# Milestone II: Implementation and Design Choices Explanation

## Analyzing Trends in Library Usage Across the Five Most Frequented San Francisco County Libraries

Group: Data Dudes

Participants: Heidi Lantz, Sadia Khan Durani, Lillian Milroy

https://www.kaggle.com/datasets/datasf/sf-library-usage-data

This project aims to analyze trends in library patron behavior using a dataset containing approximately 420,000 library patrons from the San Francisco Public Library. We intend to identify trends in how various age demographics interact with media, and how these trends may affect other aspects of library usage. Our analysis will provide insight into how different kinds of patrons engage with the library system, and how to optimize their experience going forward.

Our intended audience for this project includes library administrators, researchers, and policymakers who are interested in understanding and improving library services.

The project seeks to provide valuable insights into library patron behavior to help library administrators and stakeholders make data-driven decisions. By understanding patron preferences, usage patterns across different months, and demographics, libraries can tailor their services to better serve their communities and optimize resource allocation.

Our overarching goal is to analyze and reflect on what factors influence user engagement in the top five San Francisco library branches from 2005 to 2015.

# Cohesion

**Task 1:** *How does the duration of library membership vary by patrons' notice preference type?*

In addressing this task, we used the data attributes: Year Patron Registered, Last Circulation Year, Age Range, and Notice Preference Definition. This visualization supports the overall goal of analyzing the factors which influence user engagement, providing insight into the number of years a patron from a certain age group is active with the library for. This visualization allows the user to analyze the duration of individual patron activity over the top five libraries and draw comparisons about the activity lengths across age ranges. In addition, the user gains insight into what notice preference an age group is most likely to opt for at the time of registration. This visualization benefits the audience since it allows library administrators and researchers who work to improve library engagement to become conscious of which age groups tend to remain active for shorter or longer periods. If library administrators and researchers learn who the most active members are across all the library branches, they can move forward to make improvements to increase engagement across other age groups as well. Furthermore, they can reason as to why some age groups were more active during certain years from 2005 to 2015 depending on the types of services they provided in their libraries.

**Task 2:** *What is the variation in total checkouts across the years, and how do corresponding renewals per checkout differ across various age groups?*

This task used the following data attributes: Year Patron Registered, Sum of Total Checkouts, Renewals per Checkouts, and Age Range. The visualization is relevant to our overarching theme as it allows the user to learn which year from 2005 to 2015 had the highest total checkouts, and additionally learn how much the different age groups were involved in those trends. This visualization allows the user to learn about the overall engagement with San Francisco library branches and see which years caused an increase or decrease in total checkouts over the 10 year period. For example, we found there was a significant dip in total checkouts amount in 2007, but an increase shortly after. This prompts researchers to recall major events that might have led to a decrease in patron activity with the libraries. We will further analyze how we could infer this could be related to the 2007-2008 financial crash. The interactive dot chart implemented with the area chart visualizes the corresponding renewals per checkout for each patron. It is of benefit to the audience as it lends insight for which groups are likely to have more renewals per checkout and also facilitates comparisons of renewals between age groups. For example, when considering the lowest recorded value of renewals per checkouts (occurring in 2007), the primary age group engaging with libraries was 25-34 year old patrons, with the highest number of renewals per checkout being only 2.15 per patron. Moreover, with the peak for total checkouts occurring in 2010, a corresponding patron from the 35-44 age group had an average of 6 renewals per checkout. This major change that we observe from the visualization showcases the importance and benefit of creating this task.

**Task 3:** *Which months influence the engagement level of patrons across the top 5 libraries?*

The data attributes: Circulation Active Month, Home Library Definition, and Average Total Checkouts were utilized in this task. This visualization contributes to our main theme by recognizing the impact of the month on patrons' engagement levels with the library. For instance, our findings indicate that January sees a higher average total checkout compared to other months. One hypothesis is that individuals, driven by New Year's resolutions, express a larger motivation for reading during that time of year but eventually give them up. Additionally, the winter months might have increased library use as people spend more time indoors and choose to entertain themselves by reading. You could also make  inferences including if the home library location was more children-oriented then that contributes to the increased level of

checkout activity over the summer break. The audience benefits from this insight as it enables library administrators to better understand their visitors, including how frequently they visit the library and potential reasons motivating these patterns. For instance, if January sees increased visits due to resolutions, planning related events in February could help to sustain engagement. Recognizing these trends allows libraries to make informed changes that improve overall library usage. This information also tells us what the audience can learn about the topic. It emphasizes that the time of the year could influence people's motivation to engage with the library, and that this could differ per home library. Understanding these seasonal variations allows for the development and implementation of personalized strategies to maintain and boost library engagement throughout the year.

*Task 4:* *Which home libraries have the most average total checkouts over the number of years a patron is active in the library?*

This task used the attributes Years Active, Average Total Checkouts, and Home Library Definition. This visualization fits to our overall theme by comparing the differences in engagement over time between home libraries. It helps us inquire as to why some libraries might attract more engagement, whether it's due to events hosted, better book selection, staff per location, or technological amenities. Understanding these factors can guide us in making improvements to boost library engagement and identify factors that might be discouraging people from staying connected. What's interesting is that the visualization highlights trends, primarily Chinatown's library seeing increased usage over time, while the Main library shows a decrease. The other libraries appeared to have more constant overall checkouts across their years active.

