

Why we put GPUs in the database

Postgres Conf 2024

A black hole with a glowing accretion disk against a starry background. The text is centered over the black hole.

**AI isn't the future,
it's here.**

IT'S BEEN A LONG TIME COMING



AI is table stakes

Successful companies will incorporate AI into their core UX, if they haven't already

Personalization

Real-time personalization, knowledge based chat

Search

Proprietary data made findable and effective

Anomaly detection

Rapid detection, response command and control

Fraud detection

Fraud reduced, countermeasures activated instantly





MSFT
+46.38%



META
+137.47%



BRK.B
+35.23%



JPM
+51.95%



AMZN
+77.53%



WMT
+20.04%



HD
+20.97%



+4.04%



-4.10%



+34.51%



+53.79%



+80.80%



+31.14%



+51.54%



COST
+42.13%



Target



Target



Target



Target



Target



Target



Target



+30.03%



IBM



Intuit



Now



+16.82%



Alphabet



Alphabet



Alphabet



-6.32%



Alphabet



Alphabet



Alphabet



Goldman Sachs



Citi



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



Cherubini



GOOGL
+48.49%



+111.96%



+86.29%



Qualcomm



Intel



GE



Texas Instruments



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



LLY
+120.41%



JNJ
-1.73%



MRK
+19.50%



+10.48%



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



CAT



Caterpillar



Caterpillar



Caterpillar



Caterpillar



Caterpillar



Caterpillar



Caterpillar



Caterpillar



Caterpillar



Caterpillar



Caterpillar



Caterpillar



UNH
-6.75%



UnitedHealth



UnitedHealth



UnitedHealth



UnitedHealth



UnitedHealth



UnitedHealth



UnitedHealth



AAPL
+1.83%



+111.96%



+86.29%



Qualcomm



Intel



GE



Texas Instruments



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



Micron



LLY
+120.41%



JNJ
-1.73%



MRK
+19.50%



+10.48%



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



AbbVie



CAT



Caterpillar

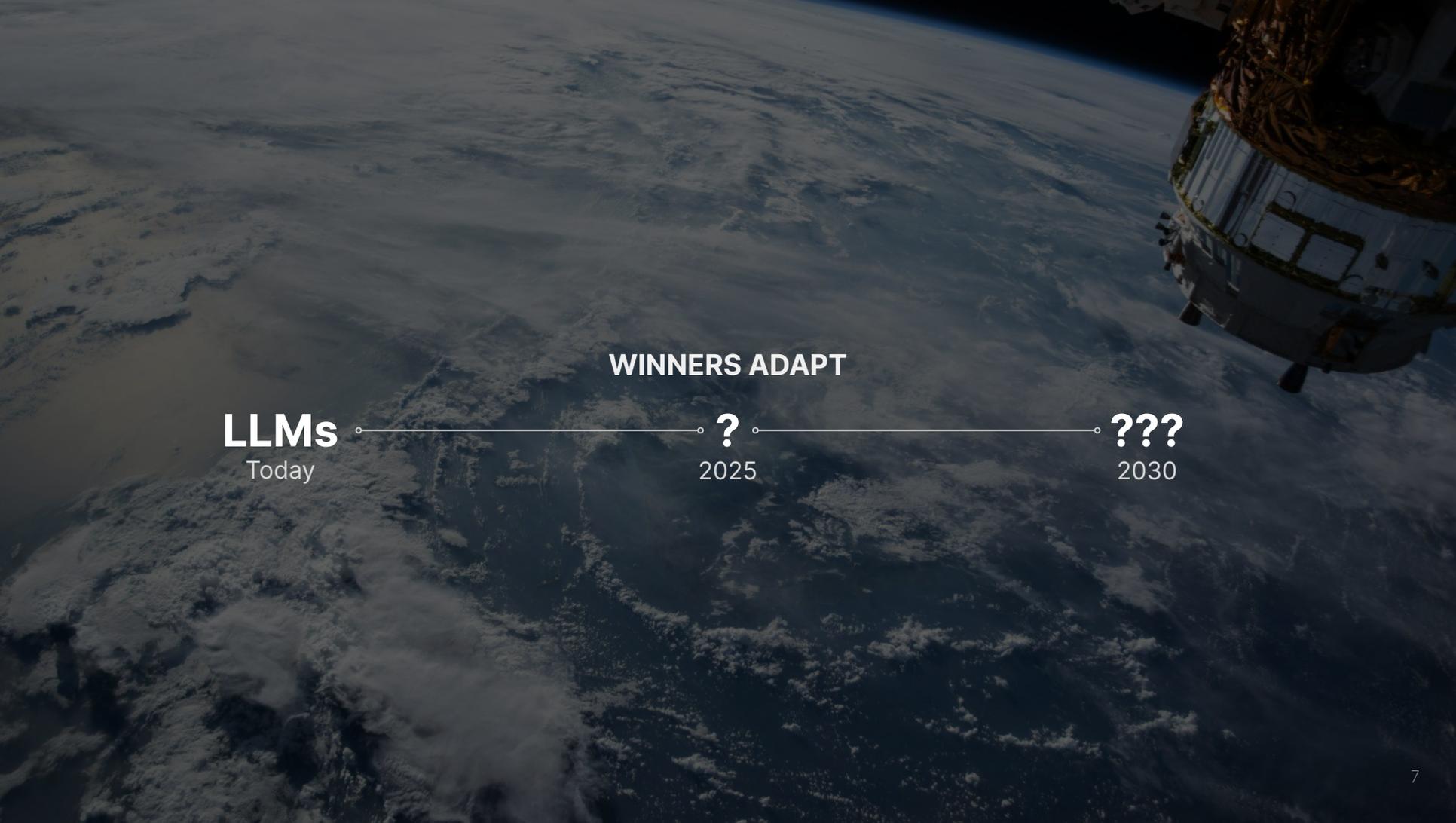


Caterpillar



Caterpillar





WINNERS ADAPT

LLMs

Today

○ ————— ○ **?** ○ ————— ○

2025

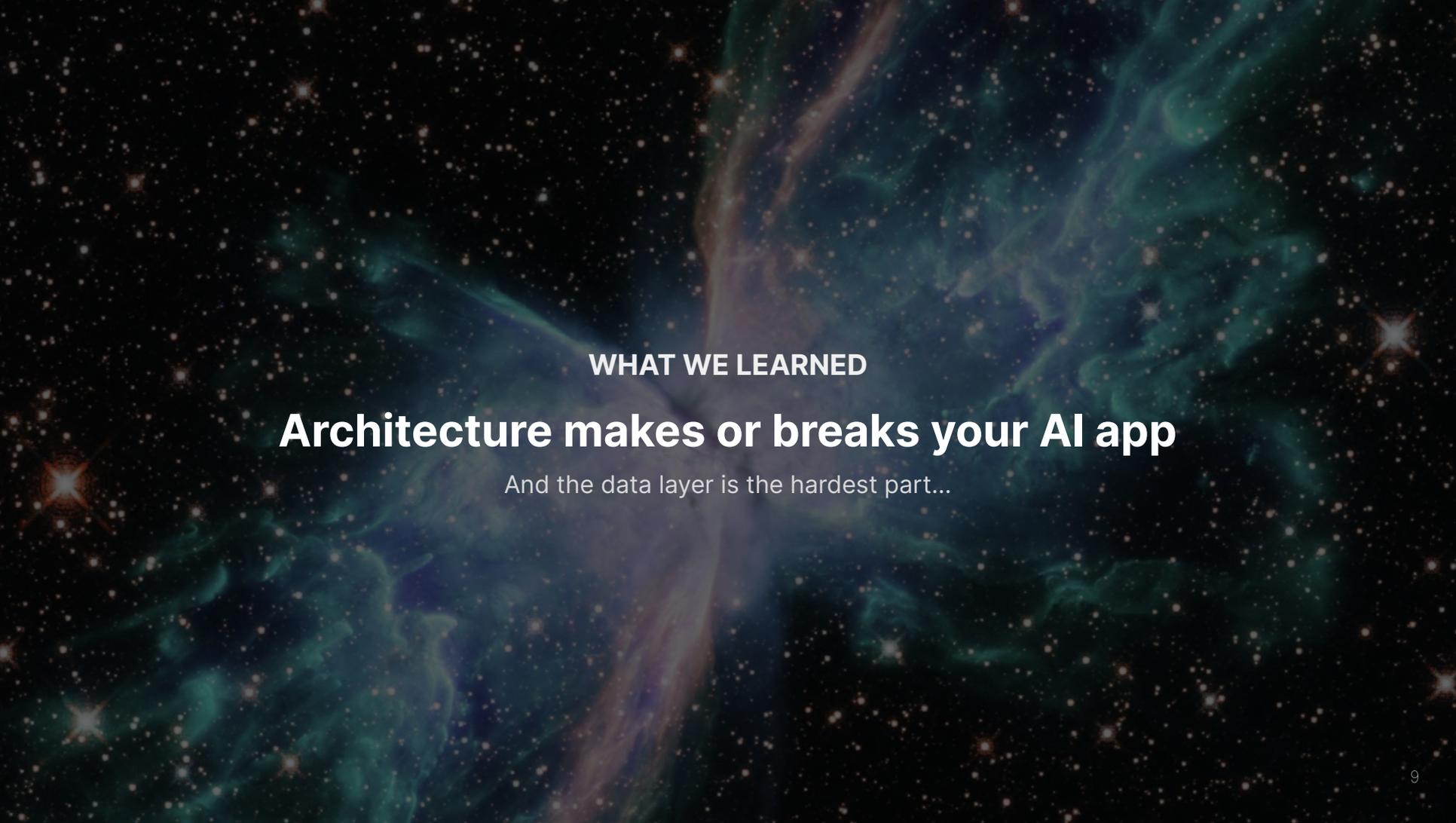
???

