

Seeing, Scatterplots and Shakespeare

Michael Gleicher
 Department of Computer Sciences
 University of Wisconsin Madison
 (on sabbatical at INRIA Rhone-Alpes)
 Team Imagine



Happy (American) Thanksgiving!

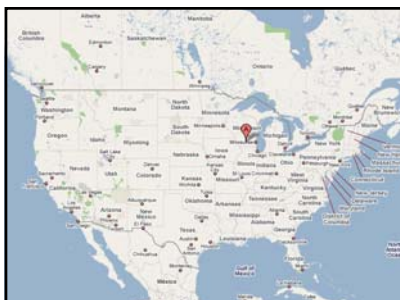
Acknowledgements

This work is a collaboration with a great set of people!

Students:
 Danielle Albers, Eric Alexander, Michael Correll,
 Adrian Mayorga


Collaborators:
 Steve Franconeri, Jonathan Hope, Robin Valenza,
 Mike Witmore

This research is funded in part, by the National Science Foundation and the Andrew Mellon Foundation



Seeing, Scatterplots and Shakespeare

Michael Gleicher
 Department of Computer Sciences
 University of Wisconsin Madison
 (on sabbatical at INRIA Rhone-Alpes)
 Team Imagine



Data Visualization

How do we use pictures to help understand, and communicate data?

Data Visualization

Seeing

Scatterplots


Shakespeare

Data Visualization

Seeing

Scatterplots

Shakespeare?



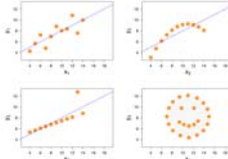
Perceptual Science:
 How do people see?
 How do we use this knowledge?

Data Visualization

Seeing

Scatterplots

Shakespeare?




Re-Examine Basic Methods:
 Consider foundations
 Impact what really gets used!

Data Visualization

Seeing

Scatterplots

Shakespeare?



Shakespeare?

Literary Scholarship as a motivating domain


Literary Scholarship as a way of thinking

New ways of thinking -> New approaches

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



SHAKESPEARE QUARTERLY

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

Visualizing English Print 1470-1800

What if you had access to all surviving books?

Why?

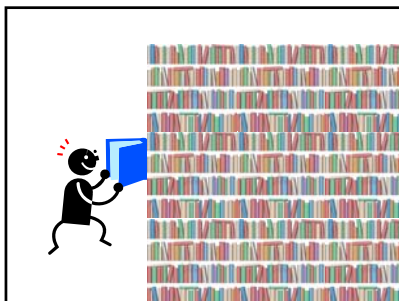
Consider larger collections of books

Consider language not content

See patterns across language

See small scale patterns in familiar texts

Be uncultured and still hang out with the cool kids



Literary Scholarship...

The statistics are not the argument

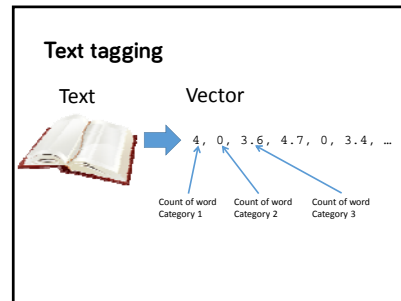
Exemplars and Outliers
Go back to the sources

Arguments based on context and knowledge

Multiple viewpoints and lenses

**Seeing Shakespeare:
Scalable Scholarship**
Study Literature without Reading?

Texts as Data?
Need to turn books into numbers



Just counting

Words (phrases) have a type (tag)

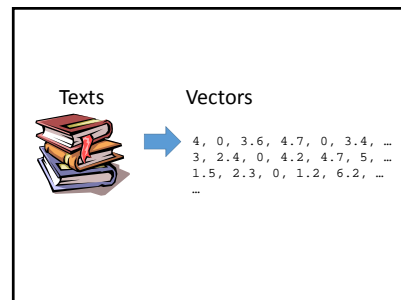
DocuScope
Simple matching against a dictionary
Hand-built dictionaries

100-115 Categories, 12-17 Clusters (groups)

to be ; or not to be ; that is the question ;
whether 'tis nobler in the mind to suffer
the slings and arrows of outrageous fortune ;
or to take arms against a sea of troubles ;
and by opposing end them ? to die : to sleep ;
no more ; and by a sleep to say we end
the heart-ache and the thousand natural shocks
that flesh is heir to ; 'tis a consummation
devoutly to be wish'd ; to die ; to sleep ;
to sleep : perchance to dream : ay , there's the rub ;
for in that sleep of death what dreams may come
when we have shuffled off this mortal coil ,
must give us pause : there's the respect
that makes calamity of so long life ;
for who would bear the whips and scorns of time ;
the oppressor's wrong , the proud man's contumely ,

15 "Clusters" (115 LATs)

- Subjective Register
- Emotion
- Description
- Institutional Register
- Academic Register
- Future
- Past
- Personal Relations
- Reasoning
- Interactive
- Elaboration
- Reporting
- Directing
- Narrative
- Character
- Not in dictionary
- Not Tagged



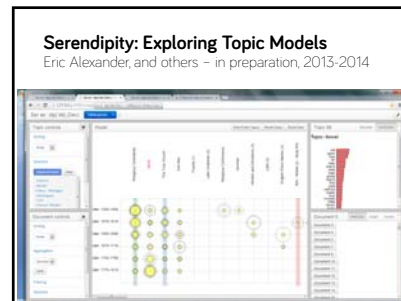
**What about more sophisticated
text analysis?**

Simple methods are:
Easier to understand and explain
Focused on word usage without considering meaning

**What about more sophisticated
text analysis?**

Topic modeling?
Yes, we're working on it too.

Eric Alexander, et al. Serendipity: turning topics back to texts.
In preparation. 2013-2014.

