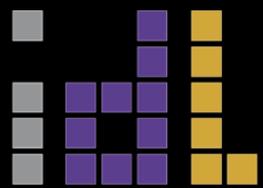


# Data to the People?

Reflections on Trying to Help People Struggle Less with Data

**Jeffrey Heer** @jeffrey\_heer

University of Washington



2010

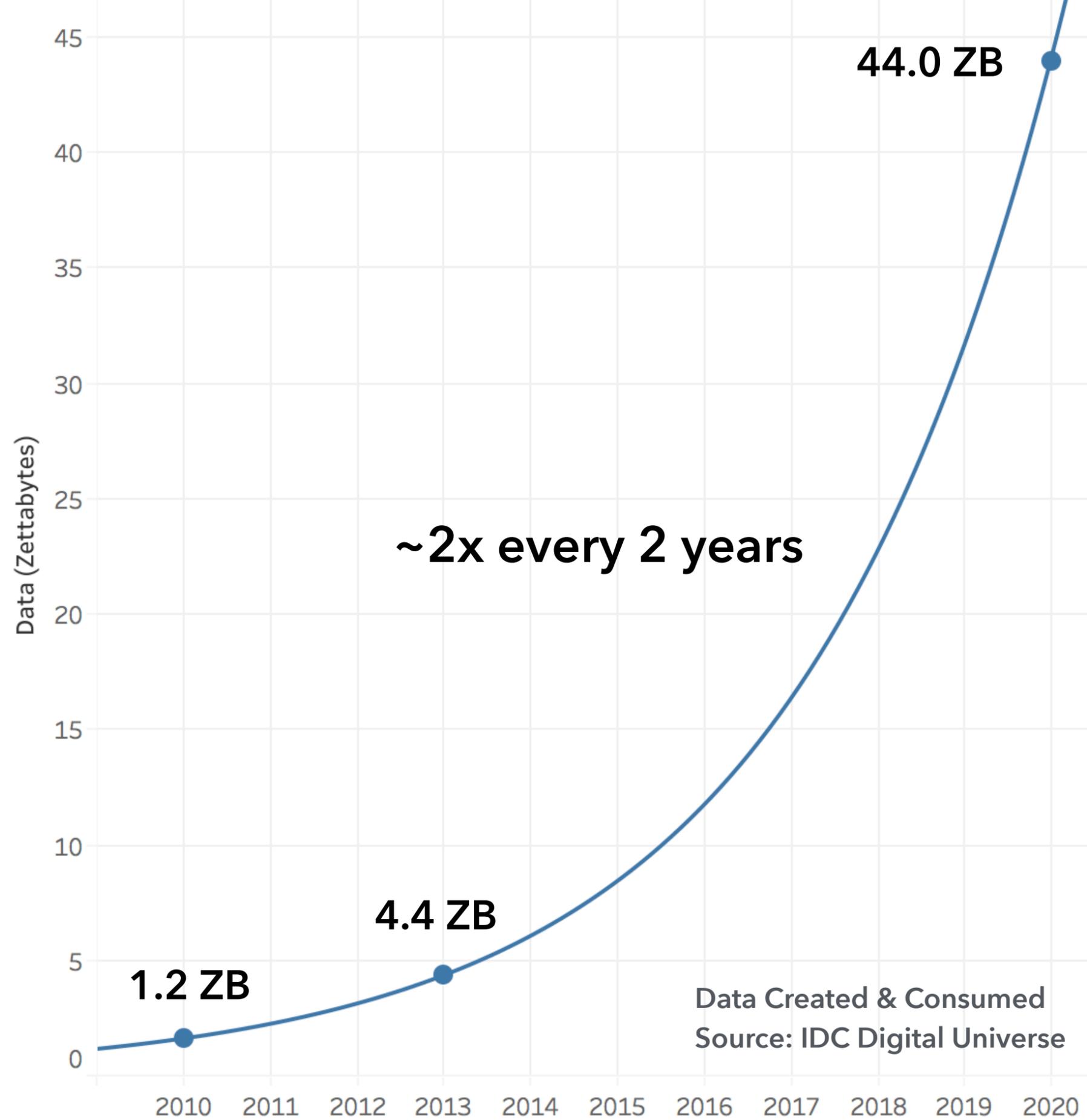


MTV  
season two

A digital tunnel of binary code (0s and 1s) receding into the distance, creating a perspective effect. The code is rendered in a light blue color against a darker blue background. The tunnel is formed by two curved walls of binary digits that meet at a vanishing point in the center.

**“Big Data”**

**2010: 1,200 exabytes  
& exponential growth**



The ability to take data—to be able to **understand** it, to **process** it, to **extract value** from it, to **visualize** it, to **communicate** it—that's going to be a hugely important skill in the next decades, ... because now we really do have essentially free and ubiquitous data. So the complimentary **scarce factor** is the ability to **understand** that data and **extract value** from it.

Hal Varian, Google's Chief Economist  
*The McKinsey Quarterly*, Jan 2009



**Four major influences** act on data analysis today:

1. The formal theories of statistics.
2. Accelerating developments in computers and display devices.
3. The challenge, in many fields, of more and larger bodies of data.
4. The emphasis on quantification in a wider variety of disciplines.



Accordingly, both approaches and techniques need to be structured so as to **facilitate human involvement and intervention.**

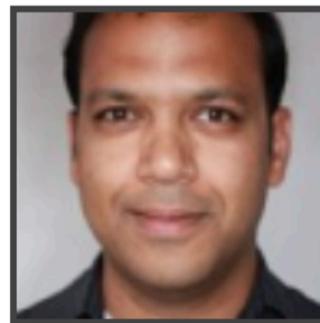
Some implications for effective analysis are: (1) it is essential to have convenience of **interaction of people and intermediate results** and (2) at all stages of data analysis, the outputs need to be **matched to the capabilities of the people who use it and want it.**

**d<sup>p</sup>** = data to the people

Facilitating interactions between people and data throughout the analytic lifecycle.



Jeff Heer  
Stanford



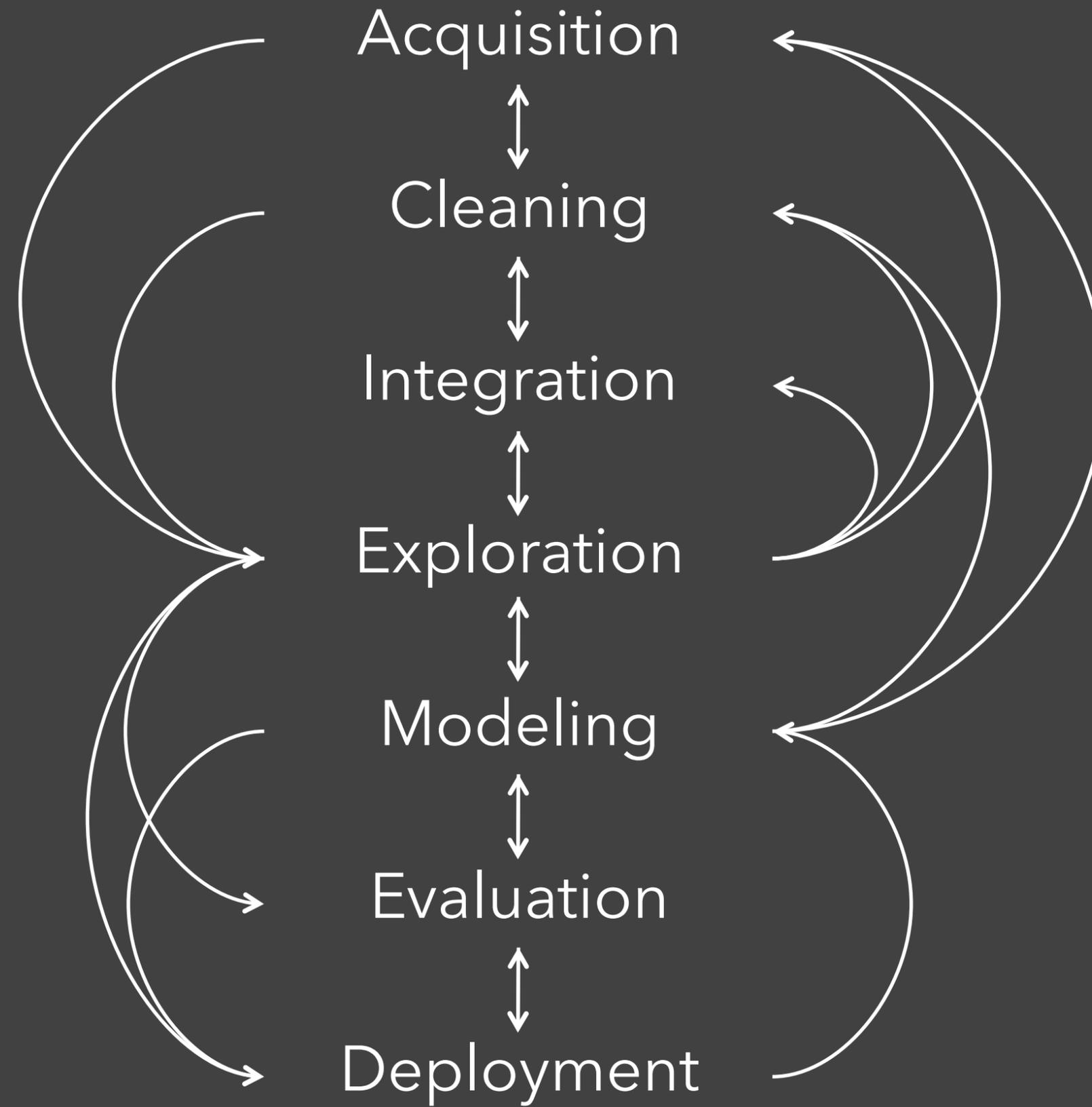
Maneesh Agrawala  
Berkeley

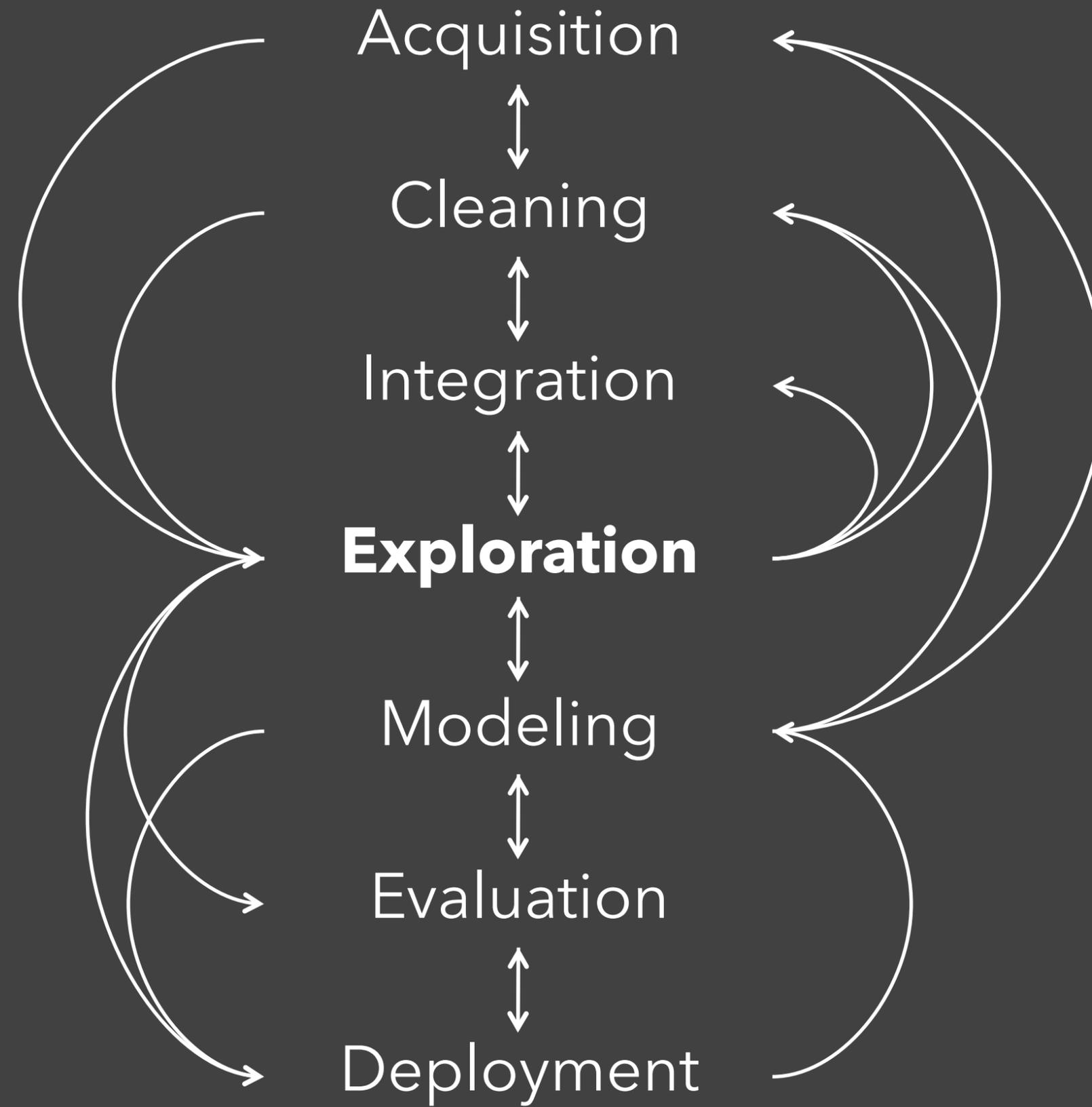


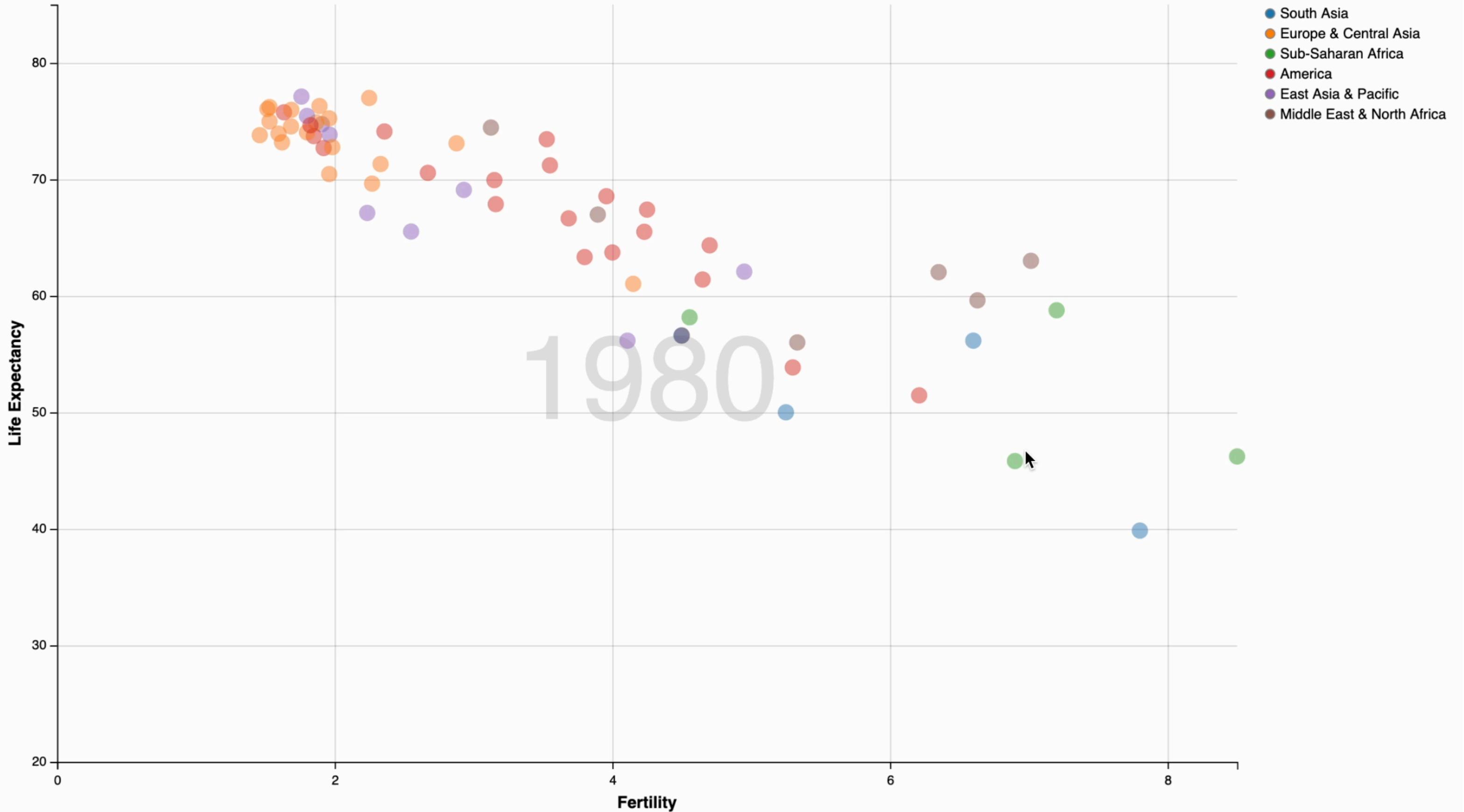
Joe Hellerstein  
Berkeley

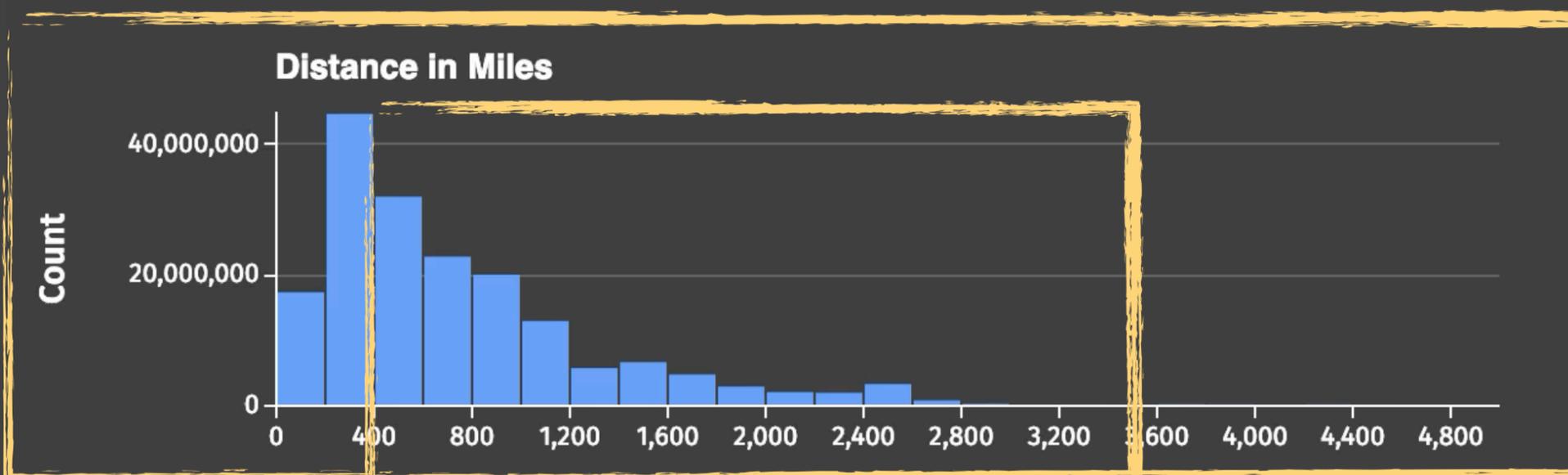
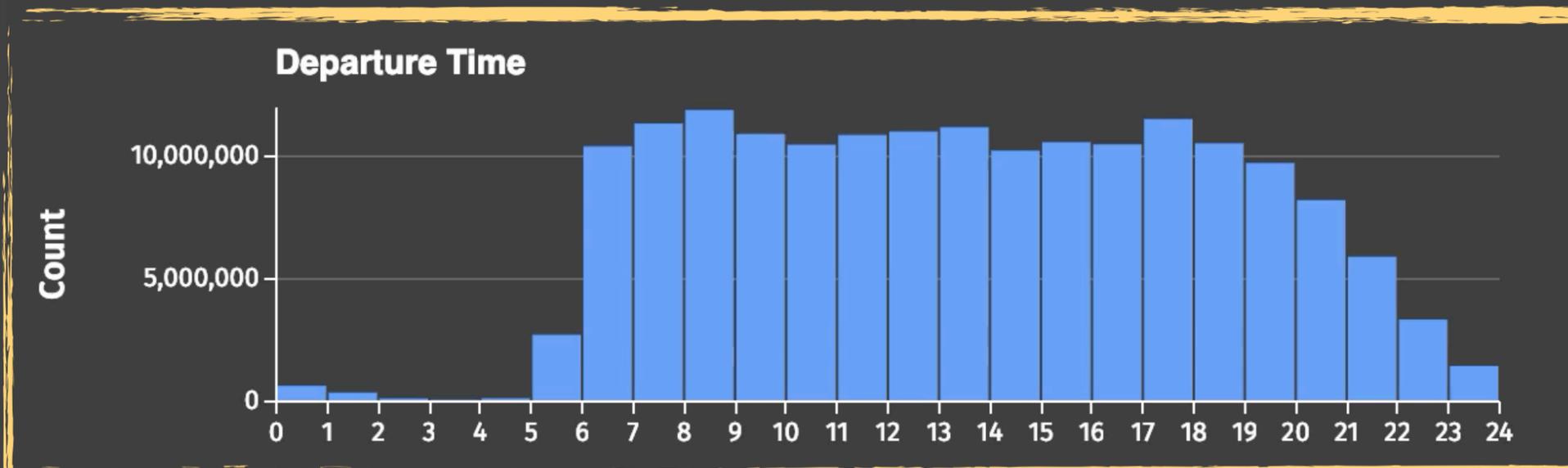
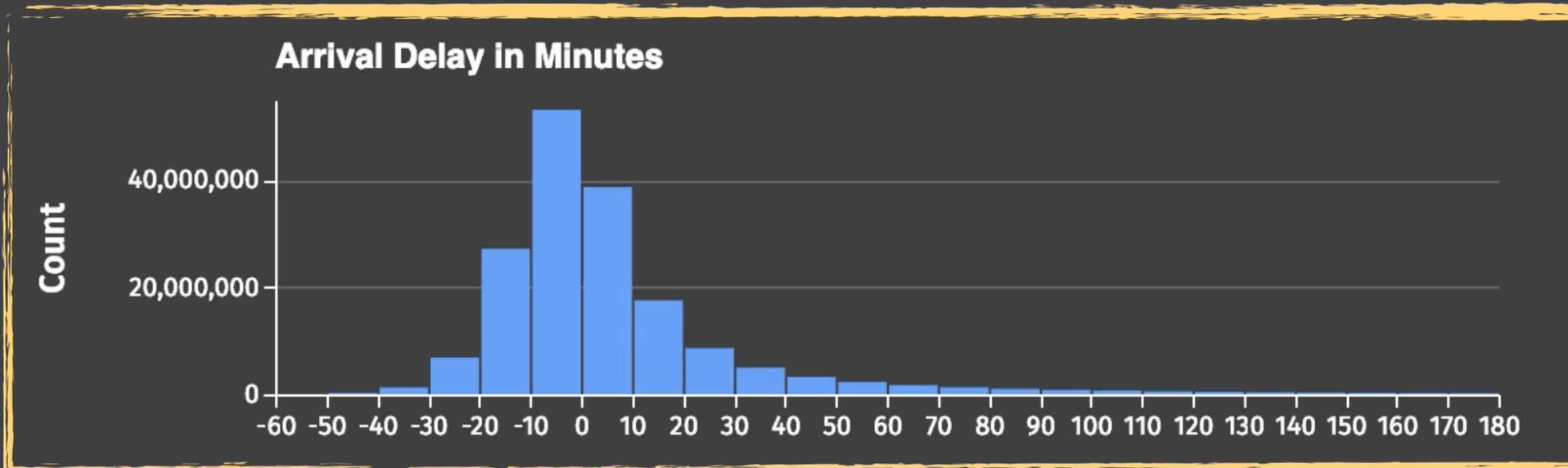


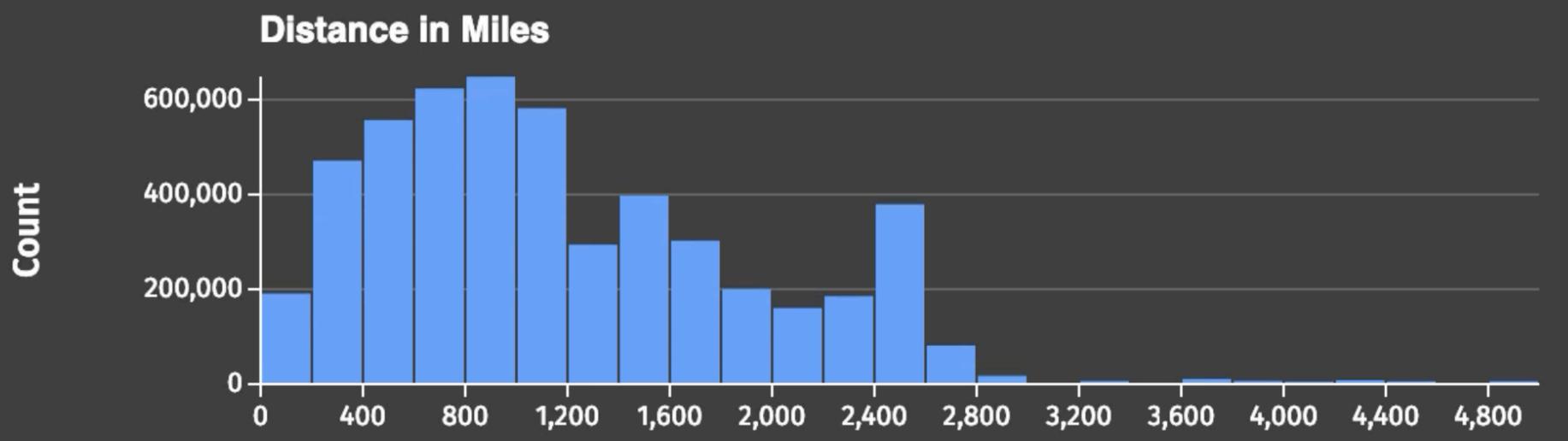
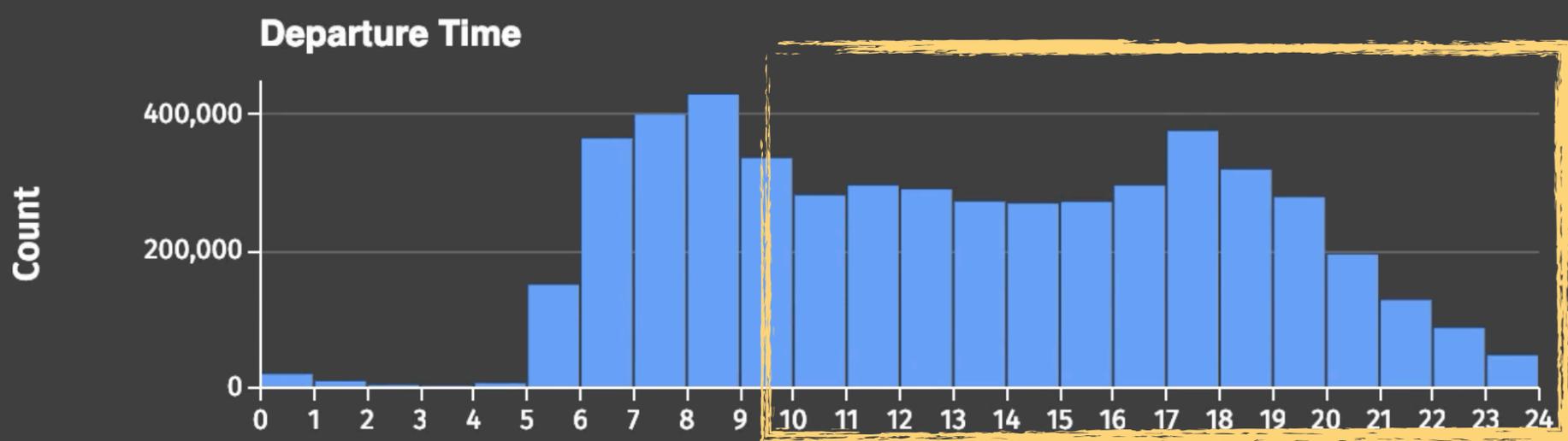
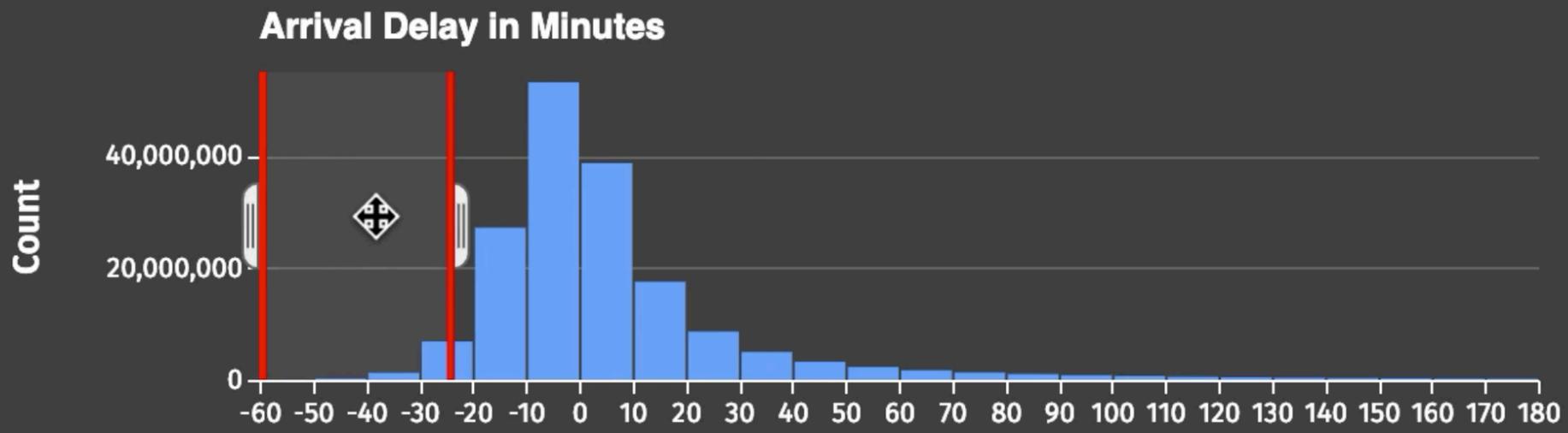
Tapan Parikh  
Berkeley



