

DataEng: Project Assignment 3

Data Integration

Genevieve LaLonde
Data Engineering Winter 2021
Bruce Irvin

Assignment date: February 16

Due date: February 28, 2021 @10pm PT

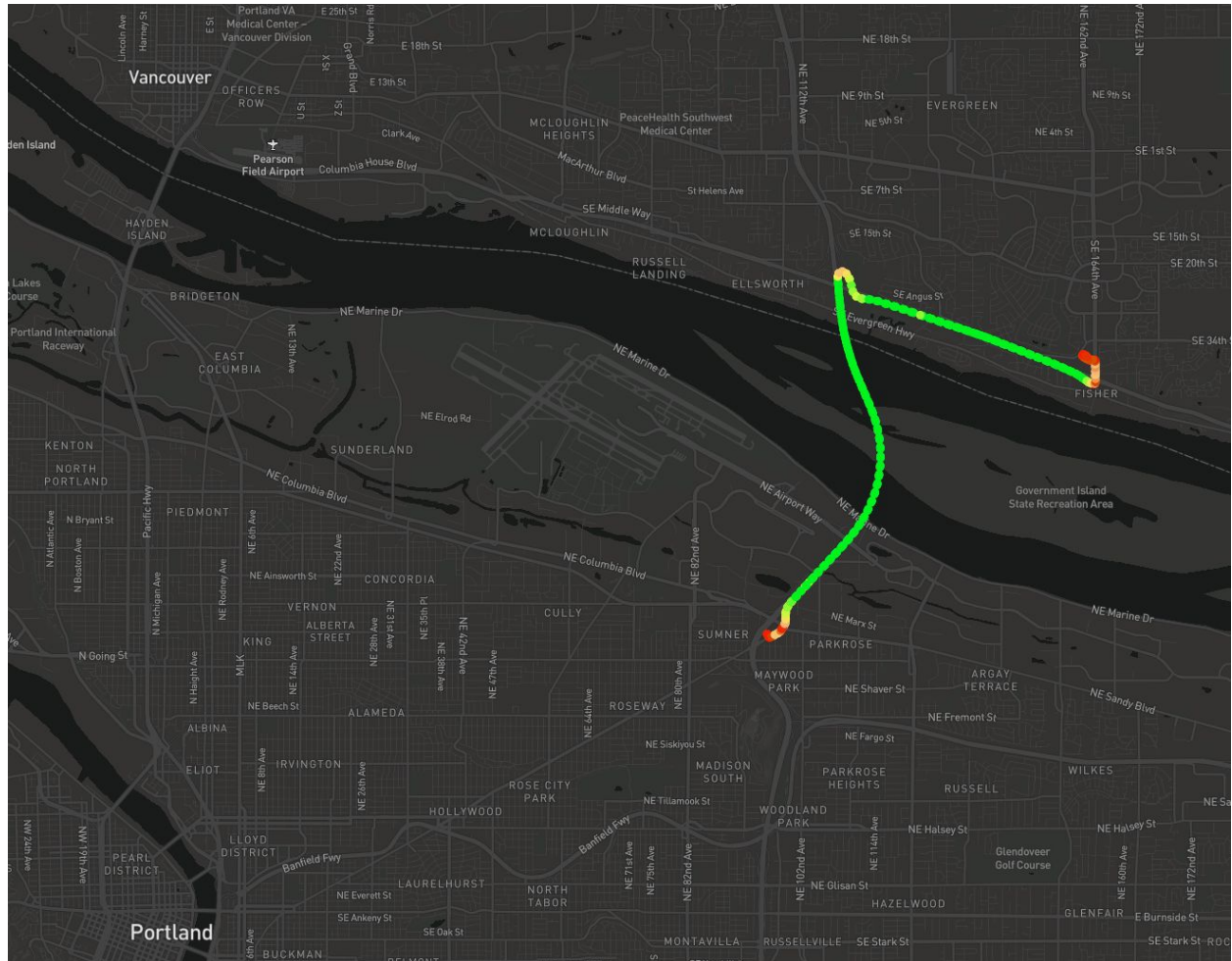
Submit: [assignment submission form](#)

Submission

Visualization 1.

A visualization of speeds for a single trip for any bus route that crosses the Glenn Jackson I-205 bridge. You choose the day, time and route for your selected trip. To find a trip that traverses this bridge, consider finding a trip that includes breadcrumb sensor points within this bounding box: [45.592404, -122.550711, 45.586158, -122.541270]. Any bus trip that includes breadcrumb points within that box either crosses the bridge or goes swimming in the Columbia river!

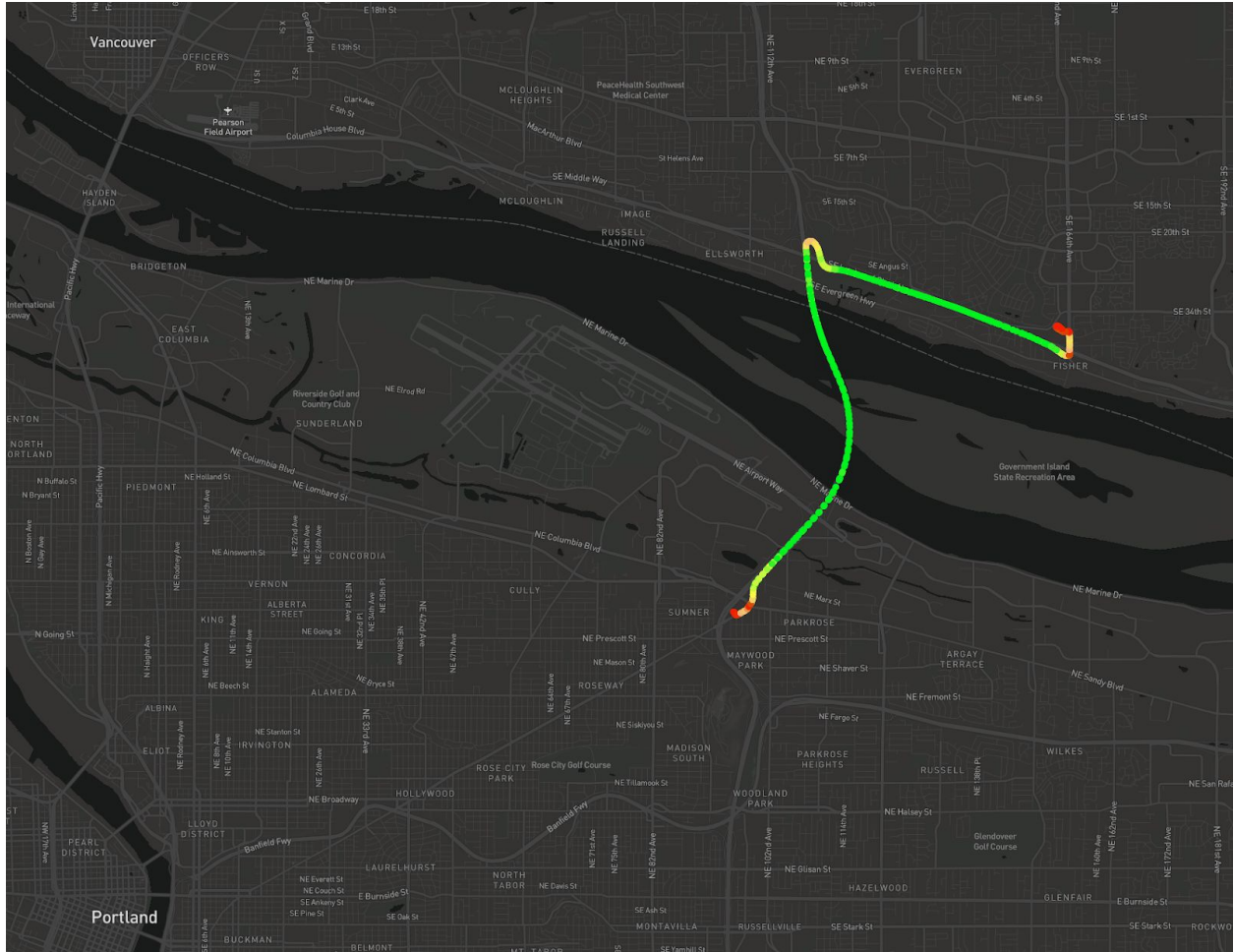
- One trip
- On route 164
- Which crosses the Glenn Jackson I-205 bridge
- On a weekday schedule
- Returning to the start
- On 2020-10-19
- Starting at 15:34:47
- Ending at 15:59:06



Visualization 3.

