# Causal analysis of network logs with layered protocols and topology knowledge

**Satoru Kobayashi**, Kazuki Otomo, Kensuke Fukuda

CNSM 2019, Halifax

Oct 23, 2019

# Outline

- Background and research goal
- Approach
  - Introduction to causal analysis of network logs
  - Proposed method for using domain knowledge in causal analysis
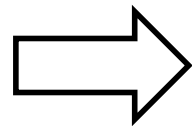- Evaluation
- Conclusion

# Difficulty of leveraging system log in network management

- Huge dataset
  - Large scale and complicated systems
  - 150,000 lines / day in SINET 5
  - Automated analysis required

- Difficulty in automated analysis
  - Free-format and sparse data
  - Contextual information required for troubleshooting

# Causal analysis in operational data

- Causal analysis: A popular approach for extracting contextual information
  - More reliable than correlation-based approach
- Problem:
  - Efficiency (large processing time)
  - No consideration of network knowledge

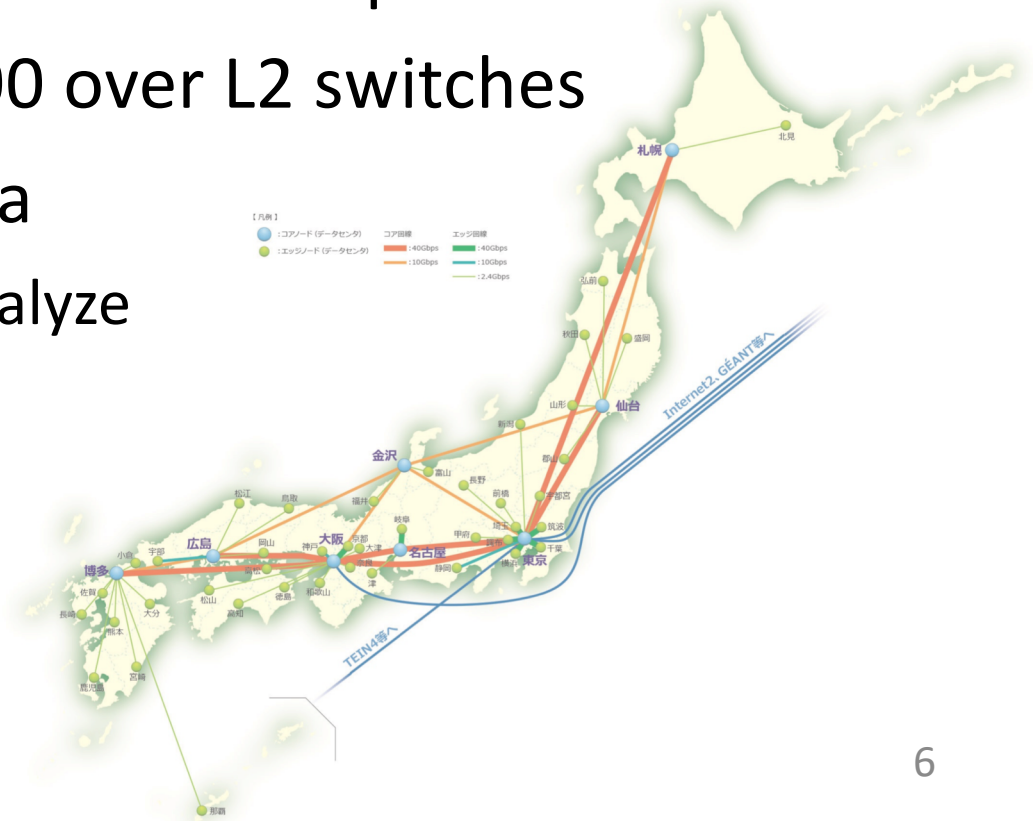$\Longrightarrow$ Causal analysis with network domain knowledge

