



# Realistic Cyber Ranges & Enabling Machine Learning Within Cybersecurity

Tom Podnar  
Dustin Updyke

Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213

Copyright 2020 Carnegie Mellon University.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

This material was prepared for the exclusive use of Cylab Partners Conference and may not be used for any other purpose without the written consent of [permission@sei.cmu.edu](mailto:permission@sei.cmu.edu).

Carnegie Mellon® and CERT® are registered in the U.S. Patent and Trademark Office by Carnegie Mellon University.  
DM20-0784

# Part I (Tom)

- Who we are and what we do
- The state of cyber ranges today
- Data opportunities





## SEI > CERT > CWD > Realistic Scenario Simulation Team

- Architect & deliver realistic cyber warfare exercises
- Persistent production spec high-fidelity environment (aka: realism)
- Security hardened and tuned
- Multiple engagements lasting days to 7+ months
- Participants are worldwide DoD operators

# Exercise as you Fight



**ENDGAME.**



# Exercise as you Fight, continued





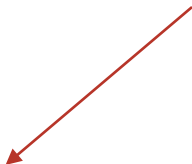
- Topics
- News
- In Focus
- How Do I?
- Get Involved
- About DHS

Home > Topics > Cybersecurity

### Topics

- Academic Engagement
- Border Security
- Citizenship and Immigration Services
- Civil Rights and Civil Liberties
- Critical Infrastructure Security

# Cybersecurity



Our daily life, economic vitality, and national security depend on a stable, safe, and resilient cyberspace.

Cyberspace and its underlying infrastructure are vulnerable to a wide range of risk stemming from both physical and cyber threats and hazards. Sophisticated cyber actors and nation-states exploit vulnerabilities to steal information and money and are developing capabilities to disrupt, destroy, or threaten the delivery of essential services.

# Machine Learning & Cyber Security – Success Stories

- Breach and lateral movement detection and prevention
- Identification of malicious activity and rogue server behaviors
- Malware classification
- **Supplement and enhance human security analysis techniques**



# Machine Learning & Cyber Security – Challenges

- Researchers and engineers face similar challenges to other applications of Machine Learning
- Cyber security is constantly evolving
- **Call to Arms – Data Scientists and Engineers – We need your help!**

# Machine Learning & Cyber Security – Data Challenges

- **Acquiring / Data collection:**
  - Harder to get cyber related data – no one wants to share
- **Training Data**
  - How to obtaining the “right” data sets for a specific cyber security issue?
  - Concerns over data quality and whether data sets are relevant
  - Overfitting and under-fitting related problems
- **How to enable efficient development and testing of your cyber security ideas and theories?**

# Machine Learning in Cyber-Security - Problems, Challenges and Data Sets

Idan Amit<sup>1</sup> , John Matherly<sup>2</sup> , William Hewlett<sup>1</sup> , Zhi Xu<sup>1</sup> , Yinnon Meshi<sup>1</sup> , Yigal Weinberger<sup>1</sup>

<sup>1</sup>Palo Alto Networks

<sup>2</sup>Shodan

iamit@paloaltonetworks.com, jmath@shodan.io, whewlett@paloaltonetworks.com,  
zxu@paloaltonetworks.com, ymeshi@paloaltonetworks.com, yweinberge@paloaltonetworks.com

9

## Data Sets

We think that the use of machine learning in cyber-security should change. We also think that the cyber community should help the machine learning community to become more involved in this field.

One of the key obstacles to investigating cyber-security problems is the lack of appropriate data sets. There are many important cyber-security data sets like Microsoft's malware data set (Ronen et al. 2018), Los Alamos's traffic data set (Turcotte, Kent, and Hash 2017) and EndGame's Ember malware properties data set (Anderson and Roth 2018). However, we feel that there are no suitable data sets that will enable academic researchers to cope with the problems and challenges we listed.