This benefits the audience because it prompts us to explore why these patterns exist and try to find practical insights for improvement. For example, it encourages libraries to rethink their events, activities, and experiences to keep people engaged in the long term. The value here is that it benefits both libraries and users. Libraries can learn from one another and improve their services, whereas users are more likely to stay engaged, creating a stronger connection with the library and its activities. The audience, being both people who manage libraries but also use them, can learn how to effectively engage users in a library and also what incentivizes you to use them. Overall, this visualization showcases that each library has unique factors influencing engagement. It emphasizes that success isn't just about the patrons but also about the library's initiatives and environment. If library administrators are aiming to boost involvement in library activities, the key is to look at the successful libraries and adopt strategies that work for their audience, promoting a more enriching experience for library users.

*Interactive Dashboard*

The visualizations in our dashboard that are linked through interaction include the area and dot chart from task two, along with the heatmap from task three. The linking is done through uni and  bi-directional interactions. Specifically, the area chart controls both the dot chart and heat map, while the dot chart and heat map show bidirectional interactions with each other. Through these interactions, the user is able to learn more specifics about the factors affecting patron engagement in the library branches from 2005 to 2015. For example, if a user selects over the years 2005 to 2007 on the area chart, both the dot chart and the heatmap would display the corresponding renewals per checkout across the age ranges and the monthly mean checkouts per library for those years, respectively. This enables users to analyze deeper the difference in engagement across the five library branches throughout the years and specifically, identify the months in which patrons of a certain age range were being more active.
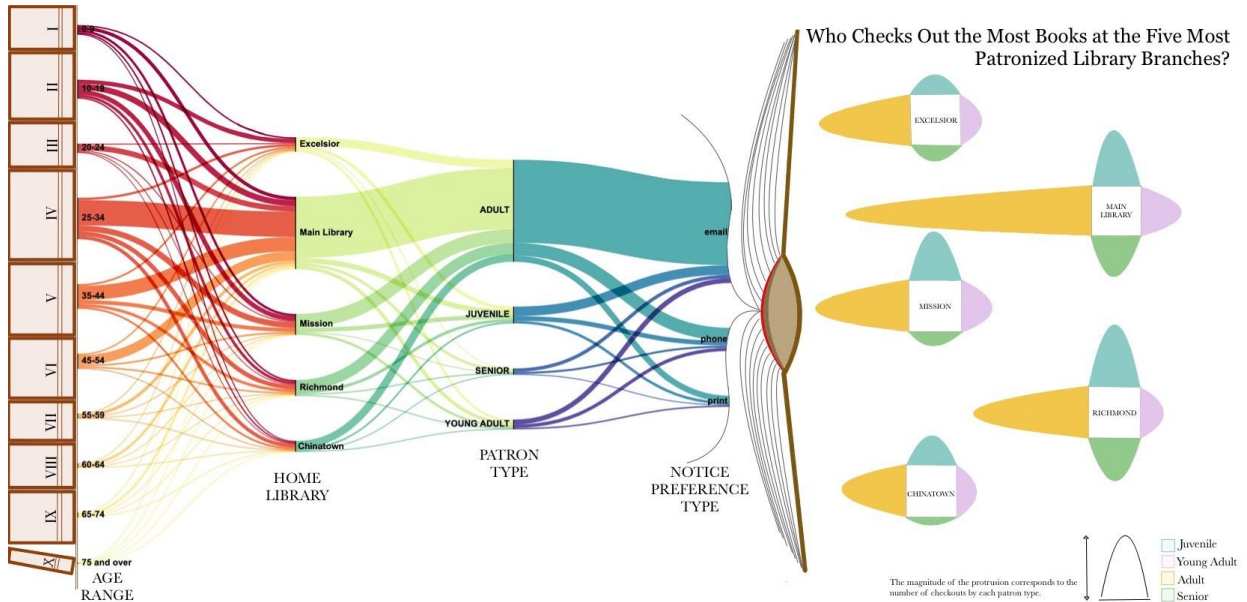
# Cohesion

## Interactive Dashboard



**Analyzing Trends in Library Usage Across the Five Most Frequented San Francisco County Libraries**
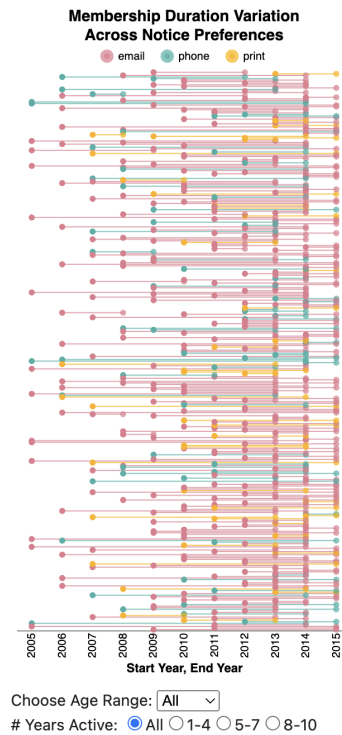
Task 1: Membership Duration Variation Across Notice Preferences

