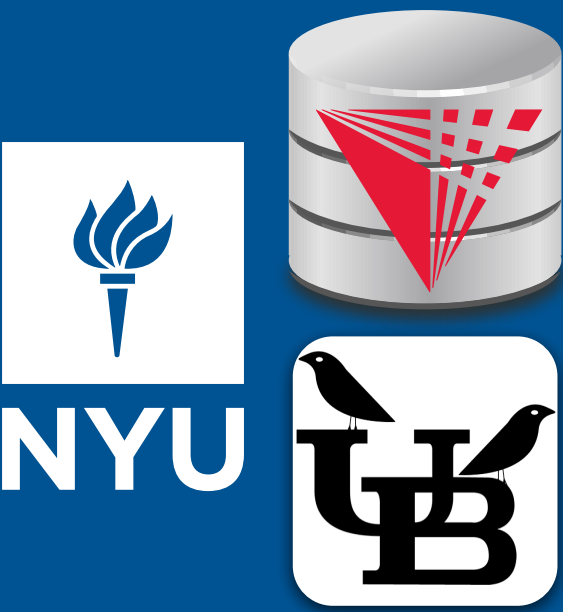


Data Debugging and Exploration with Vizier

Mike Brachmann¹, Carlos Bautista², Sonia Castelo², Su Feng³, Juliana Freire², Boris Glavic³, Oliver Kennedy¹, Heiko Müller², Rémi Rampin², William Spoth¹, Ying Yang⁴

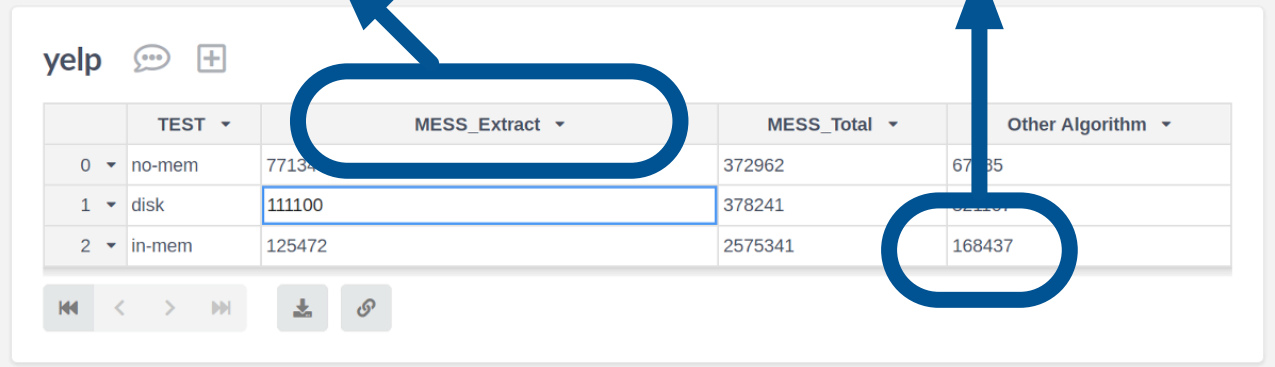
¹University at Buffalo, ²New York University, ³Illinois Institute of Technology, ⁴Oracle



