

563 Appendix

564 Here, we provide the detail about the following topics;

- 565 • Preliminaries about microservice, causality, and causal graph
- 566 • Correctness proof of RCD
- 567 • RCD for finite samples
- 568 • Implementation details
- 569 • Sensitivity analysis of γ
- 570 • Sock-shop experiment
- 571 • Related Work

572 Preliminaries

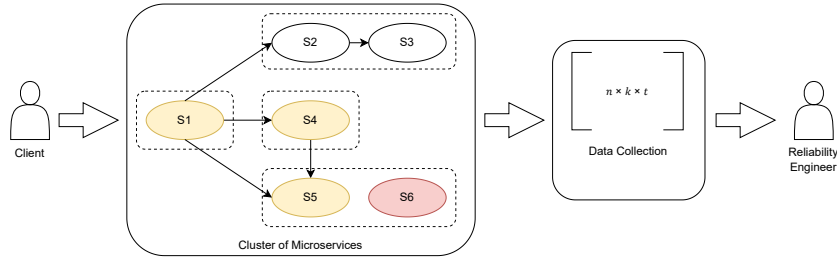


Figure 5: The workflow of a microservice-based system from the usage point of a client and SRE . Every oval is a microservice whereas every dashed rectangle is a host that hosts multiple microservices. The edges between two microservices represent the call graph. Red-colored microservices signify a failure but the effect is observed on the yellow-colored microservices. Even with the topological information about the microservices, it will be difficult for a Reliability Engineer to pinpoint the root cause just observing the failures at, say s_5 .

573 **Microservice Architecture.** Microservices are usually containerized applications that perform a
 574 small but specific set of tasks independently from other microservices. In a microservice-based
 575 application, a set of microservices work together to fulfill a client’s request. For monitoring and
 576 debugging purposes, all of the microservices expose some kind of performance metrics. There
 577 are many tools that help developers instrument the microservices to expose useful information as
 578 well as collect, store and visualize all that data [39, 33, 15, 61]. Most of these tools scrap all the
 579 microservices periodically, store the metric information, and produce a time-series dataset. SREs use
 580 this time-series metric data to monitor the performance of the overall application as well as the health
 581 of individual microservices. Figure 5 depicts a general web application with multiple microservices
 582 running on different hosts.

583 More formally, we define the microservice architecture as follows, a web application consists of
 584 a set of n microservices, $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$. At any given time, the monitoring tool collects m
 585 metrics from each of the microservices, i.e., $\mathcal{M}(i, t) = \{r_{i,1,t}, r_{i,2,t}, \dots, r_{i,m,t}\}$ is a set of m
 586 metrics of microservice i at time instance t . To combine this all together, we have a time series dataset,
 587 $\mathcal{D} = \{\mathcal{M}(1, 1), \mathcal{M}(1, 2), \dots, \mathcal{M}(n, t)\}$, that contains the metrics data of every microservice till time
 588 t .

589 **Causality and Causal Graph.** We define causality as a means to measure the dependence between
 590 two variables. One way to depict the causality between these variables is the causal graph. In the
 591 causal graph, every node represents a random variable whereas an edge between the two nodes shows
 592 the dependence between the variables. The conditional probability distribution of a node in the causal
 593 graph can be computed as

$$P(X_i | pa_i) \tag{1}$$

594 Here X_i represents the random variable i in the causal graph and pa_i is the set of parents. We can
 595 use 1 to compute the marginal probability distribution of a node by summing over the probability of
 596 all the values of pa_i .

$$P(X_i) = \sum_x P(X_i | x) P(x) \tag{2}$$

597 Moreover, the probability of the overall system with n random variables can be computed as

$$P(X_1, \dots, X_n) = \prod_i P(X_i | pa_i) \quad (3)$$

598 We can use **1** as the conditional independence (CI) test to check the (in)dependence relationship
 599 between two variables. For instance, two variables X and Y are independent only if we can find a set
 600 $S \cap \{X, Y\} = \emptyset$ such that $P(X|Y, S) = P(X|S)$.

601 Formally, a causal graph can be defined as a DAG $G(V, E)$ where V is the set of n vertices,
 602 $V = \{v_1, v_2, \dots, v_n\}$, and $E \subseteq V \times V$ is the set of the directed and undirected edges. A directed edge
 603 $\vec{e} = \{v_1, v_2\}$ signifies that v_1 causes v_2 . Here the edge \vec{e} shows that we could not find a *separating*
 604 *set*, S , such that $P(v_1|v_2, S) = P(v_1|S)$.

