# **Skip Blocks**: Reusing Execution History to Accelerate Web Scripts
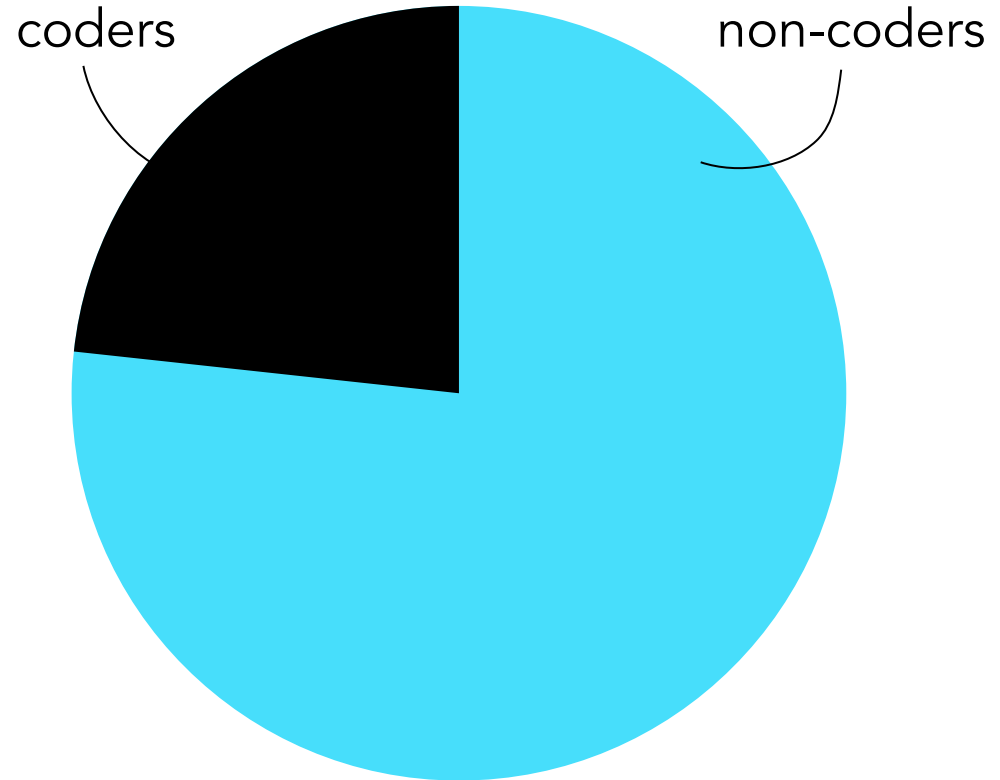
## Sarah Chasins
University of California, Berkeley
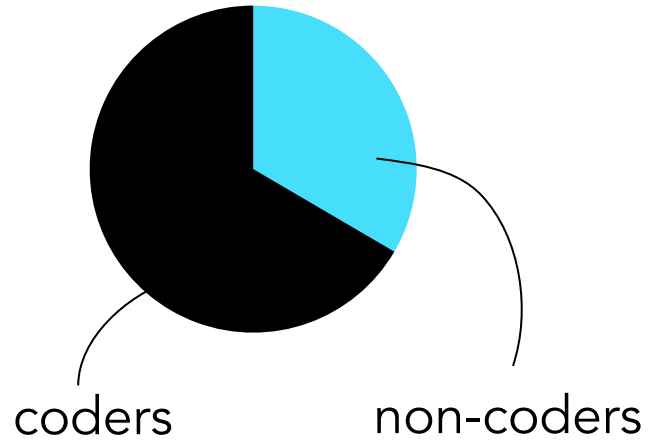
## Rastislav Bodik
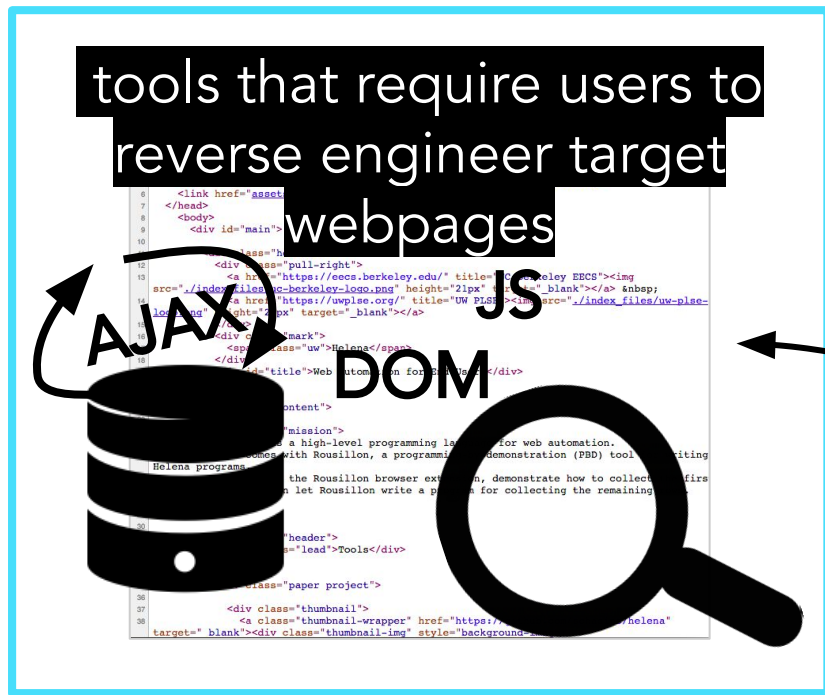University of Washington

