



# Breaking Big Data

Evading Analysis of the Metadata of Your Life

Dave Venable  
Masergy Communications  
@davevenable



hooli

# Big Data

Data sets that are so large or complex that traditional data processing applications are inadequate to deal with them.

Source: Wikipedia



# What's being collected?



An aerial photograph of a large industrial complex, possibly a refinery or chemical plant. The facility consists of several large, rectangular buildings with white roofs and grey walls. There are numerous storage tanks, some large and cylindrical, and others smaller. The surrounding area is mostly flat, with some greenery in the distance. A semi-transparent black banner with white text is overlaid across the middle of the image.

...and figure out what to do with it later

# Will Big Data Cure Cancer?

Big data technologies will play a key role in diversifying Saudi Arabia's economy away from a huge dependence on oil revenues.

## How the Chicago Cubs Used Big Data to Win the World Series

Data + Algorithms = **Big Data**

\$3.1T

Source: IBM

Bad

Data + Algorithms = Big Data



# WEAPONS OF MATH DESTRUCTION



HOW BIG DATA INCREASES INEQUALITY  
AND THREATENS DEMOCRACY

CATHY O'NEIL

Bad Biased

Data + Algorithms = Big Data

*Bad* *Biased* *Delusions*  
Data + Algorithms = Big ~~Data~~





# Encryption Works

“Properly implemented strong crypto systems  
are one of the few things that you can rely on.”

- Edward Snowden

# But so what?

Email Content

Messaging

Phone Calls

Packet Payloads

HTTPS Content

Sender

Recipient

Participants

Length

Source

Destination



# Alice and Bob



British Airways → Alice



Alice is in Seoul



“Seoul is so fun!”



Bob → Cindy



Cindy → Bob



Bob → Hotel



Cindy → Hotel



Bob - £500.00 Hotel



Hotel → Restaurant



Bob - £230.00 Restaurant



Restaurant → Hotel



Hotel → Bob's



Hotel → Cindy's



Bob → Cindy



Cindy → Bob



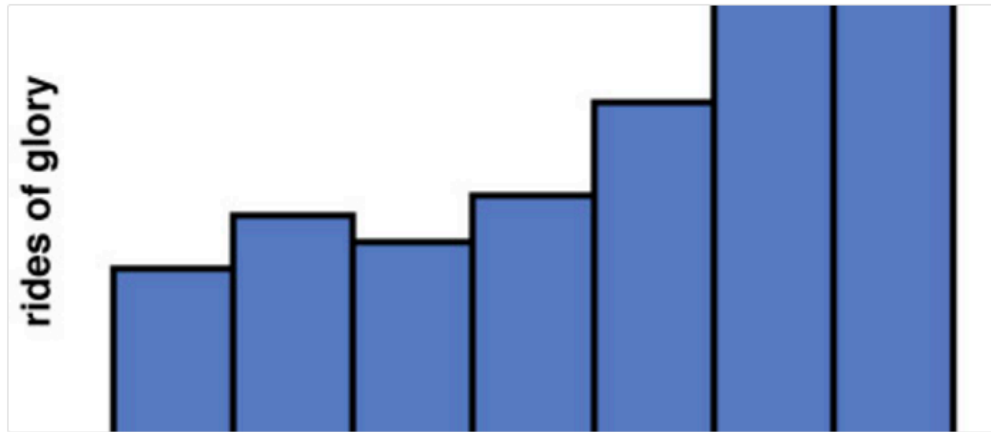
Alice is in London

#UBERDATA

# RIDES OF GLORY

MARCH 26, 2012

POSTED BY VOYTEK



# Encryption Works

“Properly implemented strong crypto systems  
are one of the few things that you can rely on.”

- Edward Snowden

# Encryption Works

“Properly implemented strong crypto systems are one of the few things that you can rely on.  
Unfortunately, endpoint security is so terrifically weak that NSA can frequently find ways around it.”

- Edward Snowden

# Endpoints Suck

# Ancient Adage

Two types of endpoints:

- Those who have been compromised
- Those who will be

# ~~Ancient~~ Adage

Two types of endpoints:

- Those who have been compromised
- Those who ~~will be~~ don't know it yet



~~New~~  
~~Old~~  
~~Ancient~~ Adage  
^

~~One~~  
~~Two~~ types of endpoints:  
^

- ~~Those who have been compromised~~
- ~~Those who will be~~ ~~don't know it yet~~

~~New~~  
~~Old~~  
~~Ancient~~ Adage

~~One~~  
~~Two~~ types of endpoints:

- ~~Those who have been compromised~~
- ~~Those who will be~~ ~~don't know it yet~~

**Everything is compromised**

...and what isn't has permission.

