

MESA WORKSHOP

PART 3

GRAMMAR OF GRAPHICS

All Mentions

★ Hadley Wickham favourited your post

 baptiste
@baptnz

19h

preparing tomorrow's talk on
"Grammar of graphics" and the accompanying workshop
"swearing with Python"

0     ...



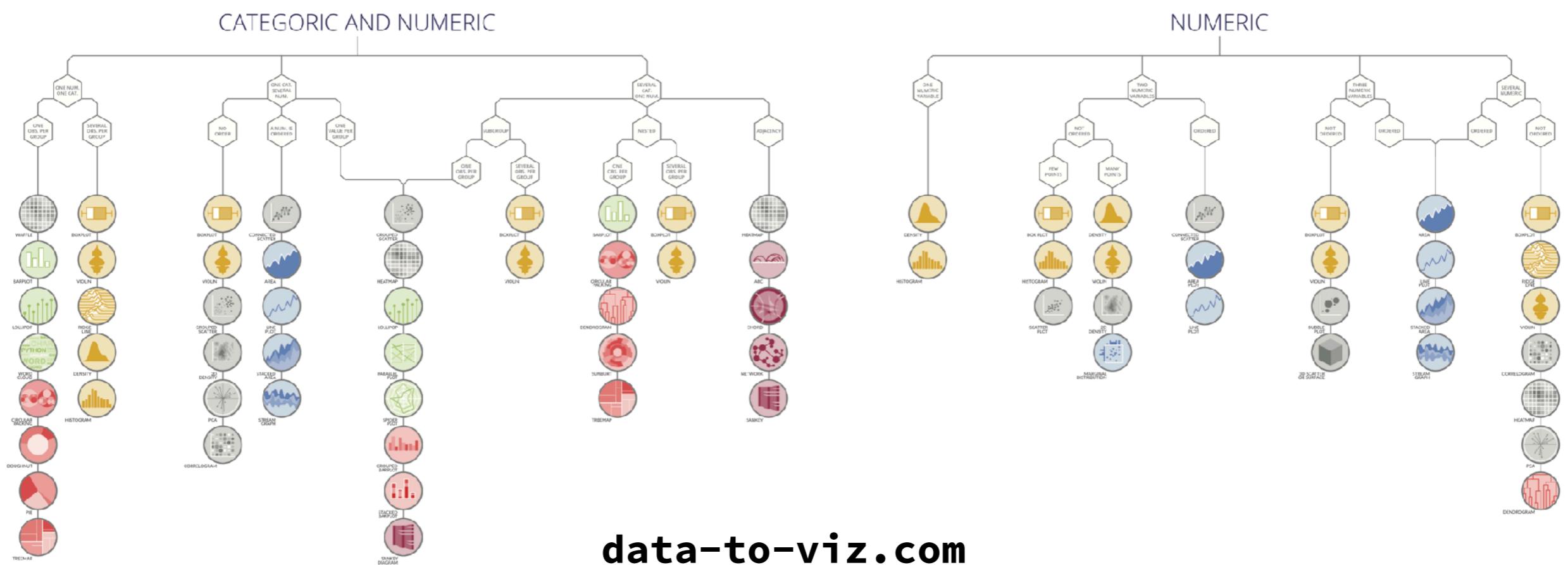
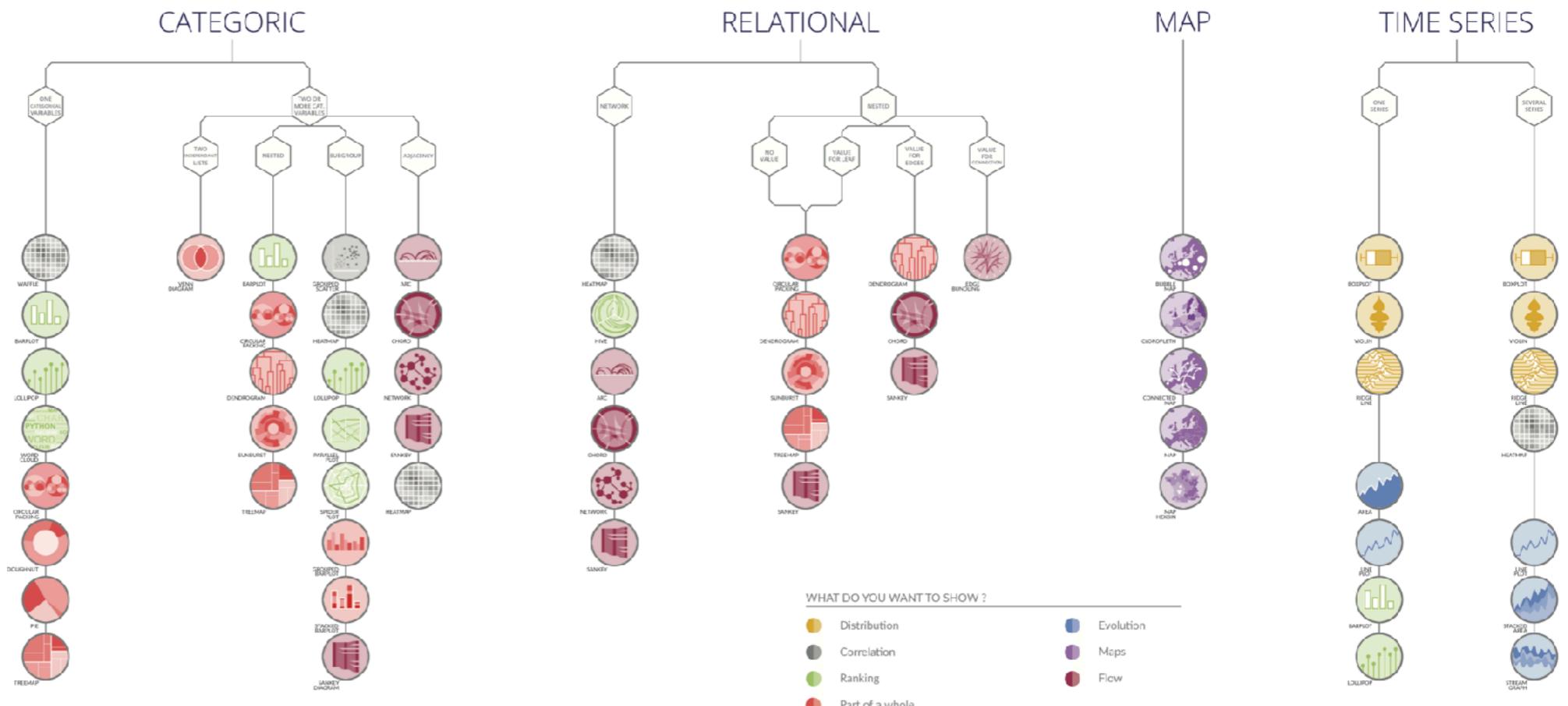
from Data to Viz

'From Data to Viz' is a classification of chart types based on input data format. It will help you find the perfect chart in three simple steps :

- 1 Identify what type of data you have.
 - 2 Go to the corresponding decision tree and follow it down to a set of possible charts.
 - 3 Choose the chart from the set that will suit your data and your needs best.

Dataviz is a world with endless possibilities and this project does not claim to be exhaustive. However it should provide you with a good starting point. For an interactive version and much more, visit:

data-to-viz.com



TAXONOMY OF GRAPHICS

Deviation

Emphasise variations ($>/<$) from a fixed reference point. Typically the reference point is zero but it can also be a target or a long-term average. Can also be used to show sentiment (positive/neutral/negative).

Example FT uses: Trade surplus/deficit, climate change

Correlation

Show the relationship between two or more variables. Be mindful that, unless you have a causal link, correlations will assume the relationships you show them to be causal (i.e. one causes the other).

Example FT uses: Inflation and unemployment, income and life expectancy

Ranking

Use where an item's position in an ordered list is more important than its absolute value. Readers will assume the relationships you show them to be causal (i.e. one causes the other).

Example FT uses: Wealth, deprivation, league tables, constituency election results

Distribution

Show values in a dataset and how often they occur. The shape (or 'skew') of a distribution is another simple way of highlighting the lack of uniformity or equality in the data.

Example FT uses: Income distribution, population, Gdp(s)/distribution, revealing inequality

Change over Time

Give emphasis to changing trends. These can be short (intra-day) movements or when longer traversing decades or centuries. Choosing the correct time period is important to provide suitable context for the reader.

Example FT uses: Share price movements, economic time series, sectoral changes in a market

Magnitude

Show size comparisons. These can be relative (just being able to see larger/smaller) or absolute to see fine differences. Usually these show a 'counted' number (for example, barrels, dollars or people) rather than a calculated rate or per cent.

Example FT uses: Commodity production, market capitalisation, volumes in general

Part-to-whole

Show how a single entity can be broken down into its component elements. If the reader's interest is solely in the size of the components, consider a magnitude-type chart instead.

Example FT uses: Fiscal budgets, company structures, national election results

Spatial

A aside from locator maps only used when precise locations or geographical patterns in space are more important to the reader than anything else.

Example FT uses: Population density, natural resource locations, natural disaster risk/impact, catchment areas, variation in election results

Flow

Show the reader volumes or intensity of movement between two or more states or conditions. These can be logical sequences or geographical locations.

Example FT uses: Movement of funds, trade, migrants, lawsuits, information/relationship graphs.

Diverging bar

A simple standard bar chart that can handle both negative and positive magnitude values.

Scatterplot

The standard way to show the relationship between two continuous variables, each of which has its own axis.

Ordered bar

Standard bar charts display ranks of values much more easily when sorted into order.

Histogram

The standard way to show a statistical distribution - keep the gaps between columns small to highlight the 'shape' of the data.

Line

The standard way to show a changing time series. If data are irregular, consider markers to represent data points.

Column

The standard way to compare the size of things. Must always start at 0 on the axis.

Bar

See above. Good when the data are not time series and labels have long category names.

Stacked column/bar

A simple way of showing part-to-whole relationships but can be difficult to read with more than a few components.

Basic choropleth (rate/ratio)

The standard approach for putting data on a map – may also be rates rather than totals and use a sensible base geography.

Column + line timeline

A good way of showing the relationship between an amount (column) and a rate (line).

Ordered column

See above.

Dot plot

A simple way of showing the change or range (min/max) of data across multiple categories.

Dot strip plot

Good for showing individual values in a distribution, can be a problem when too many dots have the same value.

Barcode plot

Like dot strip plots, good for displaying all the data in a table, they work best when highlighting individual values.

Column + line timeline

A good way of showing the relationship over time between an amount (column) and a rate (line).

Column

Columns work well for showing change over time – but usually best with only one series of data at a time.

Paired column

As per standard column but allows for multiple series. Can become tricky to read with more than 2 series.

Marimekko

A good way of showing the size and proportion of data at the same time – as long as the data are not too complicated.

Pie

A common way of showing part-to-whole data – but be aware that it's difficult to accurately compare the size of the segments.

Flow map

For showing unambiguous movement across a map.

Surplus/deficit bar

Shows a single value into two contrasting components (e.g. male/female).

Bubble

Like a scatterplot, but adds additional detail by sizing the circles according to a third variable.

