

Matthew Brehmer



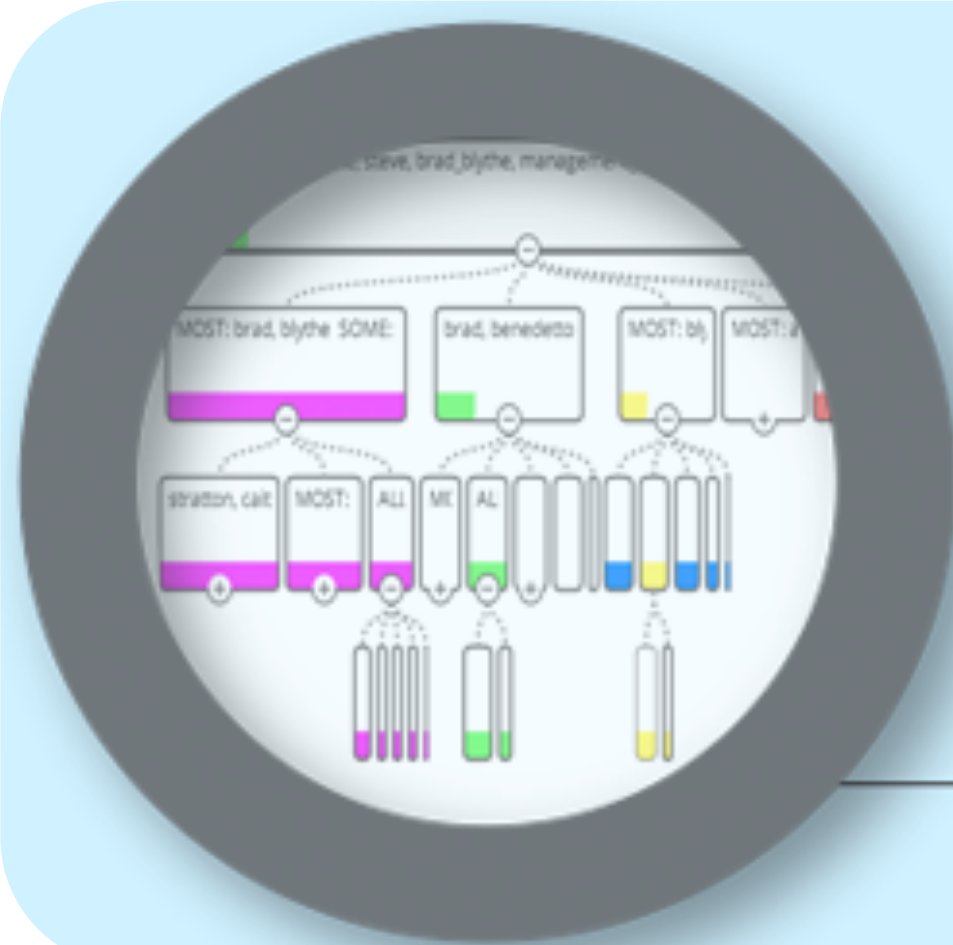
Stephen Ingram



Jonathan Stray



Tamara Munzner



# overview

the design, adoption, and analysis of a visual document mining tool for investigative journalists

IEEE InfoVis  
Nov 14, 2014



Matthew Brehmer



Stephen Ingram



Jonathan Stray



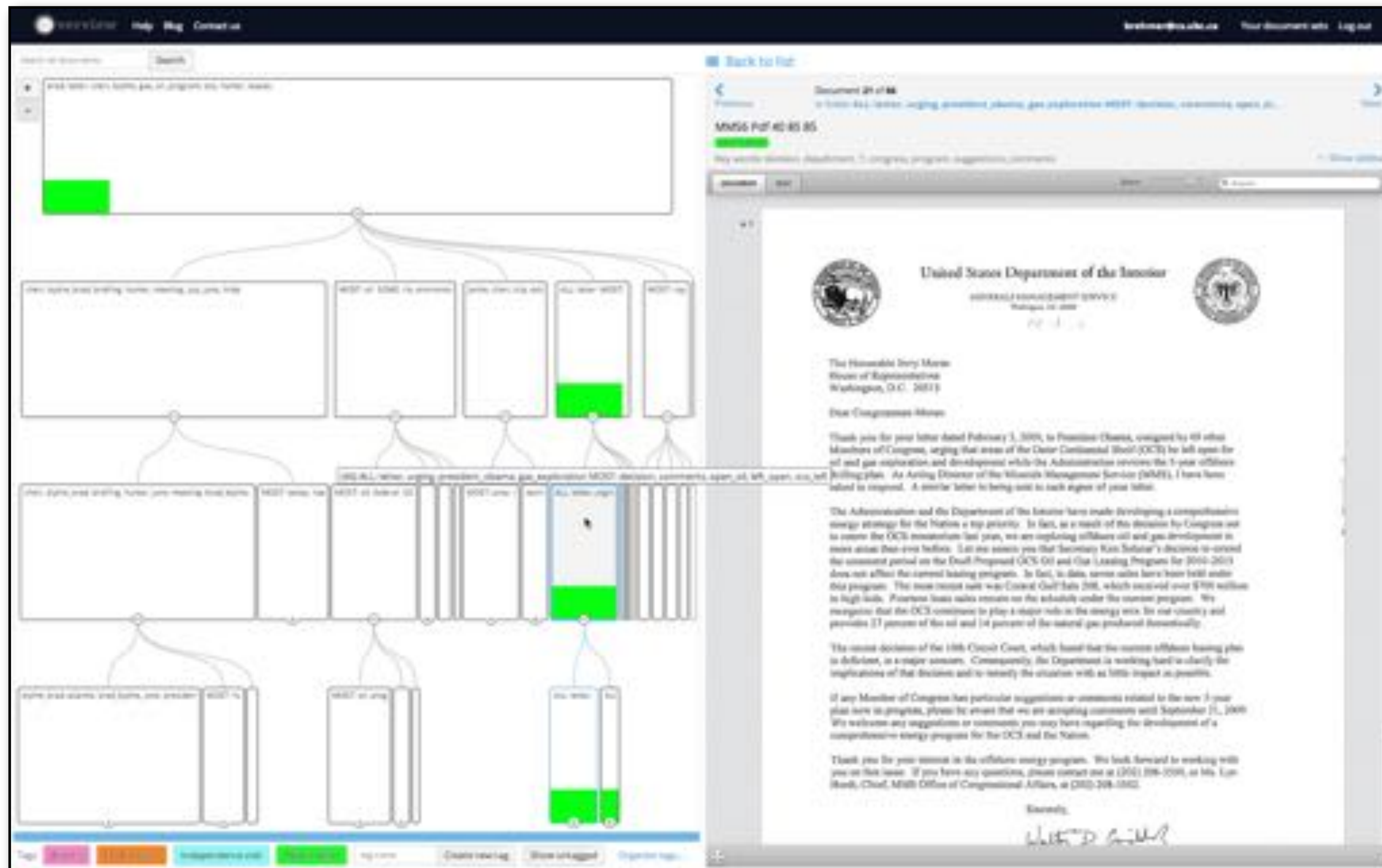
Tamara Munzner



the design, adoption, and analysis of a visual document mining tool for investigative journalists



# overview



overviewproject.org

**Police misconduct hidden from public by secrecy law, weak oversight**  
By Sandra Peddie and Adam Playford  
**Newsday**

**What did private security contractors do in Iraq?**  
by Jonathan Stray on 02/21/2012 | 3 | Edit  
**AP**

**TPD working through flawed mobile system**  
By JARREL WADE World Staff Writer on Jun 3, 2012, at 2:19 AM  
**TULSAWORLD**

**RYAN ASKED FOR FEDERAL HELP AS HE CHAMPIONED CUTS**  
By JACK GILLUM — Oct. 12 7:20 PM EDT  
**AP**

**INTERACTIVE 12.22.12**  
**Own a Gun? Tell Us Why**  
Michael Keller  
**THE DAILY BEAST**

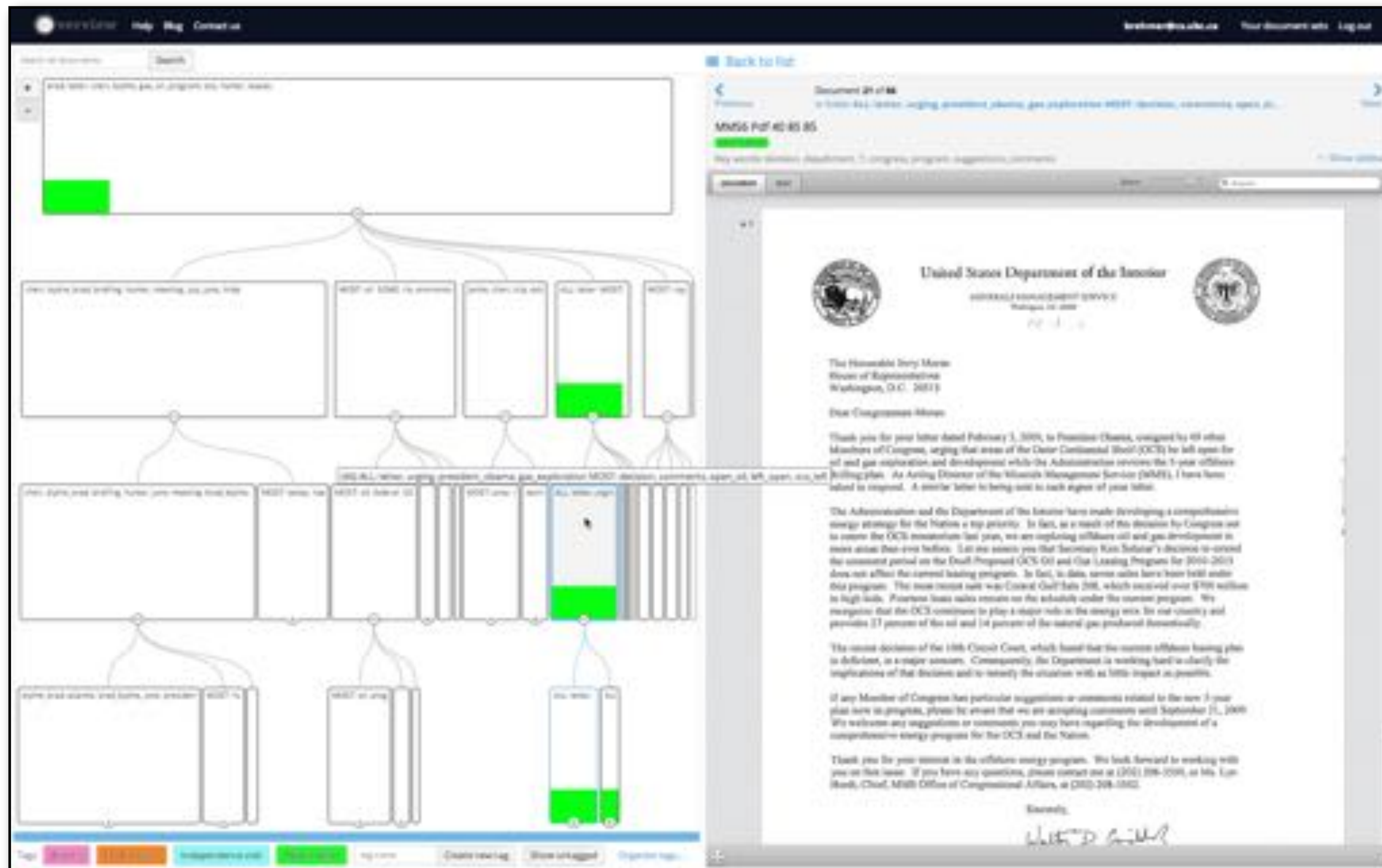
**The Brilliance of Louis C.K.'s Emails: He Writes Like a Politician**  
Where campaign strategy and comedy marketing collide  
ADRIENNE LAFRANCE | JUL 16 2014, 6:10 AM ET  
*the Atlantic*

**Surprise! Many credit card agreements allow repossession**  
Analysis: 'Security interest' clause present on 200 cards  
By Fred O. Williams  
**CreditCards.com**

Private memo reveals winding tale involving John McCain, the NRA and ... condors  
by Nancy Watzman | SEPT. 18, 2014, 1:59 P.M.  
**SUNLIGHT FOUNDATION**



# overview



overviewproject.org

**Police misconduct hidden from public by secrecy law, weak oversight**  
By Sandra Peddie and Adam Playford  
**Newsday**

**What did private security contractors do in Iraq?**  
by Jonathan Stray on 02/21/2012 | 3 | Edit  
**AP**

**TPD working through flawed mobile system**  
By JARREL WADE World Staff Writer on Jun 3, 2012, at 2:19 AM  
**TULSAWORLD**

**RYAN ASKED FOR FEDERAL HELP AS HE CHAMPIONED CUTS**  
By JACK GILLUM — Oct. 12 7:20 PM EDT  
**AP**

**INTERACTIVE 12.22.12**  
**Own a Gun? Tell Us Why**  
Michael Keller  
**THE DAILY BEAST**

**The Brilliance of Louis C.K.'s Emails: He Writes Like a Politician**  
Where campaign strategy and comedy marketing collide  
ADRIENNE LAFRANCE | JUL 16 2014, 6:10 AM ET  
*the Atlantic*

**Surprise! Many credit card agreements allow repossession**  
Analysis: 'Security interest' clause present on 200 cards  
By Fred O. Williams  
**CreditCards.com**

Private memo reveals winding tale involving John McCain, the NRA and ... condors  
by Nancy Watzman | SEPT. 18, 2014, 1:59 P.M.  
**SUNLIGHT FOUNDATION**



# A VISUALIZATION DESIGN STUDY



**real** journalists, real  
documents, real  
investigations...

## A VISUALIZATION DESIGN STUDY



**real** journalists, real  
documents, real  
investigations...

...**rational**e for visual  
encoding and  
interaction design  
choices...

## A VISUALIZATION DESIGN STUDY



**real** journalists, real  
documents, real  
investigations...

...**rational**e for visual  
encoding and  
interaction design  
choices...

...generalizable **lessons**  
for visualization design  
methodology

## A VISUALIZATION DESIGN STUDY



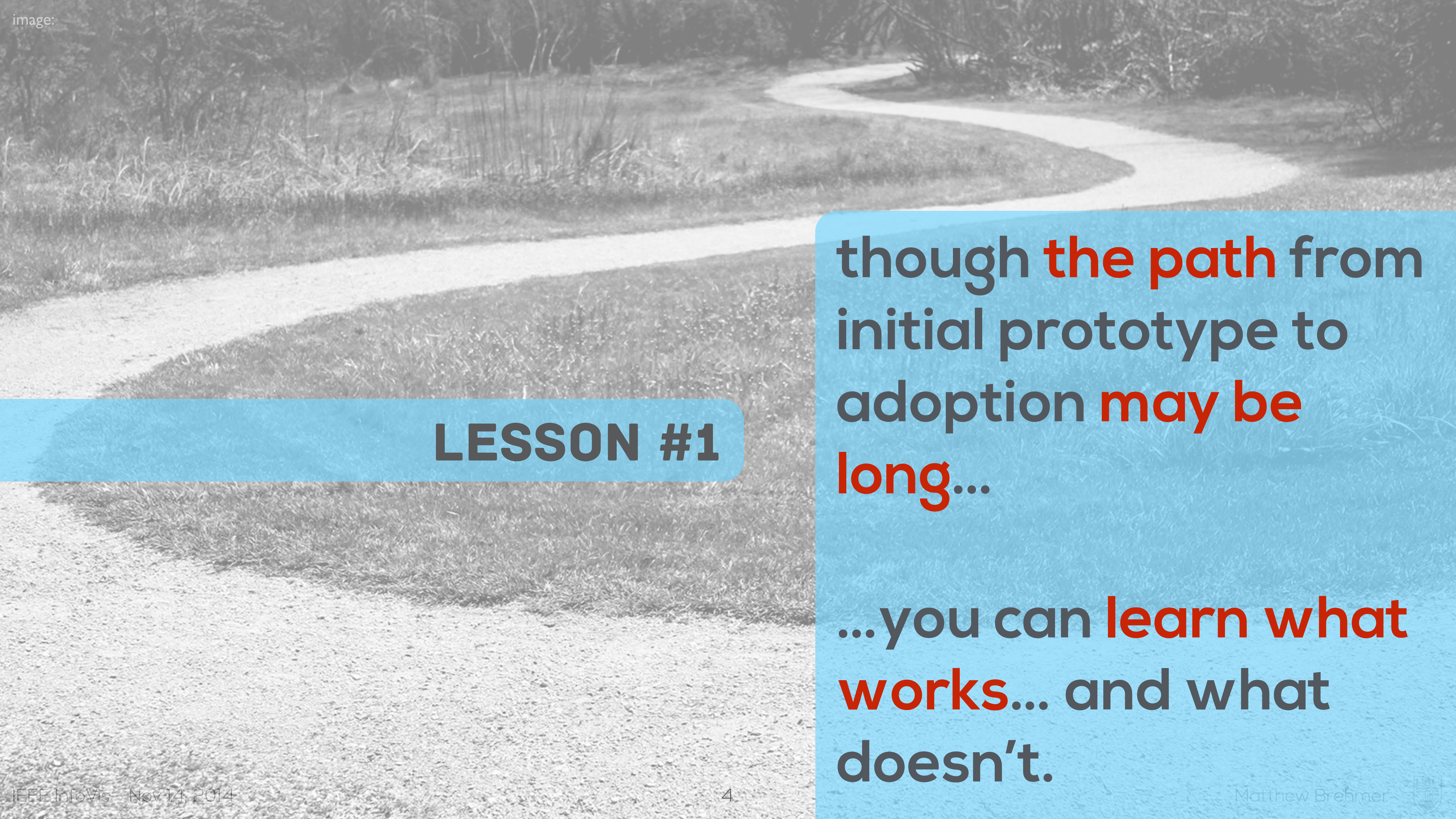


image:

## LESSON #1

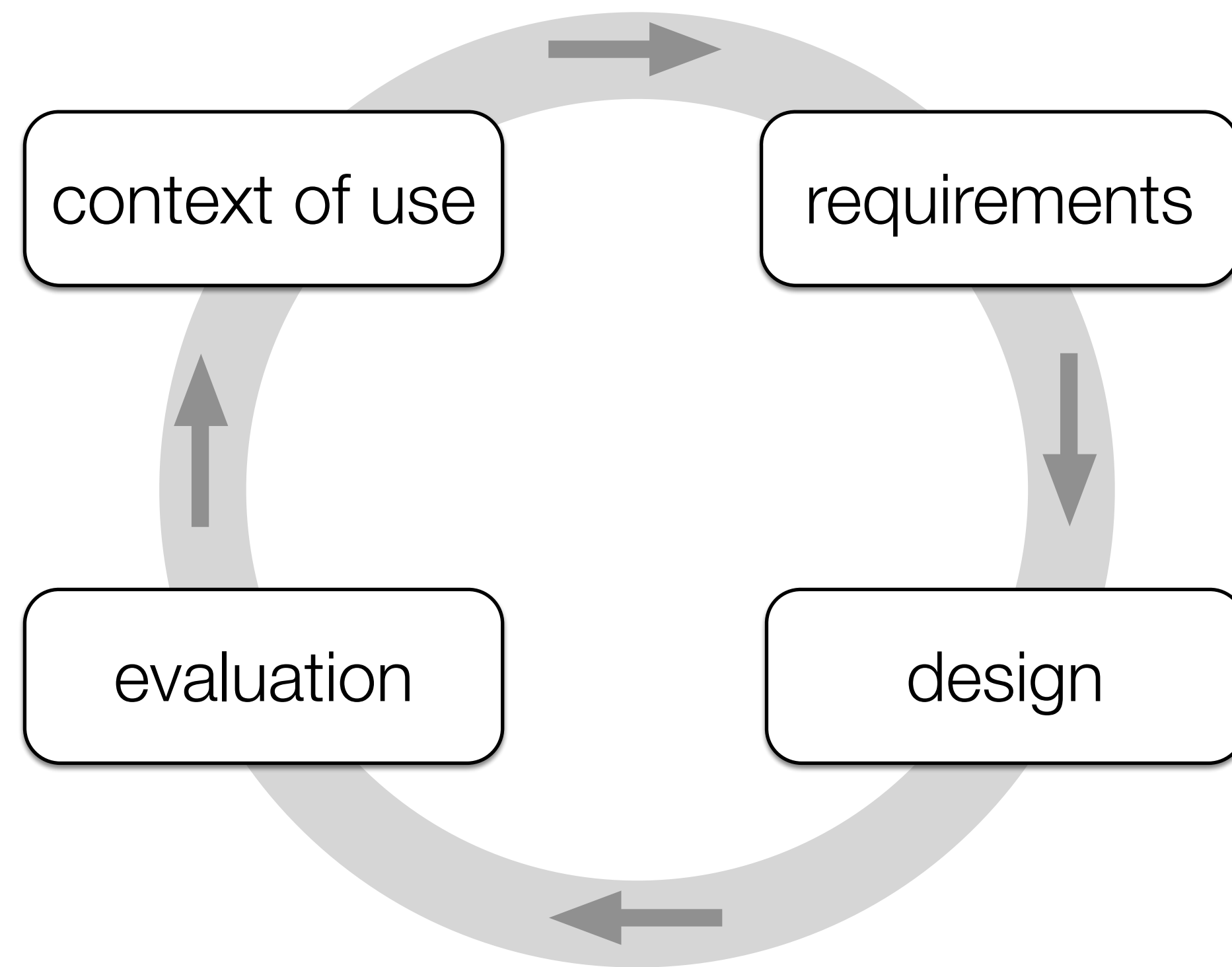
though **the path** from  
initial prototype to  
adoption **may be**  
**long...**

...you can **learn what**  
**works...** and what  
**doesn't.**





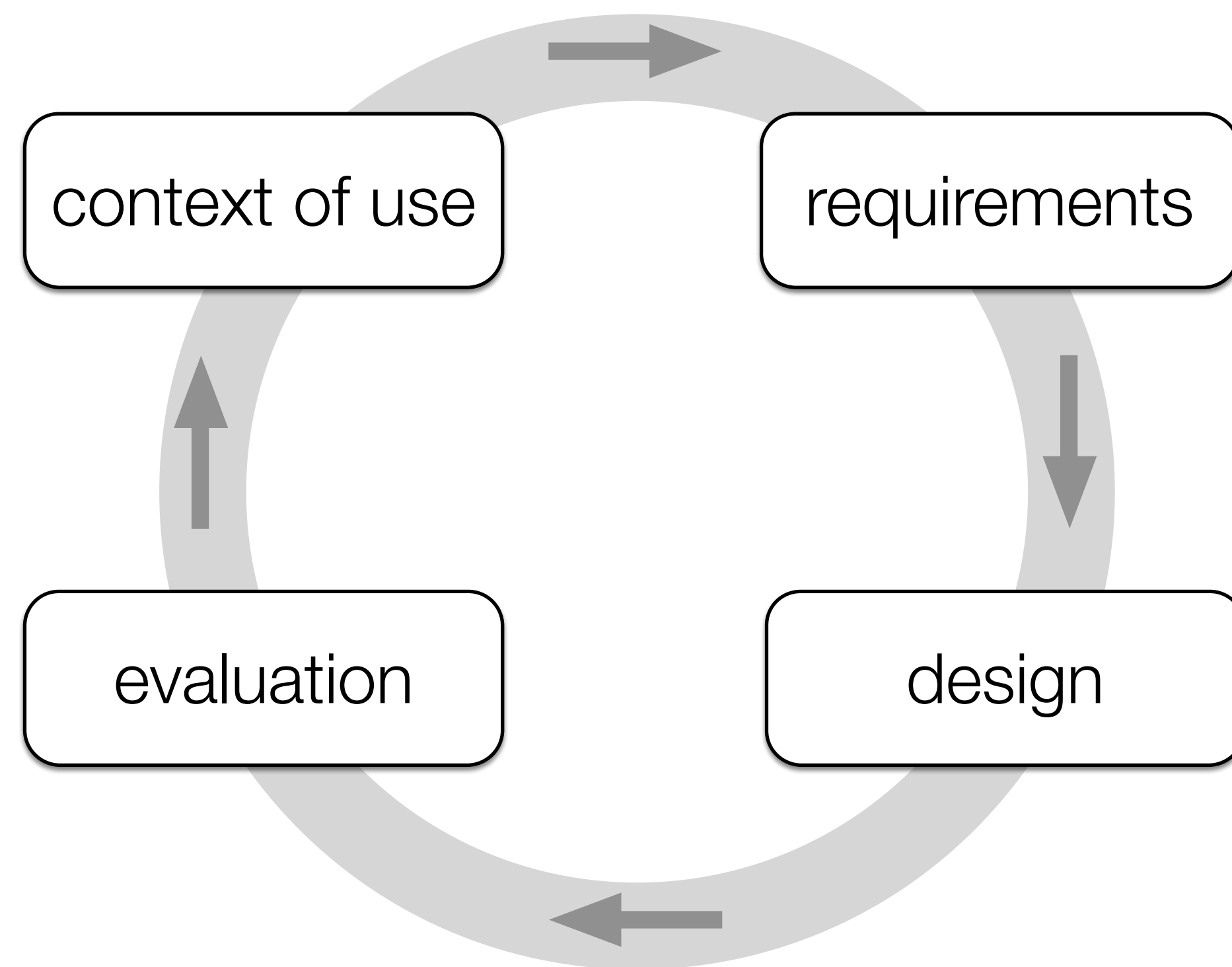
# HUMAN-CENTRED VIS DESIGN



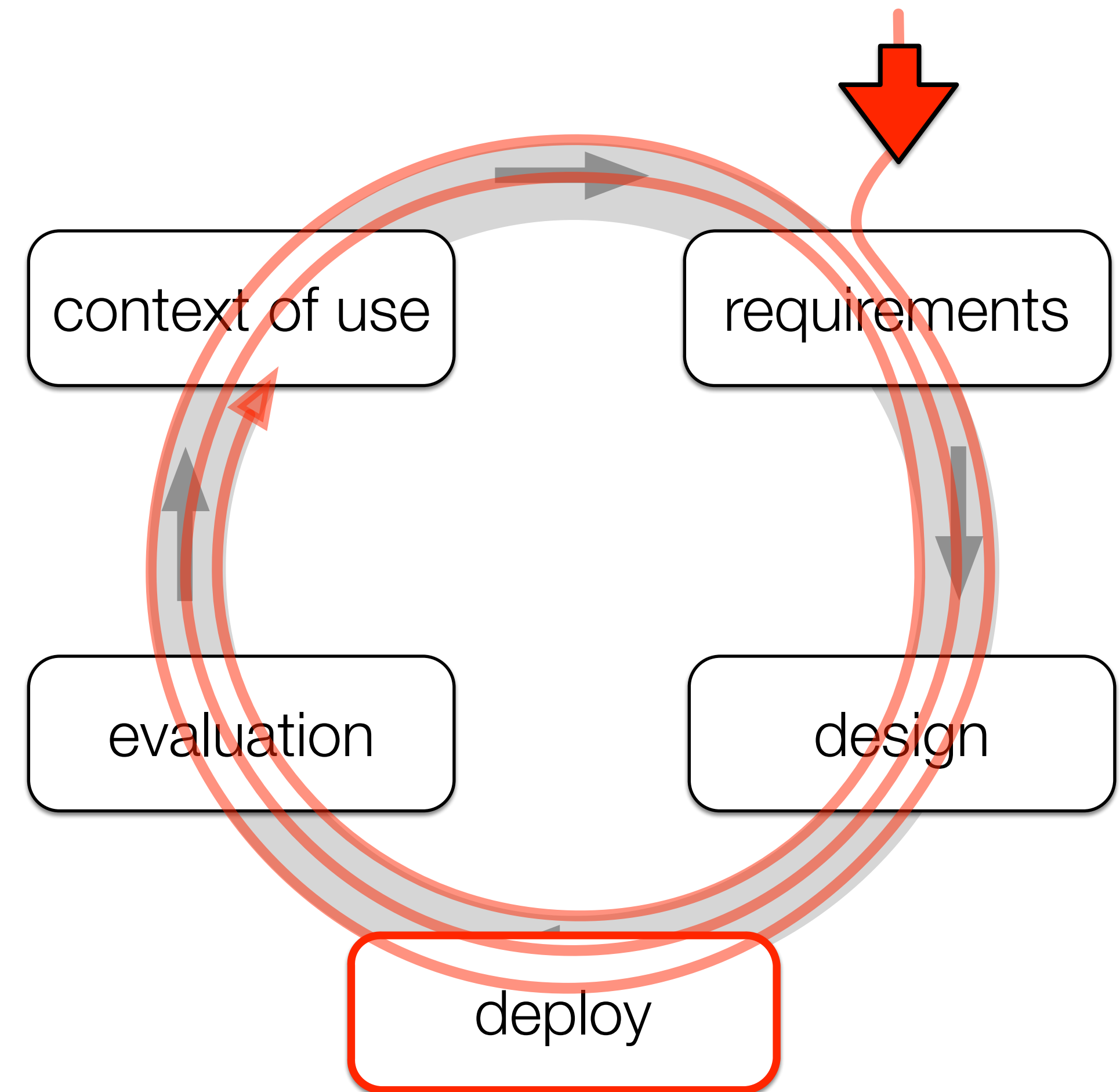
Lloyd & Dykes (IEEE TVCG / InfoVis 2011)



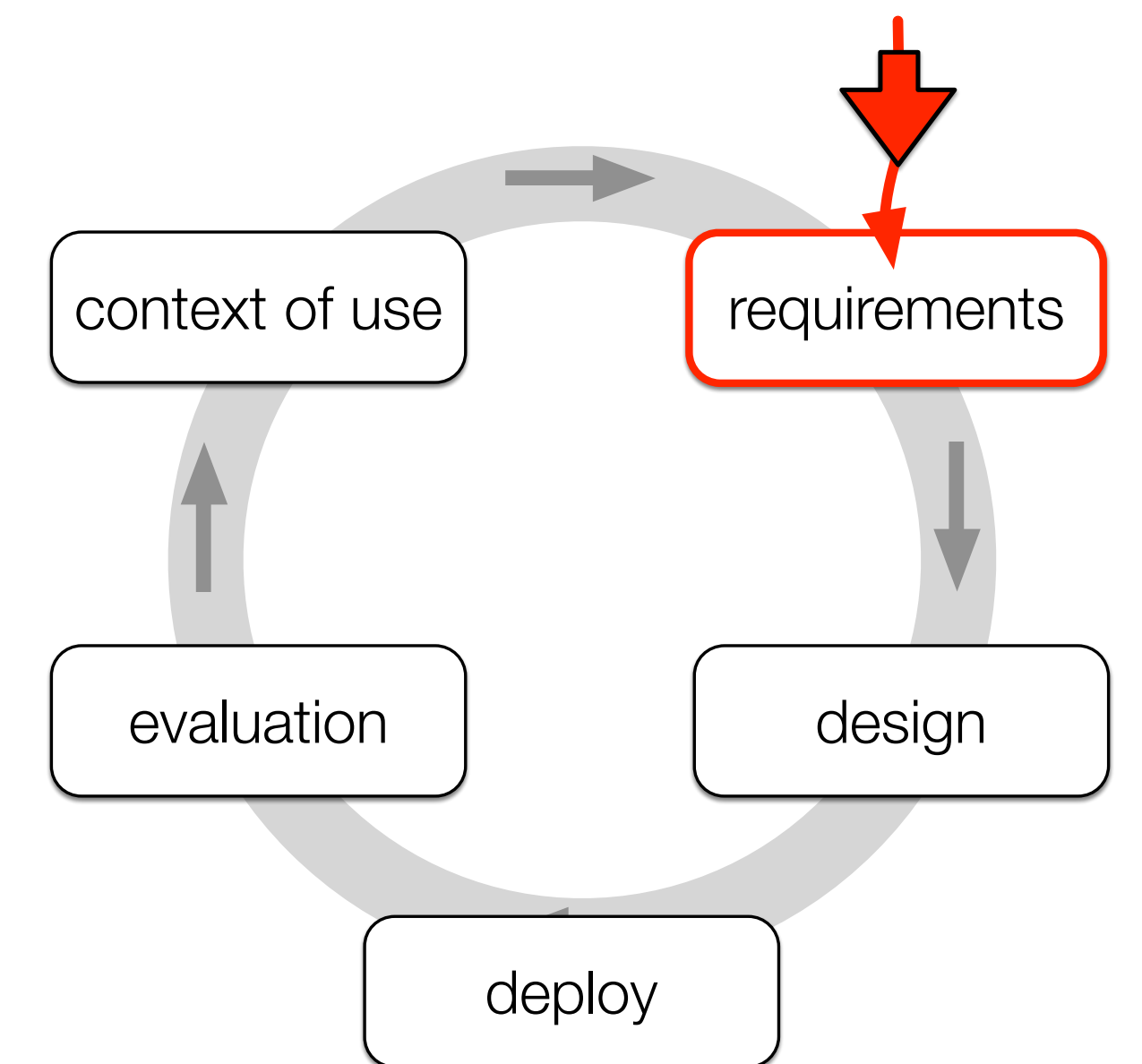
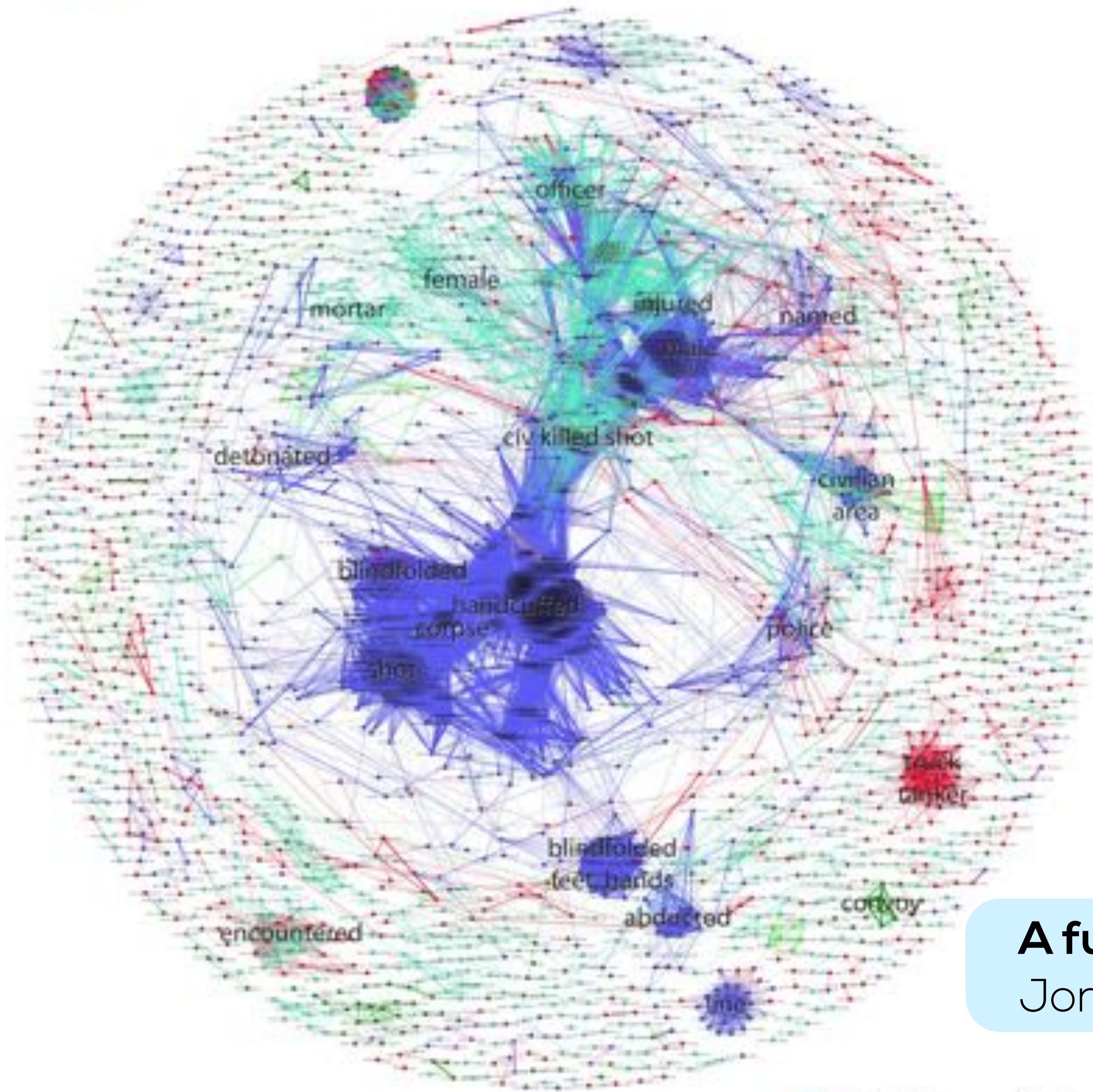
# HUMAN-CENTRED VIS DESIGN



