

Recap - AI/ML Stack  
 Dev/ML Eng / Data Scientist

- Transformer Architecture
- FM: unlabeled dataset (General purpose)
- KB (RAG?) + RAGAS Amazon Bedrock KB
- Inference Parameters (fine-tuning)
  - Temperature: low (very focused) to High (very creative) to Balanced (Mid)
  - Top P - % - Probability: 0.95 (cut dog, hot sun)
  - Top K - # - Vocab: L to M
  - max tokens # output tokens
  - Set sequences → "3" "1" "more"

Agents, Plans, Model Eval

Amazon Bedrock → (Evaluate Model, Adapt Model)

Amazon SageMaker AI

## Day 2 - Agenda

- Prompt Eng (fine-tune Models)
  - Concepts - Techniques
  - Basic Adv - Zero-shot, few-shot, CoT, ToT
  - Risks → Prompt Injection, other threats
  - Sol: Guardrails (Demo) → Model Eval, Agents, KB
- Amazon Bedrock Components - Models - IP
- Amazon Bedrock FM → API functions = invokeModel
- "Langchain" framework
  - Prompt Templates, extra features
- Architecture Patterns
  - Demo - Agents - Guardrails

Multi-modal Model = Text, images, videos, audio, web access (public) = general purpose

FM → Text generation, Text summarization, Code generation, Image generation (Diffusion Architecture), Information extraction (KB, RAG, Chatbot, Agents)

split models: Amazon
 

- Nova
- Titan

 Meta - Llama  
 Anthropic - Claude  
 AI21 - Jurassic

FM - various stages of TRG to achieve the BEST results

- Pretraining ✓ - self-supervised learning (auto-generated labels)
    - Provider RLHF (Reinforcement Learning from Human feedback)
  - fine-tuning (Supervised learning process) (add specific, smaller datasets)
  - Prompt Eng
    - Instruction fine-tuning (specific instruction) → "Prompt tuning"
    - RLHF
- FM → Text-to-Text, text to image

FM = Text to Text ⇒ LLM (NLP)  
 Natural Language Processing (Deep learning)

FM = Text to Image

Stability AI → Stable Diffusion (Model)

## Diffusion Architecture

- Deep Learning Architecture system
- 2 step process
  - forward diffusion
  - U-Net model

LLM → FM

- Clear & Concise
  - BAD: Compute the sum of total of the seq of numerals: 4, 8, 12, 16.
  - Good: What is the sum of these numbers - -- ?
- Include Context (if needed)
  - BAD: Summarize this article: [insert text]
  - Good: Provide a summary of this article to be used in a blog post: [ ]

Penalty parameters: (Text generation Model - Jurassic AI21 labs)

Help control how repetitive/varied the generated Text

- Frequency: looks at how often a token (word) has already been generated
  - effect: Too frequently, the model is less likely to generate it again
  - [very very very good]

- Presence ⇒ check whether a token is already present in the prompt
  - effect: Reduces the chance of repeating words
  - use case: Encourage new ideas rather than rephrasing the same thing

- Count Penalty = Applies a penalty based on # times a token appeared in the generated text so far.
  - Effect: Similar to frequency, more fine-grained
  - use case: keeps the generated output diverse, avoid echoing the same token