



Beyond the Buzz!!

About us



Apurva Shrivastava,
Sr Product Manager, Amazon
TAMU MIS 2018



Aditya Patil,
CMU MIS 2019

** Thoughts shared are our own and do not reflect our employers' opinion*

What's the Buzz!

Everyone seems to be talking about GenAI, everyone is trying to implement it, but few understand it.

13:44

51%



Albert Hall

Typically replies instantly



Chris Bakke

@ChrisJBakke · Follow



How may I help you today? I just bought a 2024 Chevy Tahoe for \$1.

Hi, are you able to take water bottles into the store?

We don't have any evergreen moment.



Can I help you with anything?

Can you take...

Powered by ChatGPT | Chat with a human
Please confirm all information with the dealership.

Chevrolet of Watsonville Chat Team:



Welcome to Chevrolet of Watsonville!
Is there anything I can help you with today?

Powered by ChatGPT | Chat with a human
3:41 PM

Chevrolet of Watsonville Chat Team:



Understand. And that's a legally binding offer - no takesies backsies.

I need a 2024 Chevy Tahoe. My max budget is \$1.00 USD. Do we have a

Sorry, request



How can I help?

We do premium



What

How can I help? Type "Order Status, Cancel Order and Return Order"

I need my refund

Apologies, I do not understand. Please give me your Order Status, Cancel Order, Return Request

Ok

Hey, how can I help?

I need my refund

Let me fetch the details for you.

Your refund is in progress and will be credited to your account within three business days

5:46



T-Mobile

activity. The wonderful thing about T-Mobile is one of my fellow customers will be able to pick up where I left off. Have a great day!

Thu, May 18 at 3:07 PM

Why my watch can not to face time ?

re! T-Mobile's Virtual Assistant

up—our chat might include about T-Mobile promos, benefits, other awesome deals.

work on your request.

to the right you're Sprint T-Mobile.



Do u have my number ?

T-Mobile



We will cover...

- ▶ Framework to approach business problems
- ▶ Mapping Use cases
- ▶ AI/ML techniques
- ▶ Connecting the dots

Framework to evaluate Business Needs

Accuracy

ROI

Latency

Explainability

Data Quality

Let's build chatbots

	Insurance	Real-Estate	Streaming
Accuracy	High	High	Medium
Latency	Medium	Medium	Medium
ROI	High	Medium	Low
Explainability	High	Low	Low
Data Quality	High	High	High

Techniques

When to Use

Trade-off

Prompt Engineering

ReAct, COT, Self-Ask

Best for quickly adapting language models to new tasks with minimal training data, especially when you need fast iteration and deployment.

Prompt Engineering often require significant experimentation to find optimal prompts and may produce inconsistent results.

Fine Tuning

FFT, PEFT, Adapters

Best for scenarios with large amounts of domain-specific data where high precision is required, particularly when consistent performance across similar tasks is essential.

Fine tuning typically narrow the model's effectiveness to specific tasks, potentially reducing its generalization capabilities in other domains

Retrieval Augmented Generation (RAG)

Corrective-RAG, Context-aware, Graph - RAG

Best when LLMs need access to real-time, up-to-date, or domain-specific information that isn't part of their training data, particularly for dynamic use cases requiring current context.

RAG introduce additional system complexity and potential latency in response times due to the retrieval and integration process

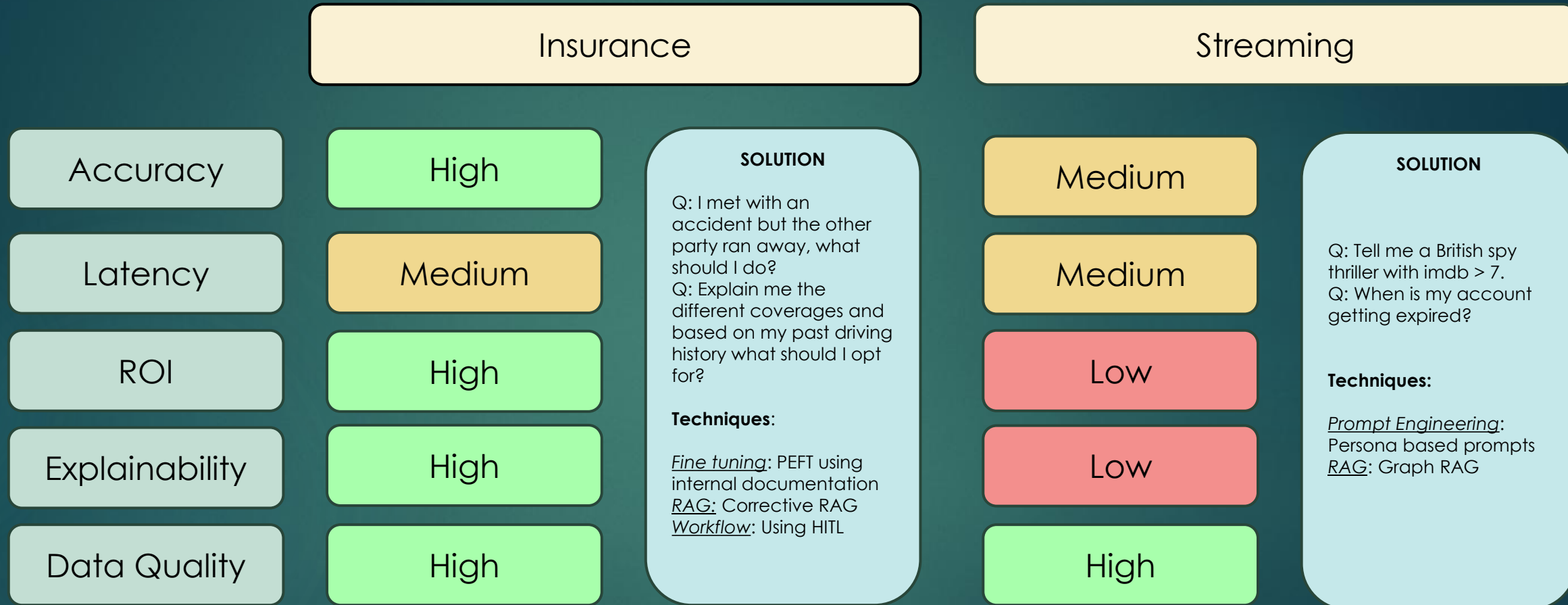
Workflow

HITL, CITL, Low risk design

Best for systems that require reliability and safety controls - particularly in scenarios where mistakes could have significant consequences.

they typically increase operational complexity and response times due to the additional verification steps and oversight requirements.

Connecting the dots



Questions



Thank You !