XY heatmap

A good way of showing the patterns between 2 categories of data, less effective at showing fine differences in amounts.

Bubble

Perfect for showing how ranks have changed over time or vary between categories.

Slope

Perfect for showing changing rankings across multiple dates. For large datasets, consider grouping lines using colour.

Lollipop

Lollipops draw more attention to the data points than standard bar/column and can also show rank and value effectively.

Bump

Effective for showing changing rankings across multiple dates. For large datasets, consider grouping lines using colour.

Cumulative curve

A standard way for showing the age and sex breakdown of a population distribution; effectively back to back histograms.

Population pyramid

A good way of showing how unequal a distribution the y axis is always cumulative frequency, x axis is always a measure.

Frequency polygons

For displaying multiple distributions of data. Like a regular line chart, best limited to a maximum of 3 or 4 datasets.

Beeswarm

Use to emphasise individual points in a distribution. Points can be sized to an additional variable. Best with medium-sized datasets.

Scatterplot

A good way of showing changing data for two variables; however there is a relatively clear pattern of progression.

Connected scatterplot

A good way of showing changing data for two variables; however there is a relatively clear pattern of progression.

Calendar heatmap

A great way of showing temporal patterns (daily, weekly, monthly) – at the expense of showing precision in quantity.

Radar

A space-efficient way of showing value of multiple variables – but make sure they are organised in a way that makes sense to the reader.

Isotype (pictogram)

Excellent solution in some instances – use only whole numbers (do not slice off an arm to represent a decimal).

Fan chart (projection)

Use to show the uncertainty in future projections – usually this grows the further forward to projection.

Lollipop

Lollipop charts draw more attention to the data points than standard bar/column – does not have to start at zero (but preferable).

Gridplot

Good for showing % data – especially when working with whole numbers and work well in small multiple layout form.

Voronoi

A way of turning points into areas – any point within each area is closer to the central point than any other centroid.

Arc

A hemisphere, often used for visualising parliamentary composition by number of seats.

Dot density

Used to show the location of individual events/locations – make sure to annotate any patterns the reader should see.

Bullet

Good for showing a measurement against the context of a target or performance range.

Vertical timeline

Presents time on the Y axis. Good for displaying detailed time series that work especially well when scrolling on mobile.

Grouped symbol

An alternative to bar/column charts when being able to count data or highlight individual elements is useful.

Waterfall

Can be useful for showing part-to-whole relationships where some of the components are negative.

Seismogram

Another alternative to the circle timeline for showing series where there are big variations in the data.

Flow map

Grid-based data values mapped to a map using a color scale. As choropleth map – but not snapped to an admin/political unit.

Contour map

For showing areas of equal value on a map. Can use deviation colour schemes for showing $>/<$ values.

Equalised cartogram

Converting each unit on a map to a square and equally-sized shape – good for representing voting regions with equal value.

Scaled cartogram (value)

Stretching and shrinking a map so that each area is sized according to a particular value.

Dot map

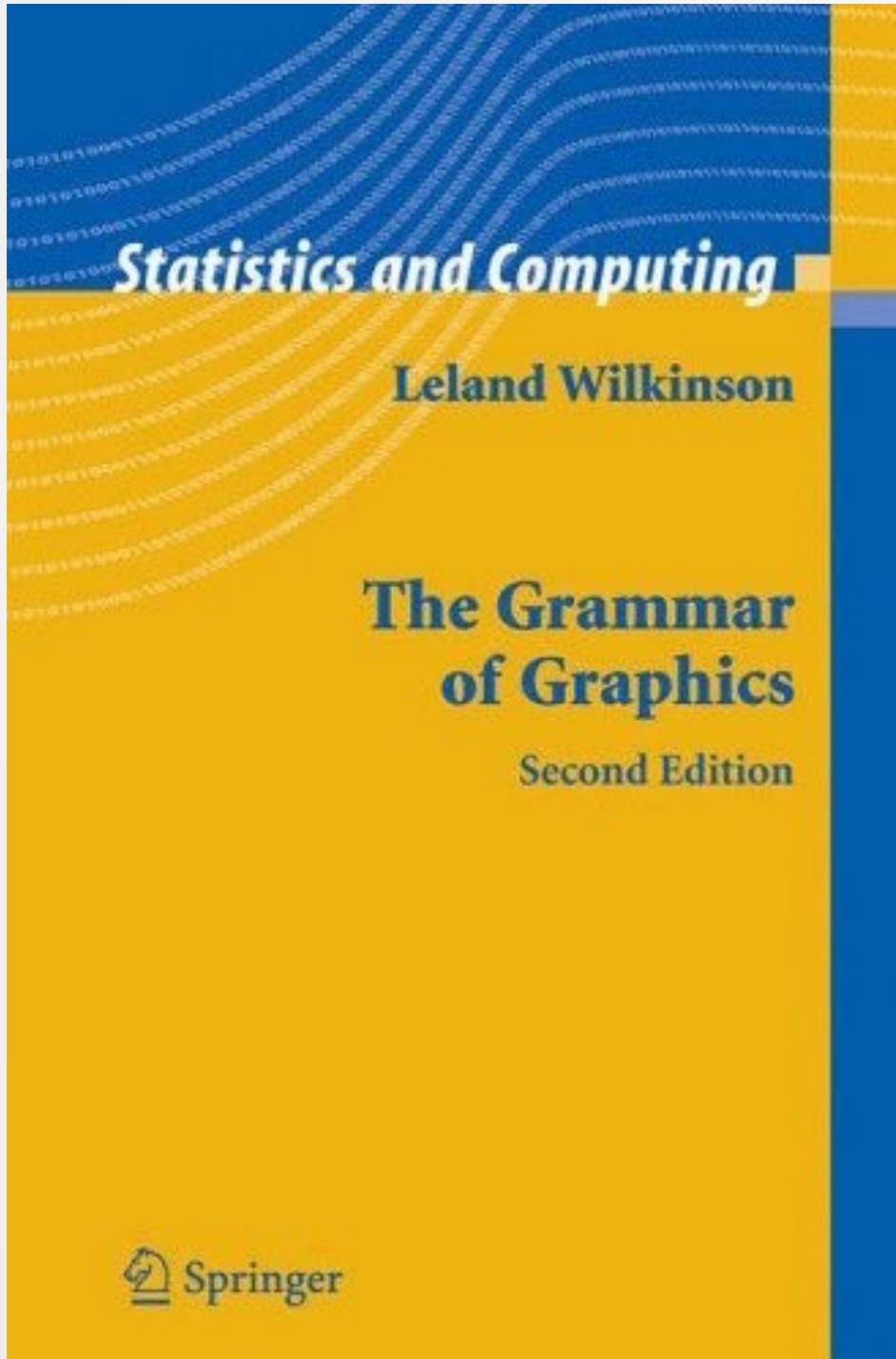
Used to show the location of individual events/locations – make sure to annotate any patterns the reader should see.

Heat map

Grid-based data values mapped to a map using a color scale. As choropleth map – but not snapped to an admin/political unit.

