

The background of the slide features a large, semi-transparent watermark of the University of California Berkeley seal. The seal is circular and contains a five-pointed star at the top, an open book in the center, and a banner at the bottom with the motto "LET THERE BE LIGHT". The text "THE UNIVERSITY OF CALIFORNIA BERKELEY" is written around the perimeter of the seal, and the year "1868" is at the bottom. The seal is rendered in a light gray color.

Data Wrangling and Predictive Interaction

Brief lessons from a decade of R&D

Joe Hellerstein

Data Platforms



Databases



Spreadsheets



Log Files



IoT Sensors



Apps

“The hardest part of AI is the **data wrangling.**”

—SWAMI SIVASUBRAMANIAN, VP AWS MACHINE LEARNING

#reInvent2018

80%

“It’s impossible to overstate this:
80% of the work in any data project
is in **cleaning the data.**”

— DJ Patil, Data Jujitsu, O’Reilly Media 2012

Analysis



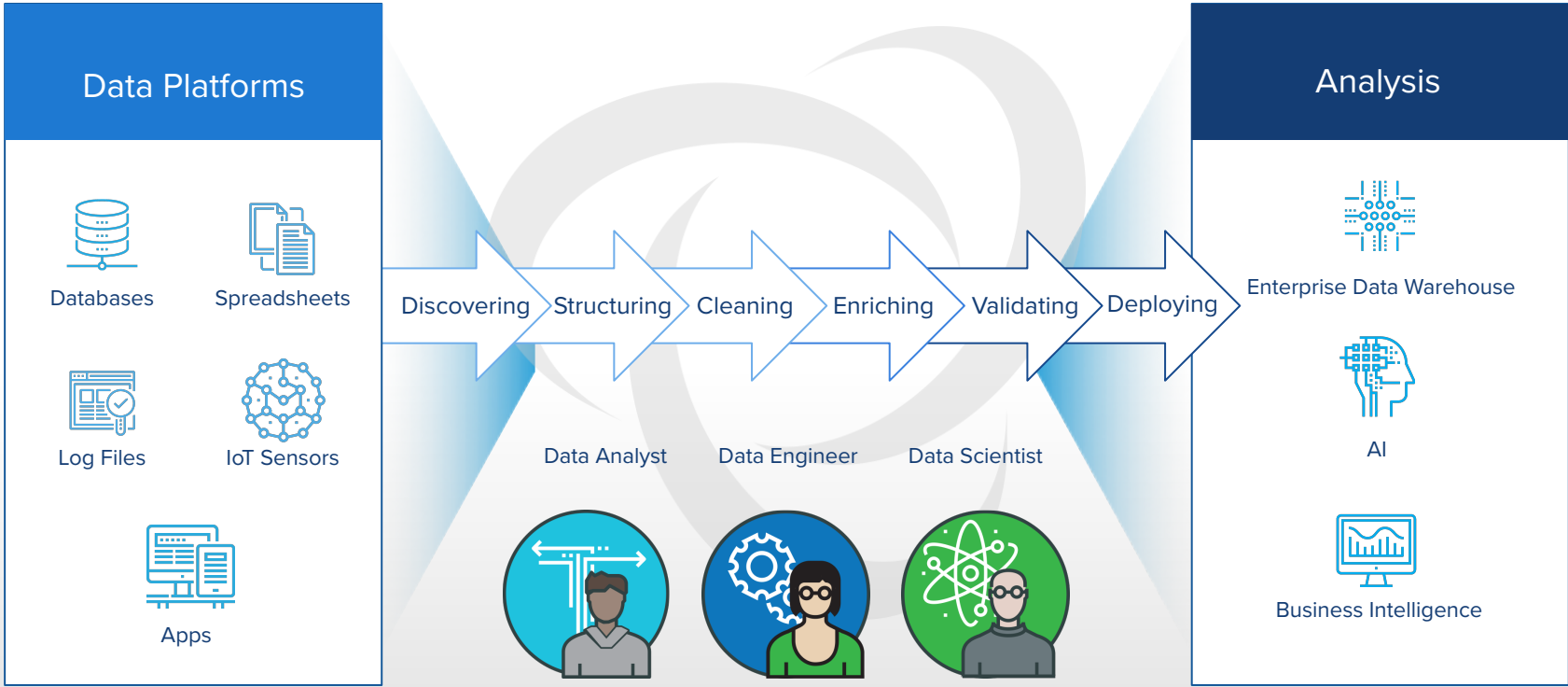
Enterprise Data Warehouse



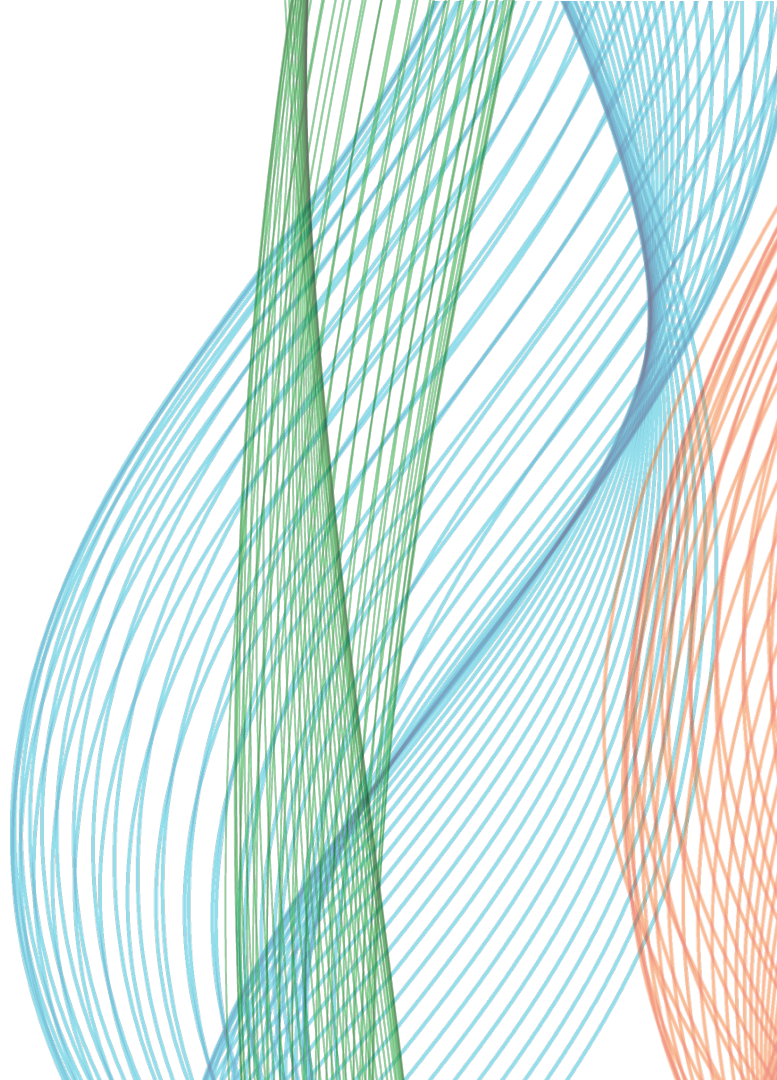
AI



Business Intelligence



Mini Demo



Methodology: Interview Studies!

IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS, VOL. 18, NO. 12, DECEMBER 2012

2917

Enterprise Data Analysis and Visualization: An Interview Study

Sean Kandel, Andreas Paepcke, Joseph M. Hellerstein, and Jeffrey Heer

Abstract—Organizations rely on data analysts to model customer engagement, streamline operations, improve production, inform business decisions, and combat fraud. Though numerous analysis and visualization tools have been built to improve the scale and efficiency at which analysts can work, there has been little research on how analysis takes place within the social and organizational context of companies. To better understand the enterprise analysts' ecosystem, we conducted semi-structured interviews with 35 data analysts from 25 organizations across a variety of sectors, including healthcare, retail, marketing and finance. Based on our interview data, we characterize the process of industrial data analysis and document how organizational features of an enterprise impact it. We describe recurring pain points, outstanding challenges, and barriers to adoption for visual analytic tools. Finally, we discuss design implications and opportunities for visual analysis research.

