

DNS backscatterの紹介

国立情報学研究所

福田 健介

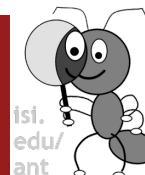
Detecting malicious activity with DNS backscatter

Kensuke Fukuda (NII/Sokendai)

John Heidemann (USC/ISI)

Appeared in IMC2015

http://www.fukuda-lab.org/publications/FH_imc2015.pdf





- 2014: Heartbleed
 - A bug in critical millions of Internet hosts
 - Security researchers scanned to find unpatched servers
- Question: Who else was scanning?
 - Criminals?
 - Black hats?
 - Others?

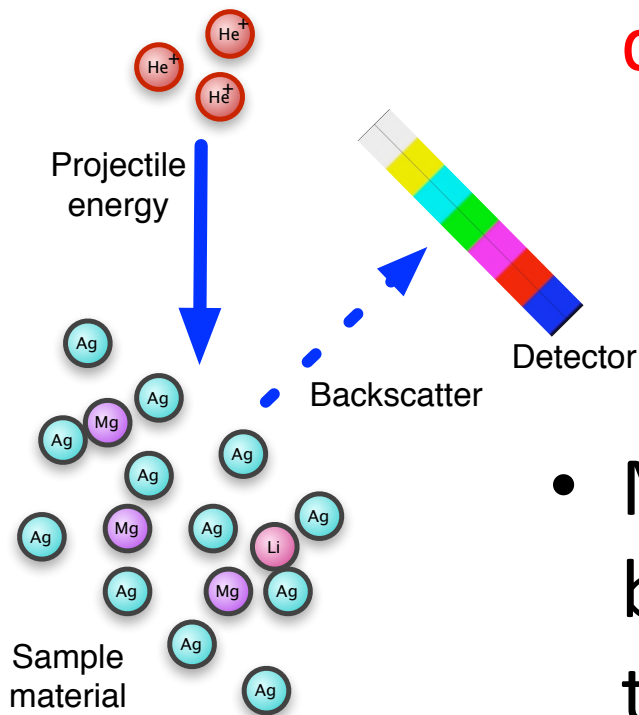
Goal: Finding Originators of Network-wide events

- Large-scale event involves many Internet hosts
 - Malicious: Scan, Attacks, Spams
 - Benign: CDN, Web crawler, DNS, NTP, Updates
 - Border: Ad tracking
- Importance of monitoring those events
 - Malicious: security consideration
 - Benign: stability of infrastructure

