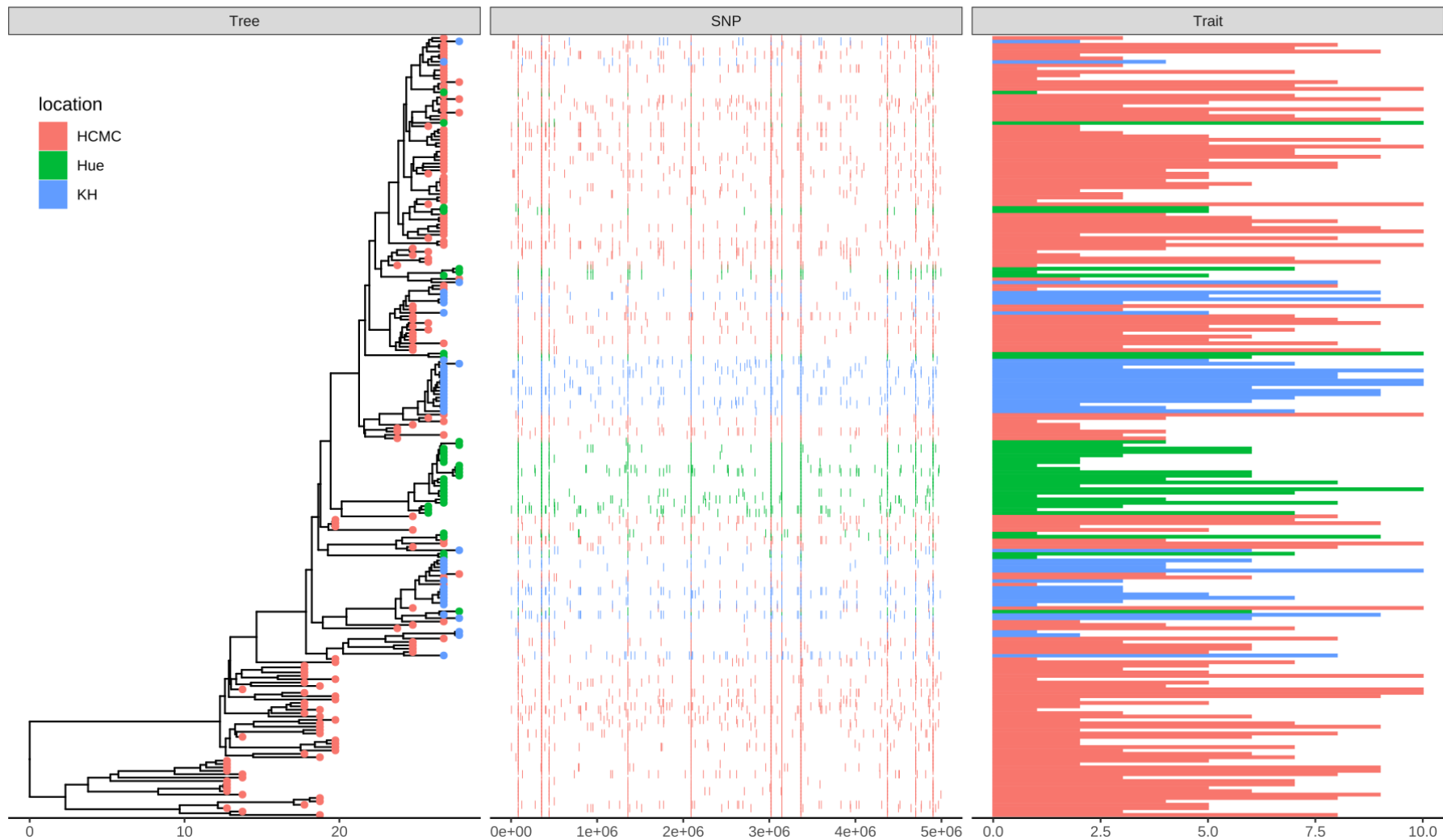# CM 515: Data Visualization Module I

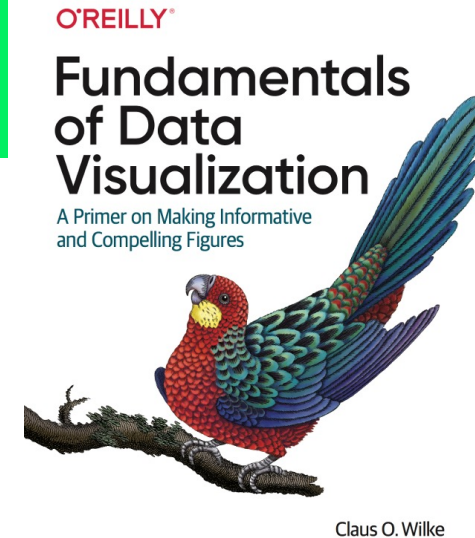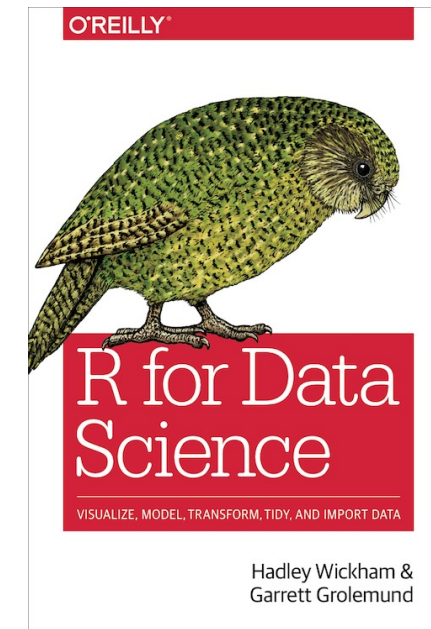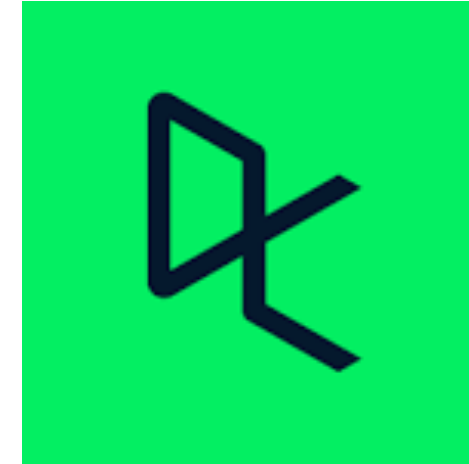**Dan Sloan**
Department of Biology
dan.sloan@colostate.edu