Index Terms—Data, analysis, visualization, enterprise.



IEEE Vis Test of Time Award 2022

A Decade of R & D

DataWrangler^{alpha}

 TRIFACTA

Google Cloud  Dataprep
by TRIFACTA

See also:

- OpenRefine
- MS PowerQuery
- Tableau Dataprep
- AWS Databrew
- Zoho Dataprep
- ...

Wrangler: Interactive Visual Specification of Data Transformation Scripts

Sean Kandel^{*}, Andreas Paepcke^{*}, Joseph Hellerstein[†] and Jeffrey Heer^{*}
^{*} Stanford University [†] University of California, Berkeley
skandel, paepcke, jheer@cs.stanford.edu hellerstein@cs.berkeley.edu

ABSTRACT

Though data analysis tools continue to improve, analysts still expend an inordinate amount of time and effort manipulating data and assessing data quality issues. Such “data wrangling” regularly involves reformatting data values or layout, correcting erroneous or missing values, and integrating multiple data sources. These transforms are often difficult to specify and difficult to reuse across analysis tasks, teams, and tools. In response, we introduce *Wrangler*, an interactive system for creating data transformations. *Wrangler* combines these capabilities of structured data with

data warehousing projects [4]. Such “data wrangling” often requires writing idiosyncratic scripts in programming languages such as Python and Perl, or extensive manual editing using interactive tools such as Microsoft Excel. Moreover, this hurdle discourages many people from working with data in the first place. Sadly, when it comes to the practice of data analysis, “the tedium is the message.”

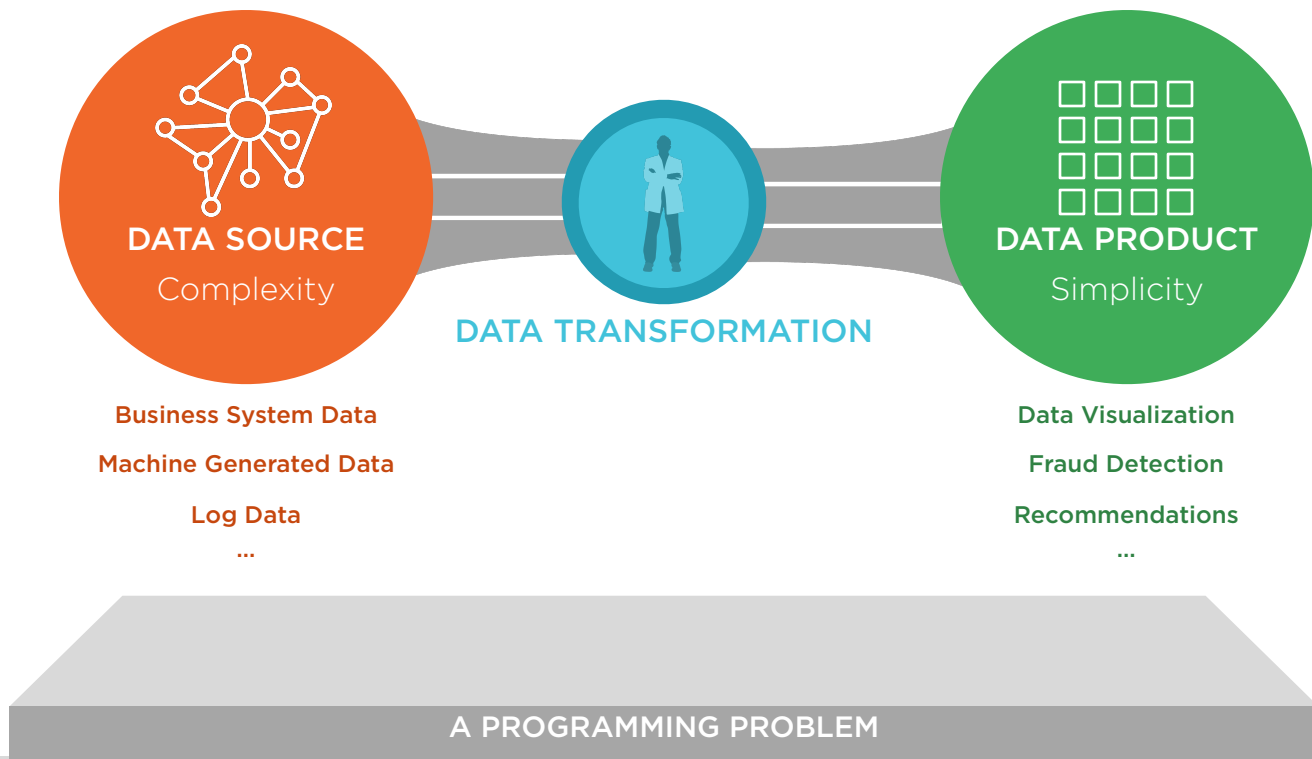
Part of the problem is that reformatting and validating data requires transforms that can be difficult to specify and evaluate. For instance, analysts often edit data into meaning-

