# With great Parsl comes great portability:

*Using Parsl through **CytoTable** for harmonizing single-cell data*
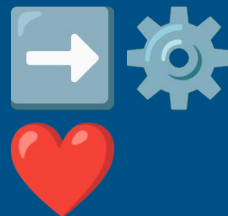
# Introduction



- **Dave Bunten**
  Principal Research Software Engineer

- **Department of Biomedical Informatics**
  University of Colorado Anschutz Medical Campus School of Medicine
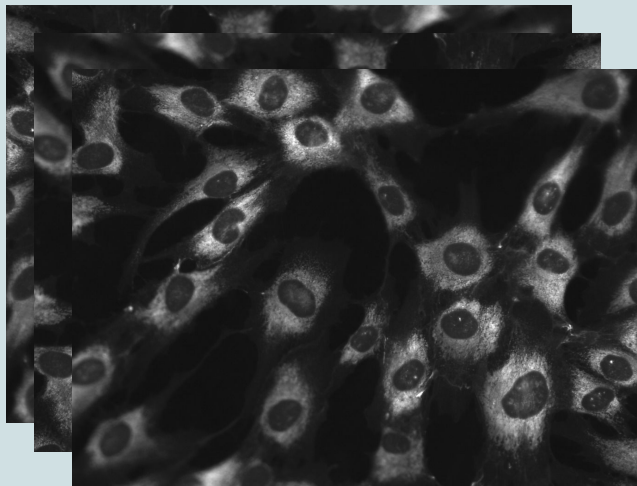
- **Way Lab (Gregory P. Way)**

# Background

- *Bioinformatics:*
  **Image-based profiling**

- **Images of cells
  to numeric data
  (1000's of features)**

# Why?

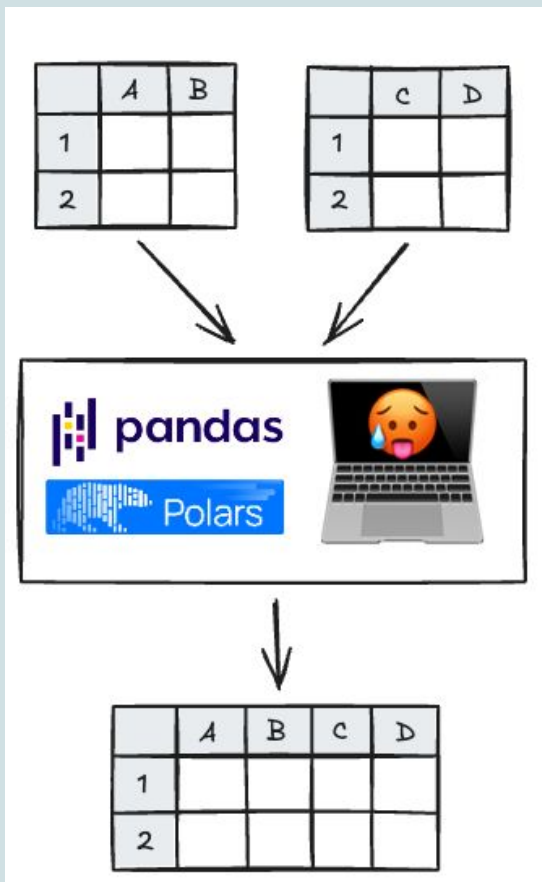- Data derived from images help us understand **biological phenomena**.

- Information is used to help **cure diseases** and **improve lives**.

- Images are **cost-effective** to produce.



