

AMBUJ KUMAR TRIPATHI

AI Prompt Engineer | GenAI Developer | RAG & Conversational AI

Available Immediately | Remote/Hybrid | Open to Relocation

Email: kumarambu8@gmail.com | Mobile: +91 9431 801363 | Location: Gorakhpur, Uttar Pradesh, India

Portfolio: <https://ambuj-portfolio-v2.netlify.app/> | GitHub: <https://github.com/Ambuj123-lab>

PROFESSIONAL SUMMARY

Generative AI Engineer and Automation Architect specialized in commercial-grade LLM optimization, RAG pipelines, and Conversational AI. **Architecture-first builder:** define system design, APIs, data flows, guardrails, and evaluation—then use **AI-assisted coding** to implement and iterate fast. Leveraging **5+ years of technical experience** in Telecom and Tech sectors to build cost-efficient, production-ready AI solutions. Successfully transitioned from legacy automation frameworks to GenAI-driven architectures, managing publicly deployed AI agents processing **525,000+ tokens** in live environments with **Langfuse observability** and **Redis** for real-time user analytics.

KEY TECHNICAL ACHIEVEMENTS

- Stakeholder & Client Management:** Expert in translating complex client requirements into technical roadmaps. Proven track record of managing expectations for high-value clients (e.g., MoD, British Telecom), conducting technical feasibility workshops, and delivering demos to non-technical stakeholders
- Production-Scale GenAI Deployment:** Deployed **5+ scalable AI agents** on Render and Streamlit Cloud, successfully processing **525,000+ tokens** and **10,000+ conversation turns** tracked via Langfuse traces and OpenRouter API routing
- Production Observability and Monitoring:** Integrated **Langfuse** to monitor production metrics including request traces, response latency, and token usage patterns, providing visibility into model performance and API costs
- Real-Time User Analytics:** Implemented **Redis** to track active user sessions and concurrency in real-time, providing instant visibility into system load and usage spikes without database latency
- Enterprise Data Privacy:** Architected **Microsoft Presidio + SpaCy PII masking layer** to detect and redact sensitive information (Names, Emails, Phone Numbers) before LLM processing, ensuring **GDPR-compliant inference pipelines**
- User Engagement Analytics:** Integrated **Google Analytics 4 (GA4)** on the portfolio UI to monitor user sessions, tracking click-through rates (CTR) on AI demos and project dwell time to analyze real-world engagement patterns
- Systematic Instruction/Prompt Optimization:** Built reusable prompt libraries and comprehensive documentation for diverse creative requirements, reducing prompt engineering time by **40%** through systematic optimization methodology and team feedback-driven iteration
- Commercial AI Production Experience:** Engineered prompts for commercial-grade generative AI models at **Hogarth Worldwide (WPP)** with strict technical constraints

Core Strengths: GenAI System Architecture | RAG Pipelines | LLM API Integration | Conversational AI Development | Prompt Engineering | Model Validation | Red Teaming | Python/JavaScript/Node.js | Cloud Platforms | Production AI Workflows | Data Analysis | Automation | Quality Assurance

GENERATIVE AI PROJECTS (ALL LIVE & ACCESSIBLE)

Production Systems:

Citizen Safety & Awareness AI - Full-Stack RAG System | January 2026

- Architected production-grade **full-stack RAG chatbot** for answering citizen safety queries from **8 government legal documents** (RBI advisories, POSH Act, POCSO Act, Cyber fraud alerts) using **dual-indexing architecture** with metadata tagging for surgical deletion
- Implemented **real-time PII detection and masking** using **Microsoft Presidio + spaCy NLP** with custom Indian phone number regex patterns before LLM processing
- Engineered **batched embedding with rate-limit-aware delays** (15 docs/batch, 10s intervals) to stay within Google AI's 100 RPM free-tier quota, processing documents efficiently without API throttling
- Designed **fault-tolerant error handling** with Circuit Breaker pattern and rate limiting using **pybreaker and SlowAPI** middleware, ensuring 99% service reliability

- Integrated **observability and analytics stack** with Langfuse for trace monitoring and Upstash Redis for real-time user session tracking
- Deployed with **OAuth 2.0 authentication** and JWT-based session management, featuring automated CI/CD pipelines on Render (Backend) and Vercel (Frontend)
- **Live Demo:** <https://citizen-safety-ai-assistant.vercel.app/>