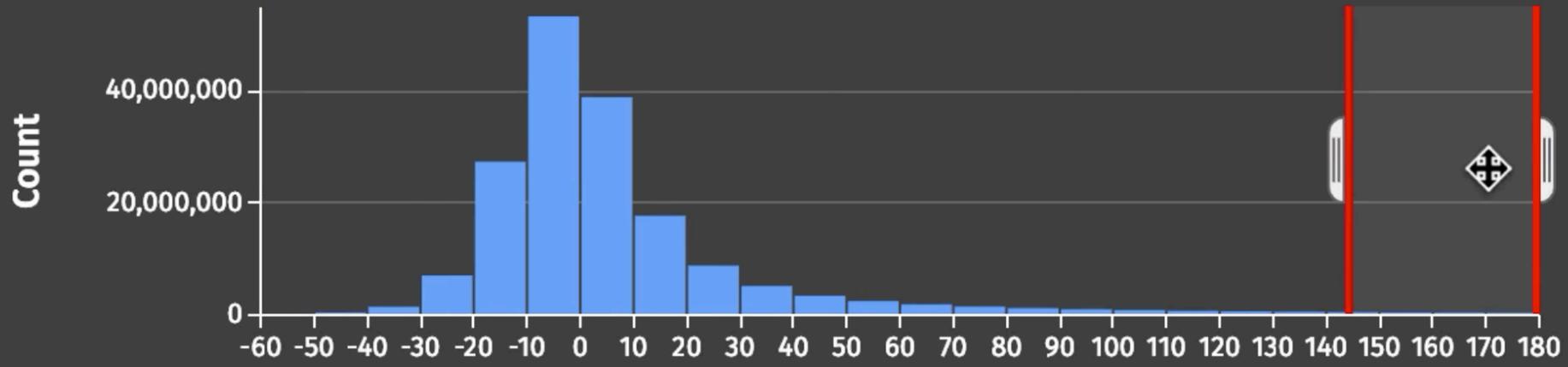




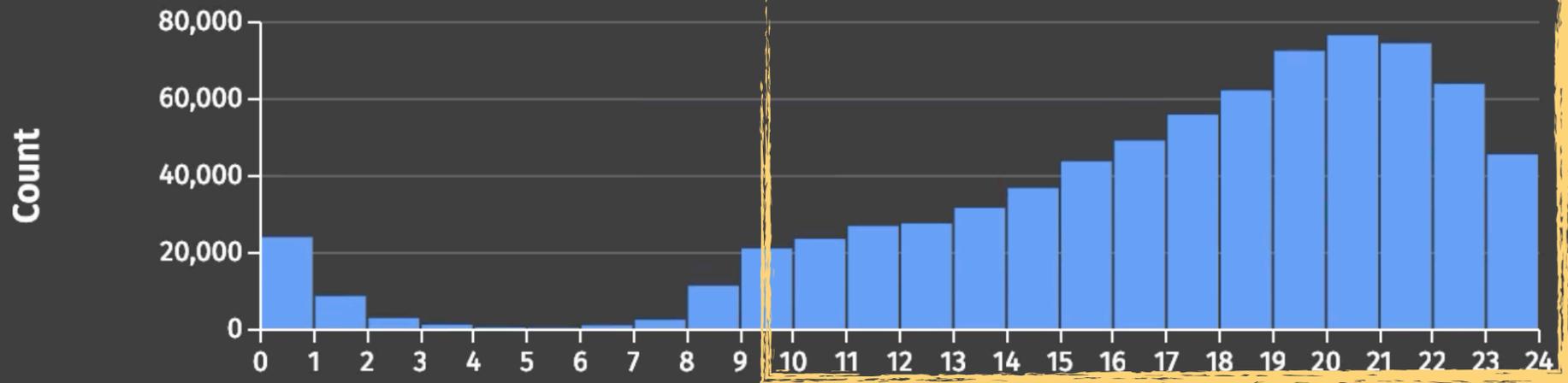




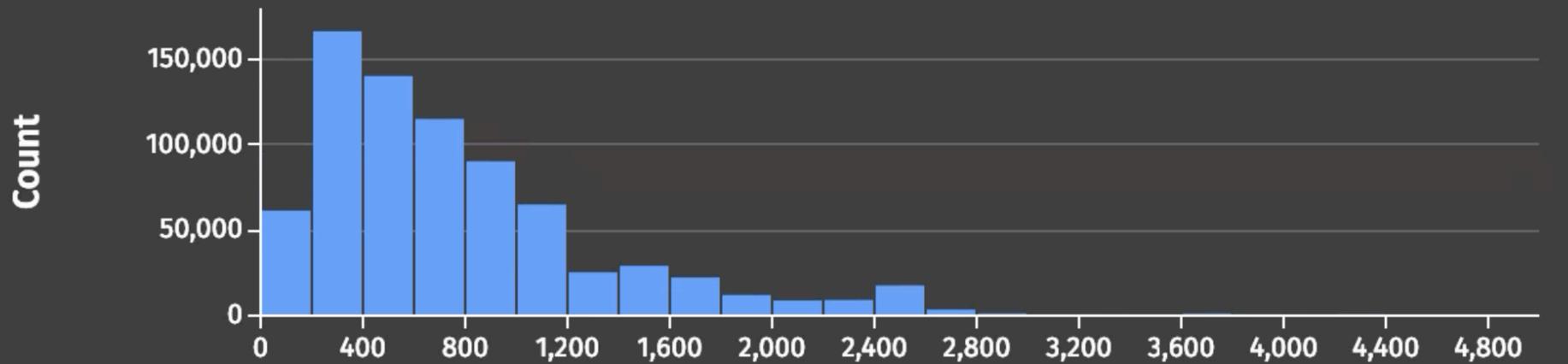
### Arrival Delay in Minutes

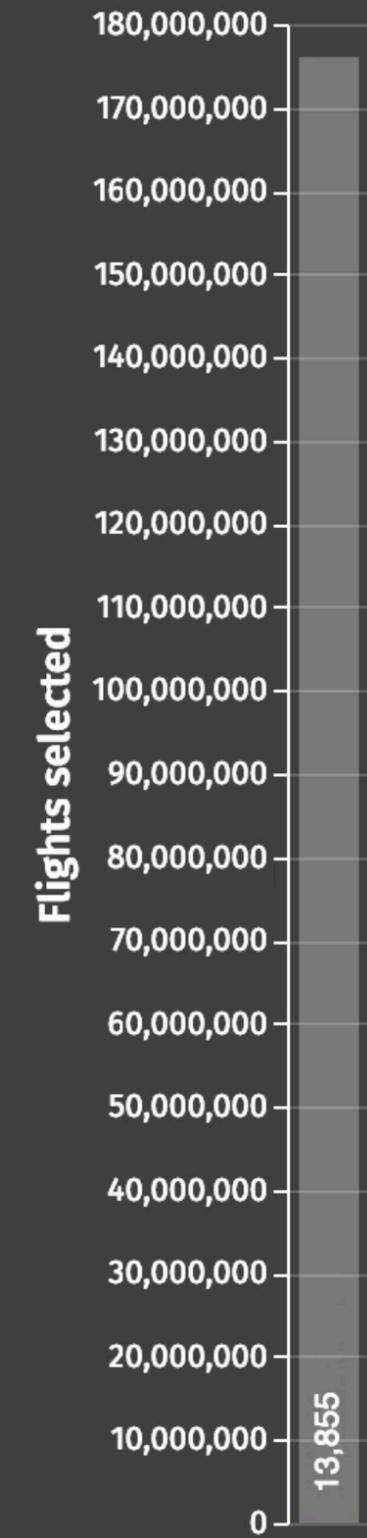
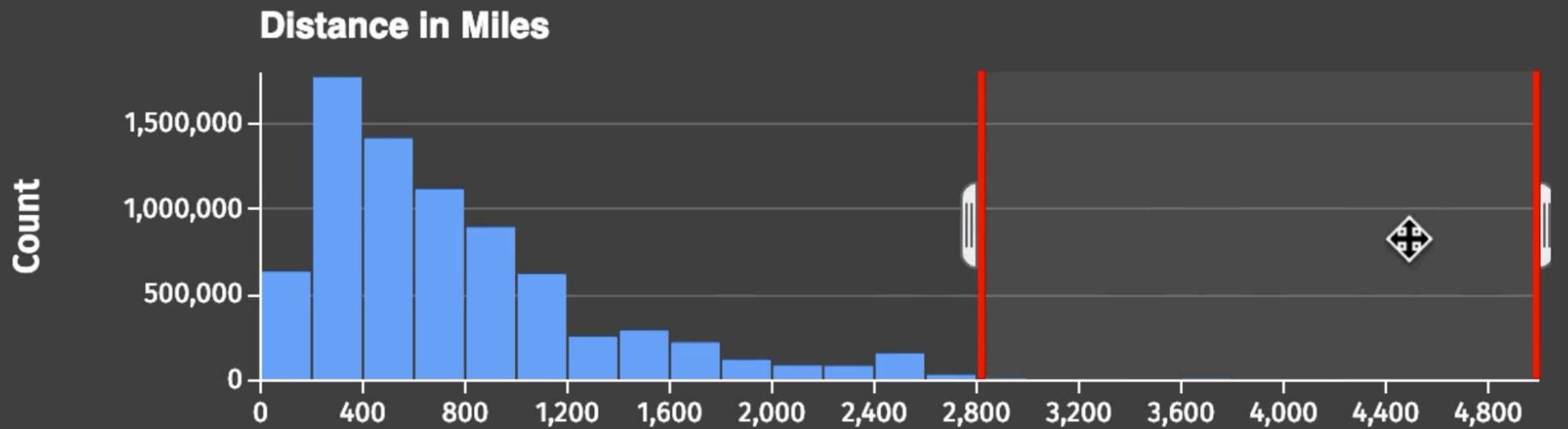
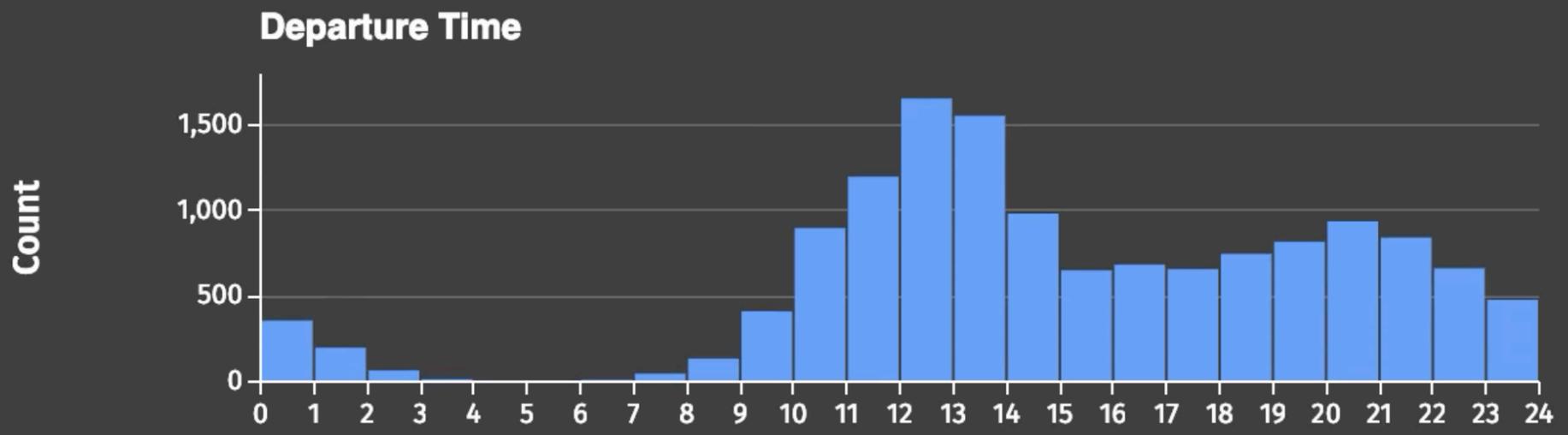
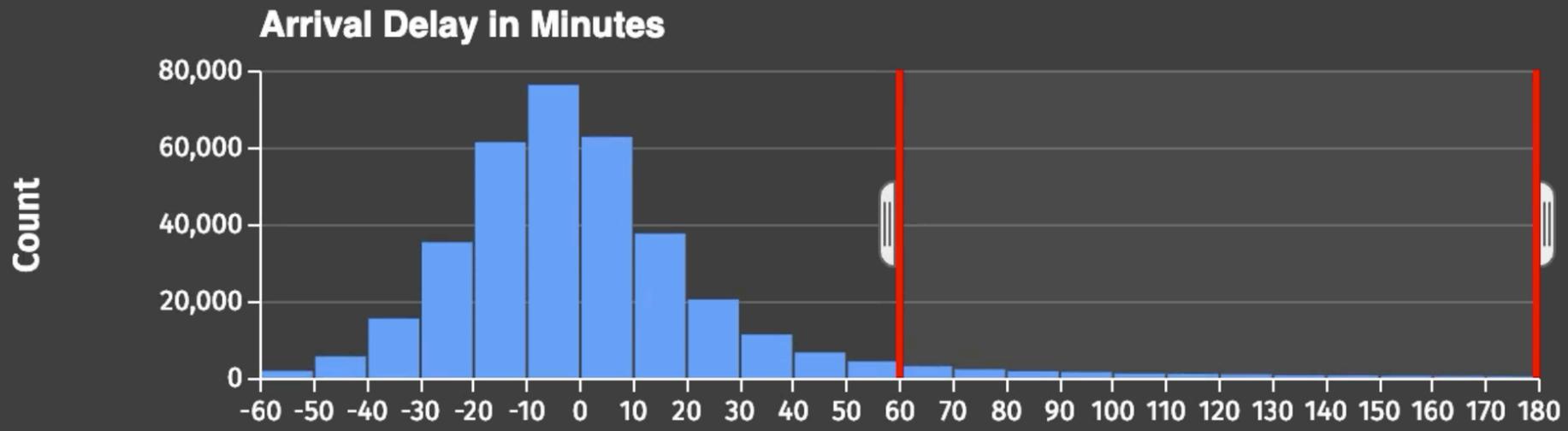


### Departure Time



### Distance in Miles





How might we flexibly author  
interactive visualizations?



# Visualization System Evolution

**2002-09** Object-Oriented Libraries (Prefuse, Flare)

Requires software engineering to use and extend

**2009-12** JavaScript Embedded DSLs (Protovis, D3)

Rich composable primitives, but specific to JavaScript

Interaction via standard event callbacks

**2012-??** Declarative DSLs (Vega, Vega-Lite)

Formal specification of visualization and interactions

APIs in various programming languages

# DB Concepts

Columnar storage

Query optimization

Closed language

Data join between

marks and data

Declarative language

Query plan / runtime

Stream management

# Query Evaluation Techniques for Large Databases

GOETZ GRAEFE

Portland State University, Computer Science Department, P.O. Box 75

Database management systems will continue to manage efficient algorithms for accessing and manipulating large required to provide acceptable performance. The advent of database systems will not solve this problem. On the contrary, exacerbate the problem: In order to manipulate large efficiently as today's database systems manipulate simple algorithms and software will become more complex, an algorithm and architectural issues is essential for the design of software.

This survey provides a foundation for the design of execution facilities in new database management systems, practical query evaluation techniques for both relational systems, including iterative execution of complex query sort- and hash-based set-matching algorithms, types of their implementation, and special operators for emerging

Categories and Subject Descriptors: E.5 [Data]: Files and Systems—query processing

General Terms: Algorithms, Performance

Additional Key Words and Phrases: Complex query evaluation plans; extensible database systems, iterative systems, operator model of parallelization, parallel systems, set-matching algorithms, sort-hash duality

## INTRODUCTION

Effective and efficient management of large data volumes is necessary in virtually all computer applications, from business data processing to library information retrieval systems, multimedia applications with images and sound, computer-aided design and manufacturing, real-time process control, and scientific computation. While database management systems are standard tools in business data processing, they are only slowly being introduced to all

the other areas.

In most cases, we have traditional reasons. and many applications nance u into ser models : grammi ing this to accep

Permission to copy without fee all or part of this material is granted or distributed for direct commercial advantage, the ACM copyright notice and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment

© 1993 ACM 0360-0300/93/0600-0073 \$01.50

# TelegraphCQ: Continuous Dataflow Processing in an Uncertain World<sup>†</sup>

Sirish Chandrasekaran, Owen Cooper, Amol Deshpande, Michael J. Franklin, Wei Hong\*, Sailesh Krishnamurthy, Sam Madden, Vijayshankar Ramakrishnan

University of California, Berkeley  
\*Intel Berkeley Laboratory  
\*\*IBM Almaden Research Center  
<http://telegraph.cs.berkeley.edu>

## Abstract

Increasingly pervasive networks are leading towards a world where data is constantly in motion. In such a world, conventional techniques for query processing, which were developed under the assumption of a far more static and predictable computational environment, will not be sufficient. Instead, query processors based on adaptive dataflow will be necessary. The Telegraph project has developed a suite of novel technologies for continuously adaptive query processing. The next generation Telegraph system, called TelegraphCQ, is focused on meeting the challenges that arise in handling large streams of continuous queries over high-volume, highly-variable data streams. In this paper, we describe the system architecture and its underlying technology, and report on our ongoing implementation effort, which leverages the PostgreSQL open source code base. We also discuss open issues and our research agenda.

## 1 INTRODUCTION

The deployment of pervasive communications infrastructure ranging from short-range wireless ad hoc sensor networks to globe-spanning intra- and internets has enabled new applications that process, analyze, and react to disparate data in a near real-time manner. Examples include: event-based business processing, profile-based data dissemination, and query processing over streaming data sources such as network monitors, sensors, and mobile devices. Such applications present challenges that cannot be met by existing database and data management technology. These challenges stem from their large scale,

<sup>†</sup> This work was funded in part by the NSF under grants IIS-0086057, SI-0122599, and EIA-0207603, by the IBM Faculty Partnership Award program, and by research funds from Intel, Microsoft, and the UC MICRO program.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment

Proceedings of the 2003 CIDR Conference

their deeply-networked environment, and the

In emerging commodity of information is realized only when data can be as locations, data is moving and change management for processing that information as the

### 1.1 Data Management

Traditional inappropriate a processing support Streaming systems where

“pulling” data our target application data is continuous subtle differences of the query traditional query and handling must instead data streams constraints (data must be spooled to a

Further detailed statistics of Finally, it time and/or over stream than query

Control applicatic which query arrives at of active

# The Design of the Borealis Stream Processing Engine

Daniel J. Abadi<sup>1</sup>, Yanif Ahmad<sup>2</sup>, Magdalena Balazinska<sup>1</sup>, Uğur Çetintemel<sup>2</sup>, Mitch Cherniack<sup>3</sup>, Jeong-Hyon Hwang<sup>2</sup>, Wolfgang Lindner<sup>1</sup>, Anurag S. Maskey<sup>3</sup>, Alexander Rasin<sup>2</sup>, Esther Ryzkina<sup>3</sup>, Nesime Tatbul<sup>2</sup>, Ying Xing<sup>2</sup>, and Stan Zdonik<sup>2</sup>

<sup>1</sup>MIT  
Cambridge, MA

<sup>2</sup>Brown University  
Providence, RI

<sup>3</sup>Brandeis University  
Waltham, MA

## Abstract

Borealis is a second-generation distributed stream processing engine that is being developed at Brandeis University, Brown University, and MIT. Borealis inherits core stream processing functionality from Aurora [14] and distribution functionality from Medusa [51]. Borealis modifies and extends both systems in non-trivial and critical ways to provide advanced capabilities that are commonly required by newly-emerging stream processing applications.

In this paper, we outline the basic design and functionality of Borealis. Through sample real-world applications, we motivate the need for dynamically revising query results and modifying query specifications. We then describe how Borealis addresses these challenges through an innovative set of features, including revision records, time travel, and control lines. Finally, we present a highly flexible and scalable QoS-based optimization model that operates across server and sensor networks and a new fault-tolerance model with flexible consistency-availability trade-offs.

## 1 Introduction

Over the last several years, a great deal of progress has been made in the area of stream processing engines (SPE). Several groups have developed working prototypes [1, 4, 16] and many papers have been published on detailed aspects of the technology such as data models [2, 5, 46], scheduling [8, 15], and load shedding [9, 20, 44]. While this work is an important first step, fundamental mismatches remain between the requirements of many streaming applications and the capabilities of first-generation systems.

This paper is intended to illustrate our vision of what second-generation SPE's should look like. It is driven by our experience in using Aurora [10], our own prototype, in several streaming applications including the Linear Road Benchmark [6] and several commercial opportunities. We present this vision in terms of our own design considerations for Borealis, the successor to Aurora, but it should

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.

Proceedings of the 2005 CIDR Conference

be emphasized that the issues raised here represent general challenges for the field as a whole. We present specifics of our design as concrete evidence for why these problems are hard and as a first cut at how they might be approached. We envision the following three fundamental requirements for second-generation SPEs:

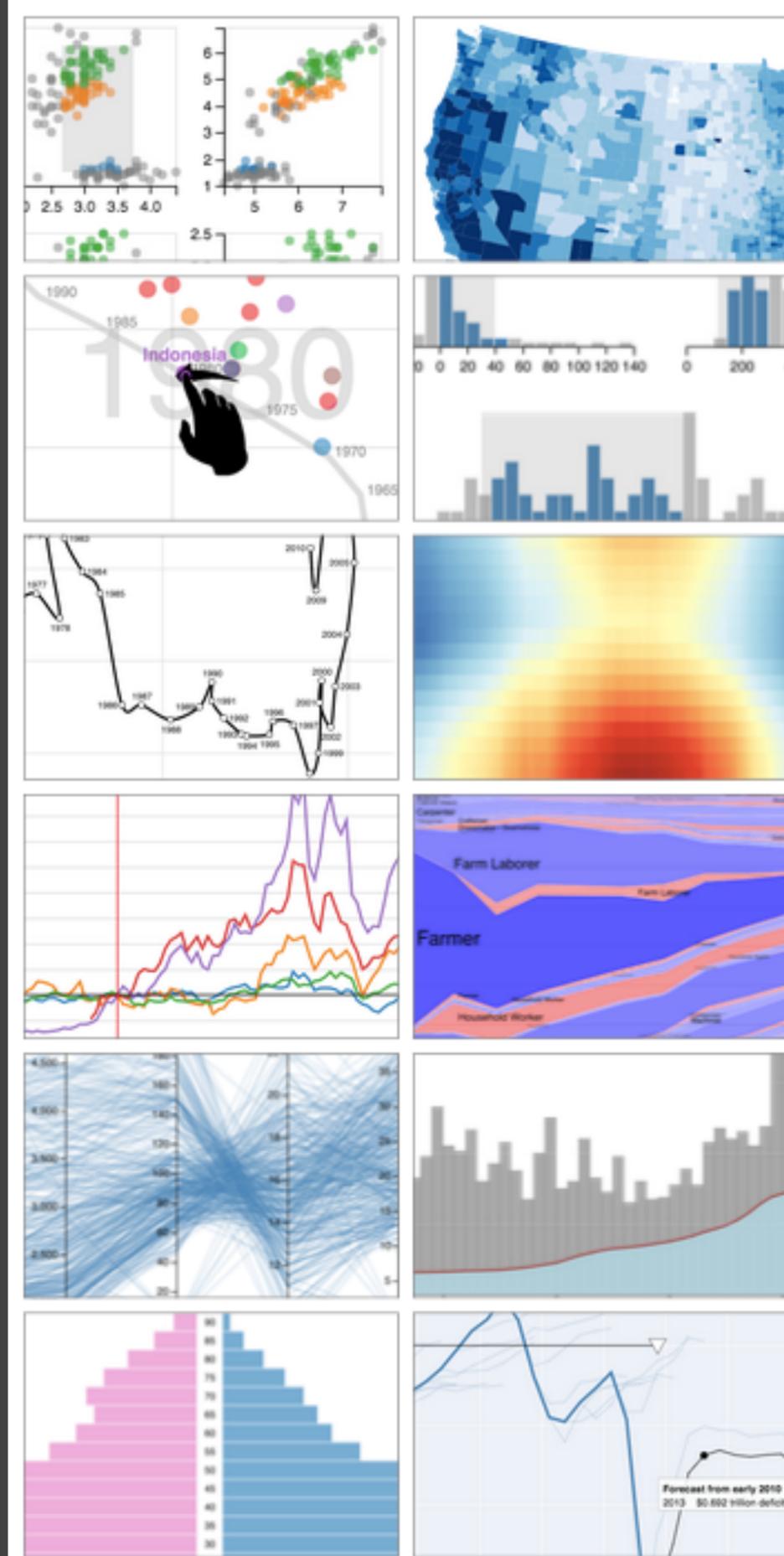
**1. Dynamic revision of query results:** In many real-world streams, corrections or updates to previously processed data are available only after the fact. For instance, many popular data streams, such as the Reuters stock market feed, often include so-called *revision records*, which allow the feed originator to correct errors in previously reported data. Furthermore, stream sources (such as sensors), as well as their connectivity, can be highly volatile and unpredictable. As a result, data may arrive late and miss its processing window, or may be ignored temporarily due to an overload situation [44]. In all these cases, applications are forced to live with imperfect results, unless the system has means to revise its processing and results to take into account newly available data or updates.

**2. Dynamic query modification:** In many stream processing applications, it is desirable to change certain attributes of the query at runtime. For example, in the financial services domain, traders typically wish to be alerted of *interesting* events, where the definition of “interesting” (i.e., the corresponding filter predicate) varies based on current context and results. In network monitoring, the system may want to obtain more precise results on a specific sub-network, if there are signs of a potential Denial-of-Service attack. Finally, in a military stream application from Mitre, they wish to switch to a “cheaper” query when the system is overloaded. For the first two applications, it is sufficient to simply alter the operator parameters (e.g., window size, filter predicate), whereas the last one calls for altering the operators that compose the running query. Although current SPEs allow applications to substitute query networks with others at runtime, such manual substitutions impose high overhead and are slow to take effect as the new query network starts with an empty state. Our goal is to support low overhead, fast, and automatic modifications.

Another motivating application comes again from the financial services community. Universally, people working on trading engines wish to test out new trading strategies as well as debug their applications on historical data before they go live. As such, they wish to perform “time travel” on input streams. Although this last example can be supported

# Vega is a *Visualization Grammar*

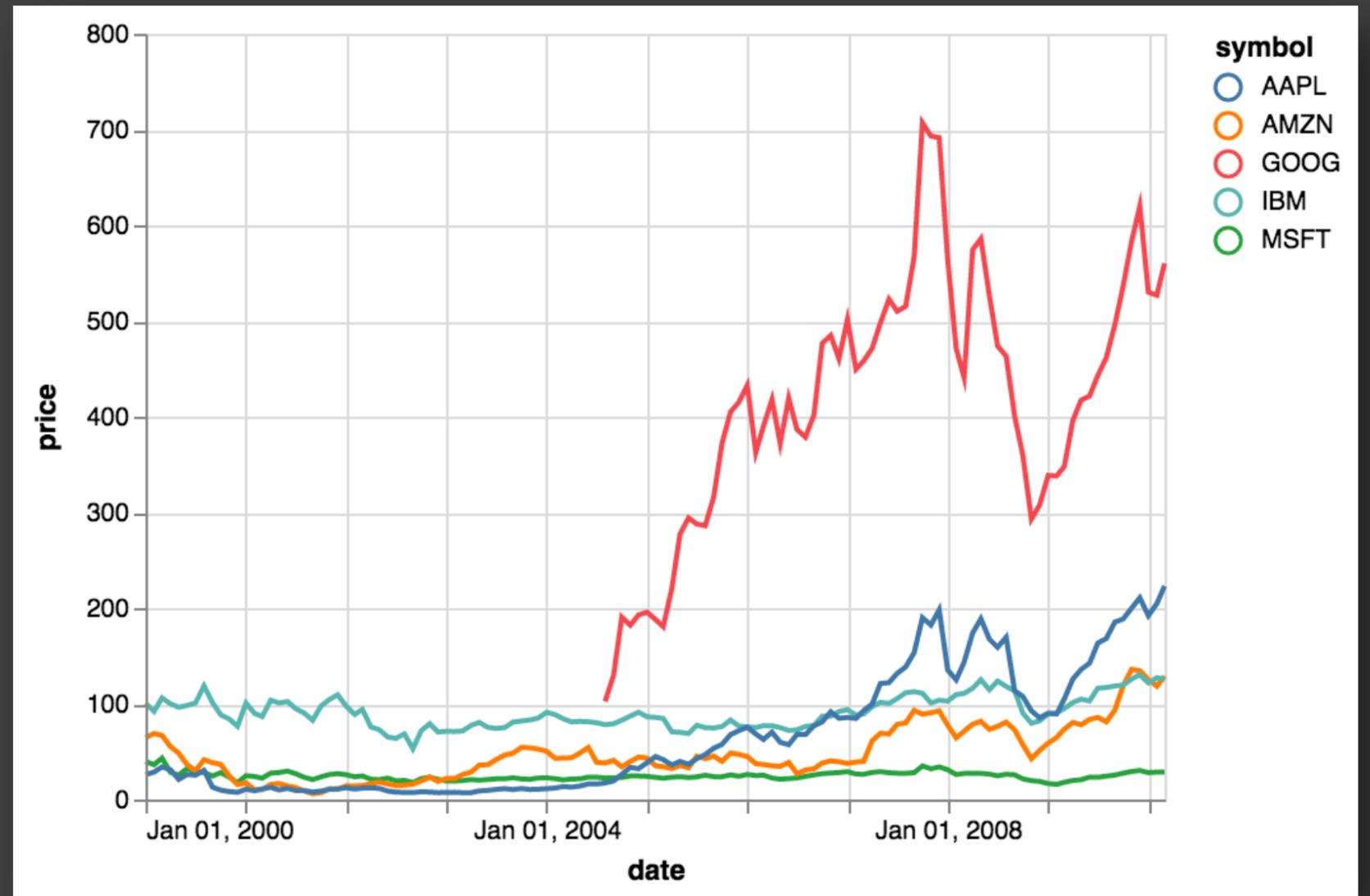
Similar in spirit to how SQL provides a language for expressing database queries, Vega is a high-level language for describing visualizations.



# Vega is a *Visualization Grammar*

## Vega Altair Python API

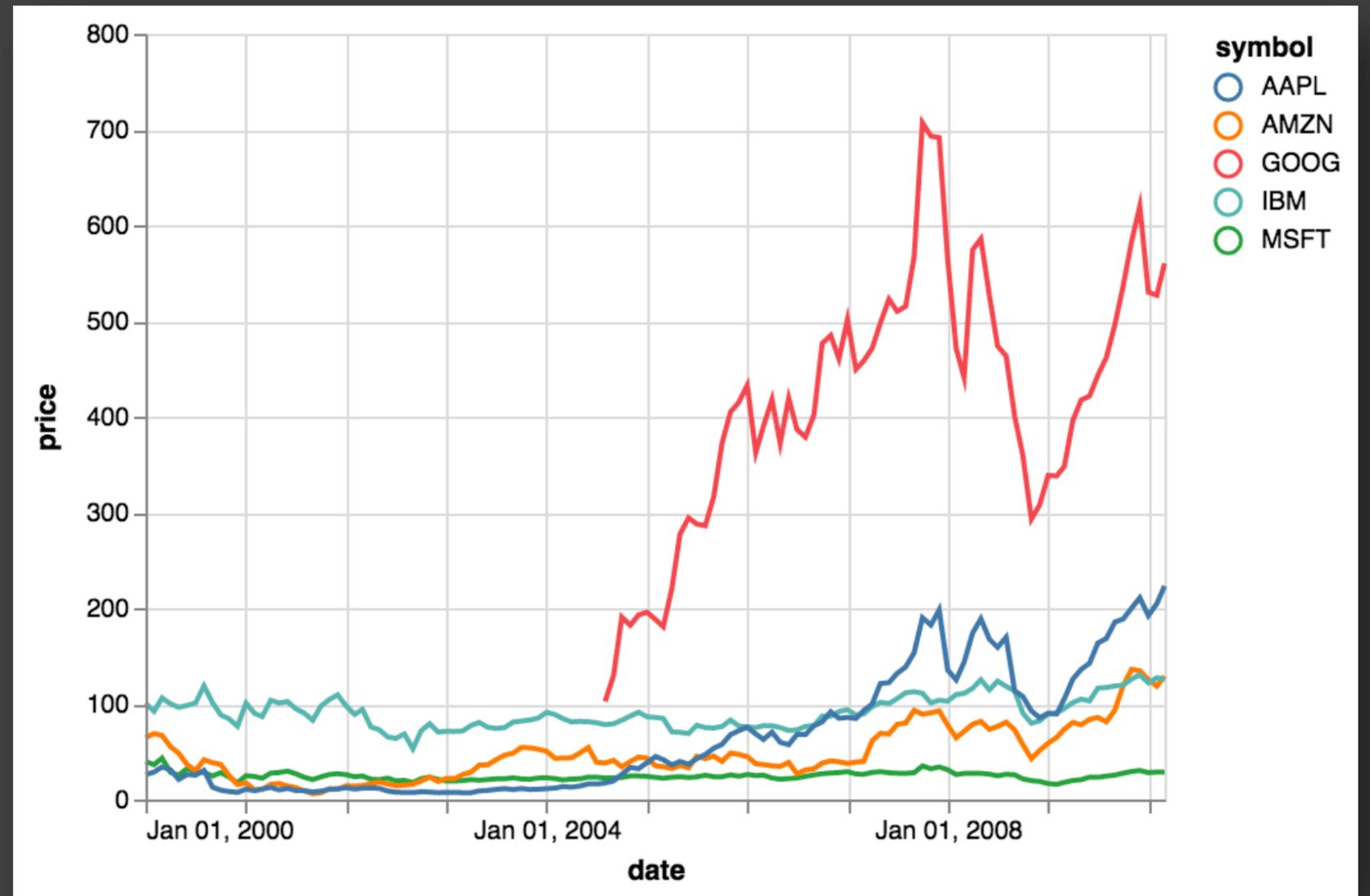
```
alt.Chart('stocks.csv')  
  .mark_line()  
  .encode(  
    x='date',  
    y='price',  
    color='symbol'  
  )
```



# Vega is a *Visualization Grammar*

## Vega-Lite JavaScript API

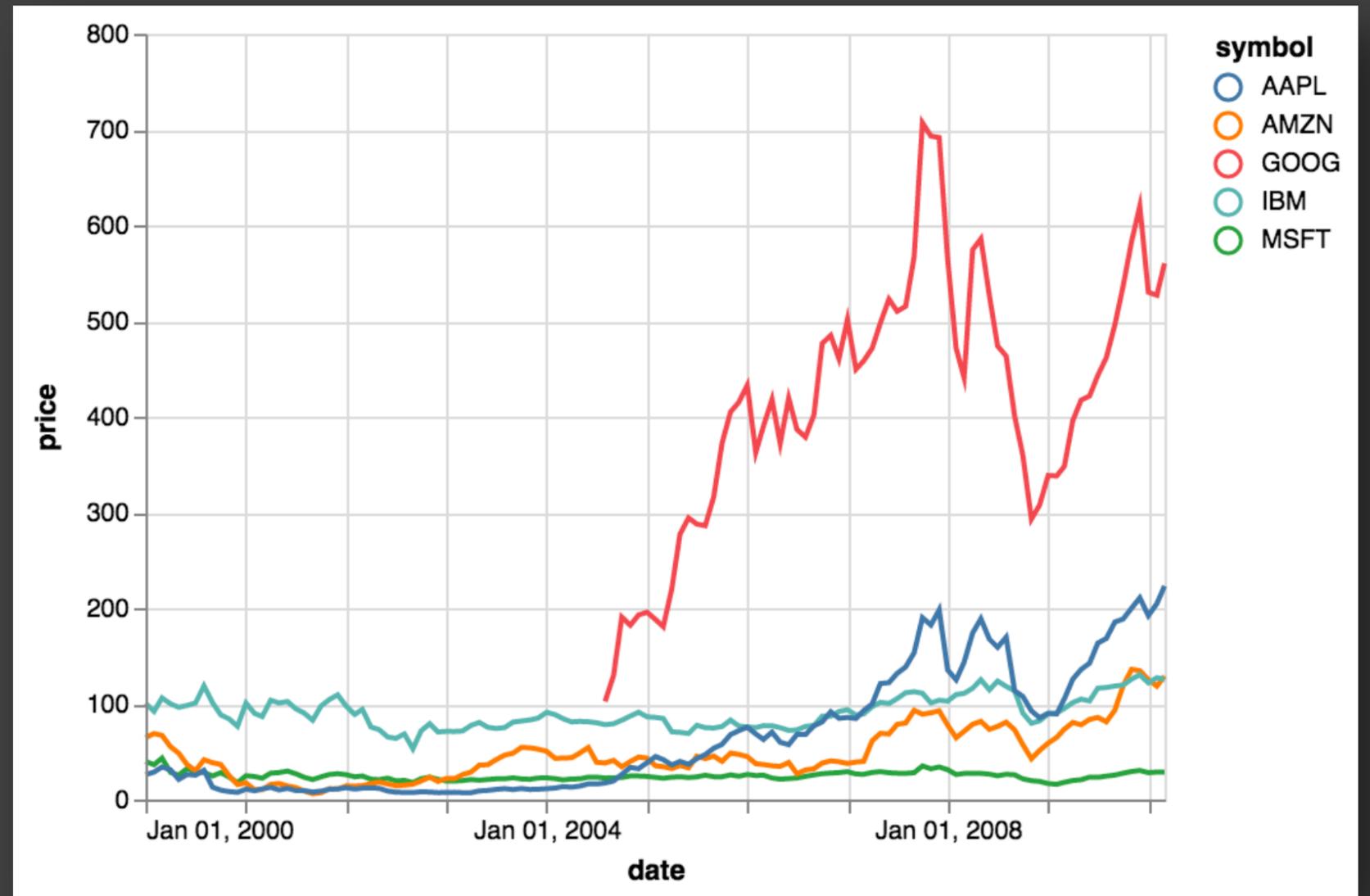
```
vl.data('stocks.csv')  
  .markLine()  
  .encode(  
    vl.x().fieldT('date'),  
    vl.y().fieldQ('price'),  
    vl.color().fieldN('symbol')  
  )
```



# Vega is a *Visualization Grammar*

## Vega-Lite JSON Specification

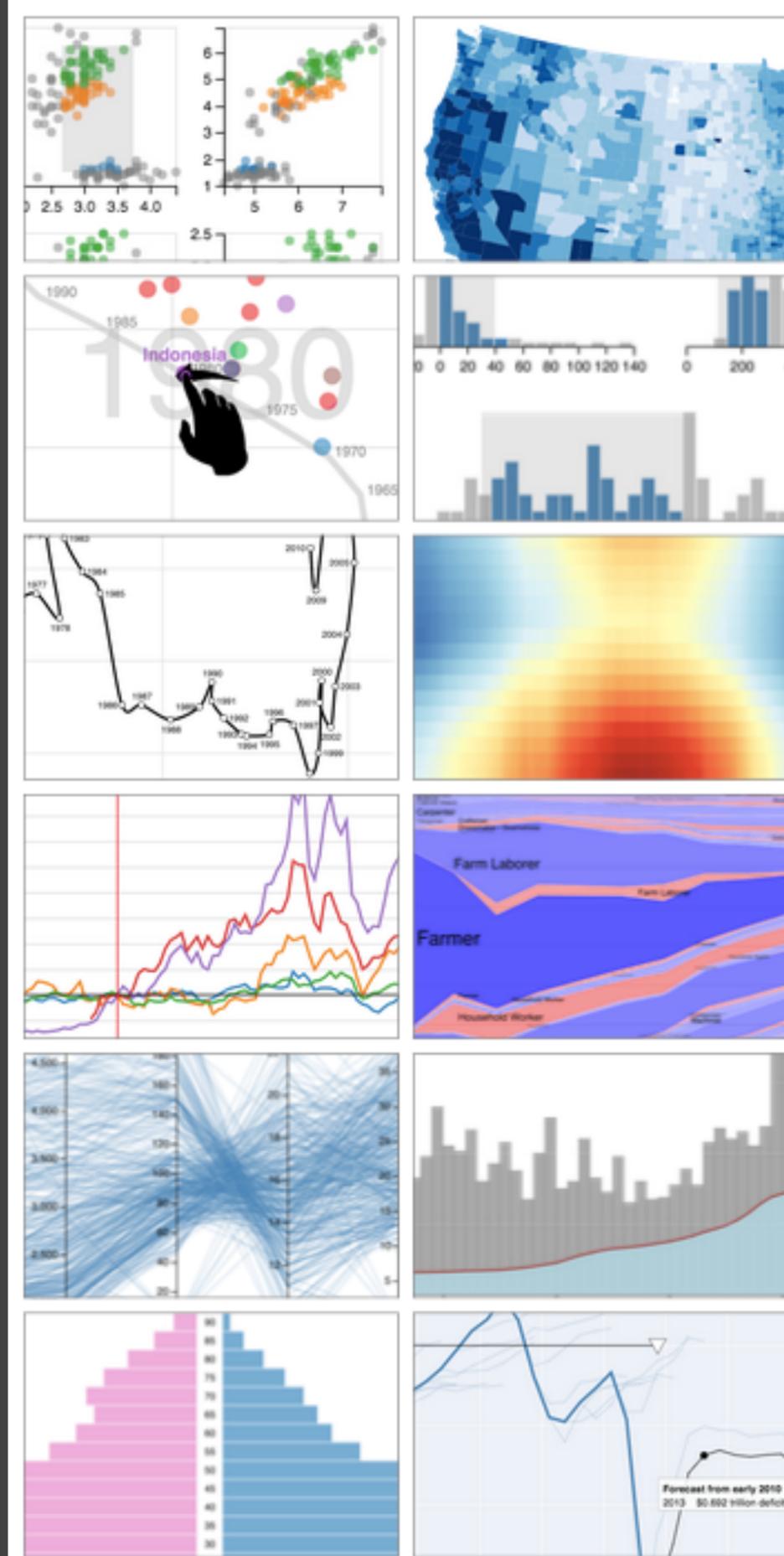
```
{
  data: {url: "stocks.csv"},
  mark: "line",
  encoding: {
    x: {
      type: "temporal",
      field: "date"
    },
    y: {
      type: "quantitative",
      field: "price"
    },
    color: {
      type: "nominal",
      field: "symbol"
    }
  }
}
```



# Vega is a *Visualization Grammar*

Similar in spirit to how SQL provides a language for expressing database queries, Vega is a high-level language for describing visualizations.

Vega provides a formal model for enumerating and reasoning about visualization designs.

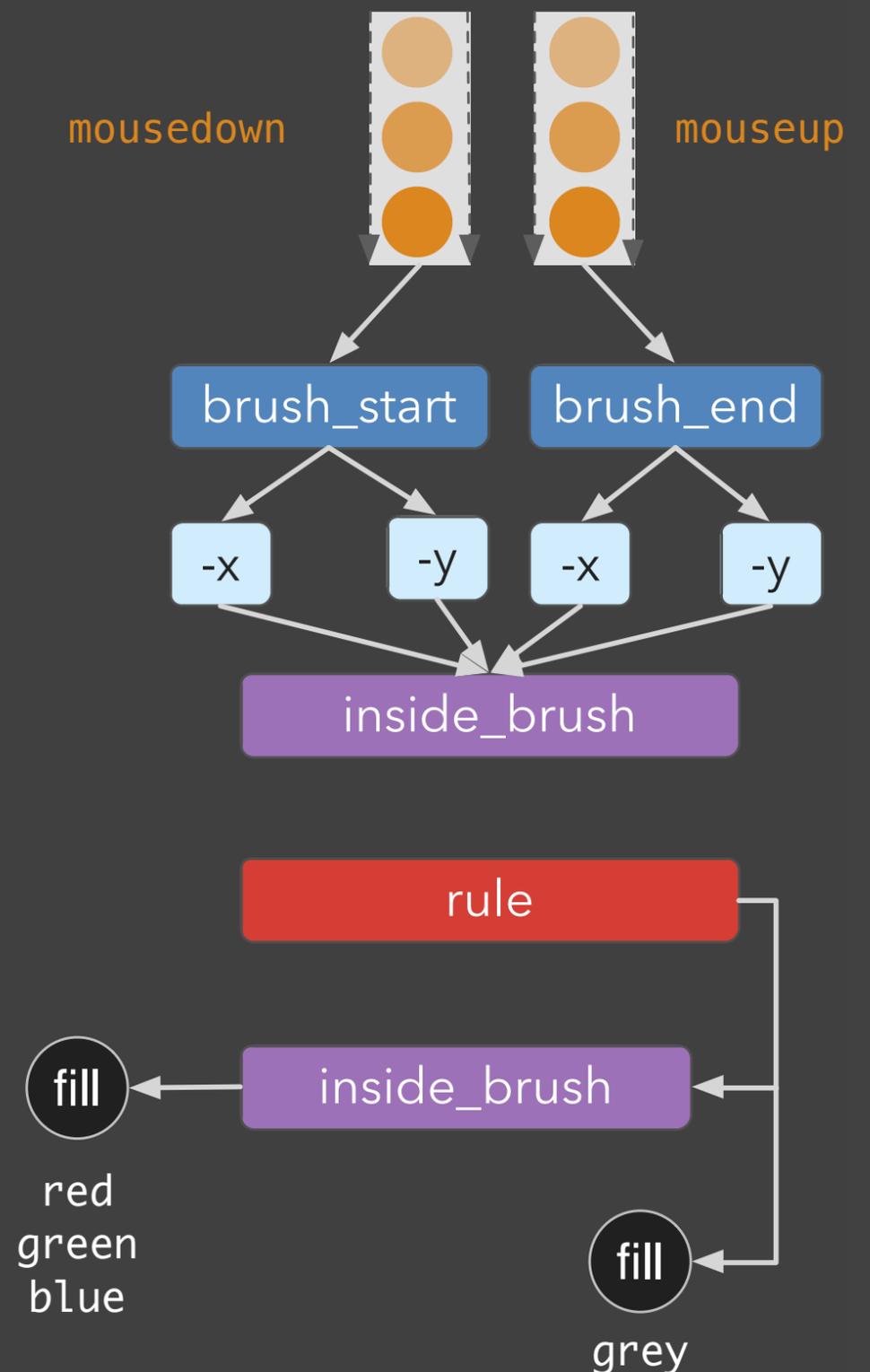


# Vega is a *Visualization Grammar*

Similar in spirit to how SQL provides a language for expressing database queries, Vega is a high-level language for describing visualizations.

