# Knock It Off:
## Profiling the Online Storefronts of Counterfeit Merchandise
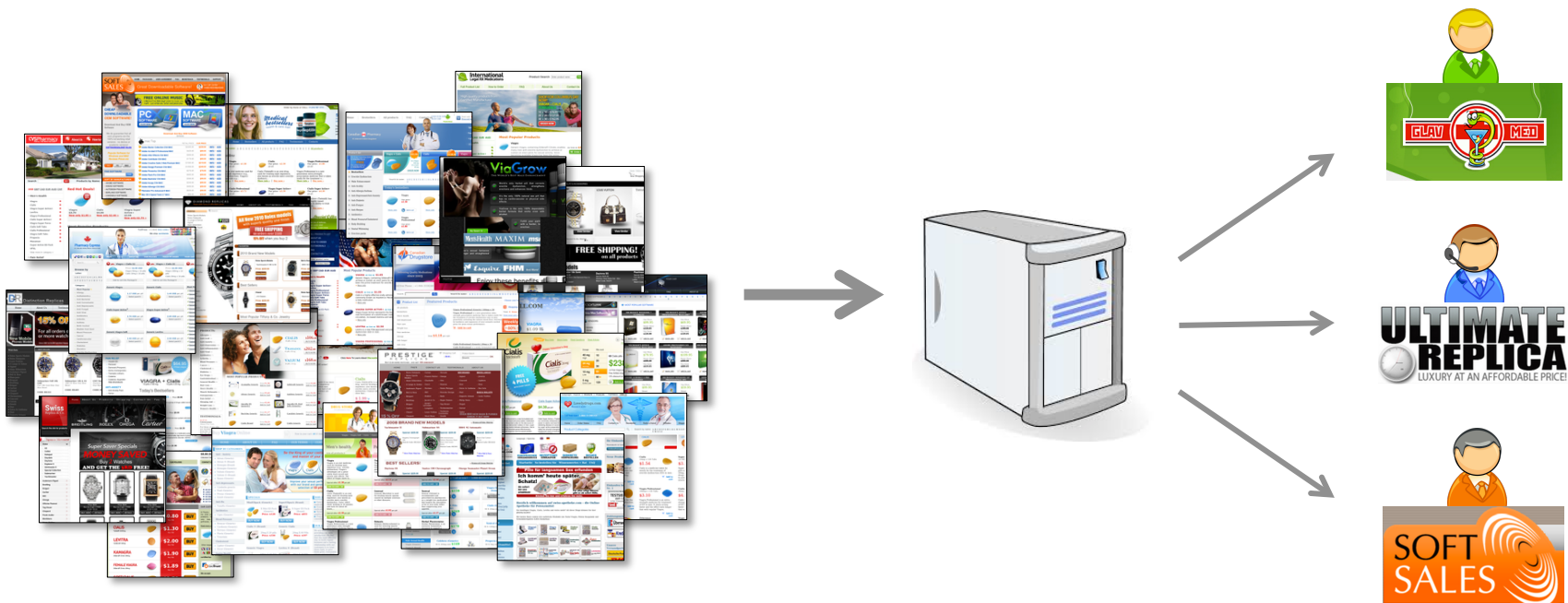
Matt Der
Lawrence Saul
Stefan Savage
Geoff Voelker

**UCSDCSE**
Computer Science and Engineering

*cns*
center for networked systems

# Problem in a nutshell

Behind the many online storefronts for counterfeit goods lurk a small handful of sophisticated criminal operations.

How can automated, data-driven methods help to identify and target them?

# Counterfeit online storefronts

# Counterfeit online storefronts

# Counterfeit online storefronts

# Who is running the store?



"affiliate programs"

# What is an affiliate program?

- Illegitimate business that sells counterfeit goods: millions of $$$ of revenue per month

- Manage Web sites that serve as online storefronts

- Enlist spammers to advertise their storefronts via bulk email

- Contract out payment & fulfillment services

# Click Trajectories



- Bottleneck is payment processing

    95% of spam-advertised pharmaceutical, replica, and software products are monetized through only a handful of merchant banks

- If affiliate program can't process credit cards, business collapses

Figure from Levchenko et al, 2011

# Key insight

**100s of thousands of storefronts**

# Key insight

**100s of thousands of storefronts**

**_dozens_ of affiliate programs**

# Our work

- **Goal**:  classify storefronts by affiliate program; disrupt their operation to undermine spam business model



