

Managing Messes in Computational Notebooks

Anonymized for submission

ABSTRACT

Data analysts use computational notebooks to write code for analyzing and visualizing data. Notebooks help analysts iteratively write analysis code by letting them interleave code with output, and selectively execute cells. However, as analysis progresses, analysts leave behind old code and outputs, and overwrite important code, producing cluttered and inconsistent notebooks. This paper introduces *code gathering tools*, extensions to computational notebooks that help analysts find, clean, recover, and compare versions of code in cluttered, inconsistent notebooks. The tools archive all versions of code outputs, allowing analysts to review these versions and recover the subsets of code that produced them. These subsets can serve as succinct summaries of analysis activity or starting points for new analyses. In a qualitative usability study, 12 professional analysts found the tools useful for cleaning notebooks and writing analysis code, and discovered new ways to use them, like generating personal documentation and lightweight versioning.

CCS CONCEPTS

- Human-centered computing → Interactive systems and tools;

KEYWORDS

Computational notebooks, messes, clutter, inconsistency, exploratory programming, code history, program slicing

ACM Reference Format:

Anonymized for submission. 2019. Managing Messes in Computational Notebooks. In *Proceedings of ACM SIGCHI Conference on Human Factors in Computing Systems (CHI '19)*. ACM, New York, NY, USA, 12 pages. https://doi.org/10.475/123_4

1 INTRODUCTION

Data analysts often engage in “exploratory programming” as they write and refine code to understand unfamiliar data, test hypotheses, and build models [19]. For this activity, they frequently use computational notebooks, which supplement the rapid iteration of an interpreted programming language with the ability to edit code in place, and see computational results interleaved with the code. A notebook’s flexibility is also a downside, leading to messy code: in recent studies, analysts have called their code “ad hoc,” “experimental,” and “throw-away” [16], and described their notebooks as

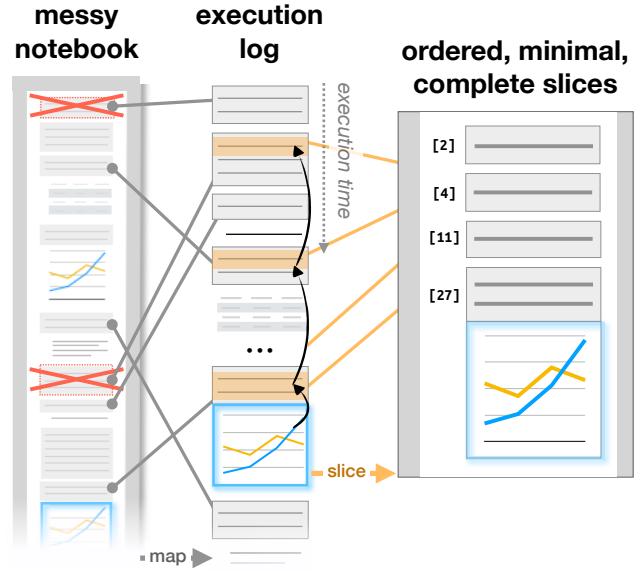


Figure 1: Code gathering tools help analysts manage programming messes in computational notebooks. The tools map selected results (e.g., outputs, charts, tables) in a notebook to the ordered, minimal subsets or “slices” of code that produced them. With these slices, the tools help analysts clean their notebooks, browse versions of results, and discover provenance of results.

“messy” [20, 29], containing “ugly code” and “dirty tricks” in need of “cleaning” and “polishing” [29].

In essence, a notebook’s user interface is a collection of code editors, called “cells.” At any time, the user can submit code from any cell to a hidden interpreter session. This design leads to three types of messes common to notebooks: disorder, where the interpreter runs code in a different order than it is presented in the cells; deletion, where the user deletes or overwrites the contents of a cell, but the interpreter retains the effect of the cell’s code; and dispersal, where the code that generates a result is spread across many distant cells. Such messes make it difficult for an analyst to understand the code, and to recall how results (i.e. charts, tables, and other code outputs) were produced, impeding their ability to reuse, extend, and share their analysis code [20, 29].

In this paper, we aim to improve the state of the art in tools for managing messes in notebooks. We introduce a suite of interactive tools, code gathering tools, as an extension to computational notebooks. The tools afford analysts the ability to find, clean, and compare versions of code in messy notebooks. They build on a static program analysis

107 technique called program slicing [33], which answers queries
 108 about the dependencies among a program's variables. With
 109 code gathering tools, an analyst first selects a set of analysis
 110 results, which can be any cell output (e.g., charts, tables, con-
 111 sole output) or variable definition (e.g., data tables, models).
 112 Then the tool searches the execution log—an ordered history
 113 of all cells executed—to find an ordered, minimal subset or
 114 "slice" of code needed to compute the selected results.

115 This paper makes two contributions. The first contribution
 116 is the design and implementation of *code gathering tools*.
 117 These tools offer interactive affordances for managing messes
 118 in notebooks. Specifically, the tools highlight dependencies
 119 used to compute results, to help analysts find code they wish
 120 to understand, reuse, and rewrite in cluttered notebooks.
 121 They provide ordered, minimal code slices that can serve as
 122 succinct summaries of analysis activity or starting points for
 123 branching analyses. Additionally, they archive past versions
 124 of results and allow analysts to explore these versions, and
 125 the code slices that produced them. Code gathering tools
 126 are implemented as an extension to Jupyter Notebook [2], a
 127 popular notebook with millions of users [17]. The extension
 128 is available for use as a design artifact and as practical tool
 129 for exploratory data analysis in notebooks.¹

130 The most important idea behind the interaction design
 131 of code gathering tools is *post-hoc mess management*—that
 132 tools should allow analysts to easily find, clean, and compare
 133 versions of code in notebooks, regardless of whether they
 134 have followed a disciplined strategy to organize and version
 135 their code. Past tools for cleaning code often require effort:
 136 annotating cells with dependency information [3], folding
 137 and unfolding cells [27], marking and tagging lightweight
 138 versions of snippets [18], etc. With code gathering tools, his-
 139 tory is stored silently, and tailored slices of code are recalled
 140 on-demand with two or fewer clicks.

141 Our second contribution is a qualitative usability study
 142 providing insight into the uses and usability of code gather-
 143 ing tools for managing exploratory messes. 12 professional
 144 data analysts used the tools in an in-lab study to clean note-
 145 books and perform exploratory data analysis. We found that
 146 affordances for gathering code to a notebook were both val-
 147 ued and versatile, enabling analysts to clean notebooks for
 148 multiple audiences, generate personal reference material,
 149 and perform just-in-time branching. We also refined our
 150 understanding of the meaning of "cleaning," and how code
 151 gathering tools supports an important yet still incomplete
 152 set of tasks analysts consider as part of code cleaning. This
 153 study confirmed that analysts thirst for tools that help them
 154 manage exploratory messes, and that code gathering tools
 155 provide a useful tool for managing these messes.

156
 157 ¹To protect the anonymity of this paper, the extension will be made available
 158 at the conclusion of blind review.

2 BACKGROUND AND RELATED WORK

Messes in computational notebooks

Lackluster code quality appears to be intrinsic to exploratory programming. Analysts regularly prioritize the speedy discovery of solutions over writing high-quality code [19]. They clutter their programs by saving old versions of their code in comments [18, 35]. In notebooks in particular, poor code quality takes on a spatial dimension. Messes accrue and disappear in an iterative process of expansion and reduction of code: analysts write code in many small cells to try out different approaches to solve a problem, view intermediary output from their code, and debug their code; and then combine and eliminate cells as their analyses reach completion [20].

