

# หลักสูตร "พัฒนา AI/LLMs App แบบ On-Premise สำหรับใช้ในองค์กร"

วันที่ 29-30 มิ.ย. และ 1-3 ก.ค. 2569 เวลา 09.00-16.00 น.

(5 วันเต็ม)

## วัตถุประสงค์ของหลักสูตร

- เพื่อให้ผู้เรียนเข้าใจแนวคิด ประโยชน์ และการวางแผนพัฒนา AI/LLMs Application แบบ On-Premise สำหรับองค์กร รวมถึงการเลือกใช้ LLMs Open Source และเครื่องมือแบบ Local เช่น Ollama ให้เหมาะสมกับงานและทรัพยากร
- เพื่อให้ผู้เรียนพัฒนา AI Chat Bot, AI Agent, Multi-Agent System และ Workflow ด้วย LangChain.js ได้อย่างมืออาชีพ
- เพื่อให้ผู้เรียนใช้งาน Tool Calling และสร้าง MCP (Model Context Protocol) Server เพื่อเชื่อมต่อระบบภายในองค์กร (RESTful API, Database) ได้อย่างปลอดภัยและมีประสิทธิภาพ
- เพื่อให้ผู้เรียนสร้างระบบ RAG และ Agentic RAG ที่รองรับเอกสารหลายรูปแบบ (Text, Docx, PDF, Image) พร้อมจัดการ Vector Database สำหรับการค้นหาข้อมูลภายในองค์กร
- เพื่อให้ผู้เรียนได้เรียนรู้หลักการ Prompt Engineering, Context Engineering, Harness Engineering เพื่อควบคุม LLM ให้เหมาะสมกับบริบทองค์กร
- เพื่อให้ผู้เรียนได้เรียนรู้หลักการของ Agent Skills และนำมาใช้กับ LLM Application ขององค์กรได้
- เพื่อให้ผู้เรียนสามารถสร้าง Chat UI ด้วย Next.js พร้อมรองรับ Real-time Streaming, Authentication, และ Role-Based Access Control (RBAC)
- เพื่อให้ผู้เรียนเข้าใจและป้องกันความเสี่ยงด้าน AI Security รวมถึงการสร้าง Guardrails (PII Detection, Prompt Injection Prevention เป็นต้น)
- เพื่อให้ผู้เรียนเรียนรู้การใช้เครื่องมือสำหรับจัดการค่าใช้จ่าย และเพิ่มประสิทธิภาพระบบด้วยเทคนิค Cost Optimization, Semantic Caching, Fallback Models และ Prompt Compression เป็นต้น
- เพื่อให้ผู้เรียนจัดการ Error Handling, Retry Logic Pattern และ Fallback Strategy เพื่อรักษาความเสถียรของระบบ AI ใน Production
- เพื่อให้ผู้เรียน Deploy AI Application แบบ Production-grade ด้วย Docker, vLLM, AI Gateway (LiteLLM Proxy) และเซิร์ฟเวอร์ในองค์กร
- เพื่อให้ผู้เรียนติดตั้ง Observability Stack (Langfuse, Monitoring, Logging) สำหรับติดตามผลประเมินประสิทธิภาพ และทำ A/B Testing
- เพื่อให้ผู้เรียนสามารถประยุกต์ใช้ ถ่ายทอดองค์ความรู้ และสร้าง Use Case ในองค์กรตนเอง เช่น HR Bot, IT Support Bot, Policy QA Bot ได้อย่างเป็นรูปธรรม

