# Pretraining and Finetuning LLMs
# from the Ground Up

| | **Workshop topics** |
|---|---|
| 1 | Introduction to LLMs |
| 2 | Understanding LLM input data |
| 3 | Coding an LLM architecture |
| 4 | Pretraining LLMs |
| 5 | Loading pretrained weights |
| 6 | Finetuning LLMs |

# Developing an LLM

https://mng.bz/lrp2

https://github.com/rasbt/LLMs-from-scratch

**(Source for most figures and code)**

# Developing an LLM

**STAGE 1: BUILDING**

1) Data preparation & sampling

2) Attention mechanism

3) LLM architecture

Building an LLM

4) Pretraining →

**STAGE 2: PRETRAINING**

5) Training loop

6) Model evaluation

7) Load pretrained weights

Foundation model

8) Finetuning

**STAGE 3: FINETUNING**

Dataset with class labels

Classifier

9) Finetuning

Personal assistant

Instruction dataset

**WITH LITGPT**

## ⚡ LitGPT

**20+ high-performance LLM implementations with recipes to pretrain, finetune, deploy at scale.**

✅ From scratch implementations   ✅ No abstractions        ✅ Beginner friendly
✅ Flash attention                ✅ FSDP                    ✅ LoRA, QLoRA, Adapter
✅ Reduce GPU memory (fp4/8/16/32) ✅ 1–1000+ GPUs/TPUs      ✅ 20+ LLMs

python 3.8 | 3.9 | 3.10 | 3.11    CPU tests passing    License Apache 2.0    chat 988 online

Lightning AI · Quick start · Models · Finetune · Deploy · All workflows · Features · Recipes (YAML) · Tutorials

Get started

https://github.com/Lightning-AI/litgpt

Lightning AI

Home   Studio templates   Agents   Teamspaces   Community   Docs

**Source**

Lightning AI   Public

**Explore**

⭐ Featured
📈 Trending
🕐 Recent
▦ All studios
👤 My studios

**Educational**

📄 Blogs
📄 Papers
📄 Tutorials

**Workflows**

▦ Data processing
📡 Endpoints
⏱ Training
📶 Serving
◯ Other

**Model types**

🎤 Audio
🖼 Image
✦ Multimodal
💬 Text
▦ Tabular

---

RAG 102
Chat with Documents
★ Featured

**Document Chat Assistant using RAG**
aniket   🚀 269   👁 6.54 K

---

Improve LLMs via Proxy-Tuning
★ Featured

**Improve LLMs With Proxy-Tuning**
sebastian   🚀 47   👁 7.88 K

---

Embed Wikipedia English under 5 dollars
★ Featured

**Embed English Wikipedia under 5 dollars**
thomasgridai   🚀 26   👁 2.73 K

---

Finetune Hugging Face BERT with PyTorch Lightning
★ Featured

**Finetune Hugging Face BERT with PyTorch Lig...**
justin   🚀 97   👁 1.98 K

---

Ingest documents (text, pdf, markdown, docx) in a vector database for Retrieval Augmented Generation (RAG)

**Document Search and Retrieval using RAG**
aniket   🚀 676   👁 7.10 K

---

Data streaming benchmarks for ImageNet
★ Featured

**Benchmark cloud data-loading libraries**
thomasgridai   🚀 23   👁 1.05 K

---

SlimPajama & Starcoder   1 trillion tokens
★ Featured

**Prepare the TinyLlama 1T token dataset**
thomasgridai   🚀 38   👁 1.64 K

---

LoRA from Scratch
★ Featured

**Code LoRA from Scratch**
sebastian   🚀 229   👁 24.66 K

---

Optimized Inference API for Mistral 7B with vLLM
★ Featured

**Optimized LLM inference API for Mistral 7B usi...**
aniket   🚀 50   👁 7.63 K

# Contact

@rasbt    in/sebastianraschka

https://sebastianraschka.com/contact/

https://lightning.ai