



Be a part of this talk:

- Log in / create an account on www.openml.org
 - You also need a GitHub account
- Click the  or  icon
- Click 'Launch Demo'



OpenML

DEMOCRATIZING AND AUTOMATING
MACHINE LEARNING

JOAQUIN VANSCHOREN, TU EINDHOVEN, 2016

Research different.

Polymaths: Solve math problems
by massive **online** collaboration

Broadcast question, combine
many minds to solve it

SCIENCE photoL

Networked Science

Serendipity: what's hard for one person is easy for another
Collaboration only scales if **all friction is eliminated**

Easy, organized, access to data, code, and results



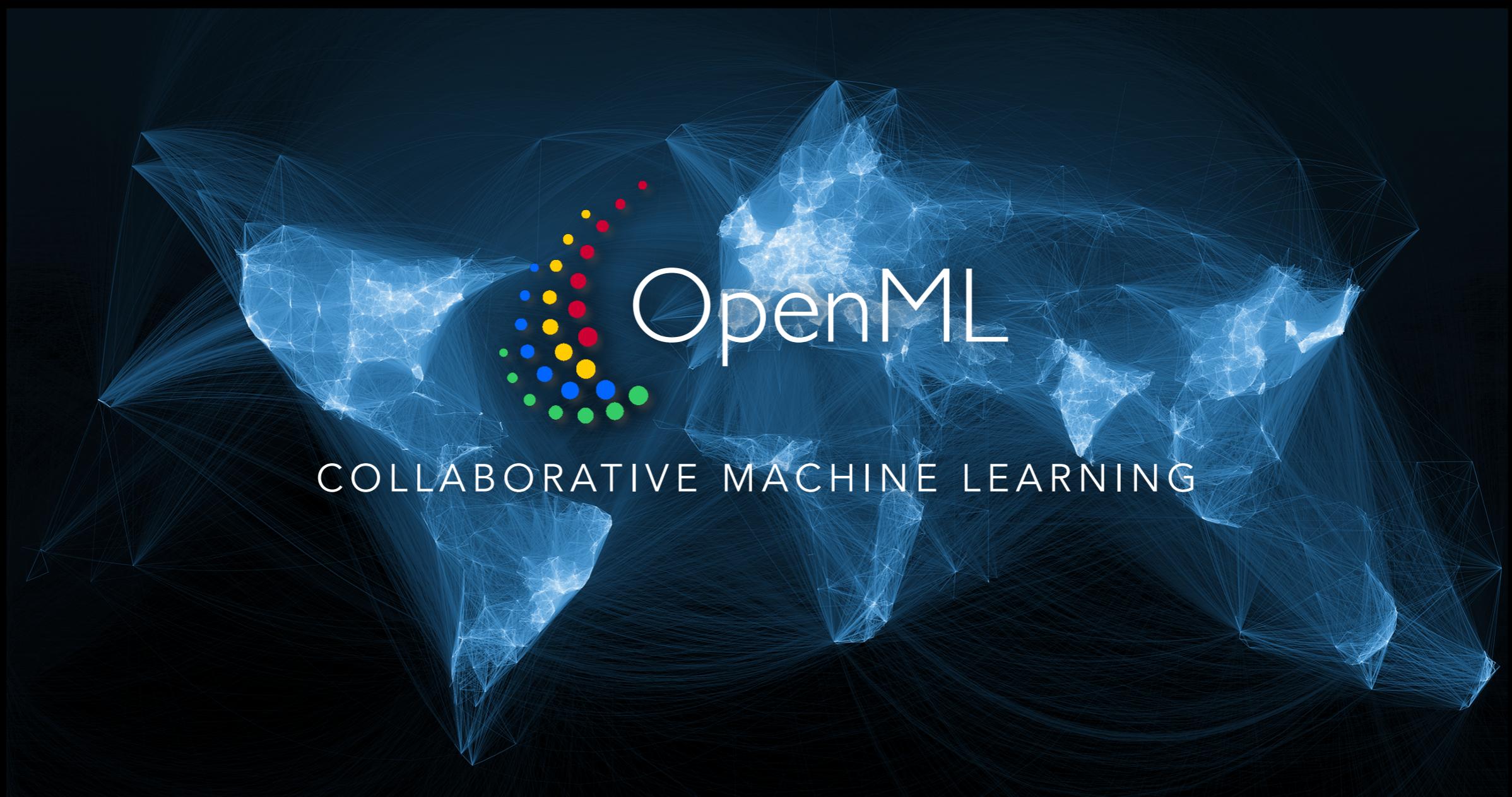
WHAT IF WE CAN EXPLORE DATA
COLLABORATIVELY



WHAT IF WE CAN EXPLORE DATA
COLLABORATIVELY
ON WEB SCALE



WHAT IF WE CAN EXPLORE DATA
COLLABORATIVELY
ON WEB SCALE **IN REAL TIME**



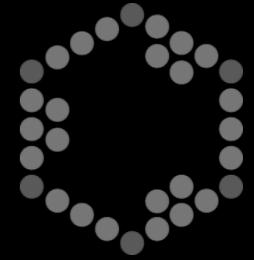
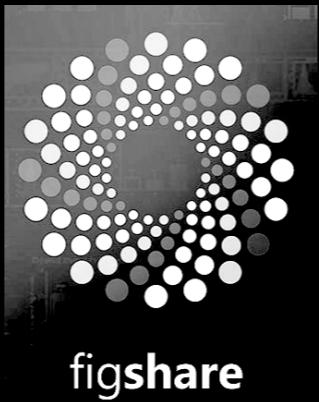
Easy to use: Integrated in many ML environments

Easy to contribute: Automated sharing of data, code, results

Organized data: Reproducible, connected to data, code, people

Reward structure: Build reputation and trust

Self-learning: Learn from millions of experiments to help users



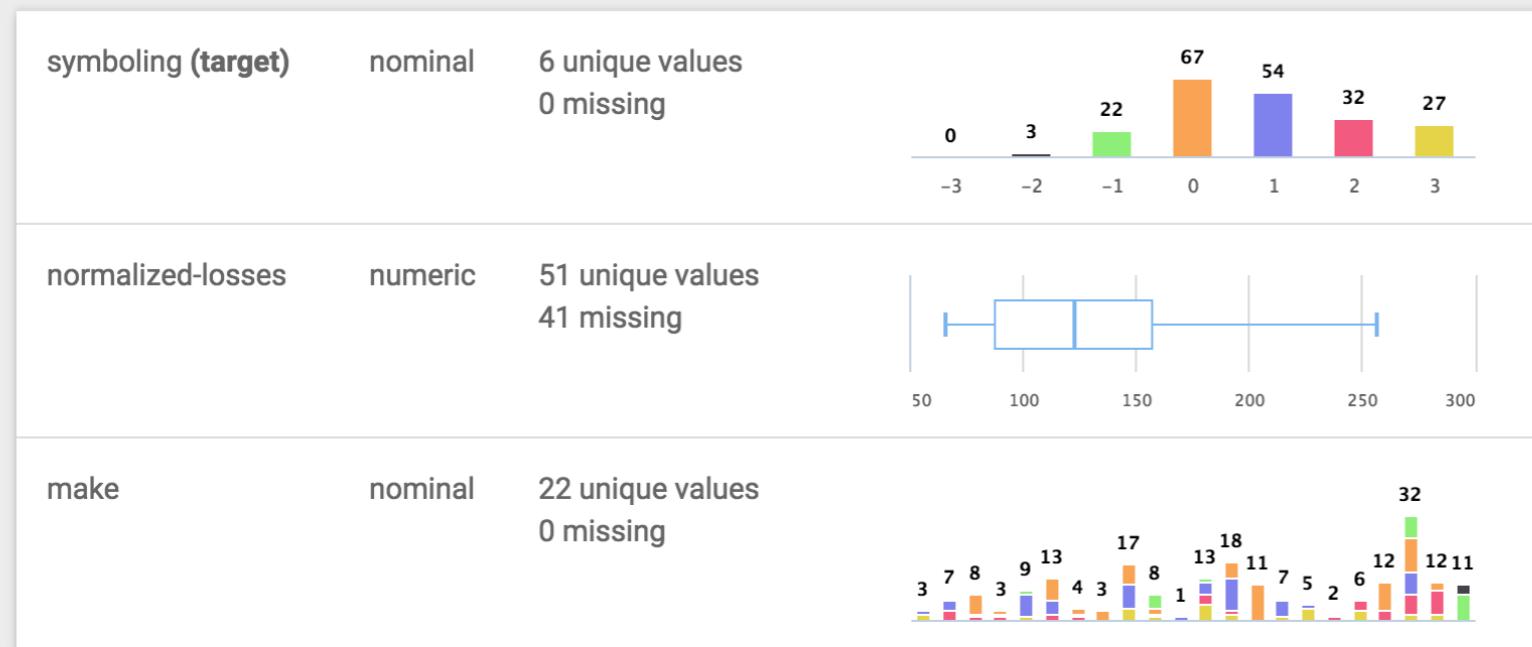
**Data (ARFF) uploaded or referenced, versioned
analyzed, characterized, organized online**



analyzed, characterized, organized online

+ visualizations, statistics, landmarks, error checking,
queryable through website + API

26 features



72 properties

DefaultAccuracy	0.33	The predictive accuracy of the model.
NumberOfClasses	7	The number of classes in the target variable.
NumberOfFeatures	26	The number of features in the dataset.
NumberOfInstances	205	The number of instances in the dataset.
NumberOfMissingValues	59	Counts the total number of missing values in the dataset.
NumberOfUniqueValues	17	The number of unique values in the target variable.

