

Efficient, Effective Interactive Visualizations

eugene wu

Columbia University

This presentation contains some animation.

See the animation in the powerpoint file at:

<https://cudbg.github.io/sigmod19tutorial//files/wu.pptx>

More Great Tutorials/Workshops!

Evaluating Interactive Data Systems: Workloads, Metrics, & Guidelines.

SIGMOD18. Jiang, Rahman, Nandi

Overview of Data Exploration Techniques.

SIGMOD15. Idreos, Papaemmanouil, Chaudhuri

HILDA: Human in the Loop Data Analysis

SIGMOD Workshop. <http://hilda.io>

DSIA: Data Systems for Interactive Analysis

VIS Workshop. <https://www.interactive-analysis.org>

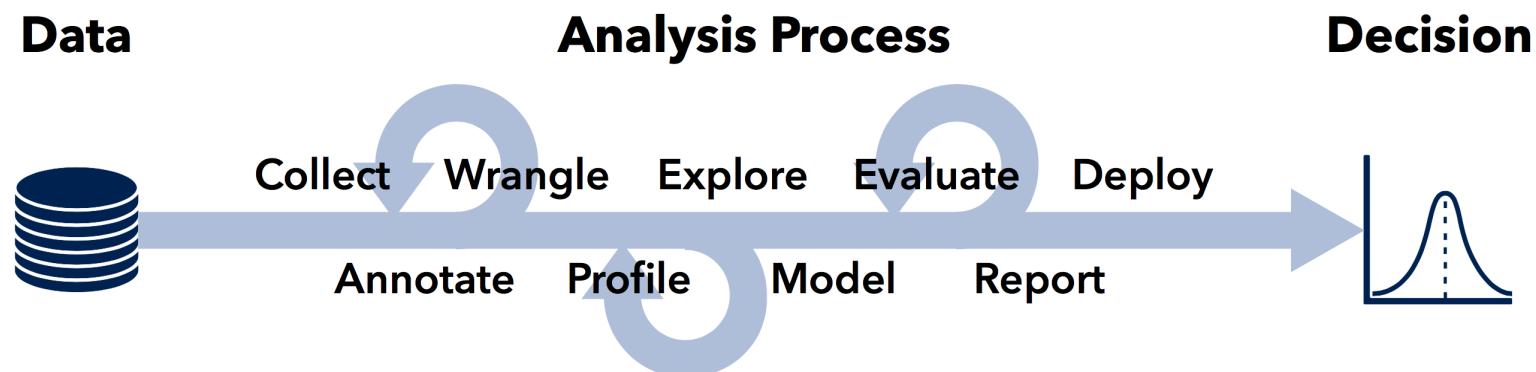
Connecting Visualization and Data Management Research

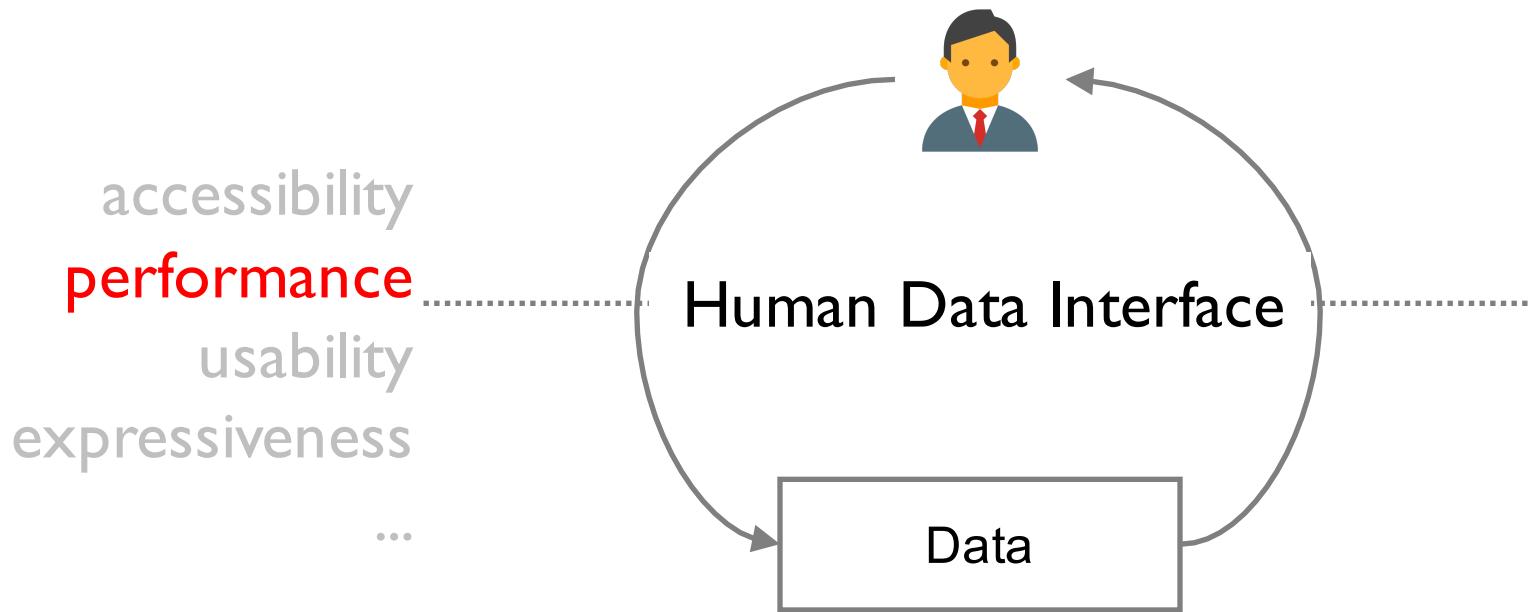
Dagstuhl Workshop. Chang, Fekete, Freire, Scheidegger.

Road Map

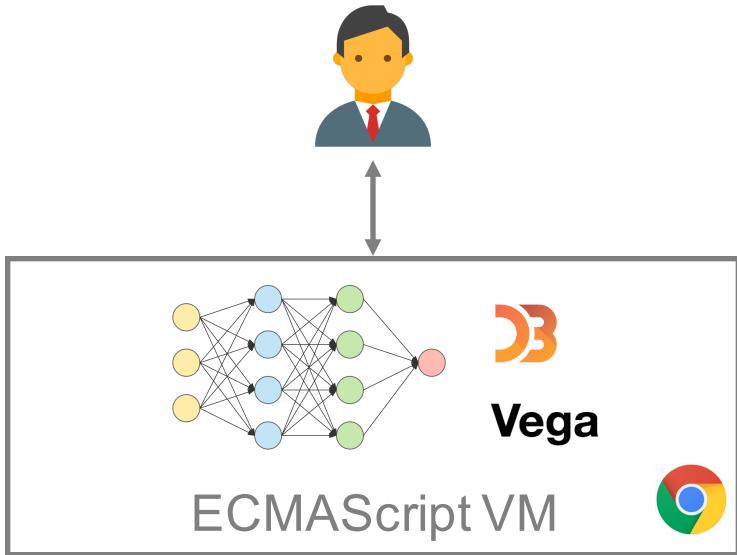
**Background
Mechanisms
A Relational Story**

Visualization is Ubiquitous





Remove computing bottleneck
Tutorial: focus on performance



SVD architecture

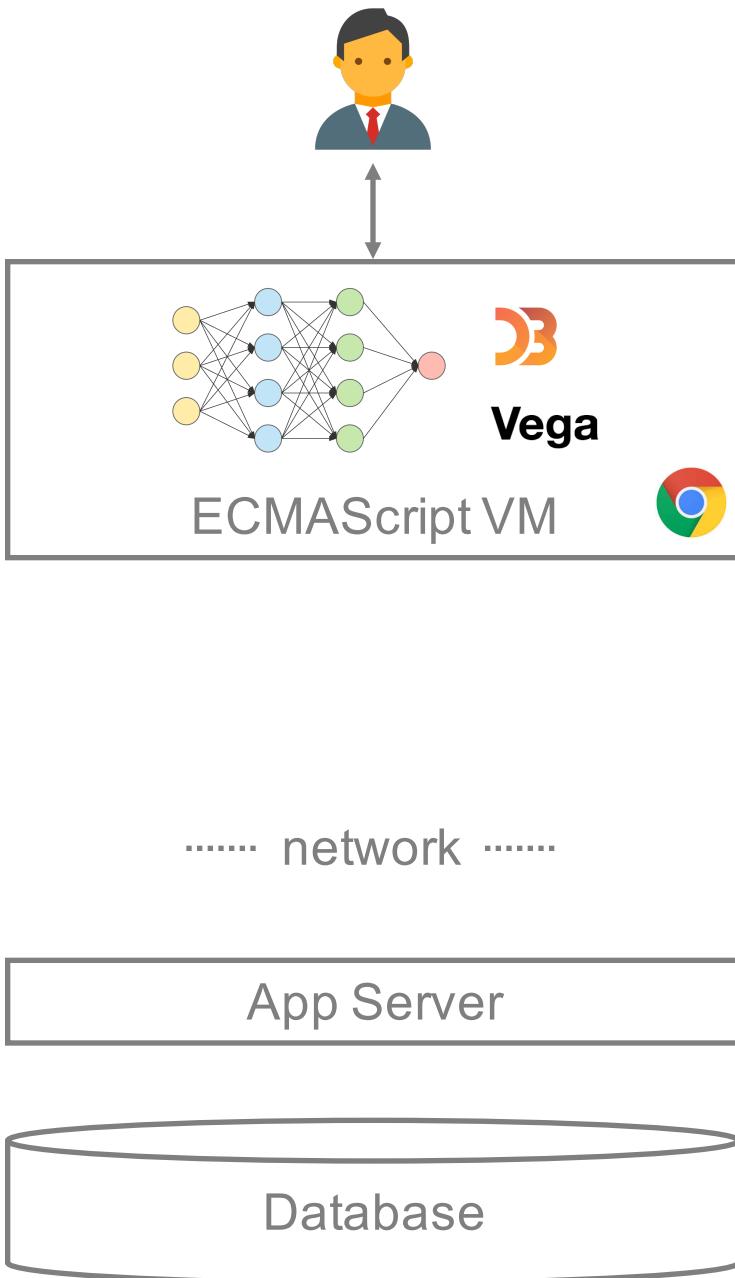
SQL,Vis,DB

Client-side vis libraries very robust

Borrows data flow, event processing, incremental updates from databases.

Outputs scene graph

Great for small datasets



SVD architecture

SQL,Vis,DB

Client-side vis libraries very robust

Borrows data flow, event processing, incremental updates from databases.

Outputs scene graph

Great for small datasets

Bigger dataset? Much messier!

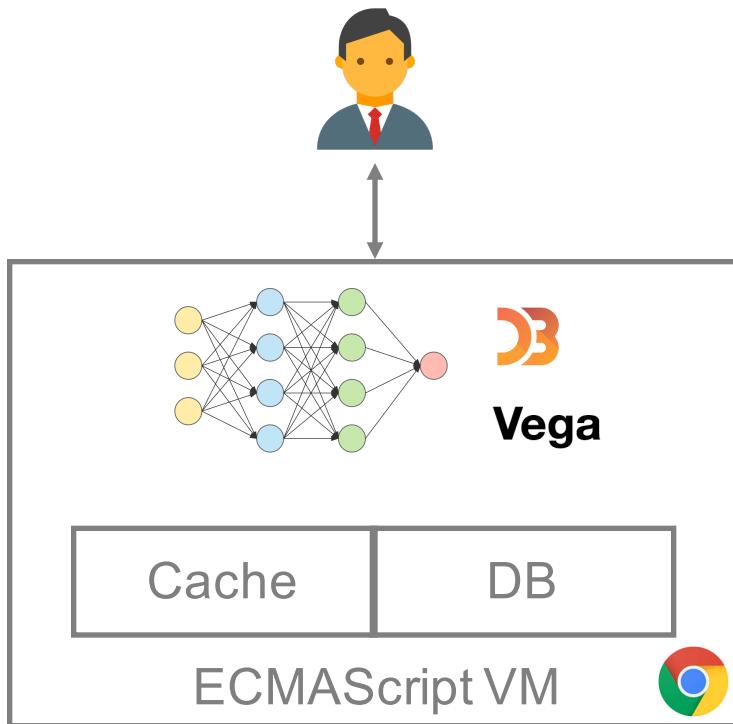
Communication overheads

DBs slower

Queries slower

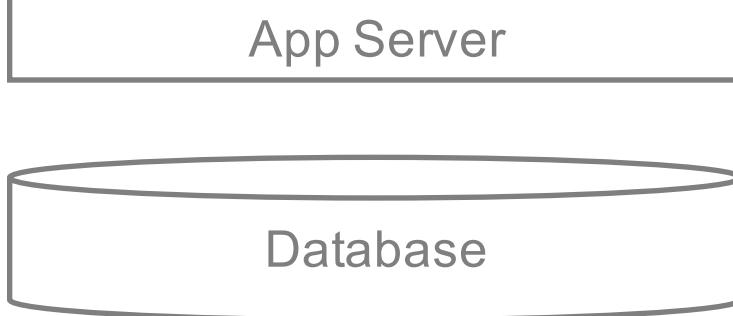
SVD architecture

SQL,Vis,DB



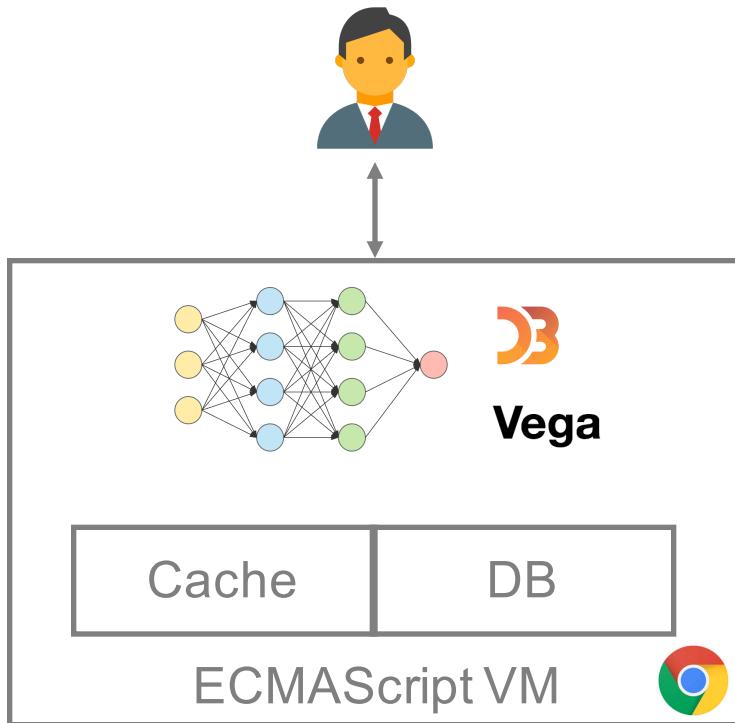
Client-side vis libraries very robust
Borrows data flow, event processing, incremental updates from databases.
Outputs scene graph
Great for small datasets

Bigger dataset? Much messier!
Communication overheads
DBs slower
Queries slower
Caches everywhere
Smart caches (DBs) everywhere



SVD architecture

SQL,Vis,DB



Latency is cumulative
Separate components/code lines
Manual implementation
Manual optimization

Akin to writing queries
before Rel Alg.

Vis Memory Hierarchy



ECMAScript

Local Machine

..... network

Big Server

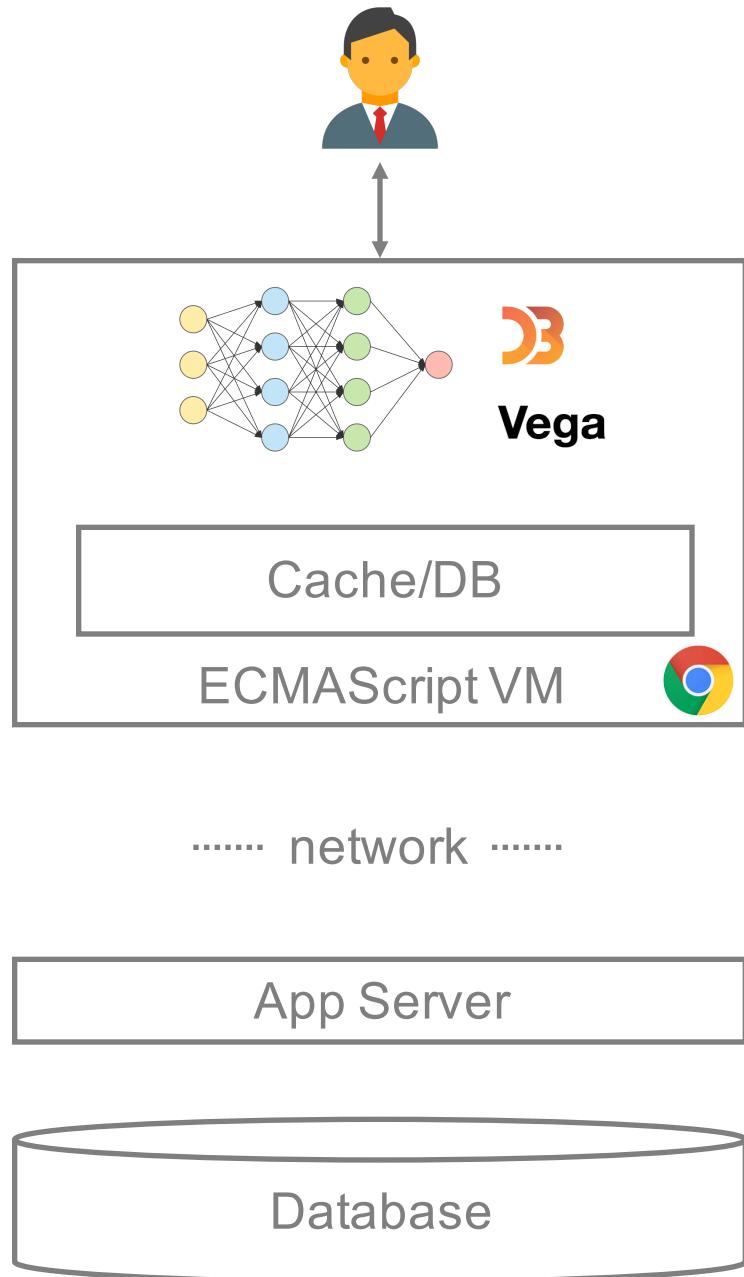
Big Database

Absolute latency expectations

Complex analysis workloads

Data dense

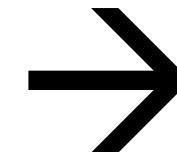
Analysis (usually) already known



The Bottom Line

Simple programming API for
interactive applications over
large-scale/in-DB data is

unsolved



Massive opportunity

Mechanisms

that leverage vis semantics

Overview of Mechanisms

Semantics	Idea
Interaction	Scale cube dimensionality to interactions
Interaction	Network Pre-fetch as prediction
Render	Aggregate results to reduce network cost Push rendering logic into query processing
Perception	Approximation Push perceptual inaccuracies into query processing
Task	Know the Task Ensure Vis API provides rich optimization hints

