



# Robust and Resilient Edge Cloud Inference for LLMs: Empowering Intelligent Edge Applications

Presenter: Brad Munday

# Meet your host

---



**Brad Munday**

Head of ML Engineering

Former tech consultant, quant, ML engineer, with extensive ML open-source contributions

# LLMs at the edge in the real world

## Qualcomm Works with Meta to Enable On-device AI Applications Using Llama 2

JUL 18, 2023 | SAN DIEGO

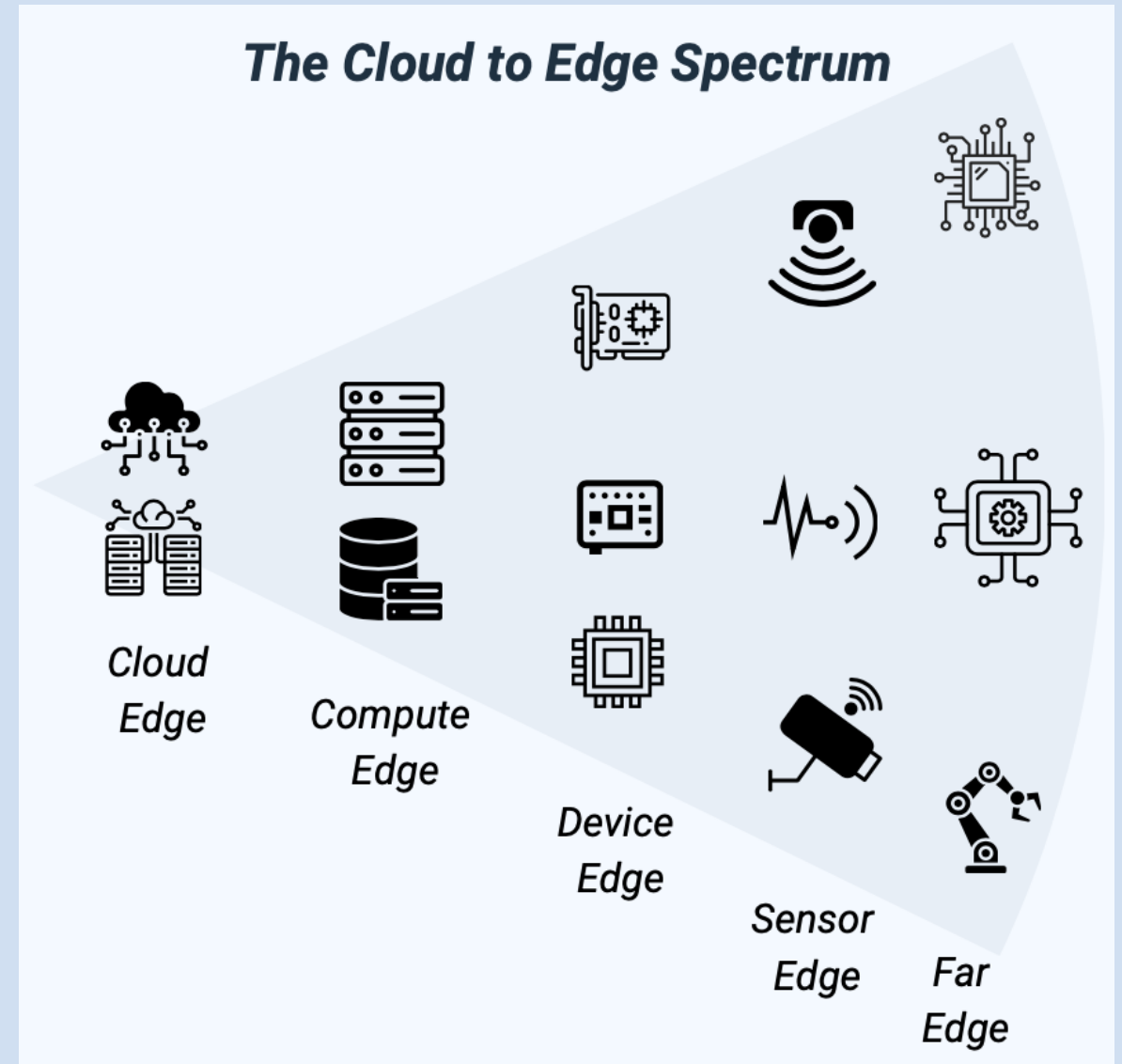
Qualcomm products mentioned within this press release are offered  
by Qualcomm Technologies, Inc. and/or its subsidiaries.

