

# **Adaptiq vs GPT**

## Comprehensive Performance Analysis

*Generated on: 2025-08-13 17:24:04*

# Benchmark Configuration

## BENCHMARK CONFIGURATION

```
=====
LLM/VLM Model   : GPT-4.1
Image Model      : black-forest-labs/flux-1.1-pro
Test Run ID      : 20250813_100138
Generated        : 2025-08-13 17:24:04
```

## SAMPLE SIZES

```
=====
Total Results      : 199 (AdaptiQ + GPT pairs)
Image Pairs        : 99 pairs
Cost/Latency Data  : N = 99 pairs (199 individual results)
CLIP Quality Data   : N = 93 pairs
```

## DATA QUALITY NOTES

```
=====
Missing CLIP scores: 6 pairs
Reason: API errors or invalid image formats
```

Notes on sample sizes:

- Cost/Latency/Token data captured for 100 image pairs.
- CLIP quality scores were successfully computed for 93 pairs due to intermittent API errors or invalid image formats.
- All per-pair comparisons are limited to the 93 pairs with complete data.

# Performance Summary & Key Insights

## Performance Statistics

Adaptiq:  
Avg Latency: 11.85±2.69s  
Avg Tokens: 7459±457  
Avg Cost: \$0.0086  
Total Cost: \$0.8578

GPT:  
Avg Latency: 13.94±3.33s  
Avg Tokens: 8347±1278  
Avg Cost: \$0.0099  
Total Cost: \$0.9826

## Key Metrics Comparison

ADAPTIQ vs GPT ( $\Delta$  = AdaptiQ - GPT):

Latency: -2.09s (-15.0%)

Cost: \$-0.0013 (-13.6%)

Tokens: -887 (-10.6%)

Interpretation:

Negative cost/latency = AdaptiQ better

Positive quality = AdaptiQ better

## Image Quality Summary

CLIP Score Winners (93 images):  
Adaptiq Wins: 31 (33.3%)  
GPT Wins: 41 (44.1%)  
Ties: 21 (22.6%)

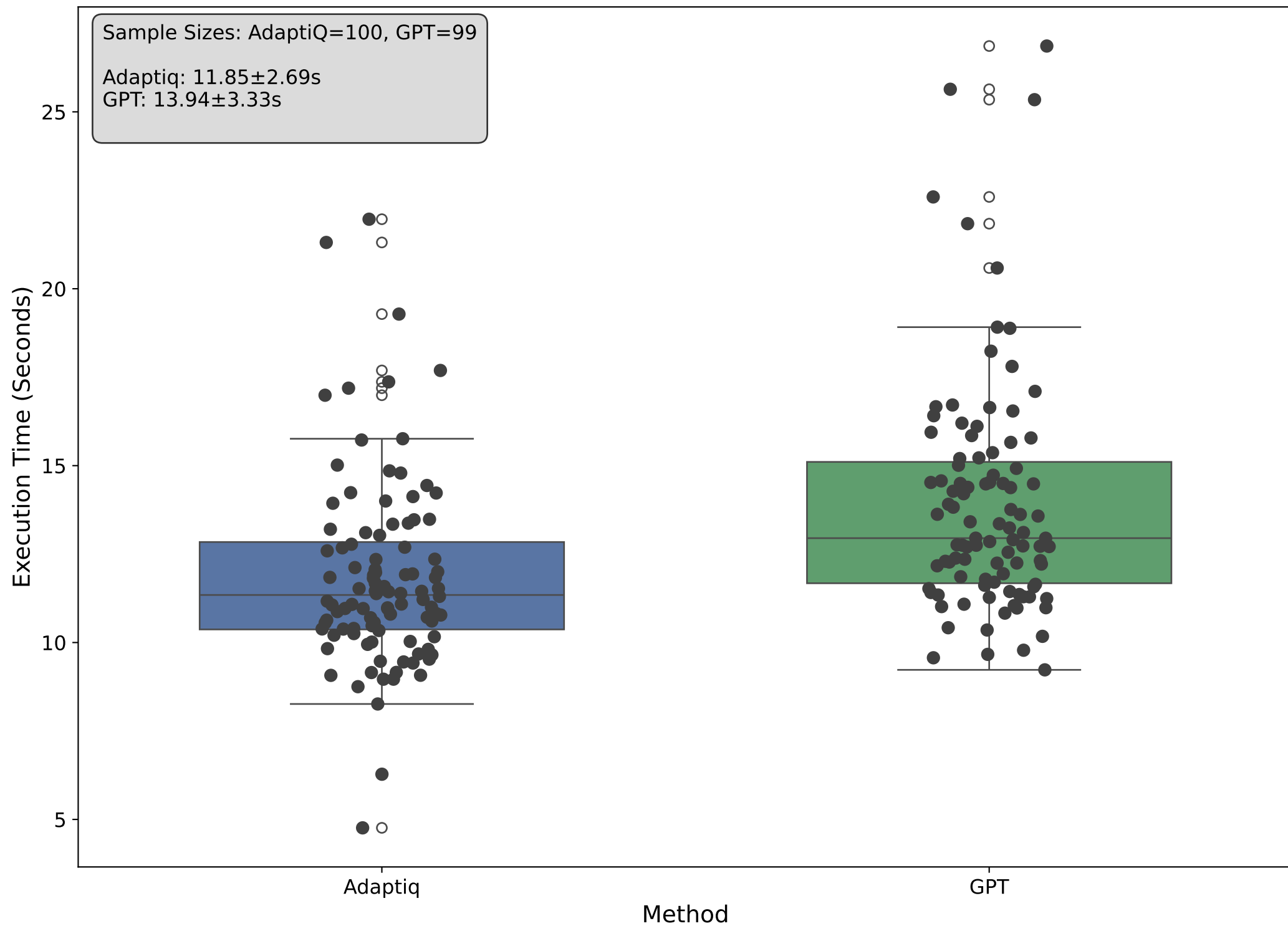
Average CLIP Scores:  
Adaptiq: 91.010  
GPT: 91.182

