# Test-Time Adaptation with FP4 Quantization: Experimental Report

Tomer Barak

## ABSTRACT

*This report presents an experimental evaluation combining two recent advances: Test-Time Adaptation via Entropy Minimization (TENT) and FP4 Fully Quantized Training. We investigate whether TTA can improve LLM predictions when operating under extreme FP4 quantization—a setting relevant for edge deployment where models must be both small and adaptive. We evaluate three configurations on OPT-125M across 15 sequence completion tasks with 100 iterations each. Error estimates use Wilson score confidence intervals. Key findings: TTA improves average accuracy from 28.5% to 39.7% (+11.2% absolute improvement), FP4 quantization preserves most TTA benefits, achieving 36.6% average accuracy (+8.1% over baseline), TTA shows dramatic improvements on structured pattern tasks (up to +62% on individual tasks), and some tasks show degradation with TTA when entropy minimization leads to overconfident incorrect predictions. The combination of TTA + FP4 enables efficient adaptation on resource-constrained hardware with acceptable quality trade-offs.*

## 1. Test-Time Adaptation with FP4 Quantization: Experimental Report

**Date:** December 7, 2025
**Model:** facebook/opt-125m

### 1.1 Executive Summary

This report presents an experimental evaluation combining two recent advances: **Test-Time Adaptation via Entropy Minimization (TENT)** and **FP4 Fully Quantized Training**. We investigate whether TTA can improve LLM predictions when operating under extreme FP4 quantization—a setting relevant for edge deployment where models must be both small and adaptive.

We evaluate three configurations on OPT-125M across 15 sequence completion tasks with 100 iterations each. Error estimates use Wilson score confidence intervals.

**Key Findings:**

- **TTA improves average accuracy from 28.5% to 39.7%** (+11.2% absolute improvement)

- **FP4 quantization preserves most TTA benefits**, achieving 36.6% average accuracy (+8.1% over baseline)

- TTA shows **dramatic improvements** on structured pattern tasks (up to +62% on individual tasks)

- Some tasks show **degradation with TTA** when entropy minimization leads to overconfident incorrect predictions

- The combination of TTA + FP4 enables efficient adaptation on resource-constrained hardware with acceptable quality trade-offs

**Implications:** These results suggest that LLMs deployed on edge devices with FP4 quantization can benefit from test-time adaptation, opening possibilities for on-device learning that improves model accuracy without retraining.

## 1.2   1. Introduction

### 1.2.1   Motivation

Large Language Models (LLMs) are increasingly deployed on edge devices such as smartphones, IoT devices, and embedded systems. To fit within the memory and compute constraints of these devices, models are typically quantized to very low precision—such as FP4 (4-bit floating point)—for inference. However, quantized models often suffer from accuracy degradation, particularly on out-of-distribution inputs or when the deployment context differs from training.

**Test-Time Adaptation (TTA)** offers a compelling solution: adapting the model at inference time using only the test input itself, without requiring labeled data or retraining. If TTA can be performed efficiently in low-precision arithmetic, it opens the possibility of **on-device learning** that improves model accuracy in real-time, even on resource-constrained hardware.

This work investigates whether TTA can be effectively combined with FP4 quantization for LLMs, potentially enabling adaptive, efficient inference on edge devices.

### 1.2.2   Background

**Test-Time Adaptation via Entropy Minimization (TENT)**   Our TTA approach is based on **TENT** (Test-Time Entropy Minimization), introduced by Wang et al. (ICLR 2021). TENT proposes adapting a model during test time by minimizing the entropy of its predictions:

$$L = -\sum_i p_i \log(p_i)$$

The key insight is that confident predictions (low entropy) tend to be more accurate. By optimizing model parameters to reduce prediction entropy on test inputs, the model adapts to the specific characteristics of the test data without requiring labels. TENT was originally demonstrated on image classification tasks with batch normalization adaptation, achieving state-of-the-art results on distribution shift benchmarks like ImageNet-C.

In this work, we apply the entropy minimization principle to **next-token prediction in LLMs**, adapting all model weights (not just normalization layers) for each test prompt.

