

MACHINE LEARNING

Running GenAI on Intel AI laptops and Simple LLM Inference on CPU and fine-tuning of LLM models using Intel OpenVINO

Abin Sabu Philip, Abinand M.

Saintgits Group of Institutions, Kottayam, Kerala

Abstract: This project is designed to help newcomers get started in the world of Generative Artificial Intelligence (GenAI) by engaging in a series of hands on activities. Participants will gain foundational knowledge in GenAI, perform inference on a text generation, Large Language Model (LLM) using the TinyLlama model, convert the model to ONNX format, optimize it with the Intel OpenVINO toolkit and evaluate the inference of the optimized model. The main goal of this project is to generate a chatbot by segmenting the input text into tokens, generating text from these tokens using sampling approaches. The goal is to get people better acquainted with GenAI and implementing improvements in a model leading to the creation of an interactive chatbot. This report discusses the method and execution process emphasizing the benefits and constraints of utilizing TinyLlama along with the complexities involved in converting and optimizing models.

Keywords: Generative AI (GenAI), Large Language Models (LLMs), Intel OpenVINO Toolkit, Model Optimization, Inference, Model Fine-Tuning, Hugging Face Transformers, Custom Chatbot Development, TinyLlama.

1 Introduction

Given the possibility and the speed at which Artificial Intelligence (AI) is being developed, new possibilities have been made available especially new advances in the Natural Language Processing (NLP). Software applications such as Generative AI (GenAI), and

Large Language Models such as GPT-3 the BERT, and other related models can all explain and even generate human-like texts. However, the Flow models are problematic in terms of theoretical computation where efficiency and costs are a major issue.

With the rapid advancement of Artificial Intelligence (AI), new opportunities have emerged, particularly in the realm of Natural Language Processing (NLP). Generative AI (GenAI) applications and Large Language Models (LLMs) like GPT-3 and BERT have demonstrated remarkable capabilities in generating and understanding human-like text. However, these powerful models often come with significant computational demands, making them expensive and resource-intensive to deploy.

This project explores the potential of using the TinyLlama model, a compact and efficient LLM that offers an excellent balance between performance and resource consumption for inference tasks on CPUs. While GPT-3 and similar models deliver powerful outputs, they require high-performance CPUs that are costly to acquire and maintain. In contrast, TinyLlama provides a more accessible approach to experimenting with GenAI technology, allowing users to gain valuable insights without the need for expensive infrastructure.

OpenVINO (Open Visual Inference and Neural Network Optimization) plays a crucial role in enhancing the efficiency of LLMs, enabling their easy deployment on Intel AI laptops. The primary goal of this project is to develop a straightforward solution for executing LLM inference on CPUs, particularly for applications like chatbots. By combining the advanced NLP capabilities of Hugging Face libraries with the optimization tools provided by OpenVINO, this project aims to deliver efficient and cost-effective AI applications.

The report outlines the step-by-step processes undertaken to achieve these goals, including environment setup, model optimization, and deployment strategies. Additionally, it showcases the practical application of these fine-tuned models through the creation of a custom chatbot, demonstrating the potential of integrating GenAI techniques with Intel's advanced optimization tools.

2 Libraries Used

In the project for various tasks, following packages are used.

```
torch
transformers
openvino
optimum[openvino]
nncf
gradio
psutil
numpy
```

3 Methodology

The methodology of this project involves several key steps to ensure a comprehensive understanding and practical experience in GenAI.

To start with, the model needed to be chosen. The reason for choosing the TinyLlama model was its suitability for environments with low hardware support. We chose TinyLlama as our model because of its excellent computable performance on less powerful machines. Then we had to test tinyLlama to see how well it works in producing text.

After testing TinyLlama, we converted it into a format called ONNX. This involved taking the TinyLlama model that was already trained and using a tool called PyTorch to change it into ONNX format. This new format was checked to make sure everything was right, and it gave us a file called model.onnx that we could work with.