Lloyd & Dykes (IEEE TVCG / InfoVis 2011)



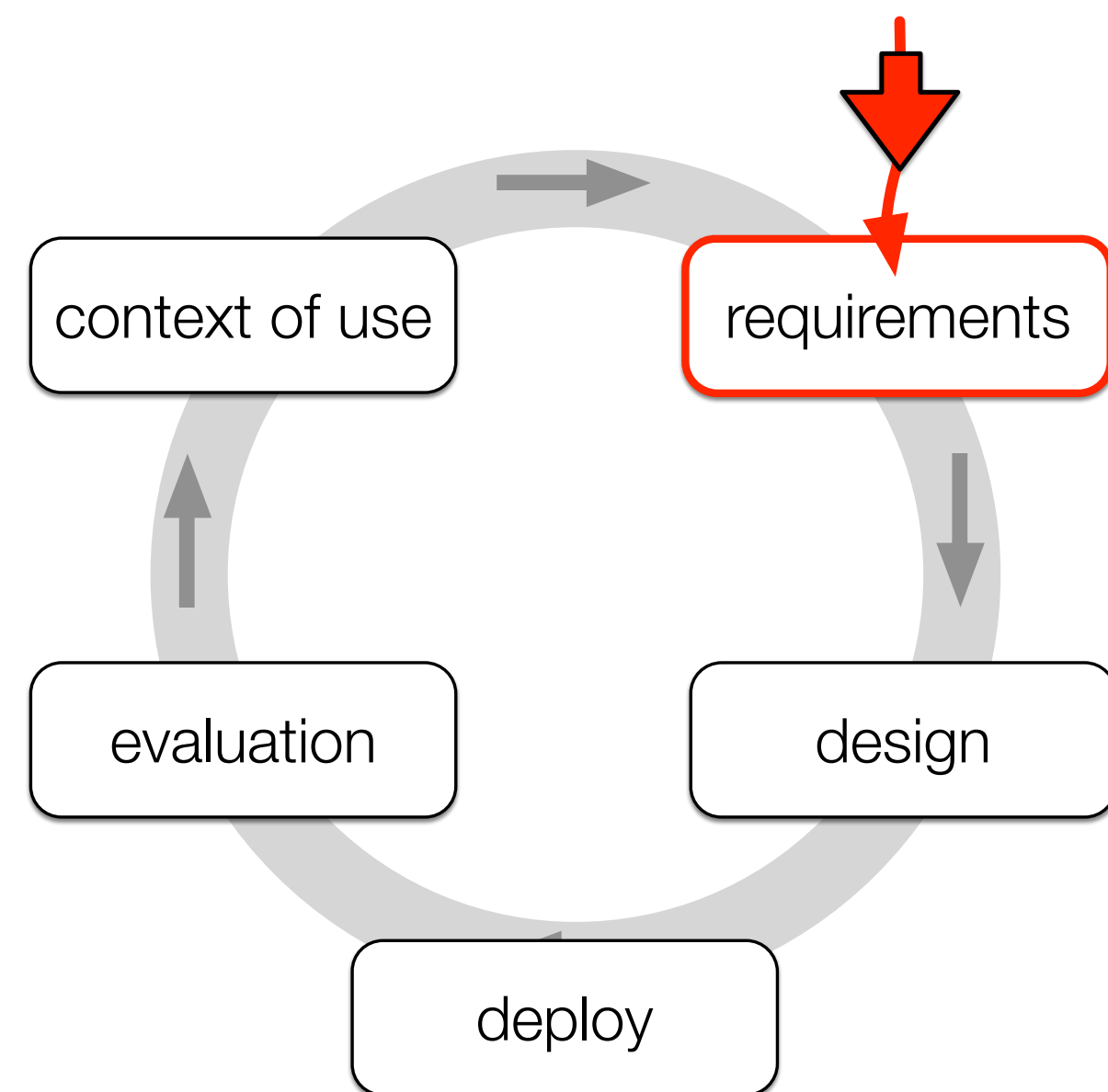
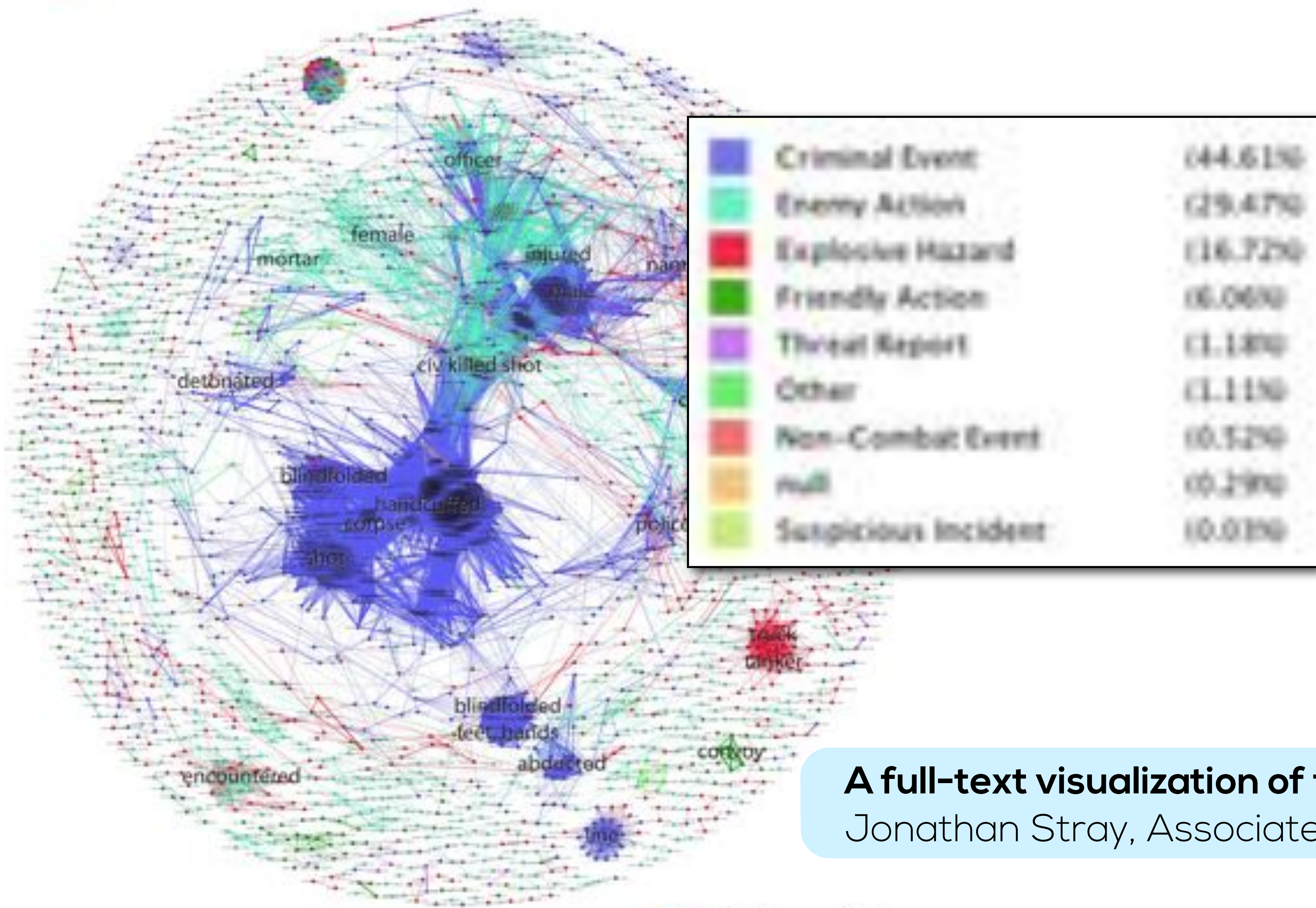




**A full-text visualization of the Iraq War Logs.**  
Jonathan Stray, Associated Press (Dec. 2010)



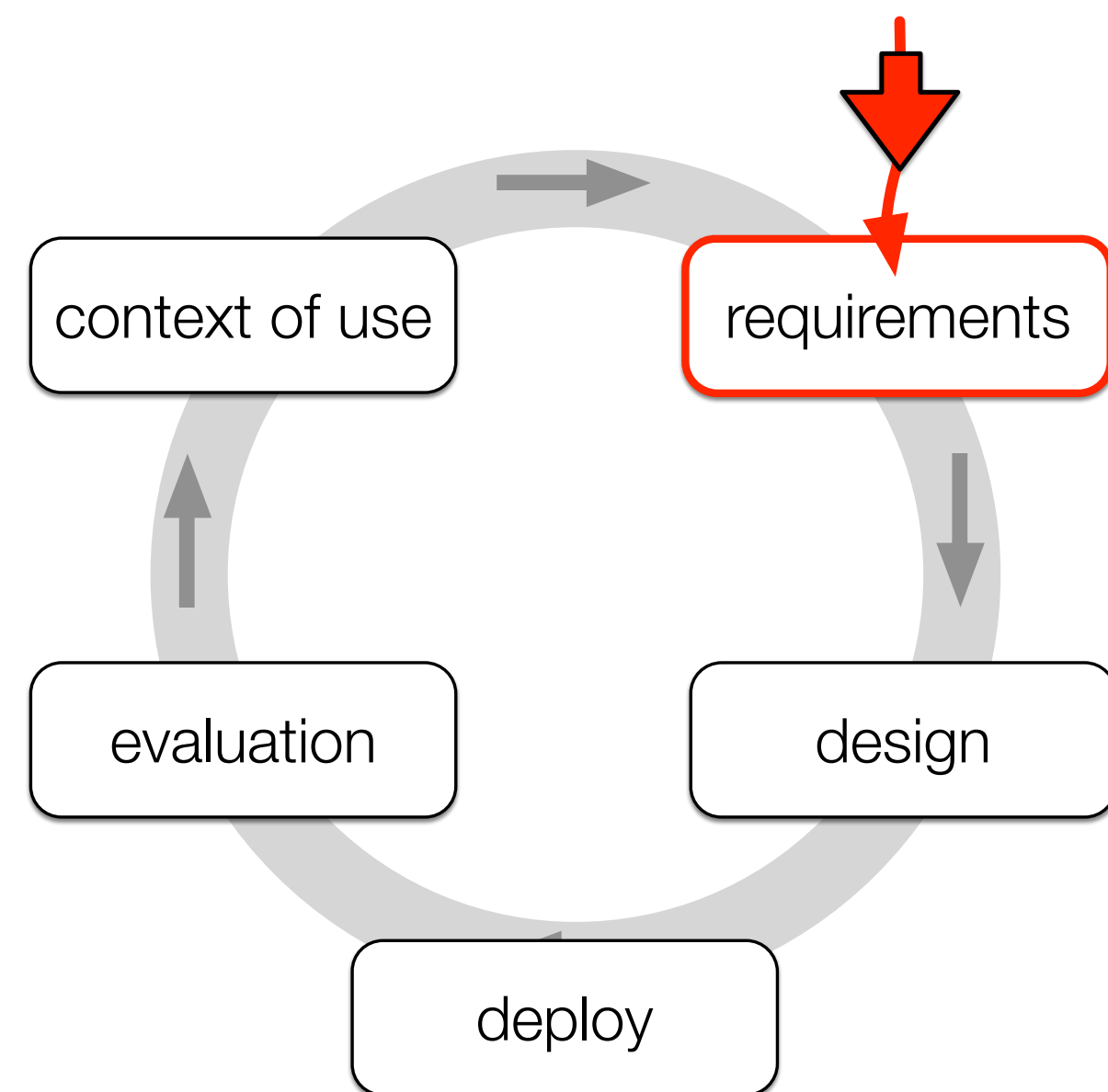
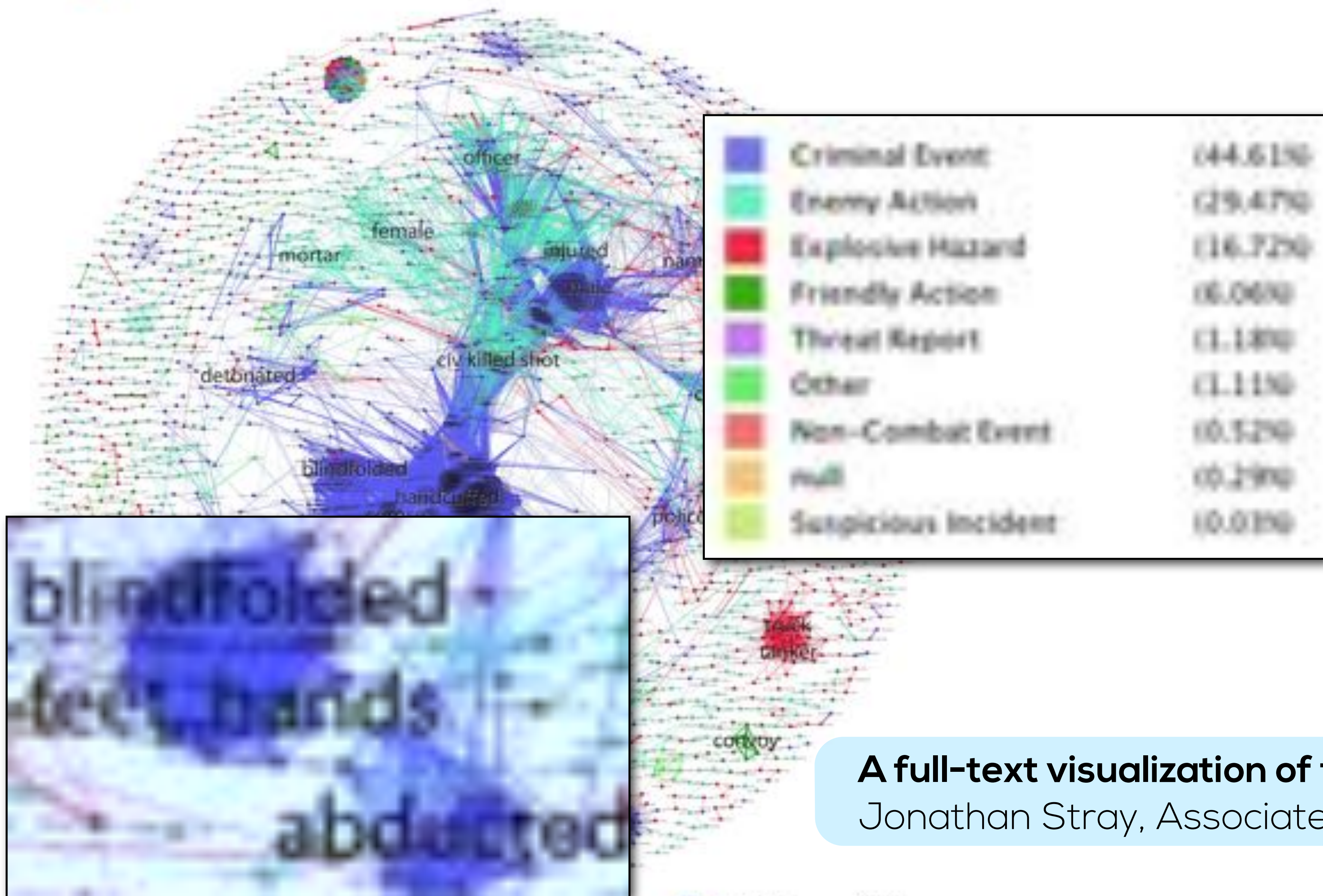




**A full-text visualization of the Iraq War Logs.**  
Jonathan Stray, Associated Press (Dec. 2010)



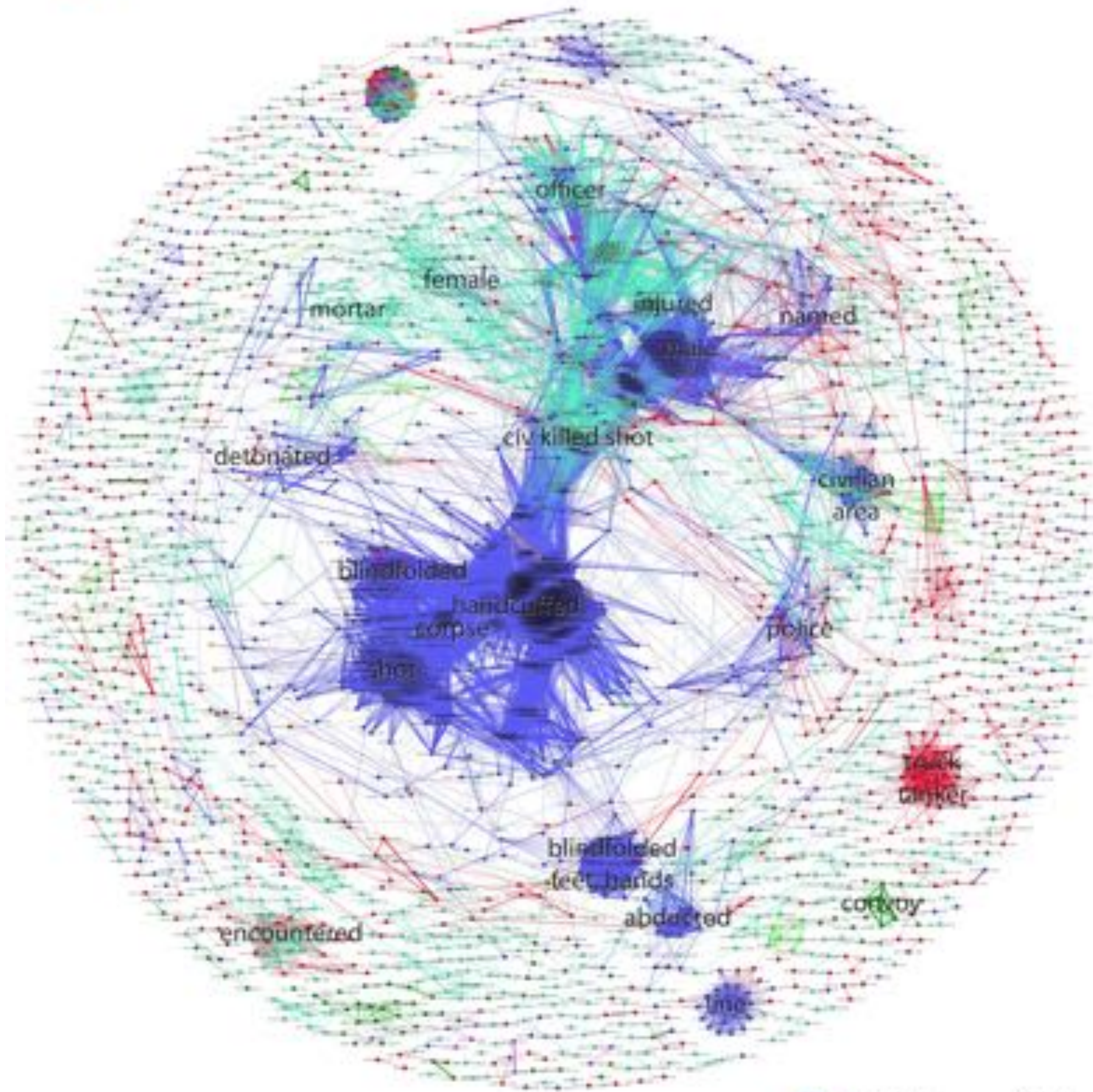




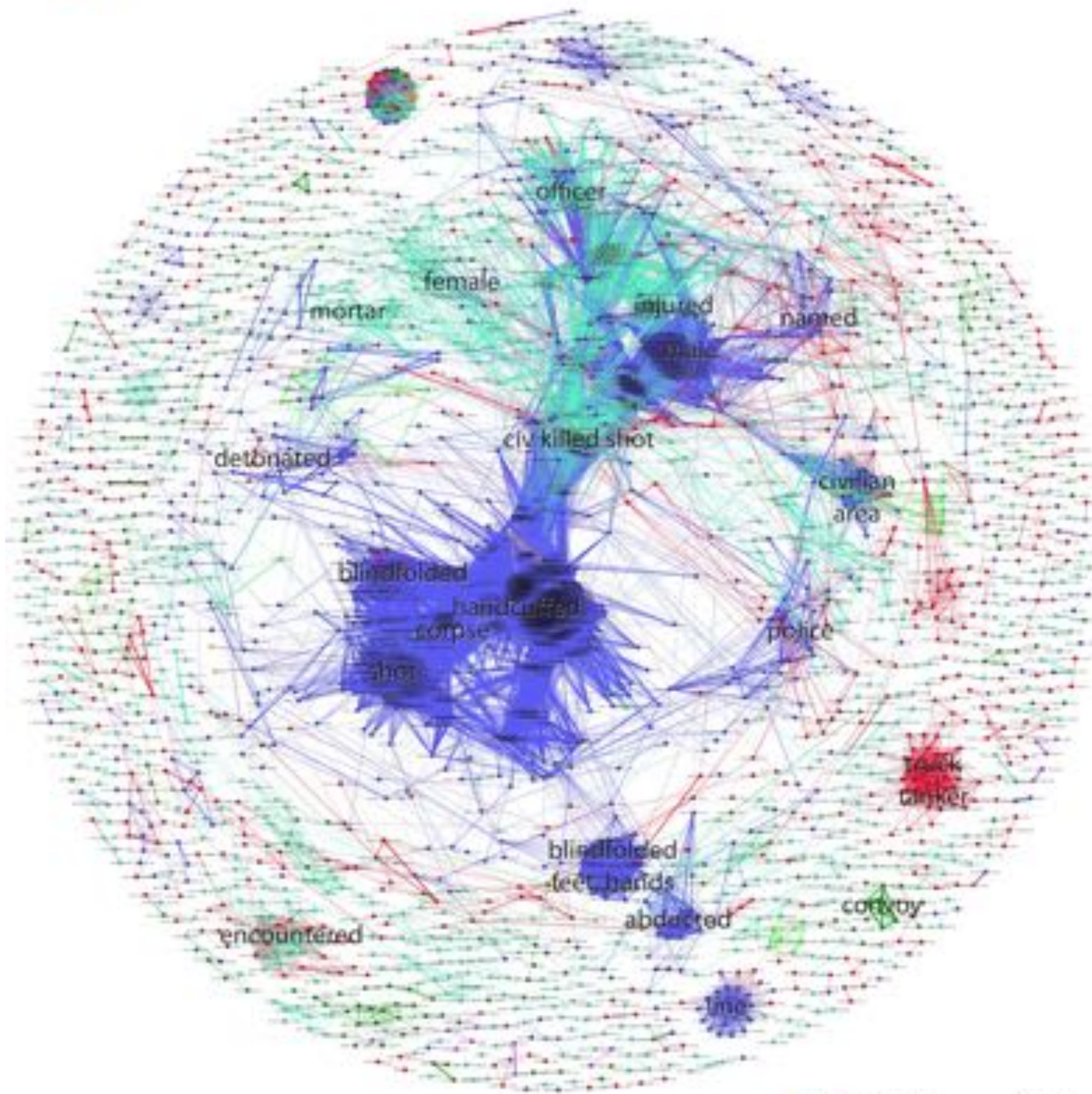
**A full-text visualization of the Iraq War Logs.**  
Jonathan Stray, Associated Press (Dec. 2010)





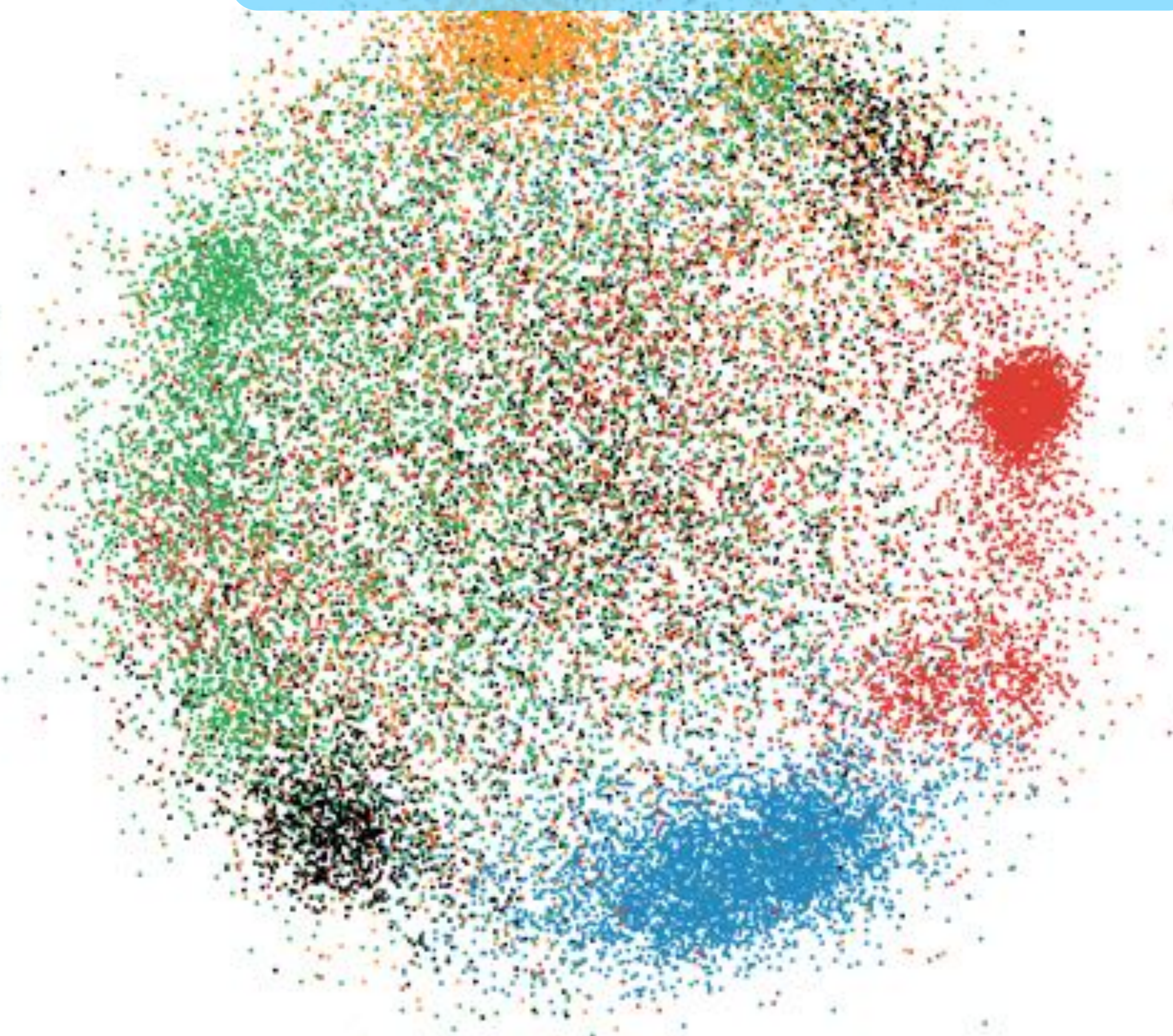






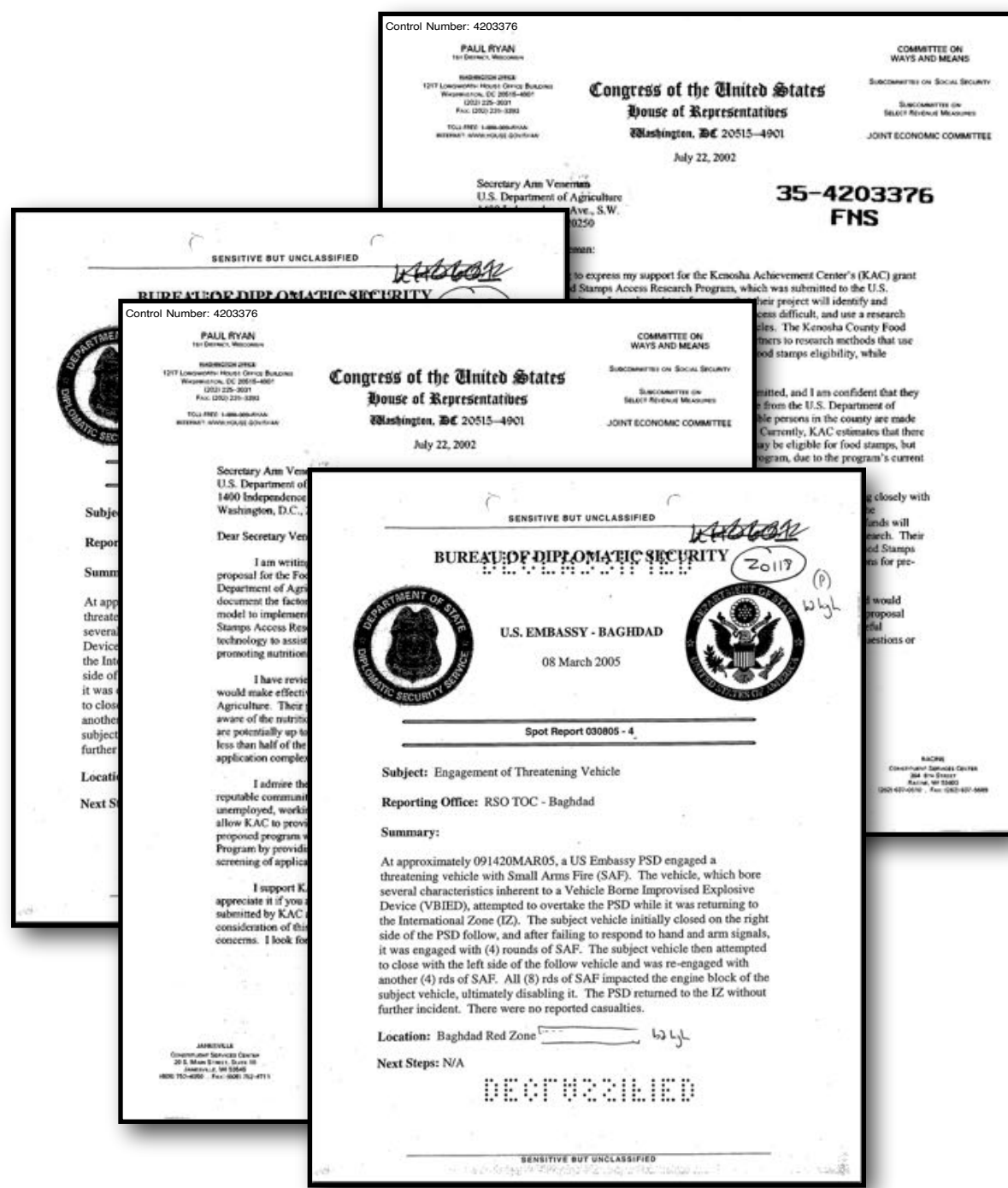
# Glimmer MDS: fast & scalable DR

Ingram et al. InfoVis 2009





# DATA ABSTRACTIONS

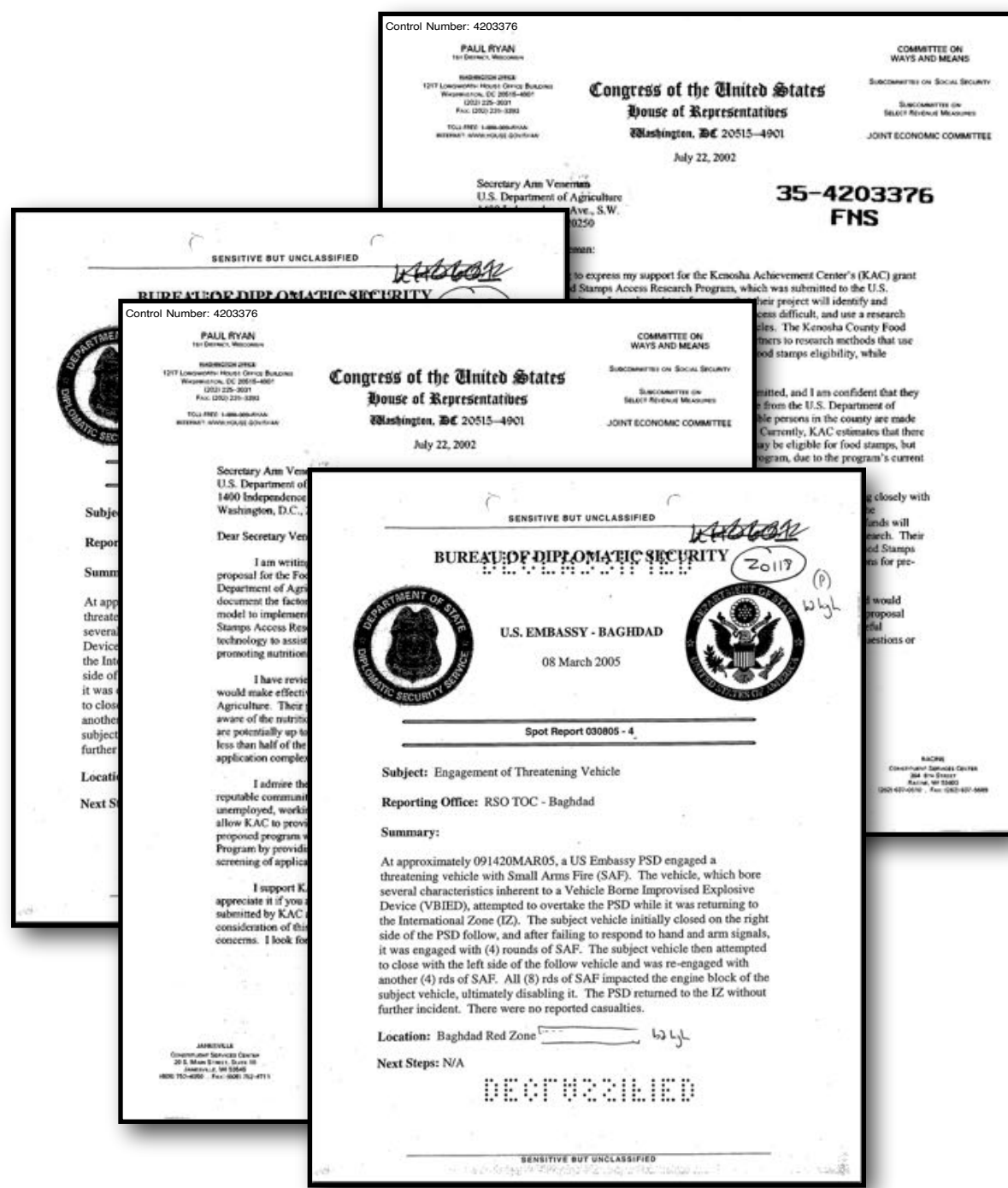


TF-IDF  
vectors





# DATA ABSTRACTIONS



TF-IDF  
vectors

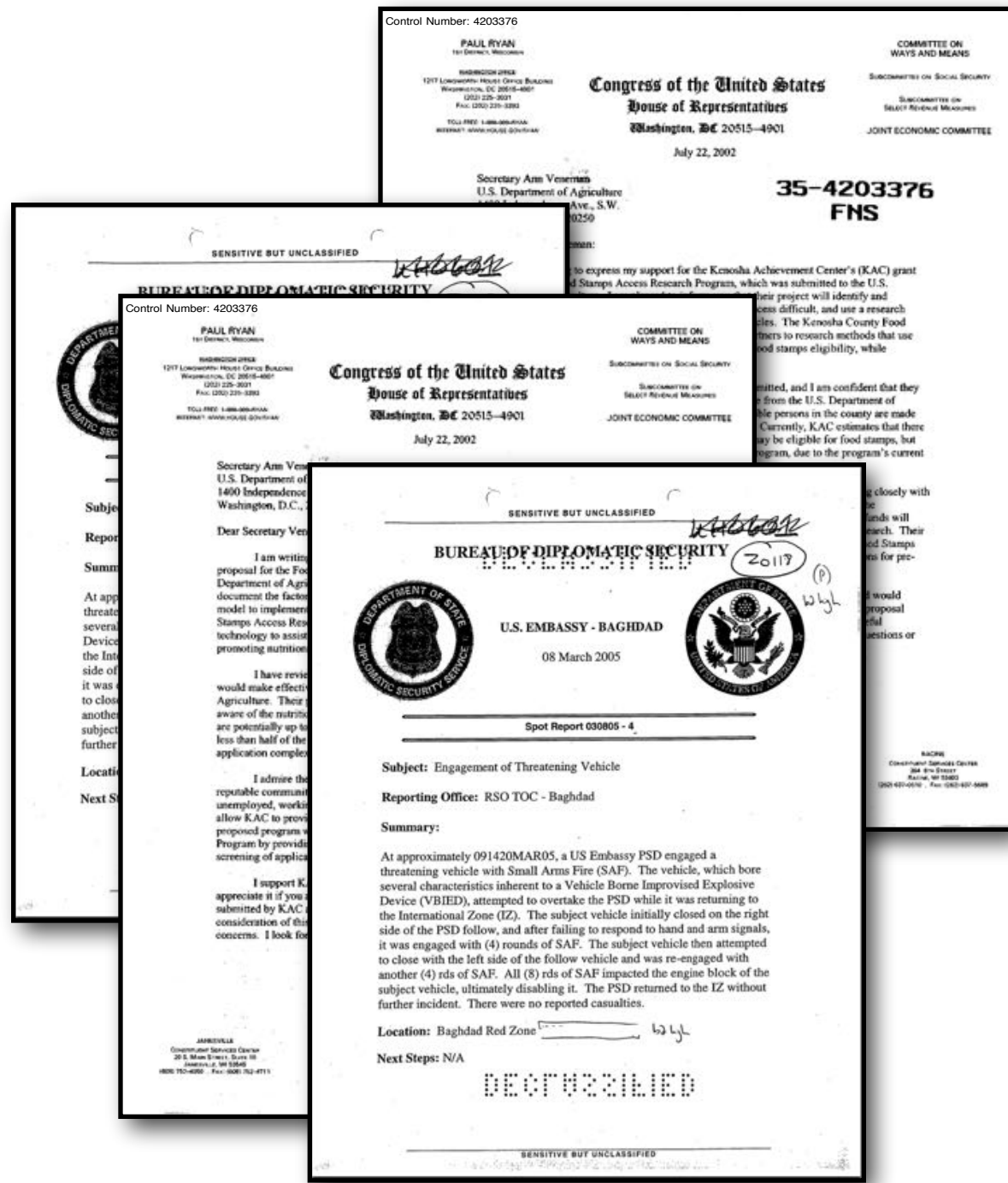
no document metadata!





