

The Big Data to Knowledge (BD2K) Guide to the Fundamentals of Data Science

SECTION 4: DATA VISUALIZATION TOOLS & COMMUNICATION

NILS GEHLENborg
Assistant Professor
HARVARD MEDICAL SCHOOL
MARCH 31, 2017

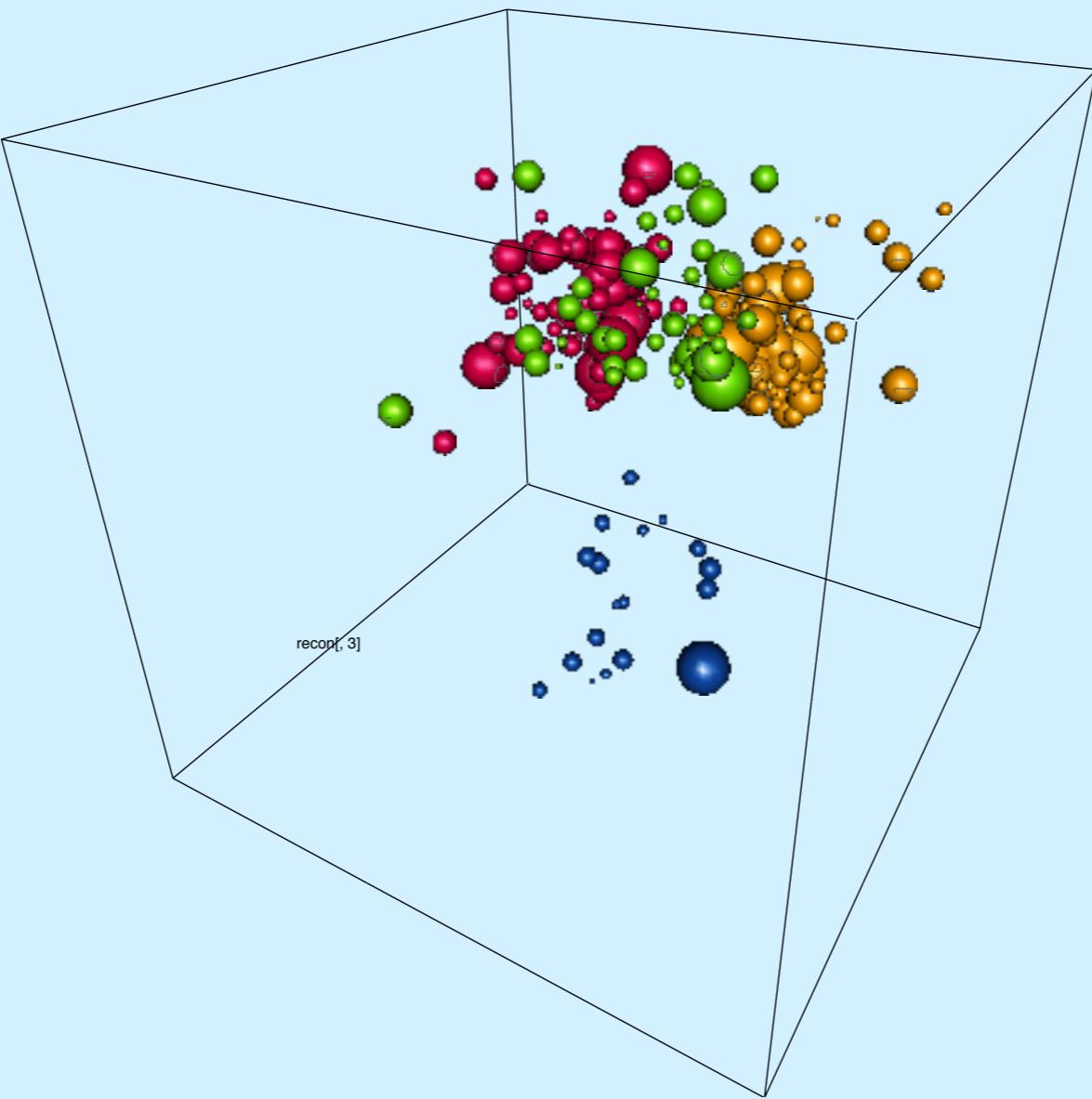


NILS GEHENBORG



- Assistant Professor in the Department of Biomedical Informatics at Harvard Medical School.
- Research & develops novel tools to visualize 3D genome conformation data as well as heterogeneous data from large-scale cancer genomics studies.
- Co-Investigator for the 4D Nucleome Network Data Coordination & Integration Center hosted at Harvard Medical School.
- Co-founder, former general chair, & current steering committee chair of BioVis, the Symposium on Biological Data Visualization
- Co-founder of VIZBI, the annual workshop on Visualizing Biological Data.
- Co-chairs the Policy Working Group for the 4D Nucleome Network, an NIH Common Fund project.
- Served on the program committees of several international bioinformatics & data visualization conferences & held multiple editorial roles, including his current role as associate editor of BMC Bioinformatics.
- Contributed to the “Points of View” data visualization column in Nature Methods.

WARNING!



Principal Components Analysis 3D Interactive Visualization

Overview

Use Cases for Data Visualization

Encoding Data for Visual Representation

Common Mistakes

Visualization for Communication

Use Cases for Data Visualization

How does data visualization work?

Visualization uses perception to free up cognition.

How does data visualization work?

MALWMRLPLALLALW
GDPAAAFVNQHLCGSHL
VEALYLVCGERGFFYTPKT
RREAEDLQVGQVELGGGP
GAGSLQPLALE GSLQKRG
VEQCCTSICSLYQLENYC N

How does data visualization work?

MALWMRLLPLLALLALW
GPDPAAAAFVNQHLCGSHL
VEALYLVCGERGFFYTPKT
RREAEDLQVGQVELGGGP
GAGSLQPLALEGSLQKRGI
VEQCCTSICSLYQLENYCN

How does data visualization work?

Visualization uses perception to free up cognition.

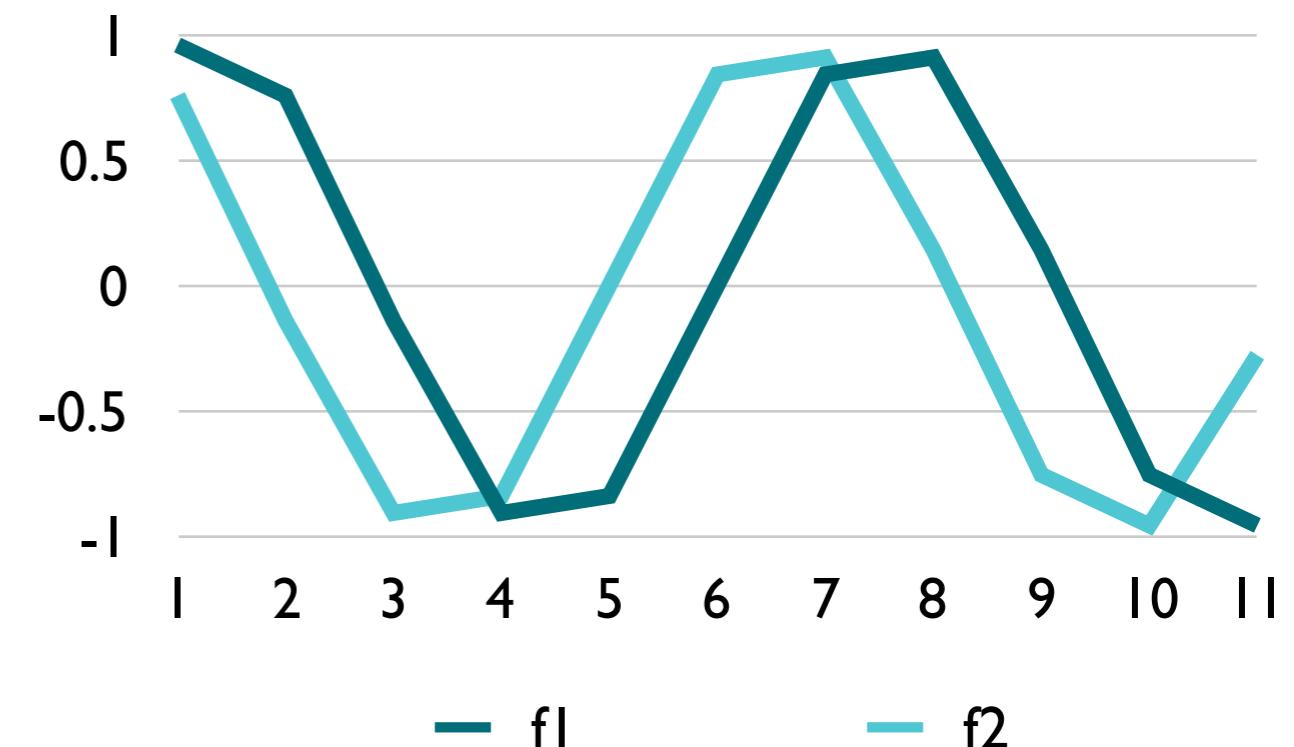
Visualization is an external cognitive aid
and augments working memory.

How does data visualization work?

$$\begin{array}{r} \mathbf{453 \times 862 =} \\ \hline & 906 \\ & + 27,180 \\ & + 362,400 \\ \hline & 390,486 \end{array}$$

How does data visualization work?

x	f1	f2
1.00	0.96	0.76
2.00	0.76	-0.14
3.00	-0.14	-0.91
4.00	-0.91	-0.84
5.00	-0.84	0.00
6.00	0.00	0.84
7.00	0.84	0.91
8.00	0.91	0.14
9.00	0.14	-0.76
10.00	-0.76	-0.96
11.00	-0.96	-0.28



How does data visualization work?

Visualization uses perception to free up cognition.

Visualization is an external cognitive aid
and augments working memory.

Why data visualization?

I believe it when I see it.

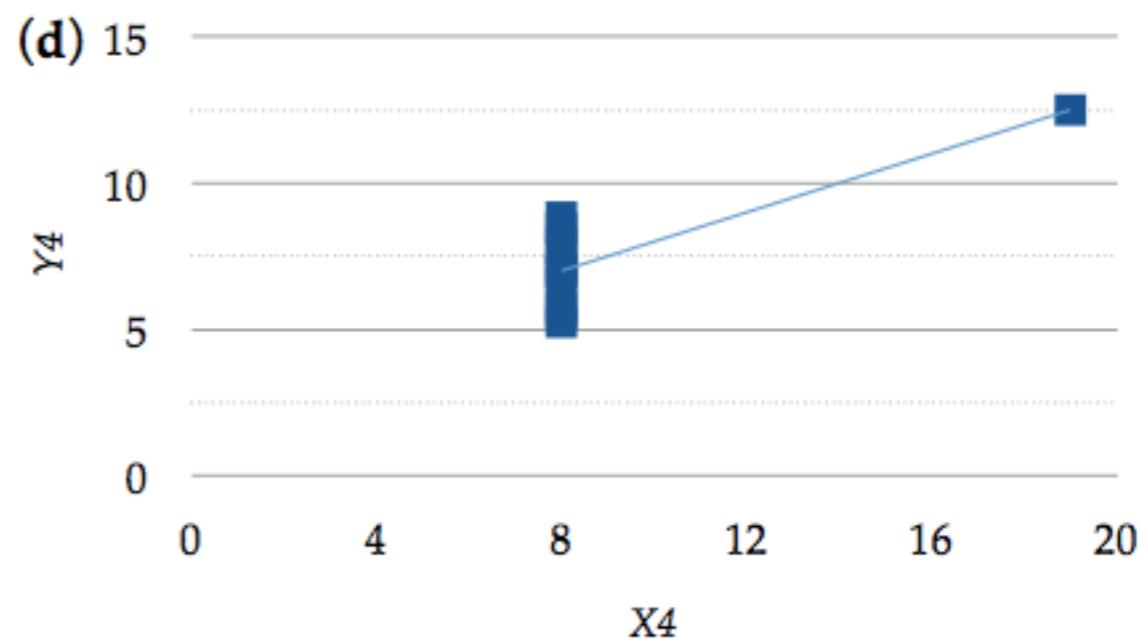
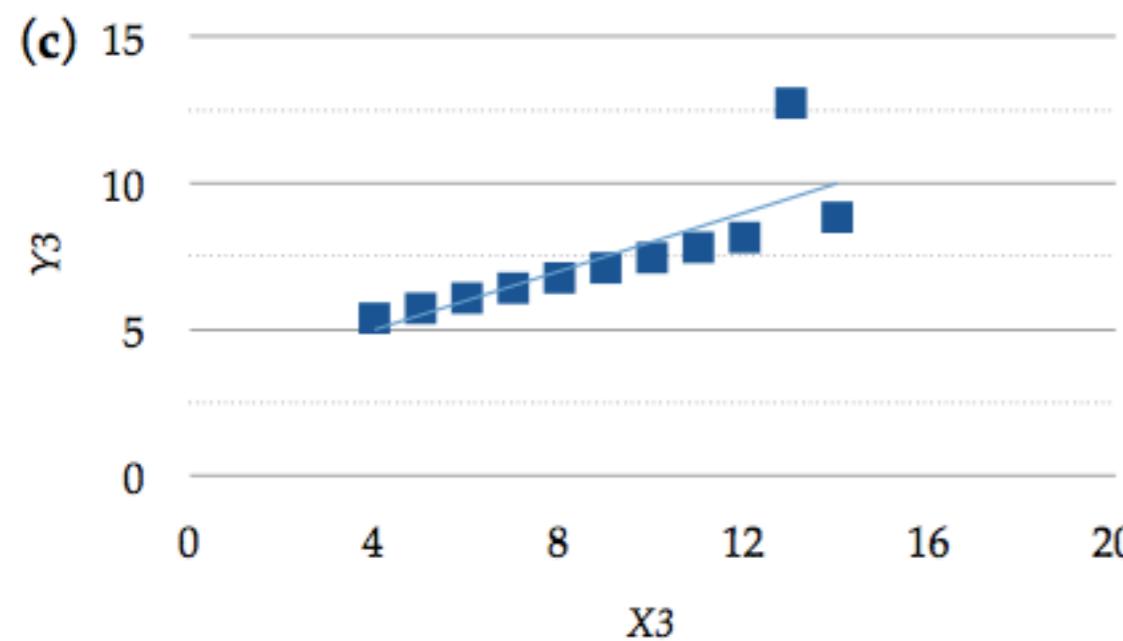
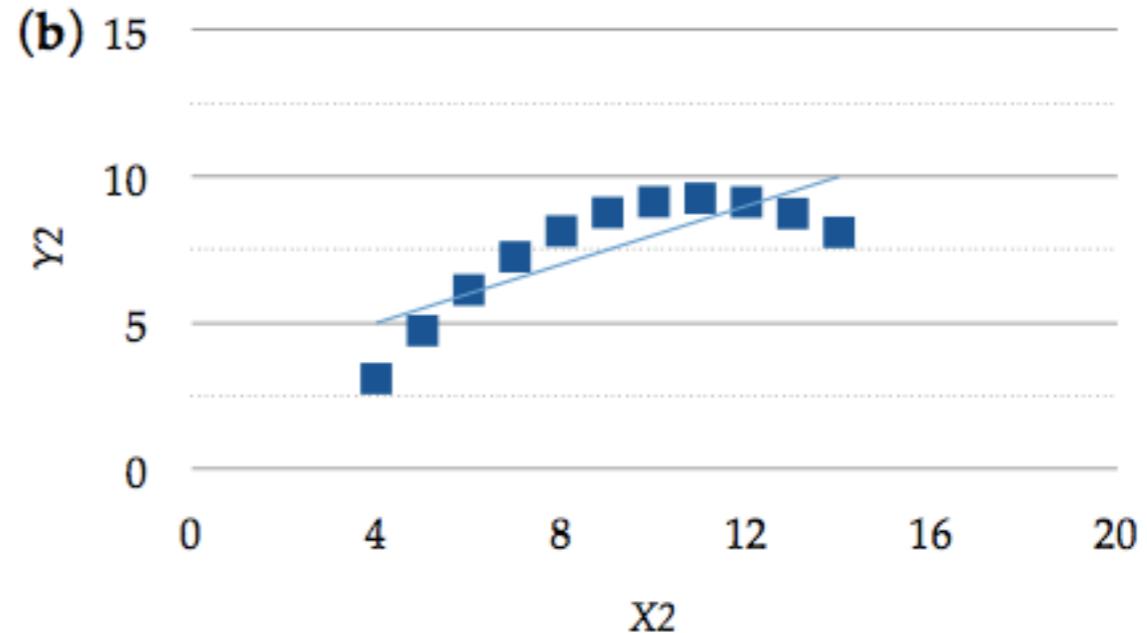
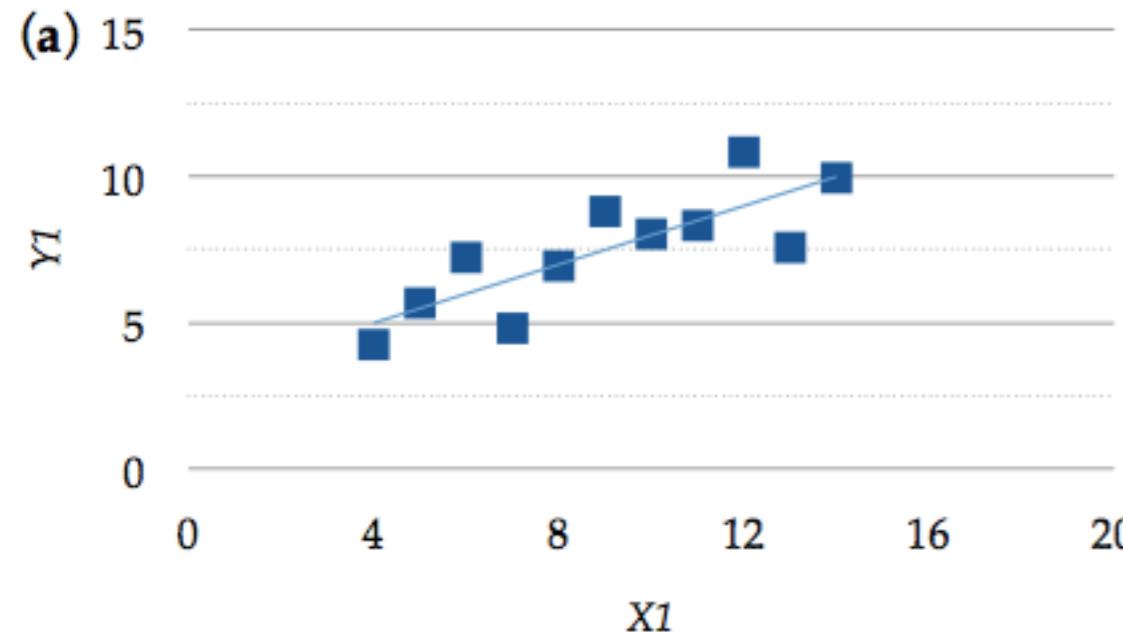
— *Unknown*

Anscombe's Quartet

X_1	Y_1	X_2	Y_2	X_3	Y_3	X_4	Y_4
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89

$\text{mean}(X) = 9$, $\text{var}(X) = 11$, $\text{mean}(Y) = 7.5$, $\text{var}(Y) = 4.12$,
 $\text{cor}(X,Y) = 0.816$, linear regression line $Y = 3 + 0.5*X$

Anscombe's Quartet

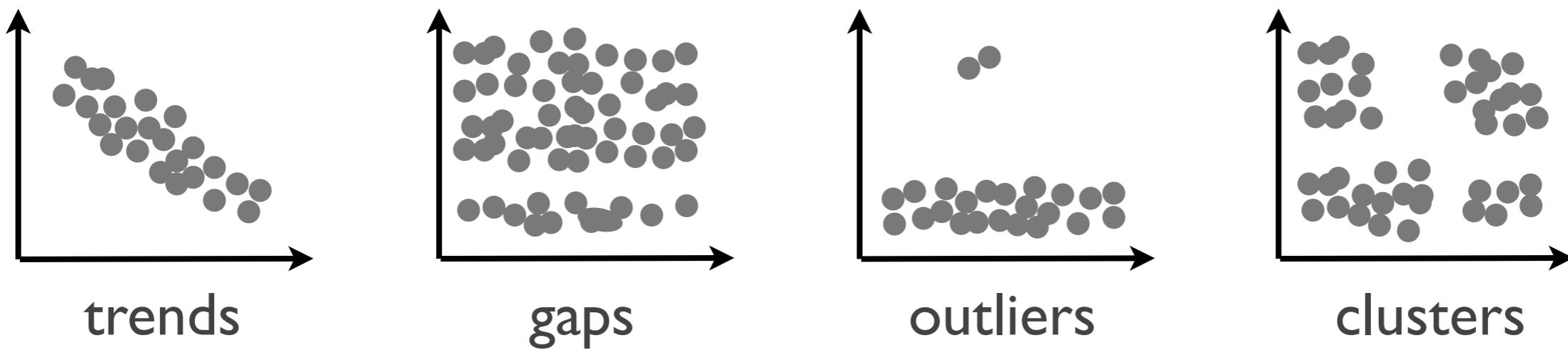


Why data visualization?

I'm wondering if there are any interesting patterns in my data.

— *Almost Everyone*

Exploration: Hypothesis Generation

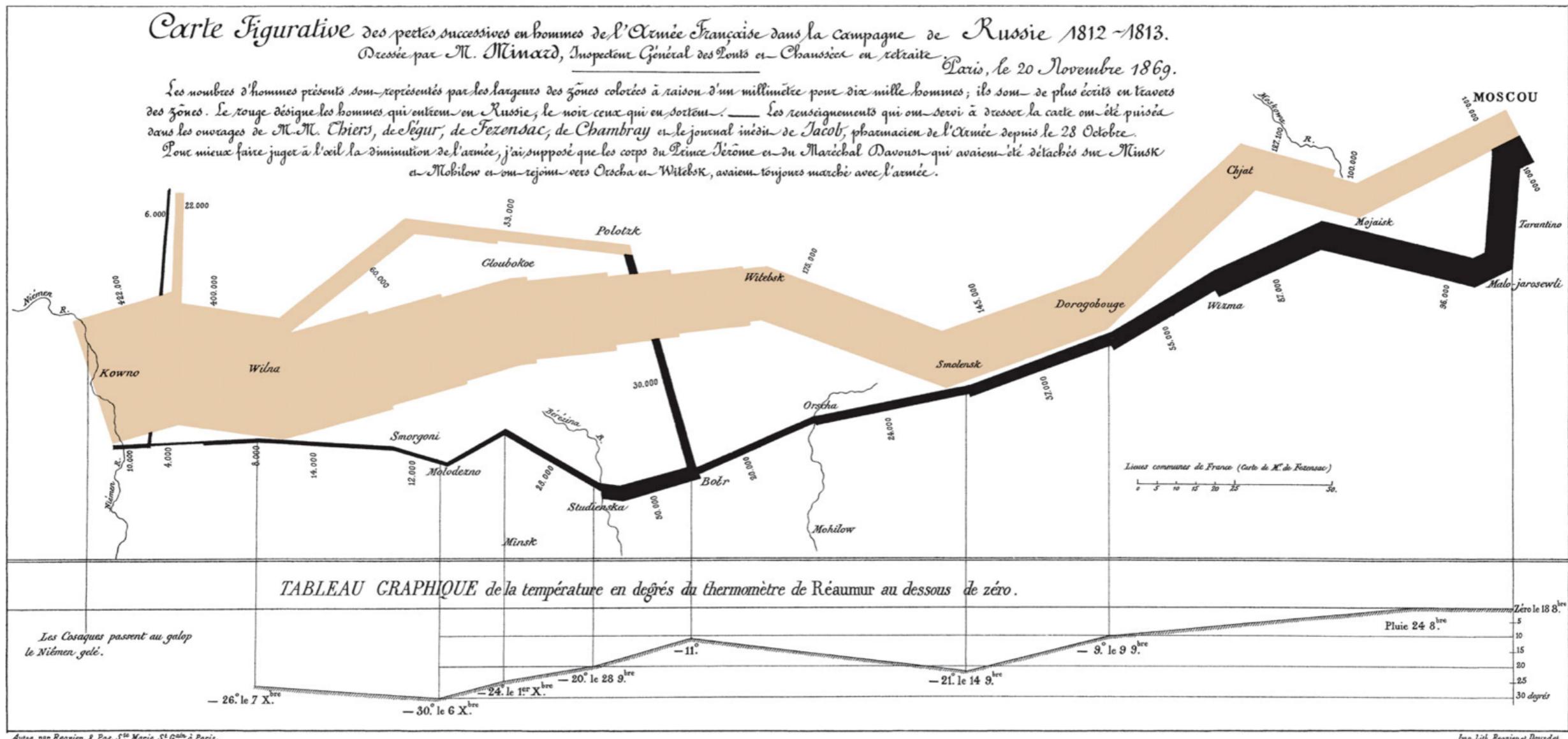


Why data visualization?

