

Removing Danger From Data

Kevin Scott, Senior Consultant @ CloverDX

Businesses face increasing risk around safety of data.

- Data is increasingly taking central role in many businesses.

... while at the same time

- Government regulations around safeguarding practices are rapidly emerging and are being actively enforced.



Rising Regulations

🟢 **European Regulation – GDPR**

Enforcement began in May 2018

Fines already levied against Google, British Airways, Marriot

🟢 **US Regulation status**

No national regulation

California CCPA regulation takes effect January 2020

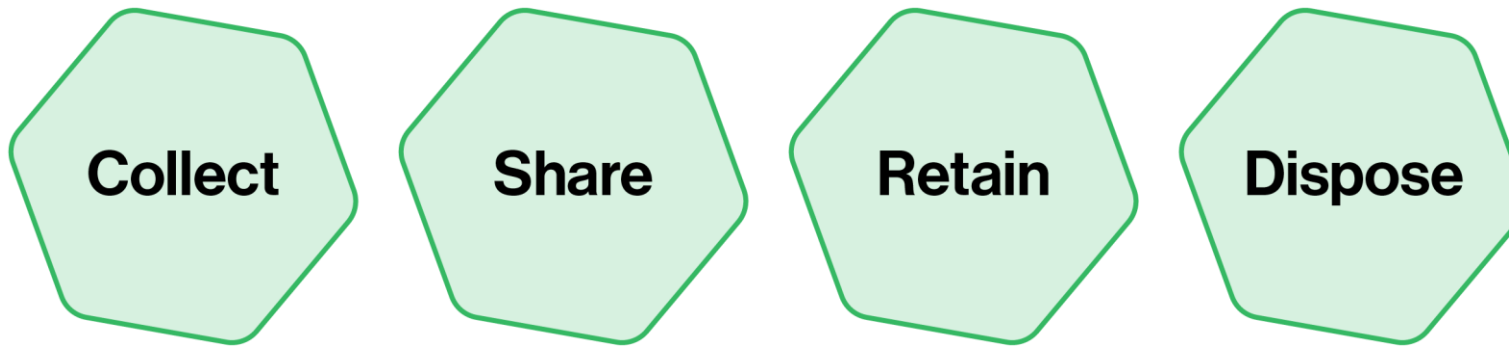
Microsoft will comply with CCPA nationwide

Other states following California model

It's not only PIs that needs safeguarding

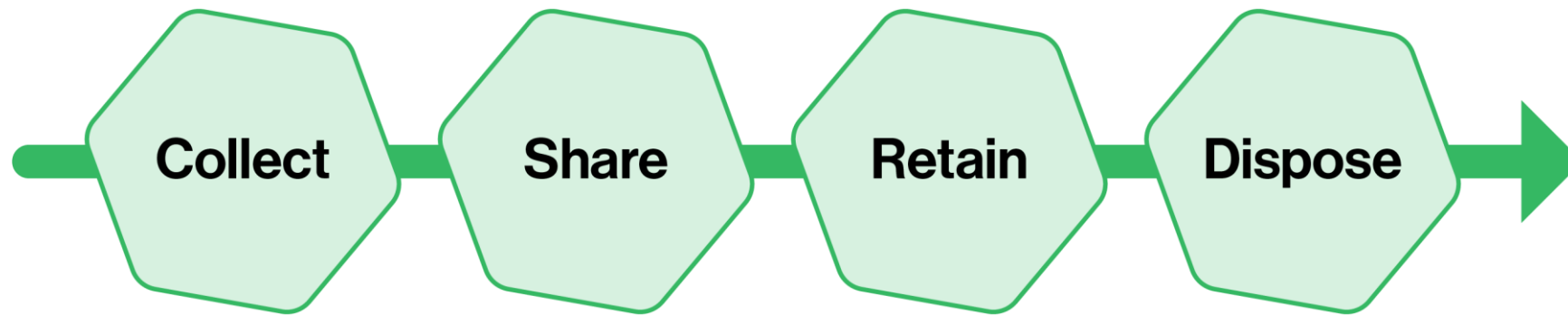
- **PII danger more obvious and includes potential legal costs**
name, birthdate, SSN, account number
- **Non-PII data is also potentially dangerous.**
sales forecasts, product plans, KPIs