# Data Visualization Resources

**Datacamp**
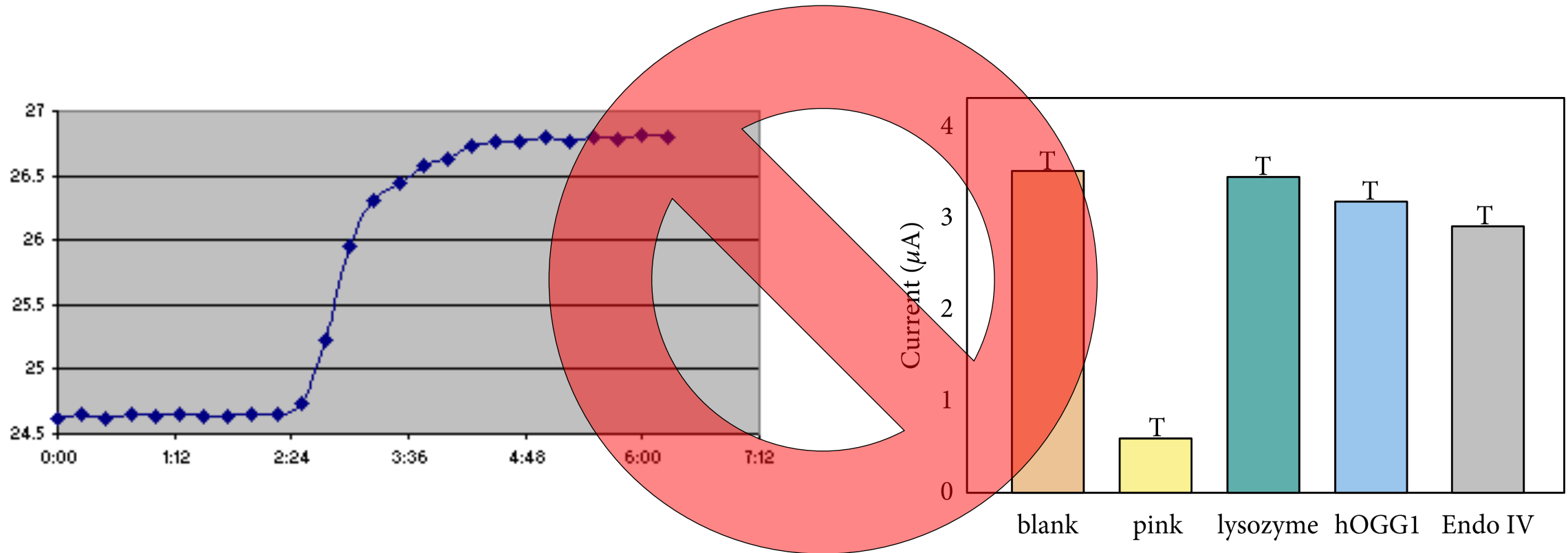
- Introduction to Data Visualization with ggplot2

- Intermediate Data Visualization with ggplot2

- Visualization Best Practices in R

**Claus Wilke Data Visualization in R Course (U Texas)**

**R for Data Science**

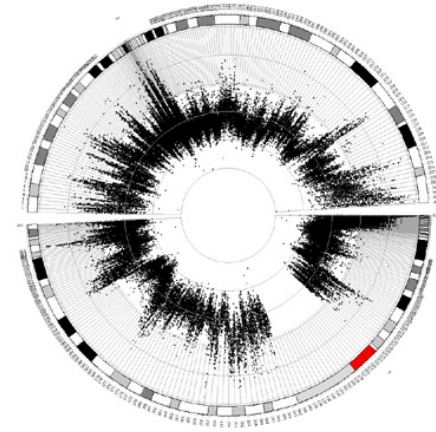# Scientific Communication and Professionalism



- Clear, accurate, and complete representation of your data

- Efficient, reproducible, and automated methods

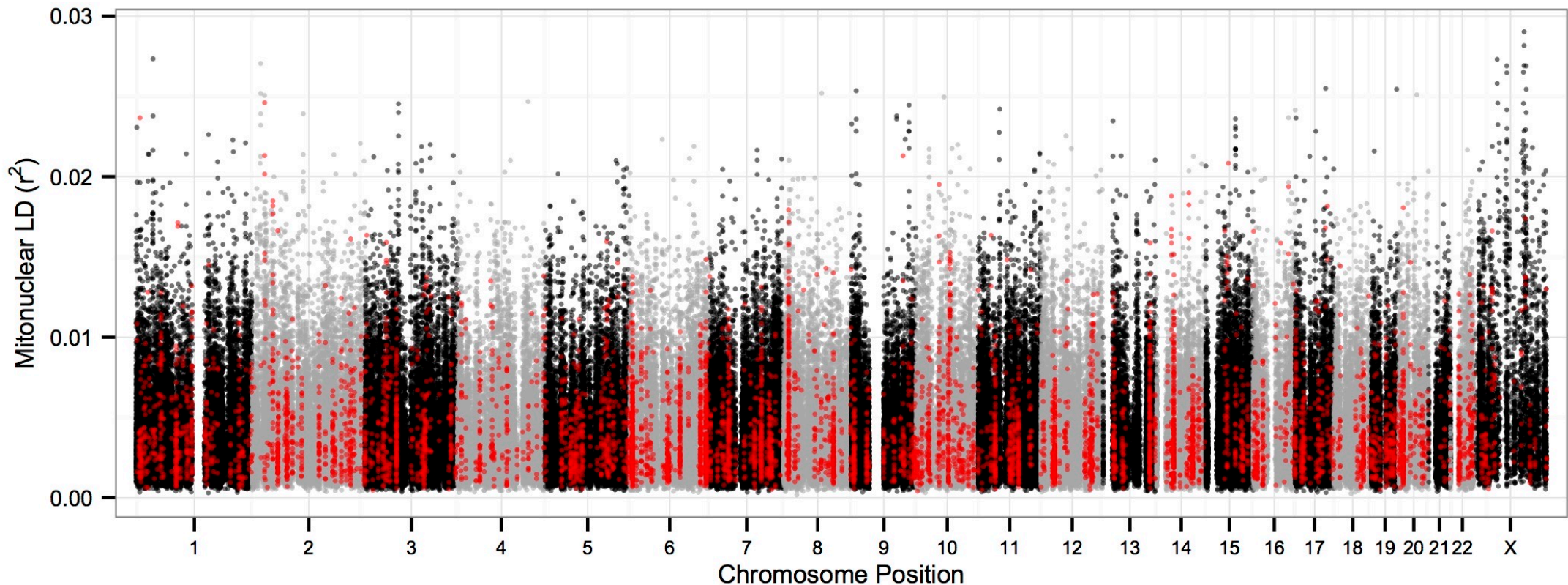- Clean, professional, and aesthetically pleasing appearance

