

RAMP

DATA CHALLENGES WITH

MODULARIZATION AND CODE SUBMISSION

LESSONS LEARNED

BALÁZS KÉGL

Université Paris-Saclay / CNRS

WHO AM I?

Balázs Kégl

- Directeur de recherche **CNRS**
 - machine learning (20 years)
interfacing with particle physics (10 years)
- Director of the **Paris-Saclay Center for Data Science**
 - interfacing with biology, economy, climatology, chemistry, etc. (4 years)
- Data science **consulting and training** (4 years)

OUTLINE

- A **short history** of RAMPs
 - **motivations**, **design principles**, and the **current tool**
- Three data challenges
 - **anomaly detection** in the LHC ATLAS detector
 - **classifying** and **quantifying** drug preparations for cancer therapy
 - **time series forecasting** of El Niño
- How can **you use it?**
 - in a **classroom**: to **teach ML**
 - as a **domain science researcher**: to **crowdsource your predictive problem**
 - as a **data science researcher**: to **benchmark your new techniques**

UNIVERSITÉ PARIS-SACLAY

19 *fondateurs*

60 000 *étudiants*

6 000 *doctorants*

15 000 *étudiants
en master*

8 *Schools*

11 000 *chercheurs
et enseignants-chercheurs*

300 *laboratoires*

8 000 *publications /an*

15 % *de la recherche
publique française*

10 *départements*

+ horizontal **multi-disciplinary** and **multi-partner**
initiatives to create cohesion

A multi-disciplinary initiative, **building interfaces**, **matching people**, helping them launching projects

345 affiliated **researchers**, **50 laboratories**

Biology & bioinformatics

IBISC/UEvry
LRI/UPSud
Hepatinov
CESP/UPSud-UVSQ-Inserm
IGM-I2BC/UPSud
MIA/Agro
MIAj-MIG/INRA
LMAS/Centrale

Chemistry

EA4041/UPSud

Earth sciences

LATMOS/UVSQ
GEOPS/UPSud
IPSL/UVSQ
LSCE/UVSQ
LMD/Polytechnique

Economy

LM/ENSAE
RITM/UPSud
LFA/ENSAE

Neuroscience

UNICOG/Inserm
U1000/Inserm
NeuroSpin/CEA

**Particle physics
astrophysics &
cosmology**

LPP/Polytechnique
DMPH/ONERA
CosmoStat/CEA
IAS/UPSud
AIM/CEA
LAL/UPSud

Machine learning

LRI/UPSud
LTCI/Telecom
CMLA/Cachan
LS/ENSAE
LIX/Polytechnique
MIA/Agro
CMA/Polytechnique
LSS/Supélec
CVN/Centrale
LMAS/Centrale
DTIM/ONERA
IBISC/UEvry
LIST/CEA

Visualization

INRIA
LIMSI

Signal processing

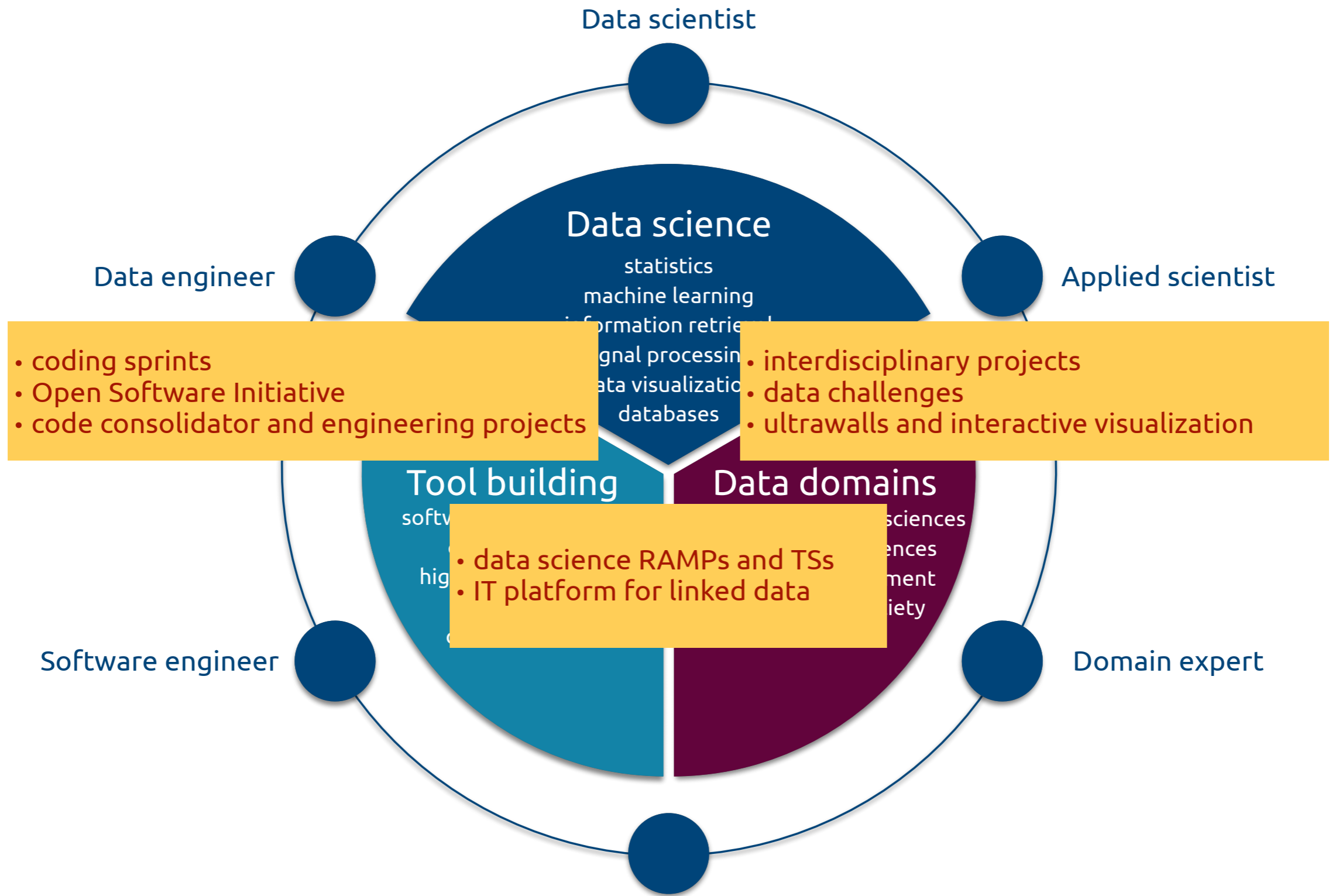
LTCI/Telecom
CMA/Polytechnique
CVN/Centrale
LSS/Supélec
CMLA/Cachan
LIMSI
DTIM/ONERA

Statistics

LMO/UPSud
LS/ENSAE
LSS/Supélec
CMA/Polytechnique
LMAS/Centrale
MIA/AgroParisTech

THE DATA SCIENCE ECOSYSTEM

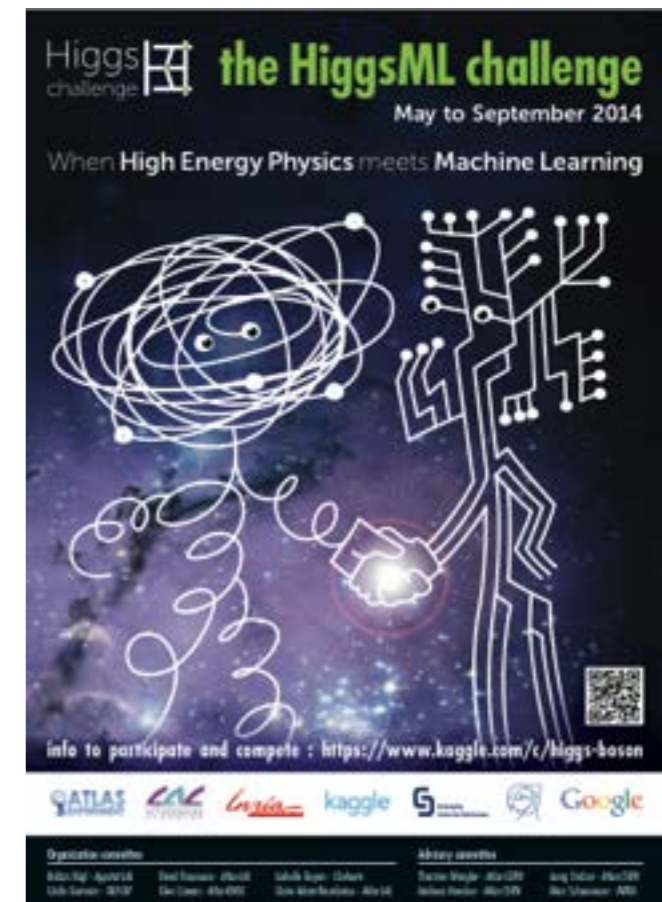
<https://medium.com/@balazskegl>



DATA CHALLENGES



- The **HiggsML** challenge on **Kaggle**
- <https://www.kaggle.com/c/higgs-boson>



HUGE PUBLICITY

kaggle

Customer Solutions

Competitions

Community ▾

Sign up

Login



Completed • \$13,000 • **1,785 teams**

Higgs Boson Machine Learning Challenge

Mon 12 May 2014 – Mon 15 Sep 2014 (21 days ago)

Dashboard ▾

Private Leaderboard - Higgs Boson Machine Learning Challenge

This competition has completed. This leaderboard reflects the final standings.

See someone using multiple accounts?
[Let us know.](#)

#	Δ1w	Team Name <small>‡ model uploaded * in the money</small>	Score <small>🔍</small>	Entries	Last Submission UTC (Best – Last Submission)
1	↑4	Gábor Melis ‡ *	3.80581	110	Sun, 14 Sep 2014 09:10:04 (-0h)
2	↓1	Tim Salimans ‡ *	3.78913	57	Mon, 15 Sep 2014 23:49:02 (-40.6d)
3	—	nhlx5haze ‡ *	3.78682	254	Mon, 15 Sep 2014 16:50:01 (-76.3d)

SIGNIFICANT IMPROVEMENT OVER THE BASELINE

#	Δ1w	Team Name <small>‡ model uploaded * in the money</small>	Score	Entries	Last Submission UTC (Best - Last Submission)
1	↑4	Gábor Melis ‡ *	3.80581	100	Sun, 14 Sep 2014 09:10:04 (-0h)
2	↓1	Tim Salimans ‡ *	3.78682	57	Mon, 15 Sep 2014 23:49:02 (-40.6d)
3	—	nhlx5haze ‡ *	3.78682	254	Mon, 15 Sep 2014 16:50:01 (-76.3d)
4	↑55	ChoKo Team 🏆	3.77526	216	Mon, 15 Sep 2014 15:21:36 (-42.1h)
5	↑23	cheng chen	3.77384	21	Mon, 15 Sep 2014 23:29:29 (-0h)
6	↓2	quantify	3.77086	8	Mon, 15 Sep 2014 16:12:48 (-7.3h)
7	↑73	Stanislav Semenov & Co (HSE Yandex)	3.76211	68	Mon, 15 Sep 2014 20:19:03
8	↓1	Luboš Motl's team 🏆	3.76050	589	Mon, 15 Sep 2014 08:38:49 (-1.6h)
9	↓1	Roberto-UCIIM	3.75864	292	Mon, 15 Sep 2014 23:44:42 (-44d)
10	↑5	Davut & Josef 🏆	3.75838	161	Mon, 15 Sep 2014 23:24:32 (-4.5d)
990	↓65	sandy	3.20546	5	Fri, 29 Aug 2014 18:14:30 (-0.7h)
991	↓65	Rem.	3.19956	2	Mon, 16 Jun 2014 21:53:43 (-30.4h)
📍		simple TMVA boosted trees	3.19956		
992	↓65	Xiaohu SUN	3.19956	3	Tue, 03 Jun 2014 13:14:47
993	↓65	Pierre Boutaud	3.19956	10	Fri, 25 Jul 2014 15:25:07 (-30d)