- **Approach**:  HTML bag-of-words, nearest neighbor classification (**automated** system)

- **Takeaway**:  highly accurate — even with simple classifier & limited labeled examples

# Challenges

# Challenges

1. Web pages that render very differently are often linked to the same affiliate program



GlavMed

EvaPharmacy

# Challenges

2. Difficulty in acquiring training data

# Challenges

2. Difficulty in acquiring training data

2. Difficulty in acquiring training data

2. Difficulty in acquiring training data

# Challenges

2. Difficulty in acquiring training data

# Challenges

2. Difficulty in acquiring training data

# Challenges

2.   Difficulty in acquiring training data



expert labeling is
slow & tedious!

# A role for machine learning

- Security experts labeled 178k storefronts
  - Estimated **~200 person-hours**
  - Painstaking manual process
    - Inspect HTML source for signals, encode with regular expressions

- **NOT** once-and-for-all effort
  - Storefronts change over time, regexs go stale

- Ripe opportunity for machine learning — a more automated approach to aid security practitioners

# Feature extraction



HTML src

```
<html>
…
</html>
```

screenshot

DNS records

DNS

# Feature extraction



HTML src          screenshot          DNS records

```
<html>
…
</html>
```

- Affiliate programs use in-house software engines to generate storefront templates
- HTML contains distinctive signatures
- Bag-of-words on HTML – automated!

# Data set



- classes:            44
- labeled exs:      178k
- largest class:    58k
- smallest class:  2

Data is high-dimensional & sparse

# Proof-of-concept experiment

- **Question**: are these HTML features enough to distinguish affiliate programs:
  - From one another?
  - From the noise in the Web crawl?
    - 43k unlabeled Web pages → "other" class

- Favorable setting: plenty of labeled data
- 45-way 1-nearest neighbor classification,
  10 random 70/30 train/test splits

# Proof-of-concept experiment

- **Question**:  are these HTML features enough to distinguish affiliate programs:
  - From one another?
  - From the noise in the Web crawl?
    - 43k unlabeled Web pages → "other" class

- Favorable setting:  plenty of labeled data
- 45-way 1-nearest neighbor classification, 10 random 70/30 train/test splits

Avg accuracy  =  99.95%
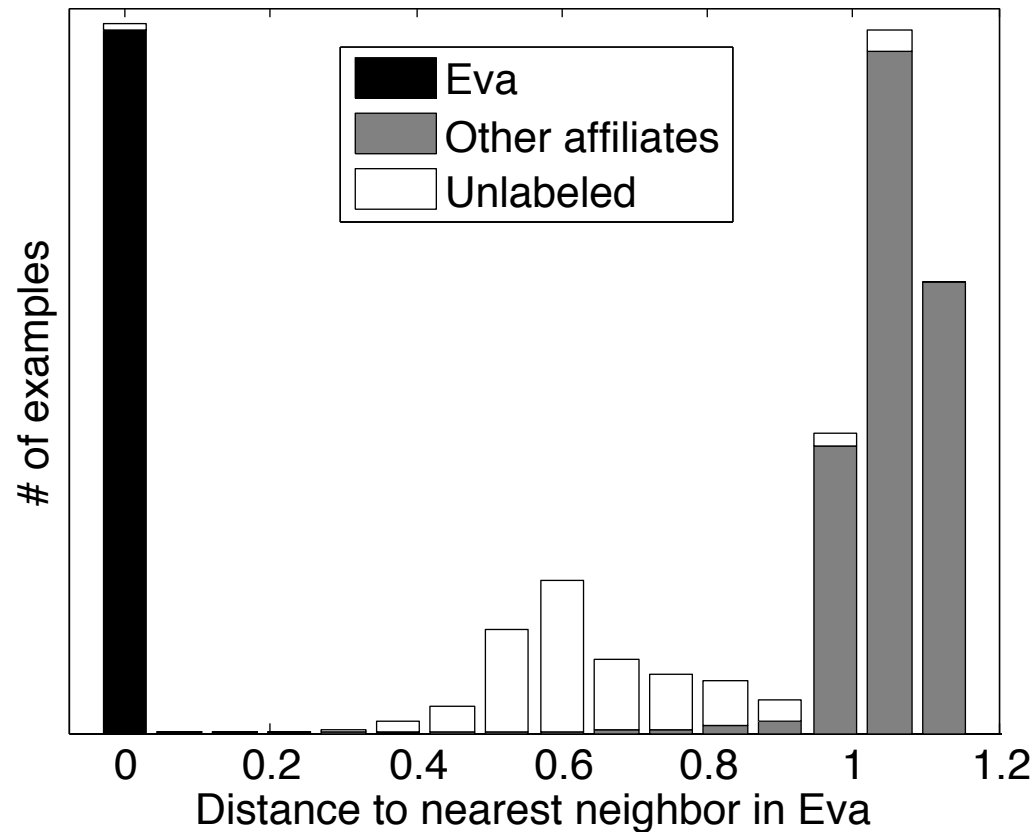
# Proof-of-concept experiment