Danger lurks in all phases of the Data Life Cycle



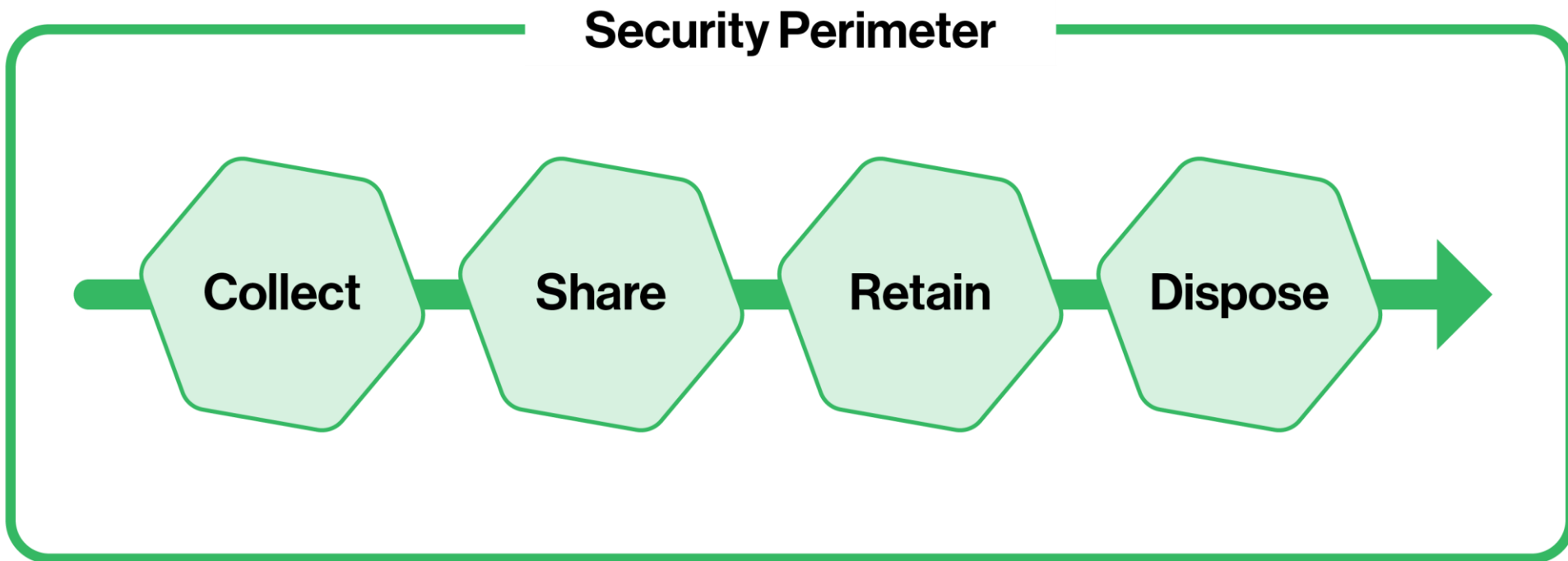
Danger lurks in all phases of the Data Life Cycle

Typically there's **processing** required to move data between stages

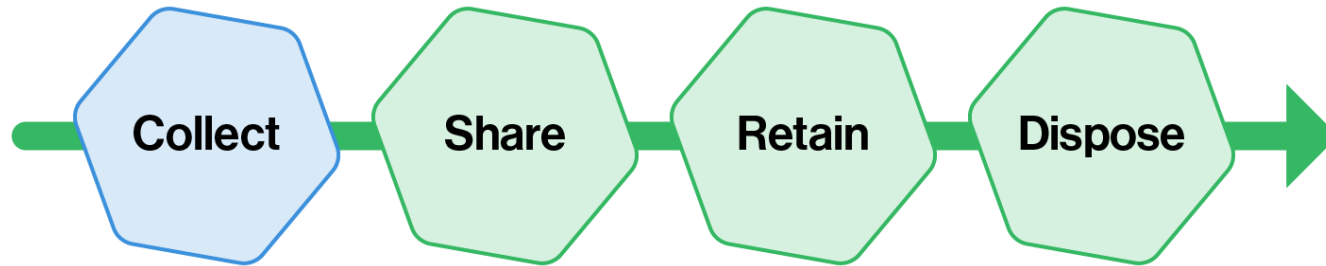


This **processing** can include **actions to remove danger**

Controlling **danger** from within



Collection Phase



- ◆ **Danger originates in collection phase**
- ◆ **Not knowing if PII data has been introduced into system**

Collection Phase:

Look for dangerous
data everywhere

expected



TIMELYCALL
Marketing Support System

Custom 1 Custom 2 Custom 3 Custom 4 Custom 5

• **Call Script**

• **Contact Summary**

Mr. Patrick LaJuett

Address: 123 West River Bend
Hilton, NY
14468

Primary #: 1-877-490-1200

Employer: Stamina Web Solutions

Title: Owner

Caller Time: 9:25 AM

Prospect Time: 9:25 AM Eastern

Best Contact Time: 9:00 AM Eastern

Calling: Hilton, NY

Dataset: IYP AOL 1_Initial Sales

List: New Leads DB

Script: 1_IYP Selling Script

• **Disposition / Call Outcome**

Sale

[Call Back](#) [Hang Up](#)

[Call History](#) [Done](#)

Answer Machine Busy Signal

Disconnected Dup. Phone #

Left Message with Gatekeeper No Answer

Privacy Manager Detected Wrong Number

• **Add Call Notes**

Submit

• **Response**

[Rebuttal](#) [Questions](#)

