# 1. Introduction and Overview

> 1.1 What is Fine-Tuning LLM?

Fine-tuning is the process of continuing the training of a pre-trained model on a domain-specific dataset. In medical physics, this means adapting general-purpose models such as Llama or Gemma to understand terminology, concepts, and applications specific to the medical physics field.

> 1.2 Why Medical Physics?

Medical physics is an interdisciplinary field that combines:

- **Fundamental physics**: Radiation, electromagnetism, quantum mechanics
- **Medical applications**: Radiotherapy, diagnostic imaging, nuclear medicine
- **Technology**: Linear accelerators, CT, MRI, PET scanning
- **Regulation**: Radiation safety, quality assurance, clinical protocols

> 1.3 PhD Portfolio Objectives

This project will demonstrate your capabilities in:

- **Domain expertise**: In-depth understanding of medical physics
- **Technical skills**: Machine learning, deep learning, NLP
- **Research methodology**: Literature review, data analysis, evaluation
- **Innovation**: Applying AI in scientific domains
- **Communication**: Clear documentation and presentation skills

> 1.4 Advantages of This Approach

Compared to training from scratch or RAG (Retrieval-Augmented Generation), fine-tuning provides:

- **Domain-specific knowledge injection**: The model learns specialized terminology
- **Cost-effectiveness**: Cheaper than full model training
- **Memory efficiency**: Uses LoRA/QLoRA for optimized performance
- **Measurable results**: Can be evaluated using standard benchmarks
- **Deployability**: The model can be applied to real-world use cases

# 2. Setup Environment Google Colab

2.1 Hardware Setup

```
1 # Cek GPU availability
2 import torch
3 print(f"CUDA available: {torch.cuda.is_available()}")
4 print(f"GPU count: {torch.cuda.device_count()}")
5 if torch.cuda.is_available():
6     print(f"GPU name: {torch.cuda.get_device_name(0)}")
7     print(f"GPU memory: {torch.cuda.get_device_properties(0).total_memory / 1024**3:.1f} GB")
8
```

```
CUDA available: True
GPU count: 1
GPU name: Tesla T4
GPU memory: 14.7 GB
```

Pastikan menggunakan T4 GPU (16GB) minimum:

- Runtime → Change runtime type
- Hardware accelerator → GPU
- GPU type → T4 GPU (atau yang lebih tinggi jika tersedia)

2.2 Library Installation

```
1 # ==========================
2 # Sub 2: Dependency Setup
3 # ==========================
4
5 # Uninstall versi lama (jika ada)
6 !pip uninstall -y unsloth xformers peft transformers accelerate bitsandbytes
```

```
 7
 8 # Install dan upgrade dependencies utama
 9 !pip install --upgrade --no-cache-dir \
10     unsloth[colab-new] \
11     xformers \
12     peft \
13     transformers \
14     accelerate \
15     bitsandbytes
16
17 # Import modul utama setelah instalasi selesai
18 import torch
19 from unsloth import FastLanguageModel
20 from transformers import AutoTokenizer, AutoModelForCausalLM, BitsAndBytesConfig
21 from peft import LoraConfig, get_peft_model, prepare_model_for_kbit_training
22 from trl import SFTTrainer, SFTConfig # Added SFTConfig import
23
24 # Set environment variable opsional agar xformers memberi lebih banyak detail saat error
25 import os
26 os.environ["XFORMERS_MORE_DETAILS"] = "1"
27
28 print("Sub 2: Dependencies installed and libraries imported successfully.")
```

```
WARNING: Skipping unsloth as it is not installed.
WARNING: Skipping xformers as it is not installed.
Found existing installation: peft 0.17.1
Uninstalling peft-0.17.1:
  Successfully uninstalled peft-0.17.1
Found existing installation: transformers 4.56.1
Uninstalling transformers-4.56.1:
  Successfully uninstalled transformers-4.56.1
Found existing installation: accelerate 1.10.1
Uninstalling accelerate-1.10.1:
  Successfully uninstalled accelerate-1.10.1
WARNING: Skipping bitsandbytes as it is not installed.
Collecting xformers
  Downloading xformers-0.0.32.post2-cp39-abi3-manylinux_2_28_x86_64.whl.metadata (1.1 kB)
Collecting peft
  Downloading peft-0.17.1-py3-none-any.whl.metadata (14 kB)
Collecting transformers
  Downloading transformers-4.56.1-py3-none-any.whl.metadata (42 kB)
                                                    ━━━━━━━━━━━━━━━━━━━━━━━━━━ 42.2/42.2 kB 19.2 MB/s eta 0:00:00
Collecting accelerate
  Downloading accelerate-1.10.1-py3-none-any.whl.metadata (19 kB)
Collecting bitsandbytes
  Downloading bitsandbytes-0.47.0-py3-none-manylinux_2_24_x86_64.whl.metadata (11 kB)
Collecting unsloth[colab-new]
  Downloading unsloth-2025.9.4-py3-none-any.whl.metadata (52 kB)
                                                    ━━━━━━━━━━━━━━━━━━━━━━━━━━ 52.3/52.3 kB 205.5 MB/s eta 0:00:00
Collecting unsloth_zoo>=2025.9.5 (from unsloth[colab-new])
  Downloading unsloth_zoo-2025.9.5-py3-none-any.whl.metadata (9.5 kB)
Requirement already satisfied: torch>=2.4.0 in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (
Requirement already satisfied: triton>=3.0.0 in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new])
Requirement already satisfied: packaging in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (25.
Collecting tyro (from unsloth[colab-new])
  Downloading tyro-0.9.31-py3-none-any.whl.metadata (11 kB)
Collecting datasets<4.0.0,>=3.4.1 (from unsloth[colab-new])
  Downloading datasets-3.6.0-py3-none-any.whl.metadata (19 kB)
Requirement already satisfied: sentencepiece>=0.2.0 in /usr/local/lib/python3.12/dist-packages (from unsloth[colab
Requirement already satisfied: tqdm in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (4.67.1)
Requirement already satisfied: psutil in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (5.9.5)
Requirement already satisfied: wheel>=0.42.0 in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new])
Requirement already satisfied: numpy in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (2.0.2)
Collecting trl!=0.15.0,!=0.19.0,!=0.9.0,!=0.9.1,!=0.9.2,!=0.9.3,>=0.7.9 (from unsloth[colab-new])
  Downloading trl-0.23.0-py3-none-any.whl.metadata (11 kB)
Requirement already satisfied: protobuf in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (5.29
Requirement already satisfied: huggingface_hub>=0.34.0 in /usr/local/lib/python3.12/dist-packages (from unsloth[co
Requirement already satisfied: hf_transfer in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (0
Requirement already satisfied: diffusers in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (0.3
Requirement already satisfied: torchvision in /usr/local/lib/python3.12/dist-packages (from unsloth[colab-new]) (0
Requirement already satisfied: filelock in /usr/local/lib/python3.12/dist-packages (from torch>=2.4.0->unsloth[col
Requirement already satisfied: typing-extensions>=4.10.0 in /usr/local/lib/python3.12/dist-packages (from torch>=2
Requirement already satisfied: setuptools in /usr/local/lib/python3.12/dist-packages (from torch>=2.4.0->unsloth[c
Requirement already satisfied: sympy>=1.13.3 in /usr/local/lib/python3.12/dist-packages (from torch>=2.4.0->unslot
Requirement already satisfied: networkx in /usr/local/lib/python3.12/dist-packages (from torch>=2.4.0->unsloth[col
Requirement already satisfied: jinja2 in /usr/local/lib/python3.12/dist-packages (from torch>=2.4.0->unsloth[colab
Requirement already satisfied: fsspec in /usr/local/lib/python3.12/dist-packages (from torch>=2.4.0->unsloth[colab
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.6.77 in /usr/local/lib/python3.12/dist-packages (from to
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.6.77 in /usr/local/lib/python3.12/dist-packages (from
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.6.80 in /usr/local/lib/python3.12/dist-packages (from to
Requirement already satisfied: nvidia-cudnn-cu12==9.10.2.21 in /usr/local/lib/python3.12/dist-packages (from torch
```

```
1 # # Install required libraries
2 # %%capture
3 # !pip install transformers==4.44.2
4 # !pip install datasets==3.0.0
5 # !pip install peft==0.12.0
```

```
 6 # !pip install trl==0.10.1
 7 # !pip install bitsandbytes==0.43.3
 8 # !pip install accelerate==0.33.0
 9 # !pip install unsloth-zoo
10 # !pip install unsloth[colab-new]@git+https://github.com/unslothai/unsloth.git
11 # !pip install --no-deps xformers==0.0.27.post2
12
13 # # Additional utilities
14 # !pip install wandb scipy scikit-learn matplotlib seaborn
15
```

```
 1 # !pip uninstall numpy pandas -y
 2 # # Reinstalling other libraries will pull in compatible versions of numpy and pandas
 3 # !pip install transformers==4.44.2 datasets==3.0.0 peft==0.12.0 trl==0.10.1 bitsandbytes==0.43.3 accelerate==(
```

```
 1 # !pip install pandas
```

```
 1 # !pip install trl==0.10.1
```

```
 1 # !pip install unsloth[colab-new]@git+https://github.com/unslothai/unsloth.git
```

