

# Meeting of the Technical Advisory Council (TAC)

February 25, 2021

 **DLF** AI & DATA

# Anti-Trust 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 Updegrave of the firm of Gesmer Undergone LLP, which provides legal counsel to the Linux Foundation.

# Recording of Calls

## **Reminder:**

TAC calls are recorded and available for viewing on the [TAC Wiki](#)

# Reminder: LF AI & Data Useful Links

- › Web site: [lfaidata.foundation](https://lfaidata.foundation)
- › Wiki: [wiki.lfaidata.foundation](https://wiki.lfaidata.foundation)
- › GitHub: [github.com/lfaidata](https://github.com/lfaidata)
- › Landscape: <https://landscape.lfaidata.foundation> or <https://l.lfaidata.foundation>
- › Mail Lists: <https://lists.lfaidata.foundation>
- › Slack: <https://slack.lfaidata.foundation>
- ›
- › 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](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>

# Agenda

- › Roll Call (5 mins)
- › Approval of Minutes from Jan 28 and Feb 11 (5 mins)
- › Incubation proposal (40 minutes)
  - › Flyte (Ketan Umare)
- ›
- › LFAI General Updates (5 minutes)
- › Open Discussion (5 minutes)

# TAC Voting Members

\* = still need backup specified on [wiki](#)

Board Member	Contact Person	Email
AT&T	Anwar Atfab	<a href="mailto:anwar@research.att.com">anwar@research.att.com</a>
Baidu	Ti Zhou	<a href="mailto:zhouti@baidu.com">zhouti@baidu.com</a>
Ericsson	Rani Yadav-Ranjan*	<a href="mailto:rani.yadav-ranjan@ericsson.com">rani.yadav-ranjan@ericsson.com</a>
Huawei	Huang Zhipeng*	<a href="mailto:huangzhipeng@huawei.com">huangzhipeng@huawei.com</a>
IBM	Susan Malaika	<a href="mailto:malaika@us.ibm.com">malaika@us.ibm.com</a>
Nokia	Jonne Soininen*	<a href="mailto:jonne.soininen@nokia.com">jonne.soininen@nokia.com</a>
SAS	Nancy Rausch	<a href="mailto:nancy.rausch@sas.com">nancy.rausch@sas.com</a>
Tech Mahindra	Nikunj Nirmal	<a href="mailto:nn006444@techmahindra.com">nn006444@techmahindra.com</a>
Tencent	Bruce Tao	<a href="mailto:brucetao@tencent.com">brucetao@tencent.com</a>
Zilliz	Jun Gu	<a href="mailto:jun.gu@zilliz.com">jun.gu@zilliz.com</a>
ZTE	Wei Meng	<a href="mailto:meng.wei2@zte.com.cn">meng.wei2@zte.com.cn</a>
Graduate Project	Contact Person	Email
Acumos	Nat Subramanian	<a href="mailto:natarajan.subramanian@techmahindra.com">natarajan.subramanian@techmahindra.com</a>
Angel	Bruce Tao	<a href="mailto:brucetao@tencent.com">brucetao@tencent.com</a>
Egeria	Mandy Chessell	<a href="mailto:mandy_chessell@uk.ibm.com">mandy_chessell@uk.ibm.com</a>
Horovod	Travis Addair*	<a href="mailto:taddair@uber.com">taddair@uber.com</a>
ONNX	Jim Spohrer (Chair of TAC)	<a href="mailto:spohrer@us.ibm.com">spohrer@us.ibm.com</a>
Pyro	Fritz Obermeyer*	<a href="mailto:fritz.obermeyer@gmail.com">fritz.obermeyer@gmail.com</a>

# Approval of January 28th, 2021 Minutes

Draft minutes from the January 28<sup>th</sup> TAC call were previously distributed to the TAC members via the mailing list

## **Proposed Resolution:**

- › That the minutes of the January 28<sup>th</sup> meeting of the Technical Advisory Council of the LF AI & Data Foundation are hereby approved.

# Approval of February 11th, 2021 Minutes

Draft minutes from the February 11<sup>th</sup> TAC call were previously distributed to the TAC members via the mailing list

## **Proposed Resolution:**

- › That the minutes of the February 11<sup>th</sup> meeting of the Technical Advisory Council of the LF AI & Data Foundation are hereby approved.



# Incubation Proposal - Flyte

Ketan Umare <ketan.umare@gmail.com>

# Project Contribution Proposal Review & Discussion: Flyte

Flyte is a container-native, type-safe workflow and pipelines platform optimized for large scale processing and machine learning written in Golang. Workflows can be written in any language, with out of the box support for Python, Java and Scala.

**Presenter:** Ketan Umare <ketan.umare@gmail.com>

## **Resources:**

Github: <https://github.com/flyteorg/flyte>

Project Level: Incubation

Proposal: <https://github.com/lfai/proposing-projects/blob/master/proposals/flyte.adoc>

# Flyte Overview



# Agenda

- Problem, Motivation & Goal
- What are the challenges for the users?
- Introducing Flyte
- Overview of Flyte's architecture
- Future
- Case for Contributing

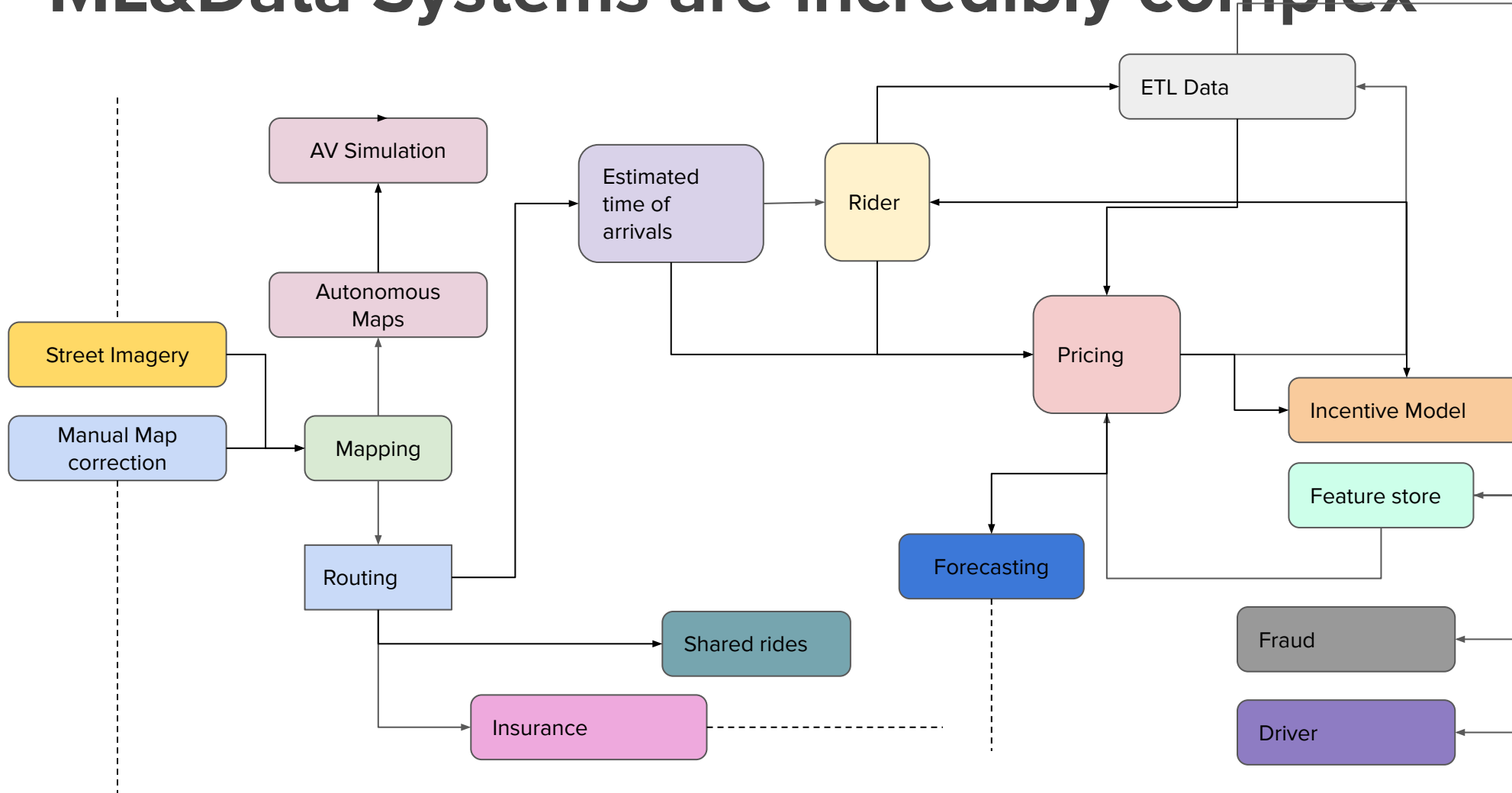
---

# Motivation & Goal

---

## Motivation & Goal

# ML&Data Systems are incredibly complex

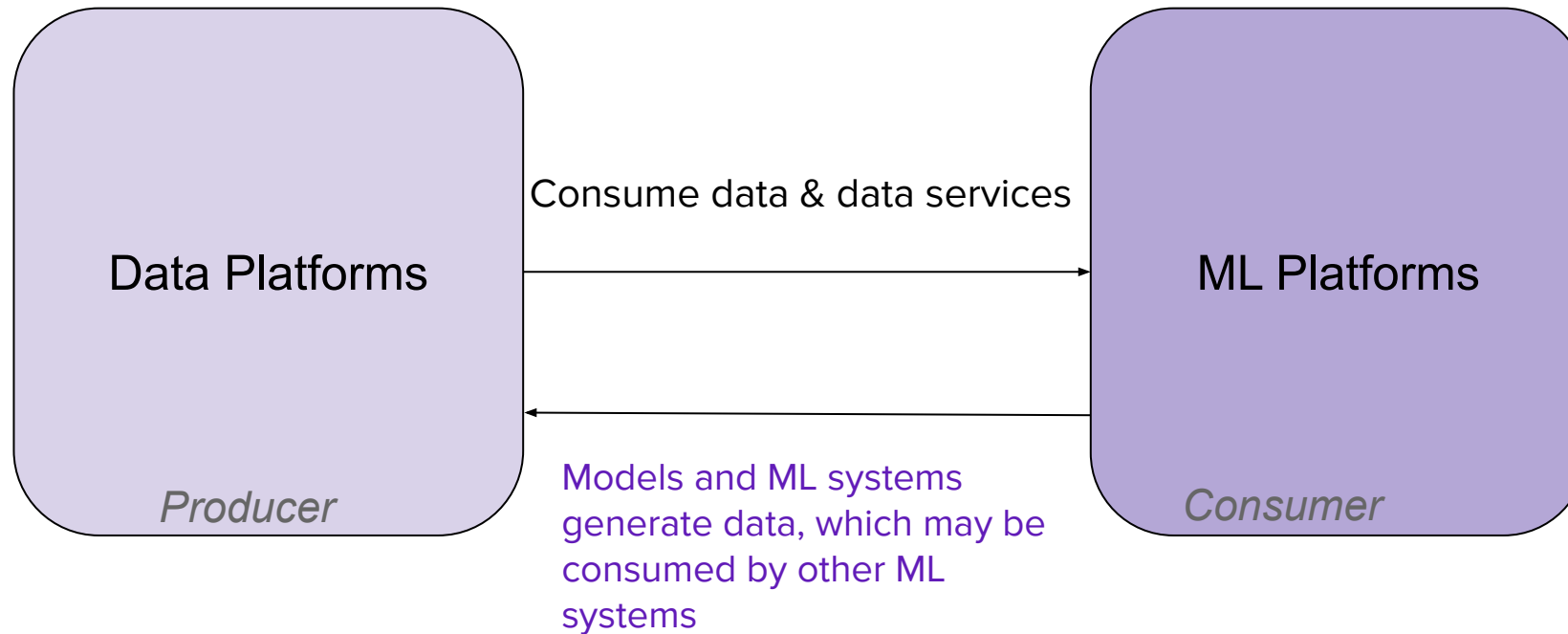


Data and ML processes often interact.