# Bob

9:00 pm / 21:00



Home → Pub



Bob is driving to the pub (too fast)



Check in at pub

2:00 am / 02:00



Bob - £100.00 @ pub



Pub → Home



Bob is driving home (too slow)

6:30 am / 06:30



Driving to work

(too carelessly)



# Your life has patterns

...that produce metadata



B B C



Your habits



#irc





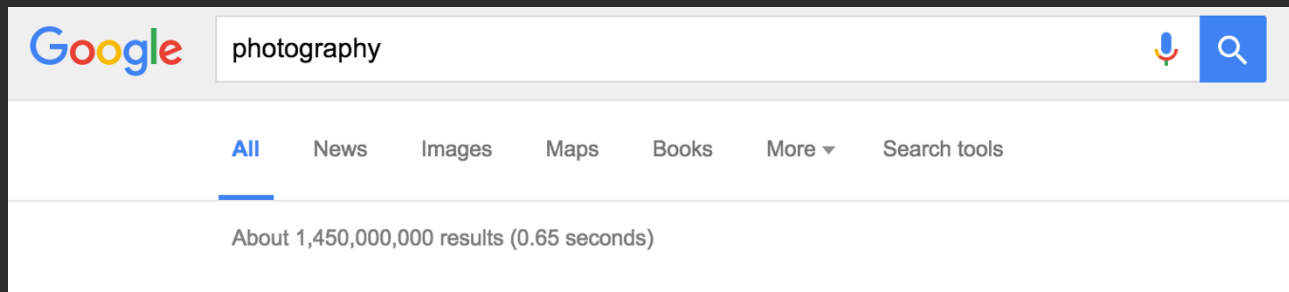
# Your locations



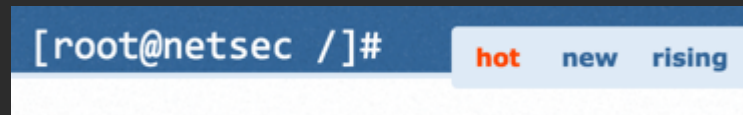
  
**black hat**  
EUROPE 2016







# Your interests





Bob ↔ Cindy



# Your friends



Same place, same time



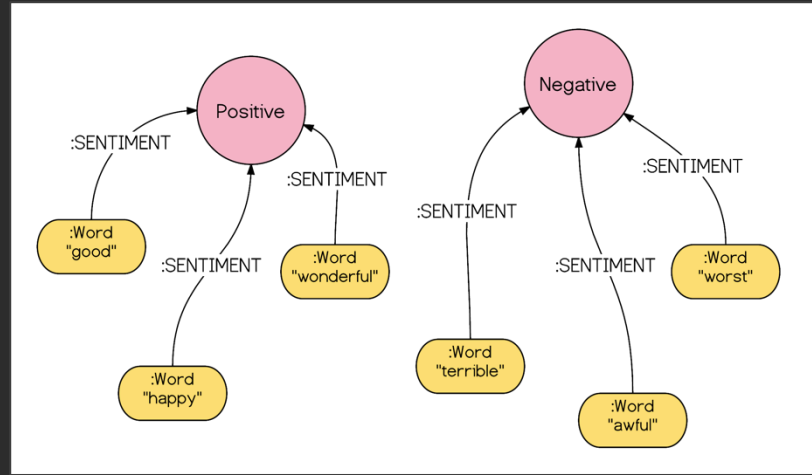


[Candidate 1] sucks



Loving this coffee shop

# Your feelings





### Anti Face

This face is unrecognizable to several state-of-art face detection algorithms.



### Face

Once computer vision programs detect a face, they can extract data about your emotions, age, and identity.

[See how a face is detected](#)

# Camouflage from face detection.

CV Dazzle explores how fashion can be used as camouflage from face-detection technology, the first step in automated face recognition.

*From all appearances, deception has always been critical to daily survival—for human and non-human creatures alike—and,*

