

Data Visualization for Biomedical Applications

Lecture 3

BMI706 - 5 April 2018

Nils Gehlenborg, PhD

Administrative

- Submission deadline for Homework 2 assignment: Sunday, April 9 at 11:59 pm Eastern
- Next assignment is part of the project, it doesn't count towards your homework grade
- Nils will be traveling on April 18: Peter Kerpedjiev and remote lecture

Talk Announcements

- Friday, April 6, 1:00pm to 6:00pm - Data Science in Biomedicine Symposium

Fernanda Viegas & Martin Wattenberg (Google)

Visualization: The secret weapon for machine learning

Dana-Farber Cancer Institute

- <http://bcb.dfci.harvard.edu/index.php/calendar/dfci-fstrf-marvin-zelen-symposium>

Talk Announcements

- Friday, April 20, 12:30pm to 2:00pm

Ben Shneiderman (University of Maryland)

Interactive Visual Discovery in Event Analytics: Electronic Health Records and Other Applications

Harvard Geological Museum (Room 100)
24 Oxford Street, Cambridge MA 02138

- <https://www.seas.harvard.edu/calendar/event/111371>

Review of Session 2

Visualization Pitfalls

- Color
- Interactions between encodings
- 3D Visualization
- Animation

Single View Interactions

Manipulate

⌚ Change over Time



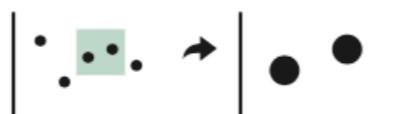
🔍 Select



ניווט (Navigate)

→ Item Reduction

→ Zoom
Geometric or Semantic



→ Pan/Translate



→ Constrained



→ Attribute Reduction

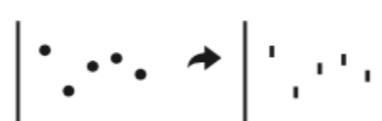
→ Slice

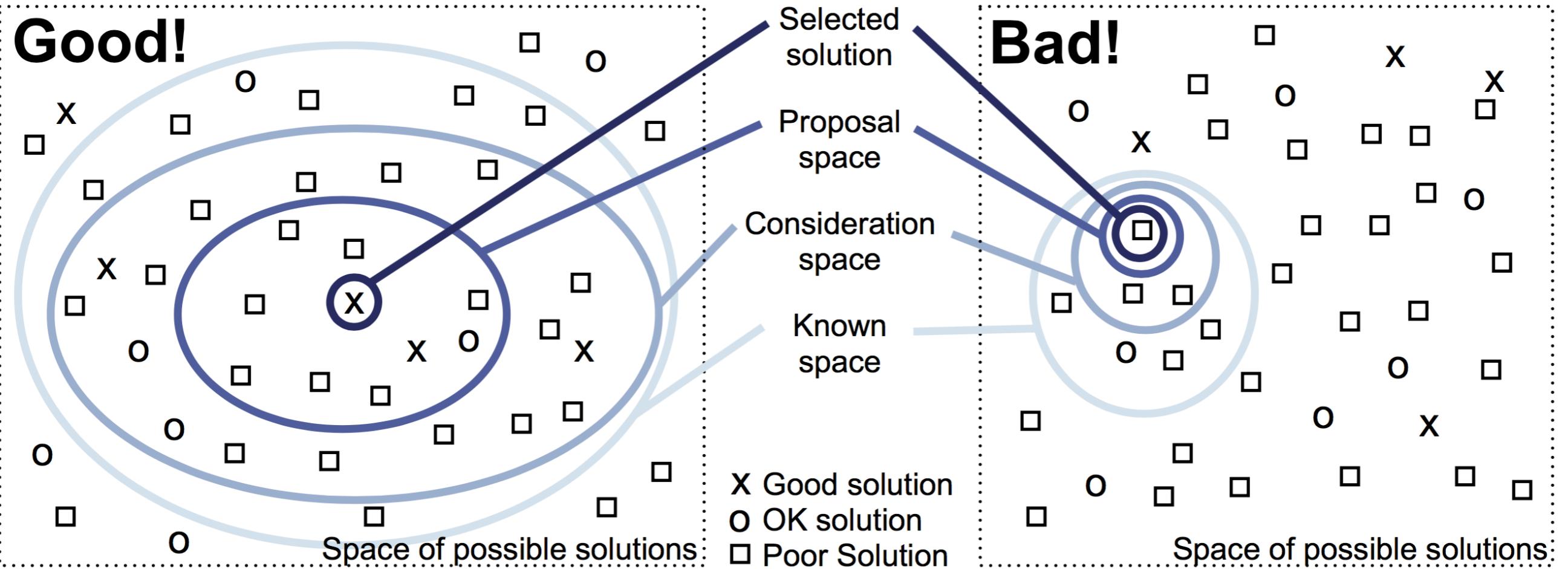


→ Cut

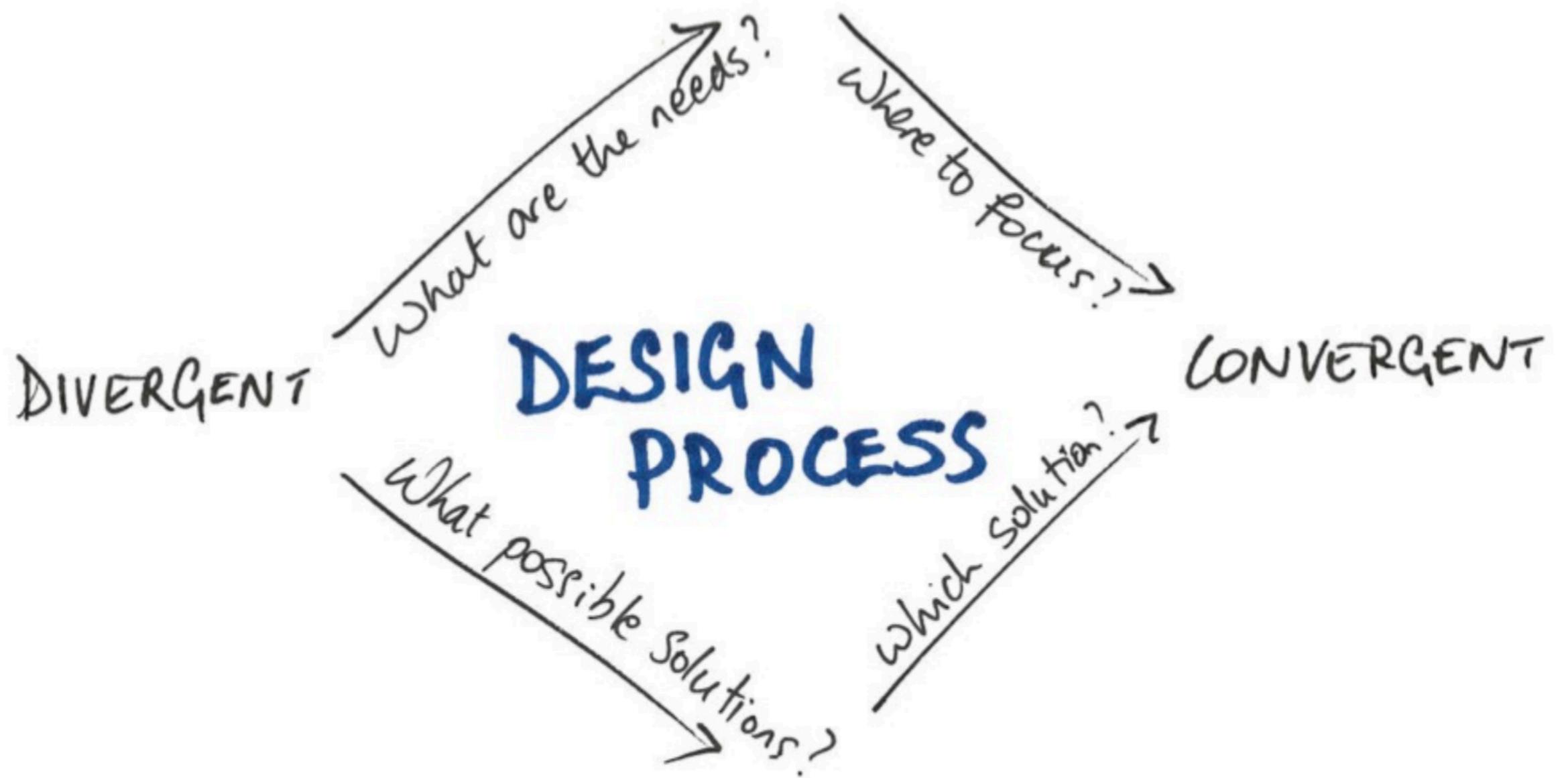


→ Project

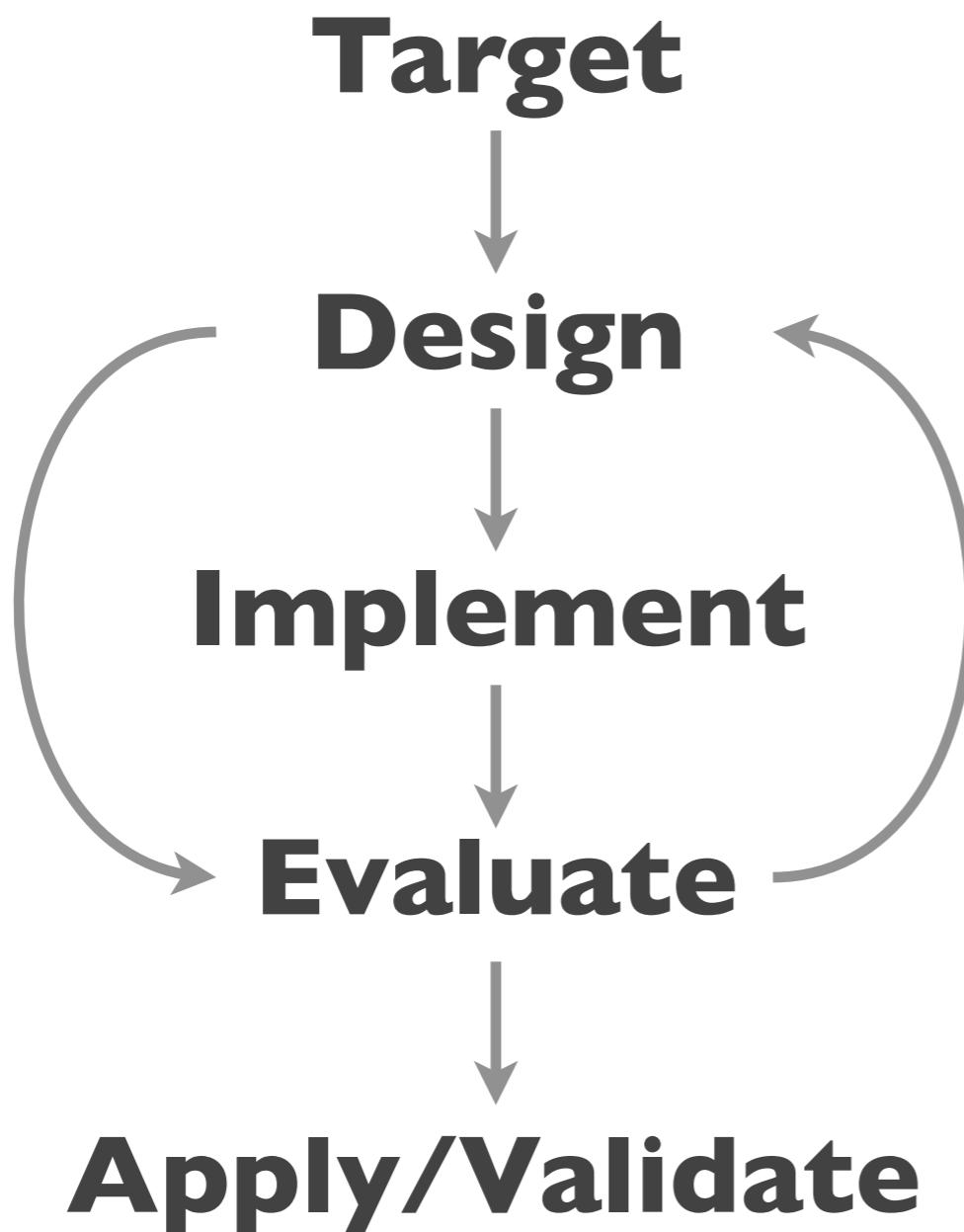




Five Design-Sheet Methodology



User-Centered Participatory Design



user-centered design
usability engineering
participatory design

Tasks?

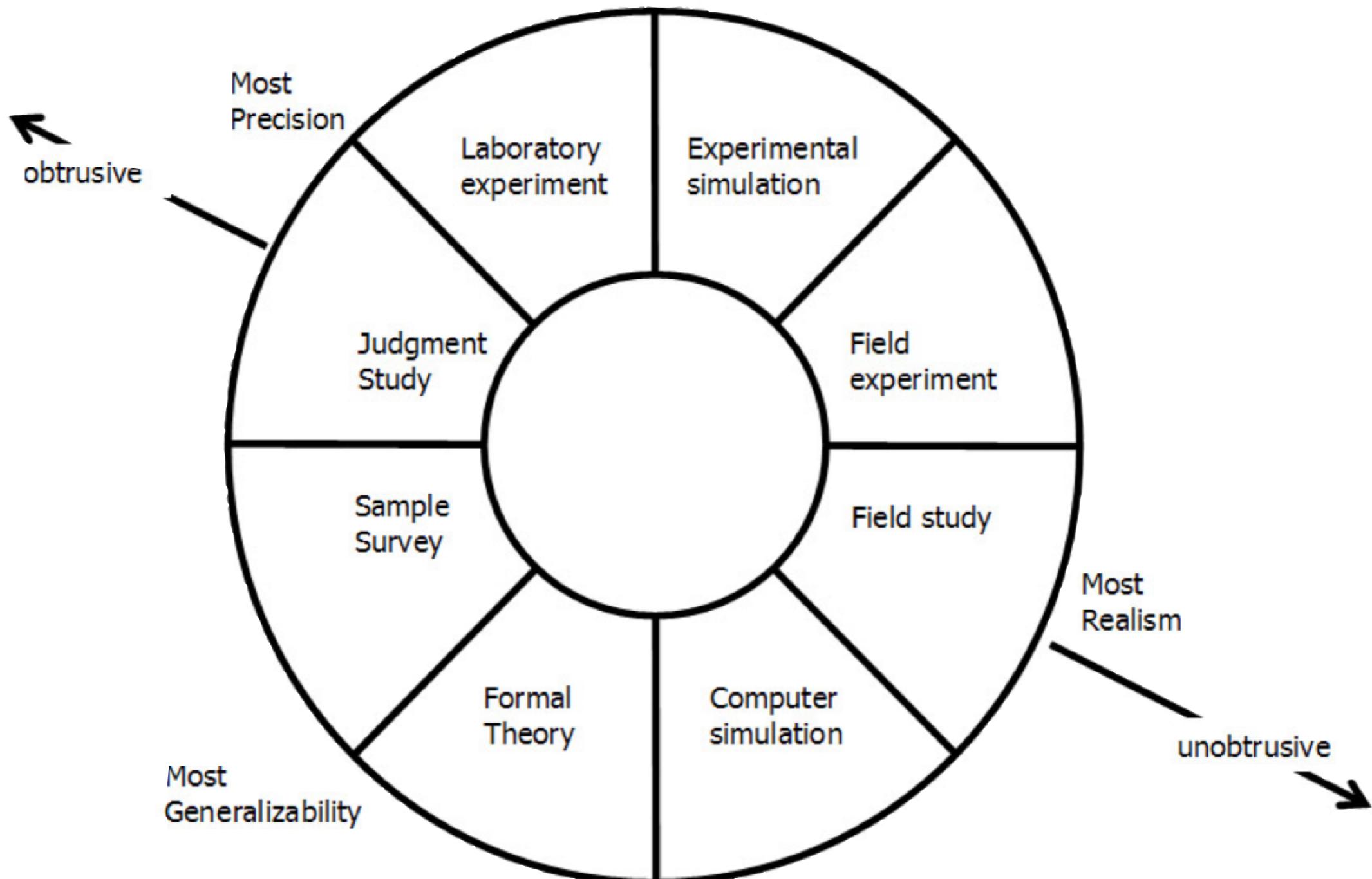
Evaluate != Validate?

Lecture 3: Learning Goals

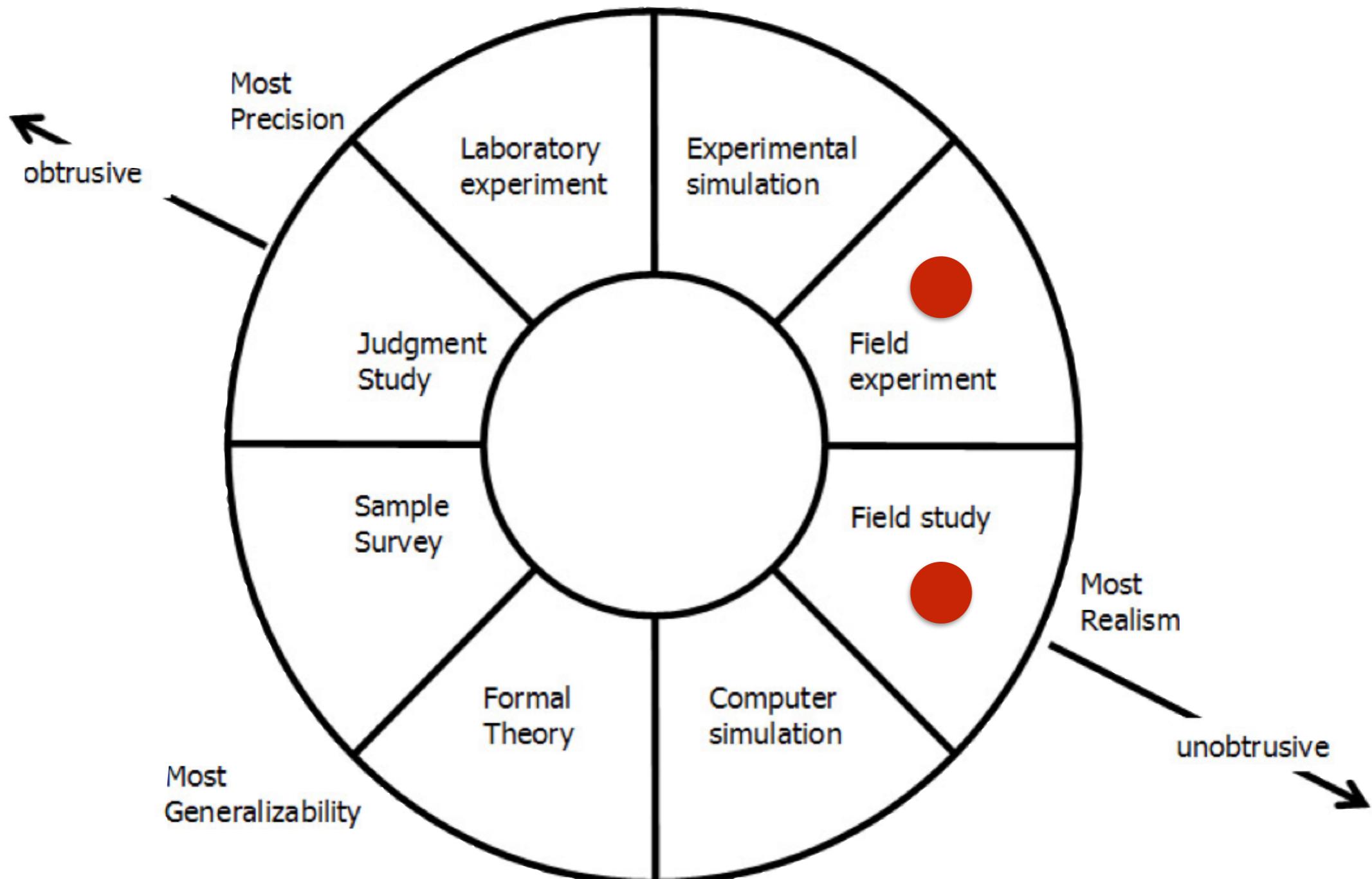
- How can we evaluate a visualization method or tool?
- What are common techniques for multivariate data visualization?
- What is the long tail of data visualization?

