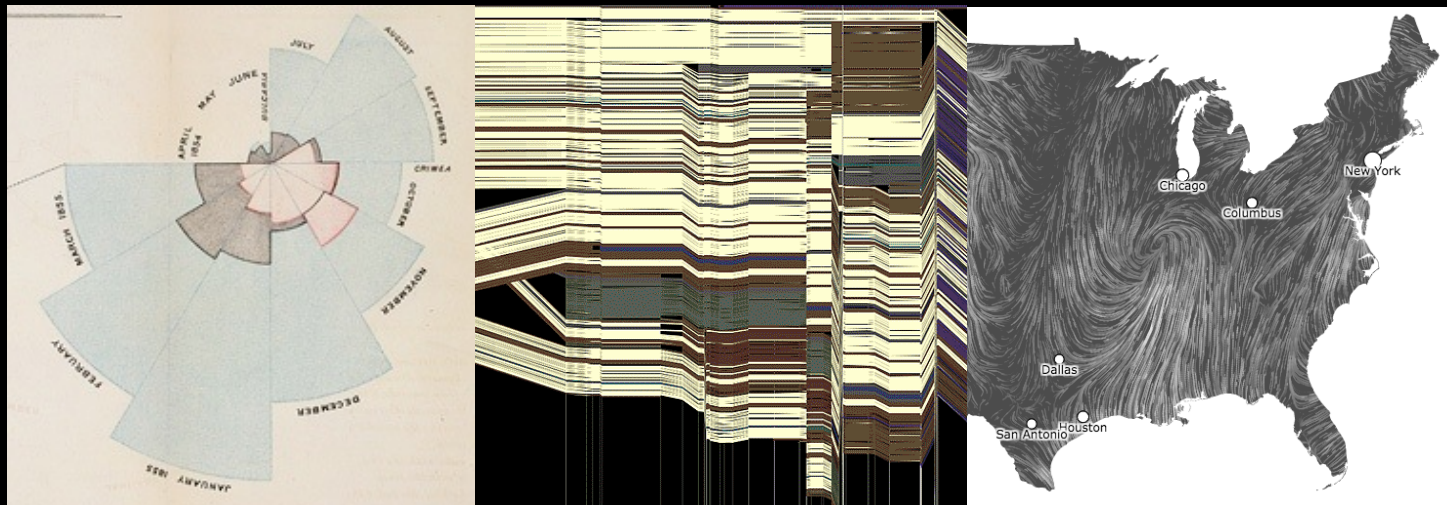


CSE 512 - Data Visualization

Scalability



Leilani Battle University of Washington

How can we visualize and
interact with **billion+ record**
databases in real-time?

Session Outline

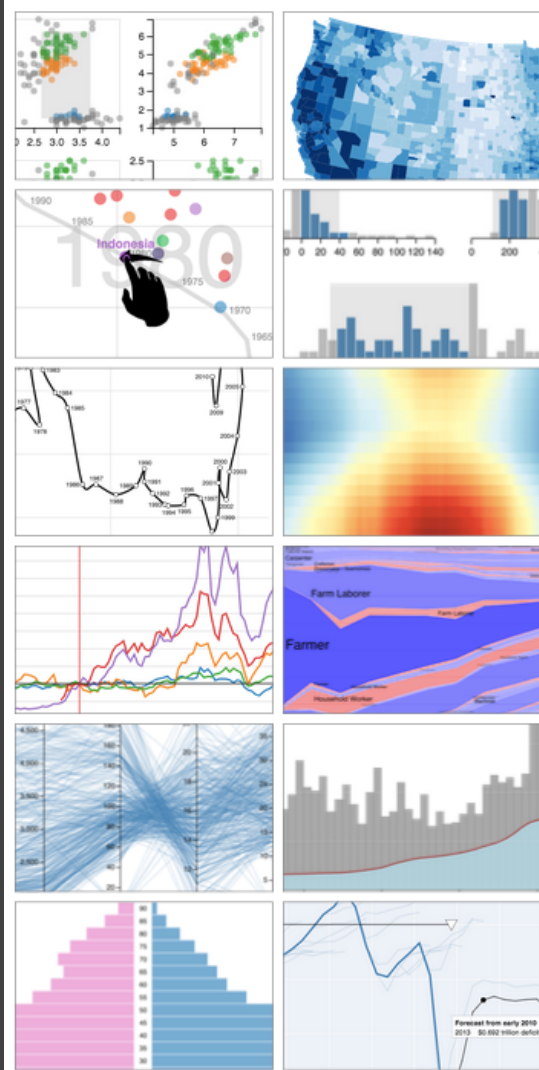
The Varieties of “Big Data”

Scalable Plotting Techniques

Scalable Interaction

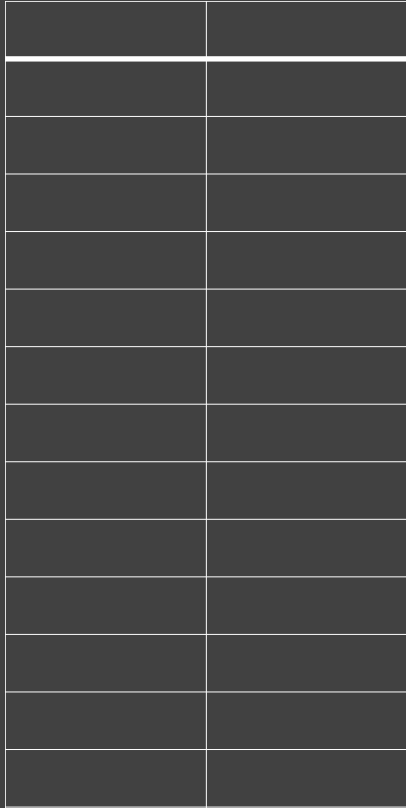
Why Latency Matters

Sampling Methods

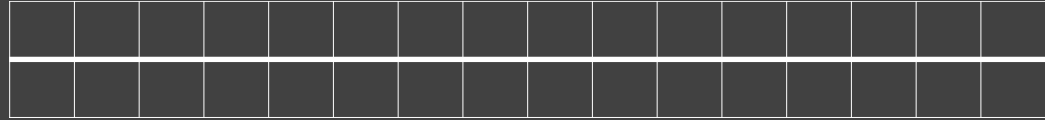


The Varieties of “Big Data”

Tall Data



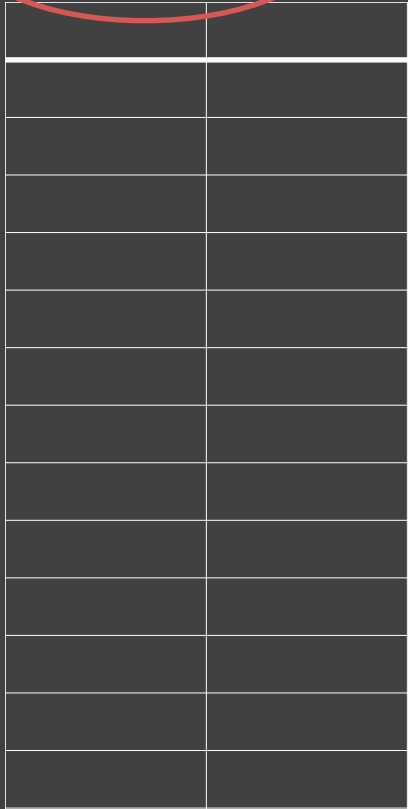
Wide data



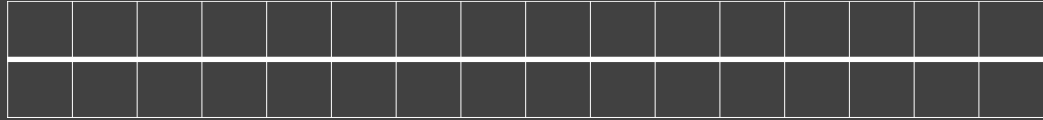
Diverse data



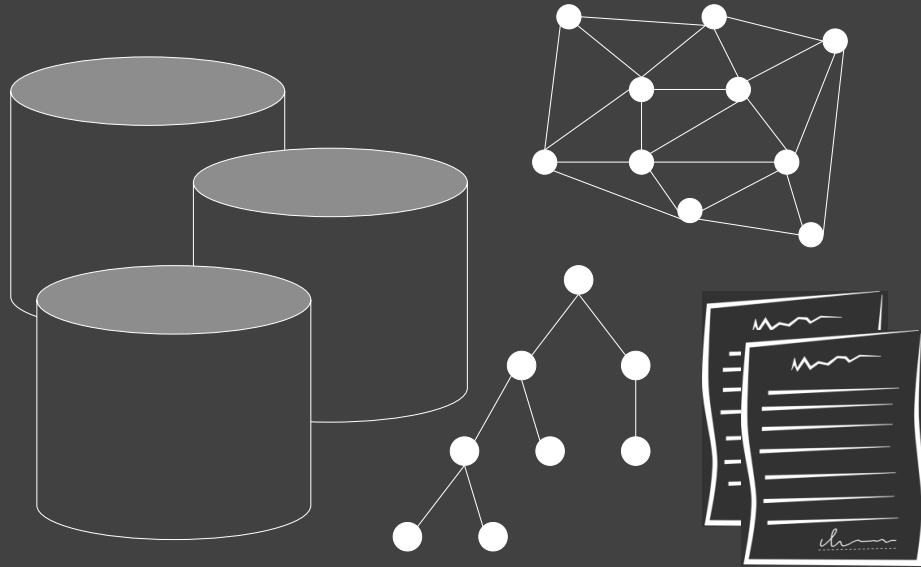
Tall Data



Wide data



Diverse data



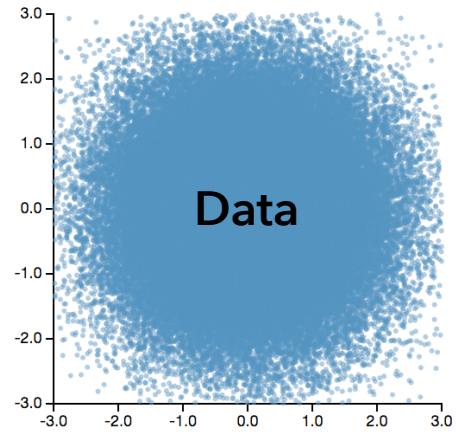
How can we visualize and
interact with **billion+ record**
databases in real-time?

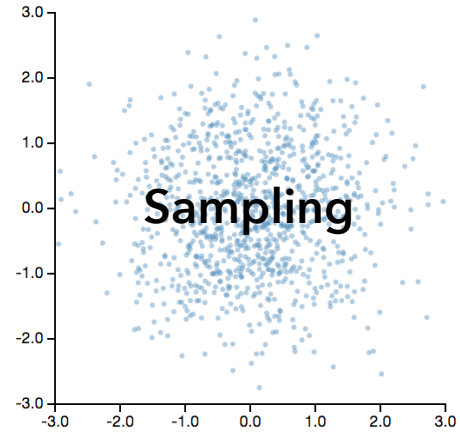
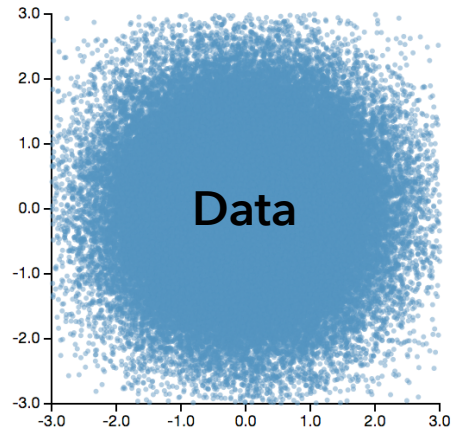
Two Challenges:

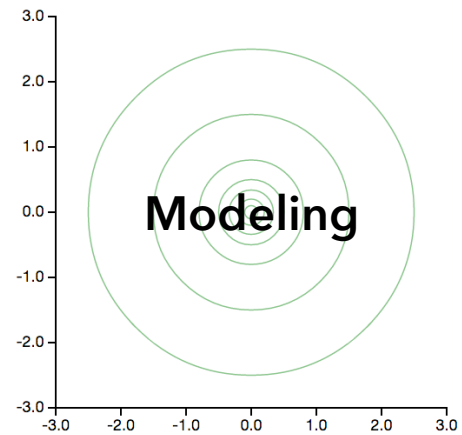
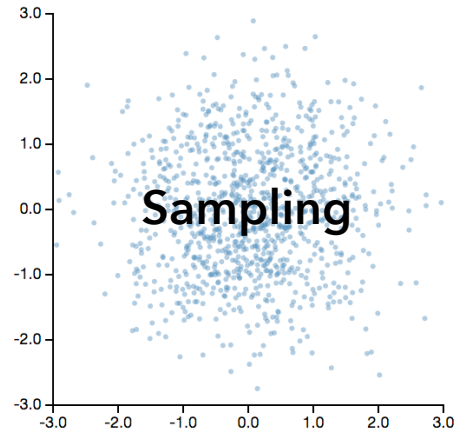
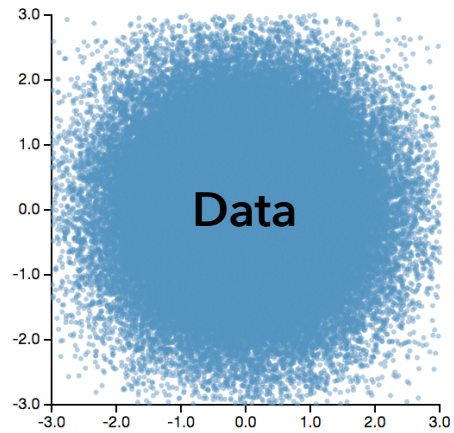
1. Effective **visual encoding**
2. Real-time **interaction**

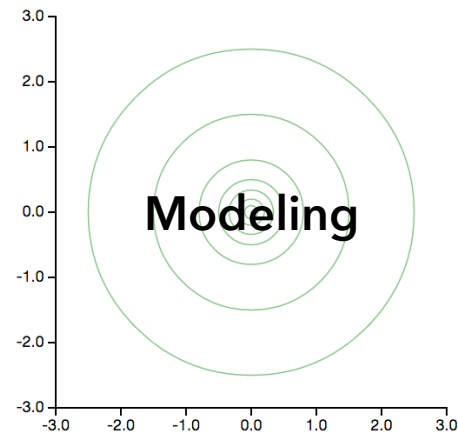
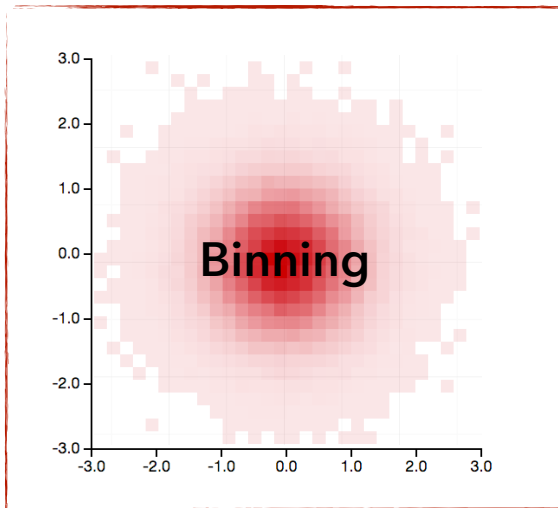
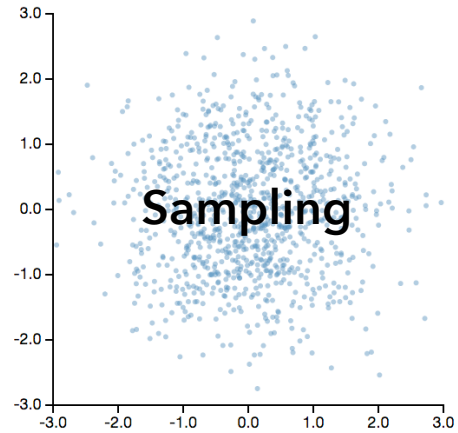
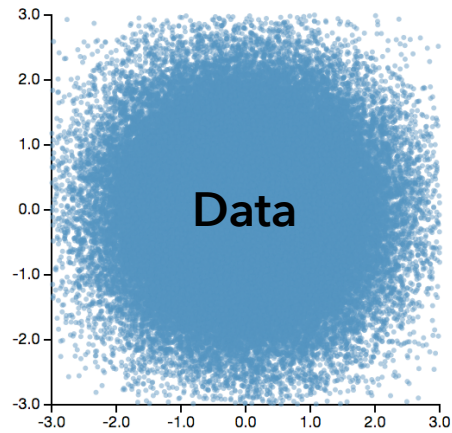
Perceptual and interactive scalability should be limited by the **chosen resolution** of the visualized data, not the number of records.

Scalable Plotting Techniques

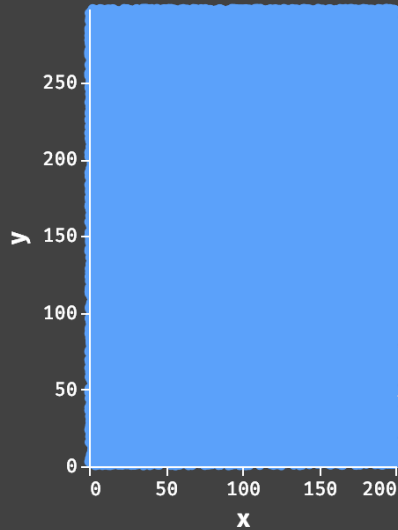




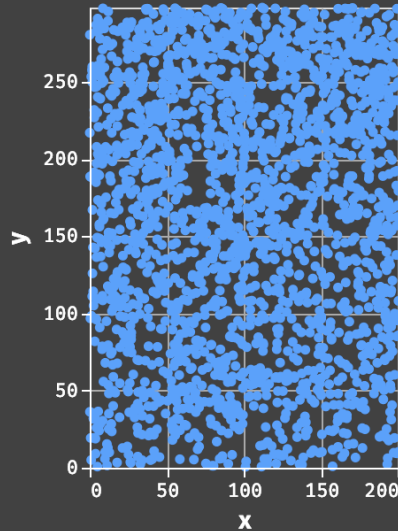




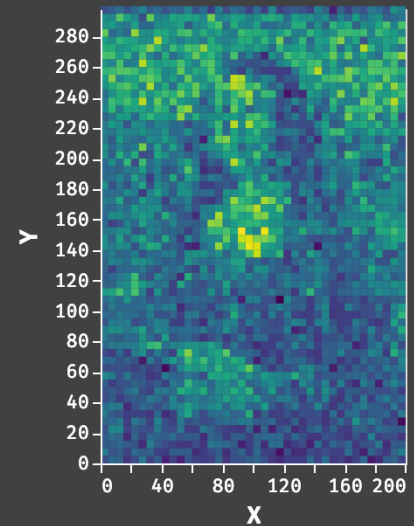
How to **Visualize** a Billion+ Records



Data



Sampling



Binned Aggregation

Decouple the visual complexity from the raw data through aggregation.

Bin > Aggregate (> Smooth) > Plot

1. Bin Divide data domain into discrete “buckets”

Categories: Already discrete (but watch out for high cardinality)

Numbers: Choose bin intervals (uniform, quantile, ...)

Time: Choose time unit: Hour, Day, Month, etc.

Geo: Bin x, y coordinates *after* cartographic projection

Bin > Aggregate (> Smooth) > Plot

1. Bin Divide data domain into discrete “buckets”

Categories: Already discrete (but watch out for high cardinality)

Numbers: Choose bin intervals (uniform, quantile, ...)

Time: Choose time unit: Hour, Day, Month, etc.

Geo: Bin x, y coordinates *after* cartographic projection

2. Aggregate Count, Sum, Average, Min, Max, ...

Bin > Aggregate (> Smooth) > Plot

1. Bin Divide data domain into discrete “buckets”

Categories: Already discrete (but watch out for high cardinality)

Numbers: Choose bin intervals (uniform, quantile, ...)

Time: Choose time unit: Hour, Day, Month, etc.

Geo: Bin x, y coordinates *after* cartographic projection

2. Aggregate Count, Sum, Average, Min, Max, ...

3. Smooth Optional: smooth aggregates [Wickham '13]

Bin > Aggregate (> Smooth) > Plot

1. Bin Divide data domain into discrete “buckets”

Categories: Already discrete (but watch out for high cardinality)

Numbers: Choose bin intervals (uniform, quantile, ...)

Time: Choose time unit: Hour, Day, Month, etc.

