

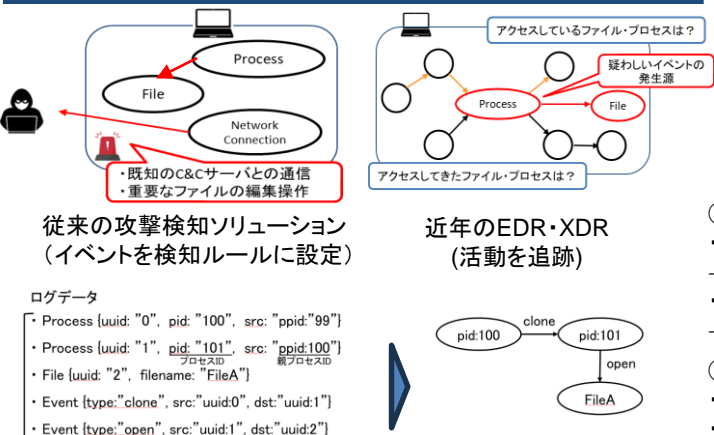
# Data Provenance を用いた悪性活動特定に向けた 良性活動抽出手法の提案

A Benign Activity Extraction Method for Malignant Activity Identification using Data Provenance

齋藤太新・マネジメント分科会・情報セキュリティ大学院大学

In order to gain a complete picture of cyber attacks and to identify the source of cyber attacks, a method is being developed to identify malicious activities by tracing data probabilities and connecting the dependencies of a series of related events. However, the problem is the dependency explosion, in which a large number of non-malicious, normal computer system actions are included in the dependencies, creating a huge graph that makes it difficult to identify malicious activities. To cope with the dependency explosion, we propose a method to reduce the search space for malicious activities by extracting and removing frequently occurring benign activities through natural language processing of log data and analysis of activities in the computer system using similarity judgments. In an evaluation experiment, we used a large public dataset to evaluate the effectiveness of the proposed method on the dependency explosion problem, and found that benign activities can be used to reduce dependencies, and that 6.8% to 39% of activities in a computer system can be defined as benign activity patterns. In addition, we showed that removing benign activities extracted from a portion of log data can significantly reduce the search space in large data sets.

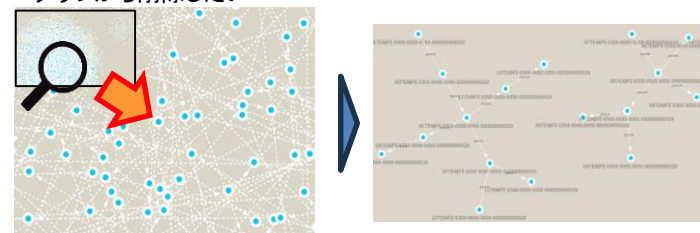
## 1. 背景: グラフ用いた悪性活動の分析



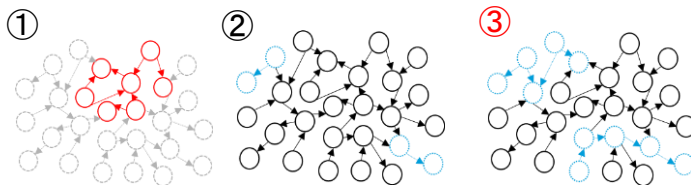
ログデータからコンピュータシステム内で発生した活動をグラフとして可視化  
→ 攻撃の分析(全体像の把握, 影響範囲・発生源特定)に活用

## 2. 課題: 依存関係の爆発

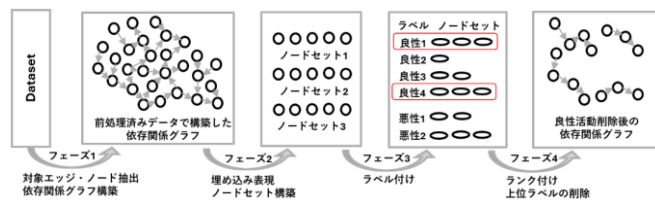
大量のノードを含む巨大なグラフは分析に役立たない  
→ 攻撃に関係のない(正規のユーザ・システムの動作)依存関係はグラフから削除したい



## 3. 既存の対処法と提案手法



- 攻撃のパターンは増え続ける(新しい攻撃の登場)  
→ 攻撃の変化に対処するための再学習が必要
- 悪性活動をピンポイントで抽出するのは困難  
→ まずは良性部分を削って、分析に使う探索空間を減らすべき
- グラフの削減率に限界がある
- 本当に除外して良いイベントなのか? (環境・分析者の違い)
- 分析対象環境から収集したデータから良性活動を抽出・グラフから削除
- 良性を見るため, ①よりも攻撃の変化に対応可
- 様々な良性活動パターンを抽出することで②の限界に対処
- 分析対象環境のログデータから良性活動を抽出するため, 環境の違いに対応可



- グラフからノード5個分の活動(ノードセット)を全て抽出
- コサイン類似度でノードセットにラベル付け
- 頻出の良性ラベルを良性活動と定義、グラフから削除

## 4. 実験結果と考察・今後の展望

- 3種類のデータセット(Darpa TC Data)を使用して実験
- 各データセットからそれぞれ一部を切り取って良性活動抽出(計9パターンの実験)
- 最大52.3%のノードを削減(半分以上削減)
- 各データセットで約11%~39%はコンピュータシステム内頻出の良性活動として定義可能
- データセットの3%程度のデータから抽出した良性活動で大規模データから依存関係を削減可能
- 汎用的な用途(E3 Theia)より, 同じ動作を繰り返すコンピュータシステムのログデータ(E5 Theia, E5 Marple)で高い削減率を発揮

データセット	データサイズ	評価用データに対する割合
E3 Theia-A	3.8GB	13.4%
E3 Theia-B	3.8GB	13.4%
E3 Theia-C	3.8GB	13.4%
E5 Theia-A	4.0GB	1.35%
E5 Theia-B	4.0GB	1.35%
E5 Theia-C	4.0GB	1.35%
E5 Marple-A	3.6GB	2.98%
E5 Marple-B	3.6GB	2.98%
E5 Marple-C	3.8GB	3.15%

データセット	平均削減率 (%)		平均実行時間 (sec)	
	最小 (n=3)	最大 (n=1500)	最小 (n=3)	最大 (n=1500)
E3 Theia	6.82	10.5	7,163	94.97
E5 Theia	27.3	27.9	792.4	3,061
E5 Marple	2.62	39.1	541.7	8,563

削除対象ラベル数 (n種)	ノード数	FN	FP	ノード削減率 (%)	実行時間 (sec)
3	11,803,667	0	11,803,657	3.03	558.7
10	11,545,745	0	11,545,735	5.15	717.1
100	8,956,381	0	8,956,371	26.4	2,091
500	6,600,627	0	6,600,617	45.8	5,768
1000	6,048,300	0	6,048,290	50.3	9,130
1500	5,800,663	0	5,800,653	52.3	11,690

今後  
パラメータの最適化, 多種類データに対応, 侵入検知システムへの応用