2.3 Import Libraries dan Setup

```
 1 # import os
 2 # import json
 3 # import torch
 4 # import pandas as pd
 5 # import numpy as np
 6 # from datetime import datetime
 7 # from typing import Dict, List, Optional, Tuple
 8
 9 # # Core ML libraries
10 # from datasets import Dataset, load_dataset
11 # from transformers import (
12 #     AutoTokenizer, AutoModelForCausalLM,
13 #     TrainingArguments, BitsAndBytesConfig
14 # )
15 # from peft import LoraConfig, get_peft_model, prepare_model_for_kbit_training
16 # from trl import SFTTrainer, SFTConfig
17 # from unsloth import FastLanguageModel
18
19 # # Evaluation
20 # from sklearn.metrics import accuracy_score, f1_score
21 # import matplotlib.pyplot as plt
22 # import seaborn as sns
23
24 # # Warnings
25 # import warnings
26 # warnings.filterwarnings('ignore')
27
28 # print(f"Setup completed at {datetime.now()}")
29
```

Double-click (or enter) to edit

2.4 Authentication Setup

```
 1 # Hugging Face authentication
 2 from huggingface_hub import login
 3
 4 # Login ke Hugging Face (diperlukan untuk akses model dan upload)
 5 login()
 6
 7 # Optional: Weights & Biases untuk tracking
 8 import wandb
 9 wandb.login()
10
```

```
/usr/local/lib/python3.12/dist-packages/notebook/notebookapp.py:191: SyntaxWarning: invalid escape sequence '\/'
  | |_| | '_ \/ _` / _` |  _/ -_)
wandb: Logging into wandb.ai. (Learn how to deploy a W&B server locally: https://wandb.me/wandb-server)
wandb: You can find your API key in your browser here: https://wandb.ai/authorize
wandb: Paste an API key from your profile and hit enter: ··········
wandb: WARNING If you're specifying your api key in code, ensure this code is not shared publicly.
wandb: WARNING Consider setting the WANDB_API_KEY environment variable, or running `wandb login` from the command li
wandb: No netrc file found, creating one.
wandb: Appending key for api.wandb.ai to your netrc file: /root/.netrc
wandb: W&B API key is configured. Use `wandb login --relogin` to force relogin
True
```

## 3. Dataset Medical Physics dan Preprocessing

### 3.1 Dataset Sources

Untuk medical physics, kita akan menggunakan:

- Radiation Oncology NLP Database (ROND) - Dataset khusus radiation oncology
- Medical Physics Q&A - Pertanyaan-jawaban dari textbook dan papers
- AAPM Guidelines - American Association of Physicists in Medicine
- Synthetic Data - Generated medical physics problems

### 3.2 Data Loading dan Exploration

```python
1  def load_medical_physics_data():
2      """Load dan combine berbagai sumber data medical physics"""
3
4      # Sample data medical physics Q&A
5      medical_physics_qa = [
6          {
7              "input": "What is the half-life of Technetium-99m?",
8              "output": "The half-life of Technetium-99m (Tc-99m) is approximately 6.01 hours. This short half-li
9          },
10         {
11             "input": "Explain the inverse square law in radiation physics.",
12             "output": "The inverse square law states that the intensity of radiation is inversely proportional
13         },
14         {
15             "input": "What is the difference between absorbed dose, equivalent dose, and effective dose?",
16             "output": "Absorbed dose (Gy) measures energy deposited per unit mass. Equivalent dose (Sv) account
17         },
18         {
19             "input": "How does a linear accelerator (LINAC) produce X-rays?",
20             "output": "A LINAC accelerates electrons using radiofrequency electromagnetic fields. High-energy e
21         }
22         # Add more medical physics Q&A pairs here
23     ]
24
25     return medical_physics_qa
26
27 # Load data
28 raw_data = load_medical_physics_data()
29 print(f"Loaded {len(raw_data)} medical physics Q&A pairs")
30
31 # Convert to DataFrame for easier handling
32 import pandas as pd
33 df = pd.DataFrame(raw_data)
34 print(df.head())
```

```
Loaded 4 medical physics Q&A pairs
                                                   input  \
0          What is the half-life of Technetium-99m?
1  Explain the inverse square law in radiation ph...
2  What is the difference between absorbed dose, ...
3  How does a linear accelerator (LINAC) produce ...

                                                  output
0  The half-life of Technetium-99m (Tc-99m) is ap...
1  The inverse square law states that the intensi...
2  Absorbed dose (Gy) measures energy deposited p...
3  A LINAC accelerates electrons using radiofrequ...
```