Technologies: React, Vite, FastAPI, LangChain, ChromaDB, Llama 3.3 70B (OpenRouter), Google Gemini Embeddings, Microsoft Presidio, spaCy, MongoDB Atlas, Upstash Redis, Langfuse, pybreaker, SlowAPI

Secure Enterprise RAG Assistant with Auto-Recovery Architecture | January 2026

Core RAG Pipeline & Auto-Recovery:

- Architected production-grade RAG system using **LangChain** and **Llama 3.3 70B**, featuring **auto-recovery mechanism** that automatically rebuilds vector database (ChromaDB) from source upon corruption, ensuring **98% system reliability**
- Developed **Custom Confidence Scoring Algorithm** using Vector Distance Normalization to quantify retrieval quality and flag potential hallucinations in real-time, enhancing user trust through transparency
- Optimized indexing latency by **98% (2 minutes to under 2 seconds)** through incremental document processing, utilizing **RecursiveCharacterTextSplitter (1000/200)** for semantic context preservation

Session Management & Data Persistence:

- Implemented **Multi-Turn Conversational Memory** using sliding window (last 3 messages) with **MongoDB Atlas** for session persistence, enabling context-aware responses across email-based authenticated sessions
- Configured **MongoDB TTL Index for automatic 30-day data expiration**, ensuring GDPR-compliant data retention and preventing storage overflow in production environments
- Integrated real-time analytics using **Upstash Redis** to track user engagement and session metrics securely

Enterprise Security & Compliance:

- Integrated **Enterprise Privacy Layer** using **Microsoft Presidio** and **SpaCy**, detecting and masking PII (Names, Phones, Emails) before LLM processing to ensure GDPR compliance
- Developed **Content Moderation System** with regex-based abusive language filter and **zero fabrication policy** through strict system prompts, ensuring safe and reliable user interactions

Monitoring & Deployment:

- Integrated **LangFuse Observability Platform** to monitor production-grade metrics including Traces, Latency, and Token Cost, providing full visibility into model reasoning processes and enabling data-driven optimization
- Optimized User Experience with **Interactive Dashboard** featuring real-time latency tracking (Streamlit metrics) and Source Viewer component, following enterprise-grade UI standards for production AI systems
- Engineered **OS-agnostic deployment pipelines**, resolving cross-platform path inconsistencies (Windows/Linux) and enabling seamless cloud deployment on Streamlit Cloud
- **Live Demo:** <https://ambuj-rag-chatbot.streamlit.app/>

Technologies: Python, LangChain, LangFuse, Llama 3.3 70B, ChromaDB, Streamlit, Microsoft Presidio, SpaCy, MongoDB Atlas, Upstash Redis, RecursiveCharacterTextSplitter

LLM-Integrated Resume Chatbot | Meta Llama 3.3 70B | August 2025

- **Stateless Backend Architecture:** Developed a lightweight Flask REST API where conversation history is managed **client-side** and passed via JSON payloads, eliminating the need for server-side databases (Redis/SQL) and reducing latency
- **Context-Aware Prompt Engineering:** Implemented **Dynamic Context Injection** by loading structured resume data (JSON) directly into the LLM's System Prompt, ensuring factual accuracy without complex RAG pipelines
- **AI Guardrails and Moderation:** Coded custom Python logic (**is_abusive function**) to filter toxic keywords before API calls, ensuring interaction safety and professional output constraints
- **Production Deployment:** Configured **Gunicorn (WSGI)** as the application server on Render to handle concurrent requests efficiently, utilizing python-dotenv for secure API key management
- **Error Resilience:** Built robust exception handling for OpenRouter API to gracefully manage timeouts and upstream errors, ensuring continuous service availability
- **Live Demo:** <https://ambuj-resume-bot.onrender.com>

Technologies: Meta Llama 3.3 70B, OpenRouter API, Python, Flask, Gunicorn, JSON Context Injection

Additional AI Projects

- **Geo AI:** Quick prototype exploring Gemini Flash 2.0's multimodal capabilities with Leaflet.js for interactive map analysis
- **Smart AI Prompt Builder:** Bilingual tool (Hindi/English) with keyword matching for 8 problem categories and 6 tone styles
- **IBM Watson DialogFlow:** Rule-based chatbot achieving 95% intent recognition, deployed on Netlify CDN with 99.9% uptime
- **Resume-as-Code Pipeline:** Automated PDF generation using Puppeteer, achieving 100% formatting consistency
- **Task Management PWA:** Privacy-first app with LocalStorage, Service Workers, Web Audio/Speech API integration

