# Reconstructing Alert Trees for Cyber Triage

**Eric Ficke**

**PhD Dissertation Defense**
**The University of Texas at San Antonio**
**April 13, 2022**

**Committee Co-Chairs**
**Prof. Ravi Sandhu**
**Prof. Shouhuai Xu**

**Committee Members**
**Prof. Xiaoyin Wang**
**Prof. Greg White**
**Prof. Mimi Xie**

**UTSA®**

# Publications

- Published

    1. S. He, **E. Ficke**, M. Pritom, H. Chen, Q. Tang, Q. Chen, M. Pendleton, L. Njilla, and S. Xu. *Blockchain-based Automated and Robust Cyber Security Management*, Journal of Parallel Distributed Computing, 163: 62-82 (2022)

    2. **E. Ficke** and S. Xu. APIN: *Automatic Attack Path Identification in Computer Networks*, IEEE International Conference on Intelligence and Security Informatics (ISI), 2020. **[Dissertation Chapter 2]**

    3. **E. Ficke**, K. Schweitzer, R. Bateman, and S. Xu. *Analyzing Root Causes of Intrusion Detection False-Negatives: Methodology and Case Study*. IEEE Military Communications Conference (MILCOM), 2019.

    4. J. Mireles, **E. Ficke**, J. Cho, P. Hurley, and S. Xu. *Metrics Towards Measuring Cyber Agility*. IEEE Transactions on Information Forensics and Security (IEEE T-IFS), 14(12): 3217-3232 (2019).

    5. **E. Ficke**, K. Schweitzer, R. Bateman, and S. Xu. *Characterizing the Effectiveness of Network-Based Intrusion Detection Systems*. IEEE Military Communications Conference (MILCOM), 2018.

- Manuscripts to be submitted for review

    1. **E. Ficke**, R. Bateman, and S. Xu. AutoCRAT: Automatic Cumulative Reconstruction of Alert Trees. **[Dissertation Chapter 3]**

    2. R. Garcia-LeBron, **E. Ficke**, W. Wu, S. Xu. Characterizing Cyber Attack Reconnaissance Trajectories.

    3. **E. Ficke**, R. Bateman, and S. Xu. Alert Tree Reduction and Visualization. **[Dissertation Chapter 4]**

# Dissertation Outline

- Introduction

- APIN: Alert Path Identification in Computer Networks

- AutoCRAT: Automatic Cumulative Reconstruction of Alert Trees

- Alert Tree Reduction and Visualization

- Conclusion

**UTSA**®

# Introduction

Chapter 1

# Background – Alert Trees

- Cyber Triage (Network-level)

  - Alert prioritization

  - Alert correlation

  - Attack lifecycle

- Attack Prediction

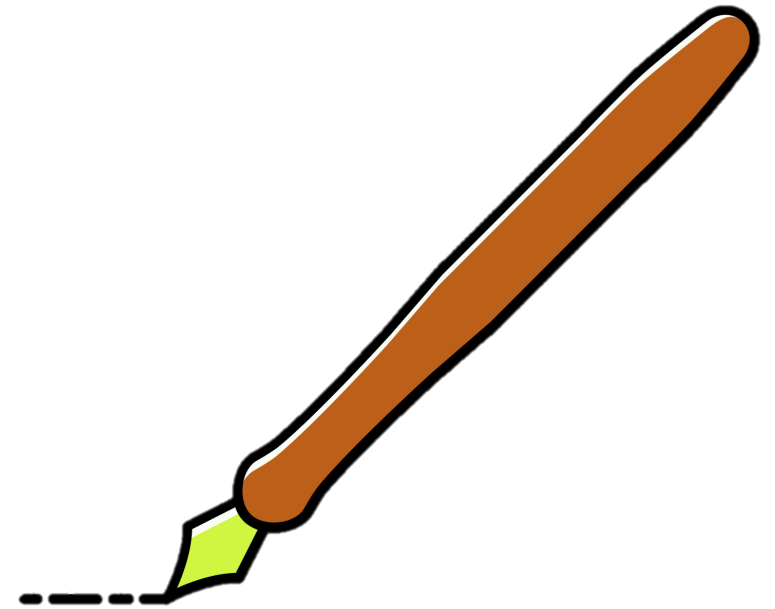  - Attack graphs / trees / paths

  - Vulnerability graphs

# Motivation

- Alert volume

  - Unrealistically low in ad hoc datasets

  - **Overwhelms human analysis** in real data

- Alert graph / tree / path formalization

  - Varies by usage

  - Depends on **spatial and temporal dependencies**

# Chapter Themes

2. Alert Path Identification (**APIN**)
   - Alert path reconstruction
   - Threat score (TS) ranking

3. Cumulative Reconstruction (**AutoCRAT**)
   - Alert tree reconstruction
   - Alternative path reconstruction method
   - Asymptotic and real analysis

4. Reduction and Visualization
   - Mitigates emergent problem of tree size

UTSA.