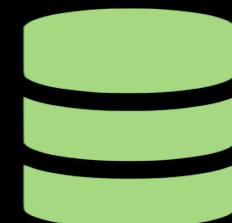


Tasks contain data, goals, procedures.
Readable by tools, automates experimentation
All results organized online: **realtime overview**

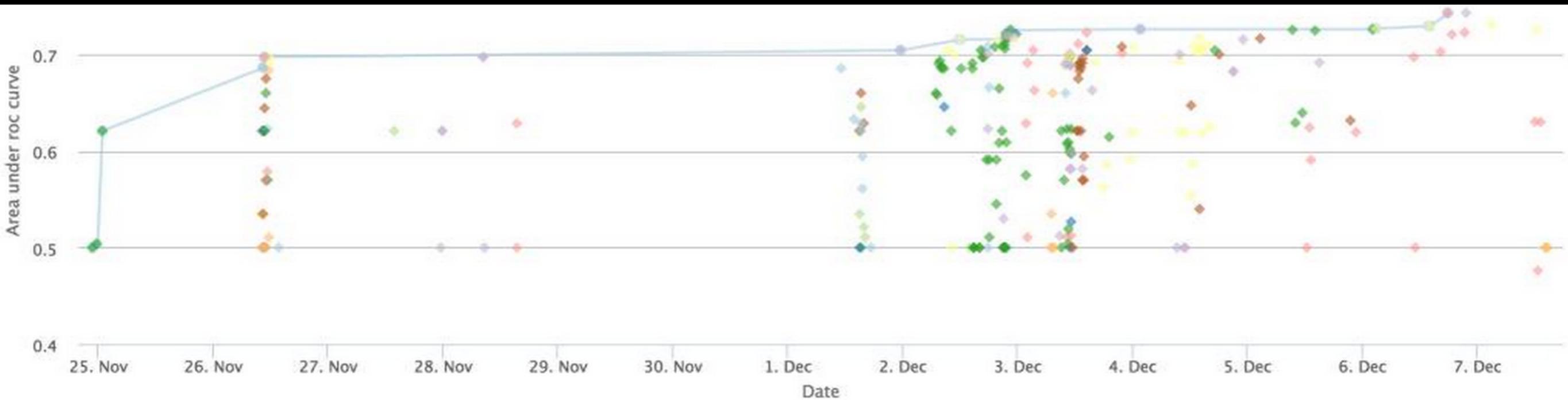


Train-test
splits

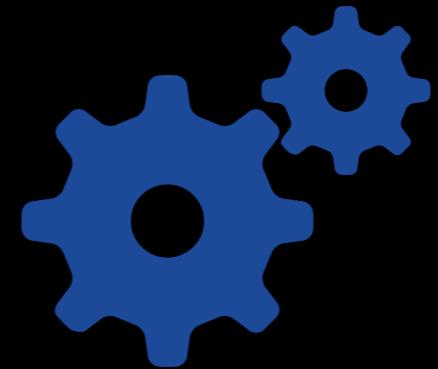
Classify target X



All results organized online: **realtime overview**

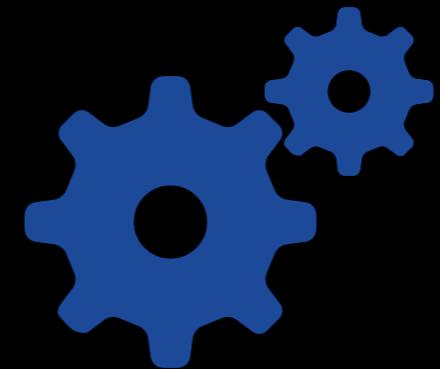


frontier Joaquin Vanschoren Perry van Wesel Jose Melo Jos Mangnus Daan Peters Tom Becht Kevin Jacobs Koen Engelen
Olav Bunte Stephan Oostveen Roy van den Hurk Sylwester Kogowski Ky-Anh Tran Edgar Salas Thomas Tiel Groenestege
Jorn Engelbart Mathijs van Liemt Henry He Richie Brondenstein Hugo Spee Stanley Clark Christoforos Boukouvalas Rogier Beckers
Stefan Majoor

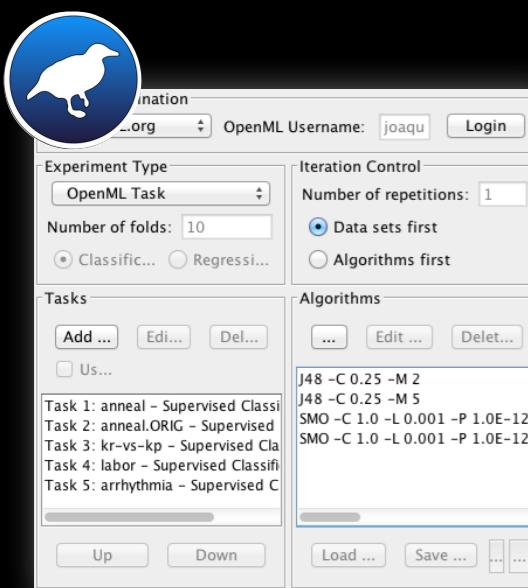


Flows (code) run anywhere, using your favorite tools

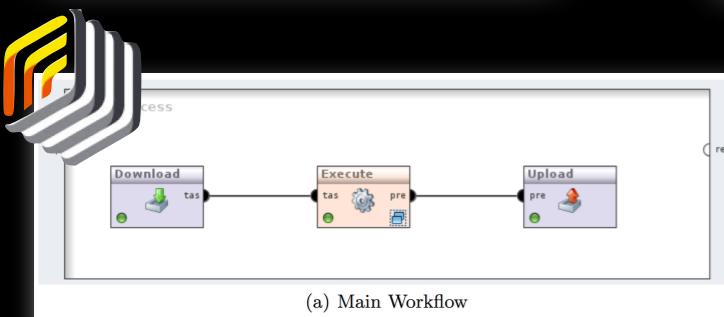
Integrations + APIs (REST, R, Python, Java,...)



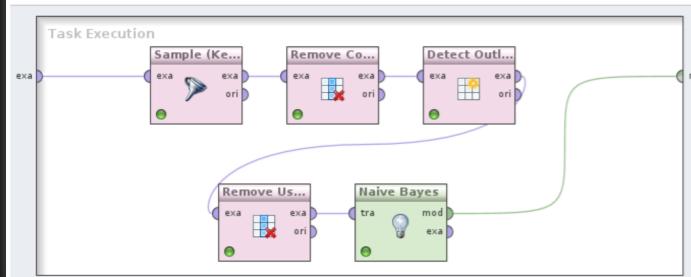
Integrations + APIs (REST, R, Python, Java,...)