- **Question**:  are these HTML features enough to distinguish affiliate programs:
  - From one another?
  - From the noise in the Web crawl?
    - 43k unlabeled Web pages → "other" class

- Favorable setting:  plenty of labeled data
- 45-way 1-nearest neighbor classification, 10 random 70/30 train/test splits

Avg accuracy  =  99.95%

How so good?!

# HTML distances are highly predictive

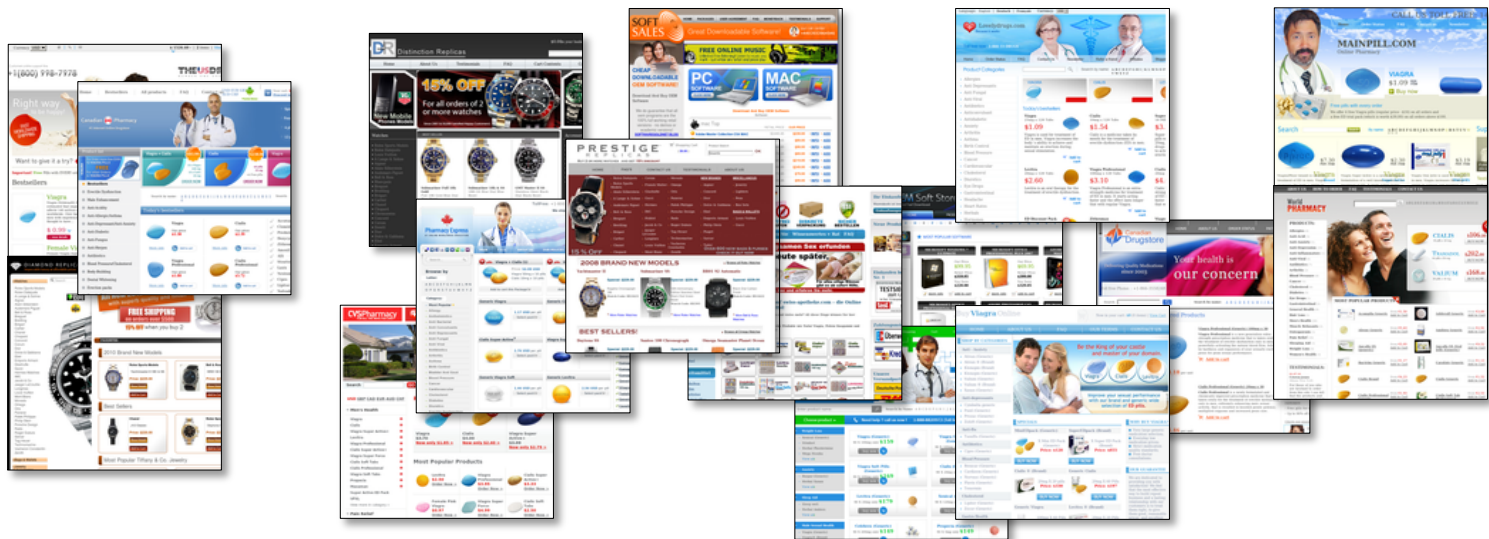Distances from *every* point to nearest neighbor in EvaPharmacy

- Experts must label **some** storefronts, but how many?
- Learning from scratch:  only small initial seed of labeled storefronts