# Dissertation Outline

- Introduction

- APIN: Alert Path Identification in Computer Networks

- AutoCRAT: Automatic Cumulative Reconstruction of Alert Trees

- Alert Tree Reduction and Visualization

- Conclusion

# APIN: <u>A</u>lert <u>P</u>ath <u>I</u>dentification in Computer <u>N</u>etworks

Chapter 2

**UTSA**®

# Motivation: Cyber Triage

- Time sensitive

- Resource intensive
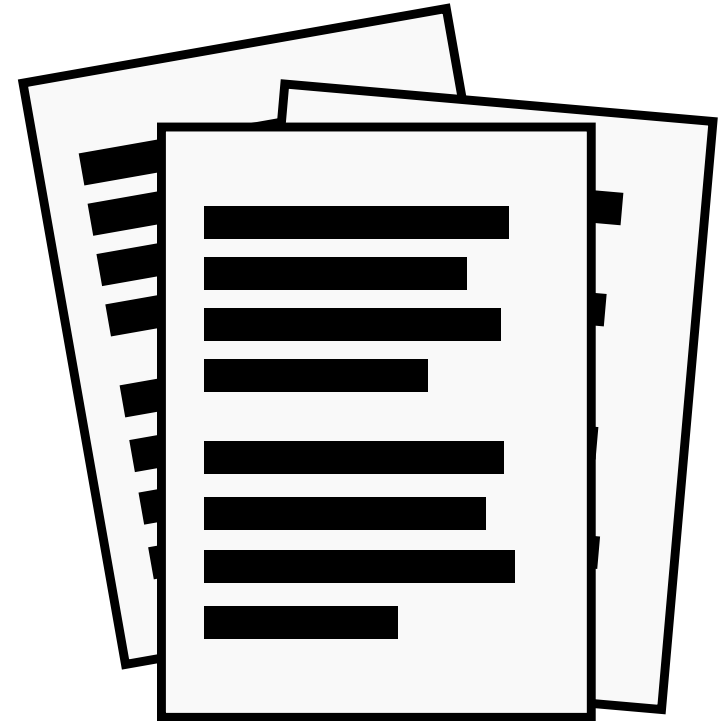
- Error prone

- Large search space

# Contributions

- Attack Tracking

  - Alert paths show footprints between victim computers

  - Spatiotemporal path reconstruction method

- Heuristics

  - Threat score shows attacker effort

  - (Actual compromise may vary)

# APIN Framework
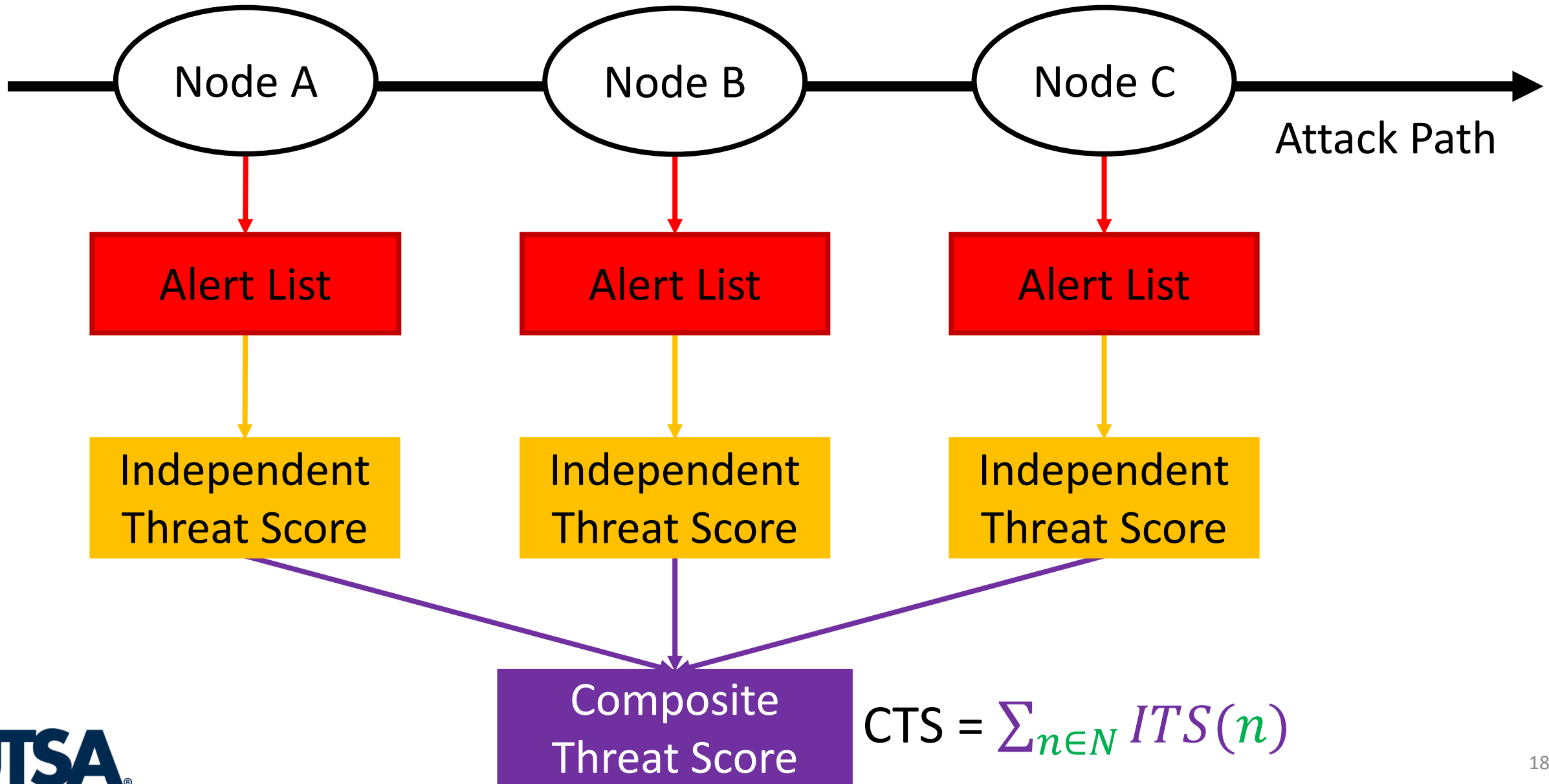
# Metric: Independent Threat Score

- Input:
  - $A_{in}$ alert types (inbound)
  - $A_{out}$ alert types (outbound)
- Terms:
  - $D_{in} = |A_{in}|$
  - $D_{out} = |A_{out}|$
  - $S_{in} = \sqrt[|A_{in}|]{\prod_{a \in A_{in}} |a|}$

$$ITS = \sqrt[3]{D_{in} \cdot D_{out} \cdot S_{in}}$$

"$D$" represents alert diversity
"$S$" represents alert scale (by type)

UTSA.

