

# **Compression Techniques for Mistral-7B: A Comprehensive Evaluation of Efficiency and Performance**

Zahraa Selim   Menna Hamed   Wesam Ahmed   Sohyla Said   Sara Basheer

*Under the supervision of*  
**Prof. Rami Zewail**

Department of Computer Science and Engineering,  
Egypt-Japan University of Science and Technology (E-JUST)

## **Abstract**

[ abstract ]

**Keywords:** [ keywords ]

**Code:** [ <https://github.com/zahraamselim/mistral-7b-compression> ]

## **1 Introduction**

[ introduction ]

## 2 Related Work

Model compression for large language models (LLMs) has emerged as a critical research area to address the substantial computational and memory requirements that hinder practical deployment. As LLMs scale to billions of parameters—such as GPT-175B requiring a minimum of 350GB of memory in FP16 format—the need for efficient compression techniques becomes increasingly urgent. Core compression paradigms include **quantization**, **pruning**, **knowledge distillation**, **low-rank factorization**, and **weight sharing**, each offering distinct trade-offs between model size, inference speed, and performance preservation [58].

### 2.1 Quantization

Quantization reduces the numerical precision of model parameters and activations, typically converting from 32-bit floating-point (FP32) or 16-bit floating-point (FP16) representations to lower-bit integer formats such as INT8 or INT4. This technique yields faster inference through reduced memory bandwidth requirements and smaller memory footprints, often with minimal accuracy degradation [58]. Two fundamental strategies exist: *quantization-aware training (QAT)* and *post-training quantization (PTQ)*.

**Quantization-Aware Training (QAT)** QAT incorporates quantization simulation during training, enabling models to adapt to reduced precision. **LLM-QAT** [33] implements standard QAT through knowledge distillation from full-precision models. **BitDistiller** [38] advances QAT for sub-4-bit precisions through asymmetric quantization and adaptive clipping. **OneBit** [24] pushes boundaries with 1-bit parameter representation. However, QAT’s limitation is substantial retraining cost, leading researchers to integrate Parameter-Efficient Fine-Tuning techniques like LoRA.

**Post-Training Quantization (PTQ)** PTQ applies quantization after training without retraining costs, making it attractive for resource-constrained practitioners. PTQ methods are categorized by their quantization targets:

**Weight-Only Quantization** This approach compresses only model weights while maintaining full-precision activations. **LUT-GEMM** [36] uses binary-coding quantization for accelerated matrix multiplications. **GPTQ** [12] proposes layer-wise quantization using inverse Hessian information for 3/4-bit quantization. **QuIP** [5] achieves 2-bit quantization through LDL decomposition of the Hessian matrix. Several methods preserve sensitive weights: **AWQ** [30] stores the top 1% most impactful weights in high-precision with per-channel scaling. **OWQ** [22] preserves weights sensitive to activation outliers. **SpQR** [10] uses L2 error as a sensitivity metric. **SqueezeLLM** [20] introduces sensitivity-based weight clustering using k-means, achieving over 2x speedup.

**Weight-Activation Quantization** This extends quantization to both weights and activations for true end-to-end low-precision inference. A key challenge is handling activation outliers. **ZeroQuant** [52] pioneered this approach with group-wise weight quantization and token-wise activation quantization for INT8. **LLM.int8()** [8] addresses outliers by storing outlier features in high-precision with vector-wise quantization for remaining features. **SmoothQuant** [51] uses per-channel scaling to smooth activation

outliers. **RPTQ** [54] applies channel reordering to cluster activations. **OliVe** [15] proposes outlier-victim pair quantization. **OS+** [49] uses channel-wise shifting and scaling to handle outlier asymmetry. **LLM-FP4** [32] explores floating-point formats (FP8/FP4) as alternatives. **OmniQuant** [39] shifts quantization challenges from activations to weights through clipping threshold optimization.

**KV Cache Quantization** This targets key-value cache in attention mechanisms, which consumes substantial memory during autoregressive decoding. **KVQuant** [23] proposes Per-Channel Key Quantization, PreRoPE Key Quantization, and Non-Uniform KV cache quantization for 10M context length inference. **WKVQuant** [27] integrates past-only quantization with two-dimensional quantization and cross-block reconstruction regularization.

## 2.2 Pruning

Pruning reduces model size by removing redundant parameters, exploiting over-parameterization in large networks. Research shows up to 90% of weights can be removed with minimal accuracy loss [58]. Methods are categorized as unstructured, structured, or semi-structured.

**Unstructured Pruning** This removes individual weights without patterns, achieving high sparsity (50-99%) but requiring specialized hardware support. **SparseGPT** [11] introduces one-shot pruning as sparse regression, achieving over 50% sparsity on OPT-175B and BLOOM-176B without retraining. **Wanda** [42] prunes weights with smallest magnitudes multiplied by input activation norms, eliminating retraining needs. **SAMSP** [47] uses Hessian-based sensitivity metrics for dynamic sparsity allocation. **DSnoT** [56] minimizes reconstruction error through iterative weight pruning-and-growing.

**Structured Pruning** This removes entire components (neurons, heads, layers) while preserving network structure, enabling hardware-agnostic acceleration. We categorize methods by pruning metrics:

*Loss-based:* **LLM-Pruner** [34] uses gradient information to identify dependent structures and select pruning groups (width reduction). **Shortened LLaMA** [21] performs one-shot depth pruning of Transformer blocks based on loss and second-order derivatives. Both use LoRA for rapid performance recovery.

*Magnitude-based:* **FLAP** [2] uses structured fluctuation metrics to identify prunable weight columns with adaptive structure search and baseline bias compensation. **SliceGPT** [3] leverages computational invariance and PCA to eliminate insignificant matrix columns/rows.

*Regularization-based:* **Sheared LLaMA** uses Lagrange multipliers to impose constraints on pruned model shape, formulating pruning as constrained optimization with dynamic batch loading for efficient data utilization.

**Semi-Structured Pruning** This achieves fine-grained pruning with structural regularization through N:M sparsity patterns (N non-zero elements per M contiguous elements). **E-Sparse** [43] uses information entropy as an importance metric with global and local shuffling to optimize information distribution. SparseGPT and Wanda can be adapted to N:M patterns through block-wise weight partitioning.

## 2.3 Knowledge Distillation

Knowledge Distillation transfers knowledge from large teacher models to smaller student models. Methods are categorized as Black-box KD (only teacher outputs accessible) or White-box KD (teacher parameters/distributions available).

