

November 2, 2022

Trolls, Pirates, Zombies, and Clowns: Data Mapping to Vanquish Privacy Foes

Brandon Schneider, IT Data Privacy Attorney, IAPP FIP

Rick Newbold, Data Governance and Privacy SME, LL.M., PSDGP, CIPP-US/G

Chris Lee, Principal Privacy Engineer, MITRE, MS, JD, CIPP-US/G

Aaron Burciaga, Chair, Analytics Certification Board

Additional Contributors:

Dr. Margaret Leary, Ph.D., CISSP, CEH, CRISC, Professor of Cybersecurity

John “Major” Nelson, DHS S&T

Speaker

Mr. Chris Lee

Principal Privacy Engineer, MITRE
MS, JD, CIPP-US/G
Washington, DC







Mr. Newbold and Mr. Schneider are sharing knowledge from the public sector privacy practitioner perspective.

Mr. Burciaga is sharing knowledge from the data analytics perspective.

Dr. Leary's material is from a cybersecurity and AI perspective.

Any reference to a private sector entity does not constitute an endorsement.

Views expressed do not reflect an official position of the U.S. Government.

NIST supplemental guidance is not copyrighted, and graphic owners are credited if graphics are not in the public domain.

The workshop presenters are still learning about the commercial market landscape and about the various vendor PII data mapping offerings.



Data Mapping is the process of inventorying personal data within business systems and is a means of creating a data inventory.

Data Inventory and Data Mapping will be used interchangeably for purposes of this workshop.

Given that a Personal Information Map can reflect a “snapshot” in time, PII data mapping may more accurately reflect PII at rest than PII in transit.

Everything that is declared through a PIA or other documentation is not necessarily grounded in truth.

Trust but verify. Obtain evidentiary proof of the truth through evidentiary artifacts and by examining, observing, and testing.

We assume the audience:

Understands the difference between security risk and privacy risk.

Appreciates presenters may have a bias toward Federal government support relative to the privacy sector.

Privacy and Security: Risk vs. Reward



	Privacy Risk	Security Risk
Confidentiality	Information having the ability to harm a person when released.	Only authorized users and processes should be able to access or modify data.
Integrity	Accuracy and consistency (validity) of data. PII accurate, relevant, timely, and complete.	Securing against improper information modification or destruction.
Availability	The data is relevant and necessary to accomplish the specified purpose(s) and only retained for as long as is necessary to fulfill the specified purpose(s). Individual Access	Information is accessible and usable upon demand by an authorized person.



Problem Statement

How do one manage, monitor, and protect personal information – minimizing risks and maximizing rewards?

Challenges:

Where is your data?

Is the data consistent and accurate?

Who has access?

What are they using it for?

Are there new uses?

How are retention and archival/destruction requirements enforced?

Opportunities:

PII mapping helps protect data, minimize risks, and maximizes the potential monetization and benefits.

Answers the questions above.



Privacy
Continuous
Monitoring
Definition:

Maintaining **ongoing awareness** of privacy risks and assessing privacy controls at a frequency sufficient to ensure compliance with applicable privacy requirements and to manage privacy.

Sources:

OMB Circular A-130

NIST Special Publication 800-137

NIST Special Publication 800-53, Rev. 5

Mr. Aaron Burciaga
Senior Data Analyst
Washington, DC

- Where is your PII, and where is your organization's PII?
- No single or exact way to map but better to have a map which can adapt to organization's needs
 - PII data flow diagram or privacy information map (PIM) or
 - Mapping PII rich locations reflected in a heat map
 - Diagram PII at rest and in transit

The Troll = The Hoarder

- The Troll does not share.
- Detection of users who do not reciprocally share.
- Observation that a collaboration site user does the relatively least sharing.
- Vanquish = target → increase awareness → discipline → terminate

The Pirate = The Rule Breaker

- The Pirate steals your data, uses your data with reckless abandon.
 - Primarily external unauthorized access
- Audit/audit logs can detect stealing.
- User and administrator analytics can detect stealing. Mapping can supplement this detection.
- If data deficiencies or losses are constantly located in a particular location, mapping can detect this location activity and detect the stealing.

The Zombie = Brainless

- The Zombie does not have their own data;
 - will consume all of your data without regard
- This consumption without reciprocity can imply zombies have access either:
 - to either privately held data; or
 - to access to publicly available data.
- Where a data exchange exists with particular data takers who never contribute data = zombie consumption is detected!
- Mapping can supplement administrator/user analytics; auditing/audit logging.

The Clown = Chaos

- The Clown represents no rules and operates in absolute chaos.
 - data everywhere and solutions packed into a clown's car.
- Where a user touches/accesses data and no rules or governance and data must be carefully managed back in order, clowns are the likely suspects.
- Data mapping can help detect data chaos or anarchy.
- If a segment or map quadrant—location is chaotic relative to other map/most other map locations, it's likely a clown.
- Rogue Databases: the Clown is often the culprit.

- Fulfills privacy regulatory compliance objectives.
 - Data mapping enables us to fact check PTAs/PIAs
- Facilitates business processes.
- Aids with managing risk and associated strategic risk-based decision making.
- Supports/facilitates/conducts PII redaction on data selected for public distribution and consumption.
- Achieves better privacy data governance through greater visibility of data holdings.

- Provides a bigger picture understanding of which systems interact/interoperate both on an enterprise-wide scale and on an interagency basis.
- Identifies where unnecessary/redundant SPII proliferation exists.
- Enables better detection of which downstream stakeholders have a share of an organizations PII data holdings. Attain a clearer picture of information sharing through a map of PII in transit (as the many Federal PTA templates requires an analysis of internal and external information sharing).
- Helps detect privacy foes.

- Drives annual budget planning based on privacy—security vulnerability risk need (resources reallocation).
 - With SPII-dense locations identified, Privacy can work with Security to analyze whether these locations are potentially vulnerable and present significant risk. If at risk, leadership can devote additional resources to these locations for a greater defense in depth posture.
- Measures the damage potential of PII data leakage.
- Supports a response to an audit (e.g., IG Federal Information Security Modernization Act (FISMA) audit or an audit pertaining to OMB regulations, NIST controls or SSN accountability/reduction).

- Provides a holistic understanding of PII data flows within an organization.
- Calls Attention to SPII (e.g., SSNs) locations within an organization which could, in turn, reveal SPII exposure risk (e.g., the existence of SPII where security controls are not implemented adequately).
- Helps (automated) implementation of SSN alternate identifier.
- Fulfills inventory requirement: Identification and tagging all PII data types Identify areas for improved (or new) efficiencies with an organizations business processes, IT systems, and IT controls.

Breakout Discussion Session I

Discussion Questions, Session I

Of the four characters (Troll, Pirate, Zombie, and Clown):

- Which character do you personally identify with?
- Which character might be easier to detect?
- Which character would you most want to vanquish?



Have you encountered vendor data mapping tools (within the last three years) that incorporate our proposed requirements?

Do you think that the Federal guidance and best practices to have a data inventory or map will gradually become a lighter lift and more commonplace? (show of hands, and let's discuss)

Mr. Brandon Schneider

Senior Privacy Lead

IAPP FIP

Washington, DC

PII Data Mapping: the process of inventorying PII through illustrating major interconnections for PII in transit and major storage locations for PII at rest.

- Go back to where you started to remember who you are.
- PII data mapping was not prominent in IT security strategic plans.
- Traditionally, PII maps (or Privacy Information Maps (PIMs)) mapping
 - reflected a “snapshot” in time
 - reflected PII at rest more commonly than PII in transit



Source: The Open-Air Mission

A Look Back – Risk Determination

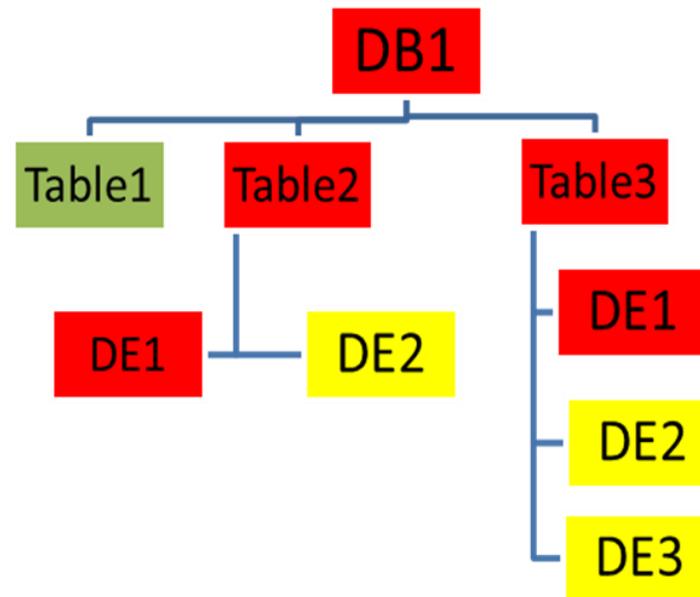
Risk of data breach is the product of Probability and Impact.

