

# M2 DATA CAMP 2017

**BALÁZS KÉGL**

Université Paris-Saclay  
CNRS

**ALEXANDRE GRAMFORT**

Université Paris-Saclay  
INRIA

Most classical data challenges are  
**HR** and **publicity** events

# We decided to turn them into a **tool** for

1. Collaborative prototyping
2. Teaching aid
3. Data science process management

# Funded by Université Paris-Saclay

## Team



Balázs Kégl



Alex Gramfort



Akin Kazakçi



Mehdi Cherti



Yohann Sitruk



Guillaume Lemaître



Alexandre Boucaud



Joris Van den Bossche

## Alumni




Djalel Benbouzid



Camille Marini


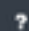
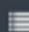

# RAMP.STUDIO

## DATA CHALLENGE WITH *CODE SUBMISSION*



≡

RAMP



Hi Balázs! ▾

Sandbox

You can either edit and save the code in the left column or upload the files in the right column. You can also import code from other submissions when the [leaderboard](#) links are open.


Edit and save your code!

classifier

```
1 from sklearn.base import BaseEstimator
2 from sklearn.ensemble import RandomForestClassifier
3
4
5 class Classifier(BaseEstimator):
6     def __init__(self):
7         pass
8
9     def fit(self, X, y):
10         self.clf = RandomForestClassifier(
11             n_estimators=2, max_leaf_nodes=3, random_state=61)
12         self.clf.fit(X, y)
13
14     def predict(self, X):
15         return self.clf.predict(X)
16
17     def predict_proba(self, X):
18         return self.clf.predict_proba(X)
```

Upload your files!

