



Scalable, Distributed, and Reproducible Machine Learning

Daniel Whitenack, [@dwhitena](https://twitter.com/dwhitena)
Data Scientist and Advocate, [@pachydermIO](https://twitter.com/pachydermIO)

Outline

1. ML/AI challenges
2. Demo with Python and Pachyderm
3. Resources

Reproducibility

generate_figs_good.py

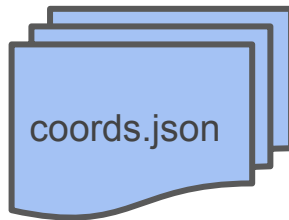
second_try_user_events.py

first_try_with_low_conf3.ipynb

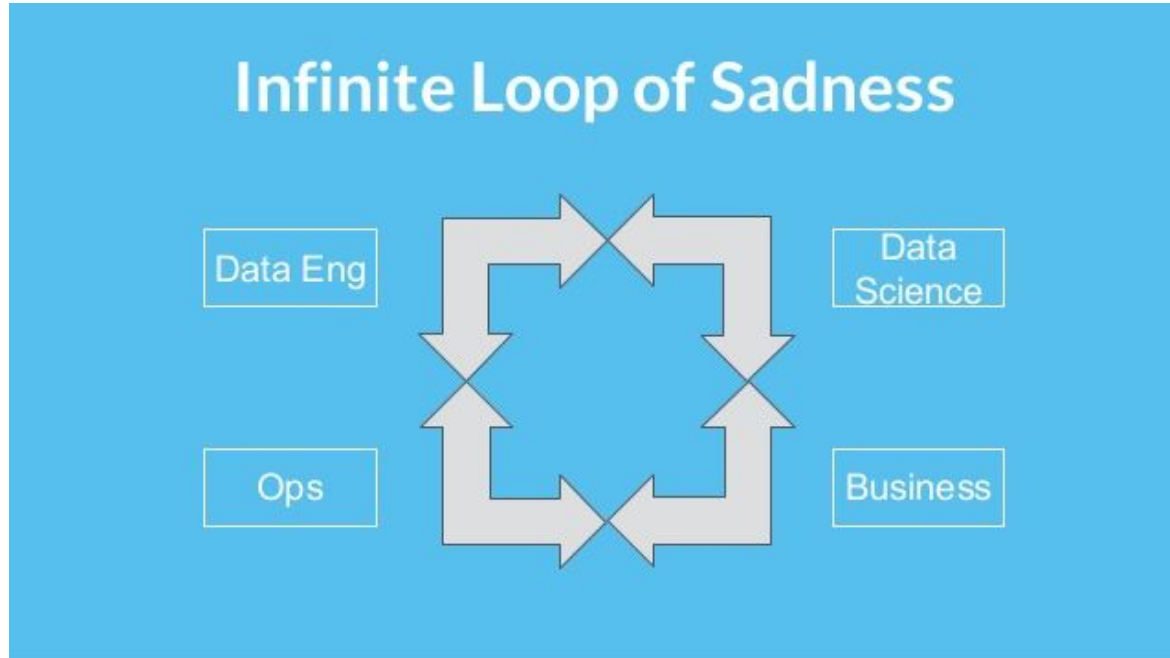
join_after_fix.json

opt_funcs_2016-11-12.m

scaling_patch.cpp



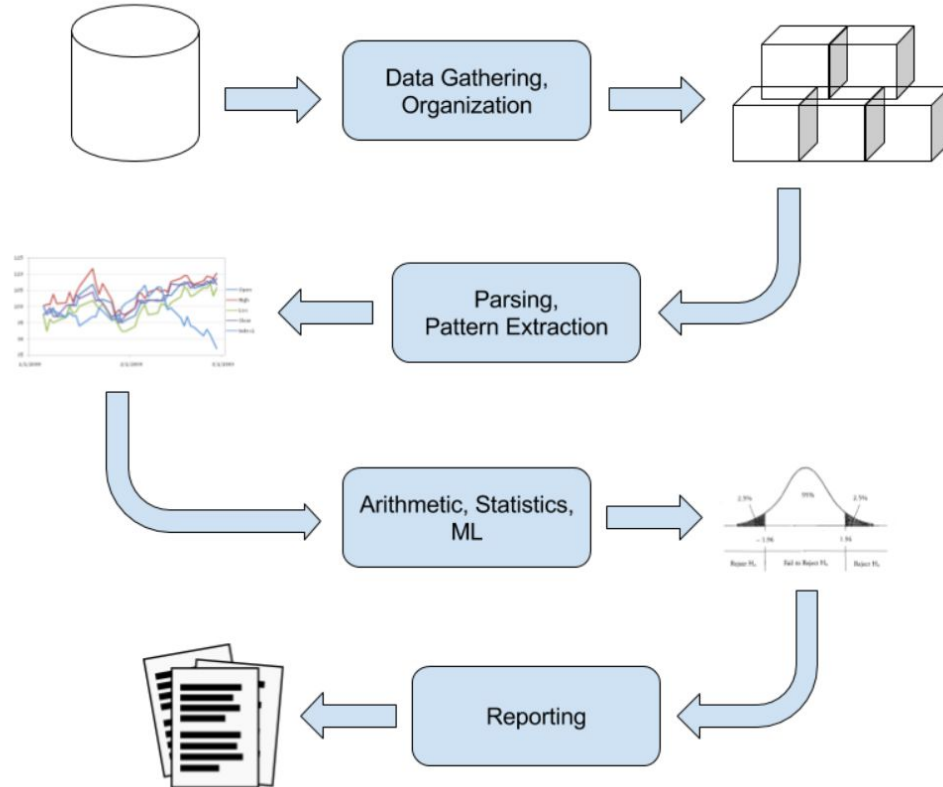
Workflow Management

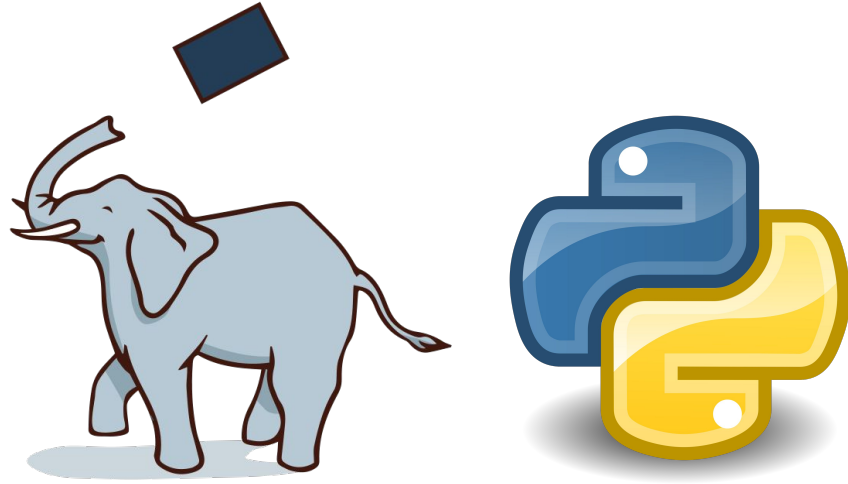


From [Josh Wills](#)

