# SOUBHAGYA SRIVASTAVA

Patna, Bihar, India • soubhagyasrivastava240@gmail.com • +91 8967010103 • [LinkedIn](#) • [GitHub](#)

## PROFESSIONAL EXPERIENCE

**AI/ML Developer Intern**                                                                                              OCT 2024 - Present
*GarudaX*                                                                                                                                    **Patna, Bihar**

- Assisted in developing and deploying DigiAstra, a Python-based vulnerability identification tool using 10+ security modules, achieving high accuracy in detecting SQL injection, security headers, tech stack, and subdomain enumeration.
- Contributed to building an AI-powered Data Loss Prevention (DLP) system with 87%+ accuracy across 15+ file types, including malware detection and multi-format data classification.
- Contributed to developing an Intrusion Detection System (IDS) and network anomaly detection models, achieving 92% accuracy on large-scale datasets using neural networks.
- Participated in internal presentations and product demos, contributing to product improvements and cross-functional discussions on AI-driven cybersecurity solutions.

**Data Science Team Lead**                                                                                                   **Gujarat, India**
*Purezza Technologies*                                                                                                       JAN 2024 - Present

- Assisted in developing an AI-powered invoice generation chatbot using NLP techniques with automatic spelling correction and dynamic product price retrieval from databases.
- Contributed to building LLM-based content generation models using Cohere API for personalized, multi-slide presentation creation.
- Supported image generation models with Stable Diffusion and Hugging Face Transformers for high-quality AI-driven visual content.
- Devised and deployed 50+ predictive models using optimized ML algorithms, ensuring high accuracy, scalability, and real-world applicability across diverse domains.
- Contributed to developing and deploying 50+ predictive models with optimized data pipelines, ensuring high accuracy, scalable deployment, and production-level performance optimization.

## SKILLS AND TECHNOLOGIES

- **AI & Model Development:** Machine Learning, Deep Learning, LLMs, Generative AI, Agentic AI, NLP, Transformers, Data Preprocessing, Feature Engineering, Model Deployment.
- **Tech Stack & Tools:** TensorFlow, PyTorch, Scikit-learn, SpaCy, Hugging Face, LangChain, OpenCV, FastAPI, Flask, Git, GitHub, SQL, AWS, Ollama, DialogFlow.
- **Development & Systems Experience:** Unix/Linux environments, distributed systems, information retrieval, TCP/IP networking, DBMS, Computer Networks, Google Colab, Jupyter, VS Code, PyCharm.
- **Programming Languages:** Proficient in Python and Java, with experience in C++ and C.

## EDUCATION

**AMITY UNIVERSITY**                                                                                                   **Patna, Bihar, India**
*Bachelor of Technology in Computer Science and Engineering*                                             **2022 - 2026**

- Appointed as Club Representative of AmiKoders under AIKYAM for contributions.
- Led an AI/ML Bootcamp, mentoring 100+ juniors in machine learning.

**DAV SCHOOL**                                                                                                   **Siliguri, West Bengal, India**
*Higher Secondary Education*                                                                                                   **2012 - 2022**

- Completed higher secondary education in PCM-CS, developing foundations in Physics, Chemistry, Mathematics, and Computer Science with practical applications.

## PROJECTS

**MAS-HIRE (MULTI-AGENT SYSTEM FOR HIRING)** [GitHub]
- This AI-driven job screening system comprises three agents built using the Gemma-2 model via Ollama. The first agent, leveraging Zero-Shot Learning, extracts key role-specific information from job descriptions. The second agent applies Few-Shot Learning to parse CVs, extracting entities like name, qualifications, and predicting top-2 suitable roles. The third agent compares job and CV data to compute match scores, enabling precise, automated shortlisting with candidate names, emails, and best-fit job roles.

**AI POWERED DOCUMENT ANALYSIS SYSTEM** [GitHub]
- An AI-powered document analysis system designed for researchers, featuring advanced plagiarism detection using TF-IDF, cosine similarity, n-gram, and LSH techniques. It classifies documents into financial, legal, or healthcare domains using a SpaCy NLP model trained on quality datasets. The system allows seamless PDF navigation and integrates a RAG-based QnA chatbot powered by LLaMA-3, all-MiniLM embeddings, and FAISS database.

**SHODH-AI: UNIFIED AI PLATFORM FOR INTELLIGENT CONTENT INTERACTION** [GitHub]
- ShodhAI is an AI-powered platform that transforms digital content, including videos, websites, PDFs, and images, into interactive learning and research experiences. It offers features like YouTube transcript analysis, real-time QnA for websites, PDF chatbots, document similarity checks, and high-precision OCR, streamlining content engagement and improving productivity for students, researchers, and professionals. Utilizing HuggingFace Transformers, Ollama Gemma2, PaddlePaddle OCR, LangChain, and Llama Models.

**YOUTUBE VIDEO TRANSCRIPT BASED CHATBOT** [GitHub]
- The YouTube Transcript ChatBot is an AI-powered tool that extracts YouTube video transcripts and allows users to ask context-aware questions using the LLAMA3 - 70B model via Groq Inference API. Built with Flask, it provides features like transcript downloads, chat history navigation, and intelligent responses, enhancing productivity for research and learning.

**DOCCUPY: AI-DRIVEN CHATBOT FOR SEAMLESS PDF INTERACTION** [GitHub]
- Doccupy is an AI-powered PDF chatbot that allows users to upload PDFs and interact with the document through an intelligent, context-aware chatbot. Utilizing Ollama models, FAISS vector database, and the LangChain framework, it provides precise, real-time answers, transforming static PDFs into dynamic, conversational learning tools.

**TEXTIFY OCR TOOLS**
- **Textify-GPU** is a high-performance OCR tool for handwritten text extraction, using transformer based stepfun-ai/GOT-OCR2_0 model and llama-3.3-70b-versatile for refined, accurate results. Users can copy the extracted text or download it in multiple formats. [GitHub]
- **Textify-CPU** is a high-performance OCR tool using PaddleOCR to extract text from scanned images, allowing users to copy or download results in PDF, DOCX, or TXT formats for convenience and flexibility. [GitHub]

**AI FACE DETECTION SYSTEM** [GitHub]
- This project implements an efficient face detection model in Python using OpenCV, leveraging a pre-trained Haar Cascade Classifier. It accurately detects faces in real-time via webcam feed, displaying bounding boxes around faces, with support for multiple detections.

## CERTIFICATIONS & ACHIEVEMENTS

- Intern of the Month Certificate from Purezza Technologies — JULY 2024
- Deep Learning Course Certificate from Scaler Academy — JAN 2024
- Machine Learning Summer Training Certificate from Analytics Vidhya — OCT 2023
- Data Science Internship Certificate from CodersCave — SEPT 2023
- Data Analytics Certificate from CISCO — JULY 2023
- Artificial Intelligence of Things Apprenticeship Certificate from CodroidHub Ltd. — JUN 2023