

Not Safe Anymore

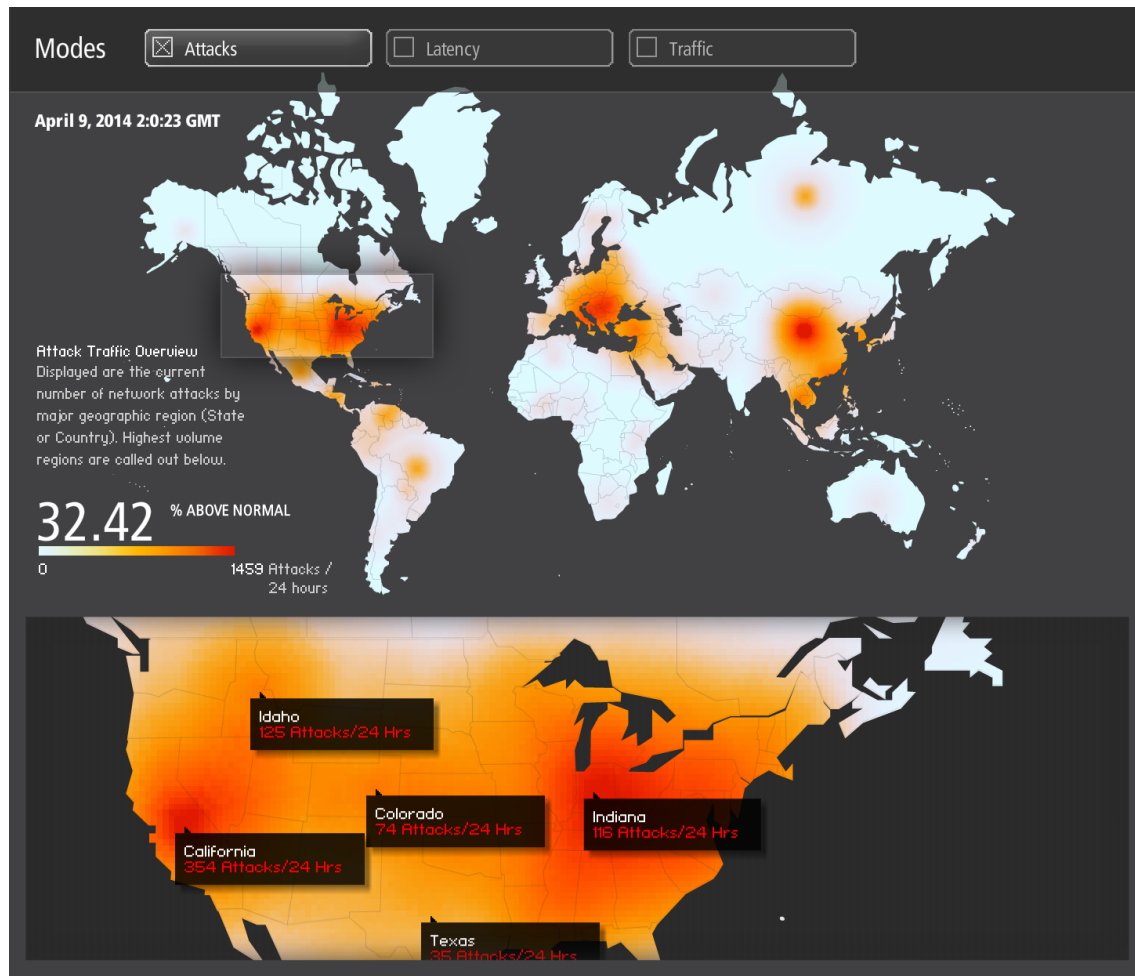
Visualizing network attacks detected by Akamai firewalls

Overview and Motivation

Web Security is an on-going concern for all organizations and individuals. It is also a time and resource consuming exercise that needs to be balanced against other priorities. Given the recent online attacks on big banks¹ and retailers², we are interested in visualizing the attributes of such attacks, including which countries originate the most attacks, what times of the day are popular for such attacks, and what networks (ISPs) are they carried out on. Having a better understanding of these attacks should help us better understand how to protect ourselves.

Related Work

One of the core inspirations was the Akamai Real-time Web Monitor, which can be seen here:



¹ Cyber attacks against banks more severe than most realize, <http://www.reuters.com/article/2013/05/18/us-cyber-summit-banks-idUSBRE94G0ZP20130518>

² Target cyber breach hits 40 million payment cards at holiday peak, <http://www.reuters.com/article/2013/12/19/us-target-breach-idUSBRE9BH1GX20131219>

This site provides an easy to read high-level overview of attack data using a global map and detail view. While this is useful for seeing a general overview, it is difficult to explore the data to surface patterns and trends. By providing more extensive filtering tools, and additional linked plots, we expect to be able to help individuals learn more about this data.

Questions

Our initial intent was to explore the following questions:

- Where do the most attacks originate?
- At what times of day do most attacks originate?
- Which ISP's are used by most attackers?

We also hoped to understand the interplay between these questions, such as if there is a correlation between a given time range and the origin of the attacks.

As we've learned more about the data we are also hoping to be able to extract additional data, such as:

- What user agent's are used for most attacks?
- What operating systems are used for the most attacks?
- Could we identify what portion of the attackers were bots?

Data

We were able to obtain 24 hours of raw data directly from Akamai to assist with this project. This dataset is a sample taken from the firewalls running on Akamai's³ global network⁴ of servers that serve web content all over the world. These firewalls log and report each intrusion attempt to Akamai's monitoring systems. This dataset was scrubbed of sensitive information by dropping dimensions that identify customers and only the keeping the dimensions needed for this project. Although Akamai collects data points on each firewall intrusion, given the volume of data and our visualization goals, we have decided to aggregate the samples on a per 5 minute basis. Finally, we are planning use data collected for a 24 hour time period for this project.

This raw data is a simple CSV file, a sample of which is shown here:

```
timestamp,country,state,city,lat,lng,network,useragent
1393976466191,US,VA,HERNDON,38.9807,-77.3799,,Mozilla%2F5.0+%28Windows+NT+6.1%3B+WOW64%3B
+rv%3A23.0%29+Gecko%2F20100101+Firefox%2F23.0
1393976513886,US,VA,HERNDON,38.9807,-77.3799,,python-requests%2F2.0.1+CPython%2F2.7.3+Lin
ux%2F3.2.0-40-virtual
1393976530853,US,VA,HERNDON,38.9807,-77.3799,,Mozilla%2F5.0+%28Windows+NT+6.1%3B+WOW64%3B
+rv%3A23.0%29+Gecko%2F20100101+Firefox%2F23.0
1393976531621,US,VA,HERNDON,38.9807,-77.3799,,Mozilla%2F5.0+%28Windows+NT+6.1%3B+WOW64%3B
+rv%3A23.0%29+Gecko%2F20100101+Firefox%2F23.0
```

³ Akamai, <http://www.akamai.com/>

⁴ Akamai's global network of ~ 137k servers,
http://en.wikipedia.org/wiki/Akamai_Technologies#The_Akamai_Network:_Edge_Platform

To be able to use the data in a useful manner we created two tools. First, a Ruby script was created that parsed the user agent to extract the browser, browser version, and operating system. Next, a Java application was created that took the parsed data and aggregates it into buckets based upon the timestamps. This tool was used to create data files for bucket sizes of 5, 10, 15, 30, and 60 minutes. As mentioned, we are currently using the 5 minute data for our work. The other datasets were created so that we can experiment with different sampling intervals to see which one works best.

