

## Overview

### Idea:

Centrally store and organize all data and analysis results in a relational database  
Keep results of **all** intermediate stages of analysis  
[e.g. fitting a correlator – keep fits to all fit ranges, not just the best fit range]  
No black boxes: structure DB so that analysis can be traced from raw data to final output

### Why?

Enforces standardization for collaborations – everyone uses a common format, all data and results live in one place. [Prevents errors converting formats!]  
Structure helps avoid “stupid mistakes” like associating data from different trajectories, working with an incomplete dataset, using outdated results, etc.  
Reproducibility & provenance – entire analysis kept in DB, easy to “check your work”.  
Can keep metadata (e.g. code version, run date) associated without it getting in the way.  
Automated analysis becomes straightforward

## Relational Databases

In practice, this means: SQL databases  
We use PostgreSQL [open source!]

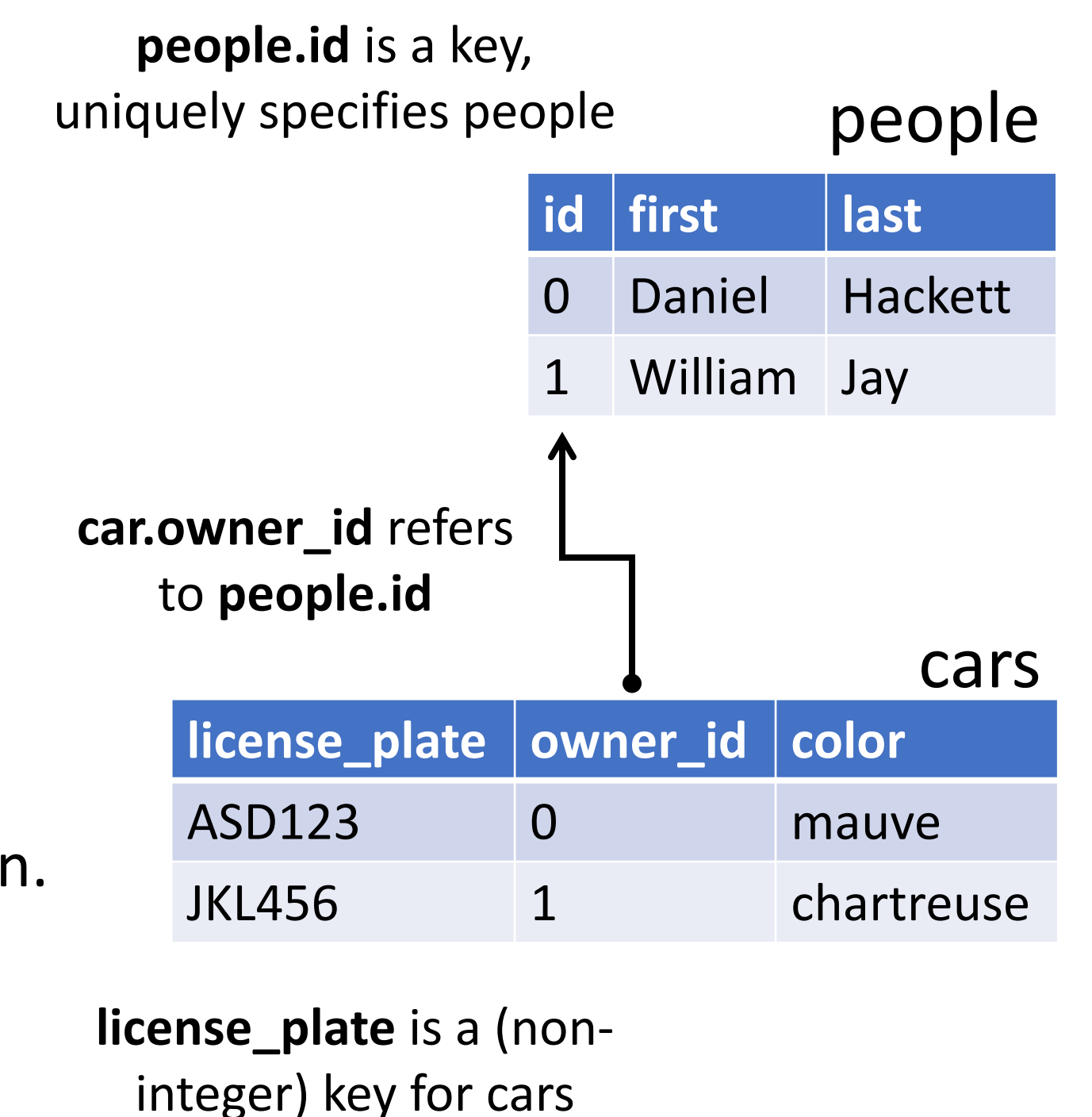


Data is stored in **tables**

**Key** columns provide unique specifiers for each entry  
Tables can refer to each other’s key columns  
→ Use this to encode structure in data

Must specify table structure and relations at initialization.  
After, any data put in must conform to these types and relations [Feature, not a bug: enforces conventions!]

