



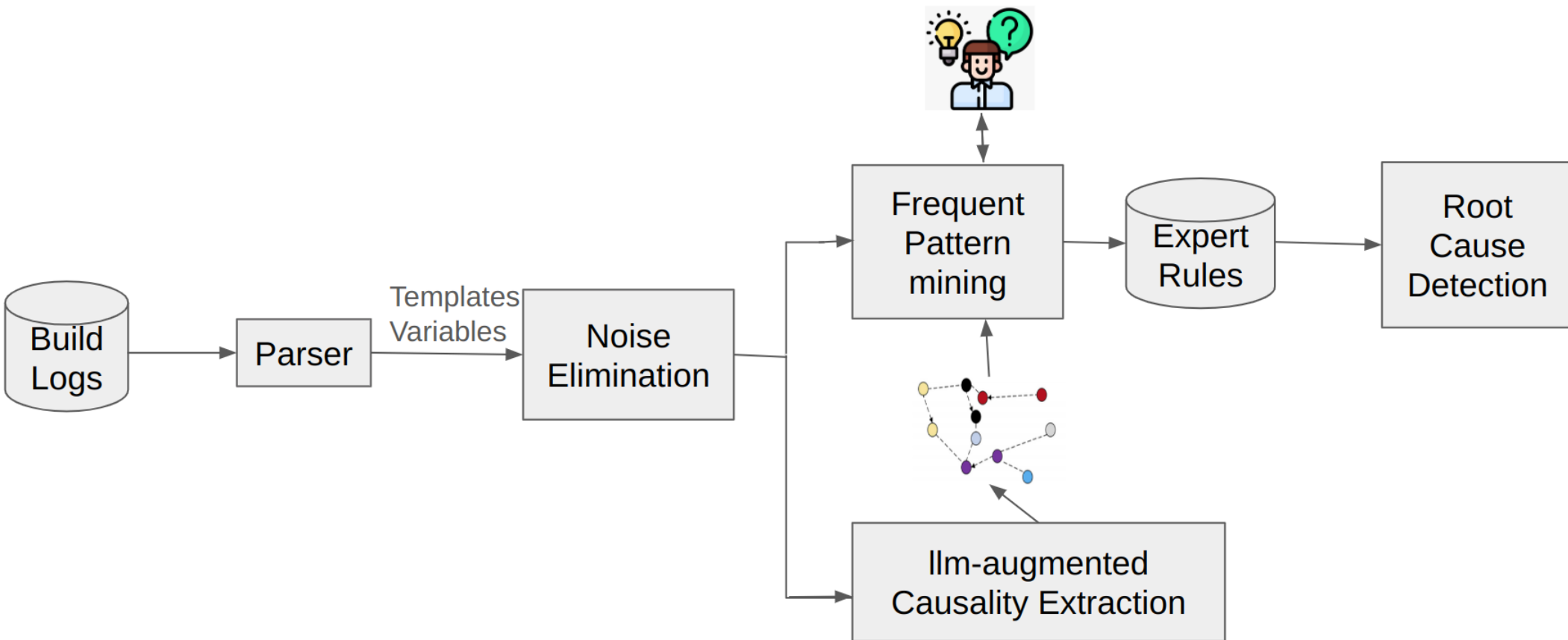
# Log Grouping and Causality Analysis

Fateme Faraji Daneshgar

Polytechnique Montréal

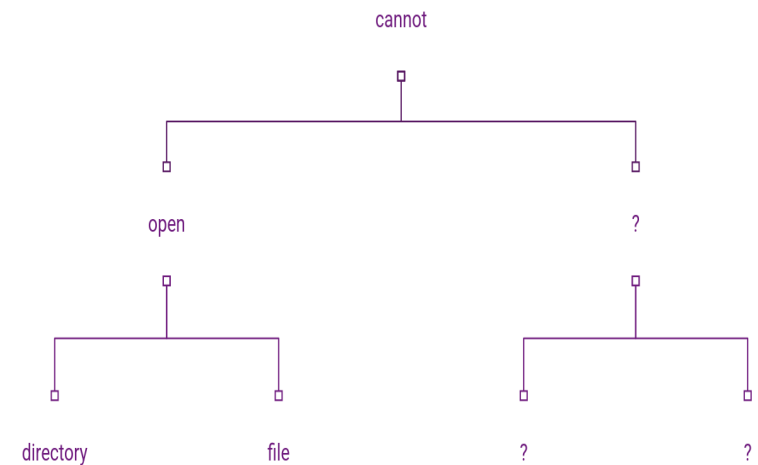
DORSAL Laboratory

# Root Cause Detection pipeline



# Parser

- Regular expression
  - Cannot open directory X
- Drain approach
  - Clustering logs using a prefix tree
  - Incremental prefix tree



