



**AlignConsulting**

It's 5 o'clock. Do  
you know where  
your info risks are?

Kate Pugh, Columbia, AlignConsulting

Sandra Serkes, Columbia, Valora Technologies

KM World November 5, 2019

# We love love our furry friends...but we *really* love metadata



Sandra Serkes, CEO, Valora Technologies  
Guest Lecturer, Columbia University  
Information and Knowledge Strategy  
Master's Program



Katrina Pugh, CEO, AlignConsulting  
Faculty and former Academic Director,  
Columbia University Information and  
Knowledge Strategy Master's Program

# What is content bloat costing us?

1 hour, 41 minutes/week spent searching  
(average worker)

13.2 billion hours across all US workers

\$628 billion opportunity  
cost for the US

Assumes 2000 hr/year/worker, 158M workers, \$23.63/hr for wages and overhead, 4.2% lost productivity searching for content. Sources: Nucleus Research "Content Bloat Drains Productivity by 8 percent," June, 2016. Trading Economics US hourly wages. <https://tradingeconomics.com/united-states/wages>

# What could info bloat cost us?



Breach of  
customer data



Proposed fine:  
**£183M**



Failing to protect  
customer data



Proposed fine:  
**£99M**

GDPR introduced hefty penalties per instance:  
up to **€20M** or **4% of global turnover** – *whichever is higher*

# Why else do we “take control of” unstructured data?

## **Compliance:**

Retention, Legal Hold, GDPR, Privacy, Security

## **Identification:**

ROT, Dupes, PII, DocType, Relevance, Privilege, Custodian, Date, IP potential and more

## **Disposition:**

Quarantine, Defensible Deletion, Redaction, Migration, Remediation

Do any of  
these look  
like you?



## Large Oil & Gas Manufacturer

- Frequent acquisitions of other companies
- Millions of documents
- Not enough SMEs
- Risky content..Legal hold...Data Privacy... ROIT



## Multi-store, multi-country fast-food

- Independent, franchise & retail sites
- Managers autonomously hire, fire
- GDPR concerns



## Info hound (bankruptcy data vault)

- Out run by content volume
- Regulatory reporting obligations
- Sensitive data
- Missing value-added revenue...?

# Laying the groundwork for today

## 1. Metadata: Frame

Metadata tells us:

- What it is (e.g., contract, whitepaper)
- What it is about (e.g. topics, insights)
- Inherited labels (e.g., file type, date, author, last modified date)
- Impact, risk

*Designing, contextualizing*

## 2. Autoclassification: Analyze and Design

Unitize, translate, transform, analyze, assess, triage groups of documents, based on metadata and context

*Processing, tagging*

## 3. Info Gov: Control

Structures, policies, procedures, processes and controls implemented to help manage information at an enterprise level over time

*Planning for future*

# 1. Metadata: Framing



# Rich Metadata for Enterprise Content

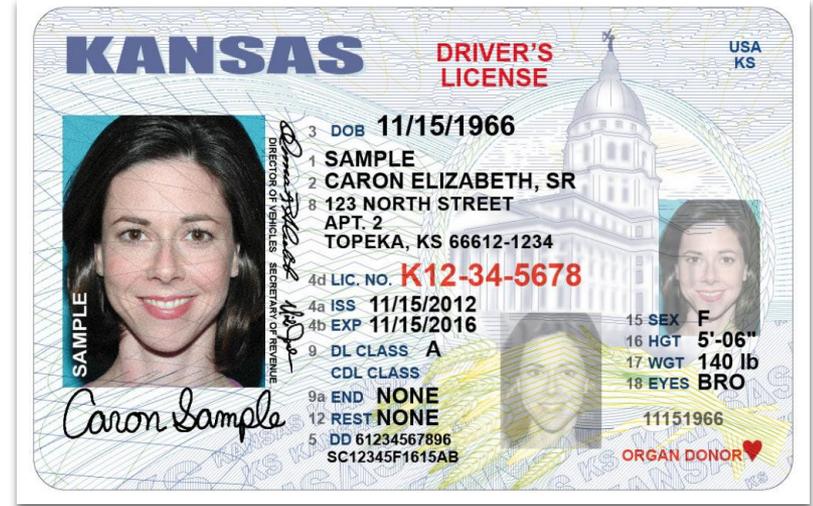


What  
attributes &  
tags make  
sense for your  
organization?



+ "hidden" metadata

# Metadata is “portable”

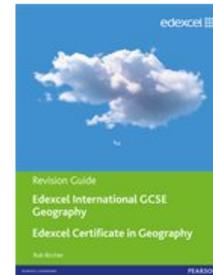


But what about similarly rich metadata (tags) for our enterprise content?

# Isn't rich metadata there...or can't we just add it?

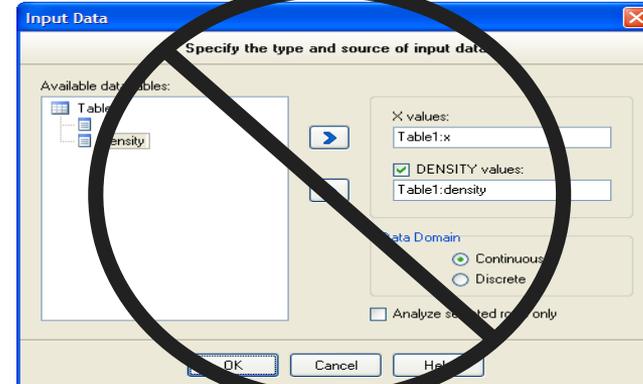
## Sort of..

- Rich metadata tends to be common with structured data (databases), and uncommon with unstructured data (email, slide decks, reports, etc.)
- Rich metadata is typically input by hand (via forms & wizards) and suffers from negligence, inconsistency, laziness, error, etc.



