

Visualizing Boston's Subway System

September 13th, 2014

project: mbtaviz.github.io notes: mbtaviz.github.io/datacon-handout.pdf

Mike Barry [@msb5014](#)

Brian Card [@bmcard](#)



Edward Tufte (@EdwardTufte) Jun 10
Graphical timetables
[edwardtufte.com/bboard/q-and-a...](#)
once handmade
now Barry+Card [mbtaviz.github.io](#) !
#Visualization

Mike Bostock (@mbostock) Jun 10
@mbtaviz Beautiful work!

Background

Visualizing MBTA Data is an interactive report of the performance and behavior of Boston's subway system over the month of February 2014. The visualizations focus on the performance and behavior of the subway system which differs from traditional train visualizations that focus on the train schedule. This started as the term project for a graduate course in data visualization. We collected the data in February, made mockups and prototypes in March and built project in April. After the course was over we spent May fine-tuning and published the end product on June 10, 2014.

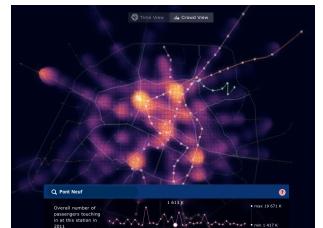
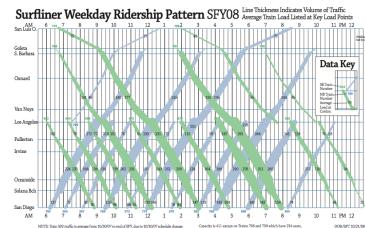
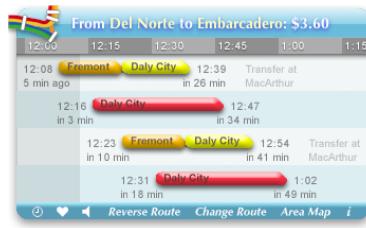
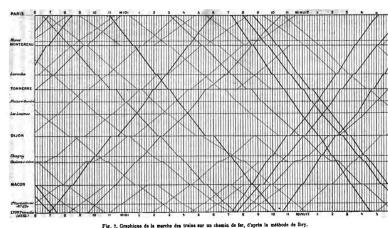
The Design Process

We started by trying to understand what information would be interesting to people and then narrowing down our focus based on the data that we could gather. We outlined three questions to focus on: **When and where are the trains crowded or delayed? How do snowstorms or special events affect the train system? How congested or delayed is my route?**

The MBTA's realtime data feed provides data for the red, orange and blue lines, but not green or silver, so the report focuses on those lines. Data is provided in JSON format and contains the current location of each train as well as the predicted time to the next stop. The real time data allowed us to see how delayed the trains were, the MBTA also provided station entry and exit data that we used to understand where the system was congested.

Existing Works

Below is a sampling of existing train visualizations. From left to right, Étienne-Jules Marey's schedule from 1885 [7], Bret Victor's Bart Widget for trip scheduling ([worrydream.com/MagicInk](#)), State of California Department of Transportation report on ridership [8], metropolitan.io a visualziation of the Paris train system ([metropolitan.io](#)).

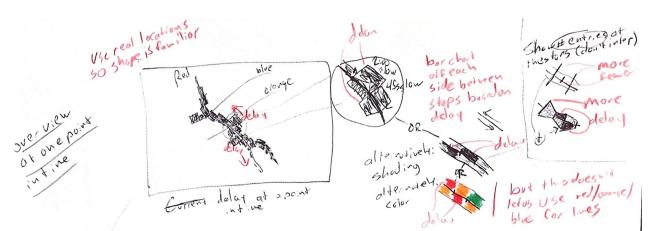


Both Marey and Victor focus on scheduling and not on train system performance. Tufte [9] provides many visualizations of train schedules, all based off of Marey's original design. The State of California also applies Marey's layout and additionally encodes the number of riders on each segment over the year. Metropolitan.io is based of of similar data to what the MBTA provides but uses different visualization techniques.

Prototypes

The project started with more than a dozen mock-ups, of these 6 to 8 made it into the prototype phase where we made samples with real data. Prototypes helped us validate ideas and find trends in the data to explore further. Several ideas such as the congestion and delay visualization only came after significant investment in prototypes.

