

Popper
falsifiable.us

Practical Reproducible Evaluation of Computer Systems

Ivo Jimenez, Michael Sevilla, Noah Watkins,
Sina Hamedian, Pete Wilcox, Carlos Maltzahn,
Jay Lofstead, Kathryn Mohror, Adam Moody,
Andrea Arpaci-Dusseau, Remzi Arpaci-Dusseau

CROSS



BIG WEATHER WEB



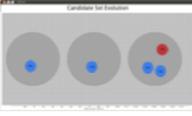
U.S. DEPARTMENT OF
ENERGY

Problem of Reproducibility in Computation and Data Exploration

The UI is capable of graphing the ratio $\frac{\text{optWIT}(OPT)}{\text{optWIT}(WFIT)}$ (Figure 2) in parallel with the analysis of the query stream (A high level of parallelism guarantees good performance). It also makes available the recommendations that are generated at each step, as well as the internal bookkeeping that the algorithm maintains. We will show some of this information as part of this scenario.

Scenario #2 We delve a little bit more into the details of our tool by allowing the candidate-index set to be automatically maintained but again disabling the feedback feature. This means that the candidate index set can automatically grow/shrink and be repartitioned over time based on the calculations of index interactions associated with each statement. This brings the tool into a completely online mode where it can operate autonomously without any user intervention.

Figure 3: Evolution of the candidate set with respect to partitioning (by calculating index interactions at each step). Each set corresponds to phases 1, 2 and 3 respectively.



We will see again how the algorithm generates a configuration such as shown in this figure, the partitioning of the candidate set will evolve for each of the three phases of the workload (Figure 3). We will show that this feature actually improves the quality of the recommendations.

Scenario #3 We complete the picture and show the effect that feedback has on the performance of WFIT by demonstrating one of the key contributions of our work: a principled feedback mechanism that is tightly integrated with the on-line algorithm (WFA).

By inspecting the recommended set of indices at any point in time, the DBA can decide whether to up- or down-vote any candidate index according to her criteria (or not vote at all). For the small test workload, it is easy to come up with reasonable “good” and “bad” votes that the algorithm can use to send as feedback to the physical design engine. We will execute three instances of WFIT concurrently with distinct feedback (good, bad, and no-feedback) and show the difference in performance for each (Figure 4).

The audience will see how, in the case of “good” feedback, the performance of WFIT increases in relation to the performance of “no-feedback” (the performance of OPT as baseline). In contrast, with “bad” feedback, the performance of WFIT will decrease; however, and more importantly, we will witness how WFIT is able to recover from poor feedback. This recovery mechanism is another important feature of the WFIT algorithm.

Scenario #4 The last scenario executes the *Reflex* workload suite of the *Online Index Selection Benchmark* [10] on Kansen. This is a complex workload consisting of approximately 1600 statements (queries and updates) that refer-

ence several datasets (TPC-C, TPC-DS, TPC-E, TPC-H and NREF).

We will show two WFIT variants: one with a stable and fixed candidate set partitioning; another whose candidate set is allowed to be automatically maintained. Similarly to scenario #1, we will graph the OPT vs. WFIT ratio in real-time as the workload is processed (Figure 5).

Figure 4: Multiple instances of WFIT running in parallel. The vote for the “good” and “bad” instances is done at step 1, causing the divergence in their behavior with respect to the “no-feedback” instance.

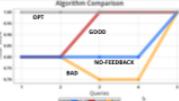
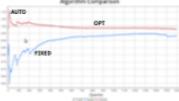


Figure 5: Two instances of WFIT running the Online Index Selection Benchmark. One with a fixed and stable candidate set (FIXED); another one with an automatically maintained candidate set (AUTO).



- What compiler was used?
- Which compilation flags?
- How was subsystem X configured?
- How does the workload look like?
- What parameters can be modified?
- What if I use input dataset Y?
- And if I run on platform Z?
- ...

5. REFERENCES
- [1] A. Bruno, S. Chaudhuri, L. Kollar, A. Marathe, V. Narasayya, and M. Syamsi. Database tuning advisor for microsoft SQL server 2005. In *SIGMOD conference*, pages 930–932, 2005.
 - [2] S. Chaudhuri and R. Krishnamurthy. An approach to physical design tuning: workload as a sequence. In *SIGMOD Conference*, pages 483–494, 2006.
 - [3] A. Bruno and S. Chaudhuri. Online computation and competitive analysis. Cambridge University Press, 1998.
 - [4] N. Bruno and S. Chaudhuri. An on-line approach to physical design tuning. In *VLDB Endowment*, pages 1–15, 2008.
 - [5] N. Bruno and S. Chaudhuri. Constrained physical design tuning. In *ICDE*, pages 1161–1164, 2009.
 - [6] D. Zillo et al. DB2 Design Advisor: Integrated Automatic Physical Optimization. In *VLDB*, pages 109–120, 2004.
 - [7] B. Dageville, D. Das, K. Dua, K. Vagnat, M. Zait, and M. Zeadman. Automatic SQL Tuning in Oracle 10g. In *VLDB*, pages 109–120, 2004.
 - [8] K. Schanatter, S. Abiteboul, T. Milo, and N. Polyzotis. On-line index selection: index shifting workloads. In *ICDE*, pages 459–468, 2009.
 - [9] K. Schanatter and N. Polyzotis. A Benchmark for Online Index Selection. In *ICDE*, pages 1761–1768, 2009.
 - [10] K. Schanatter and N. Polyzotis. On-line index tuning: a benchmark for Oracle 11g. In *PVLDB*, 5(5):478–489, 2012.
 - [11] K. Schanatter, N. Polyzotis, and L. Gotovac. Index interactions in physical design tuning: modeling, analysis, and applications. *PVLDB*, 2(1):1234–1245, 2009.