Task 2: Library Engagement Over Time and Across Age Groups

Task 3: Monthly Mean Total Checkouts Per Library

Task 4: Checkouts Over Time Per Library

---

## Novel Visualization



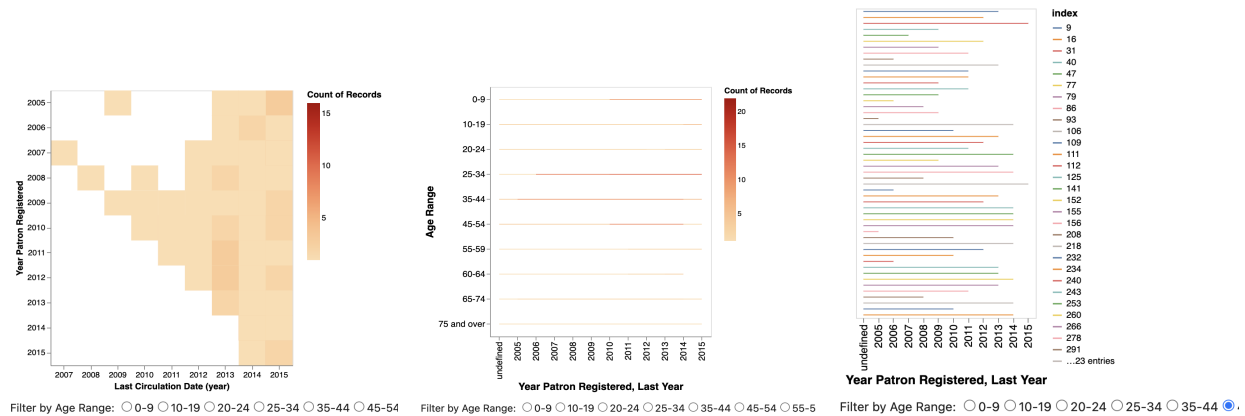Who Checks Out the Most Books at the Five Most Patronized Library Branches?

# Justifications for Each Visualization

**Visualization 1**

Associated Task: *How does the duration of library membership vary by patrons' notice preference type?*



Previous Iterations:



This visualization utilizes the following data attributes: Year Patron Registered, Last Circulation Year, Age Rang and Notice Preference Definition.

The visualization was created using a line mark, and encodes two attributes on the x channel. The first attribute encoded is Year Patron Registered (temporal) and the second is the Last Circulation Year (temporal). The y channel encodes the index number, therefore each individual line across the vertical layout of the visualization represents a patron. The end points for each line segment were created by layering an additional chart over top, encoded the same way but instead using the circle mark. The tooltip

channel also encodes Year Patron Registered and the Last Circulation Year to ensure users can easily determine the years by hovering over a line segment. Lastly, the color channel encodes Notice Preference Definition. This visualization uses the interaction type of selection through the implementation of 2 UI widgets. The first UI widget allows the user to select one of the 10 possible age ranges using a drop down menu, as the cardinality of the Age Range attribute is 10. The second UI widget is the radio buttons which allow the user to select one of the three options for the number of years active. This is encoded with radio buttons since there are only 4 options to select from.

Through these interactions, the user is able to answer the following questions as this visualization's functionality includes filtering based on two attributes. These questions would be difficult to answer through a static version as it would be a very cluttered visualization with many line segments:

- What age groups are active for longer time periods with the library?
- What age group are only engaged with the library for 1-4 years?
- What notice preference method is used more when patrons are active for shorter periods?
- How many patrons are registered for 8-10 years?

We chose to use a line chart to visualize the duration of individual patron activity with the library branches because we were working with two ordered attributes, Year Patron Registered and Last Circulation Year, which we were able to encode both on the x-axis. These attributes correctly use point marks to connect the line segments created. However, line charts also plots a quantitative attribute using spatial position; but in our line chart, this doesn't apply. On the y-axis channel we have the index number encoded so that each line can represent a single patron activity with the library branches. This may violate the Effectiveness Principle as the visualization could give a sense of ordering for the line segments as they are placed vertically. Ordering was not our goal, since index number is not an ordinal attribute part of our data. With the use of lines connecting the respective start and end years for each patron, the channels can show perceptual grouping through connection, making it easier for users to understand the implications of the line segments and connections between the data points. Because the Notice Preference Definition attribute has a cardinality of three, encoding it on the colour channel was the most effective choice. It is a categorical attribute with a small cardinality - hence, its levels are easily discriminable. The visualization adheres to the Expressiveness Principle, as everything encoded expresses the information included in the attributes.
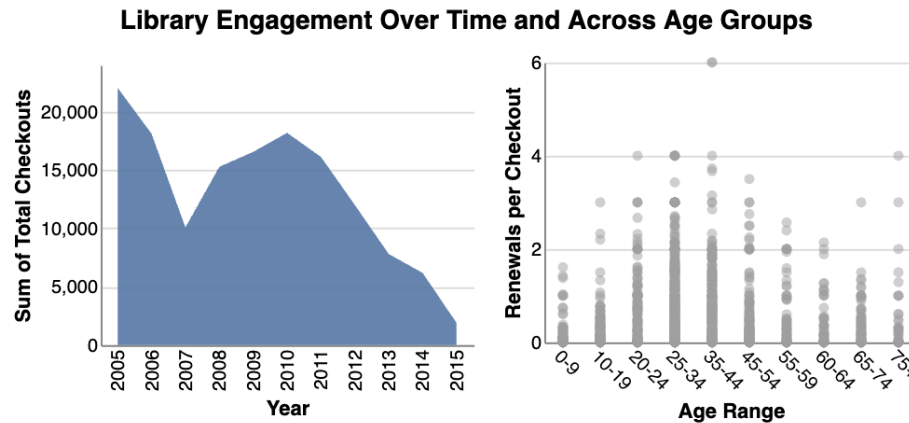
One disadvantage to the visualization is that it appears cluttered at first glance, because there are a lot of line segments due to the dataset consisting of many patrons. However, as a user can interact with the visualization and utilize its functionality of filtering based on Age Range and Years Active, the visualization becomes more understandable.