Cannot open directory X  
Cannot open directory Y  
Cannot open directory Z

Cannot open directory x

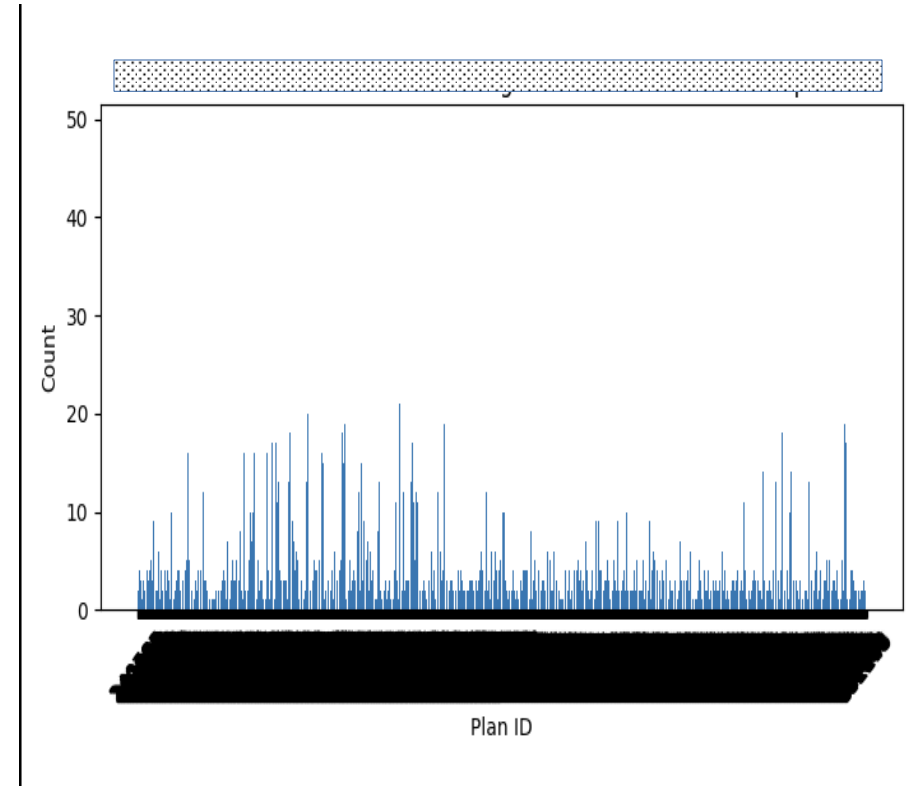
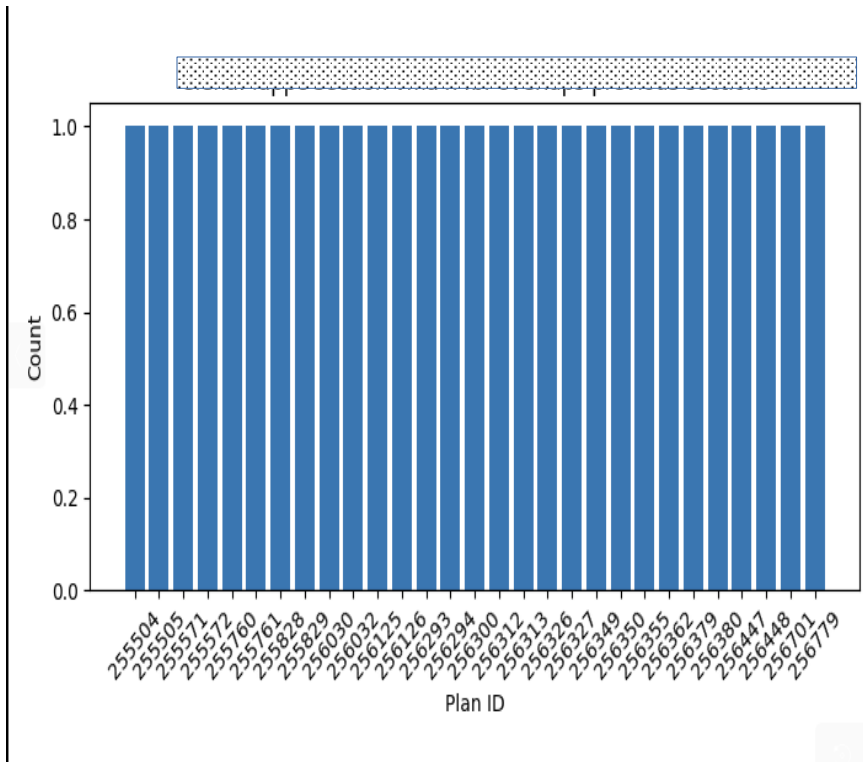
## Parser Results

