

Information Visualization I

School of Information, University of Michigan

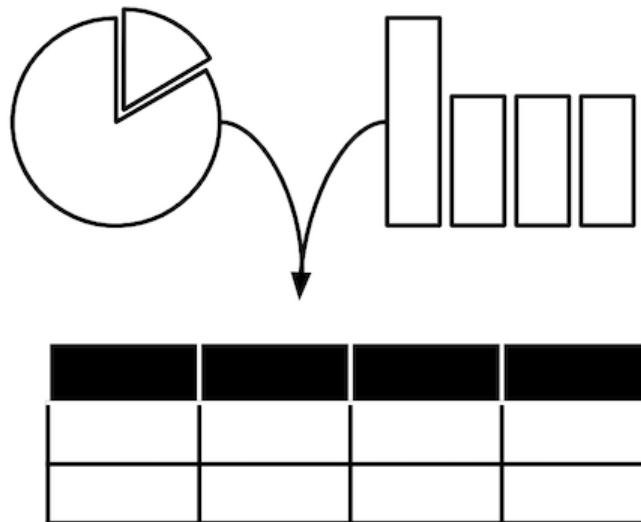
Week 2:

- Expressiveness and Effectiveness
- Grammar of Graphics

Assignment Overview

Our objectives for this week:

- Review, reflect, and apply the concepts of encoding. Given a visualization recreate the data that was encoded.
- Review, reflect, and apply the concepts of Expressiveness and Effectiveness. Given a visualization, evaluate alternatives with the same expressiveness.



Two visualizations, same expressiveness

- Review and evaluate an implementation of Grammar of Graphics using [Altair \(https://altair-viz.github.io/\)](https://altair-viz.github.io/).

The total score of this assignment will be 100 points consisting of:

- Case study reflection: Next Bechdel Test (30 points)

- Altair programming exercise (70 points)
- Bonus (5 points)

Resources:

- This article by [FiveThirtyEight](https://fivethirtyeight.com) (<https://fivethirtyeight.com>), available [online](https://projects.fivethirtyeight.com/next-bechdel/) (<https://projects.fivethirtyeight.com/next-bechdel/>) (Hickey, Koeze, Dottle, Wezerek 2017)
- Datasets from FiveThirtyEight, we have downloaded a subset of these datasets available in the folder for your use into [./assets](#) ([./assets](#))
 - The original dataset can be found on [FiveThirtyEight Next Bechdel Dataset](https://github.com/fivethirtyeight/data/tree/master/next-bechdel) (<https://github.com/fivethirtyeight/data/tree/master/next-bechdel>)

Important notes:

- 1) Grading for this assignment is entirely done by a human grader. They will be running tests on the functions we ask you to create. This means there is no autograding (submitting through the autograder will result in an error). You are expected to test and validate your own code.
- 2) You should guide your answer by the look of our examples. It doesn't need to be pixel perfect (e.g., you may not always know what our example is scaled by), but it should be pretty close.
- 3) Keep your notebooks clean and readable. If your code is highly messy or inefficient you will get a deduction.
- 4) Pay attention to the return types of your functions.
- 5) Follow the instructions for submission on Coursera. You will be providing us a generated link to a read-only version of your notebook and a PDF. When turning in your PDF, please use the File -> Print -> Save as PDF option **from your browser**. Do **not** use the File->Download as->PDF option. Complete instructions for this are under Resources in the Coursera page for this class. If you're having trouble with printing, take a look at [this video](https://youtu.be/PiO-K7AoWjk) (<https://youtu.be/PiO-K7AoWjk>).

Part 1. Expressiveness and Effectiveness (30 points)

Read the following article [Creating the next Bechdel Test](https://projects.fivethirtyeight.com/next-bechdel/) (<https://projects.fivethirtyeight.com/next-bechdel/>) and answer the following questions:

1.1 Recreate the table (by hand or excel) needed to create the following visualization (7 points)

You *should* consider the interactive parts of the visualization in your answer. Take a picture or screenshot of your table and add it to the answer below

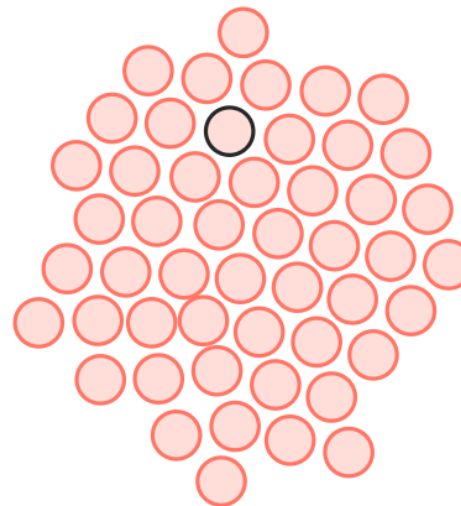
THE UPHOLD TEST

Rory Uphold: writer and actress on “This Is Why We’re Single”



A movie passes if:

- The on-set crew is 50 percent women



0 passed

All 50 failed

➤ SHOW TEST NOTES

○ THE CONJURING 2 BY THE NUMBERS

Using names to estimate gender, “The Conjuring 2” was one of the worst offenders, with men accounting for about 90 percent of the crew.

An easy way to upload images is to jump into the [./assets](#) ([./assets](#)) directory (or use the Coursera notebook explorer and navigate to it) and then use the upload button to save your image:

Files
Running
Clusters
Formgrader
Assignments

Duplicate
Rename
Move
Download
View
Edit

1

/

assignments

readonly

resources

secret

Upload

New

Notebook:

Python 3

Other:

Text File

Folder

Terminal

Once you have the image, you can link to it using the markdown command: `![answer1.2](assets/my_image_1.2.png)`

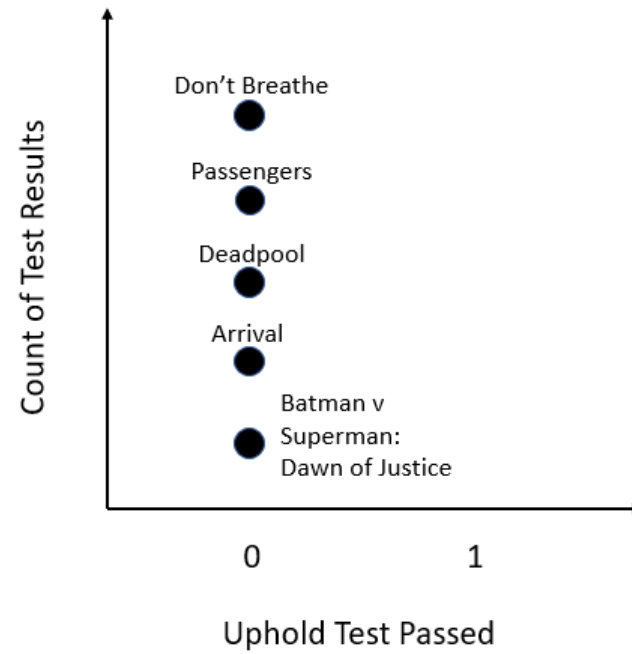
1.1 Answer

Unique Person ID	First Name	Last Name	Gender	Designation	Movie	On-Set Crew is 50% Women	Uphold Test Passed
12	Zack	Snyder	Male	Director	Batman v Superman: Dawn of Justice	FALSE	0
89652	Jeff	Markwith	Male	Set Designer	Batman v Superman: Dawn of Justice	FALSE	0
245	Bradford	Young	Male	Cinematographer	Arrival	FALSE	0
33	Julian	Clarke	Male	Editor	Deadpool	FALSE	0
6437	John	Collins	Male	Art Director	Passengers	FALSE	0
522	Carla	Vicenzino	Female	Makeup Designer	Don't Breathe	FALSE	0













1.2 Sketch an alternative visualization with the same expressiveness (7 points)

By hand is fine, but you can also use a tool. This is a sketch, the data need not be perfectly accurate or to scale. Again, upload a picture or screenshot below. Make sure there is enough annotation so it's clear why your picture has the same expressiveness.

