

# BlackSwan: The Case for a Practical Reproducibility Platform

RESCUE-HPC'18 (SC'18)

Ivo Jimenez and Carlos Maltzahn (*UC Santa Cruz*)

[ivo.jimenez|carlosm]@ucsc.edu

twitter: @ivotron

# Problem of Practical Reproducibility in Computer, Computational and Data Science

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.

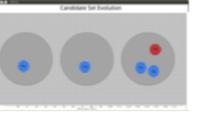


Figure 4: Multiple instances of WFIT running in parallel. The vote for the "good" and "bad" instances is done at step 1, causing the divergence in their behavior with respect to the "no-feedback" instance.

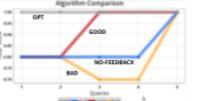
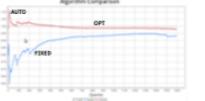


Figure 5: Two instances of WFIT running the Online Index Selection Benchmark. One with a fixed and stable candidate set (FIXED); another with an automatically maintained candidate set (AUTO).



The UI is capable of graphing the ratio  $\frac{\text{outW} - \text{O}(\text{OPT})}{\text{outW} - \text{O}(\text{WFIT})}$  (Figure 2) in parallel with the analysis of the query stream (A highly visualized interface generates good reports and publications.) 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 bypassing 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.

We will see again how the algorithm generates a configuration such as, however, at this point, 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 logic of the online algorithm (WFIT).

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" (we use 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-

[1] S. Schanitter, N. Polyzotis, L. Kollar, A. Marathe, V. Narasayya, and M. Syamas. Database tuning advisor for microsoft SQL server 2005. In *SIGMOD conference*, pages 930–932, 2005.  
[2] S. Schanitter, N. Polyzotis, and L. Gotsova. Online index selection design tuning: workload as a sequence. In *SIGMOD Conference*, pages 483–494, 2006.  
[3] A. Chaudhuri and R. Krishnamurthy. Online computation and competitive analysis. Cambridge University Press, 1998.  
[4] N. Bruno and S. Chaudhuri. An online approach to physical design tuning. *VLDB Endowment*, 1(4–15), 2008.  
[5] N. Bruno and S. Chaudhuri. Online automatic physical design tuning. In *ICDE*, pages 1161–1164, 2009.  
[6] D. Zillo et al. DB2 Design Advisor: Integrated Automatic Physical Design Tuning. In *VLDB*, pages 109–120, 2004.  
[7] B. Dageville, D. Das, K. Dasgupta, K. Vagnari, M. Zett, and M. Ziauddin. Automatic SQL Tuning in Oracle 10g. In *VLDB*, pages 109–120, 2004.  
[8] K. Schanitter, S. Abitbol, T. Milo, and N. Polyzotis. On-line index selection: shifting workloads. In *ICDE*, pages 459–468, 2009.  
[9] K. Schanitter and N. Polyzotis. A Benchmark for Online Index Selection. In *ICDE*, pages 1761–1768, 2009.  
[10] K. Schanitter and N. Polyzotis. Online index tuning: a benchmark. In *VLDB*, 34(5):478–489, 2012.  
[11] K. Schanitter, N. Polyzotis, and L. Gotsova. Index interactions in physical design tuning: modeling, analysis, and applications. *VLDB*, 2(1):1234–1245, 2009.