---

- Same Templates
- Different Templates
  - Improvement (61.5%)
  - Errors(2.8%)
  - Low frequency and invariant variables (34.8%)

# Noise Elimination

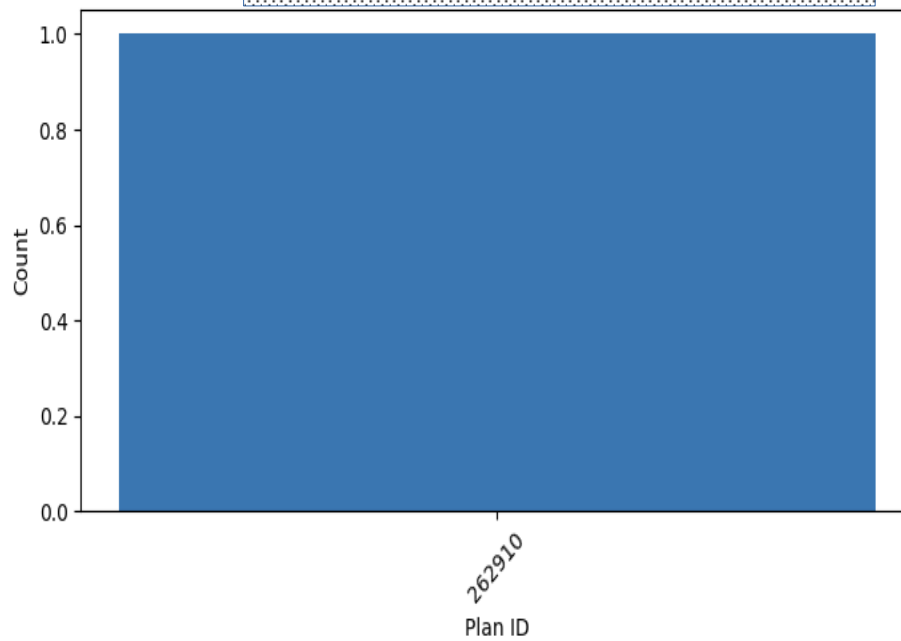
- Criterion
  - Fluctuation of the frequency of log templates