[Last Question](#) [Miscellaneous](#)

• **Other Links**

[Stats Database](#) [MapQuest Manager](#) [Custom 1](#) [Custom 2](#)

Previous Notes:

06-7-2018 at 10:32 am
by Agent 1667

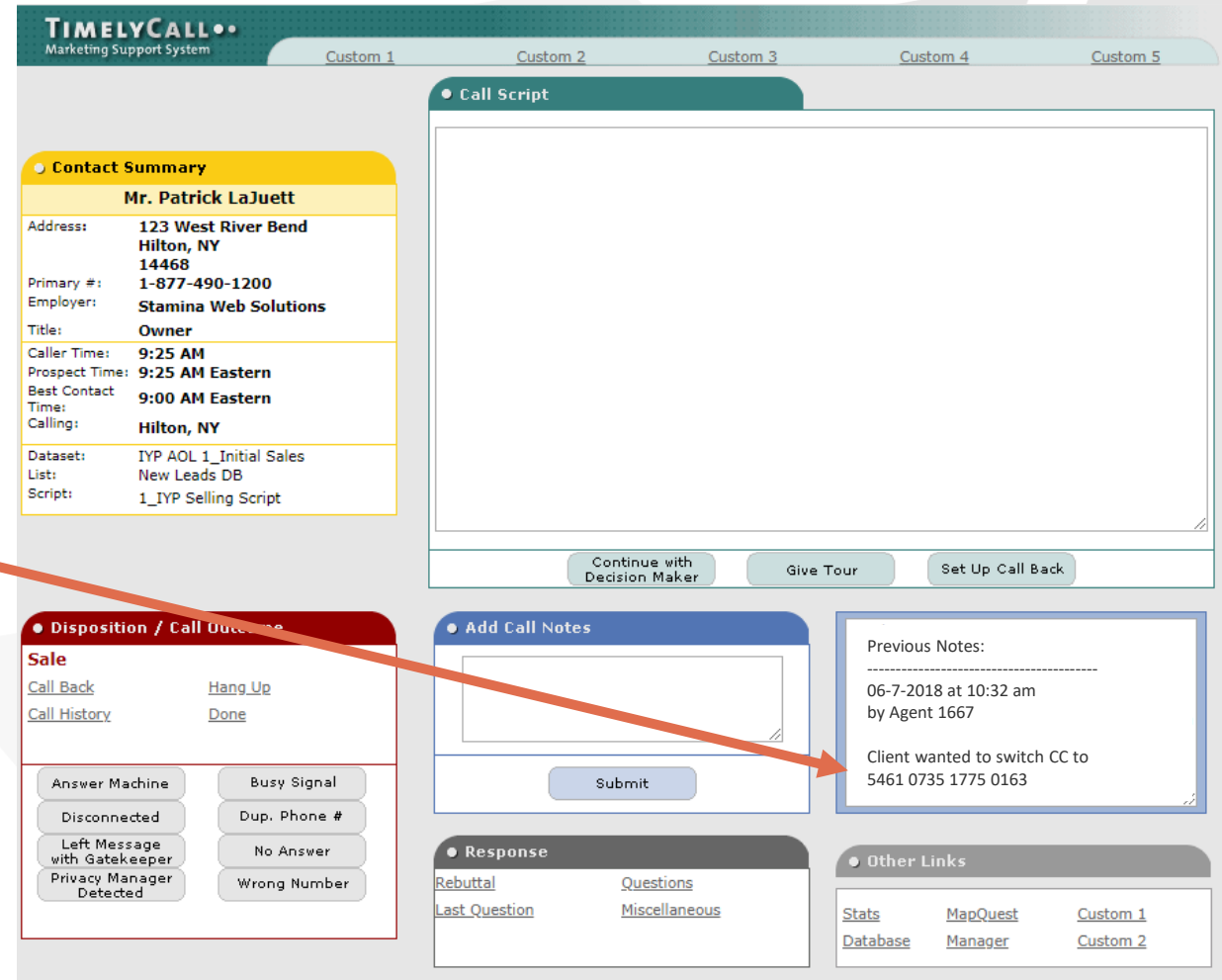
Client wanted to switch CC to
5461 0735 1775 0163

Continue with Decision Maker Give Tour Set Up Call Back

Collection Phase:

Look for dangerous
data everywhere

unexpected

The screenshot shows the TIMELYCALL Marketing Support System interface. It has a top navigation bar with tabs for Custom 1 through Custom 5. The main content area is divided into several sections: a yellow 'Contact Summary' box for Mr. Patrick LaJuett with contact details; a 'Disposition / Call Outcome' section with a 'Sale' status and various call status buttons; a 'Call Script' section with a large text area and buttons for 'Continue with Decision Maker', 'Give Tour', and 'Set Up Call Back'; an 'Add Call Notes' section with a text area and a 'Submit' button; a 'Response' section with links for 'Rebuttal', 'Questions', 'Last Question', and 'Miscellaneous'; and an 'Other Links' section with links for 'Stats Database', 'MapQuest Manager', and 'Custom 1 Custom 2'. An orange arrow points from the word 'unexpected' to the 'Add Call Notes' section, specifically highlighting the 'Previous Notes' area which contains a timestamp and a note about a client wanting to switch CC.