Eventually, messes get in the way of data analysis. It becomes difficult to understand analyses separated across many cells of a notebook, and long notebooks become time-consuming to navigate [20]. Important results accidentally get overwritten or deleted [29]. Data analysis is often fundamentally collaborative as analysts share notebooks, code snippets, documents, slides, images, and web pages [16, 20, 29]. Though when notebooks become too messy, analysts don't feel comfortable sharing their insights as is with others without substantial cleanup effort [27, 29].

To manage these messes, analysts have diverse strategies to clean their code. Many delete cells they no longer need, consolidate smaller cells into larger cells, and delete full analyses that did not turn out helpful. Long notebooks are abandoned for "fresh" ones with only a subset of successful parts from the long ones. Analysts organize code as they build it, some coding from top to bottom, some adding cells where they extend old analyses, some placing functions at the top, and some placing them at the bottom [20]. They add tables of contents, number sections, limit the size of cells, and split long notebooks into shorter ones [29].

Because of the challenges and tedium of managing messes, data analysts have given clear indications they need better tools to support the management of messes. In prior studies, analysts have asked for tools that let them collect scripts that can reproduce specific results; to compare outcomes from different versions of an analysis; and to recover a copy of a notebook that produces a version of a result [20]. They have also asked for tools that recover the history of how data was created, used, and modified [27]. This paper answers the call of these analysts, offering code gathering tools to help analysts manage messes in their notebooks.

Tools for cleaning messy code

Messiness is pervasive problem for code written in any language or domain. As such, researchers have designed tools to help coders—here, meaning software developers, end-user

213 programmers, analysts, and everyone else that writes code—
 214 rewrite unclear, bloated code. Recent tools have provided
 215 mixed-initiative interactions where coders work with their
 216 code editor to extract minimal, executable code examples
 217 that can reproduce a behavior [12] and large components of
 218 complex software projects [15]. Other tools let programmers
 219 clean application code by demonstrating by interacting with
 220 that application, and preserving only the executed code [23].
 221 Such tools, like this work, leverage program slicing [34] to
 222 trace from code a coder indicated they want to preserve to
 223 other lines that should be included in a clean program. Code
 224 gathering tools uniquely aid with cleanup of notebook code
 225 by slicing a history of executed cell code, using the outputs
 226 of analysis code as a slicing criterion.

227 Another step of cleaning code is annotating it so others
 228 can understand, run, extend, and maintain it. Beyond
 229 text editing tools for writing comments and documentation,
 230 tools can help coders annotate code by helping them
 231 create walkthroughs visiting locations of interest in code
 232 projects [24, 32] and, in the context of notebooks, by helping
 233 coders fold and annotate cells with descriptive headers [27].
 234 Tools such as these complement the code-centric cleaning
 235 utilities that code gathering tools currently provide.

239 Locating the causes of program output

240 Finding the code that is responsible for producing a result
 241 is complex when analysis produced many variants of code,
 242 generated long code files, or took place long ago. Amid these
 243 challenges, researchers have studied how tools can support
 244 analysts in locating the causes of program outputs. For exam-
 245 ple, tools can instrument programming languages (e.g., [25]),
 246 data analysis applications (e.g., [7]), and operating systems
 247 (e.g., [10]) to capture what results get produced, and the
 248 parameters, data, and code that produced them. Thereafter, this
 249 history is available and can be queried when analysts seek
 250 to recover the provenance of a result.

251 For programs more generally, prototype tools have been
 252 designed to help coders understand the causes of many
 253 types of program output, including graphics, console output,
 254 and bugs [21], interactive web page behavior [5, 6, 13, 14],
 255 user interfaces [4], and animations [9]. These tools help
 256 coders query for code by directly selecting outputs of inter-
 257 est [6, 9, 21], search for root causes of errors in code [21],
 258 record and replay ephemeral program behavior [5, 6], high-
 259 light executing lines of code [4], and filter code on a sliding
 260 scale of its relevance to a demonstrated behavior [14]. Code
 261 gathering tools similarly helps analysts find the code produc-
 262 ing specific results from analysis, highlighting potentially
 263 disordered code dependencies in notebooks.

264

265

266 Managing and reviewing code versions

267 As analysts write code, they produce many versions of their
 268 code, from the resolution of snippets to files. These versions
 269 may be informally stored as comments [18] or notebook
 270 cells [20], leading to drawbacks like the difficulty of recalling
 271 promising versions of code, and code bloat. To help ana-
 272 lists manage versions, researchers have designed tools to
 273 make informal versioning of code easier, by enabling the ver-
 274 sioning of small blocks of code directly within editors [18],
 275 automatically saving versions of cells upon each computa-
 276 tion [26], and enabling exploratory programming on several
 277 simultaneous, linked versions of a code file [11].

278 These tools help analysts manage code versions amid fast-
 279 paced iterative programming. In dialogue with this work,
 280 code gathering tools support post-hoc mess management, by
 281 helping coders find, clean, and compare versions of code even
 282 if they spent no up-front effort on managing code versions or
 283 organizing their code. This approach is inspired by Janus' au-
 284 tomatic versioning of code cells upon execution [26], and Var-
 285 iolite's archiving of runtime configurations [18]. Uniquely,
 286 code gathering tools enable analysts to recover and compare
 287 complete slices that produce results, even if they span many
 288 cells that appear out of order in a notebook.

289 In its intention to help analysts manipulate their program-
 290 ming history, code gathering tools builds on a tradition of
 291 research that shows that coders often consult programming
 292 histories to understand, debug, and extend their code [8] and
 293 rely on code's output, changelogs, and source code to find
 294 code to reuse [31]. Prior tools have supported in selectively
 295 undoing history [36], viewing the evolution of selected code
 296 snippets [30], and completing code with legacy names [22].
 297 Code gathering tools seek to bring the mastery and under-
 298 standing of history embodied by many of these prior tools
 299 into data analysts' computational notebooks.

3 A DEMO OF CODE GATHERING TOOLS

300 To convey the experience of using the code gathering tools
 301 in Jupyter Notebook, we describe a short scenario.² Consider
 302 an analyst, Dana, who is performing exploratory data analy-
 303 sis to understand variation and determiners of quality of a
 304 popular consumer good—chocolate. This section shows how
 305 code gathering tools could help her find, clean, and compare
 306 versions of code during data analysis.

307 Prologue: A proliferation of cells

308 Dana starts her analysis by loading a dataset, importing
 309 dependencies, and filtering and transforming the data. She
 310 writes code to display tables so she can preview the data,
 311 and displays charts and statistics that summarize the data's
 312 properties. To better understand key features of the data, she

²See also this paper's video figure.

319 builds a model to predict chocolate quality from the other
320 features. Through experimentation, she tailors the model
321 parameters to learn more about the features. Throughout the
322 analysis, she makes messes, overwriting old code, deleting
323 code that appears irrelevant, running cells out-of-order, and
324 accumulating dozens of cells full of code and results. Dana
325 starts to have trouble finding what she needs in the notebook.

Finding the code that produces a result

After several hours building and testing models, Dana is satisfied with a version of the model, but then recalls a peculiarity of the data: one of the columns is a percentage multiplied by 10 to avoid decimal points (e.g. 73.5% is represented with the integer 735). Although Dana wrote some code to check for triple-digit numbers, she cannot remember if she ran this code on the dataset before training her model. Hours have passed since she touched any code at the top of her notebook, where she loaded and cleaned the data.