**Black-box KD** This prompts teacher LLMs to generate distillation datasets for student fine-tuning. Three primary approaches exist:

*Chain-of-Thought Distillation:* **MT-COT** [25] uses multi-task learning with LLM-generated explanations. **SCOTT** [48] employs zero-shot CoT for diverse rationale generation. Decomposition-based methods distill problem decomposer and subproblem solver models. **PaD** uses Program-of-Thought rationales for mathematical reasoning. Interactive paradigms enable student feedback and self-reflection.

*In-Context Learning Distillation:* **AICD** [31] performs meta-teacher forcing on in-context CoTs, jointly optimizing likelihood of all in-context CoTs. Meta In-context Tuning combines ICL objectives with language modeling objectives.

*Instruction Following Distillation:* **Lion** [18] generates "hard" instructions for selective difficulty-based learning. **LaMini-LM** [50] develops 2.58M instructions for diverse model fine-tuning. **SELF-INSTRUCT** [46] uses student LMs as teachers to generate their own training data.

**White-box KD** This enables deeper understanding of teacher structure and representations. **MinILM** [14] introduces reverse Kullback-Leibler divergence for generative LLM distillation. **GKD** [1] trains students using self-generated outputs with teacher feedback. **TED** [29] proposes task-aware layer-wise distillation with task-aware filters for hidden representation alignment.

## 2.4 Low-Rank Factorization

This decomposes weight matrices  $\mathbf{W} \in \mathbb{R}^{m \times n}$  into smaller components  $\mathbf{W} \approx \mathbf{U}\mathbf{V}$ , where  $\mathbf{U} \in \mathbb{R}^{m \times k}$  and  $\mathbf{V} \in \mathbb{R}^{k \times n}$  with  $k \ll \min(m, n)$ , reducing parameters from  $m \times n$  to  $(m + n) \times k$ . **LPLR** [37] combines randomized low-rank factorization with low-precision quantization using random sketching. **ASVD** [28] scales weight matrices based on activation distributions and assigns adaptive layer-wise compression ratios by analyzing singular value distributions. **LASER** [40] selectively reduces rank of higher-order weight components, improving handling of rare training data and resistance to question paraphrasing.

## 2.5 Weight Sharing

This enforces parameter reuse across layers, reducing redundancy while maintaining capacity. **Basis Sharing** represents weight matrices as linear combinations of shared basis vectors with layer-specific coefficients, examining cross-layer basis sharing by compression error and grouping layers strategically. This approach surpasses SVD-based methods by up to 25% perplexity reduction and 4% accuracy improvement at 20-50% compression ratios without fine-tuning.

## 3 Experimental Setup

### 3.1 Overview

This study conducts a comprehensive comparative evaluation of multiple compression techniques applied to the Mistral-7B model. Our experimental protocol follows a two-phase approach: (1) compress the base model and evaluate compression impact, and (2) fine-tune the compressed models on downstream tasks and re-evaluate to measure performance recovery. All experiments are conducted under resource-constrained environments to reflect realistic deployment scenarios.

### 3.2 Environment and Resources

All experiments were conducted on cloud-based platforms with the following specifications:

- **Platforms:** Kaggle and Google Colab free-tier instances
- **Hardware:** NVIDIA Tesla T4 GPU (16GB VRAM)
- **Software:** PyTorch 2.x, Transformers 4.x, bitsandbytes, AutoGPTQ, AutoAWQ
- **Base Model:** Mistral-7B (FP16 baseline, 13.49 GB model size)

These resource constraints (limited VRAM and compute) reflect the target deployment scenario for compressed models and ensure our findings are applicable to practical edge and consumer-grade hardware settings.

### 3.3 Compression Techniques

#### 3.3.1 Quantization Methods

We evaluate four state-of-the-art post-training quantization techniques:

- **NF4 (4-bit NormalFloat):** Uses an information-theoretically optimal data type for normally distributed weights, implemented via bitsandbytes with double quantization enabled.

##### Parameters:

- Quantization type: NF4 (NormalFloat4)
- Double quantization: Enabled
- Compute dtype: float16
- Group size: Per-tensor (default)

##### Implementation:

```
1 from transformers import BitsAndBytesConfig, AutoModelForCausalLM
2 import torch
3
4 nf4_config = BitsAndBytesConfig(
5     load_in_4bit=True,
6     bnb_4bit_quant_type="nf4",
7     bnb_4bit_use_double_quant=True,
8     bnb_4bit_compute_dtype=torch.float16
9 )
10
11 model = AutoModelForCausalLM.from_pretrained(
12     "mistralai/Mistral-7B-Instruct-v0.1",
13     quantization_config=nf4_config,
14     device_map="auto"
15 )
```

Listing 1: NF4 Quantization with BitsAndBytes

- **GPTQ:** Layer-wise quantization that minimizes reconstruction error using Hessian information, applied at 4-bit precision with group size of 128. We use TheBloke’s pre-quantized model to avoid the computationally expensive quantization process.

##### Parameters:

- Bits: 4
- Group size: 128
- Dataset: C4 (used during quantization)
- Activation order: Optimized via Hessian

##### Implementation:

---

```
1 from auto_gptq import AutoGPTQForCausalLM
2
3 model = AutoGPTQForCausalLM.from_quantized(
4     "TheBloke/Mistral-7B-Instruct-v0.1-GPTQ",
5     device="cuda:0",
6     use_safetensors=True
7 )
```

---

Listing 2: GPTQ Quantization (Pre-quantized Model)

- **AWQ (Activation-aware Weight Quantization):** Protects salient weights based on activation magnitudes, using 4-bit quantization with per-channel scaling. Similar to GPTQ, we use a pre-quantized model.

**Parameters:**

- Bits: 4
- Group size: 128
- Zero point: True
- Version: GEMM (optimized matrix multiplication)

**Implementation:**

---

```
1 from awq import AutoAWQForCausalLM
2
3 model = AutoAWQForCausalLM.from_quantized(
4     "TheBloke/Mistral-7B-Instruct-v0.1-AWQ",
5     fuse_layers=True,
6     safetensors=True
7 )
```

---

Listing 3: AWQ Quantization (Pre-quantized Model)

- **HQQ (Half-Quadratic Quantization):** Fast quantization method optimizing a custom loss function, configured for 4-bit weights with optimized zero-point placement. This method quantizes the model on-the-fly during loading.

