

# Automatically Inferring the Evolution of Malicious Activity on the Internet

**Shobha Venkataraman**

AT&T Research

**David Brumley**

Carnegie Mellon University

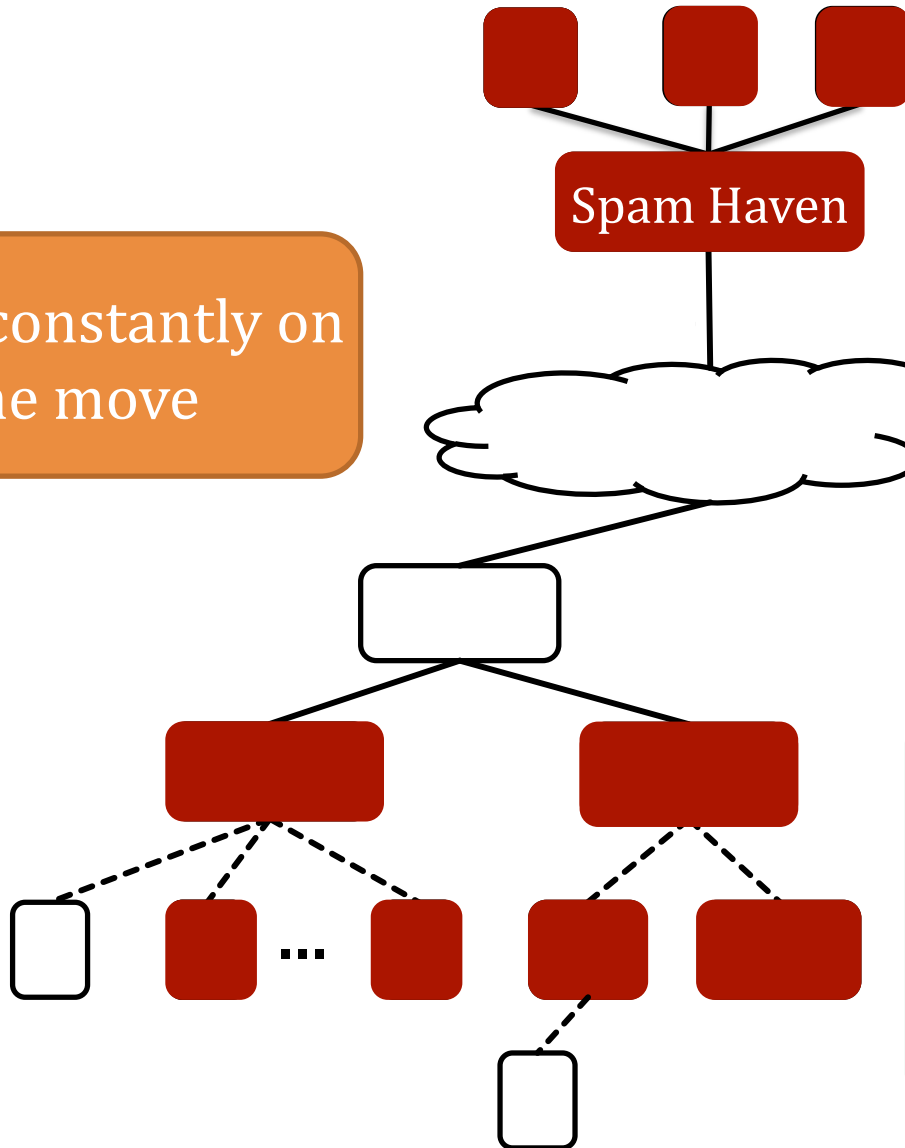
**Subhabrata Sen**

AT&T Research

**Oliver Spatscheck**

AT&T Research

Evil is constantly on the move



Labeled IP's from spam assassin, IDS logs, etc.

Tier 1

Our Goal:  
Characterize regions *changing* from bad to good ( $\Delta$ -good) or good to bad ( $\Delta$ -bad)

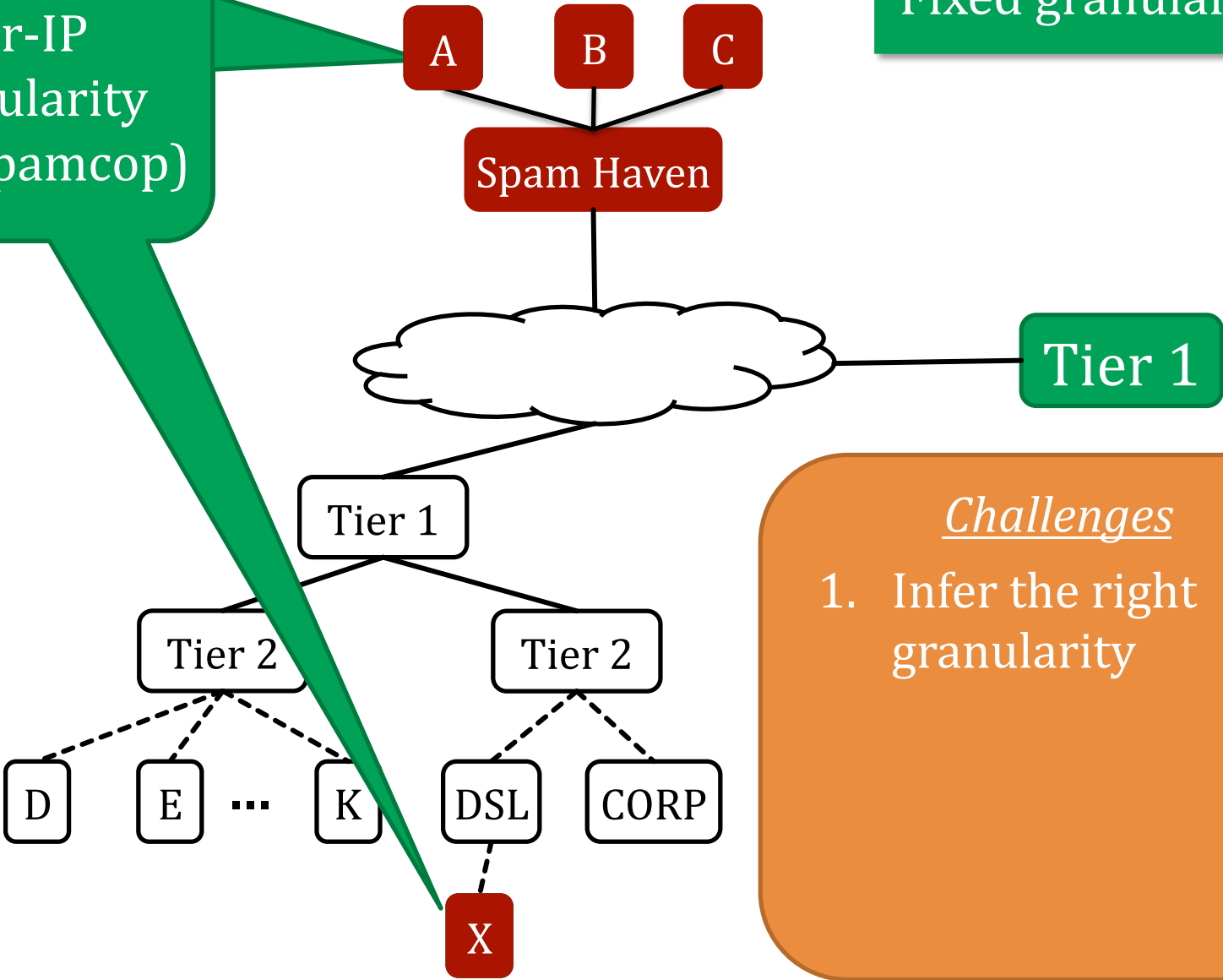
# Research Questions

Given a sequence of labeled IP's

1. Can we identify the specific regions on the Internet that have changed in malice?
2. Are there regions on the Internet that change their malicious activity more frequently than others?