Data Flow is very complex and machine learning is more than just model code.

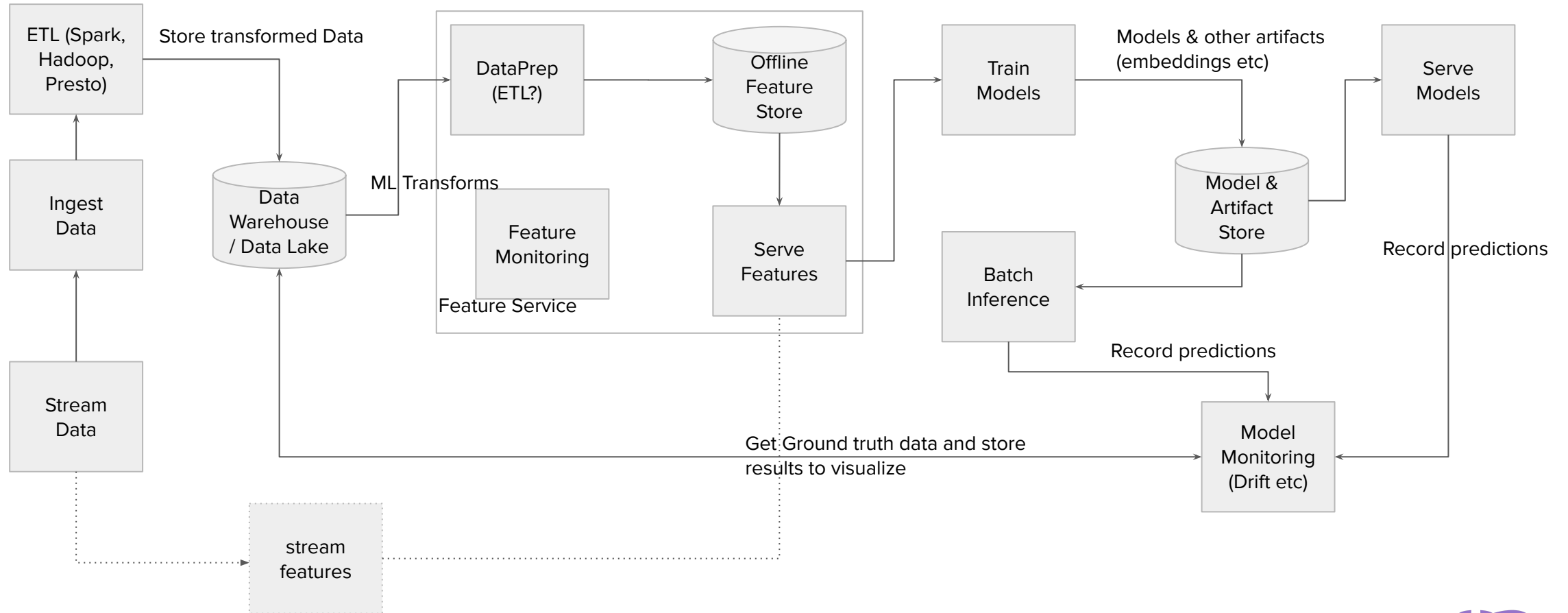
## Motivation & Goal

# The Line is getting blurred



## Motivation & Goal

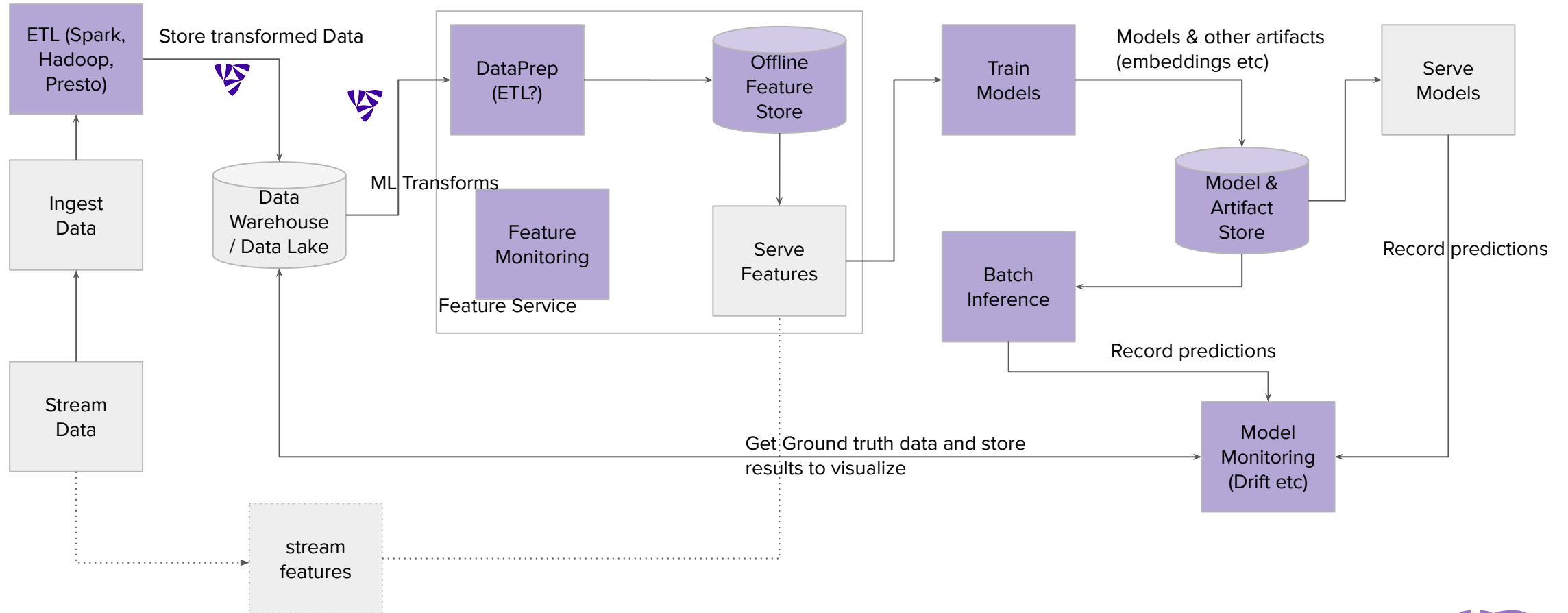
# View: ML as superset of Data





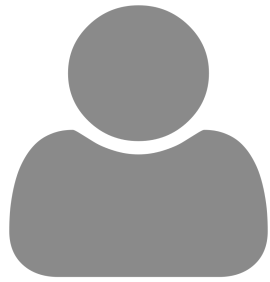
## Motivation & Goal

# Where Flyte fits in...



## Motivation & goal

# Who are the users - ML is part of the product!



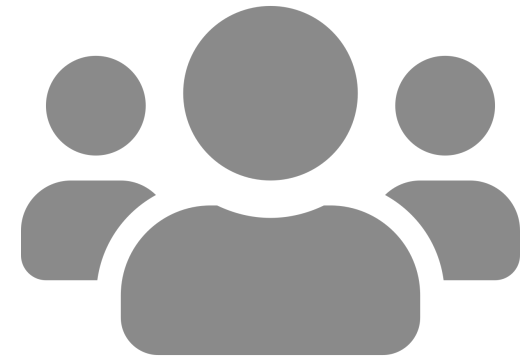
Data Scientists



Data Engineers



ML SWE's



SWE's

# Use Cases

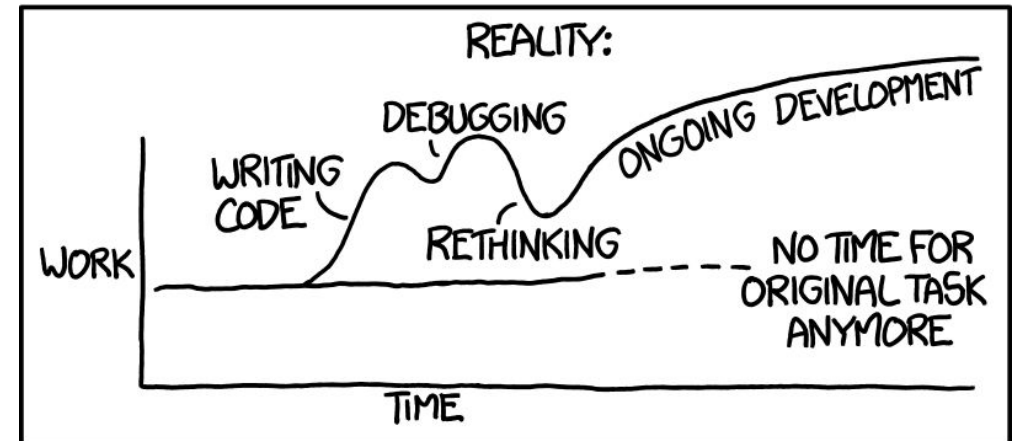
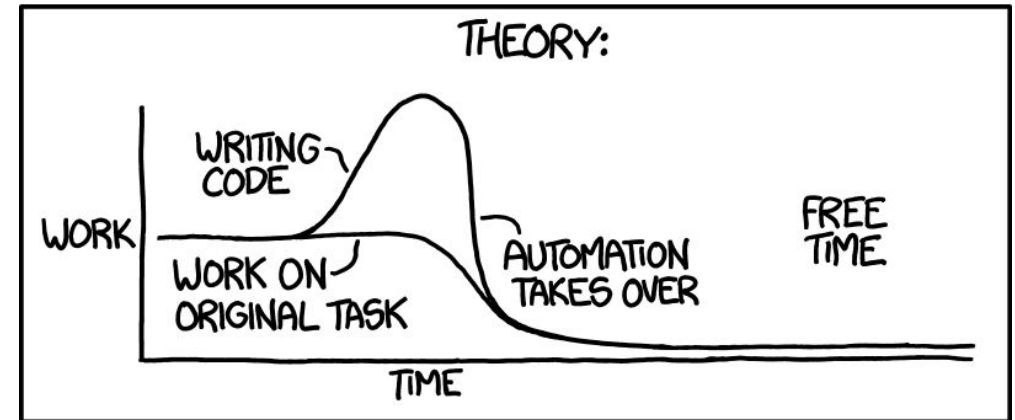
- Infrastructure
  - Match users workflow
  - Dynamic & Parameterizable
  - Simplified Ops
  - Collaboration
  - Flexible
-

## Use Cases

# Serverless experience

- No Infrastructure
- Isolated development and management
- Access resources - CPU/GPU/Mem etc
- Framework/Library independence
- Multi-tenancy unaware
- Seamless scalability
- Freedom to access
- Control costs

"I SPEND A LOT OF TIME ON THIS TASK.  
I SHOULD WRITE A PROGRAM AUTOMATING IT!"