Geo: Bin x, y coordinates *after* cartographic projection

2. Aggregate Count, Sum, Average, Min, Max, ...

3. Smooth Optional: smooth aggregates [Wickham '13]

4. Plot Visualize the aggregate values

Binned Plots by Data Type

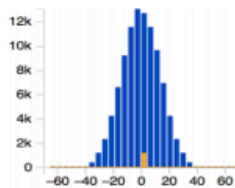
Numeric

Ordinal

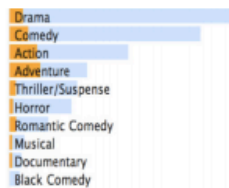
Temporal

Geographic

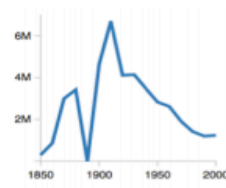
1D



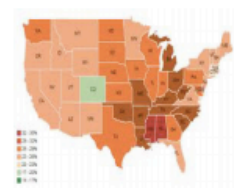
Histogram



Bar Chart

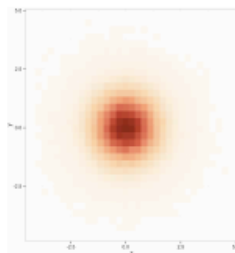


Line Graph /
Area Chart

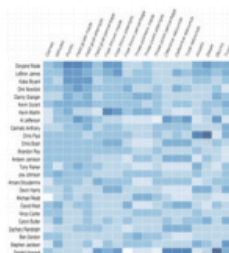


Choropleth Map

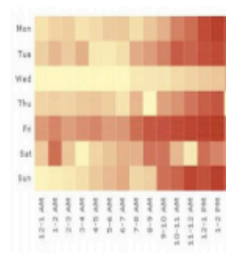
2D



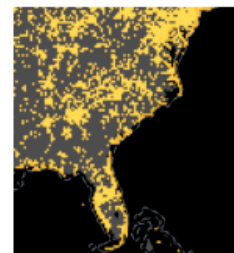
Binned
Scatter Plot



Heatmap



Temporal
Heatmap

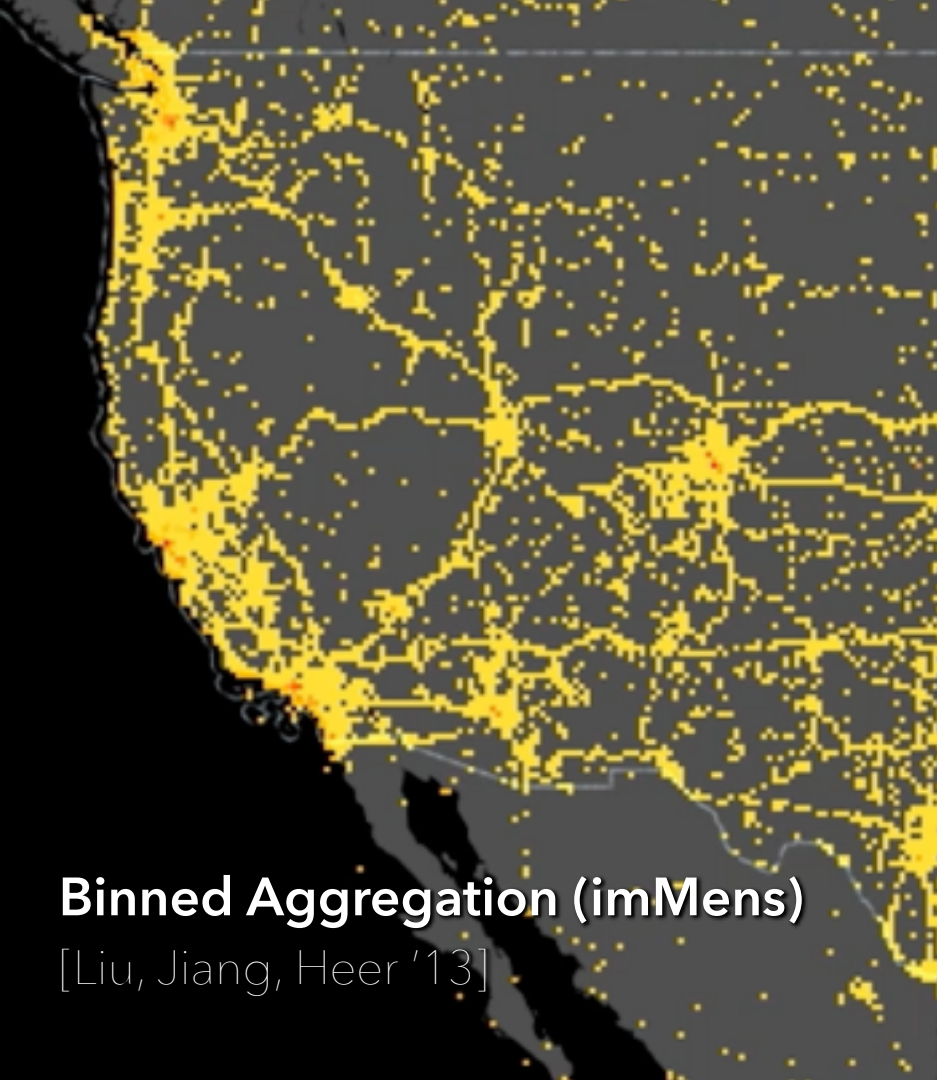


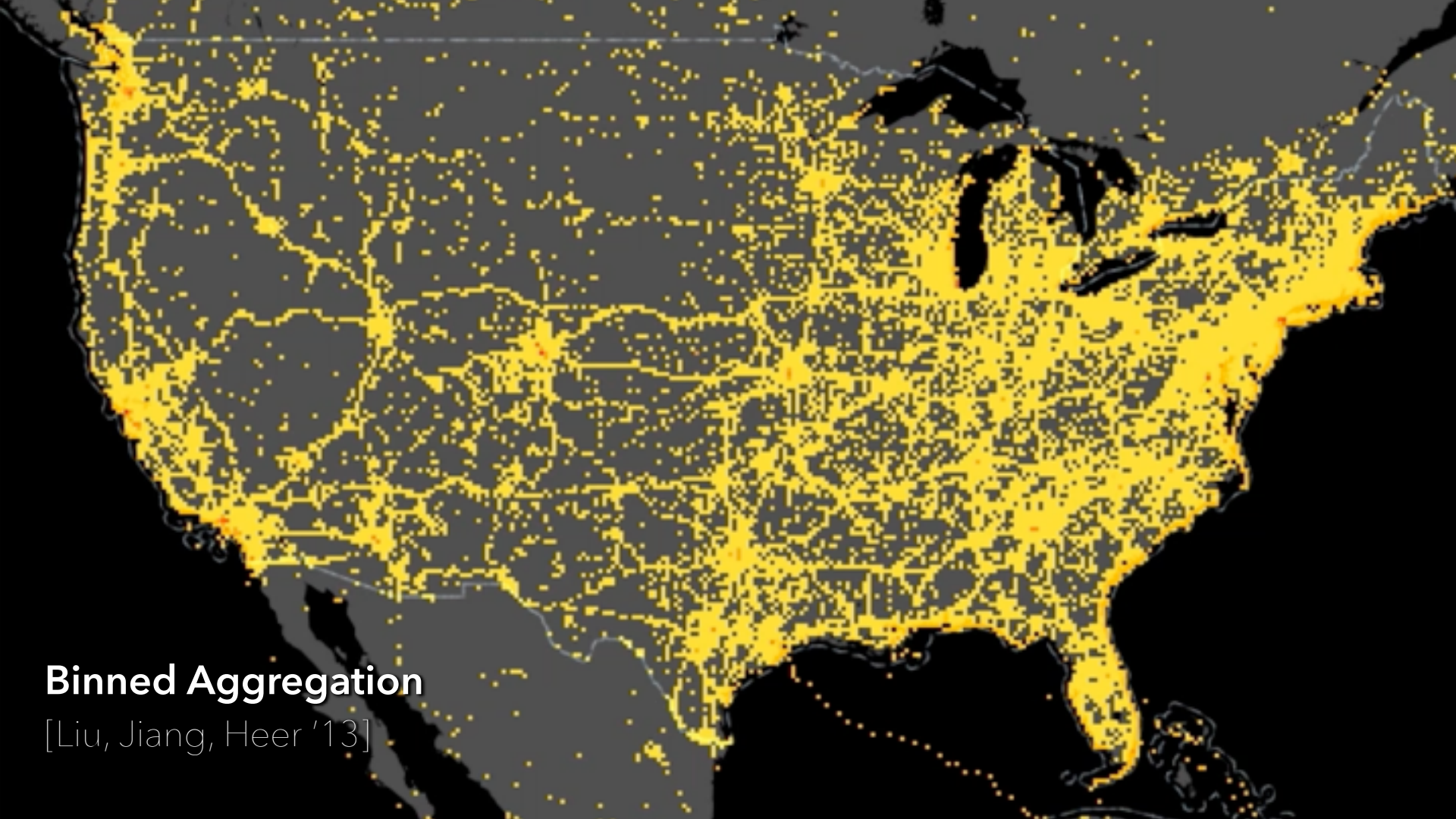
Geographic
Heatmap

Examples



Sampling Google Fusion Tables

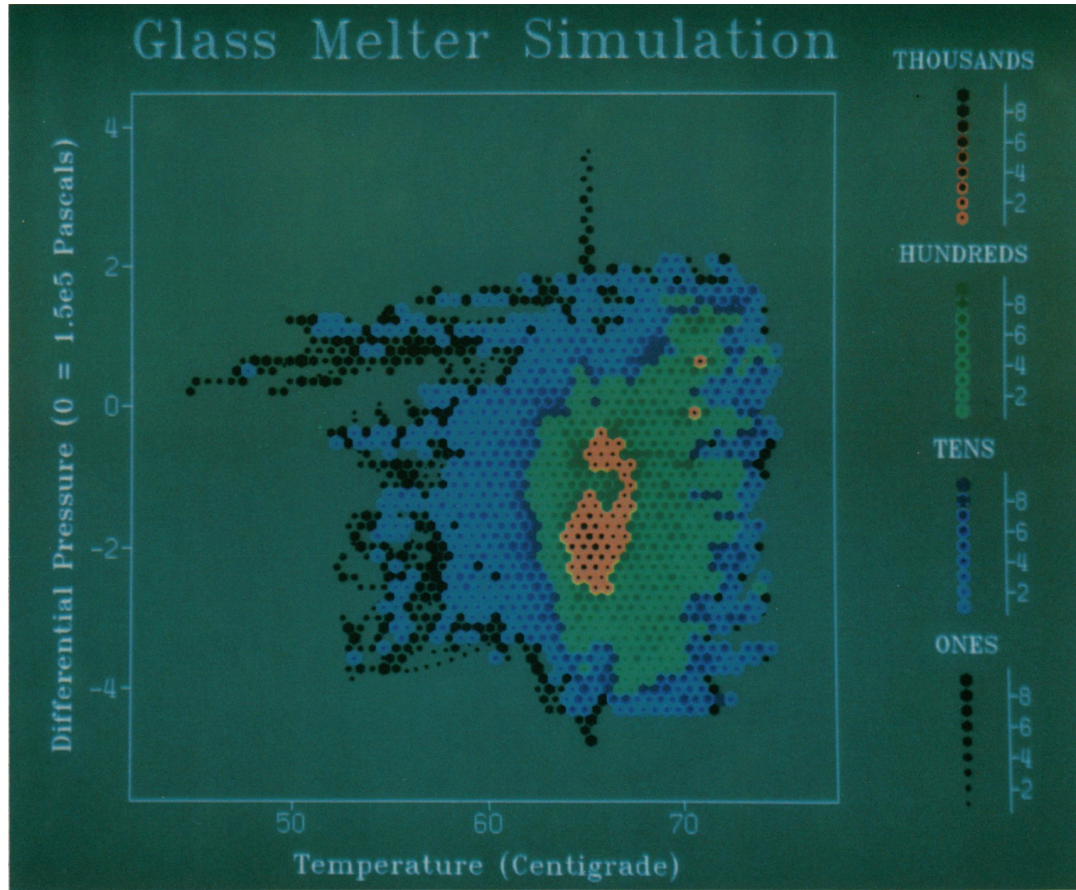




Binned Aggregation

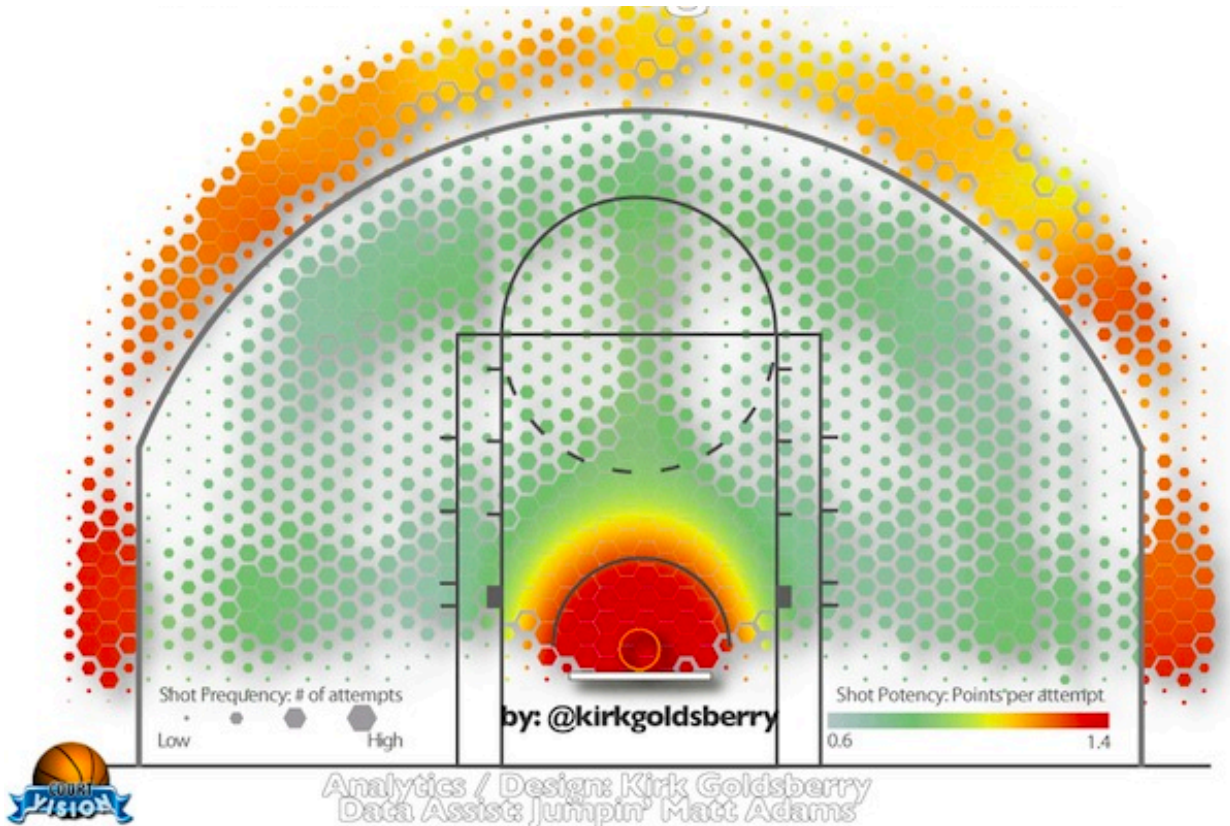
[Liu, Jiang, Heer '13]

Example: Binned Scatter Plots



Scatterplot
Matrix
Techniques
for Large N
[Carr et al. '87]

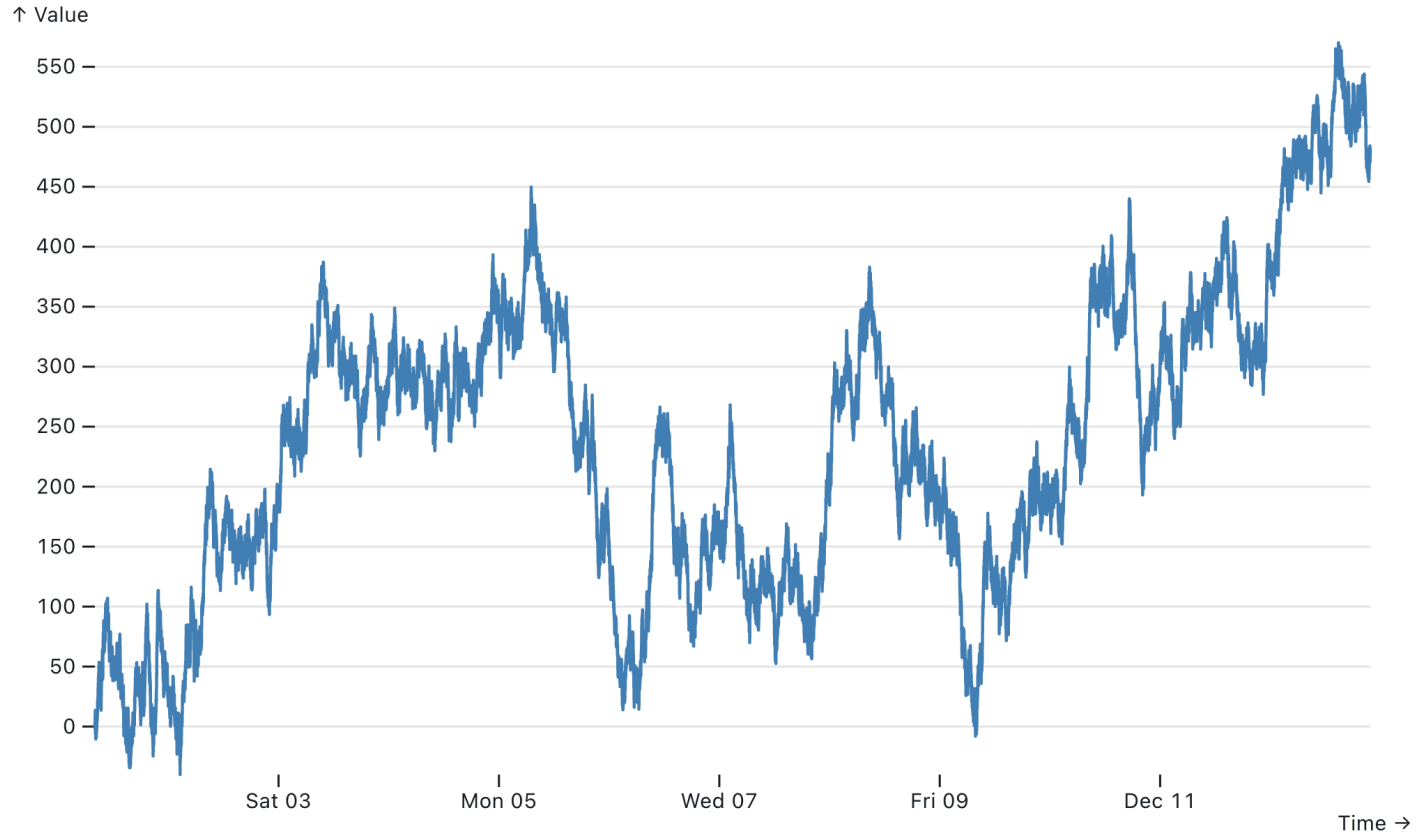
Example: Basketball Shot Chart



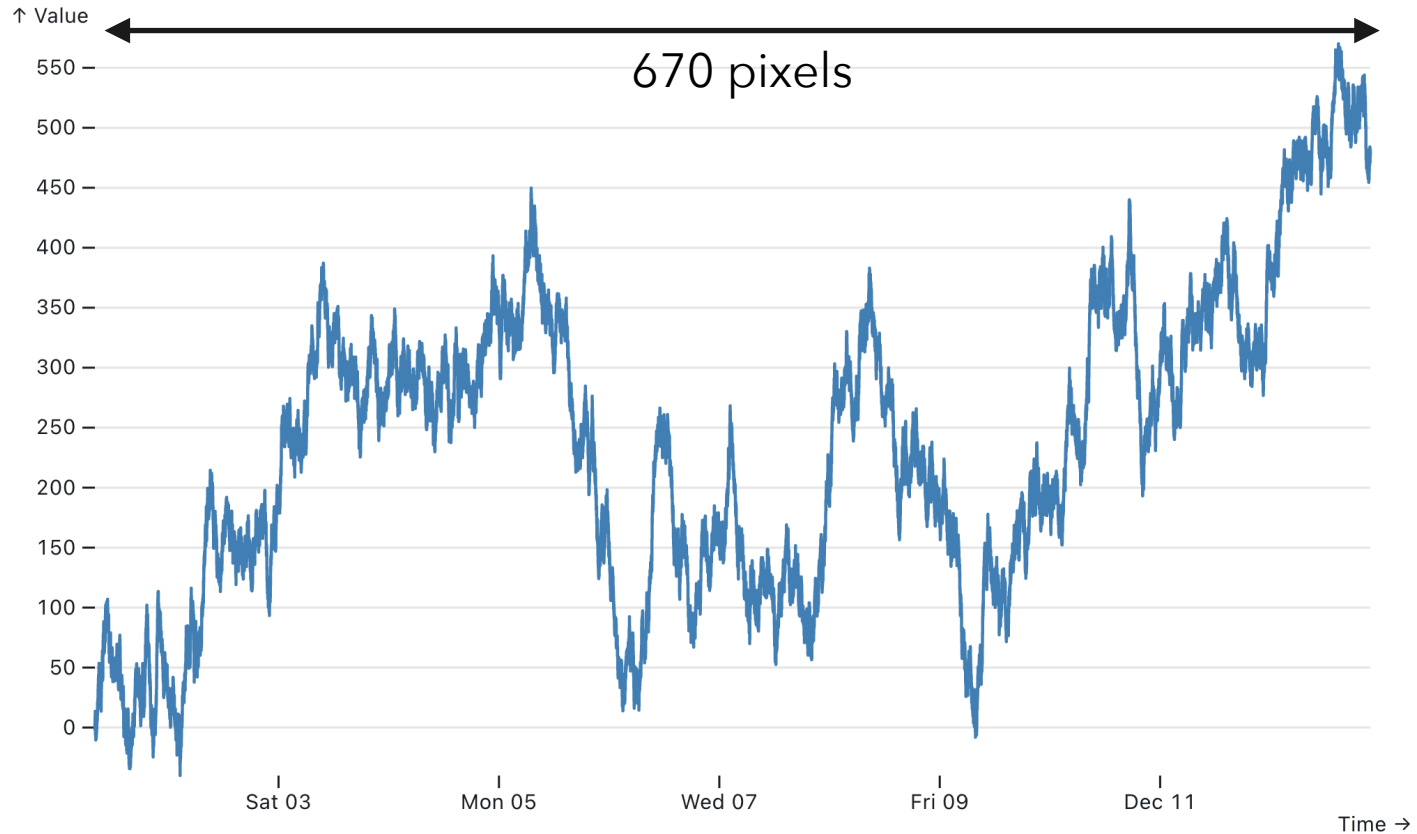
NBA Shooting 2011-12
[Goldsberry]

Time Series

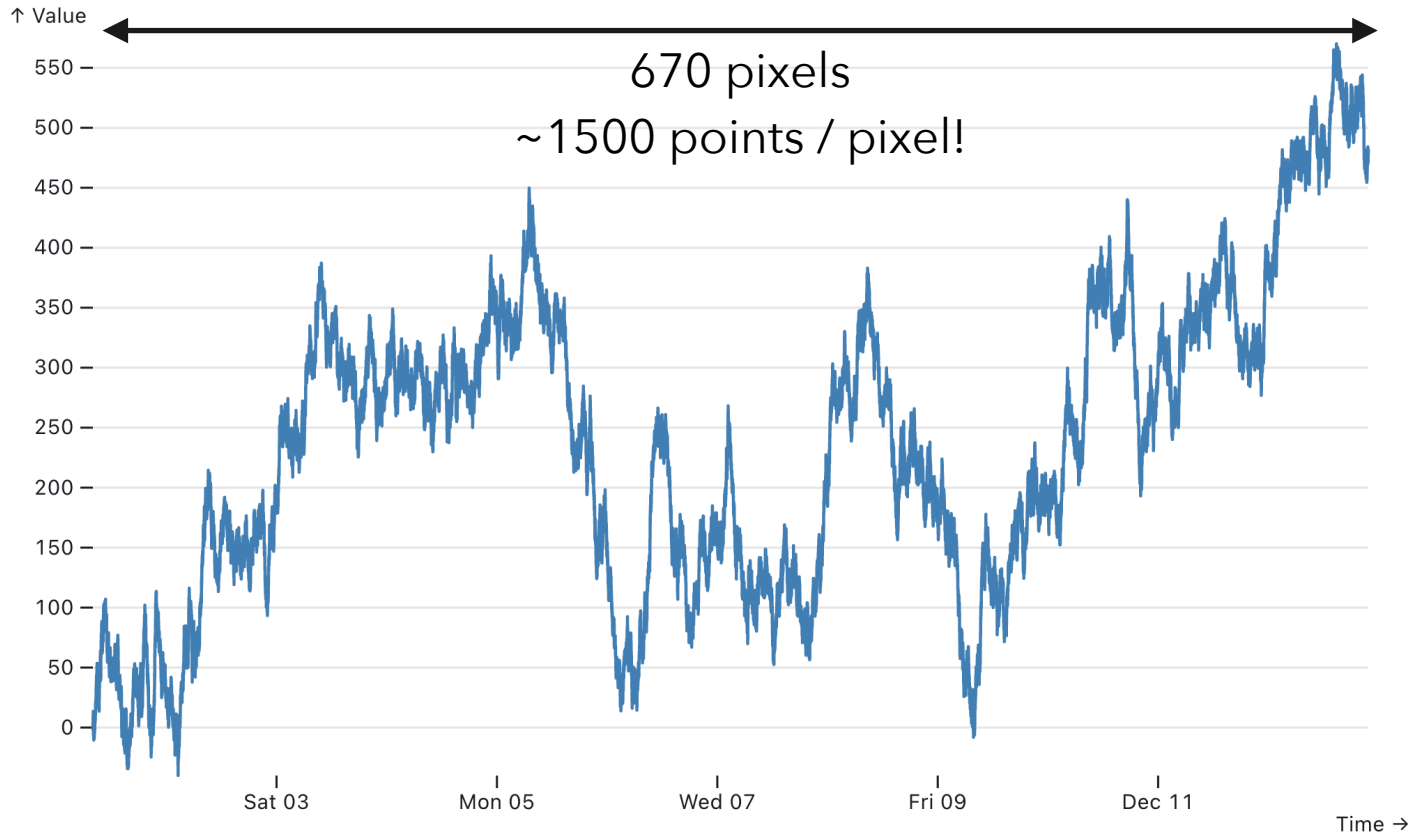
Time Series: 1M samples, 1 sample/second



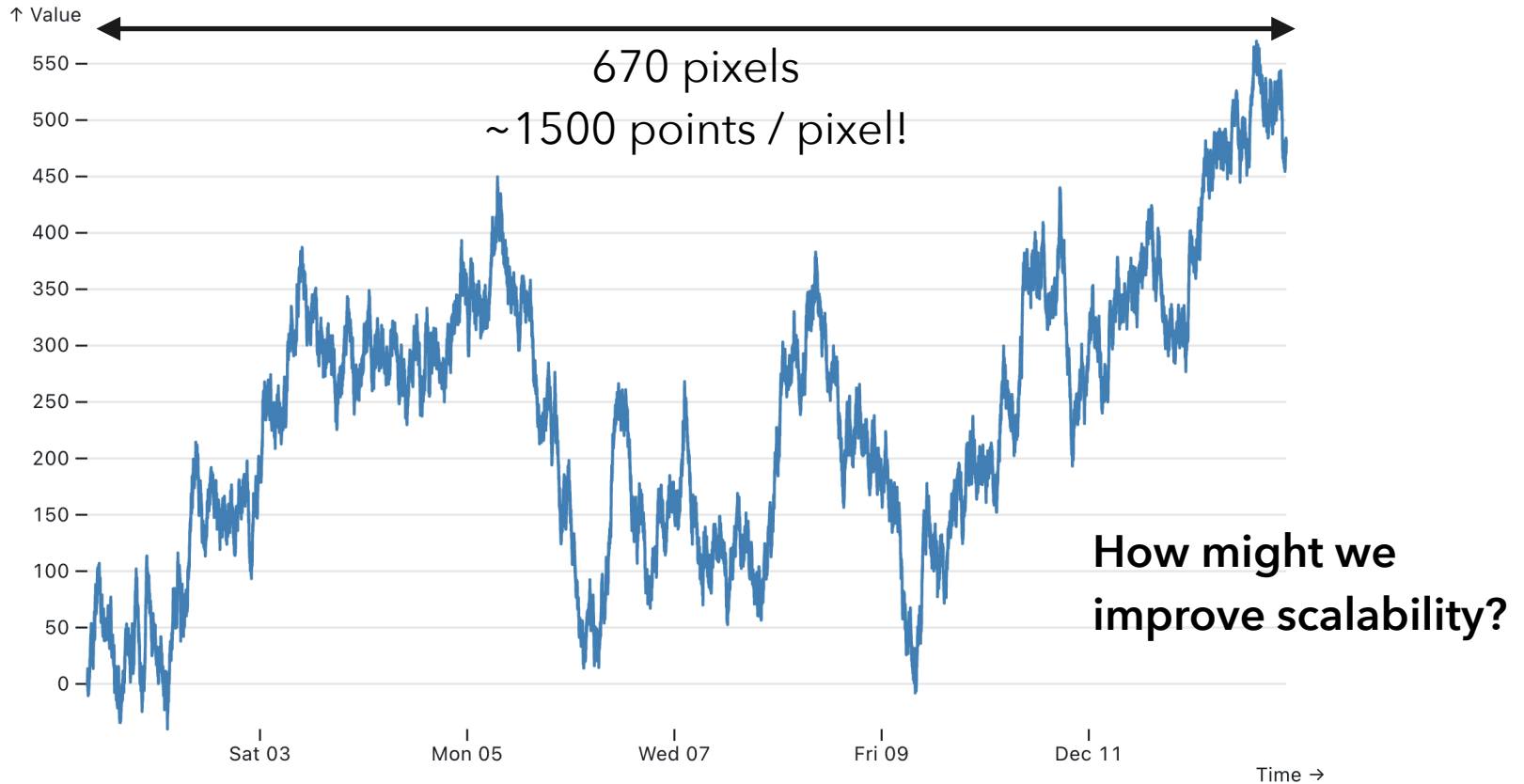
Time Series: 1M samples, 1 sample/second



Time Series: 1M samples, 1 sample/second



Time Series: 1M samples, 1 sample/second



Time-Series Aggregation [Jugel'14]



Insight: the resolution is bound by the number of pixels.

Time-Series Aggregation [Jugel'14]



Insight: the resolution is bound by the number of pixels.

1. Compute average value per pixel (1 point/pixel)
...this may miss extreme (min, max) values



Time-Series Aggregation [Jugel'14]



Insight: the resolution is bound by the number of pixels.

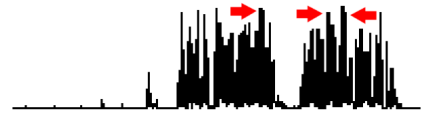
1. Compute average value per pixel (1 point/pixel)

...this may miss extreme (min, max) values



2. Plot min/max values per pixel (2 points/pixel)

...this does better, but still misrepresents



Time-Series Aggregation [Jugel'14]



Insight: the resolution is bound by the number of pixels.

1. Compute average value per pixel (1 point/pixel)

...this may miss extreme (min, max) values



2. Plot min/max values per pixel (2 points/pixel)

...this does better, but still misrepresents



3. [M4](#): min/max values & timestamps (4 points/pixel)

...this provides provable fidelity to the full data!



M4 Data Reduction in the Database

```
SELECT min(t), arg_min(v,t) FROM Q GROUP BY $pixel UNION  
SELECT max(t), arg_max(v,t) FROM Q GROUP BY $pixel UNION  
SELECT arg_min(t,v), min(v) FROM Q GROUP BY $pixel UNION  
SELECT arg_max(t,v), max(v) FROM Q GROUP BY $pixel
```

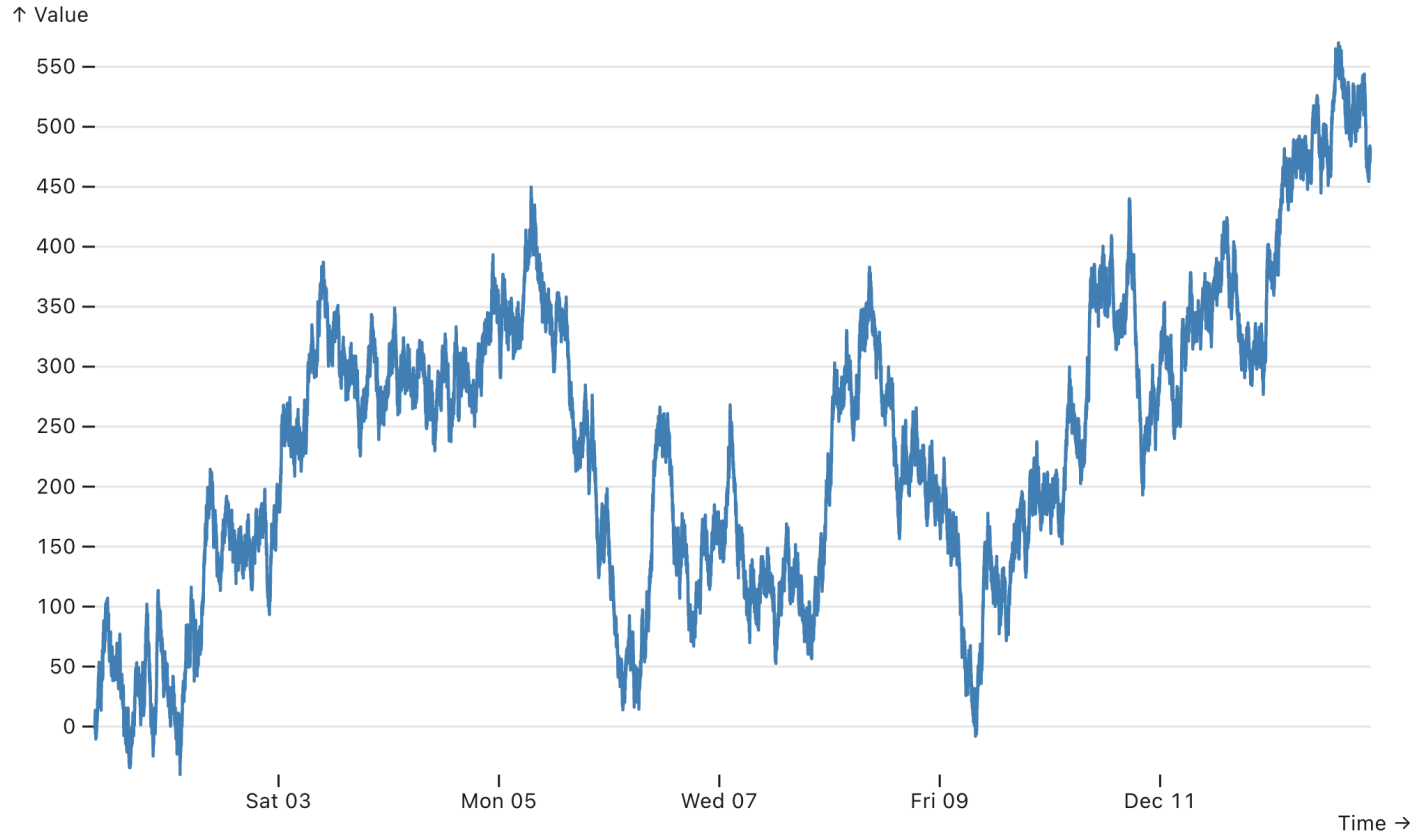
Q: query that returns a time series (t, v)

\$t1, \$t2: global min/max timestamps

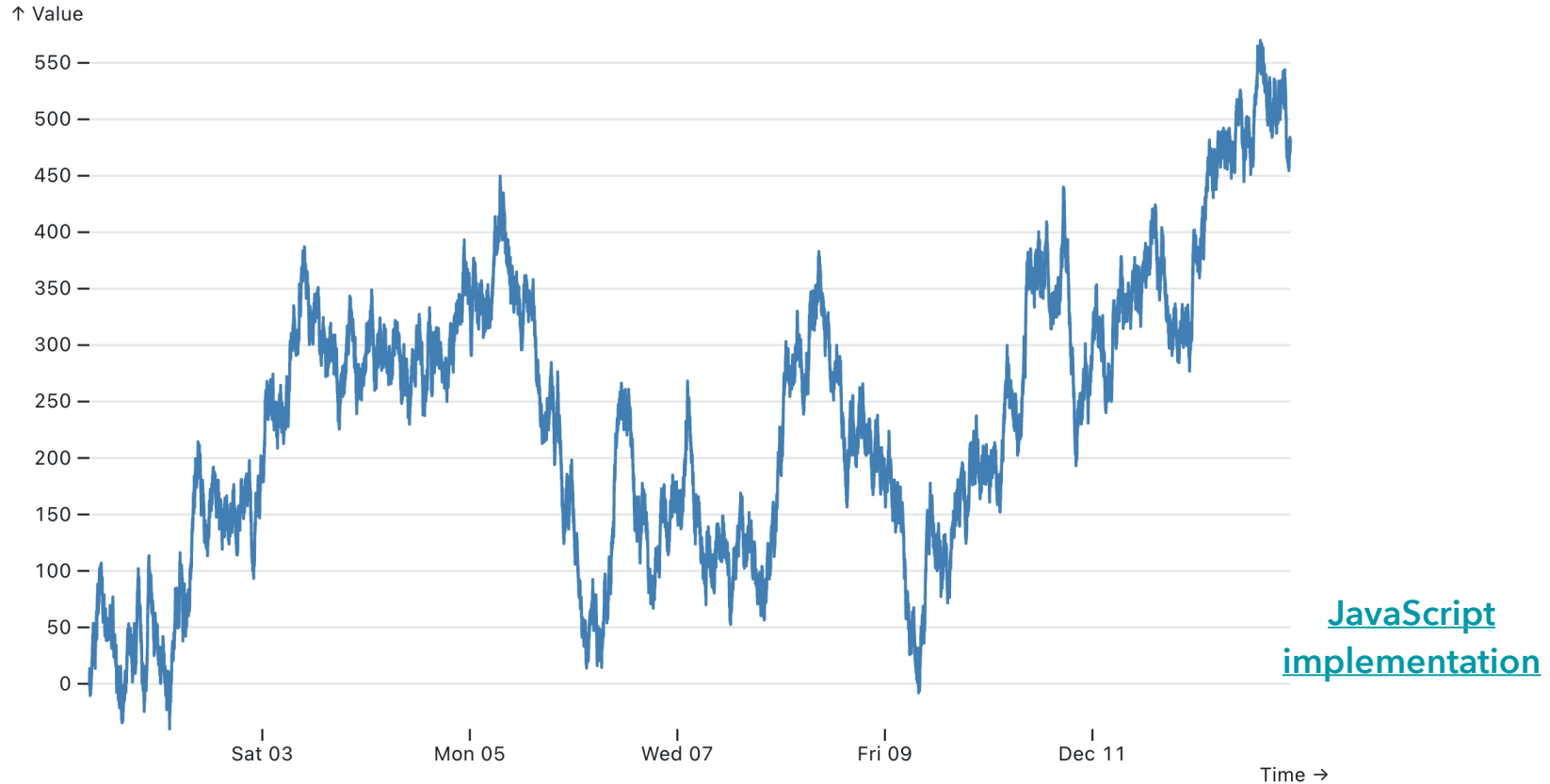
\$w: chart width in pixels

$\$pixel = \text{floor}(\$w(t - \$t1) / (\$t2 - \$t1))$

Time Series: 1M samples, 1 sample/second

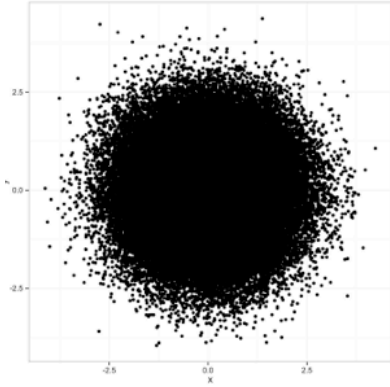


M4: 1M samples -> 2,653 plotted points

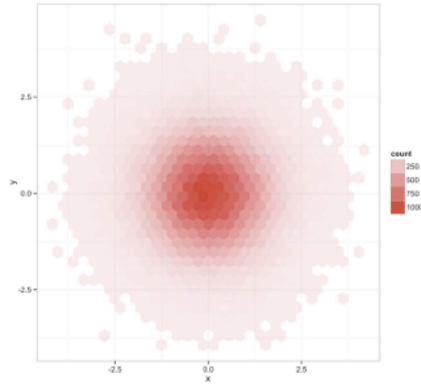


Design Subtleties

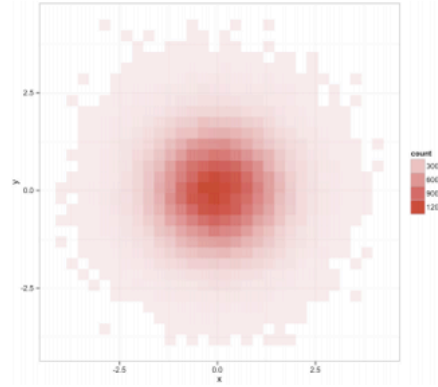
Hexagonal or Rectangular Bins?