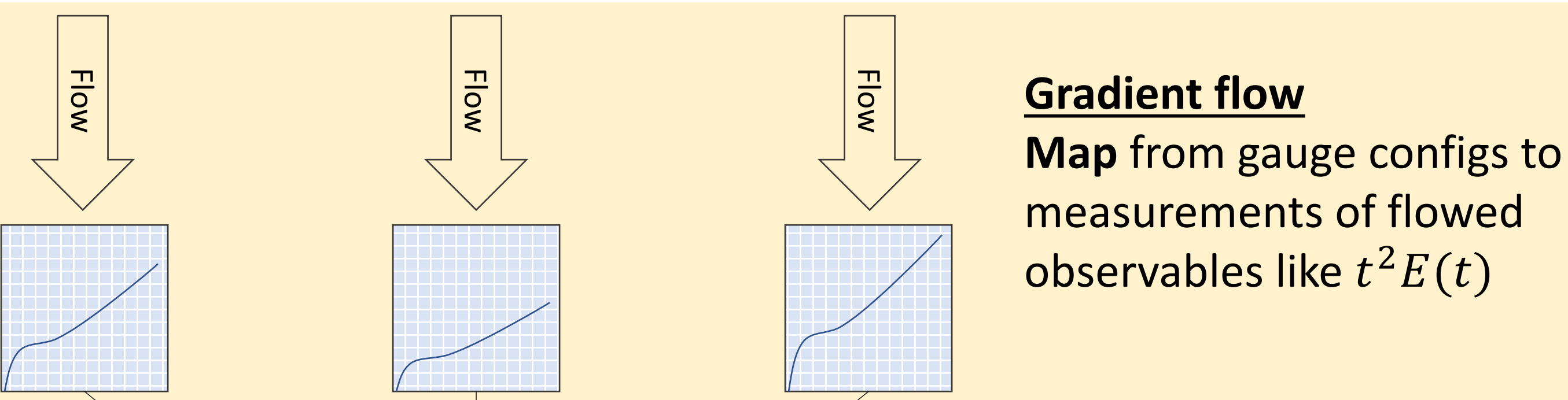
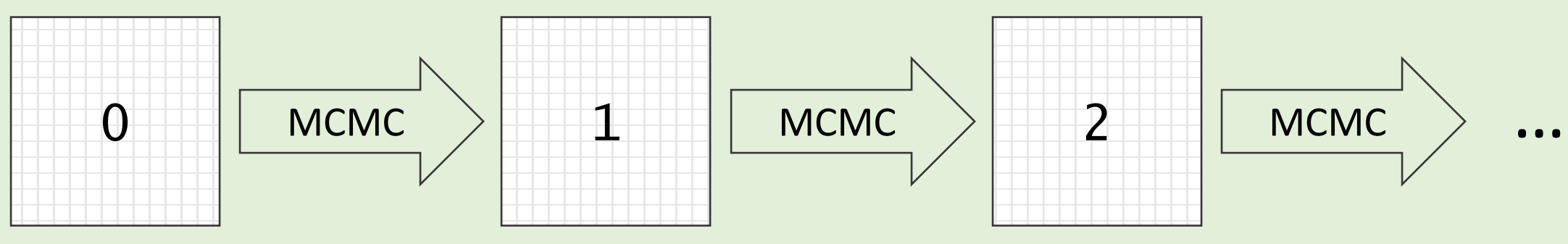
Retrieve data from DB using SQL query language



## Abstract Structure

### MCMC

Map from gauge config to next gauge config  
⇒ Ensemble ~ configs generated with same physical & MCMC parameters