This is what a sample of our processed data looks like:

```
timestamp,country,state,city,lat,lng,network,browser,browserVersion,os,count
1393976400,CN,SC,CHENGDU,30.67,104.07,chinanet,Internet Explorer,6.0,Windows,2
1393976400,CN,SC,CHENGDU,30.67,104.07,chinanet,Internet Explorer,7.0,Windows,1
1393976400,GB,EN,LONDON,51.5,-0.12,-,Internet Explorer,6.0,Windows,957
1393976400,GB,EN,LONDON,51.5,-0.12,bskyb,Safari,7.0.1,Macintosh,15
```

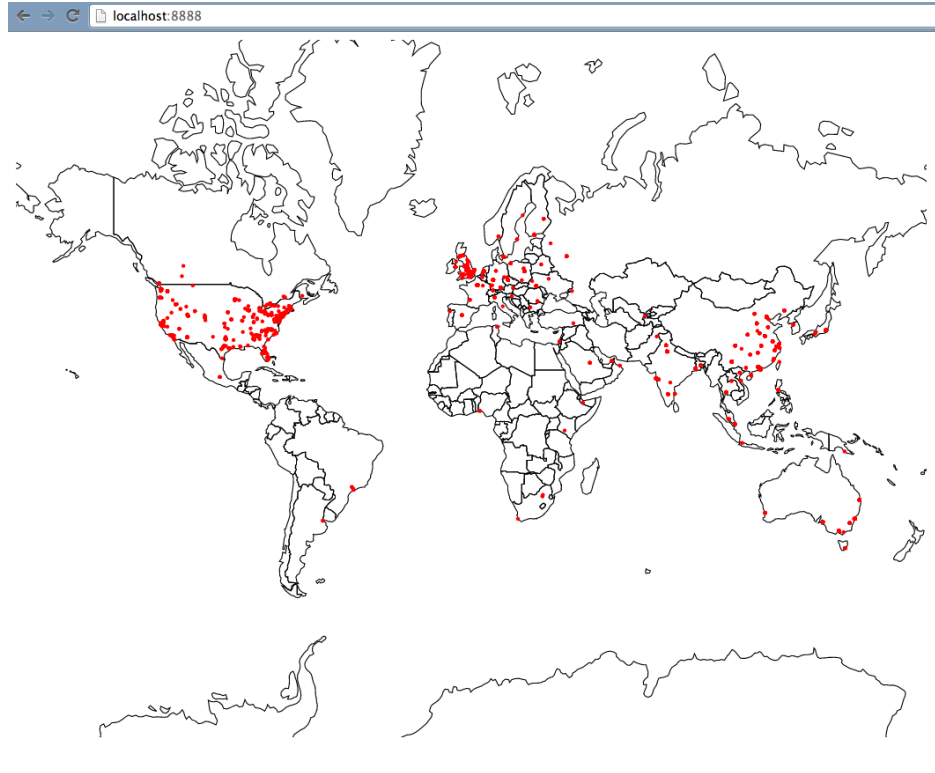
Following is our data model:

Field	Data Type
timestamp	Quantitative (interval)
country	Categorical
state	Categorical
city	Categorical
lat, lng	Quantitative (ratio)
network	Categorical
browser	Categorical
operating_system	Categorical
number_of_attacks	Quantitative (ratio)

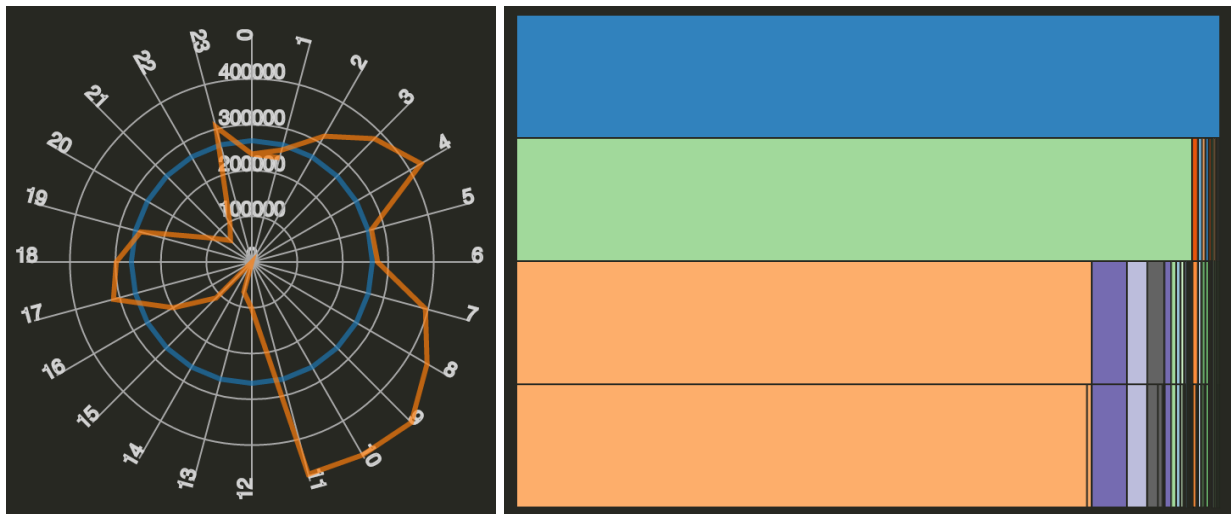
In addition to the attack data, we are using the country code data and world map data from Homework 4.

Exploratory Data Analysis

We started our process by simply visualizing the data on a map to ensure that the data was being parsed and read properly. This resulted in the visualization that you see here:



Once we had confirmed that we could load the data we decided to create several visualizations that would give us a quick understanding of what the dataset contained. These included:



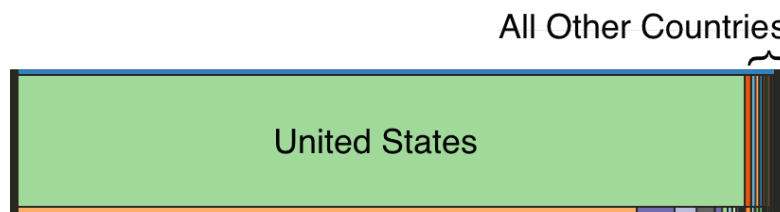
- a polar plot that showed attacks aggregated by hour of the data
- an icicle plot that compared the number of attacks at the state and city level

This allowed us to quickly compare the magnitudes of different facets of the data. We also created filters for country and network, to allow us to look at subsections of the data in more detail.

By performing this early exploration work, we were able to note several key trends that helped shape our plan:

- Most of the attacks were coming from the United States.
- There were highs and lows in the time series that implied that there were patterns in the data to understand.
- There were strong correlations between countries and networks.