100,000 Data Points



Hexagonal Bins

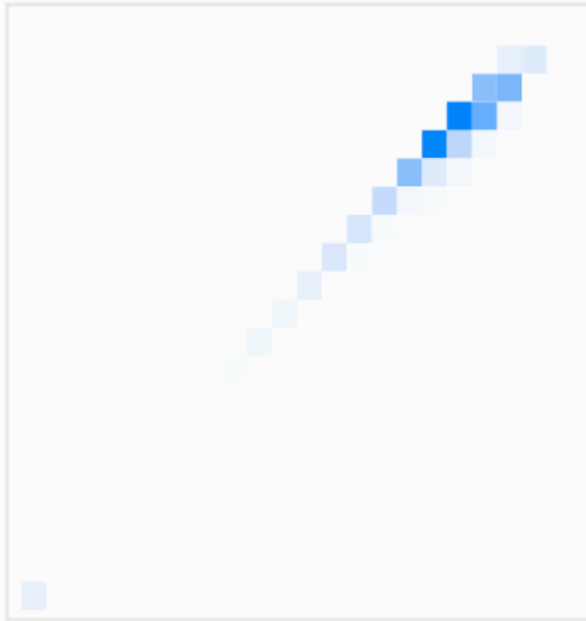


Rectangular Bins

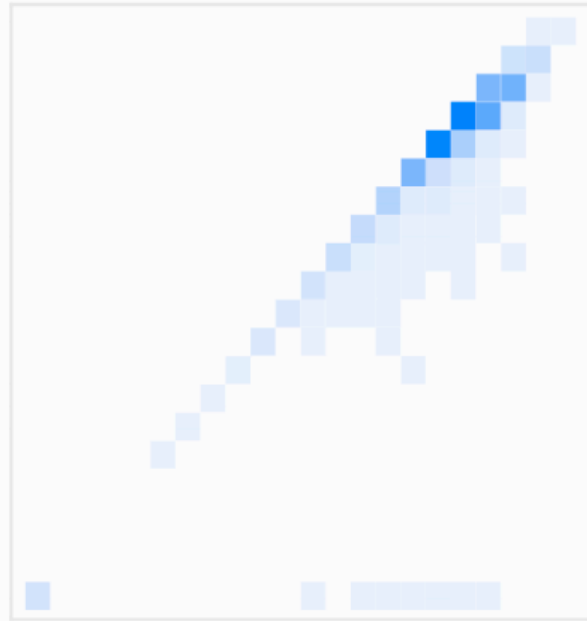
Hex bins better estimate density for 2D plots,
but the *improvement is marginal* [Scott 92].

Rectangles support *reuse* and *visual queries*.

Color Scale: Discontinuity after Zero

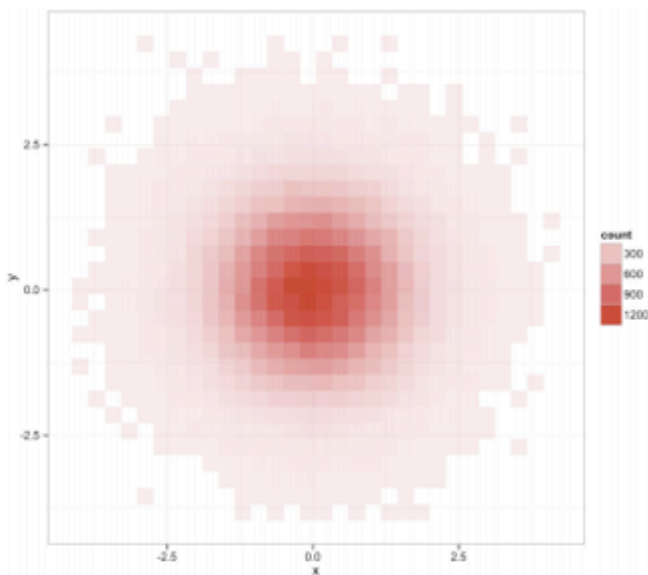


Standard Color Ramp
Counts near zero are white.

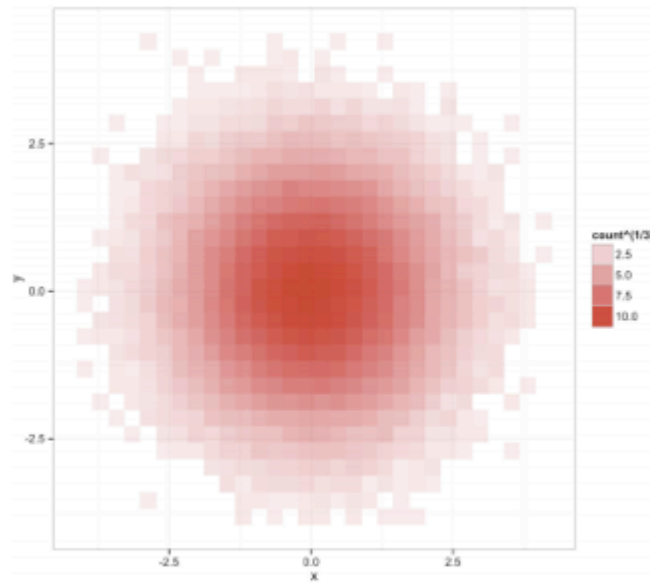


Add Discontinuity after Zero
Counts near zero remain visible.

Color Ramps / Scale Transforms



Linear interpolation in RGBA
is not perceptually linear.



Perceptual color spaces
approximate perceptual linearity.

Questions?

Administrivia

Final Project Schedule

<i>Proposal</i>	Fri May 15
<i>Prototype</i>	Wed May 27
<i>Demo Video</i>	Tue Jun 2
<i>Video Showcase</i>	Thu Jun 4 (in class)
<i>Deliverables</i>	Mon Jun 8

Logistics

Final project description posted online

Work in groups of up to 4 people

You should be well on your way at this point!

Milestone Prototype

Publish work to GitLab pages for others to examine and share feedback. You **are not** expected to have complete, polished content at this point.

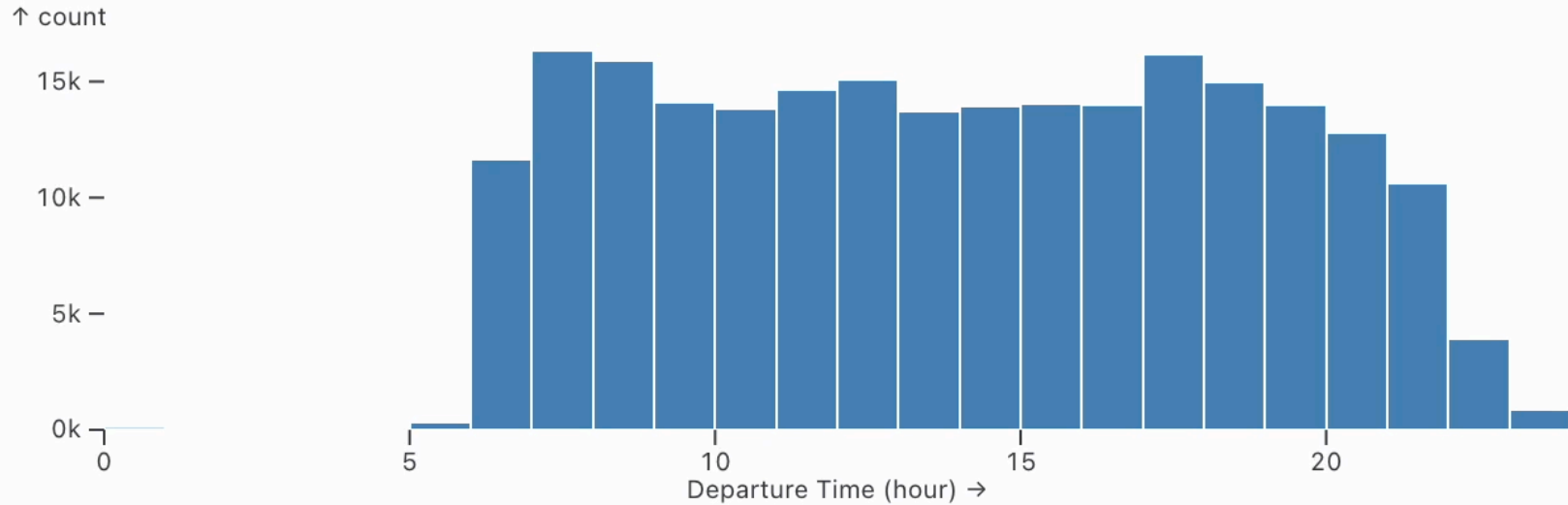
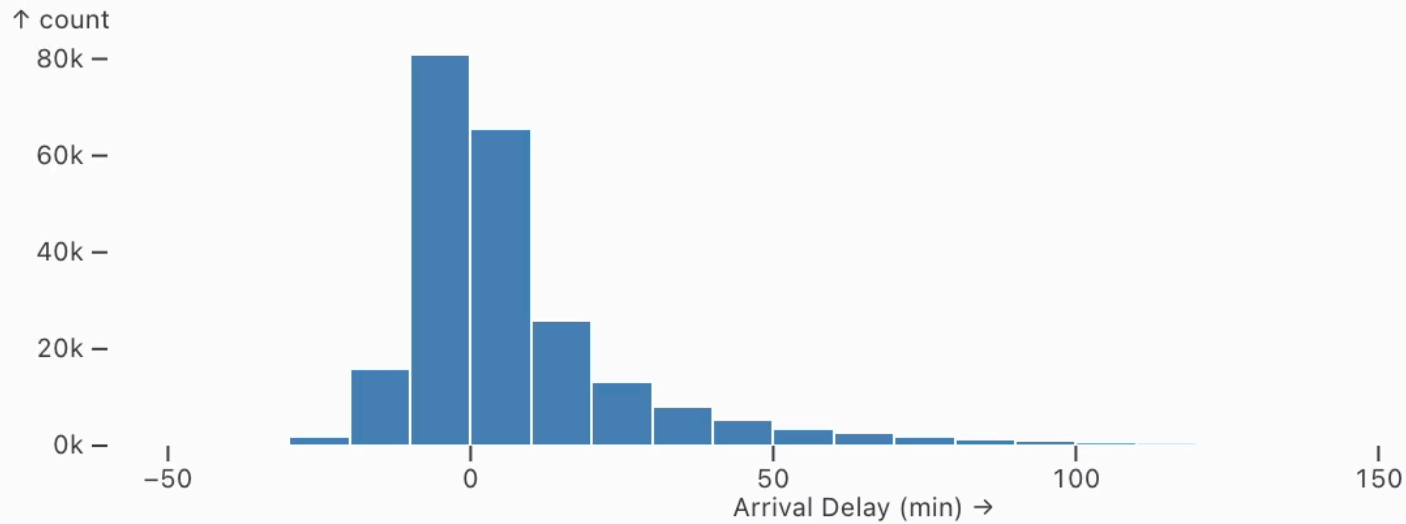
You **are** expected to provide prototype work that communicates your design goals. For example: initial visualizations, sketches, storyboards, and text annotations / idea descriptions.

*One should get a sense of what you intend to submit!
Also feel free to submit questions for us.*

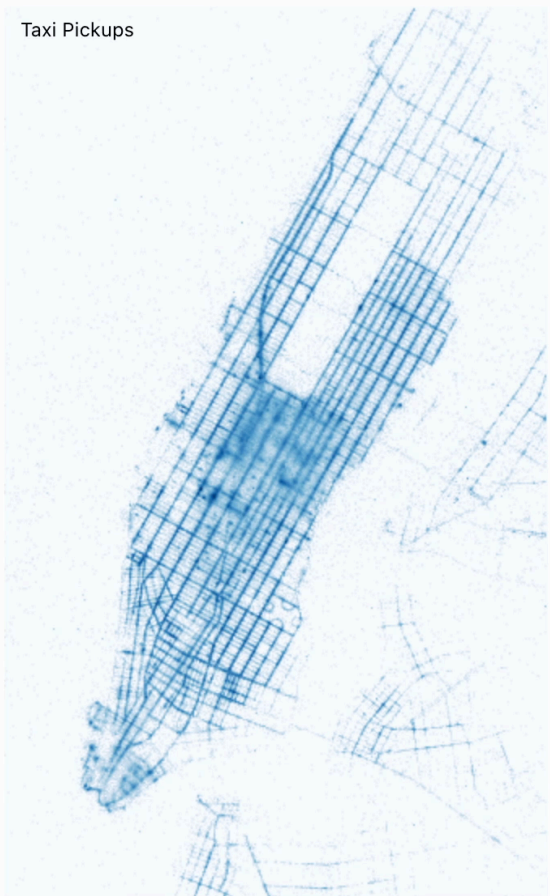
Scalable Interaction

Flight Delays

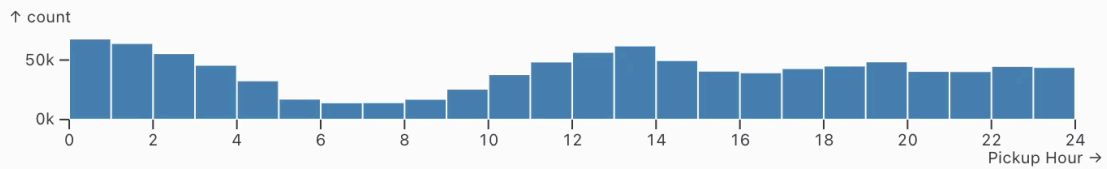
250k Records



Taxi Pickups

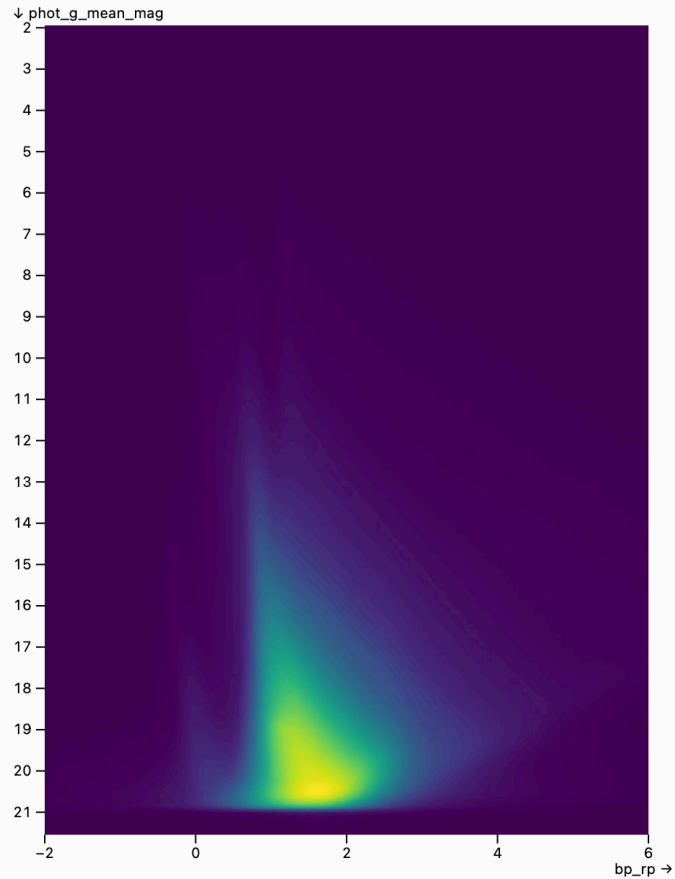
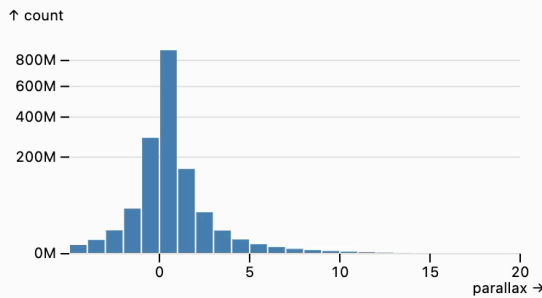
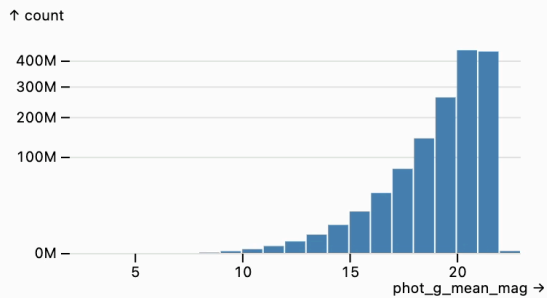
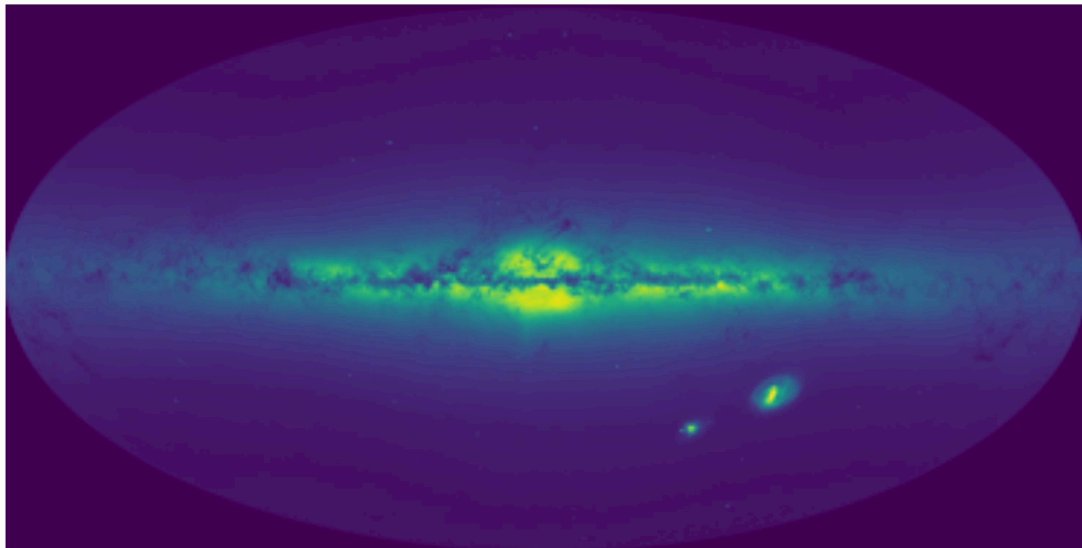


Taxi Dropoffs



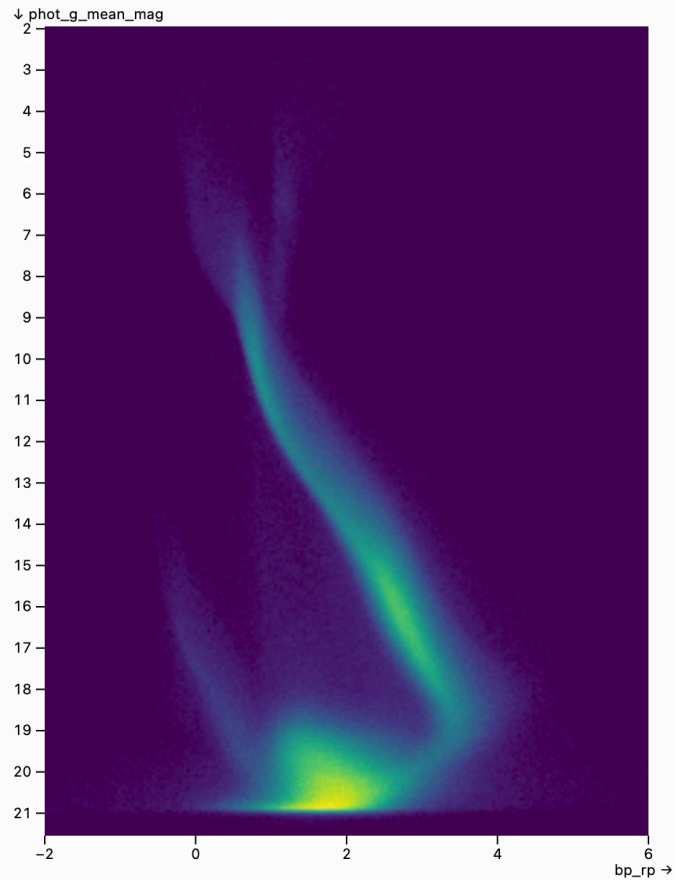
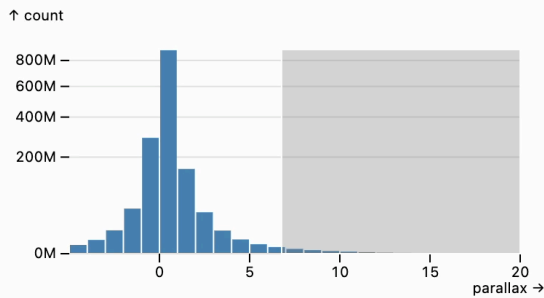
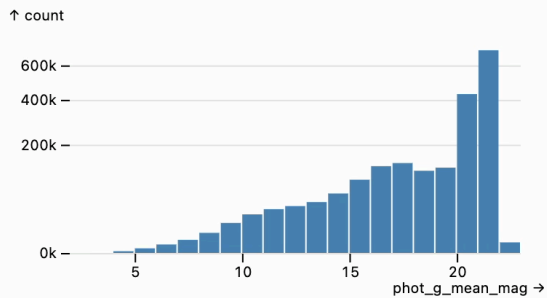
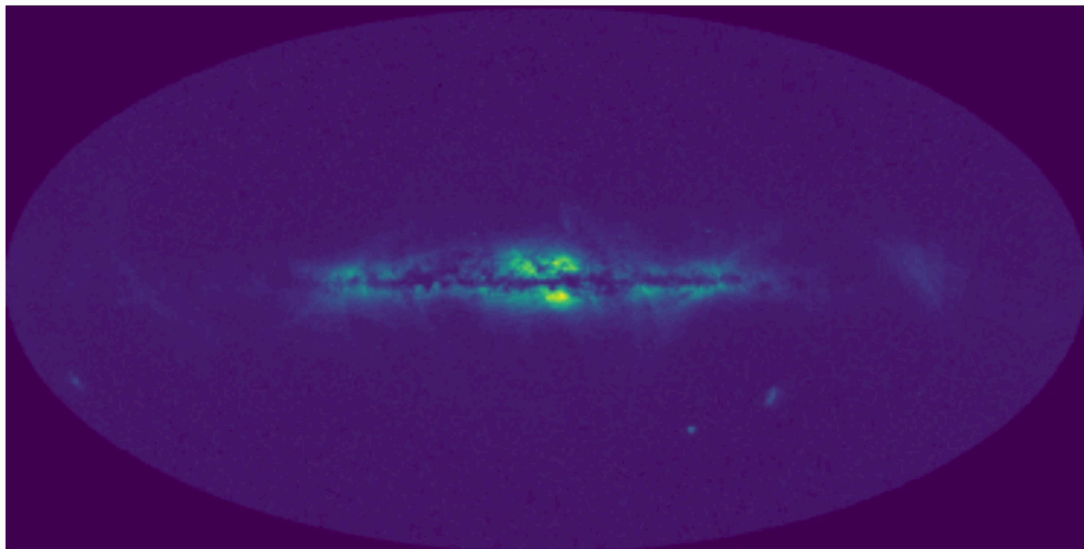
NY Taxi Rides
1M Records
Jan 1-3, 2010

Sample Size Bin Width Color



Gaia Star Catalog · 1.8B Records

Sample Size Bin Width Color



Gaia Star Catalog · 1.8B Records

Interactive Scalability Strategies

1. Query Database
2. Indexing / Preaggregation
3. Prefetching
4. Approximation

Interactive Scalability Strategies

1. Query Database Offload to a scalable backend...

Tableau, for example, issues aggregation queries.

Analytical databases are designed for fast, parallel execution.

But round-trip queries to the DB may still be too slow...

2. Indexing / Preaggregation

3. Prefetching

4. Approximation

Interactive Scalability Strategies

1. Query Database ...or alternative data frame implementation

Python: [Polars](#), [Vaex](#), [Modin](#), [cuDF](#)

R: [dbplyr](#)

All: [DuckDB](#)

2. Indexing / Preaggregation

3. Prefetching

4. Approximation

Interactive Scalability Strategies

1. Query Database

2. Indexing / Preaggregation Query data summaries

Build sorted indices or pre-aggregated data to quickly re-calculate aggregations as needed on the client.

3. Prefetching

4. Approximation

Interactive Scalability Strategies

1. Query Database

2. Indexing / Preaggregation

3. **Prefetching** Request data *before* it is needed

Reduce latency by speculatively querying for data before it is needed. Requires prediction models to guess what is needed.

4. **Approximation**

Interactive Scalability Strategies

1. Query Database

2. Indexing / Preaggregation

3. Prefetching

4. **Approximation** Give fast, approximate answers

Reduce latency by computing aggregates on a sample, ideally with approximation bounds characterizing the error.