Lab Notebook

torpor-popper (branch: master)

| Subject | Author | Date |
|--|--------------------|----------------------|
| master | Ivo Jimenez | 2016-10-05... |
| Adds torpor as a submodule | Ivo Jimenez | 2016-08-25... |
| Results of running base-vs-targets for stressng on 4 machines | Ivo Jimenez | 2016-08-25... |
| Makes use of 'baseliner' role for... | Ivo Jimenez | 2016-08-25... |
| Removes unused line | Ivo Jimenez | 2016-08-21... |
| Will add later | Ivo Jimenez | 2016-08-21... |
| Re-adds baseliner ansible role | Ivo Jimenez | 2016-08-21... |
| renames ansible folder to exper... | Ivo Jimenez | 2016-08-21... |
| benchmarks now have their own... | Ivo Jimenez | 2016-08-21... |
| Adds latest version of baseliner... | Ivo Jimenez | 2016-06-12... |
| Adds experiment on all stress-... | Ivo Jimenez | 2016-06-12... |
| Adds same-platform variability... | Ivo Jimenez | 2016-06-10... |
| WIP on using baseliner role | Ivo Jimenez | 2016-06-10... |
| Checks first if docker is already... | Ivo Jimenez | 2016-05-16... |
| Tunning cpu-quota of 5 machines | Ivo Jimenez | 2016-05-16... |
| Fixes creation of parameters an... | Ivo Jimenez | 2016-05-16... |
| Adds logic to install statically li... | Ivo Jimenez | 2016-05-16... |

Gist it

SHA: af13570c1d700f85d2a9de348beb37215eb978c4
Author: Ivo Jimenez <ivo.jimenez@gmail.com>
Date: Thu Aug 25 2016 01:29:21 GMT-0700 (PDT)
Subject: Results of running base-vs-targets for stressng on 4 machines
Parent: [7a13185872bc73719048a1bd8a76f68a06798e13](#)