Painless Data Ingest

Detects Headers

Infers Types



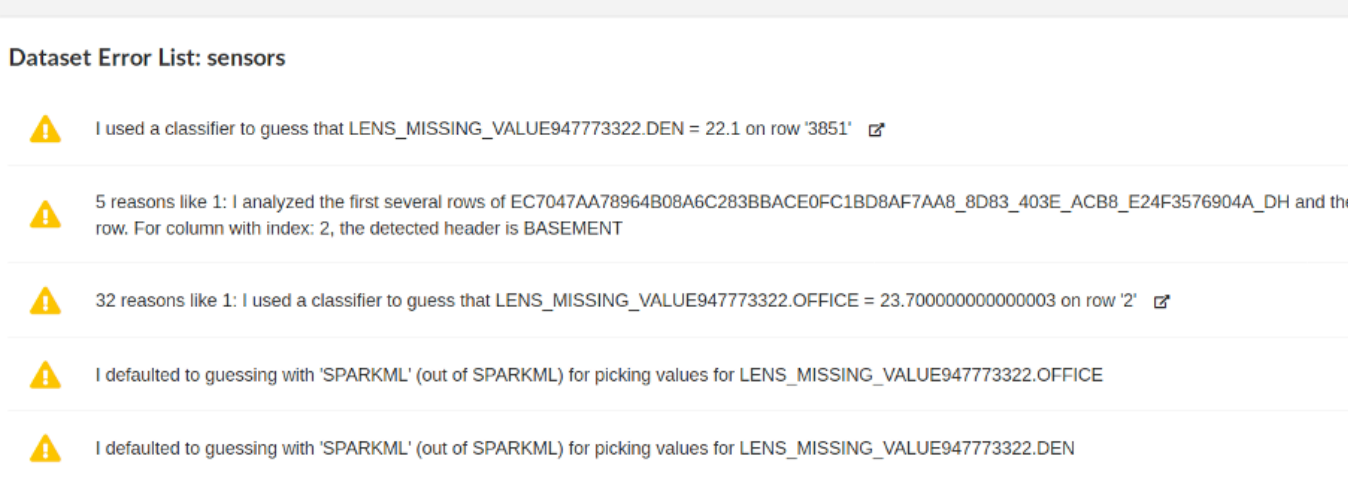
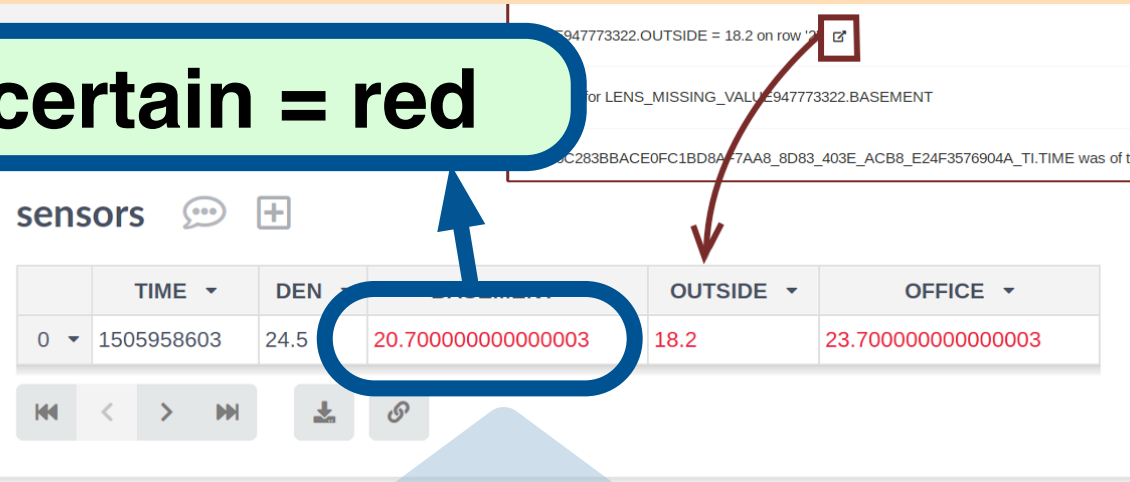
Supports your favorite data format and source

Data Testing & Debugging

Automated Curation/Cleaning Ops

Error tracking & annotations

uncertain = red



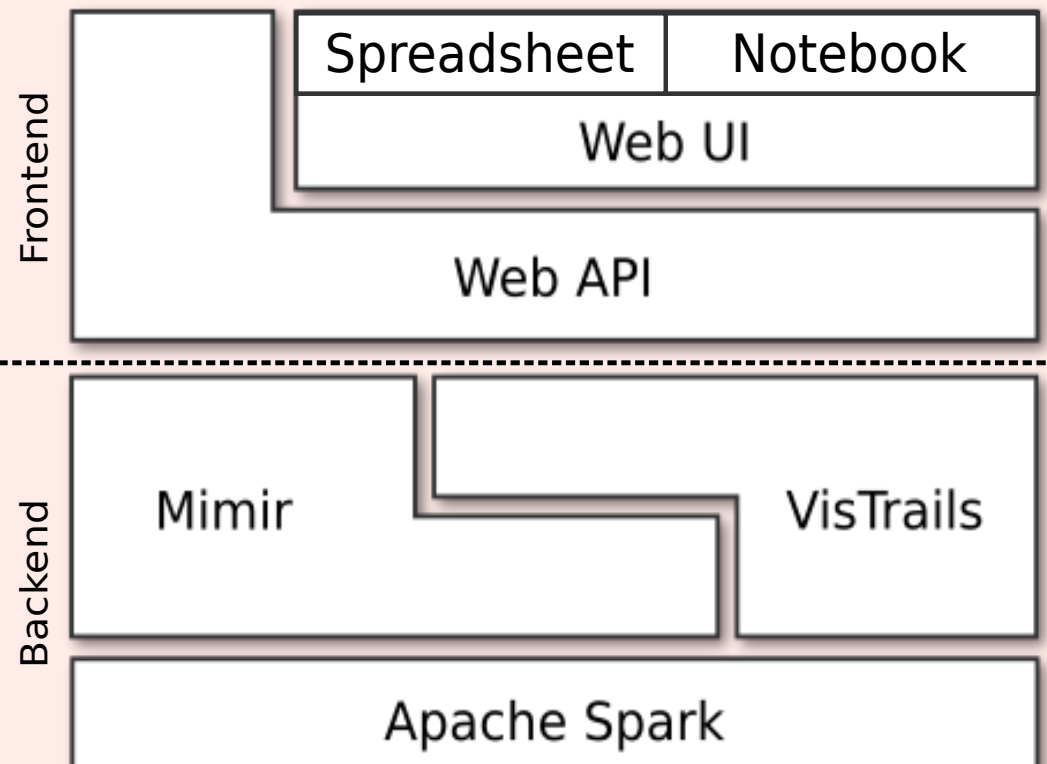
What is Vizier?

