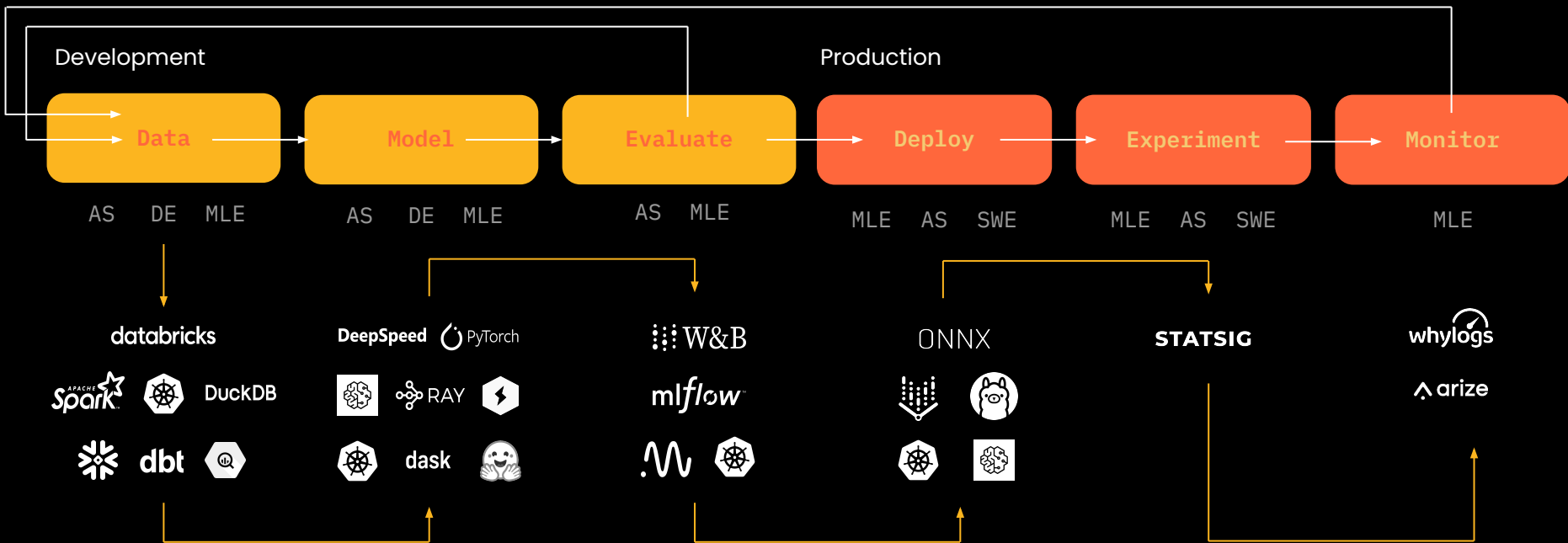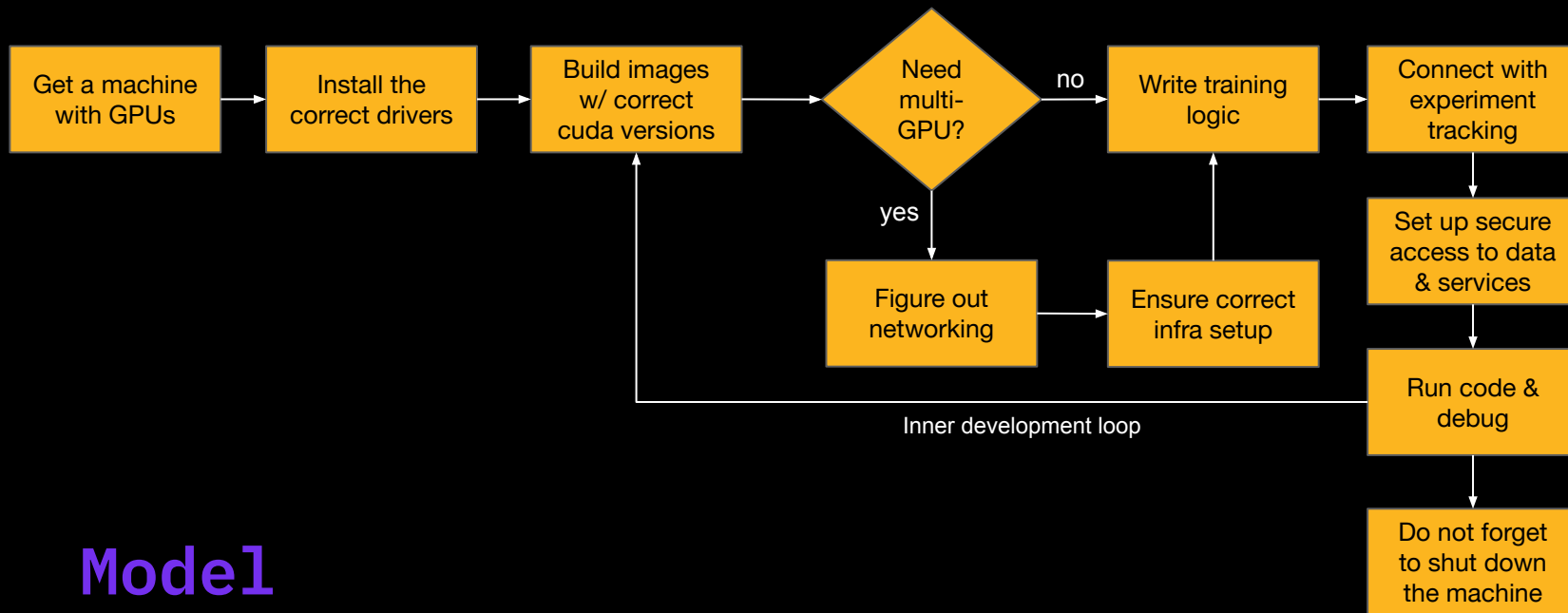# UNION

Orchestrate Your AI

# In the Era of AI every company must become an AI company

# But, AI product development is **chaotic**
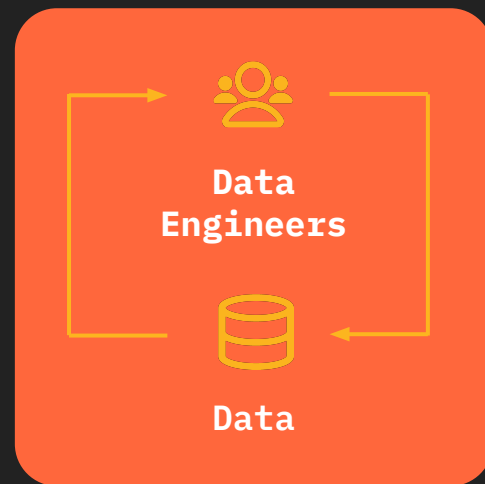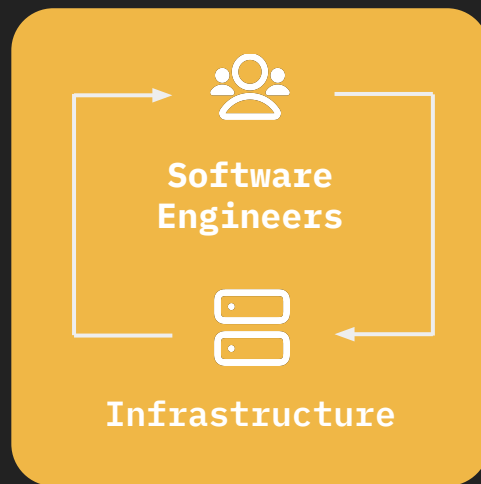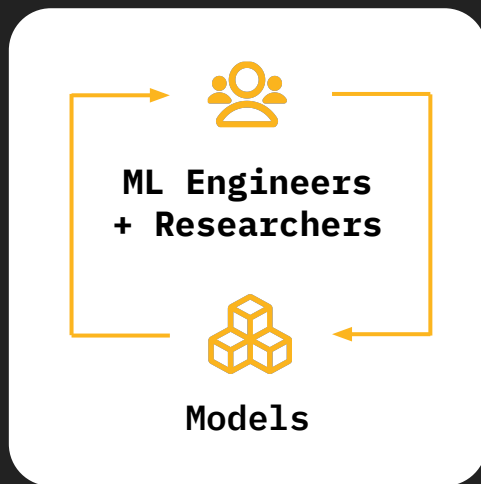
# And, solving this problem is hard