Interactive Scalability Strategies

1. Query Database
2. Indexing / Preaggregation
3. Prefetching
4. Approximation

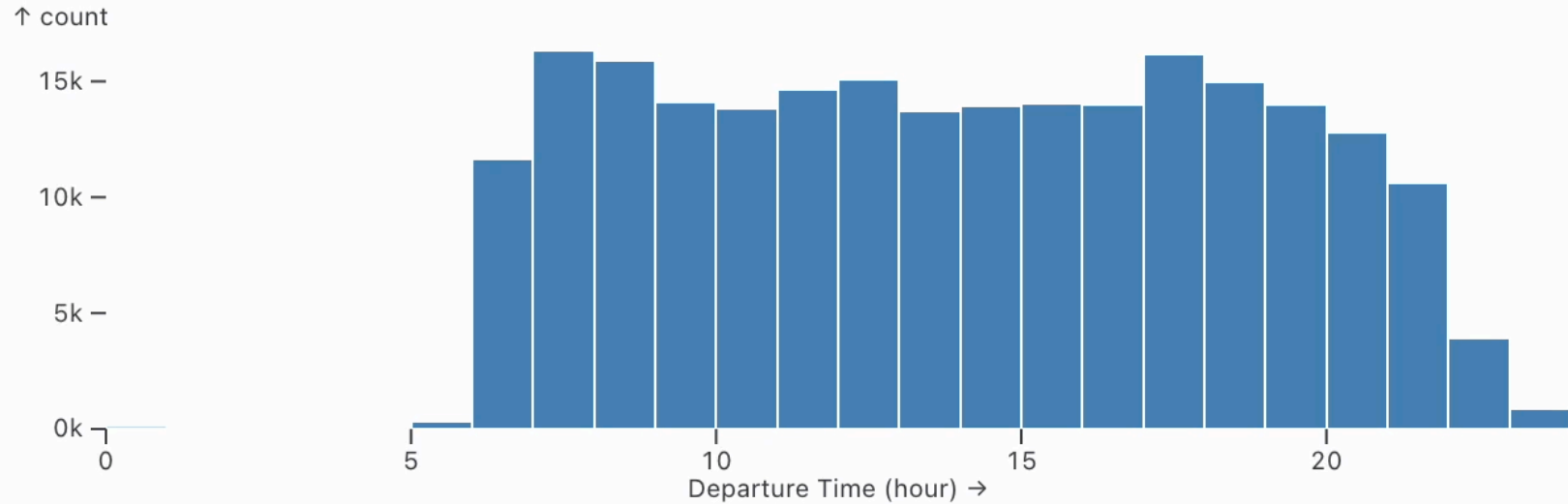
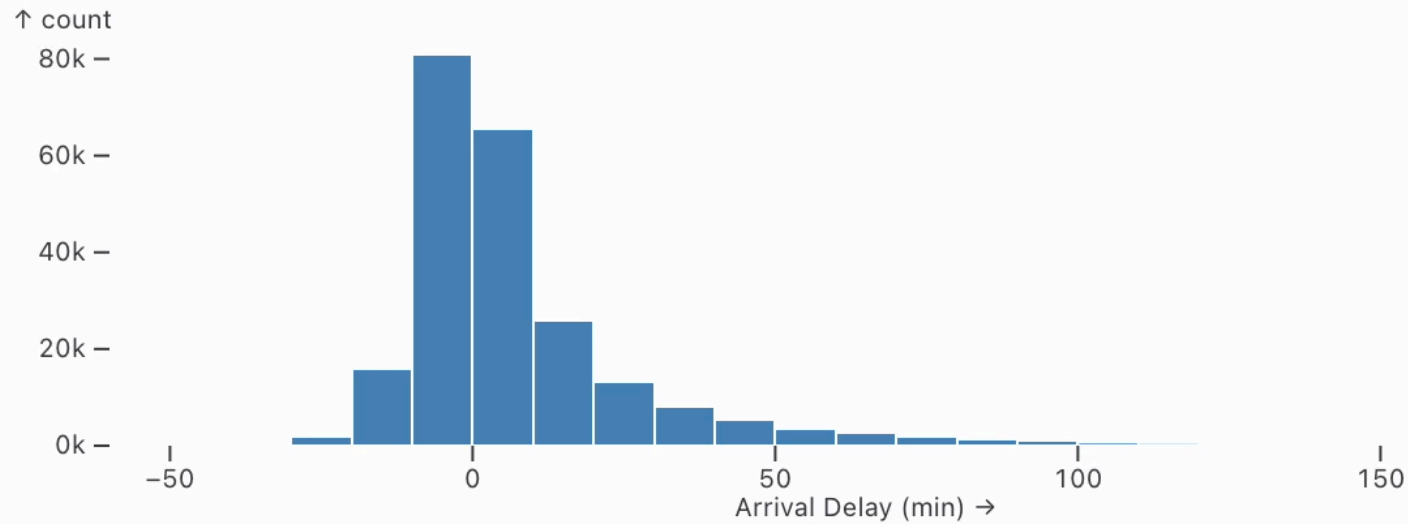
These strategies are **not** mutually exclusive!

Systems can apply them in tandem.