HUGE PUBLICITY

SIGNIFICANT IMPROVEMENT OVER THE BASELINE

yet partially **missing the objectives**

LIMITATIONS OF DATA CHALLENGES

- Organizers have **no direct access to solutions**
- Emphasize **competition**: participants **cannot build on each other's solutions**
- **No modularization**: **ideas go unnoticed** unless packaged into a top submission

We decided to design something better:

Data challenge with **code submission**

GOAL

- Design a **crowdsourcing** and **teaching** tool that
 - **hides heavy computational details** and provides a **simple interface** to data scientists to **experiment with algorithms**
 - promotes **collaboration** and the rapid **propagation of ideas**
 - **modularizes complex pipelines** so the different expertise can be applied without having to understand all the details of the full workflow
 - allows the challenge organizer to walk away with a **working prototype**

RAPID ANALYTICS AND MODEL PROTOTYPING (RAMP)

<http://www.ramp.studio>

Team



Balázs Kégl



Alex Gramfort



Akin Kazakçi



Mehdi Cherti



Yohann Sitruk

Alumni



Djalel Benbouzid



Camille Marini

RAMP

RAPID ANALYTICS AND MODEL PROTOTYPING

- Roughly two formats
 - **single day hackatons** with 20-50 participants, **open leaderboard**, 15 minute timeout
 - 1-3 week **classroom challenges** up 150 students (but no limit really): **closed phase** followed by an **open phase**
- **800+ users**, **5000+** predictive models

CURRENT RAMPs

www.ramp.studio/problems

RAMP Hi Balazs!

- **Pollenating insect classification (209 classes)**
 - La Paillasse / Futur en Seine, number of participants = **28**, number of submissions = **13**, combined score = **0.831**, [click here for score vs time plot](#)
- **Titanic survival classification**
 - DSSP 6 2016/17 2, number of participants = **31**, number of submissions = **34**, combined score = **0.87**, [click here for score vs time plot](#)
 - Entry exam to deep learning tutorial, number of participants = **35**, number of submissions = **21**, combined score = **0.86**, [click here for score vs time plot](#)
 - Ecole des Mines 2016/17, number of participants = **125**, number of submissions = **144**, combined score = **0.89**, [click here for score vs time plot](#)
- **Pollenating insect classification (18 classes)**
 - Polytechnique MAP583/MAP542 2016/17, number of participants = **166**, number of submissions = **114**, combined score = **0.959**, [click here for score vs time plot](#)
 - DSSP5 2017, number of participants = **15**, number of submissions = **24**, combined score = **0.93**, [click here for score vs time plot](#)
- **Particle tracking in the LHC ATLAS detector**
 - Initial single-day RAMP 2017, number of participants = **55**, number of submissions = **60**, combined score = **0.97**, [click here for score vs time plot](#)
- **El Nino forecast**
 - single-day RAMP at Climate Informatics Workshop 2015; Saclay Data Camp 2016/17, number of participants = **160**, number of submissions = **138**, combined score = **0.389**, [click here for score vs time plot](#)
- **Arctic sea ice forecast**
 - single-day RAMP at Climate Informatics Workshop 2016, number of participants = **46**, number of submissions = **83**, combined score = **0.31**, [click here for score vs time plot](#)
 - Polytechnique MAP542 2016/17, number of participants = **20**, number of submissions = **52**, combined score = **0.268**, [click here for score vs time plot](#)
 - Polytechnique MAP583 2016/17, number of participants = **123**, number of submissions = **252**, combined score = **0.259**, [click here for score vs time plot](#)
- **Number of air passengers prediction**
 - DSSP4/5 2016, number of participants = **95**, number of submissions = **242**, combined score = **0.236**, [click here for score vs time plot](#)
 - DSSP6 2017, number of participants = **23**, number of submissions = **59**, combined score = **0.268**, [click here for score vs time plot](#)
- **Drug classification and concentration estimation from Raman spectra**
 - Polytechnique MAP583 2016/17, number of participants = **125**, number of submissions = **258**, combined score = **0.048**, [click here for score vs time plot](#)
 - initial single-day RAMP 2016; Saclay Data Camp 2016/17, number of participants = **242**, number of submissions = **554**, combined score = **0.027**, [click here for score vs time plot](#)
 - Ecole des Mines 2016/17, number of participants = **124**, number of submissions = **560**, combined score = **0.023**, [click here for score vs time plot](#)
- **Detecting anomalies in the LHC ATLAS detector**
 - Polytechnique MAP542 2016/17, number of participants = **29**, number of submissions = **47**, combined score = **0.865**, [click here for score vs time plot](#)
 - Polytechnique MAP583 2016/17, number of participants = **133**, number of submissions = **275**, combined score = **0.899**, [click here for score vs time plot](#)
 - initial single-day RAMP 2016, number of participants = **49**, number of submissions = **19**, combined score = **0.677**, [click here for score vs time plot](#)
- **Epidemium cancer mortality rate prediction (2nd RAMP)**
 - initial single-day RAMP 2016, number of participants = **39**, number of submissions = **46**, combined score = **21.79**, [click here for score vs time plot](#)
 - Polytechnique MAP583 2016/17, number of participants = **128**, number of submissions = **192**, combined score = **18.59**, [click here for score vs time plot](#)
 - Polytechnique MAP542 2016/17, number of participants = **22**, number of submissions = **57**, combined score = **19.31**, [click here for score vs time plot](#)

DATA SCIENCE THEMES

Data science themes

- **classification**

- Iris classification
- Detecting anomalies in the LHC ATLAS detector
- Drug classification and concentration estimation from Raman spectra
- Titanic survival classification
- Pollenating insect classification (18 classes)
- Pollenating insect classification (209 classes)

- **convolutional networks**

- Pollenating insect classification (18 classes)
- Pollenating insect classification (209 classes)

- **external data**

- Number of air passengers prediction

- **feature engineering**

- El Nino forecast
- Arctic sea ice forecast
- Drug classification and concentration estimation from Raman spectra
- Detecting anomalies in the LHC ATLAS detector

- **forests**

- Iris classification
- Detecting anomalies in the LHC ATLAS detector
- Titanic survival classification
- Boston housing price regression
- El Nino forecast
- Arctic sea ice forecast
- Number of air passengers prediction
- Epidemium cancer mortality rate prediction (2nd RAMP)

- **functional data**

- Drug classification and concentration estimation from Raman spectra

- **image data**

- Pollenating insect classification (18 classes)
- Pollenating insect classification (209 classes)
- El Nino forecast

- **missing data**

- Epidemium cancer mortality rate prediction (2nd RAMP)
- Titanic survival classification

- **neural networks (deep learning)**

- Drug classification and concentration estimation from Raman spectra
- Pollenating insect classification (18 classes)
- Pollenating insect classification (209 classes)

- **regression**

- Boston housing price regression
- El Nino forecast
- Arctic sea ice forecast
- Number of air passengers prediction
- Drug classification and concentration estimation from Raman spectra
- Epidemium cancer mortality rate prediction (2nd RAMP)

- **small data**

- Drug classification and concentration estimation from Raman spectra
- Epidemium cancer mortality rate prediction (2nd RAMP)
- Detecting anomalies in the LHC ATLAS detector
- El Nino forecast
- Arctic sea ice forecast
- Number of air passengers prediction
- Particle tracking in the LHC ATLAS detector

- **supervised clustering (unsupervised classification)**

- Particle tracking in the LHC ATLAS detector

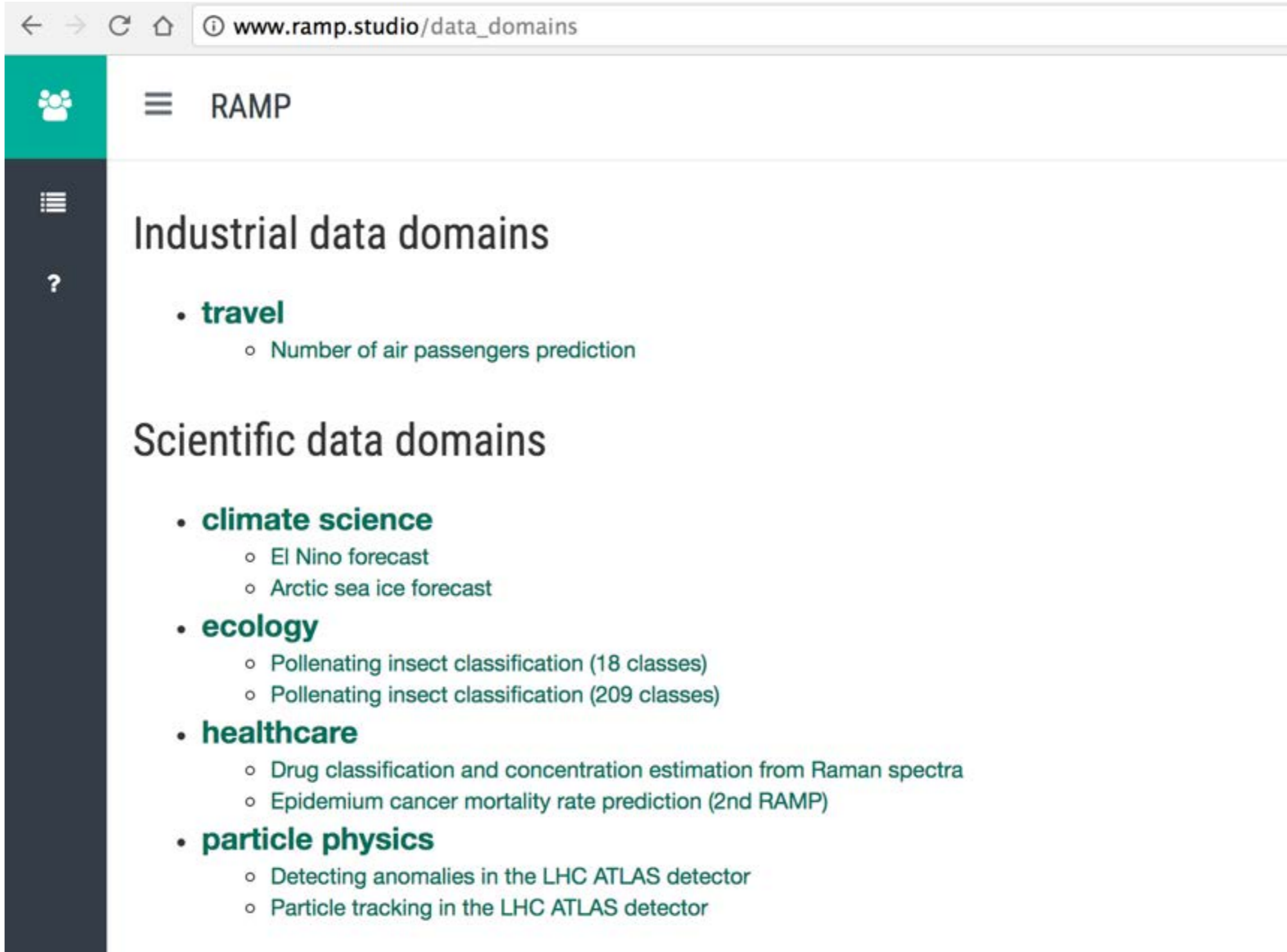
- **tabular data**

- Iris classification
- Detecting anomalies in the LHC ATLAS detector
- Titanic survival classification
- Boston housing price regression
- Number of air passengers prediction
- Epidemium cancer mortality rate prediction (2nd RAMP)

- **time series forecasting**



- El Nino forecast
- Arctic sea ice forecast

DATA DOMAINS



The screenshot shows a web browser at the URL www.ramp.studio/data_domains. The page features a dark sidebar with a teal header containing a group icon and the text 'RAMP'. Below the sidebar, the main content is organized into two sections: 'Industrial data domains' and 'Scientific data domains'. The 'Industrial data domains' section lists a single domain: 'travel', which includes the sub-domain 'Number of air passengers prediction'. The 'Scientific data domains' section lists four domains: 'climate science' (with sub-domains 'El Nino forecast' and 'Arctic sea ice forecast'), 'ecology' (with sub-domains 'Pollenating insect classification (18 classes)' and 'Pollenating insect classification (209 classes)'), 'healthcare' (with sub-domains 'Drug classification and concentration estimation from Raman spectra' and 'Epidemium cancer mortality rate prediction (2nd RAMP)'), and 'particle physics' (with sub-domains 'Detecting anomalies in the LHC ATLAS detector' and 'Particle tracking in the LHC ATLAS detector').