[See larger version of cover](#)

Price	£7.99 + £0.16 UK VAT
ISBN	9781446905777
Publication Date	March 2013
Format	BOOK
Jump to:	<a href="#">Customer reviews</a>



# What does rich metadata get you?



## Ability to know what you have, where it is, and how to manage it

- Lifecycle management tags
- Protected status tags (legal hold, privacy, access, etc.)
- Context for searches & decision making



## Ability to see across data repositories

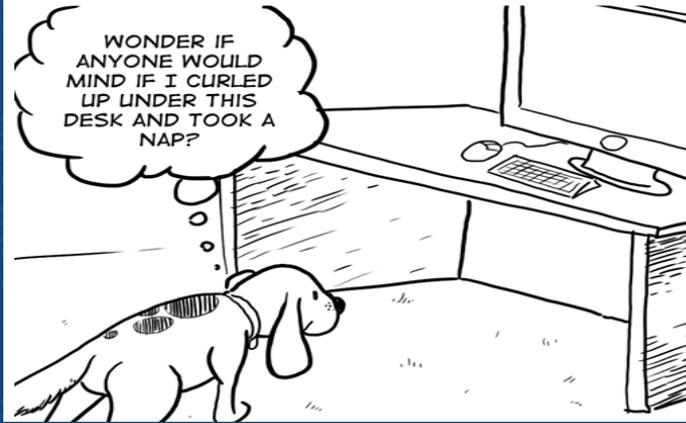
- Normalized tags
- Data visualization, presentation and reporting



## Ability to operate by facet

- Manage by content, not by time or other “simplistic” approaches

“Wonder if anyone would mind...?”



Context  
matters as  
we define  
metadata!

# The slow way: Hand classifying Cuban Missile Crisis transcript

Ingestion (OCR) - Automated

Coding – 1 Hour/page (= 84 hours)

Preliminary analytics (= 1 hr. Messy, not repeatable)

Meeting Item 40.3, Reference Reel 1

Cuba Tape

(This is a transcription of tapes recorded in the Saturday, October 27, 1962. The tape begins with an inconclusive discussion of questions, including plans to stop a ship (Grozny) and daylight surveillance missions, morning and afternoon interrupted a few minutes into the meeting as the President:)

JFK: (Reading) Premier Khrushchev told President Kennedy would withdraw offensive missiles from Cuba if the United States would withdraw its rockets from Turkey. (voices unclear.)

Voice: He didn't really say that, did he?

JFK: That may not be -- he may be putting out a statement. (Mixed voices. Calls for Pierre [Salinger].)

JFK: That wasn't in the letter we received, was it?

Voice: No. (Voices unclear)

JFK: Is he supposed to be putting out a letter he's putting out a statement?

Salinger: Putting out a letter he wrote to you.

JFK: Let's just -- uh -- keep on going (words unclear)

Voice: It's in a different statement.

Rusk: Well, I think we better get -- uh -- (words unclear) check and be sure that the letter that's coming in on the ticker is the letter that we were seeing last night. (mixed voices)

JFK: What's the advantage of the second mission?

McNamara: It creates a pattern of increasing intensity for the President. We believe that we should do this. Now, I don't recommend, although, we don't need...

JFK: What's the advantage of the second mission?

McNamara: It creates a pattern of increasing intensity for the President. We believe that we should do this. Now, I don't recommend, although, we don't need...

Participant	Reply	Coding
JFK:	(Reading) Premier Khrushchev told President Kennedy yesterday he would withdraw offensive missiles from Cuba	[Reading]
	(voices unclear.)	
Voice:	He didn't really say that, did he?	*Anti-Integrity*
JFK:	That may not be -- he may be putting out a statement. (Mixed voices. Calls for Pierre [Salinger].)	
JFK:	That wasn't in the letter we received, was it?	
Voice:	No. (Voices unclear)	
JFK:	Is he supposed to be putting out a letter he's putting out a statement?	
Salinger:	Putting out a letter he wrote to you.	
JFK:	Let's just -- uh -- keep on going (words unclear)	
Voice:	It's in a different statement.	
Rusk:	Well, I think we better get -- uh -- (words unclear) check and be sure that the letter that's coming in on the ticker is the letter that we were seeing last night. (mixed voices)	
JFK:	What's the advantage of the second mission?	
McNamara:	It creates a pattern of increasing intensity for the President. We believe that we should do this. Now, I don't recommend, although, we don't need...	

Row Labels	(truncated)	question mark	*anti-courtesy*	(interrupt)	*anti-integrity*	in disguise	*Anti-Integrity* (this is a statement in...)	*courtesy (Agreeing/Inclusion)	compliance	(clarification)	*integrity Q*	*Integrity-Q*	*Integrity-Translat on*	(Reading)	*anti-courtesy* (interrupt)	*Integ
Ball:											4					
Bundy:			1					1			12		4			
Bundy:											1					
Bundy:															1	
Bundy:			1													
Dillon:												1				
Dillon:													1			
JFK:											1		2			
JFK:			1					1	1		17		4			4
JFK:			1								2		1			
Low:																
McCon:																
McNamara:			1										16			
Nitze:													5			1
Particip:																
Rusk:											9		1			1
Rusk:			1								10		3			
Rusk:			1								2					
Salinger:																

<https://catalog.archives.gov/id/193723>

# AutoClassifying an attachment (patent application)

(9) **United States**  
 (12) **Patent Application Publication**  
 Mousseau et al.  
 (10) **Pub. No.: US 2007/0242809 A1**  
 (4) **Pub. Date: Oct 18 2007**  
 (54) **ADVANCED VOICE AND DATA OPERATIONS IN A MOBILE DATA COMMUNICATION DEVICE**  
 continuation-in-part of application No. 10/095,603, filed on Mar. 11, 2002.  
 (60) Provisional application No. 60/274,408, filed on Mar. 9, 2001.  
 (75) Inventors: **Gary Mousseau, Waterloo (CA); Mike Lazaridis, Waterloo (CA); David Yach, Waterloo (CA); Raymond Vander Weem, Waterloo (CA); Harry Major, Waterloo (CA); Atul Asthana, Waterloo (CA)**  
**Publication Classification**  
 (51) **Int. Cl. H04M 11/00**  
 (52) **U.S. CL. 379/88.16**  
 (57) **ABSTRACT**  
 A system and method for integrating voice and data operations into a single mobile device capable of simultaneously performing data and voice operations. The mobile device working in a network capable of exchanging both cell phone calls and data items to the mobile device. By wearing an earphone or an ear-bud device the user is capable of dealing with voice conversations while working with data centric information related to the current caller. By providing a data-centric device with voice capabilities there is a new range of features that allow incoming data events to trigger outgoing voice events.  
 Correspondence Address: **DOCKET CLERK PO BOX 12608 DALLAS, TX 75225 (US)**  
 (73) Assignee: **Research in Motion Limited, Waterloo (CA)**  
 (21) Appl. No.: **11/458,843**  
 (22) Filed: **Jul. 20, 2006**  
**Related U.S. Application Data**  
 (63) Continuation of application No. 10/890,824, filed on Jul. 14, 2004, now Pat. No. 7,096,009, which is a