605 **Proof of Theorem 1.**

606 Note that Ψ -PC concatenates two datasets with an F-NODE and runs PC algorithm on all the variables
 607 by treating F-NODE as another variable. Even though with more than two datasets this would create
 608 faithfulness violations among multiple F-NODE's as pointed out in [18, 30], no such issue arises
 609 when there are two datasets creating only a single F-NODE. Accordingly, Ψ -PC can avoid the careful
 610 treatment between F-NODE's that Ψ -FCI requires. Same for RCD which uses Ψ -PC

611 We will use p_N to represent the probability distribution of variables under normal operating conditions
 612 (observational), and p_A to represent the one under anomaly (interventional). Let us introduce the
 613 probability distribution p^* that is defined as $p^*(V|F = 0) = p_N(V)$ and $p^*(V|F = 1) = p_A(V)$,
 614 where V is the set of observed variables and F is the F-NODE representing the effect of intervention.

615 First, note that if Ψ -PC was run on the full graph, it would output a graph where the F-NODE is
 616 adjacent to only the true root causes. This is because any variable S that is not a root cause can be
 617 separated from the F-NODE as $S \perp\!\!\!\perp F | Pa_S$, since $p_N(s|pa_s, F = 0) = p_A(s|pa_s, F = 1)$ and
 618 Pa_S is observable for all S due to the causal sufficiency assumption.

619 RCD algorithm runs Ψ -PC on a subset of variables in each stage. F-NODE will then remain adjacent to
 620 those that are not separable from it by conditioning on any subset of the variables in that subset. Due
 621 to the extended faithfulness assumption, we have that $p_N(s|u, F = 0) \neq p_A(s|u, F = 1)$ for any root
 622 cause node S , and for any subset U of the observed variables, or equivalently $S \not\perp\!\!\!\perp F | U$. Therefore,
 623 root-causes cannot be separated using any subset of the observed variables. This establishes that,
 624 whichever subset the true root causes fall into, they will all remain adjacent to the F-NODE of that
 625 subset. This is true for all stages of the algorithm. This establishes the soundness of the algorithm,
 626 that at any stage, including the final stage, root-causes will remain adjacent to the F-NODE.

627 RCD however is not complete.⁴ We have the following simple example to demonstrate that the
 628 returned set might contain additional nodes that are not root causes. Consider the augmented graph
 629 $F \rightarrow R \rightarrow S, T \rightarrow R, T \rightarrow S$. Suppose we partition the nodes into subsets $\{\{T\}, \{R, S\}\}$. F-NODE
 630 will not be adjacent to T since $F \perp\!\!\!\perp T$. For the second subset, F-NODE will be adjacent to both
 631 R, S . This is because one needs to condition on both R and T in order to d-separate F-NODE from S .
 632 However, T was removed in the first subset and is not included in the next stage.

633 Please note that one can easily resolve this issue by running one more stage at the end of RCD
 634 algorithm by including all the nodes $\{U : F \perp\!\!\!\perp U\}$. These are the nodes that are non-descendants of
 635 the F-NODE, which contains the parents of the root causes. Including them enables finding a valid
 636 separating set for any node that is not a root cause. This is because R, Pa_R blocks all the backdoor
 637 and frontdoor paths from the root causes that would d-connect F-NODE and any non root cause node.

638 In our experiments we only implemented the vanilla RCD which is sound but not complete as it was
 639 sufficient to achieve the precision of Ψ -PC. Furthermore, existence of false positives is not a problem
 640 for root cause analysis application as long as most of the variables are eliminated quickly as potential
 641 root causes in an automated manner.

Algorithm 2 Ψ -PC: Algorithm for learning the causal graph and interventional targets from observational and interventional probability distributions [18]

Input: Observational distribution P_N , interventional distribution P_A , and set of variables V
Output: A causal graph on \mathcal{G} .

