

CoWRANGLER: Recommender System for Data-Wrangling Scripts



Bhavya Chopra, **Anna Fariha**, Sumit Gulwani, Austin Z. Henley, Daniel Perelman, Mohammad Raza, Sherry Shi, Danny Simmons, Ashish Tiwari

Data wrangling is time consuming

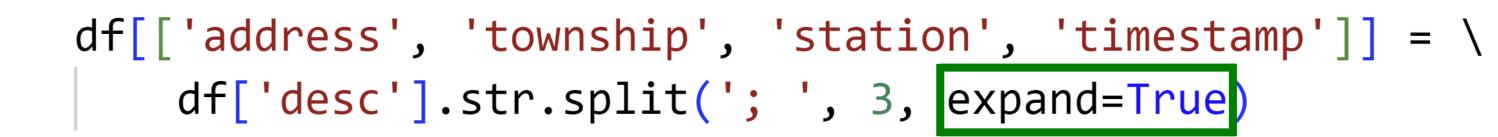
emergency_call_description REINDEER CT & DEAD END; NEW HANOVER; Station 332; 2015-12-10 @ 17:10:52 HAWS AVE; NORRISTOWN; Station STA27; 2015-12-10 @ 14:39:21 BLUEROUTE & RAMP 1476 NB TO CHEMICAL RD; PLYMOUTH; ; 2015-12-10 @ 17:35:41



I want to split the values using semicolons. But how do I deal with the third row, which has fewer splits?

```
df[['address', 'township', 'station', 'timestamp']] = \
   df['desc'].str.split('; ', 3)
ValueError: Columns must be same length as key
```

My script worked after searching the web! Setting the expand parameter to True lets each missing split occupy a column.





It took me 10 minutes to wrangle just one column. There are 8 more!

How do I know if my script is working properly?

Is my data clean now?

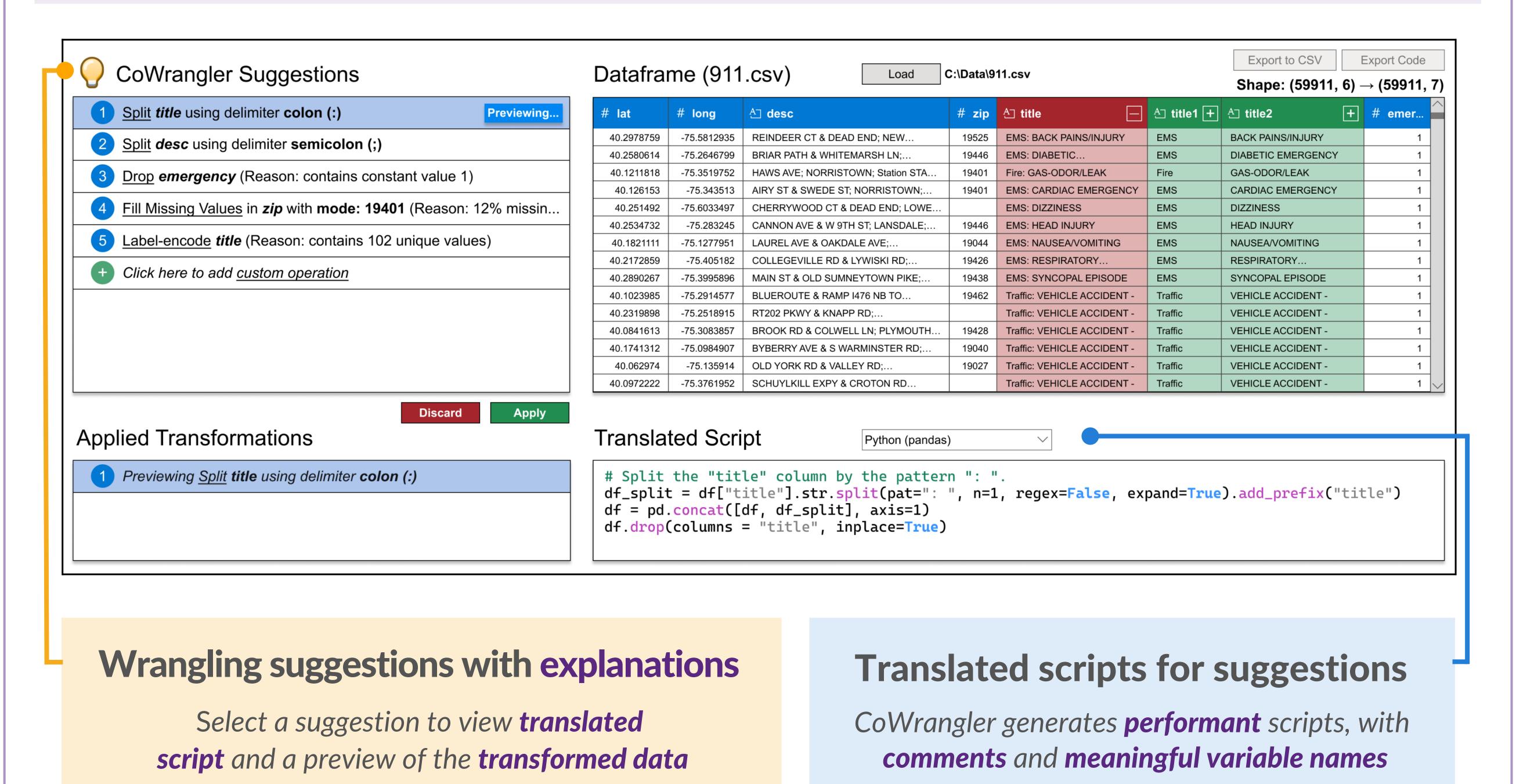
How will I navigate this API jungle? There are so many APIs and parameters!