DocType = Patent Application

Date = 10/18/2007

Date Format = US

Author = Patent Authors, Author City, Author Country

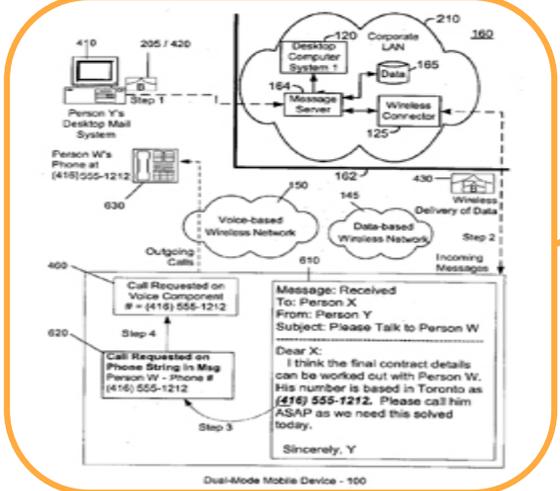
Assignee = RIM

Tone = Neutral to slightly positive

Embedded Graphic with Title

- Other Data Capturable Data Elements:
- Patent Number
  - Filing Date
  - Key Phrases & Terms
  - Managing PTO
  - Implied/Attached Docs
  - Bar Code Present
  - And many more . . .

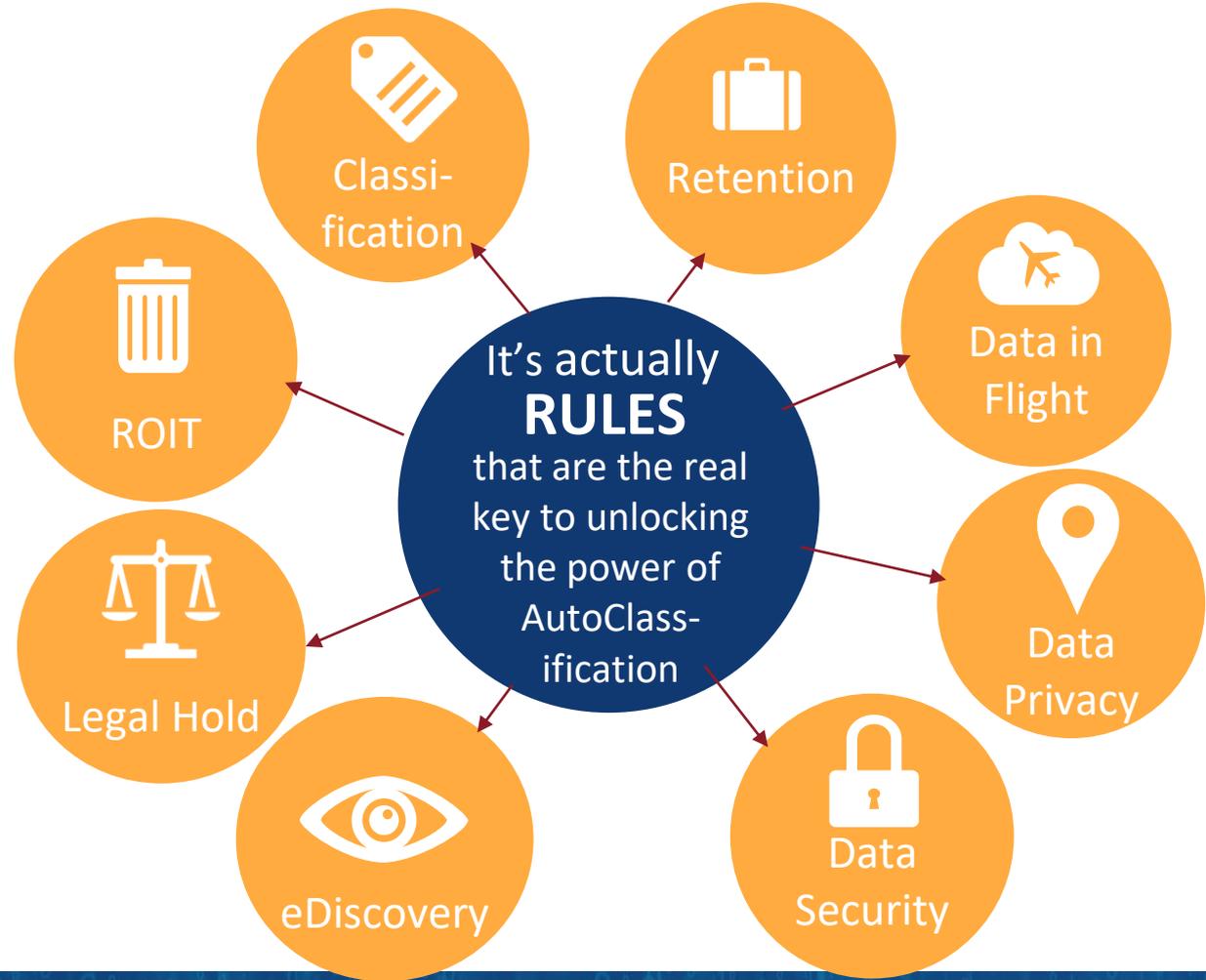
▶ **Implied status:**  
**Responsive, Nonprivileged**



## 2. Auto- Classification: Analyze and Design



# What is the secret sauce of auto-classification design?





## What are Rules?

- Software algorithms that determine the disposition of content
- Typically follow an IF-THEN format, often nested
  - Ex: IF the *DocumentType* = Contract AND the *Date* < 1/1/2020, THEN Mark File as ACTIVE (Retain)



So, if you could do any of these rules, which ones should you do first?  
[Hint: could you do multiple *simultaneously*?]

# Rules Engine illustration

Folders for organization

Ruler represents single rule

String test (can have multiple per rule)



The screenshot shows the PowerHouse Rules Editor interface. On the left is a tree view of rules organized into folders. The main area shows a configuration panel with three tabs: Operation, Fields, and Values. The Values tab is active, displaying a list of values and their weights. Callout boxes provide detailed explanations of various features.

