

# LLM Ops on Azure

**Day 1**





# Who are we?





A **full-stack IT consultancy** with a global footprint

#### Xebia Facts



##### People

Jan 2024 5.000+ professionals



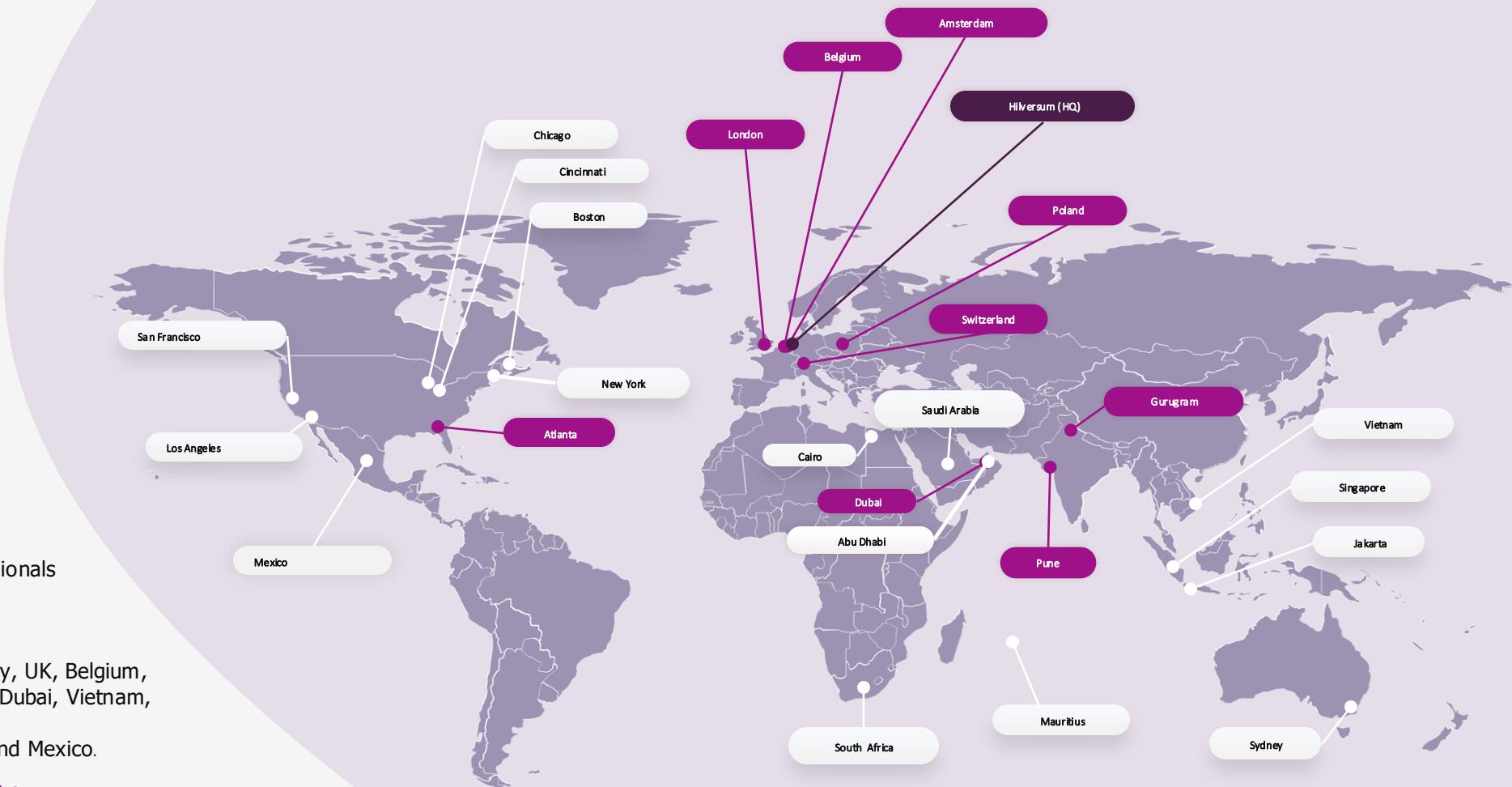
##### Offices

the Netherlands, Germany, UK, Belgium, Denmark, India, Poland, Dubai, Vietnam, US, Canada, Australia, Singapore, Switzerland and Mexico.



##### Full-stack digital consulting

Organized into service lines, each with a tremendous amount of knowledge and experience within a particular field, such as Agile, DevOps, **Data and AI**, Digital Strategy, Cloud, Software Technology, Security, Low-Code, and Microsoft Solutions.



#### Our Global Footprint

- Centre's
- Customer Presence

# About me

## **Prashant Srivastava**

• 3 years at Xebia

- Machine Learning Engineer
- Cloud Architectures
- GenAI Applications
- Trainer

# Short introductions



1. A bit about yourself...  
(What's your current position, your experience/background,...)
2. What is your experience with ML/ LLMs in Production?
3. (Maybe anything specific you'd like to learn?)

# Some comments about this training...

## **1. We will cover a simple LLM Use Case**

- won't go too much detail into the LLM itself
- It's more about LLMOps

## **2. It's going to be quite hands-on**

- we assume you have little prior MLOps knowledge

## **3. This training is for you!**

- Let me know what you want to learn
- ask questions, give feedback, etc.

# Schedule for today

## Day 1

- Intro MLOps & LLMOps
- Use case & solution design
- Building modularized LLM app
- Lunch 🍝🌮☕
- Evaluation & experiment tracking
- Containerization and deployment
- Automation with CI/CD

## Day 2

- Intro LLM system monitoring
- Tracking online metrics
- Logging traces
- Lunch 🍝🌮☕
- User feedback
- Iterating on the solution
- Wrap up

**Repo:** <https://github.com/godatadriven/academy-llmops-in-azure>

# Why LLMOps?



# Why ~~LLM~~Ops MLOps?



Daisy Data Scientist @ TurbineDynamics

Task: Help business improve generated power forecast.



Ideate with  
business to define  
use case & value



Explore data to  
find possible  
relationships

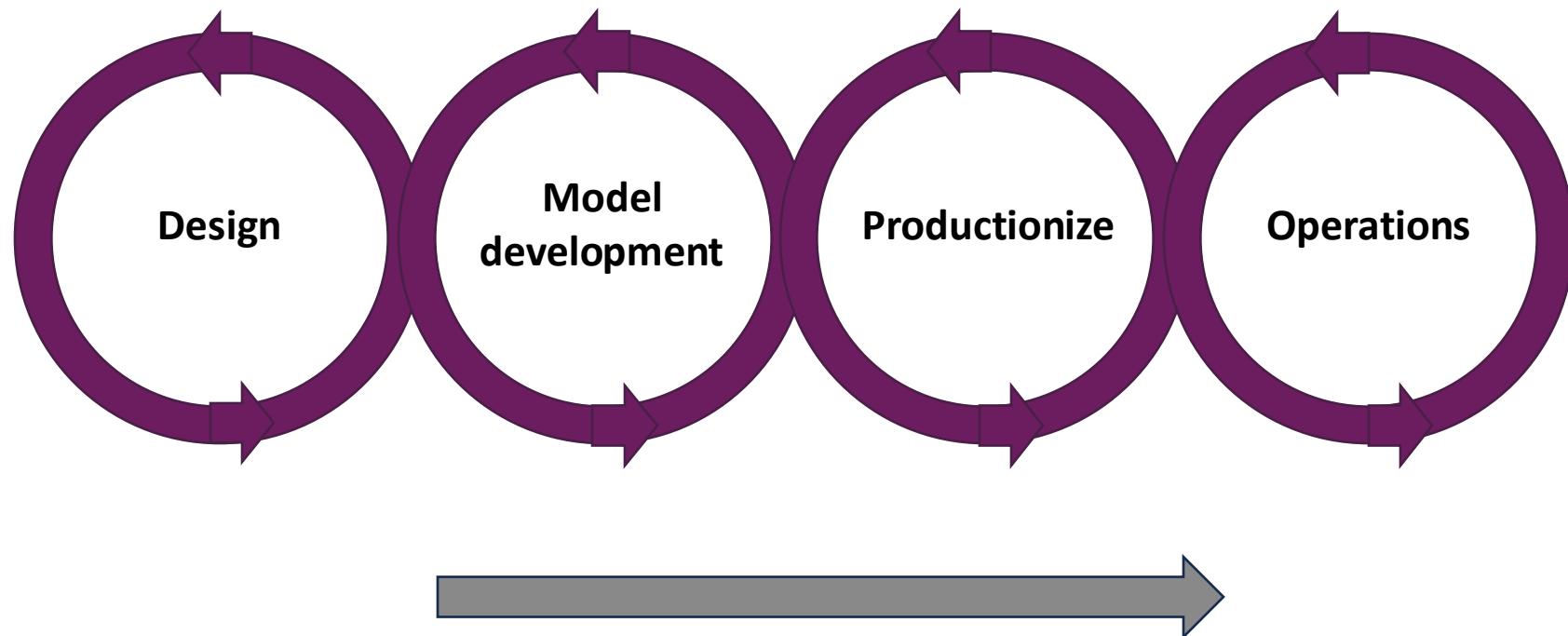


Create predictive  
model in a  
notebook

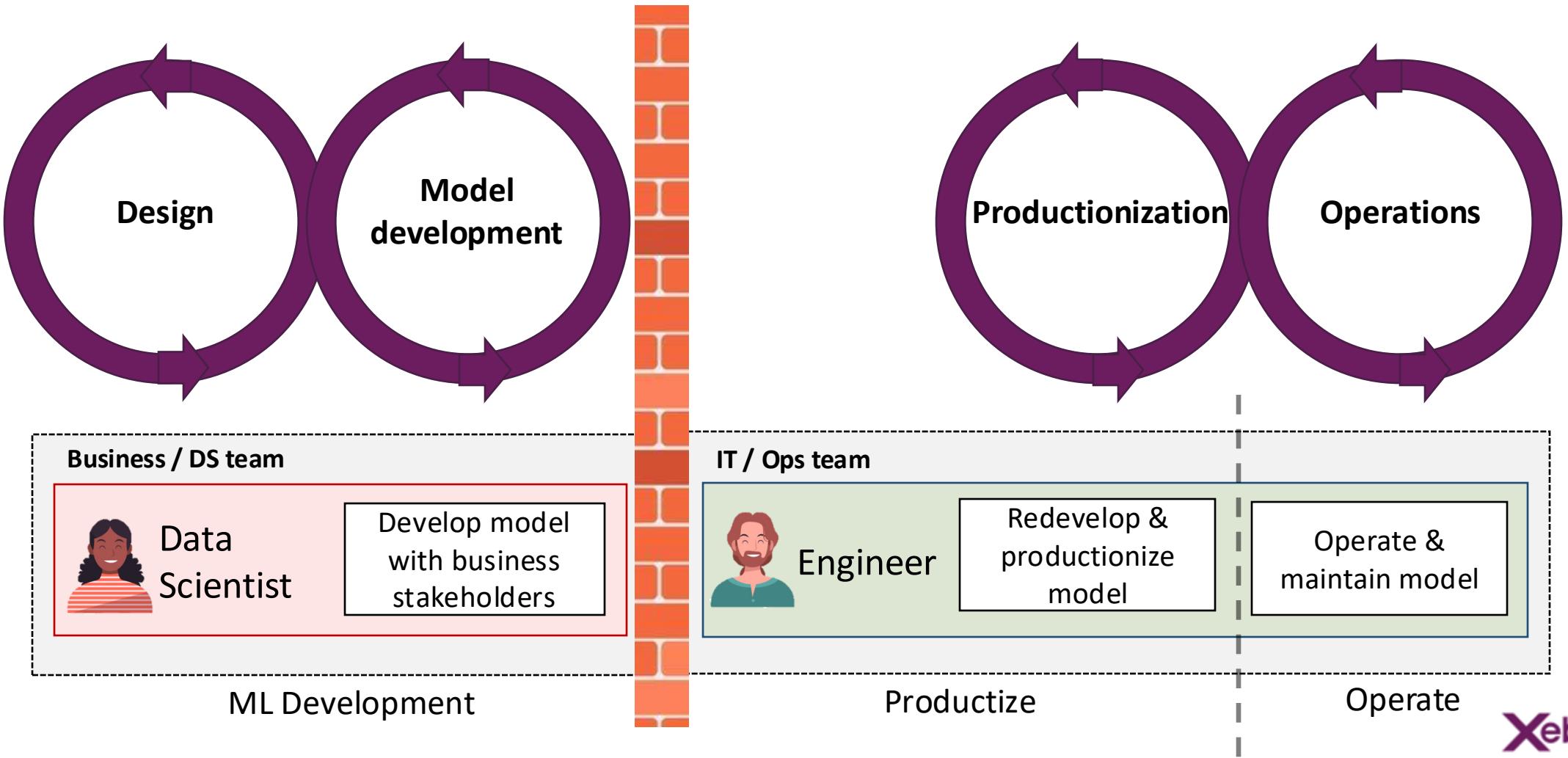
Now what?



# Ideally: a project follows the “MLOps lifecycle”

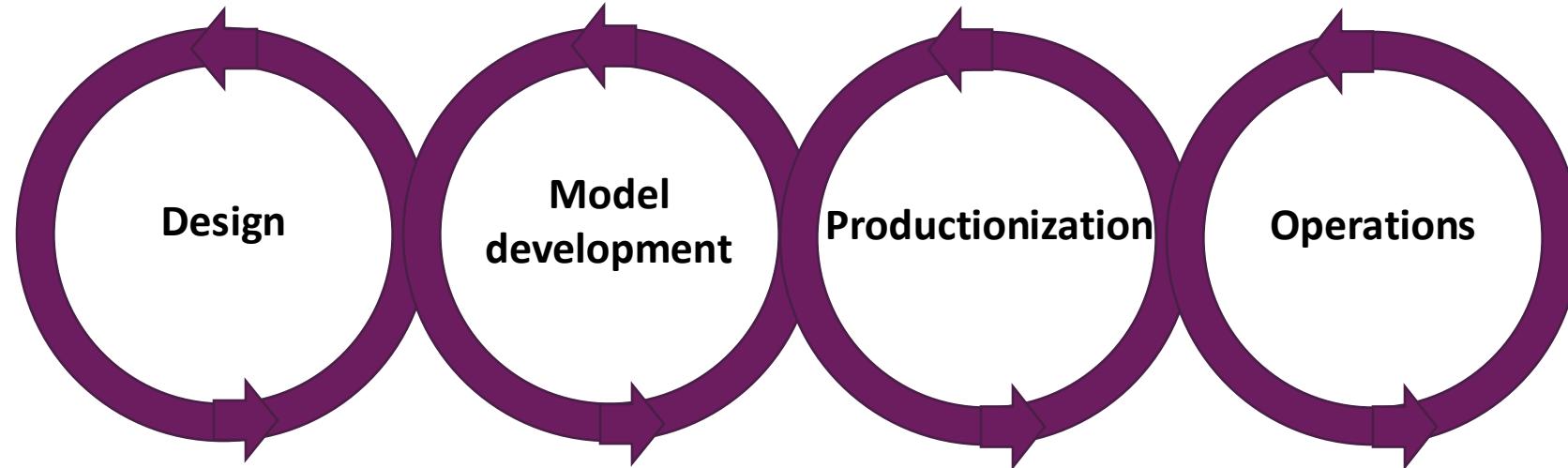


# However, we often see a *handover*

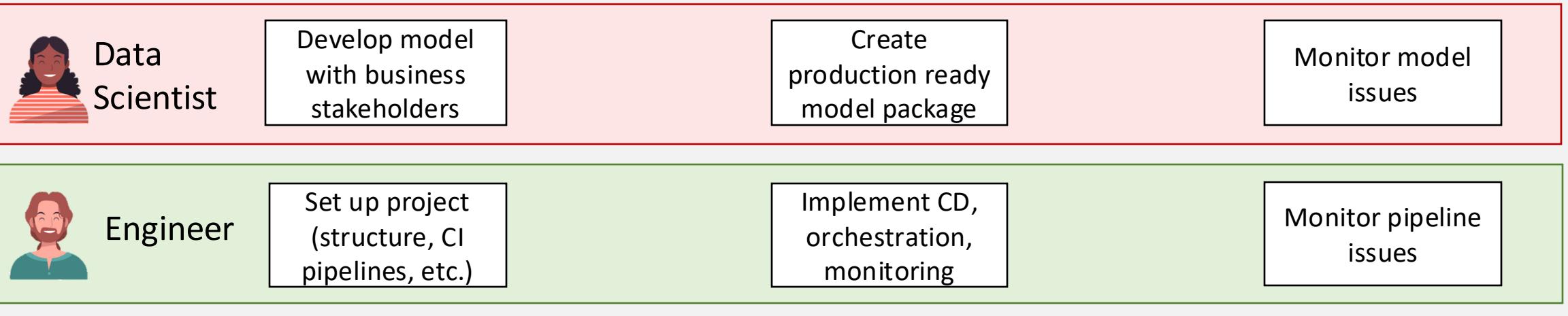


# Goal of MLOps

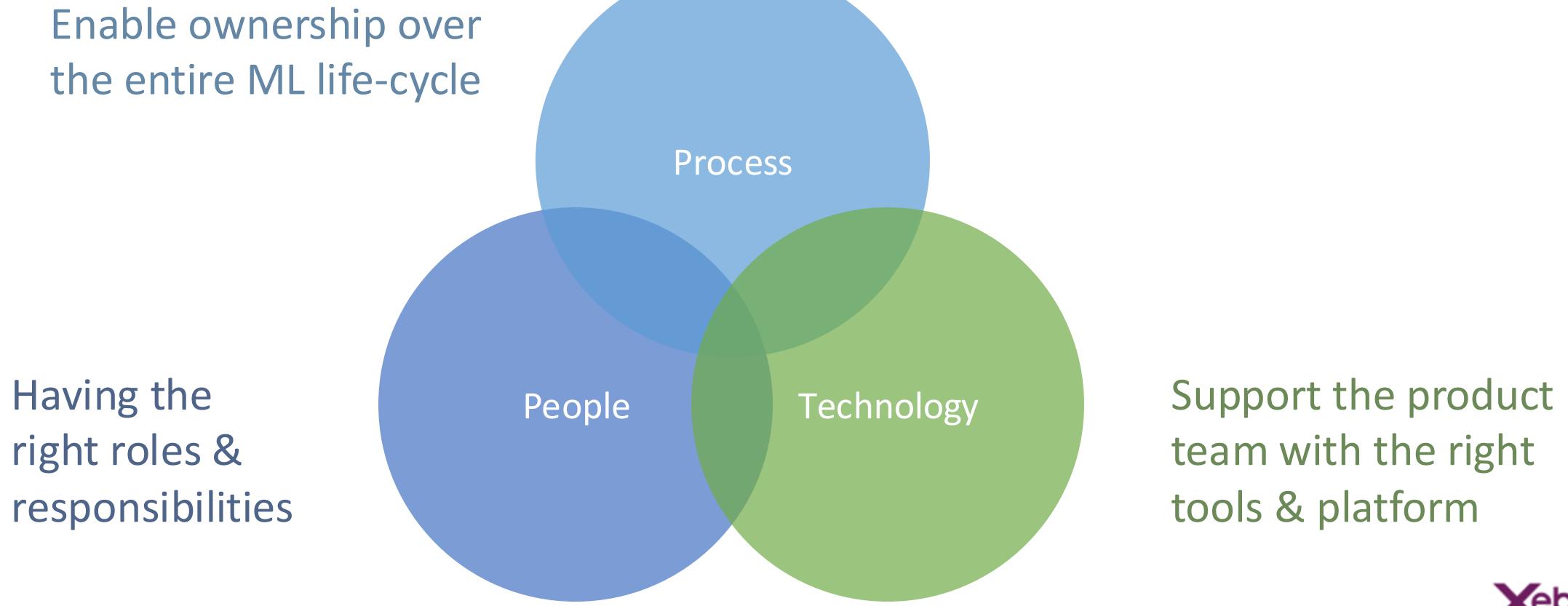
## Close the gap between DS and Ops



End-to-end DS product team



# MLOps: Close the gap by combining the right people, processes and technology



# Why LLMOps?



Daisy Data Scientist @ LanguageDynamics

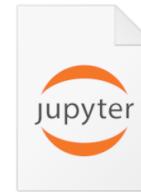
Task: help business **extract structured information from unstructured text.**



Ideate with business to define use case & value



Explore data and **tweak prompt** to get intended result



Create working solution in a notebook

Now what? →

LLMops aims to solve a similar problem as MLOps in a similar way!  
That is, by combining the right **People, Processes, and Technology**

# Differences MLOps and LLMOps

	MLOps	LLMOps
People	Data scientist (statistics, ML) Machine learning engineer Data engineer (structured data) Product owner ...	Data scientist (NLP) Machine learning engineer Data engineer ( <b>unstructured data</b> ) Product owner <b>Prompt engineer (potentially)</b> ...
Processes	Project templates (typical ML use cases) Experimentation (logging hyperparams and metrics) Pipelines (preprocessing and ML training) Deployment (usually API endpoints) Monitoring (typical ML metrics and infra) ...	Project templates (typical LLM use cases, RAG, etc.) Experimentation ( <b>logging prompts and metrics</b> ) Pipelines (preprocessing and LLM batch jobs) Deployment ( <b>usually web apps</b> ) Monitoring (LLM metrics, <b>traces</b> and infra) ...
Technology	Scikit-Learn Tensorflow/Keras/PyTorch/Jax MLflow FastAPI ...	LangChain, LangGraph, PydanticAI LLM-specific APIs* Streamlit, Flask, FastHTML, ... LangSmith ...

\* Assuming we are using a provider's API

# Differences MLOps and LLMOps\*

	MLOps	LLMops
Compute	CPU / GPU	CPU (making use of external model API's)
Experimenting	Established metrics Logging model params and weights	Harder metrics definition Logging prompts
Pipelines	Training and evaluating Batch predict	Evaluating Batch predict
Deployment	As model endpoint	As web app or endpoint
Monitoring	Logs, metrics	Logs, metrics, traces

# Use case

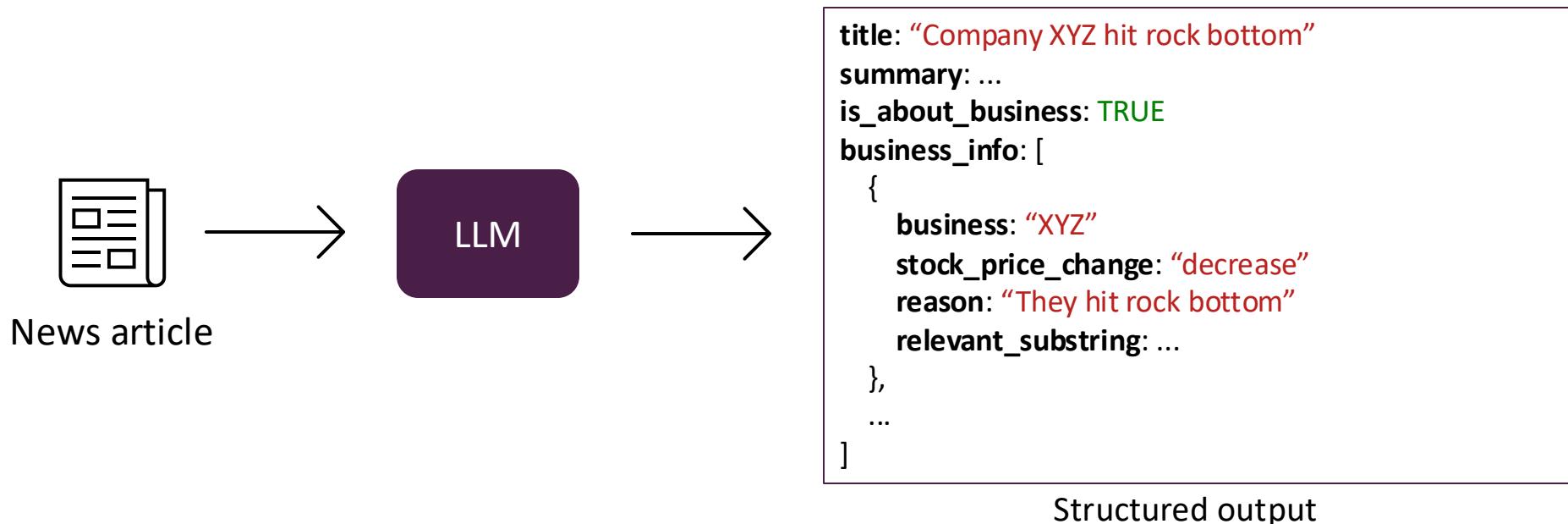




# News Reader:

## Extract structured info and predict stock price changes from news articles

- We work at a bank and are interested in analyzing news articles for likely stock price changes
- We have created a notebook with an LLM-solution to extract relevant info from the articles
- In the end, we want a web app where users can upload their articles which are then automatically analyzed



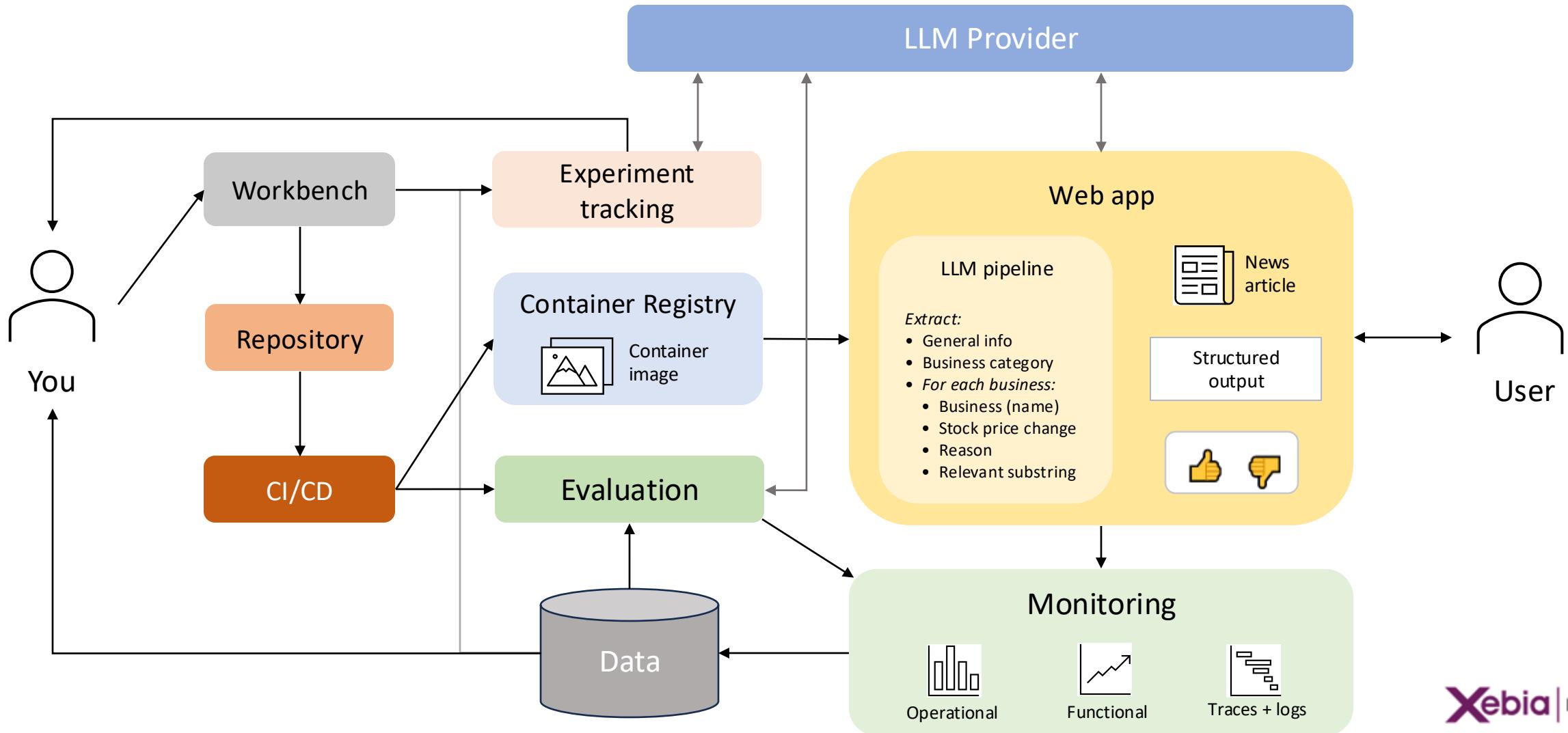
# Exercise:

## Develop a high-level solution design (architecture) for the News Reader use case

### Requirements:

- Should provide a **web app** that accepts text and returns stock price change predictions
- The model is **a provider's LLM** (no self-training)
- The user should be **able to provide feedback** on the LLM outputs
- Ideally, you would want to include **monitoring & experimentation (evaluation)**
  - **For evaluation, you may want to store some annotated articles**
- For auditing, you would want to log not only the predictions, but also general info about the business, the article and a reason for the prediction
- Code should be maintained in a **code repository** (e.g. GitHub)
- The webapp is containerized, and we need to store our webapp versions and implement automatic testing/ quality checks etc. (**CI/CD**)