Previous work:  
Fixed granularity

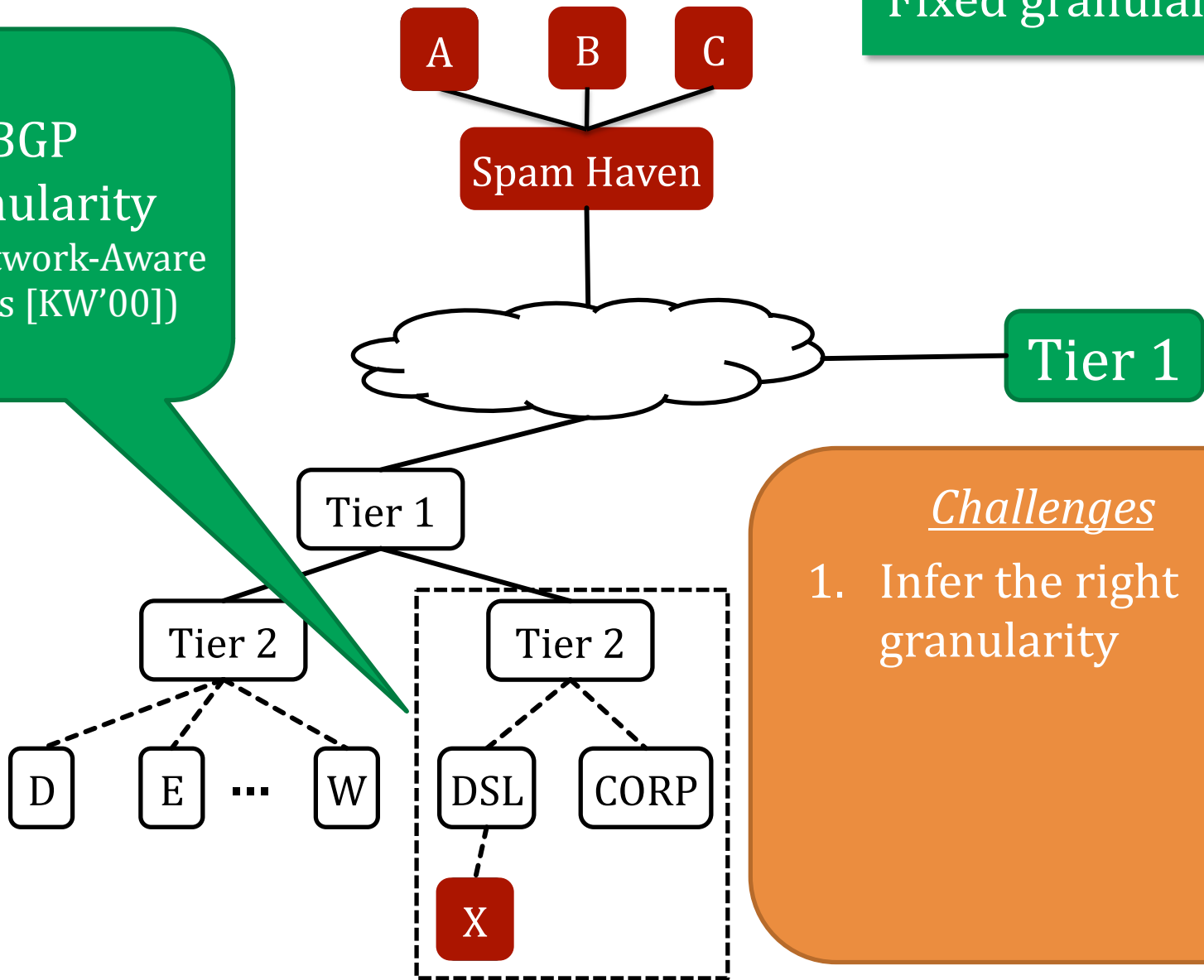
Per-IP  
Granularity  
(e.g., Spamcop)



Challenges  
1. Infer the right granularity

Previous work:  
Fixed granularity

BGP  
granularity  
(e.g., Network-Aware  
clusters [KW'00])



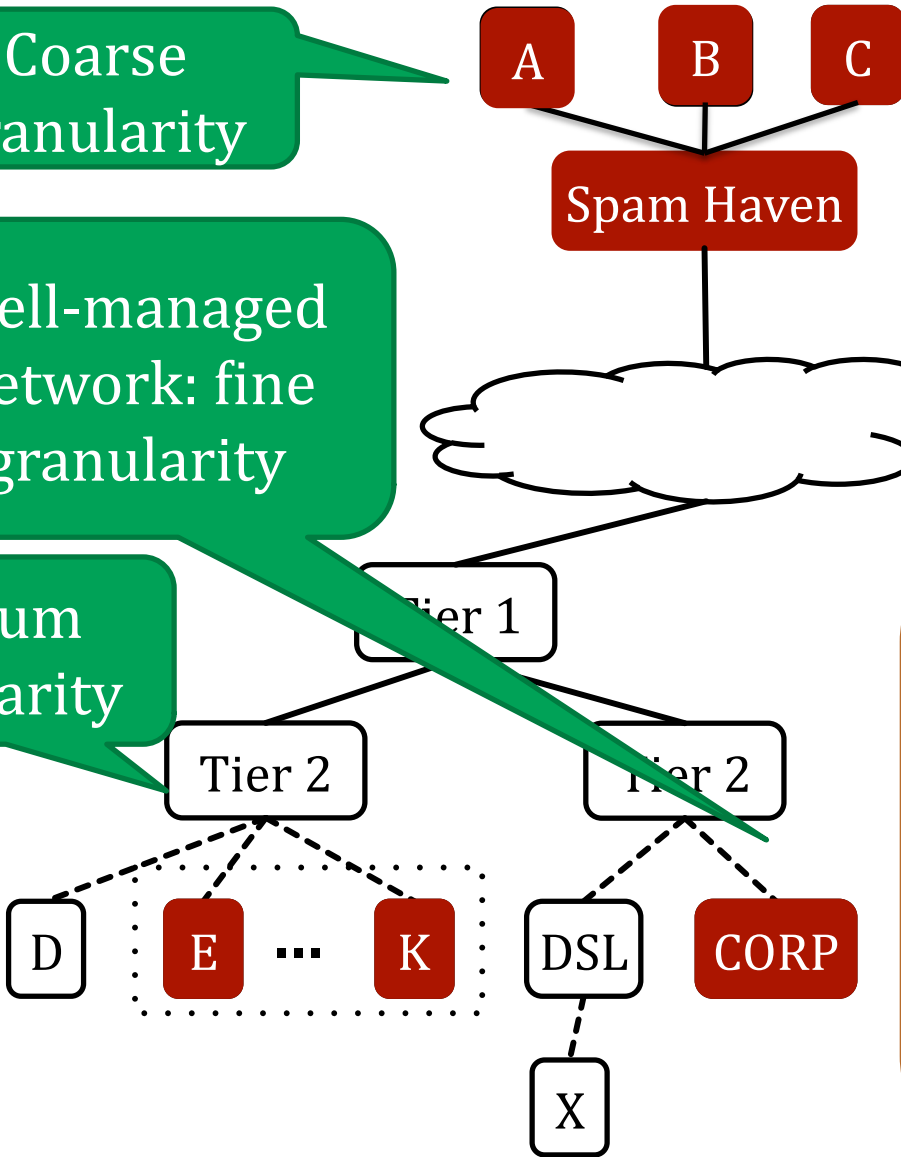
*Challenges*  
1. Infer the right  
granularity

