

Building Next-Generation LLM Applications with Private Models

Using the Reasoning Workbench to Better Select Private LLM Models

Josh Patterson

Introduction

In the wake of the groundbreaking announcements at AWS Re:Invent 2023, the landscape of Generative AI has witnessed a transformative shift, ushering in a new era of possibilities. Among the notable revelations, AWS Bedrock has swiftly emerged as the premier platform for deploying cutting-edge Generative AI applications. A pivotal requirement for operating private applications in this dynamic space lies in the utilization of private models, distinct from widely accessible options such as those provided by openAI.

As businesses strive to harness the potential of Generative AI, the imperative to seamlessly integrate private models into their infrastructure becomes paramount. In this rapidly evolving scenario, the interplay between AWS Bedrock and private models is poised to redefine the way enterprises approach artificial intelligence. This article delves into the intricate details of this intersection, unveiling the strategic implications for CXO-level executives in Fortune 500 companies. As we navigate through the intricacies of deploying private models, the significance of leveraging AWS Bedrock for next-generation AI applications will become abundantly clear.

AWS Bedrock and Private LLMs

The synergy between AWS Bedrock and the realm of Generative AI extends beyond a mere platform-provider relationship; it encapsulates a transformative paradigm shift. At the heart of this synergy lies a sophisticated integration of enterprise data within AWS—a treasure trove of knowledge that serves as the foundation for next-level AI applications. AWS Bedrock, with its robust capabilities, seamlessly accommodates the complex landscape of private models, providing a secure and scalable environment for their deployment.

To navigate the nuanced landscape of Generative AI effectively, it is essential to

adopt a strategic framework. The “knowledge vs reasoning” paradigm offers a lens through which enterprises can conceptualize the relationship between their proprietary data and the reasoning power of Large Language Models (LLMs). This framework is exemplified through the innovative use of the RAG (Retrieval-Augmented Generation) design pattern. By bridging the gap between enterprise knowledge stored in AWS and the reasoning prowess of LLMs, organizations can propel themselves into the forefront of AI innovation. The RAG design pattern serves as the connective tissue that enables the construction of next-generation AI applications on AWS Bedrock, marking a pivotal evolution in the strategic deployment of Generative AI.

Challenges with Using Private LLMs

While the marriage of AWS Bedrock and private models ushers in a new era of AI possibilities, selecting the right private model poses a nuanced challenge for enterprises. Despite the allure of openAI’s formidable models like GPT-3 and GPT-4, their complexity makes them less amenable to efficient operation as private LLM model instances. Conversely, opting for smaller private models may enhance operational efficiency, but they often lack the horsepower required for certain types of complex reasoning—a critical aspect for robust AI applications.

The conundrum faced by organizations revolves around finding a delicate balance—choosing a model that aligns with the specific requirements of their tasks and data. In the pursuit of efficient and cost-effective AI applications, there arises a pressing need for a systematic approach to compare different models accurately. Enterprises grapple with the challenge of evaluating these models against their unique datasets, seeking a method that ensures optimal performance without compromising efficiency. In this landscape of varied choices, the key lies in identifying a solution that facilitates a nuanced comparison tailored to the specific needs of the enterprise.

Enter: The Reasoning Workbench

Addressing the challenge of selecting the right private model, Patterson Consulting has spearheaded the development of a groundbreaking open-source tool known as the “Reasoning Workbench.” This tool stands as a testament to Patterson Consulting’s commitment to empowering enterprises in their quest for optimal AI model selection.

The Reasoning Workbench provides a straightforward and effective means to generate scores for comparing different models, ensuring that enterprises can make informed decisions tailored to their unique requirements. Its intuitive interface allows users to not only assess general model performance but also to generate comparative scores for prompts specifically crafted for their enterprise. This bespoke approach ensures a nuanced evaluation, aligning the chosen model

seamlessly with the intricacies of individual AI applications.

With the Reasoning Workbench, enterprises can swiftly gain a relative sense of how different models perform in the context of their specific tasks and datasets. This tool proves invaluable in the decision-making process, enabling organizations to select the model that best aligns with their use case. The Reasoning Workbench, thus, emerges as a powerful solution, streamlining the model selection process and paving the way for the rapid and confident deployment of AI applications on AWS Bedrock.

Benefits of Evaluating LLMs

The adoption of the Reasoning Workbench by enterprises brings forth a cascade of benefits, fundamentally transforming the landscape of AI application deployment. Rapid and informed decision-making regarding the choice of models allows teams to act with confidence, facilitating the swift deployment of cutting-edge AI applications. This efficiency not only enhances operational agility but also translates into tangible cost savings, as better-performing models typically require fewer tokens.

Moreover, the streamlined model selection process reduces the need for extensive prompt engineering, minimizing the complexity of interfacing with private models. This, in turn, leads to fewer inferences to private knowledge, optimizing the overall efficiency of AI applications. As a result, end-users experience faster and more responsive interactions, contributing to an enhanced user experience.

In essence, the adoption of the Reasoning Workbench not only optimizes costs but also streamlines operations, making it a pivotal tool for enterprises aiming to leverage the full potential of AWS Bedrock and private models. The net result is a transformative shift towards a more efficient, cost-effective, and user-friendly AI ecosystem.

Taking the Next Step with Private LLMs and AWS Bedrock

To embark on the journey of empowered AI application deployment, we invite CXO-level executives in Fortune 500 companies to seize the opportunity and explore the capabilities of the Reasoning Workbench. Sign up today for a private demonstration and a collaborative brainstorming session, tailored to unveil the immense potential this tool holds for your organization.

Patterson Consulting, the driving force behind the Reasoning Workbench, extends an exclusive offer to select AWS customers—a complimentary use case analysis session. This bespoke consultation will delve into your organization's specific needs and challenges, guiding you towards a strategic integration of private models with AWS Bedrock. The future of AI application deployment is at your

fingertips; take the first step towards a more agile, cost-effective, and user-centric AI ecosystem. Sign up now and revolutionize the way your organization leverages Generative AI on AWS Bedrock.