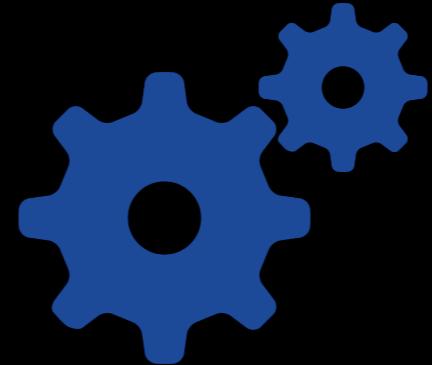
```
from sklearn import tree
from openml import tasks, runs
task = tasks.get_task(14951)
clf = tree.DecisionTreeClassifier()
run = runs.run_task(task, clf)
return_code, response = run.publish()
```



(a) Main Workflow



```
library(OpenML)
library(mlr)
task = getOMLTask(10)
lrn = makeLearner("classif.rpart")
res = runTaskMlr(task, lrn)
run.id = uploadOMLRun(res)
```



Integrations + APIs (REST, R, Python, Java,...)



```
from sklearn import tree
from openml import tasks, runs
task = tasks.get_task(14951)
clf = tree.DecisionTreeClassifier()
run = runs.run_task(task, clf)
return_code, response = run.publish()
```





Experiments auto-uploaded, evaluated online
reproducible, linked to **data, flows, authors**
and **all other experiments**



Experiments auto-uploaded, evaluated online

Result files



Description

XML file describing the run, including user-defined evaluation measures.



Model readable

A human-readable description of the model that was built.



Model serialized

A serialized description of the model that can be read by the tool that generated it.



Predictions

ARFF file with instance-level predictions generated by the model.

Area under ROC curve

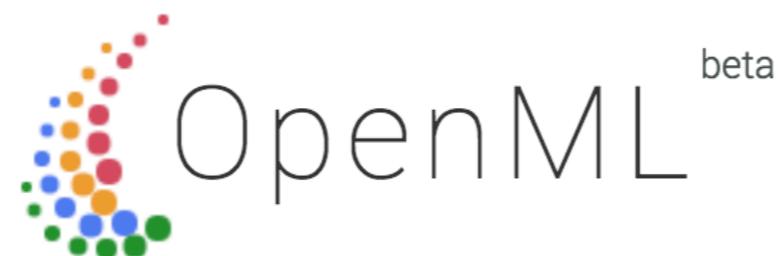
0.7007 \pm 0.0023

Per class

0	1
0.7007	0.7007

Cross-validation details (10-fold Crossvalidation)





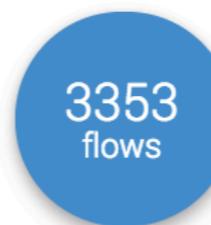
Exploring machine learning better, together



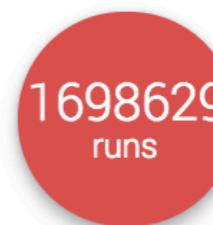
Find or add **data** to analyse



Download or create scientific
tasks

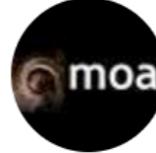


Find or add data analysis **flows**



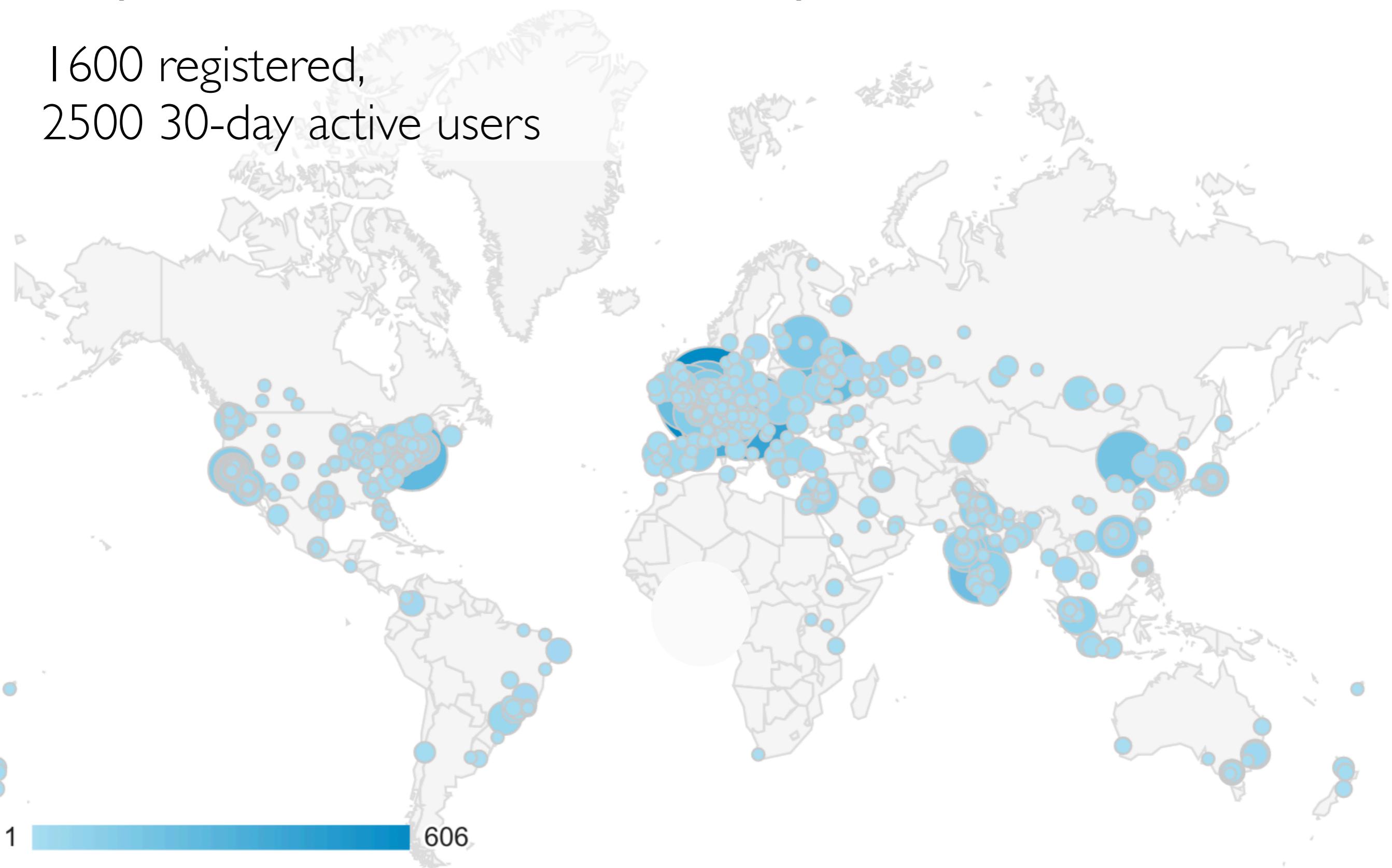
Upload and explore all **results**
online.

Download and share data, flows and runs through:



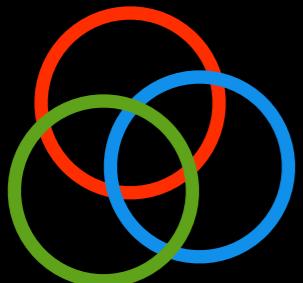
OpenML Community

| 600 registered,
2500 30-day active users



Jul-Nov 2016

Collaboration tools (in progress)



Circles

Create collaborations with trusted researchers



Studies (e-papers)

Online counterpart of a paper, linkable



Reputation

Auto-tracking of your activity, reach, impact

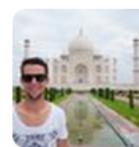


Notebooks

Easy online collaboration on data analysis scripts

Classroom challenges

Results a
Students
Simple re
Hidden t



Rogier Beckers

@RogierBeckers



+ Follow

Het bewijs dat ik studeer op zondag!

“@joavanschoren: #Machinelearning students on a #collaborative data mining ”

[View translation](#)

Lauradorp, Landgraaf



...