The next action is to optimize the ONNX model using the Intel OpenVINO toolkit. OpenVINO presents a range of tools that streamline the model inference on Intel hardware. The optimization process entails converting ONNX model into OpenVINO's intermediate representation (IR) format that will produce two files model.xml and model.bin. These files have been optimized for efficient inference by considering Intel hardware acceleration capabilities.

Having our optimized model in place, we divided the input text into smaller entities referred to as tokens. Each token was given an ID that could be understood by the model. That is when these token IDs are used to make textual material through numerous sampling methods. Temperature and sampling controls the randomness of the predictions, while top-k and top-p sampling limit the model's predictions to the most probable tokens, enhancing the quality and coherence of the generated text.

Finally, we put our optimized model into a chatbot setup. We made a user interface using a tool called Streamlit, implementing backend logic to handle user inputs, tokenize the inputs, generate responses using the optimized model, and display the responses in the chatbot interface. The performance of the chatbot undergoes evaluation and improvement so that optimal results can be achieved.

The process flow of the project is shown below:

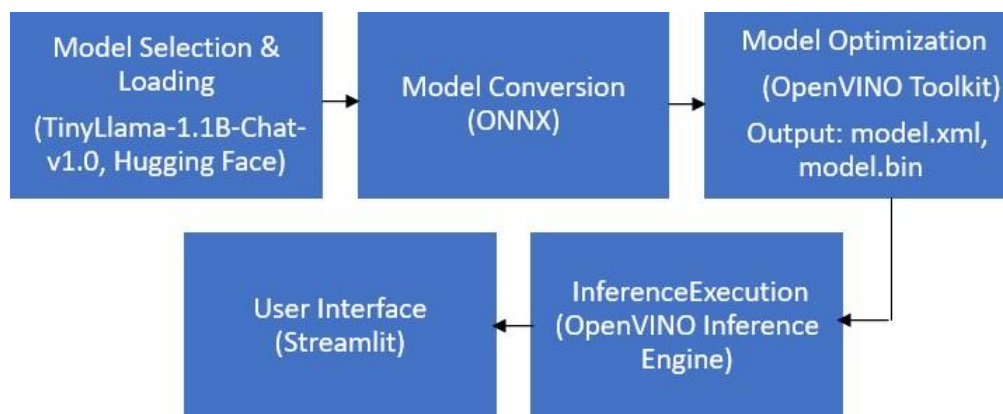


Figure 1: Process Flow

4 Implementation

We implemented our project in a step by step manner:

4.1 Model Selection and Initial Inference

We began by selecting a Large Language Model (LLM) using Hugging Face Transformers, specifically TinyLlama-1.1B-Chat-v1.0, due to its small size and lower computational requirements. This model has approximately 1.2 billion parameters, making it a balanced choice between performance and accessibility. Initial tests were conducted to assess the model's performance by feeding in sample input texts and analyzing the outputs. These tests helped us understand the model's strengths and weaknesses for future improvements.

4.2 Model Conversion to ONNX Format

At the VS code software, we ran this model and then converted this file into ONNX initiating a model conversion. The process flow had some space requirements, ie this model occupies around 2GB which was difficult. The pre-trained TinyLlama model was loaded and by using PyTorch's `torch.onnx.export` function, the model was exported into ONNX format. A check was made on integrity of resulting model.onnx file thereby confirming whether the model is ready for Intel OpenVINO optimization.

4.3 Model Optimization with Intel OpenVINO

The ONNX model was improved using the Intel OpenVINO toolkit. Therefore, the Model Optimizer transformed this to the IR format of OpenVino which gave us model.xml and model.bin files. We did that so as to make them work as fast as possible on Intel hardware by leveraging OpenVino for higher speed and efficiency. This will enable faster and smoother processing.

4.4 Tokenization and Text Generation

We applied a Tokenizer designed for TinyLlama model on input text in order to break it into smaller pieces called tokens. Each token was then given an ID number which was different from all other tokens that were used in making up the sentences. For this reason, we implemented several sampling methods in order to generate meaningful and high-quality answers such as adjusting the temperature to 0.7 and using top-k and top-p sampling. These techniques helped us generate coherent and natural-sounding text from the tokens.