care about data today

care about data tomorrow

coders

non-coders

coders

non-coders

# What web data collection tools do we have?

tools that require users to reverse engineer target webpages

JS DOM

AJAX

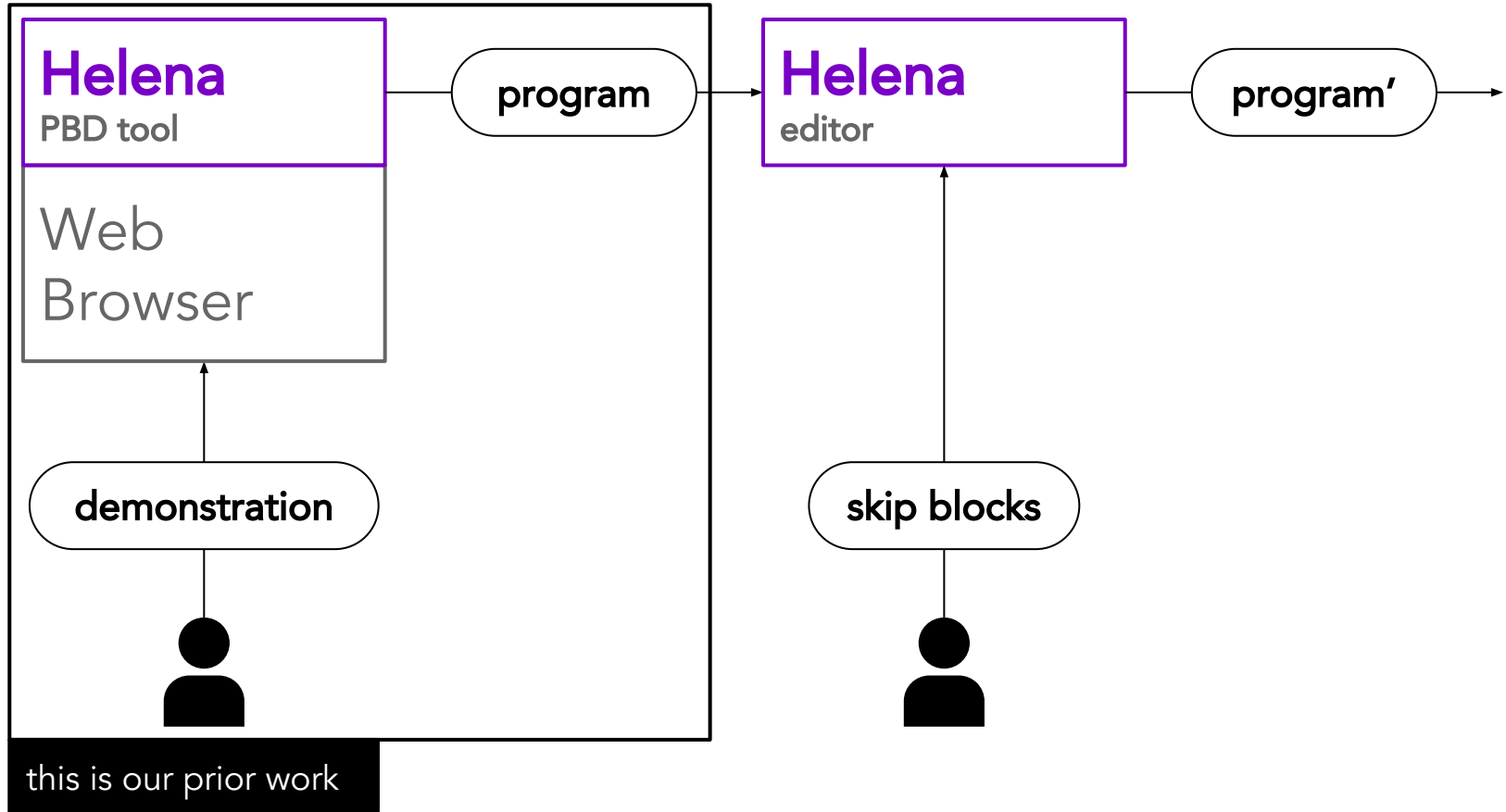**coders**

- hire a human to copy & paste
- hire a coder to use one of these

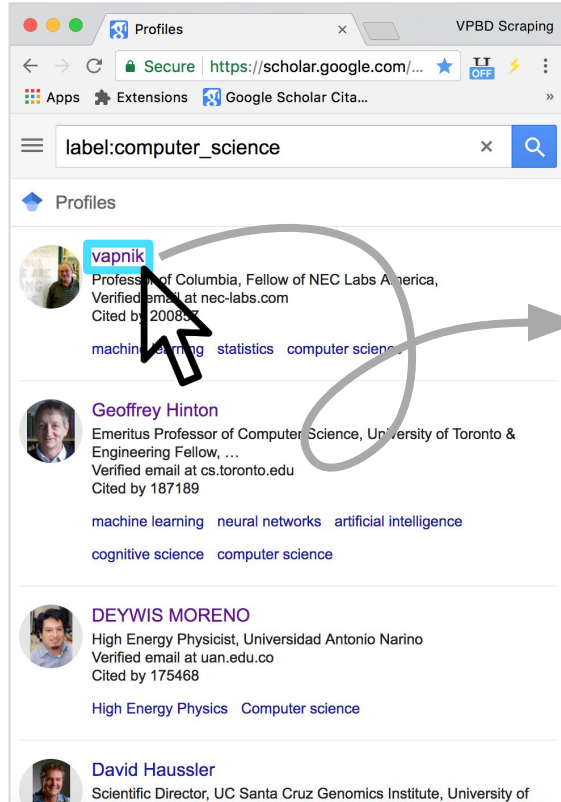- **Helena**
  WEB AUTOMATION
  FOR END USERS