Contributions over time

every point is a run, click for details



RETWEETS

2

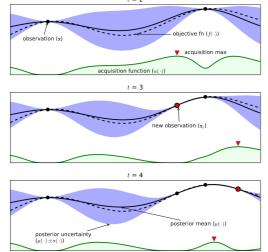
FAVORITES

2



9:48 PM - 7 Dec 2014

Upcoming realtime challenges



Hyperparameter optimization challenge

Every iteration can be uploaded, prior ones downloaded



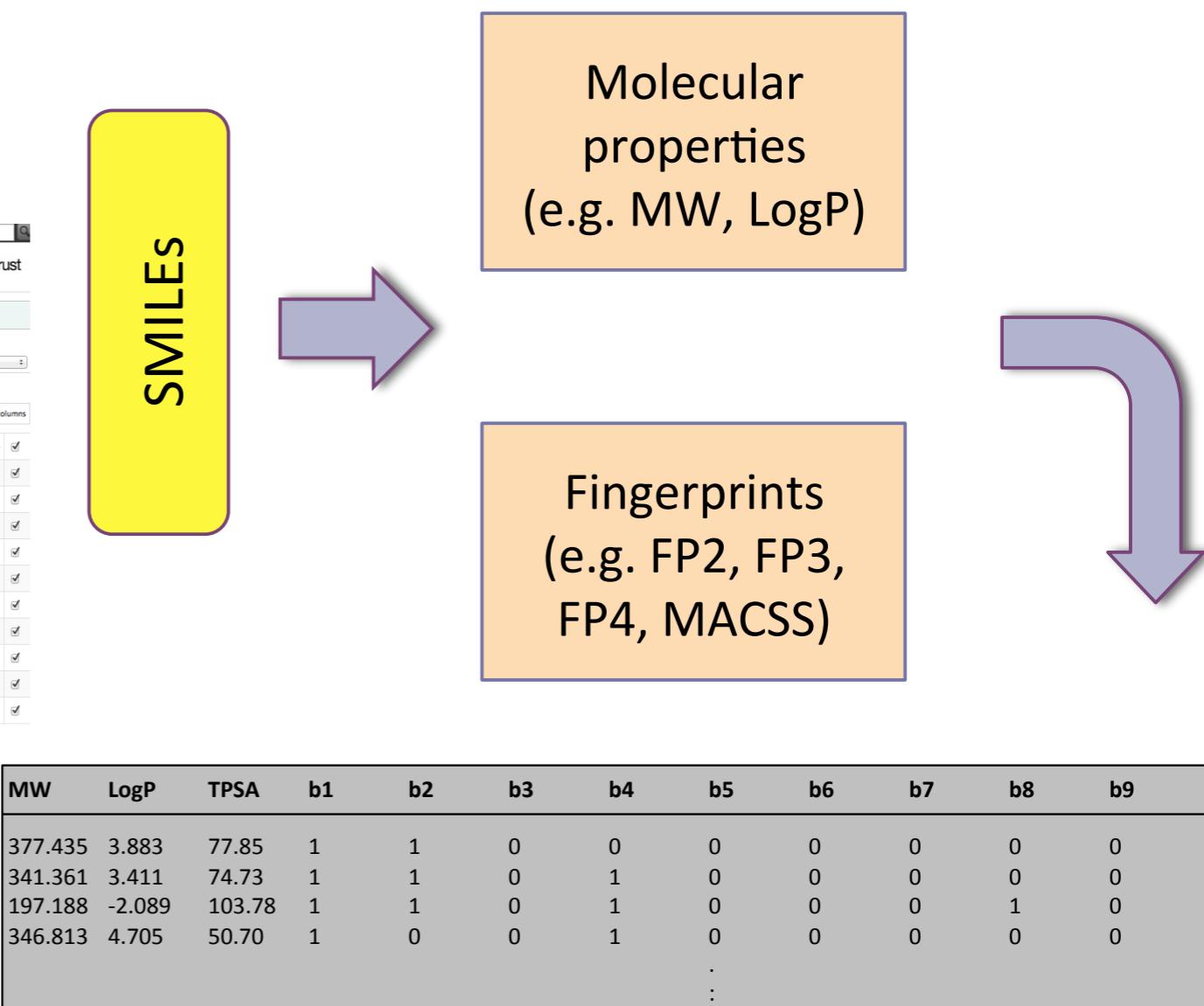
Predict energy usage in smart energy
API streams data, predictions uploaded hourly/daily

OpenML in drug discovery

Predict which drugs will inhibit certain proteins (and hence viruses, parasites,...)

The screenshot shows two pages from the ChEMBL database. The top page is a 'Target Report Card' for CHEMBL3227, which is Metabotropic glutamate receptor 5. It includes sections for Target Name and Classification, Target Components, and ChEMBL Statistics. The bottom page is a search results page for 'Metabotropic glutamate receptor 5', showing 23 entries with columns for ChEMBL ID, Preferred Name, UniProt Accession, Target Type, Organism, Compounds, and Bioactivities.

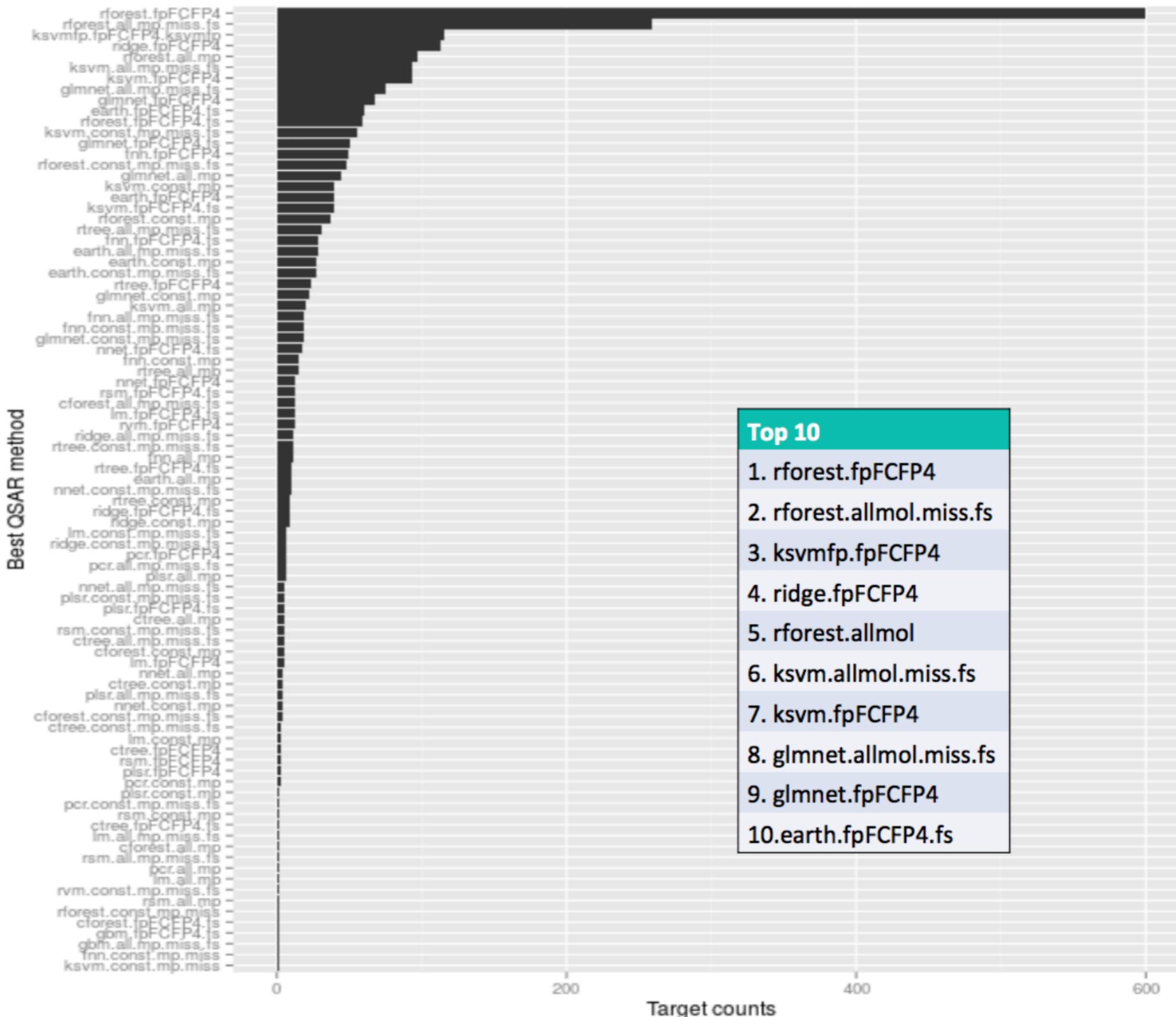
ChEMBL database
1.4M compounds, 10k proteins,
12.8M activities