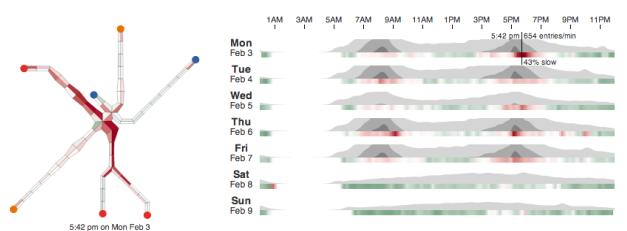
Mockup



Prototype

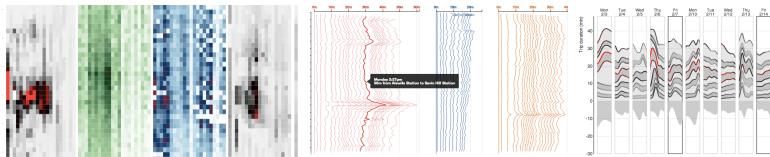


Final Implementation



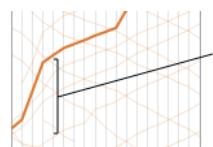
Iterate, Iterate, Iterate

We went through several iterations of each prototype before settling on the versions shown in the report. Several other prototypes were useful in exploring the data but couldn't be integrated cleanly. The image on the left shows congestion and transit times together over time but was hard to follow. The middle show delay from each line, a first attempt at the Your Commute visualization. The one on the right shows a more complex Your Commute visualization that we left out because it didn't reveal much more than a scatter plot, and the scatter plot was more obvious.



Layout and Publishing

Once we discovered the interesting features of our dataset we created the narrative and organized the visualizations into the different sections. We added the interactive text and annotations and iterated on the finished reported until we covered all of the major areas. After we were comfortable with the content we optimized the website and published it online.



	Average Number of Tr
Weekdays	Satu
Red	450
Orange	320
Blue	380
Total	1150
Notice how so	870

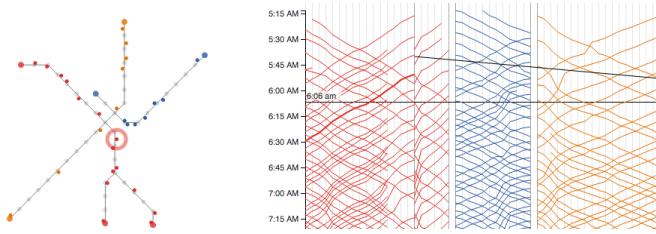
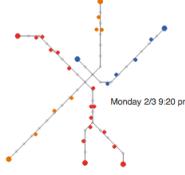
Our Design Concepts

Levels of Abstraction

Bret Victor's *Up and Down the Ladder of Abstraction*

(worrydream.com/LadderOfAbstraction) influenced many parts of the report. The model for our system is train locations on subways lines and we represent this using a small map that closely resembles the subway map seen on most trains www.mbta.com/schedules_and_maps/subway.

An abstraction one level up is the trains on subway lines moving over time. The Marey diagram represents this abstraction, but is difficult to read and interpret. Victor's advice is to tie these two together using lightweight interaction which lets the user understand the abstraction in terms of the model. This way the higher level behaviors and patterns that are clear in the abstraction can be more easily related back to the model that they represent. In our case it's easy to see the delay caused by the trains with mechanical failures in the Marey diagram because the slope of the lines change drastically. However it's difficult to see that the trains bunch up and slow down in the diagram. The tight feedback loop of hovering over a time on the Marey diagram and seeing where the trains are at that time (as well as lightweight interaction to link the highlight corresponding parts on hover) reinforces the connection the levels of abstraction and helps make the Marey diagram more understandable.



Avoiding Administrative Debris

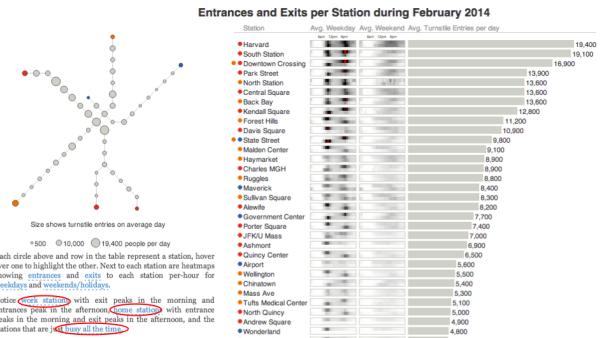
A few things you won't find in our visualization are checkboxes, sliders, or buttons (aka "administrative debris").

At first we were not going to include the Marey Diagram, we just wanted to show trains on a map and let the user control time using a slider. We realized though that the Marey Diagram could serve as an information-rich interface for controlling the visualization while avoiding typical administrative widgets.

