

# Weakly Supervised Extraction of Computer Security Events from Twitter

**Alan Ritter**, Evan Wright, William Casey, Tom Mitchell



THE OHIO STATE  
UNIVERSITY

Carnegie Mellon

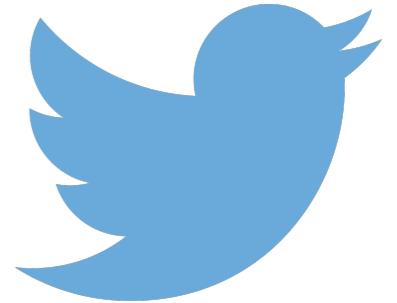
# Natural Language Processing on Newswire

- Lots of previous work:
  - MUC & ACE competitions
  - Timebank
  - Penn Treebank
  - Topic Detection and Tracking
  - Etc...





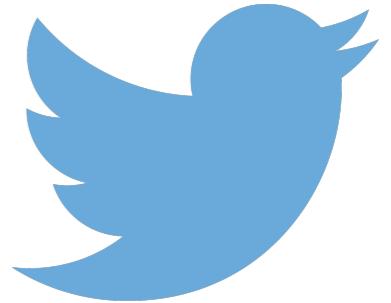
# Microblogs vs. Newswire



- **Opportunities:**
  - Fresher Information
  - Broader Coverage
- **Challenges:**
  - Irrelevant and redundant messages
  - Unreliable information



# Microblogs vs. Newswire



- **Opportunities:**

- Fresher Information
- Broader Coverage

- **Challenges:**

- Irrelevant and redundant messages
- Unreliable information



Information Overload

A red speech bubble shape points from the bottom left towards the text "Information Overload".

1 May 2011

3:58 PM EST



**Sohaib Athar**

@ReallyVirtual



+ Follow

Helicopter hovering above Abbottabad at 1AM (is a rare event).

10:24 PM EST



**Keith Urbahn**

@keithurbahn



+ Follow

So I'm told by a reputable person they have killed Osama Bin Laden. Hot damn.

11:35 PM EST

Presidential  
Address  
Broadcast



# Long tail of Smaller Events



**Luc Rossini**

@LucRossini



 Follow

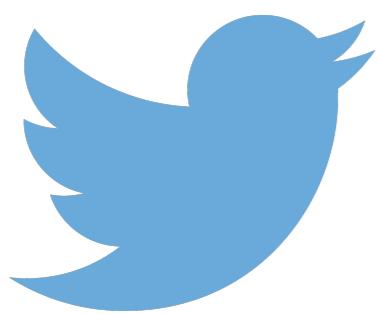
Spamhaus is currently under a DDoS attack against our website which we are working on mitigating. Our DNSBLs are not affected.

[#spamhaus](#)