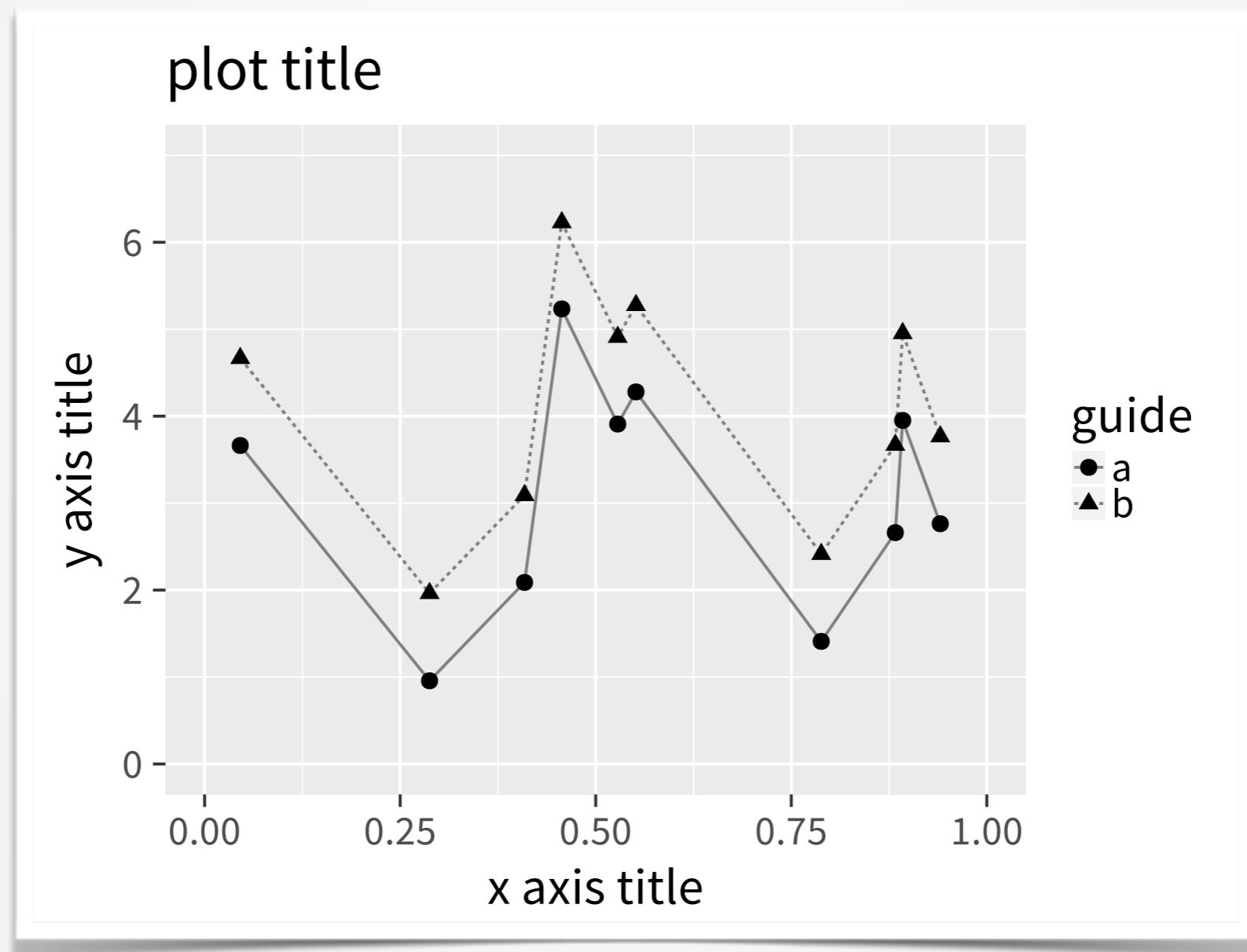
Source:
ft.com/vocabulary

GRAMMAR OF GRAPHICS

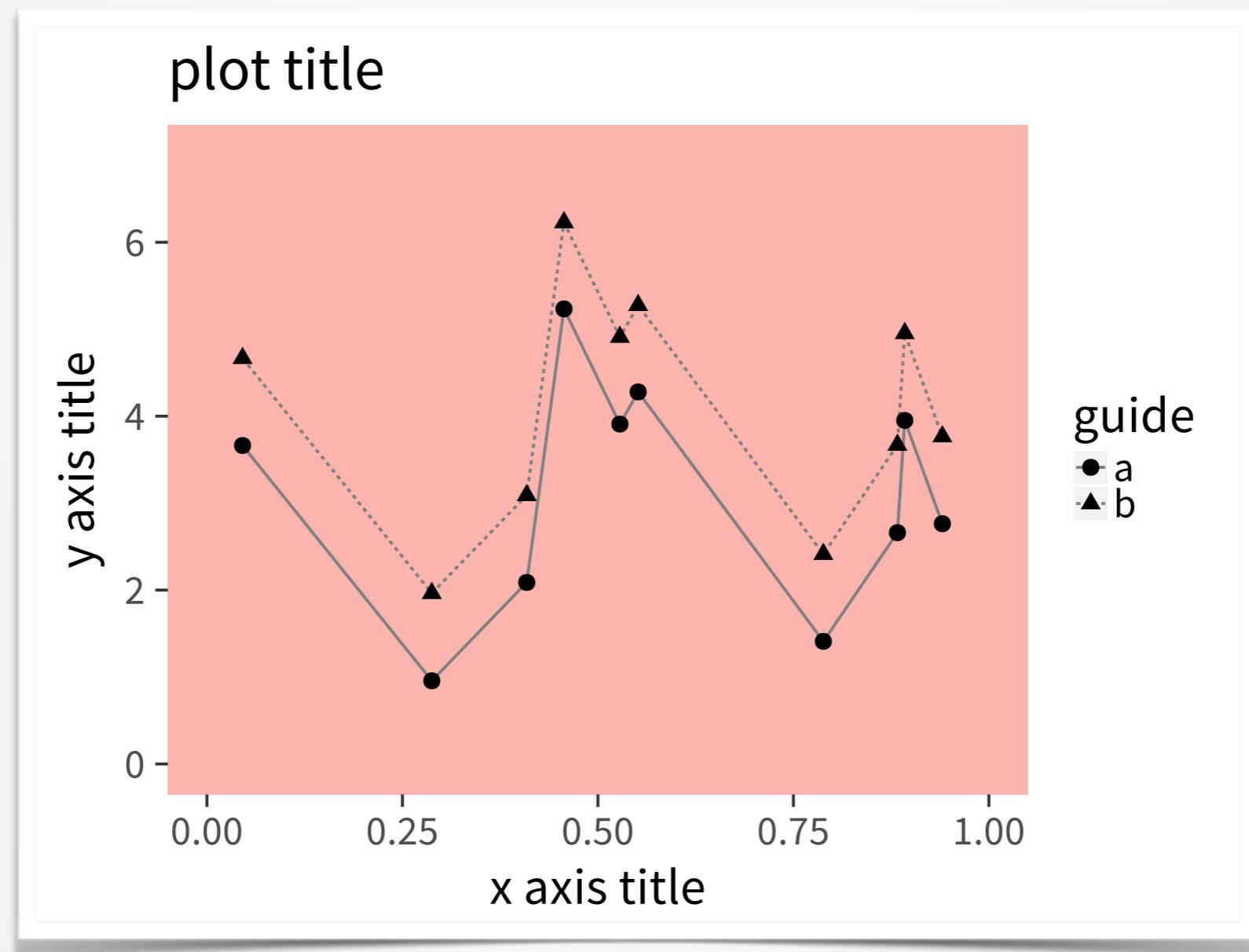


If charts are maps
of abstract worlds
the guiding principles
of graphics usage
could be derived
from the psychology
of perception. ”

