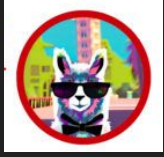# Open-Source Large Language Models in Radiology:
A Review and Tutorial for Practical Research and Clinical Deployment
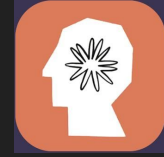
Savage et al. Radiology 2025

*A review*

By Zineb El yamani
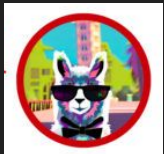August 11, 2025

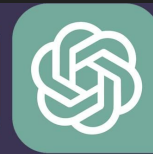# Open source  VS  Proprietary

Performance

Customizability

Cost

Licensing /innovation

Data security

Safety

# Open source  VS  Proprietary

Performance  **P**

# Open source VS Proprietary

Customizability

**Open-source:**
   you own the model

**Proprietary:**
   Always under platform control
   Risk of workflow disruption if a model version is deprecated

# Open source  VS  Proprietary

Cost

**Open-source:**

    Smaller models

    No API costs (only hardware and maintenance expenses)

**Proprietary:**

    Large parameter counts

    Pay-per-token for every API request

# Open source VS Proprietary

**Open-source:**

Encourages institutional and entrepreneurial innovation.

**Proprietary:**

- Some platforms allow revenue sharing, but always within the company's ecosystem.

Licensing /innovation

**Open-source**:

avoids sending sensitive data to third parties.
Security depends on local IT infrastructure, which may be weaker than big cloud providers.

**Proprietary**

Cloud vendors often have advanced, **HIPAA-compliant** security.

Data security

# **O**pen source **VS** **P**roprietary

**Open-source**

Less consistent testing for harmful outputs.

Developers often have fewer resources for adversarial robustness testing (prompt injection defenses).

**Proprietary**

Large-scale safety evaluation (red teaming, multiple safeguard layers).

Better resilience to adversarial prompts.

Safety

**P**

# I am a radiologist, how can I implement an open source LLM ?

Step 1:

Choose a Model

🤗 **Hugging Face**

😆 **Hugging Face**

Step 1:

Choose a Model

Metrics do not *always* align with radiology specific needs

| Open LLM | Chatbot Arena | OpenCompass |
|---|---|---|
| Open-source LLMs only | Proprietary and open-source LLMs | Proprietary and open-source LLMs |

**Performance on Benchmarks**

Performance can be approximated using public leaderboards

Further are needed !!

-complex multistep instructions
-Use case specific metrics
-Depend whether we consider
end point to be fully automated
or collaborative

**RaLEs**

**PubMedQA**
A Dataset for Biomedical Research Question Answering

🤗 **Hugging Face**

Step 1:

Choose a Model

Metrics do not *always* align with
radiology specific needs

| **Open LLM** Open-source LLMs only | **Chatbot Arena** Proprietary and open-source LLMs | **OpenCompass** Proprietary and open-source LLMs |

**Performance on Benchmarks**

Performance can be approximated
using public leaderboards

Step 1:

Choose a Model

| Open LLM | Chatbot Arena | OpenCompass |
|---|---|---|
| Open-source LLMs only | Proprietary and open-source LLMs | Proprietary and open-source LLMs |

**Performance on Benchmarks**

Performance can be approximated using public leaderboards

7B
13B
35B
70B

**Parameter Count**
↑Count = ↑Memory

**Model Size**

Memory requirements can be approximated by parameter count and reduced with quantization

**Quantized Models**
↓ Bit size
↓ Memory
↓ Performance

| exl2 | Any |
| AWQ | 4-bit |
| GPTQ | 4-bit | 8-bit |
| GGUF | 2-bit | 4-bit | 6-bit | 8-bit |

AWQ = Activation-aware Weight Quantization,
GGUF = GPT-Generated Unified Format,
GPTQ = Post-Training Quantization

Compressing LLMs without sacrificing performance

Step 1:

Choose a Model

**Parameter Count**
↑Count = ↑Memory

7B
13B
35B
70B

**Model Size**

Memory requirements can be approximated by parameter count and reduced with quantization

**Quantized Models**
↓ Bit size
↓ Memory
↓ Performance

exl2 — Any
AWQ — 4-bit
GPTQ — 4-bit — 8-bit
GGUF — 2-bit — 4-bit — 6-bit — 8-bit

**Open LLM**
Open-source LLMs only

**Chatbot Arena**
Proprietary and open-source LLMs

**OpenCompass**
Proprietary and open-source LLMs

**Performance on Benchmarks**

Performance can be approximated using public leaderboards

**Model Type**

Defined by how the training data was organized

**Instruct**
Predominantly trained to follow user instructions

**Chat**
Trained to behave like an assistant (e.g. chatbot)

**Chat-Instruct**
Relatively more training devoted to having conversations

**Pretrained**
"Base" model that has not been trained to perform specific tasks or actions

MoE ?
Merged models ?