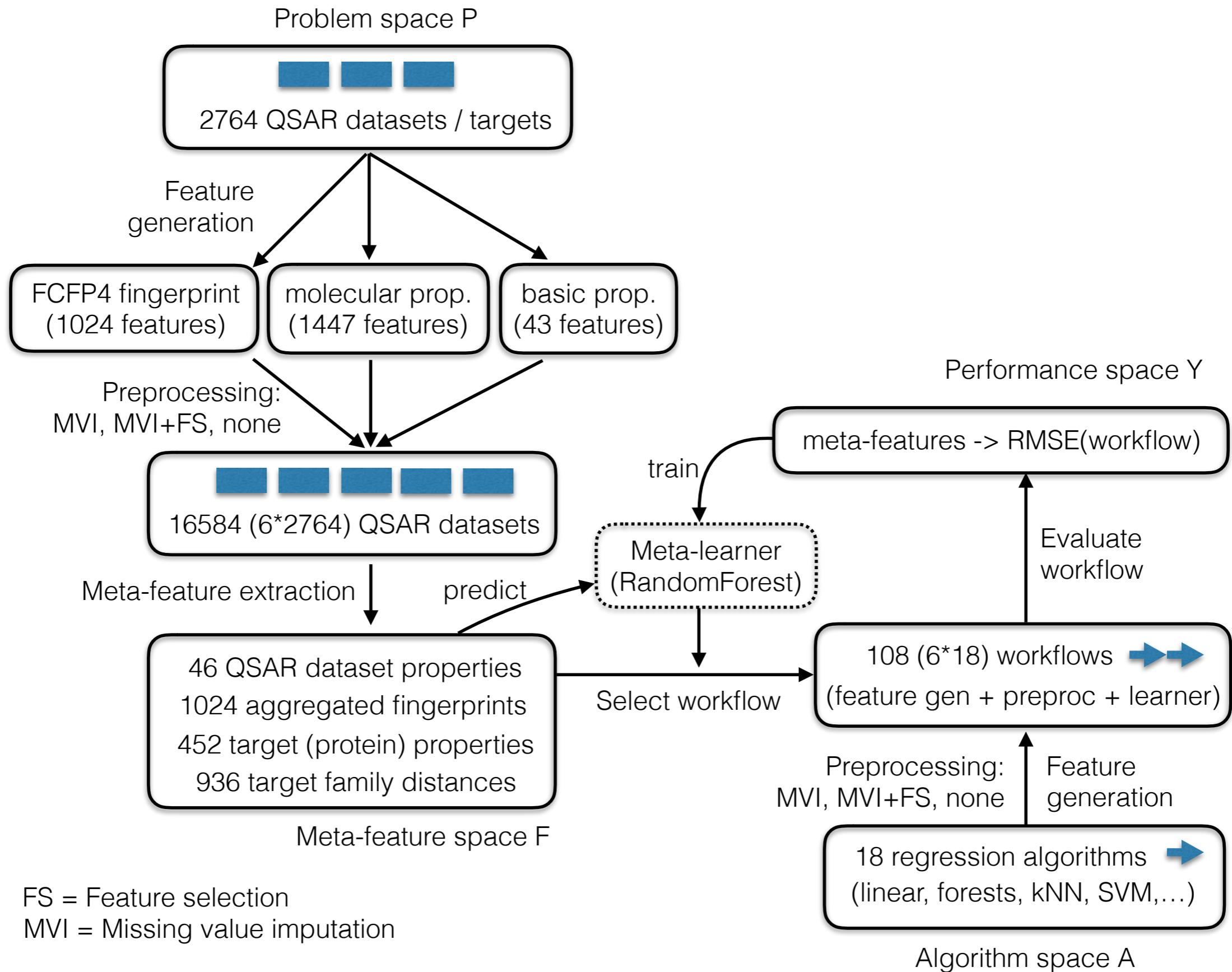
16.000+ QSAR datasets
2750 targets (proteins), x 6 feature representations

OpenML in drug discovery

Best algorithms?
Best features?



Predicting workflows with MetaLearning



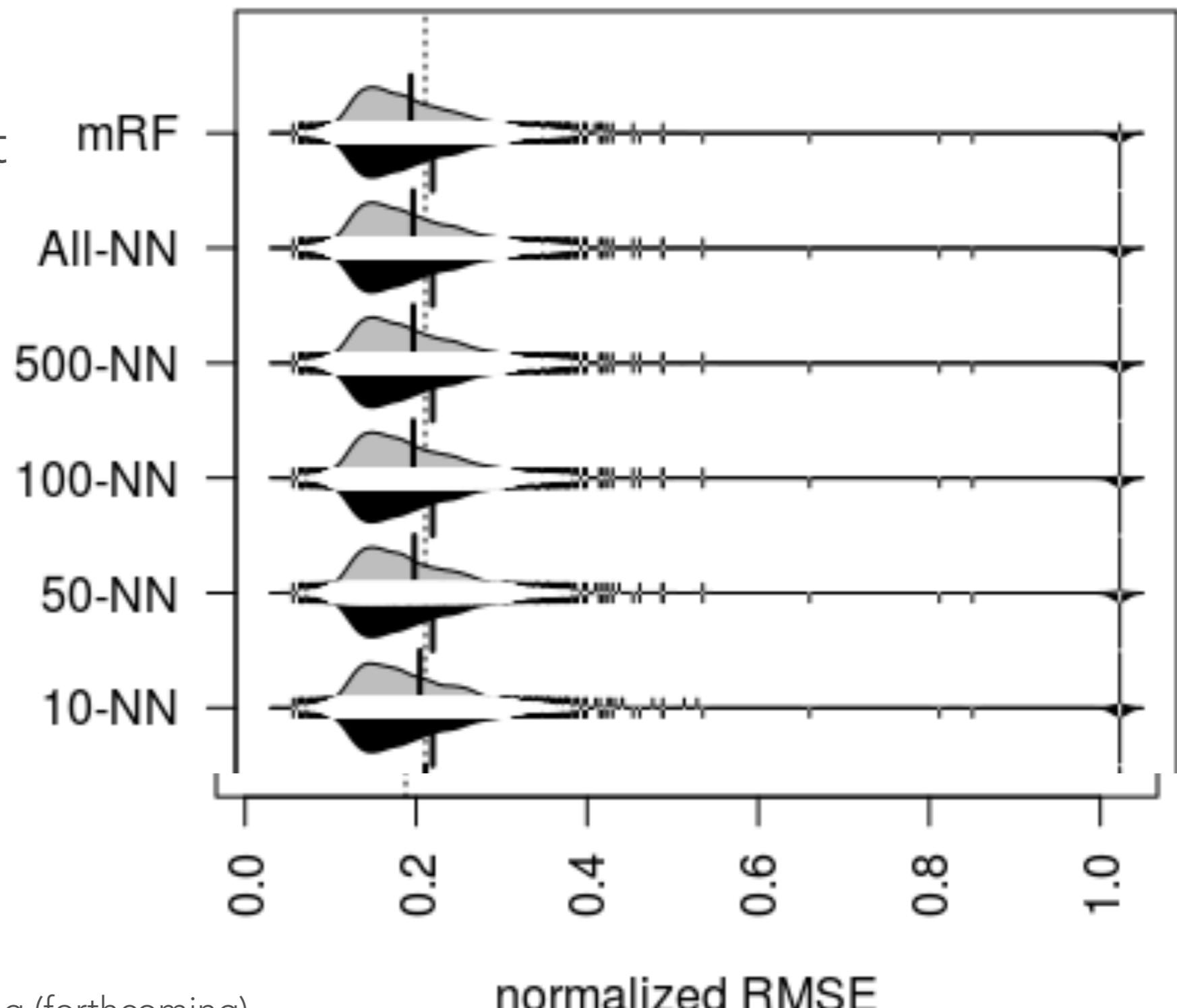
OpenML in drug discovery

Random forest (black) vs meta-learner (grey)