Evaluation

Validation Techniques



Validation Techniques



Qualitative Usability Testing

- **Field study**
 - Intended to be unobtrusive
 - Conducted in the actual situation (a user with their own analysis tasks and data)
 - Watching someone use your tool can be very informative; you may be surprised!

Qualitative Usability Testing

- **Field experiment**
 - Design user tasks to simulate real analyses
 - Recruit group of users and arrange one-on-one sessions
 - Encourage thinking aloud and note top usability issues

Qualitative Usability Testing

- Learn through observation
- Learn about behavior, not just opinions
- Not about proof, but about insight and context
- Test early and test often
- **Important that your users are assured that you are testing the software, not them**

Quantitative User Testing

- Error Rate
- Time (see LineUp paper next week)
- Insights
 - **Oh boy!**

See required reading for Week 5

Gremlin: an interactive visualization model for analyzing genomic rearrangements

O'Brien TM, Ritz AM, Raphael BJ, Laidlaw DH

IEEE Trans. Vis. Comput. Graph., 2010

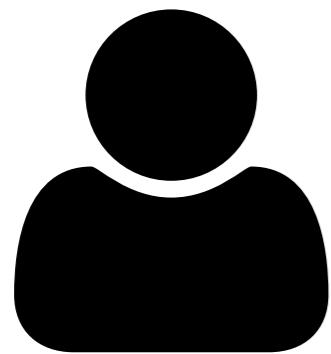
Evaluation

What is an effective visualization and what is not?

Considerations - e.g. measuring effectiveness
quantitatively is extremely problematic

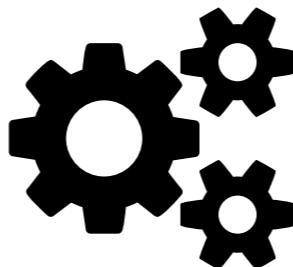
Evaluation

1



human aspect

2



system complexity

3

? vs !

outcome is question,
not answer

Evaluation

What is an effective visualization and what is not?

Considerations - e.g. measuring effectiveness quantitatively is extremely problematic

Impact - e.g. on developers (are we doing good work?), users (is this a good tool?), and reviewers (should this be published/funded?), e.g. in evaluating visual against analytical approaches

Tabular Data

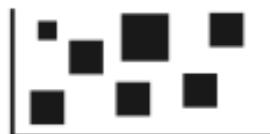
Arrange Tables

→ Express Values



→ Separate, Order, Align Regions

→ Separate



→ Order



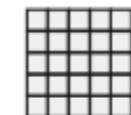
→ Align



→ 1 Key
List



→ 2 Keys
Matrix



→ 3 Keys
Volume



→ Many Keys
Recursive Subdivision



→ Axis Orientation

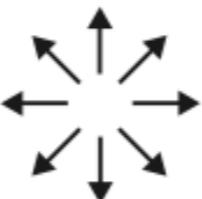
→ Rectilinear



→ Parallel

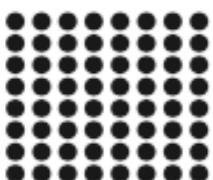


→ Radial



→ Layout Density

→ Dense



→ Space-Filling



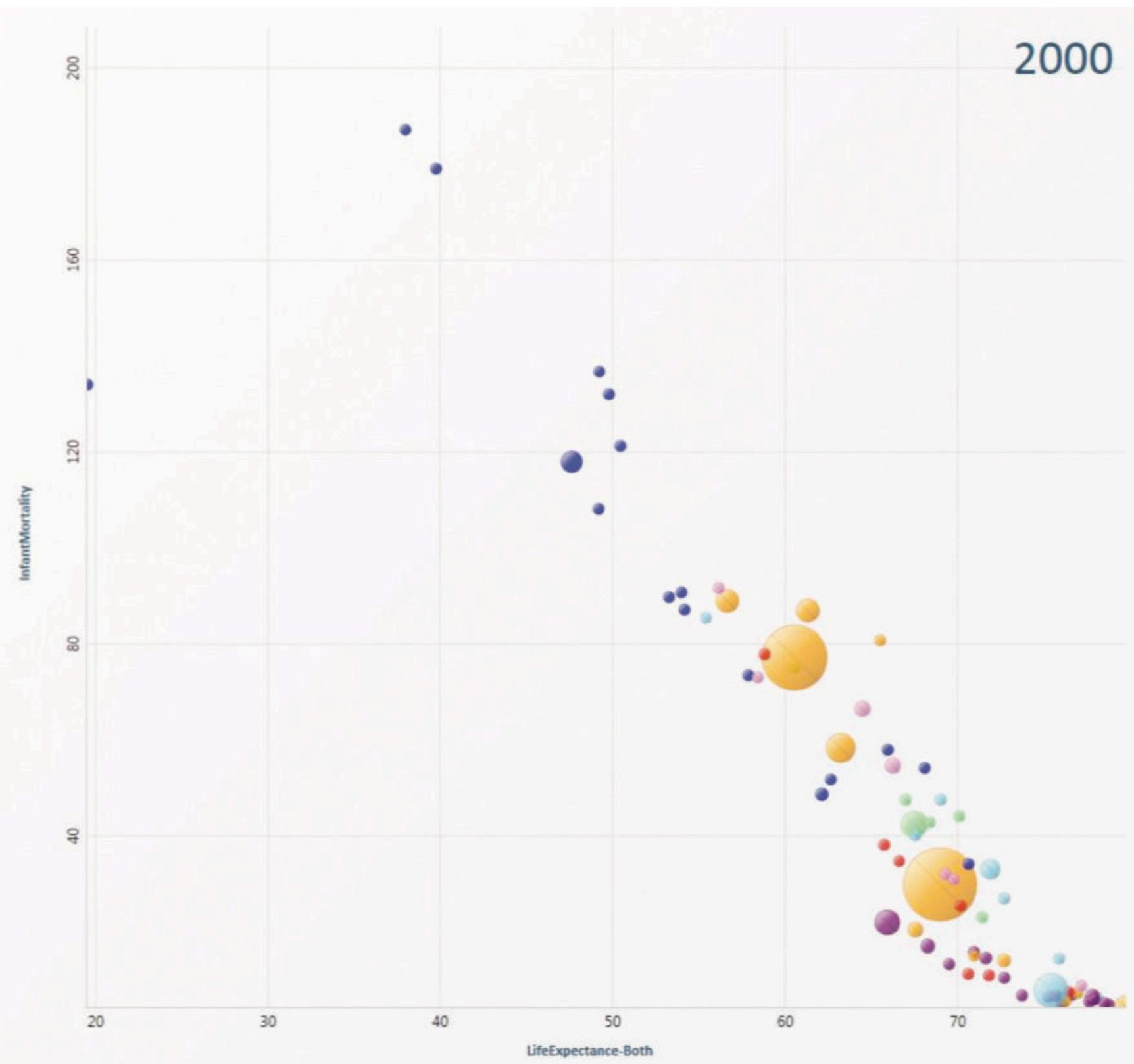
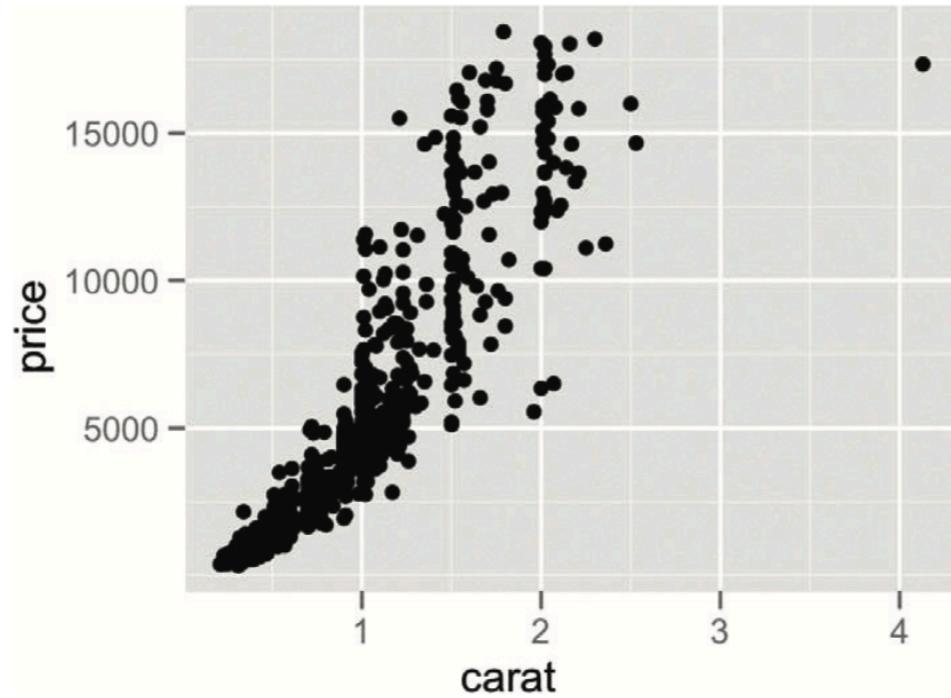
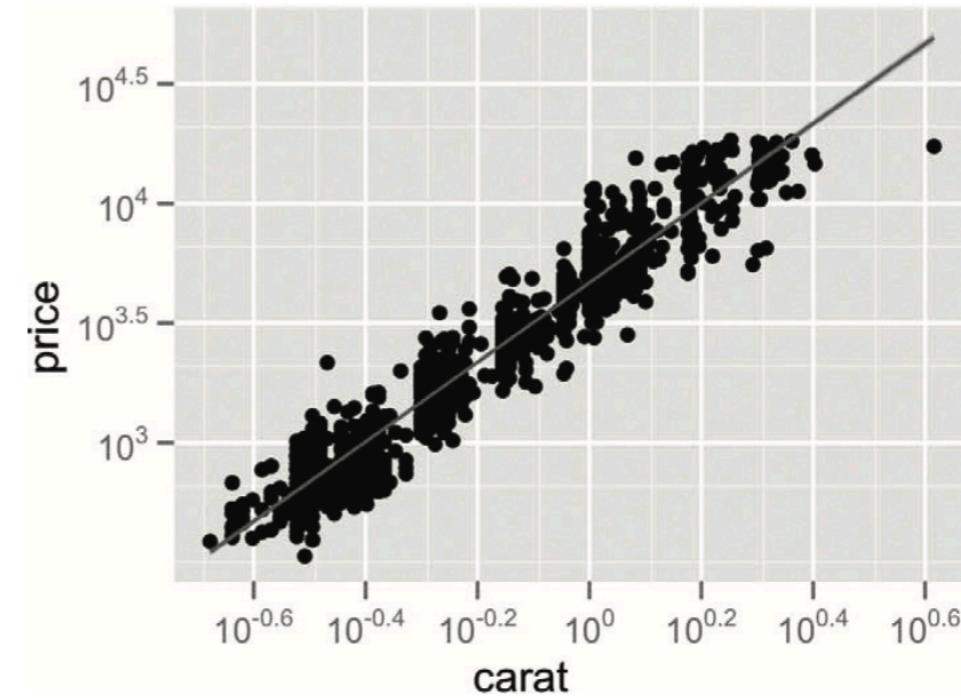


Figure 7.2. Scatterplot. Each point mark represents a country, with horizontal and vertical spatial position encoding the primary quantitative attributes of life expectancy and infant mortality. The color channel is used for the categorical country attribute and the size channel for quantitative population attribute. From [Robertson et al. 08, Figure 1c].

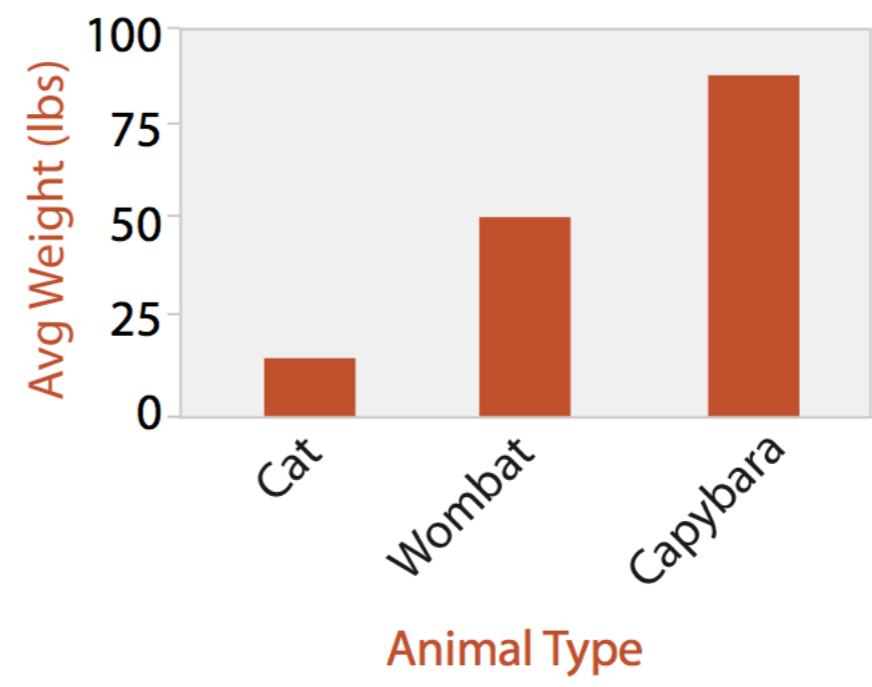
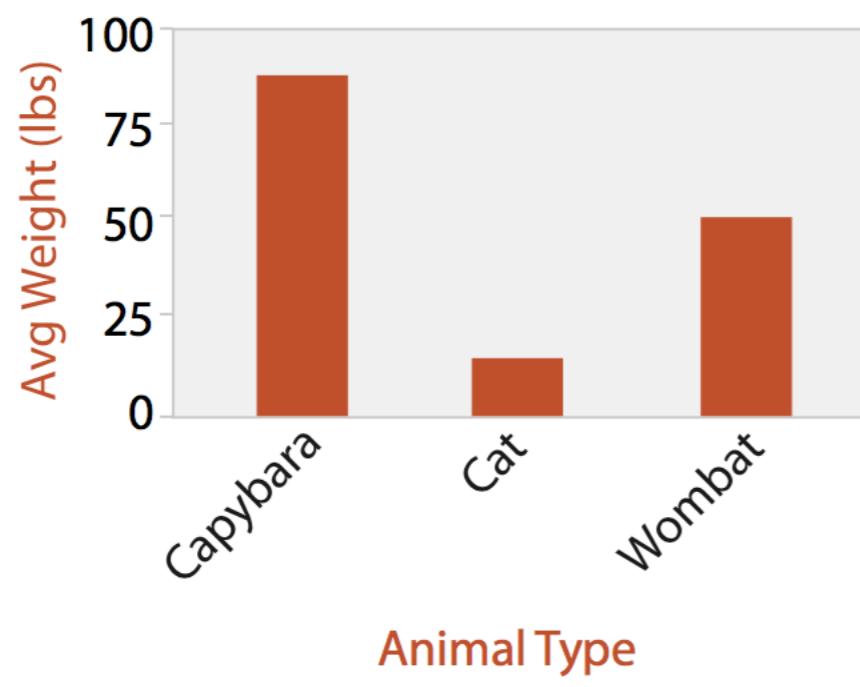


(a)



(b)

Figure 7.3. Scatterplots. (a) Original diamond price/carat data. (b) Derived log-scale attributes are highly positively correlated. From [Wickham 10, Figure 10].



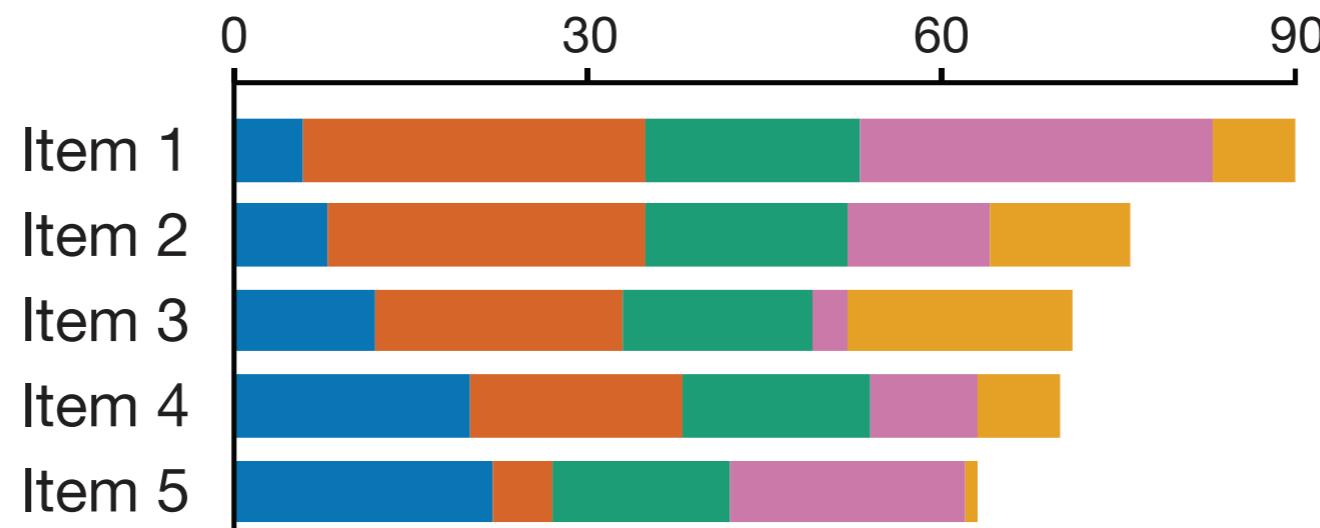
Bar Charts for Items & Categories

	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 Items & 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

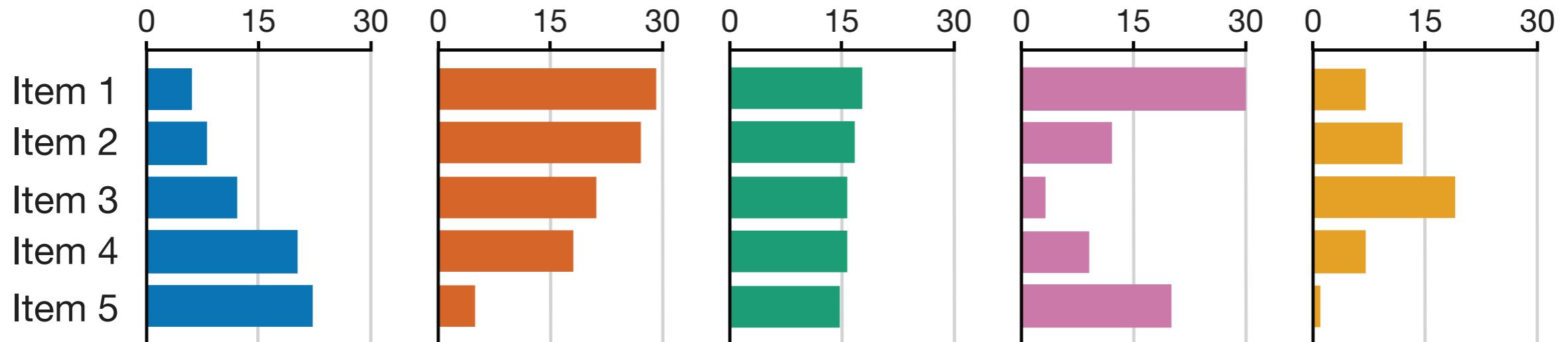
Stacked Bar Chart



Bar Charts for Items & 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 Items & 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

