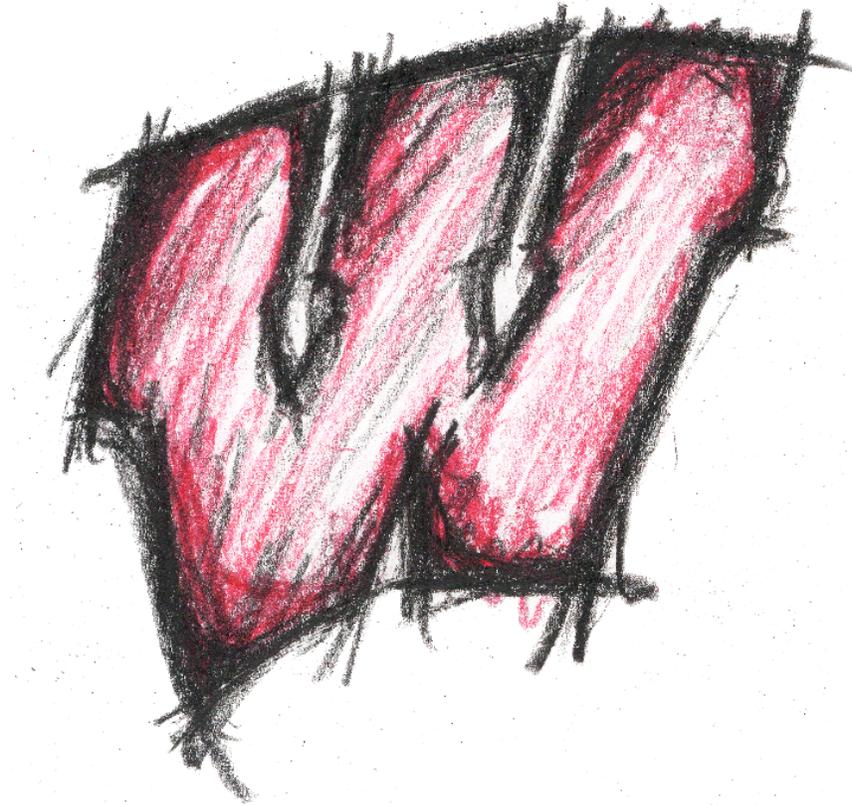


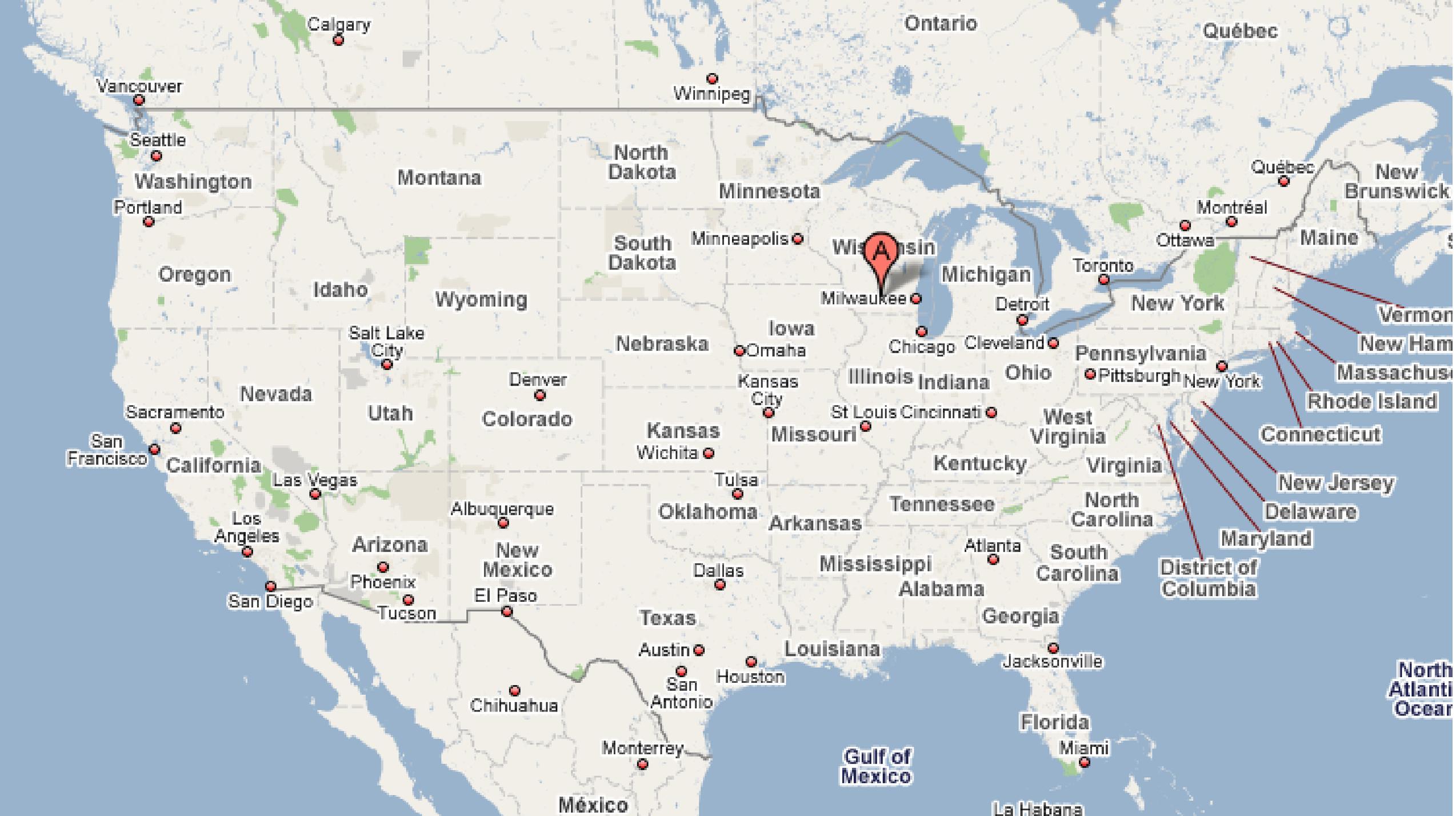
# What Shakespeare Taught Us About Visualization and Data Science

**Michael Gleicher**

Department of Computer Sciences

University of Wisconsin Madison





A

- Vermont
- New Hampshire
- Massachusetts
- Rhode Island
- Connecticut
- New Jersey
- Delaware
- Maryland
- District of Columbia

# Acknowledgements

## Visualizing English Print Project

*Co-Pis*

Michael Witmore (Folger Library)

Jonathan Hope (Strathclyde University)

Robin Valenza (UW English)

*VEP Students and Staff*

Eric Alexander, Deidre Stuffer, Erin Winter,  
Hao Fu, Joe Kohlmann, ...

Jonathan's students (Heather, Beth, ...)

Robin's students (Mattie, Cathy, ...)

## Other Vis Collaborators

*My Students (past and present)*

Danielle Szafir, Michael Correll, Alper Sarikaya, Kendall Park, ...

*Perception and HCI Collaborators*

Steve Franconeri, Bilge Mutlu,  
Chuck Hansen, and their students

*Other Domain Collaborators*

George Phillips, Dave O'Connor, Colin Dewey, ...

Lesson 1

**It's a team effort**

Lesson 1b

**I'm not the Shakespeare Expert**

So... why am I here?

Michael Gleicher  
Visual Computing Group

Human **Graphics** Interaction

authoring pictures, videos, animations

Human **Robot** Interaction

robots!

Human **Data** Interaction

visualization, visual analytics, interactive learning

~~Data Visualization~~

~~Visual Analytics~~

Human Data Interaction

## Lesson 2

**Visualization is not (just) about making pictures**

**Data Science is not (just) about doing statistics**

**They are (also) about **Human-Data Interaction****

# What Shakespeare Taught Us About Visualization and Data Science

**Michael Gleicher**

Department of Computer Sciences

University of Wisconsin Madison



**Caveats**

# This is not a talk about Shakespeare

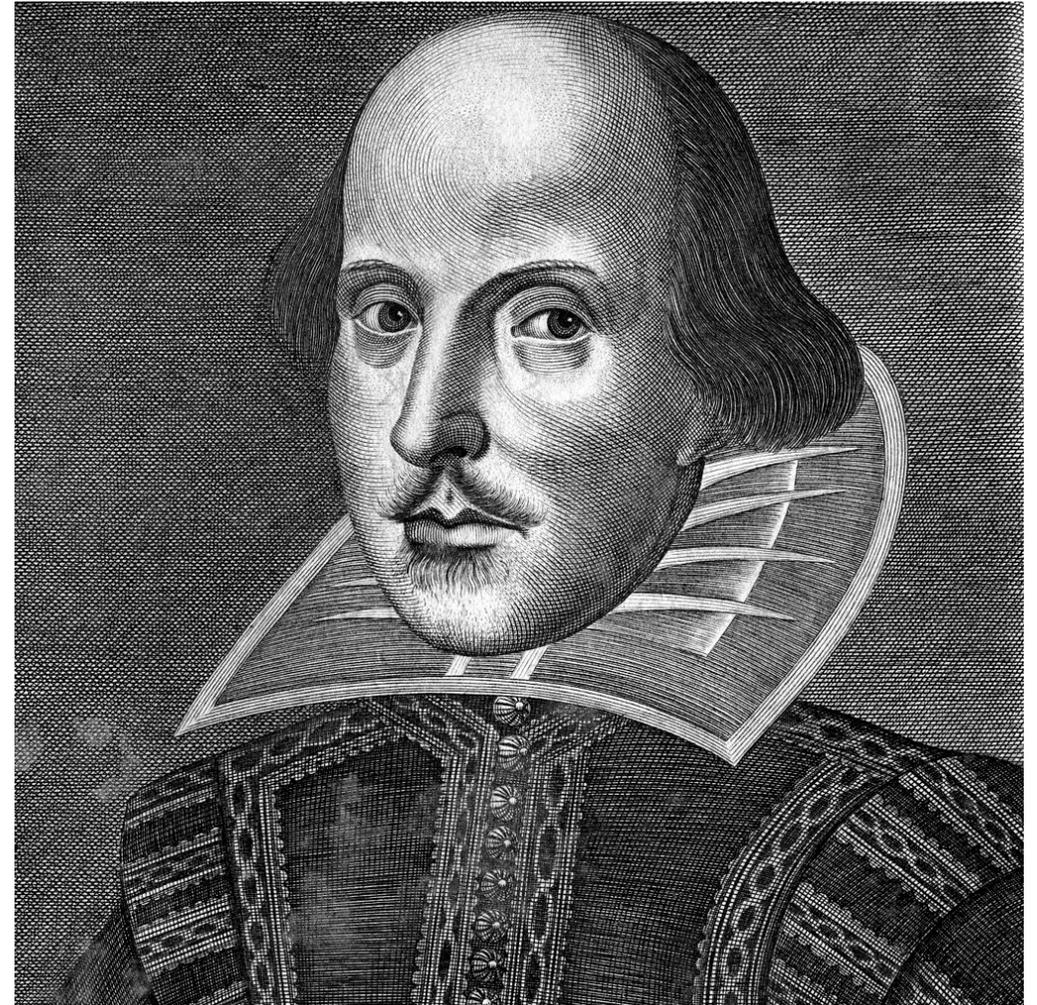
*But he makes for a catchy title*

It is a talk about a collaboration  
with **literary scholars**

With the goal of getting beyond  
Shakespeare

(Early Modern Period 1470-1700)

(Or is that 1470-1660? Or 1800?)



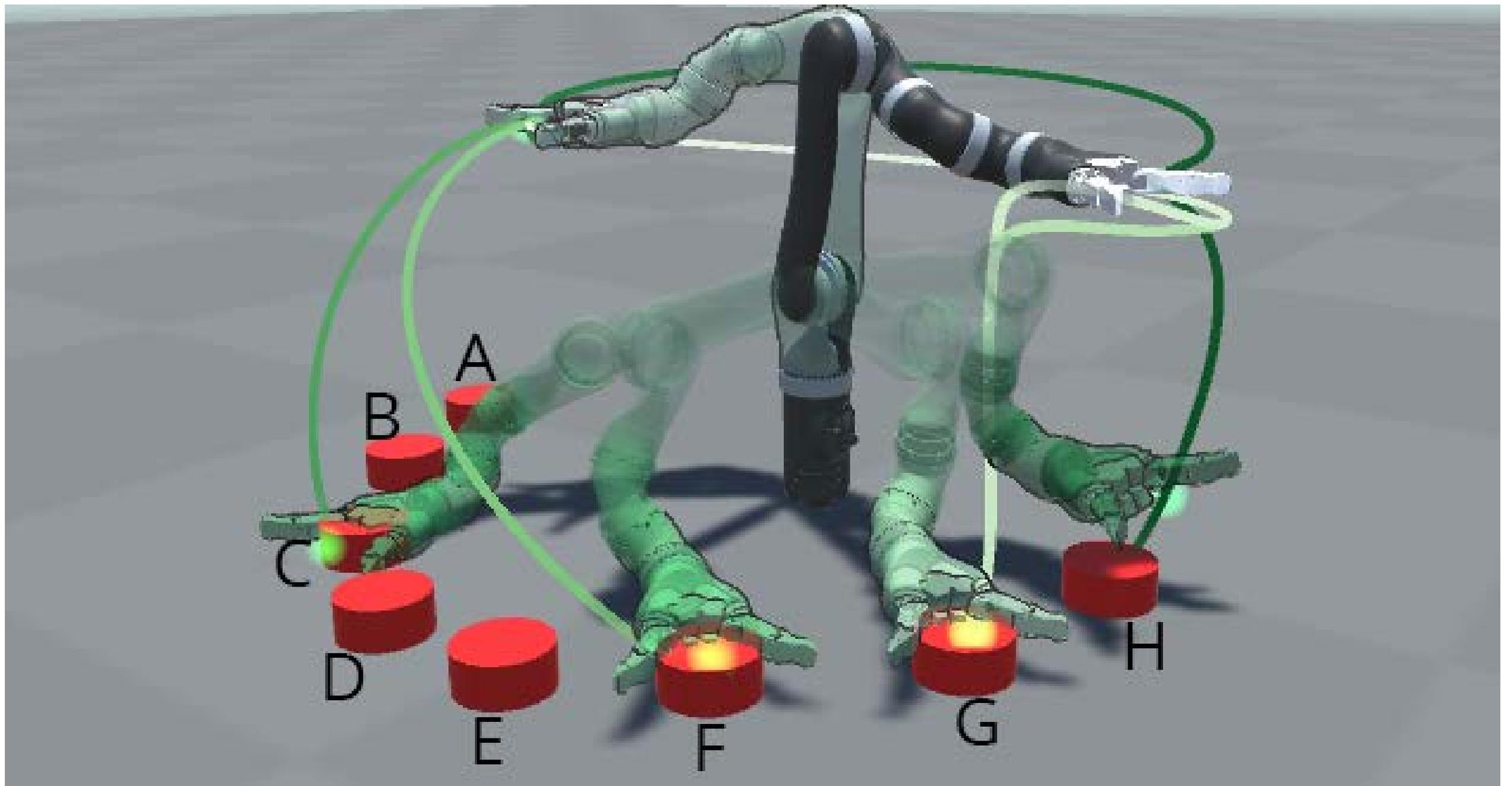
# **This is not a talk about Digital Humanities**

This is about a collaboration with literature scholars.

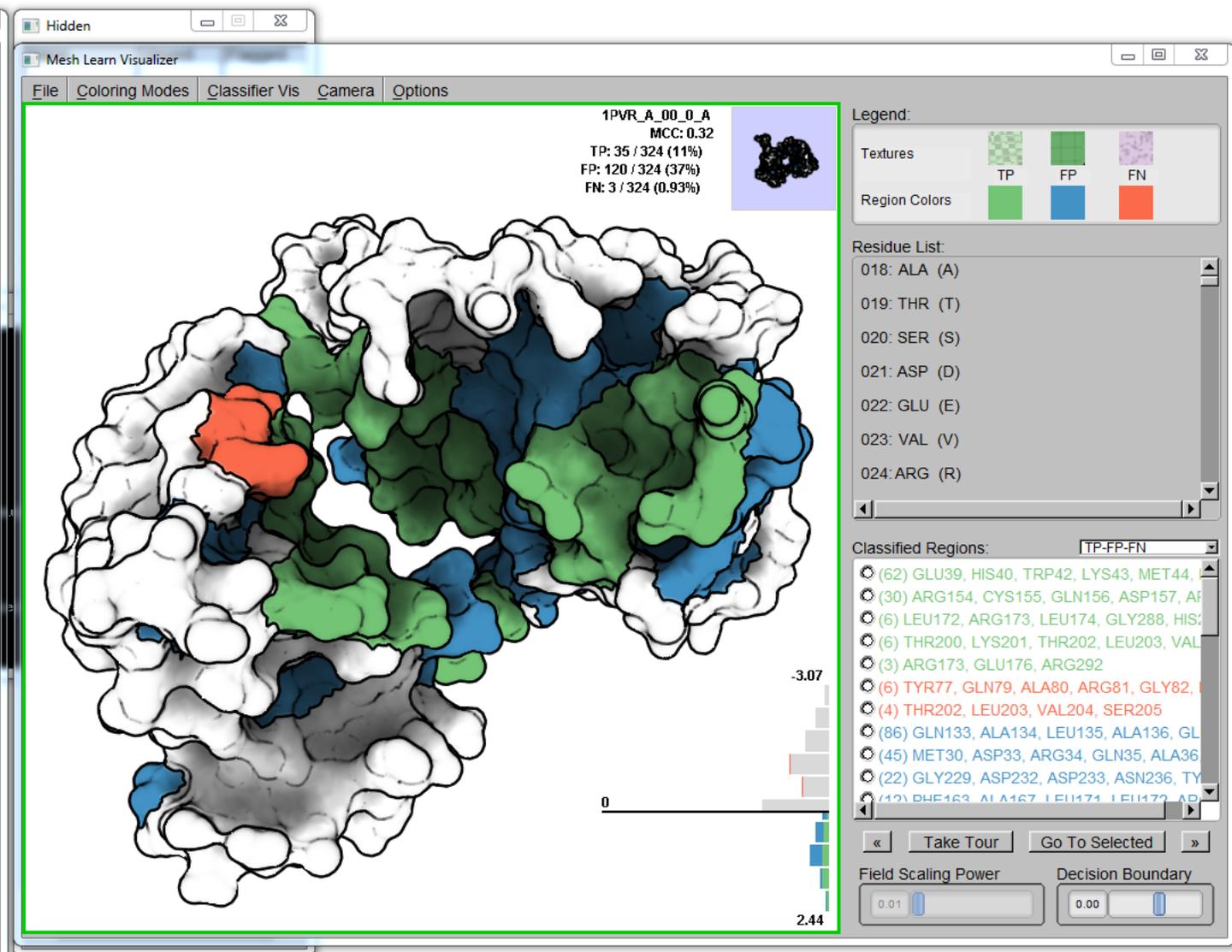
What can the rest of us learn?

# Why experience with Literary Scholarship?

*Why not experience with Biochemistry, Virology, Robotics, ... ?*



Rakita, Mutlu and Gleicher. Motion Synopsis for Robot Arm Trajectories, IEEE ROMAN 2016.



Sarikaya, Albers, Mitchel and Gleicher. *Visualizing Validation of Protein Surface Classifiers*. Computer Graphics Forum. Proceedings EuroVis 2014

## Lesson 3

**It's nice to have an application  
that people are interested in**

# Why experience with Literary Scholarship?

Some lessons we could have gotten elsewhere

*but this project made them more clear*

Some lessons came from unique aspects of the project

*but I think they are more general*

Some lessons come from the unique collaboration

*Humanist\* thinking in data analysis!*

\* I dislike the term “Humanist” because of what it implies about the rest of us. But it is how “they” self-identify.

# Mike's Mantra

Until we take the time to learn about how the other side thinks, we can't really work **together**.

Once we learn how each other thinks, our ways of thinking can infuse each other's.