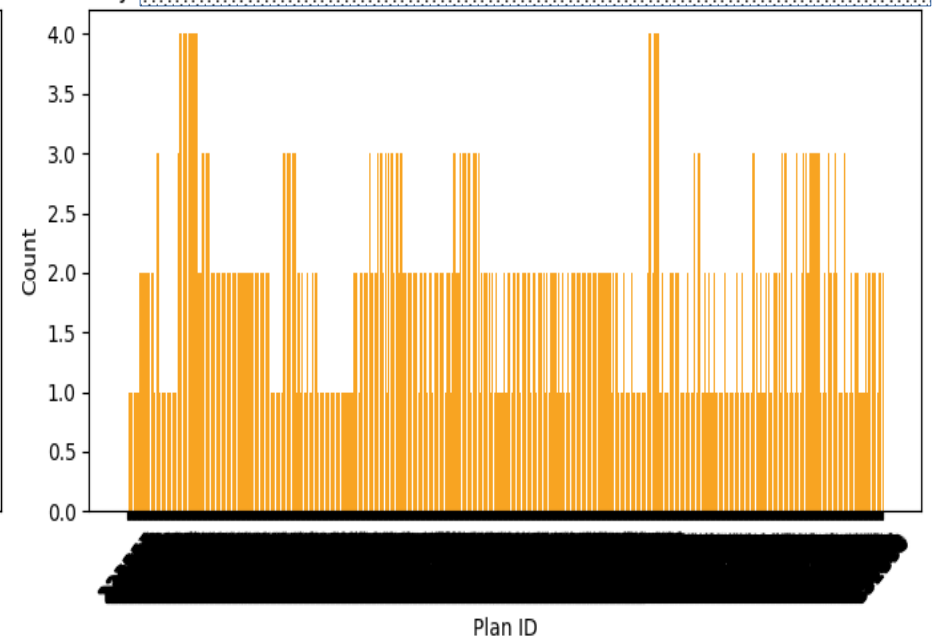
# Noise Elimination

Index 596 - Product: 6500

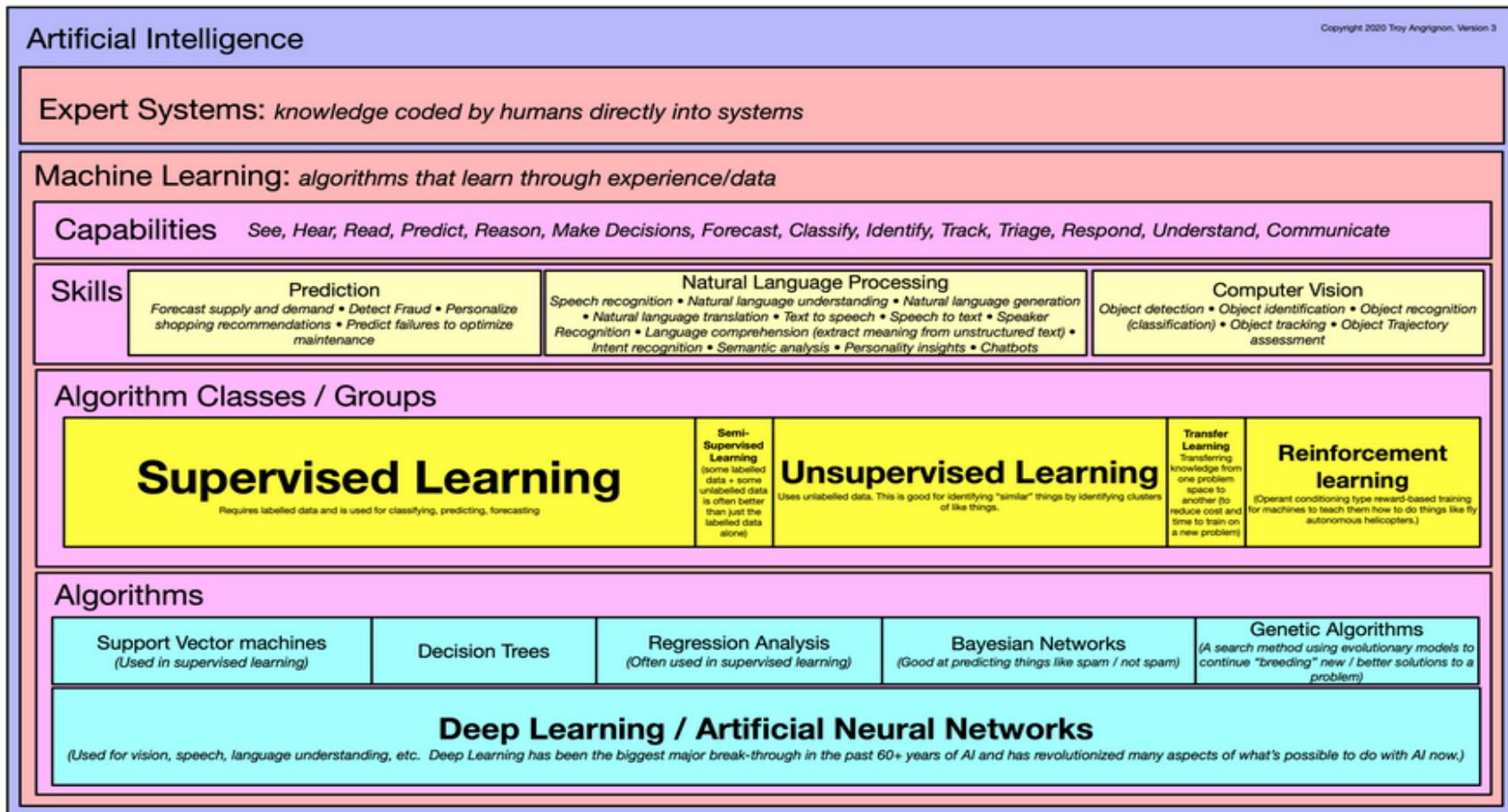
Ciena:



Poly:



# Expert Knowledge Extraction



# Frequent Pattern mining

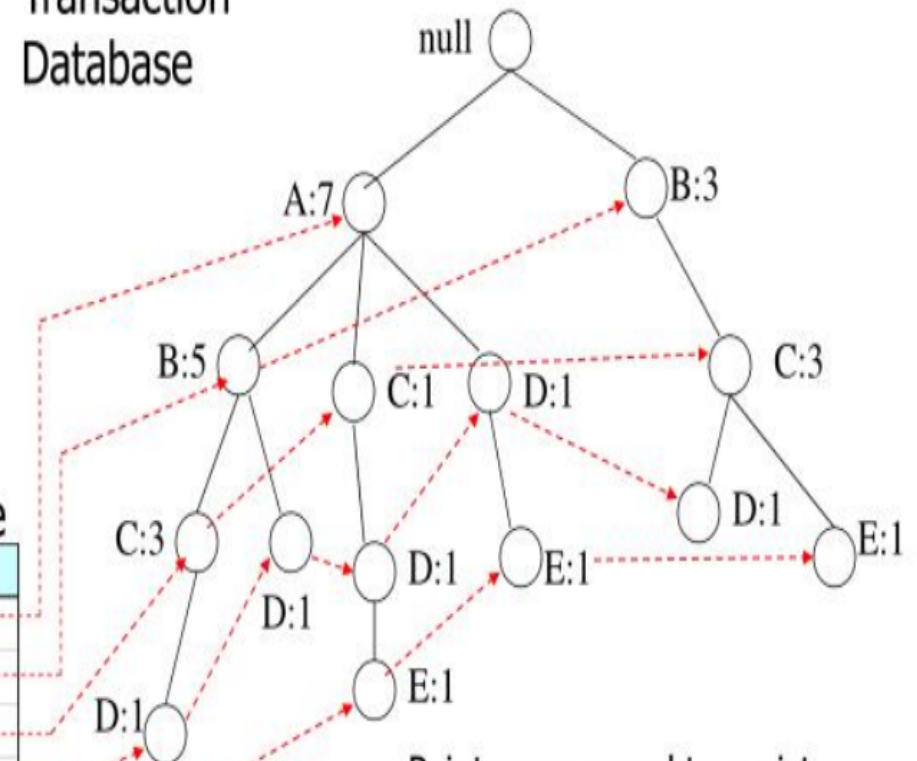
- Time consuming
- Candidate generation
  - Frequent Pattern
  - Conditional Pattern Growth

| TID | Items     |
|-----|-----------|
| 1   | {A,B}     |
| 2   | {B,C,D}   |
| 3   | {A,C,D,E} |
| 4   | {A,D,E}   |
| 5   | {A,B,C}   |
| 6   | {A,B,C,D} |
| 7   | {B,C}     |
| 8   | {A,B,C}   |
| 9   | {A,B,D}   |
| 10  | {B,C,E}   |

Header table

| Item | Pointer |
|------|---------|
| A    |         |
| B    |         |
| C    |         |
| D    |         |
| E    |         |