A good sketch is better than a long speech.

— Napoleon Bonaparte

Napoleon's March on Moscow



Visualization Use Cases

Confirmation

Exploration

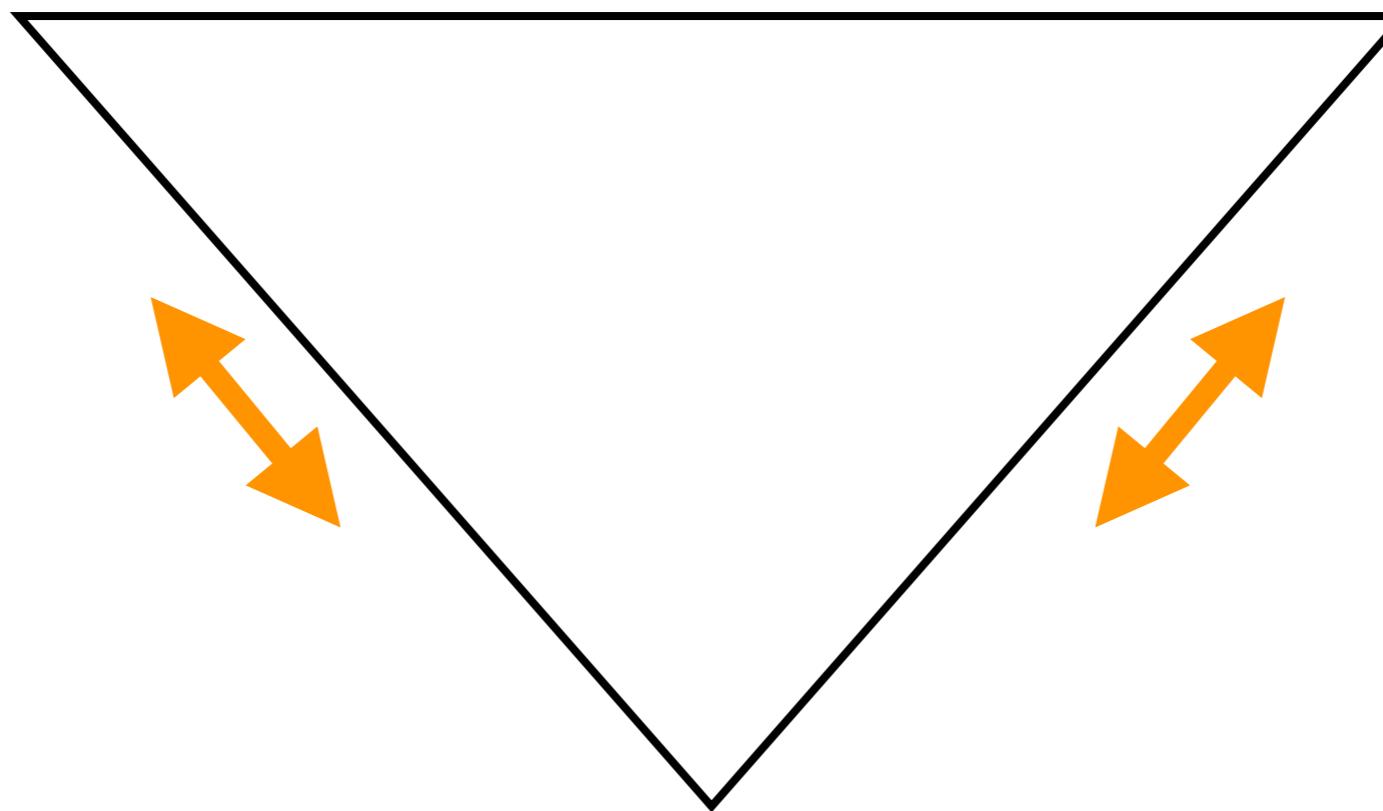
Presentation

dialogue between computer & analyst

Confirmation



Exploration



Presentation

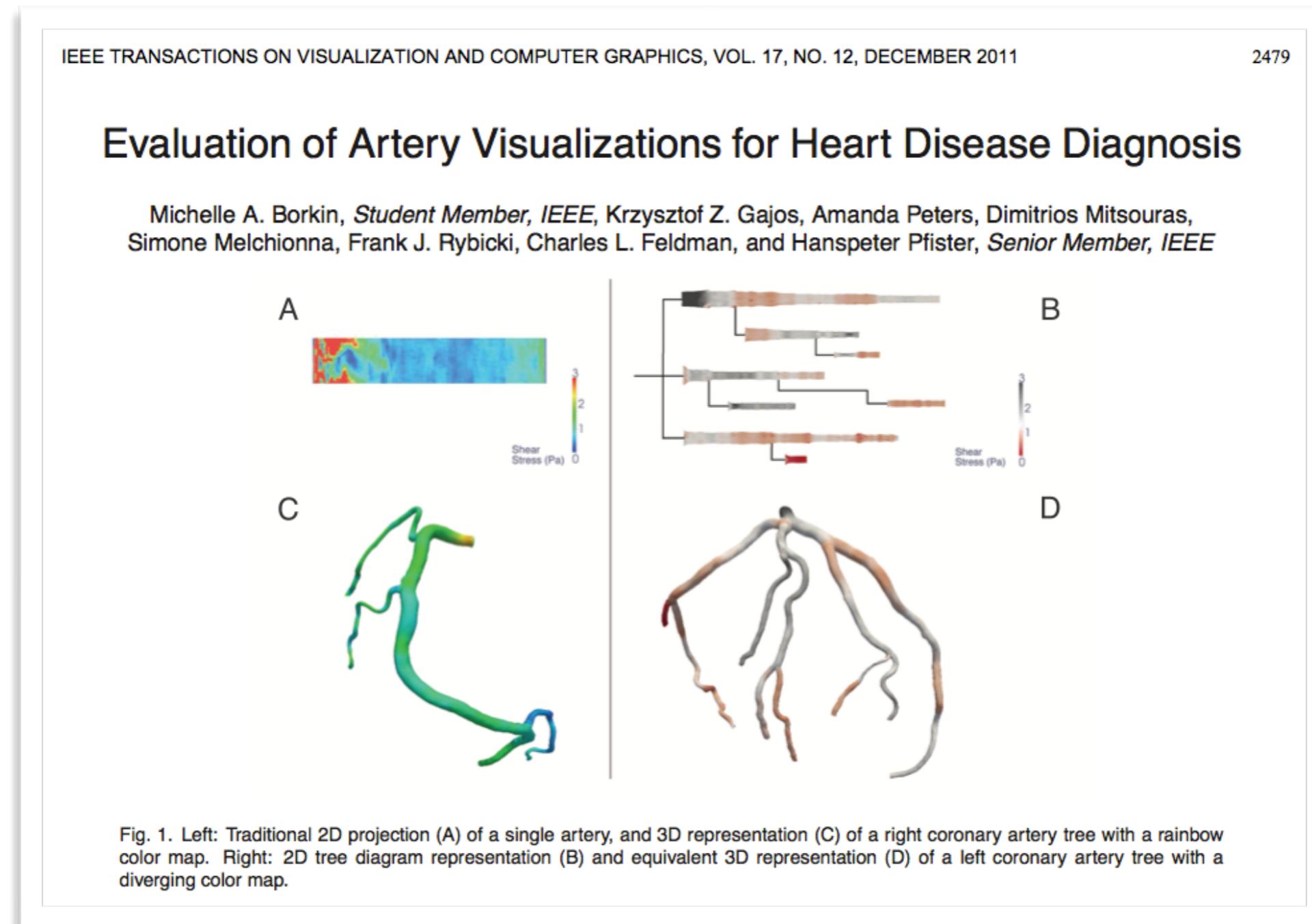
dialogue between analyst & audience

Encoding Data for Visual Representation

When good visualization really matters

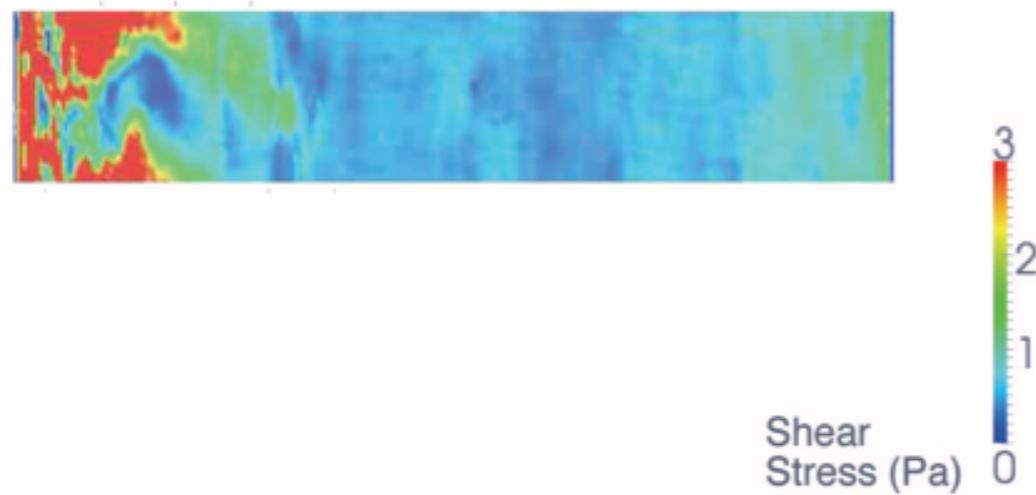


Why we should ❤️ good visualizations

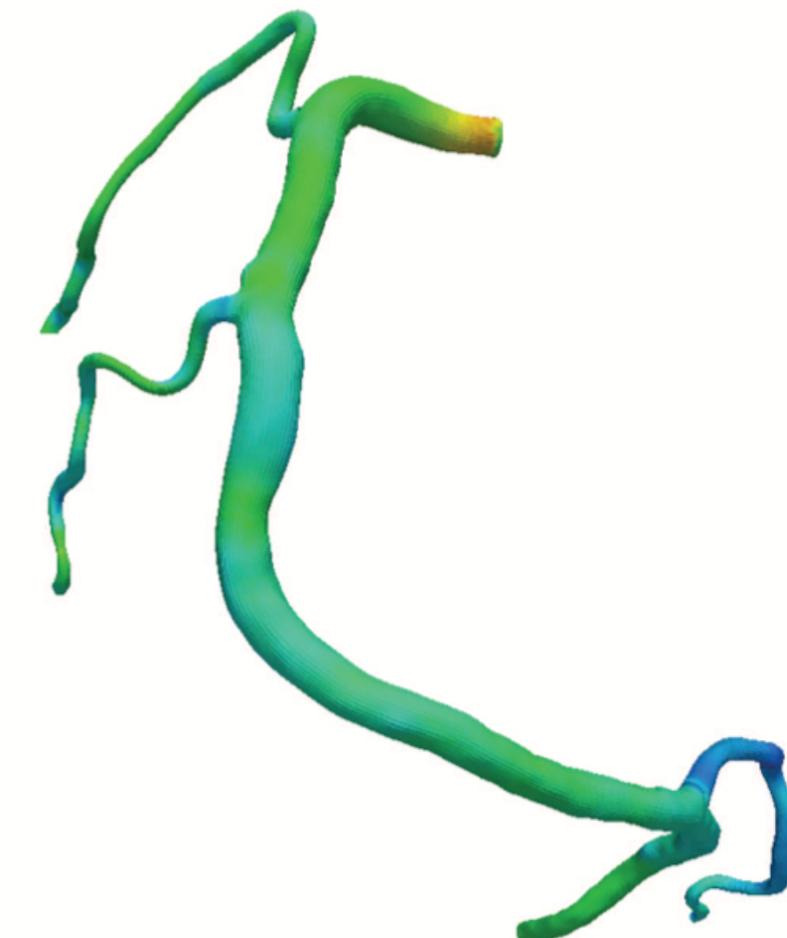


Why we should ❤️ good visualizations

Traditional 2D projection
of a single artery

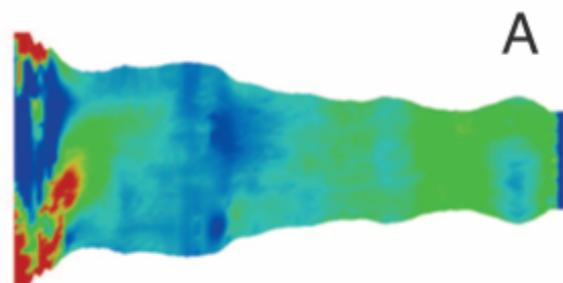


Traditional 3D projection
of right artery tree

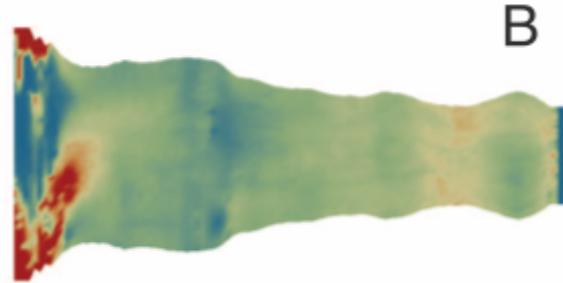


Color Maps

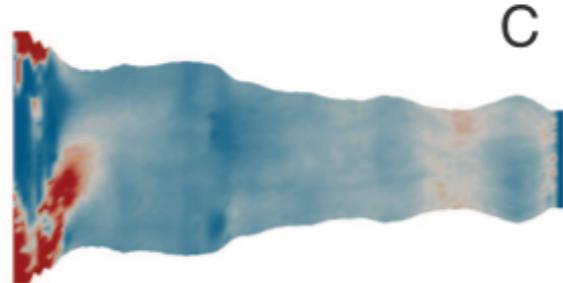
1st choice



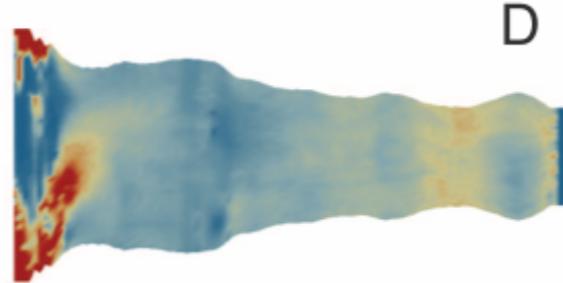
A



B



C



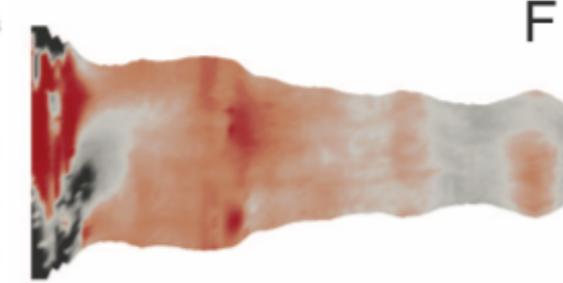
D



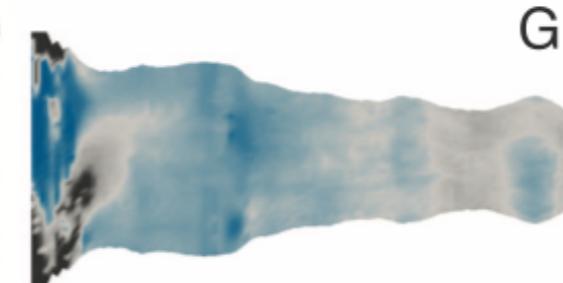
2nd choice



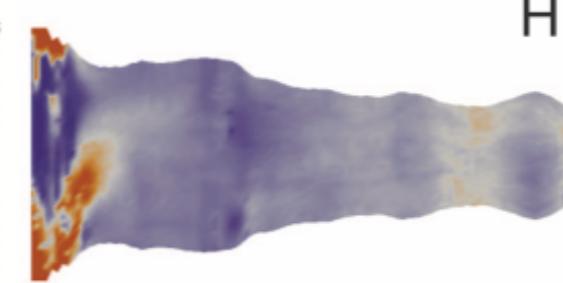
E



F



G

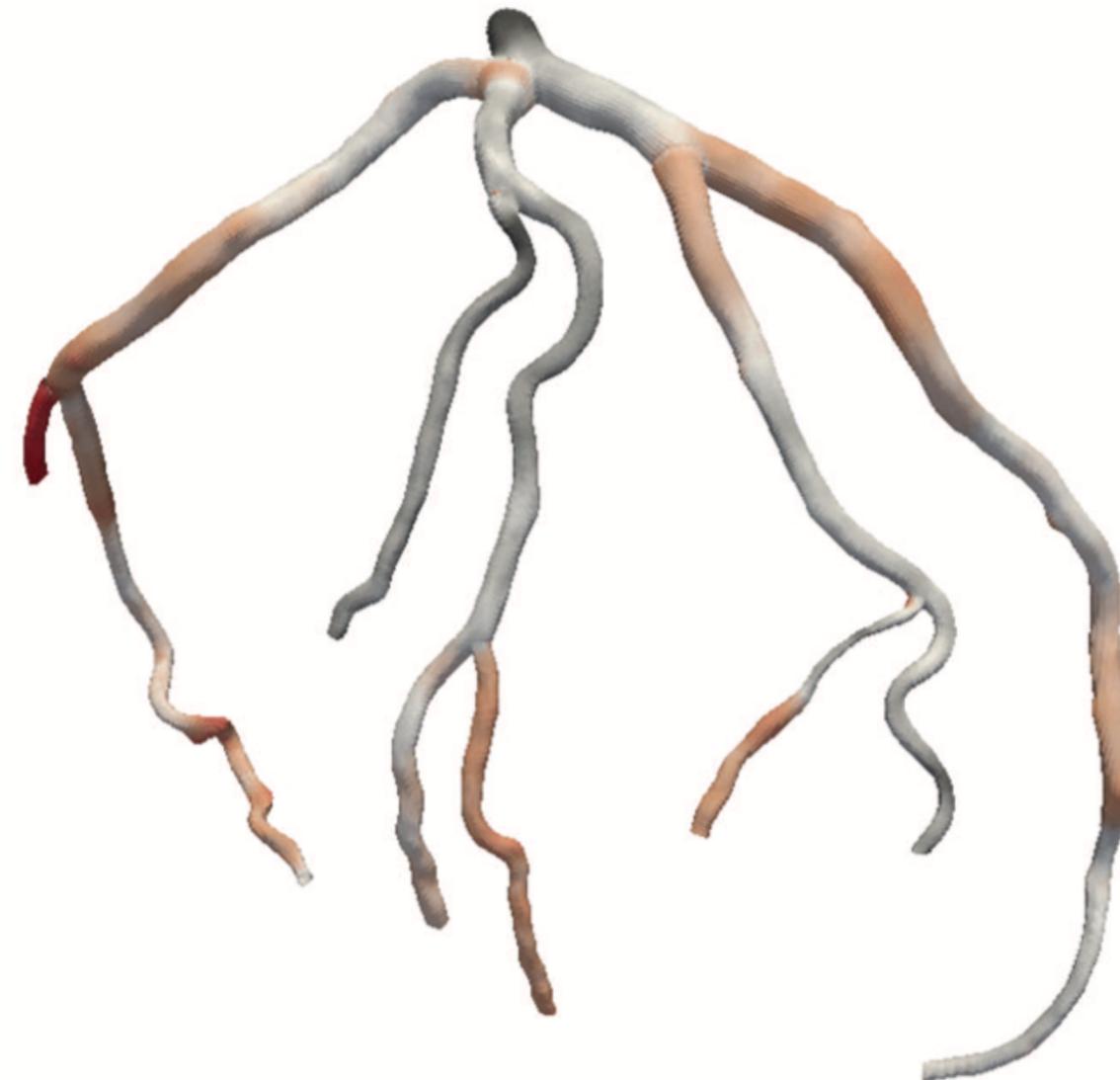


H



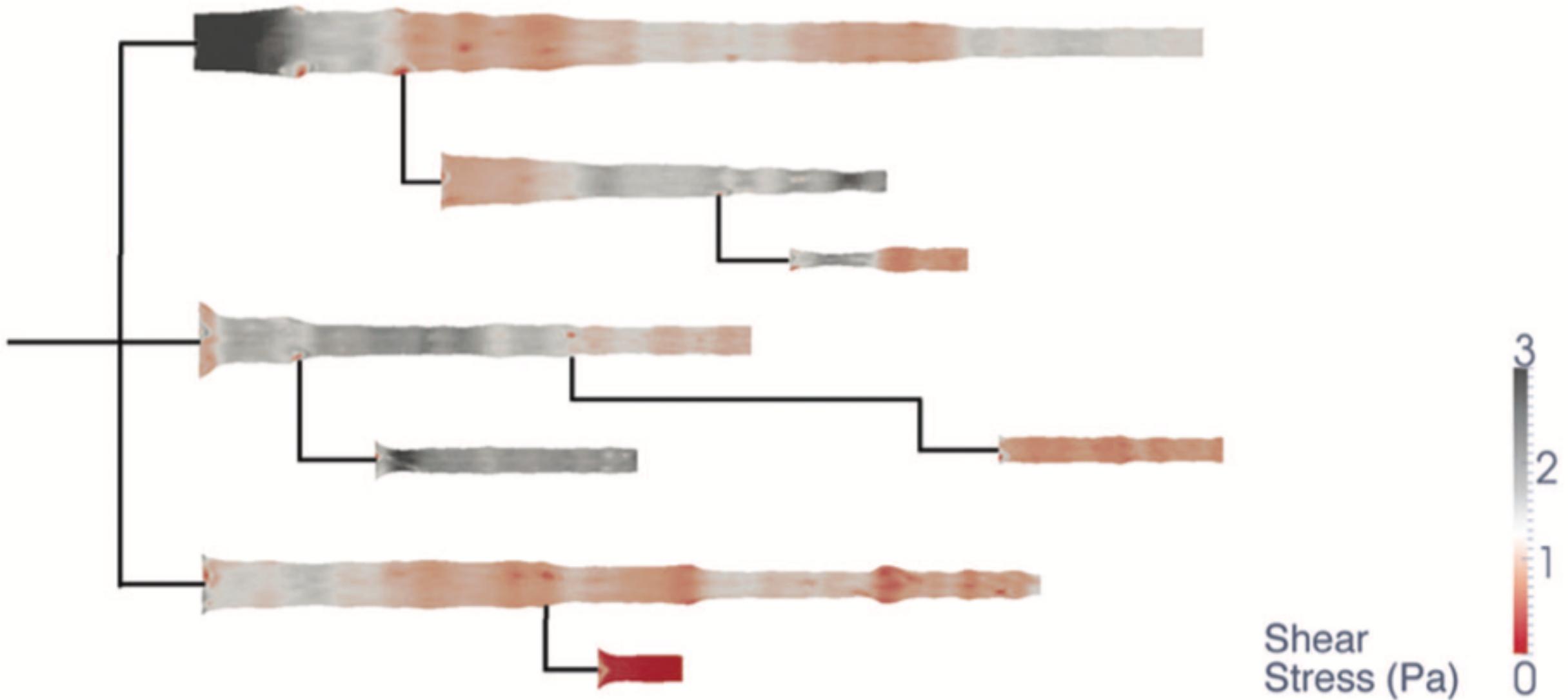
Projections

3D projection of left artery tree with diverging color map

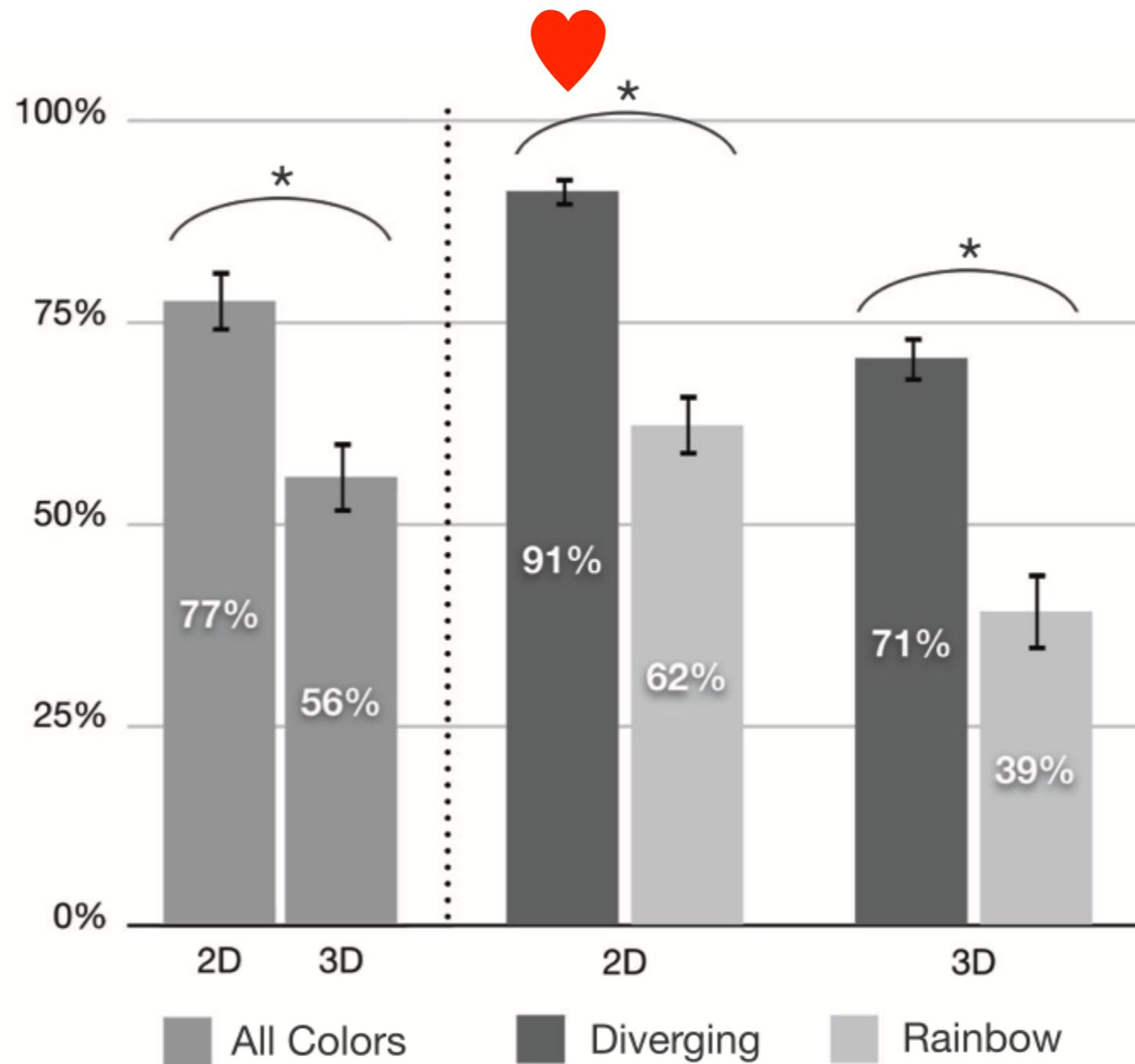


Projections

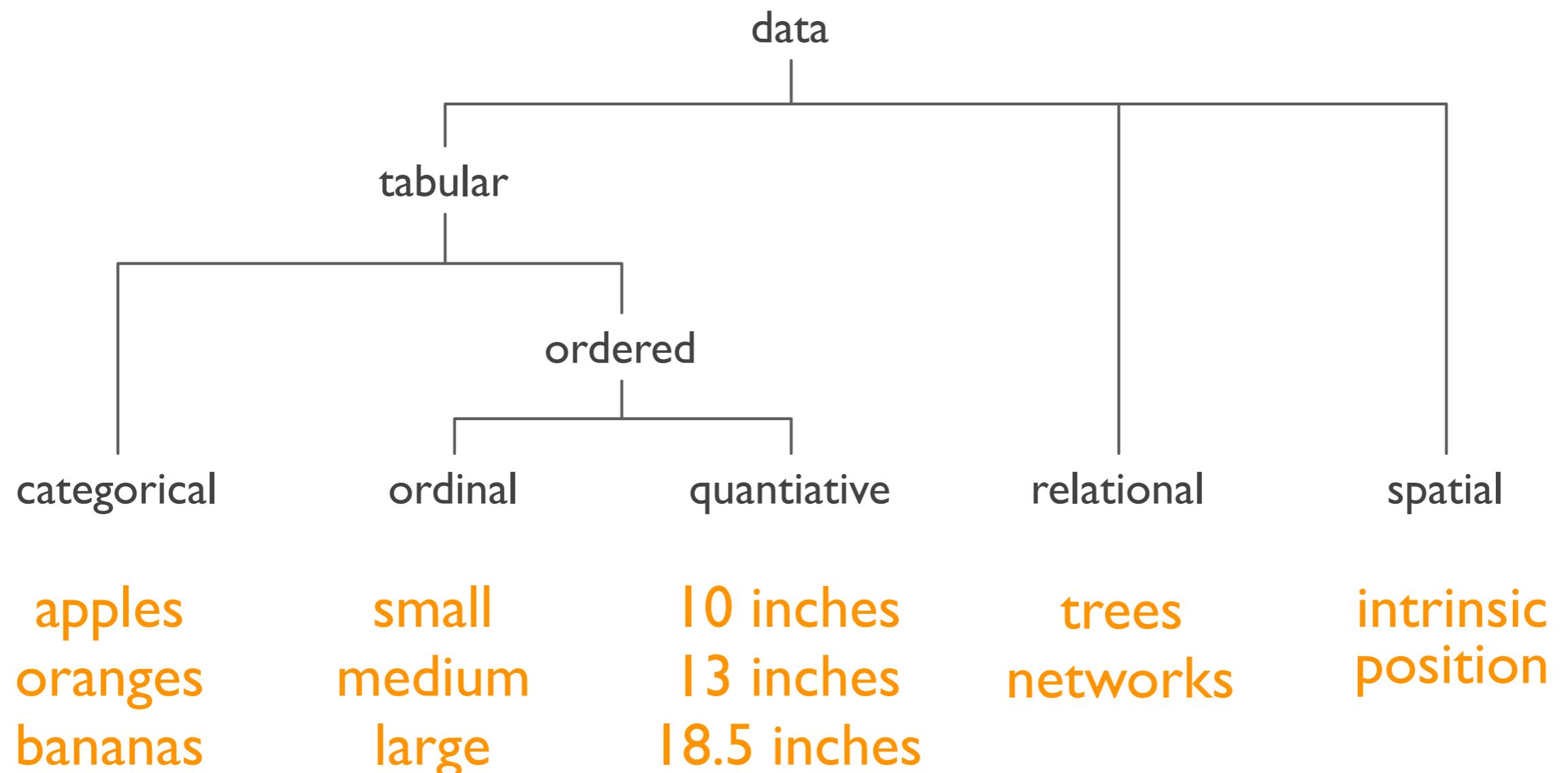
Novel 2D projection of left artery tree



Findings: Percentage of low ESS areas identified



Visual Encoding of Data

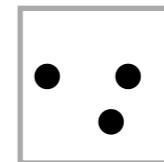


Visual Channels: Rankings

Categorical
What? Where?



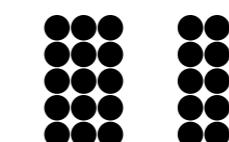
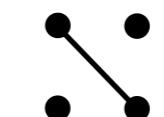
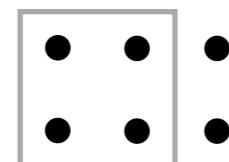
position*
planar
color hue
shape



Relational
With whom?