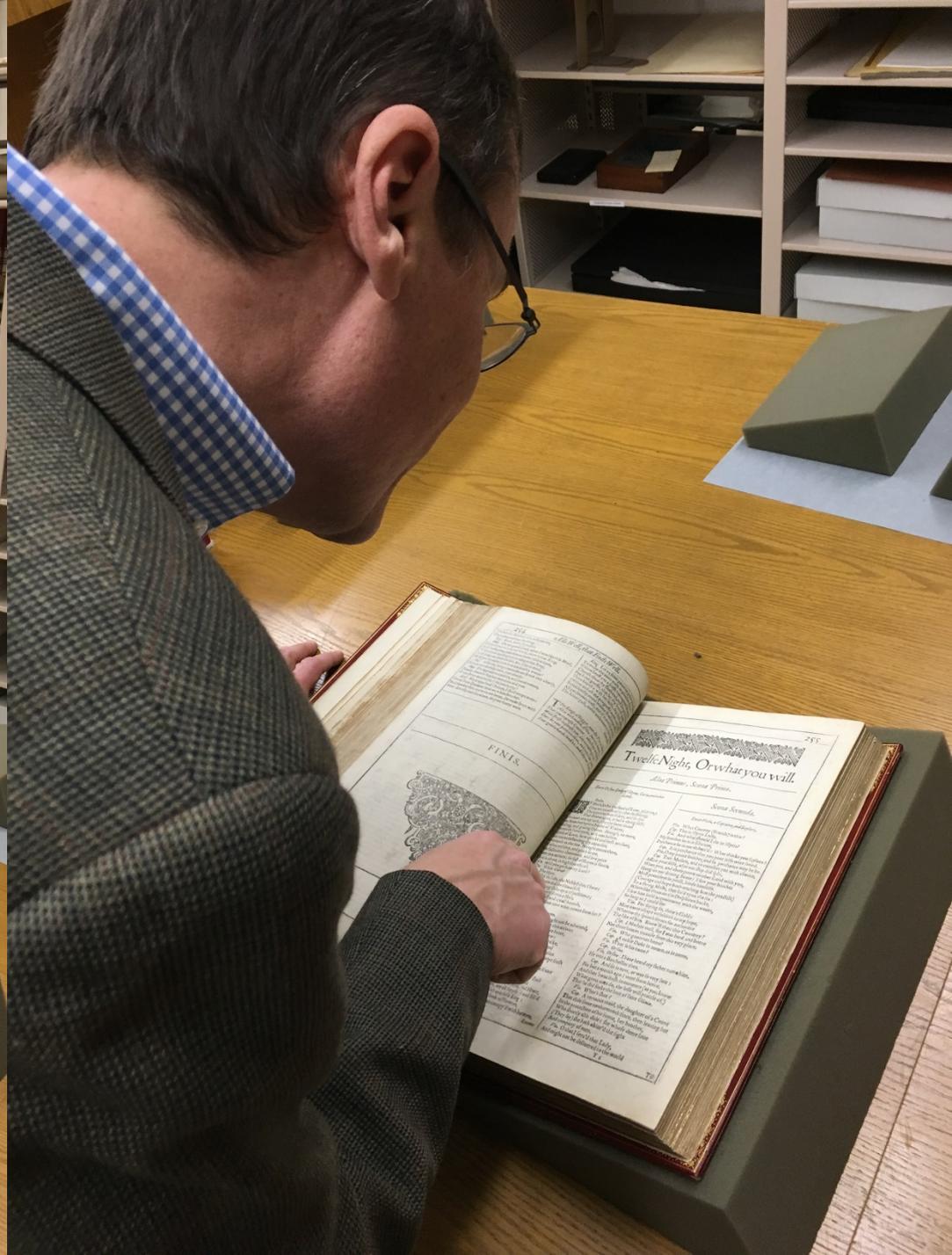
This is not just building tools for our friends.

It's a **lot** more fun and interesting

# Working with Literary Scholars



Folger Shakespeare Library





# THE TRAGEDIE OF MACBETH.

*Actus Primus. Scœna Prima.*

*Thunder and Lightning. Enter three Witches.*

**W**hen shall we three meet againe?  
In Thunder, Lightning, or in Raine?  
2. When the Hurley-burley's done,  
When the Battaille's lost, and wonne.  
3. That will be ere the set of Sunne.  
1. Where the place?

2. Vpon the Heath.  
3. There to meet with *Macbeth*.  
1. I come, *Gray-Malkin*.  
*All. Paddock* calls anon: faire is foule, and foule is faire,  
Honet through the fogge and filthie ayre. *Exeunt.*

*King.* O valiant Cousin, worthy Gentleman,  
*Cap.* As whence the Sunne 'gins his reflection,  
Shipwracking Stormes, and direfull Thunders:  
So from that Spring, whence comfort seem'd to come,  
Discomfort swells: Marke King of Scotland, marke,  
No sooner Iustice had, with Valour arm'd,  
Compell'd these skipping Kernes to trust their heeles,  
But the Norway Lord, surneying vantage,  
With furbusht Armes, and new supplies of men,  
Began a fresh assault.

*King.* Dismay'd not this our Captaines, *Macbeth* and  
*Banquoh*?

*Cap.* Yes, as Sparrowes, Eagles;  
Or the Hare, the Lyon:  
If I say sooth, I must report they were  
As Cannons ouer-charg'd with double Cracks,  
So they doubly redoubled stroakes vpon the Foe:  
Except they meant to bathe in reeking Wounds,  
Or memorize another *Golgotha*,  
I cannot tell: but I am faint,  
My Gashes cry for helpe.

*King.* So well thy words become thee, as thy wounds,

When shall we three meet againe?  
In Thunder, Lightning, or in Raine?

# Some Lessons from “Humanist” Thinking

*They got along fine before “data-centric” thinking*

Importance of **exemplars** and **outliers**

Importance of going back to the **source** (specific passages)

Value of **multiple** points of view

Contextualized arguments with lots of **background**

**Editing** and **curation** are scholarly activities

# Visualizing English Print

c 1470-~~1800~~ (1700)

What if you had access to all surviving books?

# Visualizing English print (VEP)

Large collections of texts becoming available to scholars

How to enable (traditional) literary scholars to use it?

## **Problem Factory:**

What questions should/could one ask?

How do you ask those questions?

What tools to provide to help?

*What can we learn about Visualization / Data Science?*

# Why? (what is in it for literature)

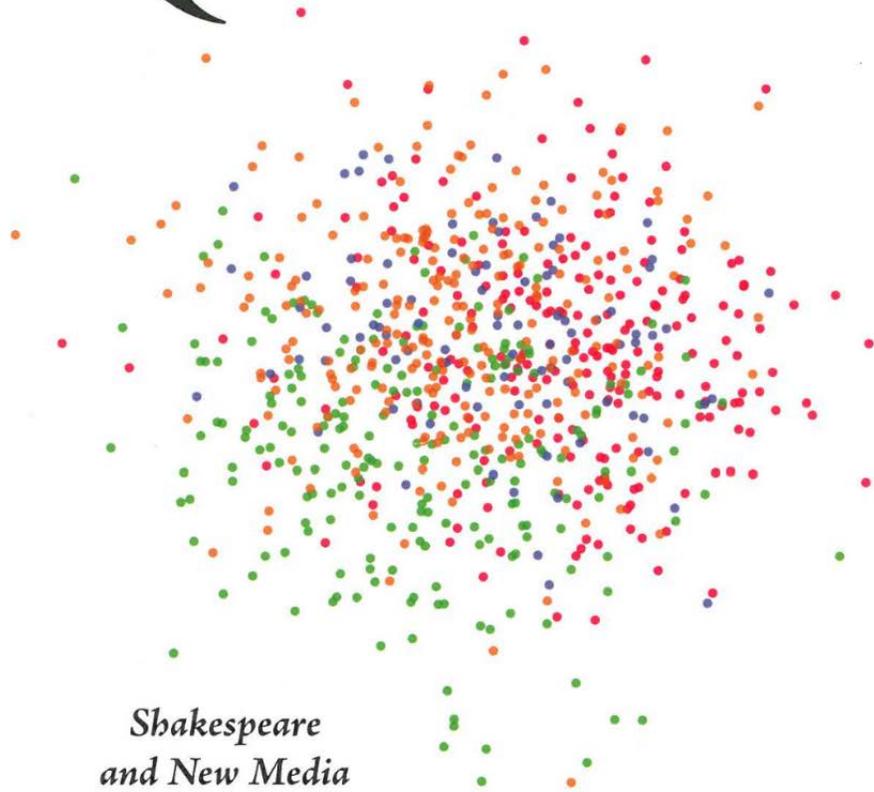
Consider larger collections of books  
beyond the “canon”

See language develop (independent of content)

See developments over decades

See patterns too subtle for people to pick up

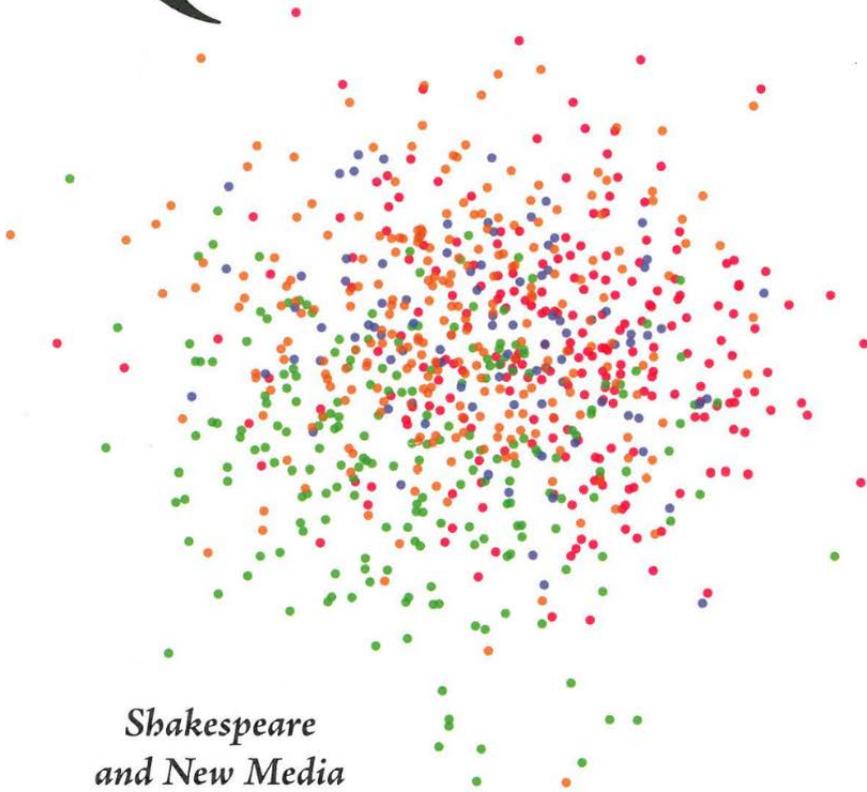
# SHAKESPEARE QUARTERLY



*Shakespeare  
and New Media*

Published for the Folger Shakespeare Library  
in association with  
The George Washington University  
by The Johns Hopkins University Press

# SHAKESPEARE QUARTERLY



*Shakespeare  
and New Media*

Published for the Folger Shakespeare Library  
in association with  
The George Washington University  
by The Johns Hopkins University Press

**One journal cover image  
leads to (at least) three challenges**

The axes are meaningless!

Explainers – crafted projections  
VAST 2013

Can people interpret this?

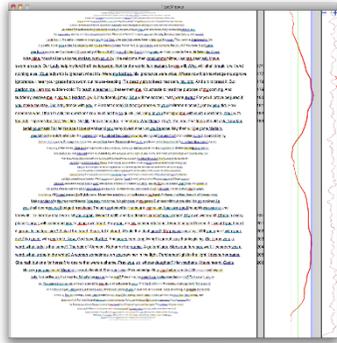
Perception of average value in scatterplots  
InfoVis 2013

The scatterplot has too many points!

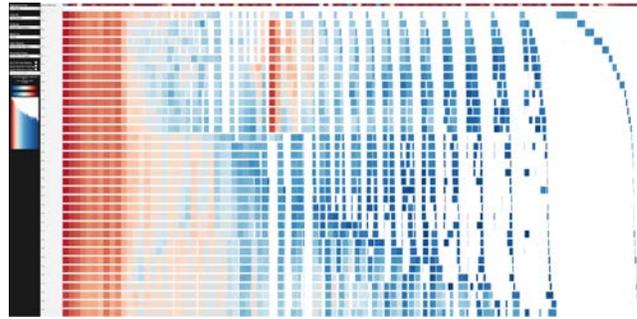
Splatterplots – scalable scatterplots  
TVCG 2013

**And these lead to many more...**

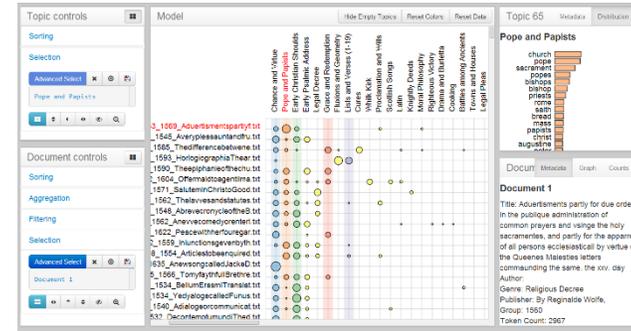
# Vis research success stories



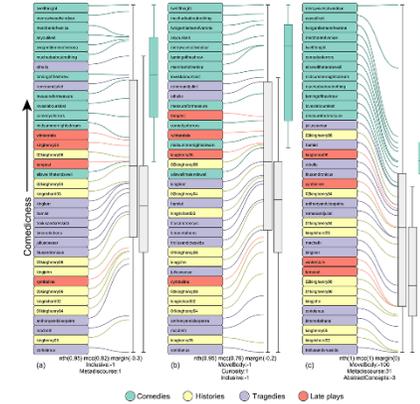
Text Viewer  
Correll et al, EuroVis 2011



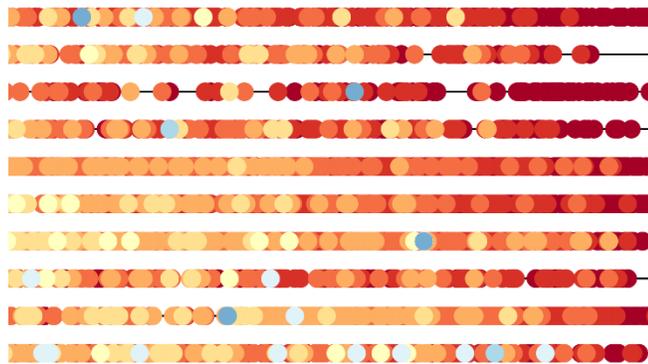
TextDNA  
Szafir et al, EuroVis 2016



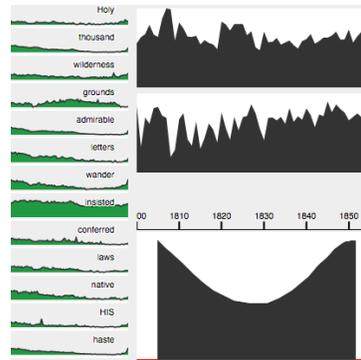
Serendip  
Alexander et al, VAST 2014



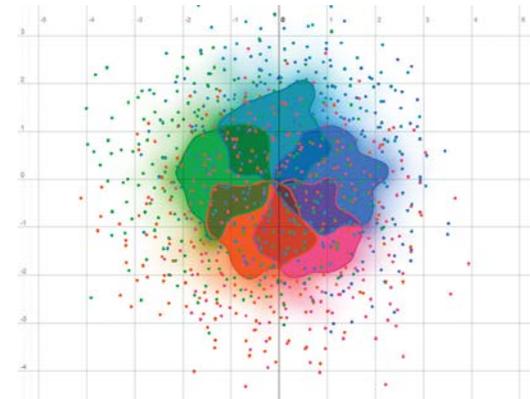
Explainers  
Gleicher, VAST 2013



Topic Model Comparison  
Alexander et al, VAST 2015



Sketch-based search  
Correll et al, VAST 2016



Splatterplots  
Mayorga & Gleicher 2013  
Sarikaya & Gleicher 2015

And others...

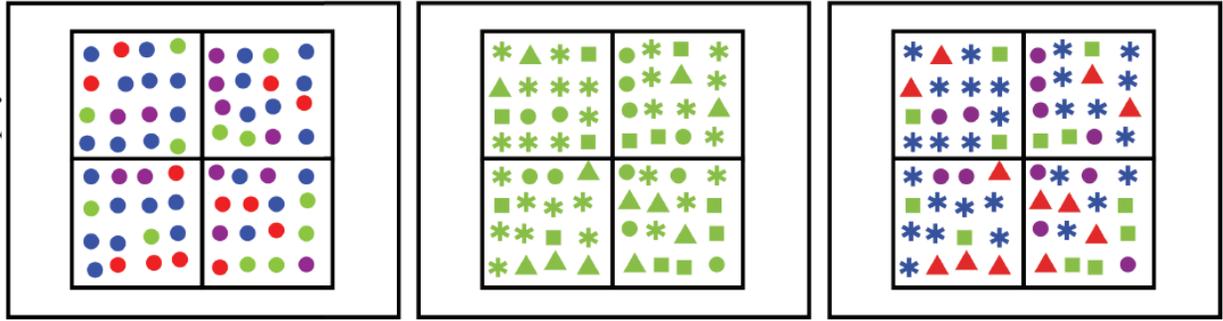
**What is visualization research?**

# Perceptual Studies

Target Preview



Test Display



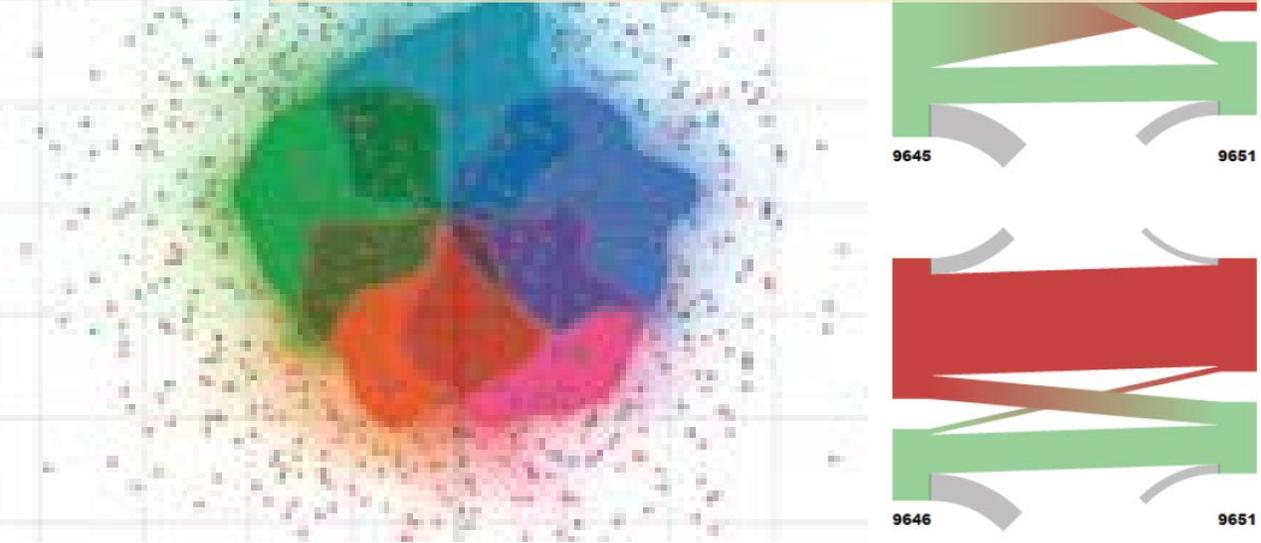
Serendip Word Rankings

Enter words separated by a space:  purple · Add

counsel x  
 punishment x  
 churches x  
 doctors x  
 council x

# Systems

# Technique Designs

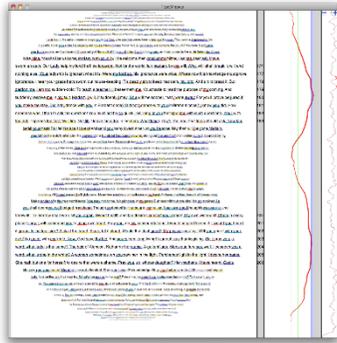


### Visual Aggregation Task

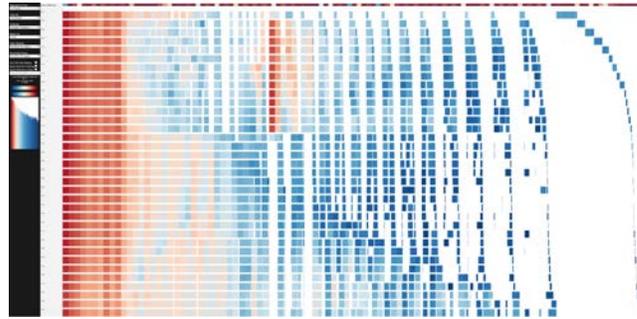
Visual Feature	Structure			
	Identification (Outlier)	Summary (Mean)	Segmentation (Clustering)	Estimation (Trends)
Position				
Size				
Orientation				