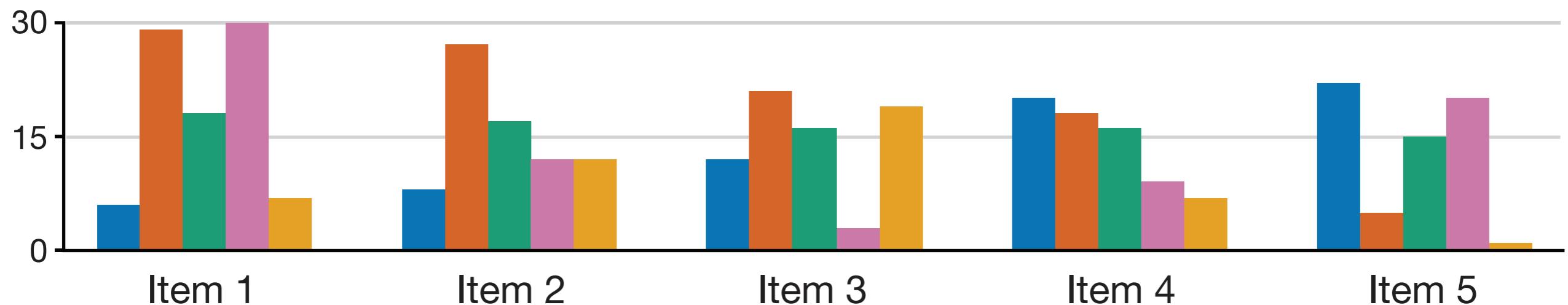


Table Lens

Table Lens: Baseball Player Statistics

Calculate: "Hits" / "At Bats" = "Avg"

	Avg	Career Avg	Team	Salary 87
Larry Herndon	0.24734983	0.27282876	Det.	225
Jesse Barfield	0.2886248	0.27268818	Tor.	1237.5
Jeffrey Leonar	0.27859238	0.27260458	S.F.	900
Donnie Hill	0.28318584	0.2725564	Oak.	275
Billy Sample	0.285	0.2718601	Atl.	NA
Howard Johnson	0.24545455	0.25232068	N.Y.	297.5
Andres Thomas	0.250774	0.2521994	Atl.	75
Billy Hatcher	0.25775656	0.25211507	Hou.	110
Omar Moreno	0.2339833	0.2518029	Atl.	NA
Darnell Coles	0.2725528	0.25153375	Det.	105

Row 304: Mike Lavalliere; Column 20: Put Outs Value: 468 810 -- 2163

Multivariate Data

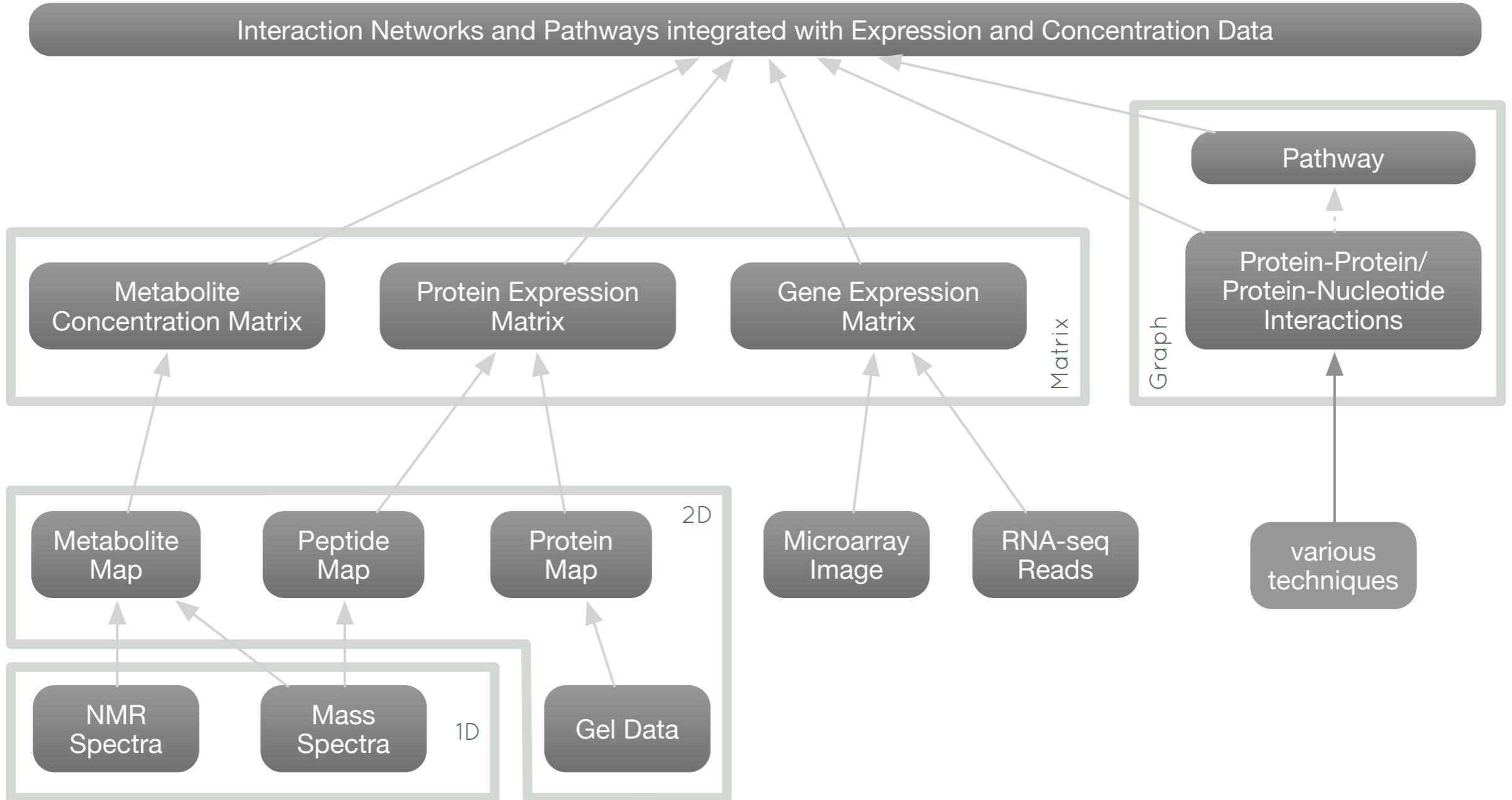
Homogeneous Tables

Multivariate Data

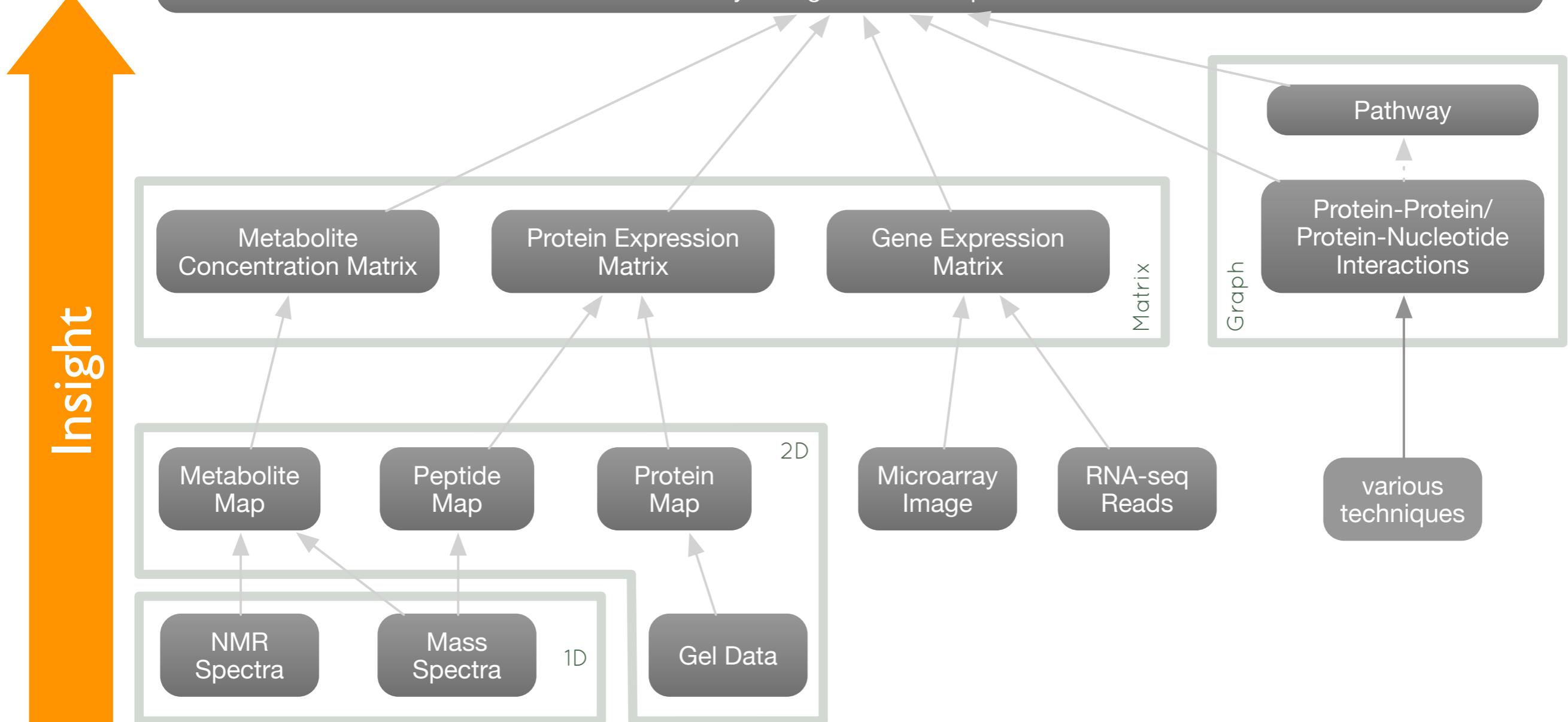
- typical “omics” data: transcriptomics, proteomics, metabolomics
- expression/concentration levels of many biological entities (transcripts, proteins, etc.) across many different conditions/time points
- entity levels measured per sample on a “genome-wide” scale
- often entities are not measured directly