@dwhitena, @PyDataAnnArbor

Audit Trails, Debugging, Maintenance

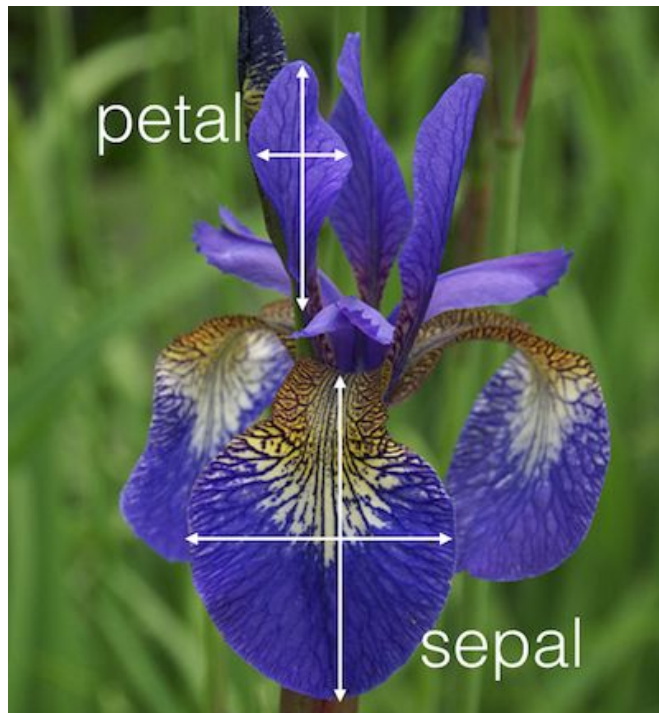




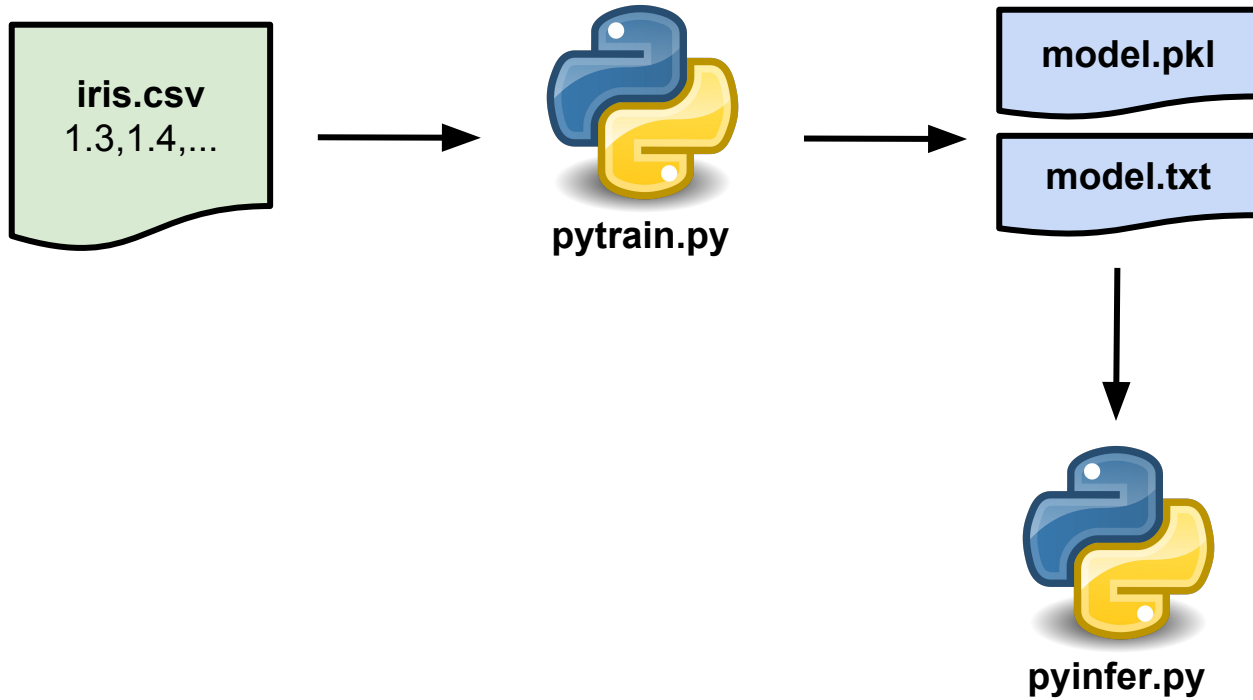