- 1: **procedure** Ψ -PC(P_N, P_A, V)
- 2: $P^*(V|\text{F-NODE} = 0) \leftarrow P_N(V), P^*(V|\text{F-NODE} = 1) \leftarrow P_A(V)$ # Concatenate
- 3: $P^*(\text{F-NODE} = 0) \leftarrow 0.5, P^*(\text{F-NODE} = 1) \leftarrow 0.5$ # Equal number of samples
- 4: **Phase-1: Learn Skeleton**
- 5: $\mathcal{G} \leftarrow$ a complete undirected graph over $V \cup \{\text{F-NODE}\}$
- 6: **for all** every pair $X, Y \in V \cup \{\text{F-NODE}\}$ **do**
- 7: **for all** $S \subseteq V \setminus \{X\}$ **do**
- 8: **if** $P^*(X|Y, S) = P^*(X|S)$ **then**
- 9: $\text{SepSet}(X, Y) \leftarrow S$ and remove the edge between X and Y from \mathcal{G}
- 10: **Phase-2: Edge Orientation**
- 11: \mathcal{R}_0 : For any triplets $\langle X, Z, Y \rangle$, such that $X \cap Y = \emptyset$, orient $X \rightarrow Z \leftarrow Y$ iff $Z \notin \text{SepSet}(X, Y)$
- 12: \mathcal{R}^* : Orient all adjacent edges of F-NODE as outgoing.
- 13: Apply the four Meek rules from [28] ($\mathcal{R}_1 - \mathcal{R}_4$) until none applies.
- 14: **return** \mathcal{G}

642 Hierarchical Learning Algorithm

643 In this section, we provide a modified version of RCD when the number of samples are finite. Similar
644 to Algorithm 1, the algorithm for finite number of samples is divided in two phases as well. In the first
645 phase, we divide the data in small subsets of size γ and run Ψ -PC to learn the potential interventional
646 targets. In the second phase, we take the union of all the potential interventional targets and narrow
647 the list down to k interventional targets. Complete pseudocode of this algorithm is provided in
648 Algorithm 3.

Algorithm 3 RCD for finite number of samples

Input: Normal dataset D , anomalous dataset D' , γ, α, k : Max. no. of root causes
Output: A list of root causes \mathcal{U} .

- 1: **procedure** RCD(D, D', γ, α, k)
- 2: $\mathcal{U} \leftarrow$ Set of variables of D, D' .
- 3: **while** $|\mathcal{U}| > k$ & $|\mathcal{U}| > \gamma$ **do**
- 4: $\mathcal{S} \leftarrow$ A random partitioning of $|\mathcal{U}|$ into subsets of size γ .
- 5: $R \leftarrow \emptyset$
- 6: **for all** $S \in \mathcal{S}$ **do**
- 7: $G \leftarrow \Psi$ -PC($D[S], D'[S], \alpha$) # Ψ -PC constructs a graph on $S \cup \{\text{F-NODE}\}$.
- 8: $R \leftarrow R \cup \text{Ne}_G(\text{F-NODE})$ # Extract neighbors of F-NODE.
- 9: $\mathcal{U} \leftarrow R$.
- 10: **return** TOPK($D[\mathcal{U}], D'[\mathcal{U}], k$) # Get top- k neighbors of F-NODE.

649 In most of the RCA literature, the output of the algorithm is an ordered list of potential root causes of
650 size k , rather than a singleton. Here $k \ll n$ and n is the total number of microservices. We create
651 an ordered list of potential root causes by running Ψ -PC multiple times with different α .

652 The Ψ -PC algorithm takes a hyperparameter α as an input which it uses this to compare the p-values
653 of CI tests to decide if two variables are independent. If the p-value is greater than α then we mark
654 those two variables as independent. In this way, a strict (small) value of α will result in a sparse graph
655 whereas a relaxed (large) value of α will produce a dense causal graph. We create an ordered list of
656 k potential root causes by running Ψ -PC multiple times with different values of α and ordering the
657 neighbors of F-NODE.

658 To create an ordered list of root causes, we modified the second phase of the hierarchical learning
659 algorithm. After running Ψ -PC on all the subsets, F-NODE may have more than one neighbor. In this

⁴In the main paper, there is a typo calling the algorithm sound and complete.