## Replicate:

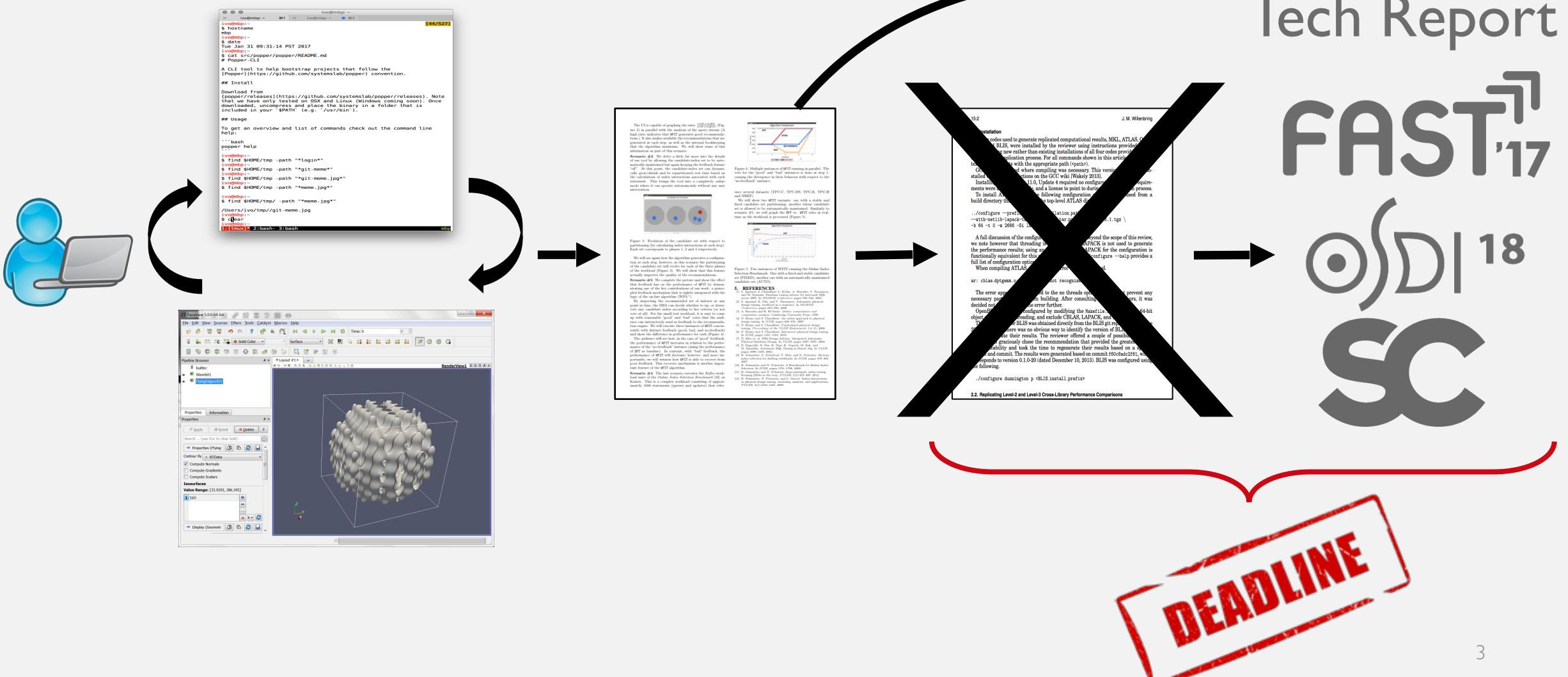
- What compiler was used?
- Which compilation flags?
- How was subsystem X configured?

## Reuse:

- What parameters can be modified?
- What if I use input dataset Y?
- And if I run on platform Z?

2

# How is a Manuscript Prepared?



# A Deadline-friendly Approach

## Typical



```
ivo@mbp: ~ 361 ivo@mbp: ~ 362 [44/527]
ivo@mbp: ~
$ 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:[(muX)* 2:bash- 3:bash
```

## DevOps



myscript.sh

```
$ bash myscript.sh
```

- Automation
  - Bash, Python, etc.
- Version-control
  - Git, Mercurial, etc.
- Portability
  - Containers, VMs, Spack, etc.

# The Popper Experimentation Protocol

1. Pick one or more tools from the DevOps toolkit.
2. Write *portable* scripts for a 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.

GitHub, Inc. popperized / popper-readthedocs-examples Watch 3 Star 0 Fork 2 Insights Branch: master popper-readthedocs-examples / pipelines / blis / Franko1307 and ivotron doc: add validations for BLIS pipeline #408 (#21) ... Latest commit d498146