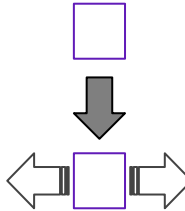
## Use Cases

# Consistent API for Jobs and Pipelines

1

**Start with one Job**

*e.g. spark job, a training job, a query etc*



2

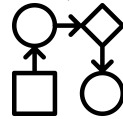
**Scale the Job**

*e.g. one region -> all regions, more GPUs*

3

**Create a pipeline**

*e.g. Fetch data -> train model -> calculate metrics*



4

**Run the pipeline on demand**

*e.g. Run a pipeline with parameters*



5

**They want to run the pipeline on schedule**

*e.g. Run every hour*



6

**Retrieve results for jobs / pipelines**



## Use Cases

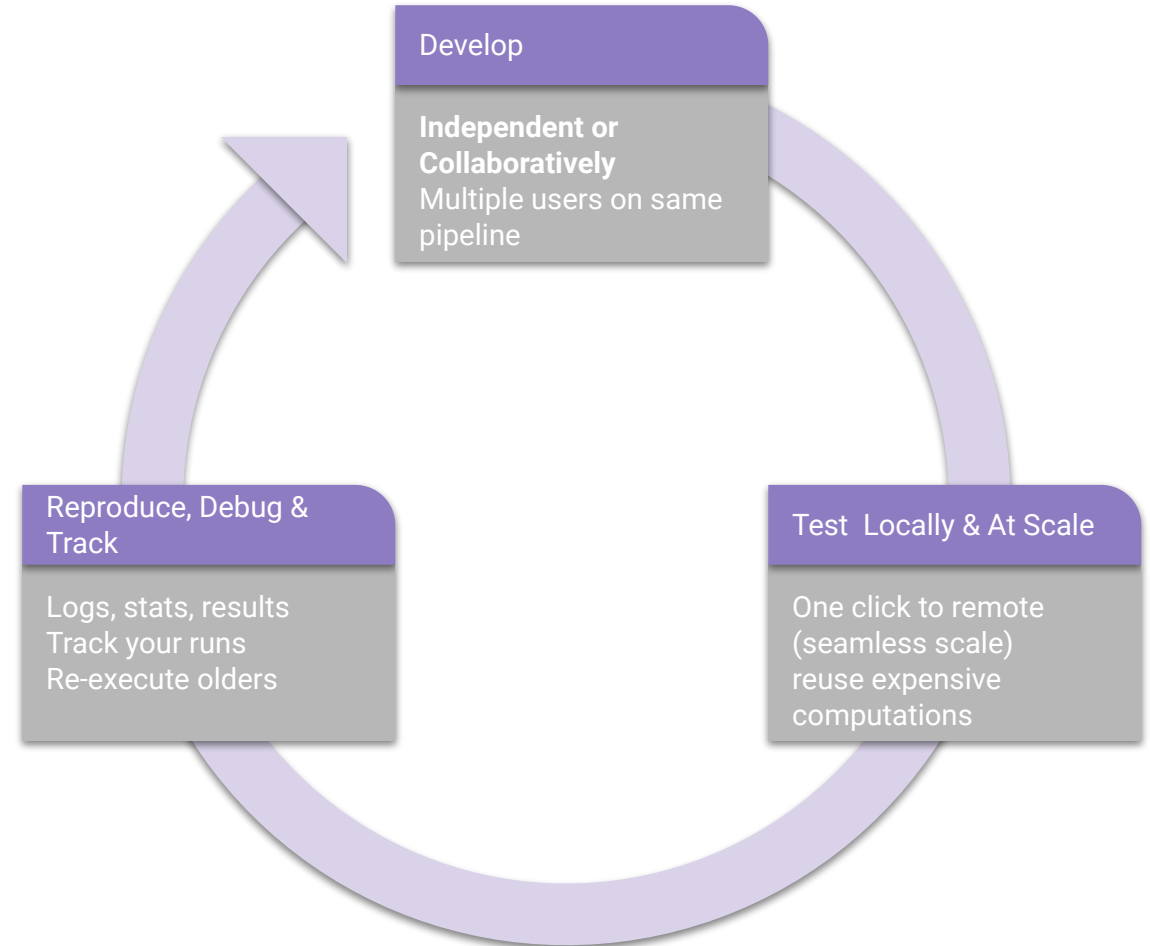
# Parameterize executions and dynamism

- **Parameterized Experiments** - Alter the coefficients, data-sets etc
- **Simultaneous Schedules** - production & shadow experiments
- **Visualize** results in the UI
- **Time** is only one parameter (time can be in the past, future or present)
- **Data lineage & Provenance** - how, where, why & when?
- **Dynamically generate pipelines** based on parameters at runtime
  - Add more regions on some days
  - Alter the shape of map generation based on data updates etc

## Use Cases

# Development & Iteration

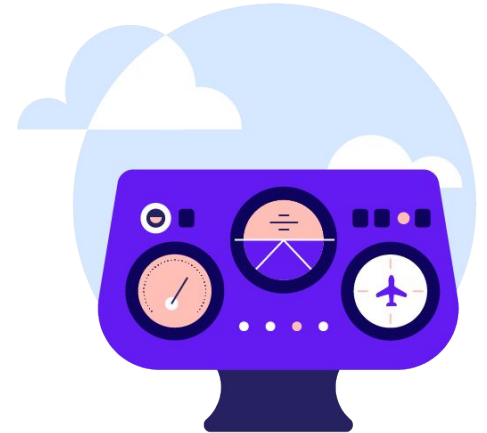
- Development should **NOT affect production** data
- Write code in multiple languages, but *Python is special “lingua-franca”*
- **Test at scale** - more data, gpus, faster training
- **Reuse** expensive computations - Fix bugs in parts of the pipeline and then run only those parts



## Use Cases

# Ops & Visibility

- **History of all executions** - parameters and debug data
- **Notifications** - events, failures (Pagerduty, Slack, emails)
- **Debug** issues in **production** using Logging, stats
- Get a **custom dashboard** to visualize stats for their pipelines
- **Track their runs independently** and may want to rewind time

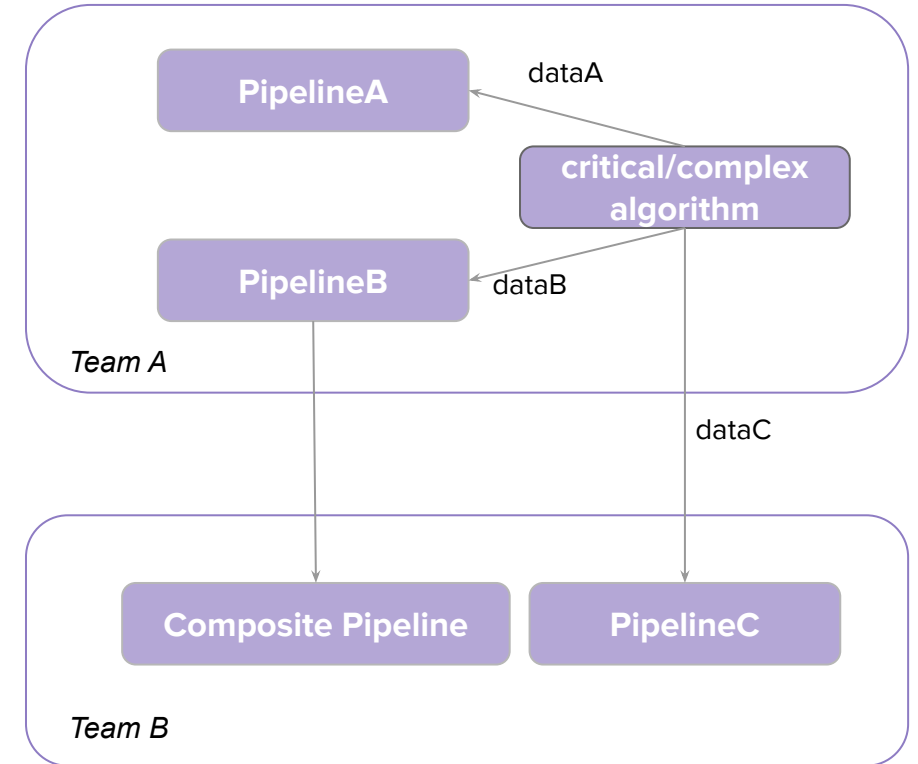




## Use Cases

# Reusability & Shareability

- **Write once Reuse code** - Feature extraction, specific deployment artifacts etc
- **Share artifacts & code** - without breaking dependencies & organizational boundaries
- **Compose pipelines** using pipelines from other teams



- *Composite pipeline is composed of TeamA.PipelineB + other tasks.*
- *PipelineC re-uses the shared critical task*

## Use Cases

# Extensibility & Flexibility

1

### Users want flexibility

*Add simple python extensions (Airflow operators)  
Maybe only for their teams*

**Flytekit** makes it easy to add new user customizations  
Flyte also allows you to run just your own containers

2

### Platform wants to keep adding new capabilities

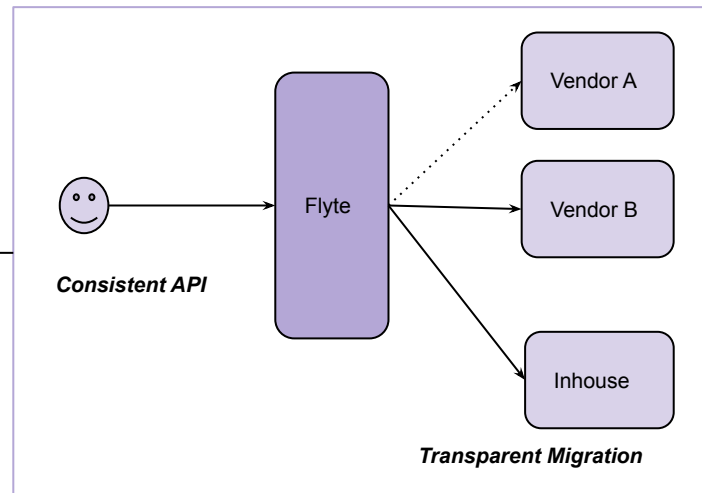
*Distributed training support, Spark, Streaming etc  
Continue adding and controlling roll-out of features*

**Flyte backend plugins** are independently deployed, maintained  
and are in the hosted service

3

### Organizations want flexibility

*Control costs (Migrate vendors, bring capabilities inhouse)  
Users velocity and existing code should **just work!***



**Flyte control plane**  
makes it possible to  
switch plugin  
associations and OSS  
makes it possible to  
migrate



- **Opinionated, scalable and hosted Orchestration Platform**
  - **Fabric that connects disparate compute technologies**
  - **Extensible, Observable & shareable**
  - **Integrates best of the breed open source solutions**
  - **Auditable, Repeatable & Secure**
-

# User Journey

## Ideate & Iterate

1. Write **business logic**
2. Test task **locally**
3. Test task **remote**
4. Orchestrate multiple tasks into a **Workflow**
5. **Execute** the workflow
6. Repeat

## Productionize

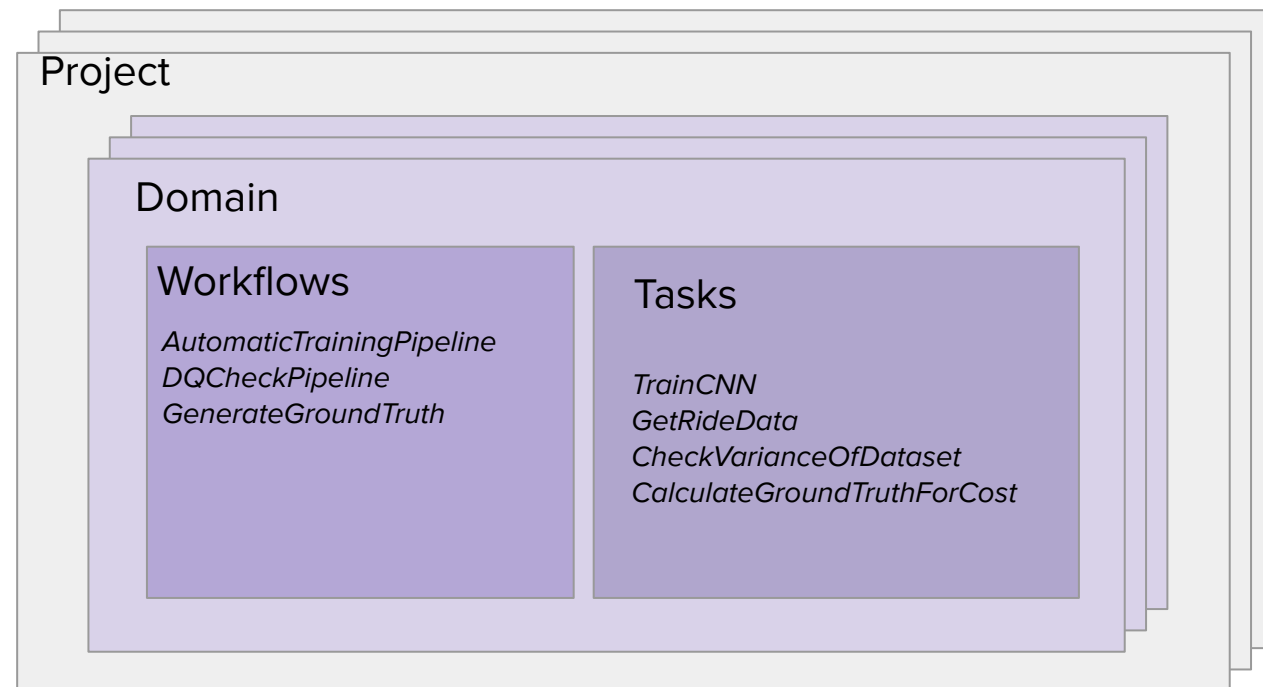
1. Promote a pipeline to **production** (CI/CD)
2. Create **one or more schedules**
3. Execute **ad-hoc**
4. **Monitor** and get notified