## FACE RECOGNITION TECHNOLOGY:

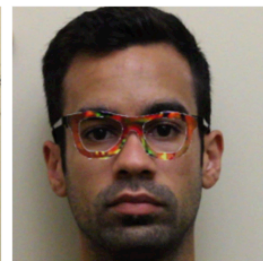
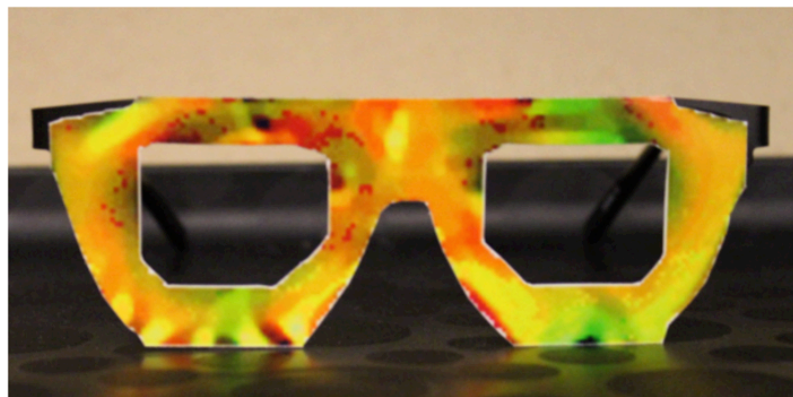
### FBI Should Better Ensure Privacy and Accuracy [Reissued on August 3, 2016]

GAO-16-267: Published: May 16, 2016. Publicly Released: Jun 15, 2016.

#### HIGHLIGHTS

#### What

The Department of Justice's Federal Bureau of Investigation (FBI) has been enforcing a return a 2015, the



ACY & SECURITY

## Facebook's Is Different

May 16, 2016 · 9:30 AM ET

NAOMI LACHANCE

## N-GRAM-BASED AUTHOR PROFILES FOR AUTHORSHIP ATTRIBUTION

VLADO KEŠELJ<sup>†</sup> FUCHUN PENG<sup>‡</sup> NICK CERCONE<sup>†</sup> CALVIN THOMAS<sup>†</sup>

<sup>†</sup>*Faculty of Computing Science, Dalhousie University, Canada*

{vlado, nick, thomas}@cs.dal.ca

<sup>‡</sup>*School of Computer Science, University of Waterloo, Canada*

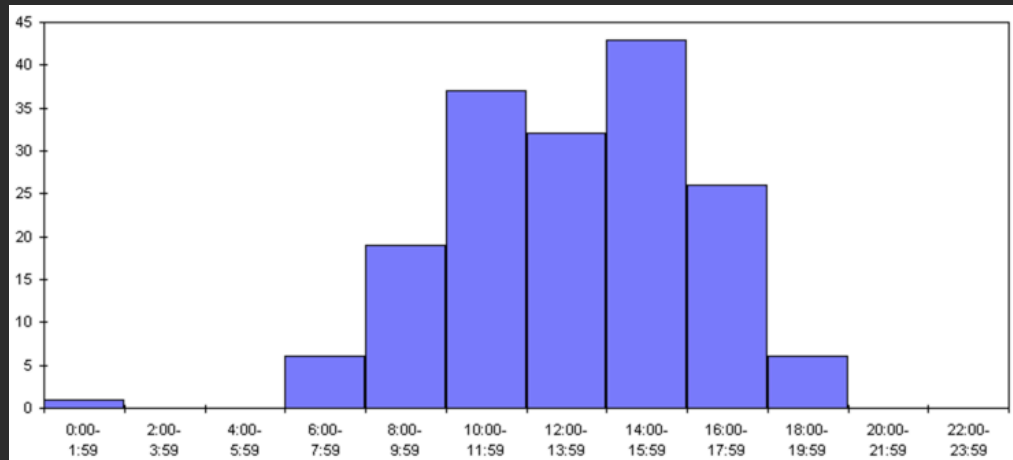
f3peng@cs.uwaterloo.ca

We present a novel method for computer-assisted authorship attribution based on character-level  $n$ -gram author profiles, which is motivated by an almost-forgotten, pioneering method in 1976. The existing approaches to automated authorship attribution implicitly build author profiles as vectors of feature weights, as language models, or similar. Our approach is based on byte-level  $n$ -grams, it is language independent, and the generated author profiles are limited in size. The effectiveness of the approach and language independence are demonstrated in experiments performed on English, Greek, and Chinese data. The accuracy of the results is at the level of the current state of the art approaches or higher in some cases.