Results of running base-vs-targets for stressng on 4 machines

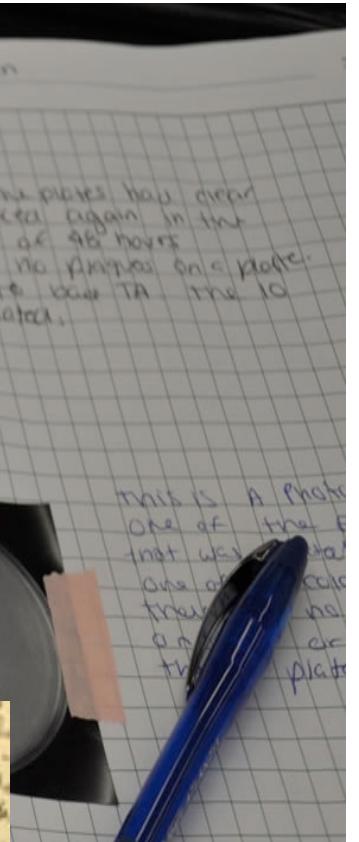
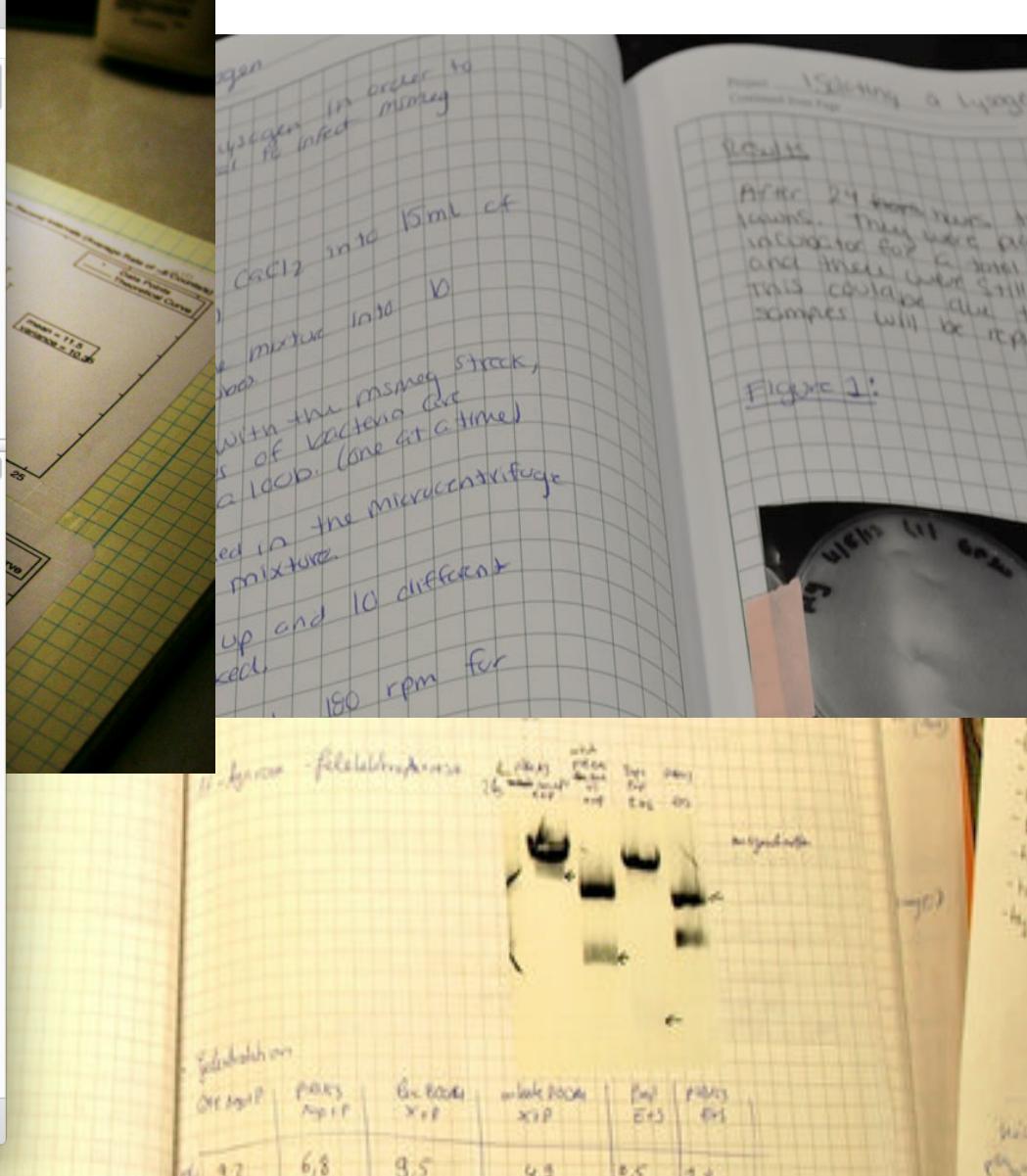
Base machine is 'issdm-12' and targets were tunned with the 'crafty' and 'c-ray' benchmarks.

The main difference between these results and the ones that appear on our VarSys '16 paper is that we are reflecting the speedup function for x=1, for both (1) tuning targets and (2) displaying results. This allows us to show better the reduction in variability without having to deal with different scales (slowdowns that lie between the [0-1] range are instead reflected and treated as speedups).

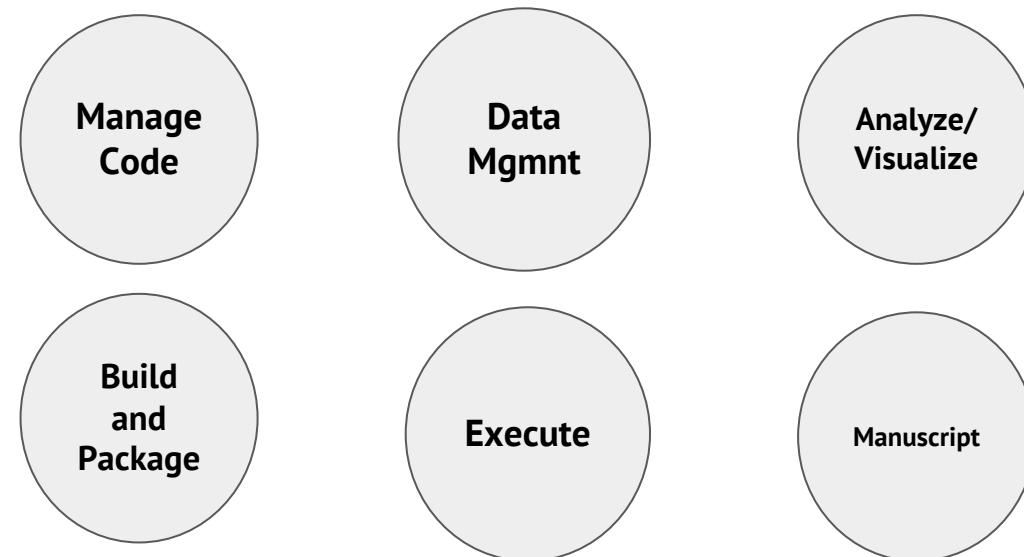
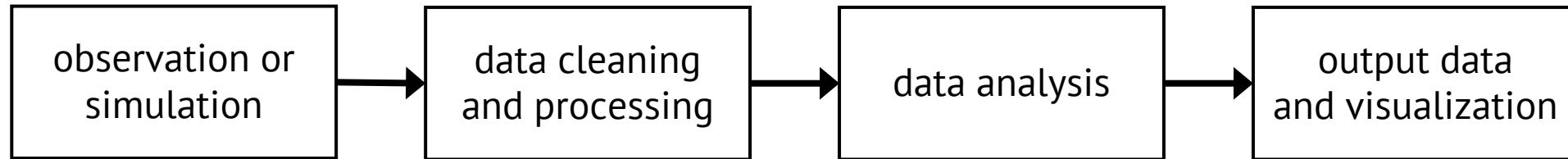
In short, we now unambiguously observe reduced variability when targets are limited. A couple of outliers, in particular stressng's memfd stressor, when limited, is very slow on target machines.