## พื้นฐานผู้เรียน

- เคยเขียนโปรแกรมด้วยภาษาใดภาษาหนึ่งมาก่อน เช่น JavaScript เป็นต้น

## เนื้อหาการเรียนและตารางเรียน 5 วัน

เวลา	เนื้อหา
วันที่ 1	
09.00 – 10.30 น.	<b>LLMs &amp; On-Premise AI Strategy</b> <ul style="list-style-type: none"><li>ติดตั้ง และตั้งค่าโปรแกรมต่างๆ ที่เกี่ยวข้อง</li><li>แนวคิดของ AI แบบ Cloud</li><li>แนวคิดและประโยชน์ของการติดตั้ง AI/LLM แบบ On-Premise</li><li>LLM คืออะไร</li><li>สถาปัตยกรรม LLM และประเภทของโมเดล</li><li>หลักการ และแนวคิดการเลือกโมเดล</li><li>การเลือกขนาดโมเดล: 7B, 13B, 70B – ข้อดีข้อเสียของแต่ละขนาด</li><li>Quantization (การบีบอัดโมเดล) - GGUF, GPTQ, AWQ</li><li>ความต้องการฮาร์ดแวร์และการวางแผนค่าใช้จ่าย</li><li>ข้อดี ข้อเสีย ของ Cloud vs. On-Premise</li></ul>
10.30 – 10.45 น.	<b>พักเบรก</b>
10.45-12.00 น.	<b>Setup Local LLM Environment</b> <ul style="list-style-type: none"><li>การติดตั้ง Ollama และจัดการโมเดลต่างๆ</li><li>การใช้งาน Ollama</li><li>การตั้งค่าโมเดลและปรับแต่งประสิทธิภาพ</li><li>โมเดลที่ทำงานกับภาษาไทยได้ดี</li><li>ติดตั้งและทดสอบโมเดลหลายตัว และการใช้งาน Chat UI สำเร็จรูป</li></ul> <b>การติดตั้ง Next.js เพื่อเตรียมสร้าง AI Chatbot</b> <ul style="list-style-type: none"><li>การติดตั้ง Next.js</li><li>การสร้างโปรเจกต์ใหม่</li><li>พื้นฐานการเขียน Next.js</li><li>การใช้งาน UI สำเร็จรูป</li><li>การติดตั้งฐานข้อมูล และการเชื่อมต่อ</li><li>ทดลองสร้างหน้า dashboard และ ระบบ Authentication</li></ul>
12.00 – 13.00 น.	<b>พักรับประทานอาหาร</b>
13.00 – 14.30 น.	<b>AI Agent Fundamentals</b> <ul style="list-style-type: none"><li>การใช้งาน LLM ประเภทต่างๆ และแนวปฏิบัติที่ดี</li></ul>

เวลา	เนื้อหา
	<ul style="list-style-type: none"> <li>● สถาปัตยกรรม Agent</li> <li>● ประเภทของ Agent</li> <li>● แนะนำ Tool Calling และ Function Calling (การเรียกใช้ฟังก์ชัน)</li> <li>● ทำอย่างไรให้ LLM ตอบคำถามขององค์กร</li> </ul>
14.30 – 14.45 น.	<b>พักเบรก</b>
14.45 – 16.00 น.	<b>LangChain.js Core Concepts</b> <ul style="list-style-type: none"> <li>● แนะนำ และติดตั้ง LangChain.js</li> <li>● สถาปัตยกรรม และ Core component</li> <li>● ทำความรู้จักกับ Agents, Models, Messages, Tools, Short-term memory</li> <li>● การเปรียบเทียบ Chat Models กับ LLMs ทั่วไป</li> <li>● ทดลองสร้าง AI Agent อย่างง่าย</li> </ul>
<b>วันที่ 2</b>	
09.00 – 10.30 น.	<b>Next.js Chat UI Development</b> <ul style="list-style-type: none"> <li>● แนะนำการสร้าง UI</li> <li>● Real-time Streaming (การส่งข้อมูลแบบเรียลไทม์) ด้วย AI SDK</li> <li>● สร้างคอมโพเนนต์หน้า Chat</li> <li>● แนะนำการจัดการ Message History และการจัดการ Session</li> <li>● ทดลองเชื่อมต่อกับ Langchain</li> </ul>
10.30 – 10.45 น.	<b>พักเบรก</b>
10.45 -12.00 น.	<b>Prompt Engineering &amp; Context Engineering</b> <ul style="list-style-type: none"> <li>● รูปแบบการออกแบบและเขียน Prompt แบบต่างๆ เช่น (Zero-shot, Few-shot, Chain-of-Thought)</li> <li>● Context Engineering สำหรับองค์กร</li> <li>● Context Engineering Best Practices</li> <li>● เทคนิคการจัดการ Context อย่างมีประสิทธิภาพ</li> <li>● แนะนำ Harness Engineering</li> <li>● การสร้าง dynamic system prompt ใน Langchain</li> <li>● การเขียน Tools ด้วย Langchain</li> <li>● ทดลองสร้าง AI Chatbot ที่สามารถเรียกใช้ Tools ได้</li> </ul>
12.00 – 13.00 น.	<b>พักรับประทานอาหาร</b>
13.00 – 14.30 น.	<b>แนวคิดของ MCP</b> <ul style="list-style-type: none"> <li>● ภาพรวม MCP (Model Context Protocol)</li> <li>● ความสำคัญของ MCP</li> <li>● แนะนำการสร้าง MCP Server เพื่อเชื่อมต่อ API ภายในองค์กร</li> <li>● การเชื่อมต่อฐานข้อมูล (PostgreSQL)</li> </ul>