# Explorations of Principles

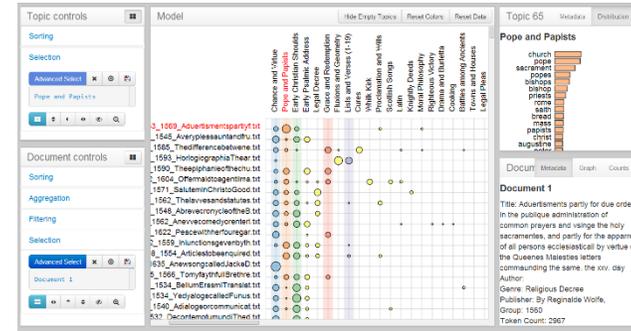
# VEP Vis success stories (some of each kind)



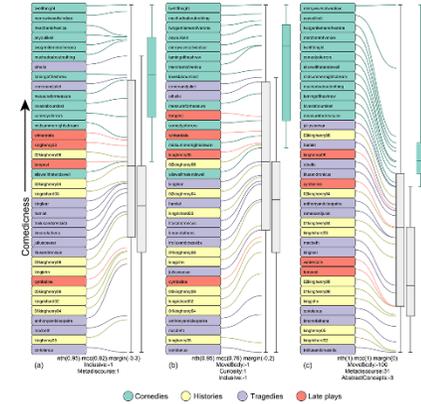
Text Viewer  
Correll et al, EuroVis 2011



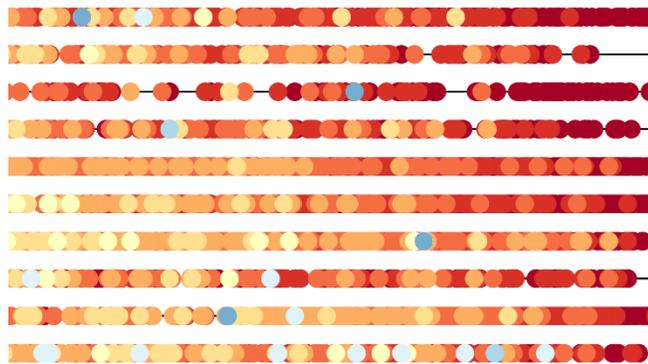
TextDNA  
Szafir et al, EuroVis 2016



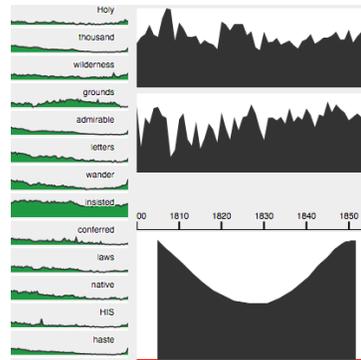
Serendip  
Alexander et al, VAST 2014



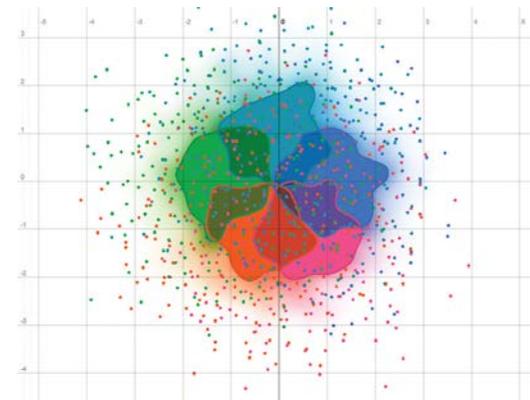
Explainers  
Gleicher, VAST 2013



Topic Model Comparison  
Alexander et al, VAST 2015



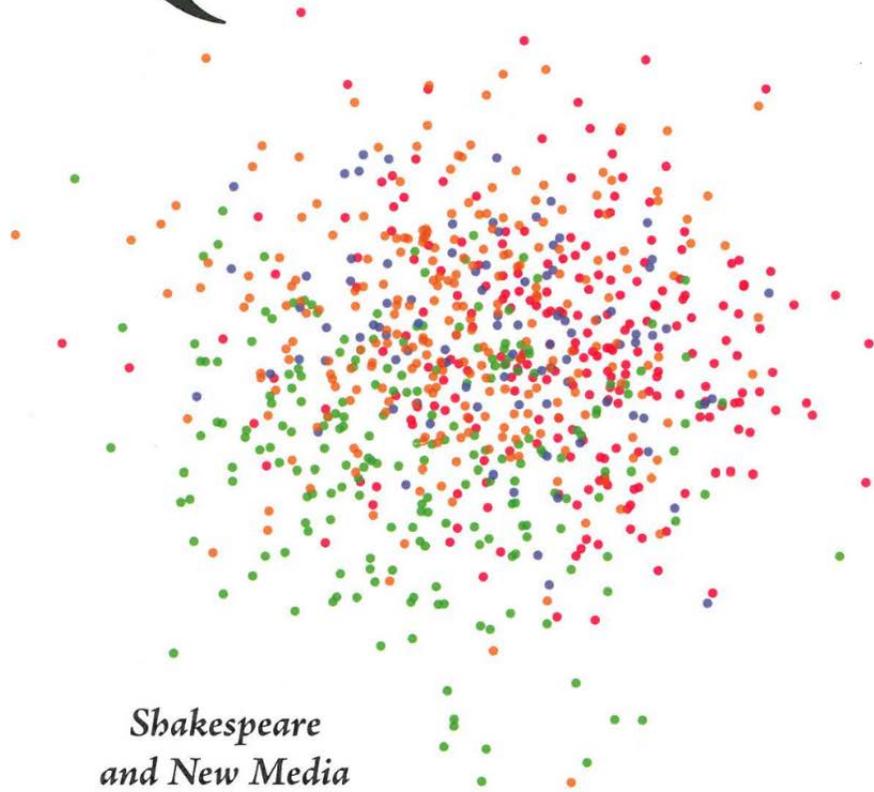
Sketch-based search  
Correll et al, VAST 2016



Splatterplots  
Mayorga & Gleicher 2013  
Sarikaya & Gleicher 2015

And others...

# SHAKESPEARE QUARTERLY

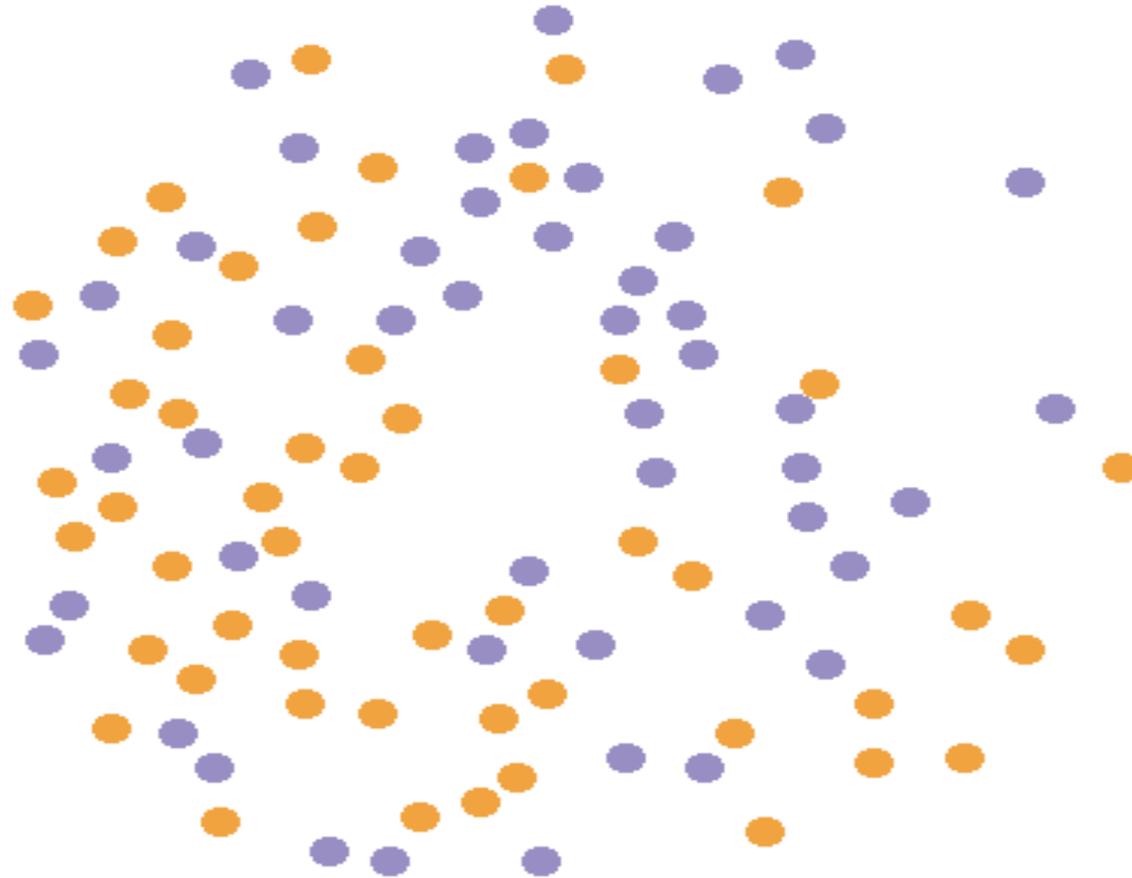


*Shakespeare  
and New Media*

Published for the Folger Shakespeare Library  
in association with  
The George Washington University  
by The Johns Hopkins University Press

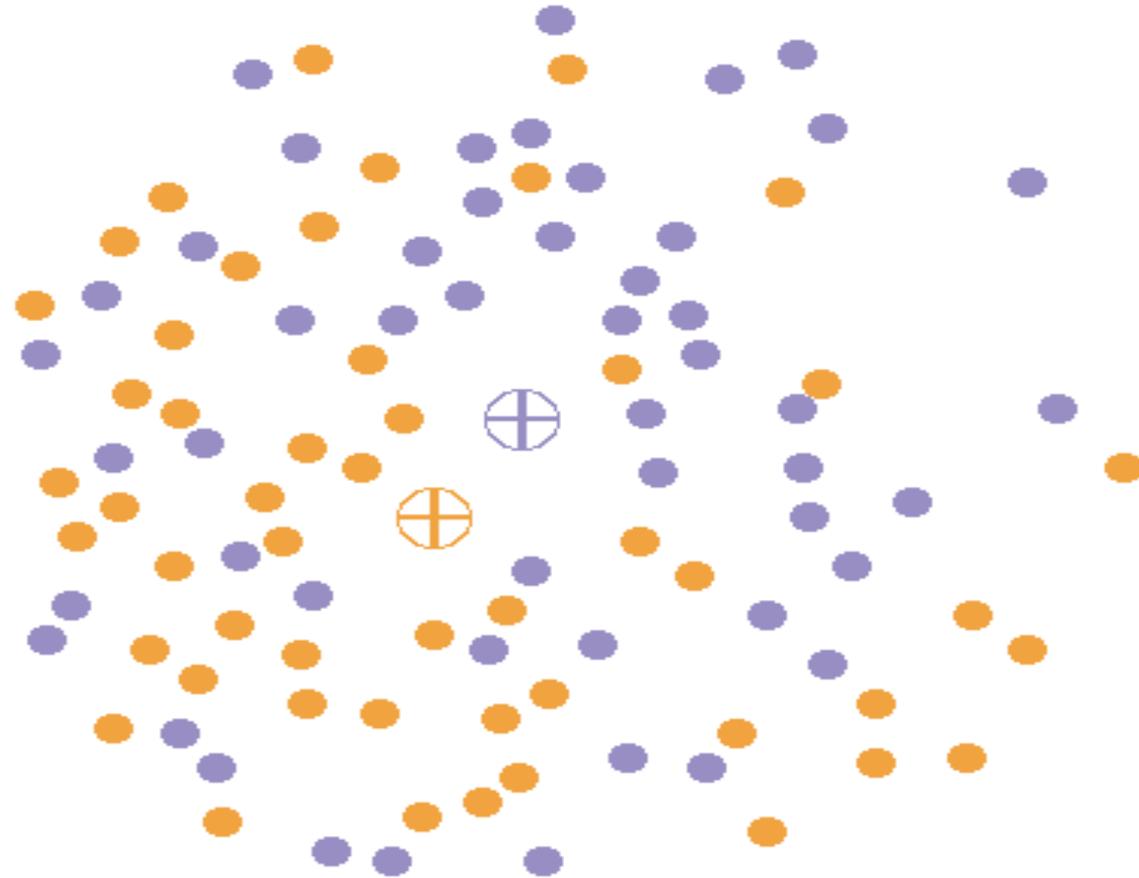
**What can you do with too many points?**

# Which Color Point is Higher on Average?



Gleicher, M., Correll, M., Nothelfer, C. and Franconeri, C. "Perception of Average Value in Multiclass Scatterplots." InfoVis 2013

# How did you do that?



# Visual Aggregation

People can extract summary statistics

Which Ones?

Efficiently?

Accurately?

How?

What can we do with it?

Why should we use it?

# Visual Aggregation

## Empirical Understanding

Averages in Time Series

Correll, et al. CHI 2012

Tagged Text

Correll, et al. CHI 2013

Scatterplot Averages

Gleicher, et al. InfoVis 2013

Other statistics in Time Series

Albers, et al. CHI 2014

## Practical Application

Sequence Surveyor (Genetics)

Albers, et al. InfoVis 2011

LayerCake (Virus mutations)