Overview of Mechanisms

Semantics Idea

Interaction Scale cube dimensionality to interactions

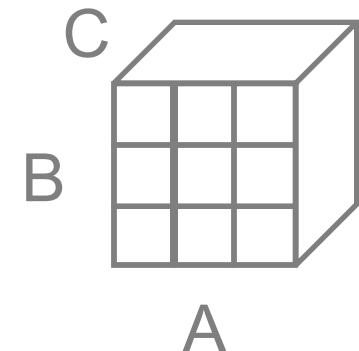
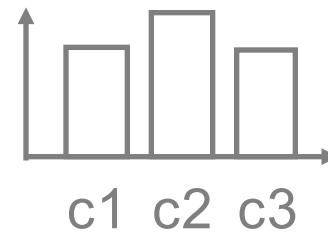
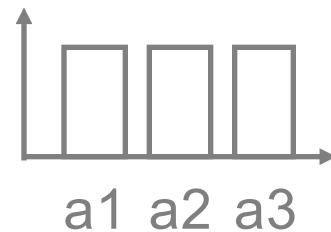
Interaction Network Pre-fetch as prediction

Render Aggregate results to reduce network cost
Push rendering logic into query processing

Perception Approximation
Push perceptual inaccuracies into query processing

Task Know the Task
Ensure Vis API provides rich optimization hints

Cubes: Immense

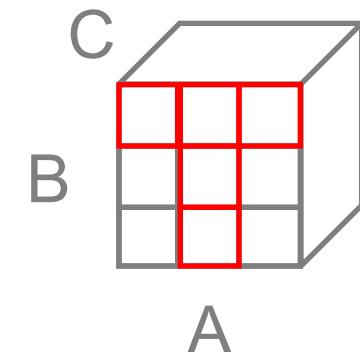
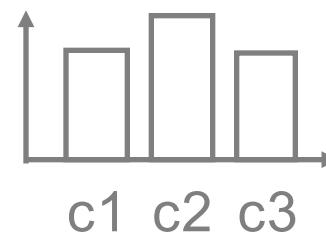
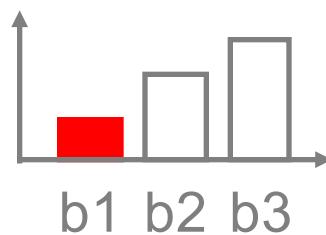
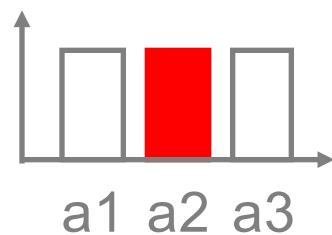


data cube

Cross filtering
interact with data in A, filter data in B

Build data cube

Cubes: Immense

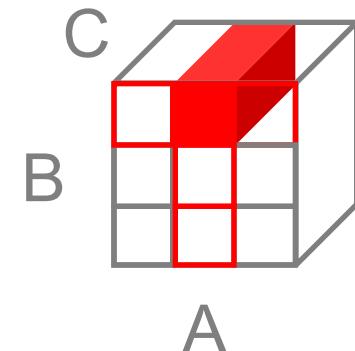
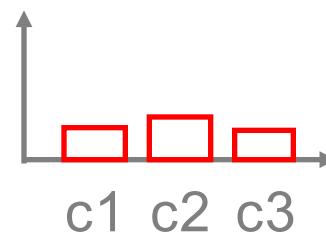
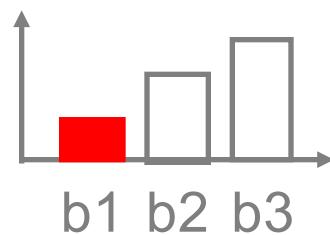
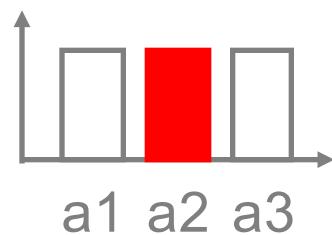


data cube

Cross filtering
interact with data in A, filter data in B

Build data cube

Cubes: Immense

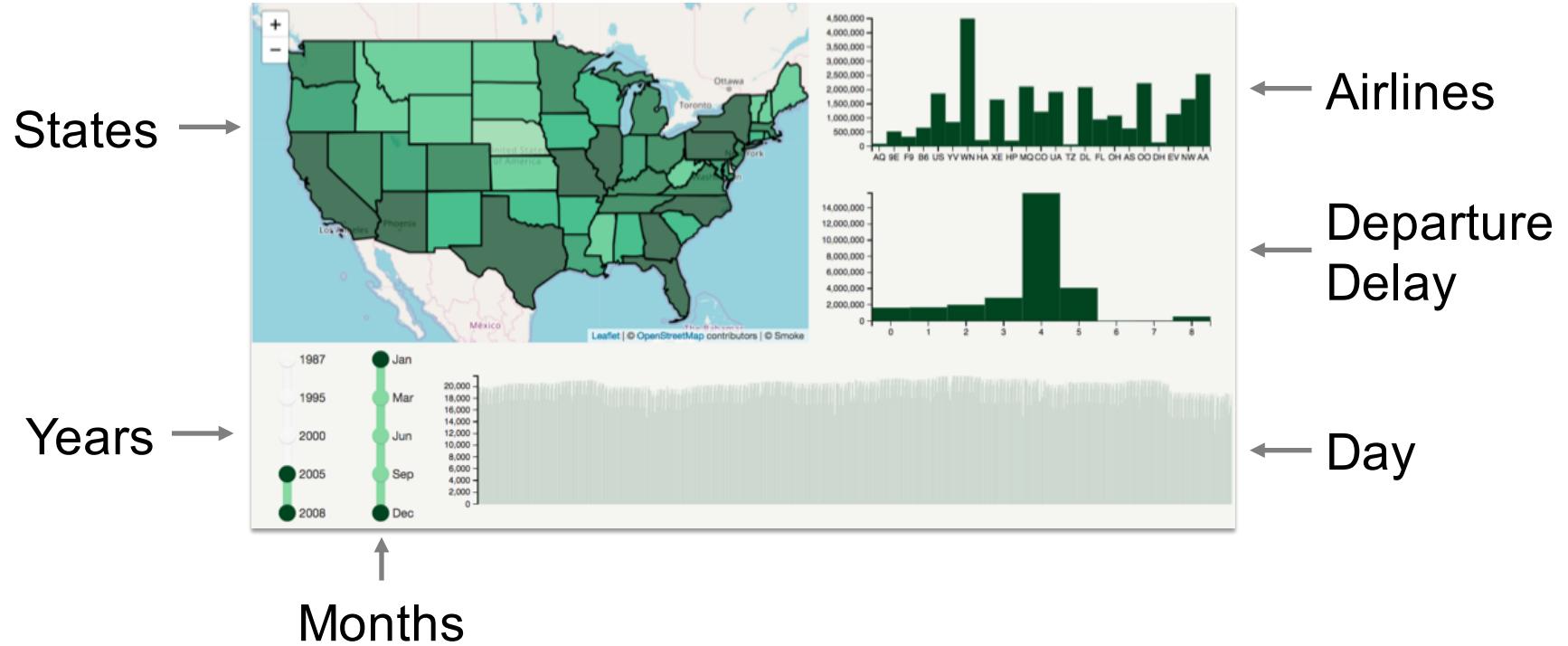


data cube

Cross filtering
interact with data in A, filter data in B

Build data cube

Cubes: Immense



6D CUBE

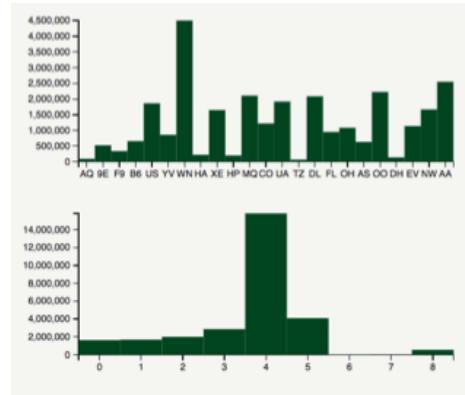
$$50 * 21 * 12 * 23 * 9 * 7000 > 18B \text{ ints} \sim 144\text{GB}$$

Idea: User doesn't express all combinations

Cubes: Immense

```
SELECT airline, COUNT
```

```
SELECT delay, COUNT
```



← Airlines

← Departure
Delay

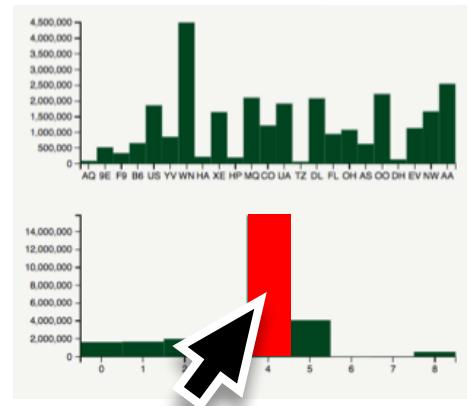
Cubes: Immens

1D

```
SELECT airline, COUNT  
      WHERE delay = 4
```

1D

```
SELECT delay, COUNT
```



← Airlines
← Departure
Delay

Interaction is 2D:
CUBE on airline, delay
Build cube for *pairs* of plots

K 1D plots, N values per attr

Naïve: N^K ~18,000,000,000

Immens: N^2K^2 809,765

Cubes

Takeaway

Scale to dimensionality of vis interactions

data is high dim

interactions can be low dim

Overview of Mechanisms

Semantics Idea

Interaction Scale cube dimensionality to interactions

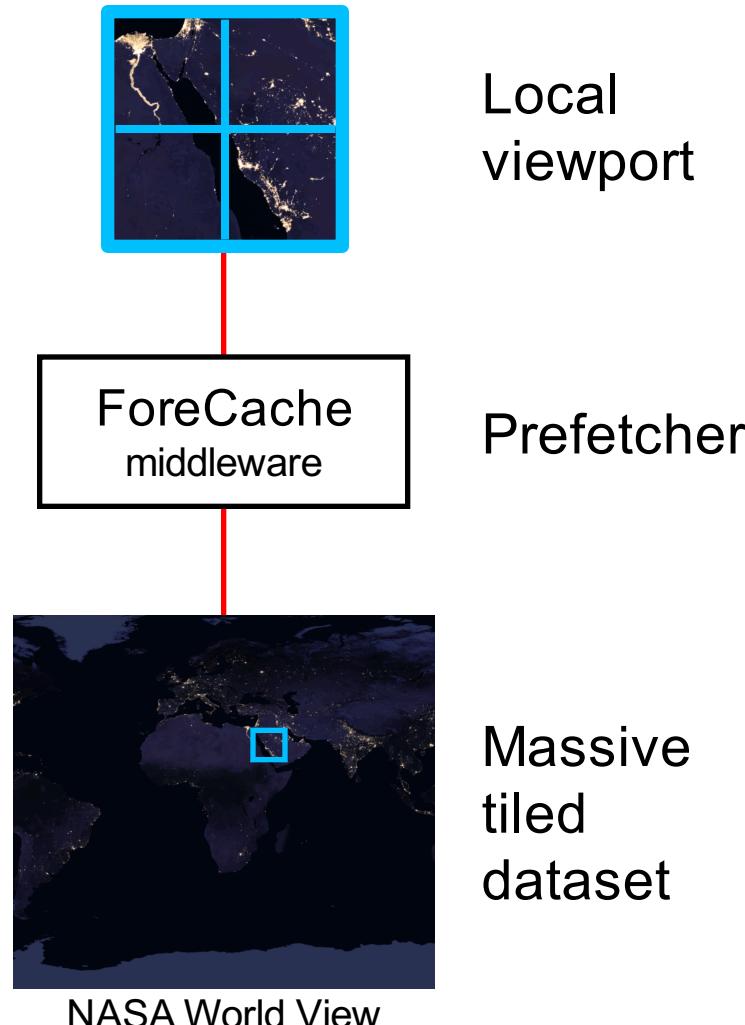
Interaction Network Pre-fetch as prediction

Render Aggregate results to reduce network cost
Push rendering logic into query processing

Perception Approximation
Push perceptual inaccuracies into query processing

Task Know the Task
Ensure Vis API provides rich optimization hints

Network: ForeCache

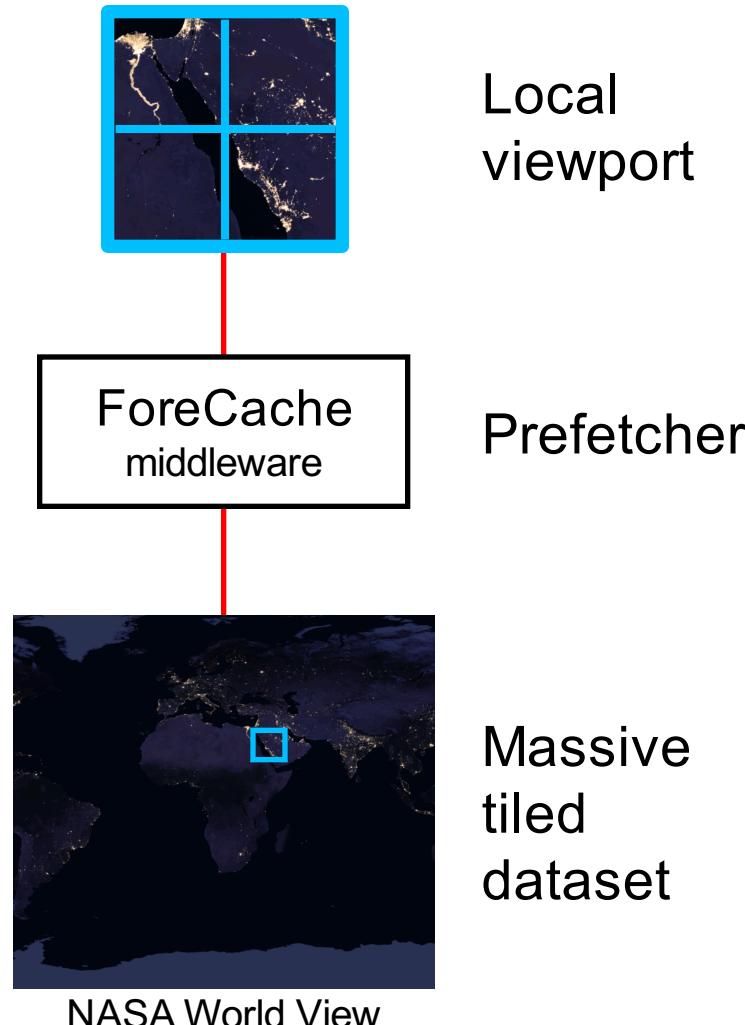


Middleware
easier to deploy
middleware has more resources

Middleware sees *inputs*
Client requests (pan/zoom)
Data content

Middleware does *actions*
make (good) predictions
pre-fetch tiles/requests

Network: ForeCache

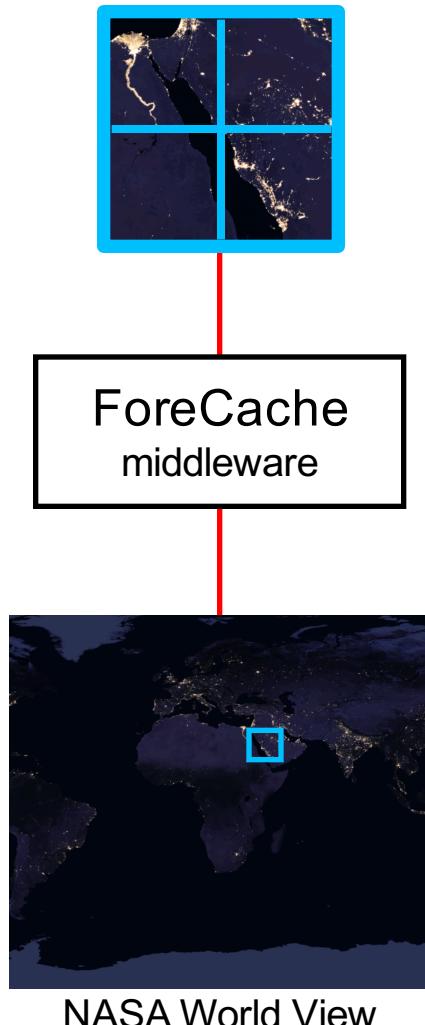


Middleware
easier to deploy
middleware has more resources

Middleware sees *inputs*
Client requests (pan/zoom)
Data content

Middleware does *actions*
make (good) predictions
pre-fetch tiles/requests

Network: ForeCache



Predict 1 of k actions
 $P(\text{action} \mid \text{state})$

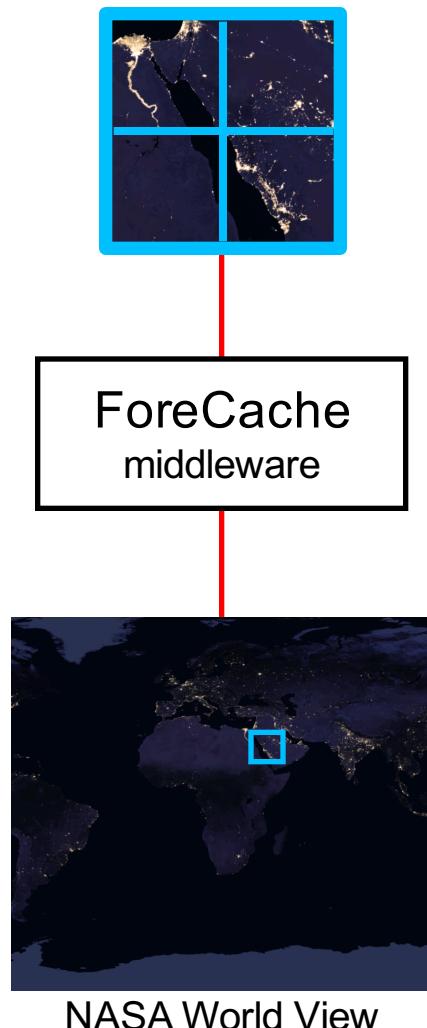
Actions:

- Pan $\leftarrow \uparrow \rightarrow \downarrow \nwarrow \nearrow \swarrow \downarrow \nwarrow$
- Zoom

Extend State using:

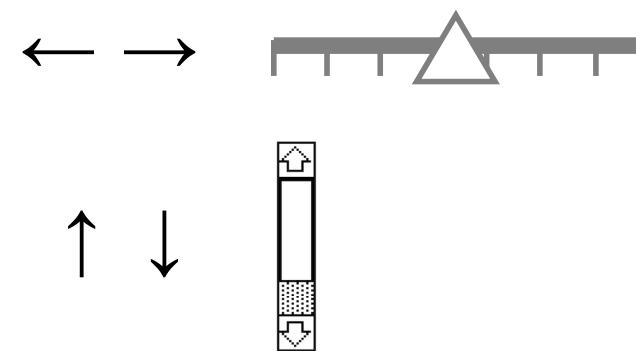
- Behavior features
 - Overview, Zoom, Details on Demand
- Map tile content features

Network: Other Examples

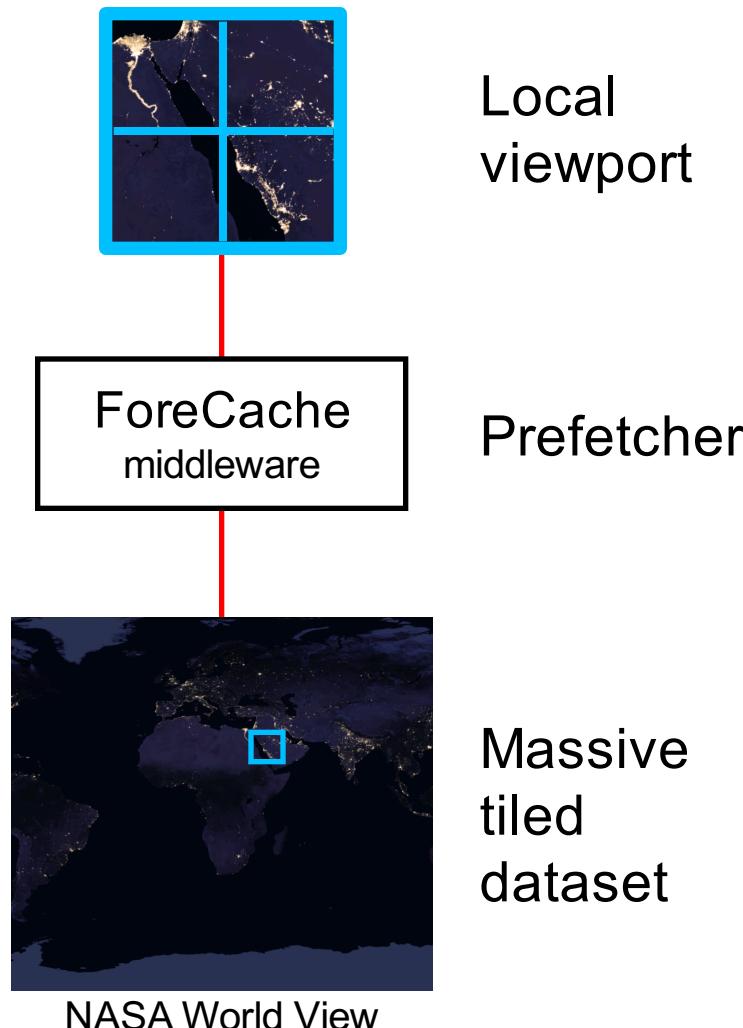


Predict 1 of k actions
 $P(\text{action} \mid \text{state})$

Restrict Actions
Semantic windows, ATLAS, SCOUT

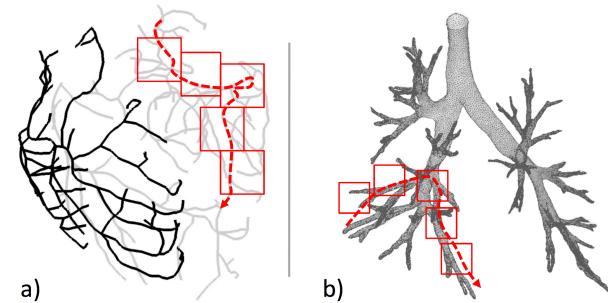


Network: Other Examples



Predict 1 of k actions
 $P(\text{action} \mid \text{state})$

Apply Domain Structure
SCOUT: graph edges instead of arbitrary 3D navigation. Prefetch from disk



DICE, Query Steering: OLAP/Markov assumptions

Network

Takeaway

Pre-fetch relies on prediction

improve predictor w/
interaction + content

Refs: Falcon (Moritz. CHI19); Momentum/Hotspot (Doshi. Thesis); ATLAS (Chan. VAST08); DICE (Kamat. ICDE14); Semantic Windows (Kalinin. SIGMOD14); Scout (Tauheed. VLDB12)

Overview of Mechanisms

Semantics Idea

Interaction Scale cube dimensionality to interactions

Interaction Network Pre-fetch as prediction

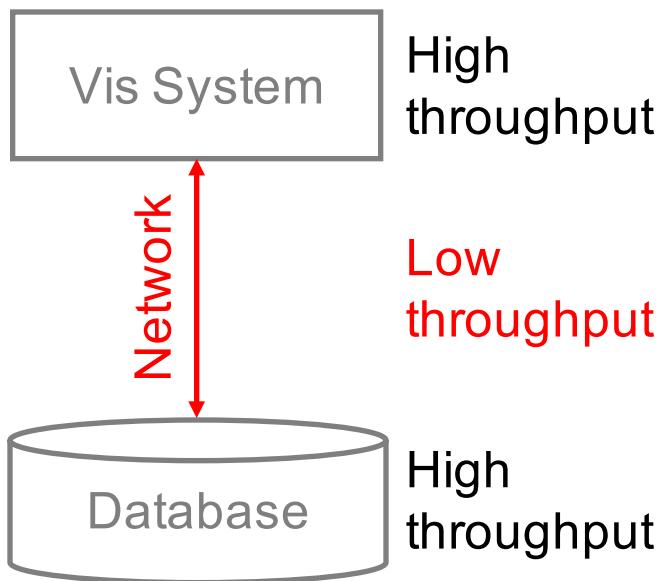
Render Aggregate results to reduce network cost
Push rendering logic into query processing

Perception Approximation
Push perceptual inaccuracies into query processing

Task Know the Task
Ensure Vis API provides rich optimization hints

Network: M4

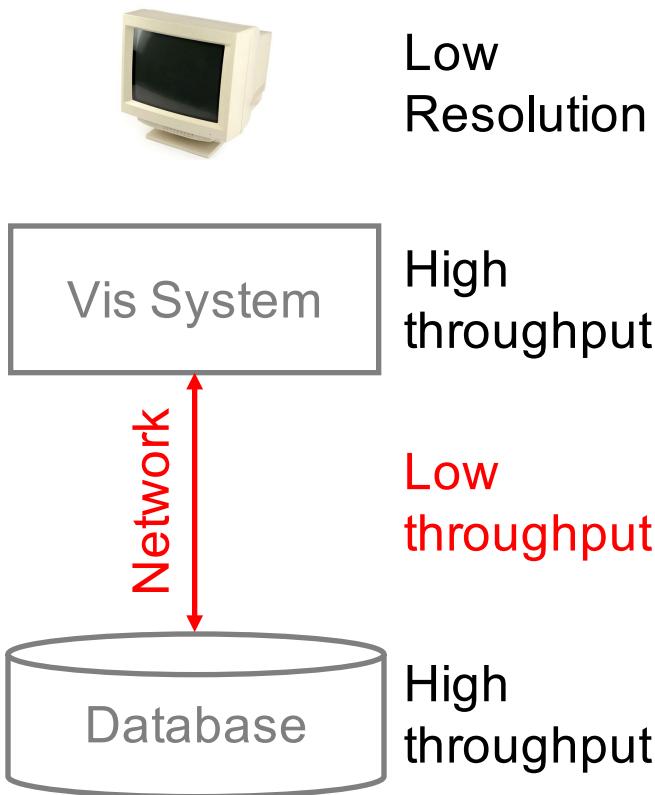
Prefetching trades bandwidth for latency



What if bandwidth limited?

Prefetch makes things worse

Network: M4



Low Resolution

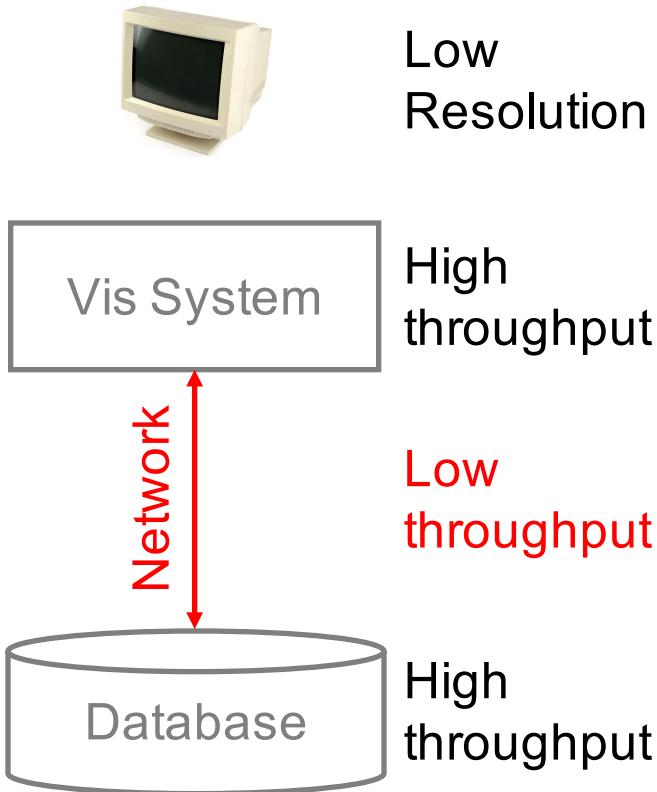
```
res = Q(database)      // on DB  
vis = render(res)     // on client
```

Network is bottleneck
Query results can be massive

Idea
vis resolution << result size

kills the network

Network: M4

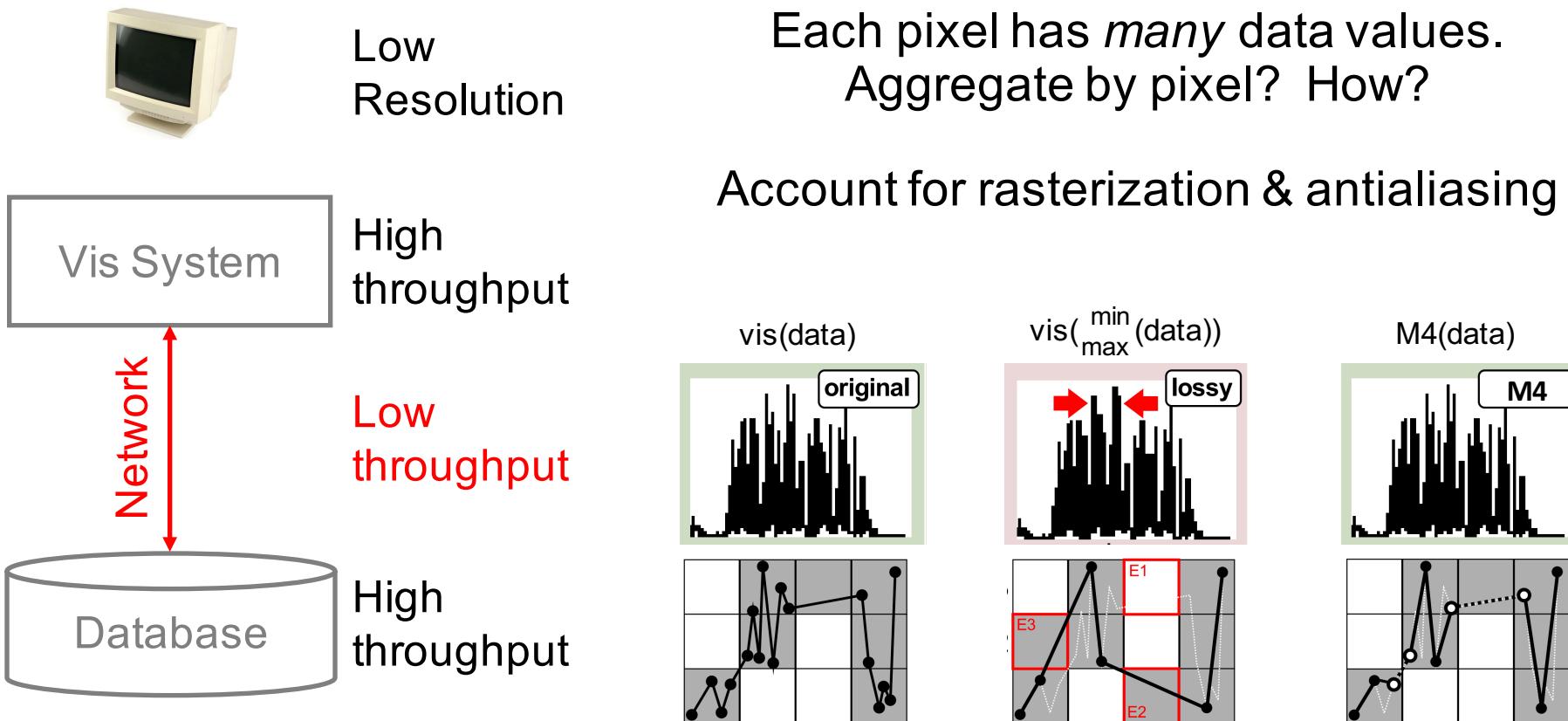


```
res' = M4_Q(database) // on DB  
vis  = render(res')    // on client
```

Network is bottleneck
Query results can be massive

Idea
vis resolution ~ M4(results)

Network: M4



Network

Takeaway

Rendering Push-down

open rendering black-box &
push into database

Overview of Mechanisms

Semantics Idea

Interaction Scale cube dimensionality to interactions

Interaction Network Pre-fetch as prediction

Render Aggregate results to reduce network cost
Push rendering logic into query processing

Perception **Progressive Vis and Approximation**
Push perceptual inaccuracies into query processing

Task Know the Task
Ensure Vis API provides rich optimization hints

Progressive Visualization

Visualization *quality* improves over time

Many forms of Progressive, such as...

Queries

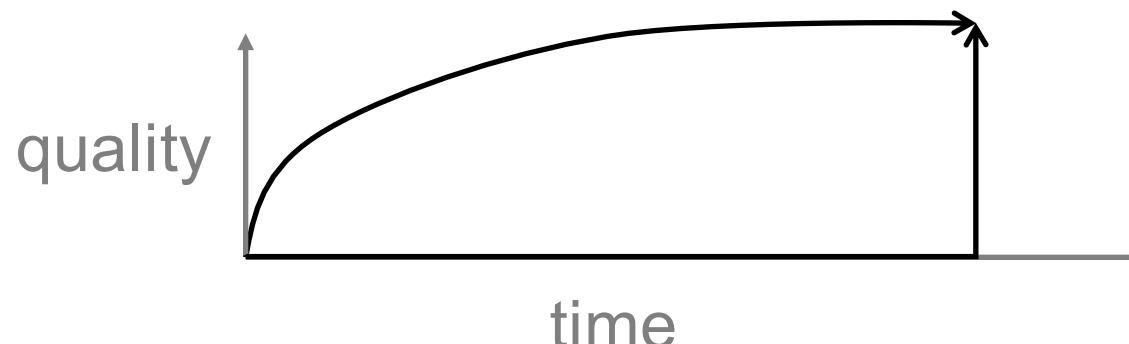
Data sampling

Rendering

Result Encoding

Interaction granularity

Interaction Design



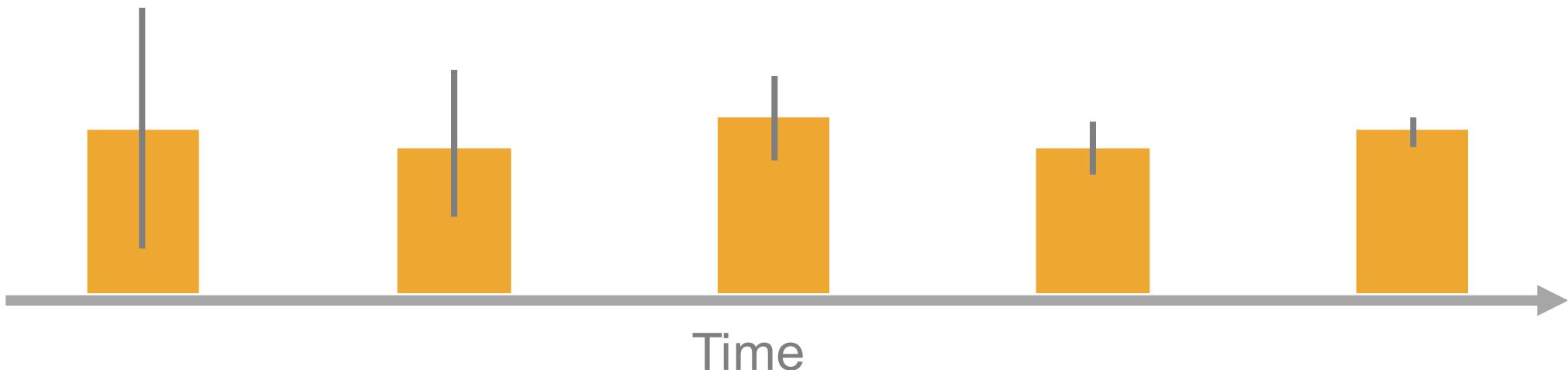
Progressive Queries

Quality ~ error bounds for $Q(D)$

Online Aggregation (Hellerstein SIGMOD97)

Wander Join (Li SIGMOD16)

Iterative sampling w/ perceptual models (P-funk Abali15)