## Key Insights

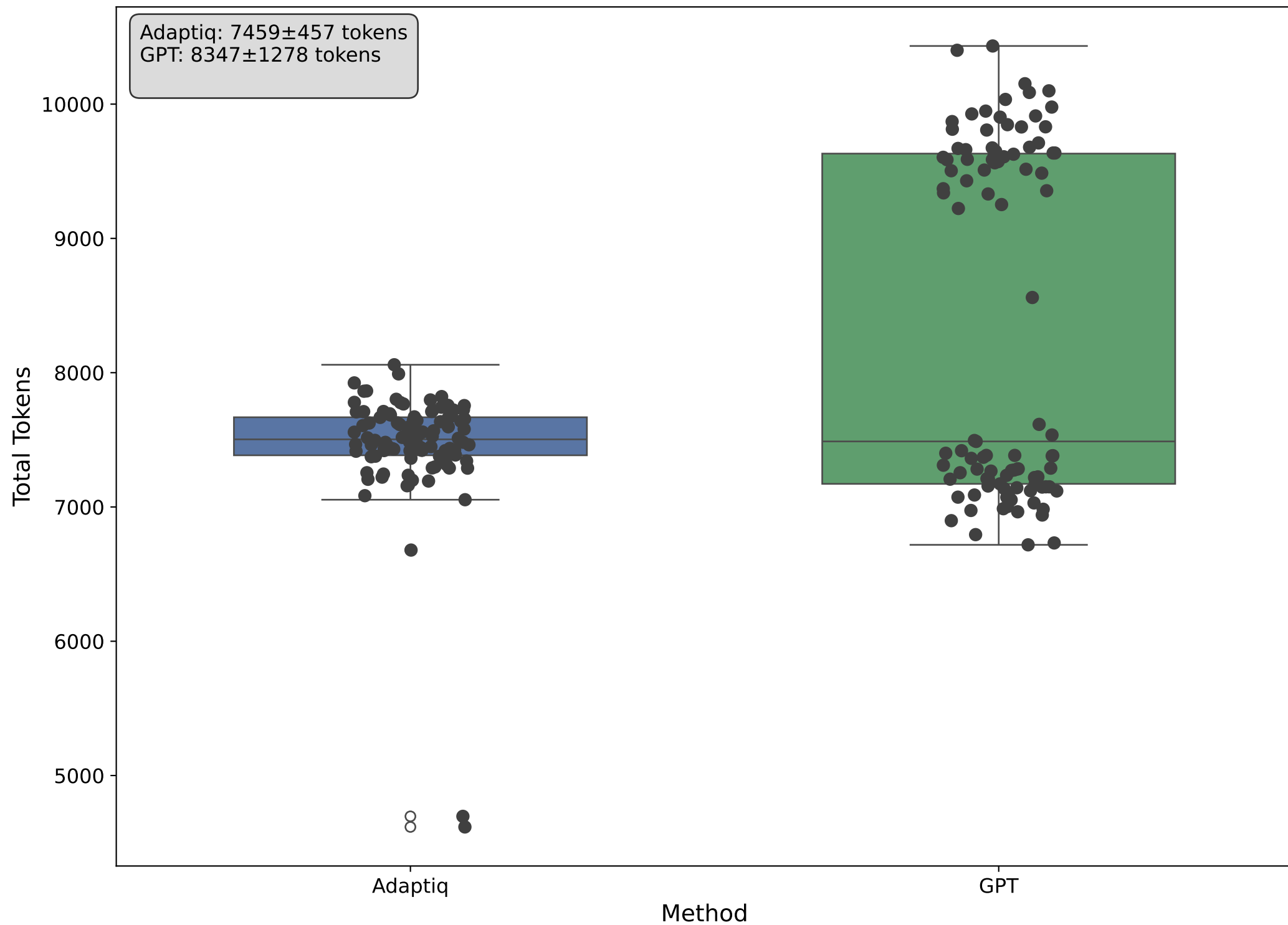
KEY FINDINGS:

- Adaptiq is 12.7% more cost-effective.
- GPT produces higher quality images more often.
- See individual pages for detailed analysis.