# Information Overload



# Information Overload



Information  
Extraction



7/4/2014

**Phishing Attack**

Victim: Bitcoins Reserve

4/25/2015

**Account Hijacking**

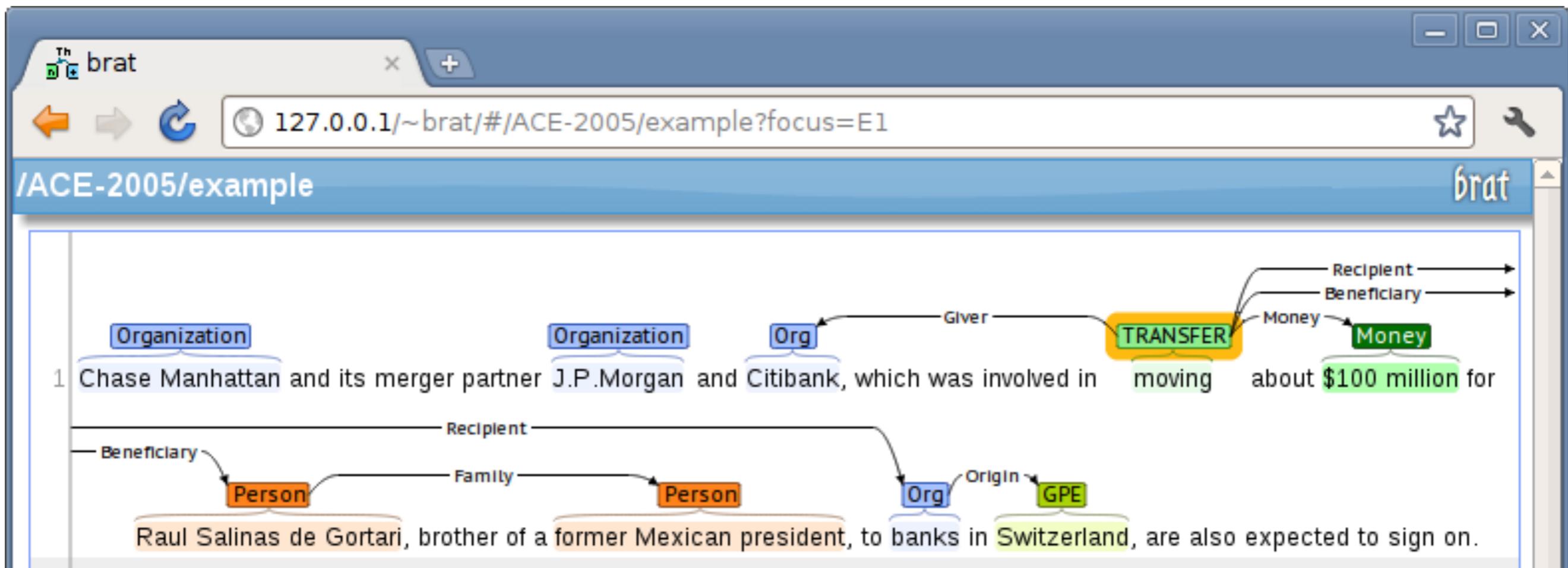
Victim: Tesla

5/16/2015

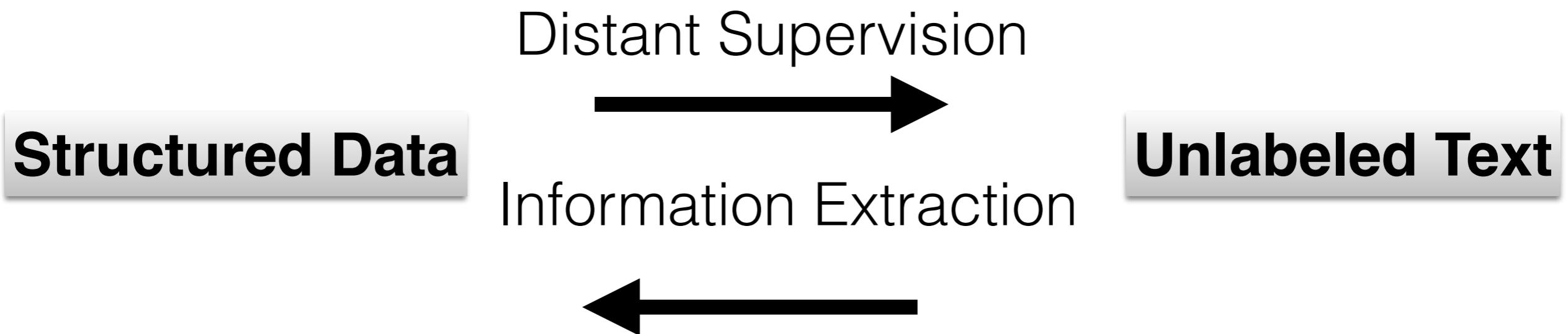
**DDOS**

Victim: PSN

# Traditional Information Extraction



# Weakly Supervised Learning



# Relation Extraction

e.g. [Mintz et. al. 2009]

 Freebase

Person	Birth Location
Barack Obama	Honolulu
Mitt Romney	Detroit
Albert Einstein	Ulm
Nikola Tesla	Smiljan
...	...

**(Albert Einstein, Ulm)**

**(Mitt Romney, Detroit)**

**(Barack Obama, Honolulu)**

“Barack Obama was born on August 4, 1961 at ...  
in the city of **Honolulu** ...”

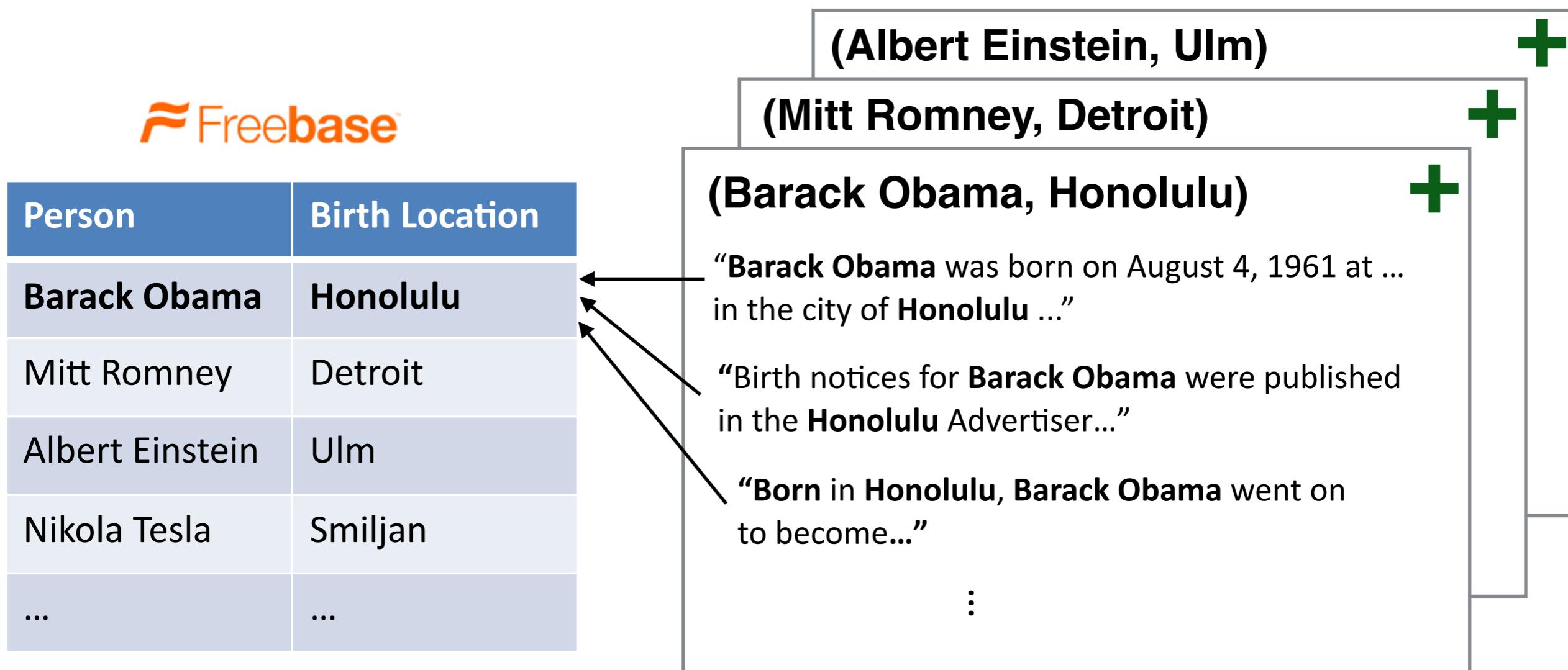
“Birth notices for **Barack Obama** were published  
in the **Honolulu Advertiser**...”

“Born in **Honolulu**, Barack Obama went on  
to become...”

:

# Relation Extraction

e.g. [Mintz et. al. 2009]



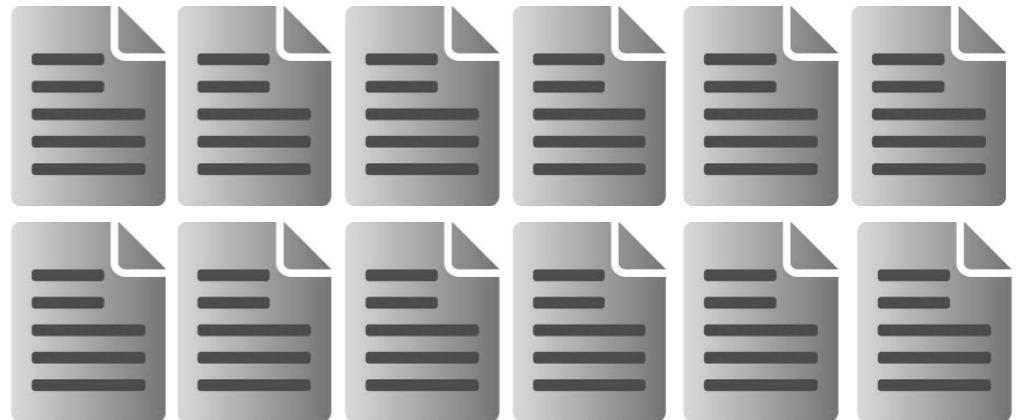
# Distant Supervision: Events vs. Relations