← → ↻ 🏠 ⓘ www.ramp.studio/data_domains

  RAMP

Industrial data domains

- **travel**
 - Number of air passengers prediction

Scientific data domains

- **climate science**
 - El Nino forecast
 - Arctic sea ice forecast
- **ecology**
 - Pollenating insect classification (18 classes)
 - Pollenating insect classification (209 classes)
- **healthcare**
 - Drug classification and concentration estimation from Raman spectra
 - Epidemium cancer mortality rate prediction (2nd RAMP)
- **particle physics**
 - Detecting anomalies in the LHC ATLAS detector
 - Particle tracking in the LHC ATLAS detector

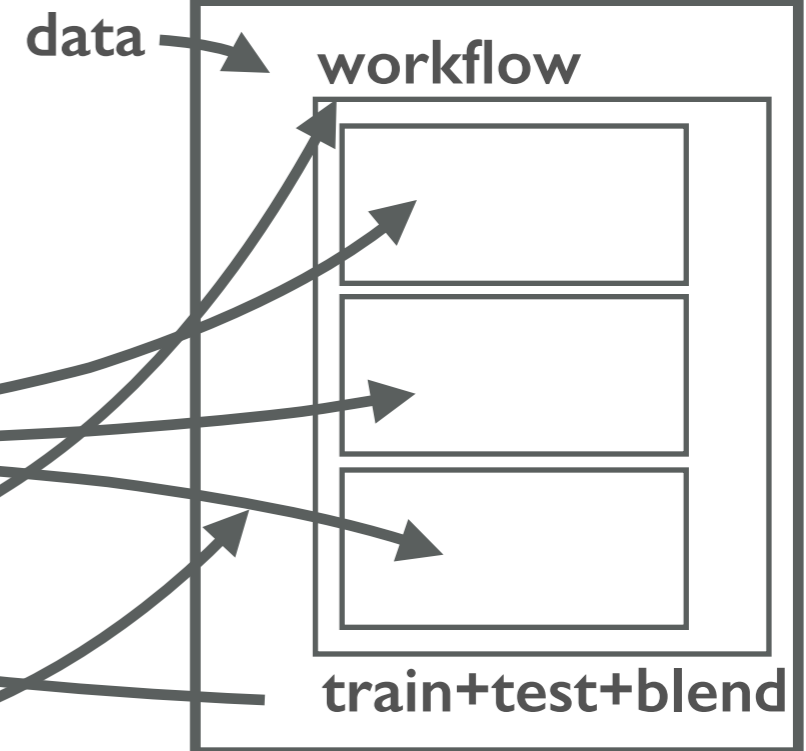
RAMP

DATA CHALLENGE WITH **CODE SUBMISSION**

frontend



backend



- users
- submissions
- score
- problems
- workflow
- starting kit
- crossval

sea_ice_M1XMAP583_201617

Description

Balázs Kégl (CNRS), Camille Marini (CNRS), Andy Rhines (UW), Jennifer Dy (NEU), Arindam Banerjee (UMN)

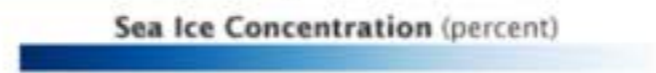
Introduction

Arctic sea ice cover is one of the most variable features of Earth's climate. Its annual cycle peaks at around 15 million square kilometers in early spring, melting back to a minimum of about 6 million square kilometers in September. These seasonal swings are important for Earth's energy balance, as ice reflects the majority of sunlight while open water absorbs it. Changes in ice cover are also important for marine life and navigation for shipping.

Arctic Minimum (September 14, 2008)



Arctic Maximum (February 28, 2009)



RAMP

www.ramp.studio/events/sea_ice_M1XMAP583_201617/sandbox



RAMP

Hi Balazs!

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!

ts_feature_extractor

```
1 import numpy as np
2 import xarray as xr
3 from sklearn.linear_model import LinearRegression
4
5 class FeatureExtractor(object):
6
7     def __init__(self):
8         pass
9
10    def transform(self, X_ds):
11        """Compute the monthly averages of the ice_area, corresponding to the month
12        The code could be simplified but in this way it is general, can be used for
13        variables as well."""
14        # This is the range for which features should be provided. Strip
15        # the burn-in from the beginning and the prediction look-ahead from
16        # the end.
17        valid_range = np.arange(X_ds.attrs['n_burn_in'], len(X_ds['time']))
18
19        # We convert the Dataset into a 4D DataArray
20        X_xr = X_ds.to_array()
21
```

regressor

Upload your files!

File list

ts_feature_extractor.py

regressor.py

Upload file

Choose File No file chosen

Upload

RAMP



sea_ice_M1XMAP583_201617

Leaderboard

Combined score: 0.268

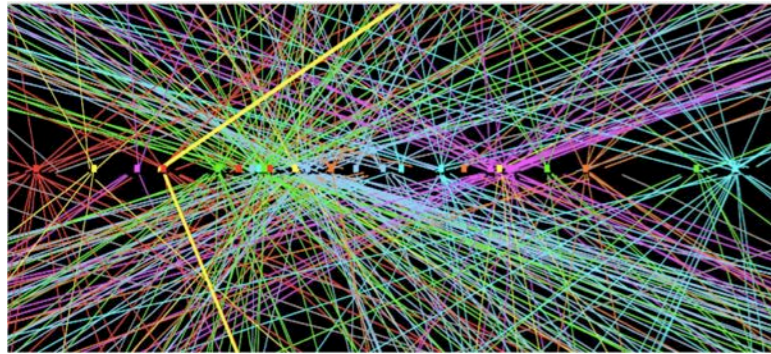
Show 10 entries

Search:

team	submission	contributivity	historical contributivity	rmse	train time	test time	submitted at (UTC)
joseph.budin	noName	26	3	0.279	286	3	2017-02-13 11:36:28 Mon
alexis.thual	timeseries	16	16	0.296	1	1	2017-02-13 17:48:47 Mon
julien.habis	try_hard3	11	8	0.300	475	3	2017-02-13 19:45:35 Mon
kangzheng.liang	thirdtry	7	7	0.291	8	1	2017-02-07 19:11:32 Tue
joseph.budin	LinReg	6	3	0.280	234	3	2017-02-13 11:25:39 Mon
gaetan.millerand	shifted+boost+nino	6	5	0.295	29	5	2017-02-04 21:04:11 Sat
thibaut.vasseur	starting_kit_help	4	4	0.289	17	9	2017-02-13 18:48:37 Mon
yu-jia.cheong	Last	3	3	0.289	18	7	2017-02-13 13:28:53 Mon
gaetan.millerand	random_test	3	3	0.295	30	5	2017-02-07 13:12:29 Tue
maxime.lapides	TestFinal	3	3	0.296	458	3	2017-02-13 17:44:37 Mon

Showing 1 to 10 of 172 entries

ANOMALY DETECTION IN THE LHC ATLAS DETECTOR



reconstruction
+ simulated anomalies

DER_mass_transverse_met_lep	1.937
DER_mass_vis	64.546
DER_pt_h	41.791
DER_deltar_tau_lep	2.301
DER_pt_tot	7.975
DER_sum_pt	105.305
DER_pt_ratio_lep_tau	0.926
DER_met_phi_centrality	1.087
PRI_tau_pt	36.259
PRI_tau_eta	-2.248
PRI_tau_phi	-2.239
PRI_lep_pt	33.582
PRI_lep_eta	-1.893
PRI_lep_phi	0.035
PRI_met	19.872
PRI_met_phi	-0.040
isSkewed	0.000

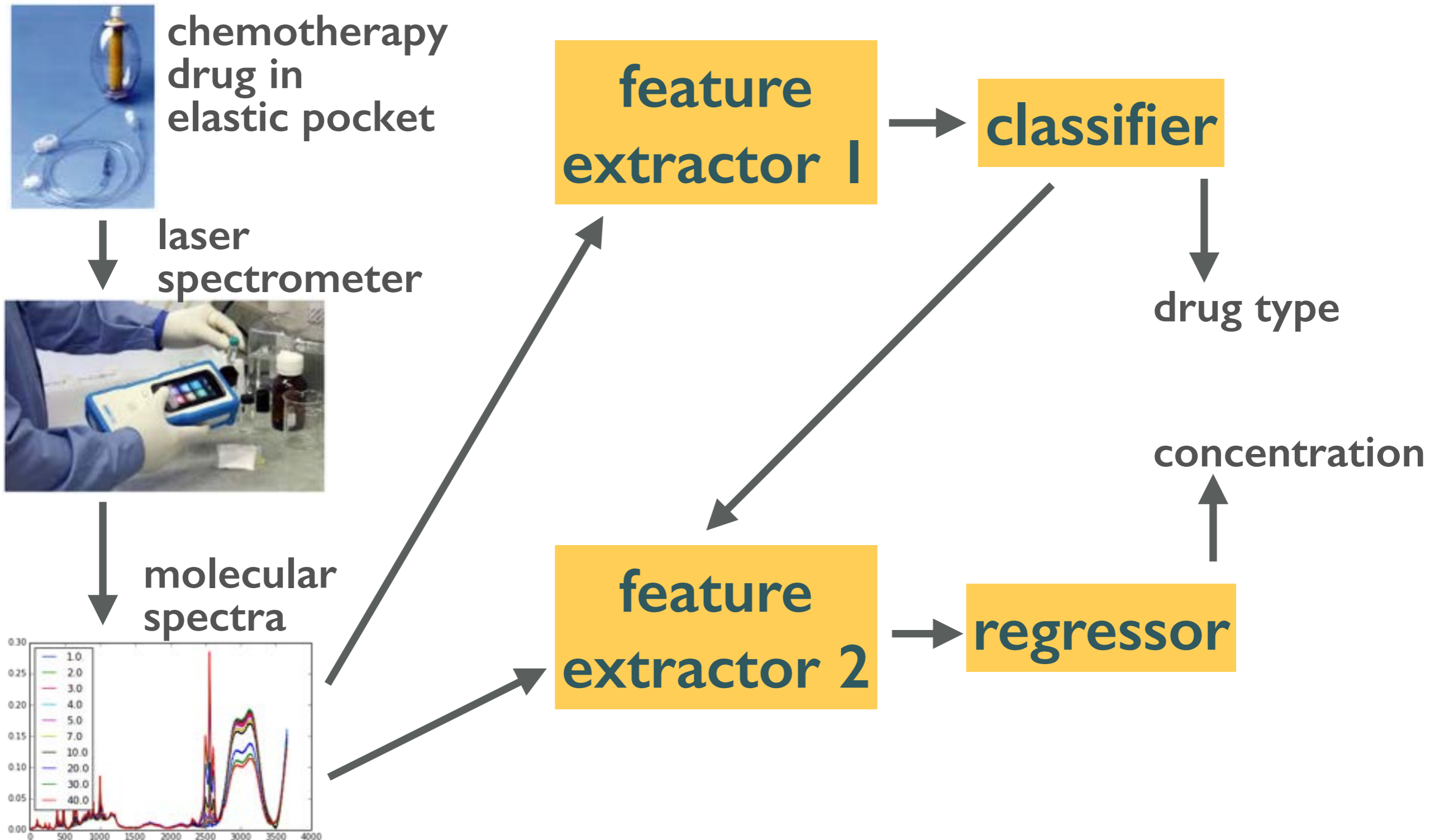
classifier

correct
(isSkewed = 0)

?