CHI 2011

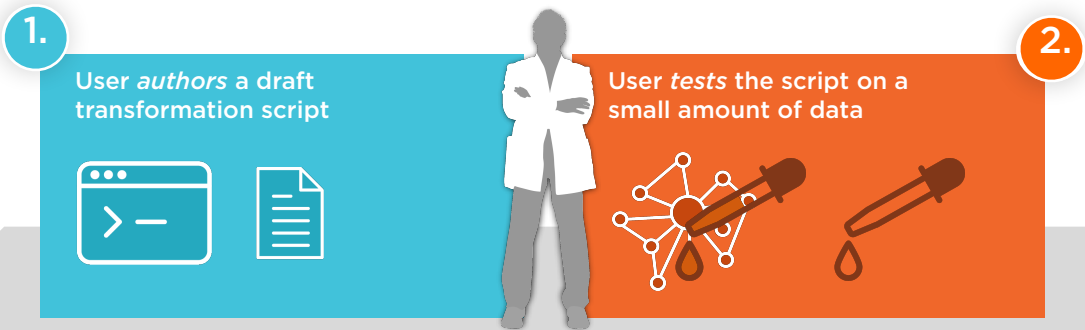
alteryx



THE DATA TRANSFORMATION PROBLEM

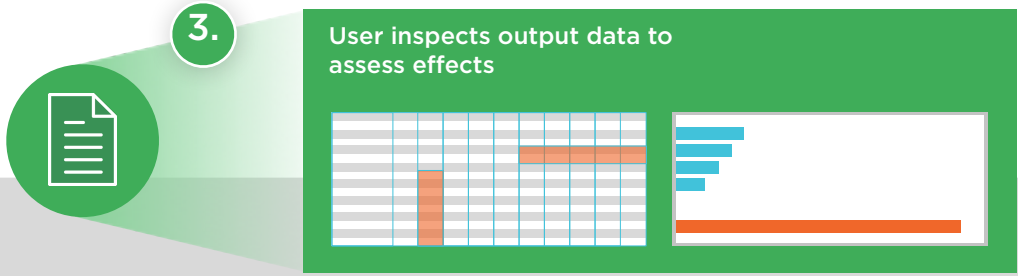