  **our tool!**

**non-coders**
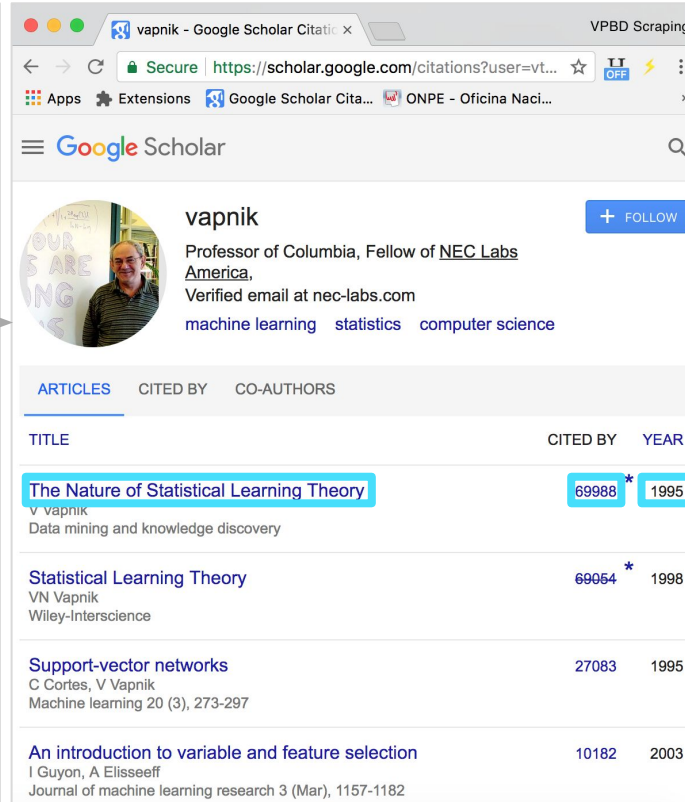
Helena
PBD tool

Web Browser

program

Helena
editor

program'

demonstration

skip blocks

this is our prior work

4

# Let's PBD a web automation script!

Goal: scrape all papers by top 10,000 CS authors from Google Scholar

DEPARTMENT OF SOCIOLOGY
UNIVERSITY *of* WASHINGTON

+

City of Seattle

How is rent changing across Seattle neighborhoods?

**Kept losing network connection**

page 1

page 2

New listings have pushed the last three listings from p1 onto p2

wasting 10+ hours scraping duplicates!

DEPARTMENT OF ECONOMICS

UNIVERSITY *of* WASHINGTON

How is the minimum wage affecting Seattle restaurants?



CIVIL & ENVIRONMENTAL ENGINEERING

UNIVERSITY *of* WASHINGTON

Can we design a better carpool matching algorithm?



EVANS SCHOOL OF PUBLIC POLICY & GOVERNANCE

UNIVERSITY *of* WASHINGTON

How do charitable foundations communicate with supporters?

# Problem Statement

(1) **Failures**: What happens when the network fails, the server fails, the computer fails?  When we lose our session with the server and have to start over?

(2) **Data changes**: What happens when the server gives the client pages produced from different (potentially conflicting) reads of the underlying data store?

not client side problems → scraping script can't prevent them, must handle them

# Solution



on the surface, seem like very different problems

**failures**                    **data changes**

"Just don't redo the same work you've already done!"

**But what's the 'same' work?  After all, data changes...**

Our answer: the skip block! User can
- tell us what makes objects the same
- associate the code that operates on an object

- If object already committed (memoized), skip block; else, run block
- No reverse engineering!  Reasoning about output data

# Recovering from Failures

```
for (aRow in p1.authors){

        scrape aRow.author_name
        scrape aRow.author_institution
        p2 = click aRow.author_name
        for (pRow in p2.papers){
            scrape pRow.title
            scrape pRow.citations
            output([aRow.author_name, pRow.title, pRow.citations])
        }

}
```

scrape stuff about the author, click the author

for the author's papers, scrape paper stuff

add a row of output with the author and paper info

13

# Recovering from Failures

**key attributes**: is the current author the same as another we've already seen?

```
for (aRow in p1.authors){
    skipBlock(Author(aRow.author_name, aRow.author_institution)){
        scrape aRow.author_name
        scrape aRow.author_institution
        p2 = click aRow.author_name
        for (pRow in p2.papers){
            scrape pRow.title
            scrape pRow.citations
            output([aRow.author_name, pRow.title, pRow.citations])
        }
    }
}
```

**block**: the code that operates on the author object

if ever, **in any run**, script has committed an object with the same key attributes, skips the block

# Recovering from Failures

```
for (aRow in p1.authors){
    skipBlock(Author(aRow.author_name, aRow.author_institution)){
        scrape aRow.author_name
        scrape aRow.author_institution
        p2 = click aRow.author_name
        for (pRow in p2.papers){
            scrape pRow.title
            scrape pRow.citations
            output([aRow.author_name, pRow.title, pRow.citations])
        }
    }
}
```

AUTHOR

PAPER

page load

skip block commit

always at least one page load per author (to load paper list), but often ≈ 40

time

a1  p1  p2  ...  p21  p22  ...  p41  p42  ...  a2  p1  p2  ...  p21  p22  ...  p41  p42  ...