Probability	Risk = Probability x Impact		
Certain 0.900	Low 0.045	High 0.180	High 0.720
Likely 0.700	Low 0.035	Medium 0.140	High 0.560
Unlikely 0.100	Low 0.005	Low 0.020	Medium 0.080
Impact	0.050 Limited	0.200 Serious	0.800 Severe

<i>PII Risk Legend</i>	High	Medium	Low
------------------------	------	--------	-----

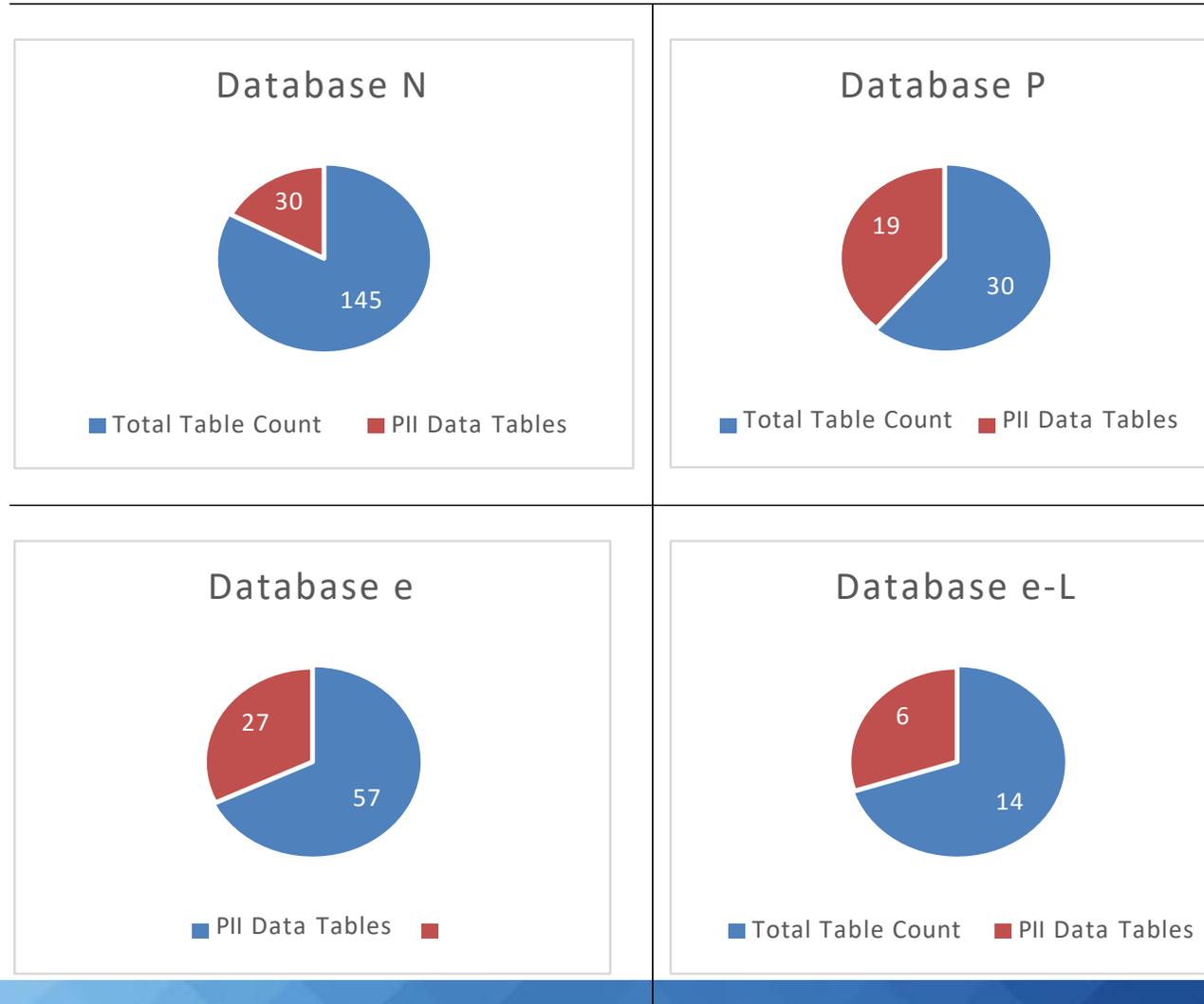
Risk Determination for Data Table Using Weighted Average Formula

Weighted Average Formula =
Sum (Element-level Risk * Number of data elements with that risk) /
Total number of PII data elements in the table



A Look Back: PII Tables Comparison Chart

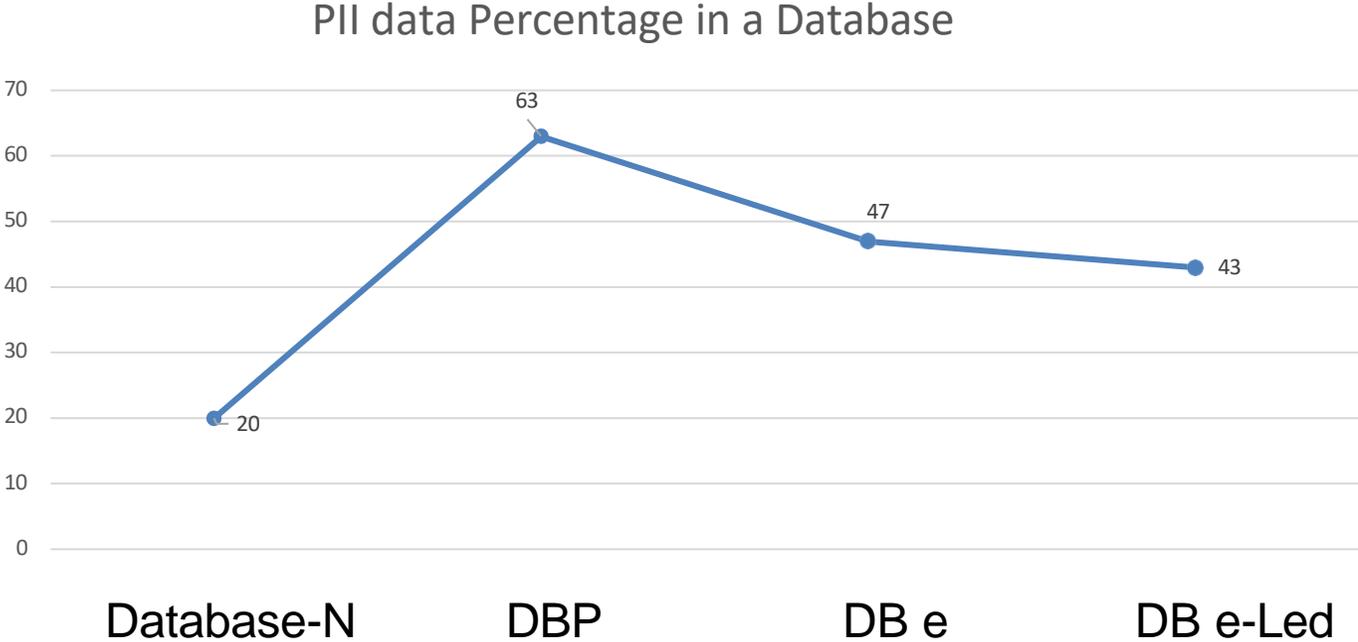
Total Tables vs. PII tables in four databases analyzed.