# Because teams are often siloed.

# Companies need a unified platform to create AI products

**Data Engineer / Software Engineer**

```python
# This code should run on spark

def transform(s: datetime, t: timedelta)
-> pyspark.DataFrame:
    sc = pyspark.Context()
    ...
    return df
```

**ML Engineer / Data Scientist / Researcher**

```python
# This code should run on one or + GPUs

def train(df: pd.DataFrame, hp:
TrainerArgs) -> nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
```

## Orchestrator

| Translate | Communicate | Translate |

**Infrastructure, Models, + Data**

# Union brings engineering, ops, + data science together... into **one efficient AI team.**

# Open source + closed source, ensures flexibility without lock-in

# Flyte open source to unlock AI

## One platform for your ML pipelines

**Build**

Write code in Python and test locally

**Iterate**

Rapidly ship to remote with a tight inner loop

**Scale**

Scale out and leverage multiple frameworks

**Visualize**

Track experiments, view plots, and use results

**Deploy**

CRON jobs and API-based executions

**Monitor**

Notifications via slack, email + pagerduty

**Extend**

Integrate with 2p and 3p external services

Flyte

# Union supercharges & extends Flyte



**Fully managed & secure**
Leverage a robust platform that meets rigorous standards for security, compliance, and operational reliability - in your cloud.

**Supercharged performance**
Run complex AI workloads with unparalleled performance, scale, and efficiency. (25x faster)

**Enhanced developer experience**
Shorten the development loop from hours to seconds while writing production-ready code.

**More efficient**
Boost ROI by enabling teams to access the resources they need while sharing underlying infrastructure.

**SLAs + customer success from the team that built Flyte**

# Our mission is to make AI products reliable, secure, + easy

# Our mission is to make AI products reliable, secure, + easy

## The Union AI Fabric

Teams collaborate to create AI applications

API

**Orchestration Engine**

| Data | → | Model | → | Evaluate | → | Deploy | → | Experiment | → | Monitor |

Inputs / Data

Compute

aws

databricks

W&B

# Workflows on Union are...

| Reproducible | Programmable | Composable | Interoperable |
| --- | --- | --- | --- |
| Scalable | Production Ready | Reliable | Efficient |

# Workflows on Union are...

## Reproducible

Automatic Versioning
Containerization
Data Immutability
Durable State & Results

## Programmable

Declarative Infra
Declarative data flow
Type-Aware
Pythonic

## Composable

Reusable components
Reusable data
Heterogeneous workloads
Accurate caching

## Interoperable

Integrations
API-driven development
Framework agnostic
Multi-language

## Scalable

Local-Remote parity
Multi-tenant
Integrated compute
Multi-cluster

## Prod-ready

Scheduling
Notifications
Observability
Dev/prod isolation

## Reliable

Retries
Checkpointing
Failure Recovery
Multi-AZ

## Efficient

Ephemeral Compute
Spot Instances
Checkpointing
Fractional GPUs

# Demo

# Demo: Stable Diffusion Fine Tuning

**Training**
Train LORA adapters on Multiple GPUs

**Optimize the model for deployment**
Use ONNX + Tensorrt to optimize the model for Nvidia hardware

**Use built-in Sagemaker deployment integration**
Automatically deploy the model

❏ Simple development and hardware targeting
❏ Caching (save money and time)
❏ Automated container builds
❏ UI based triggering
❏ Model Registry and Model cards
❏ Extensible integration system - Sagemaker
❏ Sharable building blocks - reuse optimization code!

# Demo: Spark three ways (local, open-source on k8s, and Databricks)

**Use Spark**
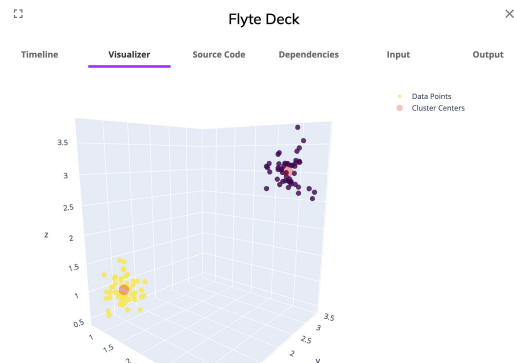Use spark local, or on K8s or on Databricks (and others)

**Simplified dependency management**
Automatically build containers and manage clusters

**Visualize**
Visualize data using FlyteDecks and debug Spark using Spark History Server

- ❏ Test spark code locally
- ❏ Declarative infrastructure
- ❏ Automated container builds
- ❏ Extensible - Databricks vs others
- ❏ Visualize data

# Demo: Integrate with Airflow

**Trigger from Airflow**
Seamless integration of Union Workflows / Tasks using [FlyteAirflowProvider](FlyteAirflowProvider)

**Migrate Airflow pipelines to Union**
Use most existing Airflow Operators and let Union handle the scaling

**Interoperate**
Use Flyte features like caching, data movement, type-safety and interoperate with Airflow operators



*"We were able to **save nine months of engineering time** by avoiding any code changes, and simply **lifting and shifting** our Airflow code and running it with Union."* — Shih-Gian Lee, Senior Machine Learning Engineer, Porch.

# Demo: Embed Wikipedia

**Scale**
Scale to multiple GPU's and cpu's efficiently

**Simple**
Code is simple, no need to learn complicated frameworks

**Efficient**
Native caching and high performance

**Model and Reuse**
Model the workflow as small tasks and reuse them

WIKIPEDIA

Shard into partitions

encode  encode  encode  encode

# Demo: Embed PDFs fast

**Scale**
Scale to multiple GPU's and cpu's efficiently
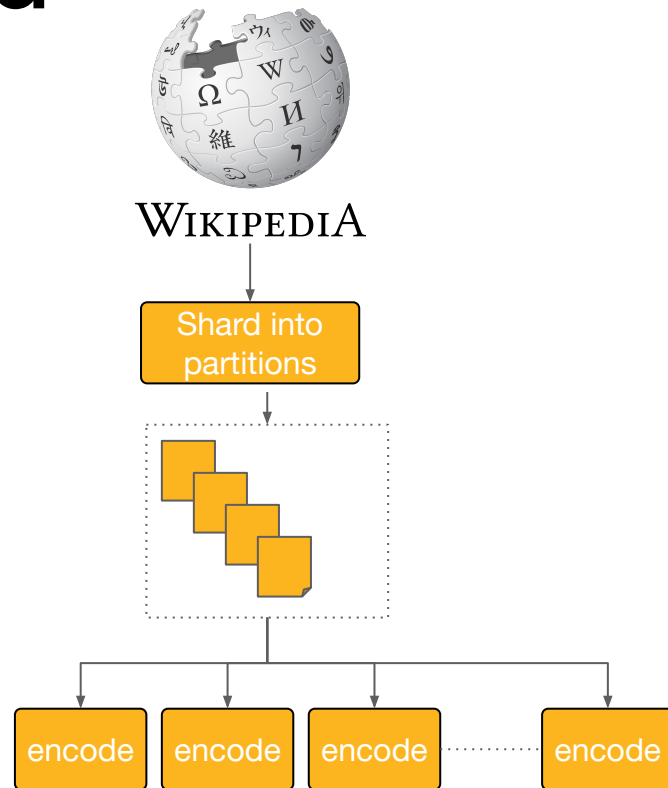
**Simple**
Code is simple, no need to learn complicated frameworks

