# Global Big Data Conference

## BIG DATA BOOTCAMP

**Denver**

September 30th, Oct 1st & 2nd 2016

Colorado Convention Center, 700 14th St, Denver, CO 80202

COLORADO CONVENTION CENTER

www.globalbigdataconference.com

Twitter : @bigdataconf

# Enabling a dialog between **People** & **Data**

Lessons in Designing for Big Data

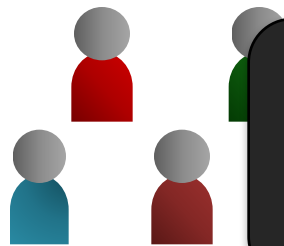College of Media, Communication and Information
UNIVERSITY OF COLORADO **BOULDER**

Danielle Albers Szafir
Assistant Professor
Department of Information Science
Danielle.Szafir@Colorado.edu

More Analysts
(& Fewer Experts)

More Heterogeneity

More Questions

Why don't we just compute the answer?

## Data Sample 1:

Mean(x) = 9
Variance(x) = 11
Correlation(x, y) = 0.816
Regression: y = 3 + 0.5x

## Data Sample 2:

Mean(x) = 9
Variance(x) = 11
Correlation(x, y) = 0.816
Regression: y = 3 + 0.5x

## Data Sample 3:

Mean(x) = 9
Variance(x) = 11
Correlation(x, y) = 0.816
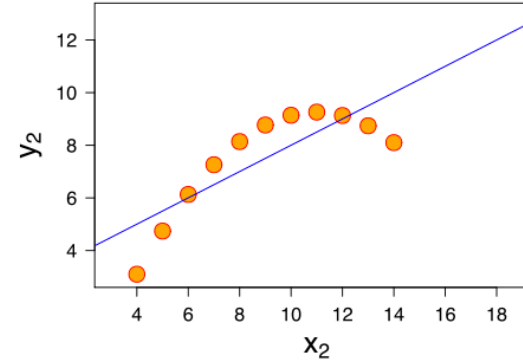Regression: y = 3 + 0.5x

## Data Sample 4:

Mean(x) = 9
Variance(x) = 11
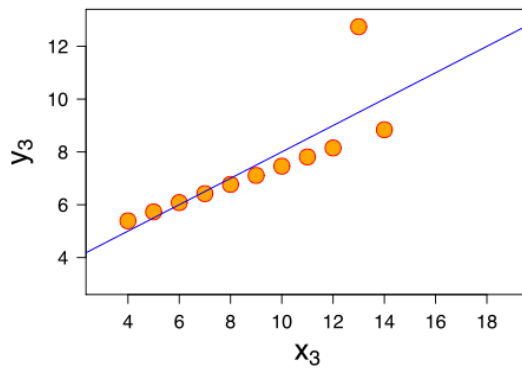Correlation(x, y) = 0.816
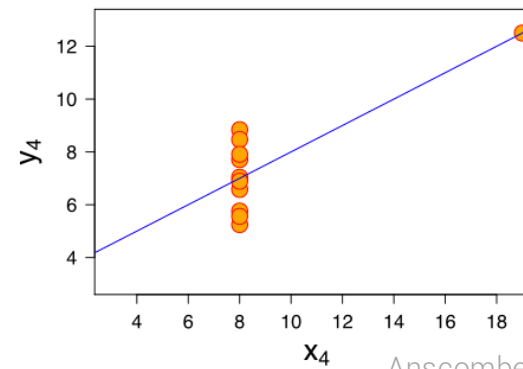Regression: y = 3 + 0.5x

Data Sample 1:

Data Sample 2:

Data Sample 3:

Data Sample 4:

Anscombe, American Statistician, 1973

© D.A. Szafir, 2016

Statistical tools are powerful, but the human brain understands patterns

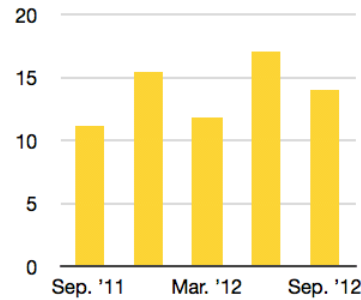# Apple Earnings Dashboard: September 2012 Quarter

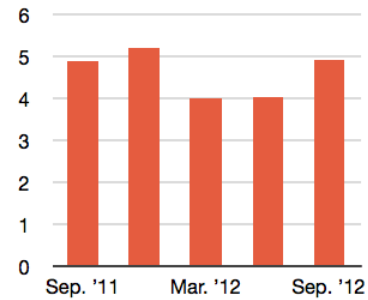## How Apple's Revenue Stacks Up (Billions)
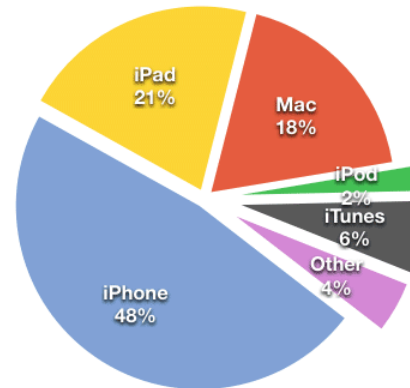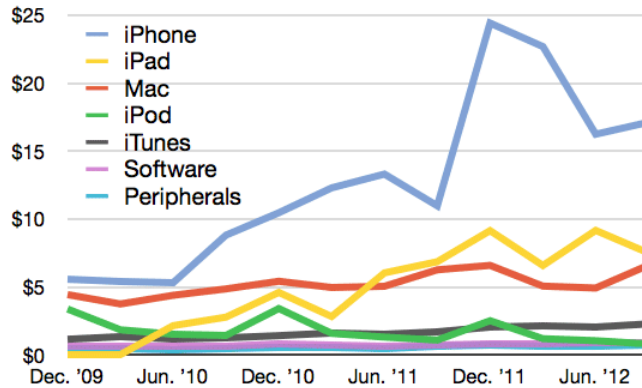


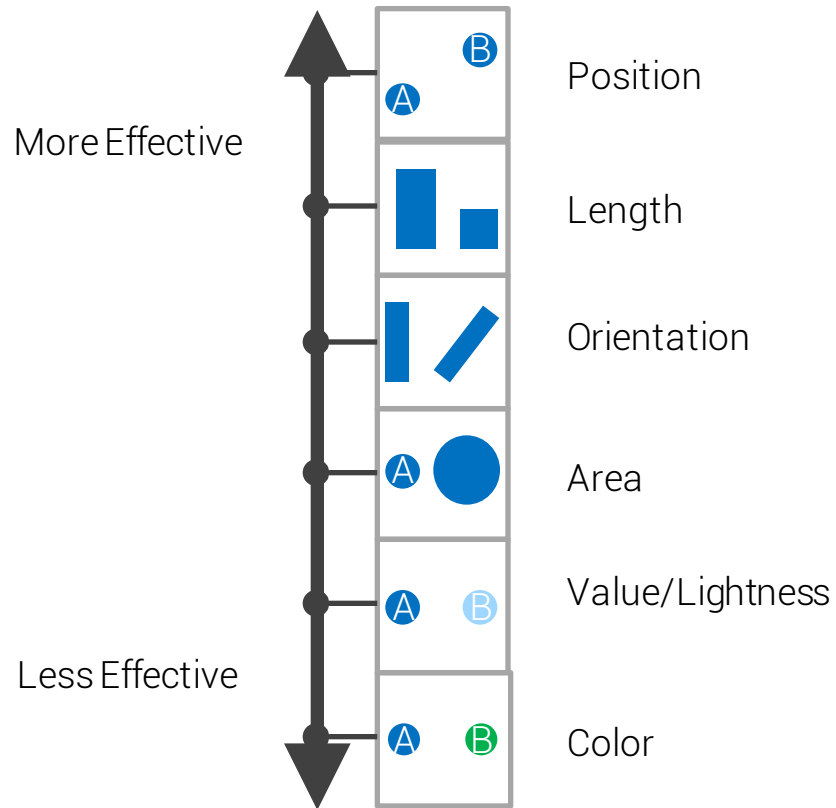## iPhone shipments (Millions)