### Scale estimation

Reduce from ensemble of measurements to estimates of observables

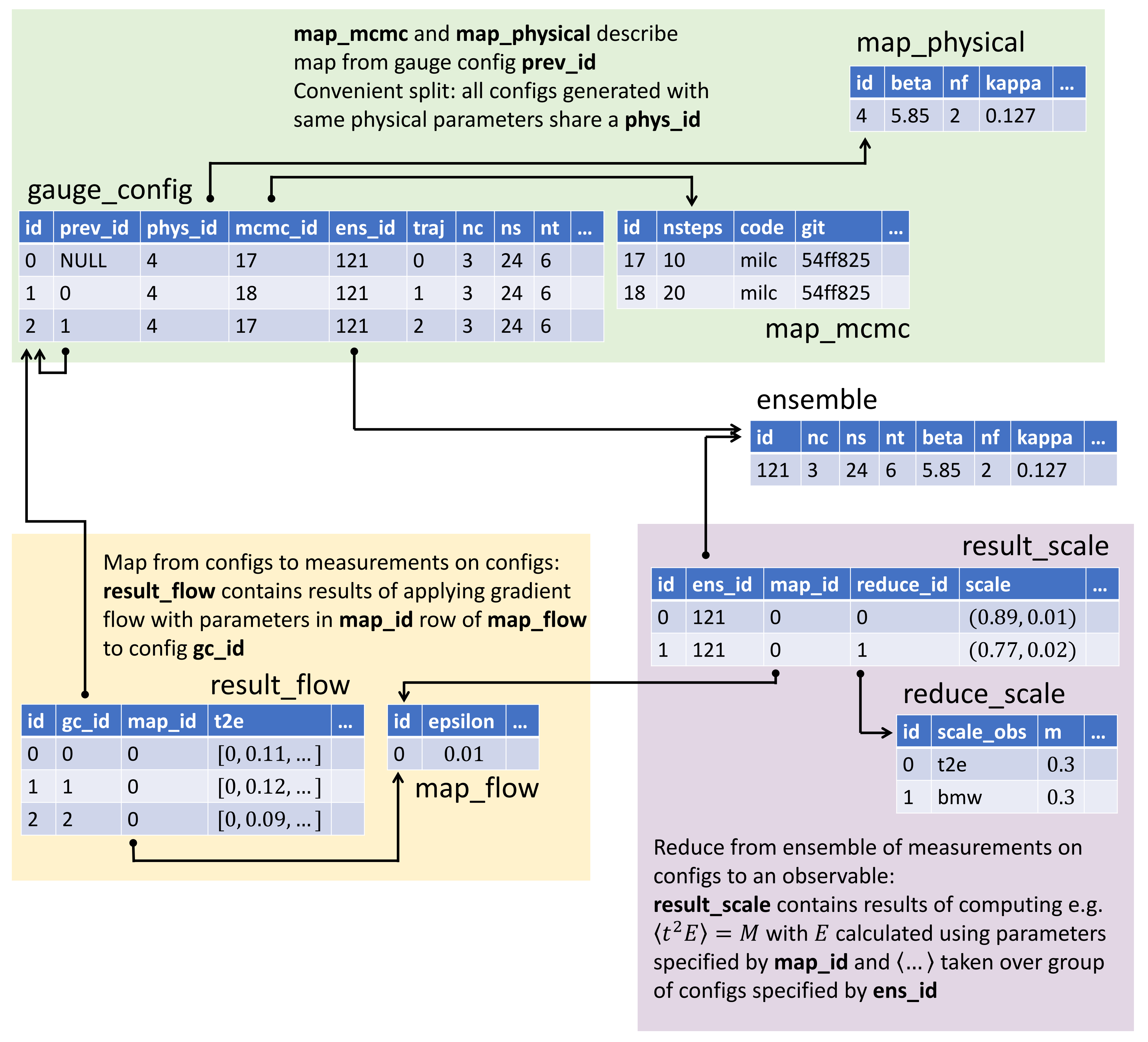
Make equilibration cut, average observables, solve e.g.  $\langle t^2 E(t) \rangle = M$

Reduction produces different estimates for different

- ensembles
- equilibration cuts
- scale observables [e.g.  $\langle t^2 E \rangle$ ,  $t \partial_t \langle t^2 E \rangle$ , ...]
- “magic number”  $M$  [e.g. 0.3 for  $t_0, w_0$ ]

$$\frac{t_0}{a^2} \quad \frac{w_0}{a^2} \quad \dots$$

## Example Lattice Database



## Smart Updating with Hashes

To start: assign a unique **output\_hash** to each gauge config [e.g. checksum, **gauge\_config.id**]

Every result has an **input\_hash** and an **output\_hash**.