เวลา	เนื้อหา
	<ul style="list-style-type: none"> <li>● การเชื่อมต่อ Authentication</li> <li>● การเรียก Tool จาก MCP</li> <li>● ทดลองสร้าง MCP Server สำหรับระบบองค์กร</li> <li>● ทดลองให้ AI Chatbot เรียก Tool จาก MCP</li> </ul>
14.30 – 14.45 น.	<b>พักเบรก</b>
14.45 – 16.00 น.	<b>AI Security &amp; Guardrails</b> <ul style="list-style-type: none"> <li>● หลักการของ Middleware ใน Langchain</li> <li>● การใช้งาน Middleware ใน Langchain</li> <li>● การป้องกัน Prompt Injection (การแทรก Prompt ที่เป็นอันตราย)</li> <li>● PII Detection (ตรวจจับและปิดบังข้อมูลส่วนบุคคล)</li> <li>● แนะนำหลักการ และการใช้งาน Guardrails</li> <li>● ทดลองสร้างระบบ Guardrails ใน Langchain</li> </ul>
<b>วันที่ 3</b>	
09.00 – 10.30 น.	<b>RAG Architecture &amp; Vector Databases</b> <ul style="list-style-type: none"> <li>● แนวคิด RAG และการป้องกันการสร้างข้อมูลเท็จ</li> <li>● การเลือก LLM ประเภท Embeddings</li> <li>● Embeddings (การแปลงข้อความเป็นตัวเลข) และความคล้ายคลึงของ Vector</li> <li>● ตัวเลือก Vector Database: Qdrant, Chroma, PGVector, Supabase Vector</li> <li>● การติดตั้ง Qdrant (Docker) และตั้งค่า</li> <li>● แนะนำ Retrieval Pipeline</li> </ul>
10.30 – 10.45 น.	<b>พักเบรก</b>
10.45-12.00 น.	<b>Production RAG Implementation</b> <ul style="list-style-type: none"> <li>● Document Loaders (การโหลดเอกสาร) - PDF, DOCX, TXT, CSV</li> <li>● กลยุทธ์การแบ่งข้อความ (Chunking Strategies)</li> <li>● Metadata Filtering และ Hybrid Search (การค้นหาแบบผสม)</li> <li>● การเพิ่มประสิทธิภาพการค้นหา (MMR, Threshold, Reranking)</li> <li>● ทดลองสร้าง จัดเก็บเอกสารและค้นหาข้อมูลใน Vector Database</li> <li>● ทดลองให้ AI Chatbot สามารถตอบคำถามจากระบบ RAG</li> </ul>
12.00 – 13.00 น.	<b>พักรับประทานอาหาร</b>
13.00 – 14.30 น.	<b>Agentic RAG แบบ Multimodal</b> <ul style="list-style-type: none"> <li>● สถาปัตยกรรม Agentic RAG</li> <li>● การทำความเข้าใจภาพ (OCR, Vision Models)</li> <li>● การเลือก Vision Model มาใช้งานอย่างเหมาะสม</li> <li>● การดึงข้อมูลจากภาพ แปลงเป็นข้อความ</li> <li>● Pipeline การประมวลผลเอกสารที่ซับซ้อน</li> </ul>