Instead of buttons with short, cryptic labels we allow you to click on or hover over linked words to control the visualization. These words are inside of sentences organized into paragraphs that put the point of interaction into context with the rest of the narrative.

The table and map below breaks down February's turnstile entries and exits by station. Hover over a row in the table to highlight the corresponding circle or the map, or vice-versa. Click on a row in the table to show a detailed heatmap for the entrances to and exits from that station over the month. Click and drag or several table rows to highlight a range of stations.

You can see the busiest stations are all along the Red Line ([Harvard](#) dropped the list, followed close to [South Station](#) and the [Downtown Crossing](#)). Next to each station are heatmaps showing [entrances](#) and [exits](#) to each station per-hour for [weekdays](#) and [weekends/holidays](#). You can see that some stations are [work stations](#) since their exits peak in the morning and entrances peak in the afternoon and that some stations are [home stations](#) since their entrances peak in the morning and exits peak in the afternoon. Some stations are just [busy all the time](#).



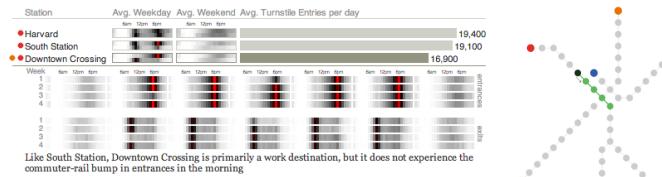
Interaction Considered Harmful

We avoided mandatory interaction as much as possible. Even though we provide details on demand through tooltips and dynamic annotations, chart keys and labels provide the same information statically. The fact that the entire visualization is readable as an image also makes it easier to explore on mobile devices.

We did compromise on this in a few places. We wanted to use small multiples for the lined-up Marey and congestion diagrams, but due to the low resolution of computer screens we could not get enough detail out of a smaller version of each. Instead of small multiples we allow you to hover over a particular time to change the data shown in the graphic



Also, the turnstile heatmaps and your commute scatterplots were expensive to render so we require the user to select which instance they want to drill down into to see more details, otherwise the page would take too long to load.



In all of these cases we start with initial default values so the user is never presented with an empty chart prior to interacting with it.

Integrating Words, Numbers, and Graphics

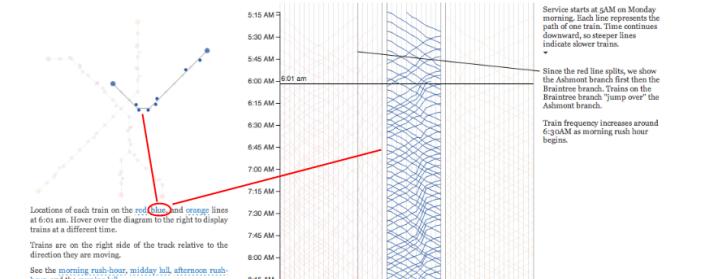
Data does not always require a visualization and visualizations do not always need to be interactive. We use sentences to show datasets with only a few numbers, tables to explain datasets with more numbers, and graphics to explain datasets with hundreds or thousands of numbers.

In a typical weekday, trains make approximately [1150 trips](#) on the red, orange, and blue lines starting at 5AM and continuing through [11PM](#) the next morning. On Saturdays trains make [870](#) trips and on Sundays they make [760](#).

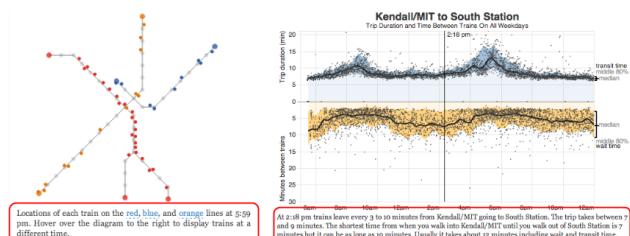
To better understand how the trains operate on a typical day, below are all trips that trains took on the red, orange, and blue lines on Monday February 3 2014. Each vertical line represents a station, and time extends from top to bottom. Steeper lines indicate slower trains. This visualization was first used by

	Weekdays	Saturdays	Sundays
Red	450	350	300
Orange	320	260	220
Blue	380	260	240
Total	1150	870	760

As a twist on this concept, we also link hovering over or clicking on words to changing associated graphics:

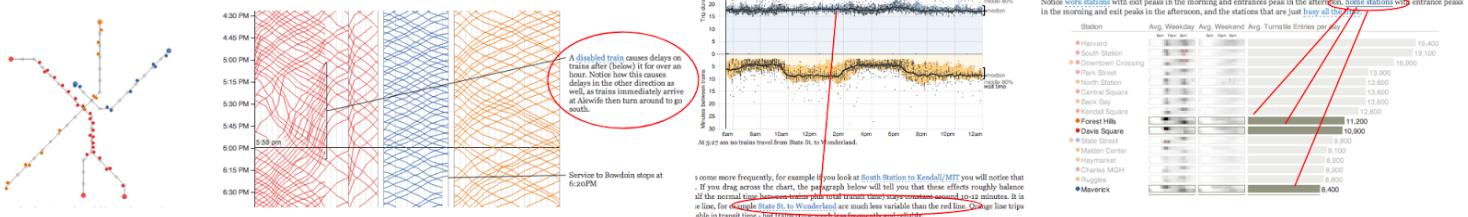


and we link hovering over graphics to changing the words and numbers that explain the graphic:



Show Me What's Interesting

We spent a lot of time analyzing the data to pick out the interesting points and showed them off with liberal use of annotations. This lets people who spend a short amount of time playing with the visualization get the highlights while users who spend more time can explore and dig into the details. If there are anomalies in the data explaining the cause of the anomalies is an easy way to show the viewer something interesting.



What We Didn't Do

We wanted to show how crowded each individual train was and also to come to a sweeping conclusion about the train system, but the data wasn't there to do either, so we ended up just showing the data. We also wanted a "grand visualization" that showed all of the variables together in a single view, but we couldn't quite pull that off. Instead we went with separate views, each of which showcased a particular aspect of the data.

Gathering Data

We wrote simple bash scripts and scheduled them using cron jobs on old laptops in parallel for the month of February to get the realtime data files. We also pulled the MBTA GTFS zip file which contained the scheduling data for all trips and worked with a contact at the MBTA to get per-minute entry and exit counts at each station from their turnstile data.

Merging and Cleaning the Data

We experimented with python, scala, and node.js for data pre-processing. Python was an order of magnitude slower than scala or node.js which were comparable. Since only one of us was comfortable with Scala and we both knew JavaScript, we chose node.js for continued work. We found the bulk of the time spent running our scripts was from opening files from disk (12 files per minute of gathered data) so our preprocessing scripts merged the two datasets, removed redundant data, and wrote the data out to hourly gzipped data files with a single json blob per line, for example `redline/yyyy/mm/dd/hh.json.gz`. This gave the optimal performance for downstream feature-extraction scripts, while still allowing us to process date ranges at hour granularity.

Extracting Visualization Data Files

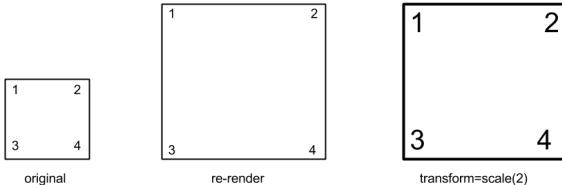
We wrote a node.js script for to generate a minimal JSON data file for each visualization embedded within our larger web page with the goals of minimizing file size and reducing the amount of work done in the browser at runtime. Almost every script ran inside a reusable helper that took a date range, opened the necessary `.json.gz` files and called a function with each line of that file.

Building the Web App

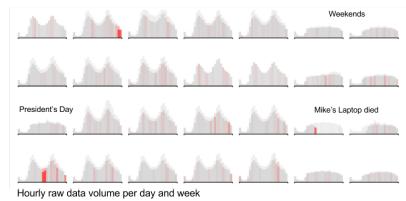
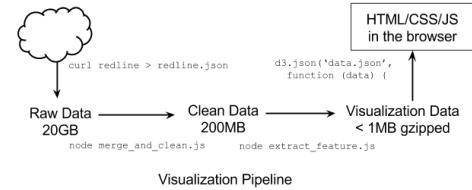
The web application is a single HTML file with descriptions and explanation surrounding placeholder empty `<div>` elements for our visualizations. Our JavaScript files load at the end, read in the data files, render SVG dynamically into the placeholders, and setup listeners to handle interaction. We used a few open source projects to make development easier including D3.js ([d3js.org](#)) for mapping our data to DOM elements that create the visualization, underscore.js ([underscorejs.org](#)) for collection manipulation, jQuery ([jquery.com](#)) for additional DOM manipulation utilities, moment.js ([momentjs.com](#)) for date and time formatting, and d3-tip ([labrat revenge.com/d3-tip](#)) for creating simple tooltips. Below we describe how we solved several problems along the way. See the github repo for the full details ([github.com/mbtaviz/mbtaviz.github.io](#)).