Our work:  
Infer granularity

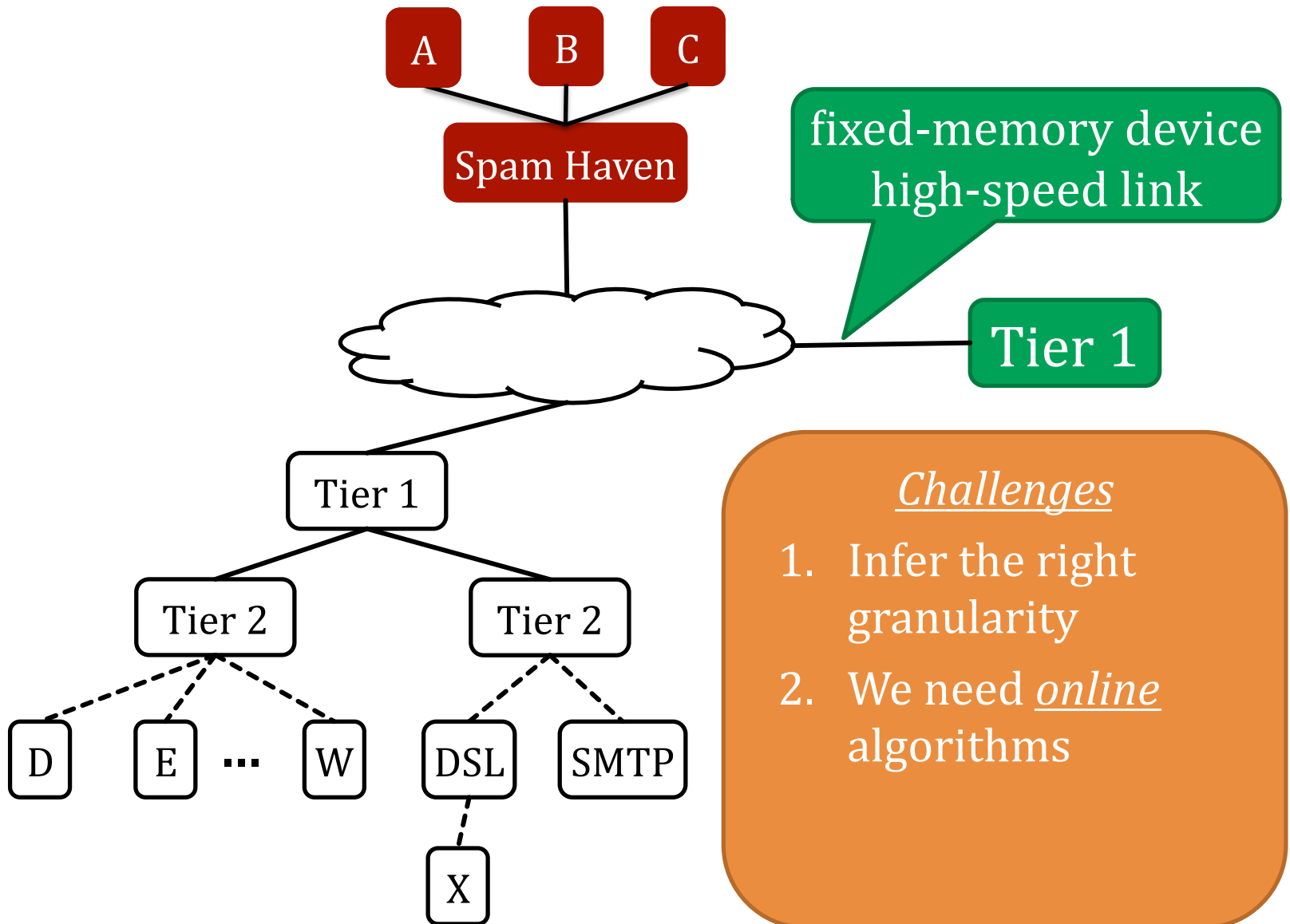
Coarse  
granularity

Well-managed  
network: fine  
granularity

Medium  
granularity



*Challenges*  
1. Infer the right  
granularity



# Research Questions

Given a sequence of labeled IP's

We Present

1. Can we identify the specific regions on the Internet that have changed in malice?

$\Delta$ -Change

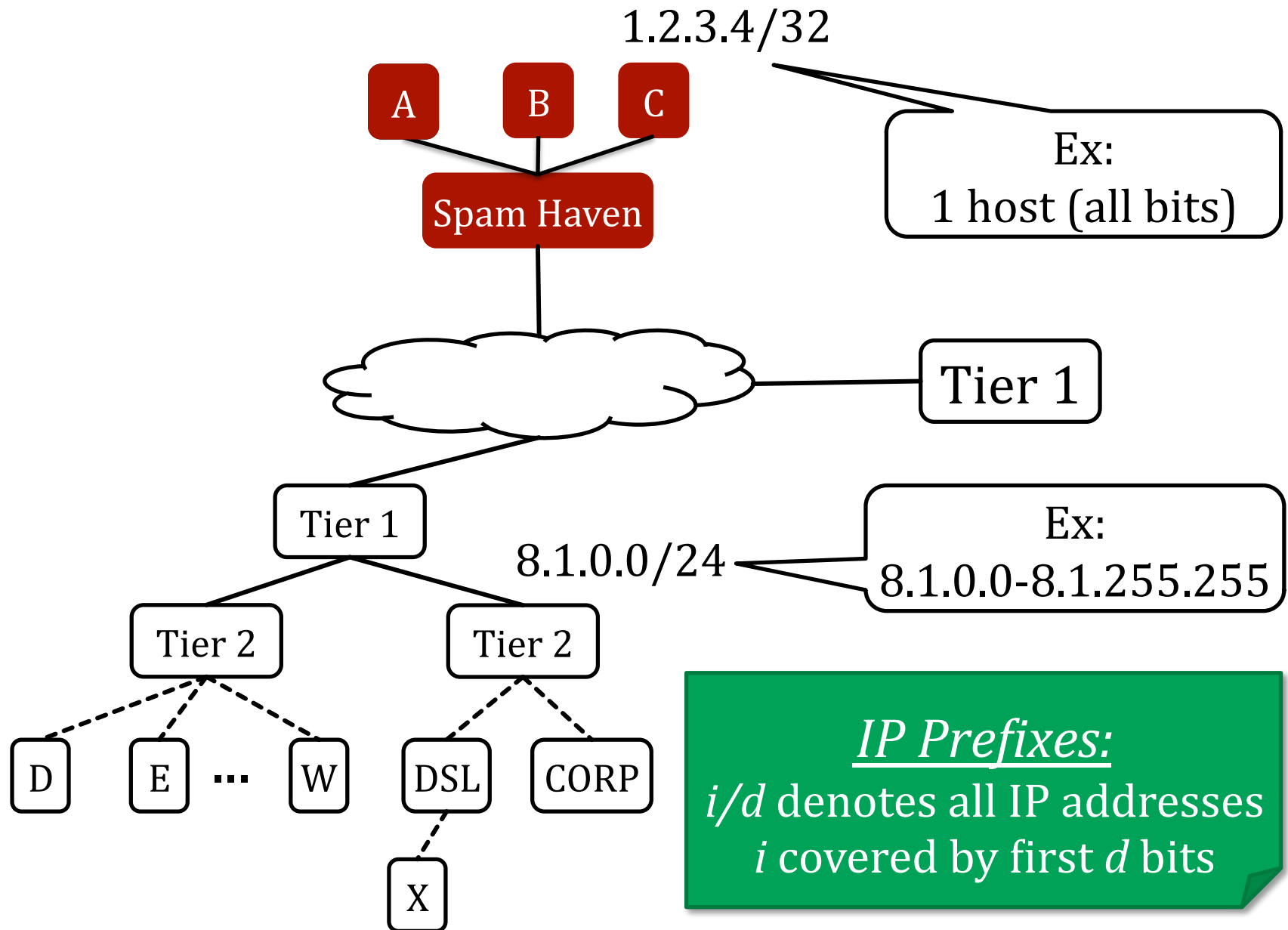
2. Are there regions on the Internet that change their malicious activity more frequently than others?