Because she has installed code gathering tools, Dana sees all variable definitions highlighted in blue and all visual outputs outlined in blue. These marks highlight all results of analysis, including tables, charts, console output, console errors, definitions of models, loaded datasets, and transformed copies of those datasets. The marks indicate that Dana can click on any of the results to see the code in the notebook computed it, directly or indirectly.

Dana wants to know whether she cleaned her percentage data before training her classifier. She clicks on the results of classification—in this case, the predictions in the variable `predictions`. After clicking this variable, all lines that were used to compute the variable’s value—including code for prediction, model training, data preparation, data loading, and importing dependencies—are highlighted in light purple (Figure 2). All other lines remain unhighlighted. With the relevant code highlighted, Dana scrolls through the sprawling notebook to browse the highlighted lines, skipping over long sections of irrelevant code and results. She finds the code that transforms the percentage data, namely, a cell defining the function `normalizeIt` and the cell after it. Because these lines are highlighted, Dana knows that she cleaned the percentage column before classification.

Removing old and distracting analysis code

362 Dana now has a notebook with a model that she likes—and
363 much more code she no longer needs. Now that Dana knows
364 what her data looks like and has a working set of data filter-
365 ing, data transformation, and model training code, the code
366 she wrote to visualize the data and debug her APIs will just
367 get in the way. Dana decides to clean her notebook to a state
368 where it only has the useful model-building code.

To clean the notebook, Dana clicks on a few results she wants to still be computed in the cleaned notebook, namely,

```
In [206]: def normalizeIt(percent):
    if percent > 100:
        percent = int(str(percent)[-1])
    return percent

In [207]: df['CocoaInPercent'] = df['Cocoa\Percent'].apply(normalizeIt)

In [209]: df['Rating'] = (df['Rating']* 100).astype(int)
df['Rating'].head(5)

In [216]: X = df.drop('Rating', axis = 1) #Features
y = df['Rating'] # target Variables
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

In [220]: dtree = DecisionTreeClassifier(max_depth=12)
dtree.fit(X_train, y_train)

In [221]: predictions = dtree.predict(X_test)
```



★

Figure 2: Finding relevant code with code gathering tools. With code gathering tools, an analyst can click on any result, and the notebook highlights in light purple just those lines that were used to compute the result or variable. The highlights appear throughout the notebook (which is condensed in the figure). Black arrows have been added in the figure to indicate the data dependencies that cause each line to be included in the highlighted set.

the classification results in the variable predictions and a bar chart showing the range of chocolate qualities used to build the classifier (Figure 3.2). Dana gets a sense of the size of the final cleaned notebook by looking at which lines in the notebook are highlighted as she selects each result. Then, Dana clicks the “Gather to Notebook” button (Figure 3.3), which opens a new notebook with the definition of predictions, the bar chart of chocolate quality, and the other code needed to produce these two results. The new, cleaned notebook has 16 cells, instead of the 47 in her original notebook. It contains the bar chart and omits 28 other visual results in the original. This reduces the overall size of the notebook from 13,044 to 1,248 vertical pixels in her browser, which is much easier to scroll through when editing the code (Figure 3.4). This cleaned notebook is guaranteed to replicate the results, as the tool reorders cells and resurrects deleted cells as necessary to produce the selected results. Dana verifies that running this notebook start-to-finish indeed replicates the chosen predictions and bar chart.

Reviewing versions of a result and the ordered, minimal code slices that produced them

To build a better predictor, Dana has been experimenting with different parameters to a decision tree classifier, like its maximum allowable depth and the minimum samples per branch. Dana remembers that she had previously created a simple, shallow decision tree with promising performance, but has not yet found a model with better performance.

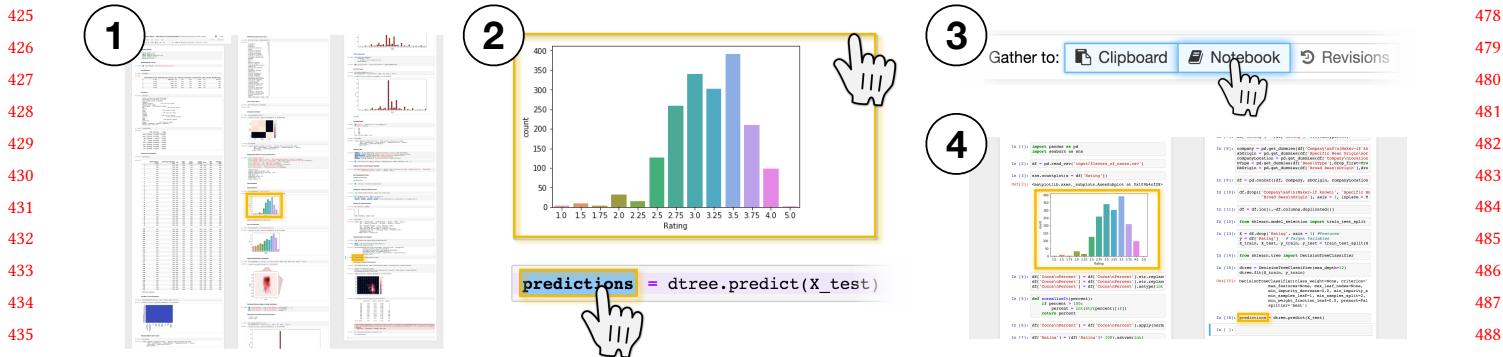


Figure 3: Cleaning a notebook with code gathering tools. Over the course of a long analysis, a notebook will become cluttered and inconsistent (1). With code gathering tools, an analyst can select results (e.g., charts, tables, variable definitions, and any other code output) (2) and click “Gather to Notebook” (3) to obtain a minimal, complete, ordered slice that replicate the selected results (4).

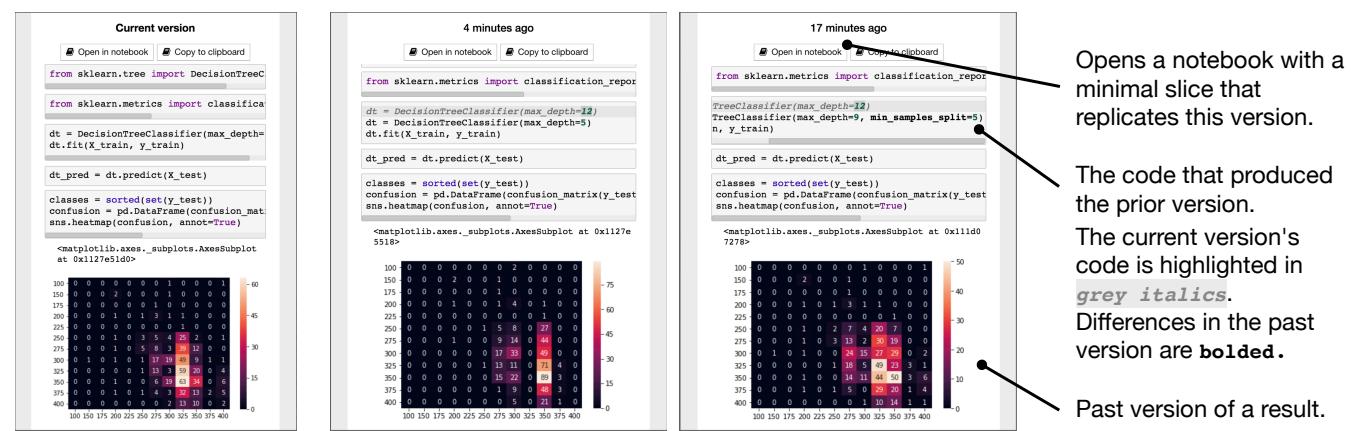


Figure 4: Comparing versions of a result with code gathering tools. When an analysts executes a cell multiple times, code gathering tools archive each version of the cell. When the analyst chooses the cell’s output—say, the confusion matrix shown above—and clicks “Gather to Revisions,” a version browser appears that lets them see all versions of that output, compare the code slices that produced each version, and load any of these slices to a new notebook, where the version’s results can be replicated.