All outbound trips for route 65 on any Sunday morning (you choose which Sunday) between 9am and 11am.

- All outbound trips
- On route 65
- On Sunday 2020-10-25
- 9:00 am - 11:00 am



Visualization 4.

The longest (as measured by time) trip in your entire data set. Indicate the date, route #, and trip ID of the trip along with a visualization showing the entire trip.

The longest duration trip I have full data for is: 170770658.

A lot of the other very long trips are from before we started ingesting Stop Event data.


```

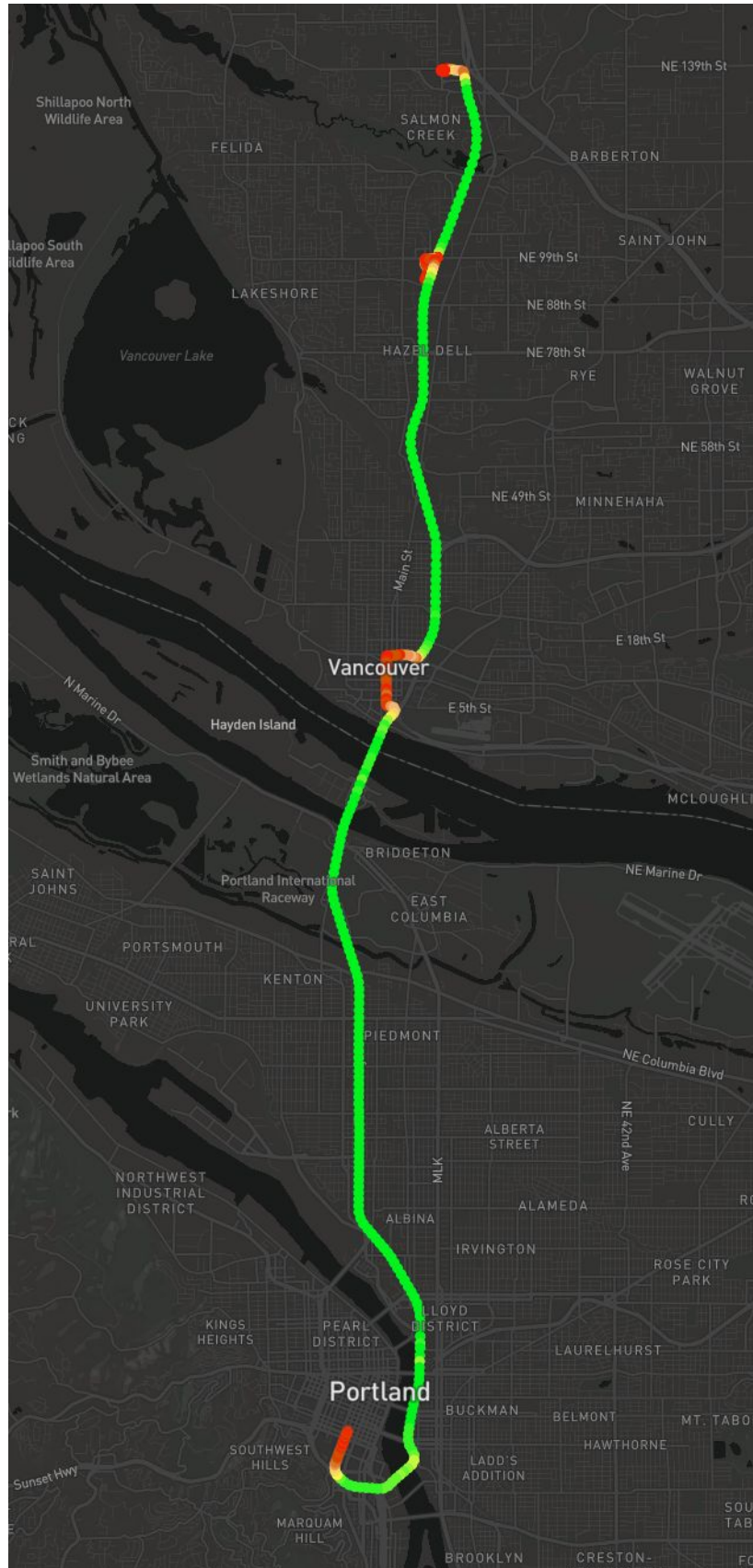
ctran=> select t.*, ib.*
ctran=> from trip t inner join (
ctran(> select max(tstamp), min(tstamp), max(tstamp) - min(tstamp) as trip_length, trip_id
ctran(> from breadcrumb b
ctran(> group by trip_id
ctran(> ) ib on t.trip_id = ib.trip_id
ctran=> order by ib.trip_length desc
ctran=> limit 20;

```

trip_id	route_id	vehicle_id	service_key	direction	max	min	trip_length	trip_id
169302880		1254260			2020-10-01 22:02:44	2020-10-01 16:30:18	05:32:26	169302880
169026975		2269			2020-09-29 10:14:19	2020-09-29 07:02:14	03:12:05	169026975
168692921		2251			2020-09-24 14:06:46	2020-09-24 11:16:41	02:50:05	168692921
167674081		4026			2020-09-11 12:34:30	2020-09-11 09:50:21	02:44:09	167674081
170770658	105	1254260	Weekday	Back	2020-10-20 19:23:25	2020-10-20 17:27:30	01:55:55	170770658
167690539		4032			2020-09-11 13:09:26	2020-09-11 11:17:18	01:52:08	167690539
167688567		4020			2020-09-11 12:51:28	2020-09-11 11:00:11	01:51:17	167688567
167859443		4029			2020-09-14 09:55:47	2020-09-14 08:04:49	01:50:58	167859443
168094965		6001			2020-09-17 23:11:42	2020-09-17 21:20:58	01:50:44	168094965
170809733	60	2291	Weekday	Back	2020-10-20 11:47:22	2020-10-20 09:59:34	01:47:48	170809733
168425778		2220			2020-09-21 17:48:10	2020-09-21 16:01:00	01:47:10	168425778
167954916		4025			2020-09-15 09:49:51	2020-09-15 08:05:25	01:44:26	167954916
167132327		4030			2020-09-04 17:41:18	2020-09-04 16:00:59	01:40:19	167132327
171423657	47	2231	Weekday	Out	2020-10-28 06:51:37	2020-10-28 05:13:25	01:38:12	171423657
168124157		2232			2020-09-17 17:38:32	2020-09-17 16:00:52	01:37:40	168124157
170239217		2243			2020-10-13 06:52:09	2020-10-13 05:14:47	01:37:22	170239217
168621282		2220			2020-09-23 18:00:06	2020-09-23 16:22:48	01:37:18	168621282
169466935		2215			2020-10-02 19:23:40	2020-10-02 17:46:55	01:36:45	169466935
171329672	47	2248	Weekday	Out	2020-10-27 06:49:11	2020-10-27 05:13:08	01:36:03	171329672
168499819		2248			2020-09-22 06:47:53	2020-09-22 05:13:11	01:34:42	168499819

(20 rows)

All datapoints on the longest duration trip are mapped below.