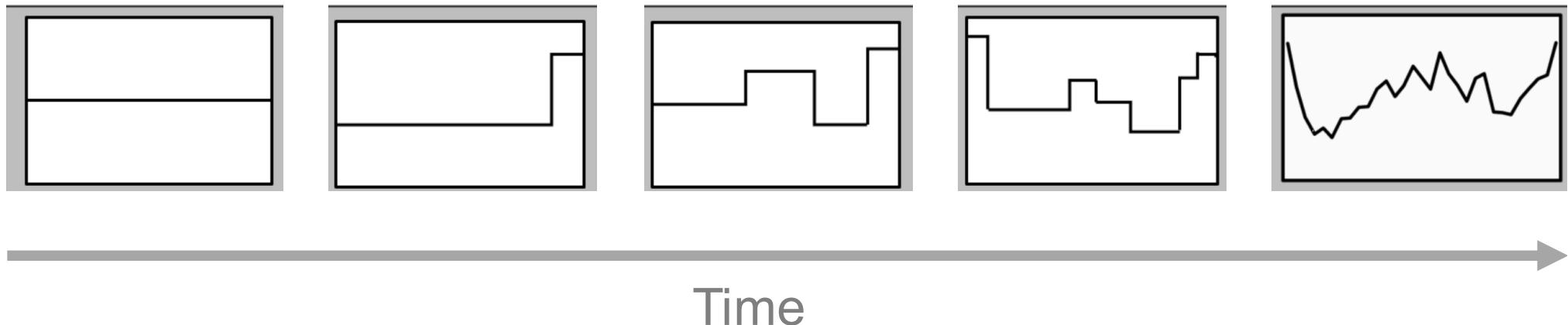
Progressive Communication

Quality ~ $E[\text{difference from full resolution}]$

Progressive image encoding (JPEG)

Speak summary, then fill in details (CiceroDB Trummer19)

Incrementally sample at higher resolutions (IncVisage Rahman17)

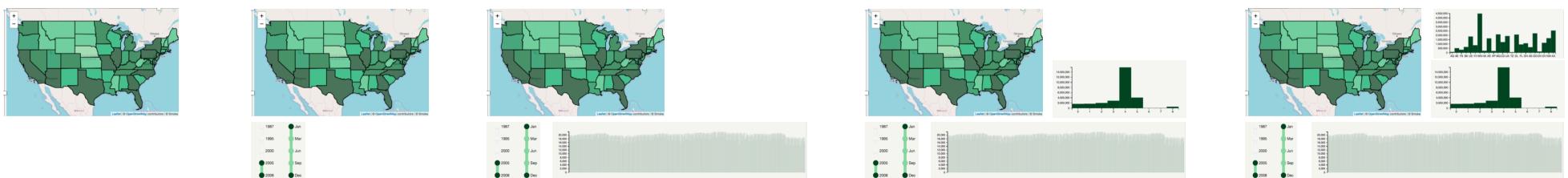


Progressive Loading

Quality ~ How much of the vis is shown

Number of charts loaded in dashboard

Webpage loading “above the fold” (Polaris netravali16)



Time

Progressive Interactions

Quality ~ Interaction granularity

Request throttling

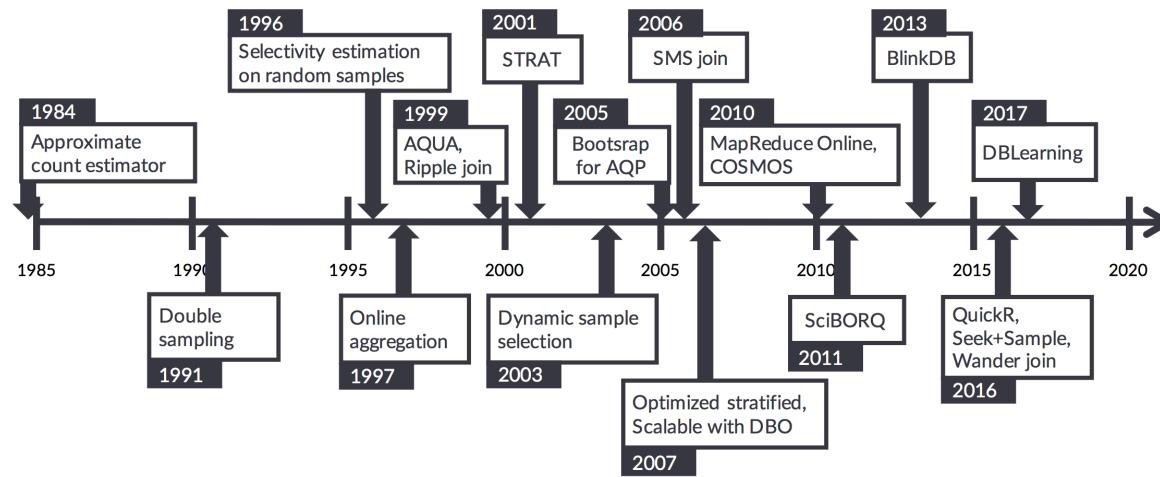
Higher brushing resolution over time (Falcon Moritz19)

Number of ticks in slider increase over time



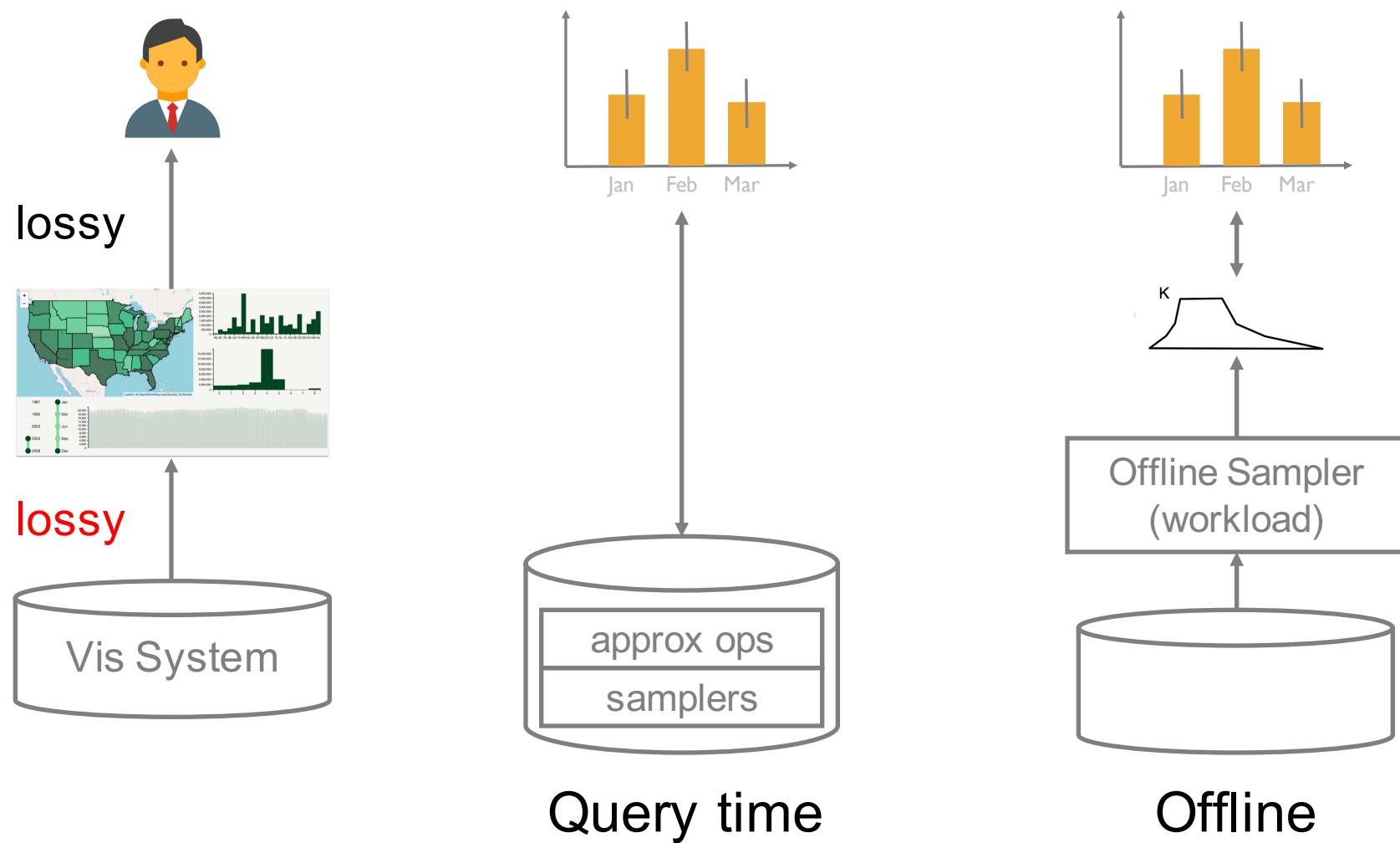
Approximate Query Processing

Progressive query processing and AQP are long-standing problems in databases

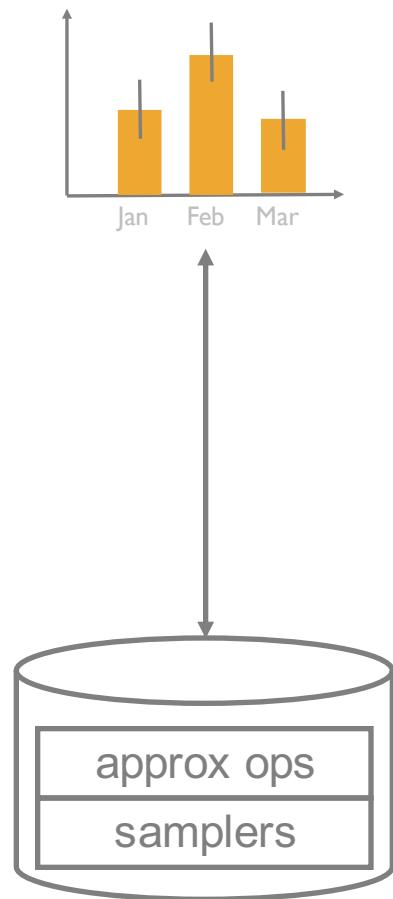


Adapting AQP to visualization also depends on visualization semantics

Approximate Query Processing



Query Time Approximation



Offline:
Do nothing

When running Q:
Choose sample operators
Draw samples to answer Q

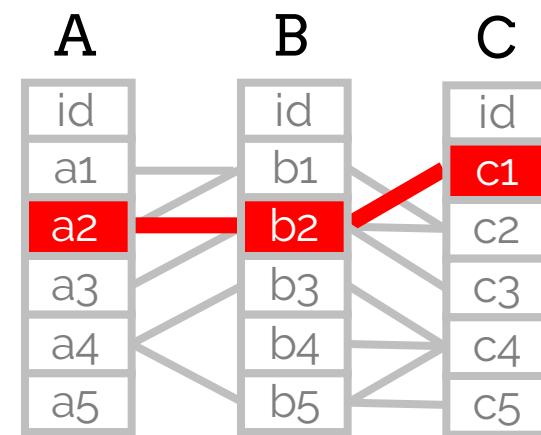
Query Time: WanderJoin

$A \bowtie B \bowtie C$

WanderJoin: leverage join indexes

edges represent join matches

Sample from A. Then path from A to C (a2-b2-c1) is a join sample
random walk is non-uniform independent sample

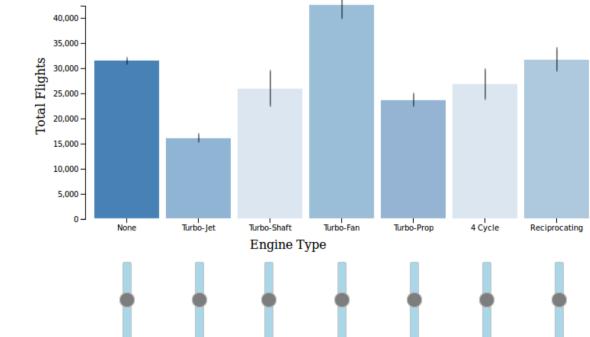
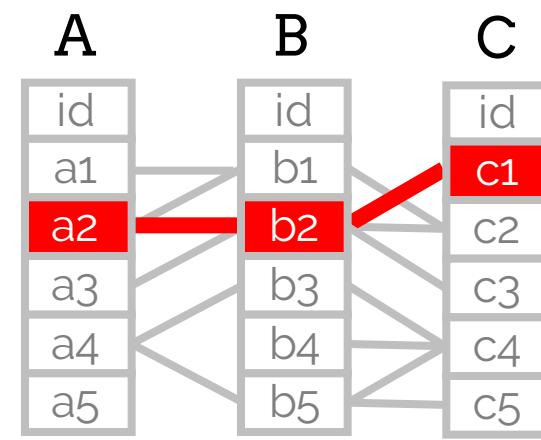


Query Time: WanderJoin

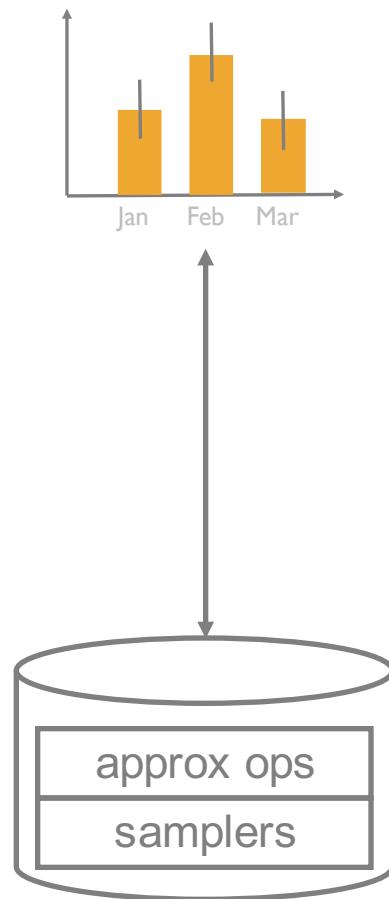
$A \bowtie B \bowtie C$

Vis is dominated by filtering and group-bys (filter by group)
Adapt WJ by biasing random walk via importance sampling for..

- Filters and dynamic selections
- User preferences



Query Time Approximation



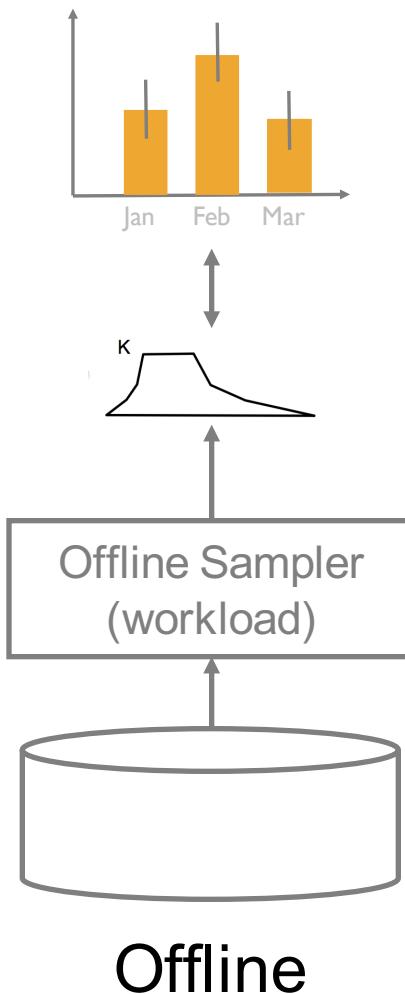
Offline:
Do nothing

When running Q:
Choose sample operators
Draw samples to answer Q

Sampling is expensive
WanderJoin uses join indexes.
Could use indexing time to build other data structs?
Can take long time for bounds to be small

Query time

Offline



Offline:

Precompute samples given workload W
Typically stratify on columns groups in W

When running Q:

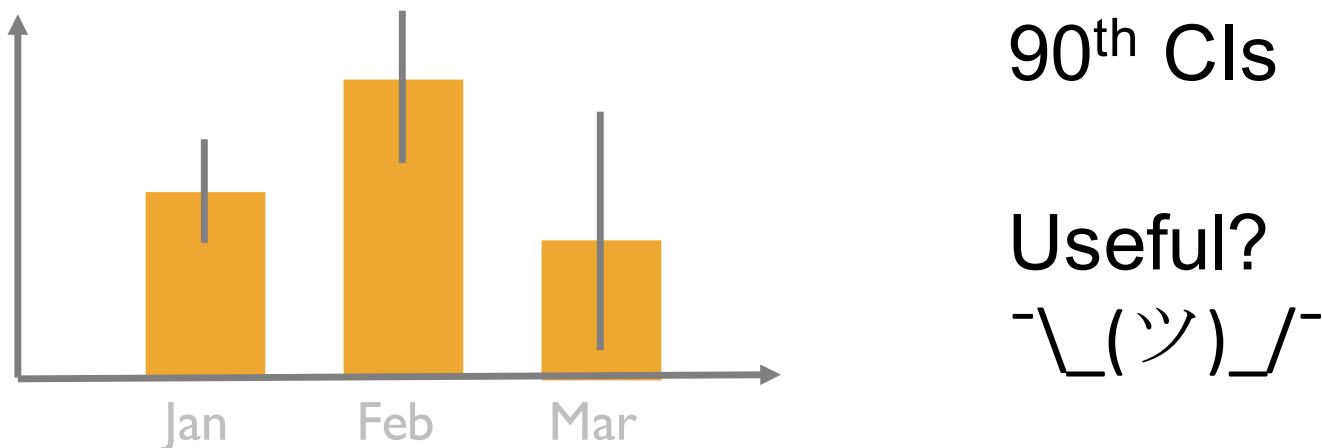
Pick precomputed samples
Use CLT/Hoeffding/bootstrap for err bounds

Offline: Sample+Seek

Adapts AQP towards visualization needs in 2 ways.