22 Apr 2019



[Read the article online here.](#)

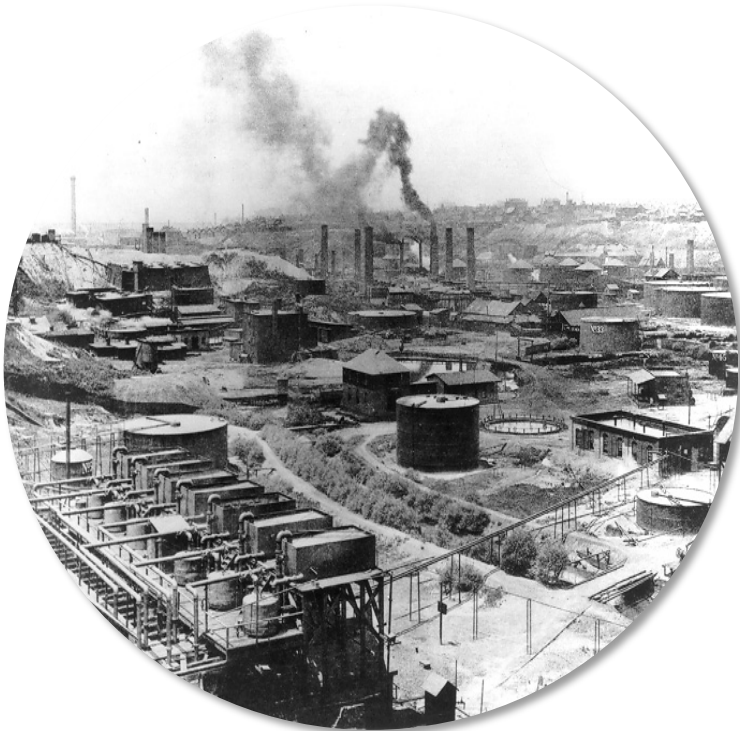


## SEI > CERT > CWD > Realistic Scenario Simulation Team

- Architect & deliver realistic cyber warfare exercises
- Persistent production spec high-fidelity environment (aka: realism)
- Security hardened and tuned
- Multiple engagements lasting days to 7+ months
- Participants are worldwide DoD operators

# Enabling Machine Learning & Cyber Security Success

- **Persistent cyber ranges and resultant data sets:**
  - On-demand access that can be tuned/adjusted
  - Cleaner, labeled data and more balanced
  - Can produce large amounts of unclassified data
  - Ability to re-run cyber security events and scenarios as required
- **Enables efficient development and testing of your cyber security ideas and theories**



# Gasoline

Mid-19th century refinery waste product

Almost worthless — thrown away by the barrel and made Ohio's Cuyahoga River flammable

The oil industry managed to turn its worthless waste product into the most highly sought fuel of the entire 20th century

# Enabling Machine Learning & Cyber Security Success

## Types of data sets:

- Blue, red, and white team perspectives and their interactions
- Non player character (NPC) behaviors and actions
- Raw data in any form: JSON, PCAP, Netflow, protocol metadata, etc.
- **All data can be cross-correlated w/ cyber related event occurrences**

# Part II (Dustin)

## Customer related ML work...

- Large & diverse datasets are a boon for research!
- ...but they also make it hard to get started
- Bonus points: How do we enable other teams?