# DATA ABSTRACTIONS

hierarchical  
clustering



TF-IDF  
vectors

dimensionality  
reduction

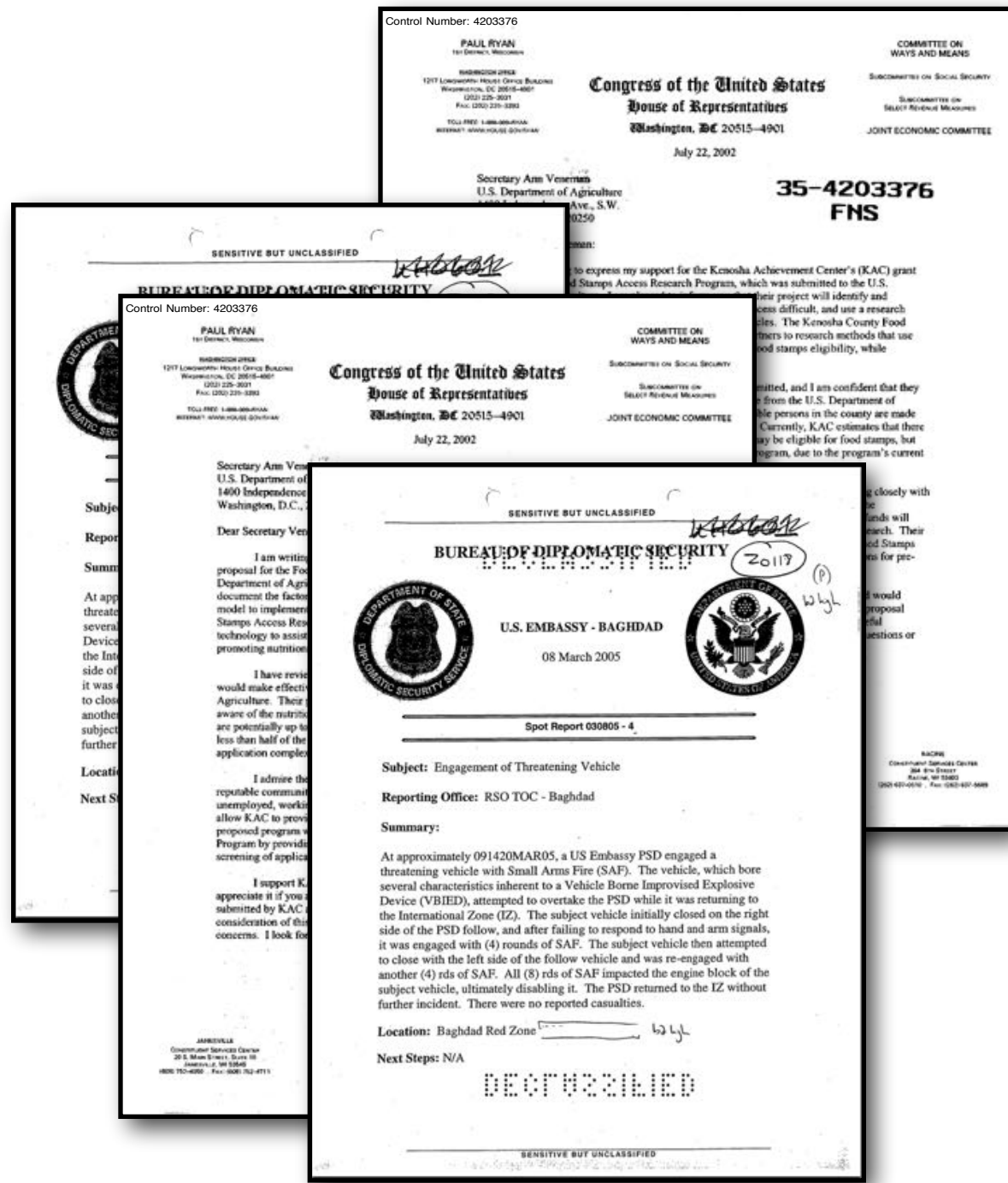
no document metadata!





# DATA ABSTRACTIONS

hierarchical  
clustering



TF-IDF  
vectors

+

user-  
defined  
tags

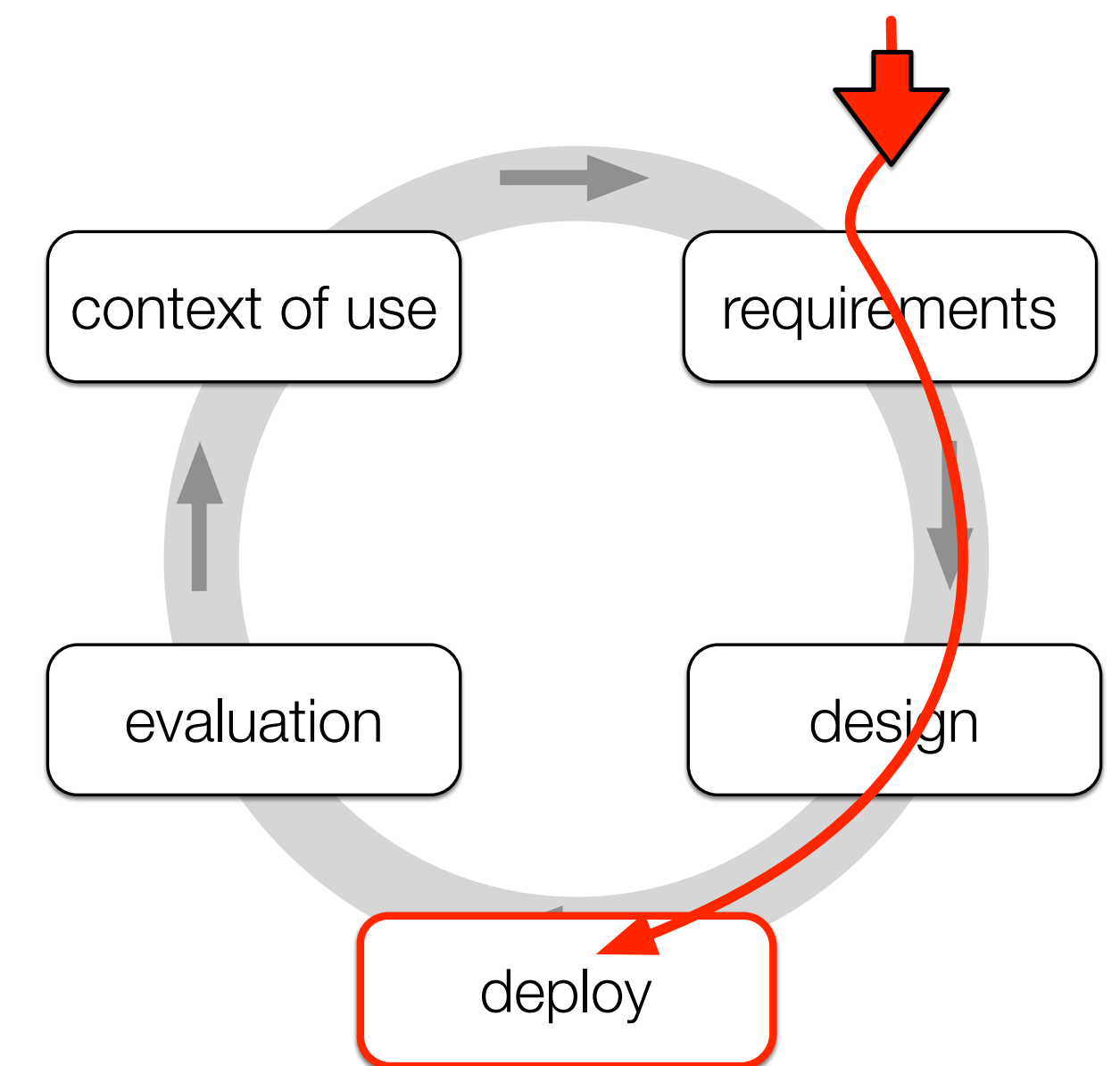
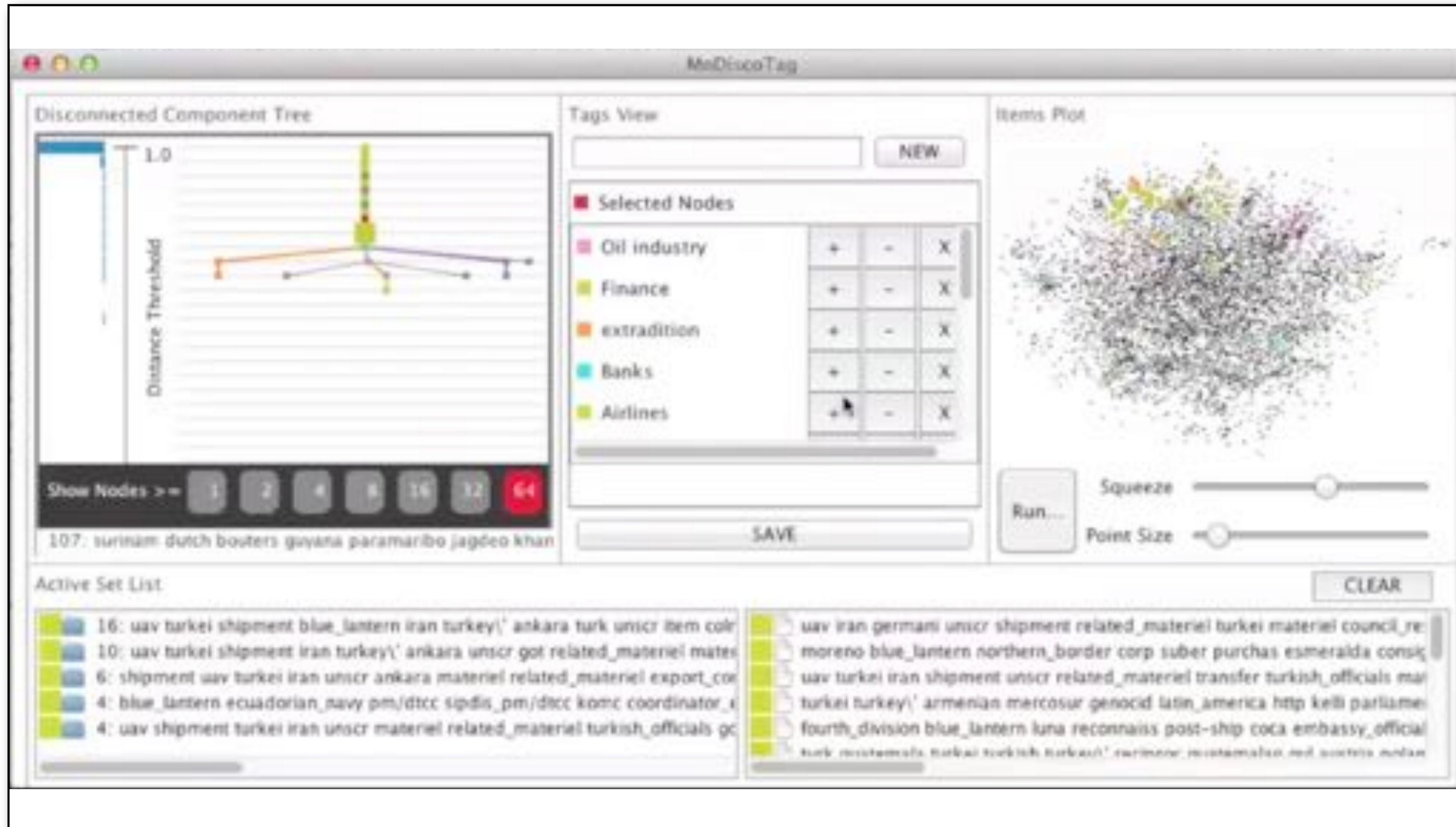
dimensionality  
reduction

no document metadata!





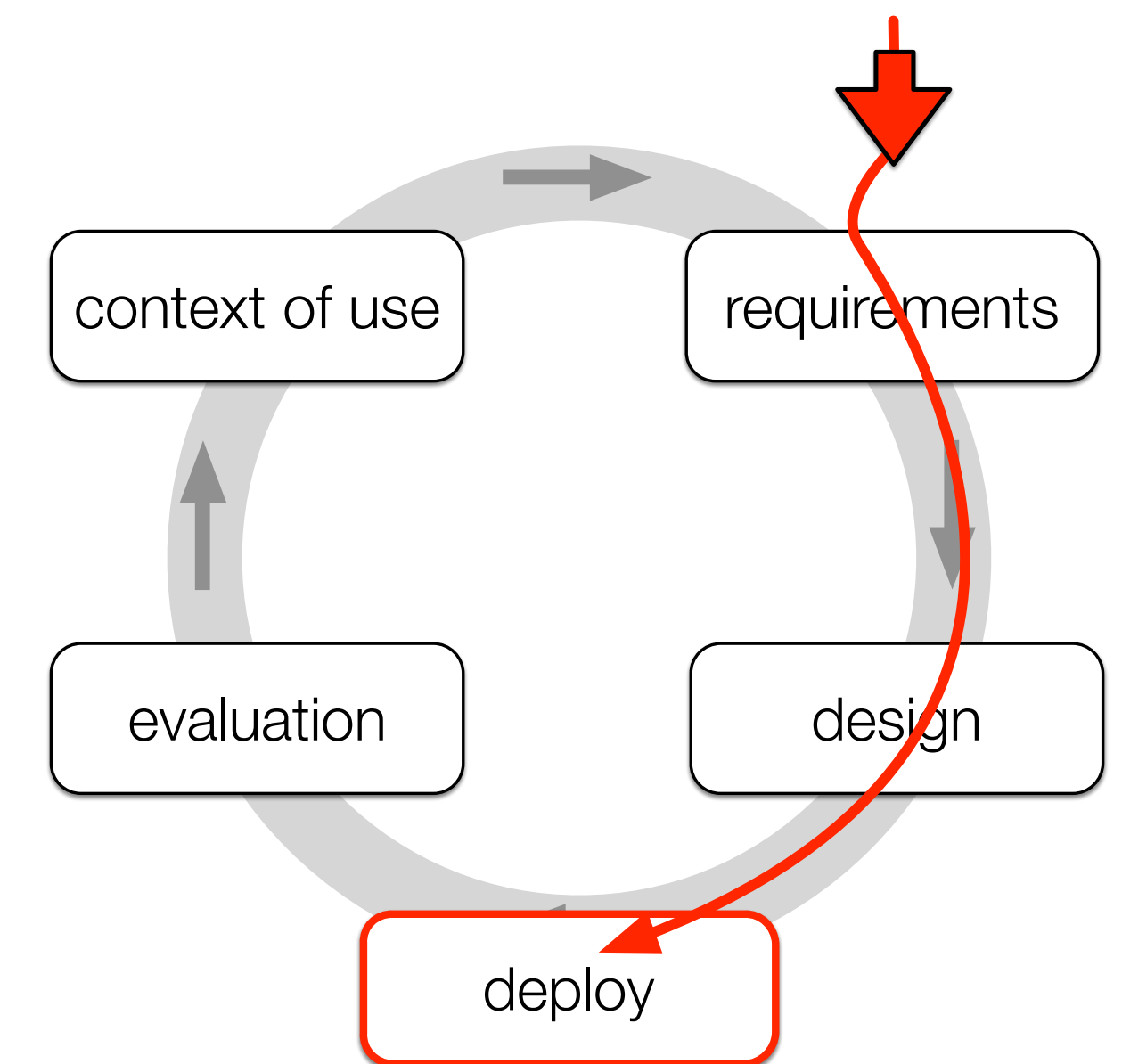
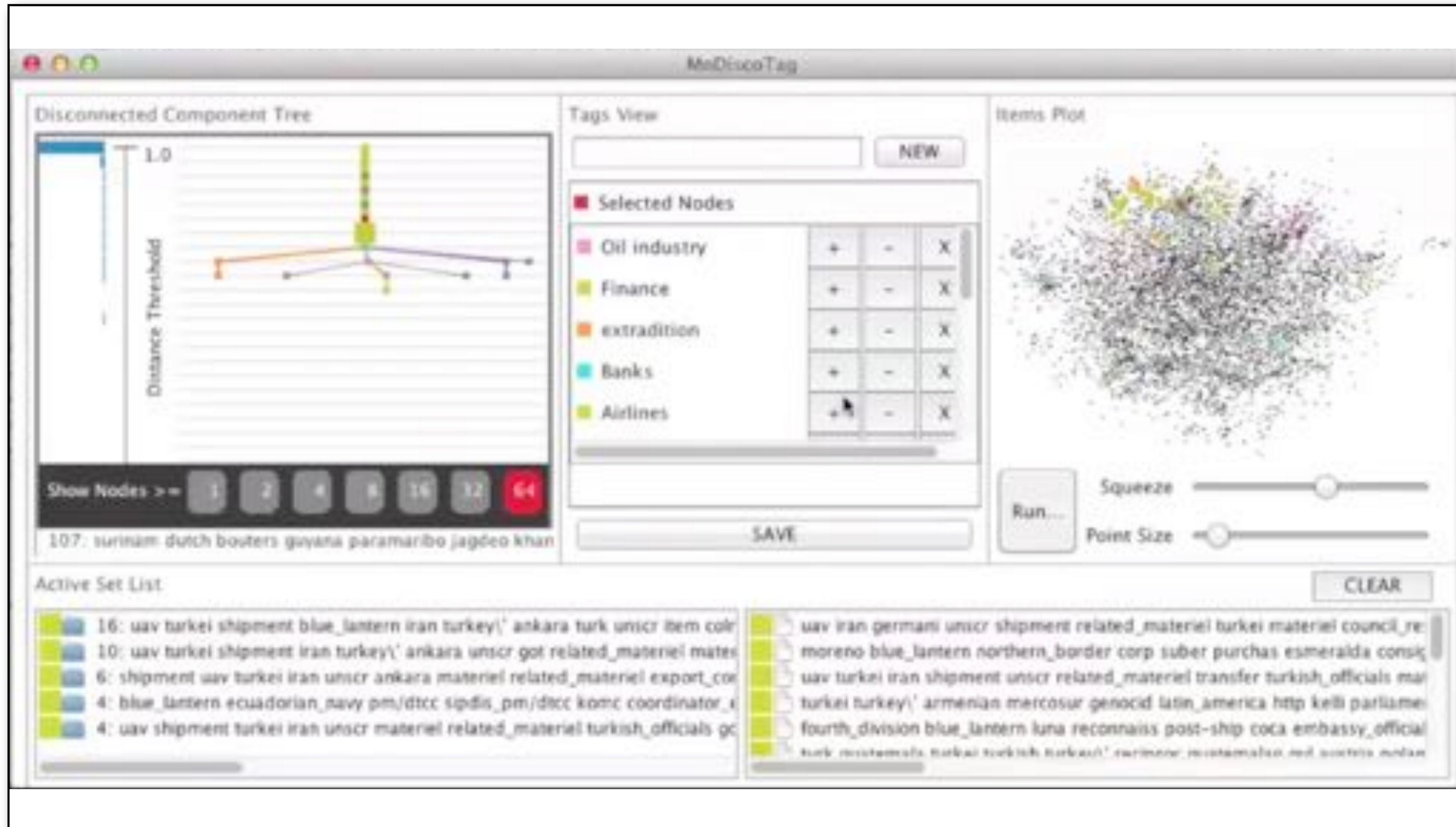
# OVERVIEW V.1 PROTOTYPE, 2011



**clockwise from TL:**  
tree visualization,  
tag controls,  
scatterplot,  
document list,  
document viewer,  
cluster list



# OVERVIEW V.1 PROTOTYPE, 2011



**clockwise from TL:**  
tree visualization,  
tag controls,  
scatterplot,  
document list,  
document viewer,  
cluster list

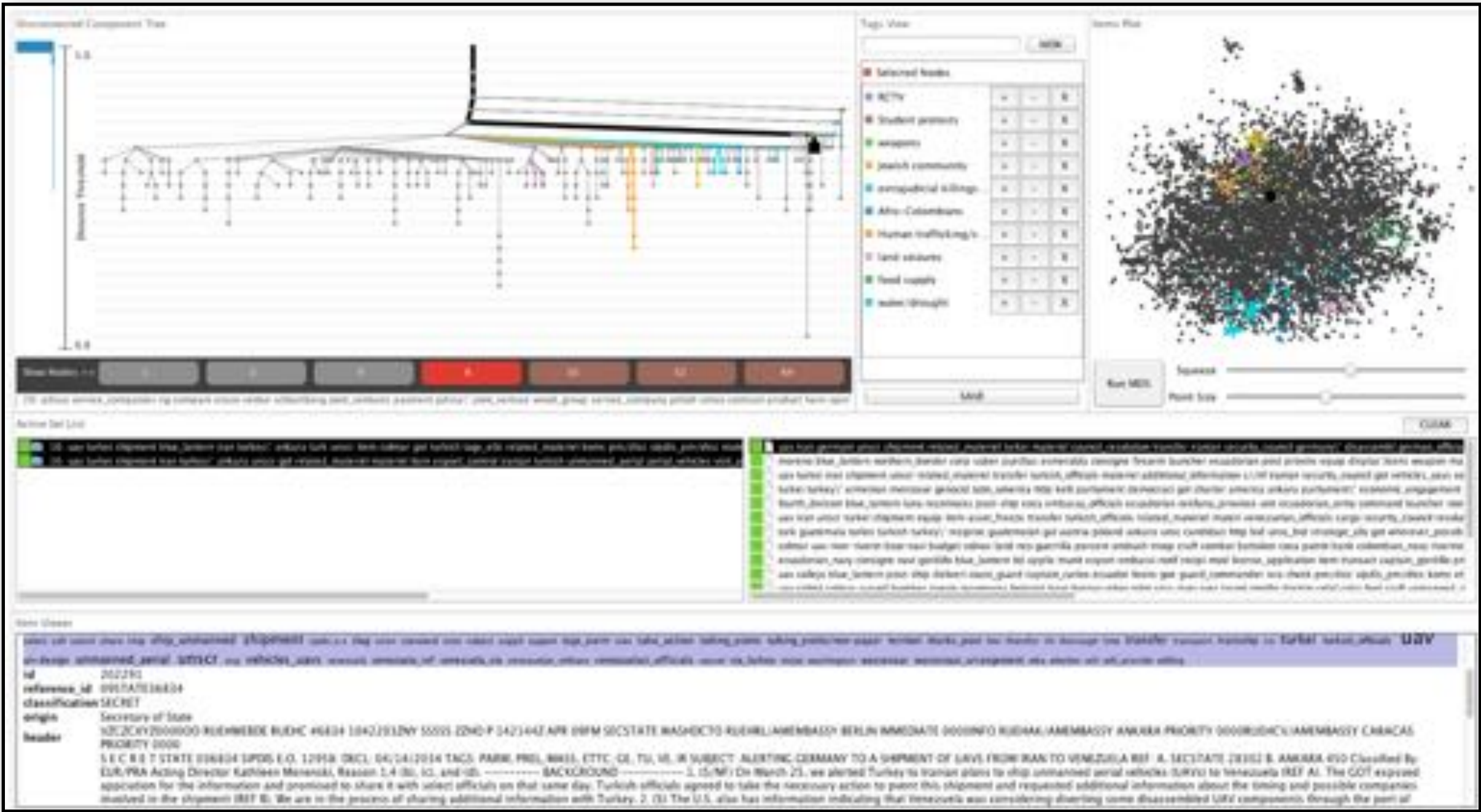
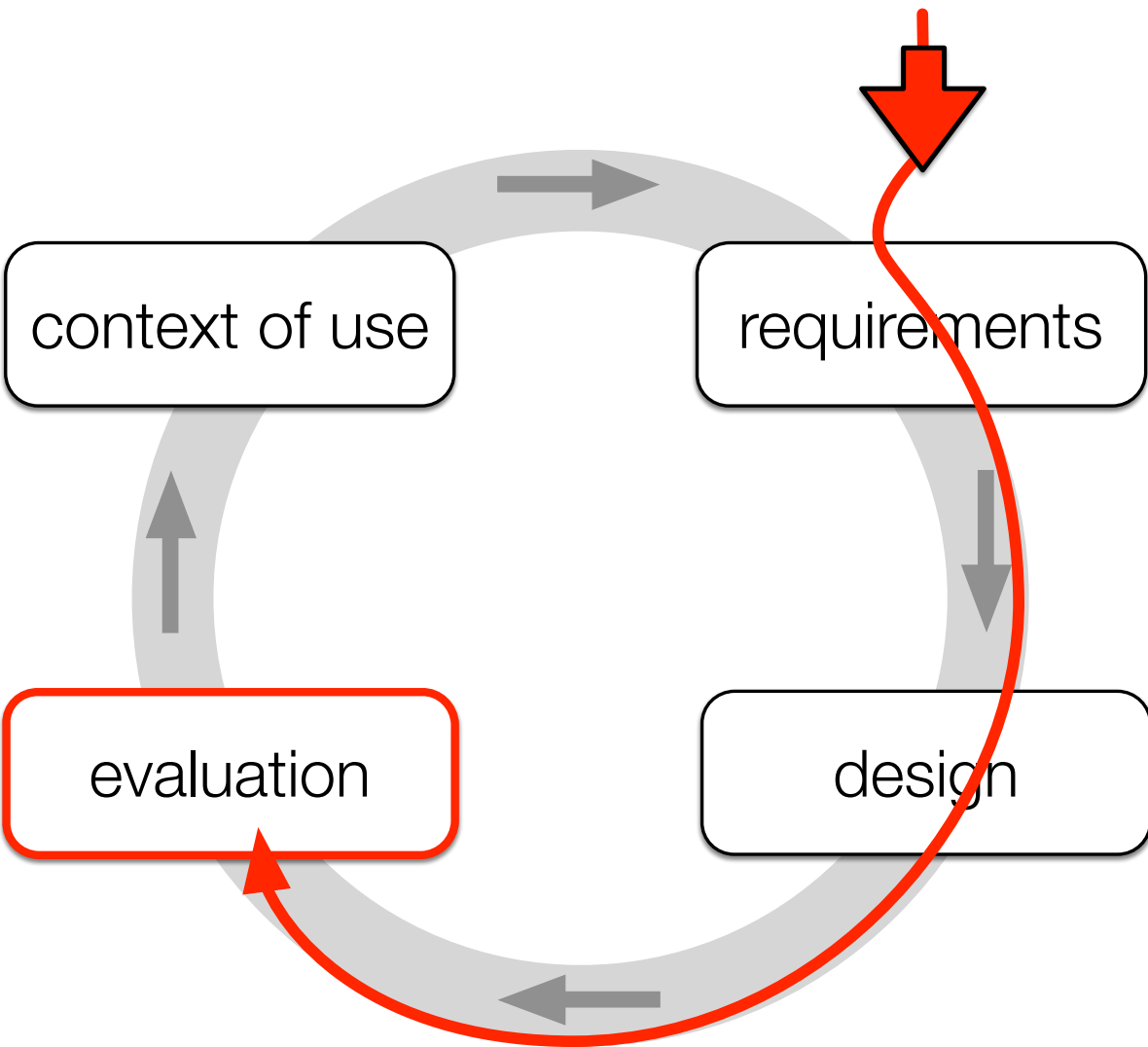


# “EXPLORING” THE WIKILEAKS CABLES

<http://bigstory.ap.org/article/venezuela-chavez-era-reporters-view>

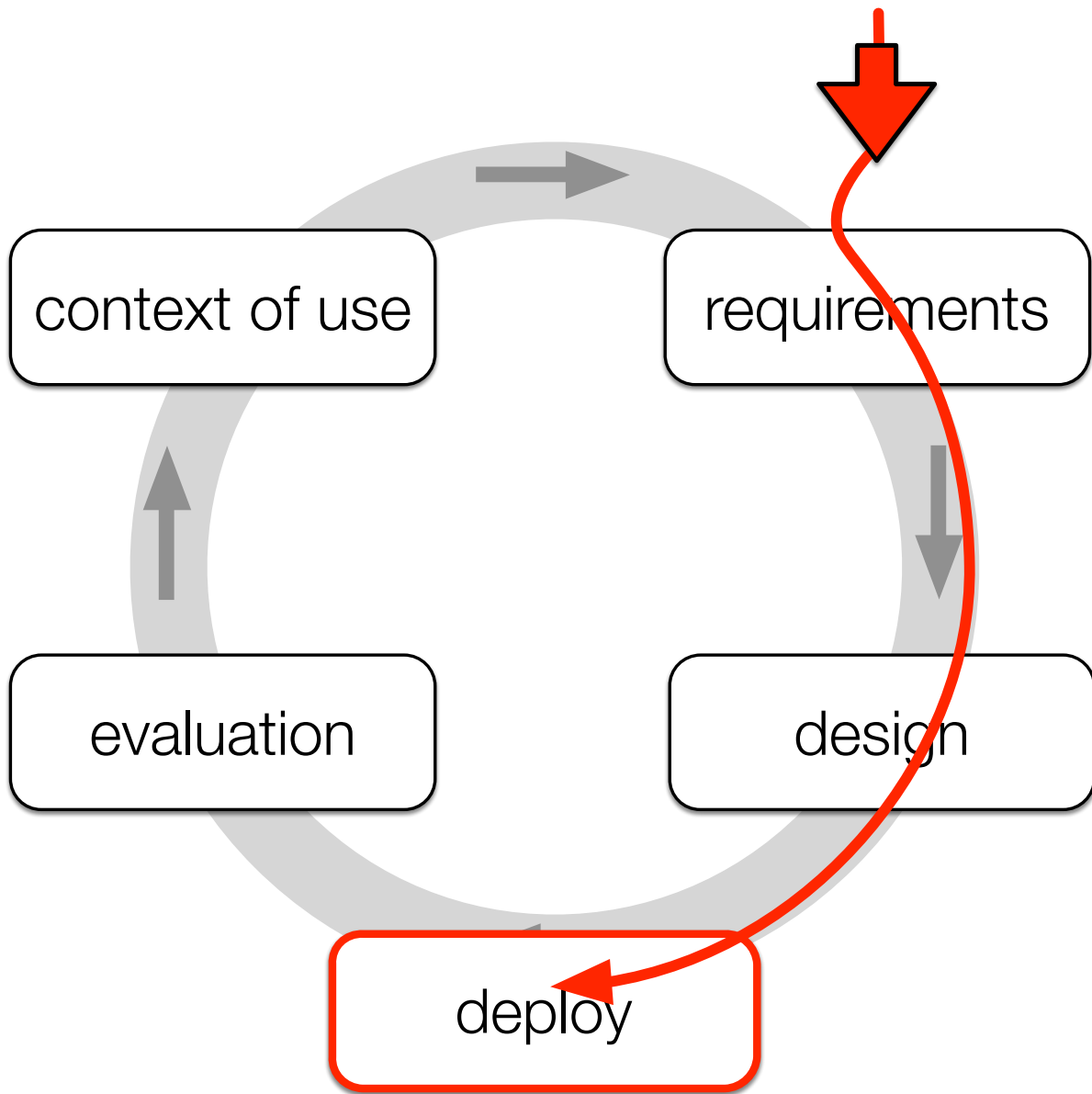
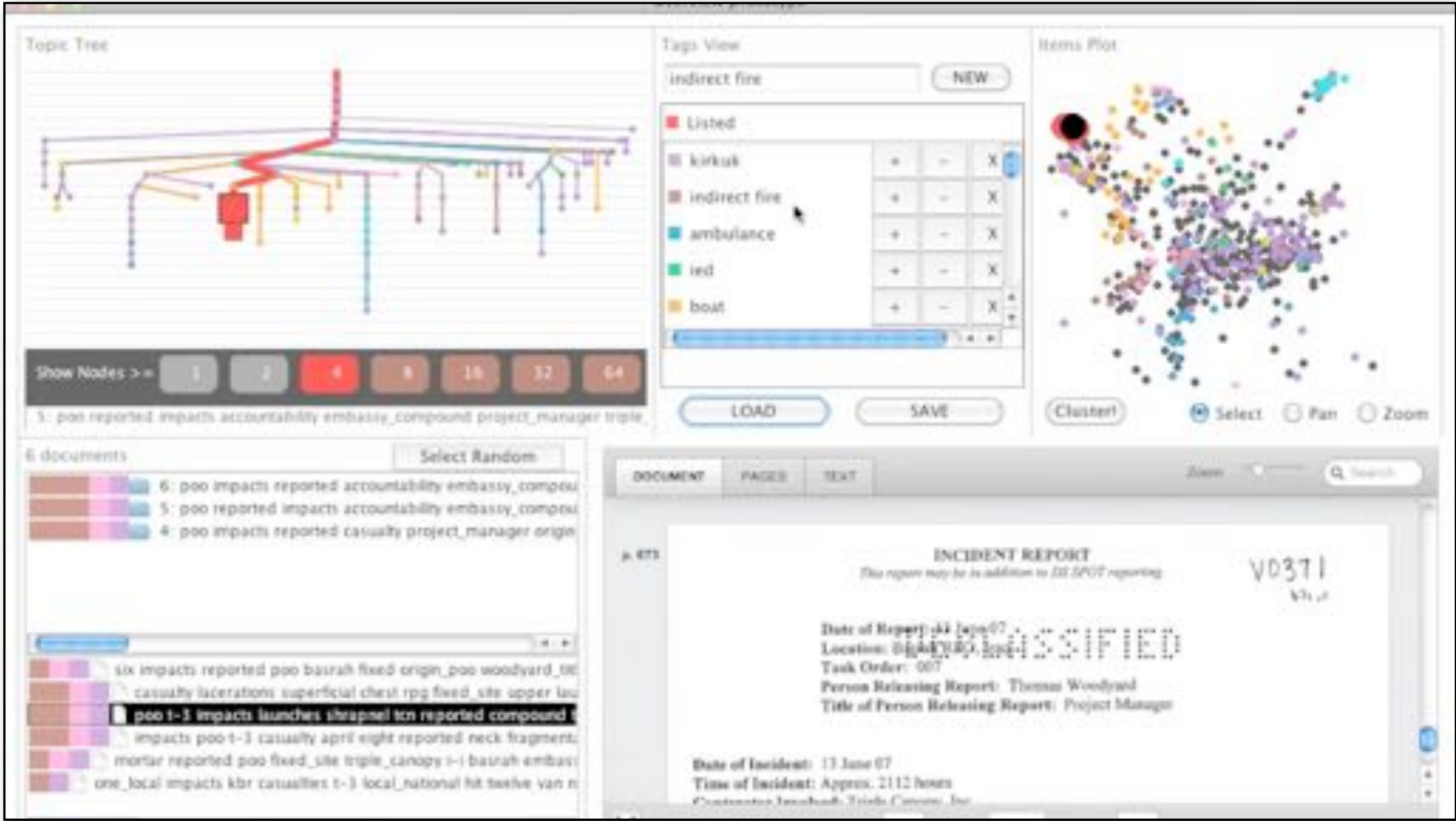