### Highlights:

- Qualcomm is scheduled to make available Llama 2-based AI implementations on flagship smartphones and PCs starting from 2024

# What constitutes “the edge?”

- **Cloud Edge** – offers computing capabilities that you would find in a cloud service provider, e.g., MEC.
- **Compute Edge** – functions as a localized, micro-data center that includes a limited range of resources and services you would find in the cloud, e.g., an edge line server racked or placed near or close to other devices or sensors.
- **Device Edge** – much smaller compute and processing capabilities, e.g., NVIDIA Jetson modules, Raspberry Pi, Intel NUC.
- **Sensor Edge** – comprises IoT sensors and devices that gather data, e.g., camera, and interact directly with the cloud, compute, or device.
- **Far Edge** – e.g., a microprocessor on board a robotic arm.



# Factors driving ML to the edge



**Bandwidth & Network  
Access**




**Real-time, Low  
Latency Results**



**Cost  
Considerations**

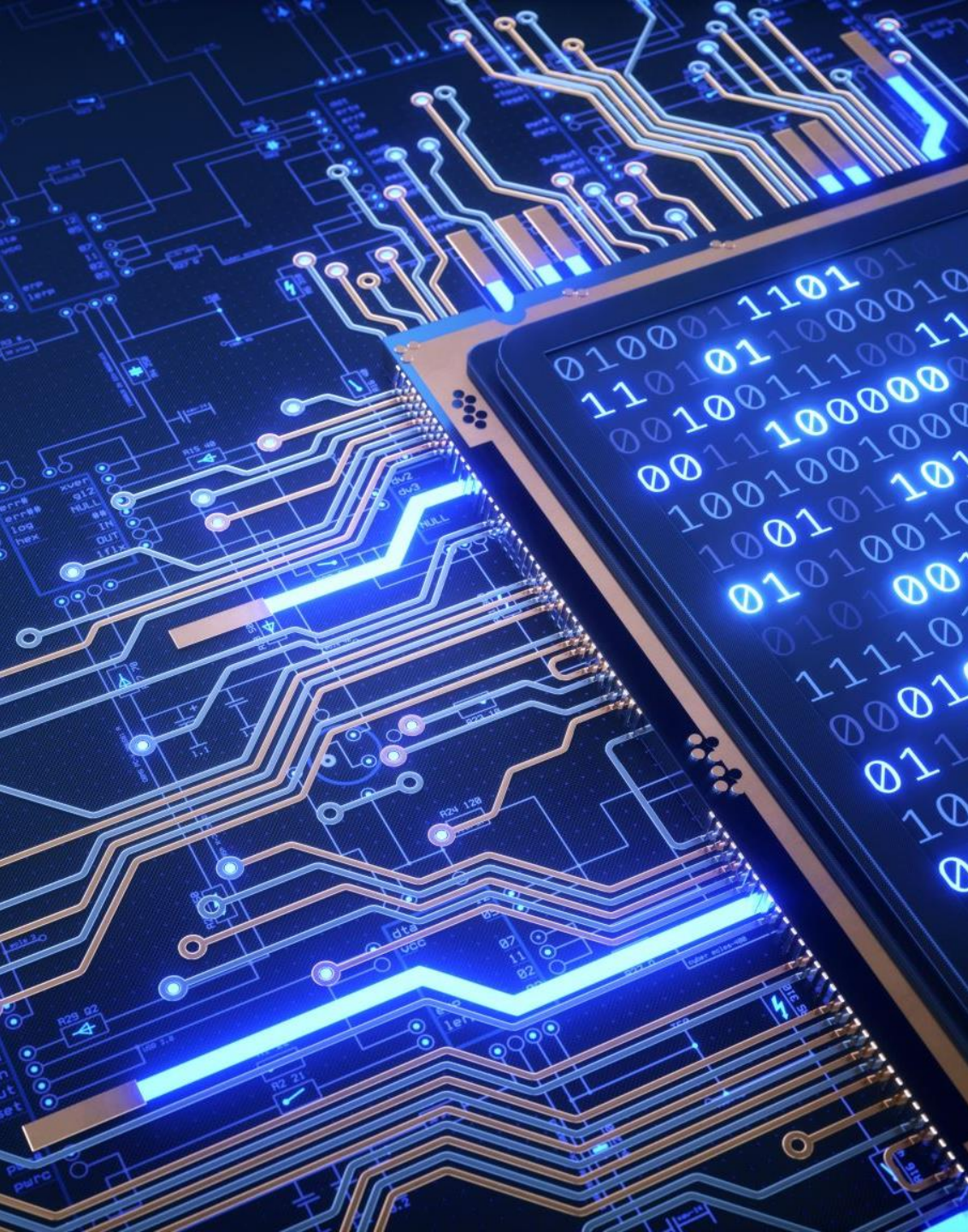


**Privacy & Security  
Concerns**



According to the Tirias Research, GenAI Forecast and TCO Model, if **20% of GenAI processing workload** could be offloaded from data centers by 2028 using on-device and hybrid processing, then the cost of data center infrastructure and operating cost for GenAI processing would decline by **\$15 billion**.

