^{t=T} = L | do(X^{t=T} = \alpha), X^{t=T}, X_F^{t=T} = H) . \tag{7}$$

Given that an anomaly was factually recorded in the data, i.e., through observing a high risk of failure $X_F^{t=T} = H$, i.e., a high count of failure/alarm events, Equation ([**?**]) estimates the probability that the risk would have been low at the Point of Incipient Failure $X_F^{t=T} = L$, had the root cause $X^{t=T}$ had the value $\alpha$, instead of the value it actually had when the anomaly was triggered. Note that this formula does not involve regular probabilistic conditioning, but the application of the Abduction-Action-Prediction method.

## 3  Experimental results and discussion

The proposed approach is primarily targeted at the railway domain, where it aims to address various operational challenges and improve overall system efficiency and reliability.

### 3.1  Data engineering pipeline : Automated Log processing

The proposed solution aims to process large machine logs in order to provide insights into the anomalous behaviors of a railway asset. It leverages two types of data: Event Variables (EV) High level, nominal qualitative data which group functions into categories, such as subsystem events. State Variables (SV) Low level, parametric quantitative or numeric data that show some kind of meaningful order or hierarchy, such as physical sensor records. The two types of data can be related, i.e., synchronized, by their timestamps.

The data engineering pipeline includes data preparation, characterization and scoring. Data preparation represents a preprocessing step that transforms raw data into usable information. The obtained structured data undergoes a transformation to be standardized into a time-series format by generating time
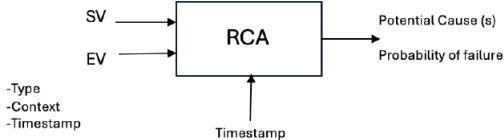
**Fig. 1.** RCA functional module diagram

windows. The training dataset is divided into time windows based on failure pattern occurrences. Each window spans 5 hours prior to a failure event. Each window is further divided into bins of 2 minutes each. The last bin of each window always contains the failure pattern.
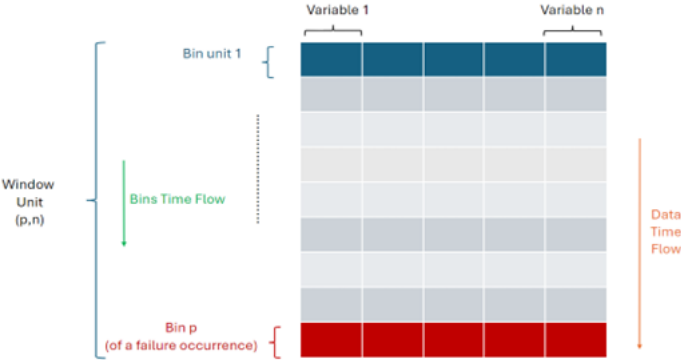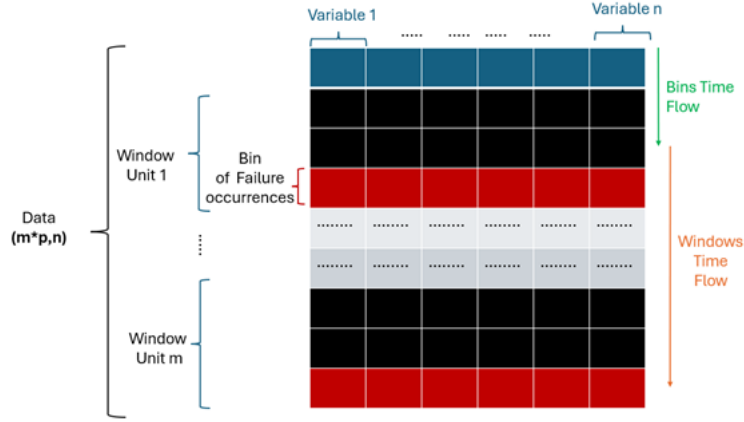


**Fig. 2.** Single window

Figure 1 shows a single window where rows represent bins and columns represent failure patterns of interest. The cell values correpond to the binary transformation by counting the variable messages occurrences within each bin therefore obtaining an integer-valued representation for all the variables. Figure 2 displays all windows combined as a tensor.