Our contribution

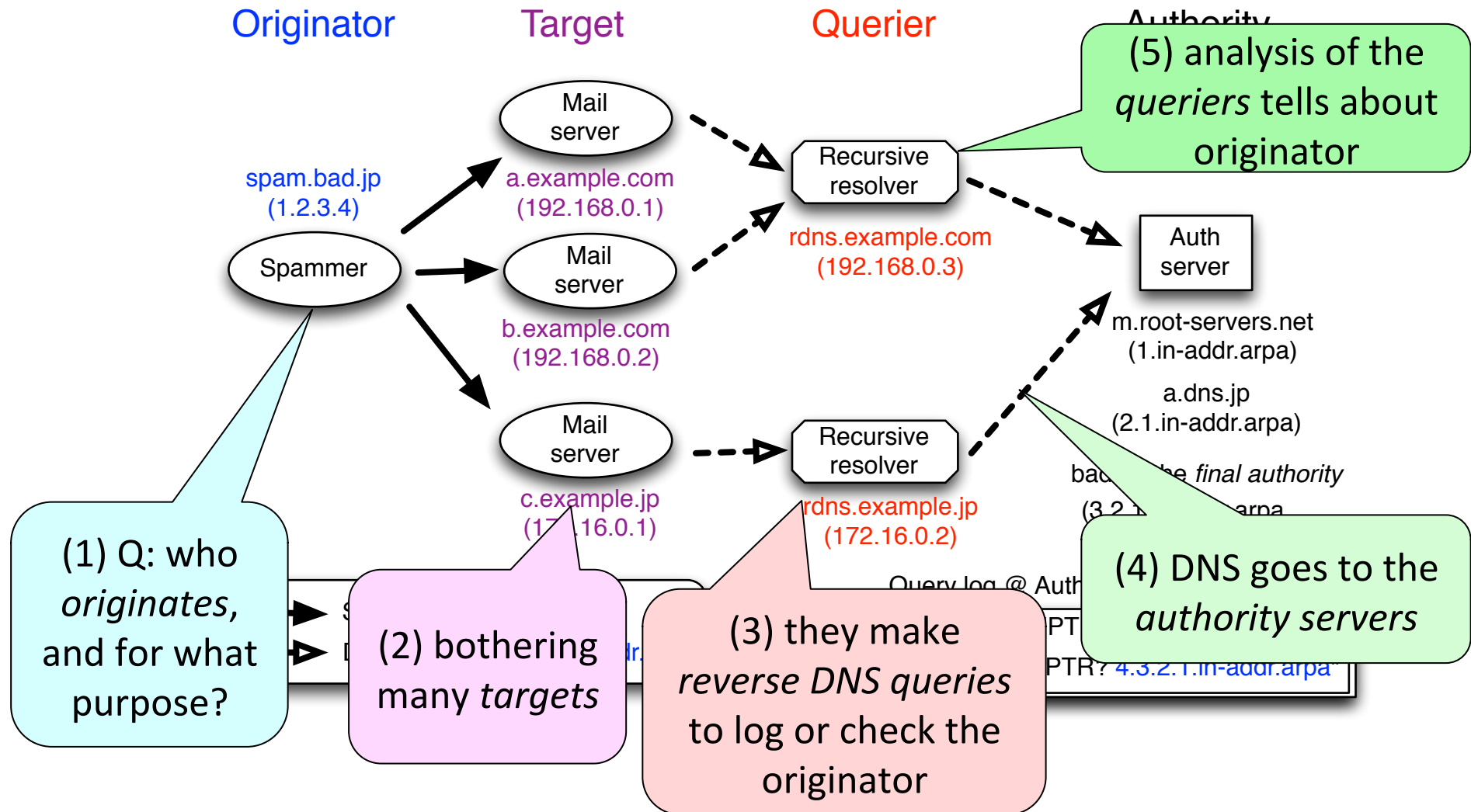
- New data source - **DNS backscatter** - to identify network-wide activity
 - Deployable
 - Privacy friendly
 - Robust against malicious source
- Validate with several DNS authoritative servers
- Evaluate over time: 6 months dataset

Key idea of DNS backscatter



- Large event triggers **reverse DNS queries** near target automatically
 - SMTP server: hostname of **spammer**
 - Firewall: hostname of **scanner**
 - Web server: hostname of **web crawler**
- Many reverse DNS queries (DNS backscatter) at **auth server** are hint to identify events

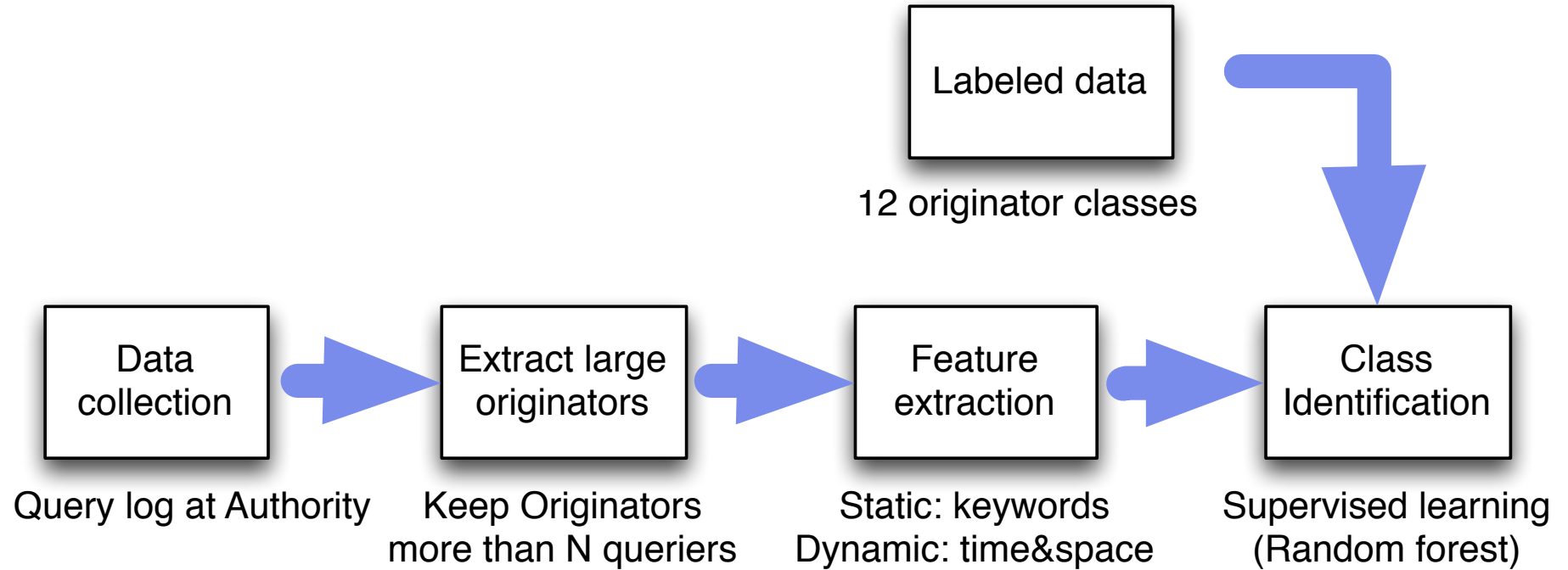
Detecting Events through DNS backscatter