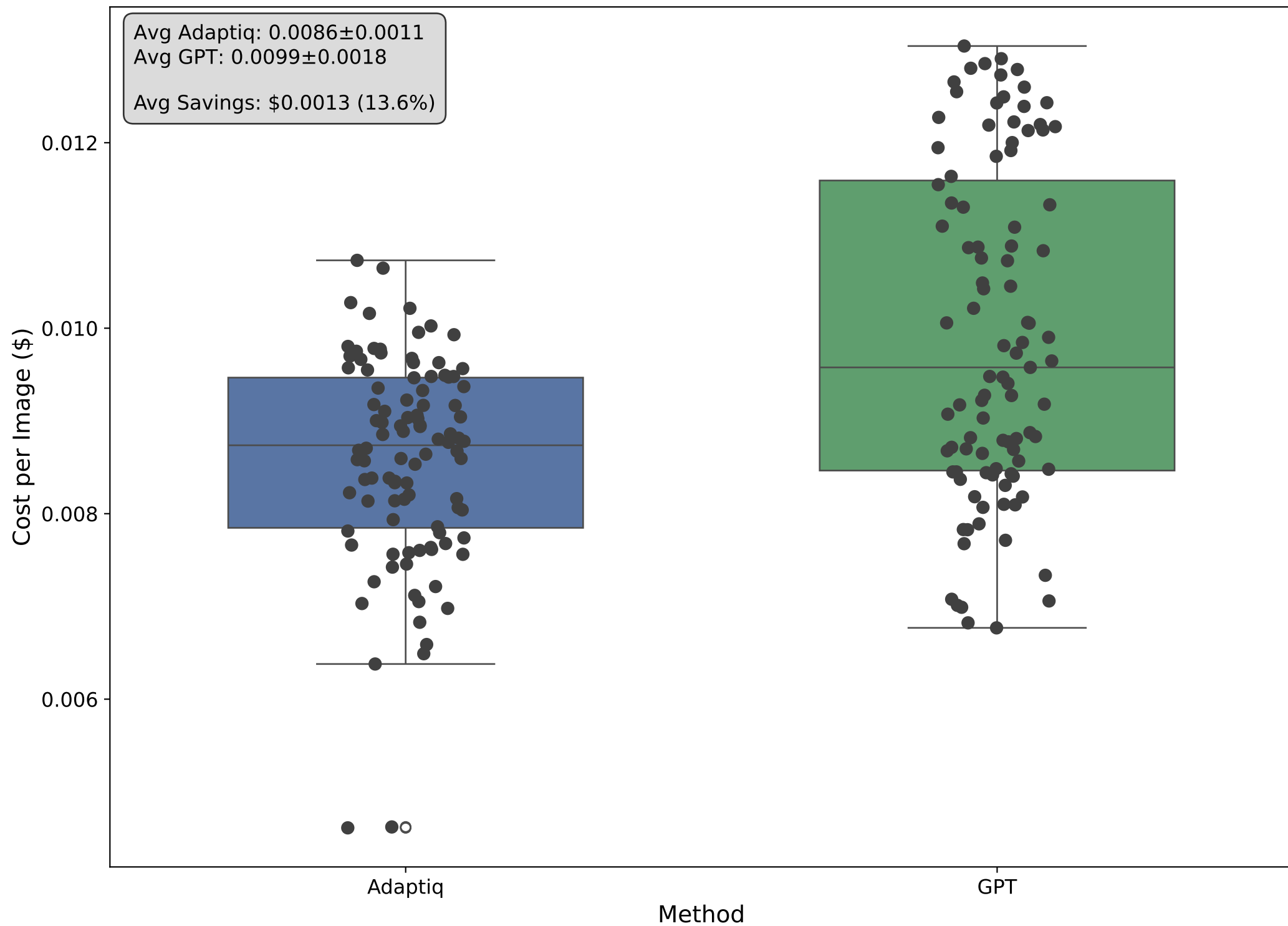
# Latency Comparison (N = 99 pairs)