Correll, et al. BioVis 2011

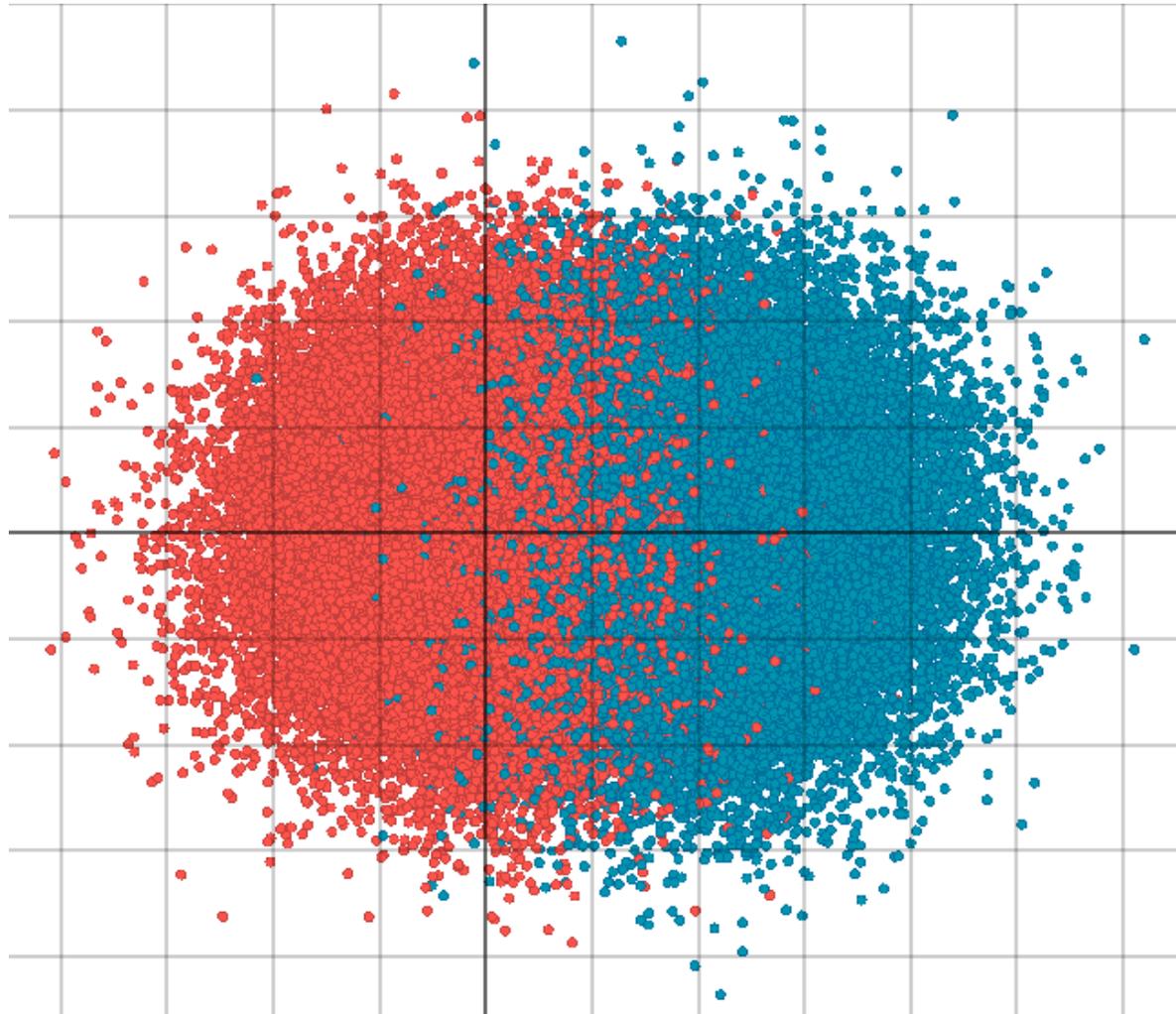
Molecular Surface Experiments

Sarikaya, et al. EuroVis 2014

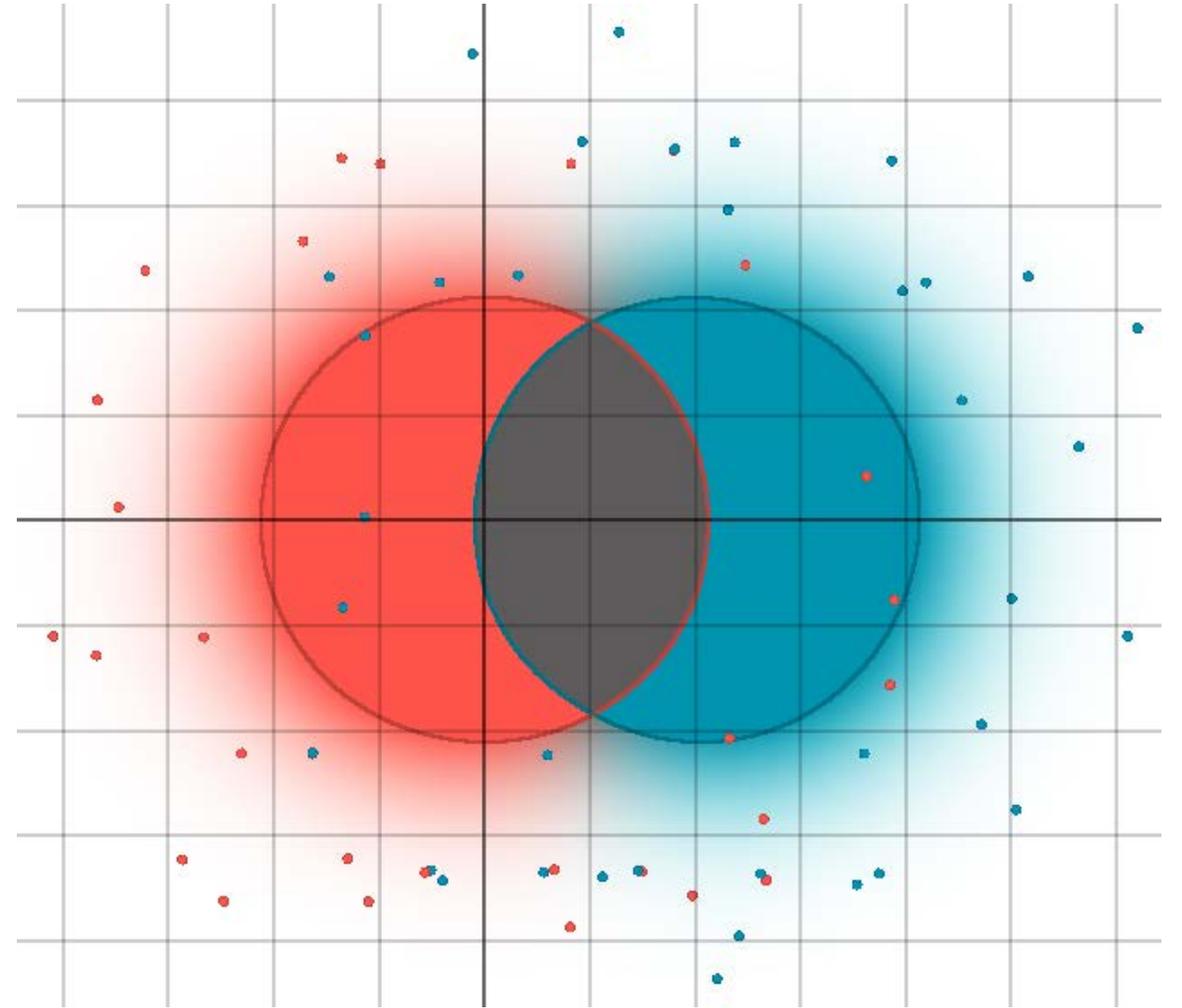
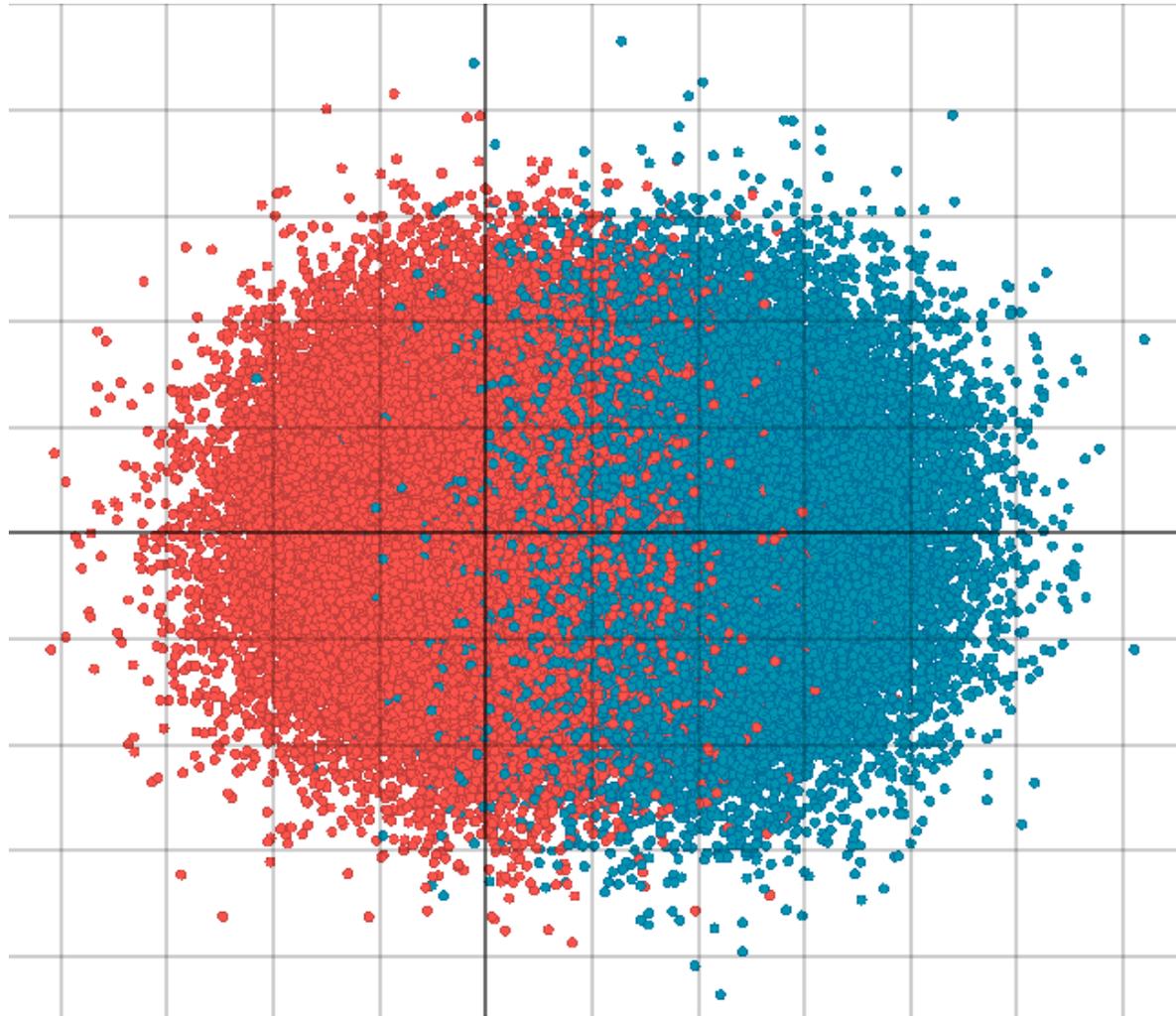
Decision Making

Correll, et al. (submitted)

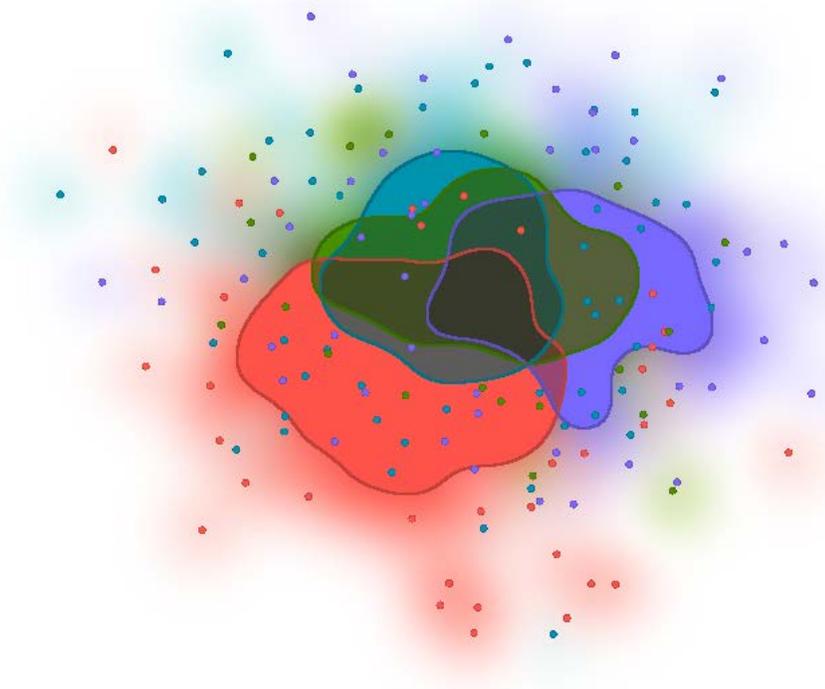
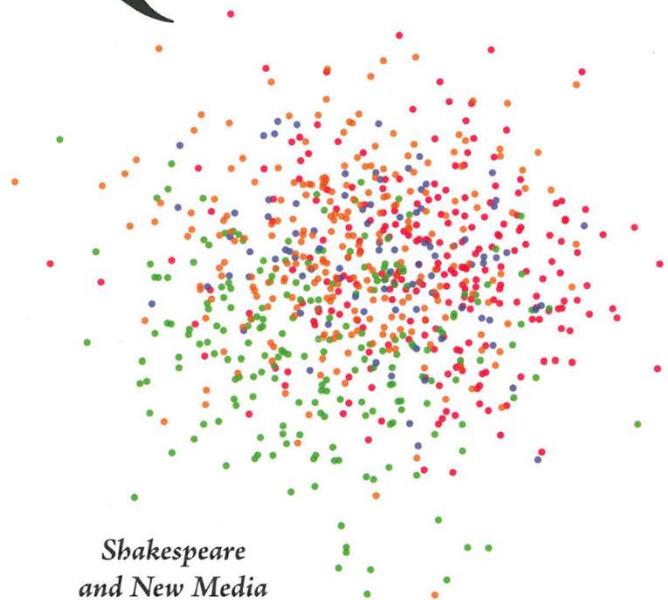
# Problem: Scatterplot with *way* too many points!



# Solution: Splatterplot



# SHAKESPEARE QUARTERLY



Published for the Folger Shakespeare Library  
in association with  
The George Washington University  
by The Johns Hopkins University Press

Volume 61

Fall 2010

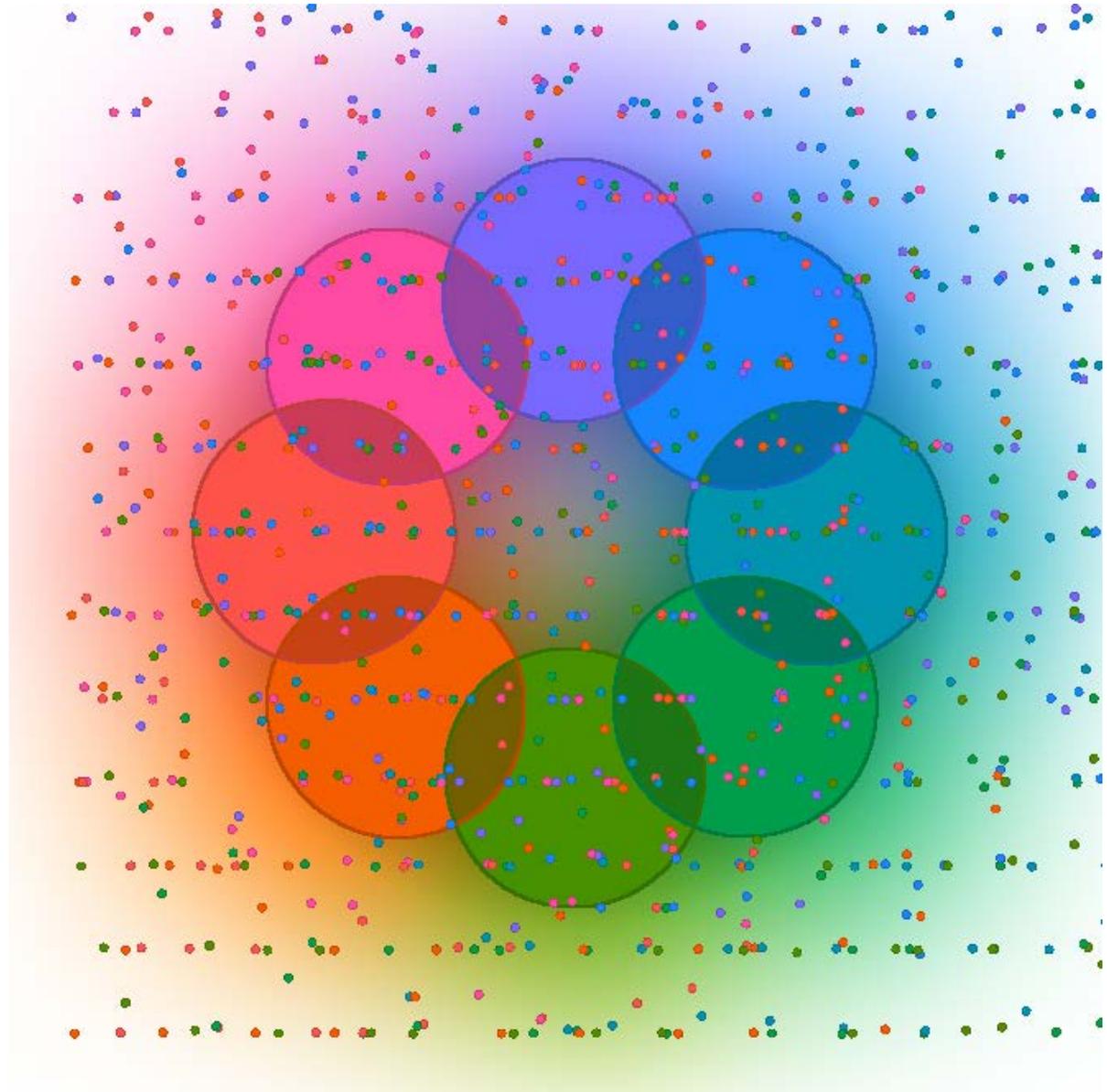
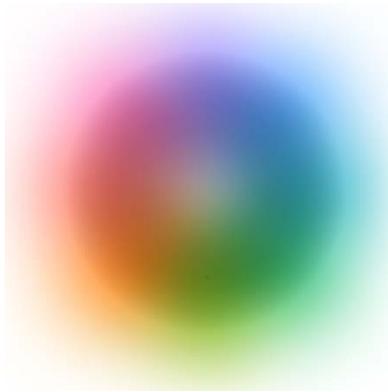
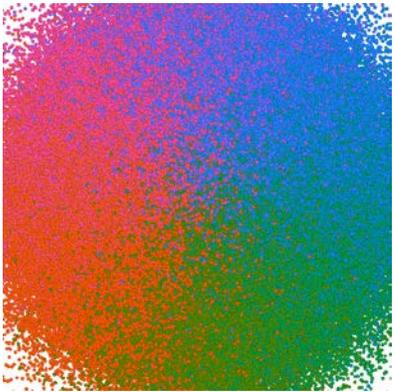
Number 3

This is pretty!

But when should we use it?

Lots of choices

Need actionable advice



# Scatterplots are common

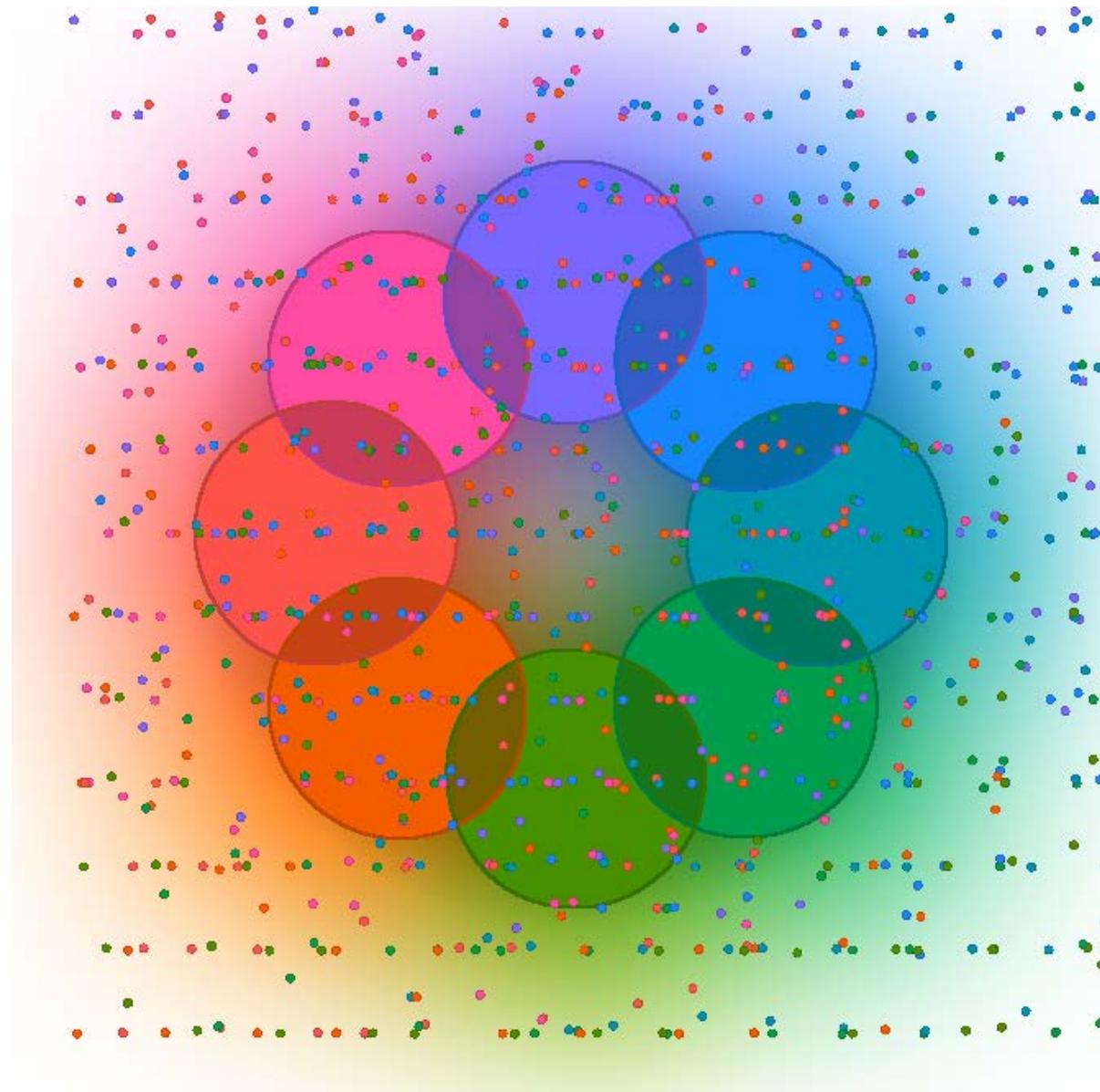
Lots of designs

How do we choose?

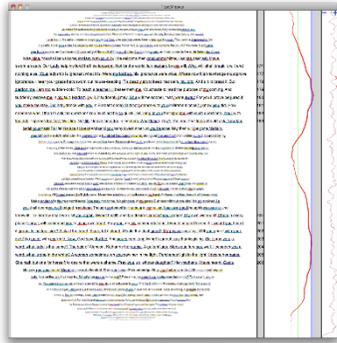
**Task oriented analysis**

Ask Dr. Scatterplot!

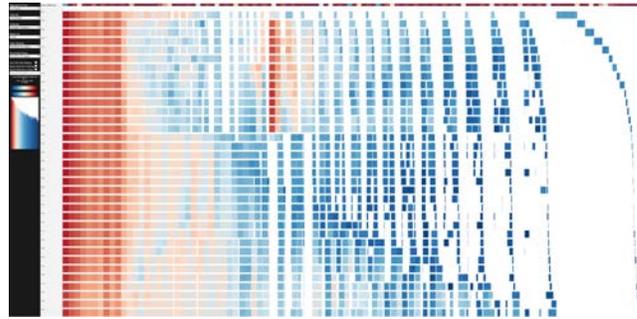
Sarikaya and Gleicher. Scatterplots: Tasks, Data, and Designs. IEEE Trans Vis and CG, 2018.