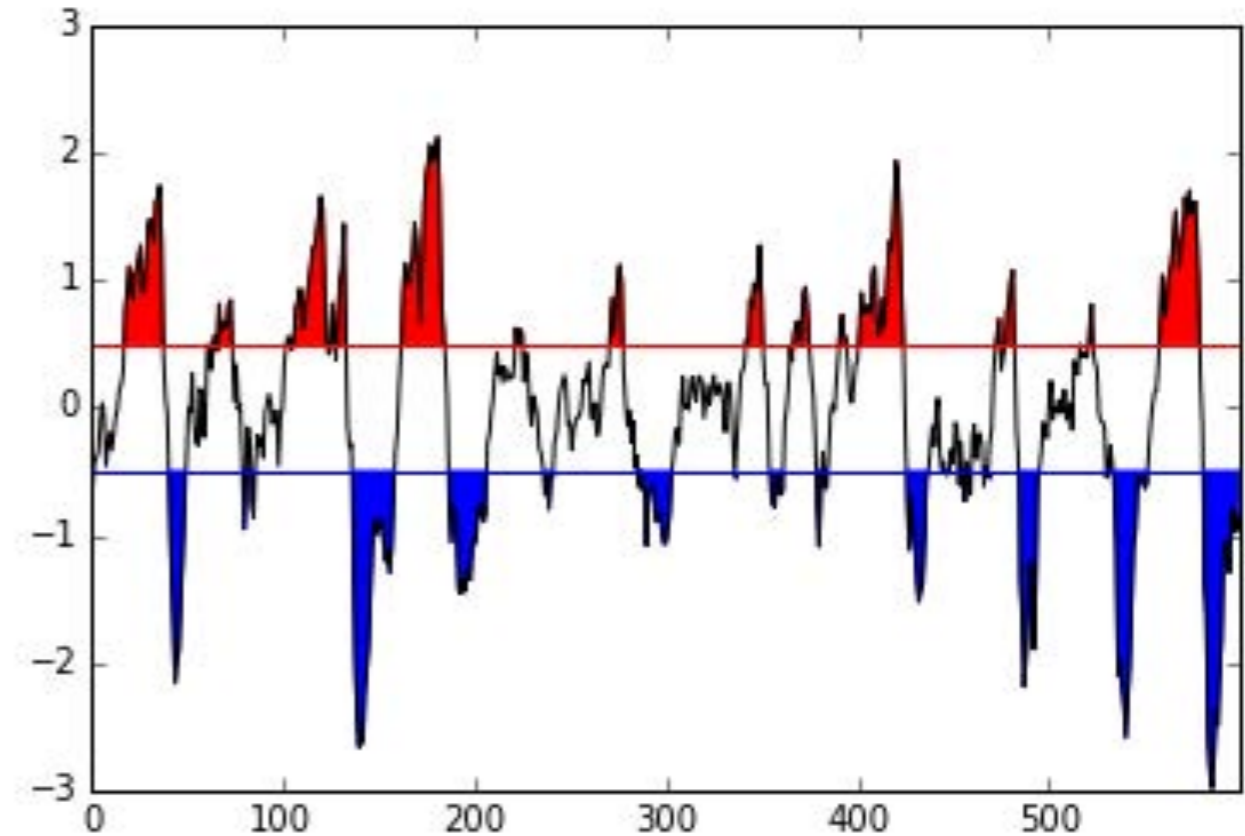
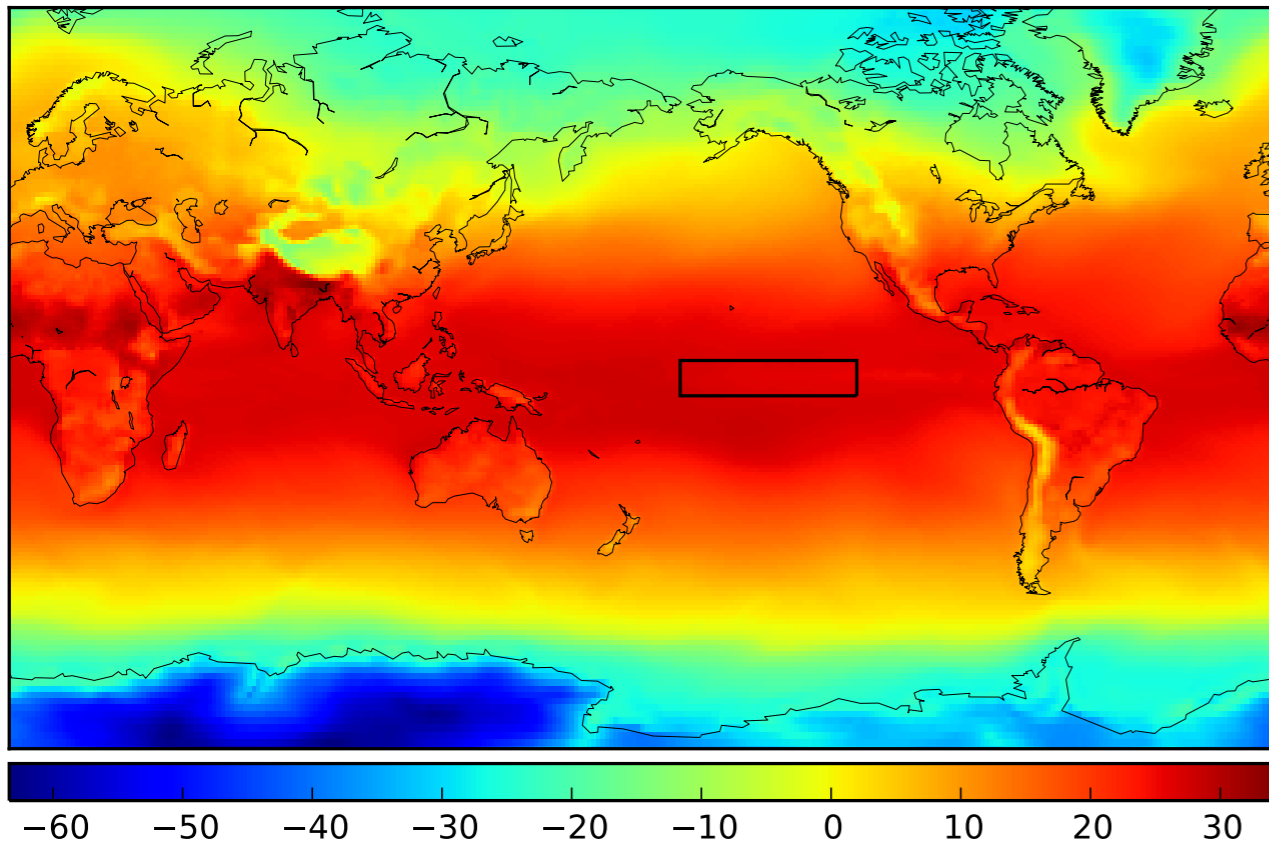
anomaly
(isSkewed = 1)