4.5 Chatbot Development and Integration

The optimized model was integrated into a chatbot framework developed using Gradio. The backend logic handled user inputs, tokenized the inputs, generated responses using the optimized model and displayed the responses in the chatbot interface. The chatbot's performance was evaluated and refined to ensure it provided accurate and relevant responses, demonstrating the practical applications of GenAI.

This project successfully introduced participants to the field of Generative Artificial Intelligence (GenAI), guiding them through the process of performing inference on an LLM, converting and optimizing the model, and developing a custom chatbot. By working with the TinyLlama model and Intel’s OpenVINO toolkit, participants gained practical experience in model optimization and text generation techniques. The resulting chatbot showcases the potential of GenAI in creating interactive and intelligent applications, providing a solid foundation for further exploration and development in this exciting field.

7 Acknowledgments

We would like to express our heartfelt gratitude and appreciation to Intel® Corporation for providing an opportunity to this project. First and foremost, we would like to extend our sincere thanks to our team mentors *Dr. Anju Pratap* for their invaluable guidance and constant support throughout the project. We are deeply indebted to our college Saintgits College of Engineering and Technology for providing us with the necessary resources, and sessions. We would also like to thank the industrial mentor *Mr Abhishek Nandy*, for taking time out of his busy schedules to provide us with training and for answering our queries. We extend our gratitude to all the researchers, scholars, and experts in the field of machine learning and natural language processing and artificial intelligence, whose seminal work has paved the way for our project. We acknowledge the mentors, institutional heads, and industrial mentors for their invaluable guidance and support in completing this industrial training under Intel® -Unnati Programme whose expertise and encouragement have been instrumental in shaping our work.

References

- [1] A. R. Borah, N. T N and S. Gupta, "Improved Learning Based on GenAI", *2nd International Conference on Intelligent Data Communication Technologies and Internet of Things (IDCIoT)*, Bengaluru, India, 2024
- [2] J. Qadir, "Learning 101 reloaded: Revisiting the basics for the GenAI era", *IEEE*, 2024
- [3] V. V. Zunin, "Intel OpenVINO Toolkit for Computer Vision: Object Detection and Semantic Segmentation", *International Russian Automation Conference (RusAutoCon)*, Sochi, Russian Federation, 2021
- [4] J. Zhang, Y. Zhang, M. Chu, S. Yang and T. Zu, "A LLM-Based Simulation Scenario Aided Generation Method", *IEEE 7th Information Technology and Mechatronics Engineering Conference (ITOEC)*, Chongqing, China, 2023