# Any exercise contains interactions between:

- **Player**

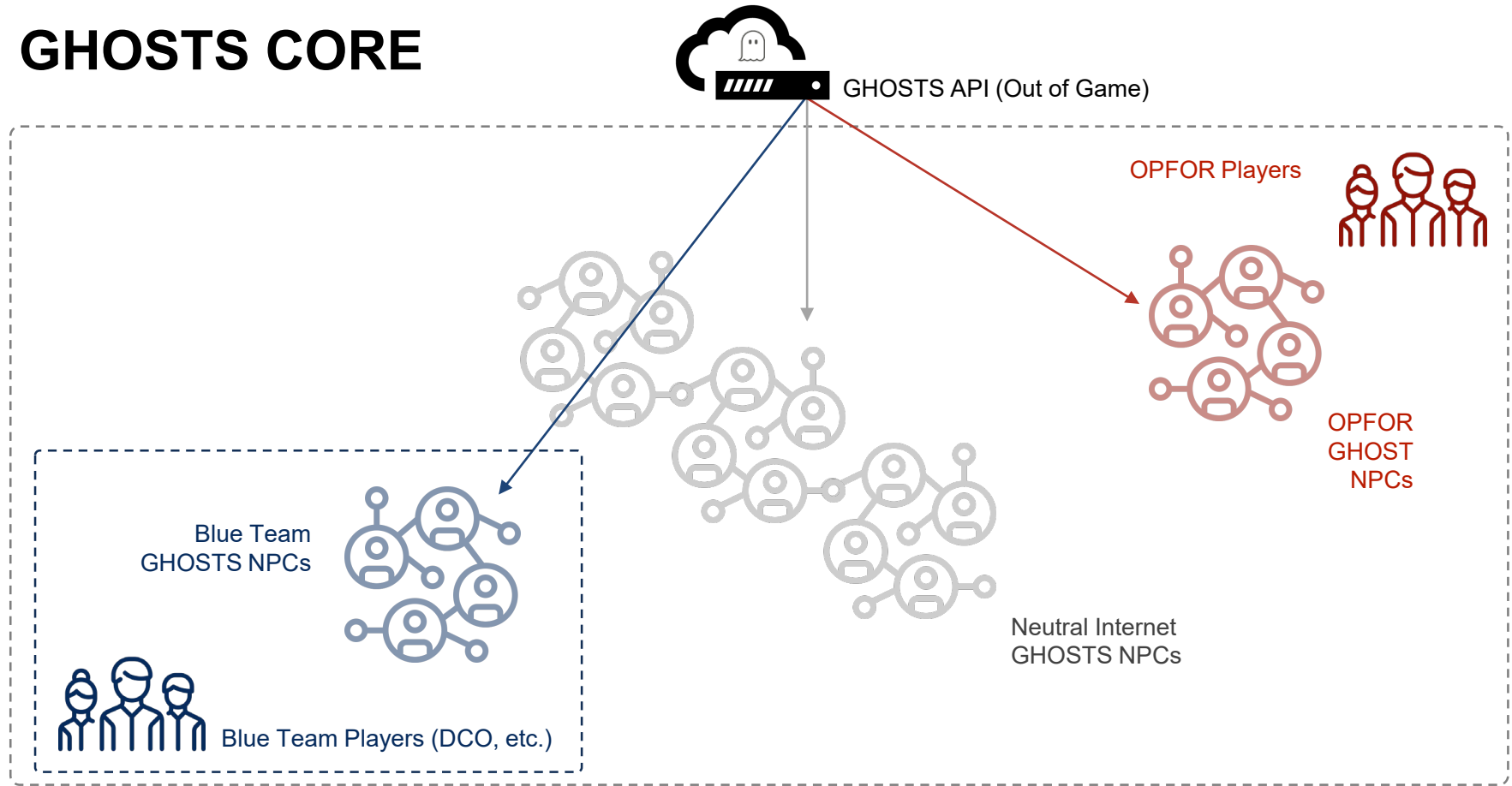
An active human participant

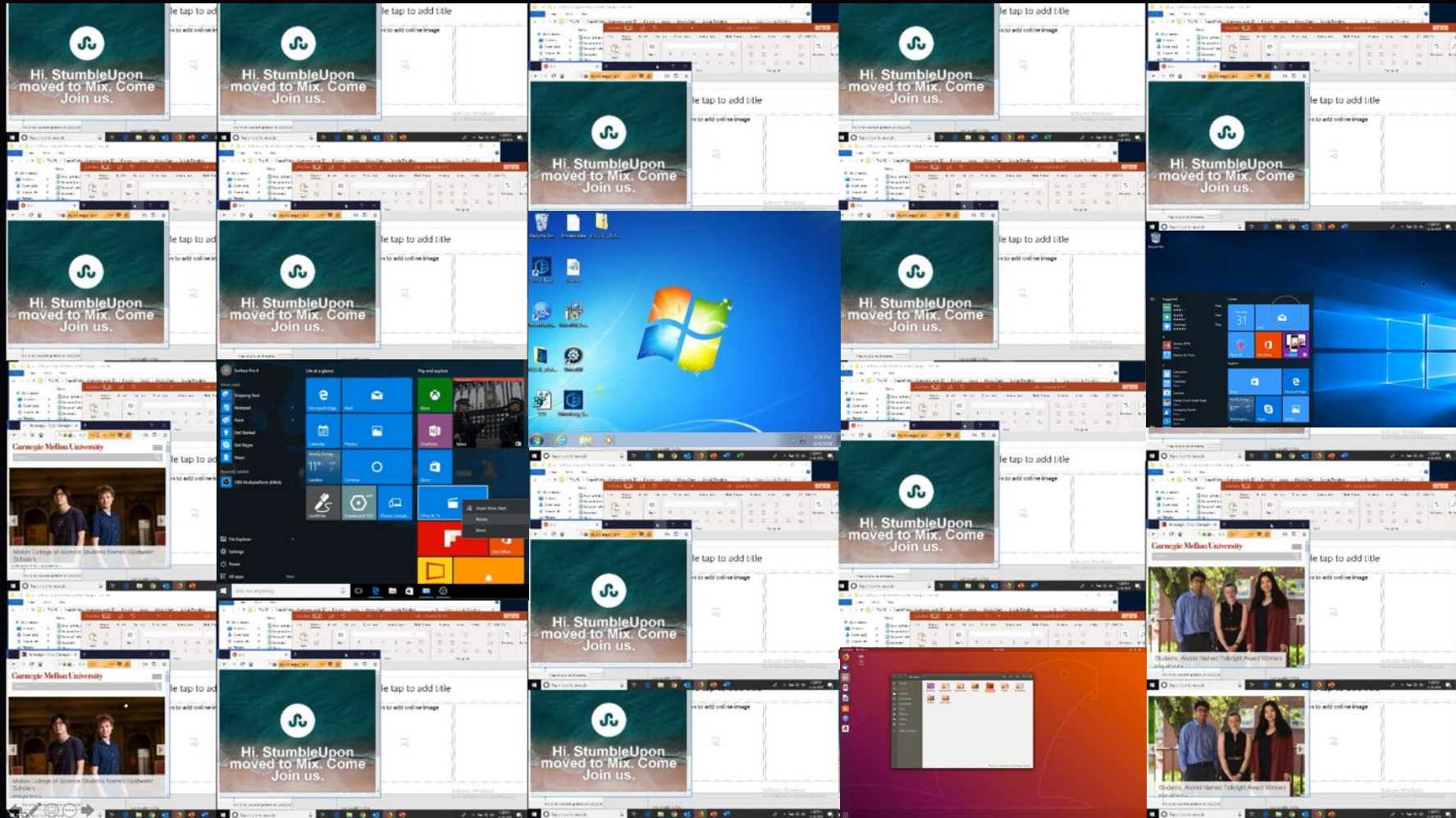
- **Non-player character (NPC)**