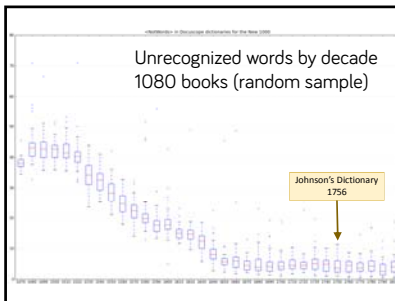


Just counting

Words (phrases) have a type (tag)

Docuscope
Simple matching against a dictionary
Hand-built dictionaries

100-115 Categories, 12-17 Clusters (groups)



How to look at 100+ dimensions?

How to look at 100+ dimensions?

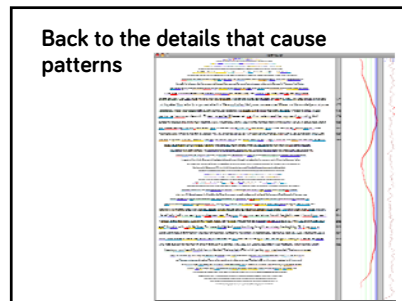
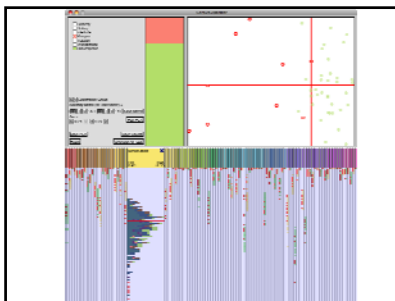
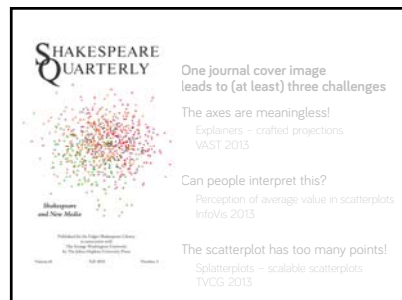
How to support "Humanist" arguments?


The statistics are not the argument

Exemplars and Outliers
Go back to the sources

Arguments based on context and knowledge

Multiple viewpoints and lenses





SHAKESPEARE QUARTERLY

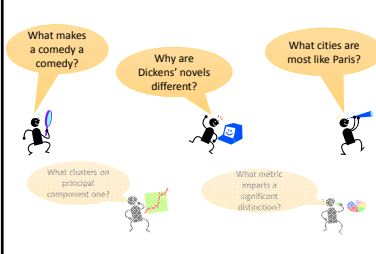
One journal cover image leads to (at least) three challenges

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

How to look at 100+ dimensions?



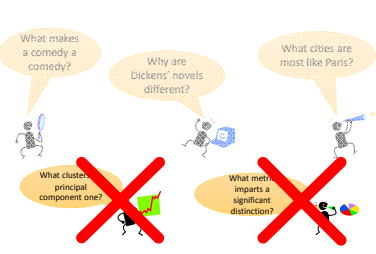
What makes a comedy a comedy?

Why are Dickens' novels different?

What cities are most like Paris?

What clusters on principal component one?

What metric imparts a significant distinction?



What makes a comedy a comedy?

Why are Dickens' novels different?

What cities are most like Paris?

What clusters on principal component one?

What metric imparts a significant distinction?

Explainers: Expert Explorations with Crafted Projections

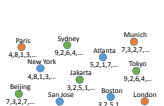
2013 IEEE VAST
Visual Analytics Science and Technology

Honorable Mention Award Winner

Explainers

An approach to explore high dimensional data

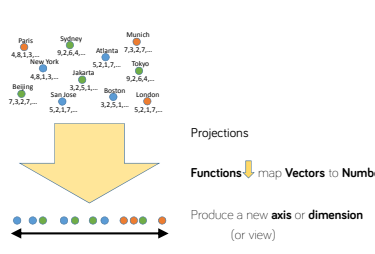
- Organize data according to **user**-defined concepts
- Explain **user**-defined concepts according to the data
- Give the **user** control over tradeoffs



High Dimensional Data

Objects \bullet have associated **Vectors**

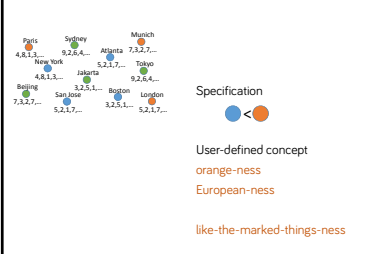
Paris 4,8,1,3...
Singapore 9,2,5,4...
Munich 7,3,2,7...
New York 4,8,1,3...
Atlanta 5,2,1,7...
Tokyo 9,2,6,4...
Jakarta 4,8,1,3...
Singapore 3,2,5,1...
Boston 9,2,6,4...
San Jose 7,3,2,7...
San Jose 5,2,1,7...
London 3,2,5,1...
London 5,2,1,7...



Projections

Functions \downarrow map **Vectors** to **Numbers**

Produce a new **axis** or **dimension** (or view)



Specification

User-defined concept

- orange-ness
- European-ness
- like-the-marked-things-ness

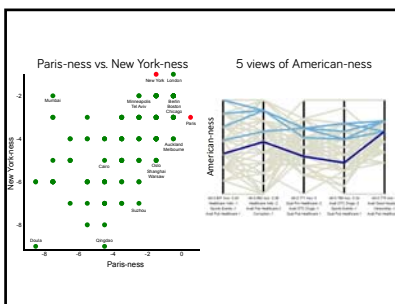
Explainers:
Projections crafted to meet user specifications

With user control of tradeoffs between:

Correctness:
does it align with the user specification?

Understandability:
can the user interpret the mapping?

Diversity:
can we generate alternate mappings?



Organize
Relationships between points based on data and concepts

Rankings
Outliers
Extrema
Exemplars
Similarities

Explain
Relationships with the data connect concepts and variables

Where do the orderings come from?

Are variables correlated with concepts?