Multivariate Data



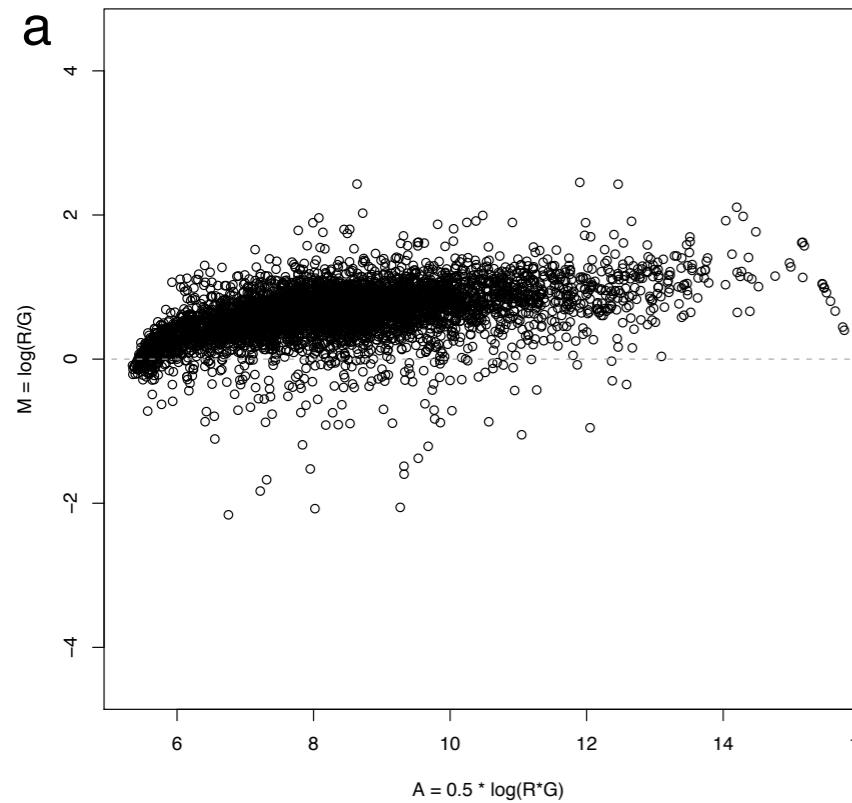
Multivariate Data



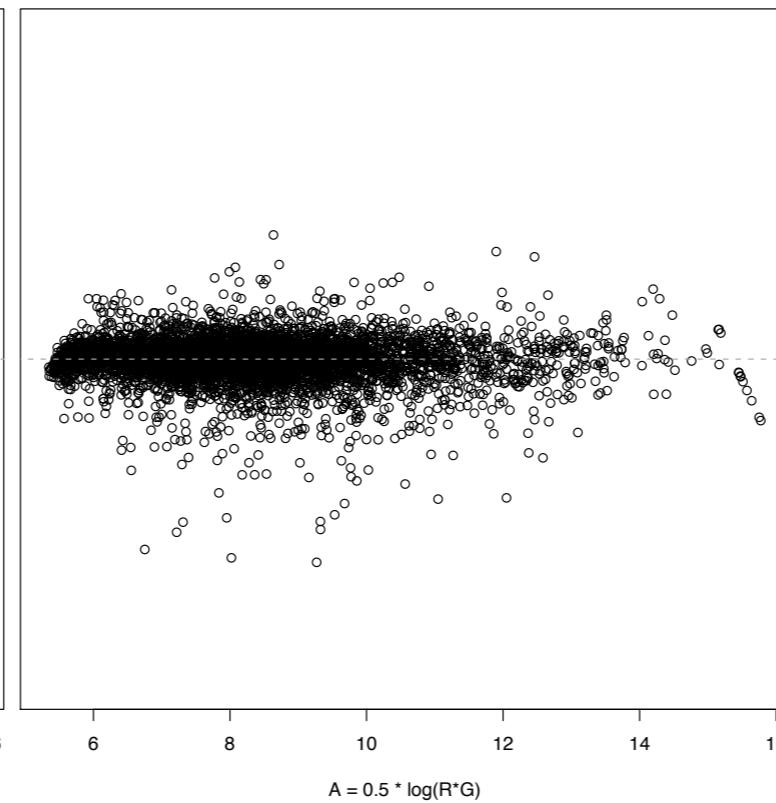
Multivariate Data: Transcriptomics

MA Plot: 1 array

before normalization

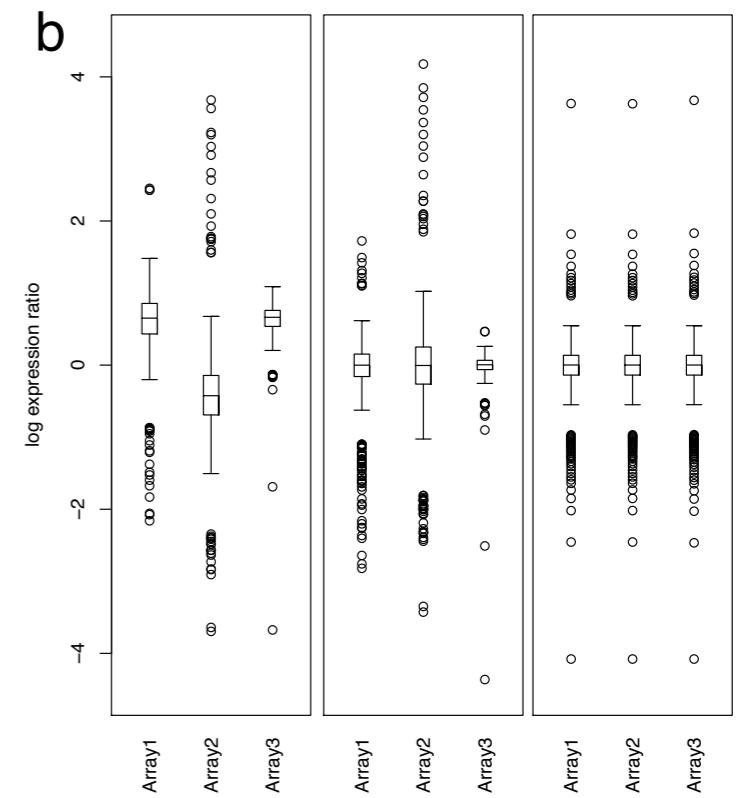


after normalization

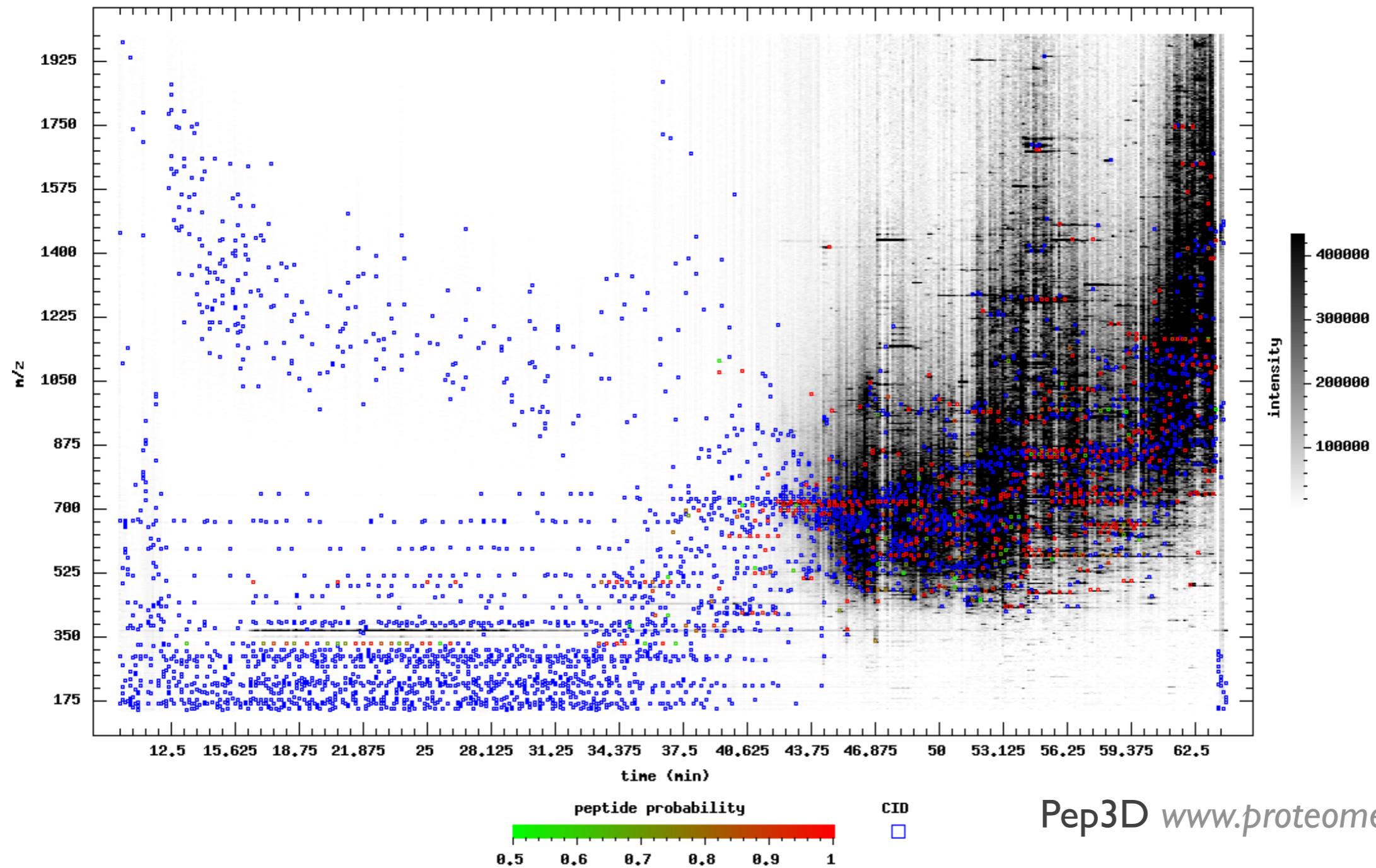


Box Plot: 3 arrays

1 2 3

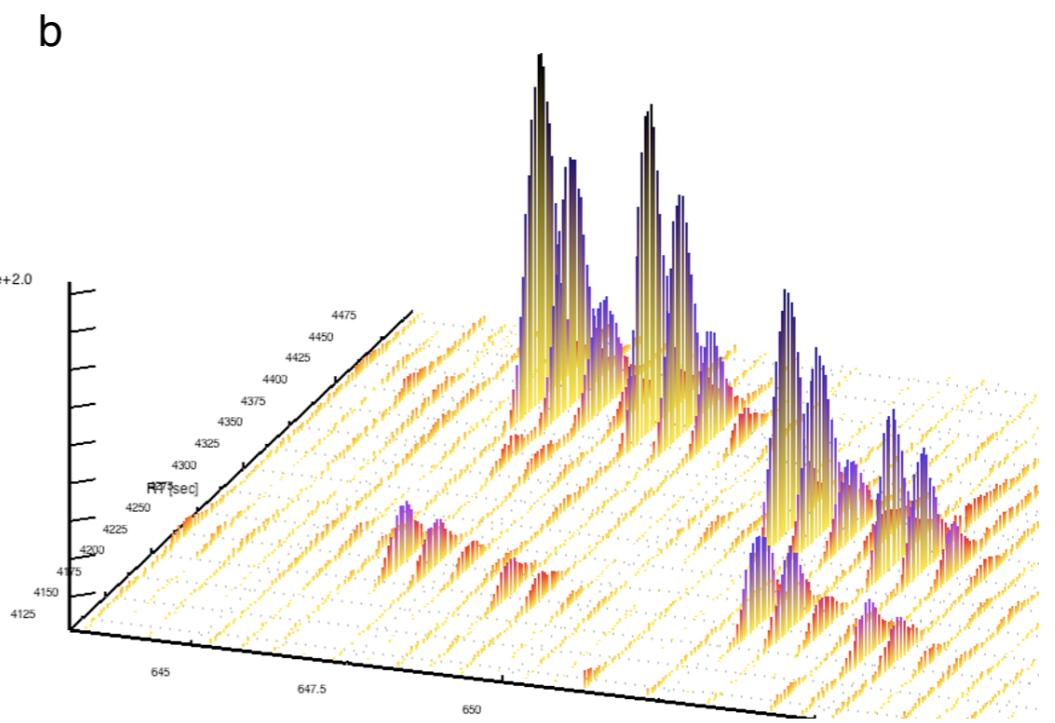
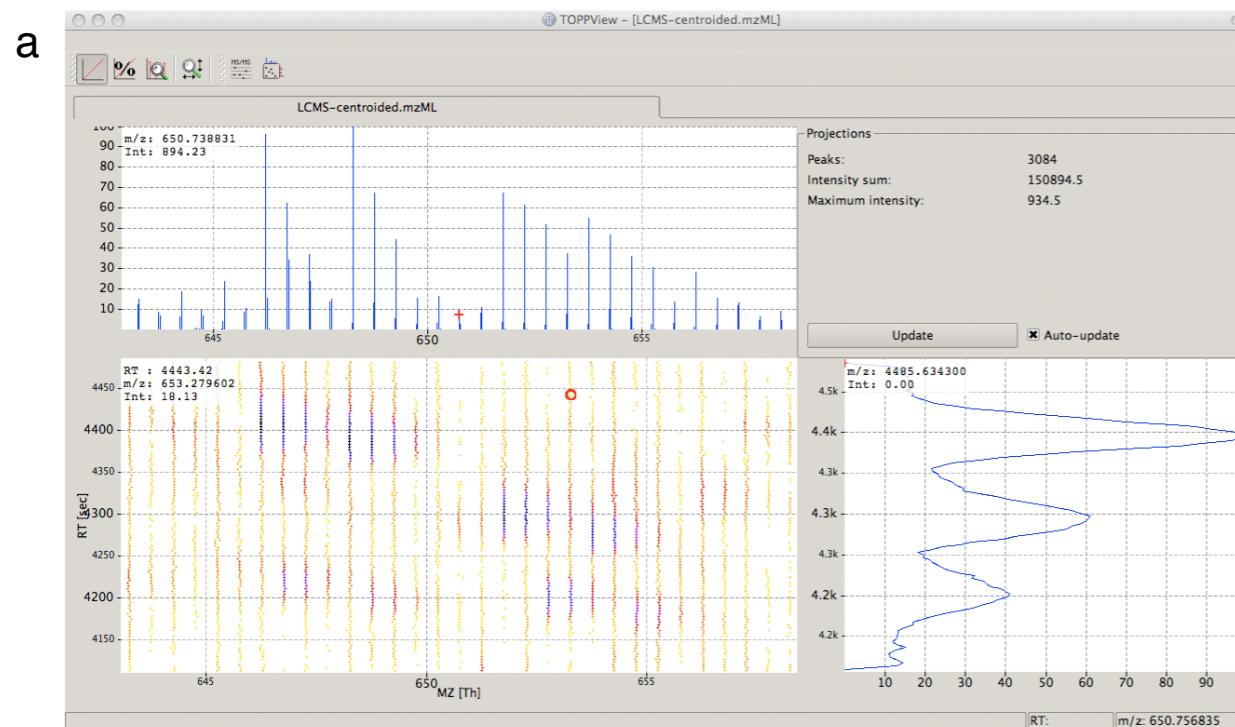


1 = before normalization
 2 = after within-array normalization
 3 = after between-array normalization



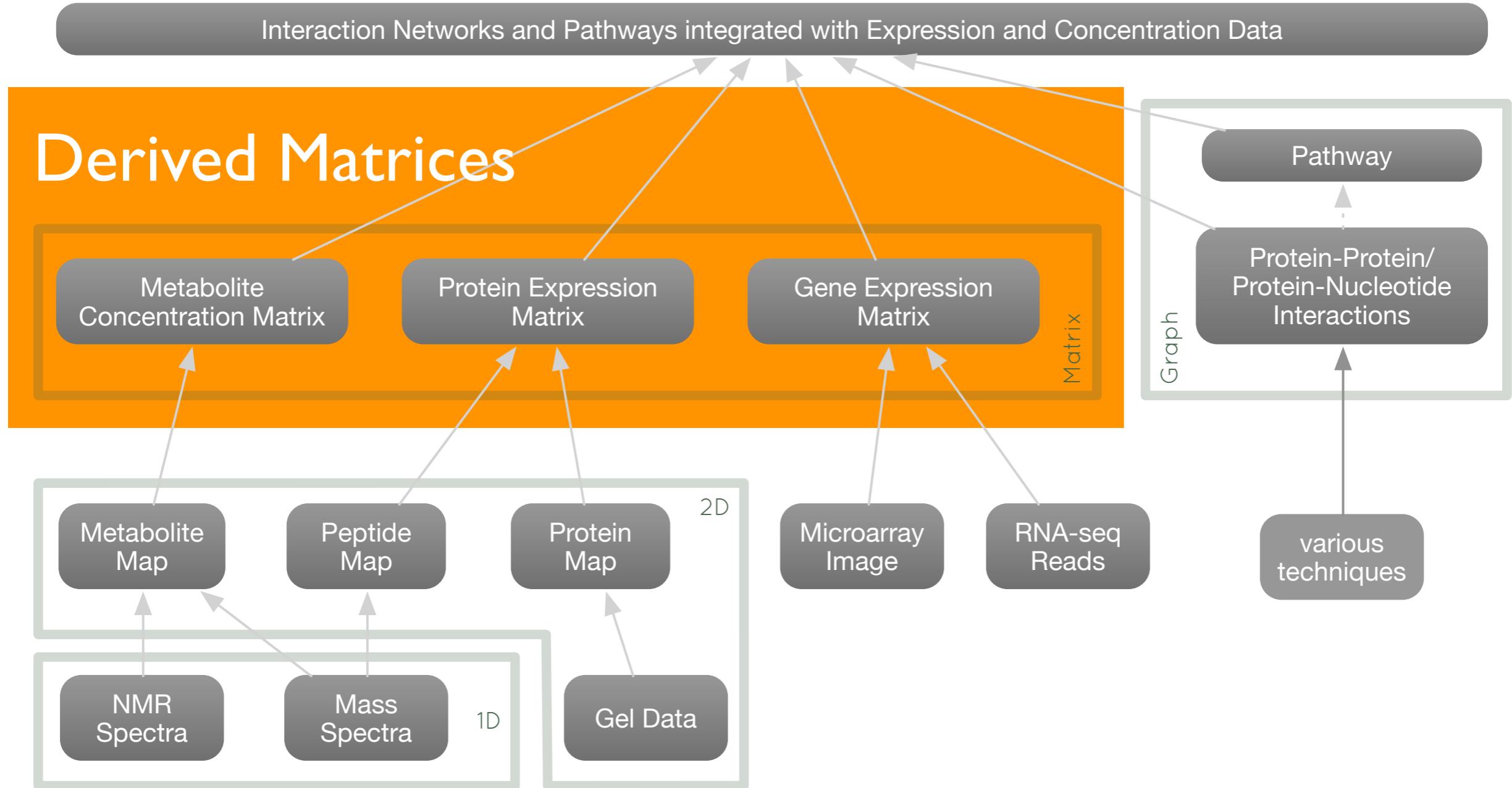
Pep3D www.proteomecenter.org

Multivariate Data: Proteomics



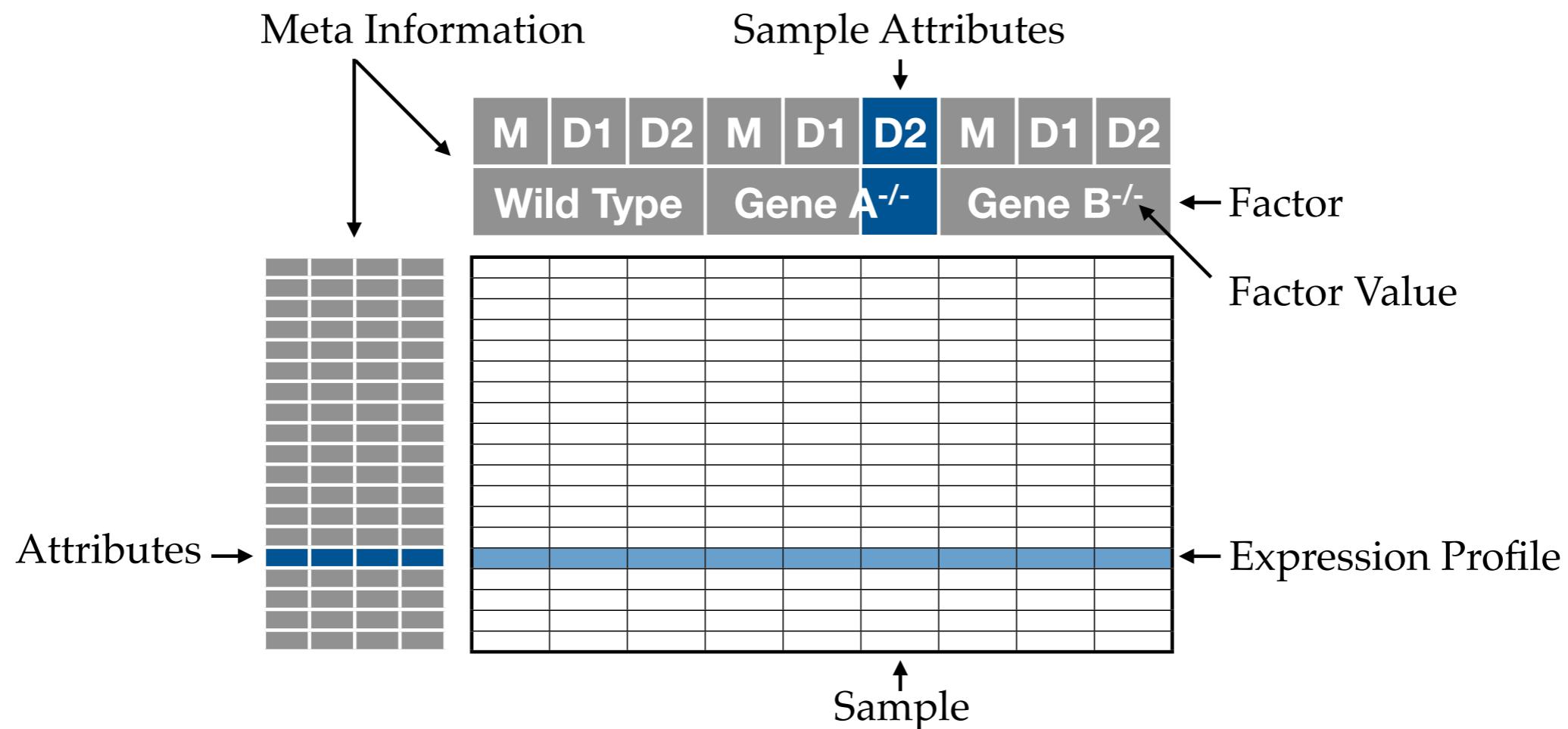
TOPPView www.open-ms.de

Multivariate Data: Derived Data



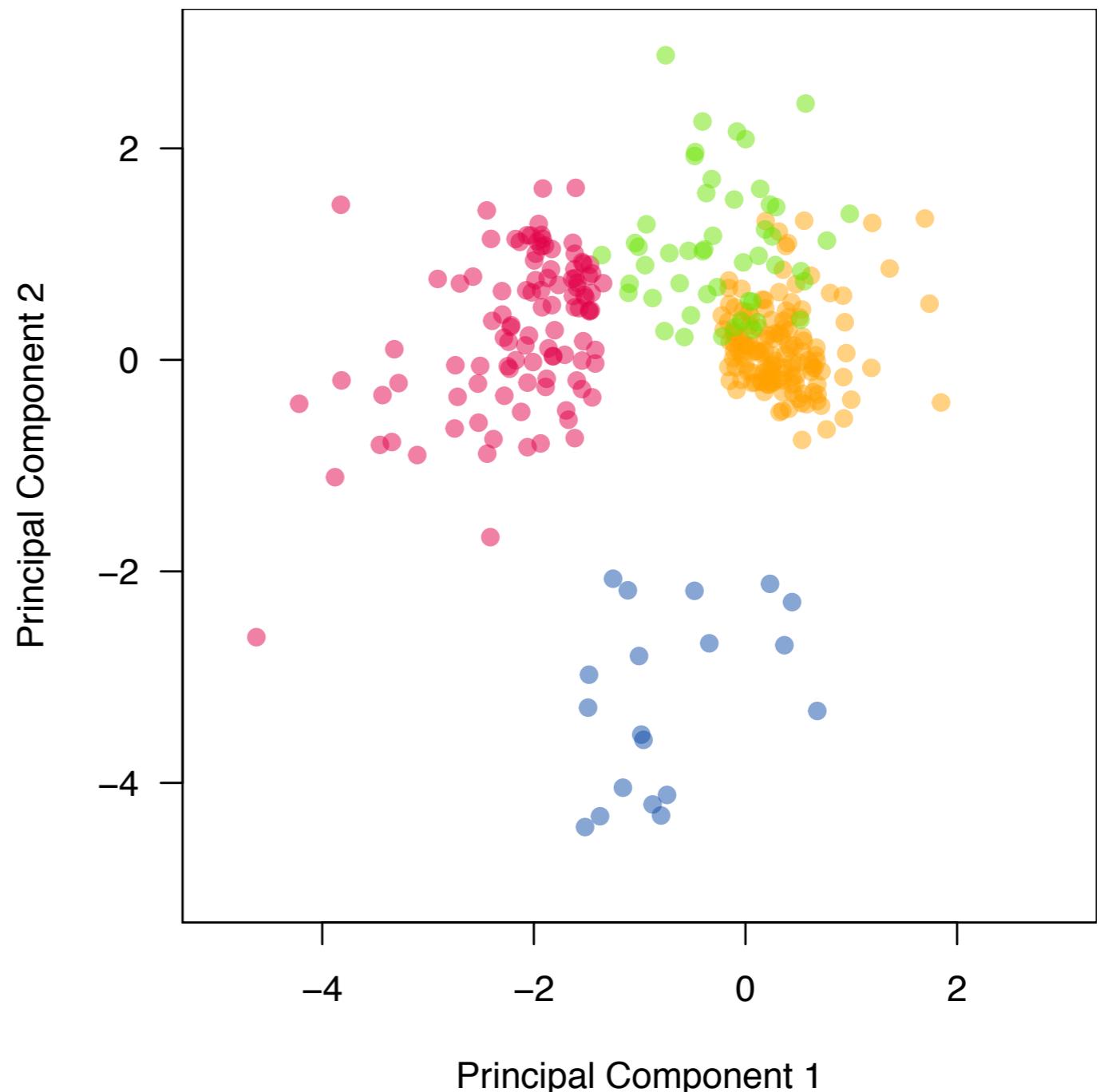
Multivariate Data: Derived Matrices

- matrices of multi-dimensional vectors
- usually abundance profiles, e.g. transcript or protein levels, metabolite concentrations



Multivariate Data: Derived Matrices

Scatter Plot

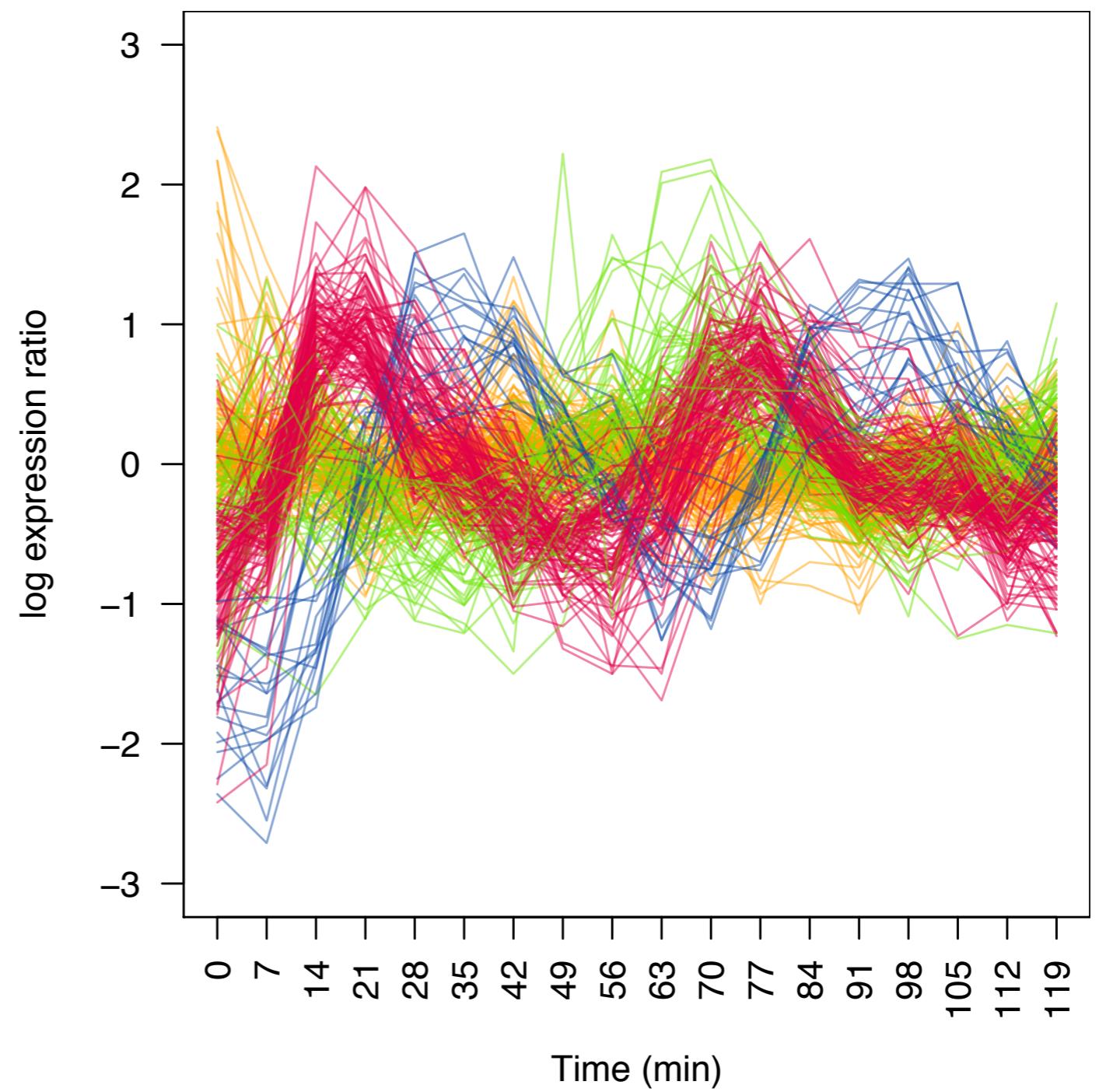


Multivariate Data: Derived Matrices

- **Scatter Plots and Dimensionality Reduction**
 - used to visualize high-dimensional profiles as projections in lower-dimensional spaces (usually 2D, sometimes also 3D ...)
 - there is always a loss of information in the process, goal is to minimize the loss of information
 - many different algorithms: Principal Components Analysis (PCA), Multi-Dimensional Scaling (MDS), Isomap, etc.
 - **Pros** - good choice to get an idea about the overall structure of the whole data set: clusters, outliers, gaps in the data
 - **Cons** - because of the dimensionality reduction the original profiles are not accessible in the visualization