Preaggregation

Flight Delays

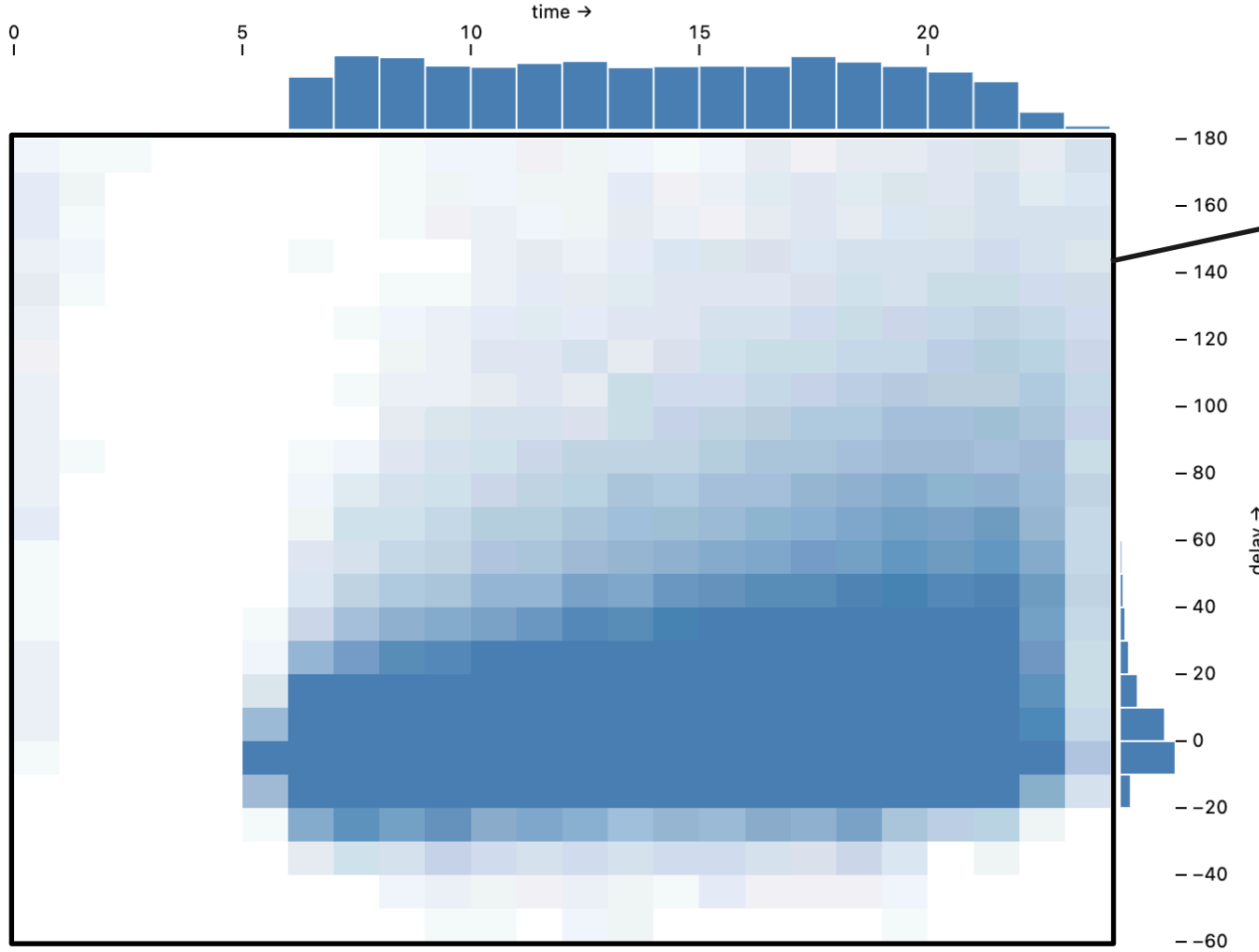
250k Records



Time Resolution Delay Resolution

Flight Delays

250k Records

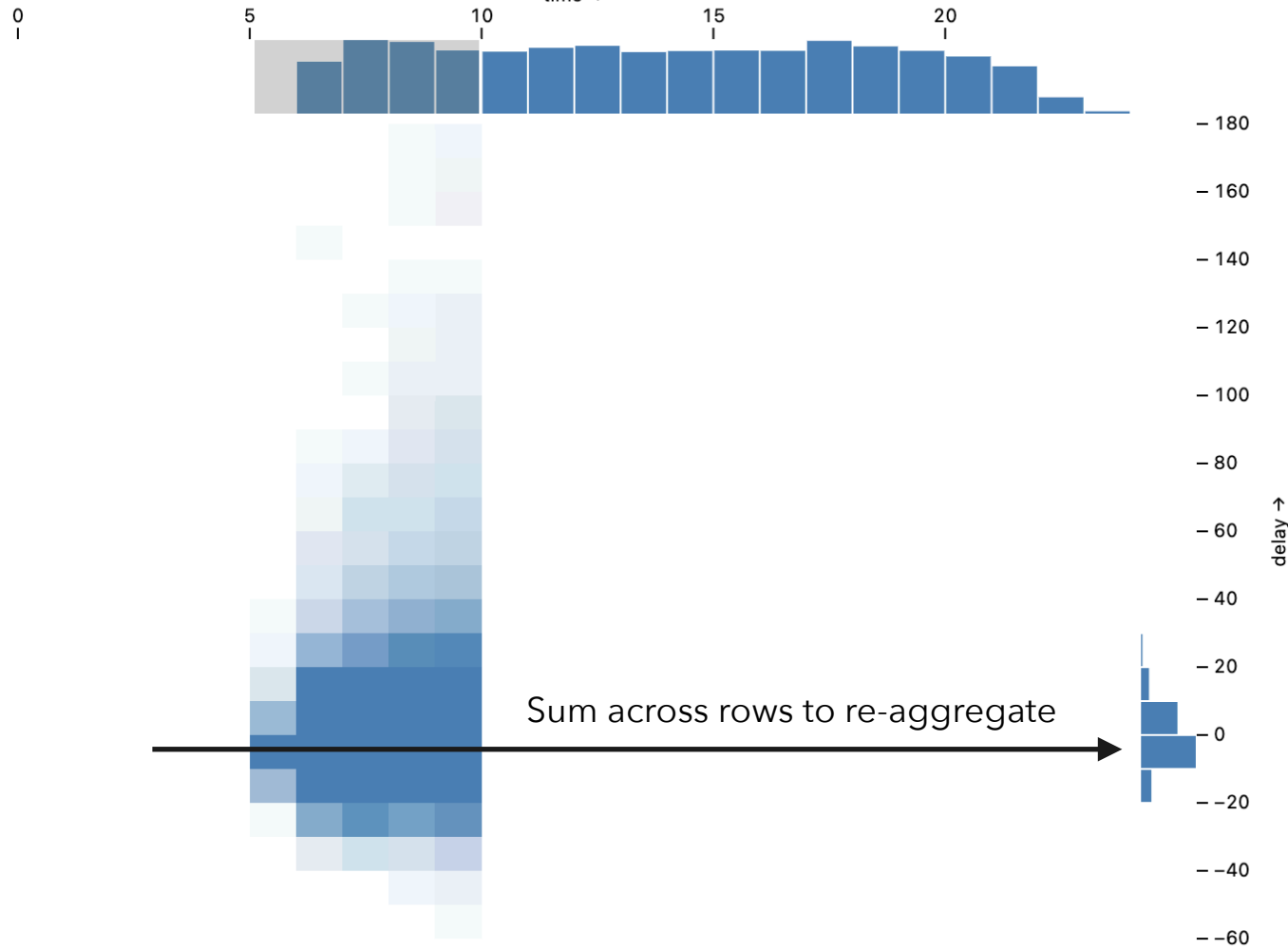


Preaggregate

Time Resolution Delay Resolution

Flight Delays

250k Records



Time Resolution

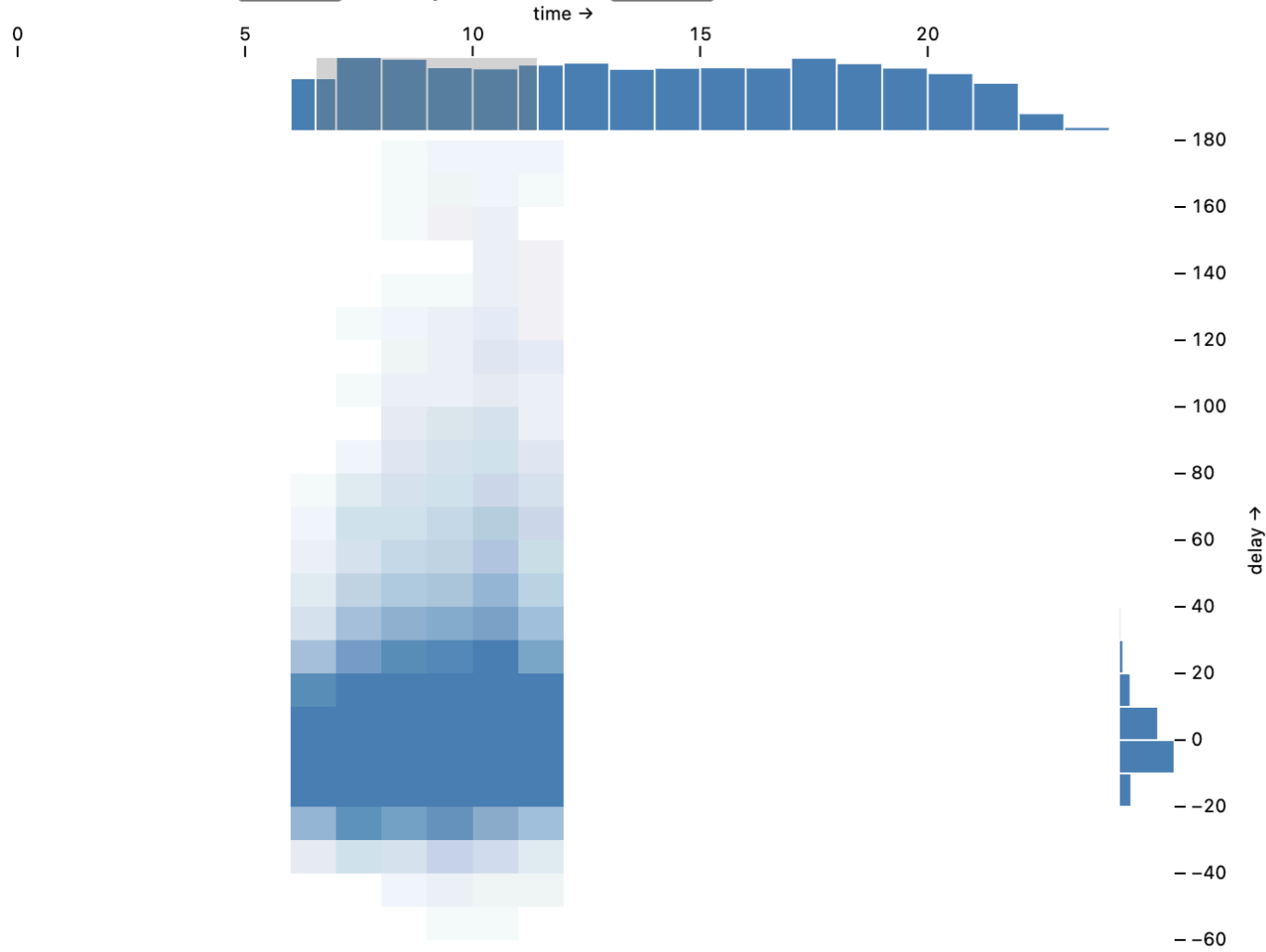
bins

Delay Resolution

bins

Flight Delays

250k Records



Time Resolution

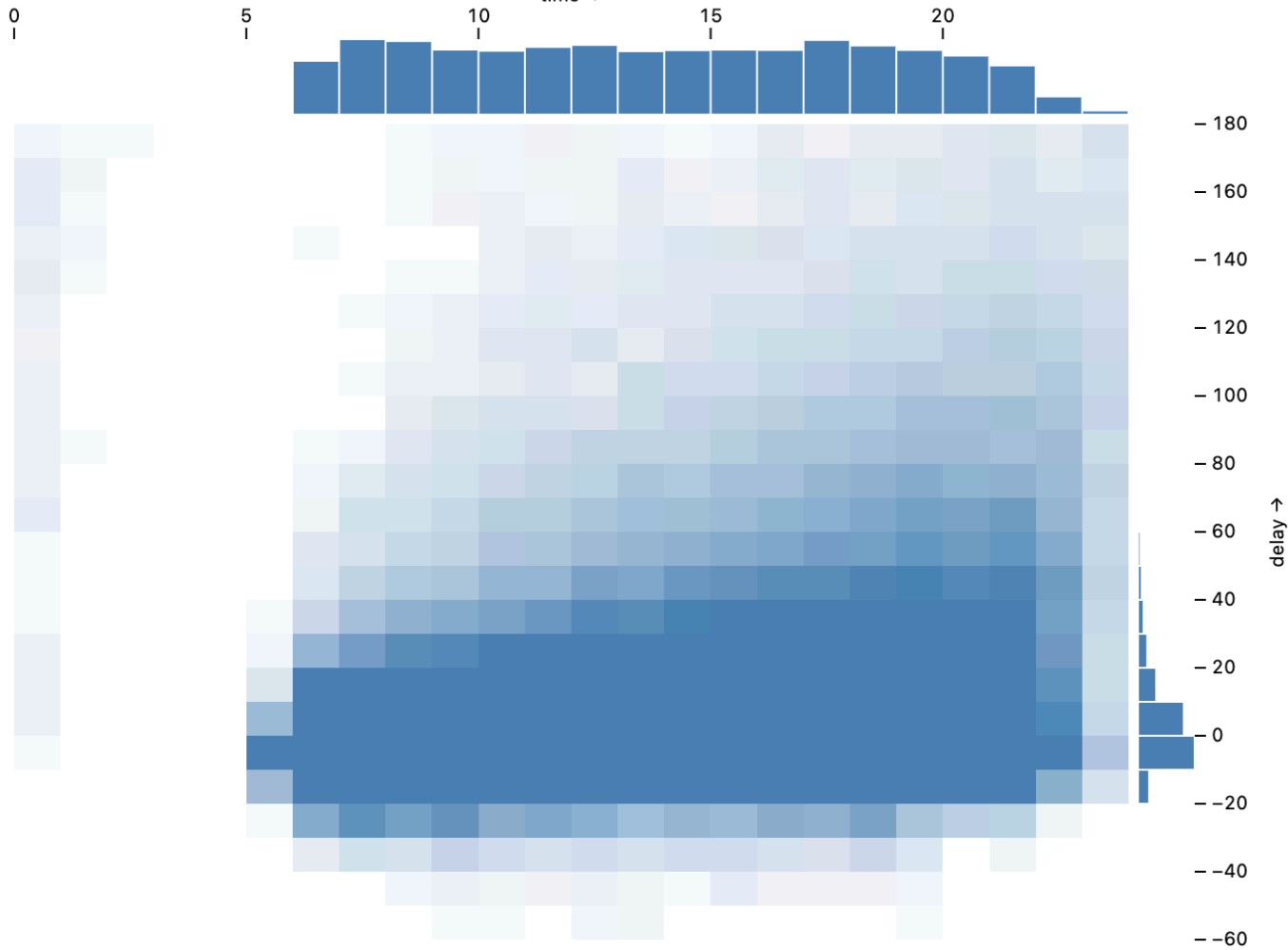
bins ▾

Delay Resolution

bins ▾

Flight Delays

250k Records



Time Resolution

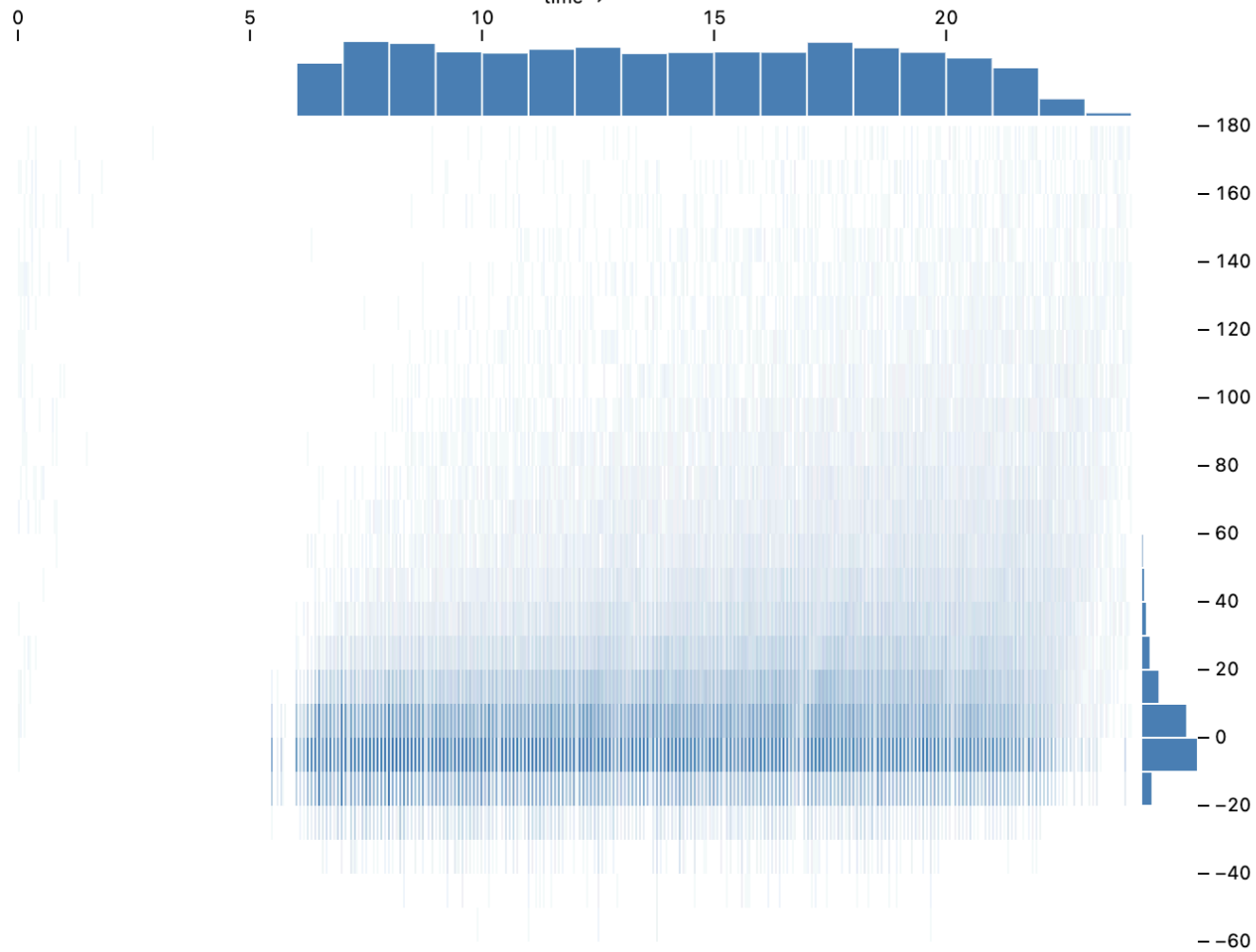
pixels ▾

Delay Resolution

bins ▾

Flight Delays

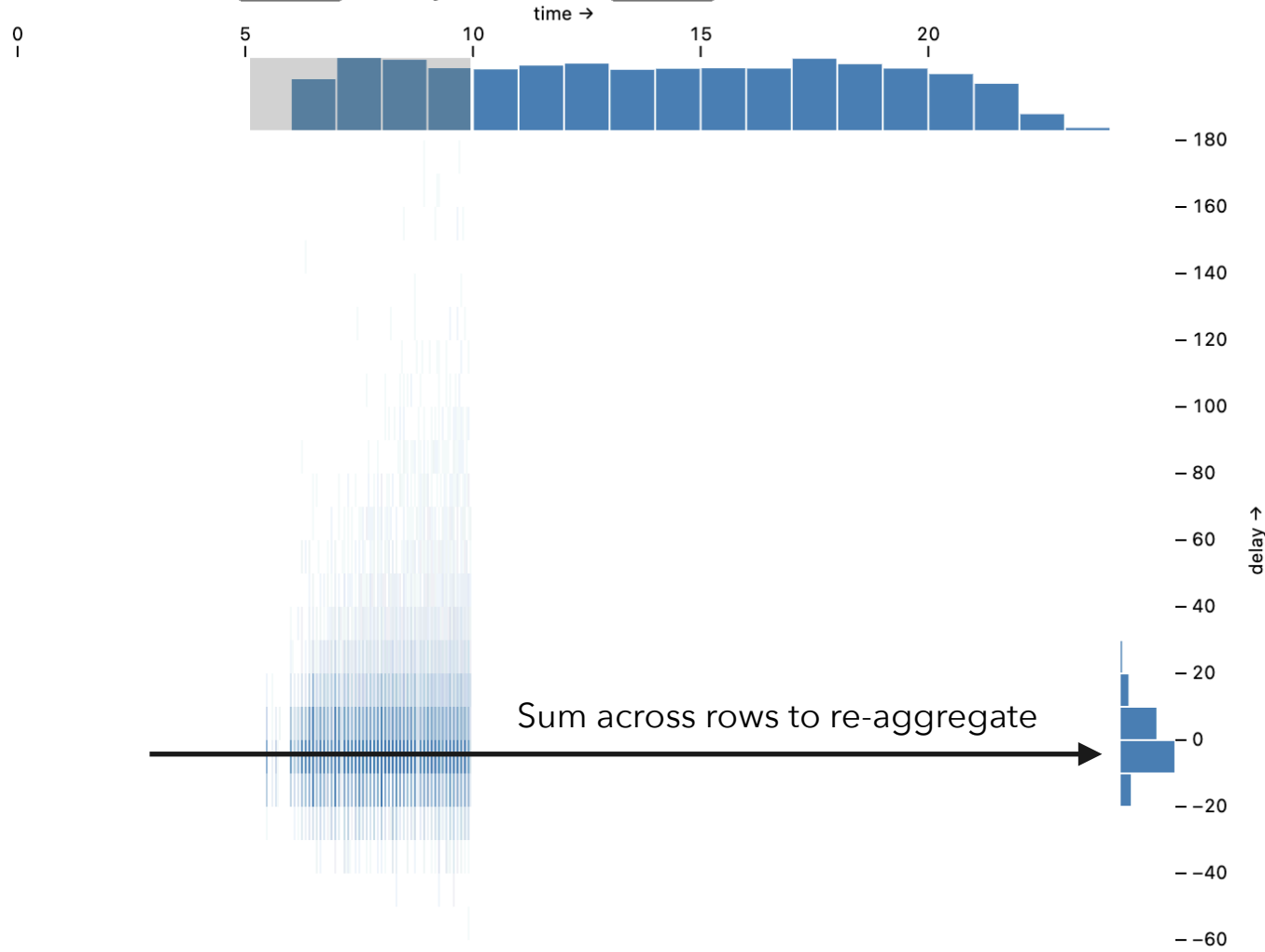
250k Records



Time Resolution Delay Resolution

Flight Delays

250k Records



Time Resolution

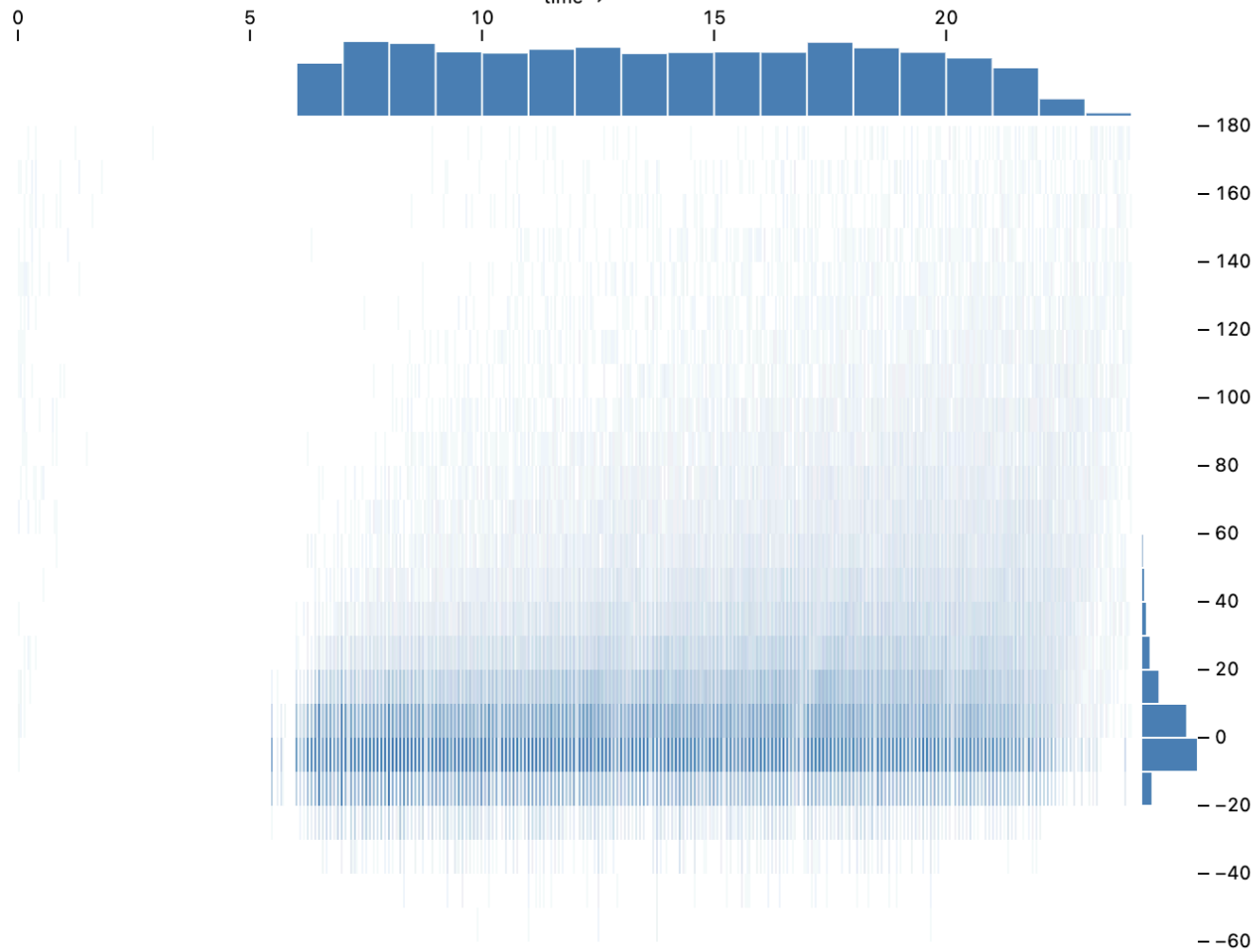
pixels ▾

Delay Resolution

bins ▾

Flight Delays

250k Records



Time Resolution

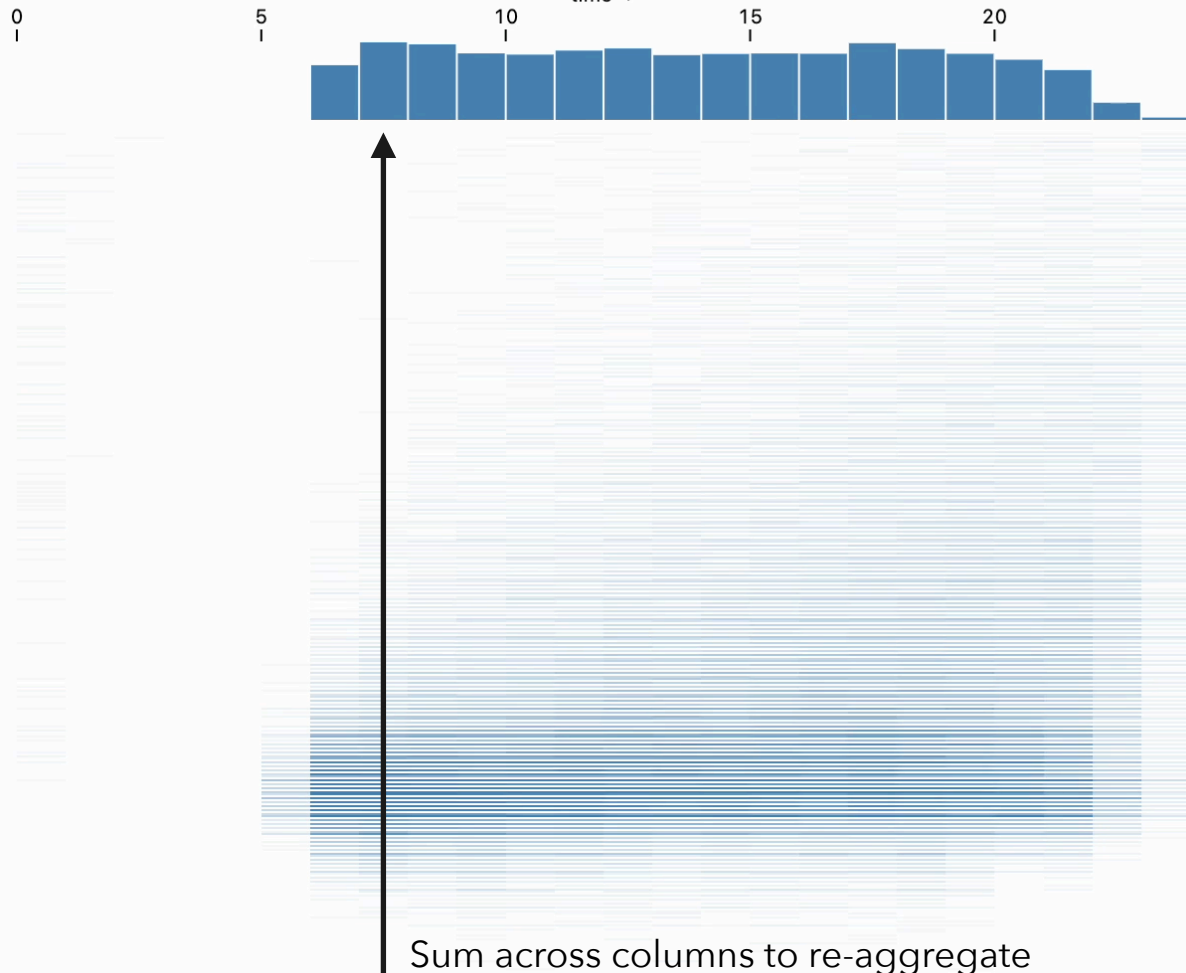
bins ▾

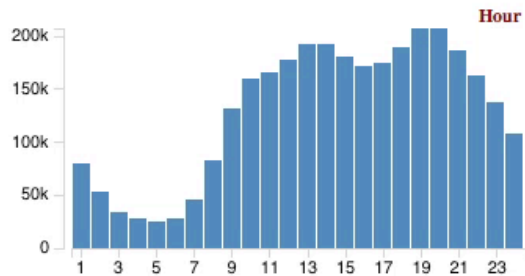
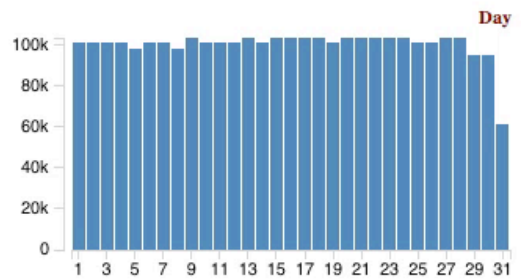
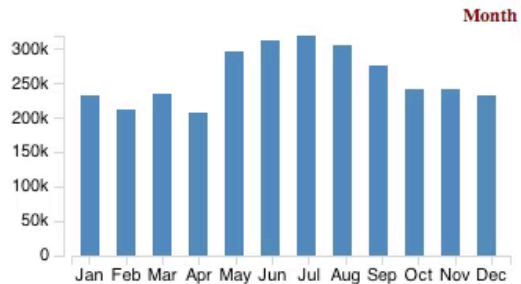
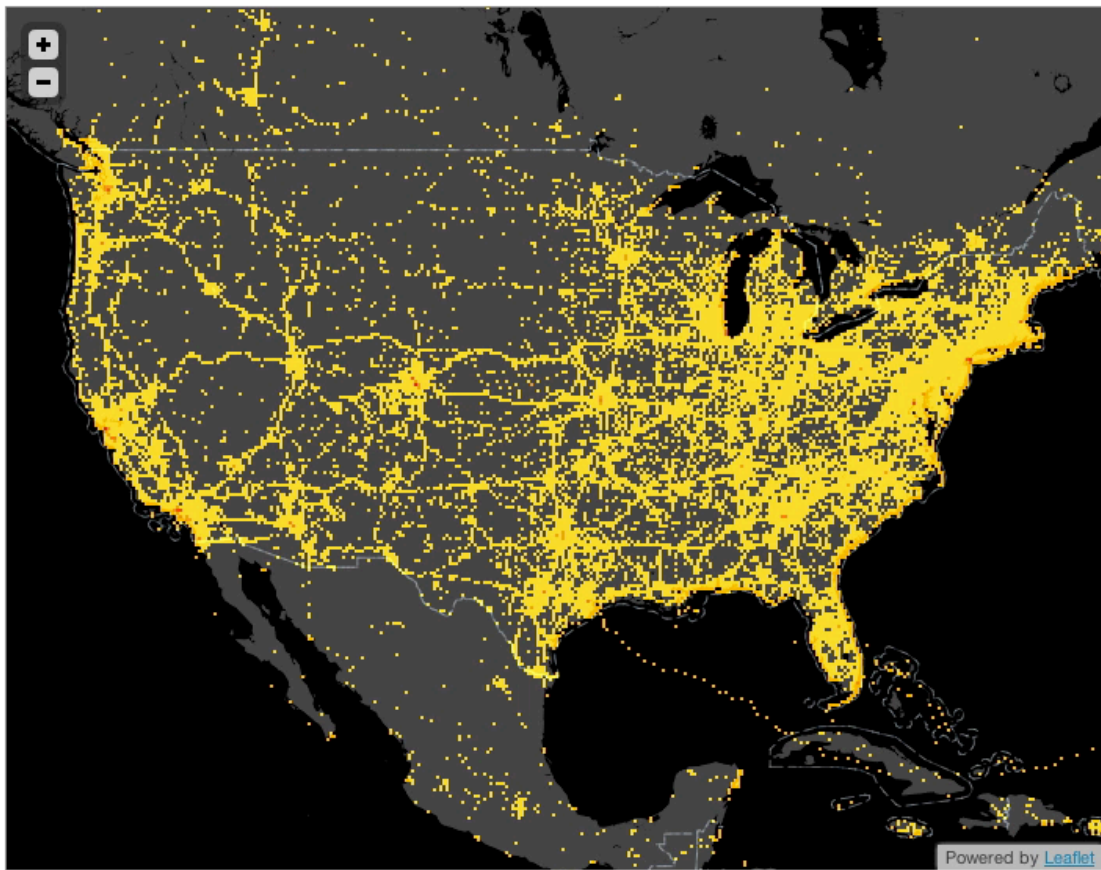
Delay Resolution

pixels ▾

Flight Delays

250k Records

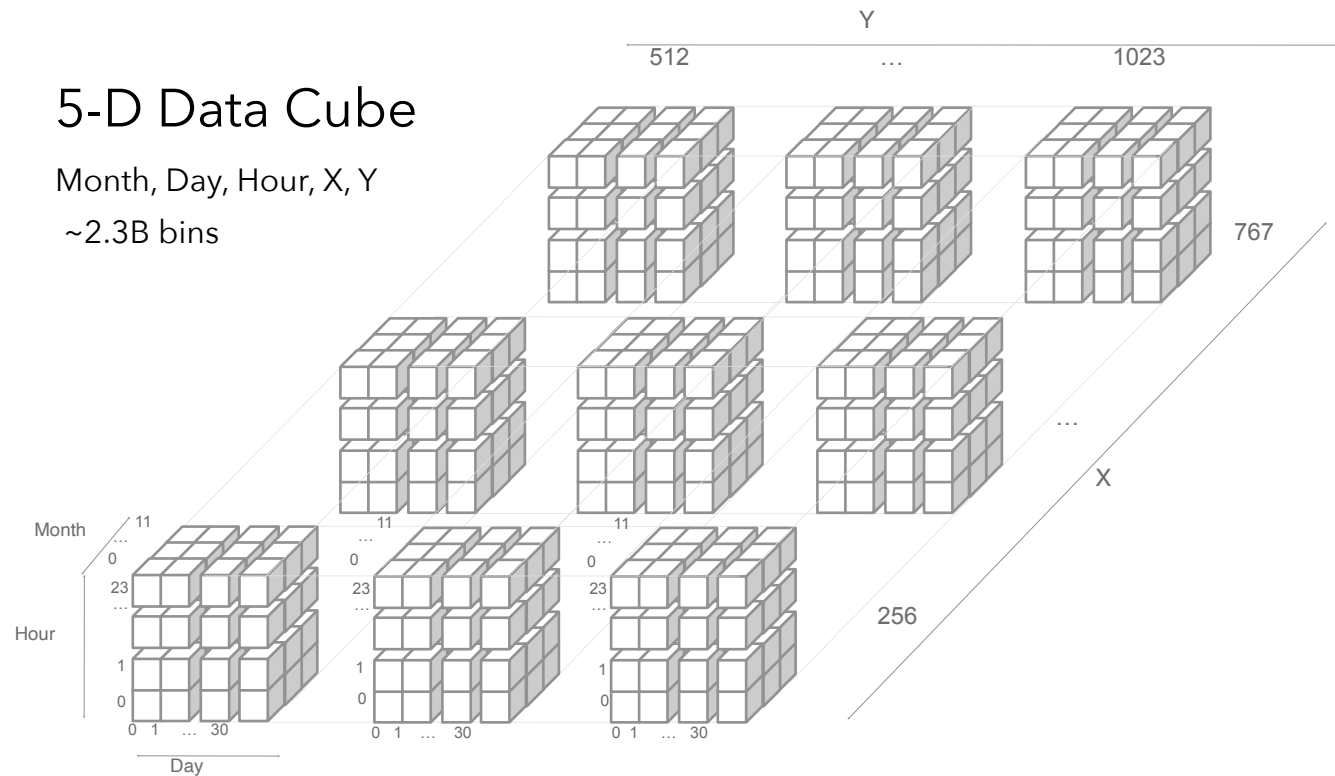




5-D Data Cube

Month, Day, Hour, X, Y

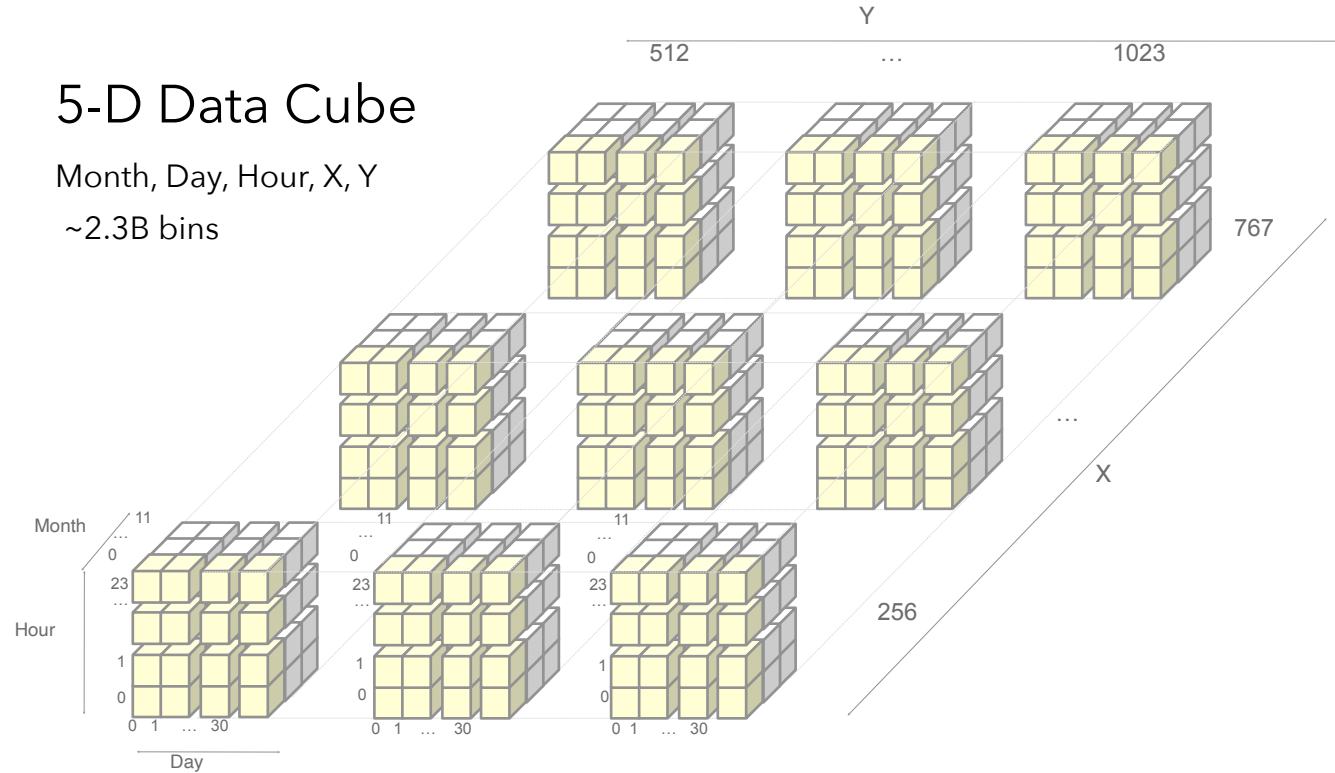
~2.3B bins

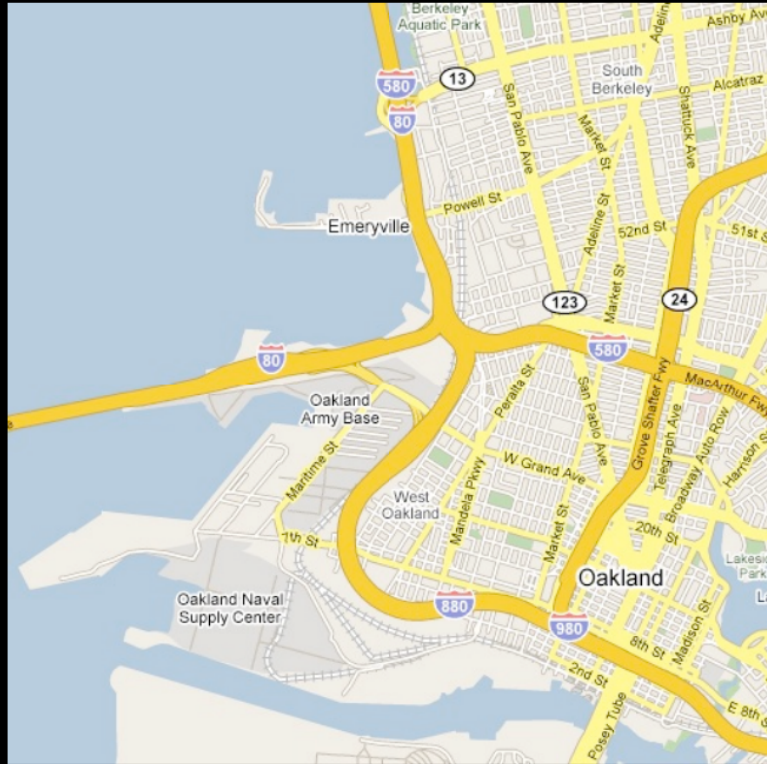


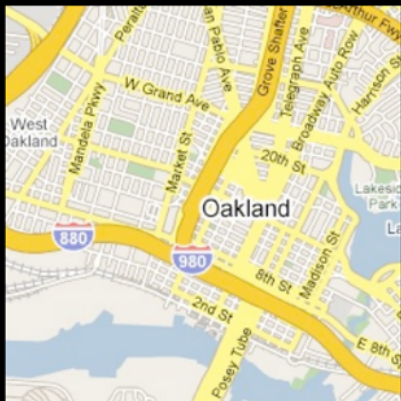
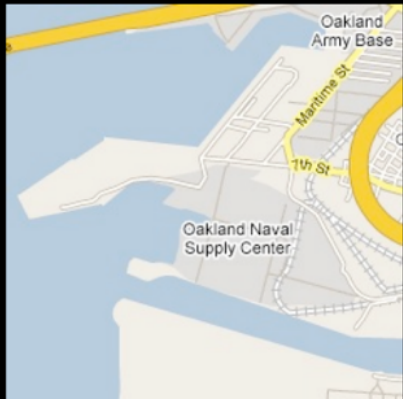
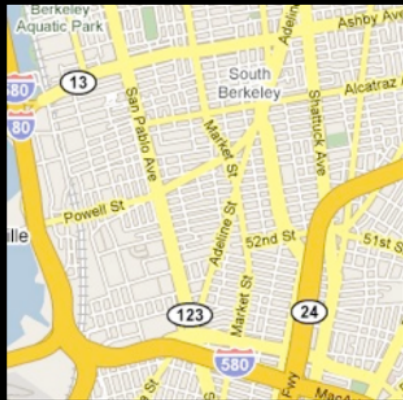
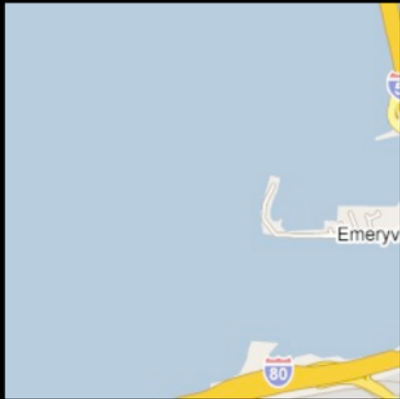
5-D Data Cube

Month, Day, Hour, X, Y

~2.3B bins

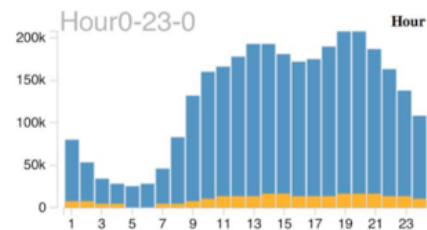
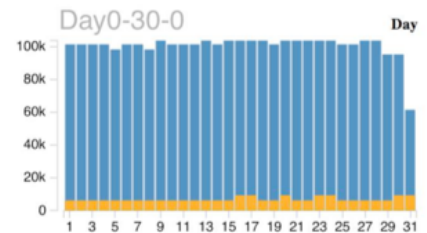
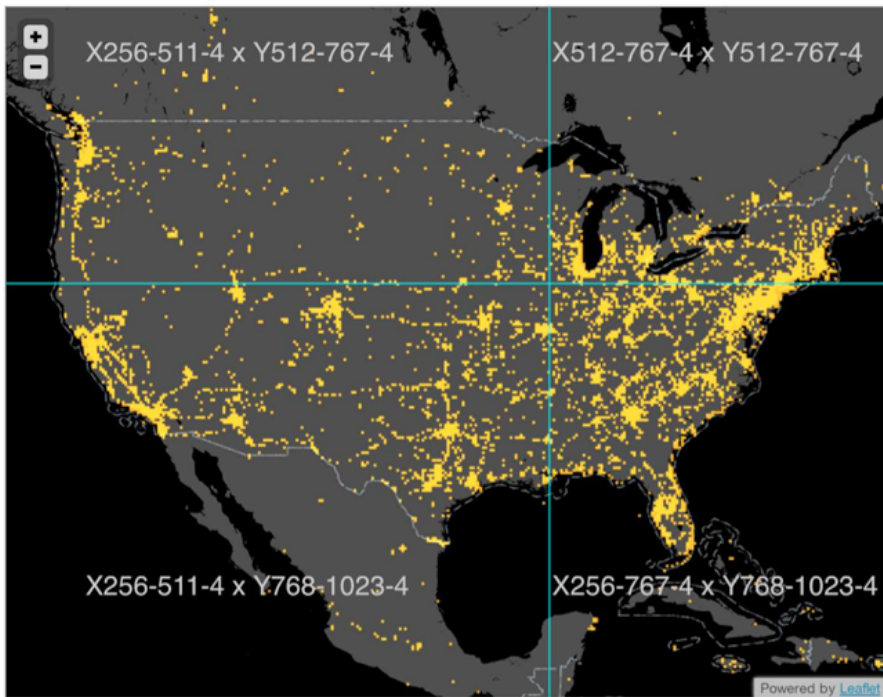


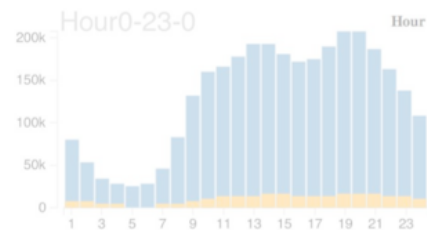
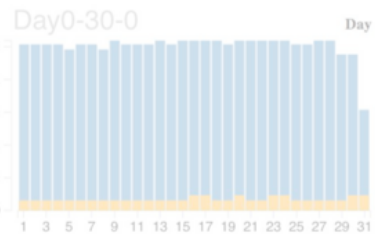
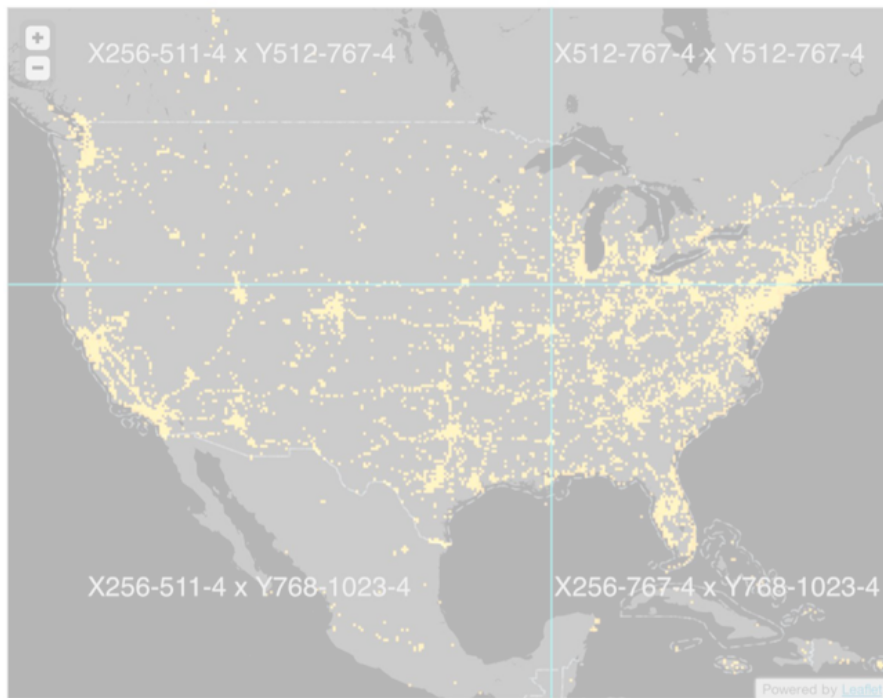


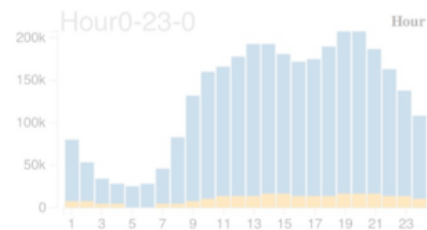
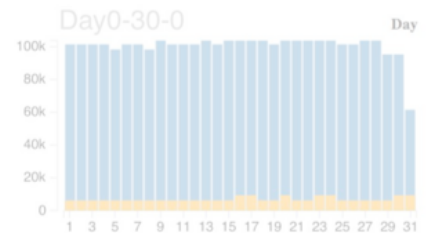
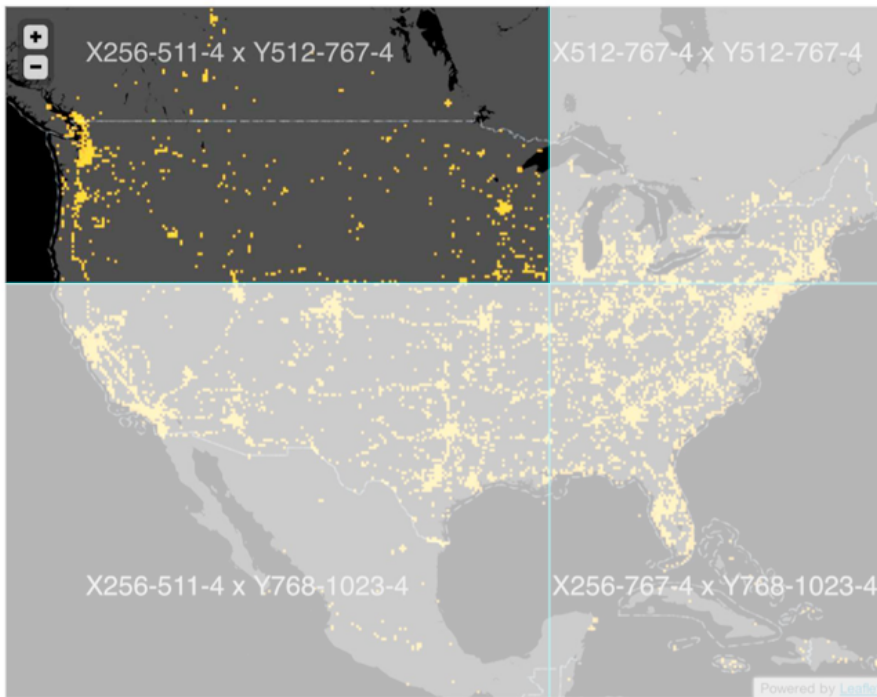


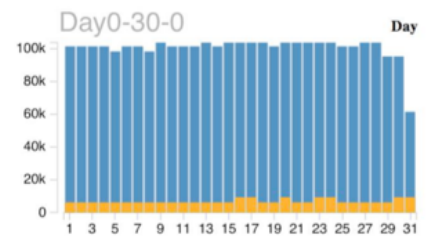
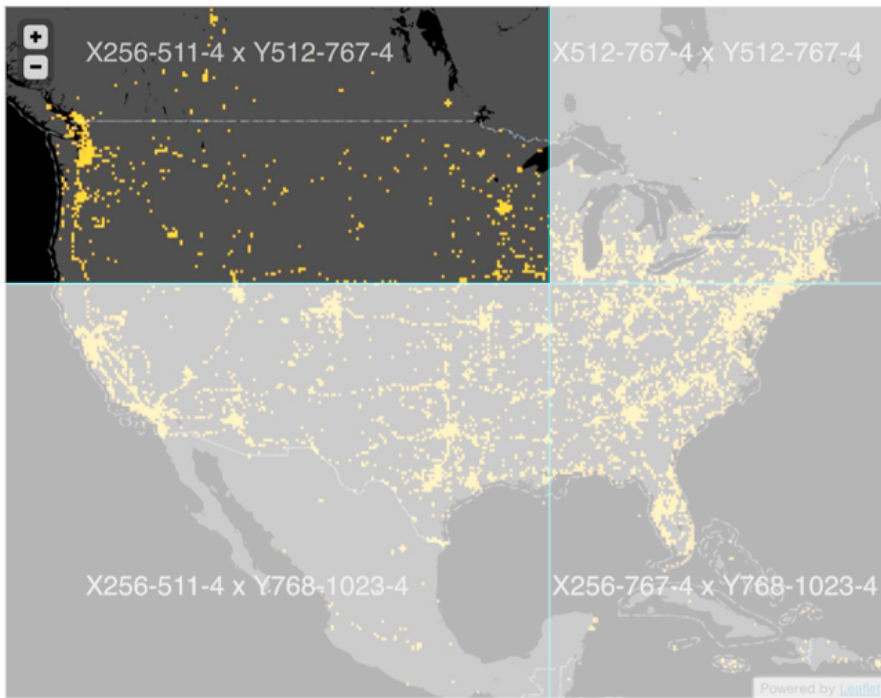
Multivariate Data Tiles