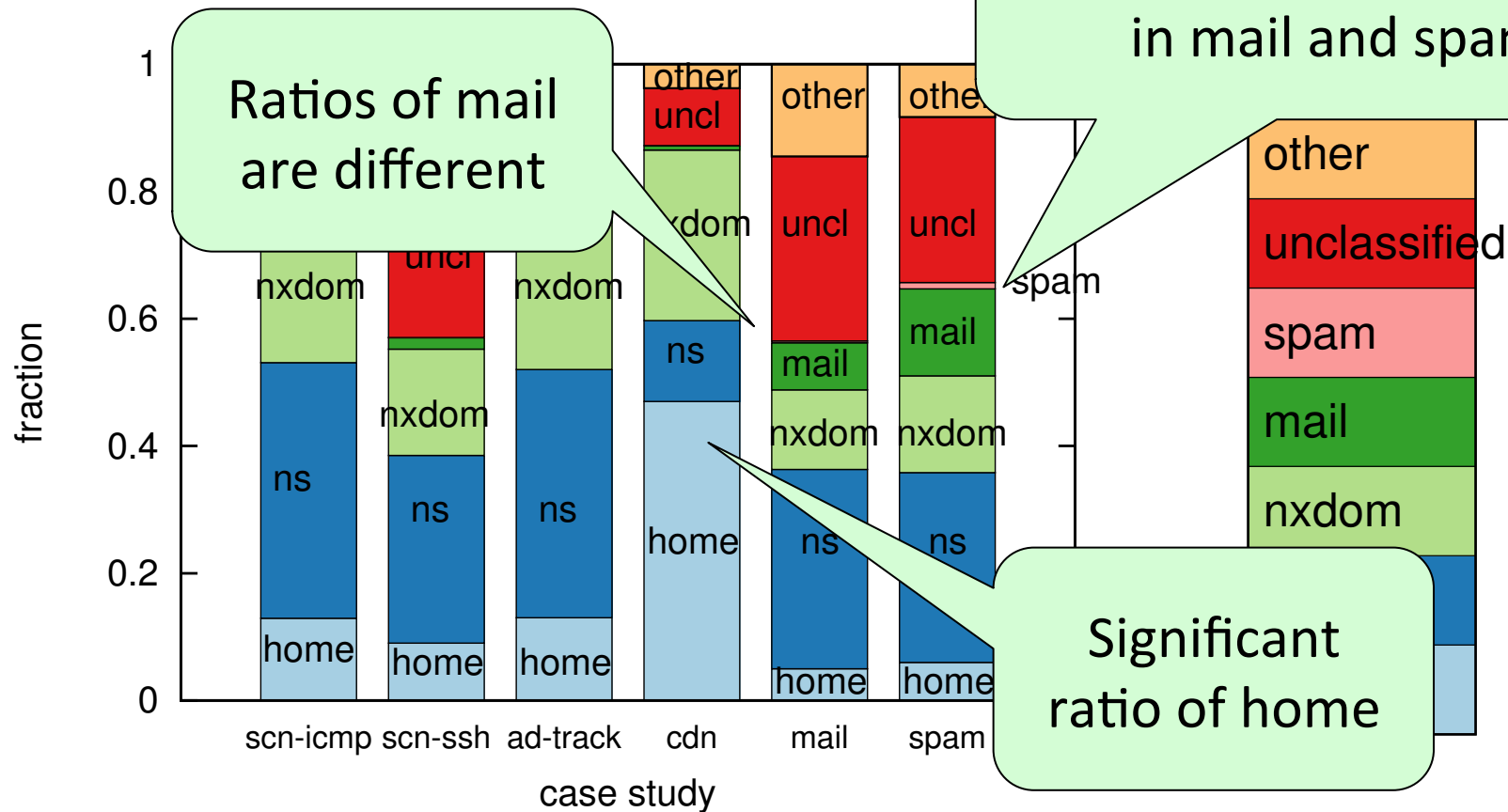
Advantages

- Deployable
 - **Centralized** monitoring at DNS authority
- Privacy friendly
 - Information is on **queriers** NOT originators
 - Reverse queries are generated **automatically**
 - Focus on **large events** (ignore small users)
- Robust against malicious originators
- Can infer different class of originator (e.g., scanner) with Machine Learning

Identification process



Discriminative power



Different mixes of features allow distinguishing different classes of events

Picking the best ML algorithm

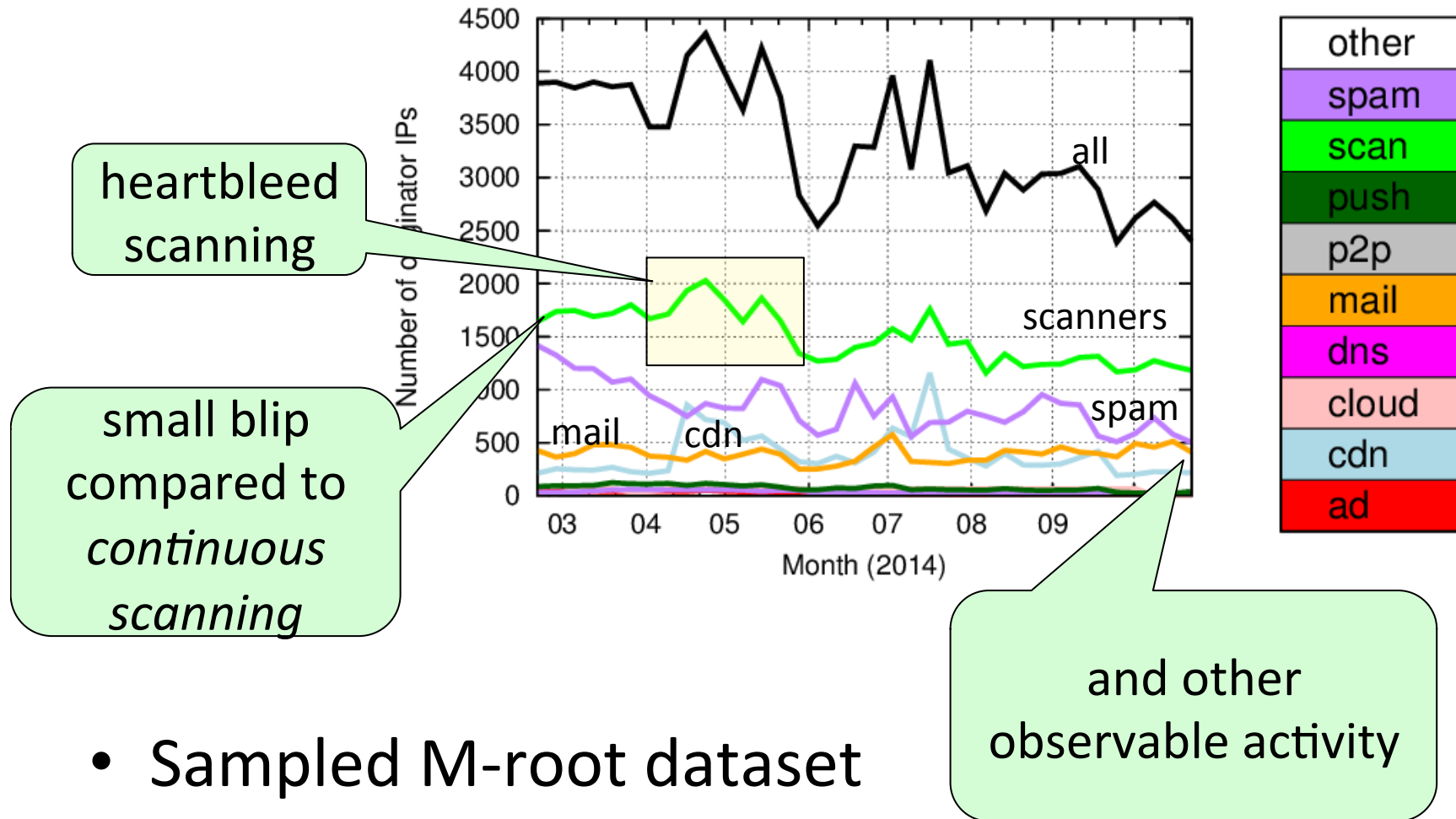
| dataset | algorithm | accuracy | precision | recall | F1-score |
|--------------------|-----------|-------------|-------------|-------------|-------------|
| JP ditl | CART | 0.66 | 0.63 | 0.60 | 0.61 |
| | RF | 0.78 | 0.82 | 0.76 | 0.79 |
| | SVM | 0.73 | 0.74 | 0.71 | 0.73 |
| B post- ditl | CART | 0.48 | 0.48 | 0.45 | 0.46 |
| | RF | 0.62 | 0.66 | 0.60 | 0.61 |
| | SVM | 0.38 | 0.50 | 0.32 | 0.40 |
| M ditl | CART | 0.53 | 0.52 | 0.49 | 0.51 |
| | RF | 0.68 | 0.74 | 0.63 | 0.68 |
| | SVM | 0.60 | 0.67 | 0.60 | 0.64 |
| M sampled | CART | 0.61 | 0.66 | 0.60 | 0.62 |
| | RF | 0.79 | 0.81 | 0.76 | 0.78 |
| | SVM | 0.72 | 0.77 | 0.71 | 0.74 |