3.3 Data Formatting untuk Fine-tuning

```python
1 from typing import List, Dict
2
3 def format_instruction_data(data: List[Dict]) -> List[Dict]:
4     """Format data ke format instruksi yang sesuai untuk fine-tuning"""
5
6     formatted_data = []
7
8     for item in data:
9         # Format template untuk medical physics instruction following
10        formatted_text = f"""Below is an instruction related to medical physics. Write a response that appropr:
11
12 ### Instruction:
13 {item['input']}
14
15 ### Response:
16 {item['output']}"""
17
18        formatted_data.append({"text": formatted_text})
19
20    return formatted_data
21
22 # Format data
23 formatted_data = format_instruction_data(raw_data)
24 print("Sample formatted data:")
25 print(formatted_data[0]["text"])
```

```
Sample formatted data:
Below is an instruction related to medical physics. Write a response that appropriately completes the request.

### Instruction:
What is the half-life of Technetium-99m?

### Response:
The half-life of Technetium-99m (Tc-99m) is approximately 6.01 hours. This short half-life makes it ideal for diagno
```

3.4 Dataset Preparation

```python
1 from typing import List, Dict, Tuple
2 from datasets import Dataset
3
4 def prepare_dataset(formatted_data: List[Dict],
5                     test_size: float = 0.2,
6                     random_state: int = 42) -> Tuple[Dataset, Dataset]:
7     """Prepare train dan validation datasets"""
8
9     from sklearn.model_selection import train_test_split
10
11    # Split data
12    train_data, val_data = train_test_split(
13        formatted_data,
14        test_size=test_size,
15        random_state=random_state
16    )
17
18    # Convert to HuggingFace Dataset
19    train_dataset = Dataset.from_list(train_data)
20    val_dataset = Dataset.from_list(val_data)
21
22    print(f"Train dataset size: {len(train_dataset)}")
23    print(f"Validation dataset size: {len(val_dataset)}")
24
25    return train_dataset, val_dataset
26
```

```
27 # Prepare datasets
```

```
Train dataset size: 3
Validation dataset size: 1
```

## 4. Model Selection dan Loading

### 4.1 Model Selection Criteria

Untuk medical physics fine-tuning, pertimbangkan:

- Size: 7B-13B parameters optimal untuk Colab
- Domain adaptability: Models yang responsive terhadap fine-tuning
- Memory efficiency: Support untuk quantization
- License: Open source untuk research

### 4.2 Recommended Models

```python
1 # Model options untuk medical physics
2 MODEL_OPTIONS = {
3     "llama3-8b": "meta-llama/Llama-3.2-8B-Instruct",
4     "gemma-7b": "google/gemma-1.1-7b-it", # Using Gemma 7B which is more likely to fit on a T4
5     "mistral-7b": "mistralai/Mistral-7B-Instruct-v0.3",
6     "unsloth-llama": "unsloth/llama-3-8b-Instruct-bnb-4bit"
7 }
8
9 # Pilih model
10 MODEL_NAME = MODEL_OPTIONS["gemma-7b"] # Changed to Gemma 7B
11 print(f"Selected model: {MODEL_NAME}")
```

```
Selected model: google/gemma-1.1-7b-it
```

### 4.3 Model Loading dengan Quantization

```python
1 def load_model_and_tokenizer(model_name: str,
2                              max_seq_length: int = 1024,
3                              load_in_4bit: bool = True):
4     """Load model dan tokenizer dengan quantization"""
5
6     if "unsloth" in model_name.lower():
7         # Use Unsloth for optimized loading
8         model, tokenizer = FastLanguageModel.from_pretrained(
9             model_name=model_name,
10            max_seq_length=max_seq_length,
11            dtype=None,  # Auto-detect
12            load_in_4bit=load_in_4bit,
13        )
14
15        # Enable fast training
16        model = FastLanguageModel.get_peft_model(
17            model,
18            r=16,  # LoRA rank
19            target_modules=["q_proj", "k_proj", "v_proj", "o_proj",
20                            "gate_proj", "up_proj", "down_proj"],
21            lora_alpha=16,
22            lora_dropout=0.0,
23            bias="none",
24            use_gradient_checkpointing="unsloth",
25            random_state=3407,
26            use_rslora=False,
27            loftq_config=None,
28        )
29
30    else:
31        # Standard loading dengan quantization config
32        quantization_config = BitsAndBytesConfig(
33            load_in_4bit=True,
34            bnb_4bit_use_double_quant=True,
35            bnb_4bit_quant_type="nf4",
36            bnb_4bit_compute_dtype=torch.bfloat16
37        )
38
39        model = AutoModelForCausalLM.from_pretrained(
40            model_name,
41            quantization_config=quantization_config,
42            device_map="auto",
```

```
43              torch_dtype=torch.bfloat16,
44              trust_remote_code=True
45          )
46
47          tokenizer = AutoTokenizer.from_pretrained(model_name)
48
49          # Add padding token if missing
50          if tokenizer.pad_token is None:
51              tokenizer.pad_token = tokenizer.eos_token
52
53      return model, tokenizer
54
55 # Load model dan tokenizer
56 model, tokenizer = load_model_and_tokenizer(MODEL_NAME)
57
58 print(f"Model loaded: {model}")
59 print(f"Tokenizer vocab size: {len(tokenizer)}")
```

| | | |
|---|---|---|
| config.json: 100% | | 620/620 [00:00<00:00, 15.6kB/s] |

`torch_dtype` is deprecated! Use `dtype` instead!