With code gathering tools, Dana can summon all past versions of her classifier’s results and compare the code she used to produce these results. To do this, she clicks on a result—namely, a confusion matrix which visualizes the accuracy of the decision tree for each class—and then on the “Gather to Revisions” button. This brings up a version browser (Figure 4). Here, Dana sees all the versions of the result, arranged from left to right, starting with the current version and ending with the oldest version. Each version includes the relative time the result was computed, the code slice that produced that version of the result, and the result itself, in this case a confusion matrix.

Scrolling horizontally to access older versions, Dana finds several examples of decision trees with comparatively good accuracy. Among these alternatives, she searches for a shallow version of the decision tree with high accuracy. Dana

looks above each result to see the code that produced it. The code is segmented into the cells that were originally executed, omitting lines that were not used to compute the result. Dana is familiar with the most recent version of the code and quickly finds the differences in past code versions by looking for bold code with a grey background. By seeing each version of the confusion matrix together with the model parameters that produced it, Dana finds the model she is looking for—a shallow tree with good performance, specifically with a depth of 3 and a minimum branch size of 5. The code that produced this version can be copied as cells or text to the clipboard, or opened as a new notebook that replicates that version; Dana opens a notebook with this version so she can refer back to it later.

425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530

531 **Cleaning finished analysis code**

532 Dana finished her data analysis and wants to share the results
 533 with an analyst on her team who can check her results and
 534 suggest improvements. However, the notebook is once again
 535 cluttered with code that would distract her colleague. While
 536 Dana wants to save her long and verbose notebook for her
 537 personal use later, she also wants a clean and succinct version
 538 of the notebook for her colleague. She chooses the prediction
 539 results of her model, clicks “Gather to Notebook,” and saves
 540 the generated notebook to a folder shared with her colleague.
 541

542 **Exporting analysis code to a standalone script**

543 After refining her analysis with her colleague, Dana wants
 544 to export a script that can be packaged with an article she
 545 is writing, so that others can replicate her results in their
 546 preferred Python environments. To do this, Dana selects
 547 the code that produces the results she wants her script to
 548 replicate, clicks “Gather to Clipboard,” and then pastes the
 549 gathered code into a blank text file. This script replicates the
 550 results Dana produced in her notebook.
 551

552 **4 DESIGN MOTIVATIONS**

553 The choice to build these specific code gathering tools was
 554 informed by formative interviews with eight data analysts
 555 and builders of tools for data analysis at a large, data-driven
 556 software company. During these interviews, we proposed
 557 several extensions to the notebook interaction model, and
 558 data scientists expressed the most enthusiasm for tools to
 559 help them clean their results, and explore past variants of
 560 their code. In this section, we highlight a few of the major
 561 design ideas motivating the current design of code gathering
 562 tools, which may be helpful for others designing tools for
 563 managing messes in computational notebooks:
 564

565 *Post-hoc management of messes.* Given analysts’ diverse
 566 personal preferences of whether and how to organize and
 567 manage versions of code, we decided code gathering tools
 568 should assist analysts regardless of such preferences. Whether
 569 they prefer to overwrite or save code versions explicitly, the
 570 tools still silently collect history, and provide access to the
 571 code that produced any visible result.

572 *Portability of gathered code.* Analysts reuse a notebook’s
 573 code in that notebook, other notebooks, and scripts [20]. To
 574 support diverse code destinations, it is equally easy to gather
 575 code into cells, new notebooks, and lines of text. Analysts
 576 can “gather to notebook,” which opens gathered cells in a
 577 new notebook, or they can “gather to clipboard,” which si-
 578 multaneously copies the cells to Jupyter Notebook’s internal
 579 clipboard, and the code text to the system clipboard.

581 *Query code via direct selection of analysis results.* Prior re-
 582 search shows that programmers frequently look to program
 583

584 output when searching for code to reuse [31]. In notebooks,
 585 visual results break up walls of monospace text, providing
 586 beacons. We anticipated that selections of results would pro-
 587 vide the most direct method for accessing relevant history.
 588

589 **5 IMPLEMENTATION**

590 A computational notebook uses an underlying language in-
 591 terpreter (like Python). At any time, an analyst can submit
 592 any cell’s code to the interpreter, in any order. The results
 593 that the interpreter produces and that the notebook displays
 594 depend on the order in which the analyst submits cell code.
 595 Hence, in the notebook context, a notebook’s “program” is
 596 not the content of the notebook’s cells, but the content of the
 597 cells that the analyst runs, in the order in which the analyst
 598 runs them. We call this the *execution log*.

599 We define code gathering as the application of program
 600 slicing to an execution log to collect ordered, minimal subsets
 601 of code that produced a given result. Program slicing is a
 602 static analysis technique wherein, given a target statement
 603 (called the slicing criterion), program slicing computes the
 604 subset of program statements (called the slice) that affect
 605 the value of the variables at the target statement [34]. In
 606 the notebook context, the variables/outputs that an analyst
 607 selects are the slicing criteria, and the gathered code is the
 608 slice. We implemented code gathering as a Jupyter Notebook
 609 extension with roughly 5,000 lines of TypeScript code. Our
 610 implementation supports notebooks written in Python 3. The
 611 details in this section could serve as a conceptual template for
 612 tool builders seeking to support code gathering for notebooks
 613 in other Pythonic languages like Julia and R.
 614

615 **Collecting and slicing an execution log**

616 To find the code that produces a result, the tools first need
 617 a complete and ordered record of the code executed in the
 618 notebook. We build such a record, the “execution log,” by
 619 saving a summary of each cell as it is executed. A cell sum-
 620 mary contains two parts: first, the cell’s code, which will be
 621 joined with the code of other cells into a temporary program
 622 used to find code dependencies; second, the cell’s results,
 623 which can be used as slicing criteria, and shown in a version
 624 browser as the output of running that cell.