CORE COMPETENCIES & TECHNICAL SKILLS

Generative AI & LLM Engineering

Prompt Engineering, Large Language Models (GPT-4, Claude, Gemini, Meta Llama 3.3 70B, IBM WatsonX), RAG Pipelines (Retrieval-Augmented Generation), Vector Databases (ChromaDB, FAISS), Semantic Search and Embeddings, LangChain Orchestration, LangFuse Observability, Multi-Turn Dialog Systems, Sliding Window Memory, Hallucination Control, Context-Aware Systems, LLM API Integration, Conversational AI Development, Chatbot Development, Multi-Modal AI, Adversarial Testing, Red Teaming, Model Validation, Model Robustness Testing, Few-Shot Learning, Zero-Shot Learning, Token Optimization, Content Evaluation, Content Moderation Systems, Data Quality Assessment, Fact-Checking and Verification, Responsible AI, AI Ethics

Programming Languages & Frameworks

Python, JavaScript (ES6+), Node.js, Express.js, HTML5, CSS3, React 18, Flask, Gradio, Streamlit, Tailwind CSS, Babel, RESTful API Development, Async/Await, Promise Handling

AI Tools & Platforms

MongoDB Atlas, Microsoft Presidio (PII Masking), SpaCy (NLP), LangChain, LangFuse, ChromaDB, Upstash Redis (Analytics), Hugging Face Transformers, OpenRouter API, Gemini API, Google Generative AI SDK, Gemini Vision API, Vertex AI, IBM WatsonX, Jupyter Notebooks, VS Code, NotebookLM, REST APIs, Puppeteer, Git, Conversation State Management, Multi-Turn Dialog Systems

Cloud Platforms

Google Cloud Platform (Vertex AI, Cloud Functions, Cloud Run, GKE), IBM (WatsonX Services, Cloud Foundry, Speech to Text, Text to Speech), Streamlit Cloud, Netlify, Render.com, Vercel, Hugging Face Spaces

Specialized Skills

Enterprise Data Privacy (GDPR Compliance), Incremental Indexing Logic, Progressive Web Apps (PWA), Service Workers, Offline Functionality, Data Processing and Visualization (Papa Parse, SheetJS, Plotly.js), Web Speech API, Web Audio API, CORS, Rate Limiting, Environment Secrets Management, Security Best Practices, Input Validation, API Rate Limiting, Session Management, PDF.js, html2canvas, jsPDF, Conversational Memory Architectures, Regex-based Content Moderation, Content Moderation Systems, Pattern-Based Guardrails, Circuit Breaker Pattern, OAuth 2.0, JWT

Additional Technical Skills

Data Analysis, Automation, Network Optimization, GIS Tools, Fiber Optic Network Design, Excel Advanced Functions, Business Process Documentation, Performance Monitoring, Quality Assurance, Cross-functional Collaboration, Technical Documentation

PROFESSIONAL EXPERIENCE

Generative AI Engineer (Creative AI Models)

Hogarth Worldwide (WPP Marketing Communications India Pvt. Ltd. via TeamLease) | Remote, India

September 2025 – October 2025 | **6-Month Contract (Concluded due to mandatory office relocation requirement - health reasons)**

- Engineered and optimized prompts for commercial-grade generative AI models (**Flux, SDXL**) within strict technical constraints (**72-77 token limits** across Base, Positive, Negative prompts) for brand-compliant visual outputs
- Conducted systematic model validation and adversarial testing, identifying critical limitations including instruction adherence failures, hallucination patterns, layout inconsistencies, and edge-case behaviors across multiple parameters

- Developed **Smart AI Prompt Builder tool** to standardize prompt creation workflows across team, featuring intelligent keyword matching algorithm detecting **8 problem categories** and recommending optimal AI platforms based on task requirements
- Collaborated with cross-functional teams to define AI integration protocols and quality assurance standards for production pipelines
- Developed prompt refinement methodology translating subjective business feedback into quantifiable technical adjustments

Key Achievement: *Created systematic prompt optimization framework reducing iteration cycles and improving output quality consistency across diverse creative requirements*

Technologies: Flux, SDXL, Prompt Engineering, Adversarial Testing, Model Validation, QA Protocols

Associate Engineer - Fibre & Network Delivery (Applied AI & Workflow Automation)