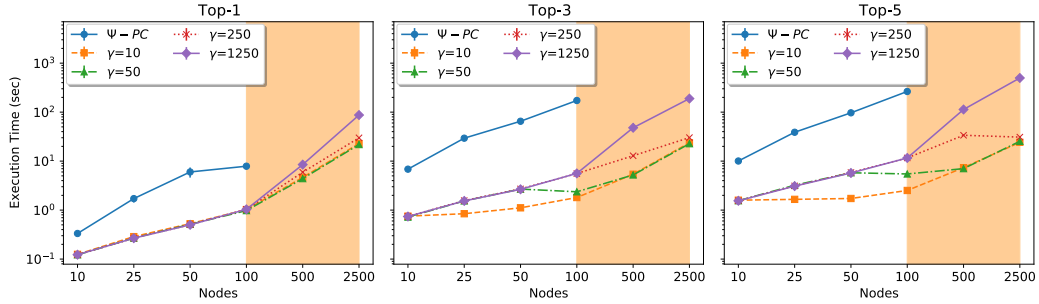


Figure 6: The execution time of Ψ -PC and RCD with different γ . It controls the size of the subsets for RCD, thus increasing γ leads to a higher execution time. However, even with higher γ and large number of nodes, RCD performs better than Ψ -PC because of its localized learning.

697 microservices to install `stress-ng` and then running it on the container of the targeted microservice.
 698 With this setup, we injected failures in all of the major microservices of Sock-shop (carts, catalogue,
 699 orders, payment, and user).

700 Sock-shop exposes a lot of metrics but not all of them are equally useful. Some of the metrics, such as
 701 the memory utilization of an inactive microservice, might remain constant during the experiment. A
 702 general solution for this problem is to remove constant variables from the data in the pre-processing
 703 step. We accomplish this by removing a metric that has a standard deviation below 1. Moreover, all of
 704 the metrics produce continuous values but the Chi-squared test only works on categorical data.
 705 Hence, we discretize the variables into a specific number of *bins* by using K-means clustering. For all
 706 our experiments with Sock-shop, we used 5 number of bins.

707 Sensitivity analysis of γ

708 RCD requires a hyperparameter, γ , which controls the size of each subset. In all of our experiments,
 709 we set γ to be 5. Here, we experiment by changing the value of γ to understand how it affects the
 710 performance of RCD. We run this experiment by controlling the number of nodes and the value of γ .
 711 The maximum in-degree of all the DAGs is 3, whereas the total number of states for every node is 6.
 712 For every experiment, we draw 10,000 samples for the normal and anomalous states of the system.
 713 We run every experiment 100 times and plot the average completion time of RCD. We also plot the
 714 execution time of Ψ -PC as the reference point. Figure 6 shows the result.

715 The first thing to note is that with a small number of nodes, different values of γ produce the result at
 716 about the same time. That is because when $n < \gamma$ where n is the number of nodes, the value of γ
 717 does not affect the number of subsets. In this case, there will only be one subset for the whole dataset.
 718 The effects of γ become more visible with a larger number of nodes and for the higher values of k .
 719 With the larger number of nodes, RCD creates more subsets if γ is small thus reducing the overall
 720 time. Even with higher values of γ , RCD performs better than Ψ -PC owing to its localized learning
 721 scheme. The time for higher values of k increases because we have to run multiple executions of
 722 Algorithm 4 to get a sorted list of all the neighbors of F-NODE.

723 Sock-shop Experiment

724 We ran the Sock-shop application for 5 minutes to collect the data for the normal state of the system.
 725 After that, we injected the failure and ran the application for 5 more minutes to collect the anomalous
 726 dataset. We repeated the experiment 5 times for two different types of failures (CPU hog and memory
 727 leak) and in total, we gathered 50 datasets. To measure the performance of the root cause detection
 728 algorithm, we ran an algorithm 50 times on every dataset and quantified their top- k precision and
 729 execution time. Table 1 shows the top- k precision and Table 2 shows the execution time of different
 730 RCA algorithms.

731 The first thing to note here is that ϵ -Diagnosis outperforms all the other algorithms. That is because it
 732 only uses a pairwise test of coefficient of variation to find the root cause metric and microservice.
 733 Another reason it takes much less time is that for Sock-shop, we only have 38 metrics. However, as

Table 2: Execution time in seconds of different algorithms for detecting the root cause on data from the Sock-shop application. ϵ -Diagnosis takes less time because of its reliance on coefficient of variation but it perfo

	CPU hog				Memory leak			
	Ψ -PC	AutoMAP	ϵ -Diagnosis	RCD	Ψ -PC	AutoMAP	ϵ -Diagnosis	RCD
Carts	8.966	17.038	0.007	0.548	17.052	14.243	0.007	0.543
Catalogue	9.410	6.139	0.006	0.393	10.512	14.572	0.005	0.326
Orders	2.964	14.361	0.004	0.739	1.972	14.579	0.007	0.659
Payment	12.403	19.832	0.005	0.670	30.857	11.807	0.006	0.490
Payment	1.338	17.259	0.005	0.490	3.763	11.645	0.006	0.697
Average	7.016	14.926	0.006	0.568	12.831	13.391	0.006	0.543

734 we show in our experiments, Figure 2, the precision of ϵ -Diagnosis drops significantly because of the
 735 weak CI test for a large number of variable.