## Relations:

HEADQUARTEREDIN(Microsoft, Redmond)



BORNIN(Barack Obama, Honolulu)

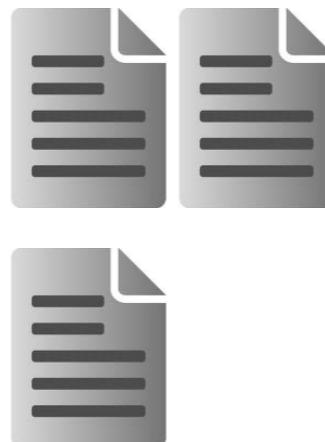


## Events:

TERRORISTATTACK(Quetta, Feb. 16 2011)



PRODUCTRELEASE(Galaxy S 6, Mar. 22 2015)



# Distant Supervision: Events vs. Relations

## Relations:

HEADQUARTEREDIN(Microsoft, Redmond)

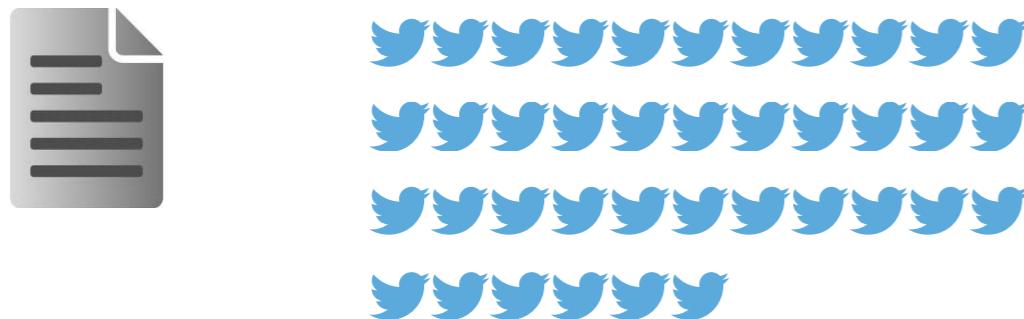


BORNIN(Barack Obama, Honolulu)

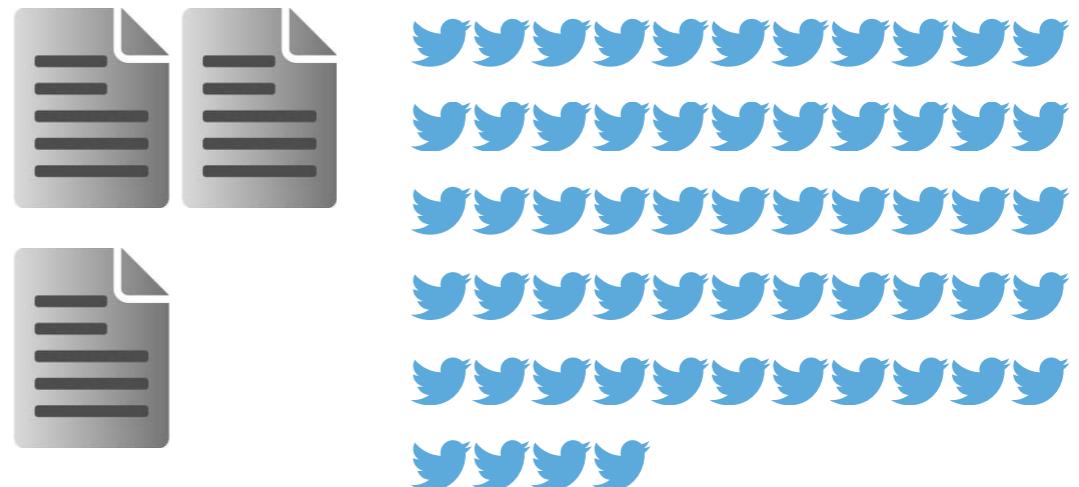


## Events:

TERRORISTATTACK(Quetta, Feb. 16 2011)



PRODUCTRELEASE(Galaxy S 6, Mar. 22 2015)



# Cybersecurity Events

**Denial of Service**

Account Hijacking

Data Breach



Luc Rossini  
@LucRossini



 Follow

Spamhaus is currently under a DDoS attack against our website which we are working on mitigating. Our DNSBLs are not affected.  
[#spamhaus](#)

# Cybersecurity Events

Denial of Service

Account Hijacking

Data Breach



**AP** The Associated Press

@AP



Following

Breaking: Two Explosions in the White House and Barack Obama is injured

# Cybersecurity Events

Denial of Service

Account Hijacking

Data Breach



Target  
@Target



Follow

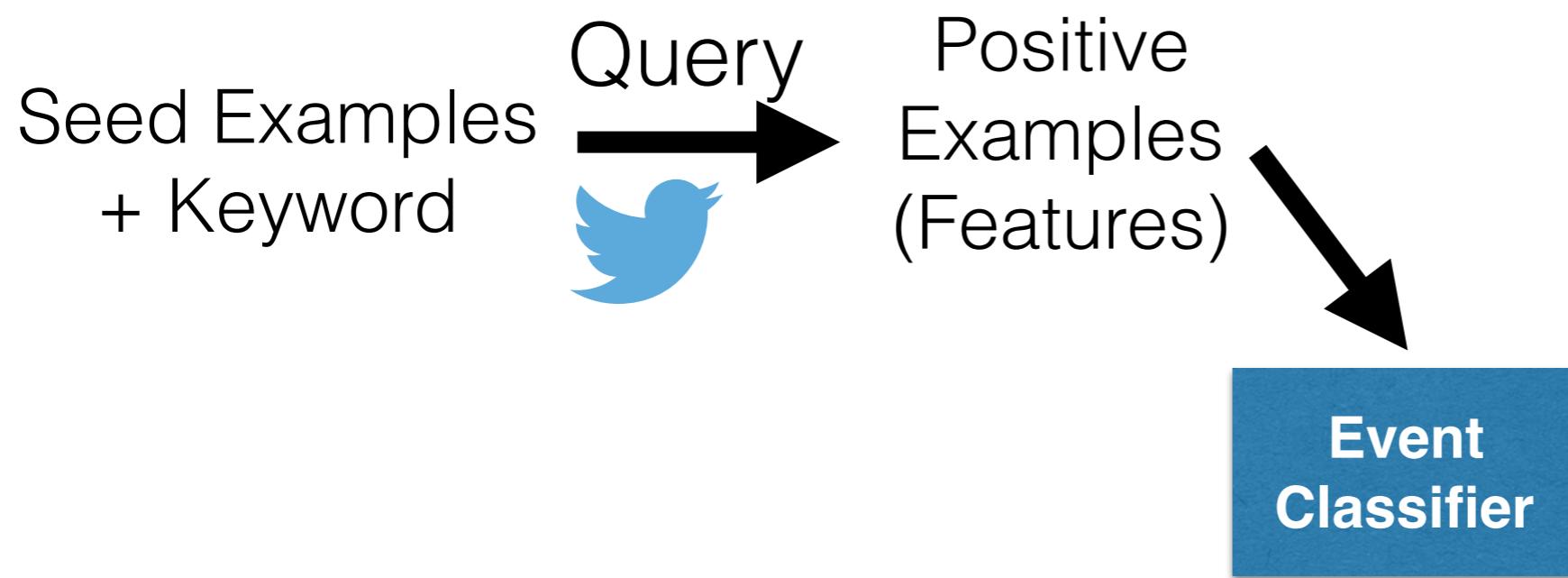
Target Confirms Unauthorized Access to Payment Card Data in U.S. Stores: Issue identified and resolved [tgt.biz/1kXzK0m](http://tgt.biz/1kXzK0m)