14

# Methods (Alert Path Identification)

Approach: breadth-first search in reverse-chronological order



Not a valid path:
Violates temporal attack dependency

8

4

16

20

5

10

Edges indexed by timestamp

# Preliminary Analysis

- Scans (high volume, low threat)
  - Deprioritize inbound alerts

- (lower threat)
  - Cause
  - Blacklist nodes
    - Restricts path identification
    - Leaves nodes unmonitored

High-granularity network segmentation improves performance significantly

# Metric: Weighted Independent Threat Score

- Input:
  - $A_{in}$ alert types (inbound)
  - $A_{out}$ alert types (outbound)

- Terms:
  - $D_{in} = |A_{in}|$
  - $D_{out} = |A_{out}|$
  - $S_{in} = \sqrt[|A_{in}|]{\prod_{a \in A_{in}} |a|}$
  - $W = w_1 + w_2 + w_3$

$$ITS = \sqrt[W]{D_{in}^{w_1} \cdot D_{out}^{w_2} \cdot S_{in}^{w_3}}$$

"$D$" represents alert diversity
"$S$" represents alert scale (by type)

# Metric: Composite Threat Score



CTS $= \sum_{n \in N} ITS(n)$

# Preliminary Results: DARPA '99

- Notable paths, using queries from top 5 nodes

| Origin | Composite Threat Score ▽ | Length (#edges) | Notable Alerts |
|---|---|---|---|
| 209.67.29.11 | | 1 | Windows 95 Malware |
| 172.16.116.194 | 3.43 | | |
| 207.25.71.141 | 3.41 | 1 | Windows 95 Malware |
| 192.168.1.30 172.16.112.5 | 3.41 | 1 | Public SNMP Access |
| 206.132.25.51 | 3.36 | 1 | Windows 95 Malware |

**Artificial datasets need multi-step attacks**

Hidden IPs are repeated from higher-ranked paths

# Results: CSE-CIC-IDS2018

| [Path Origin, Path Target] | Composite Threat Score ▽ | Length (#Edges) | Notable Alerts |
|---|---|---|---|
| 103.47.124.154 54.172.47.69 | 34.31 | 4 | EternalBlue (WannaCry) NAT Traversal |
| 172.31.67.54 52.87.201.4 | 33.60 | 3 | EternalBlue (WannaCry) NAT Traversal |
| 71.6.165.200 172.31.64.78 | 26.42 | 3 | Blacklisted IP group SQL Scan |
| 77.222.106.20 172.31.66.112 | 21.30 | 3 | EternalBlue (WannaCry) SMB Share Access |
| 172.31.64.78 172.31.0.2 | 21.11 | 1 | Suspicious DNS Query |

**UTSA**

# Results: CSE-CIC-IDS2018

**1.** 103.47.124.154 → 172.31.67.46 → 103.68.10.188 → 172.31.66.112 → 54.172.47.69 — EternalBlue NAT Traversal

**2.** 172.31.67.54 → 212.174.232.94 → 172.31.64.46 → 52.87.201.4 — EternalBlue NAT Traversal

**3.** 71.6.165.200 → 172.31.64.71 → 149.255.35.24 → 172.31.64.78 — Blacklisted IP group SQL Scan

**4.** 77.222.106.20 → 172.31.67.46 → 103.68.10.188 → 172.31.66.112 — EternalBlue SMB Share Access

**5.** 172.31.64.78 → 172.31.0.2 — Suspicious DNS Query

UTSA®

# Results: CSE-CIC-IDS2018



**1.** 103.47.124.154 → 172.31.67.46 → 103.68.10.188 → 172.31.66.112 → 54.172.47.69 — EternalBlue NAT Traversal

**4.** 77.222.106.20 — EternalBlue SMB Share Access

**2.** 172.31.67.54 → 212.174.232.94 → 172.31.64.46 → 52.87.201.4 — EternalBlue NAT Traversal

**3.** 71.6.165.200 → 172.31.64.71 → 149.255.35.24 → 172.31.64.78 — Blacklisted IP group SQL Scan

**5.** 172.31.0.2 — Suspicious DNS Query

Origin → Target

# Results: CSE-CIC-IDS2018

# Dissertation Outline

- Introduction

- APIN: Alert Path Identification in Computer Networks

- AutoCRAT: Automatic Cumulative Reconstruction of Alert Trees

- Alert Tree Reduction and Visualization

- Conclusion

# AutoCRAT: <u>Auto</u>matic <u>C</u>umulative <u>R</u>econstruction of <u>A</u>lert <u>T</u>rees

Chapter 3

UTSA®

# Motivation: Alert Tree Optimization

- Improve reconstruction

- Identify optimization tradeoffs

- Formalize alert trees

# AutoCRAT Architecture

# Methods (Path Maintenance)

Approach: maintain every path at all times, merging as they join



**Paths grow sequentially**
**Paths remain independent until linked**
**Trees form spontaneously**

# Methods (Tree Reconstruction)

Approach: maintain every path at all times, merging as they join



**Paths grow sequentially**
**Paths remain independent until linked**
**Trees form spontaneously**

# Asymptotic Comparison

|  | APIN | AutoCRAT |
|---|---|---|
| I... | $O(1)$ | $O(|A|^2)$ |
| ... | $O(|V| + |A|)$ | $O(|A| + |E|^3) \subseteq O(|A|^3)$ |
| ... | | $O(1)$ |
| Retrieve Trees | | |
| Reinsert | | |
| Database Size | $O(|V| + |A|)$ | |

$A$ – alerts
$V$ – vertices (computers)
...ndpoints

**APIN dominates insertion and storage; AutoCRAT conditionally dominates retrieval**

*APIN ranks nodes, while AutoCRAT ranks endpoints and paths.

UTSA.

# Results Comparison

| | APIN | AutoCRAT | APIN-Internal | AutoCRAT-Internal |
|---|---|---|---|---|
| | | 13h42m41s | 9s | 35s |
| | | | 0.28s | 5s |
| Retrieve Top 20 Trees | | | | |
| Coverage (Nodes) | 99.6% | 100% | | |
| Coverage (Events | 3.4% | 100% | 0.6% | |
| Database Size | 637 MB | 1.1 GB | 2.9 MB | 2.4 MB |

*APIN ranks nodes, while AutoCRAT ranks endpoints and paths.
†These ranks are inferred from their ends (for paths) or root (for trees)

*AutoCRAT dominates path retrieval and under comparable footprints, at*

**The vast majority of alerts span network borders**

UTSA®

# Dissertation Outline

- Introduction

- APIN: Alert Path Identification in Computer Networks

- AutoCRAT: Automatic Cumulative Reconstruction of Alert Trees

- Alert Tree Reduction and Visualization

- Conclusion

# Alert Tree
# Reduction and Visualization

Chapter 4

# Motivation

**Facilitate cyber triage by selectively pruning alert trees**

- Reduce visual strain

  - "Which nodes can be removed to facilitate tree interpretation?"