## Retrieve & Replay

1. Retrieve results from **executions**
2. Identify production errors
3. **Replay, reproduce** historical artifacts
4. Retrieve **artifact lineage**

# Multi-tenancy & Organization

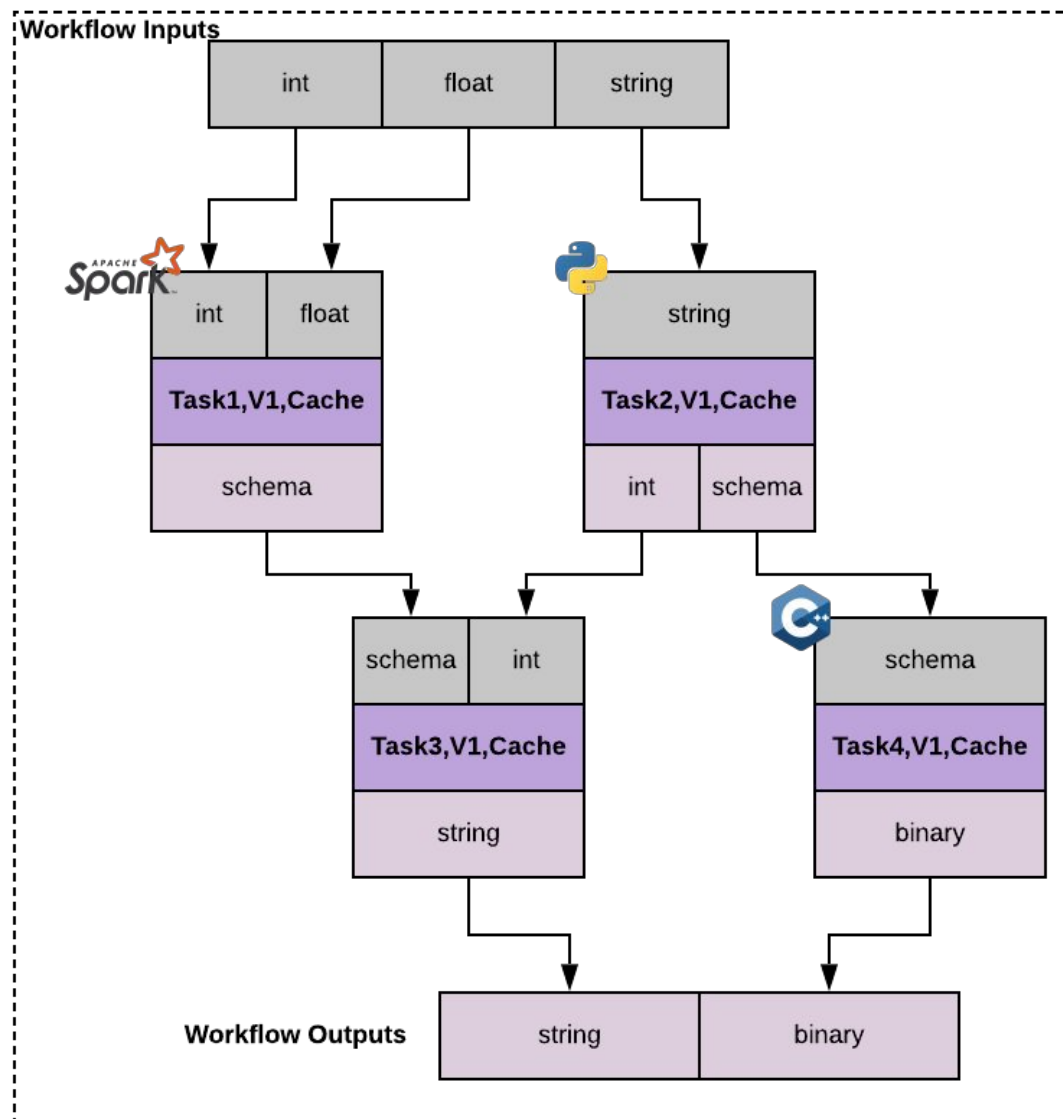
- Top tenant entity - **Project**
  - AVPerception
  - ETAModels
  - PricingModels
- Each tenant can have **Domains**
  - Development, Production
  - CI / CD
- **Workflows & Tasks**



## Flyte: Concepts

# Tasks & Workflows

- Declarative (protobuf)
- **Versioned**
- Strongly **typed interfaces**
- Models the flow of Data
- Tasks
  - Arbitrarily complex
  - Encapsulate user code
- Workflows
  - Composable
  - Dynamic
  - DSL in python (& JAVA)



## Flyte: Concepts

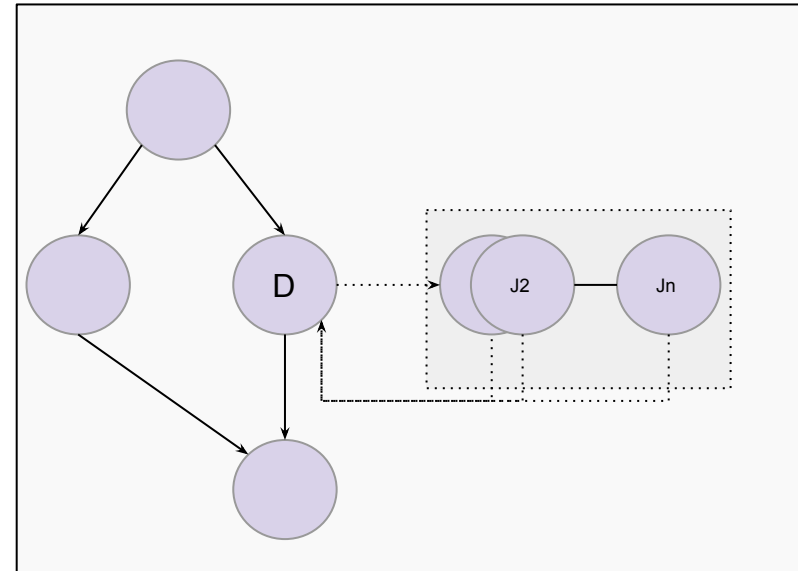
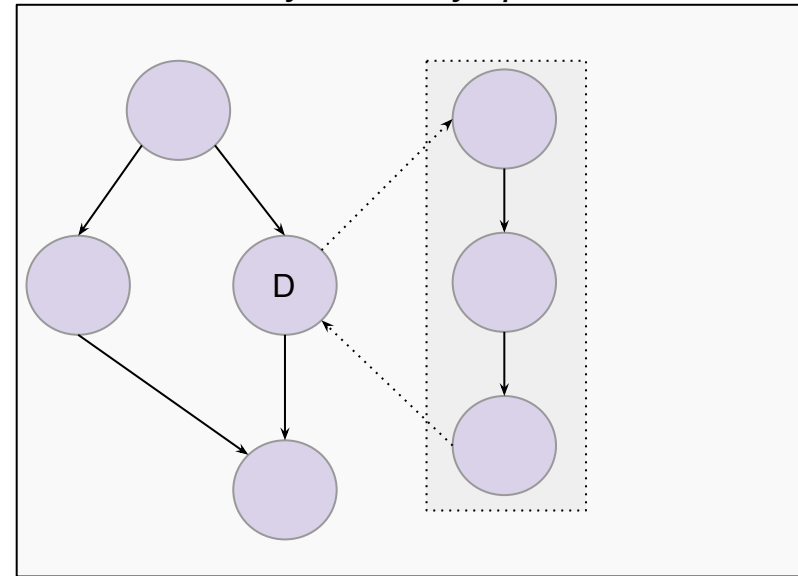
# Dynamism

Flyte allows certain nodes to alter the shape of the graph

Data parallel jobs, dynamic generation of workflows (generate logic using the available data).

*Flyte can scale to more than 10000 nodes in a graph, each with arbitrarily complex execution logic*

*Dynamically spawn a workflow*



*Dynamically spawn an array of map jobs*

## Flyte: Concepts

# DataCatalog: Lineage & Memoization

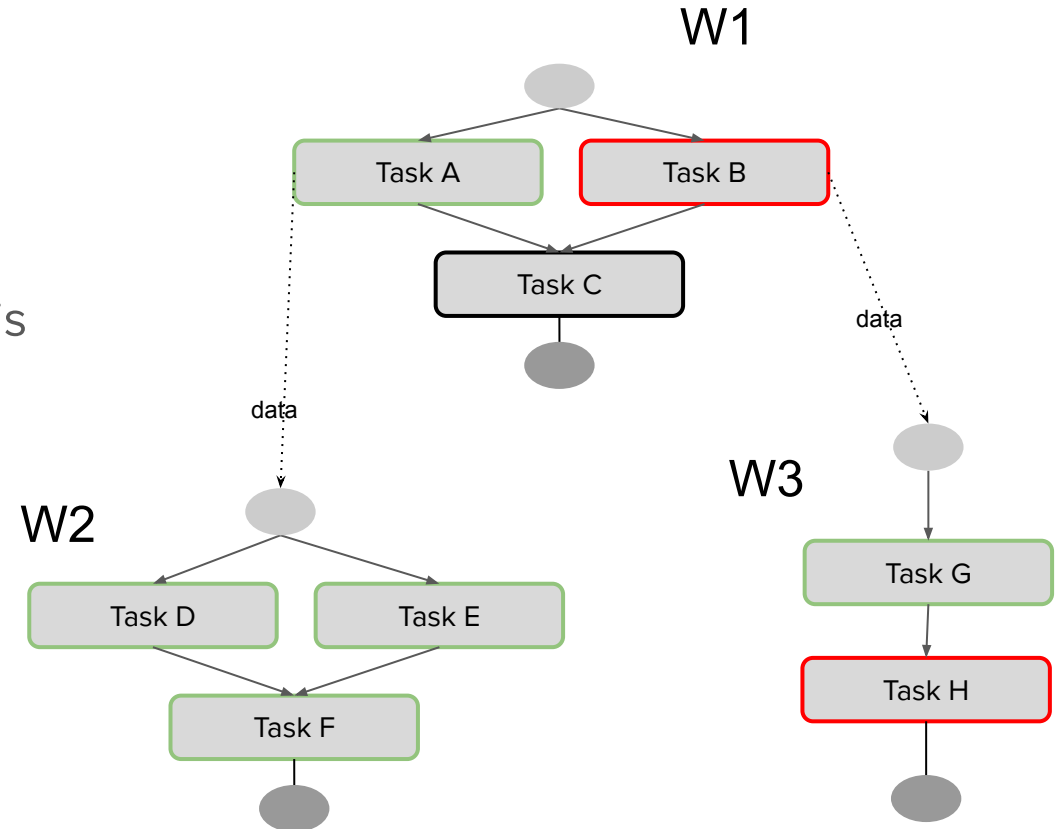
Every task execution in Flyte is **recorded** by default in Catalog Service. This enables Flyte executions to have,

### Artifact Lineage

- **Causal** dependencies between data and processes is tracked

### Memoization

- Each task execution has a **unique signature**, which includes the input values & version of code
- **Repeated** executions with matching signatures are cached





