



UNIVERSITY
OF ILLINOIS
SYSTEM

A L T O G E T H E R E X T R A O R D I N A R Y

Processing Capstone Email Using Predictive Coding

Brent M. West
University of Illinois

go.uillinois.edu/capstone

PROJECT OVERVIEW

NHPRC-funded 3 year project with Illinois State Archives

Dataset: 5.3M messages of most senior state officials from 2000-2014

Workflow:

1. File format migration (Emailchemy)
2. Appraise & sanitize (Ringtail)
3. Preservation repository (Preservica)
4. Access (EPADD)



🔄 🔍 ⚙️ 📄
Browse ▾

☰ ☰
Custom (no template) ▾
📄 List ▾

📧 P05201C012_0000...
0/0
🔍
⬇️

Review ▾
🔄 📄

		Document Title	Document Date
<input type="checkbox"/>	△ <input type="checkbox"/>	Clips 04.23.09 Thursday	4/23/2009
<input type="checkbox"/>	<input type="checkbox"/>	RE: LAC hearing	11/20/2008
<input type="checkbox"/>	<input type="checkbox"/>	Gov. Blagojevich reminds youn	5/14/2008
<input type="checkbox"/>	<input type="checkbox"/>	Procurements - Cost Cutting G	10/28/2009
<input type="checkbox"/>	<input type="checkbox"/>	[Unnamed Document]	12/5/2003
<input type="checkbox"/>	<input type="checkbox"/>	Updated: MW Weekly Update	5/30/2007
<input type="checkbox"/>	<input type="checkbox"/>	2008 N I INV LIST.doc	1/21/2009
<input type="checkbox"/>	<input type="checkbox"/>	Services Report to full OASAC.t	6/9/2009
<input type="checkbox"/>	<input type="checkbox"/>	Hearing on HIE Initiative	10/9/2009
<input type="checkbox"/>	<input type="checkbox"/>	Interview Questions for Directc	7/28/2005
<input type="checkbox"/>	<input type="checkbox"/>	Tax exemptions	7/2/2006
<input type="checkbox"/>	<input type="checkbox"/>	DocLink1.ndl	9/13/2005
<input type="checkbox"/>	<input type="checkbox"/>	Fwd: Re: Quarrance Claiborne	4/25/2007
<input type="checkbox"/>	<input type="checkbox"/>	RE: One more thing	12/1/2006
<input type="checkbox"/>	<input type="checkbox"/>	Fw: Whistleblower legislation	3/17/2003
<input type="checkbox"/>	<input type="checkbox"/>	SPSA Project	3/5/2008
<input type="checkbox"/>	<input type="checkbox"/>	Re: Sen DeMuzio	6/27/2008
<input type="checkbox"/>	<input type="checkbox"/>	SUA Recovery Act Reporting 1:	8/31/2009
<input type="checkbox"/>	<input type="checkbox"/>	Last Chance to Register: Impac	4/6/2010

Clips 04.23.09 Thursday

From: "Anderson, Alicia" <alicia.anderson@illinois.gov>
Date: Thu, 23 Apr 2009 11:33:55 -0500
Attachments: 04.23.09 Thursday.doc (71.17 kB)

CLIPS: 04.23.09

Thursday

Tribune Campaign <> finance: Money's influence not easily curbed

Tribune Illinois <> historic sites: 11 closed sites to reopen

Tribune Quinn <> orders Ill. agencies to cut waste, pollution

Tribune Ill. <> governor faces rally against his budget

Tribune Quinn won't apologize for Blagojevich connections <> (Pantagraph)

SJR Quinn calls <> for greener Illinois

SJR Holocaust <> survivor recalls Auschwitz horrors

STL Illinois' <> closed historic sites to reopen (PJ Star)

