# Demonstration of PI2: Interactive Visualization Interface Generation for SQL Analysis in Notebook

Jeffrey Tao
jat2164@columbia.edu
Columbia University
New York, NY, USA

Yiru Chen
yiru.chen@columbia.edu
Columbia University
New York, NY, USA

Eugene Wu
ewu@cs.columbia.edu
Columbia University
New York, NY, USA

## ABSTRACT

We demonstrate PI2, the first notebook extension that can automatically generate interactive visualization interfaces during SQL-based analyses.

## 1 INTRODUCTION

SQL is the dominant language for accessing and analyzing data. Along with the recent rise of notebooks, a lot of popular computational notebooks for SQL have emerged such as xeus-sqlite [4] in Jupyter, SQL notebook [5], Hex [3], Count [1], etc. Data scientists have gradually ditched their traditional SQL IDEs where they can only see one output at a time in favor of computational notebooks so that they can enjoy narrative programming benefits [11].

Traditional notebooks that can execute SQL queries [4, 5] merely render query results as tables. This is neither satisfying nor effective, as data scientists rely on interactive visualization interfaces to rapidly perform iterative analyses and to better present analyses in a compelling narrative [16].

In contrast to static tables, interactive visualization interfaces (or *interfaces*) consist of three main components: visualizations (e.g., bar and line charts), widgets (e.g. dropdown, slider), and interactions within a visualization (e.g. brushing to select points, panning, clicking). As such, numerous recent notebooks and extensions, such as Lux [12], Count Notebook [1], and Hex Notebook [3] have been designed to help users visualize data and create simple interactive visualizations during their analysis. Unfortunately, the type and complexity of interfaces that they can express are limited (Table 1). For instance, Lux [12] automatically recommends a static visualization when a notebook cell returns a dataframe, but does not support interactive analysis. Similarly, Count [1] and Hex [3] let users visualize a table and add custom widgets that manipulate simple query parameters, but these require explicit user effort. In

| | Lux | Count | Hex | **PI2** |
|---|---|---|---|---|
| *Visualizations* | ✓ | ✓ | ✓ | ✓ |
| *Widgets* | × | *Parameter* | *Parameter* | **Arbitrary** |
| *Vis. Interactions* | × | × | × | ✓ |
| *Zero Effort* | ✓ | × | × | ✓ |

**Table 1: Comparison among different tools.**



**(a) Lux.**

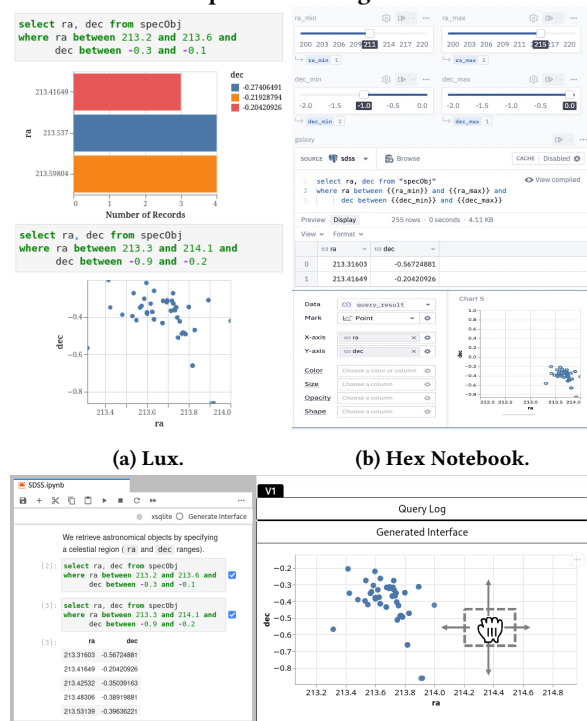**(b) Hex Notebook.**



**(c) This paper: PI2.**

**Figure 1: Different interfaces for analysis of the SDSS dataset: (a) static visualization recommendation with Lux, (b) parameterized query with widgets and visualization with Hex, (c) automatically generated interactive interface with PI2.**

short, existing notebooks have limited support for interactivity, do not support generating interactive visualizations, and require manual effort to create and lay out visualizations and widgets.