Demo Time

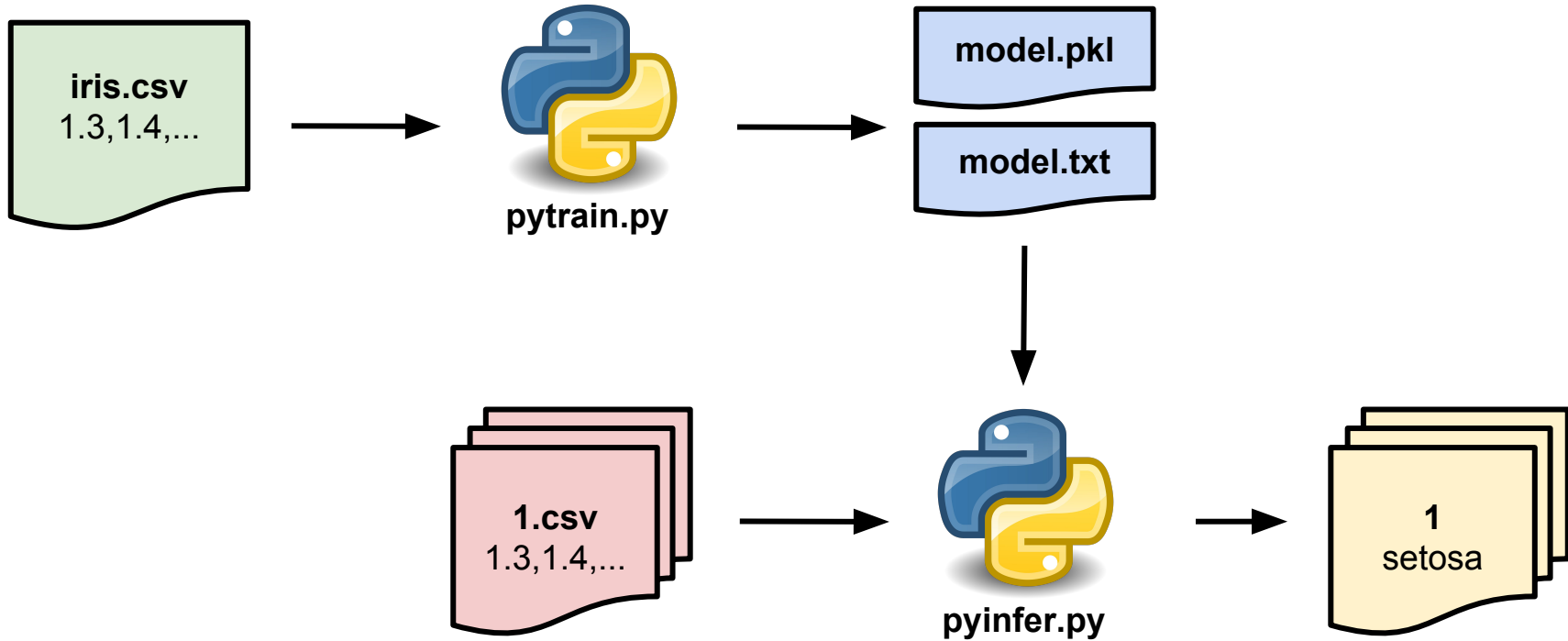


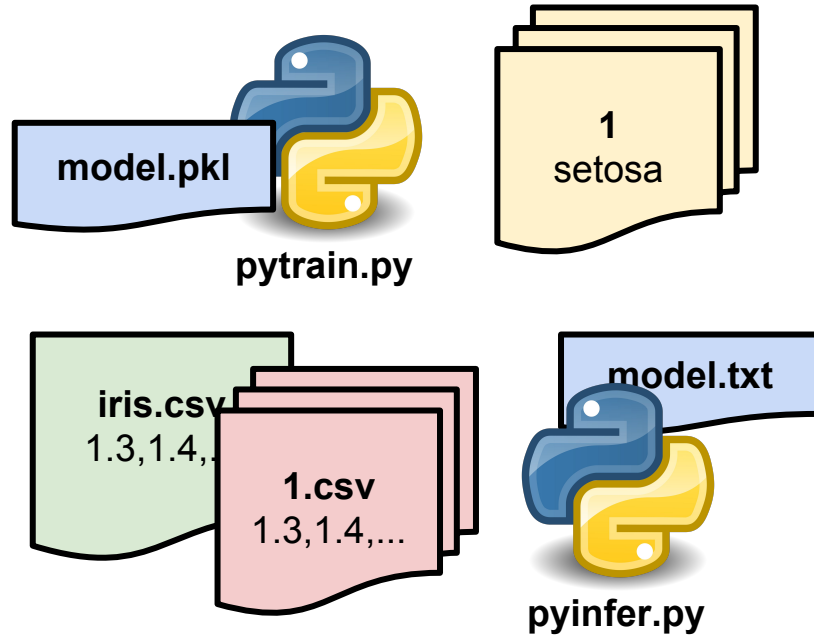
iris.csv
1.3,1.4,...

