

# **Exploratory Data Analysis & Visualization**

## **Data Storytelling**

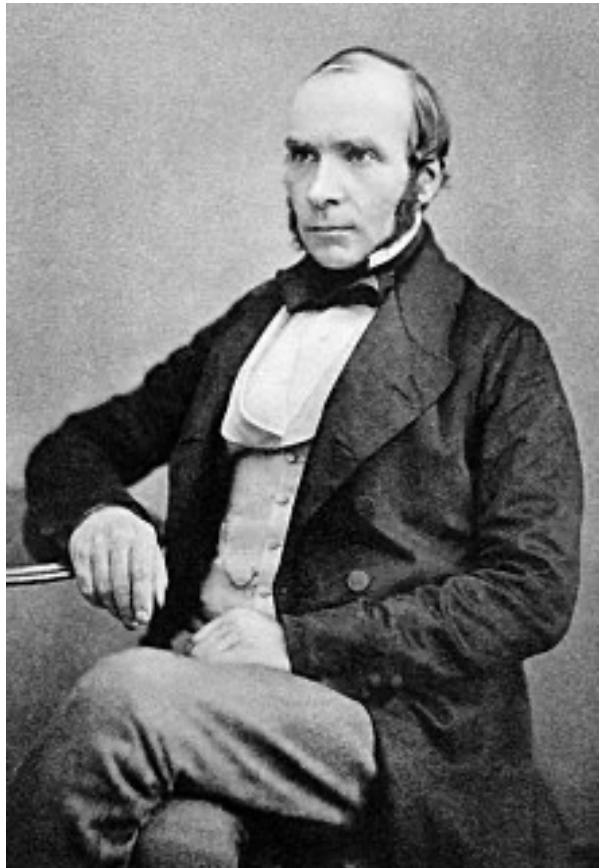
## **& Dashboards/Storyboards**

Ben Winjum

Any questions about the final project?

# Our familiar John Snow example

## John Snow



Dr. Snow was an intern physician when a cholera epidemic broke out in London in the 1830's.

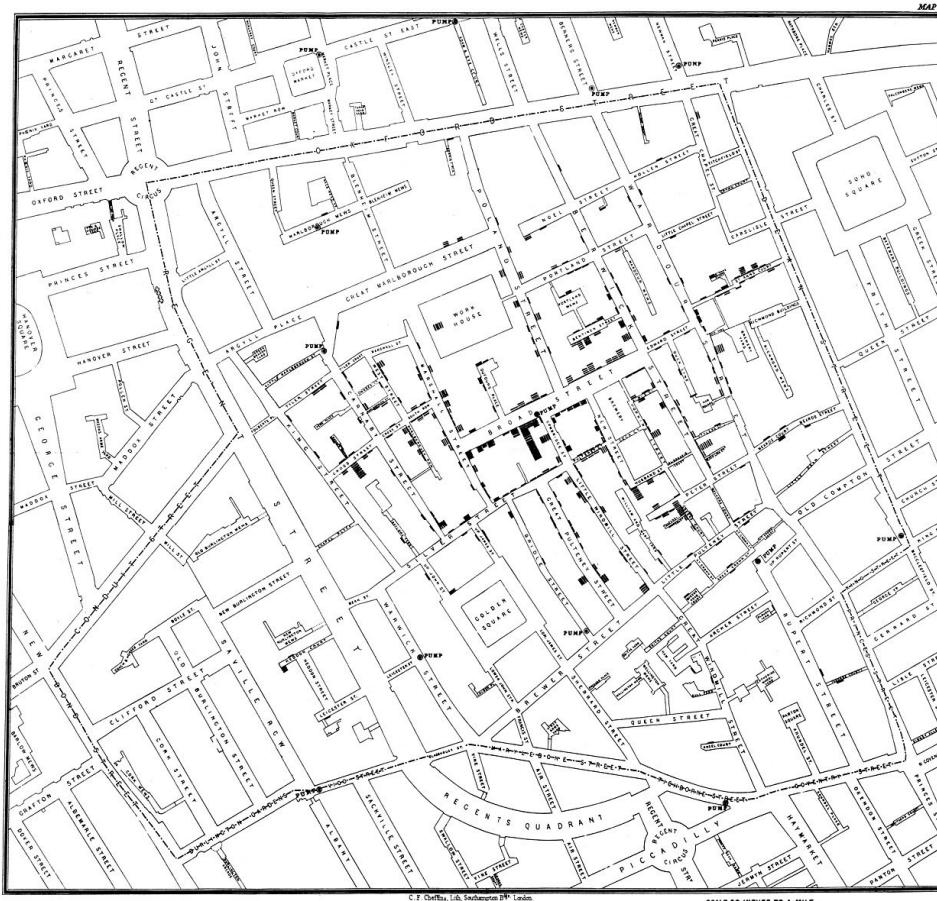
The germ theory of disease did not exist yet.

Prevailing opinion: cholera was somehow transmitted by bad air, or a "miasma," that emanated from the sick and from garbage.

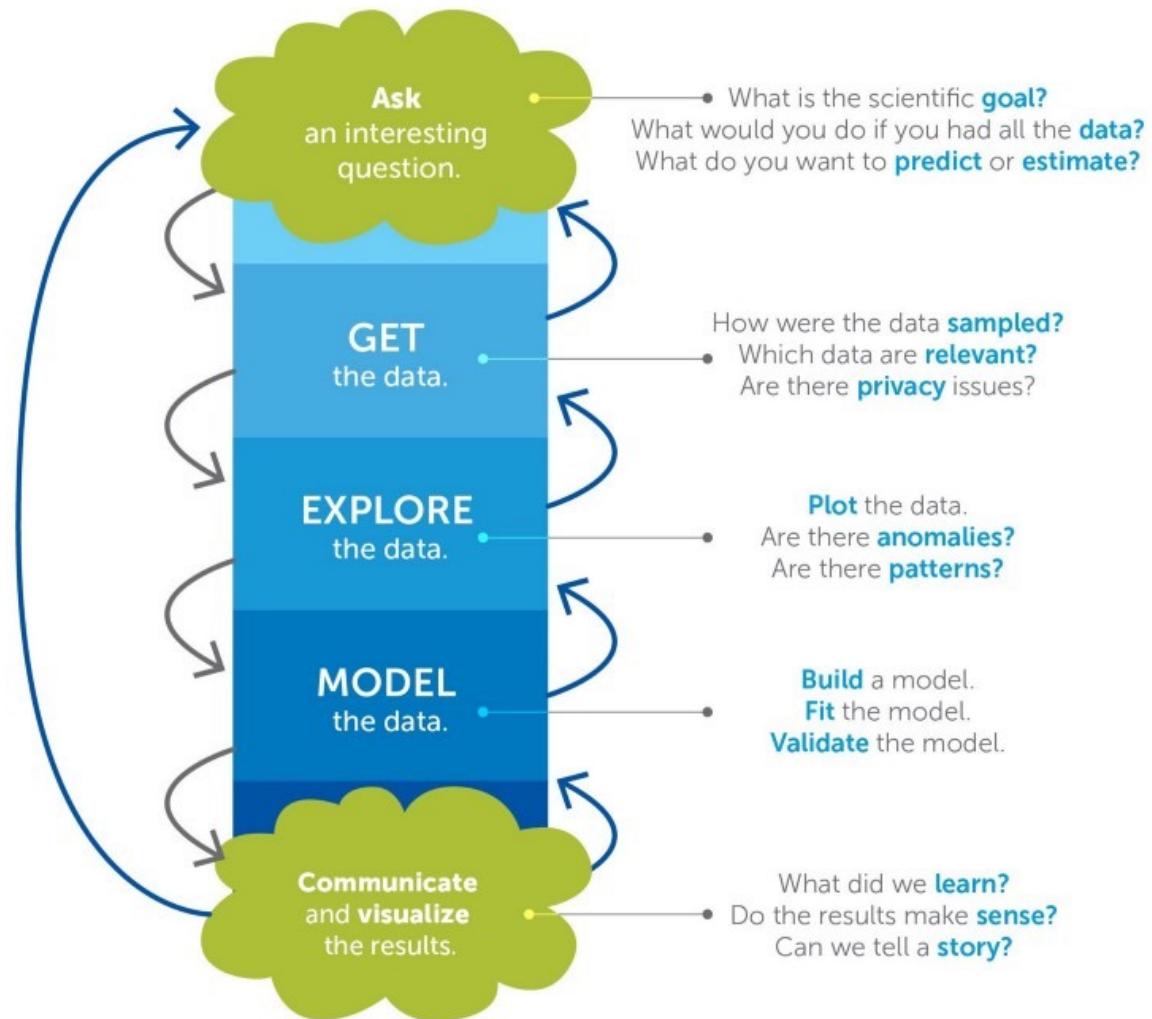
Physicians of the time gave treatments like:

- Arsenic and strychnine
- Tobacco enemas
- Wrap them in flannel soaked in turpentine
- Bleed them with leeches
- Blister them with nitric acid
- Fire cannons every hour to disperse the bad air!

# Using data visualization to draw conclusions

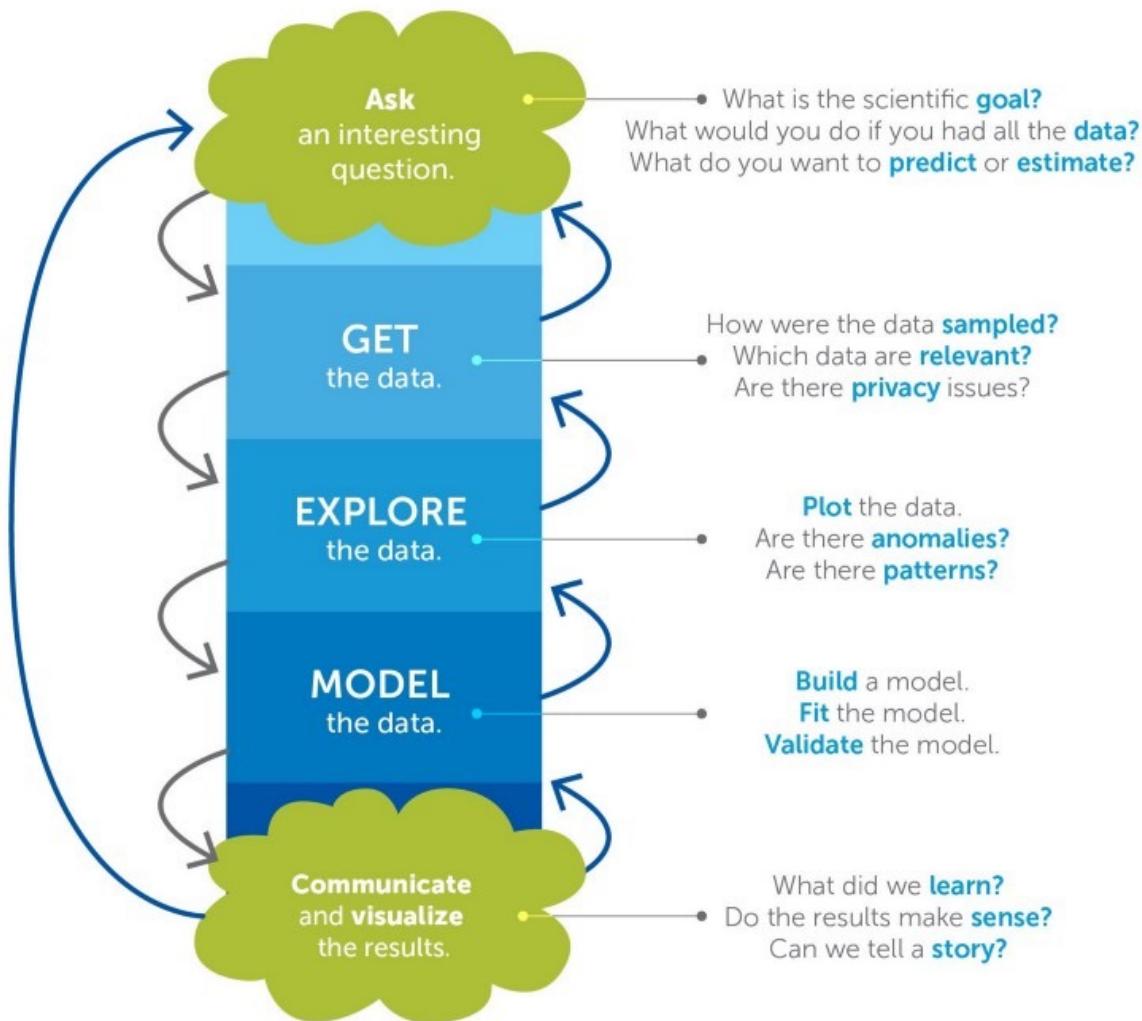


# The Data Science Process



Derived from the work of Joe Blitzstein and Hanspeter Pfister,  
originally created for the Harvard data science course <http://cs109.org/>.

# The Data Science Process



What's causing people to get cholera?

Locate cholera-stricken people

Visualize the map

What other locatable objects are relevant?

If  
(dist\_to\_broad\_st\_pump) less than (dist\_to\_other\_pumps)  
then  
(cholera\_likelihood) is greater than (random\_chance)

Present results to city officials; Enact change



Derived from the work of Joe Blitzstein and Hanspeter Pfister,  
originally created for the Harvard data science course <http://cs109.org/>.

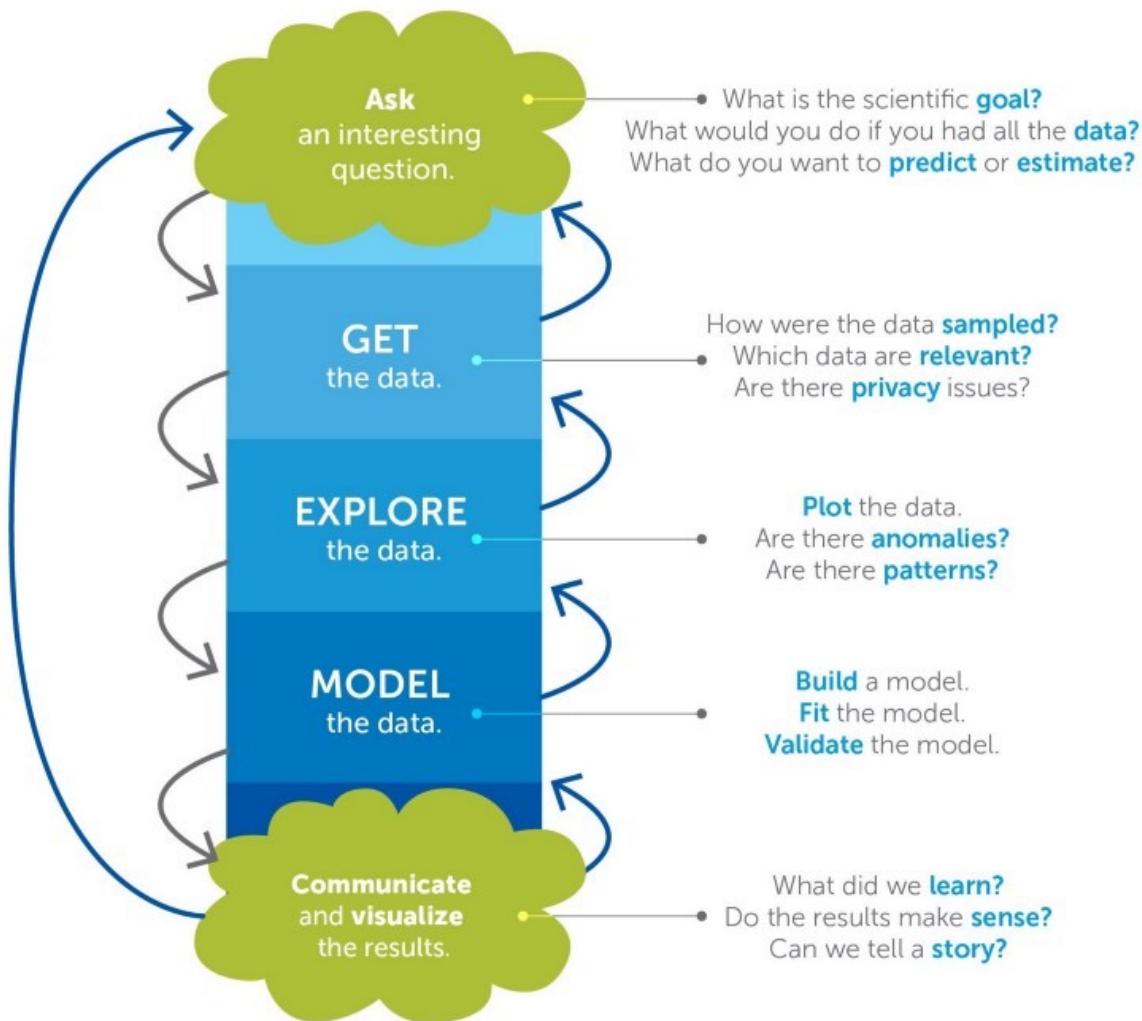
# Exploratory Data Analysis

- Get the data into a usable form / data structure
- Figure out what it contains
- Look at dimensions and values
- Clean it
  - Organize, structure, hunt for messed up values and NaNs
- Make tables
- Make plots
- Use visualization to inform our understanding of data

# Exploratory Data Analysis

- Get the data into a usable form / data structure
- Figure out what it contains
- Look at dimensions and values
- Clean it
  - Organize, structure, hunt for messed up values and NaNs
- Make tables
- Make plots
- Use visualization to inform our understanding of the data
  - It is also ***essential*** to know how to use visualization to help inform others' understanding of the data

# The Data Science Process



## Data Visualization

- Exploratory
  - Explore the data
  - Determine hypotheses
  - Decide on potential models
  - Free form, like jazz improv
- Explanatory
  - Present the data
  - Provide support for hypotheses and models
  - Enable action
  - Cohesive story and narrative, like a classical symphony



Derived from the work of Joe Blitzstein and Hanspeter Pfister,  
originally created for the Harvard data science course <http://cs109.org/>.

# Explanatory Data Analysis

- Data visualizations have a specific point

# Explanatory Data Analysis

- Data visualizations have a specific point
- They are **designed** with that point in mind

# Explanatory Data Analysis

- Data visualizations have a specific point
- They are designed with that point in mind
  - Who are you communicating to?

# Explanatory Data Analysis

- Data visualizations have a specific point
- They are designed with that point in mind
  - Who are you communicating to?
  - **What** do you want them to know or to do?

# Explanatory Data Analysis

- Data visualizations have a specific point
- They are designed with that point in mind
  - Who are you communicating to?
  - What do you want them to know or to do?
  - **How** can you use your data to achieve that?

# Who is your audience?

- Who are you targeting?
- What do they know?
- What do they want?
- The more specific, the better

# Who is your audience?

- Who are you targeting?
- What do they know?
- What do they want?
- The more specific, the better



# Who are you?

- How does your audience see you?
- What is your relationship with your audience?
- Do you need to establish credibility?

# Who are you?

- How does your audience see you?
- What is your relationship with your audience?
- Do you need to establish credibility?