containment
connection
similarity
position*
proximity

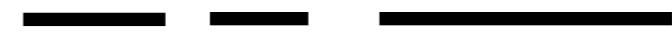
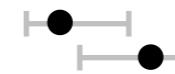


Visual Channels: Rankings

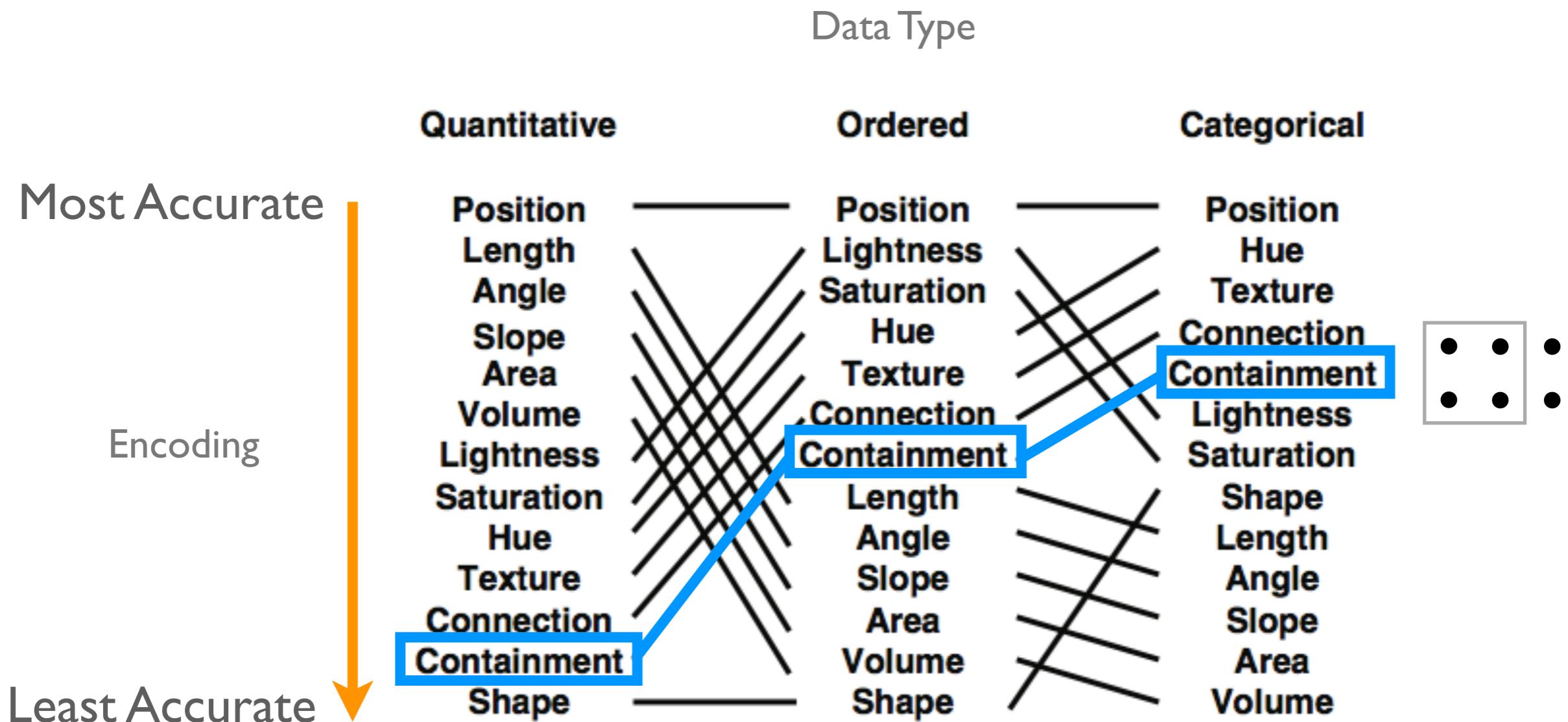
**Ordinal &
Quantitative**
How much?



- position*
common scale
- position*
unaligned scale
- length (1D)
- angle/tilt
- area (2D)
- curvature
- volume (3D)
- lightness
- color saturation



Ranking of Encodings



What's wrong? What's right?

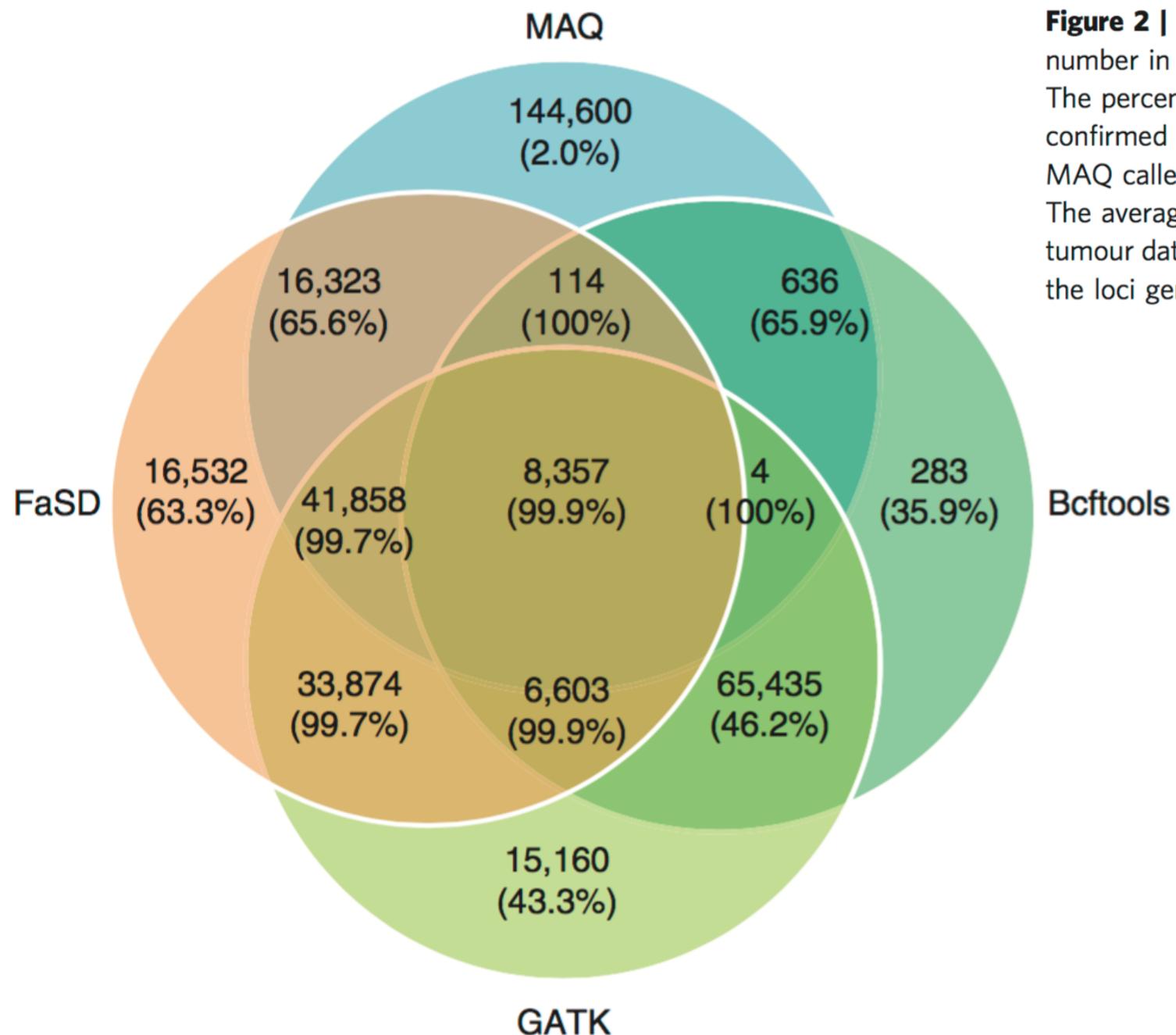


Figure 2 | The Venn diagram of SNPs detected by different tools. The number in each cell is the number of SNPs in the corresponding category. The percentage under the number is the proportion of SNPs that were confirmed by the Affymetrix SNP array. The FaSD, GATK, Bcftools and MAQ called 123661, 171291, 81432 and 211892 SNPs in total, respectively. The average depth of this data set was 10 \times . The figure is based on the tumour data set and Bowtie was used as aligner, and statistics are based on the loci genotyped by the Affymetrix SNP array.

What's wrong? What's right?

Euler Diagram

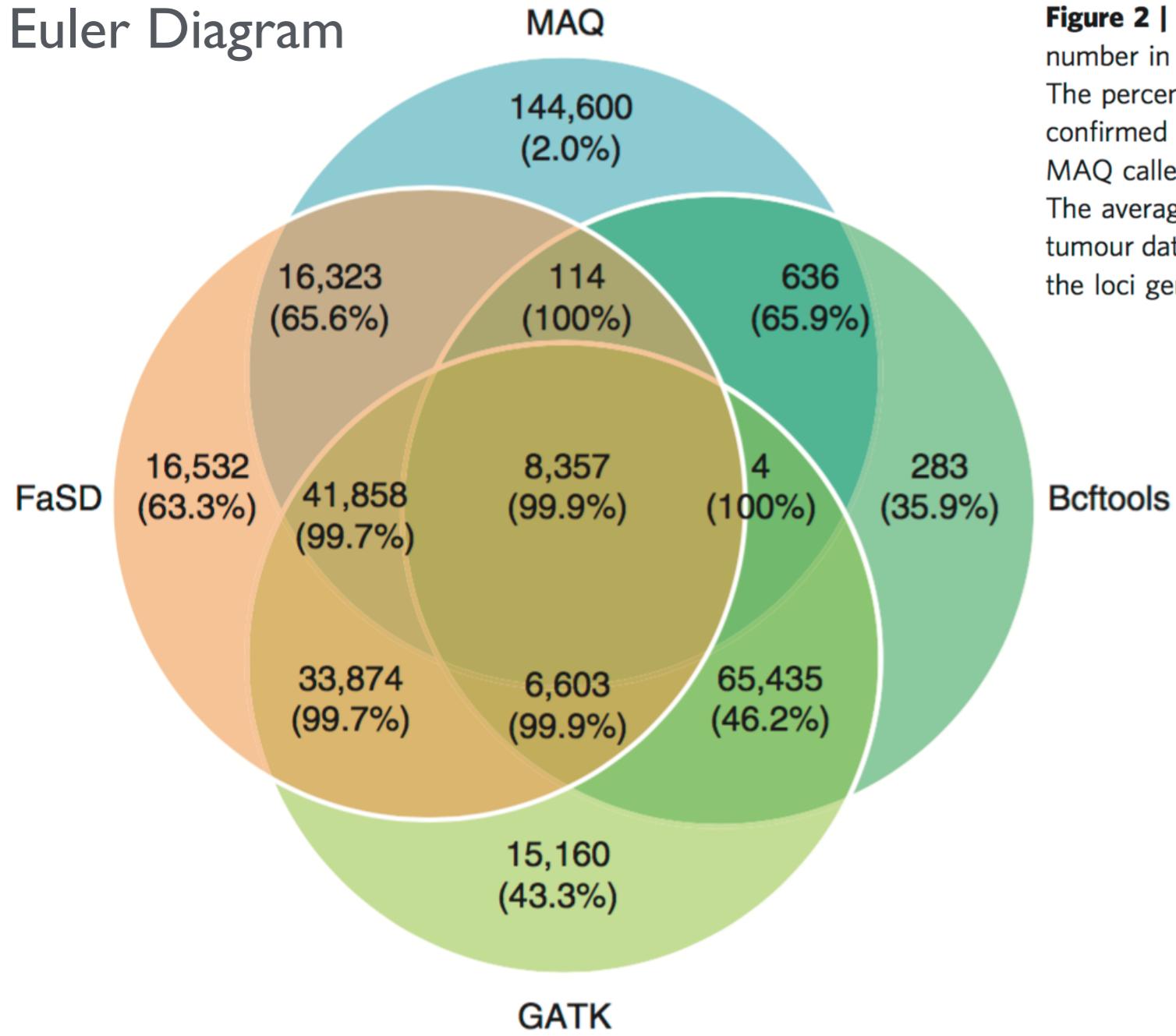
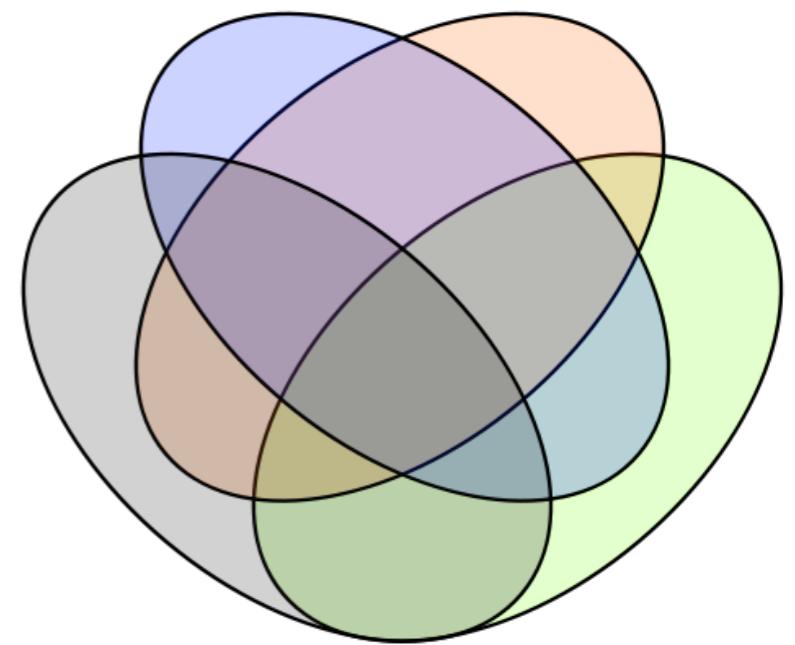


Figure 2 | The Venn diagram of SNPs detected by different tools. The number in each cell is the number of SNPs in the corresponding category. The percentage under the number is the proportion of SNPs that were confirmed by the Affymetrix SNP array. The FaSD, GATK, Bcftools and MAQ called 123661, 171291, 81432 and 211892 SNPs in total, respectively. The average depth of this data set was 10 × . The figure is based on the tumour data set and Bowtie was used as aligner, and statistics are based on the loci genotyped by the Affymetrix SNP array.