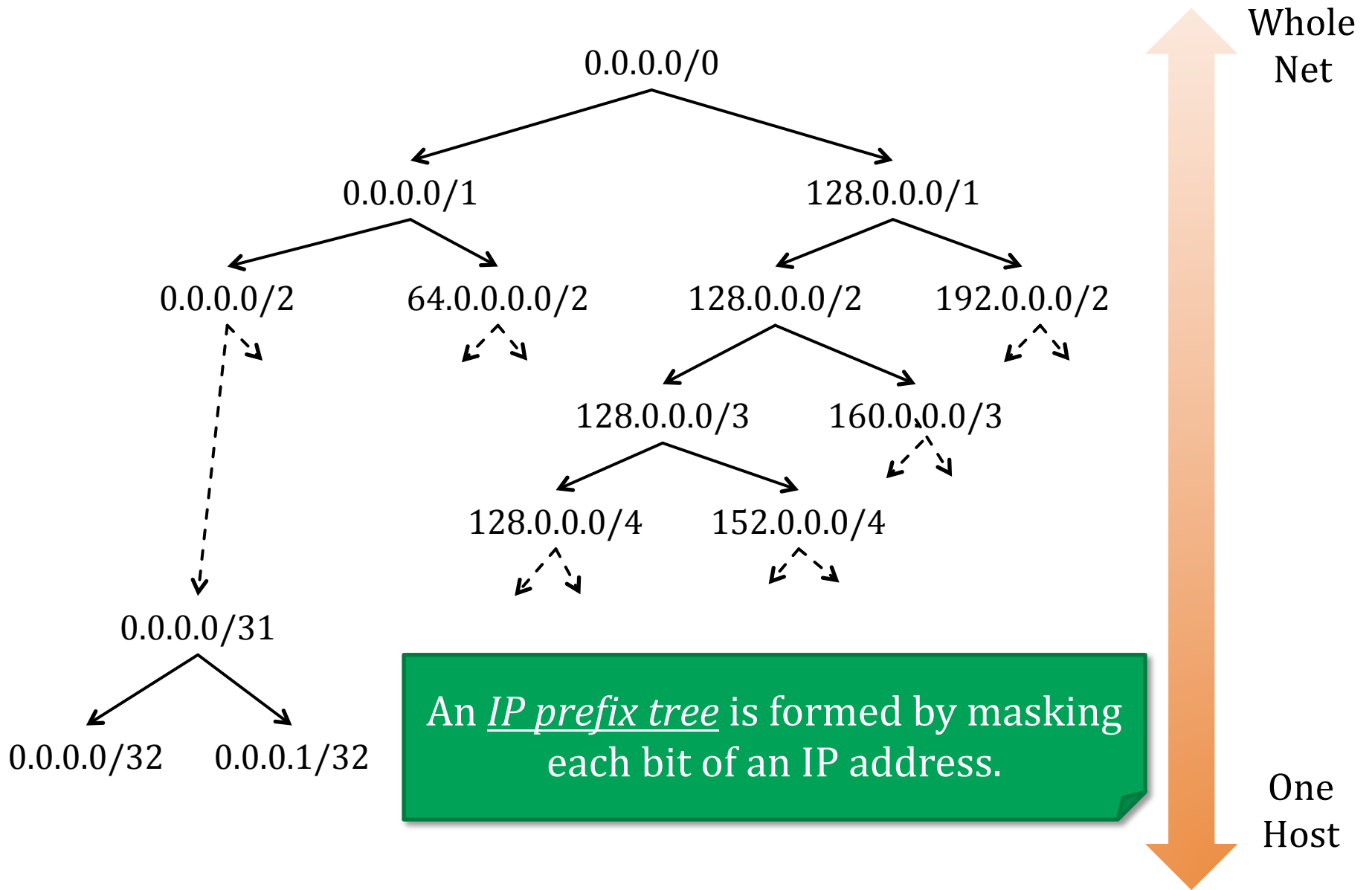
$\Delta$ -Motion



# Background

1. IP Prefix trees
2. TrackIPTree Algorithm







# TrackIPTree Algorithm

[VBSSS'09]

*In:* stream of  
labeled IPs

...  $\langle ip_4, + \rangle$   $\langle ip_3, + \rangle$   $\langle ip_2, + \rangle$   $\langle ip_1, - \rangle$

TrackIPTree

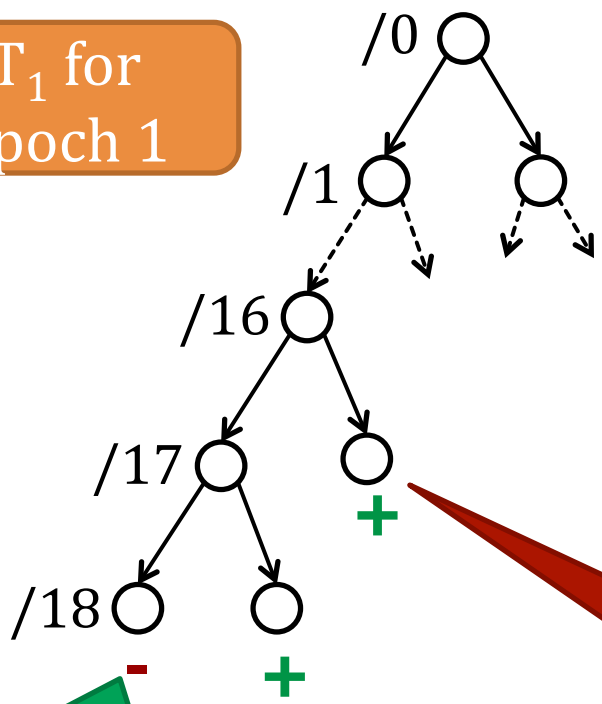
*Out:* k-IPTree

# $\Delta$ -Change Algorithm

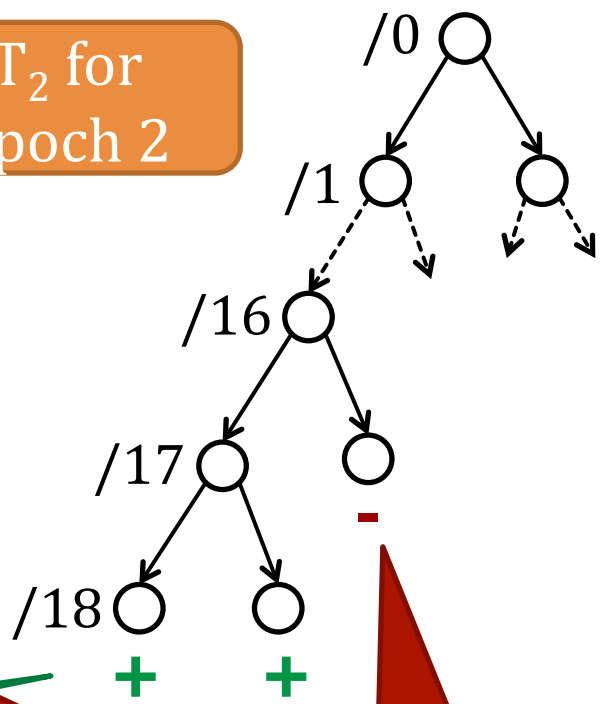
1. Approach
2. What doesn't work
3. Intuition
4. Our algorithm

Goal: identify online the specific regions on the Internet that have changed in malice.

$T_1$  for epoch 1

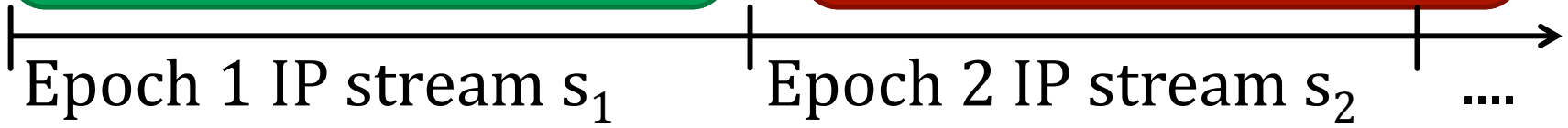


$T_2$  for epoch 2



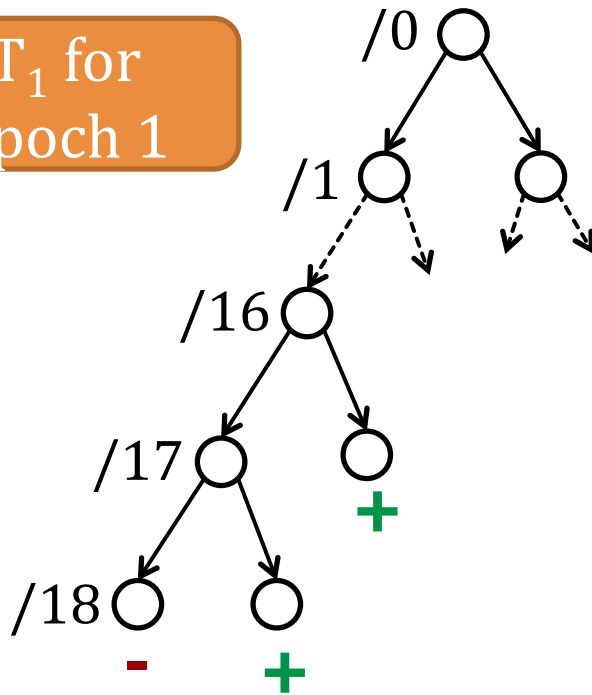
$\Delta$ -Good:  
A change from bad to good

$\Delta$ -Bad:  
A change from good to bad



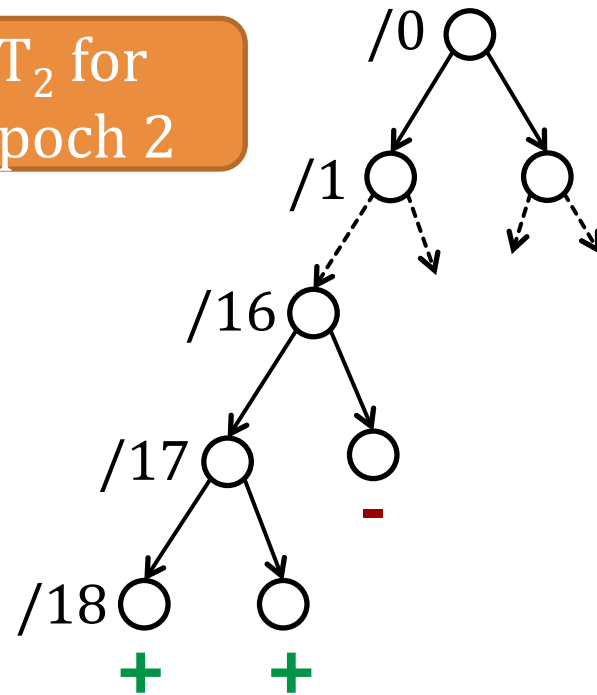
Goal: identify online the specific regions on the Internet that have changed in malice.