This paper demonstrates PI2, the first notebook extension which can automatically generate interactive visualization interfaces during SQL-based analyses. Users can select relevant queries during their analysis and invoke PI2 to synthesize a fully interactive interface with no additional effort. PI2 automatically chooses the appropriate visualizations, widgets, and visualization interactions to fully express the analysis represented by the user's selected queries.

EXAMPLE 1. *Figure 1 illustrates the types of interfaces that Lux, Hex, and PI2 will generate using queries from the Sloan Digital Sky Survey(SDSS) [15] query log. The two example queries retrieve astronomical objects by specifying a celestial region (ra and dec ranges). Lux recommends visualizations for individual tables, so it generates different visualizations for each query, despite their similarities (Figure 1a). Note that the users will need to repeatedly tweak and re-execute the queries if they continue their analysis.*

*Hex lets the user parameterize the ra and dec values in the query, create custom sliders to control each parameter, and visualize the query results in a visualization (Figure 1b). This enables more interactivity than Lux, since the user can directly use the sliders instead of editing SQL strings. However, the user still needs to configure the sliders and choose an effective visualization. Furthermore, using the interface is cumbersome, as the user needs to manipulate four separate sliders to pan and zoom.*

*In contrast, PI2 uses the same two queries to generate the interface in Figure 1c. The visualization supports panning and zooming, so the user can simply drag and scroll on the visualization to manipulate the ra and dec ranges and receive immediate visual feedback. The collapsed Query Log tab archives the input queries.*

In addition to support for visualization interactions, PI2 goes beyond existing notebooks in two ways. First, the widgets and interactions are more than simple query parameters, and can change the structure of the underlying SQL queries as well. For instance, a dropdown may choose between three subqueries, a switch may toggle a filter clause, and a tab may select between different queries to visualize. Second, PI2 takes the available screen size into account in order to select a good layout for the interface—on a large screen, the interface may show multiple visualizations side by side, whereas a small screen may show a single visualization that can be changed via interactions.

We further design the notebook extension to aid iterative analyses. The atomic unit of execution in notebooks allows users to iteratively refer back to previous cells to edit and potentially re-execute them. To adapt to edits and ensure the reproducibility of the generated interface, we take a snapshot of the queries used to generate a new interface. We also version the interfaces, so that users can go back to, or fully revert, to a previous analysis. To avoid interruption of the normal notebook workflow, we choose to lay the *Generated Interface* panel side-by-side with the notebook cells.

The next section presents a brief overview of how PI2 generates interfaces from queries, and Section 3 will illustrate these features in the context of a case study. We refer readers to the technical report [7] for complete technical details of the interface generation process.

## 2 INTERFACE GENERATION OVERVIEW

PI2 transforms an input sequence of queries into an interactive interface in four steps: it parses queries into a generalization of abstract syntax trees (ASTs) that we call DIFFTREEs; maps the DIFFTREEs to a candidate interface; estimates the interface's cost; and repeatedly transforms the DIFFTREEs to generate new candidate interfaces, optimizing according to cost. This section walks through these steps and introduces key concepts.
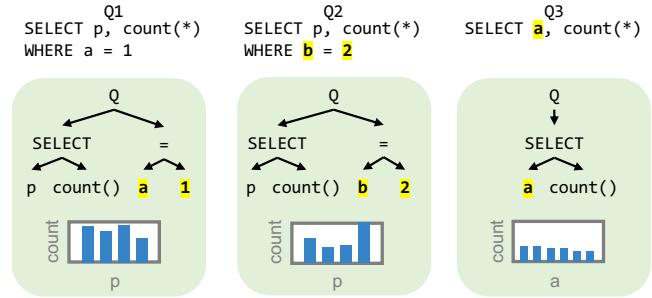