**Efficient**
Native caching and high performance

**Model and Reuse**
Model the workflow as small tasks and reuse them

**Capture artifacts**
Capture the lineage and artifacts

Split and embed pdf

Split and embed pdf

One physical container (Actor)

# Customer Innovation

**Training & Fine Tuning**

Confidently run large-scale training or fine-tuning on GPU clusters across clouds and on-premise

**Data Processing**

Seamlessly connect to your data stack. Focus on data, not infrastructure.

**Near-Line Inference**

Deliver high-throughput, reliable, and fault-tolerant inference for production AI apps

**Generative AI & LLMs**

Take generative AI applications to production faster

**Bioinformatics & Pharma**

Effortlessly carry out scientific computing workflows with must-have features out of the box

# The standard AI orchestration platform

Spotify: Luigi (2013) -> **Migrated to Flyte → Union**
LinkedIn: Pro-ML (2018) -> **Migrated to Flyte**
Flipkart: Hunch (2018) -> **Working with them**
Stripe: Railyard (2019) -> **Migrated to Flyte**
Gojek: Machine Learning Platform (2019) -> **Migrated to Flyte**
Lyft: Flyte (2020) -> **:)**
DoorDash: ML platform (2020) -> **Doordash EU migrated to Flyte**

# Deployment options

- **Managed by Union**
- **Managed by Flyte OSS Core + Union features**
- **Managed by Customer**

Union BYOC
(Bring your own cloud)

Union Control Plane

Support

Optimized Core

Extensions

BYOC

BYOC

Cloud native
enterprises

# UNION

# Deployment options

- **Managed by Union**
- **Managed by Flyte OSS Core + Union features**
- **Managed by Customer (Planned)**

## Union Serverless

| Union Control Plane |
| Support | Optimized Core | Extensions |
| Serverless Compute |

Small teams
+ individuals

## Union BYOC
(Bring your own cloud)

| Union Control Plane |
| Support | Optimized Core | Extensions |
| BYOC | BYOC |

Cloud native
enterprises

## Union On Premise

| Union Control Plane |
| Support | Optimized Core | Extensions |
| On Premise Compute |

Enterprises with
on-prem requirements

# Pricing

Utilize Union and only pay for the resources managed by the platform. Benefit from spot instances with recovery, enhanced caching, and the ability to scale to zero.

## Flyte

**Free**
(support options available)

DIY ML orchestration for teams with on-premise or bare metal deployments.

- Open source (Apache 2.0)
- Battle-tested at scale
- Deployable on-premise
- Vibrant community
- Union support plans available

## UNION Serverless

**Pay-as-you-go***
$30 in free credit

Ship your first production model in seconds without worrying about infrastructure. Ideal for individuals.

- Optimized and expanded Flyte
- Limited to 1 seat
- Scalable serverless environment
- Pay only for consumed resources
- Sign up instantly with GitHub

## UNION BYOC - Startup

**$500/mo + % of compute**
(on-demand retail pricing)

A secure and scalable platform ideal for small teams and early stage companies.

- Optimized and expanded Flyte
- Up to 5 seats
- Single cluster/cloud
- Purchase on AWS or GCP marketplace
- Standard bring-your-own-cloud (BYOC) deployment

## UNION BYOC - Enterprise

**Custom**
(w/ committed use discounts)

Built for enterprises that require large-scale, highly available, and customizable deployments.

- Optimized and expanded Flyte
- Unlimited seats
- Multi-cluster/cloud
- Purchase on AWS or GCP marketplace
- Customizable bring-your-own-cloud (BYOC) deployment

# ARCHIVE

# Use Cases

**Generative AI**
- Fine tuning
- RAG data ingestion
- Embedding
- Multimodal training & inference
Logos: Flawless, LinkedIn

**Finance/FinTech**
- AML (JPM)
- Fraud detection (Stripe)
- Time-series forecasting
- FP&A (Spotify)

**Geospatial**
- Satellite imagery
- Mapping
- Data Extraction

Logos: MethaneSAT, Muon Space

**Bioinformatics/Pharma**
- Protein Engineering
- Therapeutics
- Drug discovery
- Antibodies
- Compound discovery (Zymergen)

**Consumer**
- Recommendation Systems (HBO)
- Personalization (LinkedIn)

**Logistics**
- ETA
- Operational Research

Logos:
- Gojek

**Autonomy & Robotics**
- Computer Vision
- Perception
- SLAM
Logos:
- Tesla, Physical Intelligence, Toyota, Mercedes, Wayve, StackAV

**Retail**
- Churn prediction
- Pricing

Wallapop

# Use Cases

## Generative AI

+ Fine tuning
+ RAG data ingestion
+ Embedding
+ Multimodal training

## Finance / FinTech

+ Fine tuning
+ RAG data ingestion
+ Embedding
+ Multimodal training

## Geospatial

+ Fine tuning
+ RAG data ingestion
+ Embedding
+ Multimodal training

## Bioinformatics + Pharma

+ Fine tuning
+ RAG data ingestion
+ Embedding
+ Multimodal training

**Copy needs to be updated**

## Consumer

+ Fine tuning
+ RAG data ingestion
+ Embedding
+ Multimodal training

## Logistics

+ Fine tuning
+ RAG data ingestion
+ Embedding
+ Multimodal training

## Autonomy + Robotics

+ Fine tuning
+ RAG data ingestion
+ Embedding
+ Multimodal training

## Retail

+ Fine tuning
+ RAG data ingestion
+ Embedding
+ Multimodal training

# Why Flyte + Union

Spotify: Luigi (2013) -> **Migrated to Flyte**

Facebook: FBLearner Flow (2016) -> **We have been talking with them**

Uber: Michaelangelo (2017)

Google: TFX (2017) **-> I have a talk at Google today**

LinkedIn: Pro-ML (2018) -> **Migrated to Flyte**

Flipkart: Hunch (2018) -> **Working with them**

Netflix: Metaflow (2018)

Stripe: Railyard (2019) -> **Migrated to Flyte**

Pinterest: Galaxy (2019)

Gojek: Machine Learning Platform (2019) -> **Migrated to Flyte**

Lyft: Flyte (2020) -> **:)**

DoorDash: ML platform (2020) -> **Doordash EU migrated to Flyte**

# Rate limit velocity is impacted by org silos

ML Engineers
+ Researchers

Models

Software
Engineers

Infrastructure

Data
Engineers

Data

# But, AI product development is HARD

The ML lifecycle is iterative + collaborative with many different roles

- Applied Scientist (AS)
- Data Engineer (DE)
- ML Engineer (MLE)
- Software Engineer (SWE)