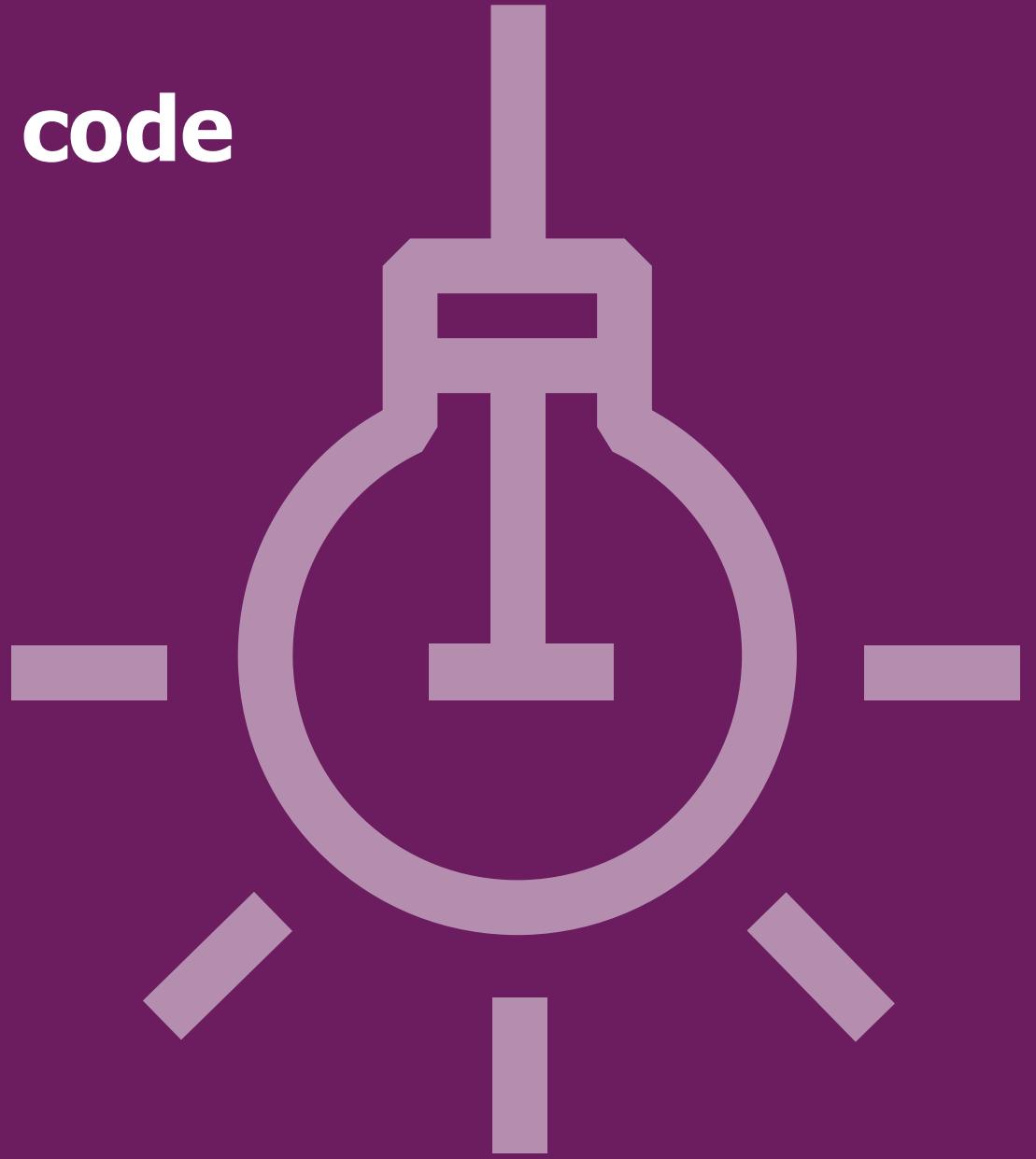
# High-level solution diagram



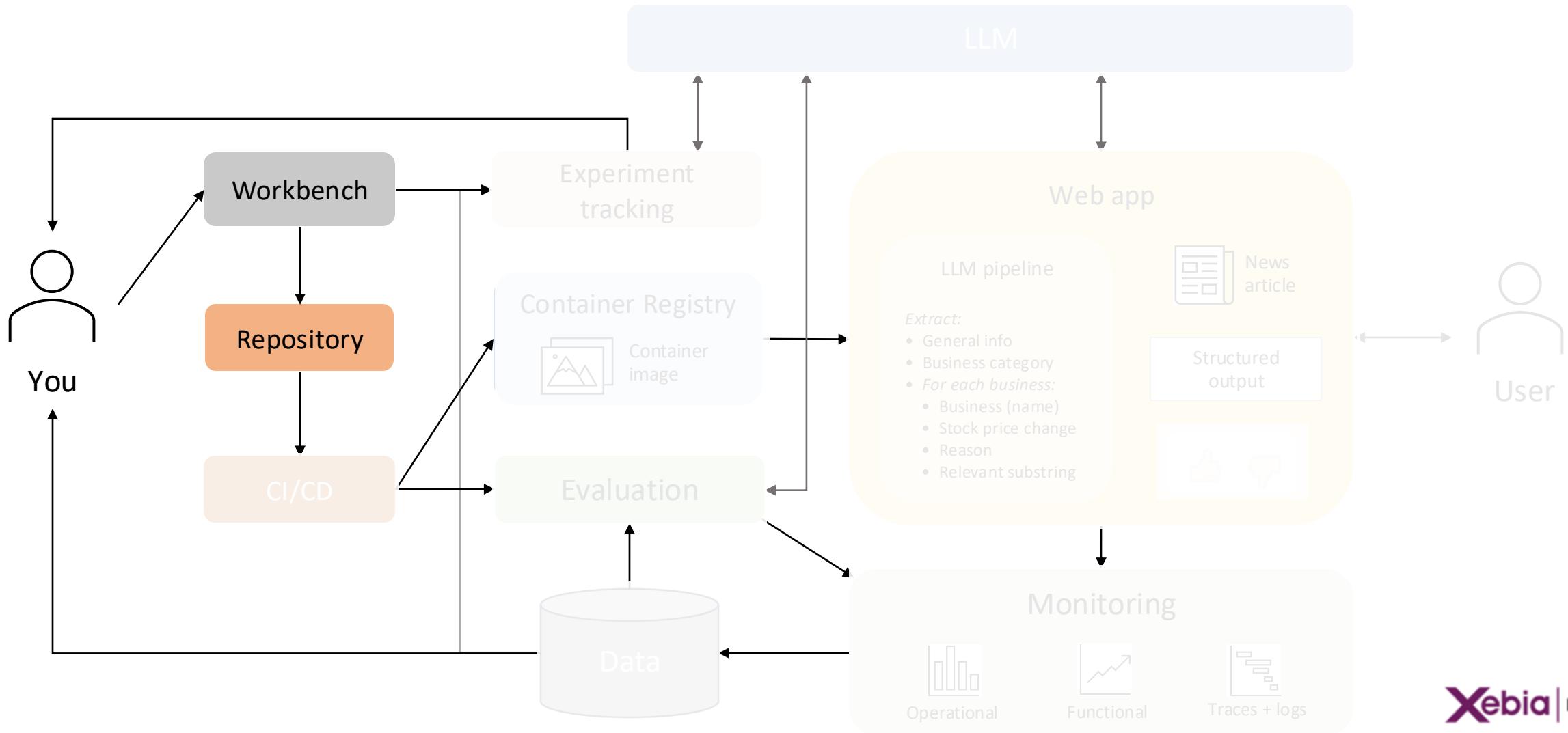
# Break time!



# Getting started with the code



# High-level solution diagram

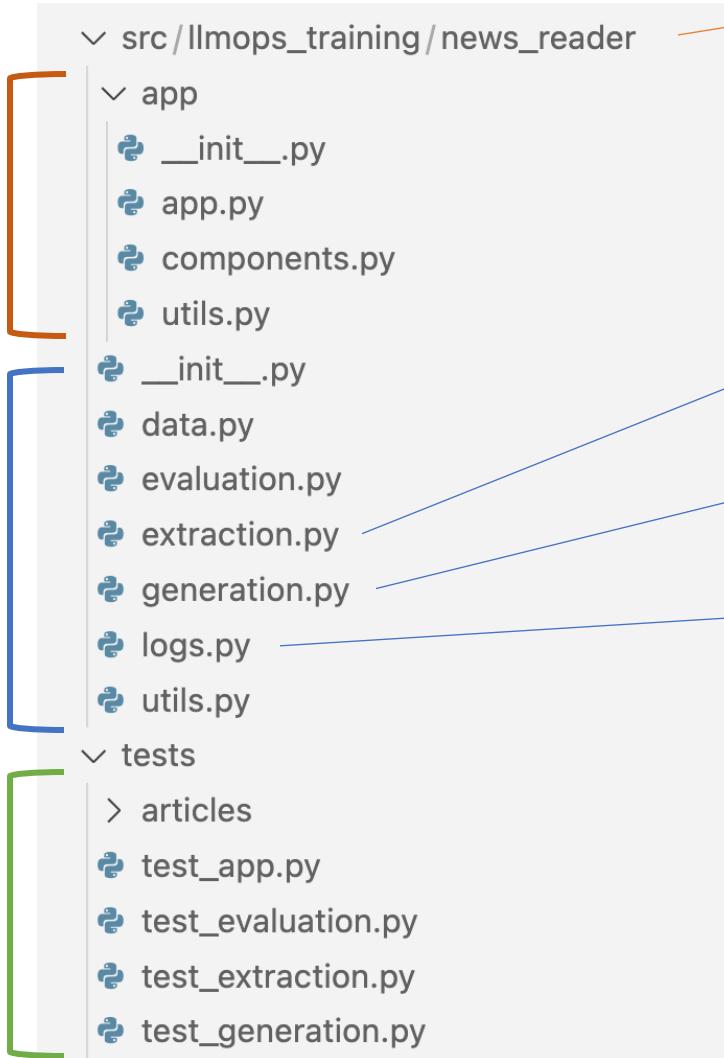


# Our code structure

Streamlit  
(web)  
app

Main  
code

Tests



Our package that we'll extend during the training

Main logic of our LLM solution

LLM API connections

Configuring logger + tracer

Where you'll spend most time:

> exercises

# How our web app looks like:

### Articles

Article Stats  
Number of articles: 1

Select Article  
Only one article available.

Extract from Articles  
Upload your articles  
Drag and drop files here  
Limit 200MB per file • TXT, MD  
Browse files

Type or paste your article  
  
Extract info

### News Reader

#### Selected Article

Senegal election: Concerns grow as President Macky Sall fails to set new election date - BBC News

President Macky Sall, elected for a second time in 2019, has said he will not run again. There is international concern about the political situation in Senegal after President Macky Sall called off this month's election citing a row over the eligibility of candidates. The West African regional bloc Ecowas and the US urged the authorities to clarify when the vote will take place, as no new date was set. France and the EU have also called for an election as soon as possible. Police have fired tear gas at pockets of demonstrators in the capital. The scenes in Dakar are a worrying sign of what may be ahead and more protests have been called for Monday. Senegal has been seen as one of the most stable democracies in West Africa, but uncertainty over when people will get a chance to vote as well as suspicion that the electoral race is not free and fair is undermining this image. President Sall's reiteration that he will not run for a third term had been welcomed, but there is worry about the consequences of a lack of an electoral timetable. The vote was supposed to have taken place in three weeks' time. Uncertainty has played a role in deadly protests in the past as has the opposition allegation that politically motivated cases were being brought against potential candidates. For months, speculation that Mr Sall was seeking a third term fuelled opposition protests that led to violence and many deaths. He finally declared he was not seeking a third term last July. An opposition candidate called the president's decision to postpone the

#### Structured Output

*An error occurred while processing this article.*

We will implement this part

# Exercise:

## Getting familiar with the use case

**First, let's get you set up, so you can start coding!**

- Open the repository in GitHub: <https://github.com/godatadriven/academy-lilmops-in-azure>
- Follow the instructions in the `readme.md` to set up your coding environment

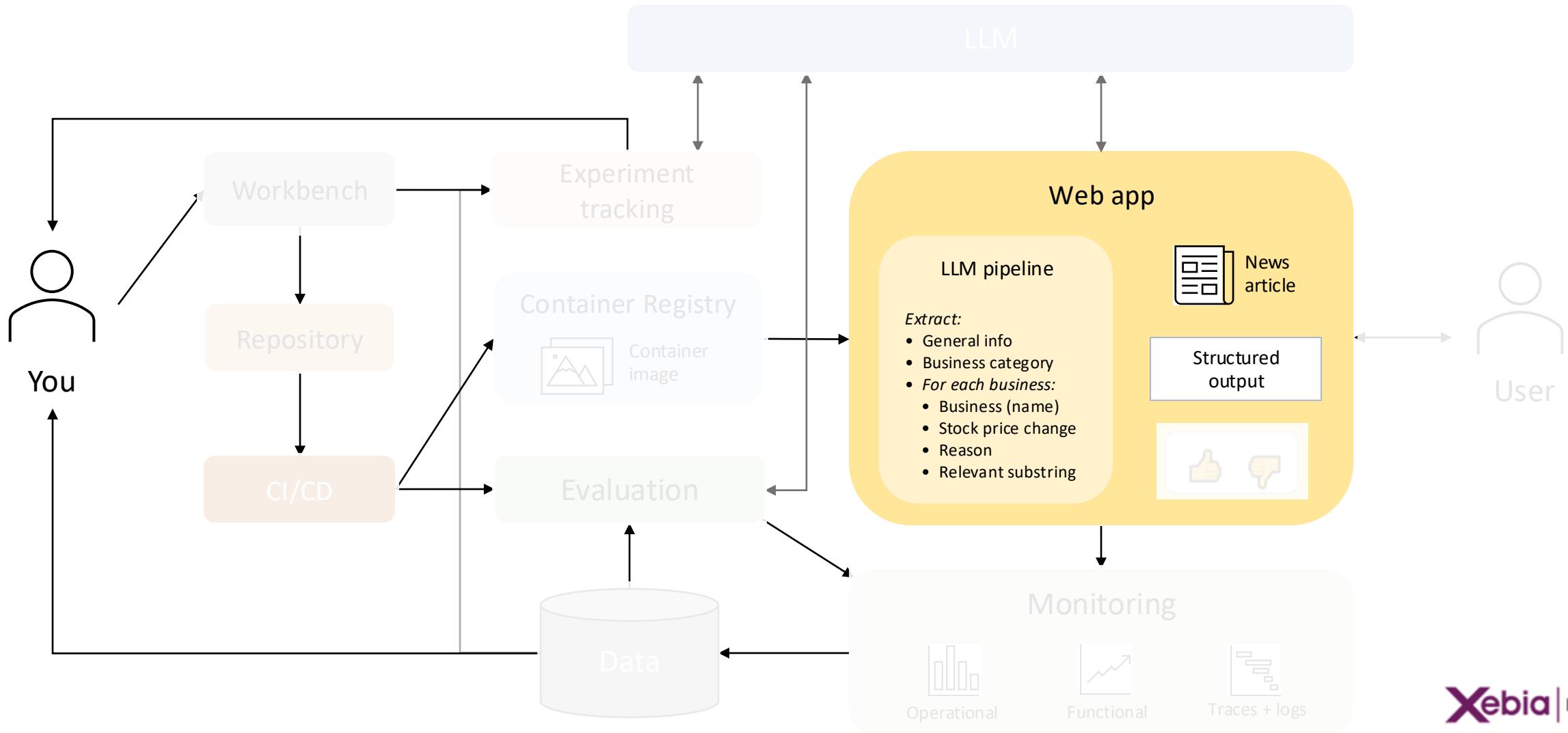
**Done setting up?**

- Run through **exercises 1 - 3** in the `/exercises` directory

# Modularizing the LLM pipeline



# High-level solution diagram



# “One prompt to rule them all” not always the best approach

```
prompt = f"""
```

You are a helpful assistant.

You extract structured info from news articles.

If a given article is **business-related**,  
**find all businesses** mentioned in the article  
and **for each business** in the article,  
**predict** whether stock price will increase

Here's a 100 articles:

```
{articles}  
"""
```

???



LLM

# Splitting complex tasks up into smaller tasks can have several benefits

- We can apply different models per task
- It allows the model to focus on smaller tasks at once, possibly improving performance
- It allows for better testing, debugging and maintaining
- In some cases, we can parallelize different tasks



However, there is a **balance** to be found!

- A single prompt can be **complex but quick** to execute
- Multiple smaller prompts can have above benefits, but may come **at the cost of speed**

# How this could look like in our case

```
prompt1 = f“Extract general info from {article}”
```

```
prompt2 = f“Classify whether {article} is about business”
```

```
prompt3 = f“Extract involved businesses from {article}”
```

```
for business in businesses:
```

```
    prompt4 = f“Predict stock price change for {business} in {article}”
```

# Exercise: Modularizing the solution

**Let's split up our LLM pipeline in smaller modules and integrate them in the app**

- Run through **exercise 4** in the /exercises folder

# Lunch break!



# Schedule

## Day 1

- Intro MLOps and LLMOps
- Use case & solution design
- Building modularized LLM app
- Lunch 
- **Evaluation and experiment tracking**
- **Containerization and deployment**
- **Automation with CI/CD**

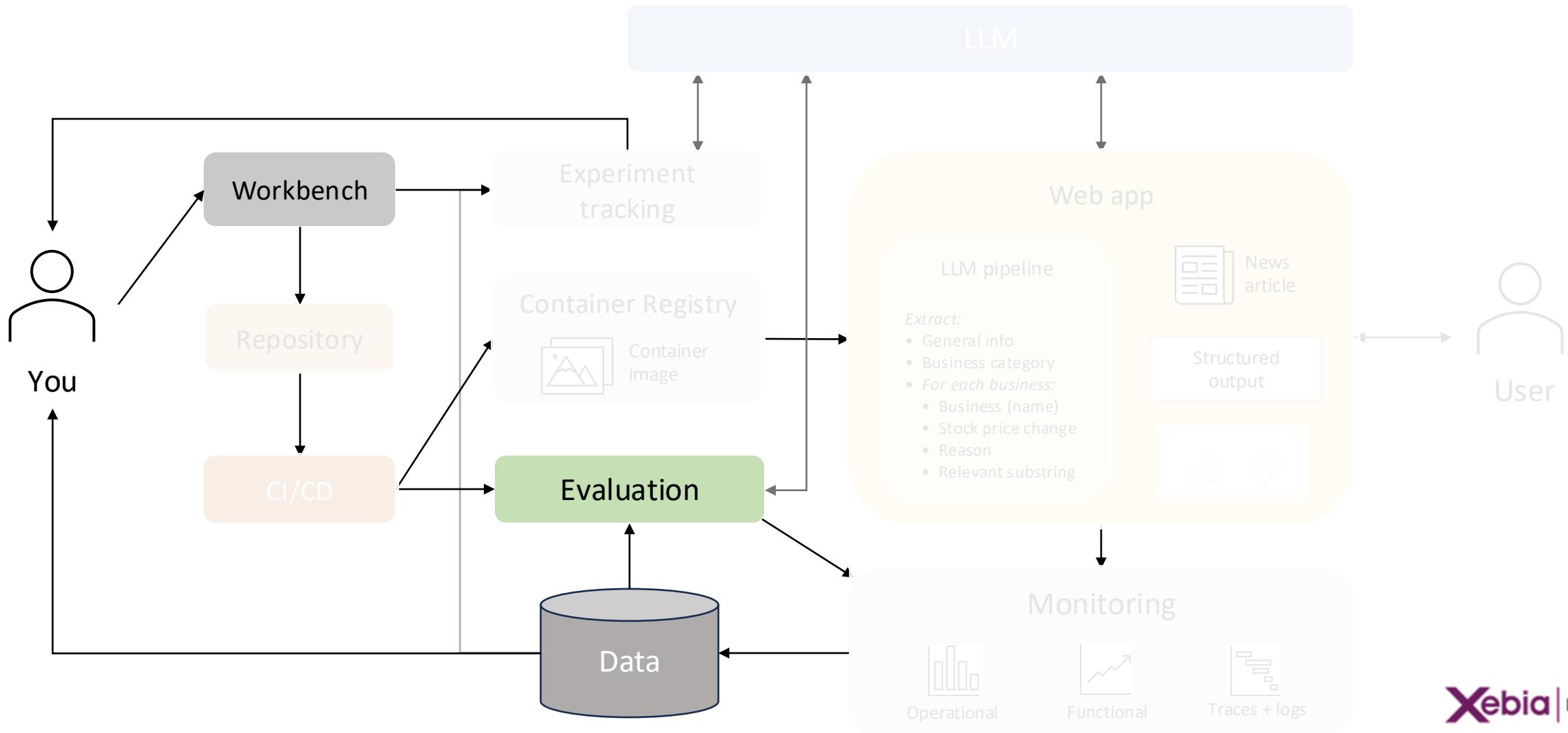
## Day 2

- Intro LLM system monitoring
- Tracking online metrics
- Logging traces
- Lunch 
- User feedback
- Iterating on the solution
- Wrap up

# Evaluating the LLM solution



# High-level solution diagram



# Different forms of evaluations

## Offline:

- Eye-balling
- Smoke tests
- Using labeled data and predefined metrics
- Letting an LLM evaluate our outputs

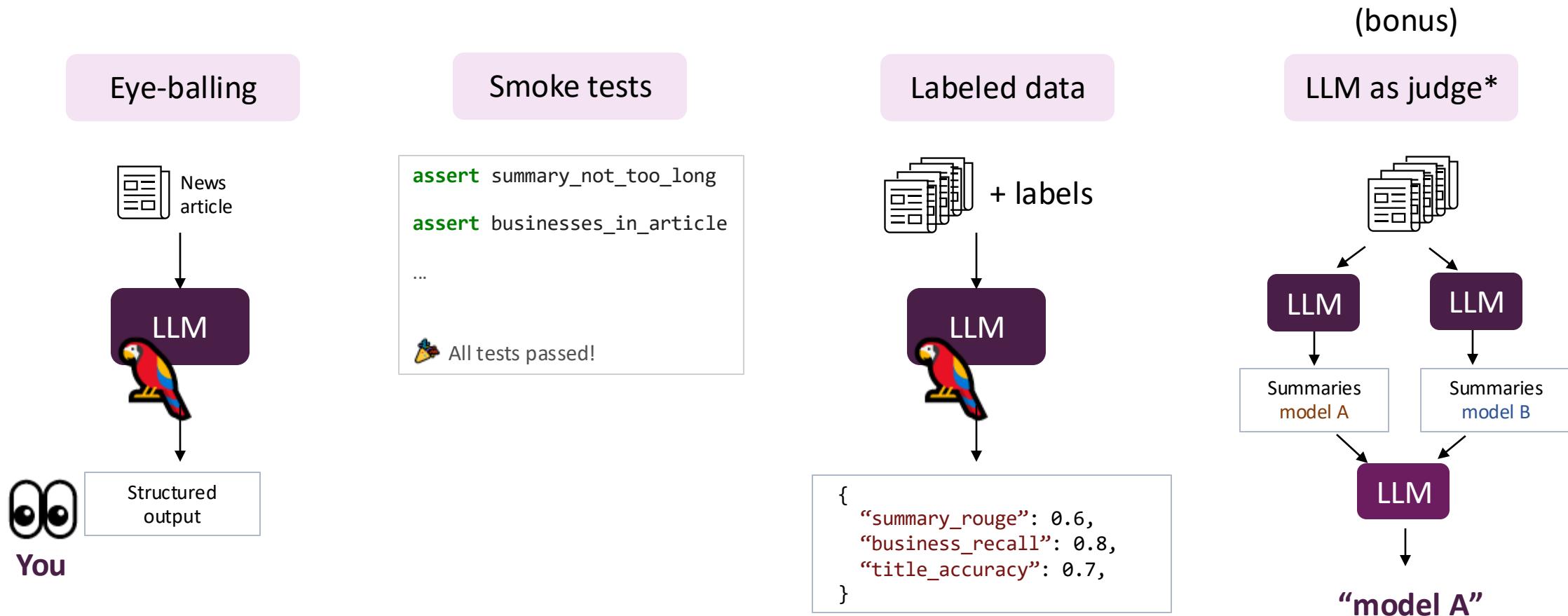
*Before we go live*

## Online:

- Metrics collected from production
- Explicit user feedback
- A/B testing

*When we are live*

# Offline evaluation for our use case



\* Only use with caution! Good alignment with human judgement essential. See [this blogpost](#)

# “You should track the correlation between model-based and human evaluation to decide how much you can rely on automatic evaluation”

“I sent my colleague Phillip the following spreadsheet every few days to grade for a different use-case”

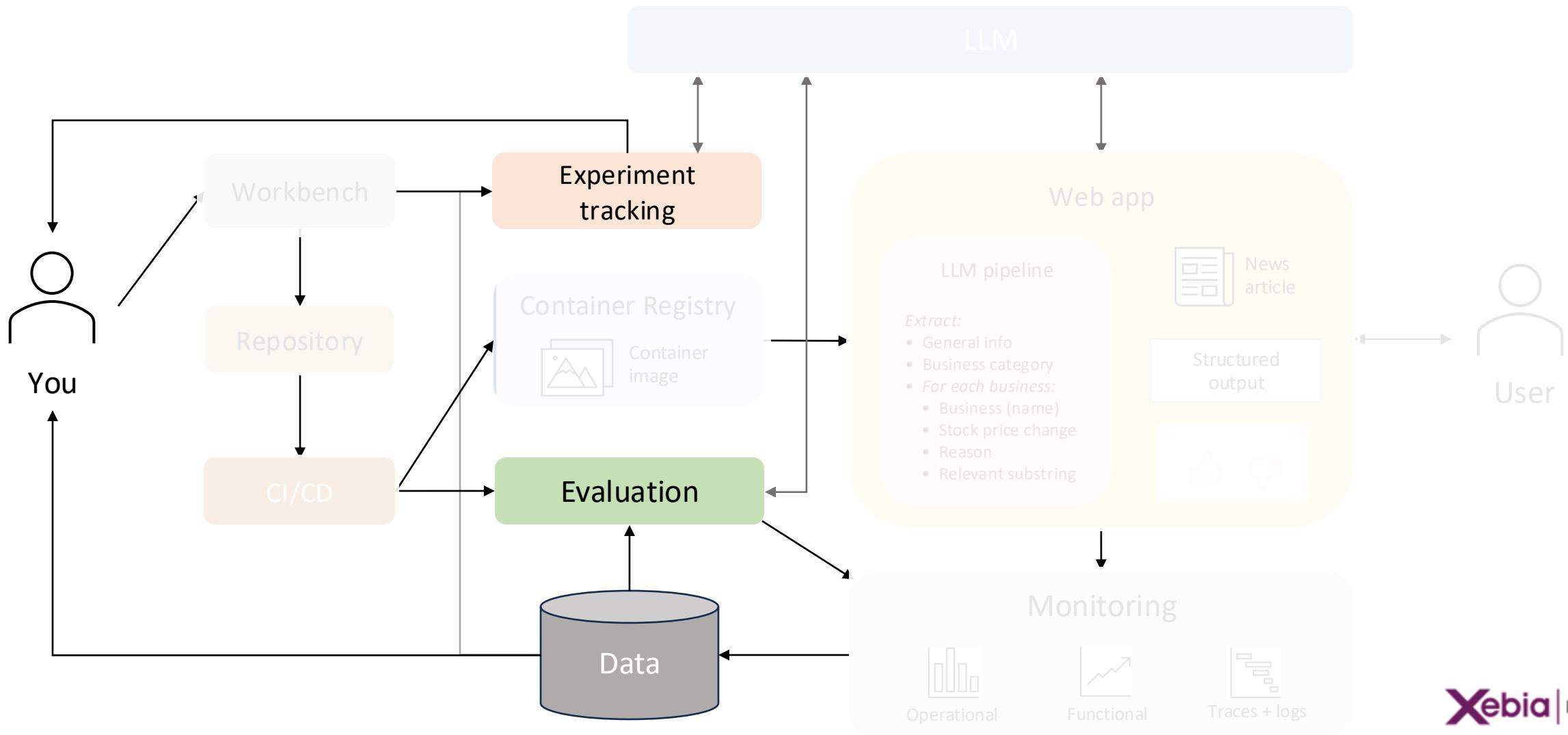
C	D	E	F	G	H	I	J
model response	model critique	model outcome	phillip critique	phillip outcome	phillip revised response	agreement	
{"calculations":[{"column":"duration_ms","op":"MAX"}],"filters":[{"column":"trace.parent_id","op":"does-not-exist","join_column":""}],"orders":[{"column":"duration_ms","op":"MAX","order":"descending"}],"limit":1,"time_range":7200}	The response is nearly correct, as it is looking for the slowest trace by using MAX(duration_ms) and ordering by duration_ms in descending order, which is appropriate for finding the 'slowest' trace. Additionally, filtering with trace.parent_id does-not-exist correctly identifies root spans. However, the query should be grouping by trace.trace_id to ensure that we identify distinct traces, not just the longest individual span. Without the correct grouping, the analysis does not guarantee that the result is a full trace, but merely the longest span. Also, specifying a limit of 1 is good as it will return the single slowest trace, as requested.	bad	The response is nearly correct, as it is looking for the slowest trace by using MAX(duration_ms) and ordering by duration_ms in descending order, which is appropriate for finding the 'slowest' trace. Additionally, filtering with trace.parent_id does-not-exist correctly identifies root spans. However, the query should be grouping by trace.trace_id to actually show the slowest trace. Without grouping, the query only shows the MAX(duration_ms) measurement over time, irrespective of which trace is responsible for that measurement.	A 1 2 3 4	B Iteration 1 2 3 3	Agreement % 68.00% 82.61% 94.00%	TRUE