625 The code for some cells, if included in the execution log,
 626 will cause errors during program slicing. Namely, if the code
 627 that contains syntax errors, the temporary program used
 628 during dependency analysis will fail to parse; if it raises run-
 629 time errors, a slice containing that cell might raise the same
 630 error. Therefore, cells with syntax errors and runtime errors
 631 are omitted from the log. Ignoring cells with parse errors
 632 is consistent with Jupyter’s semantics: if an executed cell
 633 contains any parse errors, all of its code is ignored by the
 634 interpreter. Ignoring cells with run-time errors is inconsis-
 635 tent with Jupyter’s semantics, in that the interpreter will run

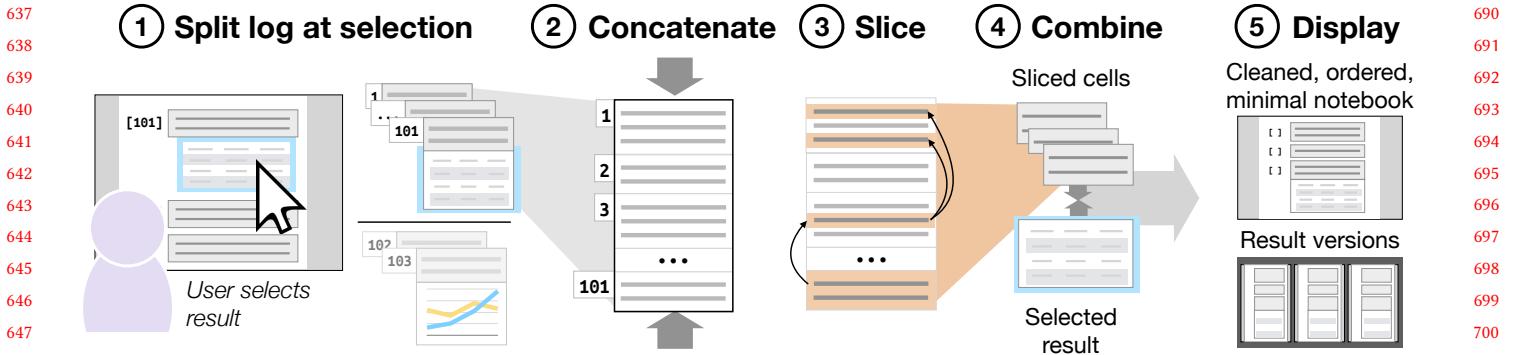


Figure 5: Implementation of code gathering. When an analyst wants to gather the code that produced a result, the code gathering backend splits the log of executed cells at the last cell where the analyst clicked a result, and discards the other cells (1), concatenates the text from the remaining cells into a program (2), slices the program using the analyst’s selections as a slicing criterion (3), combines the sliced cells with the selected results if they are code outputs (4), and displays these cells in a notebook or a version browser (5).

the statements up to the point where the error occurs. This limitation does not cause problems in practice, since analysts typically correct such errors and re-run the cells.

Next, we slice the execution log to provide code slices that replicate results. When an analyst selects results in a notebook, they specify slicing criteria. When they select a variable definition, they add the statement containing the variable definition as a slicing criterion. When they select a cell’s output, they add all statements in that cell.

To slice the execution log, there must first be a “program” to slice. We build such a program by filtering the log to the cells were executed before the cells containing slicing criteria: these cells won’t be included in the slice, and would unnecessarily slow down the slicing algorithm. Then, the program is built by joining the text of the remaining cells, in the order they were executed (Figure 5.1–2).

Finally, we slice the program (Figure 5.3). We implemented a standard program slicing workflow—parsing the program with a JSON-generated parser [1]; searching the parse tree for variable uses, definitions, and control blocks; computing control dependencies (e.g., dependencies from statements to surrounding if-conditions and for-loops) and data dependencies (e.g., dependencies from statements using a variable to statements that define or modify that variable); and slicing by tracing back from the slicing criteria to all the statements they depend on. When computing data dependencies, we determine if methods modify their arguments by looking up this information in a custom, extensible configuration file containing data dependencies for functions from common data analysis libraries (e.g., pandas, matplotlib).

Our current implementation supports interactive computation times by splitting slicing into small, reusable parts: when a cell is executed, its code is immediately parsed, and

its variable definitions and uses detected. With these precomputed pieces of analysis, gathering takes place at interactive speeds, as the most costly analyses have been performed before the analyst gathers any code.

6 IN-LAB USABILITY STUDY

We designed a two-hour, in-lab usability study to understand the support that code gathering tools can provide data analysts as they write code in computational notebooks. We were fairly confident of the ability of code gathering tools to eliminate clerical work—like the removal of irrelevant code, or recovery of dead code—given the design of the tool and evidence from several prior pilot studies. Therefore, the questions we sought to answer focused on the match between the control analysts desired over messy notebooks, and the support code gathering tools currently provide. We therefore designed our study to answer these research questions:

RQ1. What does it mean to “clean”? When we ask analysts to clean a notebook, what do they do? Could code gathering tools support the work they are doing?

RQ2. How do analysts use code gathering tools during exploratory data analysis? In our design of the tools, we hypothesized that analysts would use the tools for highlighting code, gathering to notebooks, and version browsing to find, clean, and compare versions of code. Do they?

We invited 200 randomly selected data analysts at a large, data-driven software company. The invitation stated the requirement of experience with Jupyter notebooks and Python. We recruited 12 participants altogether (3 female, aged 25–40 years, median age = 29.5 years). Participants reported the following median years of experience on a Likert scale: 6–10 years programming; 3–5 years programming in Python; and

743 1–2 years using Jupyter Notebooks. Five participants report
 744 using Jupyter Notebooks daily; three, weekly; one, monthly;
 745 and three, less than monthly. We compensated participants
 746 with a US\$50 Amazon gift card.

747 **Tasks.** To start, each participant signed a consent form
 748 and filled out a background questionnaire. The session then
 749 consisted of two cleaning tasks and a exploratory data anal-
 750 ysis task. For the two cleaning tasks, we gave participants
 751 two existing notebooks from the UCSD Jupyter Notebook
 752 archive [28], one about Titanic passengers, one about the
 753 World Happiness Index. We chose these notebooks because
 754 they are in Python, execute without errors, use popular anal-
 755 ysis and visualization libraries, involve non-technical do-
 756 mains, and are long enough to be worthy of cleaning. We
 757 counterbalanced use of the two notebooks between subjects.
 758

759 For the first cleaning task, we asked the participant to scan
 760 the notebook for an interesting result and to clean the note-
 761 book with the goal of sharing that result with a colleague (10
 762 minutes). After a brief tutorial about code gathering, we then
 763 asked the participant to repeat the cleaning task on a differ-
 764 ent notebook, this time using the code gathering features (10
 765 minutes). Finally, for the exploratory task, we gave partic-
 766 ipants a dataset about Hollywood movies and asked them
 767 to create their own movie rankings, ready for sharing (up to
 768 30 minutes). We chose this dataset as we thought it would
 769 understandable and interesting to analysts from a wide va-
 770 riety of backgrounds. During all tasks, participants could
 771 use a web browser to search the web for programming refer-
 772 ence material. After each of the three tasks, the participant
 773 filled out a questionnaire: the first about how they currently
 774 clean notebooks; the second about the usefulness of code
 775 gathering tools for notebook cleaning; and the third about
 776 the usefulness of code gathering tools for data exploration.
 777 Throughout the tasks, we encouraged participants to think
 778 aloud, and we transcribed their remarks.

779 Each participant used an eight-core, 64-bit PC with 32
 780 GB of RAM, running Windows 10, with two side-by-side
 781 monitors with 1920 × 1200 pixels. One monitor displayed
 782 Jupyter Notebooks; the other displayed our tutorial and a
 783 browser opened to a search engine.