RandomForest
is best

Hope to improve with
better training data

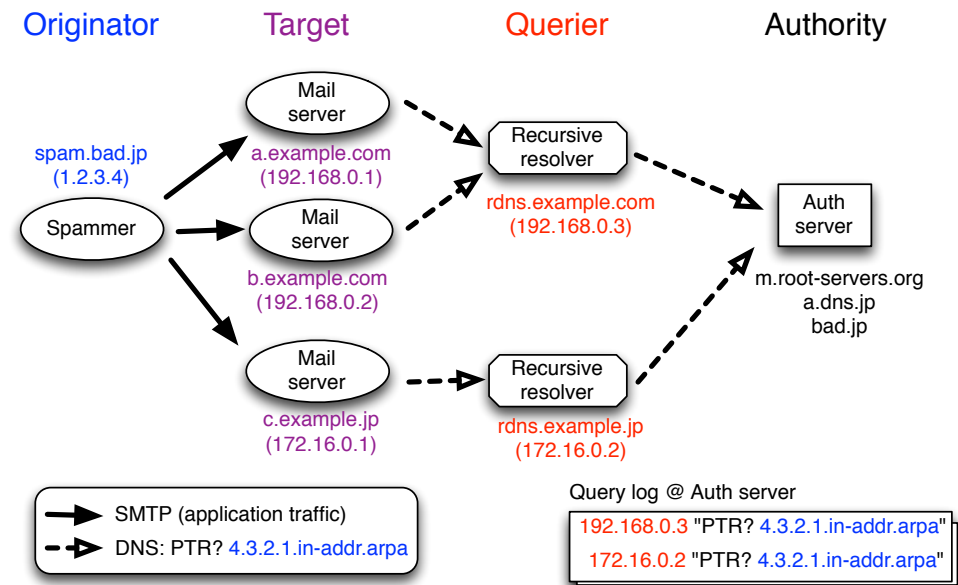
- Cross validation with 3 ML algorithms
- Num classes: 12, labeled data:200-800
- Precision: 70-80% (imbalanced dataset problem)

Finding Network-wide events over time



Conclusion

- **DNS backscatter** - a new data source for Internet-wide events
- Advantages:
 - Deployable
 - Privacy-friendly
 - Reasonable accuracy
- Longitudinal results



DNS operators may apply this to detect large events!