decoupling the visualizations We wanted to make sure our visualizations did not interact in undesired ways when placed on the same page. We used the less CSS preprocessor ([lesscss.org](#)), which allowed us to nest selectors and ensure that all rules for a visualization would apply only to that section.. Each visualization got its own copy of the data it needed and code had to be explicitly put into a common location to share between visualizations.

responsive SVG using D3? This is tricky, since D3 guides you down the path of using fixed pixel sizes for all of your elements. To make your visualization grow or shrink with the screen size you need to programmatically resize it. A few options are:



We use the re-render option in most places to preserve text size at different scales, but for the turnstile heatmap which is expensive to render, and the text feels small at the default size, we use the `transform=scale(...)` option to make our text bigger as well.



data loading Since we needed to load many data files and render several different visualizations we created a utility that requests several files asynchronously and coordinates progress and completion events, bubbling them up to user-defined listeners

```
VIZ.requiresData([
  'json!data/file1.json', 'json!data/file2.json'
]).progress(function (percent) {
  d3.selectAll(".progress").text('Loading data... ' + percent + '%');
}).onerror(function () {
  d3.selectAll(".progress").text('Failed to load data');
}).done(function (dataFromFile1, dataFromFile2) {
  // use data
});
```

screen size/responsive design We developed on 1200px wide macbooks and optimized for that, but some people we showed this to viewed on iPhones, some viewed on tablets, and Bill Shander viewed it on a large, high-resolution monitor. After this feedback we pulled in Twitter's Bootstrap ([getbootstrap.com](#)) and customized its responsive design utilities to create 3 responsive width breakpoints: 789px, 990px, and 1200px. We set `<meta name="viewport" content="width=768">` in our header to ensure that phones and tablets rendered the page using the minimum fixed width. This allowed us to not worry about rendering our visualizations any narrower than 768px and also looked acceptable on a range of screen sizes.

browser compatibility Underscore, D3, and jQuery work across modern browsers and we used an es5/6 shim to let us use newer language features in older browsers. Still, ie8 does not support SVG. Instead of bending over backwards to support them, we used phantomjs (phantomjs.org) to render our page to a png file and show that instead.

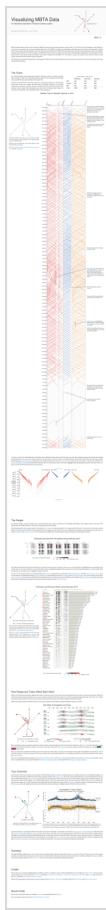
```
<!--[if lt IE 9]>

<![endif]-->
<!--[if gte IE 9]><!-->
The visualization
<!--<![endif]-->
```

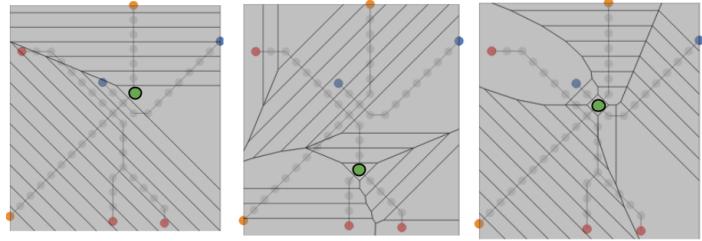
scalability & performance As our visualizations grew in size and the number of DOM nodes increased, we had to optimize a few things to get acceptable performance. Adding mouseover D3 event listeners to thousands of DOM elements was slow, so instead we added a single listener to the parent element that handled mouseover events bubbled up from its children. We pushed all calculations for each data point into the pre-processing step to avoid slowing down the browser.

To optimize further the first thing to try would be rendering parts of our visualization involving many small DOM elements that don't require interaction to canvas instead of SVG. Alternatively, we could just render those section to images and load those instead of rendering each time.