# Goal

- Provide contextual information for system management and troubleshooting from network system logs
  - Causal analysis + Network domain knowledge
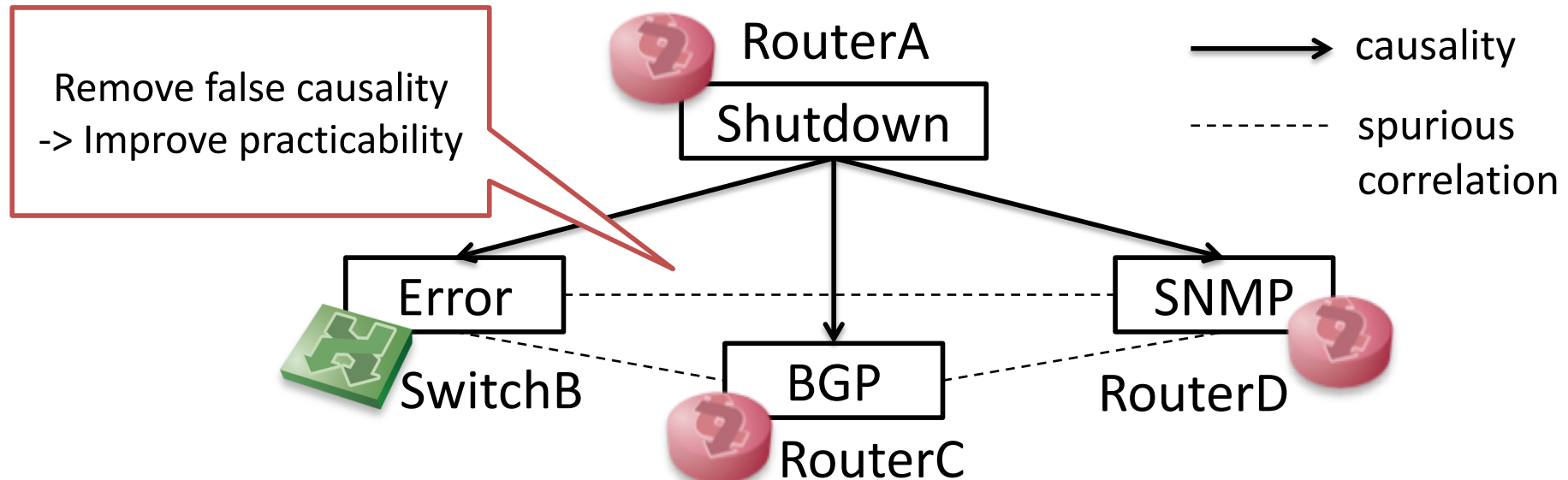  - Improve efficiency and reliability

# Dataset

- SINET4
  - https://www.sinet.ad.jp/en/top-en
  - A nation-wide R&E network in Japan
  - 8 core routers and 100 over L2 switches
  - 15 months syslog data
    - 3.5 million lines to analyze
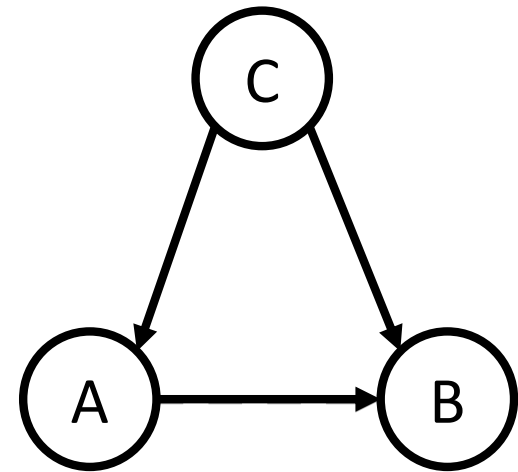
# Causal analysis of network logs[1]

Oct 17 17:00:00 routerA System shutdown by root
Oct 17 17:00:05 switchB Error detected on eth0
Oct 17 17:00:15 routerC BGP state changed from Established to Idle
Oct 17 17:00:15 routerD SNMP trap sent to routerA

......

Remove false causality
-> Improve practicability

RouterA
Shutdown

→ causality

------ spurious
correlation

Error
SwitchB

BGP
RouterC

SNMP
RouterD

[1] S. Kobayashi et al. "Mining causality of network events in log data", IEEE TNSM, vol. 15, no.1, pp. 37–67, 2018.

# Causal Inference
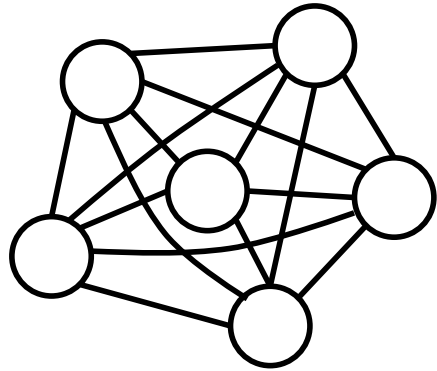
- Conditional Independence
  - A and B are independent if the effect of confounder C is excluded
  - A and B are conditionally independent given C

