# Meeting of the LF AI & Data Technical Advisory Council (TAC)

August 24, 2023

**□LF** AI & DATA

# Antitrust Policy

› Linux Foundation meetings involve participation by industry competitors, and it is the intention of the Linux Foundation to conduct all of its activities in accordance with applicable antitrust and competition laws. It is therefore extremely important that attendees adhere to meeting agendas, and be aware of, and not participate in, any activities that are prohibited under applicable US state, federal or foreign antitrust and competition laws.

› Examples of types of actions that are prohibited at Linux Foundation meetings and in connection with Linux Foundation activities are described in the Linux Foundation Antitrust Policy available at http://www.linuxfoundation.org/antitrust-policy. If you have questions about these matters, please contact your company counsel, or if you are a member of the Linux Foundation, feel free to contact Andrew Updegrove of the firm of Gesmer Undergone LLP, which provides legal counsel to the Linux Foundation.

**LF** AI & DATA

# Recording of Calls

**Reminder:**

TAC calls are recorded and available for viewing on the TAC Wiki

LF AI & DATA

# Reminder: LF AI & Data Useful Links

› Web site: lfaidata.foundation
› Wiki: wiki.lfaidata.foundation
› GitHub: github.com/lfaidata
› Landscape: https://landscape.lfaidata.foundation or https://l.lfaidata.foundation
› Mail Lists: https://lists.lfaidata.foundation
› Slack: https://slack.lfaidata.foundation
› Youtube: https://www.youtube.com/channel/UCfasaeqXJBCAJMNO9HcHfbA
› LF AI Logos: https://github.com/lfaidata/artwork/tree/master/lfaidata
› LF AI Presentation Template: https://drive.google.com/file/d/1eiDNJvXCqSZHT4Zk_-czASlz2GTBRZk2/view?usp=sharing

› Events Page on LF AI Website: https://lfaidata.foundation/events/
› Events Calendar on LF AI Wiki (subscribe available): https://wiki.lfaidata.foundation/pages/viewpage.action?pageId=12091544
› Event Wiki Pages: https://wiki.lfaidata.foundation/display/DL/LF+AI+Data+Foundation+Events

**□LF** AI & DATA

# Agenda

› Roll Call  (1 mins)

› Approval of Minutes from previous meeting (2 mins)

› SapientML – New project proposal

› Open Discussion

# TAC Voting Members - Please note

Please ensure that you do the following to facilitate smooth procedural quorum and voting processes:

- Change your Zoom display name to include your First/Last Name, Company/Project Represented
    - example: Nancy Rausch, SAS
- State your First/Last Name and Company/Project when submitting a motion
    - example: First motion, Nancy Rausch/SAS

# TAC Voting Members

Note: we still need a few designated backups specified on [wiki](wiki)

| Member Company or Graduated Project | Membership Level or Project Level | Voting Eligibility | Country | TAC Representative | Designated TAC Representative Alternates |
|---|---|---|---|---|---|
| 4paradigm | Premier | Voting Member | China | Zhongyi Tan | |
| Baidu | Premier | Voting Member | China | Jun Zhang | Daxiang Dong, Yanjun Ma |
| Ericsson | Premier | Voting Member | Sweden | Rani Yadav-Ranjan | |
| Huawei | Premier | Voting Member | China | Howard (Huang Zhipeng) | Charlotte (Xiaoman Hu), Leon (Hui Wang) |
| IBM | Premier | Voting Member | USA | Susan Malaika | Beat Buesser, Alexandre Eichenberger |
| Nokia | Premier | Voting Member | Finland | @Michael Rooke | @Jonne Soininen |
| OPPO | Premier | Voting Member | China | Jimmy (Hongmin Xu) | |
| SAS | Premier | Voting Member | USA | *Nancy Rausch | Liz McIntosh |
| ZTE | Premier | Voting Member | China | Wei Meng | Liya Yuan |
| Adversarial Robustness Toolbox Project | Graduated Technical Project | Voting Member | USA | Beat Buesser | Kevin Eykholt |
| Angel Project | Graduated Technical Project | Voting Member | China | Jun Yao | |
| Egeria Project | Graduated Technical Project | Voting Member | UK | Mandy Chessell | Nigel Jones, David Radley, Maryna Strelchuk, Ljupcho Palashevski, Chris Grote |
| Flyte Project | Graduated Technical Project | Voting Member | USA | Ketan Umare | |
| Horovod Project | Graduated Technical Project | Voting Member | USA | Travis Addair | |
| Milvus Project | Graduated Technical Project | Voting Member | China | Xiaofan Luan | Jun Gu |
| ONNX Project | Graduated Technical Project | Voting Member | USA | Alexandre Eichenberger | Andreas Fehlner, Prasanth Pulavarthi, Jim Spohrer |
| Pyro Project | Graduated Technical Project | Voting Member | USA | Fritz Obermeyer | |
| Open Lineage Project | Graduated Technical Project | Voting Member | USA | *Awaiting confirmation from Project Lead* | |