*Key words:* Authorship attribution, character  $n$ -grams, text categorization

# Your writing

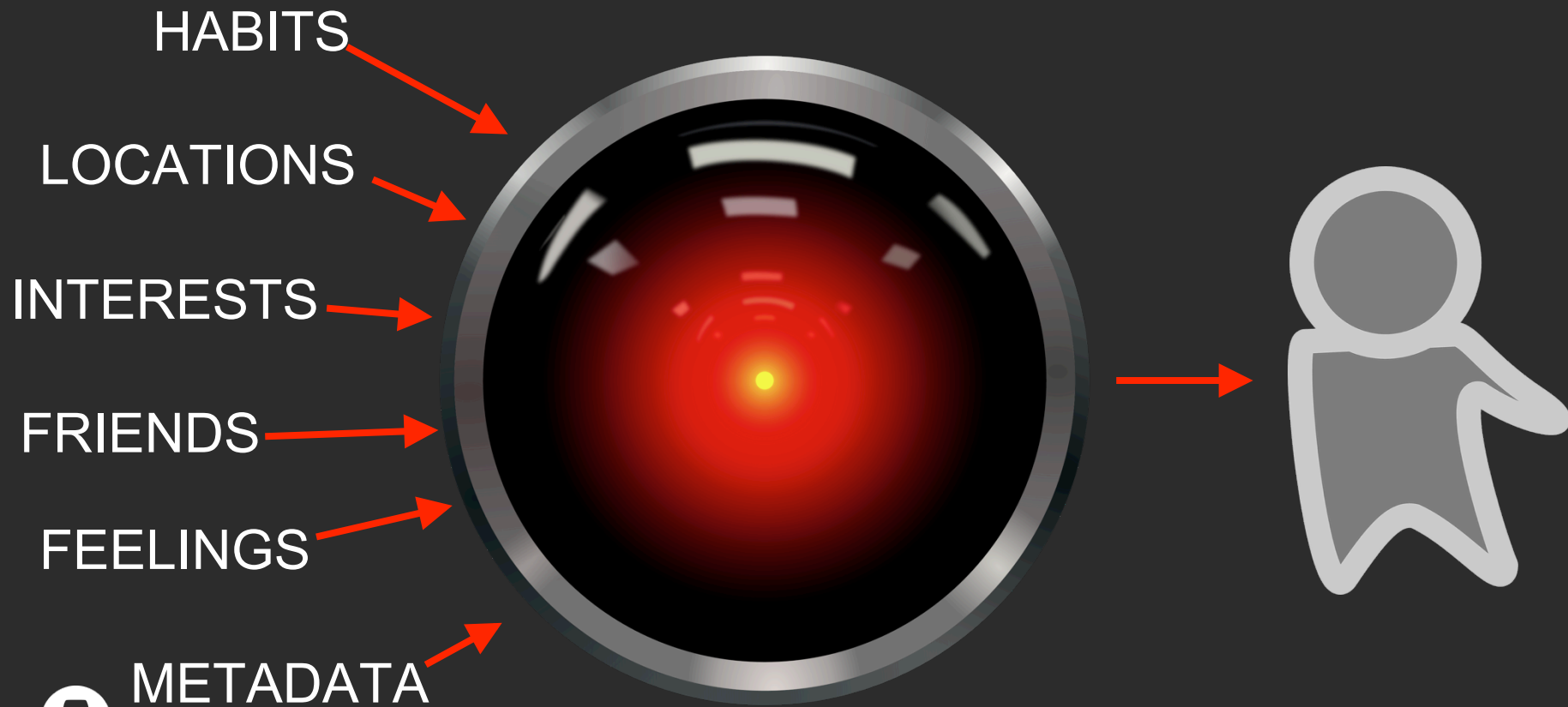




# Your timing



# Pattern of Life Analysis





# Personas

Everything you are, do, or  
use is linked to a persona.

If an activity wasn't  
intentionally done under  
an alt persona, it's linked  
to your true persona.

Anything can link two  
personas. Anything.

Links are bad.

mmkay?



# The Grugg's OPSEC: Because Jail is for wuftpd

# OPSEC Refresher

“Paranoia doesn’t work retroactively” - The Grugg

For any persona:

One account    - One persona    - One activity

One phone no.   - One persona    - One activity

One connection - One persona    - One activity

One home        - One persona    - One activity

One device      - One persona    - One activity

One location    - One persona    - One activity

One *anything else?*



One. Use.

“But Dave, what about a \_\_\_\_\_?”

One.

Use.