# The Right Tools for the Job

**A few examples**

- [R and ggplot](#)

- [Circos](#)

- [Processing](#)

- [Adobe Illustrator](#)

- [BioRender](#)

# Making Figures with Code



```
ggplot(cnld) + geom_point(aes(x=CumPos, y=r2, size=0.75, colour=as.factor(ChromPrint), alpha =
1/8)) + scale_size_identity() + theme_bw(base_size=15) +
scale_color_manual(values=c(rep(c('black', 'dark gray'),11), 'black', 'red')) +
scale_x_continuous(expand = c(0.015, 0.015),labels=c(as.character(1:chrNum), "X"),
breaks=bpMidVec) + theme(plot.margin = unit( c(0.03,0.03,0.03,0.03) , "in" ),
legend.position='none', axis.text.x = element_text(size=6), axis.text.y = element_text(size=7),
axis.title.x = element_text(size=8), axis.title.y = element_text(size=8)) + xlab('Chromosome
Position') + ylab(expression(paste("Mitonuclear LD (",r^2, ")")))
```

# The Grammar of Graphics

- **aes**: Aesthetic mapping of data to plot elements
  - position (X or Y coordinates), shapes, sizes, color, line weight/type, transparency, etc.

- **geoms**: Layers visually representing your mapped data
  - points, lines, bars, density curves, etc.

- **facets**: Dividing into subplots
  - partitions dataset based on one or more variables

- **themes**: Non-data plot elements
  - axis labels, grid lines, titles, etc.



Grammar of Graphics by L. Wilkinson

Describes non-data ink. Design elements!
The plotting space you are using
Statistical models & summaries
Rows and columns of sub-plots
Shapes used to represent your data
The scales on which the data is mapped
The actual variables to be plotted

Themes 7
Coördinates 6
Statistics 5
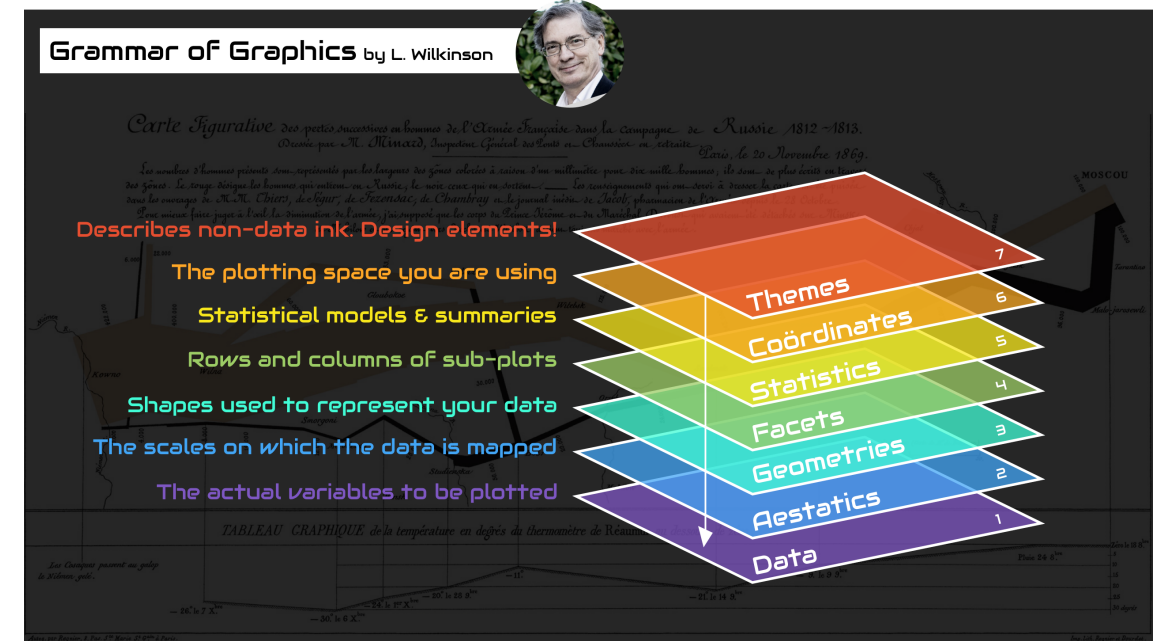Facets 4
Geometries 3
Aestatics 2
Data 1

Image: Thomas de Beus

# A Few Principles of Data Visualization

- What are you trying to communicate??

- The demise of the bar plot

- Choosing a scale: log vs. linear

- Use and choice of color palettes

https://www.sciencedirect.com/science/article/pii/S2666389920301896
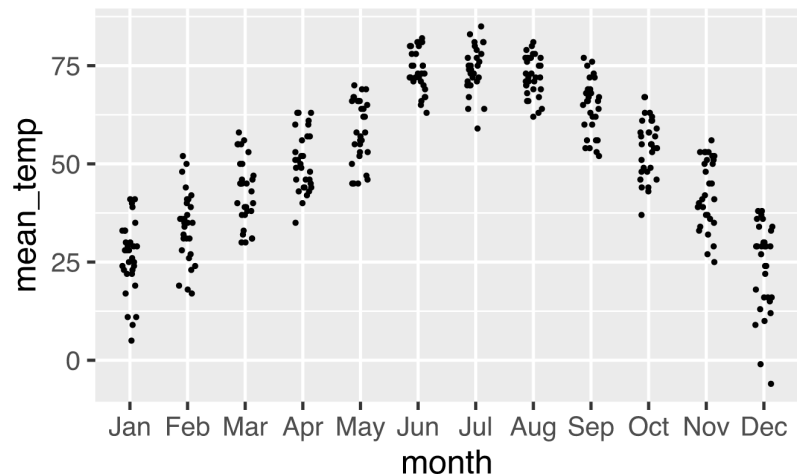
# The Demise of the Bar Plot

# When Possible… Show All the Data!

Use point size, jitter, and/or transparency to mitigate the effects of overlapping points.