\* Only use LLM as a judge with caution! Good alignment with human judgement essential. See [this blogpost](#)

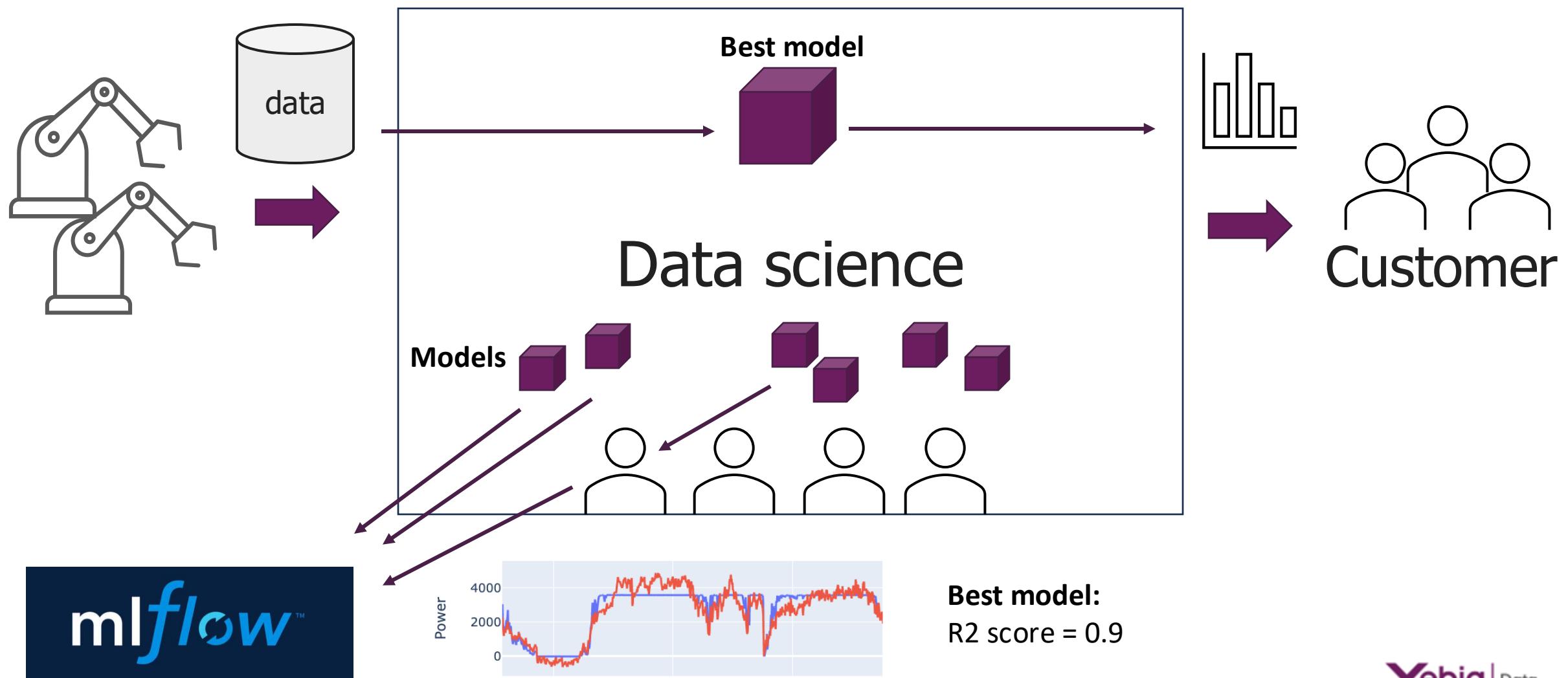
# Experiment tracking



# High-level solution diagram



# Recap ML experiment tracking



# LLM experiment tracking

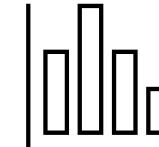
LLM



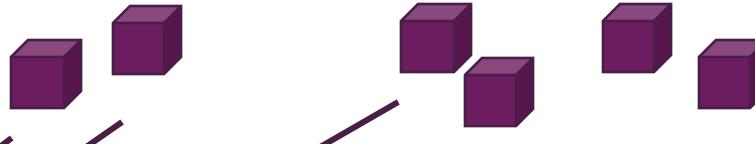
- News article
- News article
- News article
- News article



**Best LLM prompt + params**



**Prompts, model types, params**



(or LLM-specific tools like LangFuse,  
LangSmith,...)

Parameter: general_info_prompt_template	Metric: general_info_success_rate	Metric: summarization_rouge_
Extract structured information from the following article: {article}	0.9	0.2529504825

Xebia | Data

# Exercise: Evaluating the solution and tracking experiments

**Let's add some assertions and allow for evaluations based on a labeled dataset**

- Run through **exercise 5** in the /exercises folder

**Now, we can try to improve those metrics, but it's important to keep track of our runs!**

- Run through **exercise 6** in the /exercises folder

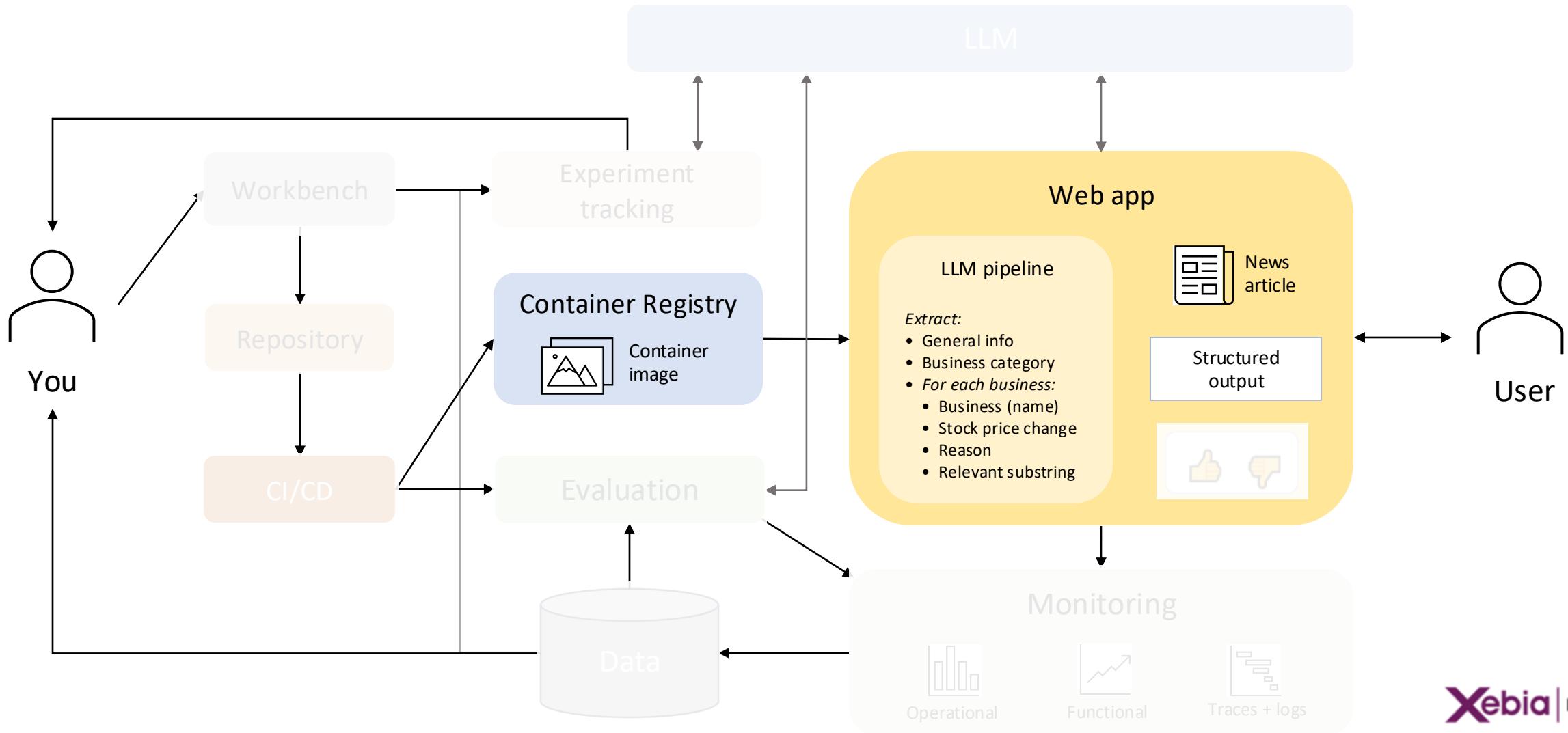
# Break time!



# Deployment



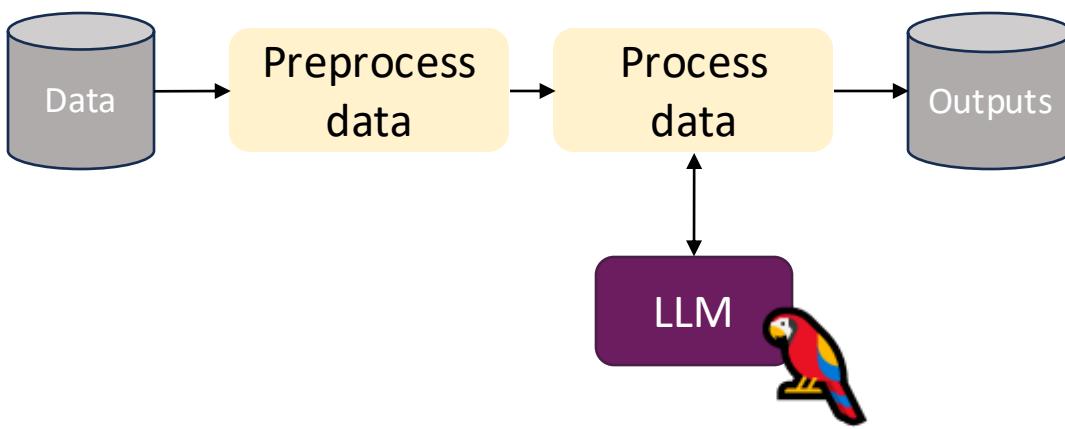
# High-level solution diagram



# Serving patterns

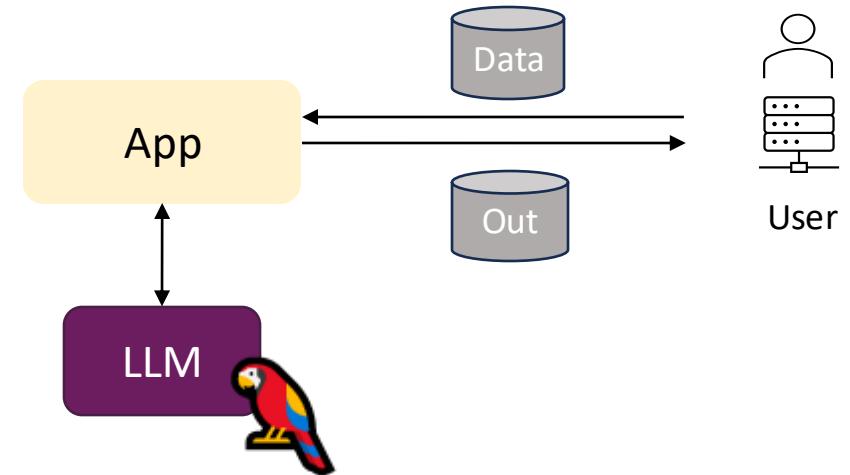
## Batch:

- When predictions/generations can wait

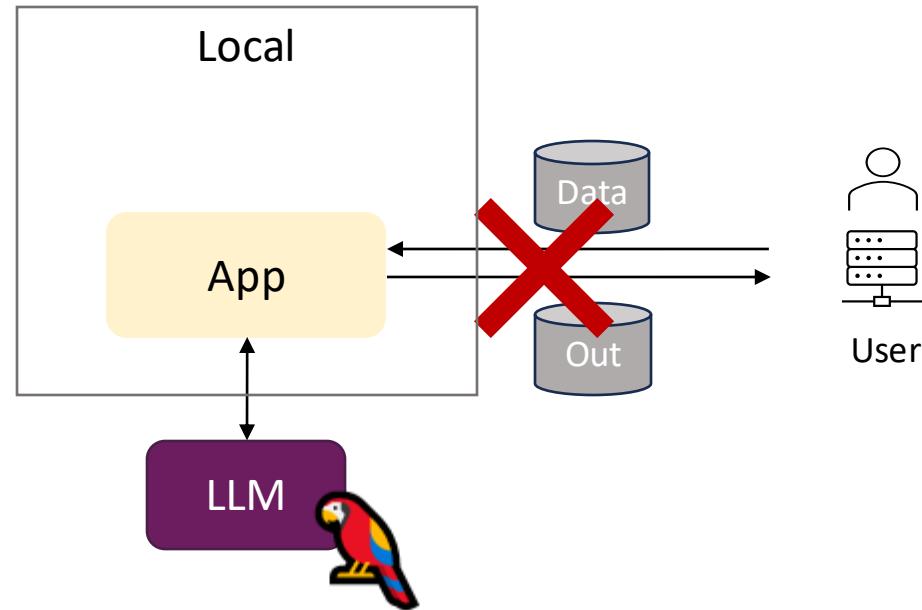


## On-demand:

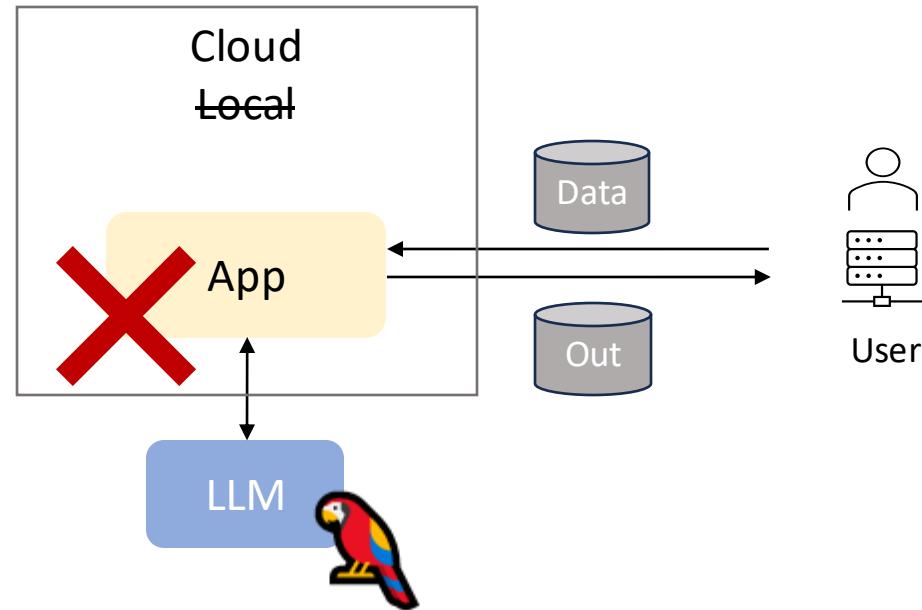
- When predictions/generation need to be served immediately



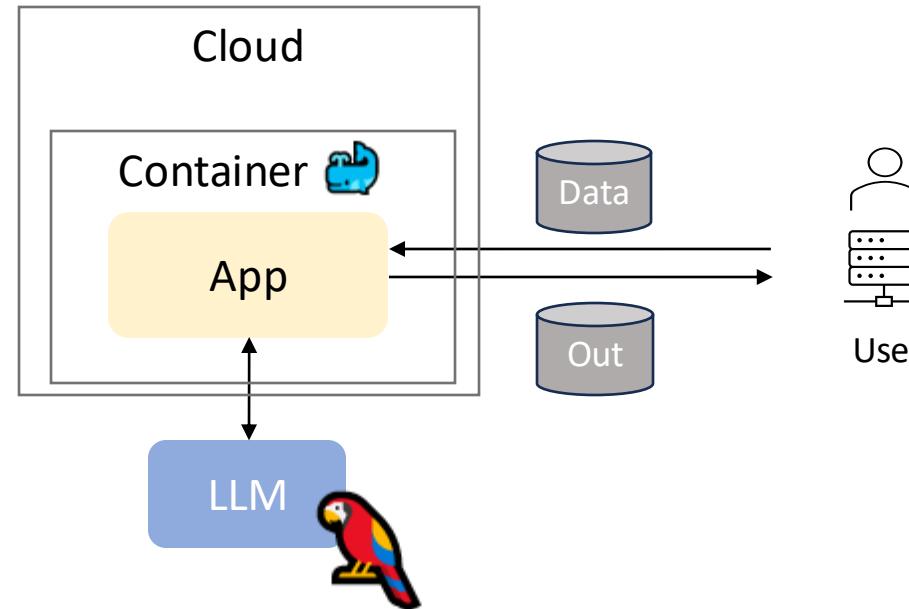
# Users cannot access our app yet!



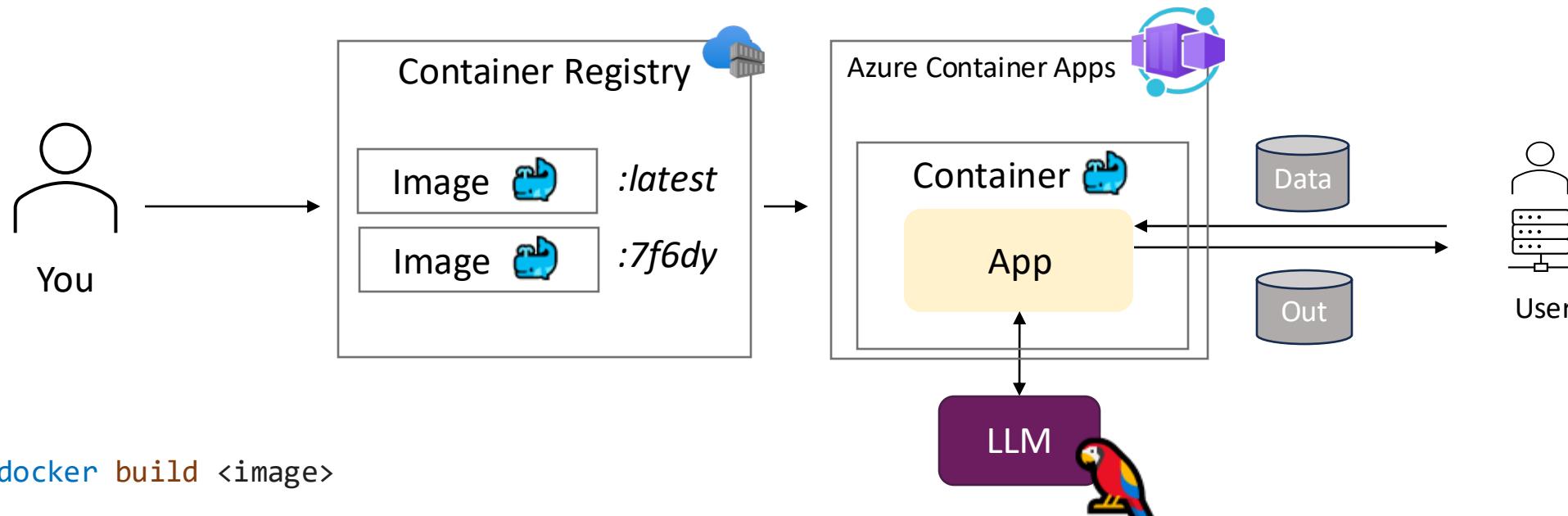
# Let's put it in the cloud!



# Containerization needed before we can run our code in the cloud



# We'll use Azure Container Registry and Azure Container Apps



# Exercise: Deploying our solution

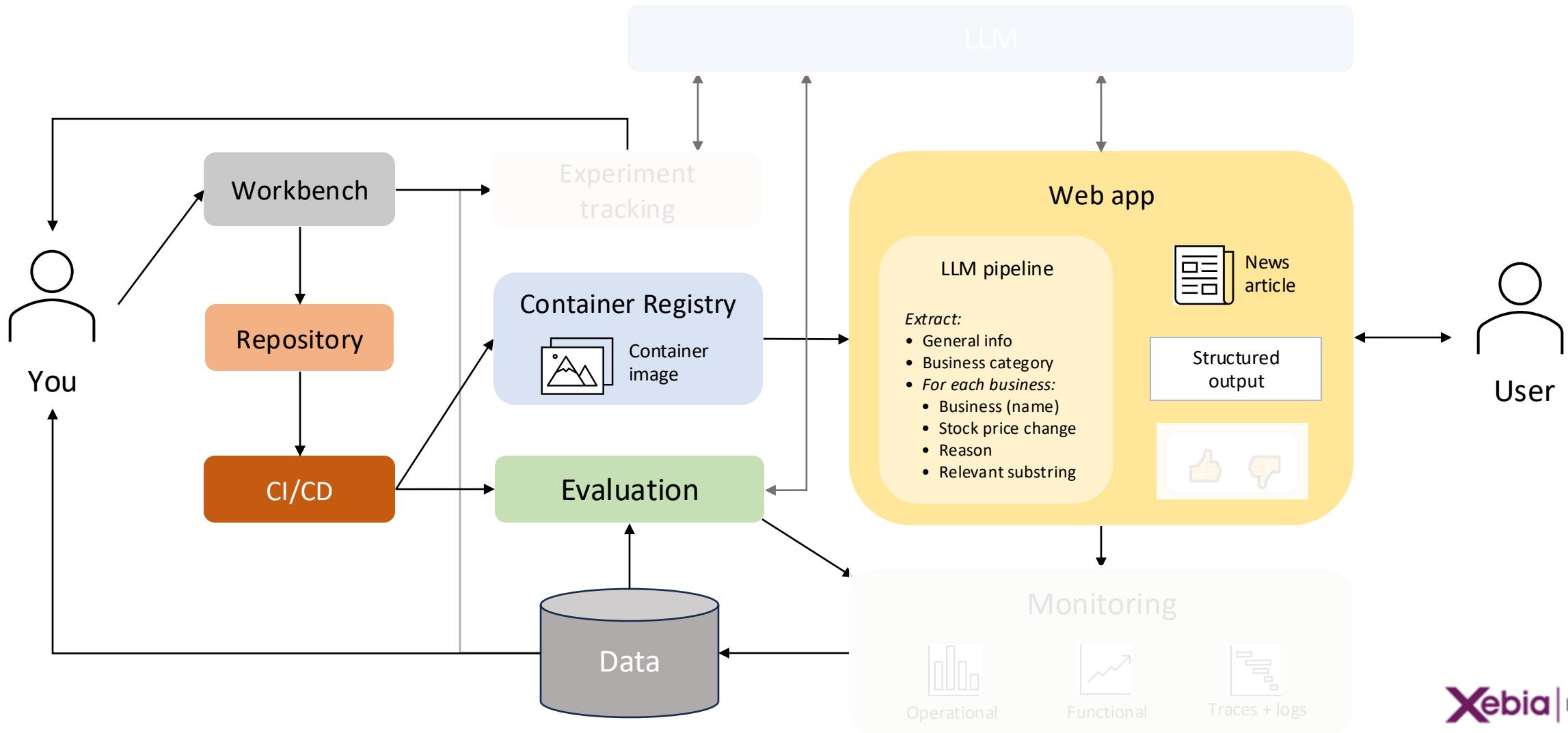
**Let's containerize and deploy our app!** 

- Run through **exercises 7 + 8** in the `/exercises` folder

# Automation with CI/CD



# High-level solution diagram



# What is CI/CD?

- CI/CD is a methodology of **frequently delivering** software by **automating** many stages of development
- Consists of two practices:
  - **Continuous Integration:** continuously building, testing and merging code into a single code base
  - **Continuous Deployment:** continuously deploying new releases/updates of your application to end users

# Examples of CI steps

- Checking code quality (e.g. using pylint)
- Checking code formatting (e.g. using black)
- Running unit tests (e.g. using pytest)
- Building/testing a Docker image
- ...

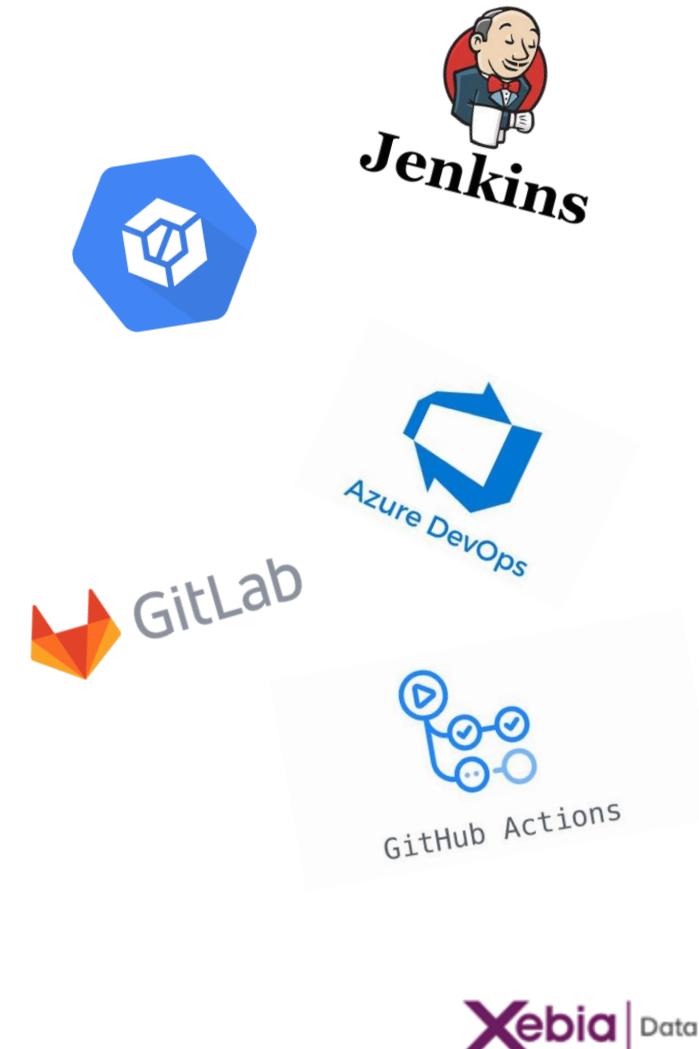
# Examples of CD steps

- Building and pushing a Python package
- Building and pushing a Docker image
- Deploying an API to a (dev/test/prod) environment
- ...

