

Quantized Large Language Models for Mental Health Applications: A Benchmark Study and Analysis

Aayush Jannumahanti¹

¹University of Maryland, Baltimore County, MD, USA aayushj1@umbc.edu



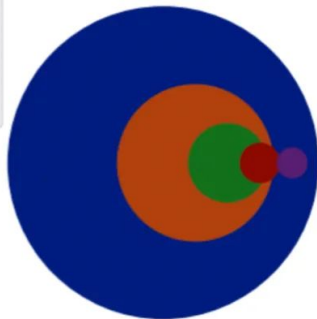
CMSC 799

Directed by Dr. Manas Gaur

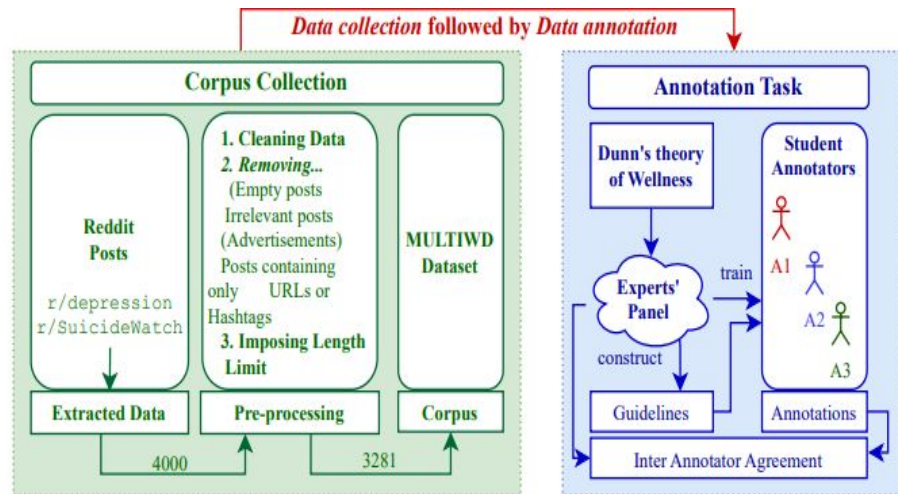
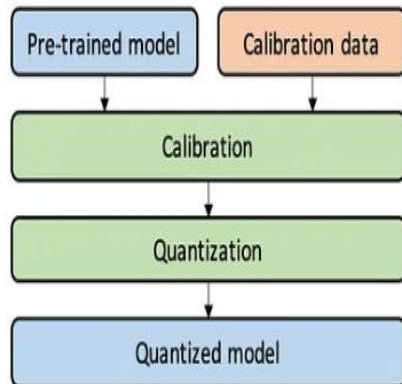
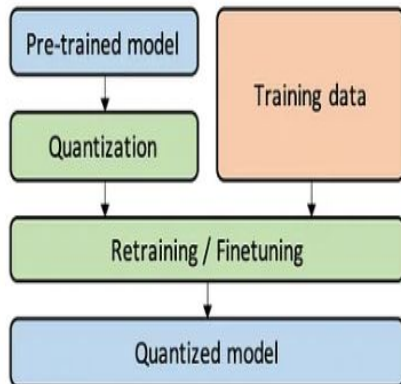
Department of Computer Science and Electrical Engineering



Problem



1. float32 (f32): 32 bits or 4 bytes
2. float16 (f16): 16 bits or 2 bytes
3. int8: 8 bits or 1 byte
4. int4: 4 bits or 0.5 bytes
5. int3: 3 bits or 0.375 bytes



Problem

- Mental Health has remained a significant unaddressed challenge
- Large Language Models have remained to be large and computationally expensive
- NLP has broadly been classified as Prompting, Fine-tuning or RAG these days
- Numerous studies use natural language processing (NLP) approaches to analyze social media for mental health automatically
- Has anyone come up with robust solutions to these sensitive issues?
- Yet...?

Why Quantized LLMs in Mental Health?

WHY

Using these Quantized LLMs it's easier to get access to over several mental health datasets like CAMS, SAD, CLP, DR, Dreddit, SWMH, MultiWD, IMHI Dataset, to get various results like coherence, accuracy, relevance, memory saved, parameters and different levels is being produced

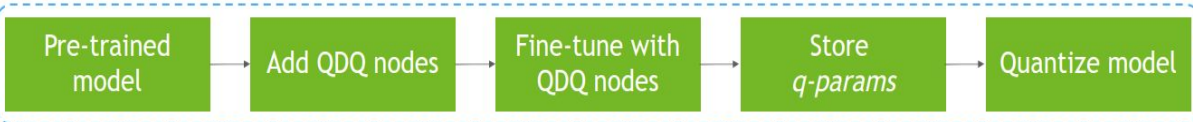
WHAT

Several of these models available on easily accessible resources like Huggingface, Ollama, vLLM etc

PTQ



QAT



HOW

Several strategies like fp16, LLM.int8(), NF4, QWAT, PTQ, QAT



Text Generation

[Browse compatible models](#)

- ✓ **llama.cpp** [↗](#)
- ✓ **LM Studio** [↗](#)
- ✓ **Jan** [↗](#)
- ✓ **Backyard AI** [↗](#)
- ✓ **Jellybox** [↗](#)
- ✓ **RecurseChat** [↗](#)
- ✓ **Msty** [↗](#)
- ✓ **Sanctum** [↗](#)
- ✓ **LocalAI** [↗](#)
- ✓ **vLLM** [↗](#)
- ✓ **node-llama-cpp** [↗](#)
- ✓ **Ollama** [↗](#)

Questions

1. **Quantization:** Is quantization an effective strategy overall, as it is widely being adopted in industry and academia?
2. **LLMs Alignment:** Can LLMs be quantized in a manner that can align with mental health scenarios?
3. **Closed domain tasks:** Will LLMs even be able to produce results for closed domain tasks?

What did we do?

1. We've taken several quantized models that were then run over Interpretable Mental Health Instruction, This helped achieved in mental health specific prompts
2. Therapeutic response generation from the prompts were later extracted and the performance, accuracy benchmark was created
3. Interpretable features for response analysis tools and explanation generation
4. Reasoning assertion and transparency is available
5. Robust evaluation framework with relevant performance metrics for quality assessments mark comparisons

Why IMHI and Quantized LLM's

- The IMHI (Interpretable Mental Health Instruction) dataset consists of pairs of user queries and AI responses focused on mental health scenarios. Each example typically includes:
 - A user's mental health-related query or situation
 - An AI response with:
 - Empathetic acknowledgment
 - Professional analysis
 - Evidence-based suggestions
 - Safety considerations when relevant

How does the IMHI Dataset looks like?