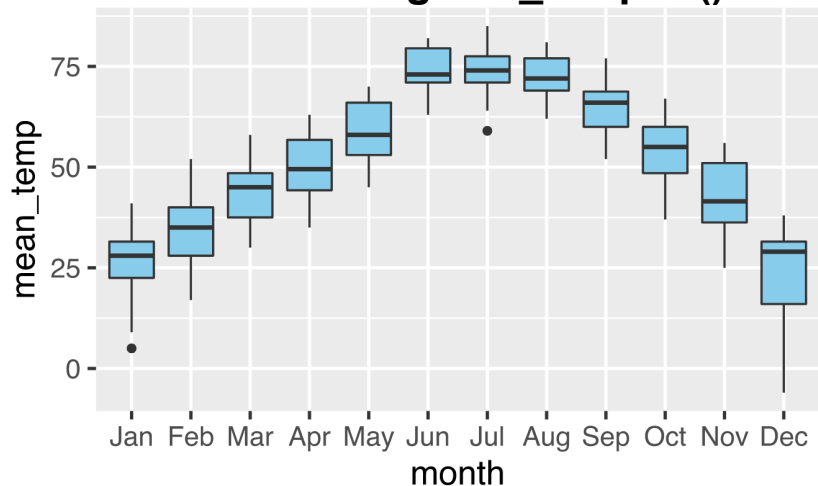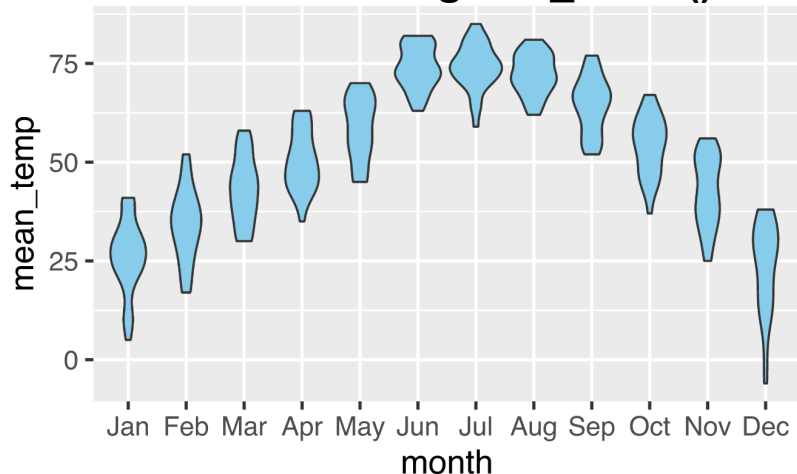# Better Ways of Comparing and Summarizing Distributions



Examples from Claus Wilke: https://wilkelab.org/SDS375/slides/visualizing-distributions-2.html#1
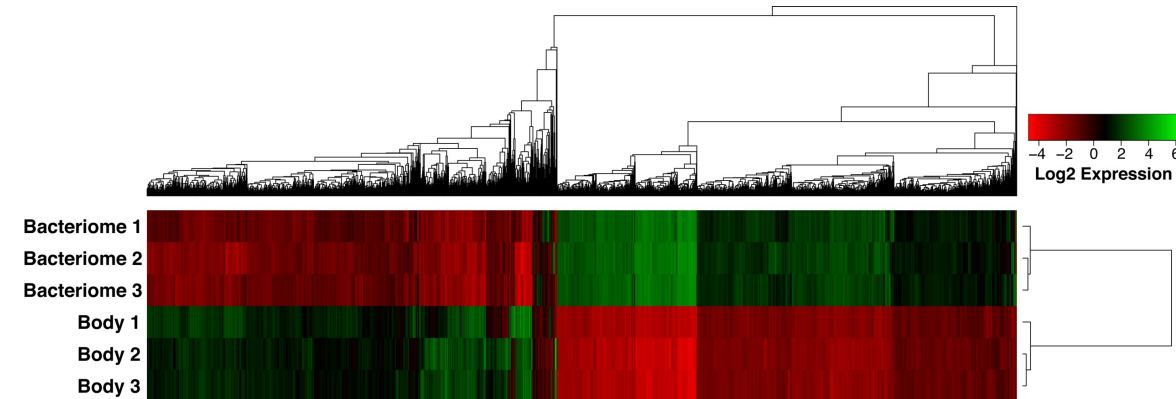
# Linear vs. Log Scales



- Use **linear** scales to emphasize **absolute** differences.

- Use **log** scales to emphasize **proportional** differences.

# Accessibility

Color is a powerful tool for visualizations, but it will not be perceived in the same way by everyone in your audience. Tips for making your visualizations accessible to color bind individuals….
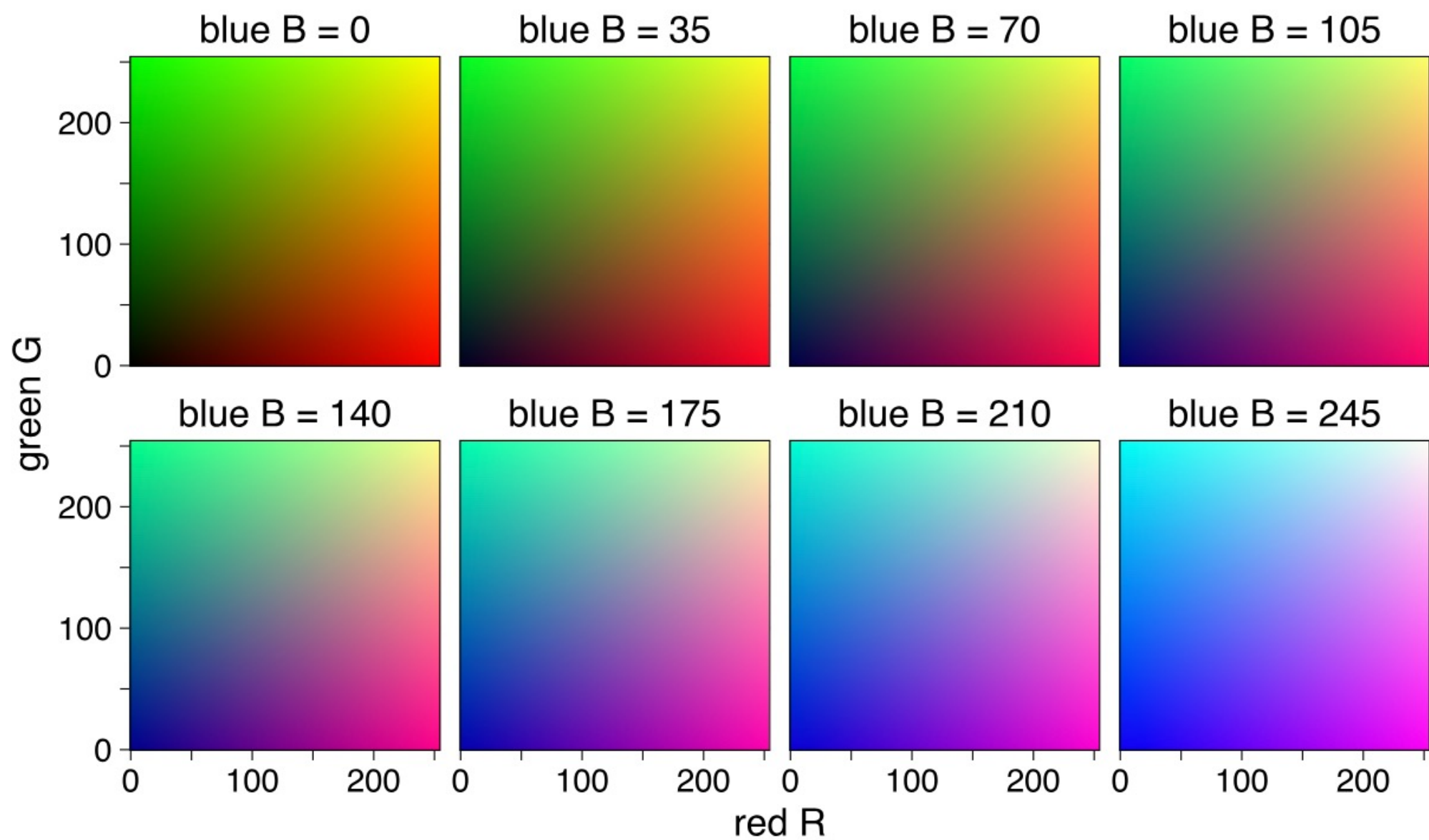