We opted to include this visualization because compared to other chart types we considered, this option best represented our unique data attributes, while also complementing their semantics. The heatmap also encodes attributes on the axes as categorical data and one quantitative attribute, but we believe that lines are better fit to represent time intervals.
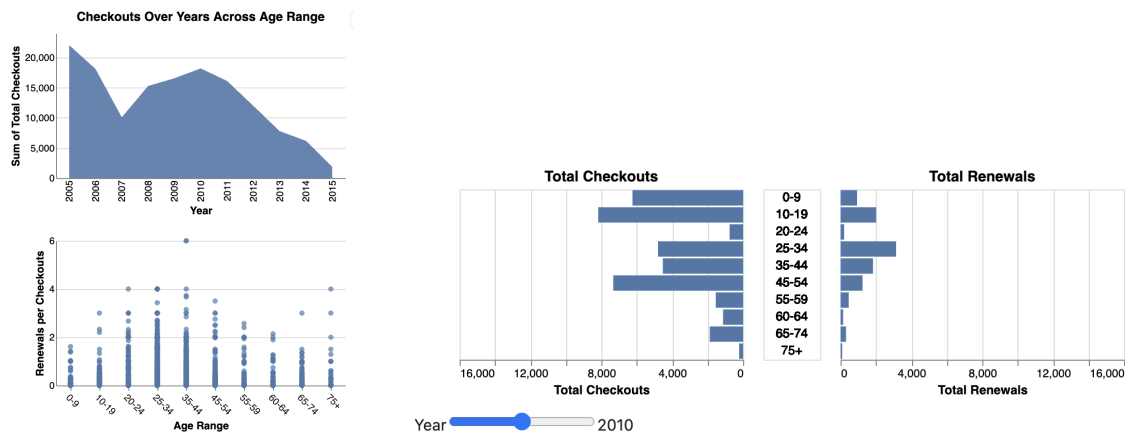
# Justifications for Each Visualization

## Visualization 2

Associated Task: *What is the variation in total checkouts across the years, and how do corresponding renewals per checkout differ across various age groups?*



Previous Iterations:



The task utilizes the following data attributes: Year Patron Registered, Sum of Total Checkouts, Renewals per Checkouts and Age Range.

This visualization consists of two plots. The first is an area chart (using the area mark), with Year (from Year Patron Registered) encoded on the x channel, and sum of total checkouts on the y channel. The second is a dot chart using the circle mark, with Age Range encoded on the x channel and Renewals per Checkouts encoded in the y channel. The area chart is linked unidirectionally to both the dot chart and the heatmap addressing task 3, while the dot chart and heatmap are linked bidirectionally to one another. The area chart essentially serves as a UI widget, facilitating the implemented selection interval, allowing the user to choose the time range from which data will be displayed in the dot chart and heatmap.

The incorporation of the selection interval permits users to explore data relevant to their time period of interest. Linking the area chart to both the dot chart and the heatmap enables users to conceptualize how each age range in the dot chart contributes to the patron activity described in the heatmap. The interactive

nature of the visualization facilitates the exploration of the following lines of inquiry, that a static visualization would not:

- Are there specific time periods (months or years) where patron activity spikes or drops? Which age range(s) are responsible for the majority of patron activity at these times?
- Are there months of the year that consistently exhibit higher or lower mean checkouts? Is there variation in the age ranges that contribute most to patron activity at these times?
- Are there any periods where checkouts are high, but renewals per checkouts are low (or vice versa)?

We chose to use an area chart because it is well-suited for communicating the progression of a quantitative variable (in this case, sum of total checkouts) over time. The simplicity offered by the area chart complemented its primary functionality as a UI widget. While theoretically a line chart could have communicated the same information, we opted to use the filled area chart for better visual clarity of the aggregated y-axis values to align with the effectiveness principle, and as well to provide contrast to make the brush more visible with the expressiveness principle. We implemented a dot chart because it represents patron behavior at an individual level as a discrete value, providing contrast to the area chart and heatmap. It is a simple and effective way to show the distribution of a quantitative variable (Renewals per Checkouts) across levels of an attribute (Age Ranges) in a population. The dot chart strikes an appropriate balance between expressiveness and effectiveness, utilizing position to represent both the discrete nature of the data, and the magnitude of renewals per checkout. Dot charts are uniquely able to capitalize on popout, clearly displaying the outlier values. Both charts demonstrate good separability; they are enhanced by the interactions implemented, but still communicate information sufficiently and interpretably at an individual level.