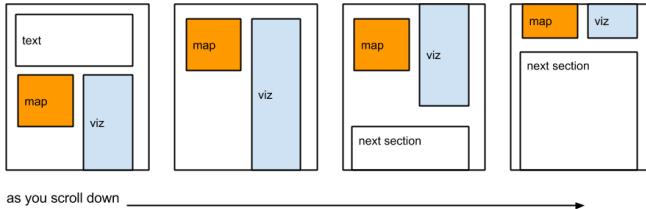
fonts We chose Helvetica for our headings to match the MBTA's sign style. For our body text we wanted to go with Garamond but after some reading learned that Georgia is better since it is optimized for computer screens, whereas Garamond is optimized for print.



voronoi picker Our final visualization allows you to click and drag between two stops to see a scatter plot of delay between them. To implement our rules for valid destinations from a source efficiently as you drag, we used D3's voronoi utility which takes a set of points and renders a polygon for each where all points inside that polygon are closest to the point. We render voronoi polygons for the set of valid destinations every time you click and start dragging, and add our mouseover listeners to those polygons as shown below.



fixing maps in place as you scroll We wanted some fixed context while scrolling to navigate the visualization, so we made a utility that listens on scroll events and programmatically sets the `top` and `left` css attributes of the map as you scroll to make the map appear to stay fixed to the top. Since scroll events fire on iOS devices only when the scroll finishes, we disable this entirely for iPhones and iPads.



Performance was slow on firefox and safari at first, but we added `transform: translate3d(0,0,0)` CSS style which causes browsers to hardware accelerate animations and the improvement was noticeable

Lessons Learned

While trying to understand our large dataset, we avoided using complex statistical techniques to over-reduce our data and instead did the minimal reduction necessary to fit it into high-resolution images. When we could see the data and all of its intricacies, our eyes drew their own conclusions.

Mockups and prototypes helped us formulate ideas but we found that nothing beat iterating on working code. When we had something concrete to talk to we came up with new ideas that we would never have thought of in a mockup or prototype. We experienced first hand that higher levels of abstraction allow for powerful reasoning, but you need to start at a level the audience is used to. We had success using words, numbers, and images to explain our data with no preference towards one or the other. We mostly avoided gratuitous interaction and color but found that subtle use of both allowed us to better link words to images. We wanted to use more small multiples, but had to compromise since the resolution of computer screens is still much lower than paper and instead used interaction to change a single larger graphic. When necessary, lightweight interaction with words and images served as a much richer alternative to typical administrative debris like buttons, checkboxes and sliders.

We learned that the T experiences the worst congestion and delay from mechanical failures followed by congestion from the increase in riders from rush hour. The MBTA does a good job of counteracting the increase in riders by increasing the trains during rush hour on the red line which keeps the commute time fairly consistent on average. The blue line has the fewest riders and experiences the least amount of delay, but also has the newest trains. The trains on the orange line break down frequently which can cause highly variable wait times. The system performs better during snow storms which might be because far fewer people ride the trains, and special events like basketball games have little impact on the overall system.

References

- [1] Roberto Scalese. "Students Create Insanely Awesome MBTA Data Maps." Online. 2014. www.boston.com/news/local/massachusetts/2014/06/10/students-create-insanely-awesome-mbta-data-maps/UIGfqPynVXU4BMMsfqLN9K/story.html
- [2] "Capturing Boston's Subway System In Stunning Detail." Online. 2014. www.citylab.com/commute/2014/06/capturing-bostons-subway-system-in-stunning-detail/372649
- [3] Mona Chalabi. "Ctrl + ← The World Cup, the World Cup And the World Cup." Online. 2014. fivethirtyeight.com/datalab/ctrl-%E2%86%90-the-world-cup-the-world-cup-and-the-world-cup
- [4] Andy Kirk. "10 significant visualization developments: January to June 2014." Online. 2014. www.visualisingdata.com/index.php/2014/08/10-significant-visualisation-developments-january-to-june-2014
- [5] Nathan Yau. "A visual analysis of the Boston subway system." Online. 2014. flowingdata.com/2014/06/11/a-visual-analysis-of-the-boston-subway-system
- [6] "Best of DataIsBeautiful Results [January - July 2014]" Online. 2014. www.reddit.com/r/dataisbeautiful/comments/2e6l6j/best_of_dataisbeautiful_results_january_july_2014
- [7] Marey, Étienne Jules. "La Méthode Graphique Dans Les Sciences Expérimentales Et Principalement En Physiologie Et En Médecine. 2. tirage augm. d'un supplément sur le Développement de la méthode graphique par la photographie; avec 383 figures dans le texte." Paris: G. Masson, 1885. archive.org/details/lamthodegraphiq00maregoog
- [8] Nguyen, Lam. "Surfliner Weekday Ridership Pattern." Memorandum to Chairs and Commissioners regarding Follow-up on September and December 2008 Rail Items. State of California Department of Transportation. January 14, 2009.
- [9] Tufte, Edward. "Envisioning Information." Graphics Press, Cheshire, CT, USA. 1990.