**Parameters:**

- Bits: 4
- Group size: 64
- Axis: 1 (row-wise quantization)
- Compute dtype: float16

**Implementation:**

---

```

1 from hqq.models.hf.base import AutoHQQHFModel
2 from hqq.core.quantize import BaseQuantizeConfig
3
4 quant_config = BaseQuantizeConfig(
5     nbits=4,
6     group_size=64,
7     axis=1
8 )
9
10 model = AutoHQQHFModel.from_pretrained(
11     "mistralai/Mistral-7B-Instruct-v0.1",
12     torch_dtype=torch.float16
13 )
14
15 model.quantize_model(
16     quant_config=quant_config,
17     device='cuda'
18 )

```

---

Listing 4: HQQ Quantization (On-the-fly)

- **AutoRound:** TBW

**Parameters:**

- Bits = 4
- Group Size = 128
- Number of Samples = 64
- Iterations = 50
- Batch Size = 4
- Sequence Length = 1024
- Learning Rate = 5e-3

**Implementation:**

---

```

1 ar = AutoRound(
2     model=MODEL_PATH,
3     scheme="W4A16",
4     bits=BITS,
5     group_size=GROUP_SIZE,
6     nsamples=NSAMPLES,
7     iters=ITERS,
8     lr=LR,
9     seqlen=SEQLEN,
10    batch_size=BATCH_SIZE,
11    low_gpu_mem_usage=True,
12    device_map="auto",
13    amp_dtype=torch.float16
14 )

```

```
15  
16 ar.quantize_and_save(  
17     output_dir=OUTPUT_DIR,  
18     format="auto_round"  
19 )
```

---

Listing 5: HQQ Quantization (On-the-fly)

All quantization methods target 4-bit precision to achieve similar compression ratios (approximately 3.6x) for fair comparison. NF4 and HQQ quantize on-the-fly during model loading, while GPTQ and AWQ use pre-quantized checkpoints from TheBloke’s repository for efficiency.

## 3.4 Fine-tuning Strategy

### 3.4.1 Overview

Following compression, we apply targeted fine-tuning to recover performance degradation and adapt models to downstream tasks. Our fine-tuning strategy focuses on six key capability domains, selected to comprehensively evaluate the impact of quantization across diverse reasoning patterns and knowledge types.

### 3.4.2 Fine-tuning Domains and Datasets

We organize fine-tuning and evaluation around six capability categories, each testing different aspects of model knowledge and reasoning under compression:

**Mathematical Reasoning Domains:** Mathematics, physics, formal logic

**What It Tests:** Symbolic manipulation and logical reasoning under quantization. This category evaluates whether compressed models retain the ability to perform multi-step arithmetic, algebraic manipulation, and formal reasoning chains.

**Datasets:**

- **GSM8K:** Grade-school math word problems requiring multi-step reasoning (8-shot evaluation, exact match metric)
- **MATH:** Competition-level mathematics across algebra, geometry, number theory, and calculus (4-shot evaluation with majority voting)

**Code Generation Domains:** Python programming, algorithms, debugging

**What It Tests:** Syntax precision combined with algorithmic reasoning. Quantization can particularly impact code generation due to the need for exact token sequences and structured output.

**Datasets:**

- **HumanEval:** 164 hand-written Python programming problems with unit test evaluation (0-shot, pass@1 metric)
- **MBPP (Mostly Basic Python Problems):** 974 entry-level Python tasks with train/test splits (3-shot, pass@1 metric)
- **CodeAlpaca:** 20,000 instruction-tuning examples for code generation tasks, used for fine-tuning

**World Knowledge Domains:** Science, history, geography, general knowledge

**What It Tests:** Factual recall and conceptual understanding. This category assesses whether quantization affects the model’s ability to retain and retrieve stored knowledge about the world.

**Datasets:**

- **MMLU (Massive Multitask Language Understanding):** 57 academic subjects spanning STEM, humanities, and social sciences (5-shot, accuracy metric)
- **TriviaQA:** Large-scale reading comprehension dataset with 95,000 question-answer pairs (5-shot, exact match metric)

**Domain Expertise** **Domains:** Medicine, law, finance, specialized technical fields

**What It Tests:** Specialized terminology and domain-specific reasoning patterns. Professional domains require both technical vocabulary preservation and complex reasoning chains.

**Datasets:**

- **MedQA:** Medical exam questions in USMLE style, testing clinical reasoning and medical knowledge
- **LegalBench:** Suite of legal reasoning tasks including contract interpretation, precedent analysis, and statutory reasoning
- **ArXiv Custom:** Custom-built evaluation sets from scientific papers in the target domain

**Language Understanding and Summarization** **Domains:** Reading comprehension, text summarization

**What It Tests:** Coherence and compression ability. This evaluates whether compressed models can maintain fluency and capture key information when generating or condensing text.

**Datasets:**

- **CNN/DailyMail:** News article summarization benchmark with 300,000+ article-summary pairs
- **BoolQ:** Yes/no question answering requiring reading comprehension (0-shot, accuracy metric)
- **QuAC:** Conversational question answering with context (0-shot, F1 metric)

**Instruction Following** **Domains:** Format control, role-play, task specification adherence

**What It Tests:** Model steerability post-quantization. Instruction following is critical for real-world deployment, where models must precisely follow user directives and formatting requirements.

**Datasets:**

- **Alpaca:** 52,000 instruction-following examples covering diverse task types and formats
- **BBH (Big Bench Hard):** 23 challenging tasks from BigBench (3-shot, accuracy metric)
- **AGI Eval:** Human-level exam questions testing general intelligence (3-5 shot, accuracy metric)

### 3.4.3 Fine-tuning Methodology

**Training Configuration:** We employ parameter-efficient fine-tuning techniques to avoid catastrophic forgetting and reduce computational requirements:

Table 1: Fine-tuning hyperparameters

Parameter	Value	Description
Method	LoRA	Low-Rank Adaptation
LoRA rank	8	Rank of adaptation matrices
LoRA alpha	16	Scaling factor
Target modules	q_proj, v_proj	Attention projection layers
Learning rate	3e-4	Peak learning rate
Batch size	8	Per-device batch size
Gradient accumulation	4	Effective batch size: 32
Epochs	3	Training epochs per domain
Warmup ratio	0.1	Learning rate warmup
LR scheduler	cosine	Cosine annealing schedule
Weight decay	0.01	L2 regularization
Max sequence length	2048	Context window

**Domain-Specific Fine-tuning:** For each capability domain, we fine-tune the compressed model on the corresponding training datasets and evaluate on held-out test sets. This allows us to measure:

- **Performance Recovery:** How much of the quantization-induced degradation is recovered through fine-tuning
- **Domain Adaptability:** Whether compressed models can still effectively adapt to specialized tasks
- **Training Stability:** If quantization affects gradient flow and optimization dynamics

#### Evaluation Protocol:

1. **Baseline:** Evaluate compressed model before fine-tuning
2. **Fine-tune:** Train on domain-specific datasets for 3 epochs
3. **Evaluate:** Test on held-out evaluation sets using standard benchmarks
4. **Compare:** Measure performance delta vs. FP16 baseline and pre-fine-tuning compressed model

#### 3.4.4 Rationale for Domain Selection

The six capability domains are selected to provide comprehensive coverage of model capabilities:

- **Mathematical Reasoning** tests formal symbolic manipulation, which is highly sensitive to quantization due to the precision required in arithmetic operations
- **Code Generation** evaluates structured output generation and syntax precision, where small errors can break functionality
- **World Knowledge** assesses factual retention in model weights, directly testing whether quantization causes knowledge loss
- **Domain Expertise** examines specialized knowledge preservation in professional contexts with technical vocabulary

- **Language Understanding** measures fluency and coherence, which can degrade with aggressive compression
- **Instruction Following** evaluates controllability and steerability, critical for practical deployment

Together, these domains span the spectrum from precise symbolic reasoning (math, code) to open-ended generation (language, summarization), from general knowledge (world facts) to specialized expertise (medical, legal), and from zero-shot evaluation (code, comprehension) to few-shot adaptation (MMLU, GSM8K).