1.2 Answer

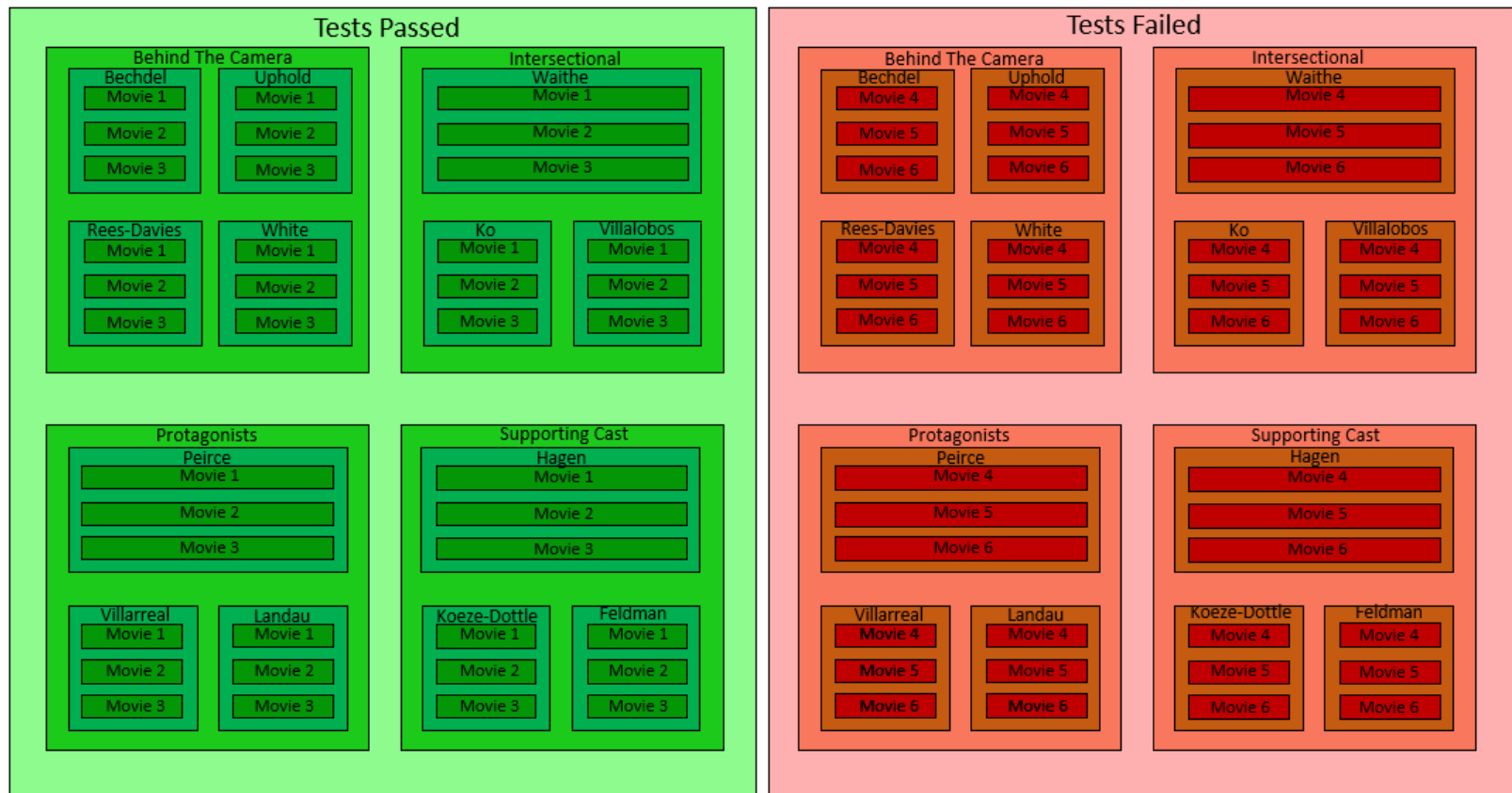


1.3 Sketch an alternative visualization with the same expressiveness of the following visualization (10 points)

MOVIES	TEST											
	BEHIND THE CAMERA				INTERSECTIONAL			PROTAGONISTS			SUPPORTING CAST	
												
Bad Moms	✓				✓	✓		✓	✓	✓	✓	✓
Hidden Figures	✓				✓	✓		✓	✓	✓		✓
Independence Day: Resurgence	✓				✓	✓		✓	✓		✓	✓
Finding Dory	✓		✓					✓	✓		✓	✓
Ghostbusters	✓					✓		✓	✓	✓		✓
Allegiant	✓					✓		✓	✓			✓
Arrival	✓							✓	✓		✓	✓
Ice Age: Collision Course	✓		✓			✓				✓	✓	

Same deal as last question: by hand or with a tool is fine. The data need not be perfectly accurate or to scale. Make sure there is enough annotation so it's clear why your picture has the same expressiveness. Again, upload a picture or screenshot below.

1.3 Answer



1.4 Reflect on which visualization you think is more *effective* and why? (6 points)

You are comparing the original figure in 1.3 and the one you created.

1.4 Answer

The figure I created in lieu of the original figure in 1.3 is not drawn to scale. It is a treemap diagram, which, if drawn to scale, would not only show which movies passed or failed various tests (Bechdel, White, etc.) in various test categories (Behind The Camera, Intersectional, etc.) but would also show how many movies fall into those categories based on the size of the rectangles. The bigger the rectangle for a group, the more values it contains. This makes the figure I created more effective than the original one since the original one just shows which movies passed or failed under each category using a tick mark.

While the movies in the original figure are ordered with the ones having most passes at the top and the ones with the least passes at the bottom, this isn't immediately obvious like it would be in a treemap diagram, which would make these differences apparent based on the size of the rectangles. For instance, the size difference would make it easier to compare between two different movies rather than counting the tick marks. Moreover, an interactive treemap diagram allows you to go deeper into the diagram and get more details on clicking the rectangles of your choice, offering it a powerful drill-down functionality which can be very useful for understanding and explaining the elements in the diagram better.

Part 2. Altair programming exercise (70 points)

We have provided some code to create visualizations based on the following datasets:

1. [all_tests \(assets/nextBechdel_allTests.csv\)](#) Is a collection of different Bechdel test results for the 50 top-grossing films [at the domestic box office in 2016 \(https://www.the-numbers.com/box-office-records/domestic/all-movies/cumulative/released-in-2016\)](https://www.the-numbers.com/box-office-records/domestic/all-movies/cumulative/released-in-2016).
2. [cast_gender \(assets/nextBechdel_castGender.csv\)](#) Is the gender for all the cast member of each movie in the Bechdel rankings
3. [top_2016 \(assets/top_2016.csv\)](#) Is the date, box office and theater count for each top 2016 movie.

Complete each assignment function and run each cell to generate the final visualizations

```
In [1]: import pandas as pd
import numpy as np
import altair as alt
```

```
In [2]: # enable correct rendering
alt.renderers.enable('default')
```

```
Out[2]: RendererRegistry.enable('default')
```

```
In [3]: # uses intermediate json files to speed things up
alt.data_transformers.enable('json')
```

```
Out[3]: DataTransformerRegistry.enable('json')
```

```
In [4]: def load_bechdel_data(alldata='assets/nextBechdel_allTests.csv',
                             castgenderdata='assets/nextBechdel_castGender.csv',
                             top2016data='assets/top_2016.csv'):
    # read all the tables
    all_tests_df = pd.read_csv(alldata)
    cast_gender = pd.read_csv(castgenderdata)
    top_2016 = pd.read_csv(top2016data)

    # set up the tables for use
    act_movies = top_2016.set_index('Movie').join(cast_gender.set_index('MOVIE')).join(all_tests_df.set_index('movie')).reset_index().dropna()
    mov_order = top_2016.sort_values(by=['Rank'])['Movie'].tolist()
    return(act_movies,mov_order)
```

```
In [5]: actors_movies,movies_order = load_bechdel_data()
```

2.1 Variables encoded (5 points)

Warmup: how many variables are encoded in the following visualization?