# Minutes approval

LF AI & DATA

# Approval of August 10, 2023 Minutes

Draft minutes from the August 10 TAC call were previously distributed to the TAC members via the mailing list

**Proposed Resolution:**

› That the minutes of the August 10 meeting of the Technical Advisory Council of the LF AI & Data Foundation are hereby approved.

# Fujitsu's Proposal to Host SapientML project in LF AI & Data

Hiro Kobashi (hkobashi@fujitsu.com)

Masahiro Fukuyori (fukuyori@fujitsu.com)

(Representing AutoML team at Fujitsu)

2023-08-24

# Why contribute SapientML to Linux Foundation

**FUJITSU**

- Neutral holding ground
  - vendor-neutral, not for profit

- Open governance model
  - Transparent and open governance model
  - Instill trust in contributors and adopters in the management of the project
  - Neutral management of projects' assets by the foundation

- Growing community
  - Increase visibility of project through LF ecosystem
  - Increase contributors by converting new & existing users
  - Opportunities to collaborate with other projects

# Agenda

- 1. Background

- 2. Challenges

- 3. How SapientML helps

- 4. Next steps

# Background

# The shortage of AI talent has not been resolved.

FUJITSU

It is still difficult to recruit AI talent.

In particular, there is a shortage of AI talent such as data scientists and data engineers.



Responses suggest that organizations are most often hiring software engineers, data engineers, and AI data scientists.

AI-related roles that respondents' organizations hired, past year, % of respondents[1]

| Role | % |
| --- | --- |
| Software engineers | 39 |
| Data engineers | 35 |
| AI data scientists | 33 |
| Machine learning engineers | 30 |
| Data architects | 28 |
| AI product owners/managers | 22 |
| Design specialists | 22 |
| Data visualization specialists | 21 |
| Translators | 8 |
| None of the above | 14 |

[1]Only asked of respondents whose organizations have adopted AI in at least one function; n = 744.

McKinsey & Company



AI high performers are much more likely than others to have hired AI data scientists, machine learning engineers, and translators in the past year.

AI-related roles that respondents' organizations hired, past year, % of respondents[1]

■ Respondents at AI high performers
□ All other respondents

| Role | Respondents at AI high performers | All other respondents |
| --- | --- | --- |
| Software engineers | 42 | 40 |
| Data engineers | 46 | 37 |
| AI data scientists | 60 | 31 |
| Machine learning engineers | 58 | 27 |
| Data architects | 33 | 31 |
| AI product owners/ managers | 40 | 21 |
| Design specialists (eg, user interface, experience design) | 29 | 22 |
| Data visualization specialists | 34 | 22 |

McKinsey & Company

**It is my personal speculation that the reason for the stagnation in the number of AI application cases is due to the shortage of AI talent.**

The state of AI in 2022—and a half decade in review by McKinsey

14

# AutoML

AutoML is a research activity that **scales AI application by automating AI utilization processes**.

# AutoML: Reality or Pipe Dream?



AutoML beginning to rival top data scientists in *some* cases

% ranking on Kaggle private leaderboard

**Performance of Top Data Scientists**

Vendor #1
Vendor #2
Vendor #3
AutoML Tables

For each product:
- Relevant tables joined by given IDs
- Some minimal preprocessing done to match input requirements
- Run until converge
- Benchmarks run between H2 2018 to today (as they became available)

KDD Cup 2014 - Predicting Excitement at Donors Choose | Allstate claims severity | Porto seguro safe driver prediction | Walmart Recruiting: Trip Type Classification | Mercari Price Suggestion Challenge | Criteo Display Advertising Challenge