## Data Analyst



What my friends think I do



What my Mom thinks I do



What my boss thinks I do



What my customers think I do



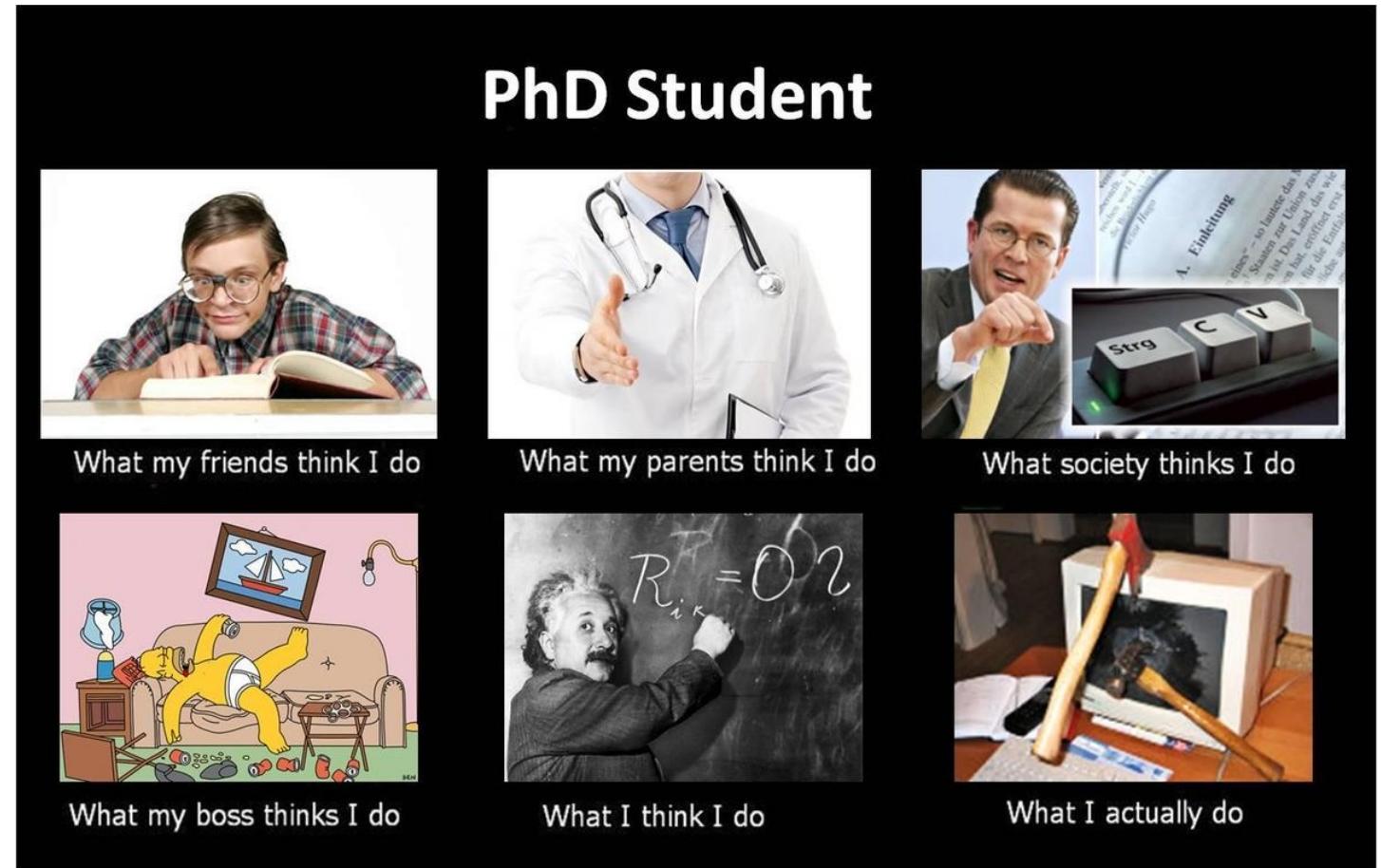
What I think I do



What I really do

# Who are you?

- How does your audience see you?
- What is your relationship with your audience?
- Do you need to establish credibility?



# What do you want your audience to know or to do?

- Always aim to either:
  - help them to know something
  - get them to do something
- If you aren't moving the audience towards a specific objective, you can still move them towards discussion or suggest next steps
- You know your data the best

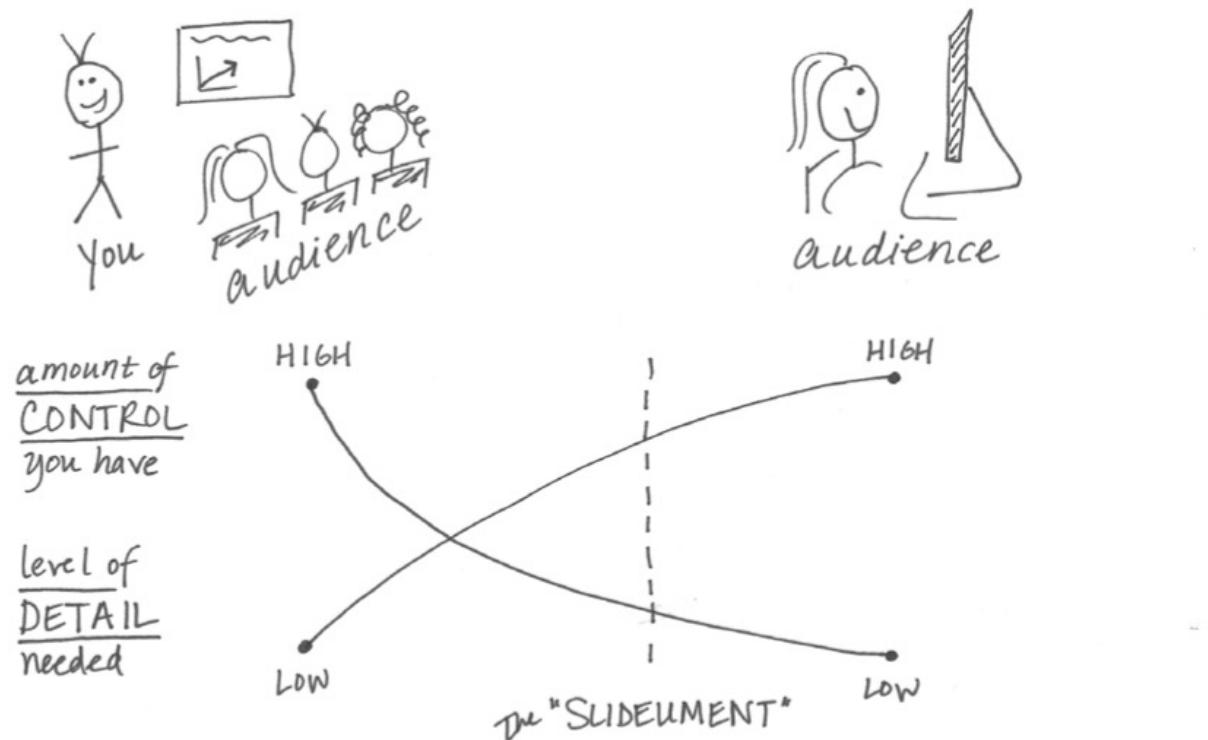
## Prompting action

Here are some action words to help act as thought starters as you determine what you are asking of your audience:

accept | agree | begin | believe | change | collaborate | commence | create | defend | desire | differentiate | do | empathize | empower | encourage | engage | establish | examine | facilitate | familiarize | form | implement | include | influence | invest | invigorate | know | learn | like | persuade | plan | promote | pursue | recommend | receive | remember | report | respond | secure | support | simplify | start | try | understand | validate

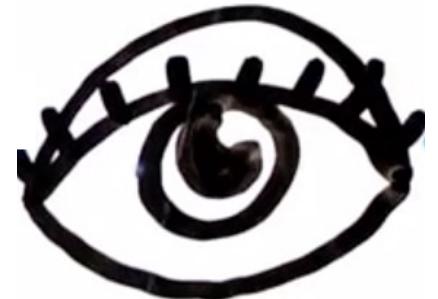
# How will you communicate?

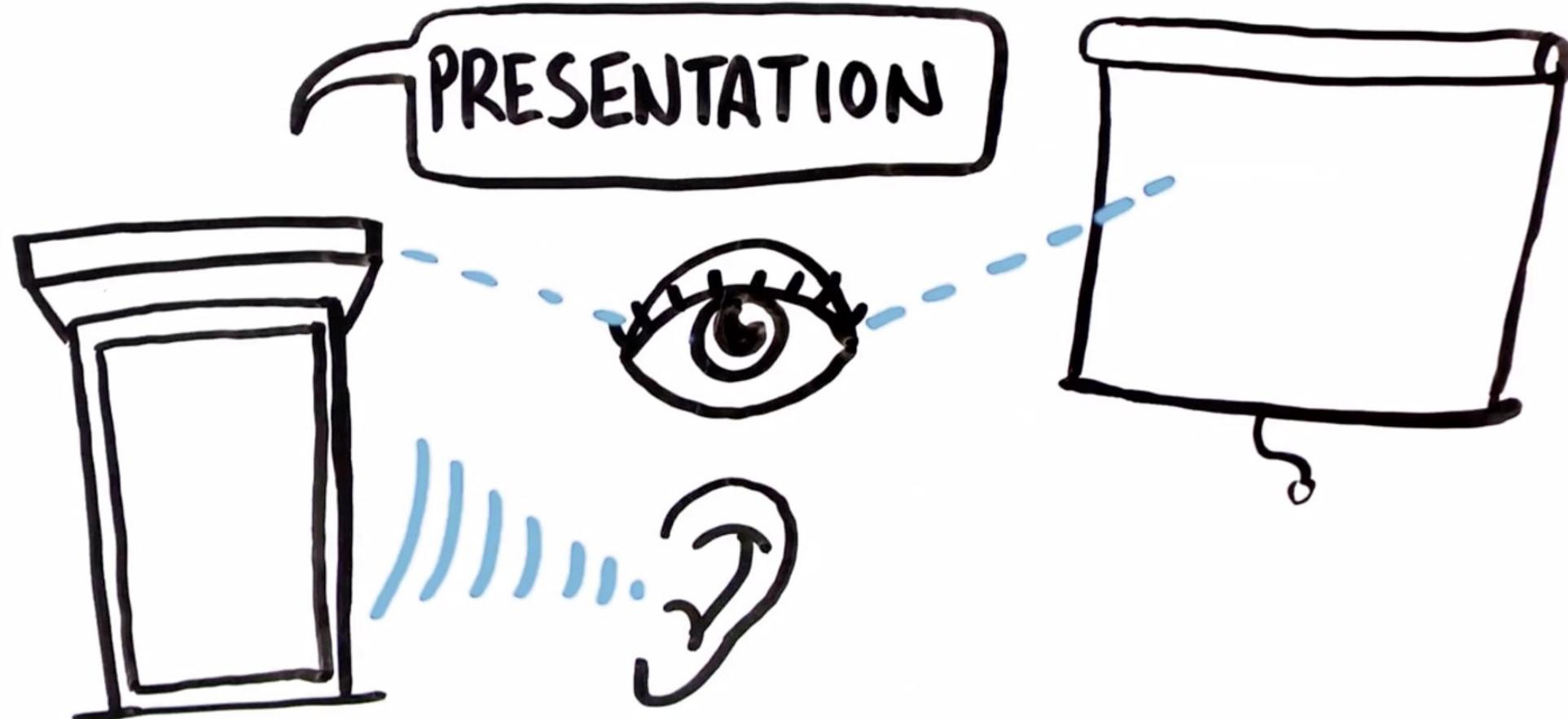
LIVE PRESENTATION ..... WRITTEN DOC OR EMAIL

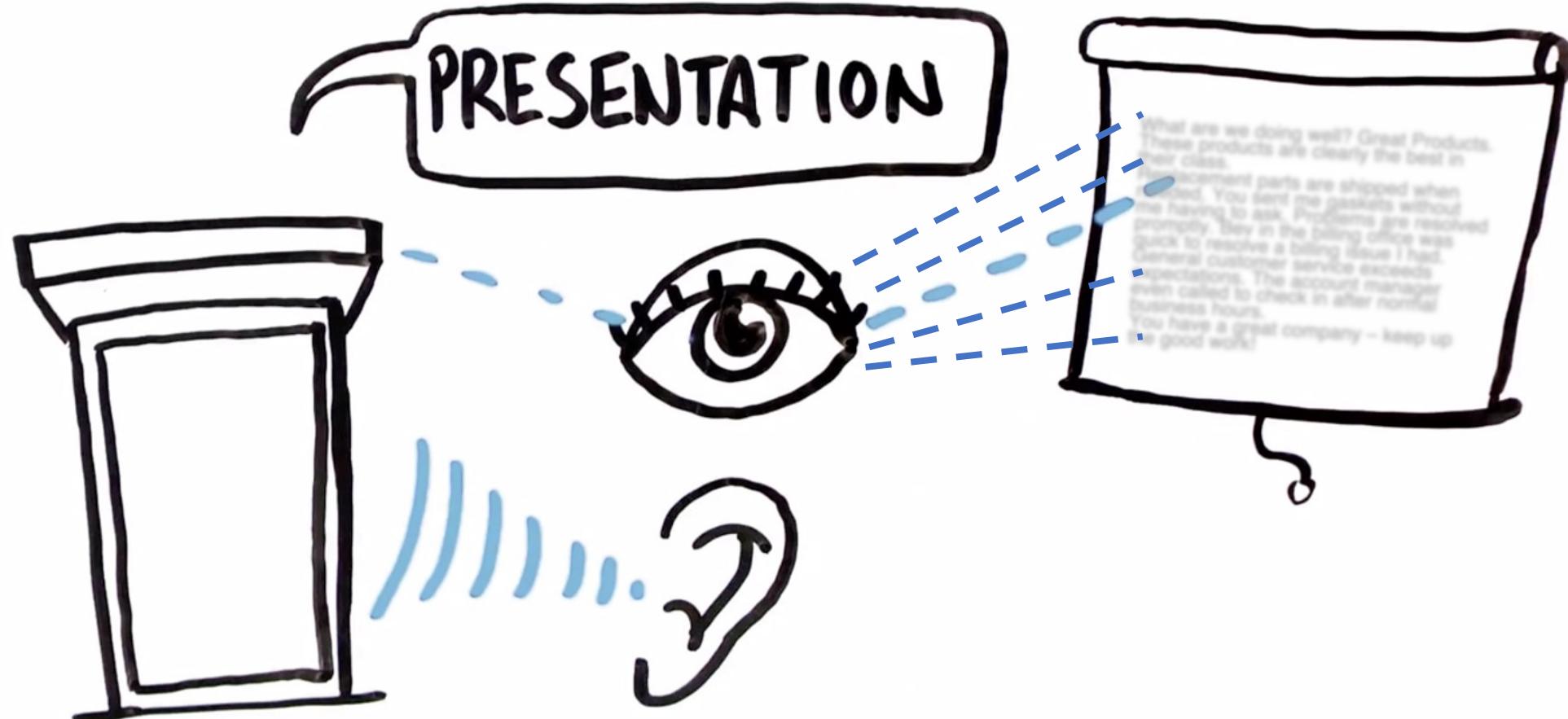


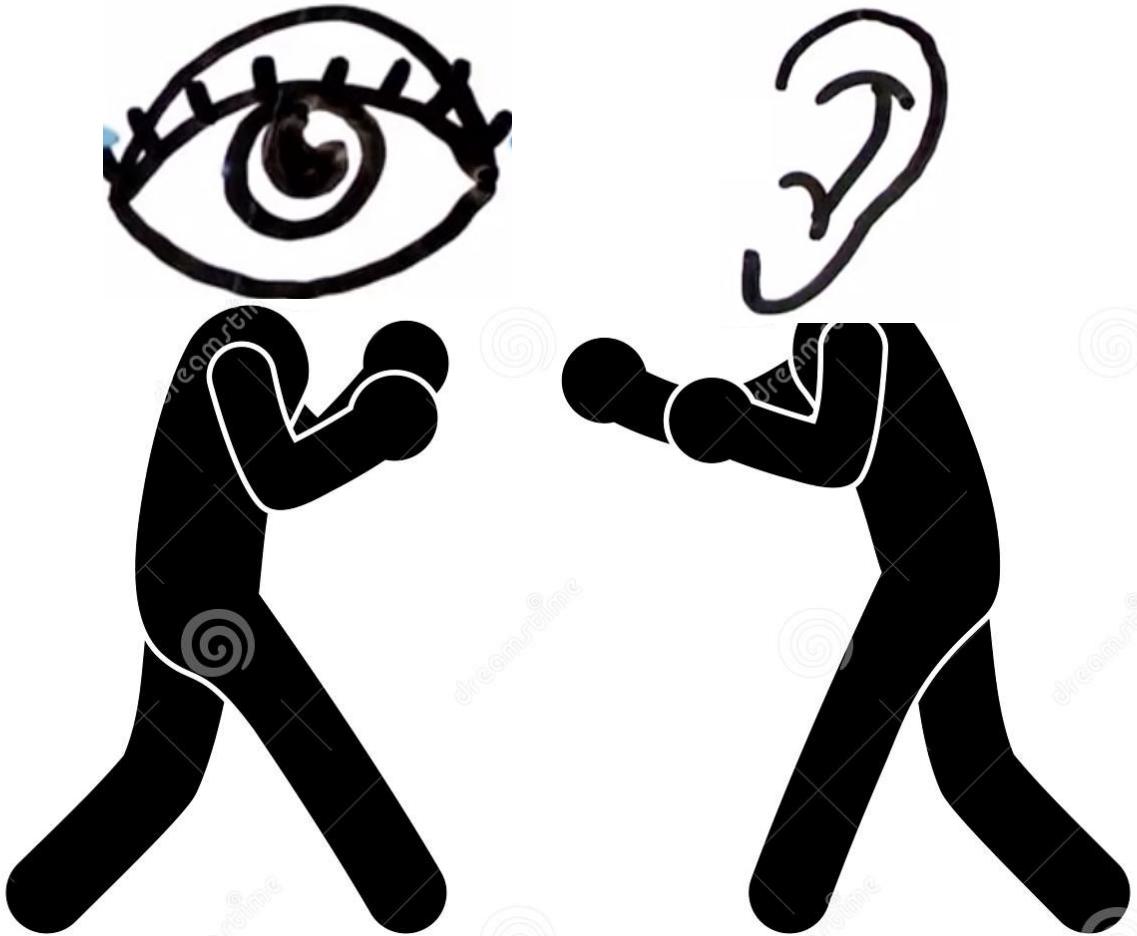
All master presenters should know:

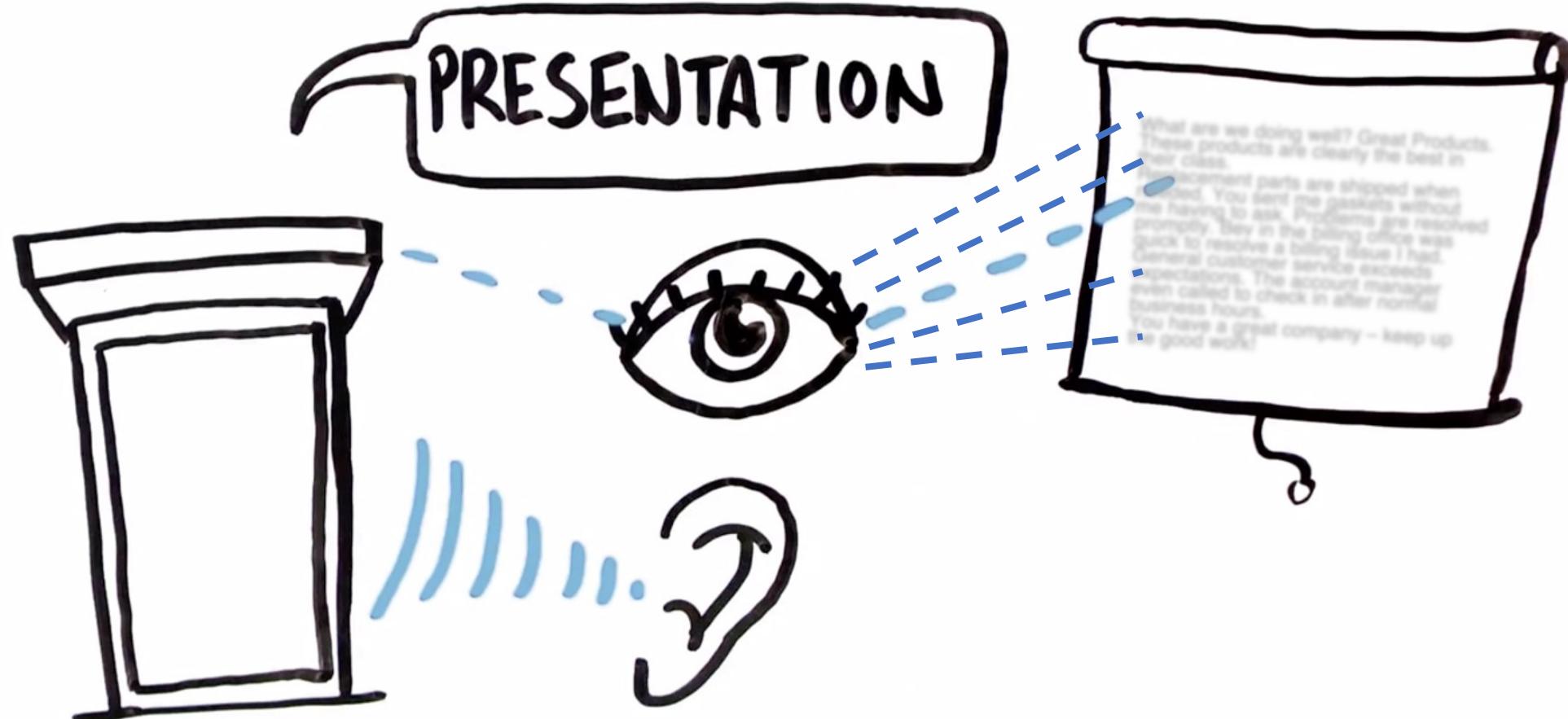
Multiple  
Sensory  
Channels  
Compete



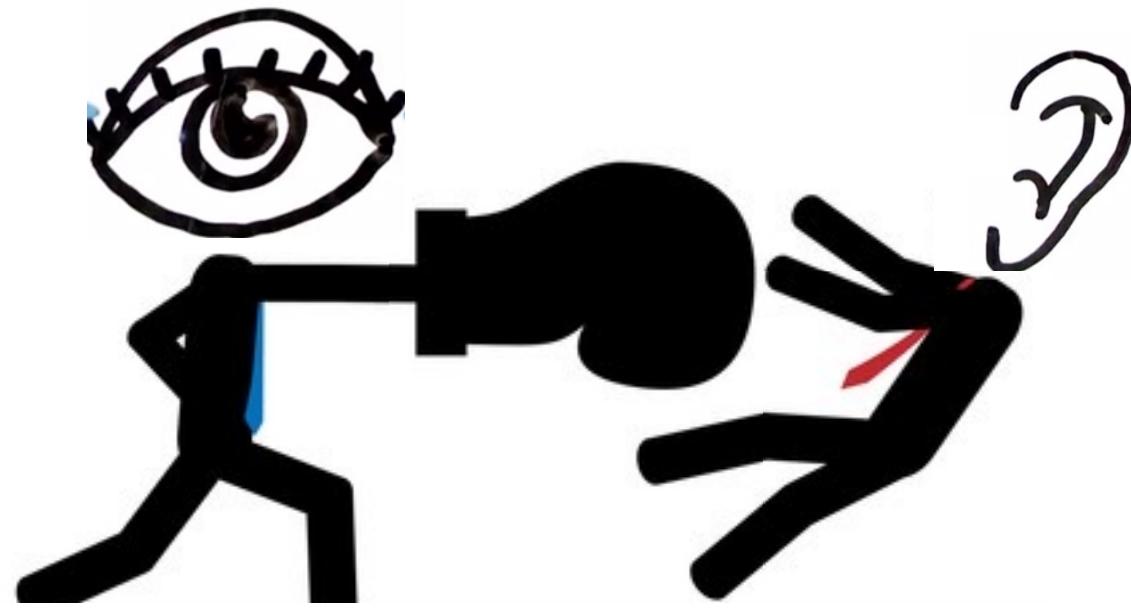


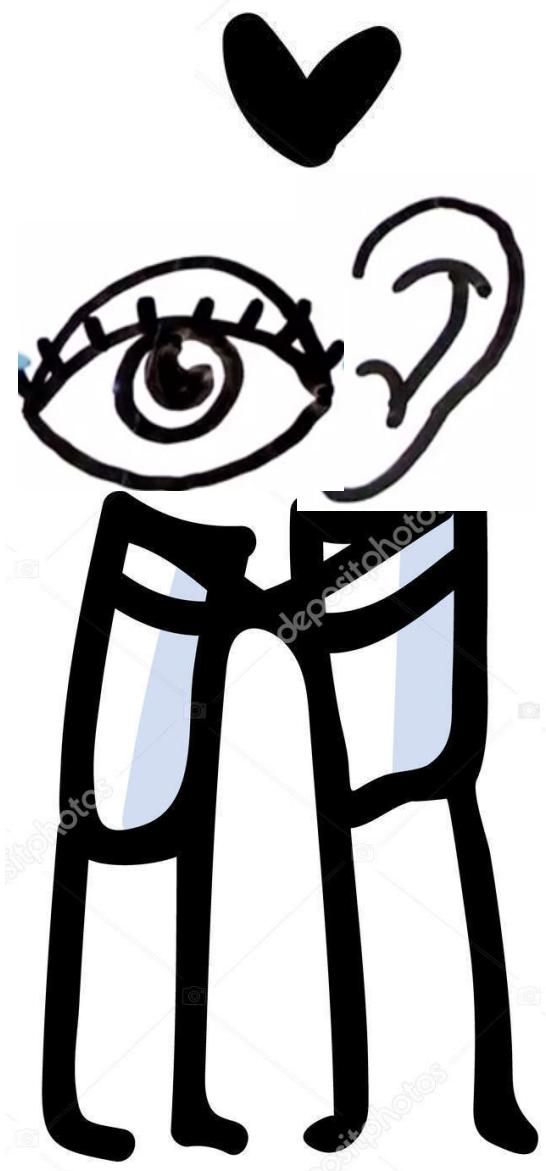






What are we doing well? Great Products. These products are clearly the best in their class. Replacement parts are shipped when needed. You sent me gaskets without me having to ask. Problems are resolved promptly. Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!





# Focusing attention



## Text – no preattentive attributes

What are we doing well? Great Products. These products are clearly the best in their class. Replacement parts are shipped when needed. You sent me gaskets without me having to ask. Problems are resolved promptly. Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!

**What are we doing well?** Great Products. These products are clearly the best in their class. Replacement parts are shipped when needed. You sent me gaskets without me having to ask. Problems are resolved promptly. Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!

What are we doing well? Great Products.  
**These products are clearly the best in their class.** Replacement parts are shipped when needed. You sent me gaskets without me having to ask. Problems are resolved promptly. Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!

What are we doing well? Great Products. These products are clearly the best in their class. *Replacement parts are shipped when needed.* You sent me gaskets without me having to ask. Problems are resolved promptly. Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!

What are we doing well? Great Products. These products are clearly the best in their class. Replacement parts are shipped when needed.

You sent me gaskets **without me having to ask**. Problems are resolved promptly. Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!

What are we doing well? Great Products. These products are clearly the best in their class. Replacement parts are shipped when needed. You sent me gaskets without me having to ask.

Problems are resolved promptly.

Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!

What are we doing well? Great Products. These products are clearly the best in their class. Replacement parts are shipped when needed. You sent me gaskets without me having to ask. Problems are resolved promptly. Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!

What are we doing well? Great Products. These products are clearly the best in their class. Replacement parts are shipped when needed. You sent me gaskets without me having to ask. Problems are resolved promptly. Bev in the billing office was quick to resolve a billing issue I had. General customer service exceeds expectations. The account manager even called to check in after normal business hours. You have a great company – keep up the good work!



Orientation



Shape



Line length



Line width



Size



Curvature



Added marks



Enclosure



Hue



Intensity



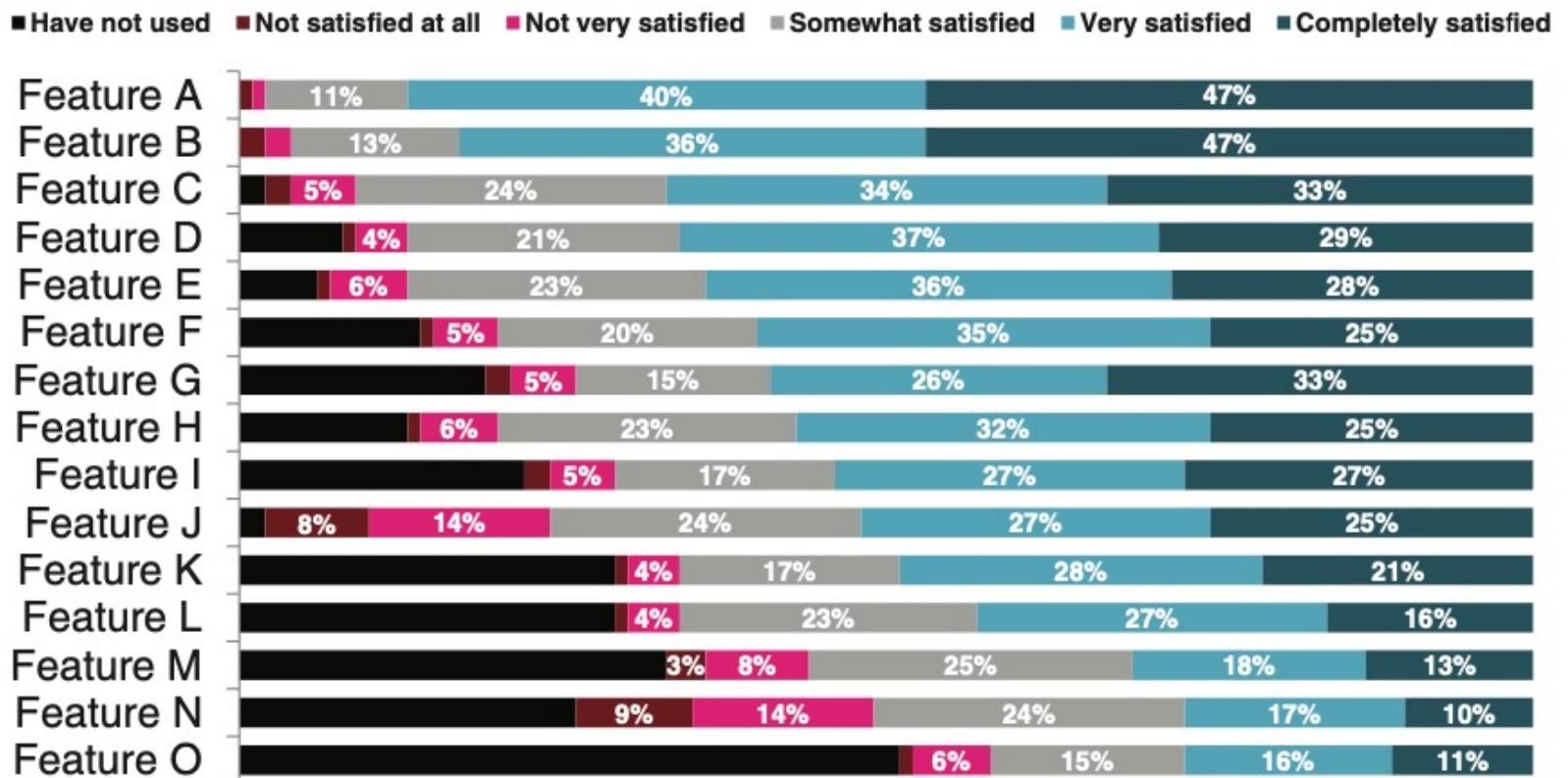
Spatial position



Motion

Focus attention on the attributes  
that best suit your purpose

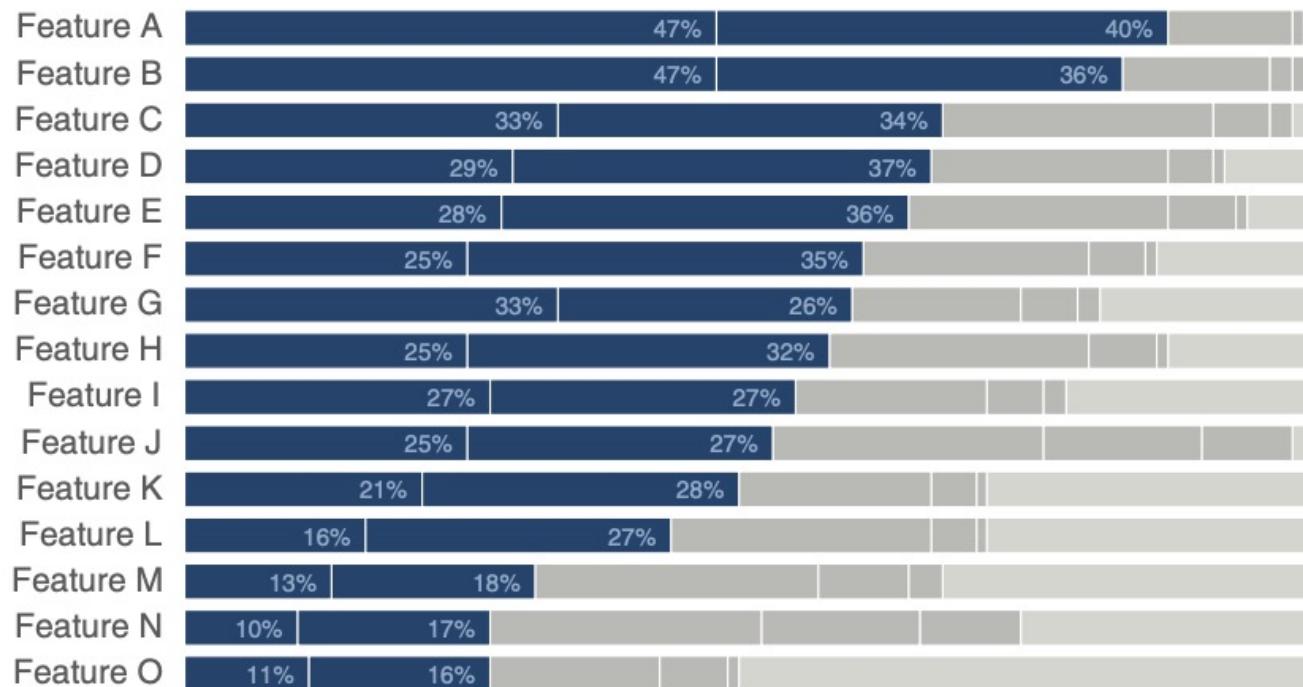
## How satisfied have you been with each of these features?



## Features A & B top user satisfaction

### Product X User Satisfaction: Features

■ Completely satisfied ■ Very satisfied ■ Somewhat satisfied ■ Not very satisfied ■ Not satisfied at all ■ Have not used



Responses based on survey question "How satisfied have you been with each of these features?".

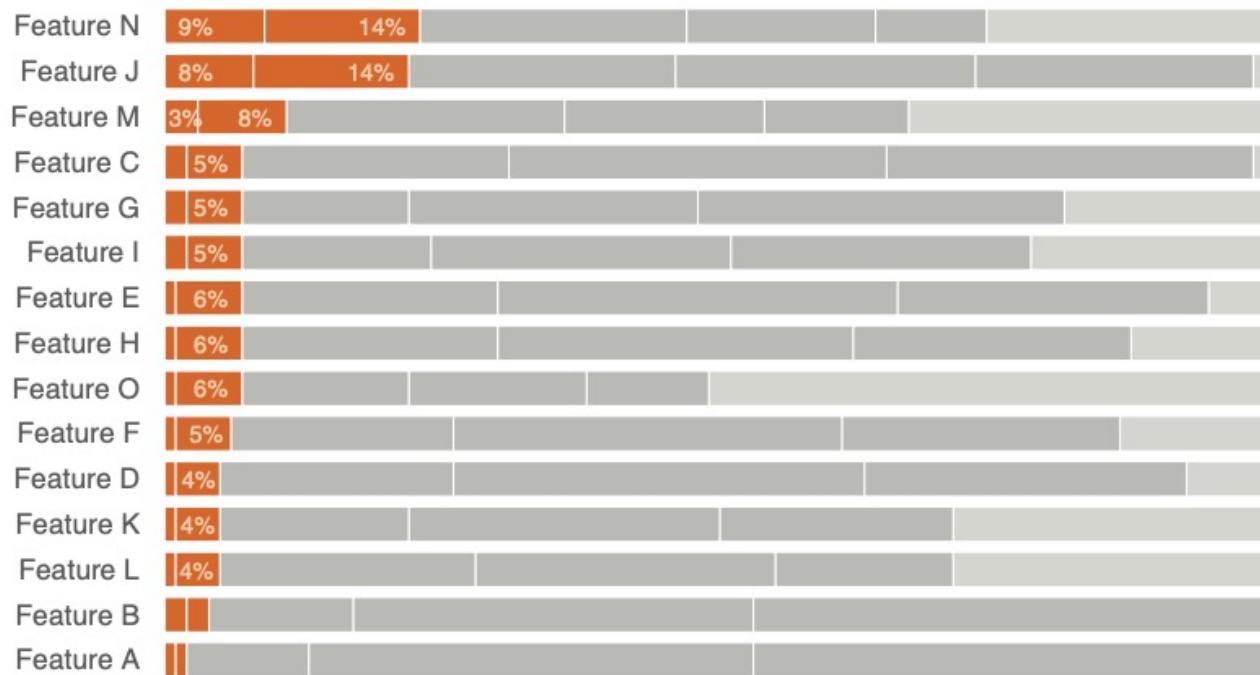
Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent?

Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

## Users least satisfied with Features N & J

### Product X User Satisfaction: Features

■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ Very satisfied ■ Completely satisfied ■ Have not used



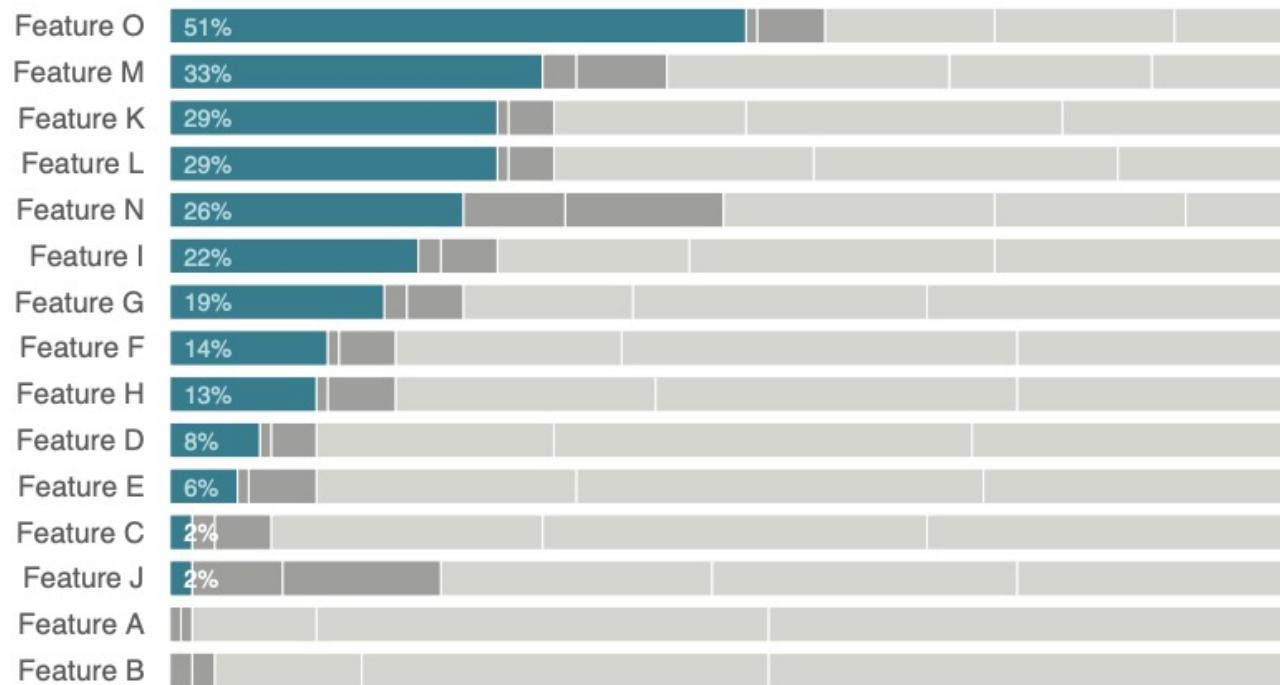
Responses based on survey question "How satisfied have you been with each of these features?".

Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent? Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

## Feature O is least used

### Product X User Satisfaction: Features

■ Have not used ■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ Very satisfied ■ Completely satisfied



Responses based on survey question "How satisfied have you been with each of these features?".

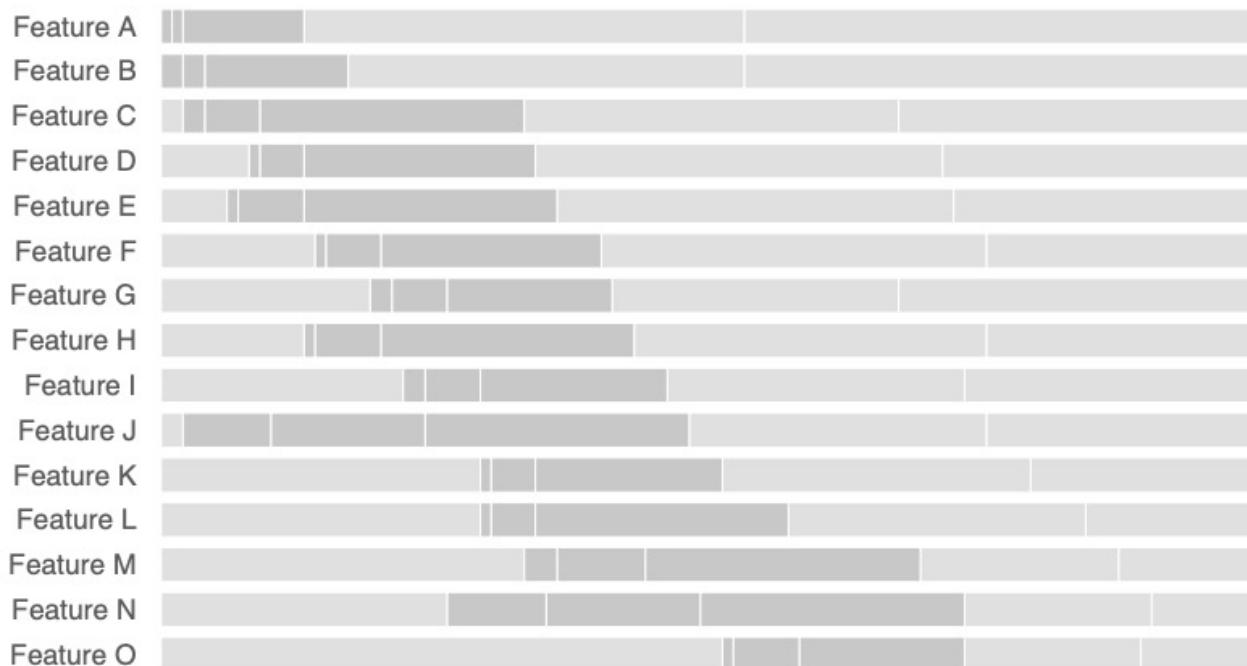
Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent?

Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

## User satisfaction varies greatly by feature

Product X User Satisfaction: **Features**

■ Have not used ■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ Very satisfied ■ Completely satisfied



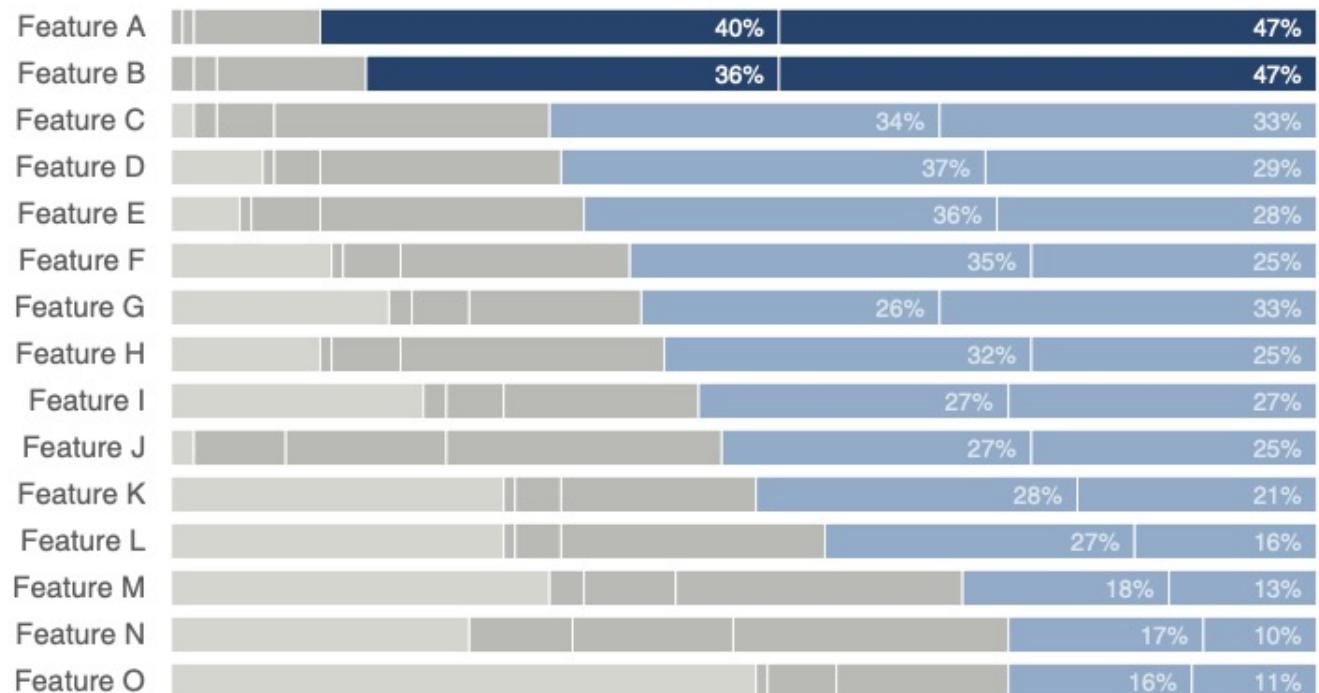
Responses based on survey question "How satisfied have you been with each of these features?".

Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent? Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

# User satisfaction varies greatly by feature

## Product X User Satisfaction: Features

■ Have not used ■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ **Very satisfied** ■ Completely satisfied



Features  
A and B  
continue  
to top user  
satisfaction

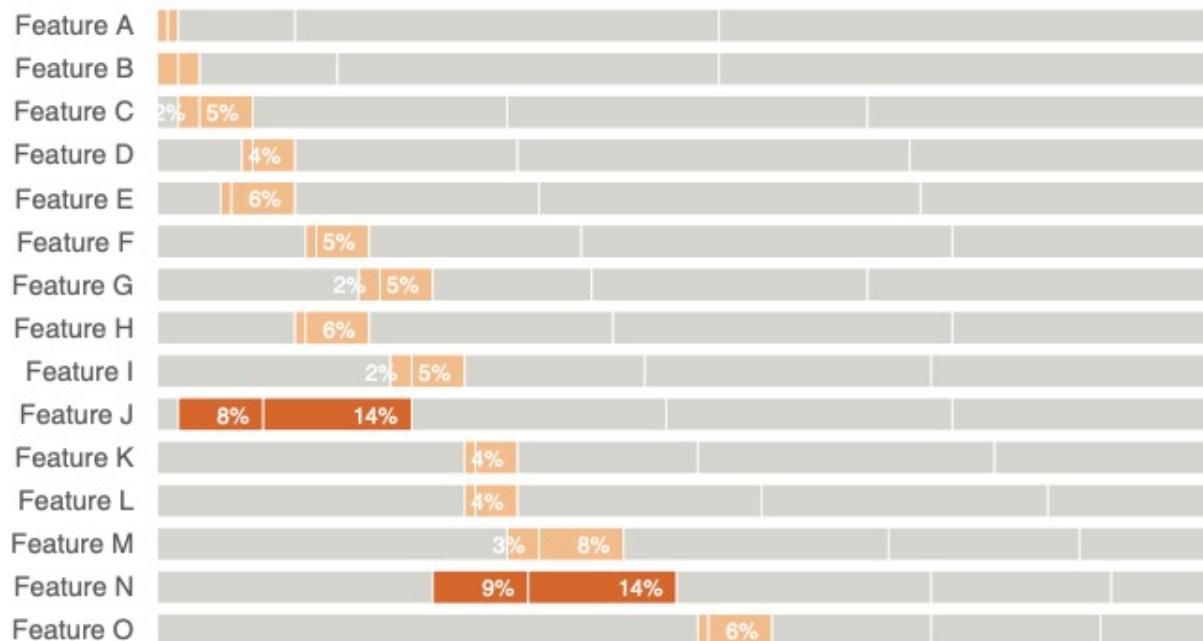
Responses based on survey question "How satisfied have you been with each of these features?".

Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent? Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

# User satisfaction varies greatly by feature

## Product X User Satisfaction: Features

■ Have not used ■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ Very satisfied ■ Completely satisfied



Users are least satisfied with Features J and N; what improvements can we make here for a better user experience?

Responses based on survey question "How satisfied have you been with each of these features?".

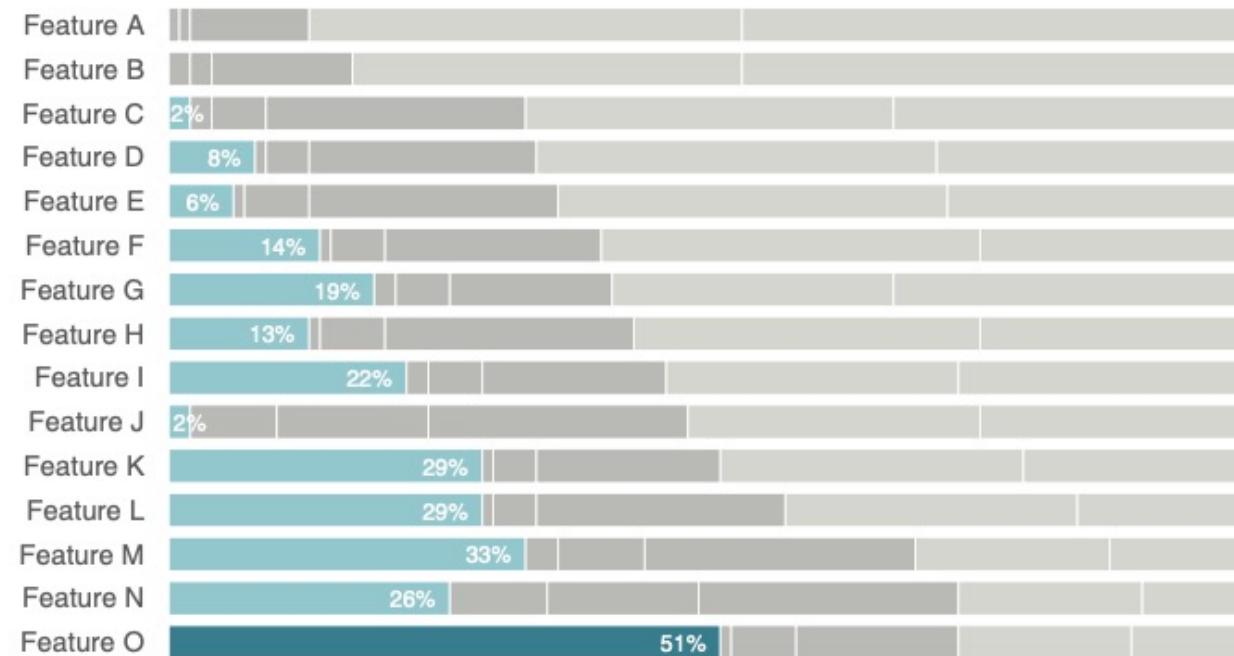
Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent?

Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

# User satisfaction varies greatly by feature

## Product X User Satisfaction: Features

■ Have not used ■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ Very satisfied ■ Completely satisfied



Feature O is least used. What steps can we proactively take with existing users to increase utilization?

Responses based on survey question "How satisfied have you been with each of these features?"

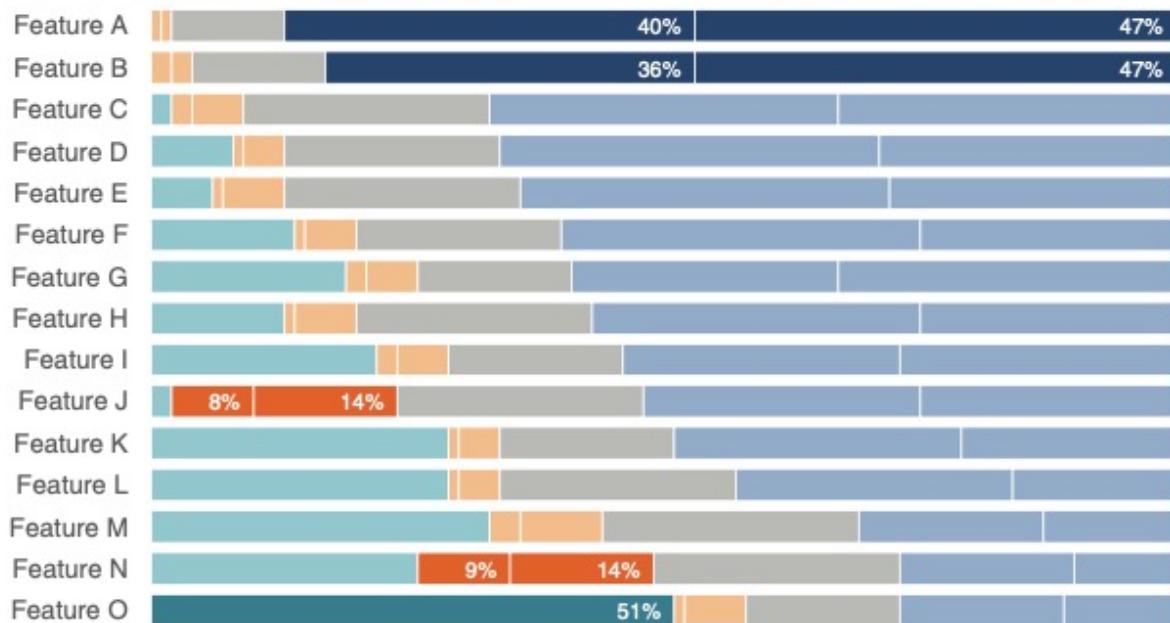
Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent?

Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

# User satisfaction varies greatly by feature

## Product X User Satisfaction: Features

■ Have not used ■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ Very satisfied ■ Completely satisfied



Features A and B continue to top user satisfaction

Users are least satisfied with Features J and N; what improvements can we make here for a better user experience?

Feature O is least used. What steps can we proactively take with existing users to increase utilization?

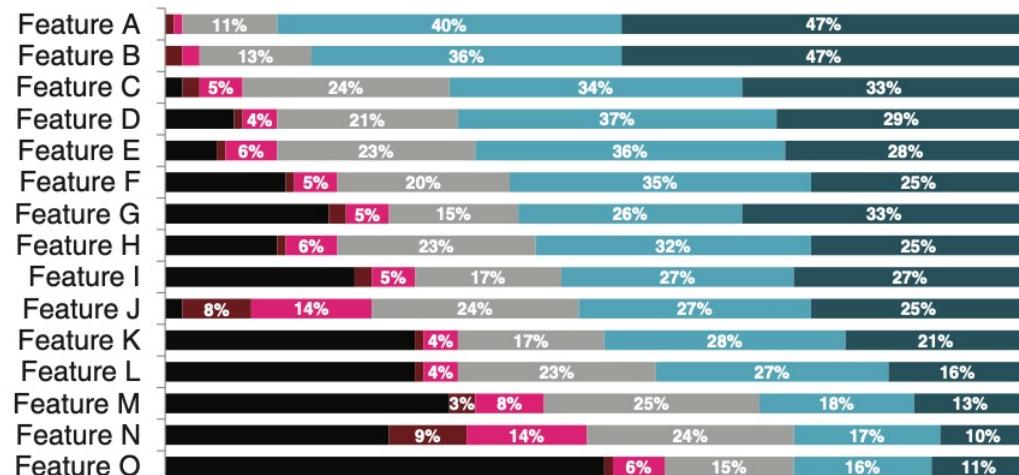
Responses based on survey question "How satisfied have you been with each of these features?".

Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent?

Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

### How satisfied have you been with each of these features?

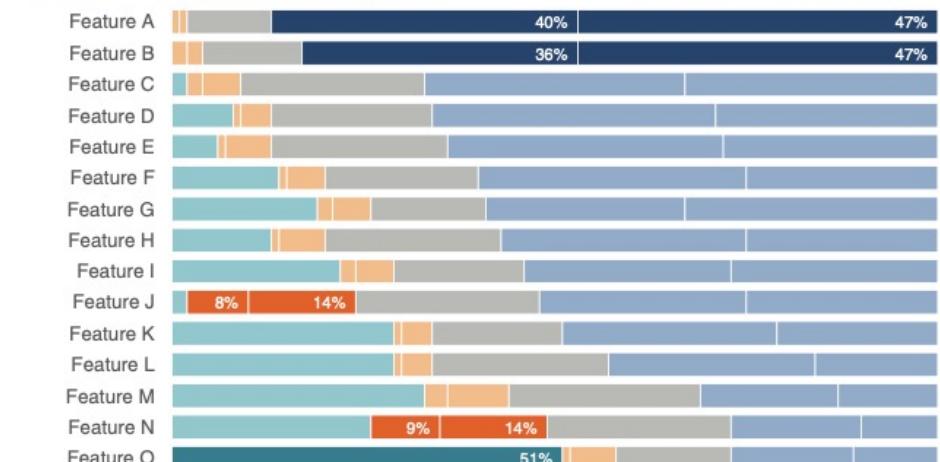
■ Have not used ■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ Very satisfied ■ Completely satisfied



## User satisfaction varies greatly by feature

### Product X User Satisfaction: Features

■ Have not used ■ Not satisfied at all ■ Not very satisfied ■ Somewhat satisfied ■ Very satisfied ■ Completely satisfied



Features A and B continue to top user satisfaction

Users are least satisfied with Features J and N; what improvements can we make here for a better user experience?

Feature O is least used. What steps can we proactively take with existing users to increase utilization?

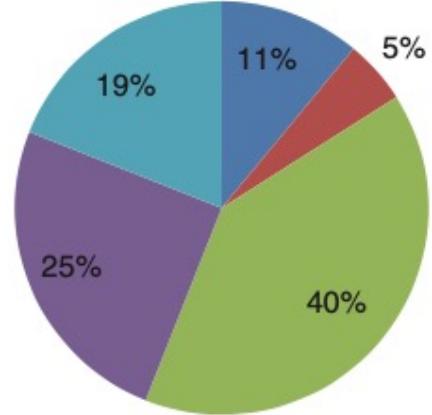
Responses based on survey question "How satisfied have you been with each of these features?".  
Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent?  
Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

Choose the type of plot  
to best suit your purpose

## Survey results: summer learning program on science

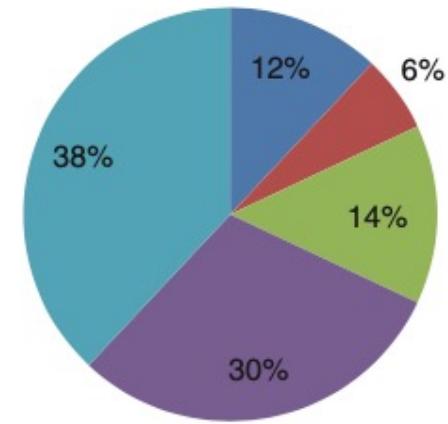
### PRE: How do you feel about doing science?

■ Bored ■ Not great ■ OK ■ Kind of interested ■ Excited



### POST: How do you feel about doing science?

■ Bored ■ Not great ■ OK ■ Kind of interested ■ Excited



## Pilot program was a success

After the pilot program,

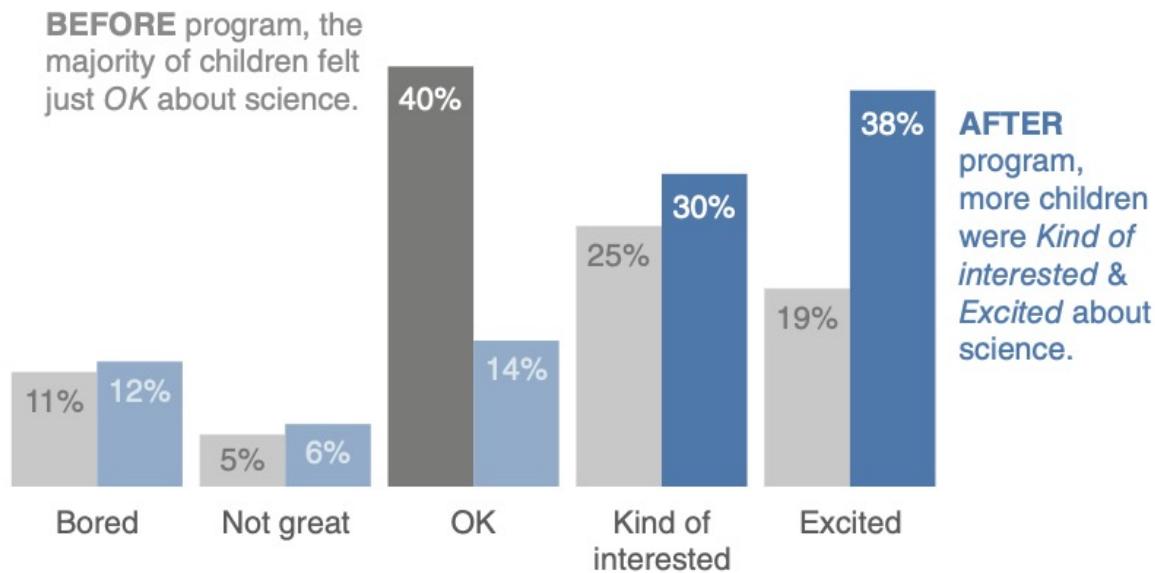
**68%**

**of kids expressed interest towards science,**  
compared to 44% going into the program.

Based on survey of 100 students conducted before and after pilot program (100% response rate on both surveys).