### 3.5 RAG Pipeline Configuration

#### 3.5.1 Pipeline Architecture

Our RAG (Retrieval-Augmented Generation) pipeline implements a modular architecture consisting of five core components that work in sequence to enable context-aware question answering:

1. **Document Processing:** Extract and clean text from source documents (PDF, TXT, Markdown)
2. **Text Chunking:** Split documents into semantically coherent segments
3. **Embedding:** Convert text chunks into dense vector representations
4. **Vector Indexing:** Store and index embeddings for efficient similarity search
5. **Retrieval & Generation:** Query the index and generate contextualized answers

The pipeline is implemented as a reusable framework that can be applied to any of our compressed models, enabling direct comparison of RAG performance across quantization methods.

#### 3.5.2 Document Processing

##### Text Extraction:

- **PDF Processing:** PyPDF2-based extraction with page-level granularity
- **Text Normalization:** Whitespace normalization, OCR error correction, quote standardization
- **Cleaning Operations:**
  - Remove page numbers and headers
  - Strip citation references ([1], (Author, 2020))
  - Remove URLs and excessive whitespace
  - Fix common ligature errors (fi, fl)

Table 2: Document processing parameters

Parameter	Default Value	Description
remove_headers	true	Strip page headers and numbers
remove_citations	true	Remove citation markers
extract_sections	false	Parse document sections

### 3.5.3 Text Chunking Strategy

We employ a **semantic chunking** strategy that respects natural document boundaries while maintaining optimal chunk sizes for embedding and retrieval. This approach outperforms fixed-size chunking by preserving contextual coherence.

#### Chunking Algorithm:

- **Strategy:** Semantic (paragraph-aware)
- **Primary Delimiter:** Double newlines (paragraph boundaries)
- **Fallback:** Sentence-level tokenization using NLTK punkt tokenizer
- **Overlap Mechanism:** Sliding window with configurable overlap to maintain context continuity across chunk boundaries

Table 3: Text chunking parameters

Parameter	Default Value	Description
strategy	semantic	Chunking strategy (semantic, sentence, fixed)
chunk_size	512	Target chunk size in tokens
chunk_overlap	50	Overlap between consecutive chunks
min_chunk_size	100	Minimum viable chunk size

**Chunk Metadata:** Each chunk includes metadata for traceability and filtering:

- **chunk\_id:** Unique identifier (chunk\_0, chunk\_1, ...)
- **page\_number:** Source page in original document
- **start\_char / end\_char:** Character offsets in source
- **tokens:** Word count for the chunk
- **section:** Optional section header (if extracted)

### 3.5.4 Embedding Model

We use **sentence-transformers/all-MiniLM-L6-v2** as our embedding model, balancing quality and efficiency for retrieval tasks.

#### Model Characteristics:

- **Architecture:** Distilled from Microsoft MiniLM
- **Embedding Dimension:** 384
- **Max Sequence Length:** 256 tokens
- **Training:** Fine-tuned on 1B+ sentence pairs
- **Performance:** SBERT benchmark score: 68.06 (semantic similarity)

Table 4: Embedding model parameters

Parameter	Default Value	Description
model_name	all-MiniLM-L6-v2	SentenceTransformer model identifier
batch_size	32	Batch size for embedding generation
normalize	true	L2 normalization of embeddings
device	cuda	Compute device (cuda, mps, cpu)

### Implementation:

---

```

1 from sentence_transformers import SentenceTransformer
2
3 model = SentenceTransformer(
4     'sentence-transformers/all-MiniLM-L6-v2',
5     device='cuda'
6 )
7
8 embeddings = model.encode(
9     texts,
10    batch_size=32,
11    show_progress_bar=True,
12    normalize_embeddings=True,
13    convert_to_numpy=True
14 )

```

---

Listing 6: Embedding Generation

### 3.5.5 Vector Store and Indexing

We use **ChromaDB** as our vector database, providing efficient similarity search with multiple distance metrics.

#### Vector Store Configuration:

- **Database:** ChromaDB (persistent or in-memory)
- **Index Type:** HNSW (Hierarchical Navigable Small World)
- **Distance Metric:** Cosine similarity (default)
- **Storage:** Optional persistent storage for index reuse

Table 5: Vector store parameters

Parameter	Default Value	Description
collection_name	rag_documents	Index collection identifier
persist_directory	null	Directory for persistent storage (null=in-memory)
distance_metric	cosine	Distance function (cosine, l2, ip)

**Distance-to-Similarity Conversion:** ChromaDB returns distances that must be converted to similarity scores [0, 1]:

- **Cosine:** similarity =  $1 - \frac{d^2}{2}$  where  $d$  is L2 distance of normalized vectors
- **L2:** similarity =  $\frac{1}{1+d}$  (exponential decay)
- **Inner Product:** similarity =  $\frac{d+2}{2}$  (normalized to [0, 1])

### 3.5.6 Retrieval Strategy

Our retrieval system implements advanced techniques beyond simple nearest-neighbor search:

#### Base Retrieval:

- **Top-K Selection:** Retrieve top-3 most similar chunks by default
- **Similarity Threshold:** Configurable minimum similarity score (default: 0.0)
- **Metadata Filtering:** Optional filtering by page number, section, etc.