**build-image**

```

14
15/15 2.1. Installation J.M. Willenbring
ivo@scruffy: /home/ivo

```

**ivo@scruffy:~**  
\$ cd tmp/mypaper/  
**ivo@scruffy:~/tmp/mypaper**  
\$ popper add popperized/popper-readthedocs-examples/blis  
Downloading pipeline blis as blis...  
Updating popper configuration...  
**Pipeline blis has been added successfully.**  
**ivo@scruffy:~/tmp/mypaper**  
\$ ls -l pipelines/blis/  
total 40K  
-rwxrwxr-x 1 ivo ivo 232 Aug 7 21:28 **analyze.sh**  
-rwxrwxr-x 1 ivo ivo 120 Aug 7 21:28 **build-docker-image.sh**  
drwxrwxr-x 2 ivo ivo 4.0K Aug 7 21:28 **docker/**  
drwxrwxr-x 2 ivo ivo 4.0K Aug 7 21:28 **docker-blis/**  
drwxrwxr-x 2 ivo ivo 4.0K Aug 7 21:28 **docker-octave/**  
-rw-rw-r-- 1 ivo ivo 1.5K Aug 7 21:28 README.md  
drwxrwxr-x 3 ivo ivo 4.0K Aug 7 21:28 **results/**  
-rwxrwxr-x 1 ivo ivo 93 Aug 7 21:28 **run.sh**  
-rwxrwxr-x 1 ivo ivo 4.6K Aug 7 21:28 **validate.sh**  
**ivo@scruffy:~/tmp/mypaper**  
\$ popper run blis  
Executing pipeline: blis  
[#####----] run

**1:bash 2:bash 3:bash- 4:python\* 5:sudo** scruffy

# Popper Compliant Pipelines

On every execution:

1. Checkout code and data dependencies.
2. Build, install and deploy system under test.
3. Allocate resources (statically or dynamically).
4. Expose pipeline parameters in a plain text file.
5. Capture environment info (hardware, software stacks).
6. Validate results.



Self-verifiable  
Experimentation Pipeline

# Community Building, Outreach and Teaching



The collage includes the following elements:

- U.S. DEPARTMENT OF ENERGY ENERGY**: Logo of the U.S. Department of Energy.
- Google Summer of Code**: Logo for Google's Summer of Code program.
- EXASCALE COMPUTING PROJECT (ECP)**: Logo for the Exascale Computing Project.
- SPASE**: Logo for the SPASE library.
- BETTER SCIENTIFIC SOFTWARE**: Logo for Better Scientific Software International.
- GITTER**: A screenshot of a Gitter messaging interface showing a conversation about the Popper experimentation protocol.
- BIDS**: Logo for the Berkeley Institute for Data Science (BIDS).
- UNIVERSITY OF CALIFORNIA MERGED**: Logos for the University of California, Merced.
- UCDAVIS**: Logo for the University of California, Davis.
- Los Alamos NATIONAL LABORATORY**: Logo for Los Alamos National Laboratory.
- Science Lab OPENCON 2018**: Logos for Mozilla Science Lab and OPENCON 2018.
- THE CARPENTRIES**: Logo for The Carpentries.
- URSI**: Logo for the Union of Radio Instrumentation Societies of India (URSI).
- GitHub Gitter screenshot**: A screenshot of a GitHub Gitter channel for the Popper package, showing a discussion thread.
- LinkedIn screenshot**: A screenshot of a LinkedIn group page for the American Meteorological Society (AMS), titled "American Meteorological Society | 98th Annual Meeting | AUSTIN | 2018".

**Key points from the slide:**

- \$\$\$ for teaching Popper
  - Organize one-day hands-on tutorials
  - Attend hackathons.
- Engage with OpenScience community.



**Popper**  
falsifiable.us

- Improved experimentation practices.
- Increased confidence on results.
- Efficient use of time, both at personal and collaboration levels.



- Steep learning curve; radical paradigm shift.
- Submission deadline pressure is too big.
- Lack of incentives for taking the leap; no group-wide adoption.