A Look Back: Data Tables Sensitivity Comparison Chart

Sim Disability	Sim Ethnicity	Sim Race	RES Contract
RES Contract Participant	Sim Congressional District	Sim Conservation District	Sim County
Sim Huc	RES Contract Account Type	RES Contract Item	RES Contract Modification
RES Contract Resource Concern	Sim Crop	Sim Livestock	Sim Office
Sim State	RES Contract Crop	RES Contract Item Resource Concern	

A Look Back: PII Data Percentage Comparison of Databases



PII Data Mapping Must Haves

- ✓ Targeting or tracking SSNs and or 9-digit recognition is a MUST
- ✓ Mapping other sensitive data (e.g., SPII) is also important
- ✓ Data at rest (databases) and data in transit is illustrated.
 - include which applications are accessing PII.
- ✓ Internal and external sharing interconnections are diagrammed
 - where sensitive information transits through and beyond the organization (using color-coded transmission lines).



PII Data Mapping Must Haves

Evolving and Adaptable Maps:

- ✓ Ability to readily:
 - modify map as new information is revealed (e.g., after PIAs, privacy risk assessments, or new projects/programs/systems arise);
 - maintain a current/updated map (i.e., at least quarterly) is important (as known information evolves/transits/ages).
 - Consider automating continuous monitoring.
 - Be Forward Thinking.

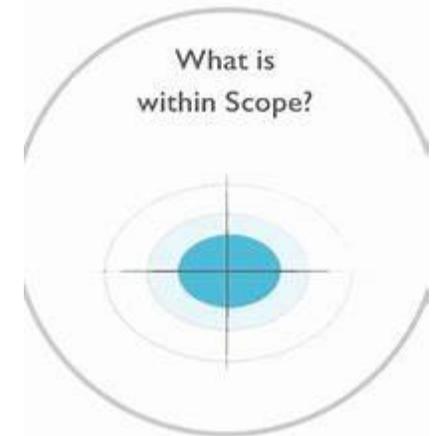


- ✓ Illustrate basic security architecture (e.g., behind additional firewall).
- ✓ Include standard administrative functions which transcend multiple major organizational divisions, especially if notoriously PII dense/sensitive (e.g., human resource data).
- ✓ Note touch points where major organizational unit data could be properly commingled.

- ✓ Note stages in the PII lifecycle, including where sensitive information is (i) created, (ii) altered, and (iii) destroyed.
- ✓ Good working relationships with CISO, OCIO, data architects, ITPMs, and system owners.
- ✓ In a later phased upgrade, the ability to automate with less reliance on manual (less efficient) mapping.

Effective preliminary scoping is critical.

- ✓ Do not map enterprise-wide initially
- ✓ Scoping Considerations:
 - Federal Information Processing Standard (FIPS) “Moderate” Applications
 - in production
 - are high profile or subject to political sensitivities/pressures.
 - prominent or higher magnitude of enterprise importance



Source: h3uni.org

- ✓ Scope to systems which generally have a higher magnitude of enterprise importance due to:
 - the high number of users (high usage rate),
 - person-centric data holdings, and
 - security categorizations (e.g., High for Confidentiality).

- ✓ Don't "boil the ocean" – no need to initially include:
 - all system subcomponents (e.g., SharePoint subsites)
 - PII within emails
 - expired data past scheduled disposal date
 - unstructured data

Break

Discussion Session II

Discussion Questions, Session II

Non-privacy professionals may not fully appreciate the full context in which to identify aggregate PII or SPII.

- Have you experienced a discrepancy in the ways privacy professionals vs. non-privacy professionals identify aggregate PII?
- Do you believe that an automated data mapping tool can fully replace the trained “human eye” in the PII data mapping process?

Have you ever wondered what a map of all your consumer related PII (or SPII in particular) across the commercial sector would look like?

- Where would most of your PII be located?
- Where would you most sensitive PII be?



Discussion Questions, Session II

Are both the public sector and the commercial sector interconnected?
For example, a commercial airline and the FAA, or TSA?



Are you more confident that your PII is safe in the public sector or the commercial sector (e.g., OPM, Equifax)?

Do you feel that the challenge of corralling the truthful complete picture necessary to map out data is due to:

- i. lack of project team knowledge;
- ii. desire by project team to hide the ball to be protective of their system/application from a security risk standpoint;
- iii. due to project team not wanting to create more work for themselves;
- iv. a combination of the above?
- v. (show of hands, and let's discuss)

- TRUST BUT VERIFY.

- Obtain evidentiary proof of the truth through obtaining evidentiary artifacts and through
 - examining
 - observing
 - testing
- a) Trust but verify the content to be inserted into the map.
- b) The map = an evidentiary artifact, supporting both:
 - privacy risk assessment process and
 - PIA declarations



Source: QuotesBook

- Initiate the Mapping Process when:
 - your organization has the time and resources
 - you experience a trend/indications of losing track of data, particularly SPII
- Maintain and share the map among trusted colleagues
- Cross-Collaborative Team effort:
 - data architect/Enterprise Architecture team
 - security engineer
 - privacy professional



Source: PeelMaster

- Key: Privacy Professional should analyze PII because non-Privacy professionals may not fully appreciate aggregate PII (“untrained eye”)

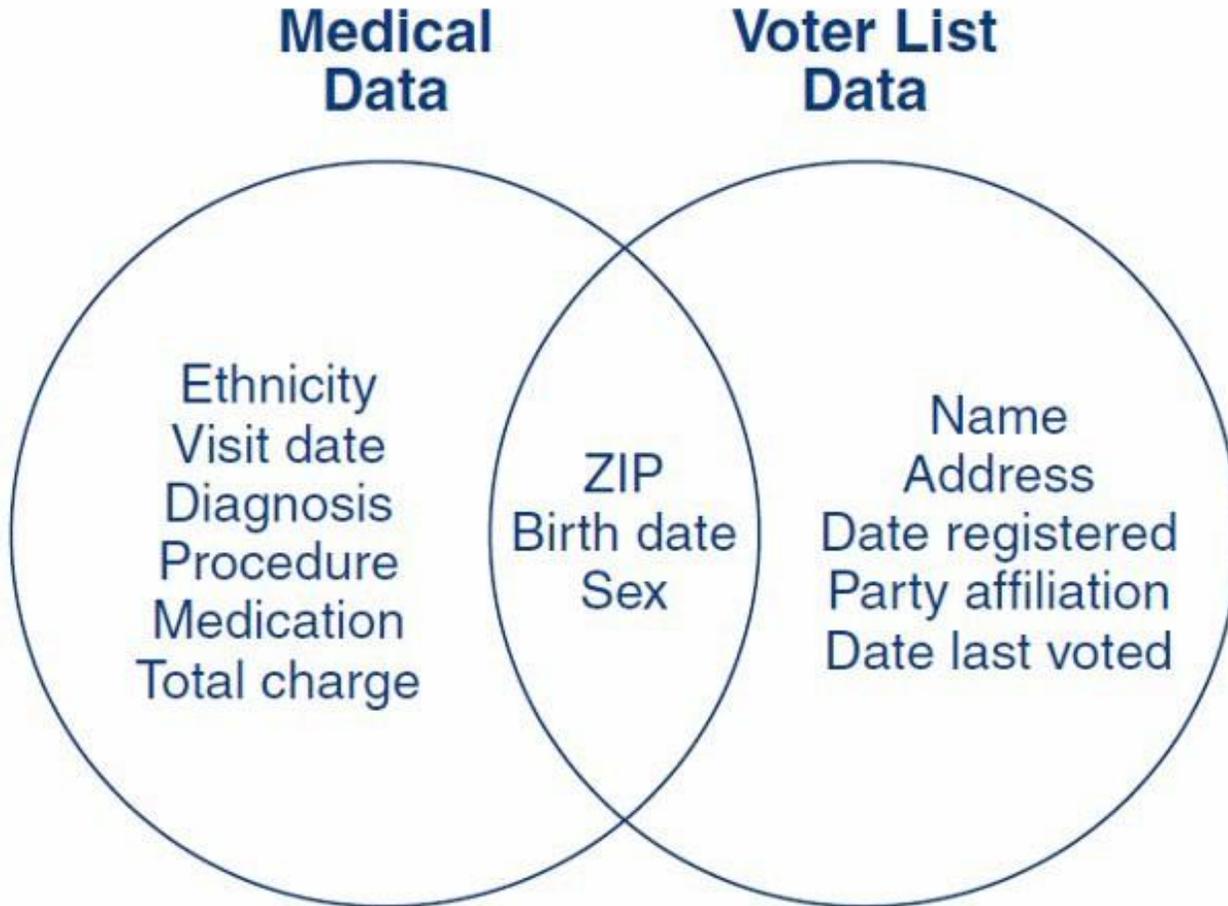
- Understand the distinction between:
 1. apparent PII: PII which by itself is singularly highly linkable to an individual.
 2. aggregate PII: A combination of PII such that, taken together, the combination increases the chances of correctly identifying an individual.