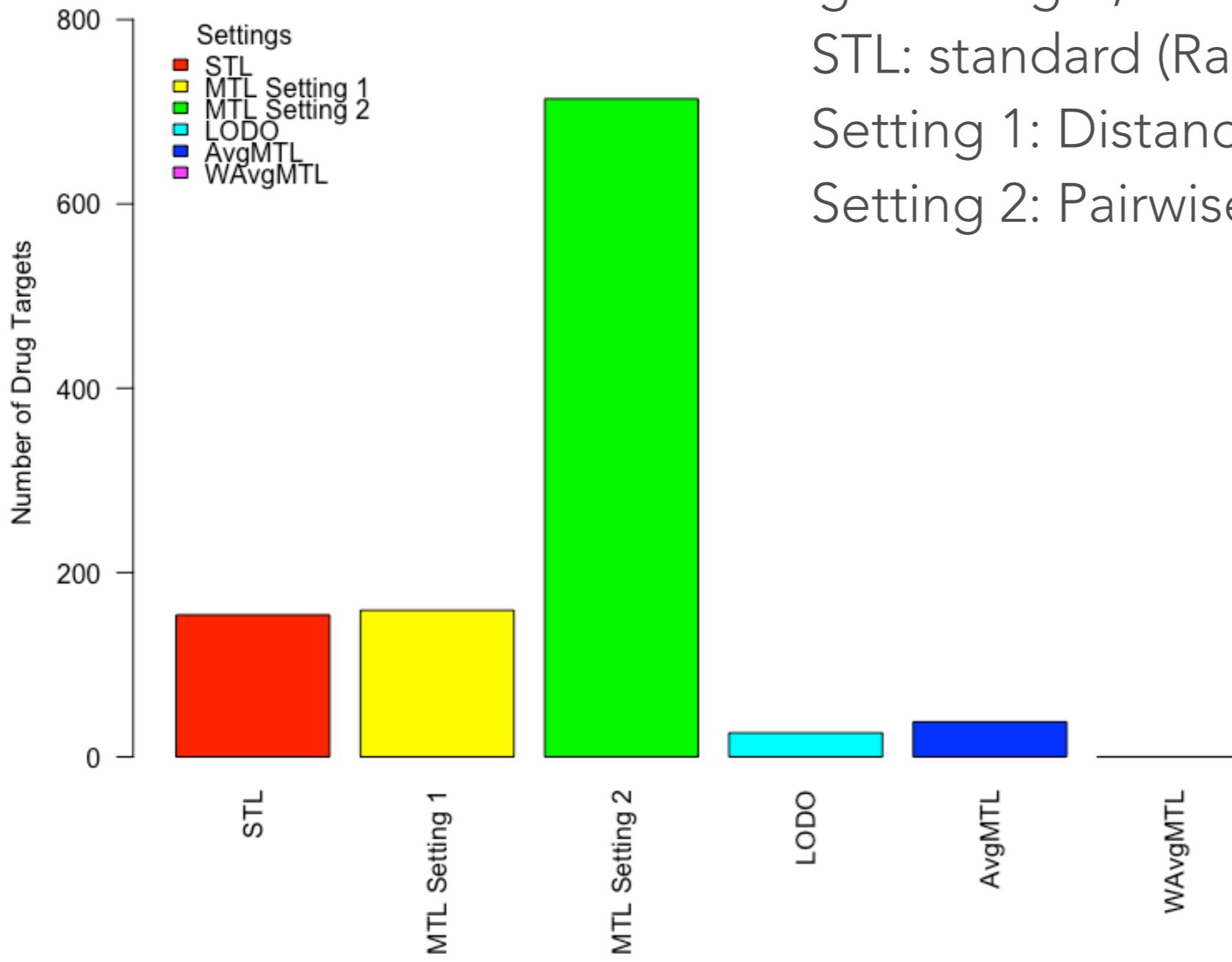
Meta-learner:

mRF: Random forest

k-NN: kNN



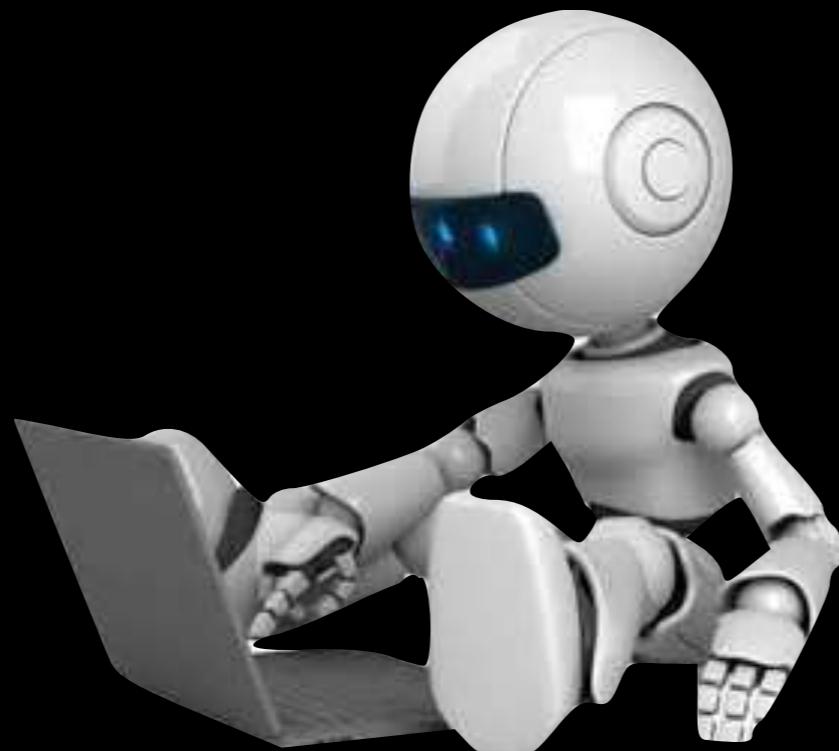
OpenML in drug discovery



Multi-target learning: if few drugs are tested on a given target, include data on 'related' targets:
STL: standard (Random Forests)
Setting 1: Distance is taxonomy (ChEMBL)
Setting 2: Pairwise sequence alignment

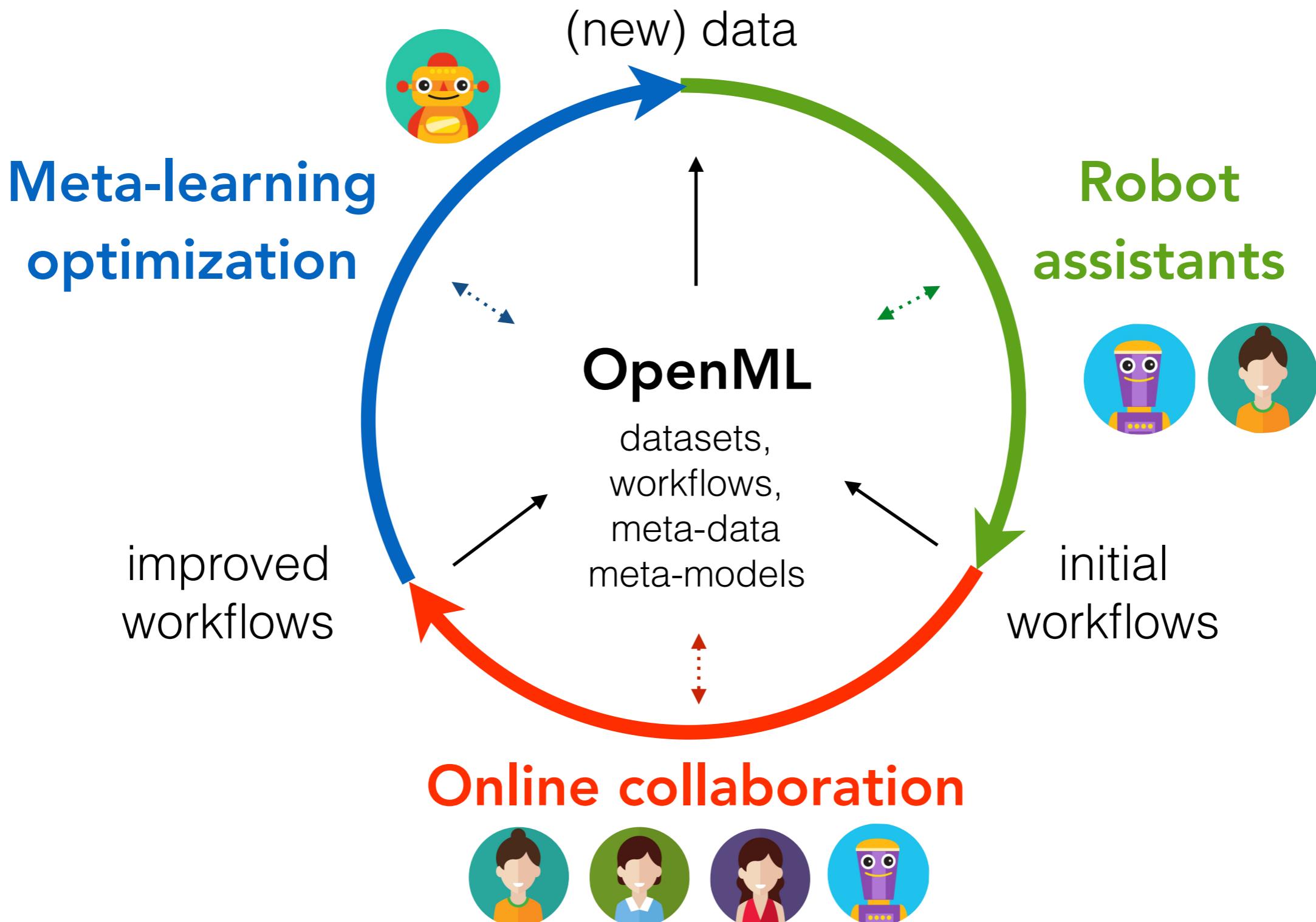
We just scratched the surface. All data is available on OpenML.