- Preserve salient information

  - "What nodes must be kept based on relevant metrics?"

# Motivating Example

- This tree (from real data) has **3090 nodes**.

- Graphviz is forced to render it at **6% of its original resolution**.*



**\*If you have difficulty reading any of the node labels, that's exactly the problem we need to solve**

# Motivating Example

- After reduction, 3090 nodes becomes 40 nodes **(98.7% reduction)**



So how do we do it???

# Alert Tree Reduction Architecture

# Terminology



**Siblings:**

Duplicate labels may exist in a tree but not in a path

UTSA.

41

# Terminology



**Branch:**

# Terminology (Graph Theory vs Data Structures)

# Terminology (Graph Theory vs Data Structures)
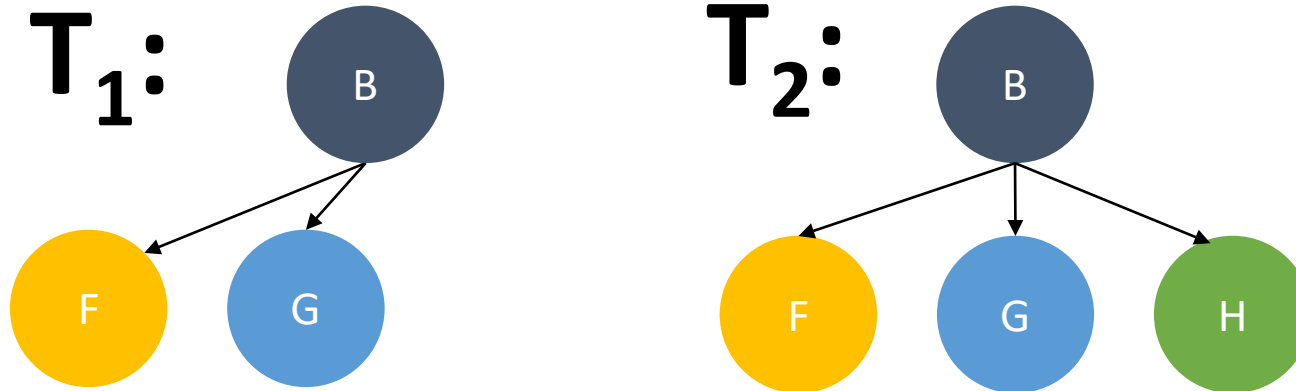
# Terminology (Graph Theory vs Data Structures)

**$T_1$:**



**$T_2$:**



In Graph Theory:

- $T_1$ is a subtree of $T_2$
- F is a subtree of $T_1$ (or $T_2$)

In Data Structures:

- F is a subtree of B (in either tree)

**We need a new term for the relationship $T_1:T_2$ that eliminates ambiguity**

# Hypotree and Hypertree



- Designate: $T_1$ is a hypotree of $T_2$ ($T_1 \lhd\!\!= T_2$)
- Designate: $T_2$ is a hypertree of $T_1$ ($T_2 \rhd\!\!= T_1$)
- Every tree is both a hypotree and a hypertree of itself
- We also designate proper hypotree ($\lhd$) and proper hypertree ($\rhd$)

# Hypotree and Hypertree



**Definition:** A tree $T_{hypo}$ is a hypotree of a tree $T_{hyper}$ if:

$\forall\ n \in T_{hypo}, \exists\ n' \in T_{hyper}:$

$\qquad \forall\ i \in \{0, 1, ..., |n.\text{ancestors}|\},\ n.\text{ancestors}_i = n'.\text{ancestors}_i$

*Hypertree is derived from hypotree. Refer to the paper for exact detail

# Merging Sibling Leaves [MSL(A)]



Node labels represent IP addresses

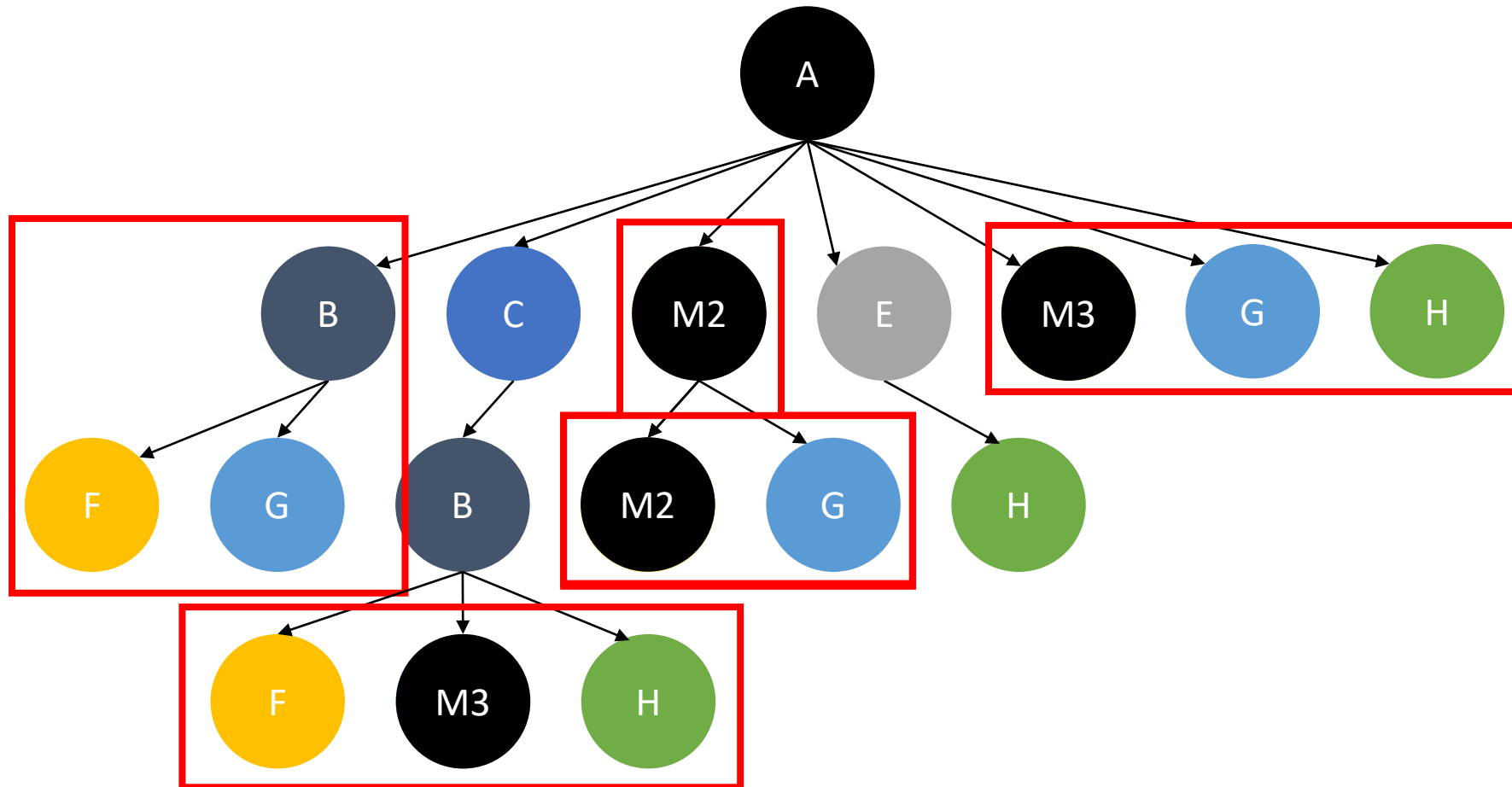Duplicate labels may exist in a tree but not in a path within that tree