STL Illinois rally supports tax hikes <>

Southern Four Southern <> Illinois historic sites reopen today

Review

Review

Archival

Non-archival

Archival

Not Coded

Restricted

Public

Restricted

Not Coded

Flag for second review

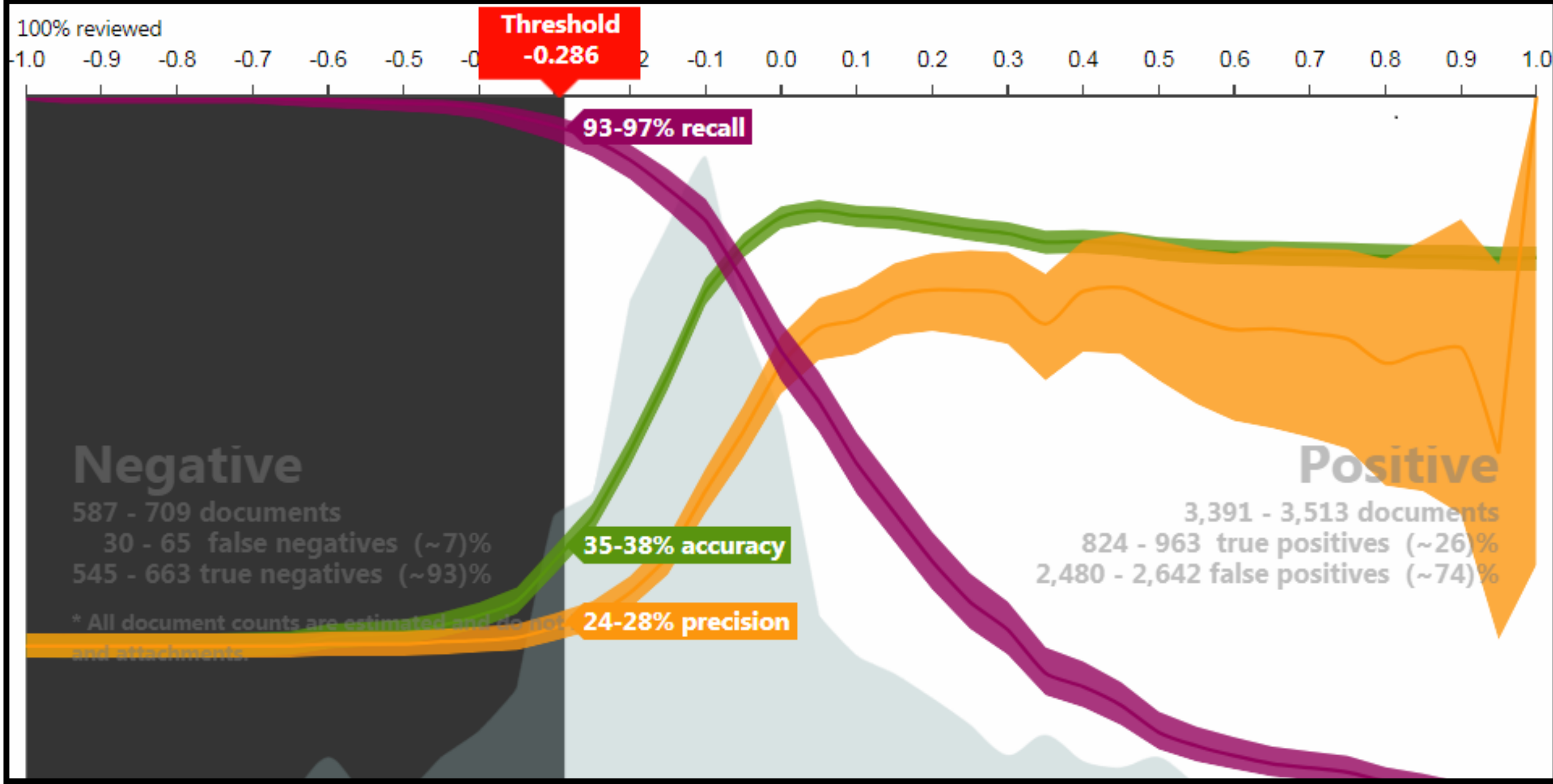
Archival re-review

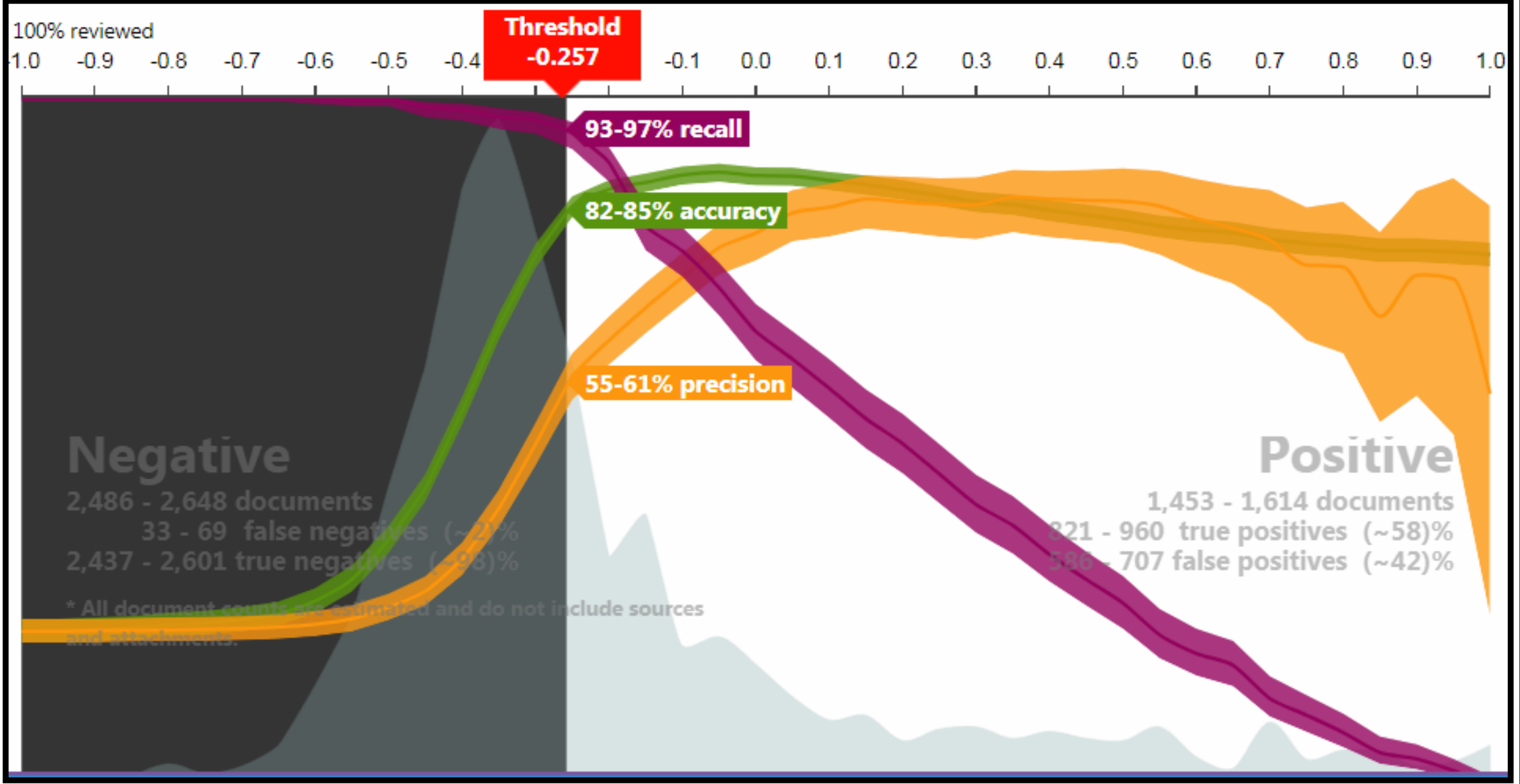
Restricted re-review

Issues

Comments

Metadata





LESSONS LEARNED

1. Augmented not automated appraisal
2. Huge labor savings (99.4%)
3. Capstone \neq all archival (41%)
4. PII small but important ($<2\%$)
5. Quantify and reduce risk
6. “Block box” commercial solutions
7. Forgiving to mistakes but expertise required
8. Trade-offs between privacy and transparency
9. Proceed deliberately: document decisions and context
10. Integrate into a lifecycle workflow

Pre Processing	OCR
Creating RDX Guid Tables	OCR
Batching	Indexing
Pre Processing	Update Field Counts
Processing	Update Field Counts
Processing	Transfer Suppressed Files
Post Processing	File Copy Batching
Standardizing Data	File Transfer
Import Files and Metadata	File Transfer Confirmation
Load Data	Gathering Report Data
Cleanup	Gathering Report Data
Hashes	Submit Index and Enrichment Job
Hashes	Submit Index and Enrichment Job
Group Coding	Submit Ingestions Cleanup Job
Group Coding	Submit Ingestions Cleanup Job
Data Filtering	Finalize Job
Data Filtering	Cleanup
Transfer Unsuppressed Files	
File Copy Batching	
File Transfer	
File Transfer Confirmation	