Value	Weight
Aarhus	Medium
Abidjan	Medium
Abkhazia	Medium
Abomey-Calavi	Medium
Abu Dhabi	Medium
Accra	Medium
Adamstown is the only city	Medium
Addis Ababa	Medium
Addu City	Medium
Aden	Medium
Afghanistan	Medium
Airai	Medium
Akhalgori	Medium
Alotini and Dhekela	Medium
Akureyri	Medium

Support for values files

Weights & confidences create hierarchical decisioning options

Ability to test on family ranges and attachments (include & exclude)

Match type options:  
Exact, Word Stemmed, OCR, RegEx, Numeric, Date

Operation tab allows initial definition

Field tab denotes which PowerHouse fields rule is executed against

Values tab shows which values the rule should identify on



# Identifying & AutoClassifying PII

### Uniform Residential Loan Application

This application is designed to be completed by the Applicant(s) with the Lender's assistance. Applicants should complete this form as "Borrower" or "Co-Borrower" as applicable. Co-borrower information must also be provided (and the appropriate box checked) when  the income or assets of a person other than the Borrower (including the Borrower's spouse) will be used as a basis for loan qualification or  the income or assets of the Borrower's Spouse will not be used as a basis for loan qualification, but his or her liabilities must be considered because the Borrower resides in a community property state, or the security property is located in a community property state, or the Borrower is relying on other property located in a community property state as a basis for repayment of the loan.

If this is an application for joint credit, Borrower and Co-Borrower each agree that we intend to apply for joint credit (sign below):

Borrower		Co-Borrower		Interest Only	
<b>I. TYPE OF MORTGAGE AND TERMS OF LOAN</b>					
Mortgage Applied for:	<input type="checkbox"/> V.A. <input type="checkbox"/> Conventional <input type="checkbox"/> Other (explain):	Agency Case Number	Lender Case Number		
<input checked="" type="checkbox"/> FHA	<input type="checkbox"/> USDA/Rural Housing Service				
Amount	Interest Rate	No. of Months	Amortization Type:	Other (explain):	
\$ 625,000.00	6.5%	360	<input checked="" type="checkbox"/> Fixed Rate <input type="checkbox"/> Other (explain):	ARM (type):	
			GPM		
<b>II. PROPERTY INFORMATION AND PURPOSE OF LOAN</b>					
Subject Property Address (street, city, state, ZIP)					No. of Units
5421 Talahassee Place, Perry FL 32347					
Legal Description of Subject Property (attach description, if necessary)					Year Built
					1984
Purpose of Loan	Property will be:				
<input type="checkbox"/> Purchase <input checked="" type="checkbox"/> Construction <input type="checkbox"/> Other (explain):	<input checked="" type="checkbox"/> Primary Residence <input type="checkbox"/> Secondary Residence <input type="checkbox"/> Investment				
<input type="checkbox"/> Refinance <input type="checkbox"/> Construction-permanent					
Complete this line if construction or construction-permanent loan					
Year Acquired	Original Cost	Amount Existing Liens	(a) Present Value of	(b) Cost of improvements	Total (+ b)
Complete this line if this is a refinance loan					
Year Acquired	Original Cost	Amount Existing Liens	Purpose of Refinance	Describe Improvements	made <input checked="" type="checkbox"/> to be made
1996	\$ 375,000	\$ 402,500	Cash-Out/Home Improvement	Cost: \$20,000	
Title will be held in what Name(s)			Manner in which Title will be held	Estate will be held in:	
Brady and Wilma Schleshenger			Tenants in common	<input checked="" type="checkbox"/> Fee Simple <input type="checkbox"/> Leasehold (show expiration date)	
Source of Down Payment, Payment Changes and/or Subordinate Financing (explain)					
<b>Borrower</b>			<b>Co-Borrower</b>		
Borrower's Name (include Jr. or Sr. if applicable)			Co-Borrower's Name (include Jr. or Sr. if applicable)		
Brady G. Schleshenger			Wilma B. Schleshenger		
Social Security Number (no dashes)	DOB (MM/DD/YYYY)	Year of school completed	Social Security Number (no dashes)	DOB (MM/DD/YYYY)	Yes School
562-55-5592	(1953) 986-9985	04/11/1969 2	548-32-9555	(1970) 986-9985	02/16/1978 15
<input checked="" type="checkbox"/> Married <input type="checkbox"/> Single <input type="checkbox"/> Widowed <input type="checkbox"/> Divorced <input type="checkbox"/> Married (no. of dependents)	<input type="checkbox"/> Married <input type="checkbox"/> Single <input type="checkbox"/> Widowed <input type="checkbox"/> Divorced (no. of dependents)	Dependents (not listed by Co-Borrower) no. 1 ages 17	<input checked="" type="checkbox"/> Married <input type="checkbox"/> Single <input type="checkbox"/> Widowed <input type="checkbox"/> Divorced (no. of dependents)	Dependents (not listed by Borrower) no. 0	ages

Clear PII:  
SSN

Likely PII:  
("warning sign")

Clear PII: Home  
Phone Number

Not PII:  
Interest Rate

Implied classification:  
Active PII, needs  
protection &  
redaction

### 3. Information Governance: Control



# The information governance lifecycle

New content



## Processing

Ingests data or processes files in place

- Creating OCR for scanned images
- Transcribing voice and video
- Extracting text for native files & email
- Speech to text for audio/video files
- Translating content to English
- Re-ordering or re-aligning pages
- Applying redactions
- Unitization



## Tagging

Extracts key Info& attributes about each document (aka Coding, Indexing)

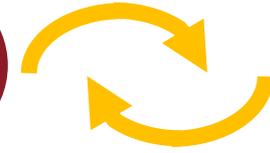
- Document Type, Important Dates
- Key Names & Phrases
- Topics, Keywords & Themes
- File, Content and DocType attributes
- Relation to other documents (duplicate, related, attached, contradictory, etc.)



## Disposition

(Rules) create a destination, action or status for each document

- Retention status & duration
- Privacy & security protections
- Folder (taxonomy) location
- Labelling & keywords display
- Activity notification
- Legal Hold