- PC algorithm [2]
  - Directed acyclic graph (DAG)
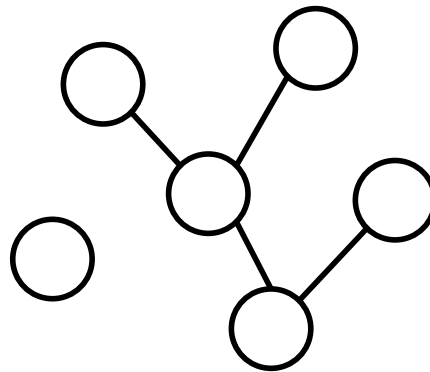  - Explore conditional independence and remove false edges



$$P(A|C)P(B|C) = P(A, B|C)$$

[2] P. Spirtes, et al. "Causation, prediction, and search". MIT press, 2000.

# Flow of PC algorithm

Complete graph (initial)

Skeleton graph

Directed acyclic graph



- Remove edges of conditional independence
- Statistical test for conditional independence ops)
  - G2 test (for binary or multi-level data) [3]
  - Fisher-Z test (for continuous data) [3]

[3] R. E. Neapolitan. "Learning Bayesian Networks." Prentice Hall Upper Saddle River, 2004.
[4] T. Verma, et al. "An algorithm for deciding if a set of observed independencies has a causal explanation". In Proceedings of UAI'92, pp. 323–330, 1992.

# Causal analysis with network logs [1]

Original log messages

Template generation

Time-series preprocessing

PC algorithm

Provided Information

Oct 23 13:00:25 sv1 interface eth1 down
Oct 23 13:00:26 rt2 connection failed to 192.168.1.4
**Oct 23 13:02:16 sv1 user sat logged in from 192.168.1.15**
Oct 23 13:02:29 sv1 su for root by sat
Oct 23 13:02:58 sv1 interface eth1 up
...

Original log

Event 1 : user ** logged in from **
2019-10-23 13:02:16
2019-10-23 14:25:00
...

Time-series event

Challenge:
How to use domain knowledge in causal analysis?

[1] S. Kobayashi et al. "Mining causality of network events in log data", IEEE TNSM, vol. 15, no.1, pp. 37–67, 2018.

# Approach: Pruning initial graph

- PC algorithm usually starts with complete graph
  - Takes large processing time if network structure is large and complex

- Prune edges in initial graph of PC algorithm
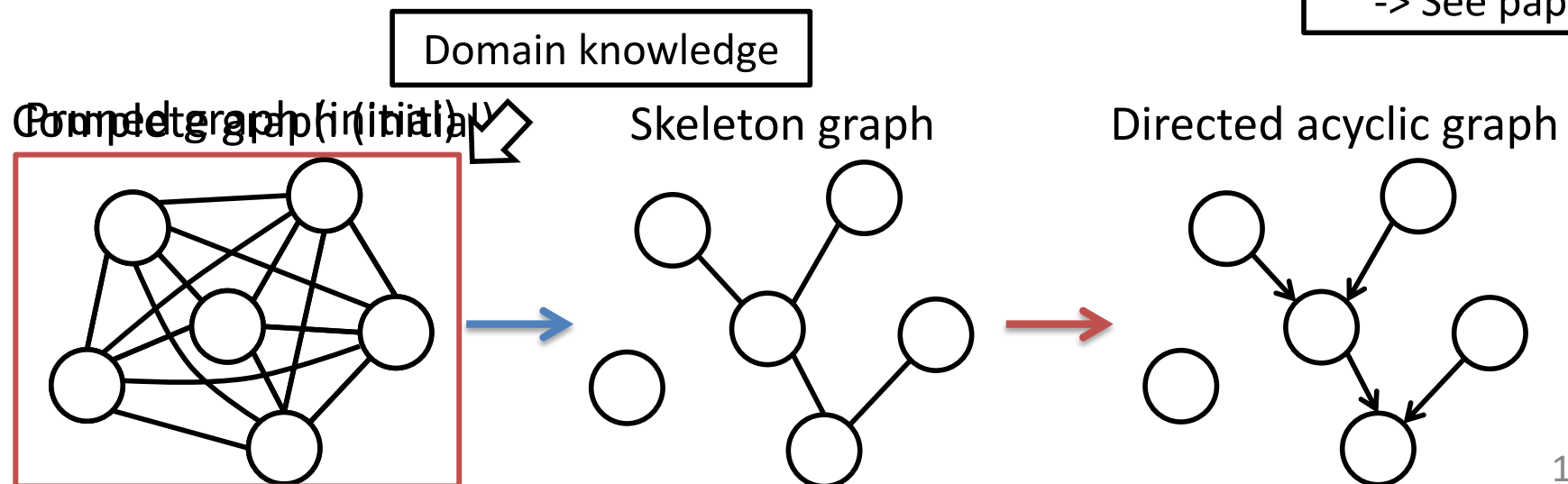  - Complete graph -> Pruned graph