2030



Our story

Founded by the engineers that built (then re-built) and scaled Instacart's machine learning platform

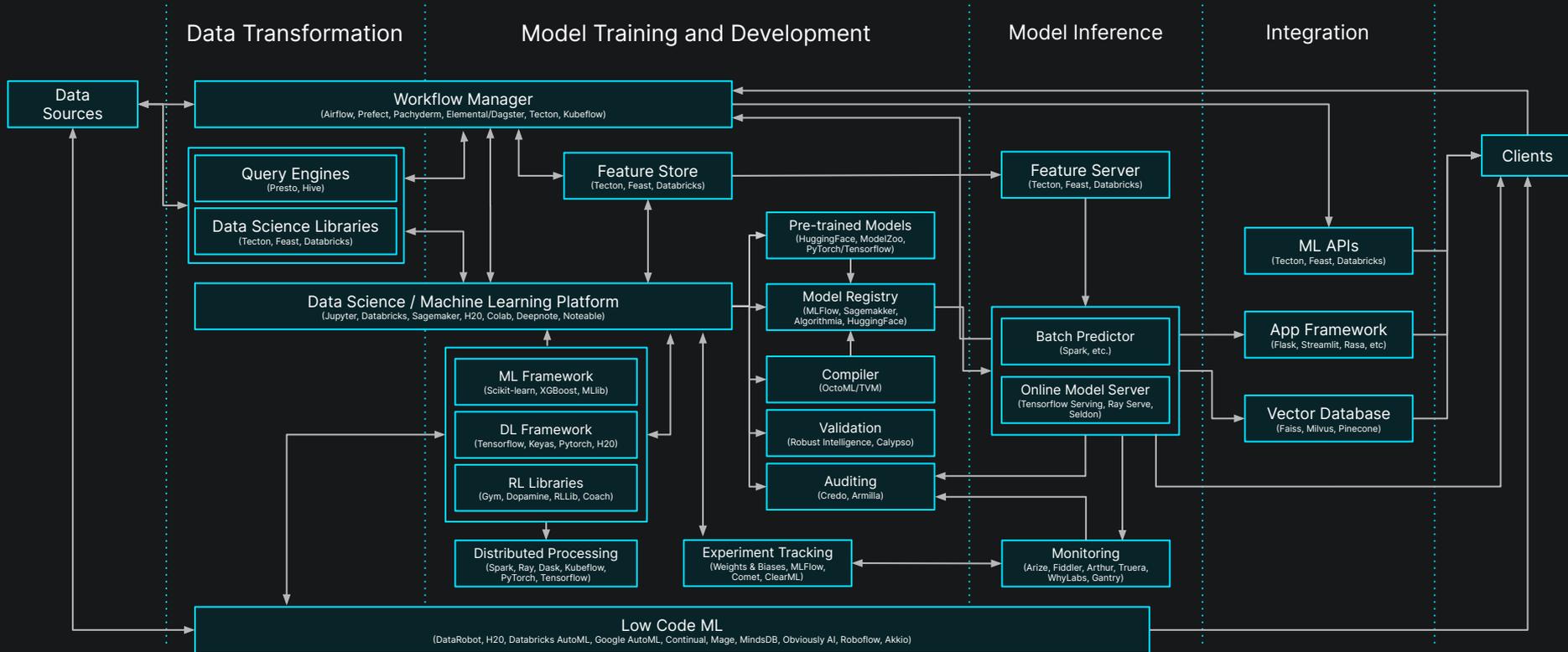


WHAT WE LEARNED

Architecture makes or breaks your AI app

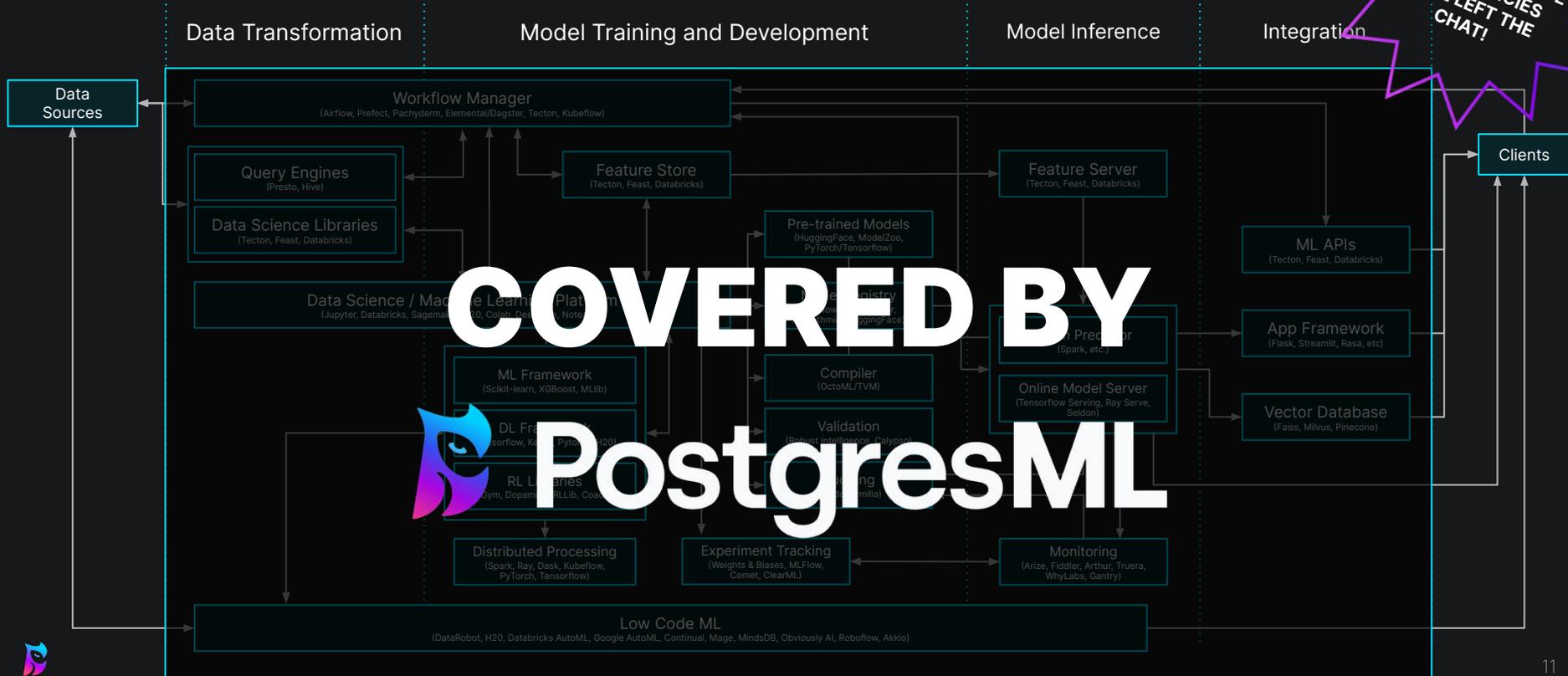
And the data layer is the hardest part...

Machine learning infrastructure



Machine learning infrastructure

DATA AND SERVICE
DEPENDENCIES
HAVE LEFT THE
CHAT!



PostgresML

Serverless | Dedicated | VPC

Use Cases

- Fraud Detection
- Chatbots
- Forecasting
- Search
- & More

Capabilities

- Unified RAG
- NLP
- Regression
- Classification
- Clustering

Serving & Storage

- Embeddings
- Inference
- Training
- Vector Database
- Feature Store
- Model Store

Postgres Extensions

- pgml
- pgvector
- pg_partman

Integrated Libraries

- Huggingface
- LLaMA
- Mistral
- Torch
- Tensorflow
- Flax
- scikit-learn
- XGBoost
- LightGBM
- CatBoost

DATA SOURCES

- SQL
- NoSQL
- Data Warehouse
- Data lake
- Streaming Data
- Connectors

INTERFACES

- Python SDK
- JavaScript SDK
- Apps
- SQL Clients
- Notebooks

SCALE
PgCat

HARDWARE



OSS ECOSYSTEM



Move the algorithm to the data

4x faster

than 🤖 HuggingFace + 🌲 Pinecone
for a RAG chatbot

10x faster

than 🌀 OpenAI for embedding
generation

10-20x faster

than 🏠 MindsDB
(not really a DB, it's a microservice)

8-40x faster

than Python + Redis for an
dmk XGBoost XGBoost microservice

Don't take our word for it

"...with [@postgresml](#), we're not just streamlining processes but revolutionizing data handling and insights generation in one fell swoop."

Co-Founder @ MedPiper Technologies, Inc

TRUSTED BY ENGINEERS AT



Join us

- Attend the pgml RAG session to see SOTA design
- Contribute to our open-source projects, including pg-cat

We're hiring:

Email: Montana@postgresml.org