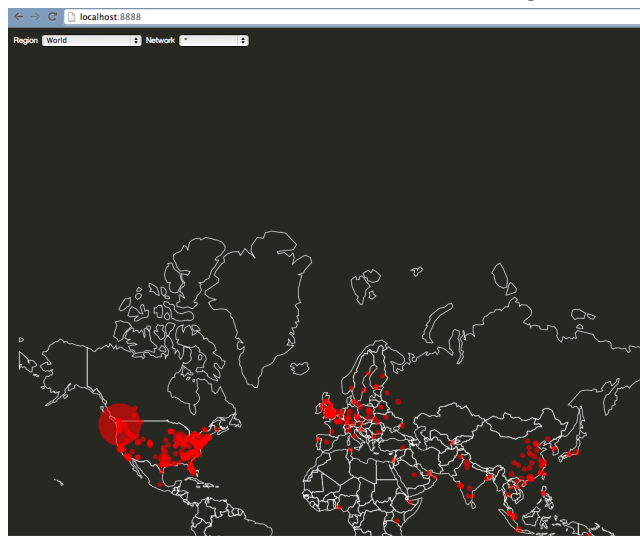
By noting that attacks in our dataset were overwhelmingly coming from the US, it became apparent that the icicle plot would not be suitable. The sizes of the rectangles were so disparate that the other countries simply displayed as a blur of colors. Also, the plot was littered with colors that were not tied to any other elements in the visualization, which was more distracting than informative.



The other patterns noted in the data indicated that there would be value in keeping both the maps and the polar plots, and in adding tables that showed the regional values. There also seemed to be value in being able to filter data by time. This led to the addition of an histogram of hourly data that could be used to filter the dataset using brushes.

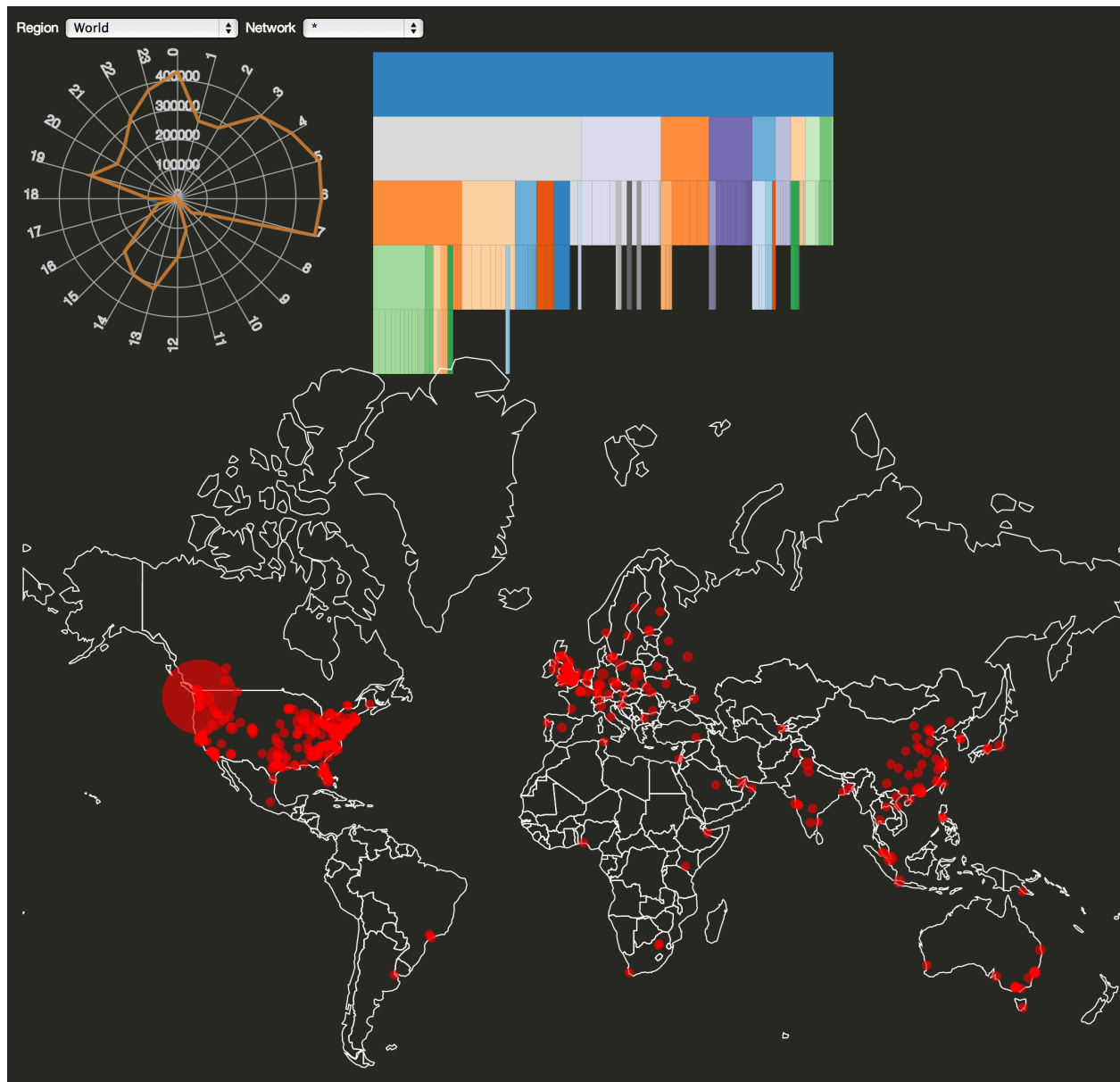
Design Evolution

We started by considering the aesthetics and figured out the color scheme to use. It made sense to use the color red to indicate attacks and a dark background to create a good contrast.



Our initial designs focused on tools that allowed the user to quickly explore the data. This was helpful as we began to tease apart and understand the data. As mentioned above, the tools

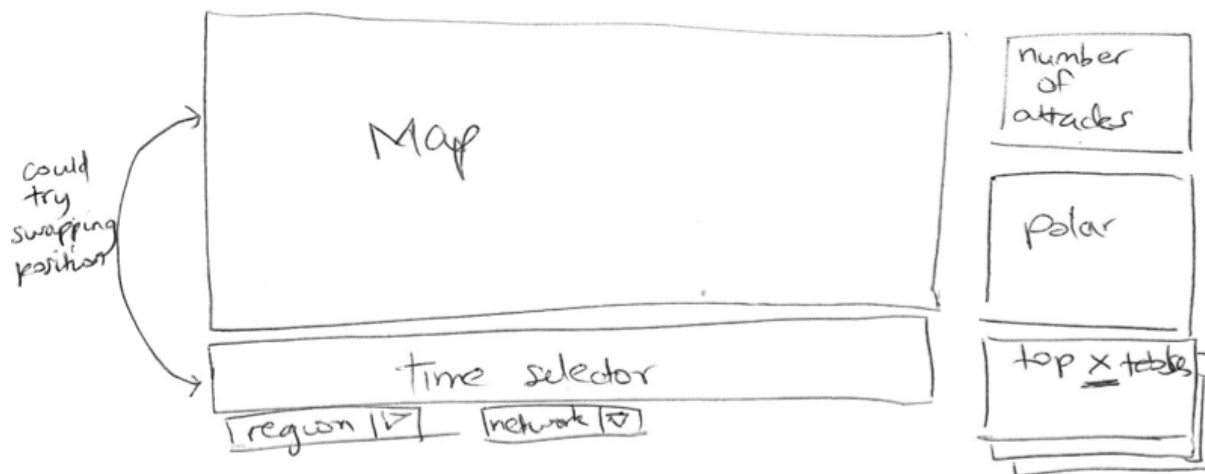
included a map, a polar plot of the aggregated hourly data, and an icicle plot of the the data. It also included drop-down boxes to quickly filter the data. It can be seen here:



The advantage to this view was that most of the data that we were exploring was clustered into the upper-left corner of the screen. That made it very easy to view while reviewing and debugging the code. There were also several disadvantages to this approach:

- the view was cluttered and unclear
- there was not a compelling flow for a user to follow through the data
- colors were used haphazardly and without defined meaning

These thoughts led to changes that significantly improved the experience for our users:

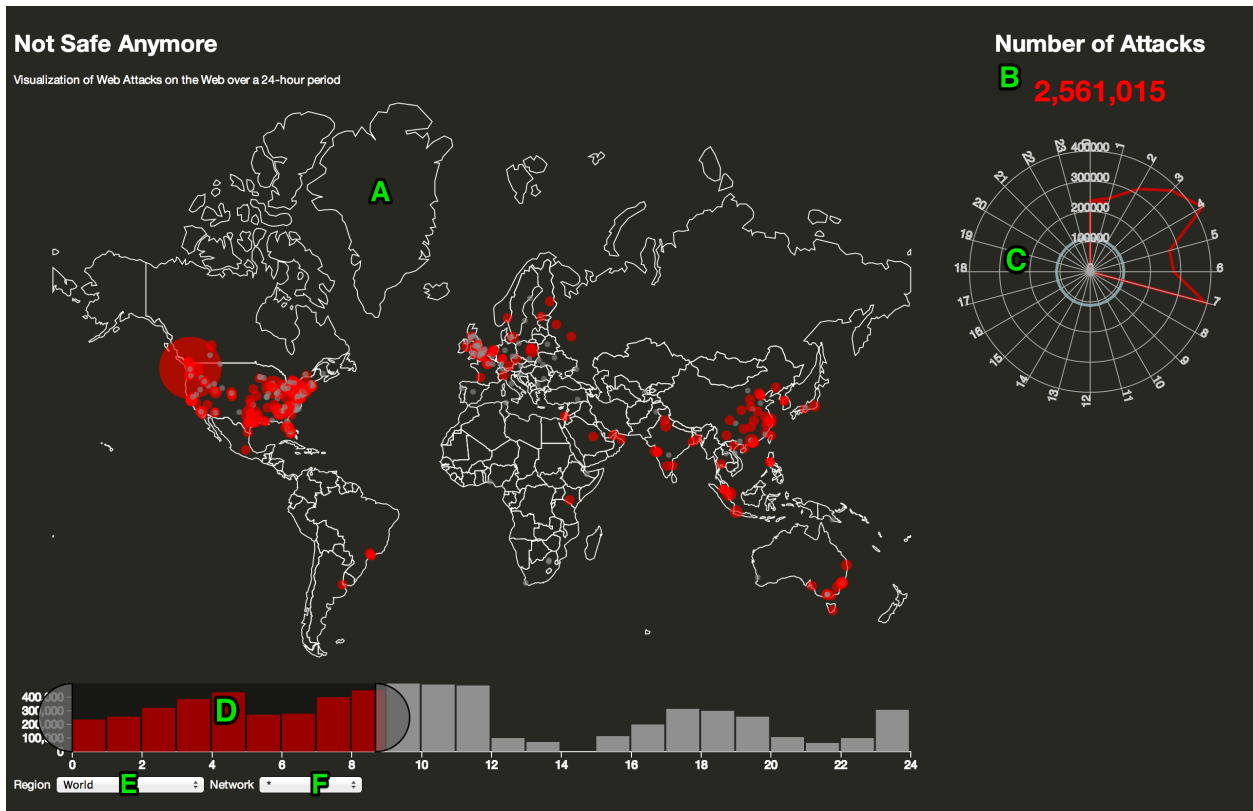


- The map was moved to the upper left to be the dominant screen element, as it took the least cognitive capacity to understand. The upper left was used as our primary consumers of this data are in the Western hemisphere and expect a left-to-right, top-to-bottom flow of information. This eases the user into the experience as it is a comfortable starting point.
- The timeline was added below the map, and large time range selectors were planned. These large selectors convey the idea that the the user should begin here to interact with the visualization.
- The region and network filters were moved below the timeline, to continue leading the users' eyes through the interactive elements.
- All of the read-only visualizations were moved to the right of the page, to supplement the primary ones. This further extended the idea of starting at the left and progressing to the right.
- All data was colored red, and all supporting and inactive elements colored in shades of grey. The choice of a bright color made the data especially dominant on the page, and also showed the correlation between the linked parts of the page. In addition, by using just one color, red, it became easy for folks that are color-blind to view the page.

By focusing on the the flow of the users' eyes, and the sparse use of color for important screen elements, the visualization should tell a much more clear and compelling story.

Implementation

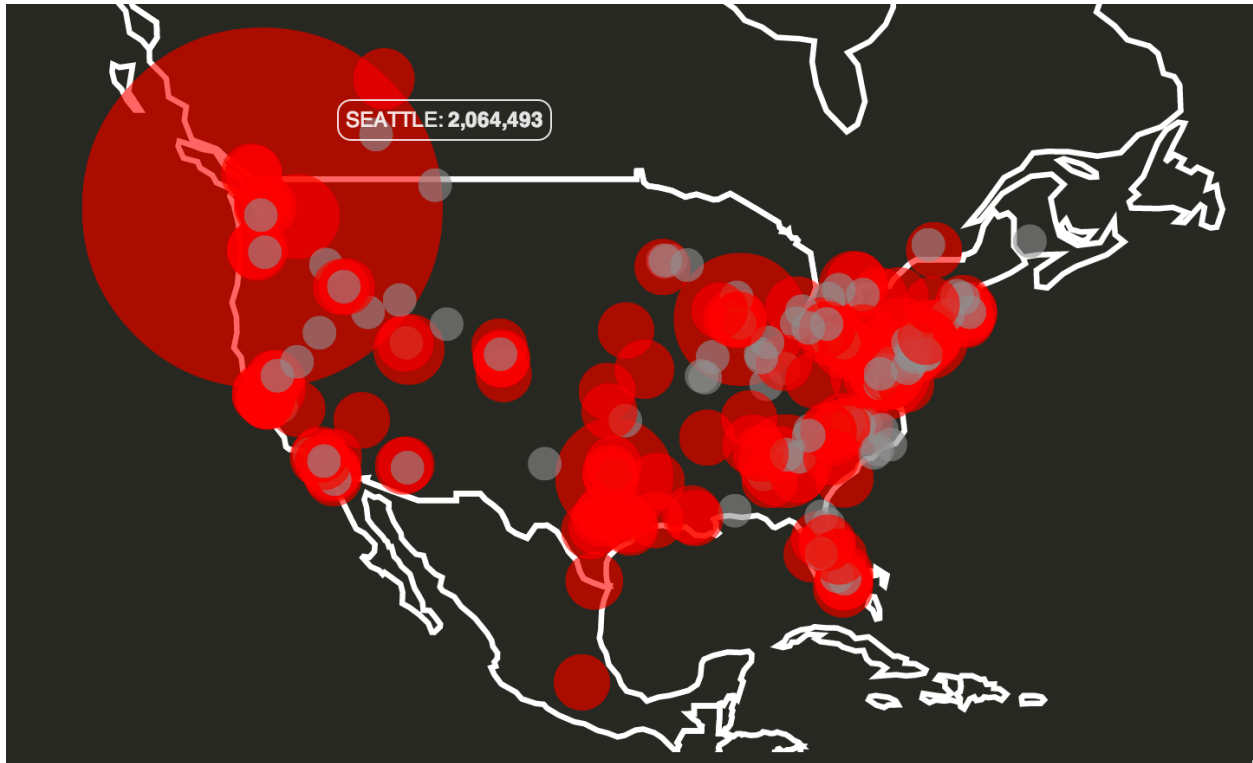
This is what the visualization looks like for the first milestone. All widgets that are part of this visualization are generated as part of a user-defined query. A query consists of the currently selected time range (D), the selected region (E), and the selected network (F). This query is run when the page loads, and can be easily modified by the user at any time.