Step 1:

Choose a Model

LLM downloaded

Step 2:

Choose LLM Generator Platform and Model Loader

**LLM Generation Platforms**

Community-built tools streamline implementing, inferencing, and modifying generation settings of LLMs

vLLM

Text Generation Web UI

SillyTavern

Or local code

**Model Loaders**

Different loaders are used depending on quantization and file type

**Hugging Face's Transformers Library**
Loads full 16-bit models

**Full-size Models**

**Quantized Models**

**llama.cpp**
Loads GGUF models
Can use GPU and CPU memory

**ExLlamav2**
Loads exl2 and GPTQ models.
GPU memory only

**AutoAWQ**
Loads AWQ models
GPU memory only

```
Step 1:

Choose a Model
```

```
Step 2:

Choose LLM Generator
Platform and Model
Loader
```

```
*Step 3:

Deploy the LLM
```

**Dataset Processing**

Iteratively process
information from a dataset

**Chatbot**

Answer questions and have
conversations with users

**Local API**
(e.g. Text Generation Web UI)

**Front-End Chatbot Interface**
(e.g. SillyTavern)

# Troubleshooting Performance issues

Prompt engineering

Retrieval-augmented Generation

Fine-tuning

**Problem :**
LLMs may exhibit deficiencies in complex reasoning, defined as low performance on tasks that require **multistep reasoning** (eg, generating a differential diagnosis)

**Common techniques:**

➔ <u>Chain-of-Thought (CoT)</u>: Tell the model to "think step by step," so it breaks reasoning into intermediate steps.
➔ <u>Reflexion</u>: The model simulates an *evaluator* that critiques its own first answer, then revises it based on that feedback.
➔ <u>Few-shot prompting:</u> Give a few solved examples in the prompt before the real question.

**Retrieval-augmented Generation**

**Problem :**

LLMs can have an **insufficient** knowledge base that can potentially lead to hallucinations. (eg, constantly changing medical guidelines)

**Solution:**

➔    Supplement the input prompt with information from other data source without the need for fine-tuning.
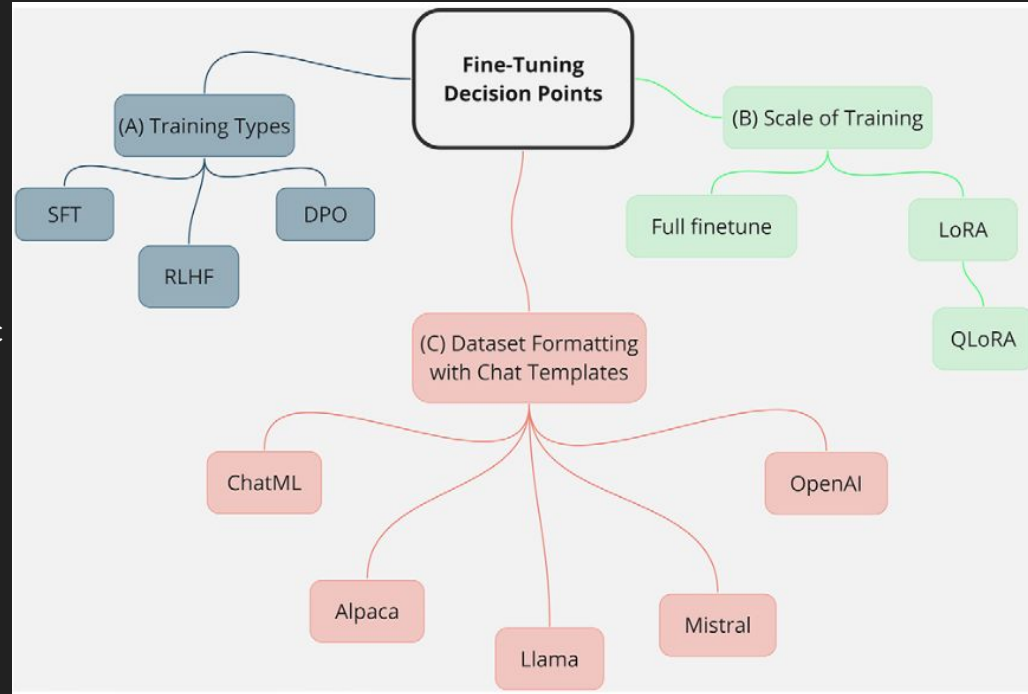
# Fine-tuning

**Problem :**
LLMs can exhibit **poor performance** in instruction following.

**Solution:**

➔ Retrain the model with additional domain-specific data so it internalizes new knowledge or skills.