# BlackSwan: Practical Reproducibility Platform

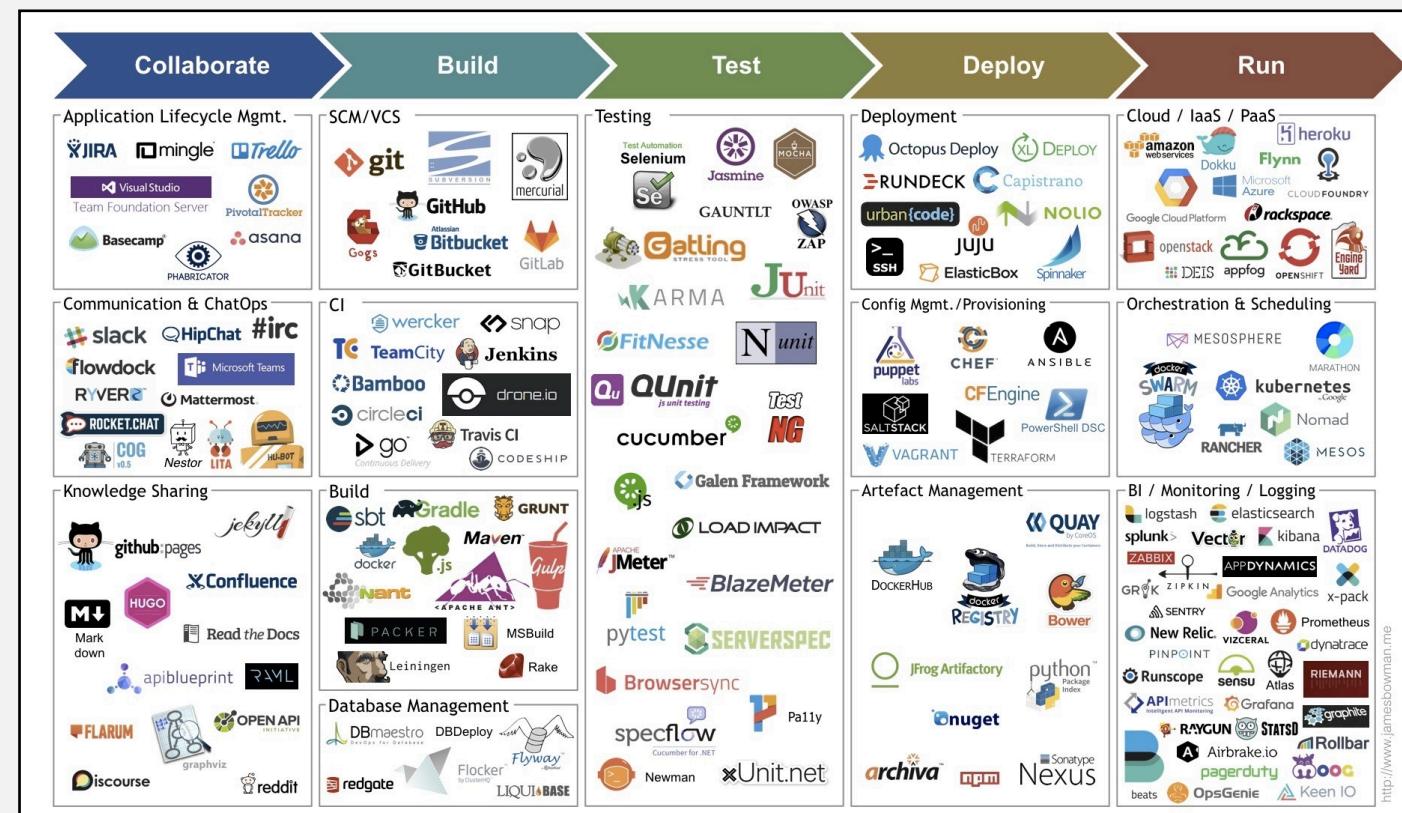


open source

## Community



## DevOps Technology





# Pipeline as the Central Concept

## Search:

spark

## Results:

spark-craill



**spark-tensorflow**



**Description:** This pipeline benchmarks tensorflow embedded in Spark workers.

**Runs on:** CloudLab

Status: blackswan passing

Import

spark-tpch



## Pipeline:



**Stage:** cluster-setup

**Description:** This stage allocates a set of nodes on CloudLab and writes a YAML file with metadata about each of the nodes in the cluster.

**Required variables:**

- CLOUDLAB\_USER\_NAME
- CLOUDLAB\_PASSWORD
- CLOUDLAB\_CERT\_PATH

**Dependencies**

- Docker +17.03
- Python 2.7

Import

# Automated Compliance Validation



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

```
when code  
not net_saturated  
expect  
throughput >= raw_bw * 0.9
```

```
expect MAPE <= 20
```

Self-verifiable Experimentation Pipeline (SEP) criteria:

1. Checkout code and data dependencies.
2. Build, install and deploy system under test.
3. Expose pipeline parameters in a plain text file.
4. Allocate resources (statically or dynamically).
5. Capture environment info (hardware, software stacks).
6. Validate results.

parameters.yml  
software.yml  
[true] system scales linearly  
infra-manifest.json  
[false] MAPE not within expected CI