A. Bubble Map showing attacks globally

This is what a user typically sees first. This map shows the attacks all over the world as bubbles/dots based on the latitude/longitude of the cities the attacks originated in. The size of the bubble is based on the number of attacks originating in that city. Cities that don't have attacks for the selected time range (see D) are greyed out.

The map can be zoomed in and out, and can also be dragged around, to focus on areas of interest. The tooltip displayed when hovering over a bubble display the total number of attacks for that city.



B. Total number of attacks for current query

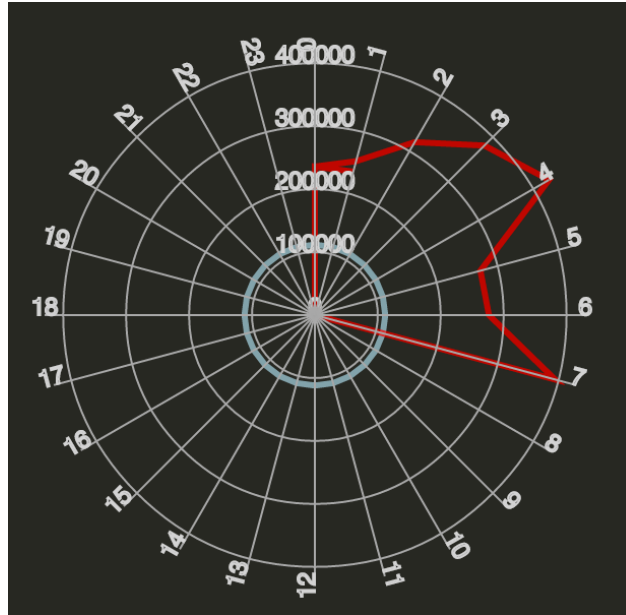
This displays the total number of attacks based on the current query. The data is shown in red, consistent with our other data elements on the page. This is a read-only widget and is updated as the query updates.

Number of Attacks

2,561,015

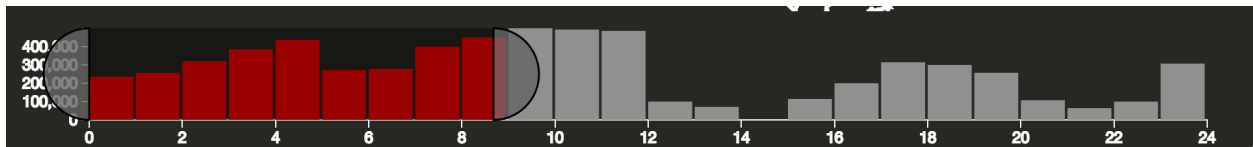
C. Polar plot showing distribution of those attacks over time

The polar plot displays the distribution of attacks based on the current query. The data is aggregated hourly and shown in red. The mean value is displayed using an overlaid circle. The hours are based on the timezone of the viewer, as opposed to the timezone of the attacker. The data is read-only, and updated as the query updates.



D. Time range selector to modify time selection

The time range selector lets users modify the visualization by focusing on specific time ranges. The height of the bars encode the number of attacks based on the current region and network.



Users can drag and select the time range they are interested in. The hours selected in this time range are shown in red while those outside the range are greyed out. The use of the large selectors provide a strong visual indicator that they are interactive.

E. Region filter

The region filter lets the user filter the data and look at attacks originating in specific regions in the world. This updates all widgets (A - D) based on the selection.

Region

F. Network filter

The network filter lets the user filter the data and look at specific networks. This updates all widgets (A - D) based on the selection.

Network

Evaluation

What did we learn:

- We can argue that we see correlation between the time of the day and the origin of the attacks. For instance, countries in the East seem to have a higher number of attacks relative to the ones in the West early in the day. This can be attributed to the differences in timezones around the world, and when businesses are operating.
- Some cities can stand out from the rest (e.g. Seattle) and that can be an indication of a real attack or a misconfigured Akamai server. We would need additional data to prove whether Seattle itself was really attacking the network or whether it was a network configuration issue.

How did we answer our questions:

- The bubble map tells us where the most attacks originate.
- The polar chart, along with the time range selector, tells us which hours of the day have the highest number of attacks.
- The region (country) and network (ISP) selectors further let the user slice-and-dice the data, and visualize data for different parts of the world.

How well does our visualization work:

Overall the visualization has been working very well. The performance and rendering of the visualization has been great. All transitions are working as expected, and show the changes in data as filters change. The flow of the page, and sparse use of color have really improved the user experience.

How could we improve it further:

- By incorporating external data sources that publish news around networking security issues, we could annotate the time range selector and provide additional context for the spikes in attacks.
- By incorporating data for multiple days, we could enhance the visualization by presenting the average (mean/median) number of attacks for any given hour.
- We could also show the hours based on the user's timezone instead of using GMT.
- We could normalize the hours of the attack data by using the timezone of the attack origin for each attack, rather than using the user's timezone.
- We could add tables to show the top 10 countries, networks, operating systems, etc. for the filtered data.
- The time range selector currently only allows for time ranges that do not cross midnight.
- The polar plot could be updated use an arc, rather than a circle, when displaying the median value for time range filtered data.