Is it OK in theory?
-> See paper

Domain knowledge

Complete graph (initial)

Skeleton graph

Directed acyclic graph

# Pruning edges with domain knowledge

- Basic idea
  - Some edge candidates are clearly not causality
    - Compared with domain knowledge of operators
  - Ignore in calculating causality

# Difficulty in pruning

- **Unobserved events** mediate causality
  - Pruning mediated causality breaks causal flow

-> How to determine the criteria?

Root cause

Unobserved

Failure

**L2 switch C**          **Router D**

...                                                              ...

**Unobserved**

Interface event

Interface error

Routing event

Really unrelated?

# Proposed method: 2 criteria

**Rule 1**. Events **in same device**, or **in same functional layer and in connected devices**



**Rule 2**. A causal edge can be mediated with **1 (or 0) unobserved event**

# Example: Good causality candidate

L2 switch **A**        L2 connection        Router **B**

| | | |
|---|---|---|
| Layer 3 | | Routing event |
| Layer 2 | Interface event | ? |
| Others | | |

Need 1 hop to follow Rule 1

Edge candidate For causal analysis

L2 switch **A** Interface event — Router **B** Routing event

# Example: Bad causality candidate

L2 switch **A**          L2 connection          Router **B**

| | | |
|---|---|---|
| Layer 3 | | Routing event |
| Layer 2 | ? | ? |
| Others | Hardware event | |

Need **2 hops** to follow Rule 1

**Violating Rule 2 !**

Edge candidate

L2 switch **A** Hardware event   ✖   Router **B** Routing event

To be pruned

16

# Algorithm to classify causality candidate

- Keep a causal edge if satisfying 1 or 2
    1. 2 events appear in same device
    2. At least 1 end node (event) is on a functional layer that connects the devices

**Condition 2 Example**

L2 switch **A**

**L2 connection**

Router **B**

Layer 3

Routing event

Layer 2

Interface event

Keep this candidate!

**L2 event**

Others

17

# Analysis in SINET4 data
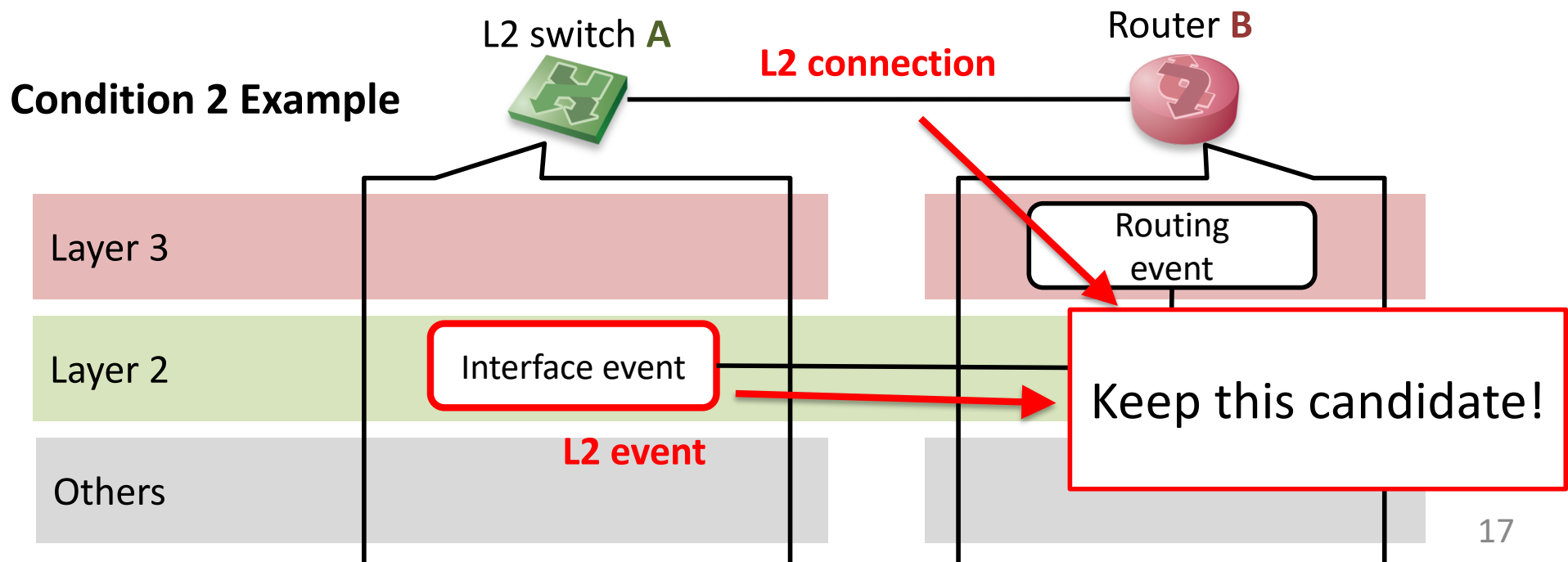
- ## Domain knowledge for pruning
  - – Network topology (L2, L3)
  - – Functional layer definition of events (L2, L3, others)
    - Manually labeled 9 classes for log templates
    - Layer definition for the classes ↓