Challenges with confidence intervals (CIs)

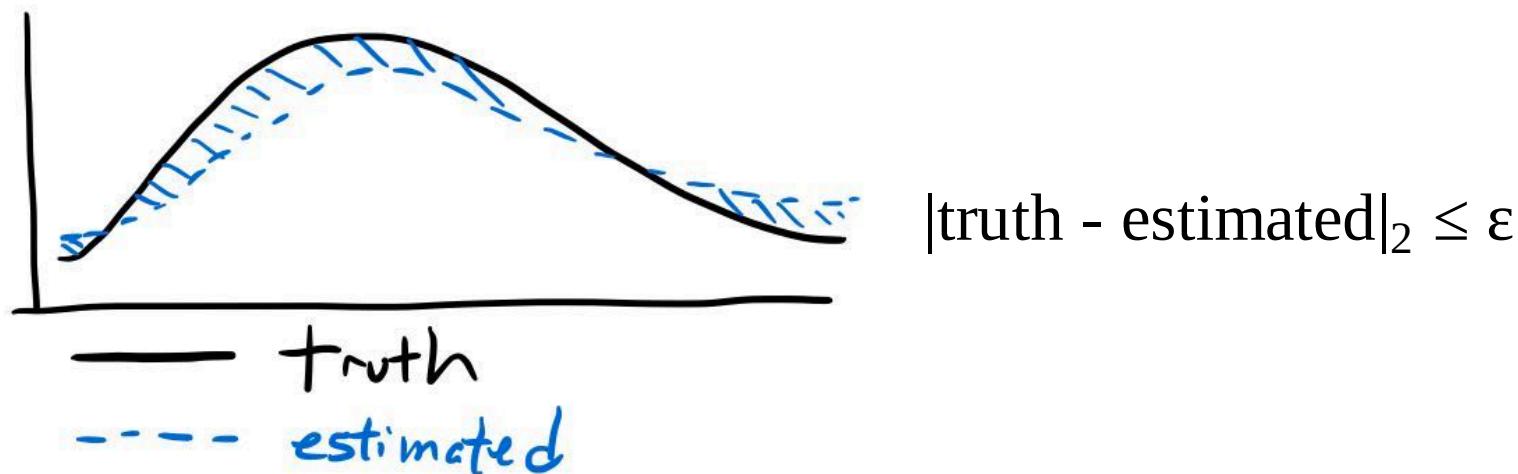
- $CI \sim \underbrace{\text{std}(n \text{ samples})}_{\text{data dependent}} / \sqrt{n}$ → sensitive to outliers
- CIs unintuitive



Offline: Sample+Seek

Proposes *Distributional Guarantee*

- Result modeled as normalized distribution
- Offline pre-computation will bound L_2 distance $\leq \varepsilon$
- Closer to understandable semantics



Offline: Sample+Seek

Don't stratify by column groups

- Column groups may be different later on

Compute **measure-biased samples** for aggregated attrs

- Need to know aggregated measures up front

```
SELECT a1, a2, ..., SUM(v1), COUNT(v2)
      FROM ...
GROUP BY a1, a2, ...
```

Offline: Sample+Seek

Measure-biased samples for $\text{SUM}(\text{val})$: proportional to val
100 rows

val	p_{naive}	$p_{\text{s+s}}$
1	$1/100$	$1/199$
\dots	\dots	\dots
1	$1/100$	$1/199$
100	$1/100$	$100/199$

If we sample 1
 val (**100**), what
is our estimate?

$$\begin{aligned} & \frac{100}{(1/100)} \\ &= 10,000 \end{aligned} \qquad \begin{aligned} & \frac{100}{(100/199)} \\ &= 199 \end{aligned}$$

Addresses err bound's *Data Dependency*

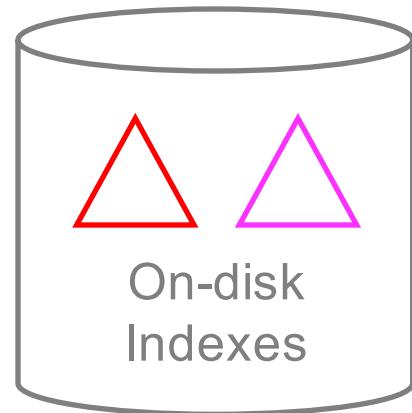
Offline: Sample+Seek

Sample probability based on value

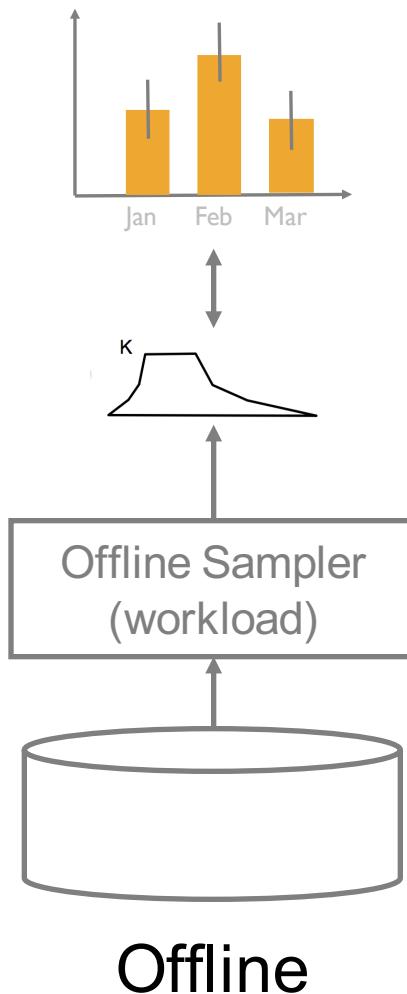
```
Q = SELECT a, SUM(b)  
      WHERE c=1 ...  
      GROUPBY b
```

```
result = Q(in-memory sample)  
if enough samples in mem:  
    return result  
if very low selectivity:  
    lookup rows directly  
else:  
    use measure-augmented index to  
    draw sample biased by b
```

In-mem
Samples



Offline



Offline:

Precompute samples given workload W
Stratify on columns groups in W

When running Q:

Pick precomputed samples
Use CLT/Hoeffding/bootstrap for err bounds

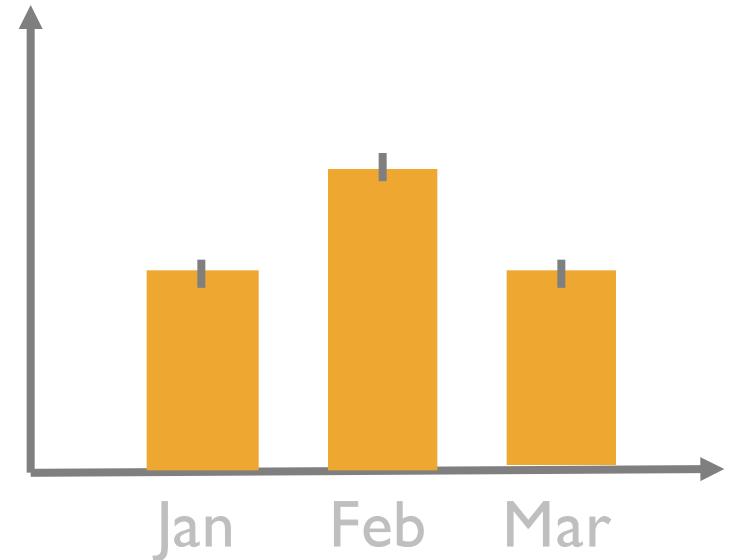
Hard to guarantee bounds are small if Q uses unseen col group

Measure-biased sampling helps, works if *aggregation function* is over sampled attributes!

What is Quality?

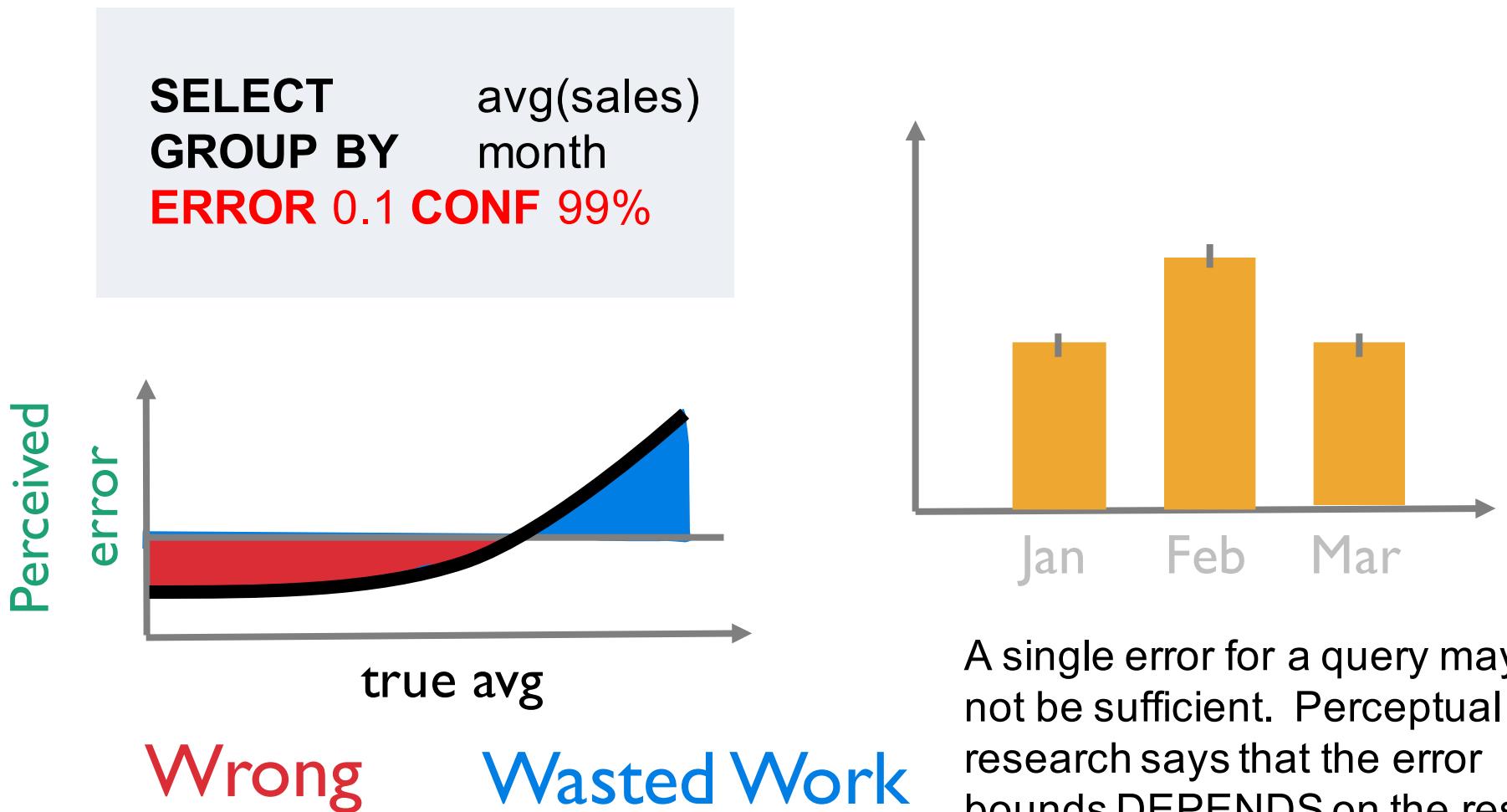
```
SELECT      avg(sales)
GROUP BY    month
ERROR 0.1 CONF 99%
```

Confidence interval per record
Error 0.1 Conf 99%



Stepping back, a bigger question is what quality should mean! Tricky even for classic error specifications in AQP

What is Quality?



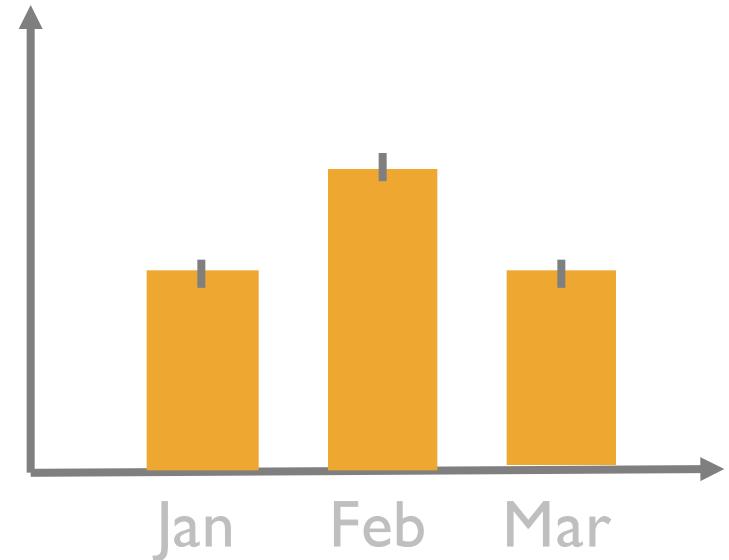
What is *Quality*?

```
SELECT      avg(sales)  
GROUP BY    month  
ERROR 0.1 CONF 99%
```

Confidence interval per record
Error 0.1 Conf 99%

Pairwise statistical test
Pairwise CI don't overlap too much

Distributional guarantee
 $E[\text{distance from true distribution}]$



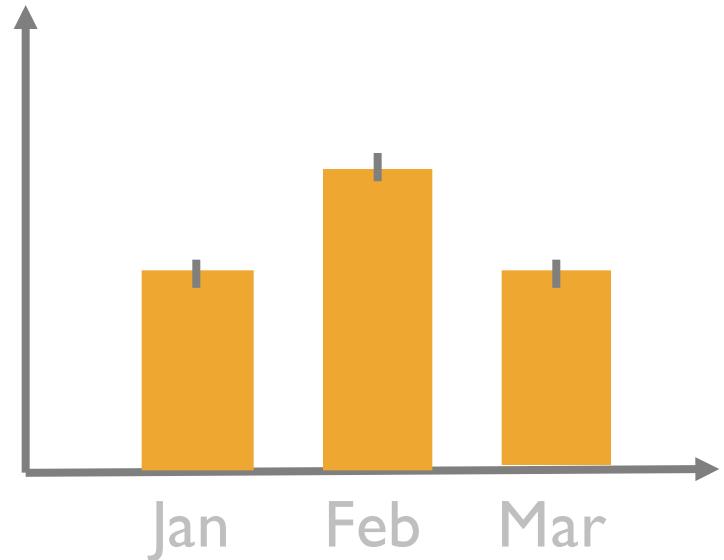
CI isn't even the only notion of quality to begin with!

What is *Quality*?

```
SELECT      avg(sales)  
GROUP BY    month  
ERROR 0.1 CONF 99%
```

Q: Who decides Quality?

A: Perception Science



The Human Side (a sample of works)

At a Glance: Approximate Entropy as a Measure of Line Chart Visualization Complexity

Ryan et al. InfoVIS19

The Human User in Progressive Visual Analytics

Micallef et al. EuroVIS (Dagstuhl report)

What Users Don't Expect about Exploratory Data Analysis on Approximate Query Processing Systems

Moritz et al. HILDA17

Why Evaluating Uncertainty Visualization is Error Prone

Jessica Hullman BELIV16

Approx + Progressive

Takeaway

Perception Push-down

Model perceptual inaccuracy &
push into database

Overview of Mechanisms

Semantics Idea

Interaction Scale cube dimensionality to interactions

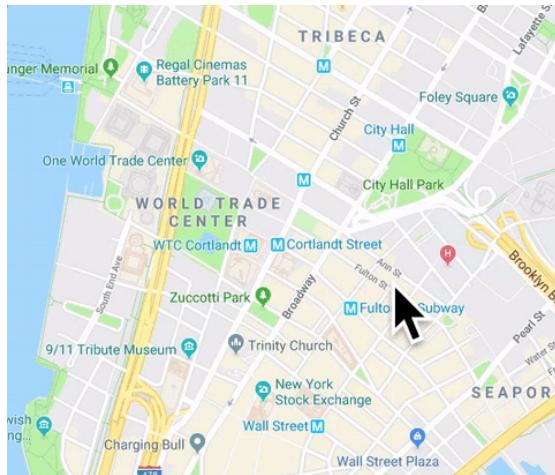
Interaction Network Pre-fetch as prediction

Render Aggregate results to reduce network cost
Push rendering logic into query processing

Perception Approximation
Push perceptual inaccuracies into query processing

Task Know the Task
Ensure Vis API provides rich optimization hints

Tasks: Kyrix



Kyrix: visualization as map

Layers render rows (map tiles, pins)

User sees through viewport

Interaction = change bounding boxes

→ Pan viewport to see more data

→ Zoom/click to switch layers

Tasks: Kyrix



Kyrix: visualization as map

Layers render rows (map tiles, pins)

User sees through viewport

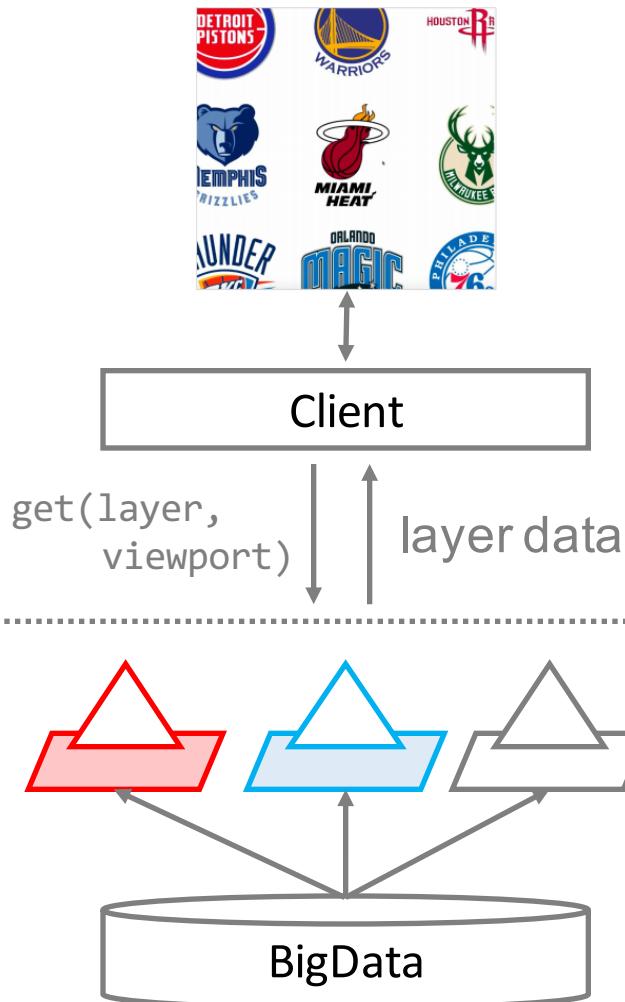
Interaction = change bounding boxes

→ Pan viewport to see more data

→ Zoom/click to switch layers



Tasks: Kyrix



simplified pseudocode

```
logos = vis.canvas(w, h)  
// xform1 computes logo positions  
// be rendered in layer  
logos.newlayer(xform1, render1)
```

Easy Developer API



Tasks

Takeaway

Know the Task

Task-based Programming API

Leverage richer semantics



An End-to-end Relational Story

Overview of Mechanisms

Semantics	Idea
Interaction	Scale cube dimensionality to interactions
Interaction	Network Pre-fetch as prediction
Render	Aggregate results to reduce network cost Push rendering logic into query processing
Perception	Approximation Push perceptual inaccuracies into query processing
Task	Know the Task Ensure Vis API provides rich optimization hints

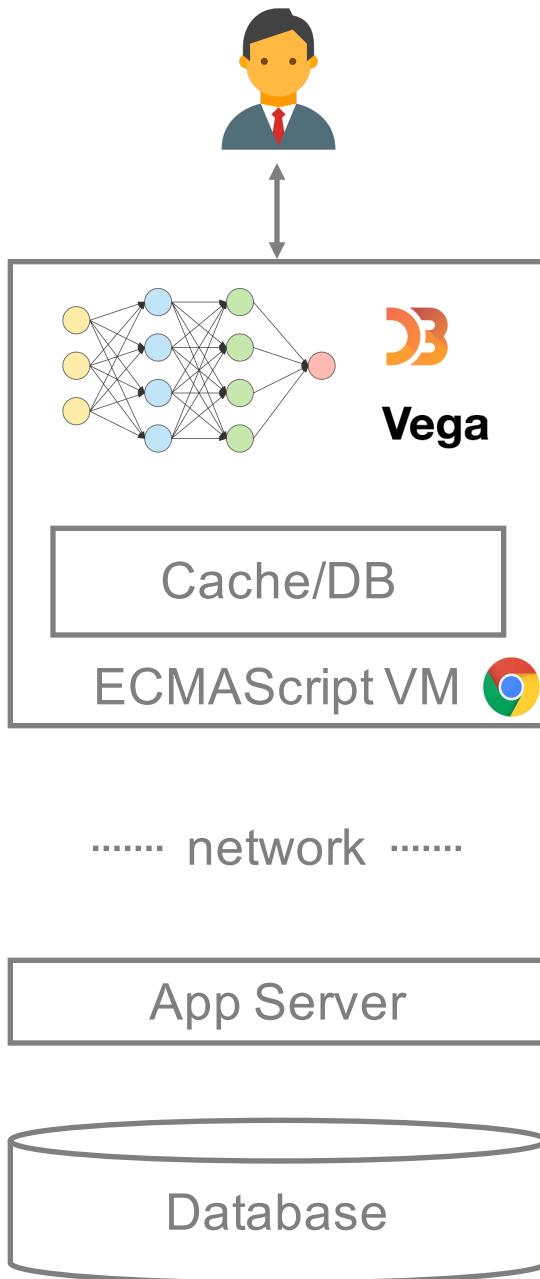
Many Disparate Optimizations

How to choose?