- Use palettes consisting of colors that are more distinguishable for individuals with common forms of color blindness.

- Use color and shape of points redundantly to distinguish among groups in plot.

```
> ggplot(MouseData, aes(x=age, y=weight, color=genotype, shape=genotype))
        + geom_point()
```

- Make use of figure labeling and legend descriptions to make the plot accessible even if colors are difficult to distinguish.
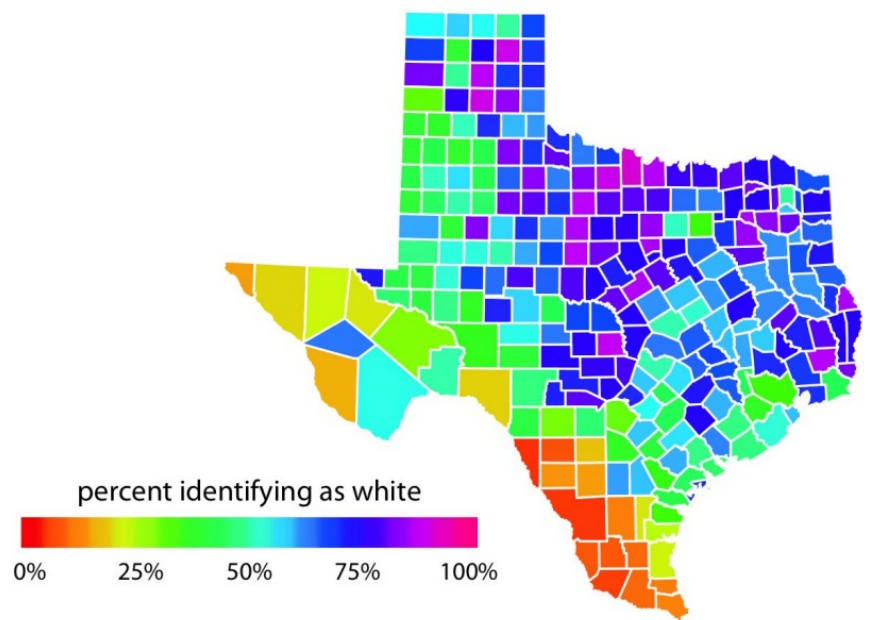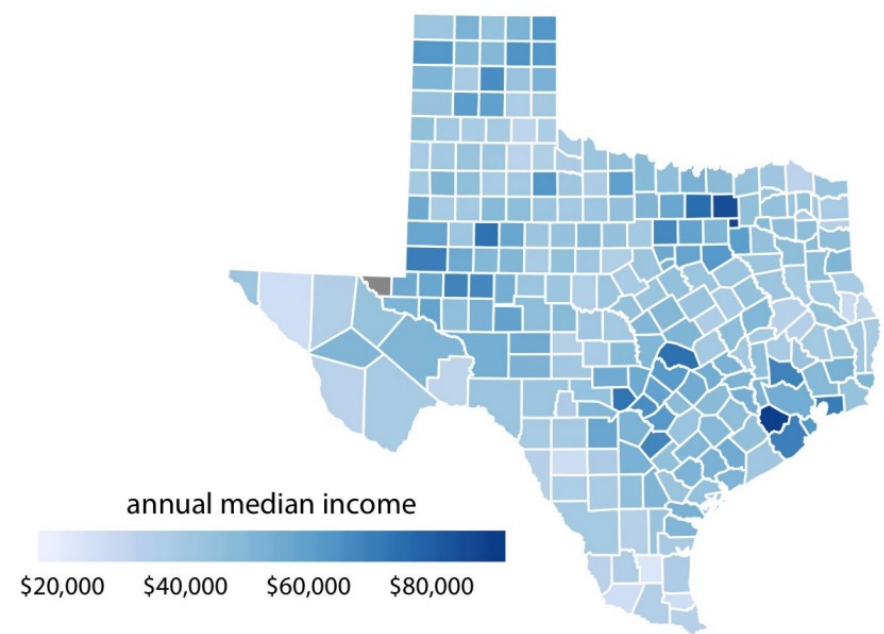
# Choosing Colors