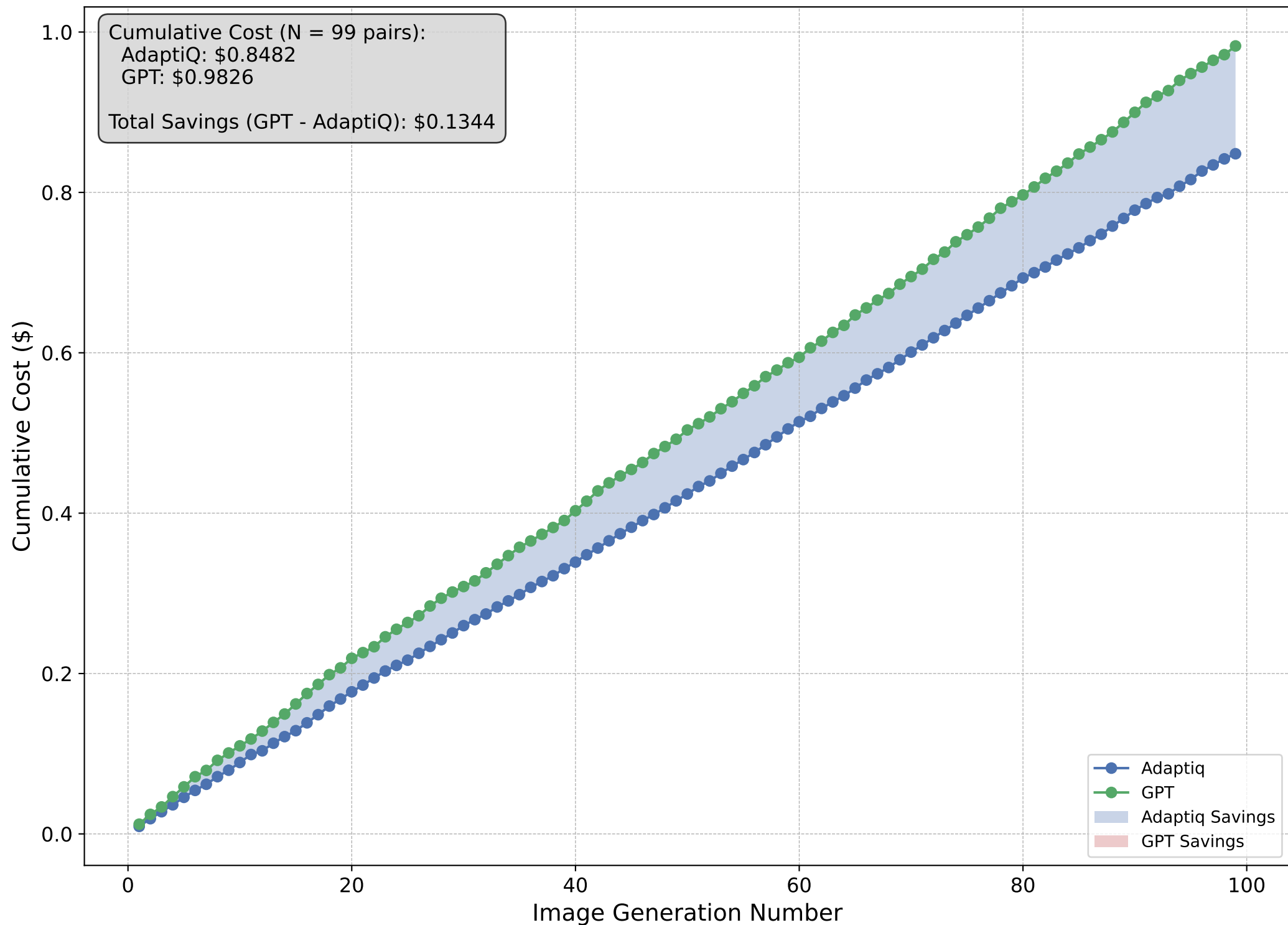
# Token Usage Comparison



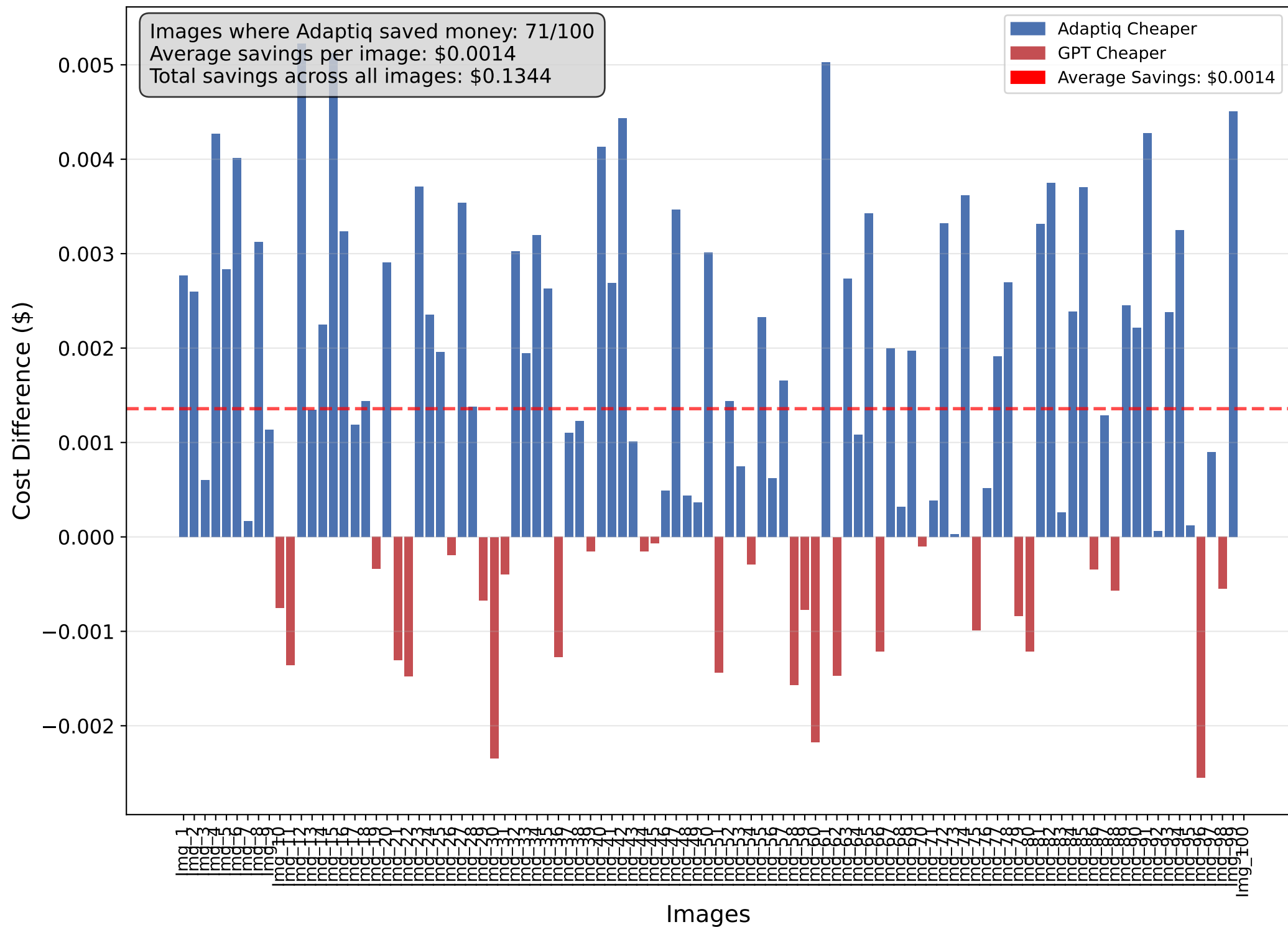