Flyte

# Real Production Scale @Lyft

9k Unique Workflows defined

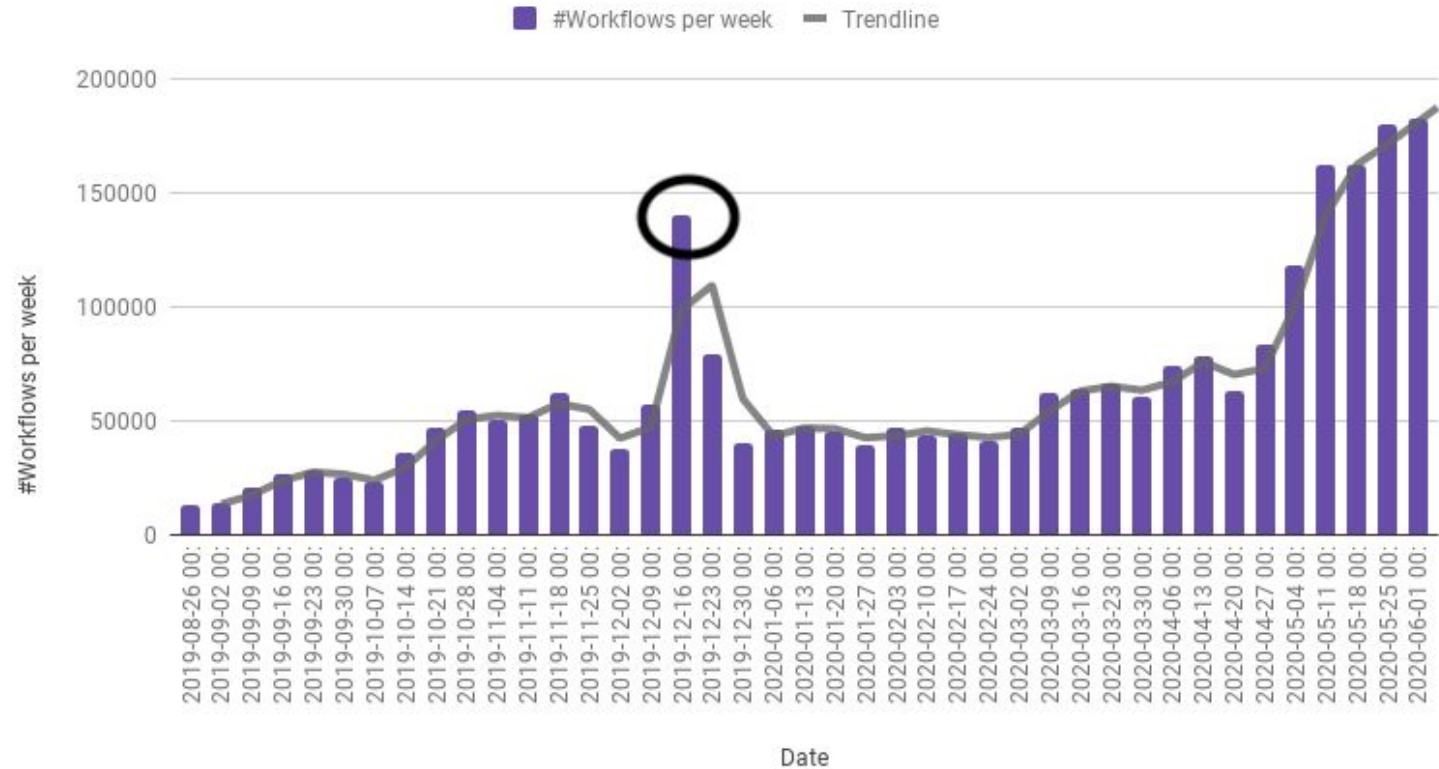
54k Unique Task definitions

1M+ Workflow executions per month

10M+ task executions per month

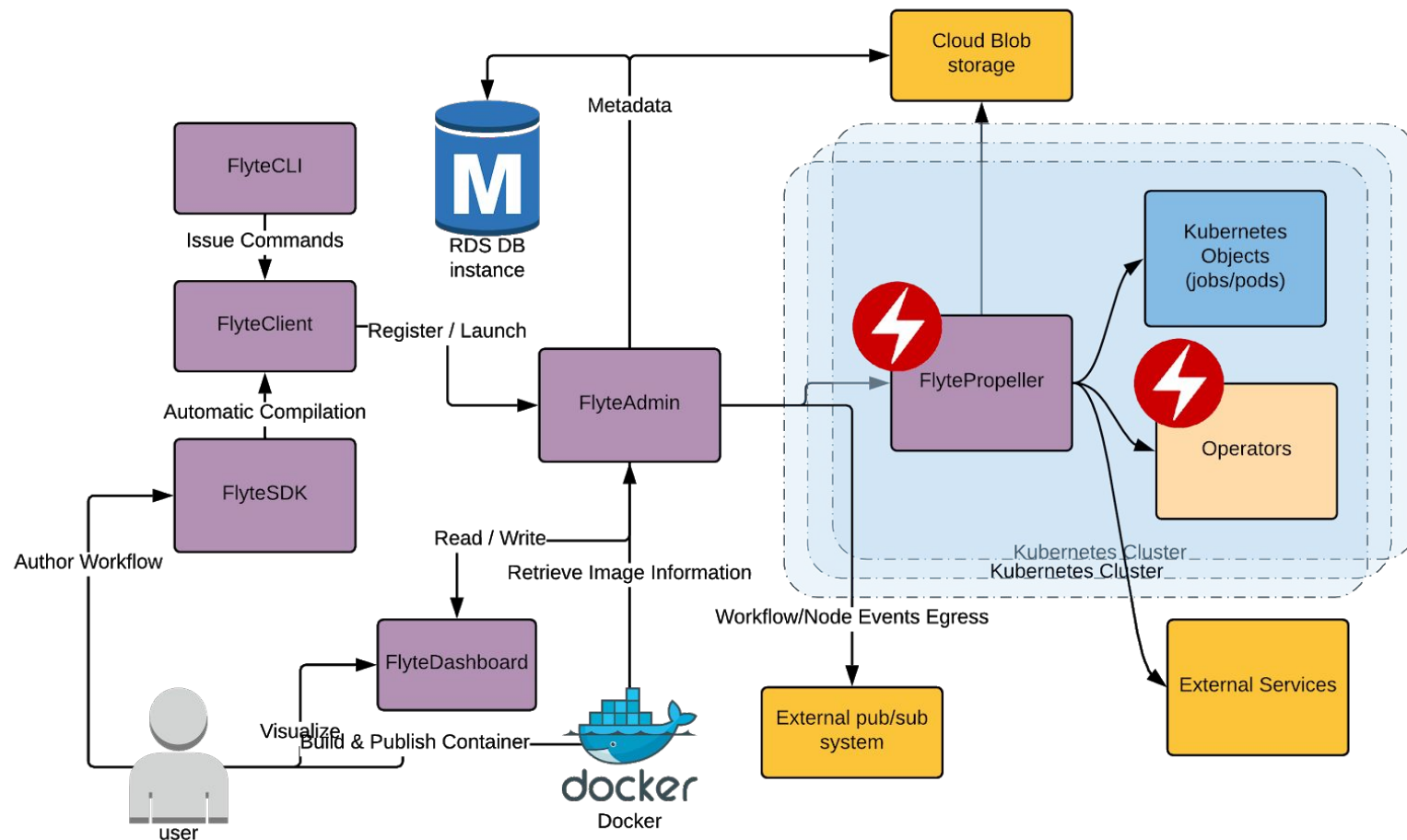
40M+ containers executed per month

#Workflows per week



# Scale Out!

- **Multi K8s cluster** out of the box
- **Highly** optimized for scale and hosted
- Kubernetes makes it possible to orchestrate containers
- Operators make it possible to have K8s services.



# Flyte

## Code Sample & UI execution

Python

```
@task
def t1(a: int) -> pandas.DataFrame:
    return pandas.DataFrame(data={"col1": [a, 2],
                                   "col2": [a, 4]})

@task
def t2(df: pandas.DataFrame) -> pandas.DataFrame:
    return df.append(pandas.DataFrame(data={"col1": [5,
                                                  10], "col2": [5, 10]}))

@workflow
def my_wf(a: int) -> pandas.DataFrame:
    return t2(df=t1(a=a))
```

Scala

```
case class SumTaskInput(a: Long, b: Long)
case class SumTaskOutput(c: Long)

class SumTask
  extends SdkRunnableTask(
    SdkScalaType[SumTaskInput],
    SdkScalaType[SumTaskOutput]
  ) {

  override def run(input: SumTaskInput): SumTaskOutput =
  {
    SumTaskOutput(input.a + input.b)
  }
}
```

The screenshot displays the Flyte UI for a workflow execution. At the top, a purple header bar shows the status 'SUCCEEDED' in a green box, the workflow path 'bulkmapmatching/development/AlignGroundTruth/wsansqr1f', and a 'Relaunch' button. Below the header, a table provides details for the workflow: Domain (development), Version (b45d67f40b61655a16f3c934f647357c291bc7b3), Cluster (flyte2), and Time (7/8/2020 4:39:31 AM UTC). The main content area is divided into two sections: 'Nodes' and 'Graph'. The 'Nodes' section contains a table with columns for Node, Status, Type, and Start Time. The 'Graph' section shows a tree view of the workflow nodes, including 'ensure-ta', 'distinct-gt', 'get-groun', 'f15a1', 'f5fkoz', 'f5ksor', and 'get-groun'. The 'distinct-gt-ids-and-ds-task' node is highlighted, and its details are shown on the right. This details panel includes the task name, a close button, the task ID, a 'SUCCEEDED' status box, the task type (Hive Batch Task), and a table for 'Inputs' showing 'ds\_string: 2020-07-02'. A 'Back to parent' button is also visible.

Domain	Version	Cluster	Time
development	b45d67f40b61655a16f3c934f647357c291bc7b3	flyte2	7/8/2020 4:39:31 AM UTC

Node	Status	Type	Start Time
> ensure-ta	SUCCEEDED	Hive Batch Task	7/8/2020 4:39:31 AM UTC 7/7/2020 9:39:31 PM PDT
distinct-gt	SUCCEEDED	Hive Batch Task	
▼ get-groun	SUCCEEDED	Dynamic Task	7/8/2020 4:41:38 AM UTC 7/7/2020 9:41:38 PM PDT
> f15a1	SUCCEEDED	Dynamic Task	7/8/2020 4:42:16 AM UTC 7/7/2020 9:42:16 PM PDT
> f5fkoz	SUCCEEDED	Dynamic Task	7/8/2020 4:42:16 AM UTC 7/7/2020 9:42:16 PM PDT
> f5ksor	SUCCEEDED	Dynamic Task	7/8/2020 4:42:16 AM UTC 7/7/2020 9:42:16 PM PDT
> get-groun	SUCCEEDED	Dynamic Task	7/8/2020 4:46:13 AM UTC 7/7/2020 9:46:13 PM PDT

Executions	Inputs	Outputs	Task
	ds_string: 2020-07-02		

https://flyte.lyft.net/console/projects/bulkmapmatching/domains/development/workflows/AlignGroundTruth

# Tasks are standalone entities

The screenshot shows the Flyte console interface. At the top, there is a header with the Flyte logo and the email 'kumare@lyft.com'. Below the header, the 'PROJECT' is set to 'KubeconDemo 2019'. The environment is 'Development'. A search bar is present. The 'Tasks' section is expanded, showing a list of tasks:

- workflows.classifier\_evaluate\_workflow.analyze\_prediction\_results**  
(No description)  
inputs: `ground_truths (integer[]), predictions (float[][])`  
outputs: `result_blobs (file/blob[]), result_files_names (string[])`
- workflows.classifier\_evaluate\_workflow.evaluate\_on\_datasets**  
(No description)  
inputs: `evaluation_clean_mpblob (file/blob), evaluation_dirty_mpblob (file/blob), model (file/blob)`  
outputs: `ground_truths_out (integer[]), predictions_out (float[][])`
- workflows.classifier\_evaluate\_workflow.fetch\_model**  
(No description)  
inputs: `model (file/blob)`  
outputs: `model_blob (file/blob)`

The screenshot shows the 'Create New Execution' dialog for the task `workflows.classifier_evaluate_workflow.analyze_prediction_results`. The dialog is open over the task details page. The task details page shows the breadcrumb `development / workflows.classifier_evaluate_workflow.analyze_prediction_results` and a 'Launch Task' button. The dialog contains the following fields:

- Task Version:** A dropdown menu showing the version `b61e3eb5a7780fb3d6b9828c9b75f105abd67201`.
- ground\_truths (integer[])\*:** A text input field containing `[1,2,3]`.
- predictions (float[][])\*:** A text input field containing `[1.0,2.0]`.