Here, node colors show labels (rather than threat score) for ease of understanding
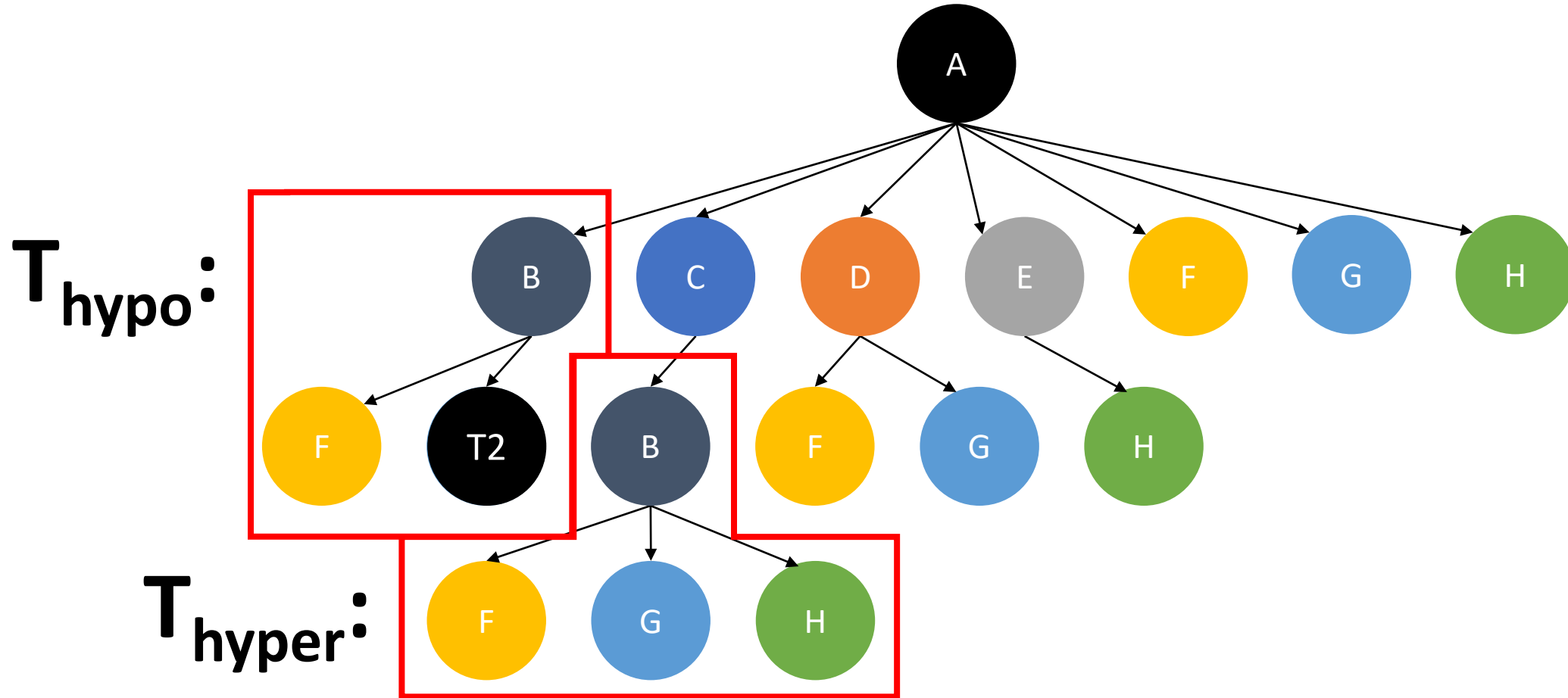
48

