



INTELLIHACK 5.0

MACHINE LEARNING HACKATHON

---

## Q3 - Training process with code

---

**The Code Rushers**

---

# 01. Instructions for Running the Program

## 1.1 Installation of Dependencies

```
1 pip install torch transformers peft bitsandbytes accelerate
   datasets huggingface_hub pymupdf pandas docx
2 pip install ragas
3 pip install python-docx
```

## 1.2 Loading the Model

```
1 from transformers import AutoModelForCausalLM, AutoTokenizer
2
3 model_name = "Qwen/Qwen2.5-3B-Instruct"
4 tokenizer = AutoTokenizer.from_pretrained(model_name)
5 model = AutoModelForCausalLM.from_pretrained(model_name,
   device_map="auto", torch_dtype="auto")
```

## 1.3 Data Extraction from PDFs

```
1 import fitz # PyMuPDF
2
3 def extract_text_from_pdf(pdf_path):
4     text_data = []
5     try:
6         doc = fitz.open(pdf_path)
7         for page in doc:
8             text_data.append(page.get_text("text"))
9     except Exception as e:
10         print(f"Error: {e}")
11     return "\n".join(text_data)
```

## 1.4 Fine-Tuning the Model with LoRA

```
1 from peft import LoraConfig, get_peft_model, TaskType
2
3 config = LoraConfig(task_type=TaskType.CAUSAL_LM, r=8,
   lora_alpha=32, lora_dropout=0.05)
4 model = get_peft_model(model, config)
```

---

## 1.5 Training Execution

```
1 from transformers import TrainingArguments, Trainer
2
3 training_args = TrainingArguments(
4     output_dir="./results",
5     num_train_epochs=3,
6     per_device_train_batch_size=8,
7     save_steps=10_000,
8     save_total_limit=2,
9     evaluation_strategy="epoch",
10    learning_rate=5e-5,
11    weight_decay=0.01,
12    logging_dir="./logs",
13 )
14
15 trainer = Trainer(
16     model=model,
17     args=training_args,
18     train_dataset=dataset,
19     eval_dataset=eval_dataset,
20 )
21
22 trainer.train()
```