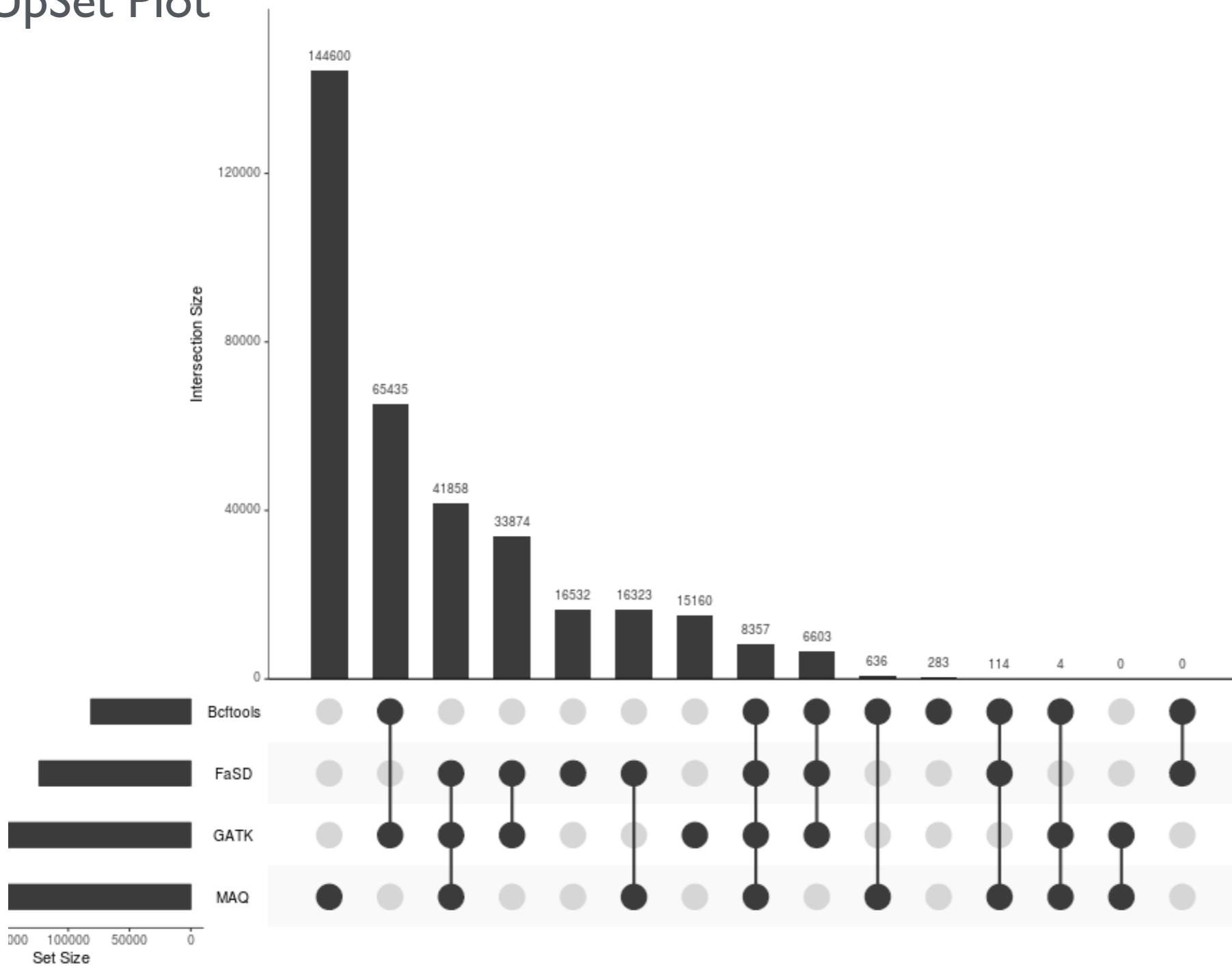
Bcftools

Venn Diagram

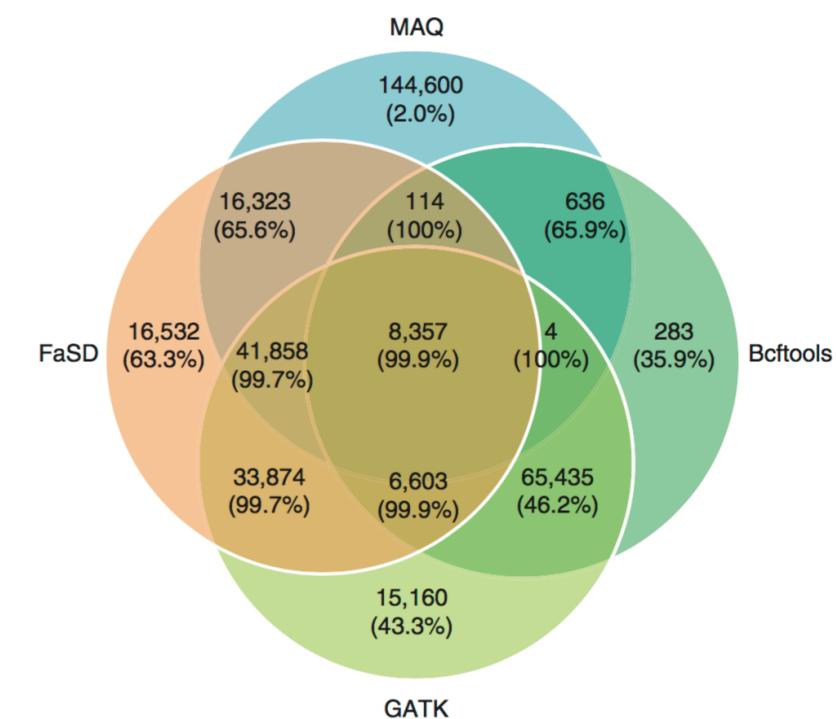


What's wrong? What's right?

UpSet Plot

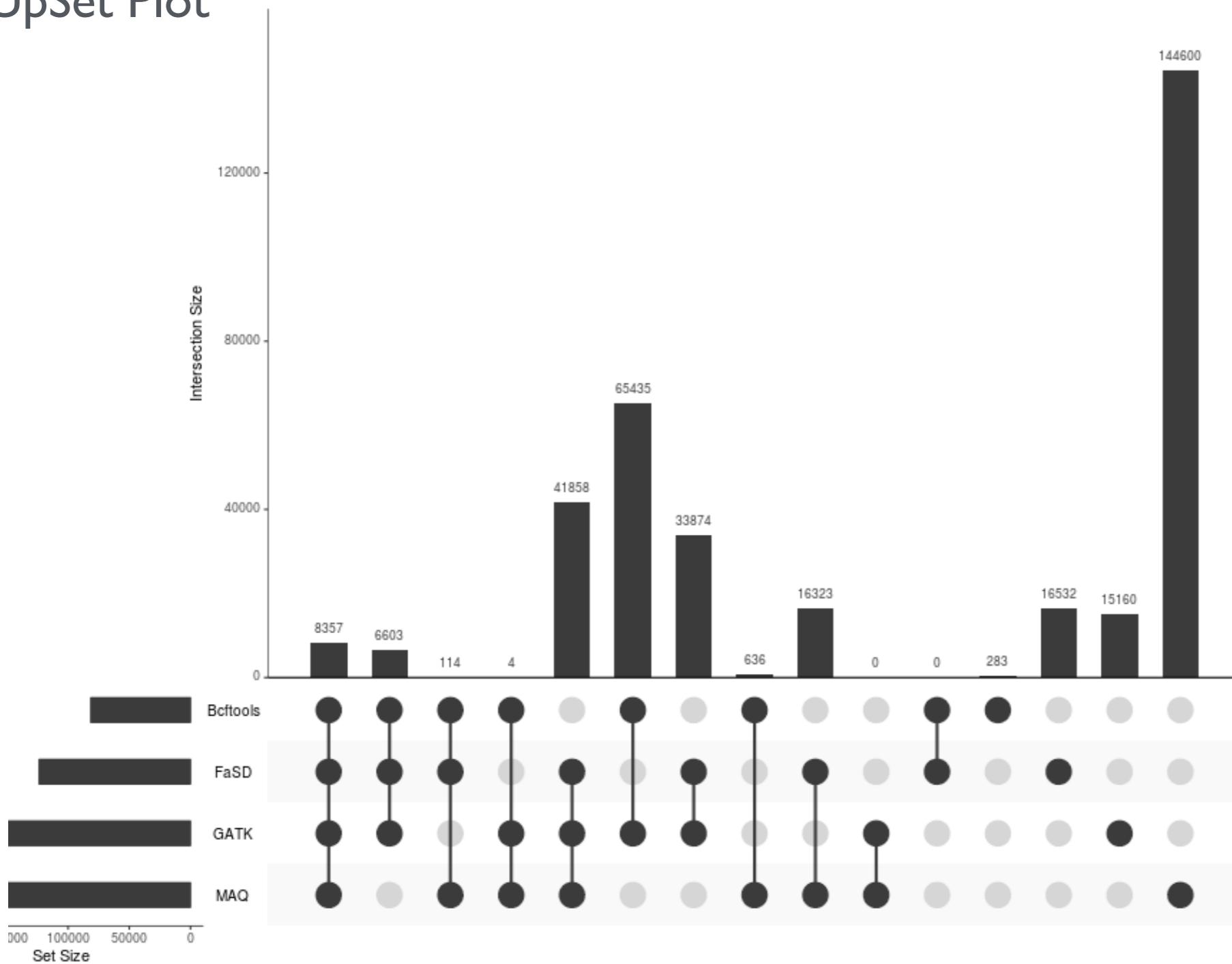


Euler Diagram

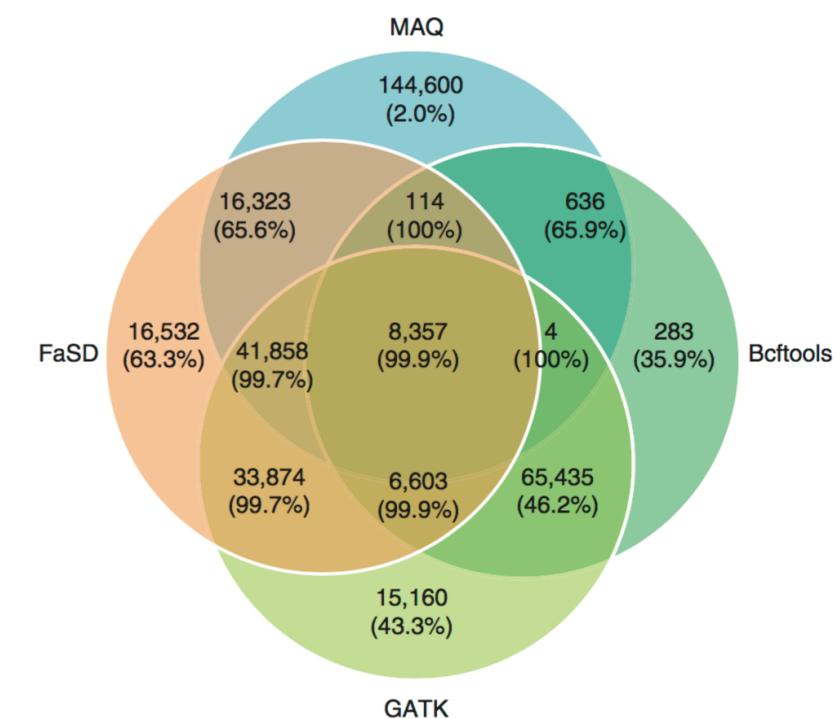


What's wrong? What's right?

UpSet Plot



Euler Diagram



Ranking of Encodings

Principle of Importance Ordering (Mackinaly 1986):

Encode more important information more effectively.

Ranking of Encodings

- How accurately can the data be read from the visualization?

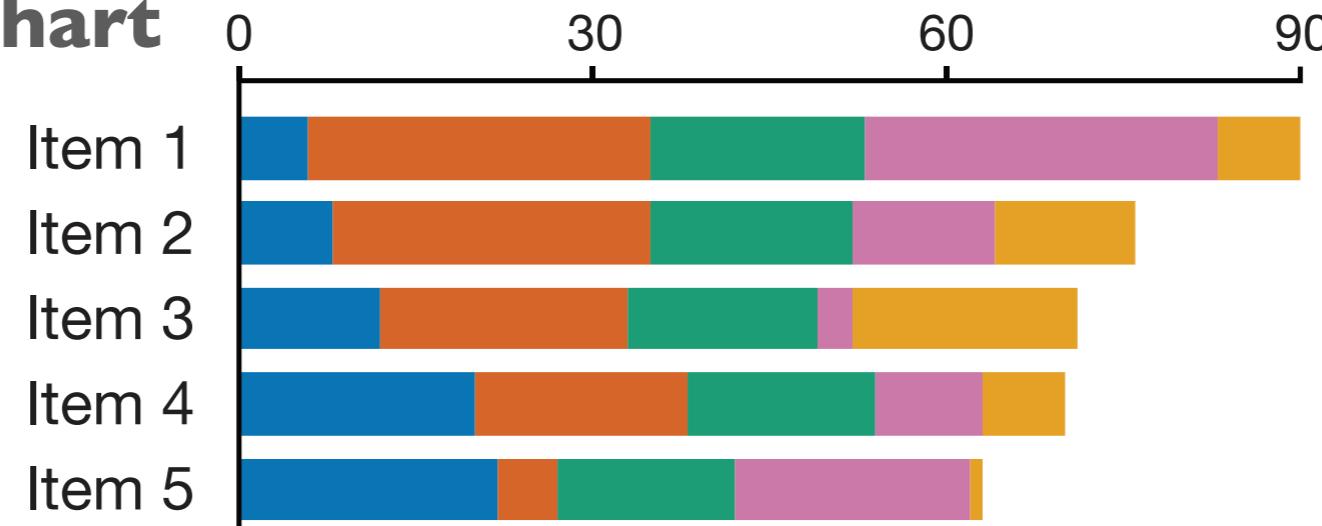
Bar Charts for Multiple Items and Categories

	1	2	3	4	5
Item	6	29	18	30	7
Item	8	27	17	12	12
Item	12	21	16	3	19
Item	20	18	16	9	7
Item	22	5	15	20	1

Bar Charts for Multiple Items and Categories

	1	2	3	4	5
Item	6	29	18	30	7
Item	8	27	17	12	12
Item	12	21	16	3	19
Item	20	18	16	9	7
Item	22	5	15	20	1

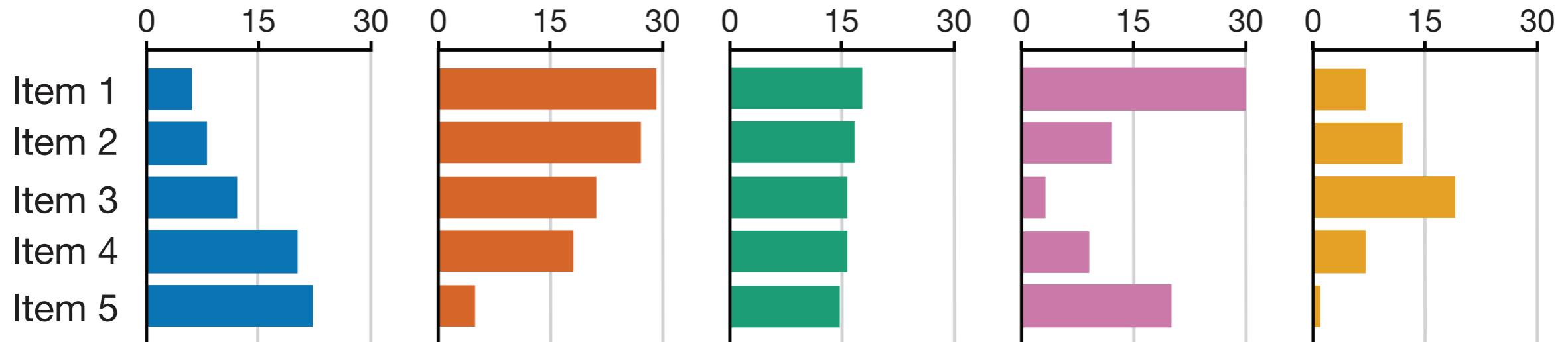
Stacked Bar Chart



Bar Charts for Multiple Items and Categories

	1	2	3	4	5
Item	6	29	18	30	7
Item	8	27	17	12	12
Item	12	21	16	3	19
Item	20	18	16	9	7
Item	22	5	15	20	1

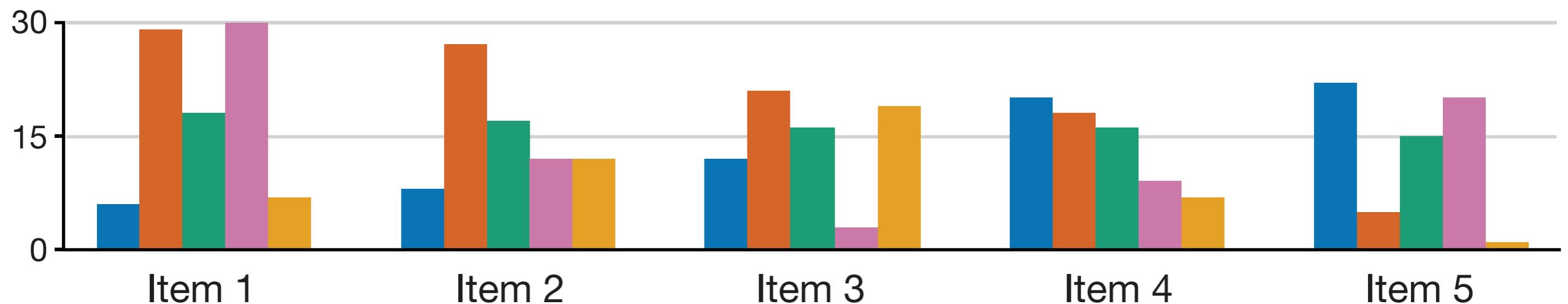
Layered Bar Chart



Bar Charts for Multiple Items and Categories

	1	2	3	4	5
Item	6	29	18	30	7
Item	8	27	17	12	12
Item	12	21	16	3	19
Item	20	18	16	9	7
Item	22	5	15	20	1

Grouped Bar Chart

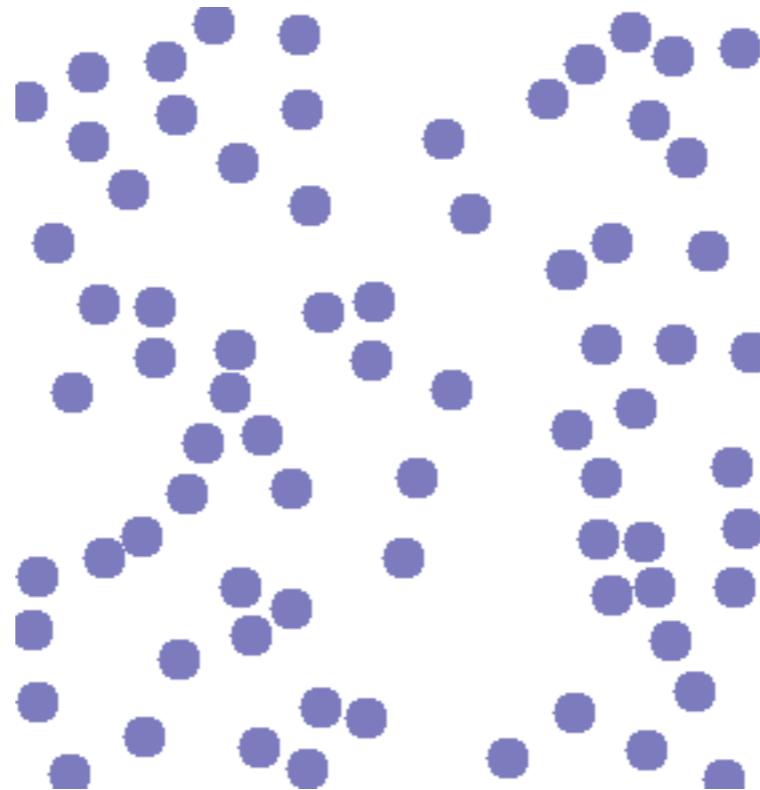


Ranking of Encodings

- How accurately can the data be read from the visualization?
- Which channels are processed preattentively?

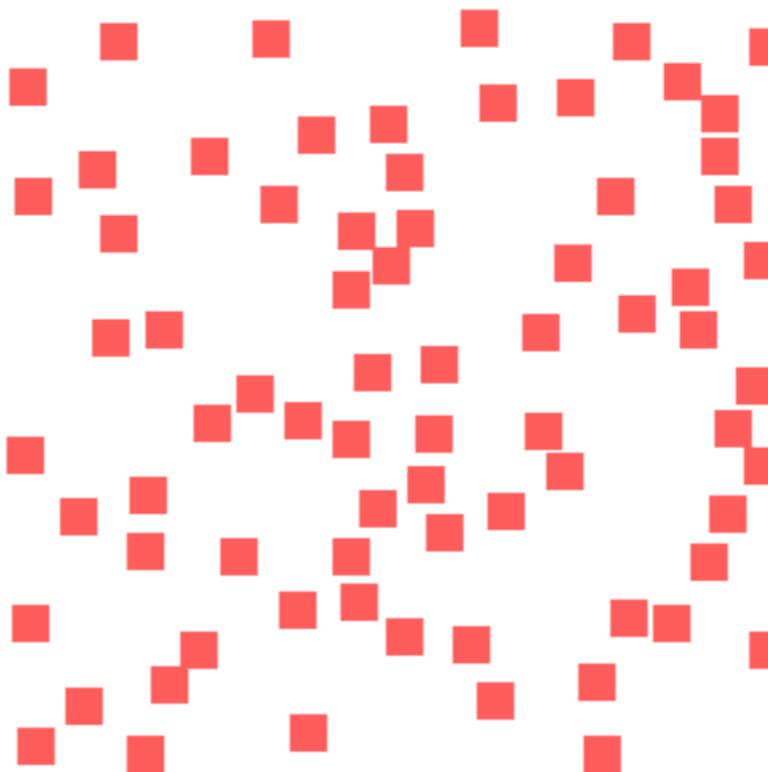
Preattentive Processing: Color

Can you spot this: ●?



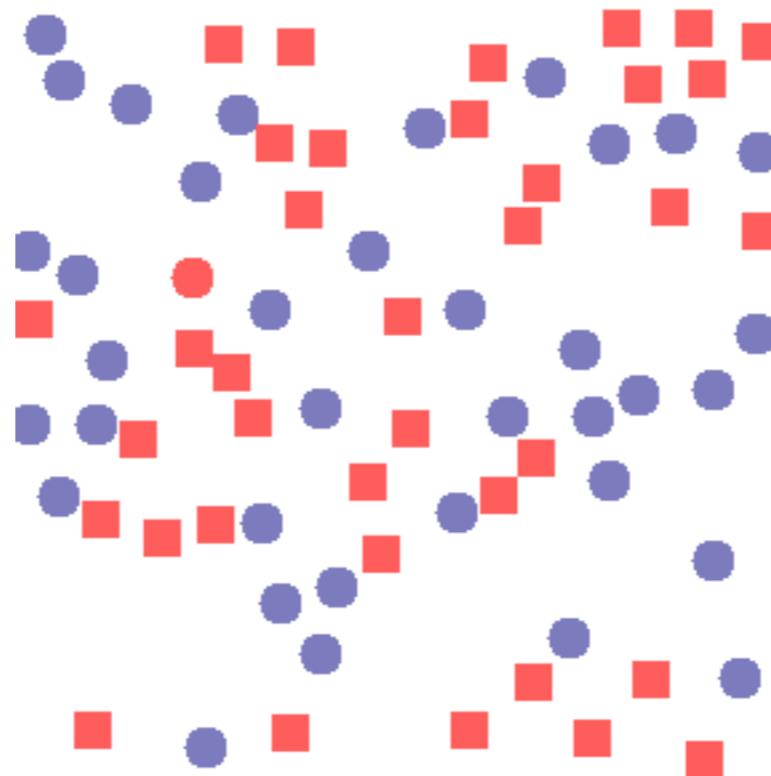
Preattentive Processing: Shape

Can you spot this: ●?



Preattentive Processing: Shape and Color

Can you spot this: ●?



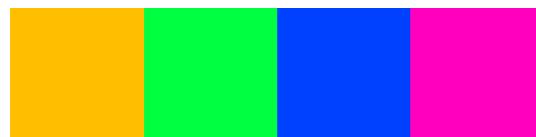
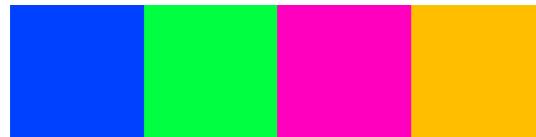
Ranking of Encodings

- How accurately can the data be read from the visualization?
- Which channels are processed preattentively?
- How many classes can be distinguished?
- Can the channels be separated from each other?