Source: Google Cloud Next 2019 talk

# AutoML: Reality or Pipe Dream?

**FUJITSU**

## AI model creation is rapidly automated (AutoML)

**Issue in AI model creation**
- Many try-and-error are needed (taking long time)
- Model quality is highly depending on individual skill

**AutoML evolution**
- The emergence of AutoML, that performs the same level of top data scientist

### California Design Den
Online sale of bedding (*1)



Quickly create AI model with AutoML without data scientists

**Reduce 50% of stocked items in factories**

### G5
Real Estate Customer Acquisition Business (*2)



AI model generated by AutoML is as accurate as its by top data scientists

**Reduce 80% time and record 95% accuracy**

### Lenovo
Device sales (*3)



AutoML enables rapid deployment of accurate models

**Reduced model build time from 4 weeks to 3 days**

# Challenges

# Challenges

- AutoML tools have not yet been adopted widely. Why?

1. **Time-consuming**/Needs intensive computational resources
2. **Blackbox nature** of the generated AI models

# Key factors to success

1. ## **<u>Speed (in AI model creation)</u>**

    - How quickly can we create an AI model to fit in the data scientist's routine

2. ## **<u>Code (with explanation)</u>**

    - Code synthesis with explanation can enable iterations (try&error), which are the matter in AI process.

3. ## **<u>Accuracy</u>**

    - No doubt, accuracy is always matter in AI model creation.

# SapientML

Meta-Learning based AutoML to bring successes to AI model creation process.

# Approach

**1** Extract expert knowledge from codes

**Break the limit (1)**
Current AutoML searches parameters from pre-defined search space. We utilize expert experience, which are stored in codes.

**2** Be interactive with data scientist

**Break the limit (2)**
Current AutoML creates AI model only. There is no flexibility (ex. No chance to modify, hard to understand why the model is good)

## Meta-Learning based AutoML

- Synthesis high accurate **AI model creation code** by using expert knowledge
- Provide the modifiable code to interact with data scientist

22

# Generative AutoML: SapientML

- AutoML that harnesses expertise of DS encapsulated in large corpora of existing ML pipelines, e.g., Kaggle



- Focus on *structured data - most abundant data in enterprise and has many use cases.*

*SapientML: An AutoML approach harnessing the wisdom (sapere) of human (sapien) data scientists.*

# How works SapientML?

- Input Table Data and ML Spec. □ Output AI model with code

# SapientML at ICSE2022



SAPIENTML: Synthesizing Machine Learning Pipelines by Learning from Human-Written Solutions
*Ripon Saha, Akira Ura, Sonal Mahajan, Chenguang Zhu, Linyi Li, Yang Hu, Hiroaki Yoshida, Sarfraz Khurshid, Mukul Prasad*

May 21-29, 2022
Pittsburgh, PA, USA

**S-rank, CORE: A\***
Flagship SE Conference

**Acceptance rate:**
26% (197/751)

- SapientML out-performs all academic AutoML tools on incl. SoTA  AL (MIT)
  - 41 benchmarks, including 20 Kaggle competitions
- No failures, highest number of wins
- SapientML synthesis is quite robust to variations in training data – pipelines should generalize well

FUJITSU

## Example: IEEE-CIS-Fraud-Detection (Kaggle Competition)

Rows: e-commerce transactions, represented by 394 features, device type, product features, etc.

Dataset

| | TransactionID | isFraud | TransactionDT | TransactionAmt | ProductCD | card1 | card2 | card3 | card4 | card5 | ... | V330 | V331 | V332 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3460022 | 0 | 12233710 | 107.950 | W | 13623 | 585.0 | 150.0 | visa | 226.0 | ... | NaN | NaN | NaN |
| 1 | 3293299 | 0 | 7602727 | 97.000 | W | 2722 | NaN | 150.0 | visa | 226.0 | ... | NaN | NaN | NaN |
| 2 | 3284094 | 0 | 7337050 | 57.950 | W | 15372 | 241.0 | 150.0 | visa | 226.0 | ... | NaN | NaN | NaN |
| 3 | 3066993 | 1 | 1724715 | 76.023 | C | 9633 | 296.0 | 185.0 | visa | 138.0 | ... | NaN | NaN | NaN |
| 4 | 3390641 | 1 | 10186470 | 23.926 | C | 14276 | 177.0 | 185.0 | mastercard | 137.0 | ... | NaN | NaN | NaN |
| 5 | 3479148 | 1 | 12848316 | 77.000 | W | 6174 | 490.0 | 150.0 | visa | 226.0 | ... | NaN | NaN | NaN |