A next generation data science platform that...

- supports iterative construction of pipelines
- provides expressive interaction paradigms without sacrificing ease-of-use
- aids data debugging & testing
- scales
- provides tracebility and provenance
- aids collaboration and reuse

data-centric!

Architecture



Versioning & Collaboration

Automatic commit on edit!

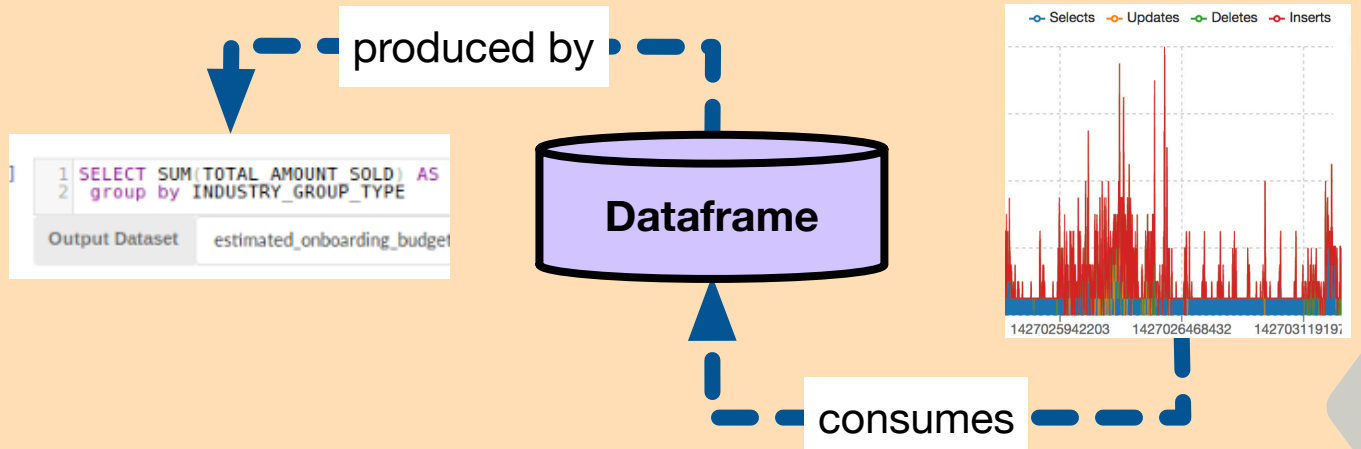
Browse old versions and create branches

Reuse your workflows!