To effectively manage the expansive event-variable space, it is recommended to filter relevant variables using Mutual Information measures and independence testing, while also removing unrelated periodic events, resulting in a refined set of significant integer-valued time series variables that reflect event count evolution.

**Fig. 3.** Combined windows : Tensor

### 3.2   Causal discovery results

The features exploited for causal discovery are the 119 event variables resulting from windowing along the timestamps present in the windowed data.

since, the PC algorithm ignores the time factor, it ignores caus_and_effect relationships that unfold over time. In our work, and in order to compute the causal time series analysis, the *PC-MCI* [18] algorithm was implemented using the following hyper-parameters:

**Table 1.** Hyperparameters used in the PC-MCI algorithm

| Parameter | Value |
|---|---|
| dataframe | df |
| cond_ind_test | CMIsymb |
| significance | 'fixed_thres' |
| n_symbs | 3 |
| verbosity | 1 |
| tau_min | 0 |
| tau_max | 5 |
| pc_alpha | 0.01 |

MMI_COMMUNICATION_ERROR_IN_CAB_A_APPLI_ERROR

2,3,4

1

1

DRM_FUNCTION_SCREEN_1_FAILURE
1,2,3,4,5

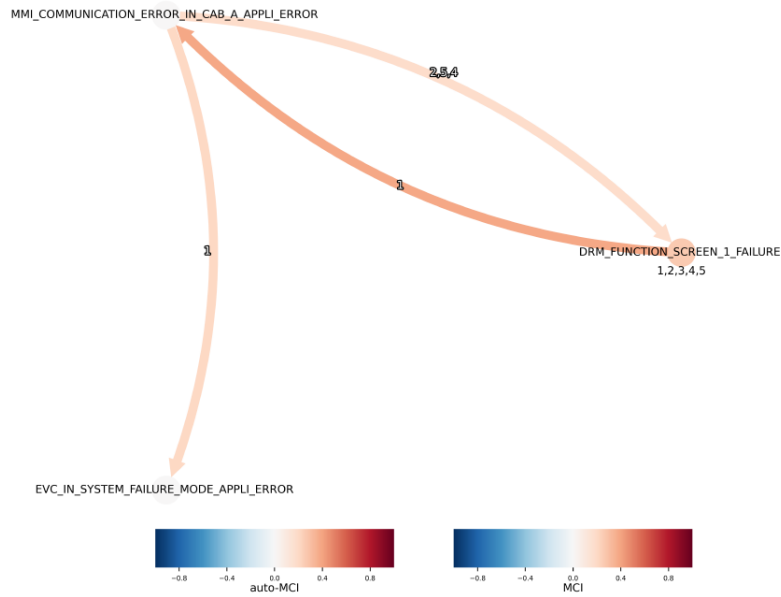EVC_IN_SYSTEM_FAILURE_MODE_APPLI_ERROR

auto-MCI

MCI

**Fig. 4.** Causal discovery on an excert from the event failure events

## 4    conclusion

This research paper introduces a counterfactual Root Cause Analysis (RCA) approach grounded in the principles of causal inference, adhering to railway industrial development standards. By utilizing multivariate time series data, the study centers on the critical task of identifying the Point of Incipient Failure. The findings suggest that this temporal milestone is pivotal for effective root cause identification, as it is likely to harbor the origin of the fault. Expanding the analysis to include various types of multivariate time series data could enhance the robustness of the findings and allow for comparisons across different industrial applications. Additionally, examining the integration of machine learning techniques with the counterfactual methodology may improve the detection and prediction of the Point of Incipient Failure, allowing for more proactive measures to address potential system failures.