# Recovering from Failures

external
failure point

didn't reach
this commit

recovery **without** the author skip block

recovery **with** the author skip block

"fast-forwarding" over prior work

skips 40
page loads

10 authors per page, so
just 1 page load by this
point, 200 skipped

page
load

skip block
commit

always at least one page load per author
(to load paper list), but often ≈ 40

time

a1    p1  p2  ...    p21  p22  ...    p41  p42  ...    a2    p1  p2  ...    p21  p22  ...    p41  p42  ...

16

# Nested Skip Blocks



| City 1 | — | Restaurant A | — | Review i |
| | | Restaurant B | | Review ii |
| | | ... | | ... |

| City 2 | — | Restaurant C | — | Review iii |
| | | Restaurant D | | Review iv |
| | | ... | | ... |

In authors vs. papers, authors is clearly the right level for the skip block. But here?

skip block only at `city` → scraping a whole city takes many hours, so scraping half a city also takes hours

skip block only at `restaurant` → iterating through a city's restaurant list takes a long time, and now we have to go through all of Seattle, San Francisco before we can resume in the middle of Vancouver

skip block at `city` & `restaurant` → adjustable granularity skipping

# Nested Skip Blocks

```
for (aRow in p1.authors){
    skipBlock(Author(aRow.author_name, aRow.author_institution)){
        scrape aRow.author_name
        scrape aRow.author_institution
        p2 = click aRow.author_name
        for (pRow in p2.papers){
            skipBlock(Paper(pRow.title, pRow.year)){
                scrape pRow.title
                scrape pRow.citations
                output([aRow.author_name, pRow.title, pRow.citations])
            }
        }
    }
}
```
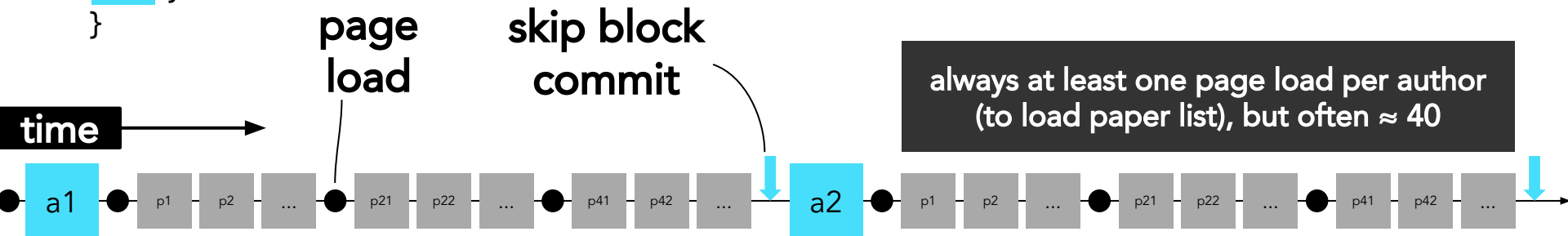
and the inner block may commit even if the outer doesn't - like a nested open transaction

