

summary_llm_int8

Here is a summary of the key points from the PDF:

1. The paper introduces LLM.int8(), a new method for 8-bit quantization of large language models that enables inference without performance degradation for models up to 175B parameters.
2. Key components of LLM.int8():
 - Vector-wise quantization: Uses separate quantization constants for each row/column vector
 - Mixed-precision decomposition: Handles outlier feature dimensions in 16-bit while keeping 99.9% of values in 8-bit
3. The method allows loading 16/32-bit checkpoints, converting to 8-bit, and using immediately without performance loss. This cuts memory usage by half compared to 16-bit.
4. Analysis shows large magnitude "outlier features" emerge in transformers beyond 6.7B parameters, which disrupt normal quantization methods. LLM.int8() handles these effectively.
5. Experiments show LLM.int8() maintains full 16-bit performance on language modeling and zero-shot tasks for models up to 175B parameters, while other methods degrade.
6. This enables running very large models like OPT-175B/BLOOM on a single server with consumer GPUs, making them more accessible.
7. The method is open-sourced and integrated into Hugging Face Transformers for broad use.
8. Limitations include focus only on Int8 (not FP8), evaluation up to 175B parameters, and application to inference but not training.

In summary, LLM.int8() enables efficient 8-bit inference for very large language models without performance loss, significantly improving their accessibility and deployability.

以下是該PDF的主要要點摘要，翻譯成繁體中文：

本論文介紹了LLM.int8()，這是一種新的大型語言模型8位量化方法，可以實現高達1750億參數的模型推理，而不會降低性能。

LLM.int8()的關鍵組成部分：

向量式量化：為每個行/列向量使用單獨的量化常數

混合精度分解：以16位處理異常特徵維度，同時保持99.9%的值在8位

該方法允許加載16/32位檢查點，轉換為8位，並立即使用而不會損失性能。與16位相比，這將內存使用量減半。

分析顯示，超過67億參數的變形器中會出現大幅度的"異常特徵"，這會干擾正常的量化方法。LLM.int8()能有效處理這些問題。

實驗表明，LLM.int8()在語言建模和零樣本任務上，對於高達1750億參數的模型，保持了完整的16位性能，而其他方法則會降級。

這使得在配備消費級GPU的單一伺服器上運行非常大的模型（如OPT-175B/BLOOM）成為可能，提高了它們的可訪問性。

該方法已開源，並整合到Hugging Face Transformers中，以供廣泛使用。

局限性包括僅關注Int8（不包括FP8），評估僅到1750億參數，以及僅應用於推理而非訓練。

總之，LLM.int8()實現了非常大型語言模型的高效8位推理，而不會損失性能，顯著提高了它們的可訪問性和可部署性。

Here's an example code snippet in Python that demonstrates the basic idea of LLM.int8() quantization:

```
import torch
import torch.nn as nn
import torch.nn.functional as F

# Define a simple transformer model
class TransformerModel(nn.Module):
    def __init__(self, hidden_size, num_heads, num_layers):
        super(TransformerModel, self).__init__()
        self.encoder = nn.TransformerEncoderLayer(d_model=hidden_size, nhead=num_heads,
```

```

dim_feedforward=hidden_size, dropout=0.1)
        self.decoder = nn.TransformerDecoderLayer(d_model=hidden_size, nhead=num_heads,
dim_feedforward=hidden_size, dropout=0.1)
        self.fc = nn.Linear(hidden_size, hidden_size)

    def forward(self, input_ids):
        encoder_output = self.encoder(input_ids)
        decoder_output = self.decoder(encoder_output)
        output = self.fc(decoder_output)
        return output

# Initialize the model, input, and quantization parameters
model = TransformerModel(hidden_size=1024, num_heads=8, num_layers=6)
input_ids = torch.randn(1, 128, 1024)
quantization_bits = 8
vector_wise_quantization = True
mixed_precision_decomposition = True

# Quantize the model using LLM.int8()
def llm_int8_quantization(model, input_ids, quantization_bits, vector_wise_quantization,
mixed_precision_decomposition):
    # Vector-wise quantization
    if vector_wise_quantization:
        model.encoder.self_attn.in_proj_weight =
torch.quantize_per_tensor(model.encoder.self_attn.in_proj_weight, scale=1.0, zero_point=0,
dtype=torch.qint8)
        model.encoder.self_attn.in_proj_bias =
torch.quantize_per_tensor(model.encoder.self_attn.in_proj_bias, scale=1.0, zero_point=0,
dtype=torch.qint8)

    # Mixed-precision decomposition
    if mixed_precision_decomposition:
        # Decompose the weight matrix into two parts: one for the top 2% magnitude dimensions and one
for the rest
        weight = model.fc.weight
        top_2_percent_dim = int(weight.shape[1] * 0.02)
        top_2_percent_weight = weight[:, :top_2_percent_dim]
        rest_weight = weight[:, top_2_percent_dim:]

    # Quantize the two parts separately
        top_2_percent_weight = torch.quantize_per_tensor(top_2_percent_weight, scale=1.0,
zero_point=0, dtype=torch.qint8)

```

```

rest_weight = torch.quantize_per_tensor(rest_weight, scale=1.0, zero_point=0, dtype=torch.qint8)

# Combine the two parts
model.fc.weight = torch.cat((top_2_percent_weight, rest_weight), dim=1)

# Convert the input to 8-bit
input_ids = torch.quantize_per_tensor(input_ids, scale=1.0, zero_point=0, dtype=torch.qint8)

# Run the quantized model
output = model(input_ids)

return output

# Run the quantized model
output = llm_int8_quantization(model, input_ids, quantization_bits, vector_wise_quantization,
mixed_precision_decomposition)
print(output)

```

Note that this is a highly simplified example and may not reflect the actual implementation of LLM.int8() in the paper. Additionally, this code snippet is for demonstration purposes only and may not be optimized for performance.