# System Overview

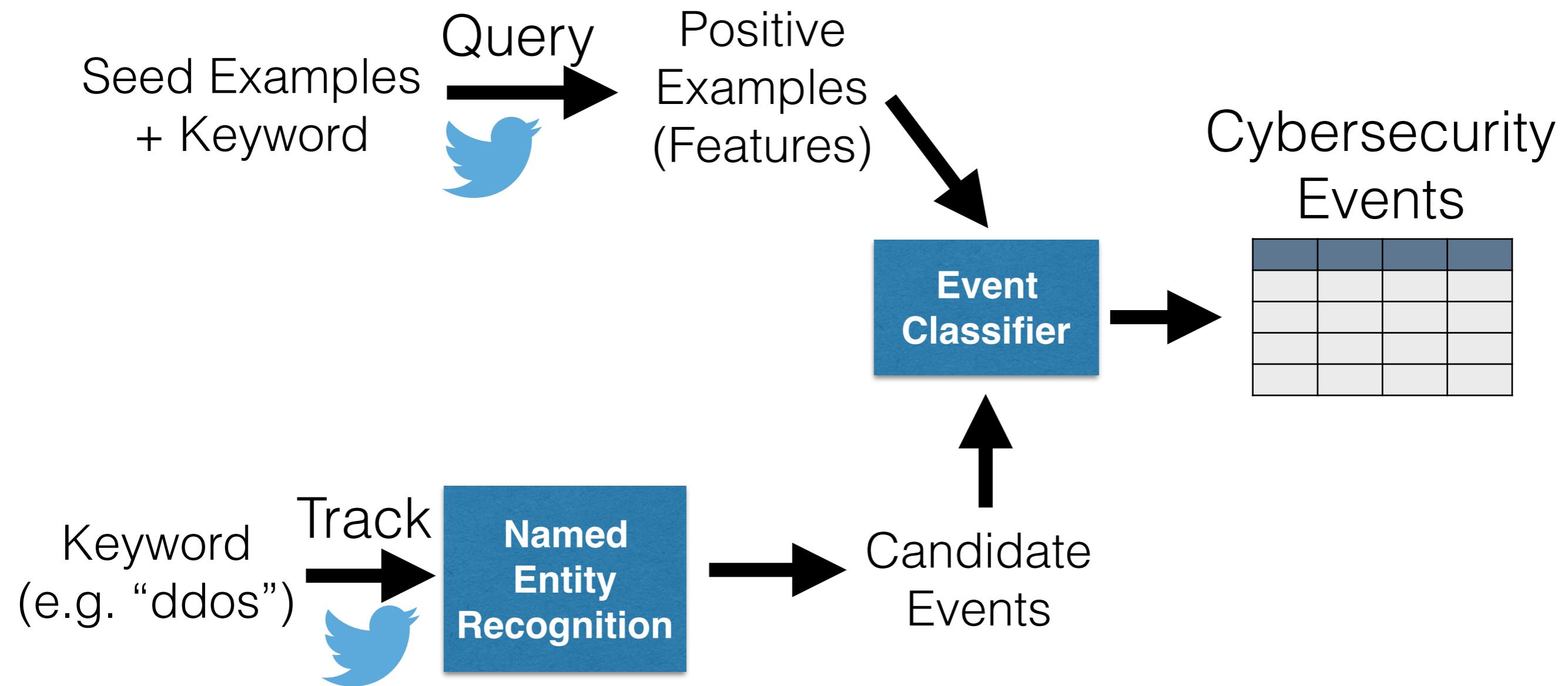
Seed Examples  
+ Keyword

**(Associated Press, 4/23/2013)**

# System Overview



# System Overview



# Account Hijacking Seed Events

Victim	Date
associated press	2013/04/23
reuters	2012/08/05
us marines	2013/09/03
sarah palin	2008/09/18
mitt romney	2012/06/05
cnn	2014/01/23
justin bieber	2012/03/27
mutunga	2013/09/27
yes scotland	2013/08/20
zuckerberg	2013/08/18

associated press hacked since:2013-4-23 until:2013-4-24

Search

Tip: use operators for advanced search.



**CyG US** @CyG\_US · Apr 23

**Associated Press** Twitter account **hacked**, false report of explosions, Obama injury at White House sent out [bit.ly/11LLrhk](http://bit.ly/11LLrhk)

Expand

Reply Retweet Favorite More



**Ambling Freest** @amblingfreest · Apr 23

**Associated Press'** Twitter **hacked**, fake tweet sends stocks plunging - Los Angeles Times: **Associated Press'** Twit... [bit.ly/12c1jwf](http://bit.ly/12c1jwf)

Expand

Reply Retweet Favorite More



**Al Stefanelli** @Stefanelli · Apr 23

**Associated Press** Twitter Accounts **Hacked**, Fake White House Bomb Report Posted [wp.me/p3oa9W-AD](http://wp.me/p3oa9W-AD)

Expand

Reply Retweet Favorite More



**DudeYouCrazy** @DudeYouCrazy · Apr 23

Wait, so was the **Associated Press** also **hacked** when its final football poll had Notre Dame ranked fourth?