**Development** →

| Data ETL | Feature Engineering | Model Training | Model Evaluation |
|---|---|---|---|
| AS    DE    MLE | AS    DE    MLE | AS    MLE | AS    MLE |

**Production** →

| Deploy + Experiment | Predict | Monitor + Debug | Retrain Model |
|---|---|---|---|
| MLE    AS    SWE | MLE    AS    SWE | MLE | MLE |

# But, AI Product development is HARD

**AI Researchers
Data Science
Applied Research**

**Models**

**ML Eng/
Software
Engineers**

**Code**  `< / >`

**Platform
Engineers**

**Infrastructure**

**Data
Engineers**

**Data**

# But, AI product development is HARD

## AI Researchers & Data Scientists

**Models**

## ML & Software Engineers

**Code** < / >

## Platform Engineers

**Infrastructure**

## Data Engineers

**Data**

# Union products expand the Flyte system



Union
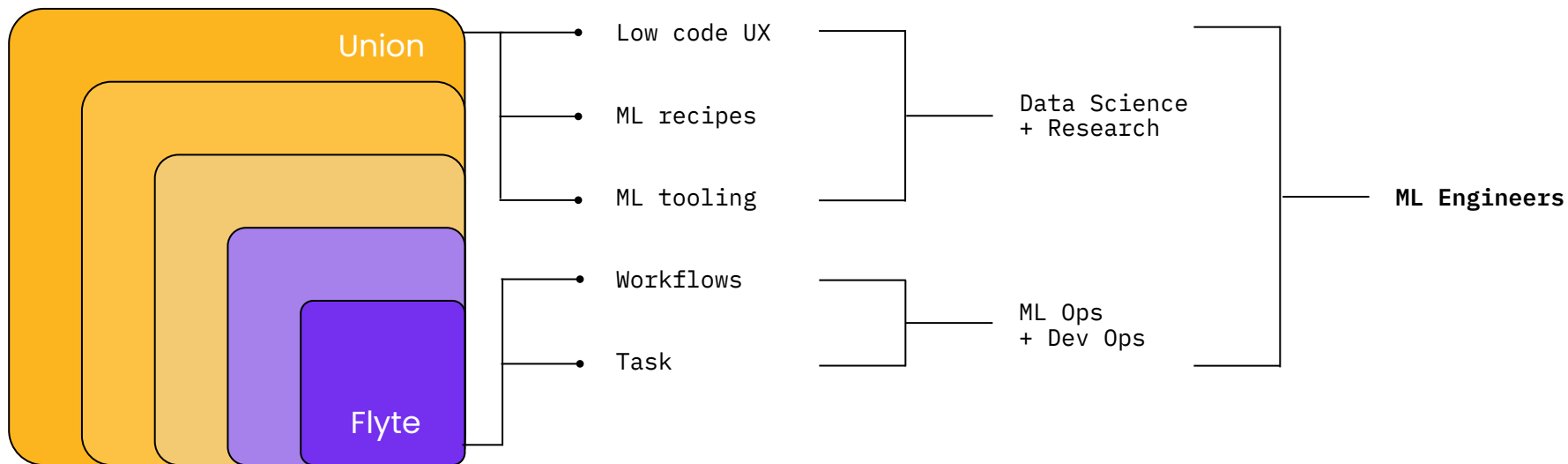Low code UX
ML recipes
ML tooling
Data Science + Research
Workflows
Task
ML Ops + Dev Ops
Flyte
ML Engineers
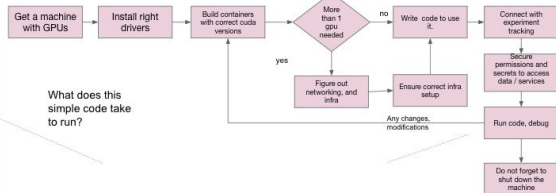
# Companies switch to Flyte & Union

# From this to this



## Today

```
# This code should run on one or more GPU's
def train(df: pd.DataFrame, hp: TrainerArgs) -> nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
```

### Applying this to a real-world example

WIP

What does this simple code take to run?

Get a machine with GPUs → Install right drivers → Build containers with correct cuda versions → More than 1 gpu needed → no → Write code to use it. → Connect with experiment tracking

yes → Figure out networking, and infra → Ensure correct infra setup → Secure permissions and secrets to access data / services

Any changes, modifications → Run code, debug → Do not forget to shut down the machine

## With Union

```
@task(task_config=Elastic(), limits=Resources(gpu=8))
def train(df: pd.DataFrame, hp: TrainerArgs) -> nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
```
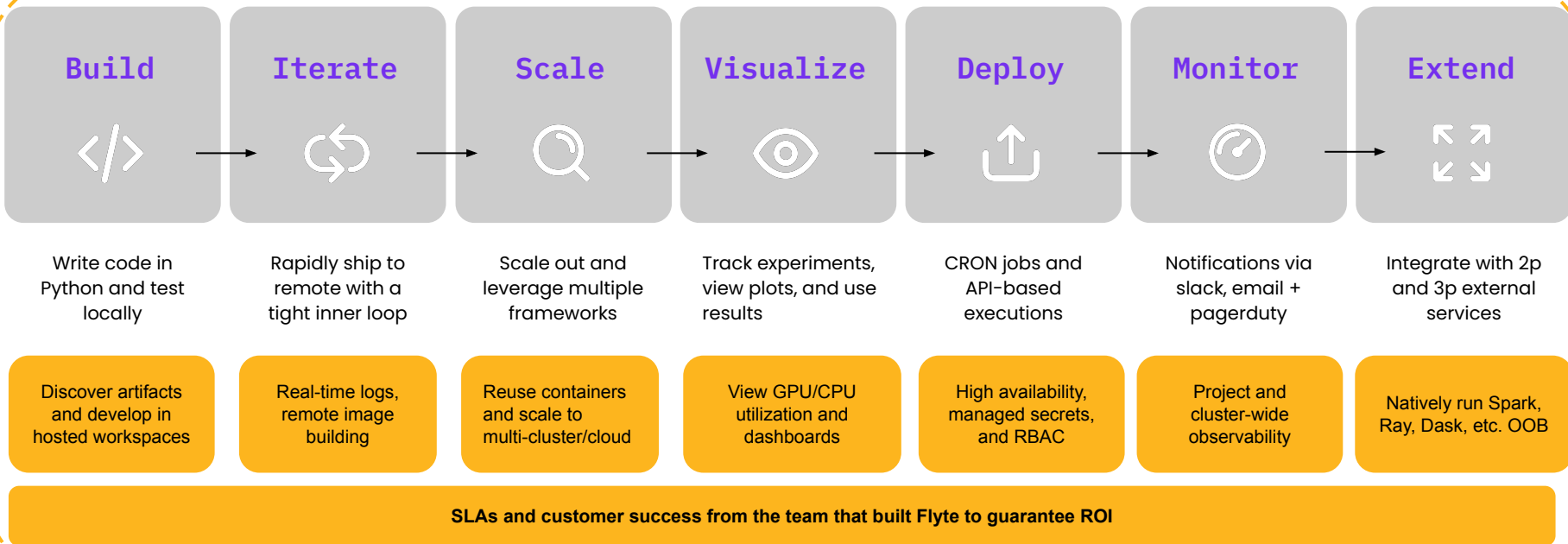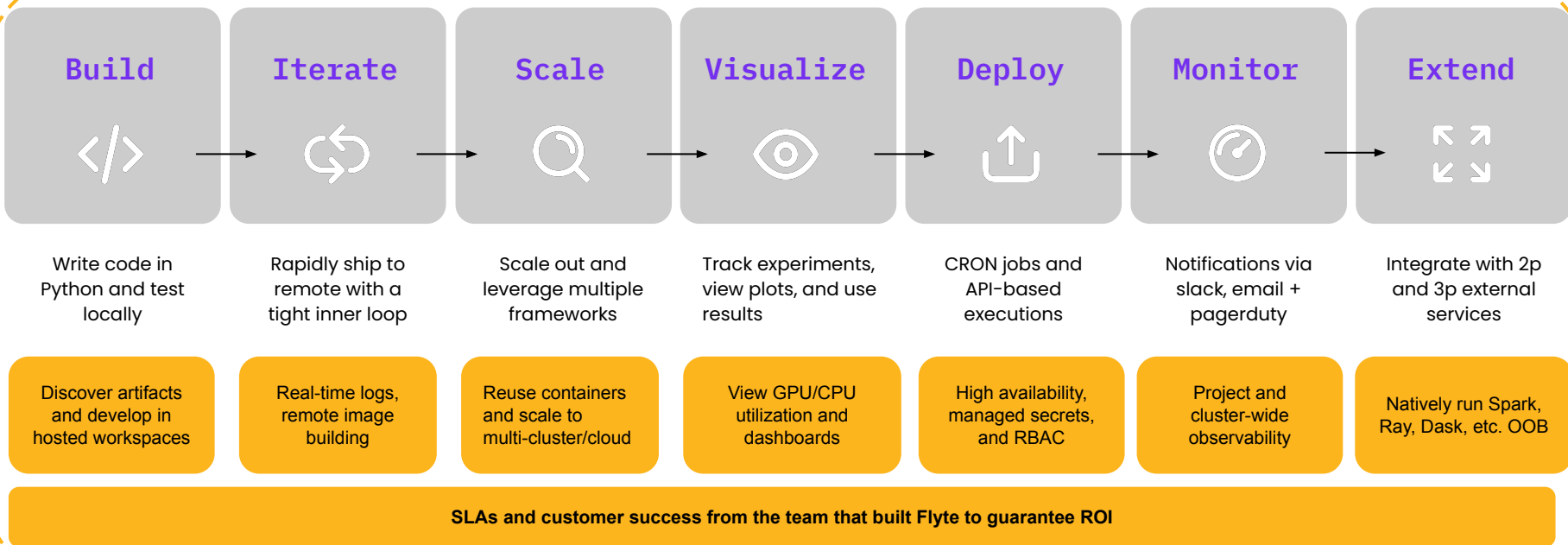
API