**input\_hash**: Hash together **output\_hashes** of all inputs & **map\_ids/reduce\_ids**

**output\_hash**: Hash together all outputs of a measurement/reduction with the **input\_hash**

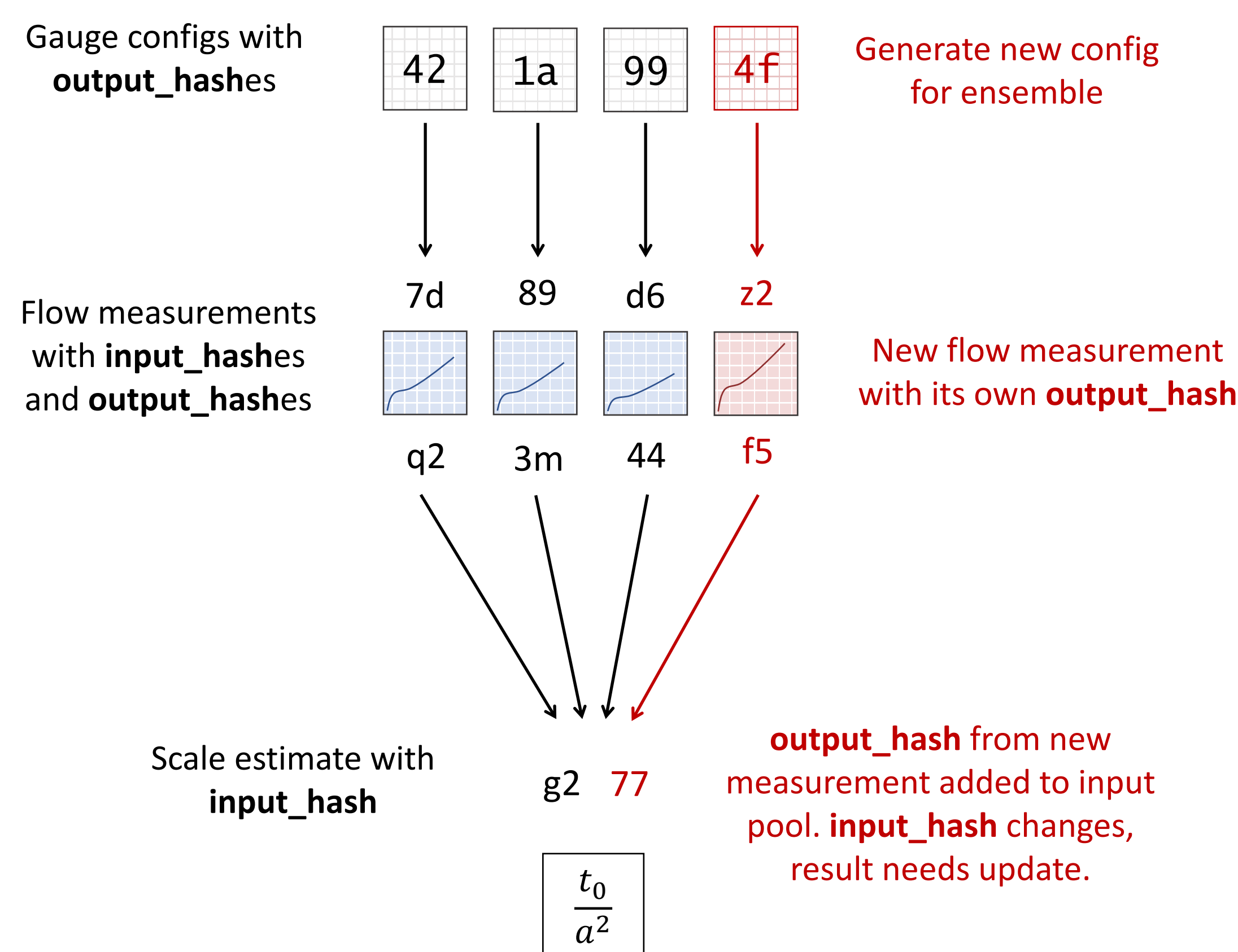
Compute hashes outwards from **gauge\_config** to determine which analyses need update

**input\_hash** changes

→ something upstream has changed – analysis needs update

**output\_hash** changes but **input\_hash** does not

→ analysis has been corrupted – easy integrity checking



## Automation

Run “bulk analysis” scripts N times daily:

- Check for newly-generated raw data and sync to DB
  - Determine which analyses need to be updated/new analyses to perform
  - Perform necessary analyses (e.g. computing scales, fitting correlators, picking best fits)
- No human intervention required to keep analysis up-to-date (wake up to fresh results every morning)

**Closing the automation loop:**

- Workflow manager** runs simulations [e.g. taxi github.com/dchackett/taxi]
- Data stored and automatically analyzed in database
- Automated run-specifier looks at analysis in DB, tells workflow manager to launch new simulations
- Applications:
  - Tuning lattice spacing,  $m_q$ , etc
  - Phase diagram exploration (ask about video!)