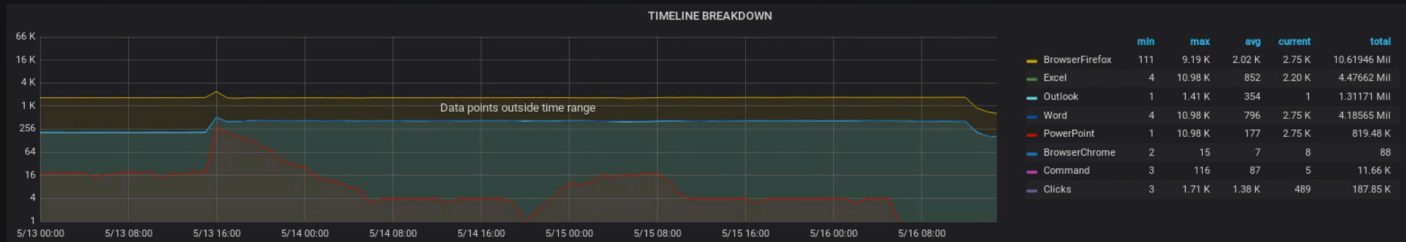
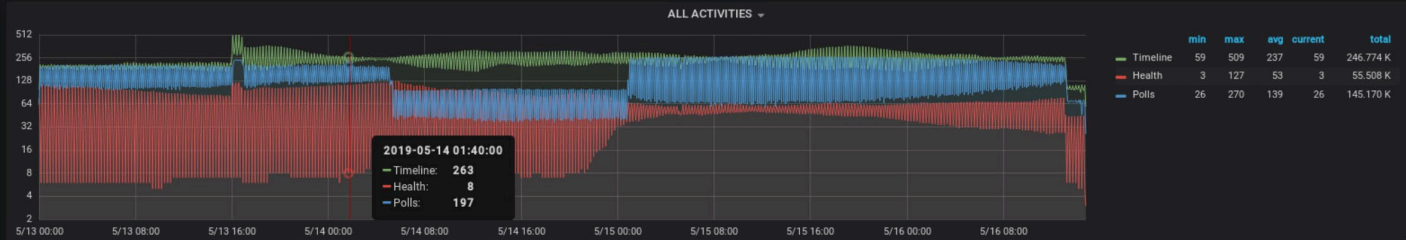
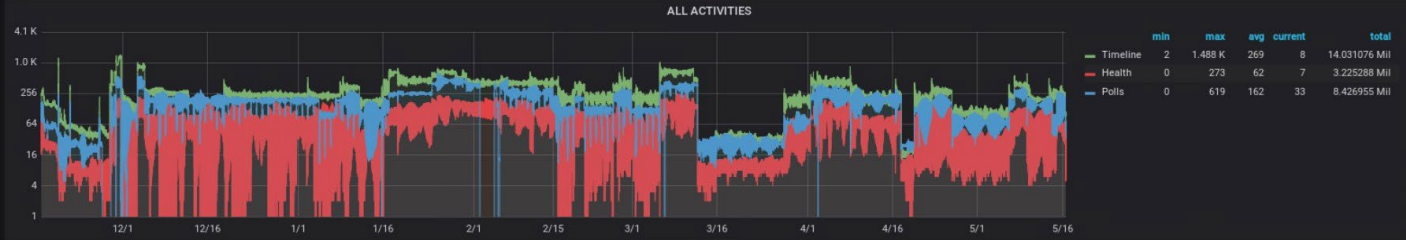
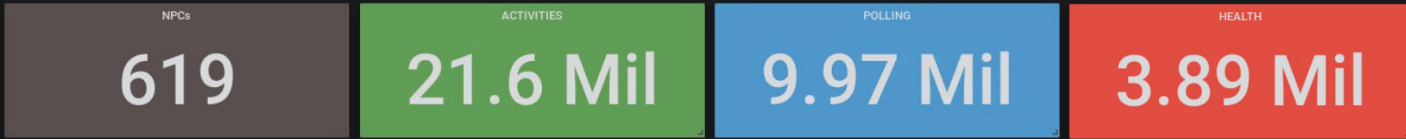
Any character not controlled by a player within the exercise



