

# Wanderson Gomes de Souza

Senior AI Engineer | Python Developer | GenAI, RAG & Agents

✉ wandersonsouza.info@gmail.com ☎ +5583988579083 🌐 linkedin.com/in/wandersongsouza

🔗 wandersonsouza.com

## 📄 PROFILE

---

Senior AI Engineer with 10+ years of experience designing, developing, and deploying AI-driven applications across Generative AI, NLP, RAG, Multi-Agent Systems, Computer Vision, and production machine learning systems. PhD with a strong background in applied machine learning, experimental design, data pipelines, and scalable AI solutions. Hands-on experience building Python applications using FastAPI, LangChain, LangGraph, CrewAI, OpenAI APIs, embeddings, vector databases, semantic retrieval, prompt engineering, and context orchestration. Strong ability to integrate LLMs with enterprise systems, optimize retrieval pipelines, and deliver maintainable AI solutions in cloud and production environments, collaborating with technical and business teams to solve complex problems.

### Main Skills

Python Development · Generative AI · LLM Systems · RAG Architectures · Agentic Workflows · Multi-Agent Systems (CrewAI, LangGraph) · LangChain · OpenAI API · LlamaParse · Hugging Face Transformers · Prompt Engineering · Context Engineering · Vector Databases & Embeddings · FAISS · Semantic Search & Indexing · API Development (FastAPI, Flask) · SQL & NoSQL Databases · Data Pipelines · Cloud Platforms (Azure, AWS) · Docker · Kubernetes/AKS · Production AI Systems

## 📁 PROFESSIONAL EXPERIENCE

---

### Tech Lead / Senior AI Specialist

07/2025 – 05/2026

Motiva

- Acted as a hands-on technical leader, developing Python-based GenAI applications with RAG, semantic search, and agentic workflows.
- Built backend services and APIs with Python and FastAPI to integrate LLM components with enterprise systems and operational workflows.
- Engineered multi-agent workflows using LangGraph, enabling dynamic, context-driven reasoning and decision-making across complex AI pipelines.
- Implemented retrieval pipelines using embeddings, vector databases, and indexing strategies to improve relevance, grounding, and accuracy.
- Applied prompt engineering, context orchestration, and LLM evaluation strategies to improve consistency, reliability, and response quality.

### Senior Machine Learning Engineer

04/2024 – 07/2025

Venturus

- Developed and architected an embedded LLM-based assistant for Smart TV environments, enabling contextual interaction over technical documentation.
- Designed modular RAG pipelines using LlamaParse for document parsing, with ingestion, chunking, semantic retrieval, context assembly, and response generation.
- Built Python-based integrations connecting LLM components, retrieval pipelines, and backend services for real-world AI applications.
- Implemented context management strategies to improve response quality under latency, memory, and limited context-window constraints.
- Evaluated LLM outputs by analyzing grounding, consistency, relevance, and response reliability in RAG-based conversational systems.

## Senior Machine Learning Engineer

04/2023 – 04/2024

*iTriad*

- Developed production-grade Python services for real-time AI systems, supporting scalable deployment, backend integration, and inference workflows.
- Built FastAPI-based inference APIs to expose ML models through maintainable services for operational and real-time applications.
- Optimized inference pipelines, reducing latency by approximately 25% while preserving model performance and system reliability.
- Designed evaluation workflows using IoU, confusion matrices, threshold tuning, and error analysis to improve model robustness.
- Delivered computer vision models for head pose estimation, achieving  $\sim 9^\circ$  MAE and improving robustness under pose variation.

## Data Scientist

01/2021 – 04/2023

*NeoPTO*

- Developed large-scale NLP pipelines processing 11M+ patent documents, enabling text analysis, indexing, classification, and retrieval.
- Built Python and PySpark ETL workflows for ingestion, cleaning, normalization, and preparation of large unstructured text datasets.
- Designed retrieval-oriented data pipelines to structure patent titles, abstracts, claims, and descriptions for downstream NLP tasks.
- Worked with SQL and NoSQL databases to manage structured metadata and unstructured text used in scalable NLP and retrieval pipelines.

## Machine Learning Researcher

03/2020 – 01/2021

*GPICEEMA*

- Led development of a computer vision system automating analog densimeter readings, enabling accurate measurements in industrial environments.
- Designed validation protocols improving robustness under illumination changes and optical distortions.
- Integrated computer vision and OCR into an end-to-end pipeline, enabling automated monitoring and data extraction.

## Machine Learning Researcher

03/2015 – 03/2020

*CBTU*

- Led deployment of a real-time computer vision system for passenger counting and flow estimation on edge devices in production environments.
- Developed computer vision models robust to occlusion, perspective distortion and environmental variability.
- Built predictive models combining temporal and visual features to estimate passenger flow and forecast peak demand patterns.

## EDUCATION

---

### PhD in Mechanical Engineering

09/2019

*Federal University of Paraiba (UFPB)*

### MSc in Computer Science

01/2014

*Federal University of Paraiba (UFPB)*

### BSc in Computer Science

12/2009

*State University of Paraiba (UEPB)*

## CERTIFICATIONS AND ADDITIONAL COURSES

---

### AI Engineer Agentic Track: The Complete Agent & MCP Course

2025

*Udemy*