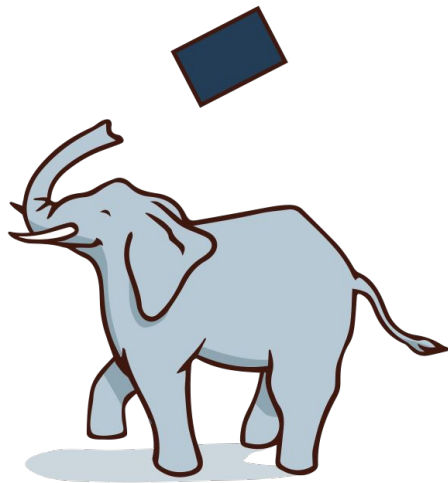








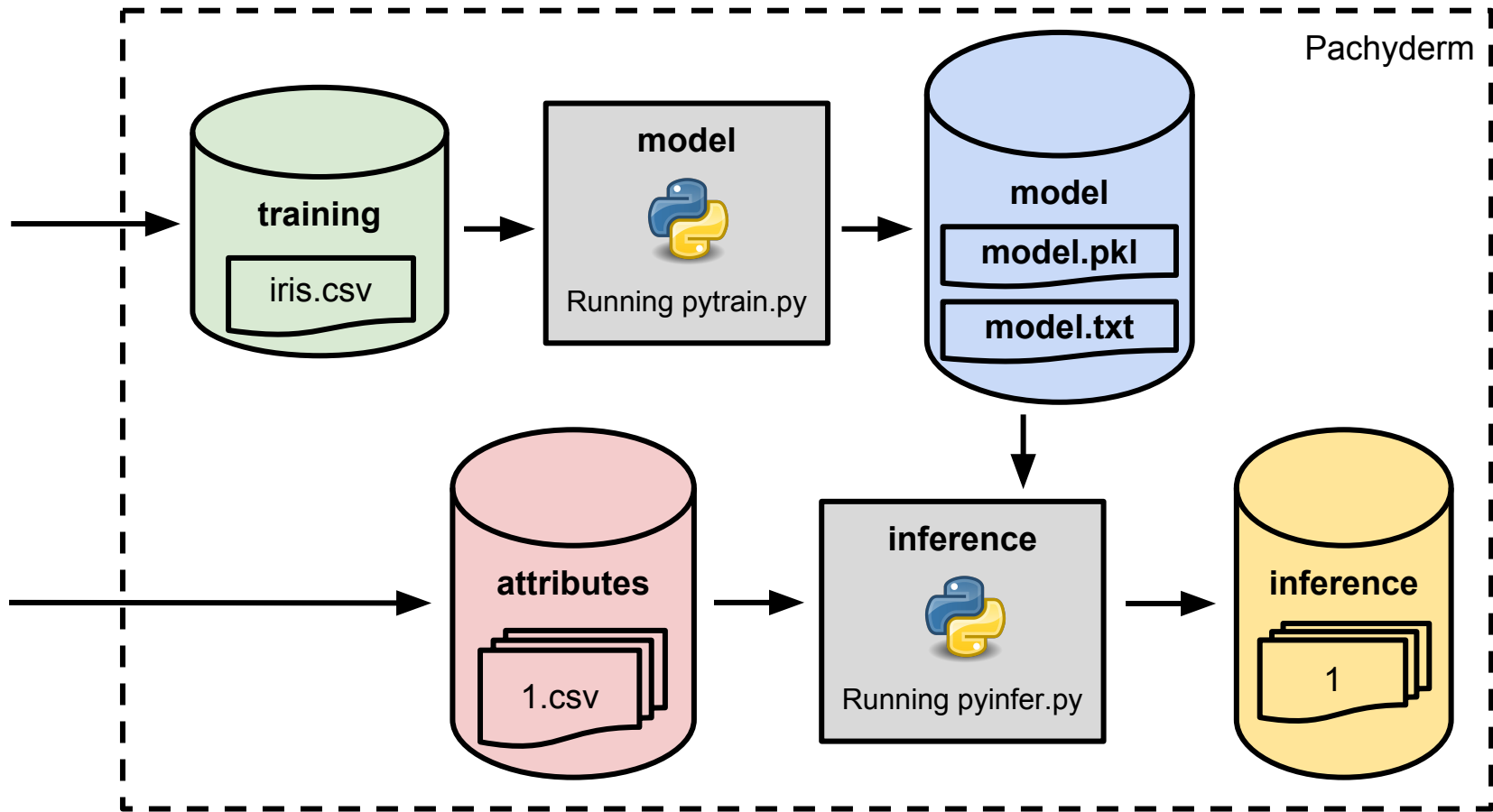