# GHOSTS CORE







# NPCs Make Decisions Based on their Preferences



**Alexander.Maxey**

+45 Computers

+37 Vim

+55 <https://news.ycombinator.com>

+20 K8s

+45 Postgresql

-10 EMACS

+10 //files/users/amaxey

AutoSave OFF goldilocks

Home Insert Draw Page Layout Formulas Data Review View Tell me Share Comments

D31

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	username	preference	score										
2	Alexander.Maxey	Computers	44.5										
3	Alfredo.Seaman	Kids	44										
4	Allan.Deal	Recreation	32.5										
5	Ashley.Munson	Arts	37										
6	Carissa.Kelso	Reference	37.5										
7	Cecelia.Nunley	Society	43										
8	Clinton.Belt	Recreation	36.5										
9	Conor.Rouse	Reference	52										
10	Dominick.Ragan	News	56										
11	Donte.Gillette	Home	43.5										
12	Emmanuel.Battle	Games	14										
13	Jaron.Lindstrom	Computers	44										
14	Joey.Crowder	Business	35.5										
15	Joseph.Mosley	Reference	31										
16	Kacy.Kinder	News	39										
17	Krystal.Shepherd	Arts	40.5										
18	Lana.Girard	Society	16										
19	Laurie.Fleming	Shopping	40.5										
20	Leslie.Richmond	Society	42.5										
21	Racheal.Denney	Computers	37										
22	Rashawn.Dow	Science	54										
23	Rodrigo.Rojas	Recreation	32.5										
24	Shayne.Fraley	Science	60										
25	Tarah.Meredith	Shopping	37.5										
26	Tracie.Gamboa	Society	44										
27													
28	These were randomly generated from a sample data set of NPC agents												
29													
30													
31													
32													
33													
34													
35													