$T_1$  for epoch 1



False positive:  
Misreporting that a change occurred

$T_2$  for epoch 2

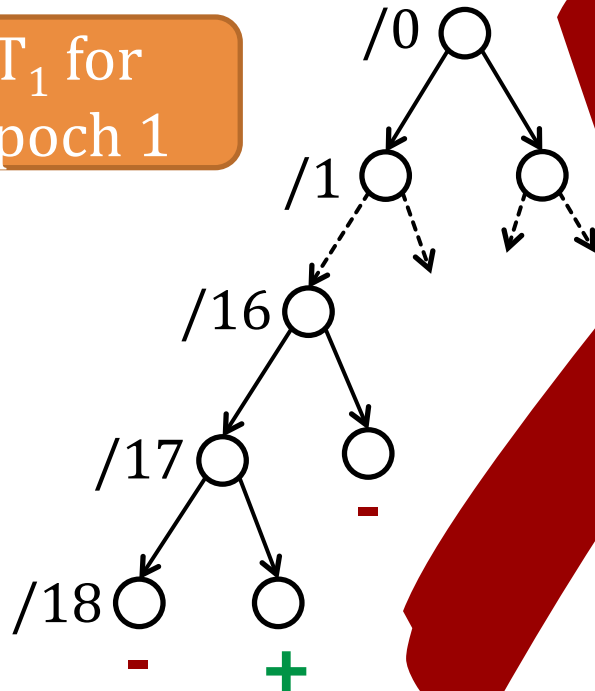


False Negative:  
Missing a real change

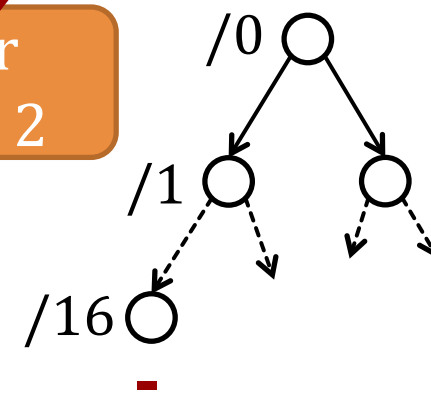


Goal: identify online the specific regions on the Internet that have changed in malice.

$T_1$  for epoch 1



$T_2$  for epoch 2



Different Granularities!

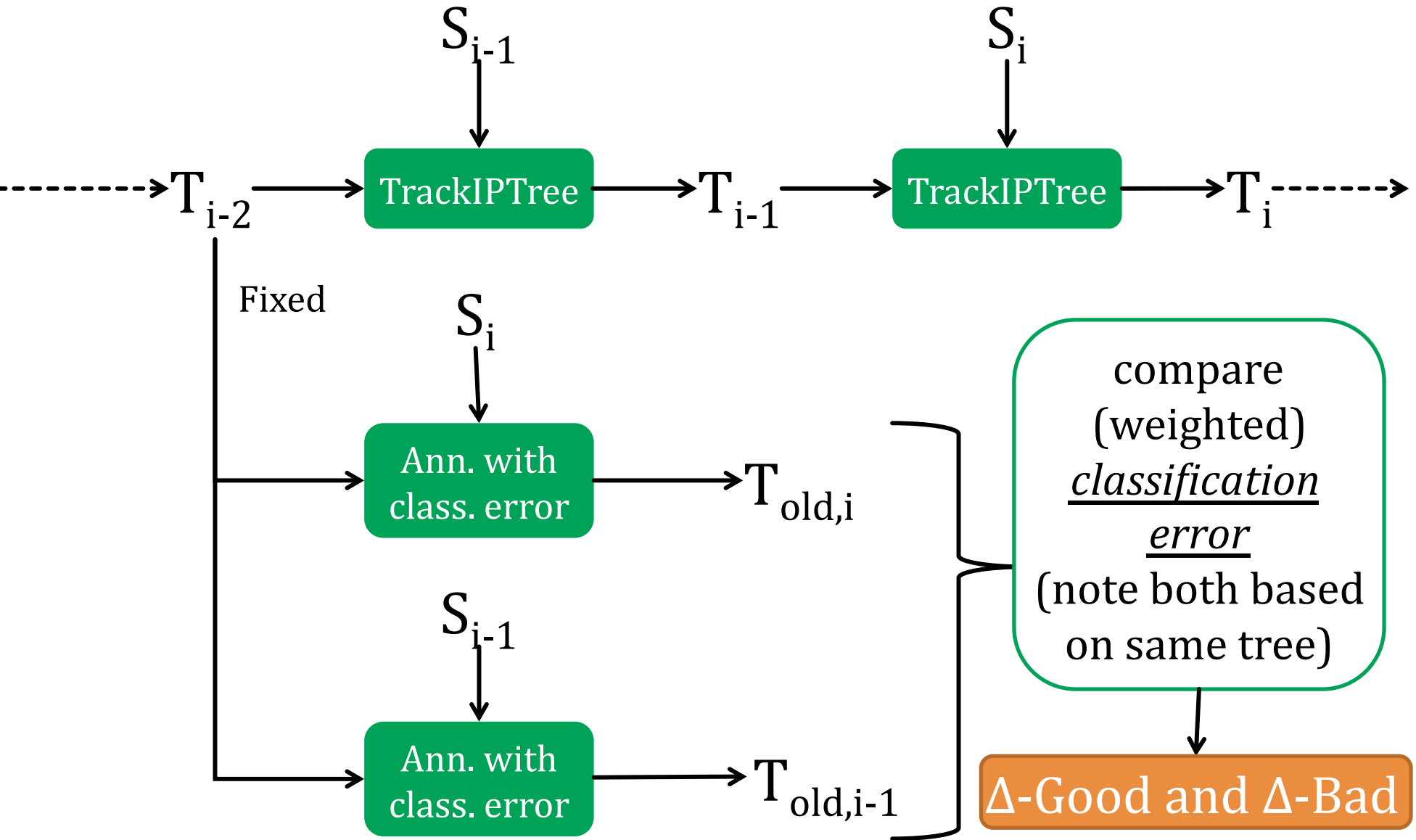
Idea: divide time into epochs and diff

- Use TrackIPTree on labeled IP stream  $s_1$  to learn  $T_1$
- Use TrackIPTree on labeled IP stream  $s_2$  to learn  $T_2$
- Diff  $T_1$  and  $T_2$  to find  $\Delta$ -Good and  $\Delta$ -Bad