7 RESULTS

786 In the section below, we refer to the 12 analysts from the
 787 study with the pseudonyms P1–12.

The meaning of “cleaning”

790 Before giving analysts the tutorial about code gathering tools,
 791 we first asked them to describe their cleaning practice and to
 792 clean a notebook in their usual way. This allowed us to under-
 793 stand their own interpretation of “cleaning” before biasing
 794 them with our tool’s capabilities. Many analysts explained

795 “cleaning” in a way that is compatible with code gathering,
 796 namely keeping a desired subset of results while discarding
 797 the rest (P8, P10–12). Indeed, one analysts’s description of
 798 cleaning is surprisingly close to the code gathering algorithm:
 799 “So I picked a plot that looked interesting and that’s maybe
 800 something I would want to share with someone and then, if
 801 you think of a dependency tree of cells, sort of walked back-
 802 wards, removed everything that wasn’t necessary” (P10).

803 In their everyday work, some analysts clean by deleting
 804 unwanted cells, but most copy/paste desired cells to a fresh
 805 notebook. (One analysts who cleans by deletion initially
 806 found the non-destructive nature of code gathering to be
 807 unintuitive, but adjusted after practice (P4).) Many described
 808 the process as error-prone and frequently re-execute the
 809 cleaned notebook to check that nothing is broken.

810 Every analyst reported that choosing a subset of cells is
 811 part of the cleaning process. However, for several analysts,
 812 “cleaning” includes additional activities. Several analysts re-
 813 ported that cleaning involves a shift in audience from oneself
 814 to other stakeholders, like peers and managers (P1, P5–7,
 815 P11). Hence, cleaning involves adding documentation (com-
 816 ments or markdown) (P1, P5, P7, P10, P11) and polishing
 817 visualizations (e.g., adding titles and legends) (P1, P6). Some
 818 analysts reported that cleanup includes improving both note-
 819 book quality (e.g., merging related cells (11) and eliminating
 820 unwanted outputs (P3, P6)) and code quality (e.g., eliminating
 821 (P3, P6) or refactoring (P3, P4, P12) repeated code). Finally,
 822 for some, cleaning involves integrating the code into a team
 823 engineering process, for example, by checking the code into
 824 a repository or turning it into a reusable script (P7).

825 **Code comments.** When code is gathered, code comments
 826 and markdown are removed, as they are unnecessary to repli-
 827 cating a result. One analyst believed this made the code look
 828 cleaner (P5). Yet at least two others believed that comments,
 829 when close to the code, should be included by default (P3,
 830 P10). One of these analysts noted that including irrelevant
 831 comments would not be problematic, as “it’s easy to remove
 832 some extraneous text” (P10).

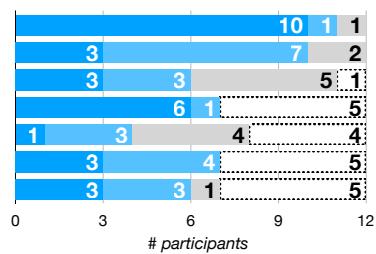
How analysts use code gathering tools to support 833 exploratory data analysis

834 After both the second notebook cleaning task and the ex-
 835 ploratory analysis task, we asked analysts to provide sub-
 836 jective assessments of code gathering, broken down into
 837 seven features (Figure 6). Gathering code to a new notebook
 838 is the clear favorite, with nearly every analyst rating it as
 839 “very useful” for both tasks. The dependency highlights were
 840 also popular. Many analysts did not find opportunities to try
 841 the version browser during the two tasks, likely due to the
 842 short duration of the lab session. Similarly, many analysts did
 843 not experience the recovery of deleted code, either because

849

	Task 1: Cleaning Notebooks				
850	Gather code to new notebook	12	1	1	
851	Highlight lines code depends on	7	3	1	
852	Paste gathered code as cells	3	6	3	
853	Reorder gathered code as cells	2	5	1	
854	Paste gathered code as text	3	4	3	
855	Review versions in revision browser	3	1	5	
856	Recover deleted / overwritten cells	4	1	7	
857		# participants			

Task 2: Exploratory Data Analysis



How useful was this feature to this task?

- Very useful
- Somewhat useful
- Not useful
- No basis to answer

Figure 6: Analysts found code gathering tools most useful for gathering code to new notebooks, when they cleaned notebooks, and when they performed exploratory analysis. Analysts also appreciated dependency highlights, especially when they were cleaning code.

no relevant code was deleted or because the user interface recovers deleted code silently.

Valued and versatile feature of gathering code to new notebooks. Nine analysts gathered code to a new notebook at least once during the exploratory task. Analysts gathered code to a notebook a median of 1.5 times ($\sigma = 3.7$) during this task, with one analyst even gathering notebooks 12 times (P3). Analysts most often gathered code to a notebook for its intended purpose of cleaning up their code as a “finishing move” after exploration (P6). Analysts clearly valued this aspect of the tool, calling it “amazing” and “beautiful” (P10), that they “loved it” (P5), it “hits the nail on the head” (P9), and will save them “a lot of time” (P11).

Analysts saw additional value in gathering code to notebooks beyond our original design intentions. During the exploratory task, one analyst used gathering to a new notebook as a lightweight branching scheme. As he explored alternatives, he would gather his preferred alternative to a new notebook to create a clean slate for further exploration (P3). Another analyst used gathering as a way to generate reference material. She created data visualizations, then gathered them to new notebooks, so she could quickly flip to the visualizations as she carried on exploring in her original notebook (P4). Finally, one analyst used gathering to support cleaning for multiple audiences. At the end of the exploratory task, he gathered many visualizations to one notebook and documented them for his peer data analysts; he then gathered his movie ranking result to a different notebook intended for those who only want to know the final answer (P2).

Analysts were eager to incorporate gathering into their data analysis workflows: seven of twelve analysts asked us when we would release the tool. One analyst envisioned gathering becoming part of code-cleaning parlance: “once this is public, people will send you bloated notebooks. I’ll say, nope, you should gather that” (P10).

Use of the dependency highlighting. During the exploratory task, 8 analysts clicked on at least one variable definition, and 9 clicked on at least one output area. Additionally, during

the cleaning tasks, as these tasks involved reading unfamiliar notebooks, a few analysts used the dependency highlights as a way to understand the unfamiliar code.

Use and disuse of the version browser. Two analysts opened the version browser at least once (P2, P3). Both copied the cells to the clipboard from a version in the version browser at least once; one analyst in fact copied cells for versions four times during their session (P2). The other analyst opened a version in a new notebook. This analyst wanted to compare versions of a cell that sorted data based on two different dimensions, and used the version browser to recover code from a prior version without overwriting the current cells, which they wished to preserve (P3).

Some analysts who did not use the version browser believed that they might eventually use it in their own work (P6, P8, P9). One analyst noted Jupyter Notebook’s implementation of “undo” is not sufficient for them, and the version browser could provide some of the backtracking functionality they want (P6). Another reported that the version browser could be useful in their current work, where they have iteratively developed an algorithm and are managing three notebooks containing different versions of analyses (P9). However, two analysts believed they wouldn’t use the version browser, as its view of versions is too restrictive. The version browser collects versions ending with multiple executions of the same cell, yet these analysts preferred to modify and re-run old analyses in new cells (P10, P11).

Downsides and gaps. A few analysts mentioned that repeatedly gathering code to a new notebook creates a different kind of mess, namely clutter across notebooks, rather than clutter within a notebook. For example, gathering multiple times typically causes initialization code (e.g. loading the dataset) to be duplicated in each generated notebook (P3, P4, P6). In effect, a notebook and the notebooks gathered from it form a parent/child relationship that the user interface does not currently recognize. Analysts suggested several improvements. First, gathering to a new notebook should create a provisional notebook, rather than being saved by default,

955 and its name should be related to the original notebook's
 956 name. One analyst suggested linked editing across this family
 957 of notebooks as a way to deal with duplicated code. For
 958 example, renaming a variable in one family member could
 959 automatically rename it in all members (P12).