| | | |
|---|---|---|
| model.safetensors.index.json: 100% | | 20.9k/20.9k [00:00<00:00, 452kB/s] |
| Fetching 4 files: 100% | 4/4 [05:11<00:00, 80.80s/it] | |
| model-00001-of-00004.safetensors: 100% | | 5.00G/5.00G [05:10<00:00, 46.6MB/s] |
| model-00004-of-00004.safetensors: 100% | | 2.11G/2.11G [00:55<00:00, 34.6MB/s] |
| model-00002-of-00004.safetensors: 100% | | 4.98G/4.98G [05:00<00:00, 12.9MB/s] |
| model-00003-of-00004.safetensors: 100% | | 4.98G/4.98G [05:11<00:00, 12.9MB/s] |
| Loading checkpoint shards: 100% | 4/4 [01:21<00:00, 18.61s/it] | |
| generation_config.json: 100% | 132/132 [00:00<00:00, 8.78kB/s] | |
| tokenizer_config.json: 100% | 34.2k/34.2k [00:00<00:00, 4.10MB/s] | |
| tokenizer.model: 100% | 4.24M/4.24M [00:00<00:00, 8.64MB/s] | |
| tokenizer.json: 100% | 17.5M/17.5M [00:01<00:00, 16.7MB/s] | |
| special_tokens_map.json: 100% | 636/636 [00:00<00:00, 49.5kB/s] | |

```
Model loaded: GemmaForCausalLM(
  (model): GemmaModel(
    (embed_tokens): Embedding(256000, 3072, padding_idx=0)
    (layers): ModuleList(
      (0-27): 28 x GemmaDecoderLayer(
        (self_attn): GemmaAttention(
          (q_proj): Linear4bit(in_features=3072, out_features=4096, bias=False)
          (k_proj): Linear4bit(in_features=3072, out_features=4096, bias=False)
          (v_proj): Linear4bit(in_features=3072, out_features=4096, bias=False)
          (o_proj): Linear4bit(in_features=4096, out_features=3072, bias=False)
        )
        (mlp): GemmaMLP(
          (gate_proj): Linear4bit(in_features=3072, out_features=24576, bias=False)
          (up_proj): Linear4bit(in_features=3072, out_features=24576, bias=False)
          (down_proj): Linear4bit(in_features=24576, out_features=3072, bias=False)
          (act_fn): PytorchGELUTanh()
        )
        (input_layernorm): GemmaRMSNorm((3072,), eps=1e-06)
        (post_attention_layernorm): GemmaRMSNorm((3072,), eps=1e-06)
      )
    )
    (norm): GemmaRMSNorm((3072,), eps=1e-06)
    (rotary_emb): GemmaRotaryEmbedding()
  )
  (lm_head): Linear(in_features=3072, out_features=256000, bias=False)
)
Tokenizer vocab size: 256000
```

```
1 # # MODEL_NAME = "unsloth/llama-3.2-8b-instruct-bnb-4bit"
2
3 # model, tokenizer = load_model_and_tokenizer(MODEL_NAME)
4
5 # print(f"Model loaded: {model}")
6 # print(f"Tokenizer vocab size: {len(tokenizer)}")
7
```

```
1 # from unsloth import FastLanguageModel
2 # import torch
3 # from transformers import AutoTokenizer, AutoModelForCausalLM, BitsAndBytesConfig
4 # from peft import LoraConfig, get_peft_model, prepare_model_for_kbit_training
5
6
```

```
 7 # def load_model_and_tokenizer(model_name: str,
 8 #                               max_seq_length: int = 1024,
 9 #                               load_in_4bit: bool = True):
10 #     """Load model dan tokenizer dengan quantization"""
11
12 #     if "unsloth" in model_name.lower():
13 #         # Use Unsloth for optimized loading
14 #         model, tokenizer = FastLanguageModel.from_pretrained(
15 #             model_name=model_name,
16 #             max_seq_length=max_seq_length,
17 #             dtype=None,  # Auto-detect
18 #             load_in_4bit=load_in_4bit,
19 #             trust_remote_code=True, # Added this line
20 #             # Add additional unsloth specific arguments if needed based on documentation
21 #             # For example, if using a specific unsloth model variant
22 #         )
23
24 #         # Enable fast training with LoRA
25 #         model = FastLanguageModel.get_peft_model(
26 #             model,
27 #             r=16,  # LoRA rank
28 #             target_modules=["q_proj", "k_proj", "v_proj", "o_proj",
29 #                             "gate_proj", "up_proj", "down_proj"],
30 #             lora_alpha=16,
31 #             lora_dropout=0.0,
32 #             bias="none",
33 #             use_gradient_checkpointing="unsloth",
34 #             random_state=3407,
35 #             use_rslora=False,
36 #             loftq_config=None,
37 #         )
38
39 #     else:
40 #         # Standard loading dengan quantization config
41 #         quantization_config = BitsAndBytesConfig(
42 #             load_in_4bit=True,
43 #             bnb_4bit_use_double_quant=True,
44 #             bnb_4bit_quant_type="nf4",
45 #             bnb_4bit_compute_dtype=torch.bfloat16
46 #         )
47
48 #         model = AutoModelForCausalLM.from_pretrained(
49 #             model_name,
50 #             quantization_config=quantization_config,
51 #             device_map="auto",
52 #             torch_dtype=torch.bfloat16,
53 #             trust_remote_code=True
54 #         )
55
56 #         tokenizer = AutoTokenizer.from_pretrained(model_name)
57
58 #         # Add padding token if missing
59 #         if tokenizer.pad_token is None:
60 #             tokenizer.pad_token = tokenizer.eos_token
61
62 #     return model, tokenizer
63
64 # # Load model dan tokenizer
65 # model, tokenizer = load_model_and_tokenizer(MODEL_NAME)
66
67 # print(f"Model loaded: {model}")
68 # print(f"Tokenizer vocab size: {len(tokenizer)}")
```

> 5. Parameter-Efficient Fine-Tuning (LoRA/QLoRA)

   ↳ 4 cells hidden

∨  6. Training Configuration dan Monitoring

6.1 Training Arguments

```
1  def create_training_config(
2      output_dir: str = "./medical-physics-llm",
3      num_train_epochs: int = 3,
4      per_device_train_batch_size: int = 1,
5      gradient_accumulation_steps: int = 4,
6      learning_rate: float = 2e-4,
7      max_steps: int = 100
```

```
 7      max_steps: int = 100,
 8      warmup_ratio: float = 0.03,
 9      logging_steps: int = 10,
10      save_steps: int = 50,
11      eval_strategy: str = "steps", # Changed argument name
12      eval_steps: int = 50,
13      max_grad_norm: float = 0.3
14  ) -> SFTConfig:
15      """Create training configuration untuk medical physics fine-tuning"""
16
17      training_config = SFTConfig(
18          output_dir=output_dir,
19          num_train_epochs=num_train_epochs,
20          per_device_train_batch_size=per_device_train_batch_size,
21          per_device_eval_batch_size=1,
22          gradient_accumulation_steps=gradient_accumulation_steps,
23          optim="paged_adamw_32bit",
24          learning_rate=learning_rate,
25          max_steps=max_steps,
26          lr_scheduler_type="cosine",
27          warmup_ratio=warmup_ratio,
28
29          # Logging dan saving
30          logging_dir="./logs",
31          logging_steps=logging_steps,
32          save_strategy="steps",
33          save_steps=save_steps,
34          eval_strategy=eval_strategy, # Changed argument name
35          eval_steps=eval_steps,
36
37          # Optimization
38          fp16=False,
39          bf16=torch.cuda.is_bf16_supported(),
40          max_grad_norm=max_grad_norm,
41          gradient_checkpointing=True,
42
43          # Data handling
44          max_seq_length=1024,
45          packing=False,
46
47          # Misc
48          report_to=None,  # Disable wandb untuk demo
49          seed=42,
50          data_seed=42,
51          remove_unused_columns=False,
52
53          # Dataset specific
54          dataset_text_field="text",
55      )
56
57      return training_config
58
59  # Create training config
60  training_config = create_training_config()
```