At the bottom of the dialog, there are two buttons: 'Cancel' and 'Launch'.

# Browse through historical executions and results

PROJECT  
KubeconDemo 2019

---

Workflows

Tasks

← development / DataPreparationWorkflow

Launch Workflow

### Description

Prepares raw videos for training/evaluation workflows. It runs a map-style job to download, breaks down streams into frames and runs luminance algorithm to pick important frames.

### Schedules

This workflow has no schedules.

### Executions

Status Version Start Time Duration

EXECUTION ID	STATUS	START TIME	DURATION
<a href="#">ffa3fe27fbaa34efc8c2</a> Last run a year ago	SUCCEEDED	11/18/2019 11:18:30 PM UTC 11/18/2019 3:18:30 PM PST	1s <a href="#">View Inputs &amp; Outputs</a>
<a href="#">f0a5fde004c1d43b6807</a> Last run a year ago	SUCCEEDED	11/17/2019 8:53:06 AM UTC 11/17/2019 12:53:06 AM PST	1s <a href="#">View Inputs &amp; Outputs</a>

# Type system - Auto launch forms!

The screenshot displays the Flyte web interface with a modal window titled "Create New Execution" for a workflow named "DataPreparationWorkflow". The modal contains several input fields for configuration:

- Workflow Version:** A dropdown menu with the value "fake11".
- Launch Plan:** A dropdown menu with the value "DataPreparationWorkflow".
- sampling\_n\_clusters (integer):** A text input field containing the value "8".
- sampling\_random\_seed (integer):** A text input field containing the value "0".
- sampling\_sample\_size (integer):** A text input field containing the value "20".
- stream\_extension (string):** A text input field containing the value "avi".

At the bottom of the modal, there are two buttons: "Cancel" and "Launch".

The background interface shows a sidebar with "PROJECT KubeconDemo 2019", "Workflows", and "Tasks" sections. The main content area is partially visible, showing a "Launch Workflow" button and a table with a "DURATION" header and rows containing "1s" and "View Inputs & Outputs".

# Get Provenance and lineage information

←
SUCCEEDED
flytesnacks/development/TrainingWorkflow/vzwqngwyvo
View Inputs & Outputs
Relaunch

Domain <b>development</b>	Version <b>28-2</b>	Cluster <b>flytestaging</b>	Time <b>7/30/2020 11:56:17 PM UTC</b>	Duration <b>1m 34s</b>
------------------------------	------------------------	--------------------------------	--	---------------------------

**Nodes**
Graph

Status ▾
Start Time ▾
Duration ▾

NODE	STATUS	TYPE	START TIME
<b>train-data</b> __main___.get_train_data2	<span style="background-color: #27ae60; color: white; padding: 2px 5px; border-radius: 3px;">SUCCEEDED</span> ↻	Unknown Task	
<b>validation-data</b> get_validation_data	<span style="background-color: #27ae60; color: white; padding: 2px 5px; border-radius: 3px;">SUCCEEDED</span> ↻	Unknown Task	
<b>train-csv</b> transform_parquet_to_csv	<span style="background-color: #27ae60; color: white; padding: 2px 5px; border-radius: 3px;">SUCCEEDED</span> ↻	Python Task	
<b>validation-csv</b> transform_parquet_to_csv	<span style="background-color: #27ae60; color: white; padding: 2px 5px; border-radius: 3px;">SUCCEEDED</span> ↻	Python Task	
<b>train</b> xgboost_hpo_task2	<span style="background-color: #27ae60; color: white; padding: 2px 5px; border-radius: 3px;">SUCCEEDED</span> ↻	Unknown Task	

**train-data**
✕

main\_.get\_train\_data2

SUCCEEDED

↻ Output for this execution was read from cache.
   
[View source execution](#)

TYPE  
 Unknown Task

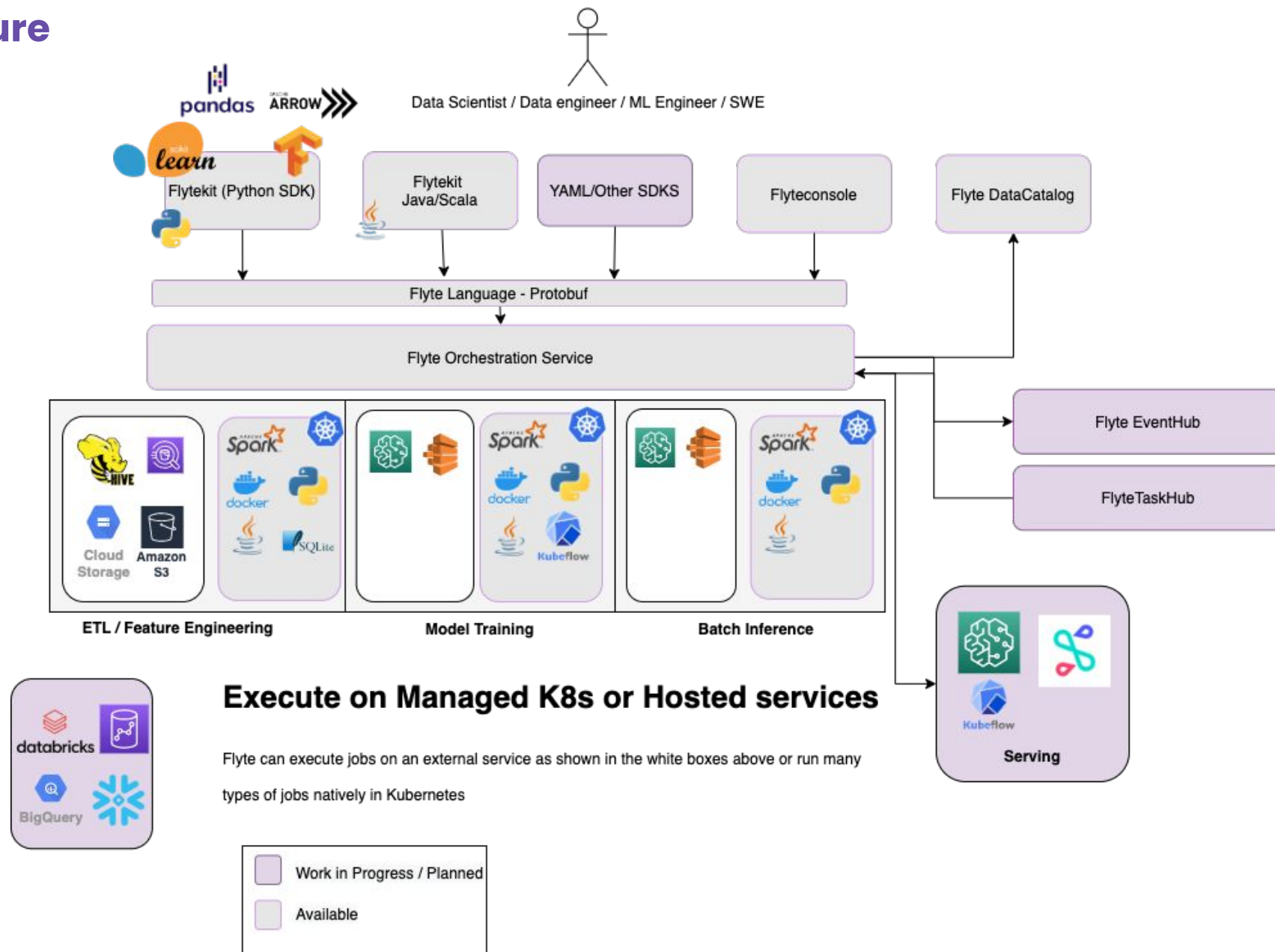
**Executions**
Inputs
Outputs
Task

Attempt 01  
succeeded

Logs  
 No logs found

started	(unknown)
run time	(unknown)

# Architecture





Flyte

# Differentiating attributes

- Fully Containerized
- Ergonomic and beautiful SDK's in python and Java/Scala
- Extensible Backend and SDK's
- Versioned and Auditable - record of all actions
- Horizontally Scalable & Battle tested - executed millions of pipelines per month
- Execute single task or a workflow, attach multiple schedules to a workflow
- Vertically integrated compute - serverless experience
- Deep understanding of data-lineage & provenance
- Operation visibility - cost, performance etc
- Pipelines portable across clouds

# Users



**& many more evaluations (Wolt, motional, gojek, universityhousing.nl, intuit, etc) in progress...**

# Contributions

There are more than 55+ unique contributors to Flyte across 16 repos ranging from Lyft, Lyft-Level5, Spotify, USU, Freenome, VMWare, Wolt and other companies. About 30 have been graduated to a permanent contributor status.

Recent major contributions by OSC

1. Flytekit JAVA (spotify) and Flytectl (OSC) are entirely open source contributed
2. Distributed Tensorflow & Pytorch Operator on K8s (entire open sourced)
3. Event Egress from Flyte (Spotify)
4. BigQuery, DataFlow and DataProc support GCP (Spotify & Freenome)
5. Flytekit plugins - Pandera, Pods etc (OSC)
6. Better onboarding experience (Freenome)
7. Documentation (Freenome and others)

**Future**

---

# Upcoming priorities

## Short Term (3-6 months)

- Role Based access controls
- Events Hub- Subscribe to Workflow / node events
- Centralized documentation
- Scalability improvements - support for extremely large DAGS - upwards of 20k nodes
- More integrations
  - AWS/GCP services
  - Flink on K8s
  - Data quality and access
  - Snowflake/Databricks etc
  - LF partners - ONNX, Feast, Mars etc
  - Serving integrations

## Long Term

- Reactive pipelines
- Complete portability across clouds
- Support for Streaming
- Low-code/No-code pipelines
- Data and model visualization plugins
- UI improvements
- Cost optimizations

# Why Contribute Flyte to LF-AI & Data?

## Neutral holding ground

- Vendor-neutral, Not for profit

## Growing community

- Increase visibility of Flyte through LF ecosystem
- Increase contributors by converting new & existing users
- Opportunities to collaborate with other hosted projects
- Flyte is unique system, which improves with collaborations and integrations. LF AI&Data - where integrations are encouraged is a natural home!

## 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

# TAC Vote on Project Proposal: Flyte

## **Proposed Resolution:**

The TAC approves the Flyte as an Incubation project of the LF AI & Data Foundation

## Next Steps

LF AI & Data staff will work with Flyte to onboard the project leading to the announcement of the project joining LF AI & Data

Explore potential integrations between the project and other LF AI & Data projects

Integrate the project with LF AI & Data operations



# LF AI & Data - General Updates

 LF AI & DATA

Machine Learning	Framework	Platform	Library	Framework	Platform	Library	Tool	Reinforcement Learning	Programming




Notebook Environment	Versioning	Store & Format	Operations	Stream Processing	SQL Engine	Feature Engineering	Visualization	Pipeline Management	Labeling and Annotation	Governance







Model	Benchmarking	Training	Parameter	Format & Interface	Marketplace	Workflow	Inference	Tool	Explainability	Adversarial	Bias & Fairness








Distributed Computing	Computing & Management	Interface	Security & Privacy	Natural Language Processing	Education

The LF AI & Data landscape explores open source projects in Artificial Intelligence and Data and their respective domains.