# Go Easy on Our Eyes
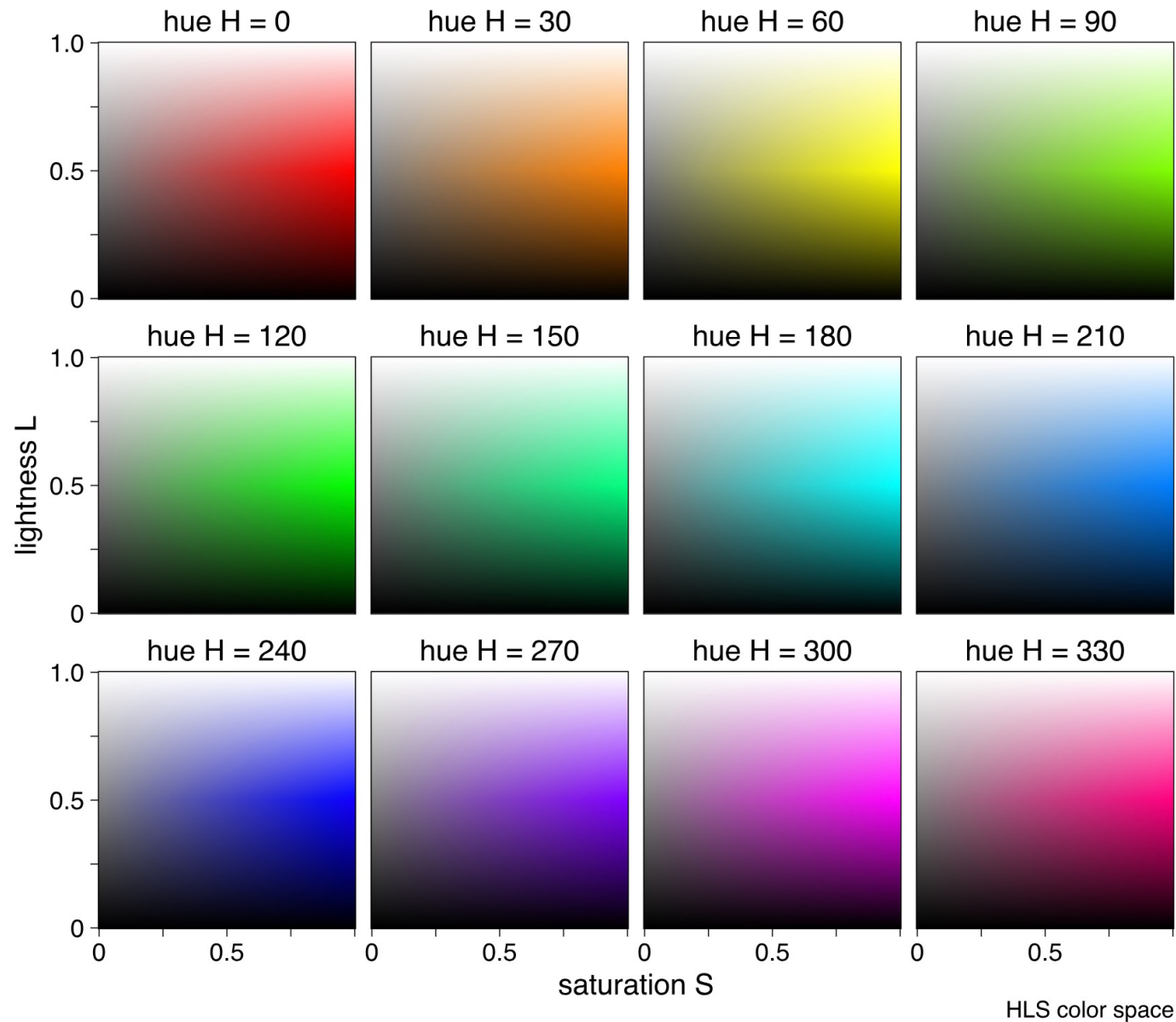
Generic named colors in R tend to have extremely high saturation.