File list

 classifier.py

Upload file

Choose File

No file chosen

Upload

RAMP

Leaderboard

Combined score: 0.899

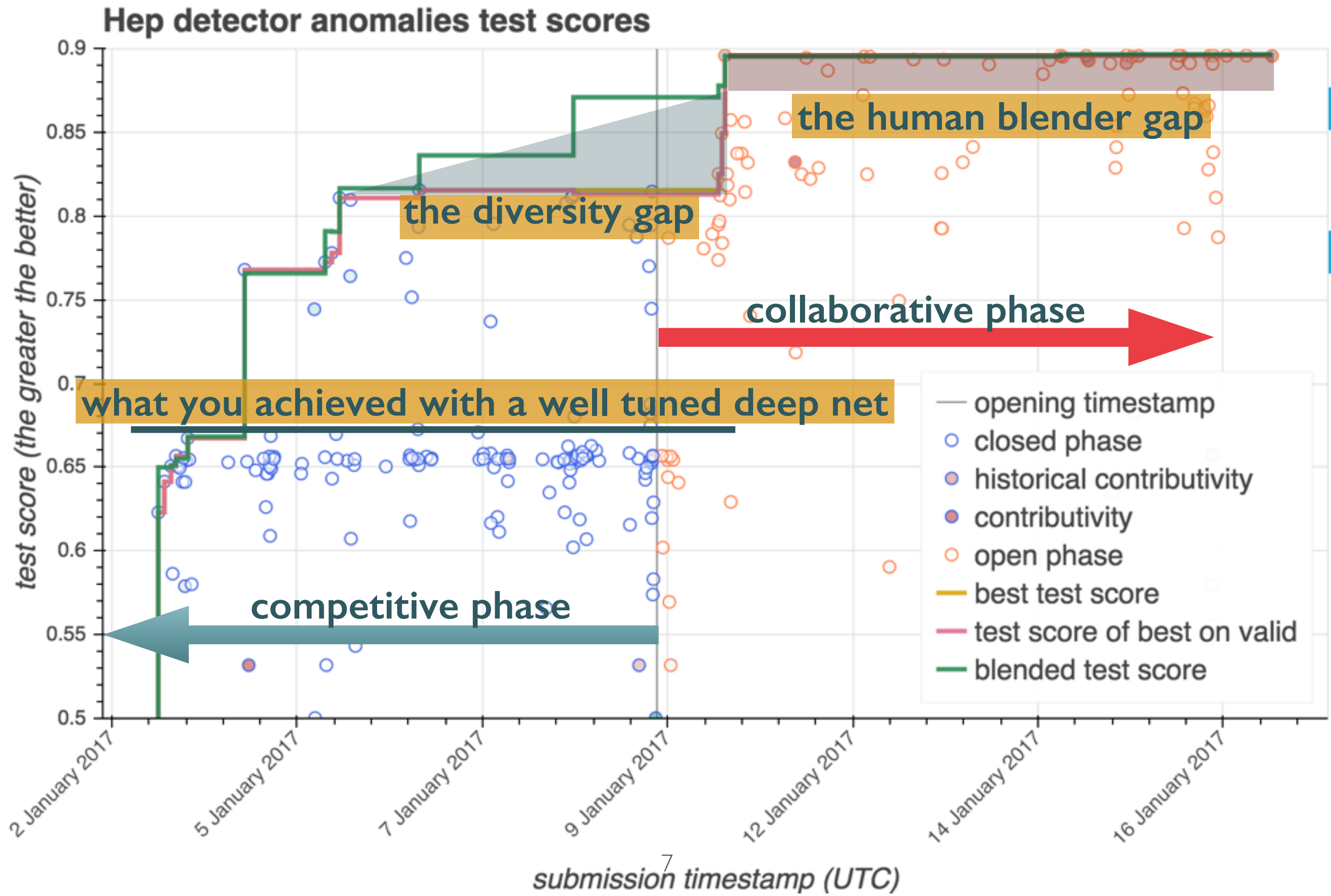
Show  entries

team	submission	contributivity	historical contributivity	auc	accuracy	nll	train time	test time	submitted at (UTC)
diego.souza	<a href="#">tuning_xgboost3</a>	9	5	0.896	0.820	0.385	3074	30	2017-01-17 11:34:53 Tue
ndeye-fatou.diop	<a href="#">kit_from_all</a>	5	1	0.896	0.819	0.382	1167	10	2017-01-14 20:03:00 Sat
diego.souza	<a href="#">tuning_xgboost2</a>	4	2	0.896	0.819	0.385	4900	17	2017-01-15 19:35:03 Sun