# Cost Distribution per Image



# Cumulative Cost Over Time (N = 99 pairs)

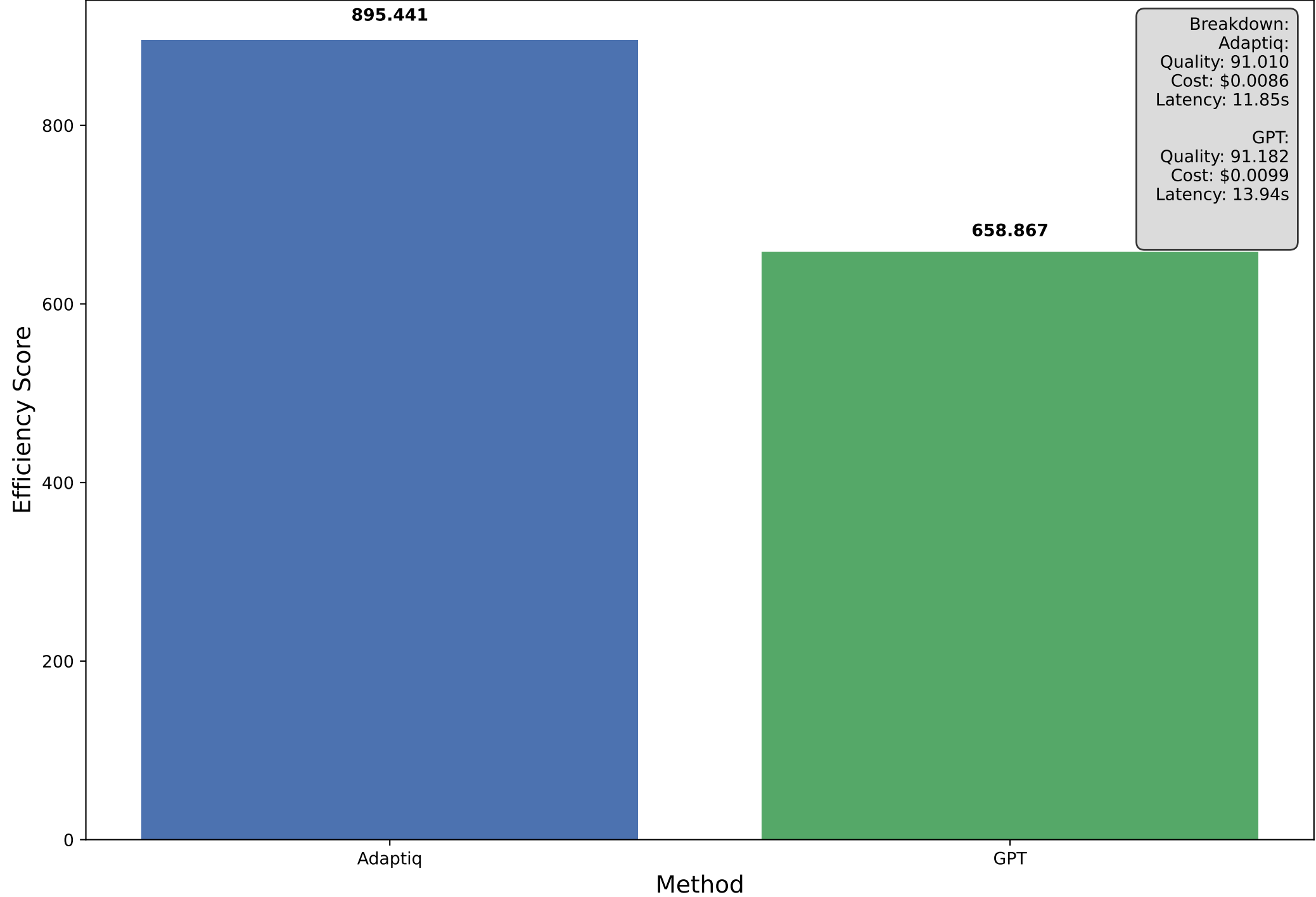