TRADITIONAL ROLE OF VISUALIZATION



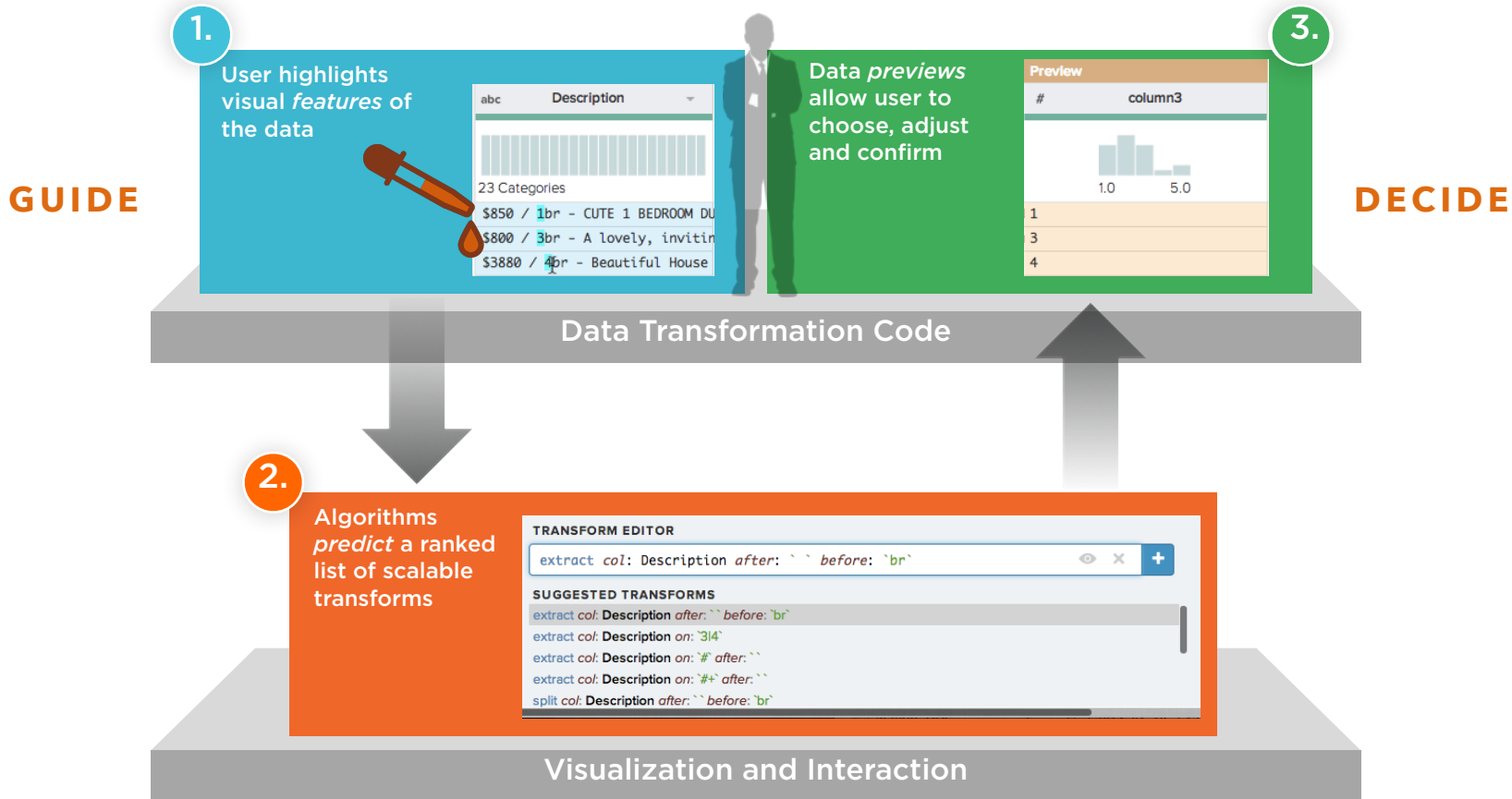
Python pandas,
R dplyr,
SQL, etc...