A Main code sections for the solution

A.1 Inference

```

1  # Imports and Setup
2  import torch
3  from transformers import LlamaTokenizer, LlamaForCausalLM
4  import os
5  import time
6  import psutil
7  import openvino.runtime as ov
8  from openvino.runtime import Core
9  import numpy as np
10 import matplotlib.pyplot as plt
11
12 # Initialize model and tokenizer
13 model_name = 'TinyLlama/TinyLlama-1.1B-Chat-v1.0'
14 tokenizer = LlamaTokenizer.from_pretrained(model_name)
15 model = LlamaForCausalLM.from_pretrained(model_name)
16
17 # Load quantized model
18 output_dir = "/content/drive/MyDrive/openvino/chatbot/FINAL"
19 quantized_ir_path = os.path.join(output_dir, "tinyllama_quantized.xml")
20 core = Core()
21 quantized_model = core.read_model(quantized_ir_path)
22 quantized_compiled_model = core.compile_model(quantized_model, "CPU")
23
24 # Define utility functions
25 def get_file_size(file_path):
26     return os.path.getsize(file_path) / (1024 ** 2) # Size in MB
27
28 def benchmark_model(model, input_text, tokenizer, num_runs=10):
29     inputs = tokenizer(input_text, return_tensors="pt")
30     start_time = time.time()
31     for _ in range(num_runs):
32         outputs = model(**inputs)
33     end_time = time.time()
34     avg_time = (end_time - start_time) / num_runs
35     return avg_time
36
37 def get_memory_usage():
38     process = psutil.Process(os.getpid())
39     return process.memory_info().rss / (1024 ** 2) # Convert to MB
40
41 def benchmark_quantized_model(compiled_model, input_text, tokenizer, num_runs=10):
42     inputs = tokenizer(input_text, return_tensors="pt")
43     input_ids = inputs["input_ids"].detach().cpu().numpy()
44
45     # Get the expected input shape from the model
46     input_shape = compiled_model.input(0).shape
47
48     # Pad or truncate the input to match the expected shape
49     if input_ids.shape[1] < input_shape[1]:
50         pad_length = input_shape[1] - input_ids.shape[1]
51         input_ids = np.pad(input_ids, ((0, 0), (0, pad_length)), mode='constant', constant_values=tokenizer.pad_token_id)
52     elif input_ids.shape[1] > input_shape[1]:
53         input_ids = input_ids[:, :input_shape[1]]
54
55     start_time = time.time()
56     for _ in range(num_runs):
57         infer_request = compiled_model.create_infer_request()
58         infer_request.infer(inputs={"input_ids": input_ids})
59     end_time = time.time()
60     avg_time = (end_time - start_time) / num_runs
61     return avg_time

```



```

63 # Calculate model sizes
64 original_model_size = sum(p.numel() * p.element_size() for p in model.parameters()) / (1024 ** 2)
65 quantized_model_size = get_file_size(quantized_ir_path) + get_file_size(quantized_ir_path.replace('.xml', '.bin'))
66
67 print(f"Original model size: {original_model_size:.2f} MB")
68 print(f"Quantized model size: {quantized_model_size:.2f} MB")
69
70 # Cell 6: Benchmark models
71 input_text = "India is Known for"
72 original_model_time = benchmark_model(model, input_text, tokenizer)
73 original_model_memory = get_memory_usage()
74
75 # Get the expected input shape from the quantized model
76 input_shape = quantized_compiled_model.input(0).shape
77 print(f"Expected input shape for quantized model: {input_shape}")
78
79 quantized_model_time = benchmark_quantized_model(quantized_compiled_model, input_text, tokenizer)
80 quantized_model_memory = get_memory_usage()
81
82 # Print comparison results
83 print(f"Original model - Average Inference Time: {original_model_time:.4f} seconds")
84 print(f"Quantized model - Average Inference Time: {quantized_model_time:.4f} seconds")
85 print(f"Original model - Peak Memory Usage: {original_model_memory:.2f} MB")
86 print(f"Quantized model - Peak Memory Usage: {quantized_model_memory:.2f} MB")
87 print(f"Original model size: {original_model_size:.2f} MB")
88 print(f"Quantized model size: {quantized_model_size:.2f} MB")
89
90
91 # Sample data
92 models = ['Original', 'Quantized']
93 inference_times = [original_model_time, quantized_model_time]
94 model_sizes = [original_model_size, quantized_model_size]
95
96 # Plot Inference Time
97 fig, ax1 = plt.subplots(figsize=(8, 6))
98 ax1.bar(models, inference_times, color=['blue', 'orange'])
99 ax1.set_ylabel('Inference Time (seconds)')
100 ax1.set_title('Inference Time Comparison')
101 for i, v in enumerate(inference_times):
102     ax1.text(i, v, f'{v:.4f}', ha='center', va='bottom')
103
104 # Adjust layout and save the plot
105 plt.tight_layout()
106 plt.savefig('inference_FINAL.png')
107 plt.close(fig) # Close the figure to avoid overlap
108
109 # Plot Model Size
110 fig, ax2 = plt.subplots(figsize=(8, 6))
111 ax2.bar(models, model_sizes, color=['blue', 'orange'])
112 ax2.set_ylabel('Model Size (MB)')
113 ax2.set_title('Model Size Comparison')
114 for i, v in enumerate(model_sizes):
115     ax2.text(i, v, f'{v:.2f}', ha='center', va='bottom')
116
117 # Adjust layout and save the plot
118 plt.tight_layout()
119 plt.savefig('model_size_FINAL.png')
120 plt.close(fig) # Close the figure to avoid overlap
121
122 print("Graphs saved successfully.")