## Business Assessment

AI and Machine learning help assess risks to reputation, ID new IP

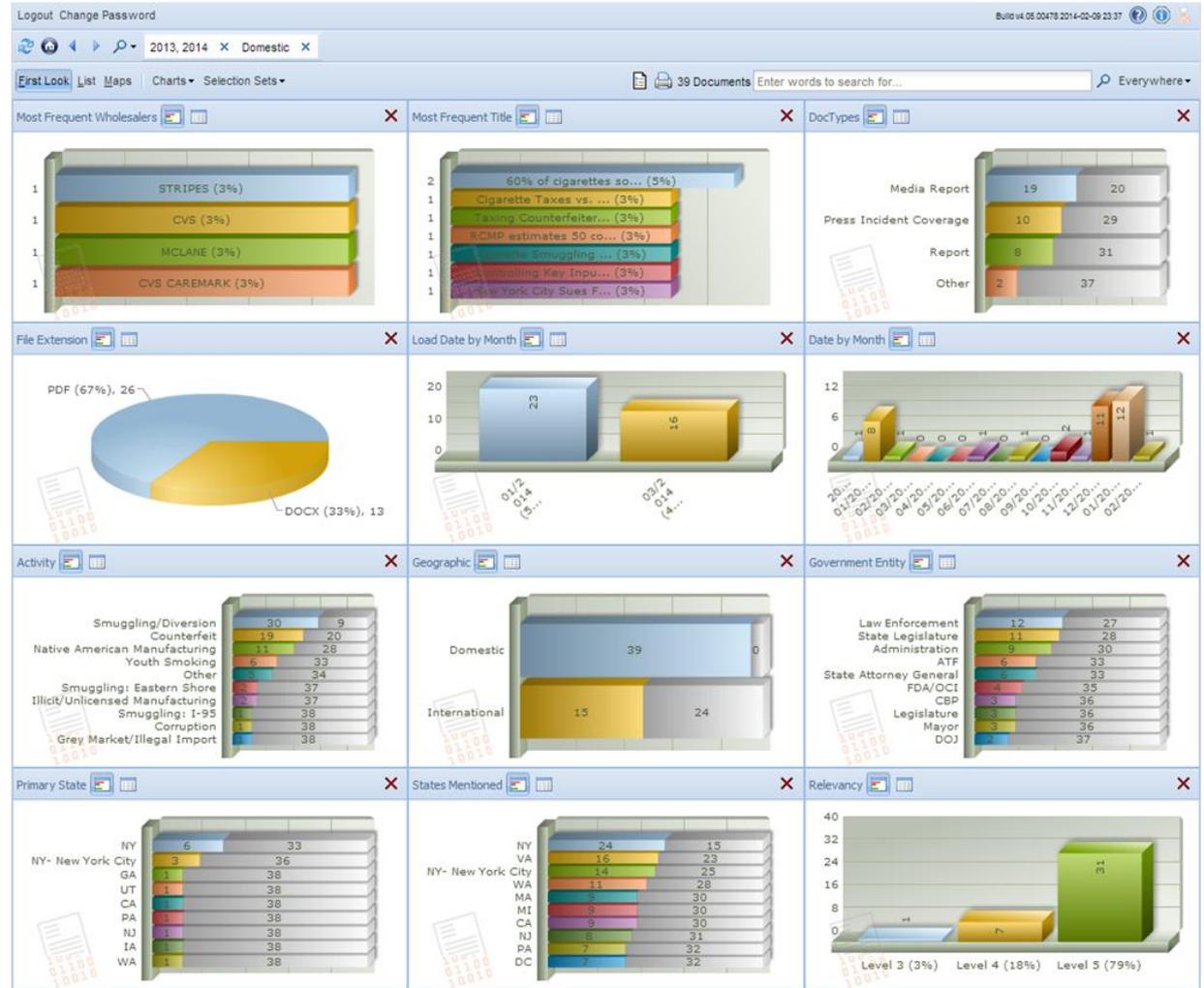
- Relevancy/ prioritization
- Contract inconsistencies
- Conflicts of interest
- New IP
- New Product or feature



Old content

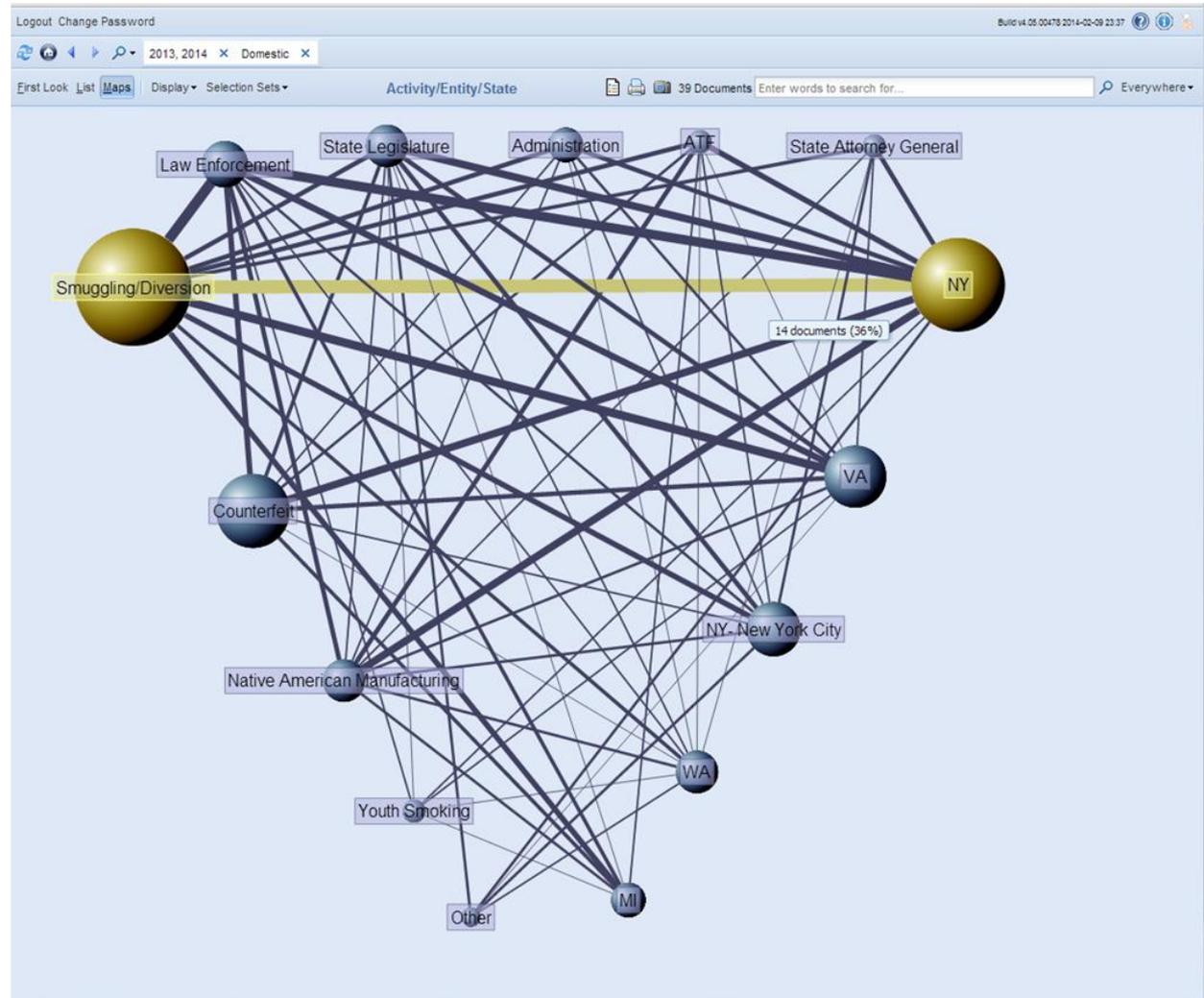


# Data Visualization of Auto Classification (Data dashboard)





# Data Visualization of Auto Classification (Graph Theory) (cont.)





# Data Visualization of Auto Classification

Logout Change Password Build v4.06.00478 2014-02-09 23:37

2013, 2014 Domestic

First Look List Maps Selection Sets 1 - 39 of 39 Documents Enter words to search for... Everywhere