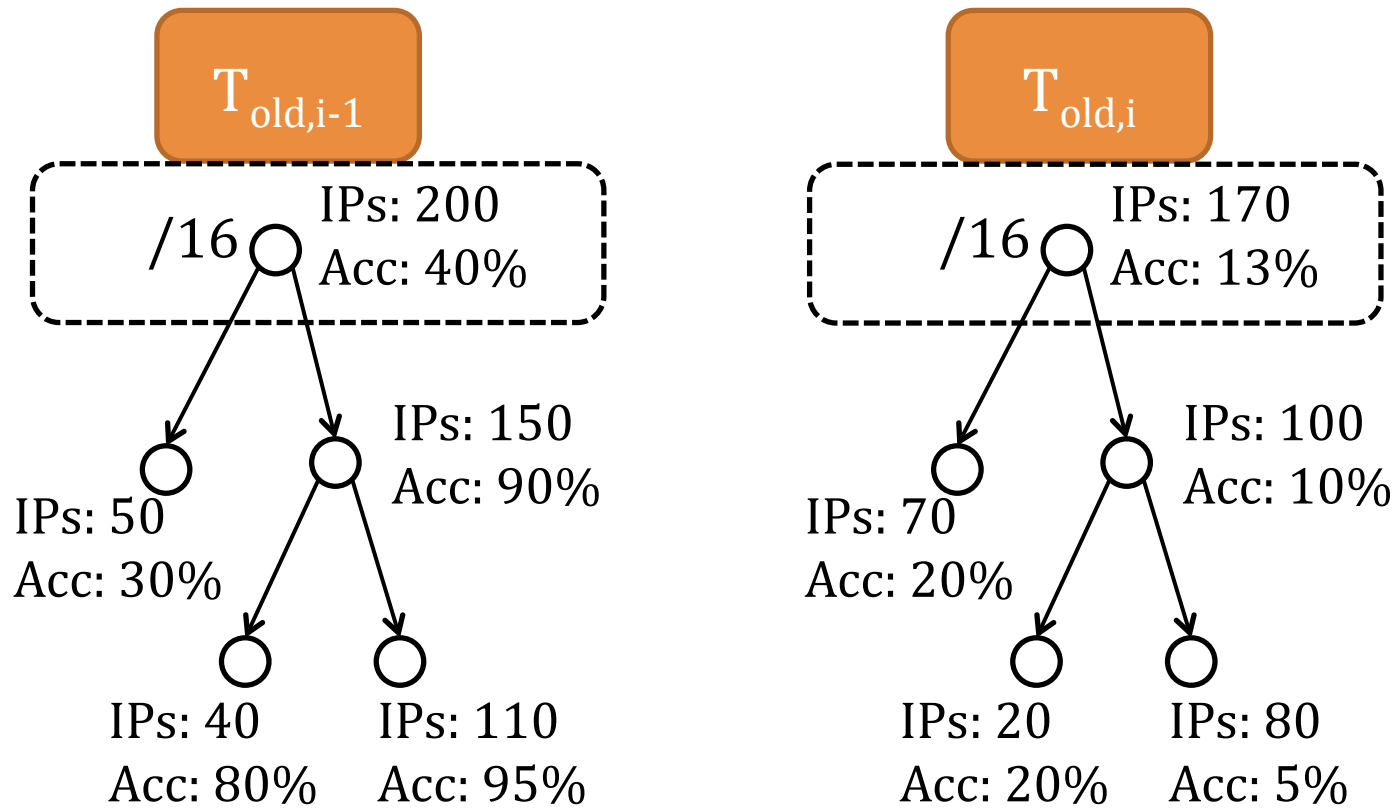
Goal: identify online the specific regions on the Internet that have changed in malice.

$\Delta$ -Change Algorithm Main Idea:  
Use classification errors between  $T_{i-1}$  and  $T_i$   
to infer  $\Delta$ -Good and  $\Delta$ -Bad