AP Caracas bureau chief  
**Ian James**  
(with Hugo Chavez, 2007)



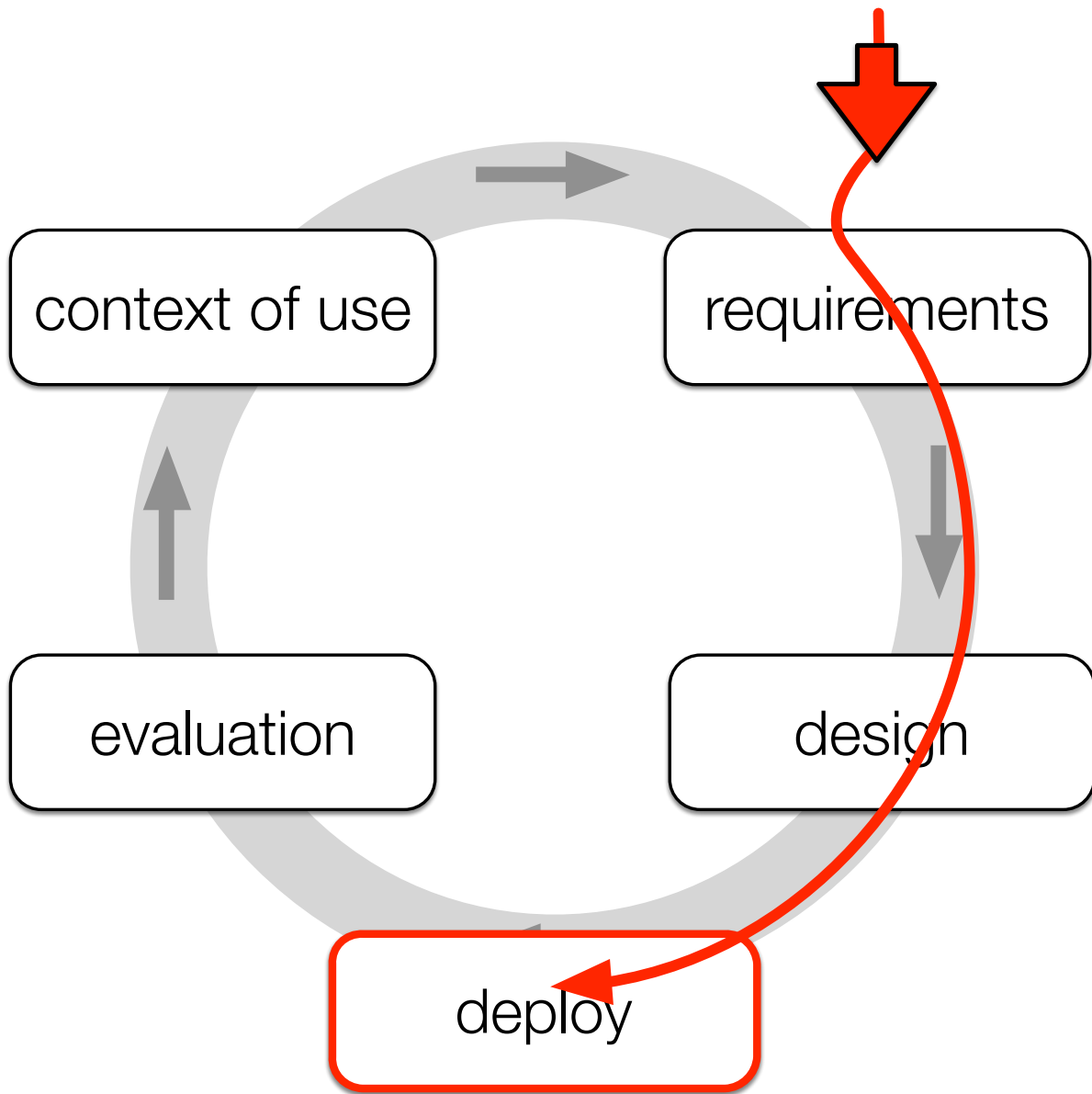
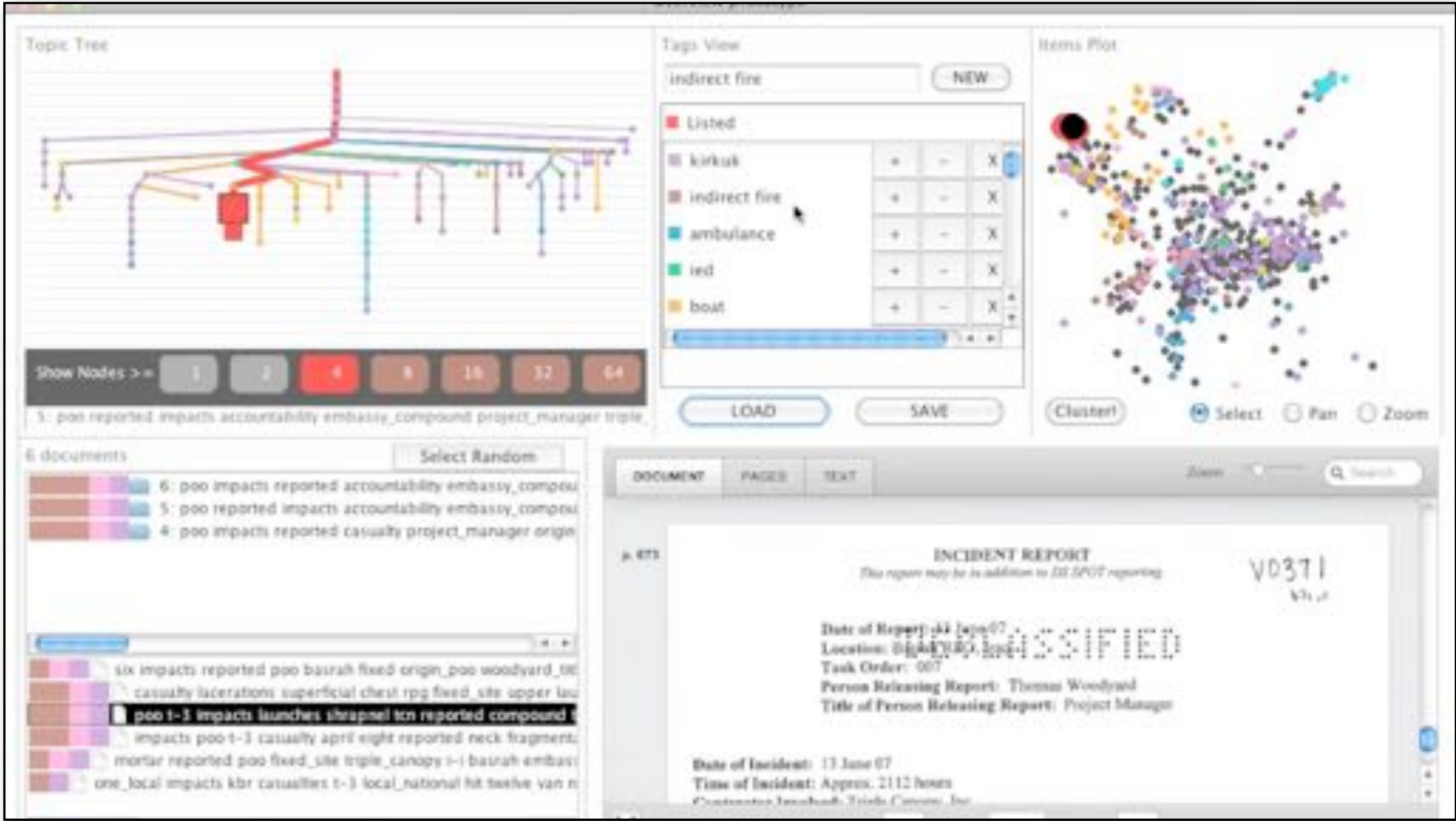


# OVERVIEW V.2 (DESKTOP APP), 2012



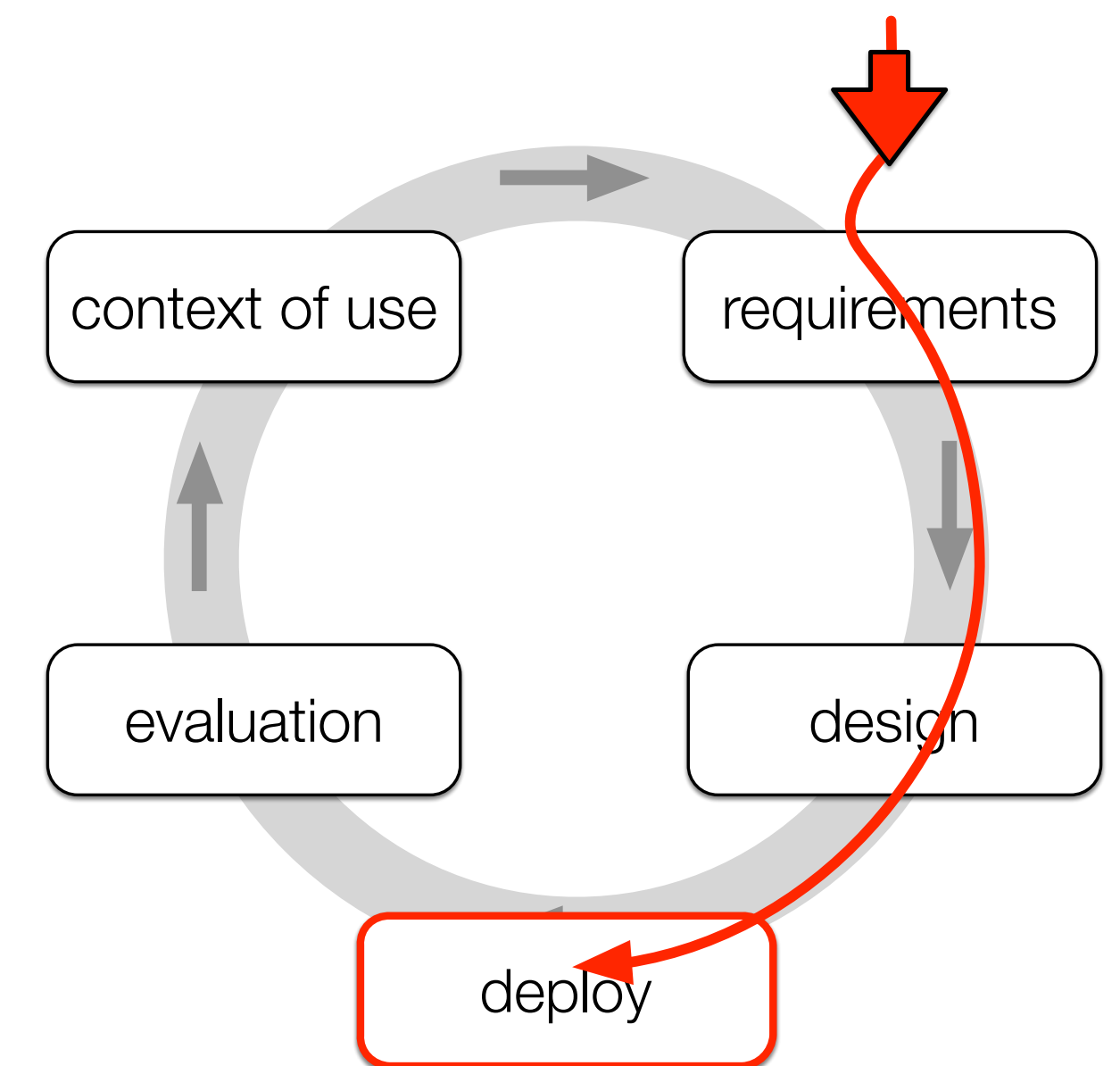
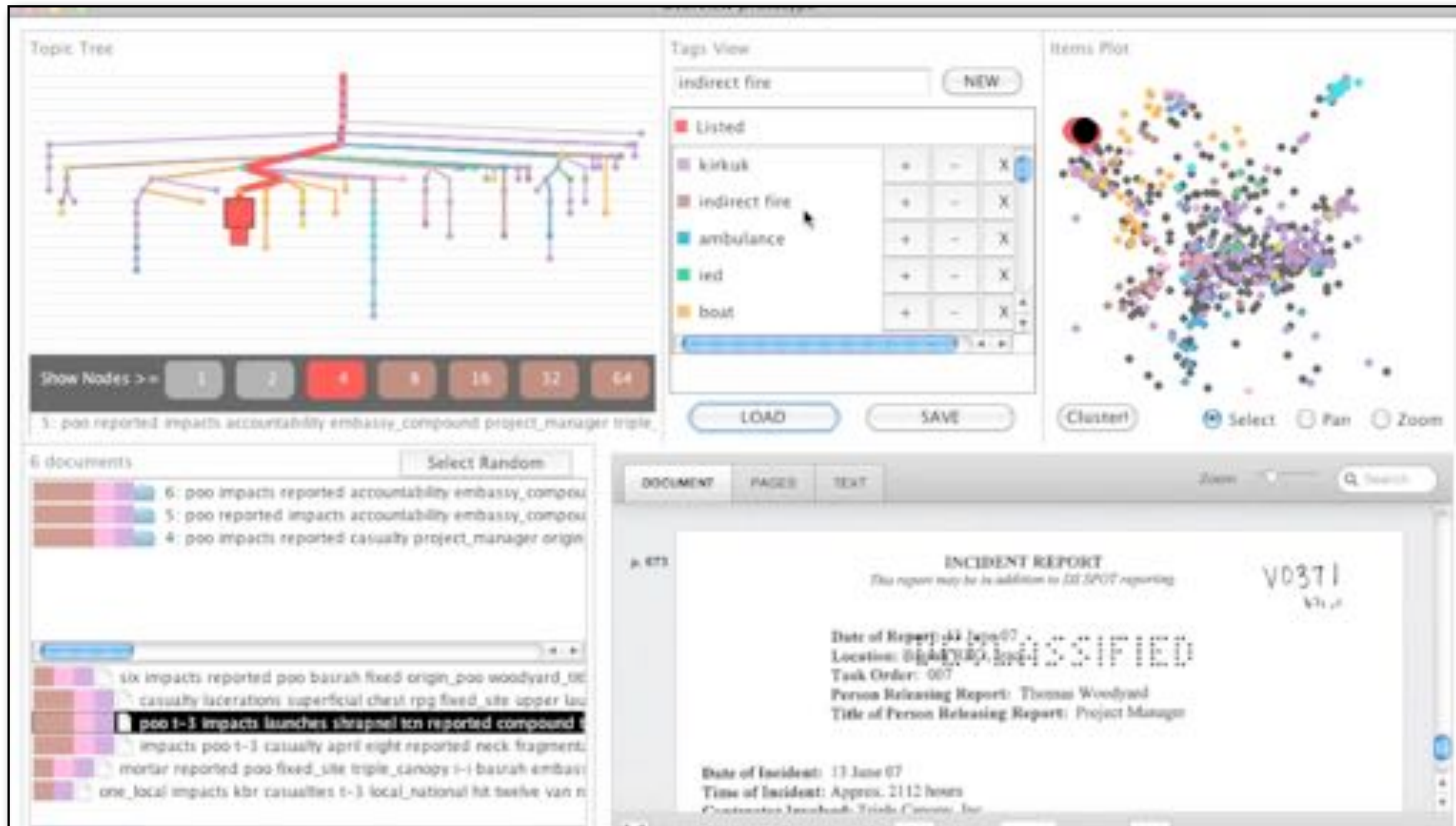


# OVERVIEW V.2 (DESKTOP APP), 2012





# OVERVIEW V.2 (DESKTOP APP), 2012



**Q:** could the use of Overview lead to a published story?

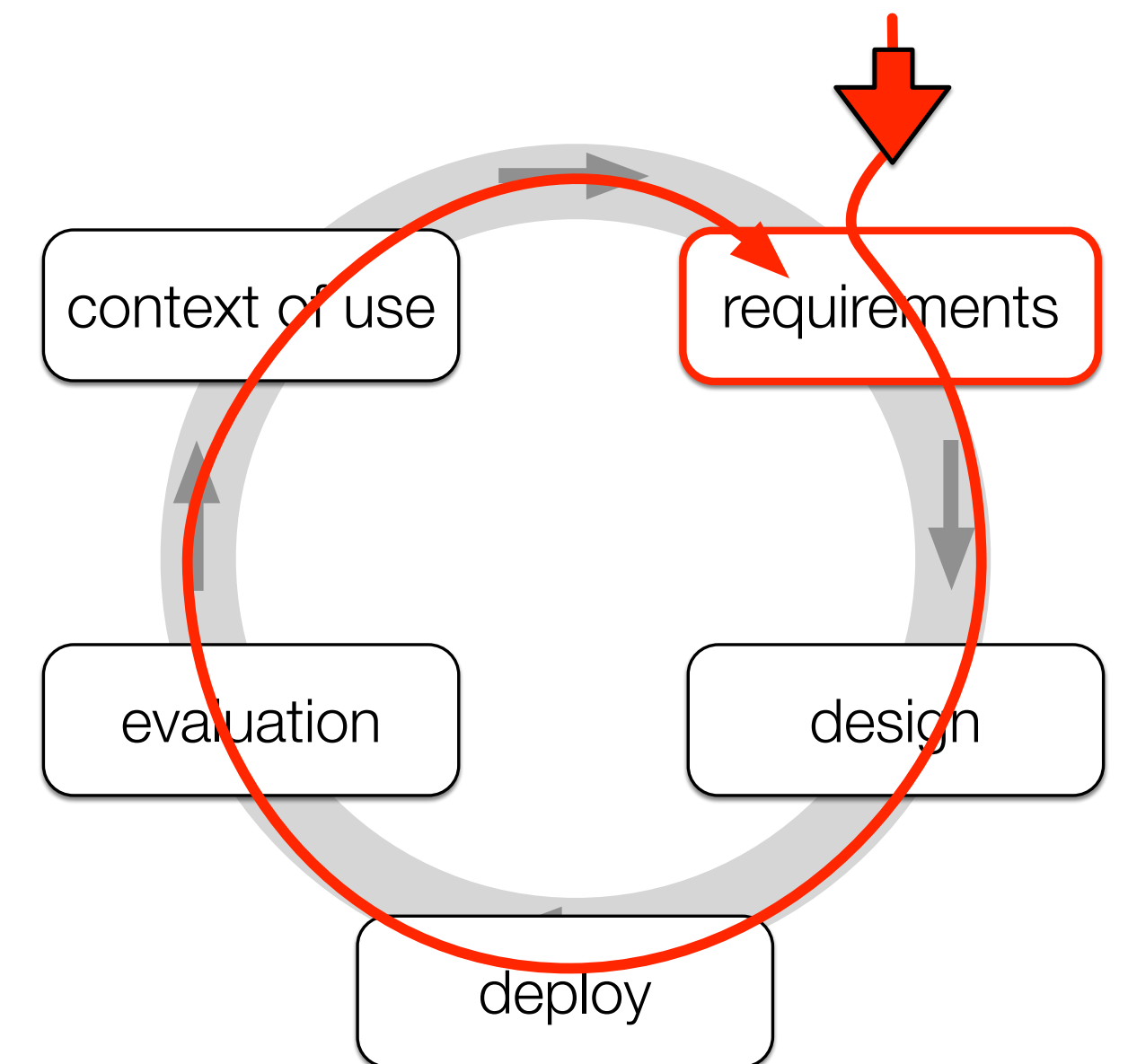


# OVERVIEW V.2 ADOPTION, 2012

**What did private security contractors do in Iraq?**

by Jonathan Stray on 02/21/2012 | 3 | Edit

**AP**





# OVERVIEW V.2 ADOPTION, 2012

## What did private security contractors do in Iraq?

by Jonathan Stray on 02/21/2012 | 3 | Edit

AP

## TPD working through flawed mobile system

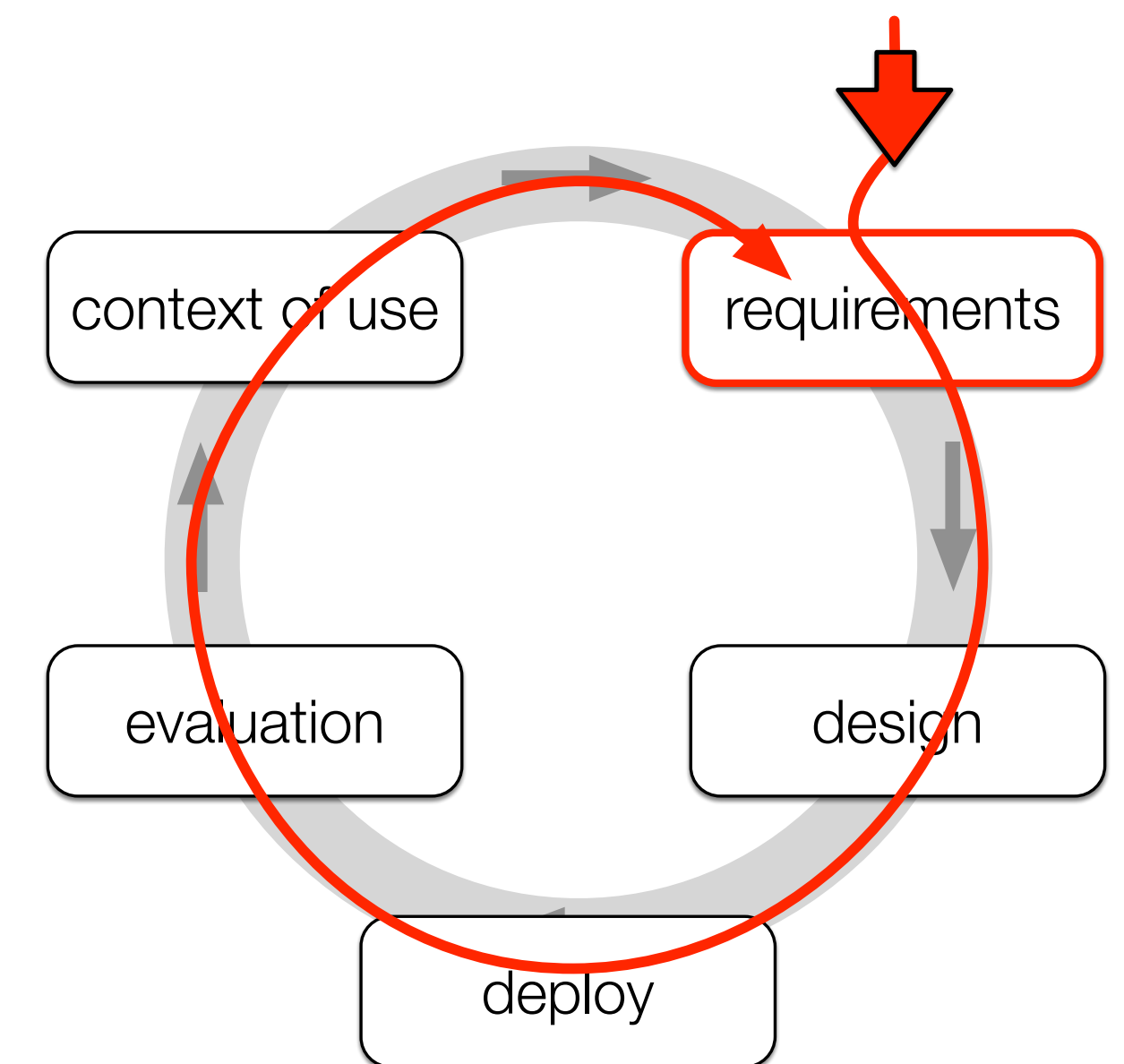
By JARREL WADE World Staff Writer on Jun 3, 2012, at 2:19 AM

TULSA WORLD

## RYAN ASKED FOR FEDERAL HELP AS HE CHAMPIONED CUTS

By JACK GILLUM — Oct. 12 7:20 PM EDT

AP





# OVERVIEW V.2 CASE STUDIES, 2012

**Document  
Collection**

**Question**

**Outcome**



# OVERVIEW V.2 CASE STUDIES, 2012

## CASE STUDY #1

### Document Collection

666 reports (4,500 pages) from FOIA

### Question

*What incidents involved security  
contractors during Iraq war?*

### Outcome

Summarized prevalence of document  
categories

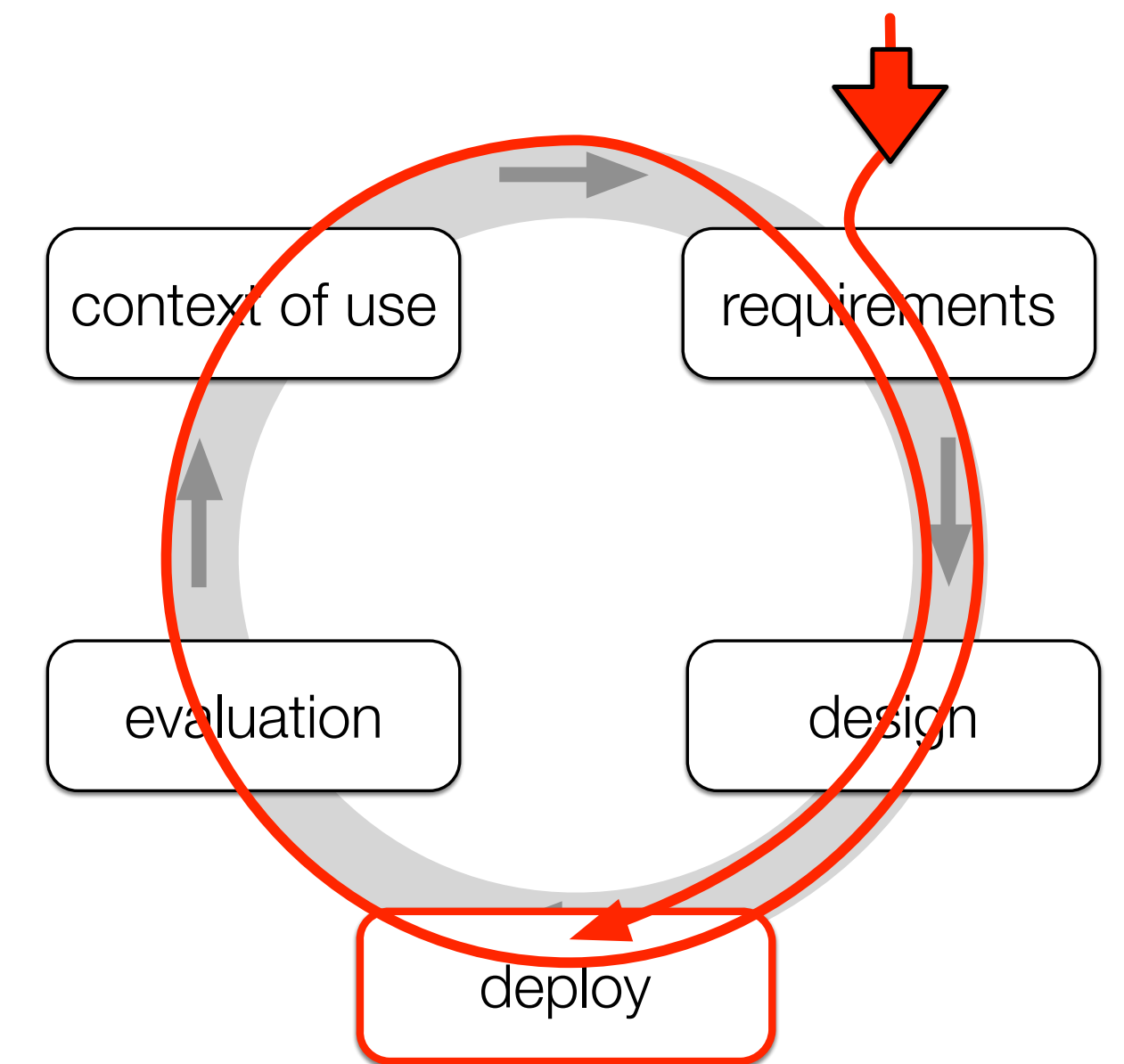
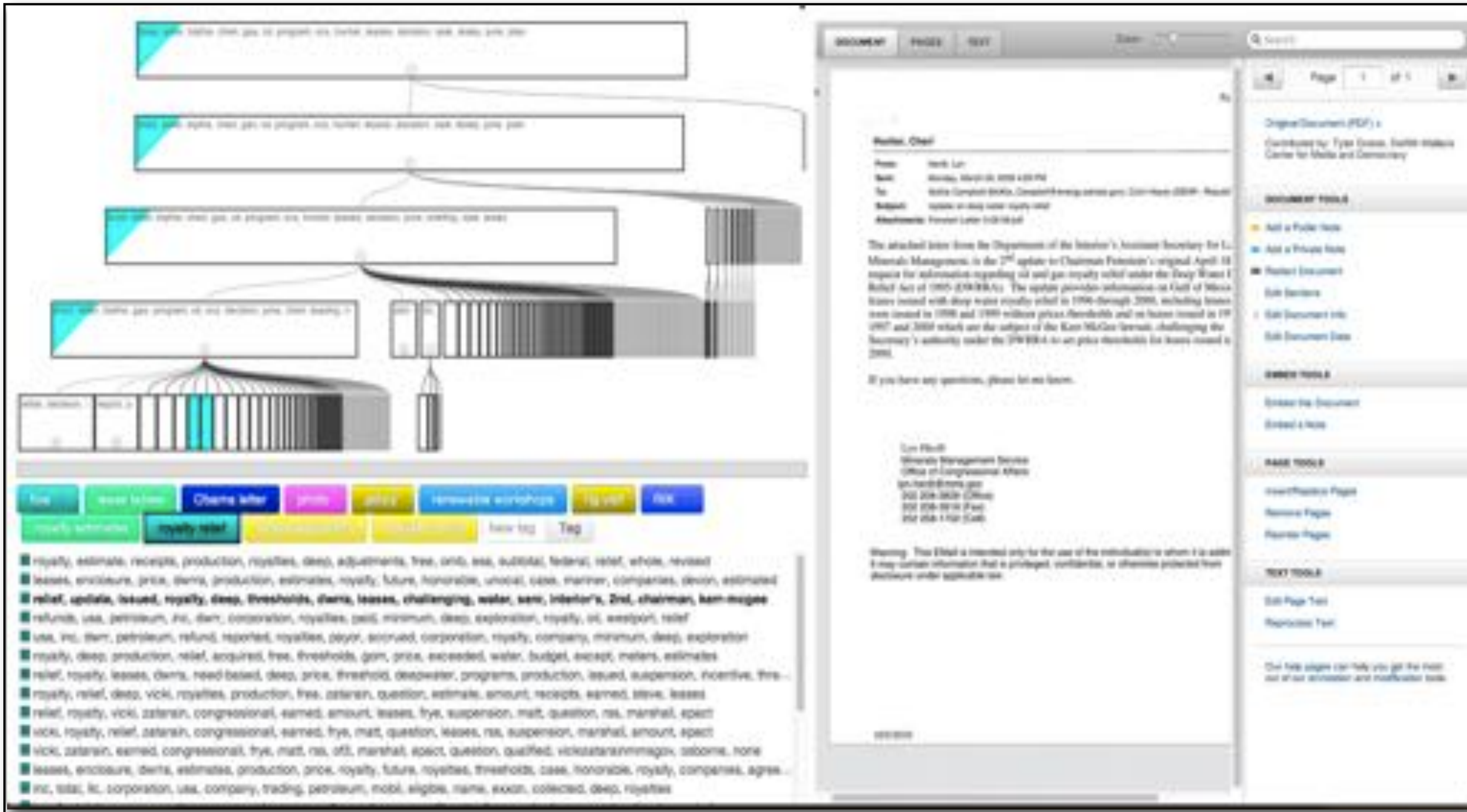