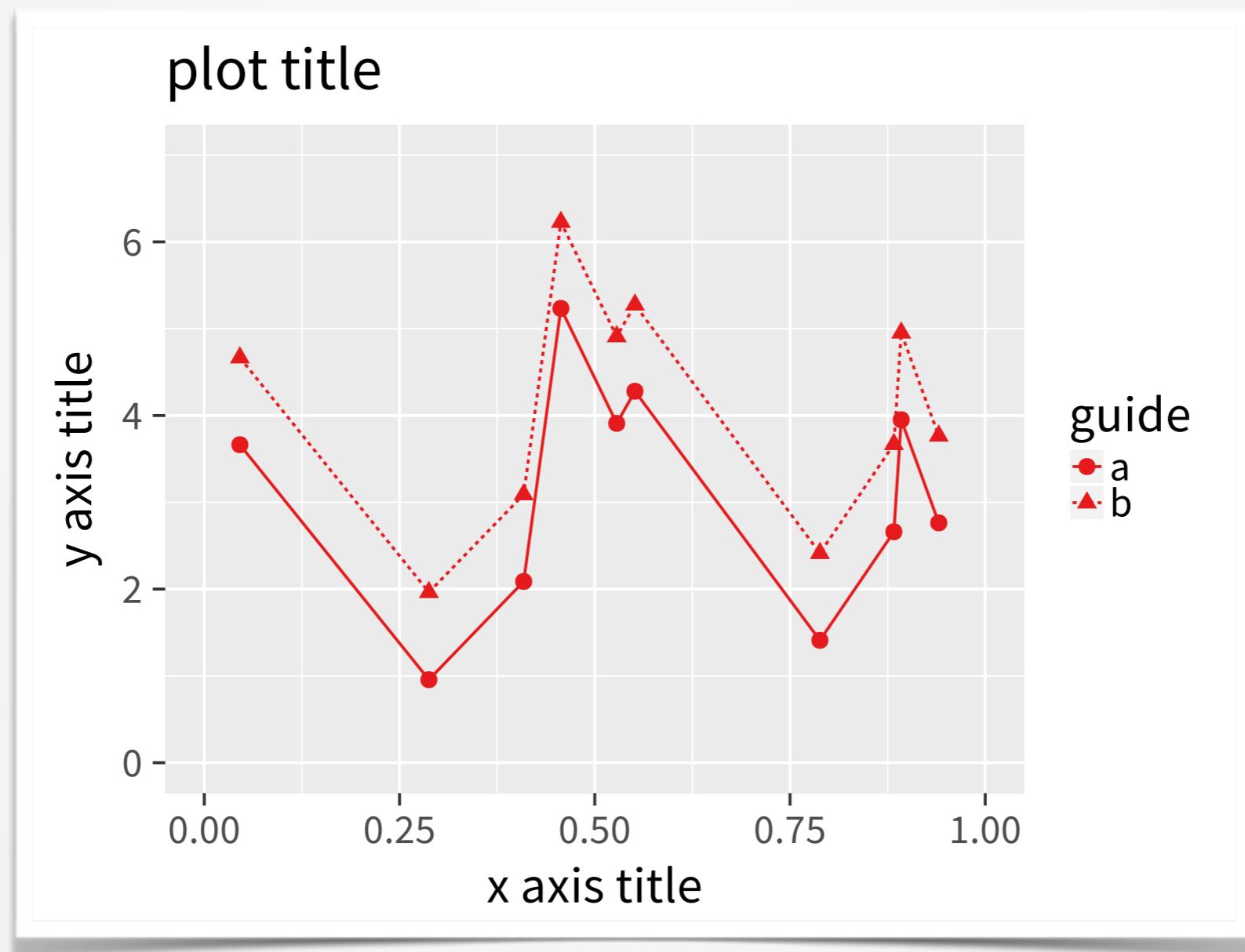
ANATOMY OF A PLOT



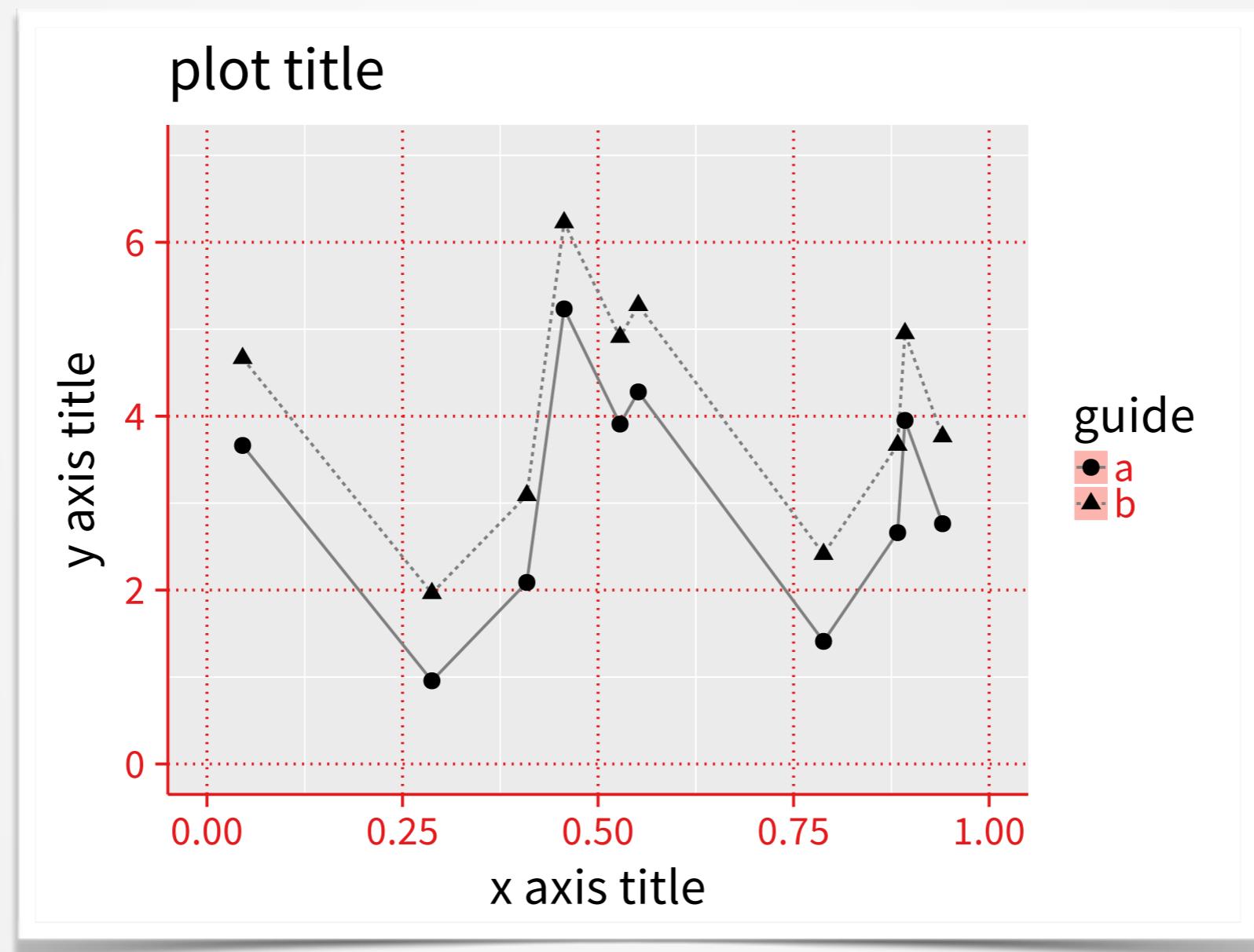
ANATOMY OF A PLOT • PLOT PANEL



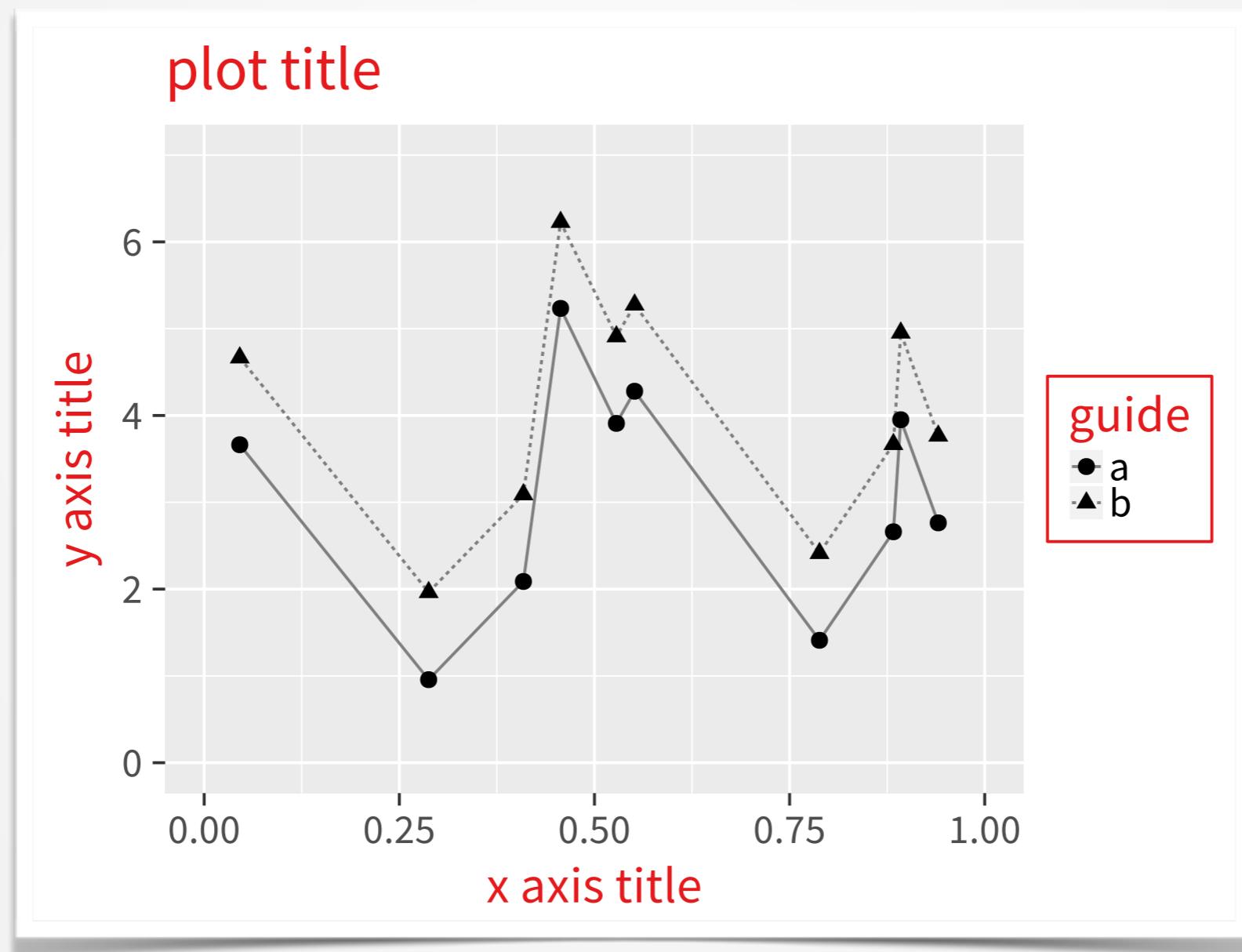
ANATOMY OF A PLOT • THE DATA



ANATOMY OF A PLOT • GUIDES



ANATOMY OF A PLOT • ANNOTATIONS

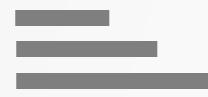


GRAMMAR OF GRAPHICS – MAPPING DATA

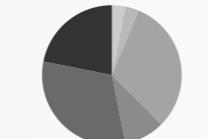
- ▶ **POSITION**



- ▶ **LENGTH**



- ▶ **ANGLE**



- ▶ **AREA**



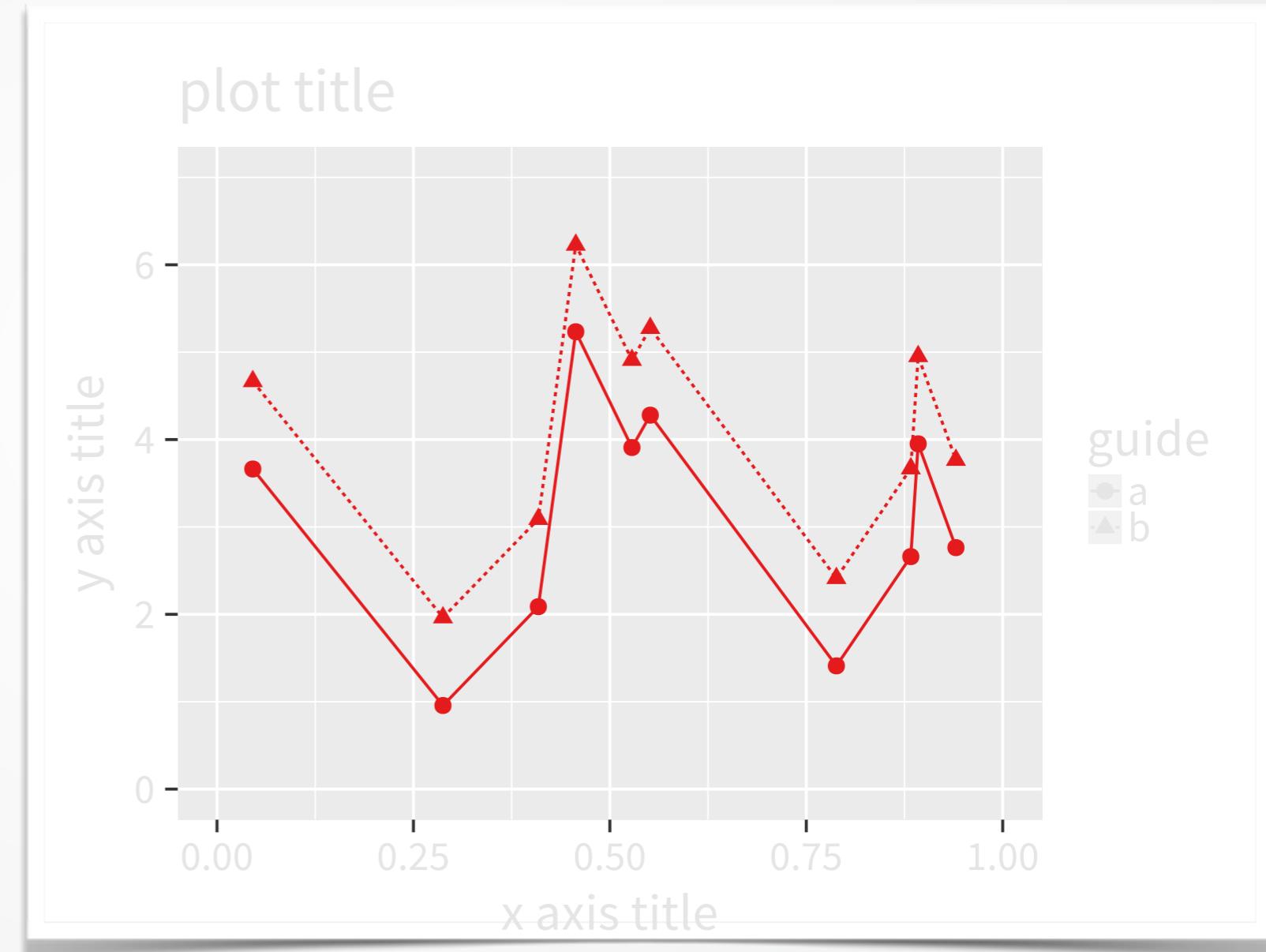
- ▶ **COLOUR**



- ▶ **SHAPE**



- ▶ **LINE TYPE, SIZE, TIME, ...**

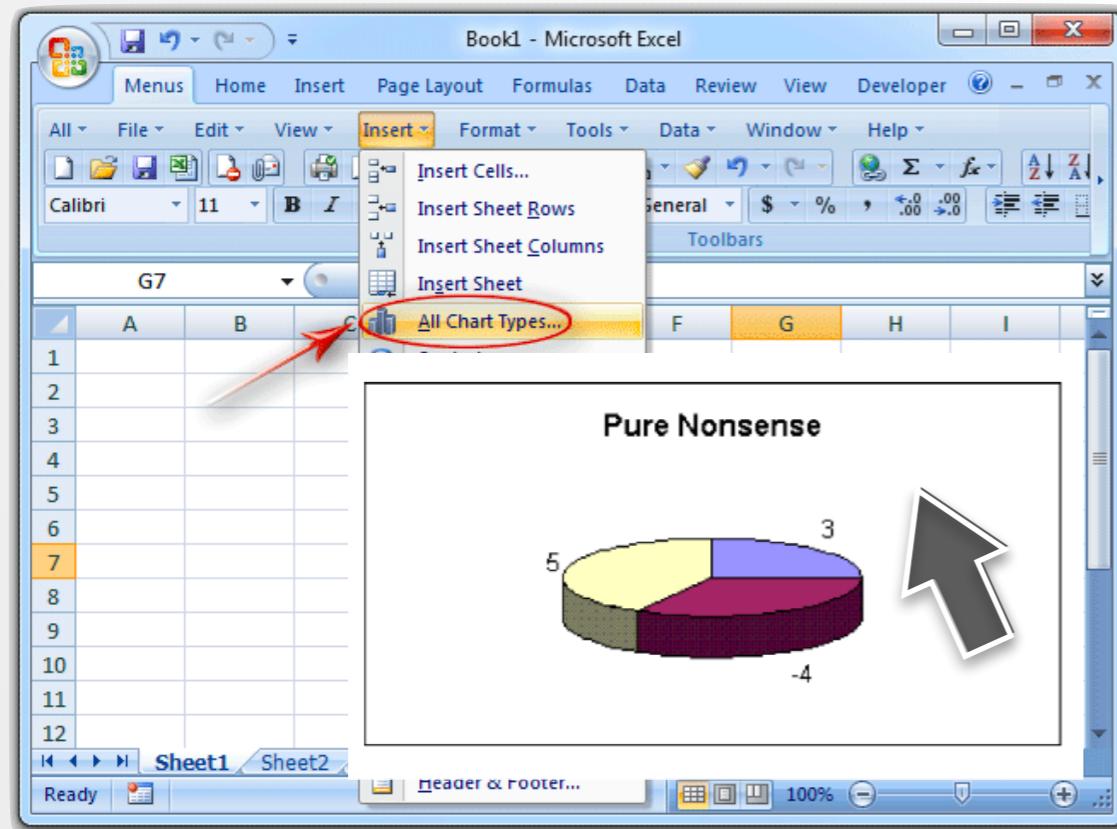


EXPRESSIVITY, LEGIBILITY, REPRODUCIBILITY

Point & Click

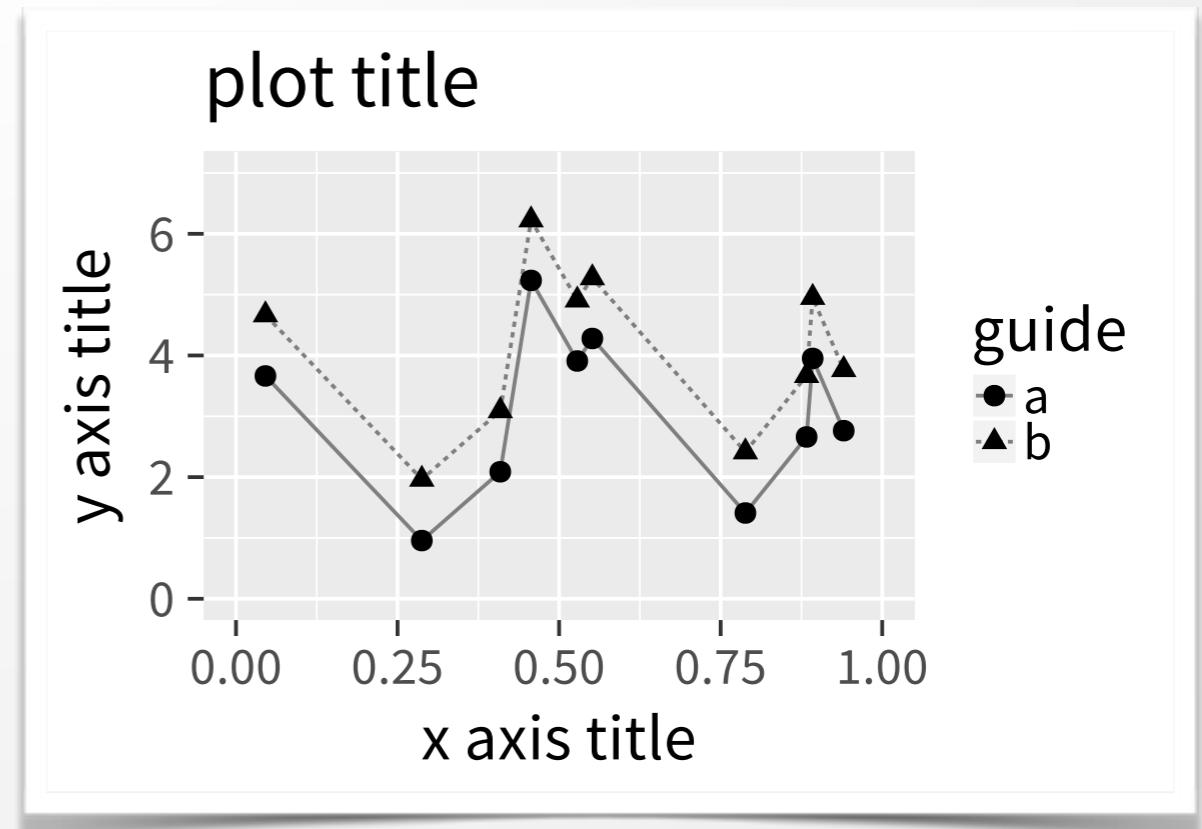
Yeah but, no but, yeah but, no but!!!