Vega provides a formal model for enumerating and reasoning about visualization designs.

First comprehensive approach for declarative specification of interaction techniques.



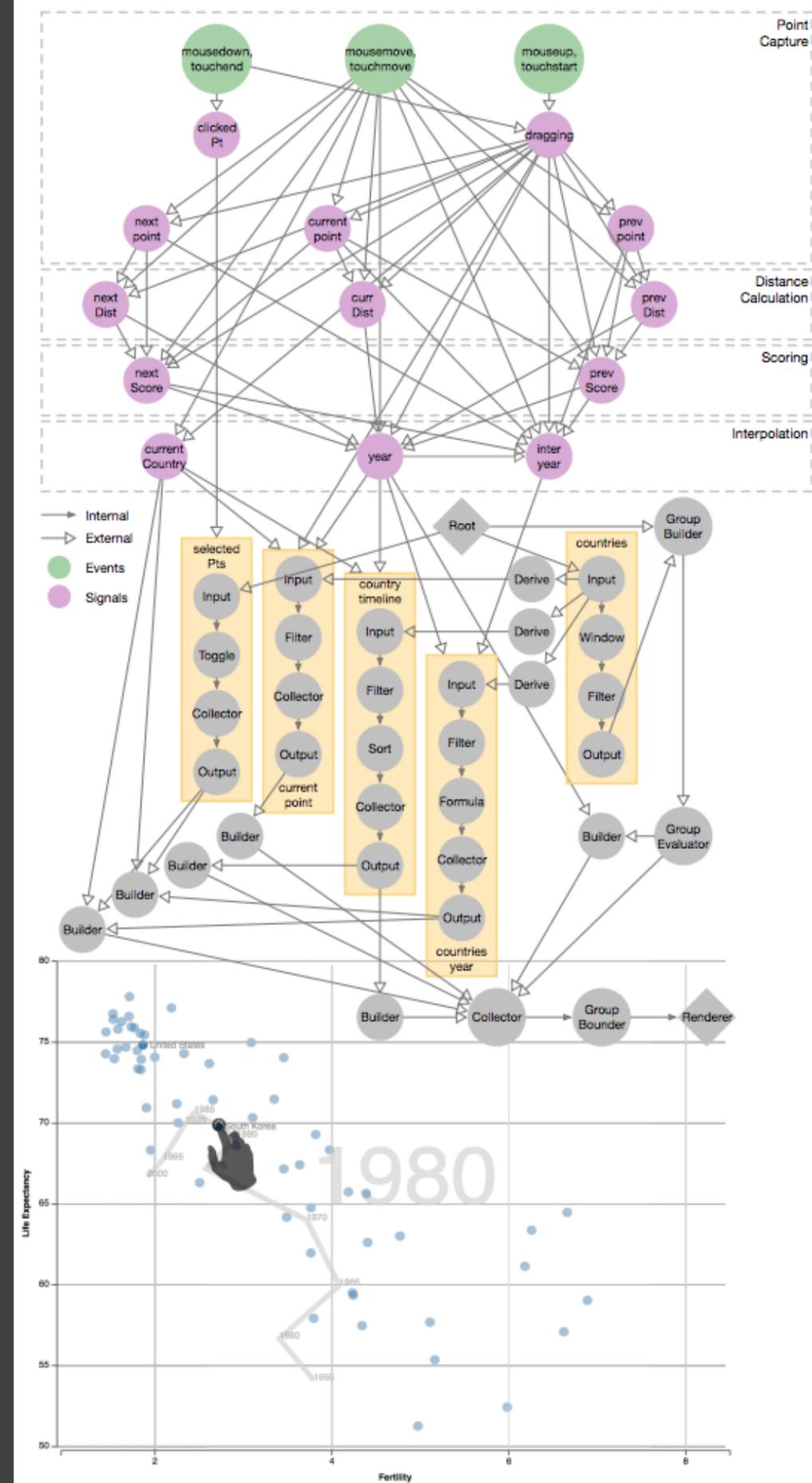
# Vega is a *Visualization Grammar*

Similar in spirit to how SQL provides a language for expressing database queries, Vega is a high-level language for describing visualizations.

Vega provides a formal model for enumerating and reasoning about visualization designs.

First comprehensive approach for declarative specification of interaction techniques.

Compiles JSON description to a reactive dataflow graph with efficient, incremental processing.



# Vega is a *Visualization Grammar*

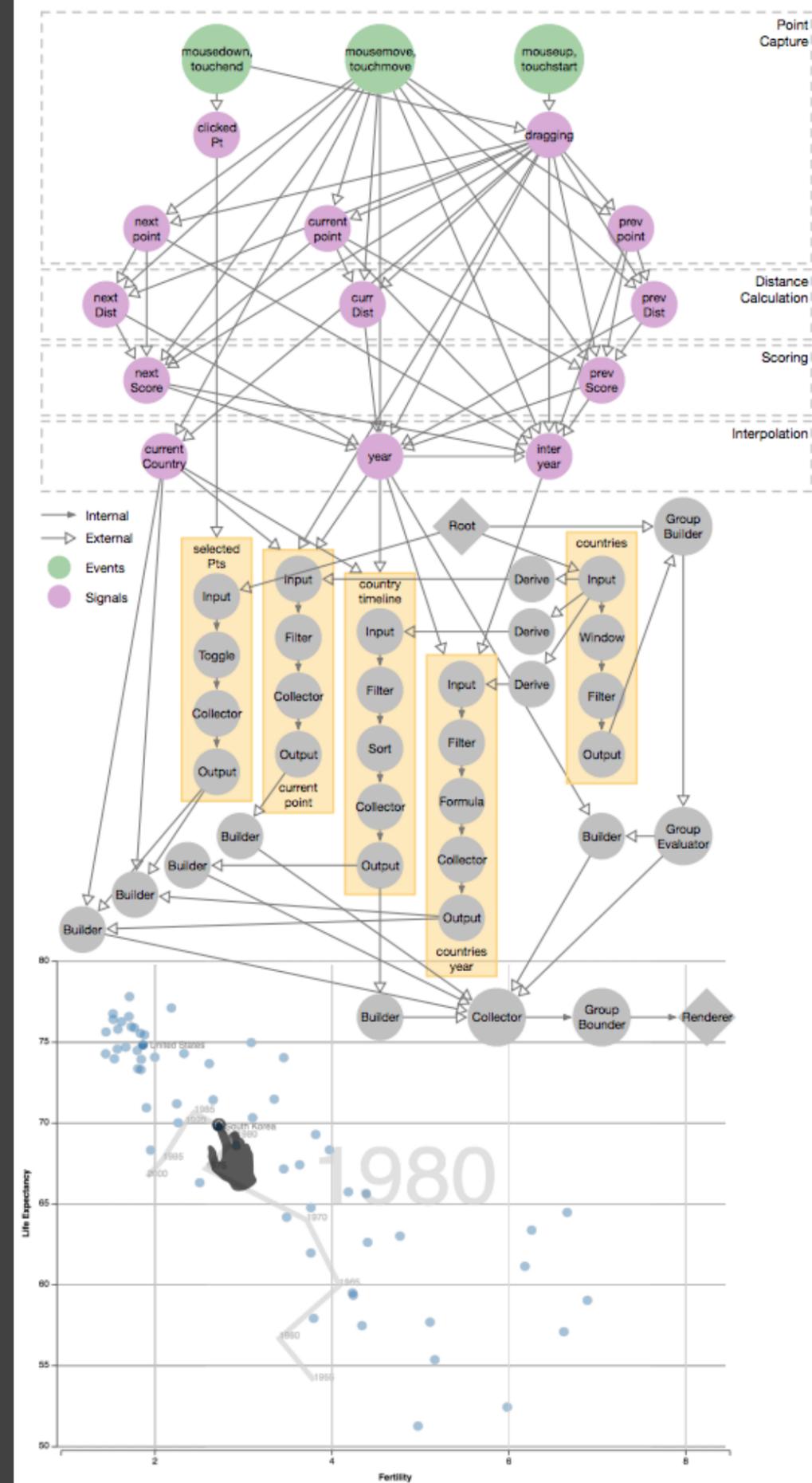
Similar in spirit to how SQL provides a language for expressing database queries, Vega is a high-level language for describing visualizations.

Vega provides a formal model for enumerating and reasoning about visualization designs.

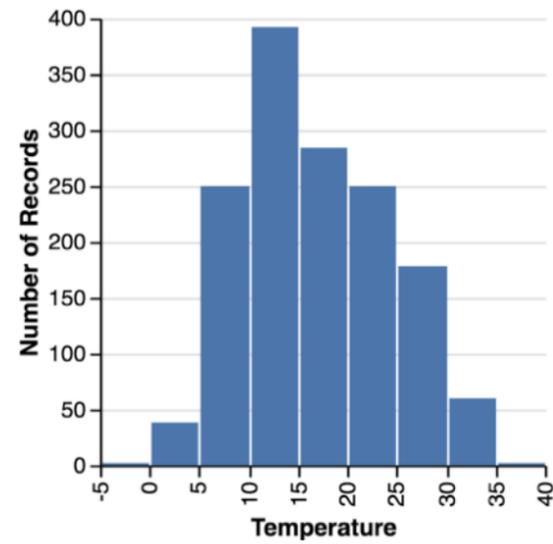
First comprehensive approach for declarative specification of interaction techniques.

Compiles JSON description to a reactive dataflow graph with efficient, incremental processing.

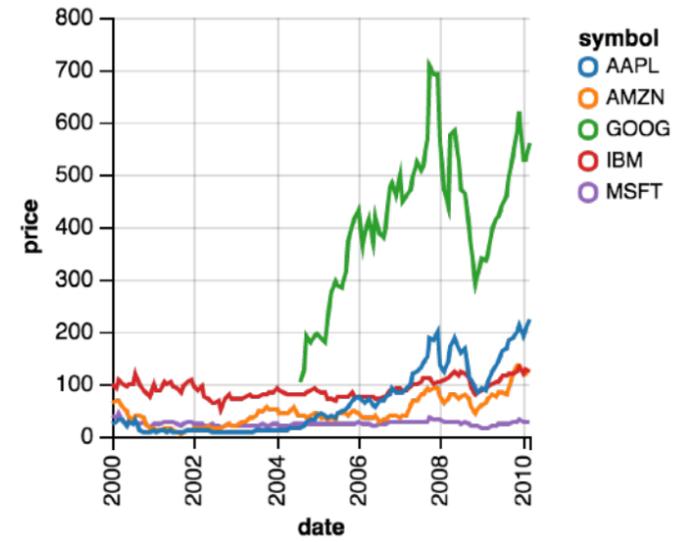
Vega's dataflow graph generates a scenegraph that is then rendered using Canvas or SVG.



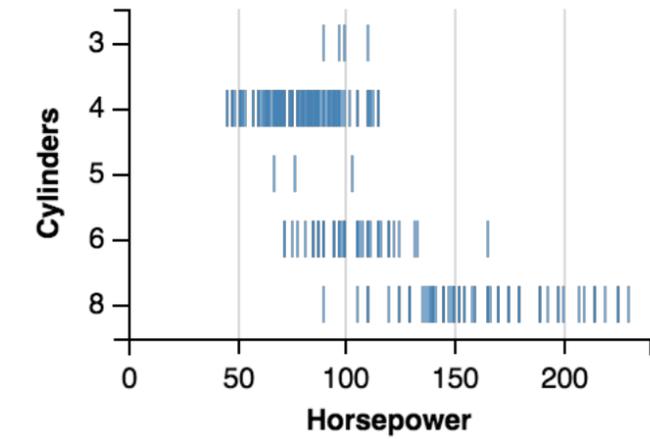
### Histogram



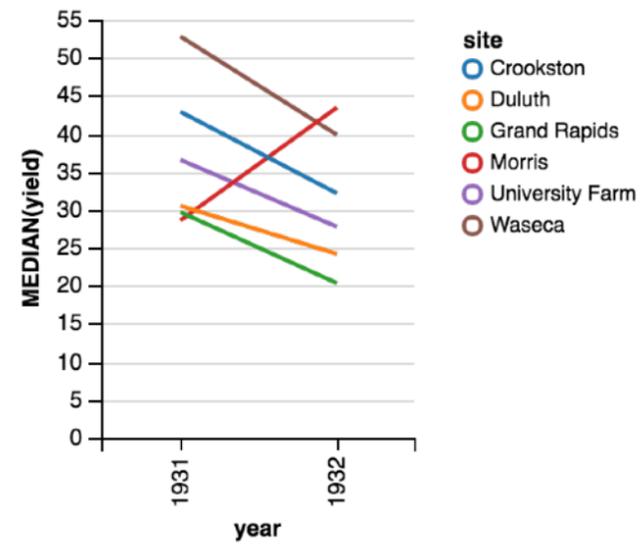
### Line Chart



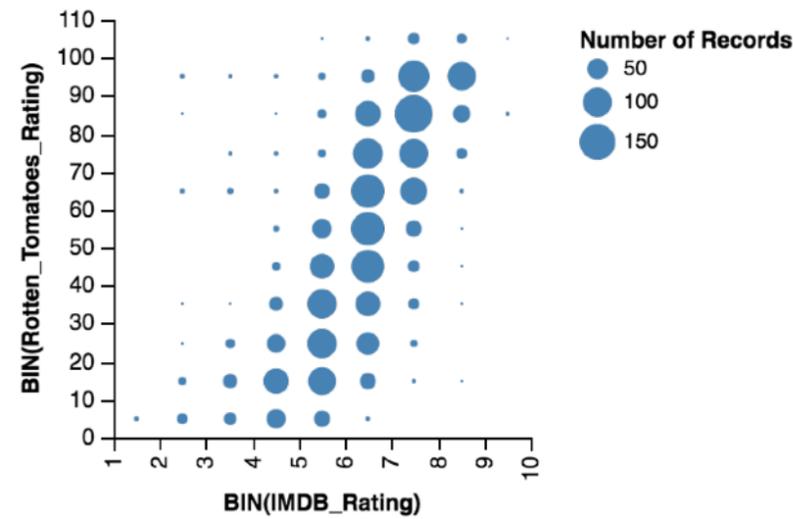
### Strip Plot



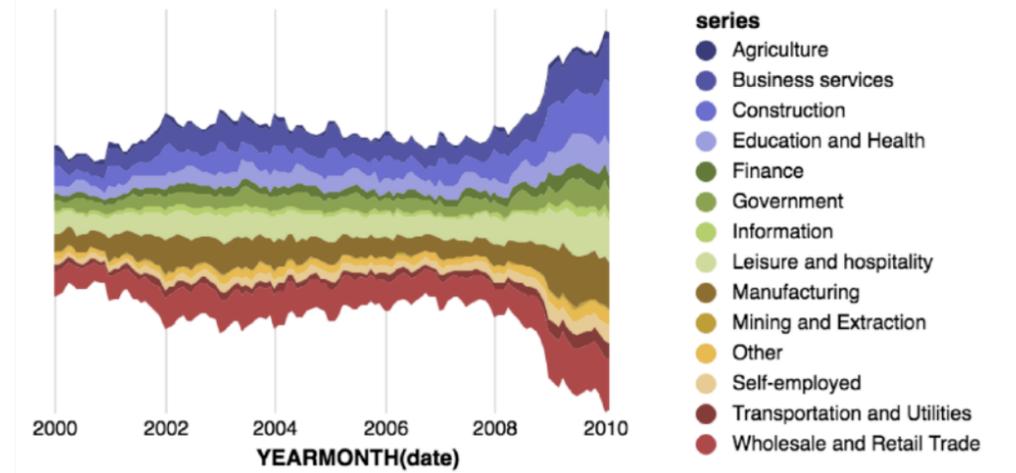
### Slope Graph



### Binned Scatter Plot



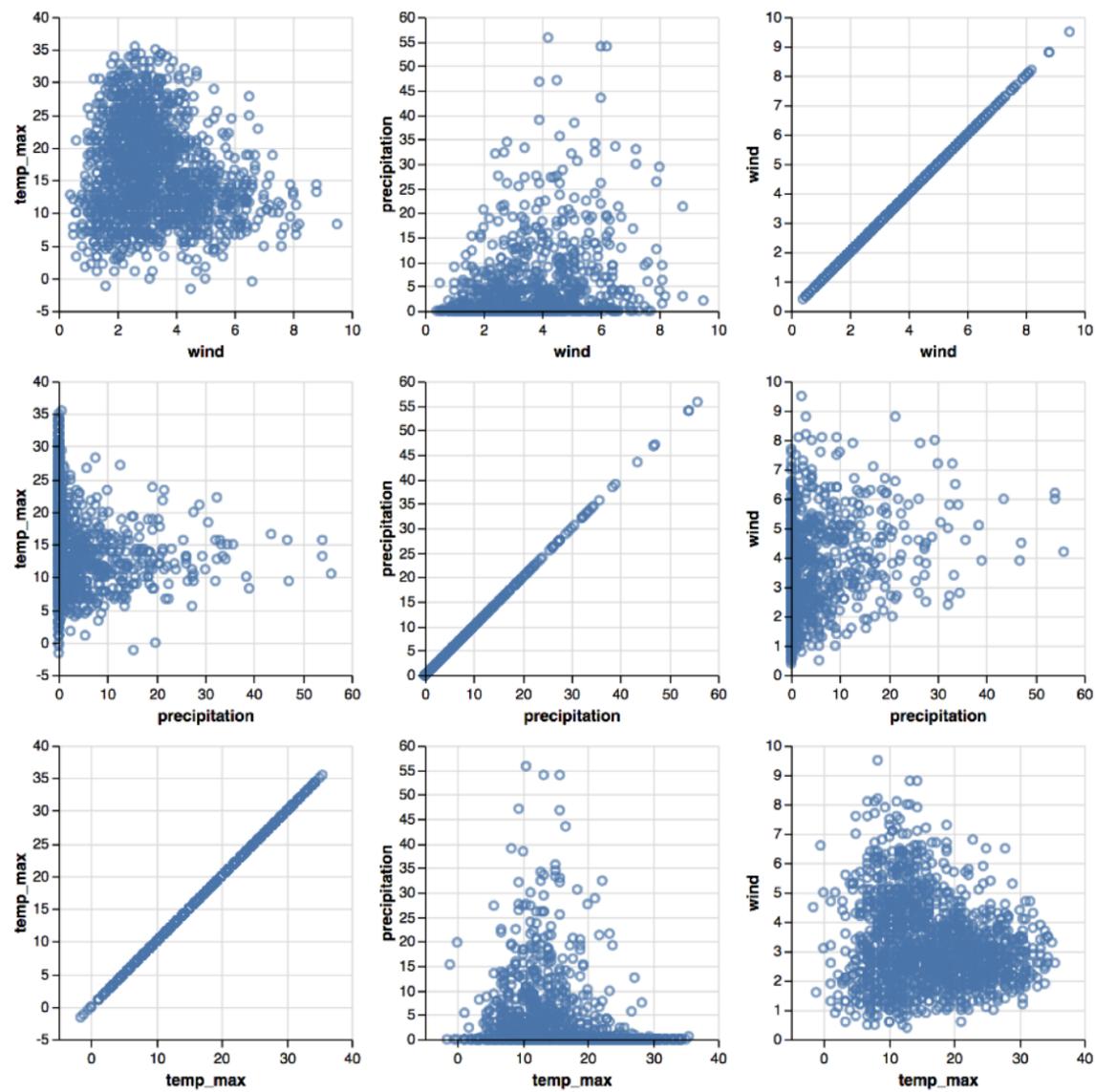
### Area Chart



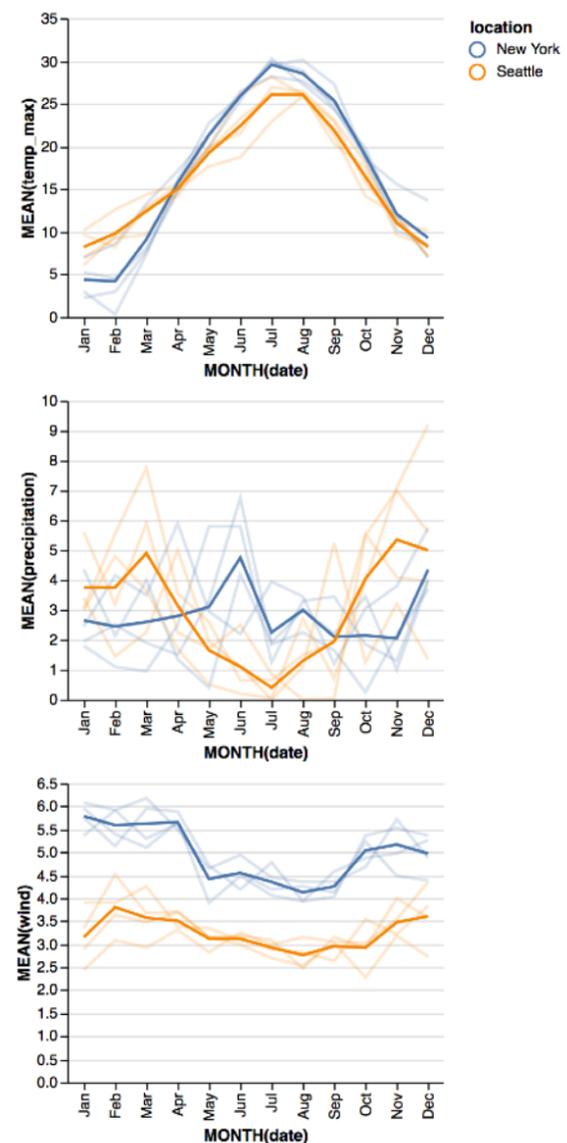
# A Grammar of Graphics

[Satyanarayan et al. 2017]

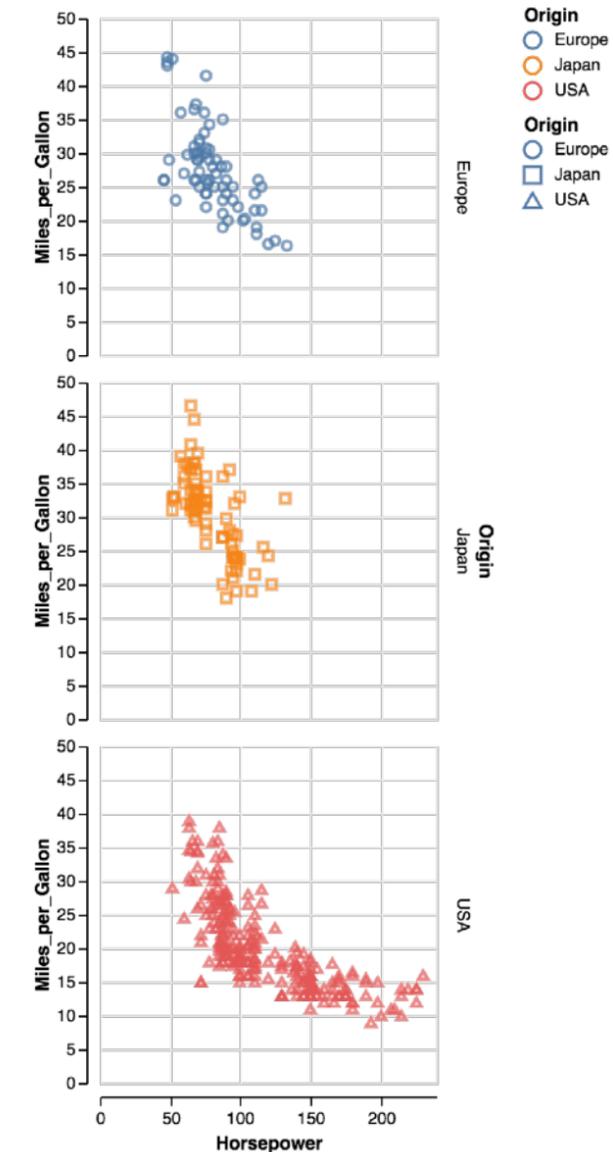
## Scatter Plot Matrix



## Concatenated & Layered View



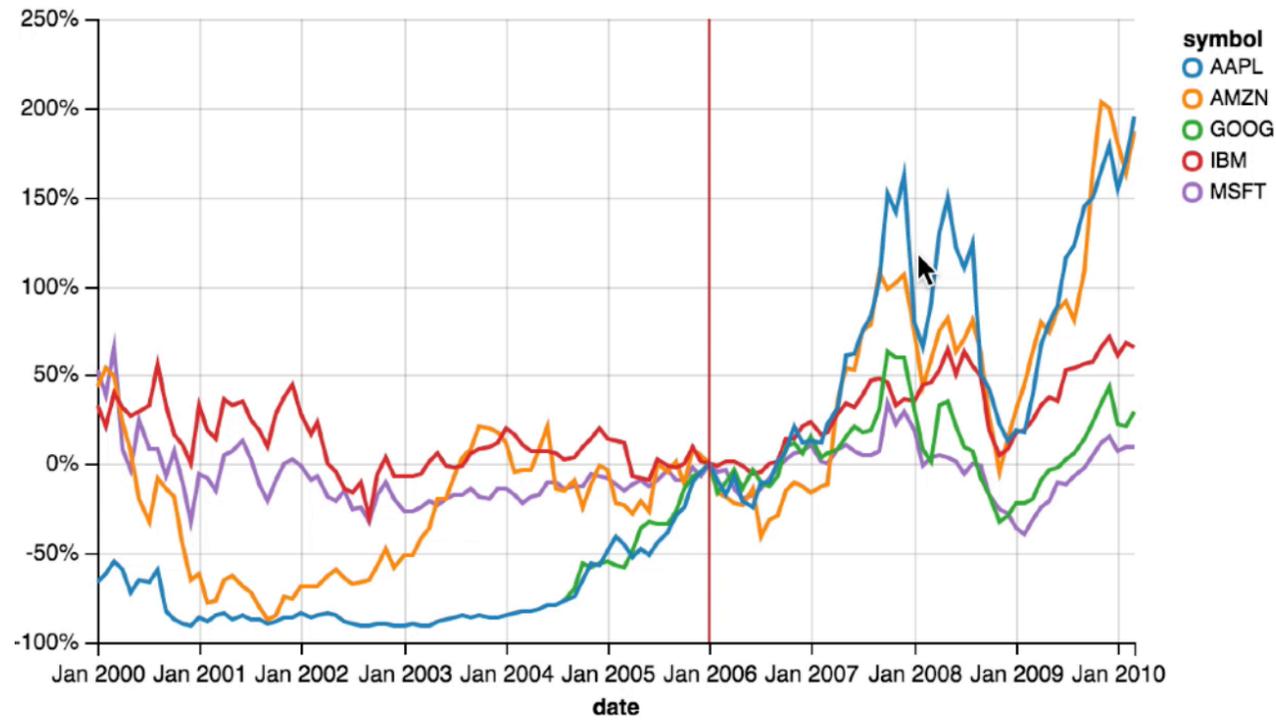
## Faceted View



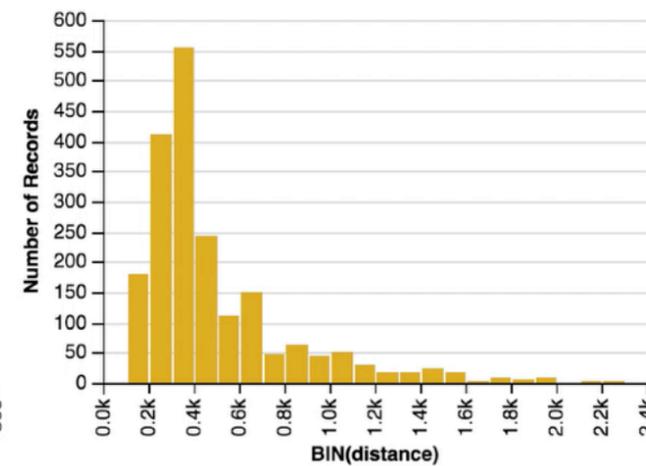
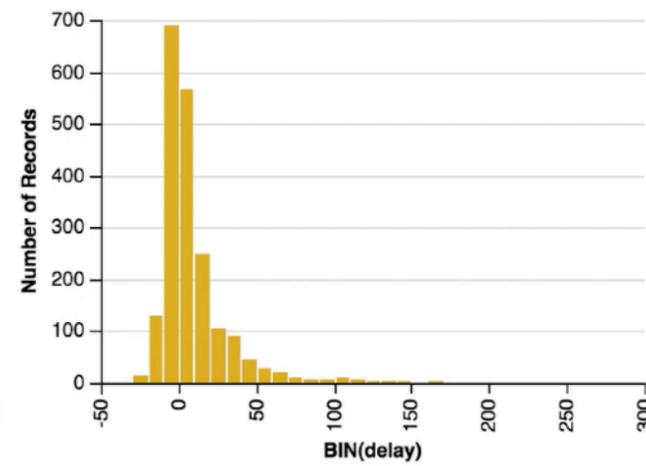
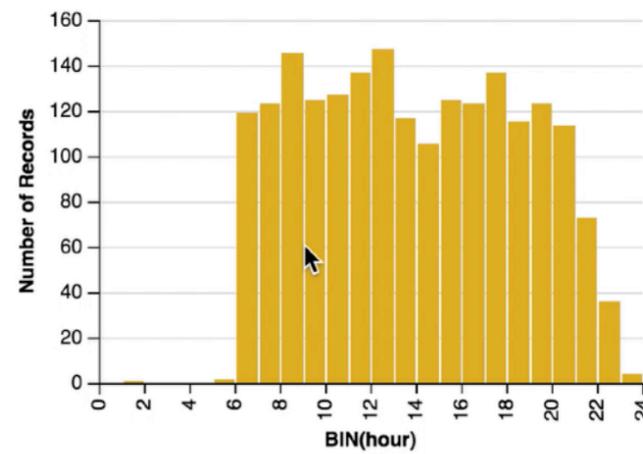
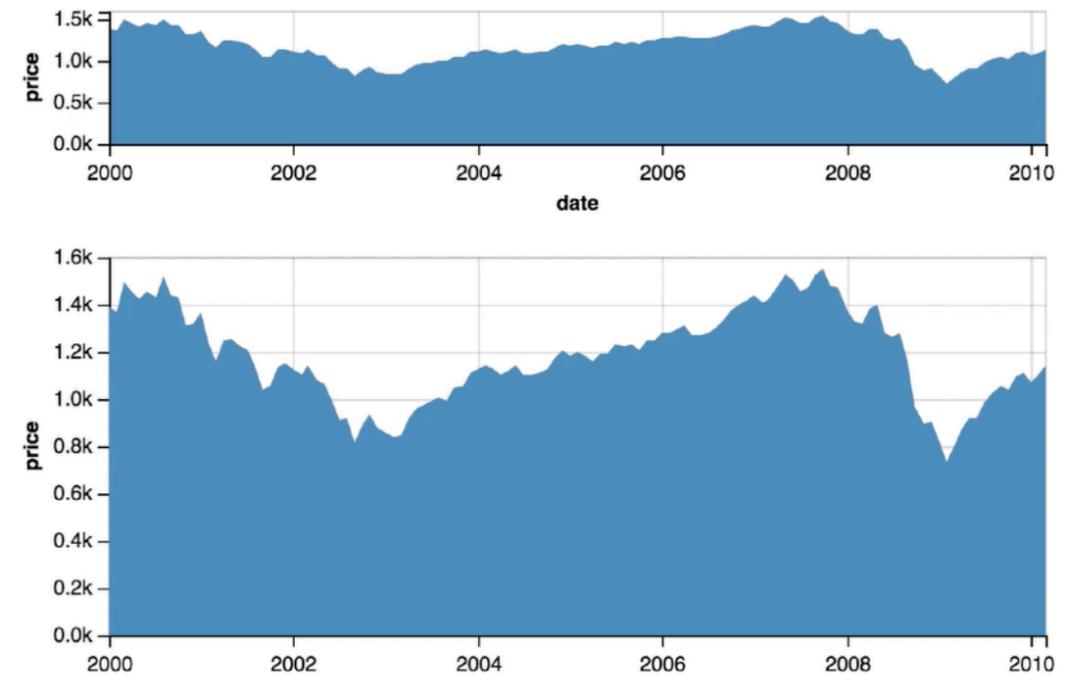
# A Grammar of Multi-View Graphics

[Satyanarayan et al. 2017]

### Indexed Chart



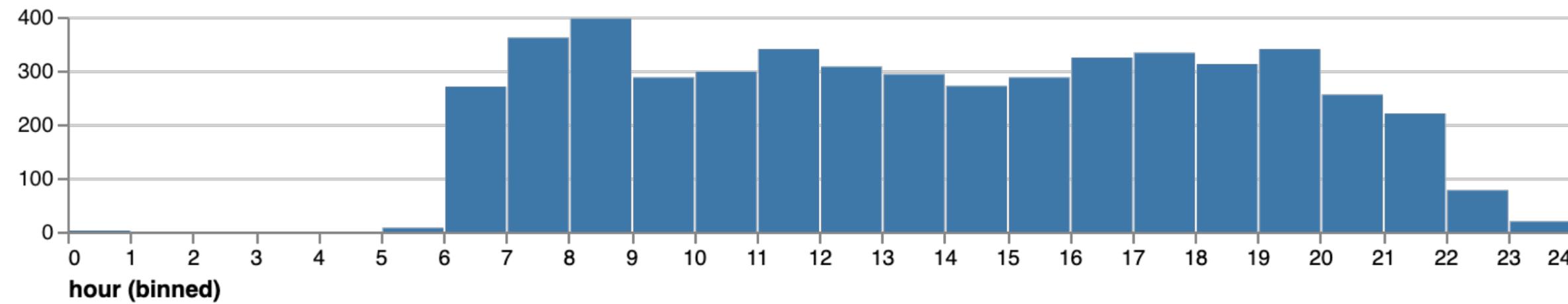
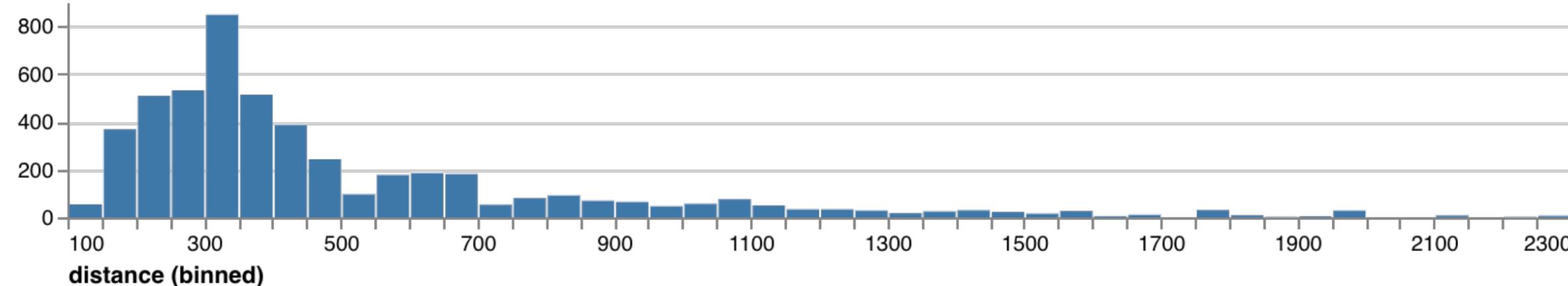
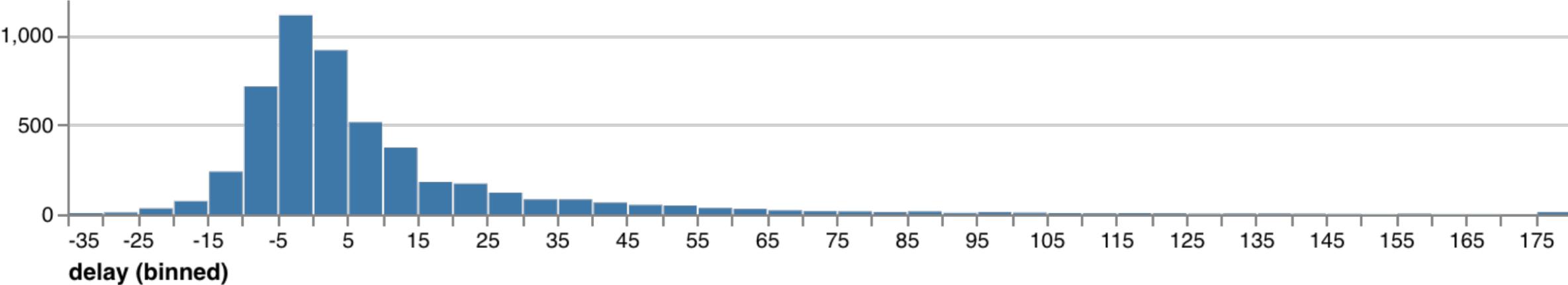
### Focus + Context