# Merging Similar Sibling Branches [MSB(A)]



Similar Branches: A set of branches for which all subtrees excluding the branch root exist in both branches
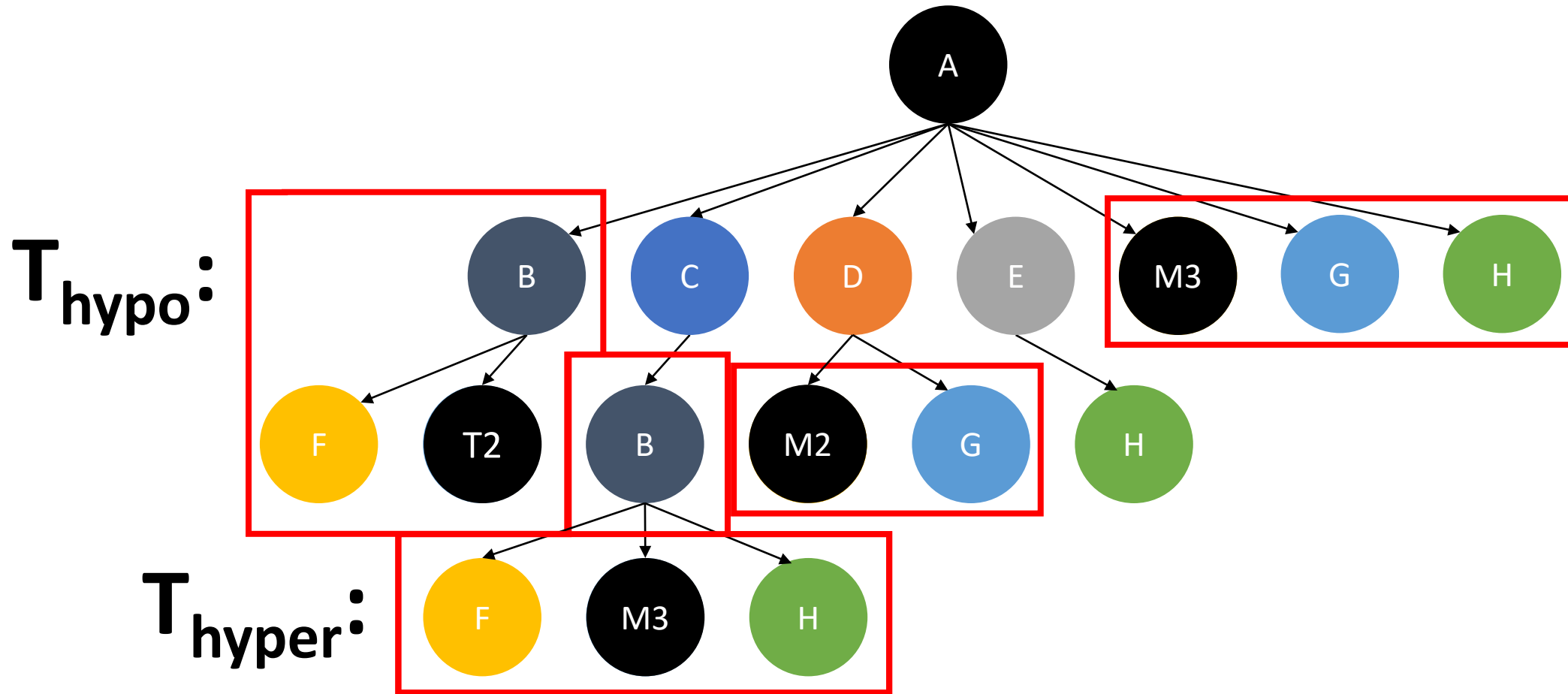
# Merging Sibling Branches & Leaves [MSL(MSB(A))]

# Truncating Hypotrees
# [TH(A)]

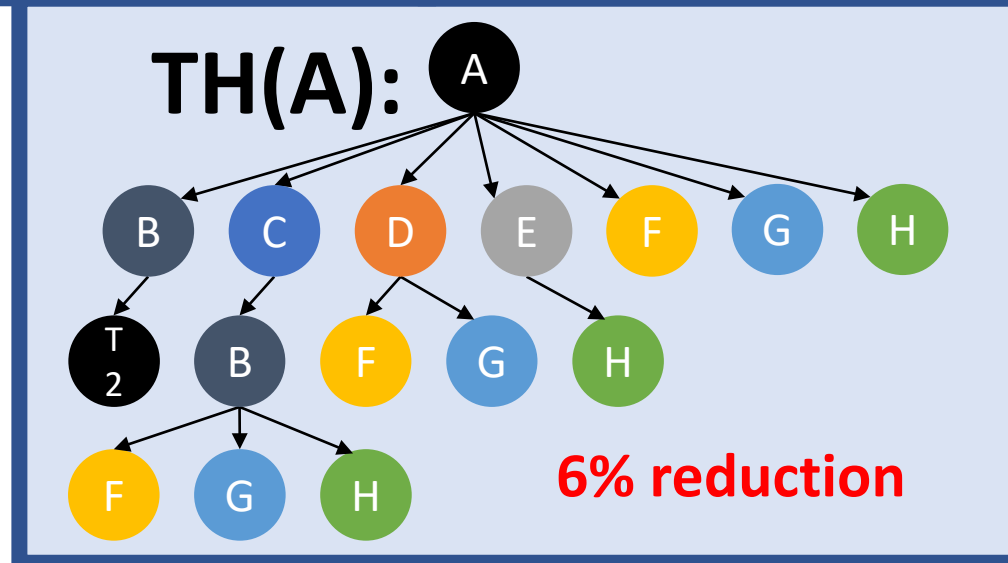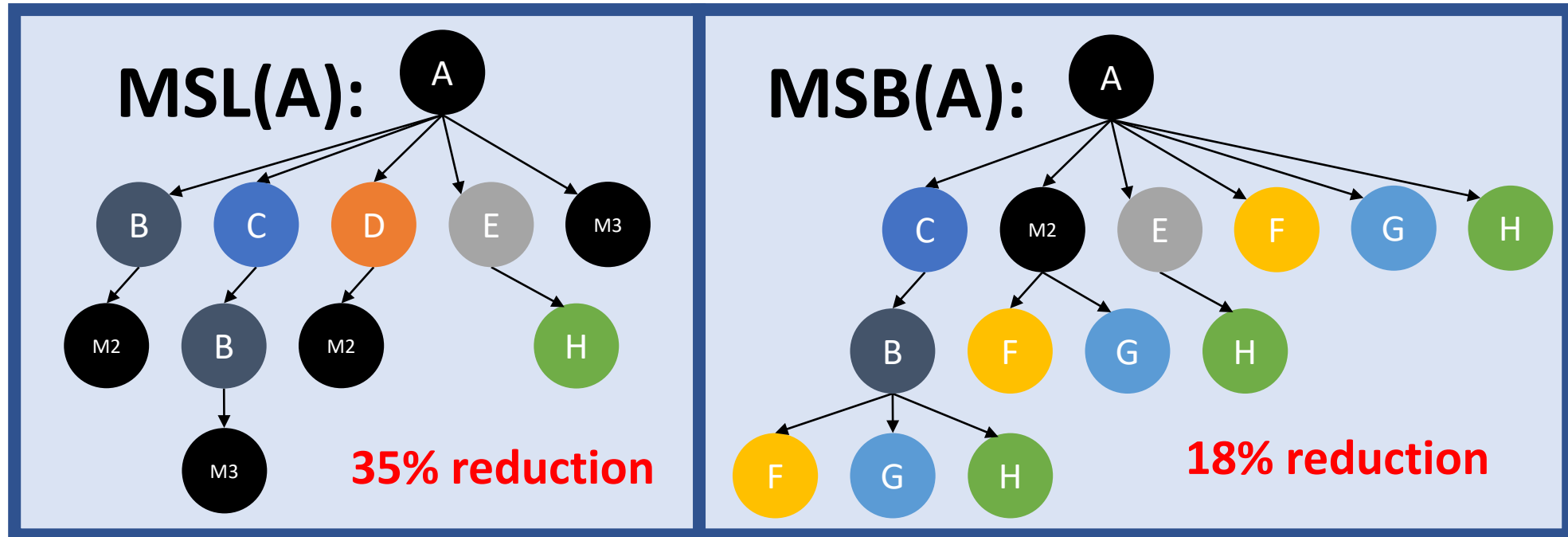# Truncating Hypotrees & Merging Sibling Leaves [MSL(TH(A))]