Expand

Reply Retweet Favorite More



**BiGxGh.Com** @bigxghdotcom · Apr 23

#MissedThisOnOMGGhana: **Associated Press'** Twitter Account **Hacked!** [bit.ly/11LBaS9](http://bit.ly/11LBaS9)

Expand

Reply Retweet Favorite More

# Features

- Binary features extracted from “bag-of-tweets”
- Context words surrounding the victim of the attack
- Context words surrounding keyword

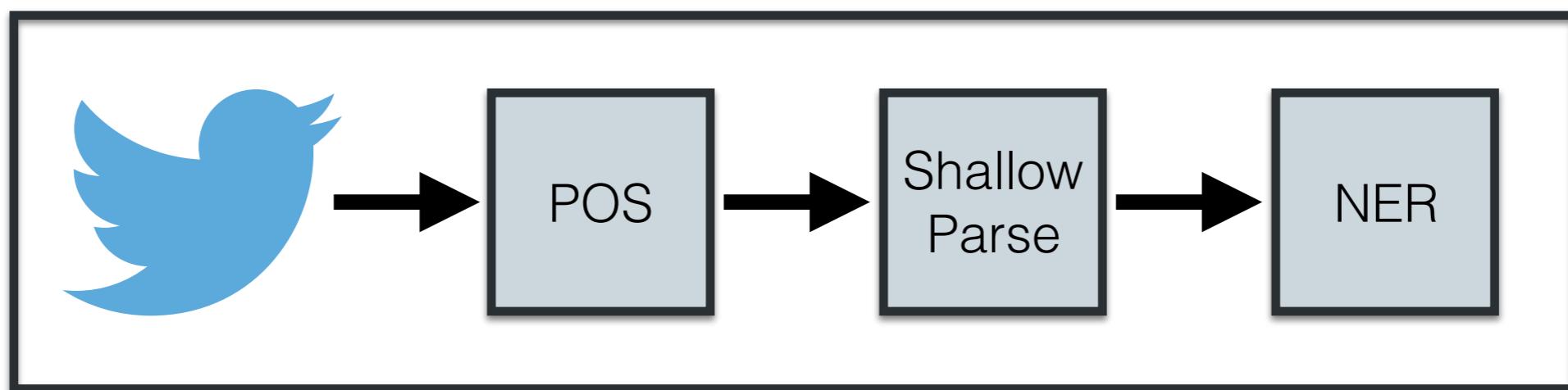
# Features

- Binary features extracted from “bag-of-tweets”
- Context words surrounding the victim of the attack
- Context words surrounding keyword

Sample Feature	Event Category
<b>security breach</b>	Data Breach
<b>X admits</b>	Data Breach
<b>X data breach affecting</b>	Data Breach
<b>DT ddos attack</b>	DDoS
<b>NNP NNP hit IN ddos</b>	DDoS
<b>X getting ddos'd</b>	DDoS
<b>PRP hacked X POS account</b>	Account Hijacking
<b>POS email was hacked</b>	Account Hijacking
<b>X POS NNP account</b>	Account Hijacking

# Gathering Unlabeled / Test Data

- Track a keyword associated with the event (e.g. **hacked, DDOS, breach**)
- Extract named entities (using Twitter-tuned NLP tools)  
**[Ritter, et. al. EMNLP 2011a]**



# Not account hijacking...



**Gavin Caney** @GavinCaney · Mar 28

PENALTY to Lynn as @JakeJonesx11 is hacked down by Hutchinson. 1-0 Blyth but for how long? #EDPsport



# Not a Data Breach...



**Blue Planet Society**  
@Seasaver



Follow

Two 40-tonne humpback [#whales](#) breach in unison in [#Hawaii](#) [dailym.ai/1FeRP75](http://dailym.ai/1FeRP75) via [@MailOnline](#)



# Learning from Unlabeled Data and Positive Seeds

Augment conditional likelihood with label regularization:

- [Mann and McCallum]

$$O(\theta) = \underbrace{\sum_i^N \log p_\theta(y_i|x_i)}_{\text{Log Likelihood}} - \underbrace{\lambda^U D(\tilde{p} || \hat{p}_\theta^{\text{unlabeled}})}_{\text{Label regularization}}$$

# KL Divergence

$$D(\tilde{p} \parallel \hat{p}_\theta) = \tilde{p} \log \frac{\tilde{p}}{\hat{p}_\theta} + (1 - \tilde{p}) \log \frac{1 - \tilde{p}}{1 - \hat{p}_\theta}$$

# KL Divergence

$$D(\tilde{p} \parallel \hat{p}_\theta) = \tilde{p} \log \frac{\tilde{p}}{\hat{p}_\theta} + (1 - \tilde{p}) \log \frac{1 - \tilde{p}}{1 - \hat{p}_\theta}$$



User-provided target expectation

# KL Divergence

Empirical expectation on unlabeled examples

$$D(\tilde{p} \parallel \hat{p}_\theta) = \tilde{p} \log \frac{\tilde{p}}{\hat{p}_\theta} + (1 - \tilde{p}) \log \frac{1 - \tilde{p}}{1 - \hat{p}_\theta}$$

User-provided target expectation

# KL-Div Gradient

$$\frac{\partial}{\partial \theta_k} D(\tilde{p} || \hat{p}_\theta) =$$

$$\frac{1}{N} \left( \frac{1 - \tilde{p}}{1 - \hat{p}_\theta} - \frac{\tilde{p}}{\hat{p}_\theta} \right) \sum_{i=1}^N p_\theta(y_i = 1 | x_i) (1 - p_\theta(y_i = 1 | x_i)) x_{i,k}$$

# KL-Div Gradient

$$\frac{\partial}{\partial \theta_k} D(\tilde{p} || \hat{p}_\theta) =$$

$$\frac{1}{N} \left( \frac{1 - \tilde{p}}{1 - \hat{p}_\theta} - \frac{\tilde{p}}{\hat{p}_\theta} \right) \sum_{i=1}^N p_\theta(y_i = 1 | x_i) (1 - p_\theta(y_i = 1 | x_i)) x_{i,k}$$