Collection Phase

Remedies:

- Data minimization
- Choose ETL tools that provide **input classification**
- Consider Software to **catalog and tag collected data**

Collection Phase:

- Data Classification tools can help track incoming sensitive data

GDPR Check Tool

Pending deletions

Value search

John Doe (settings)

Logout

Home / Reports / Data Map

Data map - list of tables and related domains

Records per page10

Export

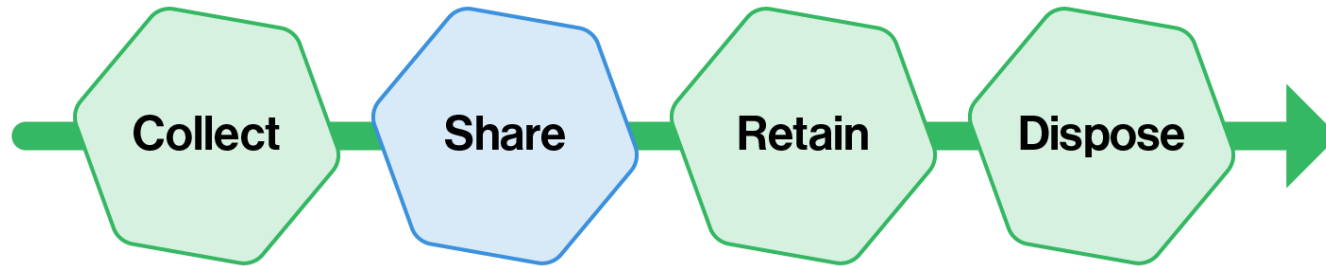
Database	Schema	Table	Column	Domains	Processes	Sample data	
Q	Q	Q	Q		Q	Q	Remove filter
firmdb	erpp	ta_person	name	First nameLast nameAddress - Street	P016P101	alex, andrej, daniel ...	Actions
firmdb	erpp	ta_person	name_alt	Last nameFirst name	P092	alex, andrej, daniel ...	Actions
firmdb	erpp	ta_person	last_name	First nameLast name	P023P101	adam, adamek, barton ...	Actions
firmdb	erpp	ta_purchord	last_name	Last nameAddress - Street	P092P101	bláha, borovička, cín ...	Actions
firmdb	erpp	ta_purchord	addrtxt	Address - StreetAddress - CityAddress - ZIP code	P016P023	brno, bayerova, dlouhá ...	Actions
firmdb	erpp	ta_purchord	addrcontxt	Address - StreetAddress - CityAddress - ZIP code	P023	60200, brno, bezručova ...	Actions
firmdb	erpp	ta_sup_contact	cvalue	E-mailLast nameFirst nameAddress - City	P016P023P092P101	aberk@firma.cz, adamec@jinafirma.cz, bures.jiri@oknaprovas.cz	Actions
firmdb	erpp	ta_supplier	city	Address - City	P016	beroun, brno, chrudim ...	Actions
firmdb	erpp	ta_supplier	txtDisplayName	First nameLast nameAddress - City	P023P092P101	adam, adamec, beron ...	Actions
firmdb	erpp	ta_supplier	conperson	Last nameFirst name	P016P023P092P101	andrej, ber, cilek ...	Actions

Showing 1 - 10 of 2160 results.

FirstPrevious12345...216NextLast

Data Classification tags can be stored and accessed in a Data Catalog

Share phase:



- **Sharing usually requires transforming data assets so they can provide desired business value.**
searching, mapping, sorting, merging categorizing, analyzing, reshaping.
- **Primary danger is oversharing**
Sharing too much data
Sharing to an unnecessarily large audience
Accidental sharing