- Experts must label **some** storefronts, but how many?
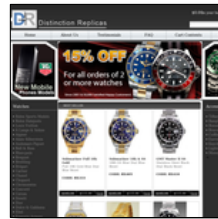- Learning from scratch:  only small initial seed of labeled storefronts

# Mimicking an operational deployment

- Experts must label **some** storefronts, but how many?
- Learning from scratch:  only small initial seed of labeled storefronts
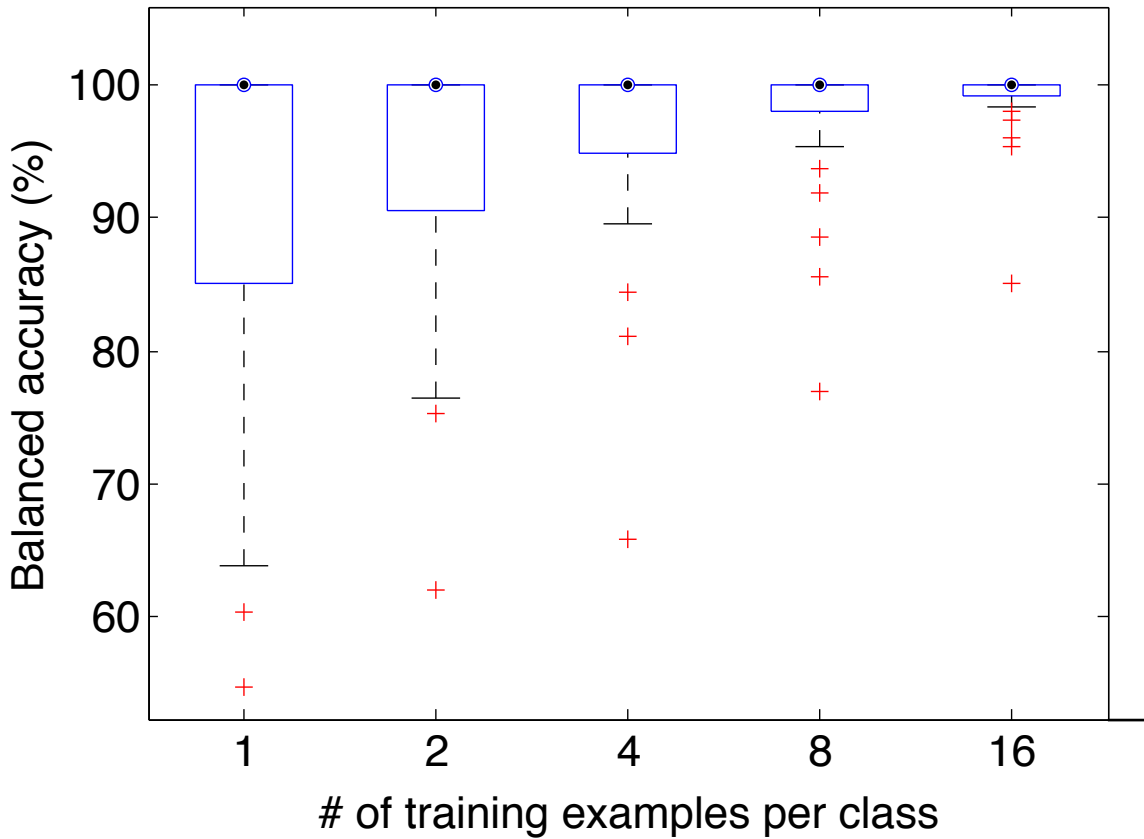
# Classification in operational setting

# One-shot learning

### Singly labeled storefront

### Correctly classified storefronts

**33drugs**



**RX-Promotions**

# Further results

- Found & labeled 3,785 additional storefront pages



- Clustered unlabeled Web pages to identify new affiliate programs

# Conclusion

- Automated system for identifying affiliate programs behind illegal online storefronts
- Simple model is highly accurate
  - Templatized storefronts, many near-duplicates
  - Affiliate programs' efforts to operate at scale make automated defense possible
- Big win for security practitioners
  - Modest labeling effort is enough to bootstrap the system
- Security a Big Data field; ML an invaluable tool for large-scale Web crawls

# Thank you!

## Questions?

UCSD CSE
Computer Science and Engineering