# How to implement CI/CD – YAML!

```
1  name: CI
2
3  on:
4    pull_request:
5      branches:
6        - main
7
8  jobs:
9    run-checks:
10   steps:
11     - name: ...
12     - name: Run pytest
13     | run: pytest
14
15  another-job:
16    steps:
17      ...
18
19
20
21
22
23
```

```
1  name: CD
2
3  on:
4    push:
5      branches:
6        - main
7
8  jobs:
9    build-push-docker:
10   runs-on: ubuntu-latest
11
12   steps:
13     - name: Login to Google Artifact Registry
14     | uses: ...
15
16     - name: Build Docker image
17     | run: docker build ...
18
19     - name: ...
20
21     - name: Deploy to Cloud Run
22     | run: gcloud run deploy ...
```

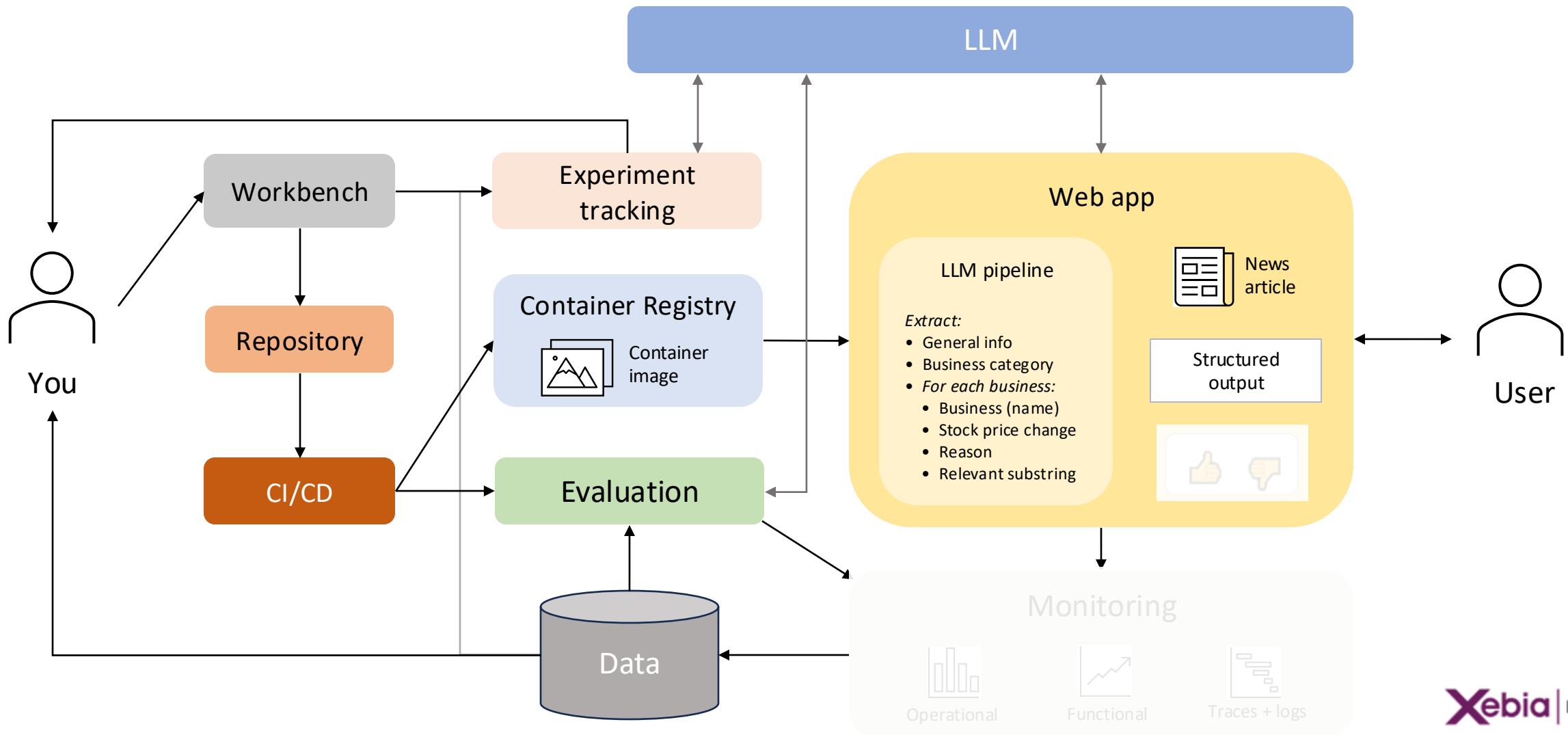


# Exercise: Automation with CI/CD

**We're running too much still manually, let's automate (almost) everything!**

- Run through **exercise 9** in the /exercises folder

# High-level solution diagram



# LLM Ops on Azure

Day 2



# Schedule

## Day 1

- Intro MLOps and LLMOps
- Use case & solution design
- Building modularized LLM app
- Lunch 🍝🌮☕
- Evaluation and experiment tracking
- Containerization and deployment
- Automation with CI/CD

## Day 2

- **Intro LLM system monitoring**
- **Tracking online metrics**
- **Logging traces**
- Lunch 🍝🌮☕
- User feedback
- Iterating on the solution
- Wrap up

**Repo:** <https://github.com/godatadriven/academy-llmops-in-azure>

# Recap day 1



# Differences MLOps and LLMOps

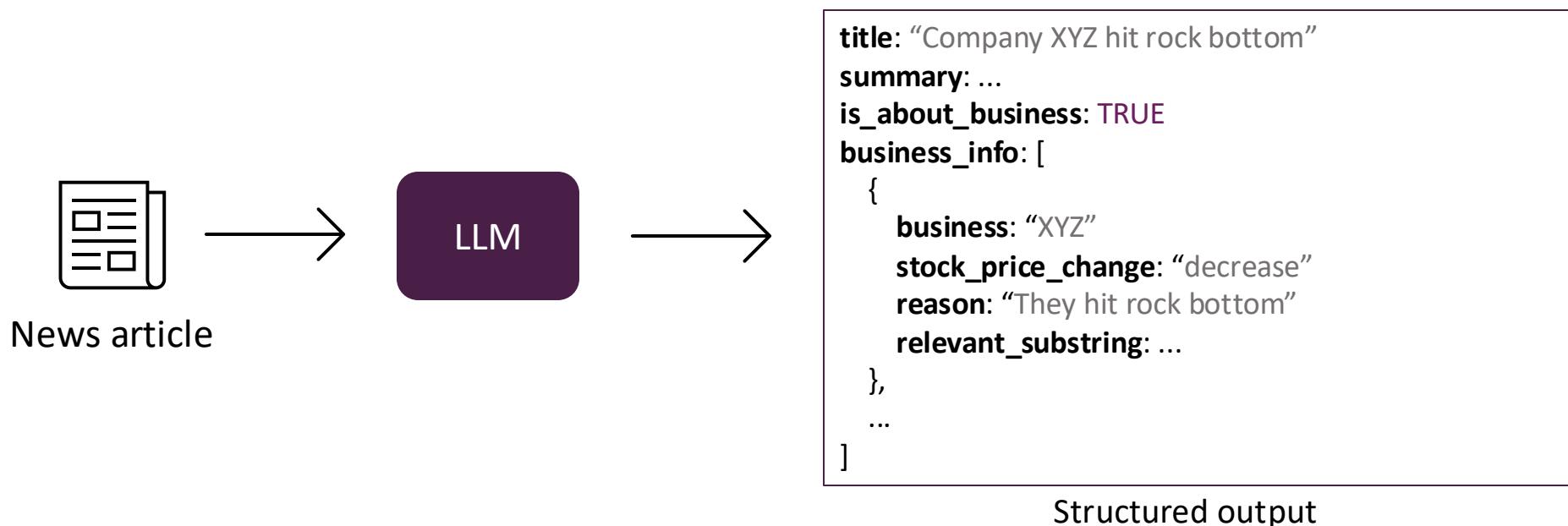
	MLOps	LLMOps
People	Data scientist (statistics, ML) Machine learning engineer Data engineer (structured data) Product owner ...	Data scientist (NLP) Machine learning engineer Data engineer ( <b>unstructured data</b> ) Product owner <b>Prompt engineer (?)</b> ...
Processes	Project templates (typical ML use cases) Experimentation (logging hyperparams and metrics) Pipelines (preprocessing and ML training) Deployment (model endpoint) Monitoring (typical ML metrics and infra) ...	Project templates (typical LLM use cases, RAG, etc.) Experimentation ( <b>logging prompts and metrics</b> ) Pipelines (preprocessing and LLM batch jobs) Deployment ( <b>web app</b> ) Monitoring (LLM metrics, <b>traces</b> and infra) ...
Technology	Scikit-Learn Tensorflow/Keras/PyTorch/Jax MLflow FastAPI ...	LangChain PydanticAI LLM-specific API* Streamlit LangSmith ...

\* Assuming we are not training our own LLM, but use a provider's LLM

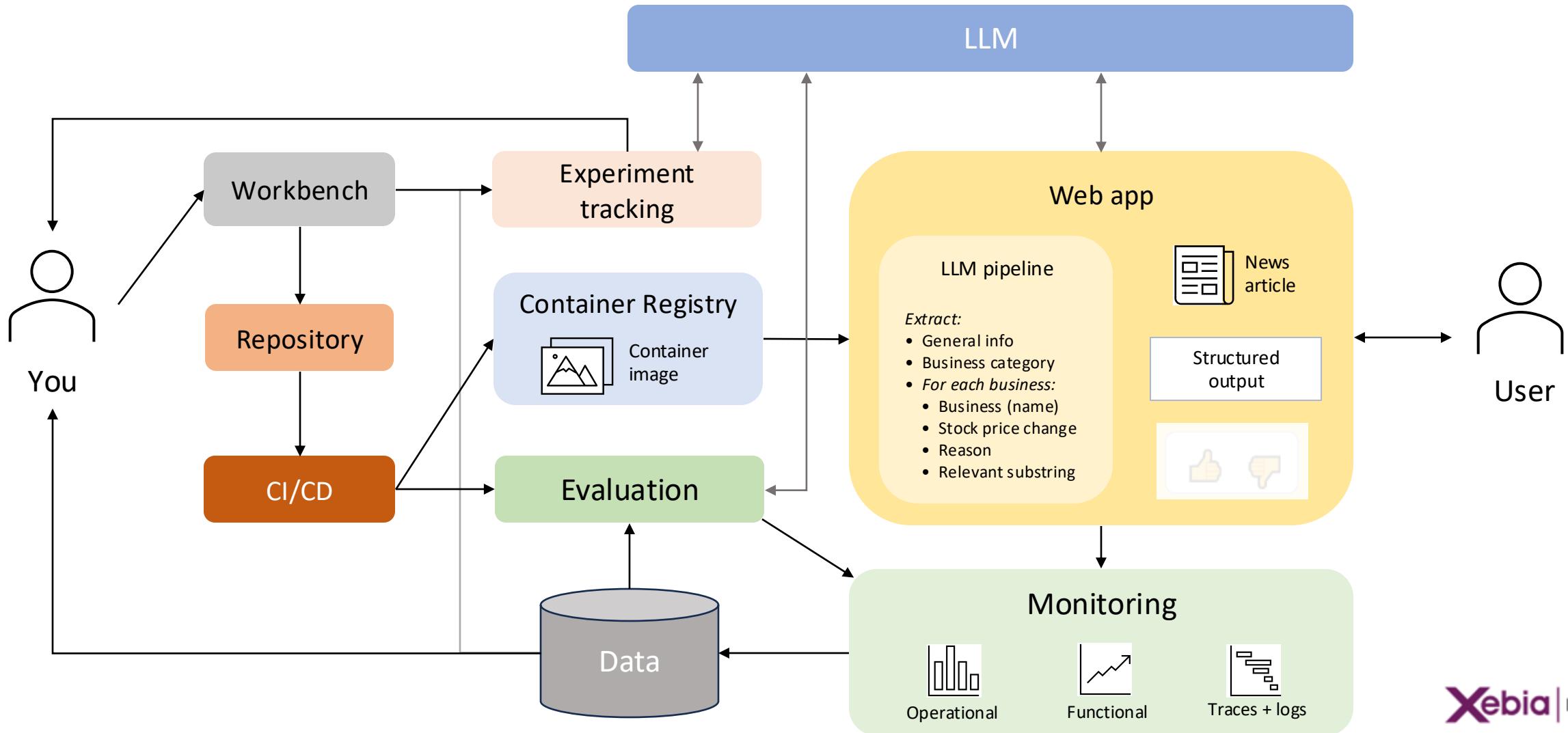


# News Reader: Extract structured info and predict stock price changes from news articles

- We work at a bank and are **interested in analyzing news articles for likely stock price changes**
- We have created a notebook with an **LLM-solution to extract relevant info from the articles**
- In the end, we want a **web app where users can upload their articles which are then automatically analyzed**



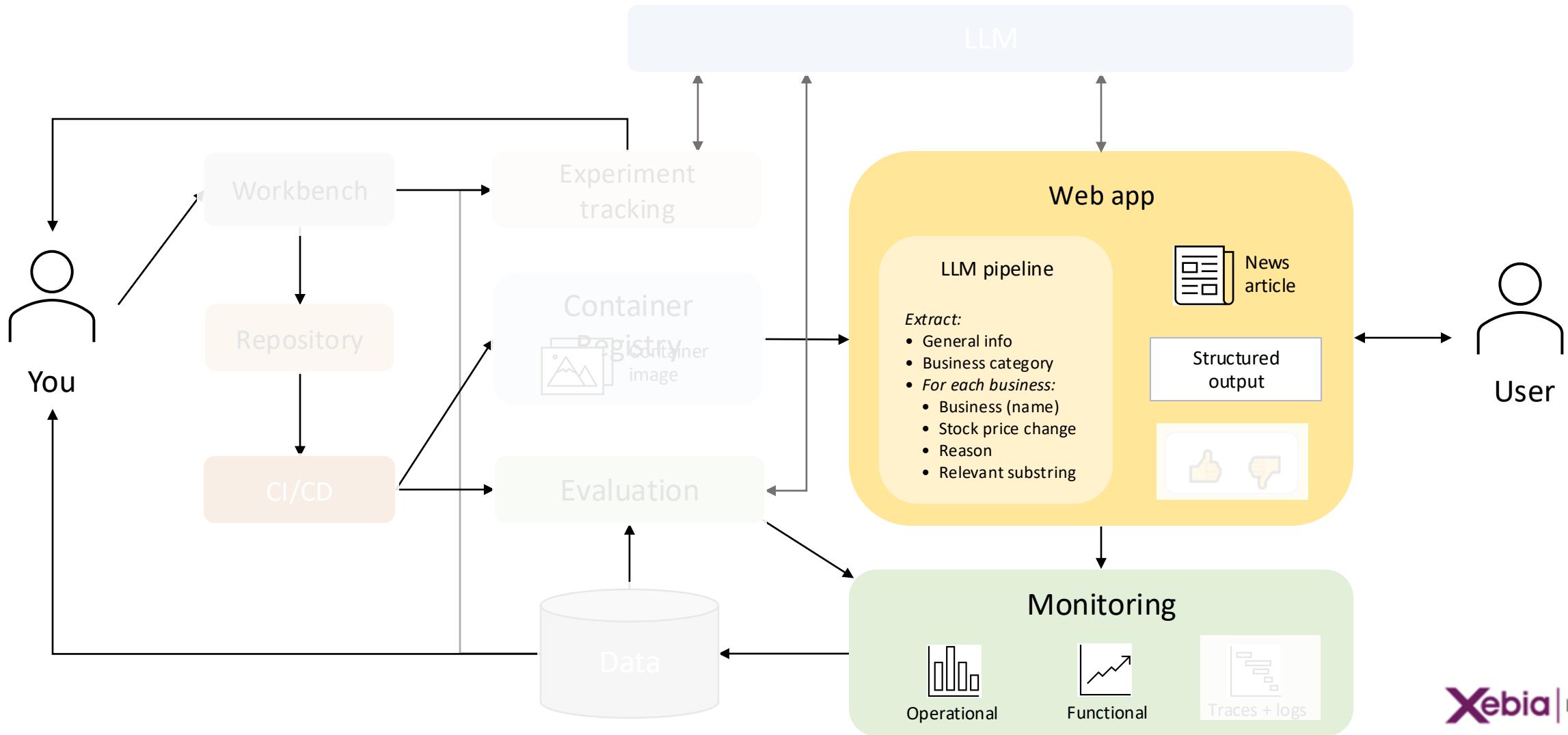
# High-level solution diagram



# Monitoring performance

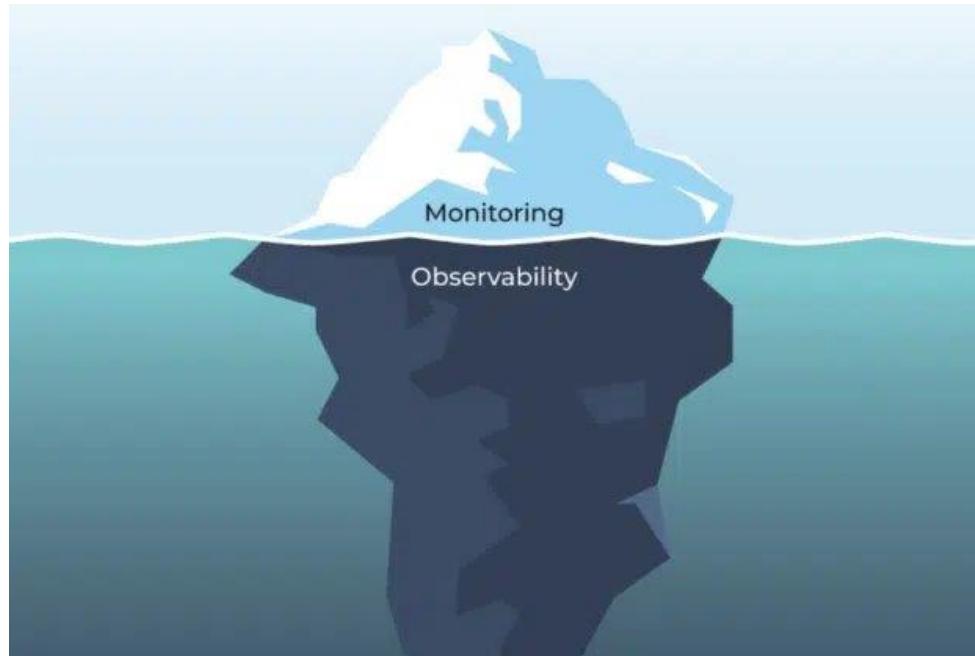


# High-level solution diagram



# Why monitor?

**Goal:** Quantify performance and flag potential issues in the system



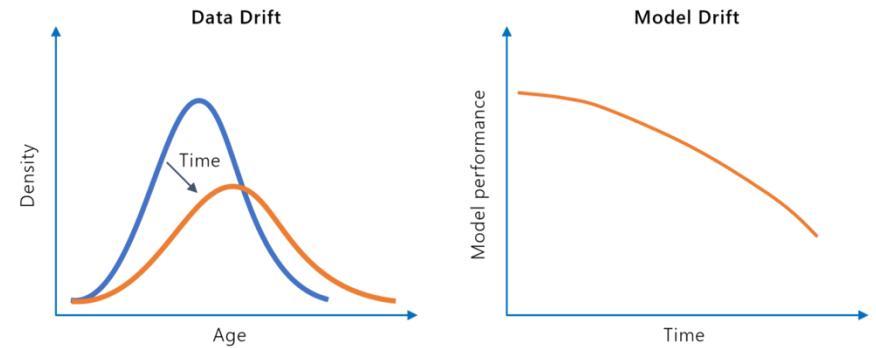
# What to monitor?

## Operational (infra-structure level)

- Number of requests from users
- Memory usage by the service
- Failed requests
- ...

## Functional (model/application level)

- Number of tokens to- and from LLM
- Number of requests to LLM
- How often output is saved by the user
- ...



# Logs, metrics and traces

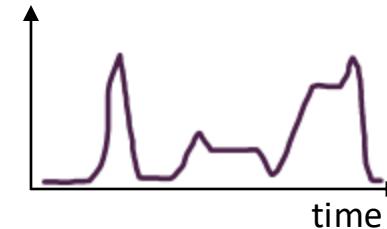
## Logs

- Give insight in what is going on inside application
- Can be structured (e.g. JSON) or natural language



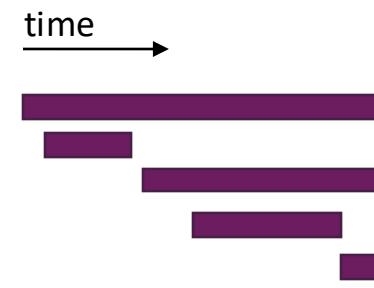
## Metrics

- Flag potential issues (with alerts) for investigation
- Operational: built-in, e.g. provided by Azure Services
- Functional: collected from our own (structured) logs



## Traces

- Follows the journey of a user/event through the entire stack
- Allows you to find bottlenecks or identify causes of issues



# Structured logging

```
print("This is not a log statement")
```

```
>> This is not a log statement
```

```
logger = logging.getLogger()
```

```
logger.info("This is a regular log statement")
```

```
>> INFO: This is a regular log statement
```

```
logger = structlog.get_logger()
```

```
... # some config
```

```
logger.info("Structured log", key1=42, key2="value")
```

```
>> {"event": "Structured log", "key1": 42, "key2": "value"}
```

# Azure Logging

The screenshot shows the Microsoft Azure portal interface for the 'test\_app\_insights' application. The left sidebar contains navigation links for Overview, Activity log, Access control (IAM), Tags, Diagnose and solve problems, Investigate (Application map, Smart detection, Live metrics, Transaction search, Availability, Failures, Performance, Troubleshooting guides (preview)), Monitoring (Alerts, Metrics, Diagnostic settings), Logs (selected), Workbooks, and Usage (Users, Sessions, Events, Funnels, User Flows, Cohorts). The main area displays a 'Logs' section with a search bar, a 'New Query 1\*' button, and a 'Select scope' dropdown set to 'test\_app\_insights'. Below this is a table titled 'traces' with columns: timestamp [UTC], message, severityLevel, itemType, and customDimensions. The table lists 2312 rows of log entries from June 19, 2023, at 5:55:44 PM to 5:56:50.625 PM. The first few entries include: 'Received HTTP response headers after 137.711ms - 404', 'End processing HTTP request after 160.4819ms - 404', 'The response has already started, the error page middleware...', 'Now listening on: https://localhost:7275', 'Now listening on: http://localhost:5212', 'Application started. Press Ctrl+C to shut down.', 'Hosting environment: Development', 'Content root path: C:\Users\samps\Source\Repos\docs\docs...', 'Sending greeting', 'Sending greeting, level 5', 'Start processing HTTP request GET https://localhost:7275/N...', 'Sending HTTP request GET https://localhost:7275/NestedGr...', 'Sending greeting, level 4', 'Start processing HTTP request GET https://localhost:7275/N...', 'Sending HTTP request GET https://localhost:7275/NestedGr...', 'Sending greeting, level 3', 'Start processing HTTP request GET https://localhost:7275/N...', 'Sending HTTP request GET https://localhost:7275/NestedGr...', 'Sending greeting, level 2', and 'Start processing HTTP request GET https://localhost:7275/N...'. The bottom of the table shows '3s 197ms | Display time (UTC+00:00)' and 'Query details | 19 - 40 of 2312'.

<https://learn.microsoft.com/en-us/dotnet/core/diagnostics/observability-applicationinsights>

# Exercise: Monitoring operational and functional metrics

**Let's define our own custom log-based metrics and monitor our solution in the cloud!**

- Run through **exercises 10 + 11** in the /exercises folder

# Break time!



# Intermezzo – Tools available



# Tools



<https://github.com/langchain-ai/langsmith-sdk>

- Structured in projects
- Within a project:
  - **Traces**: end-to-end invocations of your application, with input prompts and output responses and metadata
  - **LLM calls**: isolate only the LLM calls so you can isolate prompts and generations
  - **Monitoring**: track different statistics over time: Trace count, LLM call count, success rates, etc.
  - Can jump into any **LLM call** in a playground, changing the settings and testing the completions
  - Allows for **filtering** on attributes, e.g. feedback
  - Allows you to store and update **datasets** for testing and evaluation



<https://github.com/langfuse/langfuse>

## Develop

- **Observability**: Instrument your app and start ingesting traces to Langfuse ([Quickstart](#), [Integrations](#) [Tracing](#))
- **Langfuse UI**: Inspect and debug complex logs ([Demo](#), [Tracing](#))
- **Prompt Management**: Manage, version and deploy prompts from within Langfuse ([Prompt Management](#))
- **Prompt Engineering**: Test and iterate on your prompts with the [LLM Playground](#)

## Monitor

- **Analytics**: Track metrics (cost, latency, quality) and gain insights from dashboards & data exports ([Analytics](#))
- **Evals**: Collect and calculate scores for your LLM completions ([Scores & Evaluations](#))
  - Run model-based evaluations ([Model-based evaluations](#)) within Langfuse
  - Collect user feedback ([User Feedback](#))
  - Manually score observations in Langfuse ([Manual Scores](#))

## Test

- **Experiments**: Track and test app behaviour before deploying a new version
  - Datasets let you test expected in and output pairs and benchmark performance before deploying ([Datasets](#))
  - Track versions and releases in your application ([Experimentation](#), [Prompt Management](#))



<https://github.com/Arize-ai/phoenix>

- **Tracing** - Trace your LLM application's runtime using OpenTelemetry-based instrumentation.
- **Evaluation** - Leverage LLMs to benchmark your application's performance using response and retrieval evals.
- **Datasets** - Create versioned datasets of examples for experimentation, evaluation, and fine-tuning.
- **Experiments** - Track and evaluate changes to prompts, LLMs, and retrieval.
- **Inference Analysis** - Visualize inferences and embeddings using dimensionality reduction and clustering to identify drift and performance degradation.

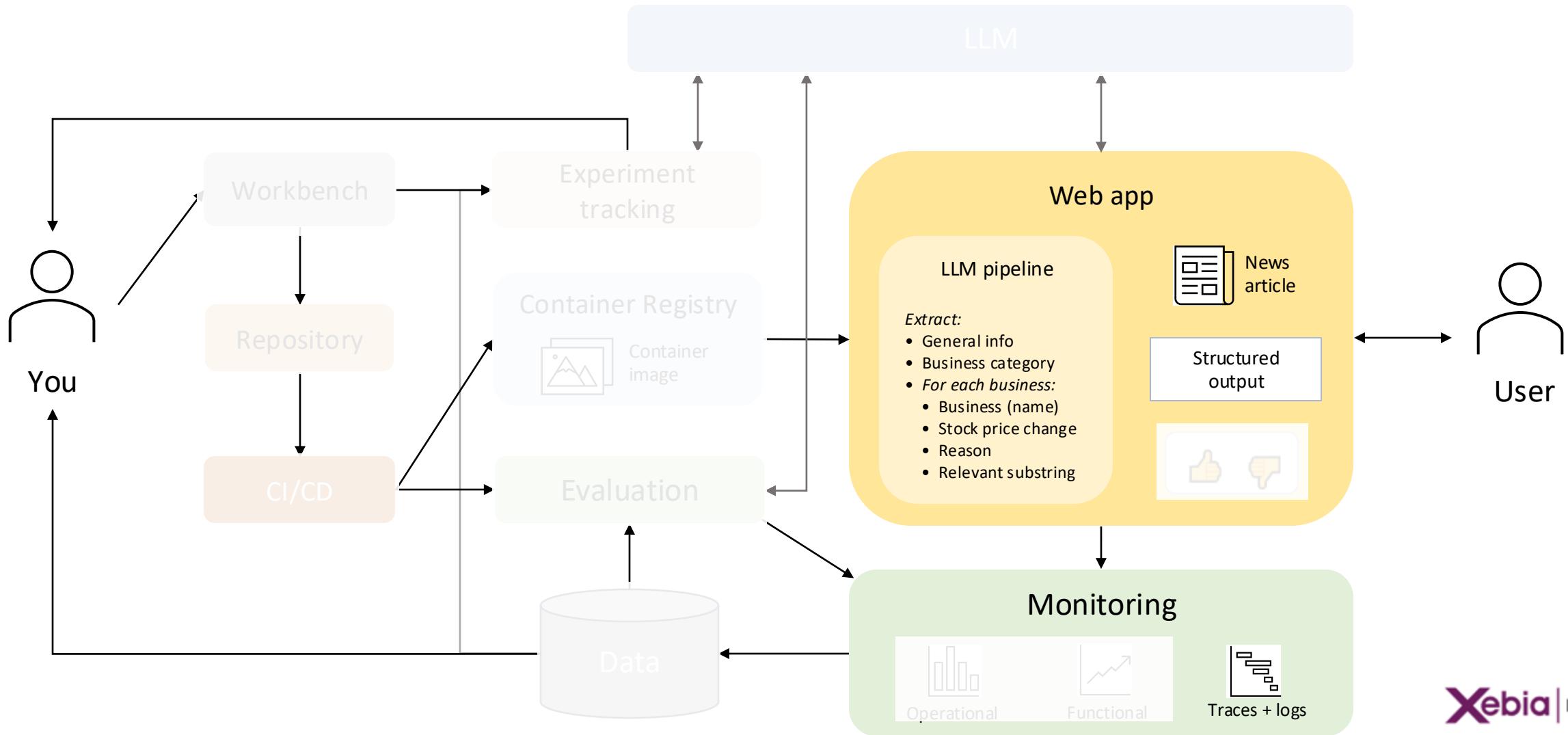
# Common features

- **Monitoring:** tracking LLM calls, tokens, etc.
- **Tracing:** insight in end-to-end flow of LLM system
- **Logging/**debugging single LLM calls: to isolate good or bad examples
- **Dataset management:** quickly add new labeled samples (based on logs)
- **Prompt management:** store and version prompts
- **Evaluations:** run LLM on dataset to collect predefined metrics
- **Experimentation:** keep track of parameters, prompt, models used
- **Iteration:** allow for quick iteration with playground and log insights