# Method Restrictions

- MSL makes some trees similar (because "M2" = "M2")
  - MSB(MSL(T)) is unsafe (but MSL(MSB(T)) is safe)
  - TH(MSL(T)) is unsafe (but MSL(TH(T)) is safe)

- MSB and TH may target the same branches
  - MSB(TH(T)) ≠ TH(MSB(T))

- The 5 valid reduction schedules:
  1. MSB(T)
  2. MSL(MSB(T))
  3. MSL(T)
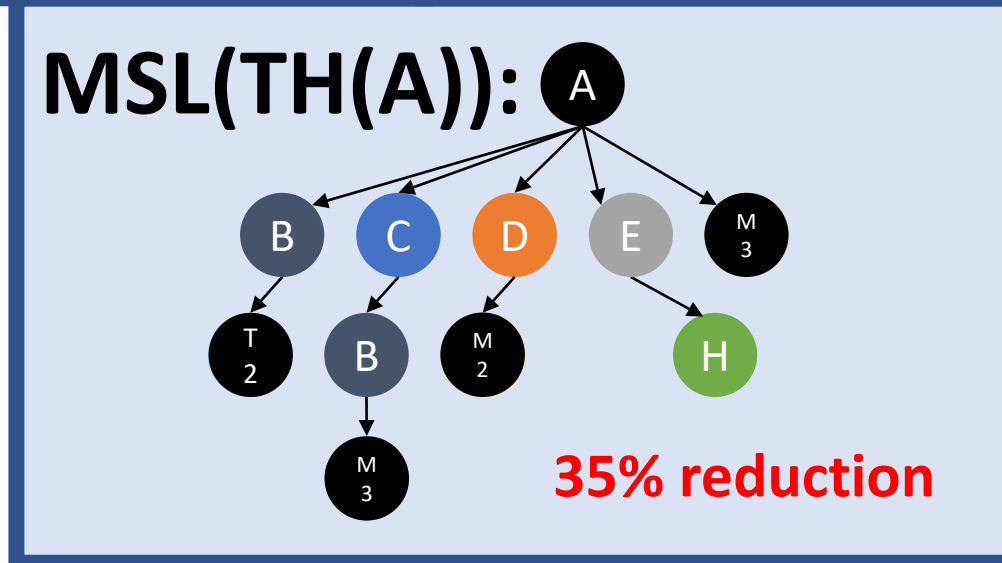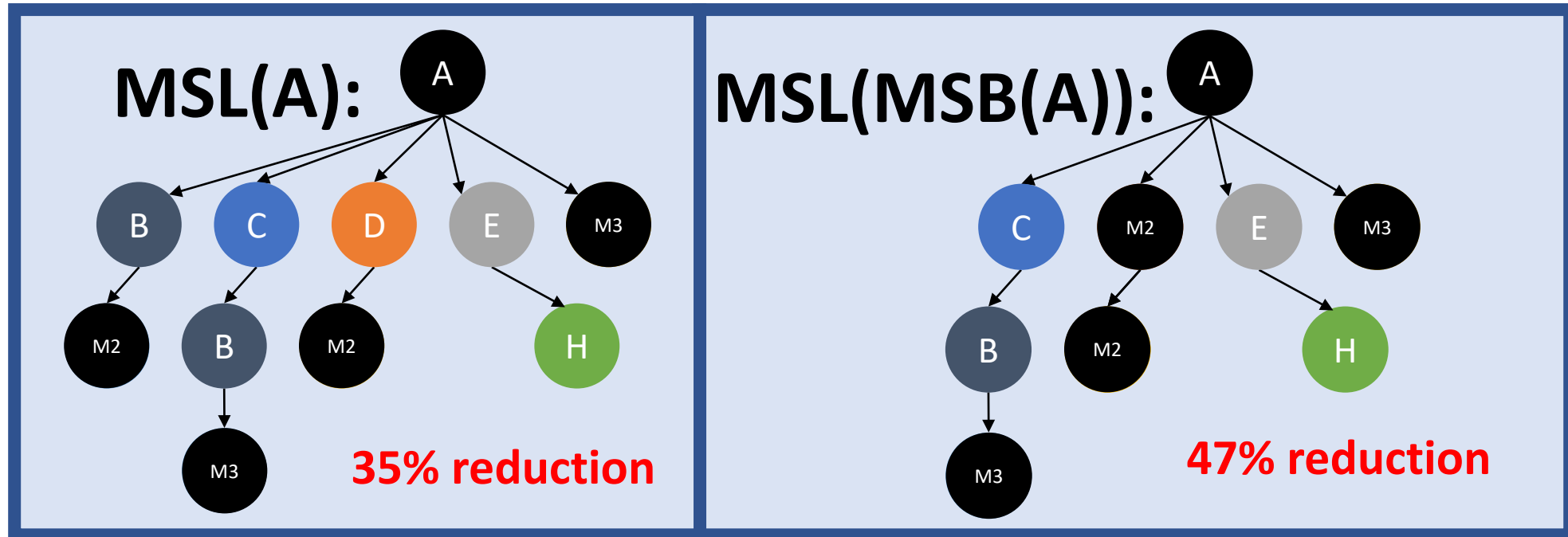  4. MSL(TH(T))
  5. TH(T)