# Cost Savings per Image (GPT Cost - Adaptiq Cost)

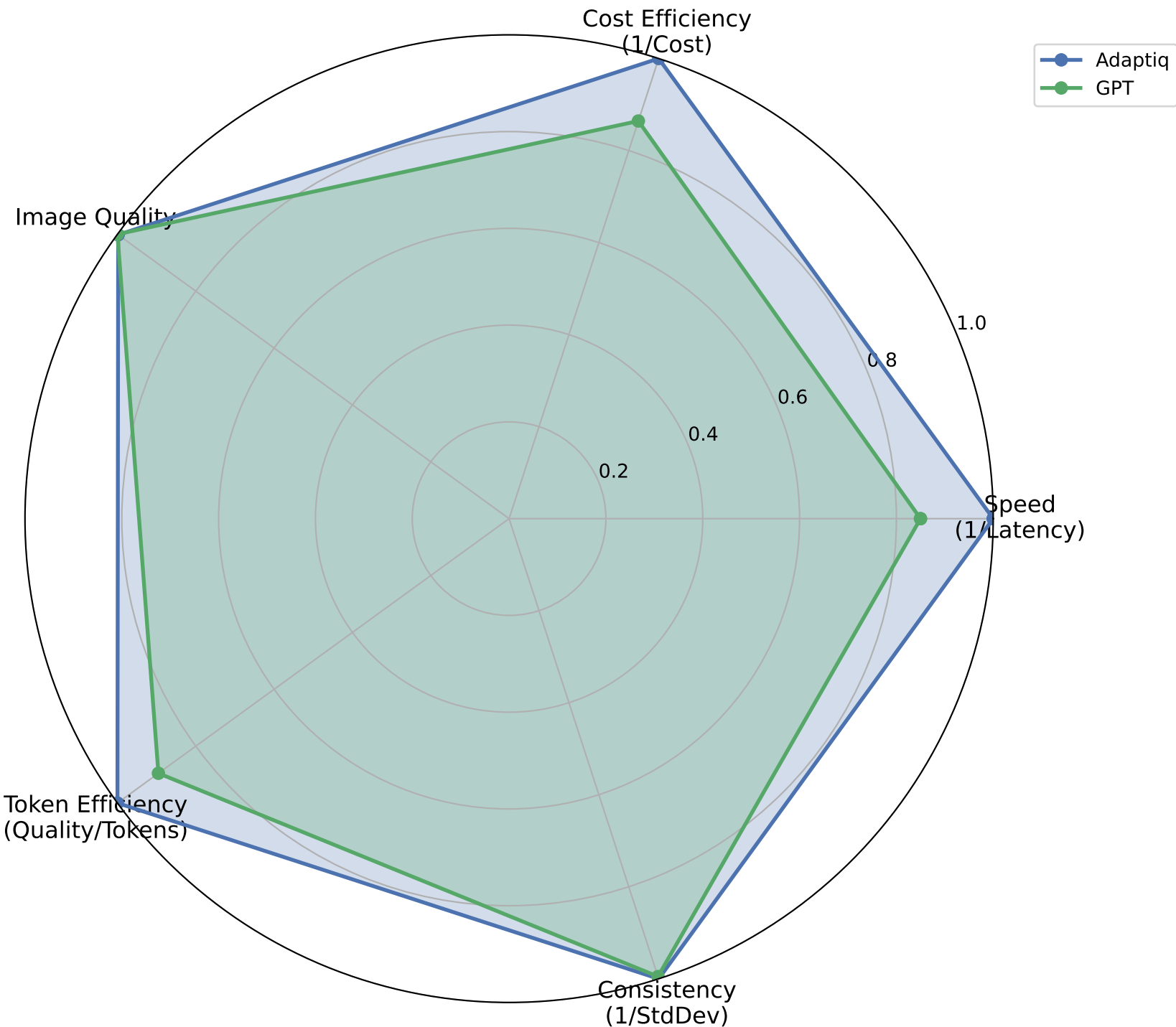