We also suggest you look closely this code and make sure you understand what we're doing. Most of the other problems below will follow this structure.


```
In [6]: def get_base_vis(indf):  
        #input: indf, the base chart dataset (i.e., actors_movies)  
  
        # this is a "base" chart -- on its own, it will not display anything because we haven't defined a mark  
        # but we can do operations that we'll want to do multiple times, in this case filtering some data  
        return alt.Chart(indf).transform_filter(  
            (alt.datum.TYPE != 'Unknown') & (alt.datum.GENDER != 'Unknown') & (alt.datum.GENDER != 'null')  
        )  
  
base_vis = get_base_vis(actors_movies)
```

```

In [7]: # this is where we build the actual chart
def gen_f_bar_vis(base, indf):
    # input: base is the 'base' vis as produced by get_base_vis
    # indf: a dataframe like actors_movies -- we *could* use this, but it's easier to
    #       use the 'base' chart so we're not repeating work

    # we will "add" to the base chart with an additional transform filter
    encoding = base.transform_filter(
        alt.datum.GENDER == 'Female'
    ).encode(
        # adding all the encoding stuff
        y= alt.Y(
            'index:N',
            sort= movies_order
        ),
        x=alt.X('count(index):Q',
            title='cast count'),
    ).mark_bar( # use a "bar" symbol
    ).properties(
        # set the title
        title='Female'
    )

    # The less efficient way (where we don't use "base" but repeat work):

    # encoding = alt.Chart(indf).mark_bar().transform_filter(
    #     (alt.datum.TYPE != 'Unknown') & (alt.datum.GENDER != 'Unknown') & (alt.datum.GENDER != 'null')
    # ).transform_filter(
    #     alt.datum.GENDER == 'Female'
    # ).encode(
    #     y= alt.Y(
    #         'index:N',
    #         sort= movies_order
    #     ),
    #     x=alt.X('count(index):Q',
    #         title='cast count'),
    # ).properties(title='Female')
    #

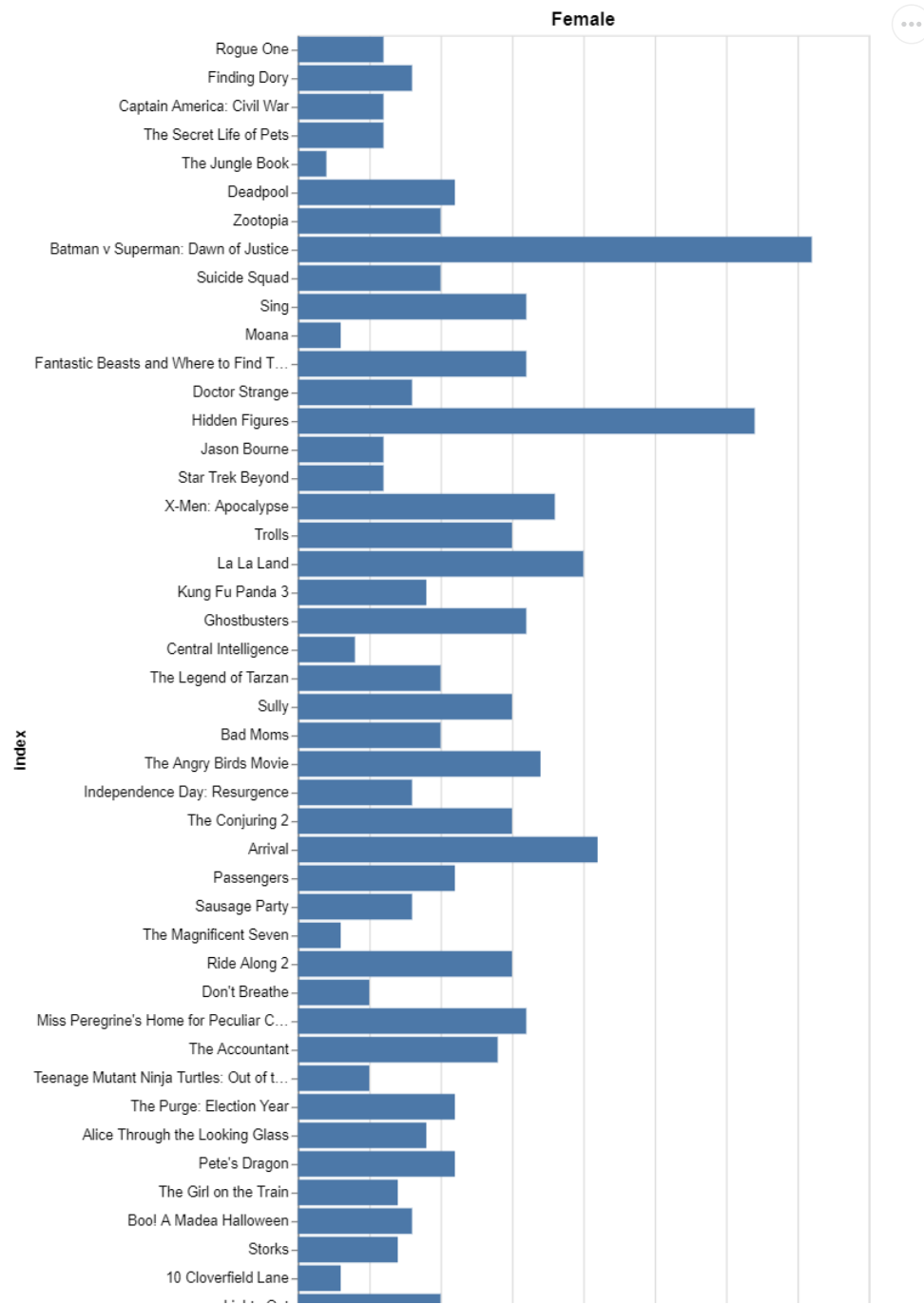
    return encoding

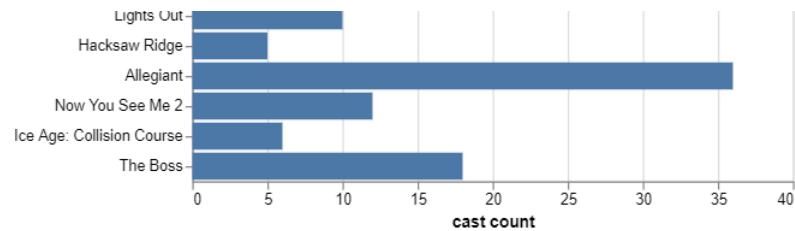
female_bar = gen_f_bar_vis(base_vis,actors_movies)

```

In [8]: female_bar

Out[8]:





How many variables are encoded in the visualization above? Add your answer below.

Answer: 2

2.2 Alternative encoding (5 points)

Complete the `gen_f_circle_vis` function. Change the encoding used in the previous example from a bar to a circle. Your visualization should look like



click [here \(assets/alt_enc_1_resized_fullImage.png\)](#) to see the full-sized image.

```

In [9]: def gen_f_circle_vis(base, indf):
        """
        input: base -- 'base' chart as defined above
        input: indf -- a dataframe like actors_movies (you can use this or 'base' to achieve the desired work)

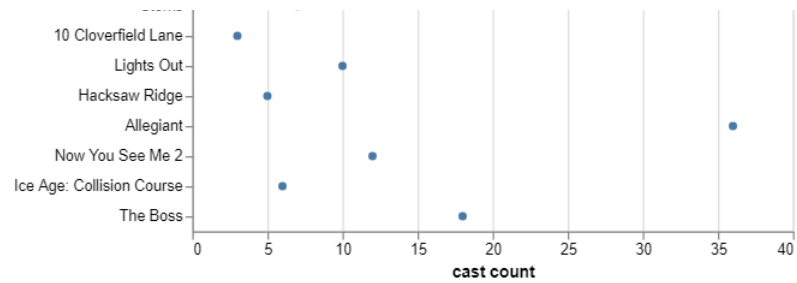
        return the call to altair function that uses the circle mark for the variables encoded in the previous example
        """
        # return encoding. (...)
        # YOUR CODE HERE

        encoding = base.transform_filter(
            alt.datum.GENDER == 'Female'
        ).encode(
            # adding all the encoding stuff
            y= alt.Y(
                'index:N',
                sort= movies_order
            ),
            x=alt.X('count(index):Q',
                    title='cast count'),
        ).mark_point(filled=True # use a "circle" symbol
        ).properties(
            # set the title
            title='Female'
        )

        return encoding

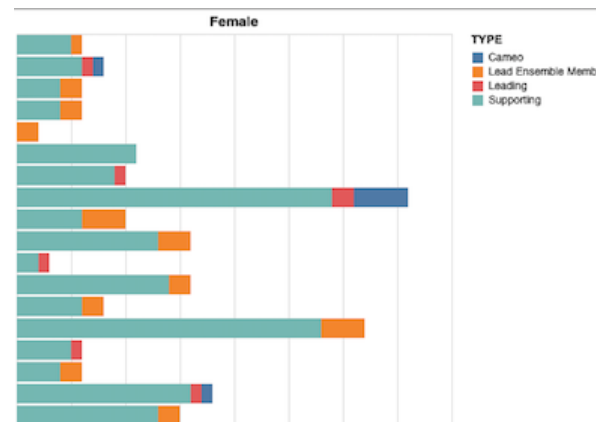
#     raise NotImplementedError()

```

2.3 Increase variables encoded (5 points)

Complete the `gen_f_stacked_vis` function. Modify the first bar chart encoding the type of the actor with the color of the bar. Your visualization should look like the following (note that we don't have labels yet):



click [here \(assets/alt_enc_2_fullSize.png\)](#) to see the full-sized image.