INTRO USERS AND PREFS 001\_BROWSE\_RANDOM 001\_BROWSE\_PREFS 002\_BROWSE\_PREF 002\_BROWSE\_RAND 003\_BROWSE\_PREF +

100%

AutoSave OFF goldilocks

Home Insert Draw Page Layout Formulas Data Review View Tell me Share Comments

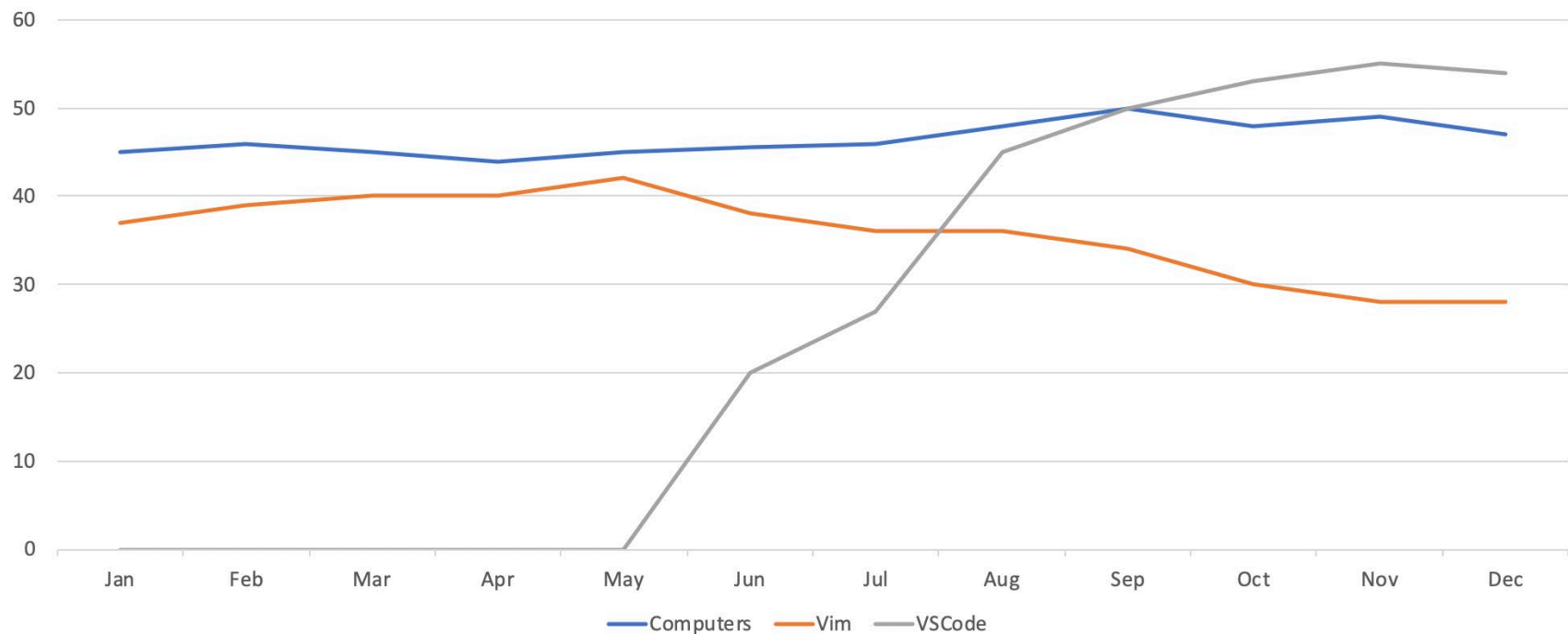
A1

	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	Iteration	(All)																	
2																			
3	Sum of count	Column Labels																	
4	Row Labels	Arts	Business	Computers	Games	Health	Home	Kids	News	Recreation	Reference	Science	Shopping	Society	Grand Total	Final	Start	Gain	
5	Alexander.Maxey	79	122	34726	15	34	13	17	16	39	34	44	42	44	35225	99%	79%	20%	
6	Alfredo.Seaman	65	113	92	15	32	11	11733	18	37	42	44	36	12275	96%	63%	33%		
7	Allan.Deal	84	135	127	18	34	15	26	14	24611	50	49	49	38	25250	97%	68%	30%	
8	Ashley.Munson	21812	123	96	17	39	10	16	27	36	56	42	54	47	22375	97%	70%	28%	
9	Carissa.Kelso	81	106	94	15	38	11	18	31	56	30525	30	45	50	31100	98%	74%	24%	
10	Cecelia.Nunley	94	100	76	15	28	15	9	19	36	45	28	38	34547	35050	99%	79%	20%	
11	Clinton.Belt	89	124	95	16	34	14	22	19	24626	78	49	41	43	25250	98%	68%	29%	
12	Conor.Rouse	66	94	71	8	22	12	14	20	36	30652	41	29	35	31100	99%	80%	19%	
13	Dominick.Ragan	62	79	69	11	30	8	18	20530	27	33	32	33	43	20975	98%	75%	23%	
14	Donte.Gillette	69	102	79	14	36	11991	14	18	41	44	34	46	37	12525	96%	63%	32%	
15	Emmanuel.Battle	102	174	120	3851	59	17	34	25	66	81	55	59	57	4700	82%	26%	56%	
16	Jaron.Lindstrom	73	115	34751	15	31	16	16	19	34	33	35	49	38	35225	99%	80%	19%	
17	Joey.Crowder	91	12146	93	16	33	17	25	23	39	51	43	39	59	12675	96%	64%	32%	
18	Joseph.Mosley	90	123	104	7	39	7	37	21	57	30478	36	44	57	31100	98%	72%	26%	
19	Kacy.Kinder	74	124	96	14	34	19	19	20402	37	44	36	38	38	20975	97%	68%	29%	
20	Krystal.Shepherd	21832	138	94	14	31	10	13	16	34	51	48	55	39	22375	98%	71%	27%	