No Change if  $\tilde{p} = \hat{p}_\theta$

Otherwise push weights up or down

# KL-Div Gradient

$$\frac{\partial}{\partial \theta_k} D(\tilde{p} || \hat{p}_\theta) =$$

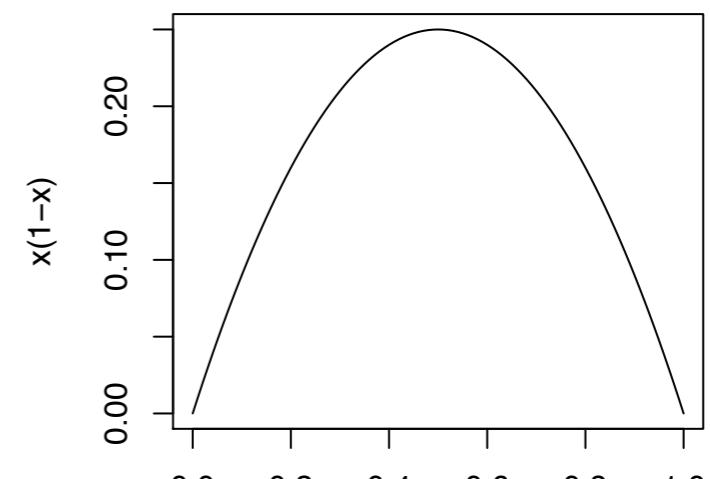
Give more weight to  
uncertain cases

$$\frac{1}{N} \left( \frac{1 - \tilde{p}}{1 - \hat{p}_\theta} - \frac{\tilde{p}}{\hat{p}_\theta} \right) \sum_{i=1}^N p_\theta(y_i = 1 | x_i) (1 - p_\theta(y_i = 1 | x_i)) x_{i,k}$$



No Change if  $\tilde{p} = \hat{p}_\theta$

Otherwise push weights up or down

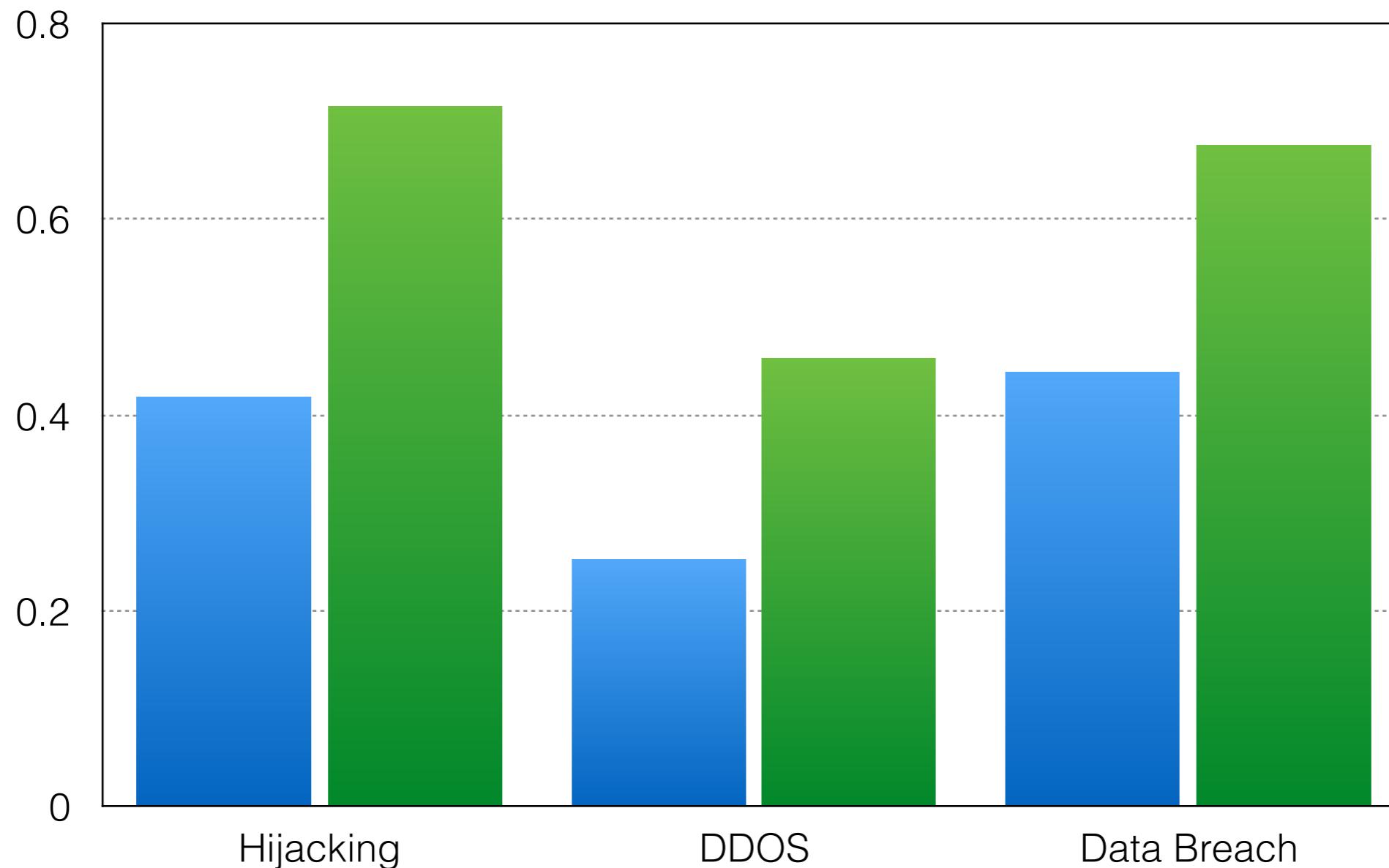


# Baselines

- Logistic Regression with Heuristic Negatives
  - used in lots of previous work
- One-class support vector machines [Schölkopf et. al.]
  - ignore unlabeled events
- Semi-supervised EM