Data Transformation Code

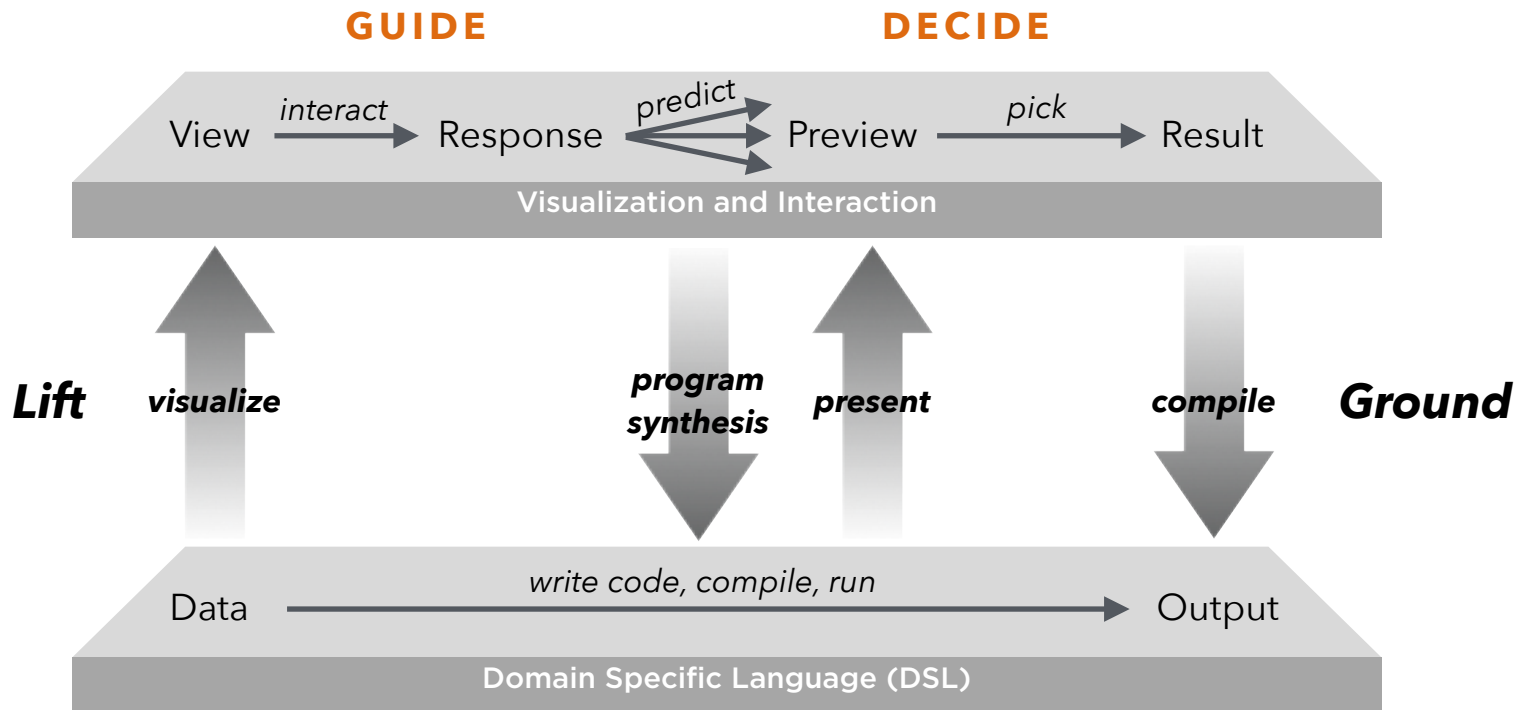


Visualization and Interaction

PREDICTIVE INTERACTION



PREDICTIVE INTERACTION



Predictive Interaction for Data Transformation

Jeffrey Heer
U. Washington & Trifacta Inc.
jheer@uw.edu

Joseph M. Hellerstein
UC Berkeley & Trifacta Inc.
hellerstein@berkeley.edu

Sean Kandel
Trifacta Inc.
skandel@trifacta.com

ABSTRACT

The human work involved in data transformation represents a major bottleneck for today's data-driven organizations. In response, we present Predictive Interaction, a framework for interactive systems that shifts the burden of technical specification from users to algorithms, while preserving human guidance and expressive power.

1. INTRODUCTION

Many of today's data management challenges stem from the increasing variety of scenarios where data is being exploited. This

data-to-metadata transformations like Pivot/Unpivot), cleaning (e.g., standardization, entity resolution, dictionary management), enrichment (e.g., joins and references), and distillation (e.g., sampling, filtering, aggregation, windowing).