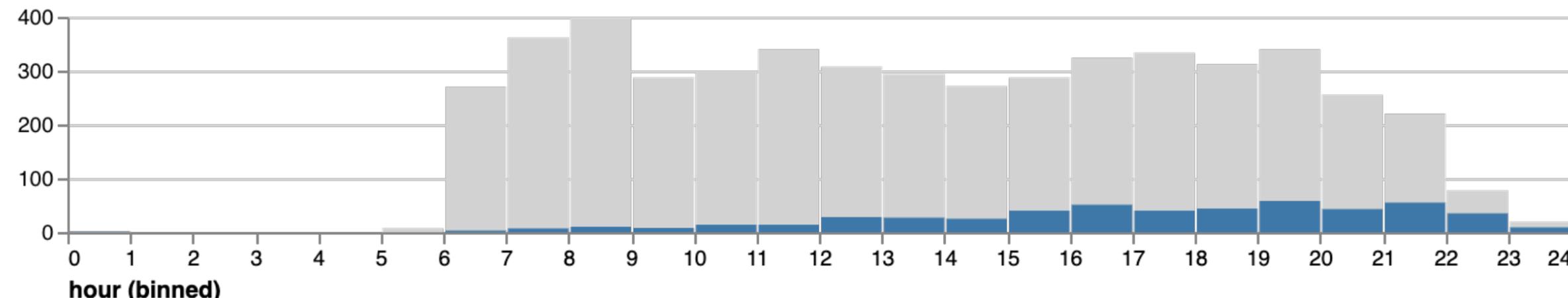
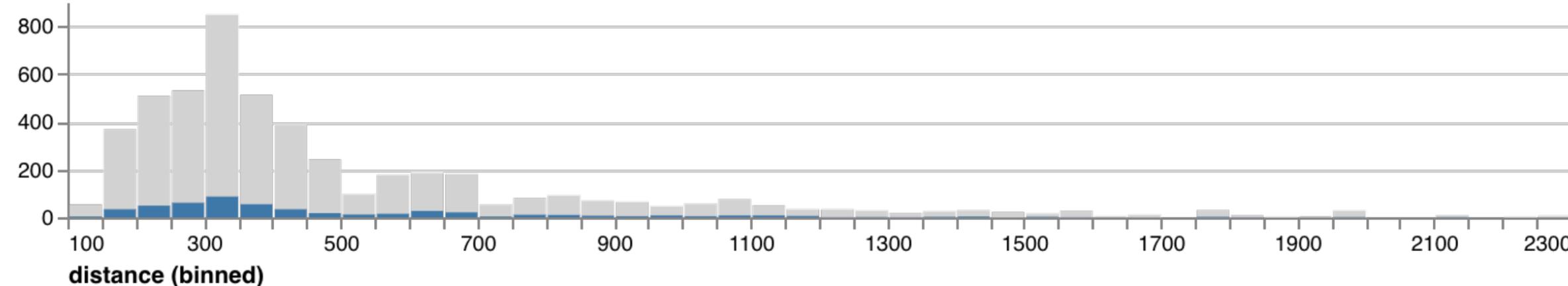
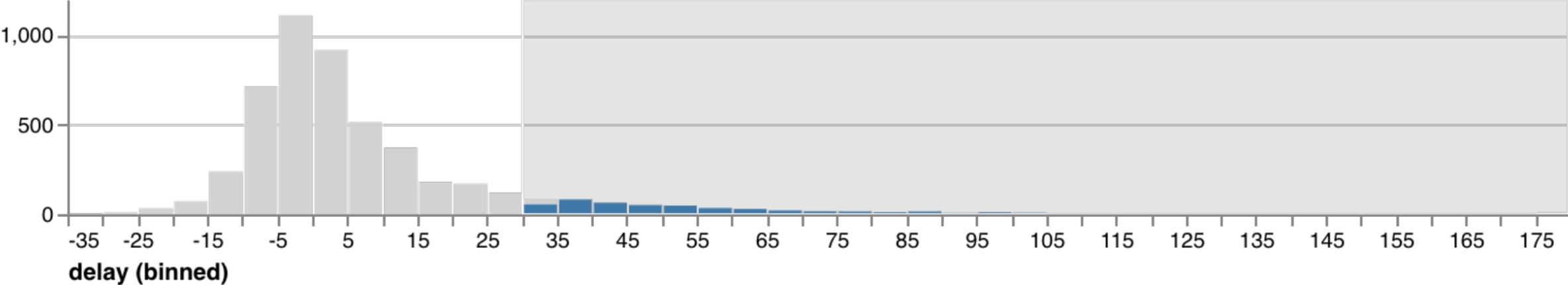
### Cross-Filtering

# A Grammar of **Interactive** Multi-View Graphics

# Cross-Filtering in Vega-Lite

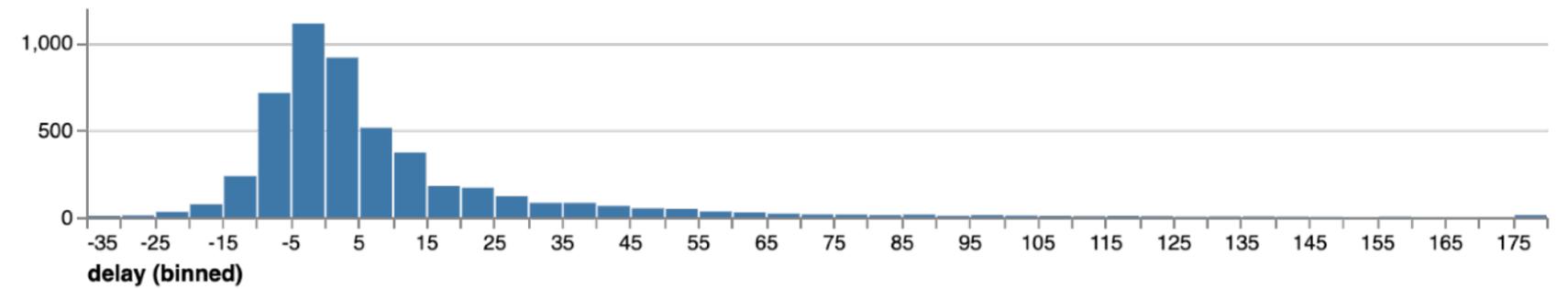


# Cross-Filtering in Vega-Lite



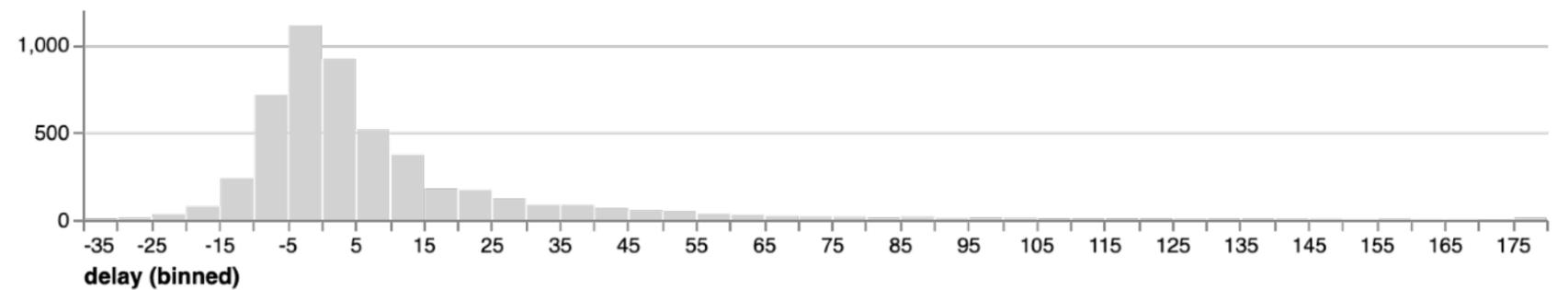
# Cross-Filtering in Vega-Lite

```
markBar().encode(  
  x().fieldQ('delay').bin(true),  
  y().count()  
)  
.data('data/flights.json')
```



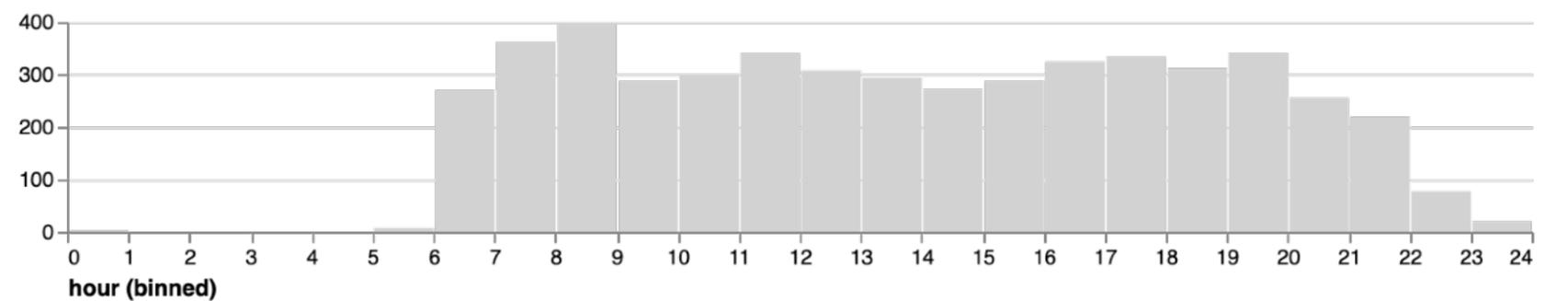
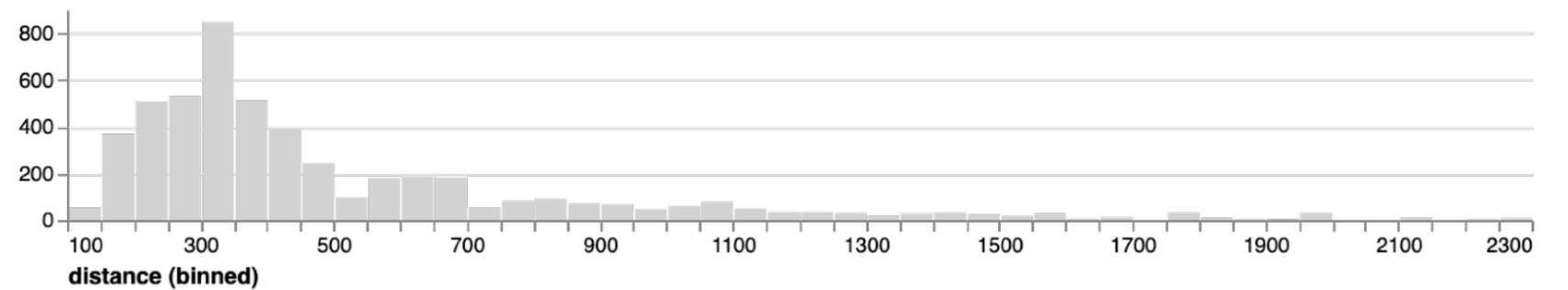
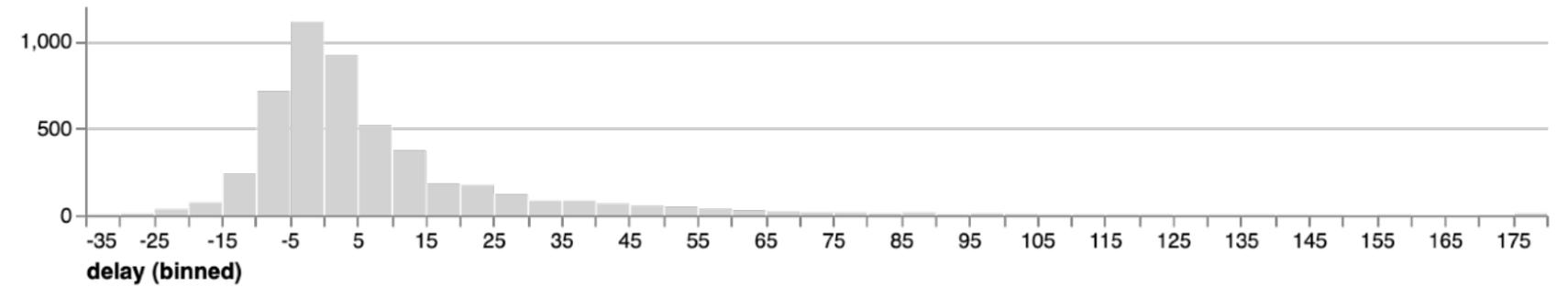
# Cross-Filtering in Vega-Lite

```
markBar().encode(  
  x().fieldQ('delay').bin(true),  
  y().count(),  
  color().value('lightgrey')  
) .data('data/flights.json')
```



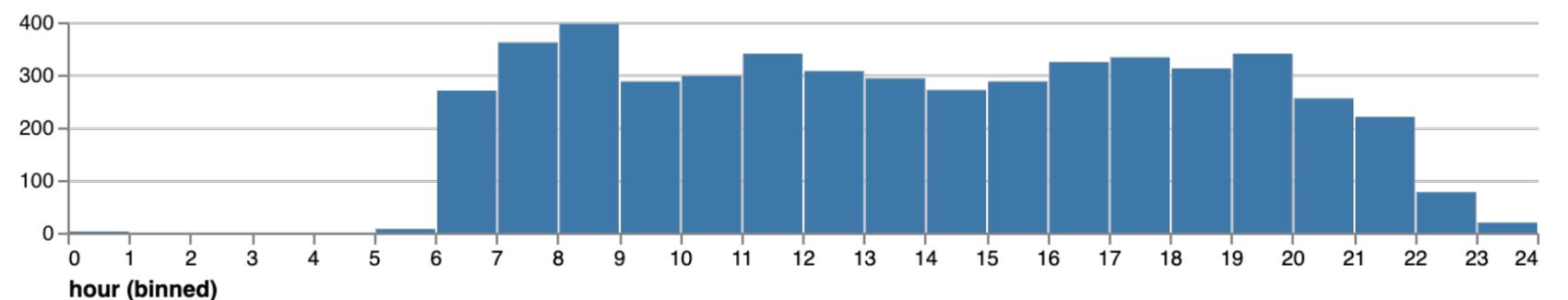
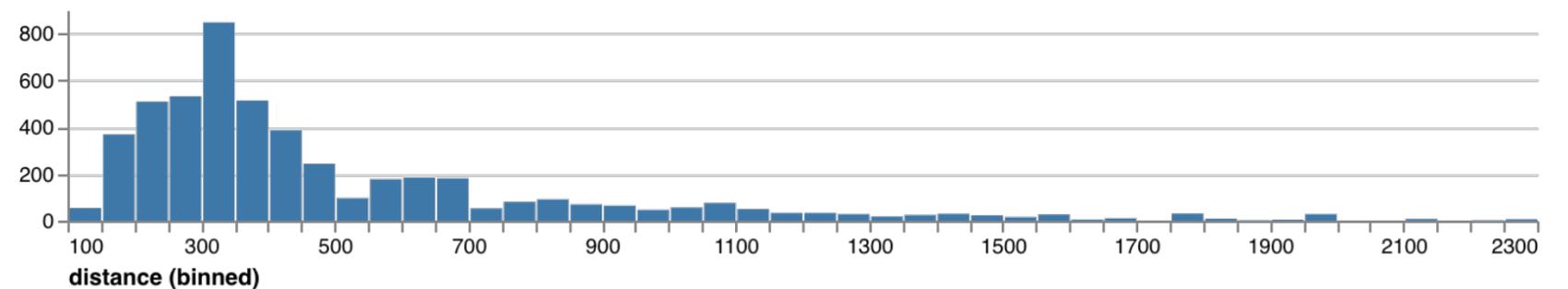
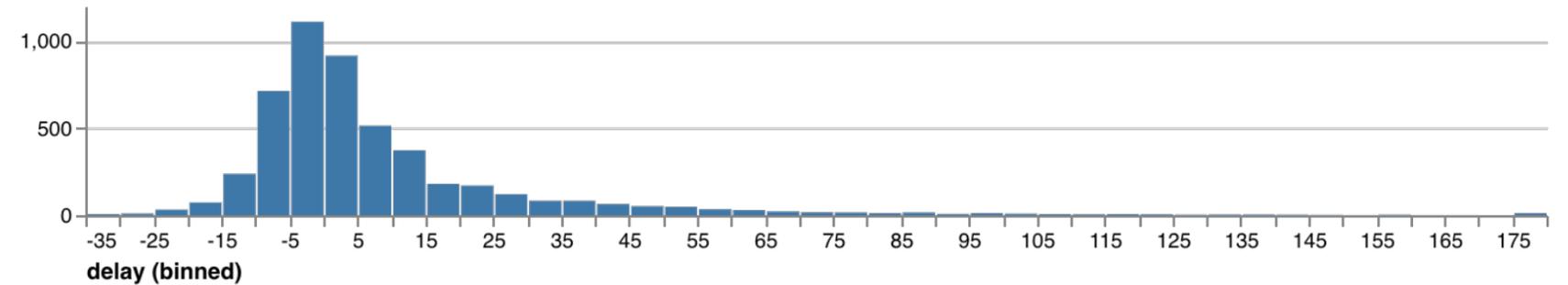
# Cross-Filtering in Vega-Lite

```
markBar().encode(  
  x().fieldQ(repeat('row').bin(true),  
  y().count(),  
  color().value('lightgrey'))  
)  
.repeat({row: ['delay', 'distance', 'hour']})  
.data('data/flights.json')
```



# Cross-Filtering in Vega-Lite

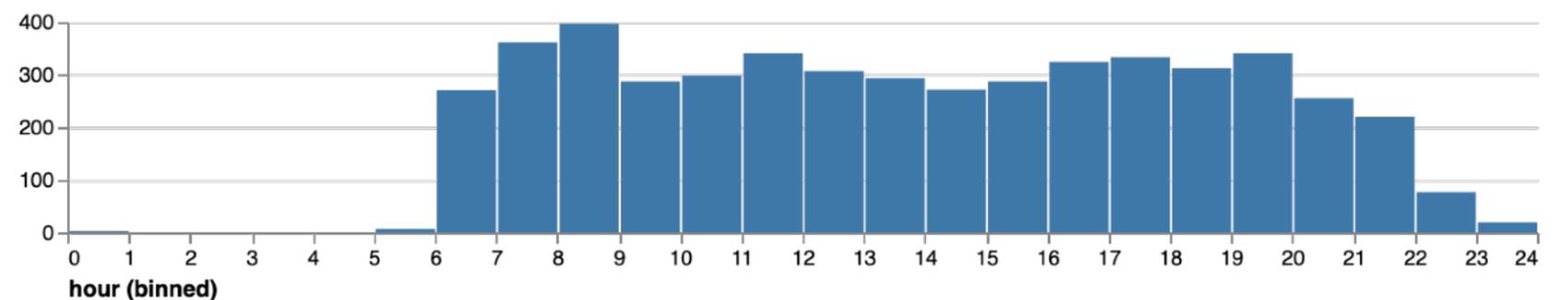
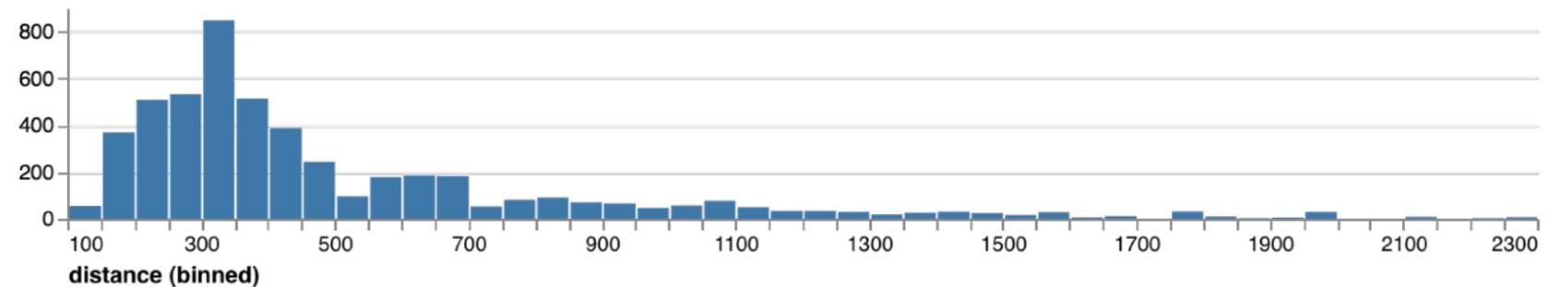
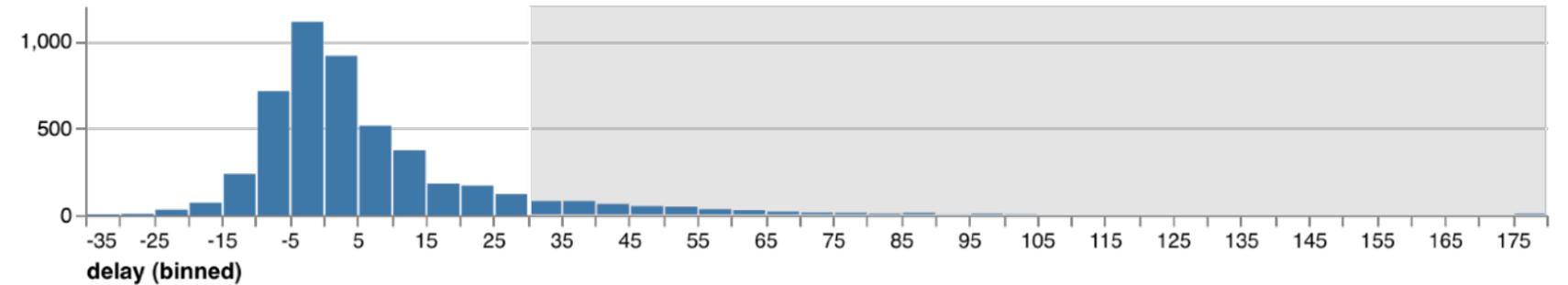
```
layer(  
  markBar().encode(  
    x().fieldQ(repeat('row')).bin(true),  
    y().count(),  
    color().value('lightgrey')  
  ),  
  markBar().encode(  
    x().fieldQ(repeat('row')).bin(true),  
    y().count()  
  )  
).repeat({row: ['delay', 'distance', 'hour']})  
.data('data/flights.json')
```



# Cross-Filtering in Vega-Lite

```
brush = selectInterval().encodings('x')

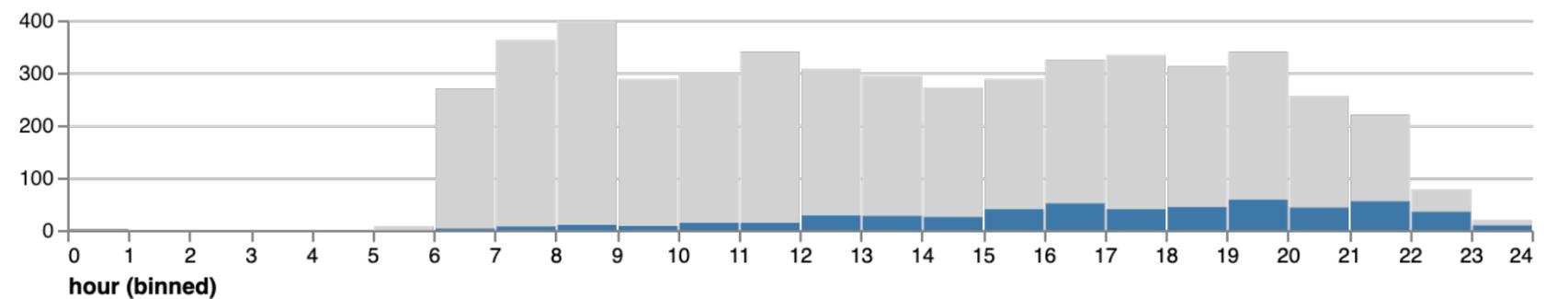
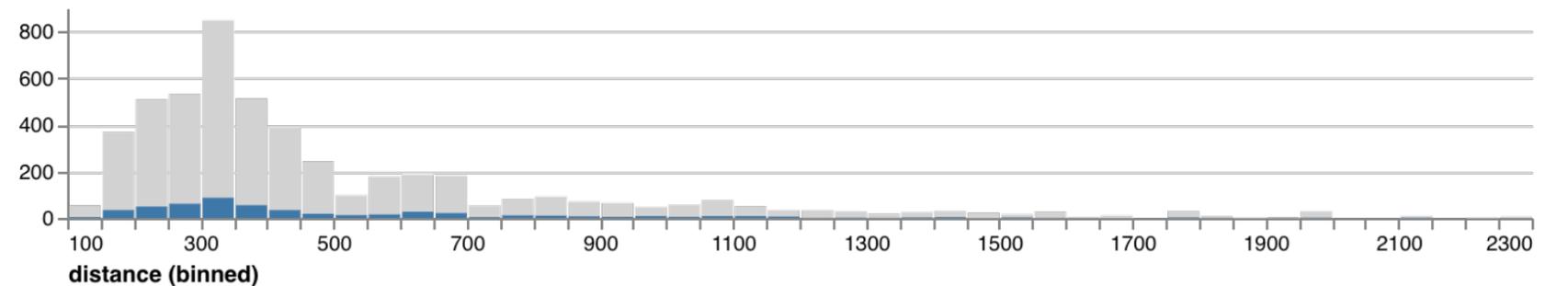
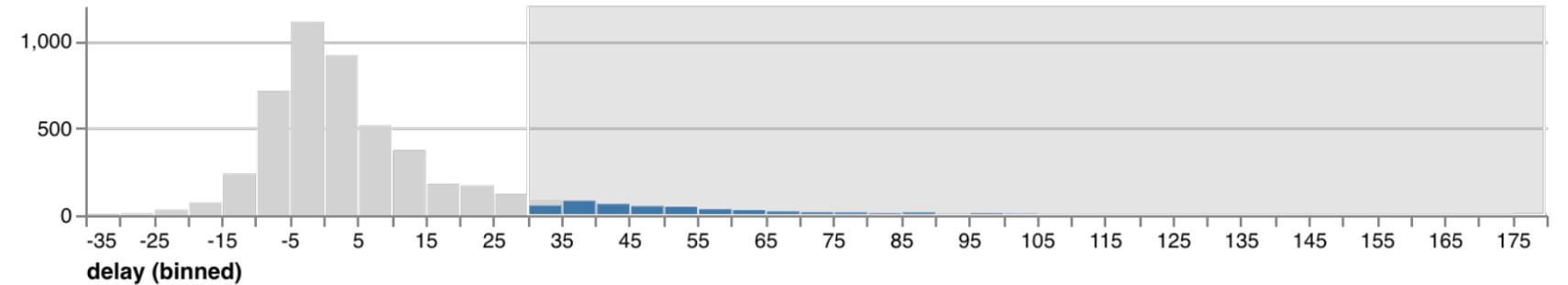
layer(
  markBar().encode(
    x().fieldQ(repeat('row')).bin(true),
    y().count(),
    color().value('lightgrey')
  ).select(brush),
  markBar().encode(
    x().fieldQ(repeat('row')).bin(true),
    y().count()
  )
)
.repeat({row: ['delay', 'distance', 'hour']})
.data('data/flights.json')
```



# Cross-Filtering in Vega-Lite

```
brush = selectInterval.encodings('x')

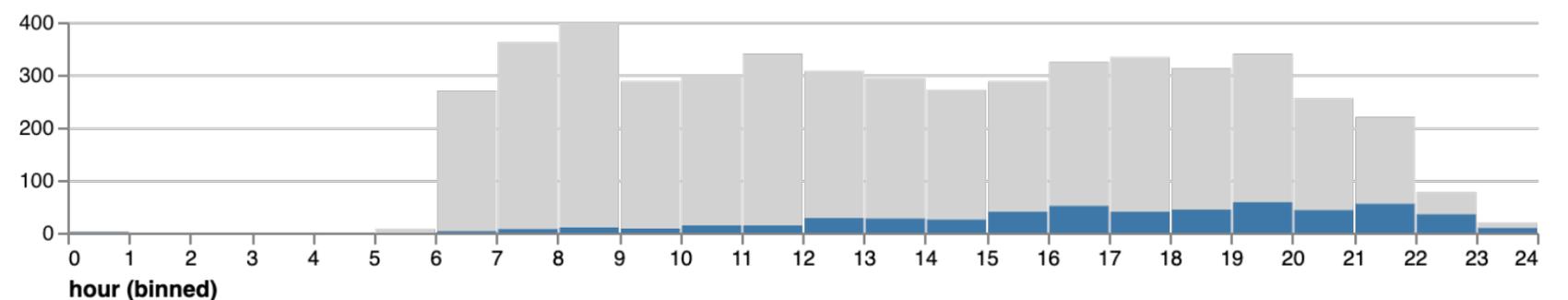
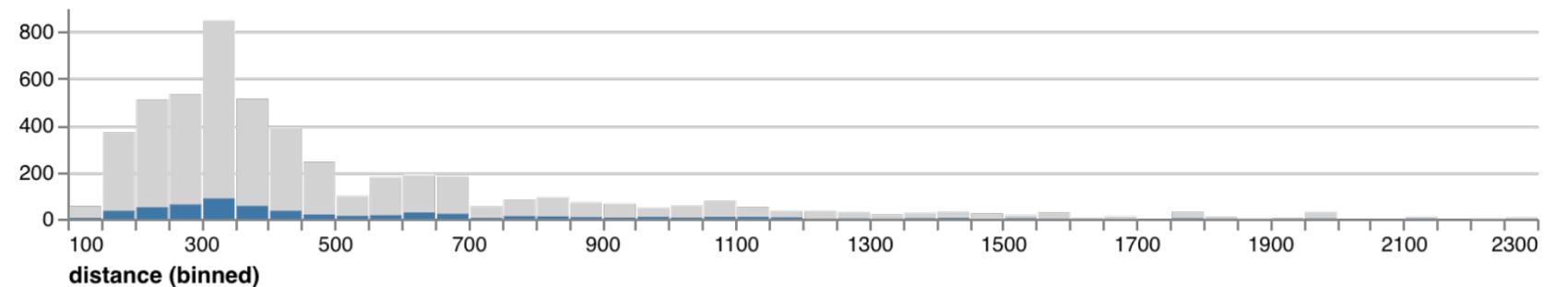
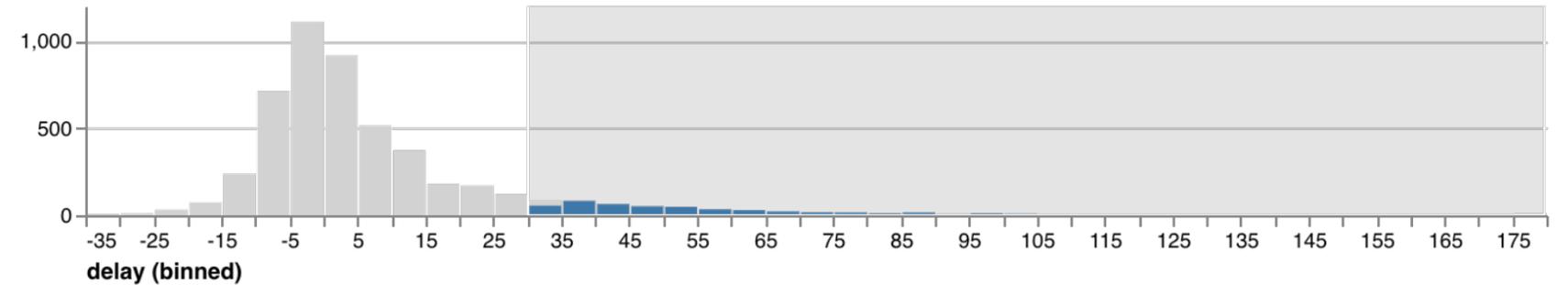
layer(
  markBar().encode(
    x().fieldQ(repeat('row')).bin(true),
    y().count(),
    color().value('lightgrey')
  ).select(brush),
  markBar().encode(
    x().fieldQ(repeat('row')).bin(true),
    y().count()
  ).transform(filter(brush))
)
.repeat({row: ['delay', 'distance', 'hour']})
.data('data/flights.json')
```



# Cross-Filtering in Vega-Lite

```
brush = selectInterval.encodings('x')

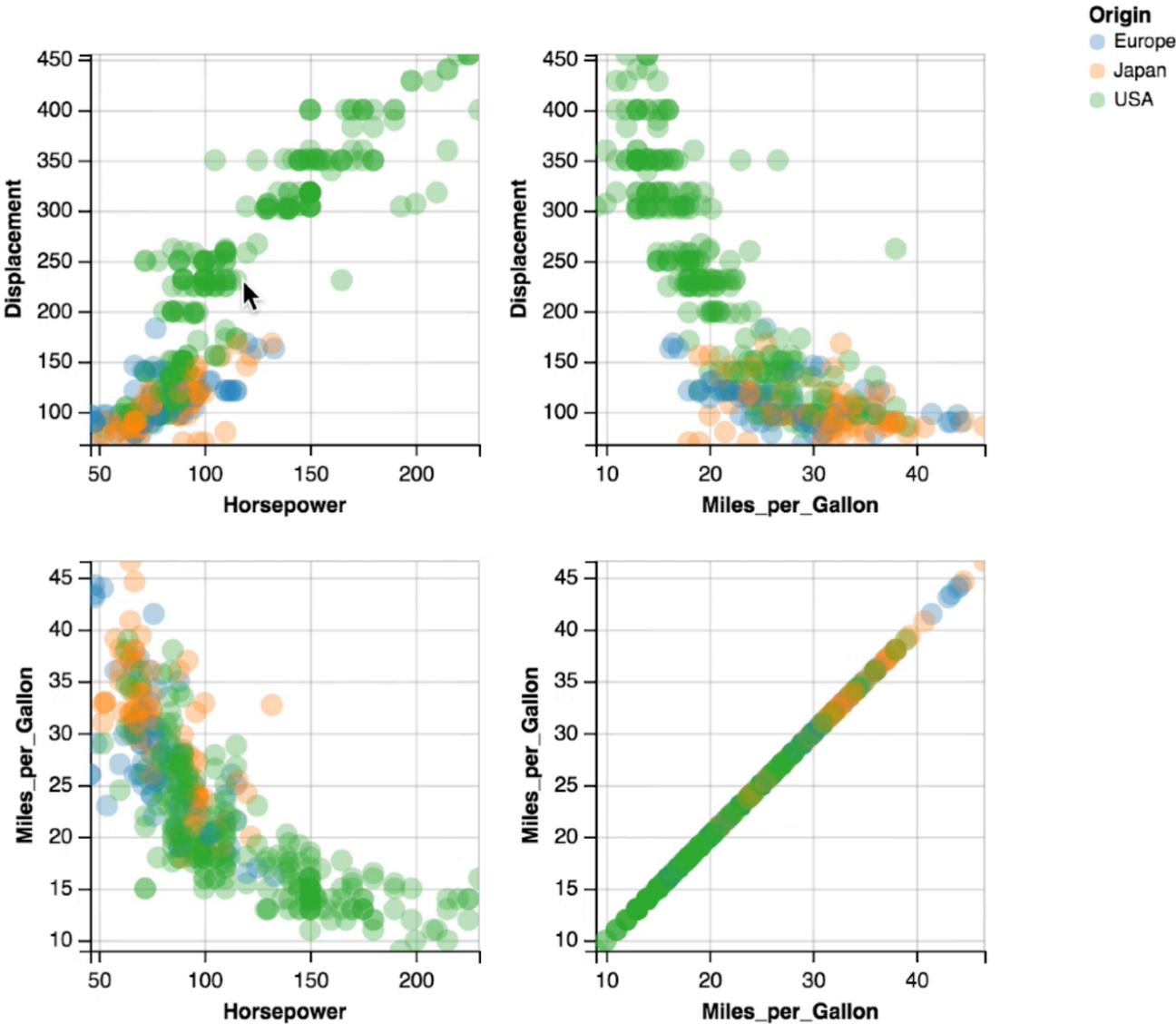
layer(
  markBar().encode(
    x().fieldQ(repeat('row')).bin(true),
    y().count(),
    color().value('lightgrey')
  ).select(brush),
  markBar().encode(
    x().fieldQ(repeat('row')).bin(true),
    y().count()
  ).transform(filter(brush))
)
.repeat({row: ['delay', 'distance', 'hour']})
.data('data/flights.json')
```