# AUC

■ Logistic Regression ■ Expectation Regularization

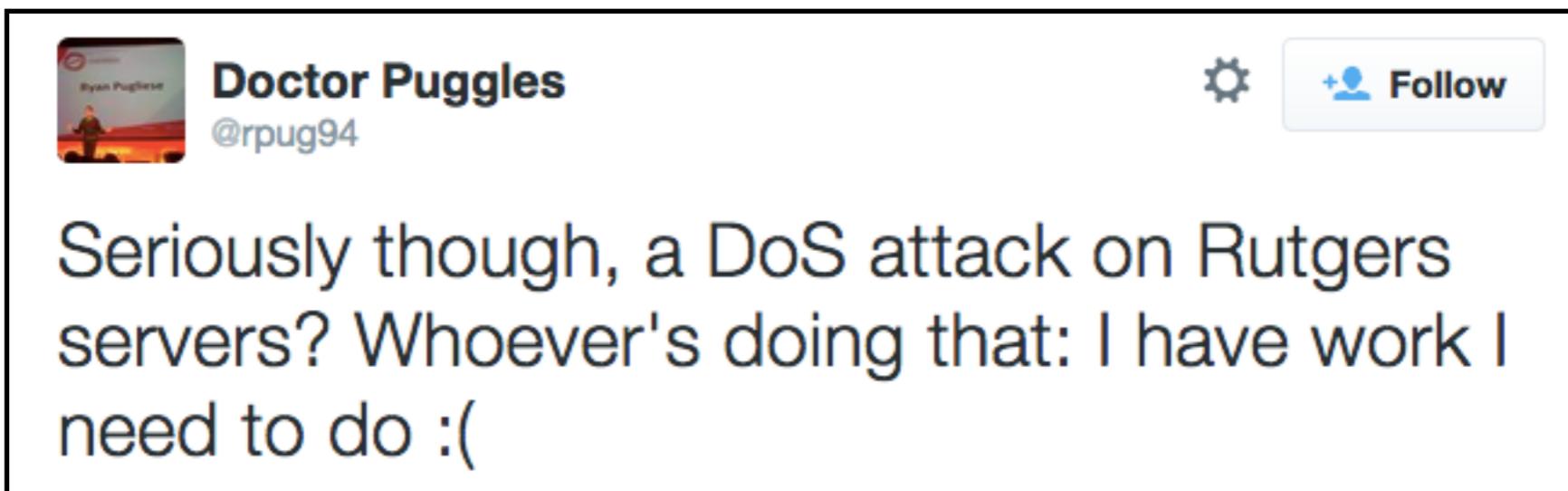


# DDOS

Mar-28-2015

china	Is China responsible for Massive DDoS Attack against GitHub? <a href="http://t.co/JDGYInB0Lv">http://t.co/JDGYInB0Lv</a>	265
rutgers	RT @MikeBrady35: That moment when your so frustrated with the bus system that you successfully ddos the absolute shit out of Rutgers servers	10
baidu	RT @jschauma: Some more info on the @GitHub DDoS, apparently by way of Baidu: <a href="http://t.co/J9HF2NijiY">http://t.co/J9HF2NijiY</a>	13
github	Is China responsible for Massive DDoS Attack against GitHub? <a href="http://t.co/JDGYInB0Lv">http://t.co/JDGYInB0Lv</a>	236
ruwireless	RT @francessssss: this kid on reddit saying he's responsible for the DDOS attack on RUwireless. bruh, I'm just trying to do my work #iwillcu...	3
github ddos	RT This GitHub DDoS is serious. <a href="https://t.co/HRmqMiOlba">https://t.co/HRmqMiOlba</a> <a href="https://t.co/PUzHbFrYWv">https://t.co/PUzHbFrYWv</a> <a href="http://t.co/59Z2O27qcz">http://t.co/59Z2O27qcz</a>	6

# DDOS



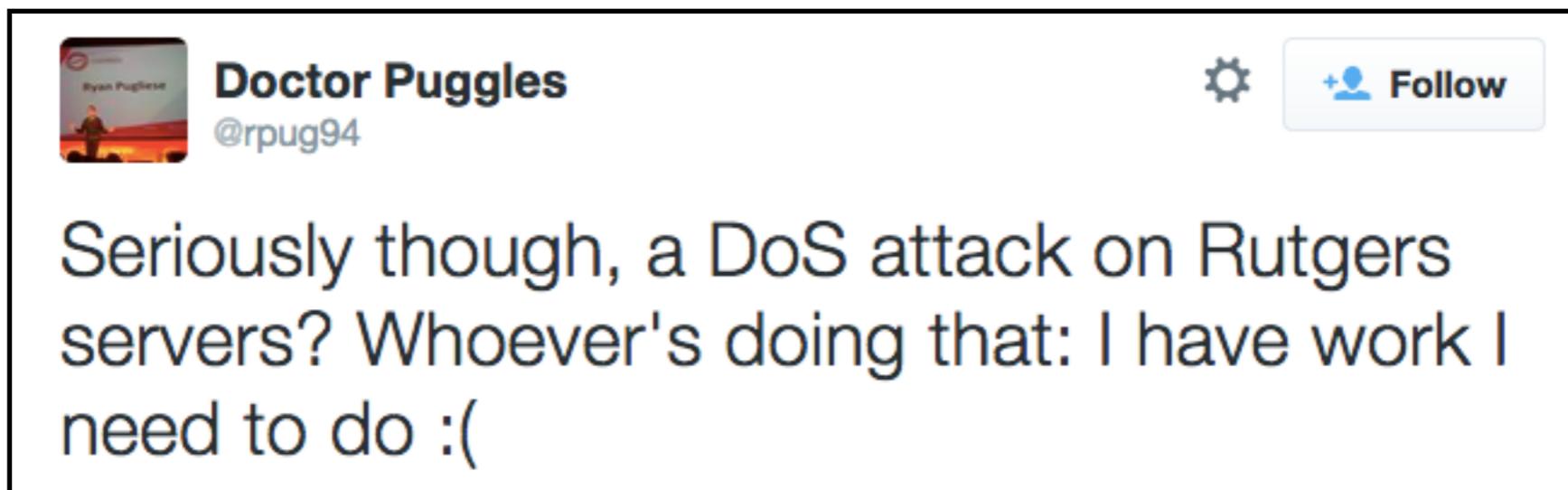
Doctor Puggles  
@rpug94

Seriously though, a DoS attack on Rutgers servers? Whoever's doing that: I have work I need to do :(

Mar-28-2015

- china Is China responsible for Massive DDoS Attack against GitHub? <http://t.co/JDGYInB0Lv> 265
- rutgers RT @MikeBrady35: That moment when your so frustrated with the bus system that you successfully ddos the absolute shit out of Rutgers servers 10
- baidu RT @jschauma: Some more info on the @GitHub DDoS, apparently by way of Baidu: <http://t.co/J9HF2NijiY> 13
- github Is China responsible for Massive DDoS Attack against GitHub? <http://t.co/JDGYInB0Lv> 236
- ruwireless RT @francessssss: this kid on reddit saying he's responsible for the DDOS attack on RUwireless. bruh, I'm just trying to do my work #iwillcu... 3
- github ddos RT This GitHub DDoS is serious. <https://t.co/HRmqMiOlba> <https://t.co/PUzHbFrYWv> <http://t.co/59Z2O27qcz> 6

# DDOS



Doctor Puggles  
@rpug94

Seriously though, a DoS attack on Rutgers servers? Whoever's doing that: I have work I need to do :(



GitHub  
@github

The attack has ramped up again, and we're evolving our mitigation strategies to match.