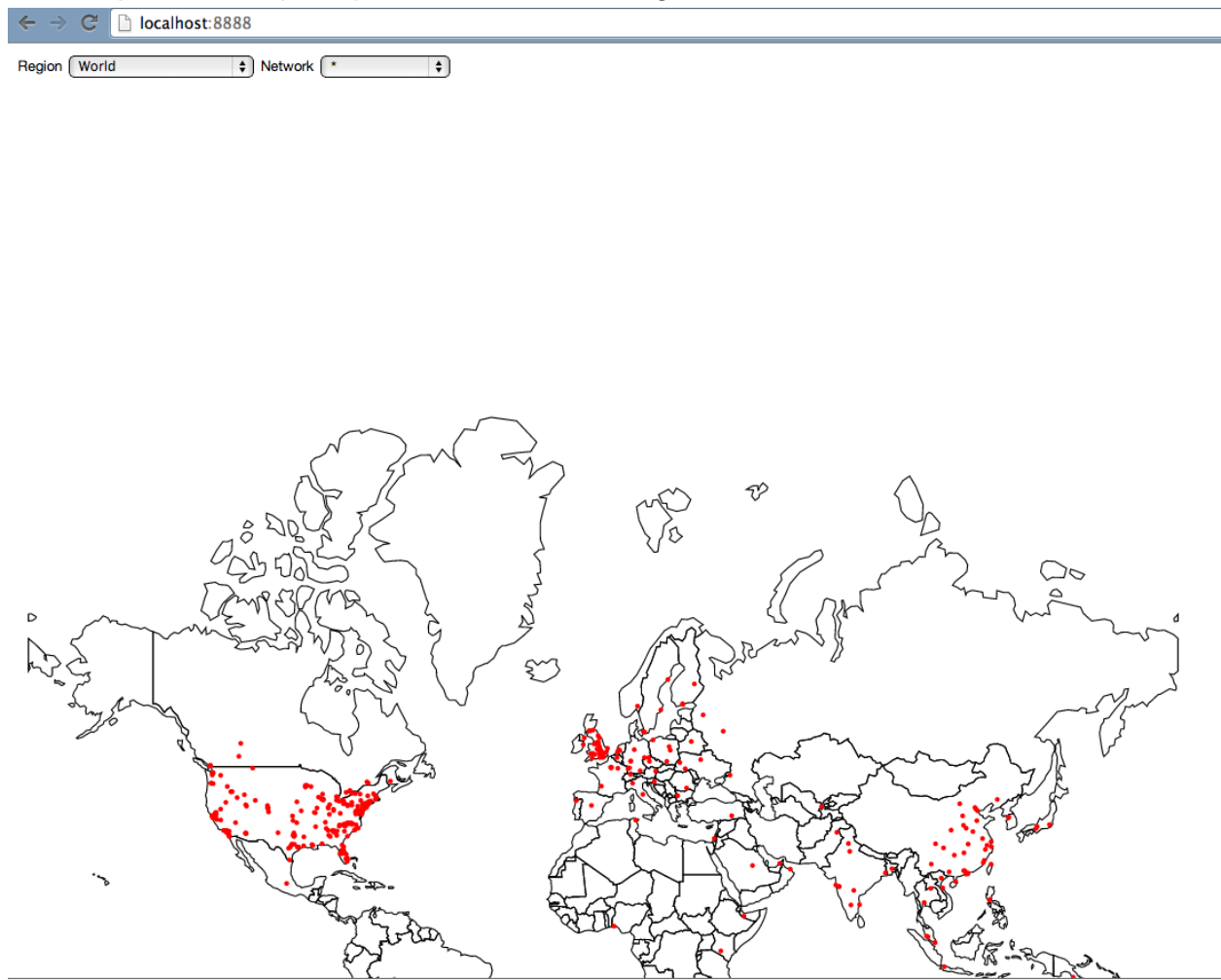
Appendix

3/30

- * Obtain the data
- * Initial work on data parsing
- * Create the initial page

3/31

- * Add basic selectors for region and network
- * Make space for the polar plot with better bounding boxes