Visualization 5a, 5b, 5c, Three or more additional visualizations of your choice. Indicate why you chose each particular visualization.

5a) Longest Route - Distance

What is the longest route in terms of estimated distance traveled?

Map the longest trip for the longest route. Use a bounding box of the movement on the x and y axis to estimate length of the route.

I've defined the route length by getting the length of the hypotenuse of the bounding box of the route's maximum positions in either axis.

$\text{length_estimate} = \sqrt{(\max(x) - \min(x))^2 + (\max(y) - \min(y))^2}$

When I first did this I was only multiplying the movement of the x and y axis. But after considering whether multiplying or adding would be a better representation (the area or half the perimeter of the bounding box), I realized the hypotenuse is much more accurate, and not weighted by the shape of the box (long and skinny vs square), nor by whether the box is more aligned with the x/y axis. However this is still just an estimate. For a more precise measurement of bus distance, I think we could convert the Geo Points of a trip to a Line String and pull its length. This would be a good area for future work.

Query Description:

Calculate the hypotenuse of the bounding box of points in the trip as the trip length. Order the result by the length of the trip, with longest first. Also return the trip info like the route id.

Query Text:

```
select t.*, ib.*
from trip t inner join (
  select
    min(latitude) S,
    max(latitude) N,
    max(latitude) - min(latitude) as y_movement,
    min(longitude) E,
    max(longitude) W ,
    max(longitude) - min(longitude) as x_movement,
    sqrt((max(latitude) - min(latitude))^2 + (max(longitude) - min(longitude))^2) as total_movement,
    trip_id
  from breadcrumb
  group by trip_id
  order by 7 desc
) ib
on t.trip_id = ib.trip_id
where route_id > 0
```

```

and total_movement > 0
order by ib.total_movement desc
;

```

Data Selection:

From the result, I picked the longest trip, for the top 3 routes occurring in the data. This redacted result is below.

trip_id	route_id	vehicle_id	service_key	direction	s	n	y_movement	e	w	x_movement	total_movement	trip_id
171076865	47	2267	Weekday	Out	45.655968	45.86684	0.21087200000000195	-122.585415	-122.385082	0.20033300000000054	0.29086132653379865	171076865
170817612	185	1254280	Weekday	Back	45.517433	45.72181	0.20435699999999457	-122.677218	-122.543822	0.13339599999999852	0.2440415461862983	170817612
170872807	190	2251	Weekday	Out	45.498038	45.721833	0.22377499999999628	-122.689352	-122.652702	0.03664999999999452	0.2267564180456881	170872807

The points of most interest to our map are:

```

trip_id | route_id
-----+-----
171076865 | 47
170817612 | 105
170872807 | 190

```

I mapped each of these separately so we can clearly see each route without overlapping.

Queries to pull the data:

```

select longitude, latitude, speed from
  breadcrumb where trip_id = 171076865;
select longitude, latitude, speed from
  breadcrumb where trip_id = 170817612;
select longitude, latitude, speed from
  breadcrumb where trip_id = 170872807;

```

Maps:

On following pages, in this order:

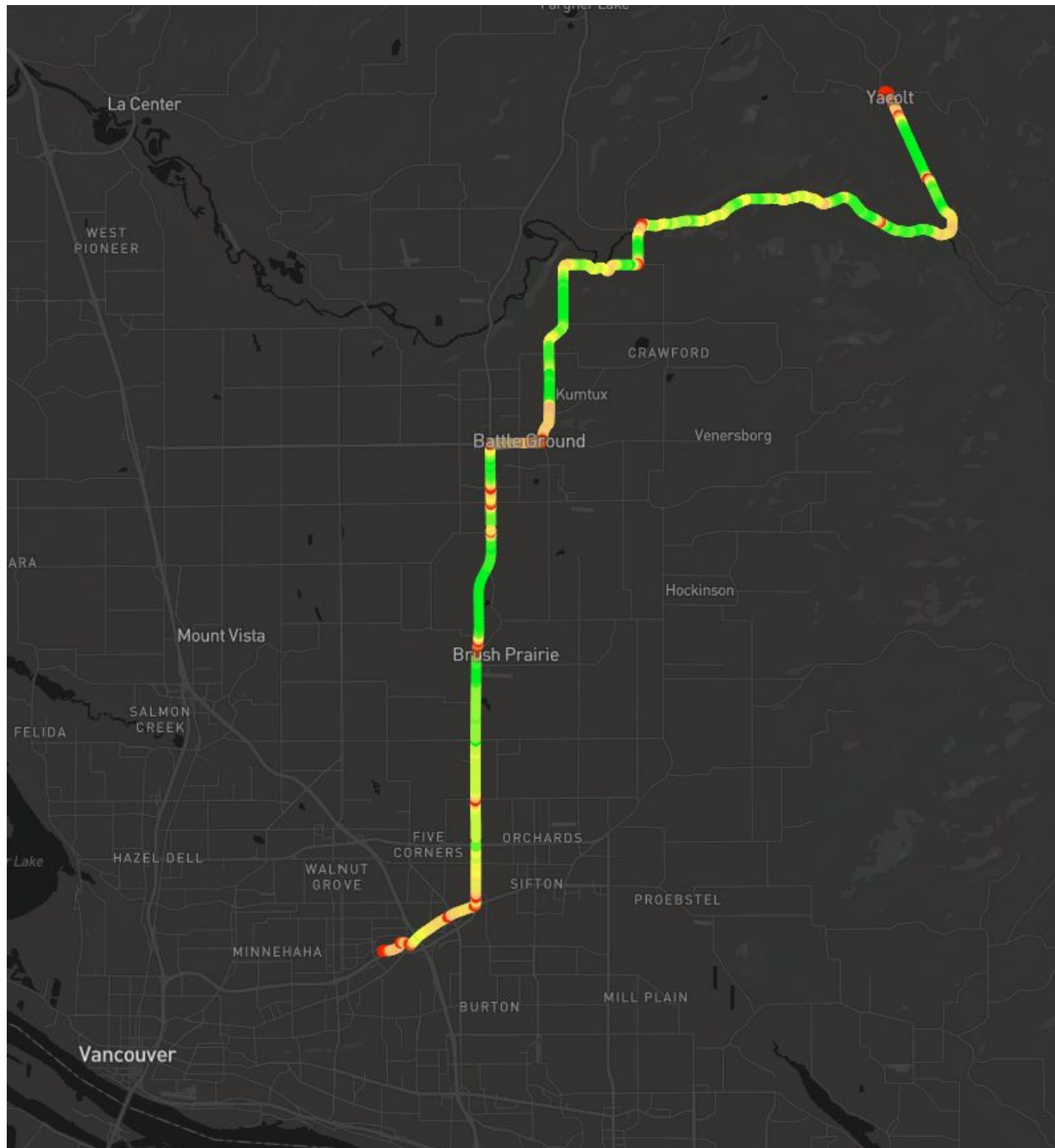
```

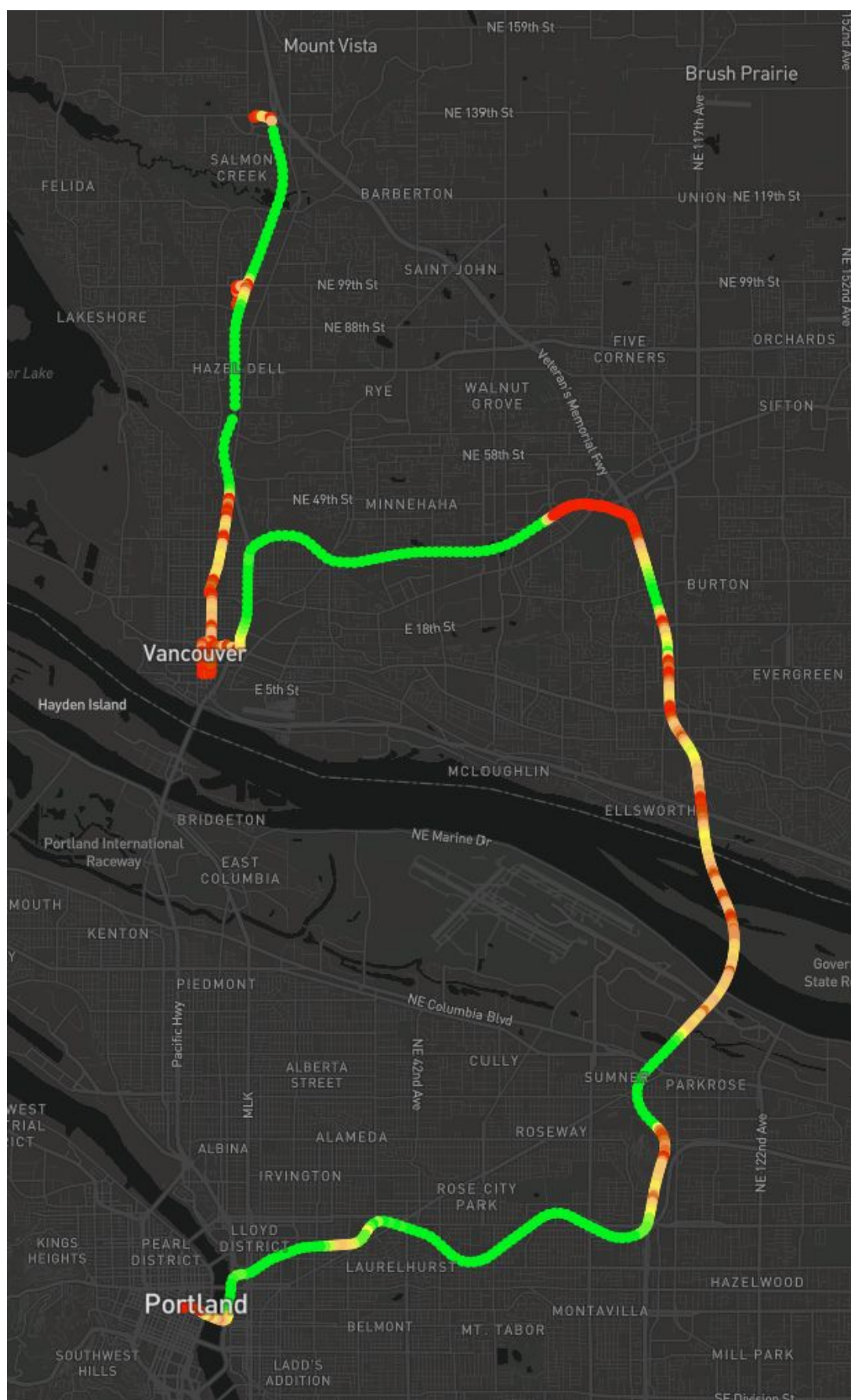
trip_id | route_id
-----+-----
171076865 | 47
170817612 | 105
170872807 | 190

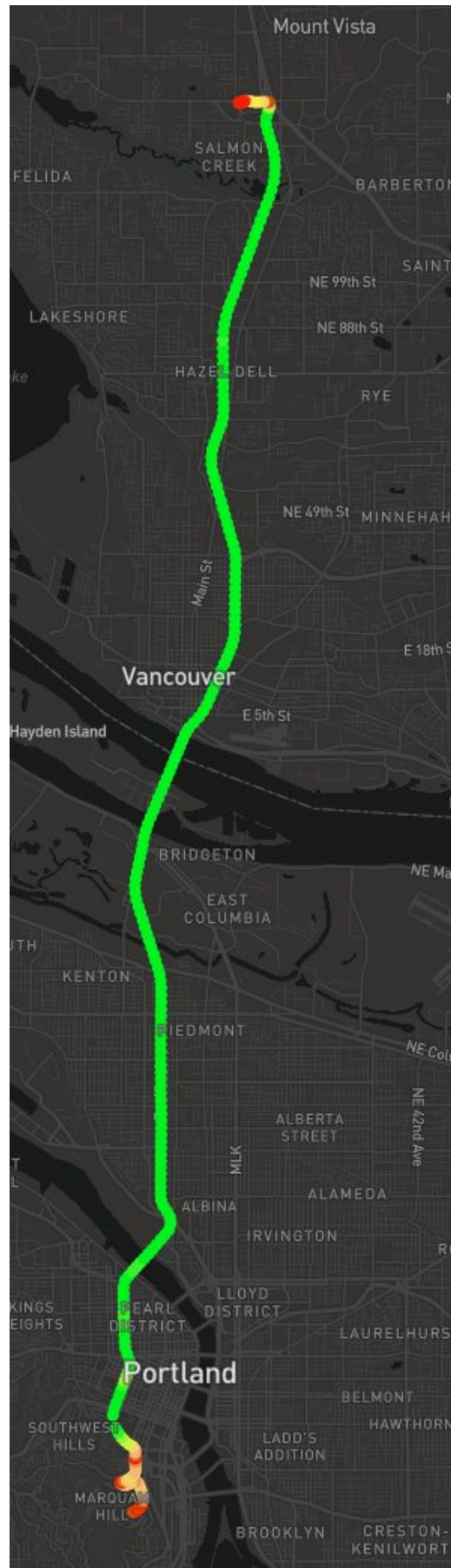
```

Improvements:

Now that I look at the maps. For most of them you could simply use the start/end point of the route instead of calculating the min/max of the x and y coordinates. That would be a suitable simplification. However this is more precise like this.







5b) Shortest Route - Distance

What is the shortest in terms of estimated distance traveled?

Map the shortest trips.

At first glance this seems easy. Simply use the same query as in 5a, just filter it in the other direction.

Query Description:

Calculate the hypotenuse of the bounding box of points in the trip as the trip length. Order the result by the length of the trip, with shortest first. Also return the trip info, like the route id.

Query Text:

```
select t.*, ib.*
from trip t inner join (
  select
    min(latitude) S,
    max(latitude) N,
    max(latitude) - min(latitude) as y_movement,
    min(longitude) E,
    max(longitude) W ,
    max(longitude) - min(longitude) as x_movement,
    sqrt((max(latitude) - min(latitude))^2 + (max(longitude) - min(longitude))^2) as total_movement,
    trip_id
  from breadcrumb
  group by trip_id
  order by 7 desc
) ib
on t.trip_id = ib.trip_id
where route_id > 0
and total_movement > 0
order by ib.total_movement
;
```

Example Result:

trip_id	route_id	vehicle_id	service_key	direction	s	n	y_movement	e	w	x_movement	total_movement	trip_id
171243872	190	2246	Weekday	Out	45.496162	45.496388	0.000146000000000006766	-122.60237	-122.60229	0.000000001118274e-05	0.00016648123817938844	171243872
171194612	41	2231	Weekday	Out	45.63264	45.63267	1.6399999999999999169e-05	-122.67083	-122.67024	0.0002639999999999999908015	0.00026354803694835207	171194612
171425201	190	2242	Weekday	Out	45.496112	45.49639	0.000270000000000015540	-122.60235	-122.60228	0.00000000172259e-05	0.00029378393421219213	171425201
171448991	50	6004	Weekday	Out	45.65625	45.656347	9.6399999999999999966307e-05	-122.58402	-122.583957	0.00032580000000000368345	0.000339166257192842	171448991
171424072	134	2232	Weekday	Out	45.721458	45.721682	0.000224000000000020056	-122.659967	-122.659703	0.0002640000000000137425	0.00034627253681378513	171424072
170990962	50	6018	Weekday	Out	45.655807	45.656182	0.0002149999999999999910966	-122.587437	-122.587052	0.000384999999999999999438656	0.00044896440512866457	170990962
170978926	72	2902	Weekday	Out	45.656055	45.656227	0.00017199999999999999997201	-122.588027	-122.58741	0.000616999999999999999912079	0.00064852556544451707	170978926
171329689	47	2240	Weekday	Out	45.775588	45.776177	0.000668999999999999999949237	-122.527842	-122.526987	0.0001358000000000021820	0.0006240516466899561	171329689
170876985	80	2289	Weekday	Out	45.656	45.65618	0.00010800000000000029184	-122.588480	-122.58767	0.000737999999999999993511	0.0007596341224547981	170876985
170681624	7	2232	Weekday	Out	45.655763	45.65607	0.00036999999999999999993911	-122.587457	-122.586732	0.000725800000000027512	0.00078732877325580256	170681624
171345414	124	2245	Weekday	Out	45.588998	45.589712	0.0007148000000000001810	-122.60709	-122.60753	0.000666999999999999999919363	0.000887790231189436	171345414
171255642	190	1254202	Weekday	Out	45.498802	45.49916	7.000000000000000000295e-05	-122.6056	-122.60464	0.000059999999999999999920783	0.00006316353750000004	171255642
171355552	74	2421	Weekday	Out	45.656813	45.65623	0.000216999999999999999924557	-122.58844	-122.587492	0.00094080000000000001645	0.0009725180944257857	171355552
171426727	72	2262	Weekday	Out	45.655998	45.656153	0.00015500000000000065655	-122.588372	-122.587377	0.00089958000000000031870	0.001080084965283678	171426727
171422221	80	2245	Weekday	Out	45.655968	45.656235	0.00026700000000000009854	-122.588498	-122.587333	0.0011650000000000003847	0.00119528458049983903	171422221
170623254	50	6009	Sunday	Back	45.655033	45.656025	0.000791999999999999999978173	-122.586123	-122.585935	0.00186000000000000036216	0.0013296194044400372	170623254
170810398	60	2292	Weekday	Out	45.638015	45.631897	0.00102800000000000038021	-122.670513	-122.66966	0.0000538000000000006432	0.00137700005808765462	170810398
170905299	60	2291	Weekday	Out	45.638720	45.631907	0.0011790000000000004052	-122.67052	-122.669802	0.00071799999999999920820	0.0013804228369110528	170905299
170802666	134	2240	Weekday	Out	45.58038	45.589752	0.0013719999999999999999361	-122.60314	-122.602513	0.0006269999999999999994382	0.0015804803611515777	170802666
171457736	60	2294	Weekday	Out	45.638575	45.631853	0.00127799999999999999992242	-122.670585	-122.669620	0.00007780000000000026012	0.001549971933286738	171457736
170704974	50	6004	Weekday	Back	45.655862	45.656125	0.000312999999999999999984534	-122.585628	-122.584025	0.00160380000000000029877	0.0016332721757283382	170704974
171345036	190	1254202	Weekday	Out	45.498325	45.499312	0.00090700000000000028696	-122.684575	-122.683847	0.001527999999999999993125	0.001819592775480593	171345036
171457732	60	2294	Weekday	Back	45.628452	45.63844	0.000193999999999999999972142	-122.670553	-122.670427	0.0001259999999999999999451938	0.001991988955578955	171457732
171338018	60	2293	Weekday	Back	45.628438	45.638438	0.0019999999999999999999539	-122.670563	-122.67045	0.000112999999999999999901952	0.002008189766438407	171338018
170905294	60	2291	Weekday	Out	45.628497	45.638987	0.00209999999999999999990513	-122.670572	-122.670447	0.0001258000000000001255	0.00201500230531409235	170905294
171522296	134	2230	Weekday	Out	45.512378	45.514165	0.001787000000000000008205	-122.681258	-122.681255	0.00093380000000000034033	0.0020159812072675796	171522296
171361154	50	2291	Weekday	Back	45.628478	45.638985	0.00202699999999999999992247	-122.67856	-122.678438	0.0001219999999999999999840756	0.0020386681166528573	171361154
171566814	50	6003	Weekday	Back	45.655813	45.656473	0.00063999999999999999993382	-122.58589	-122.583963	0.00193699999999999999990791	0.002046355852789816	171566814
170806553	72	2267	Weekday	Out	45.65573	45.656178	0.00044799999999999999996717	-122.580238	-122.585235	0.00200380000000000019754	0.0020524094648420857	170806553
170785513	80	2285	Weekday	Out	45.65573	45.656113	0.00030299999999999999993561	-122.587493	-122.58547	0.00202799999999999999994113	0.00205893613849223	170785513
171203157	60	2293	Weekday	Back	45.62841	45.638407	0.00207699999999999999990845	-122.670577	-122.670445	0.000131999999999999999935017	0.002081198284280784	171203157
171119328	74	2246	Weekday	Out	45.655738	45.655968	0.00023800000000000015807	-122.587385	-122.585233	0.002071999999999999999902574	0.0020847263689389784	171119328
171120834	60	2293	Weekday	Back	45.628203	45.63838	0.00209699999999999999991273	-122.670567	-122.67045	0.00011780000000000030135	0.0021002614123801567	171120834
170810394	60	2292	Weekday	Back	45.62020	45.638453	0.002163000000000000004226	-122.670545	-122.670515	0.00000000000000002241e-05	0.00216318083430317264	170810394
170809934	60	2291	Weekday	Back	45.628277	45.638443	0.002166800000000000005540	-122.670553	-122.670437	0.000115999999999999999991345	0.002169183962471384	170809934
170567074	60	1290380	Saturday	Back	45.626345	45.62854	0.002158000000000000008391	-122.670582	-122.670580	7.3999999999999999999987e-05	0.002196247826179297	170567074
170633158	164	2248	Weekday	Out	45.58841	45.51839	0.001908000000000000002014	-122.683132	-122.682128	0.001083999999999999999946758	0.002220803683681126	170633158
170568006	60	4015	Saturday	Out	45.628275	45.638480	0.00223800000000000000752	-122.670532	-122.670393	0.000130999999999999999991925	0.0022273414646166307	170568006
171091671	50	6003	Weekday	Back	45.655808	45.656342	0.0005340000000000000018100	-122.5862	-122.584032	0.00216880000000000011716	0.0022327964528932624	171091671
170901573	72	2268	Weekday	Out	45.655715	45.656195	0.0004799999999999999999683915	-122.588465	-122.586283	0.00218200000000000047912	0.0022341718023799357	170901573
170877257	60	2293	Weekday	Back	45.628237	45.638473	0.0002236000000000000034574	-122.670558	-122.670437	0.000128999999999999999929322	0.0022392715333370695	170877257
171395538	47	2902	Weekday	Out	45.655735	45.656083	0.000348000000000000245745	-122.587452	-122.585227	0.002249999999999999999957026	0.002252404995817896	171395538
170621802	60	4030	Sunday	Out	45.628272	45.63853	0.002237000000000000001255	-122.67057	-122.67043	0.00014800000000000008053	0.0022613370783389145	170621802
170998738	60	2292	Weekday	Back	45.628167	45.638433	0.00226599999999999999990769	-122.670555	-122.670412	0.00014299999999999999999423108	0.0022785076524848616	170998738
171435768	190	2403	Weekday	Out	45.642088	45.64435	0.0015420000000000000008984	-122.680713	-122.599013	0.001699999999999999999995987	0.002295168996531715	171435768
170969487	164	2245	Weekday	Out	45.58845	45.510532	0.00208199999999999999994366	-122.683125	-122.682142	0.000983800000000000005832	0.002302392885670579	170969487
171392220	72	2401	Weekday	Out	45.655733	45.656215	0.000482000000000000002805	-122.580277	-122.580002	0.00273980000000000011573	0.002325499731252978	171392220
170685686	72	2269	Weekday	Out	45.655748	45.656195	0.00044699999999999999994091	-122.588312	-122.585983	0.002329800000000000134	0.002371587959862626	170685686
170716794	6	2901	Weekday	Out	45.637377	45.63985	0.00167299999999999999996785	-122.655733	-122.654025	0.001707999999999999999936835	0.0023980561228077117	170716794
170684224	6	2262	Weekday	Out	45.638448	45.638673	0.000275800000000000003638	-122.658680	-122.648197	0.00241180000000000009267	0.0024214759961735835	170684224
170802137	41	1776	Weekday	Out	45.58137	45.594807	0.00074999999999999999994757	-122.586242	-122.585927	0.002349999999999999999995848	0.002435450781611453	170802137
171453993	72	2271	Weekday	Out	45.65573	45.656182	0.000458000000000000007035	-122.587377	-122.585945	0.002431999999999999999990795	0.002473646780727855	171453993
171542879	67	4029	Weekday	Out	45.594187	45.59489	0.000782999999999999999984238	-122.586723	-122.583863	0.0023680000000000001318	0.0024865813573383294	171542879