Union

# Union meets the requirements for AI dev

- **Reproducible:** Automatic Versioning, Containerization, Data Immutability, Durable State & Results
- **Programmable**: Declarative Infrastructure, Declarative Data Flow, Type-Safety & Type-Awareness, Pythonic (easy to adopt)
- **Composable**: Reusable Tasks & Workflows, Reusable Data (Artifacts), Heterogeneous Workflows
- **Interoperable**: Agents (Connect to 2p and 3p s), API-driven development, Framework-Agnostic (PyTorch, Tensorflow, etc)
- **Scalable**: Local-Remote Parity, Multi-tenant, Scalable Compute Fabric
- **Production-ready**: Scheduling, Notifications, Observability, Isolated Dev/Prod Environments
- **Reliable**: Retries, Checkpointing, Failure Recovery, No SPOF, Multi-AZ, Multi-Cluster
- **Efficiency**: Caching, Ephemeral Compute, Spot Instances, Fractional GPU

# Union is Flyte, supercharged

| Build | Iterate | Scale | Visualize | Deploy | Monitor | Extend |
|-------|---------|-------|-----------|--------|---------|--------|
| Write code in Python and test locally | Rapidly ship to remote with a tight inner loop | Scale out and leverage multiple frameworks | Track experiments, view plots, and use results | CRON jobs and API-based executions | Notifications via slack, email + pagerduty | Integrate with 2p and 3p external services |
| Discover artifacts and develop in hosted workspaces | Real-time logs, remote image building | Reuse containers and scale to multi-cluster/cloud | View GPU/CPU utilization and dashboards | High availability, managed secrets, and RBAC | Project and cluster-wide observability | Natively run Spark, Ray, Dask, etc. OOB |

**SLAs and customer success from the team that built Flyte to guarantee ROI**

# Union is Flyte, supercharged

| Build | Iterate | Scale | Visualize | Deploy | Monitor | Extend |
|-------|---------|-------|-----------|--------|---------|--------|
| Write code in Python and test locally | Rapidly ship to remote with a tight inner loop | Scale out and leverage multiple frameworks | Track experiments, view plots, and use results | CRON jobs and API-based executions | Notifications via slack, email + pagerduty | Integrate with 2p and 3p external services |
| Discover artifacts and develop in hosted workspaces | Real-time logs, remote image building | Reuse containers and scale to multi-cluster/cloud | View GPU/CPU utilization and dashboards | High availability, managed secrets, and RBAC | Project and cluster-wide observability | Natively run Spark, Ray, Dask, etc. OOB |

**SLAs and customer success from the team that built Flyte to guarantee ROI**

# Customer innovation

Spotify · TOYOTA · *LOCKHEED MARTIN* · LinkedIn · stripe · freenome

TESLA · amazon · Microsoft · NVIDIA · Cradle

## Training & Fine Tuning

Confidently run large-scale training or fine-tuning on GPU clusters across clouds and on-premise

## Data Processing

Seamlessly connect to your data stack. Focus on data, not infrastructure.

## Near-Line Inference

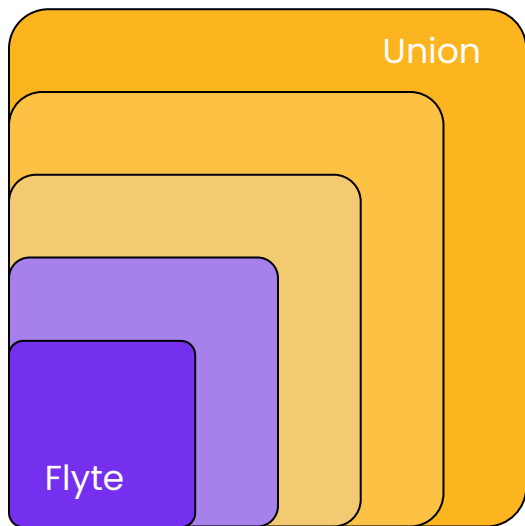Deliver high-throughput, reliable, and fault-tolerant inference for production AI apps

## Generative AI & LLMs

Take generative AI applications to production faster

## Bioinformatics & Pharma

Effortlessly carry out scientific computing workflows with must-have features out of the box

# Union supercharges and extends Flyte



**More efficient**
Boost ROI by enabling teams to access the resources they need while sharing underlying infrastructure.

**Better developer experience**
Shorten the development loop from hours to seconds while writing production-ready code.

**Supercharged performance**
Run complex AI workloads with unparalleled performance, scale, and efficiency.

**Fully managed & secure**
Leverage a robust platform that meets rigorous standards for security, compliance, and operational reliability - in your cloud.