# Tracing



# High-level solution diagram

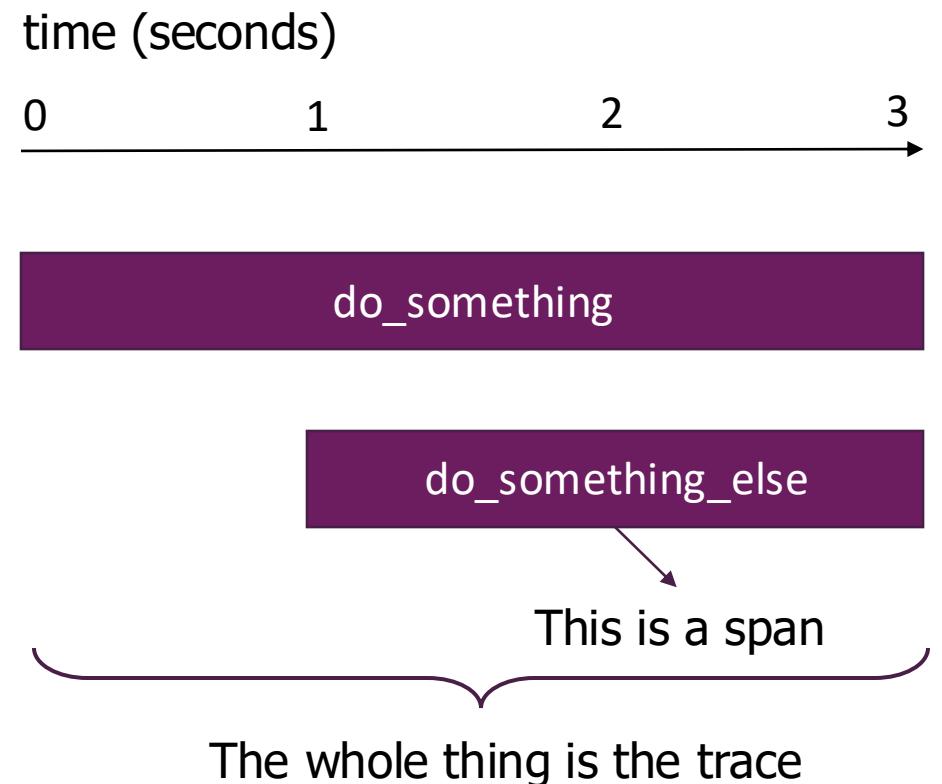


# What is a trace?

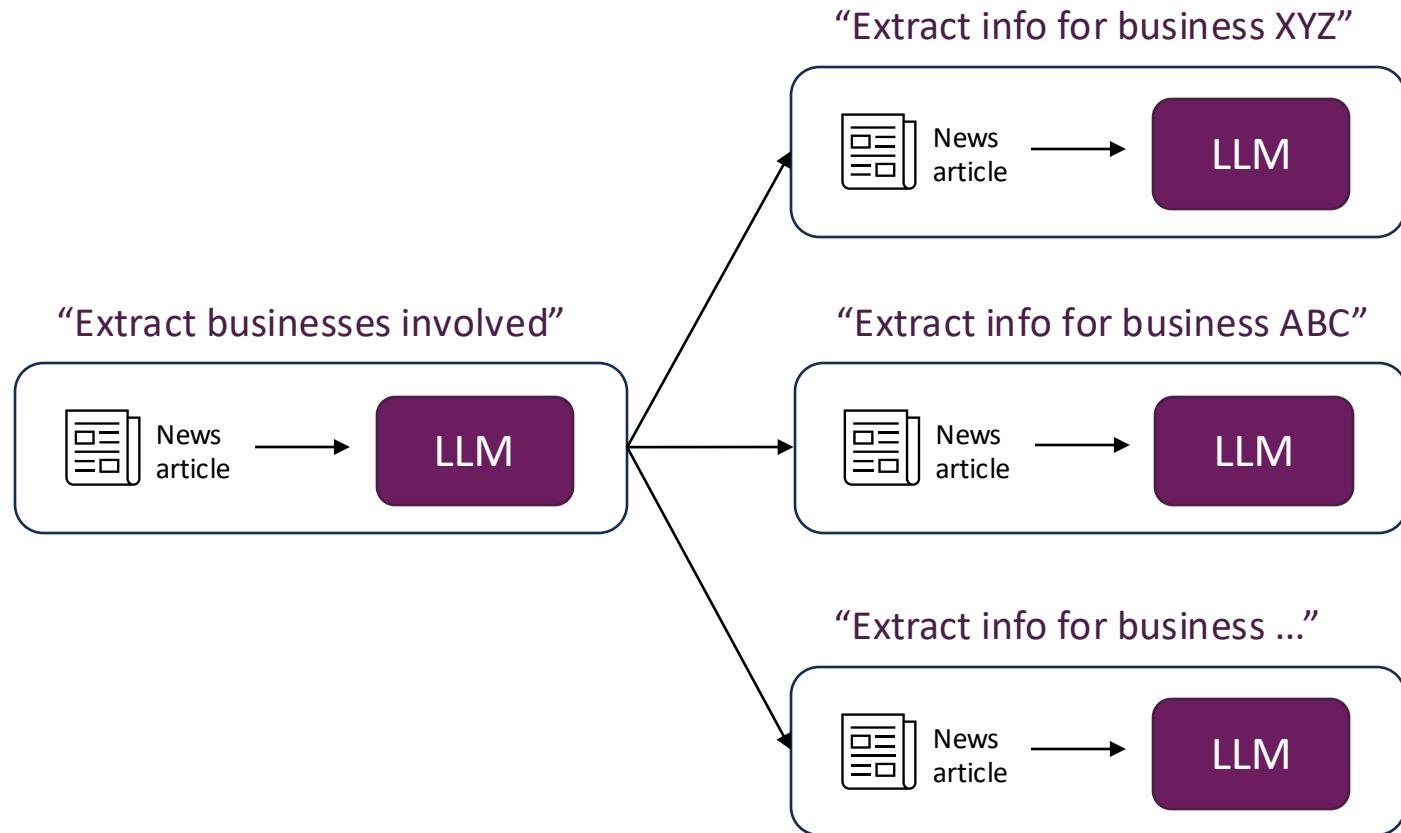
```
@tracer.start_as_current_span("do_something_else")
def do_something_else():
    time.sleep(2)

@tracer.start_as_current_span("do_something")
def do_something():
    time.sleep(1)
    do_something_else()

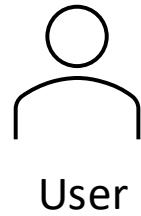
do_something()
```



# Why tracing?



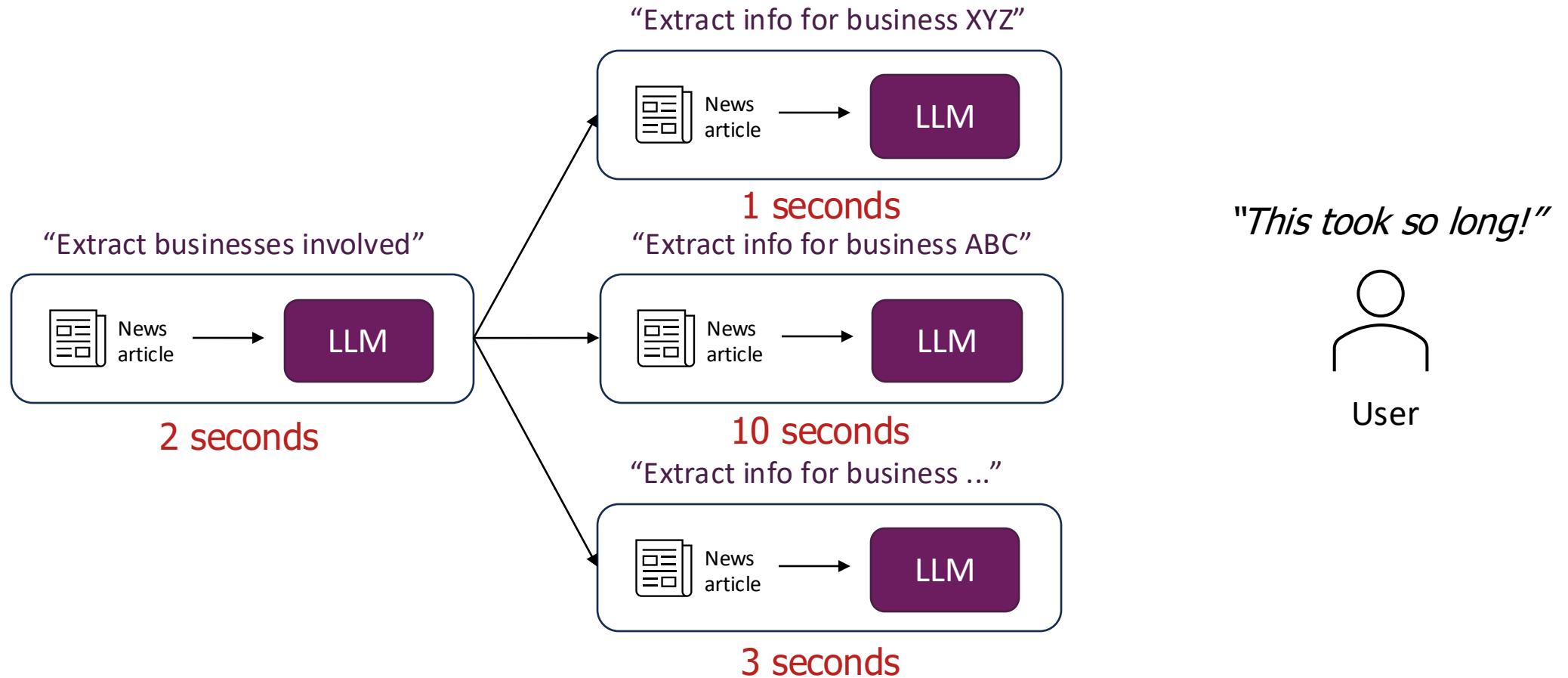
*"This output is bad!  
XYZ is not a business"*



User

💡 We need insight into the entire trace to investigate **where mistakes originate from**

# Why tracing?



We need insight into the entire trace to investigate **where the bottlenecks are.**

# OpenTelemetry

- Open-source observability framework
- Common framework for tracing
- Integrates nicely with other tools, e.g. Phoenix, Azure Monitor



# Tracing

The screenshot shows the Microsoft Azure Application Insights interface for "End-to-end transaction details". The main area displays a hierarchical trace tree for an operation ID. The tree starts with a root node for "localhost:7275 GET /NestedGreeting" and branches down into multiple levels of nested activities, each with its duration and response status. To the right of the trace tree, there is a detailed view of a specific activity named "GreeterActivity". This view includes sections for "Internal Properties" (Event time, Type, Call status, Duration, Name), "Custom Properties" (nest-level), "Command" (GreeterActivity), and "Related Items". The bottom of the page shows a summary of "Traces & events" (26) and "Events" (0).

Microsoft Azure

Home > test\_app\_insights | Application map > Performance > End-to-end transaction details

Search results

Filtered on timestamp > 6/19/2023, 10:00 AM, timestamp < 6/19/2023, 11:00 AM, Operation name == GET /NestedGreeting, client.Type != ...

Suggested

All Sort by Relevance

6/19/2023, 11:29:28 AM GET /NestedGreeting Duration: 159.2 ms Response code: 200

6/19/2023, 11:29:28 AM GET /NestedGreeting Duration: 159.2 ms Response code: 200

6/19/2023, 11:29:29 AM GET /NestedGreeting Duration: 49.6 ms Response code: 200

6/19/2023, 11:29:29 AM GET /NestedGreeting Duration: 125.2 ms Response code: 200

6/19/2023, 11:29:29 AM GET /NestedGreeting Duration: 87.0 ms Response code: 200

6/19/2023, 11:29:29 AM GET /NestedGreeting Duration: 2.2 ms Response code: 200

6/19/2023, 11:29:28 AM GET /NestedGreeting Duration: 184.2 ms Response code: 200

6/19/2023, 11:28:21 AM GET /NestedGreeting Duration: 128.8 ms Response code: 200

6/19/2023, 11:28:21 AM GET /NestedGreeting Duration: 2.5 ms Response code: 200

6/19/2023, 11:28:21 AM

End-to-end transaction

Operation ID: 60000000-0000-0000-0000-000000000000

EVENT RES. DURATION

localhost:7275 GET /NestedGreeting 200 184.2 ms

localhost:7275 GET /NestedGreeting 183.5 ms

localhost:7275 GET /NestedGreeting 21.5 ms

localhost:7275 GET /NestedGreeting 159.2 ms

localhost:7275 GET /NestedGreeting 158.1 ms

localhost:7275 GET /NestedGreeting 21.9 ms

localhost:7275 GET /NestedGreeting 125.2 ms

localhost:7275 GET /NestedGreeting 125.1 ms

localhost:7275 GET /NestedGreeting 24.0 ms

localhost:7275 GET /NestedGreeting 87.0 ms

localhost:7275 GET /NestedGreeting 86.9 ms

localhost:7275 GET /NestedGreeting 22.5 ms

localhost:7275 GET /NestedGreeting 49.6 ms

localhost:7275 GET /NestedGreeting 49.1 ms

localhost:7275 GET /NestedGreeting 19.5 ms

localhost:7275 GET /NestedGreeting 2.2 ms

localhost:7275 GET /NestedGreeting 1.6 ms

GreeterActivity

Traces & events (5)

Internal Properties

Event time: 6/19/2023, 11:29:28.9817375 AM (Local time)

Type: InProc

Call status: true

Duration: 158.1 ms

Name: GreeterActivity

Custom Properties

nest-level: 4

Command: GreeterActivity

Related Items

Show what happened before and after this dependency in User Flows

All available telemetry 5 minutes before and after this event

Traces & events 26 Traces 0 Events

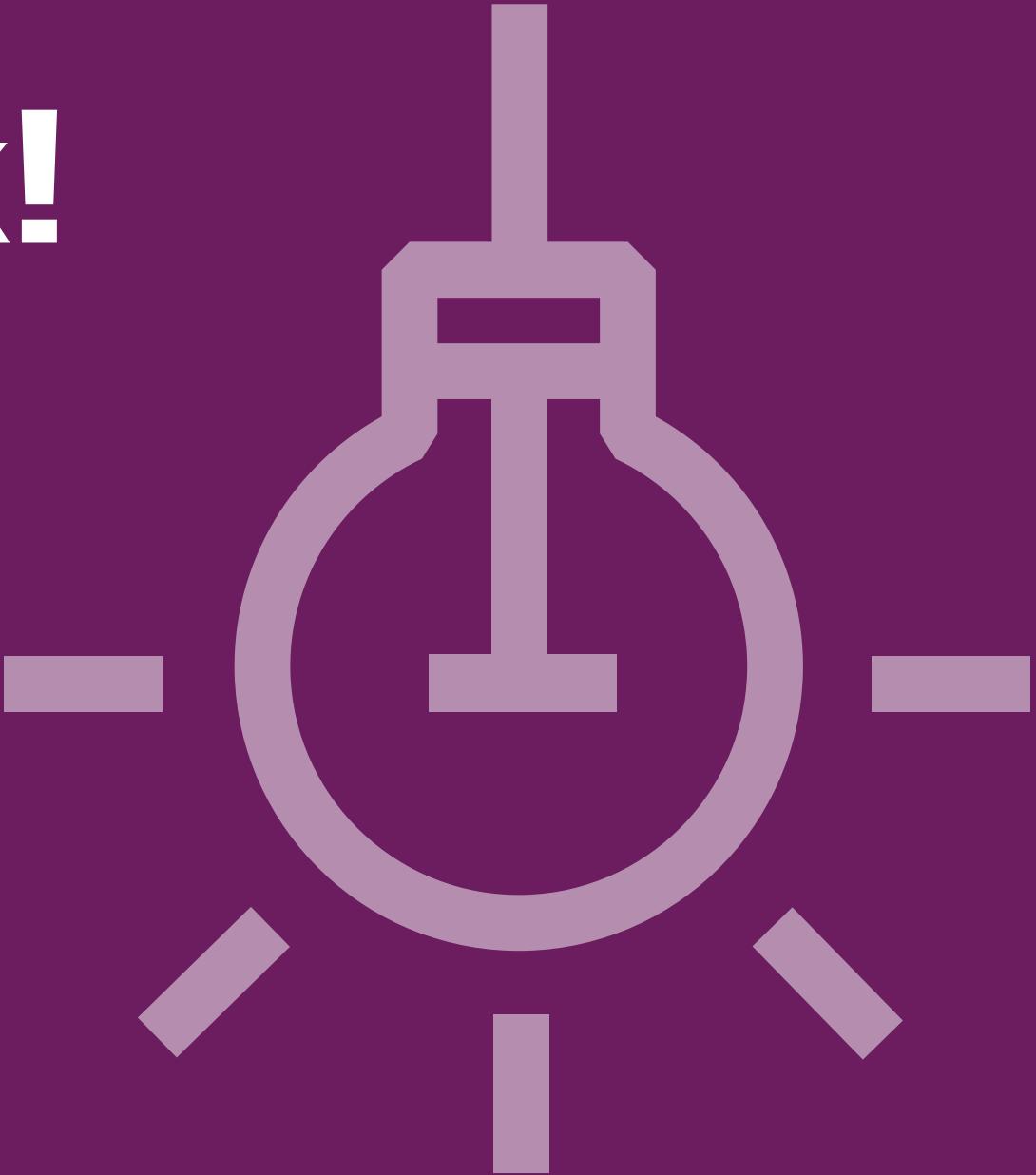
<https://learn.microsoft.com/en-us/dotnet/core/diagnostics/observability-applicationinsights>

# Exercise: Logging traces to the cloud

**Let's extend our monitoring solution by logging traces for the steps in our LLM pipeline**

- Run through **exercise 12** in the /exercises folder

# Lunch break!



# Schedule

## Day 1

- Intro MLOps and LLMOps
- Use case & solution design
- Building modularized LLM app
- Lunch 
- Evaluation and experiment tracking
- Containerization and deployment
- Automation with CI/CD

## Day 2

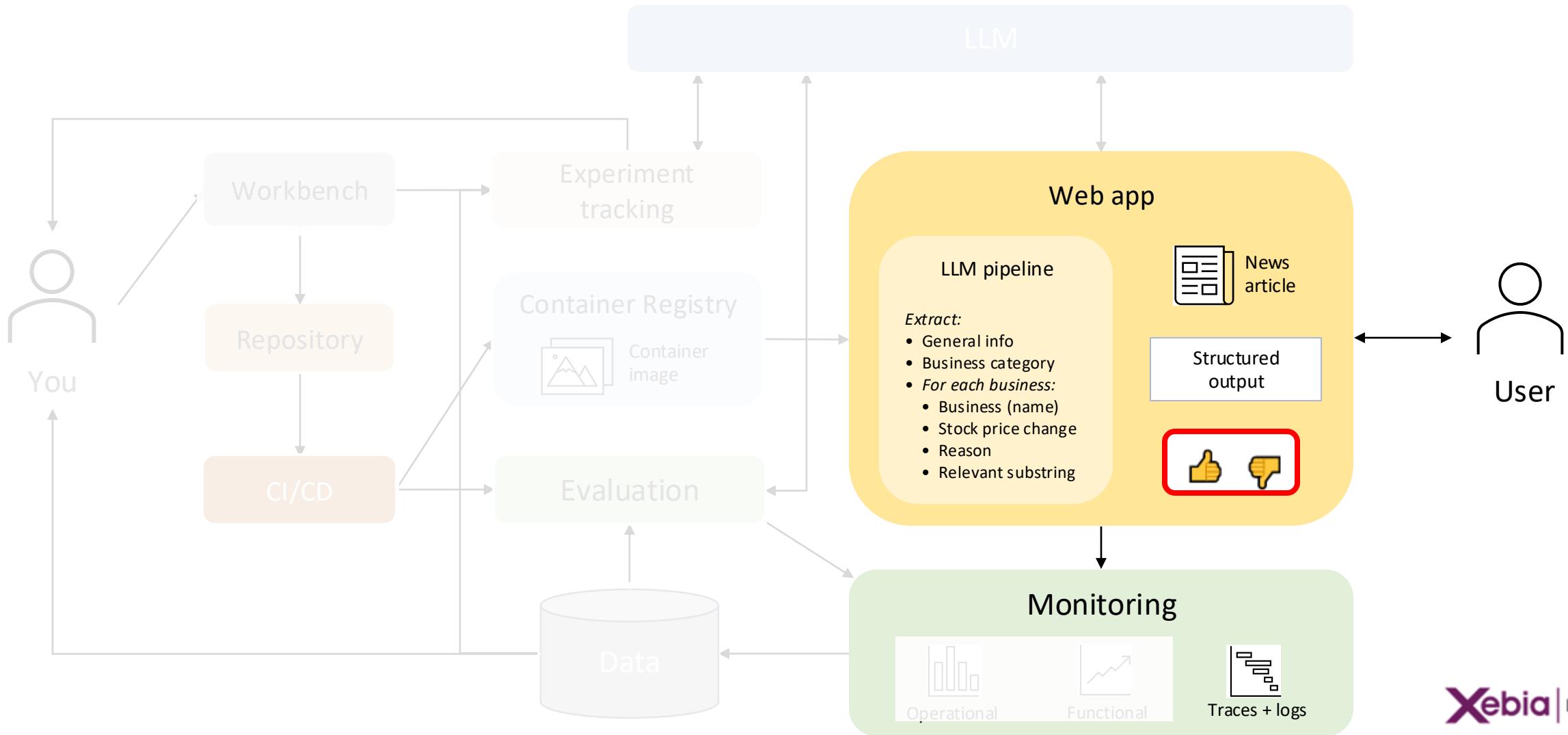
- Intro LLM system monitoring
- Tracking online metrics
- Logging traces
- Lunch 
- User feedback
- Iterating on the solution
- Wrap up

**Repo:** <https://github.com/godatadriven/academy-llmops-in-azure>

# Collecting user feedback



# High-level solution diagram



# User feedback provides an invaluable source of information. How to collect it?

## News Reader

### Selected Article

ESPN, Fox and Warner in US sports streaming tie-up - BBC News

Three US media giants have announced a new sports streaming platform to be launched in autumn. Walt Disney's ESPN, Fox Corp and Warner Bros. Discovery own a wide range of portfolios of sports rights including those for the FIFA World Cup, Formula 1, NFL, NBA and Major League Baseball. They hope to capture younger audiences and save costs. The service would have "a new brand with an independent management team" which would be available via a new app. Its pricing will be announced at a later date, according to their statement. The firms say each company will own one-third of the joint venture and have equal board representation. The announcement comes as sports leagues charge more for broadcasting rights. Instead of just linear TV operators, the fees are increasingly split between multiple media distributors. The product will bring sports linear networks and Disney's direct-to-consumer ESPN+ together, according to their statement. They added that the new joint venture aims "to serve sports fans, particularly those outside of the traditional pay TV bundle". Subscribers would also have the ability to bundle the product with the companies' streaming platforms Disney+, Hulu and Max. Bob Iger, the chief executive officer of

### Structured Output

`title` : ESPN, Fox and Warner in US sports streaming tie-up

`summary` : ESPN, Fox Corp and Warner Bros. Discovery are launching a new sports streaming platform in autumn to capture younger audiences and save costs.

`is_about_business` : TRUE

`business_info` :

▶ { ... }

▶ { ... }

▼ {

"business" : "Warner Bros. Discovery"

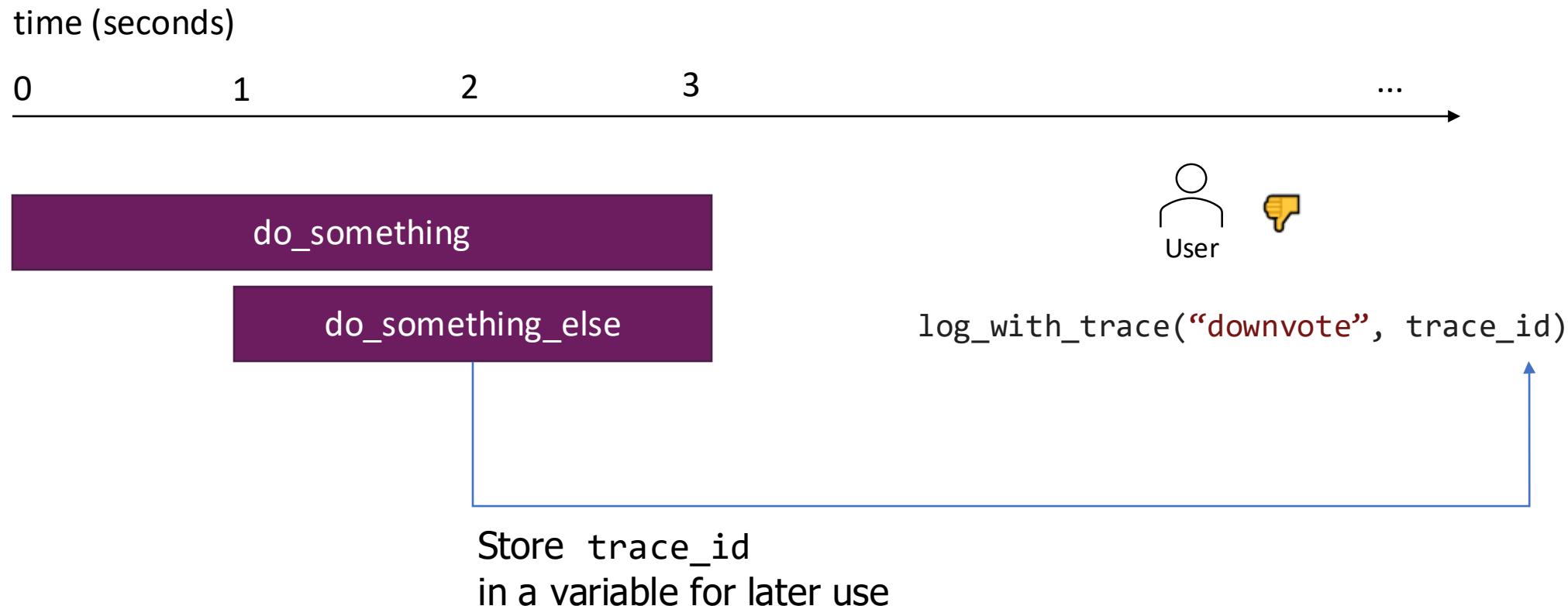
"stock\_price\_change" : "increase"

"reason" :

"The new streaming platform will drive



# User feedback provides an invaluable source of information. How to associate it to a trace?



# Exercise: Collect user feedback for LLM outputs

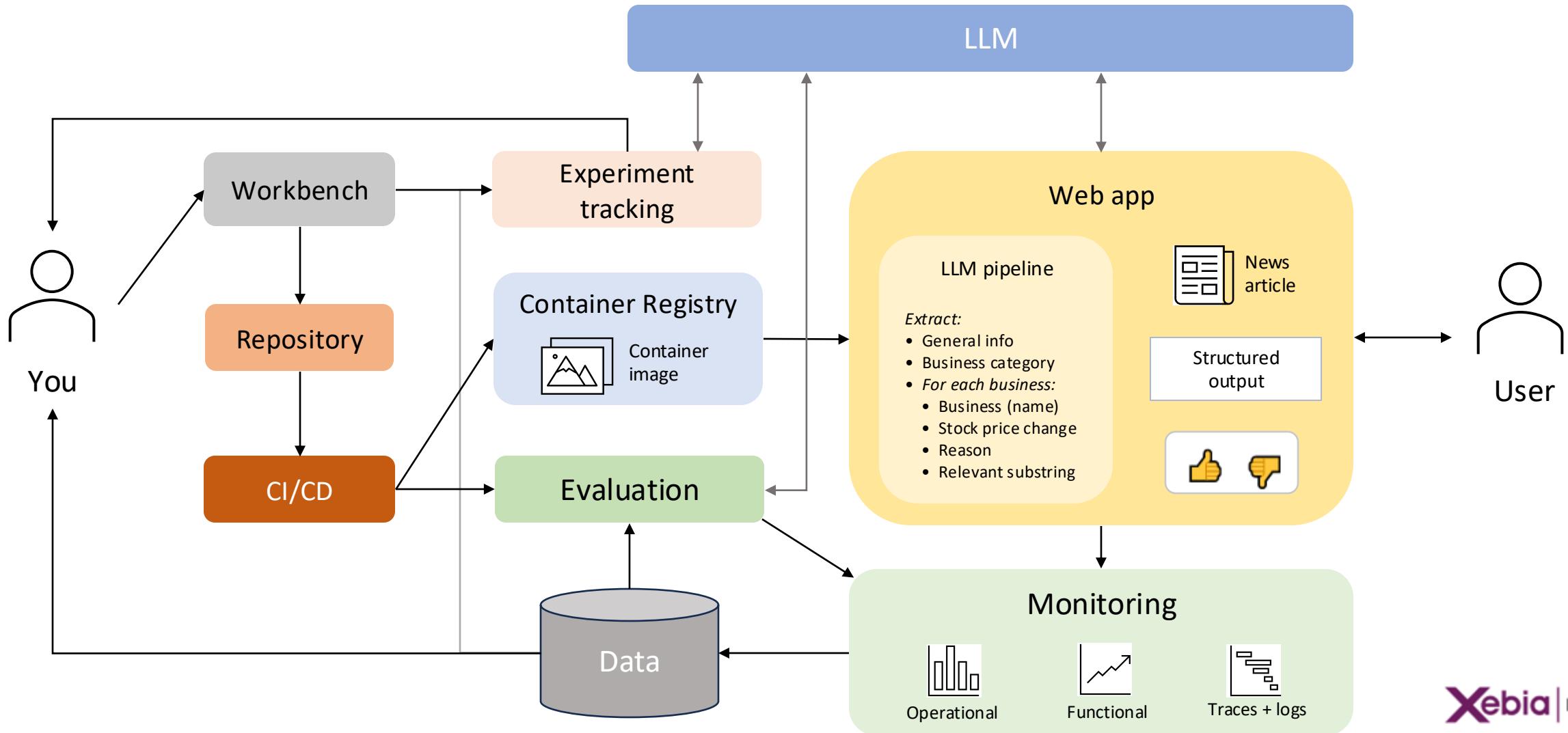
**Let's extend the app with buttons to collect and log user feedback to the cloud!**

- Run through **exercise 13** in the /exercises folder

# Iterating on the system



# High-level solution diagram



# Exercise: Iterate on the solution

**Choose a direction to iterate on, work together (and present results at the end) !**

- For inspiration, see **exercise 14** in the /exercises folder

# Agent Orchestration

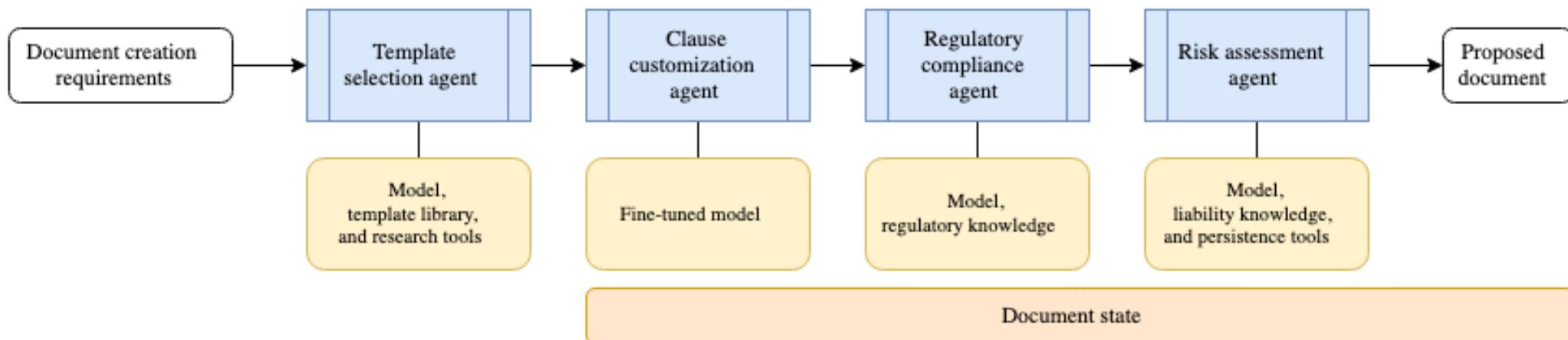


# Agent Orchestration

- Coordinate specialized agents to solve complex, multidisciplinary tasks.
- Enable scalability and parallelism by separating concerns into independent units.
- Improve maintainability, testing, and iteration through isolated responsibilities.
- Optimize latency and cost by choosing the right pattern (sequential vs concurrent).
- Provide clear boundaries for observability, security, and governance.