We can already see, something strange is going on. If this were really a representation of the shortest trips, we would see many of the same route id, the shortest route. Let's map it anyway to check it out. Map the first 50, the "shortest trips".

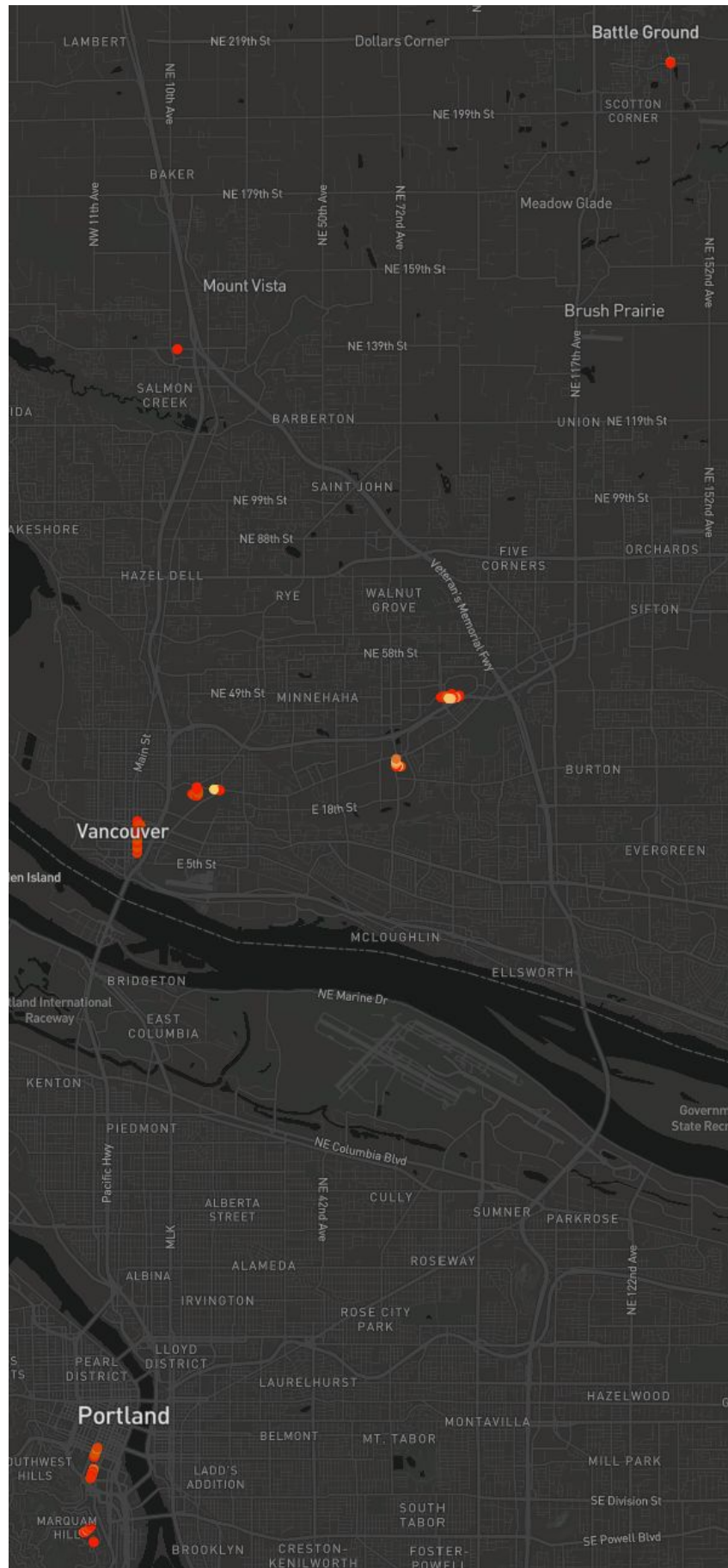
Queries to pull the data:

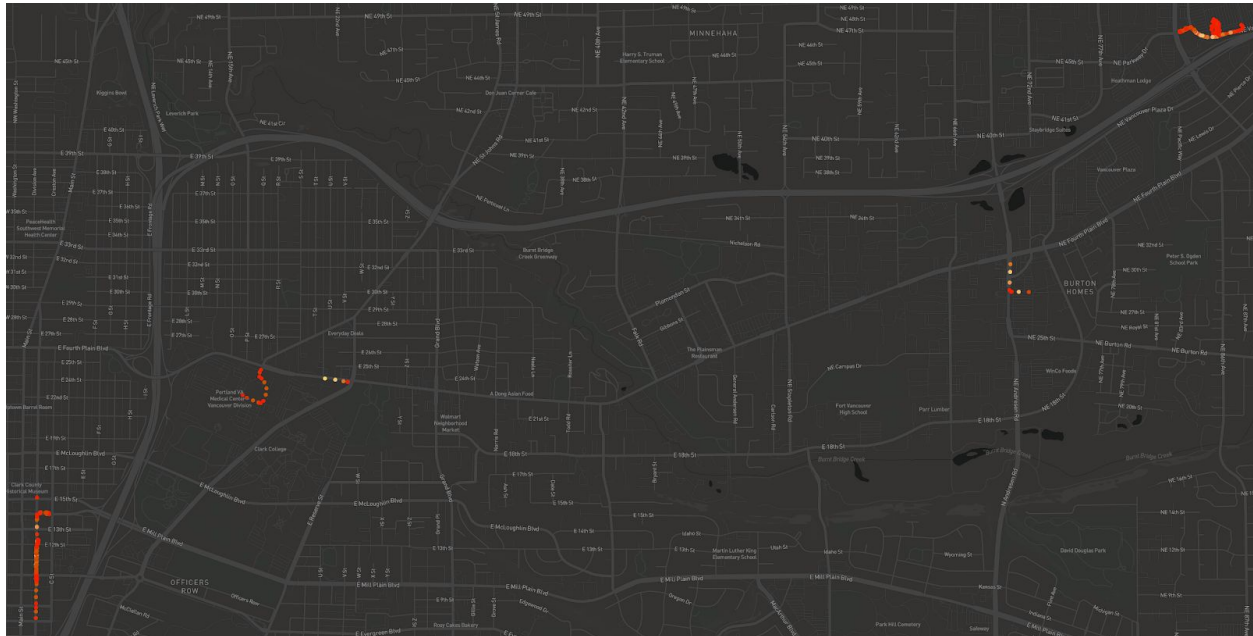
Pull the trip IDs for the first 50:

```
select longitude, latitude, speed from
breadcrumb where trip_id in
(171243872, 171194612, 171425201, 171448991, 171424072, 170990962, 170978926,
171329689, 170876985, 170681624, 171345414, 171255642, 171355552, 171426727,
171422221, 170623254, 170810398, 170905299, 170802666, 171457736, 170704974,
171345036, 171457732, 171338818, 170905294, 171522296, 171361154, 171566814,
170806553, 170785513, 171203157, 171119328, 171120834, 170810394, 170809934,
170567074, 170633158, 170566806, 171091671, 170901573, 170877257, 171339538,
170621802, 170998738, 171435768, 170969487, 171339229, 170685686, 170716794,
170684224)
```

Map:

All the data points, and a zoom in on downtown Vancouver.





