

AI Engineer Roadmap 2026

LLM apps, RAG, multimodal AI, fine-tuning, evaluation, safety and production deployment roadmap

Goal: Build production AI applications that combine models, data, product UX, backend services, evaluation and responsible deployment.

Prepared: 31 May 2026. **Use-case:** self-study, portfolio building, interview prep and job-readiness.

How to use this roadmap

- Follow the phases in order, but do not wait to build projects. Every topic should end with a public artifact: code, dashboard, demo, README, diagram or case study.
- Use YouTube for intuition and demos; use official documentation for correctness; use practice platforms for repetition; use projects for proof.
- Track progress with a weekly scorecard: hours learned, problems solved, project commits, notes written, demos recorded and applications/interviews completed.
- AI tools are allowed, but every project must show that you understand, test and can debug the output.

2026 market calibration

- AI Engineer is now a broad applied role: less pure research, more productized GenAI features, RAG, automation, evaluation, cost optimization and safety.
- The 2026 market expects strong backend/software basics plus familiarity with LLM APIs, embeddings, vector databases, model selection, evals and cloud deployment.
- Portfolio projects must show measurable quality, citations, guardrails, latency/cost thinking and user experience, not just chatbot demos.

Role title mappings

- AI Engineer, Applied AI Engineer, GenAI Engineer, LLM Engineer, AI Application Engineer, AI Product Engineer, AI Solutions Engineer

2026 hiring-ready stack

- Python/TypeScript, FastAPI/Node, LLM APIs, RAG, embeddings and vector DBs, structured outputs, fine-tuning/PEFT basics, multimodal models, evaluation, safety and guardrails, Docker/cloud deploy, observability

End-to-end roadmap overview

Phase	Focus	What to learn	Job-ready proof
0-4 weeks	Software and AI basics	Python/TS, APIs, JSON schemas, async, LLM API calls, prompt basics, Git.	Build a CLI and small web app calling an LLM API.
1-2 months	RAG and embeddings	Chunking, embeddings, vector search, reranking, citations, metadata, retrieval eval.	Build a document Q&A app with cited answers.
2-4 months	Production LLM features	Structured outputs, tool calling, streaming, caching, auth, rate limits, cost tracking.	Add a reliable AI feature to a full-stack app.
4-6 months	Evaluation and safety	Golden datasets, LLM-as-judge, regression tests, prompt injection, PII, guardrails, red teaming.	Create an eval harness and safety checklist.
6-12 months	Advanced AI engineering	Fine-tuning/PEFT, multimodal, agents, deployment, monitoring, feedback loops, governance.	Ship 3 production-style AI apps with measured quality.

Weekly operating system

- 10-12 hours learning/building: 60% project implementation, 20% docs, 10% YouTube/course intuition, 10% notes and revision.
- Minimum weekly output: 5 commits, 1 written note, 1 demo screenshot/video, 5-20 practice problems depending on role.
- Every Sunday: review blockers, update README, write what you learned, plan next week, and compare against target job descriptions.
- Every month: ship one small project or one major milestone of a bigger capstone.

Complete topic and subtopic checklist

1. LLM API development

- Models, messages, parameters, streaming, structured outputs, retries, rate limits, caching, fallbacks, cost tracking, model routing.
- Proof to create: Build API wrappers with tests and observability.

2. Prompt and product design

- Task decomposition, examples, schemas, UX constraints, error states, refusal UX, prompt versioning, user feedback.
- Proof to create: Design prompts as product assets, not magic text.

3. RAG systems

- Document parsing, chunking, embeddings, vector DBs, metadata, hybrid search, reranking, citations, permissions, freshness.
- Proof to create: Compare retrieval strategies with an eval set.

4. Fine-tuning and customization

- When to fine-tune vs prompt/RAG, dataset creation, LoRA/PEFT basics, safety risks, eval, deployment constraints.
- Proof to create: Fine-tune a small model or classifier and benchmark it.

5. Multimodal AI

- Vision inputs, OCR/document understanding, audio transcription, image generation APIs, UI considerations, evaluation.
- Proof to create: Build a document/image analysis feature with quality checks.

6. Evaluation and quality

- Golden datasets, exact/semantic scoring, LLM judge calibration, adversarial tests, regression gates, user feedback loops.
- Proof to create: Treat evals like CI for AI behavior.