Aspects of the dot chart that welcome further development are accuracy and discriminability. These channel characteristics are obfuscated because of the aggregation of the circles. The circles play an important role by representing discrete values, but alternative visualization methods that can optimize channel characteristics without sacrificing others should be investigated.
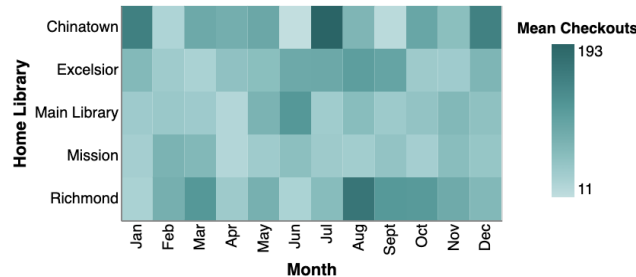
We chose to align these two charts with a single task to provide a more comprehensive, autonomous experience for the user, showing them the data they wanted to see in a way that facilitated ease of comparisons and highlighted trends. We believe it is particularly useful for visualizing how past events (social, political, or otherwise) influenced overall library usage across age ranges, what types of engagement were impacted (checkouts or renewals), and how these trends varied across libraries (when considering task two with the heatmap in task three). An example of this utility is the 2008 financial crash, which occurred at the end of September. Renewals per checkouts and the sum of total checkouts during this period plummeted as library usage stalled. Mean checkouts at most branches also dipped, with the exception of Excelsior library, which showed an increase in mean checkouts in the months leading up to the crash, hitting its peak (and one of the highest mean checkout values recorded) in September. Excelsior is a middle-class neighborhood containing the highest proportion of food service businesses in San Francisco, one of the first industries to feel the effects of the impending crash. These charts visualize how reliance on library services skyrocketed in a neighborhood uniquely impacted by the crash.
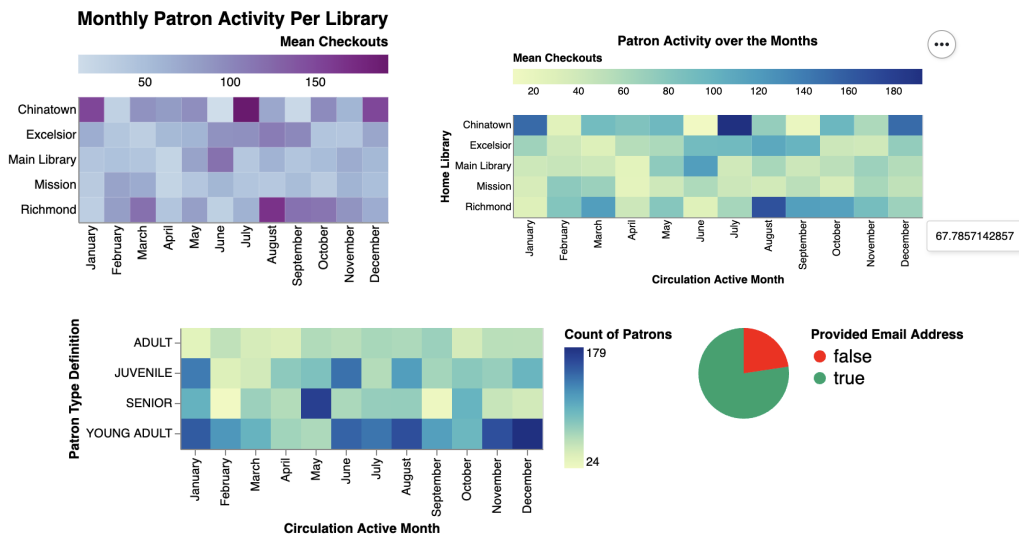
## Visualization 3

Associated Task: *Which months influence the engagement level of patrons across the top 5 libraries?*



Task 3: Monthly Mean Total Checkouts Per Library

Previous Iterations:



This task used the following data attributes: Circulation Active Month, Home Library Definition, and Average Total Checkouts.

The visualization employs the rect mark to construct a heatmap showcasing usage across months and libraries. It uses the channels of color representing the quantitative attribute mean checkouts, x being the ordinal attribute circulation active month, and y as the nominal attribute of the top five active home libraries. The tooltip channel includes the mean checkouts amount for extra clarity. The visualization incorporates an interaction feature using a selection point with the month, establishing a bi-directional relationship with task two. The selection point allows you to view the information of chosen months in the heatmap. This tells us it is an implicit interaction that has a direct focus, since you can click directly on the heatmap to select the given month. This feature allows observation regarding the factors influencing monthly library activity, and can offer insights into why some months exhibit higher usage than others. By having that bi-directional relationship, we are able to gain additional insights into how different age groups impact monthly library use. This is what the interaction makes possible that a static

version wouldn't have, by providing supplemental information and being able to analyze more relationships within the data. The reason this is important is because with a static version, it is easy to have the data blend together since a heatmap contains so much information. By being able to select certain months and choose what data you want to analyze, it makes the visualization simpler and easier to understand. We can itemize this information into the following questions:

- Which months, when selected, reveal distinct patterns in library activity, and how do these patterns contribute to the overall understanding of monthly usage?
- Are there identifiable relationships between selected months and specific age groups that significantly impact library use, and how do these relationships contribute to the observed patterns?
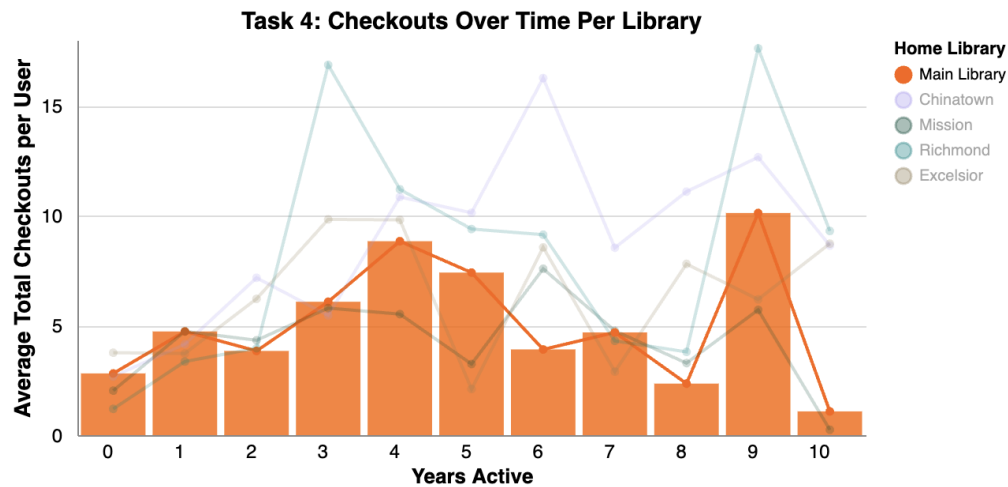- Are there any months with similar activity across the 5 libraries?