Analysis and Future Work:

There is no continuity between these points. They do not represent routes. I'm going to suspect these are error readings, where a trip was unexpectedly cut off. While my first inclination is to toss this and do a different visualization or filter these out, this query could be used for data validation, after the data has been loaded into the database. We could compare the average route length estimate for the route, and for any trips where the length is significantly smaller, those trips and their breadcrumbs should be dropped. A lot of these are in downtown Vancouver or other city centers, so it's possible they represent a bus starting up and shutting down, for example if a driver has a break scheduled at the end of a trip. This does not represent movement in traffic, and they are misrepresenting slow moving (very red/orange due to slow speeds) so we don't want them mucking up our dataset, whose purpose is to investigate traffic movement problems.

To find the actual shortest route, you could do it very precisely by dropping what may be bad datapoints as described above. Alternatively, but a similar concept of data could likely have been gathered by collecting the max length trip, for the route with the shortest average trip. Something like:

5c) Slowest 50 trips

Map the slowest 50 trips.

Slowest is defined as the average speed of the trip. Don't include very shortened trips that may have been miscalculations.

We tried to find shortest trip, we may as well try to find the slowest. First off there is a similar issue, where many datapoints show a speed of 0. So we'll average the speed across the trip. Additionally, now that we know a lot of trips are bad data, let's filter them out. Only keep the trips, where its length is as expected.

Query Description:

What is the trip with the slowest average speed per trip, where the length of the trip is between the max length of trips on that route, and the average length of trips on that route?

I used some CTEs for this to make it easier for me to conceptualize it.

Query text:

```
with greatest_cte as (
  select
    long_routes.route_id,
    max(long_routes.length_est) as max_length
  from (
    select t.*, ib.*
    from trip t inner join (
      select
        sqrt((max(latitude) - min(latitude))^2 + (max(longitude) - min(longitude))^2) as
length_est,
        trip_id
      from breadcrumb
      group by trip_id
    ) ib
    on t.trip_id = ib.trip_id
    where route_id > 0
  ) long_routes
  group by 1
),
middling_cte as
(
  select
    long_routes.route_id,
    avg(long_routes.length_est) as avg_length
  from (
    select t.*, ib.*
    from trip t inner join (
```

```

        select
            sqrt((max(latitude) - min(latitude))^2 + (max(longitude) - min(longitude))^2) as
length_est,
            trip_id
        from breadcrumb
        group by trip_id
    ) ib
    on t.trip_id = ib.trip_id
    where route_id > 0
) long_routes
group by 1
),
this_trip_cte as
(
    select t.route_id, ib.*
    from trip t inner join (
        select
            avg(speed) as trip_speed,
            sqrt((max(latitude) - min(latitude))^2 + (max(longitude) - min(longitude))^2) as
length_est,
            trip_id
        from breadcrumb
        group by trip_id
    ) ib
    on t.trip_id = ib.trip_id
    where route_id > 0
)
select
    tc.trip_id,
    tc.trip_speed,
    tc.route_id,
    tc.length_est,
    gc.max_length,
    mc.avg_length
from
    this_trip_cte tc
    join greatest_cte gc
    on tc.route_id = gc.route_id
    join middling_cte mc
    on gc.route_id = mc.route_id
where tc.route_id in (6, 60, 19)
and tc.length_est <= gc.max_length
and tc.length_est >= mc.avg_length
order by 2;

```


Example result:

trip_id	trip_speed	route_id	length_est	max_length	avg_length
171525855	4.526041666666667	6	0.05356358735186373	0.0936146677236058	0.03389456739974154
171269108	5.197530864197531	6	0.039727726464017316	0.0936146677236058	0.03389456739974154
170879378	5.276923076923077	19	0.052188512088390634	0.08867950993323866	0.051085317247486536
171357318	5.342541436464089	6	0.03994323783820918	0.0936146677236058	0.03389456739974154
171180798	5.394636015325671	19	0.05215056801224857	0.08867950993323866	0.051085317247486536
171161934	5.950920245398773	6	0.03985348800294327	0.0936146677236058	0.03389456739974154
170684028	6.099290780141844	6	0.03975776190884689	0.0936146677236058	0.03389456739974154
170972035	6.233576642335766	6	0.03971924703465229	0.0936146677236058	0.03389456739974154
171453472	6.2388059701492535	6	0.039762912732833625	0.0936146677236058	0.03389456739974154
171226796	6.318181818181818	60	0.040892009402815184	0.17028464015583014	0.040124873466136676
170612966	6.4144736842105265	6	0.039751089658018476	0.0936146677236058	0.03389456739974154
171453579	6.486486486486487	6	0.0393437341047446	0.0936146677236058	0.03389456739974154
171122635	6.494444444444445	60	0.04102605440448915	0.17028464015583014	0.040124873466136676
171453459	6.5225225225225225	6	0.03938004499744239	0.0936146677236058	0.03389456739974154
171452929	6.544776119402985	6	0.039719769309490366	0.0936146677236058	0.03389456739974154
170716771	6.598214285714286	6	0.0394041425360303	0.0936146677236058	0.03389456739974154
170880147	6.612903225806452	6	0.03972066376584783	0.0936146677236058	0.03389456739974154
170684140	6.62962962962963	6	0.03937365044035934	0.0936146677236058	0.03389456739974154
170684016	6.636363636363637	6	0.039366031715182076	0.0936146677236058	0.03389456739974154
171332096	6.724770642201835	6	0.03935969472696099	0.0936146677236058	0.03389456739974154
170684089	6.738095238095238	6	0.039751963787973965	0.0936146677236058	0.03389456739974154
170795828	6.763285024154589	19	0.05224671037491292	0.08867950993323866	0.051085317247486536
170879715	6.770642201834862	6	0.039380877593065904	0.0936146677236058	0.03389456739974154
170880135	6.773584905660377	6	0.039315969325444136	0.0936146677236058	0.03389456739974154
171250837	6.776923076923077	6	0.03973812976221975	0.0936146677236058	0.03389456739974154
170781886	6.79047619047619	6	0.03936851083035415	0.0936146677236058	0.03389456739974154
170879658	6.790909090909091	6	0.03935176166323295	0.0936146677236058	0.03389456739974154
171331978	6.803738317757009	6	0.03936080678841945	0.0936146677236058	0.03389456739974154
170972633	6.841666666666667	6	0.03977123046877937	0.0936146677236058	0.03389456739974154
171453042	6.869158878504673	6	0.039346837357035765	0.0936146677236058	0.03389456739974154
170781180	6.88785046728972	6	0.039362226982219264	0.0936146677236058	0.03389456739974154
170880210	6.895161290322581	6	0.03973975523326968	0.0936146677236058	0.03389456739974154
170683490	6.902912621359223	6	0.039364709385436035	0.0936146677236058	0.03389456739974154
171332049	6.903225806451613	6	0.039759679639048376	0.0936146677236058	0.03389456739974154
171250893	6.9105691056910565	6	0.03970971575068446	0.0936146677236058	0.03389456739974154
171617798	6.9245283018867925	6	0.039427110723462076	0.0936146677236058	0.03389456739974154
170880199	6.951456310679611	6	0.03932570366820631	0.0936146677236058	0.03389456739974154
170781311	7	6	0.03976439800876713	0.0936146677236058	0.03389456739974154
171356667	7.008064516129032	6	0.05356616323763055	0.0936146677236058	0.03389456739974154
170972684	7.009803921568627	6	0.0393609695129587	0.0936146677236058	0.03389456739974154
171269568	7.01	6	0.03939009678841408	0.0936146677236058	0.03389456739974154
171617327	7.024390243902439	6	0.039758134689646445	0.0936146677236058	0.03389456739974154
171437179	7.04	6	0.03935641229837957	0.0936146677236058	0.03389456739974154
171332152	7.049019607843137	6	0.03934743540563274	0.0936146677236058	0.03389456739974154
171617682	7.058333333333334	6	0.03972776339035063	0.0936146677236058	0.03389456739974154
171123269	7.076335877862595	60	0.04092922899346972	0.17028464015583014	0.040124873466136676
170972741	7.076923076923077	6	0.03934681732998686	0.0936146677236058	0.03389456739974154
171357131	7.079207920792079	6	0.0393681735669907	0.0936146677236058	0.03389456739974154
171269824	7.125	6	0.03975122117620888	0.0936146677236058	0.03389456739974154
170781949	7.137254901960785	6	0.039343626332615164	0.0936146677236058	0.03389456739974154
171162273	7.138613861386139	6	0.039458683632383626	0.0936146677236058	0.03389456739974154
171269704	7.142857142857143	6	0.039746454684672854	0.0936146677236058	0.03389456739974154
171269753	7.15	6	0.039379949123901484	0.0936146677236058	0.03389456739974154