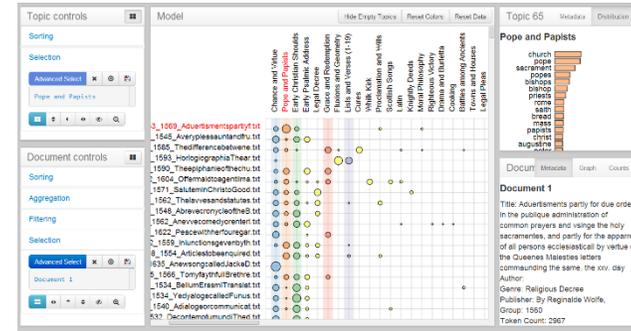
# VEP Vis success stories (some of each kind)



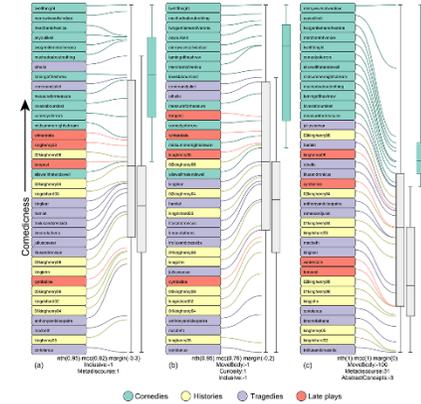
Text Viewer  
Correll et al, EuroVis 2011



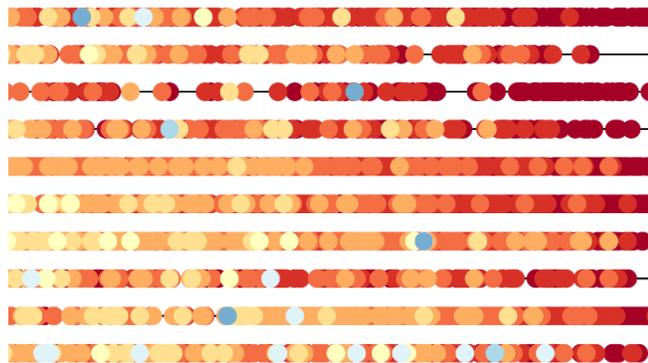
TextDNA  
Szafir et al, EuroVis 2016



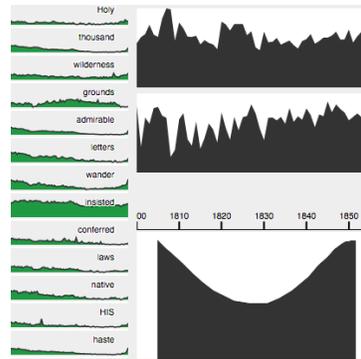
Serendip  
Alexander et al, VAST 2014



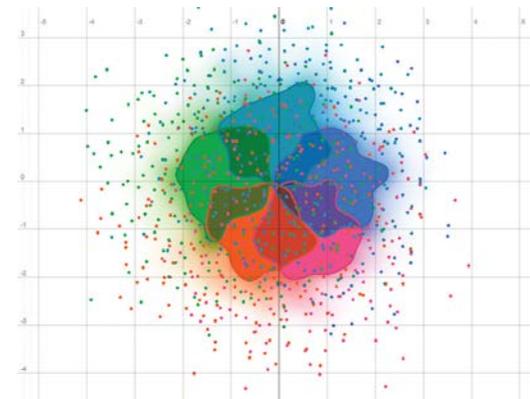
Explainers  
Gleicher, VAST 2013



Topic Model Comparison  
Alexander et al, VAST 2015



Sketch-based search  
Correll et al, VAST 2016



Splatterplots  
Mayorga & Gleicher 2013  
Sarikaya & Gleicher 2015

And others...

# But are there bigger lessons?

I learned a lot from this

# Three interesting aspects of VEP

Literary scholarship – and ways of thinking

Literary scholars – with their needs and abilities

300 years of old books – variety and quality

# 300+ years of books

A period of intense changes:  
cultural, historical, intellectual, political, ...

*How does the printed record reflect this?*

The language changed a lot

Printing changed a lot

The book business changed a lot

Libraries changed a lot

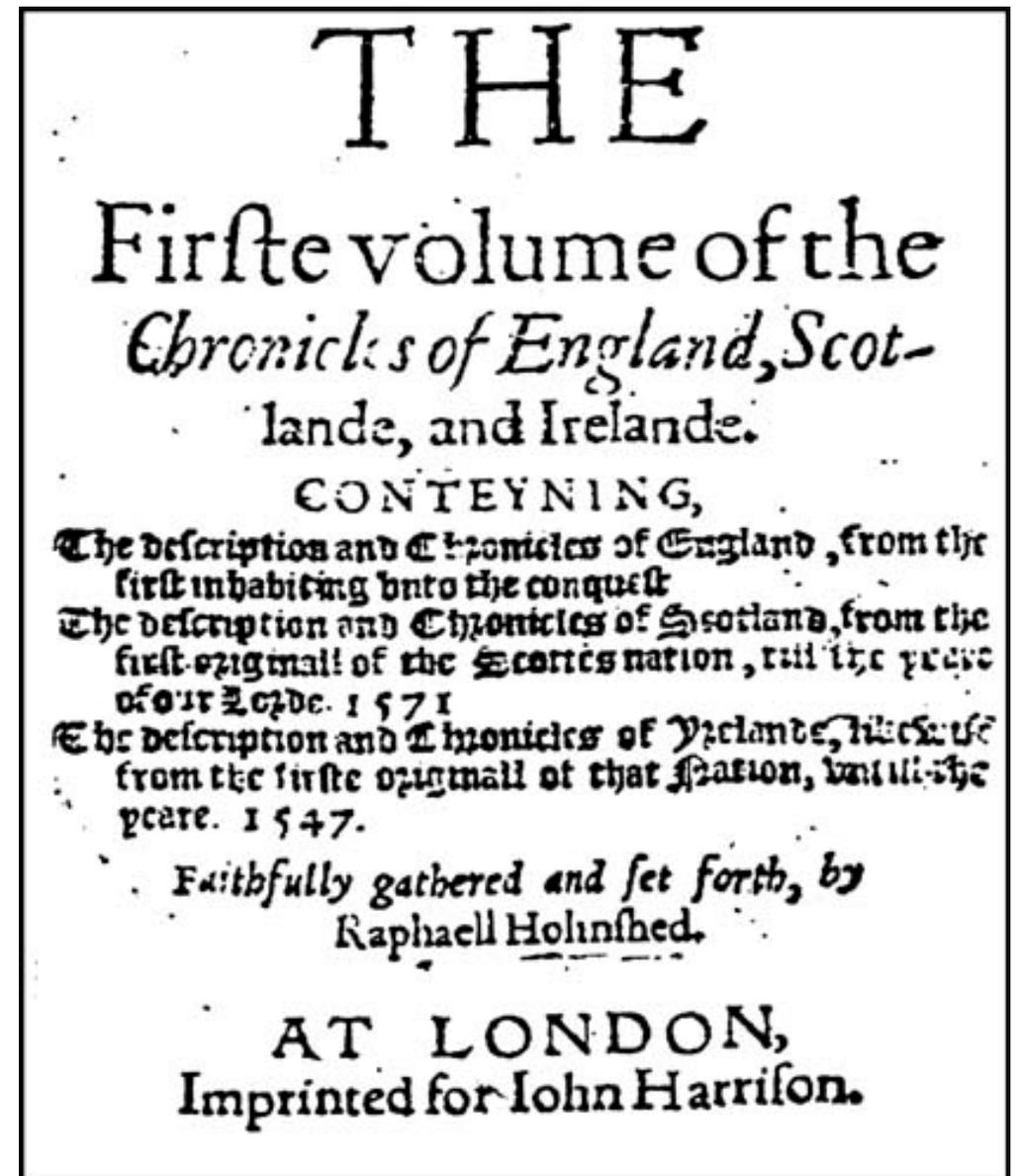
# 300+ years of books

Every known surviving book?

Huge historical efforts  
to catalog known books  
to microfilm known books  
to transcribe known books

Text Creation Partnership (TCP)

Hand-keyed about 60,000 books into SGML/XML







MACHINAMENTI NOVA COMPOSITIONE  
PER TELIS INSTRUMENTA ET  
TUBIS ANTI IMPULSIVIS ET  
DINEM AQUAM IMPULSIVIS ET  
EXTRAHI ET IMMO IN SUMMUM  
LITUM

# Sources of variation

Variety of topics

Authors' thoughts and ideas

Authors' words

Printers' practicality

Librarians' taking care of books

Transcribers translating books

Modern processing tools handling things correctly

# It's naïve to spell naïve with an i

Did the author use an i ?

Did the printer have it?

Did the dots get smudged?

Did the microfilming make it look like a dust mark?

Did the transcriber see it?

Did they encode it weirdly that year?

Did our Unicode pipeline mess it up?

Did our spelling standardizer get it wrong?

One form of variation can obscure another

But...

One person's noise is another's data

spelling obscures ideas

*or*

observe the development of spelling

rotting pages obscure content

*or*

preserved books tell us what was valued

# Data Wrangling

Make data convenient for analysis

Clean away “unwanted” variance

Leave enough signal

Hard choices!

# Data Wrangling Lessons

*Everyone knows Data Wrangling is a big deal, time consuming, ...*

Getting good meta-data requires investment

**Transparency** of data wrangling is valuable

**Comprehensibility** – the stakeholders need to understand

Curation and editing is part of scholarship.  
These are decisions – make them wisely!

# A view of Data Wrangling (or cleanup)

It's about variation:

variation is what we're interested in

in order to expose the variation you care about  
you need to clean away the variation you don't

Curation and editing are part of scholarship.  
These are decisions – make them wisely!

**A case study...**

# Counting “the”

## It is interesting at scale

Look for correlations across huge amounts of literature!

## It is transparent

we can explain – and check – every step

## It’s still hard to get right...

Just because counting is easy – doesn’t mean anything else is

# Some literature theory...

Writing about things in the world (Extra-Subjective)

vs.

Writing about abstract things (Intra-Subjective)

How does this play out over 300 years, 60,000 documents, ... ?

How does this play out over plays?

about 1500 documents- that we have lots of meta-data on

# Extra-Subjective vs. Intra-Subjective

Is there an easy way to measure it?

## Statistical analysis:

The primary variance in the collection is correlated with labels (**science**, **plays**)

The main source of this variance is the most common word

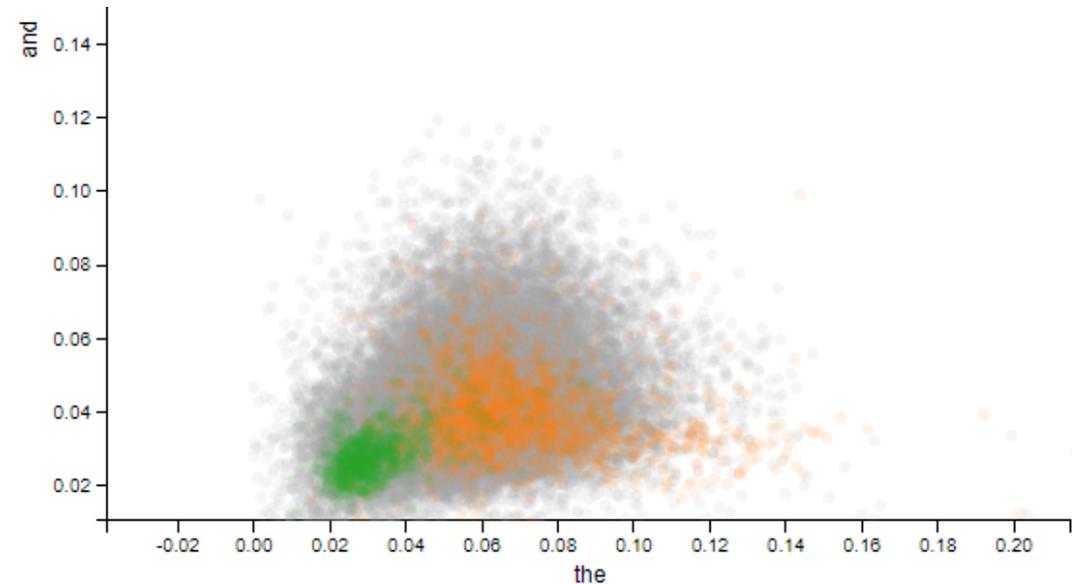
# Count “the” ?

Use definite articles to refer to objects in the world

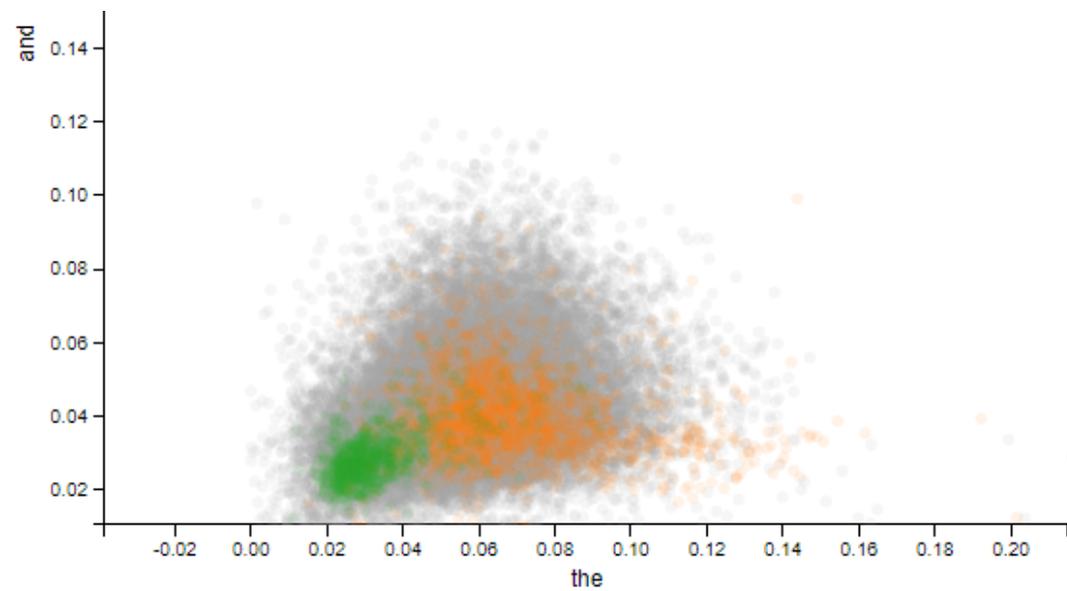
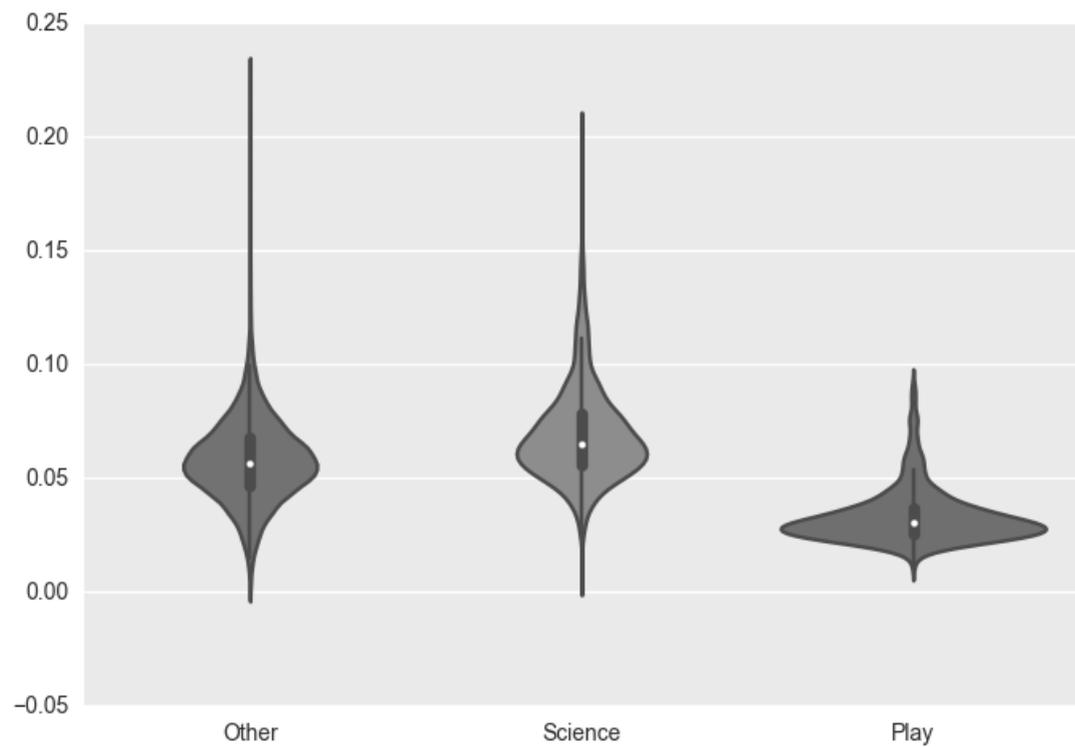
**Science** uses it more than **Plays**

We only have (extensive) labels for **science** and **plays**

Most things are unlabeled



# Count “the” ?



# 23% “the” – really?

A plain and easy rule to rigge any ship by the length of his masts, and yards, without any further trouble

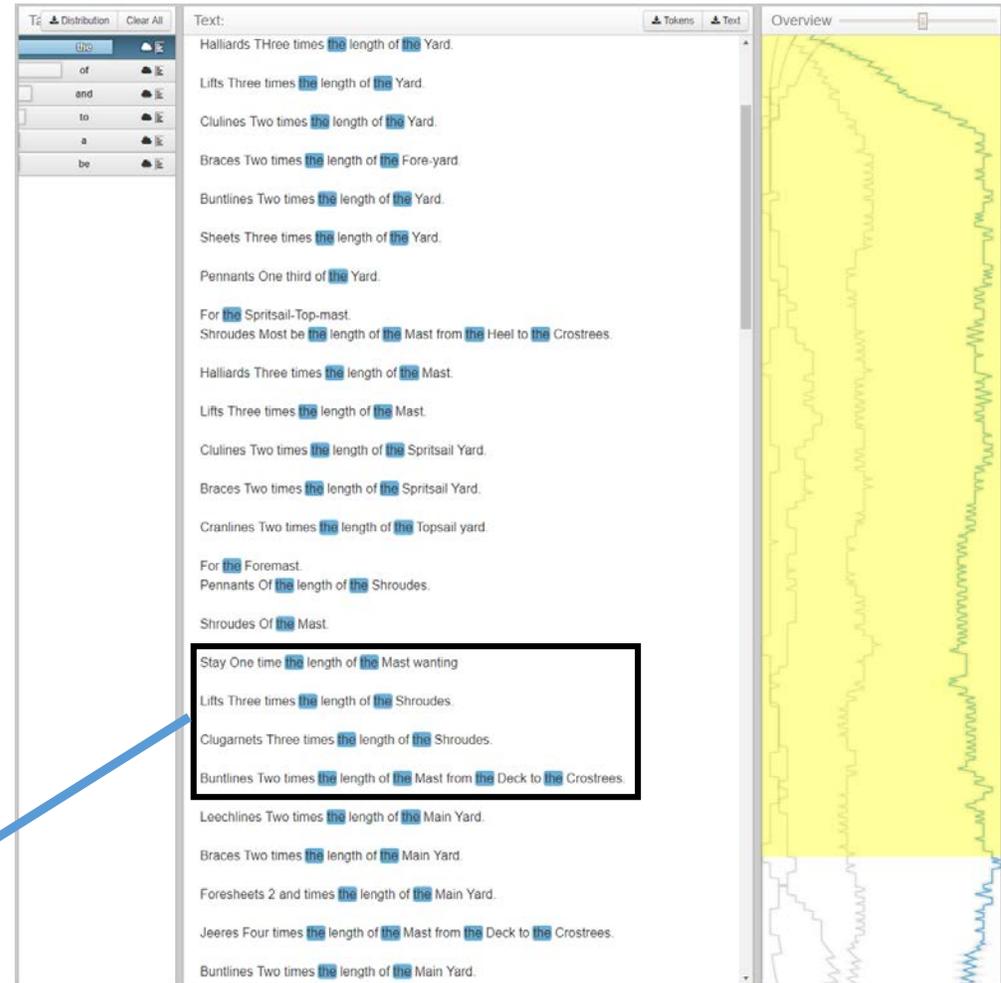
23% “the” ?

Stay One time **the** length of **the** Mast wanting

Lifts Three times **the** length of **the** Shroudes.

Clugarnets Three times **the** length of **the** Shroudes.