**Figure 2: Example of three queries and their simplified ASTs. A static interface would render one chart for each query.**
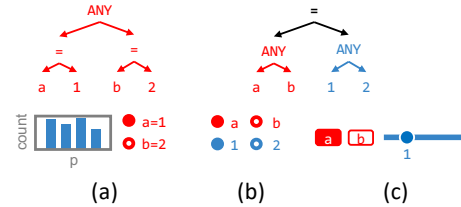


**Figure 3: Example DIFFTREEs and interfaces for Q1, Q2, focusing on the subtree for the predicate. The ANY choice node can choose one of its children. (a) the ANY node maps to a radio button that chooses between the two predicates, (b) the radio lists separately specify the left and right operands, (c) the choice nodes can instead be mapped to a button group and slider, and organized vertically.**

**Static Interfaces:** Figure 2 lists three queries[1] and their ASTs, where attributes p, a, b are integers. Q1 and Q2 change the predicate attribute and literal, and Q3 selects a instead of p. Since each AST is a DIFFTREE, a valid interface simply maps each AST's results to a static chart.

**Interactive Interfaces:** Let us focus on the differing predicate in Q1 and Q2 to show how different DIFFTREE structures can result in different interface designs. For instance, Figure 3(a) is rooted at an ANY node whose children are the two predicates. ANY is a *Choice Node* that can choose one of its child subtrees. In general, choice nodes encode subtree variations[2] that the user can control through the interface. In the example, the ANY node is mapped to two radio buttons (other widgets such as a dropdown are valid as well), where clicking on the first button would bind the ANY to its first child a=1. The DIFFTREE output is visualized as a bar chart.

*Tree Transformations:* Both of ANY's children are rooted at =, so the = can be refactored above the ANY node. This is an example of a *Tree Transformation Rule*. The resulting DIFFTREE in Figure 3(b) shows two ANY nodes that can independently choose the left and right operands. This leads to an interface with two interactions—two sets of radio buttons—and also generalizes the interface beyond the input queries. For instance, the query can now express SELECT p, count(*) WHERE b=1.

A DIFFTREE can map to many interface designs, each with different visualizations, interactions (including widgets and visualization interactions), and layouts. For instance, Figure 3(b) and (c) both express Q1 and Q2, however Figure 3(c) uses horizontal layout and

---

[1]For brevity, we omit the FROM and GROUPBY clauses and show simplified ASTs.
[2]Choice nodes generalize SQL parameterized literals to syntactic structures.
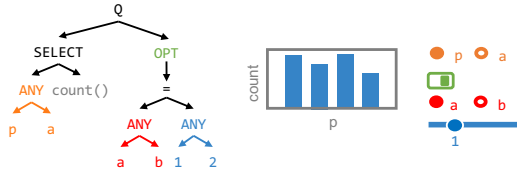
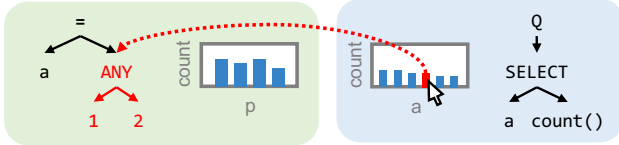**Figure 4: A DIFFTREE for Q1–3 and a candidate interface**



**Figure 5: Multi-view interface where clicking on the right-side chart updates the left chart.**

the slider can choose a continuous range of numbers that generalize beyond the radio buttons in Figure 3(b).

*All Three Queries:* Now, let us add Q3. A simple approach would be to partition the queries into two clusters, where Q3 is rendered as a static chart, and Q1 and Q2 are mapped to one of the interfaces discussed so far. We can then choose to lay these two visualizations out horizontally or vertically. Another possibility is to merge all three queries into a single DIFFTREE, which would map to an interface with a single visualization. Figure 4 illustrates one such DIFFTREE structure, where an ANY node in the SELECT clause chooses whether to project p or a. This maps to an interface similar to Figure 3(c), but with a radio button to choose the attribute to project and a toggle for the optional WHERE clause. Naturally, which of these possible interface designs (or others not discussed here) that should be generated and returned to the user depends on many factors, such as usability, layout, accessibility, and other factors that are difficult to quantify. Quantitative interface evaluation is an active area of research, and PI2 borrows current best practices to develop its cost function.