7. Safety, privacy and governance

- PII, prompt injection, data retention, model policy, guardrails, human review, audit logs, copyright, compliance basics.
- Proof to create: Create a risk register for every AI app.

8. Deployment and operations

- FastAPI/Node, Docker, serverless, queues, monitoring, tracing, latency, token costs, incident response, feature flags.
- Proof to create: Deploy with dashboards and rollback plan.

Portfolio projects and capstones

Do not treat these as toy tutorials. Each project should have a README, architecture diagram, setup steps, screenshots, demo link when possible, tests/checks, limitations and future improvements.

Project 1: RAG knowledge assistant

- Outcome: Upload docs, retrieve cited answers and track retrieval/answer metrics.
- Suggested stack: OpenAI/Anthropic/Gemini, vector DB, FastAPI, React
- Stretch goal: Add eval dashboard and permission-aware retrieval.

Project 2: AI document processor

- Outcome: Extract structured data from PDFs/images and validate results.
- Suggested stack: Multimodal LLM, OCR, Pydantic/Zod, queue
- Stretch goal: Add confidence thresholds and human review.

Project 3: AI writing/review feature

- Outcome: Generate, critique and improve a domain-specific artifact with style rules.
- Suggested stack: LLM API, evals, UI, feedback loop
- Stretch goal: Add regression tests for style and factuality.

Project 4: Fine-tuned classifier

- Outcome: Create labeled data and compare fine-tuned model vs prompt baseline.
- Suggested stack: Hugging Face, PEFT or managed fine-tuning
- Stretch goal: Add cost/latency/quality comparison.

Project 5: AI SaaS feature

- Outcome: Embed summarization, search or recommendation into a full-stack product.
- Suggested stack: Next.js, FastAPI/Node, DB, AI API
- Stretch goal: Add usage analytics, billing guardrails and safety filters.

Interview preparation checklist

- Build a small LLM feature with structured output and error handling.
- Design a RAG system and explain chunking, retrieval, permissions and evals.
- Discuss model choice, cost, latency, caching and fallbacks.
- Explain prompt injection and safety controls.
- Show an eval harness and how it prevents regressions.

Portfolio checklist before applying

- 3 AI apps with public demos, architecture diagrams and eval reports.
- One RAG project with citations and retrieval metrics.
- One multimodal or fine-tuning project.
- A clear safety/privacy section for each project.

Resume keywords to include when true

- Python/TypeScript, FastAPI/Node, LLM APIs, RAG, embeddings and vector DBs, structured outputs, fine-tuning/PEFT basics, multimodal models, evaluation, safety and guardrails, Docker/cloud deploy, observability, AI Engineer, Applied AI Engineer, GenAI Engineer

Best official docs and learning references

Prioritize these over random snippets. Read docs when implementation breaks, when preparing interviews, and when upgrading tools.