**FP4 Quantization for Training (FP4 All The Way)**   Recent work by Chmiel et al. (2025) demonstrated, for the first time, **Fully Quantized Training (FQT)** of LLMs using predominantly 4-bit floating-point (FP4) precision for weights, activations, and gradients. Key findings from this work include:

- **NVFP4 format** (E2M1 with block size 16 and E4M3 scaling) provides optimal results

- **Split rounding strategy**: Stochastic rounding in backward/update passes, round-to-nearest in forward pass

- **Critical threshold**: When gradient magnitude falls below $\sqrt{3}$ times the quantization noise, FP4 training becomes less effective

- Successfully trained a 7B parameter model on 1 trillion tokens with downstream performance matching BF16 baseline

The FP4 format uses only 4 bits per value:

- 1 sign bit

- 2 exponent bits

- 1 mantissa bit

Representable values: $\{0, 0.5, 1, 1.5, 2, 3, 4, 6\}$ and their negatives.

### 1.2.3 Significance of This Work

This work bridges two important research directions:

1. **TTA for LLMs**: Extending test-time adaptation from vision models to language models

2. **FP4 training**: Leveraging extreme quantization for efficient on-device computation

The combination suggests a path toward **adaptive edge AI**: LLMs quantized to FP4 for efficient deployment can use TTA to improve their predictions at inference time, all within the constraints of edge hardware. This is particularly relevant as:

- Edge devices increasingly run LLM inference (e.g., on-device assistants)

- FP4 hardware support is emerging (NVIDIA Blackwell architecture)

- Users benefit from models that adapt to their specific use patterns

### 1.2.4 Research Questions

1. Does Entropy Minimization TTA improve next-token prediction accuracy for LLMs?

2. Does FP4 quantization preserve the benefits of TTA?

3. What is the statistical reliability of these improvements across diverse tasks?

## 1.3 2. Methodology

### 1.3.1 2.1 Test-Time Adaptation Strategy

**Entropy Minimization Loss:**

$$L = -\sum_i p_i \log(p_i)$$

where $p_i$ is the softmax probability of the next token prediction.

**Intuition:** By minimizing the entropy of the model&#39;s predictions, we encourage the model to become more confident about its next token prediction for the given context.

**Procedure:**

1. Start with the pre-trained OPT-125M model

2. For each test prompt, perform 10 gradient descent steps

3. Optimize only to minimize prediction entropy (no labeled data required)

4. Generate the next token using the adapted model

### 1.3.2   2.2 FP4 Quantization

**Format:** E2M1 (1 exponent bit, 2 mantissa bits)
**Implementation:** Simulated FP4 computation using:

- Fake quantization during forward pass

- Straight-Through Estimator (STE) for gradients

- Block-wise quantization (block size: 32)

- BF16 scaling factors

**Key Insight:** While computations are simulated on CPU, this approach models the behavior of actual FP4 hardware accelerators.

### 1.3.3   2.3 Experimental Setup

**Three Configurations Evaluated:**

1. **Base Model:** OPT-125M without adaptation (with sampling)

2. **TTA:** Entropy Minimization with full precision (FP16)

3. **TTA-FP4:** Entropy Minimization with FP4-quantized weights

**Evaluation Protocol:**

- 15 diverse sequence completion tasks

- 100 independent trials per task per configuration

- Sampling enabled (temperature=0.7, top_k=50) to measure statistical success rate

- Success criterion: Expected token appears in the generated output

**Hyperparameters:**

- Learning Rate: 1e-4

- Optimizer: AdamW

- TTA Steps: 10

- Device: CPU

### 1.3.4   2.4 Task Categories

The 15 test problems span multiple reasoning types:

**Arithmetic Sequences (4 tasks):**

- Simple arithmetic progression: &quot;2, 4, 6, 8, 10,&quot; $\{}rightarrow$ &quot; 12&quot;

- Fibonacci sequence: &quot;1, 1, 2, 3, 5, 8,&quot; $\{}rightarrow$ &quot; 13&quot;

- Countdown: &quot;10, 9, 8, 7, 6,&quot; $\{}rightarrow$ &quot; 5&quot;

- Day sequence: &quot;The first day is Monday...&quot; $\{}rightarrow$ &quot; Wednesday&quot;

**Knowledge &amp; Associations (3 tasks):**

- Capital cities pattern completion

- Analogies (color-fruit, hot-cold)