- Partial credit can be granted for each visualization (up to 2 points) if you provide a description of what the missing piece of the function is supposed to do without need for an Altair working version

```

In [11]: def gen_f_stacked_vis(base, indf):
    """
    input: base -- 'base' chart as defined above
    input: indf -- a dataframe like actors_movies (you can use this or 'base' to achieve the desired work)

    return the call to Altair function that uses the bar mark for the variables and the color for the TYPE
    """

    # YOUR CODE HERE

    encoding = base.transform_filter(
        alt.datum.GENDER == 'Female'
    ).encode(
        y= alt.Y(
            'index:N',
            sort= movies_order,
            axis=None
        ),
        x=alt.X('count(index):Q',
            title='cast count'),
        # add color encoding
        color=alt.Color('TYPE:N'),
        order=alt.Order('TYPE:N', sort='descending')
    )

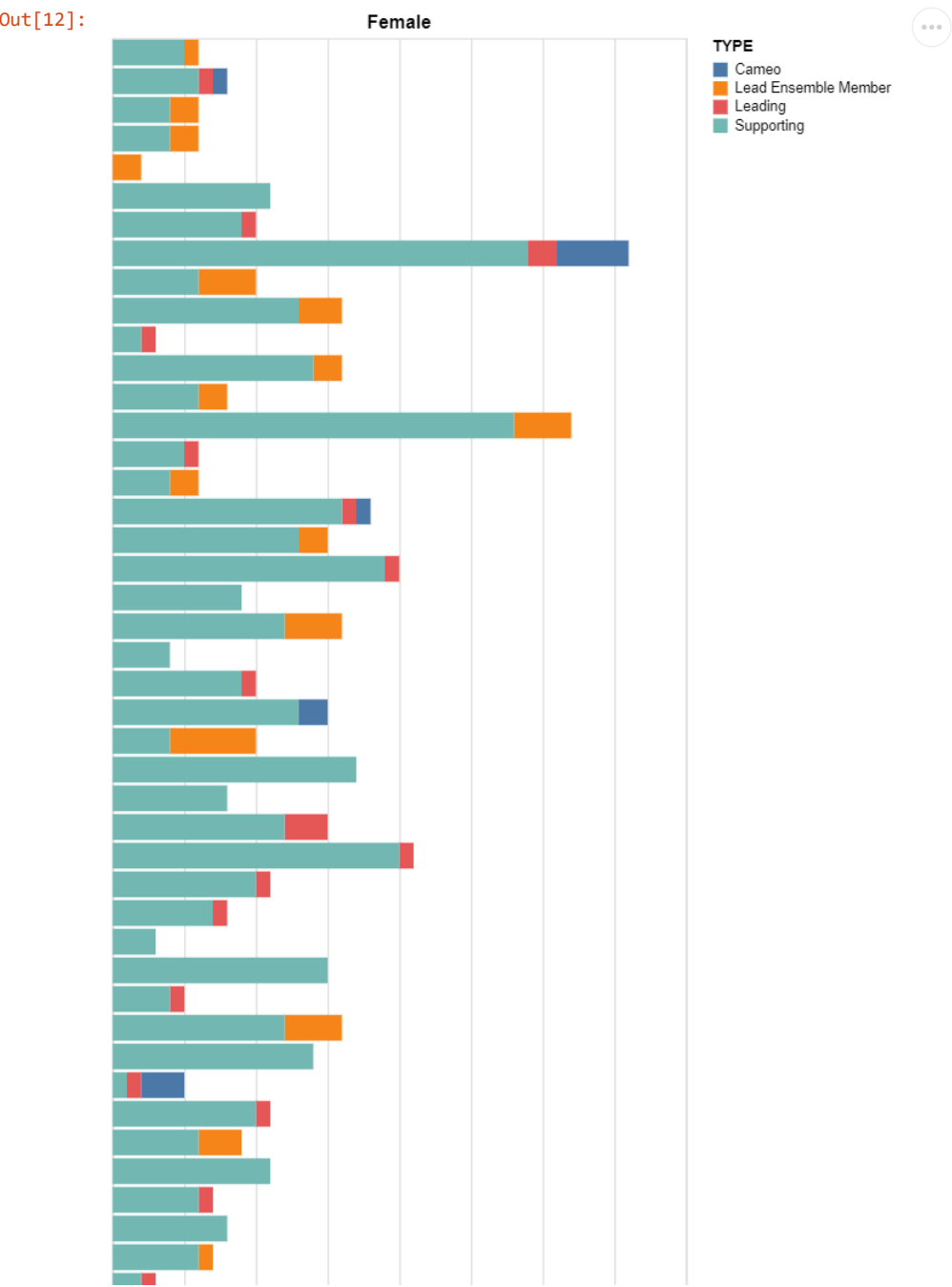
    # raise NotImplementedError()

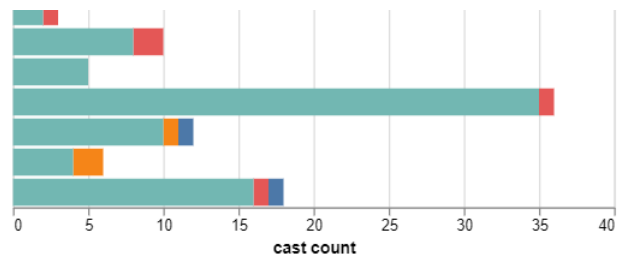
    return encoding.mark_bar().properties(title='Female')

```



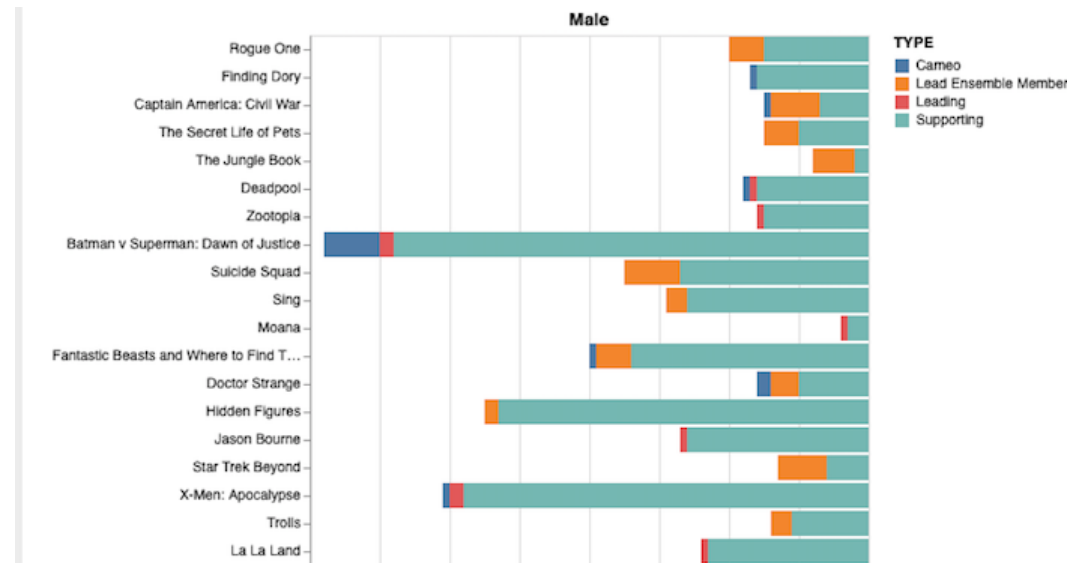
```
In [12]: female = gen_f_stacked_vis(base_vis, actors_movies)
female
```





2.4 Change filter transform (5 points)

Complete the `male_actors` function, modify the previous visualization so that the actors visualized have Male gender. Use the Altair transform function for this, not Pandas.



click [here \(assets/alt_enc_3_fullSize.png\)](#) to see the full-sized image.