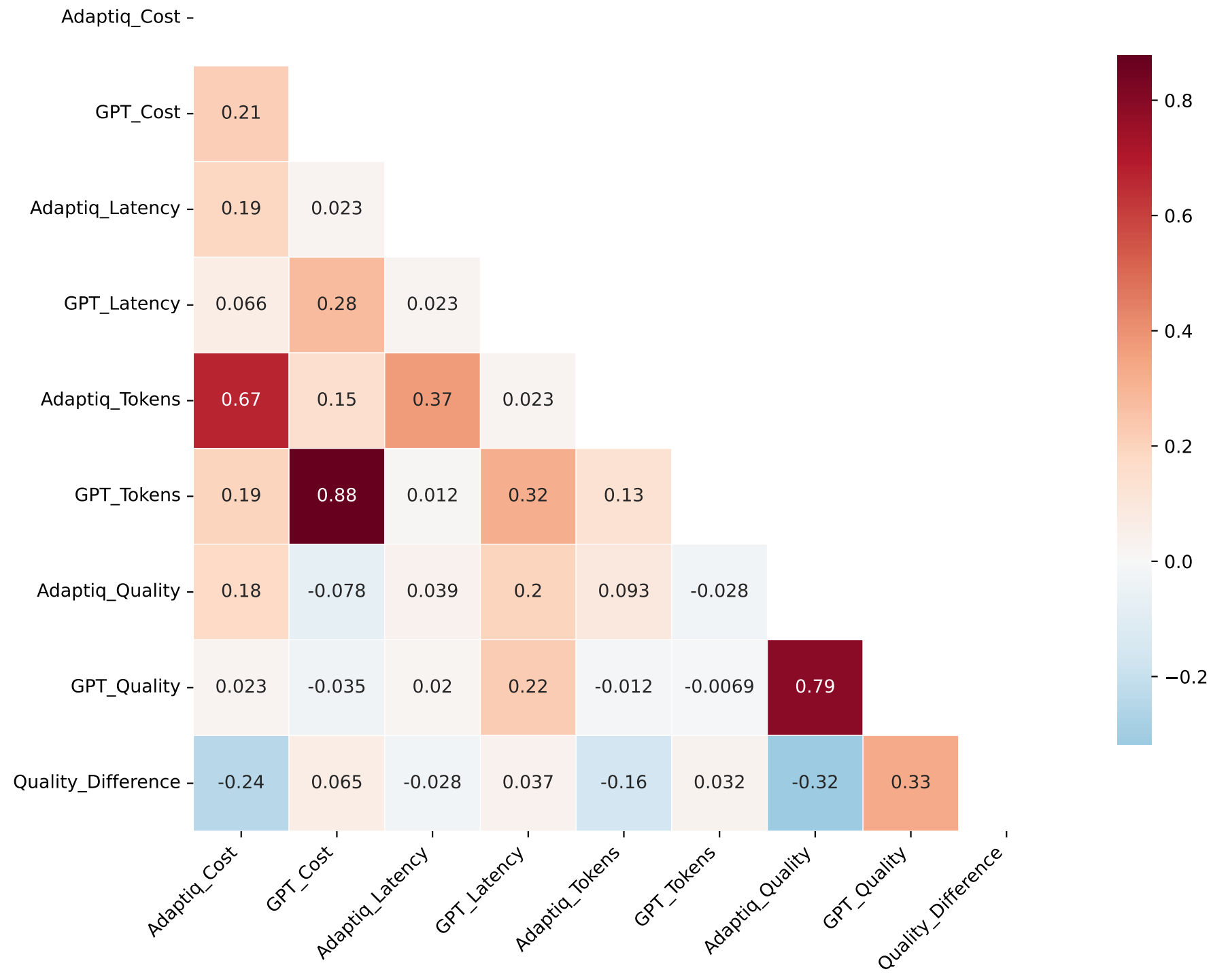
# Overall Efficiency Score (Quality per Dollar per Second)