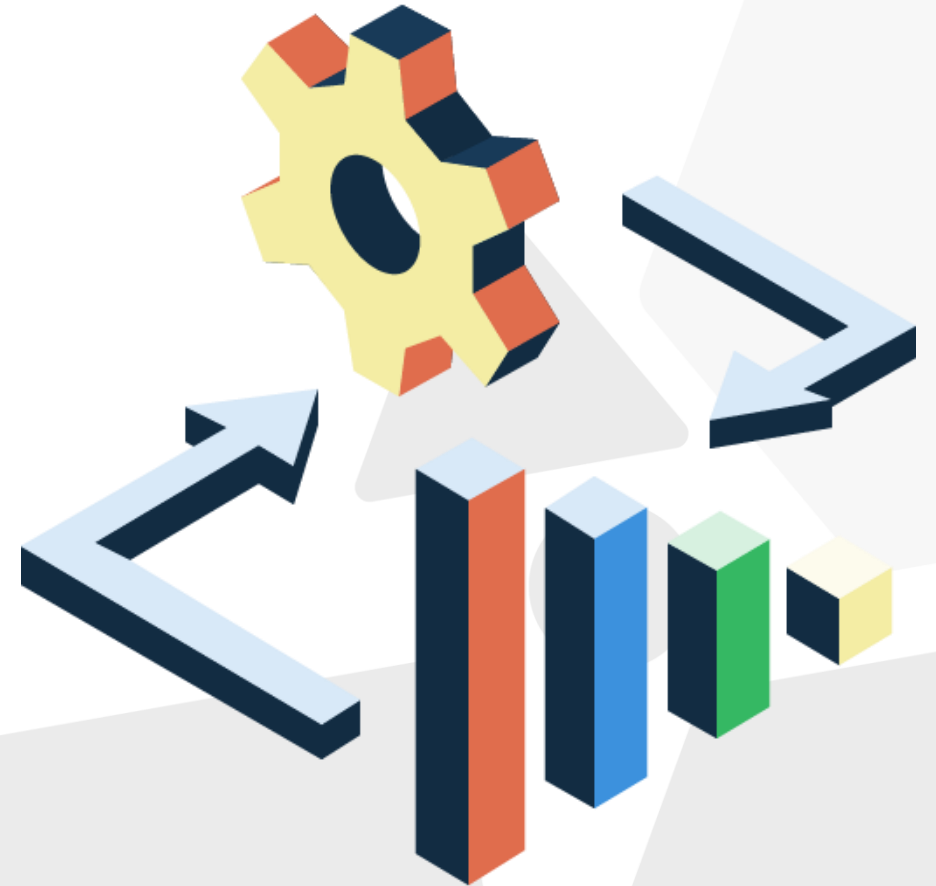
Share phase:

Processing techniques available to deliver business value without unnecessarily sharing sensitive data.

- Anonymization
- Privacy preserving analytics

Anonymization

Process of obfuscating data, so they keep some of their original attributes but not to extent they could be used to infer relation to real people or entities.



Anonymization – Shuffle

First Name	Surname	Age		First Name	Surname	Age
Alice	Smith	42		Alice	Doe	42
Bob	Johnson	21		Bob	Jackson	21
Dave	Doe	74		Dave	Chang	74
Eve	Jackson	44		Eve	Smith	44
Grace	Chang	32		Grace	Johnson	32

Anonymization – Jitter

Birth Date	Gender	Education Level		Birth Date	Gender	Education Level
7-Jul-82	Male	High School	+3 days	10-Jul-82	Male	High School
18-May-85	Male	High School	-1 day	17-May-85	Male	High School
5-Mar-87	Female	Bachelor	+3 days	8-Mar-87	Female	Bachelor
17-Jun-97	Male	Associate	-1 day	16-Jun-97	Male	Associate
27-Feb-82	Male	Graduate	-3 days	24-Feb-82	Male	Graduate
9-Nov-58	Male	Graduate	+1 day	10-Nov-58	Male	Graduate

Anonymization – Masking

Naively masked

4024 XXXX XXXX XXXX

Intelligently masked

4024 0071 4314 0399

Keeping Issuer code
VISA Credit card
Issued by Bank of America

Randomized
Account Number

Valid Luhn checksum

Preserves card types, issuers, preserves validity

Anonymization – Aggregation/Generalization

Birth Date	Gender	Education Level		Age	Gender	Education Level
7-Jul-82	Male	High School		30's	Male	High School
18-May-85	Male	High School		30's	Male	High School
5-Mar-87	Female	Bachelor		30's	Female	Bachelor
17-Jun-97	Male	Associate		20's	Male	Associate
27-Feb-82	Male	Graduate		30's	Male	Graduate
9-Nov-58	Male	Graduate		60's	Male	Graduate

Privacy Preserving Analytics

- **Methods developed at Linked-In for data mining on very large (web-scale) result sets**
- **Provide accurate analytics without inadvertently identifying individual users.**