**Multi-view Interfaces** PI2 can also generate multi-view interfaces. Figure 5 illustrates a slightly different set of queries, where the Q1 and Q2 only differ in the literal, and Q3 remains the same. Since the literal is compared to attribute a, an alternative to mapping the ANY node to a slider is to map it to a *visualization interaction* in Q3's bar chart. Specifically, each bar is derived from (a, count(*)) in Q3's result. Thus, clicking on a bar can also derive a valid value in attribute a's domain that can bind to the ANY node.

**Summary and Generation Pipeline:**

PI2 transforms an input sequence of queries into an interactive interface in four steps (Figure 6). ① It first parses the input query sequence Q into DIFFTREEs. ② PI2 maps the DIFFTREEs to an interface. An interface mapping $\mathbb{I} = (\mathbb{V}, \mathbb{M}, \mathbb{L})$ is defined by a *Visualization Mapping* $\mathbb{V}$ from DIFFTREEs results to visualizations, a *Interaction Mapping* $\mathbb{M}$ from choice nodes to interactions (including widgets and visualization interactions), and a *Layout Mapping* $\mathbb{L}$ from DIFFTREEs structures to layouts. We formulate the interface mapping problem as a schema matching problem by defining schema for both DIFFTREEs and interfaces. ③ A cost model $C(\mathbb{I}, Q)$ evaluates the interface and PI2 either returns the interface
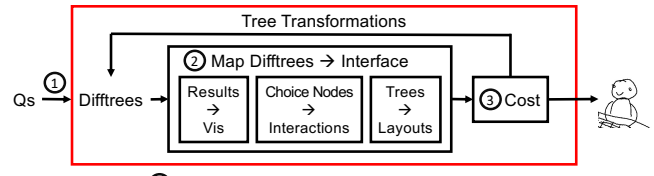


**Figure 6: PI2 interface generation pipeline.**

or chooses a valid transformation to apply to the DIFFTREEs. Informally, our problem is to return the lowest cost interface $\mathbb{I}$ that can express all queries in Q. ④ The space of possible interfaces is enormous, so we solve this problem using Monte Carlo Tree Search [? ] (MCTS). MCTS balances exploitation of good explored states (DIFFTREE structures) with exploration of new states.

## 3 DEMONSTRATION

### 3.1 Interface Design

We integrate PI2 as a Jupyter Lab [10] extension. We design the interface as shown in Figure 7. While authoring SQL queries in the Jupyter notebook, a user can check the checkbox next to each cell to include it as part of the query log for interface generation. Clicking the Generate Interface button will invoke PI2 to generate a new interface to the *Generated Interfaces* panel on the right.

The atomic unit of execution in notebooks allows users to easily refer back to previous cells to edit and potentially re-execute them. To adapt to edits and ensure reproducibility, our integration tracks interface versions in the version tabs at the top of the *Generated Interfaces* panel and archives the input query logs in the *Query Log* collapsible section for each version. PI2 lays out the interfaces and notebook cells side-by-side, so that the normal notebook analysis workflow will not be interrupted.

### 3.2 Use Case Walkthrough

We demonstrate how PI2 aids the data analysis process via a real-world scenario (shown in Figure 7): in late December 2021, an analyst named Jane at a news organization is analyzing a COVID-19 dataset of daily case counts per-state with the intent to give travel warning advice for the winter holiday season. For the sake of comparison, we show a static visualization recommendation below each cell, which is given by an existing system Lux [12].