## Pilot program was a success

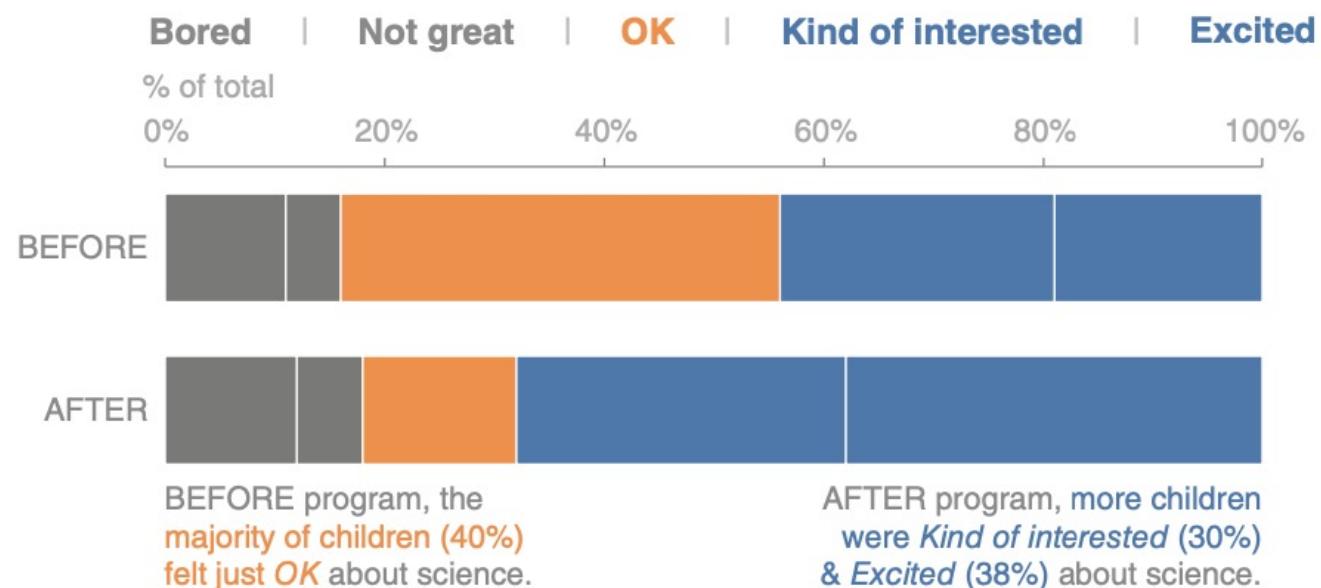
How do you feel about science?



Based on survey of 100 students conducted before and after pilot program (100% response rate on both surveys).

## Pilot program was a success

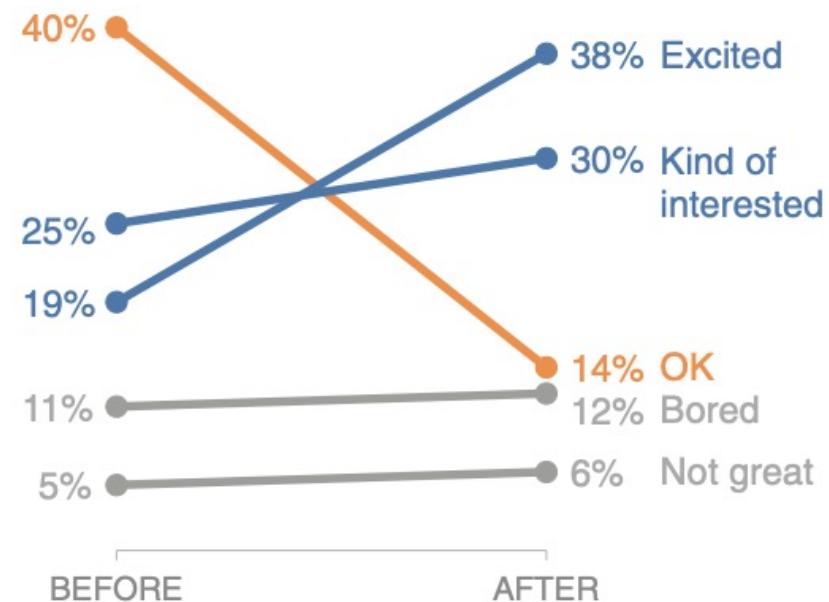
How do you feel about science?



Based on survey of 100 students conducted before and after pilot program (100% response rate on both surveys).

## Pilot program was a success

How do you feel about science?



BEFORE program, the majority of children felt just *OK* about science.

AFTER program, more children were *Kind of interested* & *Excited* about science.

Based on survey of 100 students conducted before and after pilot program (100% response rate on both surveys).

Choose the type of plot  
to best suit your purpose

... how do we identify that purpose?

# How do you know your communication's purpose?

- Use your who, what, and how
  - Who: The budget committee – they can approve funding for continuation of the summer learning program.
  - What: Approve budget of \$X to continue.
  - How: Illustrate that the summer learning program on science was a success with data collected through the survey conducted before and after the pilot program.

# How do you know your communication's purpose?

- Use your who, what, and how
  - Who: The budget committee – they can approve funding for continuation of the summer learning program.
  - What: Approve budget of \$X to continue.
  - How: Illustrate that the summer learning program on science was a success with data collected through the survey conducted before and after the pilot program.
- Sometimes your purpose is dictated to you
- Sometimes you'll dictate it to others

# Dictating your communication's purpose

- What background information is relevant or essential?
- Who is the audience or decision maker? What do we know about them?
- What biases does our audience have that might make them supportive of or resistant to our message?
- What data is available that would strengthen our case? Is our audience familiar with this data, or is it new?
- Where are the risks: what factors could weaken our case and do we need to proactively address them?
- What would a successful outcome look like?
- If you only had a limited amount of time or a single sentence to tell your audience what they need to know, what would you say?

# Your communication's purpose

- If you only had a limited amount of time or a single sentence to tell your audience what they need to know, what would you say?
  - This can be particularly effective

# Your communication's purpose

- If you only had a limited amount of time or a single sentence to tell your audience what they need to know, what would you say?
  - This can be particularly effective
- Nancy Duarte [*Resonate* (2010)] recommends identifying the Big Idea
  - It must articulate your unique point of view;
  - It must convey what's at stake; and
  - It must be a complete sentence.

# Your communication's purpose

- If you only had a limited amount of time or a single sentence to tell your audience what they need to know, what would you say?
  - This can be particularly effective
- Nancy Duarte [*Resonate* (2010)] recommends identifying the Big Idea
  - It must articulate your unique point of view;
  - It must convey what's at stake; and
  - It must be a complete sentence.
- *The pilot summer learning program was successful at improving students' perceptions of science and, because of this success, we recommend continuing to offer it going forward; please approve our budget for this program.*

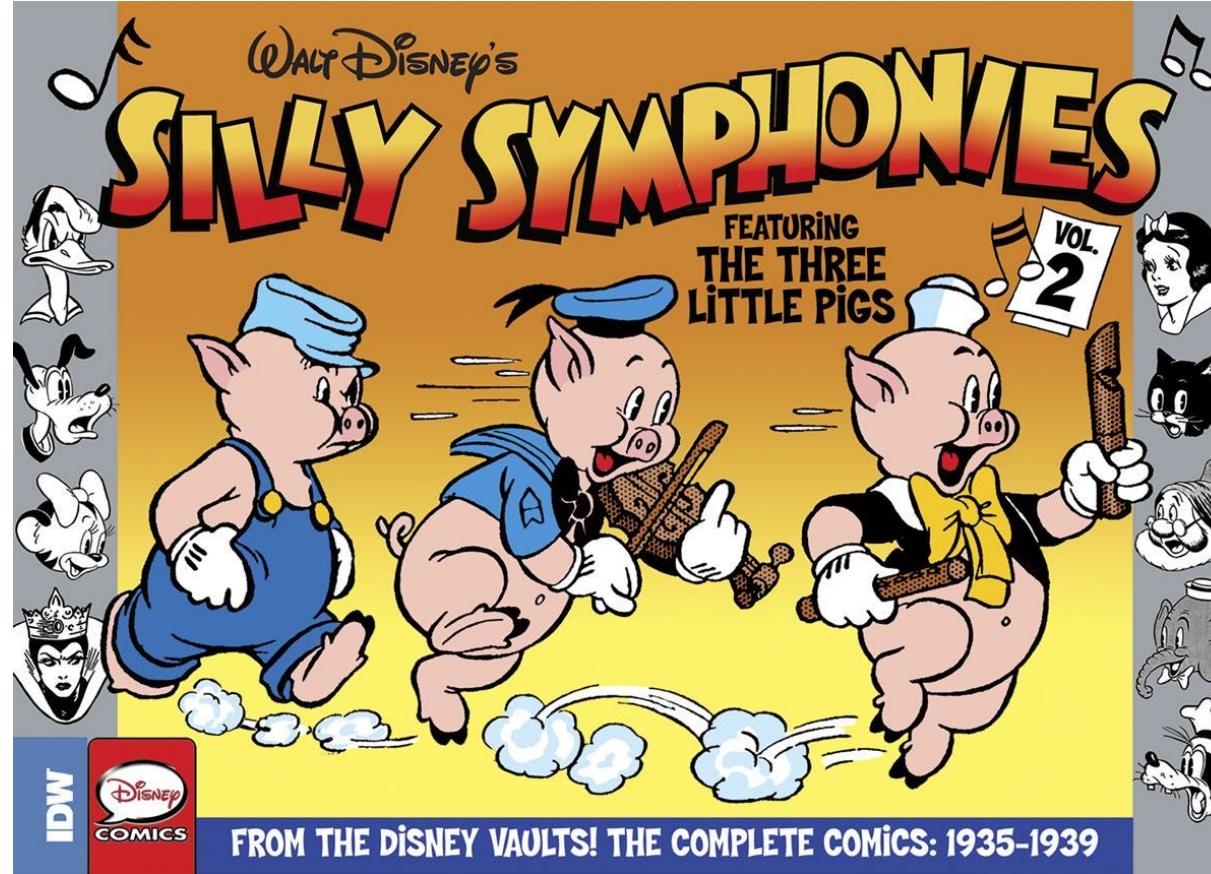
# Your communication's purpose

- If you only had a limited amount of time or a single sentence to tell your audience what they need to know, what would you say?
  - This can be particularly effective
- Alternatively....
  - The elevator pitch
  - The 3-minute pitch
  - The 3-minute story
  - ...

# Your communication's purpose

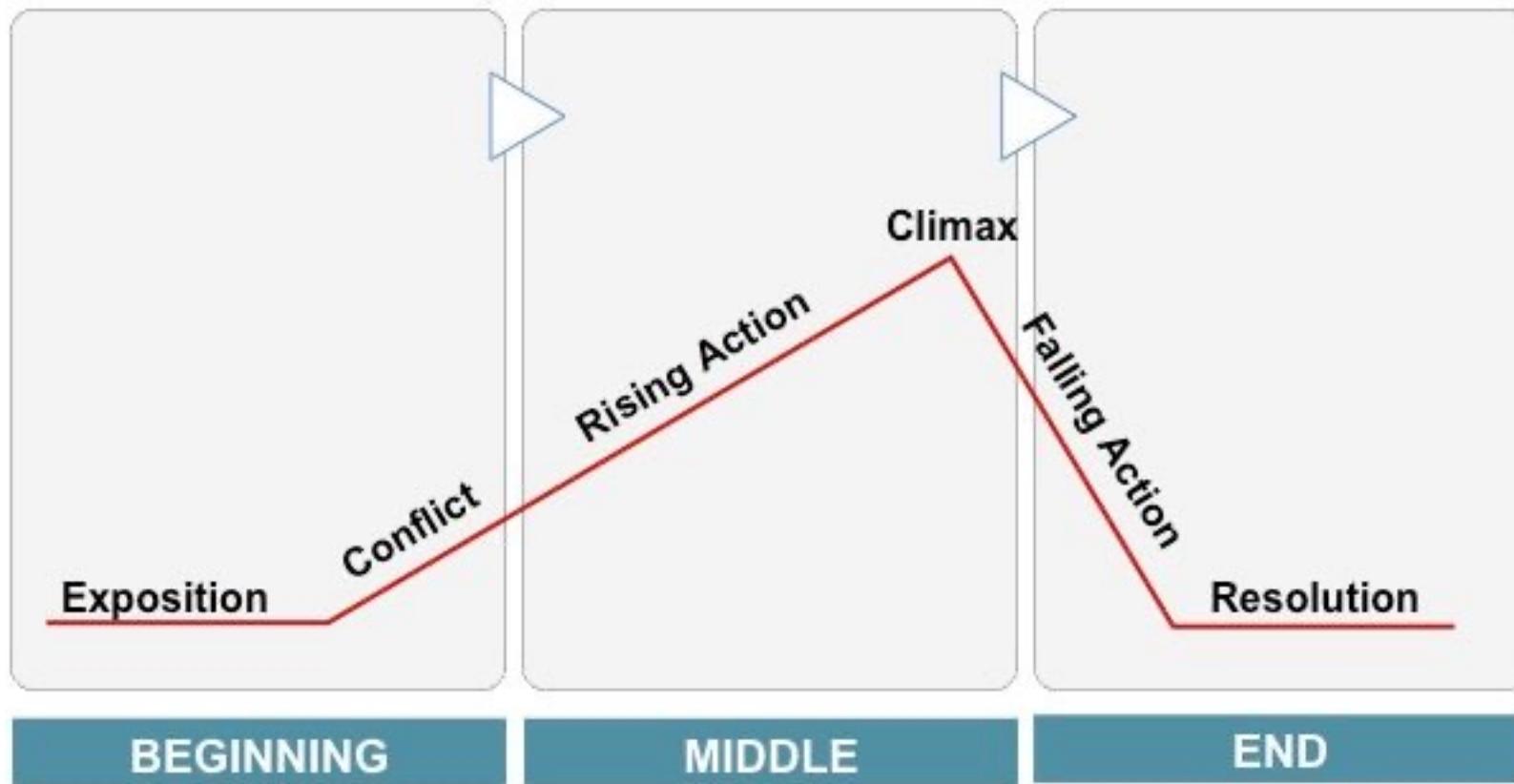
- If you only had a limited amount of time or a single sentence to tell your audience what they need to know, what would you say?
  - This can be particularly effective
- Alternatively....
  - The elevator pitch
  - The 3-minute pitch
  - The 3-minute story
  - ... -> short and story-like