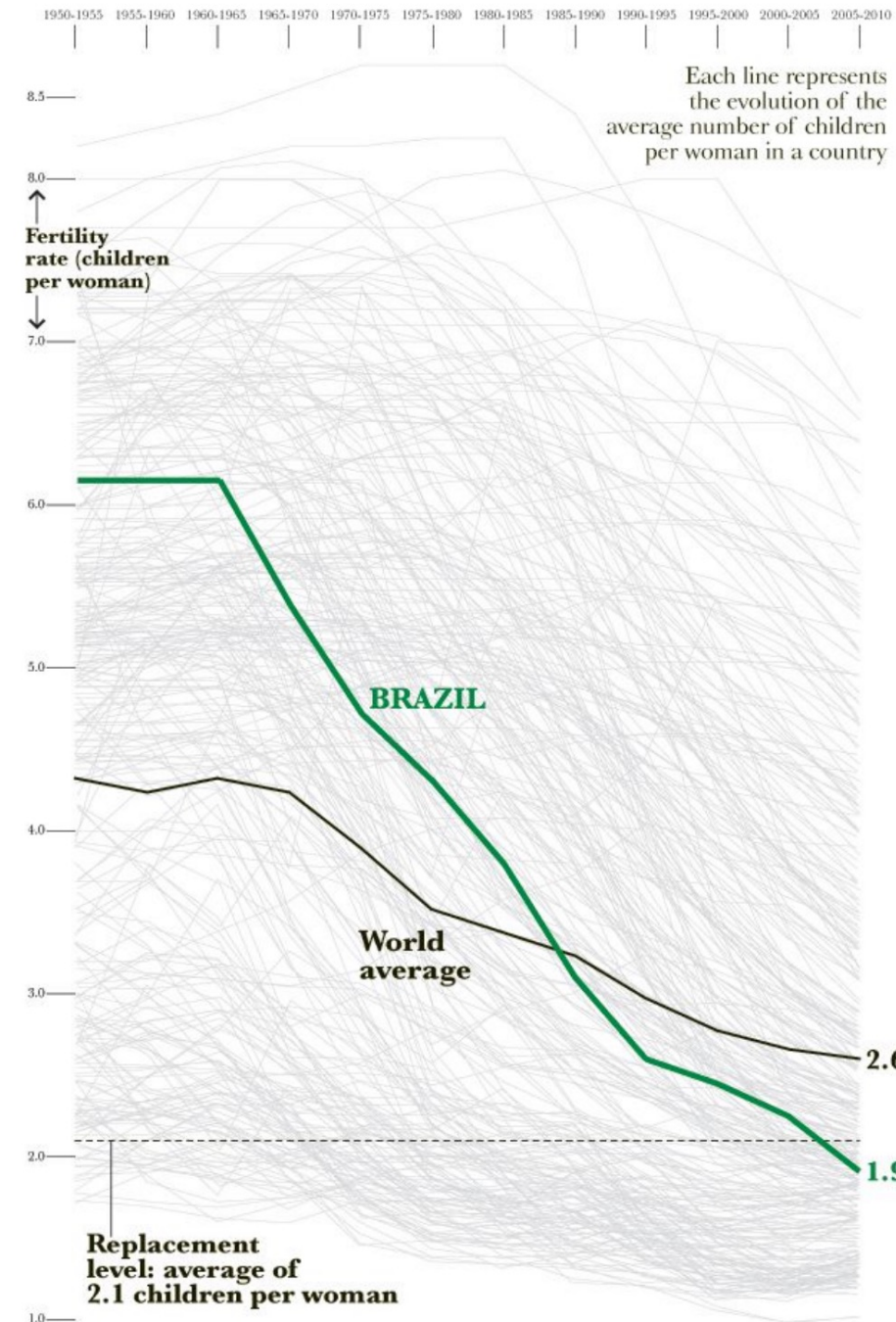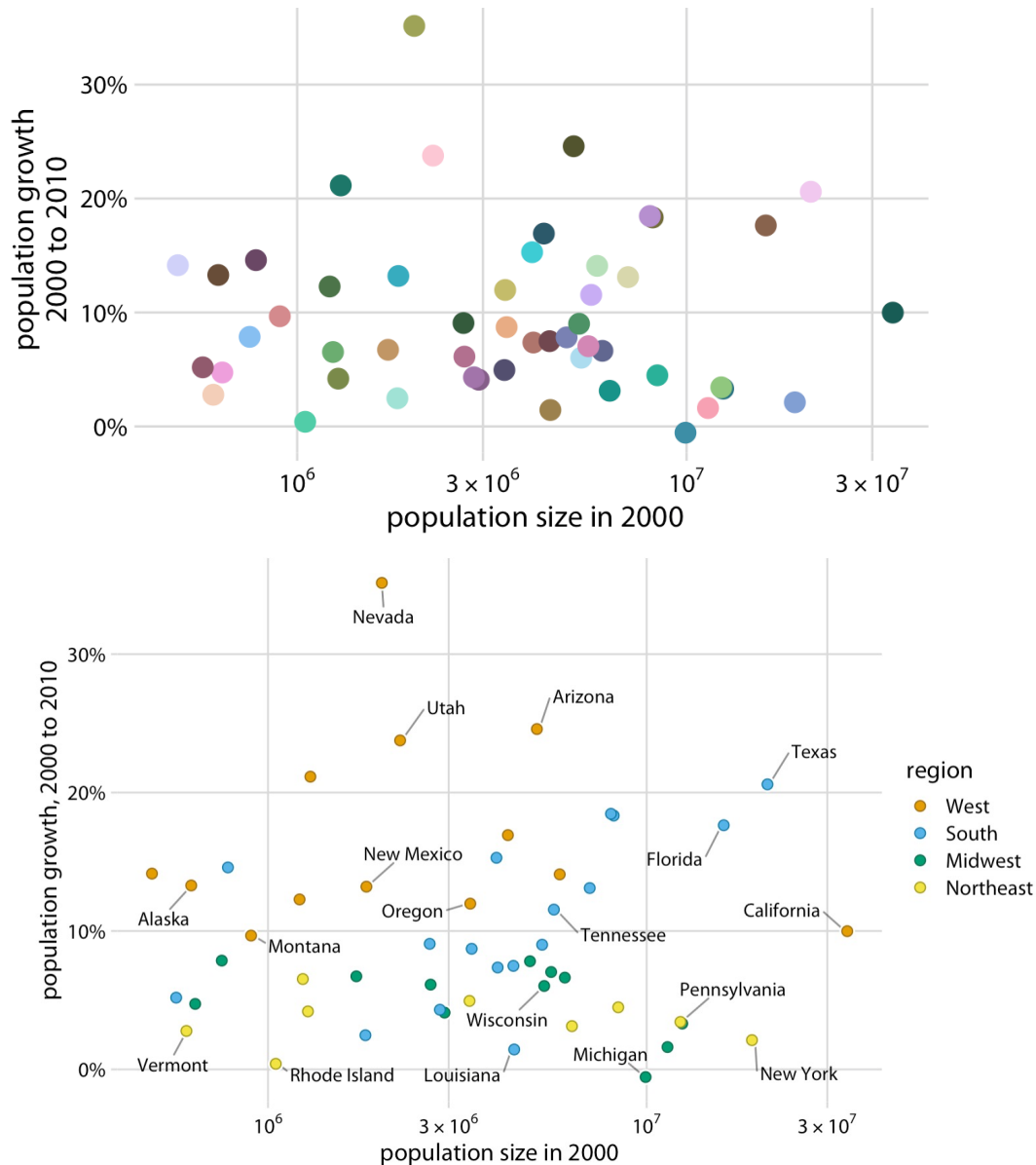| | |
|---|---|
| red |
| orange |
| yellow |
| green |
| blue |
| purple |

annual median income

$20,000    $40,000    $60,000    $80,000

percent identifying as white

0%    25%    50%    75%    100%

Jess Cohen-Tanugi

# Thinking of Colors in Terms of Hue, Lightness and Saturation



HLS color space

# Color to Emphasize, Not to Overwhelm

# Exercise and Assignment

https://dbsloan.github.io/CM515/SP24/ggplot/
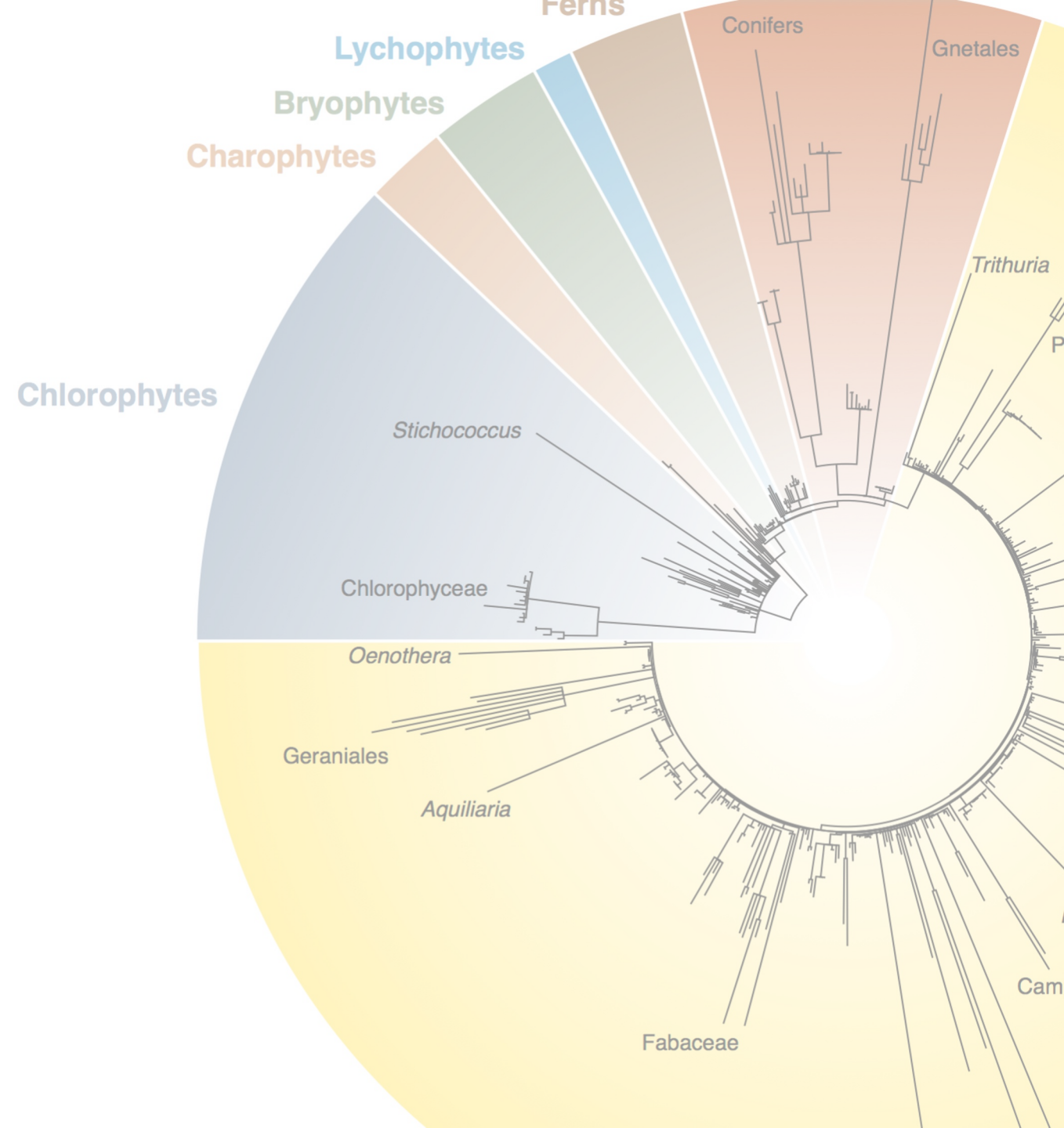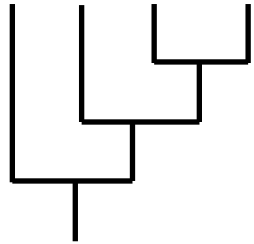
# Tree-Like Data Structures in Biology