4/1

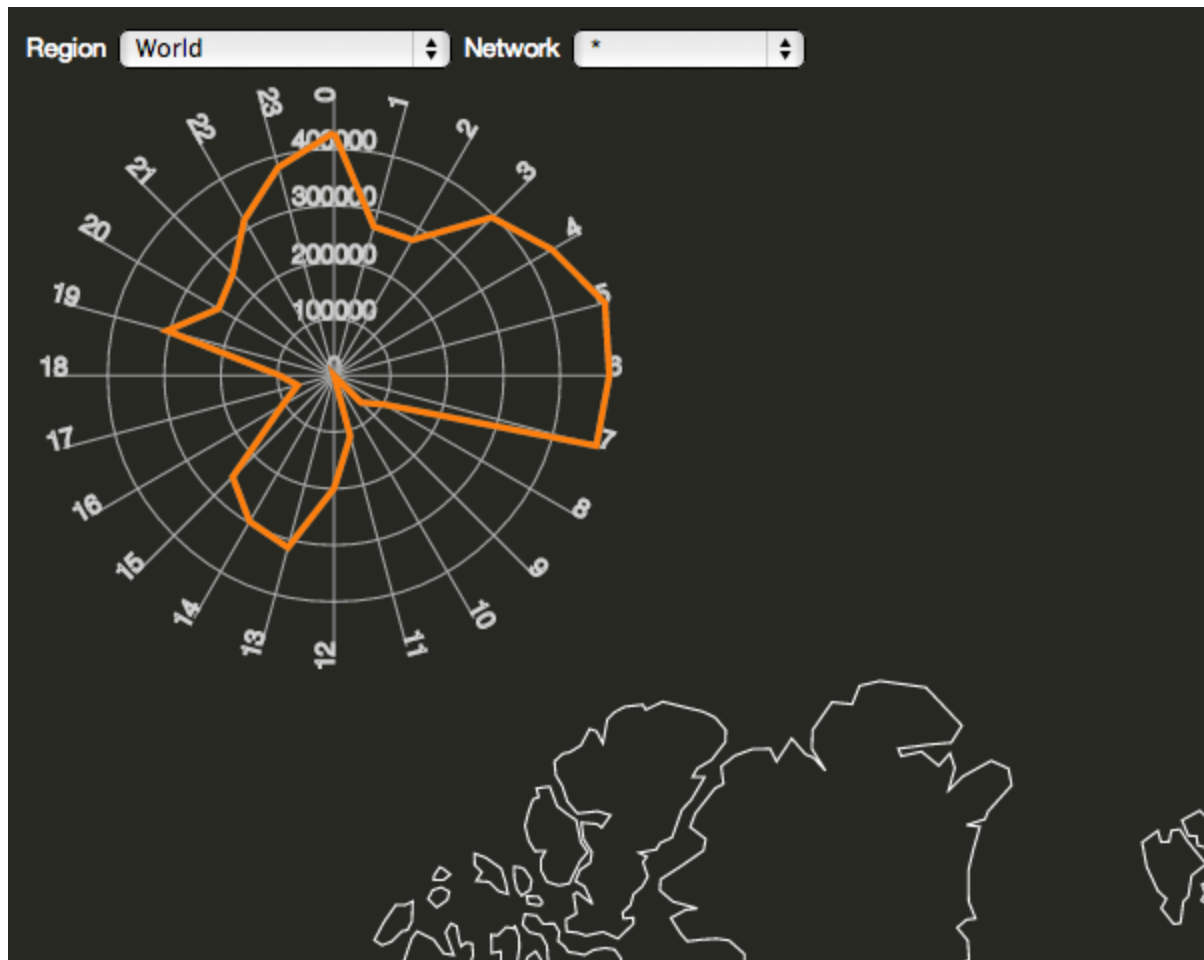
* Changing the aesthetics

* Make bubbles proportional to data



4/1

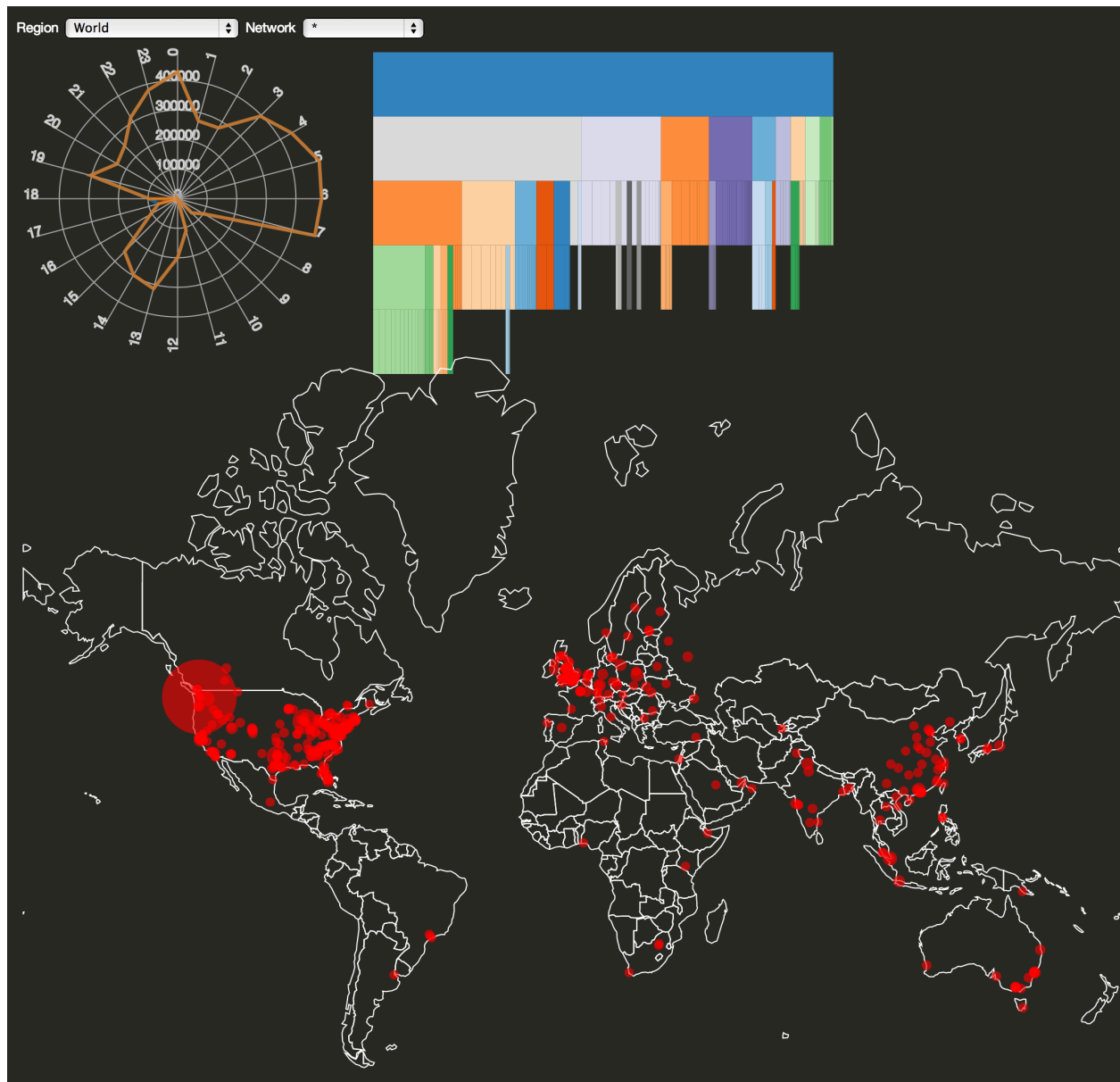
* Polar plot is now supports filtered data and visual transitions on data changes