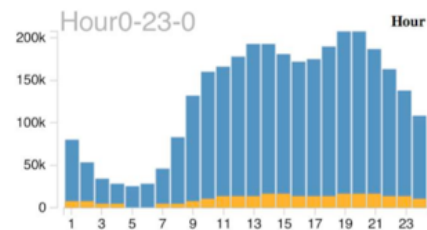
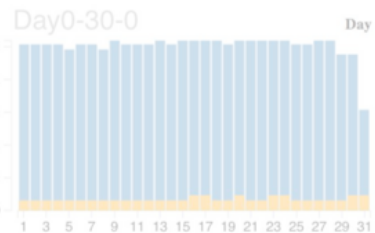
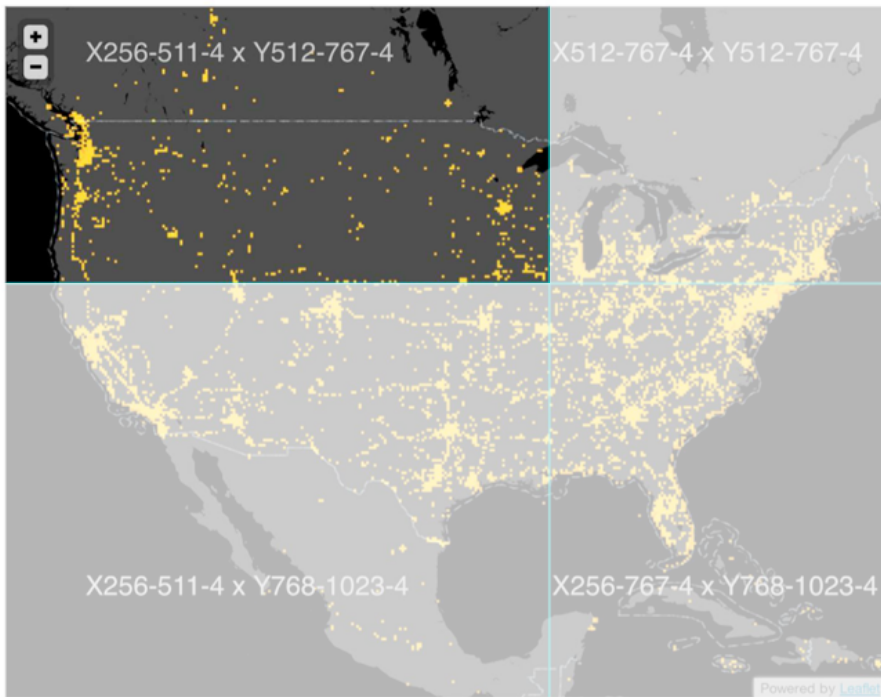
1. Send data, not pixels
2. Embed multi-dim data

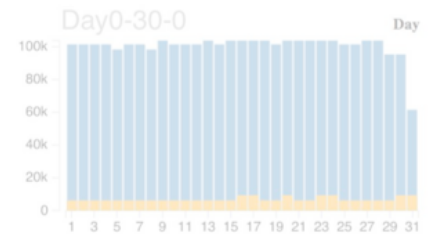
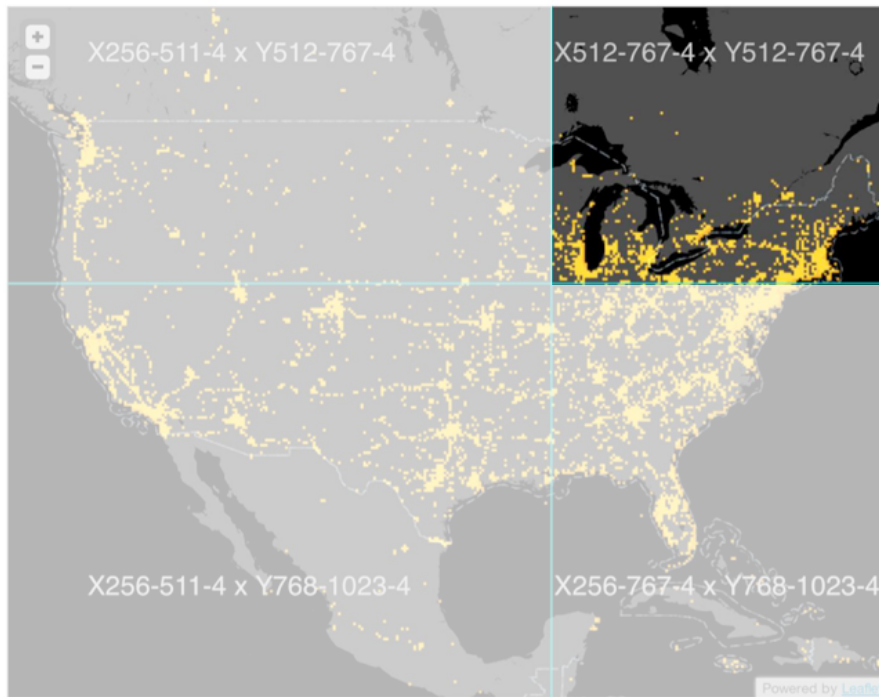


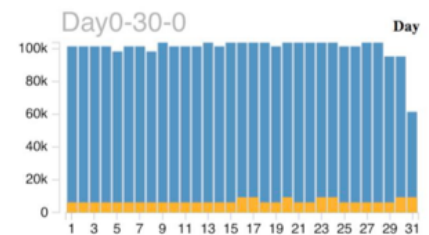
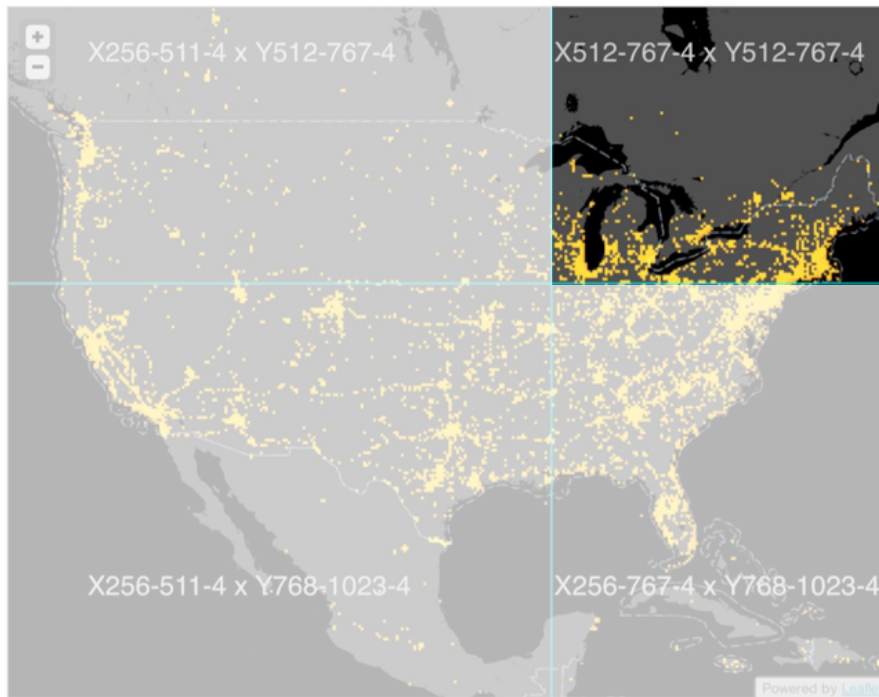


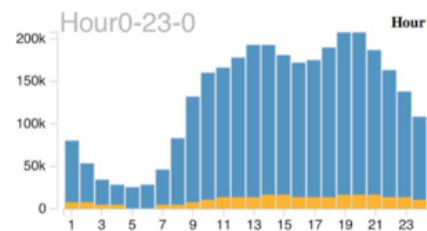
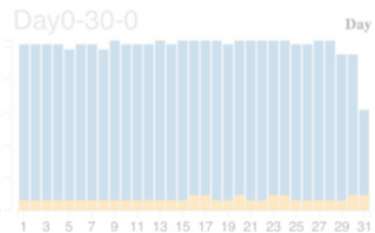
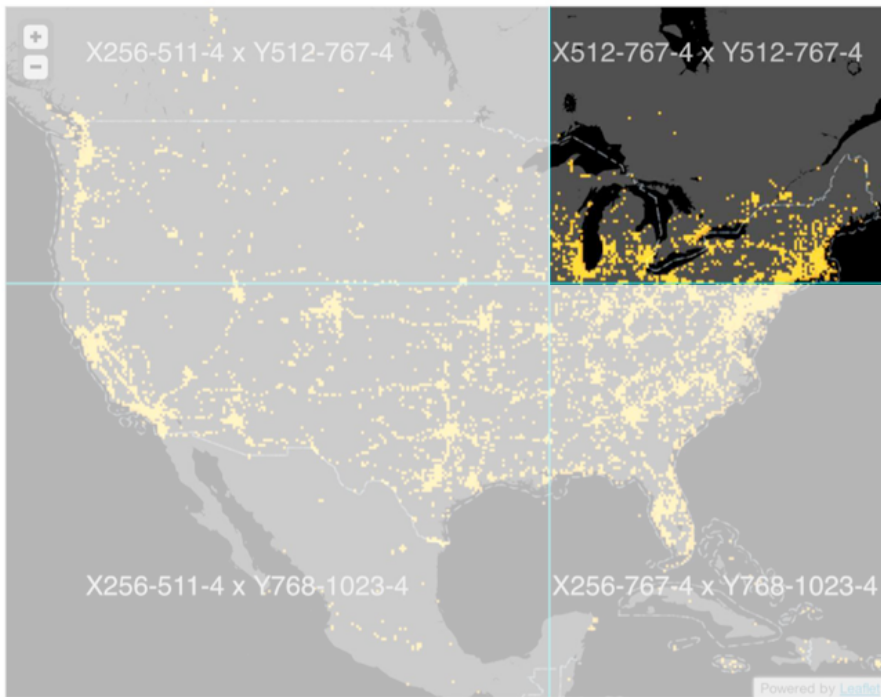


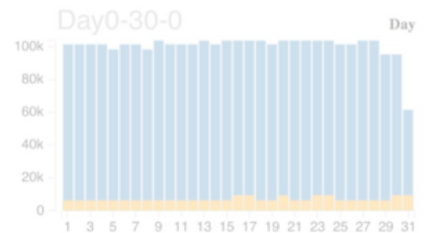
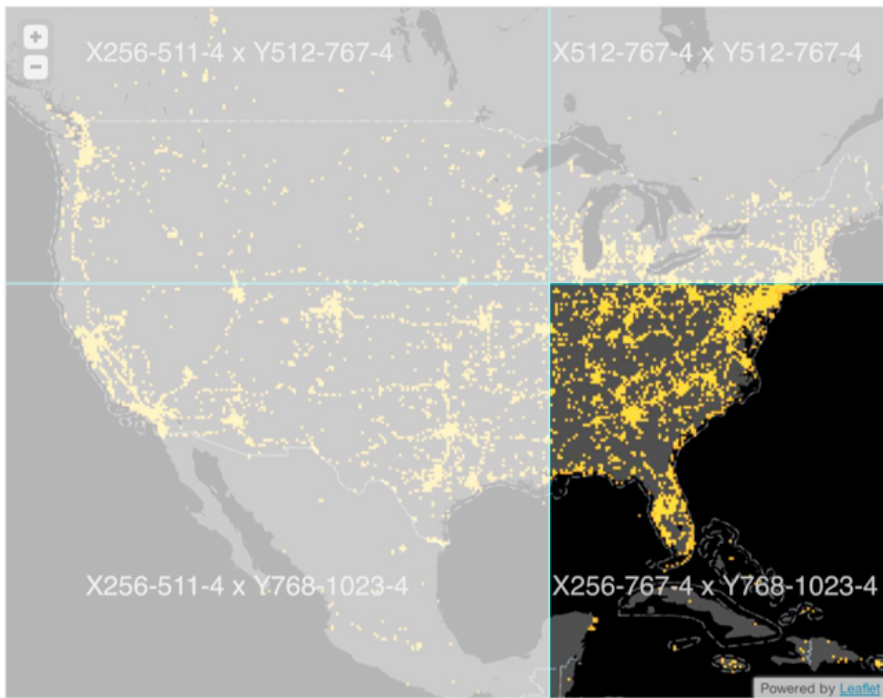


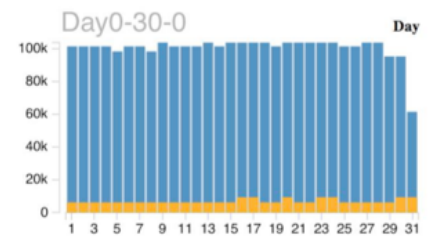
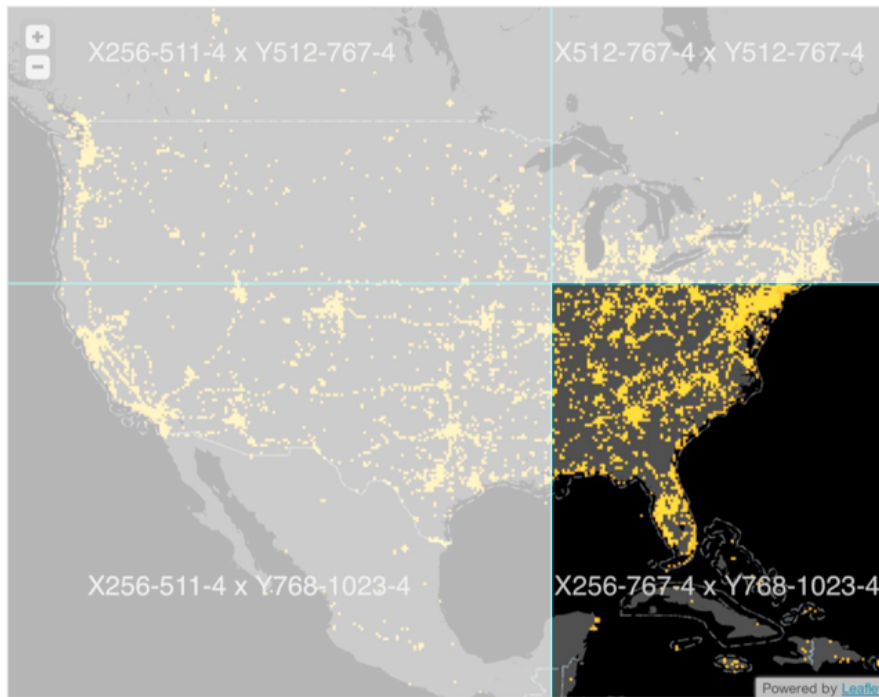


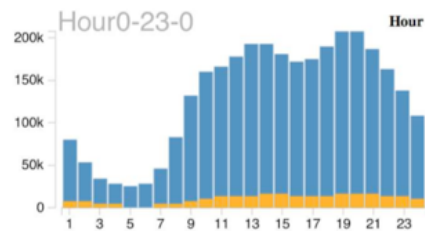
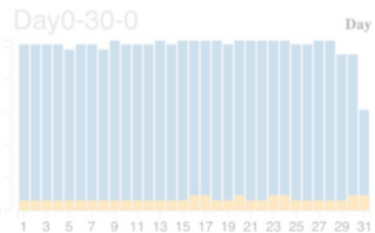
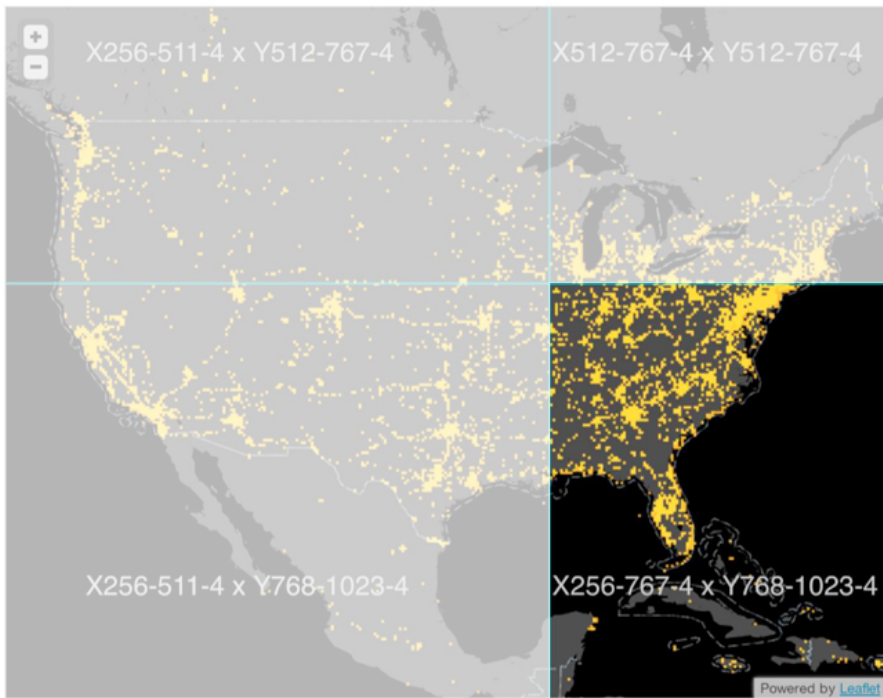


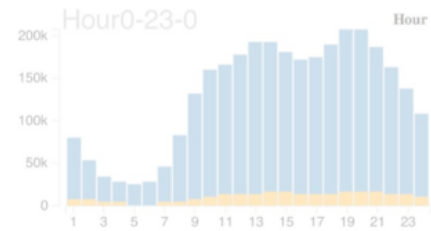
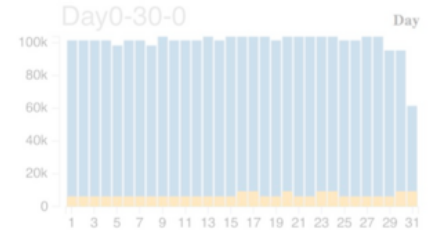
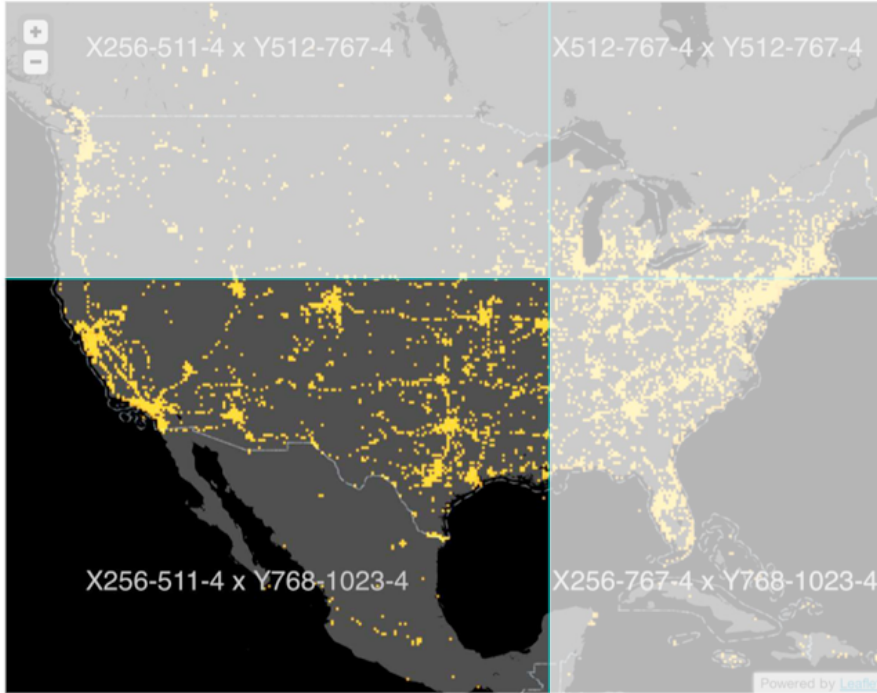


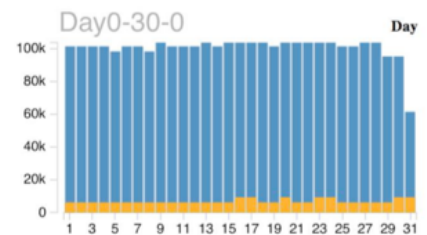
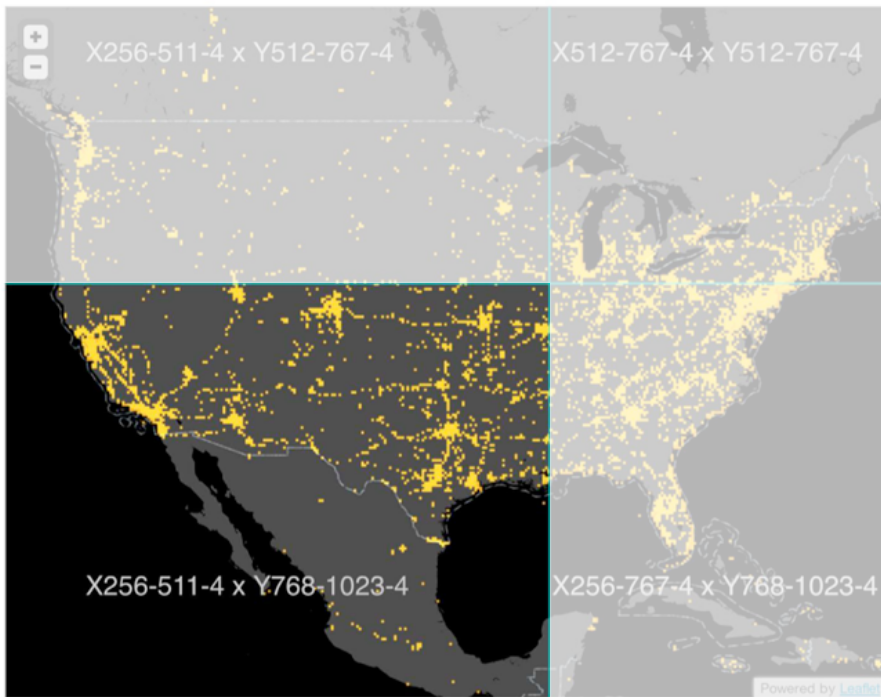


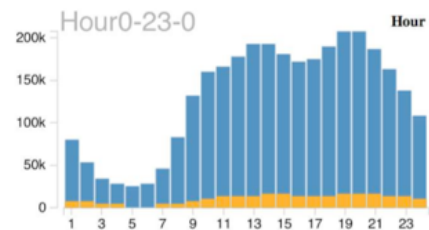
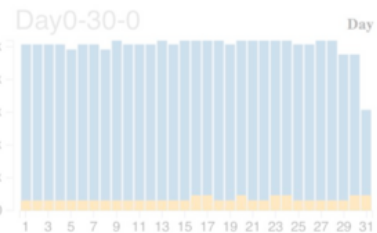
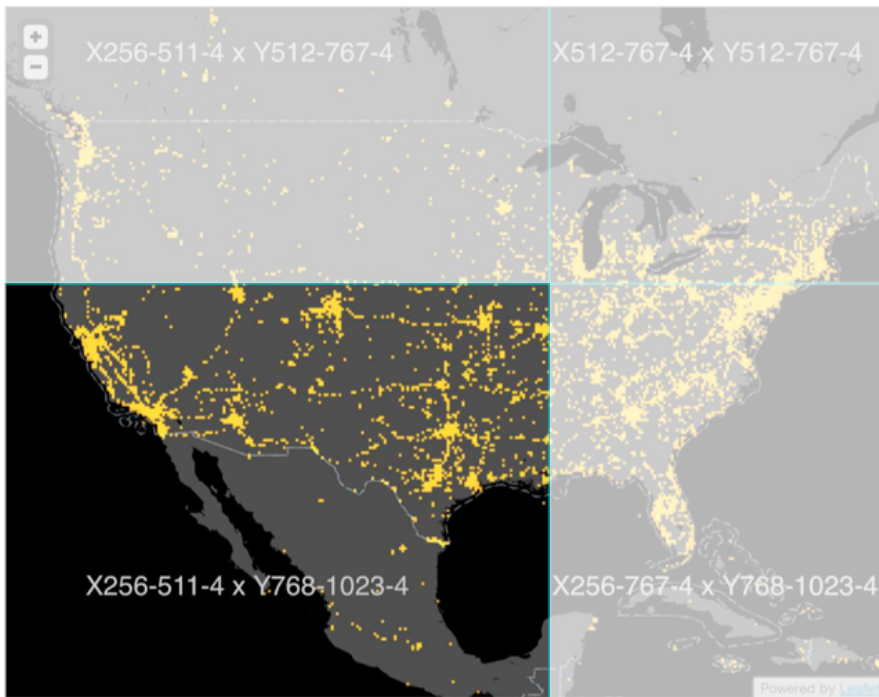


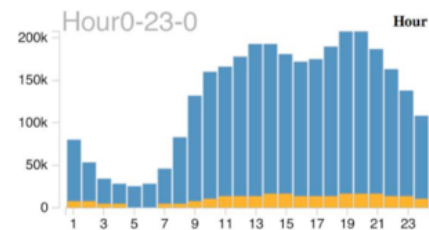
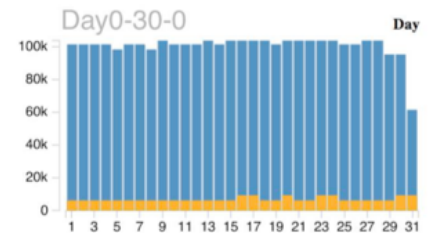
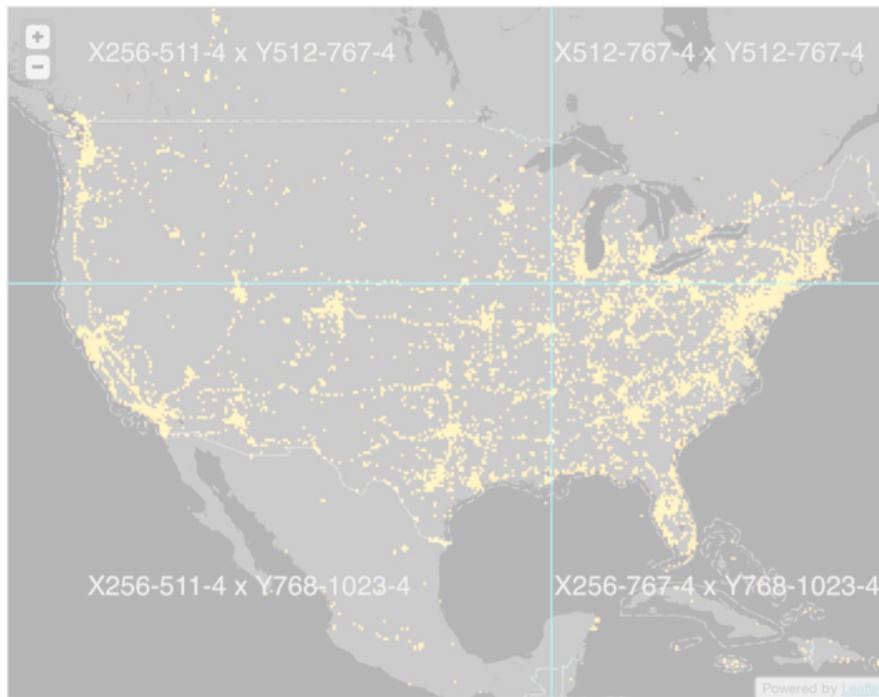


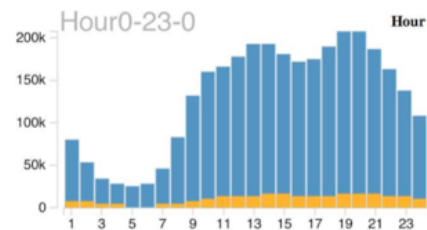
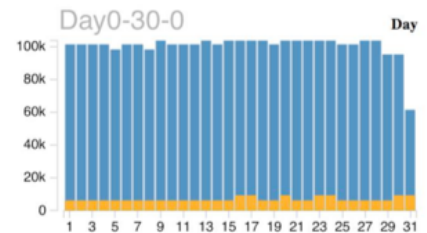
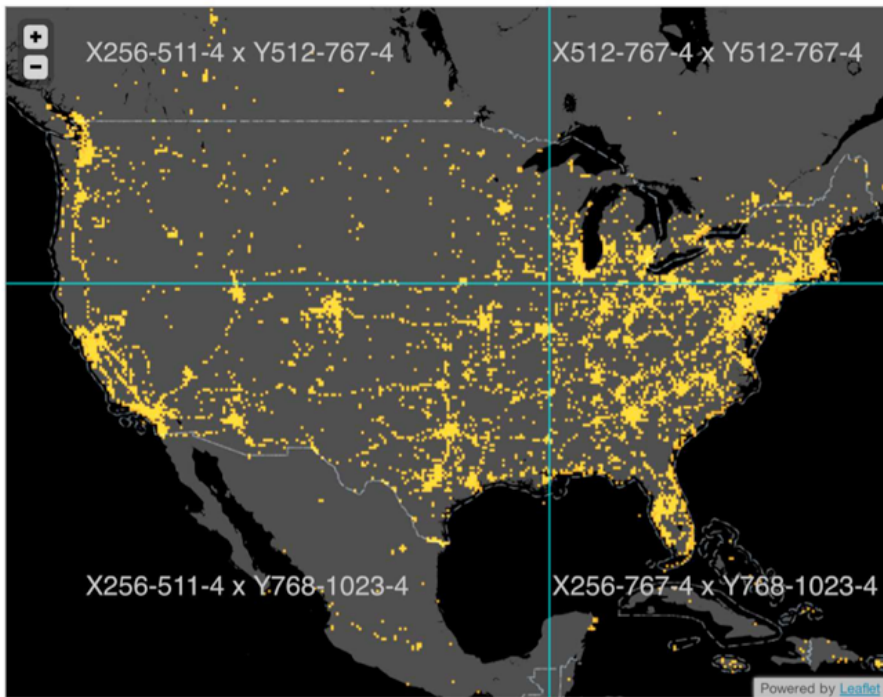