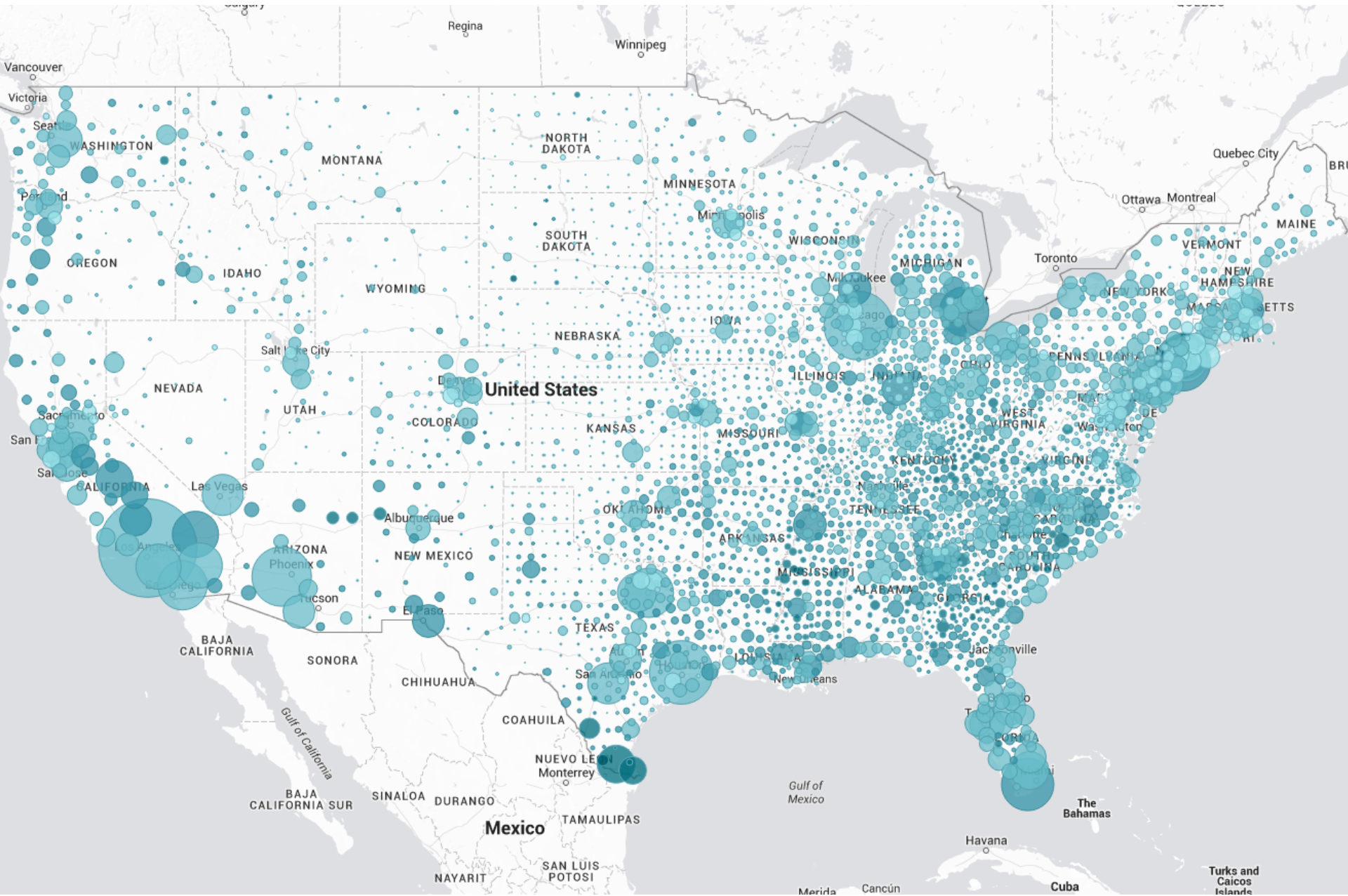## iPad shipments (Millions)



## Mac shipments (Millions)



## Revenue by Product (Billions)



- iPhone
- iPad
- Mac
- iPod
- iTunes
- Software
- Peripherals



iPhone 48%
iPad 21%
Mac 18%
iPod 2%
iTunes 6%
Other 4%

More Effective

Position

Length

Orientation

Area

Value/Lightness

Less Effective

Color

Cleveland & McGill, 1985

© D.A. Szafir, 2016

# What happens when our tools don't suit our data?

http://www.nytimes.com/newsgraphics/2014/01/05/poverty-map/?ref=multimedia
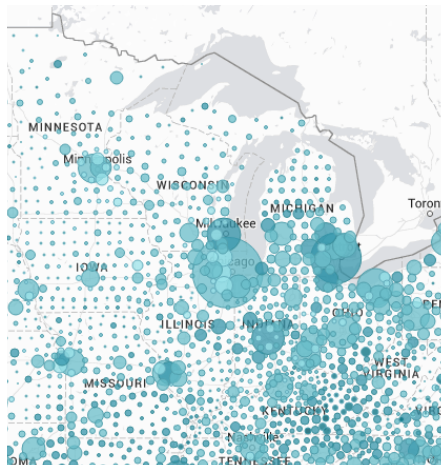
# Low-Level Tasks ➔ Individual Values
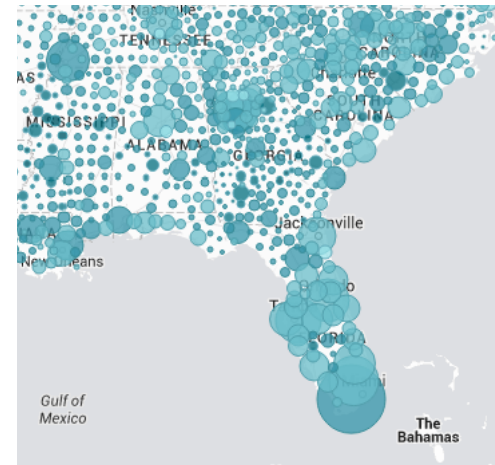
Los Angeles

Phoenix

# High-Level Tasks ➔ Combine Many Values

Midwest

Southeast

# Four types of ensemble coding in data visualizations

**Danielle Albers Szafir**

Department of Computer Sciences,
University of Wisconsin–Madison, Madison, WI, USA

**Steve Haroz**

Department of Psychology,
Northwestern University, Evanston, IL, USA

**Michael Gleicher**

Department of Computer Sciences,
University of Wisconsin–Madison, Madison, WI, USA

**Steven Franconeri**

Department of Psychology,
Northwestern University, Evanston, IL, USA

**Ensemble coding supports rapid extraction of visual statistics about distributed visual information. Researchers typically study this ability with the goal of drawing conclusions about how such coding extracts** Kahn, 2012). Other types of information can be extracted and combined in parallel from large numbers of objects at once, such as the average object size (Ariely, 2001). A growing body of work seeks to

# Binary Comparisons don't scale!

**Visual Aggregation Task**

|  | Identification (Outlier) | Summary (Mean) | Segmentation (Clustering) | Structure Estimation (Trends) |

# Big Picture Analyses

Computational Aggregation:

Visual Aggregation:



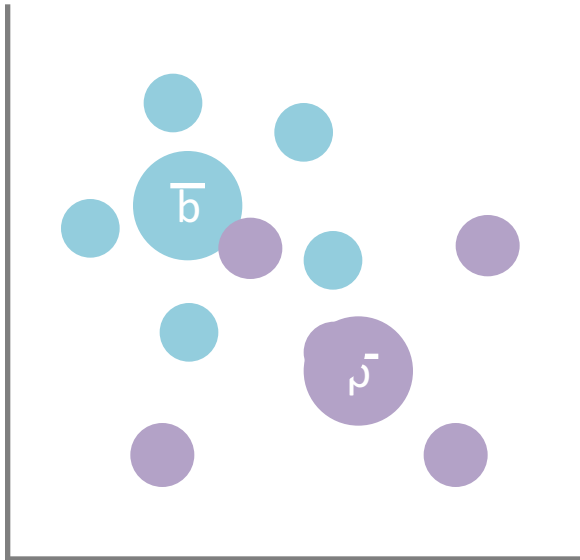Compute the answer then visualize it
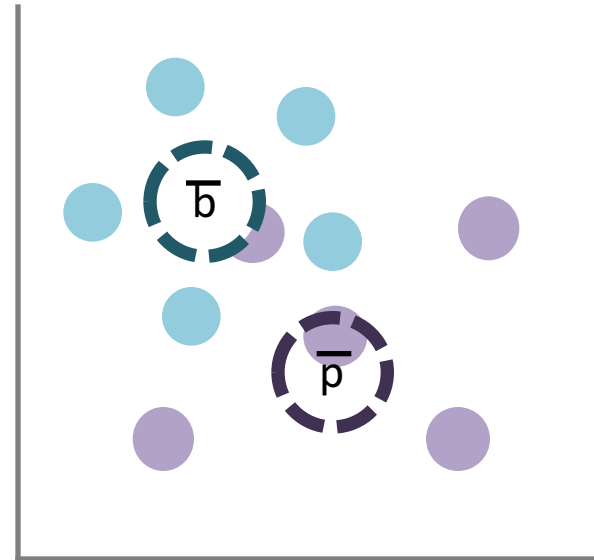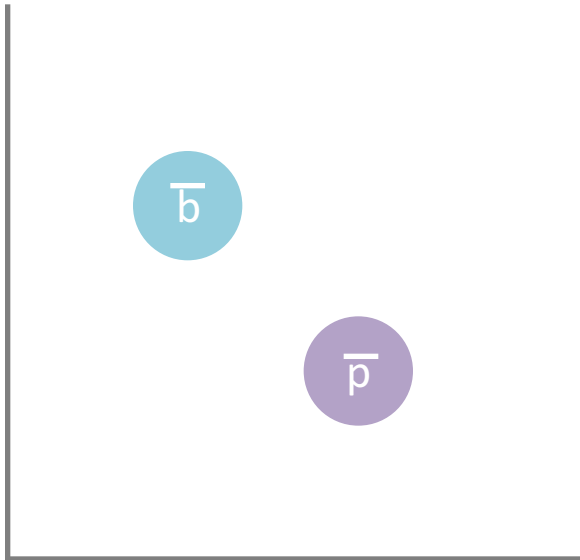
Use the visual system to estimate the answer