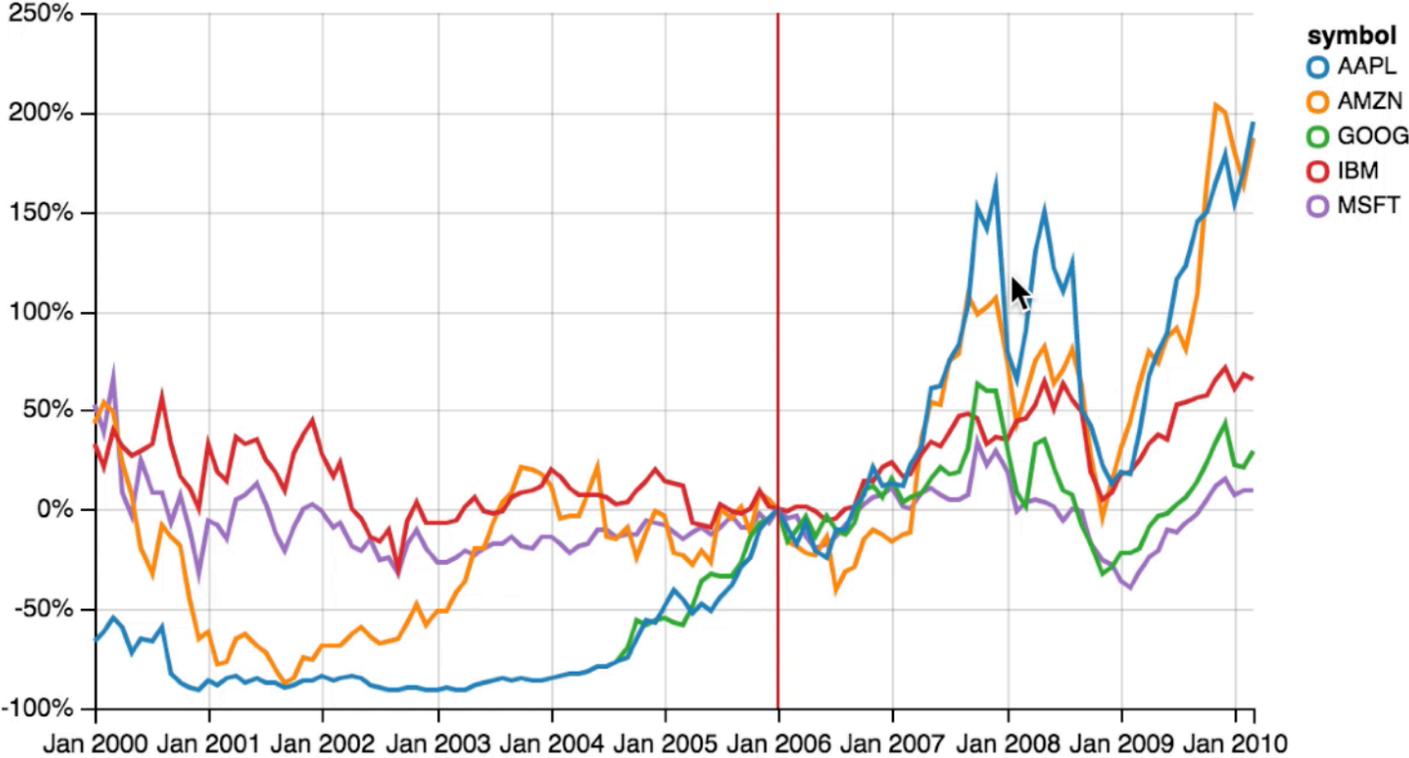
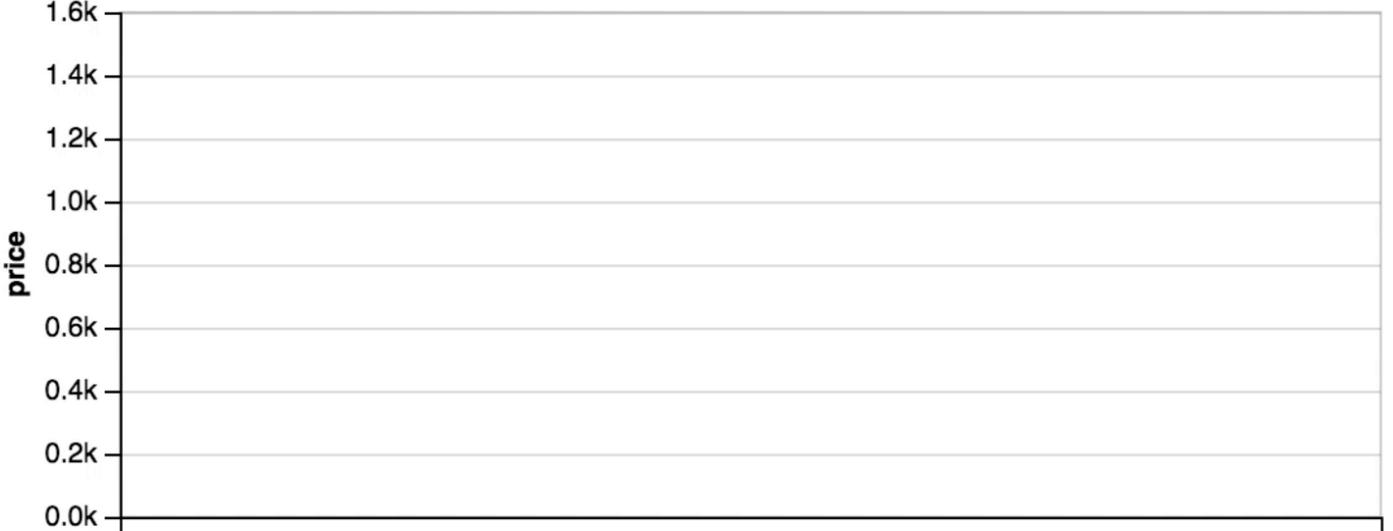
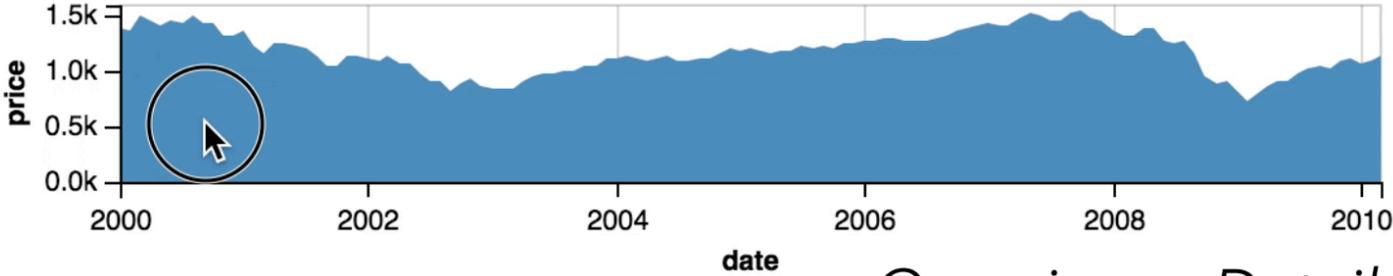
**Multi-view interactive graphics in ~10 lines of code!**

# Interactive Selections

Selections *invert* scales and *parameterize* graphics



Bind selection to scale domains:  
*Synchronized Pan & Zoom!*

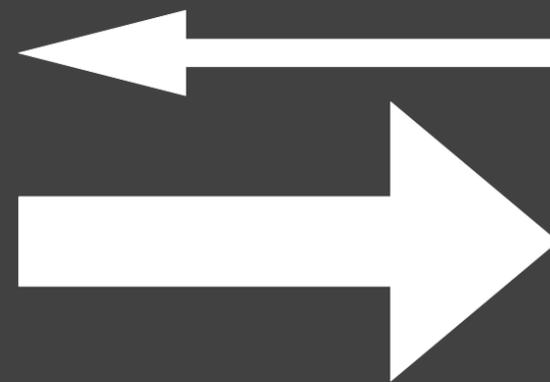


*Parameterized Transformations*

How might we enable **scalable,**  
**real-time visual queries?**



Database



**Latency?**



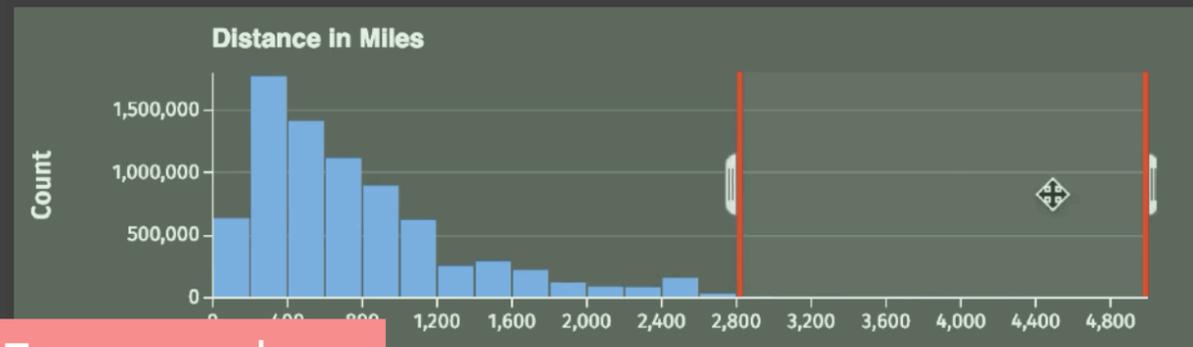
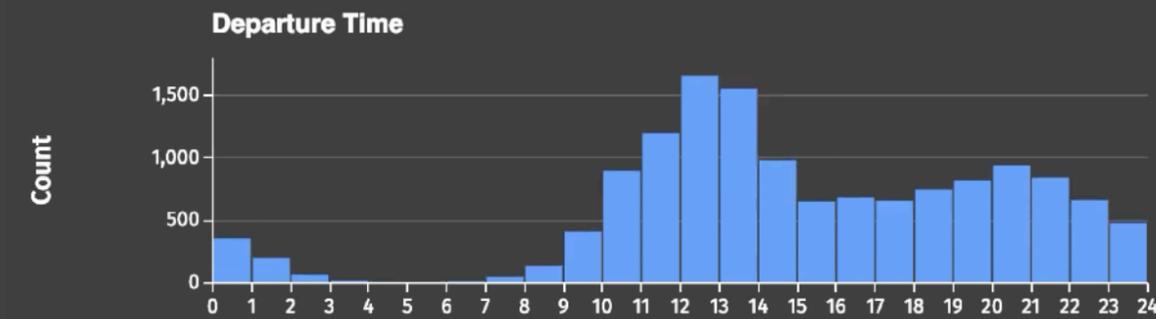
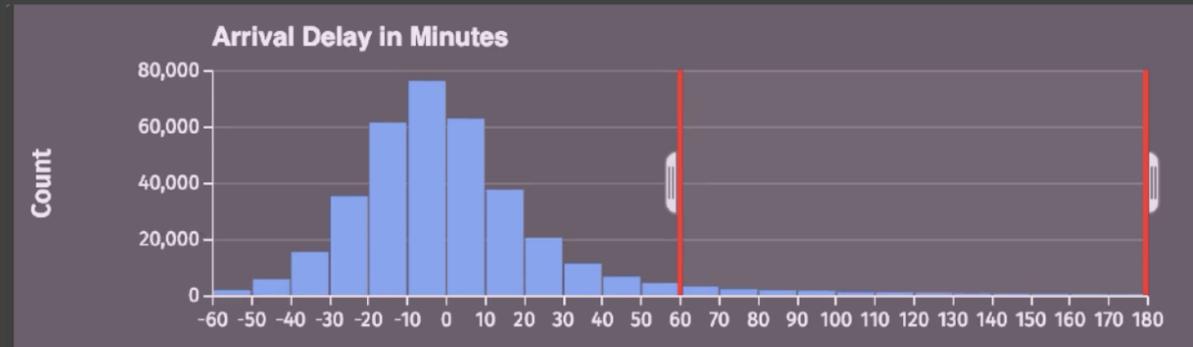
Vega

JavaScript

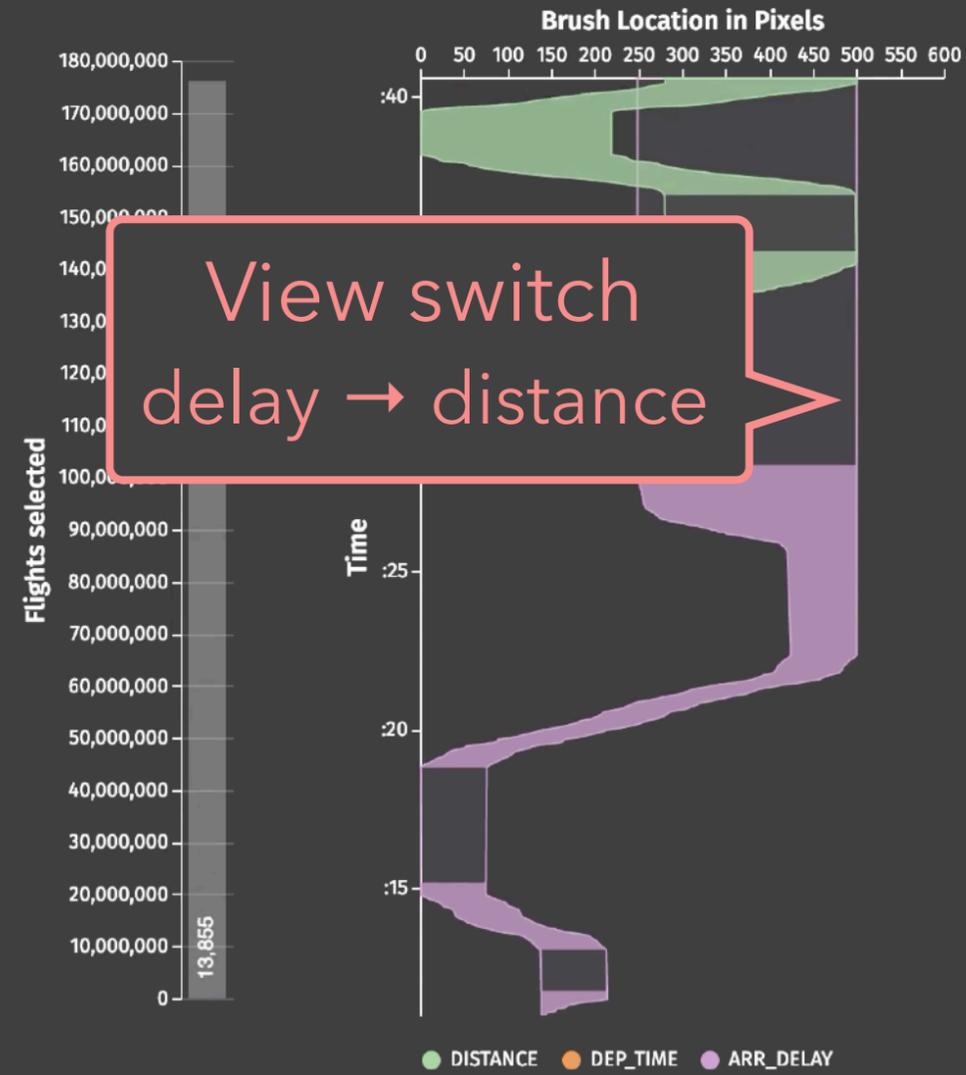
**Milliseconds matter** for latency-sensitive interactions like brushing!

[Liu & Heer 2014, Zraggen et al. 2017]

# Cross-Filtering Interaction Log

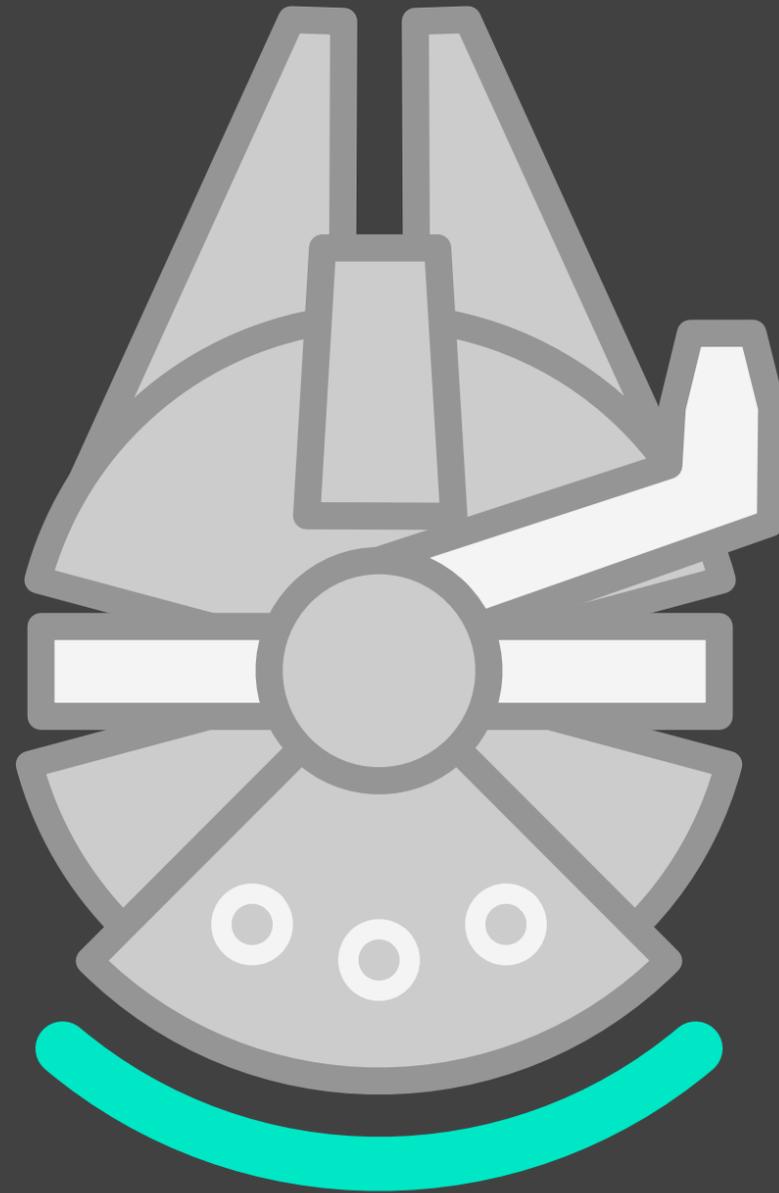


5x speedup



Brushing interactions

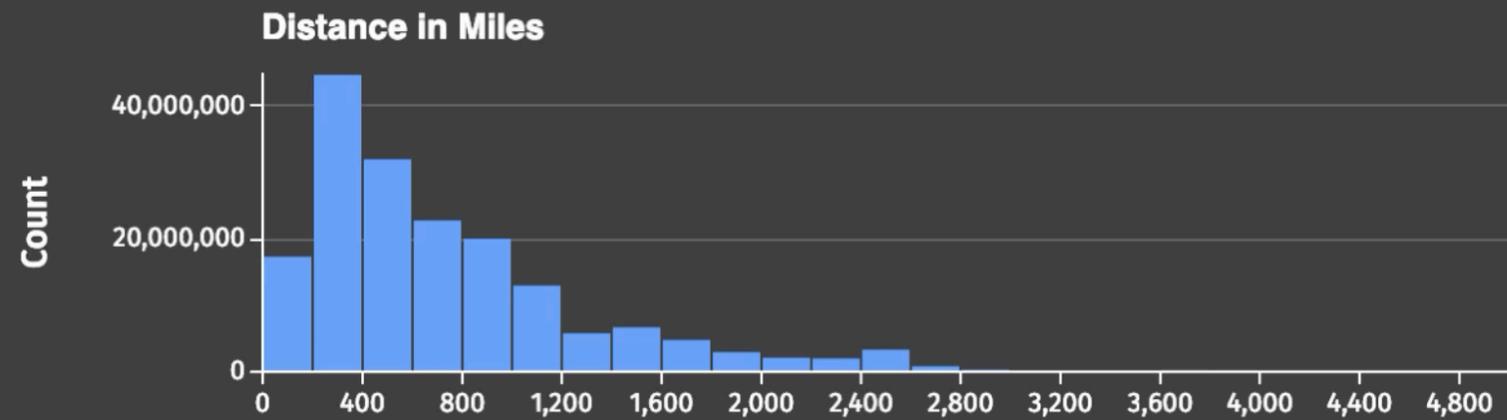
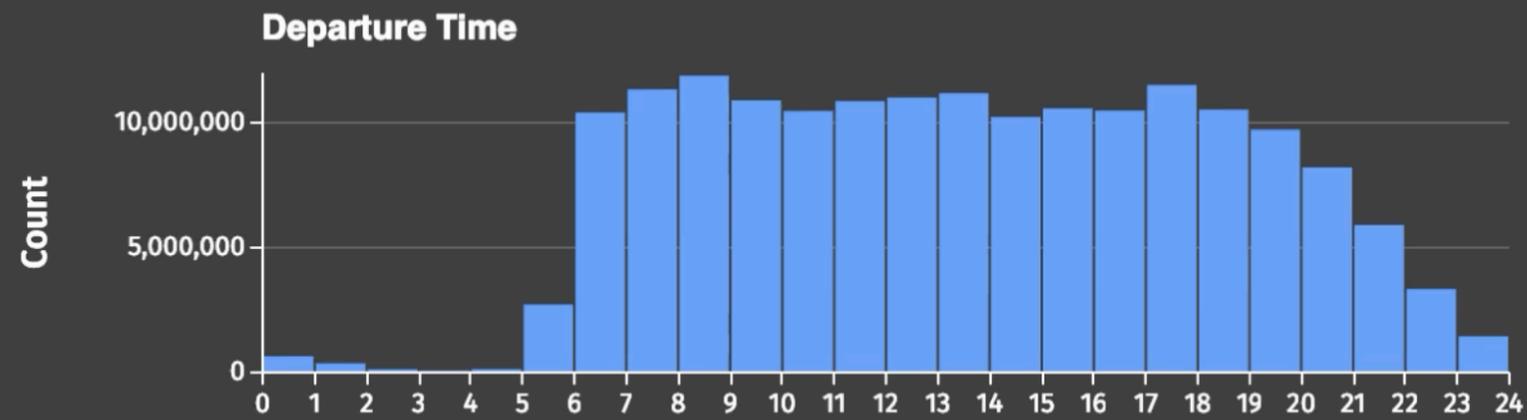
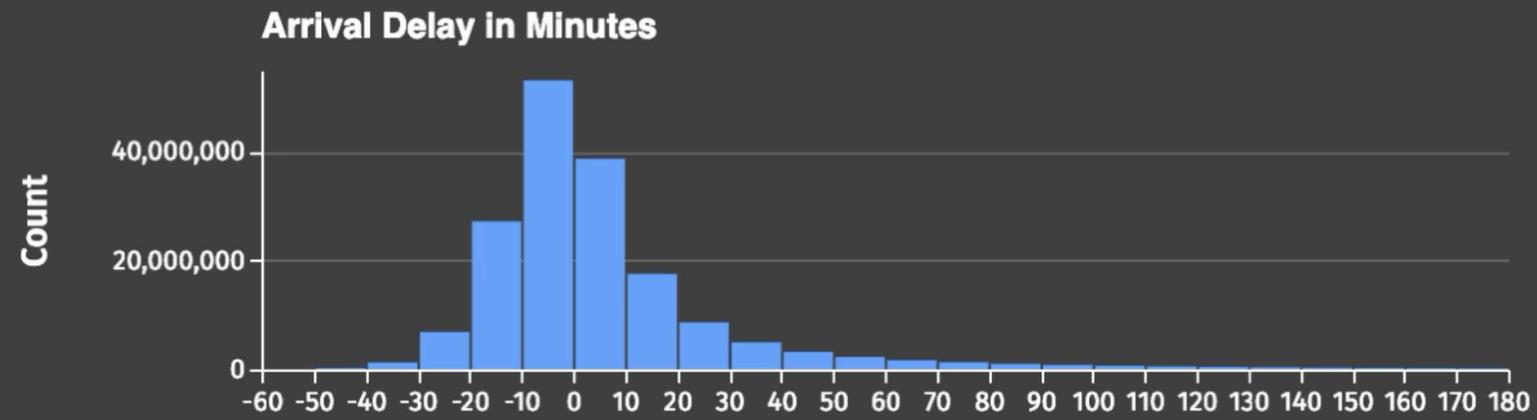
- 👁️ Brushing is more common and people are sensitive to latencies.
- 💡 Prioritize brushing latency over view switching latency.



# Falcon

[uwdata.github.io/falcon](https://uwdata.github.io/falcon)

[Moritz et al., CHI 2019]

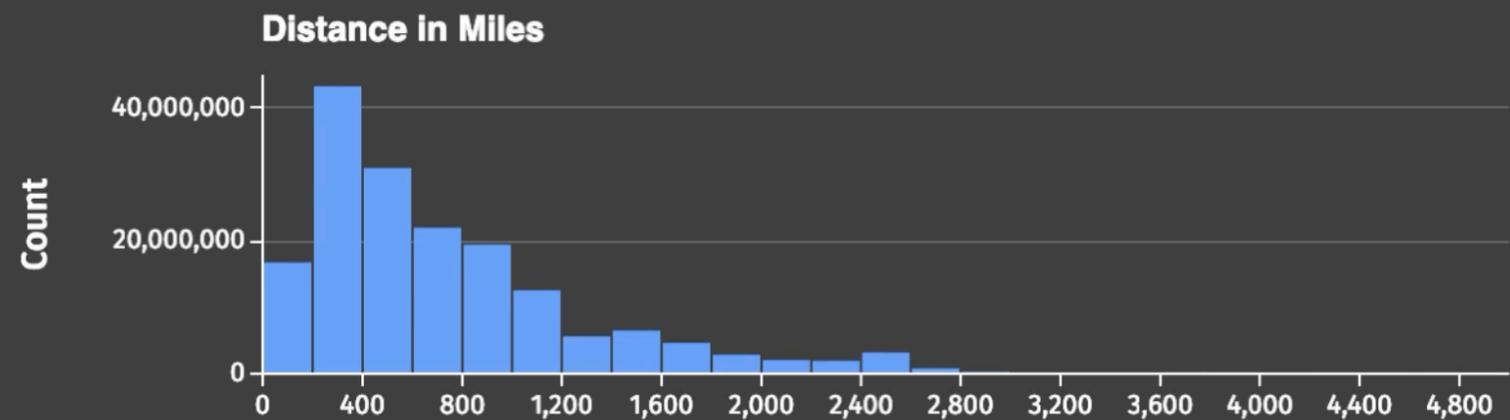
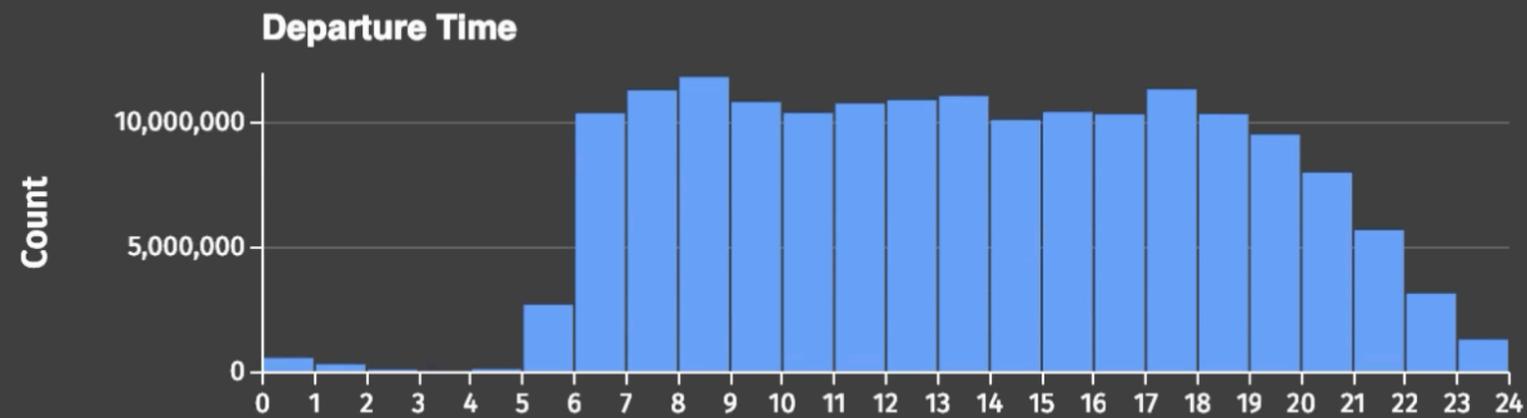
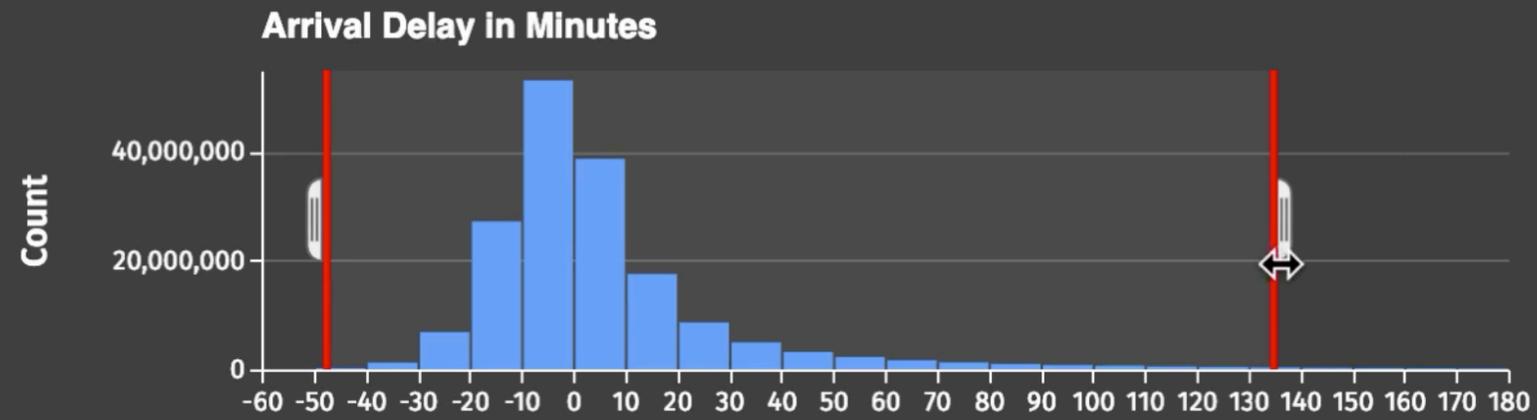


*brushes in the precomputed view*



*serves requests from a data cube*

Data Cube. Gray et al. 1997.



*brushes in the precomputed view*



*serves requests from a data cube*  
Data Cube. Gray et al. 1997.



*interacts with a new view*



*query for new data cubes*

Constant data & time.  
Client only.



*brushes in the precomputed view*



*serves requests from a data cube*  
Data Cube. Gray et al. 1997.

💡 Aggregation decouples interactions from queries over the raw data.

Requires one pass  
over the data.

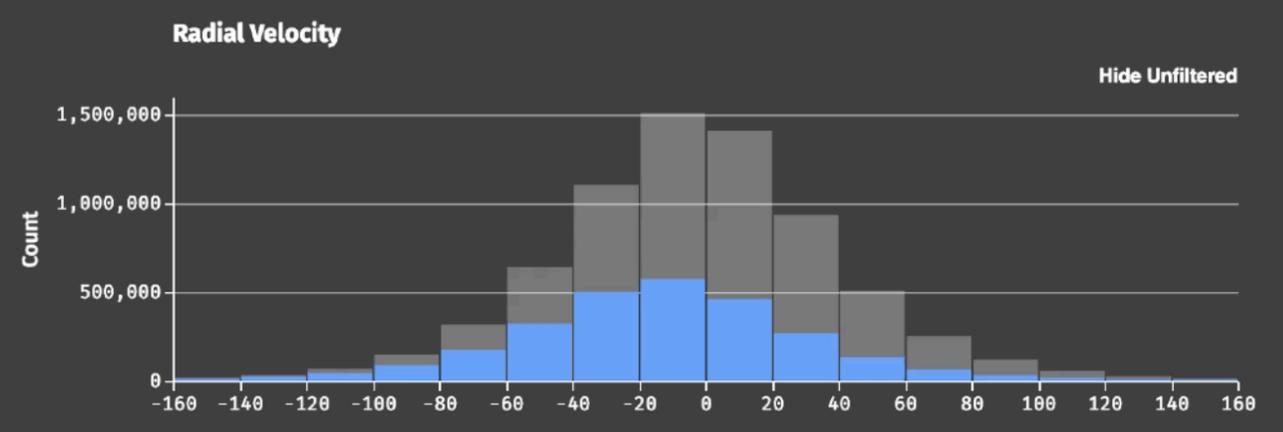
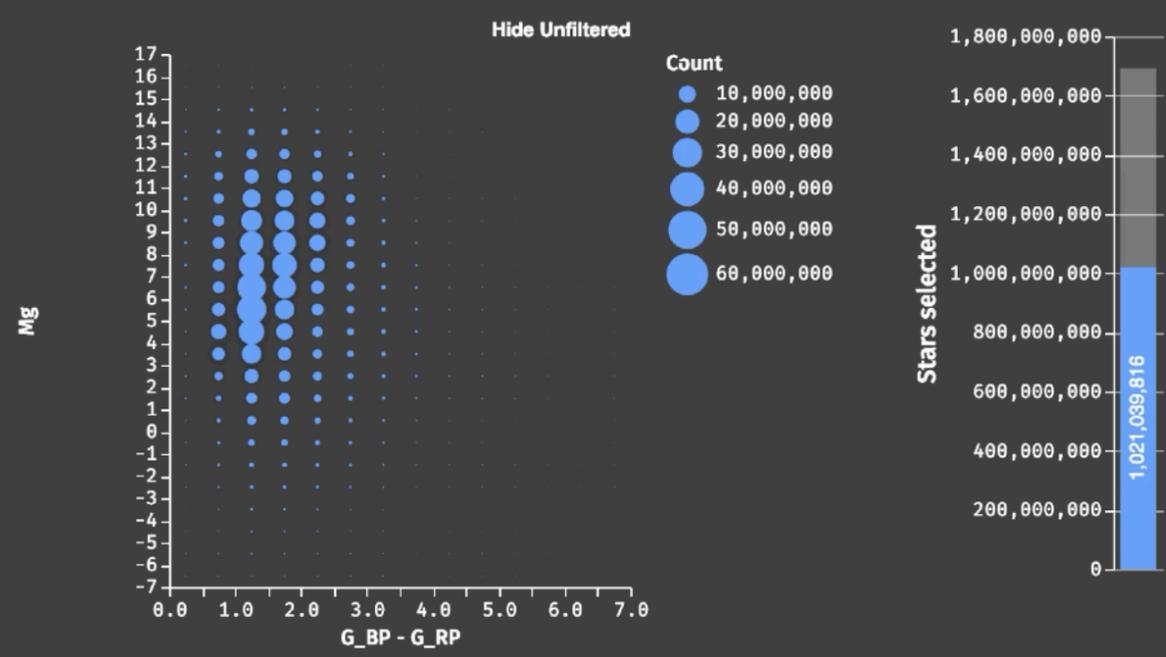
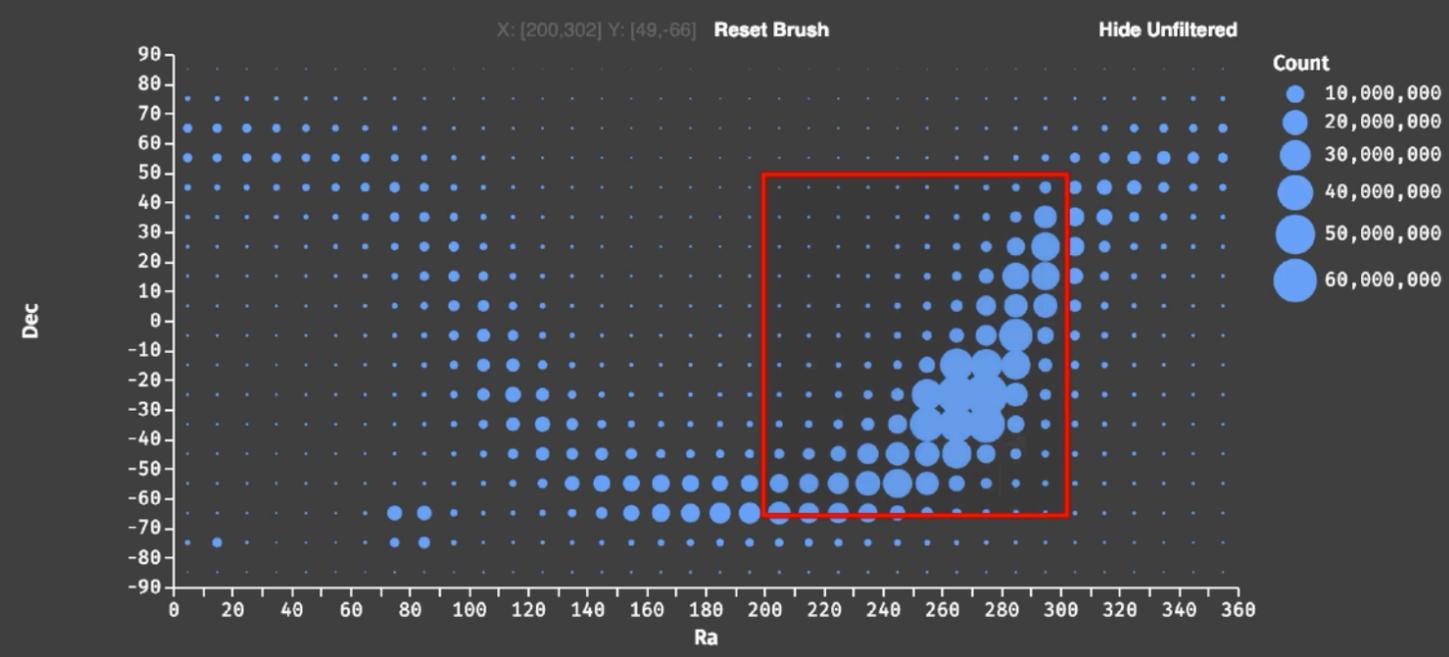
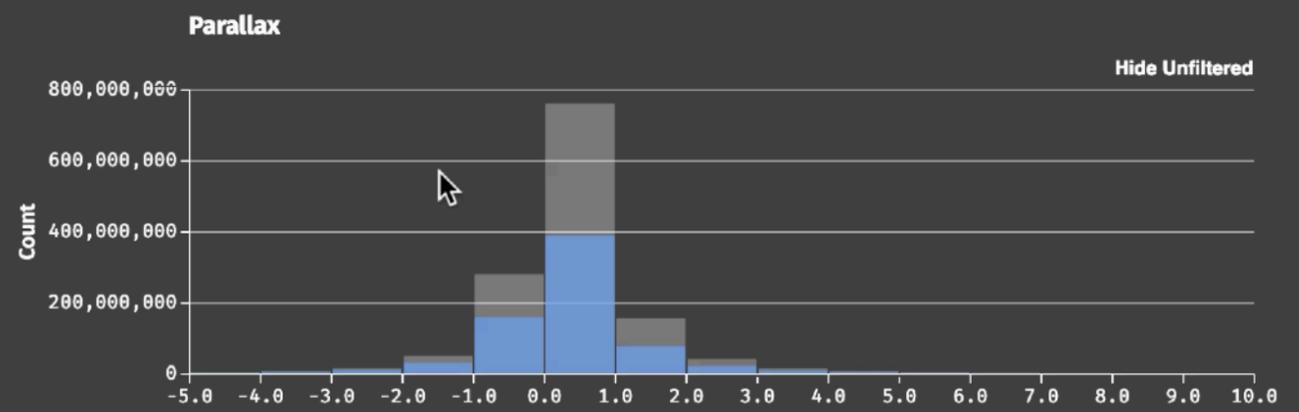
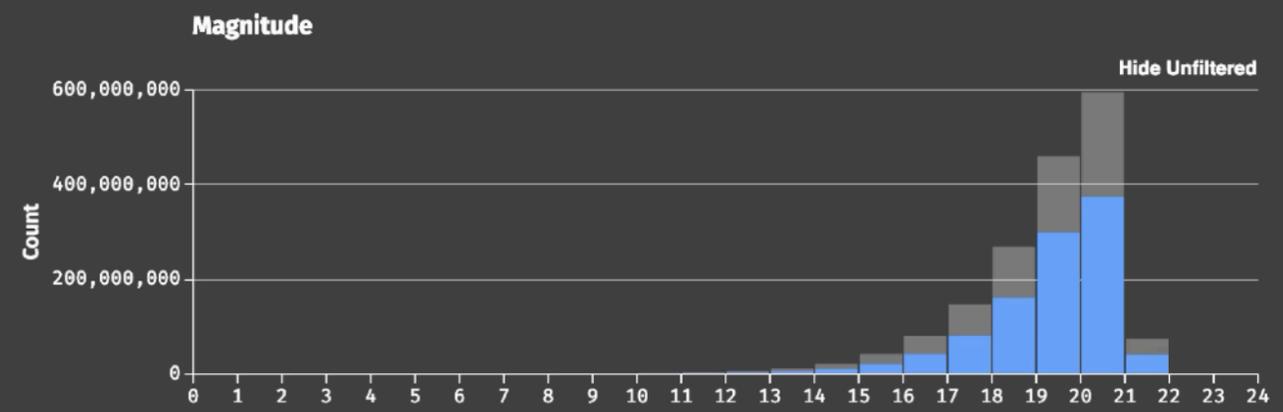


*interacts with a new view*



*query for new data cubes*

💡 View switches are **rare** and users are **not as latency sensitive** with them.