yeah but ... I swear * * * * * !!!?!!!
... but yeah _(ツ)_/ COMPUTER SAYS No

Ctrl-Z



Grammar of Graphics

```
plot(data, map(x, y)) +  
layer(point, map(shape = z)) +  
layer(line, map(linetype = t)) +  
theme(fontsize = 12)
```

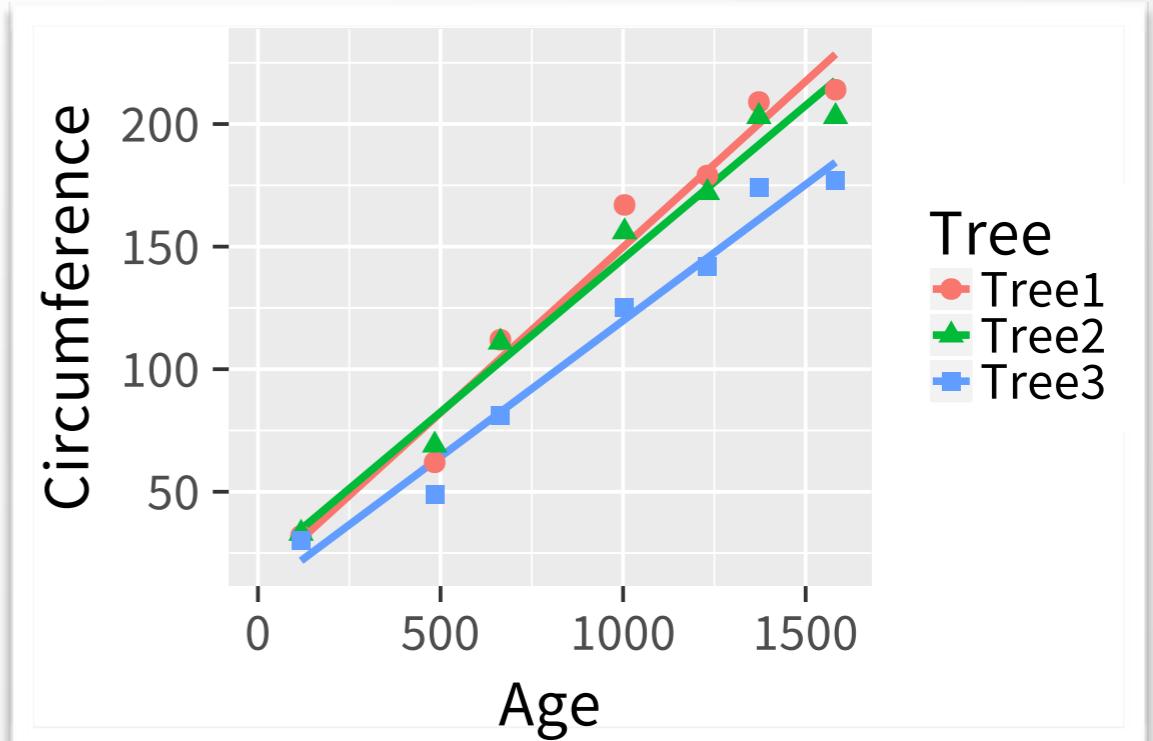


MAPPING DATA TO VISUAL ATTRIBUTES (GLYPHS)

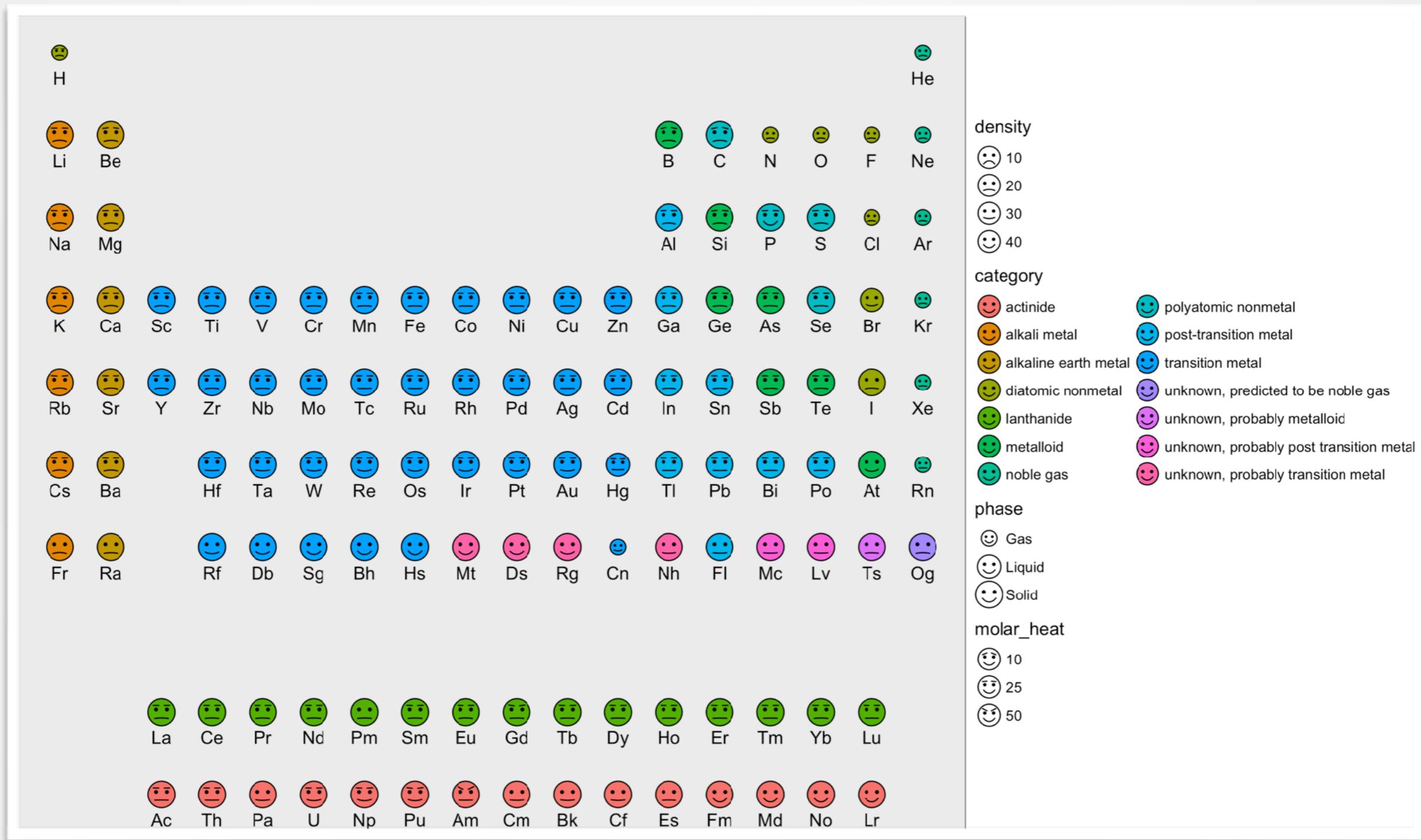
```
plot(data = d,  
      mapping = map(x = age,  
                     y = circumference)) +  
layer(type = point,  
      mapping = map(shape = Tree,  
                     colour = Tree)) +  
layer(type = line,  
      mapping = map(colour = Tree))
```

	Tree	age	circ.
1	Tree1	1582	214
2	Tree1	118	32
3	Tree2	118	33
4	Tree2	1372	203
5	Tree3	484	49
6	Tree3	1372	174
7	Tree3	1004	125