**Step 1: Overview and detailed look of the dataset.** Seeking to get an overall view of the data, Jane writes Q1 and gets a visualization recommended by Lux showing total case count over time. Looking for a more detailed view, Jane restricts the date range in Q2 to look back over two preceding half-month periods to see more recent trends. Moreover, she would like to do this over different date ranges which will result in many similar static visualizations and a lengthy notebook. With PI2, she can select these three queries via their corresponding checkboxes and automatically generate an interface. PI2 produces a unified interactive interface V1 consisting of two plots: one showing the overall timeline (G1) and the other showing just the selected date range (G2). The two plots are linked by a brushing interaction so that brushing over G1 dynamically configures the date range of G2. Whereas, none of existing notebook tools can create such visualization interactions. In this way,
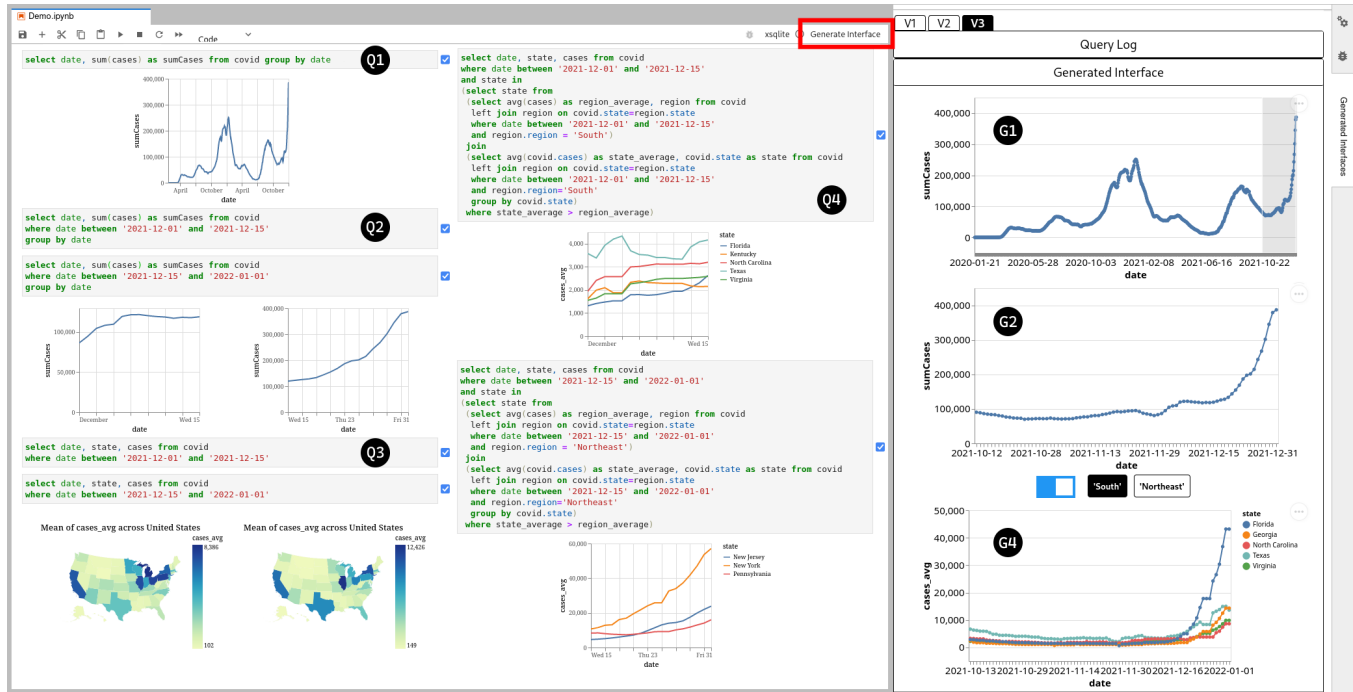
**Figure 7: PI2 Jupyter extension interface. Left: SQL-based analyses. Right: An interface automatically generated by PI2.**

Jane can use this interface to further investigate trends in different date ranges interactively and without writing more SQL queries.