Multivariate Data: Derived Matrices

Profile Plot a.k.a.
Parallel Coordinates



Multivariate Data: Derived Matrices

- **Profile Plot/Parallel Coordinate Plots**

- **Pros**

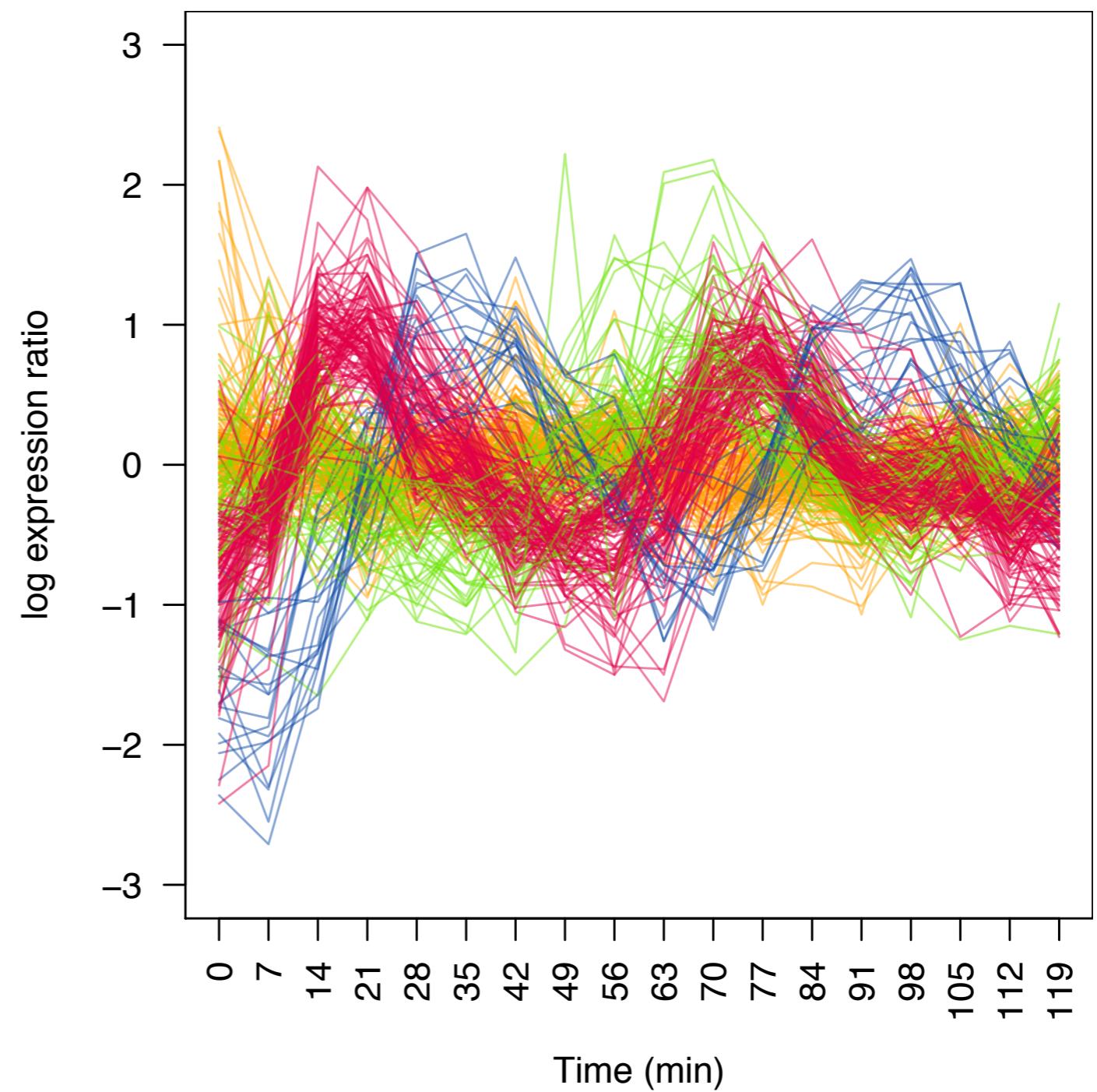
- encoding by position: profiles easy to read
 - color-coding of expression profiles (groups) very efficient

- **Cons**

- overplotting
 - grows horizontally with every additional sample

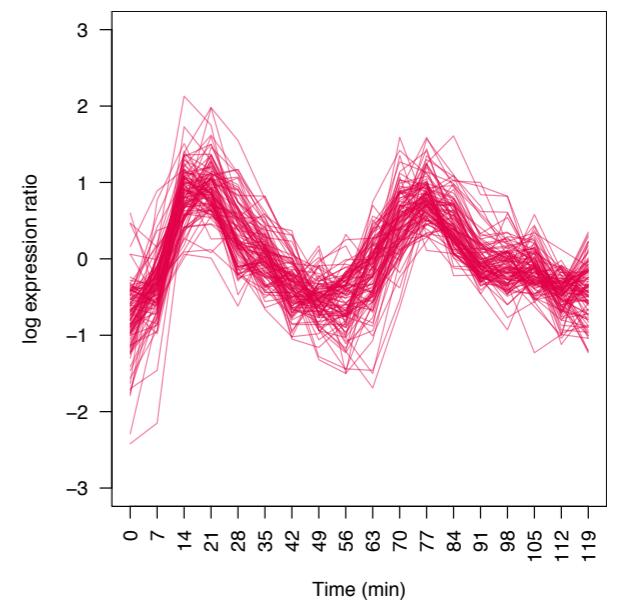
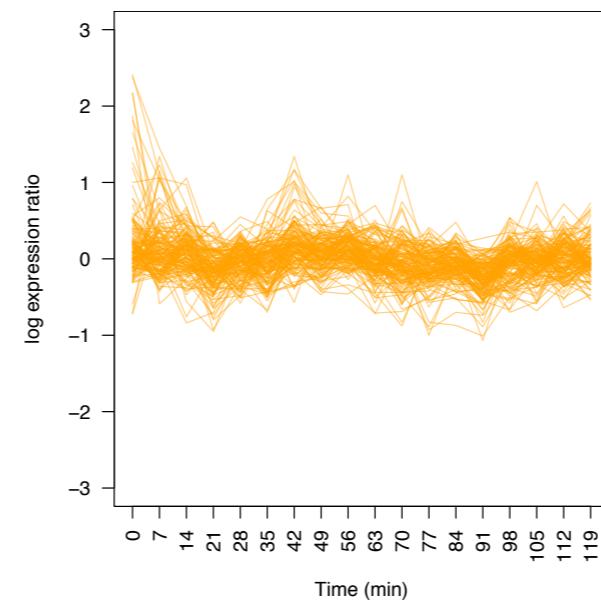
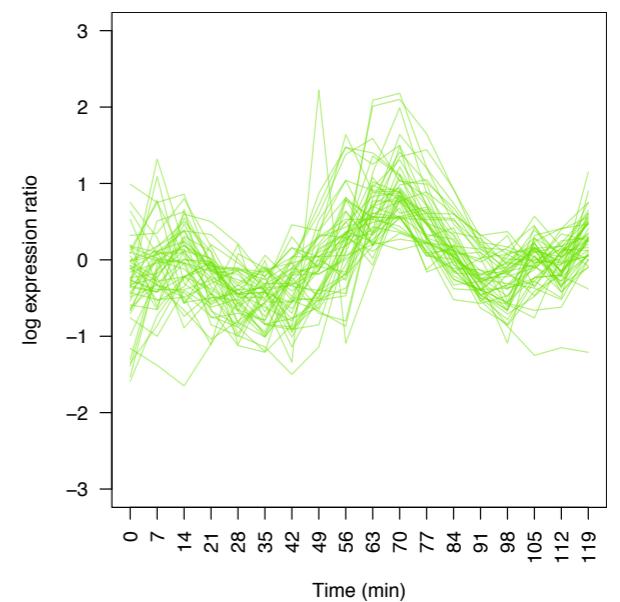
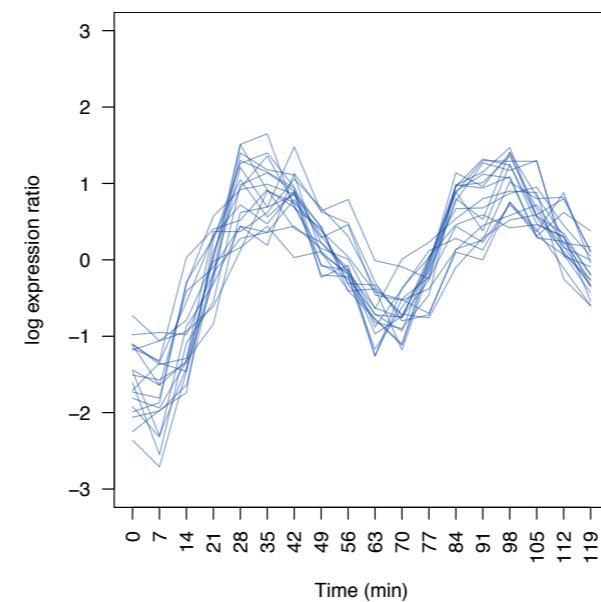
Multivariate Data: Derived Matrices

Profile Plot a.k.a.
Parallel Coordinates



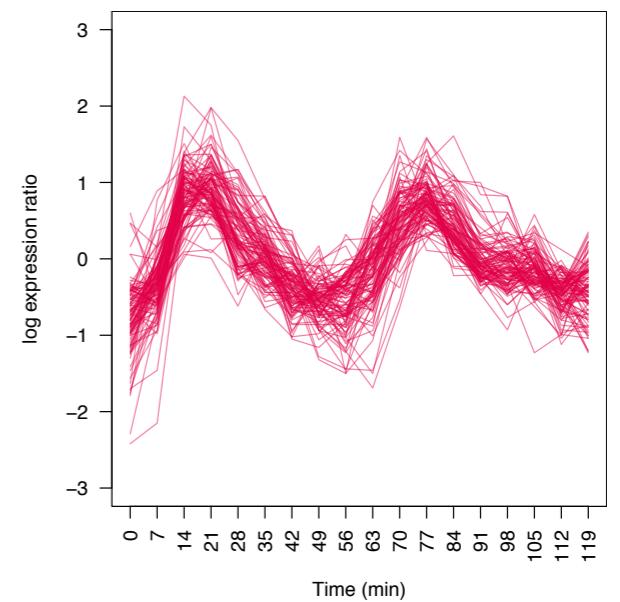
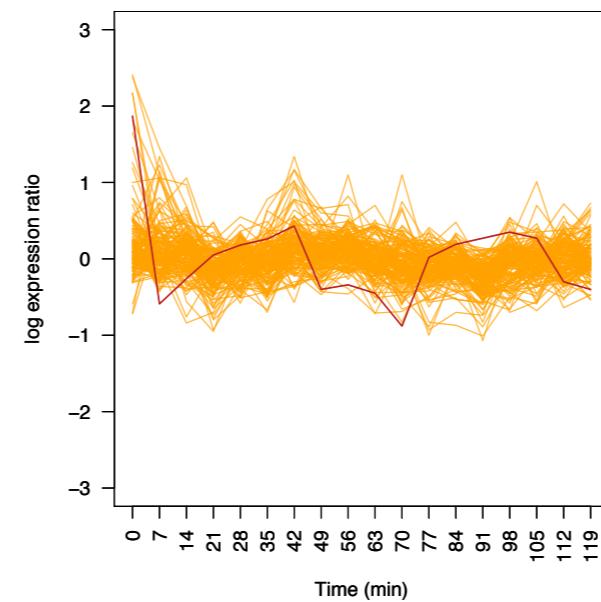
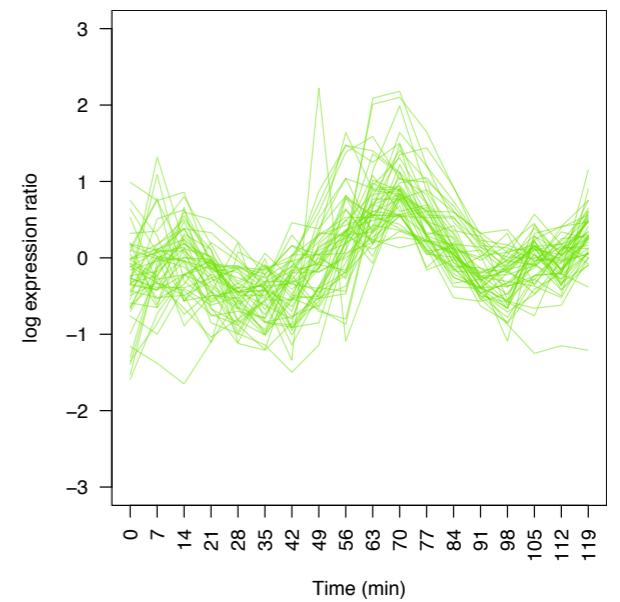
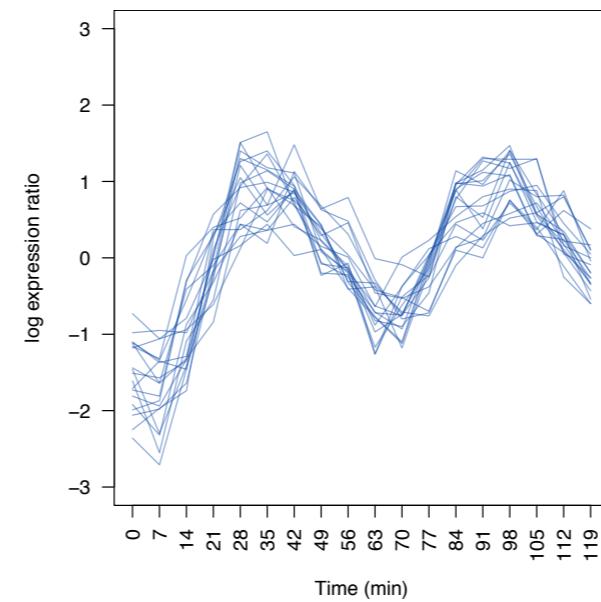
Multivariate Data: Derived Matrices

Profile Plot a.k.a.
Parallel Coordinates



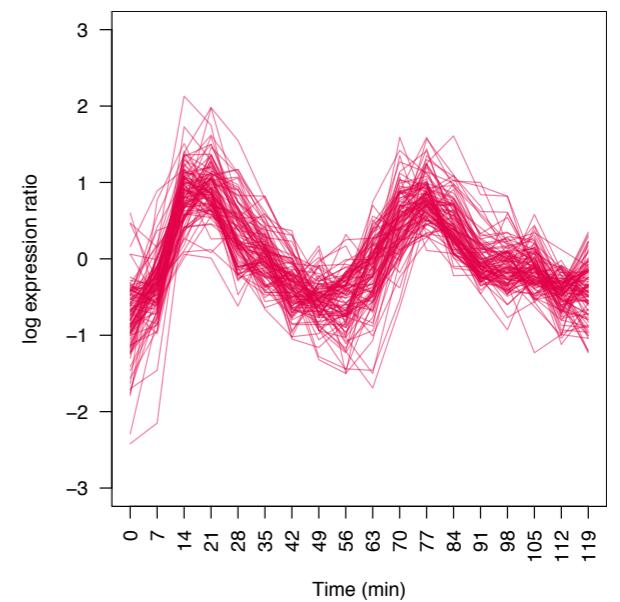
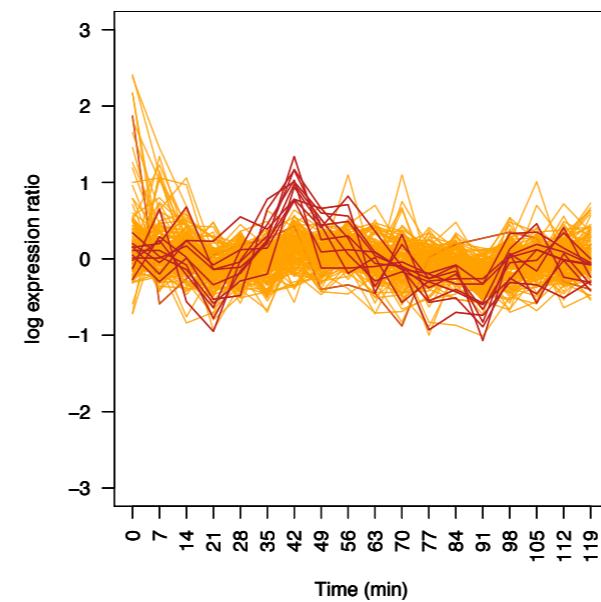
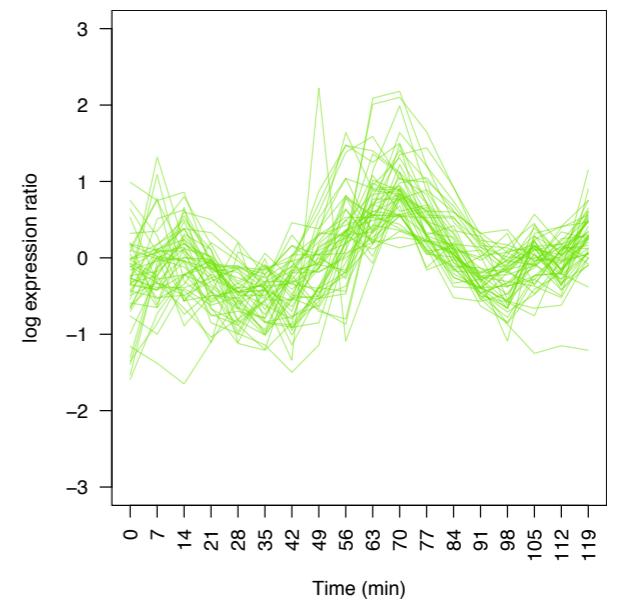
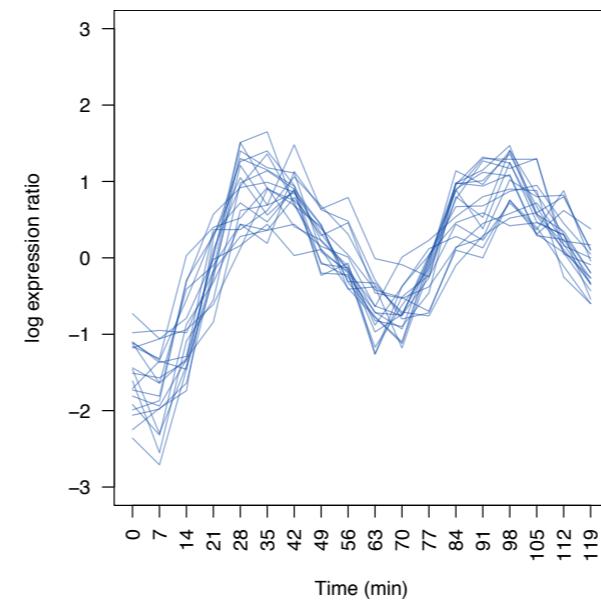
Multivariate Data: Derived Matrices