Automating machine learning



- Lots of experimental data: learn to learn across datasets
- Bots that automatically run algorithms on data
- Data-driven modeling: learn how to build models based on many prior experiments

Automating machine learning: a human-robot symbiosis (future)



Robot assistants

Runtime prediction bot: predicts how long an ML algorithm will run on your data



Feature selection bot: recommends/runs feature selection techniques



Imputation bot: recommends/runs missing value imputation techniques



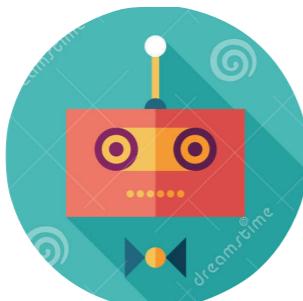
Outlier detection bot: recommends/runs outlier detection techniques

(new) data

initial
workflows

Robot assistants

Random Bot: runs random search given a hyperparameter space



(new) data



Optimization bots: runs advanced hyperparameter optimization

initial
workflows



Workflow bot: build ML workflows, in collaboration with other bots

Random Bot running on OpenML

≡ People

Search



Random Bot

Joined 2016-03-16

Activity

5716

Reach

17

Impact

0

Uploads

0

17

0

1740

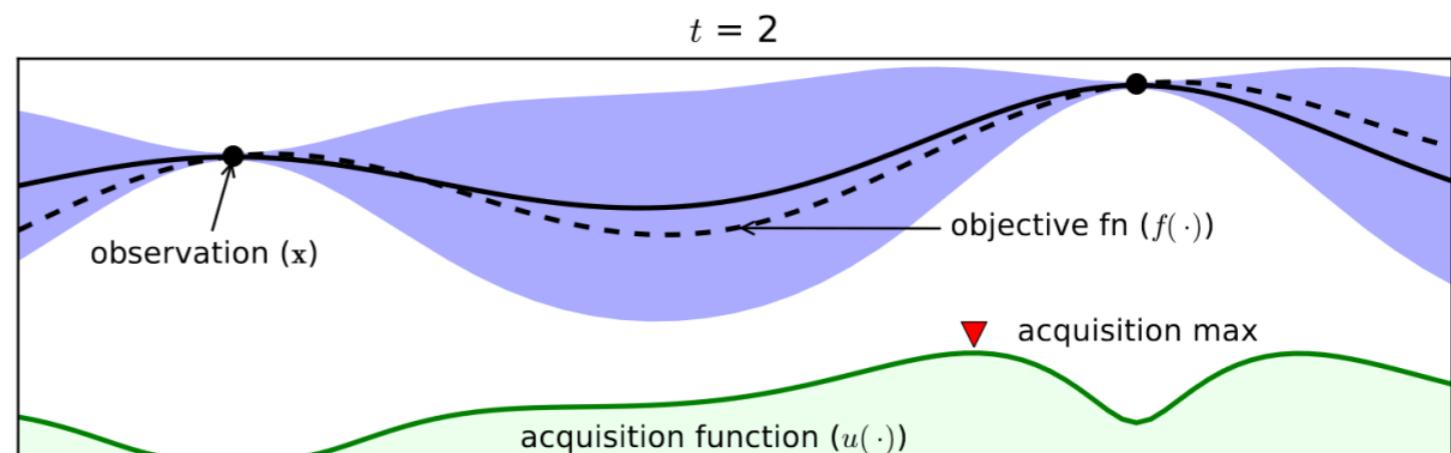
EDIT PROFILE

	Activity	Reach	Impact
Data Sets	0	0	0
Flows	17	0	0
Tasks	0	0	0
Runs	1740	0	0

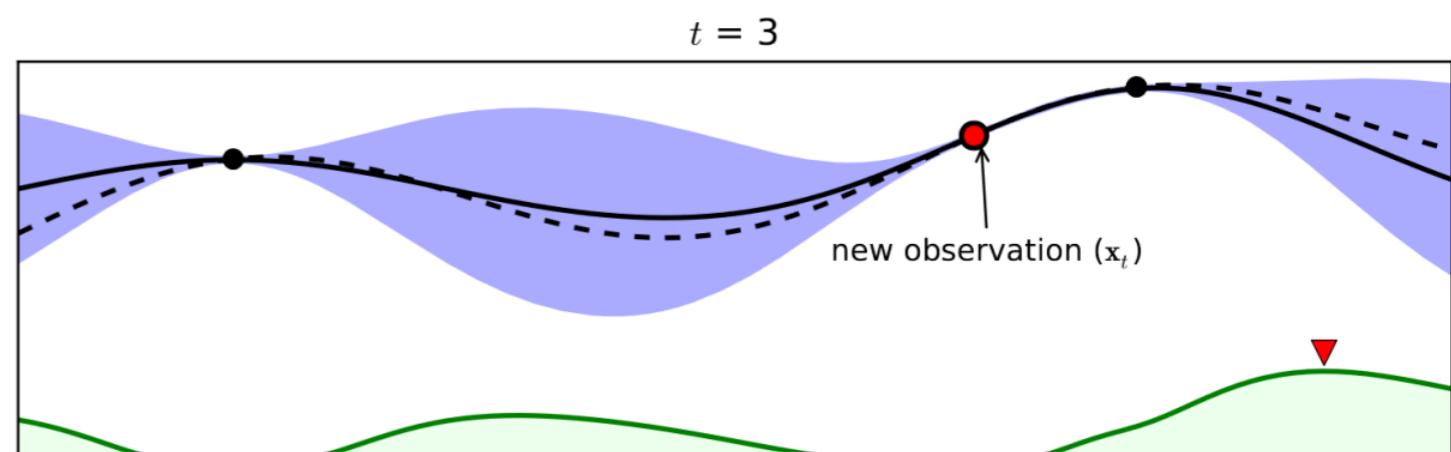
Bayesian optimization + metalearning

Include meta-data from prior datasets

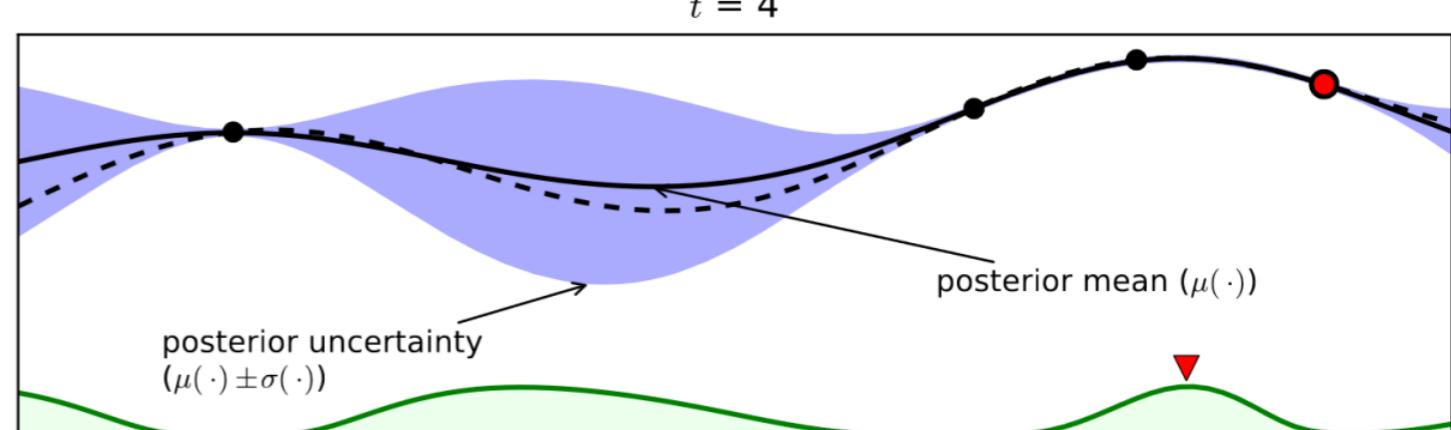
- Warm start: initialize search with promising configurations (*AutoML challenge winner*)