- **OpenAI API Docs** - LLM APIs, tools, agents and safety patterns. ([open](#))
- **OpenAI Agents SDK** - Agent workflows, tools, guardrails and tracing. ([open](#))
- **Anthropic Docs** - Claude API and tool-use patterns. ([open](#))
- **Google AI for Developers** - Gemini API and AI app development. ([open](#))
- **Microsoft Foundry Docs** - Azure AI Foundry agents, models, evaluation and monitoring. ([open](#))
- **AWS Bedrock Docs** - Managed foundation models and agents on AWS. ([open](#))
- **LangChain Docs** - LLM app building blocks. ([open](#))
- **LangGraph Docs** - Stateful, controllable agent graphs. ([open](#))
- **LlamaIndex Docs** - Data connectors and RAG frameworks. ([open](#))
- **Haystack Docs** - LLM pipelines and search. ([open](#))
- **Hugging Face Agents Course** - Free agent course. ([open](#))
- **Model Context Protocol Docs** - Tool/context protocol for AI applications. ([open](#))
- **Qdrant Docs** - Vector database. ([open](#))
- **Pinecone Docs** - Managed vector search. ([open](#))
- **Weaviate Docs** - Vector database. ([open](#))
- **pgvector** - Postgres vector similarity search. ([open](#))
- **Ragas Docs** - RAG evaluation. ([open](#))
- **DeepEval Docs** - LLM evaluation framework. ([open](#))
- **LangSmith Docs** - Tracing, testing and observability for LLM apps. ([open](#))
- **scikit-learn User Guide** - Classical ML algorithms and pipelines. ([open](#))
- **PyTorch Docs** - Deep learning framework. ([open](#))
- **TensorFlow Guides** - TensorFlow and Keras. ([open](#))
- **Hugging Face Transformers Docs** - Transformer models and pipelines. ([open](#))
- **Hugging Face Course** - Free AI courses and cookbooks. ([open](#))
- **XGBoost Docs** - Gradient boosting. ([open](#))
- **Node.js Learn** - Node runtime and backend basics. ([open](#))
- **Express Docs** - Minimal Node web framework. ([open](#))
- **NestJS Docs** - Structured TypeScript backend framework. ([open](#))
- **FastAPI Docs** - Modern Python API development. ([open](#))
- **Django Docs** - Python web framework. ([open](#))
- **PostgreSQL Docs** - Relational database reference. ([open](#))
- **MongoDB Docs** - Document database. ([open](#))
- **MDN Web Docs** - HTML, CSS, JavaScript and browser APIs. ([open](#))
- **web.dev** - Performance, accessibility and modern web guidance. ([open](#))
- **TypeScript Docs** - Typed JavaScript for production. ([open](#))
- **React Docs** - Modern React, hooks, Server Components and patterns. ([open](#))
- **Next.js Docs** - React full-stack framework, routing, data fetching, caching and deployment. ([open](#))
- **Docker Docs** - Containers and images. ([open](#))
- **Kubernetes Docs** - Container orchestration. ([open](#))
- **Helm Docs** - Kubernetes packaging. ([open](#))
- **Terraform Docs** - Infrastructure as code. ([open](#))

- **Ansible Docs** - Configuration automation. ([open](#))
- **AWS Docs** - AWS services. ([open](#))
- **Google Cloud Docs** - Google Cloud services. ([open](#))
- **Azure Docs** - Azure services. ([open](#))
- **OpenAI Cookbook** - Practical AI app examples. ([open](#))
- **Hugging Face Cookbook** - Practical model/RAG recipes. ([open](#))
- **Prompt Engineering Guide** - Prompt patterns and examples. ([open](#))
- **OWASP Top 10 for LLM Applications** - LLM security risks. ([open](#))

Best YouTube sources and how to use them

Use these for intuition, project walkthroughs and high-level context. Always verify implementation details against official documentation because tools change quickly.

- **freeCodeCamp.org** - Long-form free courses across web, data, cloud and AI. ([open](#))
- **ByteByteGo** - System design, architecture and backend concepts. ([open](#))
- **MIT OpenCourseWare** - Computer science, math and engineering fundamentals. ([open](#))
- **Stanford Online** - AI, ML and CS lecture material. ([open](#))
- **OpenAI** - OpenAI product and developer updates. ([open](#))
- **LangChain** - Agent frameworks, LangGraph, LangSmith and RAG. ([open](#))
- **Anthropic** - Claude and AI safety/product updates. ([open](#))
- **Hugging Face** - Open models, agents and deployment. ([open](#))
- **DeepLearning.AI** - Short courses on LLMs, agents, RAG and evaluation. ([open](#))
- **James Briggs** - Vector search, RAG and LLM engineering. ([open](#))
- **AI Jason** - No-code/low-code and agent workflows. ([open](#))
- **AssemblyAI** - LLM app tutorials and AI engineering concepts. ([open](#))
- **StatQuest** - Statistics and ML explained clearly. ([open](#))
- **3Blue1Brown** - Math intuition for linear algebra and calculus. ([open](#))
- **Andrej Karpathy** - Neural networks, LLMs and hands-on AI. ([open](#))
- **Krish Naik** - ML, data science and GenAI projects. ([open](#))
- **sentdex** - Python, ML and automation projects. ([open](#))
- **ArjanCodes** - Python software engineering for production-quality code. ([open](#))
- **Fireship** - Fast overviews of web tools and trends. ([open](#))
- **Web Dev Simplified** - Practical JavaScript, React and web explanations. ([open](#))
- **Traversy Media** - Project-based web development. ([open](#))
- **JavaScript Mastery** - Full-stack React/Next projects. ([open](#))

Recommended YouTube search strategy

- **YouTube search: AI Engineer roadmap 2026 project based** - Use filters for upload date and length; prefer recent, project-based content. ([open](#))
- **YouTube search: AI Engineer interview preparation 2026** - Use filters for upload date and length; prefer recent, project-based content. ([open](#))
- **YouTube search: AI Engineer portfolio project end to end** - Use filters for upload date and length; prefer recent, project-based content. ([open](#))
- **YouTube search: AI Engineer system design real world project** - Use filters for upload date and length; prefer recent, project-based content. ([open](#))

Best coding and practice platforms

Practice platforms build repetition. They do not replace projects. For each platform, maintain a solved-problems log with mistakes and revised patterns.