- Map an organization where higher dense/sensitive PII exists.
 - It is often through higher density PII that we have the higher privacy sensitivity.

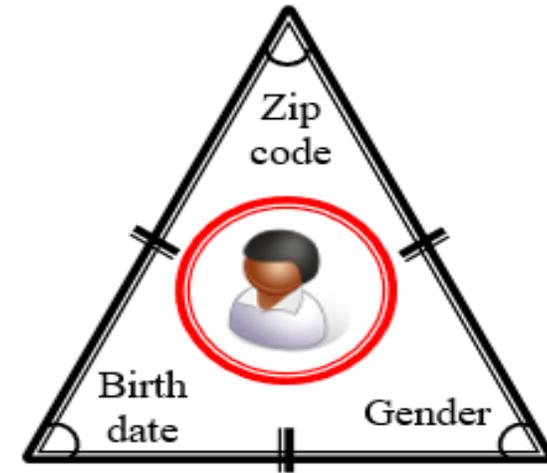


Source: Clipground 2021

The Problem with Aggregation



Information



Person = 87% chance of being identified

If you know a person's zip code, date of birth, and gender, then there is an **87%** chance you can correctly identify that person.

Latanya Sweeney, k-anonymity: a model for protecting privacy. *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 10 (5), 2002; 557-570.

Example Approaches:

- 1) Hybrid: Systems/Processes
- 2) Top down/Bottom up
- 3) Application/Database

- Analyze Both Application and Database
 1. Identify the PII data elements (PII analysis) within each:
 - a. **Application** (front end via a demo, from user perspective)
 - b. **Database** (back-end via the data architect, not visible to public)
 2. Classify the PII data elements with additional metadata; and
 3. Score the PII data elements for risk analysis.

- Analyze both apparent and aggregate PII and their privacy risk impacts.

- Maintain focused analysis on finding aggregate PII context patterns within in scope applications and databases.

Aggregate PII Context Pattern List

Column Name	Purpose
PII Context ID	ID for each Aggregate PII Context Pattern
Method of Discovery	Aggregate PII Context Pattern found through Database or Application observation
Source Database	Name of the Database where the PII data element was found
Source Table	Name of the table where the PII data element was found
Application Name	Name of the application where the PII data element was found
Report Name	Application query or data entry form where data element was found
Apparent PII	Name of the obvious PII data element
Aggregate PII	Names of the data elements that form aggregate PII
Weight	Percentage assigned to the aggregate PII context to be identified to an individual. Apparent PII can be 100% identified to a single individual therefore weight is 1.

Why Don't We Get What We Need for Data Mapping

- Colleagues may not provide needed data mapping information, because colleagues:
 - just do not know,
 - are “hiding the ball,”
 - or both of the above.
- IT Development/Project teams actually do know, but:
 - want to meet go live date (get up and running);
 - get paid to get things finished, and revealing mapping information could make a lot more work for the team;
- Finding new linkages vis-à-vis data mapping happens.
- Intentional non-disclosure happens.



Source: MLB Hidden Ball Trick (HD) – YouTube

Challenge: Data dictionaries/models are not up to date – data dictionaries/models do not reflect the current status.

- Dictionaries do not match.
 - Some of the critical information, (e.g., description of the data elements) is missing
 - Does not yield the necessary information to verify whether a particular data element is PII (e.g., “client ID” = unclear if PII by itself)
-
- Access to database structure may be limited to the Database Administrator (DBA). Mapping team members may not gain any access to the data dictionaries/models.

Recommendation: Ensure leadership support/political backing particularly when the mapping team needs access to well guarded enterprise architecture.

- Leadership may even request that the mapping team members have access, but true/full access may not be provided.

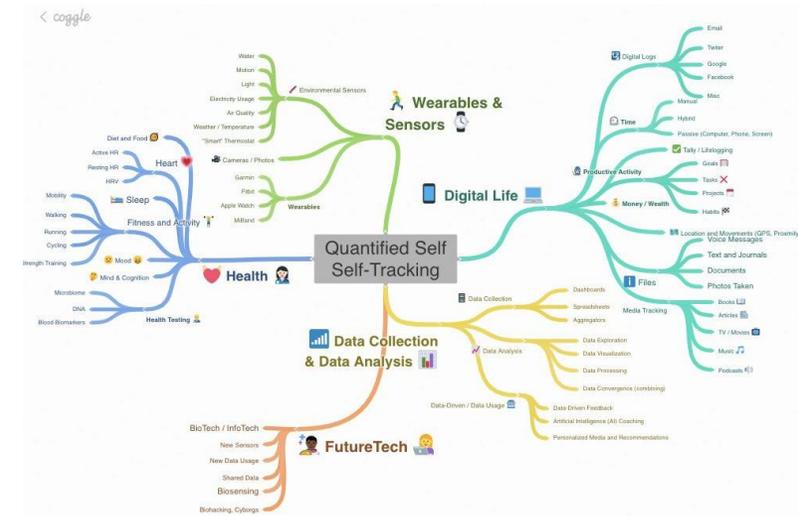
- Secure Project Teams' assistance to gain application access, and to provide support for analysis of any ambiguous data elements and application navigation.
- Have the technical project team review the map. This team may be aware of locations/location positions that other map contributors may not have considered.
- Draw Federal Information Security Modernization Act (FISMA) boundaries around proper map subsections which can help make the map more universally valuable to multiple disciplines (security, privacy, management, auditors).
- Keeping the map on Macro, or higher level, can still be effective (Metadata). Do not be overly concerned with illustrating specific data storage and flows for every individual PII element (or attribute).

- Direct access to a tool like Collibra could lesson reliance on second hand faulty data.
- Share “finished” data map phase with data governance team for broader organization-wide benefit.
- Remain Vigilant: Conduct routine spot checks for rogue databases or "without permission" sites.

Breakout Discussion Session III

Discussion Questions, Session III

- As this data could possibly become selectively absorbed into government systems, could this be viewed as data proliferation?
- Continuous monitoring may be a healthy compromise on the way to near-real time. Do you agree or disagree?
- Do you believe that real-time data mapping is more technically feasible these days?
- What have you seen out there or demo'd that gives you confidence/optimism? (show of hands and let's discuss)
- How else to you envision the data maps of the future?
- With data mapping's multiple benefits, why wouldn't you want to make the transition to more mature PII data mapping?
- Do time, personnel, money have an impact?
- Does some other factor impact this decision?