# Big Picture Analyses

Computational Aggregation:

Visual Aggregation:

Compute the answer then visualize it

Use the visual system to estimate the answer

# Encodings

# Tasks



Maxima

Minima

Range

Average

Variance

Outliers

X

© D.A. Szafir, 2016

What month has the **highest sales day**?

What month has the **highest sales on average**?

What month has the highest sales on average?

What month has the **highest sales day**?

Position for
Point Tasks

Color for
Summary Tasks

How you map the data impacts what information is
readily extracted

Can we design better visualization systems that do support these analyses?

# Two Challenges for Visualization

**Scalability**

How can we support insight across larger numbers and higher complexity?

**Comprehensibility**

How can we ensure estimates from a visualization are accurate?

# Visualization in the Age of Big Data

## Understand limits in current tools
### Large Scale Sequence Alignment

## Derive inspiration across domains
### Literary Patterns

## Link big and small
### Machine Learning & Molecules

# Visualization in the Age of Big Data

## Understand limits in current tools

What does the data look like?

## Derive inspiration across domains

Literary Patterns

## Link big and small

Machine Learning & Molecules

ACGTTT CGATGC TGCCTC AAGCTA CAACGA

Organism

Population

CGATGC ACGTTT TGCATA CAACGA CGATGC

GCATA ACGTTT CGATGC AAGCTA CGATGC

TGCCTC CAACGA ACGTTT AAGGAA CGATGC

CGATGC CAACGA CGATGC AAGCTA ACGTTT

Albers, Dewey, & Gleicher, 2011

ACGTTT CGATGC TGCCTC AAGCTA CAACGA

CGATGC ACGTTT TGCATA CAACGA CGATGC

TGCATA ACGTTT CGATGC AAGCTA CGATGC

TGCCTC CAACGA ACGTTT AAGGAA CGATGC

CGATGC CAACGA CGATGC AAGCTA ACGTTT

Albers, Dewey, & Gleicher, 2011

Gene One

Gene Two

Albers, Dewey, & Gleicher, 2011

Organism One

Organism Two

Albers, Dewey, & Gleicher, 2011

Darling et al, 2004

Limited Number of Sequences

Limited Length of Sequences

Only Reference-Dependent Analysis

Difficult to Analyze High-Level Relationships

# Visualization in the Age of Big Data

Understand limits in current tools

What does the data look like?
The Fix: Aligning patterns with tasks

Derive inspiration across domains

Literary Patterns

Link big and small

Machine Learning & Molecules

Gene One

Gene Two

Summarization

Pop-Out

Visual Search

Visual Clutter

Color better supports visual processing at scale

© D.A. Szafir, 2016

Summarization

Pop-Out

Visual Search

Visual Clutter

# Color better supports visual processing at scale

Summarization

Pop-Out

Visual Search

Visual Clutter

Color better supports visual processing at scale

Summarization

Pop-Out

Visual Search

Visual Clutter

# Color better supports visual processing at scale

Note: Due to limitations in the Census data, foreign-born populations are not available in all areas for all years.

Large
group

2%

Bubb

SecurityMax

"Sasser", "Blaster" and "MyDoom": Why Your Network Can't Stop Them
Internet Security Webinar

Monday, December 13th

Register Now

© D.A. Szafir, 2016

The New York Times

# Mapping America: Every City, Every Block

Browse local data from the Census Bureau's American Community Survey, based on samples from 2005 to 2009.

## Distribution of racial and ethnic groups

MAP KEY
One dot = 200 people
- White
- Black
- Hispanic
- Asian
- Other

Census tract 65
Whites: 2%
Blacks: 23%
Hispanics: 65%
Asians: 3%
Other groups: 5%

# Average Number of Genes per Genome



Often too many genes to display on the monitor

# Raw Sequence



# Sequence Blocks



# Aggregate Representation

# Sequence Block

# Average

# Sequence Block



# Robust Average

# Sequence Block



# Event Striping

# Sequence Block



# Color Weaving

Average

Robust Average

Event Striping

Color Weaving

# Visualization in the Age of Big Data

## Understand limits in current tools
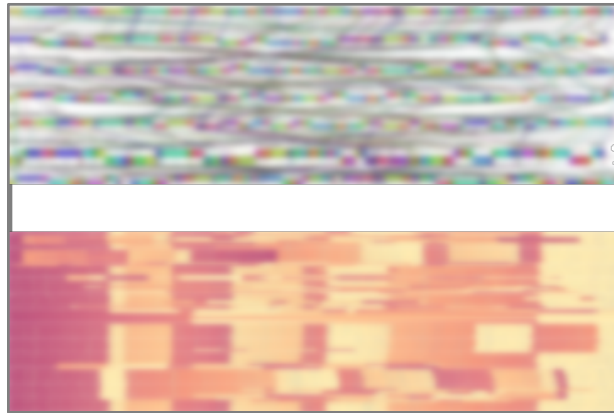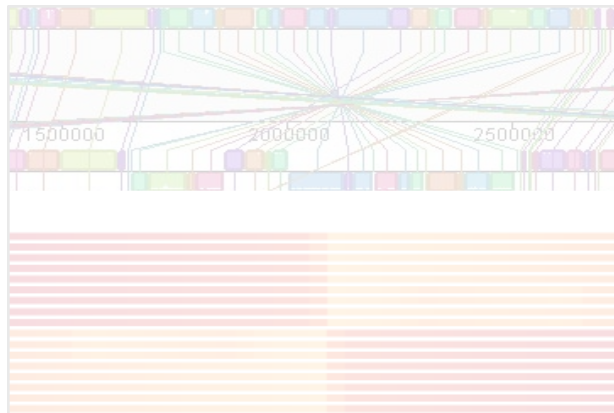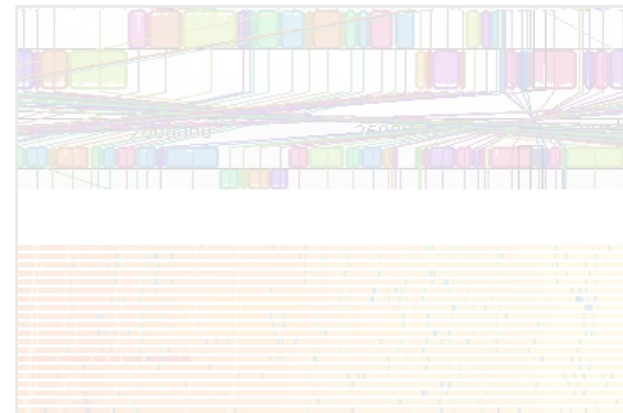
What does the data look like?
The Fix: Aligning patterns with tasks
Building The System

## Derive inspiration across domains

Literary Patterns

## Link big and small

Machine Learning & Molecules

# Sequence Surveyor

Task-Driven Sequence Aggregation
Length of Sequences

Perceptually-Driven Encoding
Number of Sequences

Dynamic Compression, Color and
Position Encoding Choices
Variety of Tasks

|  | Index | Membership Freq | Grouped Freq | Pos in Reference |
|---|---|---|---|---|
| Index | | | | |
| Grouped Freq | | | | |
| Pos in Reference | | | | |