4/2

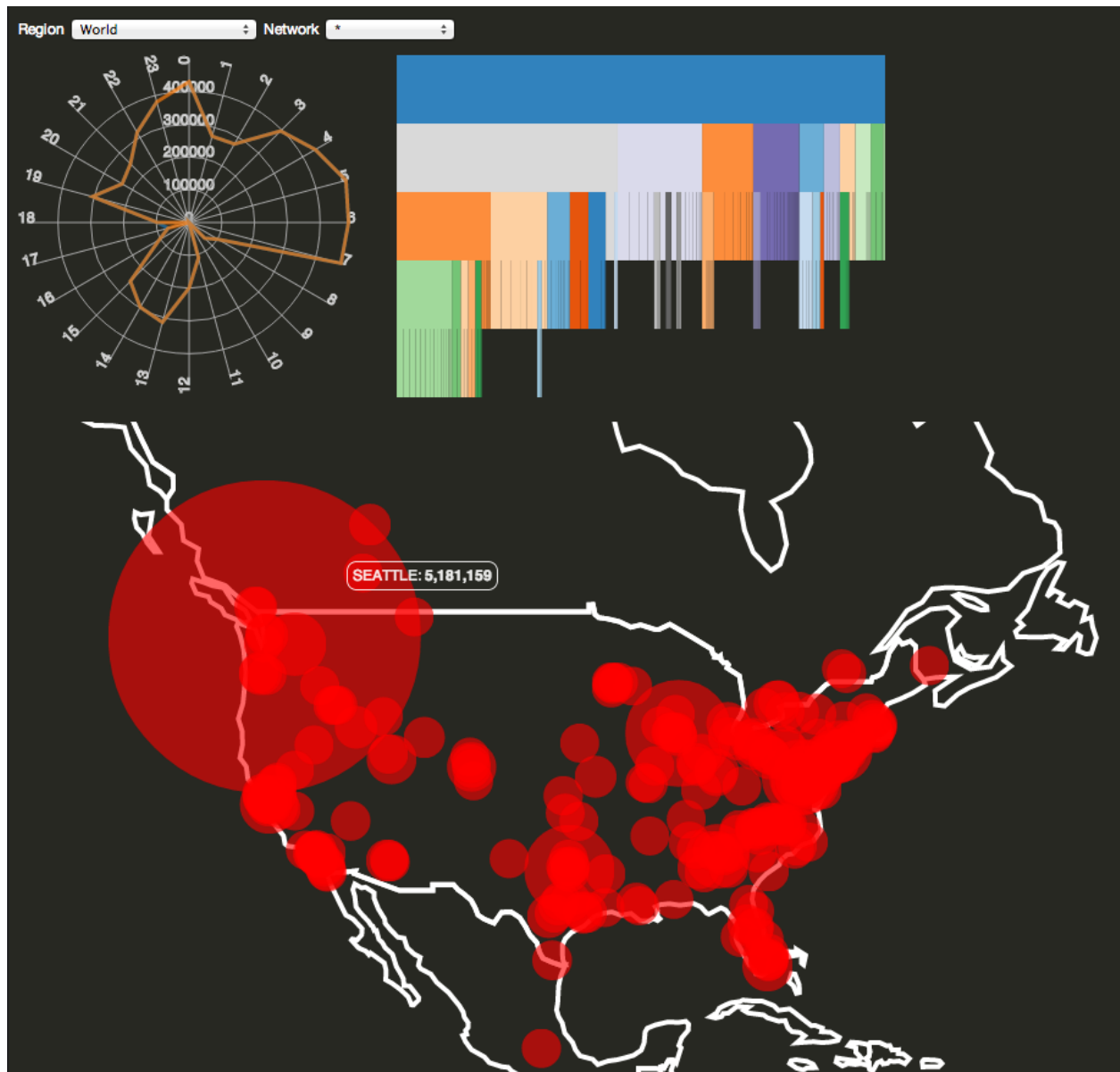
* Added icicle plot with mock data

* Refactored to use controller.js



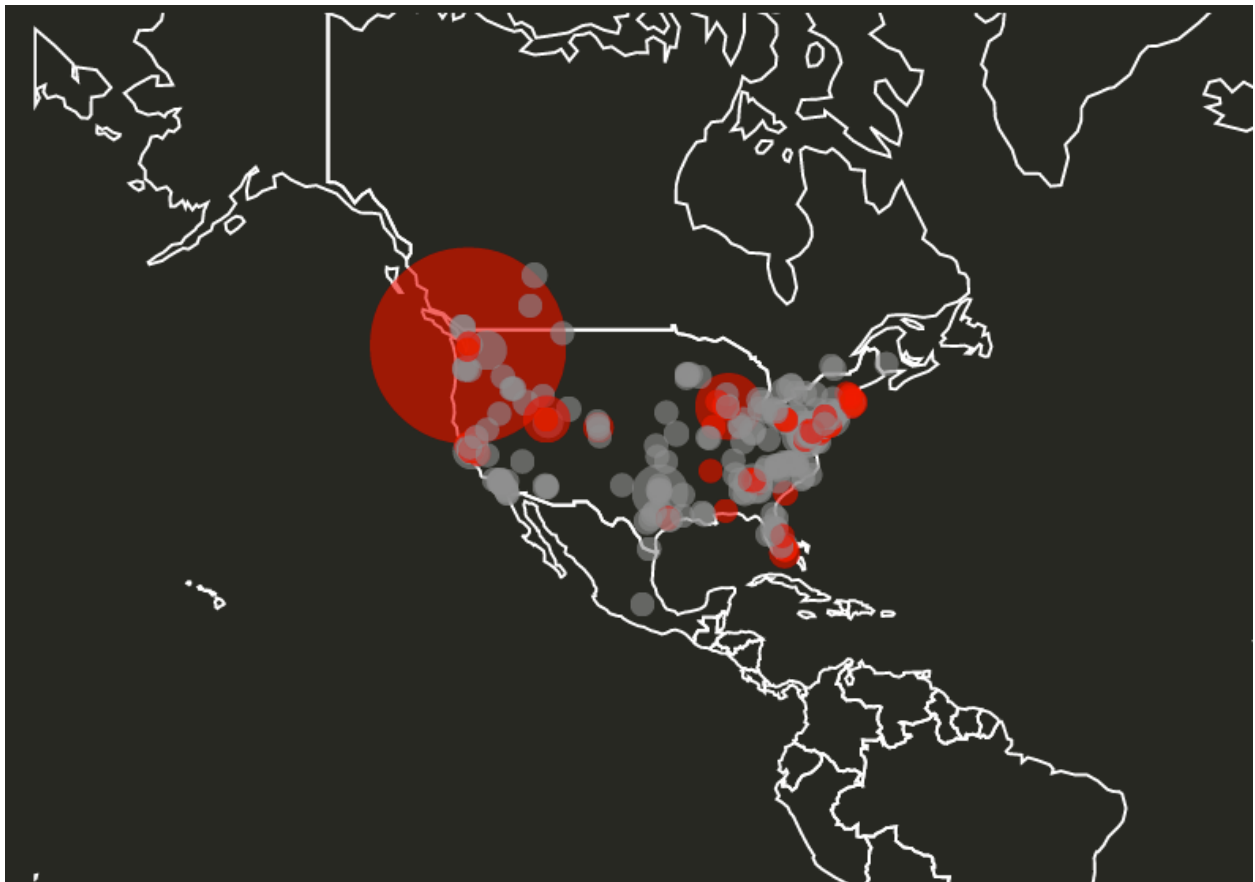
4/2

- * Map zooming, clipping, and labels are now working
- * General clean-up of polar plot
- * All live data linked to filters



4/3 Perhaps it makes more sense to show all attacks on the map and grey out the ones that are not relevant to the current selection

E.g. USA + Comcast:



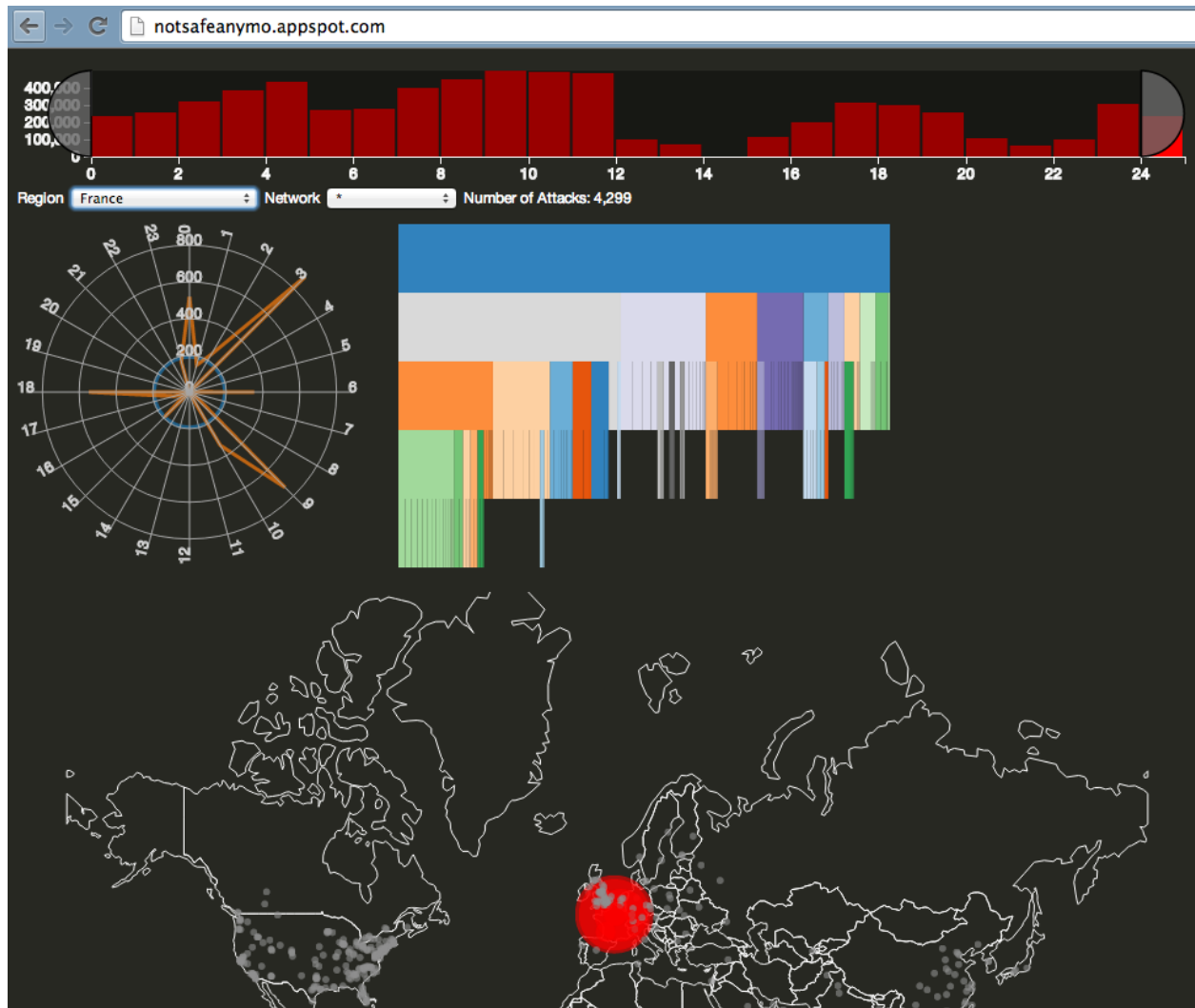
4/5 - It would be nice to display an overall count of the number of attacks



It will also be nice to update this count using animation similar to what they use to update airline schedule boards

4/6 - A bunch of updates this weekend:

- * Final fixes to polar plot
- * Addition of histogram with brushes for time of day (still connecting to data)
- * Add number of attacks in filtered dataset to page



* the time range selector is now hooked up to the drop-downs for region and network and updates as those fields update