Rick Newbold

Data Governance and Privacy SME

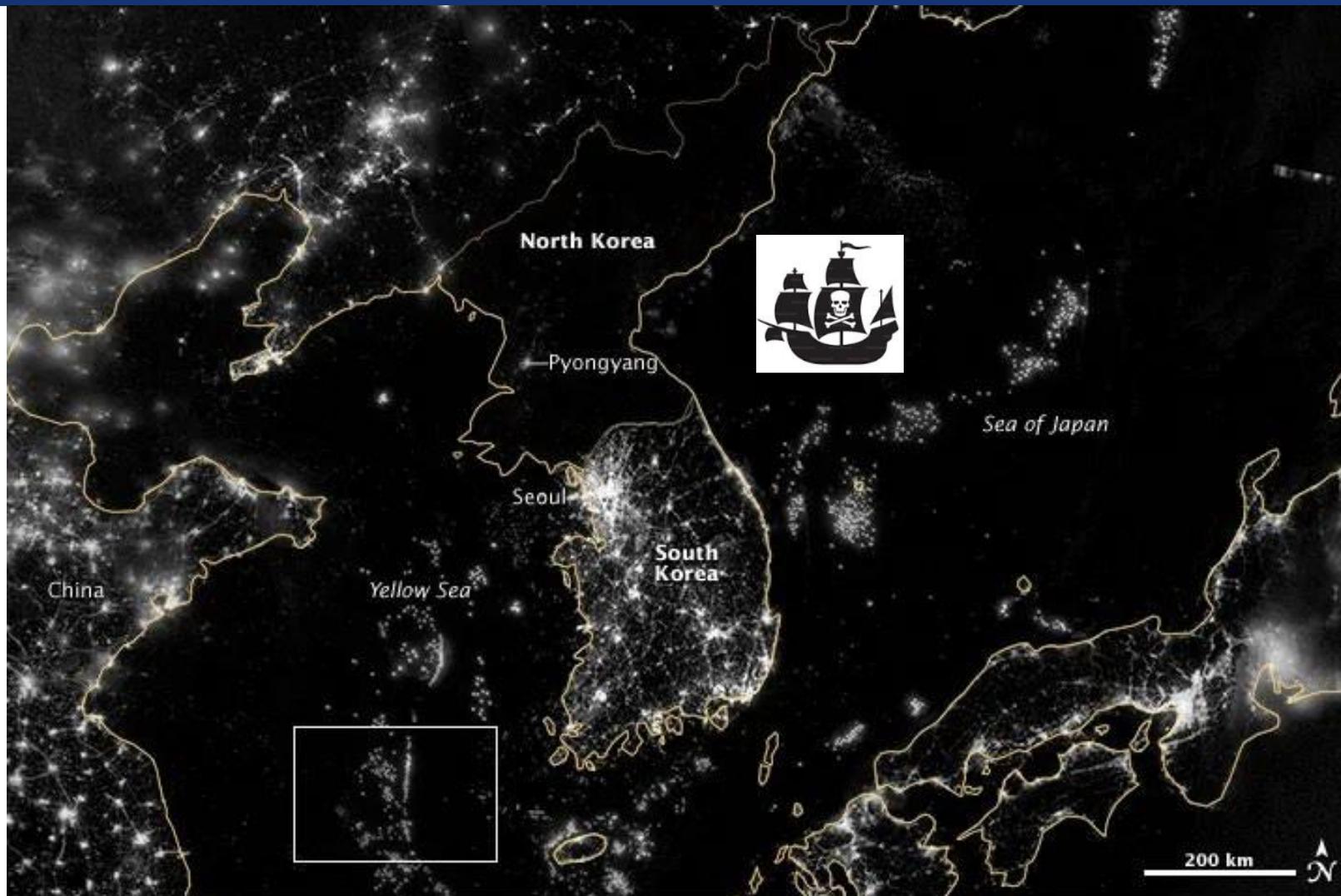
LL.M., JD, MBA, PSDGP, CIPP-US/G

DC Metro Area

Some mapping categories are already well established:

- Human/animal migration
 - Circulatory system
 - Energy usage/light emission
 - Traffic
 - Weather
 - Agriculture/forestry/mining
-
- It is often difficult to visualize PII, although we shed it everywhere we go.
 - Visualization aids resource allocation and other business decisions.
 - Pirate ship analogy.

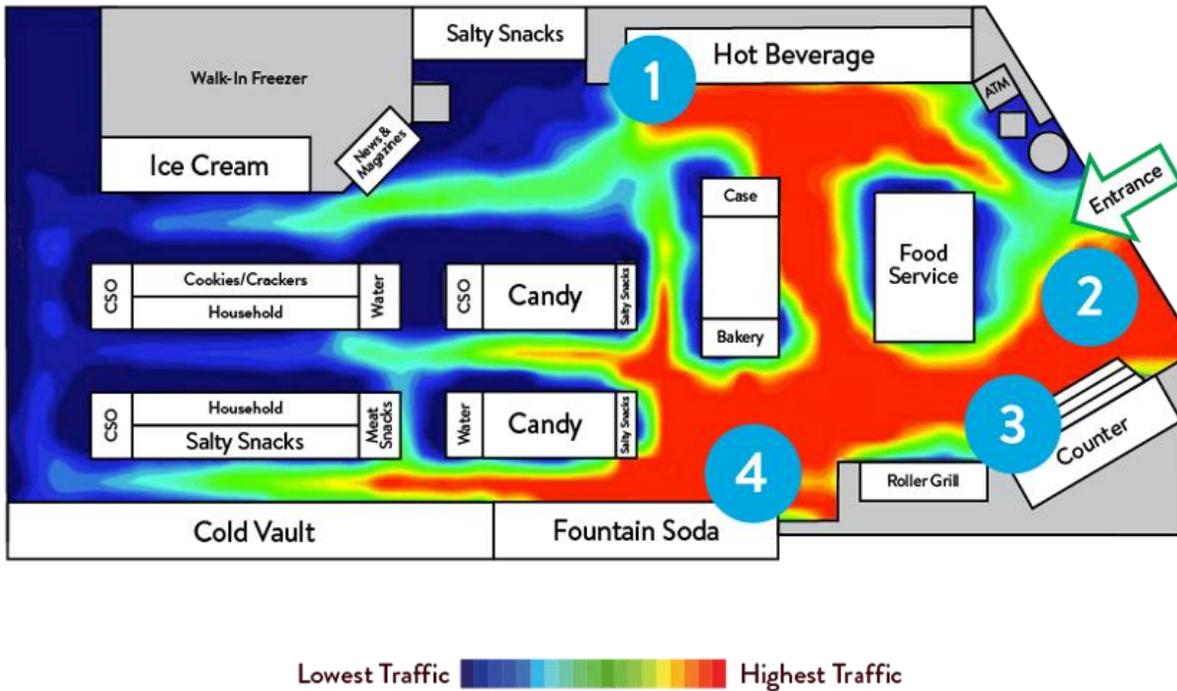
Setting Sail in Search of PII Gold



acquired September 24, 2012

Source: NASA

Foot and Auto Traffic



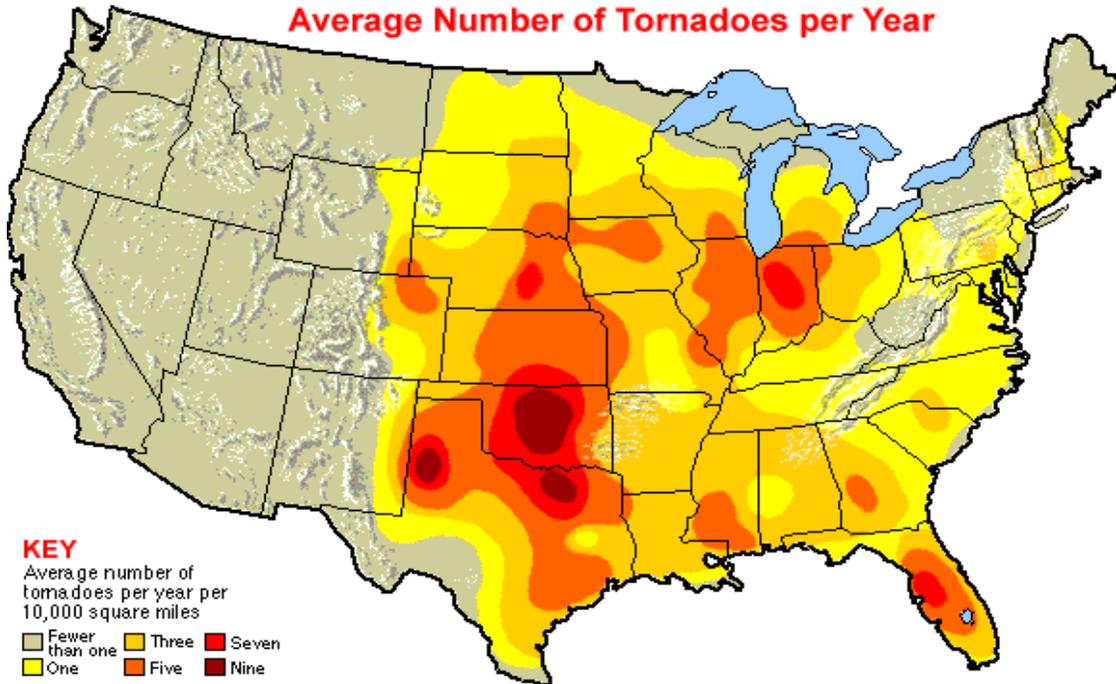
Source: skyfii.io



Source: Caliper.com

Tornado Heat Maps

Average Number of Tornadoes per Year



Copyright © 1998-1999 Oklahoma Climatological Survey. All Rights Reserved.