created [experiments/base-vs-limited-targets/all_results.csv](#)
created [experiments/base-vs-limited-targets/all_results.json](#)
created [experiments/base-vs-limited-targets/ansible.log](#)
created [experiments/base-vs-limited-targets/facts/192.168.140.81.json](#)

129 commits loaded



End-to-end Scientific Experimentation Pipelines



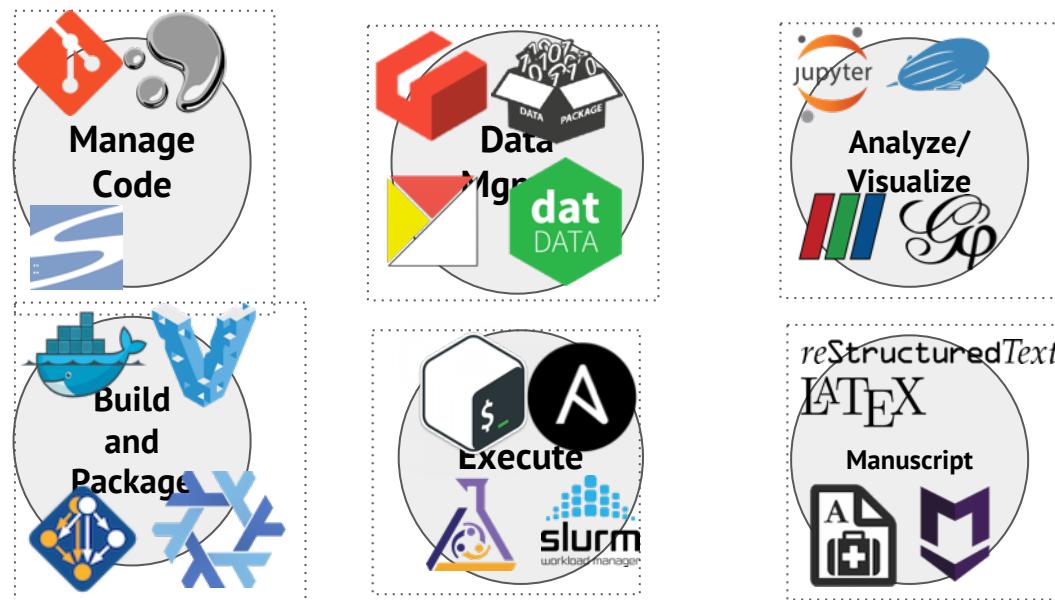
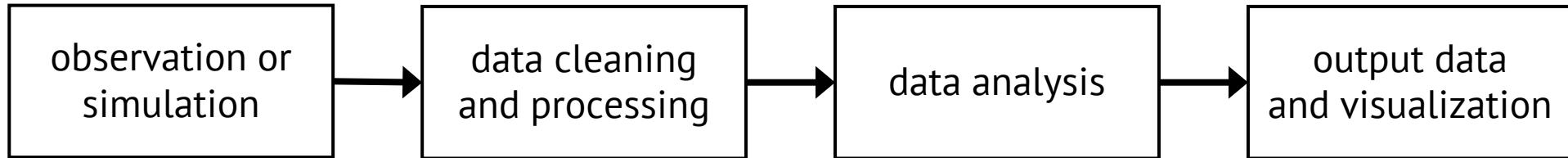
Analogies With Modern SE Practices (aka DevOps)

| Scientific exploration | Software project |
|--------------------------|-------------------------|
| Experiment code | Source code |
| Data management | Test examples |
| Analysis / visualization | Test analysis |
| Validation | CI / Regression testing |
| Manuscript / notebook | Documentation / reports |

Key Idea: manage a scientific exploration like software projects

SciOps

~~DevOps~~ View of The Experimentation Pipeline



What is DevOps?