Transaction Database

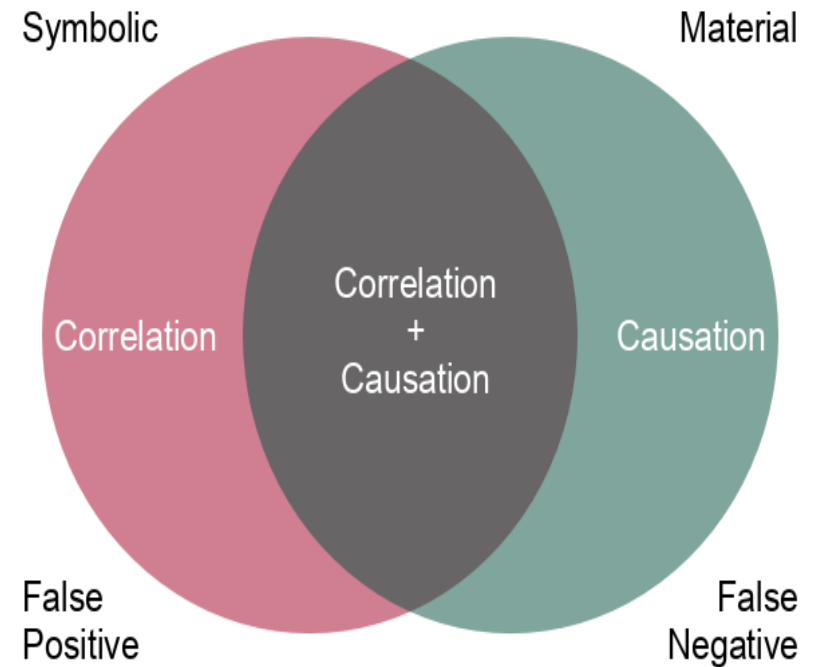


Pointers are used to assist frequent itemset generation



# Causality Analysis

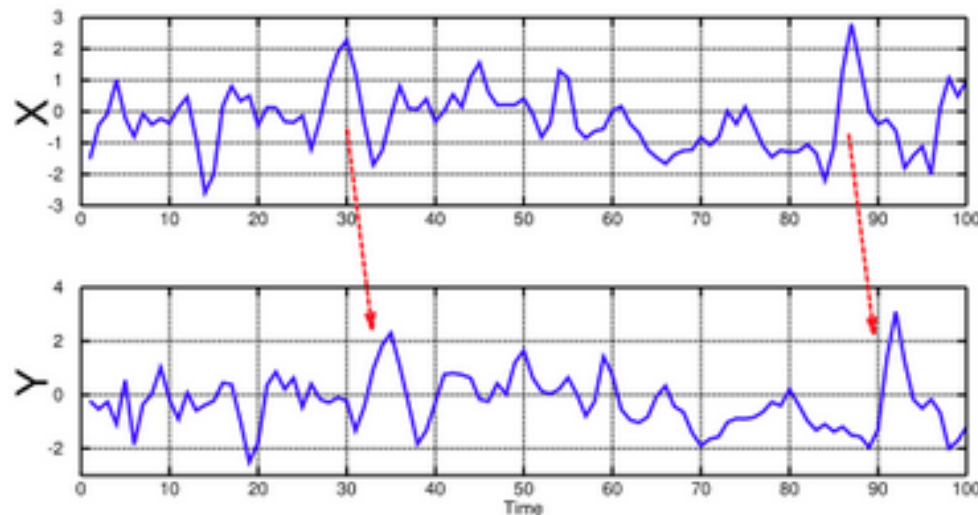
- Conditional Independence
  - Correlation does not imply Causality
    - PC, FCI, RFCI, ...
- $A \perp B|C \iff P(A|B,C)=P(A|B)$
- Variable A can predict Variable B



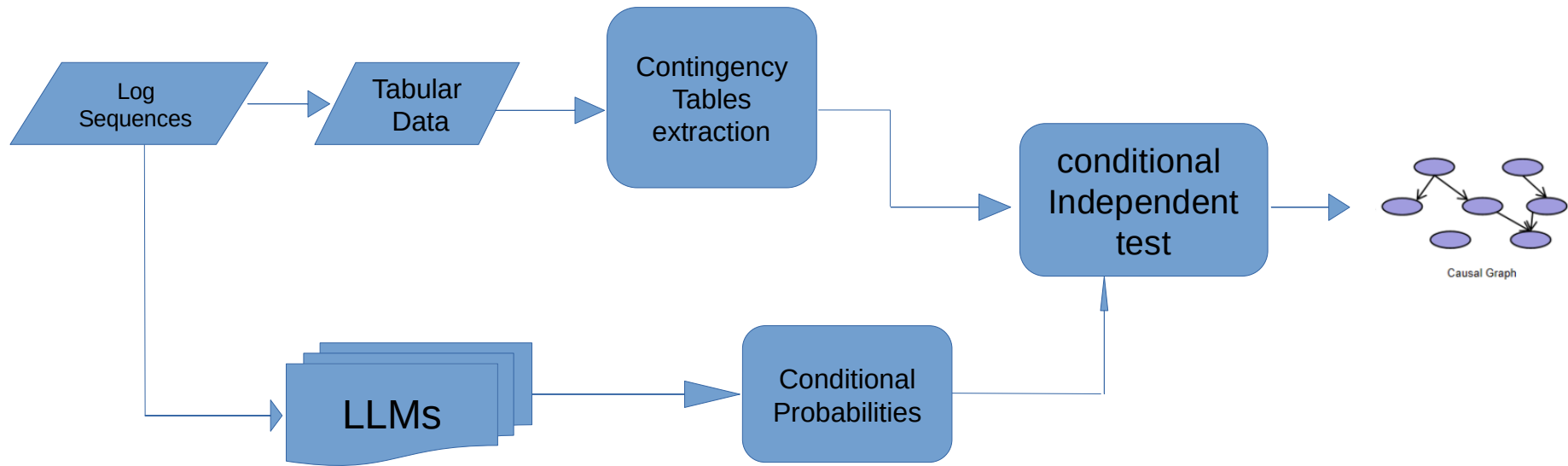
# Constraint of traditional CE for Log Analysis