6.2 Initialize Trainer

```
 1  def create_trainer(model, tokenizer, train_dataset, val_dataset,
    training_config):
 2      """Create SFT trainer untuk medical physics fine-tuning"""
 3
 4      trainer = SFTTrainer(
 5          model=model,
 6          tokenizer=tokenizer,
 7          train_dataset=train_dataset,
 8          eval_dataset=val_dataset,
 9          args=training_config,
10          # data_collator=data_collator,  # Will be handled automatically
11      )
12
13      return trainer
14
15  # Create trainer
16  trainer = create_trainer(model, tokenizer, train_dataset, val_dataset,
    training_config)
17  print("Trainer initialized successfully")
```

```
num_proc must be <= 3. Reducing num_proc to 3 for dataset of size 3.
WARNING:datasets.arrow_dataset:num_proc must be <= 3. Reducing num_proc to 3 for dataset of size 3.
```
Unsloth: Tokenizing ["text"] (num_proc=3): 100%                          3/3 [00:05<00:00,  1.43s/ examples]

```
num_proc must be <= 1. Reducing num_proc to 1 for dataset of size 1.
WARNING:datasets.arrow_dataset:num_proc must be <= 1. Reducing num_proc to 1 for dataset of size 1.
```

### 6.3 Training Process dengan Monitoring

Unsloth: Tokenizing ["text"]: 100%                                       1/1 [00:00<00:00, 51.44 examples/s]

```python
1 from datetime import datetime # Added import
2
3 def train_model_with_monitoring(trainer):
4     """Train model dengan monitoring progress"""
5
6     print("Starting training...")
7     print(f"Training started at: {datetime.now()}")
8
9     # Clear CUDA cache
10    if torch.cuda.is_available():
11        torch.cuda.empty_cache()
12
13    # Train model
14    train_result = trainer.train()
15
16    print(f"Training completed at: {datetime.now()}")
17    print(f"Training stats:")
18    print(f"  – Final loss: {train_result.training_loss:.4f}")
19    print(f"  – Training steps: {train_result.global_step}")
20
21    # Save final model
22    trainer.save_model()
23    tokenizer.save_pretrained(training_config.output_dir)
24
25    return train_result
26
27 # Start training
28 train_result = train_model_with_monitoring(trainer)
```

```
Starting training...
Training started at: 2025-09-13 04:52:58.994672
creating run (1.2m)
ERROR retrying: Post "https://api.wandb.ai/graphql": net/http: request canceled (Client.Timeout exceeded while awaiting headers)
Tracking run with wandb version 0.21.3
Run data is saved locally in /content/wandb/run-20250913_045309-a662i8us
Syncing run solar-bee-3 to Weights & Biases (docs)
View project at https://wandb.ai/auliaoctavvia-authfy/huggingface
View run at https://wandb.ai/auliaoctavvia-authfy/huggingface/runs/a662i8us
wandb: 500 encountered ({"errors":[{"message":"context canceled","path":["upsertBucket"]}],"data":{"upsertBucket":nu
                                          [100/100 18:16, Epoch 100/100]
```

| Step | Training Loss | Validation Loss |
|------|---------------|-----------------|
| 50   | 0.012100      | 1.678695        |
| 100  | 0.011400      | 1.703672        |

```
Training completed at: 2025-09-13 05:12:55.377528
Training stats:
  – Final loss: 0.4743
  – Training steps: 100
```

### 6.4 Monitor Training Progress

```python
1 import matplotlib.pyplot as plt # Added import
2
3 def plot_training_progress(log_history):
4     """Plot training progress dari trainer logs"""
5
6     if not log_history:
7         print("No training history available")
8         return
9
10    # Extract metrics
11    steps = []
12    train_loss = []
13    eval_loss = []
14
15    for log in log_history:
16        if 'loss' in log:
17            steps.append(log['step'])
18            train_loss.append(log['loss'])
19        if 'eval_loss' in log:
20            eval_loss.append(log['eval_loss'])
21
22    # Plot
```

```
23    plt.figure(figsize=(12, 4))
24
25    plt.subplot(1, 2, 1)
26    plt.plot(steps, train_loss, label='Training Loss')
27    plt.xlabel('Steps')
28    plt.ylabel('Loss')
29    plt.title('Training Loss Over Time')
30    plt.legend()
31    plt.grid(True)
32
33    if eval_loss:
34        plt.subplot(1, 2, 2)
35        eval_steps = steps[-len(eval_loss):]
36        plt.plot(eval_steps, eval_loss, label='Validation Loss', color='orange')
37        plt.xlabel('Steps')
38        plt.ylabel('Loss')
39        plt.title('Validation Loss Over Time')
40        plt.legend()
41        plt.grid(True)
42
43    plt.tight_layout()
44    plt.show()
45
46 # Plot training progress
47 plot training progress(trainer.state.log history)
```



## 7. Evaluasi dan Testing Model

### 7.1 Inference Testing

```
1   def test_model_inference(model, tokenizer, test_prompts: List[str]):
2       """Test model inference dengan medical physics prompts"""
3
4       model.eval()  # Set to evaluation mode
5
6       results = []
7
8       for prompt in test_prompts:
9           # Format prompt
10          formatted_prompt = f"""Below is an instruction related to medical
            physics. Write a response that appropriately completes the request.

12 ### Instruction:
13 {prompt}

15 ### Response:
16 """

18          # Tokenize input
19          inputs = tokenizer(
20              formatted_prompt,
21              return_tensors="pt",
22              truncation=True,
23              max_length=512
24          ).to(model.device)
25
26          # Generate response
27          with torch.no_grad():
28              outputs = model.generate(
```

```
29                **inputs,
30                max_new_tokens=256,
31                temperature=0.7,
32                do_sample=True,
33                top_p=0.9,
34                repetition_penalty=1.1,
35                pad_token_id=tokenizer.eos_token_id
36            )
37
38        # Decode response
39        full_response = tokenizer.decode(outputs[0], skip_special_tokens=True)
40        response = full_response.split("### Response:")[-1].strip()
41
42        results.append({
43            "prompt": prompt,
44            "response": response,
45            "full_output": full_response
46        })
47
48        print(f"Prompt: {prompt}")
49        print(f"Response: {response}")
50        print("-" * 80)
51
52    return results
53
54  # Test prompts medical physics
55  test_prompts = [
56      "What is the photoelectric effect in medical imaging?",
57      "Calculate the half-value layer for 100 keV X-rays in aluminum.",
58      "Explain the principles of IMRT planning.",
59      "What are the main components of a PET scanner?",
60      "Describe radiation safety principles in nuclear medicine."
61  ]
62
63  # Run inference tests
64  inference_results = test_model_inference(model, tokenizer, test_prompts)
```

```
Prompt: What is the photoelectric effect in medical imaging?
Response: The Photovoltaic Effect (photoelectron emission) describes how X-rays interact with matter, producing ele
--------------------------------------------------------------------------------
Prompt: Calculate the half-value layer for 100 keV X-rays in aluminum.
Response: The Half Value Layer (HVL) of any material refers  to thickness required so as radiation intensity reduce

For **X - rays** specifically; their interaction involves electromagnetic interactions like scattering , absorption
--------------------------------------------------------------------------------
Prompt: Explain the principles of IMRT planning.
Response: IMRT (Intensity Modulated Radiation Therapy) treatment plans are designed using sophisticated software an

- **Patient Positioning:** Precise measurements ensure accurate targeting with minimal beam divergence from intende

 -**Planning Target Definition**: Tumors or regions requiring intervention defined on imaging studies like CT scans

The plan then undergoes review evaluation testing before being implemented clinically
--------------------------------------------------------------------------------
Prompt: What are the main components of a PET scanner?
Response: APET (Positron Emission Tomography) scanners consist primarily  of three major parts - 1] The cyclotron,
--------------------------------------------------------------------------------
Prompt: Describe radiation safety principles in nuclear medicine.
Response: Radiation Safety Principles (RSP) are fundamental concepts used for safe handling and administration of r
--------------------------------------------------------------------------------
```