The channels in this visualization are fairly well-matched, making this a generally effective viz. It could use some improvement however, since color isn't always the best way to represent a quantitative variable with a large range of values. This ties in with the color hue rule, where having this many colors makes for a hard comparison. For example, it is difficult to differentiate between the Chinatown library in February versus the Mission library in October for finding which one has a higher mean checkout value. This is what makes the interaction useful, so that we can compare specific months with one another in a simplistic way. This visualization is very expressive, as the data is accurately portrayed without any unwanted inferences. The visualization excels in separability, with no visual elements interfering with each other. The interaction feature introduces a pop-out effect, directing attention to libraries with notable mean checkouts for a particular month, and the color emphasizes any standout high values. Overall, this visualization is highly expressive and showcases good separability, but is not the most effective visualization due to the wide range of color.

One drawback about this visualization is that it doesn't do a great job of highlighting trends in the data. Trends comparing 12 different months can be tough, and comparing colors on that much information is not a simple way of accomplishing this task. As well, the data does not differ that strongly between months and that could be due to the smaller sample size we chose. Even though we can notice slight differences between monthly trends, it does make it difficult to notice when you are also looking at 5 different libraries. However, it is still a good way to represent this data and is visually appealing for the user. We chose this because we knew it was a unique visualization, since we could analyze trends for the months of the year. The heatmap allows us to make a large-scale comparison that not many other visualizations would be able to accomplish. It provides valuable reflection that connects us to the main goal; library usage can be affected by different times of the year.
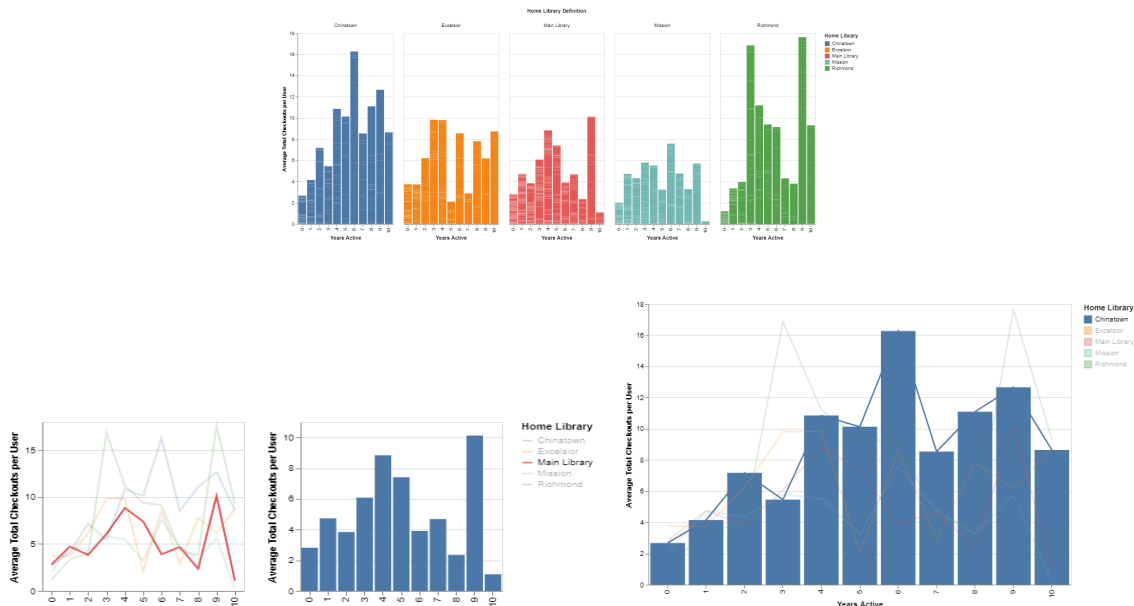
# Justifications for Each Visualization

**Visualization 4**

Associated Task: *Which home libraries have the most average total checkouts over the time a patron is active in the library?*



Previous Iterations:



This task used the following attributes: Years Active, Average Total Checkouts, and Home Library Definition.

The visualization combines the mark bar and mark line, overlapping them to let us present the data effectively. We use the mark line because it is a good way to represent information over a period of time. The mark bar is for added information and readability if you are interested in a particular library. It uses the color channel to represent the nominal attribute of different home libraries, 'x' for the quantitative years active in the library, and 'y' for quantitative average total checkouts per user. We also put an additional tooltip channel that indicates the average total checkouts per user amount, and also the home

library and years active for added clarity. It has an interaction via point selection in the legend, and by clicking, it reveals a bar chart distribution of average total checkouts for the library you selected.This characterizes the interaction as implicit, since it isn't necessarily obvious that you need to click the legend, with an indirect focus. The other libraries maintain a line plot with 0.2 opacity after you click one to allow us to compare the chosen library against the others general trends. This interaction is useful because it shows us the informative distributions that would be too much if the information was put all together in, for example, a stacked bar chart. This way, it has a popout-like effect with emphasis on the certain library you selected. It provides a simpler, more aesthetic option to understanding the differences within the libraries. This is what the interaction makes possible that a static version wouldn't have; you get a good understanding of one library of interest while still having the reference to the other libraries in the background. We can itemize this information as the following:

- How does the distribution of average total checkouts per user vary for the selected library over time?
- How does one library distribution compare to the other library trends in average total checkouts per user?
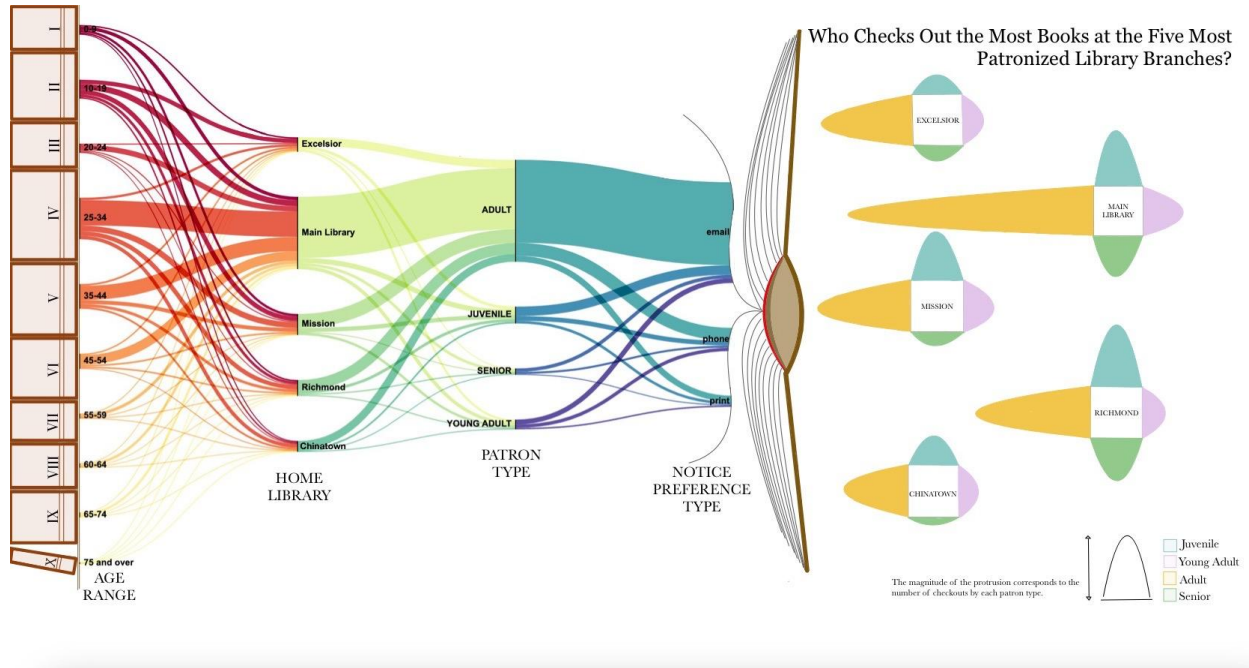
In terms of effectiveness, our visualization has optimally matched channels. We can realize this viz is highly effective since we have the use of a line chart for analyzing checkouts over a period of time, paired with the bar chart based on the same scale allows for easy year-to-year comparison. By having both graphs layered together, it does make the visualization more complex, however it is much more simpler to the alternative of representing this information in a stacked bar chart. In terms of expressiveness, the visualization accurately and faithfully represents the data. We only use five colors, ensuring that we can easily discriminate between them. This tells us that it follows the color hue rule in reference to the channel characteristics. Each part of the line for every year is distinctly marked with dots, enhancing readability and separability. The interactive feature introduces a pop-out effect, allowing a distinct focus on the selected home library. Overall, the visualization is highly effective, showcasing separability and effectiveness.

It is important to note the faults of this visualization. With a larger time range, we could likely get a better sense of overall trends; having a larger data sample to allow for better statistical accuracy could highly improve this visualization. We were able to find that Chinatown had a general upwards trend over time of the average total checkouts per user, and also that the Main Library had a downwards trend, but this may not be the case with that added information. Additionally, it would be helpful to have the mouse change when you hover over the legend to help the interaction be more explicit. This allows for users to learn more from the visualization and in a quicker amount of time. It also would be interesting to compare more libraries, but we only selected five for the simplicity of the visualization. Having a faceted bar chart could provide more information that wouldn't make the reader have to remember what the other library distributions look like, but we wanted a simplistic visualization that wouldn't take up much space. That is why we decided to have both the line and bar plot together, so that you could get a comparison easily to the other libraries. This is why we ended up choosing this viz, because it is a good expressive graph that gives us actual insights into library usage. We can notice the differences between the libraries, and it is done so in an aesthetic way. This visualization stands out for its ability to convey meaningful findings in the differences of library usage over time.
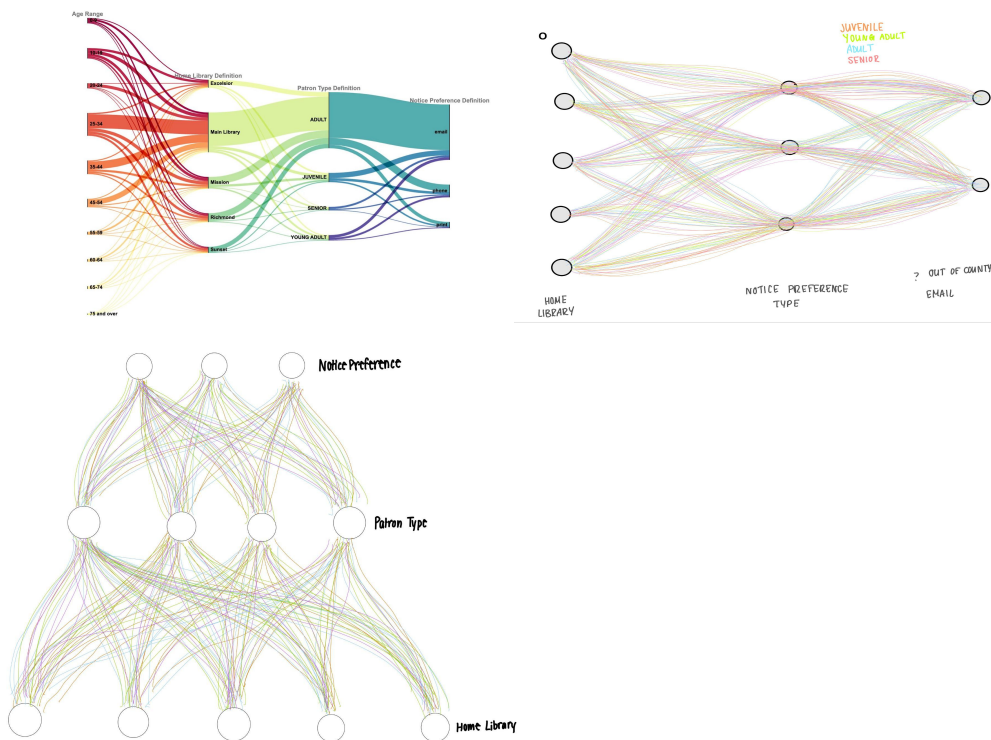
# Justifications for Each Visualization