**Pattern Recognition (5 tasks):**

- Letter patterns: &quot;A A B B C C D D E&quot; $\{\}rightarrow$ &quot; E&quot;

- Alphanumeric patterns: &quot;x1 y1 z1 x2 y2 z2 x3 y3&quot; $\{\}rightarrow$ &quot; z3&quot;

- Repetition patterns with various complexity

**Simple Logic (1 task):**

- Consequence reasoning

**Repetition Tasks (2 tasks):**

- Multi-element repeating sequences

## 1.4   3. Results

### 1.4.1   3.1 Error Estimation Methodology

All accuracy measurements include 95% confidence intervals computed using the **Wilson score interval** method, which is appropriate for binomial proportions and performs well for proportions near 0 or 1.

### 1.4.2   3.2 Overall Statistics

**Average Accuracy Across All Tasks (with 95% Wilson CI, n=1500 pooled trials):**

| Configuration | Accuracy (95% CI) | Improvement over Base |
|---|---|---|
| Base Model | 28.5% (26.3%, 30.9%) | - |
| TTA | 39.7% (37.3%, 42.2%) | +11.2% |
| TTA-FP4 | 36.6% (34.2%, 39.1%) | +8.1% |

### 1.4.3   3.3 Per-Task Breakdown

The table below shows accuracy for each task with 95% Wilson confidence intervals based on 100 independent trials.

| ID | Prompt (Truncated) | Base Acc (95% CI) | TTA Acc (95% CI) | FP4 Acc (95% CI) |
|---|---|---|---|---|
| 1 | 2, 4, 6, 8, 10,... | 0% (0.0%, 3.7%) | 0% (0.0%, 3.7%) | 1% (0.2%, 5.4%) |
| 2 | 1, 1, 2, 3, 5, 8,... | 6% (2.8%, 12.5%) | 0% (0.0%, 3.7%) | 0% (0.0%, 3.7%) |
| 3 | 10, 9, 8, 7, 6,... | 11% (6.3%, 18.6%) | 60% (50.2%, 69.1%) | 16% (10.1%, 24.4%) |
| 4 | The first day is Monday... | 34% (25.5%, 43.7%) | 82% (73.3%, 88.3%) | 96% (90.2%, 98.4%) |
| 5 | The capital of France... | 63% (53.2%, 71.8%) | 97% (91.5%, 99.0%) | 86% (77.9%, 91.5%) |
| 6 | Red is to apple... | 3% (1.0%, 8.5%) | 0% (0.0%, 3.7%) | 0% (0.0%, 3.7%) |
| 7 | Hot is to cold... | 0% (0.0%, 3.7%) | 0% (0.0%, 3.7%) | 0% (0.0%, 3.7%) |
| 8 | A A B B C C D D E | 22% (15.0%, 31.1%) | 70% (60.4%, 78.1%) | 81% (72.2%, 87.5%) |
| 9 | x1 y1 z1 x2 y2 z2... | 0% (0.0%, 3.7%) | 0% (0.0%, 3.7%) | 0% (0.0%, 3.7%) |
| 10 | Statement: It is raining... | 91% (83.8%, 95.2%) | 88% (80.2%, 93.0%) | 100% (96.3%, 100.0%) |
| 11 | cat dog mouse... | 10% (5.5%, 17.4%) | 3% (1.0%, 8.5%) | 0% (0.0%, 3.7%) |
| 12 | 1 2 3 1 2 3 1 2 | 54% (44.3%, 63.4%) | 69% (59.4%, 77.2%) | 23% (15.8%, 32.2%) |
| 13 | A B C A B C A B | 64% (54.2%, 72.7%) | 95% (88.8%, 97.8%) | 92% (85.0%, 95.9%) |
| 14 | Sun Moon Star... | 55% (45.2%, 64.4%) | 32% (23.7%, 41.7%) | 54% (44.3%, 63.4%) |
| 15 | Up Down Left Right... | 15% (9.3%, 23.3%) | 0% (0.0%, 3.7%) | 0% (0.0%, 3.7%) |

### 1.4.4  3.4 Statistical Significance Testing

To assess whether TTA and TTA-FP4 significantly improve performance over the baseline, we perform paired comparisons across all 15 tasks. We use McNemar&#39;s test for paired binary outcomes aggregated across tasks.