How do I debug cryptic error messages like ValueError?

Automatic wrangling has several challenges

- How to pick relevant wrangling transformations from the enormous space of valid transformations?
- What metric to use to rank them?
- How to translate them to human-readable scripts?
- How to make the scripts indistinguishable from human-authored scripts?
- How to make sure the scripts are performant?

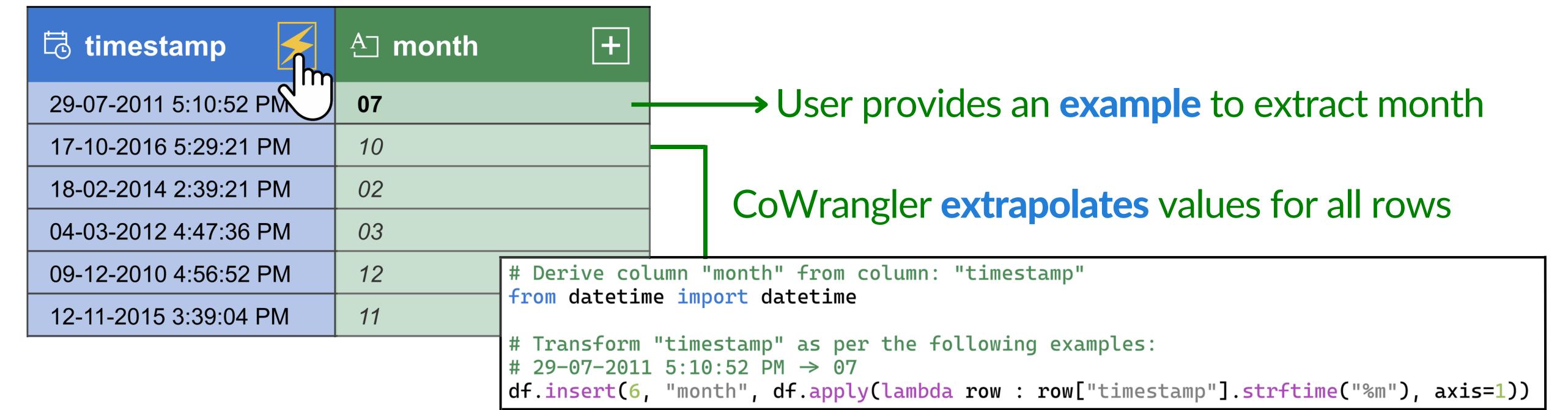
CoWrangling scripts



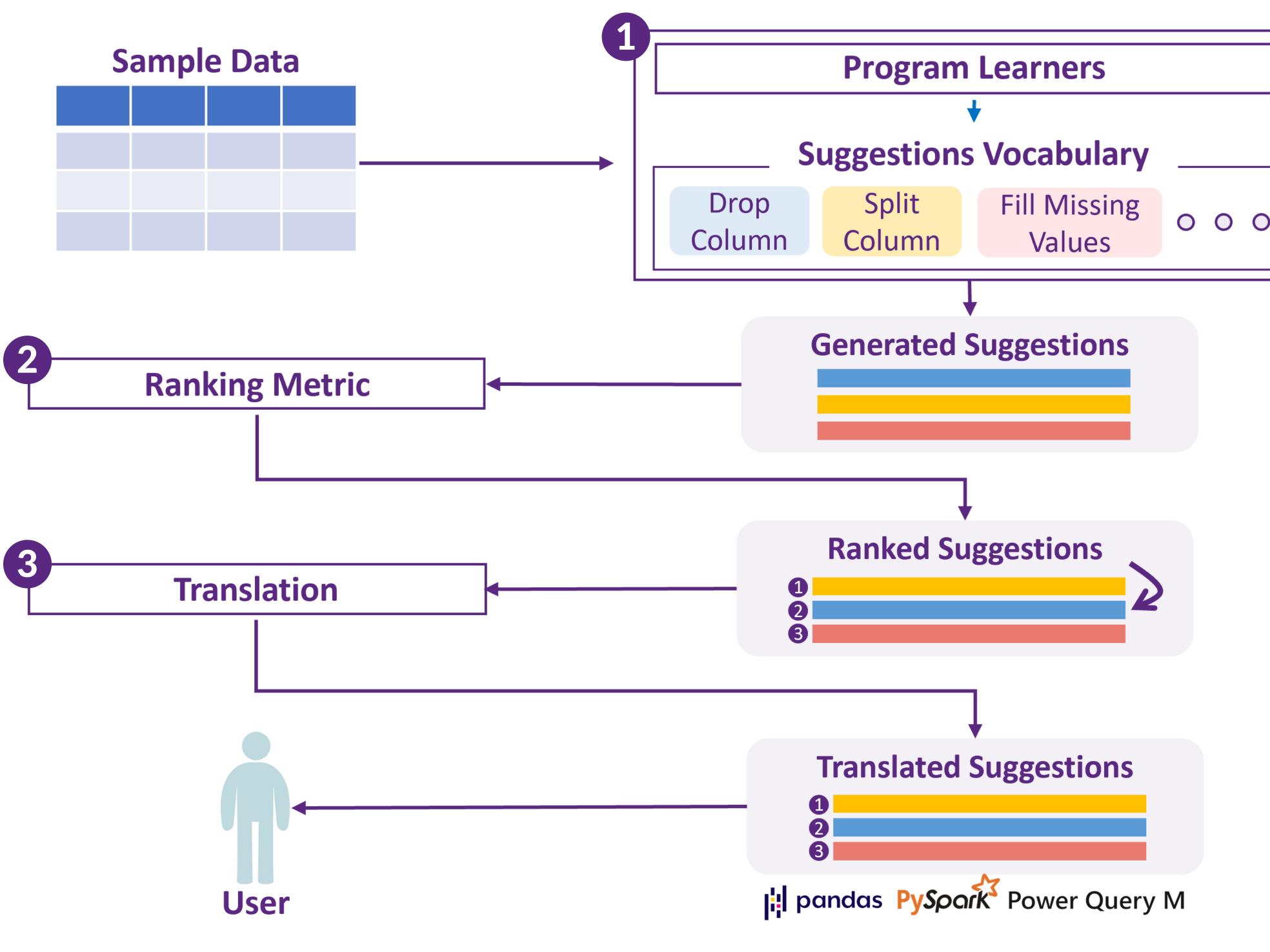
CoWrangler enables human-in-the-loop wrangling

Data scientists can customize and interact with CoWrangler suggestions by

- Editing the scripts generated by CoWrangler
- Expressing intent by example [1]



CoWrangler in a nutshell



- 1 Predictive Synthesis to learn & suggest valid transformations [2]
- 2 Ranking suggestions based on measurable data quality improvements
- 3 Smart translation to generate performant scripts, using vector APIs

CoWrangler achieves 53% accuracy

Benchmark: 2248 pandas operations from 730 Kaggle notebooks

- CoWrangler's vocabulary supports 33% operations
- Suggestions are accurate in 53.4% cases

References

- [1] S Gulwani. Automating string processing in spreadsheets using input-output examples. POPL 2011
- [2] M Raza and S Gulwani. Automated data extraction using predictive program synthesis. AAAI 2017