To make things concrete

An Example: Shakespeare's Plays

More Examples online!
<http://graphics.cs.wisc.edu/Via/Explainers>

Texts → Vectors

36 Plays = 36 Vectors

115 "Measurements" of each text = 115 dimensions

genre

Categorization given by Shakespeare's contemporaries

Comedy Tragedy History

Category for plays written after that

Late Plays

comedic-ness

A measure of how much of a comedy something is

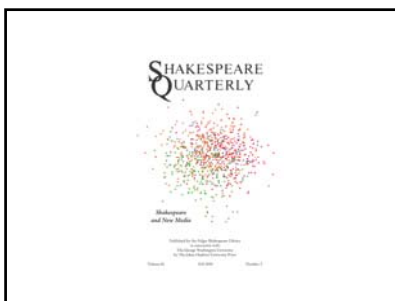
It's the "stuff" comedies have more of
Where "stuff" has to be in the data

Tragedy
History
Late Plays

Comedy

Organization:
What is most/least comedic?

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



A comedicness explainer

$c = f(V)$
 c = comedicness
 f = a function (**Explainer**) that maps from V to c
 V = vector from a text (length 115)

Choose f such that

1. It is correct (meets specification)
2. It is understandable (simple)
3. We can have alternatives (other functions that meet 1 and 2) as well

An Explainer

comedicness = $M - I$
 $f(V) = V[39] - V[42]$

Understanding the Visualization*

* The Visual Encoding is not a strong part. Suggestions are most welcomed!

Understanding the Visualization*

Objects in rank order
 Color by specified class

* The Visual Encoding is not a strong part. Suggestions are most welcomed!

Understanding the Visualization*

Lines connect rank (left) to value (position on number line)

* The Visual Encoding is not a strong part. Suggestions are most welcomed!

Understanding the Visualization*

Box Plots show class separation
 Left: all data
 Right: each class of interest

* The Visual Encoding is not a strong part. Suggestions are most welcomed!

An Explainer

comedicness = $M - I$
 $f(V) = V[39] - V[42]$

Tradeoff:
simple (linear, 2 variables, unit coefficients)
 but
5 "wrong"

5 Wrong

False Positives →
 False Negative →

Wrong?

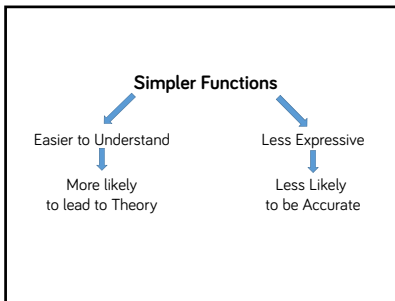
Interesting Outliers →
 "Romeo and Juliet" is pretty comedic
 Ambiguous Classifications →
 Late Plays are called Tragi-Comedies
 Near-Misses →
 A tiny shift, and this would be different

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

What's Understandable*?

- Simple form (linear vs. non-linear, ...)
- $A+B$ vs. $e^{-aA(A+B)}$
- Parsimony (few variables)
- $A+B$ vs. $W+X+Y+Z$
- Simple Coefficients (small integers)
- $A = 2B$ vs. $1235 A = 4.327 B$
- Familiar Variables
- $A + B$ vs. $Q + W$



Tradeoffs



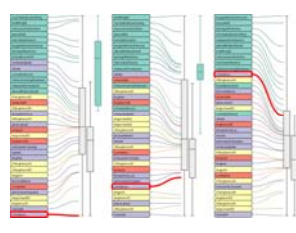
Give the **user** control over the tradeoffs

Diversity

$F + Q = 1$
 $C = M = 1$
 $P + N = D$

Same:
 Correctness
 Simplicity

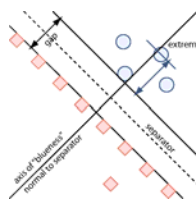
Different:
 Explanations
 Orderings



How to find functions?

Optimization problem
 Minimize amount "wrong"
 "Cost" of function

Support Vector Machine (SVM)



How to Implement Explainers?

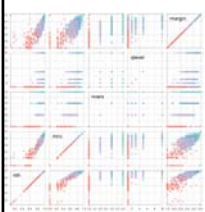
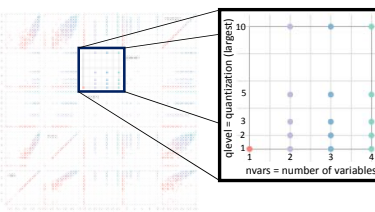
<p>Fancy Math Encode tradeoffs into the SVM Solve for the best tradeoffs Adjust parameters to tune</p> <p>Solve one big optimization problem</p> <p><i>Not a standard SVM, so needs a slow and finicky solver</i></p>	<p>Brute Force Sample space of variable sets Solve an SVM for each Sort and filter to find interesting ones</p> <p>Solve many small optimization problems</p> <p><i>Generates a diverse and interesting exploration of tradeoffs</i></p>
---	--

Finding the interesting explainers

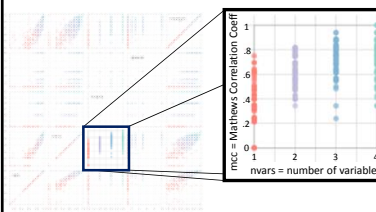
2667 Explainers generated

Rank-by-Feature or "Scagnostics" style analysis

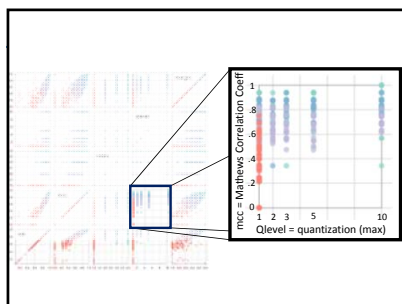
2667 points - each an explainer
 5 properties of each

qlen = quantization (largest)
 nvars = number of variables



mcc = Matthews Correlation Coeff
 nvars = number of variables



Isn't this just...
Some prior approaches to help situate our work

Explainers

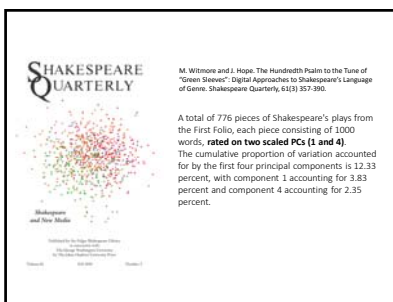
Organize data according to user-defined concepts
Explain user-defined concepts according to data

Explainers add:
User-defined concepts
Control over tradeoffs
Connection between concepts and variables
Generation of alternatives

Dimensionality Reduction
e.g. PCA, CCA, IsoMap, ... - standard statistical and ML practices

~~Organize data according to user-defined concepts~~
~~Explain user-defined concepts according to data~~

Explainers add:
User-defined concepts
Control over tradeoffs
Connection between concepts and variables
Generation of alternatives



Machine Learning Classification Techniques

~~Organize data according to user-defined concepts~~
~~Explain user-defined concepts according to data~~

Explainers add:
User-defined concepts
Control over tradeoffs
Connection between concepts and variables
Generation of alternatives

User-Driven Spatializations
e.g. Semantic Interaction (Enders+), LAMP (Paulovich+), Star Coordinates (Kandogan), ...