British Telecom Global Services Pvt. Ltd. | Gurugram, India

January 2022 – August 2024

- Led FTTP network planning using **GIS tools**, reducing deployment time by **15%** through systematic optimization
- Streamlined network estimation workflows by automating data extraction from **Piper database** to structured documentation including duct diagrams, cable diagrams, and SOC costing sheets, reducing manual effort by **60%**
- Collaborated on transition from manual processes to **Autotron-based automation pipeline**, standardizing documentation workflows across engineering teams
- Developed Excel-based automation solutions and Python scripts for repetitive data transformation tasks, improving data accuracy and processing speed
- Explored AI-powered document intelligence tools (May-August 2024) using **Gemini 1.5 Flash API** for automated extraction of technical parameters from network specifications and site survey reports
- Prototyped conversational AI interface using **Gradio** enabling stakeholders to query network deployment metrics through natural language, demonstrating potential for cross-team accessibility improvements
- Experimented with multimodal AI capabilities for automated interpretation of network design maps and infrastructure blueprints, achieving promising results in pilot testing phase
- Collaborated with MIS/Power BI teams to explore AI-driven analytics solutions and automation pipelines, receiving recognition for innovative problem-solving approaches
- Coordinated **onshore (UK) engineers** and **offshore (India) delivery teams**, achieving **98% quality assurance compliance** with streamlined processes
- Awarded "**Top Performer**" (September 2022) for innovative problem-solving and process improvement initiatives

Innovation Achievement: *Successfully automated critical network planning workflows and pioneered exploration of AI technologies for telecom operations, establishing foundation for intelligent automation initiatives*

Technologies: GIS Tools, Network Optimization, Piper Database, Autotron, Excel Advanced Automation, Python Scripting, Gemini 1.5 Flash API, Gradio

Operations & Maintenance Engineer

Lobo Staffing (Client: Tata Communications Transformation Services Ltd.) | Lucknow, India

November 2021 – January 2022

- Optimized network performance and monitoring systems, reducing downtime by **10%** through comprehensive data analysis
- Managed billing systems for OSP/ISP services with **100% accuracy** using systematic quality control processes
- Improved Field Engineer SLA compliance, achieving **98% TAT** through workflow optimization

Technologies: ServiceNow, Orion, Network Monitoring Systems, Excel Advanced Functions

Optical Fiber Execution Engineer

Annu Infra Construct India Pvt. Ltd. | Kolkata, India

December 2017 – June 2018

- Managed on-ground execution teams for high-stakes Ministry of Defense projects. Orchestrated daily operations for 15+ technical staff, ensuring 100% adherence to security protocols and project timelines
- Executed fiber optic deployment projects for **Ministry of Defense** across **3 states** ensuring complete security compliance
- Ensured **100% compliance** with security protocols and technical specifications through rigorous quality assurance
- Reduced project delays by **20%** through effective vendor coordination and strategic resource optimization

System Administrator

Teamware Solutions (Client: Tata Consultancy Services) | Noida, India

October 2013 – November 2014

- Managed revenue systems for **Delhi Jal Board**, ensuring **99% uptime** through proactive monitoring
- Trained **20+ portal users** on system optimization and best practices
- Team Capability Building:** Spearheaded training programs for 20+ portal users and junior staff, optimizing system adoption and operational efficiency for the eMigrate project.
- Supported **eMigrate portal (Ministry of Overseas Affairs)** with document verification processes

INDUSTRY CERTIFICATIONS & CREDENTIALS (ALL VERIFIED)

AI/ML & Generative AI Certifications

- NVIDIA:** Building RAG Agents with LLMs, AI on Jetson Nano
- Google Cloud:** 6 Skill Badges — Vertex AI Prompt Design, Gemini API, TensorFlow Image Classification, GenAI Apps with Streamlit, Cloud Run & GKE Deployment
- IBM AI:** AI Fundamentals, Generative AI in Action, Deep Learning Essentials, Chatbot Development, Python for Data Science
- Azure AI:** Language Models on Databricks, Responsible AI Principles
- Anthropic Academy:** Model Context Protocol, Claude with Google Vertex AI
- Linux Foundation:** Ethical Principles for Conversational AI

Enterprise Job Simulations

- BCG X:** AI-Powered Financial Chatbot Development
- Big 4 & Enterprise:** AWS Solutions Architecture, PwD Digital Assurance, Deloitte Australia Analytics, Tata iQ GenAI, Siemens

All certifications verifiable through official platforms. Verification links available in online portfolio.

EDUCATION

Post Graduate Diploma in Power Transmission & Distribution

National Power Training Institute, Ministry of Power | 2015 – 2016 | New Delhi, India

Bachelor of Technology - Electrical & Electronics Engineering

Uttar Pradesh Technical University | 2009 – 2013 | Lucknow, India