© D.A. Szafir, 2016

10x More Sequences

100x Longer Sequences

Reference-Dependent, Independent, and Metadata-Based Analyses

Explicit Support for High-Level and Low-Level Relationships

100 Bacteria
6,000 genes

50 Bacteria
5,000 genes

35 Fungi
17,000 genes

14 Pathogens
4,000 genes

8 Partial *E. Coli*
300 genes

# Explore Evolutionary Patterns in Organisms

# Explore Phylogenetic Relationships

# Explore Phylogenetic Relationships

# Explore Phylogenetic Relationships

# "At a Glance" Algorithm Debugging

Color-based aggregation better supports analyses at scale

# Visualization in the Age of Big Data

Understand limits in current tools
### Large Scale Sequence Alignment

Derive inspiration across domains
### Literary Patterns

Link big and small
### Machine Learning & Molecules

All the world's a stage,
And all the men and women merely players:
They have their exits and their entrances;
And one man in his time plays many parts,
His acts being seven ages. At first the infant,
Mewling and puking in the nurse's arms.
And then the whining school-boy, with his satchel
And shining morning face, creeping like snail
Unwillingly to school. And then the lover,
Sighing like furnace, with a woeful ballad
Made to his mistress' eyebrow. Then a soldier,
Full of strange oaths and bearded like the pard,
Jealous in honour, sudden and quick in quarrel,
Seeking the bubble reputation
Even in the cannon's mouth.

# Large Digitized Collections

*Google N-Grams: 5,195,769 books*

More Effective

Less Effective

Position

Length/Height

Orientation

Area

Value/Lightness

Color

Cleveland & McGill, 1985

© D.A. Szafir, 2016

⚠ Please don't use wordclouds

More Effective

Less Effective

Position

Length/Height

Orientation

Area

Value/Lightness

Color

Cleveland & McGill, 1985

© D.A. Szafir, 2016

# Word Usage Analysis Tasks

Characterize and compare authors

Measure shifts in an author's writing over time

Evolution of language over time

Evolution of cultural influences over time

Indicate recurring themes and topics

Characterize typographic conventions

# Word Usage Analysis Tasks

Characterize and compare  organisms

Measure shifts in  organisms  over  species

Evolution of organisms over time

Evolution of cultural influences over time

Indicate recurring  genetic material

Characterize typographic conventions

# Turning texts into sequences

```
All the world's a stage,
And all the men and women merely players:
They have their exits and their entrances,
```

```
all the world a stage
and all the men and women merely players
they have their exits and their entrances
```

| all | the | world | a | stage | and | all | the | men | and | women | merel |
|---|---|---|---|---|---|---|---|---|---|---|---|

# Text Sequence:

Present words in their
original reading order

Highlight word locations

Precise analysis for single texts

### A Midsummer Night's Dream

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| *Text Sequence:* | now | fair | Hippolyta | our | nuptial | hour | draws | on | apace | four |
| *Ranked Count:* | the | and | to | I | you | of | a | in | my | is |
| *Position* | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

# Text Sequence:
Present words in their
original reading order

Highlight word locations

Precise analysis for single
texts

# Ranked Count:
Order words by how often
they occur in a text collection

Highlight word frequency

Aggregate multiple texts

| A Midsummer Night's Dream | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| *Text Sequence:* | now | fair | Hippolyta | our | nuptial | hour | draws | on | apace | four |
| *Ranked Count:* | the | and | to | I | you | of | a | in | my | is |
| *Position* | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

| | | | | | |
|---|---|---|---|---|---|
| **King Henry IV pt. 1** | the | and | i | of | a |
| King Henry IV pt. 2 | the | and | i | of | to |
| King Henry VI pt. 1 | and | the | of | to | i |
| King Henry VI pt. 2 | the | and | to | i | of |
| King Henry VI pt. 3 | and | the | to | i | my |

© D.A. Szafir, 2016

King Henry IV pt. 1

King Henry IV pt. 2

King Henry VI pt. 1

King Henry VI pt. 2

King Henry VI pt. 3

© D.A. Szafir, 2016

5.2 million books from 1660-2009

Mitchel et al, 2011

*175,000 words over 35 decades*
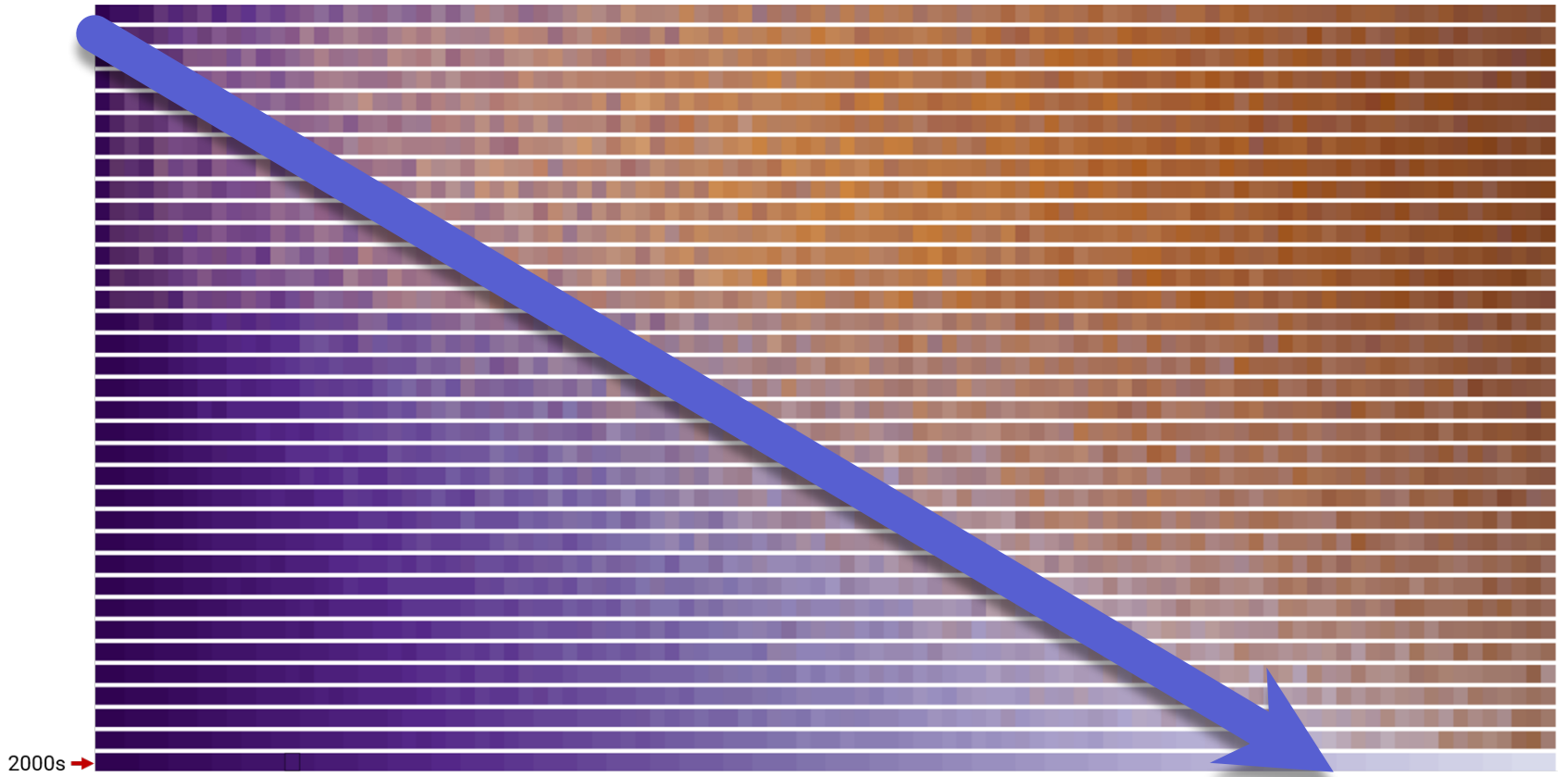
# Explore Evolutionary Patterns in Writing

Time

Popularity Rank
(High to Low)