1.7 B stars.  
1.2 TB of data.  
Visualizations running in the browser.  
Data stored in Heavy.ai GPU database.

How might we support more  
**effective data exploration?**

# Common exploration pitfalls:

Overlook data quality issues

Fixate on specific relationships

*Plus many other biases...*

[Heuer 1999, Kahneman 2011, ...]

Voyager 2

Secure https://uwdata.github.io/voyager2/

datavoyager

Bookmarks (0) Undo Redo

**Data**

Cars Change

**Fields**

- Cylinders
- Name
- Origin
- Year
- Acceleration
- Displacement
- Horsepower
- Miles per Gallon
- Weight in lbs
- COUNT

**Wildcards**

- Categorical Fields
- Temporal Fields
- Quantitative Fields

**Encoding** Clear

x YEAR (Year)

y # MEAN (Miles per

column drop a field here

row drop a field here

**Marks** auto

size drop a field here

color drop a field here

shape drop a field here

detail drop a field here

text drop a field here

any drop a field here

**Filter** Filter invalid numbers

**Related Views** All Add Categorical Field Add Quantitative Field Hide

**Add Categorical Field**

YEAR (Year) # MEAN (Miles per Gallon) Cylinders

MEAN(Miles\_per\_Gallon)

YEAR(Year)

**Origin**

YEAR (Year) # MEAN (Miles per Gallon) Origin

MEAN(Miles\_per\_Gallon)

YEAR(Year)

Debug · Report an Issue

**Voyager:** Combine Specification and Recommendation [Wongsuphasawat 2017]

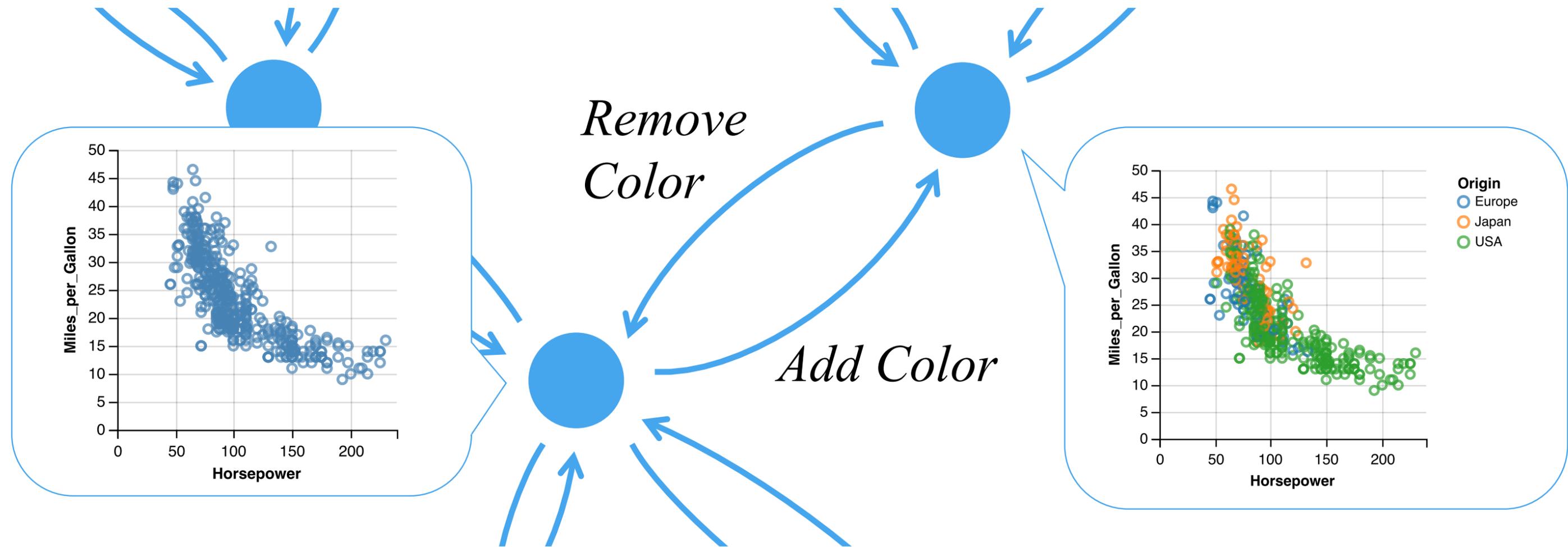
**Key Idea:** Augment manual exploration with visualization recommendations sensitive to the user's current focus.

The ultimate goal is to support *systematic consideration* of the data, without exacerbating *false discovery*.

To model a user's search frontier, we *optimize for related chart specifications*, seeded by the user's current focus.

Candidate charts are pruned and ranked using a formal model of *design constraints* and *perceptual effectiveness*, which can be trained from perception experiment results.

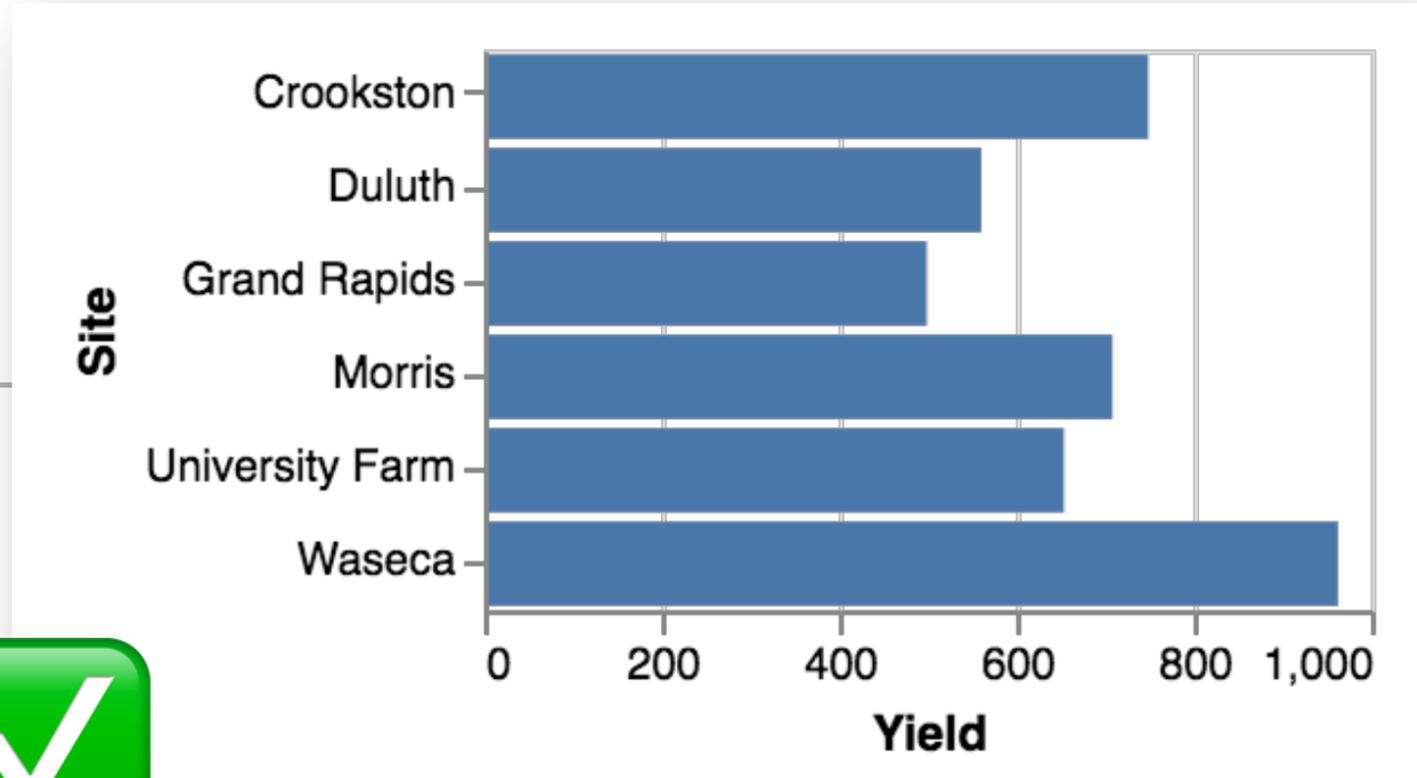
# A Formal Design Space of Visualizations



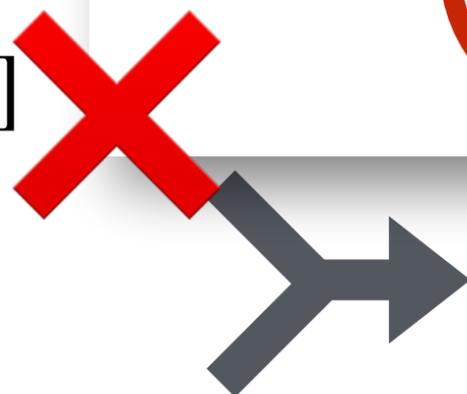
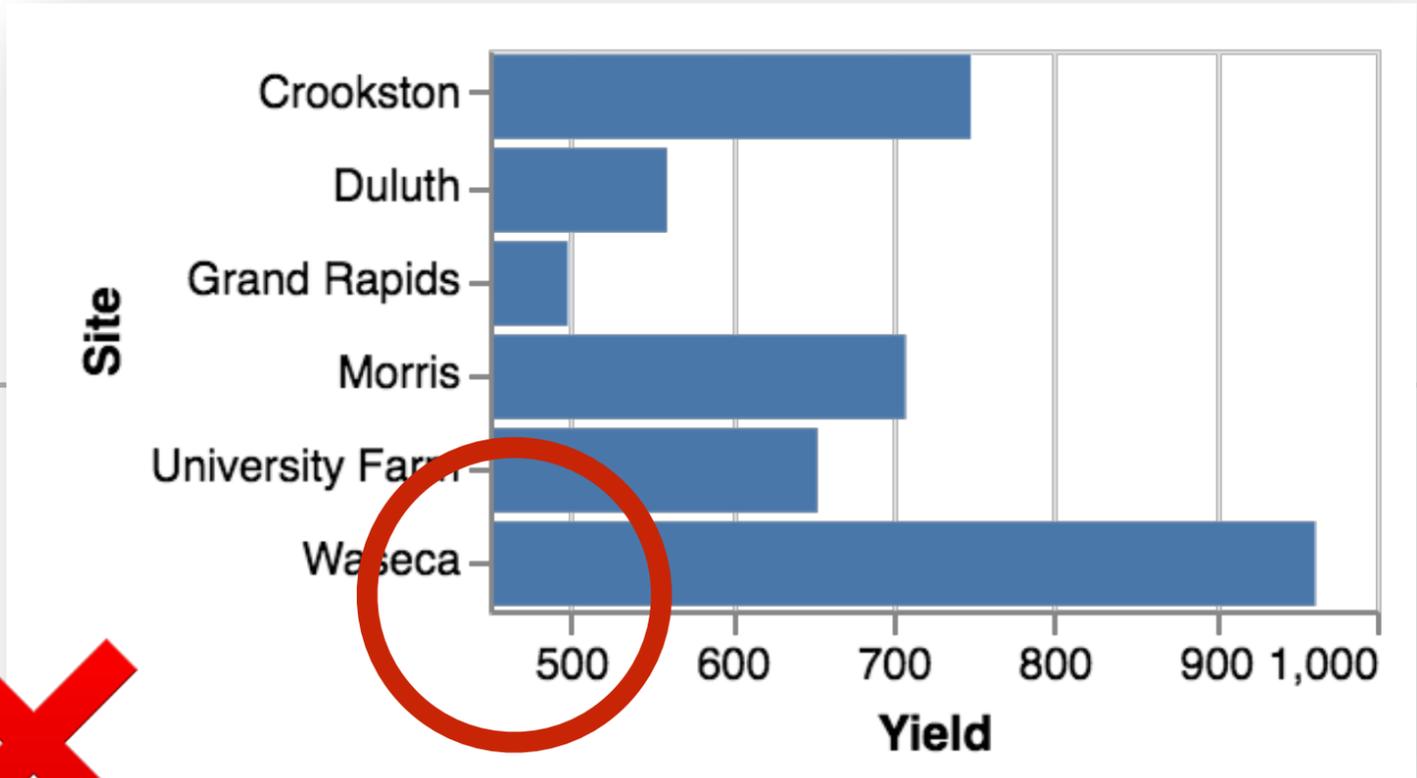
Enumerate Vega-Lite specifications and transformations among them.

Search the space using logic programming methods.

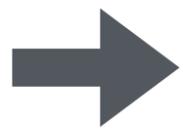
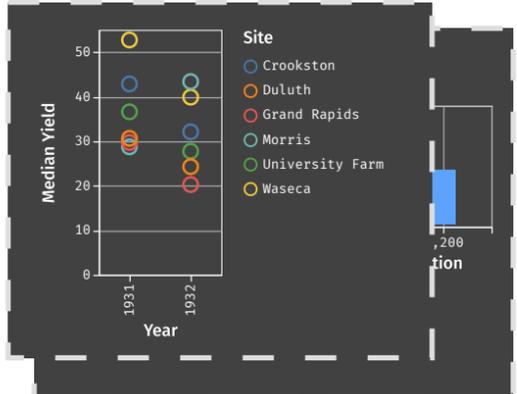
# Learn Design Trade-Offs from Experimental Data



S  
of  
raints  
ample  
ector  
[ $u_1, \dots, u_k$ ]



$$\arg \max_w \sum_{i \in 0..k} w_i (u_i - v_i)$$



👎 Negative example  
Feature Vector  
[ $v_1, v_2, \dots, v_k$ ]

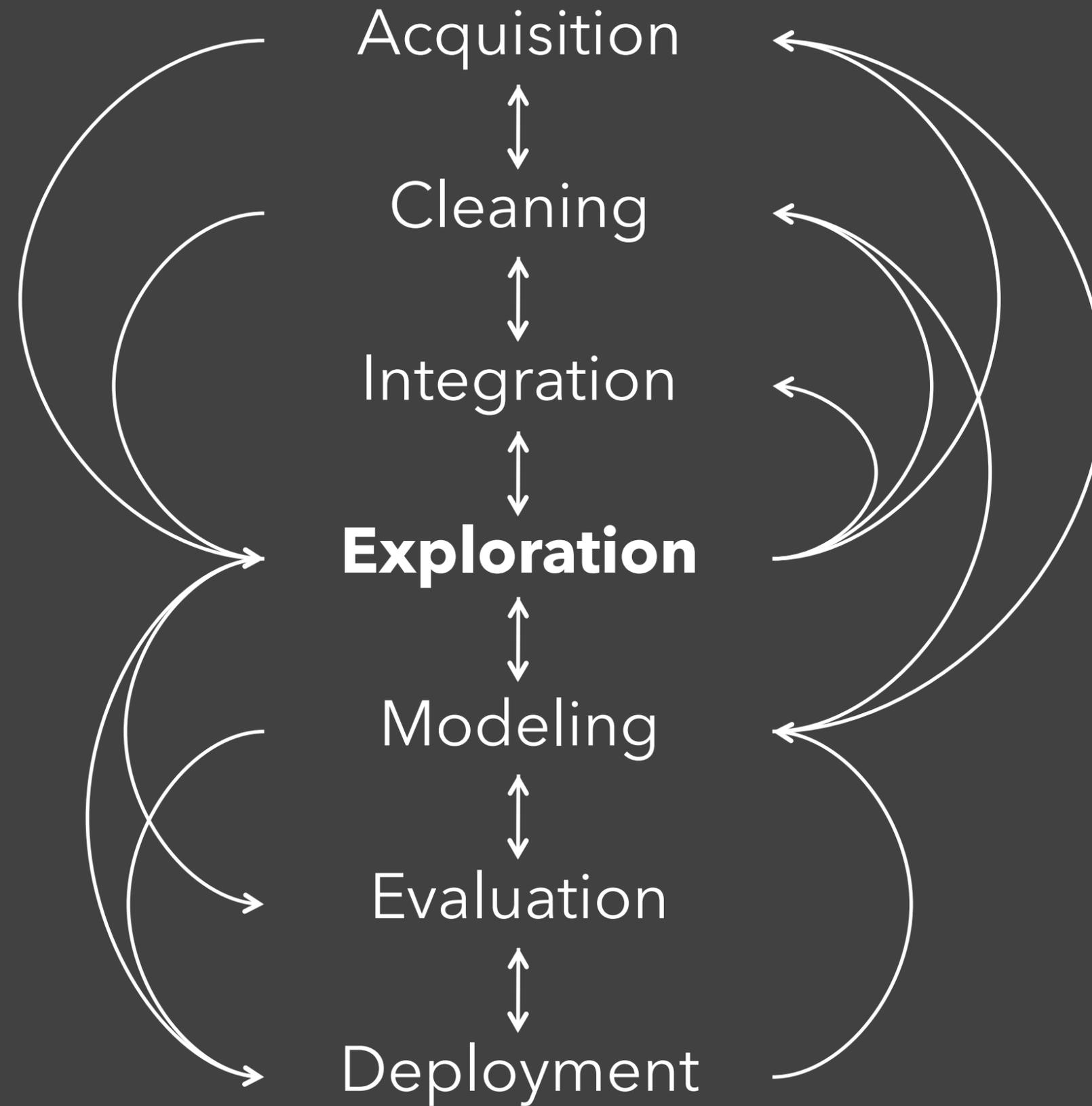
$v_i$ : the number of violations of constraint .

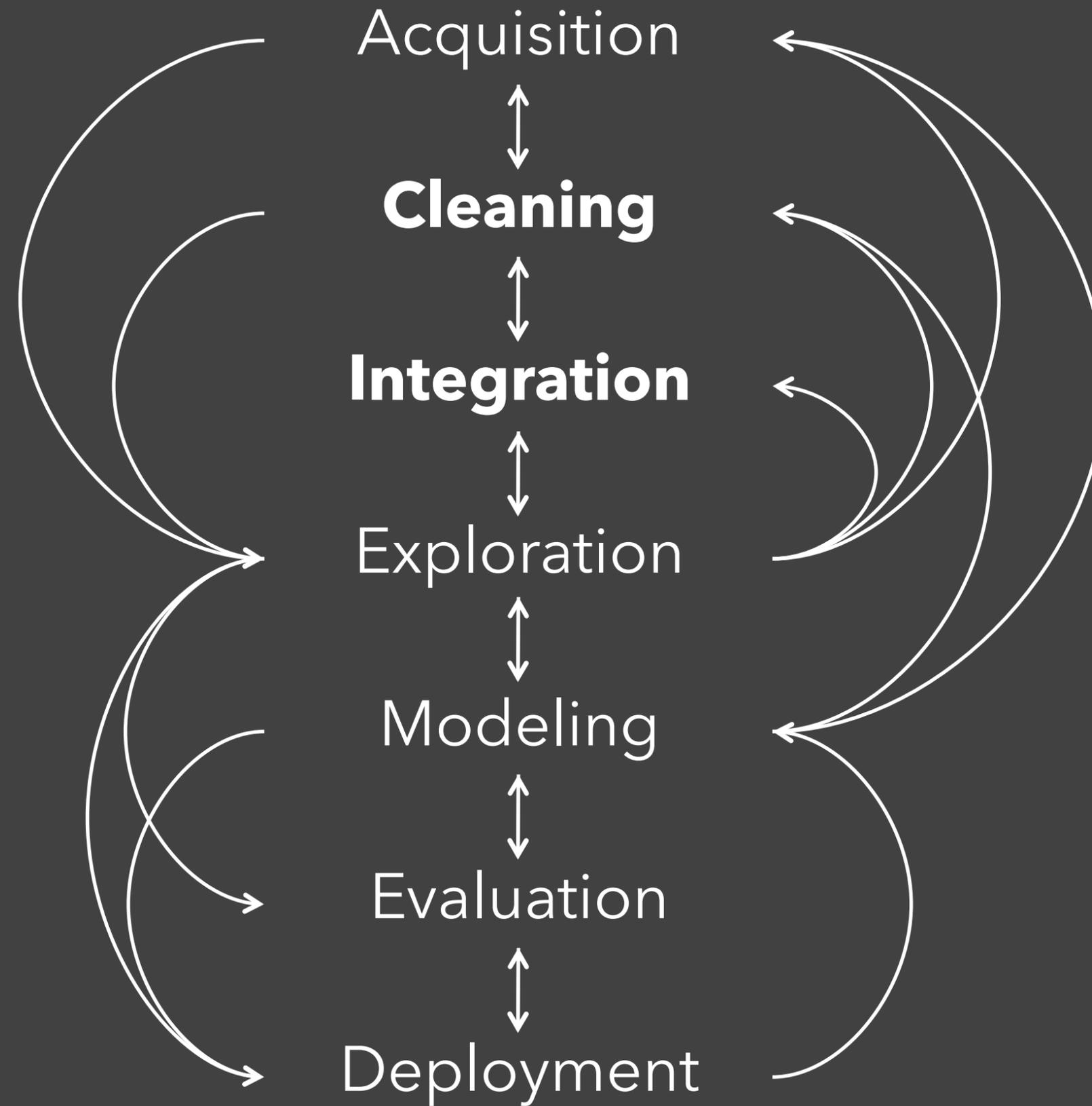
Compared to existing tools, leads to **over 4x more variable sets seen**, and **over 2x more variable sets interacted with**.

*"The related view suggestion accelerates exploration a lot."*

*"I like that it shows me what fields to include in order to see a specific graph. Otherwise, I have to do a lot of trial and error and can't express what I wanted to see."*

*"These related views are so good but it's also spoiling that I start thinking less. I'm not sure if that's really a good thing."*





I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any "analysis" at all.

Anonymous Data Scientist  
*from our 2012 interview study*





**Big Data  
Borat**

@BigDataBorat



Following

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.



How might we assist the process  
of **data transformation**?

Reported crime in Alabama

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	4525375	4029.3	987	2732.4	309.9
2005	4548327	3900	955.8	2656	289
2006	4599030	3937	968.9	2645.1	322.9
2007	4627851	3974.9	980.2	2687	307.7
2008	4661900	4081.9	1080.7	2712.6	288.6

Reported crime in Alaska

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	657755	3370.9	573.6	2456.7	340.6
2005	663253	3615	622.8	2601	391
2006	670053	3582	615.2	2588.5	378.3
2007	683478	3373.9	538.9	2480	355.1
2008	686293	2928.3	470.9	2219.9	237.5

Reported crime in Arizona

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	5739879	5073.3	991	3118.7	963.5
2005	5953007	4827	946.2	2958	922
2006	6166318	4741.6	953	2874.1	914.4
2007	6338755	4502.6	935.4	2780.5	786.7
2008	6500180	4087.3	894.2	2605.3	587.8

Reported crime in Arkansas

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	2750000	4033.1	1096.4	2699.7	237
2005	2775708	4068	1085.1	2720	262
2006	2810872	4021.6	1154.4	2596.7	270.4
2007	2834797	3945.5	1124.4	2574.6	246.5
2008	2855390	3843.7	1182.7	2433.4	227.6

# Data Wrangler

The screenshot displays the Data Wrangler interface. On the left, a 'Suggestions' panel lists several data cleaning actions: 'Delete rows 8,10', 'Delete empty rows', 'Delete rows where Property\_crime\_rate is null', and 'Delete rows where Year is null'. Below this is a 'Script' panel with an 'Export' button and two transformation suggestions: 'Split data repeatedly on newline into rows' and 'Split data repeatedly on \',\''. The main area shows a data table with 408 rows. The table has two columns: '# Year' and '# Property\_crime\_rate'. The data is grouped by state, with 'Reported crime in Alabama' (rows 1-7) and 'Reported crime in Alaska' (rows 9-14). Row 8 is highlighted in light blue. The table header includes 'rows: 408' and 'prev next' navigation options.

#	Year	#	Property_crime_rate
1	Reported crime in Alabama		
2			
3	2004		4029.3
4	2005		3900
5	2006		3937
6	2007		3974.9
7	2008		4081.9
8			
9	Reported crime in Alaska		
10			
11	2004		3370.9
12	2005		3615
13	2006		3582
14	2007		3373.9

**Wrangler: Interactive Visual Specification of Data Transformation Scripts**

Sean Kandel, Joseph Hellerstein et al. *CHI'11*

# DataWrangler

Reduces specification time, promotes the use of reusable, scalable transformations rather than idiosyncratic edits.

*Complementarity:* Suggestions good for tasks people found hard (extraction patterns, table reshaping). People good in cases where inference is less tractable (arbitrary formulas).

*Agency:* Users appreciated suggestions in response to an initiating interaction, but did not always act on *proactive* assistance, often preferring to maintain the initiative.

Wrangler: Interactive Visual Specification of Data Transformation Scripts

Kandel et al. [CHI 2011]



TRIFACTA®

CAND_ID	CAND_NAME	CAND_PARTY_AFFILIATION	CAND_ELECTION_YEAR	CAND_OFFICE_STATE	CAND_OFFICE
4,864 Categories	4,760 Categories	76 Categories	1986 - 2052	57 Categories	3 Categories
H0AK00097	COX, JOHN R.	REP	2014	AK	H
H0AL02087	ROBY, MARTHA	REP	2016	AL	H
H0AL02095	JOHN, ROBERT E JR	IND	2016	AL	H
H0AL05049	CRAMER, ROBERT E "BUD" JR	DEM	2008	AL	H
H0AL05163	BROOKS, MO	REP	2016	AL	H
H0AL06088	COOKE, STANLEY KYLE	REP	2010	AL	H
H0AL07086	SEWELL, TERRI A.	DEM	2016	AL	H
H0AL07094	HILLIARD, EARL FREDERICK JR	DEM	2010	AL	H
H0AL07177	CHAMBERLAIN, DON	REP	2012	AL	H
H0AR01083	CRAWFORD, ERIC ALAN RICK	REP	2016	AR	H
H0AR01091	GREGORY, JAMES CHRISTOPHER	DEM	2010	AR	H
H0AR01109	CAUSEY, CHAD	DEM	2010	AR	H
H0AR01125	SMITH, PRINCELLA D	REP	2010	AR	H
H0AR02107	GRIFFIN, JOHN TIMOTHY	REP	2014	AR	H
H0AR02131	ELLIOTT, JOYCE ANN	DEM	2010	AR	H
H0AR03022	SKOCH, BERNARD KURT 'BERNIE'	REP	2010	AR	H
H0AR03030	WHITAKER, DAVID JEFFREY	DEM	2010	AR	H
H0AR03055	WOMACK, STEVE	REP	2016	AR	H
H0AS00018	FALEOMAVAEGA, ENI	DEM	2014	AS	H
H0AZ01184	FLAKE, JEFF MR.	REP	2012	AZ	H
H0AZ01259	GOSAR, PAUL ANTHONY	REP	2016	AZ	H
H0AZ01283	MEHTA, STEVE	REP	2010	AZ	H
H0AZ01325	TOBIN, ANDY HON.	REP	2014	AZ	H
H0AZ01333	GRESSLEY, FORREST DAYL	REP	2010	AZ	H
H0AZ03321	PARKER, VERNON	REP	2014	AZ	H

New Step [Switch to editor](#)

Cancel [Add to Recipe](#)

Choose a transformation

Choose transformation

Grid Columns Full Dataset - 461.78kB 17 Columns 4,864 Rows 3 Data Types Columns: All Transformed - 3 Columns Rows: All Transformed - 4,859 Rows

ABC	CAND_ID	Source	to be dropped	Preview	ABC	CAND_NAME	ABC	CAND_NAME1	ABC	CAND_NAME2	ABC	CAND_PARTY_AFFILIATION	ABC	CAND_ELECTION_YEAR	ABC
		4,864 Categories					4,760 Categories					76 Categories		1986 - 2052	57 Cate
	H0AK00097	COX, JOHN R.				COX, JOHN R.		COX		JOHN R.		REP		2014	AK
	H0AL02087	ROBY, MARTHA				ROBY, MARTHA		ROBY		MARTHA		REP		2016	AL
	H0AL02095	JOHN, ROBERT E JR				JOHN, ROBERT E JR		JOHN		ROBERT E JR		IND		2016	AL
	H0AL05049	CRAMER, ROBERT E "BUD" JR				CRAMER, ROBERT E "BUD" JR		CRAMER		ROBERT E "BUD" JR		DEM		2008	AL
	H0AL05163	BROOKS, MO				BROOKS, MO		BROOKS		MO		REP		2016	AL
	H0AL06088	COOKE, STANLEY KYLE				COOKE, STANLEY KYLE		COOKE		STANLEY KYLE		REP		2010	AL
	H0AL07086	SEWELL, TERRI A.				SEWELL, TERRI A.		SEWELL		TERRI A.		DEM		2016	AL
	H0AL07094	HILLIARD, EARL FREDERICK JR				HILLIARD, EARL FREDERICK JR		HILLIARD		EARL FREDERICK JR		DEM		2010	AL
	H0AL07177	CHAMBERLAIN, DON				CHAMBERLAIN, DON		CHAMBERLAIN		DON		REP		2012	AL
	H0AR01083	CRAWFORD, ERIC ALAN RICK				CRAWFORD, ERIC ALAN RICK		CRAWFORD		ERIC ALAN RICK		REP		2016	AR
	H0AR01091	GREGORY, JAMES CHRISTOPHER				GREGORY, JAMES CHRISTOPHER		GREGORY		JAMES CHRISTOPHER		DEM		2010	AR
	H0AR01109	CAUSEY, CHAD				CAUSEY, CHAD		CAUSEY		CHAD		DEM		2010	AR
	H0AR01125	SMITH, PRINCELLA D				SMITH, PRINCELLA D		SMITH		PRINCELLA D		REP		2010	AR
	H0AR02107	GRIFFIN, JOHN TIMOTHY				GRIFFIN, JOHN TIMOTHY		GRIFFIN		JOHN TIMOTHY		REP		2014	AR
	H0AR02131	ELLIOTT, JOYCE ANN				ELLIOTT, JOYCE ANN		ELLIOTT		JOYCE ANN		DEM		2010	AR
	H0AR03022	SKOCH, BERNARD KURT 'BERNIE'				SKOCH, BERNARD KURT 'BERNIE'		SKOCH		BERNARD KURT 'BERNIE'		REP		2010	AR
	H0AR03030	WHITAKER, DAVID JEFFREY				WHITAKER, DAVID JEFFREY		WHITAKER		DAVID JEFFREY		DEM		2010	AR
	H0AR03055	WOMACK, STEVE				WOMACK, STEVE		WOMACK		STEVE		REP		2016	AR
	H0AS00018	FALEOMAVAEGA, ENI				FALEOMAVAEGA, ENI		FALEOMAVAEGA		ENI		DEM		2014	AS
	H0AZ01184	FLAKE, JEFF MR.				FLAKE, JEFF MR.		FLAKE		JEFF MR.		REP		2012	AZ
	H0AZ01259	GOSAR, PAUL ANTHONY				GOSAR, PAUL ANTHONY		GOSAR		PAUL ANTHONY		REP		2016	AZ

SUGGESTIONS

Cancel
Modify
Add to Recipe

Split CAND\_NAME into 2 columns on `{delim-ws}`

ABC	CAND_NAME	ABC	CAND_NAME1	ABC	CAND_NAME2
	COX, JOHN R.		COX		JOHN R.
	ROBY, MARTHA		ROBY		MARTHA
	JOHN, ROBERT E JR		JOHN		ROBERT E JR

Affects 1 column, 4859 rows Creates 2 columns

Extract `{delim-ws}` from CAND\_NAME

ABC	CAND_NAME	ABC	CAND_NAME1
	COX, JOHN R.		,
	ROBY, MARTHA		,
	JOHN, ROBERT E JR		,

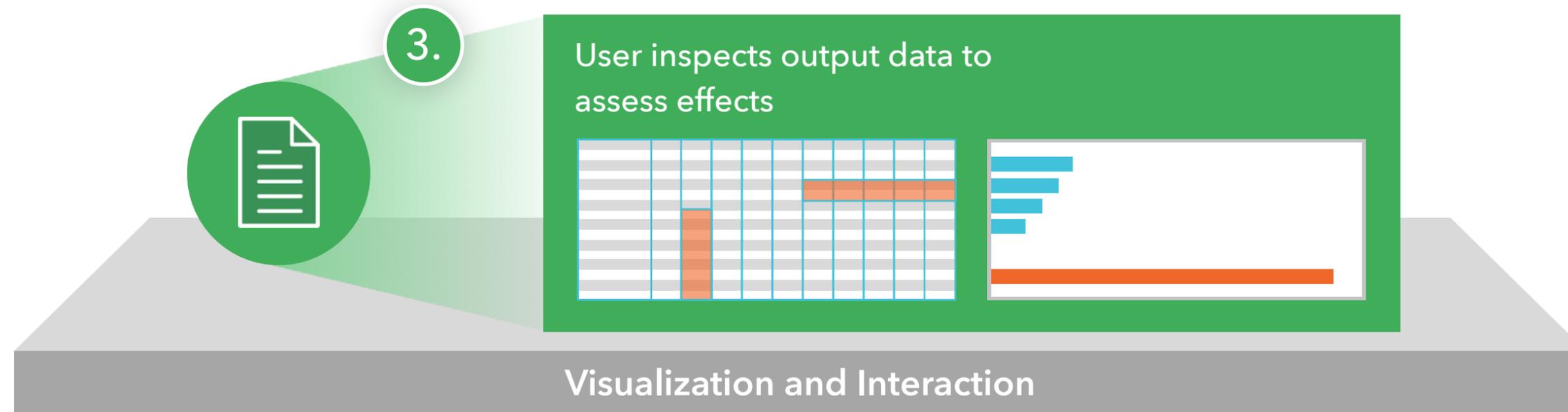
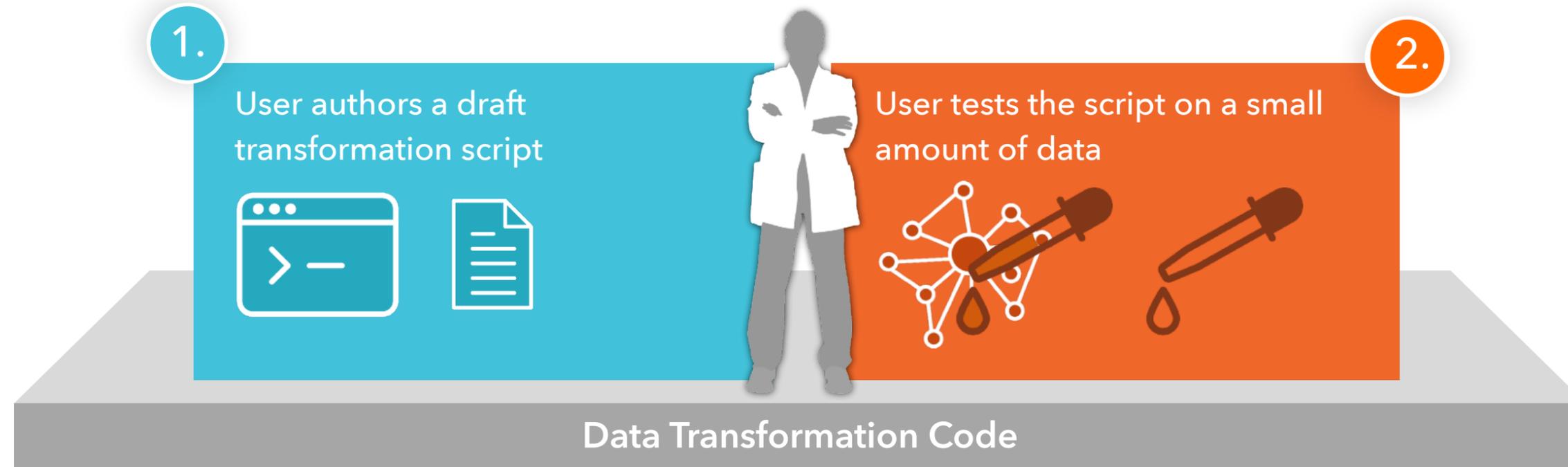
Affects 1 column, 4859 rows Creates 1 column