- Partial credit can be granted for each visualization (up to 2 points) if you provide a description of what the missing piece of the function is supposed to do without need for an altair working version

```
In [13]: def gen_m_stacked_vis(base, indf):
        """
        input: base -- 'base' chart as defined above
        input: indf -- a dataframe like actors_movies (you can use this or 'base' to achieve the desired work)

        return the call to Altair function that uses the bar mark for the variables and the color for the TYPE
        """
        # YOUR CODE HERE

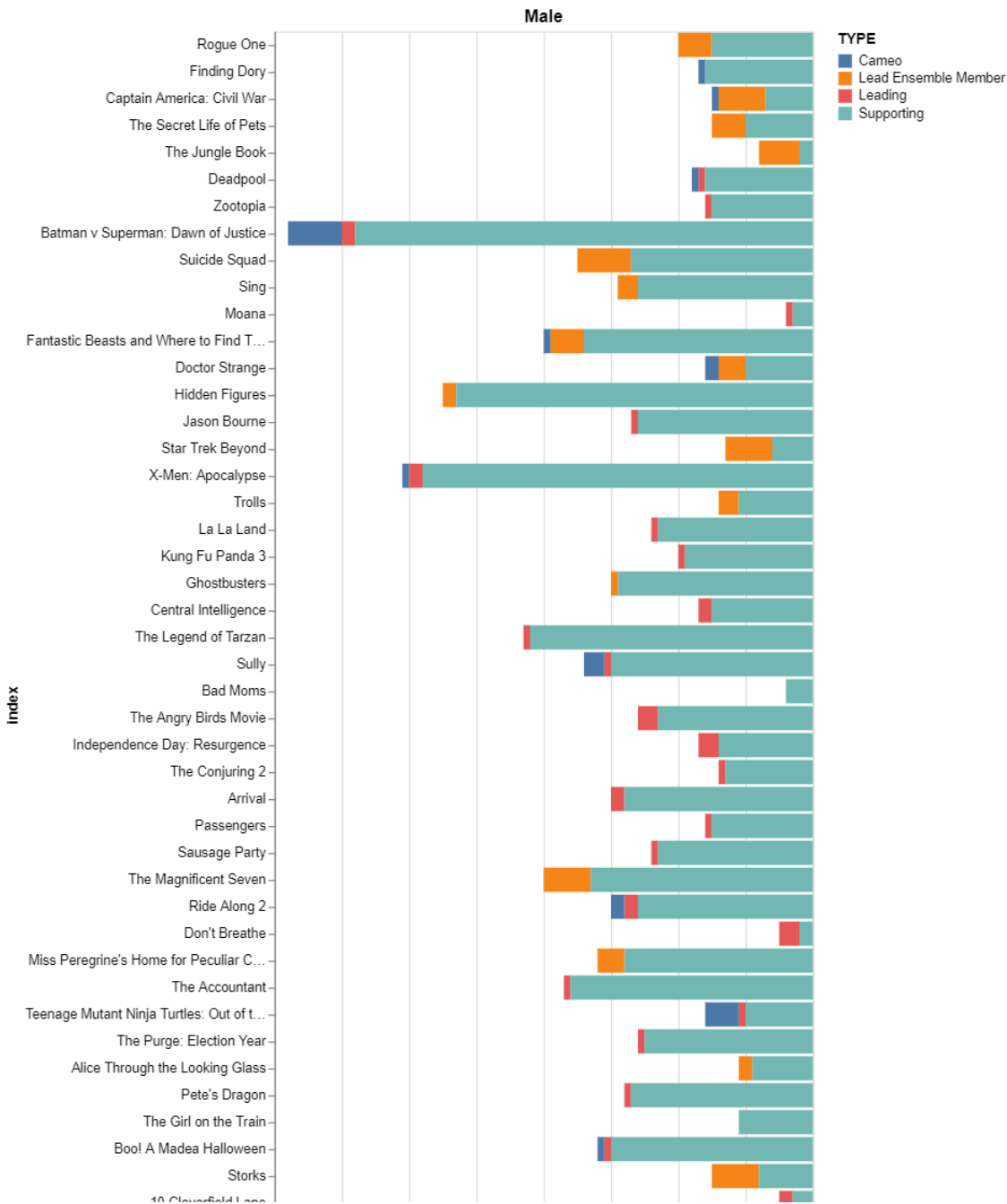
        encoding = base.transform_filter(
            alt.datum.GENDER == 'Male'
        ).encode(
            y= alt.Y(
                'index:N',
                sort= movies_order
            ),
            x=alt.X('count(index):Q',
                    sort='descending',
                    title='cast count'),
            # add color encoding
            color=alt.Color('TYPE:N'),
            order=alt.Order('TYPE:N', sort='descending')
        ).mark_bar().properties(title='Male')

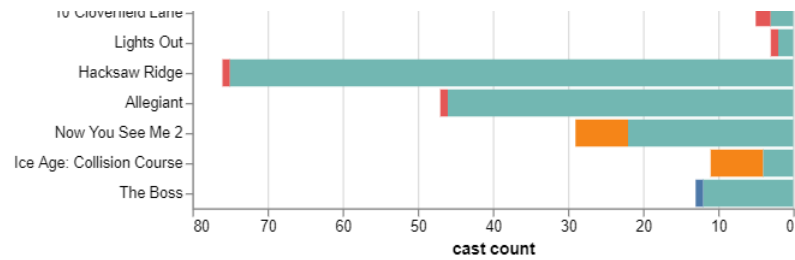
        # raise NotImplementedError()

        return encoding.mark_bar().properties(title='Male')
```

```
In [14]: male = gen_m_stacked_vis(base_vis, actors_movies)
male
```

Out[14]:





2.5 Variables encoded 2 (5 points)

Execute the following cell and determine how many variables are being encoded in the combined visualization. If you have been able to complete the previous examples, the plot should look like this