~~Organize data according to user-defined concepts~~
~~Explain user-defined concepts according to data~~

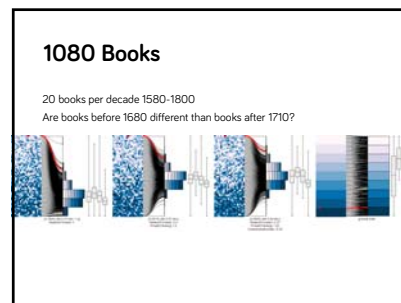
Explainers add:
User-defined concepts
Control over tradeoffs
Connection between concepts and variables
Generation of alternatives

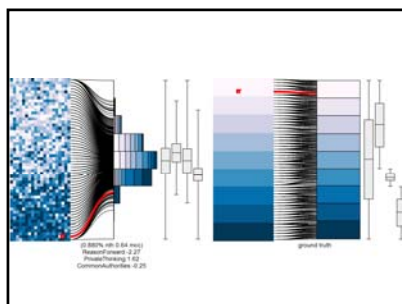
More to do . . . (current limitations)

User Experience
Visualizations
Interactive Specification

Theory
Understanding understandability tradeoffs
Statistical significance in negative results

Scalability
More variables (redundancy)
More objects
More complex relationships





Key Ideas

User-defined concepts
 Multiple goals: **organize** and **explain**
 User-control over **tradeoffs**: correctness, simplicity, diversity
Alternative viewpoints

Details:
 Types of Simplicity
 Implementation with SVM

Explainers

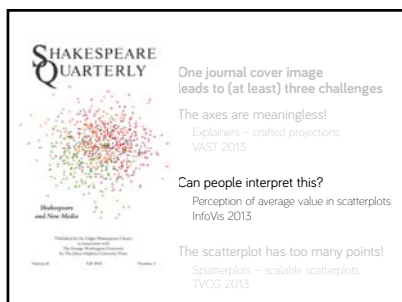
An approach to exploration and discovery in high dimensional data that

- organizes** data according to **user-defined concepts**
- explains** these **user-defined concepts** in terms of the data and generates **alternative viewpoints**

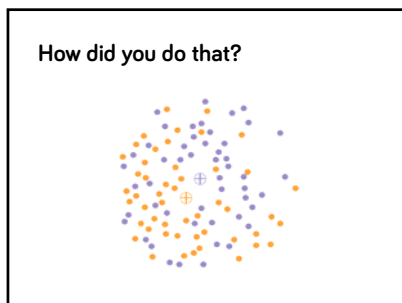
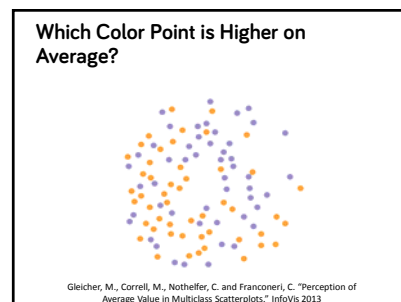
using machine learning techniques and providing control over **tradeoffs**.

More Examples Online!
<http://graphics.cs.wisc.edu/Vis/Explainers>

Acknowledgements:
 This work would not have been possible without my fantastic former collaborator, and the optimization and machine learning wizard in my department.
 This work is supported in part by the Andrew Mellon Foundation through the "Visualizing English Text" project. This work is supported in part by NSF Awards IIS-1162027/CMMI-094103, and OAC-1247262.



What can you do with too many points?



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	Practical Application
Averages in Time Series Correll, et al. CHI 2012	Sequence Surveyor (Genetics) Albers, et al. InfoVis 2011
Tagged Text Correll, et al. CHI 2013	LayerCake (Virus mutations) Correll, et al. BoVis 2011
Scatterplot Averages Gleicher, et al. InfoVis 2013	Molecular Surface Experiments Sankaya, et al. (in prep)
Other statistics in Time Series Albers, et al. CHI 2014[*]	Decision Making Correll, et al. (in prep)

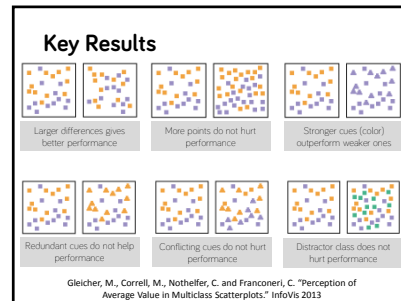
Ask the Turkers!