- **Kaggle Competitions** - ML competitions and notebooks. ([open](#))
- **DrivenData** - Social-impact ML competitions. ([open](#))
- **Papers with Code** - Benchmarks and implementations. ([open](#))
- **Hugging Face Spaces** - Deploy AI demos. ([open](#))
- **Hugging Face Agents Course** - Agent practice and benchmark tasks. ([open](#))
- **OpenAI Cookbook** - Practical LLM examples. ([open](#))
- **MLOps Zoomcamp** - End-to-end MLOps projects. ([open](#))
- **LeetCode** - DSA patterns for interviews. ([open](#))
- **NeetCode Roadmap** - Structured DSA practice. ([open](#))
- **HackerRank** - Algorithms, SQL and language tracks. ([open](#))
- **Exercism** - Language practice with mentoring. ([open](#))
- **Codewars** - Kata practice. ([open](#))
- **GitHub Skills** - Hands-on GitHub learning. ([open](#))
- **Frontend Mentor** - Pixel-to-code projects. ([open](#))
- **GreatFrontEnd** - Frontend interview practice. ([open](#))
- **BigFrontend.dev** - JS, CSS and systems problems. ([open](#))
- **DataLemur SQL Questions** - Data interview SQL. ([open](#))
- **StrataScratch** - SQL, Python and analytics interview questions. ([open](#))
- **SQLBolt** - Beginner SQL drills. ([open](#))
- **KodeKloud Free Labs** - Linux, Docker, Kubernetes and DevOps labs. ([open](#))
- **Killercoda** - Kubernetes, Linux and cloud-native labs. ([open](#))
- **Cloud Resume Challenge** - Cloud portfolio challenge. ([open](#))

30-60-90 day execution plan

Period	Primary objective	Deliverables
Days 1-30	Foundation and first small project	Finish basics, solve starter practice, publish one small project and notes.
Days 31-60	Core role skills and second project	Complete core docs, build a stronger project, add tests/checks and write a case study.
Days 61-90	Production proof and interview prep	Deploy capstone, create architecture diagram, record demo, start mock interviews and applications.

Common mistakes to avoid

- Only watching videos without building public proof.
- Building tutorial clones without explaining decisions, tradeoffs and failure cases.
- Ignoring documentation, testing, security, accessibility, observability or data quality.
- Using AI-generated code without understanding it or adding tests.
- Applying to jobs before your portfolio proves the core responsibilities of the role.

Market research sources consulted

- **World Economic Forum Future of Jobs 2025** - Macro skill demand for 2025-2030; big data, AI/ML, software, cybersecurity and tech literacy signals. ([open](#))
- **Stack Overflow Developer Survey 2025** - Developer technology, AI-tool use, documentation, languages and web framework signals. ([open](#))
- **GitHub Octoverse 2025** - Open-source and language momentum; TypeScript, Python, agents and typed-language shift. ([open](#))
- **DORA State of AI-assisted Software Development 2025** - AI in software delivery; verification, platform maturity and organizational system effects. ([open](#))
- **CNCF Annual Cloud Native Survey** - Kubernetes, cloud-native maturity and AI infrastructure adoption. ([open](#))
- **dbt State of Analytics Engineering 2025** - Analytics engineering trends: AI use, data trust, quality and investment. ([open](#))
- **LangChain State of Agent Engineering** - Agent adoption, production challenges, evaluation and observability trends. ([open](#))
- **McKinsey State of AI 2025** - Enterprise AI adoption, scaling practices, governance and value capture. ([open](#))
- **McKinsey State of AI 2025** - Enterprise AI adoption and scaling practices. ([open](#))
- **OpenAI API Docs** - AI application development documentation. ([open](#))