There are a number of recurring themes in previous work on data transformation. One is to develop new *user interfaces* for graphical specification of queries and transforms [11, 23, 26, 32], and visualization of outputs [3, 19]. A second theme is to innovate at the *data management* layer, with domain-specific languages that are well-suited to data transformation and can be executed in a scalable, high-performance manner [8, 18, 28]. And despite the

CIDR 2015

Surprising Lessons

→ UX matters far more than prediction quality

Details ✕

📞 phone_num ⋮

Unique Values

(240)966-1418	3
(301)133-0554	3
(620)784-8038	3
(915)228-6301	3
209-715-7824	3

[Show more values...](#)

Patterns

\({digit}{3}\){digit}{3}-{digit}{4}	8,750
{digit}{3} {digit}{3} {digit}{4}	2,190

[Show pattern details...](#)

Suggestions

Delete columns

phone_num

Rename

Rename phone_num to 'phone_num'

Group by

Create new columns from COUNT() grouped by phone_num

Create a new column

COUNT() grouped by phone_num



Set

Set phone_num to IFMISMATCHED(\$col, ['Phone'], '')

Surprising Lessons

- UX matters far more than prediction quality
- Embrace “read-centric” rendering of code

SUGGESTED TRANSFORMS



Source	Preview
abc	abc
Screen_Detail	Screen...
	
6 Categories	2 Categories
31 adtam_name=utarget1&adtam_source=dynamic&adtam_size=180x150	dynamic
32 adtam_name=holidaypromo1&adtam_source=dynamic&adtam_size=300x250	dynamic
33 adtam_name=utarget1&adtam_source=dynamic&adtam_size=180x150	dynamic
34 adtam_name=holidaypromo2&adtam_source=mobile&adtam_size=240x400	mobile

Surprising Lessons

- UX matters far more than prediction quality
- Embrace “read-centric” rendering of code

SUGGESTED TRANSFORMS

```
extract col: Screen_Detail on: /(?!<=adtam_source)[^\&]*(?=\&)/
extract col: Screen_Detail on: /(?!<=)[^\&]*(?=\&)/ limit: 2
extract col: Screen_Detail on: /(?!<=)[a-z]+/ limit: 2
extract col: Screen_Detail on: /[a-z]+/ limit: 4
```



Source	Preview
abc	abc Screen...
	
6 Categories	2 Categories
31 adtam_name=utarget1&adtam_source=dynamic&adtam_size=180x150	dynamic
32 adtam_name=holidaypromo1&adtam_source=dynamic&adtam_size=300x250	dynamic
33 adtam_name=utarget1&adtam_source=dynamic&adtam_size=180x150	dynamic
34 adtam_name=holidaypromo2&adtam_source=mobile&adtam_size=240x400	mobile

Surprising Lessons

- UX matters far more than prediction quality
- Embrace “read-centric” rendering of code