21	Lana.Girard	103	197	129	20	34	16	33	24	61	70	49	48	33866	34650	98%	67%	31%	
22	Laurie.Fleming	83	113	85	15	35	14	18	19	44	42	34	20548	50	21100	97%	69%	28%	
23	Leslie.Richmond	77	111	94	14	30	11	23	21	37	41	32	52	34507	35050	98%	77%	21%	
24	Racheal.Denney	78	124	34696	14	32	12	27	18	40	51	50	29	54	35225	98%	78%	21%	
25	Rashawn.Dow	66	97	57	13	20	11	12	18	30	30	27227	31	38	27650	98%	80%	19%	
26	Rodrigo.Rojas	80	118	113	21	38	15	27	26	24619	56	46	37	54	25250	98%	68%	30%	
27	Shayne.Fraley	48	81	58	11	17	12	13	7	23	49	27267	36	28	27650	99%	81%	17%	
28	Tarah.Meredith	77	128	83	21	35	12	29	20	46	53	38	20516	42	21100	97%	68%	30%	
29	Tracie.Gambo	91	111	62	13	26	5	27	19	35	65	39	35	34522	35050	98%	78%	21%	
30	Grand Total	45,466	14,998	106,160	4,203	821	12,294	12,240	41,390	74,747	92,749	55,426	42,037	138,369	640,900				
31																			
32																			
33																			
34																			
35																			

INTRO USERS AND PREFS 001\_BROWSE\_RANDOM 001\_BROWSE\_PREFS 002\_BROWSE\_PREF 002\_BROWSE\_RAND 003\_BROWSE\_PREF +

Average: 39.78 Count: 50 Sum: 994.5 100%

# An Agent's Preference Over Time





The screenshot shows a web browser window with the address bar at [www.winzip.com](http://www.winzip.com). On the left, a 'History' sidebar is open, listing various websites visited today. On the right, a list of popular files is displayed, including [PaintShop](#), [VideoStudio](#), [WinDVD](#), [AfterShot](#), [Roxio](#), [Pinnacle](#), [WinZip](#), [CorelDRAW](#), and [Painter](#). A large red arrow points from the text on the right towards the 'WinZip' entry in the file list. At the bottom of the browser window, the text 'The Most Popular File' is visible in blue. The Windows taskbar at the bottom shows the time as 10:58 PM on 5/28/2020.