# Other Options?



ReactionGifs.me

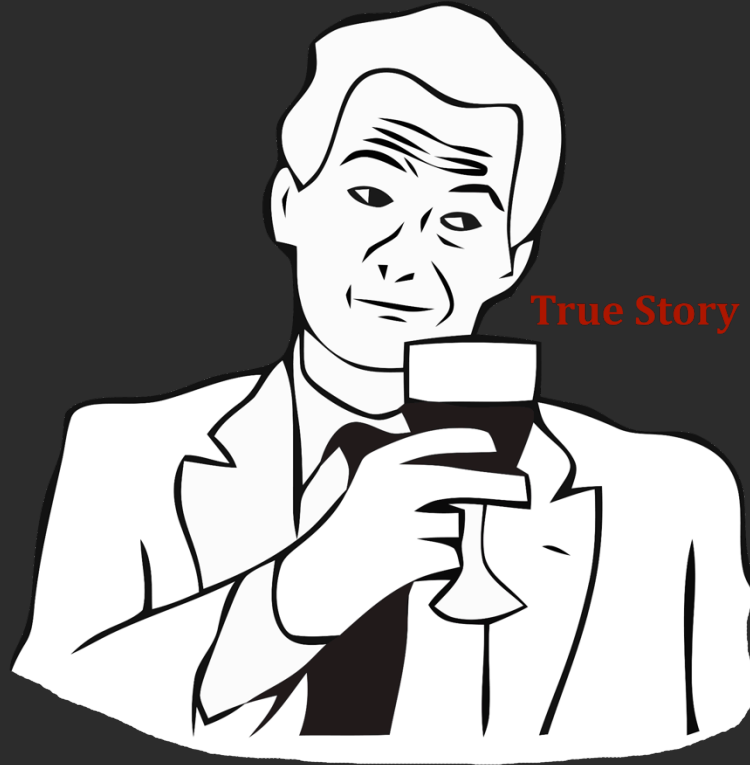
When a new American diplomat arrives for duty at the U.S. embassy in Moscow or Beijing, CIA officials say, Russian and Chinese intelligence operatives run data analytics programs that check the “digital dust” associated with his or her name. If the newcomer’s footprint in that dust – social media posts, cell phone calls, debit card payments – is too small, the “diplomat” is flagged as an undercover CIA officer.

Source: Reuters

# OPSEC Fails



# Bombing the Test



# Bombing the Test

1. Go to Harvard
2. Take difficult class
3. Don't study for final exam
4. How do you get out of taking the exam?



Email a Bomb Threat!

# Bombing the Test

1. Started up Tor
2. Emailed bomb threat (via Tor) to:
  - a. Two Harvard officials
  - b. Harvard Police Department
  - c. The Crimson (Harvard's newspaper)
3. Went to take the test
4. Evacuated with other students



# Bombing the Test

1. Police determined the email sent via Tor
2. Looked at Harvard's network logs
3. ID'd Tor users at the time email was sent
4. Knocked on doors



# Conventional Advice

VPN = Privacy

Tor = Anonymity

- TOR connection to a VPN => OK
- VPN connection to TOR => GOTO JAIL

There's no one size fits all for OPSEC

# What about Bob?



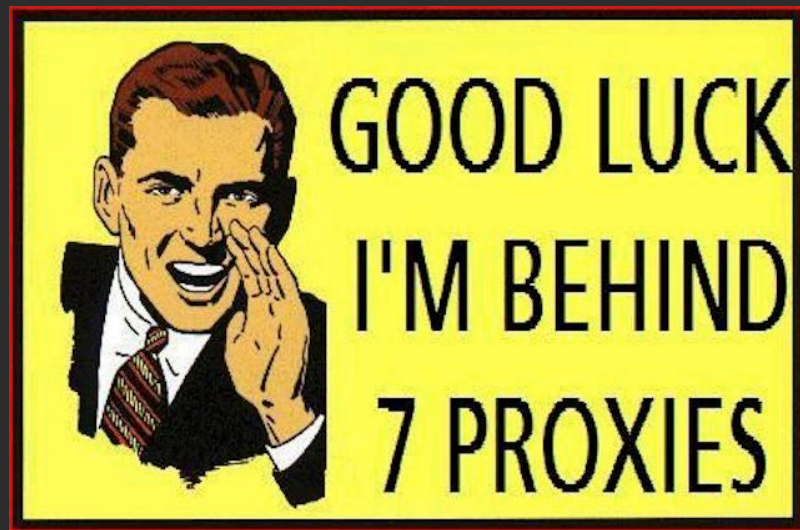
# What about Bob? #1

1. Lives upstairs from a cafe
2. Goes downstairs to use their wifi
3. VM on personal laptop
4. Does [activity]



## What about Bob? #2

1. Drives across town to cafe
2. Brings phone but turns it off
3. Uses their wifi - with a proxy
4. VM on personal laptop
5. Does [activity]



## What about Bob? #3

1. Public transport to park with wifi
2. Phone is on at home
3. Uses Tor
4. Separate hardware
5. Does [activity]

# Information Leaks





# Summary

Assume everything is compromised

Be aware of links

Look at your own pattern of life

Adjust your personas' PoLs to be different

Watch out for leaks where you least expect them



# Breaking Big Data

Evading Analysis of the Metadata of Your Life

Dave Venable  
Masergy Communications  
@davevenable