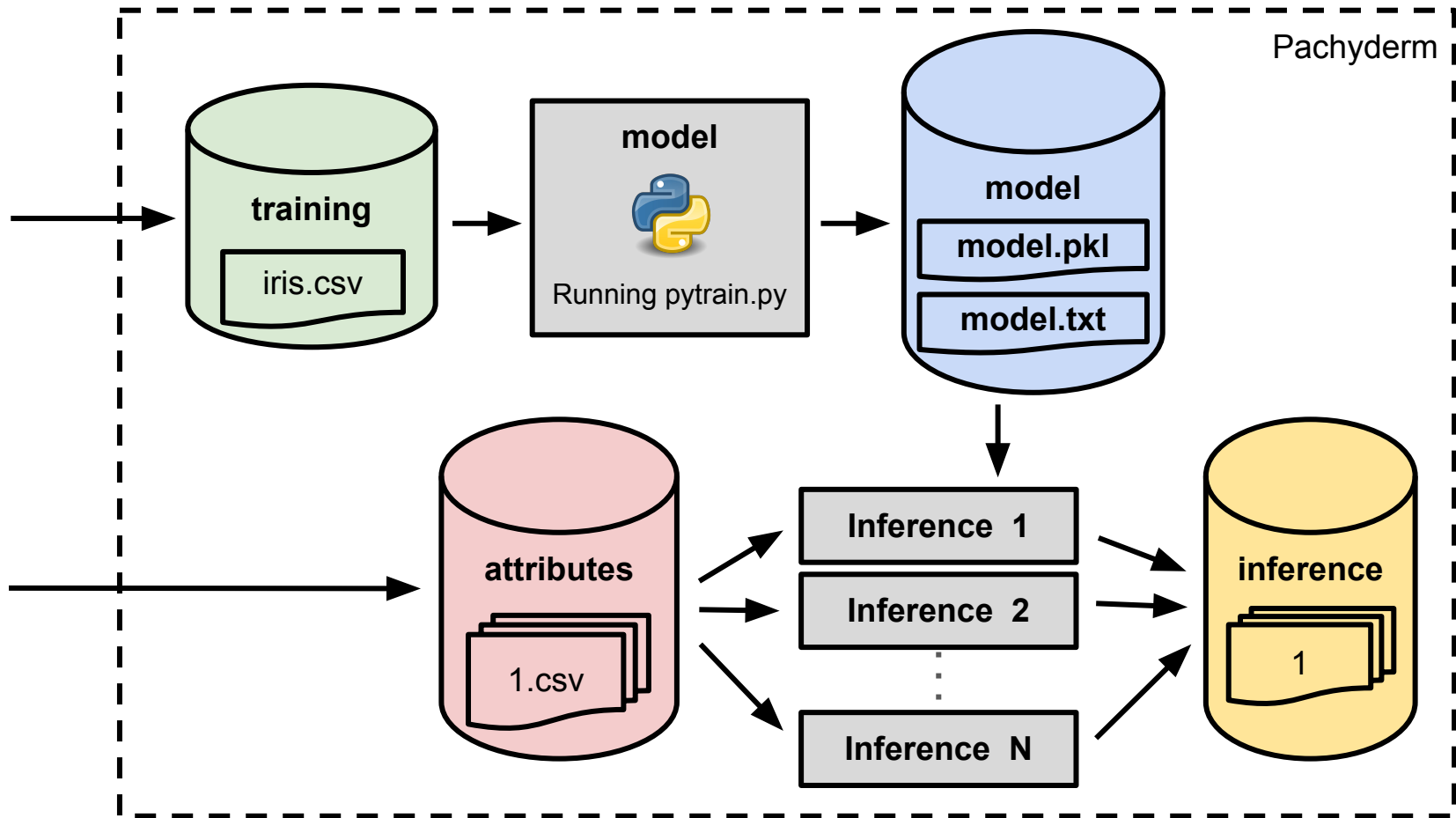


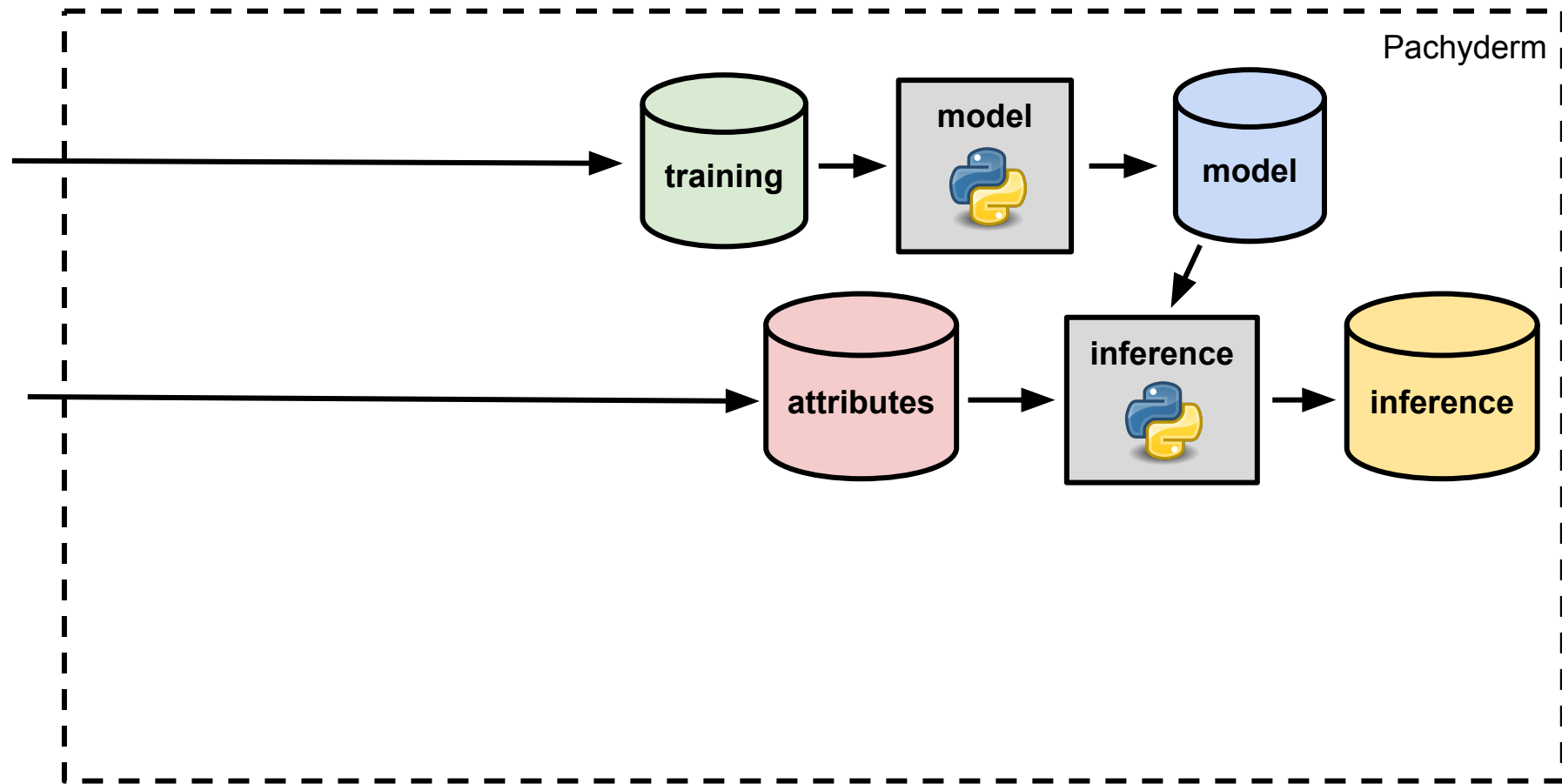
... enter Pachyderm

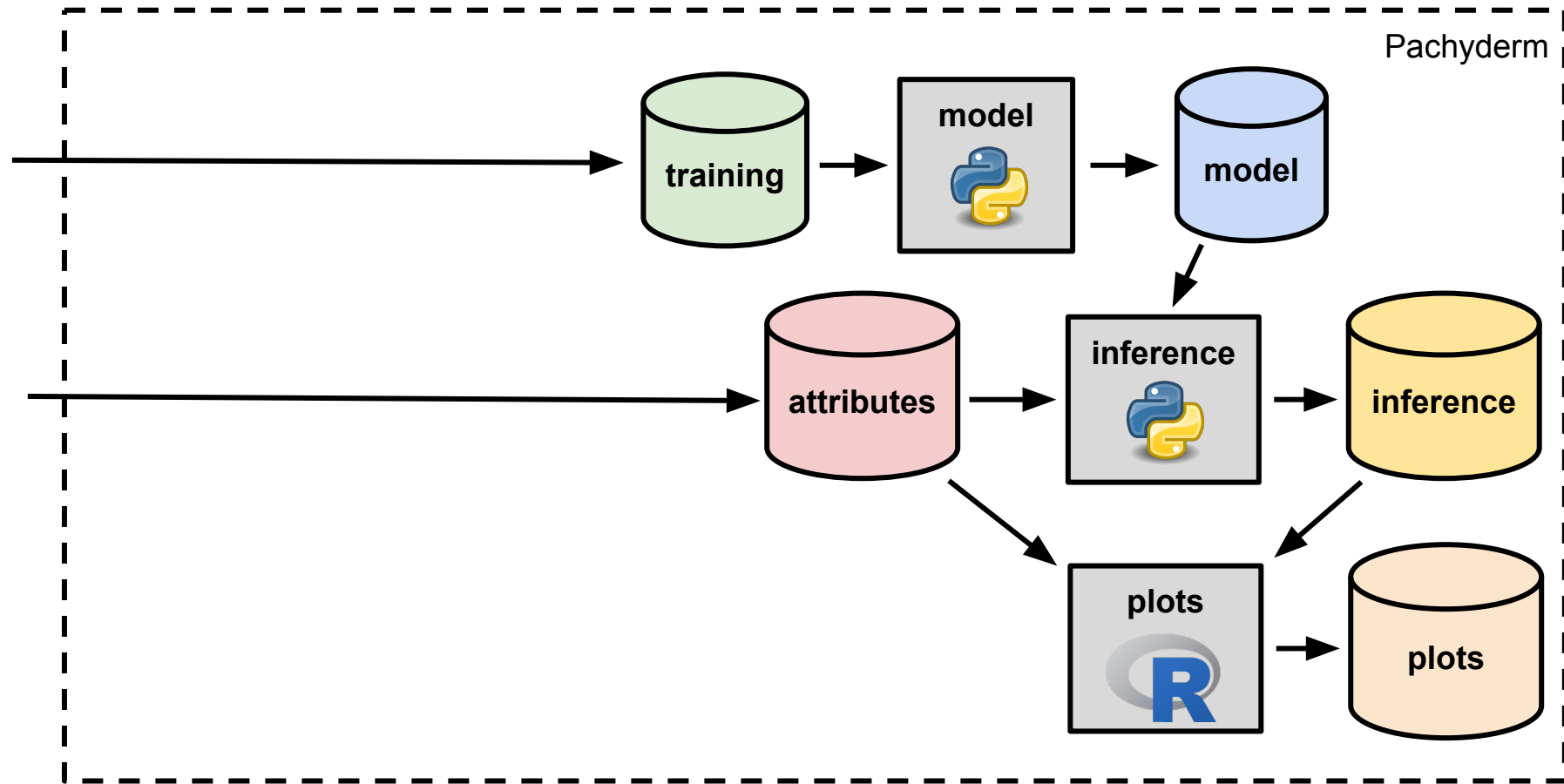
An open source, distributed processing and data versioning
framework built on containers.

