

New tech in town

amazon bedrock

is a fully managed service that lets you use foundation models (Claude, LLaMa, Mistral) without managing infrastructure

supports “agent capabilities” like RAG, guardrails and tool to use

Agent Core (Guardrails, RAG, Knowledge Bases)

guardrails: policies that filter and shape LLM output (for safety and formatting)

RAG: allows the retrieval from external knowledge databases or docs and then inject into prompts

Knowledge bases: Pre-indexed data stores that Bedrock can query for RAG

Agent Core: orchestration layer around the raw models

Amazon SageMaker

is basically a platform to build, train and deploy ML models

two layers:

- SageMaker (full platform) → train/deploy custom models with infra handled
- SageMaker AI (lower level) → Raw ML components (training jobs, endpoints, datasets)

EC2 (Elastic Compute Cloud)

is a raw compute layer - virtual machines (Linux/Windows servers) on demand

AWS SDK's

Python boto3: allows to call AWS APIs

MCP Server Protocol

MCP = Model Context Protocol is a standardized way for agents/LLMs to connect to external tools or services

- you can define tools/functions that agents can call
- you can define resources and prompts (so it can be reusable)
- any MCP-compliant client (like an LLM or orchestrator) can talk to the servers