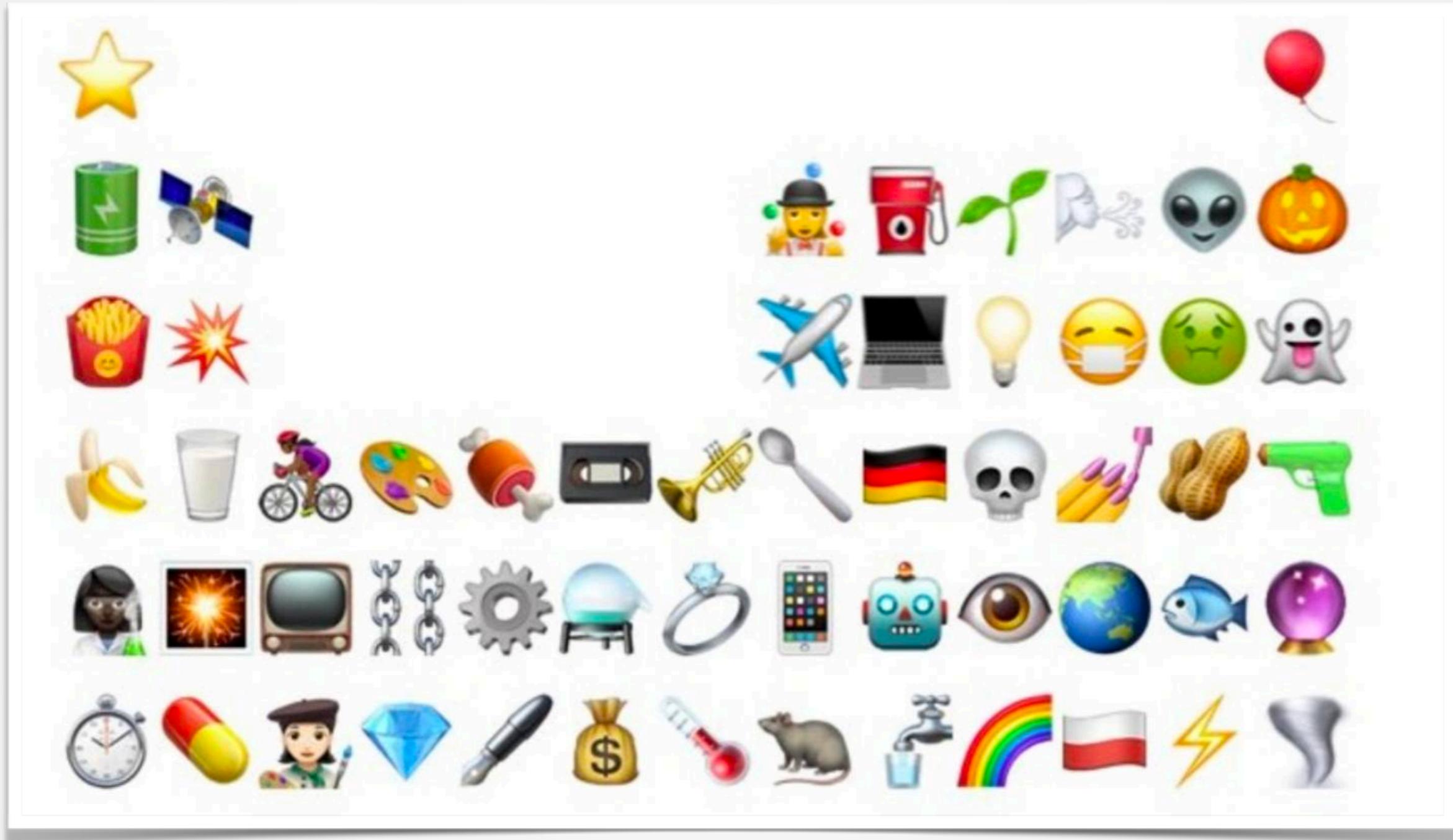
mapping
guides



MAPPING DATA TO VISUAL ATTRIBUTES



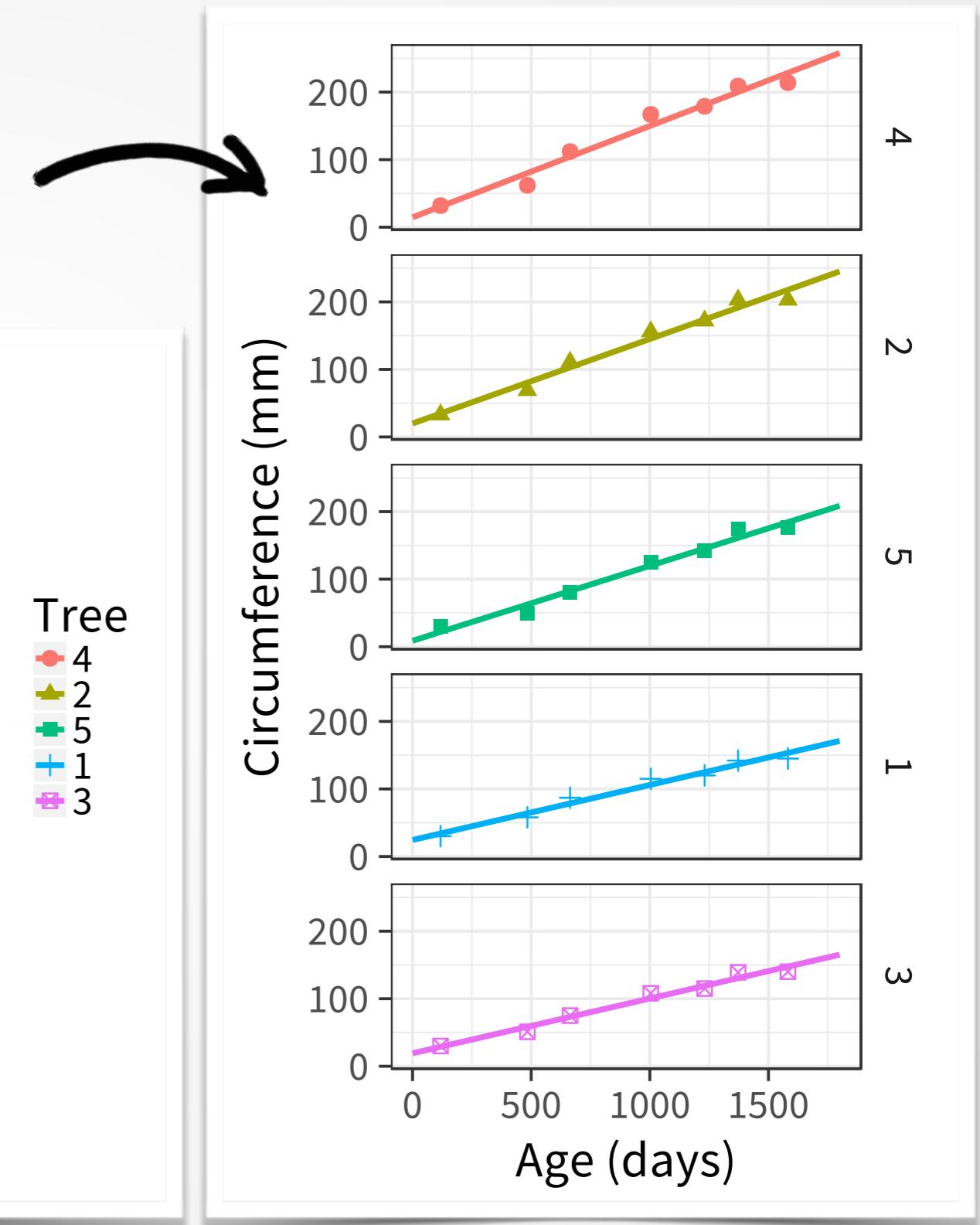
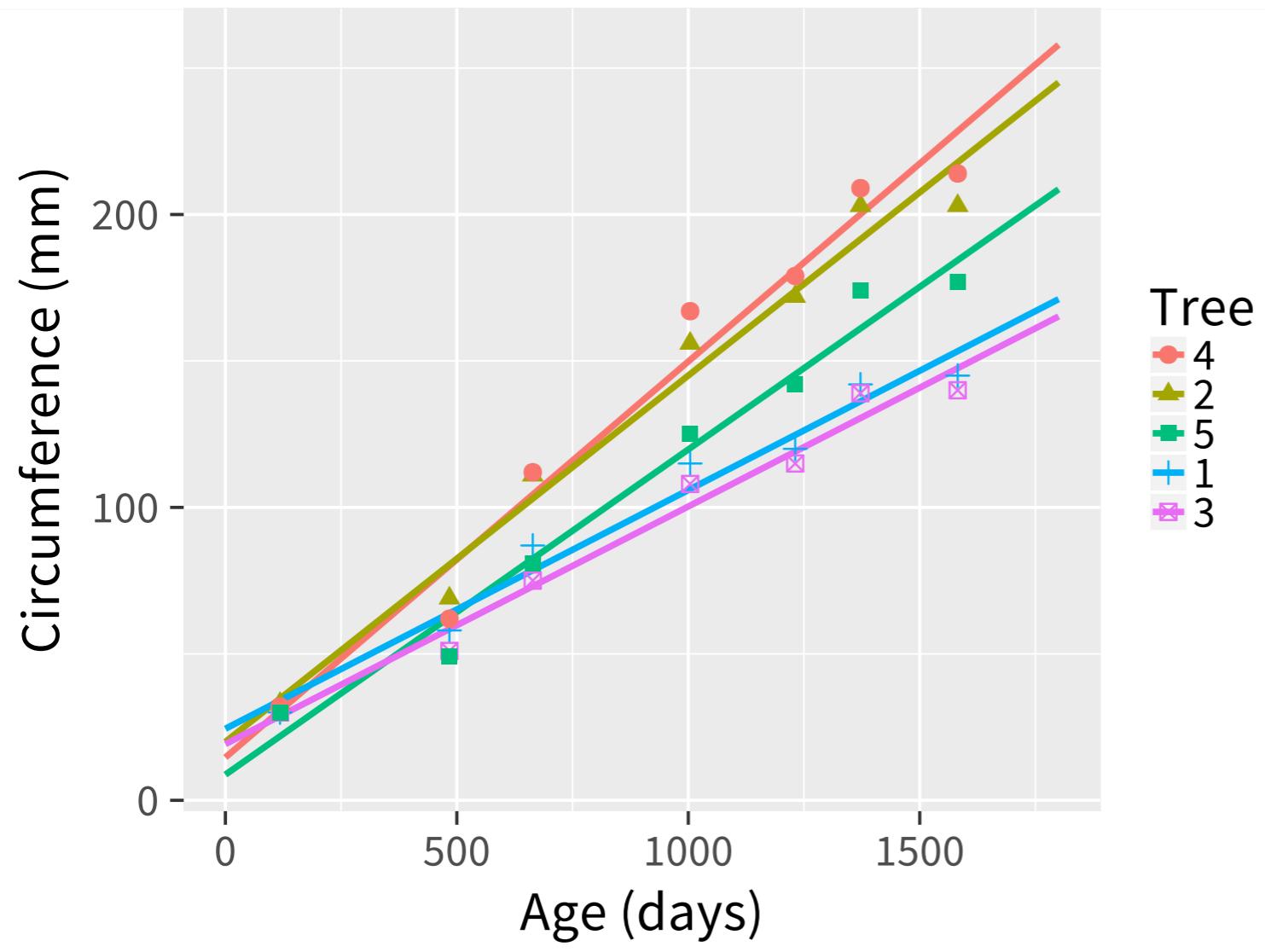
CONSIDER DIFFERENT VISUAL ATTRIBUTES



