

PartyRock



O que é: Pense no PartyRock como um playground de experimentação prático, divertido e visual.

Como funciona

É uma ferramenta no-code (sem código) que permite que qualquer pessoa crie pequenos "apps" de IA generativa em minutos.

Você descreve o que quer, arrasta e solta componentes, e ele funciona.

Tecnologia por baixo

Ele é construído em cima do Amazon Bedrock. Ou seja, quando você usa o PartyRock, ele está, nos bastidores, fazendo chamadas de API para os modelos do Bedrock.

Para quem é

Ideal para iniciantes, estudantes, gerentes de produto ou qualquer pessoa que queira aprender e prototipar o que é possível fazer com IA Generativa, sem precisar escrever código ou ter uma conta AWS.

- ❑ **Importante:** Não é um serviço para criar aplicações de produção em larga escala. É para aprendizado e descoberta.

Amazon Bedrock

O que é: Este é o serviço fundamental (o "alicerce", como o nome diz) para construir aplicações de IA Generativa na AWS.

Como funciona

É um serviço totalmente gerenciado que oferece acesso aos principais Modelos de Fundação (FMs) do mercado através de uma única API. Em vez de você ter que hospedar e gerenciar um modelo gigante, você simplesmente chama a API do Bedrock.

Principais Vantagens

Escolha de Modelos

Você não fica preso a um fornecedor. Pela mesma API, você pode acessar:

- Modelos da Anthropic (ex: Claude 3 Sonnet, Opus)
- Modelos da Meta (ex: Llama 3)
- Modelos da própria Amazon (ex: Titan Text, Titan Image Generator)
- Modelos da Cohere (ex: Command)
- E outros (ex: Stability.ai para imagens).

Personalização Privada

Você pode fazer fine-tuning (ajuste fino) desses modelos usando seus próprios dados (ex: documentos internos da sua empresa). Seus dados permanecem privados e não são usados para treinar o modelo original.

Serverless

Você não gerencia servidores. Você paga apenas pelos tokens (basicamente, a quantidade de texto) que processa na entrada e na saída.

Para quem é: Desenvolvedores e empresas que querem integrar IA Generativa em suas aplicações de forma escalável, segura e com flexibilidade de escolha.



Outros Serviços Essenciais de IA Generativa na AWS



Amazon Q 🤖 (O Assistente)

O que é: É o assistente de IA da AWS, focado em negócios e desenvolvimento. Ele é um "usuário" dos modelos do Bedrock, treinado para tarefas específicas.

Onde ele atua:

- **Amazon Q for Business:** Conecta-se às fontes de dados da sua empresa (documentos, wikis, Slack, S3) para responder perguntas específicas do seu negócio. (Ex: "Qual foi o faturamento do produto X no último trimestre?", "Resuma nossa política de férias.").
- **Amazon Q for Developers (antigo CodeWhisperer):** Funciona dentro da sua IDE (como o VS Code) para sugerir código, explicar blocos de código, fazer debug e até ajudar a migrar versões de linguagens (ex: Java 8 para 17).
- **Amazon Q in QuickSight:** Ajuda a criar dashboards e análises de BI (Business Intelligence) usando linguagem natural.



Amazon SageMaker 👩‍🔬 (A Plataforma Completa de ML)

O que é: Se o Bedrock é para usar modelos prontos via API, o SageMaker é a plataforma completa para construir, treinar e implantar seus próprios modelos de Machine Learning (incluindo os generativos) do zero.

Para IA Generativa:

- **Controle Total:** Você usa o SageMaker se quiser ter controle absoluto sobre a arquitetura do modelo, o processo de treinamento e a infraestrutura de implantação (ex: escolher tipos de GPU específicas).
- **SageMaker JumpStart:** É um "catálogo" dentro do SageMaker que oferece acesso a muitos modelos de código aberto (como o Llama) que você pode implantar e gerenciar facilmente em suas próprias instâncias.

Para quem é: Cientistas de Dados e Engenheiros de Machine Learning que precisam de controle granular para treinar ou hospedar seus próprios modelos.