960 961 **Validating the design motivations**

962 Analysts' feedback offered evidence of the role that our de-
 963 sign motivations played in the usefulness of the tools:
 964

965 *Post-hoc management of messes.* Analysts valued the abil-
 966 ity to manage messes without up-front effort to organize and
 967 version code. This was a benefit of gathering to new note-
 968 books, as analysts appreciated simple affordances to clean
 969 up their messy analysis code (P1, P2, P6, P8, P9). For one
 970 analyst, the tool encouraged them, for better or worse, to
 971 "not to care... too much about data cleaning or structure at
 972 this moment. I say it was nice in a way, that I can just kind of
 973 go on with what I want to do" (P12). For some analysts, this
 974 was the downside of the version browser, which required
 975 them to run new versions of code in the same cell (P10, P11).

976 *Portability of gathered code.* Analysts reused gathered code
 977 by opening fresh notebooks, pasting cells, and pasting plain-
 978 text. In the exploratory task, nine analysts gathered code to
 979 notebooks, and five gathered code to the clipboard, to paste
 980 as either cells or plaintext. By pasting plaintext into one cell,
 981 analysis code looked "a lot cleaner" (P3), and several ana-
 982 lysts wanted an easier ways to gather code to scripts. Others
 983 preferred pasting code as distinct cells (P5, P7). One analyst
 984 simply liked having the choice (P10).

985 *Querying code via direct selection of analysis results.* Ana-
 986 lysts appreciated the directness with which they could gather
 987 code: "It was very easy to just click, click, click on something
 988 and then grab the code that produced" a result (P10). The
 989 directness allowed analysts to clean their code by asking,
 990 "what do I need?" rather than "what do I not need?" (P3).

991 992 **Limitations**

993 Our study has two limits to external validity, which are com-
 994 mon in lab-based usability evaluations: first, the participants
 995 did not do their own work, on their own data, in their own
 996 time frame. We created realistic tasks by choosing notebooks
 997 and datasets from the UCSD Notebook Archive, itself mined
 998 from GitHub. Ideally, participants would use their own data
 999 and analyses. However, several informants in our formative
 1000 interviews said their data was too sensitive for us to observe,
 1001 so we did not pursue this option. The second limitation is the
 1002 study's short duration, which we believe accounts for the
 1003 low use of the version browser feature. As P6 commented,
 1004 "the other features will be more valued for notebooks that
 1005 have been used for a long time/long project."

8 DISCUSSION AND FUTURE WORK

1008 There are several simplifying assumptions embodied in the
 1009 prototype, which could be made more flexible in future tools.
 1010

1011 *Flexible versions of results.* The Revisions button shows dif-
 1012 ferent versions of a result over time, but this raises a question:
 1013 when should we consider two results in the execution log
 1014 to be different versions of the "same" result? Our prototype
 1015 uses the identity of the cell that produces a result to deter-
 1016 mine "sameness". This heuristic works for many styles of
 1017 exploratory programming, but lacks flexibility. Future tools
 1018 could determine "sameness" with heuristics based on text
 1019 content, rather than notebook structure.

1020 *Long-lived execution histories.* The code gathering tools'
 1021 execution log lasts only for a single programming session,
 1022 which limits the scope of the Revisions button and resurrect-
 1023 ing code from deleted or overwritten cells. Future tools could
 1024 create a persistent execution log, both for individual users
 1025 and for sharing. Further, transferring code or cells between
 1026 notebooks (e.g., copy/paste) could also transfer the relevant
 1027 execution steps. This creates both technical challenges (e.g.,
 1028 merging log segments into a coherent history) and user inter-
 1029 face challenges (e.g. scaling the Revisions display to months'
 1030 or years' worth of revisions).

1031 *Cleaning incrementally.* Given the fluid nature of notebook
 1032 programming, we designed the use of code gathering tools to
 1033 be entirely *post hoc*, i.e. a user never needs to anticipate that
 1034 cleaning will happen in the future. While gathering meets
 1035 this goal, the gathering step requires the user to select all
 1036 the results to be saved at once. Several participants wanted
 1037 to create a cleaned notebook in a series of steps, that is, the
 1038 ability for a gathering step to "patch" notebooks gathered
 1039 in a previous step. Future tools could use algorithms from
 1040 revision control systems to support this flexibility.

1041 **Conclusions**

1042 Our qualitative usability study with 12 professional data sci-
 1043 entists confirmed that cleaning computational notebooks is
 1044 primarily about removing unwanted analysis code and re-
 1045 sults. The cleaning task can also involve secondary steps like
 1046 improving code quality, writing documentation, polishing
 1047 results for a new audience, or creating scripts. Participants
 1048 find the primary cleaning task to be clerical and error-prone.
 1049 They therefore responded positively to the code gathering
 1050 tools, which automatically produces the minimal code nec-
 1051 essary to replicate a chosen set of analysis results, using
 1052 a novel application of program slicing. Analysts primarily
 1053 used code gathering as a "finishing move" to share work,
 1054 but also found unanticipated uses like generating reference
 1055 material, creating lightweight branches in their code, and
 1056 creating summaries for multiple audiences.