| Layer definition | Event group (label) |
|---|---|
| L3 | Routing-EGP, Routing-IGP, VPN |
| L2 | Interface, Network |
| Others | System, Service, Management, Monitor |

# Evaluation

- Compare 3 methods (different initial graph)
  - <u>Processing time</u> & <u>Quality of edges</u>



Proposed method

| None | Area-based [1] | **Multi-layered** |

Complete graph

Multiple complete subgraph

Pruned

Domain knowledge

PC algorithm → DAG (×3)

[1] S. Kobayashi et al. "Mining causality of network events in log data", IEEE TNSM, vol. 15, no.1, pp. 37–67, 2018.

# Processing time of PC algorithm

Average processing time for 1-day data



Decreased 74% compared with None

16% faster than existing method

Proposed method

# Quality of causal edges

- Event classes of end nodes of detected edges

Multi-Layered method:
Same distribution with None

| Type | #Nodes | #Ends of edges | | |
|---|---|---|---|---|
| | | None | Area | ML |
| System | 49,005 | 24,577 | 23,033 | 22,662 |
| Network | 10,585 | 1,402 | 1,391 | 1,355 |
| Interface | 13,562 | 1,943 | 2,062 | 2,134 |
| Service | 7,697 | 742 | 435 | 314 |
| Mgmt | 81,628 | 29,379 | 27,911 | 26,332 |
| Monitor | 2,467 | 267 | 305 | 304 |
| VPN | 4,538 | 97 | 1,171 | 155 |
| Rt-EGP | 4,738 | 1,923 | 2,063 | 2,063 |
| Rt-IGP | 870 | 18 | 19 | 17 |
| Total | 175,090 | 60,348 | 58,390 | 55,336 |

# Quality of causal edges

- Event classes of end nodes of detected edges



End node

Data center

Core router

Area

**Multi-Layered** method:
Same distribution with None

**Area-based** method
has a problem:
**Relations among multiple
areas are missed**

⇩

Fails to detect
conditional independence

⇩

**False Positive** edges

34
14
32
04
55
63
17
36

# Summary of evaluation

| Pruning methods | Processing time | Quality of edges |
|---|---|---|
| None | ✕<br>Takes 10 minutes / day | ◯<br>(Shown in previous paper [1]) |
| Area-based method | ◯<br>Decrease 69% | ✕<br>No consideration of area gaps |
| Multi-Layered method (proposed method) | ◎<br>Decrease 74% | ◯<br>Similar distribution to None |

[1] S. Kobayashi et al. "Mining causality of network events in log data", IEEE TNSM, vol. 15, no.1, pp. 37–67, 2018.

# Discussion

- ## Parallel processing?
  - Available in PC algorithm [5]

- ## Available in other causal algorithms?
  - Depends on algorithms
  - Easily available in regression-based methods or constraint-based causal methods

- ## Available in any network?
  - Effective even in full-mesh-topology network

[5] T. Le, et al. "A fast PC algorithm for high dimensional causal discovery with multi-core PCs," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 13, no. 9, pp. 1–13, 2014.

# Conclusion

- Causal inference approach with network domain knowledge for helping troubleshooting

- Pruning initial graph of PC algorithm
  - Considering unobserved events

- Improvement in terms of processing time and quality of edges
  - Decrease 74%, 16% faster than Area-based method
  - Solve area-gap problem in Area-based method

- https://github.com/cpflat/logdag