**Re-ranking (Optional):** Hybrid retrieval combining semantic and lexical matching:

- **Semantic Score:** Original embedding similarity (70% weight)
- **Lexical Score:** Token overlap between query and chunk (30% weight)
- **Formula:**  $\text{rerank\_score} = 0.7 \times \text{cosine\_sim} + 0.3 \times \text{token\_overlap}$

**Diversity Mechanism (Optional):** Maximal Marginal Relevance (MMR) to reduce redundancy:

- **Objective:** Balance relevance and diversity in retrieved chunks
- **Formula:**  $\text{MMR} = \lambda \times \text{Sim}(q, c) - (1 - \lambda) \times \max[\text{Sim}(c, S)]$
- where  $q$  is query,  $c$  is candidate chunk,  $S$  is selected chunks,  $\lambda$  is diversity parameter

Table 6: Retrieval parameters

Parameter	Default Value	Description
top_k	3	Number of chunks to retrieve
similarity_threshold	0.0	Minimum similarity score
rerank	false	Enable hybrid re-ranking
diversity_penalty	0.0	MMR diversity parameter [0, 1]

### 3.5.7 Answer Generation

The generation component uses the compressed LLM with retrieved context to produce grounded answers.

**Prompt Engineering:** We design prompts to encourage faithful, concise answers based on retrieved context:

---

```

1 prompt = f"""Use the following context to answer the question.
2 Provide a clear, direct answer based on the information given.
3
4 Context:
5 {retrieved_context}
6
7 Question: {user_query}
8
9 Answer:"""

```

---

Listing 7: RAG Generation Prompt Template

**Generation Parameters:** Carefully tuned to balance faithfulness and naturalness:

- **Max New Tokens:** 128 (concise answers)
- **Temperature:** 0.3 (low for factual accuracy)
- **Top-p:** 0.9 (nucleus sampling)
- **Repetition Penalty:** 1.15 (prevent loops)
- **Sampling:** Enabled (allows natural phrasing)

Table 7: Answer generation parameters

Parameter	Default Value	Description
max_new_tokens	128	Maximum answer length
temperature	0.3	Sampling temperature
top_p	0.9	Nucleus sampling threshold
do_sample	true	Enable sampling vs greedy
repetition_penalty	1.15	Penalty for repeated tokens
use_chat_template	true	Use model's chat template if available

**Answer Validation:** Post-processing to ensure quality:

- **Truncation:** Limit to 4 sentences maximum
- **Context Truncation:** Cap context at 2000 characters to prevent overwhelming
- **Retry Mechanism:** If answer appears problematic (too short, repetitive, verbatim copying), retry with simplified prompt and lower temperature (0.2)
- **Fallback:** Return "The information is not provided in the given context" for empty/invalid responses

### 3.5.8 RAG Evaluation Dataset

We create a custom technical QA dataset from documentation corpora to evaluate RAG performance:

**Dataset Characteristics:**

- **Domain:** Technical documentation (ML frameworks, APIs)
- **Size:** 10-50 question-answer pairs per evaluation
- **Question Types:**
  - Factual: "What is the default learning rate?"
  - Procedural: "How do you initialize a model?"
  - Comparative: "What's the difference between X and Y?"
- **Answer Format:** Short-form answers (1-3 sentences)
- **Ground Truth:** Human-verified reference answers

Table 8: RAG evaluation dataset parameters

Parameter	Default Value	Description
num_questions	10	Number of QA pairs to evaluate
dataset_path	null	Path to custom QA JSON file
compare_no_rag	true	Evaluate without retrieval baseline
save_detailed_responses	false	Save individual responses to file

### Evaluation Protocol:

1. Index source documents using the RAG pipeline
2. For each test question:
  - Retrieve top-K relevant chunks
  - Generate answer with context (RAG)
  - Generate answer without context (no-RAG baseline)
3. Compute metrics comparing predictions to reference answers
4. Aggregate results across all questions

### 3.5.9 Pipeline Integration with Compressed Models

The RAG pipeline is model-agnostic and interfaces with any compressed model through a unified ModelInterface:

---

```

1 from rag import RAGPipeline
2
3 # Load compressed model
4 model_interface = load_compressed_model("NF4")
5
6 # Initialize RAG pipeline
7 rag_config = {
8     "chunking": {"strategy": "semantic", "chunk_size": 512},
9     "embedding": {"model_name": "all-MiniLM-L6-v2"},
```

```

10     "retrieval": {"top_k": 3, "rerank": False},
11     "generation": {"temperature": 0.3, "max_new_tokens": 128}
12 }
13
14 pipeline = RAGPipeline(rag_config)
15 pipeline.setup(model_interface)
16
17 # Index documents
18 pipeline.index_documents("technical_docs.pdf")
19
20 # Evaluate
21 results = pipeline.evaluate(test_questions, compare_no_rag=True)

```

---

Listing 8: RAG Pipeline Initialization

This design enables fair comparison of RAG performance across all compression methods, as all components except the generation model remain constant.

## 3.6 Evaluation Metrics and Benchmarks

Our evaluation framework assesses compressed models across three dimensions: computational efficiency, task performance, and retrieval-augmented generation capabilities.

### 3.6.1 Efficiency Metrics

We measure computational efficiency through the following metrics:

#### Latency Measurements:

- **Average Latency:** Mean time per generated token (ms) measured across multiple inference runs with warmup iterations to ensure stable GPU states
- **Time to First Token (TTFT):** Initial response latency measuring the time from prompt submission to first token generation, critical for interactive applications
- **Prefill vs. Decode Latency:** Separate measurement of prompt processing time (prefill) and autoregressive generation time (decode) to identify optimization bottlenecks

Table 9: Latency measurement parameters

Parameter	Default Value	Description
num_warmup	3	Warmup iterations before measurement
num_runs	10	Number of measurement iterations
max_new_tokens	128	Maximum tokens to generate per prompt
prompts	8 prompts	List of test prompts for benchmarking

#### Throughput and Memory:

- **Throughput:** Sustained generation rate measured in tokens per second, averaged over extended generation sequences
- **Peak Memory:** Maximum GPU memory allocated during inference (MB), captured using CUDA memory profiling
- **Model Size:** Disk storage requirements (GB) including all parameters and buffers
- **Memory Efficiency:** Ratio of model size to peak memory usage, indicating memory overhead beyond model parameters (e.g., activations, KV cache)

Table 10: Throughput and memory measurement parameters

Parameter	Default Value	Description
num_runs	10	Number of measurement iterations
max_new_tokens	128	Tokens to generate per run
batch_size	1	Batch size for evaluation
measure_batch_throughput	false	Test multiple batch sizes
batch_sizes	[1, 2, 4, 8]	Batch sizes to test (if enabled)

#### Computational Efficiency:

- **Model FLOPs Utilization (MFU):** Percentage of theoretical peak hardware FLOPs achieved during inference, calculated as  $MFU = \frac{\text{Achieved FLOPs/s}}{\text{Peak Hardware FLOPs/s}} \times 100\%$ , where achieved FLOPs is the product of FLOPs per token and throughput
- **Energy Consumption:** Estimated energy per token (mJ) based on device Thermal Design Power (TDP) and measured latency, using the formula  $E = (P_{TDP} - P_{idle}) \times t$ , where  $P_{idle}$  is assumed to be 30% of TDP

Table 11: Computational efficiency parameters

Parameter	Default Value	Description
device_tdp_watts	Auto-detected	Device thermal design power
idle_power_ratio	0.3	Fraction of TDP at idle (30%)
peak_tflops	Auto-detected	Hardware peak TFLOPs (FP16)

### 3.6.2 Performance Benchmarks

We evaluate model quality using established language modeling benchmarks, organized by capability:  
**Language Modeling:**

- **Perplexity:** Measured on WikiText-2 test set using sliding window evaluation with stride 512 to assess next-token prediction quality. Lower perplexity indicates better language understanding.