2000s

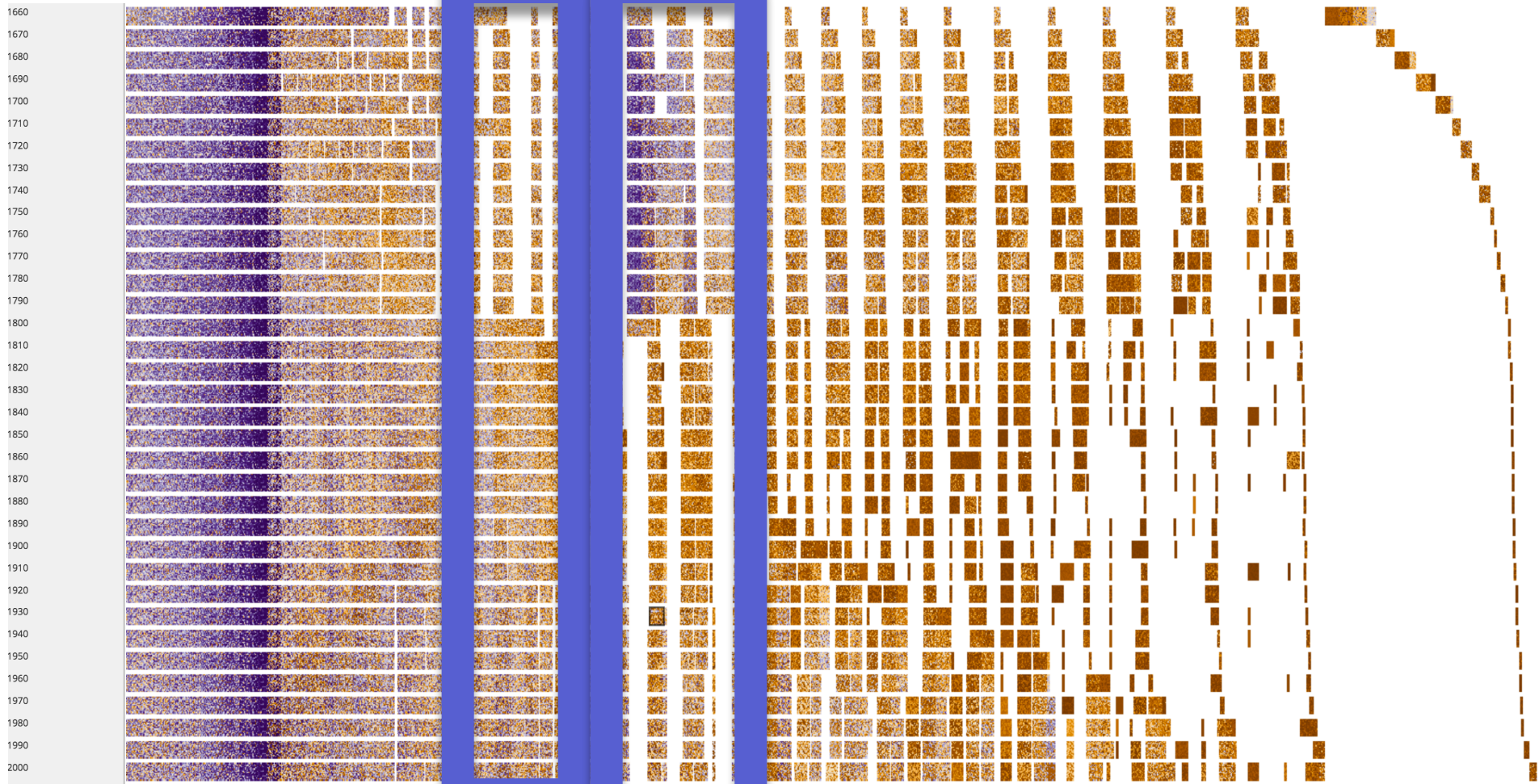© D.A. Szafir, 2016

2000s

# Confirm Prior Hypotheses

© D.A. Szafir, 2016

| Sequences | | |
|---|---|---|
| | ■ exercife | rank : 935, frequency : 15 |
| 1660 | ■ fi | rank : 1877, frequency : 15 |
| 1670 | ■ fecond | rank : 328, frequency : 15 |
| 1680 | ■ defire | rank : 600, frequency : 15 |
| 1690 | ■ pafs | rank : 556, frequency : 15 |
| 1700 | ■ unlefs | rank : 826, frequency : 15 |
| 1710 | ■ obfervation | rank : 1116, frequency : 15 |
| 1720 | ■ pleafure | rank : 415, frequency : 15 |
| 1730 | ■ neceflary | rank : 353, frequency : 15 |
| 1740 | ■ ferve | rank : 894, frequency : 15 |
| 1750 | ■ faw | rank : 484, frequency : 15 |
| 1760 | ■ fent | rank : 256, frequency : 15 |
| 1770 | ■ obferve | rank : 695, frequency : 15 |
| **1780** | ■ fupply | rank : 851, frequency : 15 |
| 1790 | ■ reafon | rank : 202, frequency : 15 |
| 1800 | ■ truft | rank : 1052, frequency : 15 |
| 1810 | ■ raifed | rank : 575, frequency : 15 |
| 1820 | ■ confent | rank : 1151, frequency : 15 |
| 1830 | ■ fuffer | rank : 981, frequency : 15 |
| | ■ folid | rank : 1766, frequency : 15 |