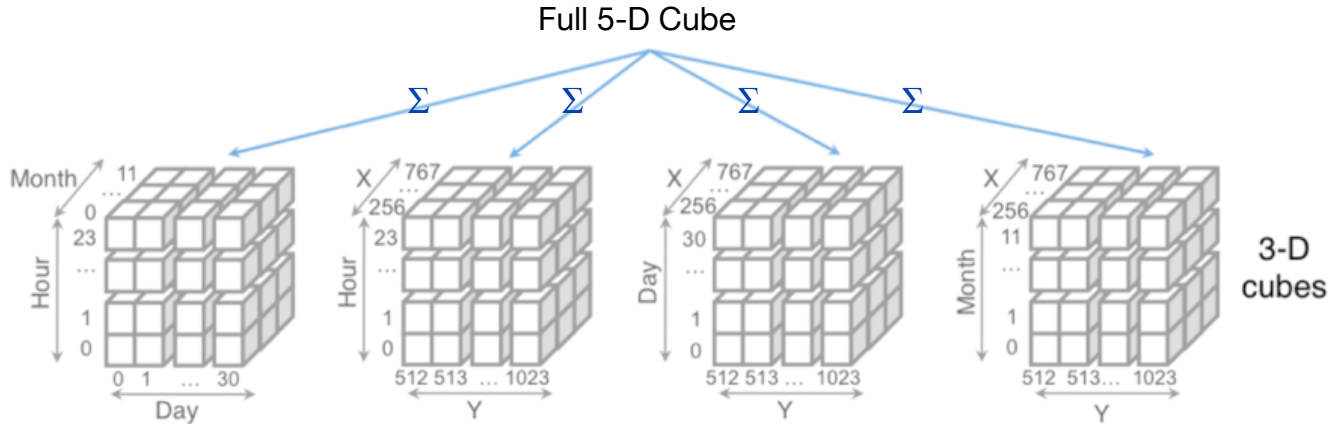






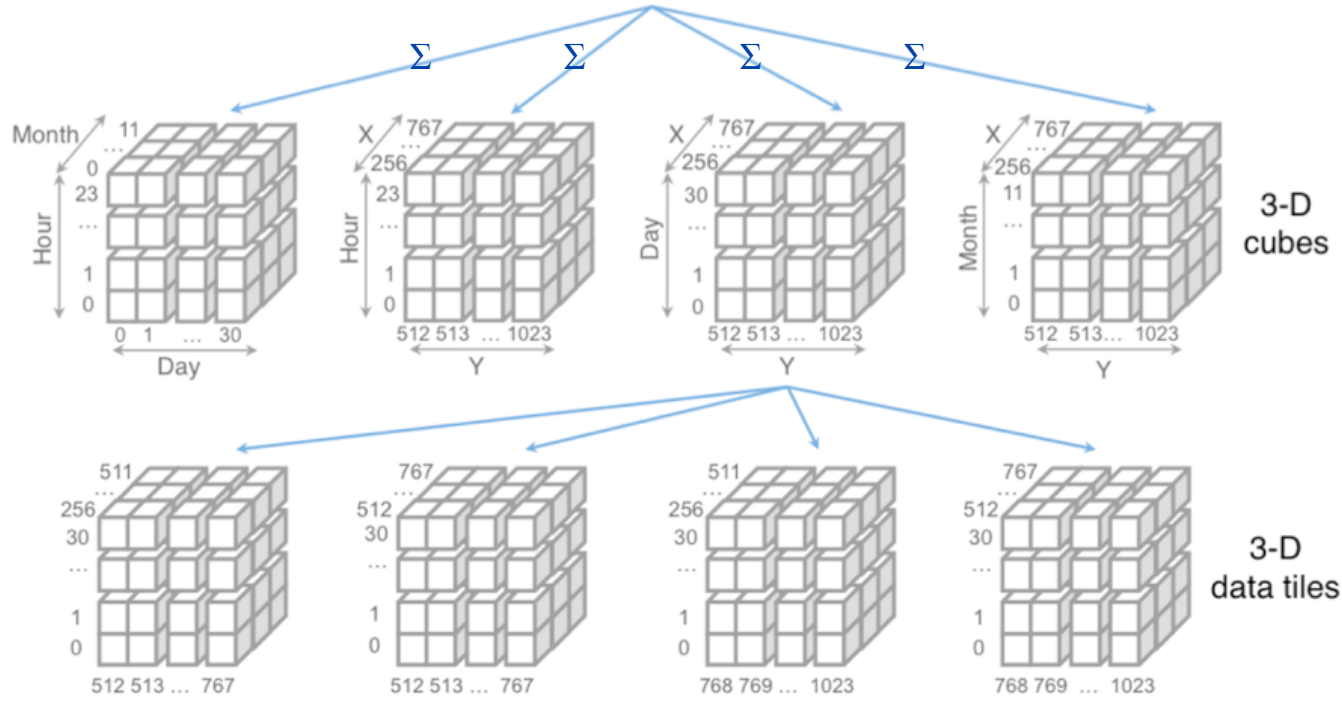


Full 5-D Cube

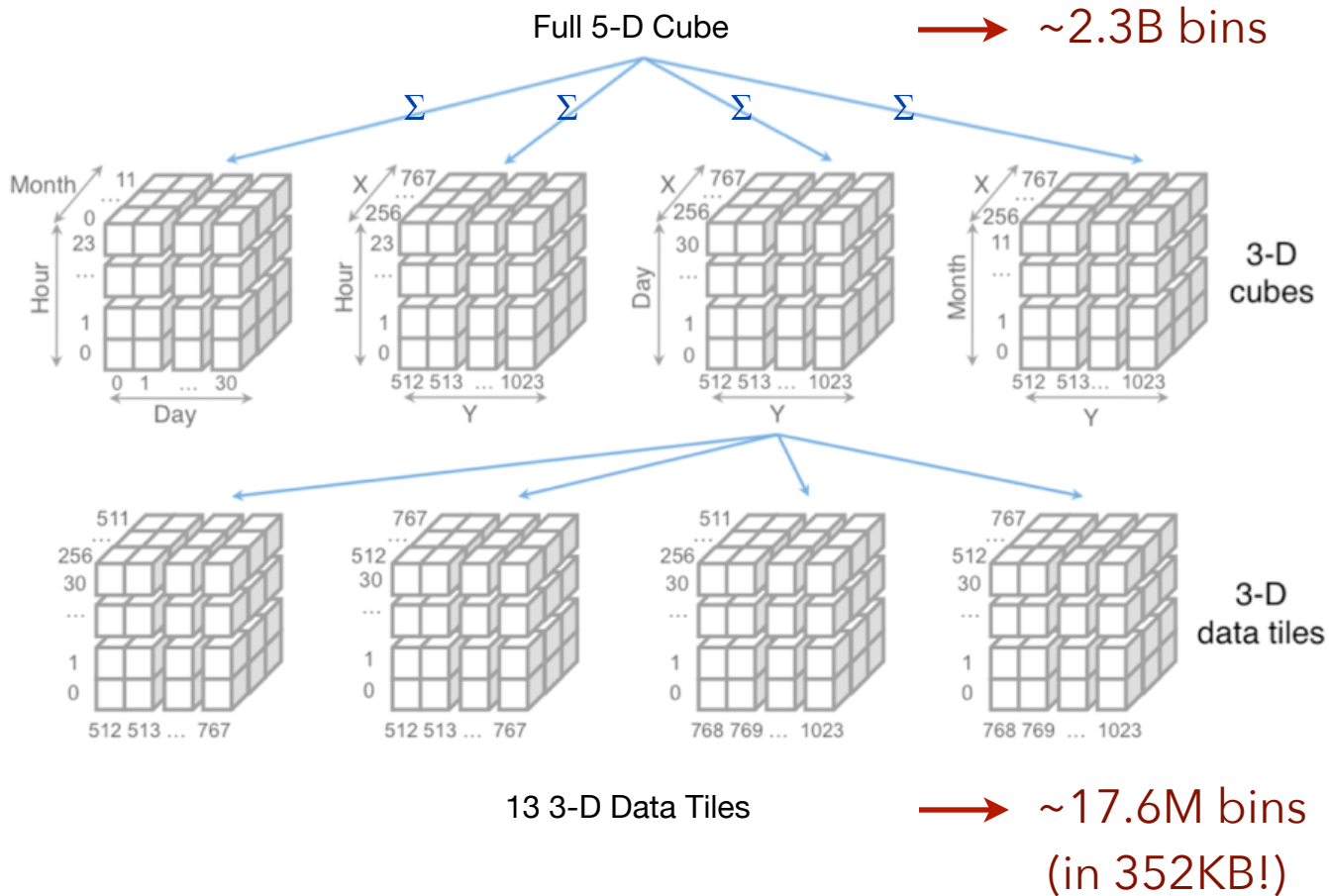


For any pair of 1D or 2D binned plots, the maximum number of dimensions needed to support brushing & linking is **four**.

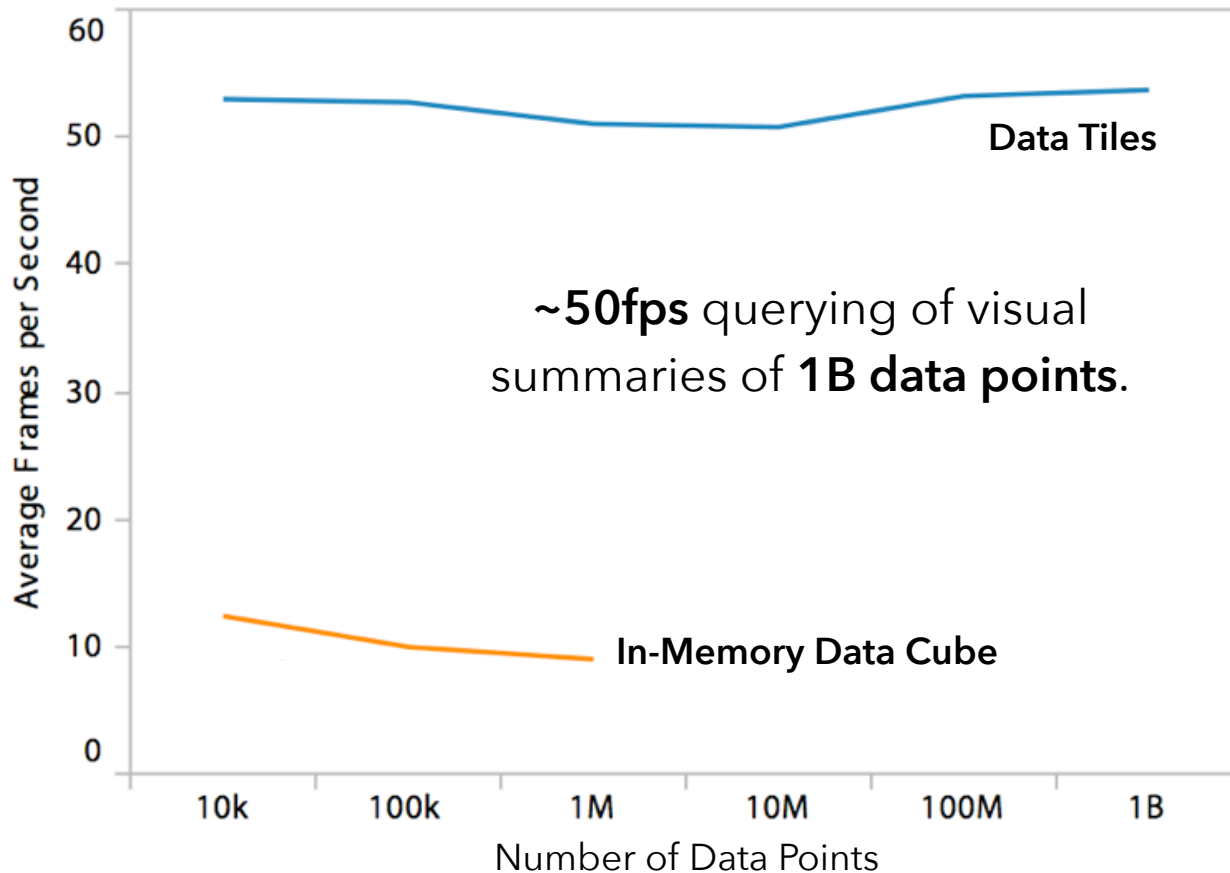
Full 5-D Cube



13 3-D Data Tiles



5 dimensions x 50 bins/dim x 25 plots



Limitations and Questions

But where do the multivariate data tiles come from?

They must be computed, either ahead of time or on-the-fly. Up to the 100M point range, an analytic database can do this on the fly. In the 1B point range, pre-computation avoids delays.

We can also *prefetch*: we can start computing new data tiles as soon as the pointer enters a chart, before a selection is made.

Does super-low-latency interaction really matter?

Is it worth it to go to all of this trouble? (Short answer: yes!)

High latency leads to reduced analytic output [Liu & Heer, InfoVis 2014]

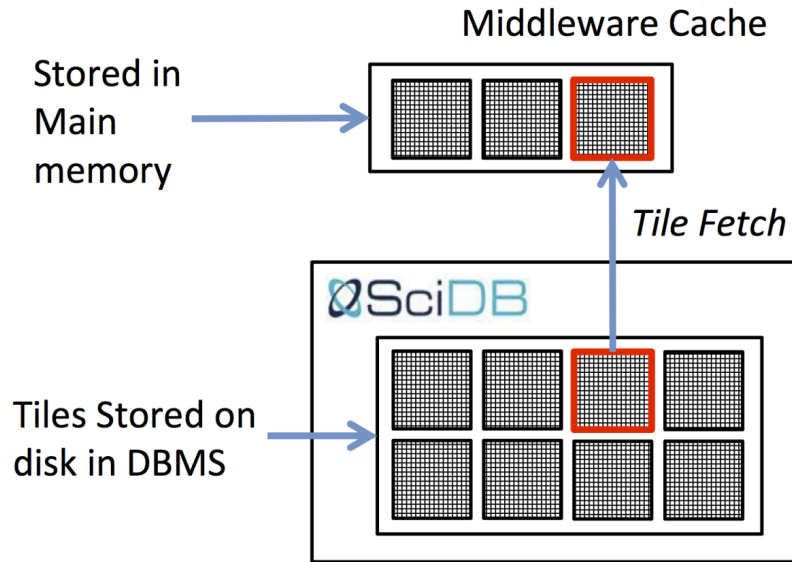
Prefetching

ForeCache

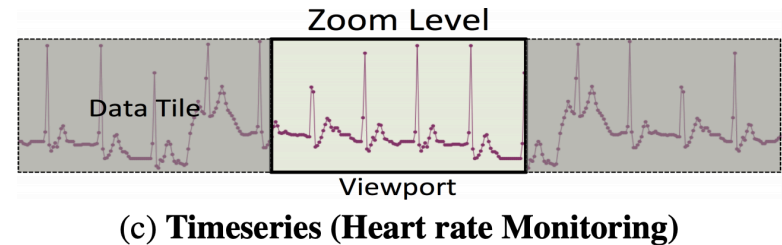
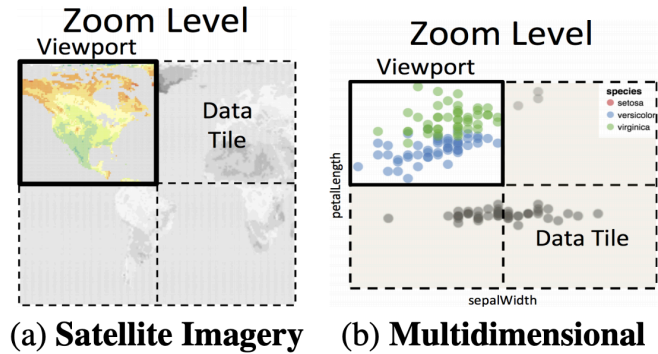
[Battle, Chang, & Stonebraker '16]

Strategies: Query Database, Data Cubes, Prefetching

ForeCache is also a Data Tile-Based System



Manage a Cache of Tiles from DBMS



Example Tile-Based Views

Key Idea: Model & Predict User Behavior

1. Classify the User's Analysis Phase

Foraging: Searching for patterns of interest

Sensemaking: Closely examine a region-of-interest (ROI)

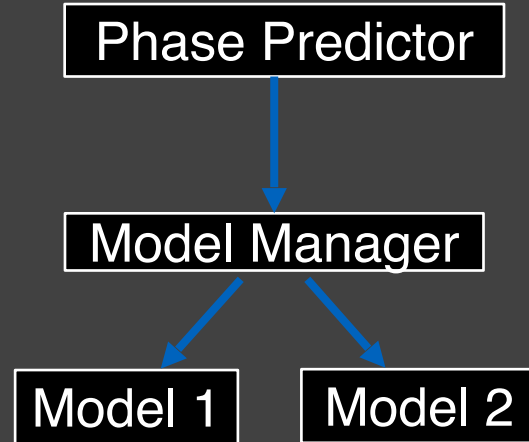
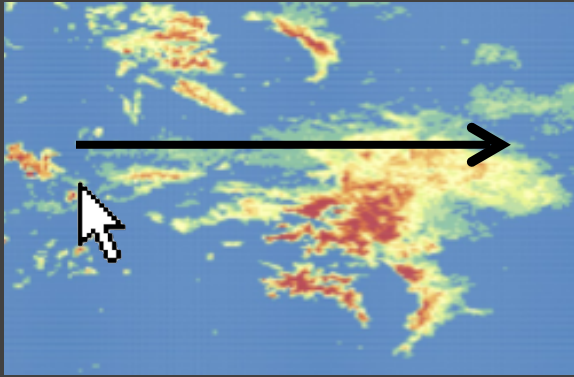
Navigation: Transition between levels of detail

2. Predict Which Data Tiles Will be Requested

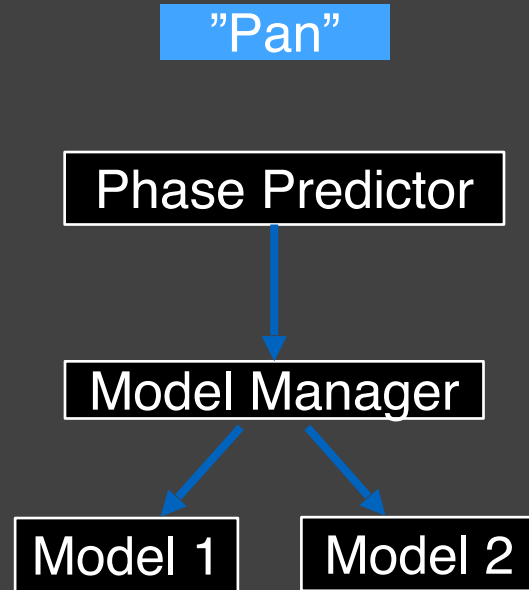
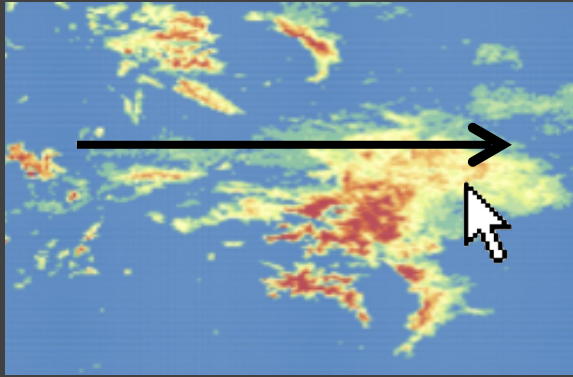
Train a machine learning classifier (SVM) to predict phase.

The input data is the activity trace of user interactions.

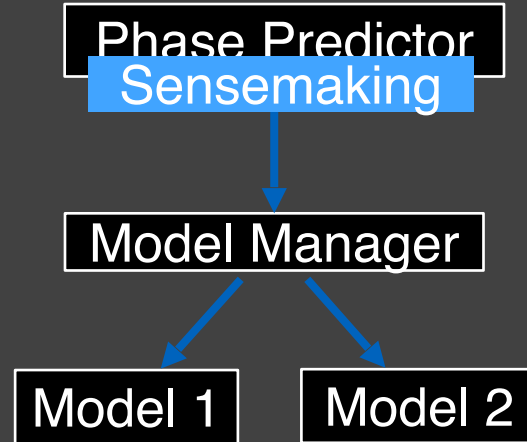
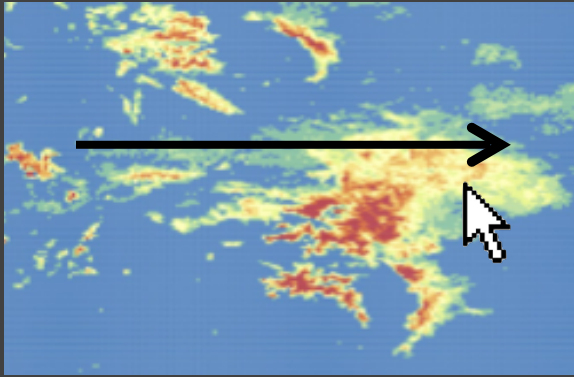
Using Phases to Predict Tiles



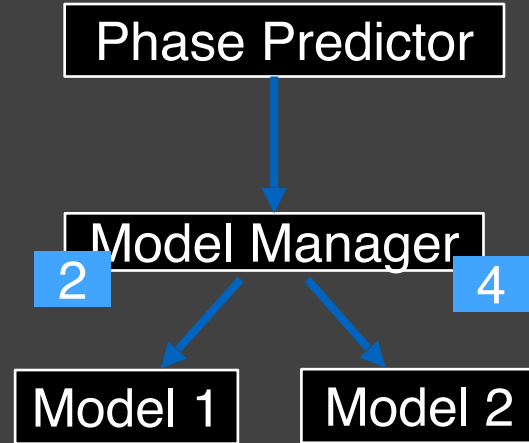
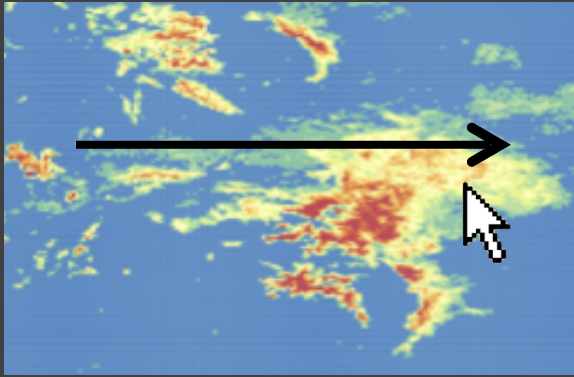
Using Phases to Predict Tiles



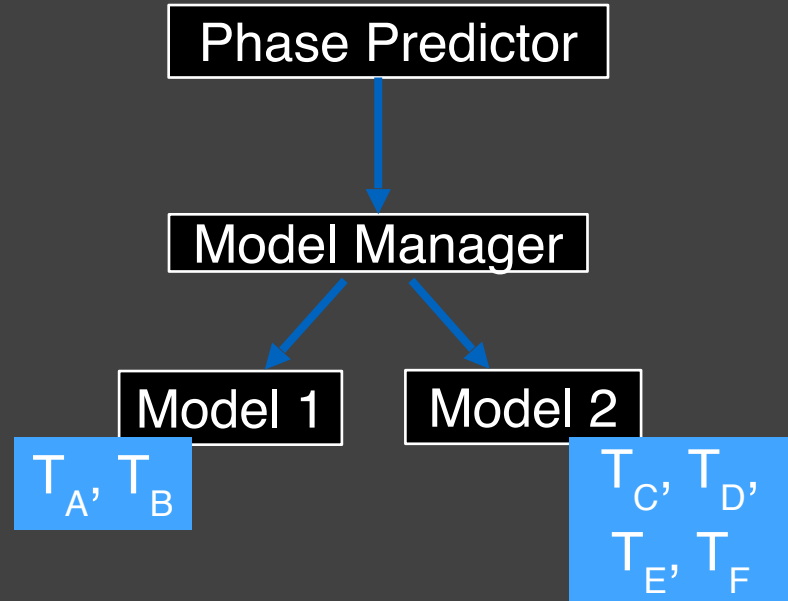
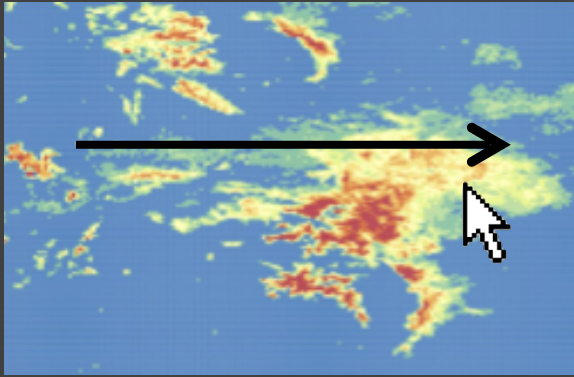
Using Phases to Predict Tiles



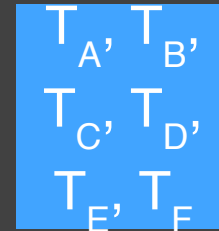
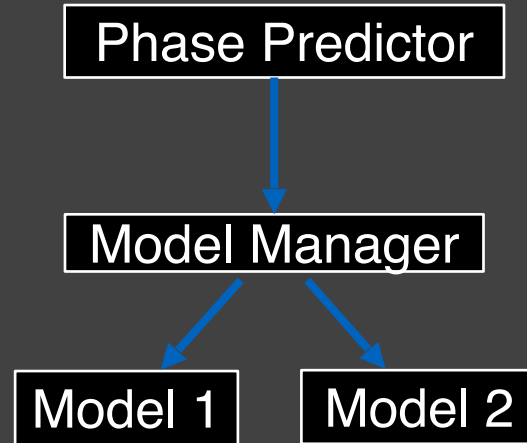
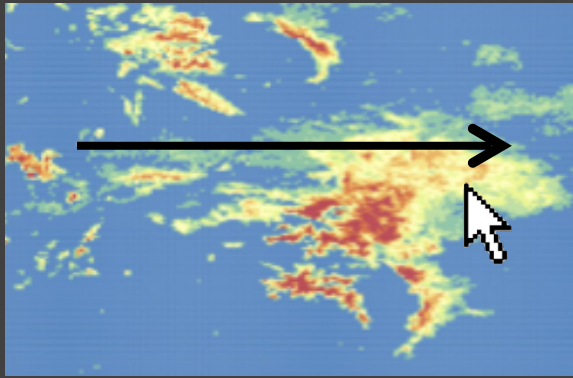
Using Phases to Predict Tiles



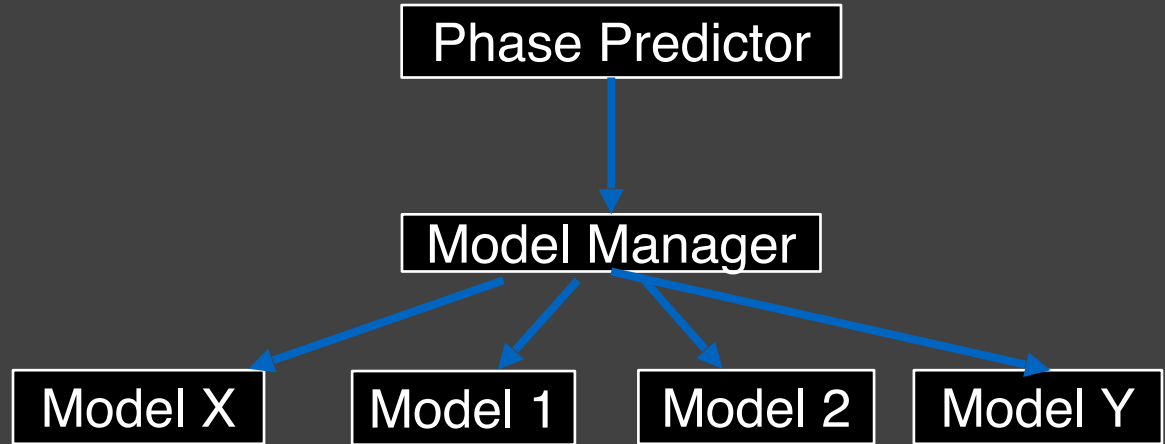
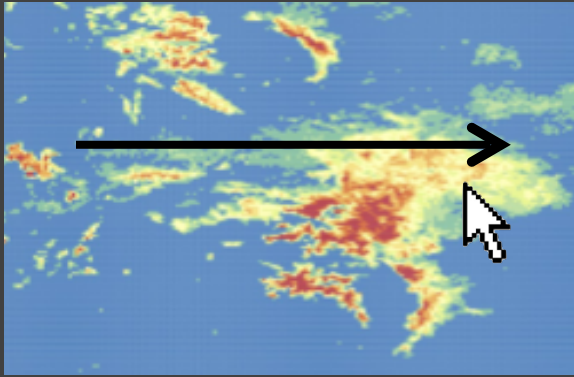
Using Phases to Predict Tiles