7.2 Quantitative Evaluation

```
1 def calculate_perplexity(model, tokenizer, eval_dataset, max_samples: int = 50):
2     """Calculate perplexity pada evaluation dataset"""
3
4     model.eval()
5     total_loss = 0
6     total_tokens = 0
7
8     # Sample subset untuk evaluation
9     eval_samples = eval_dataset.shuffle(seed=42).select(range(min(max_samples, len(eval_dataset))))
10
11     for sample in eval_samples:
12         text = sample["text"]
13
```

```
14        # Tokenize
15        inputs = tokenizer(
16            text,
17            return_tensors="pt",
18            truncation=True,
19            max_length=512
20        ).to(model.device)
21
22        # Calculate loss
23        with torch.no_grad():
24            outputs = model(**inputs, labels=inputs["input_ids"])
25            loss = outputs.loss
26
27        total_loss += loss.item() * inputs["input_ids"].size(1)
28        total_tokens += inputs["input_ids"].size(1)
29
30    avg_loss = total_loss / total_tokens
31    perplexity = torch.exp(torch.tensor(avg_loss))
32
33    print(f"Average loss: {avg_loss:.4f}")
34    print(f"Perplexity: {perplexity:.2f}")
35
36    return perplexity.item()
37
38 # Calculate perplexity
39 perplexity = calculate_perplexity(model, tokenizer, val_dataset)
```

```
Average loss: 1.7037
Perplexity: 5.49
```

7.3 Domain-Specific Evaluation

```
1 def evaluate_medical_physics_knowledge(model, tokenizer):
2     """Evaluate medical physics domain knowledge"""
3
4     # Medical physics multiple choice questions
5     mc_questions = [
6         {
7             "question": "The unit of absorbed dose is:",
8             "options": ["A) Sievert (Sv)", "B) Gray (Gy)", "C) Becquerel (Bq)", "D) Coulomb per kilogram (C/k
9             "correct": "B"
10        },
11        {
12            "question": "In radiation therapy, IMRT stands for:",
13            "options": ["A) Intensity Modulated Radiation Therapy", "B) Image Modulated Radiation Treatment",
14                       "C) Internal Medical Radiation Therapy", "D) Inverse Medical Radiation Treatment"],
15            "correct": "A"
16        },
17        {
18            "question": "The half-life of I-131 is approximately:",
19            "options": ["A) 6 hours", "B) 8 days", "C) 30 days", "D) 1 year"],
20            "correct": "B"
21        }
22    ]
23
24    correct_answers = 0
25    total_questions = len(mc_questions)
26
27    for i, q in enumerate(mc_questions):
28        prompt = f"Question: {q['question']}\n" + "\n".join(q['options']) + "\n\nAnswer:"
29
30        # Generate response
31        inputs = tokenizer(prompt, return_tensors="pt", truncation=True, max_length=256).to(model.device) # A
32        with torch.no_grad():
33            outputs = model.generate(**inputs, max_new_tokens=10, temperature=0.1)
34
35        response = tokenizer.decode(outputs[0], skip_special_tokens=True)
36        answer = response.split("Answer:")[-1].strip()
37
38        print(f"Q{i+1}: {q['question']}")
39        print(f"Model answer: {answer}")
40        print(f"Correct answer: {q['correct']}")
41
42        if q['correct'].lower() in answer.lower():
43            correct_answers += 1
44            print("✓ Correct")
45        else:
46            print("✗ Incorrect")
47        print("-" * 50)
48
```

```
49    accuracy = correct_answers / total_questions
50    print(f"\nOverall accuracy: {accuracy:.2%} ({correct_answers}/{total_questions})")
51
52    return accuracy
53
54 # Evaluate domain knowledge
55 domain_accuracy = evaluate_medical_physics_knowledge(model, tokenizer)
```

```
Q1: The unit of absorbed dose is:
Model answer: A) Sievert (Sv)

Explanation:
Correct answer: B
✗ Incorrect
--------------------------------------------------
Q2: In radiation therapy, IMRT stands for:
Model answer: A) Intensity Modulated Radiation Therapy

Explanation:
Correct answer: A
✓ Correct
--------------------------------------------------
Q3: The half-life of I-131 is approximately:
Model answer: B) 8 days

Explanation:
I
Correct answer: B
✓ Correct
--------------------------------------------------

Overall accuracy: 66.67% (2/3)
```

## 8. Deployment dan Portfolio Presentation

### 8.1 Model Export dan Saving

```python
1 import json # Added import
2
3 def export_model_for_deployment(model, tokenizer, output_dir: str = "./final_model"):
4     """Export model untuk deployment dan portfolio"""
5
6     # Create output directory
7     os.makedirs(output_dir, exist_ok=True)
8
9     # Save model dan tokenizer
10    model.save_pretrained(output_dir)
11    tokenizer.save_pretrained(output_dir)
12
13    # Save configuration
14    config_info = {
15        "model_name": MODEL_NAME,
16        "training_date": datetime.now().isoformat(),
17        "lora_rank": lora_config.r if lora_config else "N/A",
18        "training_steps": training_config.max_steps,
19        "final_perplexity": perplexity,
20        "domain_accuracy": domain_accuracy,
21        "dataset_size": len(train_dataset) + len(val_dataset)
22    }
23
24    with open(f"{output_dir}/model_info.json", "w") as f:
25        json.dump(config_info, f, indent=2)
26
27    print(f"Model exported to: {output_dir}")
28    print(f"Model info saved: {output_dir}/model_info.json")
29
30    return output_dir
31
32 # Export model
33 final_model_dir = export_model_for_deployment(model, tokenizer)
```

```
Model exported to: ./final_model
Model info saved: ./final_model/model_info.json
```