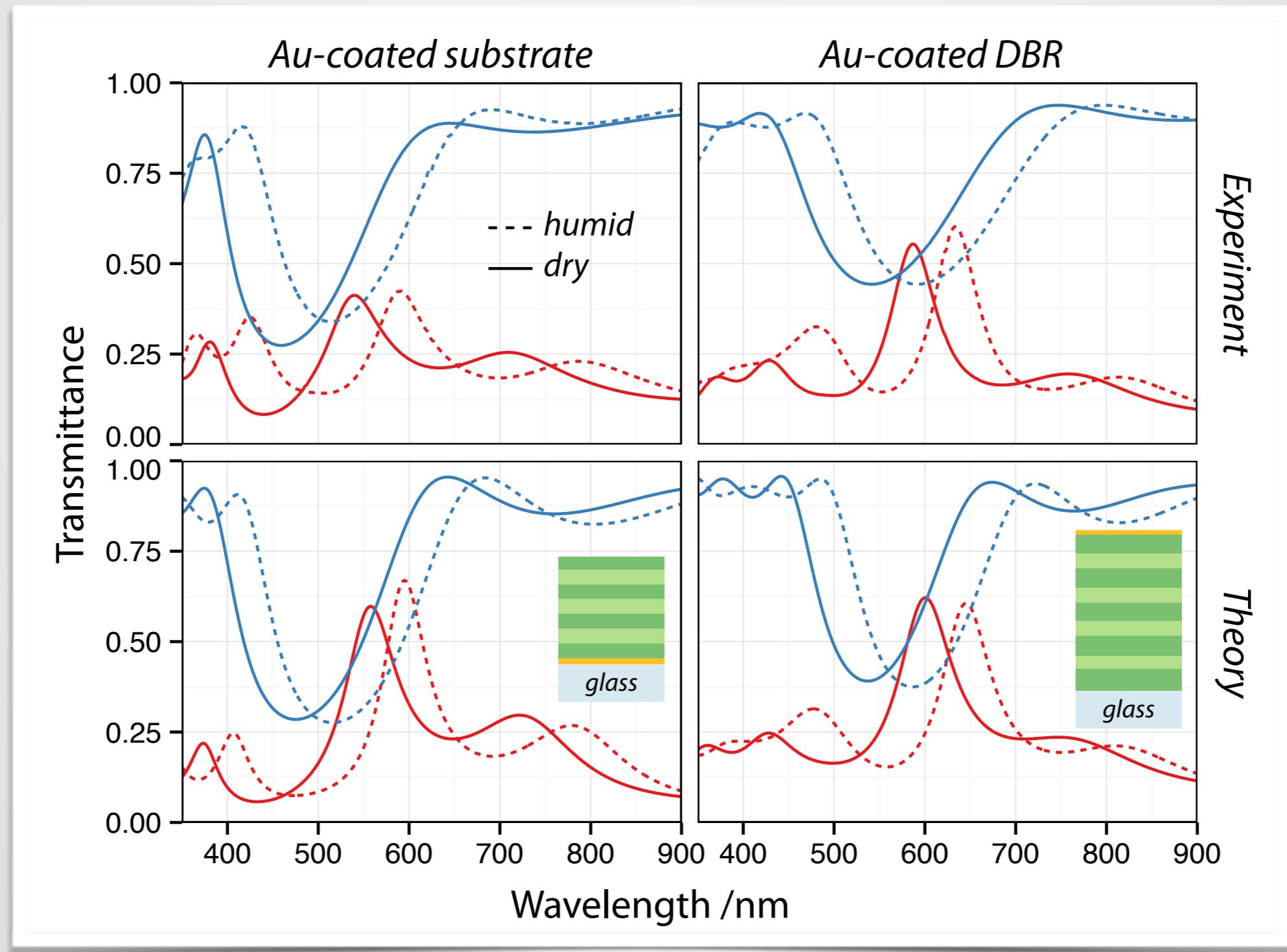
credit: @nicgaston

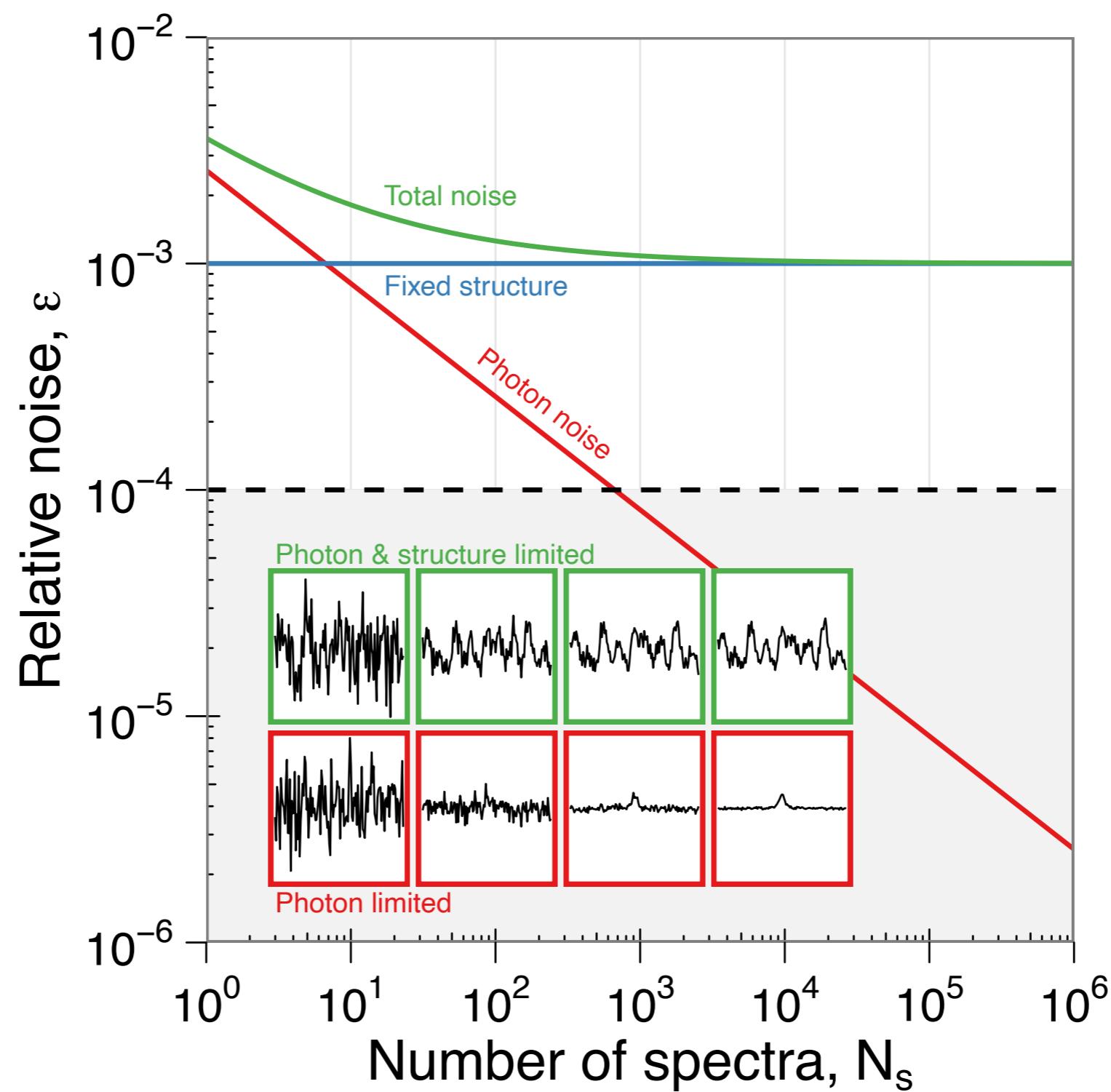
GRAPHICAL EXPLORATIONS – D.R.Y. PRINCIPLE

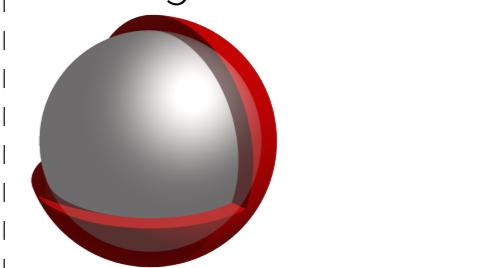
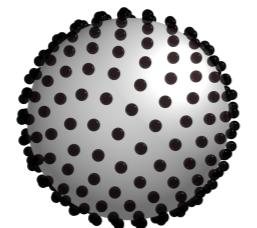
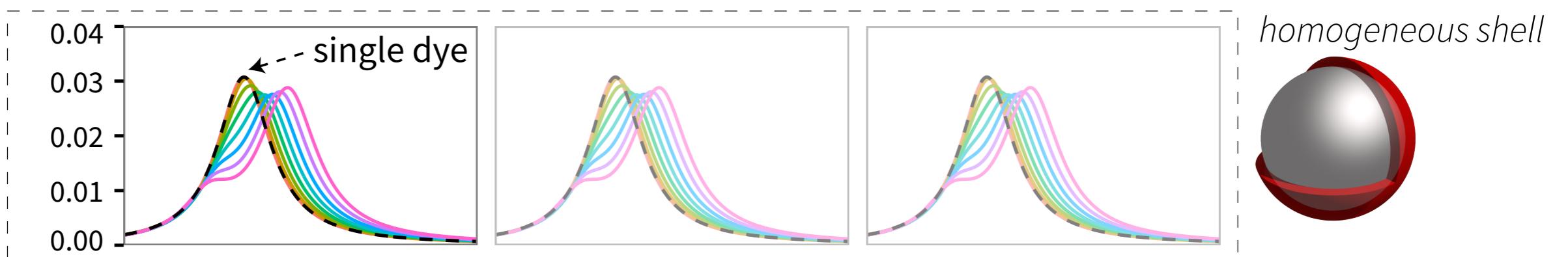
```
last_plot() +  
  facet_grid(Tree ~ .) +  
  theme_publication
```