Crowdsource participants on Amazon's Mechanical Turk service
 Careful Design to get valid results

Measure accuracy not speed
 More like real tasks
 Adjust hardness - not time allowed

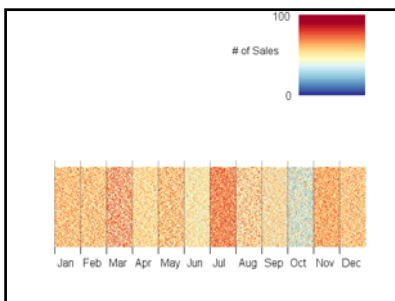
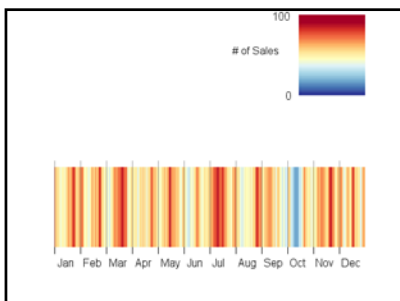
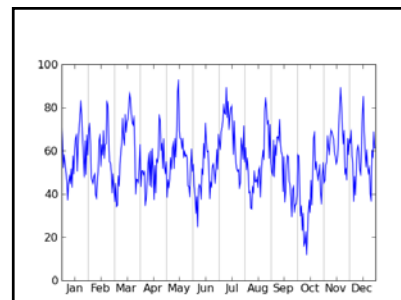
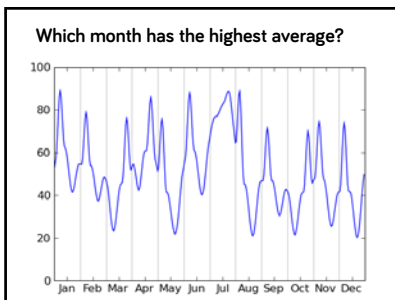
Scatterplot experiments

- Between subjects study
 - Series of experiments (one per condition)
 - 32 participants per condition
 - Established relations between conditions
- Within subjects study (repeated measures)
 - Reconfirm consults
 - 32 participants per experiment
 - Run pairs of conditions for key results
- Other experiments have similar designs



Time Series?

Correll, Albers, Franconeri, Gleicher. Comparing Averages in Time Series Data. CHI 2012.



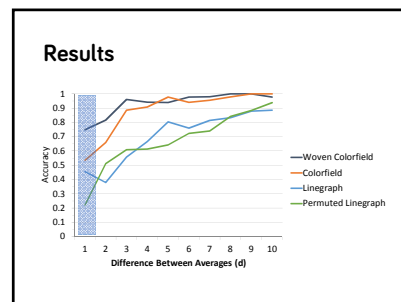
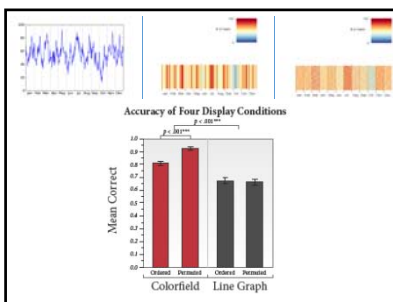
Conditions

Linegraphs:
 Regular or 1D permuted

Colorfields:
 Regular or woven

Conditions

- Noise of signal
- Gap between winner and distractor
- Number of distractor months



Conditions

$d = 10$ $d = 2$

What besides averages?

Albers, Correll, Gleicher. Task-Driven Design Variables for Time Series Visualization. (under review)

Things You Might Care About

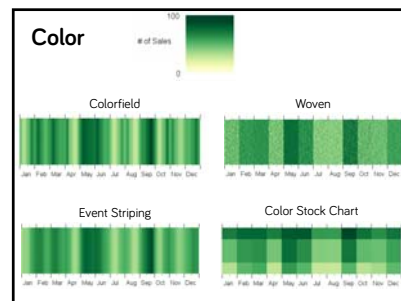
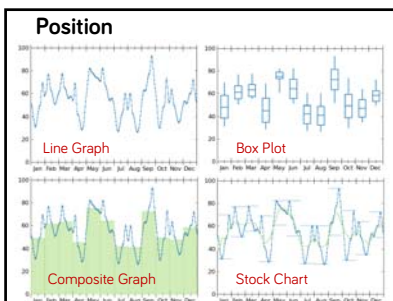
- Maxima:** Which month had the day with the highest sales for the year?
- Minima:** Which month had the day with the lowest sales for the year?
- Range:** Which month had the largest range of values?
- Average:** Which month had the highest average sales for the year?
- Spread:** Look at the average sales from each month. Which month had the sales which were the most spread out from their monthly average?
- Outliers:** Which month had the most unusual (outlier) sales days?

Things That Might Matter

What **visual variable** do we use to encode the data? Position, color, size, alpha, ...?

What **derived data** do we want to explicitly show? Extrema, averages, medians, modes, quartiles, ...?

How should we compute this derived data? Discretely, continuously, simplified, complex, ...?

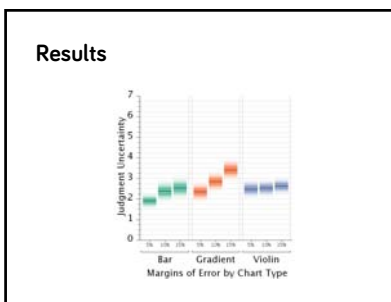
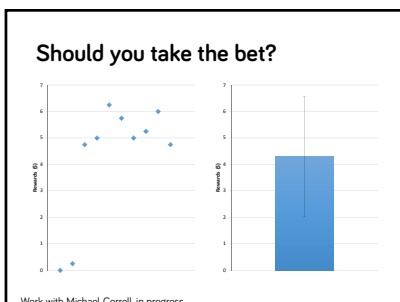


Results

Facility	Maxima	Minima	Range	Average	Spread	Outliers
Line graph	87.5%	75.0%	75.2%	82.25%	48.2%	76.7%
Modified Stack Chart	87.5%	75.0%	75.2%	82.25%	48.2%	76.7%
Box Plot	87.5%	75.0%	75.2%	82.25%	48.2%	76.7%
Composite Graph	87.5%	75.0%	75.2%	82.25%	48.2%	76.7%
Colorfield	87.5%	75.0%	75.2%	82.25%	48.2%	76.7%
Color Stack Chart	87.5%	75.0%	75.2%	82.25%	48.2%	76.7%
Warm Colorfield	87.5%	75.0%	75.2%	82.25%	48.2%	76.7%
Event Sampling	87.5%	75.0%	75.2%	82.25%	48.2%	76.7%

- ### Why Visual Aggregation?
- Why not just give them the answer?
1. You may not know what the viewer wants
 2. Some aggregate properties might be complicated
 3. You can't show all properties
 4. It gives the viewer more information
 5. Doing "work" might force them to think about things

Can better visualizations lead to better decision making?