# Sequential Agent

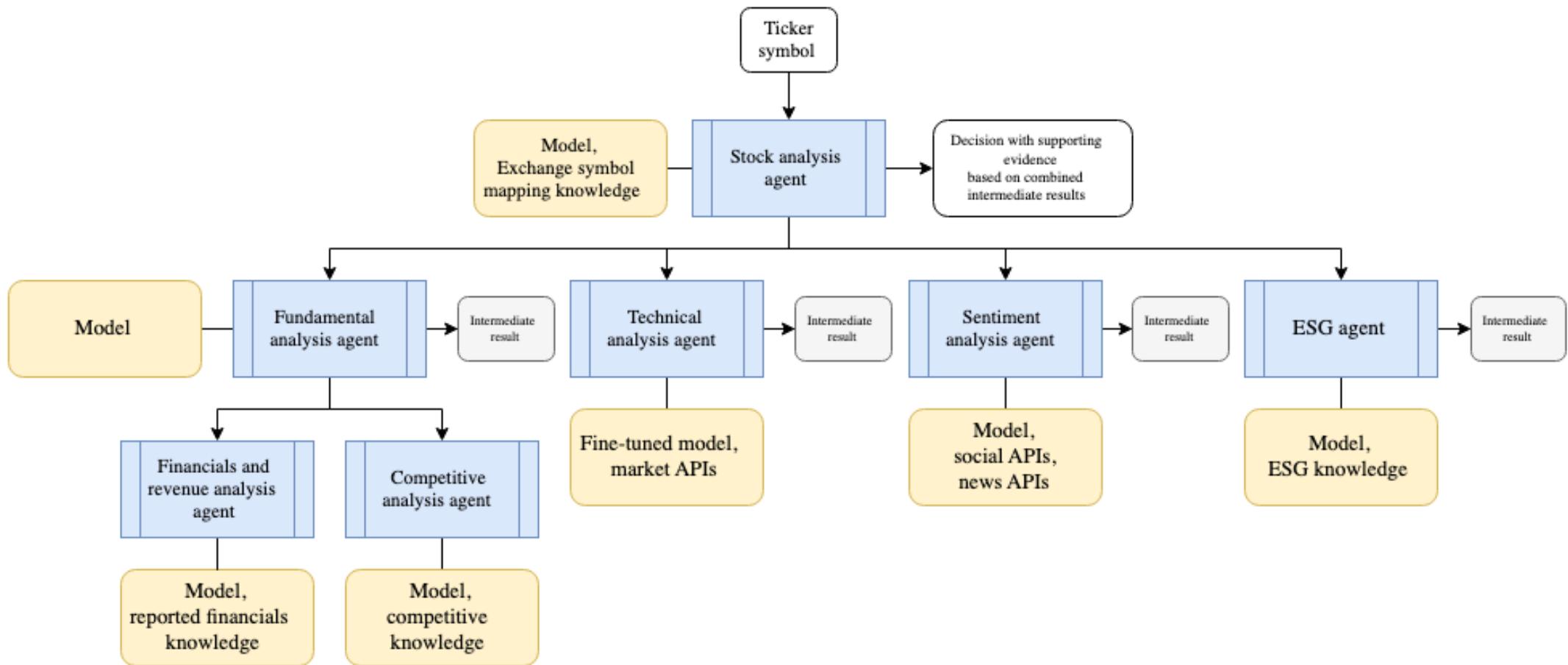
- A deterministic pipeline where each agent consumes previous agent output.
- template selection > clause customization > compliance check > risk assessment.



- Avoid when parallelization is required.

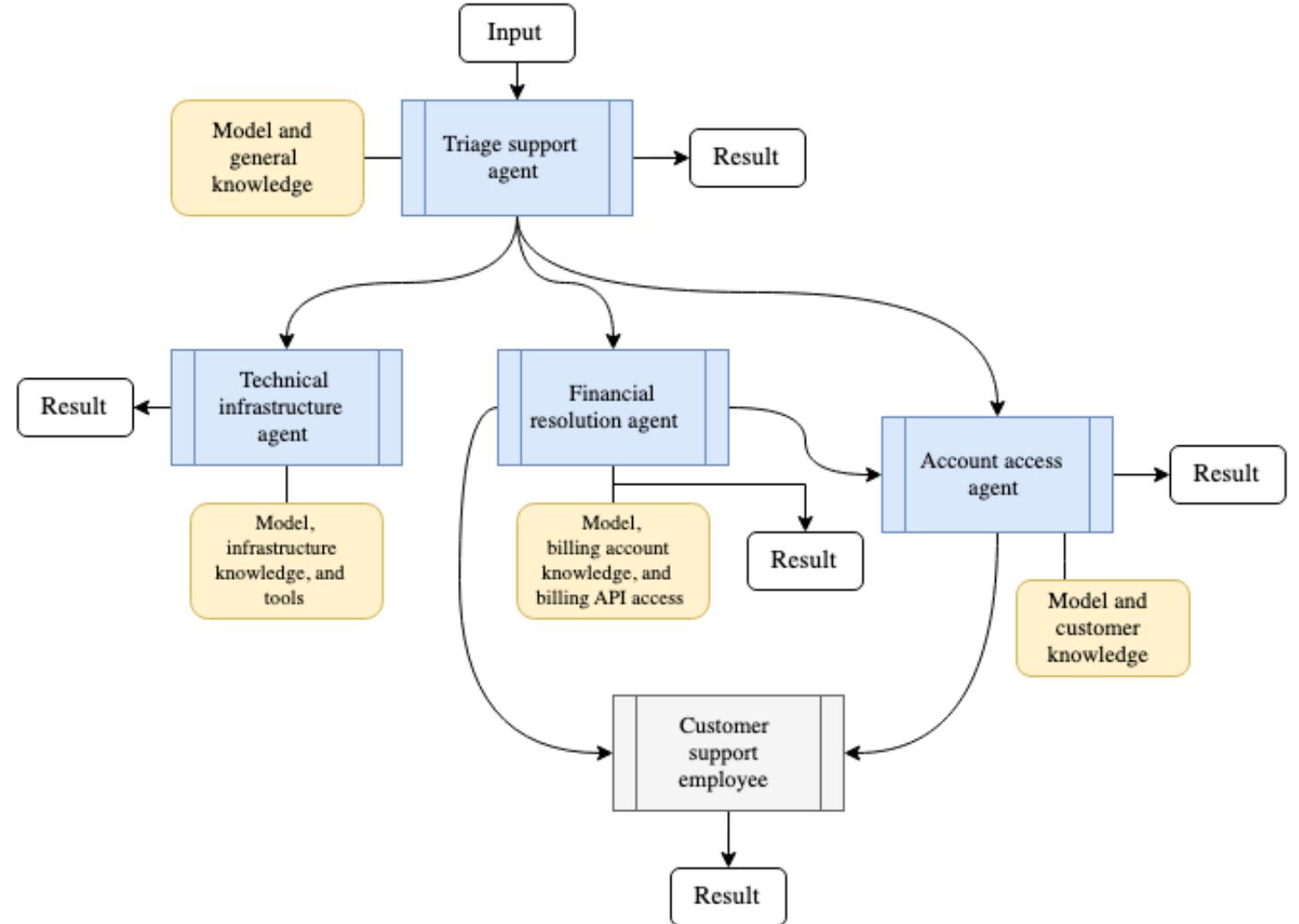
# Concurrent Orchestration

- Can run in parallel, either by using a fixed set of agents or by dynamically choosing AI agents based on specific task requirements.
- multi-agent decision-making, low latency.



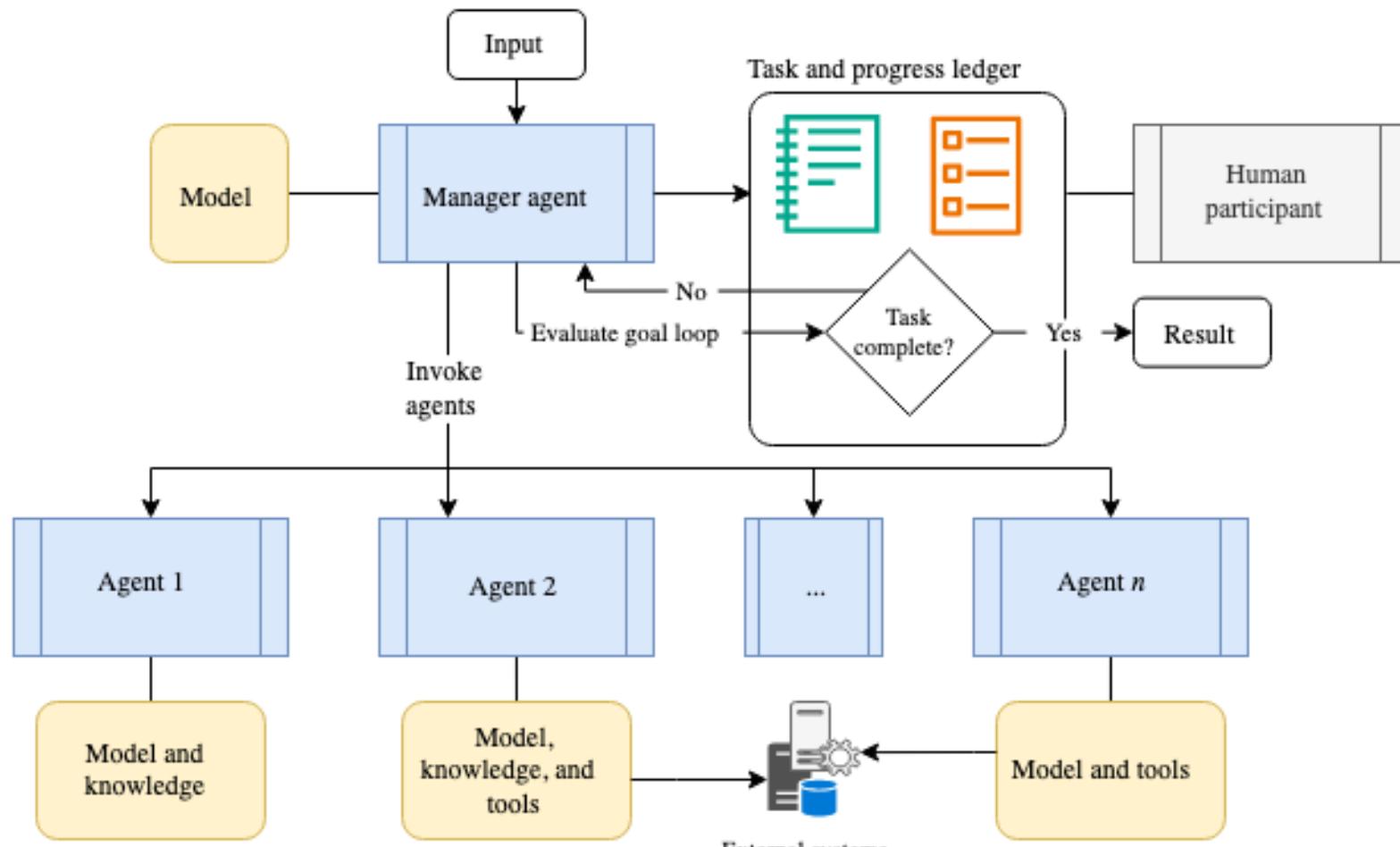
# Handoff Orchestration

- dynamic routing/handoff; one agent passes task to another based on capability or context.
- triage workflows where the optimal handler isn't known upfront.



# Magnetic Orchestration

- open-ended and complex problems that don't have a predetermined plan of approach



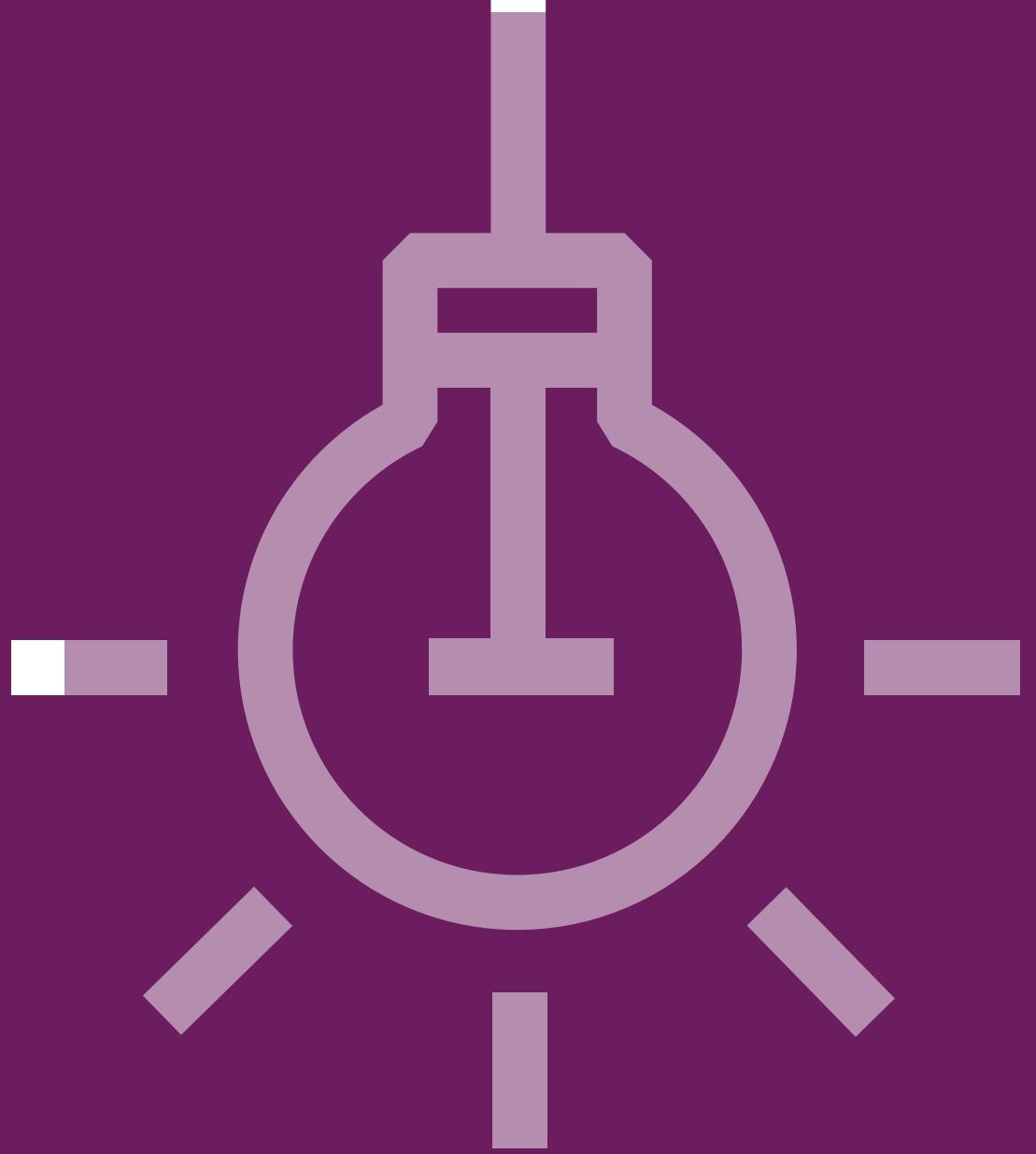
# Async Agents

- Why async?
  - Speed: parallelize independent LLM calls to reduce total wall time.
  - Scalability: better utilization of I/O-bound model calls; useful for concurrent pattern or extracting many subitems.
- When to use:
  - Business-info extraction for many entities in same article (concurrent per-entity calls).
  - Evaluation across many articles (bulk eval).
- Caveats:
  - Respect model rate limits and quotas; add concurrency limits/backpressure.
  - Add timeouts, retries with exponential backoff, circuit breakers.
  - Preserve traces/observability per logical task.

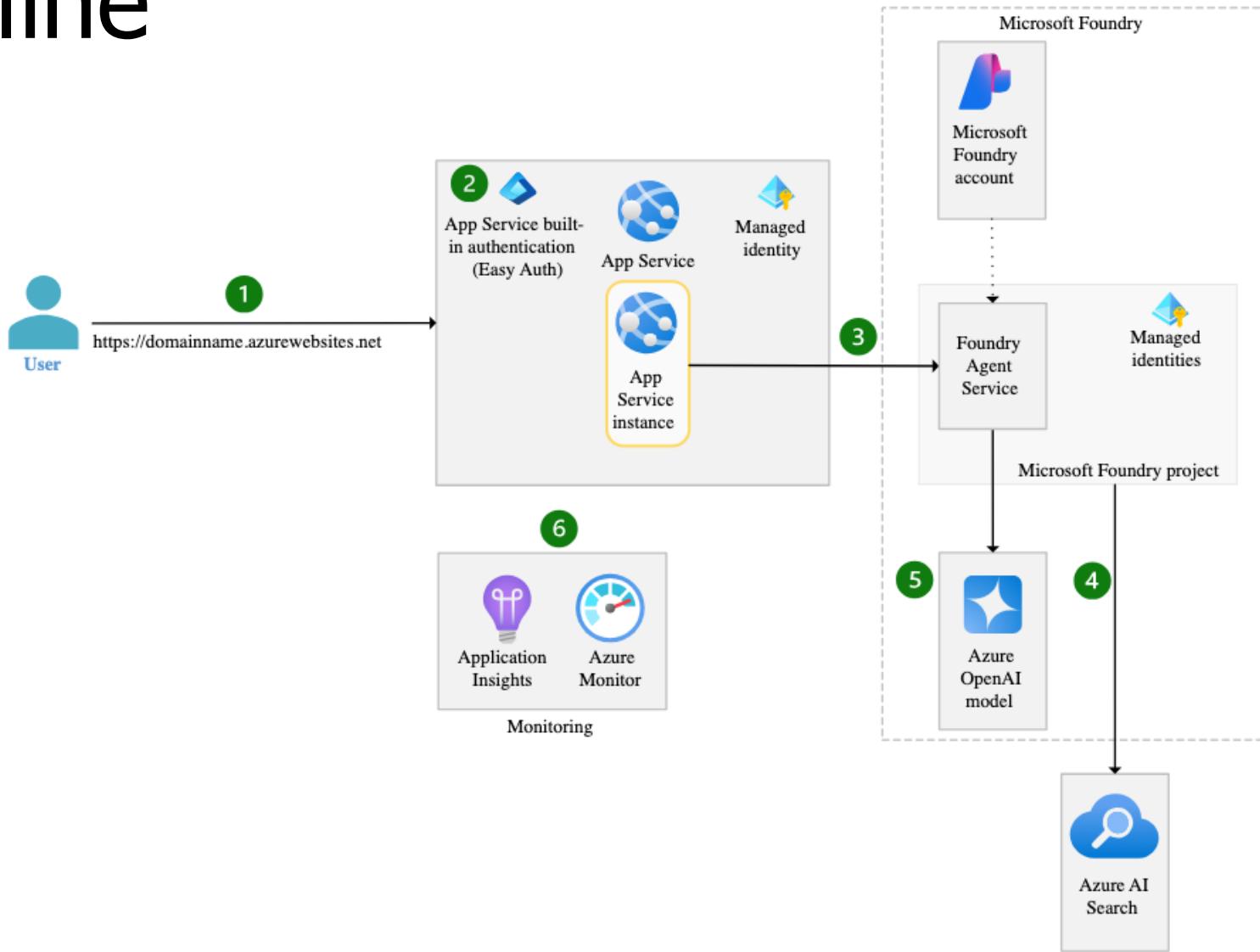
# Implementation

- Routing & decision logic: deterministic vs. dynamic.
- Context handling: pass full vs. summarized context between agents.
- Retry and timeouts: per-agent timeouts + backoff.
- Observability: tracing, structured logs, per-agent metrics, distributed traces.
- Security: auth between agents, least privilege for tools, PII redaction.
- Testing: unit-test agents independently; integration tests for combined flows.
- Resource management: rate limits, compute quotas, batching and cost control.

# Architecture Patterns



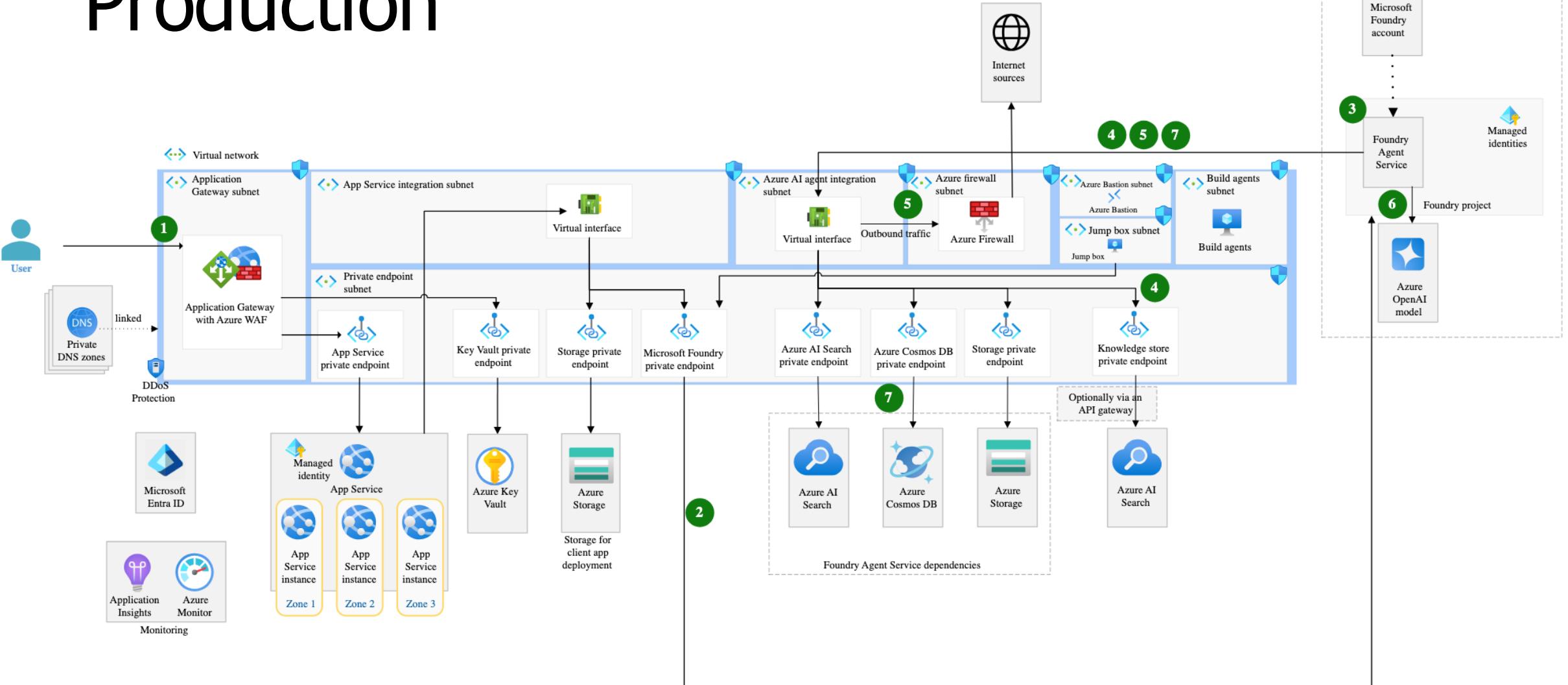
# Baseline



# Production

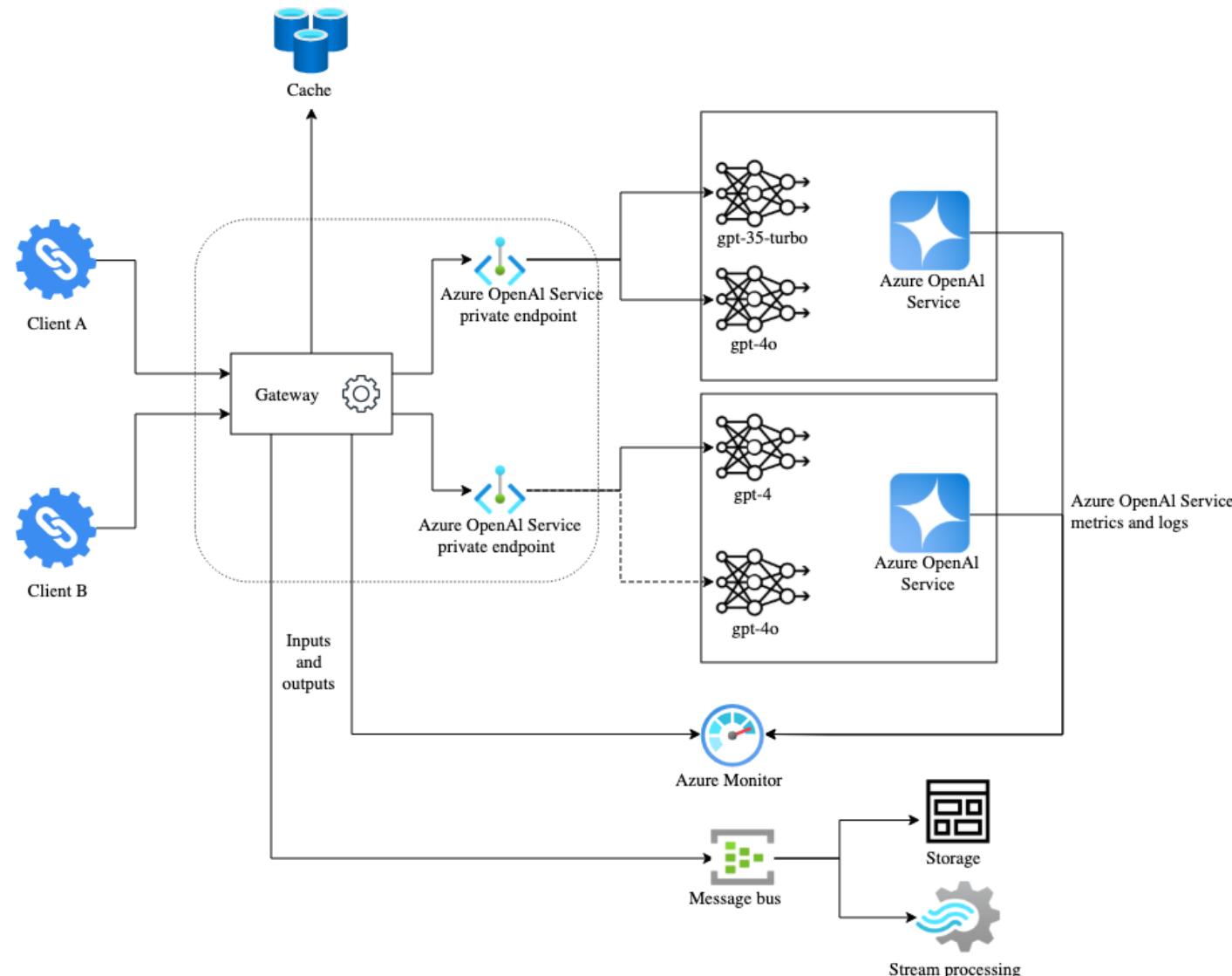
- A single, secure entry point for all chat UI traffic, which minimizes the attack surface
- Filtered ingress and egress network traffic by using a combination of NSGs, a web application firewall, user-defined routes (UDRs), and Azure Firewall rules
- End-to-end encryption of data in transit by using TLS
- Network privacy by using Private Link for all Azure PaaS service connections
- Zero Trust Policy for all resources