1057 1058 1059 1060

REFERENCES

- [1] Jison. <https://jison.org> Accessed September 8, 2018.
- [2] Jupyter. <http://jupyter.org/> Accessed September 2, 2018.
- [3] Unofficial Jupyter Notebook Extensions. <https://jupyter-contrib-nbextensions.readthedocs.io/en/latest/>. Accessed September 3, 2018.
- [4] Joel Brandt, Vignan Pattamatta, William Choi, Ben Hsieh, and Scott R. Klemmer. 2010. Rehearse: Helping Programmers Adapt Examples by Visualizing Execution and Highlighting Related Code. (2010). Stanford technical report.
- [5] Brian Burg, Richard Bailey, Andrew J. Ko, and Michael D. Ernst. 2013. Interactive Record/Replay for Web Application Debugging. In *Proceedings of the 26th annual ACM symposium on User interface software and technology*. ACM, 473–484.
- [6] Brian Burg, Andrew J. Ko, and Michael D. Ernst. 2015. Explaining Visual Changes in Web Interfaces. In *Proceedings of the 28th Annual ACM Symposium on User Interface Software & Technology*. ACM, 259–268.
- [7] Steven P. Callahan, Juliana Freire, Emanuele Santos, Carlos E. Scheidegger, Cláudio T. Silva, and Huy T. Vo. 2006. VisTrails: Visualization meets Data Management. In *Proceedings of the 2006 ACM SIGMOD international conference on Management of data*. ACM, 745–747.
- [8] Mihai Codoban, Sruti Srinivasa Ragavan, Danny Dig, and Brian Bailey. 2015. Software History under the Lens: A Study on Why and How Developers Examine It. In *Software Maintenance and Evolution (ICSME), 2015 IEEE International Conference on*. IEEE, 1–10.
- [9] Paul A. Gross, Micah S. Herstand, Jordana W. Hodges, and Caitlin L. Kelleher. 2010. A Code Reuse Interface for Non-Programmer Middle School Students. In *Proceedings of the 15th international conference on Intelligent user interfaces*. ACM, 219–228.
- [10] Philip J. Guo and Margo Seltzer. 2012. BURRITO: Wrapping Your Lab Notebook in Computational Infrastructure. In *Proceedings of the 4th USENIX Workshop on the Theory and Practice of Provenance (TaPP'12)*. USENIX Association, Berkeley, CA, USA. <http://dl.acm.org/citation.cfm?id=2342875.2342882>
- [11] Björn Hartmann, Loren Yu, Abel Allison, Yeonsoo Yang, and Scott R. Klemmer. 2008. Design as exploration: creating interface alternatives through parallel authoring and runtime tuning. In *Proceedings of the 21st annual ACM symposium on User interface software and technology*. ACM, 91–100.
- [12] Andrew Head, Elena L Glassman, Björn Hartmann, and Marti A Hearst. 2018. Interactive Extraction of Examples from Existing Code. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM.
- [13] Joshua Hirschman and Haoqi Zhang. 2015. Unravel: Rapid web application reverse engineering via interaction recording, source tracing, and library detection. In *Proceedings of the 28th Annual ACM Symposium on User Interface Software & Technology*. ACM, 270–279.
- [14] Joshua Hirschman and Haoqi Zhang. 2016. Telescope: Fine-Tuned Discovery of Interactive Web UI Feature Implementation. In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology*. ACM, 233–245.
- [15] Reid Holmes and Robert J. Walker. 2012. Systematizing Pragmatic Software Reuse. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 21, 4 (2012), 20.
- [16] Sean Kandel, Andreas Paepcke, Joseph M. Hellerstein, and Jeffrey Heer. 2012. Enterprise Data Analysis and Visualization: An Interview Study. *IEEE Transactions on Visualization & Computer Graphics* 12 (2012), 2917–2926.
- [17] Kyle Kelley and Brian Granger. 2017. Jupyter frontends: From the classic Jupyter Notebook to JupyterLab, nteract, and beyond. (2017).
- JupyterCon. A video of this talk is available at: <https://www.youtube.com/watch?v=YKmJvHjTGAM>.
- [18] Mary Beth Kery, Amber Horvath, and Brad Myers. 2017. Variolite: Supporting Exploratory Programming by Data Scientists. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*. 1265–1276.
- [19] Mary Beth Kery and Brad A. Myers. 2017. Exploring Exploratory Programming. In *Visual Languages and Human-Centric Computing (VL/HCC), 2017 IEEE Symposium on*. IEEE, 25–29.
- [20] Mary Beth Kery, Marissa Radensky, Mahima Arya, Bonnie E. John, and Brad A. Myers. 2018. The Story in the Notebook: Exploratory Data Science using a Literate Programming Tool. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM.
- [21] Andrew Ko and Brad A. Myers. 2009. Finding Causes of Program Output with the Java Whyline. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, 1569–1578.
- [22] Yun Young Lee, Darko Marinov, and Ralph E. Johnson. 2015. Tempura: Temporal dimension for ides. In *Software Engineering (ICSE), 2015 IEEE/ACM 37th IEEE International Conference on*, Vol. 1. IEEE, 212–222.
- [23] Josip Maras, Maja Štula, Jan Carlson, and Ivica Crnković. 2013. Identifying Code of Individual Features in Client-Side Web Applications. *IEEE Transactions on Software Engineering* 39, 12 (2013), 1680–1697.
- [24] Christopher Oezbek and Lutz Prechelt. 2007. JTouBus: Simplifying Program Understanding by Documentation that Provides Tours Through the Source Code. In *Software Maintenance, 2007. ICSM 2007. IEEE International Conference on*. IEEE, 64–73.
- [25] João Felipe Pimentel, Juliana Freire, Leonardo Murta, and Vanessa Braganholo. 2016. Fine-Grained Provenance Collection over Scripts Through Program Slicing. In *International Provenance and Annotation Workshop*. Springer, 199–203.
- [26] Adam Rule. 2018. *Design and Use of Computational Notebooks*. Ph.D. Dissertation. UC San Diego.
- [27] Adam Rule, Ian Drosos, Aurélien Tabard, and James D. Hollan. 2018. Aiding Collaborative Reuse of Computational Notebooks with Annotated Cell Folding. (2018).
- [28] Adam Rule, Aurélien Tabard, and James D. Hollan. Data from: Exploration and Explanation in Computational Notebooks. <https://doi.org/10.6075/J0JW8C39>.
- [29] Adam Rule, Aurélien Tabard, and James D Hollan. 2018. Exploration and Explanation in Computational Notebooks. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM.
- [30] Francisco Servant and James A. Jones. 2012. History Slicing: Assisting Code-Evolution Tasks. In *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering*. ACM.
- [31] Sruti Srinivasa Ragavan, Sandeep Kaur Kuttal, Charles Hill, Anita Sarma, David Piorkowski, and Margaret Burnett. 2016. Foraging Among an Overabundance of Similar Variants. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. ACM, 3509–3521.
- [32] Ryo Suzuki. 2015. Interactive and Collaborative Source Code Annotation. In *Software Engineering (ICSE), 2015 IEEE/ACM 37th IEEE International Conference on*, Vol. 2. IEEE, 799–800.
- [33] Mark Weiser. 1981. Program slicing. In *Proceedings of the 5th international conference on Software engineering*. IEEE Press, 439–449.
- [34] Mark D. Weiser. 1979. *Program Slices: Formal, Psychological, and Practical Investigations of an Automatic Program Abstraction Method*. Ph.D. Dissertation. <https://search-proquest-com.libproxy.berkeley.edu/docview/302949655?accountid=14496> Copyright - Database copyright ProQuest LLC; ProQuest does not claim copyright in the individual underlying works; Last updated - 2016-06-08.

1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166

- 1167 [35] YoungSeok Yoon and Brad A. Myers. 2012. An Exploratory Study of
1168 Backtracking Strategies Used by Developers. In *Proceedings of the 5th*
1169 *International Workshop on Co-operative and Human Aspects of Software*
1170 *Engineering*. IEEE Press, 138–144.
- 1171
- 1172
- 1173
- 1174
- 1175
- 1176
- 1177
- 1178
- 1179
- 1180
- 1181
- 1182
- 1183
- 1184
- 1185
- 1186
- 1187
- 1188
- 1189
- 1190
- 1191
- 1192
- 1193
- 1194
- 1195
- 1196
- 1197
- 1198
- 1199
- 1200
- 1201
- 1202
- 1203
- 1204
- 1205
- 1206
- 1207
- 1208
- 1209
- 1210
- 1211
- 1212
- 1213
- 1214
- 1215
- 1216
- 1217
- 1218
- 1219
- [36] YoungSeok Yoon and Brad A. Myers. 2015. Supporting Selective Undo
1220 in a Code Editor. In *Software Engineering (ICSE), 2015 IEEE/ACM 37th*
1221 *IEEE International Conference on*, Vol. 1. IEEE, 223–233.
- 1222
- 1223
- 1224
- 1225
- 1226
- 1227
- 1228
- 1229
- 1230
- 1231
- 1232
- 1233
- 1234
- 1235
- 1236
- 1237
- 1238
- 1239
- 1240
- 1241
- 1242
- 1243
- 1244
- 1245
- 1246
- 1247
- 1248
- 1249
- 1250
- 1251
- 1252
- 1253
- 1254
- 1255
- 1256
- 1257
- 1258
- 1259
- 1260
- 1261
- 1262
- 1263
- 1264
- 1265
- 1266
- 1267
- 1268
- 1269
- 1270
- 1271
- 1272