Profile Plot a.k.a.
Parallel Coordinates



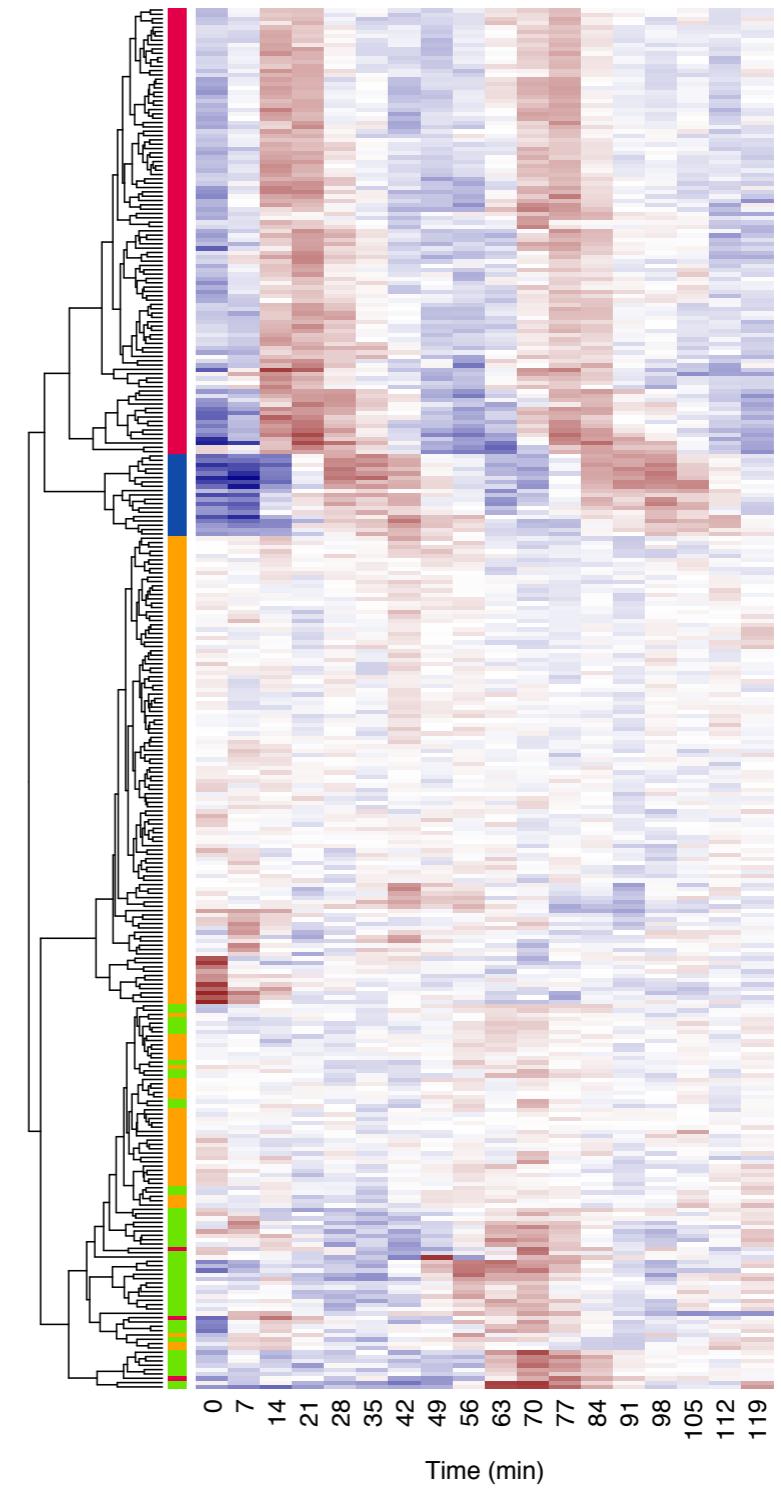
Multivariate Data: Derived Matrices

Profile Plot a.k.a.
Parallel Coordinates



Multivariate Data: Derived Matrices

Heat Map with Dendrogram



Multivariate Data: Derived Matrices

- **Heatmap**

- **Pros**

- no overplotting, yet a very dense information display
 - can be combined with dendrogram and additional information can be encoded in further columns or in the height of rows

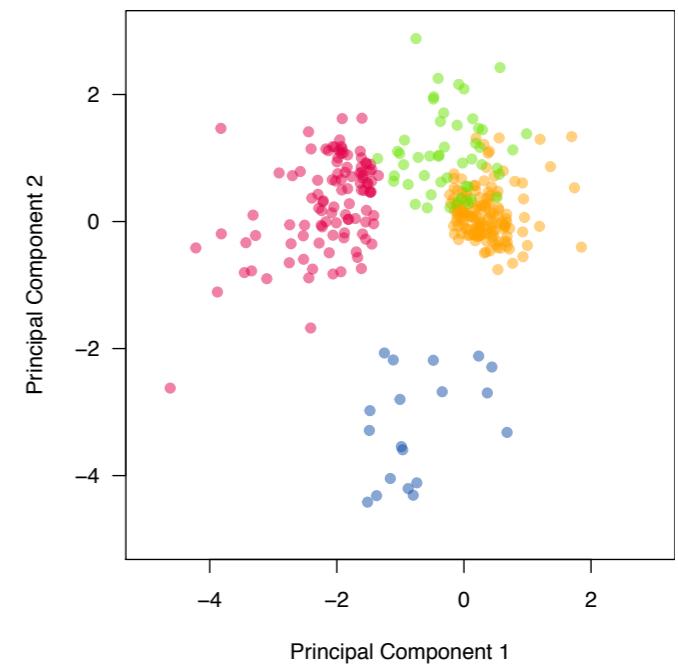
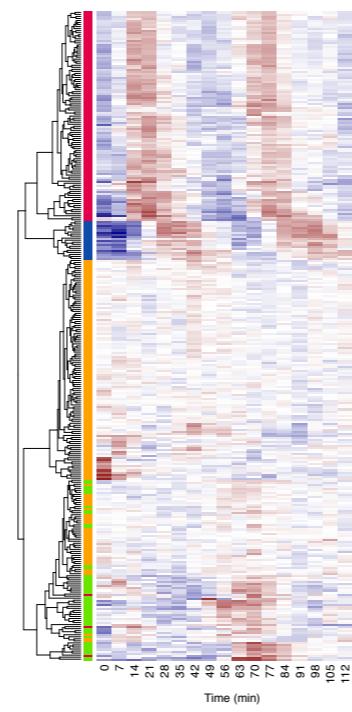
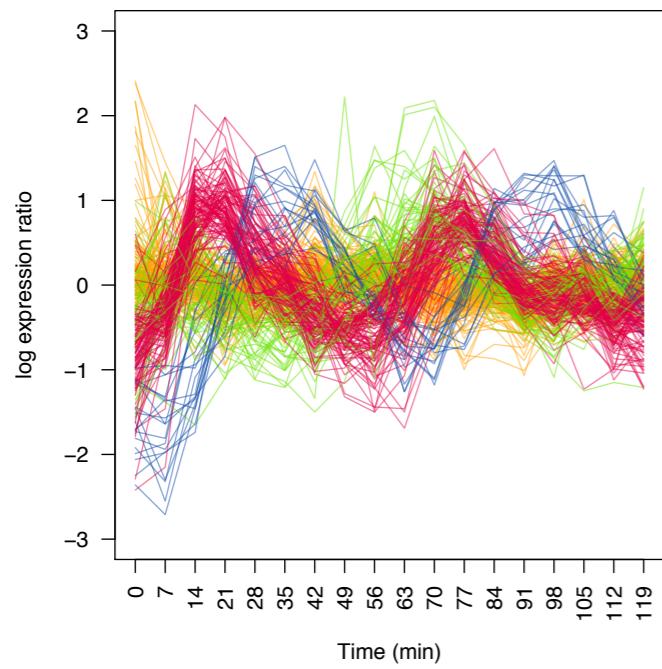
- **Cons**

- only qualitative interpretation possible due to color coding
 - grows horizontally with every additional sample and grows vertically with every additional profile