Typical



```
ivo@mbp:~ ivo@mbp:~ ivo@mbp:~ ivo@mbp:~ [44/527]
$ hostname
mbp
ivo@mbp:~
$ date
Tue Jan 31 09:31:14 PST 2017
ivo@mbp:~
$ cat src/popper/popper/README.md
# Popper-CLI

A CLI tool to help bootstrap projects that follow the
[Popper](https://github.com/systemslab/popper) convention.

## Install

Download from
[popper/releases](https://github.com/systemslab/popper/releases). Note
that we have only tested on OSX and Linux (Windows coming soon). Once
downloaded, uncompress and place the binary in a folder that is
included in your '$PATH' (e.g. '/usr/bin').

## Usage

To get an overview and list of commands check out the command line
help:

```bash
popper help
```

ivo@mbp:~
$ find $HOME/tmp -path "*login*"
ivo@mbp:~
$ find $HOME/tmp -path "*git-meme*"
ivo@mbp:~
$ find $HOME/tmp -path "*git-meme.jpg*"
ivo@mbp:~
$ find $HOME/tmp -path "*meme.jpg*"

ivo@mbp:~
$ find $HOME/tmp/ -path "*meme.jpg*"
/Users/ivo/tmp//git-meme.jpg
ivo@mbp:~
$ clear
ivo@mbp:~
1:[tmux]* 2:bash- 3:bash
```

DevOps



myscript.sh

\$ bash myscript.sh

The Popper Convention

1. Pick one or more tools from the DevOps toolkit
2. Write scripts for an experiment pipeline
3. Put all scripts in a version control repository



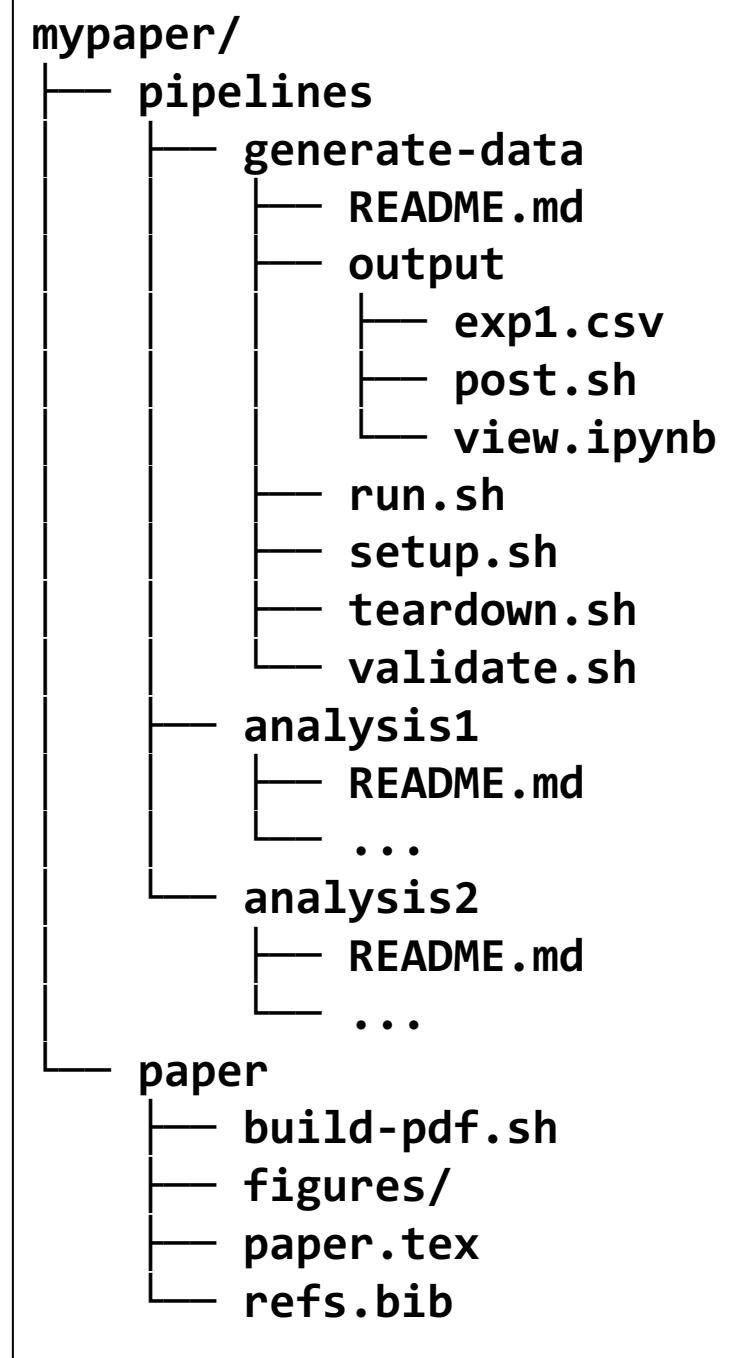
[1]: Jimenez et al. *Standing on the Shoulders of Giants by Managing Scientific Experiments Like Software*, ;login: Winter 2016, Vol. 41, No. 4.
[2]: Jimenez et al. *The Popper Convention: Making Reproducible Systems Evaluation Practical*, REPPAR 2017.