# Validation Dashboard

RUN	MESSAGE	DURATION	COMPLETED
232	Synchronised bandwidth-monitored framework	3h 21s	-
231	Rete monitorata incrementale	2h 5m 58s	22 secs ago
230	Eius accusamus qui omnis.	29 sec	1 mins ago
229	Phased zeroadministration middleware	18 min 22 sec	18 mins ago
228	Innovative 24/7 encryption. Fundamental dynamic matrix. Moratoria programmabile...	22 min 46 sec	36 mins ago
227	Public-key logistical function	20 min 59sec	43 mins ago
226	Asperiores omnis voluptatem autem voluptatem.	1h 0m 51s	1 hour ago
225	Focused national analyzer	7m 24s	3 hours ago
1	-	3m 21s	16 hours ago
3	Exclusive national success. Up-sized high-level matrix. Definizione espansa locale...	37s	22 hours ago
10	Reduced high-level application	5m 27s	1 day ago
1	Sfida switchable massimizzata	1m 46s	6 days ago
4	Future-proofed clear-thinking groupware	13m 48s	1 week ago

Compare



# Differential Analysis

```
1 {  
2   "data": [  
3     {  
4       "type": "articles",  
5       "id": "1",  
6       "attributes": {  
7         "body": "The shortest article. Ever.",  
8         "created": "2015-05-22T14:56:29.000Z",  
9         "updated": "2015-05-22T14:56:28.000Z"  
10      },  
11      "relationships": {  
12        "author": {  
13          "data": {  
14            "id": "42",  
15            "type": "people"  
16          }  
17        }  
18      }  
19    ],  
20    "included": [  
21      {  
22        "type": "people",  
23        "id": "42",  
24        "attributes": {  
25          "name": "John",  
26          "age": 80,  
27          "gender": "male"  
1 {  
2   "data": [  
3     {  
4       "type": "articles",  
5       "id": "1",  
6       "attributes": {  
7         "title": "JSON API paints my bikeshed!",  
8         "body": "The shortest article. Ever.",  
9         "created": "2015-05-22T14:56:29.000Z",  
10        "updated": "2015-05-22T14:56:28.000Z"  
11      },  
12      "relationships": {  
13        "author": {  
14          "data": {  
15            "type": "animals"  
16          }  
17        }  
18      }  
19    ],  
20    "included": [  
21      {  
22        "type": "people",  
23        "id": "42",  
24        "attributes": {  
25          "name": "John",  
26          "age": 80,  
27          "gender": "male"  
28      }  
29    }  
30  ]  
31}
```

Semantic (environment) diff:

- Hardware.
- Software dependencies.
- Parameters.
- Input datasets.
- etc..

# Black Swan – OSS Community-driven Platform

- Success depends on having multiple domain-specific communities involved in the creation, maintenance and curation of pipelines.
- Killer feature: UX – Make it *ridiculously easy* to implement and re-run an experimentation pipeline.

# Use Cases

- Research curation.
  - Institutional libraries (universities and government).
  - Synergies with existing efforts (UC3, DOE code, etc.).
- Technology transfer in R&D scenarios.
  - University / Industrial / Government labs.
  - DOE labs scenario.
- **Reproducibility Evaluation for Conferences and Journals.**
  - **SEP as complement to AD/AE.**
  - **Paper reproduced by SC's Student Cluster Competition.**



Ivo Jimenez  
@ivotron

Thrilled to announce that our "Black Swan" @cross\_ucsc incubator proposal has been accepted!! This means that, after I graduate, I get to keep working on making @getpopper better, and collaborating with more communities to bring DevOps practices to #openscience #opendata #oss

CROSS Funding Decision -- Black Swan

**S** Stephanie Lieggi to Ivo, Carlos, Doug, Karen, Sage, James Oct 24 Dear Ivo,

Congratulations! Your incubator proposal "Black Swan: The Popper Reproducibility Platform" was selected for funding. We can offer you a position to lead this incubator once you have completed your Ph.D.

12:12 PM - 24 Oct 2018

4 Retweets 23 Likes



UNIVERSITY OF CALIFORNIA  
**SANTA CRUZ**

**CROSS**  
CENTER FOR RESEARCH IN  
OPEN SOURCE SOFTWARE

- Incubation of BlackSwan OSS project.
  - Starting Spring 2019.
- Looking for collaborators!
  - Beta-test community aspects of the platform.

Thanks!