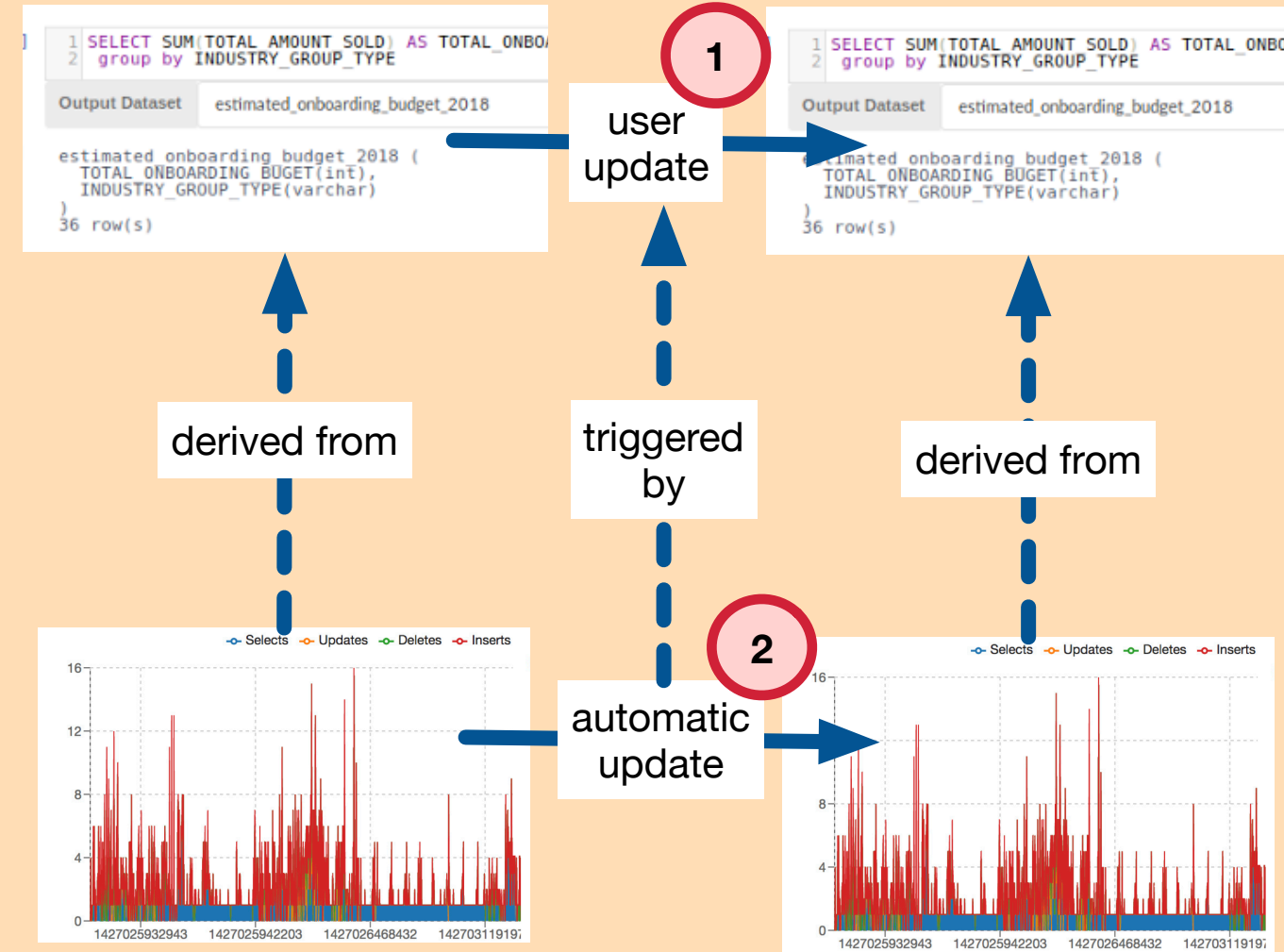
Share links to any version of ... a dataset ... a notebook ... a plot

Data-Centric

Interoperability through data



Automatic refresh of dependent objects



Supporting new languages

Implement a library for Vizier's dataset API

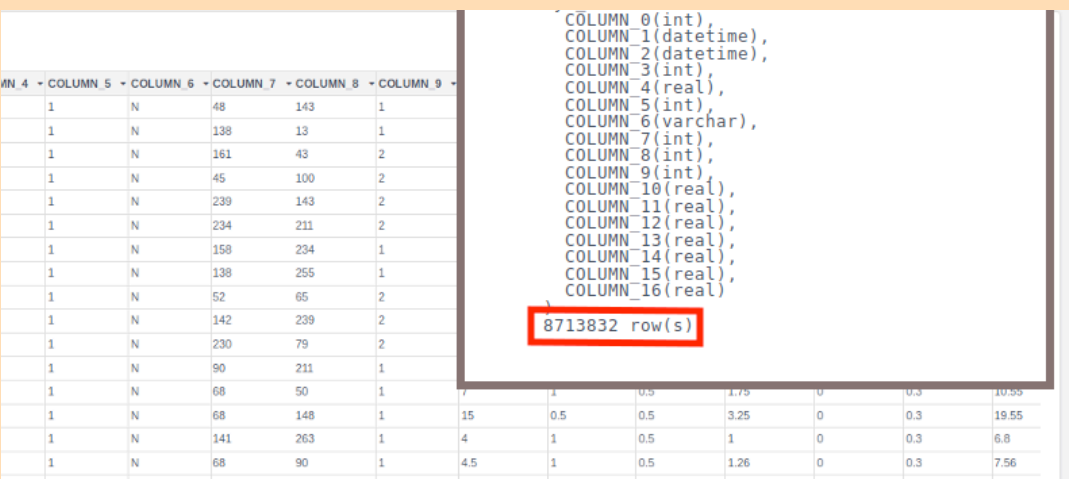
Dependency-tracking for free by monitoring dataset API calls

Languages never communicate directly!
No interoperability issues!

Scalable

Backed by Spark

Vizier dataset = Spark dataframe



Interested in how Vizier tracks uncertainty?

Come to our research talk [1]:
Thu, 11:30-12:50 SIGMOD Research 13, Administratiezaal

Multi-lingual + Multi-modal

Scala

SQL

Python

Vizual spreadsheet scripting

Spreadsheets

Notebooks

Plots

Project Page



Demo Video



Github Repo



[1] S. Feng, A. Huber, B. Glavic, O. Kennedy. Uncertainty Annotated Databases - A Lightweight Approach for Approximating Certain Answers. SIGMOD 2019.