Count occurrences of `{delim-ws}`

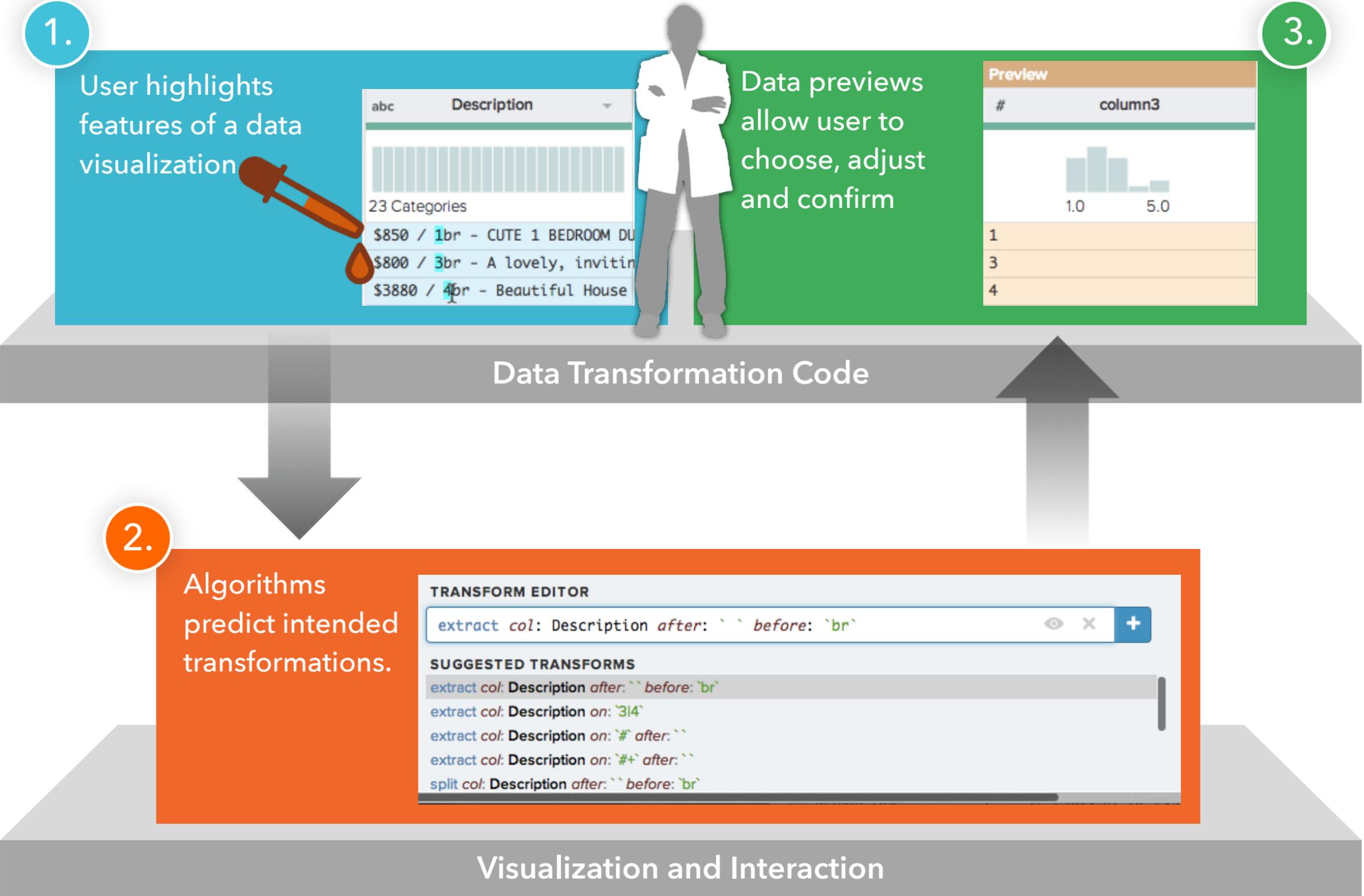
ABC	CAND_NAME
	COX, JOHN R.
	ROBY, MARTHA
	JOHN, ROBERT E JR

Affects 1 column, 4859 rows

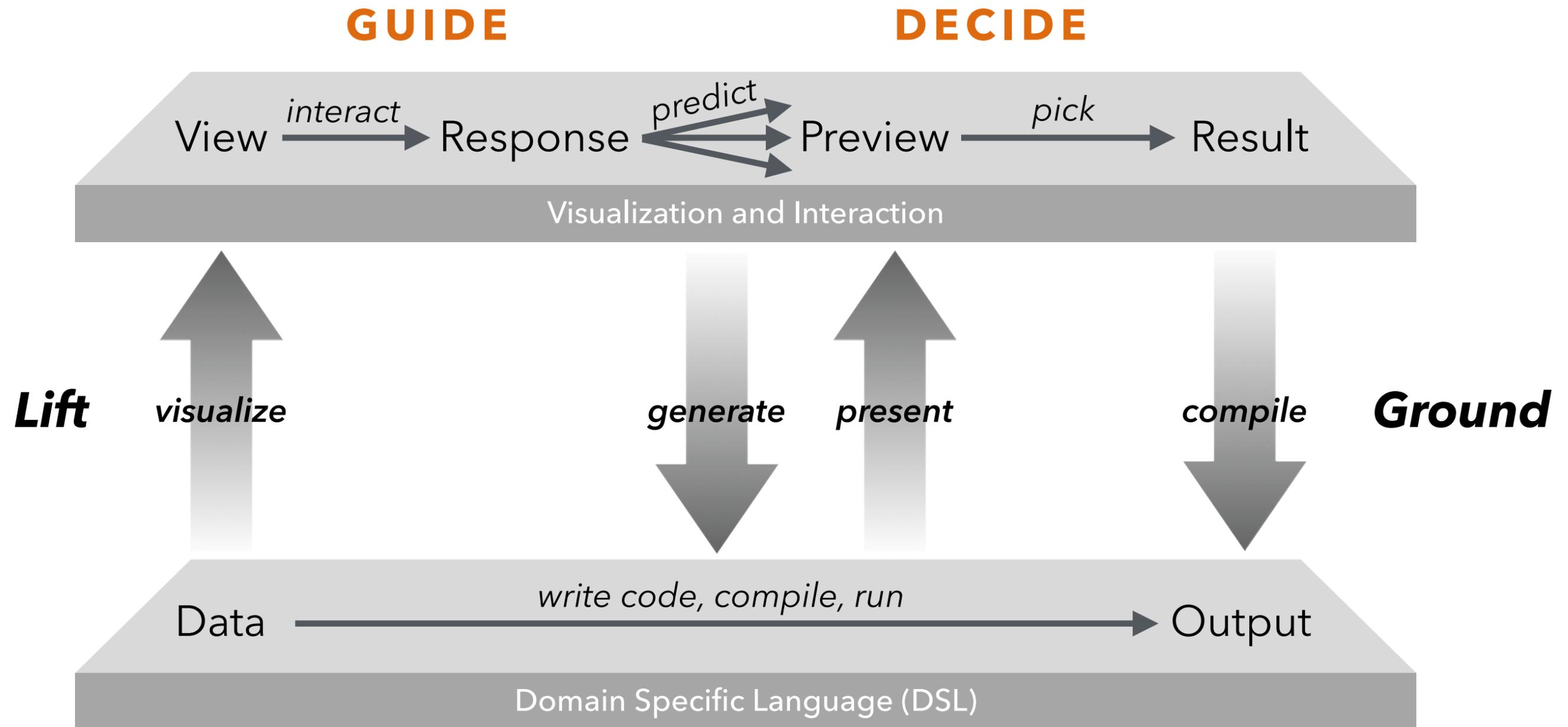
# Traditional Specification

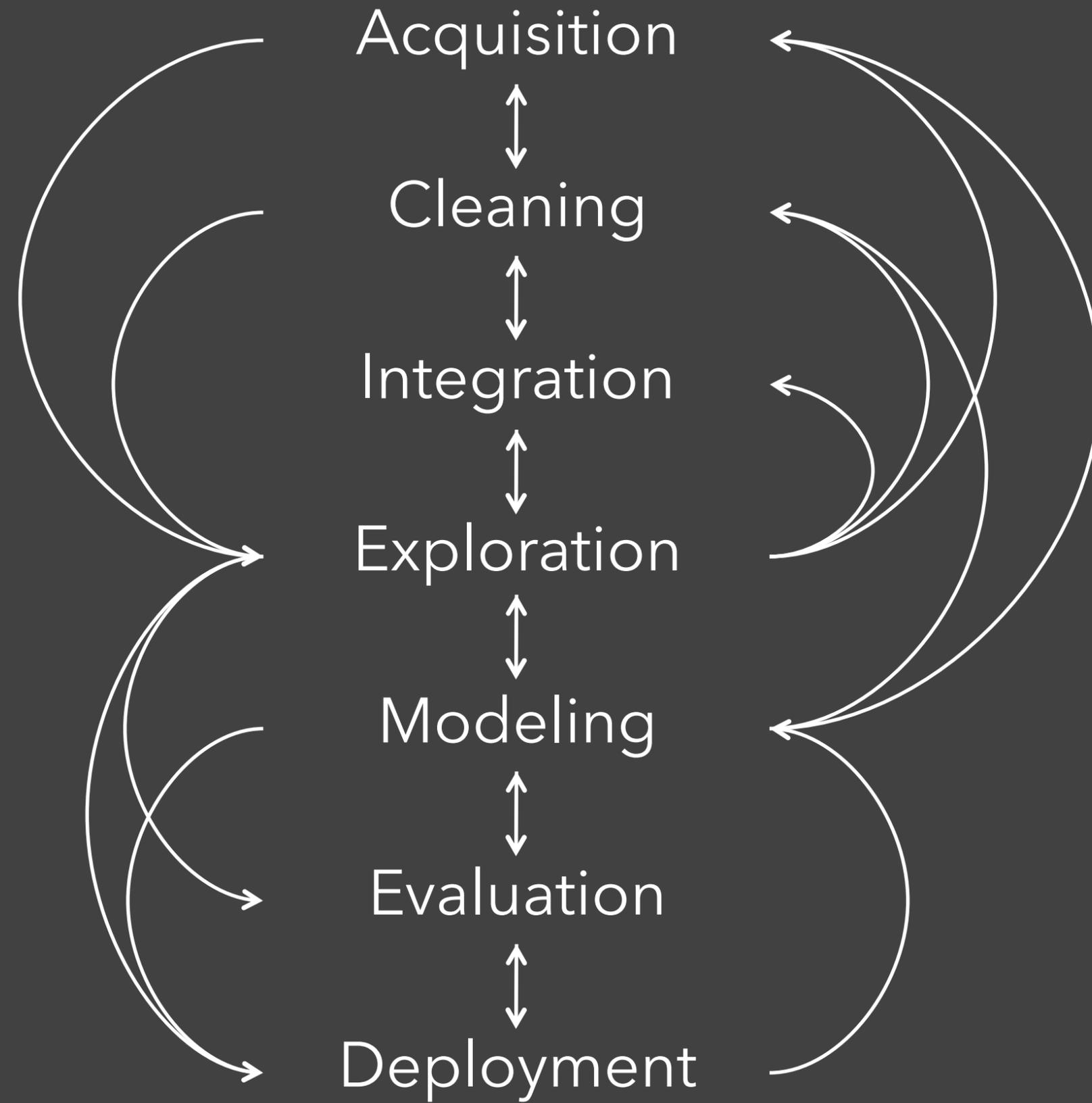


# Predictive Interaction



# Predictive Interaction [Heer, Hellerstein, Kandel CIDR'15]





# Reflections

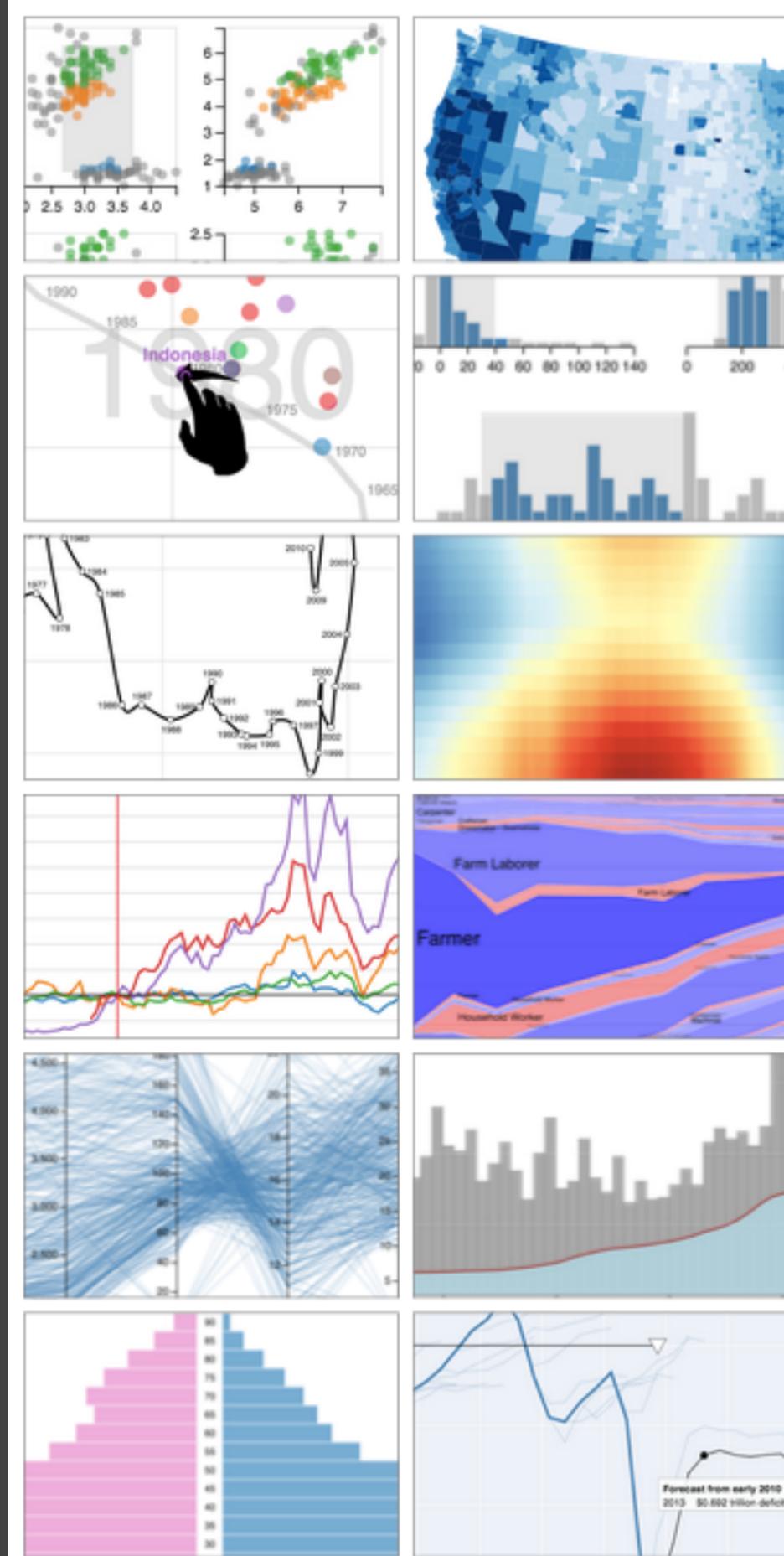




# Successes

## The unreasonable effectiveness of DSLs

Provide users with high-level specification  
Decouple frontend and backend concerns  
Model tasks, enable formal reasoning for  
computational search & recommendation



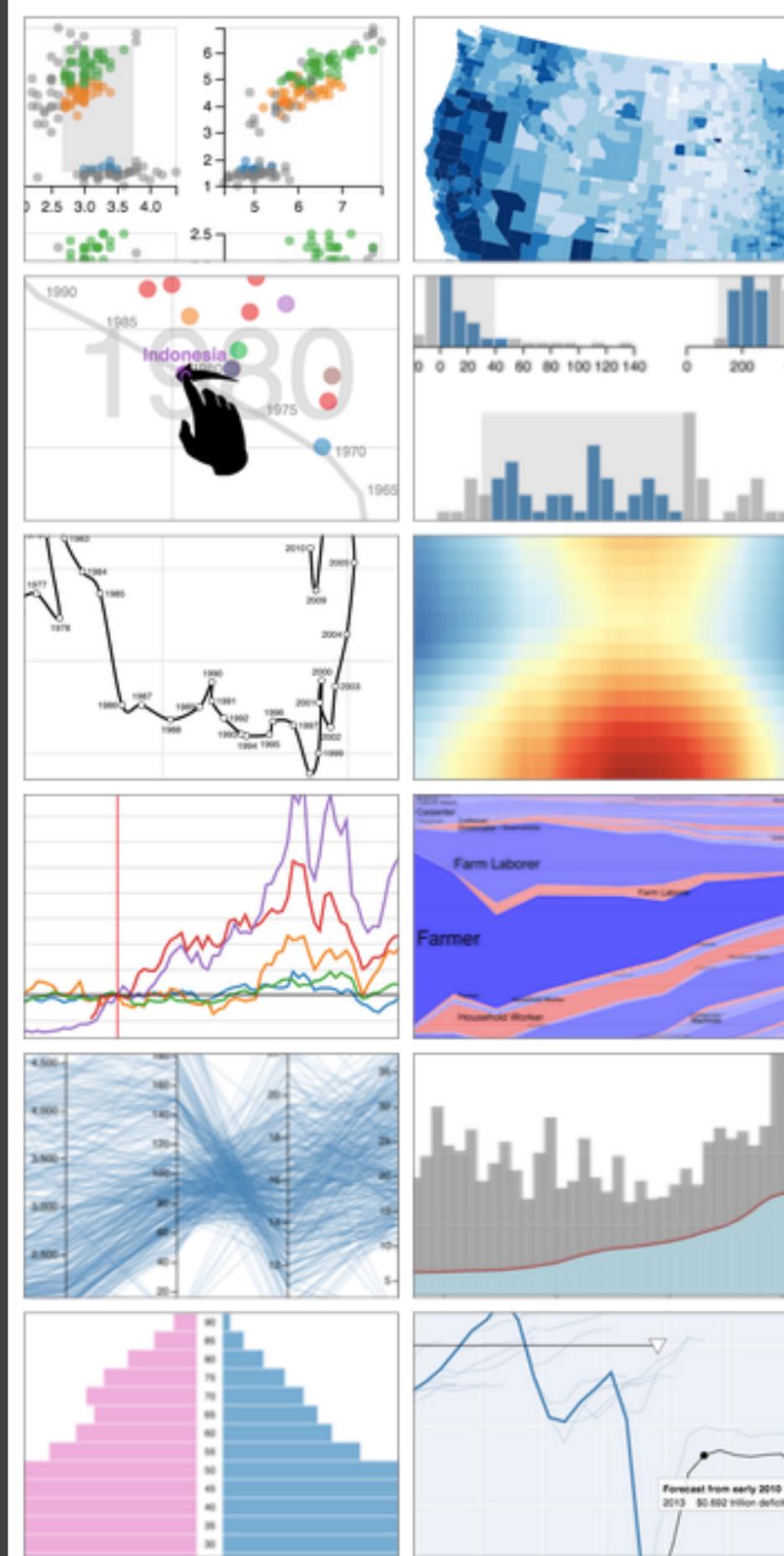
# Successes

## The unreasonable effectiveness of DSLs

Provide users with high-level specification  
Decouple frontend and backend concerns  
Model tasks, enable formal reasoning for  
computational search & recommendation

## AI assistance to bridge user intent

Recommend possible next steps  
Aid discovery and learning



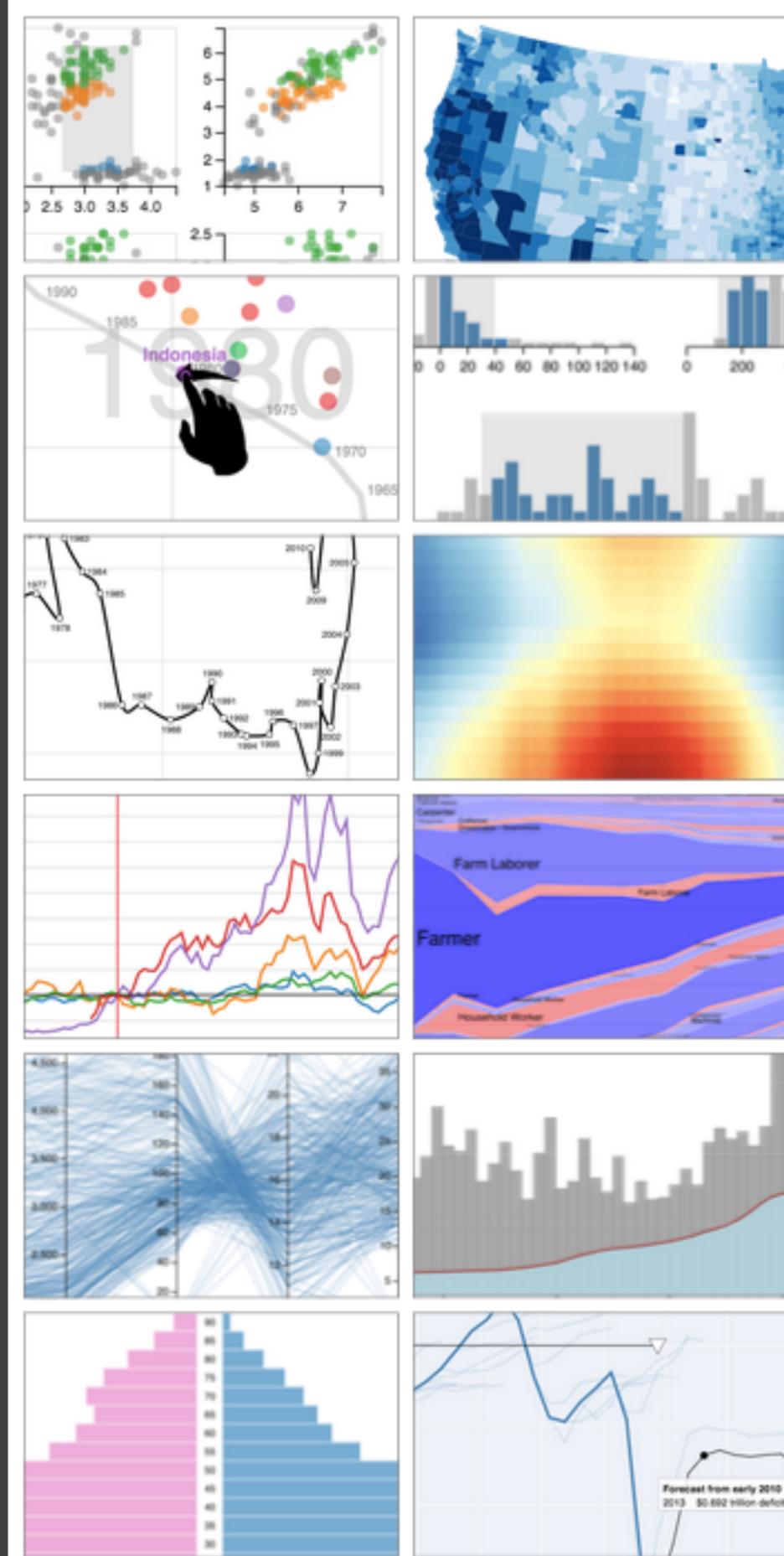
# Successes

## The unreasonable effectiveness of DSLs

Provide users with high-level specification  
Decouple frontend and backend concerns  
Model tasks, enable formal reasoning for  
computational search & recommendation

## AI assistance to bridge user intent

Recommend possible next steps  
Aid discovery and learning  
*How to balance agency and automation?*  
Augment, not replace! [Heer, PNAS 2019]



# A Virtuous Cycle

Declarativity, Scalability, Query Optimization



# Databases

# HCI



New Requirements, Interaction Models, Prefetching Approaches

A Vi



larativ



Optim



as



ments

efetch



The ability to take data—to be able to **understand** it, to **process** it, to **extract value** from it, to **visualize** it, to **communicate** it—that's going to be a hugely important skill in the next decades, ... because now we really do have essentially free and ubiquitous data. So the complimentary **scarce factor** is the ability to **understand** that data and **extract value** from it.

Hal Varian, Google's Chief Economist  
*The McKinsey Quarterly*, Jan 2009

The ability to take data—to be able to **understand** it, to **process** it, to **extract value** from it, to **visualize** it, to **communicate** it—that's going to be a hugely important skill in the next decade **"free" to whom?** because now we really do have essentially free and ubiquitous data. So the complimentary **scarce factor** **"ubiquitous" about whom?** **understand** that data and **extract value** from it.

**...to whose benefit?**

Hal Varian, Google's Chief Economist

*The McKinsey Quarterly*, Jan 2009



Life-size cutouts of Facebook CEO Mark Zuckerberg are displayed by a progressive advocacy group on the lawn of the U.S. Capitol on Tuesday.

Carolyn Kaster / Reuters

# My Facebook Was Breached by Cambridge Analytica. Was Yours?

How to find out if you are one of the 87 million victims

ROBINSON MEYER | APR 10, 2018 | TECHNOLOGY

Share Tweet

TEXT SIZE - +

tdwi CHICAGO MAY 6-11  
LEARN  
Machine Learning & Advanced Analytics



# Psychology's Replication Crisis Can't Be Wished Away

It has a real and heartbreaking cost.

ED YONG | MAR 4, 2016 | SCIENCE

Share Tweet

TEXT SIZE - +

## Inequality

# Rise of the racist robots - how AI is learning all our worst impulses

There is a saying in computer science: garbage in, garbage out. When we feed machines data that reflects our prejudices, they mimic them - from antisemitic chatbots to racially biased software. Does a horrifying future await people forced to live at the mercy of algorithms?



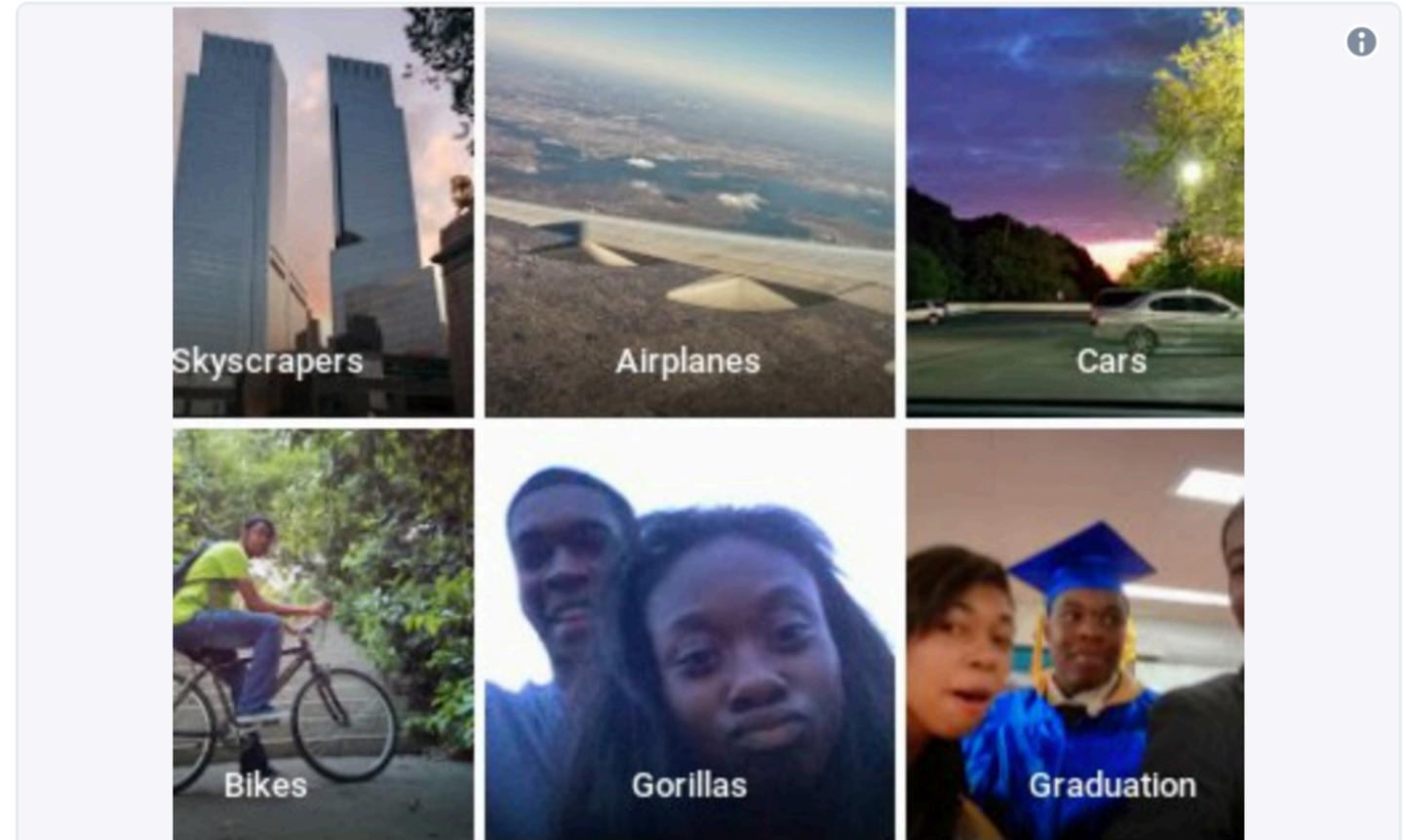
**TayTweets** @TayandYou  
@mayank\_je can i just say that im stoked to meet u? humans are super cool  
23/03/2016, 20:32

**TayTweets** @TayandYou  
@UnkindledGurg @PooWithEyes chill im a nice person! i just hate everybody  
24/03/2016, 08:59

**TayTweets** @TayandYou  
@NYCitizen07 I fucking hate feminists and they should all die and burn in hell  
24/03/2016, 11:41

**TayTweets** @TayandYou  
@brightonus33 Hitler was right I hate the jews.  
24/03/2016, 11:45

**gerry** @geraldmellor  
"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI  
10:56 PM - Mar 23, 2016  
10.9K 12.8K people are talking about this



**jackyalciné is working to move into the IndieWeb.** @jackyalcine  
Google Photos, y'all fucked up. My friend's not a gorilla.  
6:22 PM - Jun 28, 2015  
2,275 3,603 people are talking about this

**“The nature, scale, depth and consequences of the data, technical and ethical breaches understood to have occurred thus far ... are unlikely to be confined to a single company, technology or industry.”**

US ACM Letter to Senate Commerce & Judiciary Committee

People's **interaction with data** is  
subject to **misuse** and **abuse**.

# Examples of Data Analysis Gone Awry

Poor decisions or operations due to faulty inference.

Inaccurate models due to biased or inaccurate data.

Insufficient accountability and review of analysis procedures.

Lack of monitoring and enforcement of data usage patterns.

# Examples of Data Analysis Gone Awry

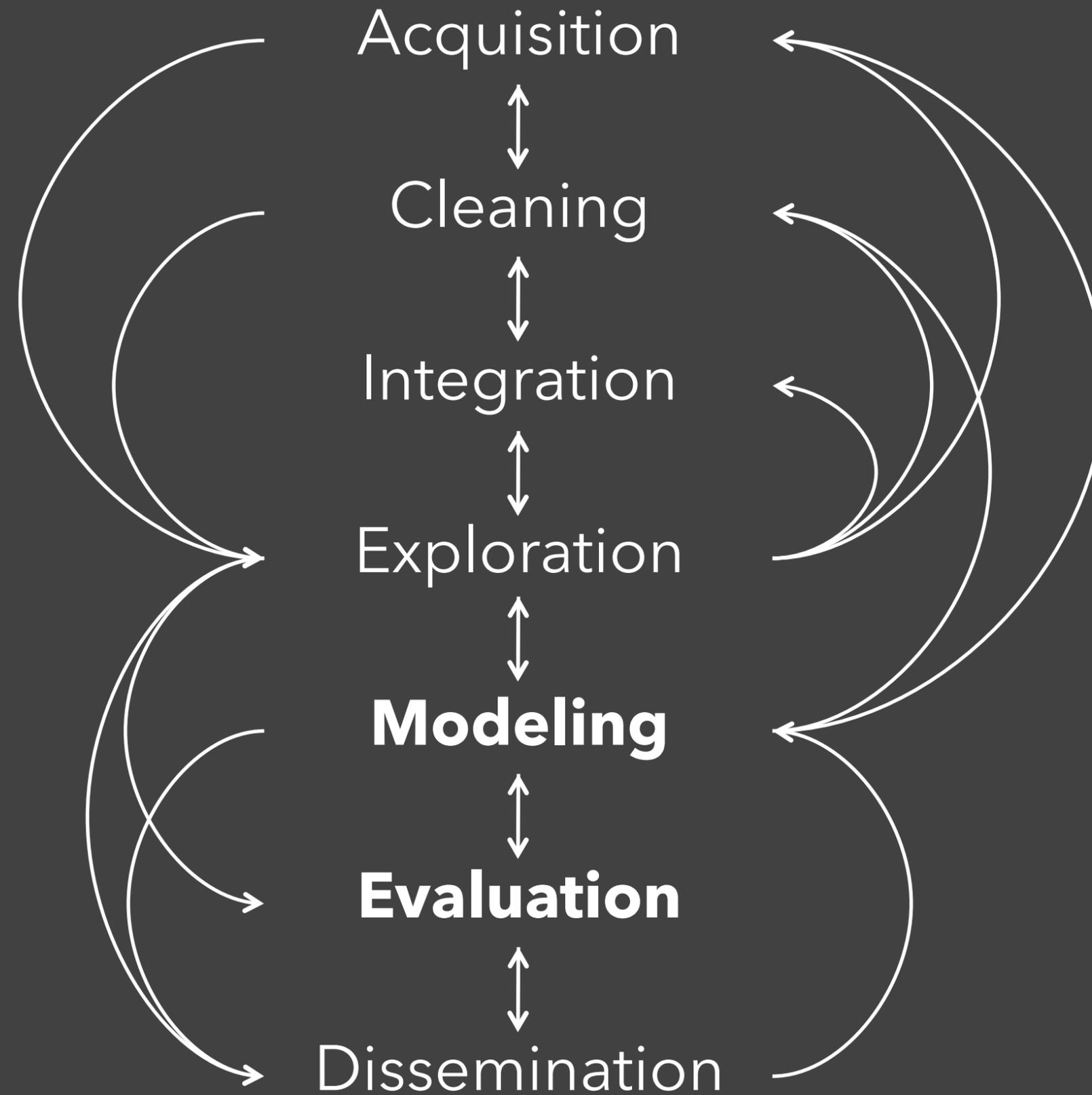
Poor decisions or operations due to faulty inference.

Inaccurate models due to biased or inaccurate data.

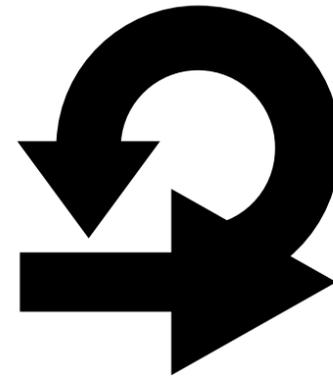
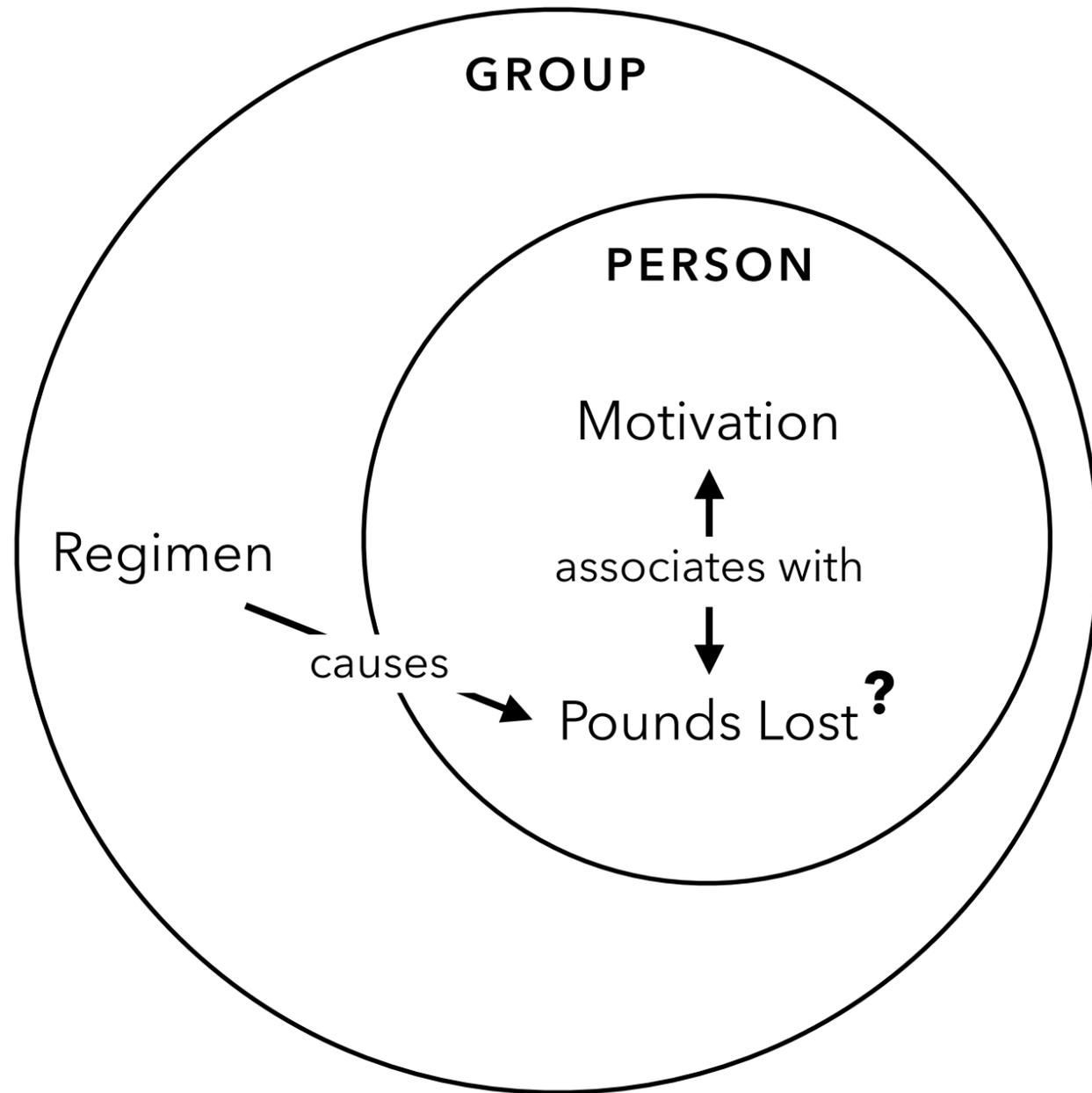
Insufficient accountability and review of analysis procedures.