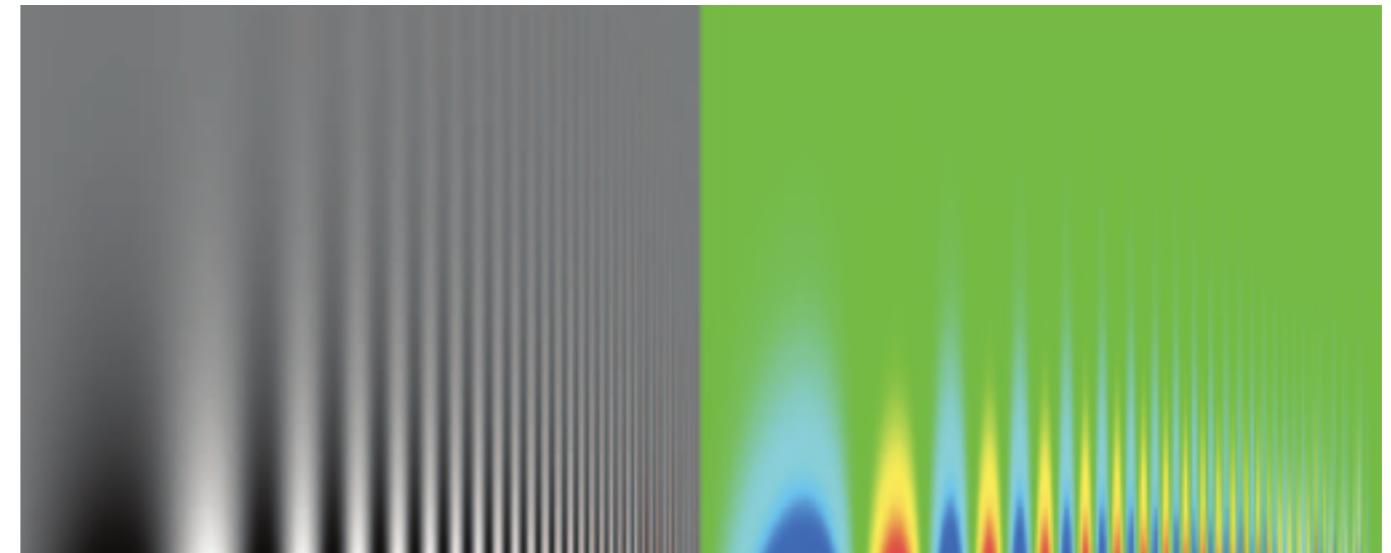
Common Mistakes

Color Pitfalls: Rainbow Color Map

hard to order



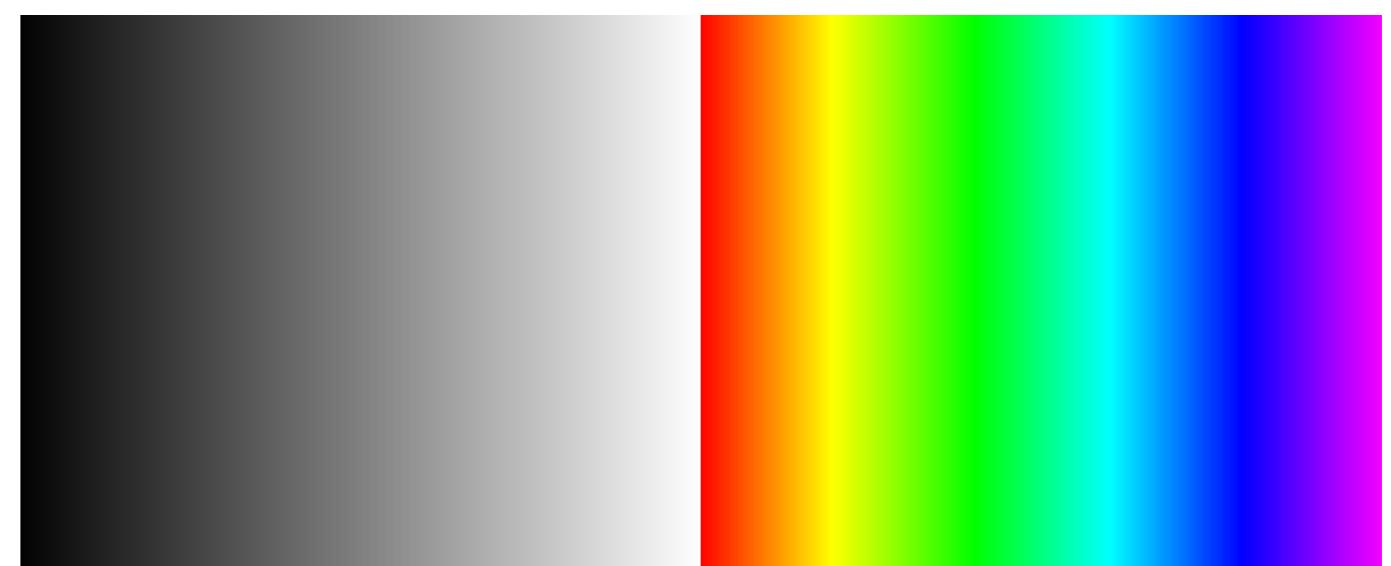
lower resolution



easy to order

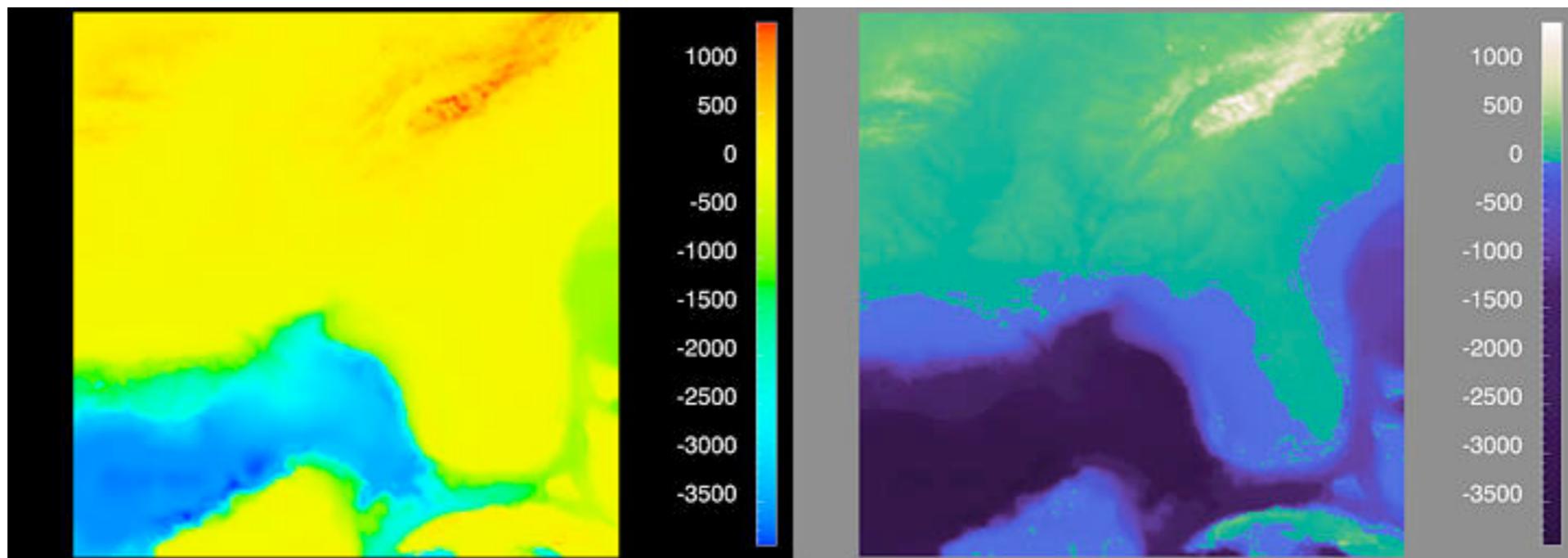


creates artifacts



Color Pitfalls: Rainbow Color Map

Southeastern United States and Gulf of Mexico

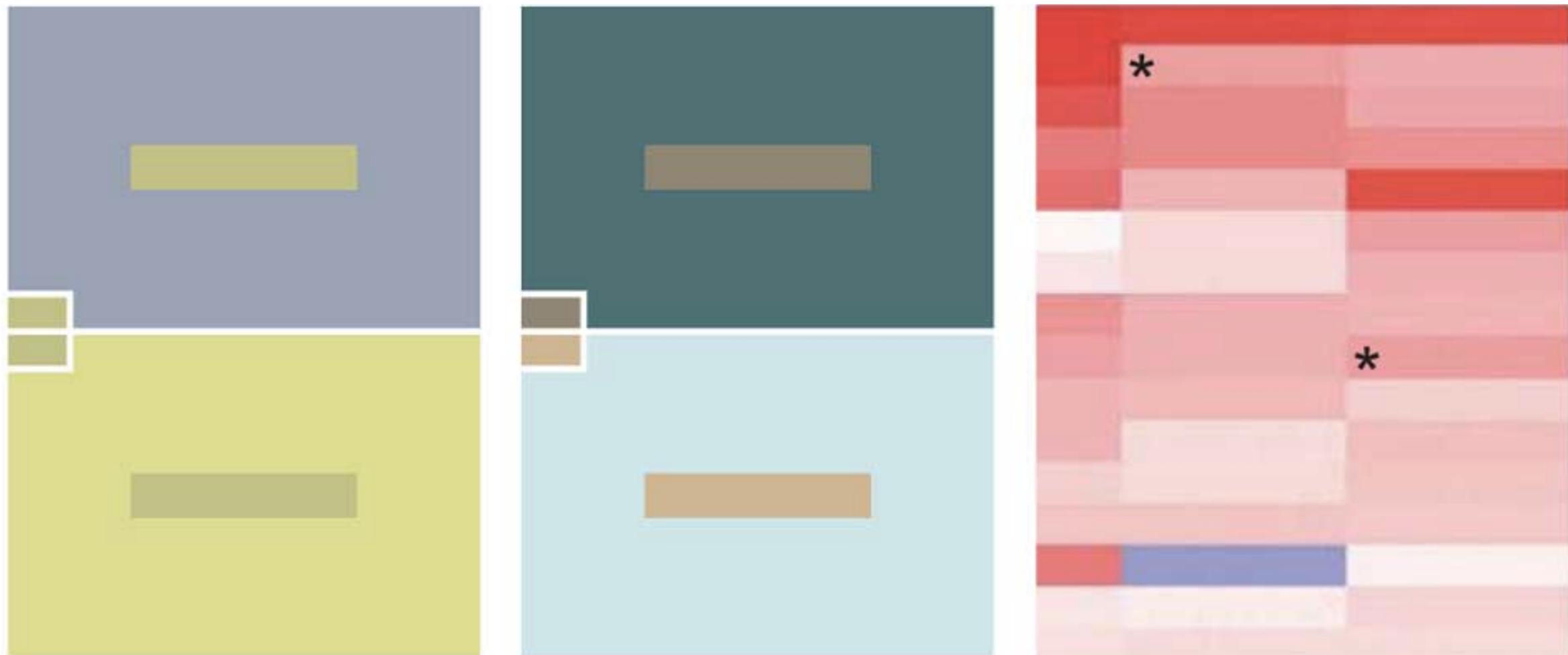


Problems:

- zero crossing not explicit
- lack of ordering of colors makes it hard to interpret the map

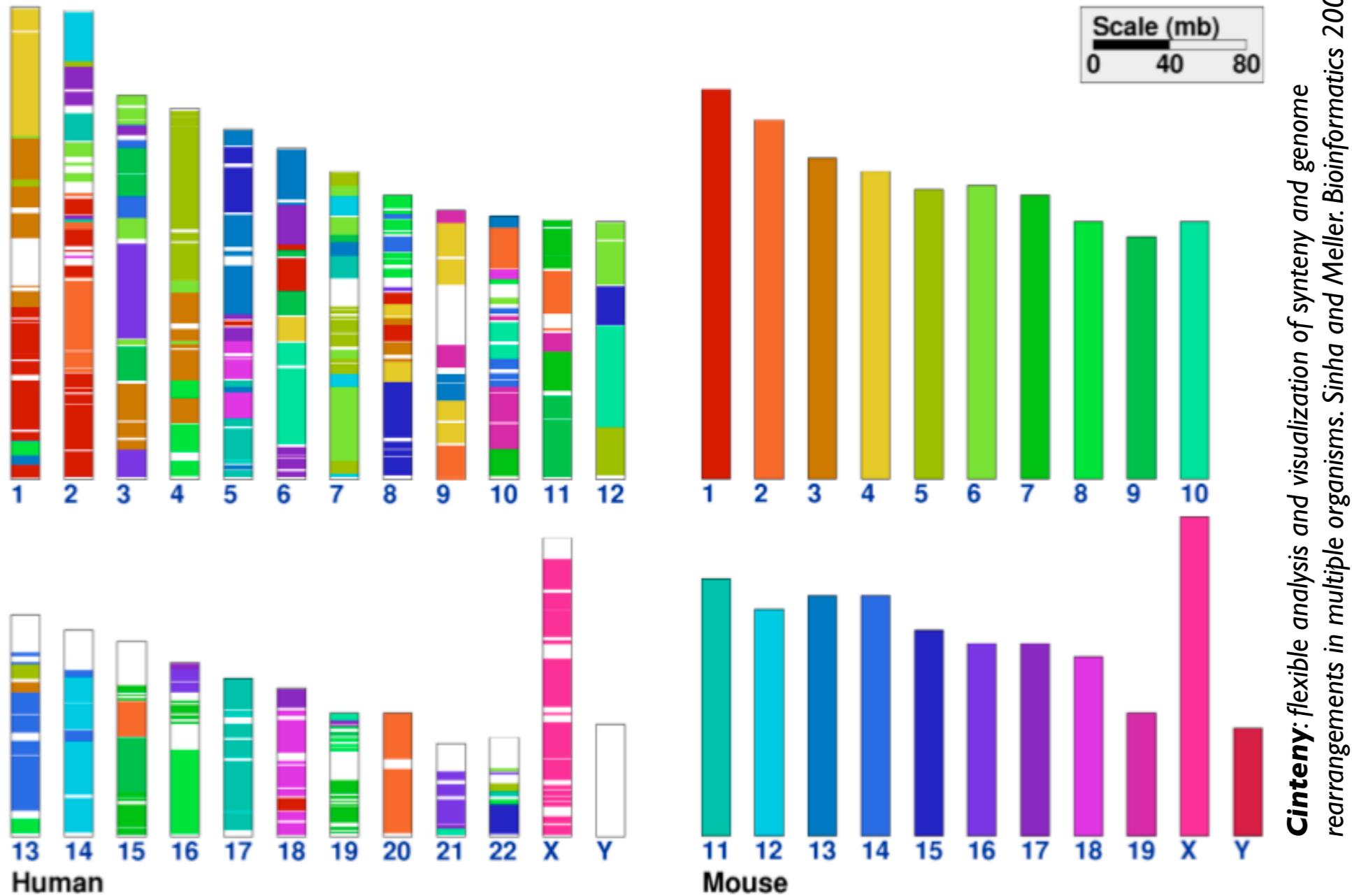
Color Pitfalls: Relativity

Color is a relative medium and context matters

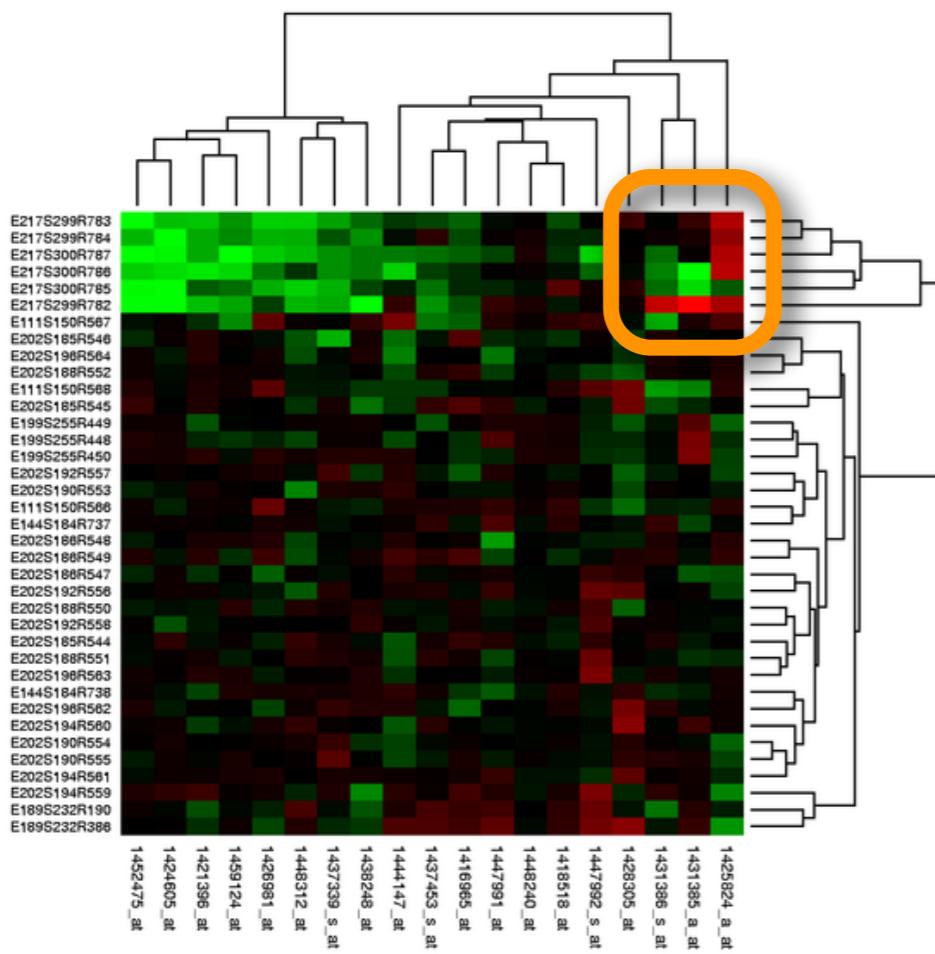


Color Pitfalls: Discriminability

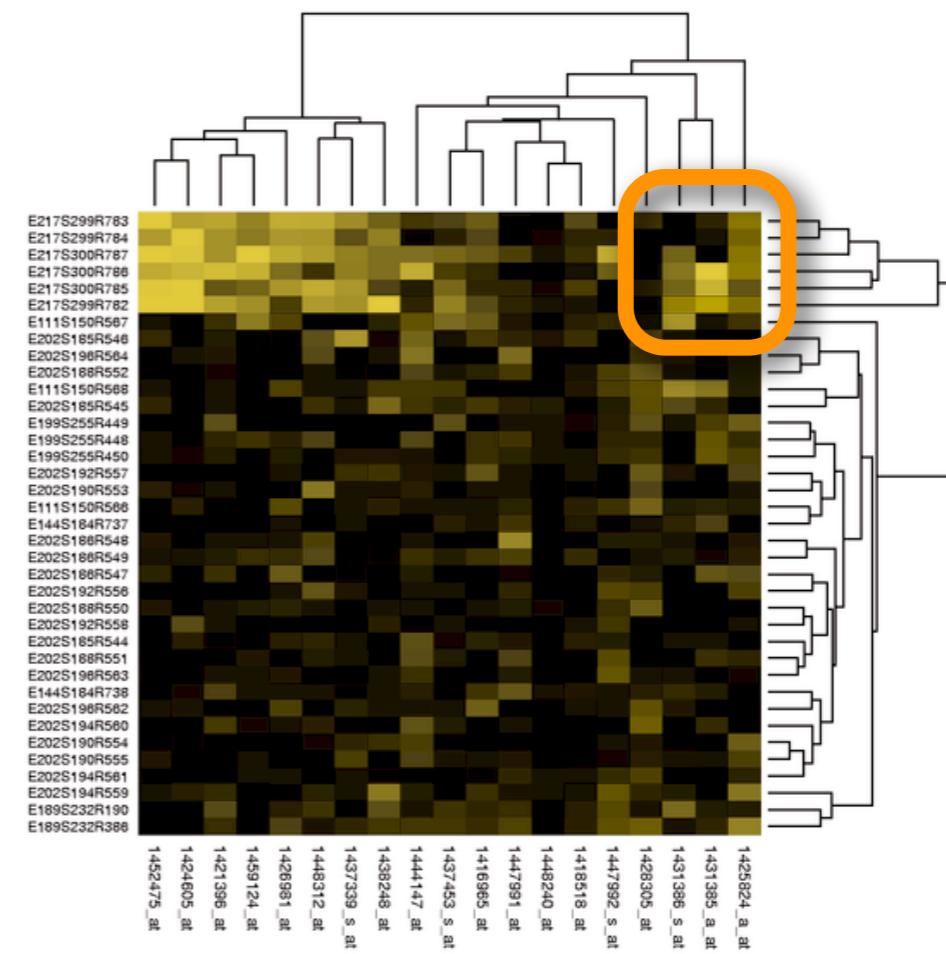
Only 6-12 colors are visually discernible



Color Pitfalls: Color Blindness



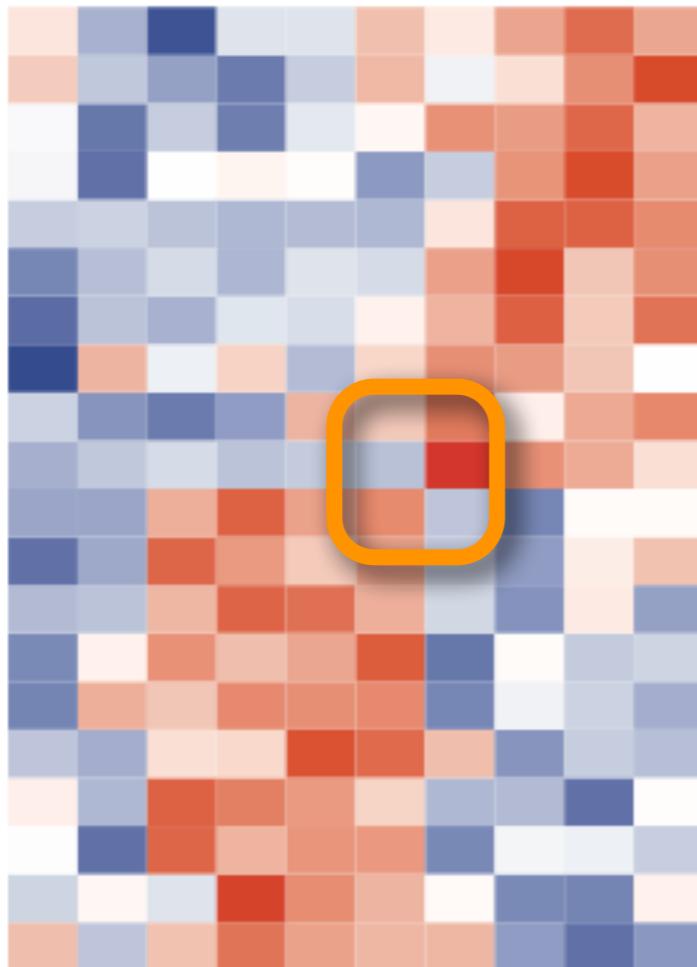
Normal Vision



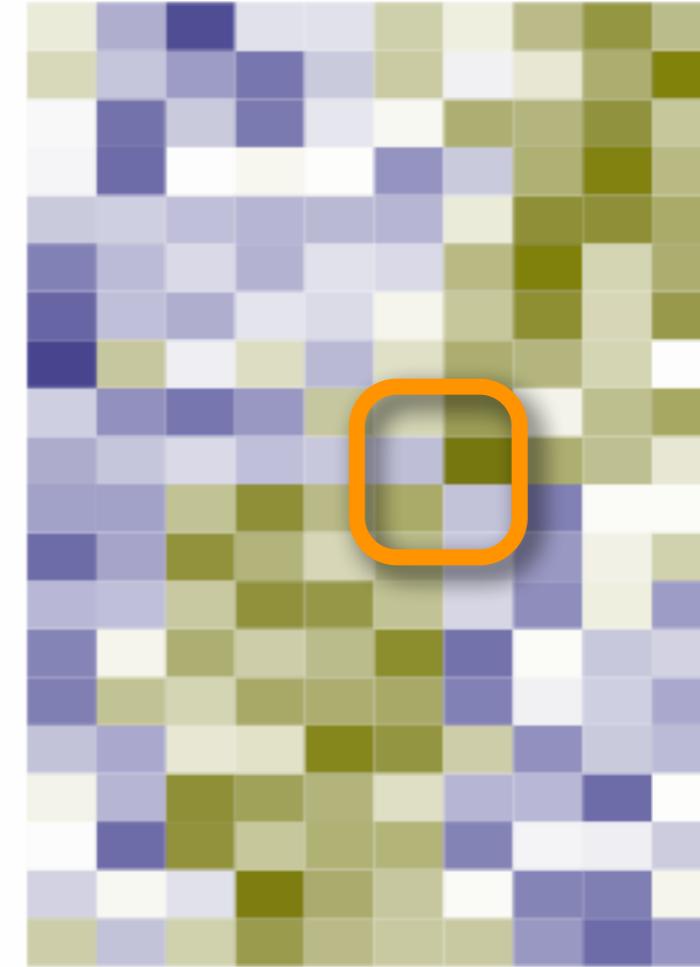
Deutanope Vision
("Red-Green Blindness")