## References

1. Alaeddini, A., and Dogan, I. Using Bayesian networks for root cause analysis in statistical process control. Expert Systems with Applications, 38:11230–11243, 2011.
2. Assaad, C. K., Devijver, E., and Gaussier, E. Survey and Evaluation of Causal Discovery Methods for Time Series. Journal of Artificial Intelligence Research, 73:767–819, 2022.
3. Bareinboim, E., Correa, J. D., Ibeling, D., and Icard, T. On Pearl's Hierarchy and the Foundations of Causal Inference. Probabilistic and Causal Inference: The Works of Judea Pearl, pages 507–556, 2022.
4. Bobbio, A., Codetta-Raiteri, D., Montani, S., and Portinale, L. Reliability analysis of systems with dynamic dependencies. Bayesian Networks: A Practical Guide to Applications, pages 225–238, 2008.
5. Bunge, M. Causality and Modern Science. Routledge, 2009.
6. Cook, R. I. How Complex Systems Fail. Cognitive Technologies Labratory, University of Chicago, 2000.
7. Glymour, C., Zhang, K., and Spirtes, P. Review of Causal Discovery Methods Based on Graphical Models. Frontiers in Genetics, 10(524):1–15, 2019.
8. Han, X., Zhang, L., Wu, Y., and Yuan, S. On Root Cause Localization and Anomaly Mitigation through Causal Inference. Proc. of the 32nd ACM International Conference on Information and Knowledge Management, 2023.
9. Karimi, A.-H., Barthe, G., Schölkopf, B., and Valera, I. A Survey of Algorithmic Recourse: Contrastive Explanations and Consequential Recommendations. ACM Computing Surveys, 55(5):95, 2022.
10. Kobayashi, S., Otomo, K., Fukuda, K., and Esaki, H. Mining causality of network events in log data. IEEE Transactions on Network and Service Management, 15(1):53–67, 2017.
11. Lilliefors, H. W. On the Kolmogorov-Smirnov Test for Normality with Mean and Variance Unknown. Journal of the American Statistical Association, 62(318):399–402, 1967.
12. Li, M., Li, Z., Yin, K., Nie, X., Zhang, W., Sui, K., Pei, D. Causal Inference-Based Root Cause Analysis for Online Service Systems with Intervention Recognition. Proc. of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pages 3230–3240, 2022.
13. Pearl, J. The Do-Calculus Revisited. Proc. of the 28th Conference on Uncertainty in Artificial Intelligence, pages 4–11, 2012.
14. Pearl, J., Glymour, M., and Jewell, N. P. Causal Inference in Statistics: A Primer. John Wiley and Sons Ltd, 2016.
15. Pearl, J. The seven tools of causal inference, with reflections on machine learning. Communications of the ACM, 62 (3):54–60, 2019.
16. Pourret, O. Introduction to Bayesian networks. Bayesian Networks: A Practical Guide to Applications, pages 1–13, 2008.
17. Runge, J., Bathiany, S., Bollt, E. et al. Inferring causation from time series in Earth system sciences. Nature Communications, 10(2553):1–13, 2019.
18. Runge, J., Nowack, P., Kretschmer, M., Flaxman, S., and Sejdinovic, D. Detecting and quantifying causal associations in large nonlinear time series datasets. Science Advances, 5(eaau4996):1–15, 2019.
19. Spirtes, P., Glymour, C., and Scheines, R. Causation, Prediction, and Search. MIT Press, 2001.

20. Van Houdt, G., Depaire, B., and Martin, N. Root Cause Analysis in Process Mining with Probabilistic Temporal Logic. Munoz-Gama, J., Lu, X. (eds) Process Mining Workshops. ICPM 2021. Lecture Notes in Business Information Processing, 433:73–84, 2022.
21. Weidl, G., Madsen, A. L., and Dahlquist, E. Decision support on complex industrial process operation. Bayesian Networks: A Practical Guide to Applications, pages 313– 328, 2008.
22. Wilson, P. F., Dell, L. D., and Anderson, G. F. Root Cause Analysis: A Tool for Total Quality Management. ASQ Quality Press, 1993.
23. Yang, W., Zhang, K., and Hoi, S. C.H. A Causal Approach to Detecting Multivariate Time-series Anomalies and Root Causes. Proc. of the International Conference on Learning Representations, 2023.
24. Trilla, A., Rajendran, R., Yiboe, O., Possamaï, Q., Mijatovic, N., & Vitrià, J. Industrial-Grade Time-Dependent Counterfactual Root Cause Analysis through the Unanticipated Point of Incipient Failure: a Proof of Concept. Proc. of the Causal Inference for Time Series Data Workshop at the 40th Conference on Uncertainty in Artificial Intelligence, 2024.