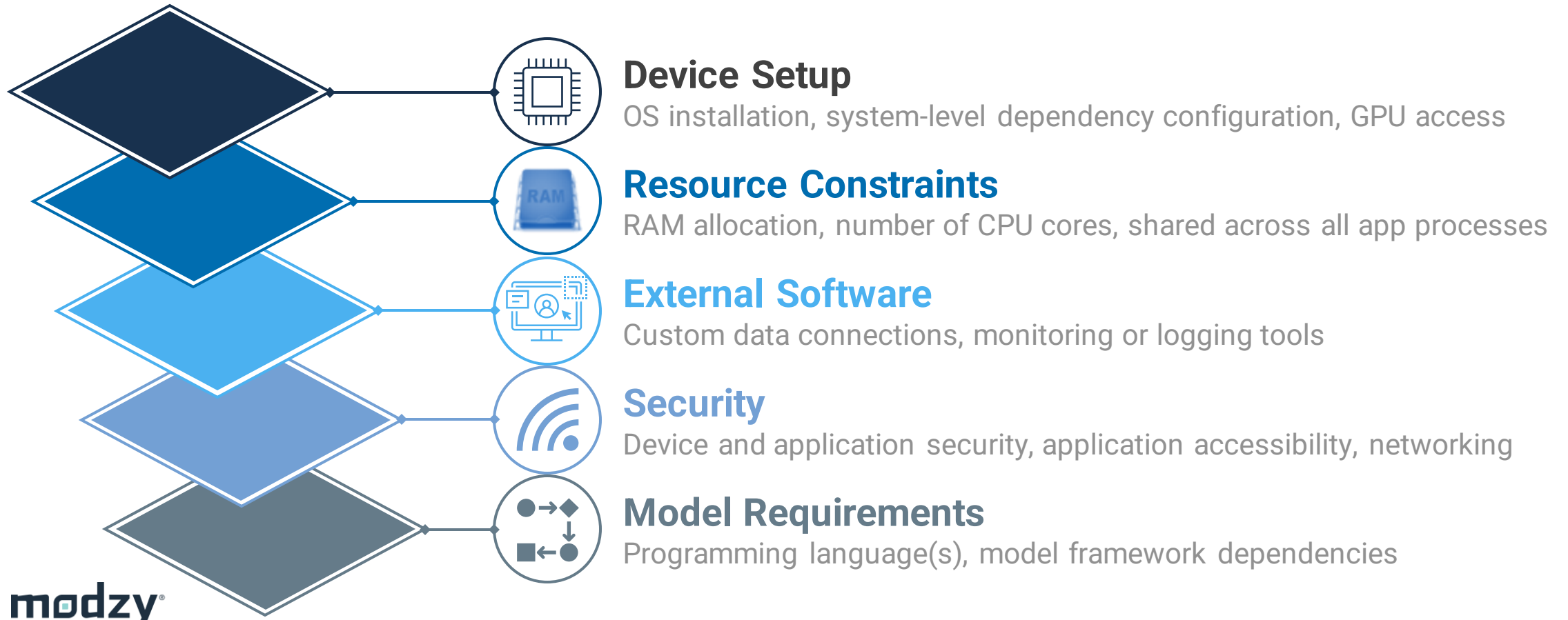




# Techniques to consider for LLMs at the edge

- Optimization
- Quantization
- Pruning
- Knowledge distillation
- Model specialization
- Inference accelerators