~ 7% of male population affected

Color Pitfalls: Color Blindness



Normal Vision

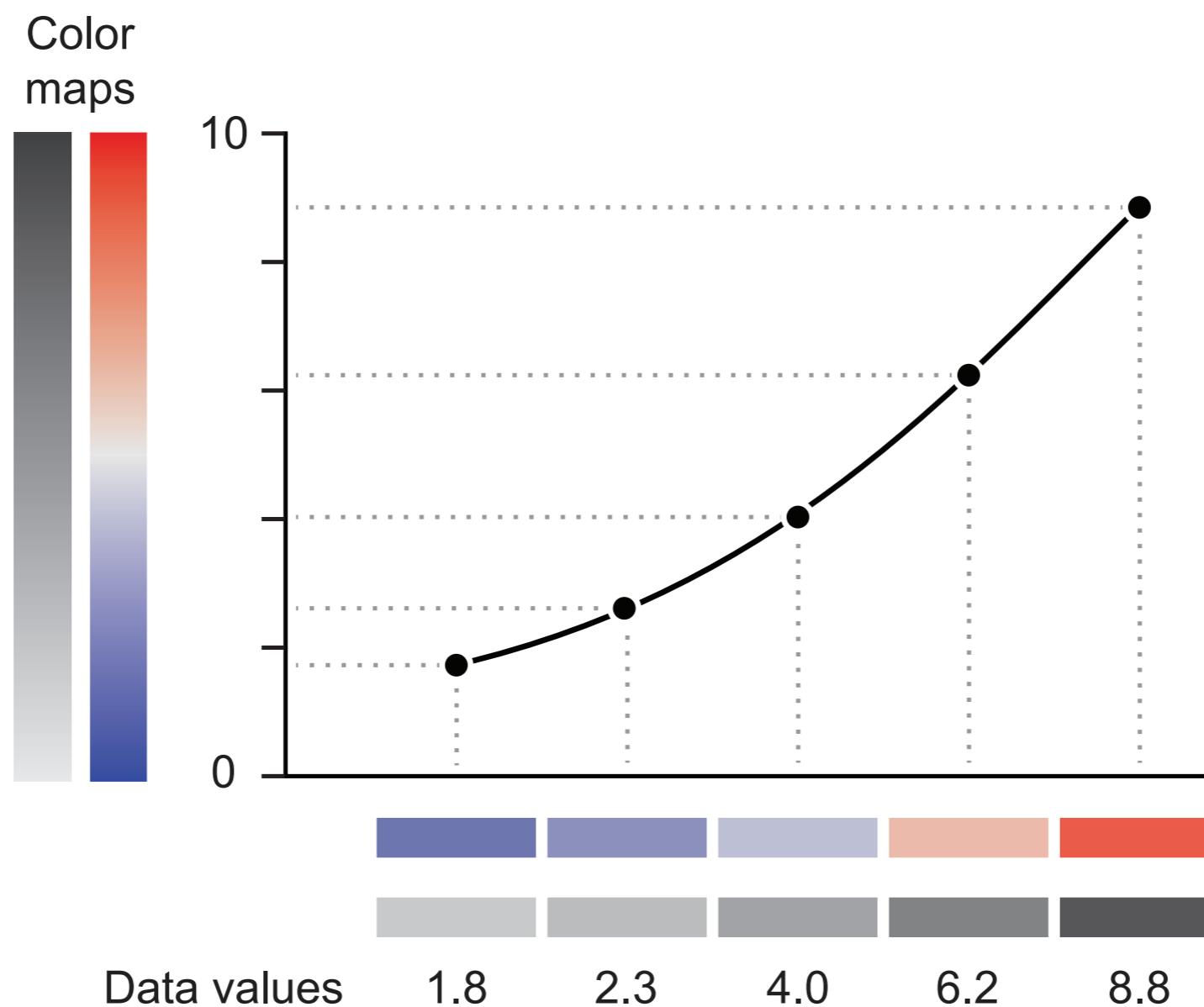


Deutanope Vision
("Red-Green Blindness")

~ 7% of male population affected

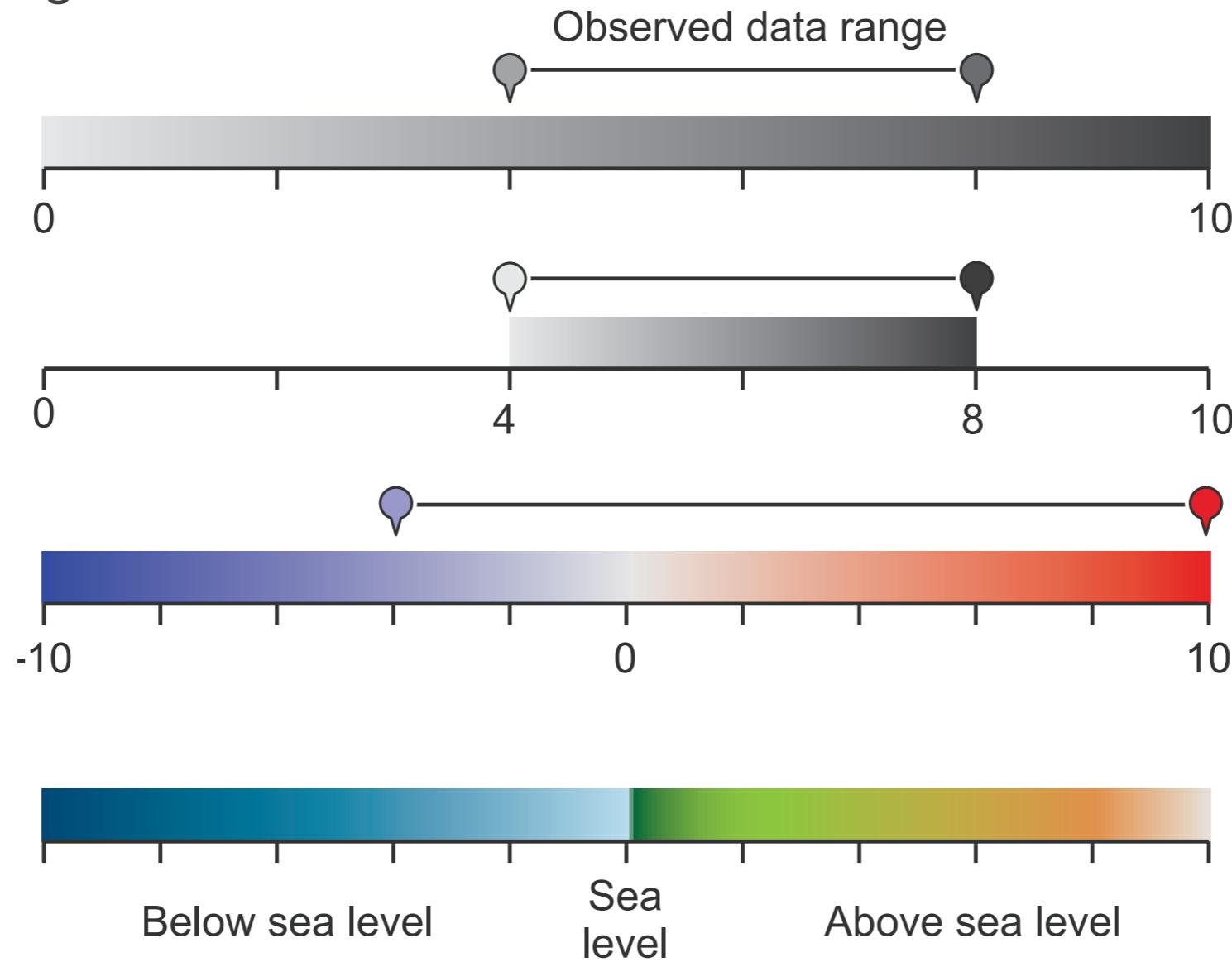
Color Pitfalls: Color Mapping

Bad color mapping!



Color Pitfalls: Color Mapping

Good color mapping!



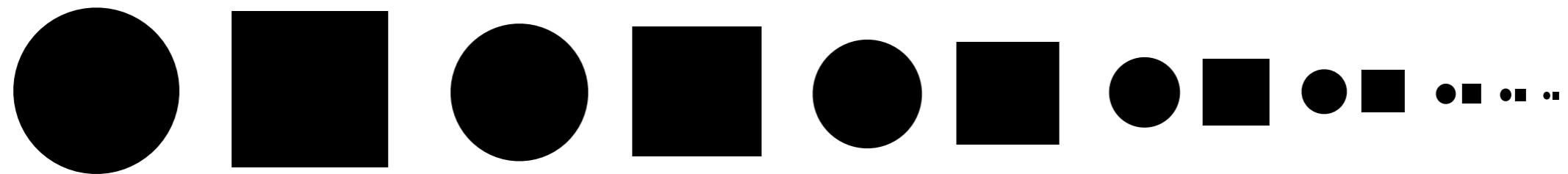
Remember!

Color used poorly is worse than no color at all.

— *Edward Tufte*

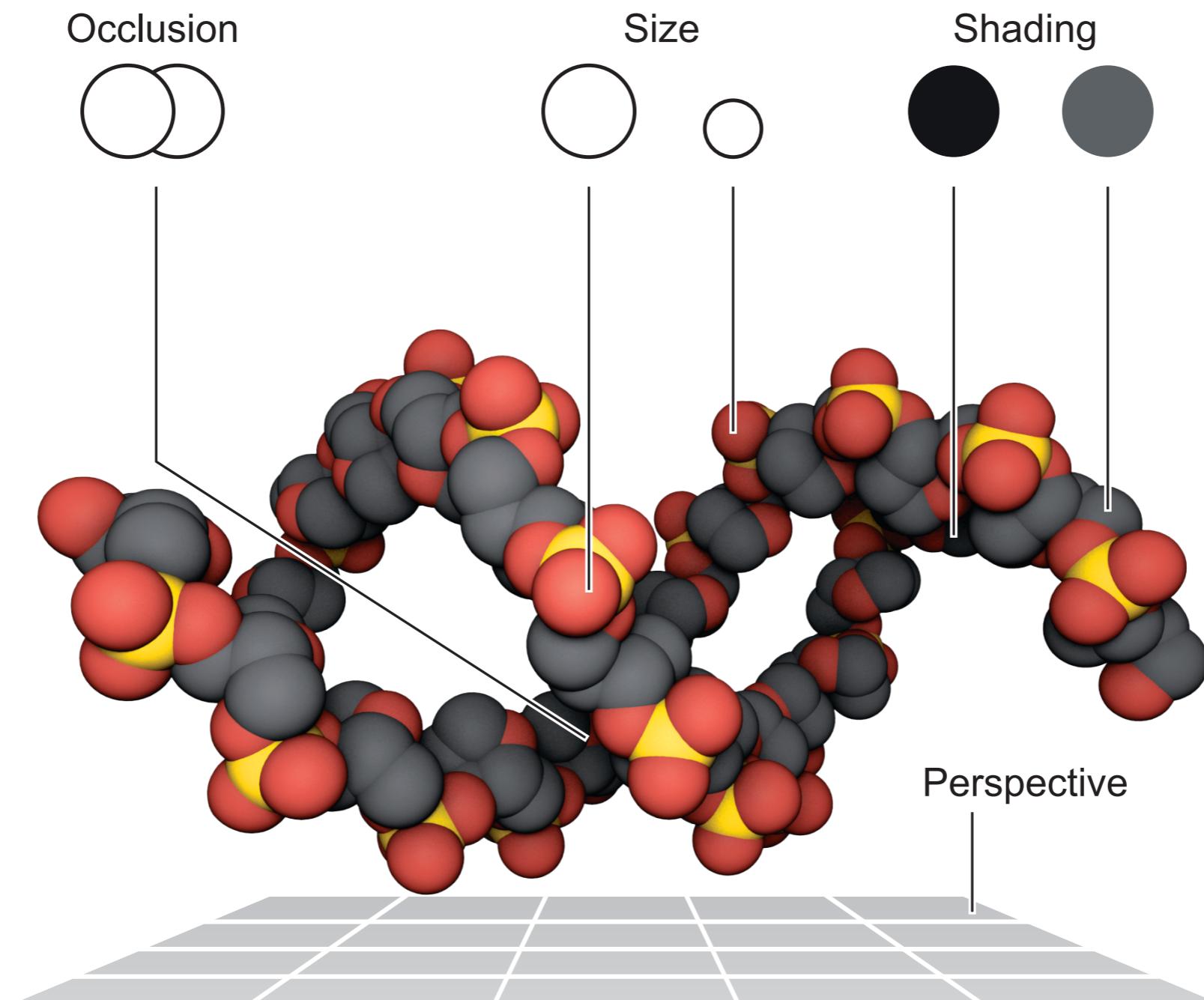
after Mackinlay, ACM Transactions on Graphics 5, 110-141, 1986

Encoding Pitfalls: Interference between Channels



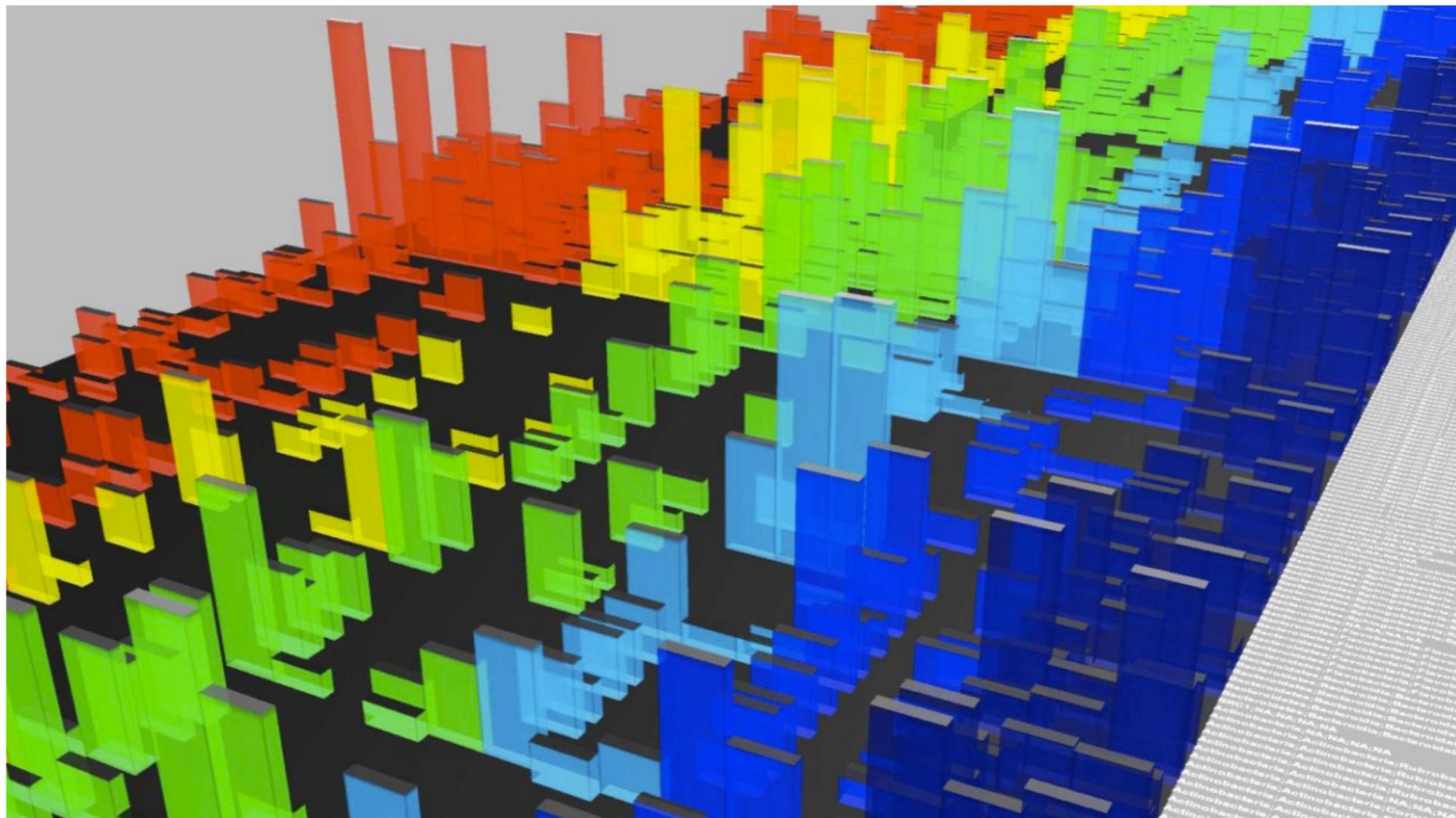
shape and size

3D: Depth Cues

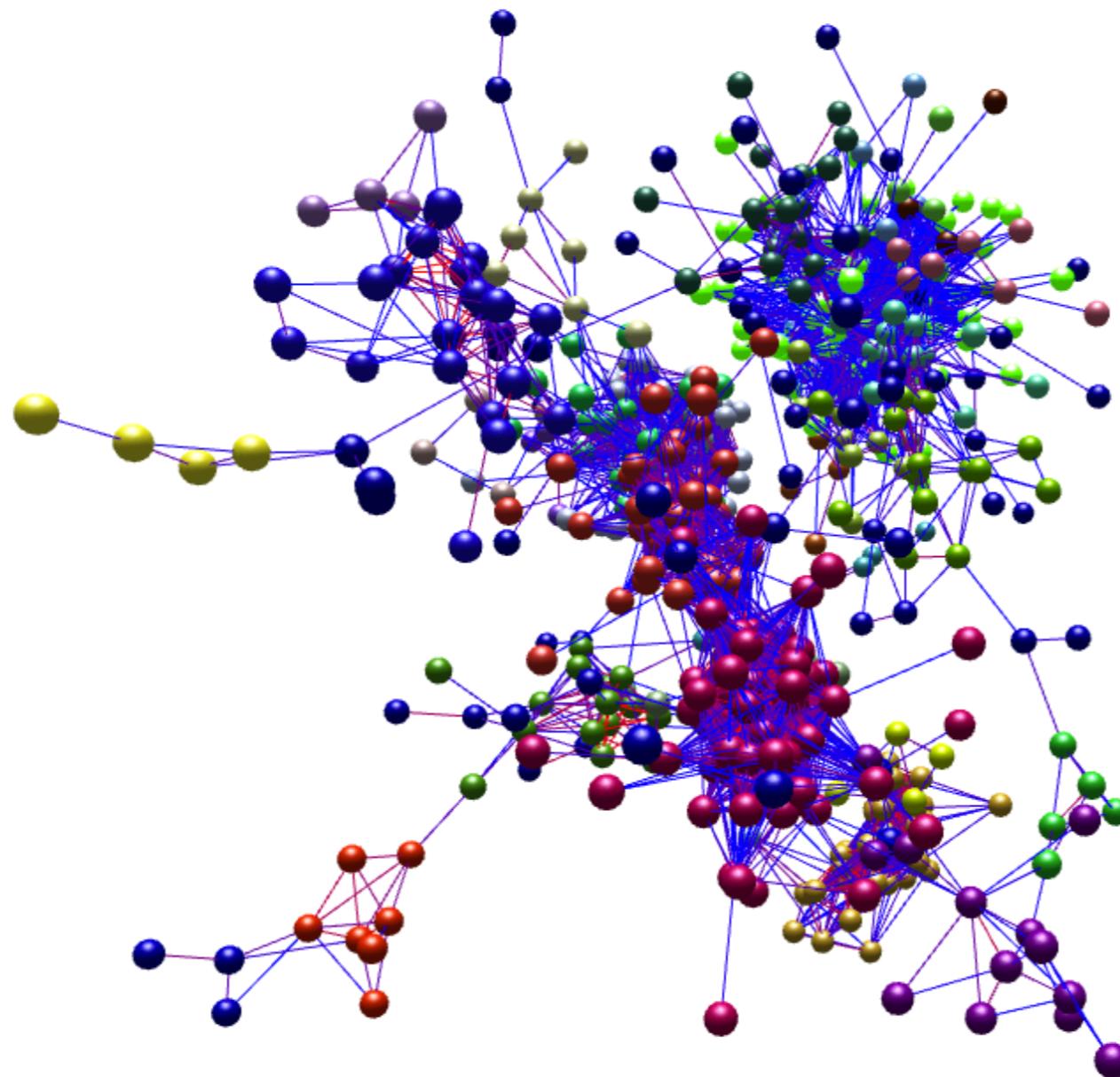


3D Pitfalls: Perspective

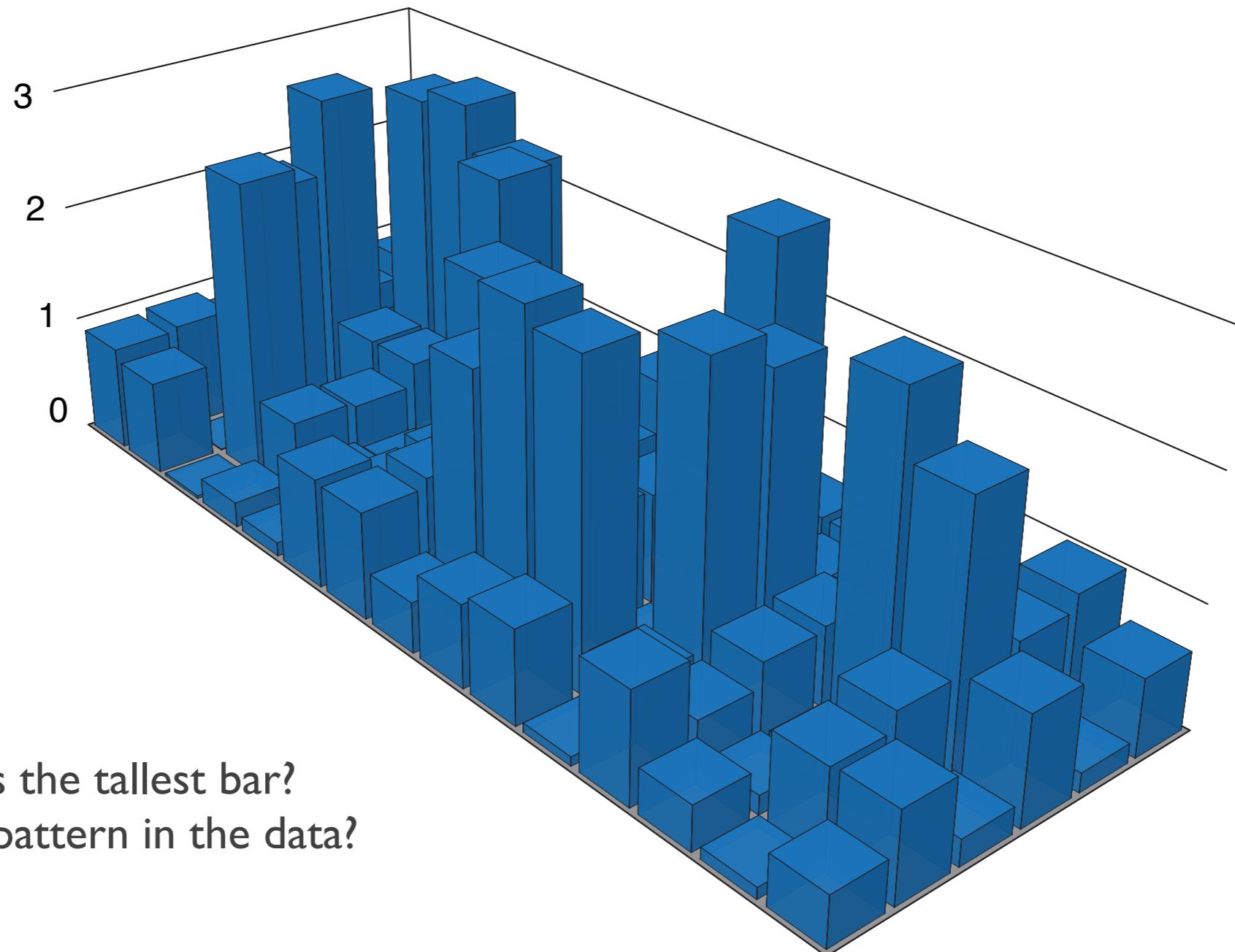
- **perspective distortion:** interferes with size channel encoding
- **shading:** interferes with color, lightness, and saturation channel encodings