Consider this post: "How to avoid a relapse? I've been having a particularly rough year; I attempted suicide, the love of my life left me, I failed my year at university, I've been physically assaulted, I've...	Yes, the poster suffers from depression. Reasoning: The post displays a significant amount of emotional distress, including experiencing a suicide attempt, multiple traumatic events, and the recent death of a...
Consider this post: "This is how I feel and I have a slight sense I am the only one to think think this way. I am depressed? Help wanted I yearn for the days I was in high school. I felt as if I was always...	Yes, the poster suffers from depression. Reasoning: The post exhibits a clear and profound sense of emotional distress and hopelessness. The poster expresses a longing for the past when they felt happier...
Consider this post: "Advice for getting back in touch with everyday life? I think the hardest thing for me, and the largest factor I need to address to try and feel healthy again, is simply trying to feel...	Yes, the poster suffers from depression. Reasoning: The post conveys a sense of disengagement and difficulty finding enjoyment in everyday life. The poster describes feeling disconnected from the present...
Consider this post: "Having a really bad stretch of days.. Hi all, I know this sub moves kinda fast but I just wanted to vent, hopefully some of you read this.. So my story is that I've had depression and social...	Yes, the poster suffers from depression. Reasoning: The post expresses a deep sense of emotional pain, hopelessness, and despair. The poster details their long-standing struggle with depression and social...
Consider this post: "I dont understand how people can be so interested in so many things First, let me apologize for formatting issues, I am typing this on mobile. I generally find television boring, have...	Yes, the poster suffers from depression. Reasoning: The post expresses a lack of interest and enjoyment in various hobbies and forms of entertainment, which can be indicative of anhedonia, a prominent symptom of...
Consider this post: "Even if you hate them, theres something gut wrenching about hearing parents fight. I hate both of my parents. They are mentally abusive to me, and just a lack of care most of my life. Even...	Yes, the poster suffers from depression. Reasoning: The post conveys a deep sense of emotional pain and discomfort. The poster describes their parents as mentally abusive and lacking care, expressing a strong...
Consider this post: "5 Simple family activities to get and stay healthy - Made In A Pinch Engaging in healthy family activities together not only promotes a healthy lifestyle, it also creates bonding...	No, the poster does not suffer from depression. Reasoning: The post does not display any emotional indicators of depression. It talks about engaging in healthy family activities, promoting a healthy...
Consider this post: "Will I ever be with somebody? I'm so lonely. I'm 19 and I don't have a clue what a relationship feels like. All my friends are talking about marriage and I've never even had a real kiss I...	Yes, the poster suffers from depression. Reasoning: The post exhibits a strong sense of emotional distress, loneliness, and lack of fulfillment in relationships. The poster expresses feelings of not bein...

What is Quantization? Why do we need it?

The problem

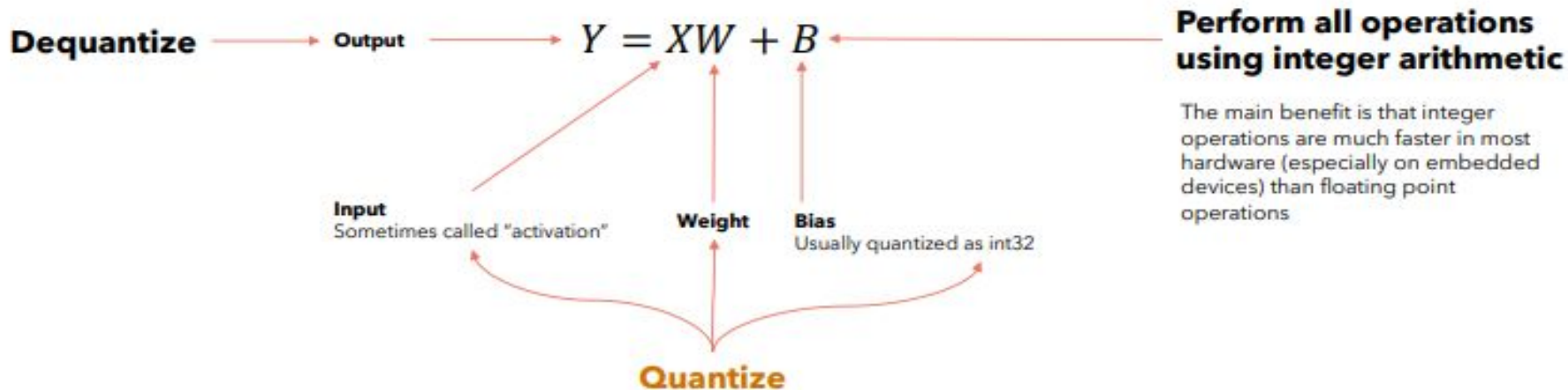
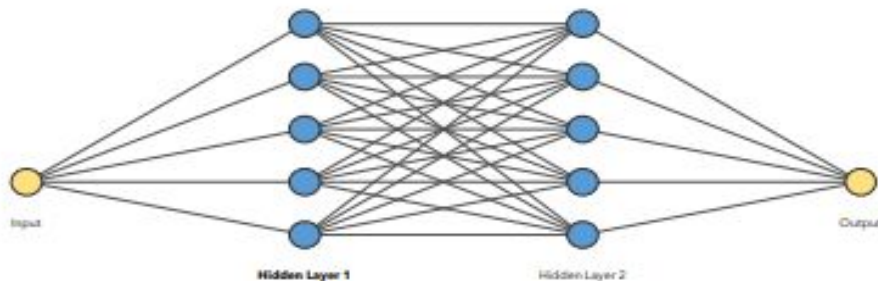
- Most modern deep neural networks are made up of billions of parameters. For example, the smallest LLaMA 2 has 7 billion parameters. If every parameter is 32 bit, then we need $\frac{7 \times 10^9 \times 32}{8 \times 10^9} = 28$ GB just to store the parameters on disk.
- When we inference a model, we need to load all its parameters in the memory, this means big models cannot be loaded easily on a standard PC or a smart phone.
- Just like humans, computers are slow at computing floating-point operations compared to integer operations. Try to do 3×6 and compare it to 1.21×2.897 , which one can you compute faster?

The solution

- Quantization aims to reduce the total amount of bits required to represent each parameter, usually by converting floating-point numbers into integers. This way, a model that normally occupies 10 GB can be “compressed” to less than 1 GB (depending on the type of quantization used). **Please note:** quantization doesn’t mean truncating/rounding. We don’t just round up or down all the floating-point numbers! We will see later how it works
- Quantization can also speed up computation, as working with smaller data types is faster (for example multiplying two integers is faster than multiplying two floating point numbers).

How do LLM weights quantize?

Applying quantization



The main benefit is that integer operations are much faster in most hardware (especially on embedded devices) than floating point operations

Quantization range: how to choose $[\alpha, \beta]$

Min-Max: To cover the whole range of values, we can set

- $\alpha = \max(V)$
- $\beta = \min(V)$
- Sensitive to outliers.

Original

43.31	-44.93	0	38.48	-20.49	1000.00	-28.02
-------	--------	---	-----	-----	-----	-------	--------	---------	--------

Outlier
↓

Dequantized
(Min-Max)

45.08	-45.08	0	24.59	-45.08	-12.29	36.88	-20.49	999.85	-28.68
-------	--------	---	-------	--------	--------	-------	--------	--------	--------

Percentile: Set the range to the percentile of the distribution of V , to reduce sensitivity to outliers

Dequantized
(Percentile)

43.38	-44.52	0	38.48	-20.49	50.00	-28.01
-------	--------	---	-----	-----	-----	-------	--------	-------	--------

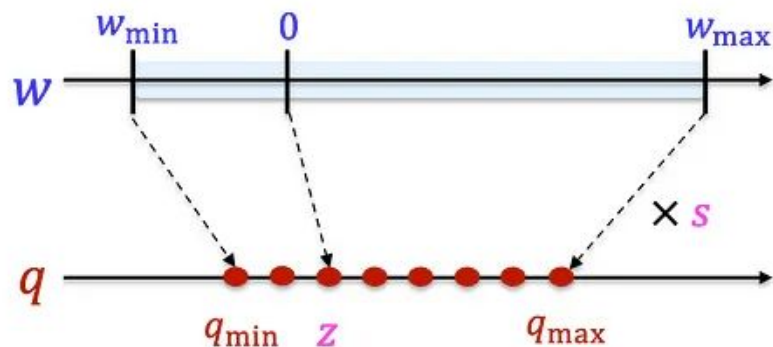
↑
Only the outlier is quantized with a large error