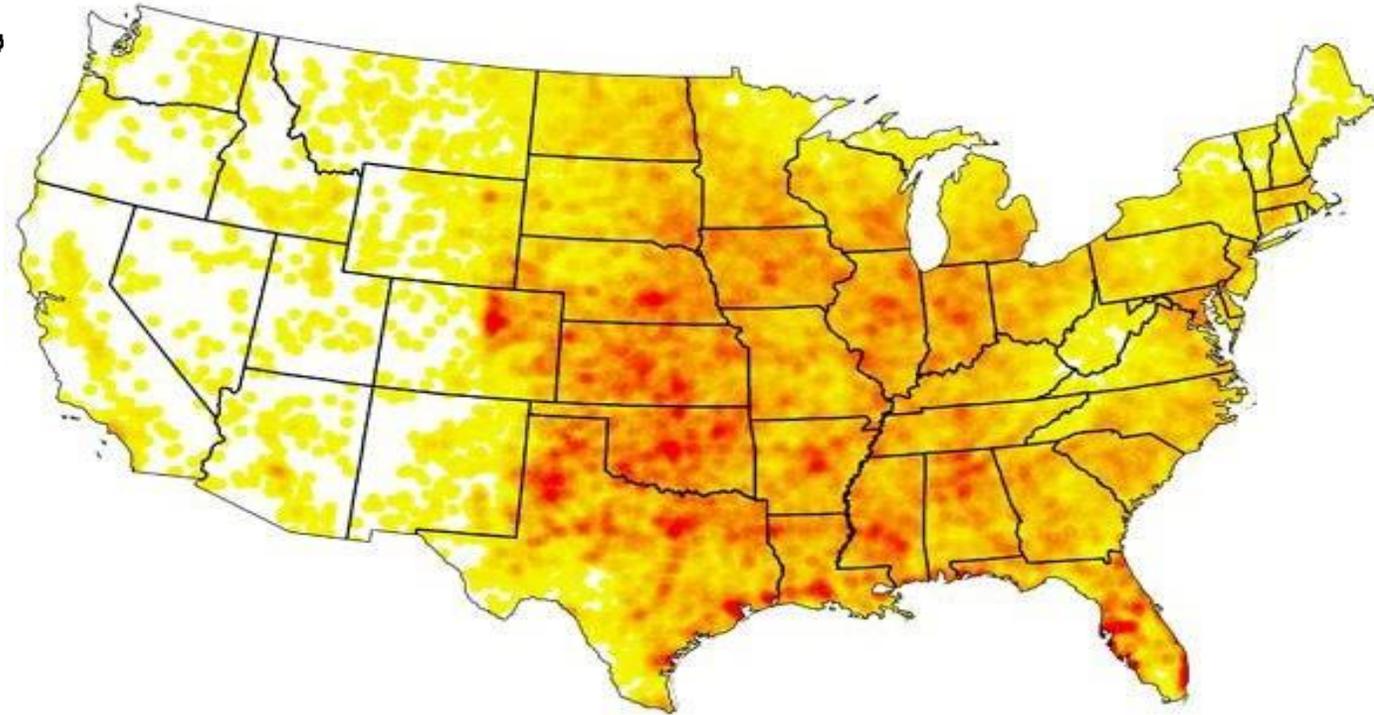
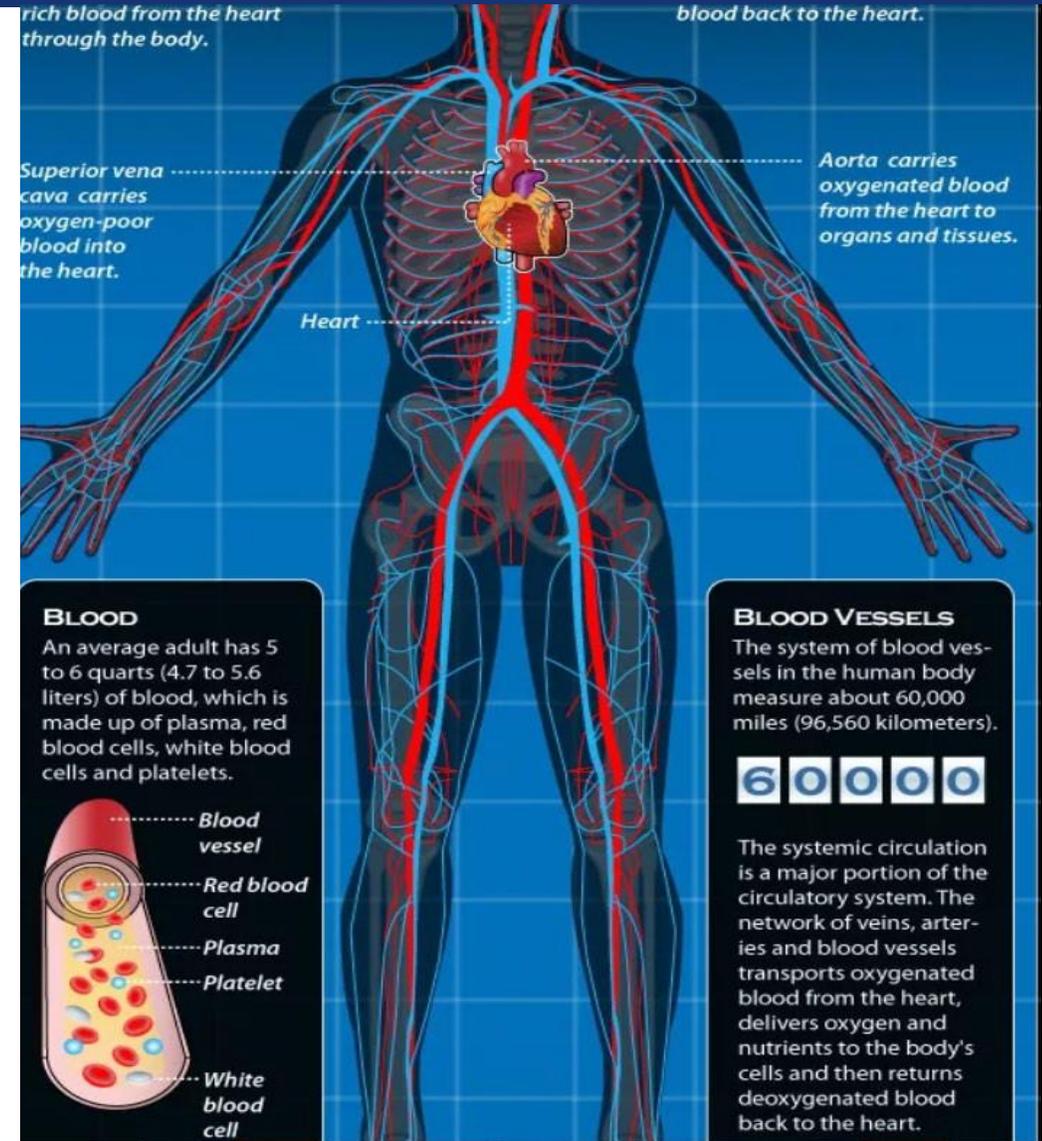
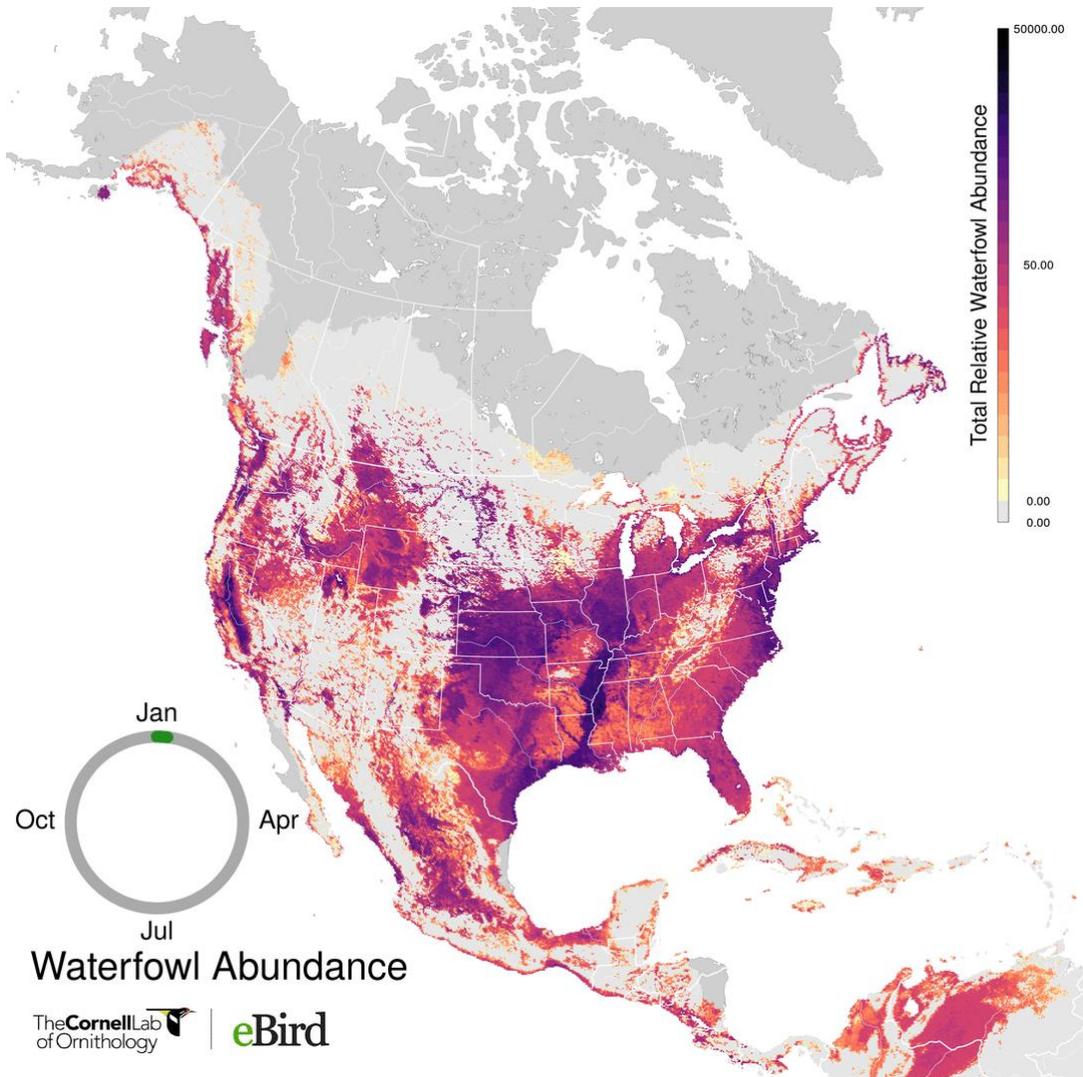