```

A.2 Conversion and optimisation


```

1 import torch
2 from transformers import LlamaTokenizer, LlamaForCausalLM
3 import os
4 import openvino.runtime as ov
5 from nncf import quantize
6 from nncf.quantization import QuantizationPreset
7 from torch.utils.data import Dataset, DataLoader
8
9 # Initialize model and tokenizer
10 model_name = 'TinyLlama/TinyLlama-1.1B-Chat-v1.0'
11 tokenizer = LlamaTokenizer.from_pretrained(model_name)
12 model = LlamaForCausalLM.from_pretrained(model_name)
13
14 # Create dummy input
15 dummy_input = tokenizer("Hello, how are you?", return_tensors="pt")
16 print("Input shape:", dummy_input['input_ids'].shape)
17
18 # Specify output directory
19 output_dir = "FINAL"
20 os.makedirs(output_dir, exist_ok=True)
21
22 # Specify the output file path for ONNX
23 onnx_path = os.path.join(output_dir, "tinylama.onnx")
24
25 # Export the model to ONNX format
26 torch.onnx.export(
27     model,
28     (dummy_input['input_ids'],),
29     onnx_path,
30     opset_version=14,
31     input_names=["input_ids"],
32     output_names=["output"]
33 )
34
35 # Load the ONNX model and convert to OpenVINO IR
36 core = ov.Core()
37 onnx_model = core.read_model(onnx_path)
38
39 # Specify the output paths for the OpenVINO IR format files
40 ir_xml_path = os.path.join(output_dir, "tinylama_ir.xml")
41 ir_bin_path = os.path.join(output_dir, "tinylama_ir.bin")
42
43 # Convert the ONNX model to OpenVINO IR format
44 ov.serialize(onnx_model, ir_xml_path, ir_bin_path)
45
46 print("Model has been successfully converted to OpenVINO IR format and saved to", output_dir)
47
48 # Create a dummy calibration dataset
49 class CalibrationDataset(Dataset):
50     def __init__(self, num_samples, batch_size=1):
51         self.num_samples = num_samples
52         self.batch_size = batch_size
53
54     def __len__(self):
55         return self.num_samples
56
57     def __getitem__(self, item):
58         # Generate a tensor with the correct shape [1, 7]
59         input_ids = torch.randint(0, 100, (1, 7), dtype=torch.int32) # Ensure the shape is [1, 7]
60         return {"input_ids": input_ids}

```

www.saintgits.org

```
56 # Sample text generation prompts
57 sample_prompts = [
58     "Write an e-mail",
59     "Advantages of tinyllama",
60     "Write a poem",
61     "Write a code for :",
62     "Best movies of all time ?"
63 ]
64
65 # Define Gradio interface
66 iface = gr.Interface(
67     fn=generate_response,
68     inputs=gr.Textbox(lines=2, placeholder="Enter your queries "),
69     outputs=gr.Textbox(label="Reply"),
70     title="Chatbot using TinyLlama and Intel OpenVINO ",
71     description="Welcome to Chatbot! Enter your questions ",
72     examples=sample_prompts,
73     theme="default",
74     allow_flagging="never"
75 )
76
77 # Launch the Gradio interface
78 try:
79     logger.info("Launching Gradio interface...")
80     iface.launch(share=True)
81     logger.info("Gradio interface launched successfully")
82 except Exception as e:
83     logger.error(f"Error launching Gradio interface: {e}")
```

B Github Repository

<https://github.com/ABINSABUPHILIP/TECHNOWIZZ>