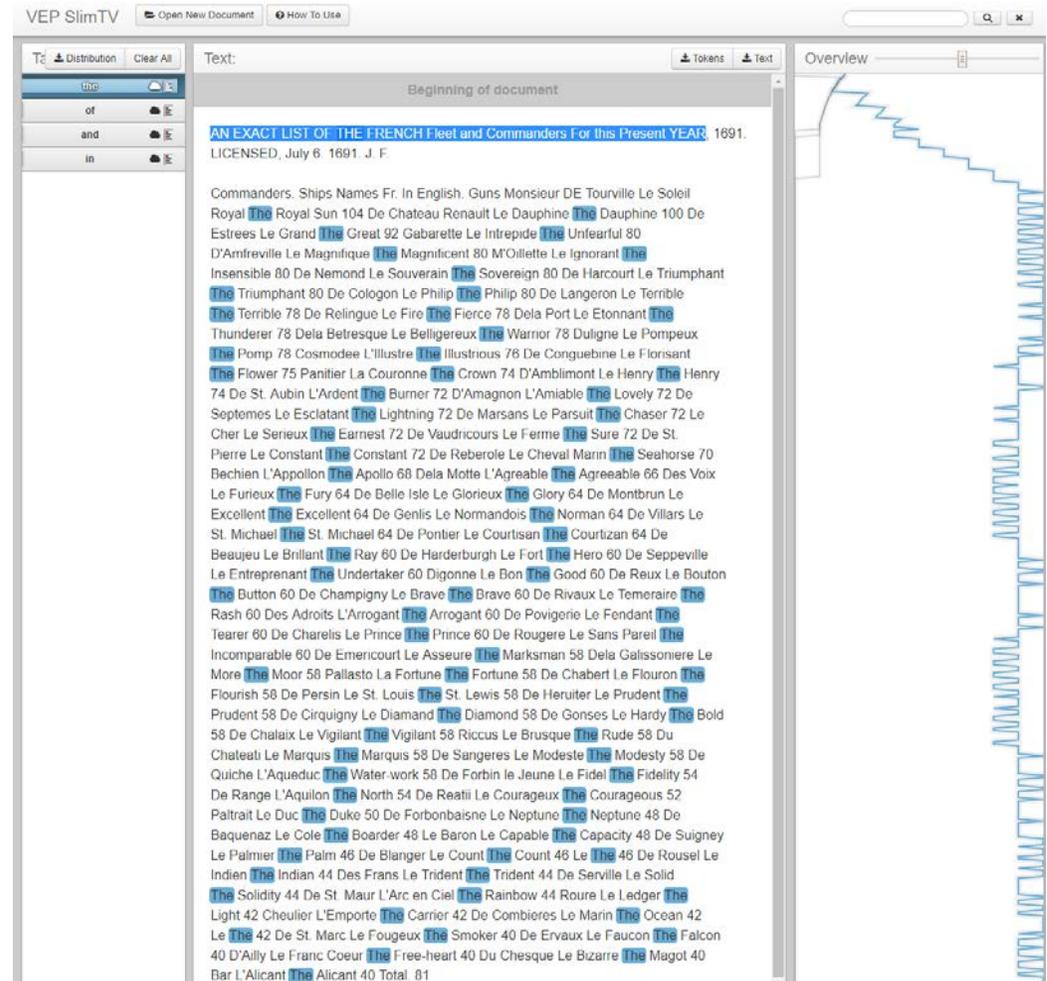
Buntlines Two times **the** length of **the** Mast from **the** Deck to **the** Crostrees.



# 51% “the” if you count differently

AN EXACT LIST OF THE  
FRENCH Fleet and  
Commanders For this Present  
YEAR

51% of “words” – only counting  
things we recognize as words



The screenshot shows the VEP SlimTV interface. The main window displays a text document titled "Beginning of document" with the following text: "AN EXACT LIST OF THE FRENCH Fleet and Commanders For this Present YEAR, 1891. LICENSED, July 6 1891. J. F. Commanders. Ships Names Fr. In English. Guns Monsieur DE Tourville Le Soleil Royal The Royal Sun 104 De Chateau Renault Le Dauphine The Dauphine 100 De Estrees Le Grand The Great 92 Gabarette Le Intrepide The Unfearful 80 D'Amfreville Le Magnifique The Magnificent 80 M'Oilette Le Ignorant The Insensible 80 De Nemond Le Souverain The Sovereign 80 De Harcourt Le Triumphant The Triumphant 80 De Cologon Le Philip The Philip 80 De Langeron Le Terrible The Terrible 78 De Relingue Le Fire The Fierce 78 Dela Port Le Etonnant The Thunderer 78 Dela Betresque Le Belligereux The Warrior 78 Duligne Le Pompeux The Pomp 78 Cosmodee L'illustre The illustrious 76 De Conguebine Le Florissant The Flower 75 Panitier La Couronne The Crown 74 D'Amblimont Le Henry The Henry 74 De St. Aubin L'Ardent The Burner 72 D'Amagnon L'Amiable The Lovely 72 De Septomes Le Esclatant The Lightning 72 De Marsans Le Parsuit The Chaser 72 Le Cher Le Serieux The Earnest 72 De Vaudricours Le Ferme The Sure 72 De St. Pierre Le Constant The Constant 72 De Reberole Le Cheval Marin The Seahorse 70 Bechien L'Appolon The Apollo 68 Dela Motte L'Agreeable The Agreeable 66 Des Voix Le Furieux The Fury 64 De Belle Isle Le Glorieux The Glory 64 De Montbrun Le Excellent The Excellent 64 De Genlis Le Normandois The Norman 64 De Villars Le St. Michael The St. Michael 64 De Pontier Le Courtisan The Courtizan 64 De Beaujeu Le Brillant The Ray 60 De Harderburgh Le Fort The Hero 60 De Seppeville Le Entreprenant The Undertaker 60 Dignon Le Bon The Good 60 De Reux Le Bouton The Button 60 De Champigny Le Brave The Brave 60 De Rivaux Le Temeraire The Rash 60 Des Adroits L'Arrogant The Arrogant 60 De Povigerie Le Fendant The Tearer 60 De Charelis Le Prince The Prince 60 De Rougere Le Sans Pareil The Incomparable 60 De Emercourc Le Assuree The Marksman 58 Dela Galissoniere Le More The Moor 58 Pallasto La Fortune The Fortune 58 De Chabert Le Flouron The Flourish 58 De Persin Le St. Louis The St. Lewis 58 De Heruiter Le Prudent The Prudent 58 De Cirquigny Le Diamand The Diamond 58 De Gonses Le Hardy The Bold 58 De Chalaix Le Vigilant The Vigilant 58 Riccus Le Brusque The Rude 58 Du Chateati Le Marquis The Marquis 58 De Sangeres Le Modeste The Modesty 58 De Quiche L'Aqueduc The Water-work 58 De Forbin le Jeune Le Fidel The Fidelity 54 De Range L'Aquilon The North 54 De Raatii Le Courageux The Courageous 52 Paltrait Le Duc The Duke 50 De Forbonbaisne Le Neptune The Neptune 48 De Baquenaz Le Cole The Boarder 48 Le Baron Le Capable The Capacity 48 De Suigney Le Palmier The Palm 46 De Blanger Le Count The Count 46 Le The 46 De Rousel Le Indien The Indian 44 Des Frans Le Trident The Trident 44 De Serville Le Solid The Solidity 44 De St. Maur L'Arc en Ciel The Rainbow 44 Roure Le Ledger The Light 42 Cheulior L'Emporte The Carrier 42 De Combières Le Marin The Ocean 42 Le The 42 De St. Marc Le Fougeux The Smoker 40 De Ervaux Le Faucon The Falcon 40 D'Ailly Le Franc Coeur The Free-heart 40 Du Chesque Le Bizarre The Magot 40 Bar L'Alicant The Alicant 40 Total. 81

# Why is this hard?

We need to carefully treat each book in a uniform way

*Which words do we extract off the pages?*

We need to carefully treat each word in a uniform way

*How to deal with bad characters and spelling variation?*

We need to analyze in a uniform and transparent way

We need good meta-data

The ***scholars*** (and their audience) must understand each step!

We need to get **back to the sources!**

Big lesson:

**Data Science is a process**  
**Challenges are everywhere!**

Big lesson:

**We must help all users understand  
all phases of the process**

Big lesson:

**We must help all users understand  
all phases of the process**

# Towards Comprehensibility in Modeling

The other title for this talk

# Data Science is a Process

There are lots of steps

Many provide challenges

Many provide opportunities

Identify Problem

Gather Data

Clean Data

Abstract

Design

Build

Evaluate/Validate

Interpret

Act

Disseminate

An Aside...

**Can we automate this process?**

**Make it possible for subject experts?**

# Data Driven Discovery of Models

Smart systems to do steps

Smart user experience to guide users

Some steps seem really hard

Identify Problem

Gather Data

Clean Data

Abstract

Design

Build

Evaluate/Validate

Interpret

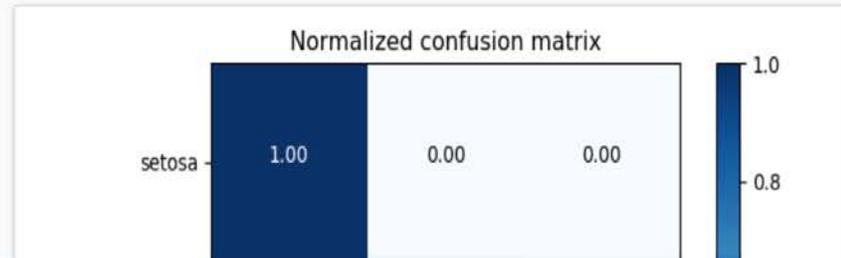
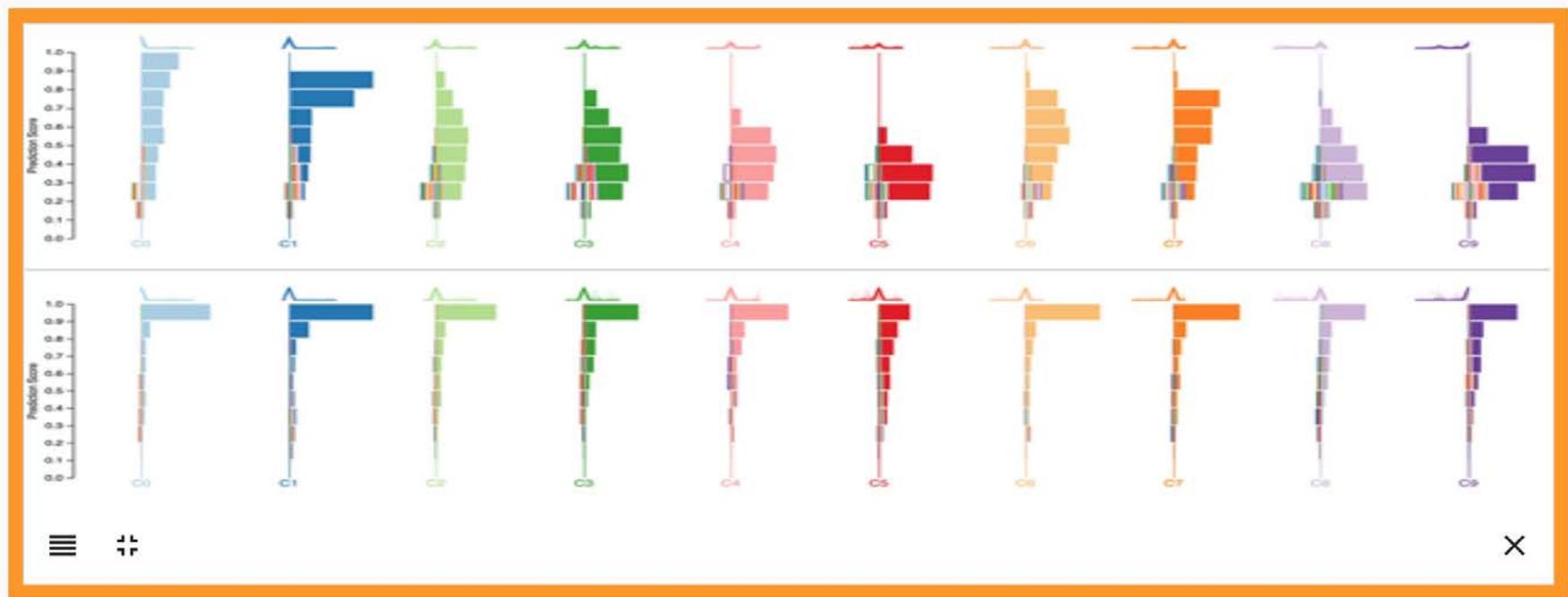
Act

Disseminate

- ✓ Approve Problems
  - ✓ Wait for Results
  - ✓ Select Metric  
Selected: specificity
  - ✓ Select Model  
Default: Model 1
  - ✓ Examine model results in detail
  - ✓ Compare models
- Model 1
- SELECT MODEL**
- 5 Create Executable
- Wanna skip to the end? **SHIP IT!**

### Problem Definition

This will be a description of the problem for the user to review. This problem is very interesting. It is perhaps the most interesting of all the problems. Can you think of a problem more interesting than it?



- ✓ Approve Problems
    - ✓ Edit/Examine the problem definition
    - ✓ Sign off on problem
    - ✓ Provision backends
  - START BACKENDS
  - 2 Wait for Results
  - 3 Select Metric  
Default: accuracy
  - 4 Select Model  
Default: Model 1
  - 5 Create Executable
- Wanna skip to the end? SHIP IT!

### Problem Definition

This will be a description of the problem for the user to review. This problem is very interesting. It is perhaps the most interesting of all the problems. Can you think of a problem more interesting than it?

☰ ☰

- ✓ Backend
- ✓ Backend

☰ ☰

**Workflow View**

- Show steps in process
- Show progress and next steps
- Connect to Visualizations
- Allow for skipping (guess for me)
- Allow for review (go back)

# Comprehensibility across the Process

## Who?

Stakeholders

Developers

Data Scientists

Domain Experts

Audience

Subjects

## Why?

Reason for Wanting

Improve Performance

Build Theory

Extend/Characterize

Build Trust

Actionability

Convince

## Where?

Phase of the **Process**

Identify Problem

Gather Data

Clean Data

Abstract

Design

Build

Evaluate

Interpret

Disseminate

# Who? Why? Where? (for How? and What? See the paper)

## Who?

Stakeholders

Developers

Data Scientists

Domain Experts

Audience

Subjects

## Why?

Reason for Wanting

Improve Performance

Build Theory

Extend/Characterize

Build Trust

Actionability

Convince

## Where?

Phase of the **Process**

Identify Problem

Gather Data

Clean Data

Abstract

Design

Build

Evaluate

Interpret

Disseminate

# Some phases for comprehensibility

Data Wrangling

understand what data the data is and what it can do

Model Building

understand classifiers to build theory / identify items

Validation Experiments

understand experiments to identify items / build trust

# Towards Comprehensibility in Modeling

This is not just comprehensible models

# genre

Categorization given by Shakespeare's contemporaries

**Comedy**

**Tragedy**

**History**

Category for plays written after that

**Late Plays**

# data

Count words

Count kinds of words – using dictionaries

DocuScope (Igarashi and Kaufer)

Lists of phrases for each “Language Action Type”

Simple and Transparent dictionary matching

(once the dictionaries are made)

Fancier methods (Topic Models, Machine Learning, ...)

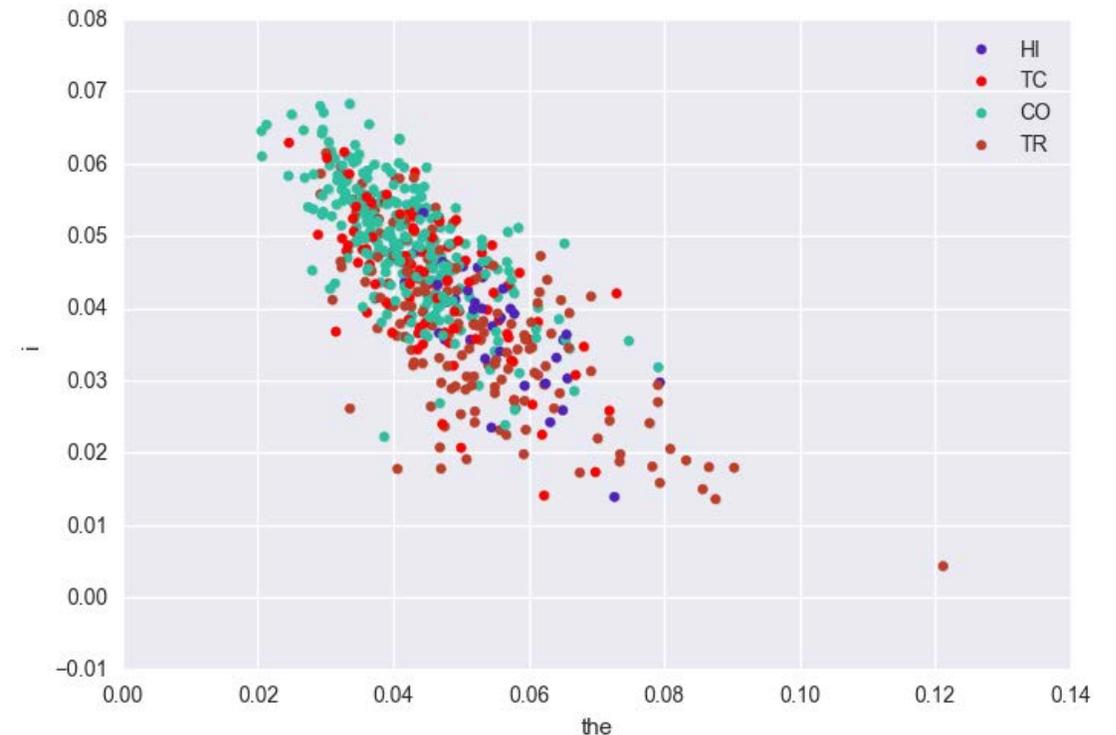
# Is “the” good enough?

**Histories** and **Tragedies** are more external

**Comedies** are more internal

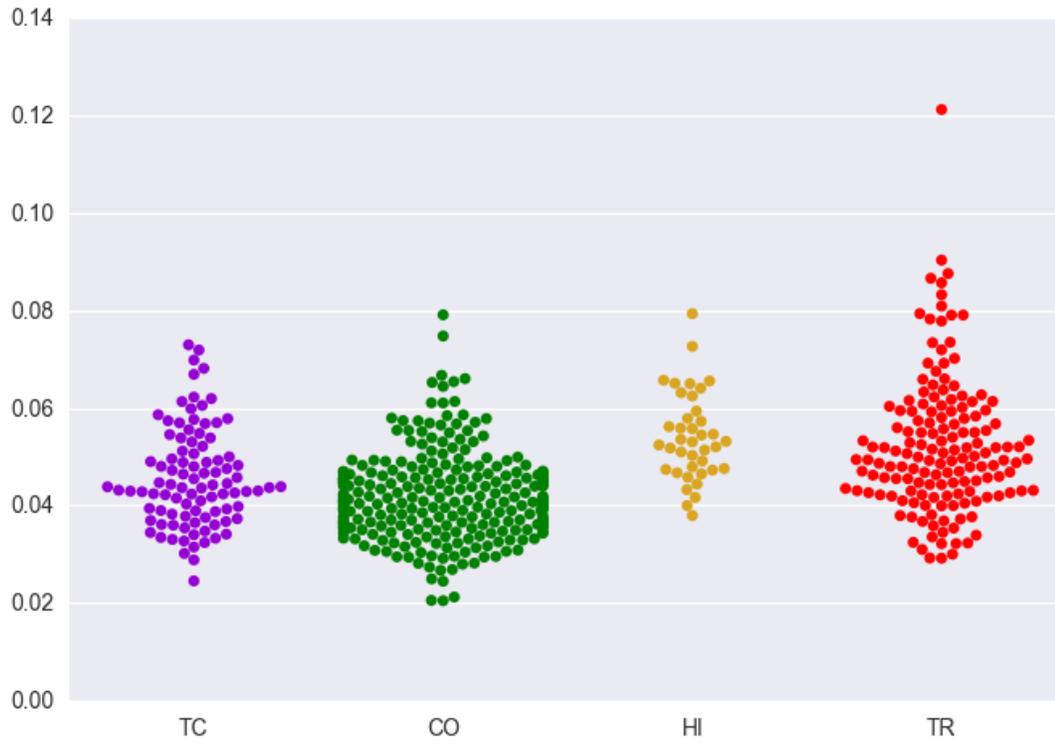
The second most common word, “I”, might be better (it is for internal)

600 plays of Shakespeare’s era

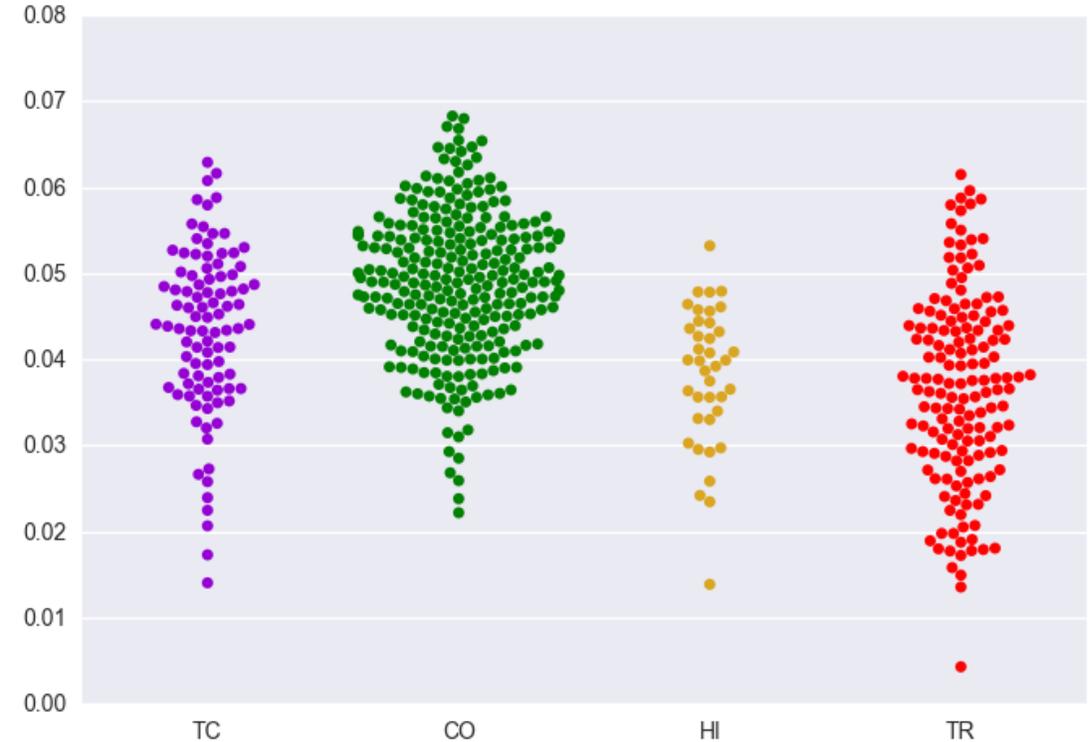


# Neither is great, maybe try combining them?

The



I



# Mathematical Models

“the” – “I” --- pretty good

(“the” + “to”) – (“I” + “a”) --- very good

How do we come up with such equations?