# Story

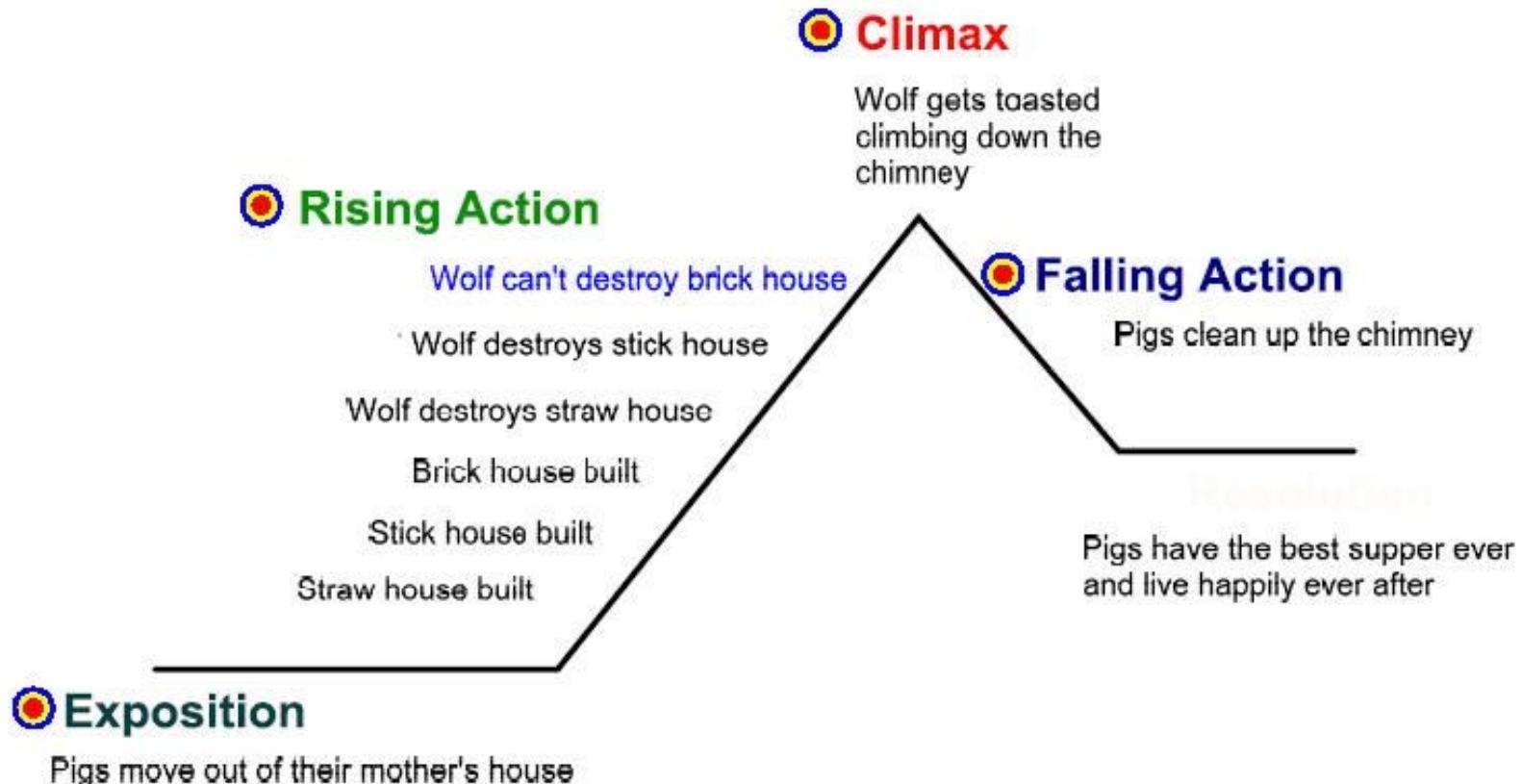


<https://www.youtube.com/watch?v=Olo923T2HQ4>

# Story Arc



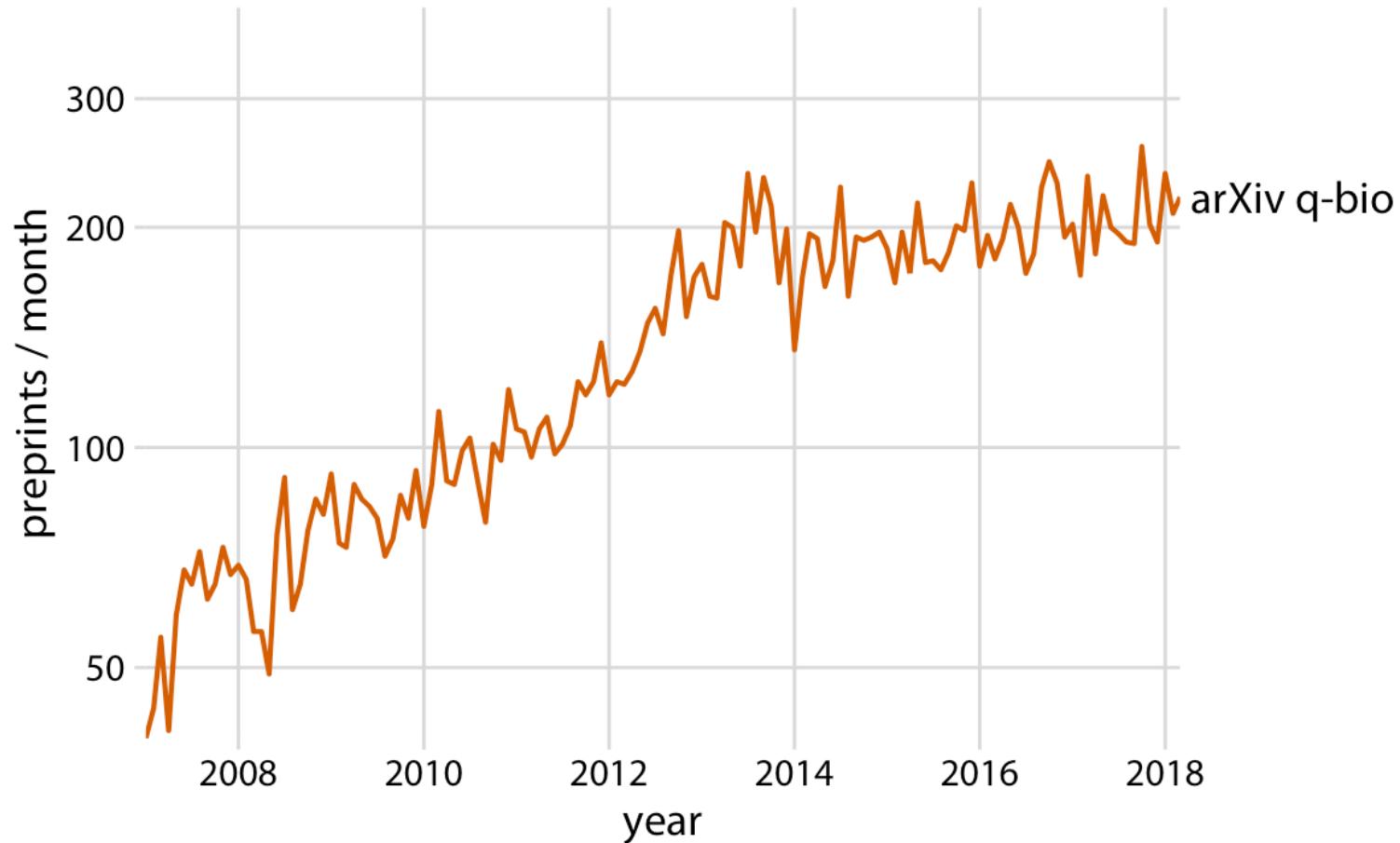
# Story Arc



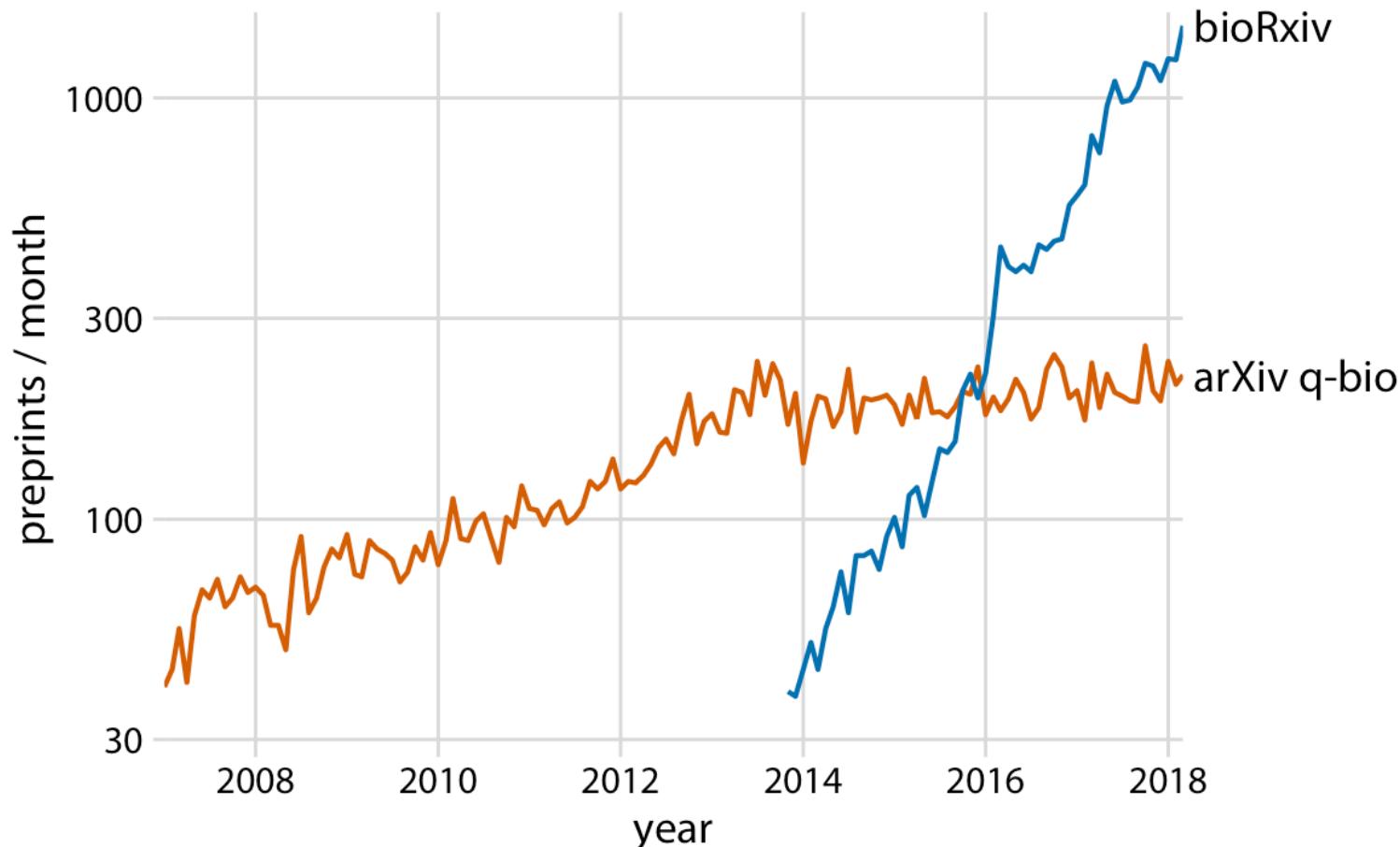
# Narrative elements

- Who are the main characters?
- Is there a conflict or issue being addressed?
- The conflict/tension/resolution can be supported in analyses in small and large ways

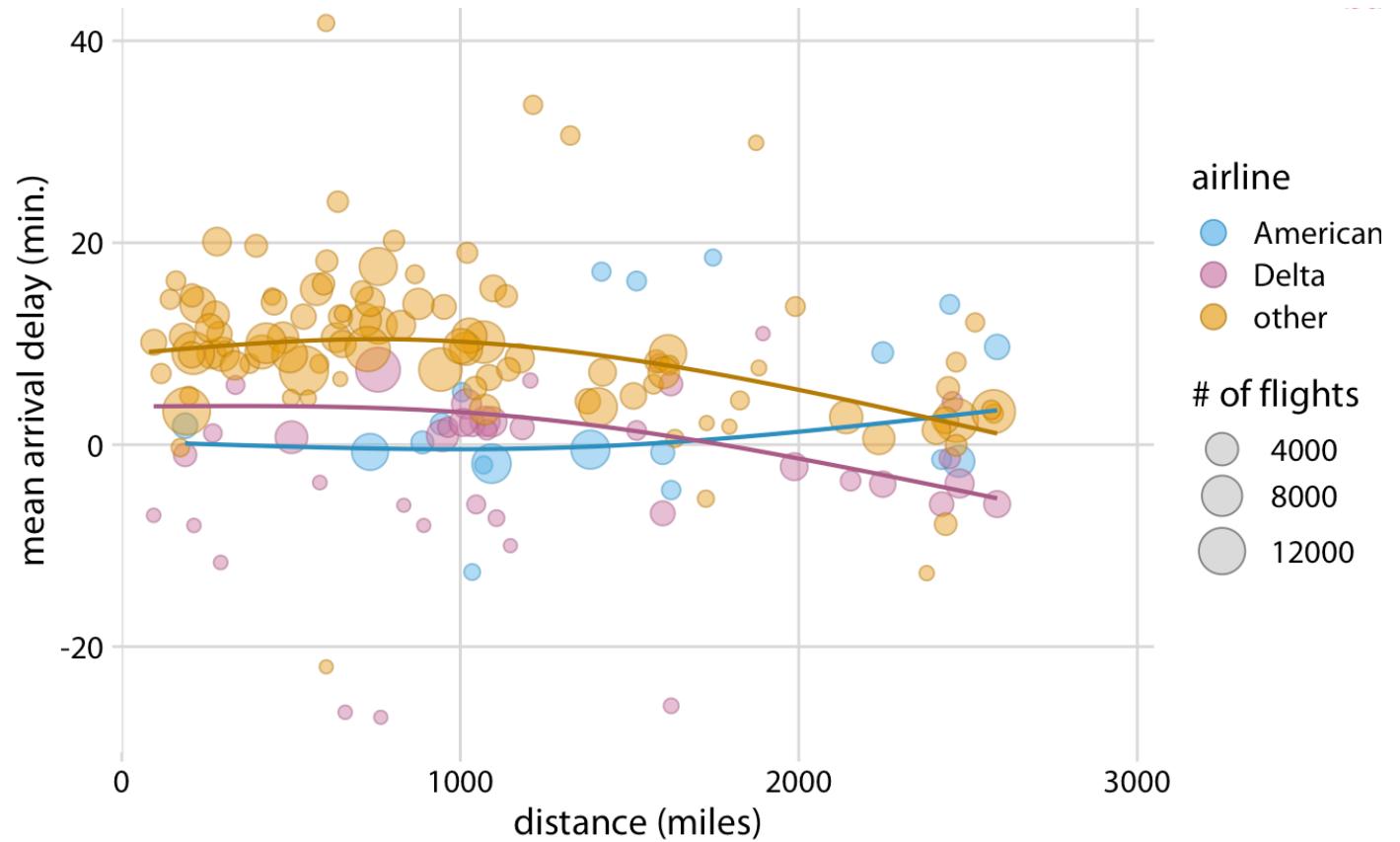
# Narrative elements in data visualization?



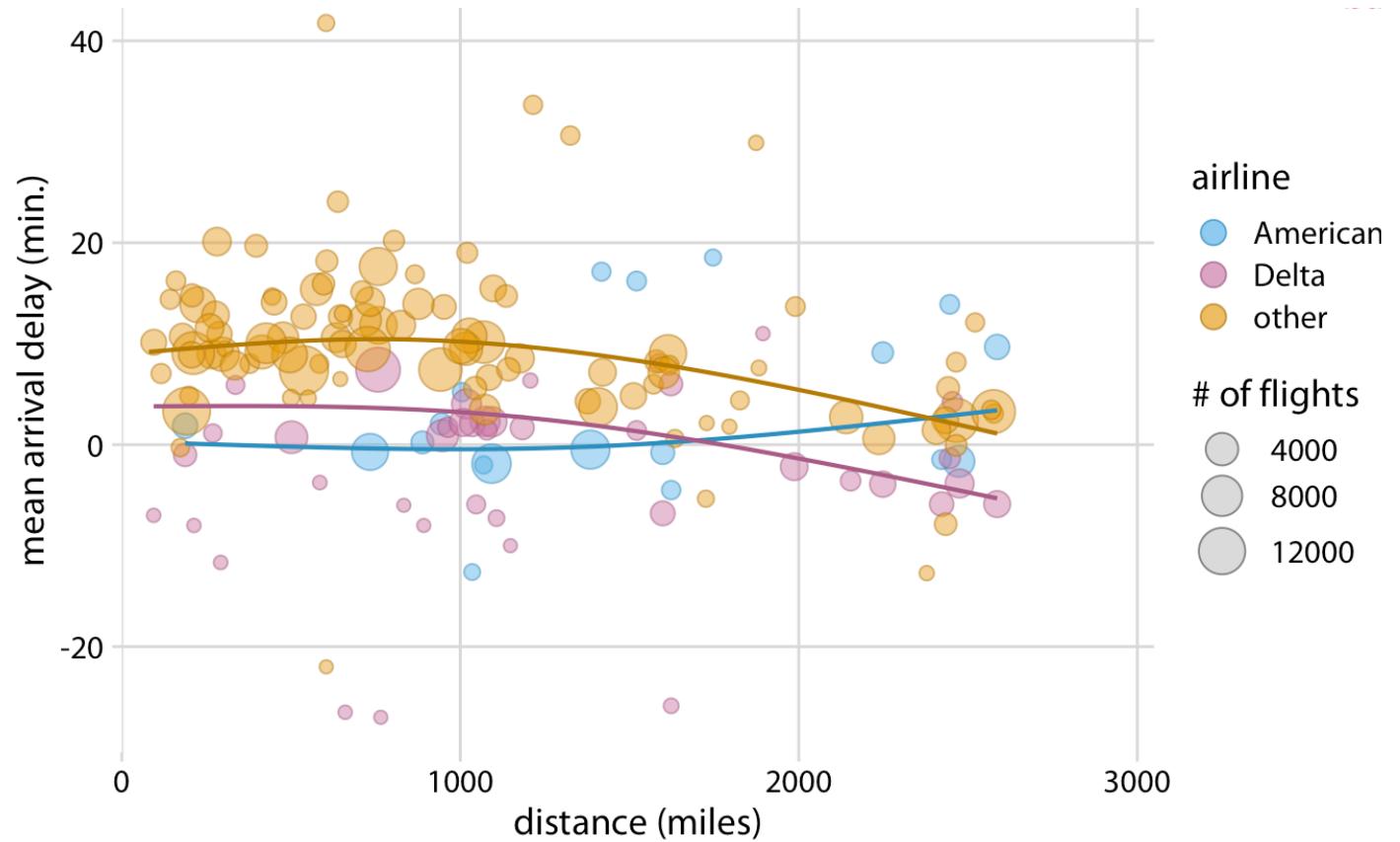
# Narrative elements in data visualization: set-up and pay-off



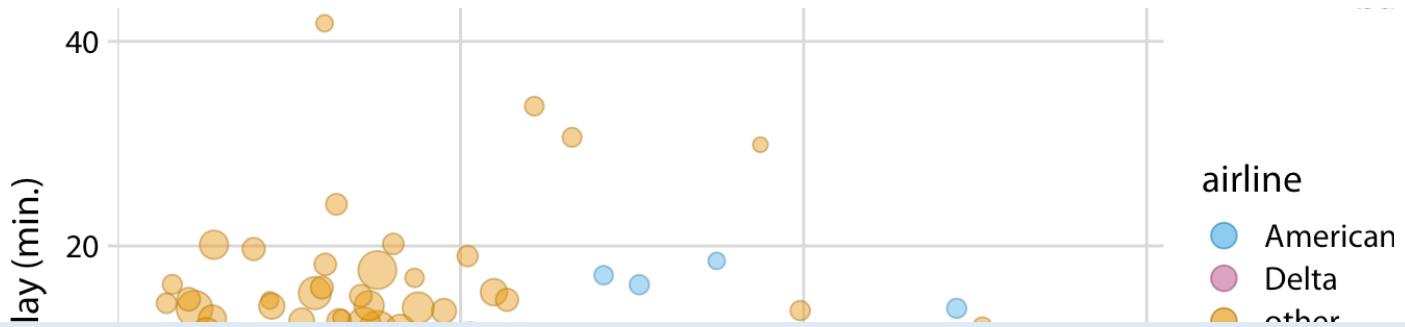
# Tell your story for the Generals / CEOs



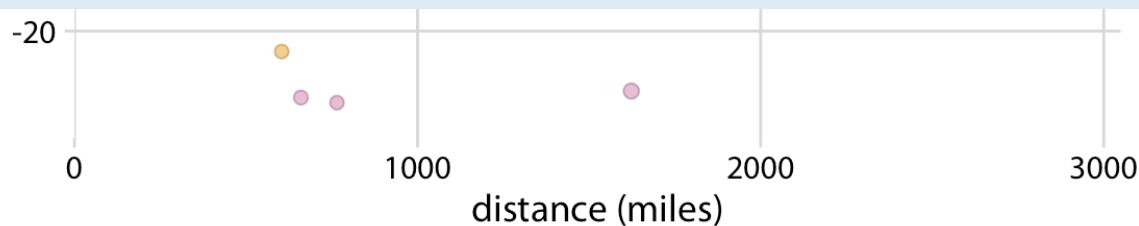
# Tell your story for the Generals / CEOs



# Tell your story for the Generals / CEOs



Never assume your audience can rapidly process complex visual displays  
(key is *rapidly*)



# Tell your story for the Generals / CEOs

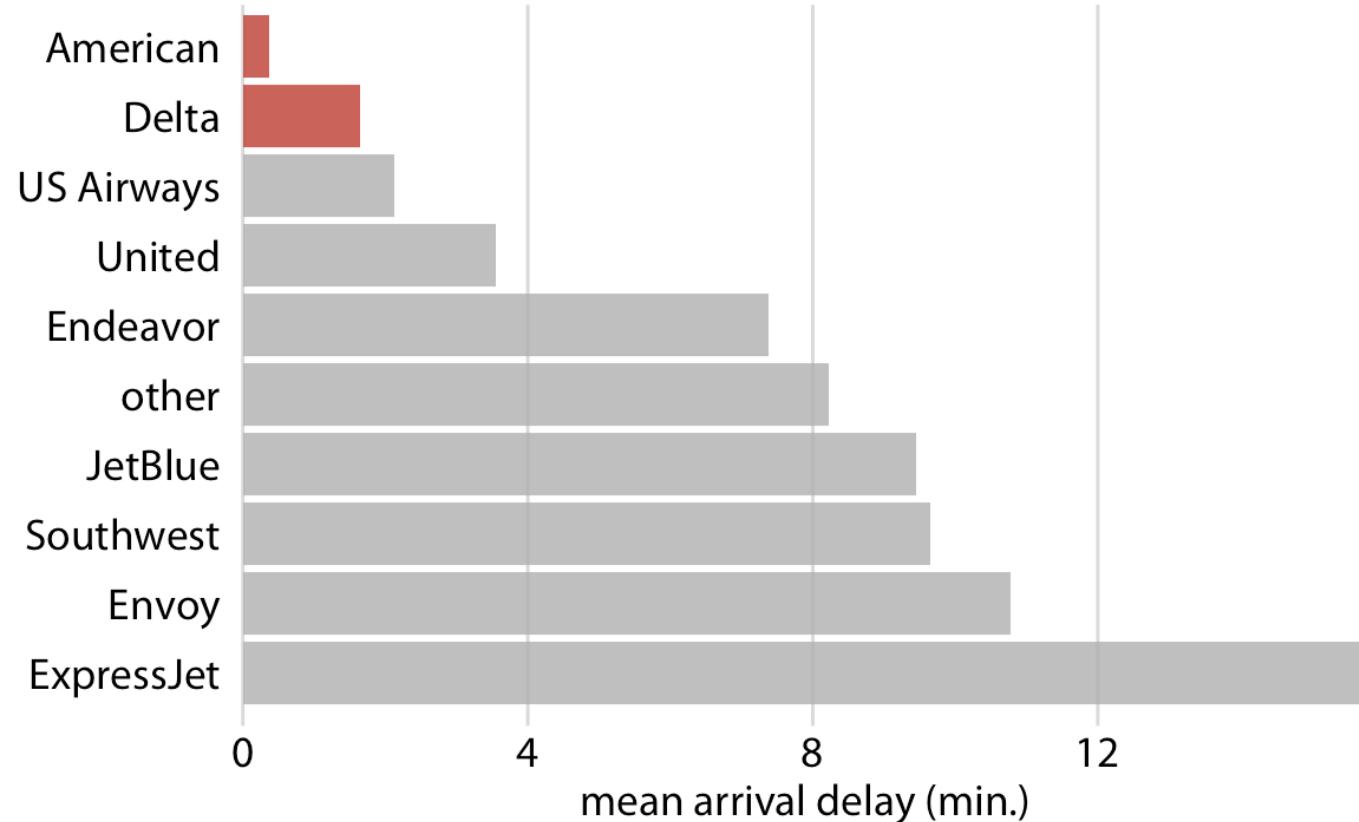


Figure 29.4: Mean arrival delay for flights out of the New York City area in 2013, by airline. American and Delta have the lowest mean arrival delays of all airlines flying out of the New York City area. Data source: U.S. Dept. of Transportation, Bureau of Transportation Statistics.