Figure 1. A heat map of tornado locations from 1950 to the present.

Migratory and Circulatory Pathways

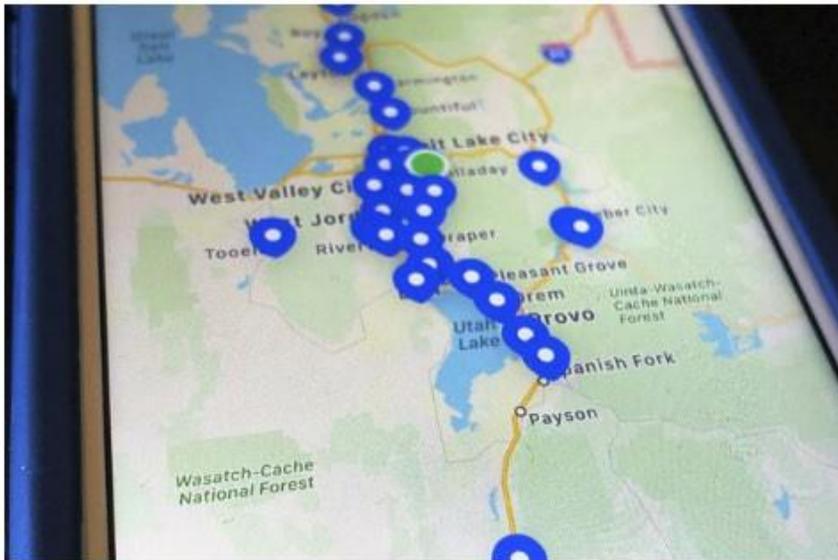


POLITICO

TECHNOLOGY

Homeland Security records show 'shocking' use of phone data, ACLU says

The civil liberties group released documents showing new details about how agencies had purchased information on people's movements throughout North America.



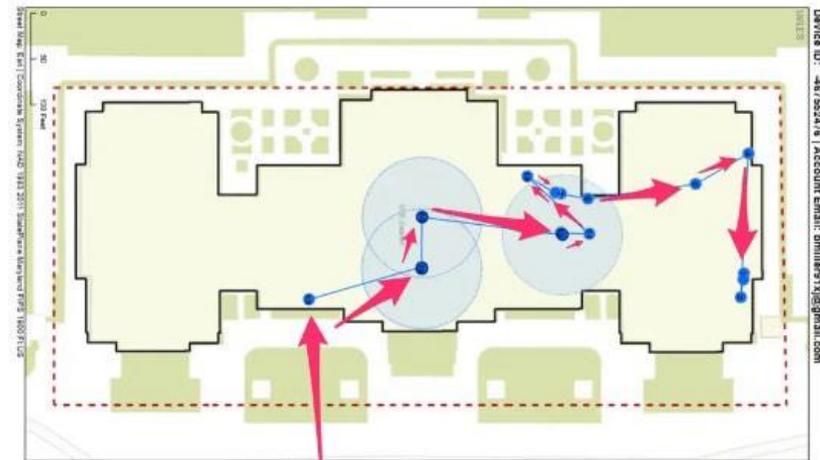
In just three days in 2018, documents show that the CBP collected data from more than 113,000 locations from phones in the Southwestern United States – equivalent to more than 26 data points per minute – without obtaining a warrant. | Lindsay Whitehurst/AP Photo

INSIDER

HOME > POLITICS

The DOJ is creating maps from subpoenaed cell phone data to identify rioters involved with the Capitol insurrection

Madison Hall Mar 24, 2021, 3:34 PM



- Radius \leq 100ft
- Radius $>$ 100ft
- Uncertainty Radius
- Capitol Building

Capitol rioter

Map of a Capitol insurrectionist created by the DOJ using subpoenaed cell phone data. US Department of Justice

Information Business Model – Then vs. Now

1974

- Public court filings
- Books and library records
- Less targeted advertising



Source: [Record Nations](#)