ndeye-fatou.diop	<a href="#">kit_from_all_clearer</a>	3	0	0.896	0.819	0.384	1175	10	2017-01-15 03:45:44 Sun
etienne.boursier	<a href="#">combine_features</a>	2	7	0.896	0.820	0.383	2712	3	2017-01-10 15:26:21 Tue
clement.vignac	<a href="#">boursier_improved_1</a>	1	0	0.896	0.819	0.385	2499	4	2017-01-16 08:21:55 Mon

# Code submission

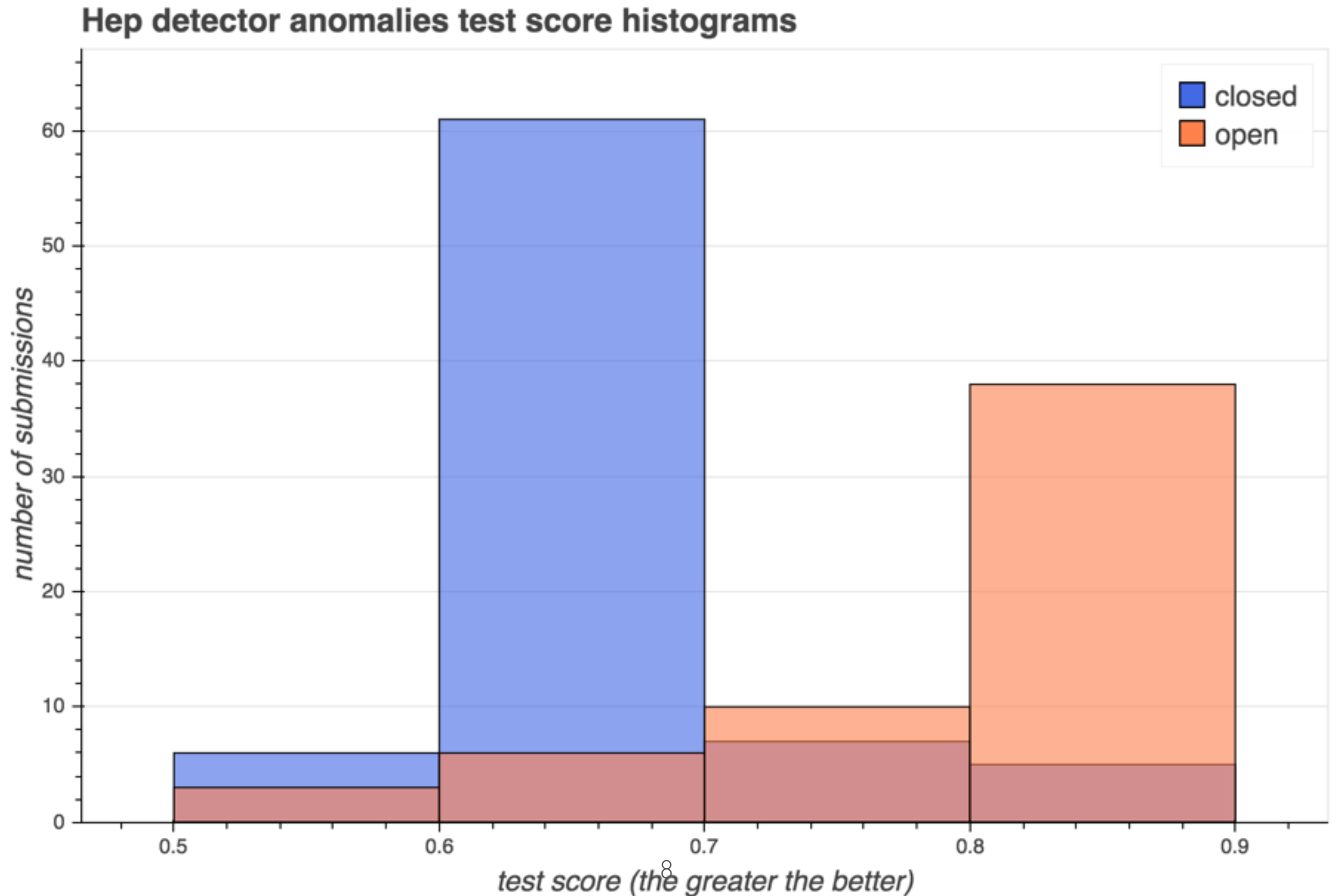
1. lets us deliver a **working prototype**
2. lets the participants **collaborate**
3. makes the **backend challenging** to run (cloud management)