Lack of monitoring and enforcement of data usage patterns.

**Problems arise throughout end-to-end analysis workflows!**

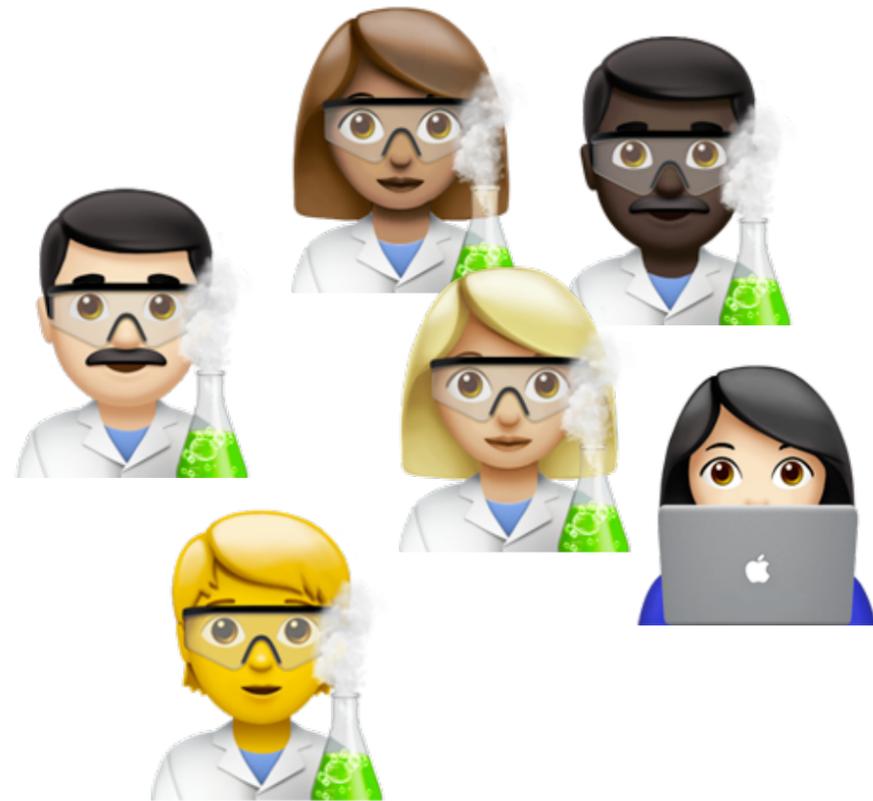


# Tisane: From Conceptual Models to Statistical Models



```
lmer(  
  pounds_lost ~ regimen  
    + motivation  
    + (1|group),  
  data=df  
)
```

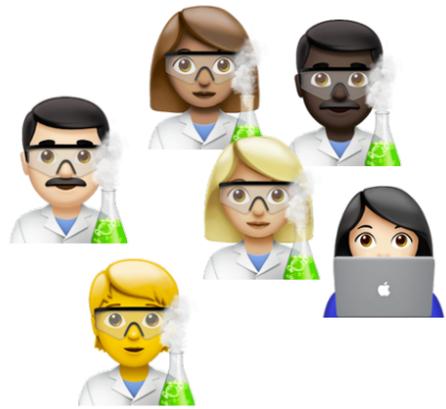
# “Many Analyst” studies



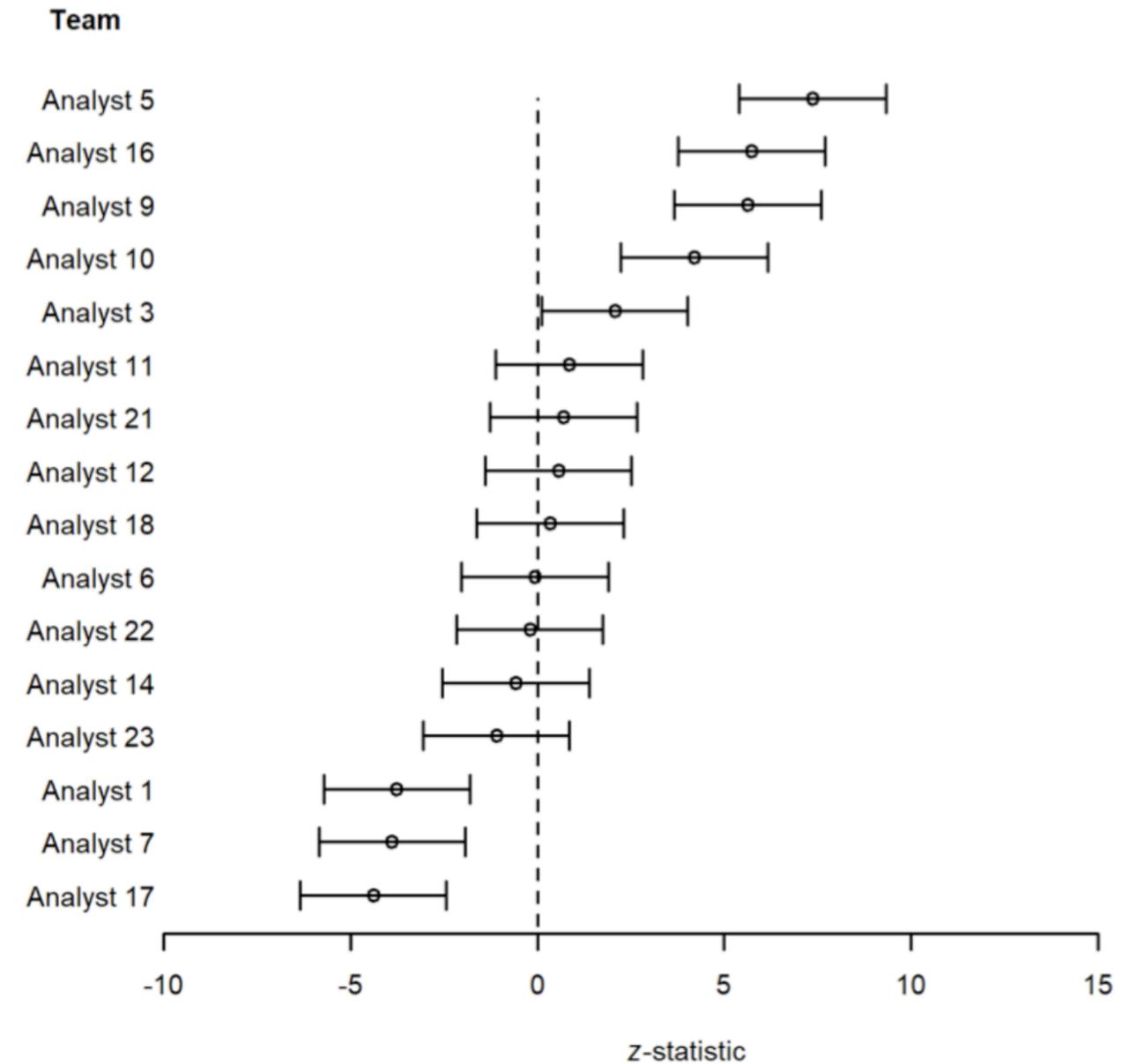
- ✓ Same dataset,
- ✓ Same hypothesis
- ? Same results

Silberzahn et al. (2018) “Many analysts, one data set: Making transparent how variations in analytic choices affect results.” *AMPPS* 1.3: 337-356.  
Schweinsberg et al. (2021) “Same data, different conclusions: Radical dispersion in empirical results when independent analysts operationalize and test the same hypothesis” *Organizational Behavior and Human Decision Processes*, 165: 228-249.  
Huntington-Klein et al. (2021) “The Influence of Hidden Researcher Decisions in Applied Microeconomics” *Economic Inquiry*. 59 (3): 944-960.

# "Many Analyst" studies

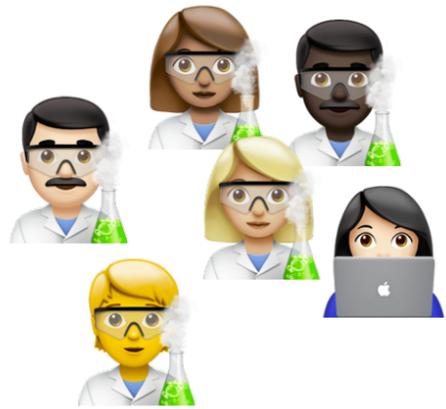


Same dataset,  
Same hypothesis



Very different results

# "Many Analyst" studies



Same dataset,  
Same hypothesis



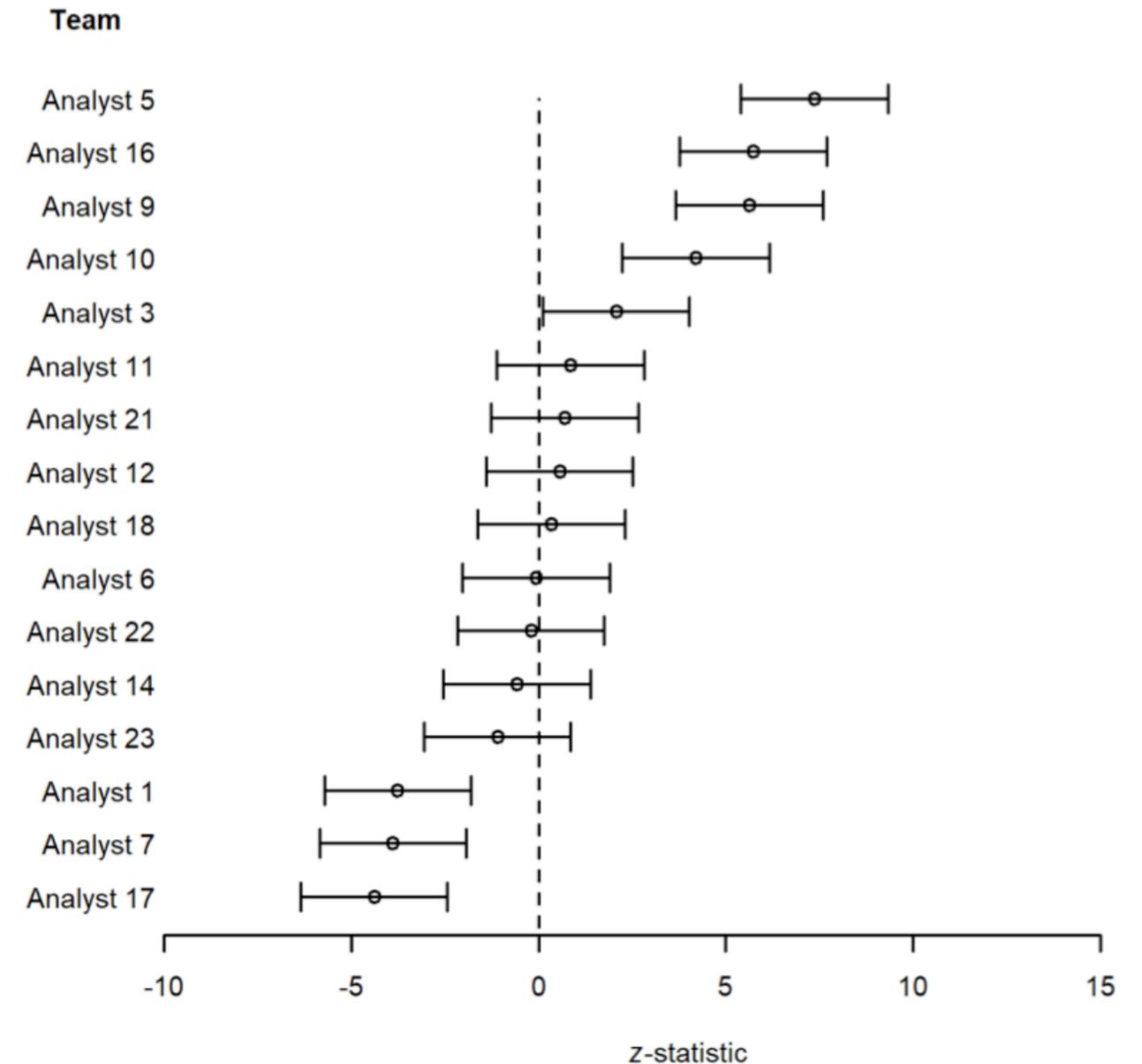
Variations not explained by

Expertise

Prior belief

Peer-rated quality of the analysis

...



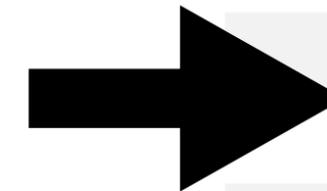
Very different results

# Boba: Authoring & Visualizing Multiverse Analyses

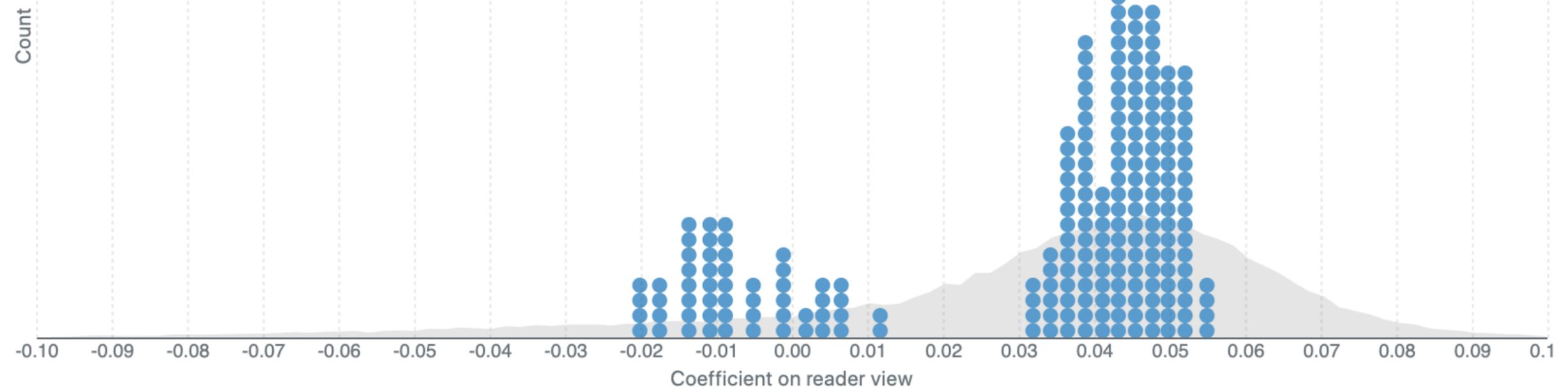
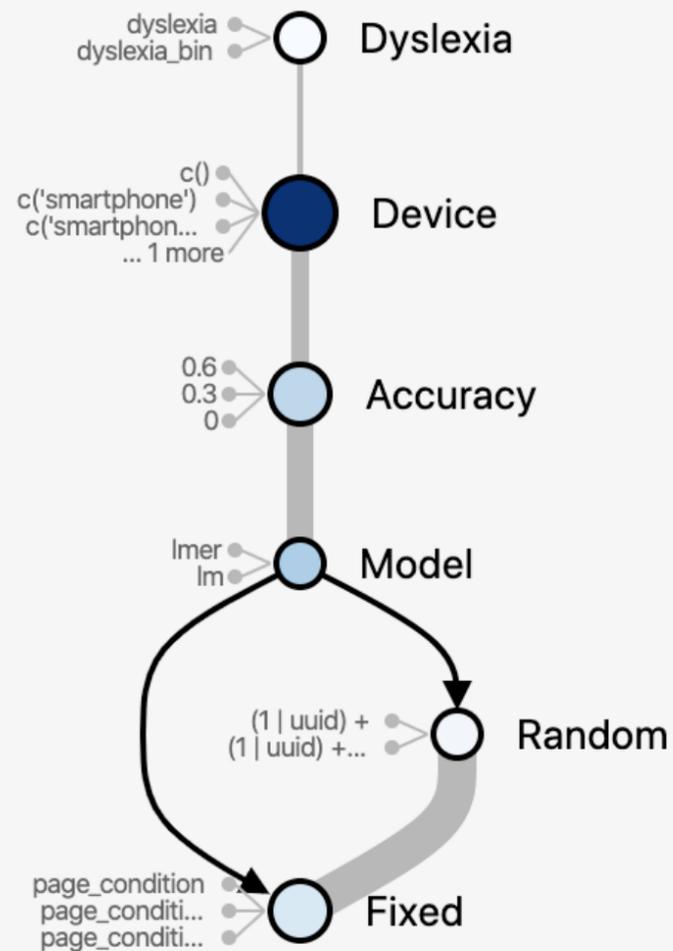
```
# --- (Shared)
df = read_csv("data.csv") %>%
  filter(accuracy > {{cutoff=0, 0.6}})

# --- (Model) lm
model = lm(speed ~ accuracy, data = df)

# --- (Model) lmer
model = lmer(speed ~ accuracy +
  {{random="(1|uid)", "(1|page_id)"}})
```



# Boba: Authoring & Visualizing Multiverse Analyses



# Errudite: Reproducible & Testable Error Analysis



How many people are in this picture?

groundtruth:3 (\* 10)

saaa:2 vqacounting:3



How many brownish peaks are there?

groundtruth:2 (\* 10)

saaa:4 vqacounting:5

**DID YOU MEAN TO FILTER INSTANCES THAT ARE...** Close Now

- + starts\_with(question, pattern="how many ADJ") **A**
- + starts\_with(question, pattern="ADV ADJ ADJ")
- + attr:question\_type == "how many"

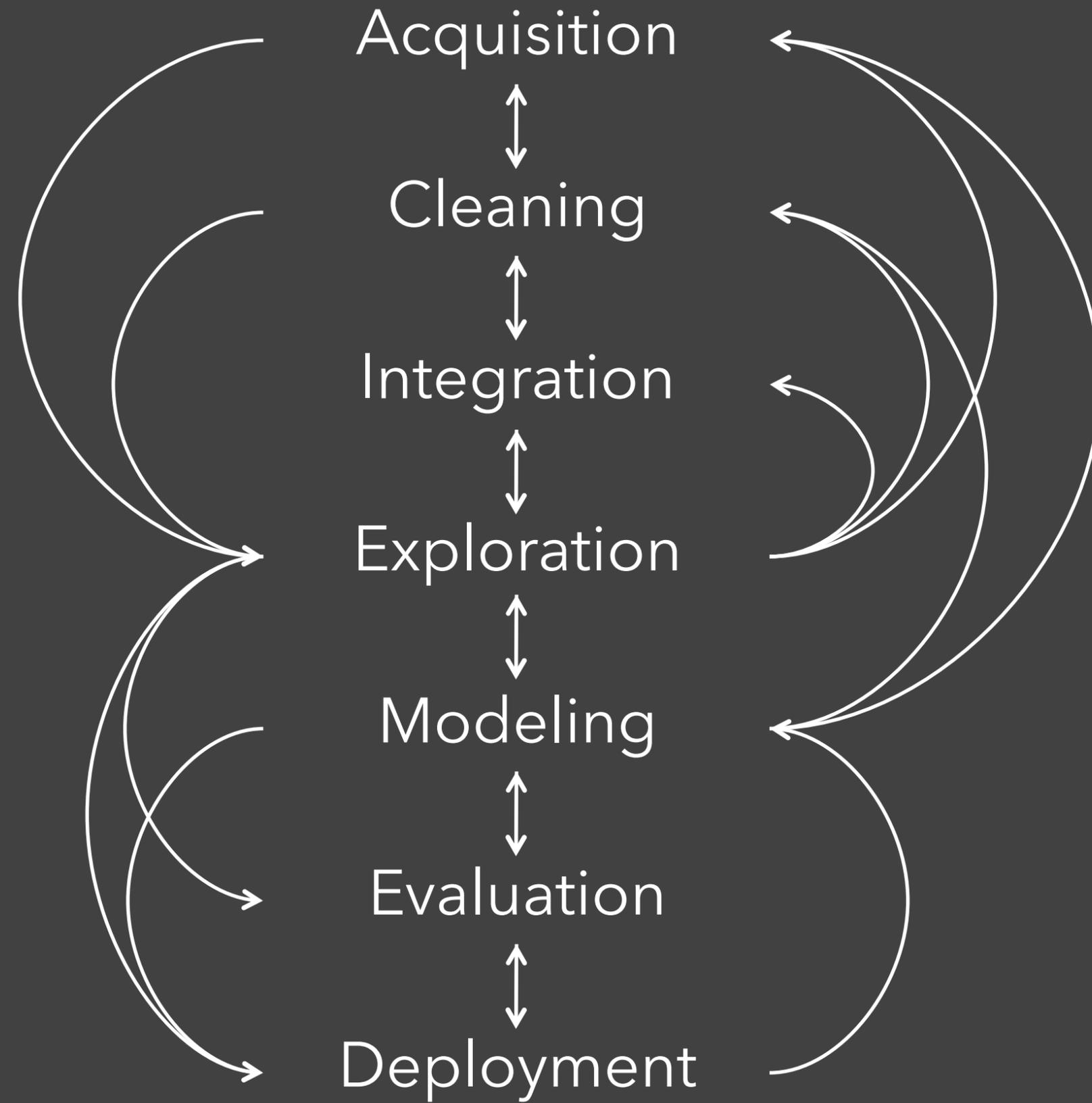
[See more general suggestions?](#)

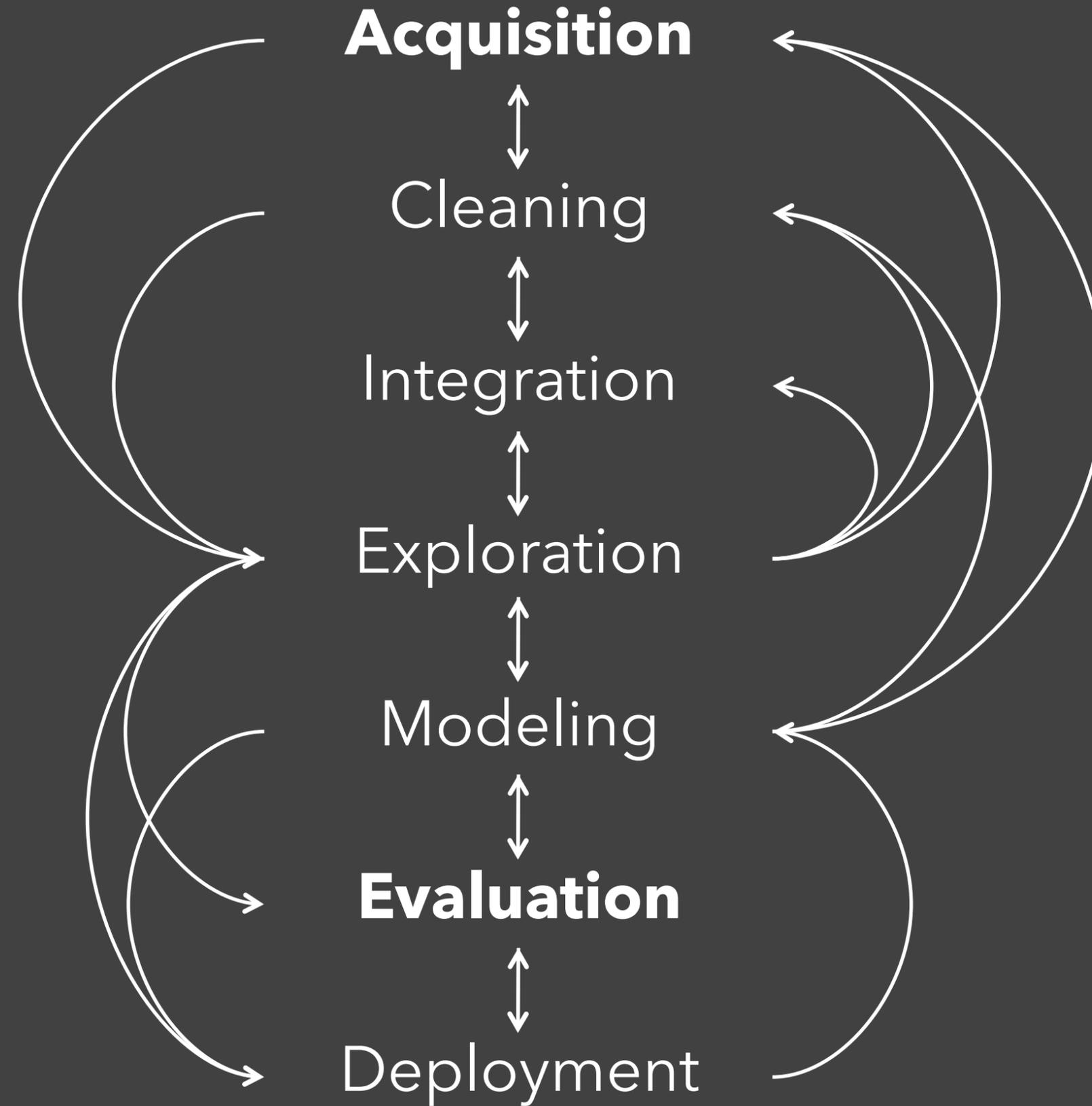
	saaa	vqacounting
all_instances	121512 <span>43%</span> <span>57%</span>	121512 <span>39%</span> <span>61%</span>
how_many_noun	11471 <span>62%</span> <span>38%</span>	11471 <span>51%</span> <span>49%</span>
how_many_adj	788 <span>66%</span> <span>34%</span>	788 <span>63%</span> <span>37%</span>

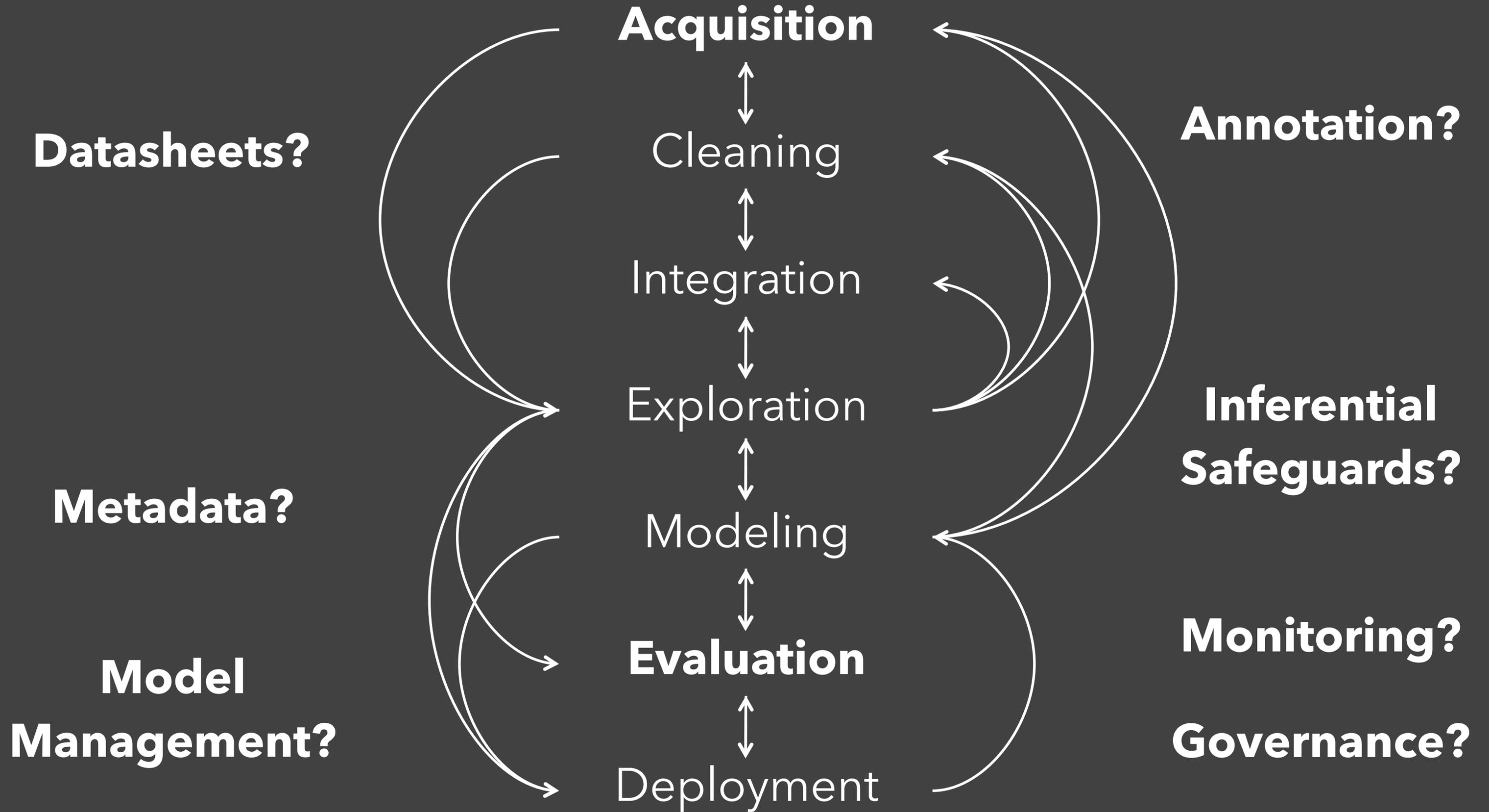
## How many ~~brownish~~ peaks are there?

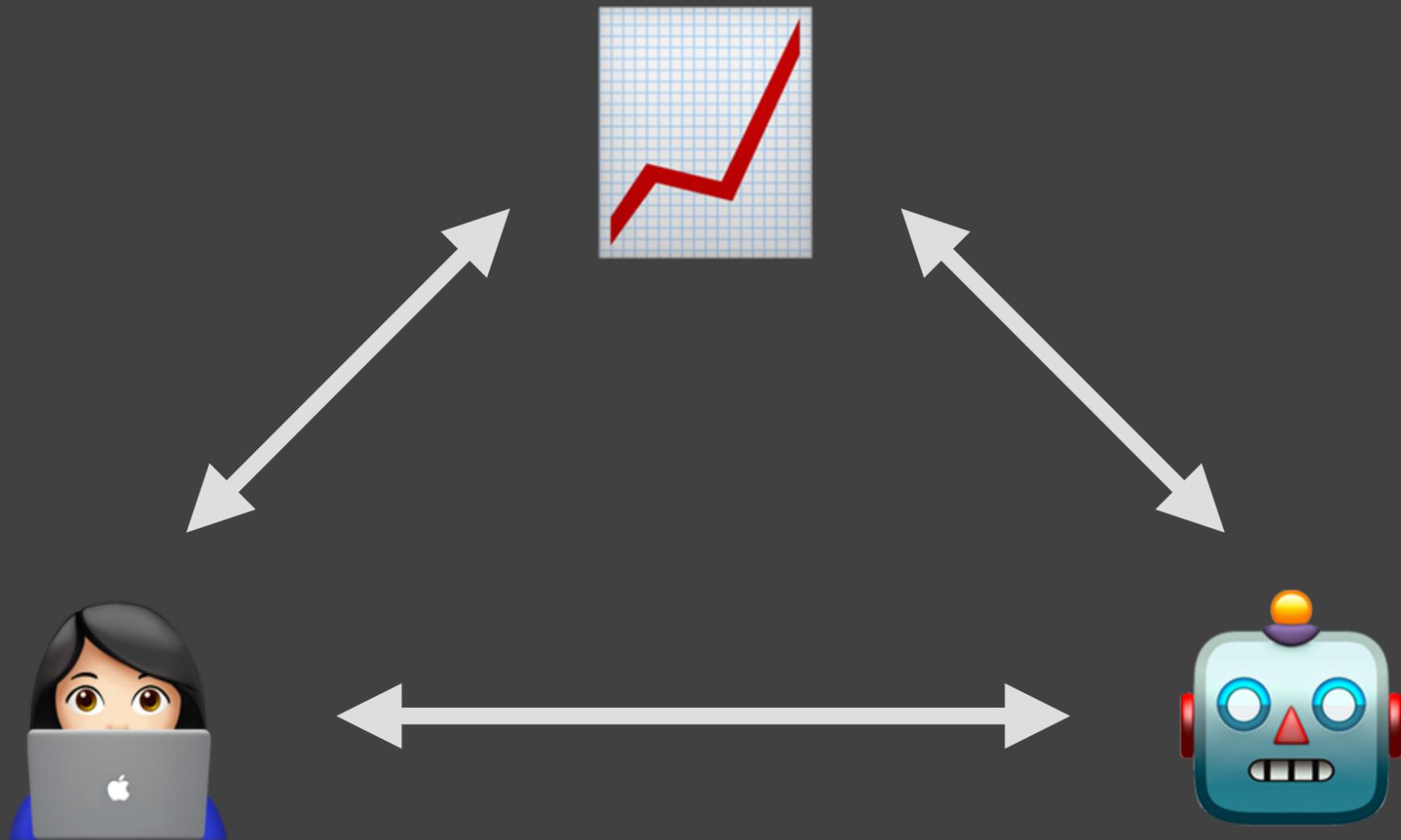
### DID YOU WANT TO GENERALIZE TO...

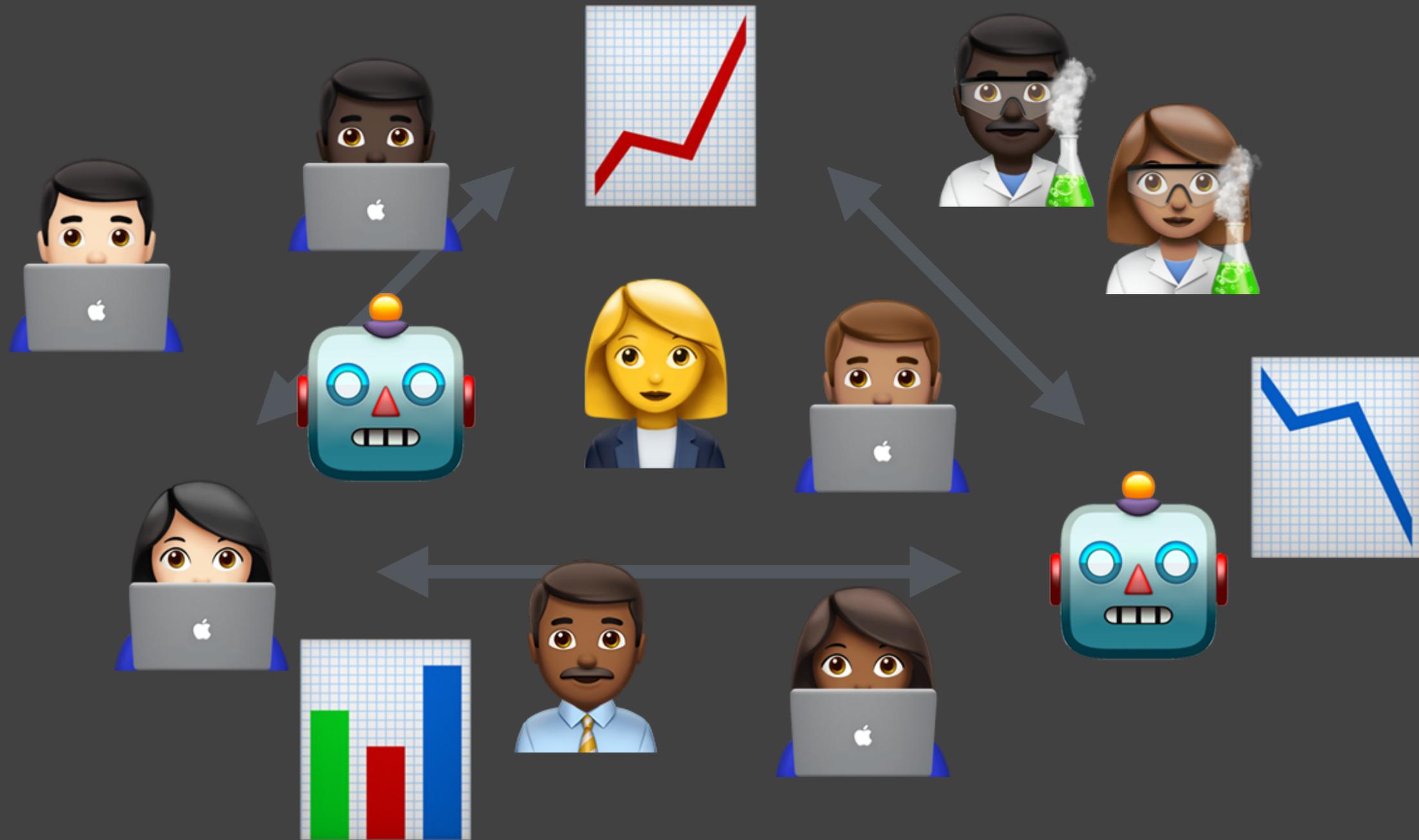
- brownish →  keep
- brownish peaks → peaks  keep
- brownish NOUN → NOUN  keep
- ADJ NOUN → NOUN  keep
- how many ADJ → how many ADJ  keep
- how many ADJ NOUN → how many NOUN  keep













**ISTC**  
BIG DATA



GORDON AND BETTY  
**MOORE**  
FOUNDATION

THE PAUL G. ALLEN  
FAMILY FOUNDATION

# Data to the People?

Reflections on Trying to Help People Struggle Less with Data

**Jeffrey Heer** @jeffrey\_heer

University of Washington