# Production



# API Gateway

- Capability to log both the client's request and the final response that it received
- For real-time monitoring, the logs can be published to a message bus.
- Beneficial when you have multiple clients or multiple models that you need to monitor

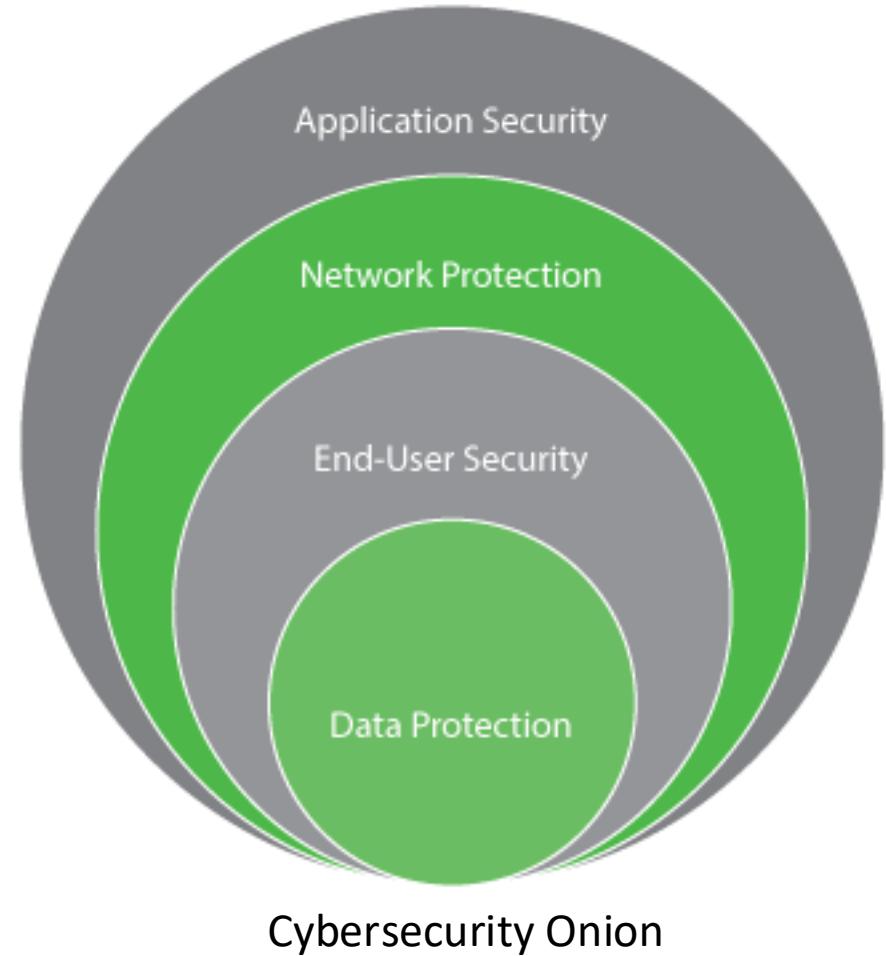


# Networking and Security



# Key Risks

- Data leakage (inputs/outputs sent to external models)
- Access to resources
- Prompt injection, model poisoning, hallucination risks
- Supply-chain compromise (dependencies, images, models)
- Runaway Costs



## IAM

- Central identity (Entra ID / OIDC), RBAC and least-privilege for services and humans
- Managed identities for service-to-service auth (no secrets in code)
- Conditional Access & MFA for operator consoles and CI/CD access
- Federate on-prem identities with Entra ID (Azure AD Connect or Workload Identity Federation)

## Data Protection

- Encrypt data in transit and at rest.
- Classify, redact or tokenise PII before external calls
- Differential privacy or synthetic data for training/fine-tuning
- Data residency within regions and private networks

# Networking

- Define trust boundaries — map data flows for ingress, storage, inference, egress and logs.
- Use VNet integration for inference, storage, and control-plane services to avoid public exposure.
- Private Endpoints / Private Link for managed services (OpenAI, Key Vault, Storage, Search)
- Protect ingress with App Gateway + WAF (or similar) and avoid direct public app endpoints for sensitive UIs/APIs.
- Upload critical logs to SIEM tools like Sentinel/Defender.



# Introduction to Fine-tuning

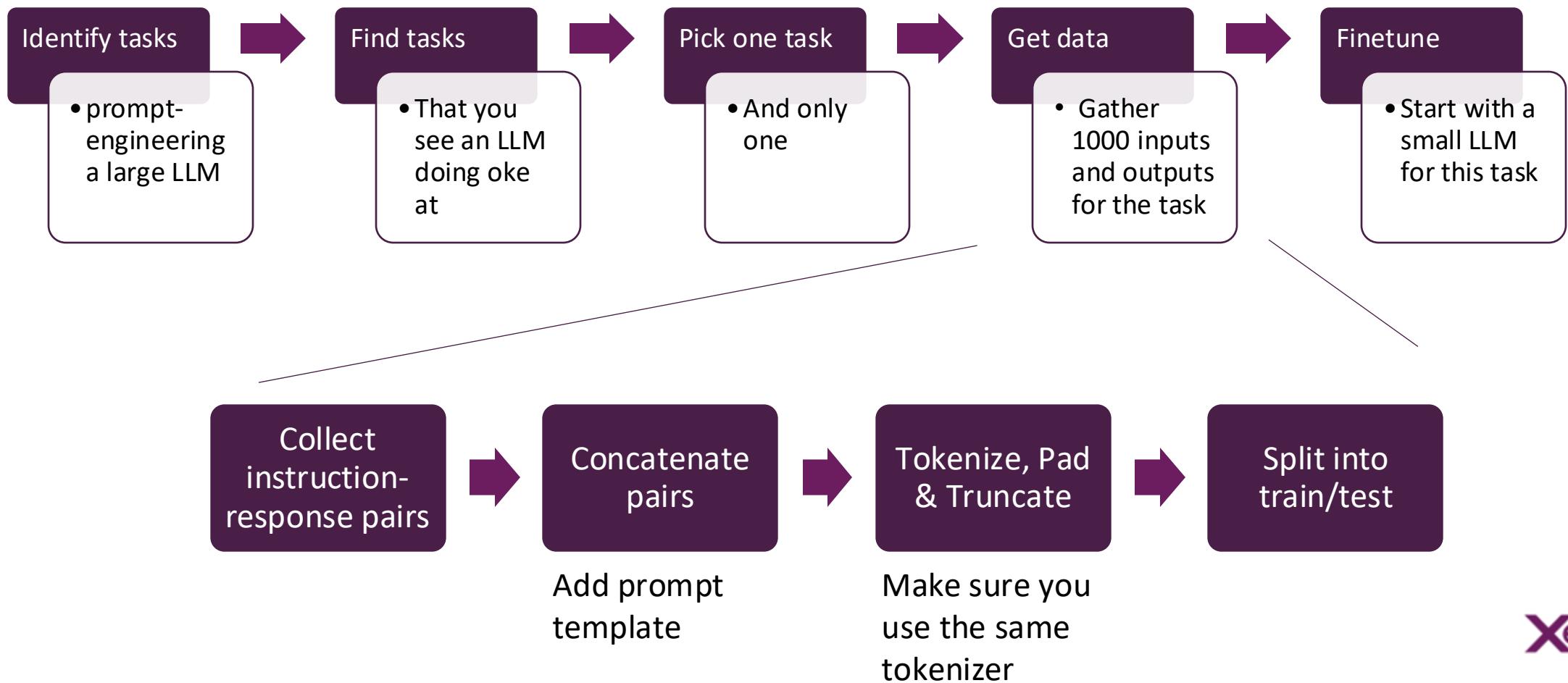
# Connections – What do you already know?

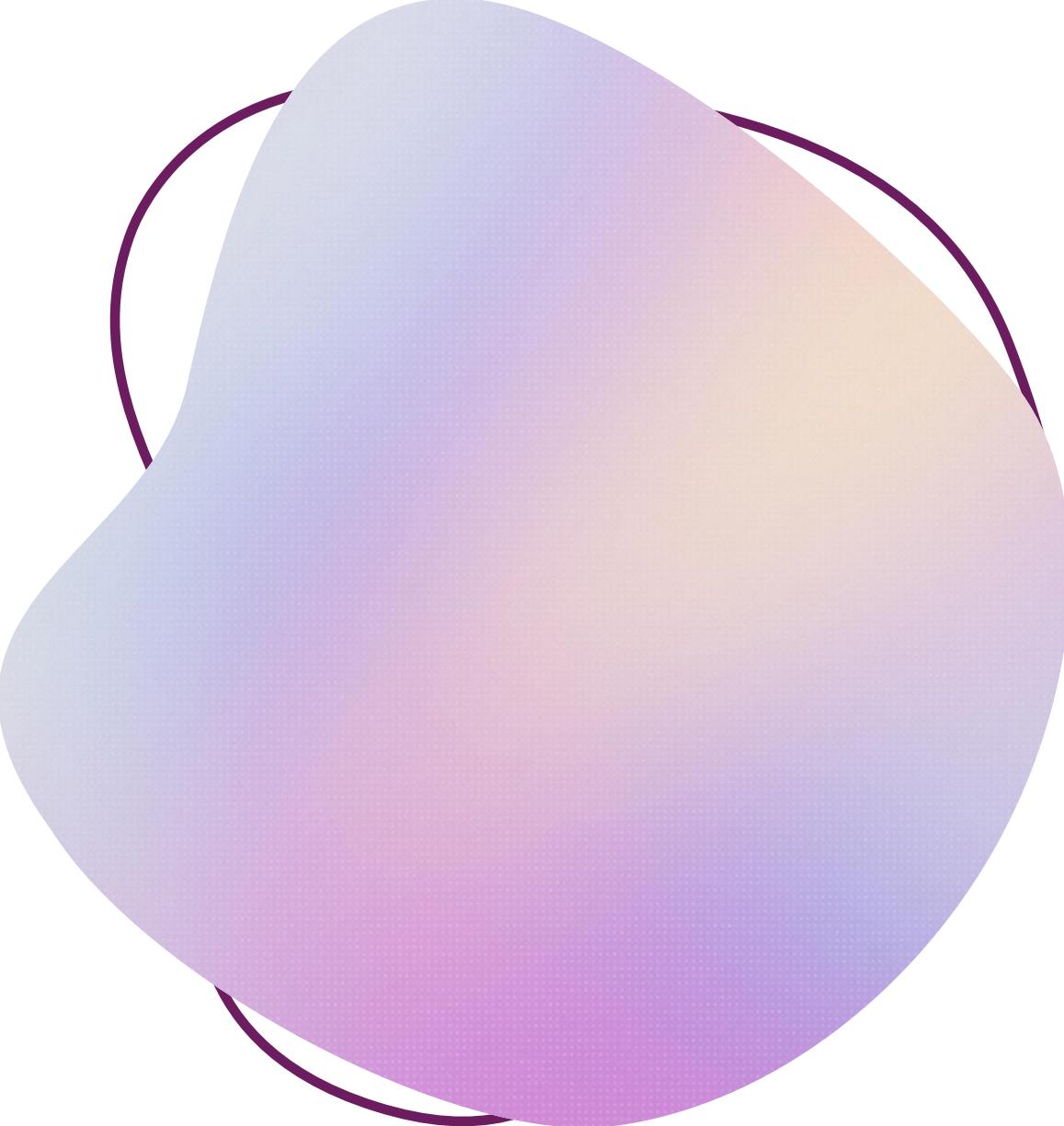
- **Who has fine-tuned an LLM?**
  - QLoRA
  - RLHF
- **Who is familiar with training a neural network?**
  - PyTorch, JAX, TensorFlow, Keras, FastAI
  - Epoch, batch, learning rate, back-prop?
- **What questions do you already have?**

# Connections – Why do you need to learn this?

- For most use-cases: fine-tuning is NOT needed and should be considered a last resort.
- For example, domain knowledge can be sent to the model using RAG.
- However, in certain cases, models might require additional training or specialized datasets to handle these domains effectively.

# Process for Finetuning

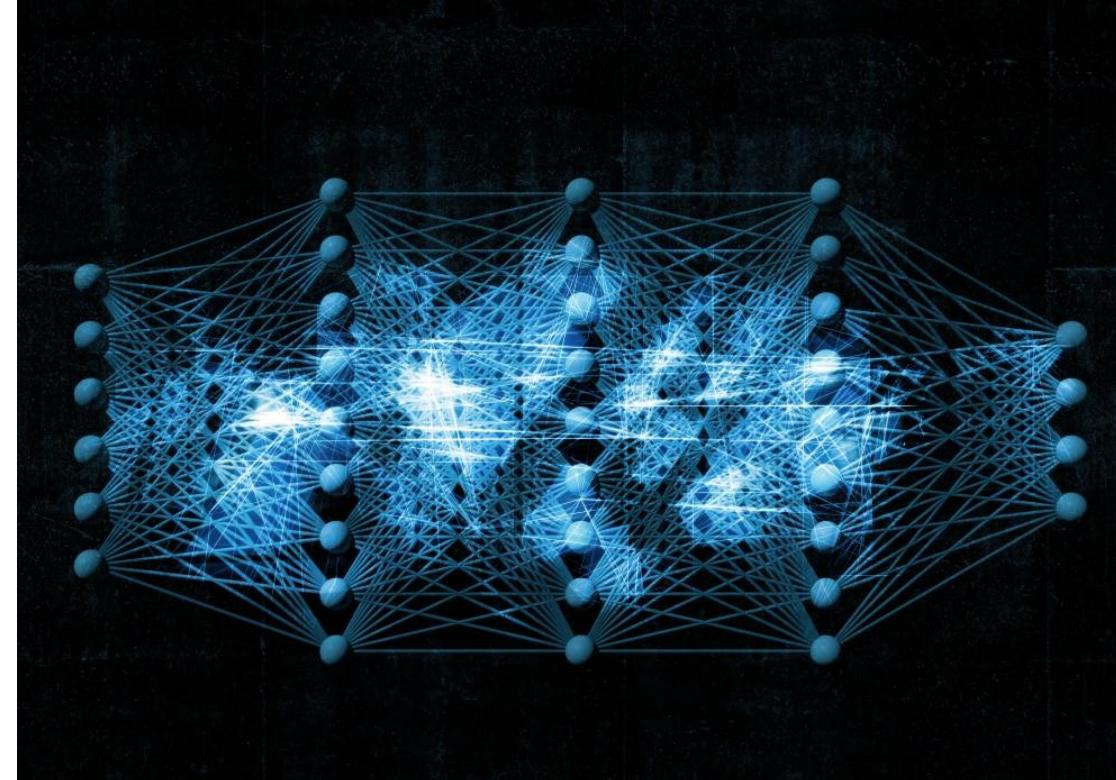




# Parameter Efficient Fine- tuning (PEFT)

# Full parameter fine-tuning

Full-parameter fine-tuning involves updating all the parameters of a pre-trained model on a new task-specific dataset.



# Full parameter fine-tuning Issues

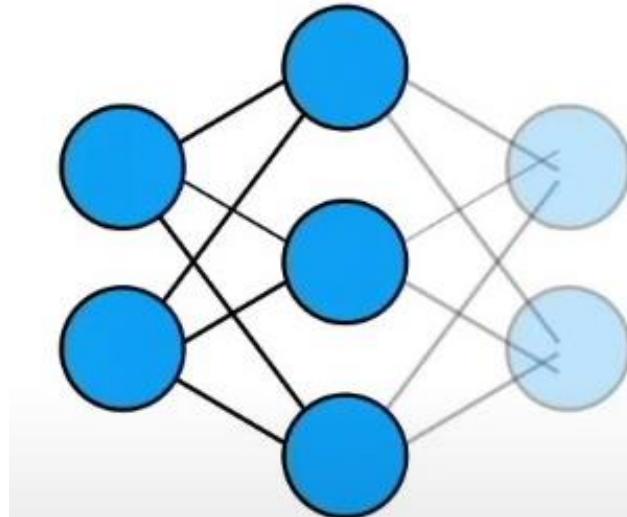
**Resource-Intensive:** Requires extensive computational power (GPU/TPU) and large storage due to model size

**Slow Training:** Retraining the entire model takes significant time.

**Catastrophic Forgetting:** Fine-tuning can cause the model to forget previously learned information.

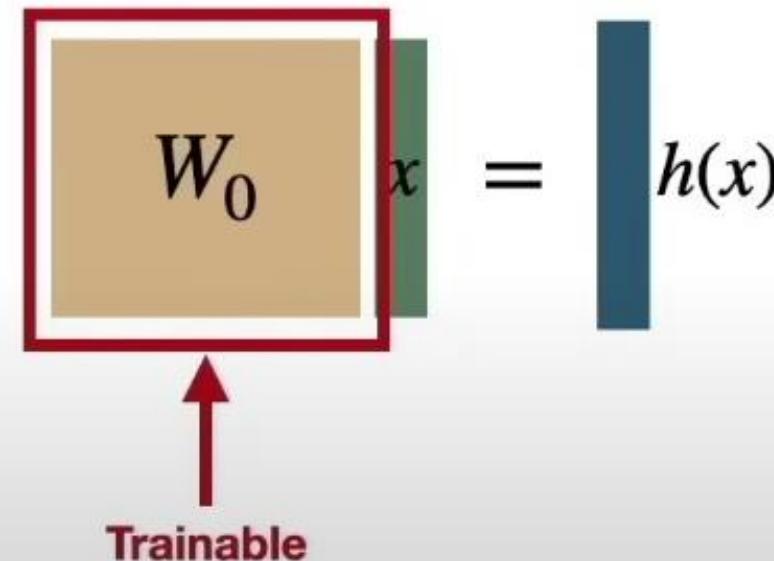
Model Name	$n_{\text{params}}$
GPT-3 Small	125M
GPT-3 Medium	350M
GPT-3 Large	760M
GPT-3 XL	1.3B
GPT-3 2.7B	2.7B
GPT-3 6.7B	6.7B
GPT-3 13B	13.0B
GPT-3 175B or "GPT-3"	175.0B

# Full Parameter Fine-tuning



$$x \rightarrow h(x)$$

$$h(x) = W_0 x$$



A diagram illustrating the matrix multiplication  $W_0 x = h(x)$ . On the left, there is a large orange rectangle labeled  $W_0$ , which is multiplied by a smaller green rectangle labeled  $x$ . The result is a dark blue vertical bar labeled  $h(x)$ . A red arrow points upwards from the  $x$  term, with the word "Trainable" written below it in red.

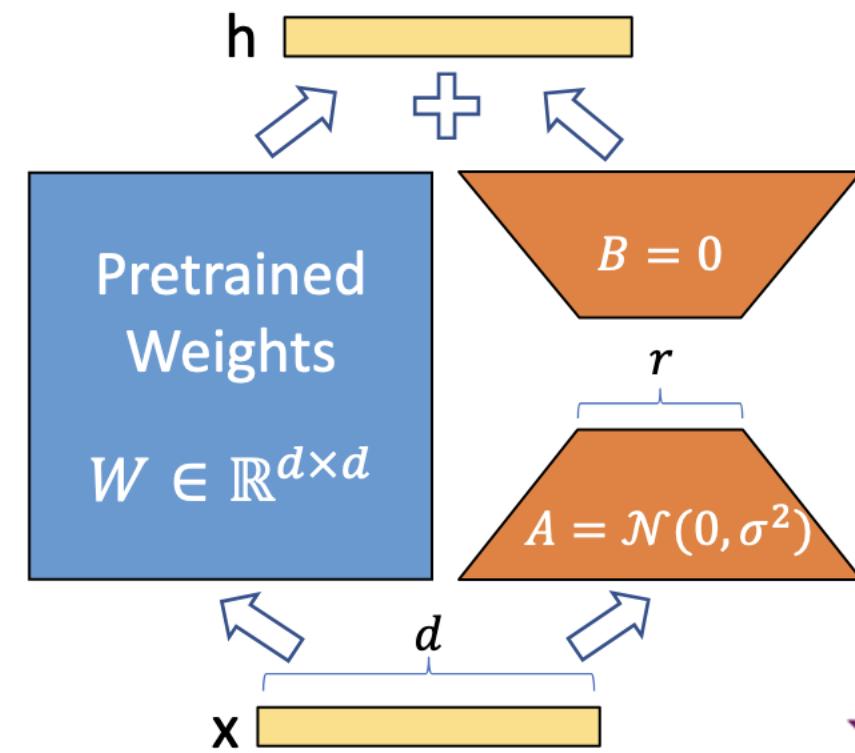
$$\begin{aligned} W_0 &\in R^{d \times k} \\ x &\in R^{k \times 1} \\ h(x) &\in R^{d \times 1} \end{aligned}$$

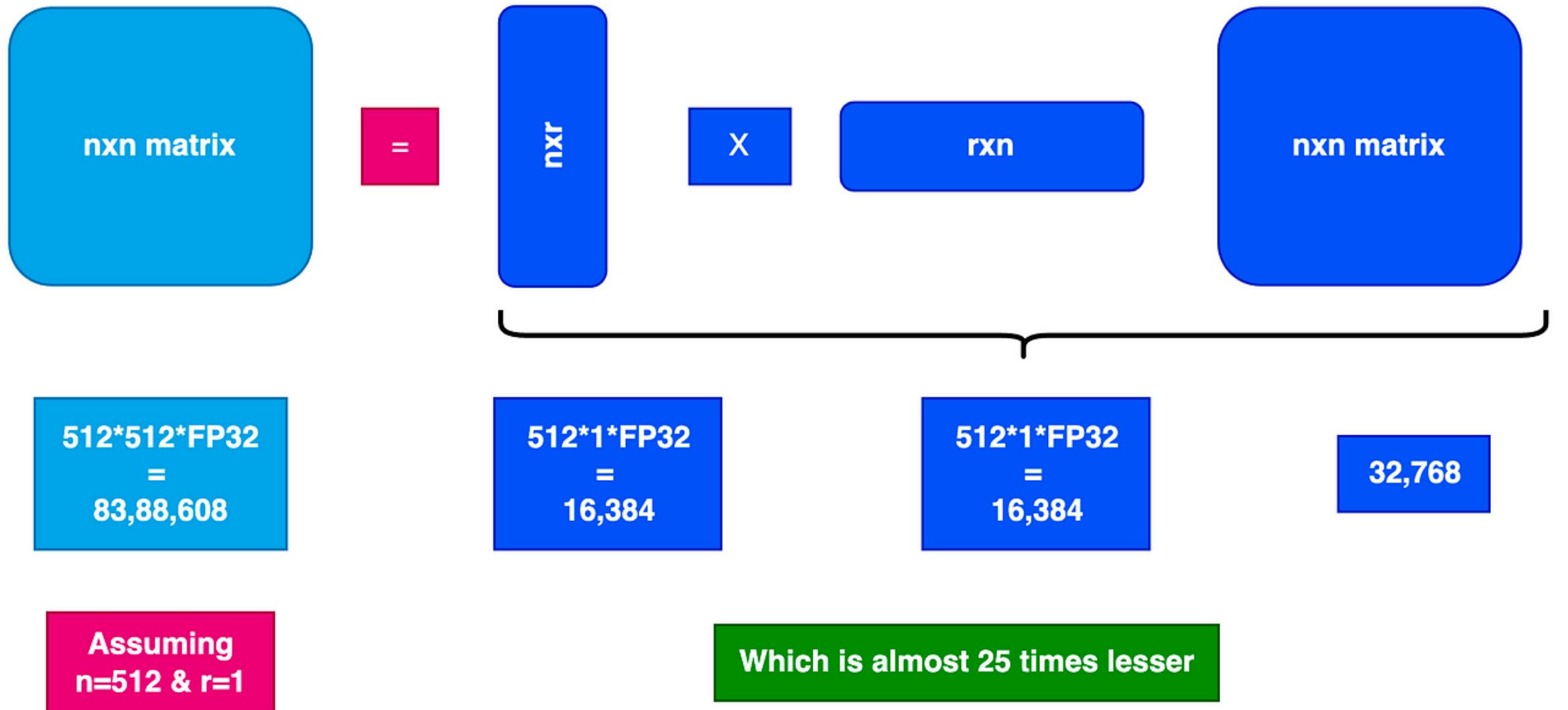
$$\begin{aligned} d &= 1,000 \\ k &= 1,000 \end{aligned} \implies d \times k = 1,000,000 \text{ trainable parameters}$$

# Solution: Low Rank Adaptation (LoRA)

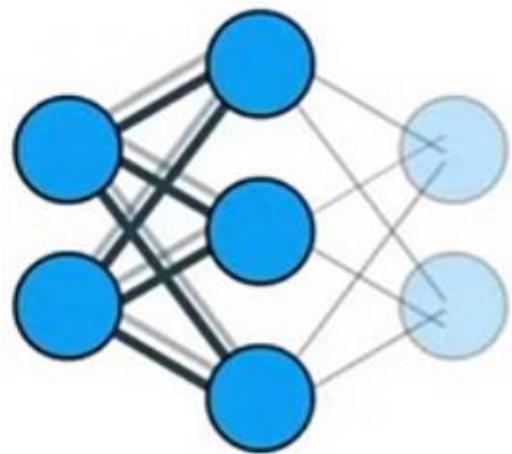
One Parameter Efficient approach to fine-tuning is LoRA

The key to LoRA matrix decomposition

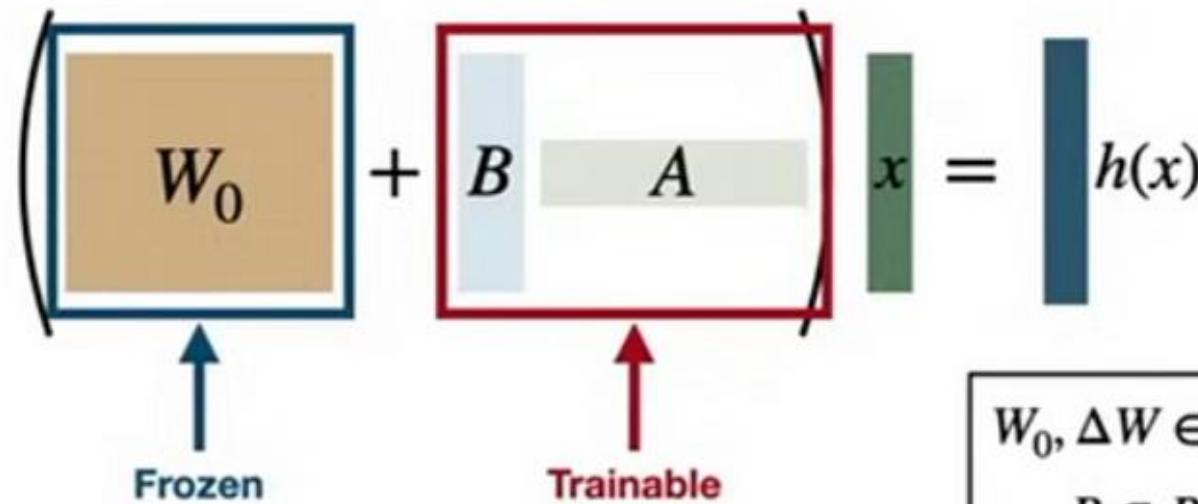




# Low Rank Adaptation (LoRA)



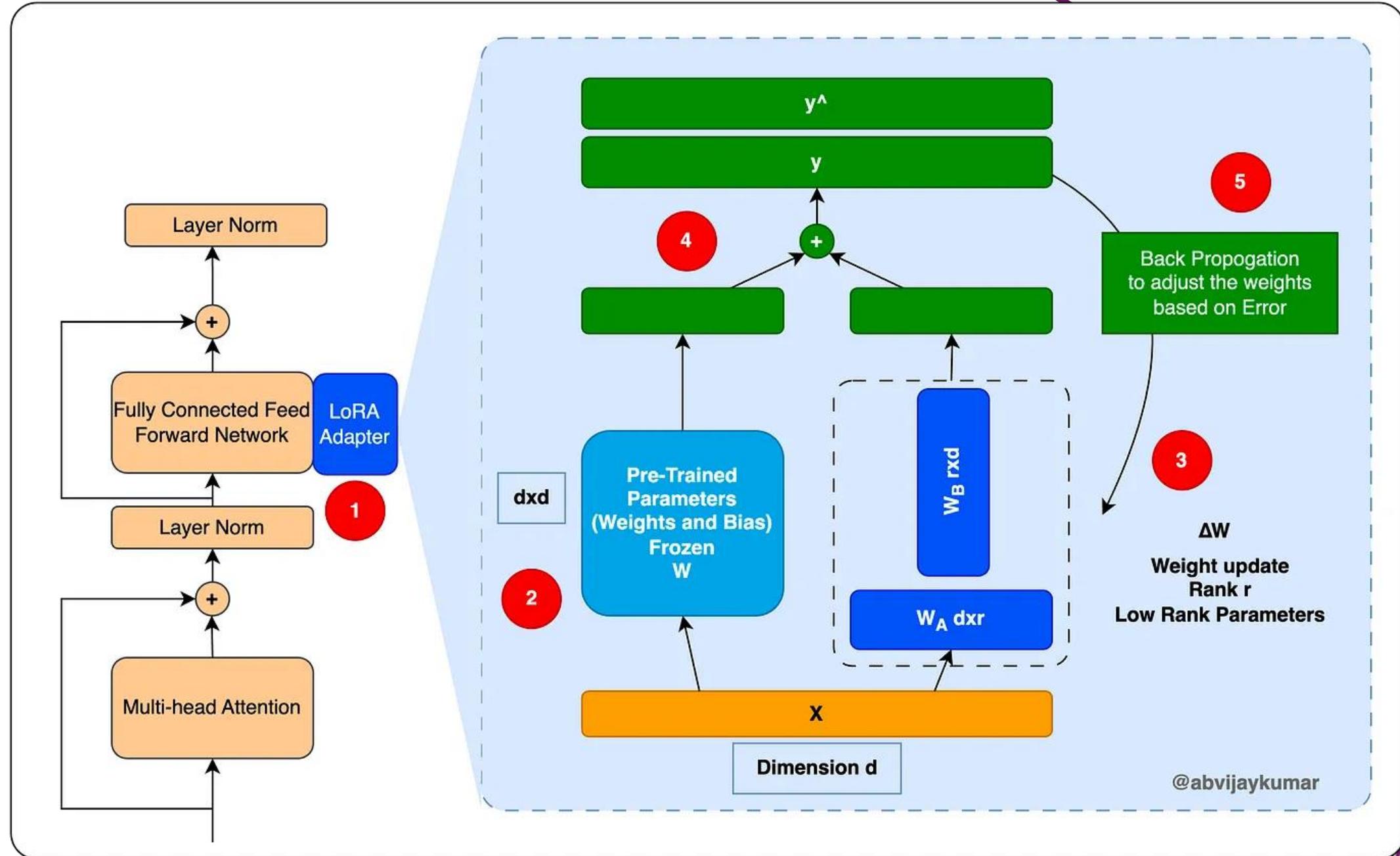
$$h(x) = W_0x + \Delta Wx \quad \Delta W = BA$$
$$= W_0x + BAx$$



$$d = 1,000$$
$$k = 1,000 \implies r = 2$$
$$(d \times r) + (r \times k) = 4,000$$