FY24 Ringtail Project_converted.7z

Description	Started via upload
Job ID	2
Started	11/9/2023 9:54 AM
Duration	1h 39m 52s
Status	Completed with exceptions

Processed files: 601,125 (94.91GB)

Total expanded documents (77.14GB)	714,252		
<i>Suppressed documents</i>			
Duplicates	19,513	Master Document	2.7%
Outside date range	Not applicable		
Excluded NIST files	23		0.0%
Excluded by category or extension	0		0
Does not match search term family	Not applicable		
Unsuppressed documents (76.07GB)	694,716		97.3%

Suppressed document ingestion exceptions	116
Renamed Extension	33
NIST Item	27
Non-Searchable PDF	23
Text Stripped	16
Multimedia	7
Unknown Binary	6
Encrypted	2
Corrupted	1
Export Failed	1

Unsuppressed document ingestion exceptions i 19,462

Renamed Extension	4,618
Non-Searchable PDF	4,614
Text Stripped	4,304
Corrupted	2,004
Unknown Binary	1,343
System File	860
Multimedia	642
Encrypted	248
NIST Item	241
Data Type Conversion Failed	228
Field Data Truncated	124
Missing Hash Value	79
Empty	49
Extracted Text Only	30
Export Failed	27
Field Data Extraction Error	25
Inaccessible Content	23
Databases	3

Included fields	Chat Start Time	Evidence ID	Office Property - Comments
	Chat Viewed History Count	Excluded File	Office Property - Company
BCC	Conversation Topic	Exif Date Time	Office Property - Date Last Printed
CC	Date Accessed	Exif Date Time Digitized	Office Property - Date Last Printed Time
Collection ID	Date Accessed Time	Exif Date Time Original	Office Property - Date Last Saved
Custodian	Date Appmt End	Extended File Path	Office Property - Keyword
Custodian ID	Date Appmt End Time	File Application	Office Property - Last Author Saved By
Document Date	Date Appmt Start	File Extension - Loaded	Office Property - Title
Document Type	Date Appmt Start Time	File Extension - Original	Original Full Path
Evidence Job ID	Date Created	File Name	PDF - Encryption Level
From	Date Created Time	File Path	PDF - Portfolio
Media ID	Date Modified	File Size	PDF Properties
To	Date Modified Time	GUID	Processing Exceptions
Appmt Location	Date Received	GUID - Parent	Processing Time Zone
Appmt Optional Attendees	Date Received Time	ISO Media Type	Revision Number
Appmt Organizer	Date Sent	Languages	Subject
Appmt Required Attendees	Date Sent Time	Latitude	Transport Message Headers
Chat Action Count	Date Top Family	Longitude	Word-Pdf-Image Page Count
Chat Attachment Count	Date Top Family Time	Mapi-DeliveredTo	DPM File ID
Chat Disclaimer Count	Document Category	Mapi-Message-Flags	Ingestion Complete
Chat End Date	Document Kind	Message Class	Ingestion Exception Detail
Chat End Time	Email Delivery Receipt Request	Multimedia Audio Codec	MD5 Hash
Chat Enter Count	Email Folder	Multimedia Duration	
Chat Event Count	Email Importance	Multimedia Video Codec	
Chat Exit Count	Email Message ID	NSF UNID	
Chat Invite Count	Email Message ID Replied To	Office Exceptions	
Chat Message Count	Email Read Receipt Request	Office Exceptions - Excel Hidden Sheet Count	
Chat Participant Count	Email Sensitivity	Office Exceptions - Notes Count	
Chat Sender Count	Entity	Office Property - Author	
Chat Start Date	EntryID		