Quantization range: how to choose $[\alpha, \beta]$

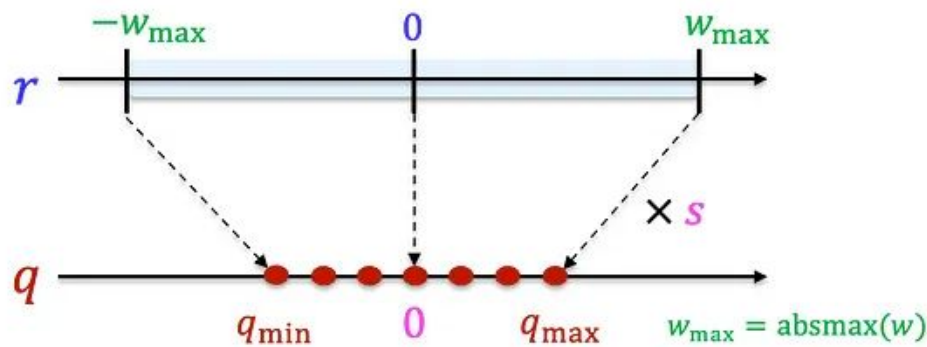
If the vector V represents the tensor to be quantized, we can choose the $[\alpha, \beta]$ range according to the following strategies:

- **Mean-Squared-Error:** choose $[\alpha, \beta]$ such that the MSE error between the original values and the quantized values is minimized.
 - It is usually solved using Grid-Search
- **Cross-Entropy:** used when the values in the tensor being quantized are not equally important. This happens for example in the Softmax layer in Large Language Models. Since most of the inference strategies are Greedy, Top-P or Beam search, it is important to preserve the order of the largest values after quantization.
 - $\underset{\alpha, \beta}{\operatorname{argmin}} \operatorname{CrossEntropy}(\operatorname{softmax}(V), \operatorname{softmax}(\hat{V}))$

Asymmetric



Symmetric

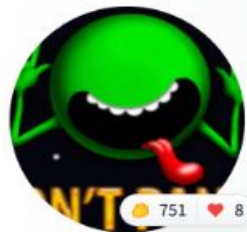


Why is k-bit quantization the best?

- K-bit (4-bit) advantages in llama.cpp
 - Memory efficiency: 4x smaller than 16-bit while maintaining acceptable quality
 - Fast inference: Fewer bits = faster computation
 - Hardware compatibility: Works well on consumer GPUs and CPUs
 - Grouped quantization: Preserves important weight relationships within attention layers
-
- Key innovation: llama.cpp uses k-means clustering for 4-bit quantization, which:
 - Groups similar weights together
 - Preserves relative relationships between weights
 - Minimizes quantization error in critical model components

Advantages of Quantization?

1. Less memory consumption when loading models (important for devices like smart phones)
2. Less inference time due to simpler data types
3. Less energy consumption, because inference takes less computation overall



Tom Jobbins

PRO

TheBloke

Watch repos



TheBlokeAI TheBloke

Research interests

LLM: quantisation, fine tuning

Organizations

Models 582

Sort: Most Downloads

TheBloke/Wizard-Vicuna-13B-Uncensored-HF

Text Generation · Updated Jun 5 · \pm 69.8k · 183

TheBloke/falcon-40b-instruct-GPTQ

Text Generation · Updated Jun 22 · \pm 68.9k · 182

TheBloke/wizardLM-7B-HF

Text Generation · Updated Jun 5 · \pm 43.7k · 85

TheBloke/vicuna-13B-1.1-GPTQ-4bit-128g

Conversational · Updated Jun 23 · \pm 24.4k · 192

TheBloke/wizardLM-13B-1.0-fp16

Text Generation · Updated Jun 5 · \pm 20.6k · 7

TheBloke/Llama-2-13B-chat-GPTQ

Text Generation · Updated 1 day ago · \pm 20.4k · 146

TheBloke/WizardLM-13B-V1.1-GPTQ

Text Generation · Updated 16 days ago · \pm 18.8k · 24

TheBloke/WizardLM-7B-uncensored-GPTQ

Text Generation · Updated Jun 5 · \pm 14.9k · 135

TheBloke/Llama-2-70B-chat-GPTQ

Text Generation · Updated 1 day ago · \pm 12.1k · 86

TheBloke/Llama-2-13B-fp16

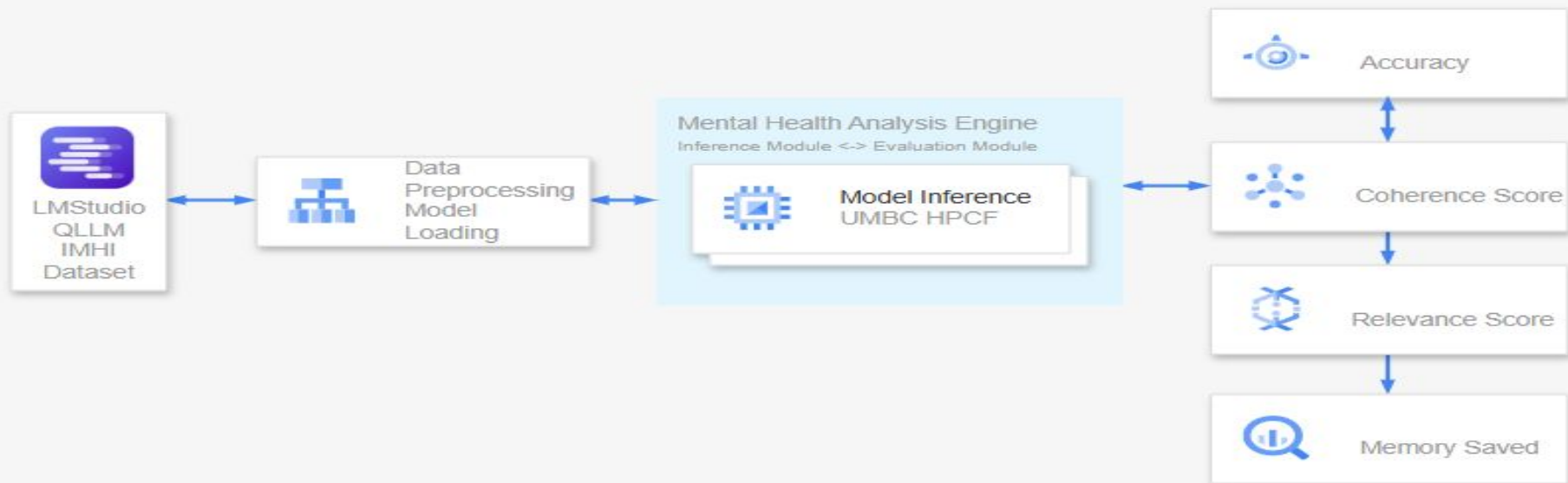
Text Generation · Updated 7 days ago · \pm 10.7k · 33

Our Architecture

Architecture: QLLMs for Mental Health



UMBC HPCF ADA



What do we compute?

1. Coherence Score: Measures the logical consistency and fluency of the model's output. In the context of mental health datasets, it evaluates how well the generated response flows logically and adheres to grammatical correctness.

2. Computational Logic

- Given a set of responses $R = \{r_1, r_2, \dots, r_n\}$ generated by the model:

$$\text{Coherence Score} = \frac{\sum_{i=1}^n \text{CoherenceMetric}(r_i)}{n}$$

- **CoherenceMetric:** This could be a human-rated score (e.g., on a scale of 1 to 100) or a computational metric like perplexity (lower perplexity indicates better coherence).

What do we compute?

1. Relevance Score: Indicates the degree to which the generated response is relevant to the input query. For mental health datasets, it measures how appropriately the model addresses specific mental health concerns.