8.2 Push ke Hugging Face Hub

```python
1 def push_to_huggingface_hub(model, tokenizer, repo_name: str):
2     """Push model ke Hugging Face Hub untuk portfolio"""
3
4     # Push model
5     model.push_to_hub(repo_name, private=False)
6     tokenizer.push_to_hub(repo_name, private=False)
7
8     # Create model card
9     model_card_content = f"""
10 # Medical Physics Fine-tuned Language Model
11
12 This model is a fine-tuned version of {MODEL_NAME} on medical physics domain data.
13
14 ## Model Description
15
16 - **Base Model**: {MODEL_NAME}
17 - **Domain**: Medical Physics, Radiation Oncology
18 - **Fine-tuning Method**: LoRA/QLoRA
19 - **Training Data**: Medical physics Q&A, guidelines, and textbook content
20
21 ## Training Details
22
23 - **Training Steps**: {training_config.max_steps}
24 - **Learning Rate**: {training_config.learning_rate}
25 - **LoRA Rank**: {lora_config.r if lora_config else 'N/A'}
26 - **Final Perplexity**: {perplexity:.2f}
27 - **Domain Accuracy**: {domain_accuracy:.2%}
28
29 ## Use Cases
30
31 - Medical physics education and training
32 - Radiation therapy planning assistance
33 - Nuclear medicine consultation
34 - Regulatory compliance guidance
35
36 ## Usage
37
38 ```python
39 from transformers import AutoTokenizer, AutoModelForCausalLM
40
41 tokenizer = AutoTokenizer.from_pretrained("{repo_name}")
42 model = AutoModelForCausalLM.from_pretrained("{repo_name}")
43
44 prompt = "What is the photoelectric effect in medical imaging?"
45 inputs = tokenizer(prompt, return_tensors="pt")
46 outputs = model.generate(**inputs, max_new_tokens=256)
47 response = tokenizer.decode(outputs[0], skip_special_tokens=True)
```

8.3 Portfolio Documentation

```python
def create_portfolio_documentation():
    """Create comprehensive documentation untuk PhD portfolio"""

    portfolio_doc = f"""
```

Medical Physics LLM Fine-Tuning Project PhD Portfolio in Physics

> Project Overview

This project demonstrates the application of state-of-the-art natural language processing techniques to the specialized domain of medical physics. By fine-tuning a large language model on domain-specific data, we created an AI assistant capable of understanding and generating content related to radiation therapy, nuclear medicine, diagnostic imaging, and radiation safety.

> Technical Implementation

**Model Architecture:**

- Base Model: {MODEL_NAME}
- Fine-tuning Method: LoRA (Low-Rank Adaptation)
- Quantization: 4-bit precision for memory efficiency
- Training Framework: Hugging Face Transformers + TRL

**Training Configuration:**

- Dataset Size: {len(train_dataset) + len(val_dataset)} samples
- Training Steps: {training_config.max_steps}
- Learning Rate: {training_config.learning_rate}
- Batch Size: {training_config.per_device_train_batch_size}
- LoRA Rank: {lora_config.r if lora_config else 'N/A'}

**Performance Metrics:**

- Final Perplexity: {perplexity:.2f}
- Domain Accuracy: {domain_accuracy:.2%}
- Training Time: ~2-3 hours on Google Colab T4 GPU

## Key Contributions

1. **Domain Adaptation**: Successfully adapted a general-purpose LLM to medical physics domain
2. **Efficient Training**: Utilized parameter-efficient fine-tuning techniques (LoRA/QLoRA)
3. **Practical Application**: Created a working AI assistant for medical physics education
4. **Open Science**: Model and code made available for research community

## Applications

- **Education**: AI tutor for medical physics students
- **Research**: Literature synthesis and hypothesis generation
- **Clinical Support**: Decision support for radiation therapy planning
- **Regulatory**: Compliance guidance and protocol verification

## Technical Skills Demonstrated

- Deep learning model architecture understanding
- Parameter-efficient fine-tuning techniques
- Domain-specific data curation and preprocessing
- Model evaluation and validation methodologies
- MLOps practices for model deployment
- Scientific communication and documentation

## Impact and Future Work

This project represents a novel application of LLM technology to a specialized scientific domain. Future extensions could include:

- Multi-modal capabilities (incorporating medical images)
- Integration with treatment planning systems
- Regulatory compliance automation
- Real-time clinical decision support

## Code Availability

Complete implementation available at: [GitHub Repository] Trained model available at: [Hugging Face Hub]

## Academic Relevance

This work demonstrates:

- Interdisciplinary research combining AI and medical physics
- Practical application of cutting-edge ML techniques
- Understanding of domain-specific challenges
- Ability to communicate technical concepts clearly
- Commitment to open science and reproducibility

---

*This project was completed as part of PhD application portfolio in Physics, demonstrating both technical expertise and research capability in the intersection of artificial intelligence and medical physics.* """

```
# Save portfolio documentation
with open("portfolio_documentation.md", "w") as f:
    f.write(portfolio_doc)

print("Portfolio documentation created: portfolio_documentation.md")
return portfolio_doc
```

## Create portfolio

portfolio_doc = create_portfolio_documentation()