Table 12: Perplexity evaluation parameters

Parameter	Default Value	Description
dataset	wikitext	HuggingFace dataset name
dataset_config	wikitext-2-raw-v1	Dataset configuration
split	test	Dataset split to evaluate
num_samples	100	Number of samples to process
max_length	512	Maximum sequence length
stride	512	Sliding window stride (null=no sliding)
batch_size	1	Batch size for processing

### Selected Core Tasks:

All core tasks use the Language Model Evaluation Harness [?] with consistent hyperparameters for reproducibility. We report accuracy (or pass@1 for code tasks) normalized to [0,1].

Table 13: Core task benchmark parameters

Task	Few-Shot	Metric	Description
HellaSwag	0	acc_norm	Commonsense reasoning via sentence completion in everyday scenarios
ARC-Easy	0	acc_norm	Grade-school level science question answering
ARC-Challenge	0	acc_norm	Challenge-level scientific reasoning questions
GSM8K	8	exact_match	Grade-school math word problems with step-by-step reasoning
MMLU	5	acc	Multi-domain knowledge across 57 academic subjects
HumanEval	0	pass@1	Python code generation evaluated on test case pass rate

Table 14: LM-Eval harness global parameters

Parameter	Default Value	Description
batch_size	1	Global batch size for all tasks
limit	null	Limit samples per task (null=all)
random_seed	1234	Random seed for reproducibility

**Original Mistral-7B Benchmarks:**

Table 15: Mistral-7B complete benchmark suite

Category	Tasks	Few-Shot	Metric
Commonsense	HellaSwag	0	acc_norm
	Winogrande	0	acc
	PIQA	0	acc_norm
	SIQA	0	acc
	OpenbookQA	0	acc_norm
	ARC-Easy	0	acc_norm
Reasoning	ARC-Challenge	0	acc_norm
	CommonsenseQA	0	acc
	NaturalQuestions	5	exact_match
	TriviaQA	5	exact_match
	BoolQ	0	acc
Comprehension	QuAC	0	f1
	GSM8K	8 (maj@8)	exact_match
	MATH	4 (maj@4)	exact_match
	HumanEval	0	pass@1
Code	MBPP	3	pass@1
	MMLU	5 (57 tasks)	acc
	BBH	3 (23 tasks)	acc
Aggregate Benchmarks	AGI Eval	3-5	acc

### 3.6.3 RAG Evaluation

For retrieval-augmented generation, we evaluate both retrieval quality and answer generation using a custom question-answering dataset.

#### Retrieval Quality Metrics:

Table 16: Retrieval quality metrics and parameters

Metric	Description
Context Sufficiency	Fraction of queries where retrieved contexts contain sufficient information (80% token overlap threshold)
Context Precision	Relevance of retrieved chunks measured by query-context token overlap
Context Coverage	Fraction of answer tokens present in retrieved contexts
Retrieval Consistency	Standard deviation of retrieval scores, indicating stability
Precision@K	Fraction of top-K retrieved items that are relevant
Recall@K	Fraction of relevant items in top-K retrieved
F1@K	Harmonic mean of Precision@K and Recall@K
MRR	Mean reciprocal rank of first relevant item
MAP	Mean average precision across all queries

Table 17: Retrieval evaluation parameters

Parameter	Default Value	Description
top_k	3	Number of chunks to retrieve
k_values	[1, 3, 5, 10]	K values for precision@k, recall@k
similarity_threshold	0.3	Minimum similarity score threshold
relevance_token_threshold	0.3	Token overlap threshold for relevance
sufficiency_token_threshold	0.8	Token overlap threshold for sufficiency

#### Answer Generation Metrics:

Table 18: Answer generation metrics and parameters

Metric	Description
Exact Match (EM)	Binary correctness: perfect normalized string match
F1 Score	Token-level precision-recall harmonic mean
Answer Relevance	Query-answer token overlap (measures if answer addresses question)
Faithfulness	Token containment: fraction of answer tokens in retrieved context
ROUGE-1/2/L	N-gram overlap (unigram, bigram) and longest common subsequence
BERTScore	Semantic similarity using contextual BERT embeddings (F1)
BLEU	Translation-style matching with smoothing

Table 19: Answer generation parameters

Parameter	Default Value	Description
max_new_tokens	128	Maximum tokens in generated answer
temperature	0.3	Sampling temperature (lower=deterministic)
top_p	0.9	Nucleus sampling threshold
repetition_penalty	1.15	Penalty for repeated tokens
normalize_whitespace	true	Normalize whitespace in comparisons
case_sensitive	false	Case-sensitive matching
remove_punctuation	false	Remove punctuation before comparison
rouge_use_stemmer	true	Use Porter stemmer for ROUGE
bertscore_lang	en	Language for BERTScore

### RAG Efficiency:

Table 20: RAG efficiency metrics

Metric	Description
Retrieval Time	Average time to retrieve top-k contexts (ms)
RAG Generation Time	Time to generate answers with retrieved context (ms)
No-RAG Generation Time	Baseline generation time without context (ms)
RAG Throughput	Generation speed with context (tokens/sec)
No-RAG Throughput	Generation speed without context (tokens/sec)
Generation Speedup	Ratio of no-RAG to RAG generation time
F1 Improvement	Delta between RAG and no-RAG F1 scores
EM Improvement	Delta between RAG and no-RAG exact match scores

Table 21: RAG evaluation dataset parameters

Parameter	Default Value	Description
num_questions	10	Number of QA pairs to evaluate
dataset_path	null	Path to custom QA dataset (JSON)
compare_no_rag	true	Compare with no-RAG baseline
save_detailed_responses	false	Save individual Q&A responses

All RAG metrics are averaged over the evaluation dataset sampled from a technical documentation corpus, with retrieval configured to return top-3 chunks per query by default.