CLASSIFYING AND REGRESSING ON MOLECULAR SPECTRA

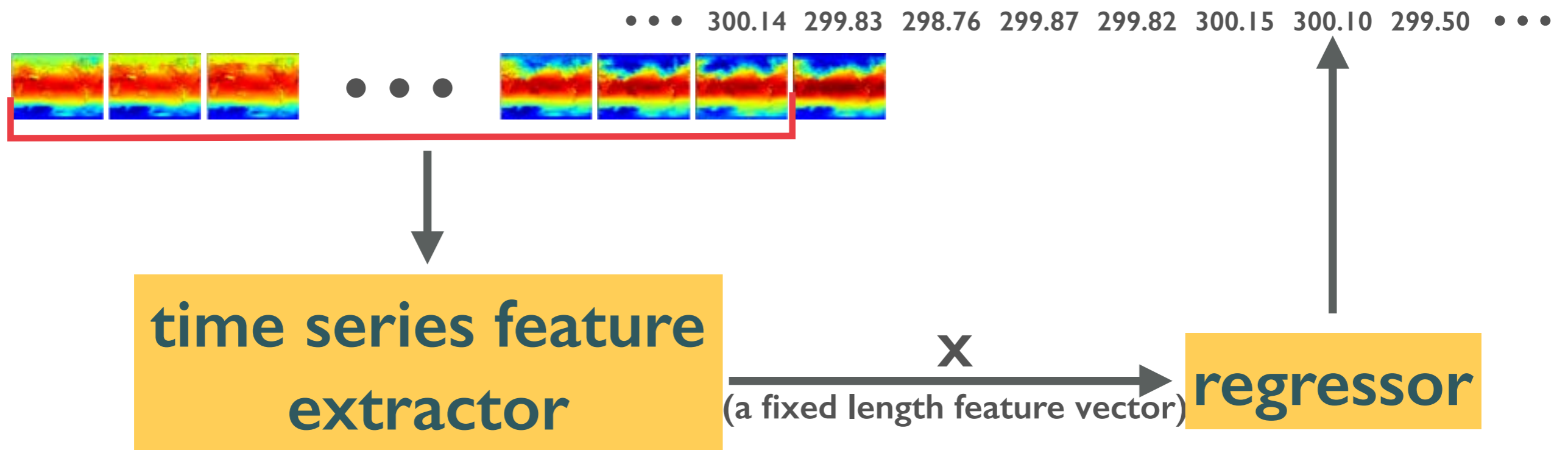


FORECASTING EL NINO SIX MONTHS AHEAD

Temperature map



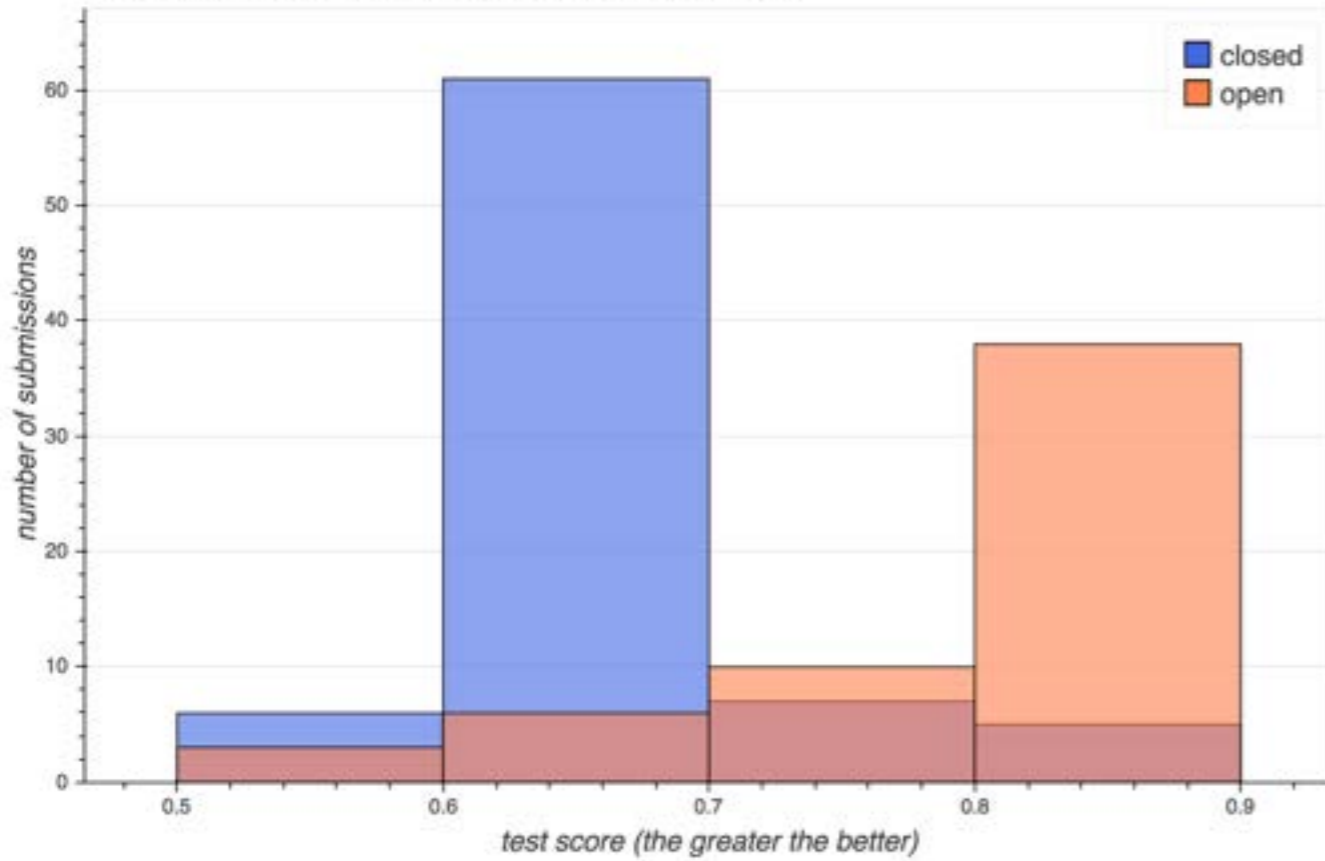
FORECASTING EL NINO SIX MONTHS AHEAD



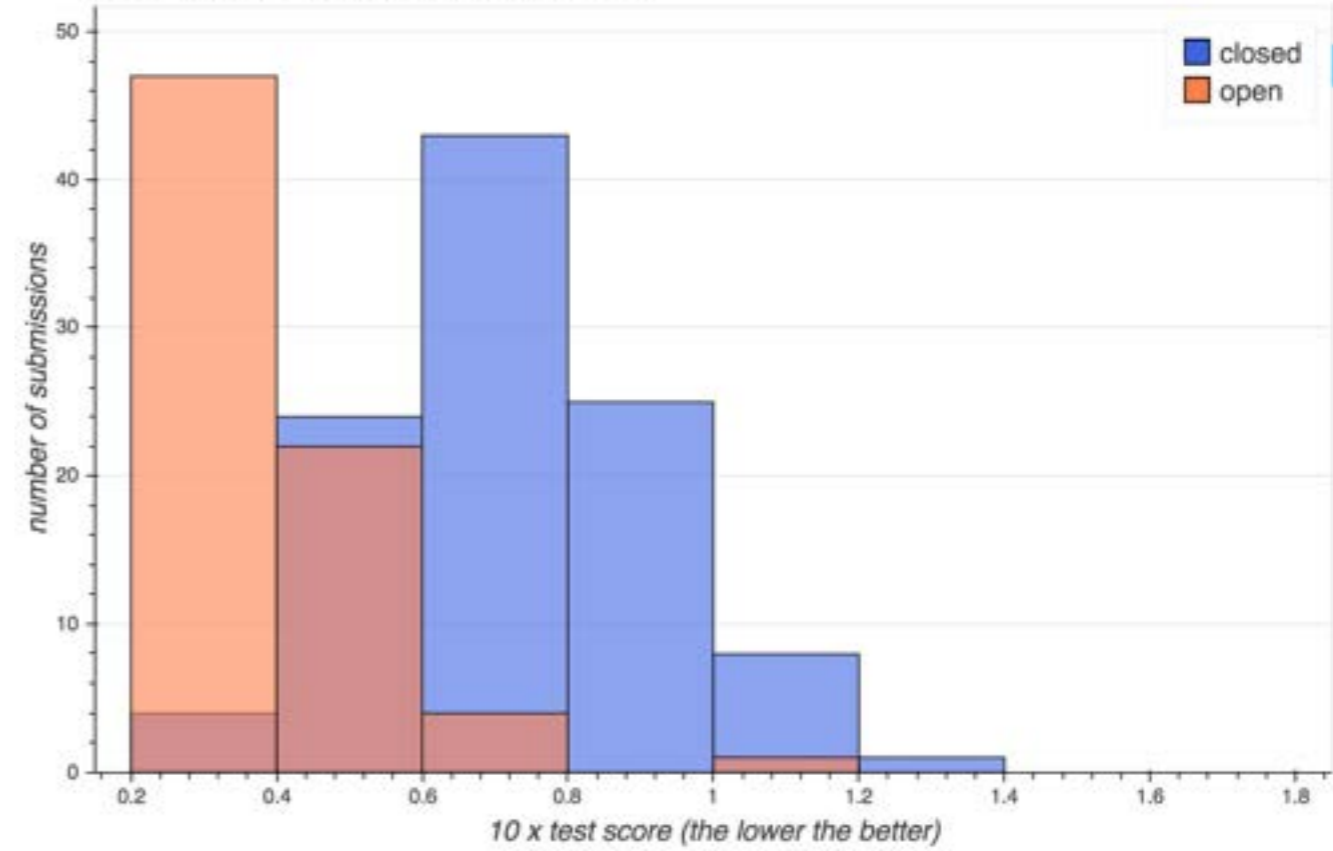
Analyzing the process

OPEN PHASE LETS PARTICIPANTS CATCH UP

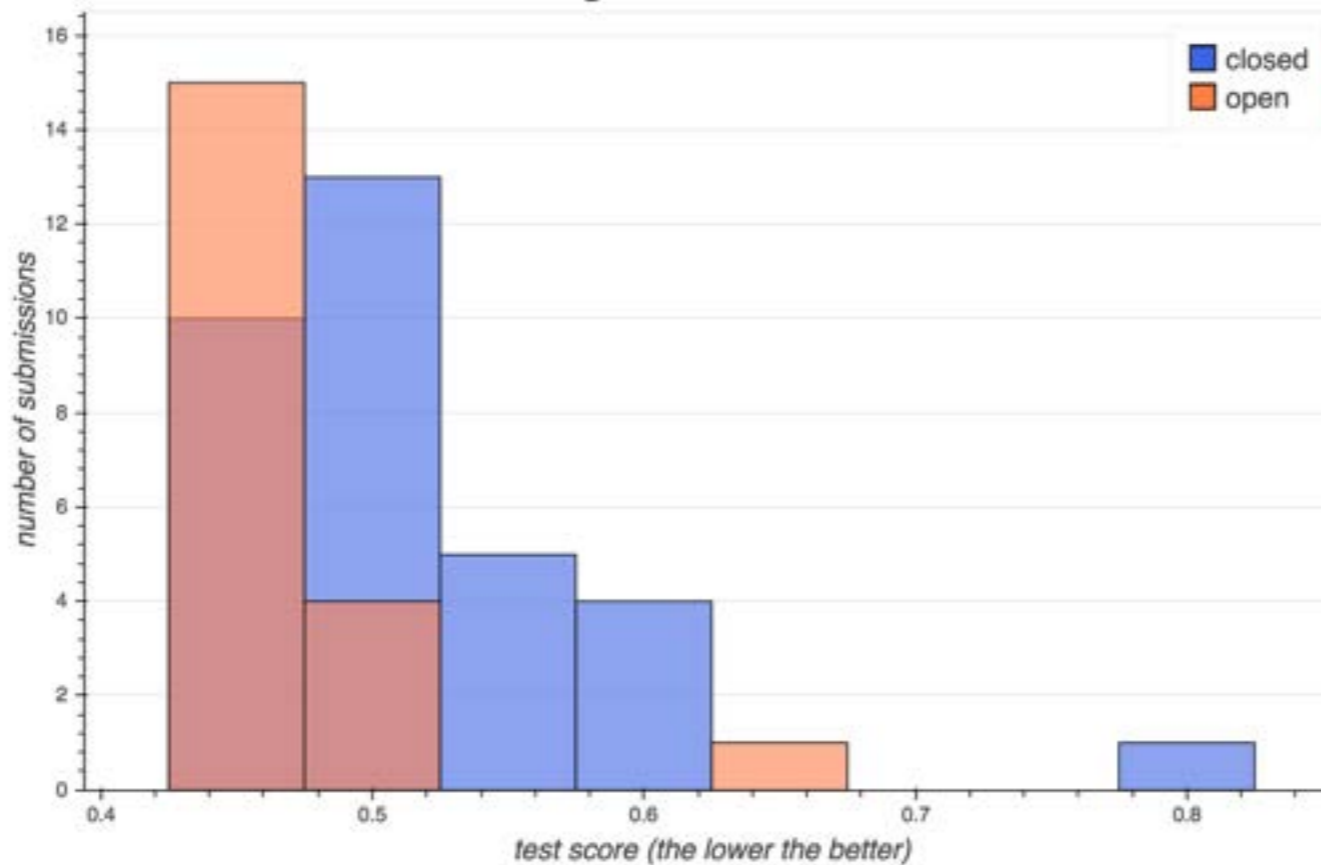
Hep detector anomalies test score histograms