# Identify Cultural Shifts

Wife

Women

1910 – 1919

© D.A. Szafir, 2016

# The Plays of William Shakespeare

*961,304 words over 36 plays*

# Author Attribution

Look for inspiration in other data domains

Henry VI pt 1
Julius Caesar
King John
Titus Andronicus

# Visualization in the Age of Big Data
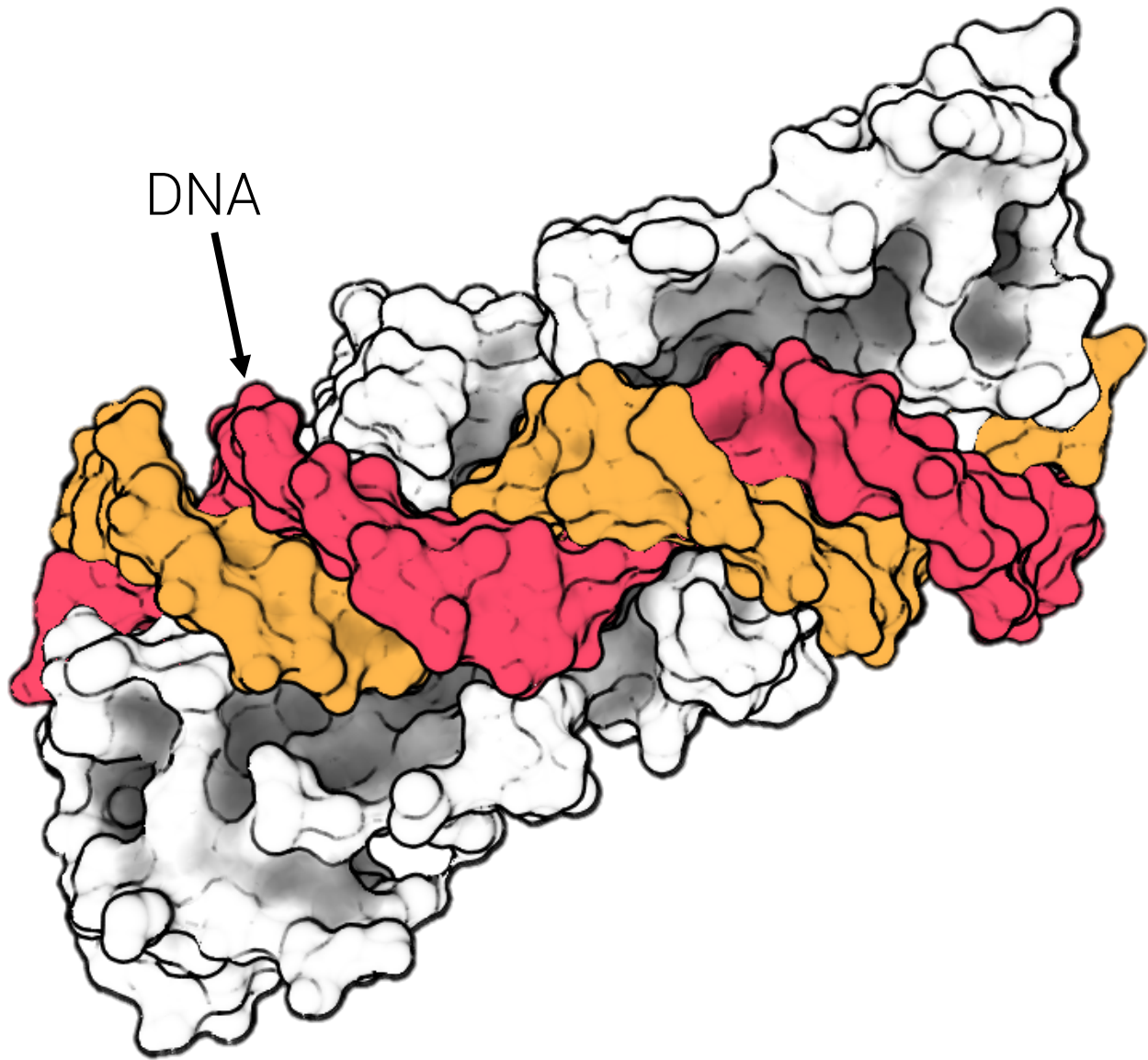
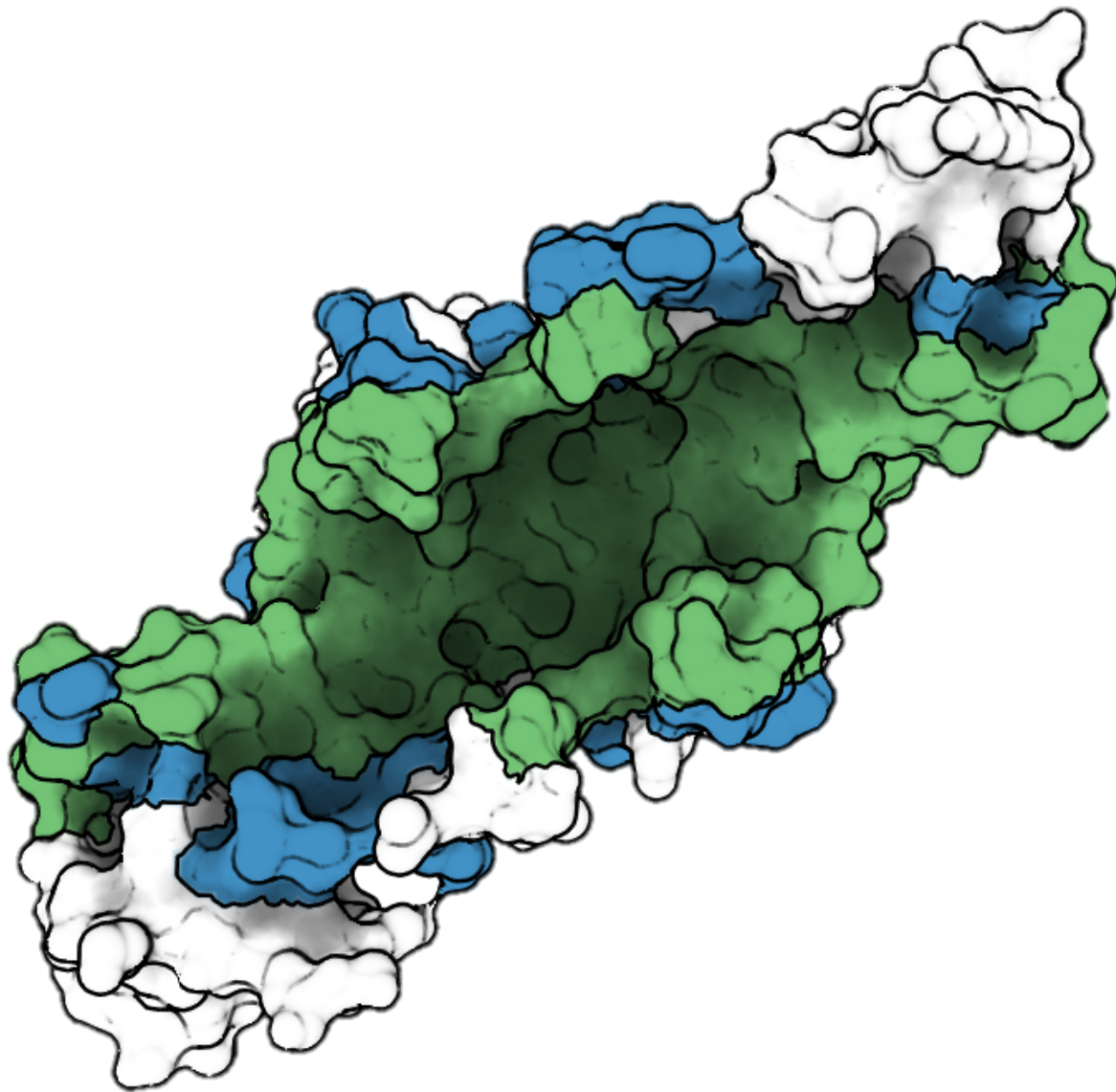Understand limits in current tools

Large Scale Sequence Alignment

Derive inspiration across domains

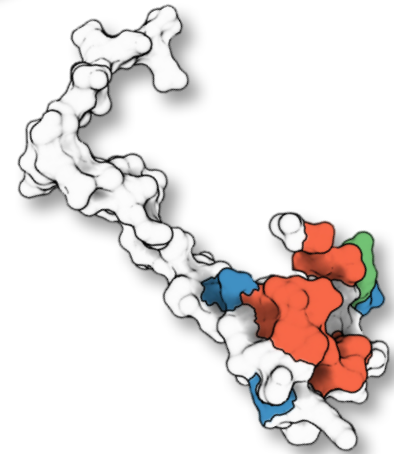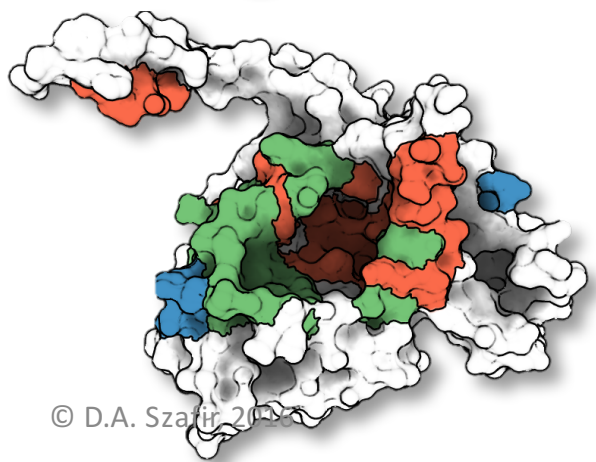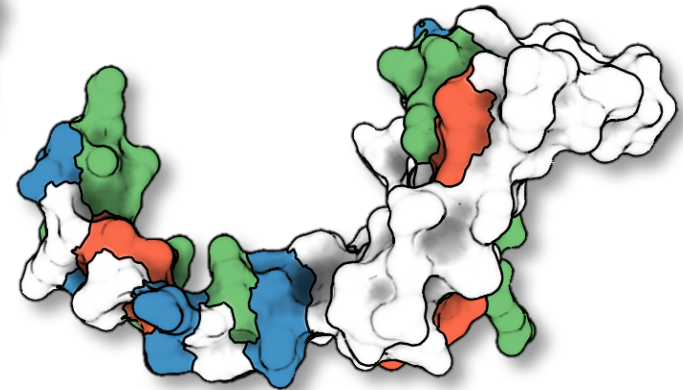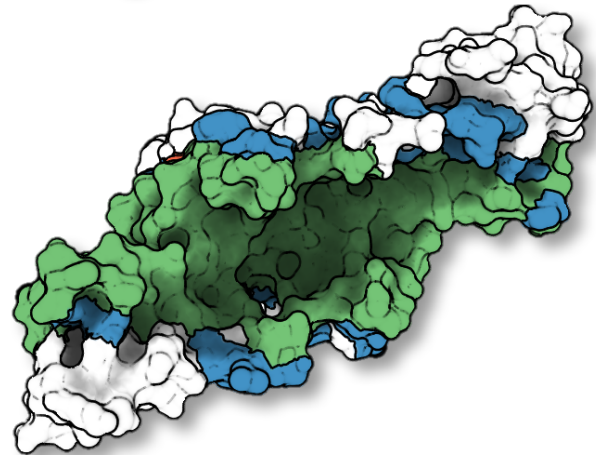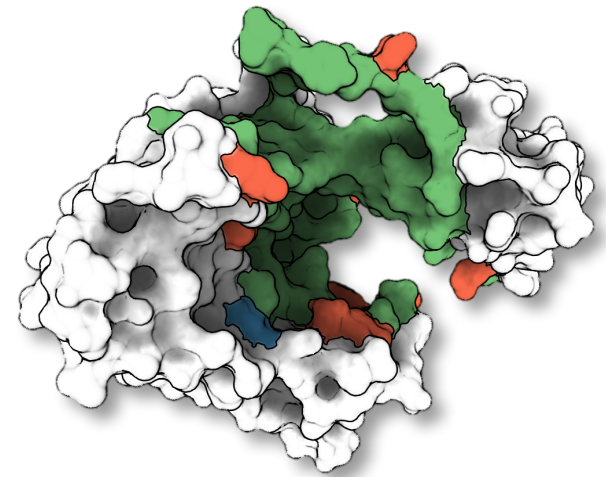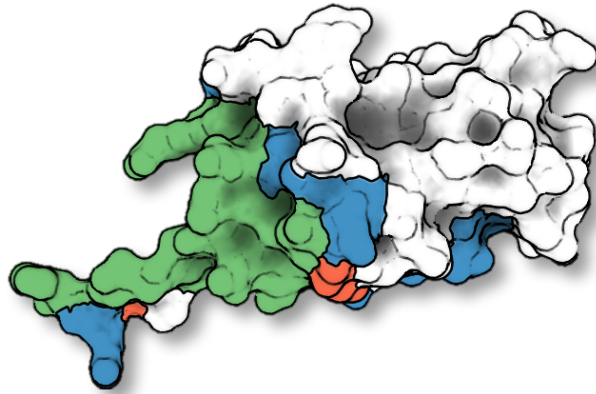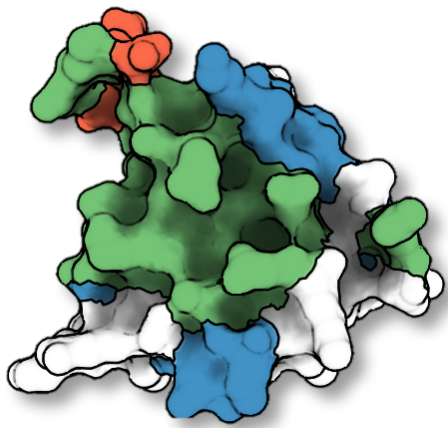Literary Patterns