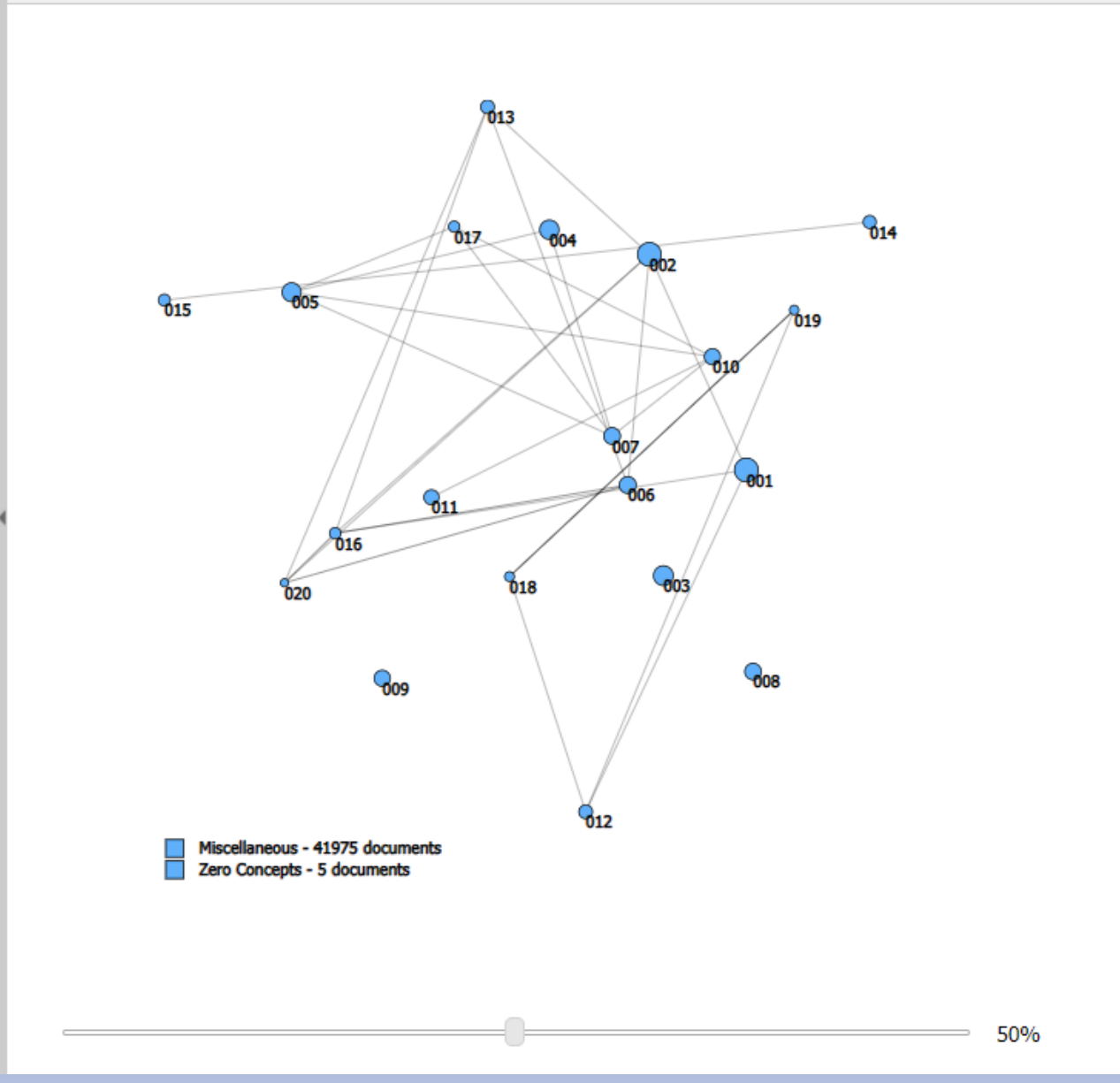
Rows: All × Date ×

	Document Kind ▾	Document Type ▾	Org ▾	People ▾	Relevance: Archival ▾
1996 ▾	45	45	22	22	0
1997 ▾	22	22	18	18	0
1998 ▾	36	36	22	23	0
1999 ▾	423	423	357	357	0
2000 ▾	7523	7523	6116	6120	0
2001 ▾	12317	12317	9169	9176	5
2002 ▾	18709	18709	12718	12798	0
2003 ▾	22988	22988	17014	17087	0
2004 ▾	734	734	570	571	0
2005 ▾	358	358	325	325	0
2006 ▾	603	603	503	503	0
2007 ▾	759	759	664	664	0
2008 ▾	1241	1241	1114	1114	0
2009 ▾	1432	1432	1292	1292	0
2010 ▾	7761	7761	6791	6795	0
2011 ▾	9497	9497	7835	7843	0
2012 ▾	10055	10055	8392	8394	0
2013 ▾	13786	13786	11094	11099	0
2014 ▾	19082	19082	14322	14338	0
2015 ▾	19391	19391	15260	15261	0
2016 ▾	16442	16442	13446	13453	0
2017 ▾	16564	16564	13474	13483	0
2018 ▾	17782	17782	14358	14363	0
2019 ▾	20687	20687	16725	16727	0
2020 ▾	204904	204904	178341	178389	0
2021 ▾	292629	292629	244332	244397	0
2022 ▾	38397	38397	32322	32334	0
2023 ▾	2584	2584	2188	2196	0