**Task:** Classifying if a transaction is *Fraud*

**Problem type:** (binary) classification

**Target column: isFraud** (*1 or 0*)

# SpaientML generated pipeline



```
# LOAD DATA
__train_dataset=pd.read_csv("training.csv", delimiter=",")
__test_dataset=pd.read_csv("test.csv", delimiter=",")

# PREPROCESSING-1
_STRING_CATG_COLM_HAS_MISSING = ['card4', 'card6', 'P_emaildomain',…]
for _col in _STRING_COLS_WITH_MISSING_VALUES:
    __si = SimpleImputer(missing_values=np.nan, strategy='most_frequent')
    __train_dataset[_col] = __si.fit_transform(__train_dataset[_col].values.reshape(-1,1))[:,0]
    __test_dataset[_col] = _si.transform(__test_dataset[_col].astype(__train_dataset[_col].dtypes).values.reshape(-1,1))[:,0]

# PREPROCESSING-2
_CAT_COLS = ['ProductCD', 'card4', 'card6', …,'M9']
_ohe = OrdinalEncoder(handle_unknown="use_encoded_value", unknown_value=-1)
__train_dataset[_CAT_COLS] = pd.DataFrame(_ohe.fit_transform(__train_dataset[_CAT_COLS]), columns=_CAT_COLS)
__test_dataset[_CAT_COLS] = pd.DataFrame(_ohe.transform(__test_dataset[_CAT_COLS]), columns=_CAT_COLS)

# PREPROCESSING-4
from sklearn.preprocessing import StandardScaler
__ss= StandardScaler()
__feature_train = pd.DataFrame(__ss.fit_transform(__feature_train.values), index=__feature_train.index, columns=__feature_train.columns)
__feature_test = pd.DataFrame(__ss.transform(__feature_test.values), index=__feature_test.index, columns=__feature_test.columns)

# PREPROCESSING-5
from imblearn.over_sampling import SMOTE
smote = SMOTE()
__feature_train, __target_train = smote.fit_resample(__feature_train, __target_train)

# MODEL
from catboost import CatBoostClassifier
__model = CatBoostClassifier()
__model.fit(__feature_train, __target_train)
__y_pred = __model.predict(__feature_test)

# EVALUATION
from sklearn import metrics
__f1 = metrics.f1_score(__target_test, __y_pred, average='macro')
print('RESULT: F1 Score: ' + str(__f1))
```

*# shortened for presentation*

Load (training, test) data
↓
Fill missing values
↓
Assign numeric encoding to categorical strings
↓
Apply Scaling
↓
Apply sampling to balance data
↓
Train CatBoostClassifier
↓
Evaluate F1 score

# GUI – Simple 3 steps to generate AI



## 1.Load dataset

## 2.Specify ML task

## 3. Request to generate

# GUI – Prediction using generated AI



It generates codes for top 3 AI Models as well as AI model performances

We can check feature importance and correlation between feature and target variable(column).

We can predict sales price for a set of unseen houses (i.e., test dataset) by clicking 「Predict」.

## Start Exploratory Data Analysis instantly



## Customize AI model easily by modifying generated code

# Information for proposal

- GH repo: https://github.com/sapientml/sapientml

- License: Apache 2.0

- Proposal: https://github.com/lfai/proposing-projects/blob/master/proposals/sapientml.adoc

- Possible Collaboration in LF AI&Data
  - Pre-processing: Amundsen, Feast, Feathr
  - Modelling: Adversarial Robustness Toolbox, AI Explainability 360, AI Fairness 360, Intersectional Fairness
  - Deployment: Acumos
  - Coding: Elyta, Kedro

# The Value of SapientML

FUJITSU

## User
(Data Scientist)

Beginner to intermediate level

Expert

## Value

✓ **No programming required. Generation of practical AI models**

✓ **Reduces the effort of trial and error by making clear model selection including code generation reasons**