## Link big and small

Machine Learning & Molecules

DNA
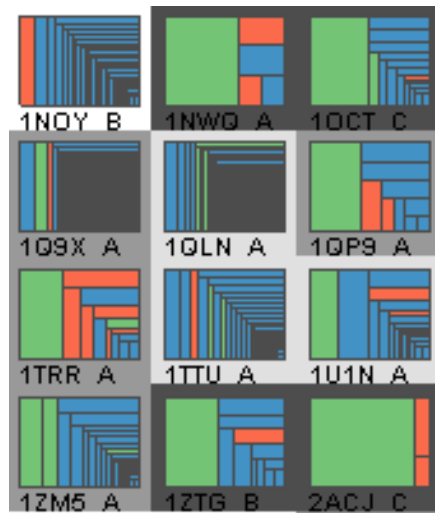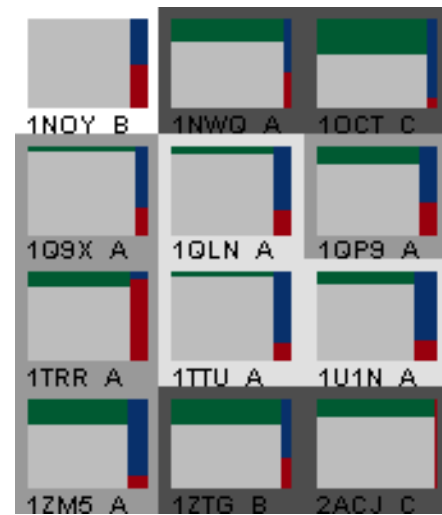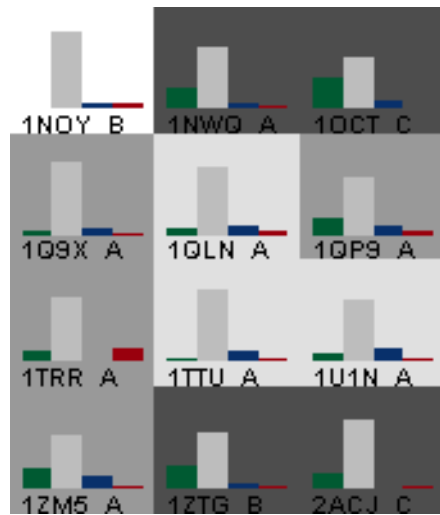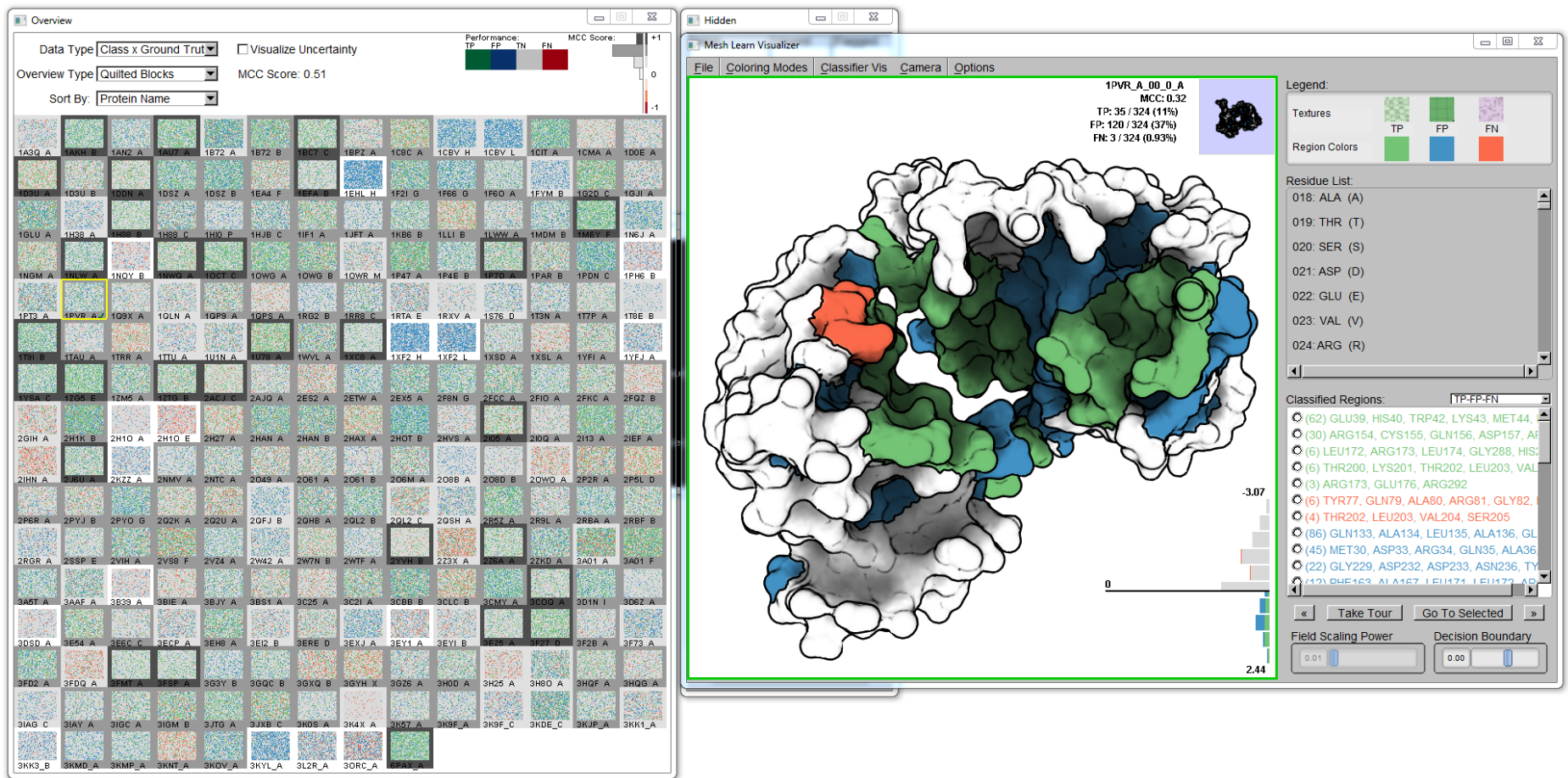
Hundreds of proteins with binding site predictions computed over hundreds of ligands
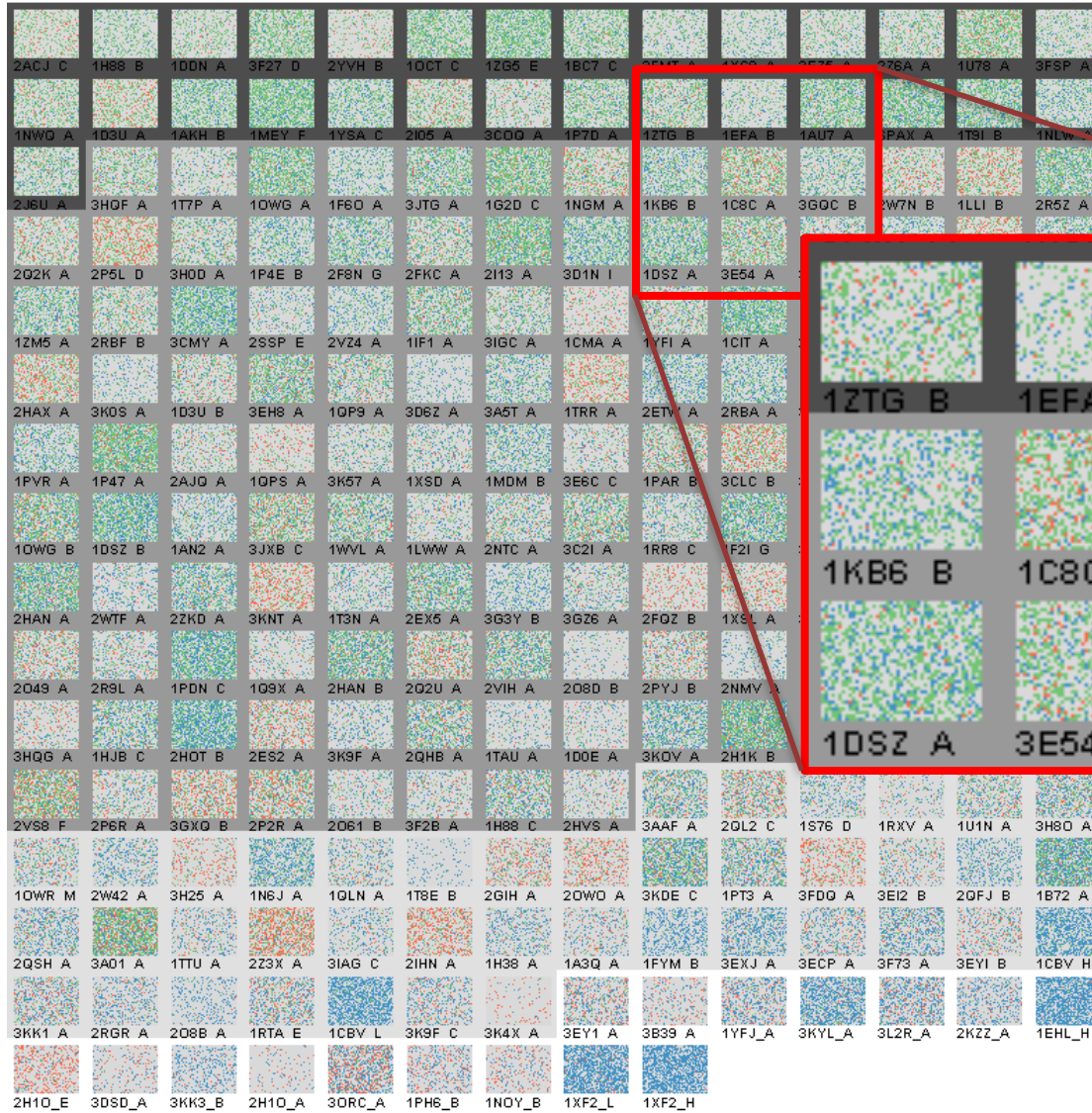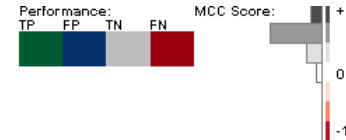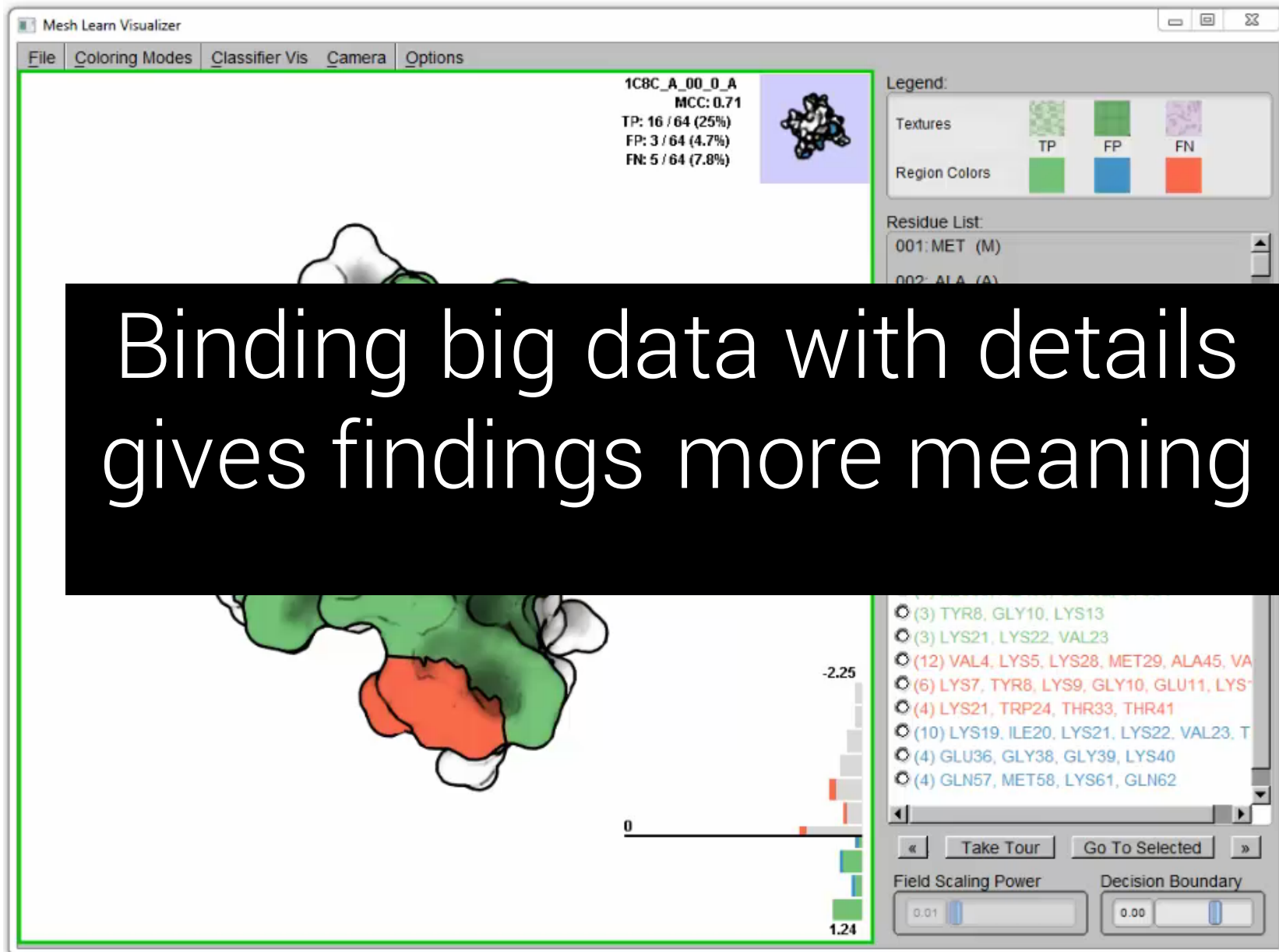
Task-driven overviews of large-scale machine learning performance data

DNA binding predictions over 216 proteins with 40 to 800 residues per protein

Binding big data with details gives findings more meaning

# Visualization in the Age of Big Data

## Understand limits in current tools

Large Scale Sequence Alignment

## Derive inspiration across domains

Literary Patterns

## Link big and small

Machine Learning & Molecules

# Designing for Big Data

Consider how the ways we communicate data support high-level tasks.

Look at parallels in the data structure and tasks associated with your data.

Don't lose sight of the details.

# Thank You!

Danielle Albers Szafir
danielle.szafir@colorado.edu
@dalbersszafir

Demos & Papers at:
http://danielleszafir.com

# Extra Slides