# Tell your story for the Generals / CEOs

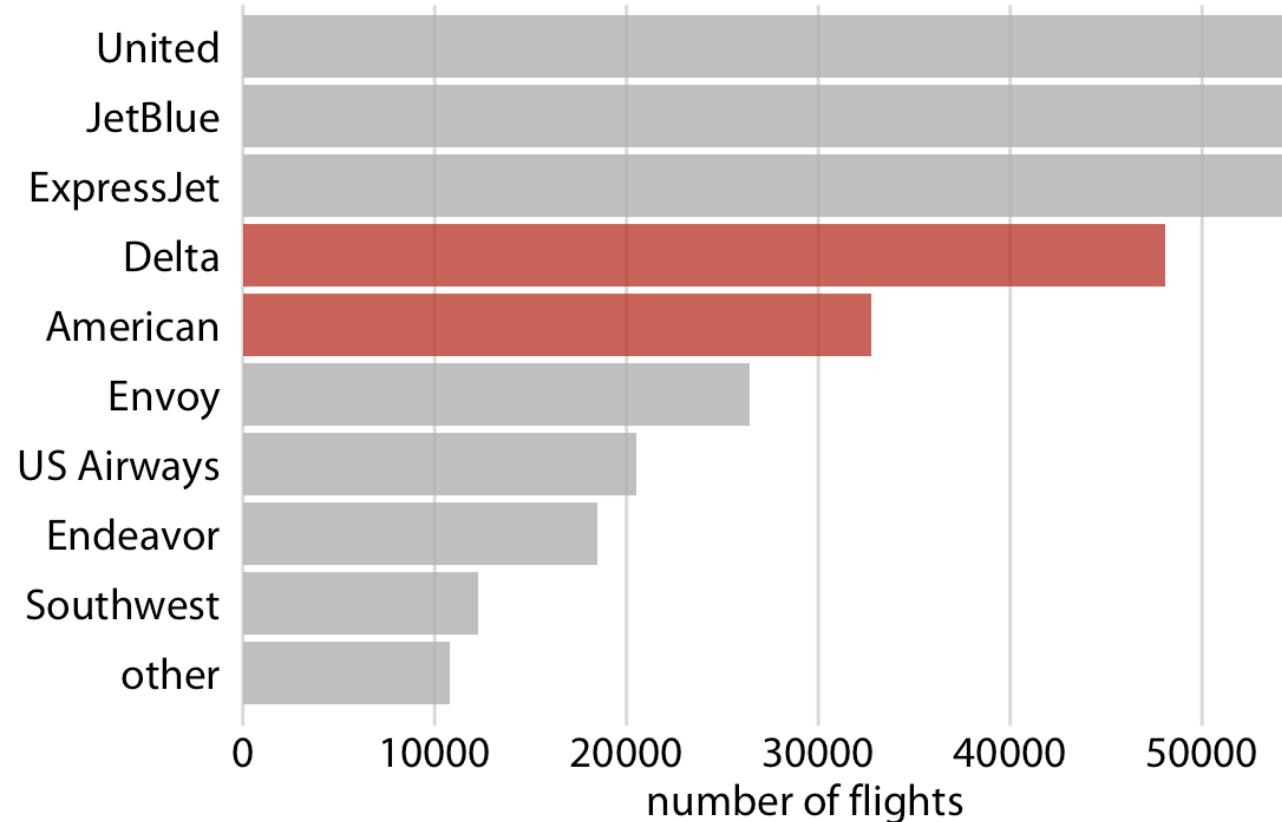
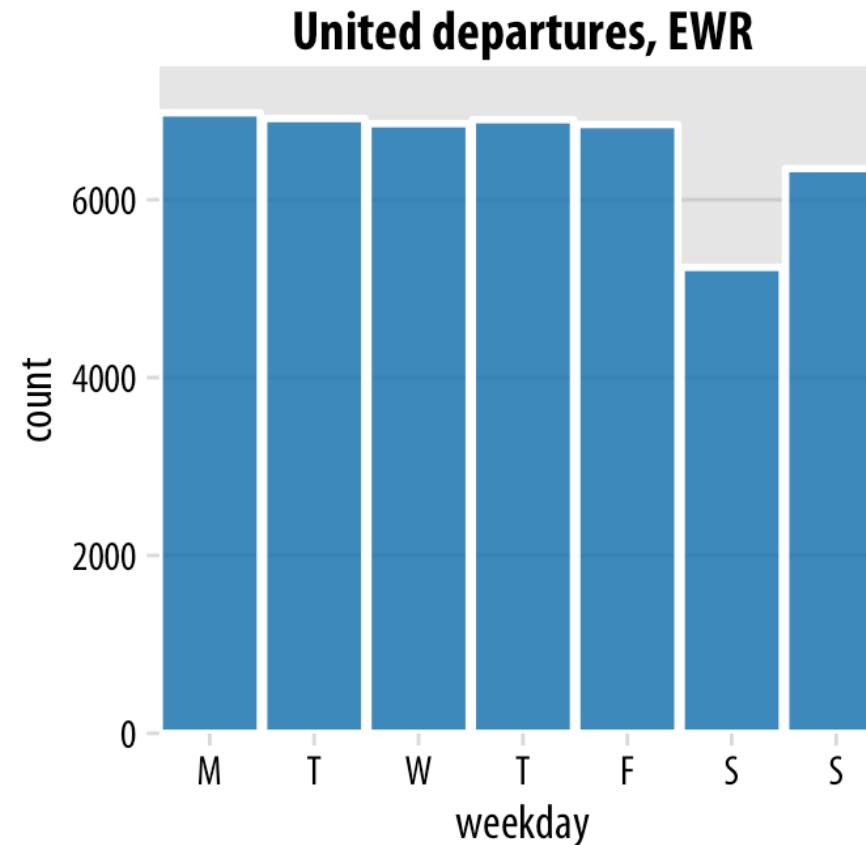
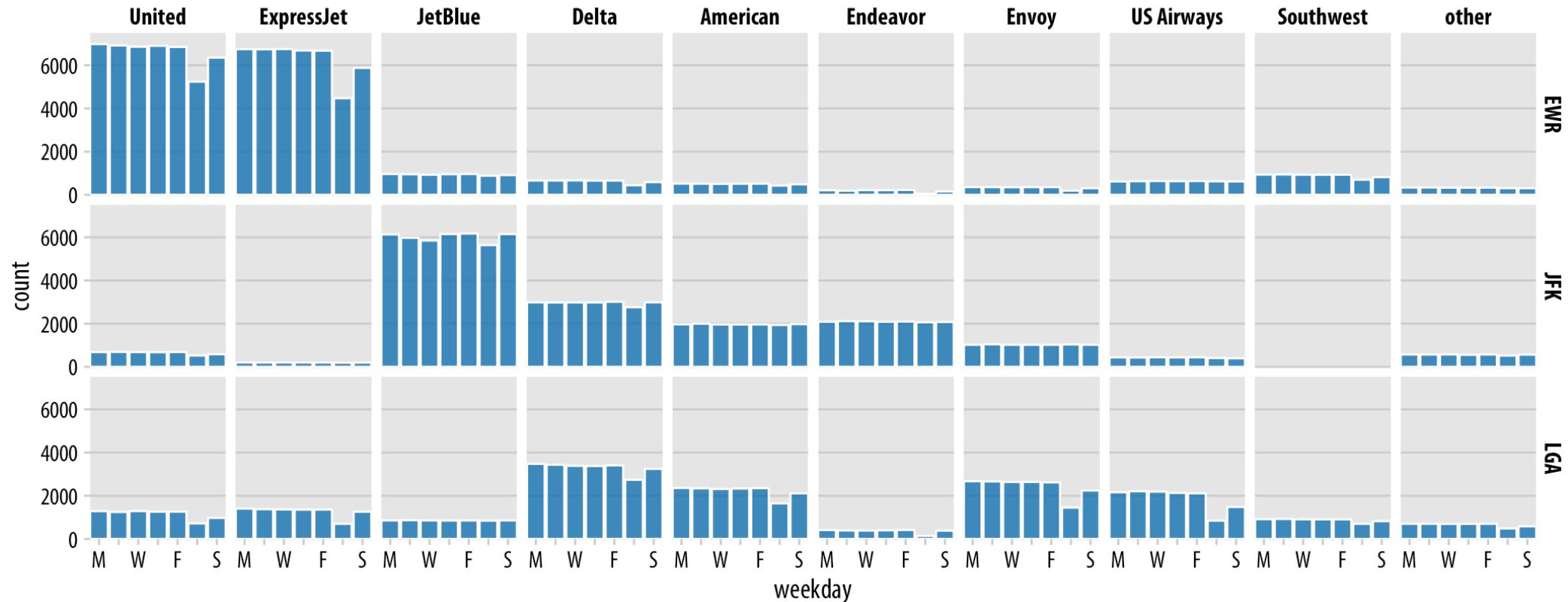


Figure 29.5: Number of flights out of the New York City area in 2013, by airline. Delta and American are fourth and fifth largest carrier by flights out of the New York City area. Data source: U.S. Dept. of Transportation, Bureau of Transportation Statistics.

# Build up towards complex figures



# Build up towards complex figures

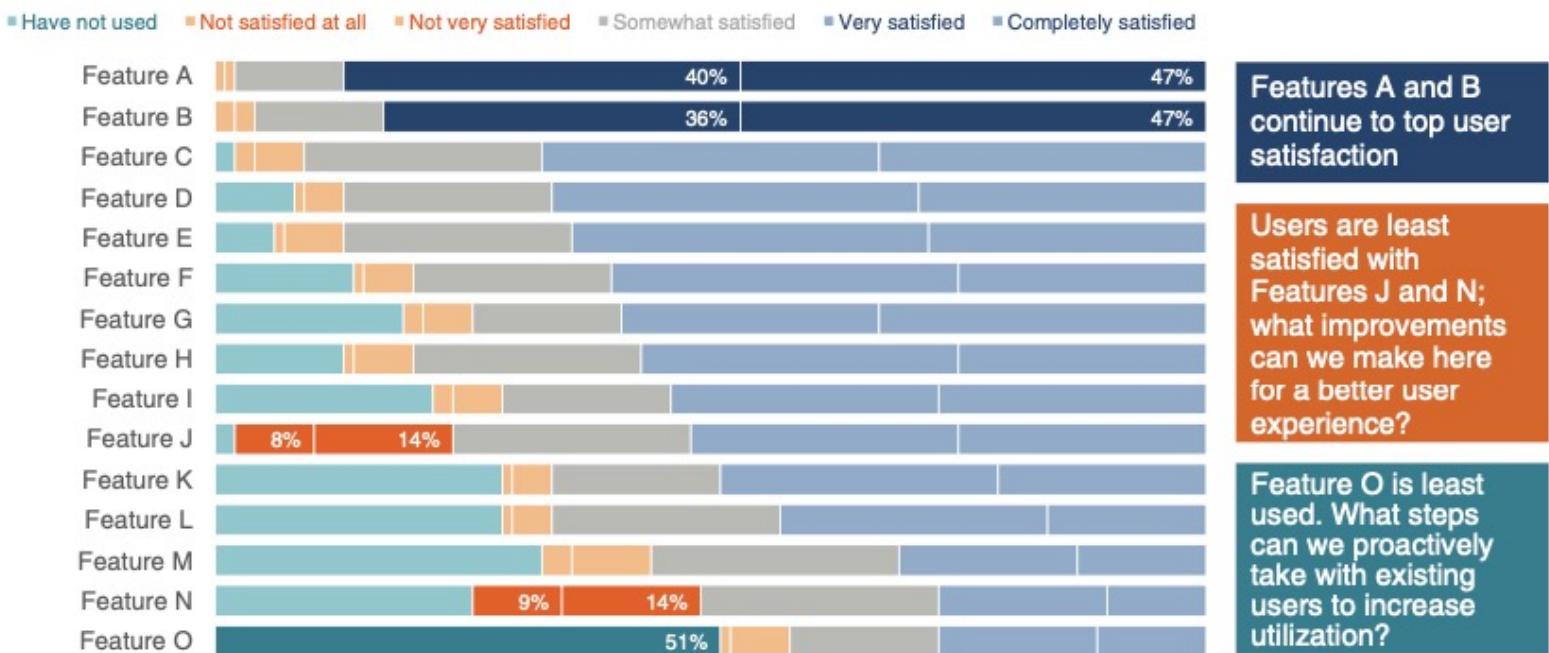


# Narrative elements

- Who are the main characters?
- Is there a conflict or issue being addressed?
- The conflict / tension / resolution can be supported in analyses in small and large ways

## User satisfaction varies greatly by feature

Product X User Satisfaction: Features



Responses based on survey question "How satisfied have you been with each of these features?".  
Need more details here to help put this data into context: How many people completed survey? What proportion of users does this represent?  
Do those who completed survey look like the overall population, demographic-wise? When was the survey conducted?

Features A and B continue to top user satisfaction

Users are least satisfied with Features J and N; what improvements can we make here for a better user experience?

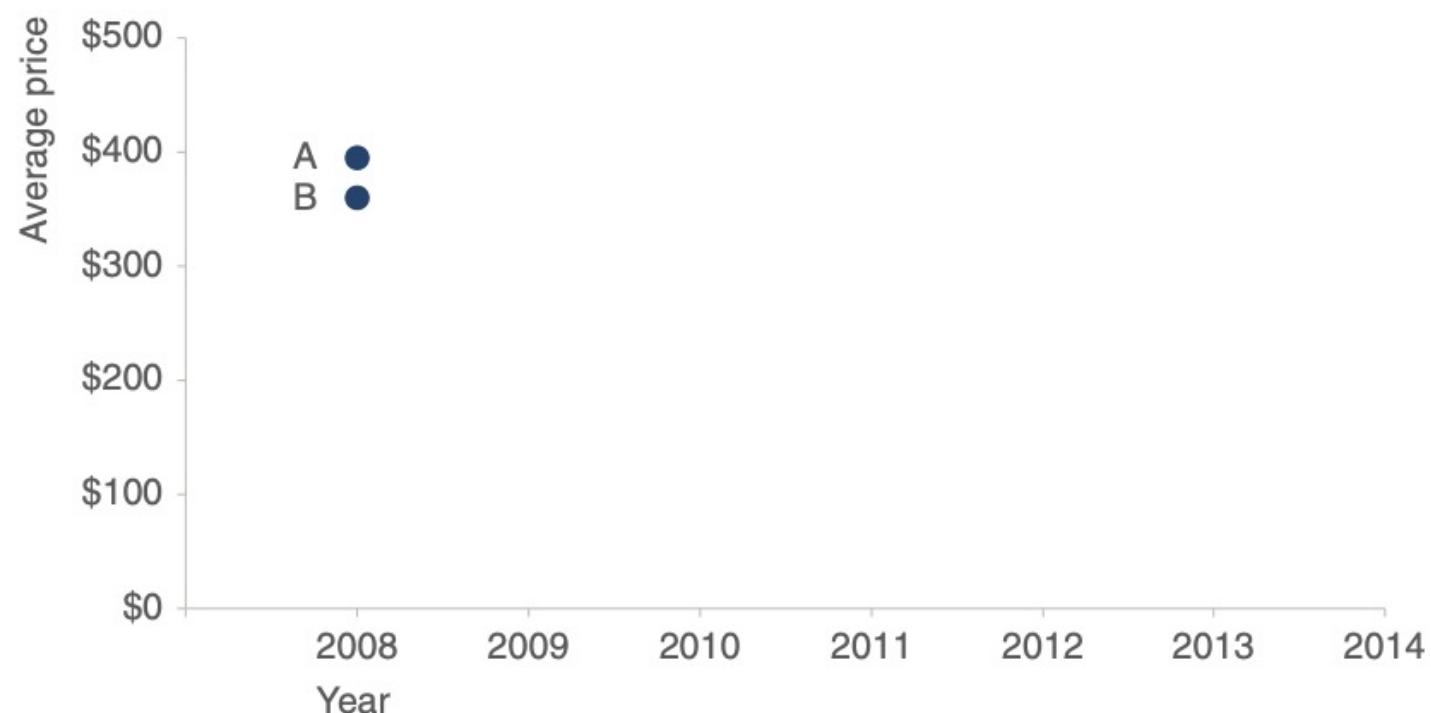
Feature O is least used. What steps can we proactively take with existing users to increase utilization?

# Telling a story through your data visualizations



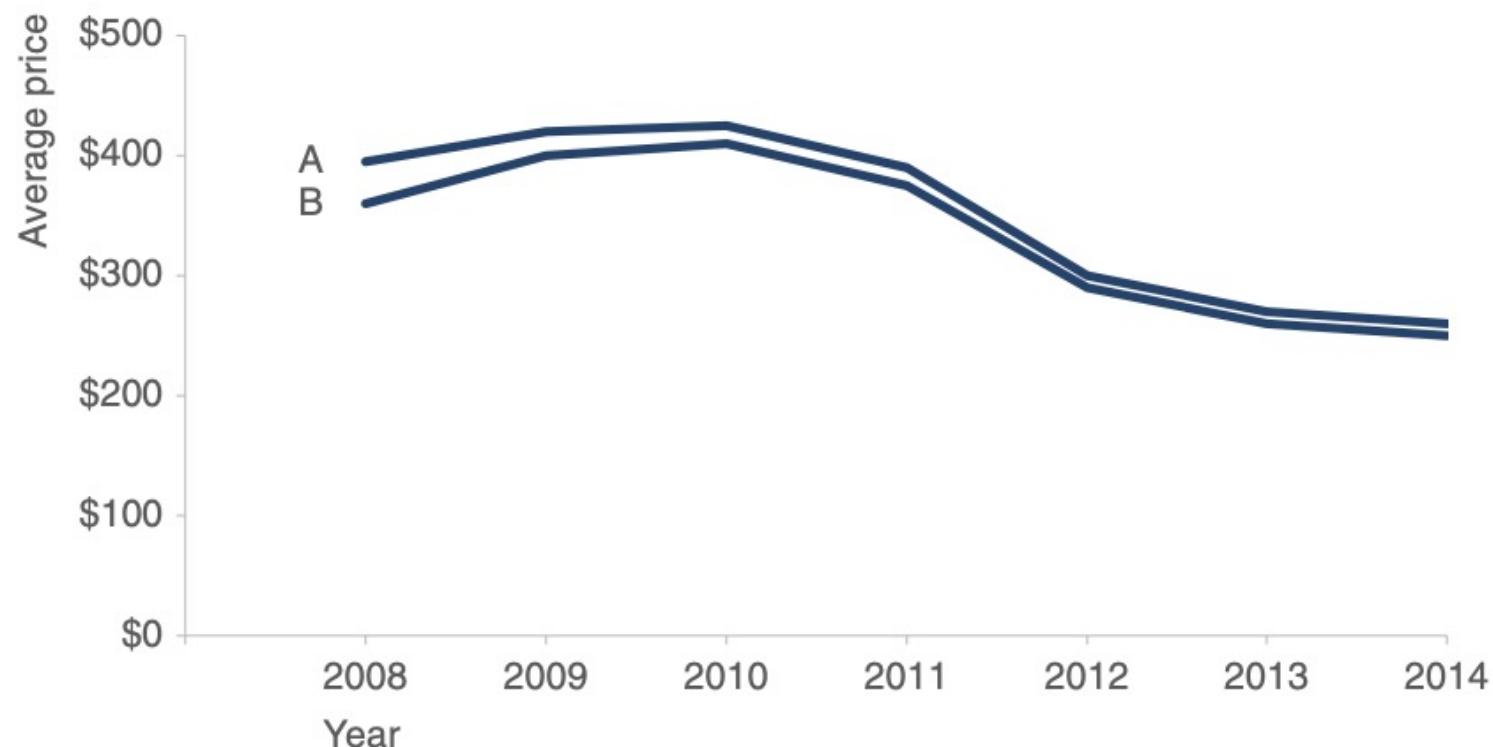
Products A and B were launched in 2008 at price points of **\$360+**

Retail price over time



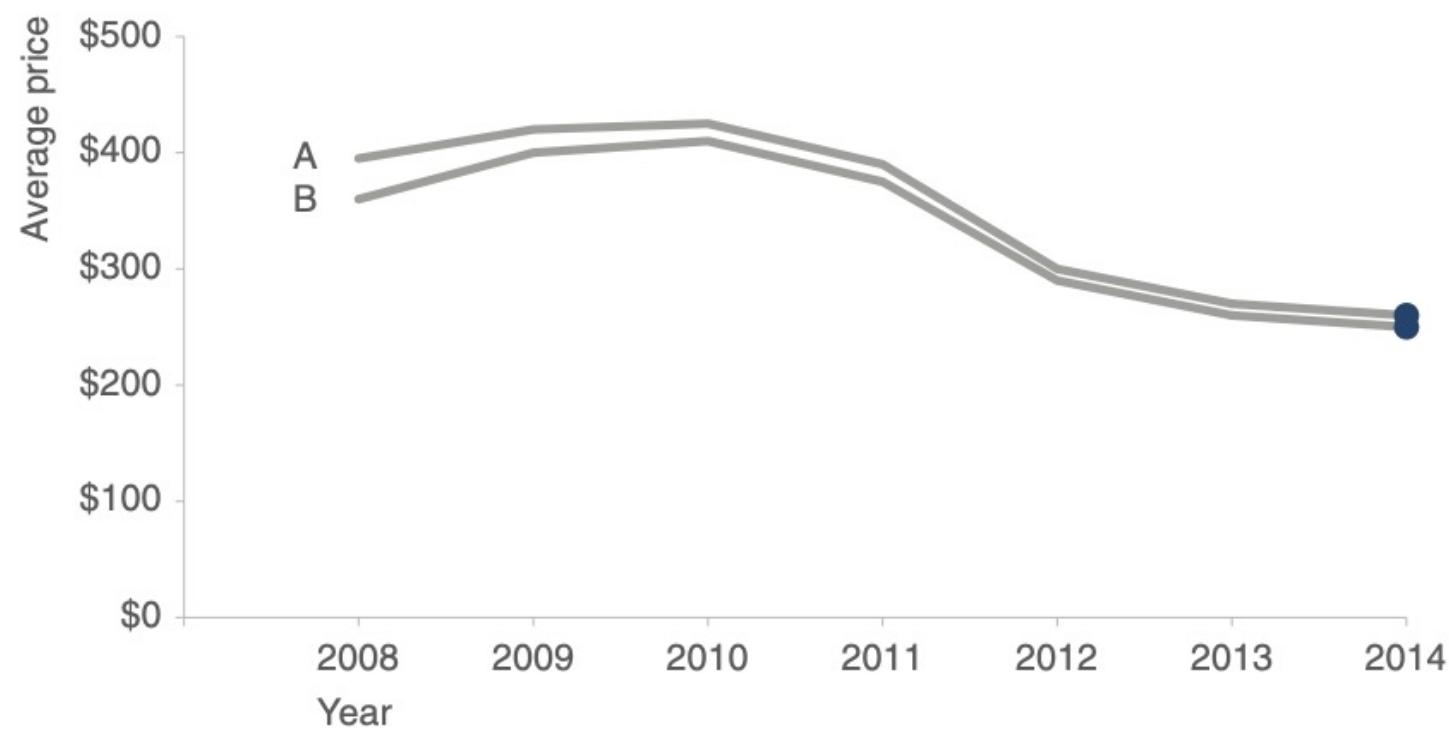
They have been priced similarly over time, with B consistently slightly lower than A