**SLAs and customer success from the team that built Flyte**

# You need a system that abstracts sharing & Infrastructure

Data engineer / Software engineer

```
# This code should run on spark
def transform(s: datetime, t: timedelta) ->
pyspark.DataFrame:
    sc = pyspark.Context()
    ...
    return df
```

ML Engineer / Data Scientists / Researcher

```
# This code should run on one or more GPUs
def train(df: pd.DataFrame, hp: TrainerArgs) -> nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
```

## Orchestrator

Translate → Communicate → Translate

## Low-level infrastructure

# Orchestrate Your AI

Bring together ML, Platform, Data and Ops teams to create AI products efficiently

## Flyte, supercharged

All of the features in flyte, optimized for speed and enhanced for dynamic execution and managed K8s

## Unified workstreams

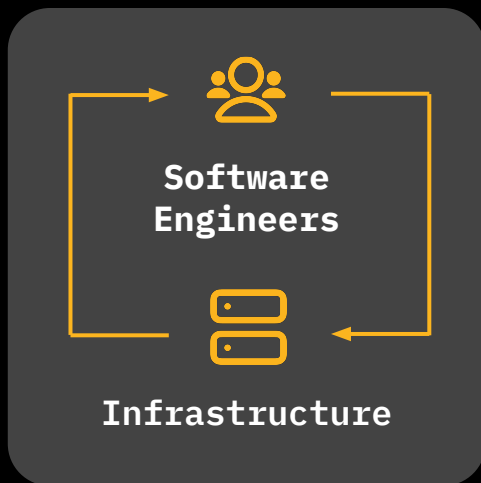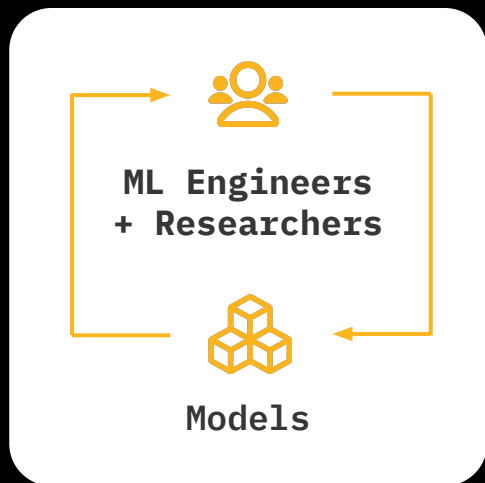Modern AI orchestration that joins teams to productionize AI apps, process and workflows

## Maximized AI ROI, derisked

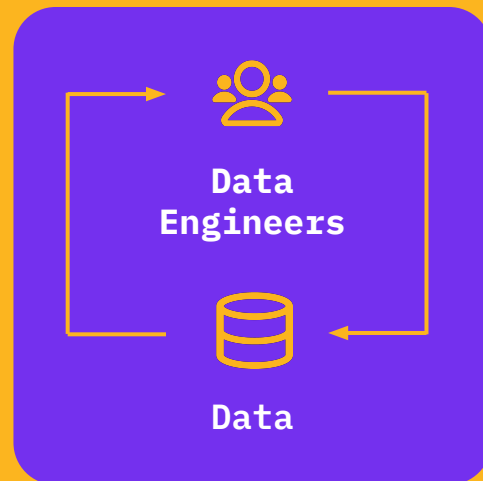Reduce operating costs with efficient resource management, while increasing velocity
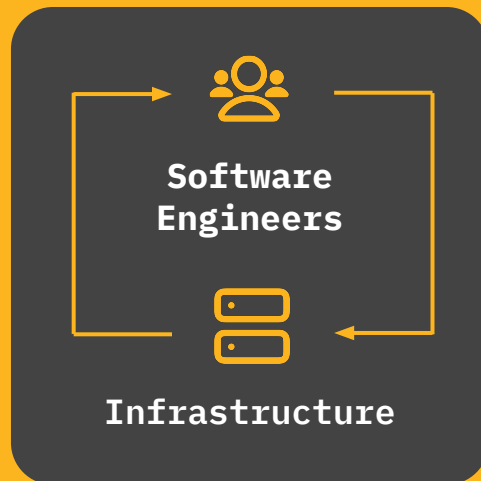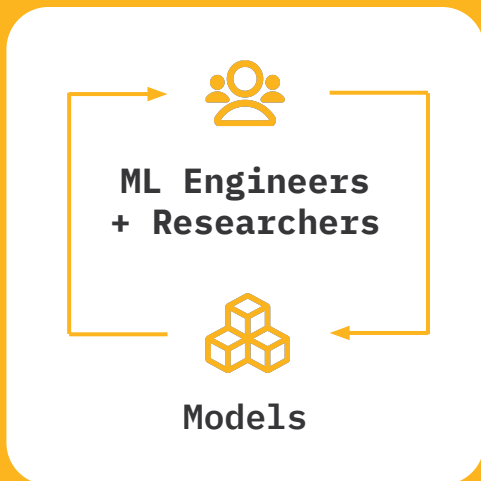
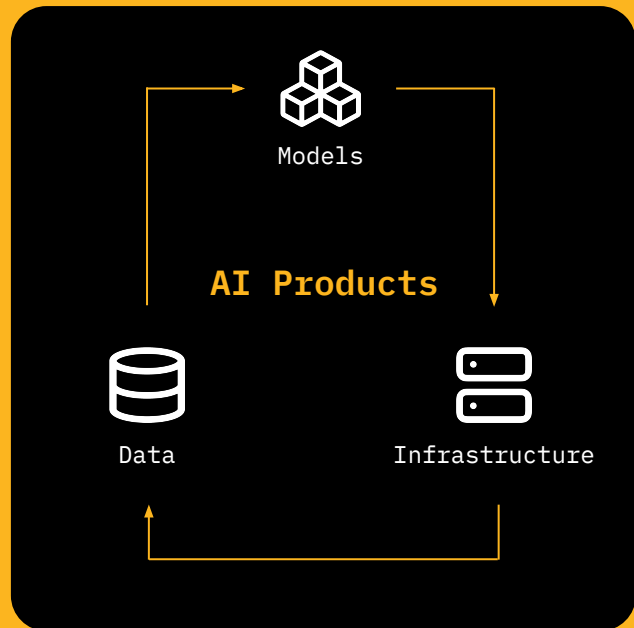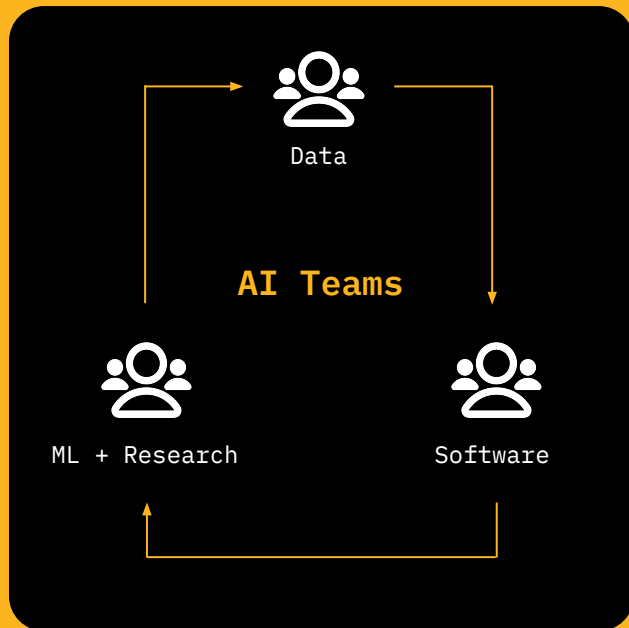All built on a foundation of trust