- Mar-28-2015
- china Is China responsible for Massive DDoS Attacks? 265
  - rutgers RT @MikeBrady35: That's what I'm talking about! 10
  - baidu RT @intrauma: Some more info 13
  - github Is China responsible for Massive DDoS Attack against GitHub? http://t.co/JDGYInB0Lv 236
  - ruwireless RT @francessssss: this kid on reddit saying he's responsible for the DDOS attack on RUwireless. bruh, I'm just trying to do my work #iwillcu... 3
  - github ddos RT This GitHub DDoS is serious. https://t.co/HRmqMiOlba https://t.co/PUzHbFrYWv http://t.co/59Z2O27qcz 6

# Account Hijacking

Mar-29-2015

british airways	British Airways says some frequent flyer accounts hacked <a href="http://t.co/rsh1qOZ7WH">http://t.co/rsh1qOZ7WH</a>	210
nigeria	"Nigeria: Election Website Hacked" #newspapers #feedly <a href="http://t.co/JngBLefyn">http://t.co/JngBLefyn</a>	37
slackster	Group Chatting Application Slacker Has Been Hacked - <a href="http://t.co/whBy5QdmCR">&gt;</a> And I always thought slackster was for streaming music. #infosec	9
uk	RT @PotCoinFan: Hacking is cool if we do it just not you. - UK Government Authorized GHCQ to Hack Any Device <a href="https://t.co/D0AbsHG9nB">https://t.co/D0AbsHG9nB</a> via @H...	11
inec	RT @OmoOduaRere: INEC's Website get hacked by Nigeria Cyber hacking group on Election day #Nigeriadecides <a href="http://t.co/nn8mOhPulp">http://t.co/nn8mOhPulp</a>	61
taylor swift	Taylor Swift Hacked, 6 Nudes for Sale? <a href="http://t.co/k6HhFe9nEa">http://t.co/k6HhFe9nEa</a>	13
inec website	INEC Website has been Hacked <a href="http://t.co/346VTvfvcd">http://t.co/346VTvfvcd</a>	4
darren	The people who think Darren was hacked are the same people who think he doesn't lie to his fans, wrote This Time for Lea, and likes vag.	12

# Account Hijacking

The Guardian   
@guardian

Follow

British Airways frequent flyer accounts hacked [trib.al/82aR1y4](http://trib.al/82aR1y4)

Mar-29-2015

British Airways says some frequent flyer accounts hacked <http://t.co/rsh1qOZ7WH> 210

nigeria "Nigeria: Election Website Hacked" #newspapers #feedly <http://t.co/JngBLefyn> 37

slackerr Group Chatting Application Slacker Has Been Hacked - [&gt;](http://t.co/whBy5QdmCR) And I always thought slacker was for streaming music. #infosec 9

uk RT @PotCoinFan: Hacking is cool if we do it just not you. - UK Government Authorized GHCQ to Hack Any Device <https://t.co/D0AbsHG9nB> via @H... 11

inec RT @OmoOduaRere: INEC's Website get hacked by Nigeria Cyber hacking group on Election day #Nigeriadecides <http://t.co/nn8mOhPulp> 61

taylor swift Taylor Swift Hacked, 6 Nudes for Sale? <http://t.co/k6HhFe9nEa> 13

inec website INEC Website has been Hacked <http://t.co/346VTfVcD> 4

darren The people who think Darren was hacked are the same people who think he doesn't lie to his fans, wrote This Time for Lea, and likes vag. 12

# Account Hijacking

Mar-29-2015

british airways British Airways says some frequent flyer accounts hacked <http://t.co/rsh1aOZ7WH> 210

nigeria "Nigeria: Election Website" 37

slackerr Group Chatting Application streaming music. #i 9

uk RT @PotCoinFan: Hacking <https://t.co/D0AbsHG9nB> 11

inec RT @OmoOduaRere: INEC <http://t.co/nm8mOhPulp> 61

taylor swift Taylor Swift Hacked, 6 Nu 13

inec website INEC Website has been H 4

darren The people who think Darren likes vag. 12

The Guardian  Follow

British Airways frequent flyer accounts hacked <trib.al/82aR1y4>

Yahoo Labs follows  Breaking News Tech @breakingbytes • Mar 28

Website of Nigerian electoral commission hacked on day of national election, not expected to affect poll - @techcabal <bit.ly/1HaJelP>

TechCabal

[!] Struck By Nigerian Cyber Army | Team NCA [!]



Sorry xD Your Site has been STAMPED by Team Nigerian Cyber Army  
FEEL SOME SHAME ADMIN!!  
Security is just an illusion

Remember US :D GREETINGS OF PEACE TO CITIZEN OF NIGERIA FROM TEAM NCA

# Data Breach

Mar-28-2015

banks crabb badger davis  
elphick breach

kreditech

uber

premra

provider reports security  
breach

google

israel

amazon

twitch

RT @L...

Kreditech Su  
http://t.co/v...

RT @EnverP...  
reported recently?

RT @DMBisson: Premera faces class action suit over data #breach http://t.co/oiKLf01Hkl via @modrnhealthcr  
#security #breach

Jack Attack

Google Ads  
http://t.co/H...

Fav if you w  
take back our land by the sword

Twitch, the video game streaming site Amazon bought for \$970 million, has been hacked (AMZN)  
http://t.co/1YMfZPT2Cx #business

New hacking group DDoS attacks Amazon's Twitch, US department websites | http://t.co/jkzZV2JH28



**briankrebs**

@briankrebs



Follow

Kreditech, a lending firm for the "unbanked,"  
acknowledges breach after customer data  
posted online



**The Drum**

@TheDrum



Follow

• **@Uber** admits #data breach exposed  
personal details of 50,000 drivers

# Contributions

- Microblogs are a valuable source of info for cybersecurity events
- Novel approach to event extraction with minimal supervision
  - Weakly supervised learning as a problem with positive and unlabeled data.
  - New approach based on expectation regularization
- Demo: <http://securitytweets.org/>

# Contributions

- Microblogs are a valuable source of info for cybersecurity events
- Novel approach to event extraction with minimal supervision
  - Weakly supervised learning as a problem with positive and unlabeled data.
  - New approach based on expectation regularization
- Demo: <http://securitytweets.org/>

**Thanks!**