[l.fai.foundation](https://l.fai.foundation)

Machine Learning	Framework	Platform	Library	Framework	Platform	Library	Tool	Reinforcement Learning	Programming
		 LF AI & Data	 LF AI & Data						 LF AI & Data

Notebook Environment	Notebook Environment	Versioning	Store & Format	Operations	Stream Processing	SQL Engine	Feature Engineering	Visualization	Pipeline Management	Labeling and Annotation	Governance
			 LF AI & Data	 LF AI & Data  LF AI & Data  LF AI & Data <small>Incubating</small>	 LF AI & Data						 LF AI & Data

Model	Benchmarking	Training	Parameter	Format & Interface	Marketplace	Workflow	Inference	Tool	Explainability	Adversarial	Bias & Fairness
		 LF AI & Data	 LF AI & Data	 LF AI & Data	 LF AI & Data		 LF AI & Data		 LF AI & Data	 LF AI & Data	 LF AI & Data

Distributed Computing	Computing & Management	Interface	 The LF AI & Data landscape explores open source projects in Artificial Intelligence and Data and their respective sub-domains. <a href="https://lfaidata.foundation">lfaidata.foundation</a>				Security & Privacy	Natural Language Processing	Education
	 LF AI & Data	 LF AI & Data	 LF AI & Data	 LF AI & DATA Landscape	 LF AI & DATA			 LF AI & Data	 LF AI & Data  LF AI & Data <small>Incubating</small>

# Suggested Additions

## Project Key

**Yellow** = not in [Landscape](#), maybe should be added

## Programming

[Numpy](#)  
[Numba](#)  
[SciPy](#)  
[Dask](#)  
[Julia](#) (\*)  
[Python](#)  
[Rstudio](#)

## Notebooks

[Flyra](#)  
[I-python](#)  
[Jupyter Notebooks](#)  
[PixieDust](#)  
[Rmarkdown](#)

## Security & Privacy

[HE-Lib](#) (\*)  
[TensorFlow Privacy](#)  
[TF-Encrypted](#)

## Distributed Computing

*Management*  
[OpenShift](#)  
[Kubernetes](#)  
[Mesos](#)  
[Ranger](#)  
[Storm](#)

*Interface*  
[Sparklyr](#)  
[Toree](#)  
[Livy](#)  
[Spark-NLP](#)

## Data

*Versioning*  
[Pachyderm](#) (\*)

*Store & Format*  
[Alluxio](#)  
[Arrow](#)  
[Avro](#)  
[Delta Lake](#) (\*)

[Druid](#)  
[JanusGraph](#)  
[Parquet](#)  
[Ceph](#)

*Stream Processing*

[Flink](#)  
[Kafka](#)  
[Logstash](#) (\*)  
[FluentD](#) (\*)

*Relational DB*

[Postgres](#)  
[MySQL](#)  
[CouchDB](#)

*SQL Engine*  
[Presto](#) (\*)

*Visualization*

[Bokeh](#)  
[D3](#)  
[Plotly](#)  
[Facets](#)  
[Grafana](#)  
[Seaborn](#)  
[Superset](#) (\*)  
[TensorBoard](#)  
[Prometheus](#)

## Data

*Governance*  
[Egeria](#)  
[CLDA](#)

*Feature Engineering*  
[Tsfresh](#)

*Operations*  
[FEAST](#) (\*)  
[Amundsen](#) (\*)  
[Hive](#) (\*)  
[Snorkel](#) (\*)

*Pipeline Management*  
[Beam](#)

*Labeling & Annotation*  
[Vott](#) (\*)

*Exploration*  
[Hue](#)  
[Kibana](#)

## Machine Learning

*Framework*  
[LightGBM](#)  
[Mahout](#)  
[Ray](#) (\*)

*Platform*  
[Kubeflow](#)  
[H2O](#)  
[SystemML](#)  
[Mlflow](#) (\*)  
[Seldon](#) (\*)  
[Marvin-AI](#) (\*)

*Library*  
[Scikit-learn](#)  
[XGBoost](#)  
[cat-boost](#)  
[SparkML](#)

## Deep Learning

*Framework*  
[TensorFlow](#)  
[PyTorch](#)  
[MX-Net](#)

*Library*  
[Keras](#)

## Reinforcement Learning

[DeepMind Lab](#) (\*)  
[OpenAI Gym](#) (\*)

## Model

*Inference*  
[TensorRT](#)  
[TensorRT Inference](#)

*Benchmarking*  
[MLPerf](#)

*Training*  
[Horovod](#) (\*)

*Parameter*  
[HyperOpt](#)  
[Katib](#)

*Format & Interface*  
[ONNX](#)

*Marketplace*  
[MAX](#) (\*)

*Workflow*  
[Kubeflow Pipelines](#)  
[Tekton](#)

[Airflow](#) (\*)  
[Nifi](#) (\*)  
[Argp](#) (\*)  
[Mleap](#) (\*)  
[Volcano](#) (\*)

*Tool*  
[KFServing](#)  
[ONNX Runtime](#)  
[TorchServe](#) (\*)  
[Clipper](#) (\*)  
[MMS](#) (\*)

## Trusted AI

*Explainability*  
[AI Explainability 360](#)  
[Alibi](#) (\*)  
[LIME](#)  
[SHAP](#)

*Bias & Fairness*  
[AI Fairness 360](#)

*Adversarial Attacks*  
[Adversarial Robustness Toolbox](#)

## Natural Language Processing

[UIMA](#)  
[BERT](#)  
[Core NLP](#)  
[Lucene](#)  
[PyText](#)  
[Spacy](#)  
[Transformers](#) (\*)

*Education*  
[OpenDS4All](#)

## 2020 TAC Meetings Summary

Jan Feb Mar	16: Milvus (Zilliz)*	13: <i>MLOps Work (LF CD)</i>  27: <i>Collective Knowledge (Coral Reef)</i>	12: NNStreamer (Samsung)*  26: ForestFlow (?)*
Apr May Jun	9: <i>Trusted AI &amp; ML Workflow (LF)</i>  23: <i>Open Data Hub (Red Hat)</i>	7: Ludwig (Uber)*  21: <i>SnapML (IBM)</i>	4: <i>Trusted AI (AI for Good, Ambianic.ai, MAIEI)</i>  18: Fairness, Explainability, Robustness (IBM)*
Jul Aug Sep	16: <i>Mindspore (Huawei)</i>  30: Amundsen (Lyft)*	16: <i>Delta (Didi)</i> <b>16: Horovod (Uber/LF)**</b>  30: <i>ModelDB (?)</i> 30: <i>Egeria, OpenDS4All, BI&amp;AI (LF ODPI)</i>	10: SOAJS (HeronTech)* 10: Delta (Didi)* 24: FEAST (Gojek)* <b>24: Egeria, (LF ODPI)**</b> 24: OpenDS4All (ODPI)* 24: BI&AI Committee (ODPI)
Oct Nov Dec	8: <i>Fairness, Explainability, Robustness (LF)</i>  22: <i>OpenLineage (DataKins)</i> 22: <i>IDA (IBM/Salesforce)</i>	5: DataPractices.Org (WorldData/LF)* 5: <i>Kubeflow-On-Prem (Google, Arrikto/Intel)</i>  19: <i>OpenDS4All, DataPractices.Org, edX Ethical AI (LF)</i>	3: TBD - JanusGraph (LF)* 3: <i>TBD - RosaeGL (?)</i>  17: TBD – Seldon Core (Seldon)* <b>17: TBD – Pyro (Uber/LF)**</b>

(Entity)\* = incubating vote

**\*\* bold = graduate vote**

*Italics = invited project presentation*

## 2021 TAC Meetings Pipeline Summary

Jan Feb Mar	14: Data Lifecycle Framework (IBM)* 28: Tentative: Verse (Seldon)	11: MARS (Aliabab) 25: Flyte (Lyft)	11: Streams (IBM) 25: Tentative: Substra Framework
Apr May Jun	8: Adlik (ZTE)** 22: Kubeflow-On-Prem (Google, Arrikto, Intel)	?: Ray (Anyscale.io) ?: Pachyderm (Pachyderm) ?: DataHub (LinkedIn)	?: Common Knowledge (Code Reef) ?: Couler (Ant Financial)
Jul Aug Sep	?: KubeflowServing (Google, Arrikto, Seldon)	?: Kubeflow Pipeline (Google, Bloomberg)	?: Open Data Hub (Red Hat)
Oct Nov Dec	?: Vespa (Verizon Media)	?: Snorkle (Snorkle) ?: Plotly (DASH) ?: Mellody (Substra) ?: mloperator (Polyaxen) ?: SnapML (IBM)	?: PMML/PFA (DMG.org) ?: Mindspore, Volcano (Huawei) ?: TransmorgrifAI (Salesforce) ?: AIMET (Qualcomm) ?: Elyra-AI (IBM)