Popper CLI tool



- Make it **super easy** to automate execution and validation of experimentation pipelines
 - easy → low-overhead → more likely it'll be used
- Common convention to organize the contents of a repo
- CLI tool that helps users to implement pipeline stages
- Provide domain-specific examples
 - Today: Genomics, MPI, Ceph, Atmospheric Science
 - Working with domain-experts to contribute more examples

Common convention to organize the contents of a repo



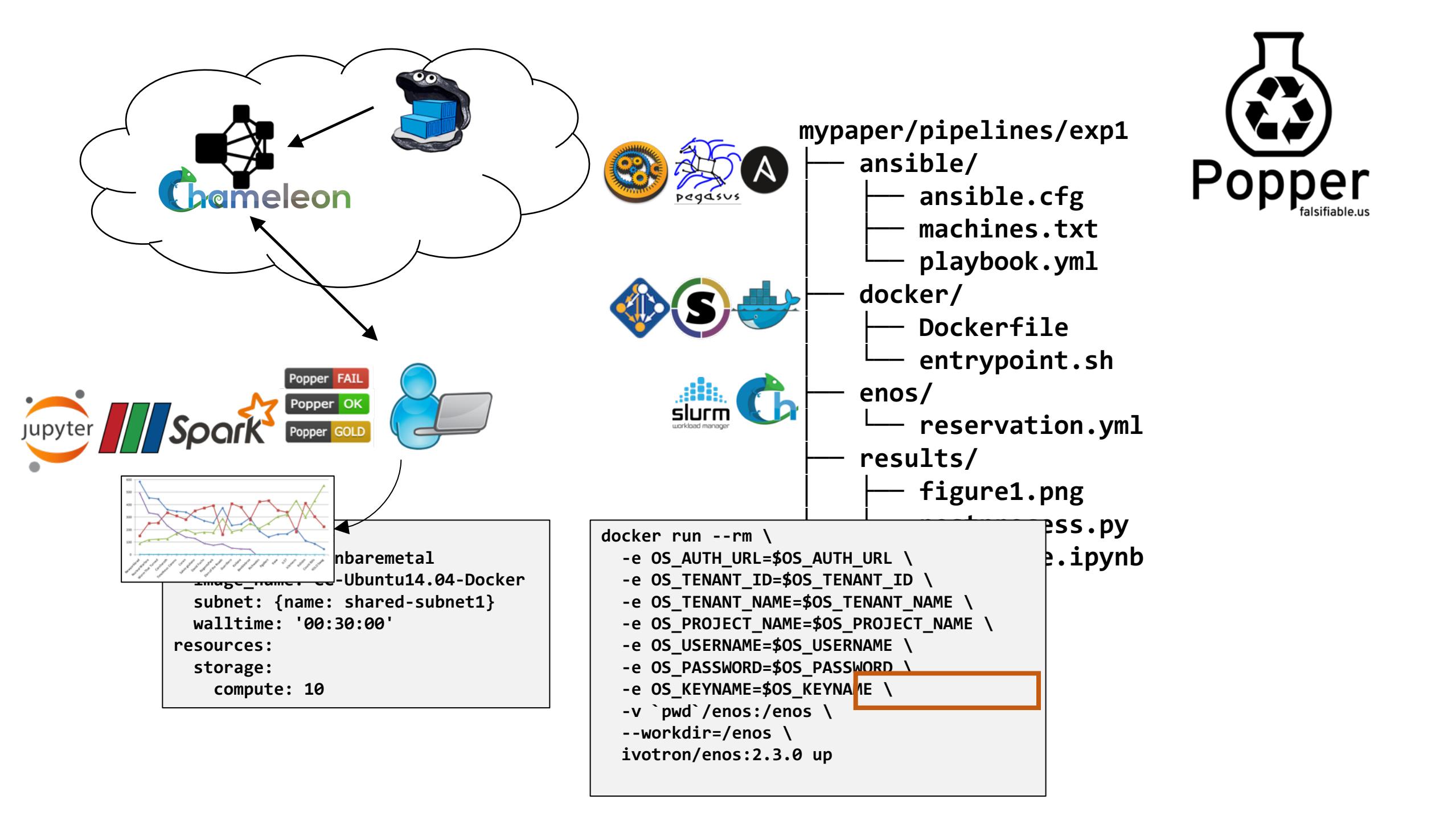
CLI tool

```
$ cd my-paper-repo
$ git init
Initialized empty Git repository in my-paper-repo/.git
$ popper init
Initialized popper repository.
$ popper pipeline init mypipeline --stages=prepare,execute,cleanup
-- Initialized exp1 pipeline.
```

```
$ ls -l pipelines/mypipeline
total 20K
-rw-r----- 1 ivo ivo 8 Apr 29 23:58 README.md
-rwxr-x--- 1 ivo ivo 210 Apr 29 23:58 execute.sh
-rwxr-x--- 1 ivo ivo 206 Apr 29 23:58 prepare.sh
-rwxr-x--- 1 ivo ivo 61 Apr 29 23:58 cleanup.sh
```

```
#!/bin/bash
#!/bin/bash

# trigger execution of experiment
docker run google/kubectl run ...
...
```



```
$ popper run exp1
```

Popper run started

Stage: setup.sh

Stage: run.sh

Stage: validate.sh .