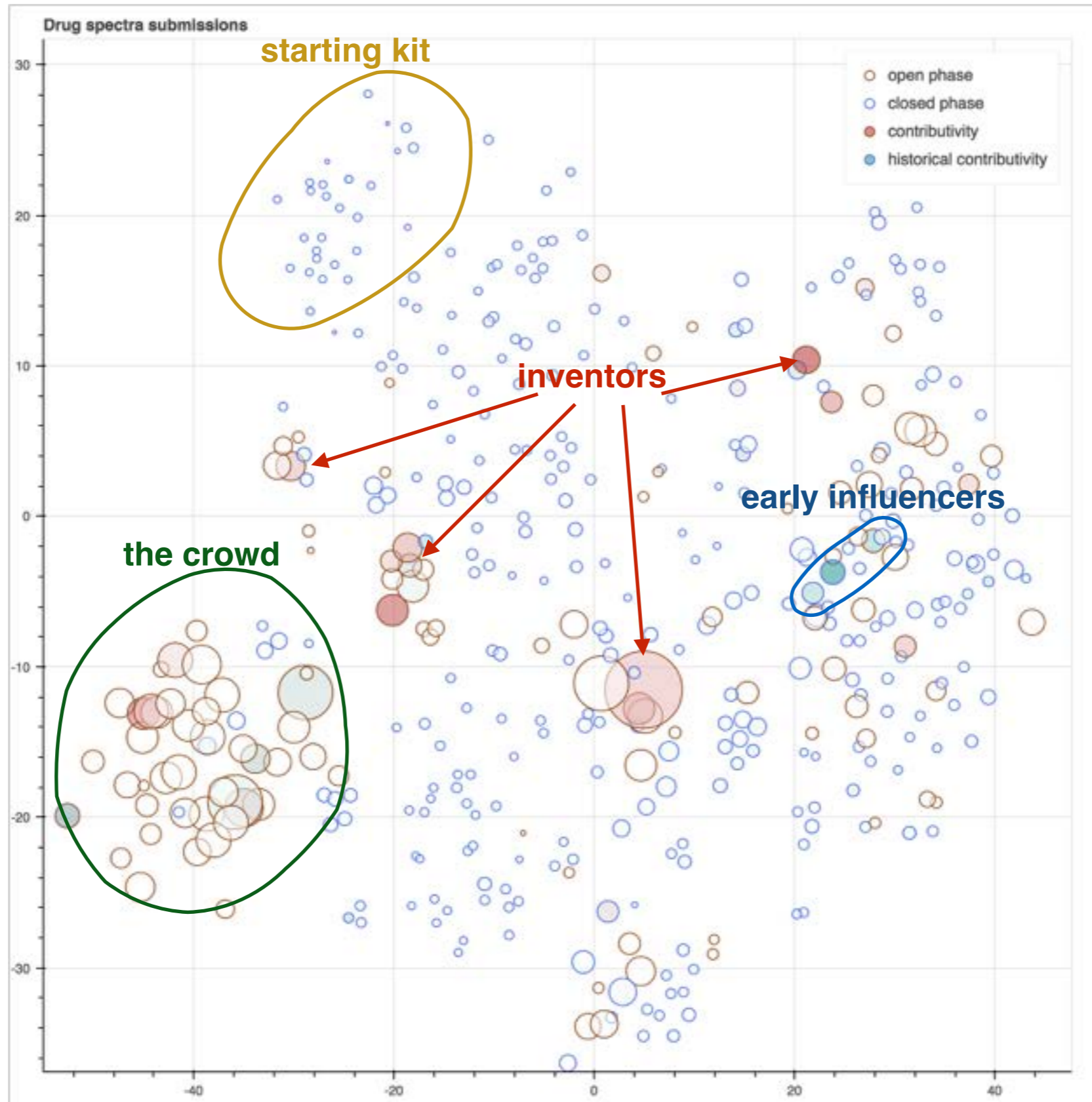
Drug spectra test score histograms



El nino forecast test score histograms

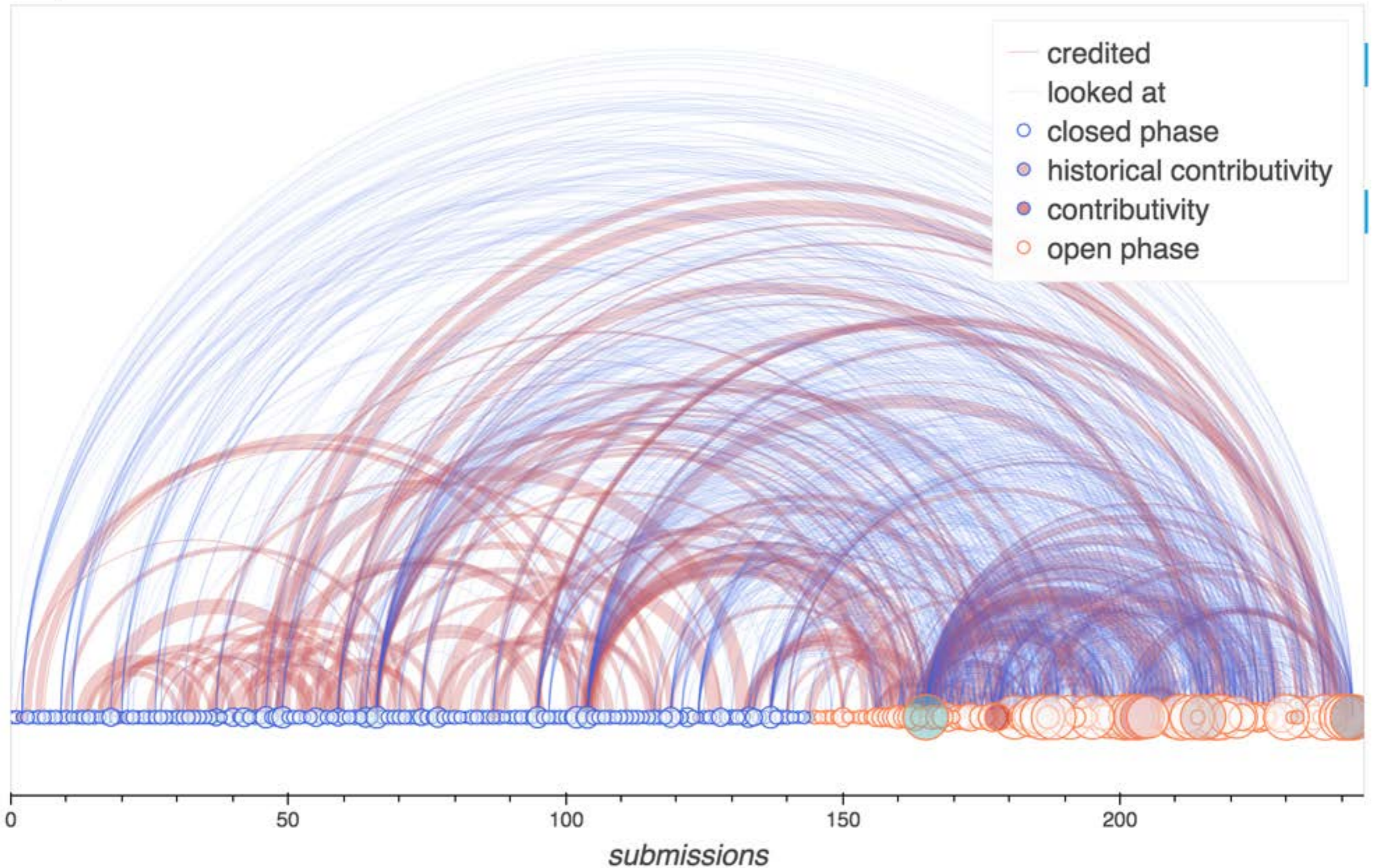


THE DYNAMICS OF COLLABORATION



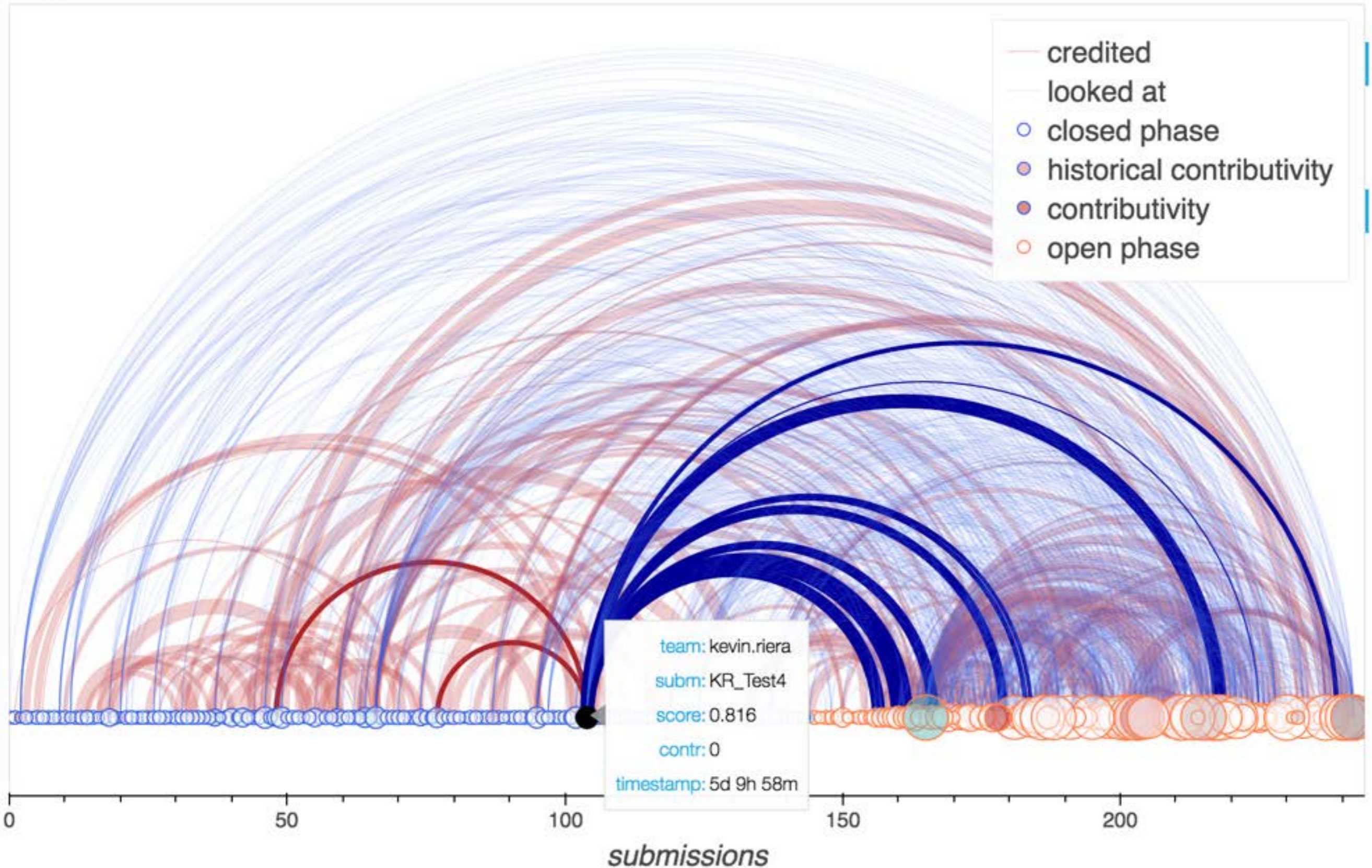
THE DYNAMICS OF COLLABORATION

Hep detector anomalies submissions



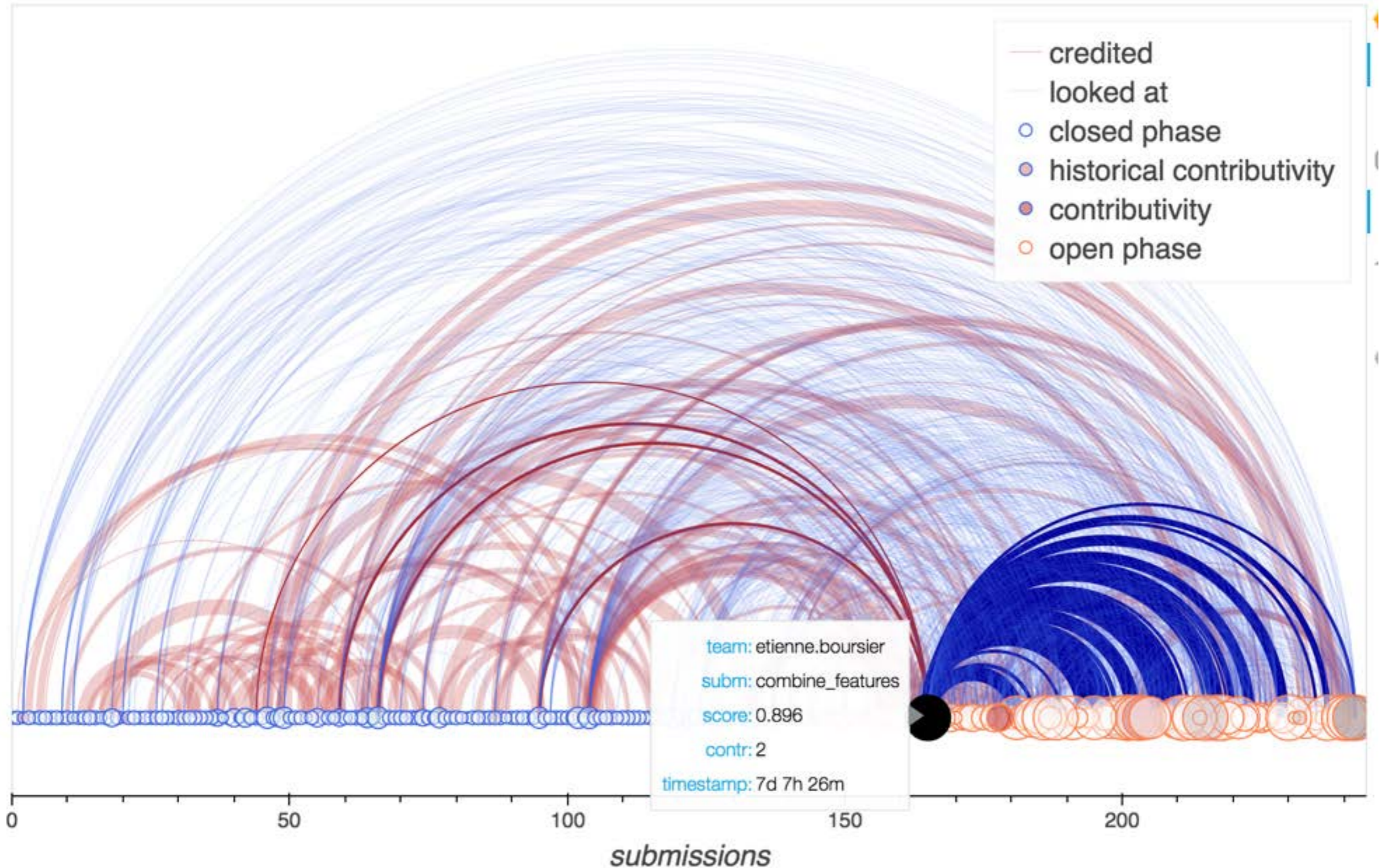
THE DYNAMICS OF COLLABORATION

Hep detector anomalies submissions

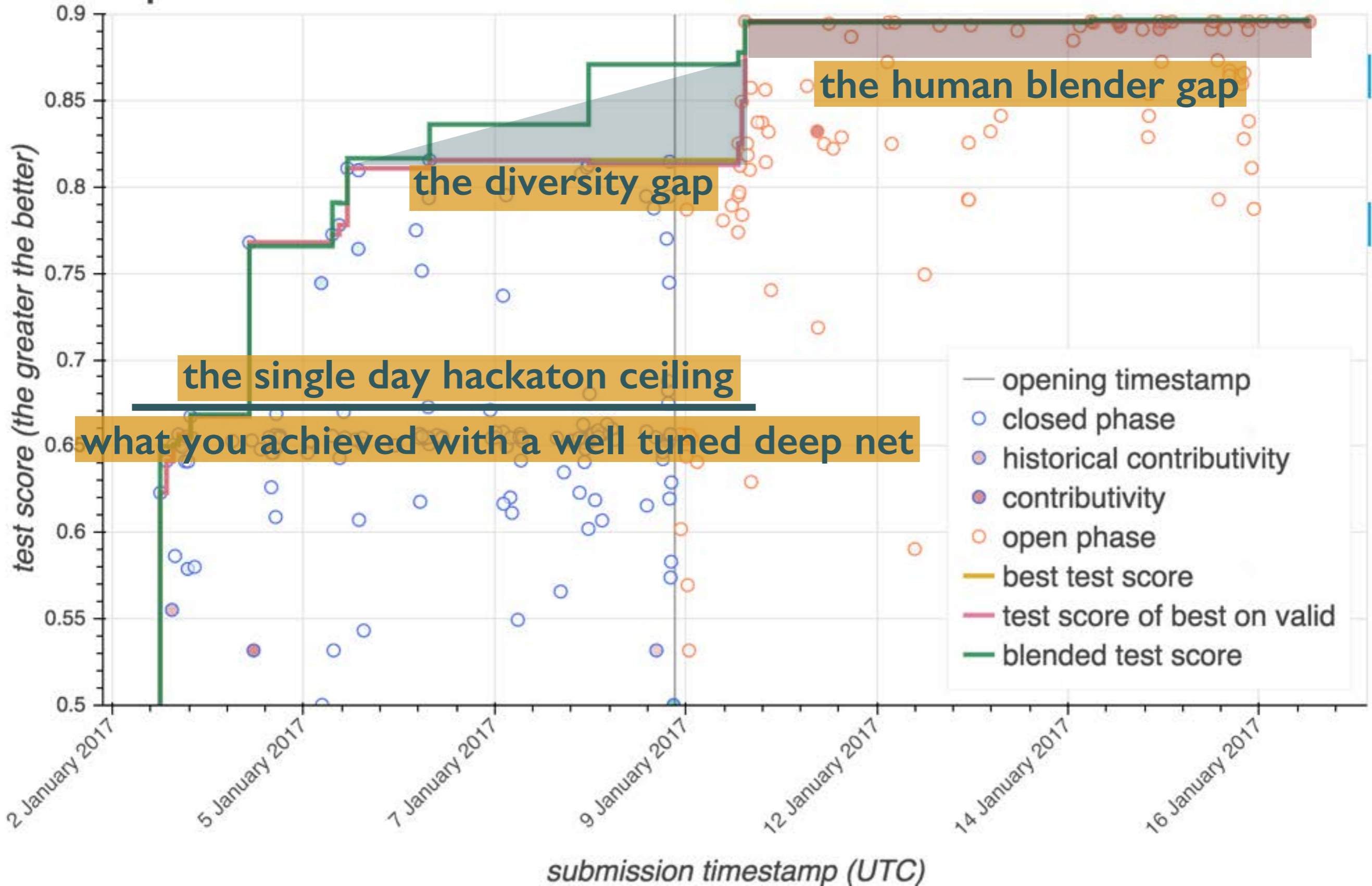


THE DYNAMICS OF COLLABORATION

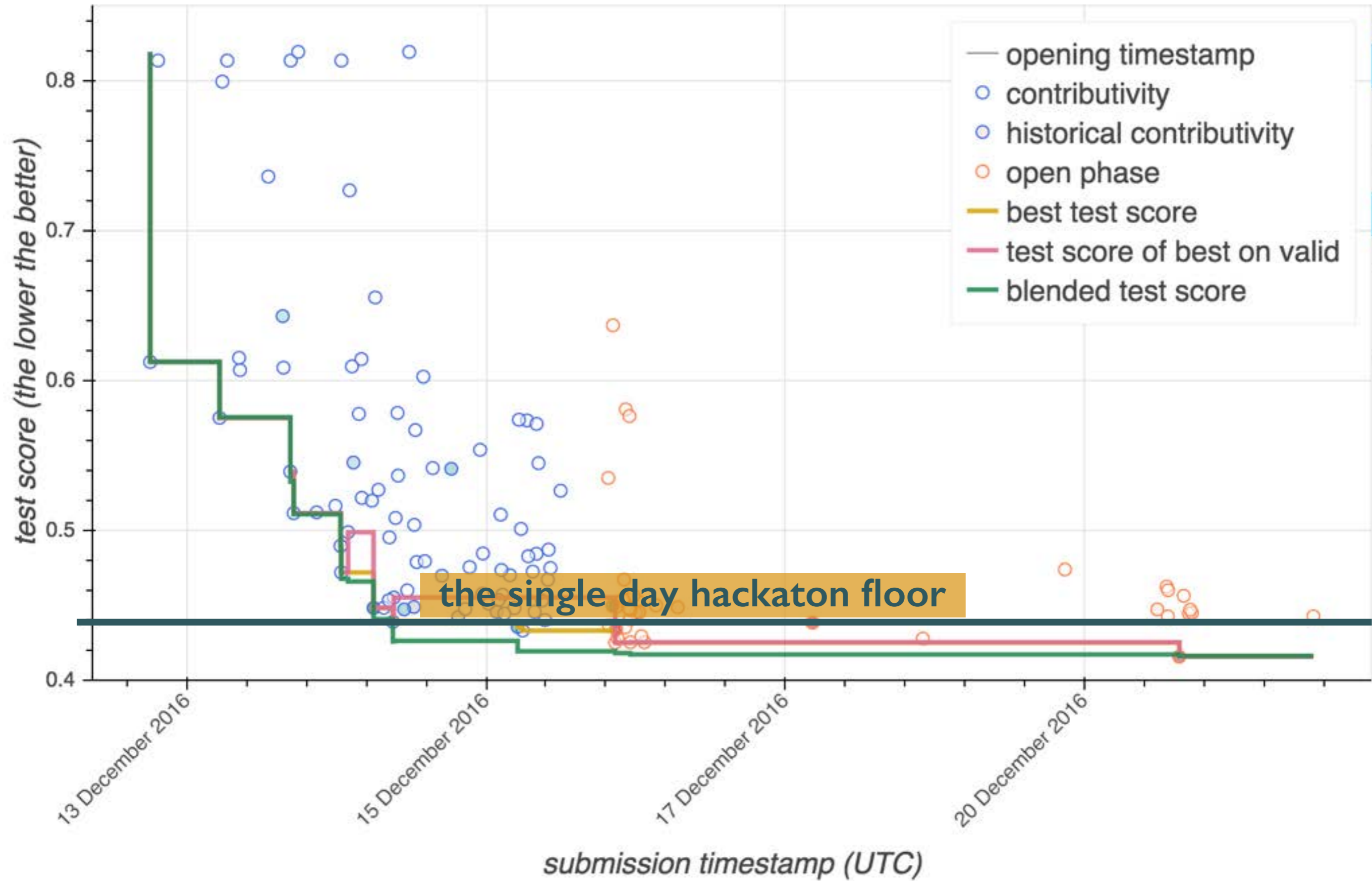
Hep detector anomalies submissions



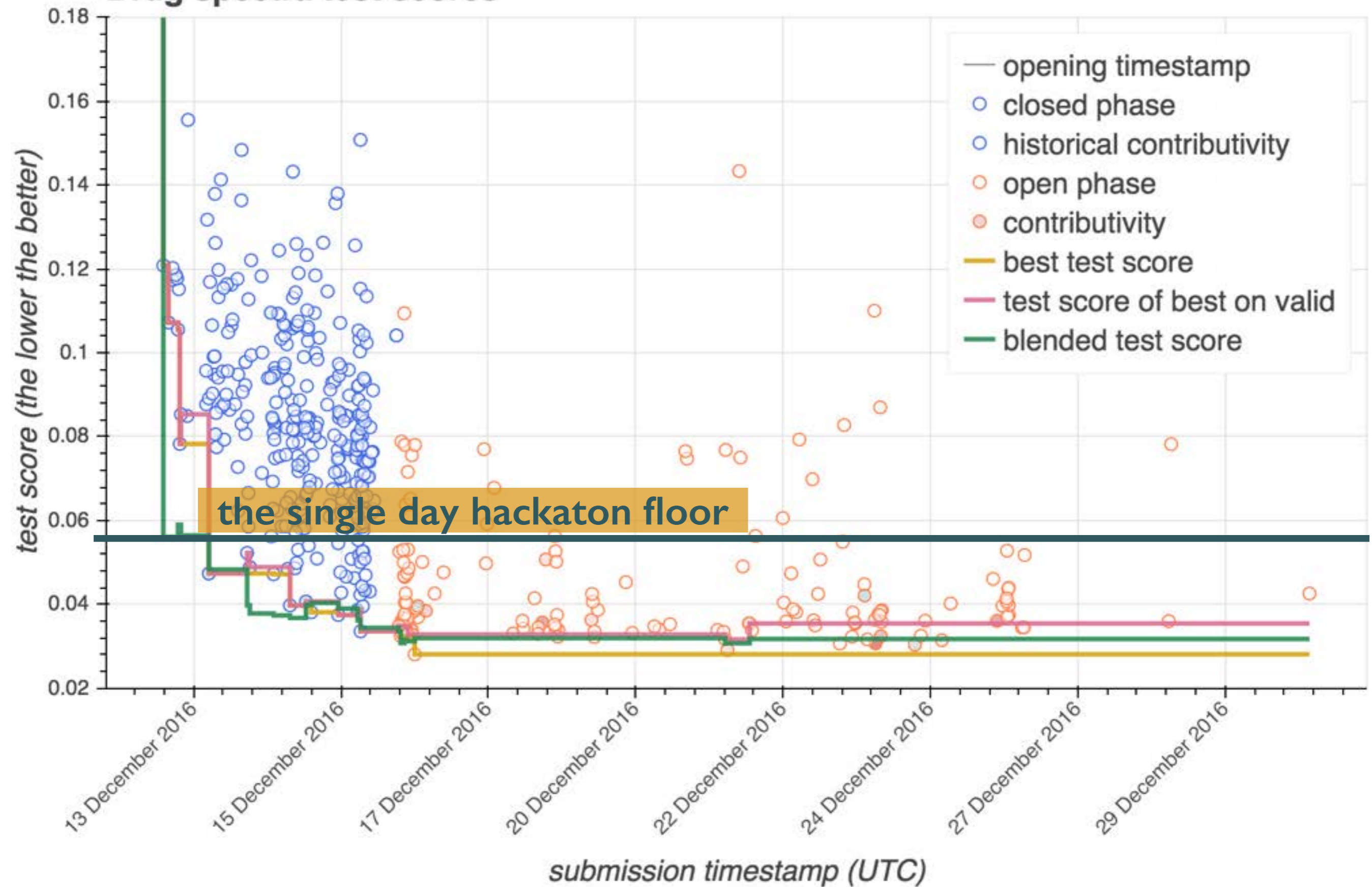
Hep detector anomalies test scores



El nino forecast test scores



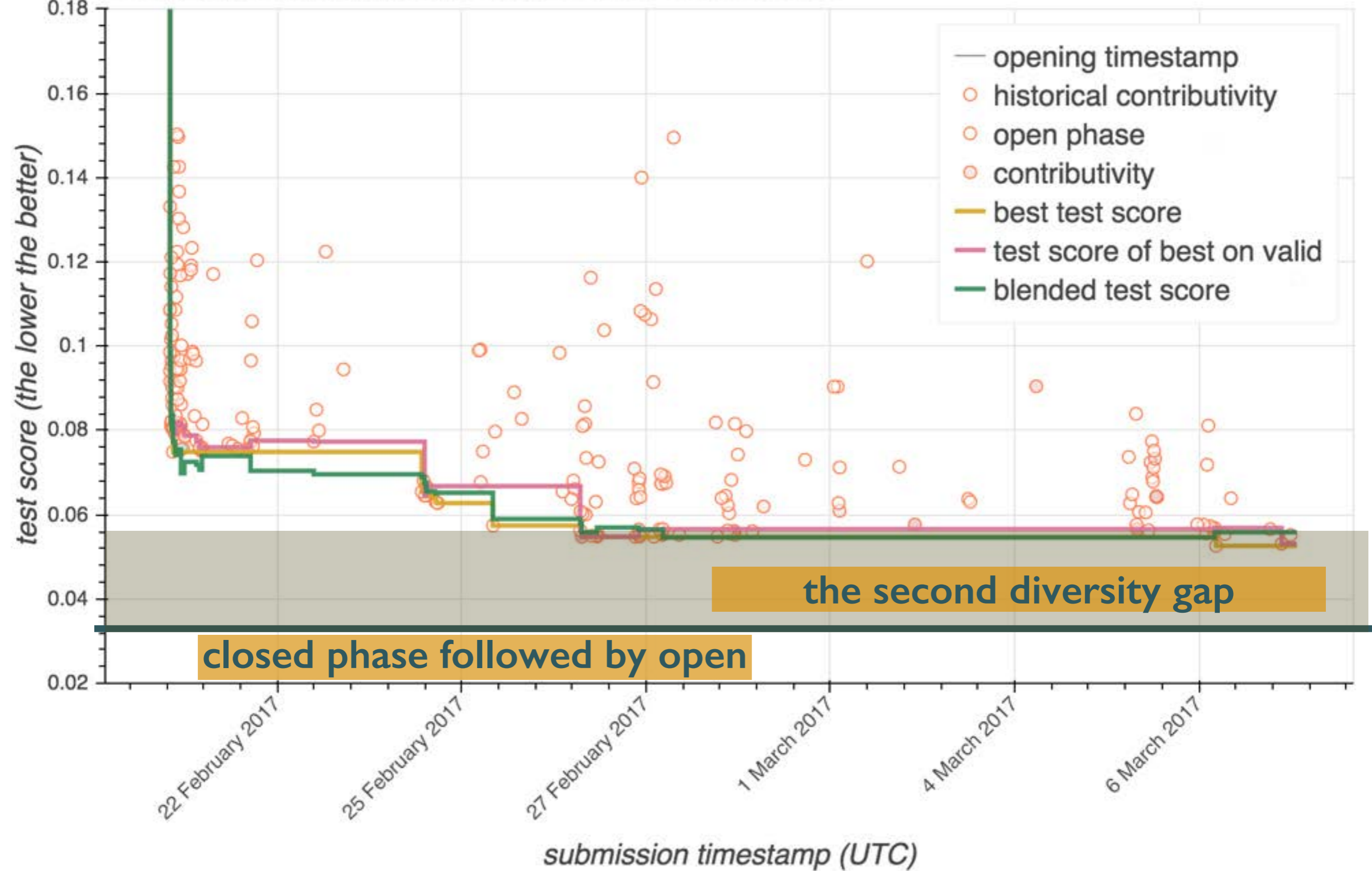
Drug spectra test scores



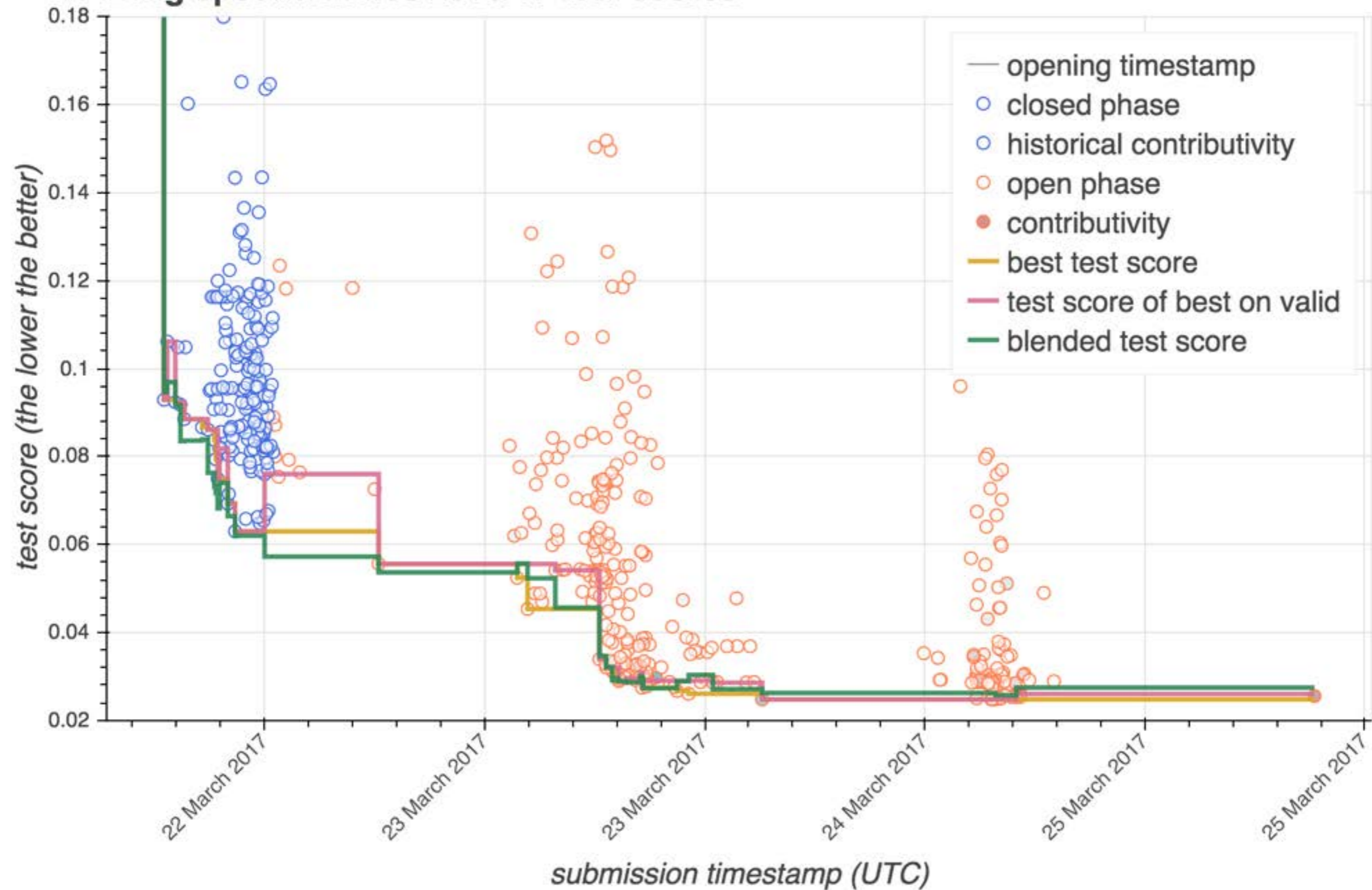
WHAT WE LEARNED

- **Open phase helps novice participants to catch up: the goal of teaching!**
 - Sometimes also makes the best and blended score better
- **Human blending often beats machine blending**