MSL: Merge Sibling Leaves
MSB: Merge Sibling Branches
TH: Truncate Hypotrees

# Method Comparisons (Toy Example)



MSL: Merge Sibling Leaves
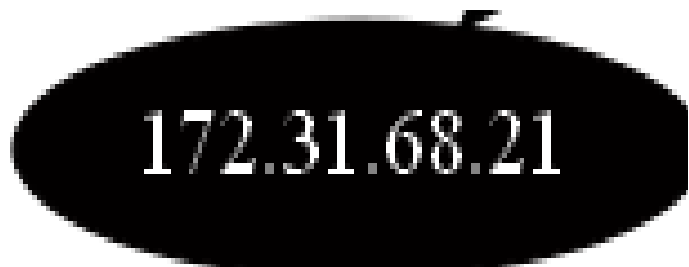MSB: Merge Sibling Branches
TH: Truncate Hypotrees

# Method Comparisons (Toy Example)



MSL: Merge Sibling Leaves
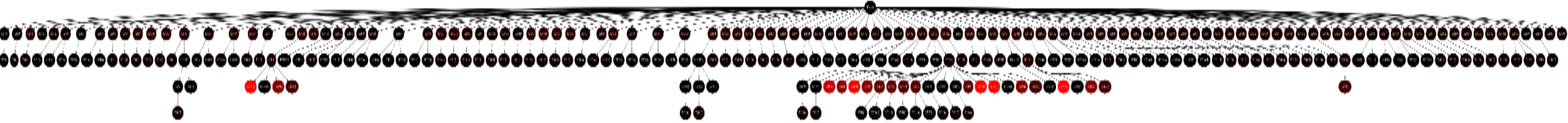MSB: Merge Sibling Branches
TH: Truncate Hypotrees

# Visualization

- Black (low threat) -> red (high threat)
  - Min-max normalized

- Merged nodes
  - Color shows highest threat of those merged
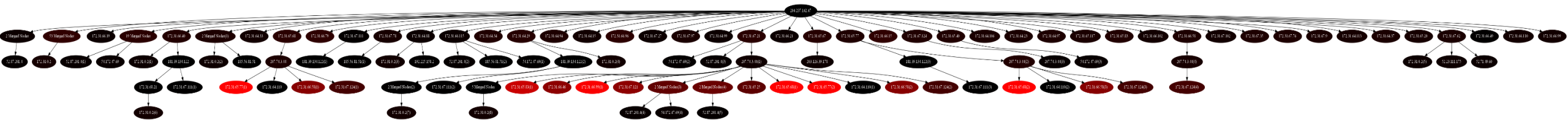
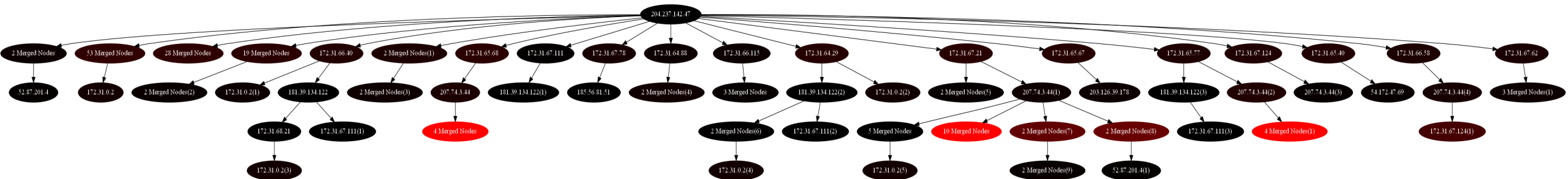# Results (Visual):
# Forward Tree 204.237.142.47

- Full Tree



- Tree with branches merged (R1)



- Tree with branches and leaves merged (R2)

# Metrics

- Visual Strain Reduction (VSR)

- Node Retention (NR)

- Threat Score Retention (TSR)

- Reduction Index (RI)
  - $RI = 3/(VSR^{-1} + NR^{-1} + TSR^{-1})$

# Results (Numerical)

| Reduction | Tree Set | VSR | NR | TSR | RI |
|---|---|---|---|---|---|
| MSB | Top 5 | 0.243 | 0.539 | 0.278 | 0.313 |
| | Random 5 | 0.352 | 0.553 | 0.254 | 0.349 |
| | Bottom 5 | 0.433 | 0.493 | 0.36 | 0.42 |
| MSL | Top 5 | 0.363 | 0.577 | 0.611 | 0.489 |
| | Random 5 | 0.282 | 0.824 | 0.799 | 0.499 |
| | Bottom 5 | 0.791 | 0.744 | 0.73 | 0.754 |
| TH | Top 5 | 0.009 | 1 | 0.999 | 0.026 |
| | Random 5 | 0 | 1 | 1 | 0 |
| | Bottom 5 | 0.037 | 1 | 0.983 | 0.103 |

MSB: Merge Sibling Branches
MSL: Merge Sibling Leaves
TH: Truncate Hypotrees
VSR: Visual Strain Reduction
NR: Node Retention
TSR: Threat Score Retention
RI: Reduction Index

UTSA.

# Dissertation Outline

- Introduction

- APIN: Alert Path Identification in Computer Networks

- AutoCRAT: Automatic Cumulative Reconstruction of Alert Trees

- Alert Tree Reduction and Visualization

- Conclusion

# Conclusion

Chapter 5

# Discussion

- APIN

  - Relies on <span style="color:red">network segmentation</span>

  - Dominates <span style="color:green">maintenance time</span>

- AutoCRAT

  - Relies on <span style="color:red">ordering assumption</span>

  - Dominates <span style="color:green">retrieval time</span>

- Reduction improves visualization

eric.ficke@utsa.edu

UTSA®