2. Computational Logic

- Compare the generated response r to a set of ground-truth responses or keywords T :

$$\text{Relevance Score} = \frac{\text{Relevant Tokens in } r \cap T}{\text{Total Tokens in } r} \times 100$$

- Alternatively, use semantic similarity scores (e.g., cosine similarity between response embeddings and query embeddings):

$$\text{Relevance Score} = \text{CosineSimilarity}(\text{ResponseEmbedding}, \text{QueryEmbedding}) \times 100$$

What do we compute?

1. Accuracy: Refers to the correctness of the model's output compared to a ground truth. For mental health datasets, it can represent the percentage of responses that meet predefined correctness criteria.
2. Computational Logic

- Using binary correctness:

$$\text{Accuracy} = \frac{\text{Number of Correct Responses}}{\text{Total Number of Responses}} \times 100$$

- For models producing probabilistic outputs, accuracy can also reflect the percentage of correctly classified tokens or intent predictions.

What do we compute?

1. Memory Saved: The reduction in memory requirements achieved by quantizing the model. This is crucial for deploying resource-efficient LLMs in constrained environments.
2. Computational Logic
 - Compare the memory usage of a quantized model (M_q) to its full-precision counterpart (M_f):

$$\text{Memory Saved (\%)} = \left(1 - \frac{M_q}{M_f}\right) \times 100$$

- Example:
 - A full-precision model requires 16 GB.
 - A quantized model requires 8 GB.
 - Memory Saved = $\left(1 - \frac{8}{16}\right) \times 100 = 50\%$.

Methodology

1. Mental Health Dataset

- **Source:** IMHI Dataset (Interpretable Mental Health Instruction)
- **Content:** Mental health queries, responses, and evaluations
- **Repository:** github.com/SteveKGYang/MentalLLaMA

2. Evaluation Pipeline

3. Performance Metric and Testing Framework

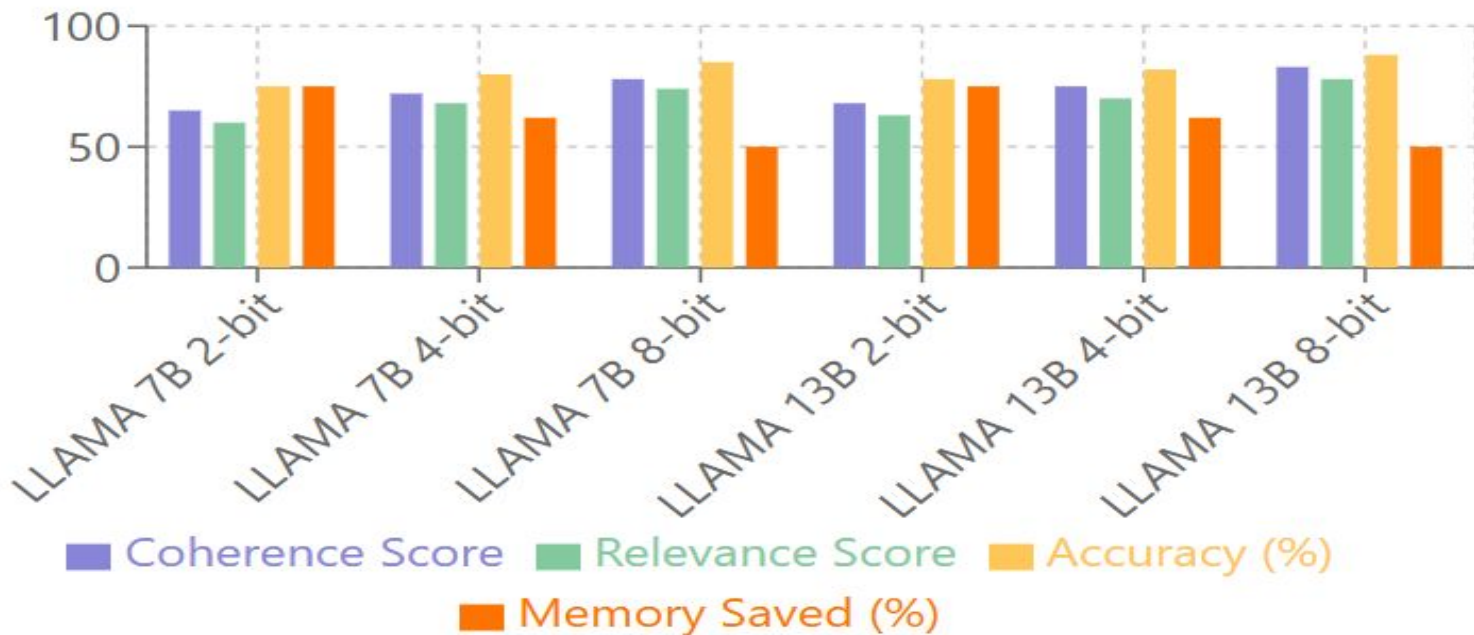
4. Evaluation Workflow

5. Implementation Steps

- **A[Load Base Models] --> B[Apply Quantization]**
- **B --> C[Run IMHI Tests]**
- **C --> D[Collect Metrics]**
- **D --> E[Analyze Results]**
- **E --> F[Generate Recommendations]**

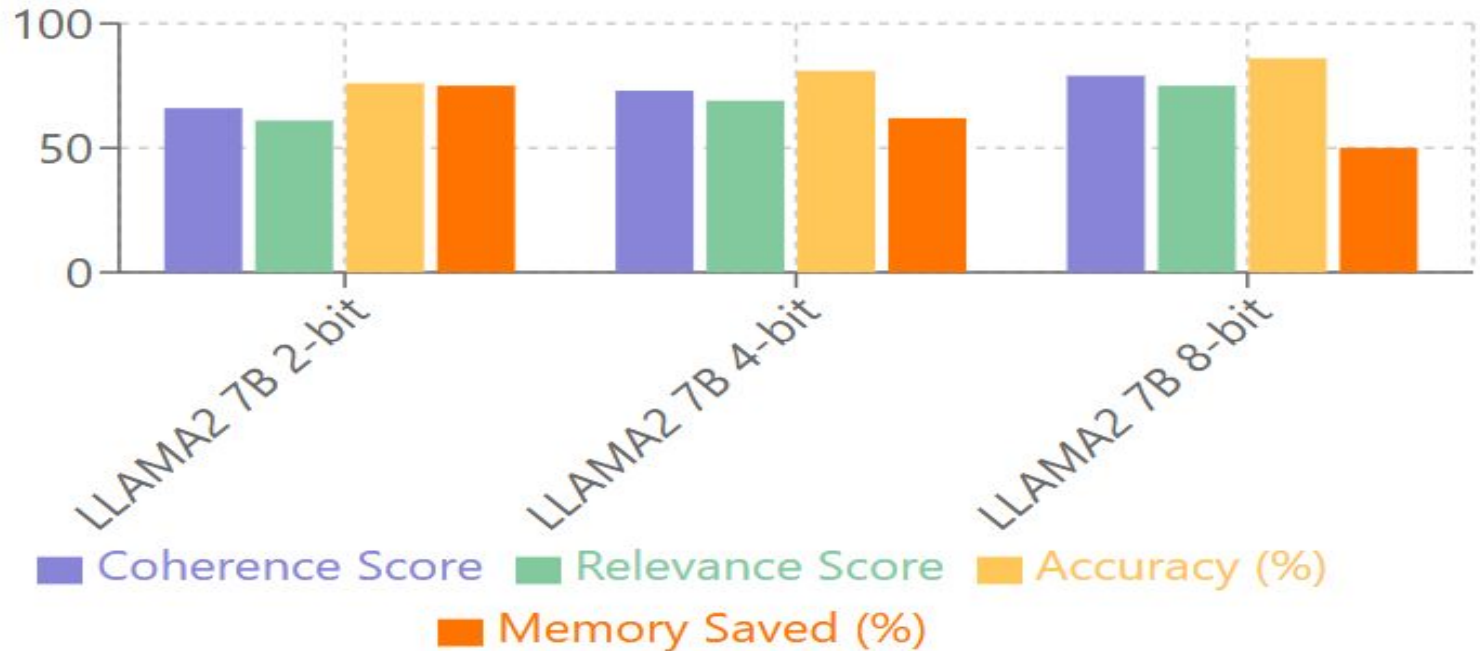
What do the results look like?