- Tabular data
  - Missing the order
- Number of variables
- Time-series based
  - Sparse data
- Large language models

| $X$ | $Y$ | $C_0$ | $C_1$ |
|-----|-----|-------|-------|
| 0   | 0   | 0     | 0*    |
| 0   | 0   | 0     | 0*    |
| 0   | 0   | 0     | 0*    |
| 0   | 0   | 0     | 0*    |
| 1   | 1   | 1*    | 1     |
| 1   | 1   | 1*    | 1     |
| 1   | 1   | 1*    | 1     |
| 1   | 1   | 1*    | 1     |



# LLM-Augmented PC



# Contingency Table

- Log Sequences :
  - {AB,CA,BC,ACB,A}
- Conditional Independence test
  - A and B given C

$$E_{i,j} = \frac{\text{Row Total} \times \text{Column Total}}{\text{Grand Total}}$$

- Chi-square test

| A | B | C |
|---|---|---|
| 1 | 1 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 0 | 0 | 0 |
| 1 | 1 | 1 |
| 1 | 0 | 0 |

Contingency Table for C = 0:

| A \ B | 0 | 1 | Total |
|-------|---|---|-------|
| 0     | 1 | 1 | 2     |
| 1     | 1 | 1 | 2     |
| Total | 2 | 2 | 4     |

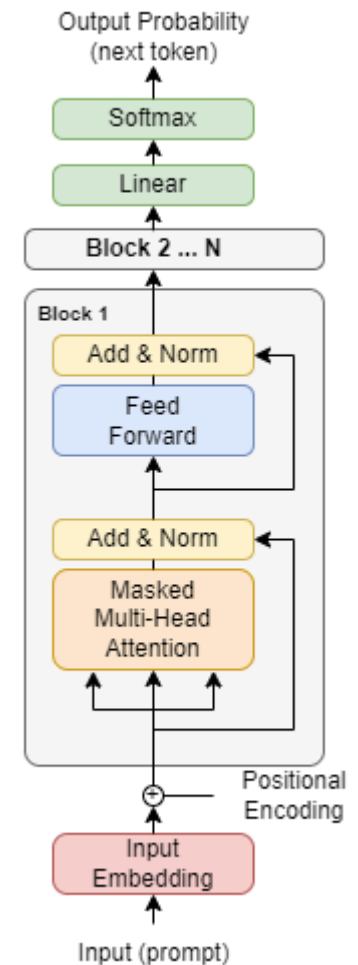
Contingency Table for C = 1:

| A \ B | 0 | 1 | Total |
|-------|---|---|-------|
| 0     | 1 | 0 | 1     |
| 1     | 1 | 1 | 2     |
| Total | 2 | 1 | 3     |

# LLM-Augmented PC

- Decoder-only transformer
- Each log template is considered as a token
- Continual pre-train
- $P(X|Y)$  and  $P(X|Y,C)$ :
- For each sequence having Y (Y and C):
  - probability of X for all positions after Y

Suppose  $s_i=Y$     $S_1 S_2 S_3 \dots S_i S_{i+1} \dots S_n$



## LLM-Augmented PC

---

- Traditional PC
  - log template seq:  $S_1 S_2 S_3 \dots S_i S_{i+1} \dots S_n$
  - 1 if we have X (**before or after**) and 0 otherwise
- LLM\_Augmented approach:
  - log template seq:  $S_1 S_2 S_3 \dots S_i S_{i+1} \dots S_n$
  - How probable is having or not having X **after** Y

Thanks!

[fateme-faraji.daneshgar@polymtl.ca](mailto:fateme-faraji.daneshgar@polymtl.ca)