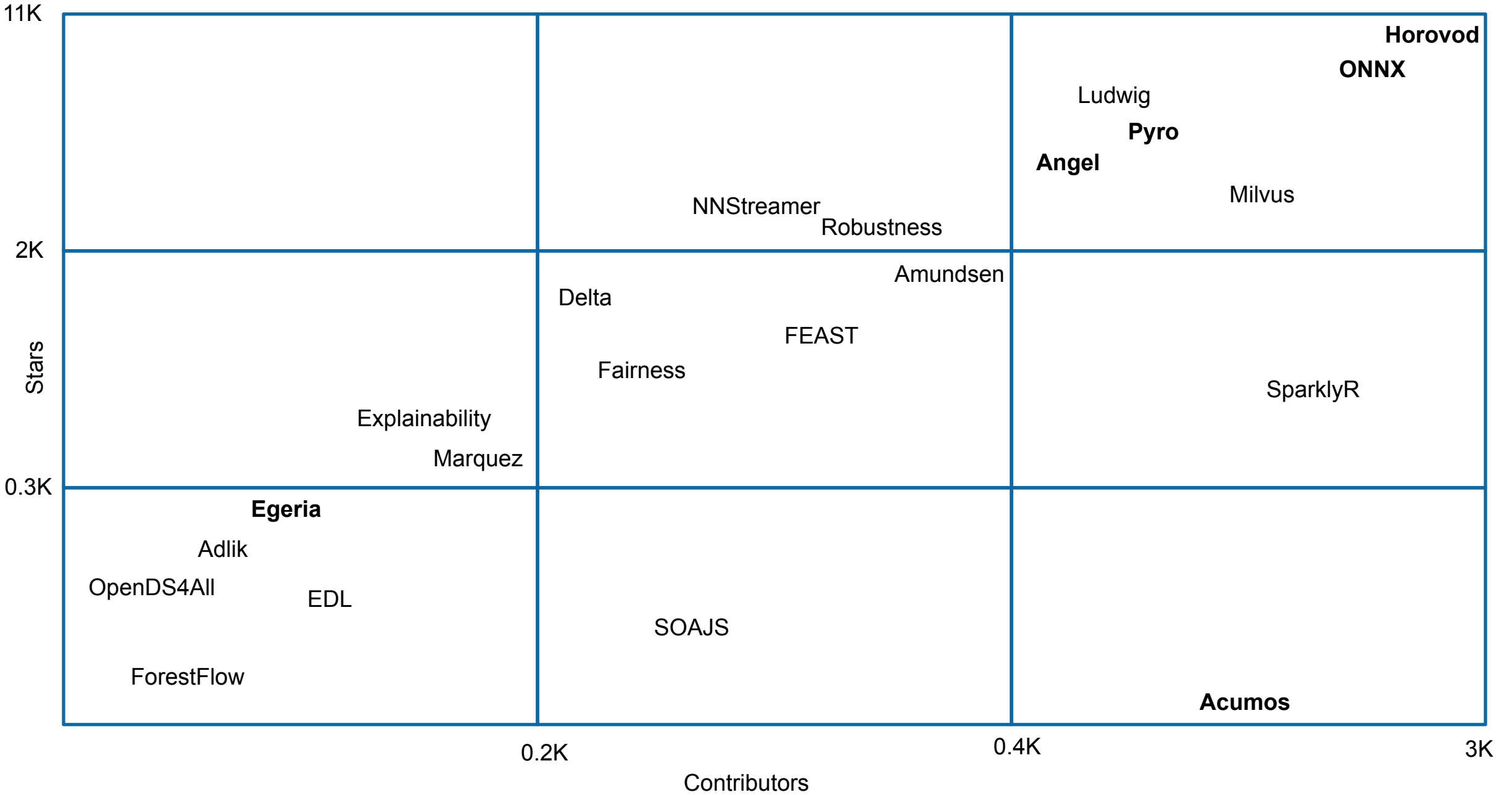
(Entity)\* = incubating vote

\*\* **bold** = graduate vote

*Italics* = invited project presentation

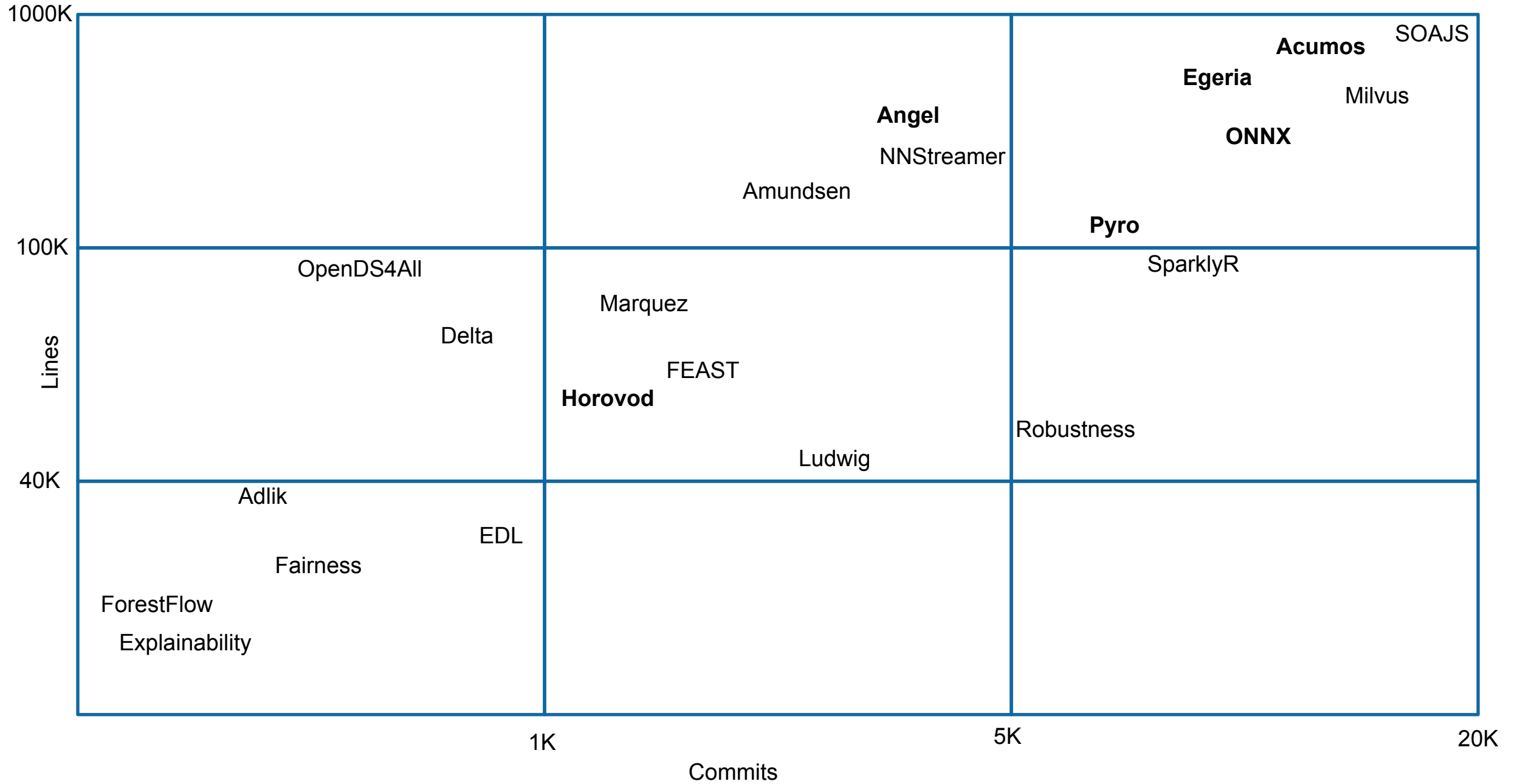
Getting to know the projects more

Data from November 23, 2020 – Stars and Contributors





Data from November 23, 2020 – Lines of Code and Commits



# Looking to host a project with LF AI & Data

- › Hosted project stages and life cycle:  
<https://lfaidata.foundation/project-stages-and-lifecycle/>
- › Offered services for hosted projects:  
<https://lfaidata.foundation/services-for-projects/>
- › Contact:  
Jim Spohrer (TAC Chair) and Ibrahim Haddad (ED, LF AI & Data)

# Promoting Upcoming Project Releases

We promote project releases via a blog post and on LF AI & Data [Twitter](#) and/or [LinkedIn](#) social channels

For links to details on upcoming releases for LF AI & Data hosted projects visit the [Technical Project Releases wiki](#)

If you are an LF AI & Data hosted project and would like LF AI & Data to promote your release, reach out to [pr@lfai.foundation](mailto:pr@lfai.foundation) to coordinate in advance (min 2 wks) of your expected release date.

# Note on quorum

As LF AI & Data is growing, we now have 16 voting members on the TAC.

TAC representative - please ensure you attend the bi-weekly calls or email Jacqueline/Ibrahim to designate an alternate representative when you can not make it.

We need to ensure quorum on the calls especially when we have items to vote on.

# Updates from Outreach Committee

# Upcoming Events

- › Upcoming Events
  - › Visit the [LF AI & Data Events Calendar](#) or the [LF AI & Data 2021 Events wiki](#) for a list of all events
  - › To participate visit the [LF AI & Data 2021 Events wiki page](#) or email [info@lfaidata.foundation](mailto:info@lfaidata.foundation)
  
- › Please consider holding virtual events

To discuss participation, please email [events@lfaidata.foundation](mailto:events@lfaidata.foundation)

# Upcoming Events

<https://lfaidata.foundation/events/>

- **March 24, 2021 - ONNX Community Virtual Meetup**
  - a. **Wednesday @ 5:00 pm - 8:00 pm PT USA**  
**Thursday @ 8:00am - 11am China Time**  
[LF AI Day: ONNX Community Virtual Meetup – March 2021](#)  
**(Virtual - Free - Asia-friendly time – Host Ti Zhou - Baidu)**
  
- **Sept 29 - Oct 1, 2021 - OSS Global**
  - a. **Mini-Summit, Booth, Track**

# LF AI PR/Comms

- › Please follow LF AI & Data on [Twitter](#) & [LinkedIn](#) and help amplify news via your social networks - Please retweet and share!
  - › Also watch for news updates via the tac-general mail list
  - › View recent announcement on the [LF AI & Data Blog](#)
- › Open call to publish project/committee updates or other relevant content on the [LF AI & Data Blog](#)
- › To discuss more details on participation or upcoming announcements, please email [pr@lfaidata.foundation](mailto:pr@lfaidata.foundation)



# Call to Participate in Ongoing Efforts

 **OLF** AI & DATA

# Trusted AI

- › **Leadership:**  
Animesh Singh (IBM), Souad Ouali (Orange), and Jeff Cao (Tencent)
- › **Goal:** Create policies, guidelines, tooling and use cases by industry
- › **Slack conversation channel:**  
#trusted-ai-committee  
<https://lfaifoundation.slack.com/archives/CPS6Q1E8G>
- › **Github:**  
<https://github.com/lfai/trusted-ai>
- › **Wiki:**  
<https://wiki.lfai.foundation/display/DL/Trusted+AI+Committee>
- › **Email lists:**  
<https://lists.lfaidata.foundation/g/trustedai-committee/>
- › **Next call:** Monthly alternating times  
<https://wiki.lfai.foundation/pages/viewpage.action?pageId=12091895>

# ML Workflow & Interop

- › **Leadership:**  
Huang “Howard” Zhipeng (Huawei)
- › **Goal:**  
Define an ML Workflow and promote cross project integration
- › **Slack conversation channel:**  
#ml-workflow  
<https://lfaifoundation.slack.com/archives/C011V9VSMQR>
- › **Wiki:**  
<https://wiki.lfaidata.foundation/pages/viewpage.action?pageId=10518537>
- › **Email lists:**  
<https://lists.lfaidata.foundation/g/mlworkflow-committee>
- › **Next call:** Monthly check calendar/slack  
<https://wiki.lfai.foundation/pages/viewpage.action?pageId=18481242>

# BI & AI

- › **Leadership:**  
Cupid Chan (Index Analytics)
- › **Goal:** Identify and share industry best practices that combine the speed of machine learning with human insights to create a new business intelligence and better strategic direction for your organization.
  
- › **Slack conversations channel:**  
**#bi-ai-committee**  
<https://lfaifoundation.slack.com/archives/C01EK5ND073>
- › **Github:**  
<https://github.com/odpi/bi-ai>
- Wiki:**  
<https://wiki.lfaidata.foundation/pages/viewpage.action?pageId=35160417>
- Email lists:**  
<https://lists.lfaidata.foundation/g/biai-discussion>
- Next call:** Monthly community call TBD

# Ongoing effort to create AI Ethics Training

Initial developed course by the LF: Ethics in AI and Big Data - published on edX platform:

<https://www.edx.org/course/ethics-in-ai-and-big-data>

The goal is to build 2 more modules and package all 3 as a professional certificate - a requirement for edX

- › **To participate:**  
<https://lists.lfaidata.foundation/g/aiethics-training>

# Upcoming TAC Meetings

# Upcoming TAC Meetings (Tentative)

- ›
- › Mar 11: Sandbox project proposal - RosaeNLG
- › Mar 25: Substra Foundation
- › April 8: Adlik (ZTE)
- › April 22: TBD
- › May 6: All project updates

›  
Please send agenda topic requests to  
[tac-general@lists.lfaidata.foundation](mailto:tac-general@lists.lfaidata.foundation)

# 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>

# Open Discussion



# Mission

To build and support an open community and a growing ecosystem of open source AI, data and analytics projects, by accelerating innovation, enabling collaboration and the creation of new opportunities for all the members of the community

# Legal Notice

- › The Linux Foundation, The Linux Foundation logos, and other marks that may be used herein are owned by The Linux Foundation or its affiliated entities, and are subject to The Linux Foundation's Trademark Usage Policy at <https://www.linuxfoundation.org/trademark-usage>, as may be modified from time to time.
- › Linux is a registered trademark of Linus Torvalds. Please see the Linux Mark Institute's trademark usage page at <https://lmi.linuxfoundation.org> for details regarding use of this trademark.
- › Some marks that may be used herein are owned by projects operating as separately incorporated entities managed by The Linux Foundation, and have their own trademarks, policies and usage guidelines.
- › TWITTER, TWEET, RETWEET and the Twitter logo are trademarks of Twitter, Inc. or its affiliates.
- › Facebook and the “f” logo are trademarks of Facebook or its affiliates.
- › LinkedIn, the LinkedIn logo, the IN logo and InMail are registered trademarks or trademarks of LinkedIn Corporation and its affiliates in the United States and/or other countries.
- › YouTube and the YouTube icon are trademarks of YouTube or its affiliates.
- › All other trademarks are the property of their respective owners. Use of such marks herein does not represent affiliation with or authorization, sponsorship or approval by such owners unless otherwise expressly specified.
- › The Linux Foundation is subject to other policies, including without limitation its Privacy Policy at <https://www.linuxfoundation.org/privacy> and its Antitrust Policy at <https://www.linuxfoundation.org/antitrust-policy>, each as may be modified from time to time. More information about The Linux Foundation's policies is available at <https://www.linuxfoundation.org>.
- › Please email [legal@linuxfoundation.org](mailto:legal@linuxfoundation.org) with any questions about The Linux Foundation's policies or the notices set forth on this slide.