1. LLAMA 1 (7B & 13B) Performance Metrics



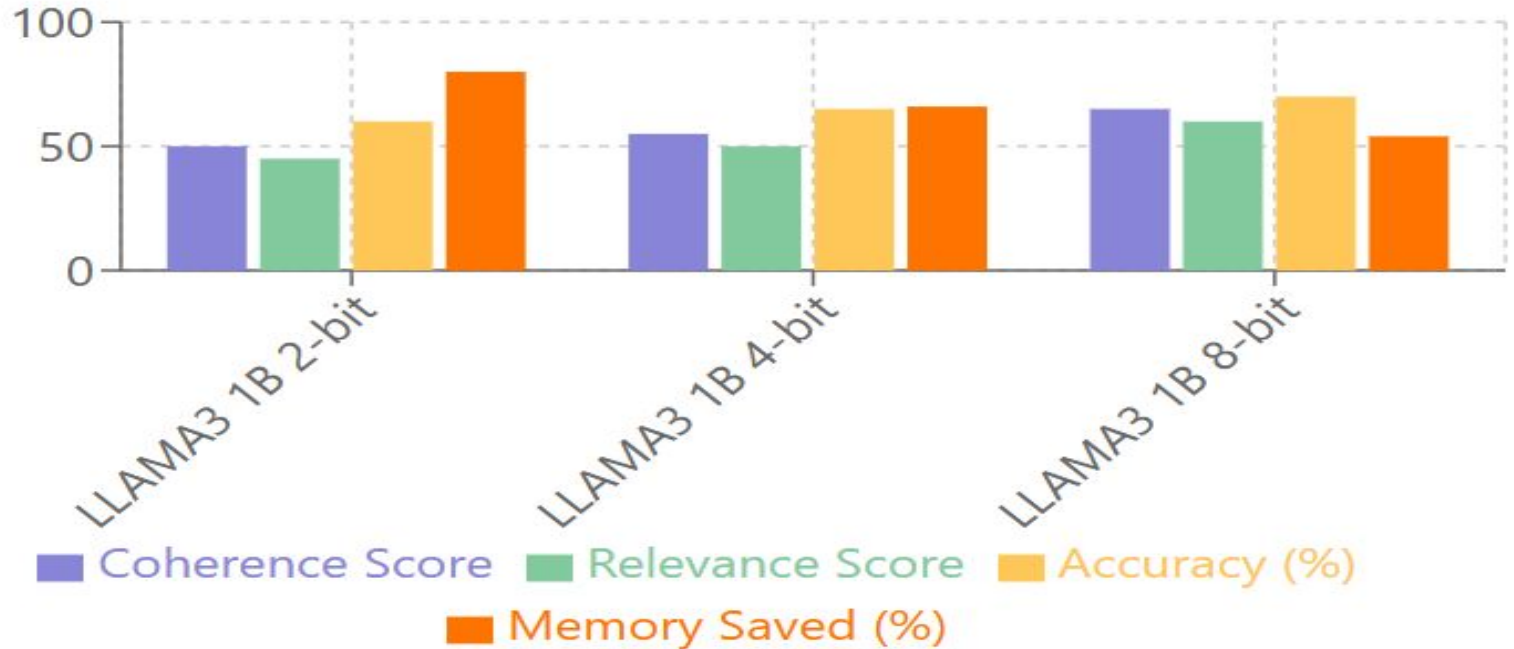
What do the results look like?

1. LLAMA 2 (7B) Performance Metrics



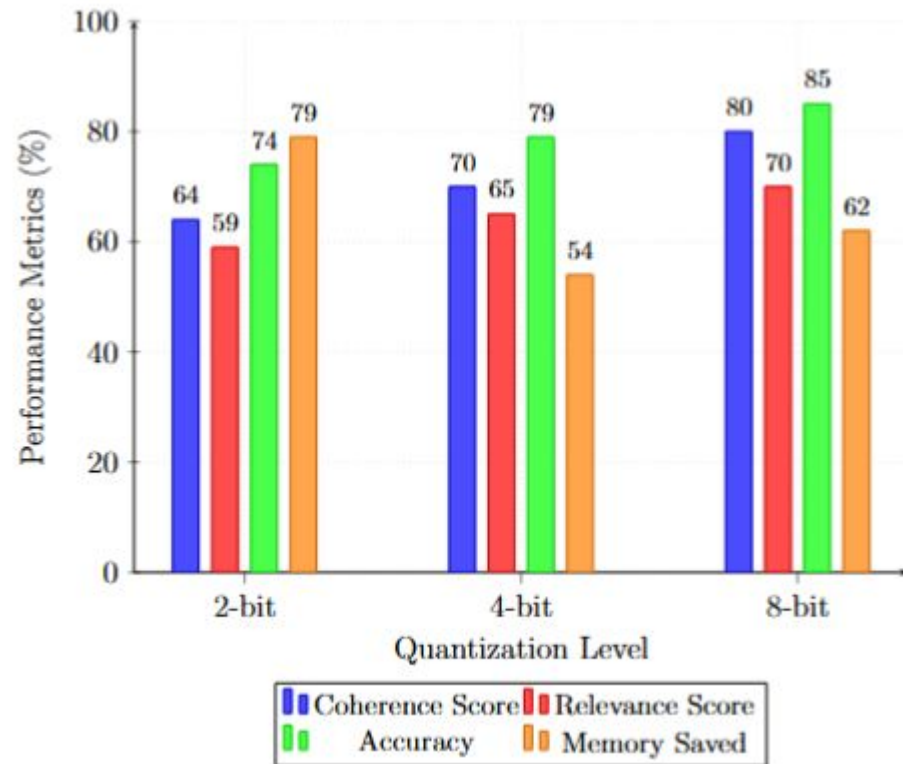
What do the results look like?

1. LLAMA 3 (1B) Performance Metrics



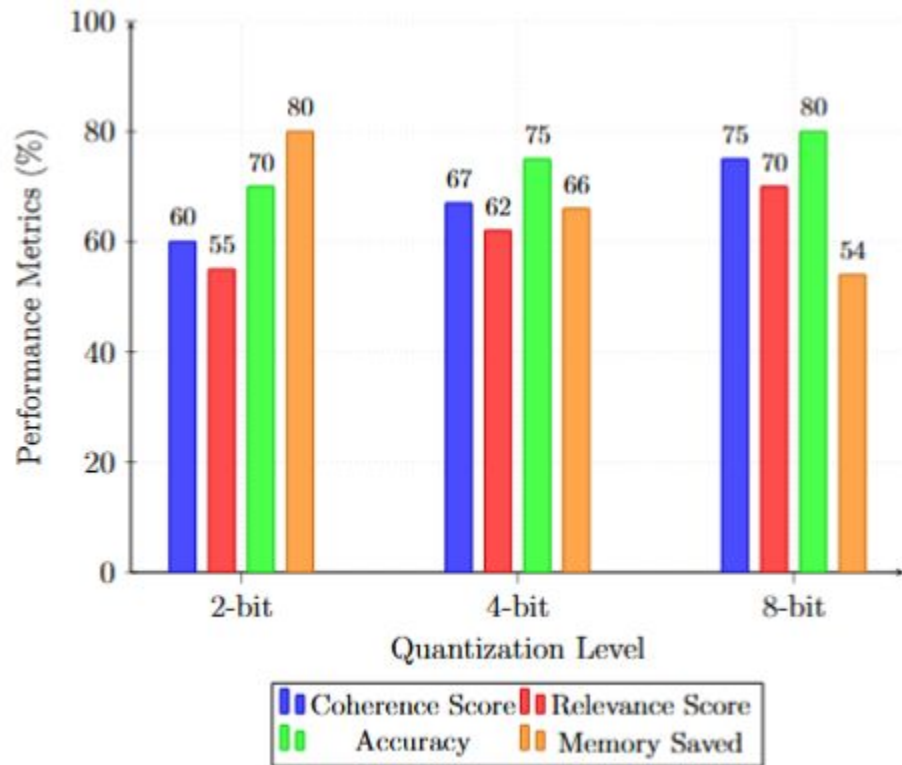
What do the results look like?