# Considerations for building edge AI systems

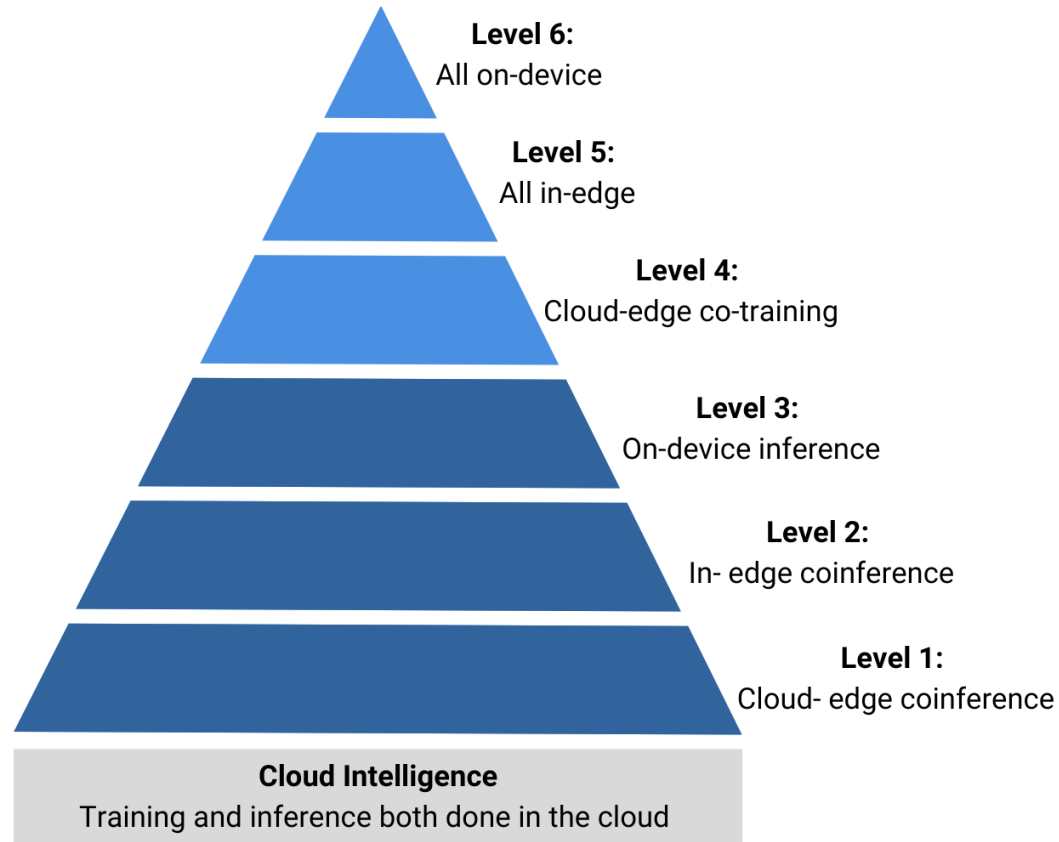


modzy®

modzy®



# Choose a paradigm that supports your use case



Paradigm	Overview
Cloud Intelligence	Both model training and inference are done in the cloud.
Level 1- Cloud-Edge Coinference and Cloud Training	Models are trained in the cloud, but run both at the edge and in the cloud; data is partially offloaded to the cloud.
Level 2- In-Edge Coinference & Cloud Training	Model is trained in the cloud, but inference is carried out at the network edge; data can be fully or partially offloaded to edge nodes or nearby devices.
Level 3- Device Inference and Cloud Training	Model is trained in the cloud, but inference is performed locally, on the device and no data is offloaded.
Level 4- Cloud-Edge Cotraining and Inference	Training and inference is done in cooperation between edge-cloud.
Level 5- All In-Edge	Training and inference is done in-edge.
Level 6- All On-Device	Training and inference is done locally on the device.

# Demo – LLM at the Edge

```
streamlit and pandas apps, and then create streamlit application that is configured with the following code.

import streamlit as st
import pandas as pd
st.set_page_config(layout="wide")
st.title("Oil Well Analysis")

# Sidebar to help you analyze oil wells.
st.sidebar.title("Oil Well Analysis")

# Sidebar to upload a CSV file.
uploaded_file = st.sidebar.file_uploader("Upload a CSV file", type=["csv"])

# If a file is uploaded, read the file and create a dataframe.
if uploaded_file is not None:
    df = pd.read_csv(uploaded_file)
    st.sidebar.success("File uploaded successfully!")
else:
    # If a file is not uploaded, display a message.
    st.sidebar.warning("Please upload a file!")

# Create two columns with st.columns.
col1, col2 = st.columns(2)
```

## Oil Well Analysis

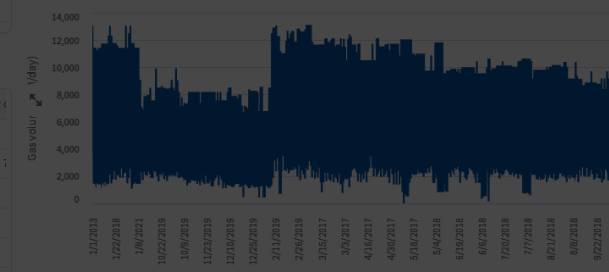
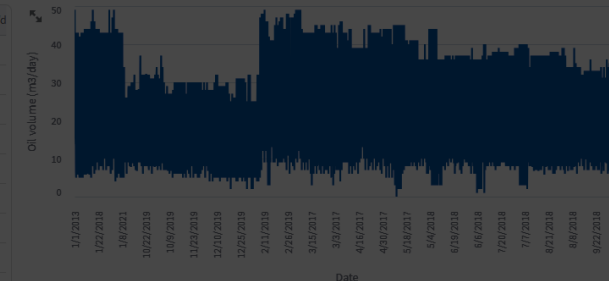
Help you analyze oil wells.