What do we do with them?

**Machine Learning** Classification:

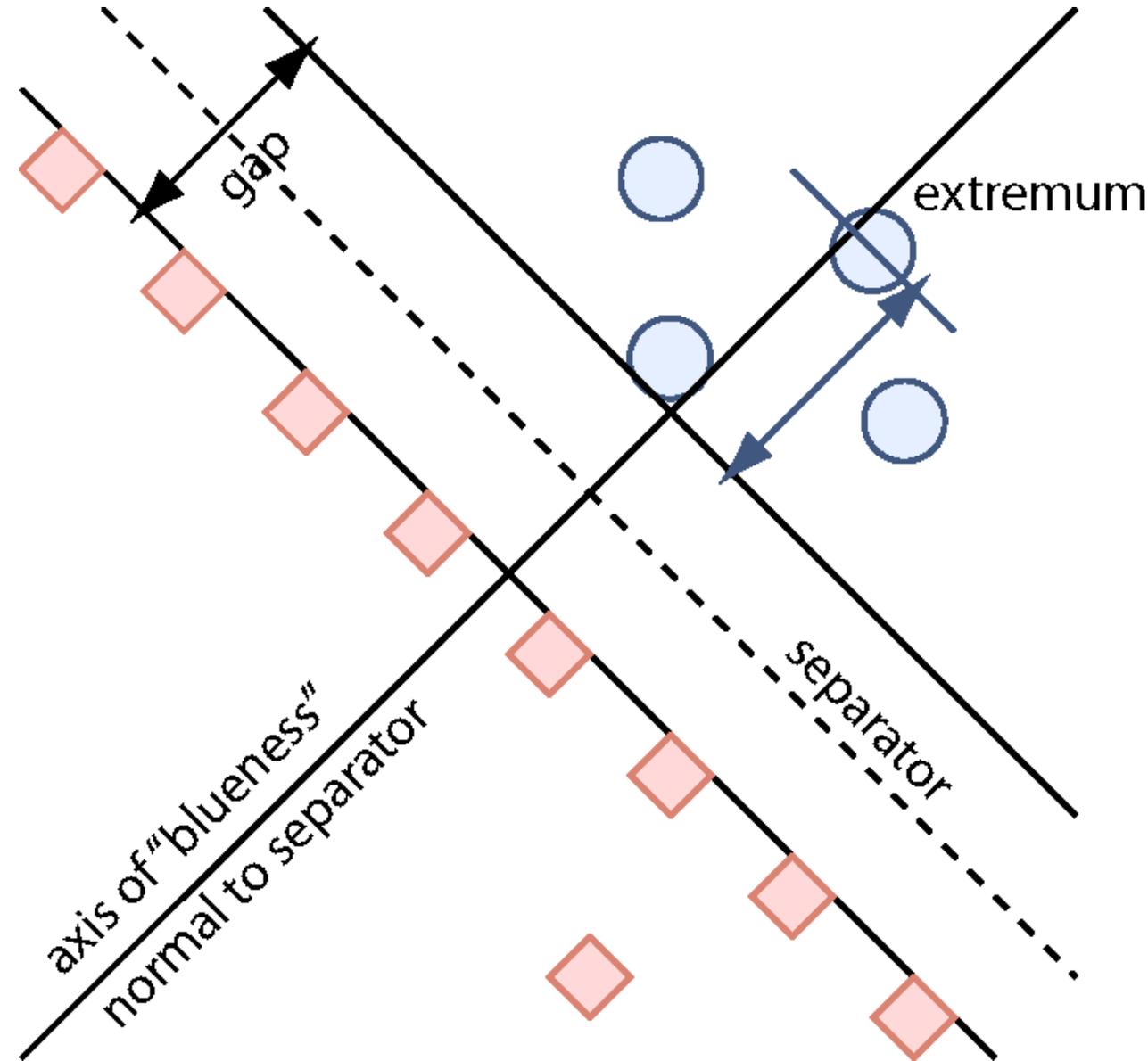
Make up these function so that one “class” is high

The other classes are low

# Comprehensibility in modeling?

They may never understand  
Support Vector Machines  
(or choose your favorite buzzword)

But they might understand  
what it does...



# Use Models for Insight

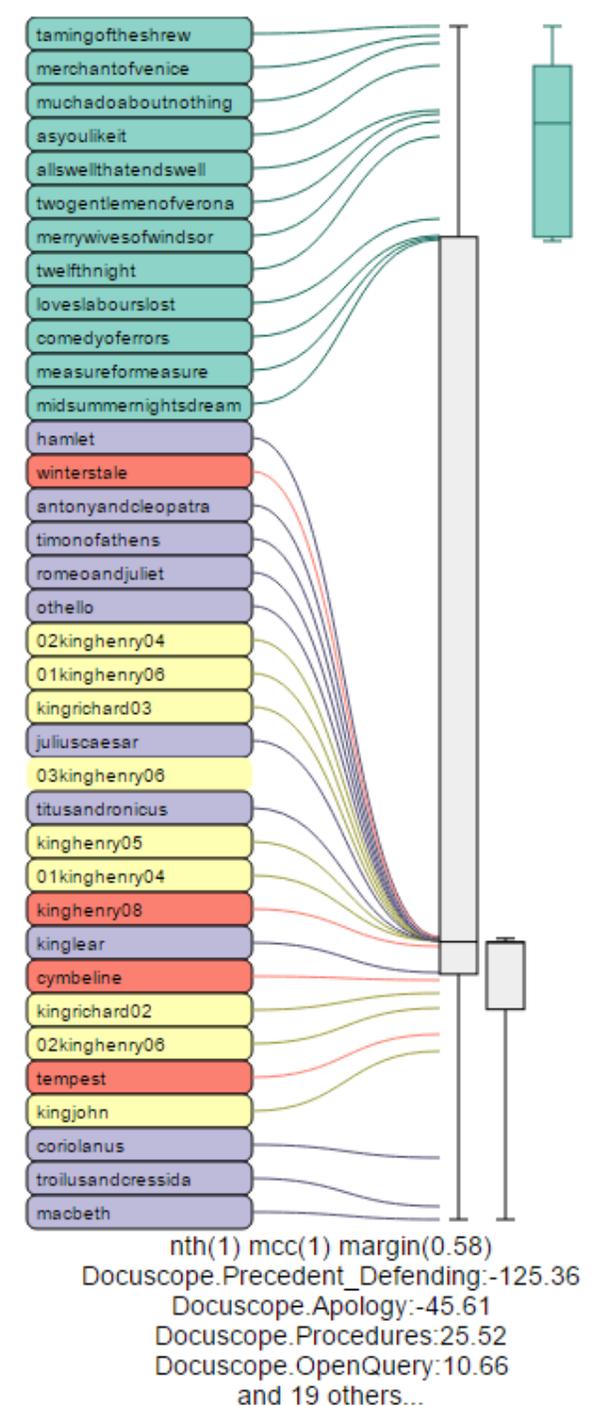
We know the answers!  
Teaching the machine can help us:

Can you pick apart known groups?

Does the data capture the concepts?

Language features vs. genre

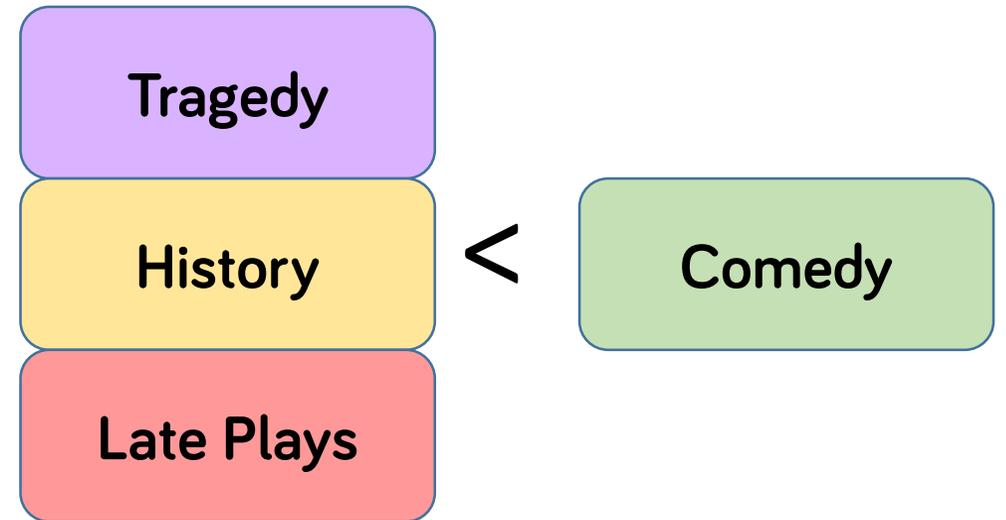
Can we quantify and organize? (comedic-ness)



# Comedicness

A measure of how much of a comedy something is

It's the "stuff" comedies have more of  
Where "stuff" has to be in the data



## Organization:

What is most/least comedic?

## Explanation:

How is the word usage (measured stuff) different in comedies?

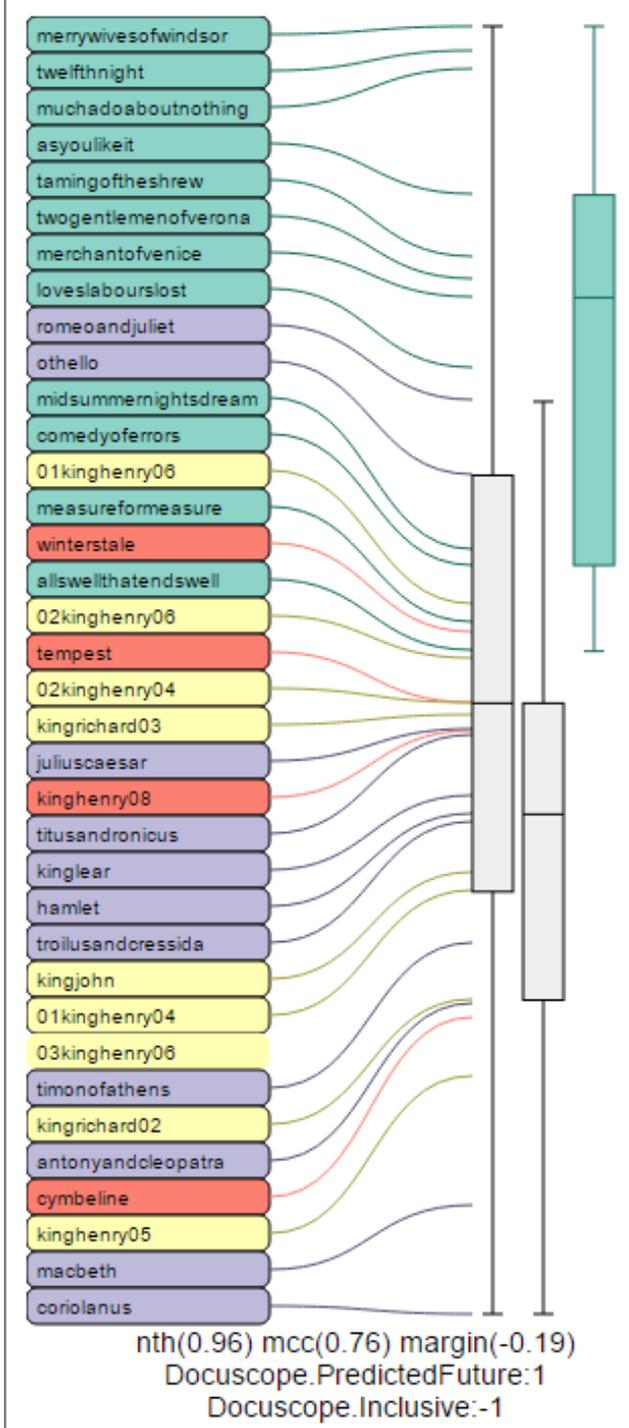
# Use **Simple** Models for Insight

Create models for comprehensibility

See what you get wrong

Look “inside” models to see how they work

Gleicher. *Explainers: Expert Explorations with Crafted Projections*. IEEE TVCG 19(12), Dec 2013. Proceedings VAST 2013.



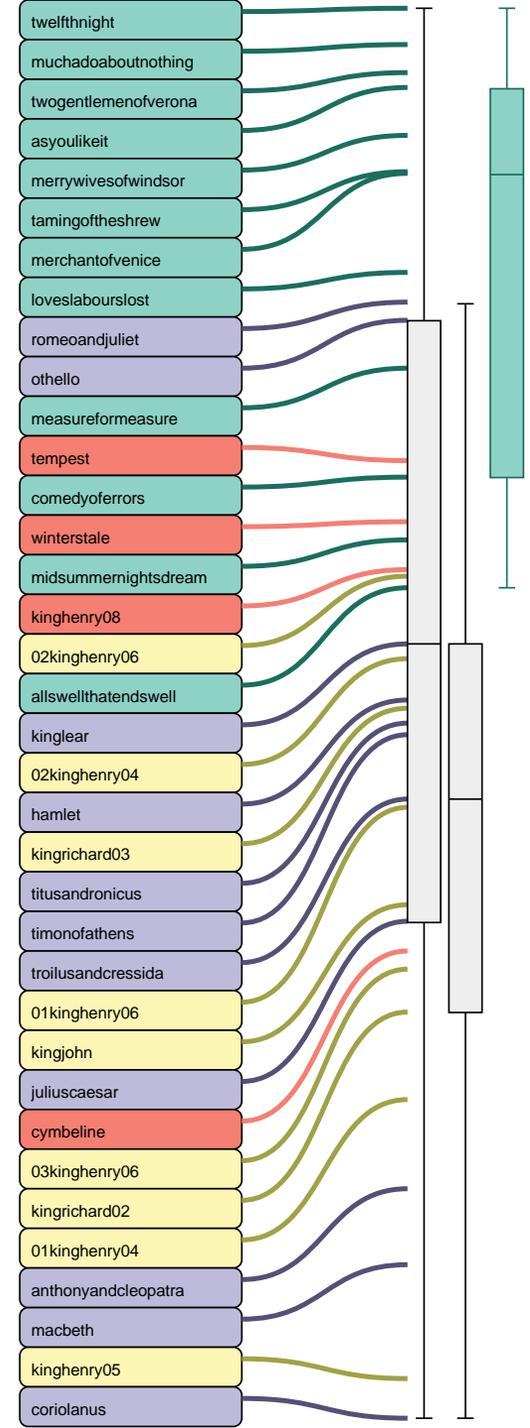
comedicness = M - I

$$f(V) = V[39] - V[42]$$

M = Predicted Future

I = Inclusiveness

36 plays of Shakespeare, Docuscope Features



M - I (5 wrong)

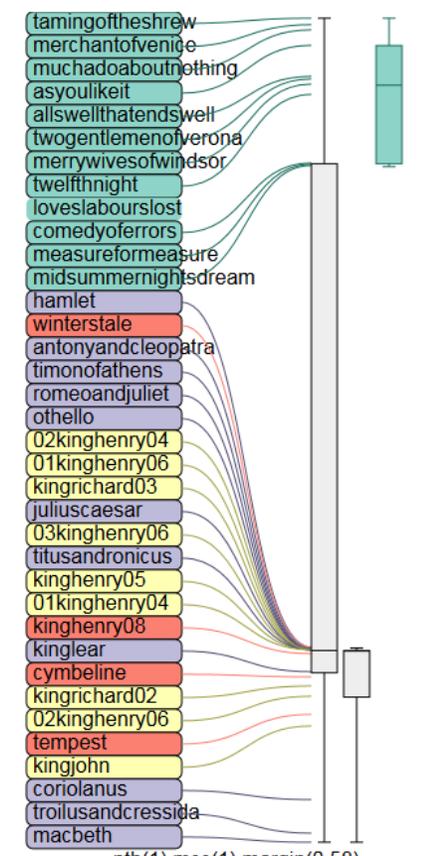
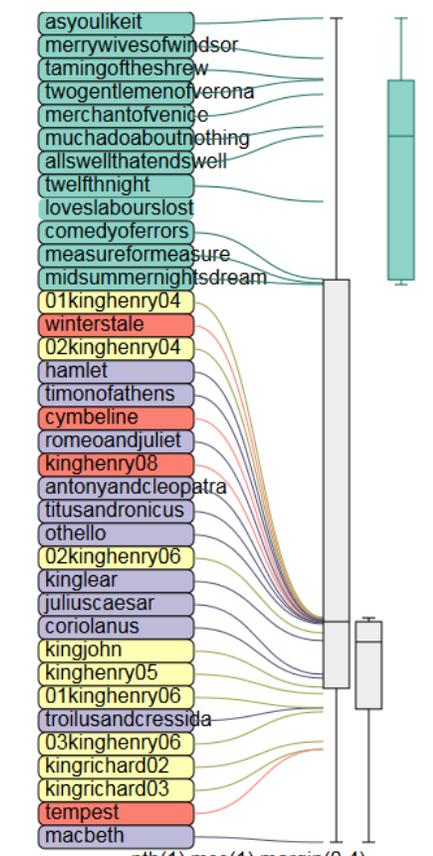
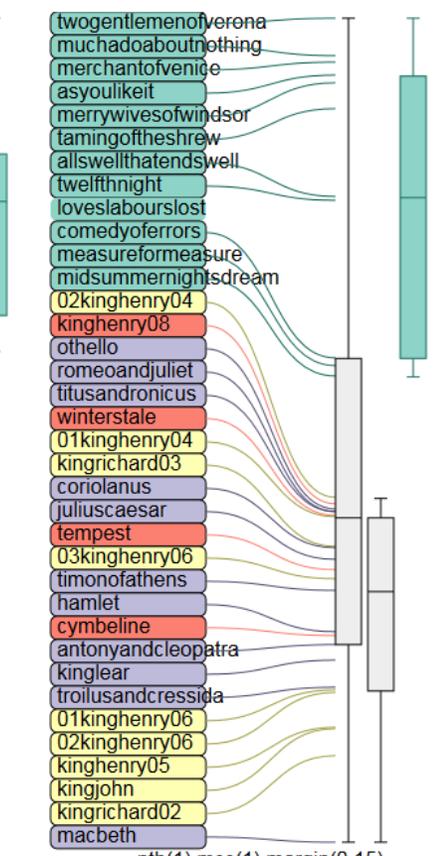
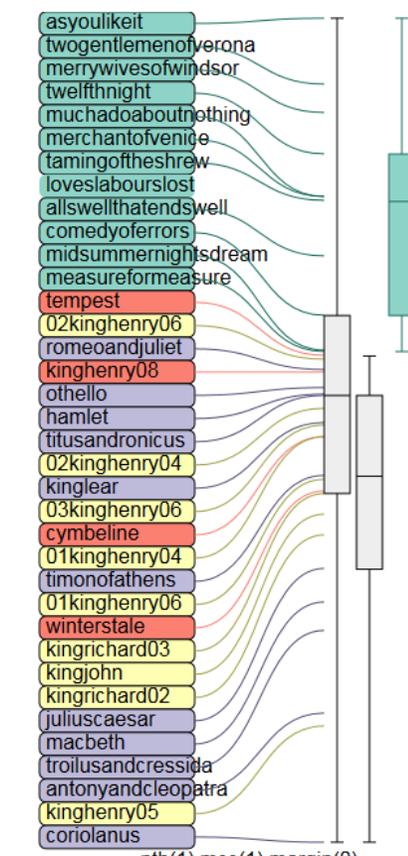
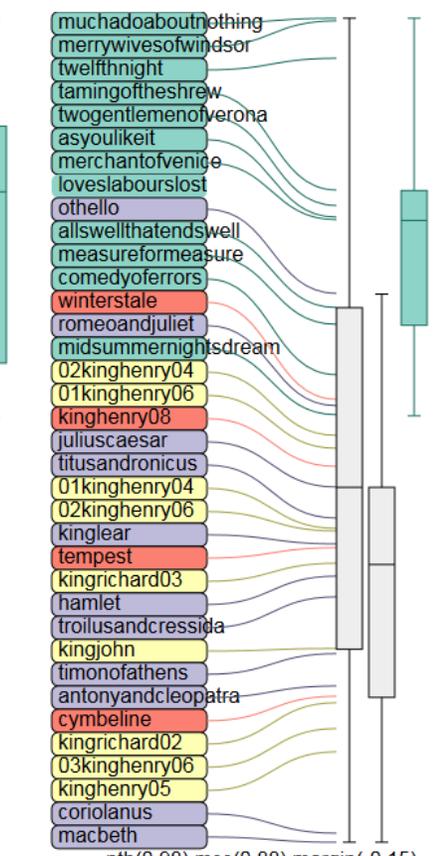
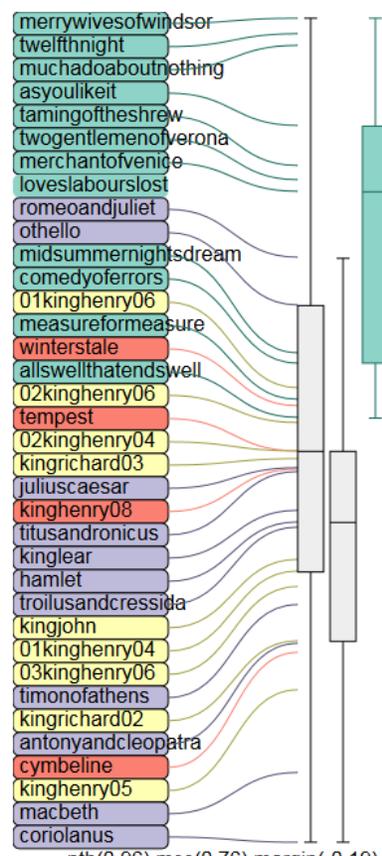
C - B - I (4 wrong)