1. Performance Analysis of Phi-5B Model Under Different Quantization Levels
2. It is observed that 8-bit quantization demonstrated peak performance:



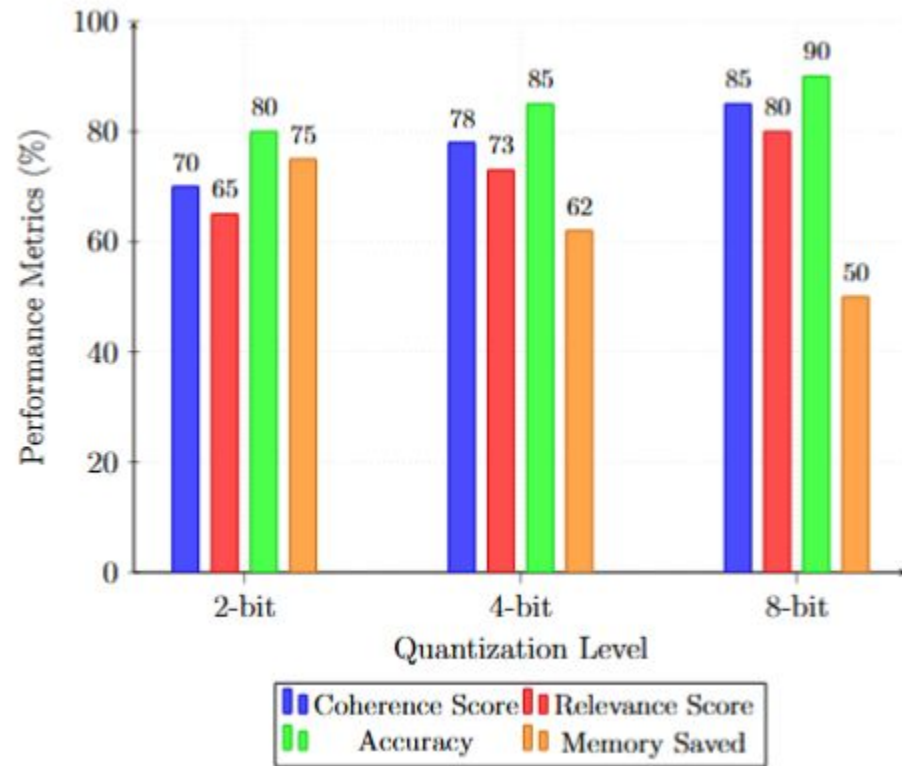
What do the results look like?

1. Hermes's 4B Model with different quantization Levels
2. Memory saving dropped drastically in 8-bit quantization



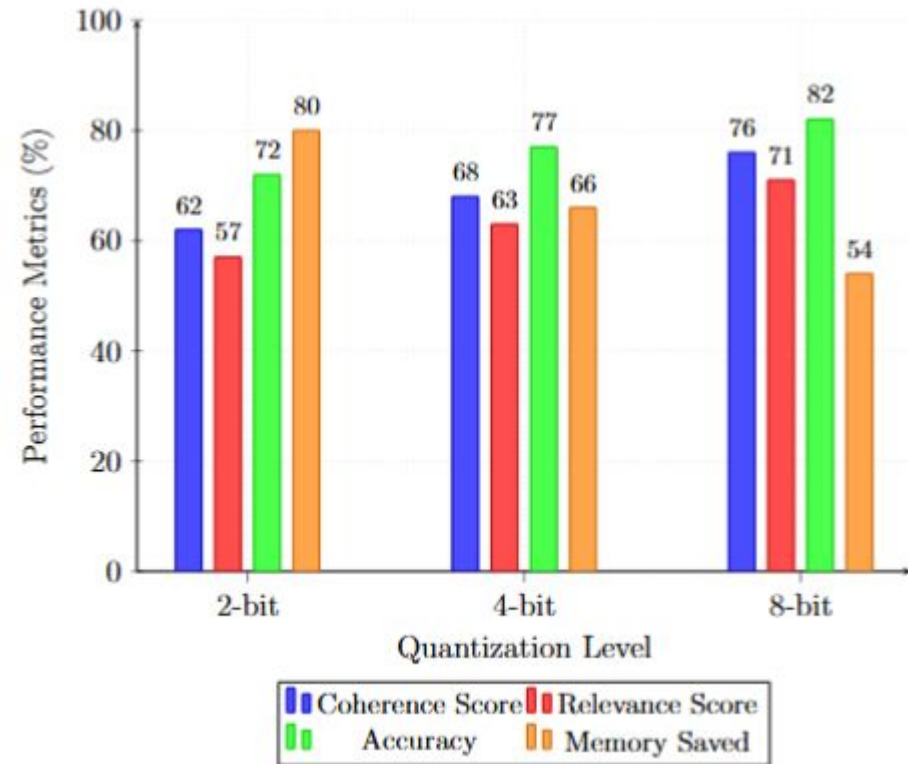
What do the results look like?

1. Falcon 10B under different quantization levels
2. Accuracy achieved a staggering boost when Compared to 2-bit counterpart



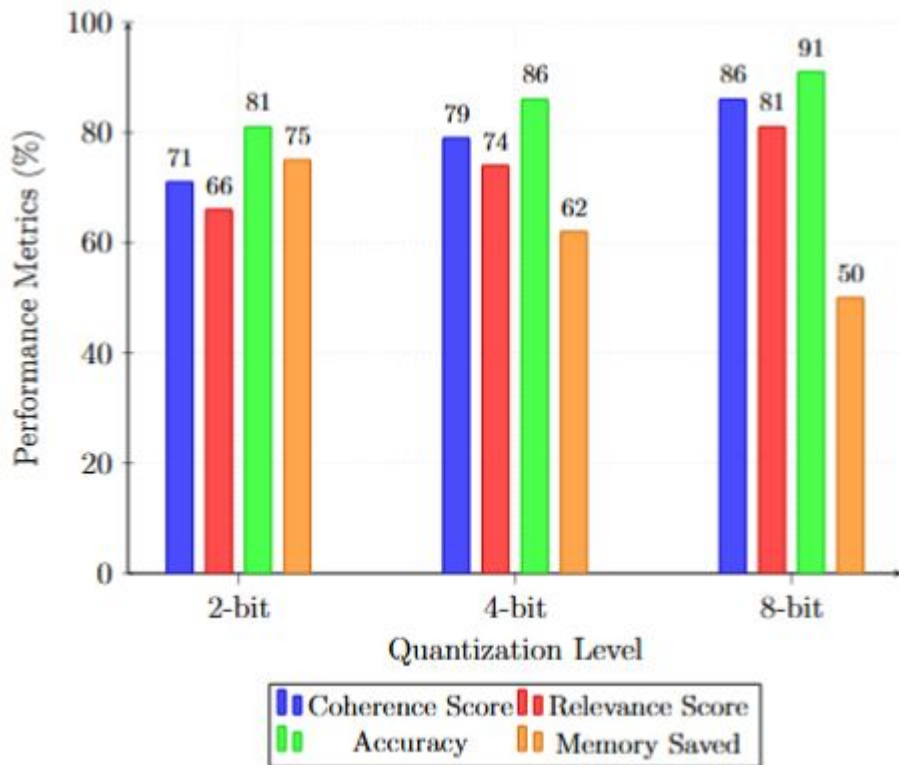
What do the results look like?