# Refreshing a Dataset

this is the default **staleness threshold**
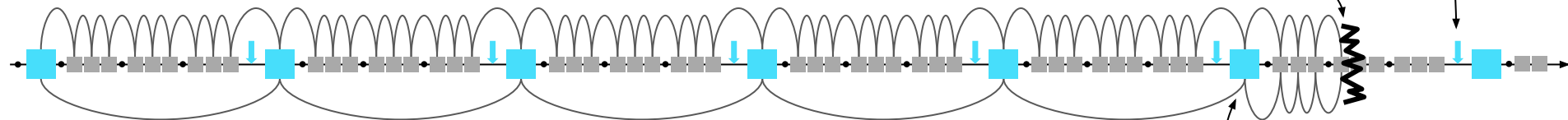
```
for (aRow in p1.authors){
    skipBlock(Author(aRow.author_name, aRow.author_institution), -∞)){
        scrape aRow.author_name
        scrape aRow.author_institution
        p2 = click aRow.author_name
        for (pRow in p2.papers){
            skipBlock(Paper(pRow.title, pRow.year)){
                scrape pRow.title
                scrape pRow.citations
                output([aRow.author_name, pRow.title, pRow.citations])
            }
        }
    }
}
```

**-∞ means skip any duplicate we've seen ever**

If we're scraping once a week, we don't want to revisit each author.  But after a year, maybe we should see what's new.

# Refreshing a Dataset

```
for (aRow in p1.authors){
    skipBlock(Author(aRow.author_name, aRow.author_institution), now - 365*24*60)){
        scrape aRow.author_name
        scrape aRow.author_institution
        p2 = click aRow.author_name
        for (pRow in p2.papers){
            skipBlock(Paper(pRow.title, pRow.year)){
                scrape pRow.title
                scrape pRow.citations
                output([aRow.author_name, pRow.title, pRow.citations])
            }
        }
    }
}
```

Also have logical time (ex: last 3 runs)

Bonus!  In addition to failure recovery and data redundancy handling, get incremental/longitudinal scraping!

# Demo time!

# Benchmark Suite

Need web data? → **Urban@UW** → benchmark suite: 7 long-running web scraping tasks

Ex: for 50 top foundations, scrape the last 1,000 tweets they tweeted

Ex: scrape all Seattle apartment listings from Craigslist

# Data Change
within one run

Measured full execution time of:
- Script with skip blocks
- Script without skip blocks

Chart shows speedup from using skip blocks

**higher is better**

**1.7x** Skipping one ad skips one page load, and pagination gives us so many duplicate ads!

**0.9x** All overhead, no gains - skipping a tweet doesn't skip any page loads!

Speedup (y-axis): 0.0, 0.5, 1.0, 1.5, 2.0, 2.5

Categories (x-axis): Twitter, Yelp Restaurants, Yelp Menus, Foundations, Yelp Reviews, (peace icon), Zimride Carpool

# Data Re-Scraping
within multiple runs

Executed script with skip blocks
**One week later,** measured full execution time of:
- Script with skip blocks
- Script without skip blocks

Chart shows speedup from using skip blocks

higher is better

**49x** Lots of benefit from last week's data. Gates Foundation doesn't post that many new tweets in a week!

**1.9x** Little additional benefit from last week's data; ≈ same speedup as first run.



Speedup axis: $10^0$, $10^1$, $10^2$, $10^3$, 0

Categories: (peace icon), Yelp Menus, Yelp Reviews, Yelp Restaurants, (twitter icon), Foundations, Zimride Carpool

# Failure Recovery
with skip block fast-forwarding

For each benchmark, for three failure locations, the execution time of:
- Script that recovers by naive restarting
- Script that recovers by skip block fastforwarding

Normalized by execution time of a script that doesn't encounter failures

lower is better

overall, performance close to ideal!

execution time if there's no failure

failure during high churn → see new data → slower recovery

Legend:
- Error Loc 1, Skip Block
- Error Loc 1, Naive
- Error Loc 2, Skip Block
- Error Loc 2, Naive
- Error Loc 3, Skip Block
- Error Loc 3, Naive

Y-axis: Execution Time (0.0, 0.5, 1.0, 1.5, 2.0, 2.5)

X-axis benchmarks: Community Foundations, Craigslist, Twitter, Yelp Menus, Yelp Restaurants, Zimride Carpool