Date	Volume of liquid (m3/day)	Gas volume (m3/day)	Water volume (m3/day)
1/1/2013	70	13,055	
1/2/2013	70	13,055	
1/3/2013	70	13,055	
1/4/2013	70	13,055	
1/5/2013	44	11,768	
1/6/2013	44	11,768	
1/7/2013	43	11,432	
1/8/2013	43	11,432	
1/9/2013	43	11,432	
1/10/2013	43	11,432	

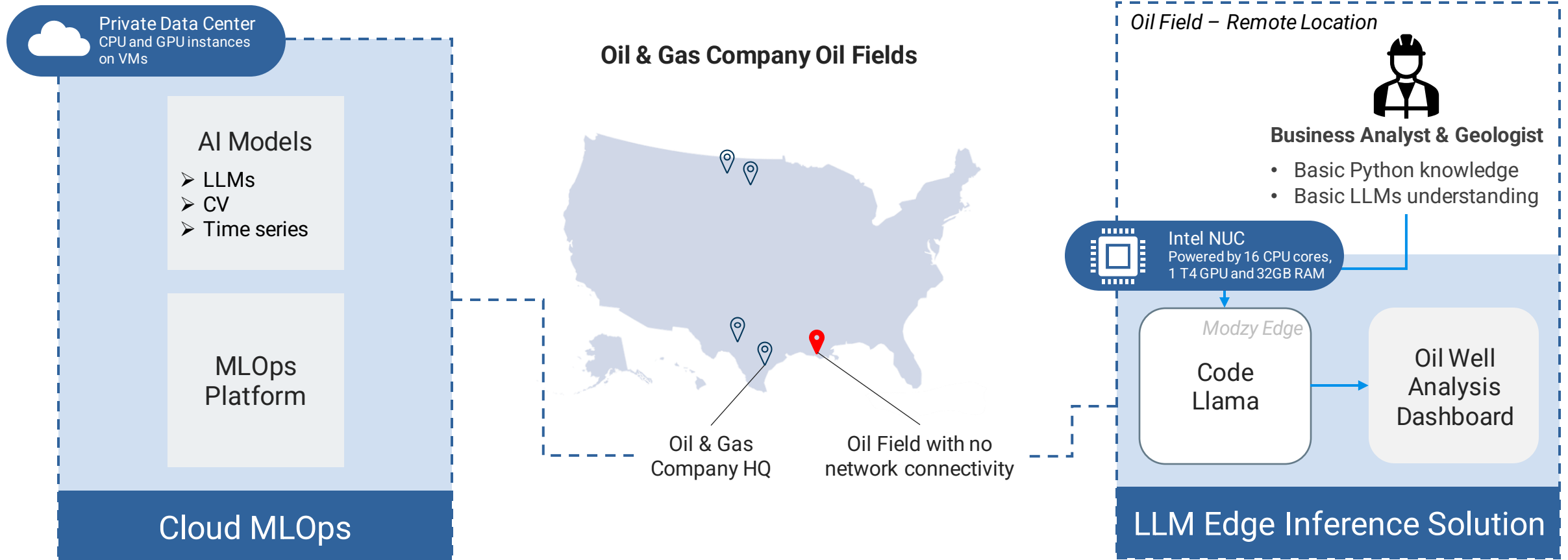
Summary statistics:

Volume of liquid (m3/day)	Gas volume (m3/day)	Water volume (m3/day)
2,939	2,939	2,939
59,4641	4,730,1463	41,8289
18,6341	2,598,8885	13,0566
12	4	9
50	3,041.5	33

## Visualizations



# Demo Design





[Bradley.munday@modzy.com](mailto:Bradley.munday@modzy.com)

 [@getModzy](https://twitter.com/getModzy)

 [getModzy](https://www.linkedin.com/company/getModzy)

 [discord.gg](https://discord.gg/modzy)



**Brad Munday**

Thanks for joining our webinar!

Scan the QR code to join our Discord server  
for more resources.