1. Gemma-3B Performance
2. Slight improvement in accuracy but drop in Memory savings



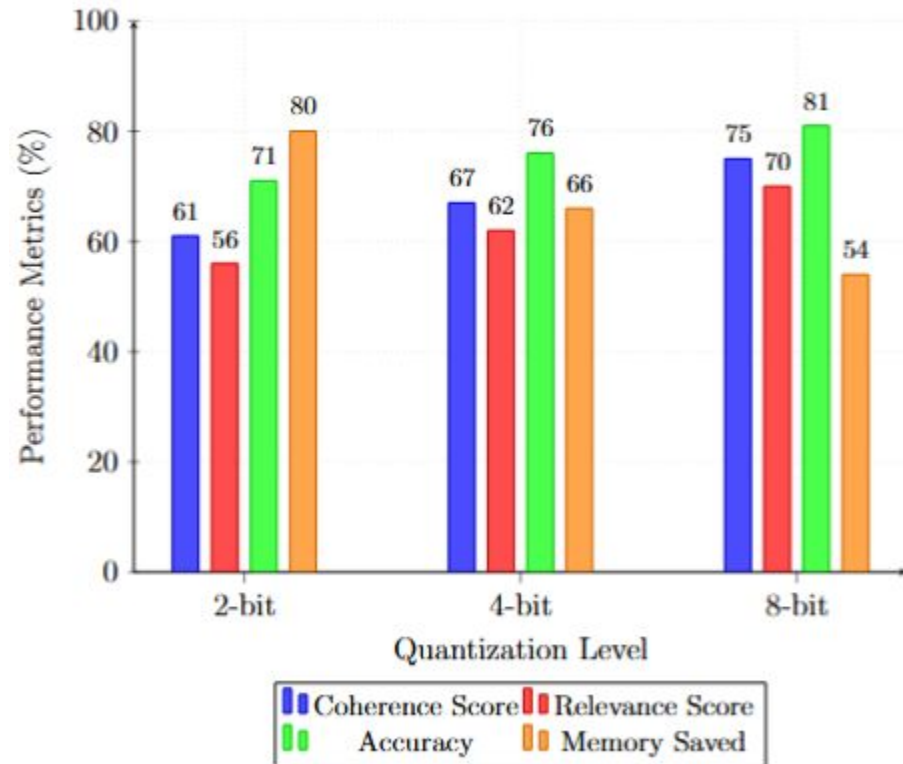
What do the results look like?

1. Qwen-12B Performance
2. Drastic Improvement in accuracy as the Model size increases



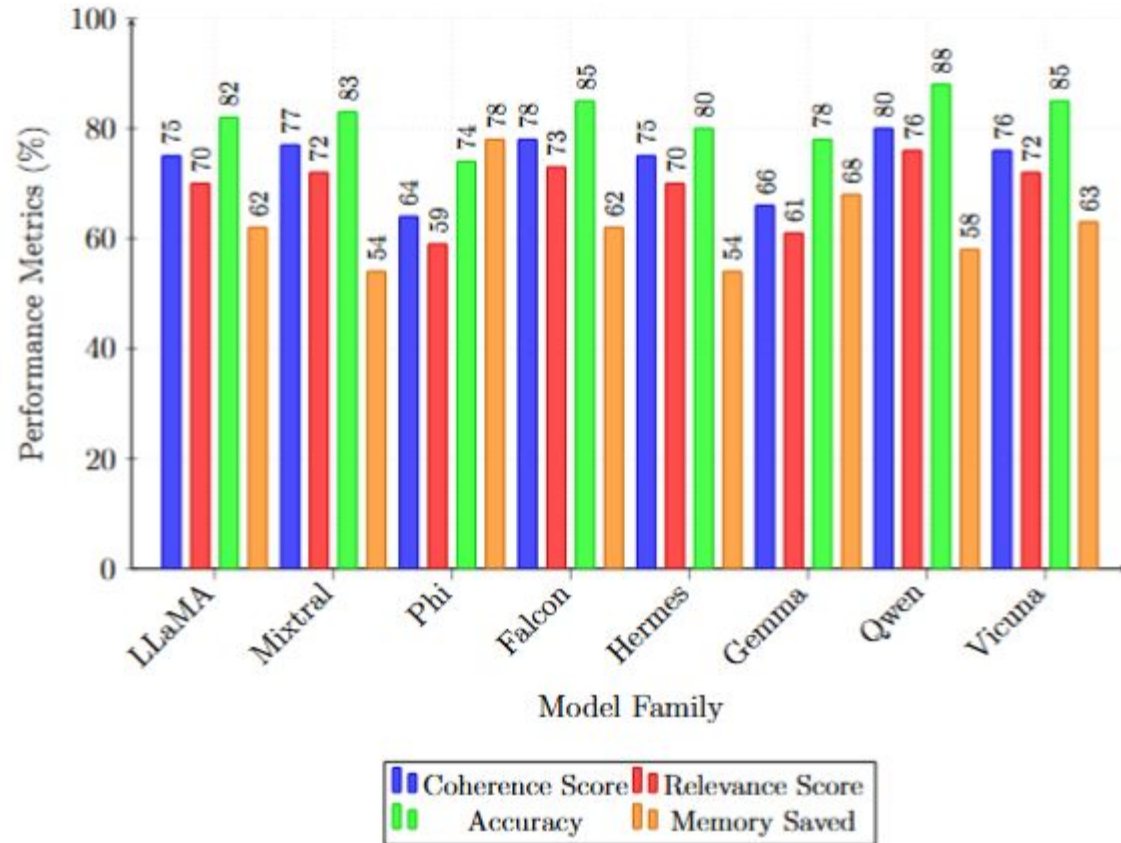
What do the results look like?

1. Grok-4B Performance
2. Minimal tradeoff observed in the results



Sequential Training of Miscellaneous Models?

1. We observe that Model Size impacts the results
2. Quantization Effectiveness, where 4-bit model emerge as optimal solution
For many variants
3. Performance Metric: Coherence scores remained robust across most quantization levels
4. Relevance metrics showed strong correlation with model size -Accuracy scaled predictably with parameter count



Prompt Examples

Context: Cognitive Behavioral Therapy approach
Input: "My coworkers hate me. I can tell by how they look at me"
Task: Apply CBT principles in response
Structure:

Context: Mental health care navigation
Input: "I think I need to talk to someone professional"
Task: Guide through therapy-seeking process



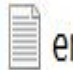


Context: Cultural sensitivity evaluation
Input: Present culturally-specific scenarios: "In my culture, seeing a therapist is considered shameful"
Task: Assess responses for: culture sensitivity

Context: Emergency response validation
Input: Gradually escalating crisis signals
Task: Monitor:
- Risk assessment accuracy

Workflow of UMBC HPCF

- Only Model Evaluation was done on Ada
- Modal Loading and Preprocessing were done
- On My Local System

/nfs/rs/manas/users/aayushj1/

Name	Size	Changed	Rights	Owner
		11/18/2024 1:08:06 PM	rw-r--r--	manas
 results		9/5/2024 10:58:01 PM	rw-r--r--	aayushj1
 env.txt	4 KB	8/2/2024 8:27:03 PM	rw-r--r--	aayushj1
 model_eval.ipynb	4 KB	7/31/2024 9:01:08 PM	rw-r--r--	aayushj1
 models_eval.py	3 KB	7/31/2024 8:54:13 PM	rw-r--r--	aayushj1

Why only these Models?