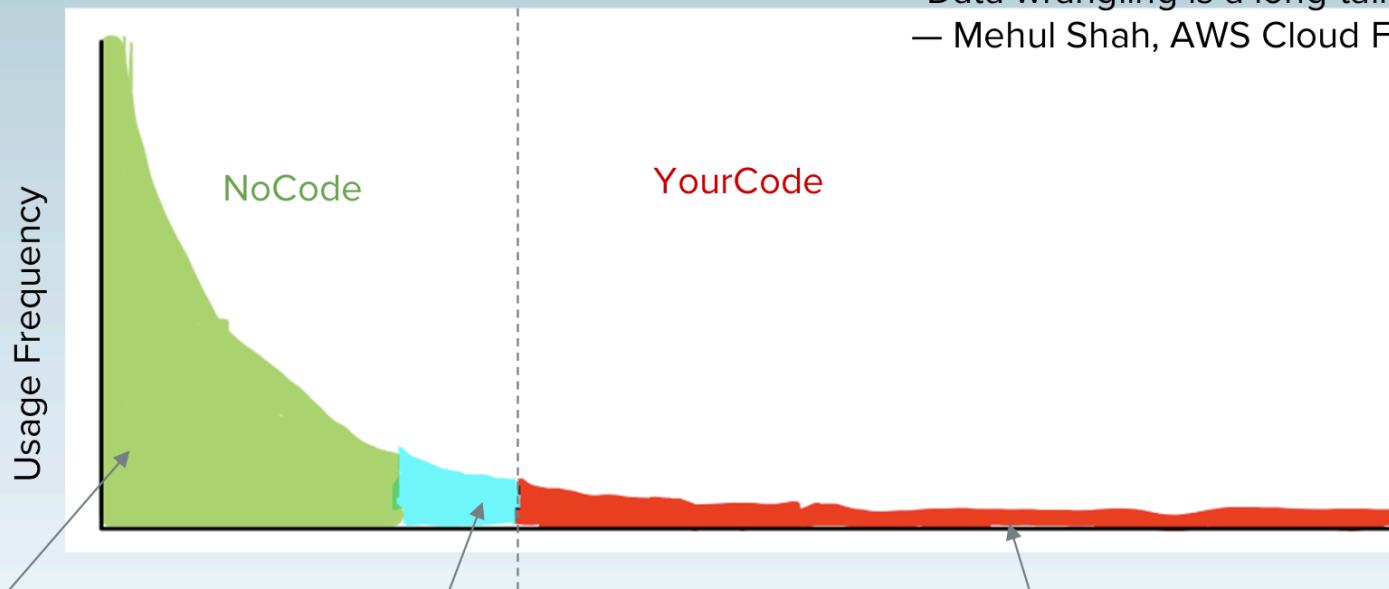
SUGGESTED TRANSFORMS

```
extract col: Screen_Detail after: `adtam_source=` before: `&`  
extract col: Screen_Detail limit: 2 after: `=` before: `&`  
extract col: Screen_Detail on: `{lower}+` limit: 2  
extract col: Screen_Detail on: `{lower}+` limit: 4
```

Source	Preview
abc	abc Screen...
	
6 Categories	2 Categories
31 adtam_name=utarget1&adtam_source=dynamic&adtam_size=180x150	dynamic
32 adtam_name=holidaypromo1&adtam_source=dynamic&adtam_size=300x250	dynamic
33 adtam_name=utarget1&adtam_source=dynamic&adtam_size=180x150	dynamic
34 adtam_name=holidaypromo2&adtam_source=mobile&adtam_size=240x400	mobile

Surprising Lessons

“Data wrangling is a long-tail business.”
— Mehul Shah, AWS Cloud Formation & Glue



Tasks that you can do point-and-click in Trifacta

Tasks that you can do in Trifacta

Tasks that you can't do in Trifacta

Surprising Lessons

- UX matters far more than prediction quality
- Embrace read-only rendering of code
- Solve for Low Code + Your Code



Selected Academic Papers

Raman, Vijayshankar, and Joseph M. Hellerstein. "Potter's wheel: An interactive data cleaning system." In *VLDB*, 2001.

Kandel, Sean, Andreas Paepcke, Joseph Hellerstein, and Jeffrey Heer. "Wrangler: Interactive visual specification of data transformation scripts." In *SIGCHI*, 2011.

Kandel, Sean, Andreas Paepcke, Joseph M. Hellerstein, and Jeffrey Heer. "Enterprise data analysis and visualization: An interview study." *IEEE Trans. on Vis. and Comp. Graphics* 18, no. 12 (2012).

Heer, Jeffrey, Joseph M. Hellerstein, and Sean Kandel. "Predictive Interaction for Data Transformation." In *CIDR*, 2015.

Hellerstein, Joseph M., Jeffrey Heer, and Sean Kandel. "Self-Service Data Preparation: Research to Practice." *IEEE Data Eng. Bull.* 41, no. 2 (2018).

@joe_hellerstein
hellerstein@berkeley.edu