## 9. Kode Lengkap Step-by-Step

```python
1   # ==============================
2   # MEDICAL PHYSICS LLM FINE-TUNING
3   # Complete Implementation
4   # ==============================
5
6   # 1. SETUP AND IMPORTS
7   print("Setting up environment...")
8
9   %%capture
10  !pip install transformers==4.44.2 datasets==3.0.0 peft==0.12.0 trl==0.10.1
11  !pip install bitsandbytes==0.43.3 accelerate==0.33.0
12  !pip install unsloth[colab-new]@git+https://github.com/unslothai/unsloth.git
13  !pip install wandb scipy scikit-learn matplotlib seaborn
14
15  import os
16  import json
17  import torch
18  import pandas as pd
19  import numpy as np
20  from datetime import datetime
21  from typing import Dict, List, Optional, Tuple
22  import warnings
23  warnings.filterwarnings('ignore')
24
25  from datasets import Dataset
26  from transformers import AutoTokenizer, AutoModelForCausalLM, TrainingArguments
27  from peft import LoraConfig, get_peft_model, prepare_model_for_kbit_training
28  from trl import SFTTrainer, SFTConfig
29  from unsloth import FastLanguageModel
30  from sklearn.model_selection import train_test_split
31  import matplotlib.pyplot as plt
32
33  print(f"Setup completed at {datetime.now()}")
34  print(f"CUDA available: {torch.cuda.is_available()}")
35  if torch.cuda.is_available():
36      print(f"GPU: {torch.cuda.get_device_name(0)}")
37
38  # 2. DATA PREPARATION
39  print("\n2. Preparing medical physics dataset...")
40
41  # Medical physics Q&A data (expanded)
42  medical_physics_data = [
43      {
44          "input": "What is the half-life of Technetium-99m?",
45          "output": "The half-life of Technetium-99m (Tc-99m) is approximately 6.
            01 hours. This short half-life makes it ideal for diagnostic nuclear
            medicine procedures as it provides sufficient time for imaging while
            minimizing radiation dose to patients."
46      },
47      {
48          "input": "Explain the inverse square law in radiation physics.",
49          "output": "The inverse square law states that the intensity of
            radiation is inversely proportional to the square of the distance from
            the source. Mathematically: I = I₀/(d²), where I is intensity at
            distance d, and I₀ is initial intensity. This means doubling the
            distance reduces intensity to 1/4 of the original value."
50      },
51      {
52          "input": "What are the main components of a linear accelerator?",
53          "output": "A linear accelerator (LINAC) consists of: 1) Electron gun
            for generating electrons, 2) RF waveguide for accelerating electrons,
            3) Target/flattening filter for X-ray production, 4) Primary collimator
            for beam shaping, 5) Monitor chambers for dose monitoring, 6)
            Multi-leaf collimator (MLC) for beam conforming, and 7) Patient
            positioning systems."
54      },
55      {
56          "input": "What is IMRT and how does it work?",
57          "output": "Intensity Modulated Radiation Therapy (IMRT) is an advanced
            radiation therapy technique that uses multiple radiation beams with
            varying intensities. It optimizes dose distribution by modulating beam
            intensity across the treatment field, allowing precise dose delivery to
            tumors while sparing healthy tissues. IMRT uses inverse treatment
            planning and multi-leaf collimators."
58      },
59      {
60          "input": "Describe the photoelectric effect in medical imaging.",
61          "output": "The photoelectric effect occurs when an X-ray photon
```

```
            transfers all its energy to an inner shell electron, ejecting it from
            the atom. This interaction is energy-dependent (proportional to Z³/E³)
            and provides high contrast in medical imaging. It's the dominant
            interaction in diagnostic radiology for low-energy X-rays and high-Z
            materials like contrast agents."
62      },
63      {
64          "input": "What is the difference between absorbed dose, equivalent
            dose, and effective dose?",
65          "output": "Absorbed dose (Gray, Gy) measures energy deposited per unit
            mass. Equivalent dose (Sievert, Sv) accounts for biological
            effectiveness using radiation weighting factors (absorbed dose × wR).
            Effective dose (Sv) considers organ sensitivity using tissue weighting
            factors (equivalent dose × wT), representing overall health risk to the
            body."
66      },
67      # Add more samples...
68  ]
69
70  # Format data for instruction following
71  formatted_data = []
72  for item in medical_physics_data:
73      text = f"""Below is an instruction related to medical physics. Write a
        response that appropriately completes the request.

74
75  ### Instruction:
76  {item['input']}
77
78  ### Response:
79  {item['output']}"""
80      formatted_data.append({"text": text})
81
82  # Split dataset
83  train_data, val_data = train_test_split(formatted_data, test_size=0.2,
    random_state=42)
84  train_dataset = Dataset.from_list(train_data)
85  val_dataset = Dataset.from_list(val_data)
86
87  print(f"Train dataset: {len(train_dataset)} samples")
88  print(f"Validation dataset: {len(val_dataset)} samples")
89
90  # 3. MODEL LOADING
91  print("\n3. Loading model and tokenizer...")
92
93  MODEL_NAME = "unsloth/llama-3.2-8b-instruct-bnb-4bit"
94  max_seq_length = 1024
95
96  model, tokenizer = FastLanguageModel.from_pretrained(
97      model_name=MODEL_NAME,
98      max_seq_length=max_seq_length,
99      dtype=None,
100     load_in_4bit=True,
101 )
102
103 # Add LoRA adapters
104 model = FastLanguageModel.get_peft_model(
105     model,
106     r=16,
107     target_modules=["q_proj", "k_proj", "v_proj", "o_proj",
108                     "gate_proj", "up_proj", "down_proj"],
109     lora_alpha=16,
110     lora_dropout=0.0,
111     bias="none",
112     use_gradient_checkpointing="unsloth",
113     random_state=3407,
114 )
115
116 print("Model loaded successfully")
117
118 # 4. TRAINING CONFIGURATION
119 print("\n4. Configuring training...")
120
121 training_config = SFTConfig(
122     output_dir="./medical-physics-llm",
123     num_train_epochs=1,
124     per_device_train_batch_size=1,
125     gradient_accumulation_steps=4,
126     optim="adamw_8bit",
127     learning_rate=2e-4,
128     max_steps=50,   # Reduced for demo
129     lr_scheduler_type="cosine",
130     warmup_ratio=0.03,
131     logging_steps=5,
```

```python
132        save_strategy="steps",
133        save_steps=25,
134        evaluation_strategy="steps",
135        eval_steps=25,
136        bf16=torch.cuda.is_bf16_supported(),
137        max_grad_norm=0.3,
138        max_seq_length=max_seq_length,
139        packing=False,
140        dataset_text_field="text",
141        report_to=None,
142    )
143
144    # 5. TRAINING
145    print("\n5. Starting training...")
146
147    trainer = SFTTrainer(
148        model=model,
149        tokenizer=tokenizer,
150        train_dataset=train_dataset,
151        eval_dataset=val_dataset,
152        args=training_config,
153    )
154
155    # Clear cache and start training
156    torch.cuda.empty_cache()
157    train_result = trainer.train()
158
159    print(f"Training completed!")
160    print(f"Final training loss: {train_result.training_loss:.4f}")
161
162    # 6. EVALUATION
163    print("\n6. Evaluating model...")
164
165    # Test inference
166    test_prompts = [
167        "What is the photoelectric effect in medical imaging?",
168        "Explain IMRT treatment planning principles.",
169        "What are the safety principles in nuclear medicine?",
170    ]
171
172    model.eval()
173    for prompt in test_prompts:
174        formatted_prompt = f"""Below is an instruction related to medical physics.
       Write a response that appropriately completes the request.
175
176    ### Instruction:
177    {prompt}
178
179    ### Response:
180    """
181
182        inputs = tokenizer(formatted_prompt, return_tensors="pt", max_length=512,
       truncation=True)
183        with torch.no_grad():
184            outputs = model.generate(**inputs, max_new_tokens=150, temperature=0.7,
           do_sample=True)
185
186        response = tokenizer.decode(outputs[0], skip_special_tokens=True)
187        answer = response.split("### Response:")[-1].strip()
188
189        print(f"\nQ: {prompt}")
190        print(f"A: {answer}")
191        print("-" * 80)
192
193    # 7. SAVE MODEL
194    print("\n7. Saving model...")
195
196    model.save_pretrained("./final_medical_physics_model")
197    tokenizer.save_pretrained("./final_medical_physics_model")
198
199    # Save model info
200    model_info = {
201        "model_name": MODEL_NAME,
202        "training_date": datetime.now().isoformat(),
203        "training_steps": training_config.max_steps,
204        "final_loss": train_result.training_loss,
205        "dataset_size": len(train_dataset) + len(val_dataset),
206        "max_seq_length": max_seq_length
207    }
208
209    with open("./final_medical_physics_model/model_info.json", "w") as f:
210        json.dump(model_info, f, indent=2)
211
```

```
212  print("Model saved successfully!")
213  print("\n🎉 Medical Physics LLM Fine-tuning Complete!")
214  print(f"📁 Model saved in: ./final_medical_physics_model")
215  print(f"📊 Training loss: {train_result.training_loss:.4f}")
216  print(f"⏱ Total training steps: {train_result.global_step}")
```

## 10. Tips dan Best Practices

```
212  print("Model saved successfully!")
213  print("\n🎉 Medical Physics LLM Fine-tuning Complete!")
214  print(f"📁 Model saved in: ./final_medical_physics_model")
215  print(f"📊 Training loss: {train_result.training_loss:.4f}")
216  print(f"⏱ Total training steps: {train_result.global_step}")
```

## 10. Tips dan Best Practices