- Developer tells the system **Hard for dev. API?**
- Special case the system **Limits Flexibility**

Do they compose?

Need to model **application semantics** **How?**



Data Visualization Management System

Want to express *visualization, interaction, tasks, perception* all together

Vis and interaction as queries

Apply relational ideas end-to-end

- to interactions
- to consistency
- to design



A Big “Query”

Data Visualization Management System

Want to express visualization, interaction, tasks, perception all together

Vis and interaction as queries

Apply relational ideas end-to-end

- to interactions
- to consistency
- to design

A Big “Query”



Roadmap

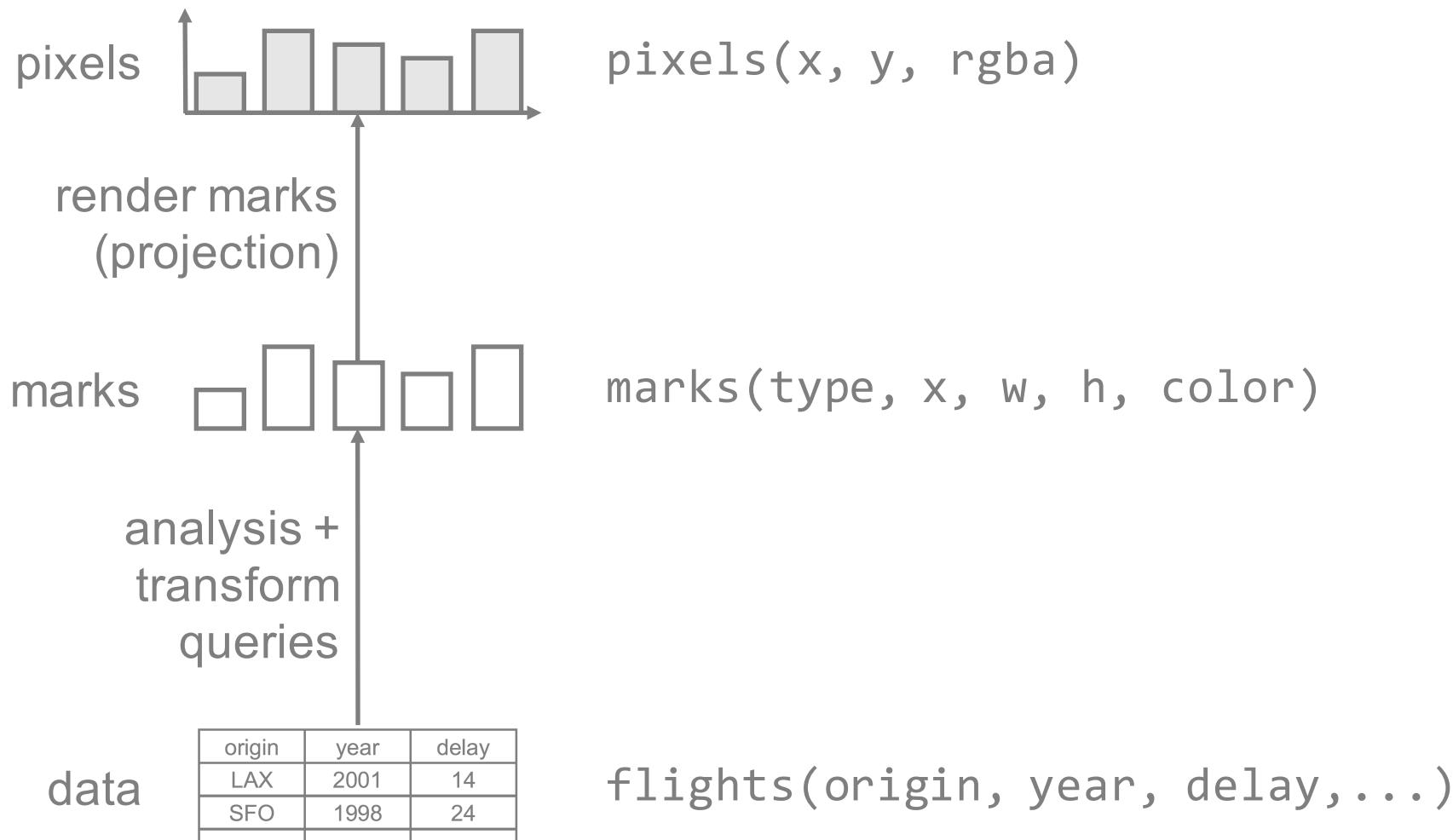
Overview of Relational Perspective

Optimizations within relational framework

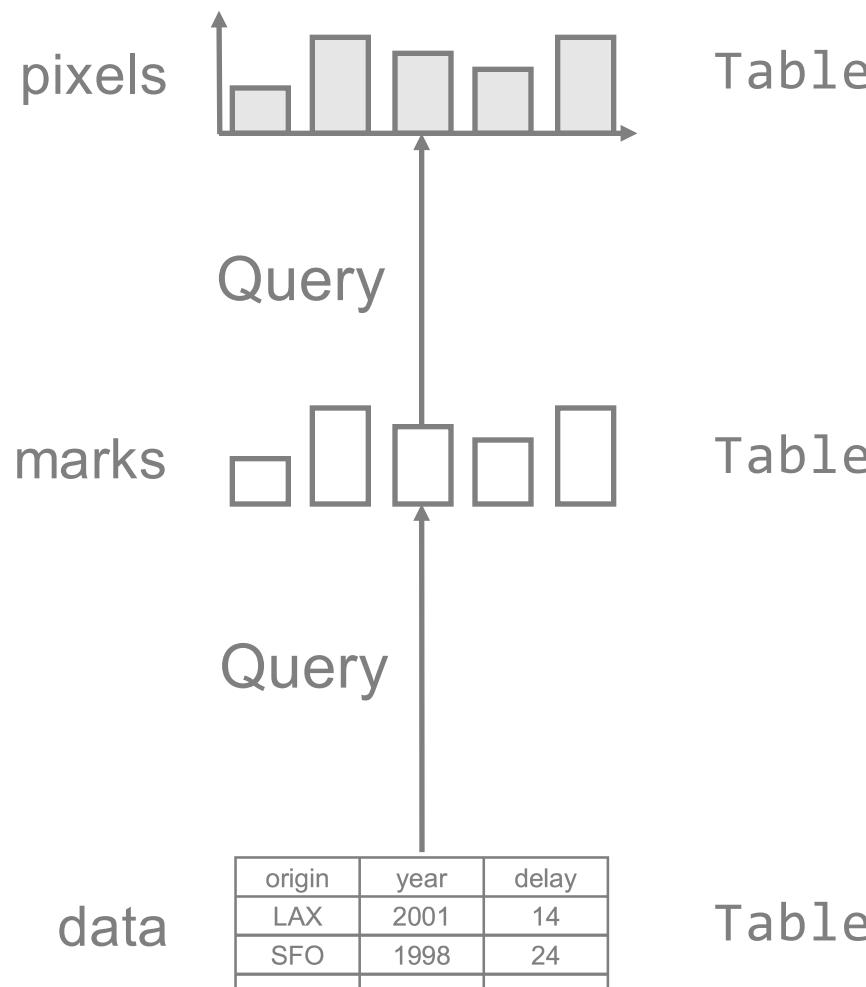
Interactions as logical expressions

- Single-view interaction
- Multi-view interactions

A Relational Perspective



A Relational Perspective

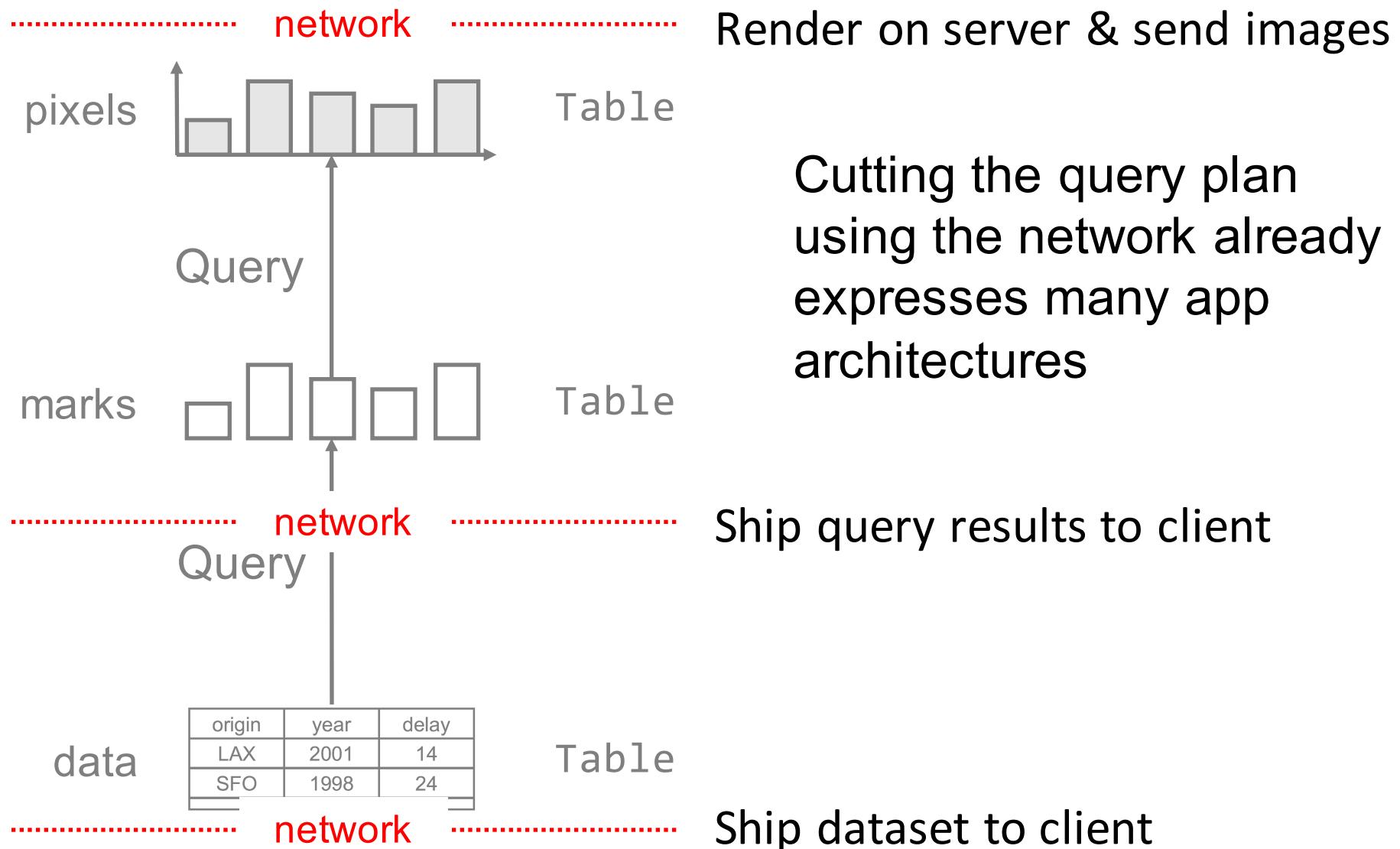


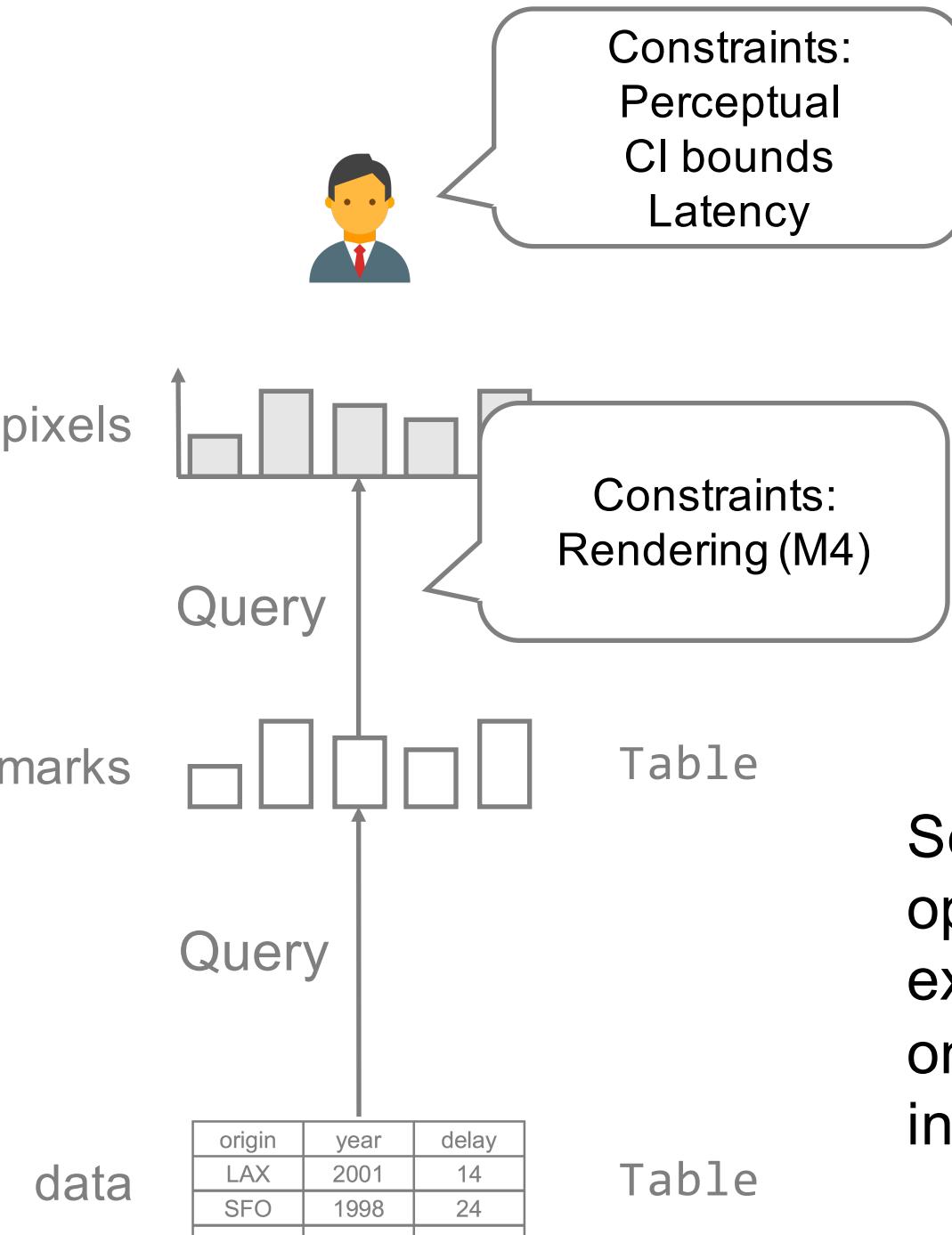
Logical model:

pipeline as a big query

input data, marks, pixels
as tables.

A Relational Perspective

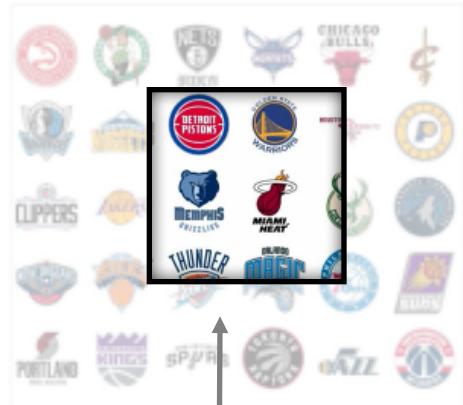




Several earlier push-down optimizations can now be expressed as constraints on the query output or intermediate results

Example: Single View Interactions

- Using Kyrix as an example



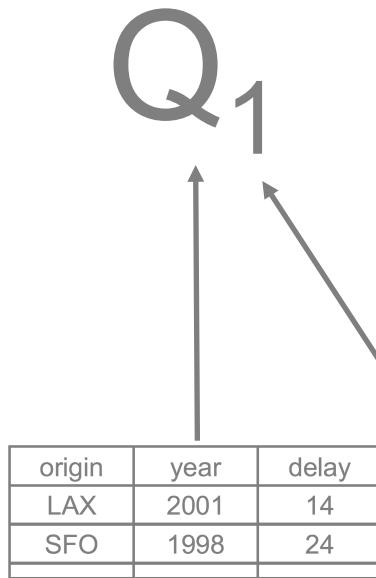
Q_1

origin	year	delay
LAX	2001	14
SFO	1998	24

```
SELECT x, y, img  
FROM NBA_icons  
WHERE x >= ? AND  
      x < ? AND  
      y >= ? AND  
      y < ?
```

At its core, Kyrix can be expressed as parameterized filter queries!