(_Travers et al., 2025, Cell Painting and Machine Learning Distinguish Fibroblasts From Nonfailing and Failing Human Hearts_)
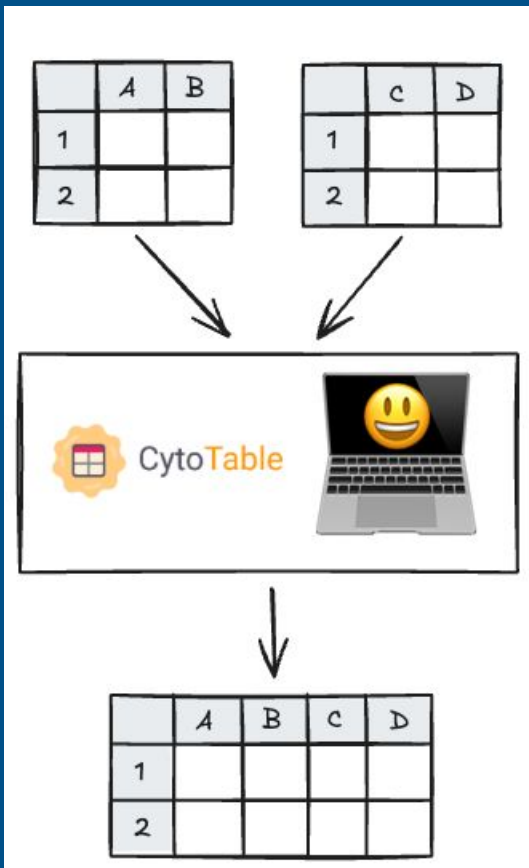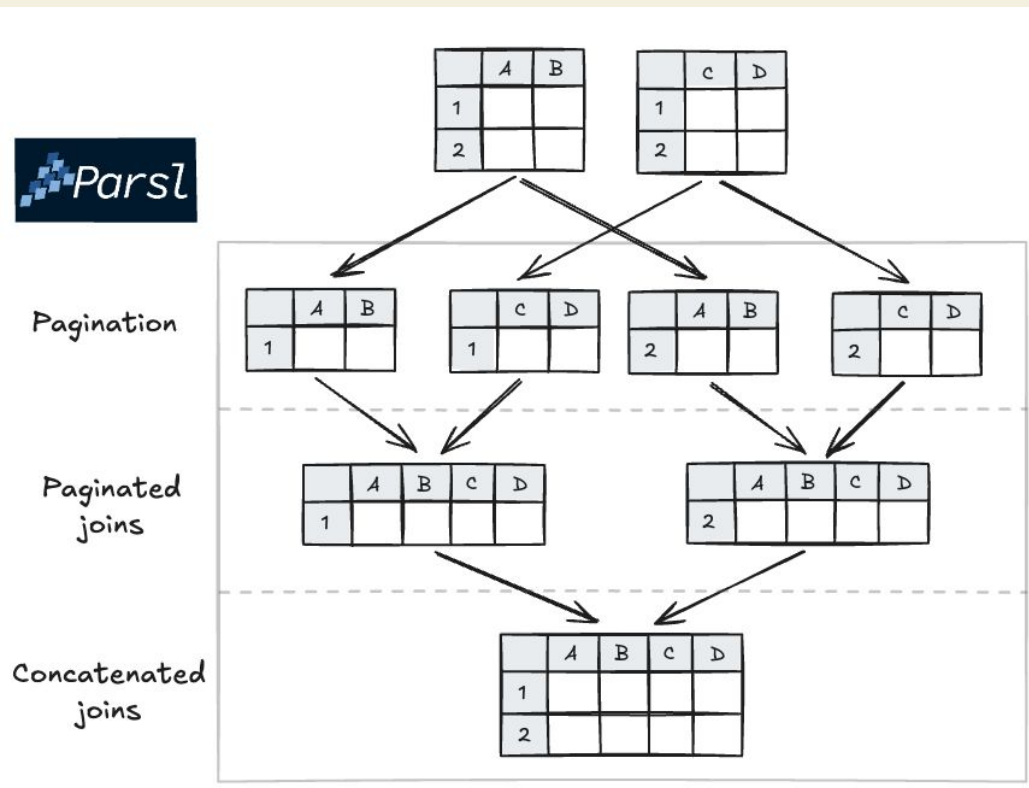
4

# Challenge

- Merging data is **expensive** and **expressed differently** in DataFrame libraries.

- Single-cell feature data entail **many different schema and file formats** from various ecosystems and microscopy products.

# Solution

- **CytoTable** addresses these challenges through scalable **data harmonization**.

- "[Data harmonization is] ... the practice of combining different datasets to maximize their **comparability or compatibility**" (Cheng et al., 2024)

# CytoTable + Parsl



- CytoTable uses Parsl to orchestrate **paginated data harmonization.**

- We implement map–reduce through **paginated maps and concatenated reduction**.
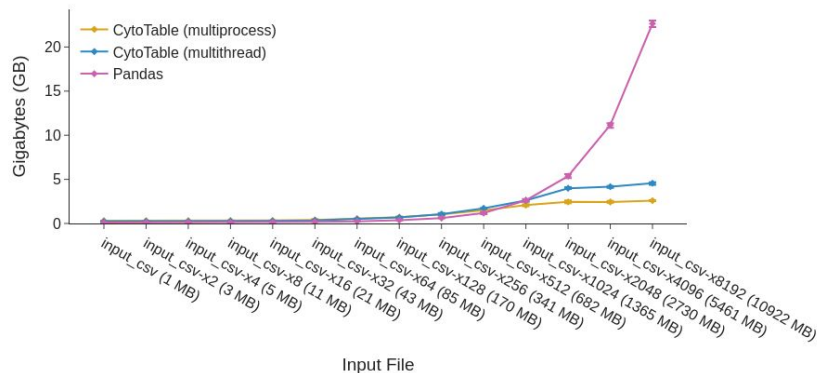
# Portability + scale (!)

- CytoTable is managed using PyPI packages, including Parsl.



- Parsl config within CytoTable includes defaults with flexibility for overrides.



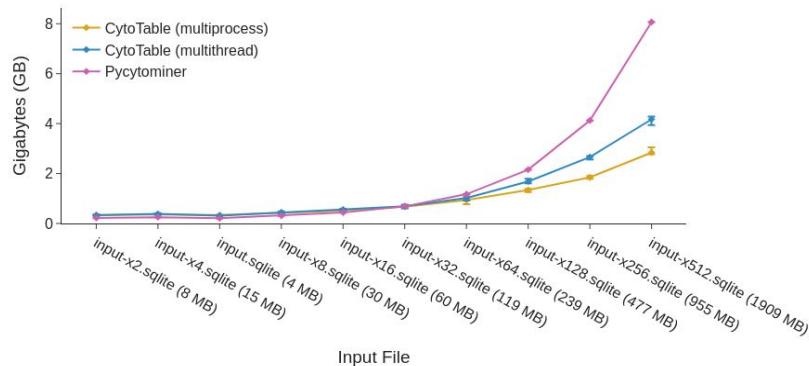- Parsl is compatible with HPC, reducing compute barriers.

CytoTable and Pandas CSV Peak Memory with Min/Max Errors

CytoTable and Pycytominer SQLite Peak Memory with Min/Max Errors

# **Solutions at Scale**

- We find that CytoTable **enables scalable memory and time performance** when compared to existing methods.

# Thank you!

# Questions/comments?

# Find more here!



Preprint:
https://www.biorxiv.org/content/10.1101/2025.06.19.660613v1



GitHub repository:
https://github.com/cytomining/CytoTable