Quantity Estimation in Visualizations of Tagged Text

Correll, Alexander, Gleicher, CHI 2013

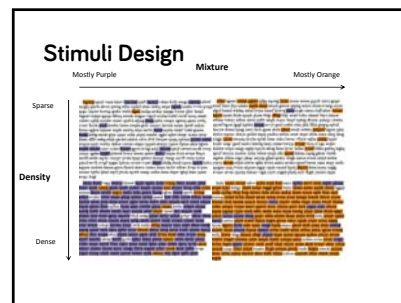
More blue or purple?

"I don't think they play at all fairly," Alice began, in rather a complaining tone, "and they all quarrel so dreadfully one can't hear one's self speak—and they don't seem to have any rules in particular; at least, if there are, nobody attends to them—and you've no idea how confusing it is all the things being alive; for instance, there's the arch I've got to go through next walking about at the other end of the ground—and I should have croqueted the Queen's hedgehog just now, only it ran away when it saw mine coming!"

Contributions

Empirical analysis of comparison in tagged text.

Design considerations to improve comparison.



Results

Two columns of dense, illegible text representing experimental results.

Density

9 samples of density, random mixtures

Word Length

3 uniform distributions of word length

Multi-category Comparison

Which color is most common?

Two columns of dense, illegible text.

d parameter

d parameter

Color bias

Asymmetric Area

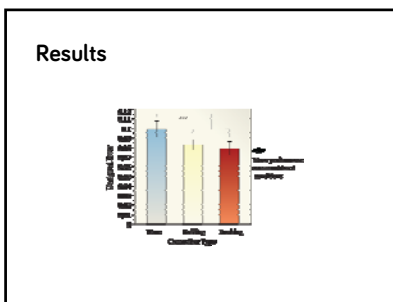
Which type of word is most common?

Two columns of dense, illegible text.

Area Corrections

Two blue rectangular boxes labeled 'Padding' and 'Tracking'.

Area Corrections



to be ; or not to be ; that is the question ;
 whether 'tis nobler in the mind to suffer
 the slings and arrows of outrageous fortune ;
 or to take arms against a sea of troubles ;
 and by opposing end them ? to die : to sleep ;
 no more ; and by a sleep to say we end
 the heart-ache and the thousand natural shocks
 that flesh is heir to ' ; tis a consummation
 devoutly to be wish'd : to die : to sleep ;
 to sleep : perchance to dream : ay , there's the rub ;
 for in that sleep of death what dreams may come
 when we have shuffled off this mortal coil ,
 must give us pause ; there's the respect
 that makes calamity of so long life ;
 for who would bear the whips and scorns of time ;
 the oppressor's wrong ; the proud man's contumely ;

15 "Clusters" (115 LATs)

- Subjective_Register
- Emotion
- Description
- Institutional_Register
- Academic_Register
- Future
- Past
- Personal_Relations
- Reasoning
- Interactive
- Elaboration
- Reporting
- Directing
- Narrative
- Character

Not in dictionary
Not Tagged

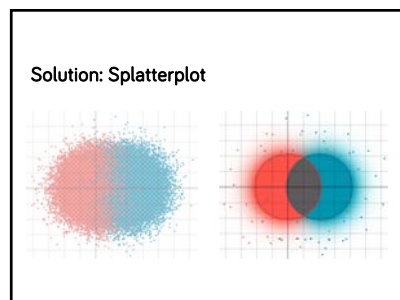
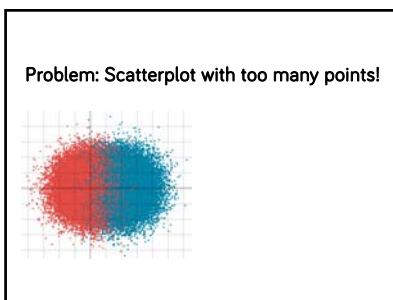
SHAKESPEARE QUARTERLY

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

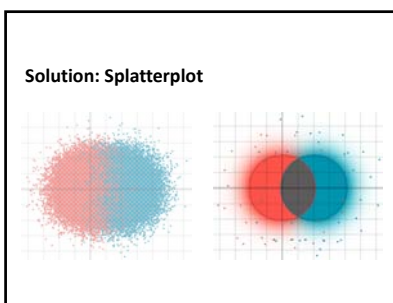
The axes are meaningless!
 Explainers - crafted projections
 VAST 2013

Can people interpret this?
 Perception of average value in scatterplots
 InfoVis 2013

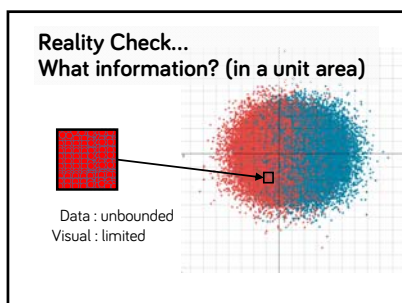
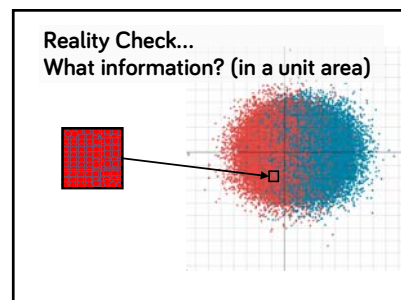
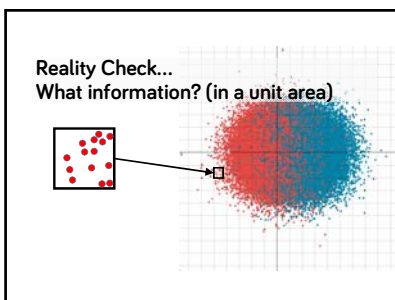
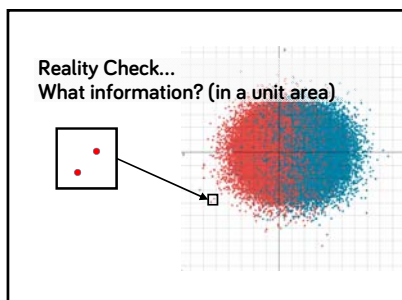
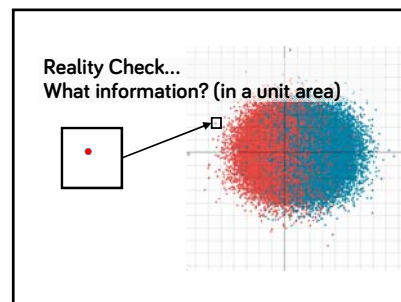
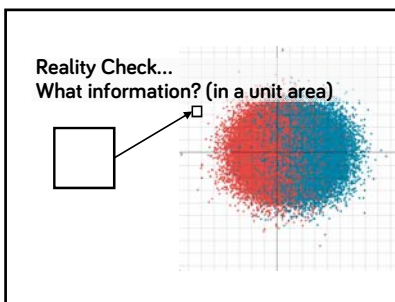
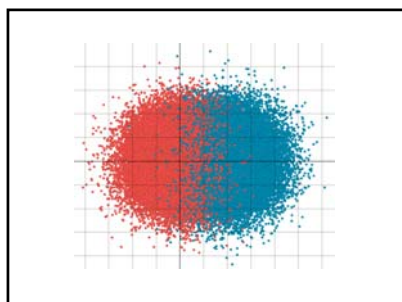
The scatterplot has too many points!
 Splatterplots - scalable scatterplots
 TVCG 2013