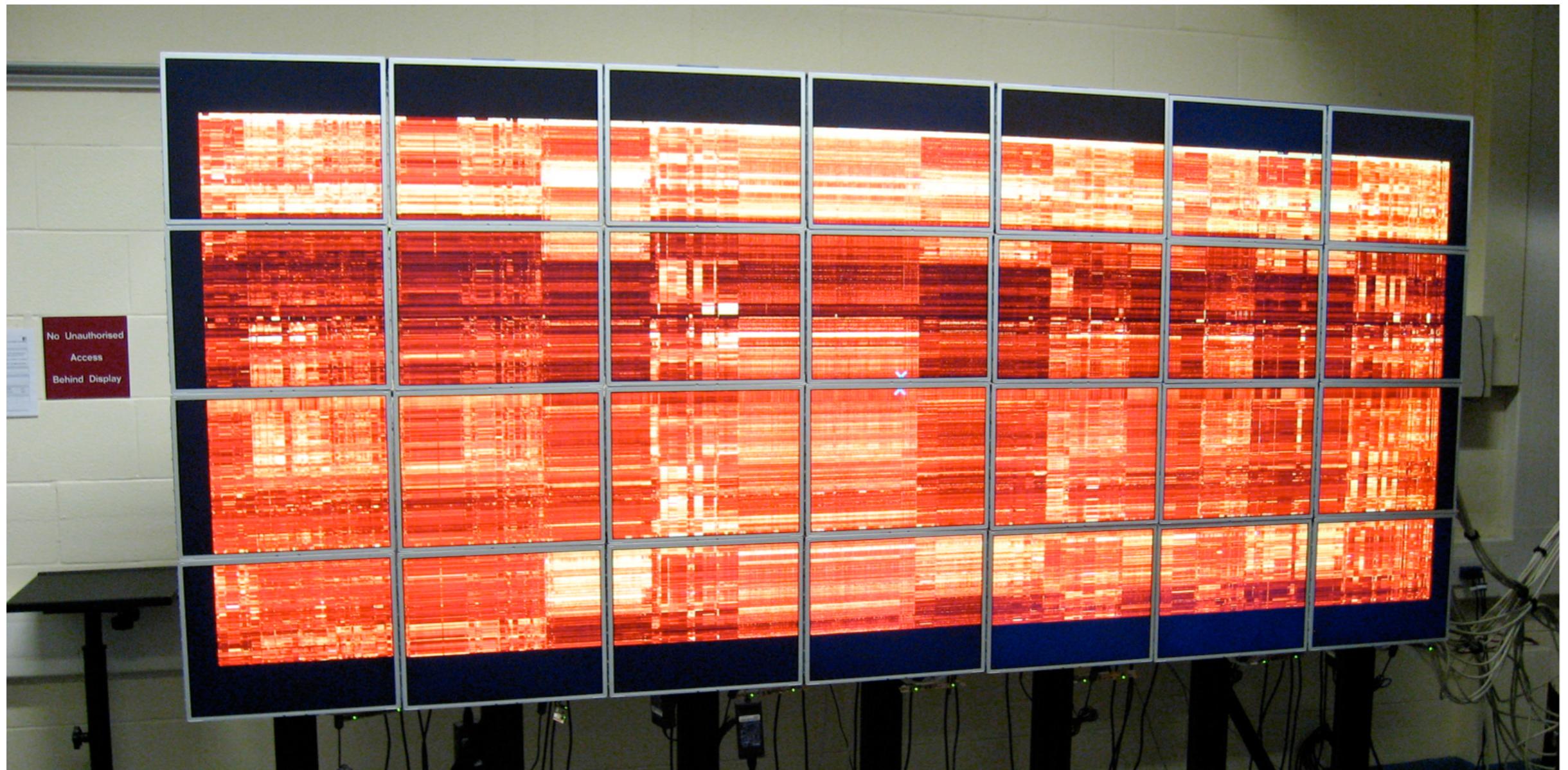
Multivariate Data: Summary

few, high-res

many, low-res



Problem: Very Large Expression Matrices

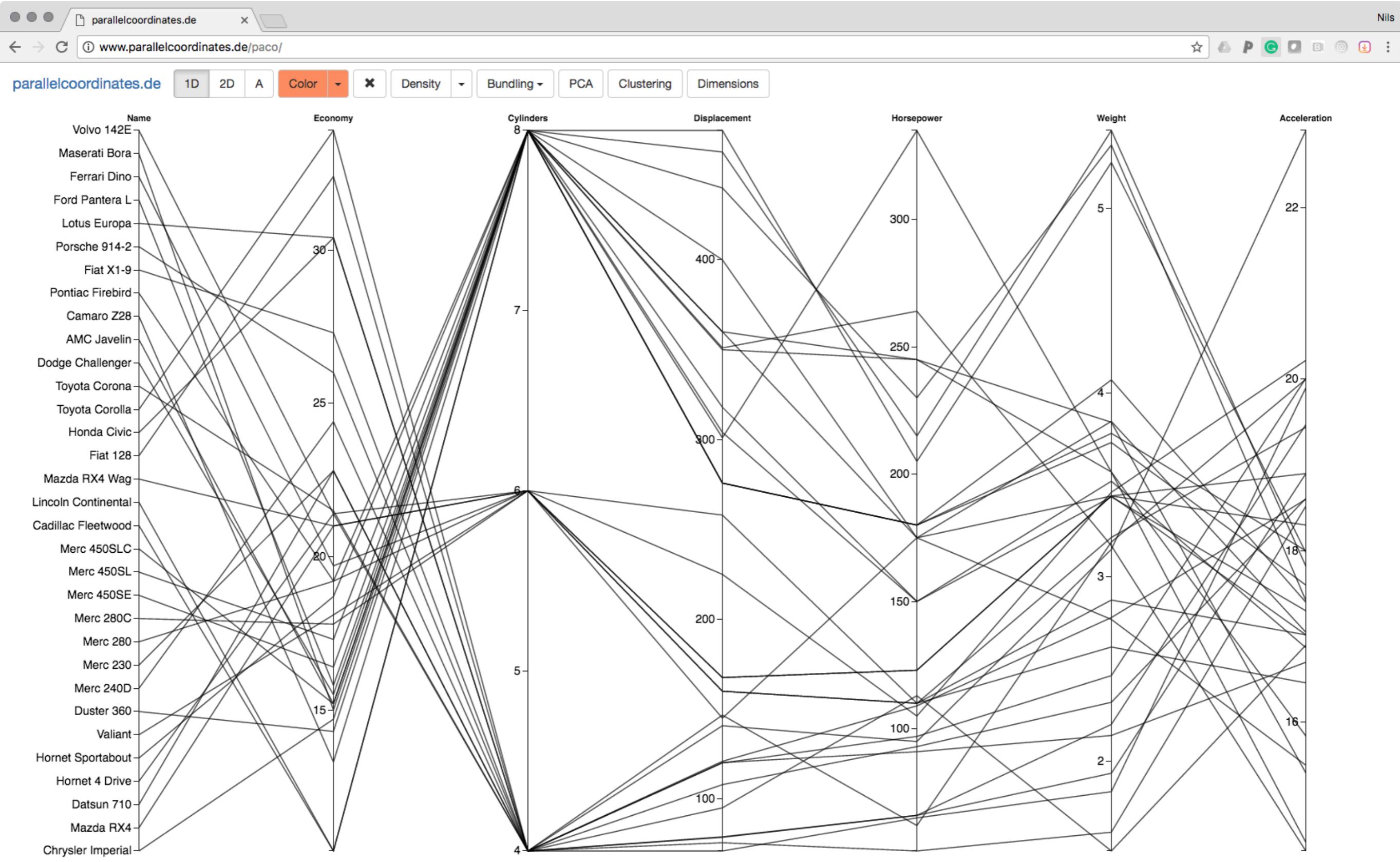


Power Wall (7x4 screens = 11,200x4,800), University of Leeds

1000 transcripts, 5372 samples

Multivariate Data

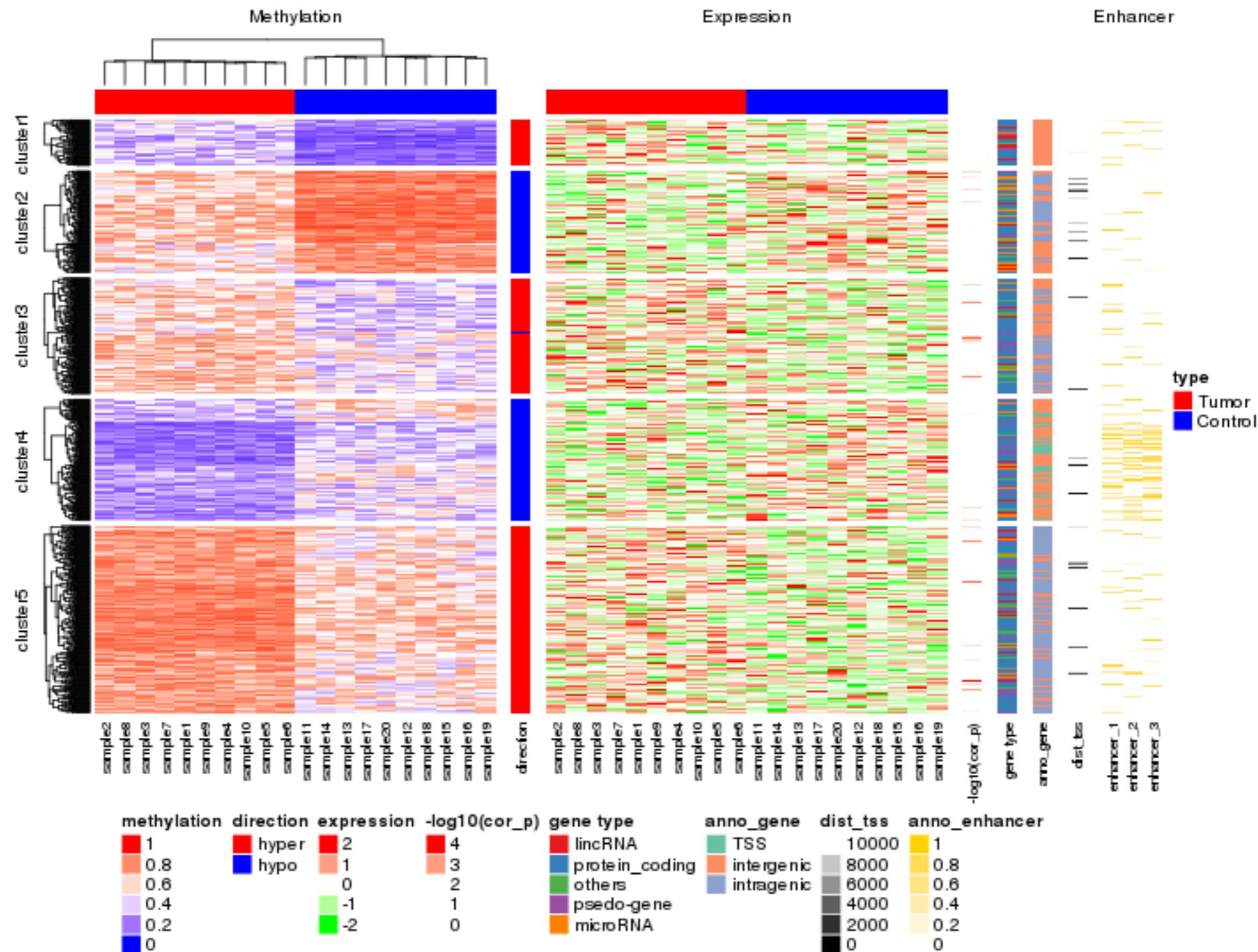
Heterogeneous Tables



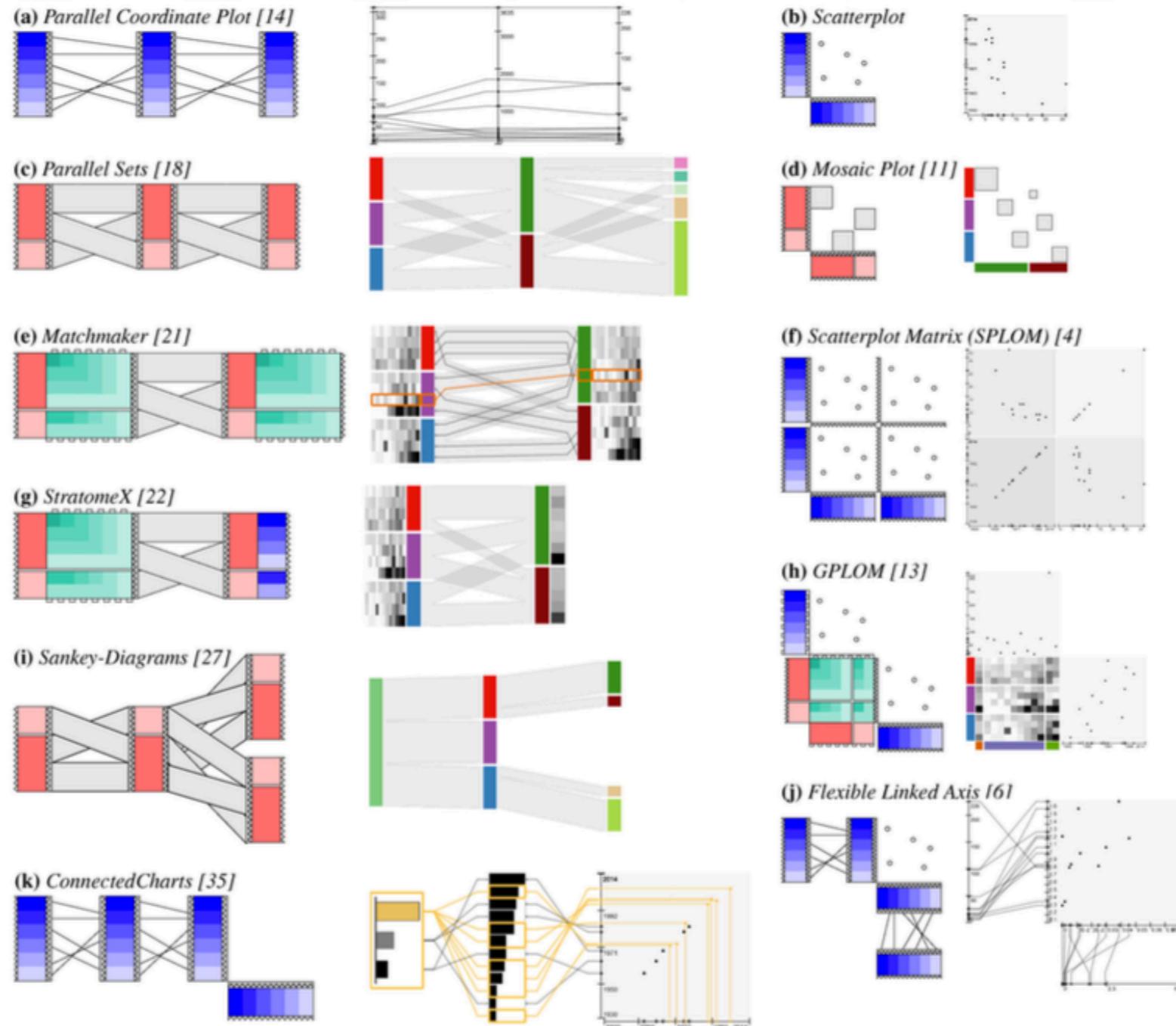
drop a csv-formatted file to load your own data. Note that the first line must describe the data scheme. See [this dataset](#) for an example.

<http://www.parallelcoordinates.de/paco/>

Correspondence between methylation, expression and other genomic features



Domino



Domino

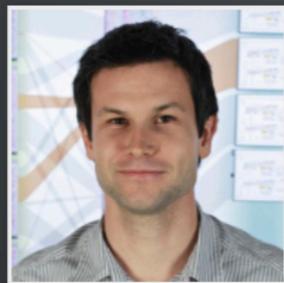
Extracting, Comparing, and Manipulating Subsets
across Multiple Tabular Datasets



JKU
JOHANNES KEPLER
UNIVERSITY LINZ



Samuel Gratzl



Marc Streit



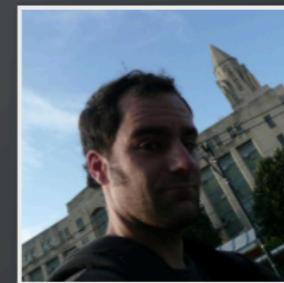
HARVARD
MEDICAL SCHOOL



Nils Gehlenborg



HARVARD
School of Engineering
and Applied Sciences



Alexander Lex

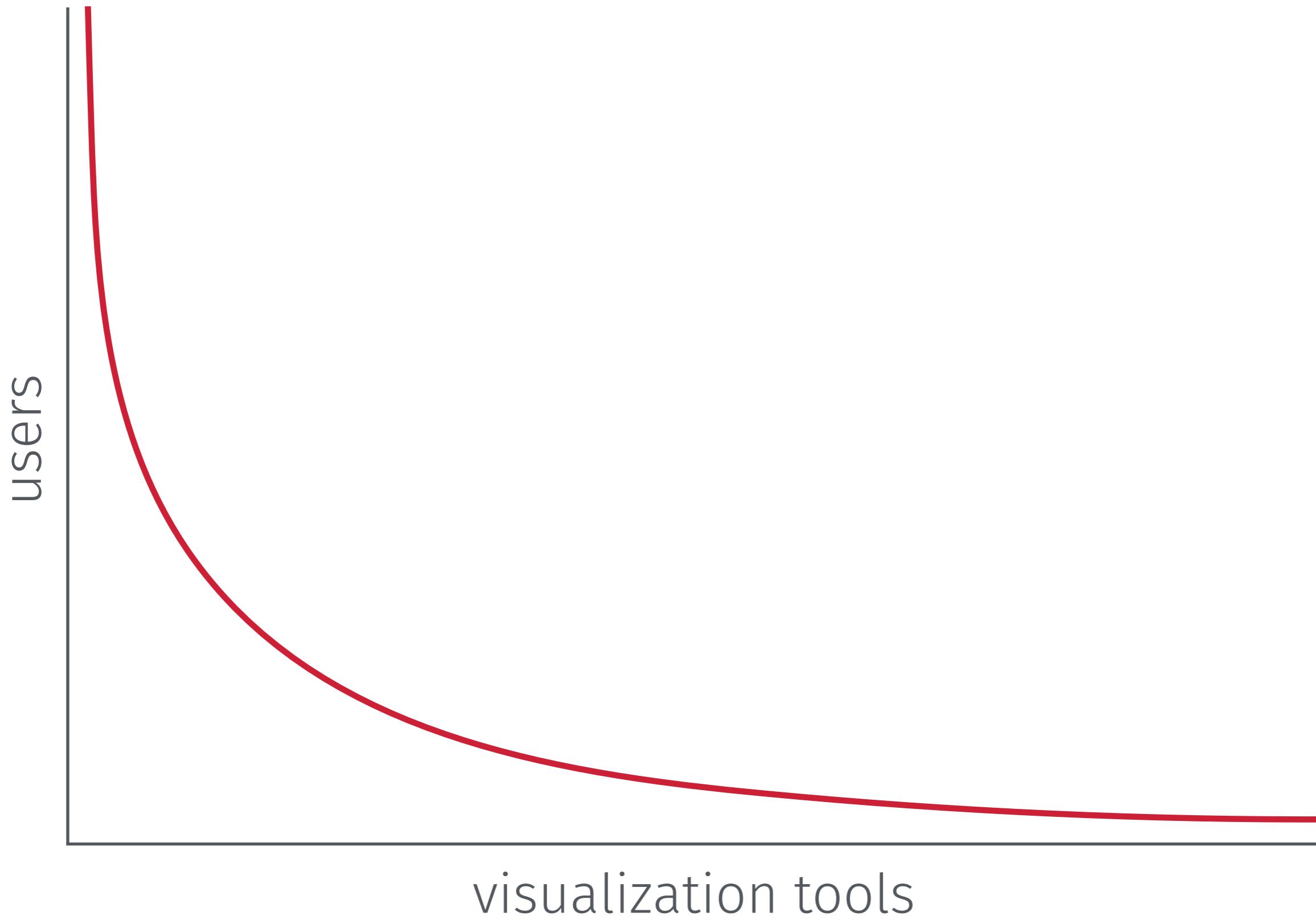


Hanspeter Pfister

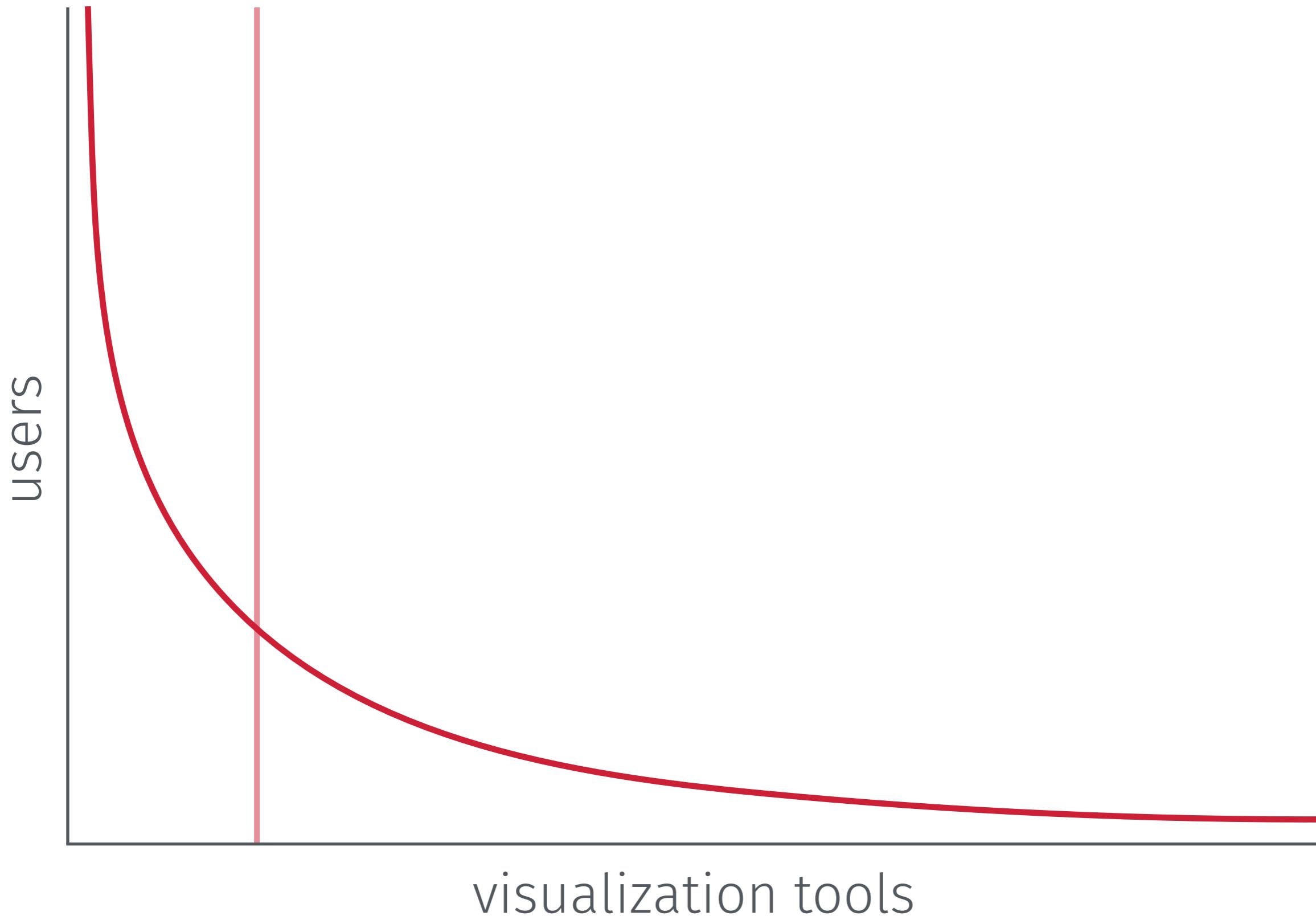
https://sgratzl.github.io/domino_vis2014/#/title

The Long Tail

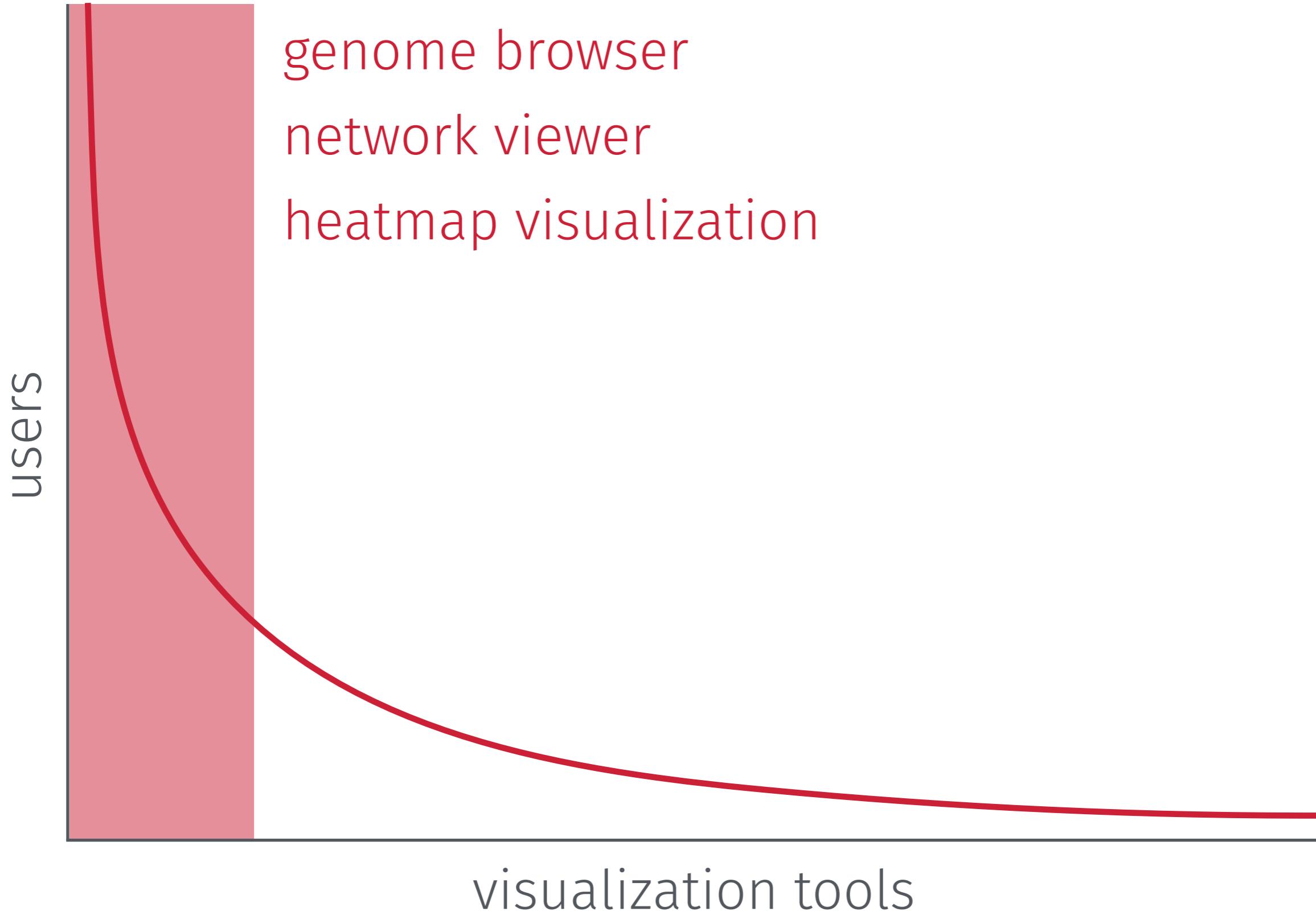
Data Visualization: The Long Tail



Data Visualization: The Long Tail



Data Visualization: The Long Tail



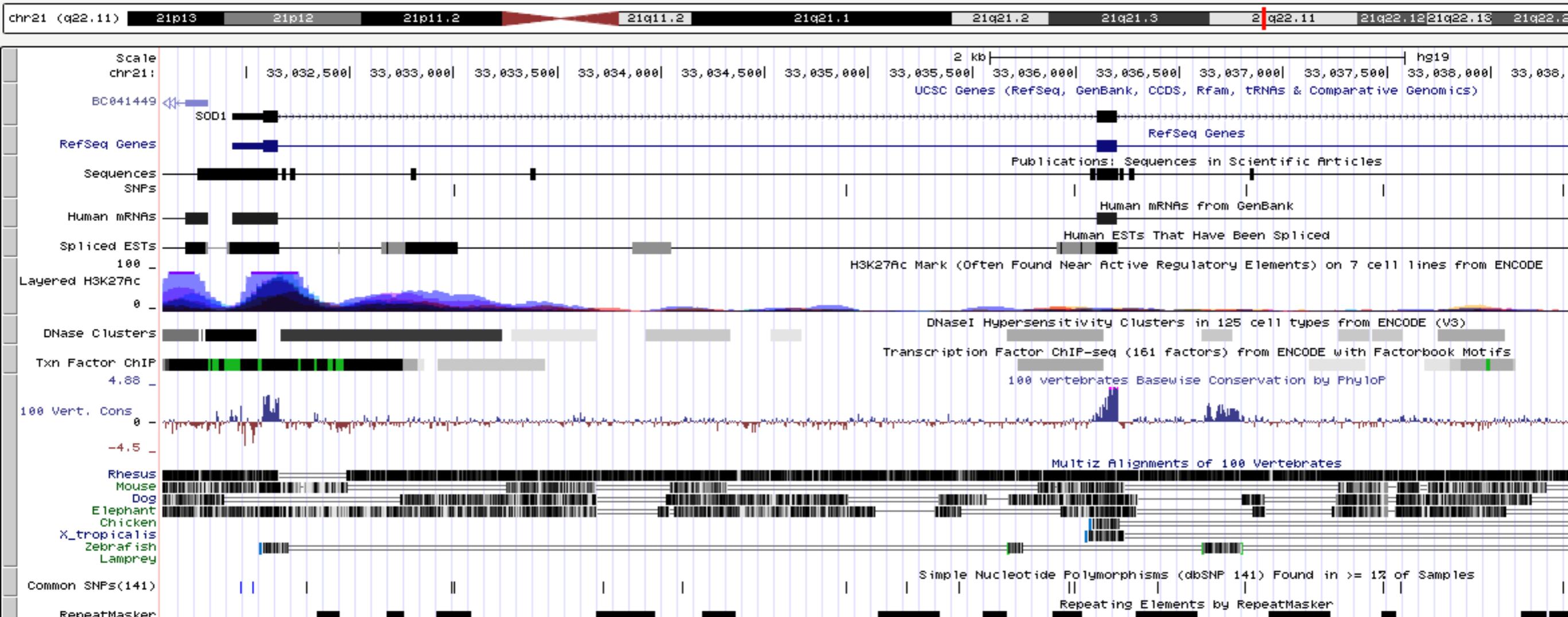

[Genomes](#)
[Genome Browser](#)
[Tools](#)
[Mirrors](#)
[Downloads](#)
[My Data](#)
[View](#)
[Help](#)
[About Us](#)

UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

move <<< << < > >> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x

chr21:33,031,597-33,041,570 9,974 bp. enter position, gene symbol or search terms

go



move start

Click on a feature for details. Click or drag in the base position track to zoom in. Click side bars for track options. Drag side bars or labels up or down to change position.