click [here \(assets/visualization_fullSize.png\)](#) to see the full-sized image.

```

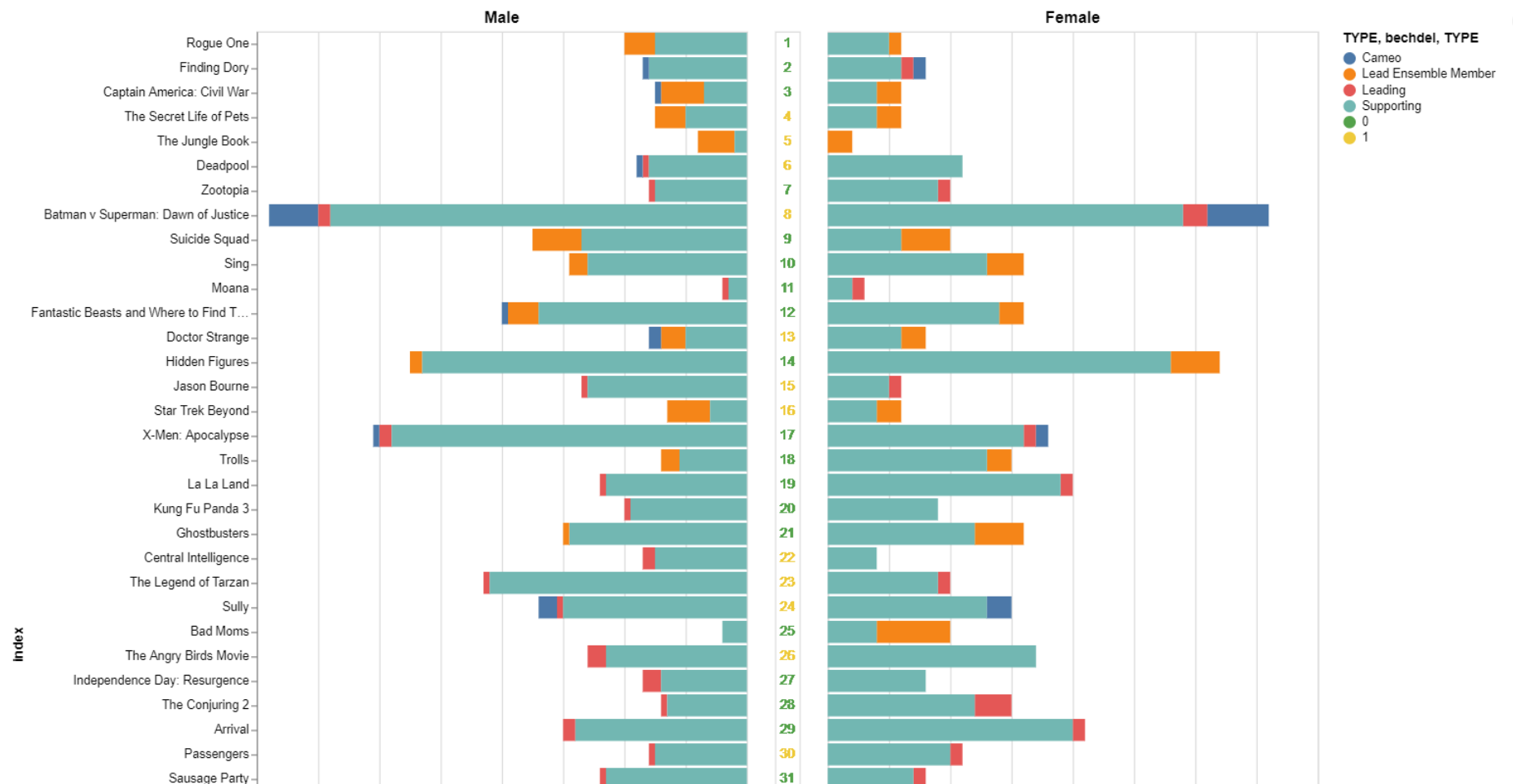
In [15]: def get_middle_vis(base, indf):
    """
    input: base -- 'base' chart as defined above
    input: indf -- a dataframe like actors_movies (you can use this or 'base' to achieve the desired work)

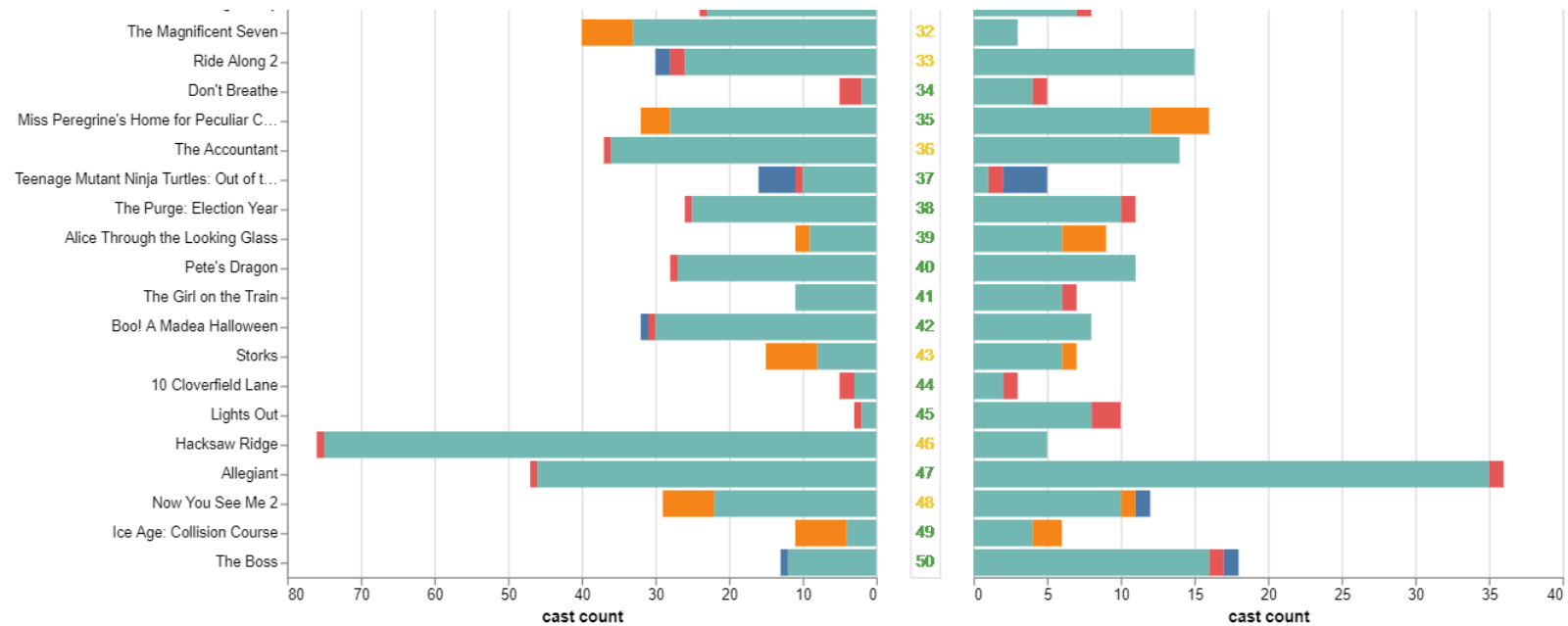
    return the "middle" column -- a text visualization
    """
    middle = base.encode(
        y=alt.Y('Rank:O', axis=None),
        text=alt.Text('Rank:Q'),
        color=alt.Color('bechdel:N')
    ).mark_text().properties(width=20)
    return(middle)

# merge together the three charts, male, middle, female
middle = get_middle_vis(base_vis, actors_movies)
male | middle | female

```

Out[15]:





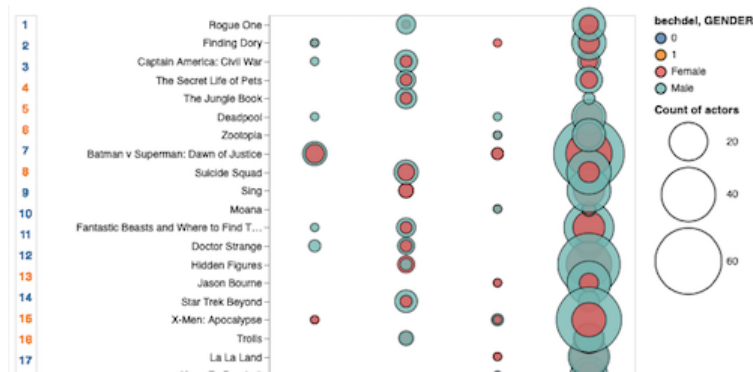
How many variables are encoded in the visualization above? (this should be an integer)

Answer: 5

2.6 Alternative encoding 1 (20 Points)

Create a new visualization within the `gen_bubble_vis` function with the following encoding:

- Use circles as the mark
- Use the scale of the circles to encode the number of actors on each category
- Use the y position of the circle to encode the movie
- Use the x position of the circle to encode the type of actor
- Use the color of the circle to encode the gender of the actor
- Match the styling of the example (it doesn't need to be pixel perfect, but should be close)
- you don't need to generate the leftmost column (it's the same as "middle" above and we'll add it at the end.



click [here \(assets/visualization_2_fullSize.png\)](#) to see the full-sized image.

- Partial credit can be granted for each visualization (up to 5 points) if you provide a description of what the missing piece of the function is supposed to do

```
In [16]: def gen_bubble_vis(base, indf):
    """
    input: base -- 'base' chart as defined above
    input: indf -- a dataframe like actors_movies (you can use this or 'base' to achieve the desired work)

    return an altair chart per the specification above
    """
    # YOUR CODE HERE

    plot = base.mark_circle(
        opacity=0.8,
        stroke='black',
        strokeWidth=0.8,
    ).encode(
        alt.X('TYPE:N'),
        alt.Y('index:N',
            sort= movies_order),
        color=alt.Color('GENDER:N'),
        size=alt.Size('count(index):Q',
            scale=alt.Scale(range=[0, 5000]),
            legend=alt.Legend(symbolFillColor='#FFFFFF')
        )
        #complete this
    ).properties(
        width=350,
        height=880
    )

    # raise NotImplementedError()
    return plot
```