Special	DocID	DocType	Date	Title	Authors	Activity	Govt. Entity	Pages	Geographic
<input type="checkbox"/>	VAL0000038	Media Report	01/10/2013	Cigarette Taxes vs. Cigarette...	New York Times; ...	Smuggling/Diversion		2	Domestic
<input checked="" type="checkbox"/>	VAL0000044	Press Incident Coverage	12/19/2013	New York City Council mulls...	Reuters; Skinner, ...	Other	Legislature; Mayor; St...	1	Domestic
<input type="checkbox"/>	VAL0000045	Media Report	02/04/2013	SMUGGLING MAY UNDERMINE...	LaFaive, Michael; ...	Smuggling/Diversion	Law Enforcement	3	Domestic; Intern...
<input type="checkbox"/>	VAL0000048	Media Report	01/10/2013	Financing Terrorism Terroris...	Cutting Edge; Wills...	Counterfeit; Native Ameri...	ATF; FBI	5	Domestic; Intern...
<input checked="" type="checkbox"/>	BI-MM-00000005	Report	01/10/2013	CIGARETTE TAX EVASION ...	WA State Dept. of ...	Native American Manufa...	Administration; State ...	1	Domestic
<input type="checkbox"/>	BI-MM-00000008	Press Incident Coverage	12/02/2013	Despite law, tribe sells 1.7 t...	Associated Press	Native American Manufa...	ATF; State Attorney ...	4	Domestic; Intern...
<input type="checkbox"/>	BI-MM-00000024	Press Incident Coverage	12/30/2013	New York City Sues FedEx ...	Dow Jones Busine...	Other		2	Domestic
<input checked="" type="checkbox"/>	BI-MM-00000026	Press Incident Coverage	12/19/2013	New York City Council mills ...	Reuters; Skinner, ...	Other	Legislature; Mayor; St...	1	Domestic

image Text Properties 1 of 2 100%

12/10/13 60% of cigarettes sold in New York are smuggled: report - Jan. 10, 2013



## 60% of cigarettes sold in New York are smuggled: report

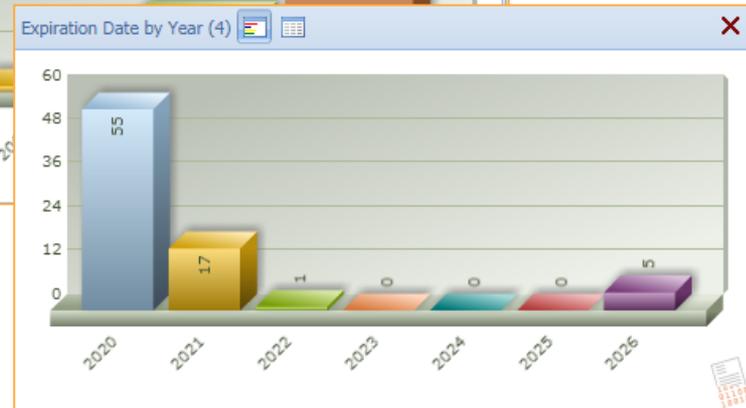
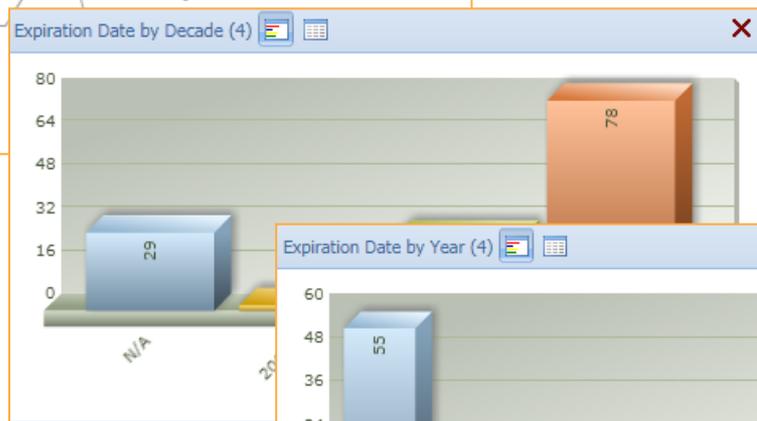
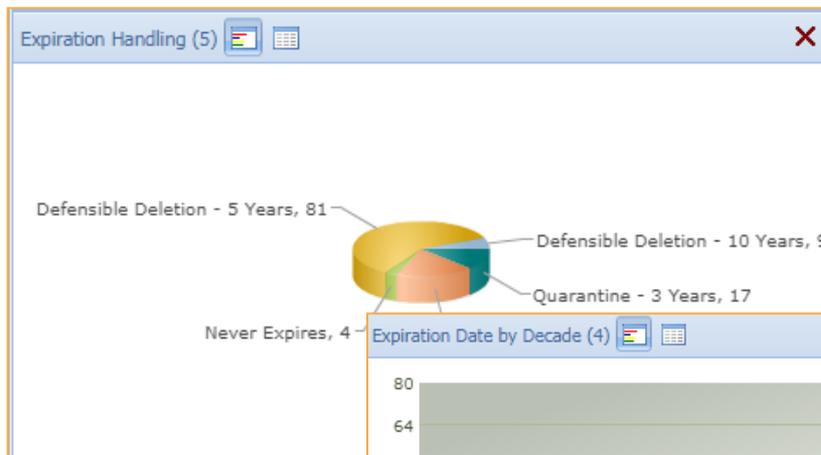
By Aaron Smith @AaronSmithCNN January 10, 2013: 3:24 PM ET