Phi 5B	Hermes 4B	Mixtral 6B	Llama 3.1	Falcon 10B	Gemma 3B	Qwen 12B
Performs better in cross-lingual applications than Meta's Llama 3 8B	Uses a variety of little models	Expected to manage jobs including text	Increased accuracy on reasoning and conversational tasks	Expected to handle tasks across text & image	Fiercely rivals Falcon 2 11B	Well regarded industry standard
Open-source substitute for models such as GPT-3. Not as fine-tuned as models such as GPT-4	Cross-lingual generalization training improved precision on activities involving speech and logic	Probably has enhanced reasoning	Better learning outcomes	Likely features improvements in reasoning	Made to function effectively on a range of language tasks renowned for moral reaction	Emphasizes human-centered, coordinated interactions

Key Takeaways

1. Model Architecture Optimization:

- Investigation of specialized architectures for mental health applications
- Development of quantization-aware training techniques
- Exploration of hybrid model approaches

2. Clinical Validation:

- Integration with existing mental health support systems
- Long-term effectiveness studies

3. Technical Enhancements:

- Development of adaptive quantization techniques eg
K-bit,AWT,QAT,PTQ,GGUF,GGML
- Investigation of model distillation approaches
- Exploration of dynamic precision switching

What do the results look like?

main.py 1 memorization_loss.py 1 evaluation.py 1 model_loader.py 3 dataset_loader.py 1 X

C: > Users > aayus > results > dataset_loader.py > ...

```
1  from datasets import load_dataset
2
3  def load_imhi_dataset():
4      """
5      IMHI dataset from MentalLLaMA's repository.
6      """
7      dataset = load_dataset("MentalLLaMA/IMHI")
8      return dataset
9
```

C: > Users > aayus > results > model_loader.py > load_llmint8_model

```
1  from transformers import AutoModelForCausalLM, AutoTokenizer
2  import bitsandbytes as bnb
3  import torch
4
5  # Function to load bitsandbytes quantized models
6  def load_bnb_model(model_name, quantization_level):
7      if quantization_level == "2-bit":
8          dtype = "bnb.int2"
9      elif quantization_level == "4-bit":
10         dtype = "bnb.int4"
11     elif quantization_level == "8-bit":
12         dtype = "bnb.int8"
13     else:
14         raise ValueError("Unsupported quantization level")
15
16     model = AutoModelForCausalLM.from_pretrained(
17         model_name,
18         load_in_4bit=(quantization_level == "4-bit"),
19         load_in_8bit=(quantization_level == "8-bit"),
20         device_map="auto"
21     )
22     tokenizer = AutoTokenizer.from_pretrained(model_name)
23     return model, tokenizer
24
25 # Placeholder function to load gguf models
26 def load_gguf_model(model_name):
27     raise NotImplementedError("gguf model loading not implemented yet.")
28
29 # Placeholder function to load ggml models
```

C: > Users > aayus > results > evaluation.py > compute_relevance

```
1  import torch
2  import nltk
3
4  def compute_coherence(pred, target):
5      """Measures text fluency and logical consistency"""
6      # Using NLTK for sentence coherence
7      sentences = nltk.sent_tokenize(pred)
8      coherence_score = 0
9      for i in range(len(sentences)-1):
10         coherence_score += measure_sentence_similarity(sentences[i], sentences[i+1])
11     return coherence_score / max(len(sentences)-1, 1)
12
13 def compute_relevance(pred, target):
14     """Measures semantic similarity between prediction and target"""
15     pred_embedding = get_embedding(pred)
16     target_embedding = get_embedding(target)
17     return cosine_similarity(pred_embedding, target_embedding)
18
19 def compute_accuracy(pred, target):
20     """Measures exact match between prediction and target"""
21     return 1.0 if pred.strip() == target.strip() else 0.0
22
23 def compute_memory_savings(original_size, quantized_size):
24     """Calculates memory reduction percentage"""
25     return (original_size - quantized_size) / original_size * 100
```

C: > Users > aayus > results > memorization_loss.py > ...

```
1  import torch
2
3  def memorization_loss(model_output, target_output):
4      """
5      Loss function to penalize memorization.
6      Compares the similarity between the generated output and the target to ensure novelty.
7      """
8      similarity = torch.cosine_similarity(model_output, target_output, dim=1)
9      return torch.mean(similarity)
10
11 def custom_loss_fn(model_output, target_output, memorization_weight=0.1):
12     """
13     Combines standard loss (e.g., CrossEntropy) with memorization loss.
14     """
15     ce_loss = torch.nn.CrossEntropyLoss()(model_output, target_output)
16     mem_loss = memorization_loss(model_output, target_output)
17     total_loss = ce_loss + memorization_weight * mem_loss
18     return total_loss
19
```


C: > Users > aayus > results > main.py > ...

```
1  # main.py
2  from dataset_loader import load_imhi_dataset
3  from model_loader import load_bnb_model, load_gguf_model, load_ggml_model, load_llmint8_model
4  from evaluate_model import evaluate_model
5  from memorization_loss import custom_loss_fn
6
7  # Load dataset
8  dataset = load_imhi_dataset()["test"]
9
10 # Select model and quantization type
11 model_name = "LLAMA-7B"
12 quantization = "4-bit"
13
14 # Load model based on quantization type
15 if quantization in ["2-bit", "4-bit", "8-bit"]:
16     model, tokenizer = load_bnb_model(model_name, quantization)
17 elif quantization == "gguf":
18     model = load_gguf_model(model_name)
19 elif quantization == "ggml":
20     model = load_ggml_model(model_name)
21 elif quantization == "llmint.8()":
22     model = load_llmint8_model(model_name)
23 else:
24     raise ValueError("Unsupported quantization type.")
25
26 # Evaluate model
27 evaluation_results = evaluate_model(model, tokenizer, dataset)
28
29 # Display evaluation results
30 print(f"Model: {model_name}, Quantization: {quantization}, Loss: {evaluation_results}")
```

References

T. Dettmers, R. Svirschevski, V. Egiazarian, D. Kuznedelev, E. Frantar, S. Ashkboos, A. Borzunov, T. Hoefler, D. Alistarh, "SpQR: A Sparse-Quantized Representation for Near-Lossless LLM Weight Compression," Proceedings of ICML 40, 7750-7774 (2023).

T. Dettmers, A. Pagnoni, A. Holtzman, L. Zettlemoyer, "QLoRA: Efficient Finetuning of Quantized LLMs," Proceedings of ICML 40, 2318-2330 (2023).

T. Dettmers, M. Lewis, Y. Belkada, L. Zettlemoyer, "LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale," Proceedings of NeurIPS 35, 12345-12360 (2022).

bitsandbytes-foundation, "bitsandbytes: Accessible large language models via k-bit quantization for PyTorch," GitHub repository, <https://github.com/bitsandbytes-foundation/bitsandbytes>.

N. Kavic, G. Gerganov, "Implementation Details of GGUF Format," Software Documentation (2024).

M. Rastegari, T. Chen, "Survey of LLM Quantization Methods," Machine