# $\Delta$ -Change Algorithm

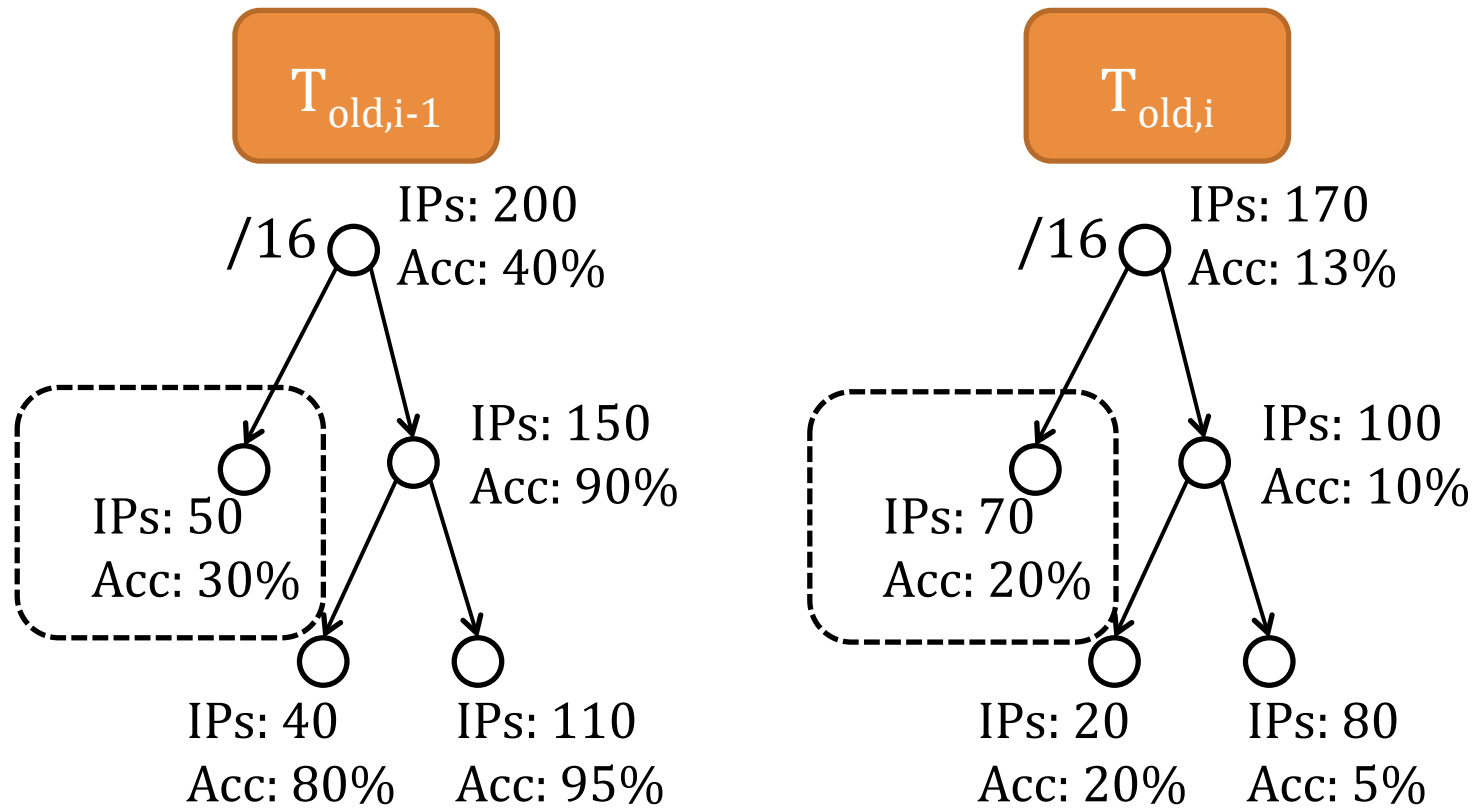


# Comparing (Weighted) Classification Error



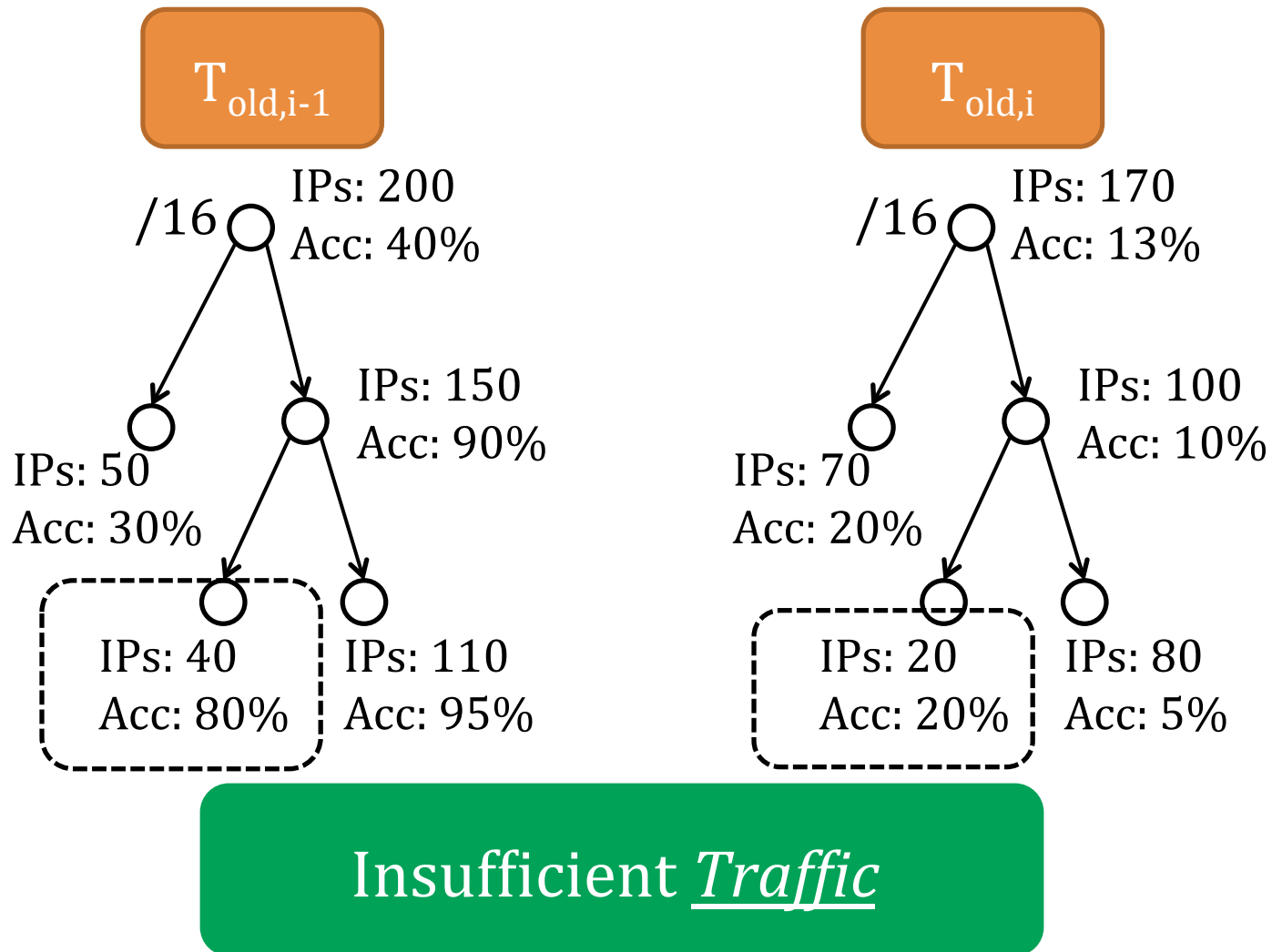
$\Delta$ -Change Somewhere

# Comparing (Weighted) Classification Error

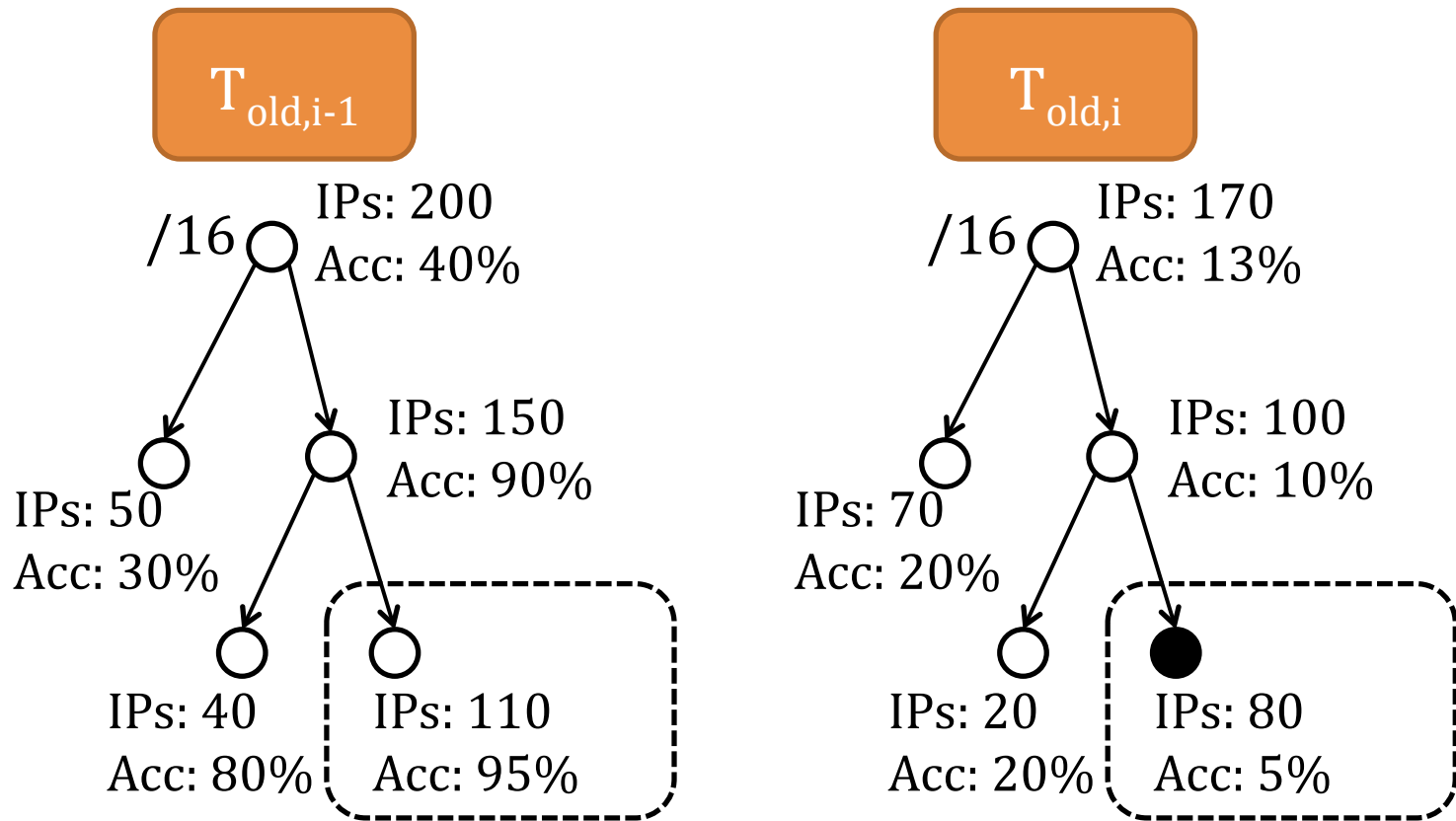


Insufficient Change

# Comparing (Weighted) Classification Error



# Comparing (Weighted) Classification Error



$\Delta$ -Change Localized

# Evaluation

1. What are the performance characteristics?
2. Are we better than previous work?
3. Do we find cool things?



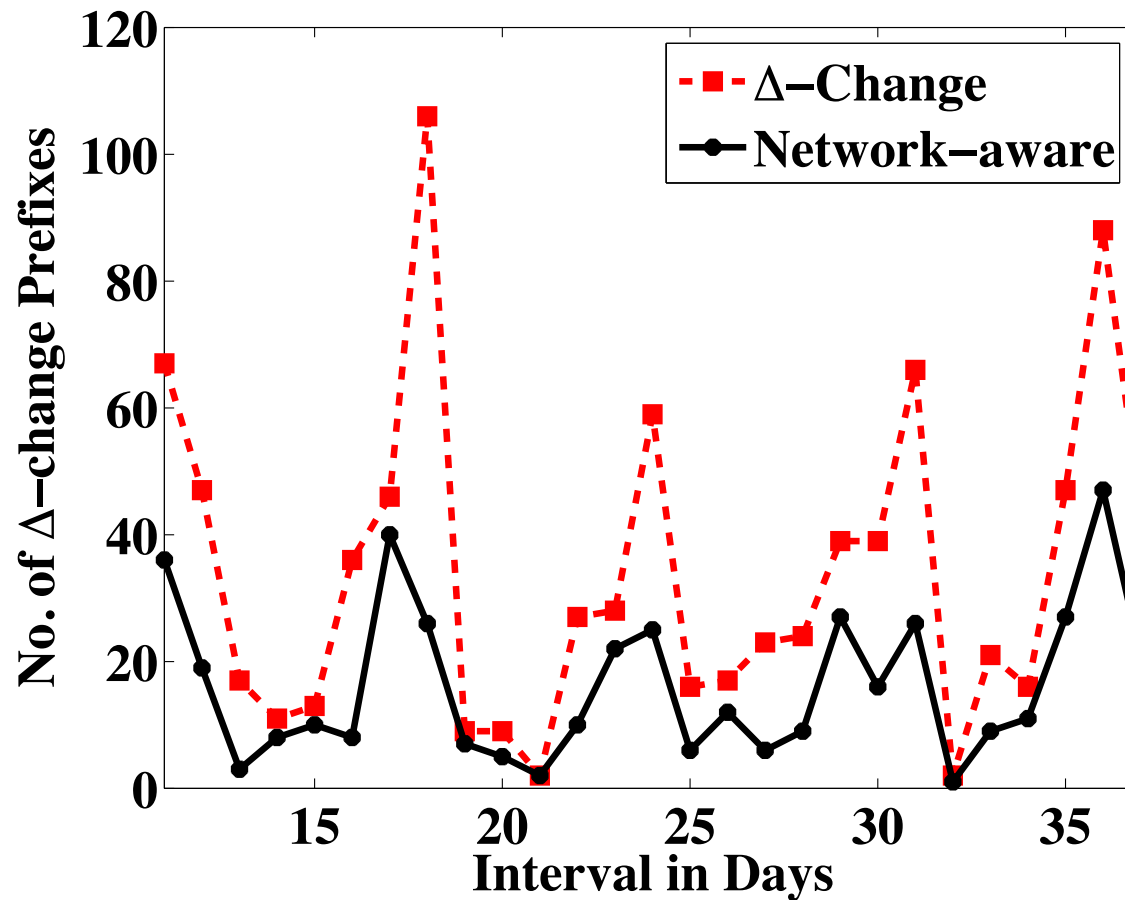
# Performance

In our experiments, we :