Recommended 2.1k





# Automated Retention Handling





Retention

# Automated Retention Handling (cont.)

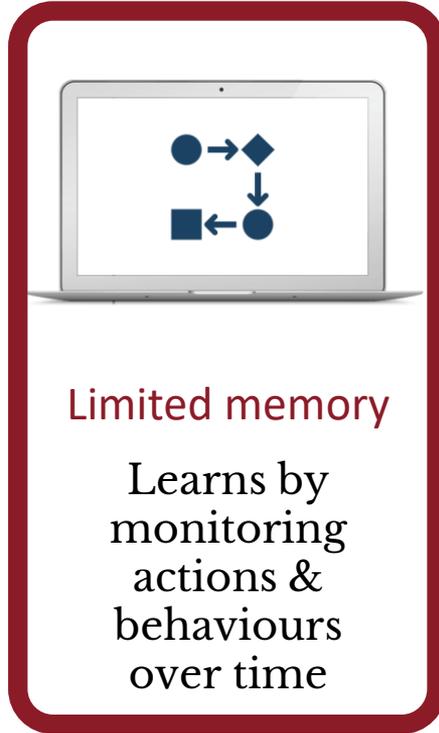
...	Special	DocID	DocType	Date	Expiration Date	Expiration Handling	Title
		IG-00028283	Purchase Document			Policy Review	CLOSING STATEMENT
		IG-00028284	Press Incident Cov...	03/09/2015	03/09/2020	Defensible Deletion - 5 Years	Counterfeit auto parts are threat to auto safety
		IG-00028287	Press Incident Cov...	07/24/2015	07/24/2020	Defensible Deletion - 5 Years	Counterfeit LED TVs from China confiscated at Miami Seaport
		IG-00028289	Press Incident Cov...	08/23/2016	08/23/2021	Defensible Deletion - 5 Years	Counterfeit pain pills likely came to Prince illegally
		IG-00028293	Press Incident Cov...	06/25/2015	06/25/2020	Defensible Deletion - 5 Years	Counterfeit products drive fashion industry to improve, study suggests
		IG-00028294	Report	03/2016		Policy Review	UNREGULATED E-WASTE EXPORTS FUEL COUNTERFEIT ELECT...
		IG-00028298	Purchase Document	04/30/2015	04/30/2018	Quarantine - 3 Years	Your Credit Card Account Statement
		IG-00028299	Press Incident Cov...	05/11/2015	05/11/2020	Defensible Deletion - 5 Years	Does Etsy Have A Problem With Fake Goods?
		IG-00028301	Report	03/2014		Policy Review	ARREST SUMMARY REPORT - January - February - March 2014
		IG-00028304	Civil Court Docume...	06/13/2014		Never Expires	UNRESOLVED ISSUE
		IG-00028314	Press Incident Cov...	07/30/2015	07/30/2020	Defensible Deletion - 5 Years	Etsy Investors Sue for Fraud over Risk of Counterfeit Goods Etsy Inc. ...
		IG-00028308	Report	03/21/2014		Policy Review	EXECUTIVE SUMMARY

# “State of the State” of Artificial Intelligence in auto-classification



## Reactive machines

Does not use past experiences to make decisions



## Limited memory

Learns by monitoring actions & behaviours over time



## Theory of mind

Has thoughts and emotions that affect their own behavior



## Self-awareness

Has consciousness, empathy and a sense of being

# Processing: Full spectrum of human communication

## Scientists have found a way to decode brain signals into speech

It's a step towards a system that would let  
people send texts straight from their brains.

by Antonio Regalado

MIT Technology Review, April 24, 2019



# Tagging: AI recognizes tone, slang

- Google Perspective API uses machine learning models to score the perceived impact a comment might have on a conversation. Developers and publishers can use this score to give real-time feedback to commenters or help moderators

**Reading Experiment**

What if it were easier to find new perspectives?

Here we've used Perspective to organize comments on three topics that are often difficult to discuss online. Use the slider to

Climate Change Brexit US Election



Showing 56 of 61 total comments based on toxicity\*

It's rubbish	idiots. backward thinking people. nationalists. not accepting facts. susceptible to lies.
Fools	
Hate it. Didn't vote for it.	They are ignorant
Dreadful: I'm a Remainer	Daft buggers
left wing wimps	I think they are wrong

Likely to be perceived as toxic (0.93) [Learn more](#) [SEEM WRONG?](#)

Eating tomatoes  
Loving lobster  
yucky  
fish pie  
soup nazi  
You are a jerk

**Writing Experiment**

What if you could see the potential impact of your writing?

<https://www.perspectiveapi.com/#/>

# Tagging: Intent recognition example: Sarcasm



**John Brock**

@jbrock5016

Follow

Well that was fun. More cold, snow, and ice please. [#sarcasm](#)



**Ellen**

@OrdinaryGirl11

Follow

Well 3 hrs of sleep was lovely. It's going to be a long day.

Humans don't always recognize it! It requires context (speaker, situation, world)

Approaches:

1. Machine learning (volume, then phrase matching / information retrieval)
2. Rules (e.g., unexpected juxtaposition, exaggeration)
3. Deep learning

Check out [Indian Institute of Technology research Sarcasm Suite](#) Source: [Quartz \(2016\)](#) and [NVDia](#)

# Tagging: Context switching

Book Flights   ONLINE MODE 

Thu Jul 12 2018 at 7:19 pm  
Hey, what's up

Thu Jul 12 2018 at 7:18 pm  
I'm good. Thank you.

Thu Jul 12 2018 at 7:19 pm  
Can you book me a flight ticket?

Thu Jul 12 2018 at 7:19 pm  
Sure, where would you like to travel from?

Thu Jul 12 2018 at 7:20 pm  
That'd be from San Francisco. By the way, what's the weather forecast there for this Sunday?

“How do you really carry context throughout a dialogue? This is the biggest challenge. A lot of how you understand what I’m saying depends on what I said maybe five sentences ago, or fifteen sentences ago. You’re building up a state of the conversation.”

Satinder Singh, Director of the  
Artificial Intelligence Lab  
at the University of Michigan.

Source: [IBM](#)  
([Hold and resume concept \(Kore.ai\)](#))

## 4. Case Studies



Recap: Do  
any of  
these look  
like you?



## Large Oil & Gas Manufacturer

- Frequent acquisitions of other companies
- Millions of documents
- Not enough SMEs
- Risky content..Legal hold...Data Privacy... ROIT



## Multi-store, multi-country fast-food

- Independent, franchise & retail sites
- Managers autonomously hire, fire
- GDPR concerns



## Info hound (bankruptcy data vault)

- Out run by content volume
- Regulatory reporting obligations
- Sensitive data
- Missing value-added revenue...?

# Case Study #1: Oil & Gas



Acquiring a company without most of the people. 200 GB of data and none of the tacit knowledge - who owns this? Should we keep it?