Sufficient to infer RTree indexes and caching



```
SELECT x, y, img  
FROM NBA_icons, viewport vp  
WHERE x >= vp.minx AND  
      x < vp.maxx AND  
      y >= vp.miny AND  
      y < vp.maxy
```

We can remove the “?”s by “relational-izing” the current viewport. This gives us freedom to redefine vp as a view

viewport(minx, maxx, miny, maxy)
user interaction as data



Q₁

origin	year	delay
LAX	2001	14
SFO	1998	24

```
SELECT x, y, img  
FROM NBA_icons, last_vp vp  
WHERE x >= vp.minx AND  
      x <  vp.maxx AND  
      y >= vp.miny AND  
      y <  vp.maxy
```

```
last_vp = SELECT *  
          FROM viewport  
          ORDER BY tstamp DESC  
          LIMIT 1
```

viewport(minx, maxx, miny, maxy, t)
user interaction as data



```
SELECT x, y, img  
FROM NBA_icons, last_vp vp  
WHERE x > vp.minx AND
```

Manipulating viewport view definition
enables historical replay, undo, ...
for free!

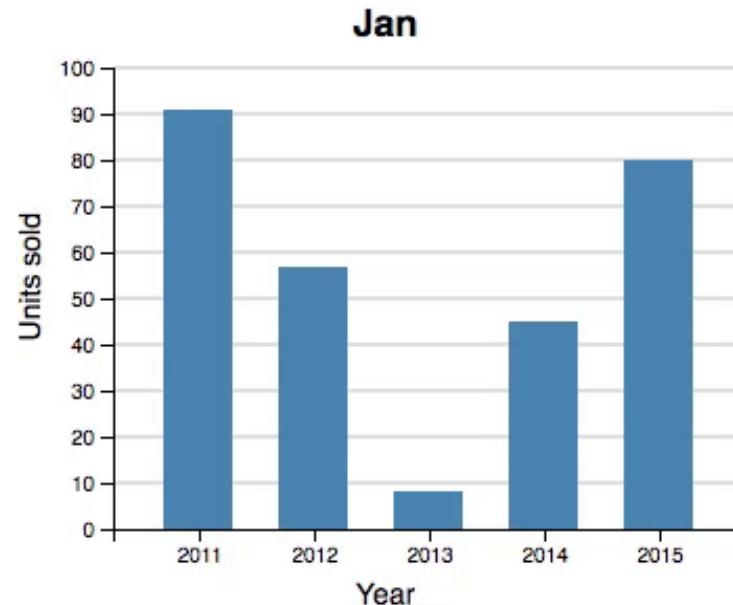
origin	year	delay
LAX	2001	14
SFO	1998	24

viewport(minx, maxx, miny, maxy, **t**)
user interaction as data

Another Benefit of Relationalizing

This slide contains animation. See powerpoint slides to see animation

Jan
Feb
Mar
Apr
May
Jun
Jul
Aug
Sep
Oct
Nov
Dec



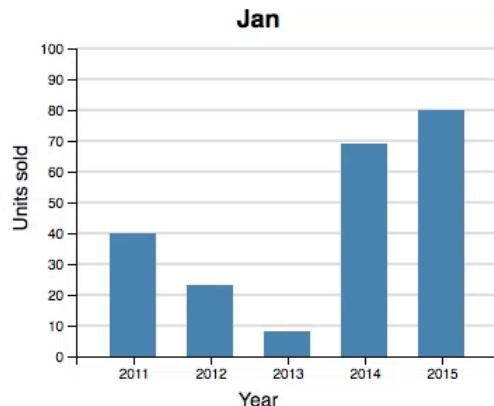
But what if
requests take
time?

Request

Response

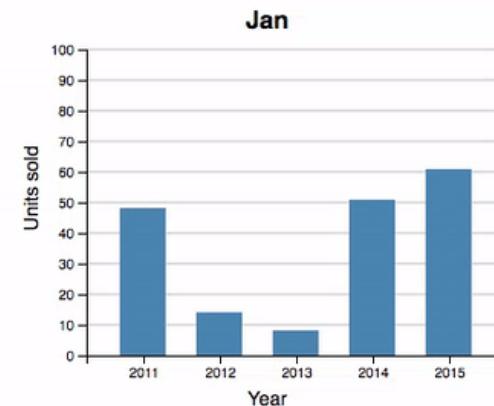
This slide contains animation. See powerpoint slides to see animation

Jan
Feb
Mar
Apr
May
Jun
Jul
Aug
Sep
Oct
Nov
Dec



No CC

Jan
Feb
Mar
Apr
May
Jun
Jul
Aug
Sep
Oct
Nov
Dec



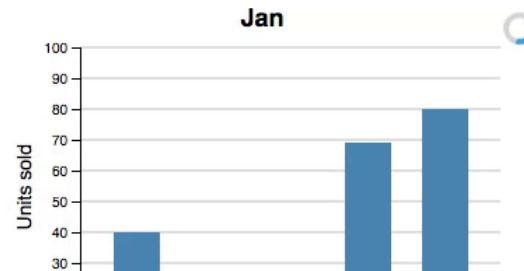
Serial Order

Jan
Feb
Mar
Apr
May
Jun
Jul
Aug
Sep
Oct
Nov
Dec

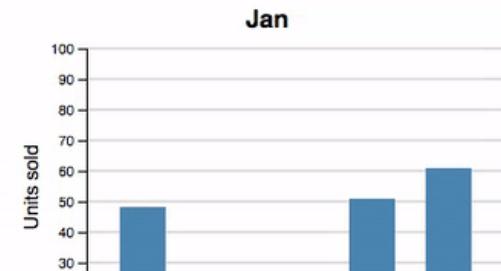
All of these are sensible choices for a designer when dealing with latencies.

Historical Small Multiples

Jan
Feb
Mar
Apr
May
Jun
Jul



Jan
Feb
Mar
Apr
May
Jun
Jul

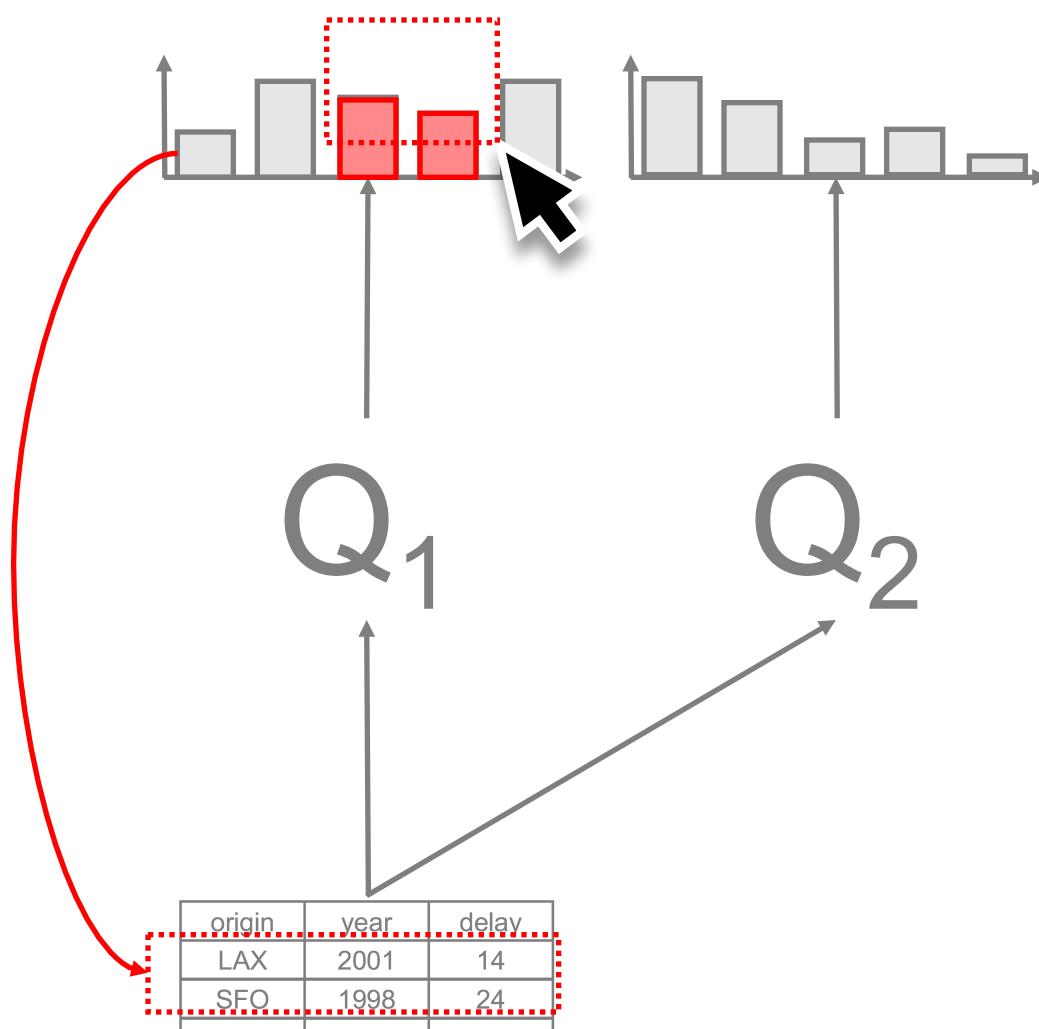


Picking concurrency control
for interaction design is simple by
relationalizing user inputs & responses

J
F
M
A
M
J
J
A
S

Historical Small Multiples

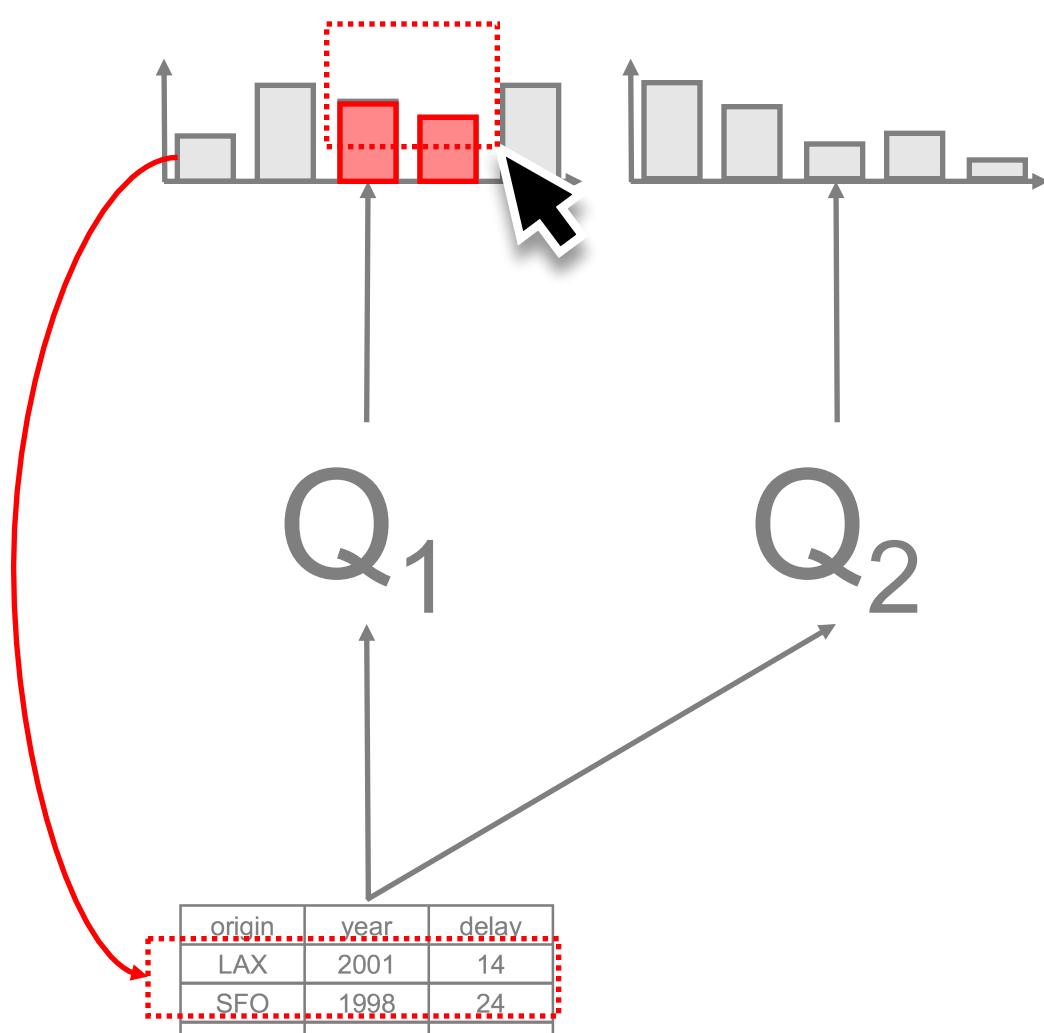
Example: Multi-view Interactions



$Q_2 = \text{SELECT } x, f(y)$
 $\text{GROUP BY } x$

User selects bars in
 Q_1 's chart...

Example: Multi-view Interactions



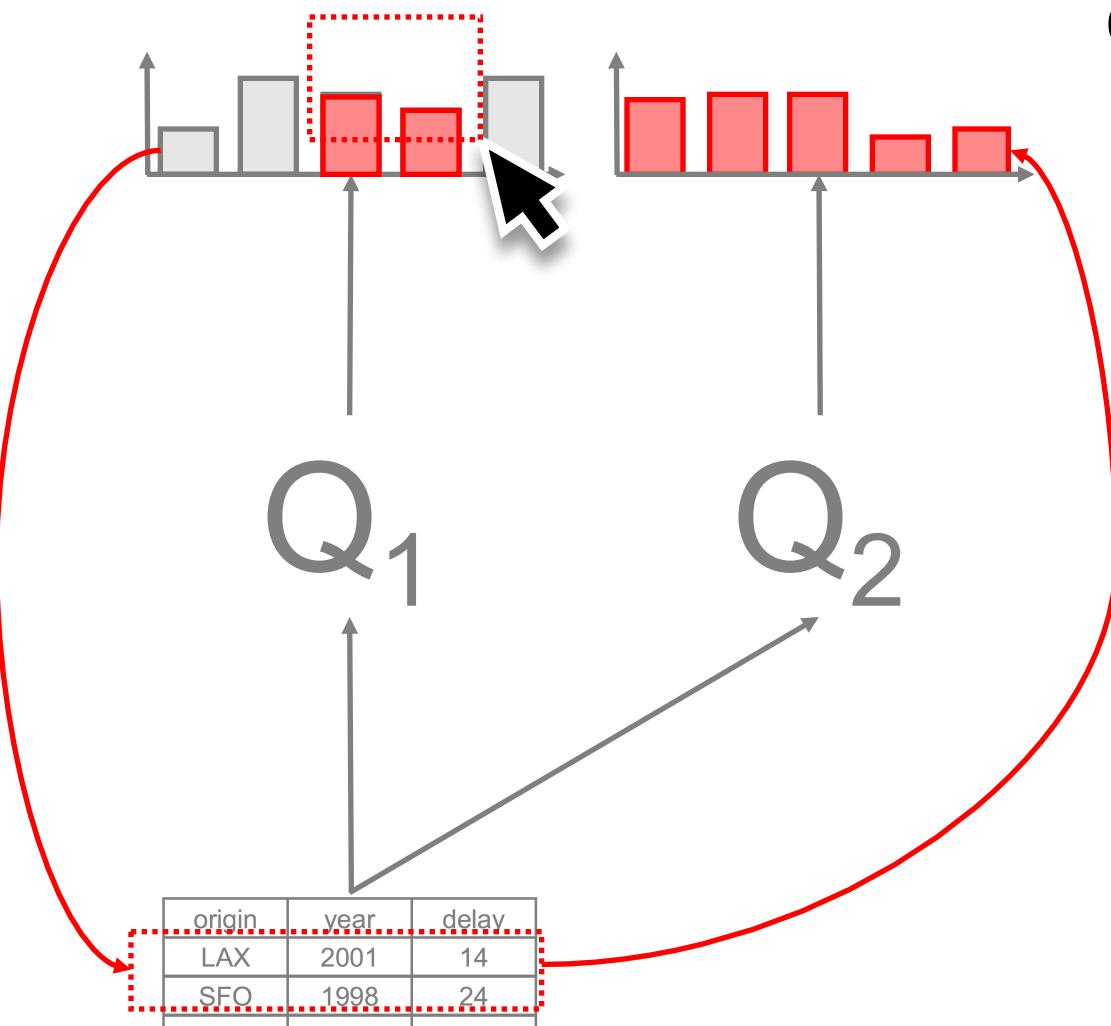
$Q_2 = \text{SELECT } x, f(y)$
WHERE $\text{id} \in$

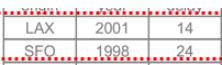
LAX	2001	14
SFO	1998	24

GROUP BY x

Update Q_2 's chart by adding the WHERE clause to Q_2 .

Many apps do this by manipulating SQL string literals to construct the query!