- let  $k=100,000$  (k-IPTree size)
- processed 30-35 million IPs (on day's traffic)
- using a 2.4 Ghz Processor

Identified  $\Delta$ -Good and  $\Delta$ -Bad  
in <22 min using <3MB memory

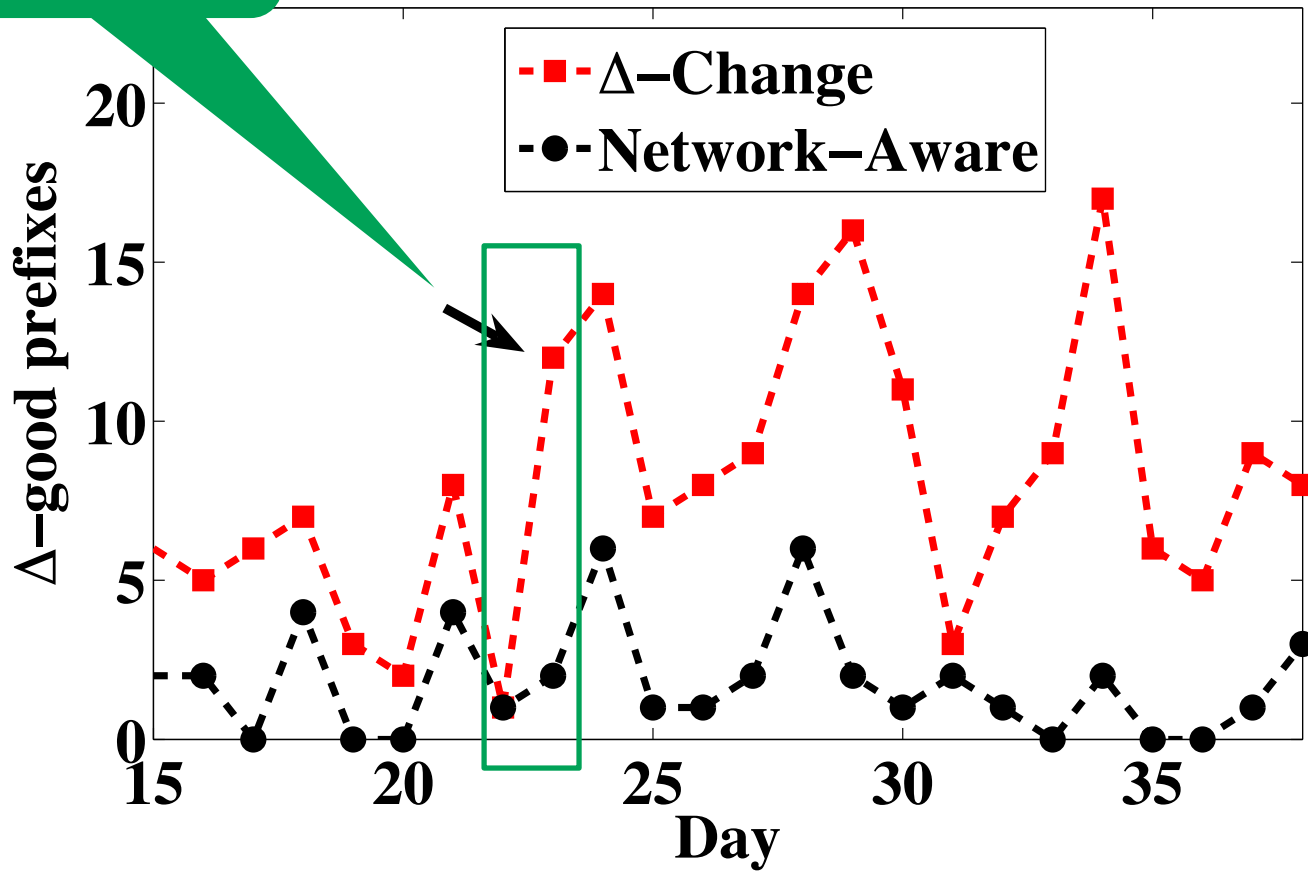
# How do we compare to network-aware clusters? (By Prefix)



2.5x as many changes  
on average!

# Spam

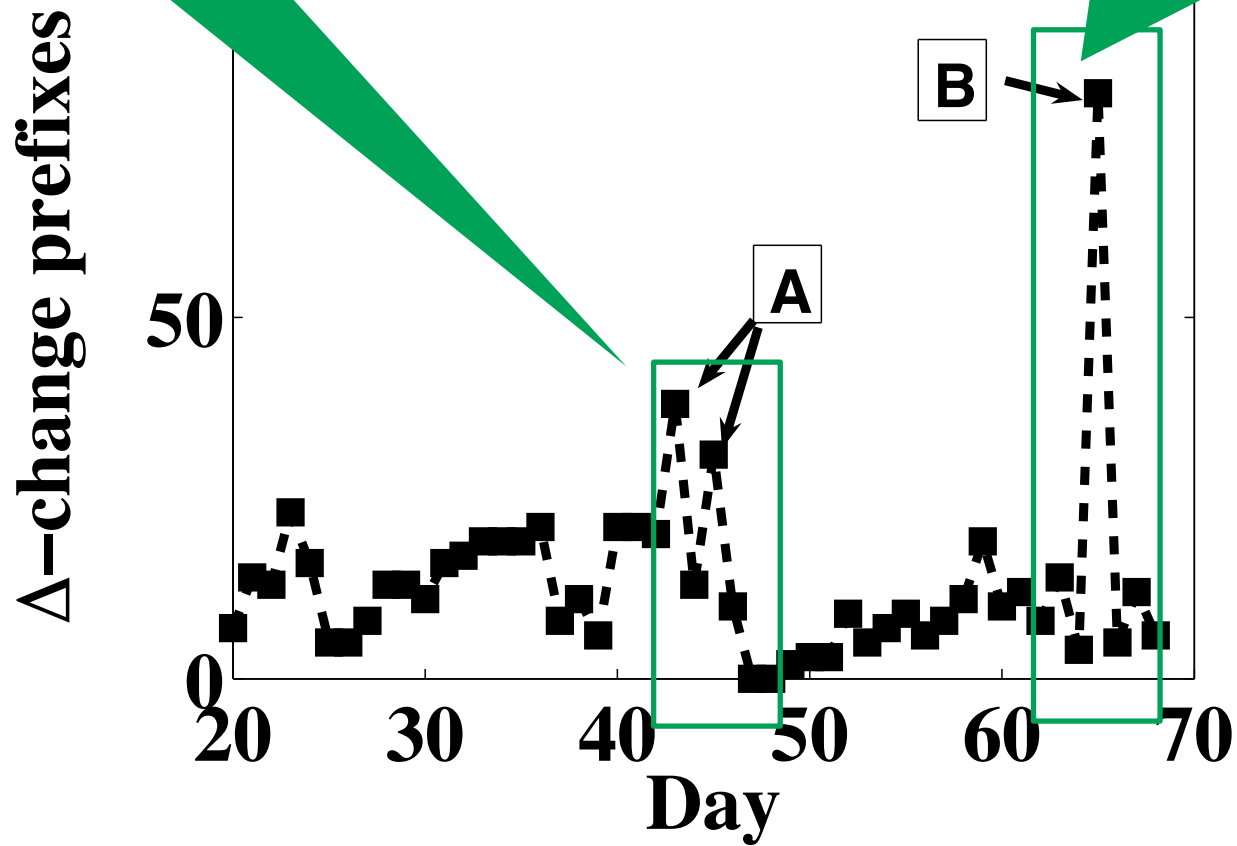
Grum botnet  
takedown



# Botnets

22.1 and 28.6 thousand new DNSChanger bots appeared

38.6 thousand new Conficker and Sality bots



# Caveats and Future Work

“For any distribution on which an ML algorithm works well, there is another on which it works poorly.”

– The “No Free Lunch” Theorem



Our algorithm is efficient and works well in practice.

....but a very powerful adversary could fool it into having many false negatives. A formal characterization is future work.

# Conclusion

$\Delta$ -Change and  $\Delta$ -Motion: two new *online* algorithms for capturing how malice evolves on the internet

- Scalable
- Discovers right IP granularity
- Finds cool changes



END