**Step 2: Drill down into state Level.** Jane's task is to give travel warning advice. She writes Q3 to study each state's trends over time. Lux generates choropleth maps visualizing the case count for each state averaged over the date window, which do not allow her to see trends over time. Alternatively, she can use PI2 to generate a new interface V2 with three plots: two are the sames as interface V1 and the third chart (G3) is a bar chart (x → date, y → cases, color → state). We omit G3 in Figure 7 due to space constraints. In this interface V2, the linked brushing interaction will configure the queries underlying both G2 and G3 such that Jane can brush over G1 to see the detailed trend and per-state breakdown trend within the selected date range at the same time.

**Step 3: Focused region investigation.** However, the state breakdown view is visually noisy due to the large number of states. Jane decides to group states into regions and only show those states whose average case counts over a period of time exceed the region's average, expressed as the complicated query Q4 consisting of joins, and correlative subqueries. Further, Jane would like to study the South and Northeast regions within different date ranges. She selects all the queries and invokes PI2. The interface V3 that PI2 generates in Figure 7 has three plots, allowing Jane to view the overall timeline (G1), detail view of the selected date range (G2), and state breakdown filtered for above average states (G4). This new interface maintains the date brushing functionality of previous versions and introduces query configuration widgets: a toggle which allows her to toggle between G3 and G4, and a pair of buttons that switch between the South and Northeast regions. Structurally, the toggle corresponds to an OPT choice node that distinguishes the existence of a complicated subquery—{and state

in...} in Q4 not present in Q3. Through this interface, Jane is able to fluidly reconfigure the date range by brushing on G1, toggle off to see the overall state breakdown of cases, and toggle on to choose to observe trends in the Northeast or South. Noticing very high rates of growth in case count, Jane makes a recommendation that travelers avoid Florida in the South and New York in the Northeast.

Through the above scenario, we show PI2's ability to consume arbitrarily complex SQL queries and automatically generate complete interfaces, including visualization interactions and widgets expressing arbitrary query differences that no other tools can.

**Demonstration engagement.** Participants will be able to write their analysis and generate interactive visual interfaces using PI2 Jupyter extension. We will prepare three datasets – COVID-19, SDSS, and S&P500 for users to explore.

## REFERENCES

[1] 2021. Count – The SQL notebook. https://count.co/.
[2] 2021. Google Covid-19 News. https://news.google.com/covid19/map?hl=en-US&gl=US&ceid=US%3Aen.
[3] 2021. Hex Technologies. https://hex.tech/.
[4] 2021. A Jupyter kernel for SQLite. https://github.com/jupyter-xeus/xeus-sql.
[5] 2021. SQL Notebook. https://sqlnotebook.com/.
[6] Jacques Bertin. 1983. *Semiology of graphics; diagrams networks maps.* Technical Report.
[7] Yiru Chen and Eugene Wu. 2021. PI2: Generating Visual Analysis Interfaces From Queries. arXiv:2107.08203 [cs.DB]
[8] Rémi Coulom. 2006. Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search. In *Computers and Games.*
[9] Krzysztof Gajos and Daniel S Weld. 2004. SUPPLE: automatically generating user interfaces. In *Proceedings of the 9th international conference on Intelligent user interfaces.* 93–100.
[10] B Granger and J Grout. 2016. JupyterLab: Building blocks for interactive computing. *Slides of presentation made at SciPy* (2016).
[11] Donald Ervin Knuth. 1984. Literate programming. *The computer journal* 27, 2 (1984), 97–111.

[12] Doris Jung-Lin Lee, Dixin Tang, Kunal Agarwal, et al. 20. Lux: Always-on Visualization Recommendations for Exploratory Data Science. *VLDB* (20).

[13] Dominik Moritz, Bill Howe, and Jeffrey Heer. 2019. Falcon: Balancing interactive latency and resolution sensitivity for scalable linked visualizations. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–11.

[14] Daniel G Murray. 2013. *Tableau your data!: fast and easy visual analysis with tableau software.* John Wiley & Sons.

[15] SDSS. 2021. Sloan Digital Sky Survey, 2017. http://www.sdss.org/.

[16] Edward Segel and Jeffrey Heer. 2010. Narrative visualization: Telling stories with data. *IEEE transactions on visualization and computer graphics* 16, 6 (2010), 1139–1148.