# OVERVIEW V.2 CASE STUDIES, 2012

	CASE STUDY #1	CASE STUDY #2
Document Collection	666 reports (4,500 pages) from FOIA	5,996 emails from FOIA
Question	<i>What incidents involved security contractors during Iraq war?</i>	<i>Were municipal police funds mismanaged?</i>
Outcome	Summarized prevalence of document categories	Located evidence supporting hypothesis



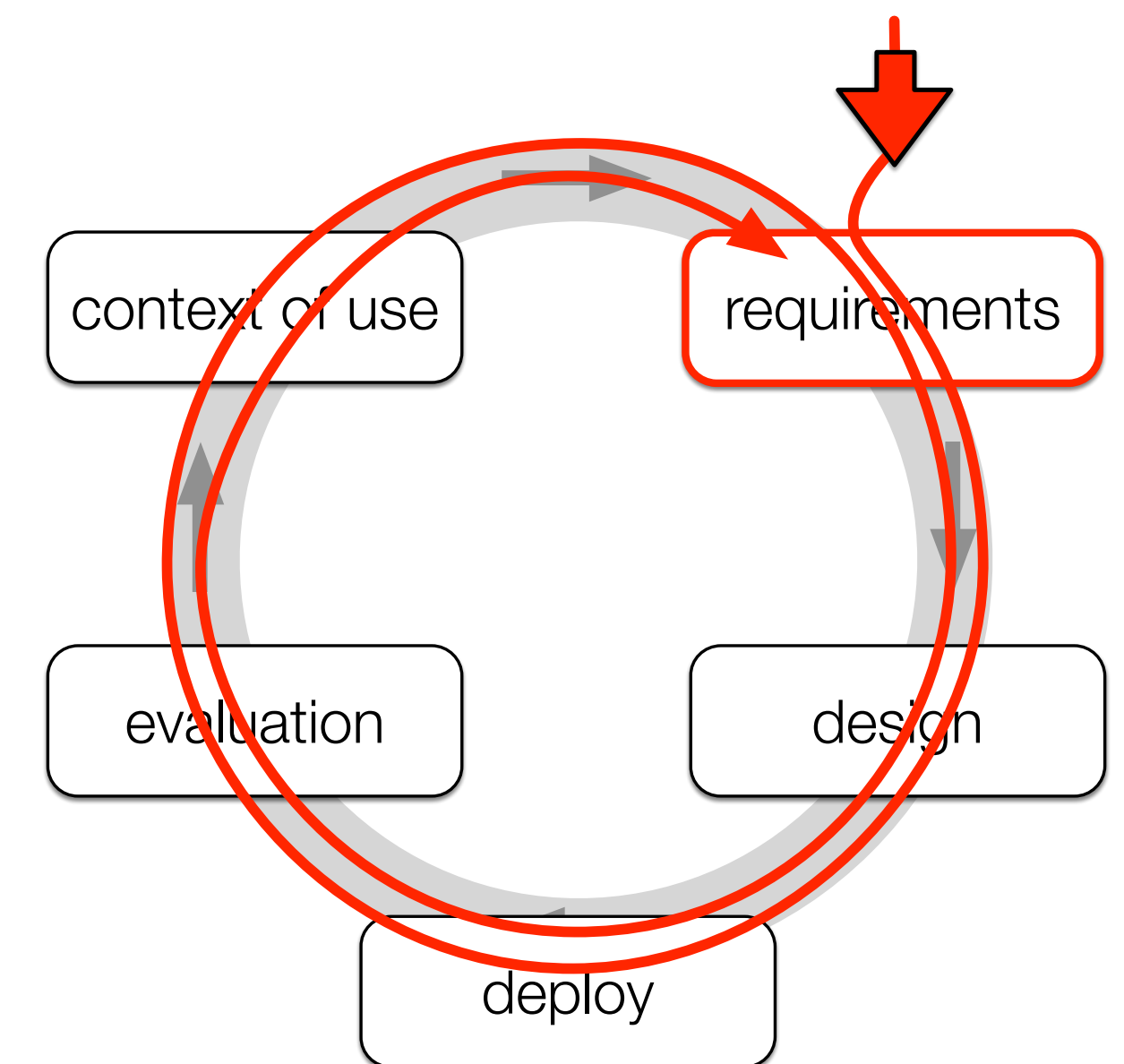
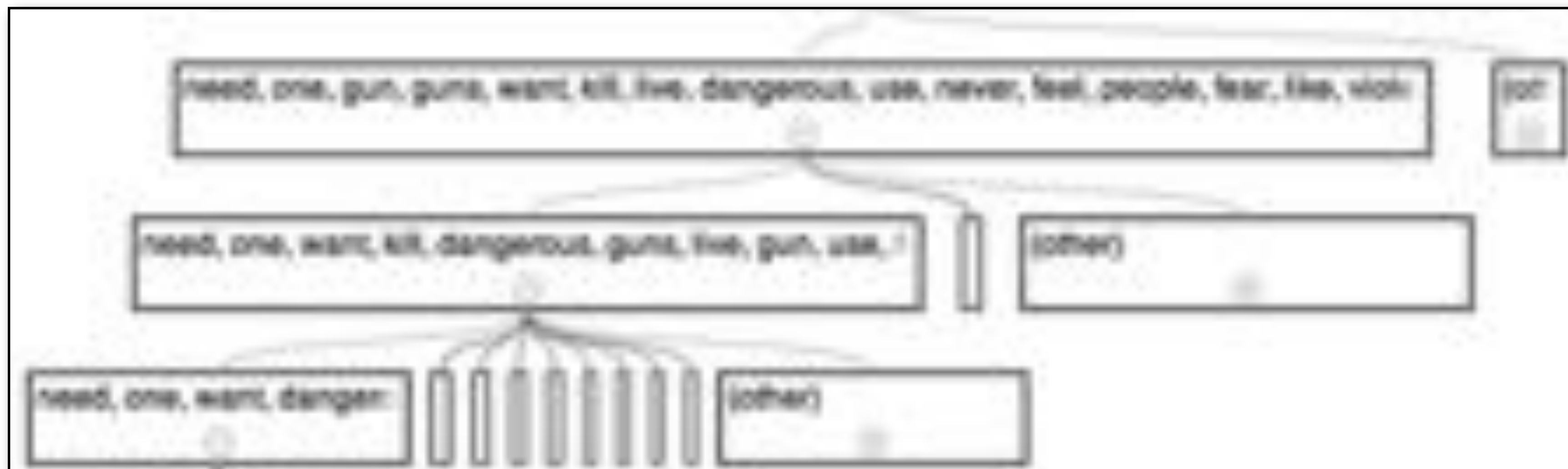
# OVERVIEW V.3 (WEB APP), FALL 2012



**clockwise from TL:**  
tree visualization,  
document viewer,  
document list,  
tag controls



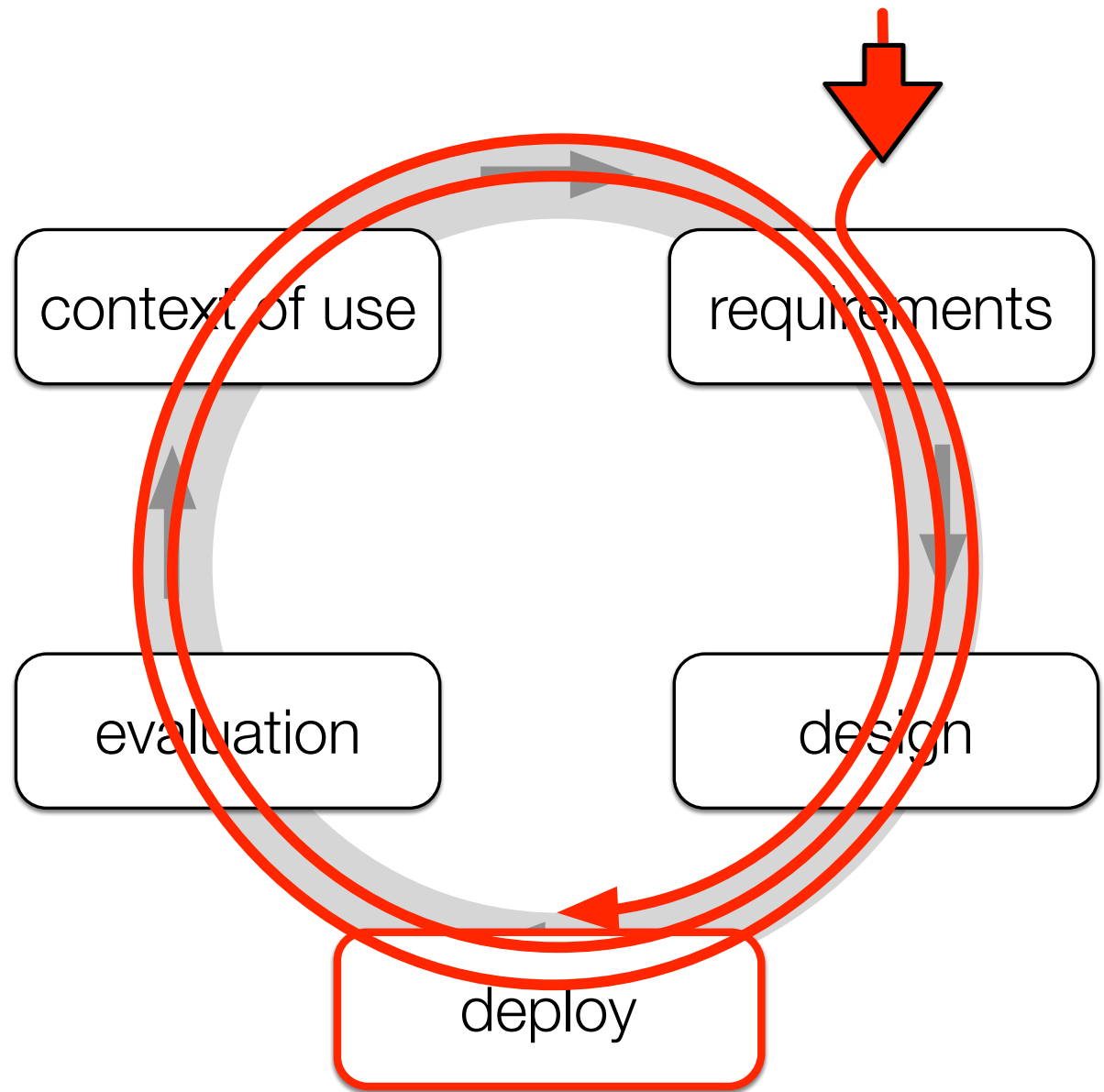
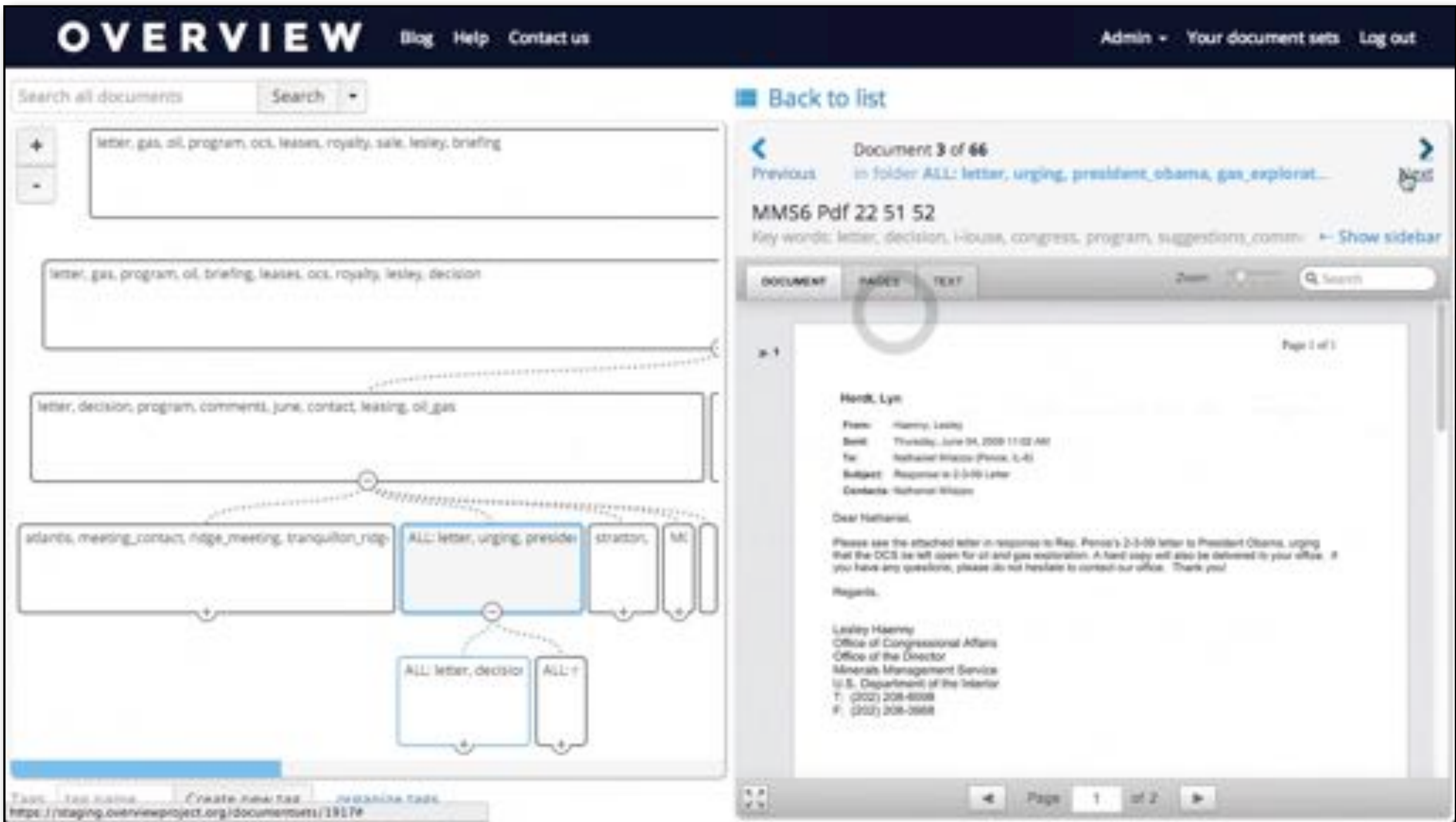
# OVERVIEW V.3 EVALUATION, 2013



**+ think-aloud  
evaluation  
with 5  
journalists**



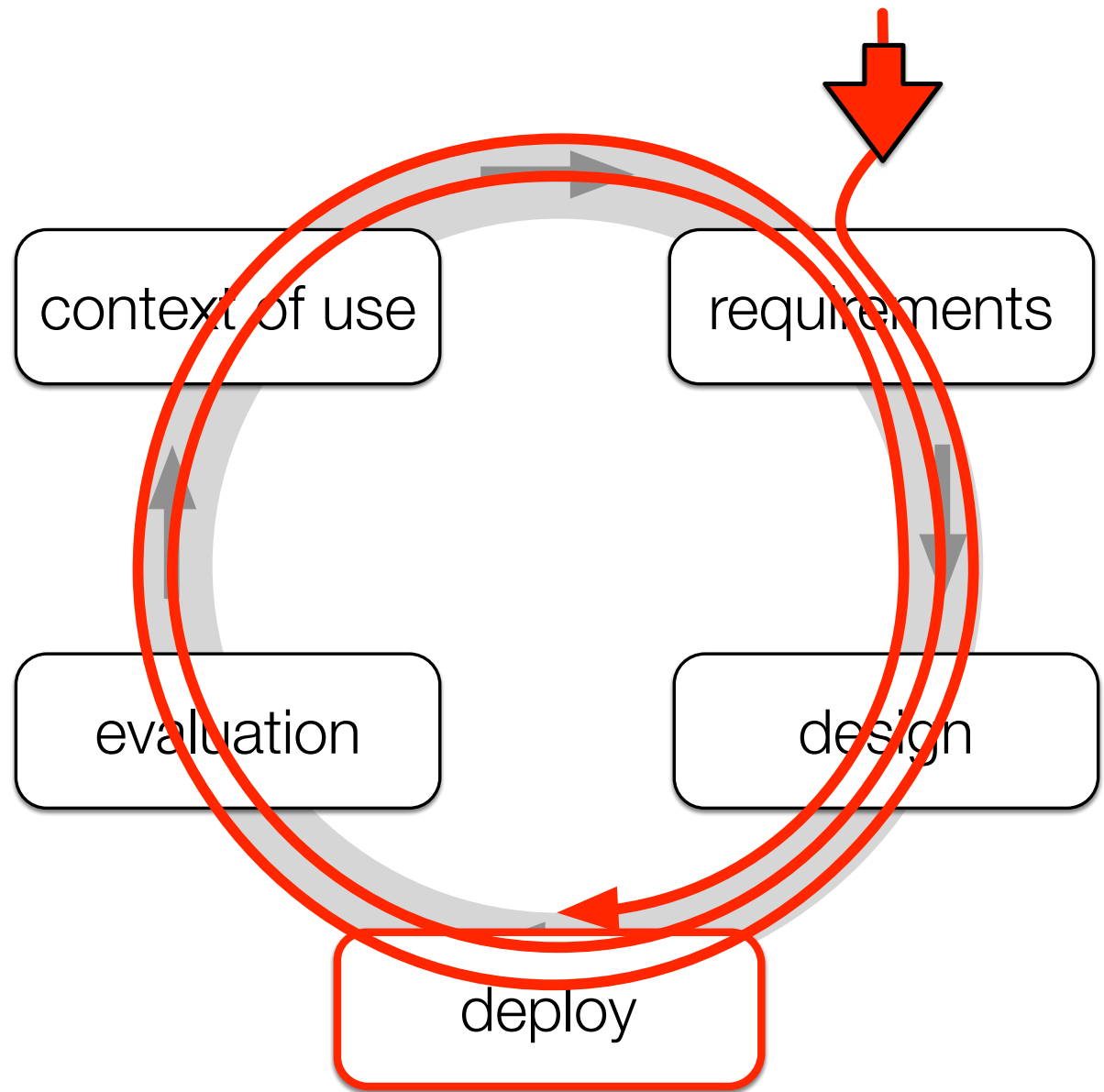
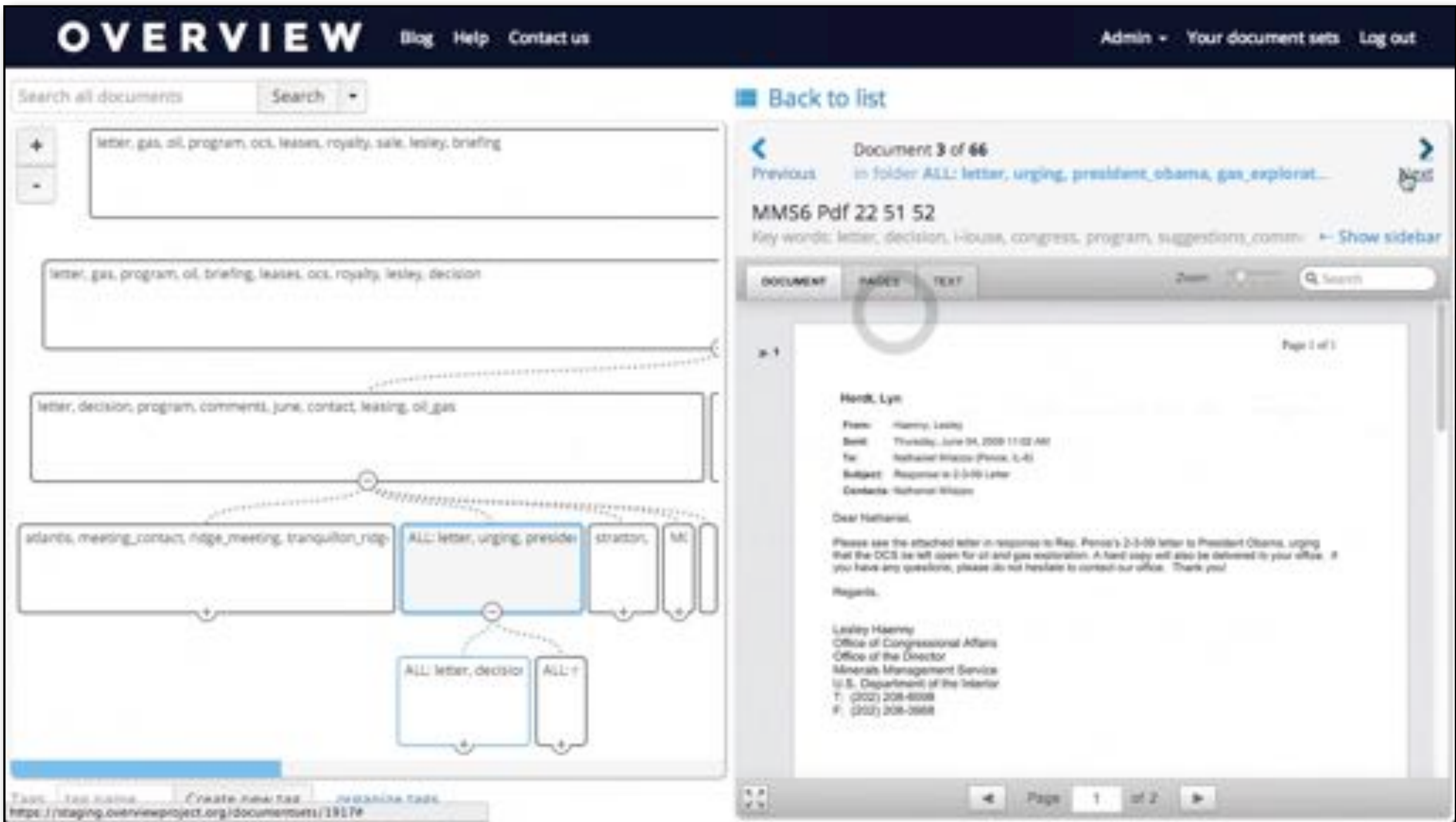
# OVERVIEW V.4, FALL 2013



**clockwise from TL:**  
full-text search,  
document list /  
viewer,  
tag controls,  
tree visualization



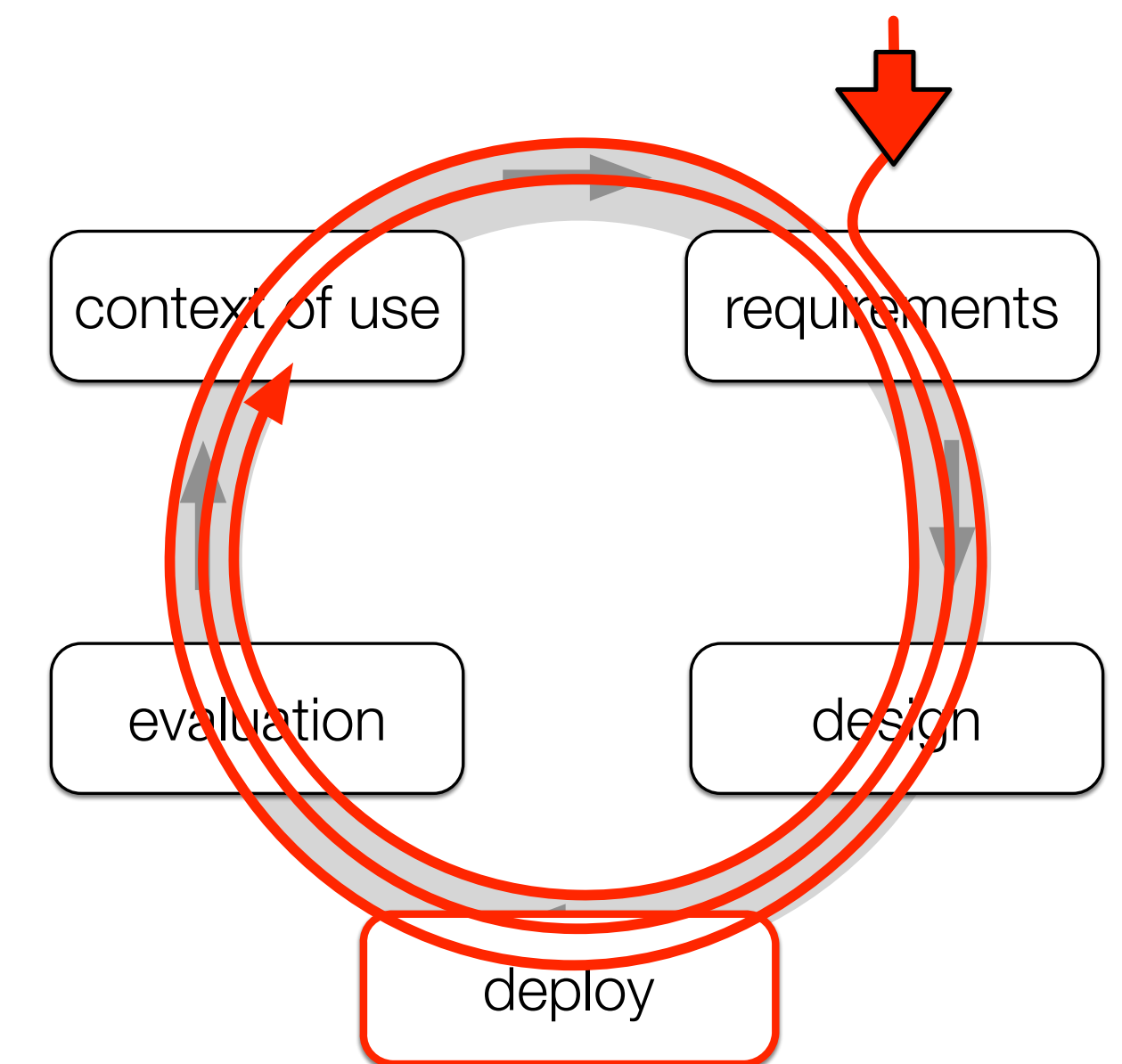
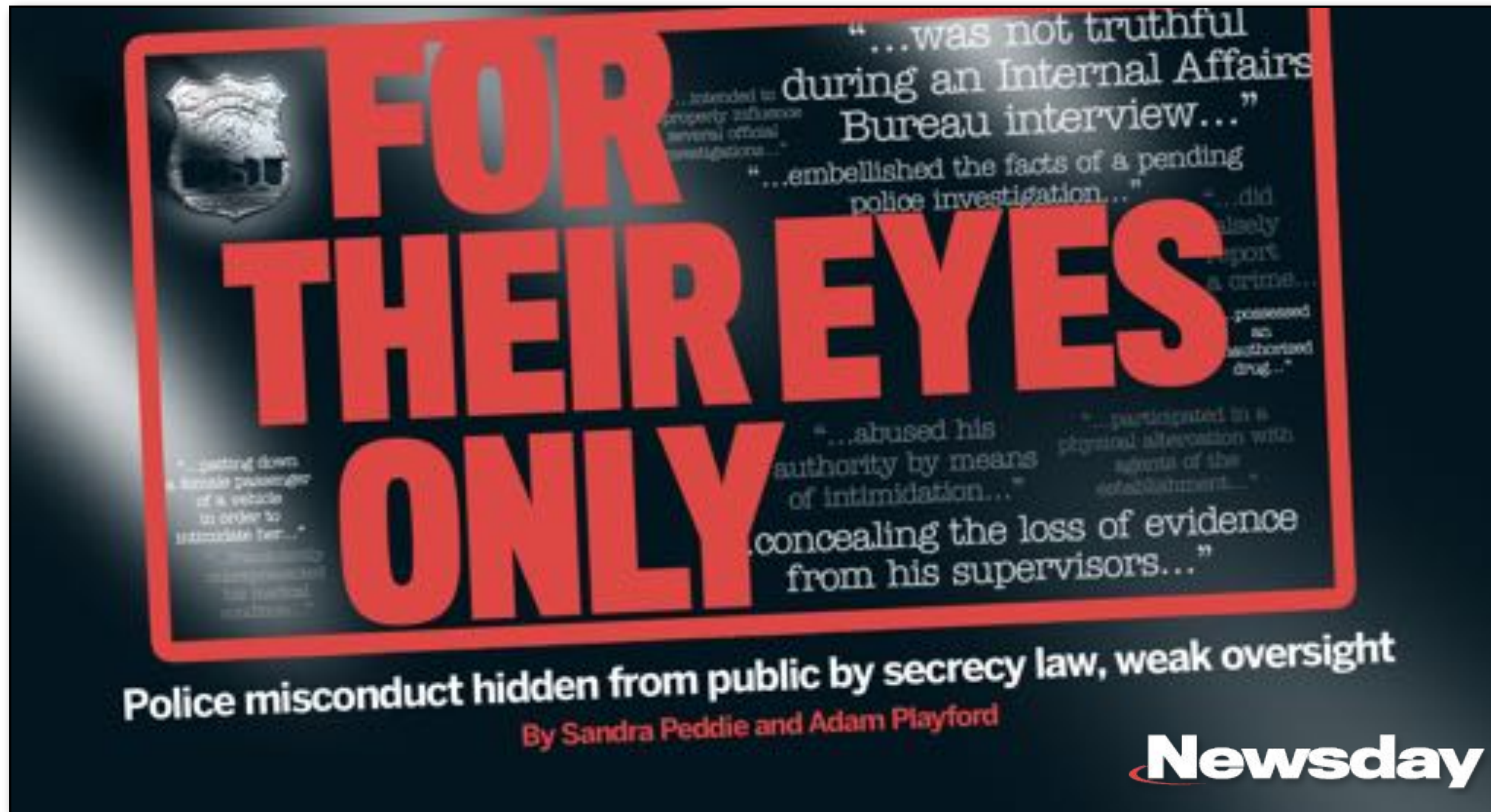
# OVERVIEW V.4, FALL 2013



**clockwise from TL:**  
full-text search,  
document list /  
viewer,  
tag controls,  
tree visualization



## OVERVIEW V.4 ADOPTION



finalist for  
2014  
Pulitzer  
Prize in  
journalism



# OVERVIEW V.4 ADOPTION

## CASE STUDY #6

### Document Collection

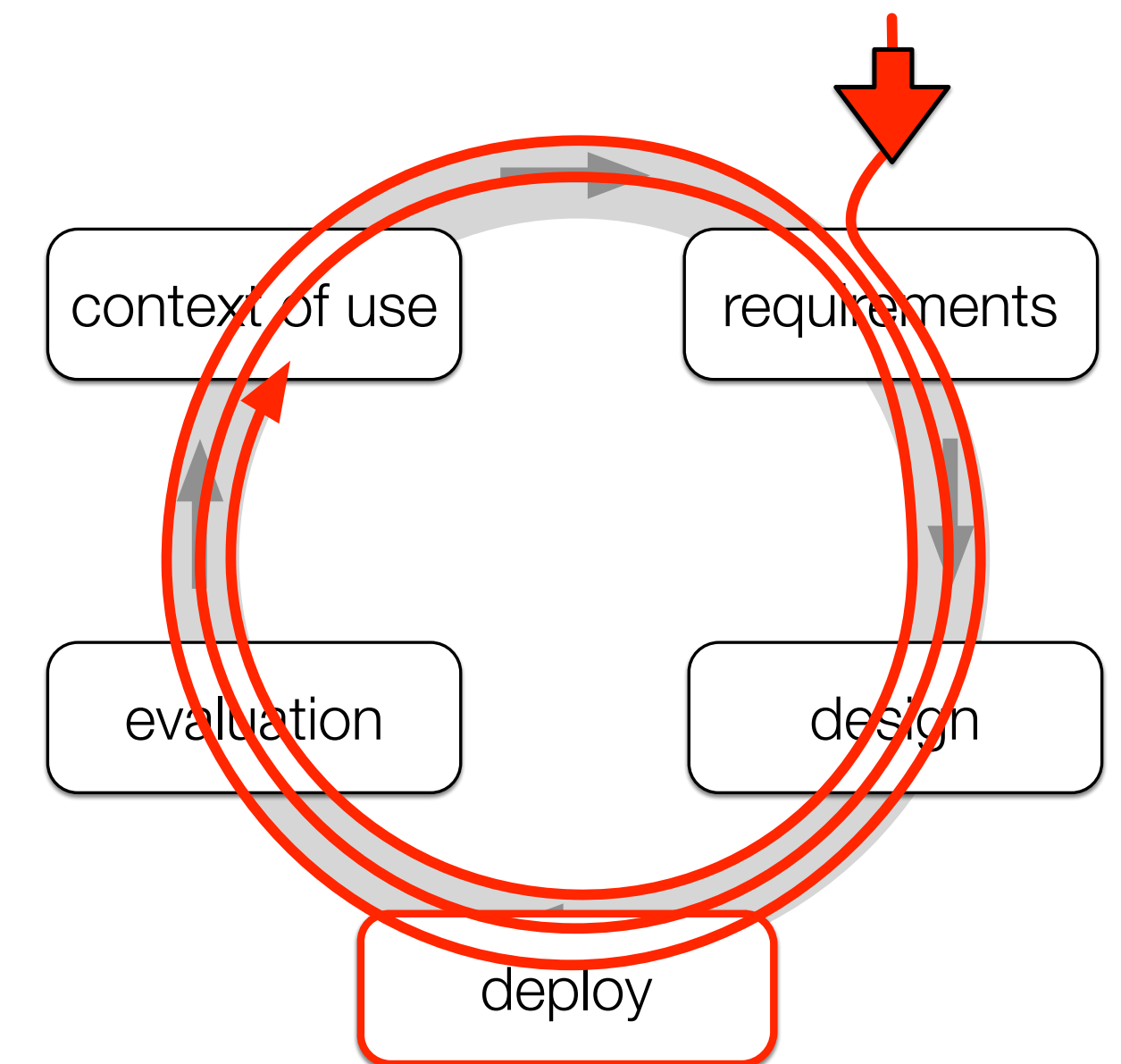
1,680 proposed and passed bills

### Question

*Did the NY state legislature fail to pass any bills addressing police misconduct?*