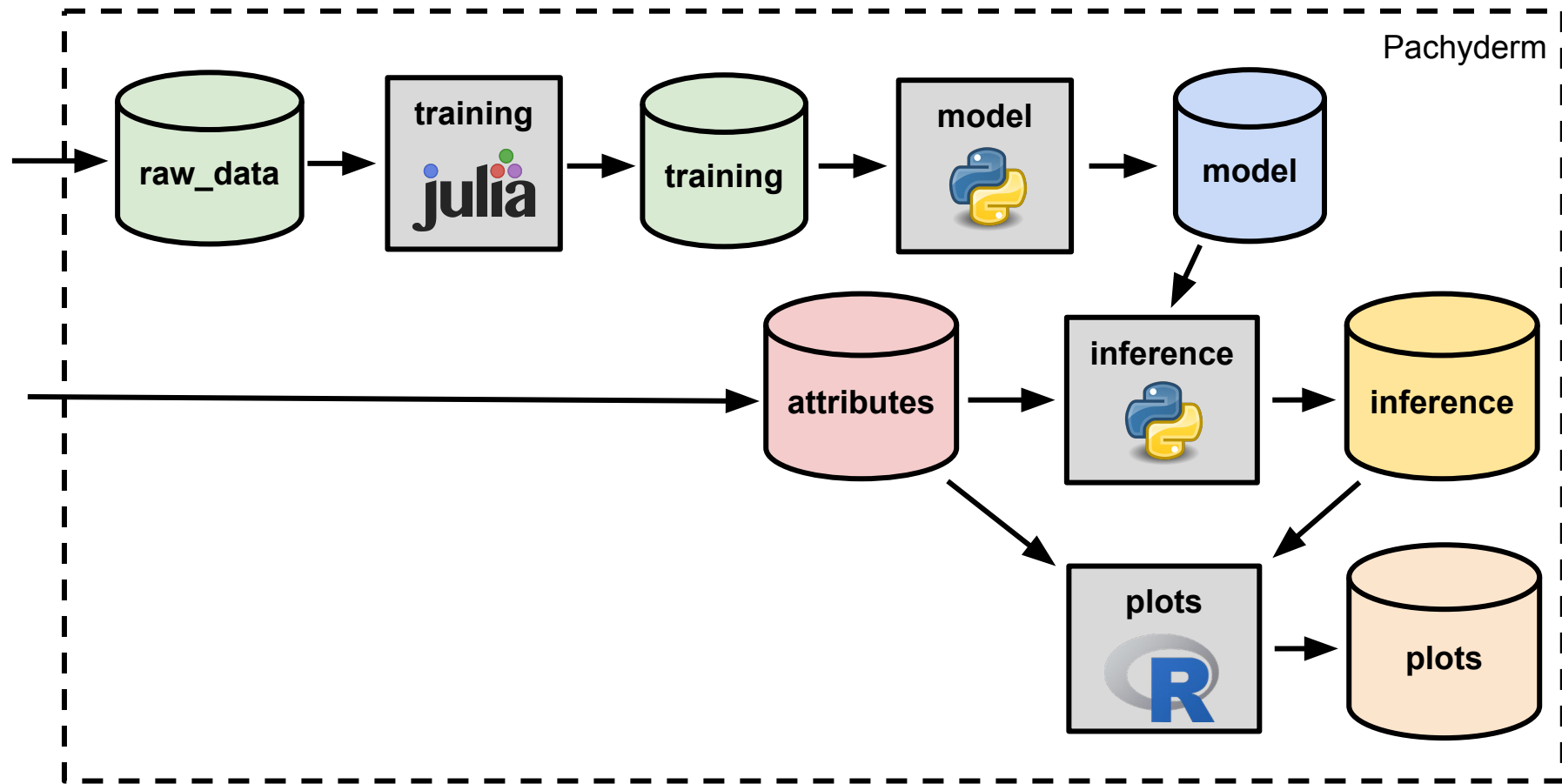
@dwhitena, @PyDataAnnArbor

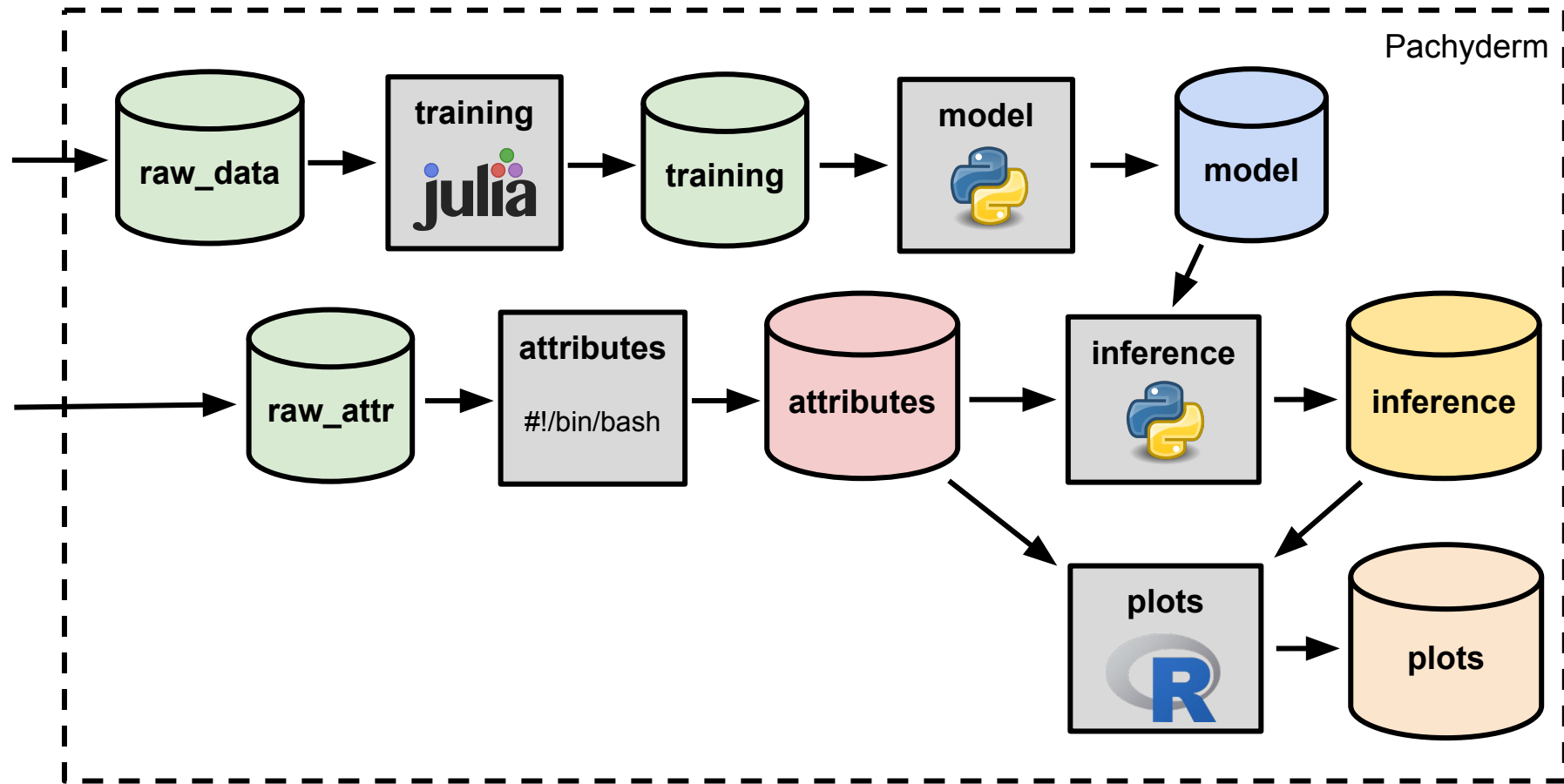


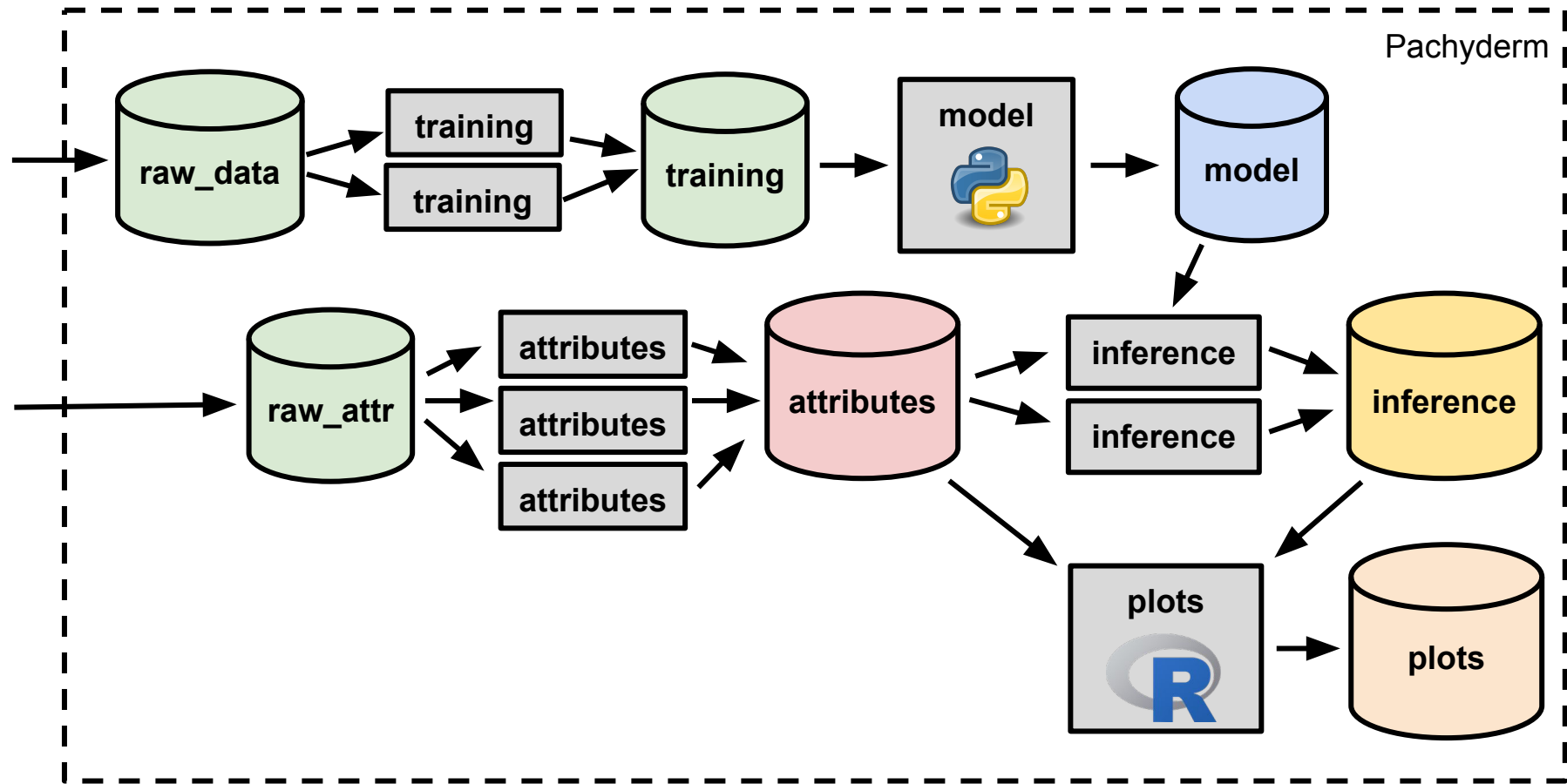












Conclusion/Resources

- Run [the code/pipeline](#)
- Run [other ML examples](#)
- Join the [Pachyderm Slack channel](#)
- Check out the [Pachyderm docs](#)
- Slack/tweet me [@dwhitena](#)
- Read a [related article](#)