# User Study

**Detecting Duplicates**

See a duplicate detection tutorial

| name text | name link | price text | price link | description text | description link | reviews text | reviews link | photos text | photos link |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ |
| Grilled Pork | https://www.yelp.com/menu/paseo-caribbean-food-fremont-seattle-2/item/grilled-pork | $8.50 | | cubed pork loin grilled over lava rocks & basted w/ paseo marinade until golden brown. | | 130 reviews | https://www.yelp.com/menu/paseo-caribbean-food-fremont-seattle-2/item/grilled-pork#menu-reviews | 18 photos | https://www.yelp.com/menu/paseo-caribbean-food-fremont-seattle-2/item/grilled-pork |

Add Annotation

**5: Menu Items**

Each row represents a menu item. Some items appear more than once. Pick columns that identify unique menu items.

| name | price | photos | photos_link | reviews | reviews_link | description |
|---|---|---|---|---|---|---|
| ☐ | ☐ | ☐ | ☐ | ☐ | ☐ | ☐ |
| Grilled Pork | $8.50 | 18 photos | https://www.yelp.com/menu/paseo-caribbean-food-fremont-seattle-2/item/grilled-pork | 130 reviews | https://www.yelp.com/menu/paseo-caribbean-food-fremont-seattle-2/item/grilled-pork#menu-reviews | cubed pork loin grilled over lava rocks & basted w/ paseo marinade until golden brown. |

You can explore the data source here if you want to see more items: https://www.yelp.com/menu/paseo-caribbean-food-fremont-seattle-2

Prev | Next | Sample Question | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | Submit
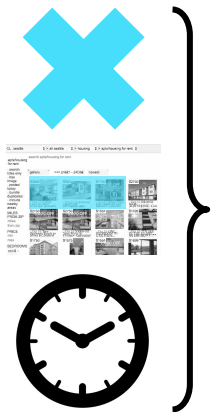
the UI in the online survey

# User Study

If a participant uses the Helena UI to add a skip block, doesn't adjust the default skip block parameters, how many rows of output data are wrong?

**difference not statistically significant**

**time to write each skip block:**
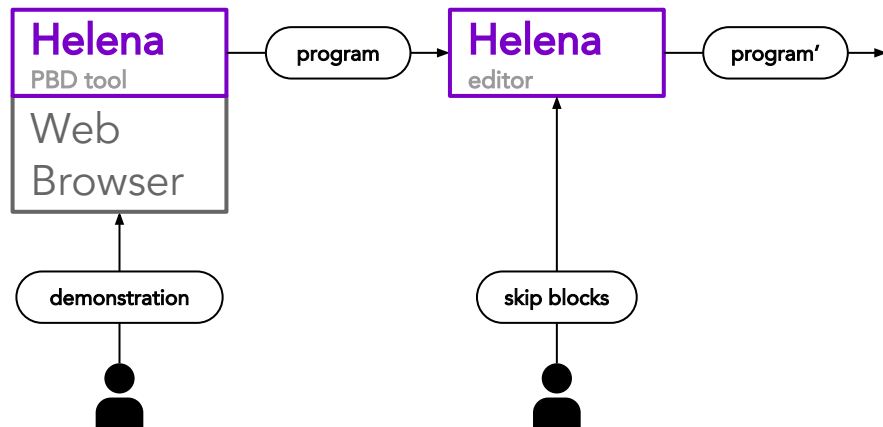**coders: 52 seconds**
**non-coders: 61 seconds**

|  | coders | non-coders |
|---|---|---|
| % of rows **kept** that should have been skipped | 0% | 0% |
| % of rows **skipped** that should have been kept | 1.3% | 2.3% |

| | Snake River Farms Kurobota Ham* | $15.00 |
| | 5 reviews   4 photos | |
| | Avocado and Roma Tomatoes* | $14.00 |
| | 10 reviews   4 photos | |

Unified handling of three apparently disparate challenges with a **single language construct**.

By keeping reasoning at the level of target output data, made skip blocks **usable by non-programmers**.

Helena
PBD tool
Web Browser

program

Helena
editor

program'

demonstration

skip blocks

helena-lang.org

contact: schasins@cs.berkeley.edu

github.com/schasins/helena