Cluster List

ID	Name	Total
001	Fenner, Melvin, MEL	45663
002	Shield, Julie, Shield Illinois	45406
003	student, community, Education	37940
004	archive, Library, Ellen	36877
005	Library, Urbana, Archivist	36326
006	Watkins, Ronald, Ron	32510
007	Library, archive, Urbana	32039
008	New, Twitter, Facebook	31825
009	Brewer, Michael, Mike	31022
010	Library, Copyright, archive	30919
011	Senate, Committee, faculty	28857
012	file, Medicat, NAME	26385
013	Yoni, Merid, lab	26235
014	Unsubscribe, privacy, Data	24872
015	SUBSCRIPTION, Service, release	22646
016	icon, vector graphics, white cross Description	21773
017	Archivist, Library, ICA	21284
018	Redwood City, Jefferson Avenue, Redwood	18748
019	Box Health, Redwood City, Redwood	18116
020	Greta, Andrew, Shield	14458
ZERO	No concepts	5

Mine Map



Concepts

Name	Density	Total
Illinois	<div style="width: 100%;"></div>	326037
University	<div style="width: 100%;"></div>	317442
information	<div style="width: 100%;"></div>	235350
business	<div style="width: 100%;"></div>	178976
Service	<div style="width: 100%;"></div>	172326
Director	<div style="width: 100%;"></div>	159205
office	<div style="width: 100%;"></div>	158612
SYSTEM	<div style="width: 100%;"></div>	152808
Urbana-Champaign	<div style="width: 100%;"></div>	152793
Urbana	<div style="width: 100%;"></div>	149917
work	<div style="width: 100%;"></div>	149743
question	<div style="width: 100%;"></div>	146692
Library	<div style="width: 100%;"></div>	142040
Shield	<div style="width: 100%;"></div>	141326

Confidence		
Confidence	Count	Density
90 - 100%	4	
80 - 90%	1934	<div style="width: 100%;"></div>
70 - 80%	11854	<div style="width: 100%;"></div>
60 - 70%	57942	<div style="width: 100%;"></div>
50 - 60%	125747	<div style="width: 100%;"></div>
40 - 50%	212786	<div style="width: 100%;"></div>

RE: Hoke out at Z bldg today or tomorrow.

From:

"Brewer, Michael K (O&M)" <mkbrewer@university-of-illinois>

To:

"Madsen, Kenneth M (O&M)" <kmmadsen@oandm.uiuc.edu>

Date:

Mon, .. Feb:06 +0000

I've got a message in to to find out what the status is.

We asked them to deliver .-. of the a long time ago. If they can not do that now, we will buy some more on Keith's money to get some in stock.

Thanks for me.

-Mike

-----Original Message-----

From: Madsen, Kenneth M (O&M)

Sent: Monday, February .., :.. AM

To: Brewer, Michael K (O&M)

Subject: RE: Hoke out at Z bldg today or tomorrow.

Mike

Are we going to get ." pipe from Hoke?

-----Original Message-----

From: Brewer, Michael K (O&M)

Sent: Monday, February .., :.. AM

To: Roberts, Joseph P (O&M); Madsen, Kenneth M (O&M)

Cc: Larson, Michael J (O&M); Erickson, Keith R (O&M)

Subject: Hoke out at Z bldg today or tomorrow.

Hoke called and will send someone up to stop the leak at the strainer today or tomorrow. They may have to shut the station down to replace the gasket.

When they do, I'll call so we can look to agree that the leak is stopped and let Mr. know he will need to do

Thank you.

Mike Brewer, PE, Utilities, O&M

Phone (...) ...-7067