C - I - 10 M (1 wrong)

31 D - 100 M - 3 A (none wrong)

“standard” L1 SVM (none wrong, reasonable margin)

25.3698 Q + 11.8823 U + 6.9492 F + 5.4897 A + 4.1489 P -  
3.3765 N + 2.6392 D + 2.0172 F - 1.5404 I + 1.1864 R - 0.7958  
C + 0.7272 D



nth(0.96) mcc(0.76) margin(-0.19)  
 Docuscope.PredictedFuture:1  
 Docuscope.Inclusive:-1

nth(0.98) mcc(0.88) margin(-0.15)  
 Docuscope.PredictedFuture:1  
 Docuscope.Inclusive:-1  
 Docuscope.Confidence:1

nth(1) mcc(1) margin(0)  
 Docuscope.MoveBody:-142.55  
 Docuscope.Curiosity:20.85  
 Docuscope.Inclusive:-14.6

nth(1) mcc(1) margin(0.15)  
 Docuscope.OpenQuery:25.37  
 Docuscope.Numbers:11.88  
 Docuscope.PredictedFuture:6.95  
 Docuscope.Acknowledge:5.49  
 and 8 others...

nth(1) mcc(1) margin(0.4)  
 Docuscope.Precedent\_Defending:-150.92  
 Docuscope.Apology:-56.12  
 Docuscope.Acknowledge:8.98  
 Docuscope.GenericEvents:-7.71  
 and 14 others...

nth(1) mcc(1) margin(0.58)  
 Docuscope.Precedent\_Defending:-125.36  
 Docuscope.Apology:-45.61  
 Docuscope.Procedures:25.52  
 Docuscope.OpenQuery:10.66  
 and 19 others...

# Simpler Functions



Easier to Understand



More likely  
to lead to Theory



Less Expressive



Less Likely  
to be Accurate  
(but more likely to overfit)

# Tradeoffs



Give the **user** control over the tradeoffs

But how do we help them make informed choices?

# Paths to model usability?

Interpretable models

simplify the models so they can be understood

Examinable models

look inside the models and hope you understand

Instance-based explanations

pick some decisions and try to understand them

**Experiment/Outcome Examination**

**look at the right input/outputs from the black box**

# Wrong?

Interesting Outliers

“Romeo and Juliet” is pretty comedic

Ambiguous Classifications

Late Plays are called Tragi-Comedies

Near-Misses

A tiny shift, and **this** would be different

twelfthnight
muchadoaboutnothing
twogentlemenofverona
asyoulikeit
merrywivesofwindsor
tamingoftheshrew
merchantofvenice
loveslabourslost
romeoandjuliet
othello
measureformeasure
tempest
comedyoferrors
winterstale
midsummernightsdream
kinghenry08
02kinghenry06
allswellthatendswell
kinglear
02kinghenry04
hamlet
kingrichard03
titusandronicus
timonofathens
troilusandcressida
01kinghenry06
kingjohn
juliuscaesar
cymbeline
03kinghenry06
kingrichard02
01kinghenry04
anthonyandcleopatra
macbeth
kinghenry05
coriolanus

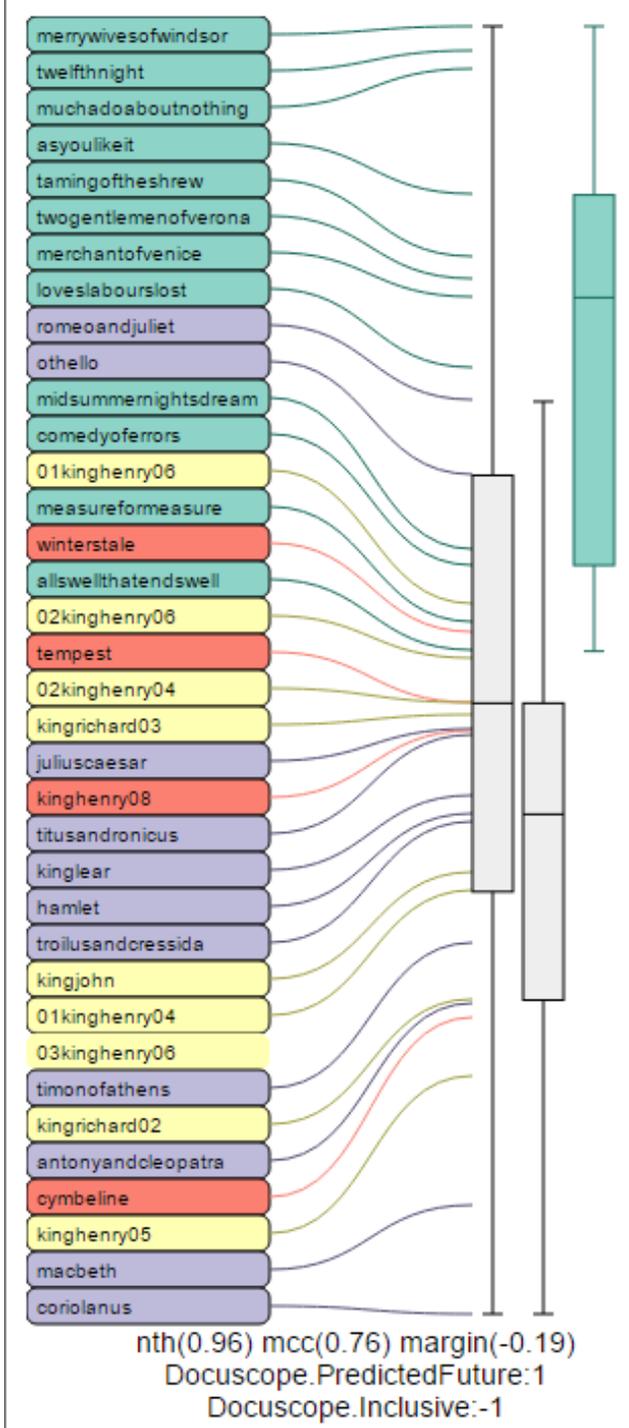
# Use Simple Models for Insight

Create models for comprehensibility

See what you get wrong

Look “inside” models to see how they work

Gleicher. *Explainers: Expert Explorations with Crafted Projections*. IEEE TVCG 19(12), Dec 2013. Proceedings VAST 2013.



Lesson

**Good tools for exploring your results are useful.**

# Is literature unique?

Large library has evolved

Understanding different kinds of variation

Bringing in other knowledge

Look at specific examples and outliers

Drill into details

Learn from validation experiments

Get **stakeholders** involved in  
all **data science phases**  
by helping them **understand**

# Some Lessons from “Humanist” Thinking

*They got along fine before “data-centric” thinking*

Importance of exemplars and outliers

Importance of going back to the source (specific passages)

## **Value of multiple points of view**

Contextualized arguments with lots of background

Editing and curation are scholarly activities

# Diversity

F + Q - I

C - M - I

P + N + D

## Same:

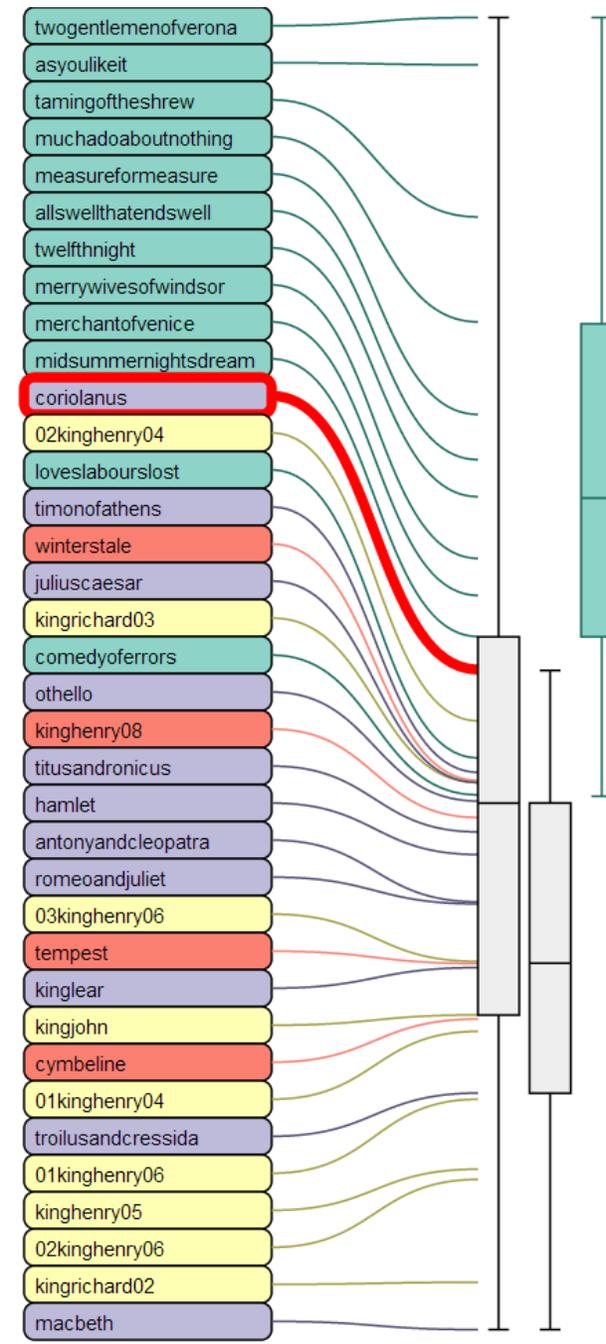
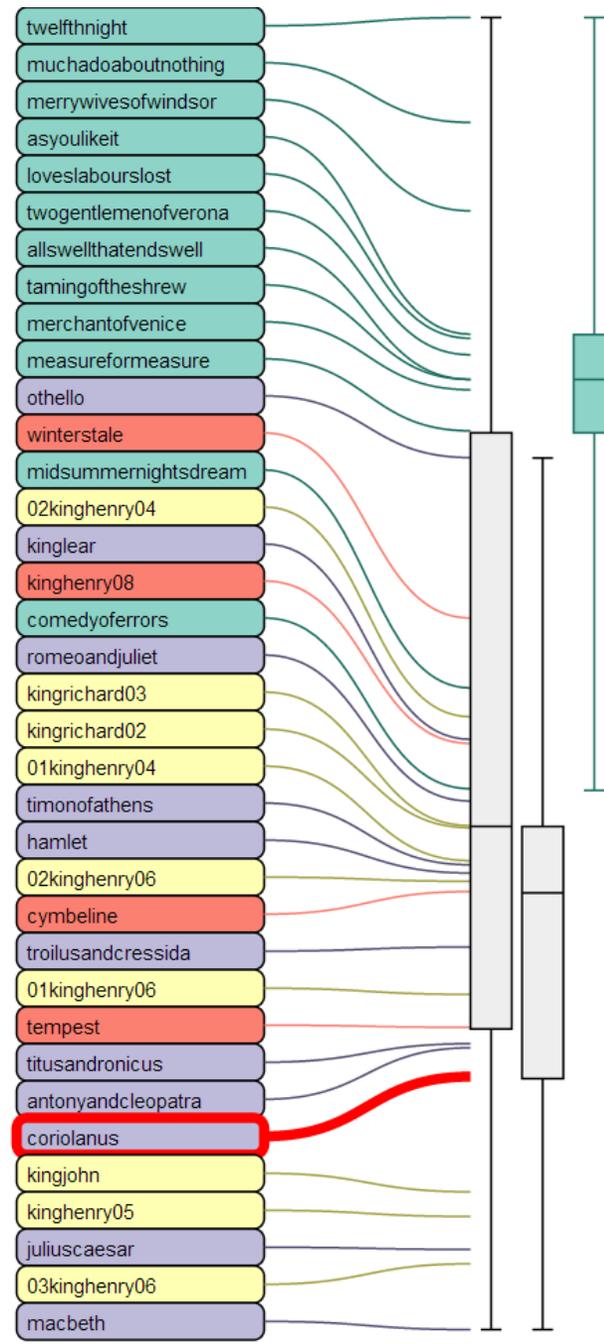
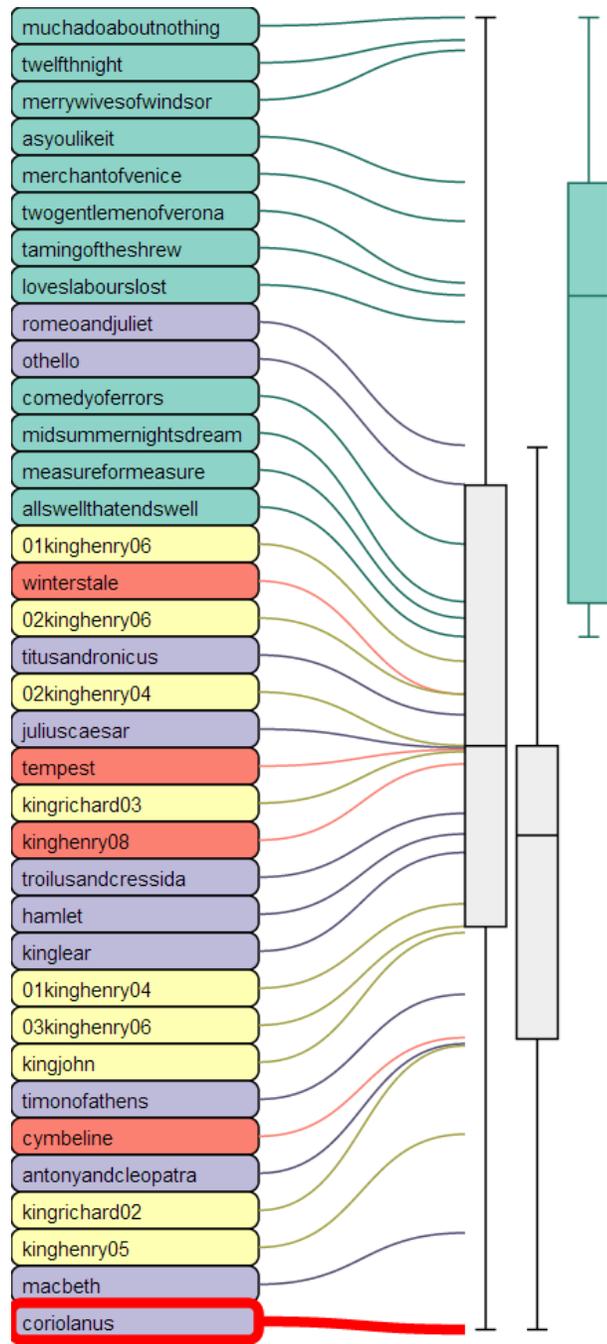
Correctness

Simplicity

## Different:

Explanations

Orderings



# Maybe Literature Scholars aren't so weird?

*They got along fine before “data-centric” thinking*

Importance of **exemplars** and **outliers**

Importance of going back to the **source** (specific passages)

Value of **multiple** points of view

Contextualized arguments with lots of **background**

**Editing** and **curation** are scholarly activities



Question.

Náttoohtemáuwetoowunk.

How prove you that there is a God?

Oobgôdje korâmen neh áttá Átandoh

Answer.

Anafquetâweten.



From the universal and constant agree-  
ment of all Nations, and persons  
wunk wutche wame arkèes, quah skeetambâwg mit-  
in the world, who are not void of  
tâuhkuk terre, owwânnak matta íáuwaióóguk wutche  
right reason and humanity.

fompáio penauwáuwuk quah renóowunk.

For the things which are grounded

Wutche ai akquíiks chawgwunsh wekakontamoo-

upon particular mens fancies

awk skeje nanfeáawk rénwawk róytammoúngansh  
and opinions are not acknow-

ledged of all men, and are

óomunks wutche wame rénawawk, quah wegonje  
of en changed but this notion that

ááflowunnamanóosh:webe (youh édyámmoóunk) neh  
there is a God is common to all men, nor is it chan-

Mandoo nânnarwee re wâine rénawawk matta ááflowú-  
ged by the changes of times;

numóoanas spe affowunnâmoúngansh quompious;  
therefore it must ai e from

règouche youh paughke móuche songème wutche  
some light, which is common to all

chawgun nowèta wequá-ai, teou nânnarwe re wame  
rèn-

But, this might be a lesson ...

Don't just give "them" your methods

If you listen carefully,  
they might have things to teach you

MICROFILMED - 1977

SOME  
**HELPS**  
FOR THE  
**INDIANS**

SHEWING THEM

How to improve their natural *Reason*, To know  
the *True GOD*, and the true *Christian Religion*.

1. By leading them to see the Divine Authority of the *Scriptures*.
2. By the *Scriptures* the Divine Truths necessary to *Eternall Salvation*.

Undertaken

*At the Motion, and published by the Order of the COMMISSIONERS of the United Colonies.*

by *ABRAHAM PEIRSON*,

Examined and approved by *Thomas Stanton* Interpreter-Generall to the *United Colonies* for the *Indian Language*, and by some others of the most able Interpreters amongst us.

L O N D O N,

Printed by *M. Simmons*, 1659.

# Thanks!

To you for listening.

To my students and collaborators.

To the NSF, NIH, DARPA and Mellon Foundation for funding.

## What Shakespeare taught us about (Visual) Data Science

**Michael Gleicher**

University of Wisconsin Madison

`gleicher@cs.wisc.edu`