Retail price over time



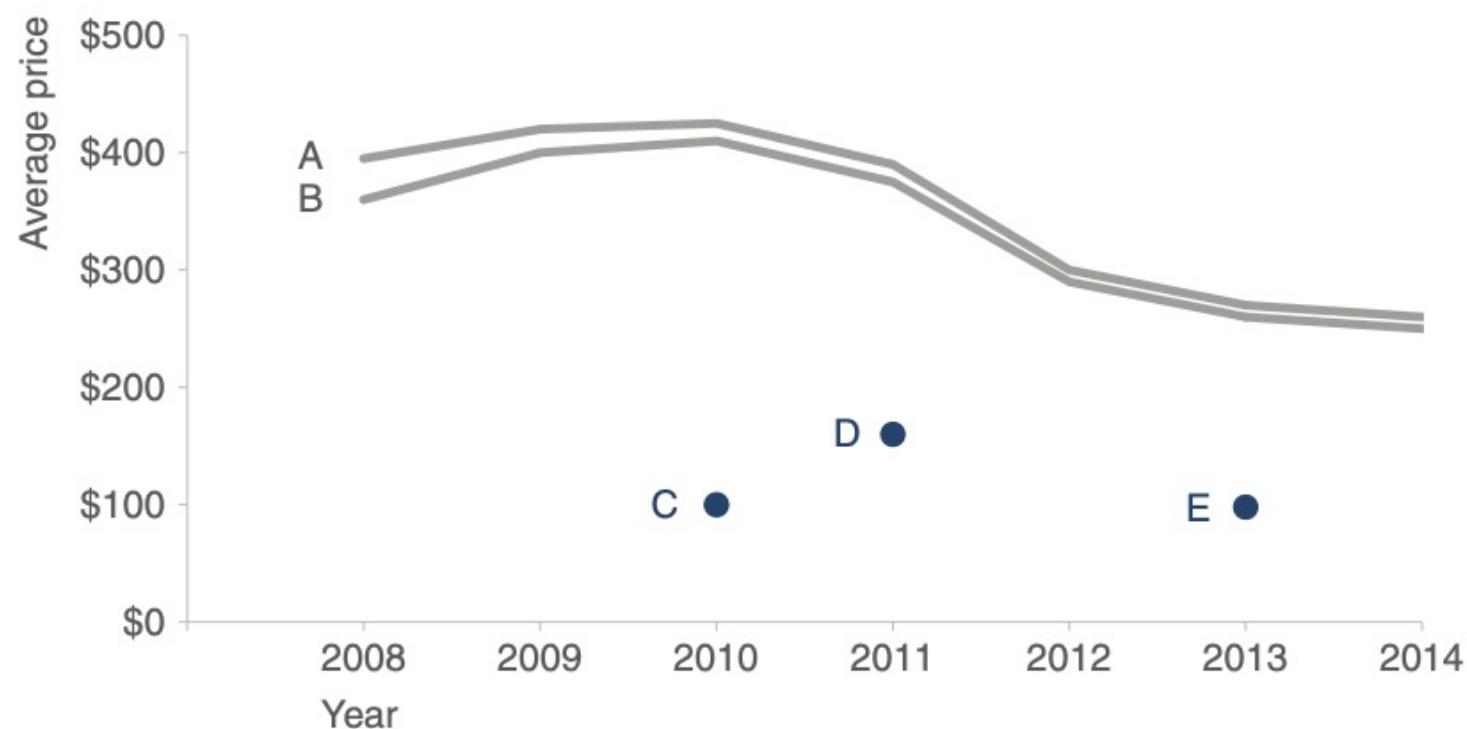
In 2014, Products A and B were priced at **\$260** and **\$250**, respectively

Retail price over time



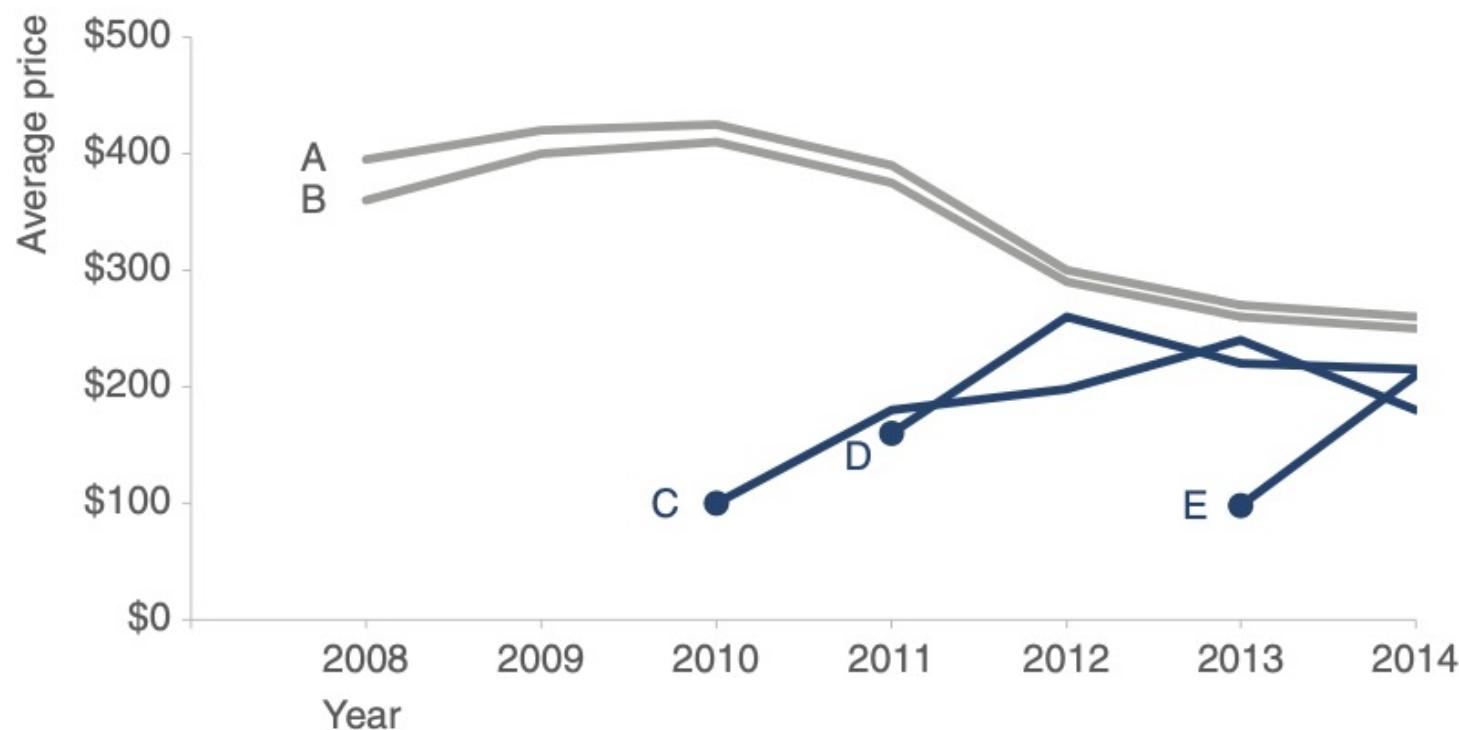
Products C, D, and E were each introduced later  
at **much lower price points**...

Retail price over time



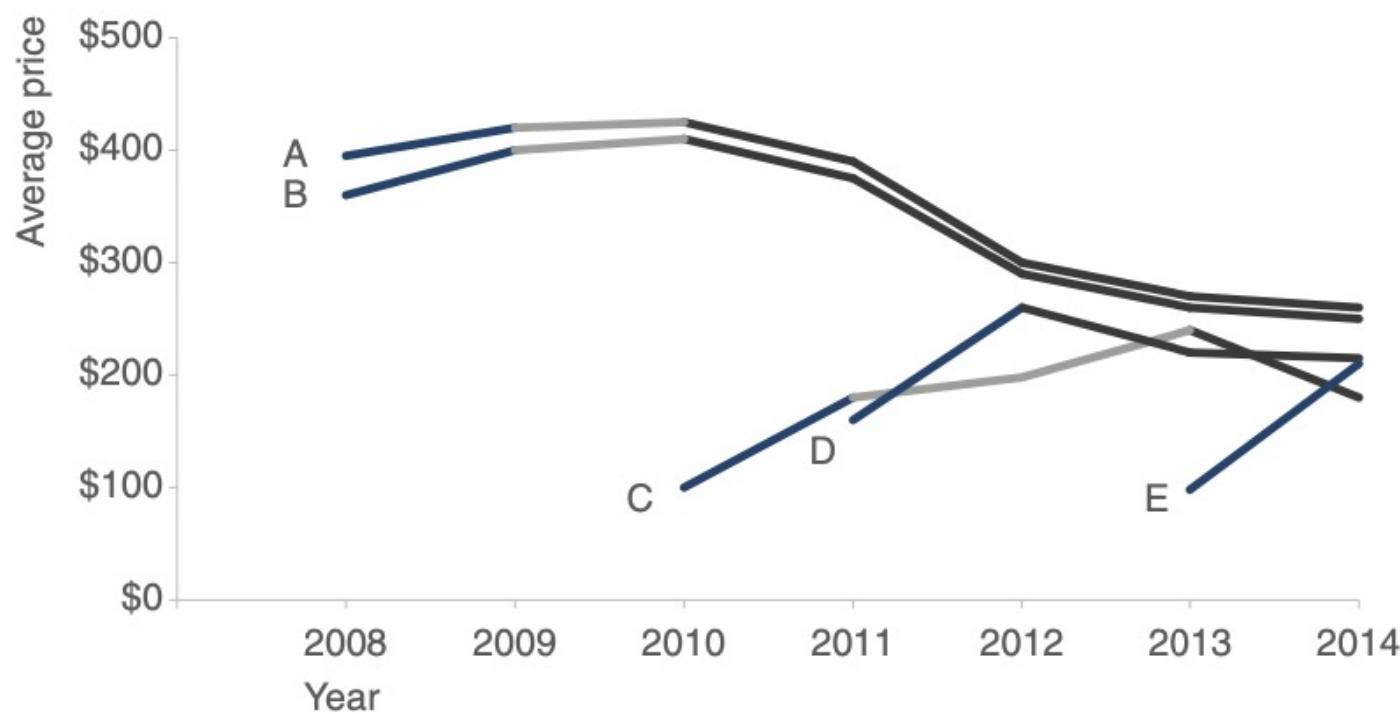
...but all have **increased** in price since their respective launches

Retail price over time



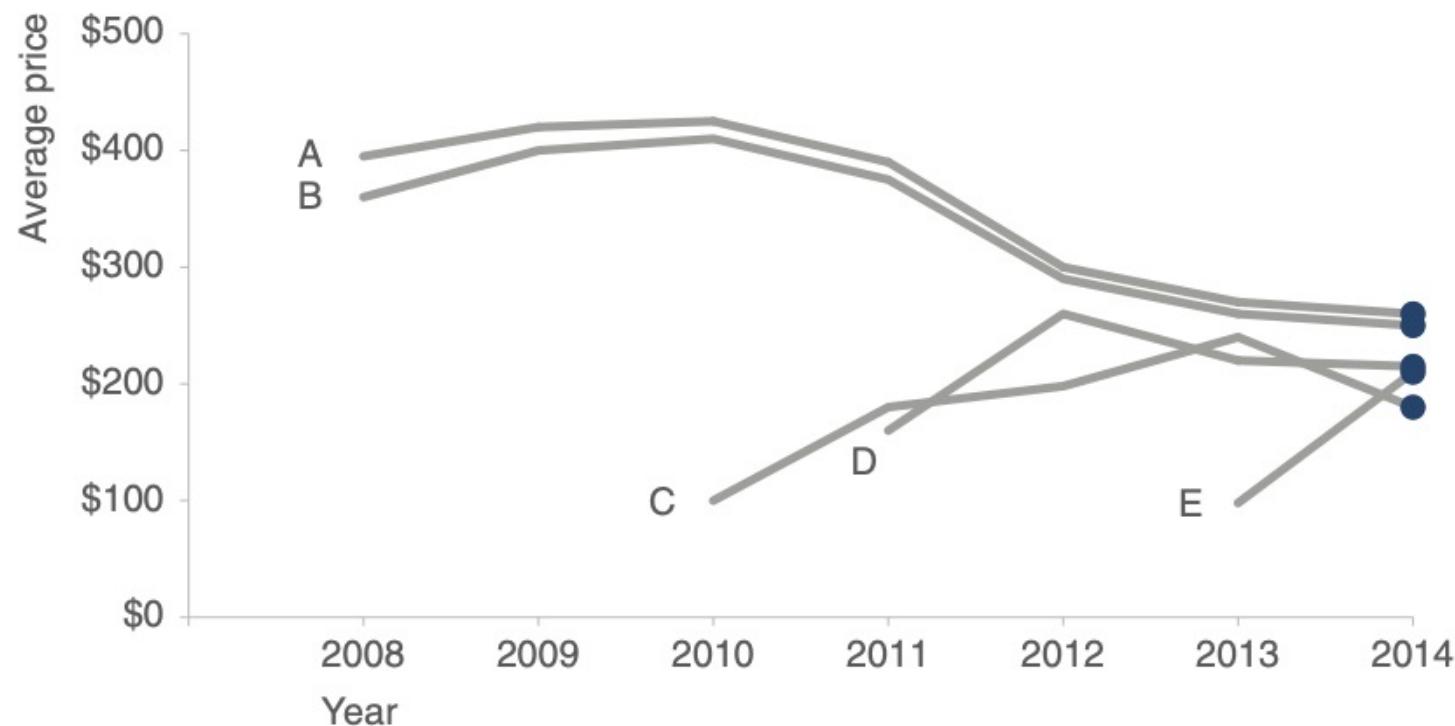
In fact, with the launch of a new product in this space, we tend to see an **initial price increase**, followed by a **decrease** over time

Retail price over time



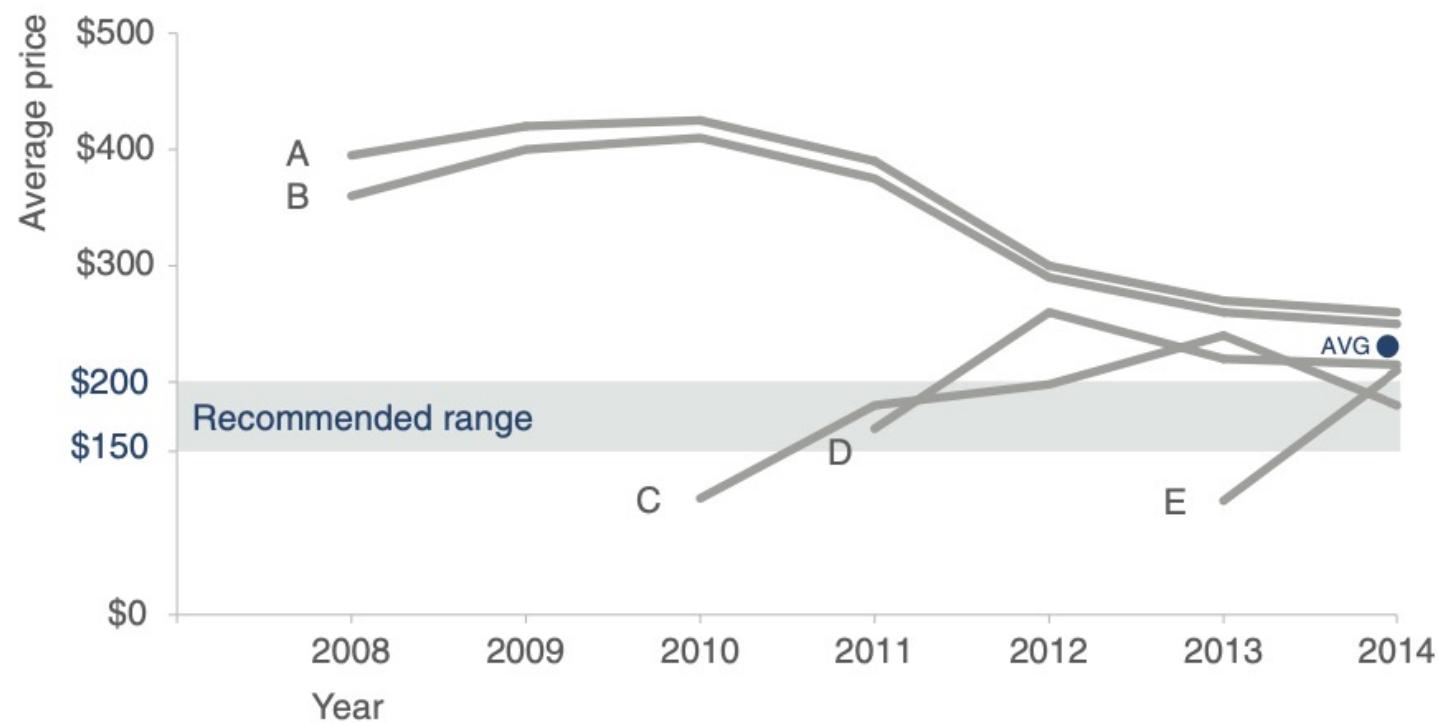
As of 2014, retail prices have converged, with an **average retail price of \$223**, ranging from a low of \$180 (C) to a high of \$260 (A)

Retail price over time



To be competitive, we recommend introducing our product *below* the \$223 average price point in the **\$150–\$200 range**

Retail price over time

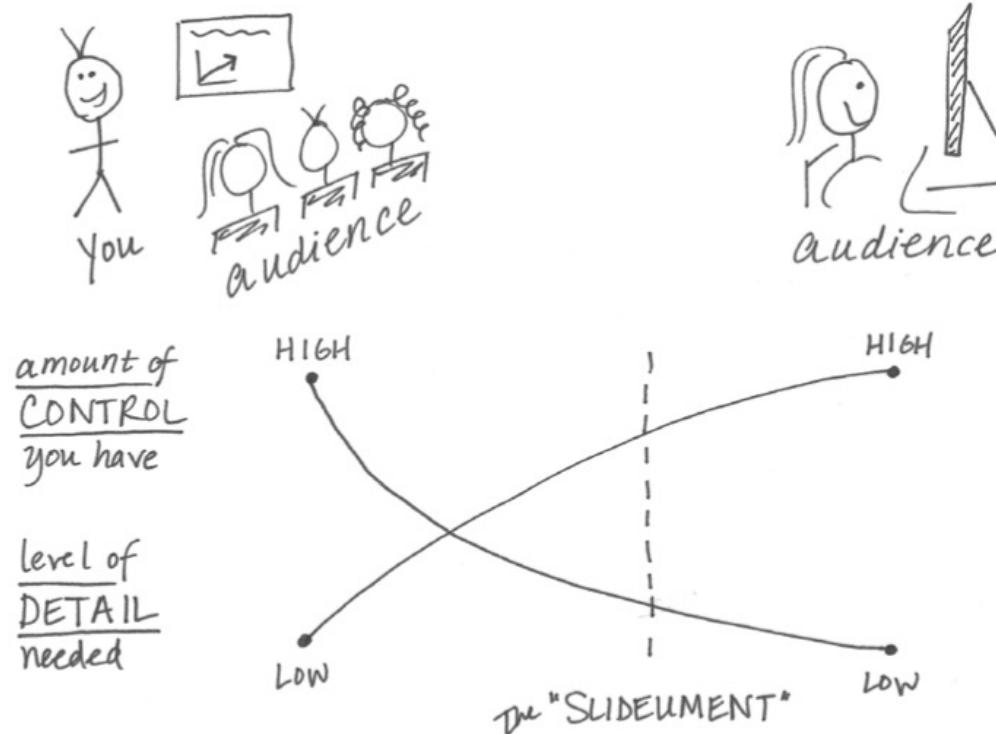


# For presentations: Don't use your slides as a teleprompter

- Write out notes
- Practice your presentation by actually talking
  - Exercise different parts of your brain
- Give mock presentation to others

# How will you communicate?

LIVE PRESENTATION..... WRITTEN DOC OR EMAIL



- For your final project:
  - Notebook
  - Presentation
- How does this frame the content and delivery?

# How to use the data

- Draw on the data visualization and analysis we have covered throughout the course
- Presentation: do not present your total analysis
  - After all the work you've done, it can be extremely tempting to describe it all... or at least many significant parts of it
  - DON'T!

Is this the result of my analysis?



Or is this the result of my analysis?



# Visualizing Data

- The bulk of our course has been focused on **how to** visualize
- The more you practice, the more comfortable and rapid you'll be
  - Both viz for exploration and viz for explanation
- Remember: Exploratory visualizations and explanatory visualizations are made for different sets of “who, what, and why/how”

Now for some more practical practice