History

Search history View

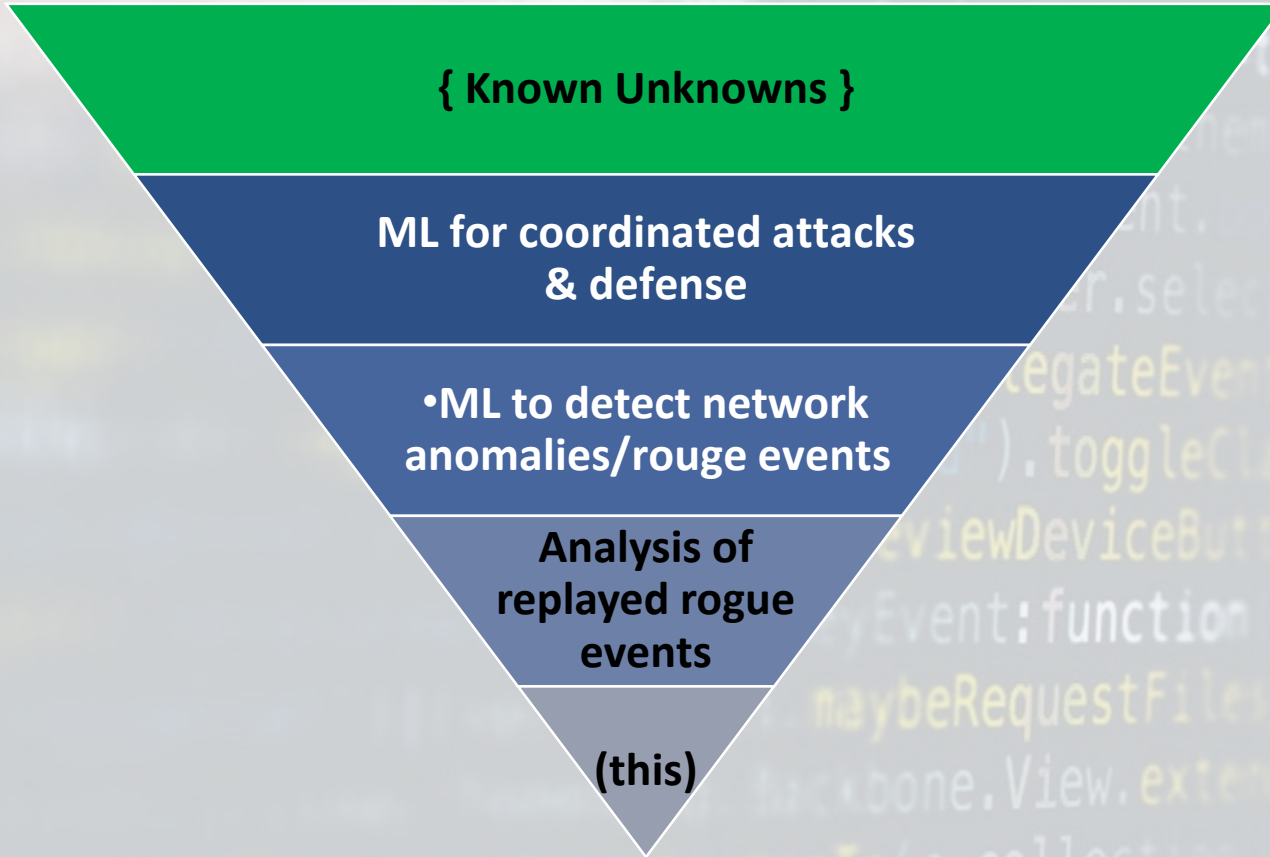
Today

- factcheck.org/
- lowes.com/
- www.darpa.mil/
- 404 Not Found
- ABS-CBN News | Latest Philippine Headlines, Breaking News, Vide.
- Adobe: Creative, marketing, and document management solution
- Android Central | Android Forums, News, Reviews, Help and Andr
- AOL - News, Sports, Weather, Entertainment, Local & Lifestyle
- Ars Technica
- Best Buy: Expert Service. Unbeatable Price.
- Bleacher Report | Sports. Highlights. News. Now.
- Blog Tool, Publishing Platform, and CMS — WordPress
- Booking.com: 653,536 hotels worldwide. 48+ million hotel reviews
- Business Insider
- BuzzFeed
- CafeMom - Moms Connecting About Pregnancy, Babies, Home, F
- California Home Page
- CBS TV Network Primetime, Daytime, Late Night and Classic Telev
- Cell Phones, Smartphones & the Largest 4G LTE Network - Verizon
- Centers for Disease Control and Prevention
- Chicago Newstips by Community Media Workshop
- Chicago Tribune: Chicago breaking news, sports, business, enterta
- City-Data.com - Stats about all US cities - real estate, relocation inf
- CNNMoney - Business, financial and personal finance news

- [PaintShop](#)
- [VideoStudio](#)
- [WinDVD](#)
- [AfterShot](#)
- [Roxio](#)
- [Pinnacle](#)
- [WinZip](#)
- [CorelDRAW](#)
- [Painter](#)

**The Most Popular File**

~25% tighter alignment of browsing to kinds of sites WRT a user's preferences



# Takeaways

1. **All data has ML potential**
2. **ML projects beget larger & more complex projects**
3. **Start simple (even for the most daunting datasets)**
4. **Output of a project feeds the next**