



# binbash®

## Integrating AWS Bedrock's LLM-Powered RAG Solution into B2Chat's Multichannel Platform

### Executive Summary

B2Chat, a multi-agent, multi-channel customer messaging platform, aims to enhance customer interactions across popular channels like WhatsApp, Facebook Messenger, and Instagram. To elevate its user experience and achieve real-time data-driven responses, B2Chat partnered with binbash, an AWS Advanced Tier Services Partner, to implement a Proof-of-Concept (PoC) for a Large Language Model (LLM) with Retrieval-Augmented Generation (RAG) using AWS Bedrock. This case study outlines how binbash successfully deployed a scalable, secure, and efficient cloud-based solution, leveraging binbash Leverage™ and Infrastructure as Code (IaC) principles to accelerate delivery.

### Customer Challenge

As B2Chat sought to expand its AI capabilities, the main challenges were:

- Integration with AI Models: Efficient integration of LLMs to deliver accurate, real-time responses.
- Data Retrieval: Implementing a robust RAG system for seamless data retrieval across diverse customer interactions.
- Scalability & Security: Ensuring a scalable and secure cloud infrastructure to support increased user demands and AI-powered features.
- Operational Efficiency: Streamlining cloud deployment processes for rapid updates and adjustments.

### Solution

To address B2Chat's needs, binbash designed a comprehensive solution using AWS Bedrock, Terraform, and binbash Leverage™. The solution included:

1. **AWS Bedrock Integration:** Deployed foundation models and embedding models for LLM capabilities, enabling accurate text

## B2Chat

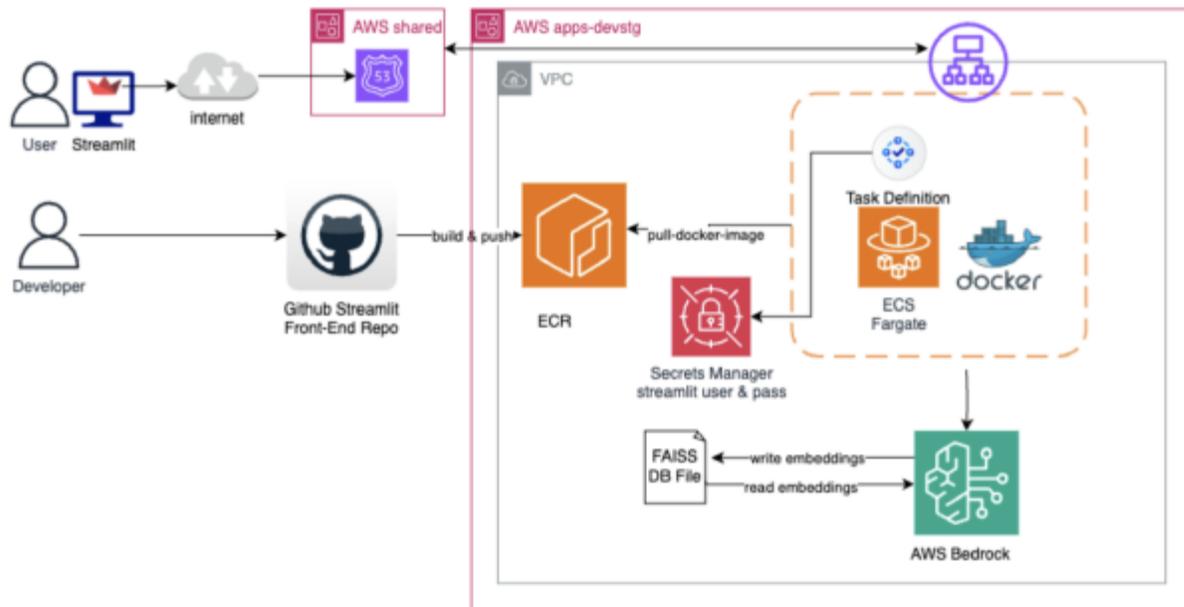
B2Chat is a multi-agent, multi-channel platform that centralizes customer interactions across various messaging services, including WhatsApp, Facebook Messenger, Instagram, Telegram, and LiveChat. It enables businesses to respond to all messages from these platforms in one consolidated inbox, facilitating efficient communication and turning every conversation into a potential sale. Designed to enhance customer satisfaction and streamline sales processes, B2Chat serves over 600 companies worldwide, helping them grow their online presence and improve client engagement.



# binbash®

processing and response generation tailored to customer queries.

2. **Retrieval-Augmented Generation (RAG):** Implemented RAG logic for efficient document indexing, retrieval, and response generation, integrating various AWS services.
3. **Infrastructure as Code (IaC):** Used Terraform to deploy cloud resources, ensuring automated, repeatable, and consistent infrastructure management.
4. **Containerized Deployment:** Developed a Docker-based deployment for the backend/frontend, streamlining the deployment of AI models and RAG components to AWS ECR.



## Key Components of the solution

- **LLM Integration with AWS Bedrock:** Configured foundation and embedding models to optimize AI performance and deliver contextually relevant responses.
- **RAG System:** Enhanced document indexing and retrieval using the RAG architecture, supporting vector-based searches with FAISS.
- **AWS ECR Deployment:** Built and pushed Docker images to AWS ECR for consistent and scalable model deployment.
- **Cost Monitoring:** Integrated AWS CloudWatch Alarms, Budgets, and Cost Explorer to manage costs and enhance cloud resource utilization



# binbash<sup>®</sup>

## Key Milestones

- **AWS Bedrock Integration:** Successfully deployed foundation models, embedding models, and RAG architecture within AWS Bedrock.
- **PoC Deployment:** Delivered a functional PoC demonstrating the LLM and RAG capabilities with real-time response generation.
- **Documentation & Knowledge Transfer:** Provided detailed documentation and knowledge transfer sessions to ensure B2Chat's team could manage and operate the AI-powered platform.

## Results

- **Enhanced AI Performance:** Integrating LLM and RAG through AWS Bedrock improved B2Chat's response accuracy and efficiency, leading to faster customer service interactions.
- **Increased Operational Efficiency:** The use of binbash Leverage™ accelerated the deployment process, achieving rapid implementation of AI models and infrastructure.
- **Improved Cost Management:** Proactive cost monitoring and optimization reduced operational expenses, enhancing resource utilization and cloud cost efficiency.
- **Scalability & Security:** The solution ensured a secure, scalable cloud infrastructure aligned with AWS Well-Architected best practices.

## Conclusion

The collaboration between B2Chat and binbash resulted in a robust, AI-powered messaging solution, leveraging AWS Bedrock for LLM and RAG capabilities. The implementation enabled B2Chat to improve customer interactions across multiple channels, setting a strong foundation for future AI-driven growth.