## Fine-tuning

Training methods: *depend on the complexity and breadth of the desired task*

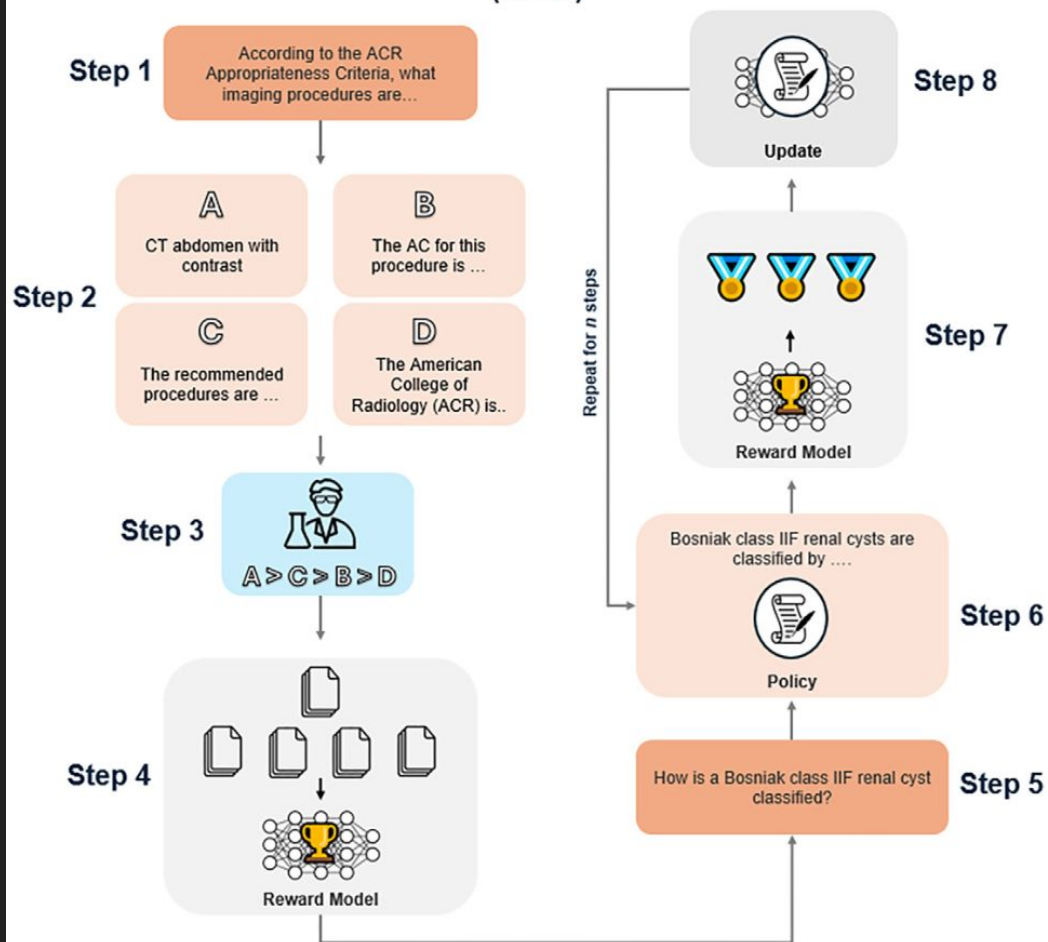**<u>SFT</u> (Supervised Fine-Tuning):** Train on prompt–response pairs.

Good for well-defined tasks with a narrow range of correct answers

**<u>RLHF</u> (Reinforcement Learning from Human Feedback):** Humans rank outputs, a reward model learns preferences, and the LLM adapts to maximize that reward.

**<u>DPO</u> (Direct Preference Optimization):** Like RLHF but skips the reward model

Simpler and needs fewer examples.