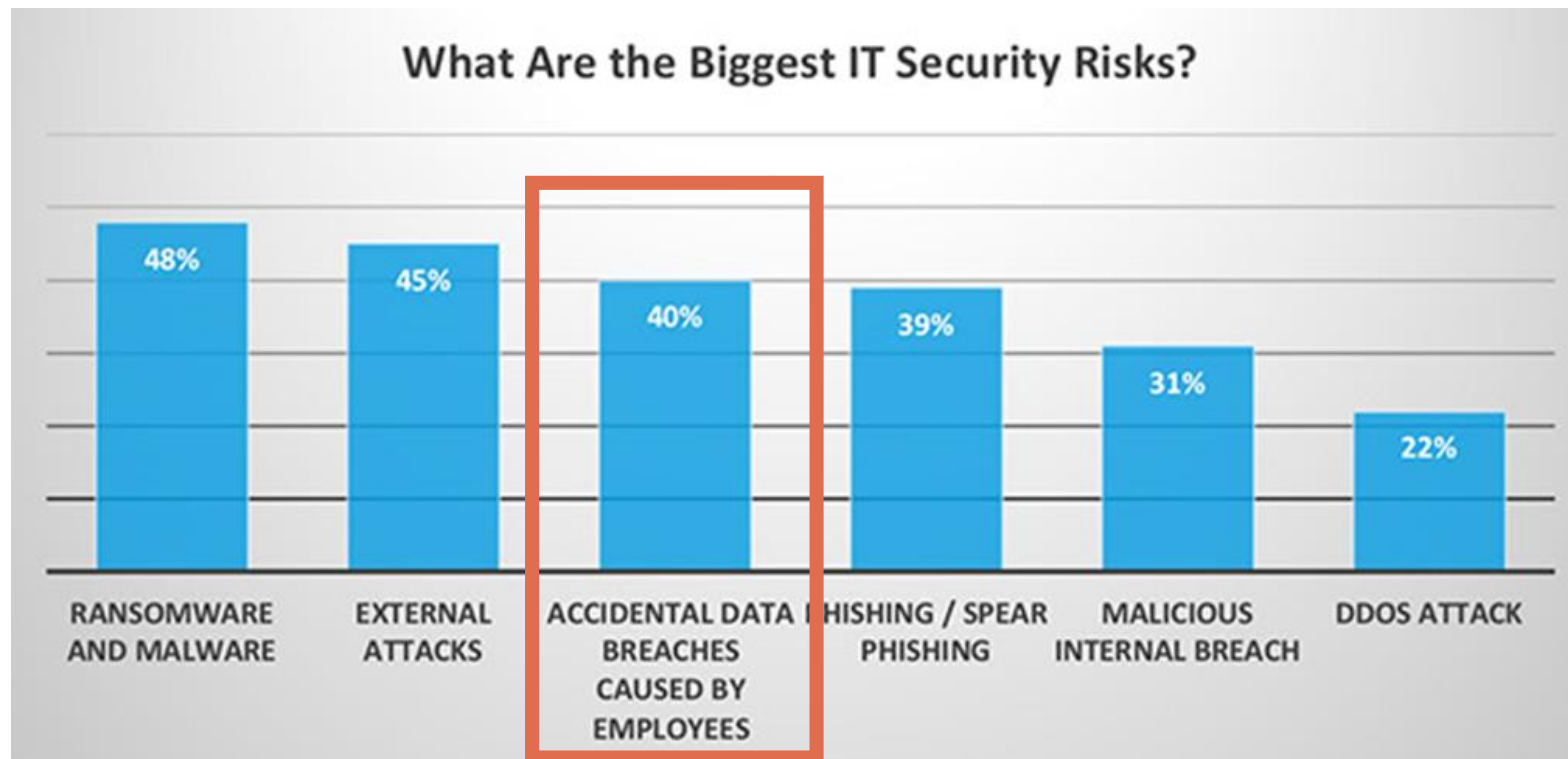
Inject small amounts of random noise to query/search

Offset noise with more processing to maintain data consistency

See [Privacy-Preserving Analytics and Reporting at LinkedIn](#)

Share Phase

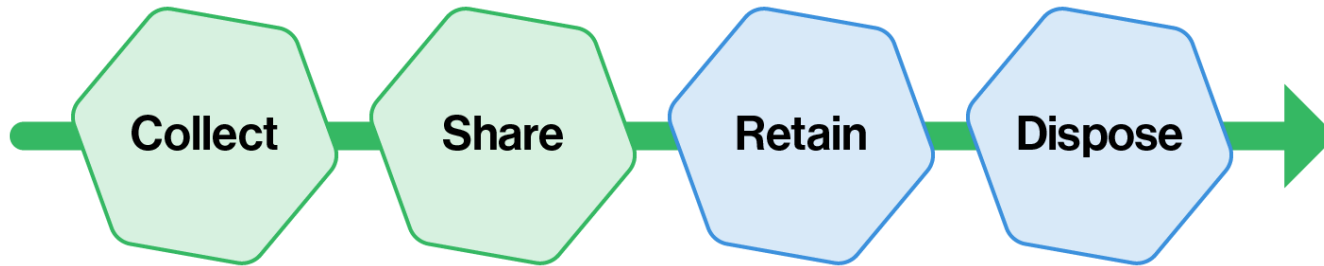
- Employees simply doing their job can be dangerous.