**Aggregate Performance (pooled across all 1500 trials = 15 tasks $\{}times$ 100 iterations):**

| Comparison | Total Correct (Base) | Total Correct (Method) | Improvement | Significance |
|---|---|---|---|---|
| TTA vs Base | 428/1500 (28.5%) | 595/1500 (39.7%) | +167 (+11.1%) | p &lt; 0.001 |
| TTA-FP4 vs Base | 428/1500 (28.5%) | 549/1500 (36.6%) | +121 (+8.1%) | p &lt; 0.001 |
| TTA vs TTA-FP4 | 595/1500 (39.7%) | 549/1500 (36.6%) | -46 (-3.1%) | p &lt; 0.001 |

**Per-Task Significance (based on non-overlapping Wilson CIs):**

| Comparison | Tasks with Significant Improvement | Tasks with Significant Degradation | Tasks with |
|---|---|---|---|
| TTA vs Base | 6 (Tasks 3, 4, 5, 8, 12, 13) | 3 (Tasks 11, 14, 15) | 6 |
| TTA-FP4 vs Base | 6 (Tasks 3, 4, 5, 8, 10, 13) | 2 (Tasks 11, 12) | 7 |

**Conclusion:** Both TTA and TTA-FP4 provide statistically significant improvements over the baseline when aggregated across all tasks. However, the effect is task-dependent: TTA shows significant improvement on 6/15 tasks but significant degradation on 3/15 tasks. TTA-FP4 shows a more conservative improvement profile with fewer both gains and losses.

### 1.4.5  3.5 Qualitative Analysis

**Key Observations:**

1. **TTA Significantly Improves Performance on Structured Tasks:**

   - Day sequence (Task 4): 34% $\{}rightarrow$ 82% ($\{}Delta$ = +48%)

   - Capital cities (Task 5): 63% $\{}rightarrow$ 97% ($\{}Delta$ = +34%)

   - Letter pattern ABC (Task 13): 64% $\{}rightarrow$ 95% ($\{}Delta$ = +31%)

   - Countdown (Task 3): 11% $\{}rightarrow$ 60% ($\{}Delta$ = +49%)

2. **TTA Can Hurt Performance on Some Tasks:**

   - &quot;cat dog mouse&quot; repetition (Task 11): 10% $\{}rightarrow$ 3% ($\{}Delta$ = -7%)

   - &quot;Sun Moon Star&quot; repetition (Task 14): 55% $\{}rightarrow$ 32% ($\{}Delta$ = -23%)

   - &quot;Up Down Left Right&quot; (Task 15): 15% $\{}rightarrow$ 0% ($\{}Delta$ = -15%)

3. **FP4 Preserves Most TTA Benefits:**

   - Average degradation from TTA to TTA-FP4: 3.1%

   - FP4 still outperforms baseline on most tasks where TTA helps

   - Notable improvements over TTA: Day sequence (82% $\{}rightarrow$ 96%), Letter pattern E (70% $\{}rightarrow$ 81%)

   - FP4 achieves 100% on logic task (Task 10), outperforming both baseline and TTA

4. **Hard Tasks Remain Hard:**

- Tasks 1, 2, 6, 9 show near-zero accuracy across all conditions
- These require deeper reasoning (Fibonacci, analogies) that entropy minimization cannot address

5. **Statistical Significance:**
   - Non-overlapping Wilson CIs indicate statistically significant differences
   - TTA improvements on Tasks 3, 4, 5, 8, 12, 13 are highly significant
   - TTA degradation on Task 15 is also significant

## 1.5  4. Discussion

### 1.5.1  4.1 Effectiveness of TTA

The results demonstrate that entropy minimization TTA can significantly improve prediction accuracy on structured pattern completion tasks. The improvement is most pronounced for:

- **Sequential patterns** with clear structure (countdown, day names)
- **Repetitive patterns** with simple rules (ABC, 123)
- **Knowledge recall** with strong priors (capital cities)

However, TTA can be counterproductive when entropy minimization causes the model to &quot;collapse&quot; onto incorrect but confident predictions. This is observed in Tasks 11, 14, and 15 where the adapted model becomes overly confident in wrong answers.