Stage: teardown.sh ...

Popper run finished

Status: OK



Codified Validations

num_nodes, throughput, raw_bw, net_saturated

```
Src,Egid,Version,Datetime,Lat,Lon,Magnitude,Depth,NST,Region  
ci,14692356,1,"Tuesday, May 4, 2010 03:21:38 UTC",32.6443,-11  
ci,14692348,1,"Tuesday, May 4, 2010 03:19:38 UTC",32.1998,-11  
ci,14692332,1,"Tuesday, May 4, 2010 03:16:56 UTC",32.6756,-11  
ci,14692324,1,"Tuesday, May 4, 2010 03:08:47 UTC",32.6763,-11  
ci,14692316,1,"Tuesday, May 4, 2010 03:08:08 UTC",32.6778,-11  
ci,14692308,1,"Tuesday, May 4, 2010 03:06:20 UTC",32.7071,-11  
ci,14692300,1,"Tuesday, May 4, 2010 03:01:52 UTC",32.1948,-11  
ak,10047267,1,"Tuesday, May 4, 2010 03:01:04 UTC",61.2695,-14  
ci,14692284,1,"Tuesday, May 4, 2010 02:58:51 UTC",32.7016,-11  
ci,14692276,1,"Tuesday, May 4, 2010 02:57:46 UTC",32.6998,-11  
ak,10047263,1,"Tuesday, May 4, 2010 02:56:28 UTC",63.5779,-11  
ak,10047261,1,"Tuesday, May 4, 2010 02:52:00 UTC",66.4986,-14  
ci,14692268,1,"Tuesday, May 4, 2010 02:48:40 UTC",32.6813,-11  
ci,14692260,1,"Tuesday, May 4, 2010 02:35:27 UTC",32.2006,-11  
nc,71392116,0,"Tuesday, May 4, 2010 02:15:24 UTC",38.8415,-12  
ci,14692244,1,"Tuesday, May 4, 2010 02:05:07 UTC",33.5248,-11  
ci,14692228,1,"Tuesday, May 4, 2010 01:57:08 UTC",32.6823,-11  
ci,14692220,1,"Tuesday, May 4, 2010 01:53:28 UTC",32.6881,-11  
ci,14692212,1,"Tuesday, May 4, 2010 01:48:53 UTC",32.6398,-11  
ci,14692188,1,"Tuesday, May 4, 2010 01:26:58 UTC",32.5003,-11  
ci,14692180,1,"Tuesday, May 4, 2010 01:19:44 UTC",32.6836,-11  
ci,14692172,1,"Tuesday, May 4, 2010 01:12:01 UTC",32.5321,-11  
ci,14692164,1,"Tuesday, May 4, 2010 01:08:24 UTC",32.6833,-11
```

```
Stage: run.sh .....  
Stage: validate.sh ....
```

```
[true] check linear scalability  
[true] check system throughput
```

Popper run finished

Status: GOLD



```
expect  
linear(num_nodes, throughput)
```

```
when  
not net_saturated  
expect  
throughput >= (raw_bw * 0.9)
```

[1]: Jimenez et al. *Tackling the reproducibility problem in storage systems research with declarative experiment specifications*, PDSW '15.

[2]: Jimenez et al. *I Aver: Providing Declarative Experiment Specifications Facilitates the Evaluation of Computer Systems Research*, TinyTOCS, Vol. 3,

Archiving/DOI service integration

```
$ popper archive --zenodo --user=ivotron --password=***  
Creating archive for repository on Zenodo.  
| #####| 100 %
```