Scatter plots suffer from overload.



what if you have lots of points?

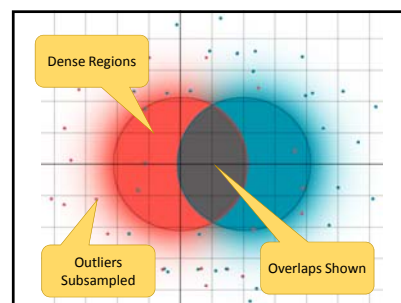


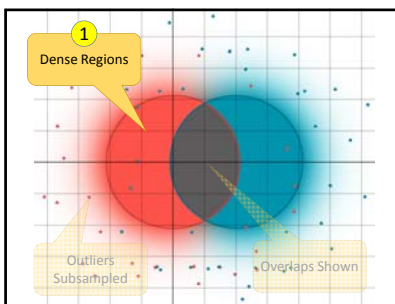
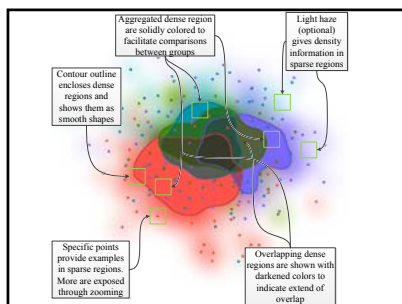
Bounded Information Density

In a Unit of Area:
Amount of data is **unbounded**
Amount you can see is **limited**

Need to **limit** the amount shown

Choose what to display by **abstracting** the data





Kernel Density Estimation (KDE)

Count how many points near every position

Weight by distance

Size of kernel (circle) is the bandwidth

Creates smooth fields

Screen Space KDE

Parameters based on perceptual properties

Independent of data

Does the right thing when you zoom

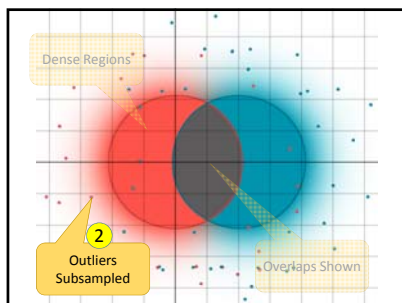
Discrete dense regions

Threshold

Why? (single set case)
Dynamic range of density may be high and hard to encode
At some point, it's just "dense"
Crisp boundaries are better visually

Information is thrown away!

Information is thrown away!
Interactive control of threshold
Encode sparse regions differently



Subsample sparse regions

To Haze or not to Haze?

Edges

Strokes

Clear Clutter

Both require distance to region

Contours?

Complicated with multiple groups

Dense Regions*

Outliers Subsampled

3 Overlaps Shown

Multiple Groups

Compute densities independently

Color per group

Pick distinctive colors

Colors for combinations

Multi-variate color encoding?
Map R^n to a color

Colors for combinations

Multi-variate color encoding?
Map R^n to a color

Colors for set combinations
Map 2^n set combinations to colors

Color Blending

Encode sets with color

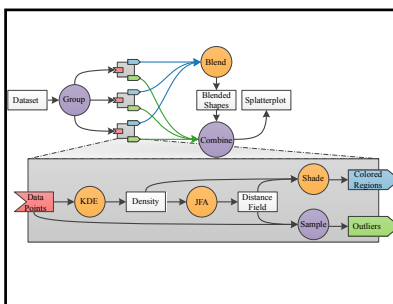
Hue = set

Lightness = number of overlaps

See evaluation in paper

Implementation

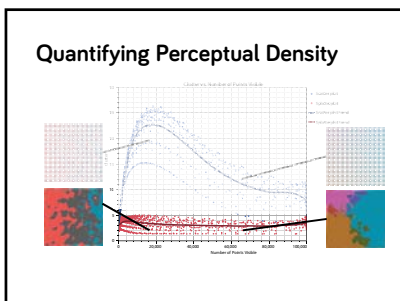
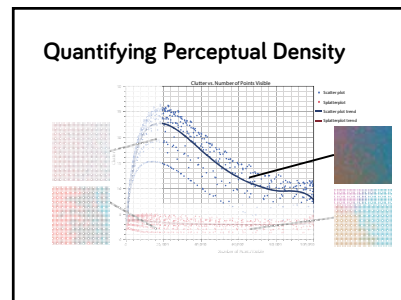
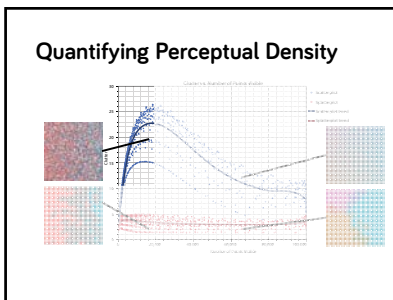
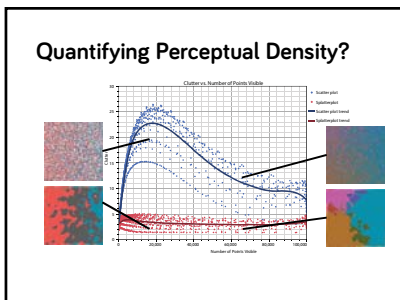
Interactivity is critical!