### Outcome

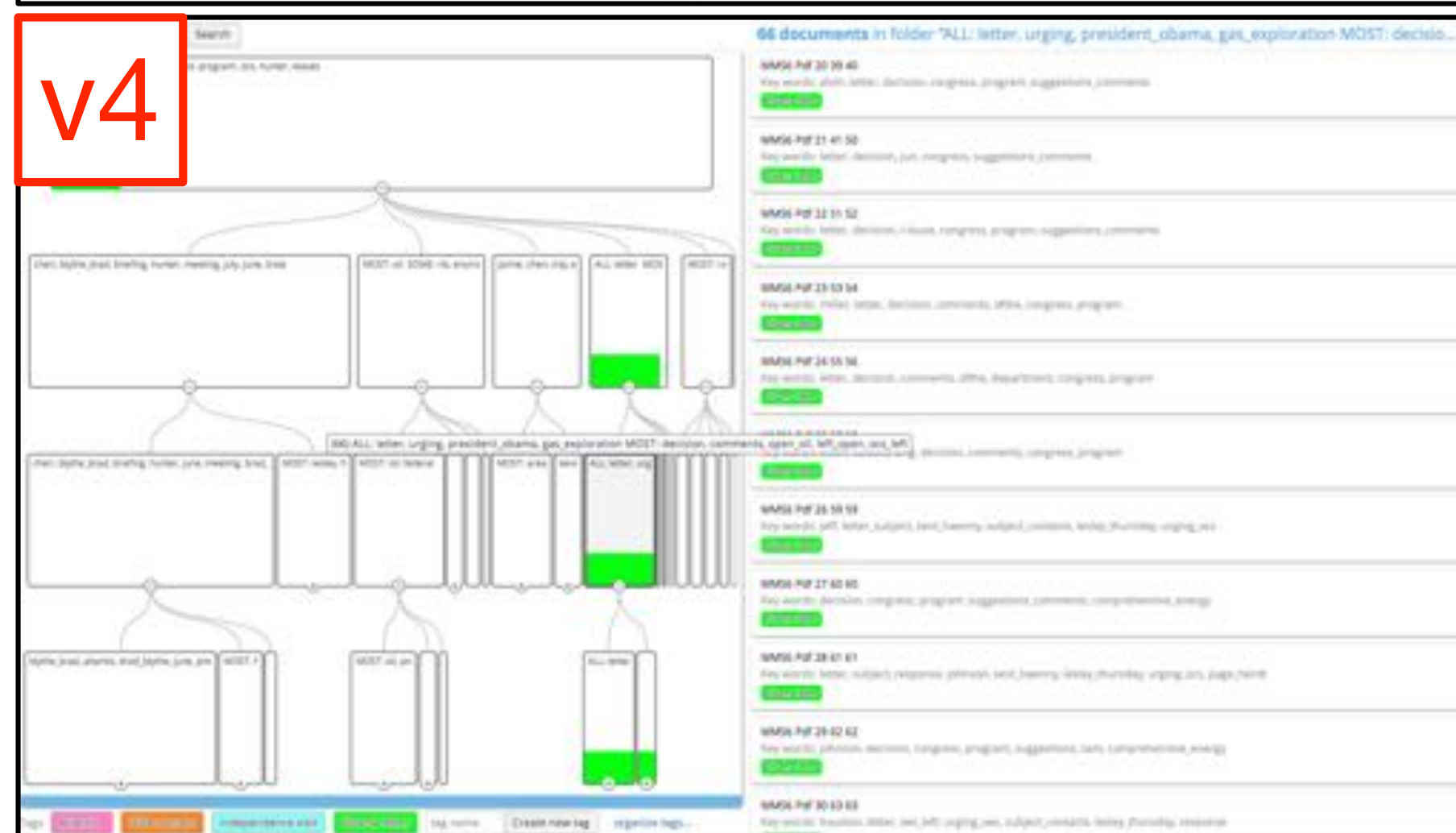
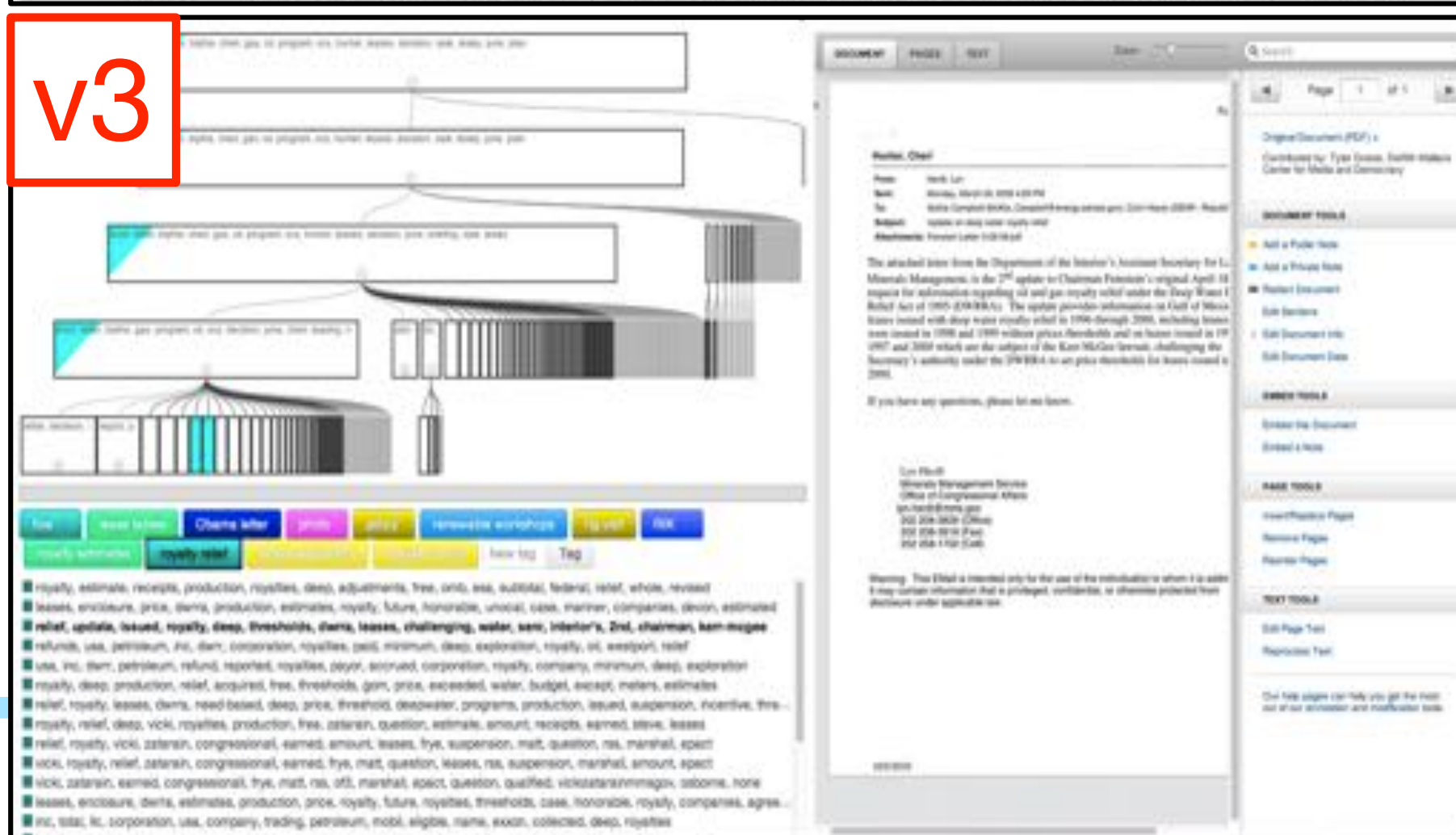
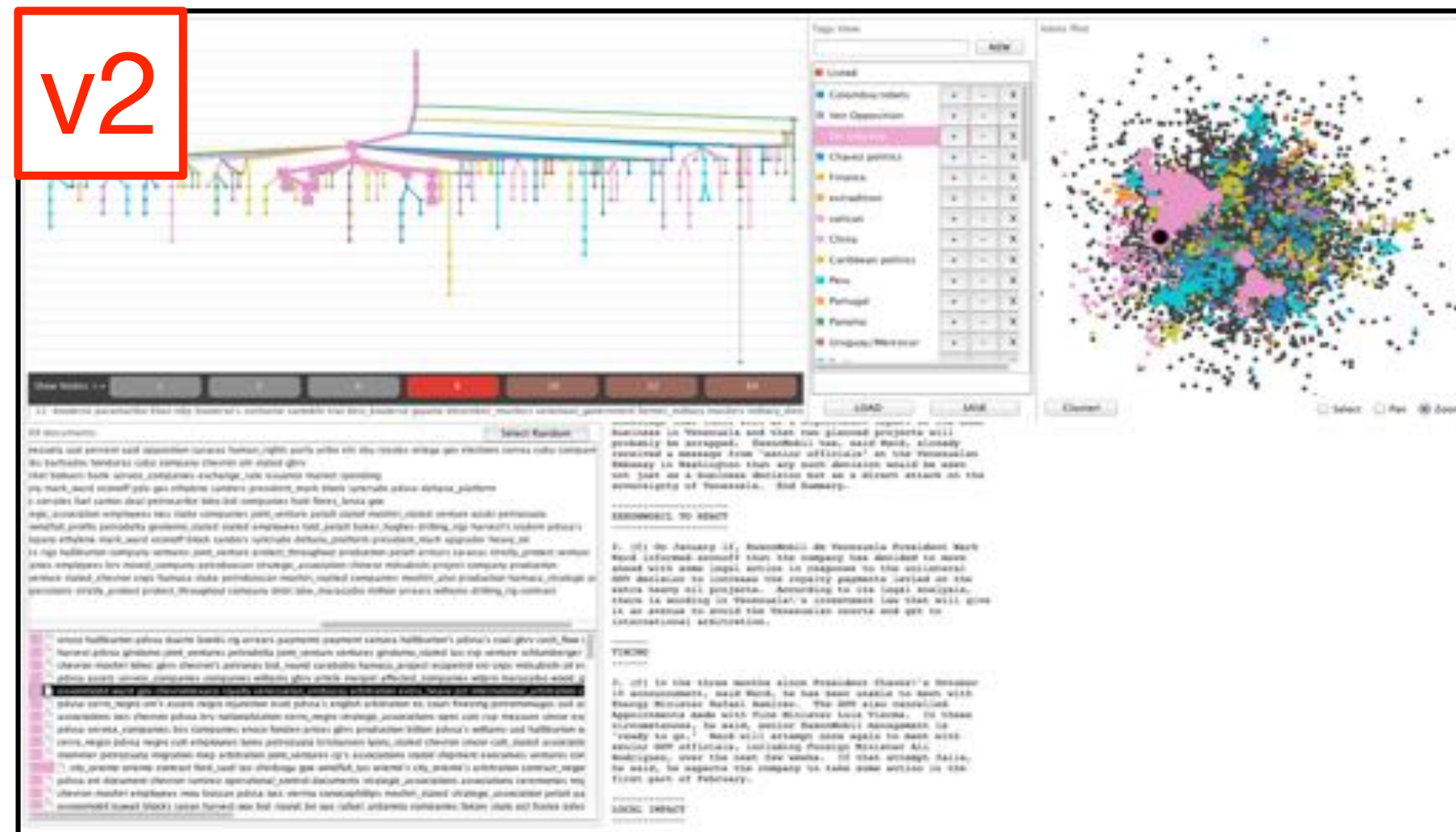
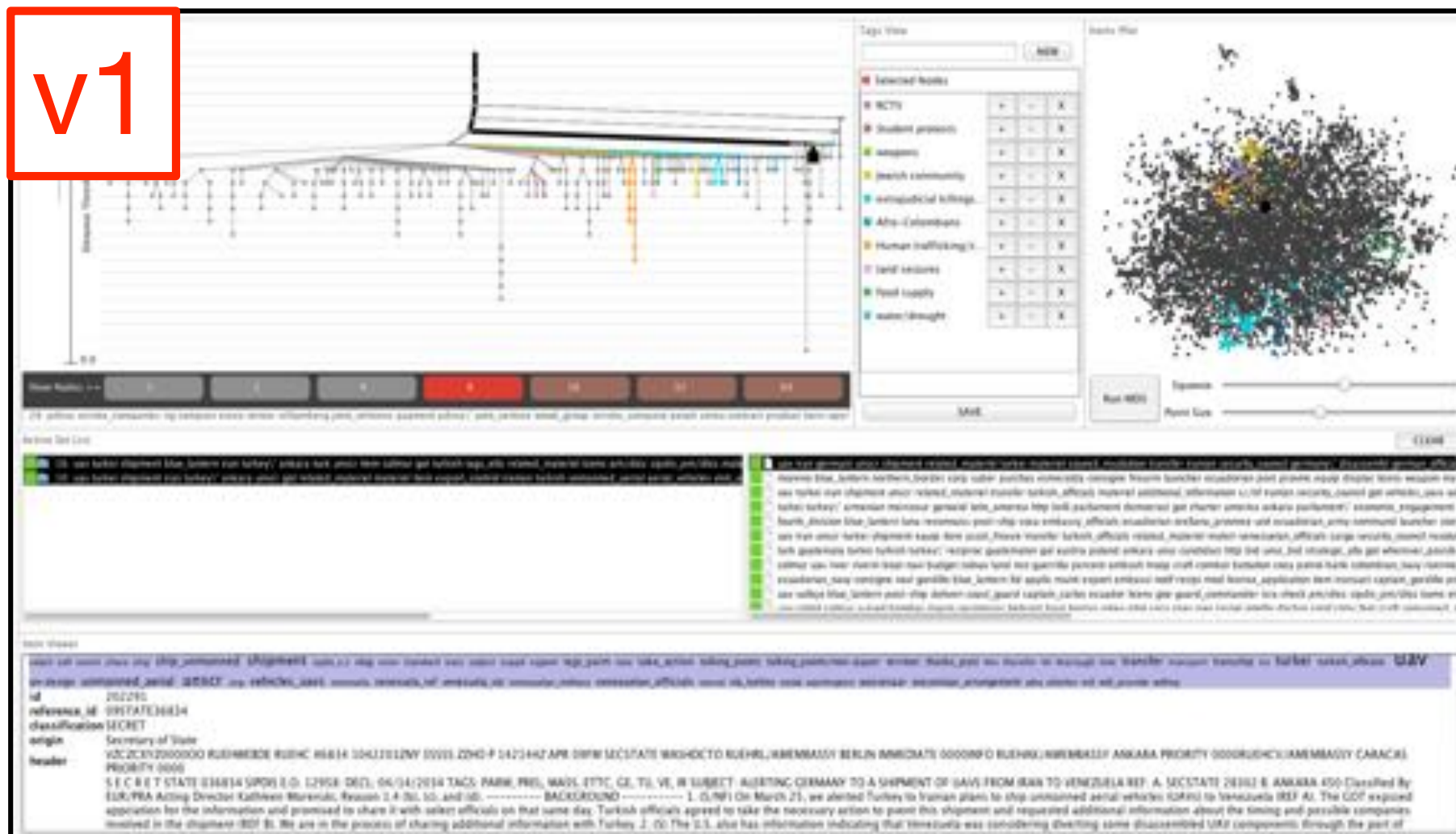
Proved non-existence of evidence



**finalist for  
2014  
Pulitzer  
Prize in  
journalism**



# 4 VERSIONS, 6 CASE STUDIES



CS1

CS2

CS3

CS4

CS5

CS6



# 4 VERSIONS, 6 CASE STUDIES

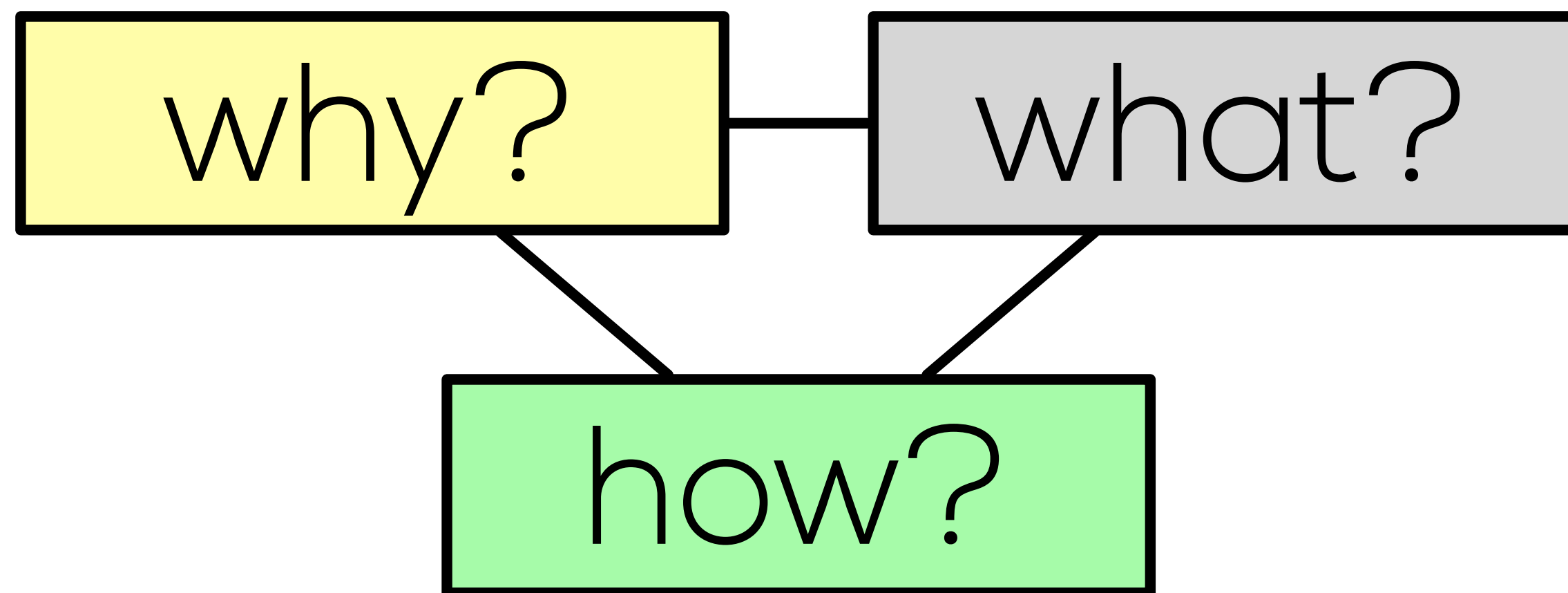
Case Study	#1	#2	#3	#4	#5	#6
Document Collection	4,500 pages from FOIA	5,996 emails from FOIA	8,680 pages from FOIA	1,278 survey comments	4,653 emails from FOIA	1,680 bills
Question	<i>What did security contractors do during Iraq war?</i>	<i>Were municipal police funds mismanaged?</i>	<i>Were Paul Ryan's campaign statements hypocritical?</i>	<i>What is the gun ownership debate about?</i>	<i>Was gov't response to emergency incident effective?</i>	<i>Did gov't pass bills addressing police misconduct?</i>



**"explore"** is not a  
well-defined task



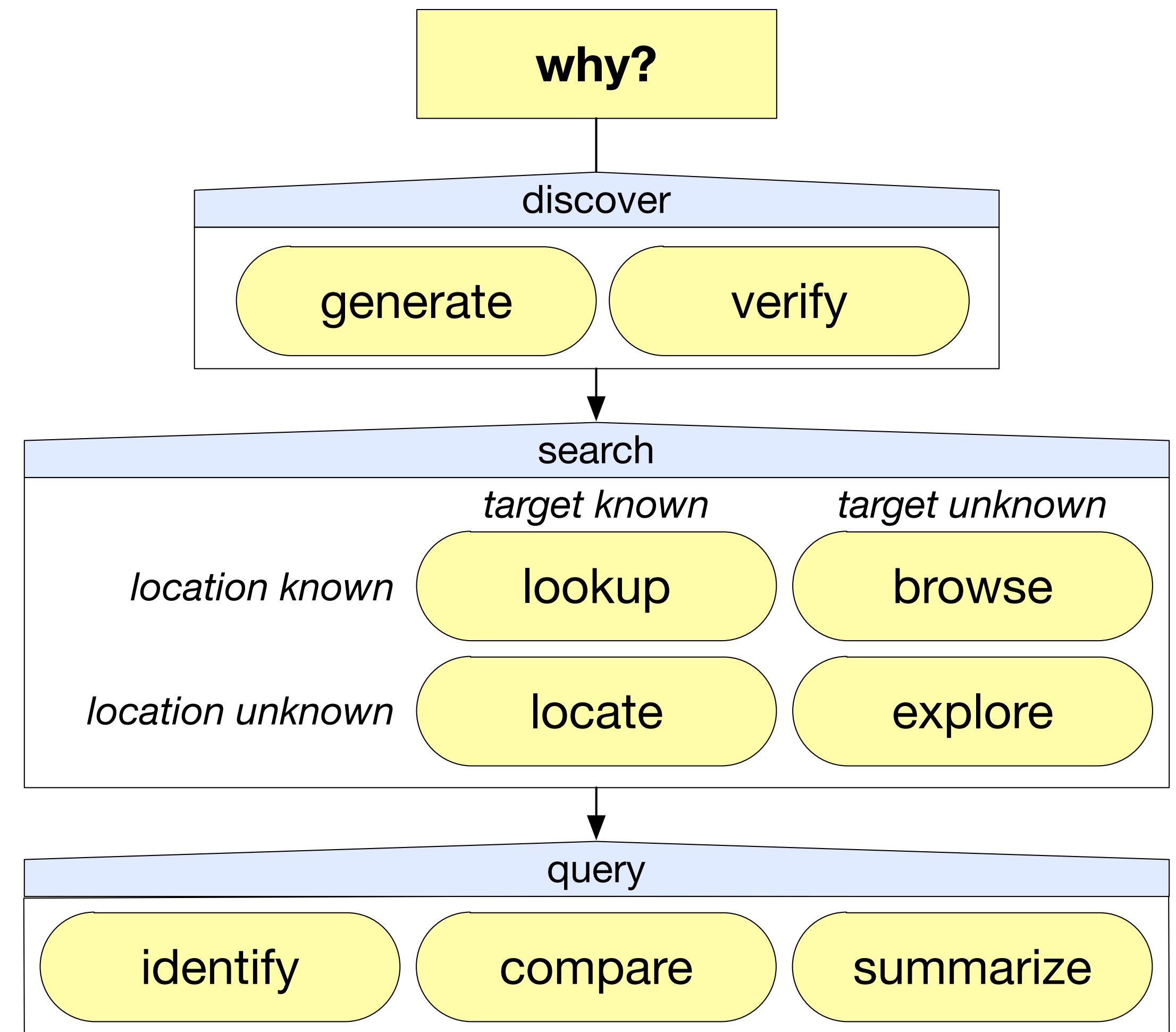
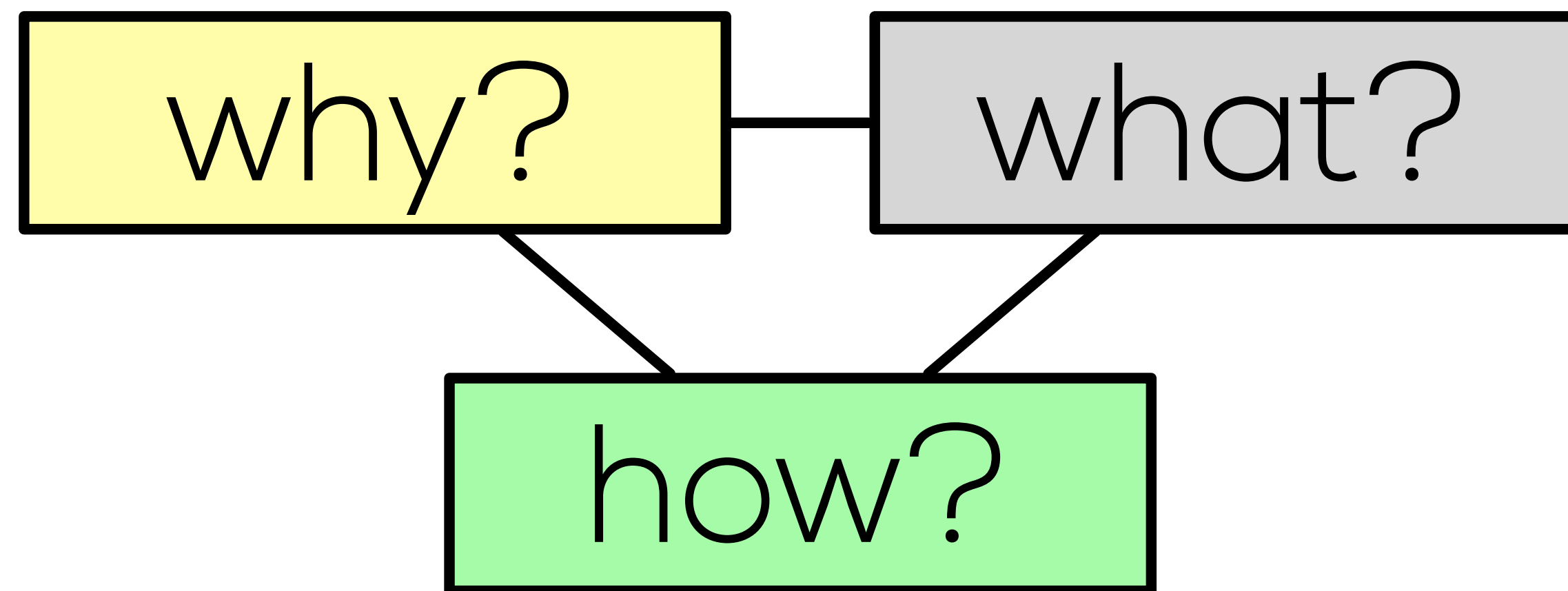
# VIS TASK ANALYSIS



Brehmer & Munzner: "A multi-level typology of abstract visualization tasks" (IEEE TVCG / InfoVis 2013).



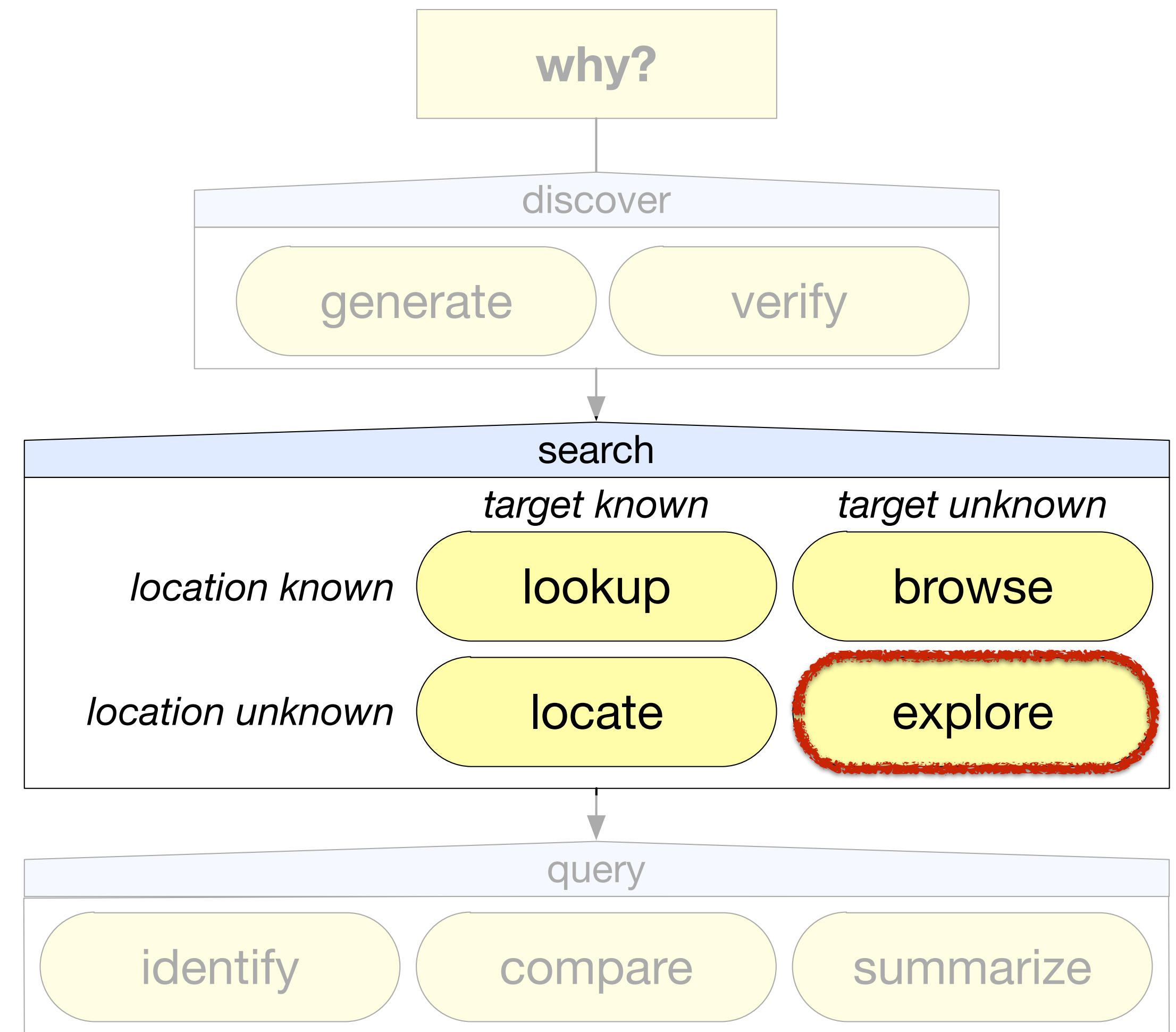
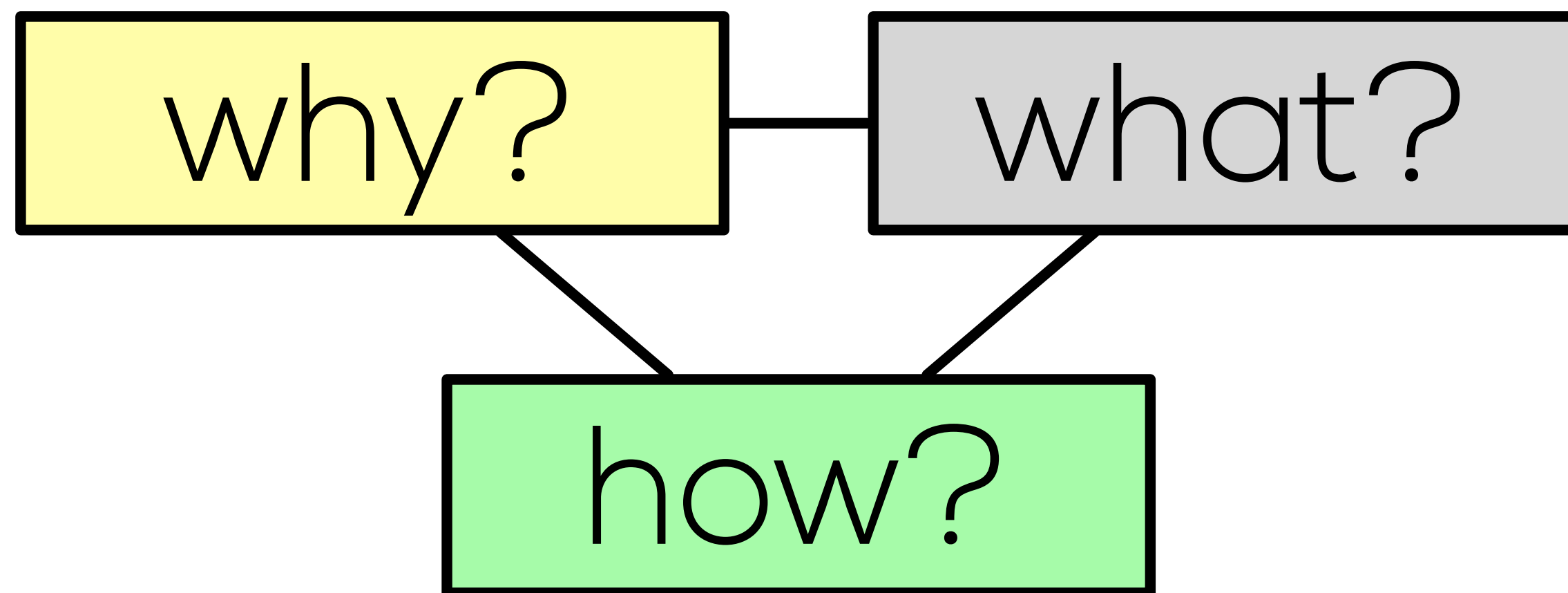
# VIS TASK ANALYSIS



Brehmer & Munzner: "A multi-level typology of abstract visualization tasks" (IEEE TVCG / InfoVis 2013).



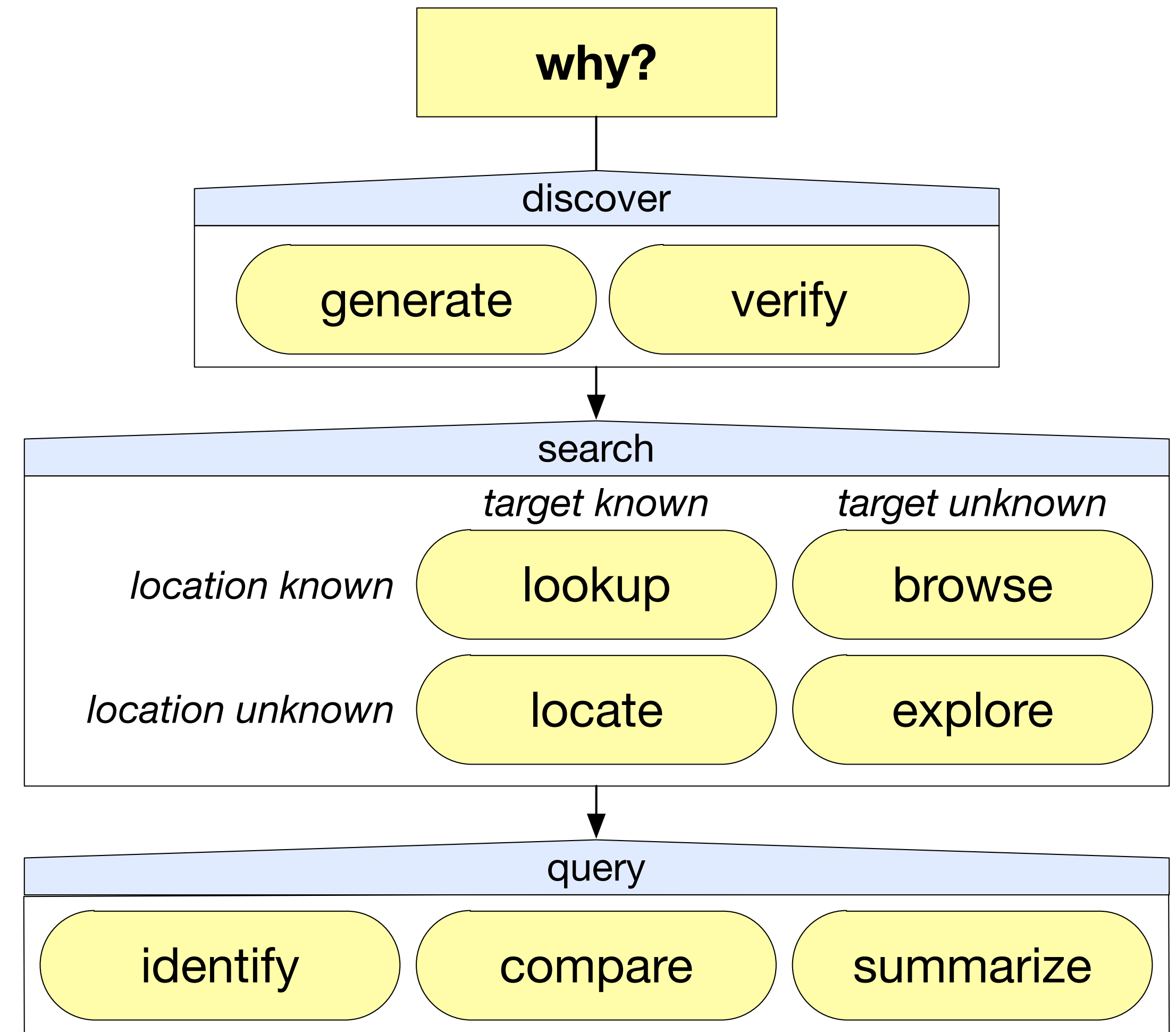
# VIS TASK ANALYSIS



Brehmer & Munzner: "A multi-level typology of abstract visualization tasks" (IEEE TVCG / InfoVis 2013).