# Rate limit velocity is impacted by org silos

# AI product development is **hard + chaotic**



WIP

## Development

| Data ETL | Feature Engineering | Model Training | Model Evaluation |
|---|---|---|---|
| AS    DE    MLE | AS    DE    MLE | AS    MLE | AS    MLE |

## Production

| Deploy + Experiment | Predict | Monitor + Debug | Retrain Model |
|---|---|---|---|
| MLE    AS    SWE | MLE    AS    SWE | MLE | MLE |

# The standard AI orchestratoin platform

Spotify: Luigi (2013) -> **Migrated to Flyte → Union**
Facebook: FBLearner Flow (2016) -> **We have been talking with them**
Uber: Michaelangelo (2017)
Google: TFX (2017) **-> I have a talk at Google today**
LinkedIn: Pro-ML (2018) -> **Migrated to Flyte**
Flipkart: Hunch (2018) -> **Working with them**
Netflix: Metaflow (2018)
Stripe: Railyard (2019) -> **Migrated to Flyte**
Pinterest: Galaxy (2019)
Gojek: Machine Learning Platform (2019) -> **Migrated to Flyte**
Lyft: Flyte (2020) -> **:)**
DoorDash: ML platform (2020) -> **Doordash EU migrated to Flyte**

Storyline

AI product development is CHAOTIC
→ need to anchor on a customer story/challenge - all the different tools and systems, the constant loop
→ the fact that the models are iterative - this is new type of development

Solving this problem is Hard
→ because it is deep tech
→ and we have the issue of silo's and people/process

We are solving this problem today (SHOW THE HOW)
- Solve the silo problem
- Solve the deep tech problem

→ the result is a simple platform where you build your AI products (Union is the fabric)

# Applying this to a real-world example

```python
@task(task_config=Spark({"spark.executors": 4}))
# This code should run on spark
def transform(s: datetime, t: timedelta) ->
pyspark.DataFrame:
    sc = pyspark.Context()
    ...
    return df



@task(task_config=Elastic(),limits=Resources(gpu=8))

def train(df: pd.DataFrame, hp: TrainerArgs) ->
nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
    ...
```
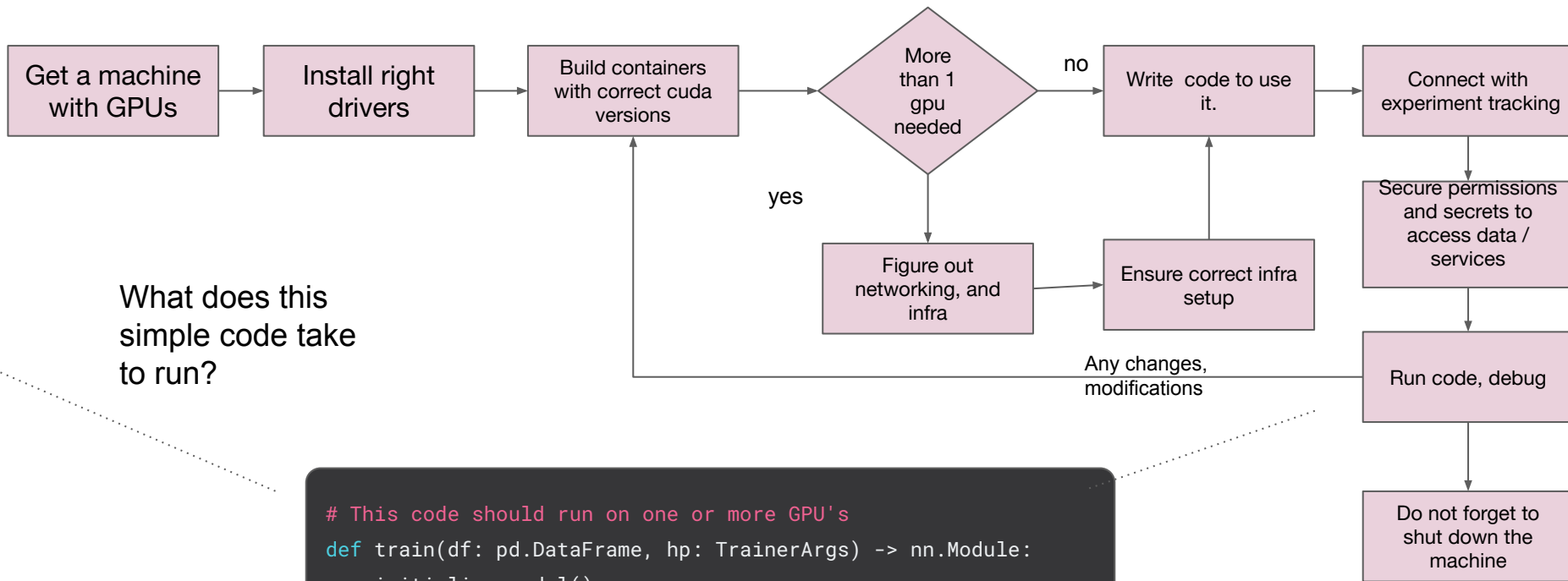
Declarative
infrastructure

Publish to
huggingface

Quantize / evaluate /
optimize

Run
predictions

Deploy

# Applying this to a real-world example



```
# This code should run on one or more GPU's
def train(df: pd.DataFrame, hp: TrainerArgs) -> nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
```

What does this simple code take to run?

Imagine every mle / data scientists has to do this. The cost, access management and scale?