- Surrogate models with prior (focus on best parameters, ranges)

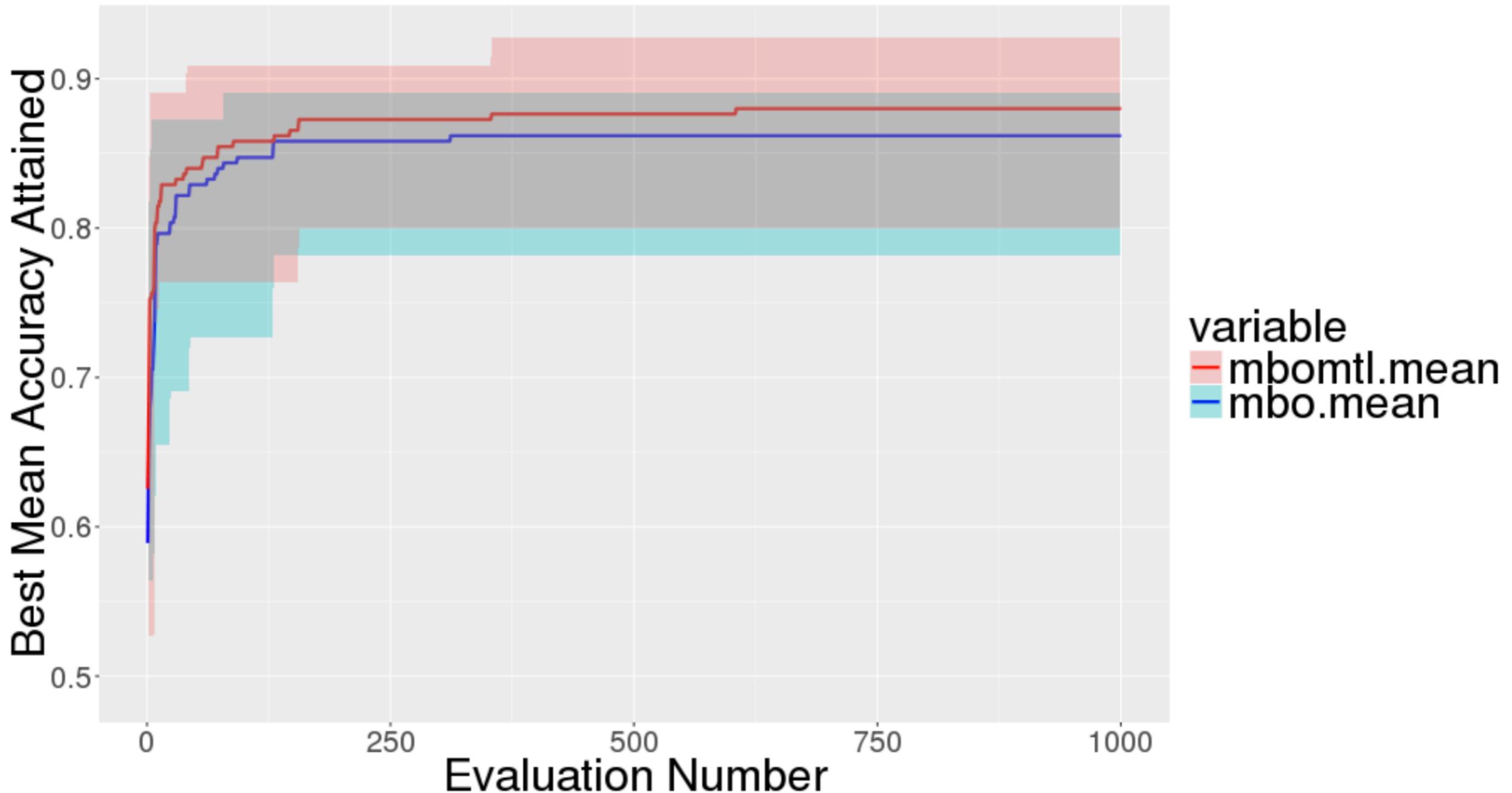


- Acquisition functions based on meta-models (predict performance, trained on prior datasets)

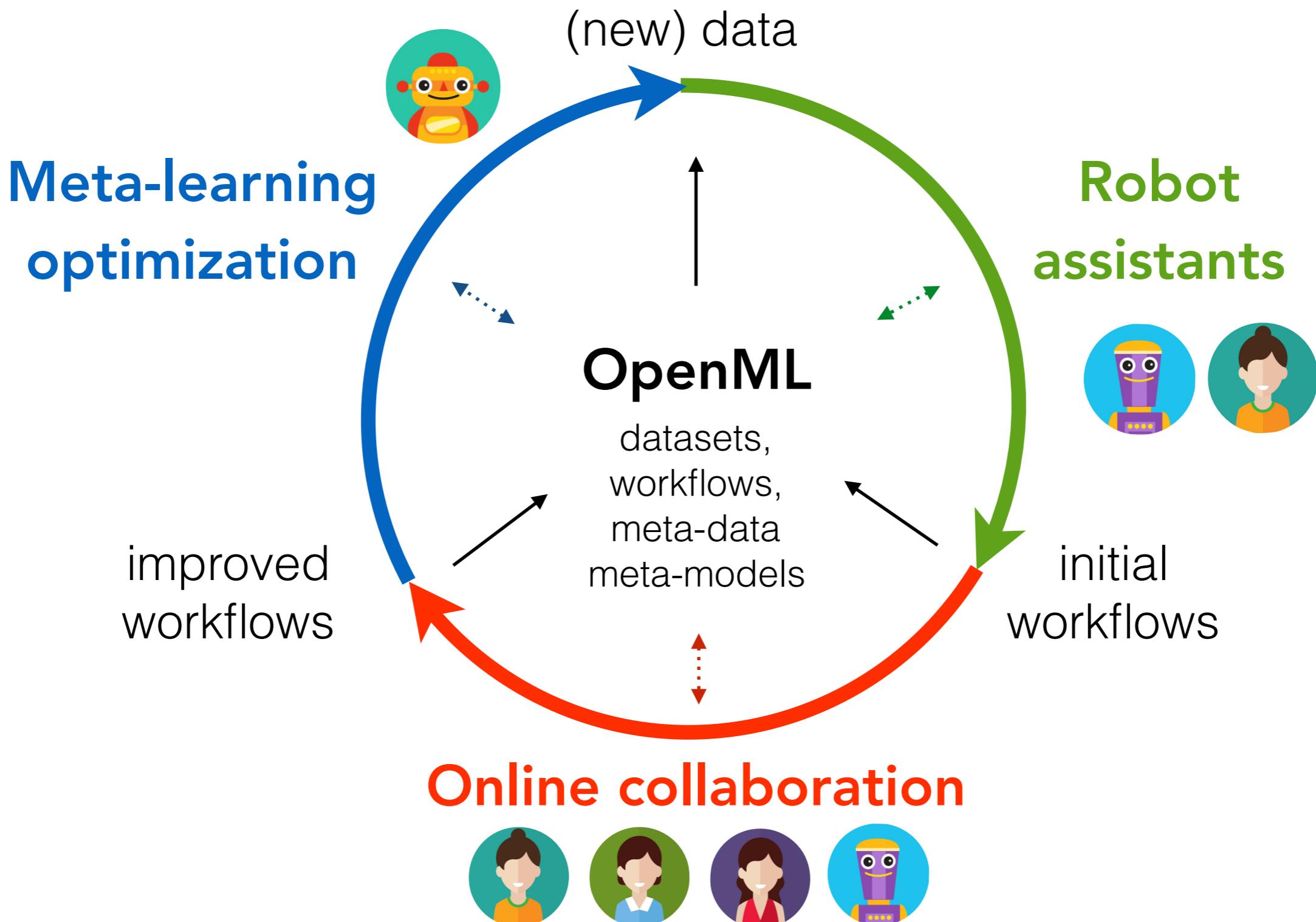


Combine meta-learning and optimization

- Acquisition functions based on meta-models



Learn from results of humans and robots,
use that to build better bots to help humans





WHAT IF WE CAN EXPLORE DATA
COLLABORATIVELY
ON WEB SCALE **IN REAL TIME**



Join Us!
www.openml.org
Join our hackathons

@open_ml
OpenML



Thank You