< 2.0 >

[track search](#) [default tracks](#) [default order](#) [hide all](#) [add custom tracks](#) [track hubs](#) [configure](#) [reverse](#) [resize](#) [refresh](#)

[collapse all](#)

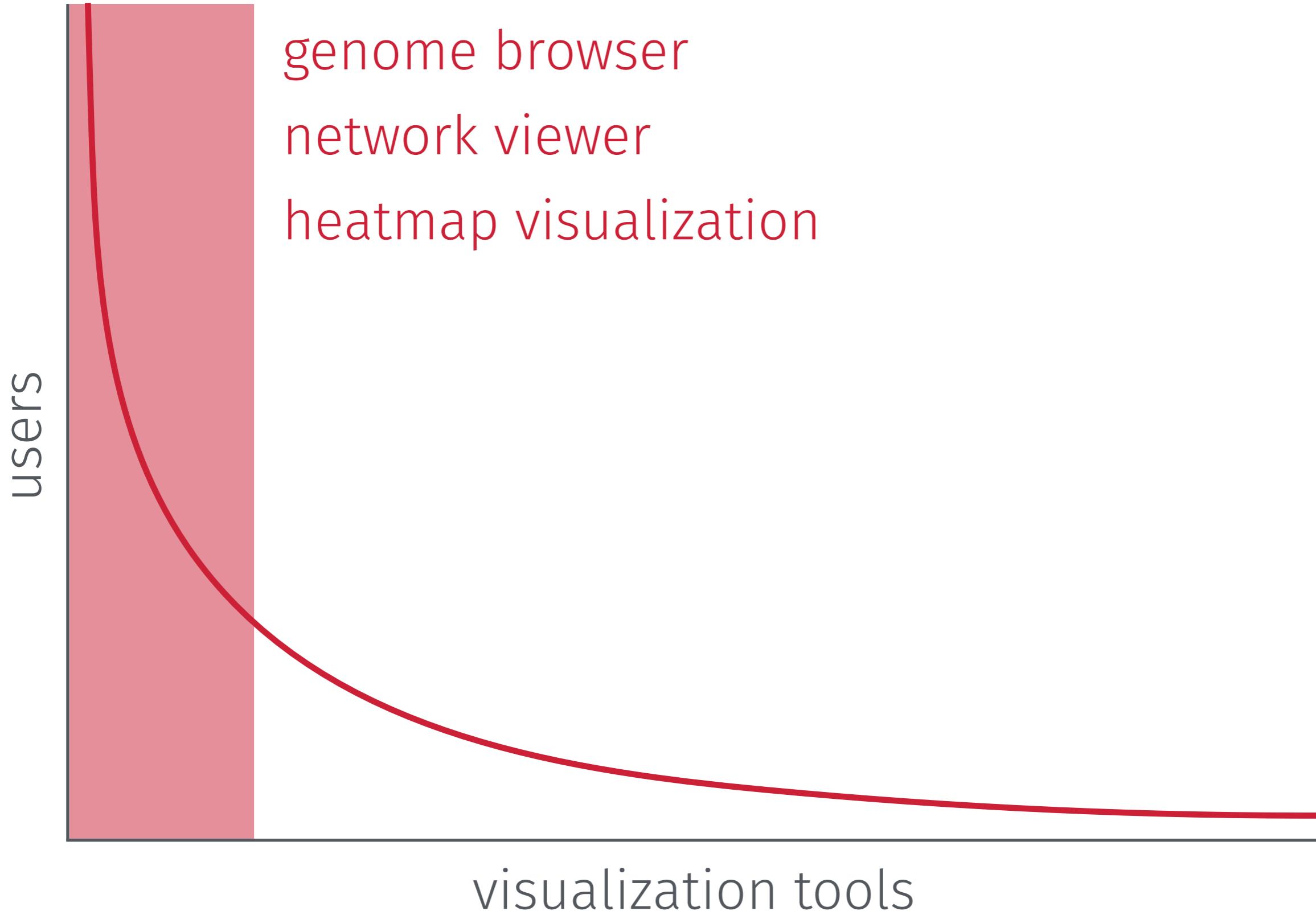
Use drop-down controls below and press refresh to alter tracks displayed.
Tracks with lots of items will automatically be displayed in more compact modes.

[expand all](#)

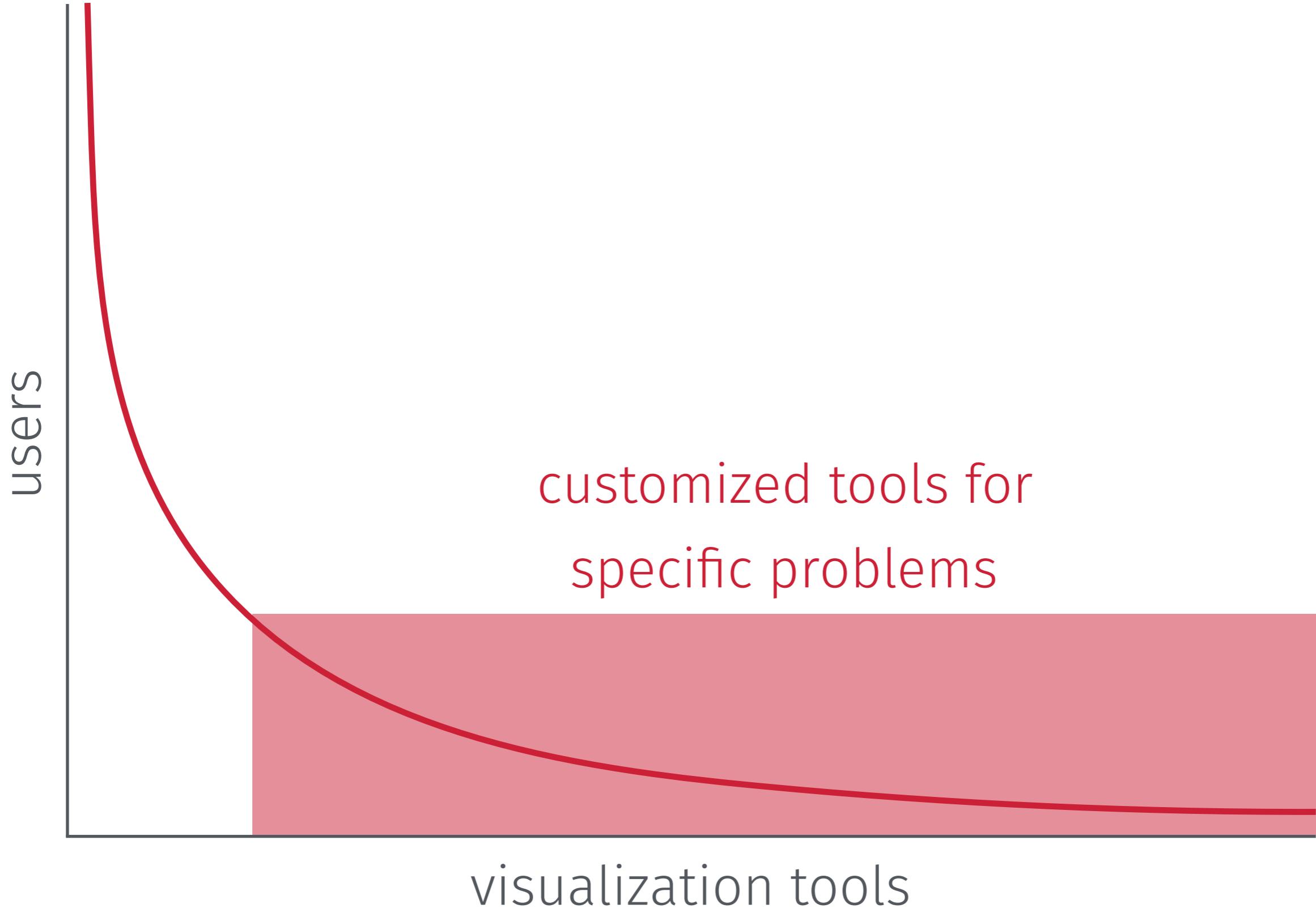
Mapping and Sequencing

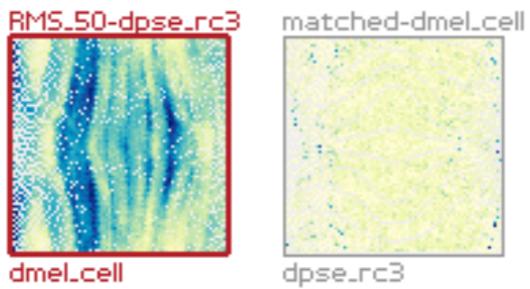
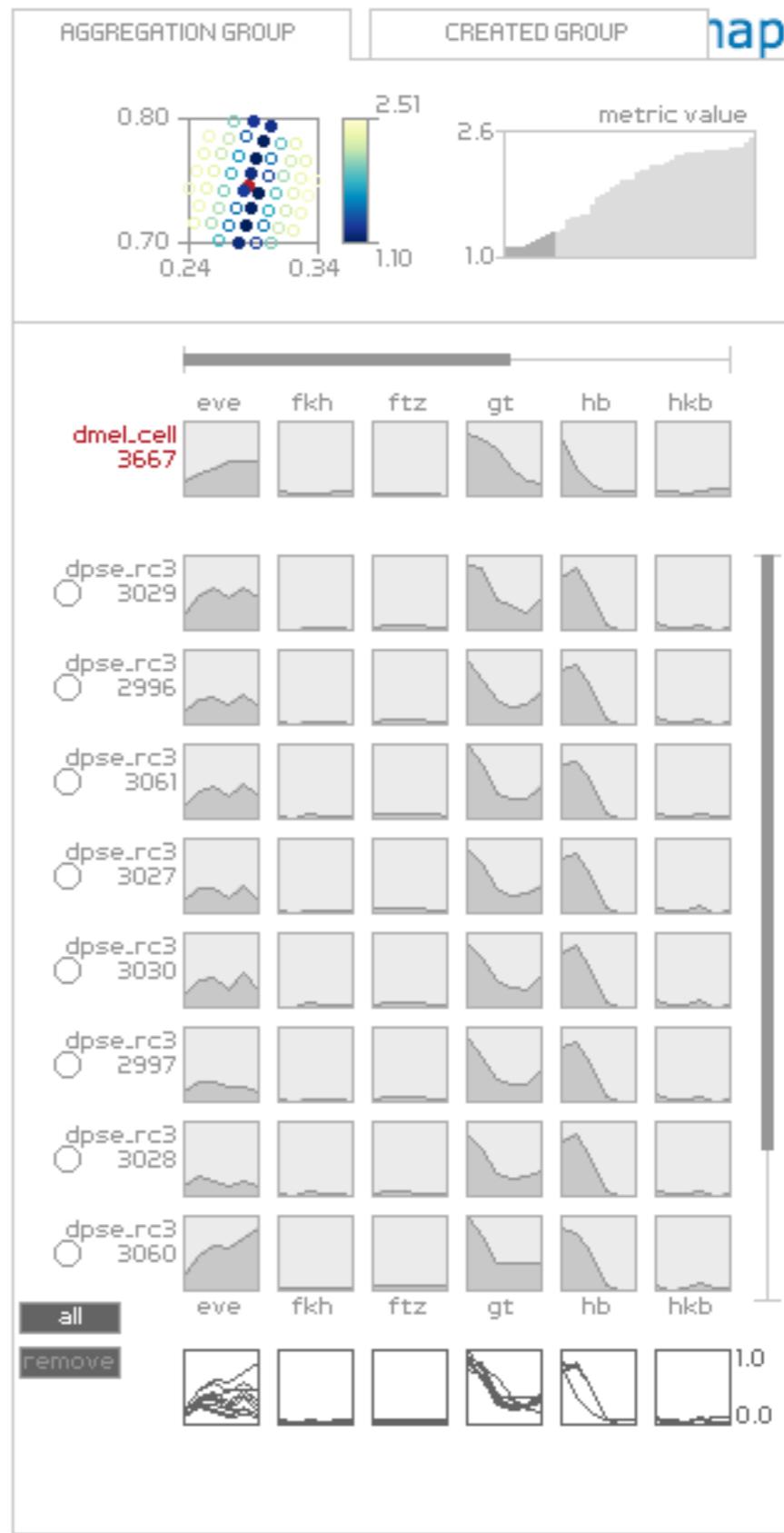
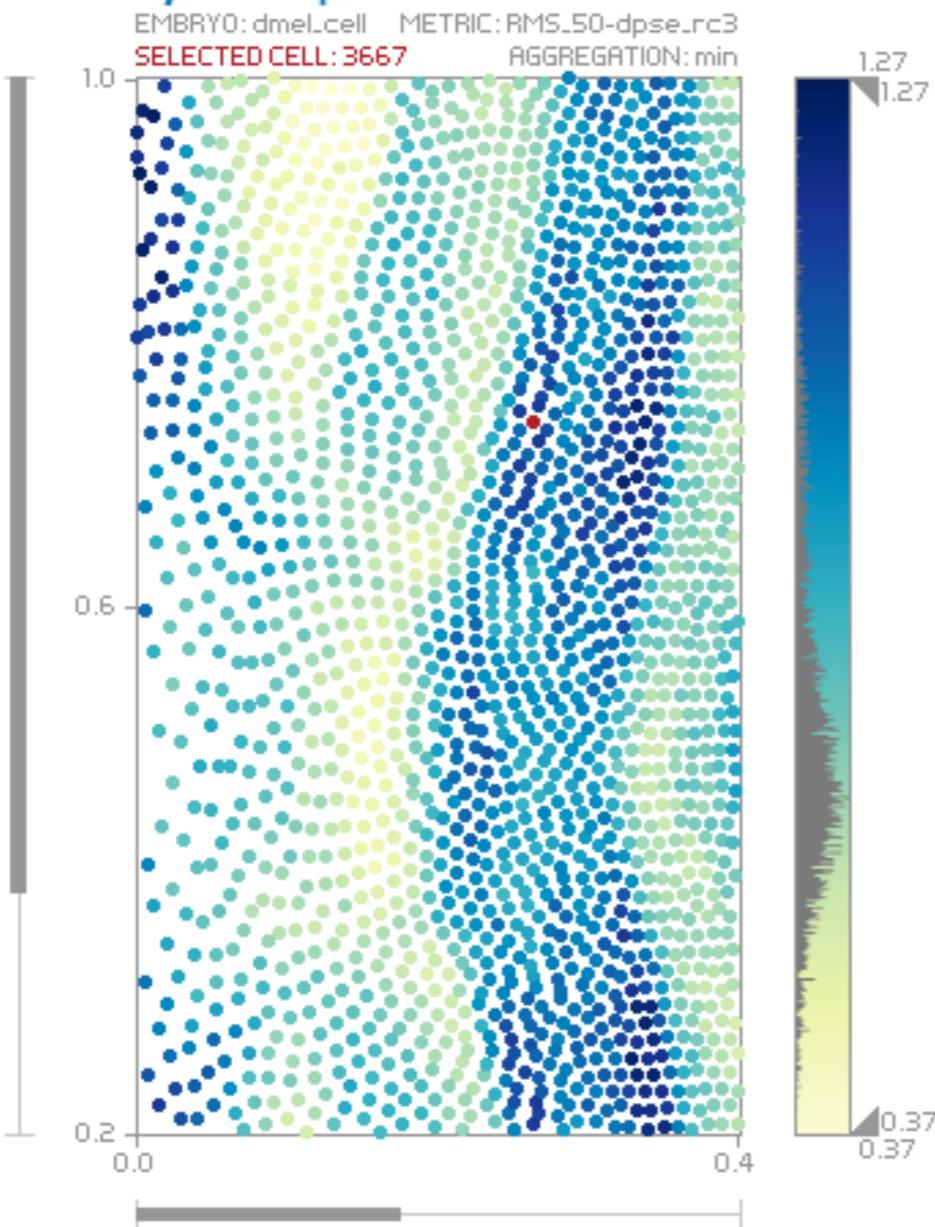
Display a menu for "http://genome-euro.ucsc.edu/cgi-bin/hgc?hgsid=200830848_Pr7vSRBAuNDdCe9FN6cl81oDcNRD&c=chr21&o=33031596&t=33041570&g=phyloP100wayAll&i=phyloP100wayAll"

Data Visualization: The Long Tail

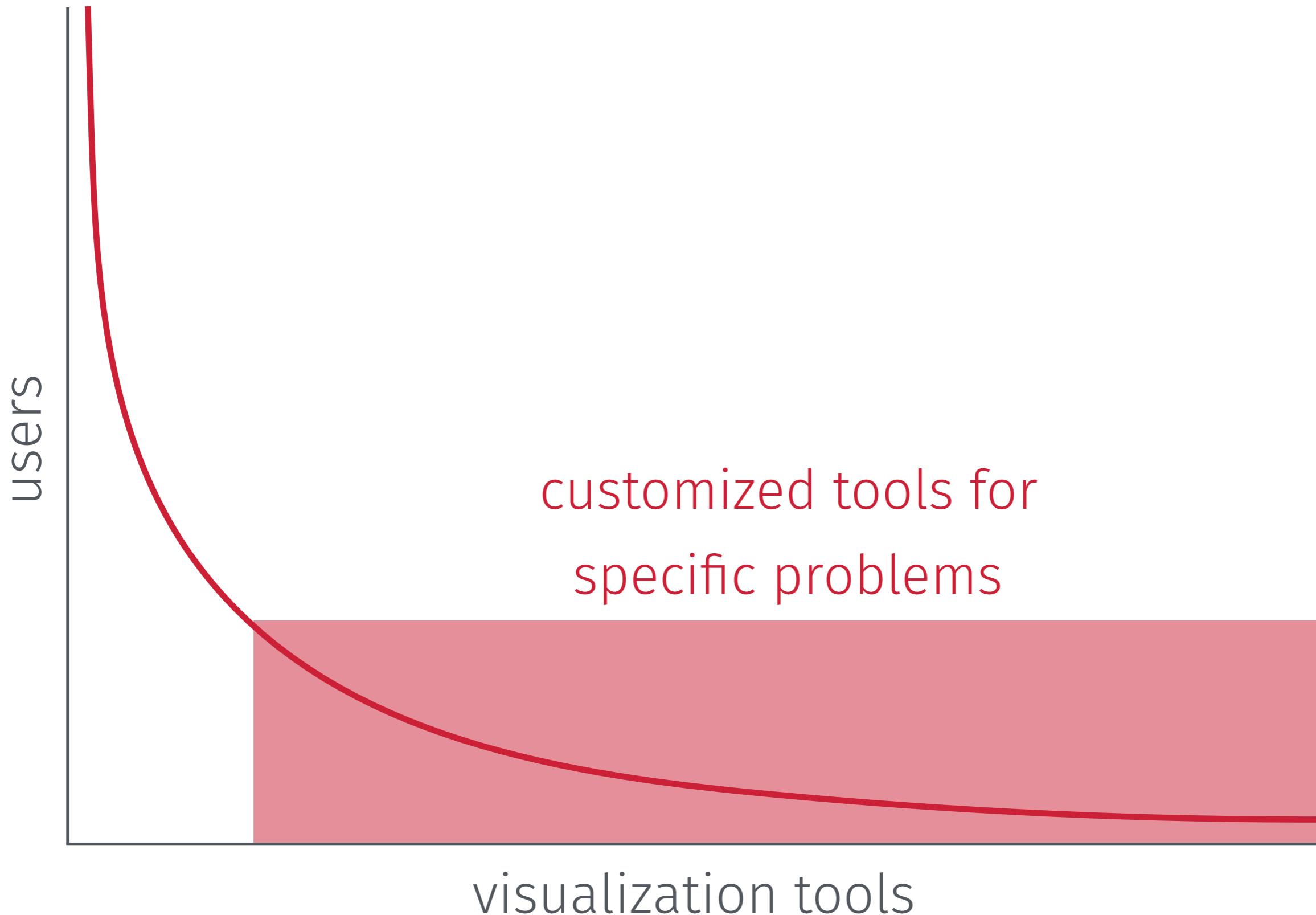


Data Visualization: The Long Tail



Summaries**Embryo Map**

Data Visualization: The Long Tail



Data Visualization: The Long Tail