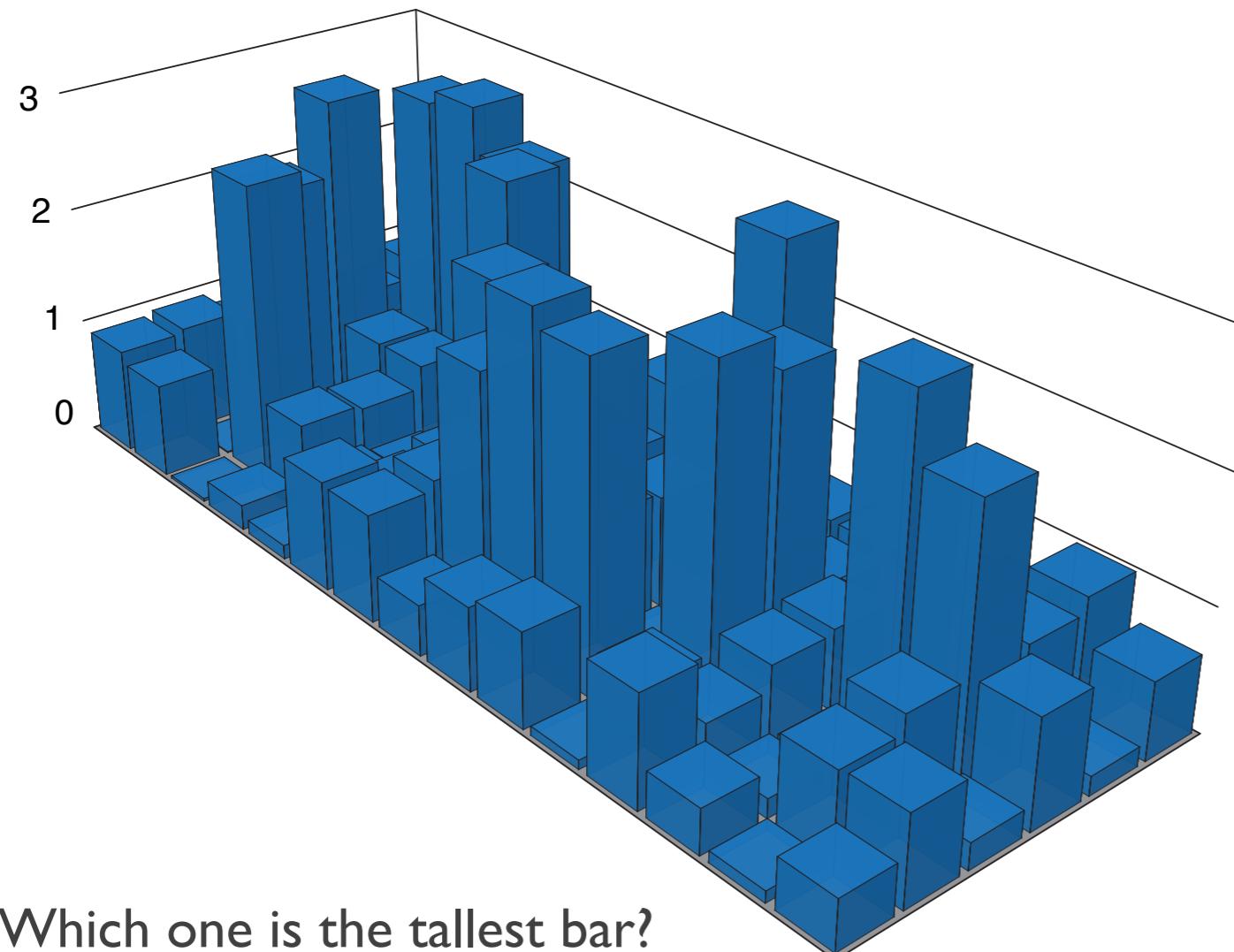
3D Pitfalls: Occlusion



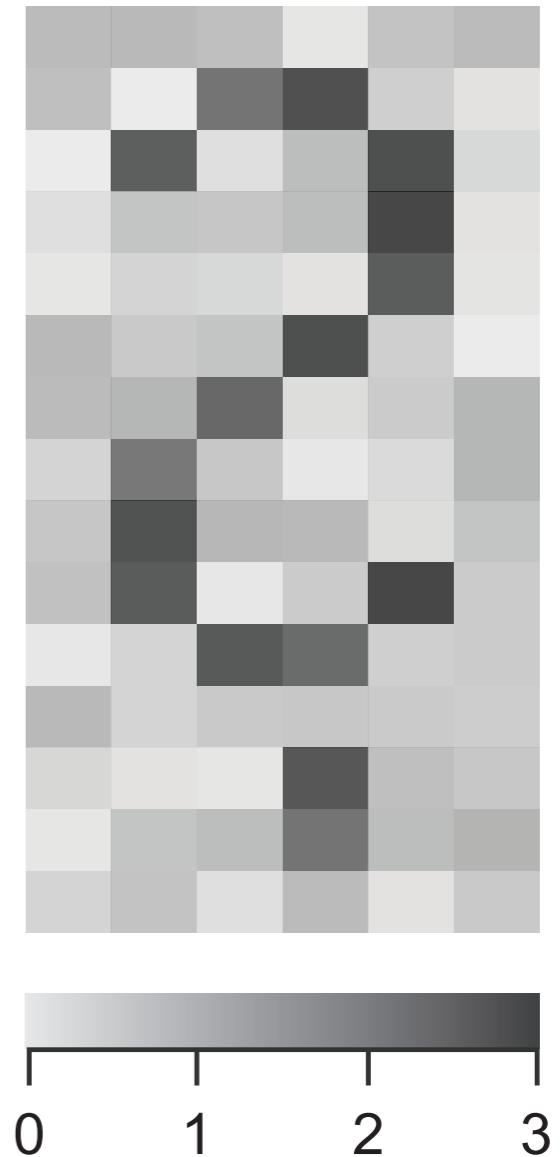
3D Pitfalls: Occlusion and Perspective



3D Pitfalls: Occlusion and Perspective



Which one is the tallest bar?
What is the pattern in the data?



Visualization for Communication

Insights into the Paper, Data into the Database

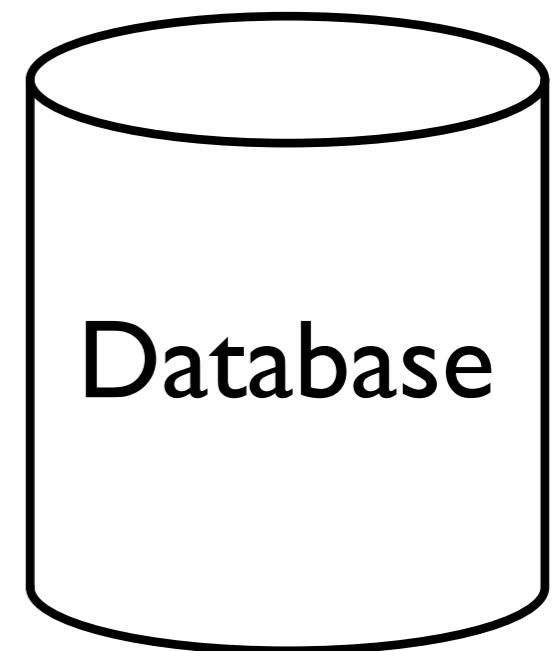
ARTICLE

doi:10.1038/nature11412

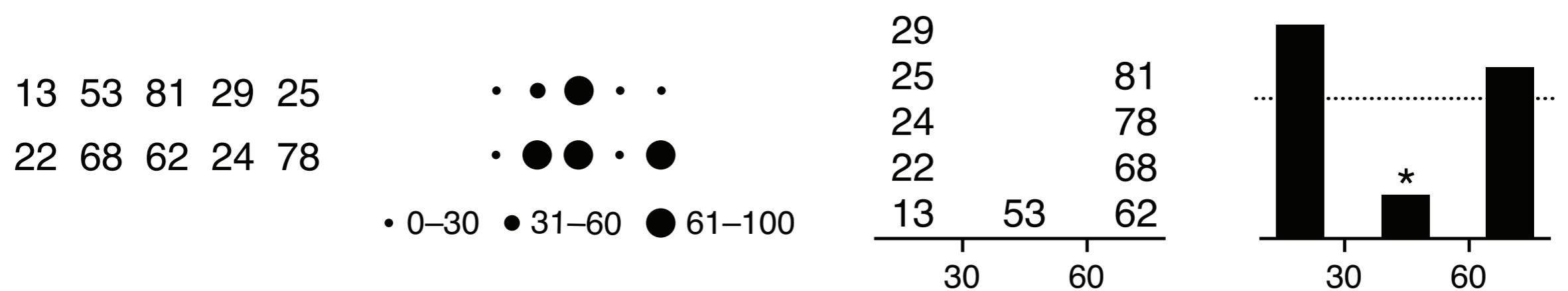
Comprehensive molecular portraits of human breast tumours

The Cancer Genome Atlas Network*

We analysed primary breast cancers by genomic DNA copy number arrays, DNA methylation, exome sequencing, messenger RNA arrays, microRNA sequencing and reverse-phase protein arrays. Our ability to integrate information across platforms provided key insights into previously defined gene expression subtypes and demonstrated the existence of four main breast cancer classes when combining data from five platforms, each of which shows significant molecular heterogeneity. Somatic mutations in only three genes (*TP53*, *PIK3CA* and *GATA3*) occurred at >10% incidence across all breast cancers; however, there were numerous subtype-associated and novel gene mutations including the enrichment of specific mutations in *GATA3*, *PIK3CA* and *MAP3K1* with the luminal A subtype. We identified two novel protein-expression-defined subgroups, possibly produced by stromal/microenvironmental elements, and integrated analyses identified specific signalling pathways dominant in each molecular subtype including a HER2/phosphorylated HER2/EGFR/phosphorylated EGFR signature within the HER2-enriched expression subtype. Comparison of basal-like breast tumours with high-grade serous ovarian tumours showed many molecular commonalities, indicating a related aetiology and similar therapeutic opportunities. The biological finding of the four main breast cancer subtypes caused by different subsets of genetic and epigenetic abnormalities raises the hypothesis that much of the clinically observable plasticity and heterogeneity occurs within, and not across, these major biological subtypes of breast cancer.

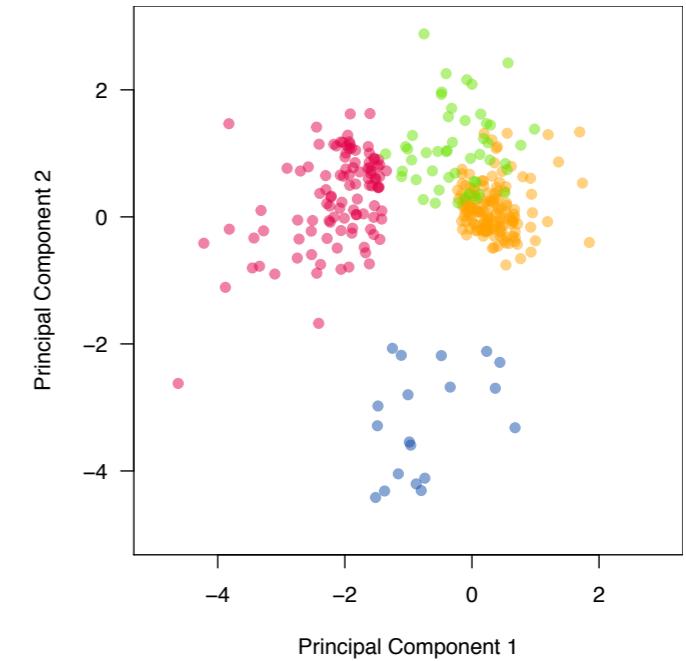
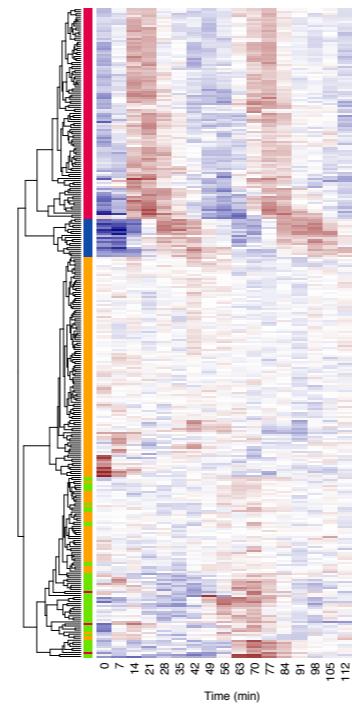
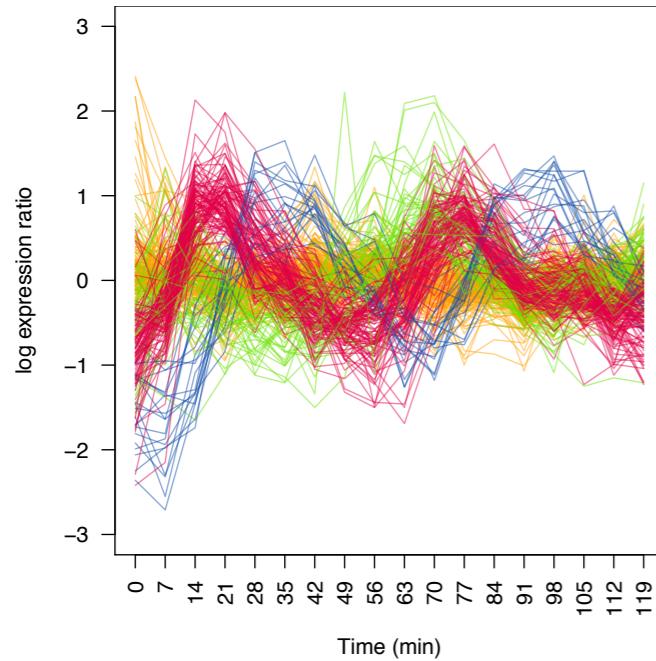