## Visualization 5

*Novel Viz*



Previous iterations:

# Justifications for Each Visualization

Our novel visualization serves to showcase the interactions between several demographic attributes with respect to patterns in library engagement at the five most patronized library branches, supplying an informative and aesthetically pleasing overview that provides context to the more focused tasks.

The alluvial chart demonstrates the flow of the relationships between four major demographic attributes, each represented by a node: Age Range, Home Library Definition, Patron Type Definition, and Notice Preference Definition. Visualizing the proportions of each attribute level in conjunction with the nodal transitions provide insight into how these attributes interact with one another, facilitating an understanding of the larger picture requiring little effort on the part of the audience. For example, you can quickly and easily observe that the age ranges making up a larger proportion of library users are 25-34 years, 35-44 years, and 10-19 years; the Main library is the most common home library; and that adults in the 25-54 year interval are more likely to utilize the Main Library, while more senior patrons are more likely to frequent one of the smaller branches.

The diagrams adjacent to the alluvial chart provide a simple comparison of user engagement (based on the metric 'checkouts') between the five libraries of focus. Each diagram represents a different library branch, and consists of a square with parabolic projections protruding from every face. The color of the projection corresponds to patron type, while the magnitude of the projections represent the number of total checkouts by that patron type. Utilizing position to delineate the five libraries promotes discriminability. The use of color to encode the four patron types and position to encode the magnitude of total checkouts is in keeping with the expressiveness principle. This aspect of the overall visualization makes good use of the channel characteristics popout and grouping; the audience can see which types of patrons engage the most at each library, and compare performance against the other libraries. More could be done in future iterations to improve the accuracy channel; while the diagrams are to scale, because it was done by hand, the values being represented are not entirely intuitive.

Deeper investigation leads to a more nuanced understanding of lapses and strengths in the provision of library services, particularly when considered alongside the alluvial chart. While older patrons tend to have a smaller, less urban library as their home library, total checkouts for this demographic are highest at the Main Library. This indicates that seniors have to travel to the Main Library because their  smaller, home branches may not have the materials they are most interested in. While checkouts for the senior demographic may be low at smaller libraries, higher membership of this demographic at those locations indicates that they may utilize the library for other purposes, like social ones.

# Reflection

**Project Strengths**

1. Through our work, we discovered and communicated trends regarding library usage, such as variation in the promotion of patron engagement across different library branches, relevant to the stakeholders and policy makers who make decisions regarding resource allocation for public services.
2. We collaboratively came up with appealing and interpretable visualizations relevant to our outlined tasks, that can be easily understood and interacted with by our intended audience.
3. Our visualizations strike an appropriate balance between addressing the given task and adhering to our primary objective; resulting in a comprehensive and cohesive final output.

**Project Weaknesses**

1. Once we cleaned our data, it became clear that there was much less usable data recorded than we thought. Our analysis was unfortunately limited in scope by the collection methods and semantics of our data.
2. The subject matter of our project may not succeed in captivating the interest and attention of audience members that do not have a vested or personal interest in library usage.

**Things we would do differently**

1. Choose a different data set, that would allow for a more robust and diverse analysis, with more attributes to factor into our visualizations
2. Utilize our time in a way that is more commensurate to the weight of the deliverables.
3. Find a better way to share code with one another.

# Work Distribution

## HIGH LEVEL OVERVIEW

| Project Milestone | Team Member 1 Sadia | Team Member 2 Heidi | Team Member 3 Lillian |
|---|---|---|---|
| Milestone 1 | 33 | 33 | 33 |
| Milestone 2 | 33 | 33 | 33 |

## BREAKDOWN

| Milestone Description | Sadia | Heidi | Lillian | Percentage of Total Work for the Deliverable |
|---|---|---|---|---|
| **Milestone 1 - Implementation & Report** | | | | |
| Task 1 - Finding Data | 33 | 33 | 33 | 10 |
| Data Abstraction | 100 | 0 | 0 | 5 |
| EDA | 100 | 0 | 0 | 20 |
| Tasks and visualization | 0 | 50 | 50 | 30 |
| Written Work | 20 | 40 | 40 | 30 |
| Next Steps | 100 | 0 | 0 | 5 |
| **Milestone 2 - Implementation & Report** | | | | |
| Task 1 | 100 | 0 | 0 | 21 |
| Task 2 | 30 | 0 | 70 | 21 |
| Task 3 | 30 | 70 | 0 | 21 |
| Task 4 | 0 | 100 | 0 | 21 |
| Final Dashboard | 100 | 0 | 0 | 11 |
| Reflection | 10 | 30 | 60 | 5 |
| Novel Vis | 0 | 0 | 100 | 21 |