- **Human feature engineering easily beats deep learning on some data**
- **Course RAMPs beat single day hackatons significantly**
 - larger number of students?
 - longer RAMPs?
 - novice and master-level students are better than data science researchers?
 - stronger incentives?
 - closed phase preceding an open phase (vs pure open RAMP) helps to create diversity?

Drug_spectra_m1xmap583_201617 test scores



Drug spectra mines 2016/17 test scores



Classifying and quantifying monoclonal antibody preparations for cancer therapy using machine learning

Laetitia Le ^{ab}, Camille Marini ^{ce}, Alexandre Gramfort ^{cfg},
David Nguyen ^a, Mehdi Cherti ^{ch}, Sana Tfaili ^b, Ali
Tfayli ^b, Arlette Baillet-Guffroy ^b, Eric Caudron ^{ab}, Balázs
Kégl ^{ch}

^a European Georges Pompidou Hospital (AP-HP), Pharmacy department, Paris, France

^b Lip(Sys) Chimie Analytique Pharmaceutique, Univ. Paris-Sud, Université Paris Saclay, F92290 Chatenay-Malabry, France (EA4041 Groupe de Chimie Analytique de Paris Sud)

^c Center of Data Science, Université Paris-Saclay

^d Université Paris-Sud

^e CMAP, Ecole Polytechnique, Palaiseau, France

^f INRIA, Parietal team, Saclay, France

^g LTCI, Télécom ParisTech

^h LAL, CNRS, France

WHAT'S NEXT

- **More RAMPs**
 - **galaxy** morphology, detecting **autism** from brain fMRI, detecting **Mars craters**, forecasting **space weather** (solar storm early warning)
- **More courses**
 - >1000 students next year
- **Build your own RAMPs**

WHAT IS IN IT FOR YOU

- If you are a **data science teacher**
 - we have been using **classroom RAMPs in different formats**: homework, final project, data camp
 - students **love it and work** their butt off, they **learn from each other** and **collaborate**
- If you are a **domain science researcher**
 - We can **solve your predictive problems** better than any single researcher in a classical project
- If you are a **data science researcher**
 - You can **benchmark your new techniques** on a variety of problems

sign up: www.ramp.studio

contact us: balazs.kegl@gmail.com