Burning Platform

Legal Hold  
Retention  
Privacy/PII  
Contract anomalies  
Animal testing compliance

Scope

Custodian  
Matter  
PII/PHI  
Record Type  
Expiry Handling  
Dates

Critical Metadata

Retention  
ROT  
Data Privacy  
Classification  
Legal Hold  
Duplication

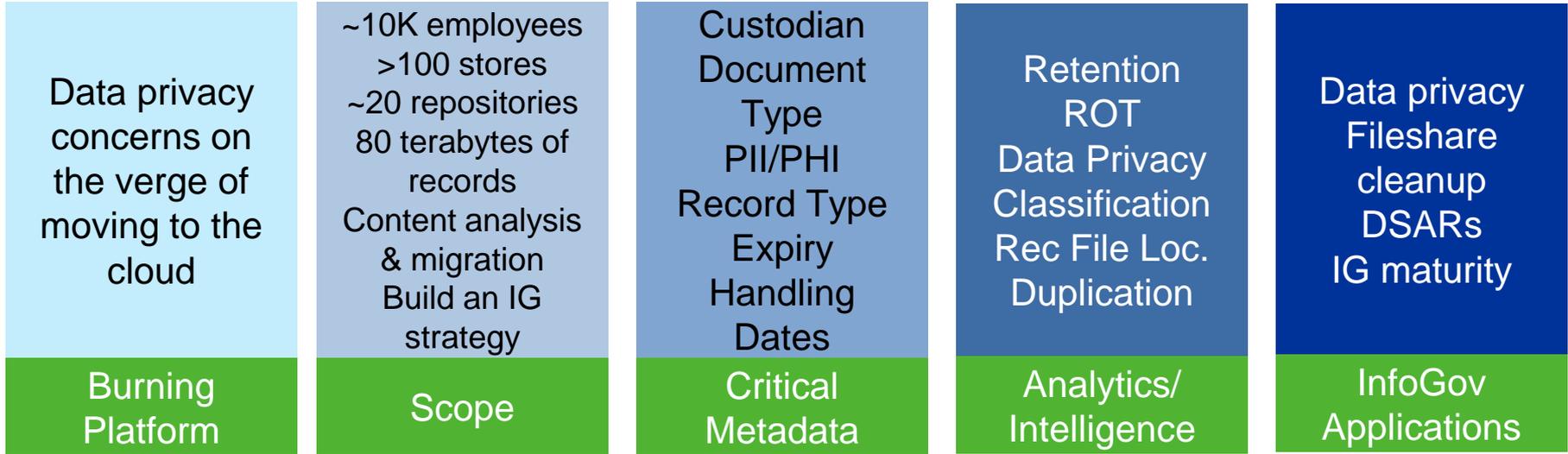
Analytics/  
Intelligence

Contract consolidation  
Compliance watch  
Fileshare cleanup  
DSARs

InfoGov  
Applications

“Don’t buy it before you know the info risk!”

# Case Study #2: Multi-store/int'l fast food retailer



“It’s not about migration. It’s about gov and access.”

# Case Study #3: Info Hound (bankruptcy data vault)



Too much data  
to assess,  
evaluate,  
prioritize in  
tight timelines

Burning  
Platform

Federal, state,  
local bankruptcy  
court filings  
1705 Juris-  
dictions  
300 Formats  
100K-200K  
cases/yr

Scope

Debtor Name &  
specifics  
Custom fields  
per DocType  
Jurisdictional  
Dates  
Financial

Critical  
Metadata

Bankruptcy terms  
violations?  
Defrauding  
regulators?  
Criminal checks?  
Redundancies?  
Expirations?

Analytics/  
Intelligence

Compliance  
reporting  
Risk analysis

InfoGov  
Applications

“The 2008 meltdown forced us all to rethink & retool.”

# Three Cases' lessons learned

1. You lay bare facts and patterns (“It’s like walking with all of the labels on your clothes on the outside”)
2. Building the rules engine – how you label it, and what you do with it -- takes longer than you’d think. (“You learn when to dig deeper and when to back off.”)
3. It’s not just “Document Intelligence,” or “eDiscovery.” It’s “info governance”

# Discussion

# Thanks for attending!



## AlignConsulting



[www.ValoraTech.com](http://www.ValoraTech.com)



[www.alignconsulting.com](http://www.alignconsulting.com)



[info@ValoraTech.com](mailto:info@ValoraTech.com)



[katepugh@alum.mit.edu](mailto:katepugh@alum.mit.edu)



+1 781.229.2265



+1 617-967-3910



# Appendix



# Rules for Protecting Consumer Data

Protecting PII is  
at the core of  
privacy  
regulations

## Privacy Regulations

require:

- ❓ Do you know **where** it is?
- ❓ Do you know **what** it is?
- ❓ Is it **sensitive**?
- ❓ Is it **protected and secure**?
- ❓ Is it **accessible** when needed?

## AutoClassification

allows you to:

- ✅ **Locate it:**  
email, shared drive, ECM
- ✅ **Identify it:**  
SSN, DOB, drivers license
- ✅ **Secure it:**  
set security access controls
- ✅ **Redact it:**  
grant permission levels
- ✅ **Action it:**  
DSAR, Right to be Forgotten



# Rules for eDiscovery

Note: If you've got a good IG practice through AutoClassification, eDiscovery is much simpler, faster & lower cost

## eDiscovery

requires us to:

- ❓ Look for content
- ❓ Filter through tons of data
- ❓ Relevant & irrelevant, Privileged & Non-priv
- ❓ Document review
- ❓ Document production

## AutoClassification

you've already:

- ✅ **Located it:**  
across all data stores
- ✅ **Identified it:**  
for Early Data Assessment
- ✅ **Cleaned it:**  
only dealing with relevant content
- ✅ **Secured it:**  
who has access to what
- ✅ **Actioned it:**  
Retention, archived, etc.



# Rules or Legal Hold

Note: You can still dynamically govern content under Legal Hold

## Legal Hold

asks us:

- ❓ Which data? Whose data?
- ❓ For how long?
- ❓ What's relevant? What's not?
- ❓ Access for inside counsel?  
Outside counsel?
- ❓ What happens when Hold lifts?

## AutoClassification

you've already:

- ✅ **Located it:**  
we haven't missed anything/where
- ✅ **Cleaned it:**  
already excluded irrelevant content
- ✅ **Identified type & custodian:**  
we know content, context & owner
- ✅ **Identified retention dates:**  
schedule already exists
- ✅ **Set it & forget it:**  
ongoing perpetual analysis