Your DOI link is: <https://zenodo.org/record/1165550>

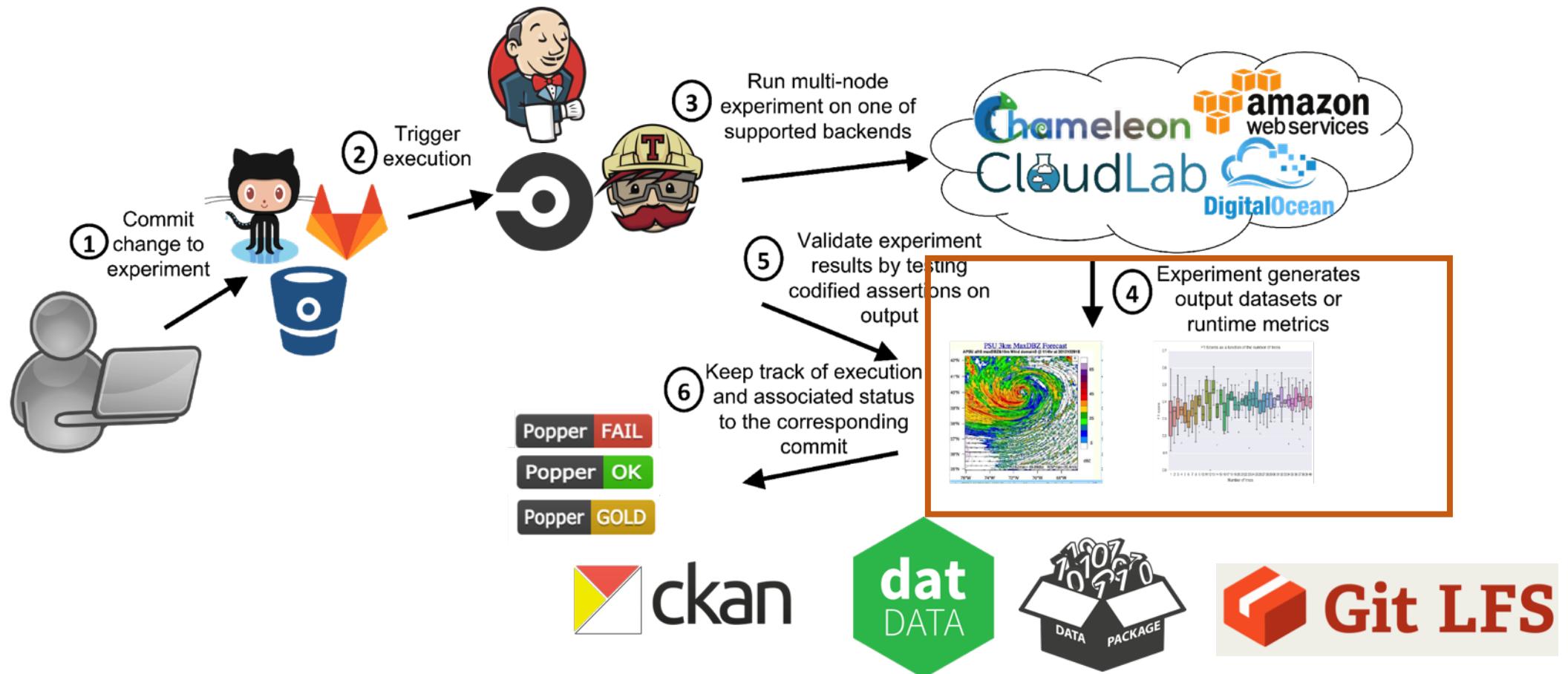


ACM/Popper Badges



| Result Status | Artifacts | Re-executed By | ACM | Popper |
|----------------------|------------------|-----------------------|---|--------------------|
| Repeatability | Original | Original Author(s) | | Popper GOLD |
| Replicability | Original | 3 rd Party |  | Popper GOLD |
| Reproducibility | Re-implemented | Anyone |  | |

Popper and CI



Push-button Reproducible *Evaluation*

Introduction to Popper

Secure | https://ivotron.github.io/popper-lesson/

Home Code of Conduct Setup Episodes Extras License Improve this page

Search...

Introduction to Popper

Popper is a convention and CLI tool for conducting experiments and writing academic articles following best OSS and DevOps practices. The open-source CLI tool helps researchers organize all the artifacts associated to a scientific exploration inside a code repository such as Git. These “popperized” repositories are self-contained and can be easily shared (cloned via Git) so that others can replicate the original results (see examples [here](#)).

★ Prerequisites

In this lesson we use the [Unix Shell](#). Some previous experience with using the shell to list directories, create, copy, remove and list files and directories, and run simple scripts is necessary. Some familiarity with [Git](#) makes it easier to follow the episodes but is not required.

★ Setup

In order to follow this lesson, you will need to install software in your computer. Please follow instructions on the [setup](#) page.

Schedule

| | Setup | Download files required for the lesson |
|-------|---|---|
| 00:00 | 1. Introduction | How can I generalize the structure of my scientific explorations? |
| 00:00 | 2. Popper Pipelines | How can I automate a scientific exploration pipeline? |
| 00:00 | 3. Visualizing Popper Pipelines | How can I quickly get a sense of what a pipeline does? |
| 00:00 | 4. Continuous Validation of Pipelines | How can I check the integrity of a pipeline over time? |
| 00:00 | 5. Portable Pipelines With Docker | How can I easily run a pipeline on other machines? |
| 00:00 | 6. Conclusion | What are the advantages and disadvantages of using Popper? |
| 00:00 | Finish | |

The actual schedule may vary slightly depending on the topics and exercises chosen by the instructor.