FACETING – SMALL MULTIPLES REVISITED







$\rho_{\text{dye}} / \text{nm}^2$

—	0
—	0.2
—	0.4
—	0.6
—	0.8
—	1
—	1.2
—	1.4

even spacing

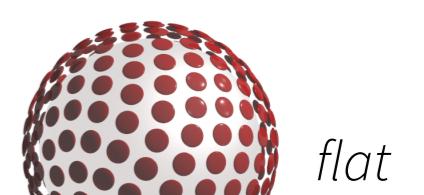
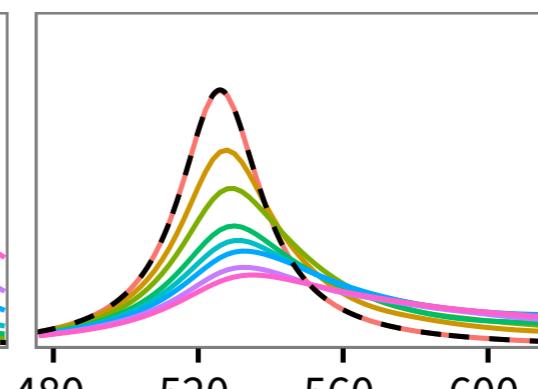
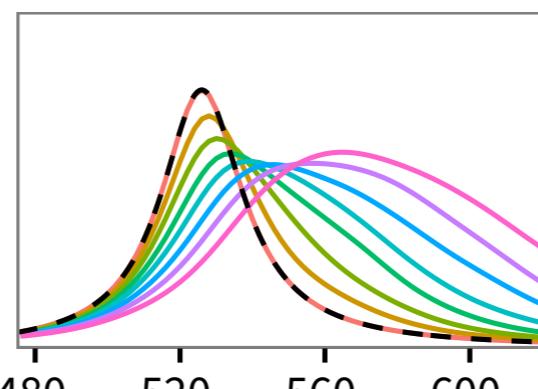
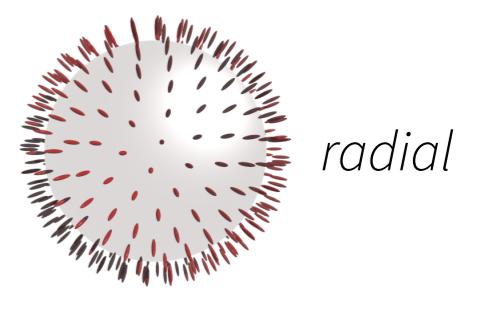
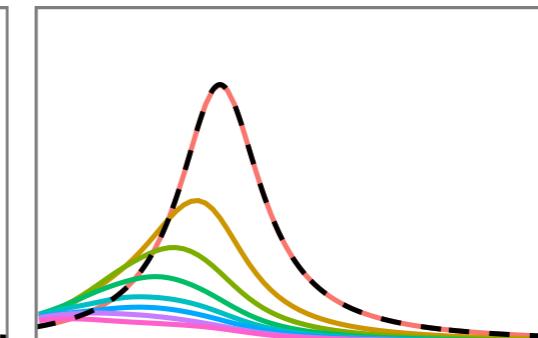
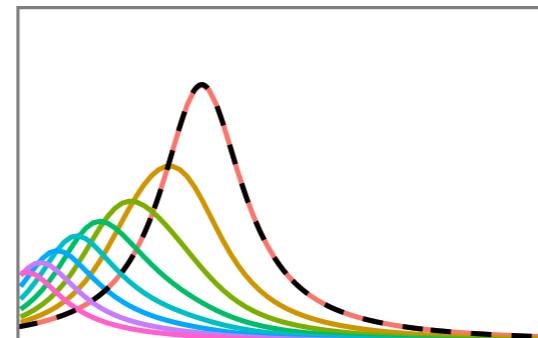
pseudo-random

random

isotropic

$\sigma_{\text{abs}} / \text{nm}^2$

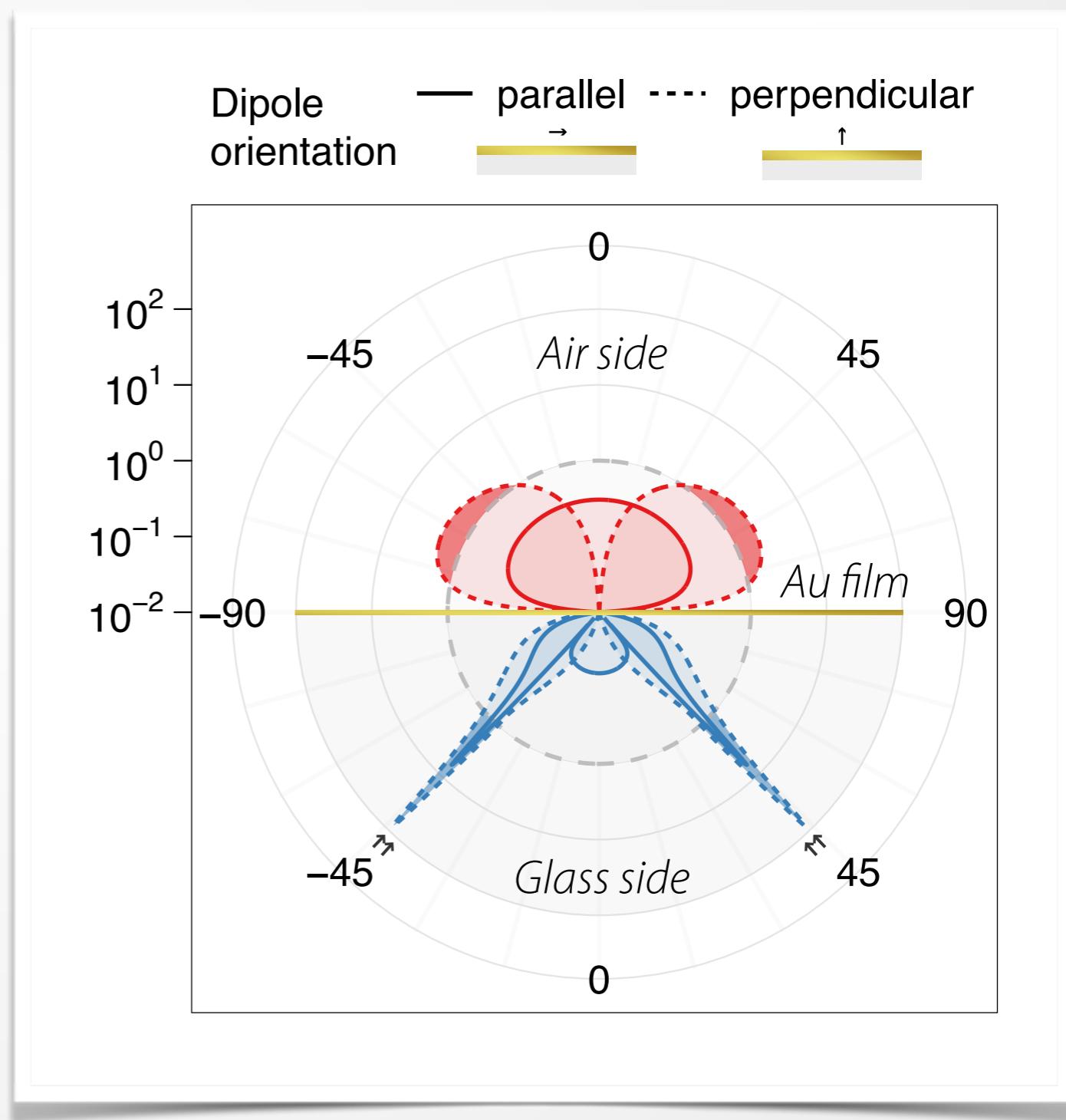
—	0
—	0.2
—	0.4
—	0.6
—	0.8
—	1
—	1.2
—	1.4



Wavelength /nm

COORDINATE TRANSFORMATIONS

```
plot(data, map(x, y)) +  
  layer(line, map(linetype = t)) +  
  coord_polar(theta = x) +  
  ...
```



But different representations of
the exact same data can lead to
different understanding and, more
importantly, to different decisions.

-R. Kosara

GRAMMAR OF GRAPHICS

Additional reading

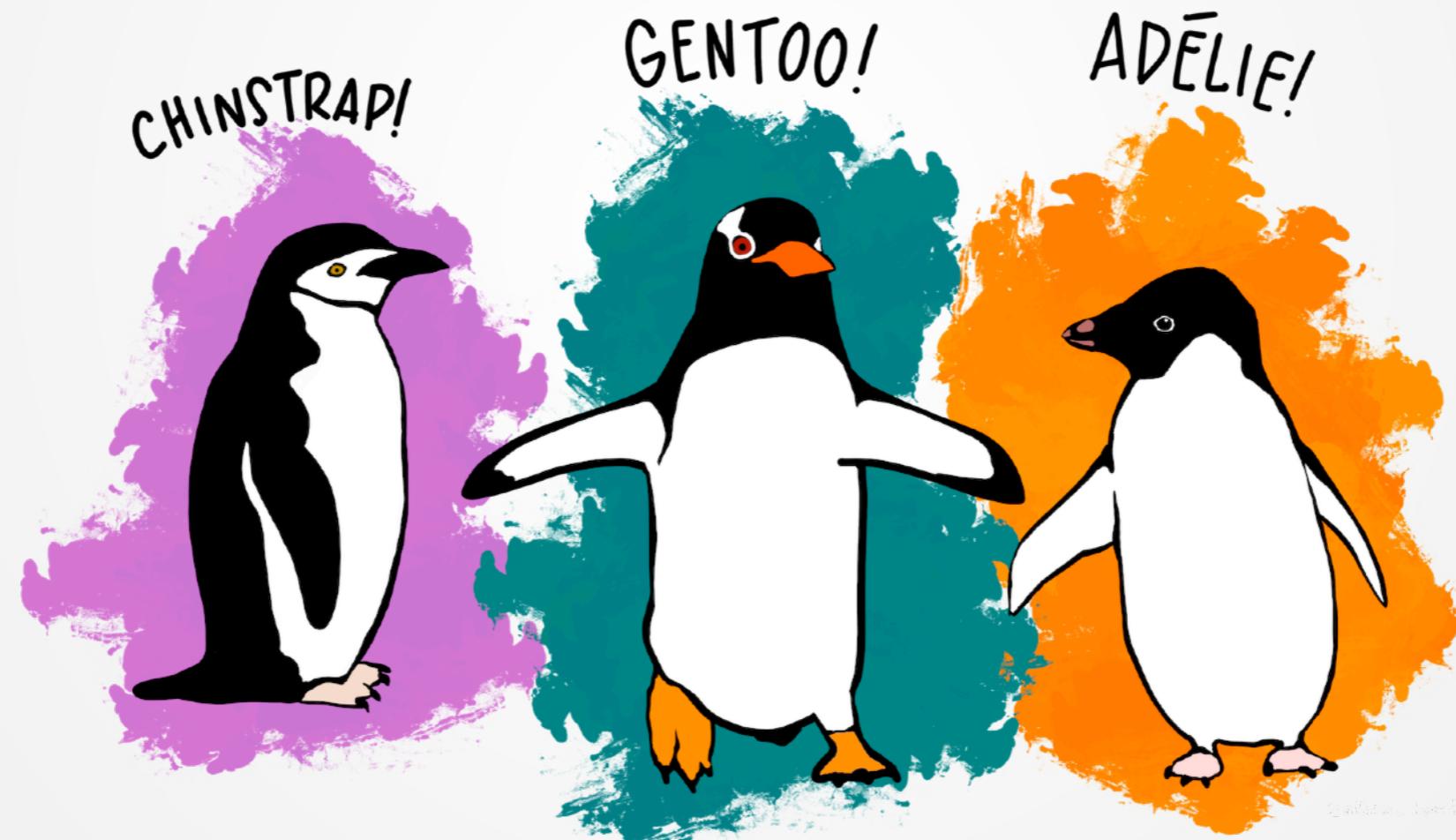
<https://vita.had.co.nz/papers/layered-grammar.html>

Python packages

- ▶ <https://plotnine.readthedocs.io>
- ▶ <https://altair-viz.github.io>
- ▶ <https://plotly.com/python>
- ▶ <https://seaborn.pydata.org>

PRACTICE WITH A TOY DATASET – PENGUINS

<https://github.com/mcnakhaee/palmerpenguins>



TOY DATASET — PENGUINS