Next I did a simple query of the breadcrumb data for the first 50 trips in that list, as seen below. Alternatively you could use the previous query as the in-list, and only select the trip ID out of that subquery. But I liked taking that intermediate step, so I can visually validate results as seen above.

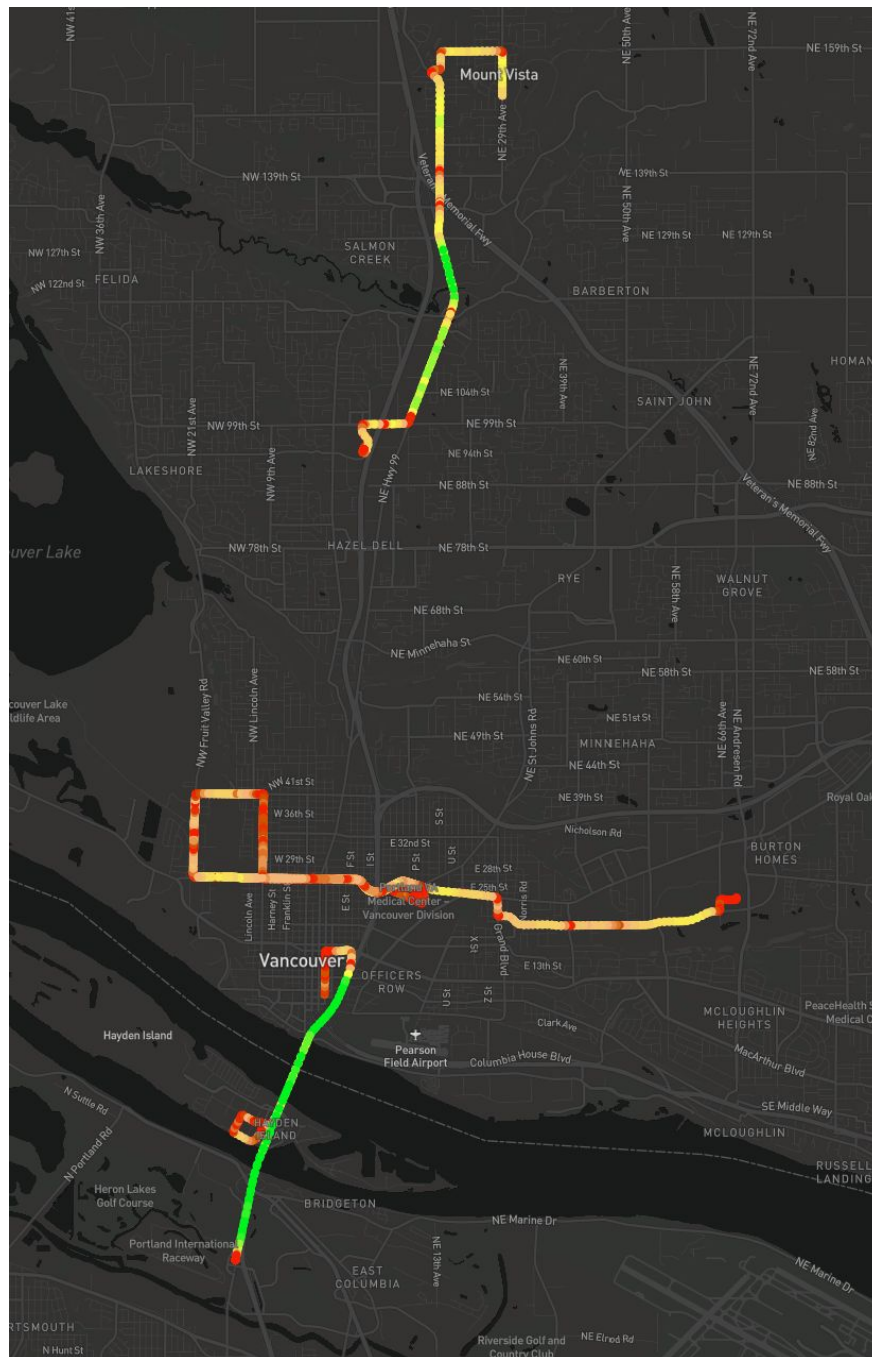
Data Generating Query:

```
select longitude, latitude, speed from
breadcrumb where trip_id in
```

```
(171525855, 171269108, 170879378, 171357318, 171180798, 171161934, 170684028,
170972035, 171453472, 171226796, 170612966, 171453579, 171122635, 171453459,
```


171452929, 170716771, 170880147, 170684140, 170684016, 171332096, 170684089, 170795828, 170879715, 170880135, 171250837, 170781886, 170879658, 171331978, 170972633, 171453042, 170781180, 170880210, 170683490, 171332049, 171250893, 171617798, 170880199, 170781311, 171356667, 170972684, 171269568, 171617327, 171437179, 171332152, 171617682, 171123269, 170972741, 171357131, 171269824, 170781949)

Map:



SQL Performance

I was interested to find that there was a significant performance difference between these 2 CTE below. I ran them on their own as a debugging step before running the full query above. The first one hangs because it takes so long. I thought I would be improving results, by selecting fewer columns so the project is not as large, and not nesting things as much, which was making it easier for me to think of this conceptually. However it seems operating on just one table at a time to pull the trip ids and calculate their length, and subsequently pulling data for those, is a much better plan.

This was very very not performant, it hung.

```
with this_trip_cte as
(
  select
    b.trip_id,
    max(route_id),
    sqrt((max(latitude) - min(latitude))^2 + (max(longitude) - min(longitude))^2) as length_est
  from breadcrumb b
  inner join trip t
  on b.trip_id = t.trip_id
  where route_id > 0
  group by b.trip_id
)
select * from this_trip_cte limit 1
```

This was very fast, comparatively.

```
with this_trip_cte as
(
  select t.*, ib.*
  from trip t inner join (
    select
      sqrt((max(latitude) - min(latitude))^2 + (max(longitude) - min(longitude))^2) as
length_est,
      trip_id
    from breadcrumb
    group by trip_id
  ) ib
  on t.trip_id = ib.trip_id
  where route_id > 0
)
select * from this_trip_cte limit 1
```

Your Code

Provide a reference to the repository where you store your code. If you are keeping it private then share it with Bruce (bruce.irvin@gmail.com), David and Aman (github references TBD).

https://github.com/coding-gen/dataeng/tree/main/project/4_final