736 Considering RCD, we observe that it outperforms Ψ -PC and AutoMAP. More specifically, it reduces
 737 the time to find the root cause by 91% and 96% as compared to Ψ -PC and AutoMAP in the case of
 738 CPU hog failures. We get this gain by using a hierarchical approach to learn the root cause and only
 739 learning the neighborhood of F-NODE rather than learning the complete causal graph.

740 Related Work

741 Performance diagnosis in a distributed system is of paramount importance due to the higher failure
 742 chances of each component, hence affecting the system in its entirety. There has been significant
 743 research work not only in academia [41, 8, 53] but in the industry [51, 52] as well devoted to the
 744 topics of identifying and localizing root causes for the system failures. Works on root cause analysis
 745 can broadly be segregated into the following methods:

746 **Log-based.** These works make use of system log information to investigate the causes of anomalies
 747 and failure of a system [12, 55, 60]. A general approach to log analysis involves parsing the raw text
 748 data to learn historical patterns, and leveraging the patterns learned to detect anomalies. [12] uses
 749 finite state automata to represent the patterns, whereas a recent study [60] performs a learning-based
 750 methodology on the features extracted from logs to detect anomalies. The log-based approaches
 751 assume that the application produces useful logs. However, in reality the quality and quantity of
 752 different modules of the same applications can be poles apart thus, making them useful only in limited
 753 space.

754 **Trace-based.** Studies like Pinpoint [7], X-Trace [11], and Microscope [22] record the execution path
 755 information and locate the culprit services. While Pinpoint and X-Trace require instrumenting the
 756 source code and a general understanding of the code by system administrators, Microscope aims
 757 at generating a service causal graph based on the trace data. Seer [13], on the other hand, uses
 758 RPC-level tracing along with low-level hardware monitoring to identify the root cause microservice,
 759 while [21, 23] leverages unsupervised trace anomaly detection and microservice localization to detect
 760 the root cause. Furthermore, GMTA [16] uses a graph-based approach to cluster the traces. It also
 761 builds an interpretable tool for understanding microservice architecture and diagnosing problems (it
 762 does not explicitly work on identifying the root causes). One common trait of most of these tools is
 763 that they require application instrumentation and expert knowledge therefore making them difficult to
 764 use in real-world applications.

765 **Metrics-based.** These tools aim at performing system diagnosis by leveraging the value of various
 766 metrics observed. MonitorRank [19] localizes the root cause API of failures by examining the
 767 historical time series of performance metrics along with the call graph of services. Grano [51]
 768 deployed at Ebay Inc. uses a dependency graph-based approach among physical resources to perform
 769 root cause analysis. With a call graph, we can analyze the faulty and the impacted services. However,
 770 to understand why a fault occurred, we need to understand how performance metrics for each
 771 component interact with each other. Works like [53, 54, 8, 29, 41] perform root cause analysis
 772 by building causal graphs using variations of the PC algorithm at the performance metrics level.
 773 The causal discovery algorithms output a CPDAG where some edges are undirected, and hence
 774 post-processing is necessary to convert the CPDAG to a DAG. After generating a DAG, the most
 775 common approach is to run a variety of random walk or graph traversal algorithms for ranking a
 776 subset of probable root causes. These tools do not require any application instrumentation or expert
 777 knowledge. However, they suffer from high execution times because of PC that tries to learn the
 778 complete underlying causal graph.

779 On the contrary, our proposed algorithm models the failures as an intervention on the failing node.
780 Thus, allowing us to leverage the state-of-the-art algorithm from causal inference to learn the
781 interventional target. Furthermore, we propose localized hierarchical learning to overcome the high
782 execution time of PC, therefore, making it applicable in real-world applications.