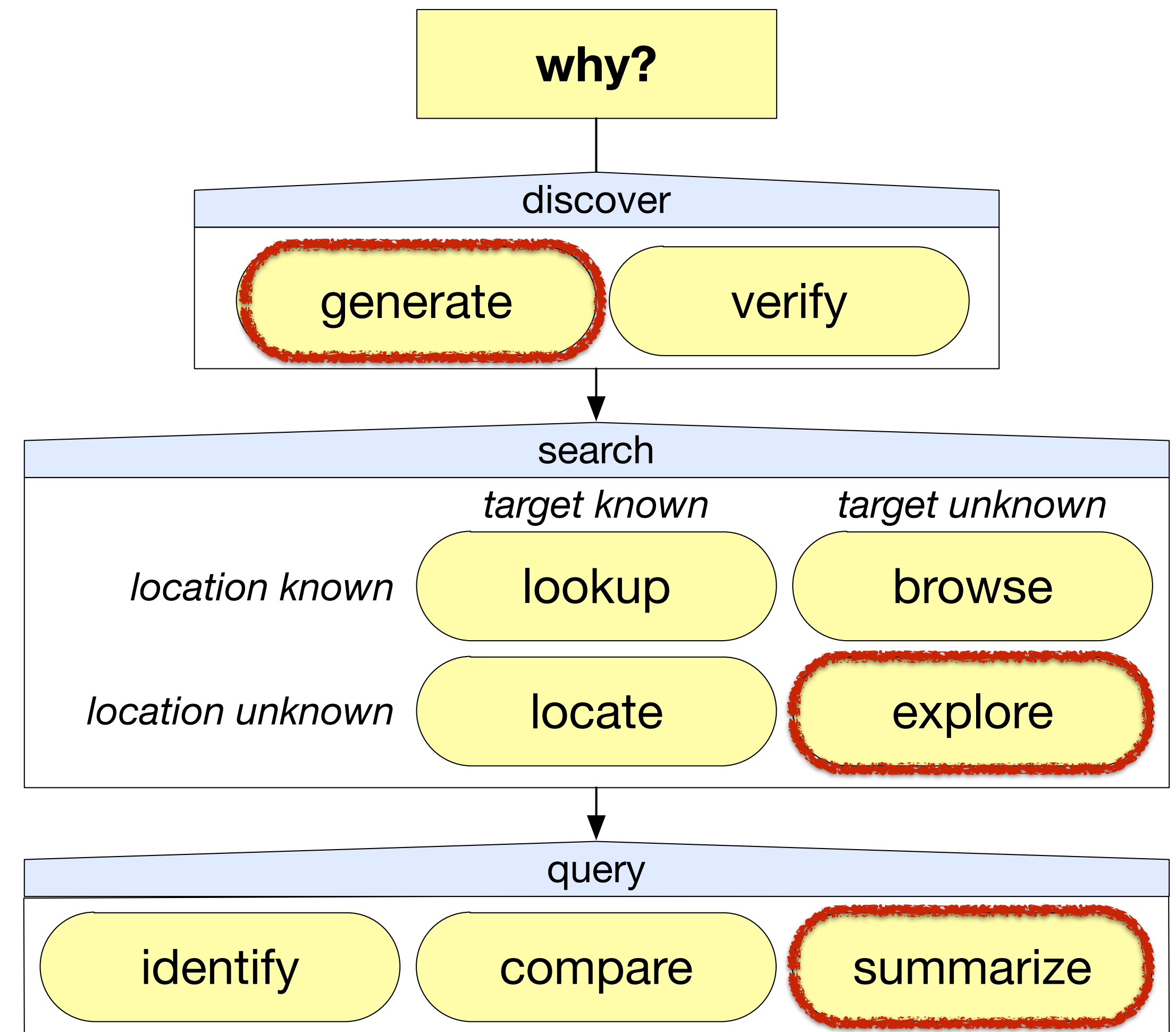
# VIS TASK ANALYSIS





# VIS TASK ANALYSIS

**T1: generate hypotheses**  
→ **explore** → **summarize**  
case studies 1, 4

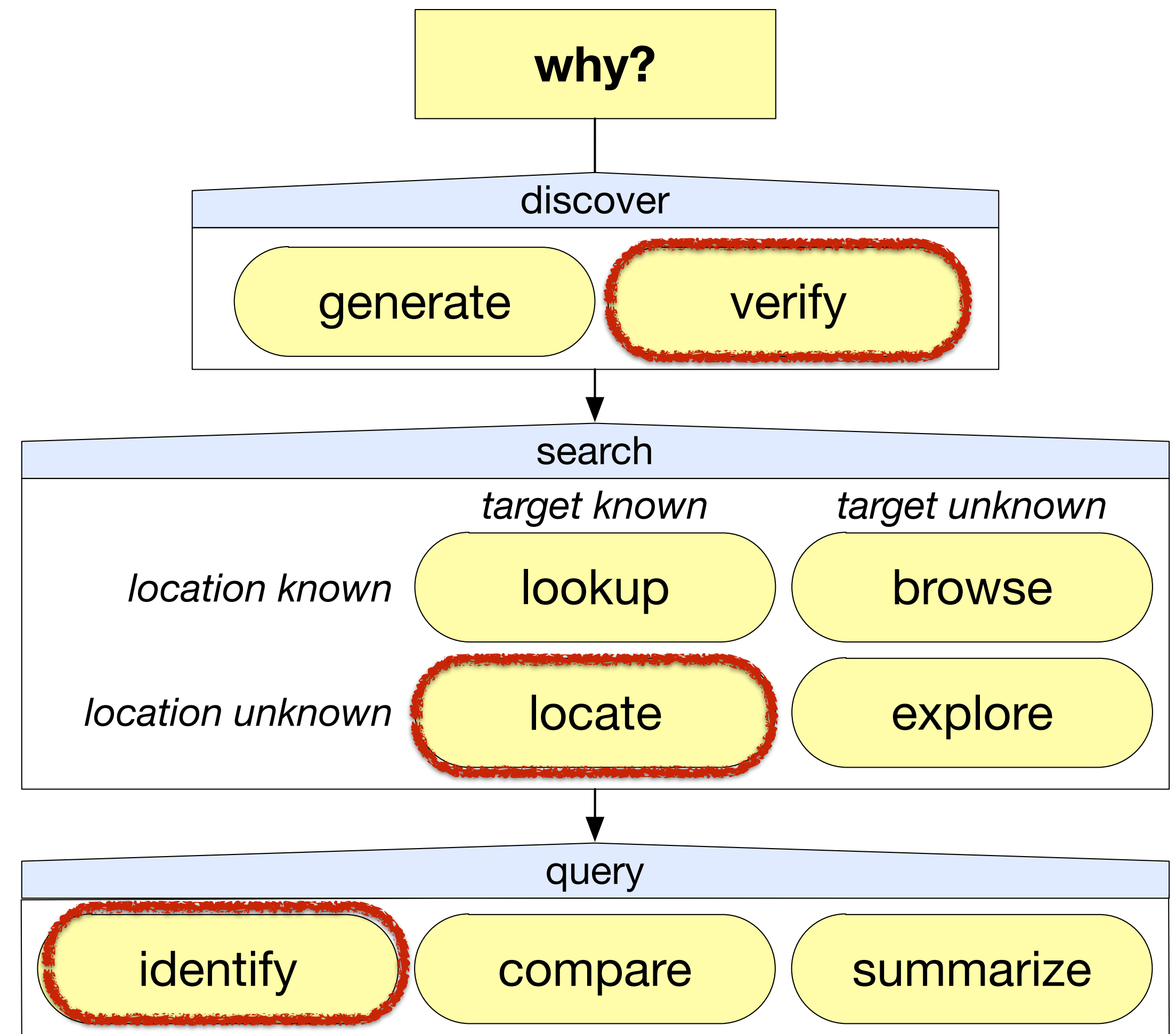




# VIS TASK ANALYSIS

**T1: generate hypotheses**  
→ **explore** → **summarize**  
case studies 1, 4

**T2: verify hypotheses**  
→ **locate** → **identify**  
case studies 2, 3, 5, 6



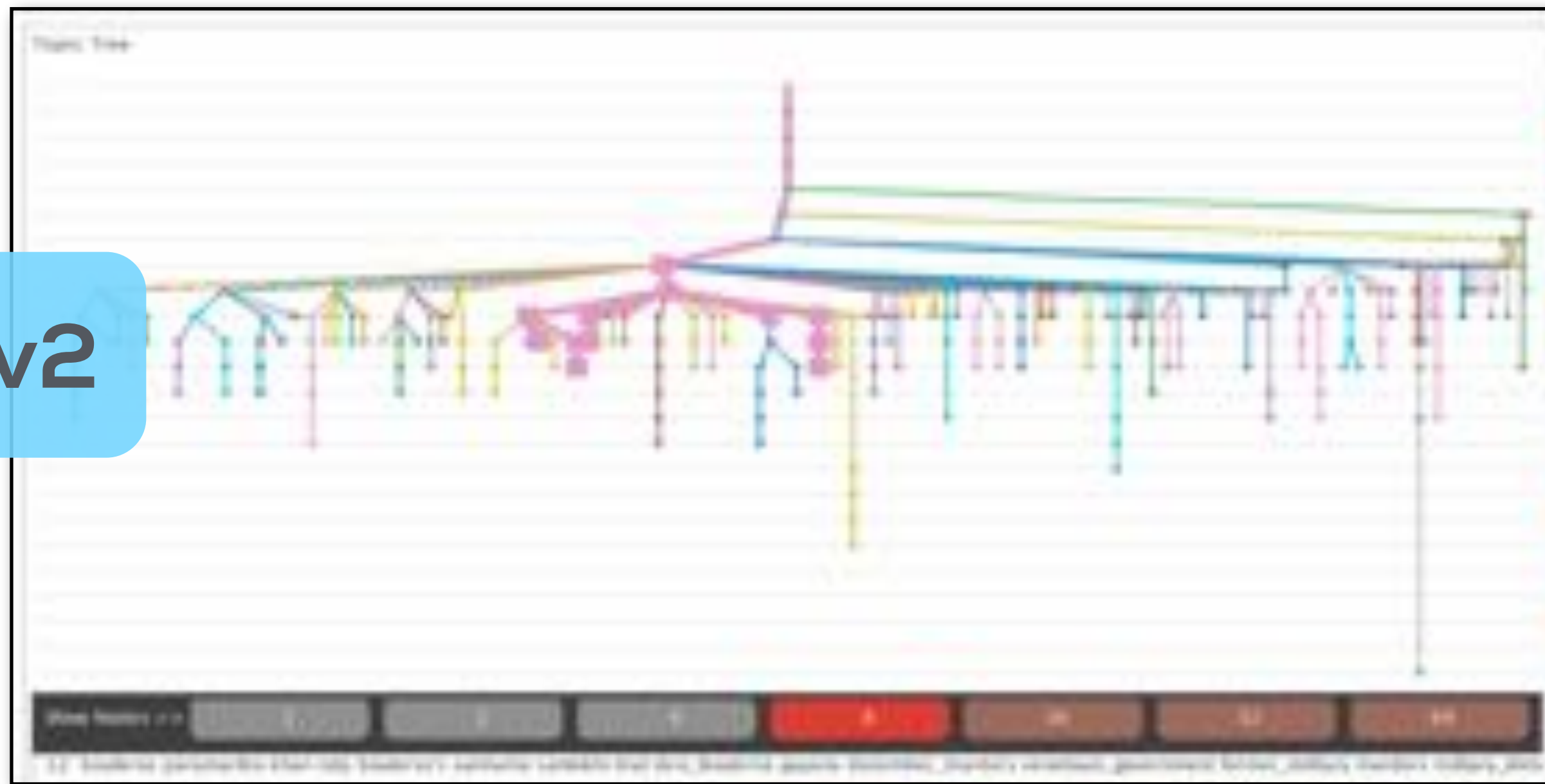




**tree vis + tagging**  
more effective  
than scatterplot

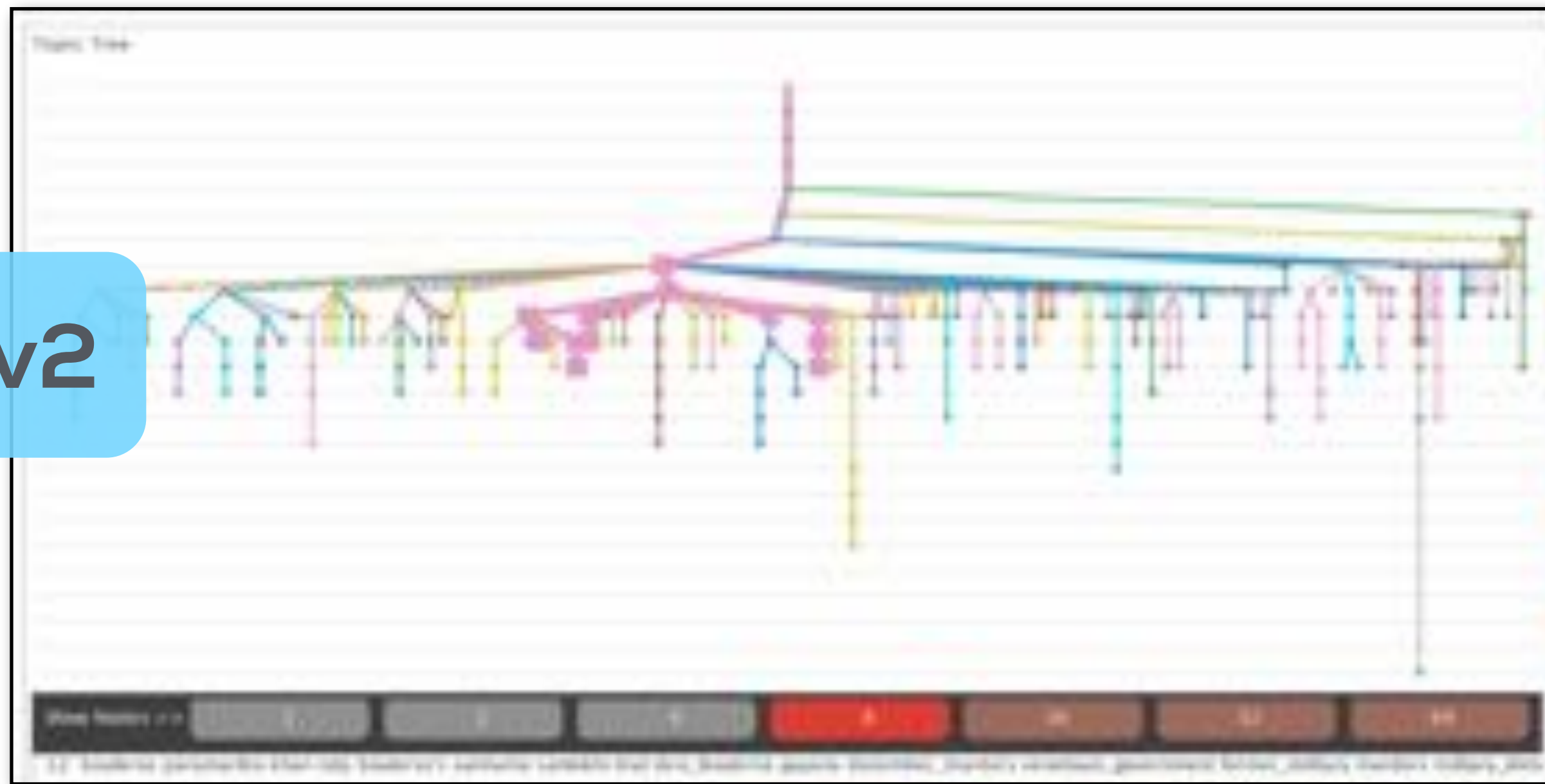


v2

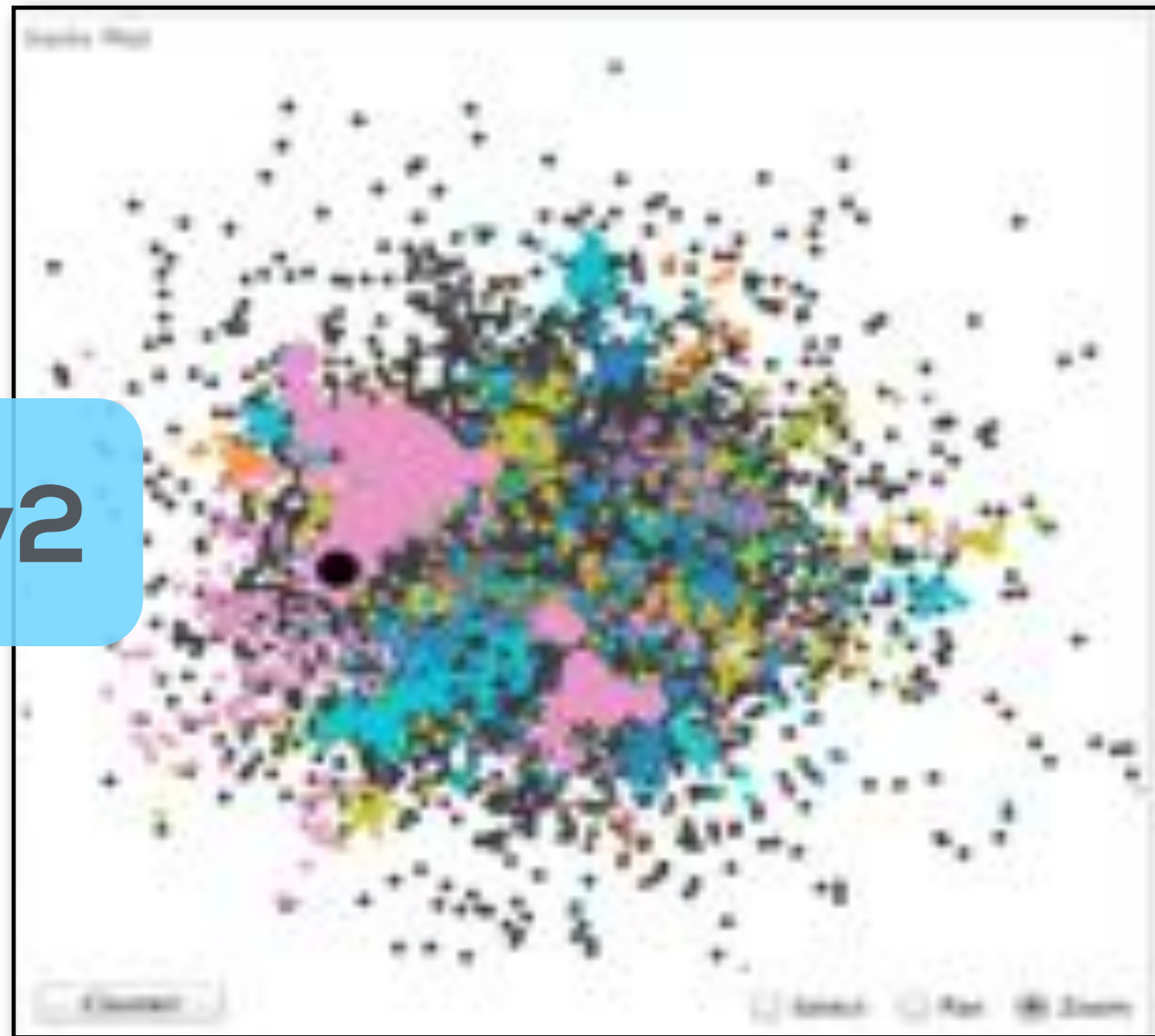




v2

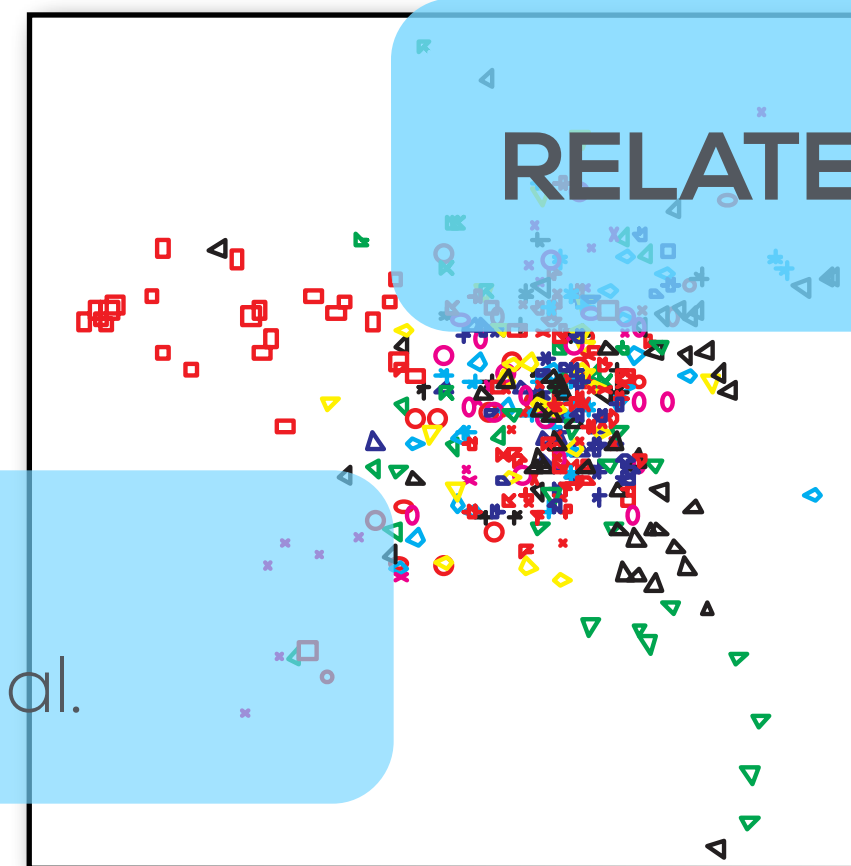


v2

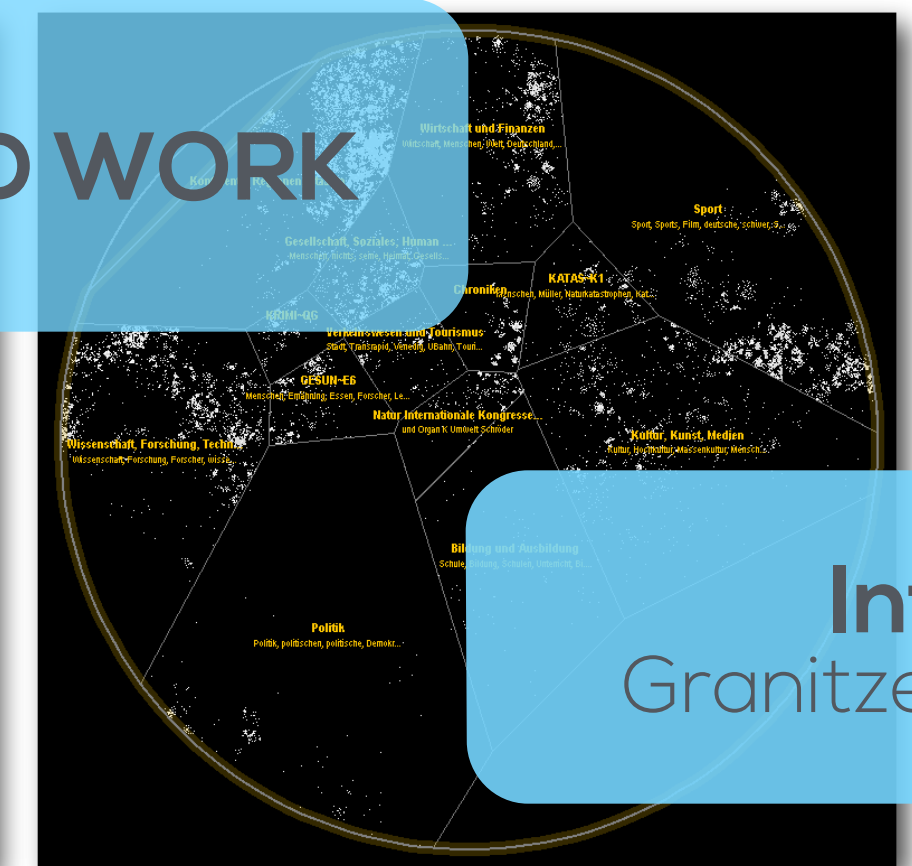


## RELATED WORK

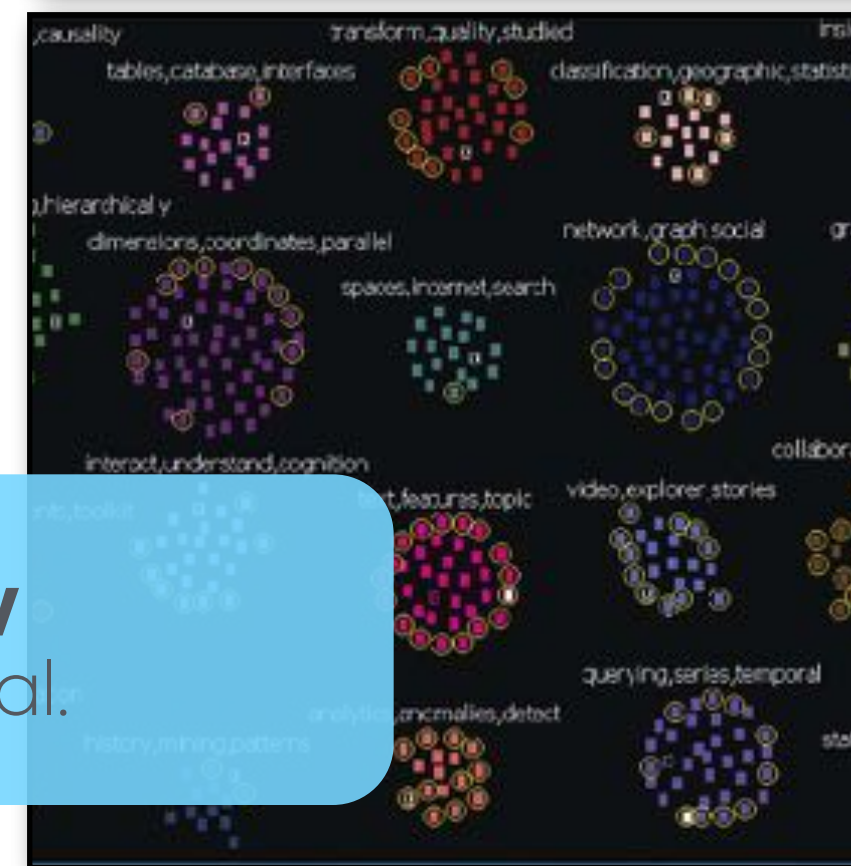
**EV**  
Chen et al.



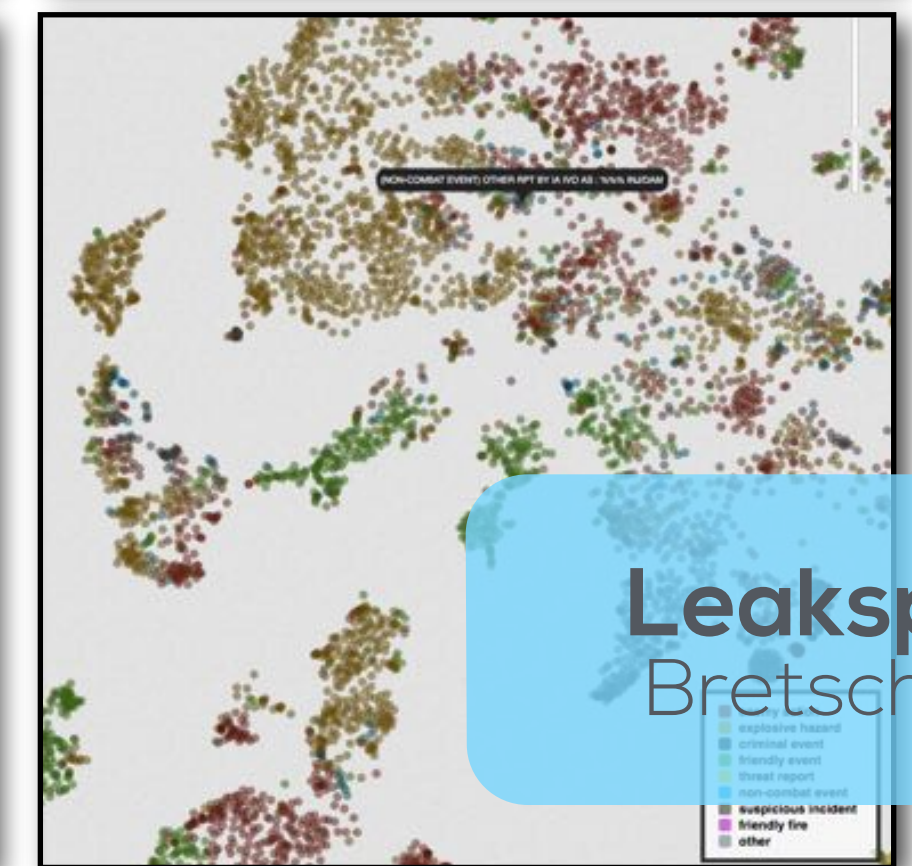
**InfoSky**  
Granitzer et al.



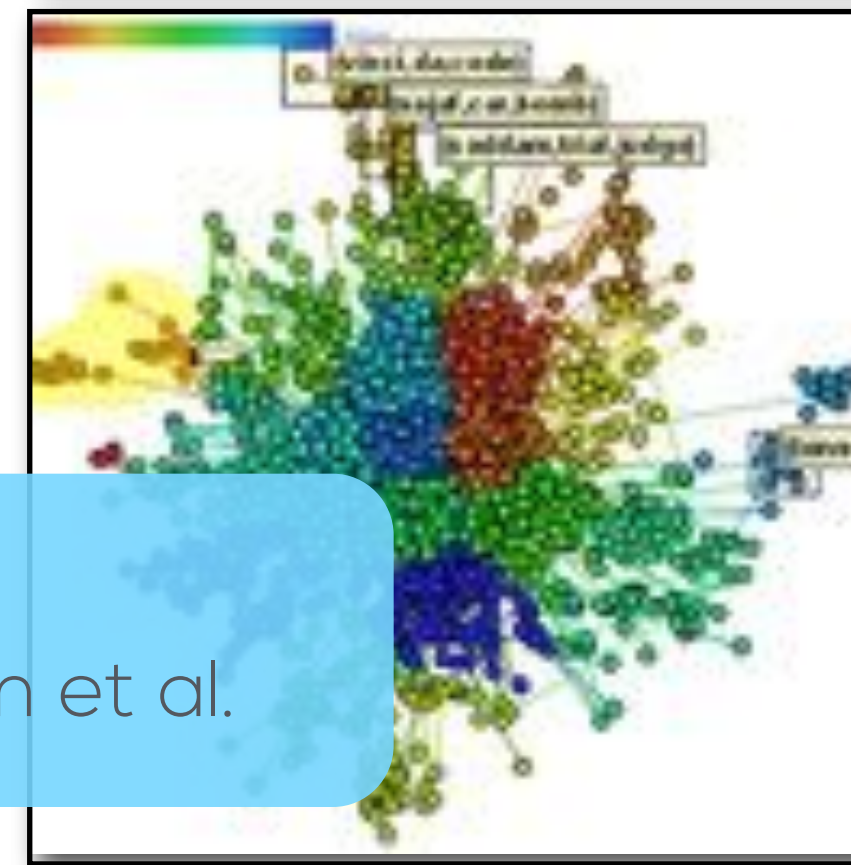
**Jigsaw**  
Görg et al.



