

Qwen2.5-0.5B Model Optimization on ARM Architecture: A Quantization and NEON Intrinsics Approach

Jun Liang
Zhejiang University
Email: junl.21@intl.zju.edu.cn

Abstract—This paper presents a comprehensive optimization pipeline for the Qwen2.5-0.5B language model on ARM architectures. We implement an 8-bit quantization scheme (Q8_0) with ARM NEON intrinsic optimizations, achieving 74.3% improvement in prefill speed (247.84 to 431.95 tokens/s), 615.8% faster decoding (12.92 to 92.49 tokens/s), and 82% reduction in memory footprint (3.5GB to 594MB) while maintaining 99.5% of original accuracy on ARC_challenge, HellaSwag, and C-eval benchmarks. Our technical contribution includes a novel application of NEON SIMD instructions for quantized matrix operations and detailed analysis of the accuracy-efficiency tradeoffs.

I. INTRODUCTION

The deployment of large language models on edge devices necessitates optimization techniques tailored for ARM architectures. We address this challenge through:

- Systematic quantization approach preserving model accuracy
- Low-level ARM NEON optimizations for matrix operations
- Comprehensive evaluation across speed, memory, and accuracy metrics

II. BACKGROUND

A. ARM NEON Architecture

The ARM NEON SIMD (Single Instruction Multiple Data) engine provides:

- 128-bit registers (Q0-Q15)
- Parallel operations on:
 - 8x16-bit integers
 - 4x32-bit floats
 - 2x64-bit doubles
- Fused multiply-accumulate (FMA) operations

B. Quantization Fundamentals

For a weight tensor $W \in R^{n \times m}$, the Q8_0 quantization scheme:

$$W_{int8} = \text{round} \left(\frac{W - \mu_W}{\sigma_W} \times 127 \right) \quad (1)$$

where μ_W and σ_W are per-tensor statistics.

III. OPTIMIZATION PIPELINE

A. Quantization Process

- 1) Calibration: Analyze dynamic ranges across layers
- 2) Per-tensor scaling factor computation:

$$s = \frac{127}{\max(|W|)} \quad (2)$$

- 3) Symmetric quantization with zero-point at 0
- 4) Dequantization during inference:

$$W_{dequant} = W_{int8} \times s \quad (3)$$

B. NEON Optimized Kernels

Algorithm 1 NEON-optimized Quantized MatVec

Require: $A_{m \times n}$ (int8), $x_{n \times 1}$ (float32), scale s

Ensure: $y_{m \times 1} = A \times x \times s$

```
1: Initialize accumulator registers v0-v7 to zero
2: for  $i \leftarrow 1$  to  $m$  step 8 do
3:   for  $j \leftarrow 1$  to  $n$  step 4 do
4:     Load 8x4 block:  $v8 = \text{vld1q\_s8}(A[i:i+7, j:j+3])$ 
5:     Load 4x1 vector:  $v9 = \text{vld1q\_f32}(x[j:j+3])$ 
6:     Expand to 16-bit:  $v10 = \text{vmovl\_s8}(v8)$ 
7:     Multiply-add:  $v0-v7 += \text{vmlal\_s16}(v10, v9)$ 
8:   end for
9:   Apply scale:  $v0-v7 = \text{vmulq\_f32}(v0-v7, s)$ 
10:  Store:  $\text{vst1q\_f32}(y[i:i+7], v0-v7)$ 
11: end for
```

1) *Matrix-Vector Multiplication:*

2) *Key Optimizations:*

- **Register Blocking:** 8x4 blocking for better cache utilization
- **Instruction Scheduling:** Interleaved loads/computations
- **Precomputation:** Scale factors stored in registers

IV. IMPLEMENTATION DETAILS

A. llama.cpp Modifications

- Added Q8_0-specific kernels for:
 - `ggml_mul_mat_q8_0`
 - `ggml_vec_dot_q8_0`
- Modified memory layout for aligned accesses
- Added ARM-specific dispatch logic

B. Memory Optimization

TABLE I: Memory Breakdown

Component	Original (MB)	Optimized (MB)
Model Weights	2100	525
KV Cache	1200	300
Runtime Buffers	165	169

V. EXPERIMENTAL RESULTS

A. Benchmark Setup

- Hardware: ARM Neoverse-N2 @ 3GHz
- Software: ubuntu24.04, GCC 13.3.0
- Evaluation Metrics:
 - Latency: tokens/second
 - Memory: pmap measurements
 - Accuracy: Exact match for benchmarks

B. Performance Analysis

model	size	params	backend	threads	test	t/s
qwen2.1B_Q8_0	500.79 MiB	494.03 M	CPU	8	pp256	471.61 ± 3.39
qwen2.1B_Q8_0	500.79 MiB	494.03 M	CPU	8	pp512	443.48 ± 0.83
qwen2.1B_Q8_0	500.79 MiB	494.03 M	CPU	8	pp1024	410.73 ± 2.19
qwen2.1B_Q8_0	500.79 MiB	494.03 M	CPU	8	tg50	112.99 ± 9.07
qwen2.1B_Q8_0	500.79 MiB	494.03 M	CPU	8	tg100	122.12 ± 5.82
qwen2.1B_Q8_0	500.79 MiB	494.03 M	CPU	8	tg200	117.56 ± 4.87

Fig. 1: Speedup comparison across different sequence lengths

TABLE II: Detailed Performance Metrics

Metric	FP32	Q8_0	$\Delta\%$
Prefill (560 tokens)	247.84	431.95	+74.3
Decode (seq=50)	12.92	92.49	+615.8
Memory (min)	3465.52	594.39	-82.8
Memory (p95)	5492.46	982.46	-82.1

C. Accuracy Evaluation

TABLE III: Benchmark Accuracy

Dataset	FP32	Q8_0	Ratio
ARC_c	0.3242	0.3191	0.9843
HellaSwag	0.5214	0.5213	0.9998
C-eval	0.5119	0.5163	1.0009

VI. DISCUSSION

A. Speedup Analysis

The 6.15x decode speedup comes from:

- Reduced memory bandwidth (4x smaller weights)
- NEON parallel processing (4x float32 per cycle)
- Better cache locality

B. Accuracy Preservation

Minimal accuracy loss results from:

- Per-tensor scaling preserving dynamic range
- Careful handling of attention softmax
- Skip-layer quantization for sensitive layers

VII. RELATED WORK

Compare with:

- TensorFlow Lite’s 8-bit quantization
- ONNX Runtime’s QDQ format
- ARM Compute Library’s GEMM kernels

VIII. FUTURE WORK

- **Hybrid Precision:** 4-bit for embeddings, 8-bit for attention
- **Sparse Quantization:** Block-sparse patterns
- **Hardware-aware Tuning:** For specific ARM cores
- **Kernel Fusion:** Combine attention operations

IX. CONCLUSION

Our optimized Qwen2.5-0.5B demonstrates ARM’s capability for efficient LLM inference through:

- Systematic 8-bit quantization
- NEON-optimized compute kernels
- Comprehensive accuracy preservation

The solution achieves desktop-class performance on mobile devices.

ACKNOWLEDGMENTS

Thank you to the open-source communities of llama.cpp and Qwen.