- The tree of life and species relationships

- Gene family evolution

- Hierarchical clustering of gene expression patterns, ecological/microbiome communities, etc.

# Newick Tree Format

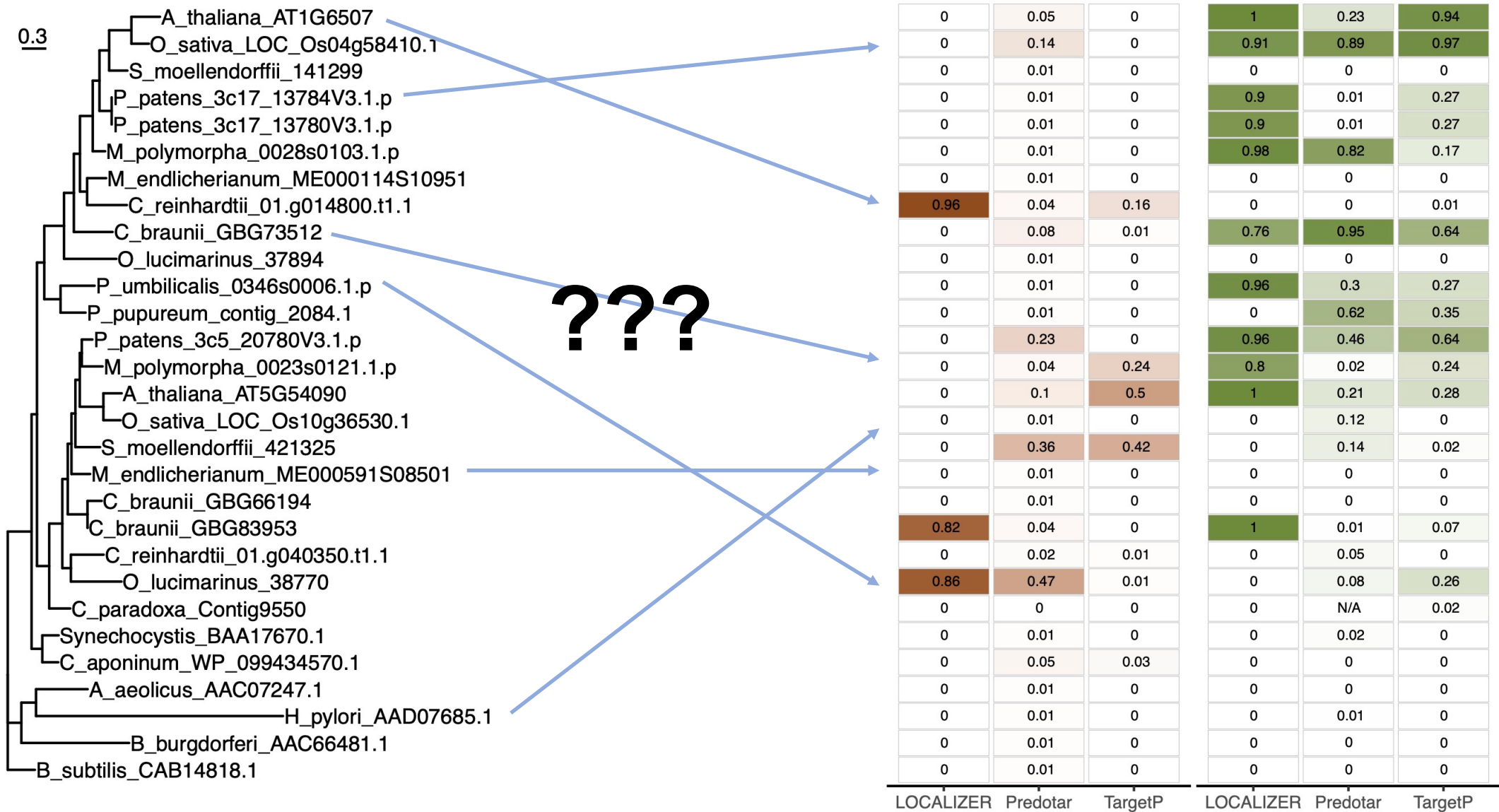- Standard parentheses/comma-based format for summarizing tree branching patterns

(A, (B, (C, D)));



- Numerical values and annotations can be added to record features such as branch lengths, statistical (e.g., bootstrap) support, node/branch labels, and other features.

# Linking Trees and Data Plots

# Exercise and Assignment

[https://dbsloan.github.io/CM515/SP24/ggtree/](https://dbsloan.github.io/CM515/SP24/ggtree/)