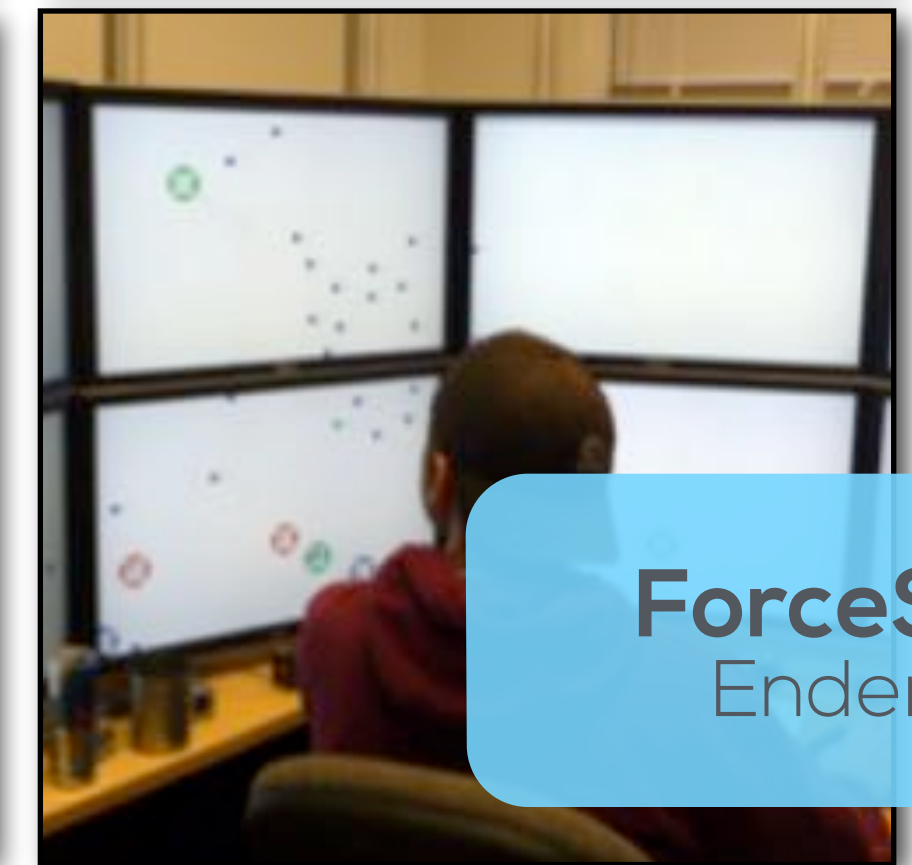
**Leaksplorer**  
Bretschneider



**PEx**  
Paulovich et al.

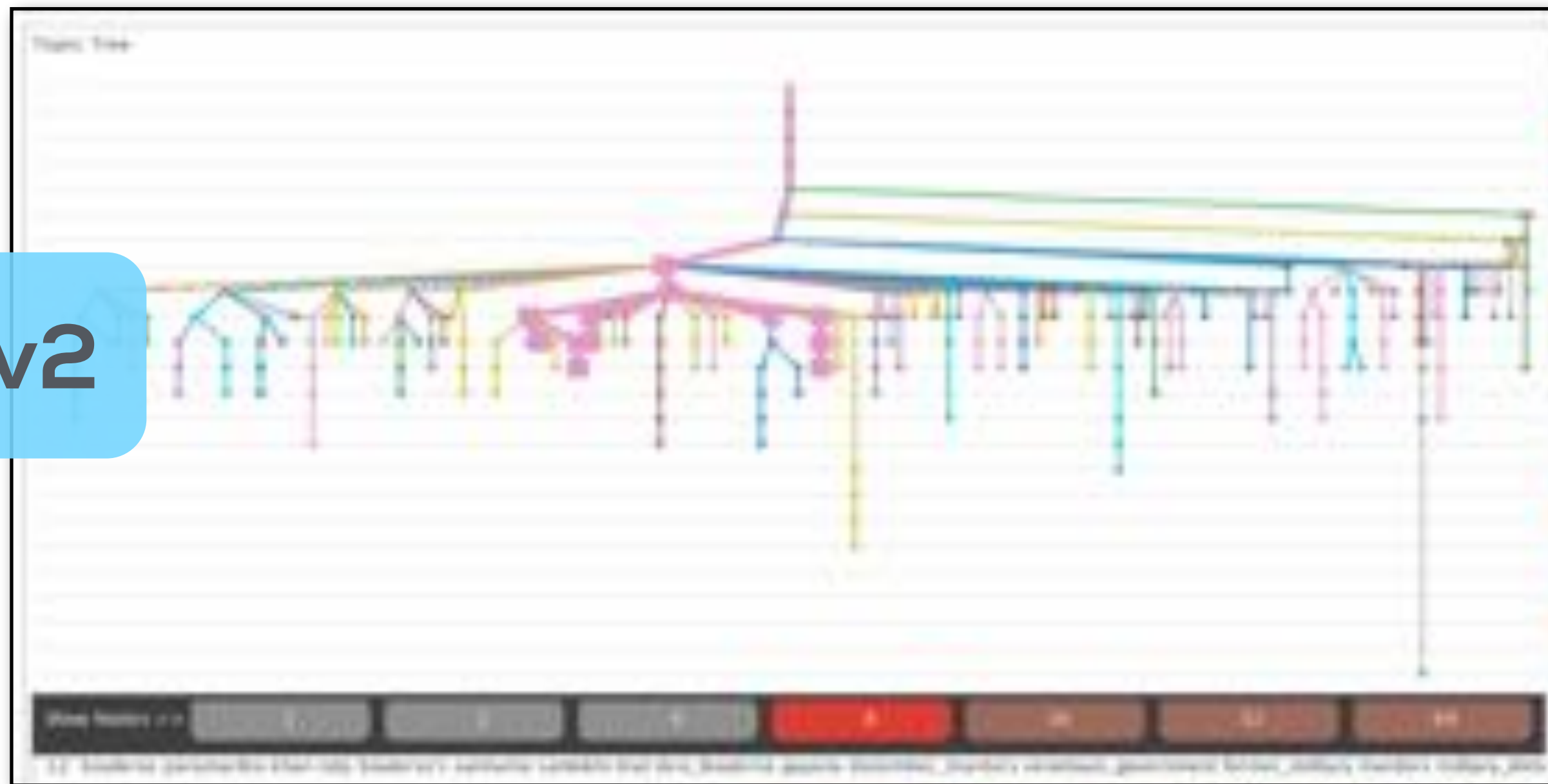


**ForceSPIRE**  
Endert et al.

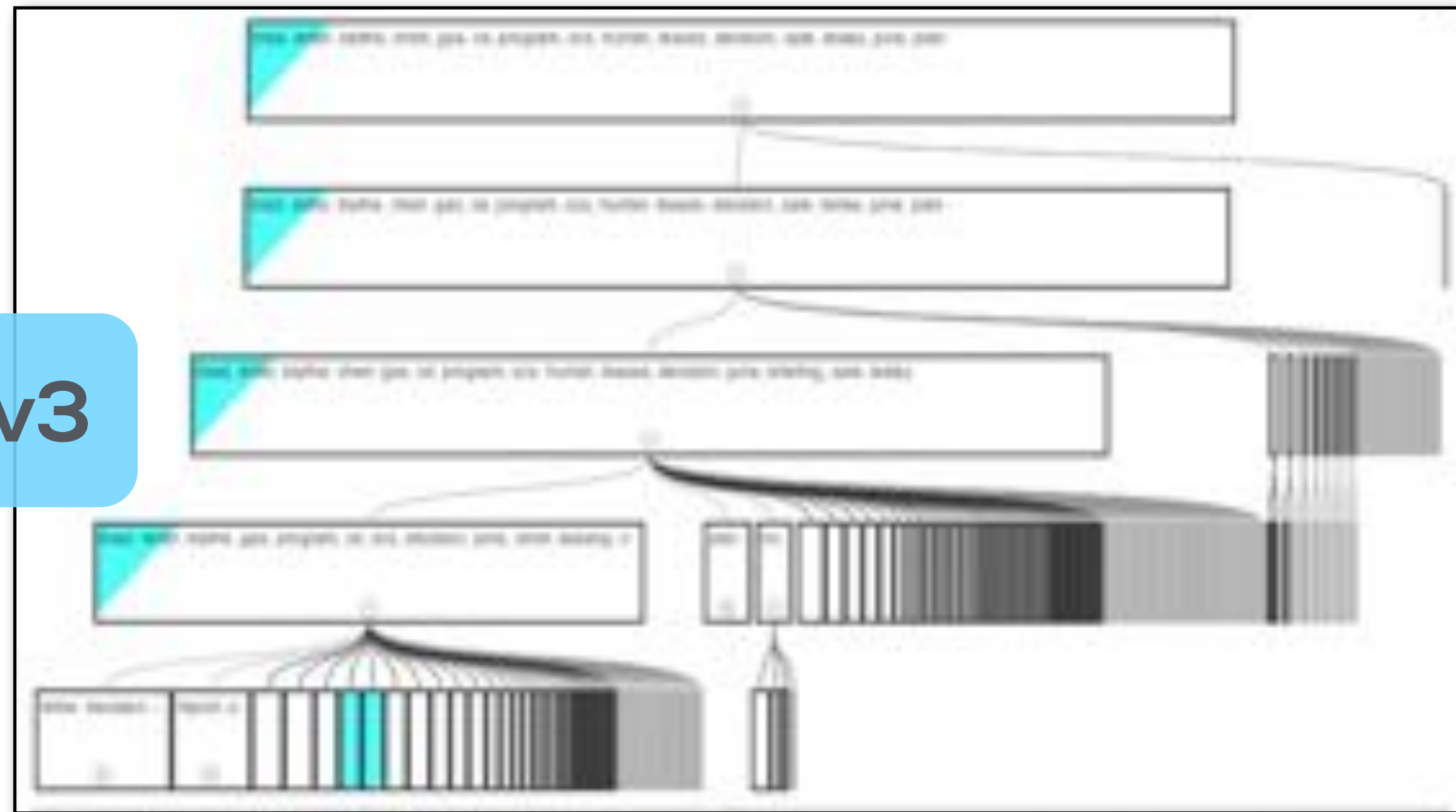




v2

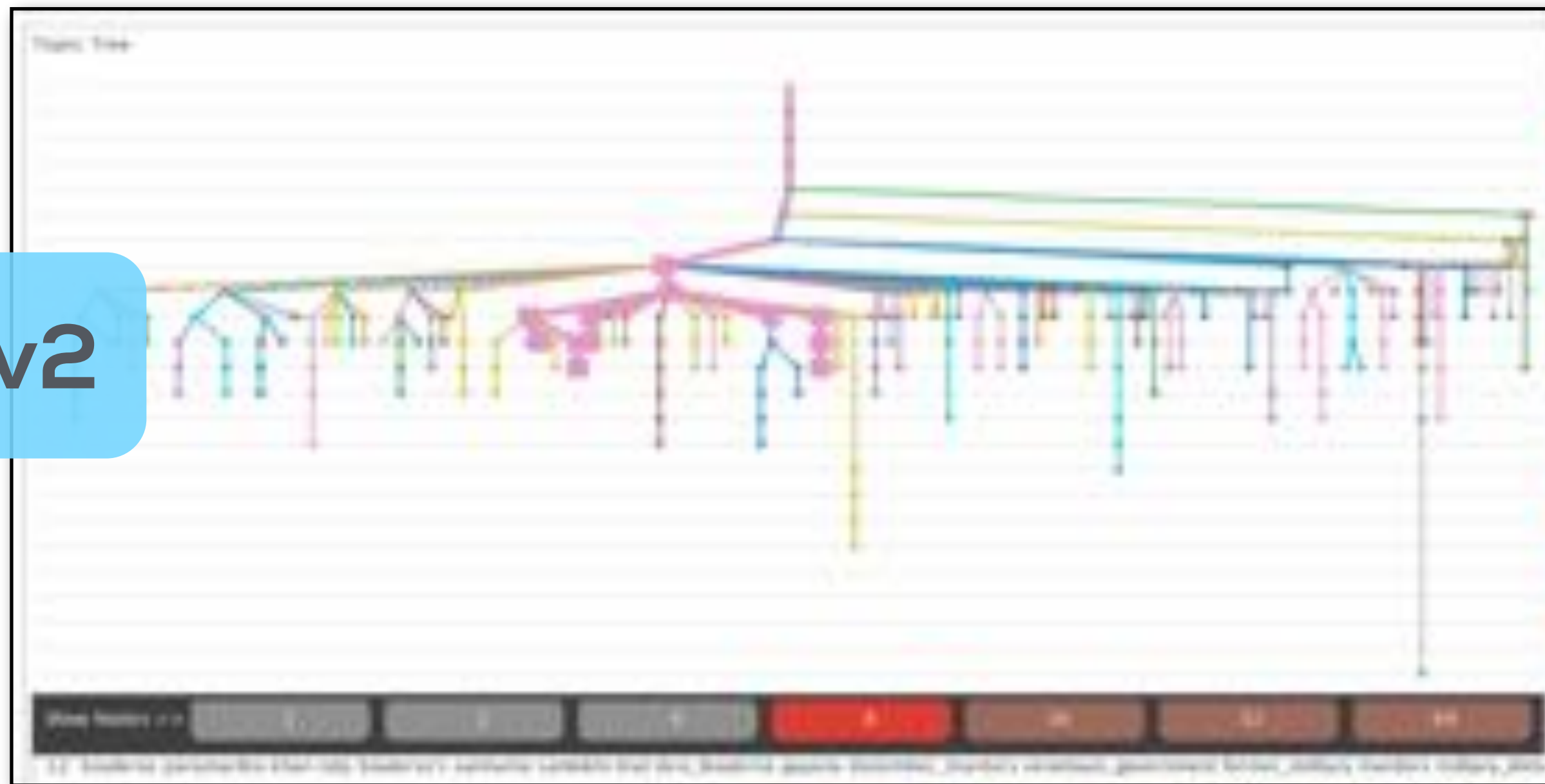


v3

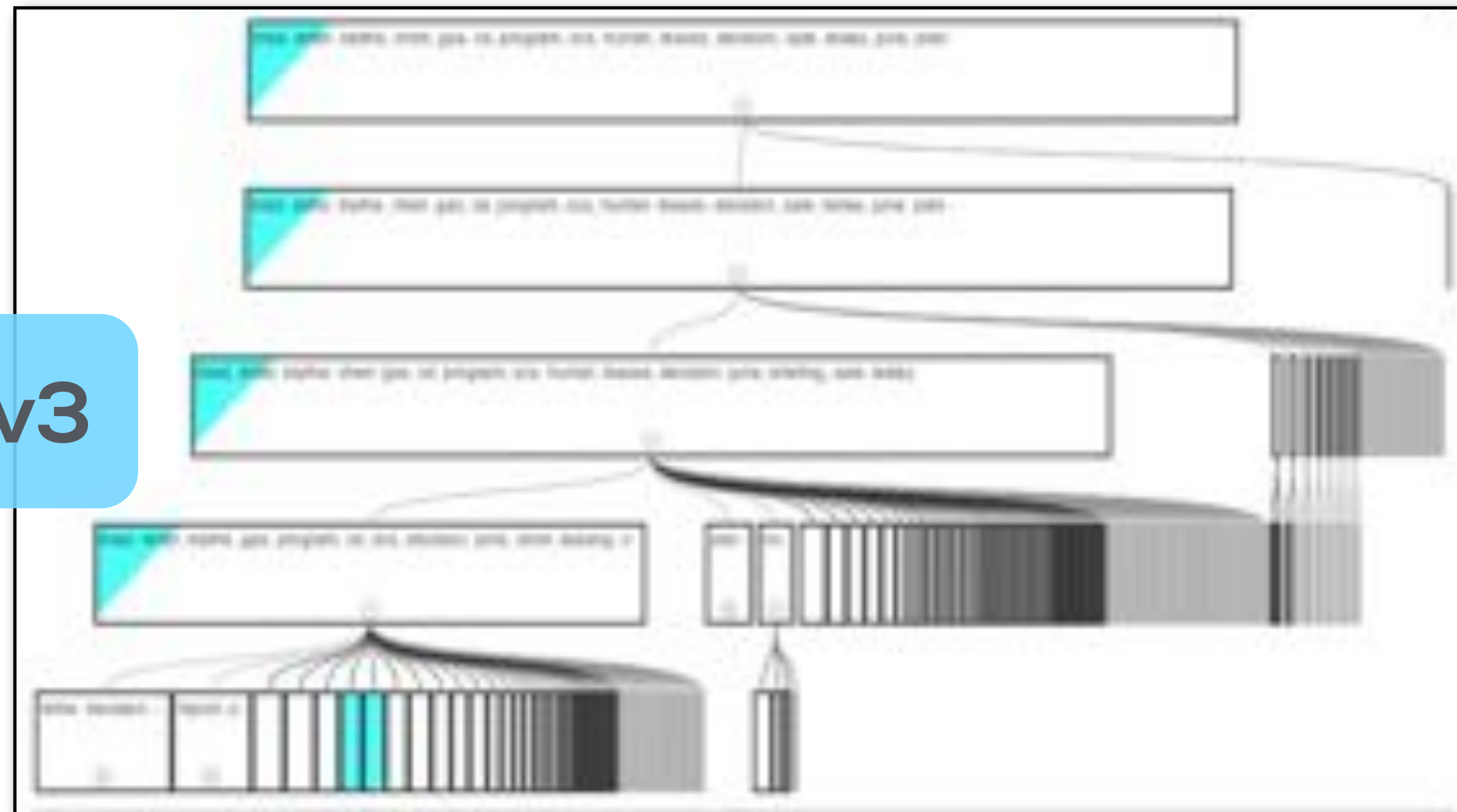




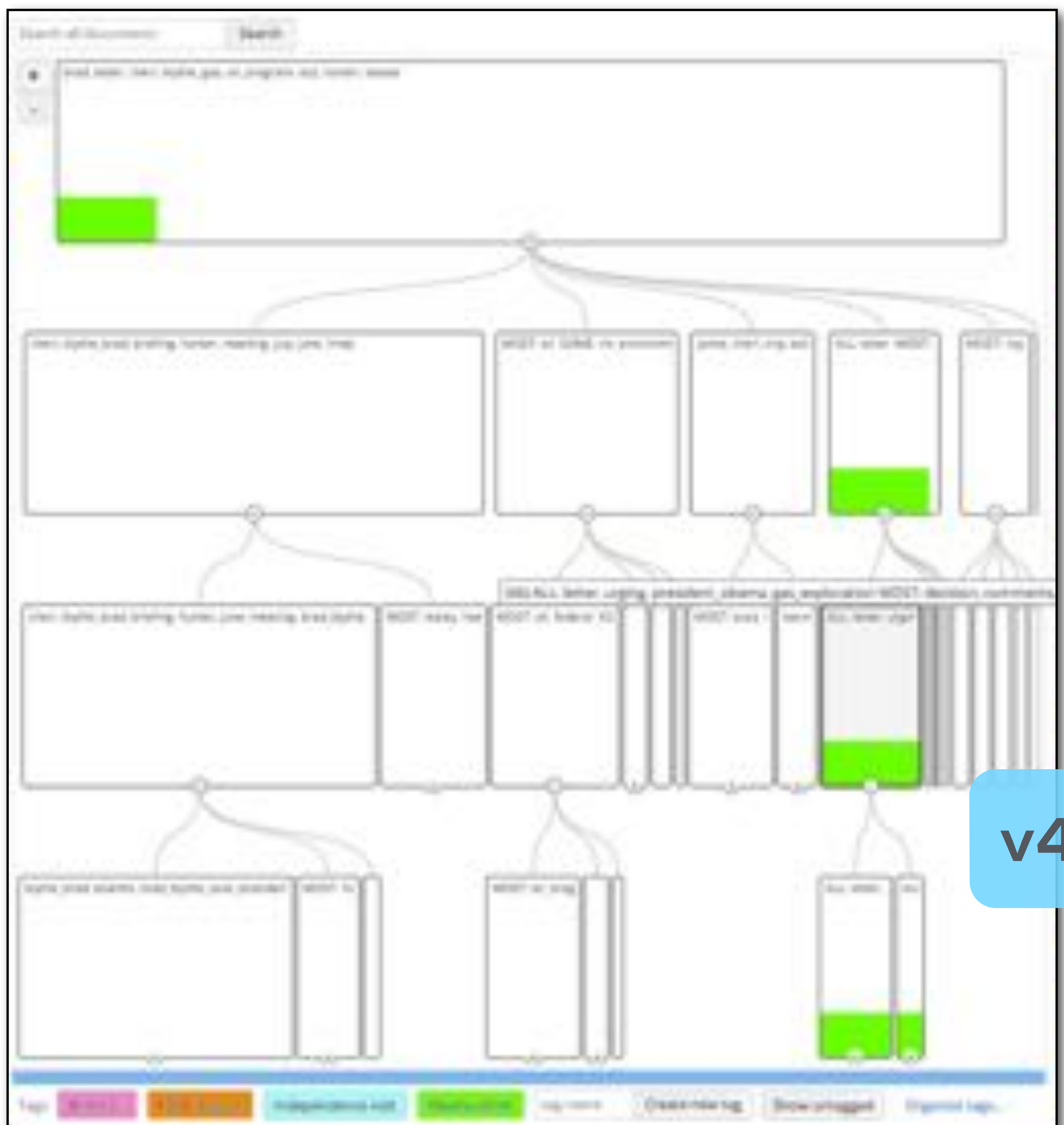
v2



v3



v4









v4 cluster  
consistency labels



'Show  
v4  
untagged'  
button

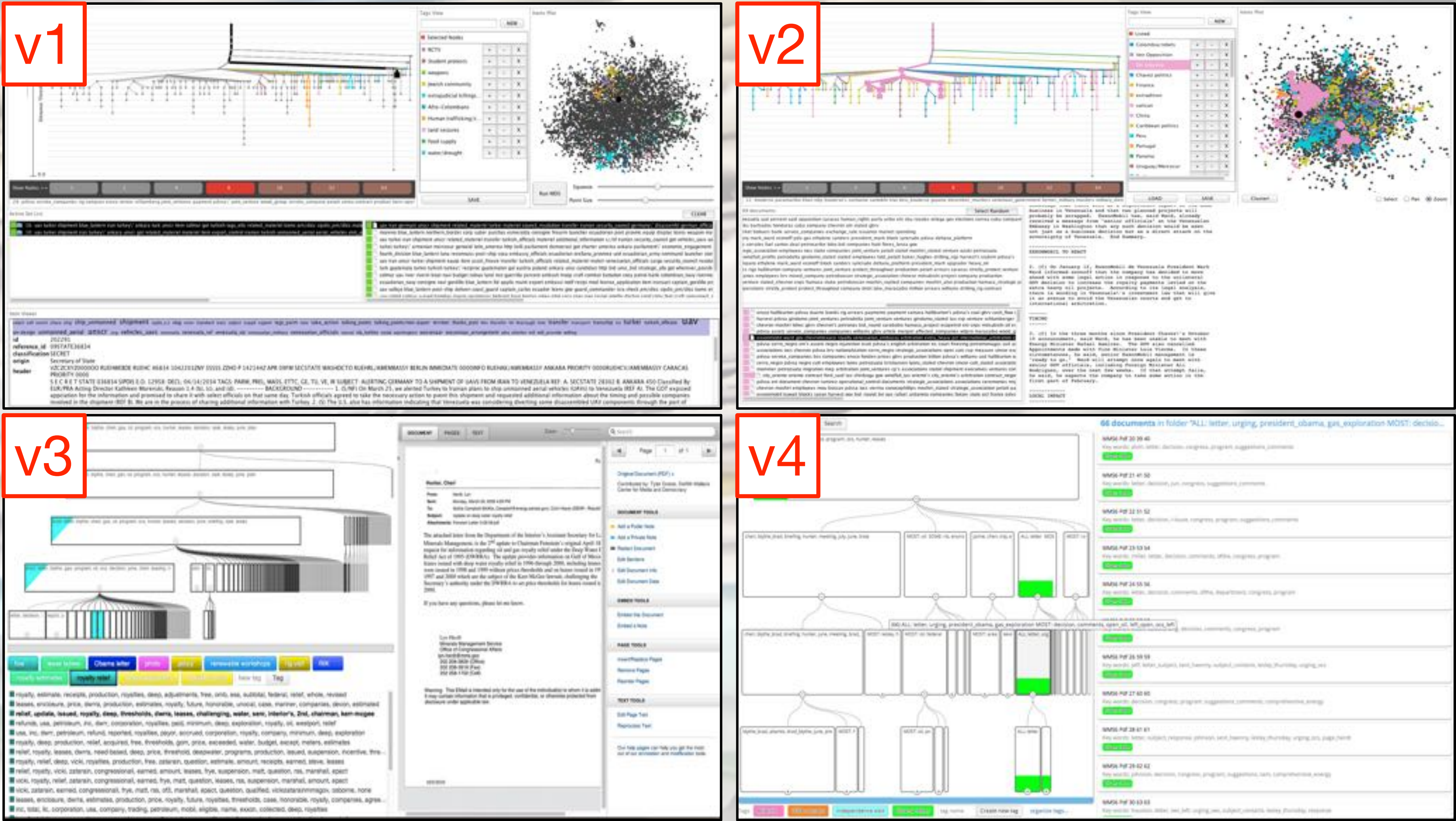


**simplify** for  
infrequent use,  
**reduce** data  
wrangling

## LESSON #4



# SIMPLIFY, REDUCE WRANGLING





# SIMPLIFY, REDUCE WRANGLING





# SIMPLIFY, REDUCE WRANGLING



[documentcloud.org](http://documentcloud.org)





image:

# LESSONS





# LESSONS

## #1: study adoption



# LESSONS

**#1:** study adoption

**#2:** (re)define the tasks



# LESSONS

**#1:** study adoption

**#2:** (re)define the tasks

**#3:** trees & tags, not points



## LESSONS

**#1:** study adoption

**#2:** (re)define the tasks

**#3:** trees & tags, not points

**#4:** simplify, less wrangling



# ACKNOWLEDGEMENTS

## research team

Jonathan



Tamara



Stephen



Matthew



## development team

Jonas Karlsson



Adam Hooper



## journalists

Jonathan Stray



Jarrel Wade



Jack Gillum



Michael Keller



(anonymous)



Adam Playford





Matthew Brehmer  
@mattbrehmer



Stephen Ingram  
@FroweFace



Jonathan Stray  
@jonathanstray



Tamara Munzner



# Private memo reveals winding tale involving John McCain, the NRA and ... condors

by *Nancy Watzman* | SEPT. 18, 2014, 1:59 P.M.



# The Brilliance of Louis C.K.'s Emails: He Writes Like a Politician

Where campaign strategy and comedy marketing collide

ADRIENNE LAFRANCE | JUL 16 2014, 6:10 AM ET



# Surprise! Many credit card agreements allow repossession

Analysis: 'Security interest' clause present on 200 cards

By Fred O. Williams



[overviewproject.org](http://overviewproject.org)  
blog: [overview.ap.org](http://overview.ap.org)

[@overviewproject](https://twitter.com/overviewproject)  
[github.com/overview](https://github.com/overview)

thanks: M. Borkin, J. Dawson, J. Ferstay, J. Fulda, S.-H. Kim, H. Lam, J. McGrenere, R. Rensink, M. Sedlmair







# SUPPLEMENTAL MATERIAL



Matthew Brehmer



Stephen Ingram



Jonathan Stray



Tamara Munzner

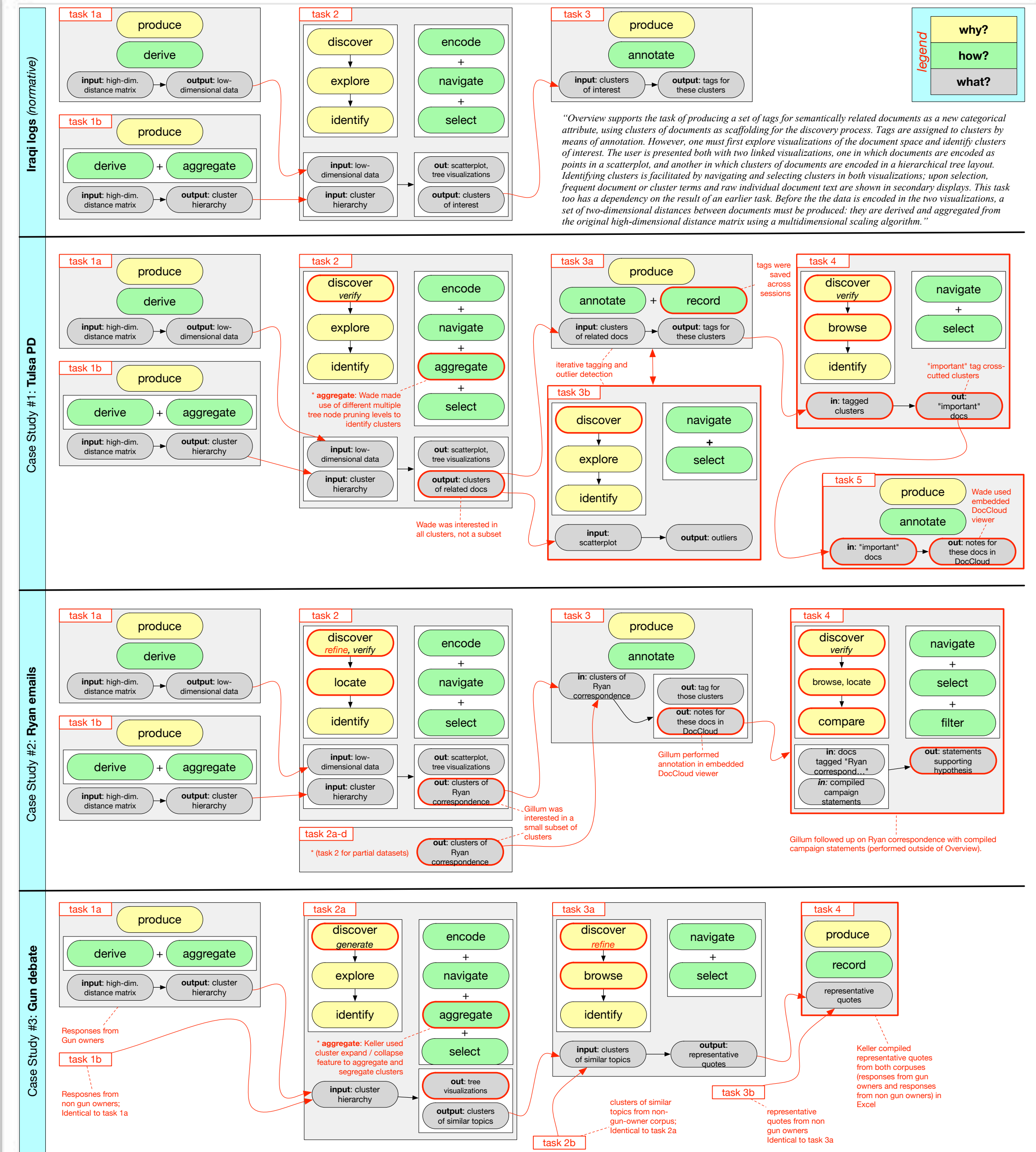
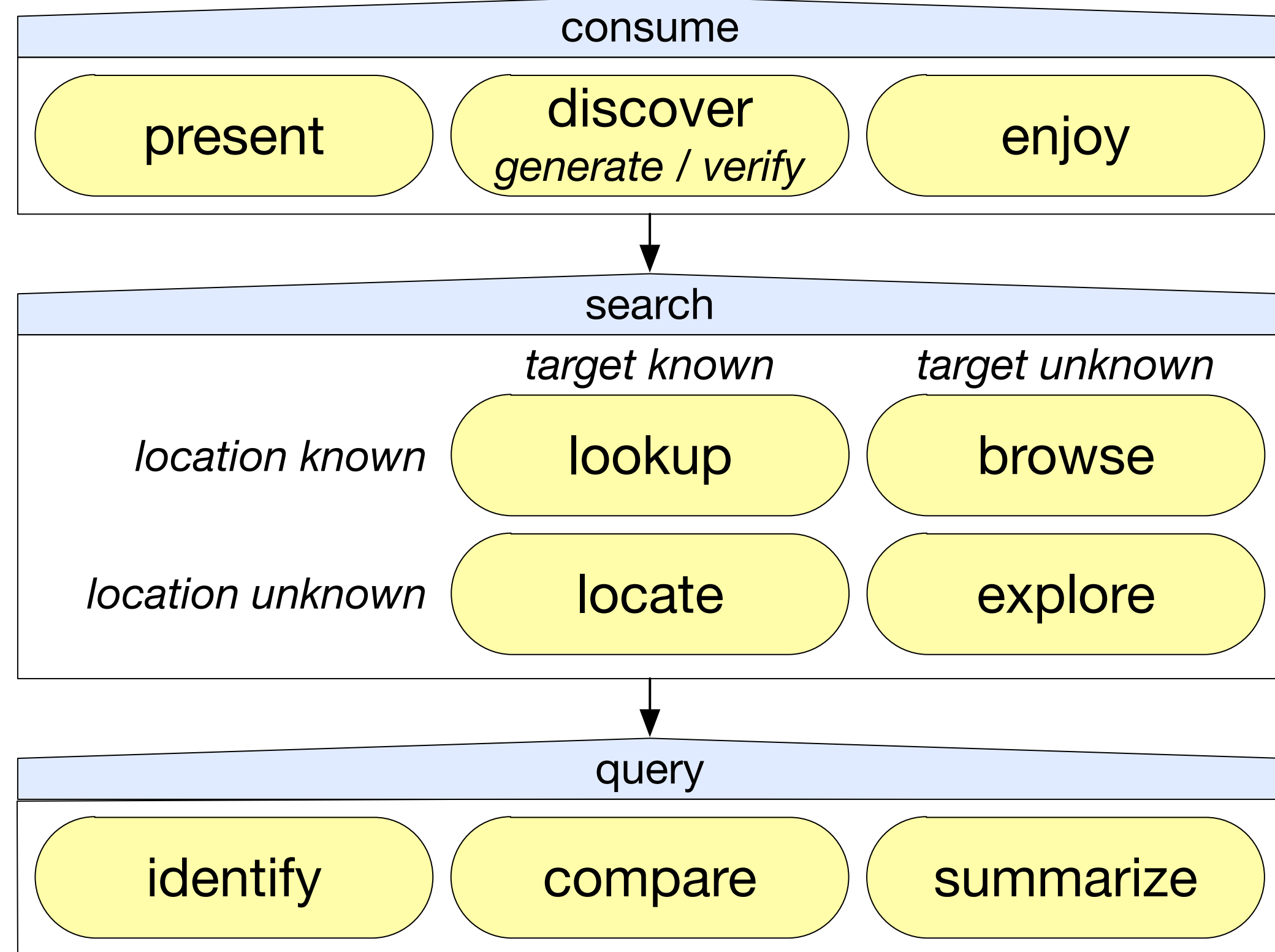


# overview

the design, **re-design, deployment, adoption, analysis, re-design, deployment, adoption, analysis, re-design, deployment, adoption, and analysis** of a visual document mining tool for investigative journalists: **a visualization design study**



# VIS TASK ANALYSIS





# OVERVIEW VS. JIGSAW

Görg et al. (IEEE TVCG 2013)  
Kang & Stasko (IEEE TVCG / VAST 2012)

**Domain**

**Data  
Abstractions**

**Scale**

**Complexity**



## OVERVIEW VS. JIGSAW

Görg et al. (IEEE TVCG 2013)  
Kang & Stasko (IEEE TVCG / VAST 2012)

### OVERVIEW

<b>Domain</b>	Investigative Journalism
<b>Data Abstractions</b>	Hierarchical Clustering of documents (raw text only), tags,
<b>Scale</b>	Tens of thousands of documents
<b>Complexity</b>	Dual view, web app



## OVERVIEW VS. JIGSAW

Görg et al. (IEEE TVCG 2013)  
Kang & Stasko (IEEE TVCG / VAST 2012)

	OVERVIEW	JIGSAW
<b>Domain</b>	Investigative Journalism	Intelligence Analysis, Law Enforcement, Academic Research
<b>Data Abstractions</b>	Hierarchical Clustering of documents (raw text only), tags,	Documents + document metadata + extracted entities +
<b>Scale</b>	Tens of thousands of documents	Hundreds to thousands
<b>Complexity</b>	Dual view, web app	Multiple view, desktop app



# OVERVIEW VS. HIERARCHICAL TOPICS

Dou et al. (IEEE TVCG / VAST 2013)

**Domain**

**Data  
Abstractions**

**Scale**

**Complexity**



# OVERVIEW VS. HIERARCHICAL TOPICS

Dou et al. (IEEE TVCG / VAST 2013)

## OVERVIEW

<b>Domain</b>	Investigative Journalism
<b>Data Abstractions</b>	Hierarchical Clustering of documents (raw text only), tags,
<b>Scale</b>	Tens of thousands of documents
<b>Complexity</b>	Dual view, web app



# OVERVIEW VS. HIERARCHICAL TOPICS

Dou et al. (IEEE TVCG / VAST 2013)

	OVERVIEW	HIERARCHICAL TOPICS
Domain	Investigative Journalism	Academic research (?)
Data Abstractions	Hierarchical Clustering of documents (raw text only), tags,	Hierarchical Clustering of documents + temporal metadata
Scale	Tens of thousands of documents	Hundreds to thousands
Complexity	Dual view, web app	Dual view, deployment details not reported





# TIMELINE & ATTRIBUTION

Jonathan



Tamara



Stephen



Matthew



Jonas Karlsson



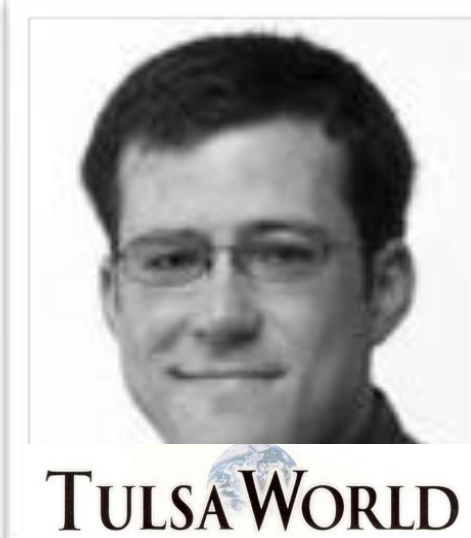
Adam Hooper



Jonathan Stray



Jarrel Wade



Jack Gillum



Michael Keller



(anonymous)



Adam Playford



2011

begin

2012

V1

V2 +  
CS 1

CS2

2013

V3

CS3

CS4

2014

V4

CS5

CS6



[overview.ap.org/completed-stories](http://overview.ap.org/completed-stories)

Private memo reveals winding tale involving John McCain,  
the NRA and ... condors

by *Nancy Watzman*

SEPT. 18, 2014, 1:59 P.M.



**Records: DHHS  
downplayed food stamp  
issues**

By Tyler Dukes



Posted: 1:00 p.m. today

**The Brilliance of Louis C.K.'s Emails:  
He Writes Like a Politician**

Where campaign strategy and comedy marketing collide

ADRIENNE LAFRANCE | JUL 16 2014, 6:10 AM ET

*the Atlantic*

**Surprise! Many credit card agreements allow repossession**

Analysis: 'Security interest' clause present on 200 cards

By [Fred O. Williams](#)