เวลา	เนื้อหา
	<ul style="list-style-type: none"> <li>ทดลองสร้าง AI Agent Bot ที่รองรับรูปภาพ</li> </ul>
14.30 – 14.45 น.	<b>พักเบรก</b>
14.45 – 16.00 น.	<b>Short-term memory &amp; Long-term memory ใน LangChain</b> <ul style="list-style-type: none"> <li>short-term memory คืออะไร</li> <li>ทดลองใช้งาน short-term memory</li> <li>ทดลองใช้ Runtime ใน LangChain</li> <li>การใช้งาน Long-term memory</li> </ul>
<b>วันที่ 4</b>	
09.00 – 10.30 น.	<b>Workshop: Short-term memory &amp; Long-term memory</b> <ul style="list-style-type: none"> <li>ทดลองสร้าง Chat History เพื่อให้ AI จำจดบทสนทนา</li> <li>ทดลองทำ Workshop สำหรับการตั้งค่า user preference เพื่อให้ระบบจดจำสไตล์และตัวตนของผู้ใช้งาน</li> </ul>
10.30 – 10.45 น.	<b>พักเบรก</b>
10.45-12.00 น.	<b>ทำความรู้จักกับ Multi-Agent Systems</b> <ul style="list-style-type: none"> <li>แนวคิด Multi-Agent และ LangChain.js</li> <li>ทำไมและเมื่อไหร่ต้องใช้ Muti Agent</li> <li>การเลือก Pattern แบบต่างๆ</li> <li>การสื่อสารระหว่าง Agent และจัดการสถานะ</li> <li>ทดลองให้ AI Agent มี Agent เพิ่ม ด้วย Subagents</li> <li>ทดลองสร้าง และใช้ Skill กับ AI Agent</li> </ul>
12.00 – 13.00 น.	<b>พักรับประทานอาหาร</b>
13.00 – 14.30 น.	<b>Workshop: Multi-Agent Systems</b> <ul style="list-style-type: none"> <li>ทดลองให้ AI Agent มี Agent เพิ่ม ด้วย Subagents</li> <li>ทดลองสร้าง และใช้ Skill กับ AI Agent</li> </ul>
14.30 – 14.45 น.	<b>พักเบรก</b>
14.45 – 16.00 น.	<b>แนะนำ vLLM</b> <ul style="list-style-type: none"> <li>แนะนำ vLLM</li> <li>พื้นฐานของ Inference Server</li> <li>กลไกการทำงานของ vLLM</li> <li>คุณสมบัติและข้อได้เปรียบของ vLLM</li> <li>ข้อแตกต่างระหว่าง vLLM กับ Ollama</li> <li>ทดลองใช้ vLLM Playground</li> </ul>
<b>วันที่ 5</b>	
09.00 – 10.30 น.	<b>แนะนำ LiteLLM</b> <ul style="list-style-type: none"> <li>แนะนำ LiteLLM</li> </ul>

เวลา	เนื้อหา
	<ul style="list-style-type: none"> <li>● การติดตั้ง LiteLLM ด้วย Docker</li> <li>● LiteLLM แก้ปัญหาอะไรบ้าง</li> <li>● LiteLLM Proxy Server คืออะไร</li> <li>● พีเจอร์เสริมและการใช้งาน</li> <li>● ทดลองติดตั้ง และสร้าง Proxy Server</li> </ul>
10.30 – 10.45 น.	<b>พักเบรก</b>
10.45-12.00 น.	<b>Production Deployment Strategy</b> <ul style="list-style-type: none"> <li>● การตั้งค่า Docker และ Docker Compose</li> <li>● Deploy vLLM สำหรับ High-Performance Inference (ประมวลผลเร็ว)</li> <li>● AI Gateway ด้วย LiteLLM Proxy</li> <li>● Load Balancing และ Rate Limiting (การกระจายโหลดและจำกัดอัตรา)</li> </ul>
12.00 – 13.00 น.	<b>พักรับประทานอาหาร</b>
13.00 – 14.30 น.	<b>Observability &amp; Monitoring</b> <ul style="list-style-type: none"> <li>● แนะนำ Langfuse</li> <li>● การเชื่อมต่อ Langfuse กับ LiteLLM สำหรับ Tracing (ติดตามการทำงาน)</li> <li>● กลยุทธ์การทำ Logging (บันทึกข้อมูล)</li> <li>● Metrics และ Dashboards (แดชบอร์ดแสดงผล)</li> <li>● Latency, Token, Error Rate</li> </ul>
14.30 – 14.45 น.	<b>พักเบรก</b>
14.45 – 16.00 น.	<b>แนะนำการนำ CI/CD เข้ามาย่างงาน</b> <ul style="list-style-type: none"> <li>● แนะนำ CI/CD Pipeline สำหรับแอป AI</li> <li>● Code Review และแนวปฏิบัติที่ดี</li> <li>● ถาม-ตอบ และปิดหลักสูตร</li> </ul>

- ติดต่อ: [codingthailand@gmail.com](mailto:codingthailand@gmail.com) หรือโทร: 085-4952624 อ.เอก