# THE POWER OF THE (COLLABORATING) CROWD



# OPEN PHASE LETS PARTICIPANTS CATCH UP

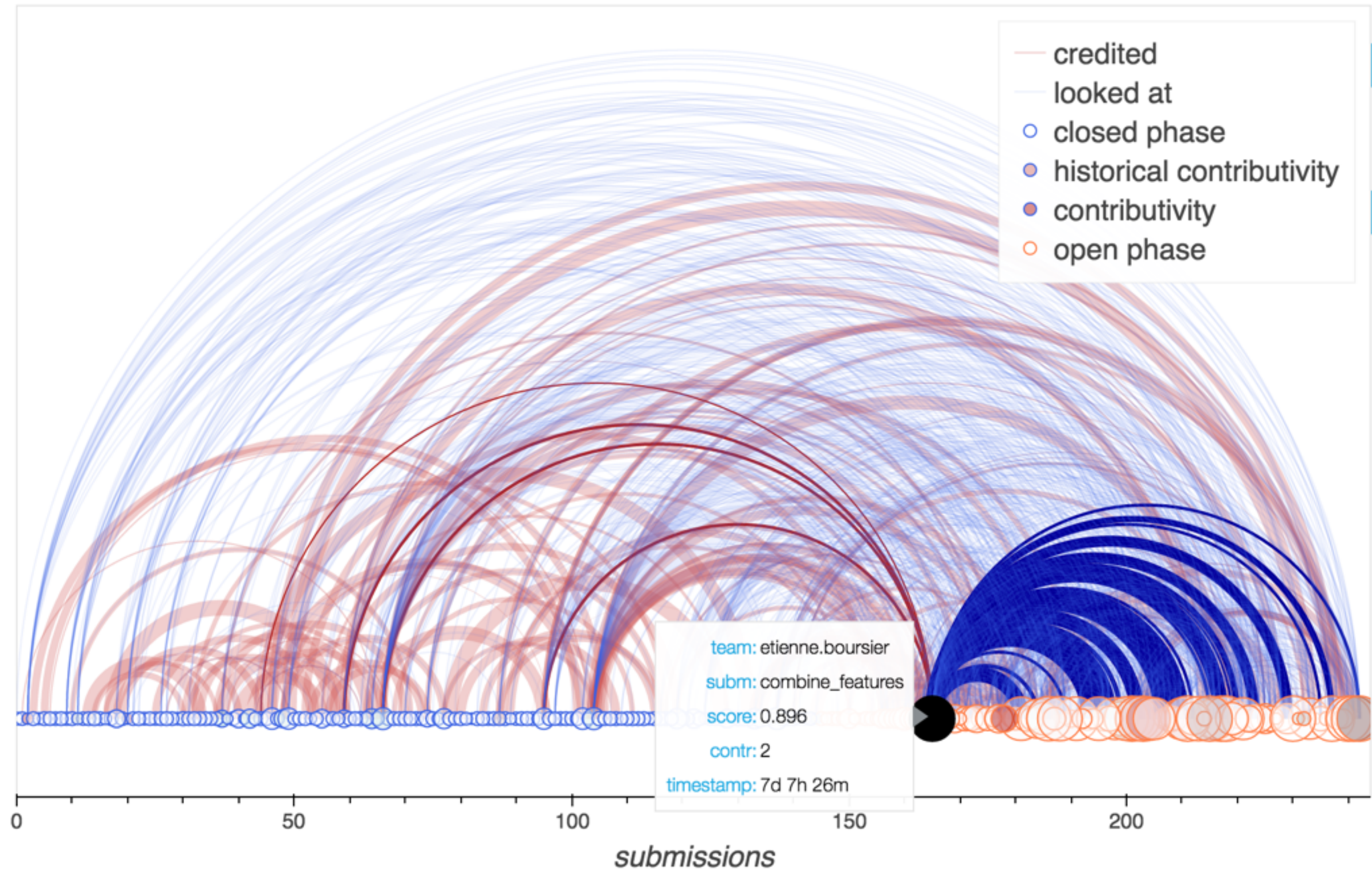
## *THE GOAL OF TEACHING*





# COMMUNICATION AND REUSE

## Hep detector anomalies submissions

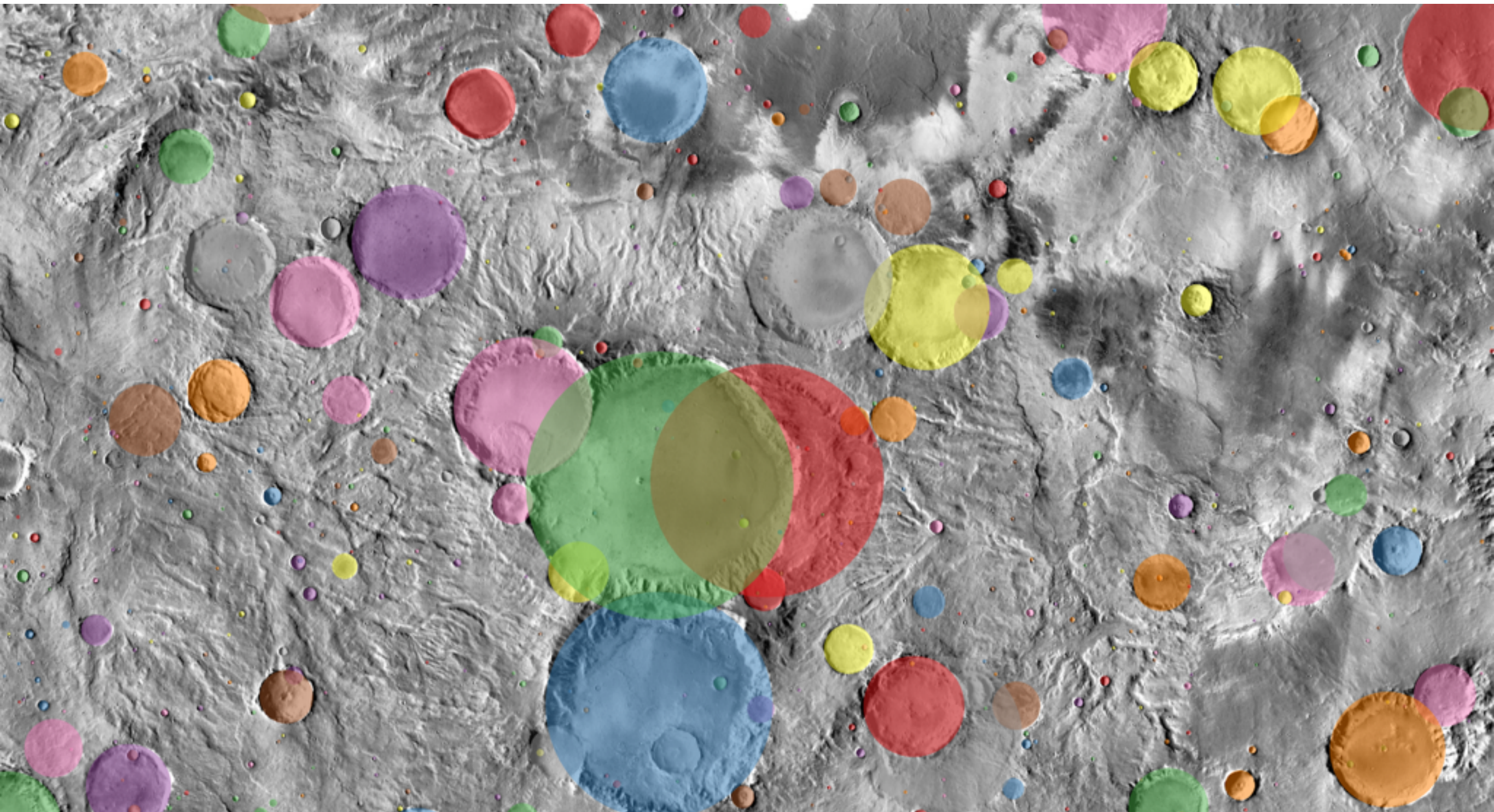


**Choose** one of three challenges



# CHALLENGE I

## DETECT MARS CRATERS





# CHALLENGE 2: FAKE NEWS

## PREDICT THE TRUTHFULNESS OF NEWS



[EDITIONS](#) ▾ [TRUTH-O-METER™](#) ▾ [PEOPLE](#) ▾ [PROMISES](#) ▾ [PANTS ON FIRE](#) [ABOUT US](#)



**Hillary Clinton, Russia, and uranium: What you need to know**



**Can Americans expect a \$4,000 "raise" from Trump tax plan?**



**The big picture: Niger and what we know about what happened**

# CHALLENGE 3: KAGGLE SEGURO

## PREDICT INSURANCE CLAIMS

Secure | <https://www.kaggle.com/c/porto-seguro-safe-driver-prediction/data>

kaggle

Search kaggle



Compe



Featured Prediction Competition

## Porto Seguro's Safe Driver Prediction

Predict if a driver will file an insurance claim next year.



Porto Seguro · 3,218 teams · a month to go (a month to go until m

# RULES

- **Choose one challenge**
  - You will be graded on your **best** score across the three challenges
  - But **you can sign up** for all three
- **Competitive phase until December 17, 20h**
  - except for Kaggle Seguro
  - you will work **on your own**
  - mainly on your computer
  - submit **max once a day** on ramp.studio
  - a total of **50h training time**
    - single CPU
    - on Mars crater, you will be able to request GPUs

# RULES

- Collaborative phase until January 30, 20h
  - You will be incentivized to look at and reuse each other's codes
  - You will be graded on your influence (credits you receive from your fellow students)
  - You will be graded on the marginal improvements of the combined score (to incentivize good but also diverse models)
  - Incentivizing the jumps also means that this phase will probably finish before Christmas
  - A fresh 50h of computational time
  - Grading of the competitive/collaborative phases will be closer to 5/5 than to 8/2 originally announced