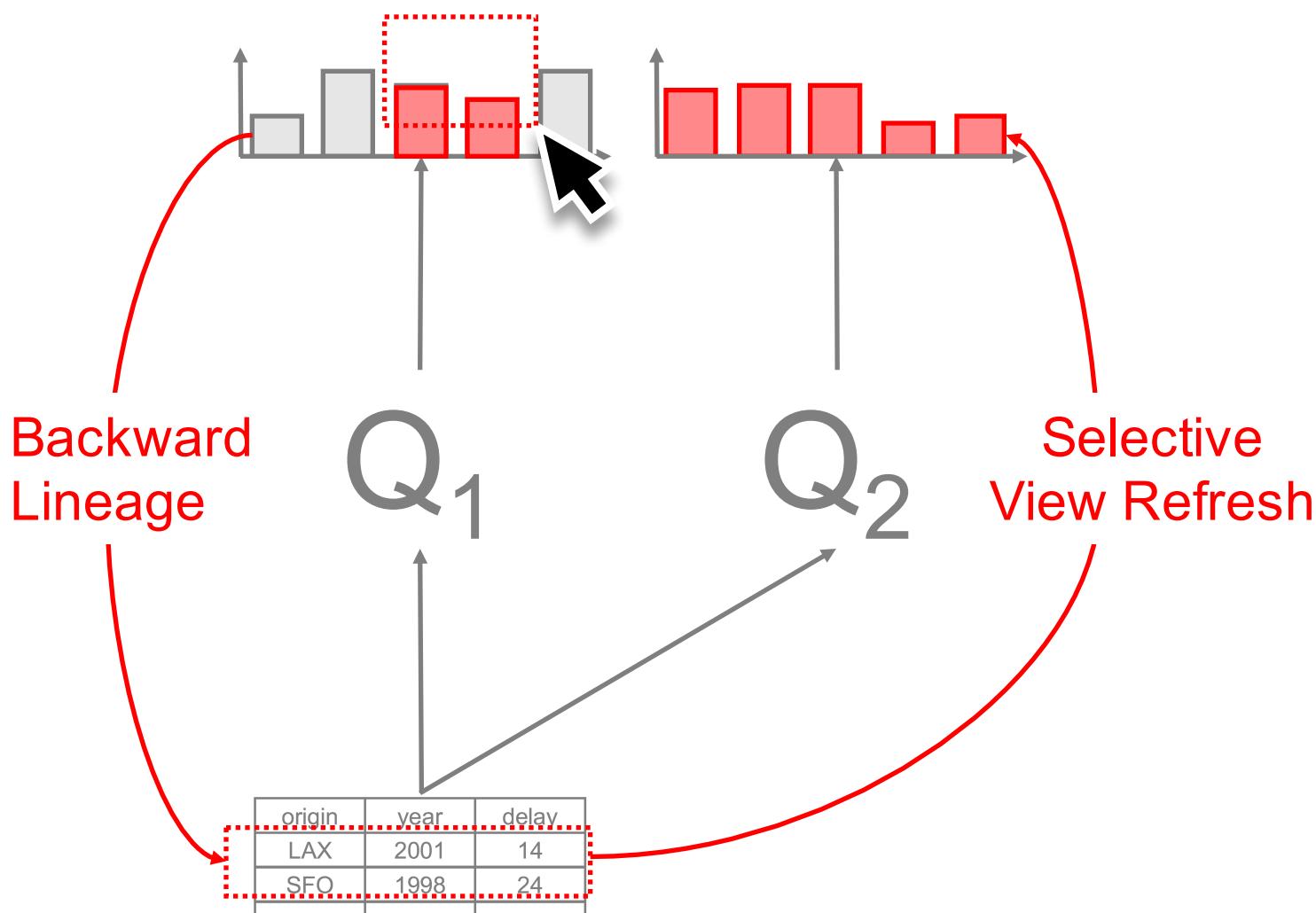
$Q_2 = \text{SELECT } x, f(y)$
 WHERE $\text{id} \in$ 
 GROUP BY x

Syntactic manipulation rather than logical spec

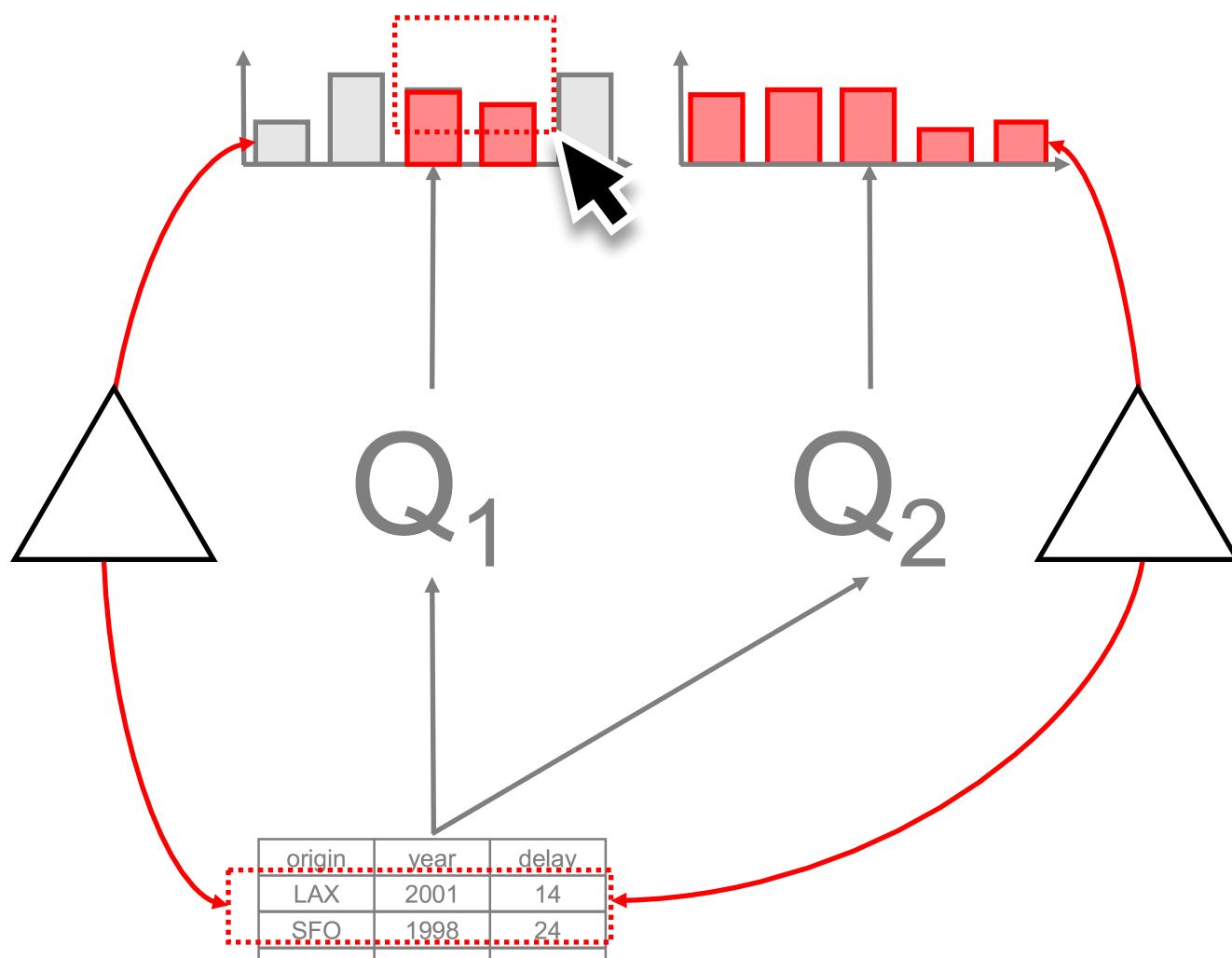


Can do better?

Smoke: Interaction as Lineage



Smoke: Interaction as Lineage

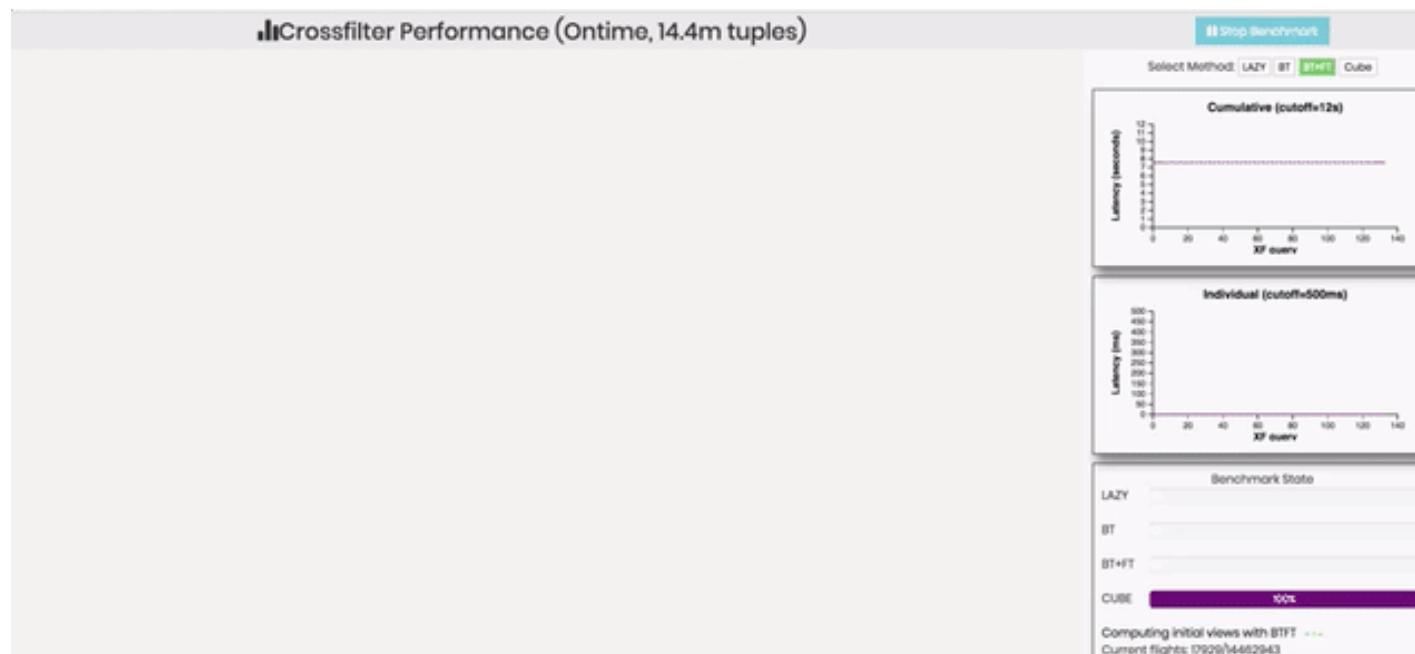


Builds lineage
indexes while
running Q_1 and Q_2
w/ low overhead

This slide contains animation. See powerpoint slides to see animation

Smoke: Interaction as Lineage

Lineage enables <100ms interactivity
Avoids data cube precomputation (mins or hrs)

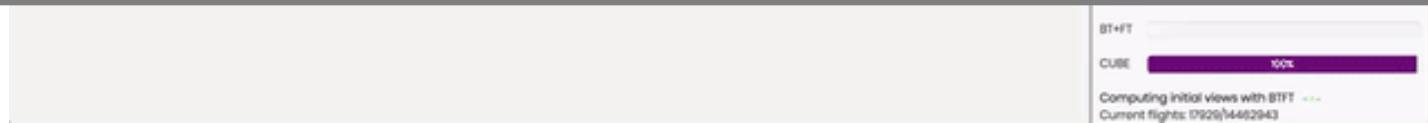


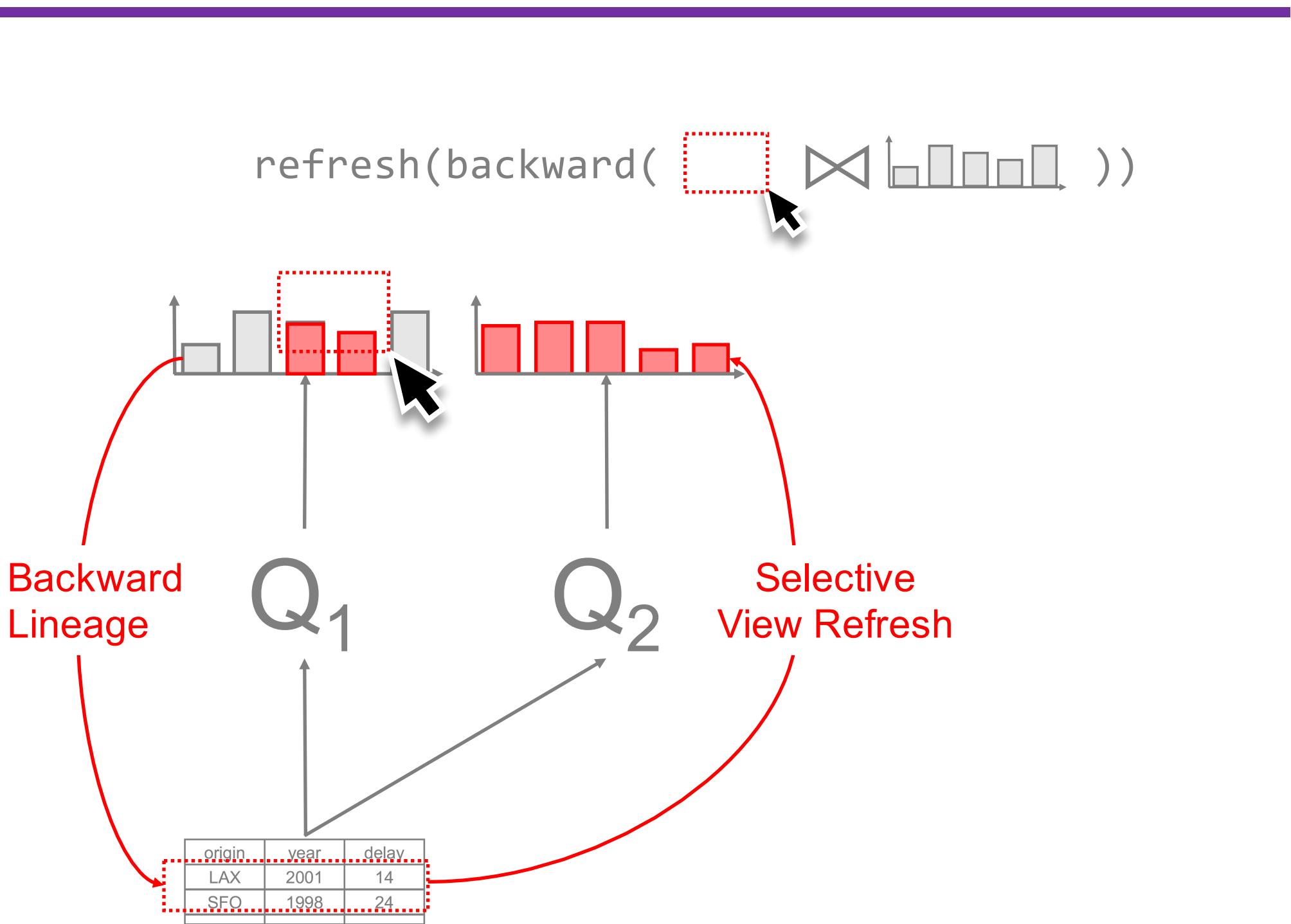
Smoke: Interaction as Lineage

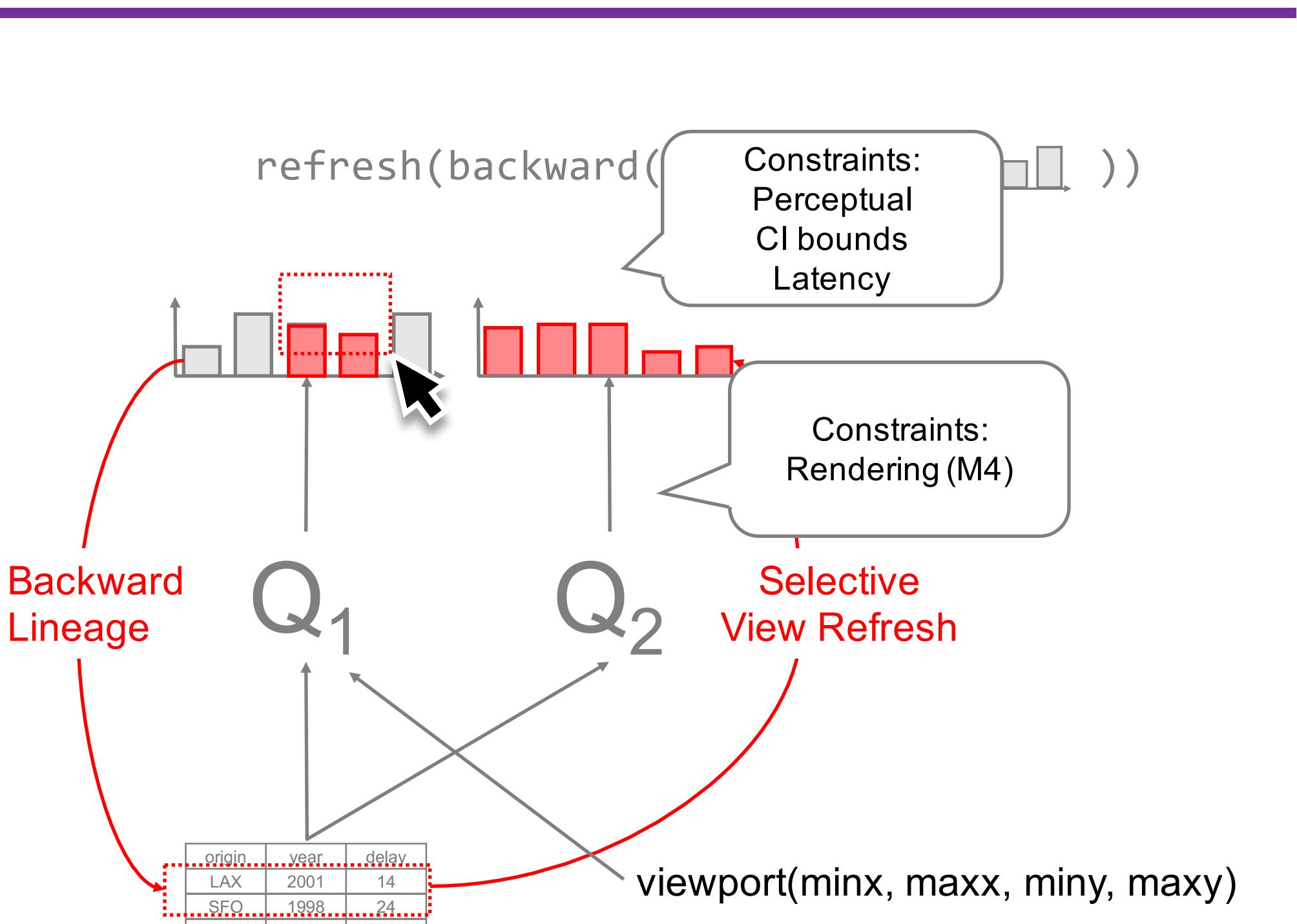
Lineage enables <100ms interactivity

Benefits

*Any visualization expressible as lineage
(most coordinated visualizations)
can be optimized *automatically**







Constraints



interaction(vis(database))

SQL(Lineage(SQL))

Stepping Back

Hacking entire SVD stack is hard

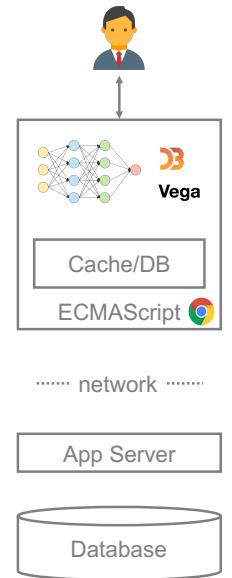
Programming API is important

Users often have existing data, analyses, designs

Wins are from moving up in semantics

Data flow mechanisms for execution

Higher level semantics for optimization



End-to-end relational approach needs to draw from...

Hierarchical models (for layout, subgraphs, etc)

Second order logic (changing group-by attrs)

Ordered relations (most vis is ordered)

...

Open Problem

**Algebra to compose
data + interaction + design + task**



eugenewu.net