# Reinforcement Learning from Human Feedback (RLHF)

**Step 1** — According to the ACR Appropriateness Criteria, what imaging procedures are…

**Step 2**

- **A** — CT abdomen with contrast
- **B** — The AC for this procedure is …
- **C** — The recommended procedures are …
- **D** — The American College of Radiology (ACR) is..

**Step 3** — A ≥ C > B > D

**Step 4** — Reward Model

**Step 5** — How is a Bosniak class IIF renal cyst classified?

**Step 6** — Bosniak class IIF renal cysts are classified by …. / Policy

**Step 7** — Reward Model

**Step 8** — Update

Repeat for *n* steps

**Direct Preference Optimization (DPO)**

**Step 1:** According to the ACR Appropriateness Criteria, what imaging procedures are…

**Step 2:**
A — CT abdomen with contrast
B — The AC for this procedure is …

**Step 3:**
A ✓ Preferred
B ✗ Unpreferred

**Step 4:**
Frozen LLM Reference Copy
Trainable LLM

**Step 5:**

Frozen LLM:
✓ .04 CT | .89 abdomen | .98 with | .45 contrast
Preferred score = .04 x .89 x .98 x .45 = .017

✗ .67 The | .85 AC | .15 for | .20 this | .06 procedure | .34 is
Unpreferred score = .67 x .85 x .15 x .20 x .06 x .34 = .0003

Trainable LLM:
✓ .95 CT | .92 abdomen | .10 with | .70 contrast
Preferred score = .95 x .92 x .10 x .70 = .061

✗ .15 The | .22 AC | .08 for | .11 this | .01 procedure | .27 is
Unpreferred score = .15 x .22 x .08 x .11 x .01 x .27 = .0000008

**Step 6:**
$$R_{Frozen} = \frac{\text{Preferred Score}}{\text{Unpreferred Score}}$$

$$R_{Trainable} = \frac{\text{Preferred Score}}{\text{Unpreferred Score}}$$

**Step 7:**
Weights are updated
$$Loss_{DPO} = \frac{R_{Trainable}}{R_{Frozen}} +$$

## Fine-tuning

Training scale: *depend on the reasoning ability required and the computational resources of the user*

**Full fine-tune**: Adjust all model parameters *(best performance but expensive)*.

**LoRA (Low-Rank Adaptation)**: Only adjust small parameter subsets *(much cheaper in memory and time)*.

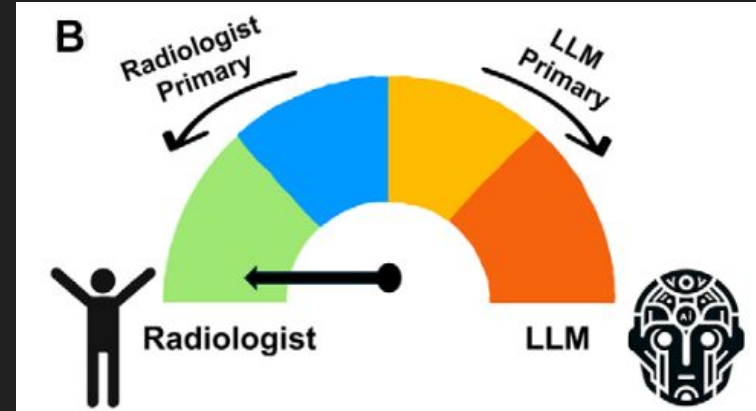**QLoRA**: LoRA + quantized LLM  *(even lower resource usage)*.

# I am a radiologist, how can I implement an open source LLM ?

## Define the Problem First

➔   What is the *exact* use case? Clinical task? Research? Administrative?

➔   What's the measurable outcome? Accuracy, speed, cost savings, reduced workload, patient safety?

## Regulatory & Risk Context

➔   Will this be used in clinical care or just for research?

➔   Is a human-in-the-loop required?

➔   Does the output have direct patient impact?

➔   What are the privacy laws that apply (HIPAA, …)?

➔   What's the risk if the model hallucinates or makes an error?

➔   Implement post-deployment surveillance for safety.

## Cost & Resource Planning

**Do you need the best performance immediately, with minimal setup?**

➜     **Proprietary models:** token-based API costs (input + output).

**Do you need full control over the model?**

➜     **Open-source models:** hardware cost (GPUs, cloud compute), IT staff time, energy consumption.

**Do you plan for scaling ?**

➜     more users = more compute or higher API spend.