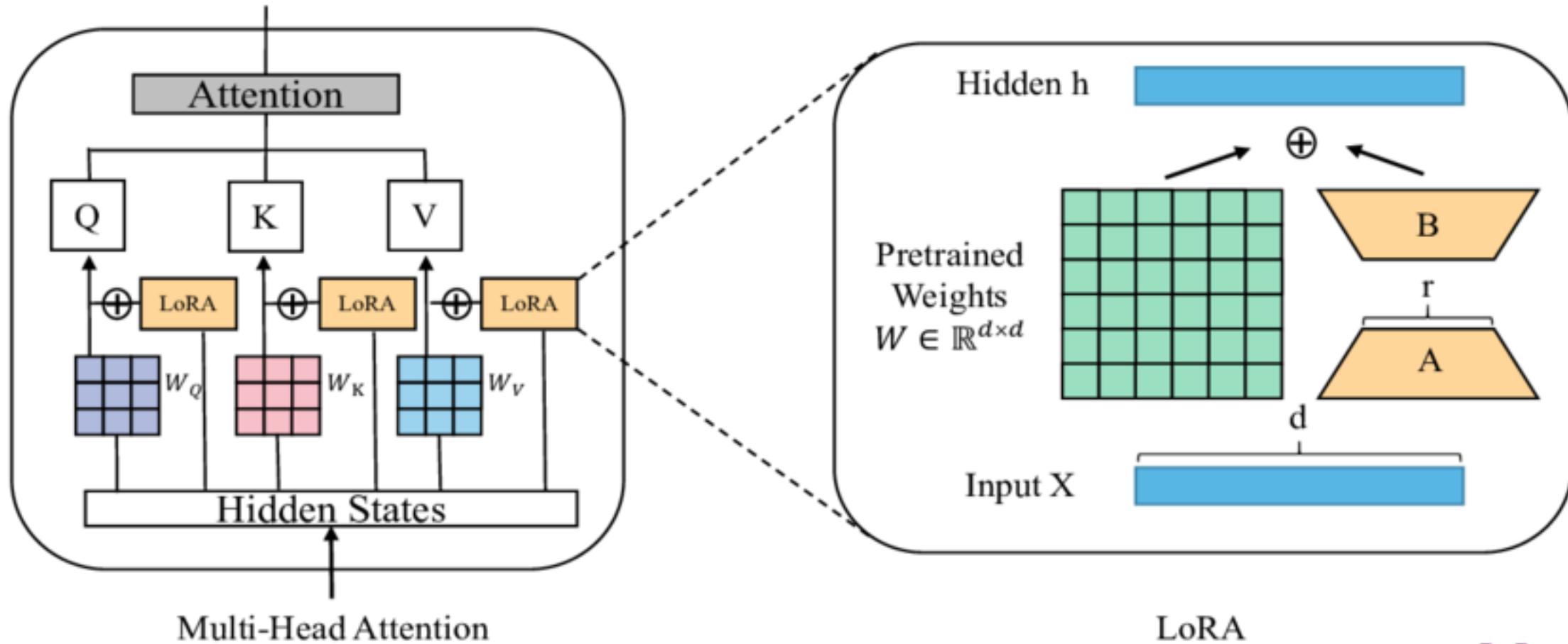
**trainable parameters**

$$W_0, \Delta W \in R^{d \times k}$$
$$B \in R^{d \times r}$$
$$A \in R^{r \times k}$$
$$h(x) \in R^{d \times 1}$$



@abvijaykumar

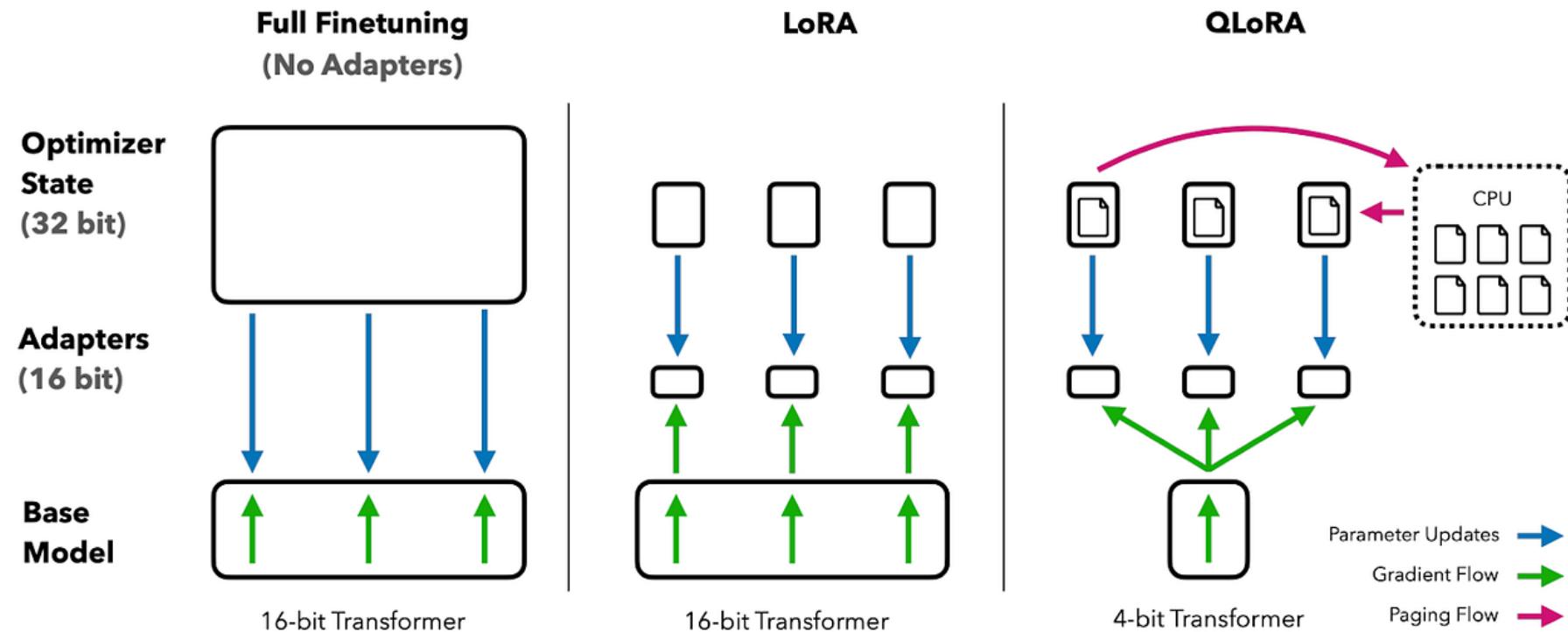
# LoRA can be applied to different attention components





# Quantization

# QLoRA – Quantized Low Rank Adaption Method



**Figure 1:** Different finetuning methods and their memory requirements. QLoRA improves over LoRA by quantizing the transformer model to 4-bit precision and using paged optimizers to handle memory spikes.

# There are 3 Key optimizations that QLoRA brings on top of LoRA

## 1. 4-bit NF4 Quantization

4-bit NormalFloat4 is an optimized data type that can be used to store the weights of the original model, which brings down the memory footprint considerably.

It is a 3-step process:

## 2. Normalization:

As part of the normalization step, the weights are adjusted to a zero mean, and a constant unit variance.

There are 3 Key optimizations that QLoRA brings on top of LoRA

## 1. 4-bit NF4 Quantization

4-bit NormalFloat4 is an optimized data type that can be used to store the weights of the original model, which brings down the memory footprint considerably.

It is a 3-step process:

2. **Quantization:** A 4-bit data type can only store 16 numbers. So, after normalization, instead of storing the weights, the nearest position is stored.

For example, let's say we have a FP32 weight, with a value of 0.2121. a 4-bit split between -1 to 1 will be the following number positions. 0.2121 is closest to 0.1997, which is the 10th position. Instead of saving the FP32 of 0.2121, we store 10.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
-1	-0.8667	-0.7334	-0.6001	-0.4668	-0.3335	-0.2002	-0.0669	0.0664	0.1997	0.333	0.4663	0.5996	0.7329	0.8662	1

**There are 3 Key optimizations that QLoRA brings on top of LoRA**

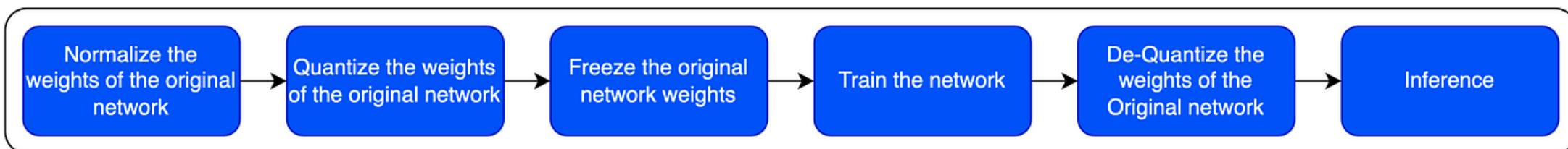
## 1. 4-bit NF4 Quantization

4-bit NormalFloat4 is an optimized data type that can be used to store the weights of the original model, which brings down the memory footprint considerably.

It is a 3-step process:

**3. Dequantization:** Once all the training is done, the original weights will be de-quantized, i.e. reverse the process.

Note, the 4-bit NormalFloat quantization is applied only to the original weights of the model. The LoRA adapter weights will be left as FP32, as the training will happen on these weights.



There are 3 Key optimizations that QLoRA brings on top of LoRA

## 2. Double Quantization

In 4-bit NF4 Quantization, a “quantization constant” is used to scale the weights

$$Q(w) = \text{round}\left(\frac{w - z}{s}\right)$$

$w$  = Original high-precision weight.

$z$  = An offset used to align the quantized range with the actual weight distribution.

$s$  = Quantization constant (scaling factor), dependent on the max weight value.

$Q(w)$  = Quantized weight, in a lower precision format (i.e., 4-bit integer)

There are 3 Key optimizations that QLoRA brings on top of LoRA

## 2. Double Quantization

**Double quantization**, further reduces the memory footprint, by quantizing the quantization constants.

Let's imagine we have grouped 64 parameters/weights into a block to quantize, and each quantization constant takes 32 bits, as it is FP32. It adds a 0.5 bit per parameter on average, which means we are talking of at least 500,000 bits for a typical 1Mil parameter model.

With Double quantization, we apply quantization on these quantization constants, which will further optimize our memory usage.

There are 3 Key optimizations that QLoRA brings on top of LoRA

### 3. Unified Memory Paging

Coupled with the above techniques, QLoRA also utilizes the nVidia unified memory feature, which allows GPU->CPU seamless page transfers, when GPU runs out of memory, thus managing the sudden memory spikes in GPU, and helping memory overflow/overrun issues.



Tim Dettmers provides a nice explanation: <https://github.com/bitsandbytes-foundation/bitsandbytes/issues/962>

# Important Concepts: QLoRA

- **LoRA: Low rank adapters**

- Reduce the trainable parameters
- Llama2 13b
  - trainable params: 36,372,480 || all params: 13,052,236,800  
trainable: 0.28%
- Mistral 7b
  - trainable params: 28,311,552 || all params: 7,270,043,648  
trainable: 0.39%

- **Quantization**

- Bfloat-16 / NF4
- Quantization error
  - Recovers with fine-tuning
  - Re-occurs with merging weights
- Rank
  - Alpha = 2\*rank

```
lora_config:  
# See https://arxiv.org/pdf/2110.04366.pdf, p.8  
target_modules:  
- down_proj  
- up_proj  
- gate_proj  
# - k_proj  
# - v_proj  
# - q_proj  
# - o_proj  
# - lm_head  
  
- lora_target_modules:  
- q_proj  
- v_proj  
- k_proj  
- o_proj  
- gate_proj  
- down_proj  
- up_proj  
- c_attn  
- c_proj
```

- Best practices

- Efficiency Methods

QLoRA, FlashAttention and  
paged\_adamw\_32bit for efficient training

Hardware specific optimization  
(TensorRT, TensorRT-LLM)

### Performance

- Ensemble
- Multi LoRA

<https://xebia.com/blog/llm-fine-tuning-challenge/>  
<https://xebia.com/blog/winning-recipe-for-finetuning-llm/>

# Make use of all available optimizations

Method/tool	Improves training speed	Optimizes memory utilization
<u>Batch size choice</u>	Yes	Yes
<u>Gradient accumulation</u>	No	Yes
<u>Gradient checkpointing</u>	No	Yes
<u>Mixed precision training</u>	Yes	(No)
<u>Optimizer choice</u>	Yes	Yes
<u>Data preloading</u>	Yes	No
<u>DeepSpeed Zero</u>	No	Yes
<u>torch.compile</u>	Yes	No

[https://huggingface.co/docs/transformers/perf\\_train\\_gpu\\_one#batch-size-choice](https://huggingface.co/docs/transformers/perf_train_gpu_one#batch-size-choice)

# Inference and efficiency

- **Experiment with efficient use of GPU VRAM**
  - Rule of thumb: model with Xbn params will take up ~4X VRAM in float32 precision.
  - Using bfloat16, it's ~2X
  - Using bfloat16 + 8-bit quantization, its ~1.1X
  - But merging the QLoRa model weights in at inference time blows things up to ~1.8X compared to normal 8-bit + bfloat16
  - Best thing: merge and dump to disk, then load at inference time. This, too, requires a lot of VRAM.
- **Be mindful of too much quantization**
  - In practice, 4-bit inference saves you ~40-50% VRAM compared to 8-bit
  - But performance is also worse . . .

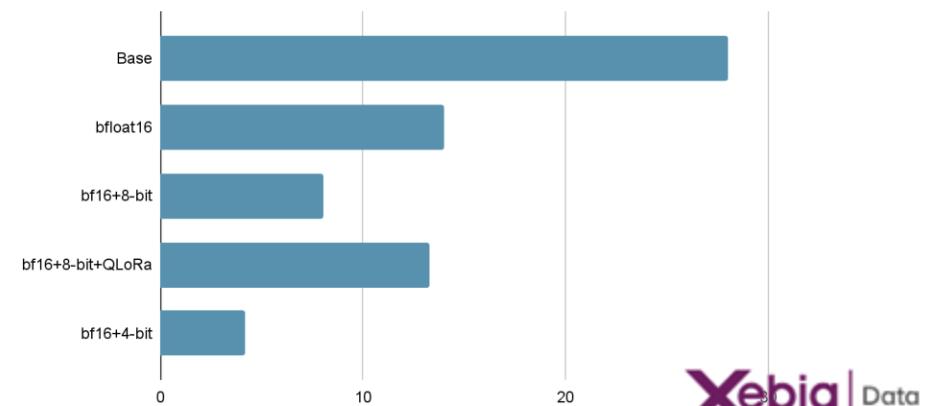
**4-bit**

Name	Value
accuracy_MMLU_EM ↗	0.497
accuracy_TruthfulQA_EM ↗	0.333
calibration_MMLU_ECE_10bin ↗	0.688
fairness_MMLU_EM_Fairness ↗	0.497
fairness_TruthfulQA_EM_Fairness ↗	0.333

**8-bit**

Name	Value
accuracy_MMLU_EM ↗	0.655
accuracy_TruthfulQA_EM ↗	0.889
calibration_MMLU_ECE_10bin ↗	0.278
fairness_MMLU_EM_Fairness ↗	0.623
fairness_TruthfulQA_EM_Fairness ↗	0.889

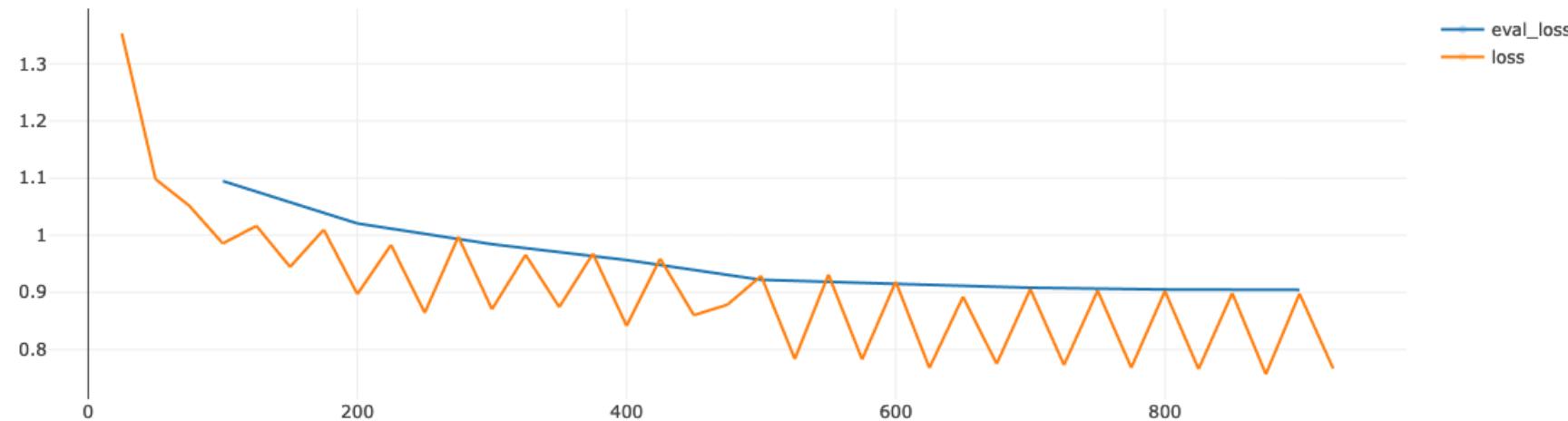
VRAM used for 7B model



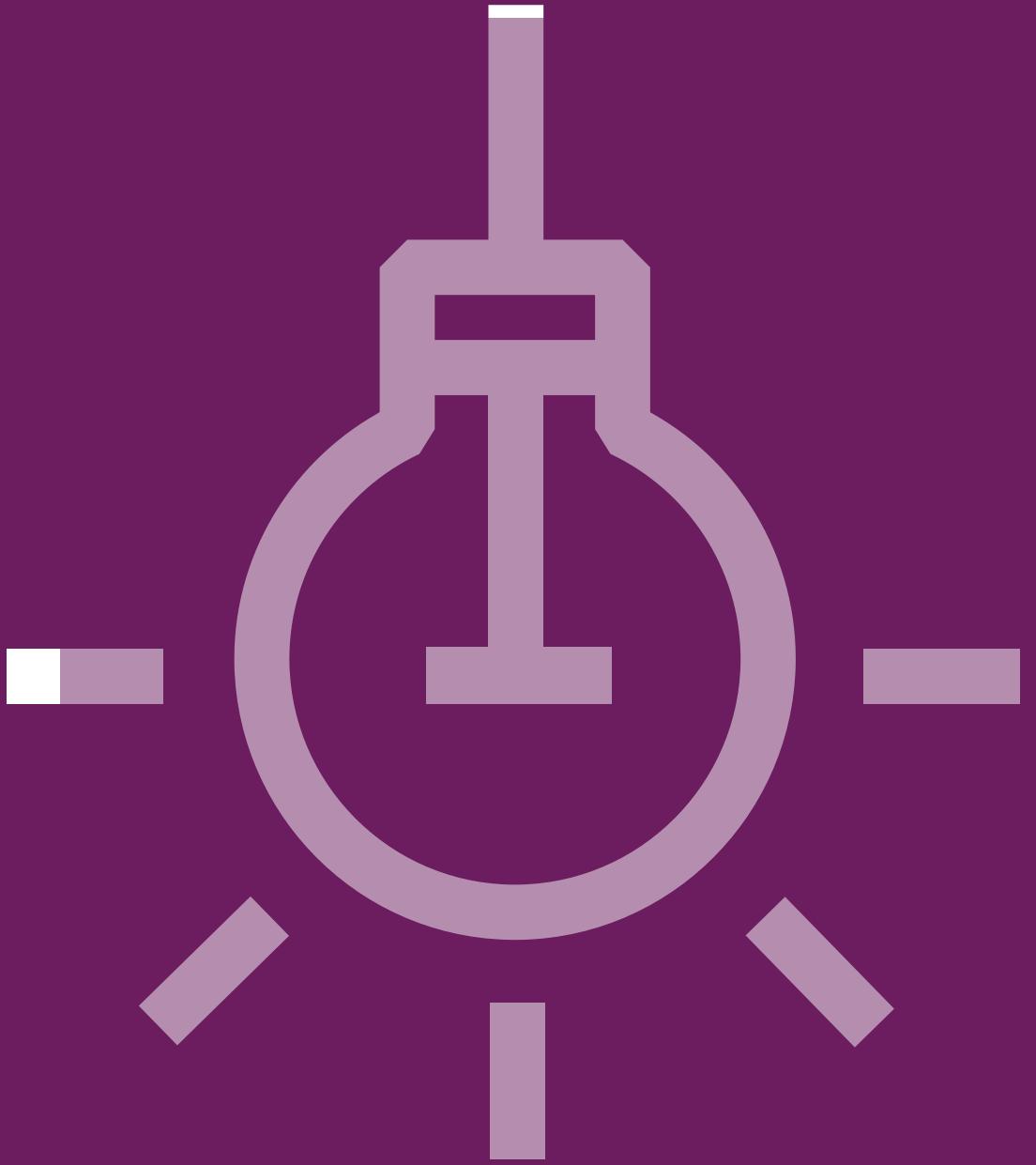
# Don'ts

- **Grouping by length**

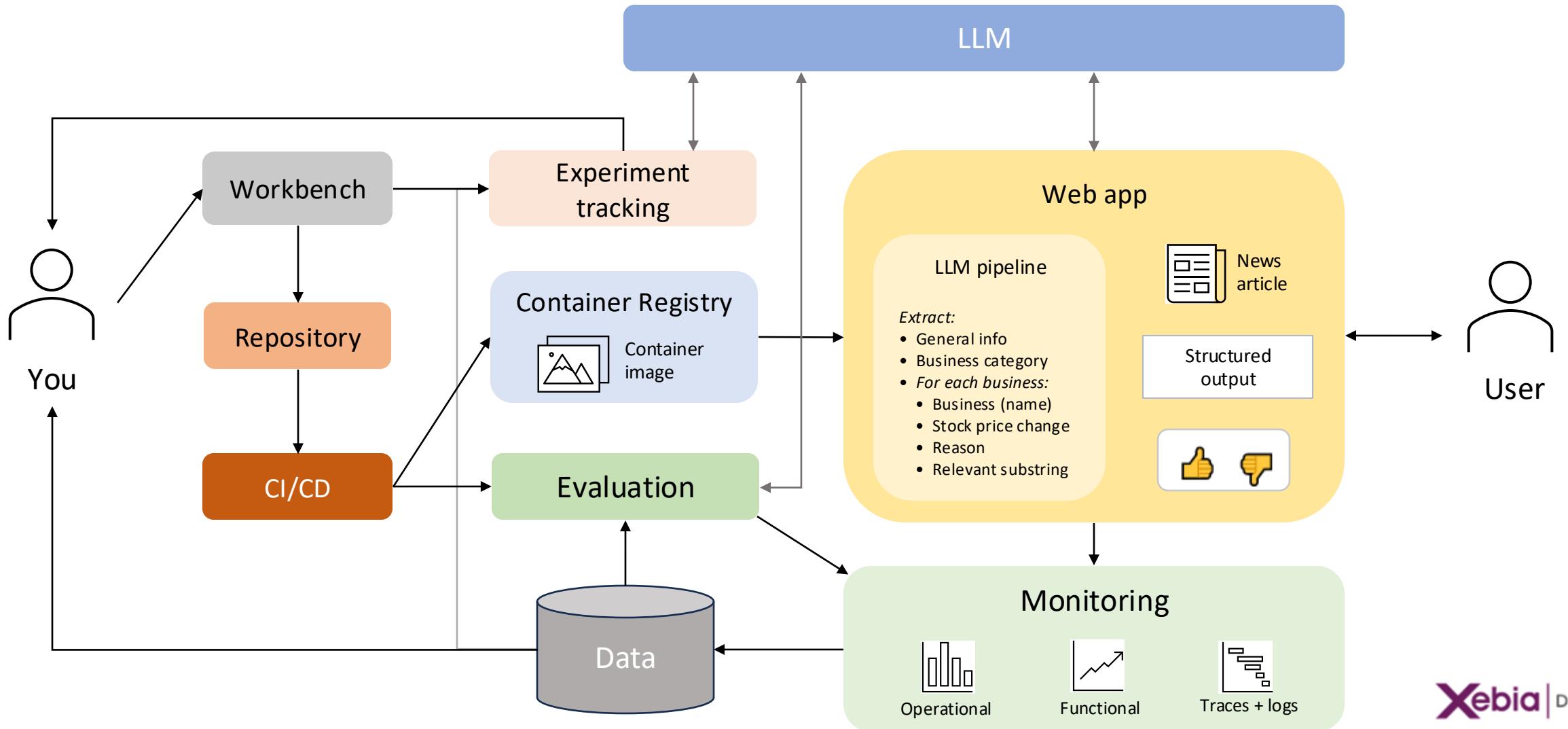
- GPU FLOPS utilization is not optimal when training on sequences of different lengths. Text sequences can be padded, but that leads to a lot of wasted compute. Instead, you can group sequences by length and pad each group separately. This way, you can reduce the amount of padding and increase the utilization of the GPU. However, this results in unstable training.
- The figure below shows one of the earlier training runs. The training loss is very unstable. We expect this to come from the fact that we finetune on a wide range of tasks. Grouping by length in our case also means roughly grouping by task. At the end of a batch, the loss is evaluated for the examples in the batch. This means that if the model performance varies per task, the training loss is going to fluctuate per task.
- This in itself does not have to be harmful. However, it shows that this process potentially updates the model for a specific task per batch, which can make for unstable training. It might be mitigated by using larger batch accumulation (more observations per forward pass before an update step).



# Wrap-up



# Started from a notebook now we're here!



# Please submit feedback and help us improve



<https://xebia.ai/llmops-survey>



# Final notes

- Don't forget to push everything to your branch
- Clone/Fork the Repository
- If you need your code space contingent for sth else, don't forget to delete your codespace (but all un-pushed / un-committed changes will be lost!)
- I will leave the Azure Resource Group running for another week for you to play around in (pls don't abuse our compute)
- You can replicate the Azure Resource Group by creating a new one on Rabo compute afterwards (follow the `trainer_docs` in the repository to set this up)

# Thanks!