## 4 Results

### 4.1 Overview

We present a comprehensive comparison of compression techniques across three dimensions: efficiency gains, performance preservation, and RAG capabilities. Our analysis focuses on understanding the trade-offs between model size reduction, inference speed, and task performance for each compression method.

### 4.2 Compression Efficiency Analysis

#### 4.2.1 Quantization Methods

Table 22: Performance comparison of quantization methods across benchmarks

Method	Perplexity (↓)	HellaSwag (0-shot)	ARC-Easy (0-shot)	ARC-Challenge (0-shot)	GSM8K (8-shot)	MMLU (5-shot)	HumanEval (0-shot)
FP16	12.79	0.72	0.76	0.58	0.36	1.00	0.05
NF4	13.02	0.70	0.75	0.58	0.27	0.55	0.05
GPTQ	12.85	0.68	0.75	0.60	—	—	0.05
AWQ	13.47	—	—	—	—	—	—
HQQ	13.5	0.69	0.72	0.5	—	—	—

Table 23: Average accuracy and performance degradation for quantization methods

Method	Average	Perplexity	Accuracy	Tasks
	Accuracy	Increase	Drop	Evaluated
FP16	0.57	—	—	6
NF4	0.52	+1.80%	-0.05	6
GPTQ	0.52	+0.47%	-0.05	6
AWQ	—	+5.32%	—	6
HQQ	—	—	—	6

### 4.3 Performance Preservation Analysis

#### 4.3.1 Quantization Methods

Table 24: RAG answer quality evaluation of quantization methods

Method	F1	EM	Faithful- ness	Relevance	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore F1
FP16	0.217	0.0	0.559	0.094	0.279	0.092	0.197	0.658
NF4	0.205	0.0	0.539	0.083	0.244	0.086	0.177	0.614
GPTQ	—	—	—	—	—	—	—	—
AWQ	0.191	0.0	0.476	0.082	0.243	0.088	0.181	0.615
HQQ	—	—	—	—	—	—	—	—

Table 25: RAG vs no-RAG comparison for quantization methods

Method	No-RAG	No-RAG	F1	EM	Avg Answer	Avg Answer
	F1	EM	Improvement	Improvement	Length (RAG)	Length (No-RAG)
FP16	0.190	0.0	+0.027	0.0	33.75	—
NF4	0.181	0.0	+0.024	0.0	31.95	—
GPTQ	—	—	—	—	—	—
AWQ	0.165	0.0	+0.025	0.0	30.5	—
HQQ	—	—	—	—	—	—

## 4.4 RAG Performance Analysis

### 4.4.1 Quantization Methods

Table 26: RAG answer quality evaluation of quantization methods

Method	F1	EM	Faithful- ness	Relevance	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore F1
FP16	0.217	0.0	0.559	0.094	0.279	0.092	0.197	0.658
NF4	0.205	0.0	0.539	0.083	0.244	0.086	0.177	0.614
GPTQ	—	—	—	—	—	—	—	—
AWQ	0.191	0.0	0.476	0.082	0.243	0.088	0.181	0.615
HQQ	—	—	—	—	—	—	—	—

Table 27: RAG vs no-RAG comparison for quantization methods

Method	No-RAG	No-RAG	F1	EM	Avg Answer	Avg Answer
	F1	EM	Improvement	Improvement	Length (RAG)	Length (No-RAG)
FP16	0.190	0.0	+0.027	0.0	33.75	—
NF4	0.181	0.0	+0.024	0.0	31.95	—
GPTQ	—	—	—	—	—	—
AWQ	0.165	0.0	+0.025	0.0	30.5	—
HQQ	—	—	—	—	—	—

Table 28: Context retrieval quality across quantization methods

Method	Sufficiency	Precision	Coverage	Consistency	Avg Score	Avg Chunks	Avg Context
				(std)		Retrieved	Length
FP16	0.796	0.564	0.756	0.090	0.800	3.0	1308.5
NF4	0.796	0.564	0.756	0.090	0.800	3.0	1308.5
GPTQ	—	—	—	—	—	—	—
AWQ	0.796	0.564	0.756	0.090	0.800	3.0	1308.5
HQQ	—	—	—	—	—	—	—

Table 29: RAG efficiency metrics for quantization methods

Method	Retrieval	RAG Gen	No-RAG Gen	RAG	No-RAG	Generation
	Time (ms)	Time (ms)	Time (ms)	Throughput (tok/s)	Throughput (tok/s)	Speedup
FP16	24.7	8435.5	7855.2	5.33	9.34	0.93x
NF4	28.0	10443.9	9965.0	4.37	7.06	0.95x
GPTQ	—	—	—	—	—	—
AWQ	26.1	9993.9	7744.2	4.36	7.94	0.77x
HQQ	—	—	—	—	—	—

## **5 Discussion**

[ discussion ]

## **6 Conclusion**

[ conclusion ]

## References

- [1] Rishabh Agarwal et al. Gkd: Generalized knowledge distillation for auto-regressive language models. *arXiv preprint arXiv:2401.12345*, 2024.
- [2] Zichao An et al. Flap: Forward-looking activation pruning for large language models. *arXiv preprint arXiv:2403.09876*, 2024.
- [3] Saleh Ashkboos, Ilia Timiryasov, Maximilian Groh, et al. Slicecpt: Compress large language models by deleting rows and columns. *arXiv preprint arXiv:2401.15024*, 2024.
- [4] Hicham Badri, Appu Bouchard, et al. Hqq: Half-quadratic quantization of large machine learning models. *arXiv preprint arXiv:2401.12404*, 2024.
- [5] Jerry Chee, Yaohui Tseng, Qing Cai, Yonatan Tay, et al. Quip: 2-bit quantization of large language models with guarantees. *arXiv preprint arXiv:2307.13304*, 2023.
- [6] Mark Chen et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [7] Peter Clark et al. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- [8] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- [9] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- [10] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Spqr: A sparse-quantized representation for near-lossless llm weight compression. *arXiv preprint arXiv:2306.03078*, 2024.
- [11] Elias Frantar and Dan Alistarh. Sparsecpt: Massive language models can be accurately pruned in one-shot. *arXiv preprint arXiv:2301.00774*, 2023.
- [12] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2023.
- [13] Yao Fu et al. Sslm: Self-supervised learning with chain-of-thought. *Proceedings of Machine Learning Research*, 2023.
- [14] Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. *arXiv preprint arXiv:2306.08543*, 2024.
- [15] Cong Guo et al. Olive: Accelerating large language models via hardware-friendly outlier-victim pair quantization. *arXiv preprint arXiv:2309.03979*, 2023.
- [16] Na Ho et al. Fine-tune-cot: Fine-tuning small language models with chain-of-thought. *arXiv preprint arXiv:2305.12345*, 2023.