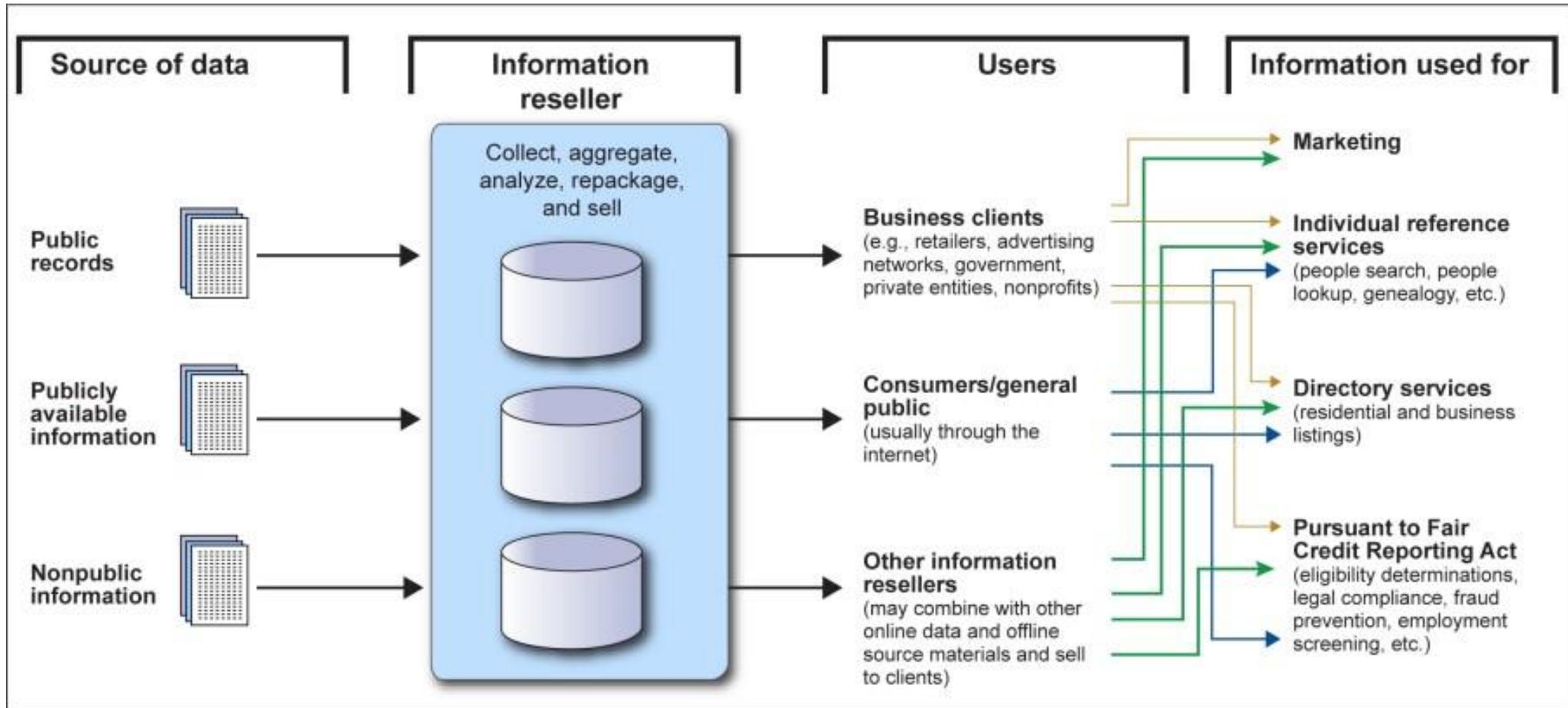
Today

- Internet browsing history
- Online/credit card purchasing history
- GPS location history
- Social media history
- Data aggregation
- Data brokers
- Subscription services
- Large and growing data market



Source: [Apple](#)

Consumer Data Flow



Source: GAO. | GAO-19-621T



“Information that has been published or broadcast for public consumption, is available on request to the public, is accessible **online or otherwise** to the public, is available to the public by **subscription or purchase**, could be seen or heard by a casual observer, is made available at a meeting open to the public, or is obtained by visiting a place or attending an event that is open to the public.”

Source: DoD Directive 3115.18 “DoD Access to and Use of Publicly Available Information (PAI)” (August 20, 2020)

Publicly Available Information

Factors to Consider in Properly Obtaining and Using Publicly Available Information:

1. Publicly Available – The information must be available to any member of the general public.

2. Lawfully Obtained – The information must be obtained lawfully under the Constitution and statutes of the United States, meaning, for purposes of this Guidance:

No special legal authorizations are needed, such as court orders or search warrants (if they're needed, then the information could still potentially be collected – it's just not "publicly available");

The collection technique is authorized/appropriate for IC professionals seeking publicly available information, rather than those IC professionals with specialized missions and authorities; and

The collection activity does not focus on a person solely because of race, ethnicity, national origin,

religion, or protected First Amendment freedom of speech or assembly.

3. Affiliation – Prior to obtaining information, the professional has identified and complied with obligations, if any, to disclose affiliation with the IC.

4. United States Person Information – If the collection includes information concerning United States persons, then:

There must be a *valid mission* requirement for the information under Executive Order 12333; and

The information must be retained and disseminated in accordance with *Executive Order 12333, the Privacy Act, and other applicable requirements.*

5. Accuracy – Safeguards are in place to ensure that the information is used in a manner that satisfies IC standards for information accuracy, quality, and reliability.



Source: LinkedIn

Source: ODNI Civil Liberties and Privacy Guidance for IC Professionals: "Properly Obtaining and Using Publicly Available Information" (Approved for public release in 2014)

- For PII data mapping, we have a long way to go...
- Progress can be gradual. Mapping continuous monitoring may occur prior to near-real time mapping.
- Imagine PII information maps of the future.
- We challenge you to take steps toward making these NextGen maps feasible for most organizations.

Identifying the challenge for our community. A historic challenge.

- Real-time (or near real time) is important because we must adapt to data moving faster and decisions being made more quickly. Enables us to pin-point where SPII is within the map.
- Continuous monitoring mapping may be a preliminary step prior to near real-time.
- Even after the initial map is produced, the map will require constant monitoring and updates. Real-time mapping helps to accomplish this, but how do we make the displayed map reflect the real-time PII?

Next Generation PII Data Mapping



Source: Interpark

Next Generation PII Data Mapping



Next Generation PII Data Mapping



Source: Desura

- Facilitates the required, healthy conversation between Project Team IT professionals and Privacy Team.
- May use “carrot” more rather than the “stick” if no firm legal obligations.
- Data mapping drives return on investment (ROI) which drives better decisions within an organization (to protect the data, to protect public image).
- Seek increased funding (through grants) for development of data mapping (real-time) for the Federal government.

How to Hold the Privacy Foes in Check?

Use an escalatory or graduated scale:

- Identify the foes through the mapping.
- Once the foes are identified, Privacy/Security professionals should counsel the foes.
 - Offer training and awareness, particularly for first time or incidental offenders.
 - Terminate access.
 - Refer to leadership for disciplinary action and/or law enforcement as appropriate.
- Shut down rogue databases. When foe (e.g., clown) activity is detected which includes “without permission” sites or rogue databases, shut down the unauthorized rogue sites/databases.

How to Hold the Privacy Foes in Check?

- Shut down rogue databases. When foe (e.g., clown) activity is detected which includes “without permission” sites or rogue databases, shut down the unauthorized rogue sites/databases.
- The four foes are our greatest challenge. Our greatest hope is the: farmers, miners, welders, and surveyors. The four "hope" skilled occupations build and give back to the community greater good.
- Root out the foes and strive to systematically replace the foes with the farmers, miners, welders, and surveyors.
- AI may be combined with data mapping to enable more efficient detection of foes.

Success – Data Mapping

- The first step in a larger process
- Identifies Risks
- Creates a Competitive Advantage



Data Mapping Benefits:

- Enables more informed decision making.
- Mitigates PII data breach risks
- Return on Investment (ROI)

The community should resolve to build a Next Generation PII Data Map which maps PII in near real time.



U.S. Federal Statutes

44 U.S.C. 3511 - Data Inventory and Federal Data Catalogue

The Privacy Act of 1974

E-Government Act of 2002

U.S. Policy & Guidance

Office of Management and Budget (OMB) Memoranda 7-16 and 13-13

NIST 800-53 App. J. Revision 4, SE-1: "Inventory of PII"

NIST 800-53 Revision 5: "PM-5(1): System Inventory/Inventory of PII"

International, State, Other

GDPR – Record of Processing Activities (ROPA)

ISO 27001 Asset Inventory

CCPA - Part of compliance process

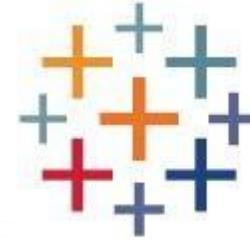
MITRE Privacy Engineering Tools

Backup Slides

Reserve Slides for Questions

 **DATA**GRAIL®

 **TIMi**



<https://ethyca.com/>



OneTrust
PRIVACY, SECURITY & GOVERNANCE

BigID Data Mapping



 **DATA**GRAIL®

Algolia



<https://sourceforge.net/software/product/Securiti/>