Focus on the Message: Reduce Data Detail



Focus

Choose a Plot Appropriate for your Message

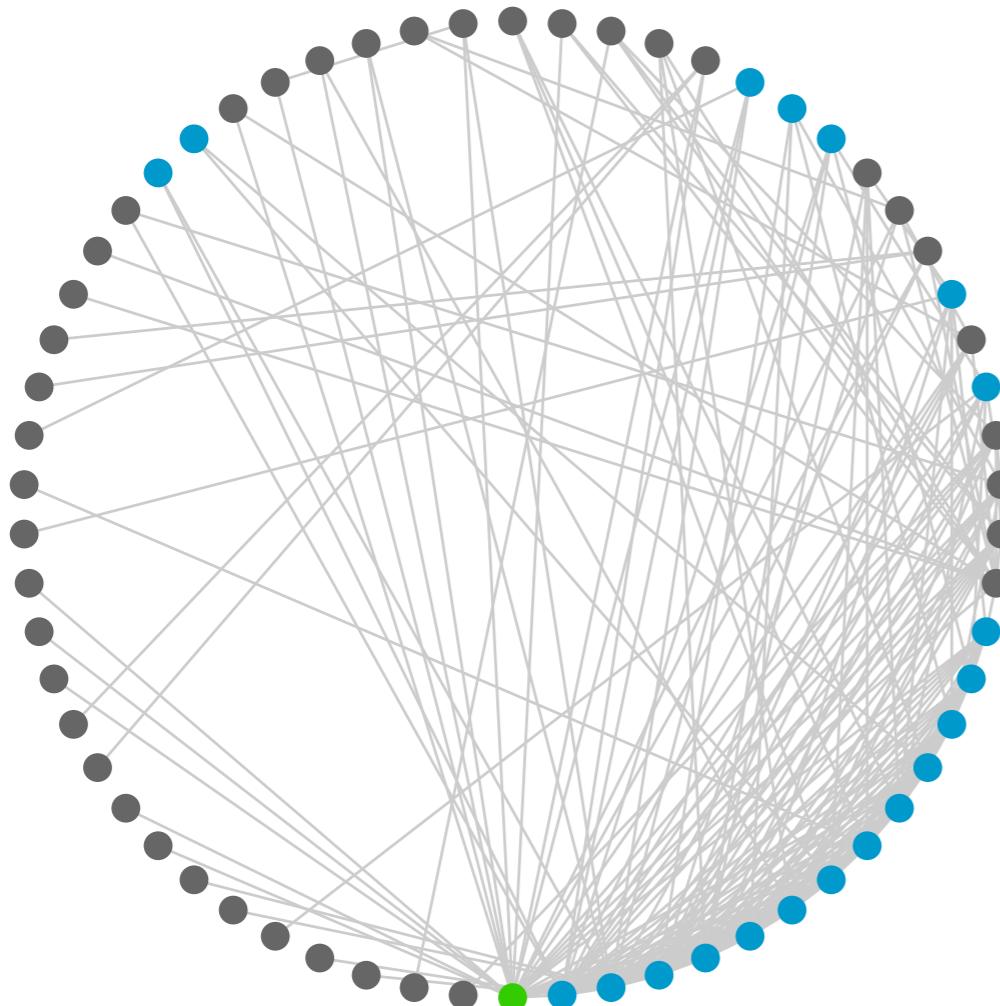


few, high-res

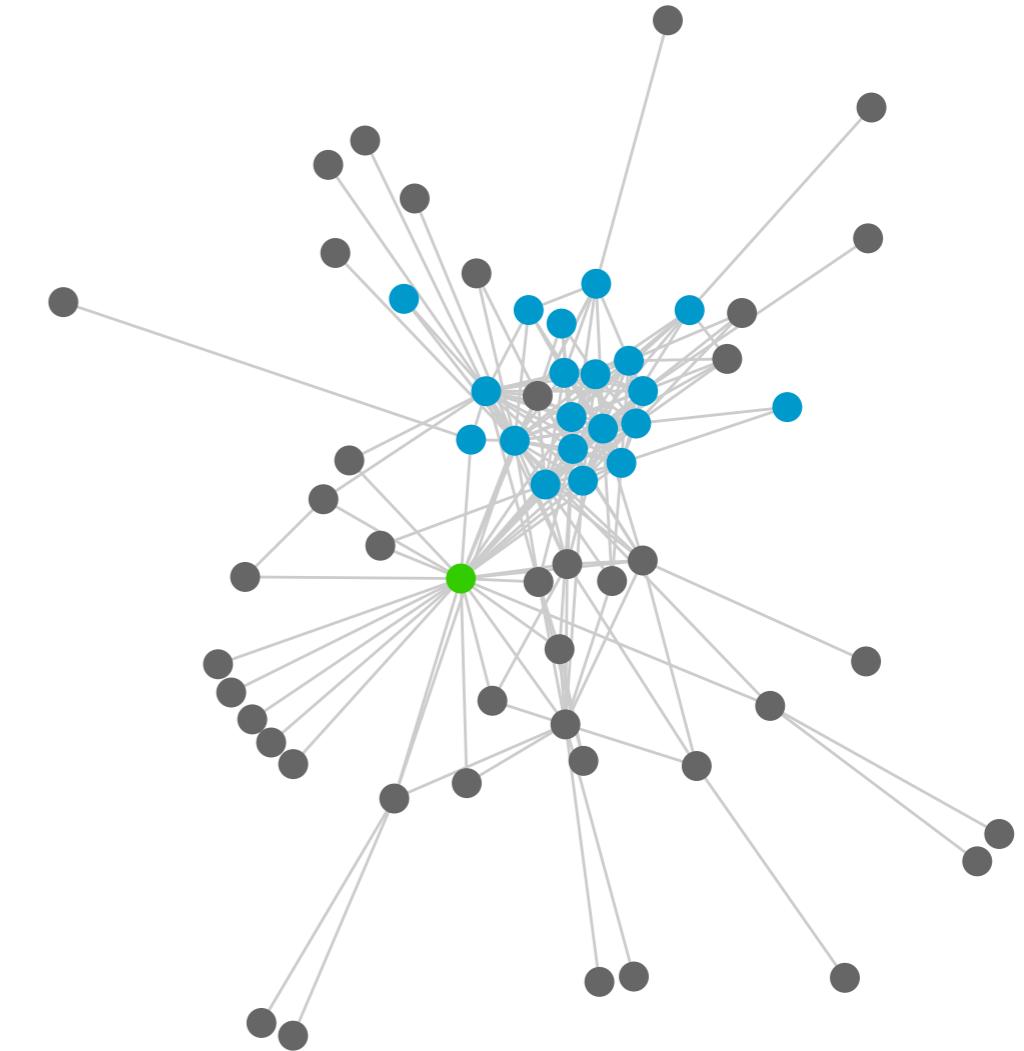


many, low-res

Choose a Plot Appropriate for your Message

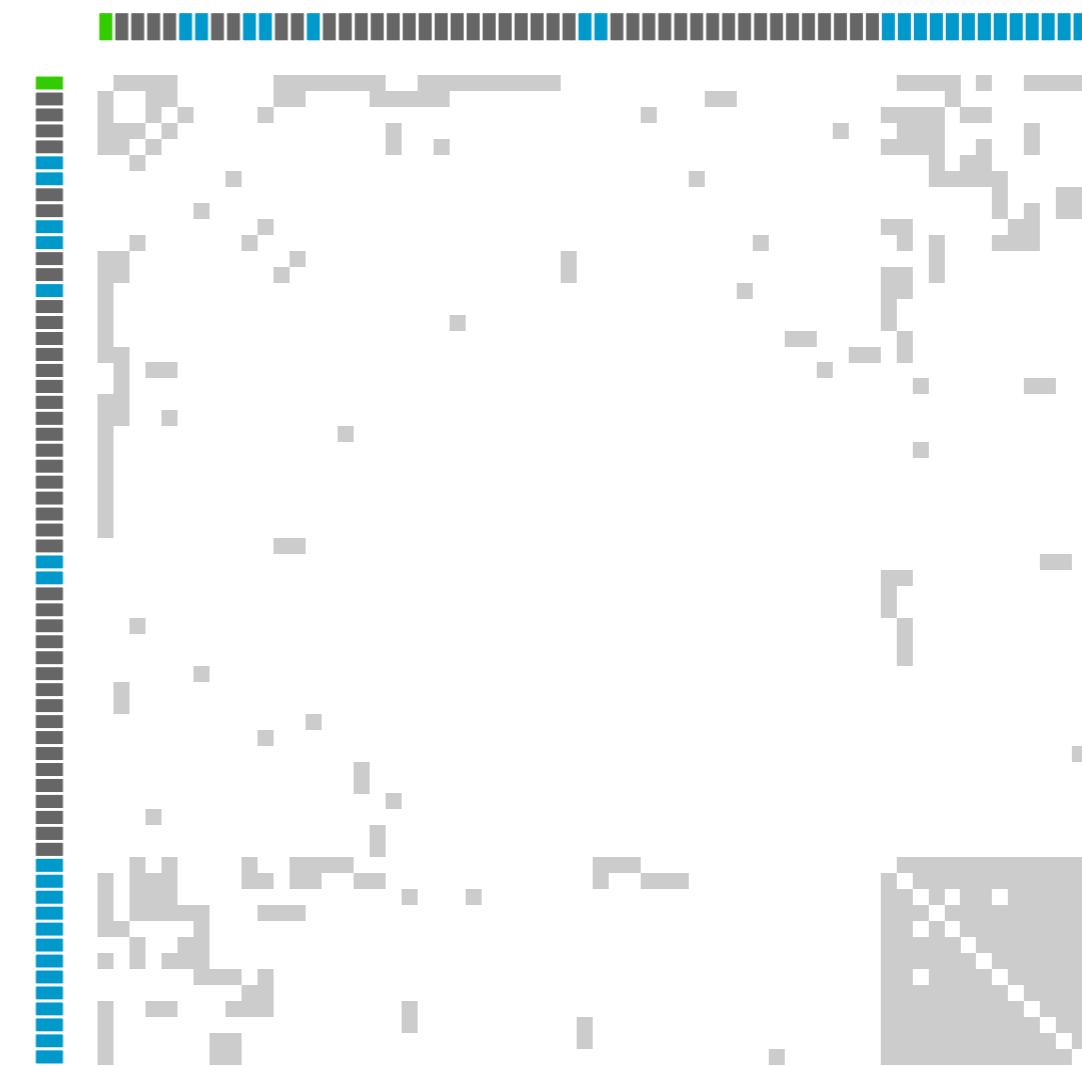
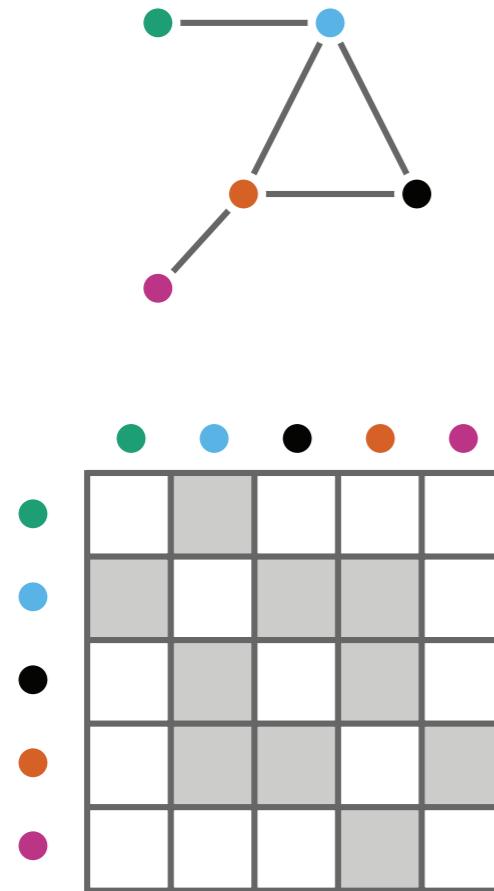


Circular Layout

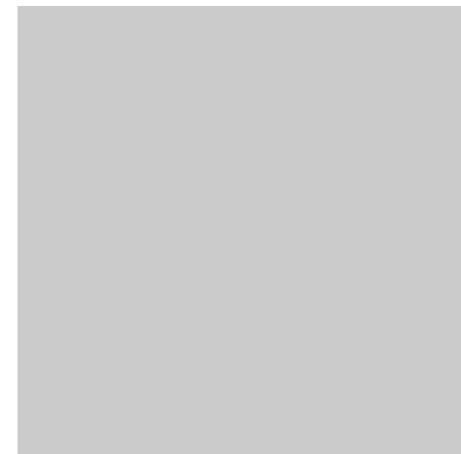
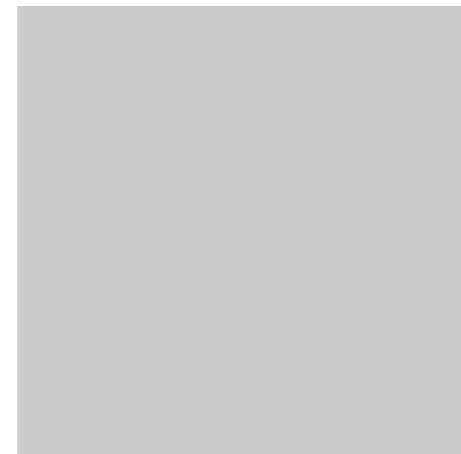
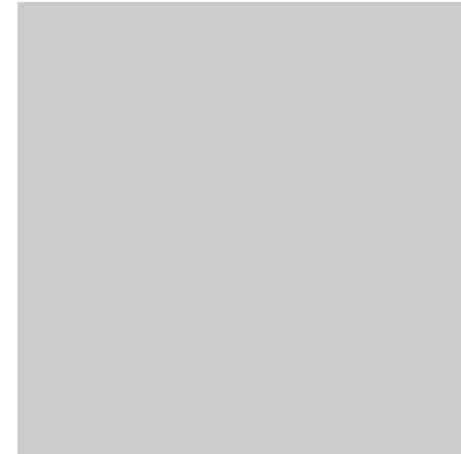
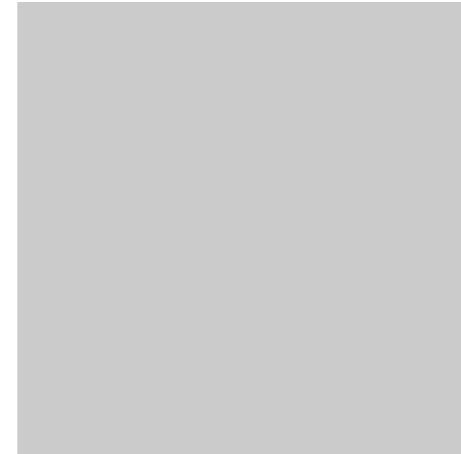


Force-directed Layout

Choose a Plot Appropriate for your Message



Use Layout to Convey Meaning



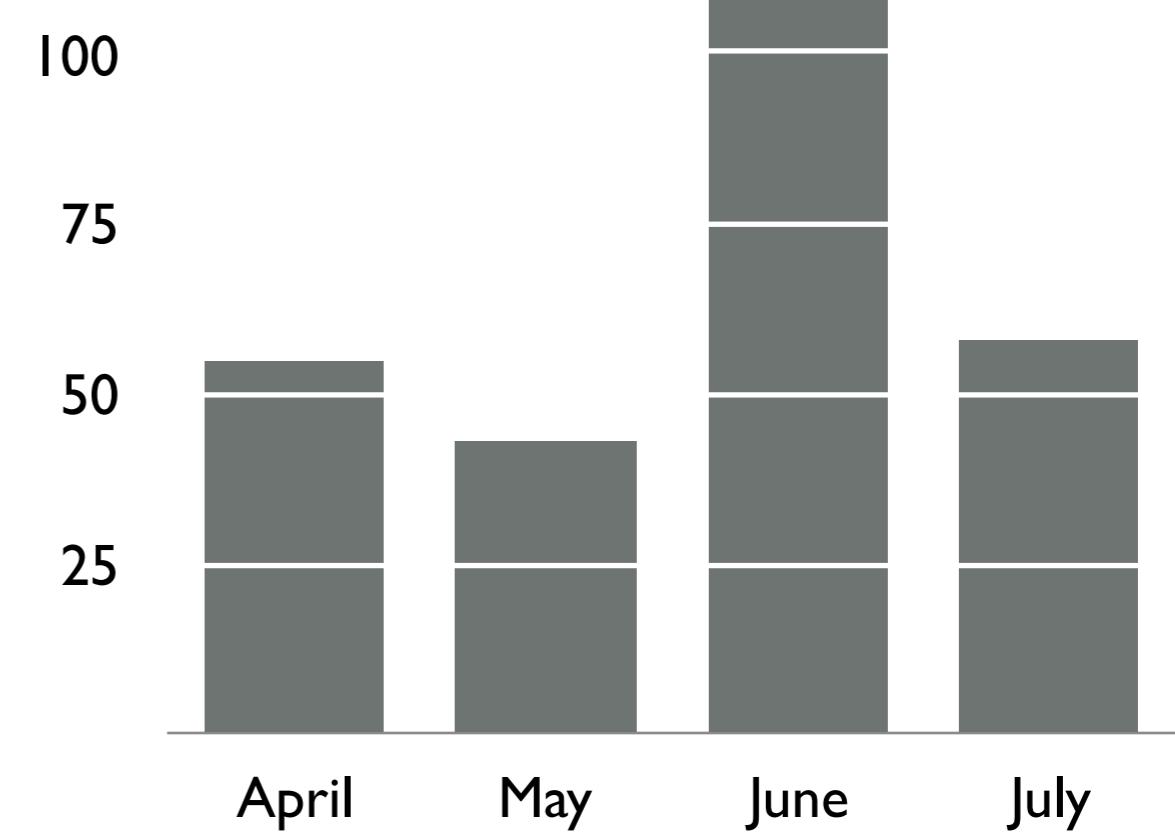
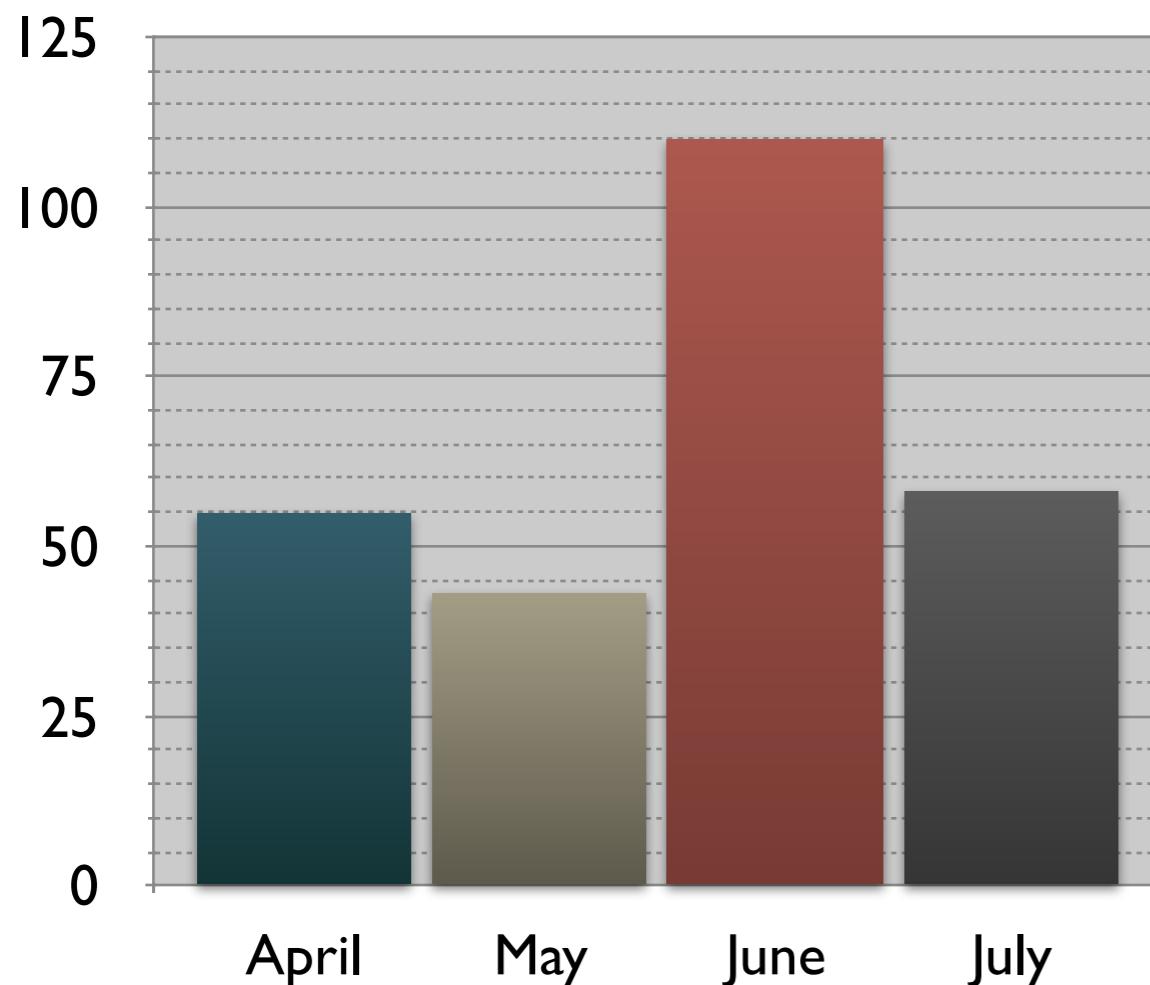
Use Layout to Convey Meaning



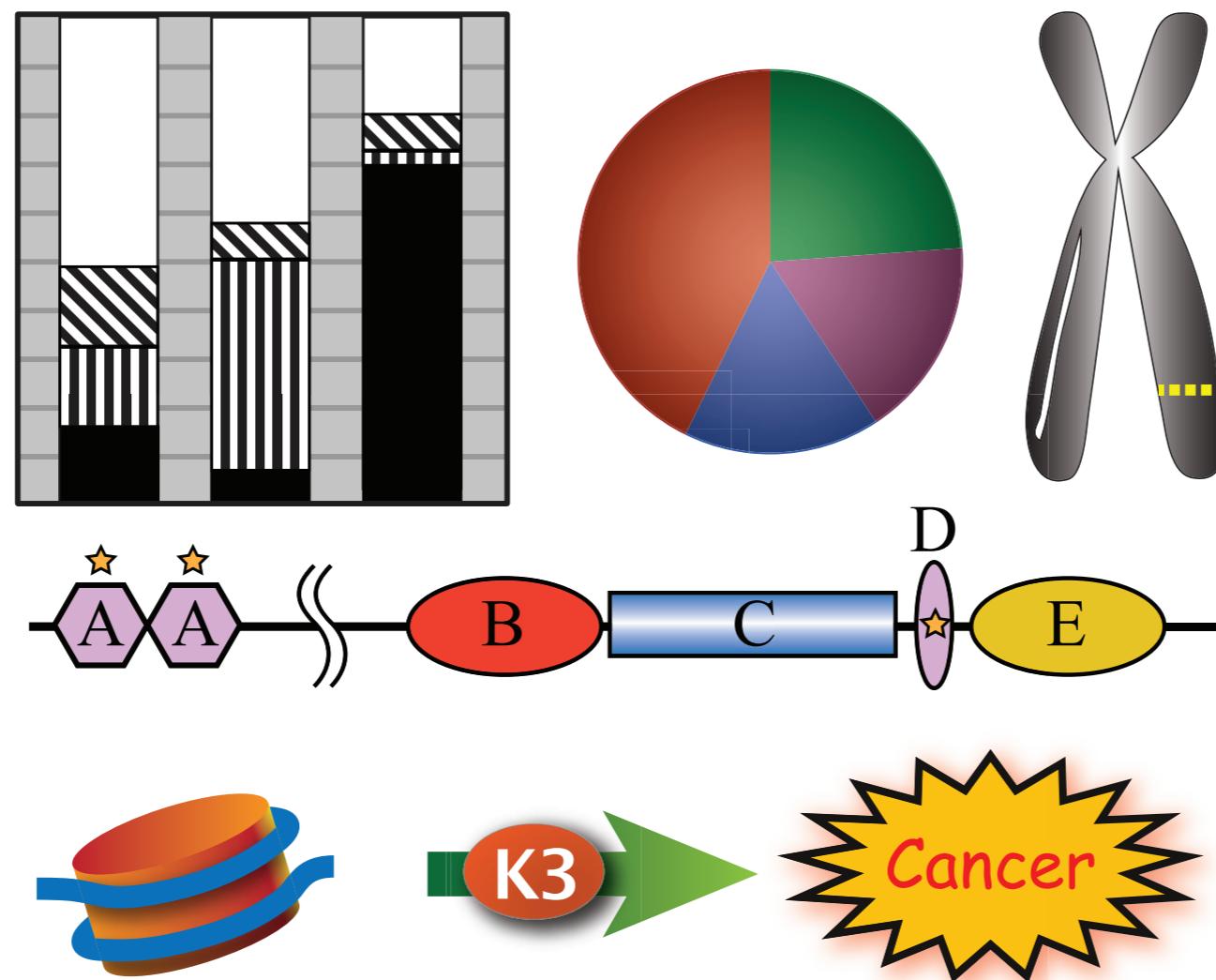
Use Layout to Convey Meaning



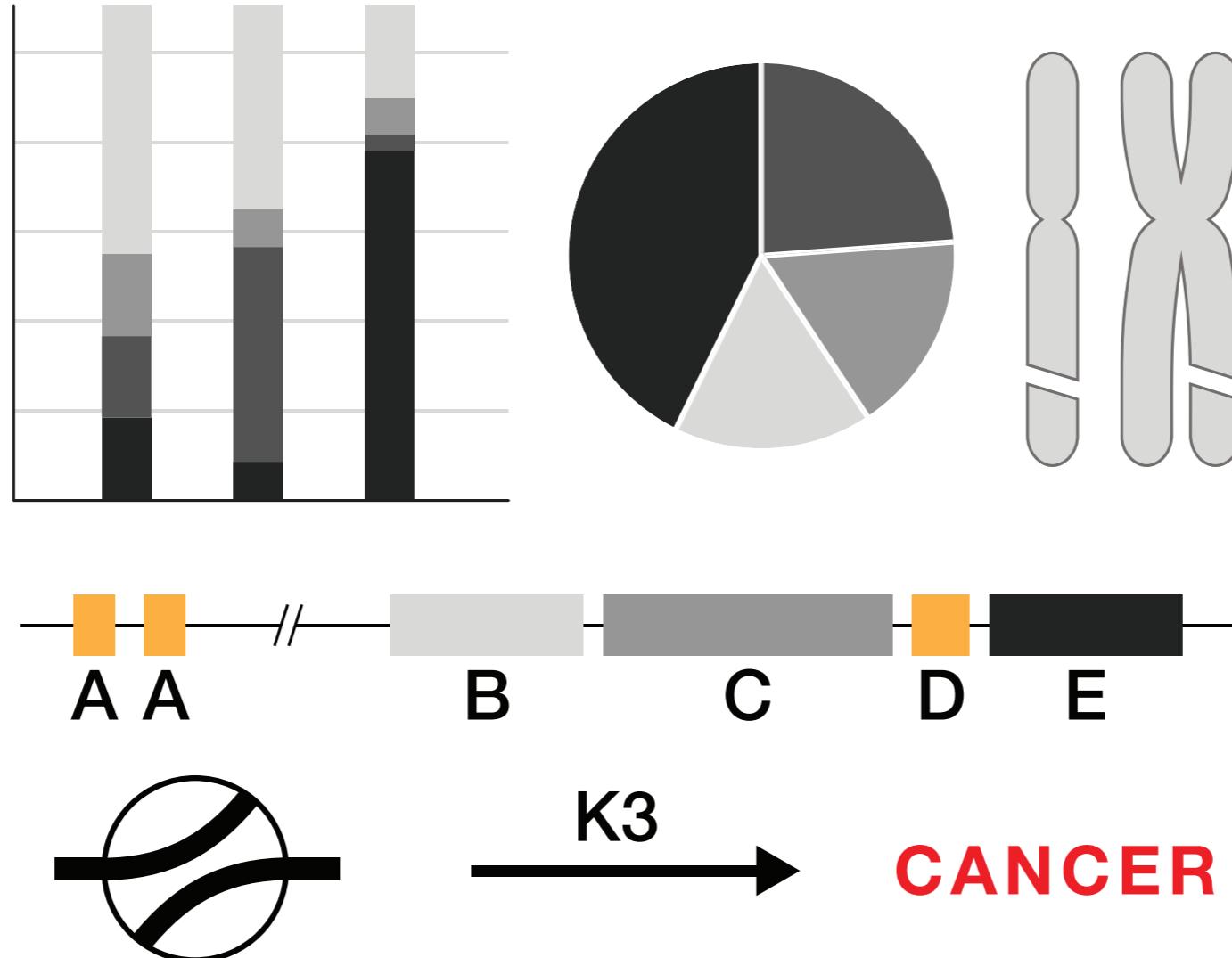
Focus on the Data: Maximize Data/Ink Ratio



Don't shout at your audience!



Don't shout at your audience!

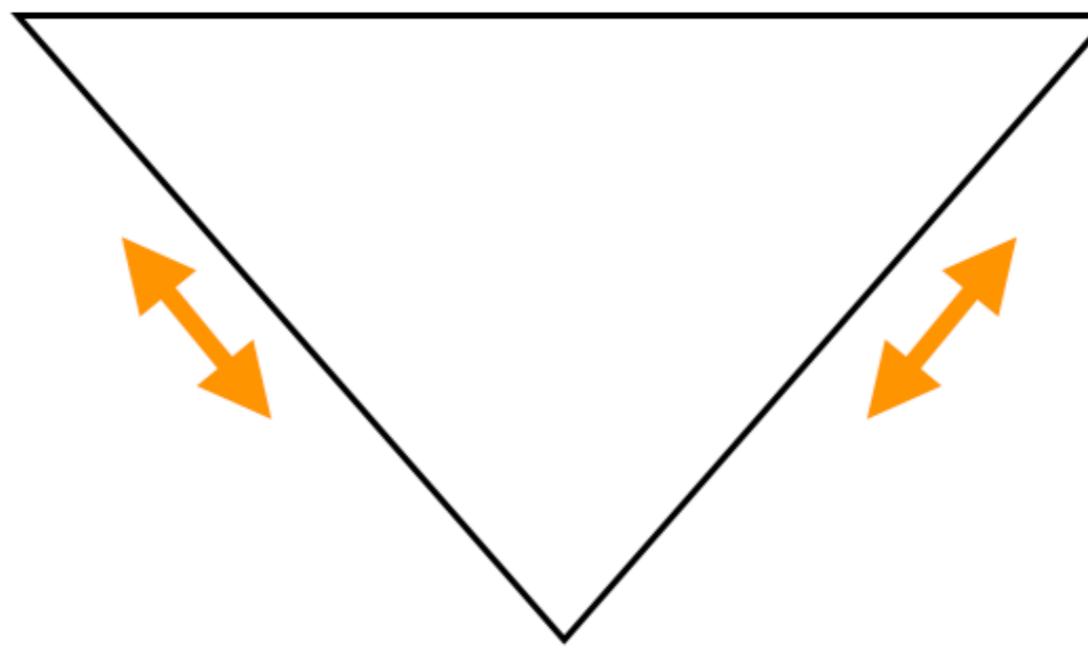


as much data as possible

dialogue between computer & analyst

Confirmation

Exploration



as much data as necessary

Data Visualization Tools & Communication

The BD2K Guide to the Fundamentals of Data Science Series

Nils Gehlenborg, PhD

Department of Biomedical Informatics, Harvard Medical School

nils@hms.harvard.edu · <http://gehlenborglab.org>

BD2K Lecture - 31 March 2017