Using Phases to Predict Tiles



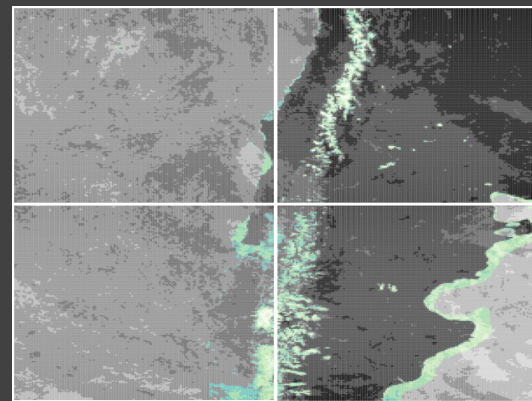
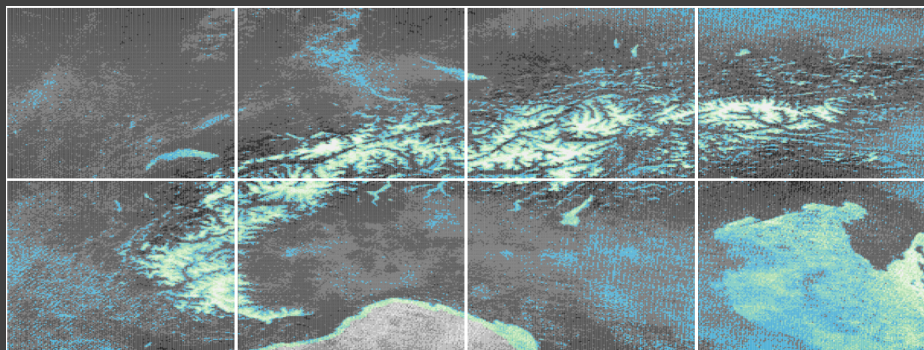
Using Phases to Predict Tiles



Evaluating ForeCache: A User Study

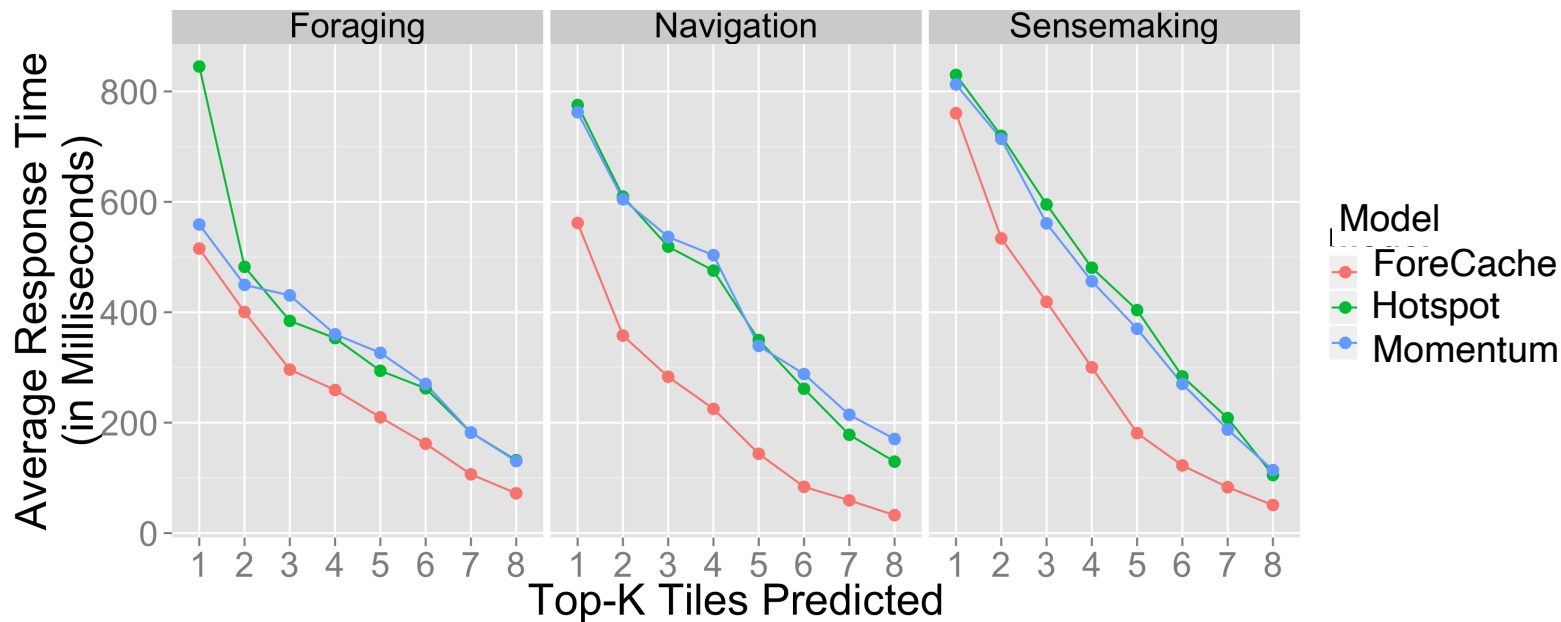
Participants: 18 earth science researchers

Explored NASA MODIS snow cover queries



Results: ForeCache was 20% More Accurate and 88% Faster than Existing Pre-fetching Methods

ForeCache vs. Existing Techniques



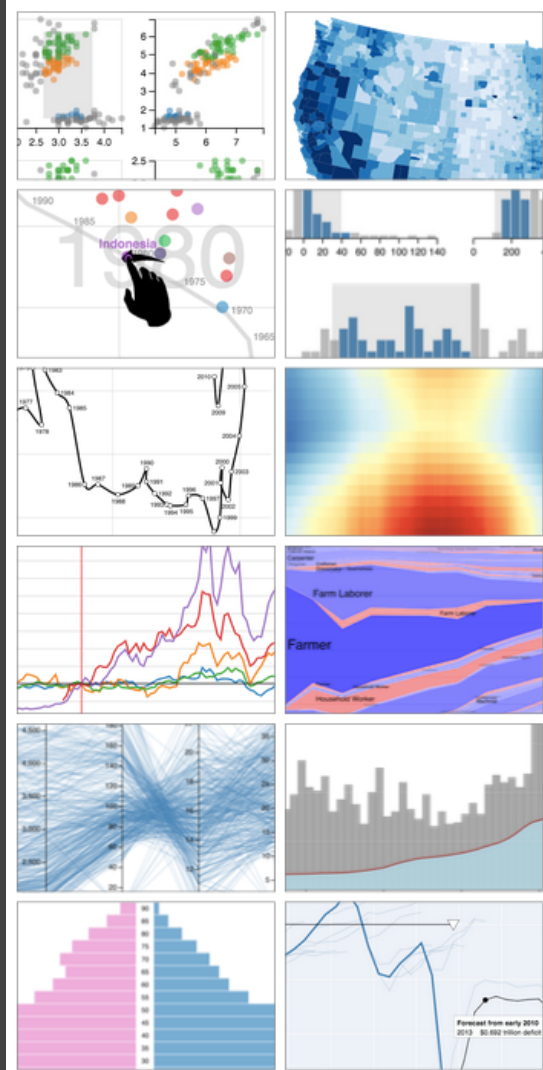
Sampling

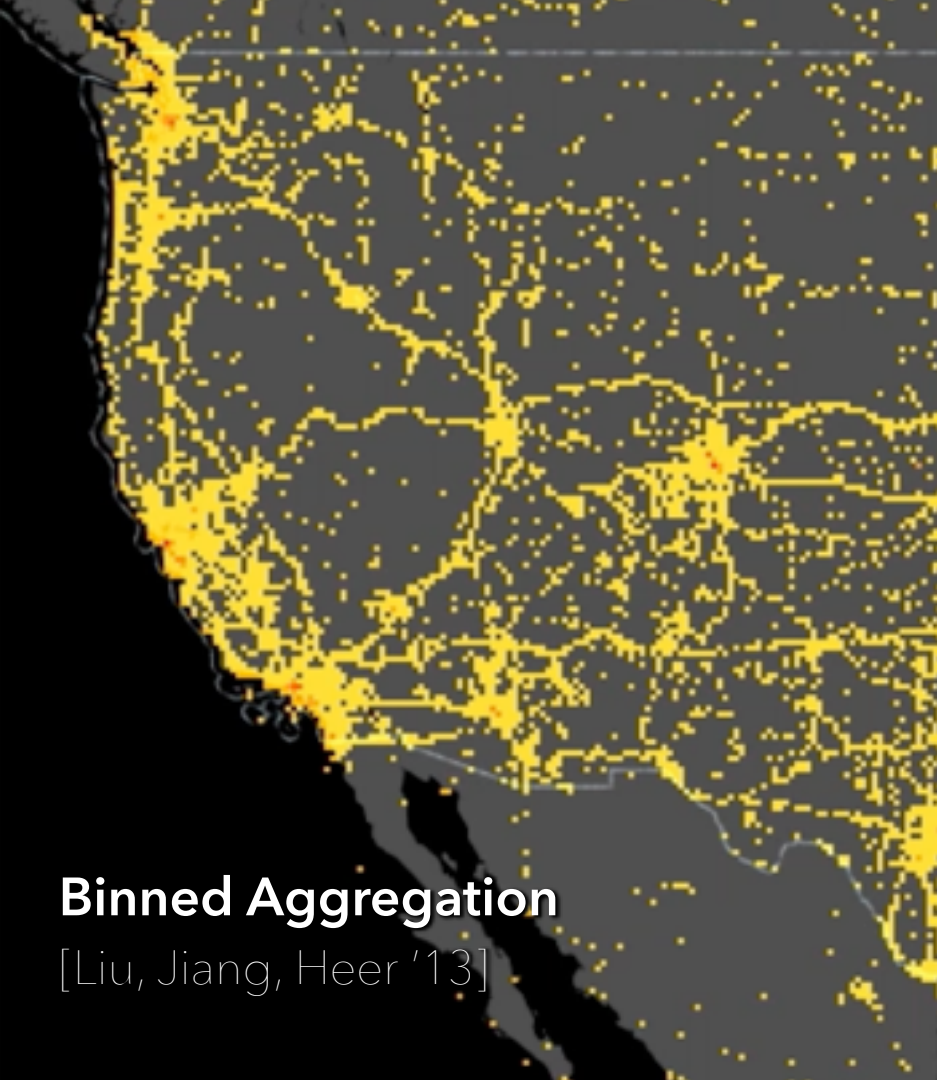
Common Sampling Methods

First-N: Useful for transformation, but not inference.

Random: Good default, but may miss features of interest. Possible in one pass via reservoir sampling, or faster if stored in randomized order.

Stratified: Sample within groups, ensure coverage and balance across those categories.





Binned Aggregation

[Liu, Jiang, Heer '13]



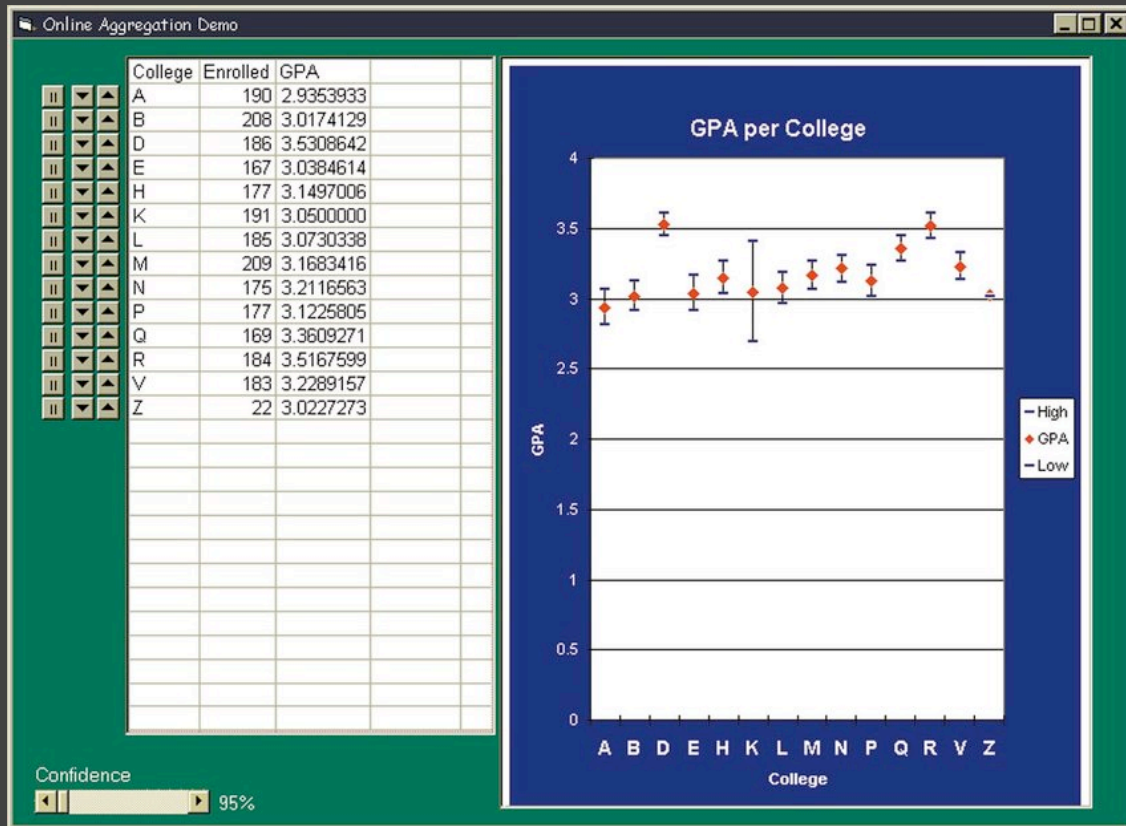
Sampling Google Fusion Tables

Online Aggregation [Hellerstein, Haas, Wang '97]

Provide dynamic, *progressive* results as queries run: see results over growing samples.

Visualize current results with confidence intervals to convey uncertainty of estimate.

Challenge: difficult to ensure truly random sampling.

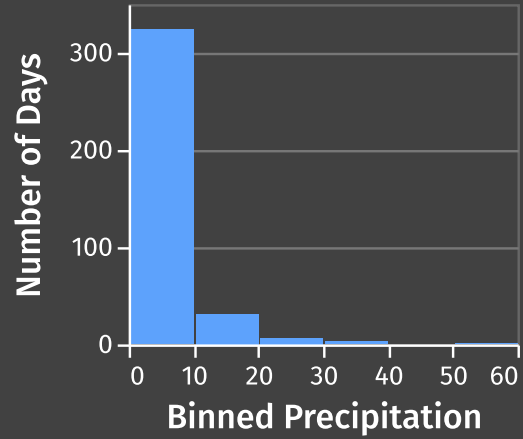
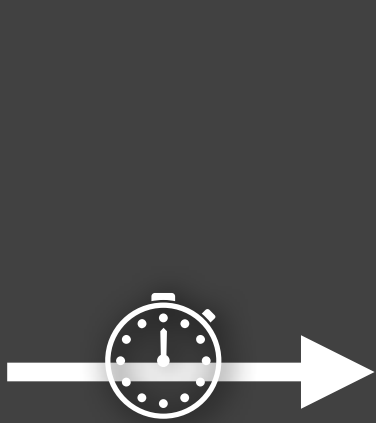


**What if data is too large to
query in a reasonable time?**

Trust, but Verify: Optimistic Vis

[Moritz, Fisher, Ding & Wang '17]

Strategies: Query Database, Approximation



An alternative perspective:

Trust the approximation for now and
verify the final results later.

Optimistic Visualization

Trust but Verify. Moritz et al. *CHI 2017*.



1. Analysts uses initial estimates.
2. Precise queries run in the background.
3. System confirms results. Analyst detects errors.

Analysts can use approximations and also trust them.

Optimistic Visualization

Data: FAAData

Heatmap

Type to filter schema...

- # Year
- # Quarter
- # Month
- # DayOfMonth
- # DayOfWeek
- FlightDate
- A UniqueCarrier
- # AirlineID
- A Carrier
- A TailNum
- # FlightNum
- # OriginAirportID
- # OriginAirportSeqID
- # OriginCityMarketID
- A Origin
- A OriginCityName
- A OriginState
- A OriginStateFips
- A OriginStateName
- # OriginWac
- # DestAirportID
- # DestAirportSeqID
- # DestCityMarketID
- A Dest
- A DestCityName
- A DestState
- A DestStateFips
- A DestStateName
- # DestWac

X-Axis
Field: DepDelay
Binning: 64
Sort by key:

Y-Axis
Field: ArrDelay
Binning: 40
Sort by key:

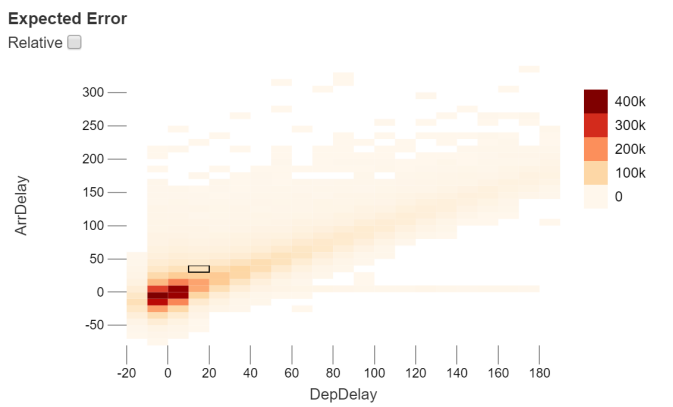
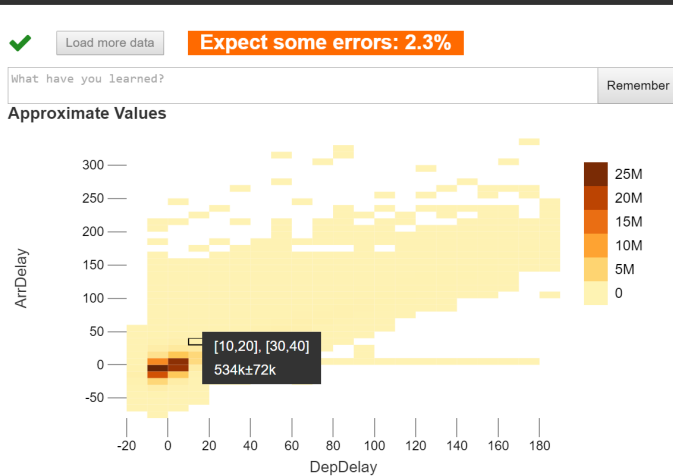
Value
Function: Count

Persistent Filters
e.g. AND(Carrier \$IN\$[ha, dl])(DepDelay>=0)

Zoom
clear Capture as Filter

(ArrDelay \$RNG\$
[[-148.80619517543857,390.49205043859655]])

(DepDelay \$RNG\$
[[-19.819658218570382,187.25649037534237]])



Massive drop off after Sep 2001

Exact data loaded (18s)

3 decades of flights

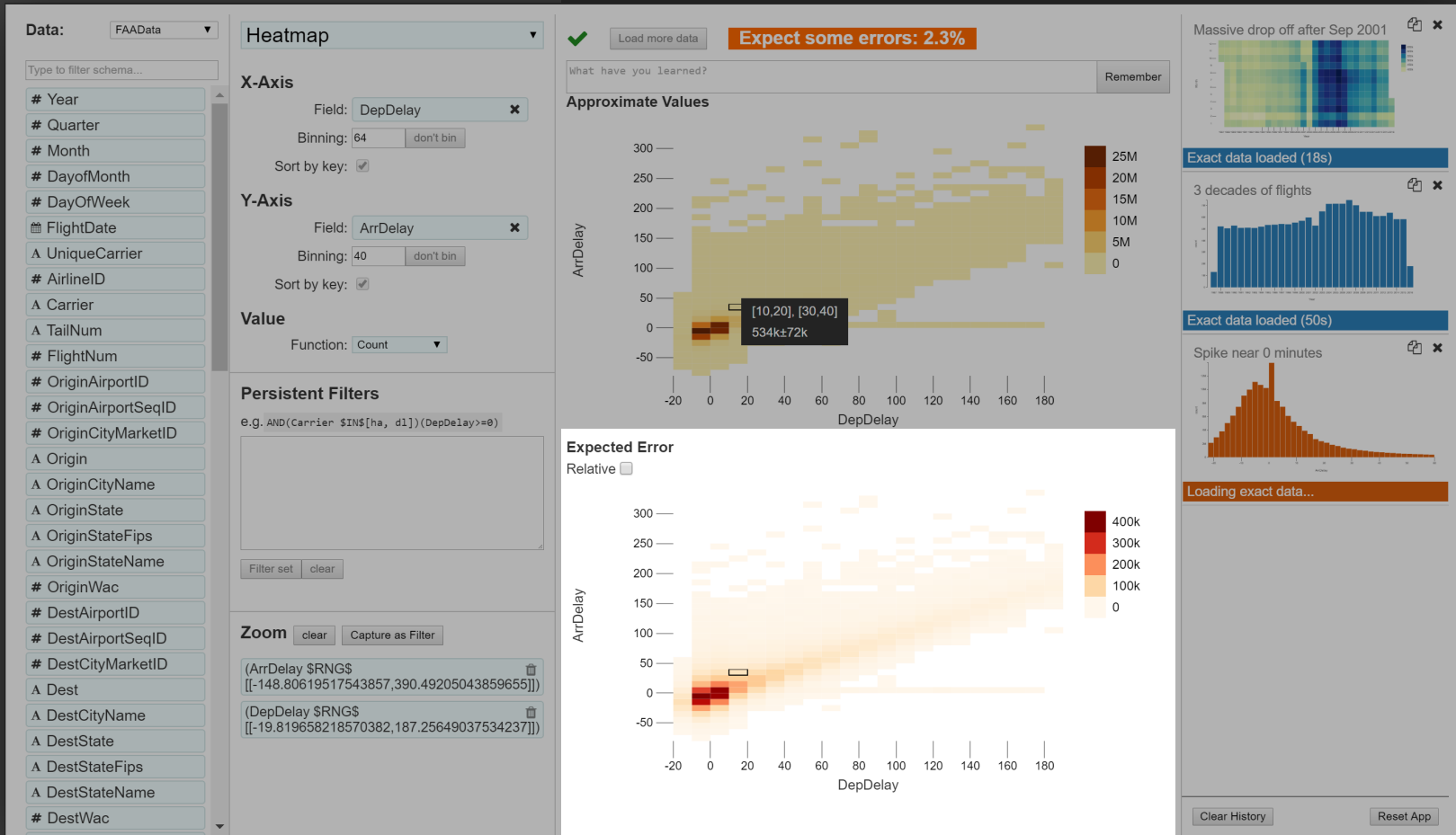
Exact data loaded (50s)

Spike near 0 minutes

Loading exact data...

Clear History Reset App

Visualize Uncertainty



Show a History of Previous Charts

Data: FAADData

Heatmap

Type to filter schema...

- # Year
- # Quarter
- # Month
- # DayOfMonth
- # DayOfWeek
- 📅 FlightDate
- A UniqueCarrier
- # AirlineID
- A Carrier
- A TailNum
- # FlightNum
- # OriginAirportID
- # OriginAirportSeqID
- # OriginCityMarketID
- A Origin
- A OriginCityName
- A OriginState
- A OriginStateFips
- A OriginStateName
- # OriginWac
- # DestAirportID
- # DestAirportSeqID
- # DestCityMarketID
- A Dest
- A DestCityName
- A DestState
- A DestStateFips
- A DestStateName
- # DestWac

X-Axis
Field: DepDelay
Binning: 64
Sort by key:

Y-Axis
Field: ArrDelay
Binning: 40
Sort by key:

Value
Function: Count

Persistent Filters
e.g. AND(Carrier \$IN\$[ha, dl])(DepDelay>=0)

Zoom
clear Capture as Filter
(ArrDelay \$RNG\$
[[-148.80619517543857,390.49205043859655]])
(DepDelay \$RNG\$
[[-19.819658218570382,187.25649037534237]])

✓ Load more data **Expect some errors: 2.3%**

what have you learned? Remember

Approximate Values

ArrDelay

DepDelay

534k±72k

25M
20M
15M
10M
5M
0

Expected Error
Relative

ArrDelay

DepDelay

400k
300k
200k
100k
0

Massive drop off after Sep 2001

Exact data loaded (18s)

3 decades of flights

Exact data loaded (50s)

Spike near 0 minutes

Loading exact data...

Clear History Reset App

Help Analysts Confirm Results

Data: FAAData

Type to filter schema...

- # Year
- # Quarter
- # Month
- # DayOfMonth
- # DayOfWeek
- 📅 FlightDate
- A UniqueCarrier
- # AirlineID
- A Carrier
- A TailNum
- # FlightNum
- # OriginAirportID
- # OriginAirportSeqID
- # OriginCityMarketID
- A Origin
- A OriginCityName
- A OriginState
- A OriginStateFips
- A OriginStateName
- # OriginWac
- # DestAirportID
- # DestAirportSeqID
- # DestCityMarketID
- A Dest
- A DestCityName
- A DestState
- A DestStateFips
- A DestStateName
- # DestWac
- A CRSDEPTime

Heatmap

X-Axis: Field: DepDelay
Binning: 64
Sort by key:

Y-Axis: Field: ArrDelay
Binning: 64
Sort by key:

Value: Function: Count

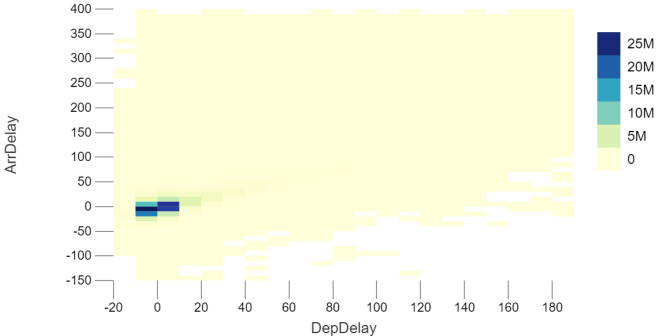
Persistent Filters
e.g. AND(CARRIER \$IN\$[ha, dl])(DEPDELAY >= 0)

Zoom
(ArrDelay \$RNG\$ [[-148.80619517543857, 390.49205043859655]])
(DepDelay \$RNG\$ [[-19.819658218570382, 187.25649037534237]])

What have you learned?

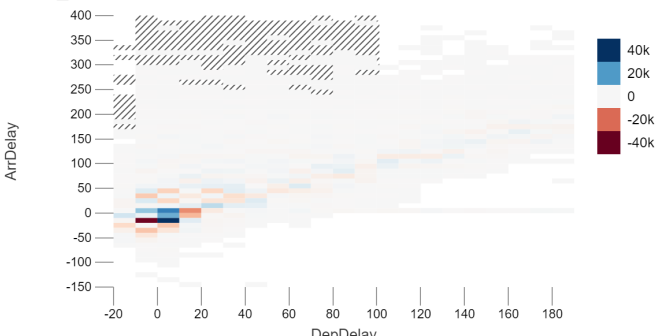
The visualization is read only because you're looking at the history. [Return to the working view](#) or make a [copy of the current chart](#).

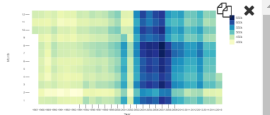
Exact Data




Difference to Approximate Data

Relative

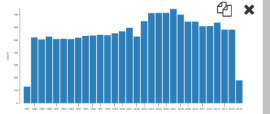




Exact data loaded (51s)

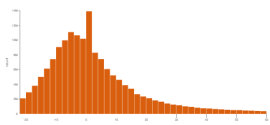


Exact data loaded (94s)



Exact data loaded (48s)

You are looking at the history and cannot make any changes.



Return to editing

Clear History Reset App

Evaluation

Case studies with teams at Microsoft who brought in their own data.

Approximation works

“seeing something right away at first glimpse is really great”

Need for guarantees

“[with a competitor] I was willing to wait 70-80 seconds. It wasn’t ideally interactive, but it meant I was looking at all the data.”

Optimism works

“I was thinking what to do next— and I saw that it had loaded, so I went back and checked it . . . [the passive update is] very nice for not interrupting your workflow.”

In Conclusion...

Two Challenges:

1. Effective **visual encoding**
2. Real-time **interaction**

Perceptual and interactive scalability should be limited by the **chosen resolution** of the visualized data, not the number of records.

Bin > Aggregate (> Smooth) > Plot

1. **Bin** Divide data domain into discrete “buckets”
2. **Aggregate** Count, Sum, Average, Min, Max, ...
3. **Smooth** *Optional*: smooth aggregates [Wickham '13]
4. **Plot** Visualize the aggregate values

Interactive Scalability Strategies

1. Query Database
2. Client-Side Indexing / Data Cubes
3. Prefetching
4. Approximation

These strategies are **not** mutually exclusive!
Systems can apply them in tandem.