Share Phase:

- **5 most common technologies that lead to accidental data breaches by employees:**
 - 1. External email services (Gmail, Yahoo!, etc.) (51 percent)**
 - 2. Corporate email (46 percent)**
 - 3. File sharing services (FTP sites, etc.) (40 percent)**
 - 4. Collaboration tools (Slack, Dropbox, etc.) (38 percent)**
 - 5. SMS / messaging apps (G-Chat, WhatsApp, etc.) (35 percent)**

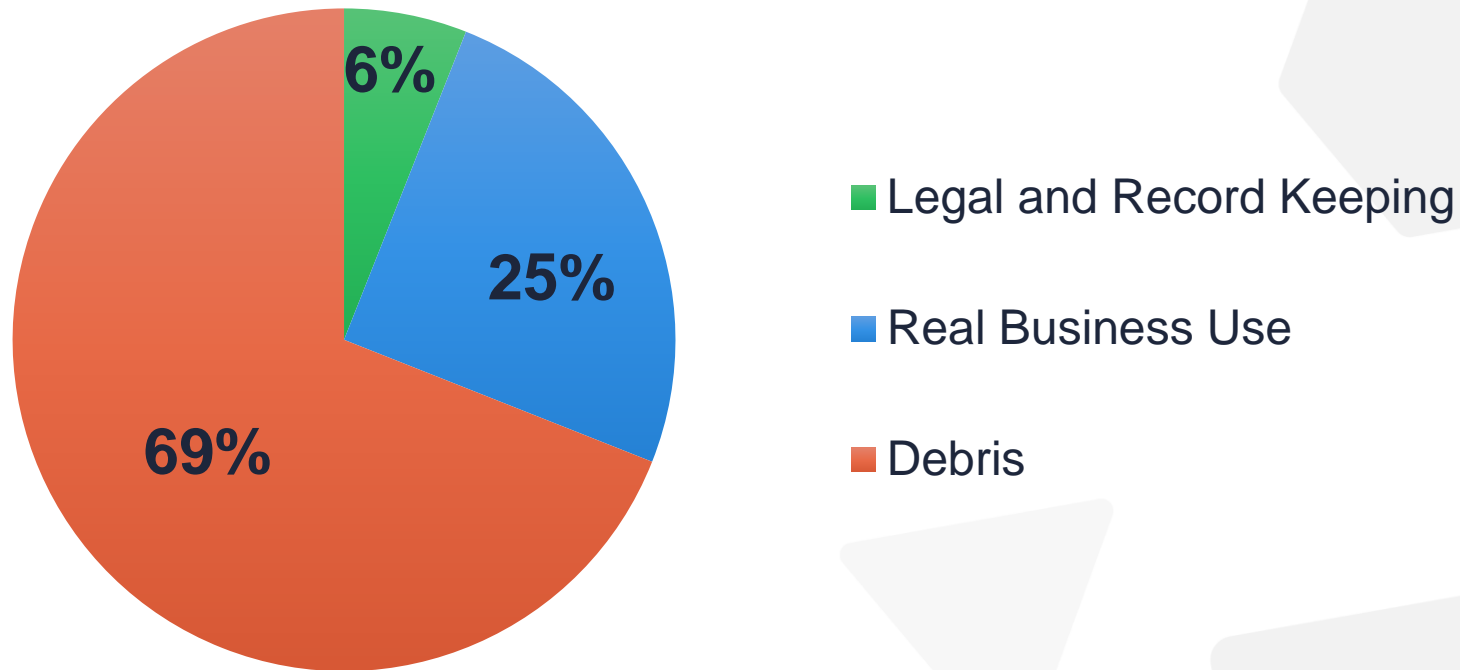
Retain and Dispose Phases



- Minimize time data is kept
- A Data Catalog can also house retention and disposal status