# SPECIAL RULES FOR THE KAGGLE CHALLENGE

- **You can only sign up if**
  - You have your Kaggle account
  - You submit a solution at Kaggle with a **public score of at least 0.28**, before **Nov 6 20h**, when the competitive RAMP phase starts
  - You should **not have more than 5 submissions at Kaggle**, but try even less
  - You agree to **enter the RAMP team on Kaggle** and stop submitting to Kaggle on your own
- **Timeline**
  - The **competitive RAMP phase** will end on **Nov 20 at 20h**
  - The **collaborative phase** will run until the **Kaggle submission deadline on November 29**
  - Grading will be similar to the normal challenges, but this challenge will be **open to non-datacamp participation**
  - In case we win money, 50% goes to the CDS, 50% will be shared among the participants, according to a similar scheme we're using to grade you (contribution to jumps and influence)



# HOW DO I CHOOSE

- Challenge 1:

- **Image detection**, **deep learning**, hot but complicated
- we will help but you will need to **manage your GPU-equipped machine** (either your own or on AWS)
- you will **learn a lot** and will be **red hot on the job market**
- potentially a **large margin of improvement over the baseline**
- we're planning a **collective ICML paper**

- Challenge 2:

- **NLP** + **categorical** variables, **small data**, computationally relatively simple
- it is possible that there is **little margin over the baseline**
- we are planning to run a **high-profile money-prize challenge on this later**

- Challenge 3:

- **tabular** but quite **big data** (~1M instances)
- **tight timeline**
- may be quite high profile if we end up close to the top

platform:

[www.ramp.studio](http://www.ramp.studio)

toolkit:

[github.com/paris-saclay-cds/ramp-workflow](https://github.com/paris-saclay-cds/ramp-workflow)

examples:

[github.com/ramp-kits](https://github.com/ramp-kits)

slack:

<https://join.slack.com/t/datacamp2017/signup>

course syllabus:

<http://bit.ly/datacamp2017>