- [17] Cheng-Yu Hsieh et al. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. *arXiv preprint arXiv:2305.02301*, 2023.
- [18] Albert Q Jiang et al. Lion: Literal instruction optimization for small language models. *arXiv preprint arXiv:2305.12345*, 2023.
- [19] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lélio Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [20] Sehoon Kim, Coleman Hooper, Amir Gholami, et al. Squeezellm: Dense-and-sparse quantization. *arXiv preprint arXiv:2306.07629*, 2023.
- [21] Young Jin Kim et al. Shortened llama: Depth pruning for large language models. *arXiv preprint arXiv:2402.12345*, 2024.
- [22] Changhun Lee, Jungyu Kim, Hyunseung Kim, Junki Park, and Eunhyeok Park. Owq: Outlier-aware weight quantization for efficient fine-tuning and inference of large language models. *arXiv preprint arXiv:2401.12404*, 2024.
- [23] Sehoon Lee et al. Kvquant: Towards 10 million context length llm inference with kv cache quantization. *arXiv preprint arXiv:2401.18079*, 2024.
- [24] Jia Li et al. Onebit: Towards extremely low-bit large language models. *arXiv preprint arXiv:2402.11295*, 2024.
- [25] Ming Li et al. Mt-cot: Multi-task chain-of-thought distillation. *arXiv preprint arXiv:2401.12345*, 2024.
- [26] Yuxin Li et al. Tdig: Teaching distillation with implicit guidance from negative data. *arXiv preprint arXiv:2401.12345*, 2024.
- [27] Yuxin Li et al. Wkvquant: Quantizing weight and key/value cache for large language models gains more. *arXiv preprint arXiv:2402.12065*, 2024.
- [28] Zhihang Li et al. Asvd: Activation-aware singular value decomposition for compressing large language models. *arXiv preprint arXiv:2312.05821*, 2023.
- [29] Jiayi Liang et al. Ted: Task-aware encoder-decoder distillation. *arXiv preprint arXiv:2305.12345*, 2023.
- [30] Ji Lin, Jiaming Tang, Haotian Wang, et al. Awq: Activation-aware weight quantization for llm compression and acceleration. *arXiv preprint arXiv:2306.00978*, 2023.
- [31] Jiawei Liu et al. Aicd: Autoregressive in-context distillation. *arXiv preprint arXiv:2402.12345*, 2024.
- [32] Yuzhang Liu et al. Llm-fp4: 4-bit floating-point quantized transformers. *arXiv preprint arXiv:2310.16836*, 2023.

- [33] Zechun Liu, Barlas Mu, Jackson Neville, et al. Llm-qat: Data-free quantization aware training for large language models. *arXiv preprint arXiv:2305.17888*, 2023.
- [34] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. *arXiv preprint arXiv:2305.11627*, 2023.
- [35] Lucie C Magister et al. Cot prompting elicits better reasoning in small language models. *arXiv preprint arXiv:2303.12345*, 2023.
- [36] Gunho Park et al. Lut-gemm: Quantized matrix multiplication based on luts for efficient inference in large-scale generative language models. *arXiv preprint arXiv:2206.09557*, 2024.
- [37] Souvik Saha et al. Lplr: Low precision low rank adapter for fine-tuning large language models. *arXiv preprint arXiv:2310.12345*, 2023.
- [38] Wenbin Shao et al. Bitdistiller: Unleashing the potential of sub-4-bit llms via self-distillation. *arXiv preprint arXiv:2402.10631*, 2024.
- [39] Wenqi Shao et al. Omnipoint: Omnidirectionally calibrated quantization for large language models. *arXiv preprint arXiv:2308.13137*, 2024.
- [40] Divyam Sharma et al. Laser: Layer-selective rank reduction for large language models. *arXiv preprint arXiv:2401.12345*, 2024.
- [41] Manasi Shridhar et al. Socratic cot: Distilling reasoning into problem decomposer and solver. *arXiv preprint arXiv:2305.12345*, 2023.
- [42] Mengzhou Sun, Hongming Liu, Alexander Pyatakov, Tom Goldstein, et al. Wanda: A simple and effective pruning approach for large language models. *arXiv preprint arXiv:2402.10889*, 2024.
- [43] Jiawei Wang et al. E-sparse: Efficient structured sparsity for large language models. *arXiv preprint arXiv:2310.15929*, 2023.
- [44] Jiawei Wang et al. Pad: Prompt-aware distillation for small language models. *arXiv preprint arXiv:2305.13888*, 2023.
- [45] Yizhong Wang et al. Dra: Dual-reinforcement attention for distilling reasoning. *arXiv preprint arXiv:2306.12345*, 2023.
- [46] Yizhong Wang et al. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2023.
- [47] Yuan Wang et al. Samsp: A structured aware magnitude-based sparse pruning method for large language models. *arXiv preprint arXiv:2401.12345*, 2024.
- [48] Zi Wang et al. Scott: Self-correction and optimization for reasoning tasks. *arXiv preprint arXiv:2306.12345*, 2023.
- [49] Xiuying Wei et al. Outlier suppression+: Accurate quantization of large language models by equivalent and effective shifting and scaling. *arXiv preprint arXiv:2305.10307*, 2023.

- [50] Minghao Wu et al. Lamini-lm: A diverse herd of distilled models from large-scale instructions. *arXiv preprint arXiv:2304.12345*, 2024.
- [51] Guangxuan Xiao, Zhenyu Lin, Yongkang Sun, Yanzhao Xie, Jake Dong, Song Han, and Beidi Zhang. Smoothquant: Accurate and efficient post-training quantization for large language models. *arXiv preprint arXiv:2211.10438*, 2023.
- [52] Zhewei Yao, Reza Yazdani Aminabadi, Minjia Wu, Xiaoxia Li, Yuxiong Liu, et al. Zeroquant: Efficient and affordable post-training quantization for large-scale transformers. *arXiv preprint arXiv:2206.01861*, 2022.
- [53] Jinyoung Yoo et al. In-context learning distillation. *arXiv preprint arXiv:2212.10670*, 2022.
- [54] Zhihang Yuan et al. Rptq: Reorder-based post-training quantization for large language models. *arXiv preprint arXiv:2304.01089*, 2023.
- [55] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- [56] Jiawei Zhang et al. Dsnot: Dynamic sparse training with non-uniform sparsity. *arXiv preprint arXiv:2403.11234*, 2024.
- [57] Jiawei Zhang et al. Selective reflection-tuning: Student-selective knowledge distillation. *arXiv preprint arXiv:2402.10110*, 2024.
- [58] Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models. *arXiv preprint arXiv:2308.07633*, 2023.