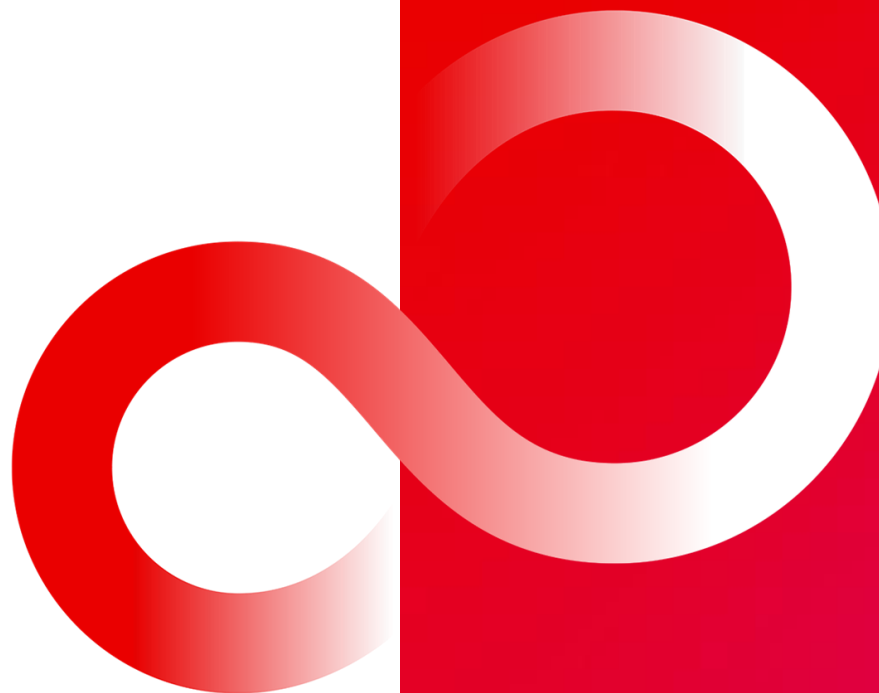✓ **Freely customizable including proprietary know-how and additional precision tuning**

## SapientML Functions/Features

High speed

High Accuracy

Code with Explanation

**Code synthesis technology using past knowledge**

· Fast code generation with AI prediction
· High accuracy using past knowledge

· Give Reason for Code Generation

PF Independent

**Platform independent because it generates Machine Learning programs (*)**
*) Major AutoML tools are platform dependent because they output AI models on their cloud services

32

© 2023 Fujitsu Limited

# Next Step

FUJITSU

# Next Step

- Support more meta-features of datasets

- Support more ML components
  - Preprocess components
  - Model


☐ Accumulate more datasets and ML pipelines from community

Thank you

# Approval on SapientML

**Proposed Resolution:**

› SapientML as a sandbox project of the LF AI & Data Foundation is hereby approved.

# Upcoming TAC Meetings

# Upcoming TAC Meetings

› September 7 – Update from the Trusted AI  committee, Update from the MLSecOps Committee

› September 21 – Marquez graduation request;  AIDA, a new project requesting Sandbox Incubation


Please note we are always open to special topics as well.


If you have a topic idea or agenda item, please send agenda topic requests to tac-general@lists.lfaidata.foundation

**LF** AI & DATA

# Upcoming Events of Interest

› 2023 AICON Middle East Summit - October 8th to 9$^{th}$ in Riyadh
https://lfaidata.foundation/blog/2023/07/18/2023-aicon-middle-east-summit-call-for-topics-from-around-the-world/

› Open Source Summit Europe in Bilbao, Spain, September 19-21 – LF AI&Data will have a booth
https://events.linuxfoundation.org/open-source-summit-europe/

**LF** AI & DATA

24AUG2023

# Open Discussion

**□LF** AI & DATA

# TAC Meeting Details

› To subscribe to the TAC Group Calendar, visit the wiki: https://wiki.lfaidata.foundation/x/cQB2 _____

› Join from PC, Mac, Linux, iOS or Android: https://zoom.us/j/430697670

› Or iPhone one-tap:

  › US: +16465588656,,430697670# or +16699006833,,430697670#

› Or Telephone:

  › Dial(for higher quality, dial a number based on your current location):

  › US: +1 646 558 8656 or +1 669 900 6833 or +1 855 880 1246 (Toll Free) or +1 877 369 0926 (Toll Free)

› Meeting ID: 430 697 670

› International numbers available: https://zoom.us/u/achYtcw7uN

**LF** AI & DATA

# Legal Notice

**LF AI & DATA**