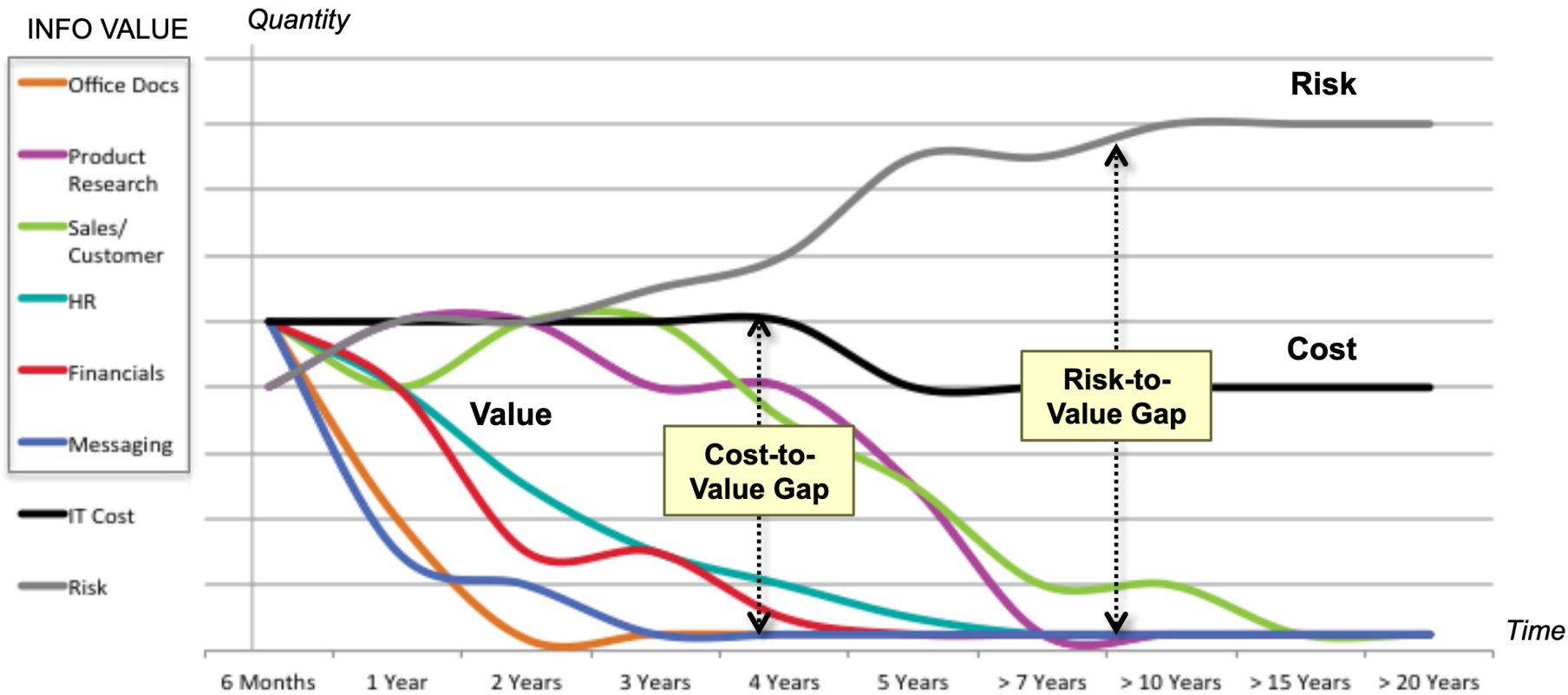
Most enterprise data is retained indiscriminately

Enterprise Data Retention



https://cedar.princeton.edu/sites/cedar/files/media/information_lifecycle_governance.pdf

Information Value Declines Over Time, Costs and Risks of retention do not.



https://cedar.princeton.edu/sites/cedar/files/media/information_lifecycle_governance.pdf

Retain and Dispose Phases


In 2019 CapitalOne Breach exposed sensitive data that included rejected credit card applications from as far back as 2005.


Remediation expected to cost \$200-300 million *without any explicit regulatory fines.*





Retain and Dispose Phases

GDPR Check Tool

 Pending deletions



 Value search

 John Doe (settings)


 Logout

Home / Reports / Pending deletions






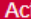
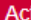

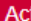
Retentions

Records per page 5   Export

Search for a person

Michael Jordan

Showing retention records for name: **Michael Jordan**

System	State	Applicable date	Comment	
				Remove filter
SalesForce CRM	Deleted	2018-01-07	Deletion requested on 20th June 2017	
AR Data	Pending deletion	2020-04-12	Legal requirement - invoices. Deletion requested on 2017-12-29 June 2017	
ERPSource	Pending deletion	2018-01-07	Legal requirement - invoices. Deletion requested on 2017-12-12 June 2017	
SugarCRM	Deleted	2018-01-07	Deletion requested on 12th December 2017	
Sales Tools	Deleted	2018-01-07	Deletion requested on 12th December 2017	

Showing 1 - 5 of 12 results.

First

Previous

1

2

3

Next

Last

Summary

- Emerging Regulations are increasing danger in data

- Removing Danger = Identification + Remediation**

Realizing the dangers exist
Software tools exist to help

- Develop Risk Intelligence**

Size of risk vs cost of remedy
Implement reasonable precautions



Thank you

About CloverDX Data Integration Platform

CloverDX is a data integration platform for designing, automating and operating data jobs at scale. We've engineered CloverDX to solve complex data movement and transformation scenarios with a combination of visual IDE for data jobs, flexibility of coding and extensible automation and orchestration features.

www.cloverdx.com

References

◆ Data minimization

<https://www.dataguise.com/gdpr-compliance-data-minimization-use-purpose/>

◆ Regulation Costs

<https://www.cnbc.com/2019/07/10/gdpr-fines-vs-marriott-british-air-are-a-warning-for-google-facebook.html>

◆ Capital One (Data EOL)

<https://www.theverge.com/2019/7/31/20748886/capital-one-breach-hack-thompson-security-data>

<https://www.nytimes.com/2019/07/29/business/capital-one-data-breach-hacked.html>

◆ Princeton Lifecycle Governance

https://cedar.princeton.edu/sites/cedar/files/media/information_lifecycle_governance.pdf