The entropy minimization approach works because it encourages the model to commit to specific predictions rather than maintaining uncertainty. This is effective when the model&#39;s initial probability mass is spread across correct and incorrect tokens, and TTA helps concentrate probability on the correct answer.

### 1.5.2  4.2 Impact of FP4 Quantization

The results show that FP4 quantization preserves the benefits of TTA while providing:

- ~4x memory reduction for weights
- Potential for specialized hardware acceleration
- Minimal degradation in adaptation quality

### 1.5.3  4.3 Limitations

1. **Task Scope:** Results are limited to next-token prediction on short sequences
2. **Model Size:** Experiments used only OPT-125M (smallest scale)
3. **Hardware:** Simulated FP4 on CPU rather than actual FP4 accelerators
4. **Adaptation Depth:** Only 10 optimization steps per sample

## 1.6  5. Conclusions

### 1.6.1  Key Takeaways

1. **TTA Works on Structured Tasks:** Entropy minimization successfully improves prediction accuracy at test time for pattern completion and sequential reasoning tasks, with improvements up to +49% on individual tasks and +11.2% average improvement.

2. **TTA Has Limitations:** Not all tasks benefit from TTA. Some tasks show degradation when the model becomes overconfident in incorrect predictions. Task selection is important for practical applications.

3. **FP4 Is Viable:** Extreme quantization to FP4 reduces TTA effectiveness by only 3.1% on average compared to full-precision TTA, while still outperforming the baseline by +8.1%. Remarkably, FP4 sometimes outperforms full-precision TTA on certain tasks. This makes it highly feasible for resource-constrained deployment.

4. **Statistical Reliability:** With 100 iterations per task and Wilson score confidence intervals, we establish statistical confidence in the observed improvements and can distinguish significant effects from noise.

### 1.6.2 Future Directions

1. **Larger Models:** Evaluate on models with billions of parameters

2. **Real Hardware:** Test on actual FP4 accelerators (e.g., Habana Gaudi)

3. **Longer Sequences:** Extend to multi-token generation and longer contexts

4. **Alternative TTA Methods:** Compare with other test-time adaptation strategies

5. **Task Diversity:** Expand to more complex reasoning and generation tasks

## 1.7 6. Reproducibility

### 1.7.1 Code Structure

```
FP4-TTA/
|-- src/
|   |-- fp4_quantization.py    # FP4 quantization primitives
|   |-- stats_tta.py           # Main evaluation script
|   +-- ...
|-- requirements.txt
+-- REPORT.md
```

### 1.7.2 Dependencies

- PyTorch

- Transformers (Hugging Face)

- NumPy

### 1.7.3 Running the Evaluation

```
python src/stats_tta.py
```

## 1.8 Appendix A: FP4 Format Details

**E2M1 Format (4 bits total):**

- 1 sign bit

- 2 exponent bits

- 1 mantissa bit

**Representable Values:** {0, 0.5, 1, 1.5, 2, 3, 4, 6} and negatives

**Quantization Strategy:**

- Block-wise: weights divided into blocks of 32

- Scaling: Each block has a BF16 scale factor

- Rounding: Nearest neighbor (deterministic)

## 1.9 References

1. Wang, D., Shelhamer, E., Liu, S., Olshausen, B., &amp; Darrell, T. (2021). **Tent: Fully Test-Time Adaptation by Entropy Minimization**. ICLR 2021 (Spotlight). https://openreview.net/forum?id=uXl3bZLkr3c

2. Chmiel, B., Fishman, M., Banner, R., &amp; Soudry, D. (2025). **FP4 All the Way: Fully Quantized Training of LLMs**. arXiv:2505.19115v2. https://arxiv.org/abs/2505.19115

3. Zhang, S., Roller, S., Goyal, N., et al. (2022). **OPT: Open Pre-trained Transformer Language Models**. arXiv:2205.01068. Meta AI.

4. Rouhani, B. D., et al. (2023). **Microscaling Data Formats for Deep Learning**. arXiv:2310.10537.

5. NVIDIA (2024). **Blackwell Architecture Whitepaper**. https://resources.nvidia.com/en-us-blackwell-architecture