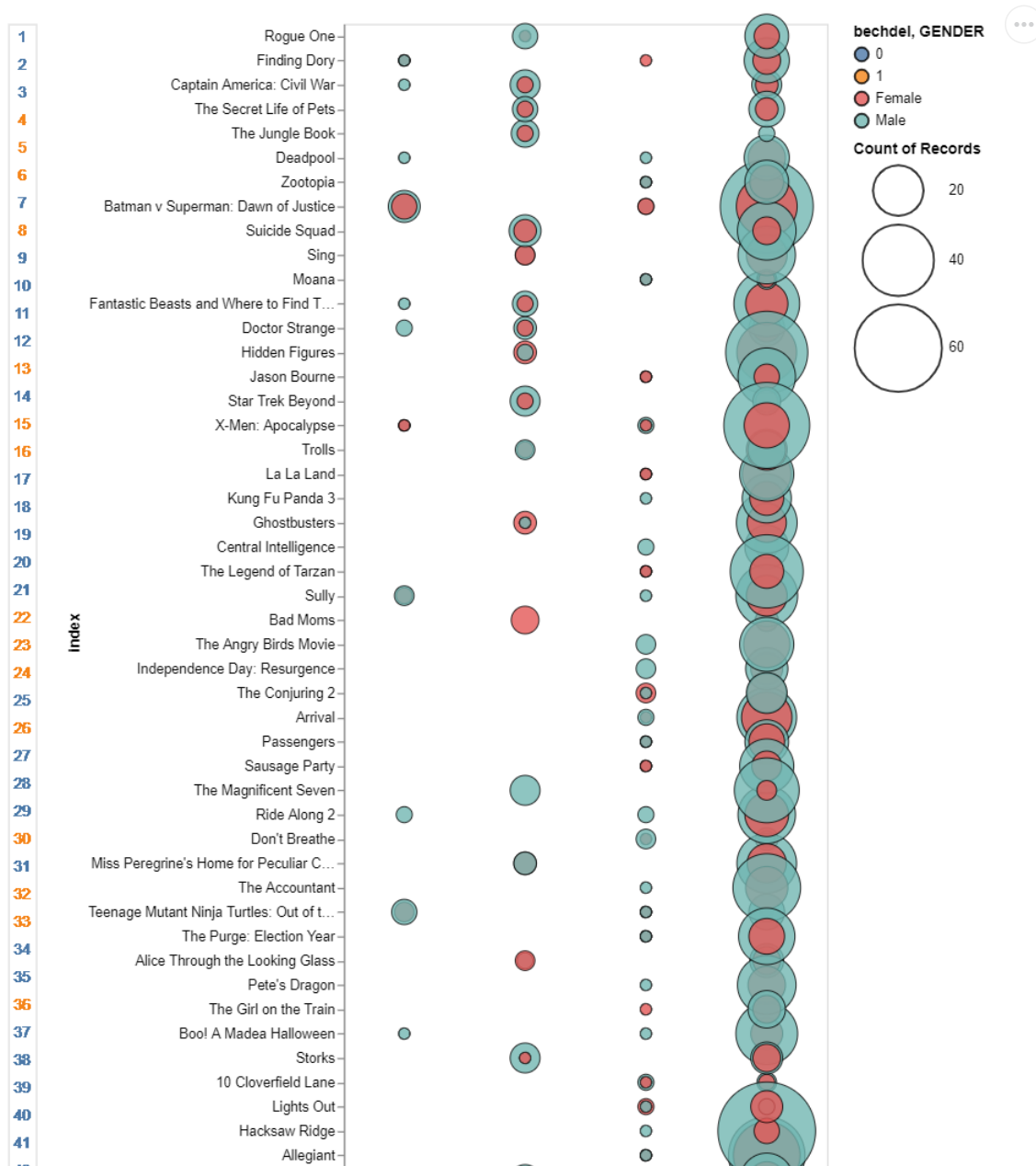
```

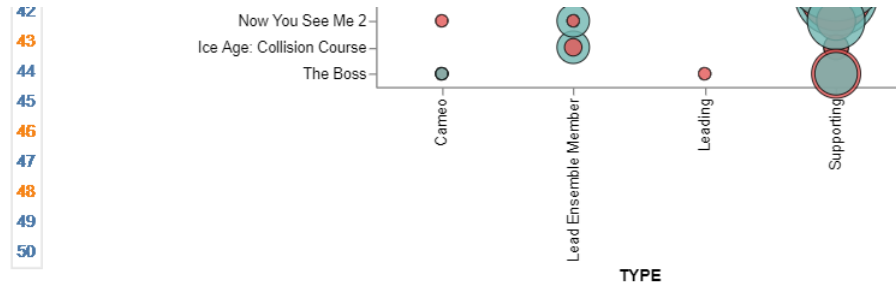
In [17]: # Let's create the bubble chart
al_enc_one = gen_bubble_vis(base_vis, actors_movies)

# add middle to the left edge and display
middle | al_enc_one

```

Out[17]:





2.7 Alternative encoding 2 (25 Points)

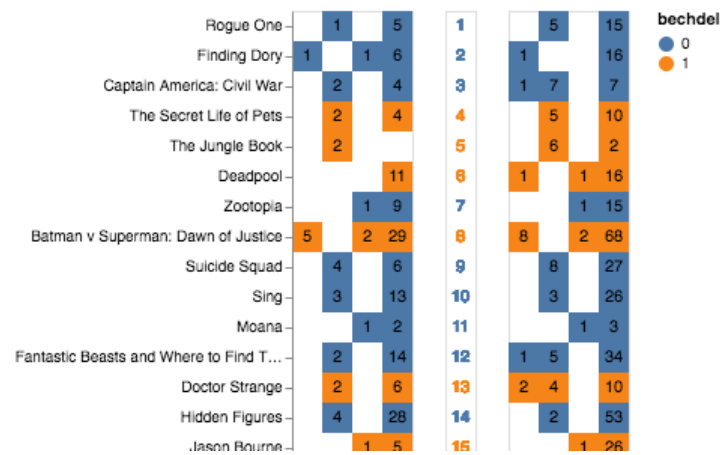
We will be completing the "heat map" style vis as specified below, but it will be easier to create a male and female subchart in two separate functions (the code will get long otherwise) and then put them together. Complete `gen_f_hm_vis()` and `gen_m_hm_vis()` functions to create a new visualization with the following encoding:

- The left and right plot should filter male and female actors respectively
- Use rectangles as the mark
- Use the text inside each rectangle to encode the count of actors one each category (gender, type and movie)
- Use the y position of the rectangle to encode the movie
- Use the x position of the rectangle to encode the type of actor
- Use the color of the rectangle to encode whether that movie passes the Bechdel test or not (bechdel variable)
- Note that only the female vis has labels on the left, the male vis does not

The top of the female plot would look like this (click [here \(assets/visualization_3_fullSize.png\)](#) to see the full-sized image):



If you've done everything correctly, your final visualization should look like the one below (click [here \(assets/visualization_3_full.png\)](#) for the full plot).



- Partial credit can be granted for each visualization (up to 4 points for each function) if you provide a description of what the missing piece of the function is supposed to do without need for an Altair working version

```

In [18]: def gen_f_hm_vis(base, indf):
    """
    input: base -- 'base' chart as defined above
    input: indf -- a dataframe like actors_movies (you can use this or 'base' to achieve the desired work)

    return an altair chart per the specification above
    """
    # modify to add filter transform
    plot = base.transform_filter(
        alt.datum.GENDER == 'Female'
    ).mark_rect().encode(
        alt.X('TYPE:N'),
        alt.Y('index:N',
            sort= movies_order),
        color=alt.Color('bechdel:N')
    )

    #modify to add filter transform
    text = base.transform_filter(
        alt.datum.GENDER == 'Female'
    ).mark_text(baseline='middle').encode(
        x='TYPE:O',
        y= alt.Y(
            'index:O',
            sort= movies_order,
            axis=None
        ),
        #change this
        text='count(GENDER):Q'
    )
    # YOUR CODE HERE
    # raise NotImplementedError()

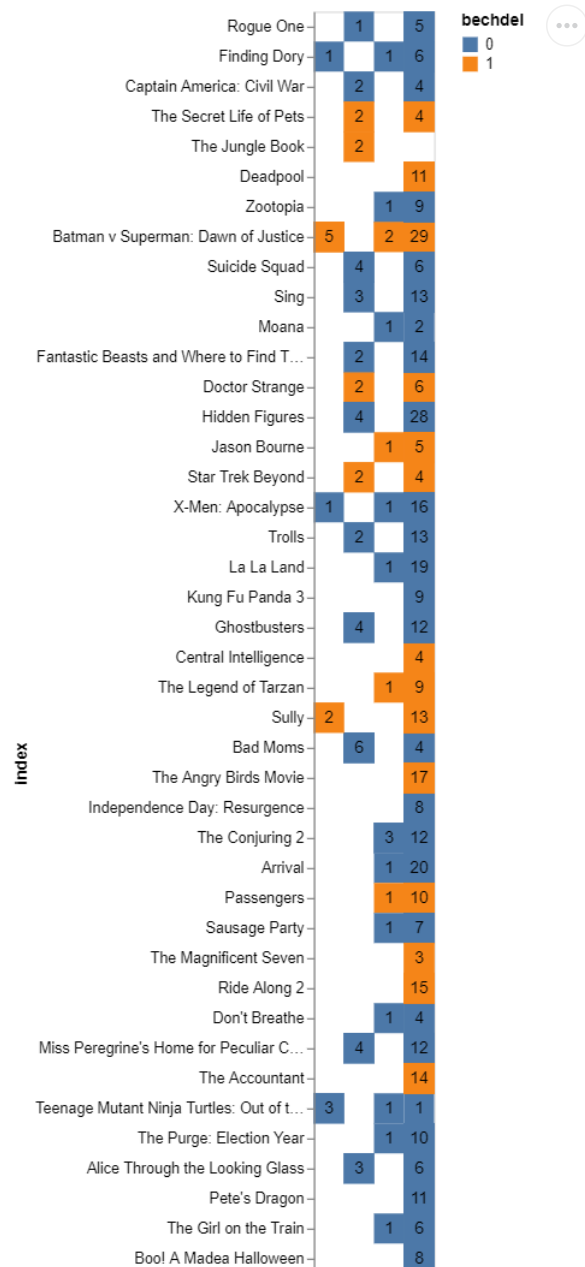
    return plot + text

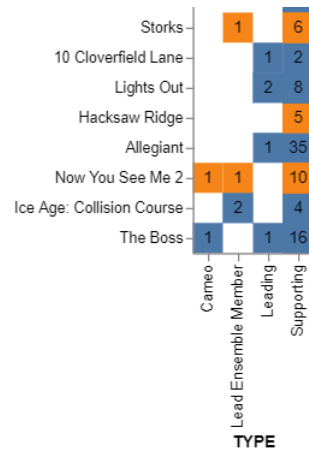
```

```
In [19]: # create the female side
f_a_1 = gen_f_hm_vis(base_vis, actors_movies)

# display it to check
f_a_1
```