# Real-world example
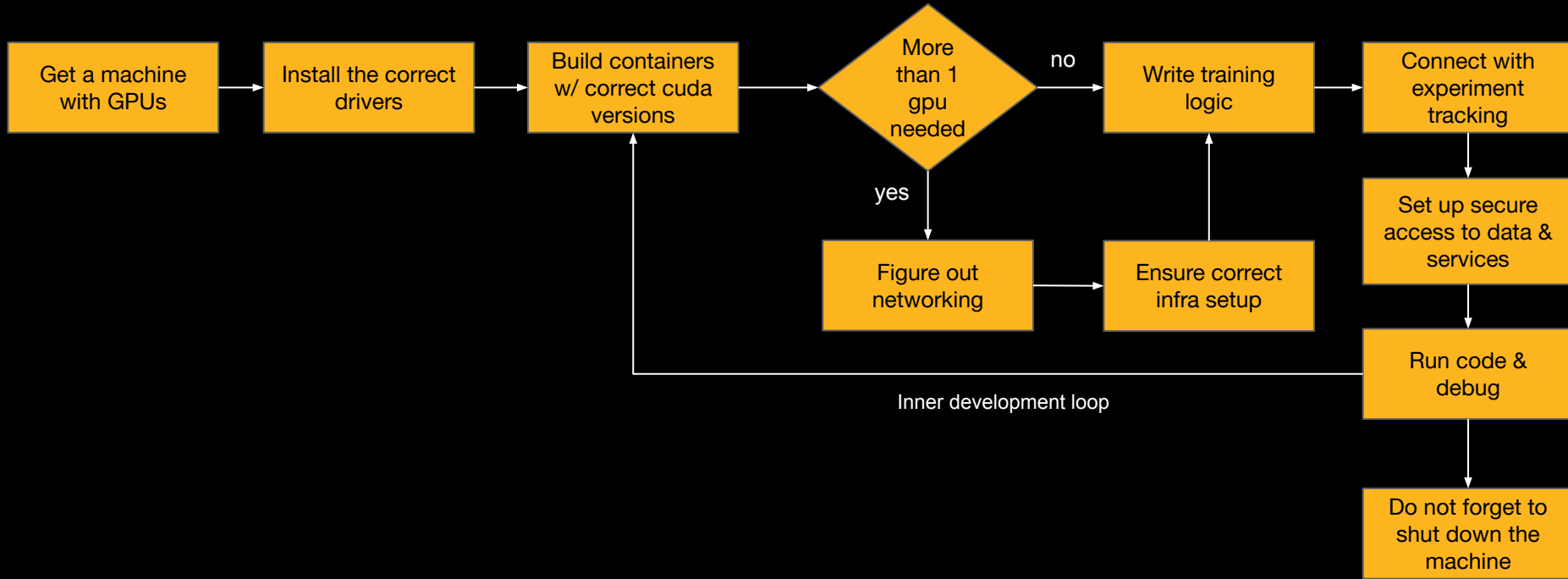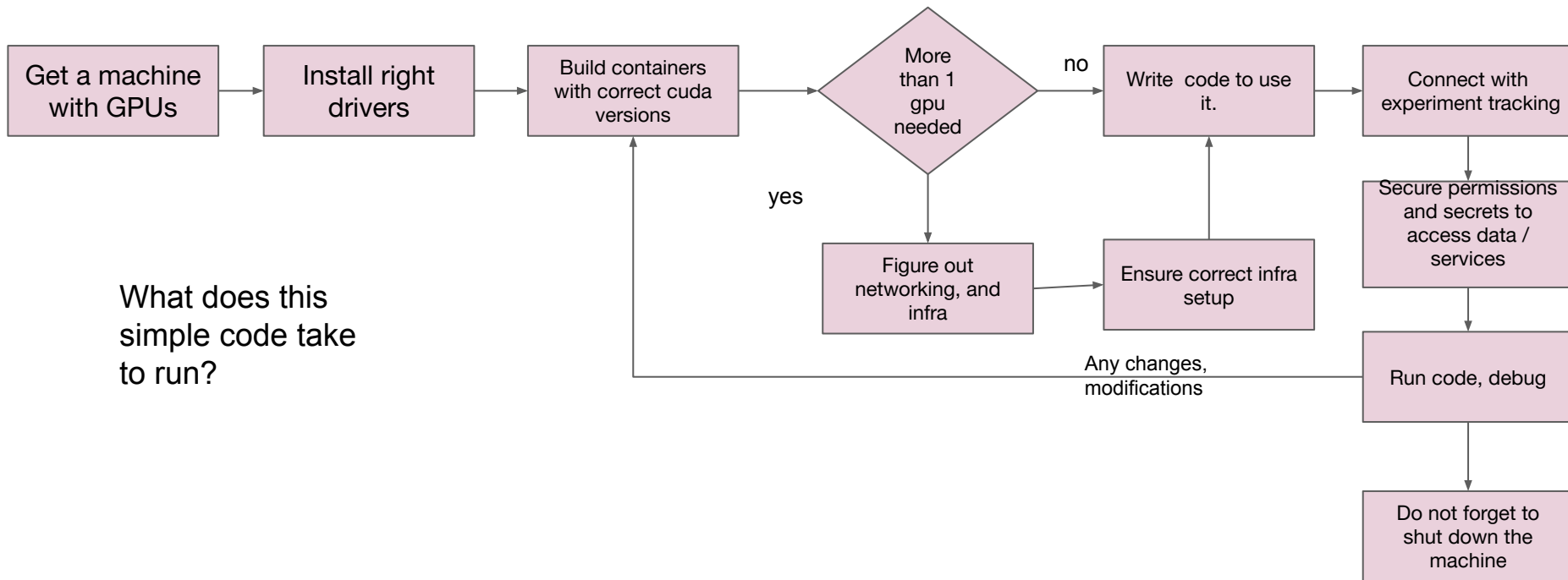
```
# This code should run on one or more GPU's
def train(df: pd.DataFrame, hp: TrainerArgs) -> nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
```

**WIP**

What does this
simple code take
to run?

Get a machine with GPUs → Install right drivers → Build containers with correct cuda versions → More than 1 gpu needed

More than 1 gpu needed —no→ Write code to use it. → Connect with experiment tracking

yes

Figure out networking, and infra → Ensure correct infra setup

Connect with experiment tracking → Secure permissions and secrets to access data / services → Run code, debug

Any changes, modifications

Run code, debug → Do not forget to shut down the machine

# Applying this to a real-world example

```python
@task(task_config=Spark({"spark.executors": 4}))
def transform(s: datetime, t: timedelta) -> pyspark.DataFrame:
    sc = pyspark.Context()
    ...
    return df
```

```python
@task(task_config=Elastic(), limits=Resources(gpu=8))
def train(df: pd.DataFrame, hp: TrainerArgs) -> nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
```

**Declarative infrastructure**

Quantize / evaluate / optimize

Publish to huggingface

Run predictions

Deploy

# You need a system that abstracts sharing + infrastructure (which is HARD)

**Data Engineer /**
**Software Engineer**

```
# This code should run on spark

def transform(s: datetime, t: timedelta) ->
pyspark.DataFrame:
    sc = pyspark.Context()
    ...
    return df
```

**ML Engineer /**
**Data Scientist / Researcher**

```
# This code should run on one or more GPUs

def train(df: pd.DataFrame, hp: TrainerArgs)
-> nn.Module:
    initialize_model()
    huggingFace.Trainer()
    ...
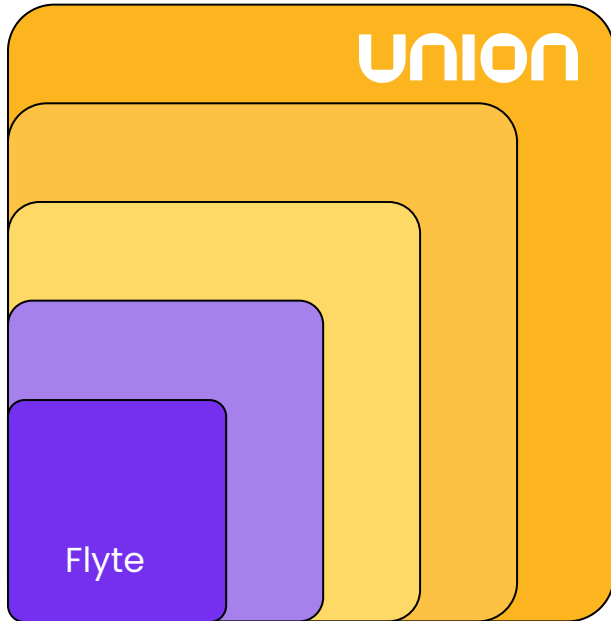```

**Orchestrator**

Translate → Communicate → Translate

Low level infrastructure

# Union supercharges + extends Flyte

UNION

Flyte

**More efficient**
Boost ROI by enabling teams to access the resources they need while sharing underlying infrastructure.

**Enterprise developer experience**
Shorten the development loop from hours to seconds while writing production-ready code.

**Supercharged performance**
Run complex AI workloads with unparalleled performance, scale, and efficiency.

**Fully managed & secure**
Leverage a robust platform that meets rigorous standards for security, compliance, and operational reliability - in your cloud.

**SLAs + customer success from the team that built Flyte**

# Our mission at Union: Make creating AI products **reliable**, **secure** and **easy**