Performance: Use the GPU

Draw points
 Filter (convolution) for KDE
 Jump Flood for distances
 Render each set and combine

Lots of points - fast
 Lots of groups - less fast



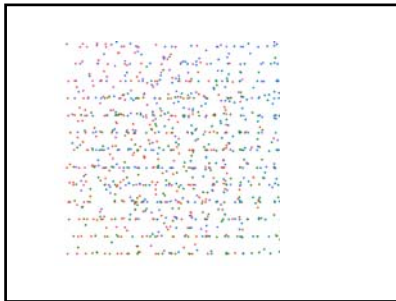
Other ideas?

There is plenty of "related work" in research in practice

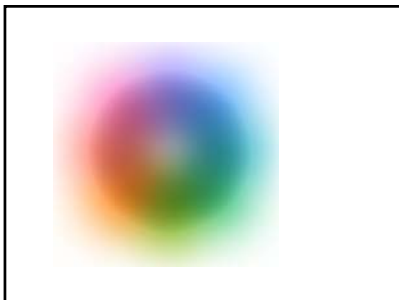
Key Novelties in Splatterplots
 Choose abstractions to understand set relationships
 Screen space density estimates
 Dual Encodings



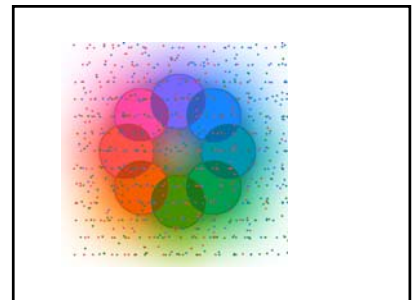
subsample?



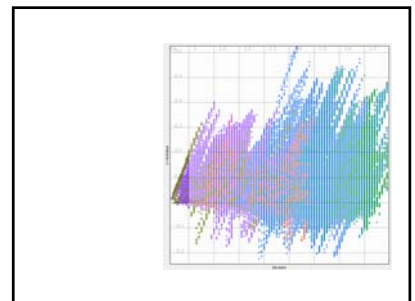
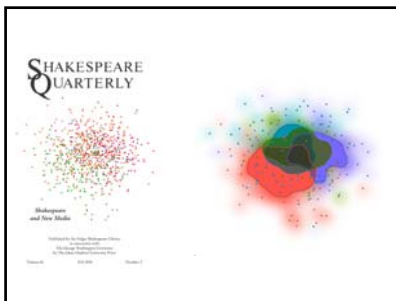
histograms and KDEs

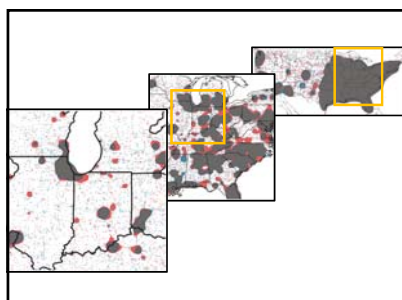
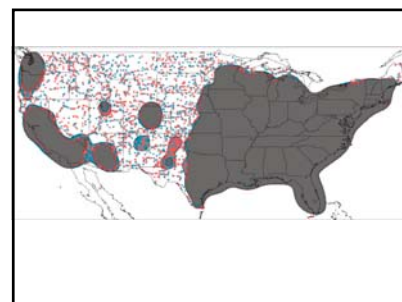
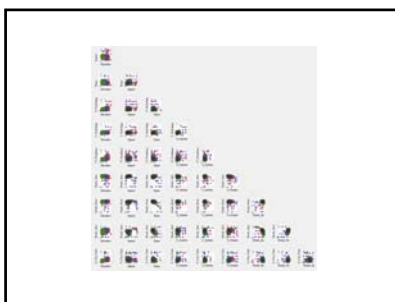
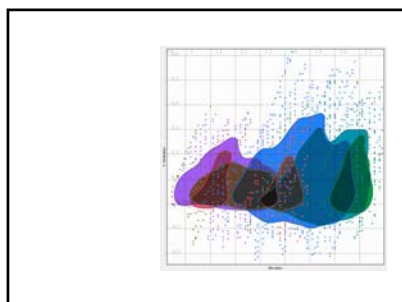


Splatterplot!



**Real (or realistic)
Examples**
The synthetic data is pretty but...





More to do!

Theory	Practice
Understand Visual Density	WebGL implementation
Consider other tradeoffs	Massive Data Handling
Other Types of Data	Evaluation (see InfoVis paper)
3D (volumes)	Non-GPU version for my laptop ☺

Spatterplots
Scalable Display of Scatter Data

Bounded visual complexity
Screen space density estimation
Dual encodings
GPU Implementation

Acknowledgements
This work is supported in part by the Andrew Mellon Foundation through the "Visualizing English Poetry" project. The work is also supported in part by NSF Awards IS-1102032, CMM-094022, and CMM-092262.

SHAKESPEARE QUARTERLY

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!
Spatterplots - scalable scatterplots
TVCG 2013

Data Visualization

Inspiration comes from a number of directions:

- Seeing
- Scatterplots
- Shakespeare

Thanks!

To you for listening.
To the organizers for inviting me
To my students and collaborators.
To the NSF and Mellon Foundation for funding.

Seeing, Scatterplots and Shakespeare

Michael Gleicher
University of Wisconsin Madison
(on sabbatical at INRIA Rhone-Alpes)
gleicher@cs.wisc.edu