Out[19]:





```
In [20]: def gen_m_hm_vis(base, indf):
    """
    input: base -- 'base' chart as defined above
    input: indf -- a dataframe like actors_movies (you can use this or 'base' to achieve the desired work)

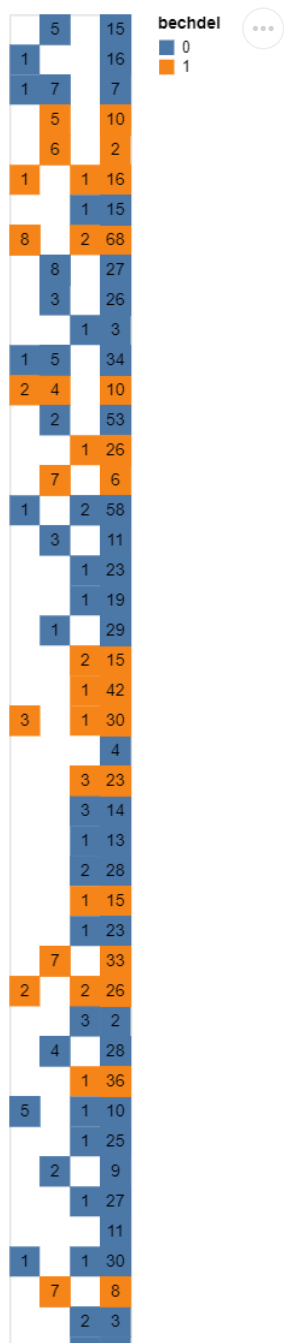
    return an altair chart per the specification above
    """
    # modify to add filter transform
    plot = base.transform_filter(
        alt.datum.GENDER == 'Male'
    ).mark_rect().encode(
        alt.X('TYPE:N'),
        alt.Y('index:N',
            sort= movies_order),
        color=alt.Color('bechdel:N')
    )

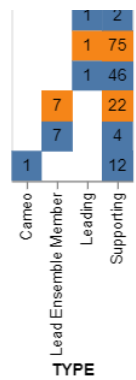
    #modify to add filter transform
    text = base.transform_filter(
        alt.datum.GENDER == 'Male'
    ).mark_text(baseline='middle').encode(
        x='TYPE:O',
        y= alt.Y(
            'index:O',
            sort= movies_order,
            axis=None
        ),
        #change this
        text='count(GENDER):Q'
    )
    # YOUR CODE HERE
    # raise NotImplementedError()

    return plot + text
```

```
In [21]: m_a_1 = gen_m_hm_vis(base_vis, actors_movies)
         m_a_1
```

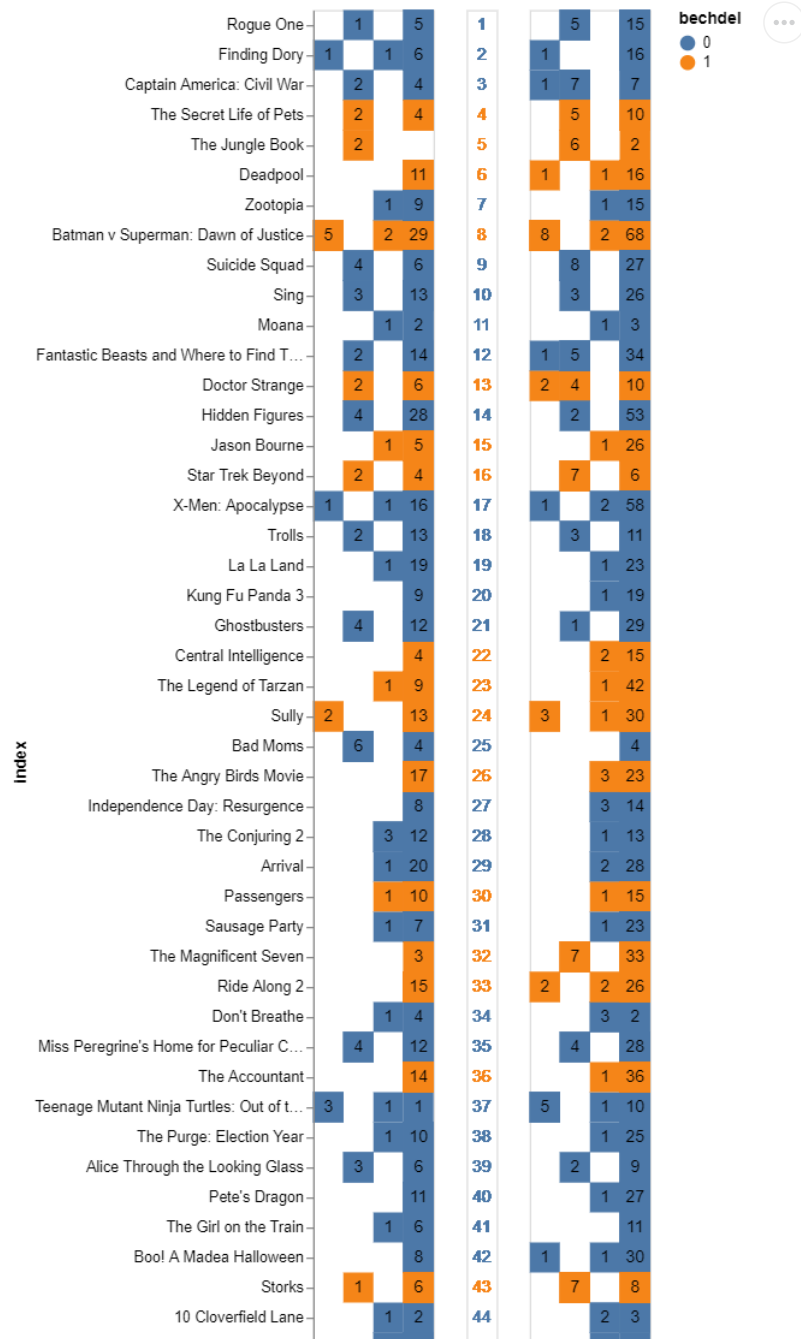
Out[21]:

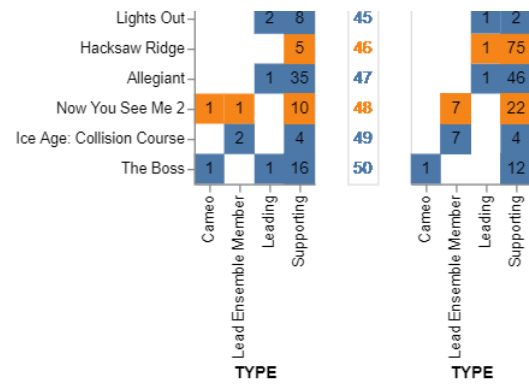





```
In [22]: # create the visualization
f_a_1 | middle | m_a_1
```

Out[22]:





2.8 (Bonus) Compare effectiveness (5 points)

Look at the visualization for question 2.7. How does this visualization compare in terms of effectiveness to the visualizations in questions 2.5 and 2.6?

2.8 Answer