# Performance Radar Chart (Normalized Scores)

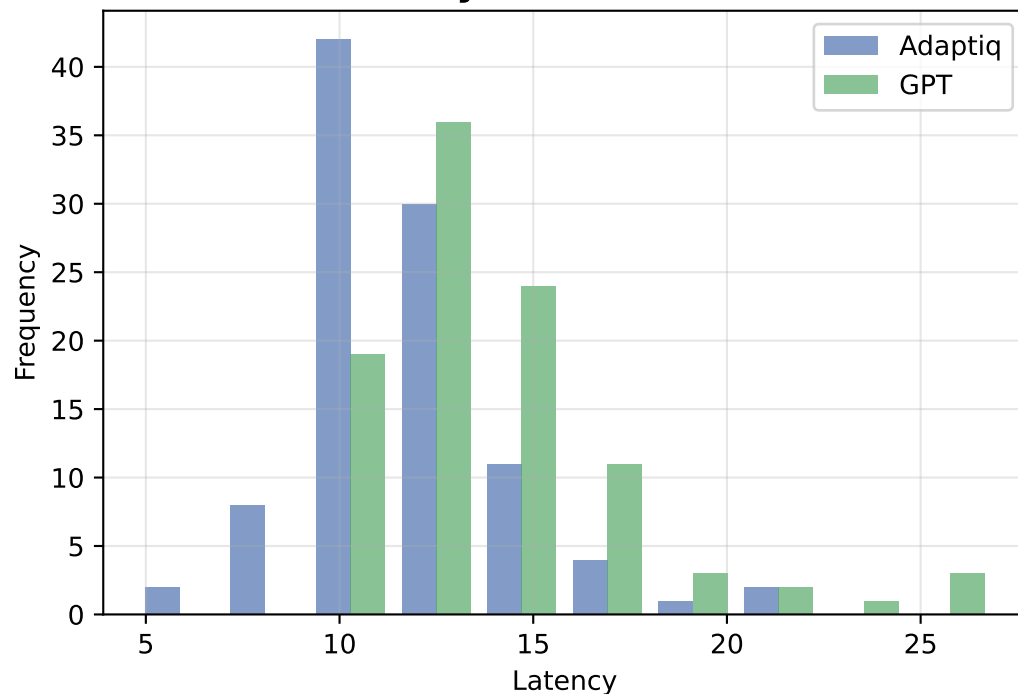


## Correlation Matrix of Performance Metrics

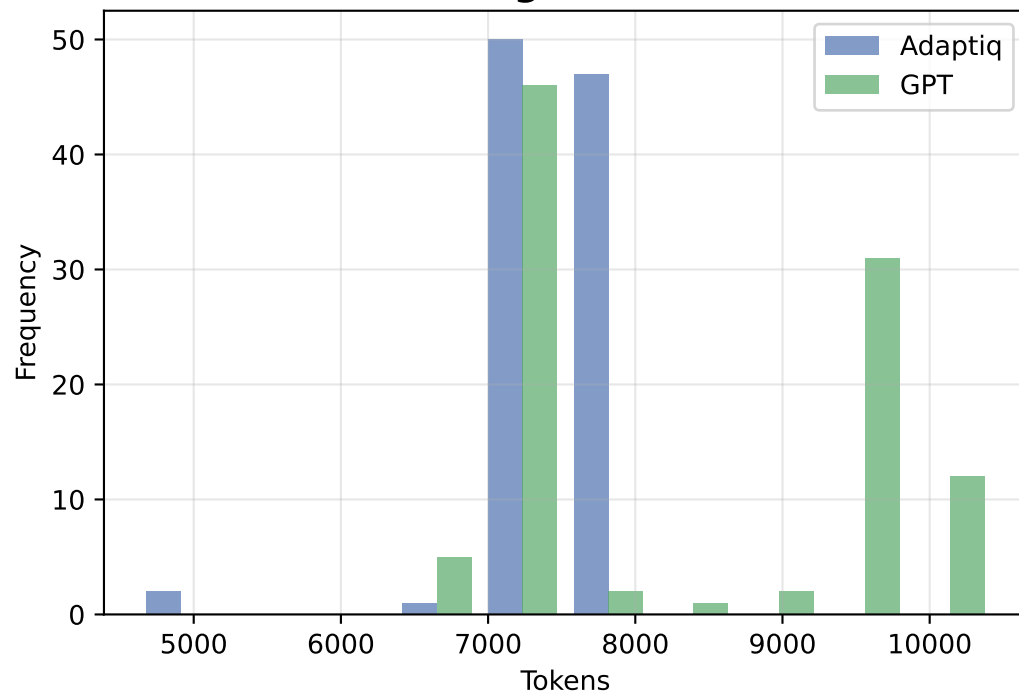


# Distribution Comparison: Adaptiq vs GPT

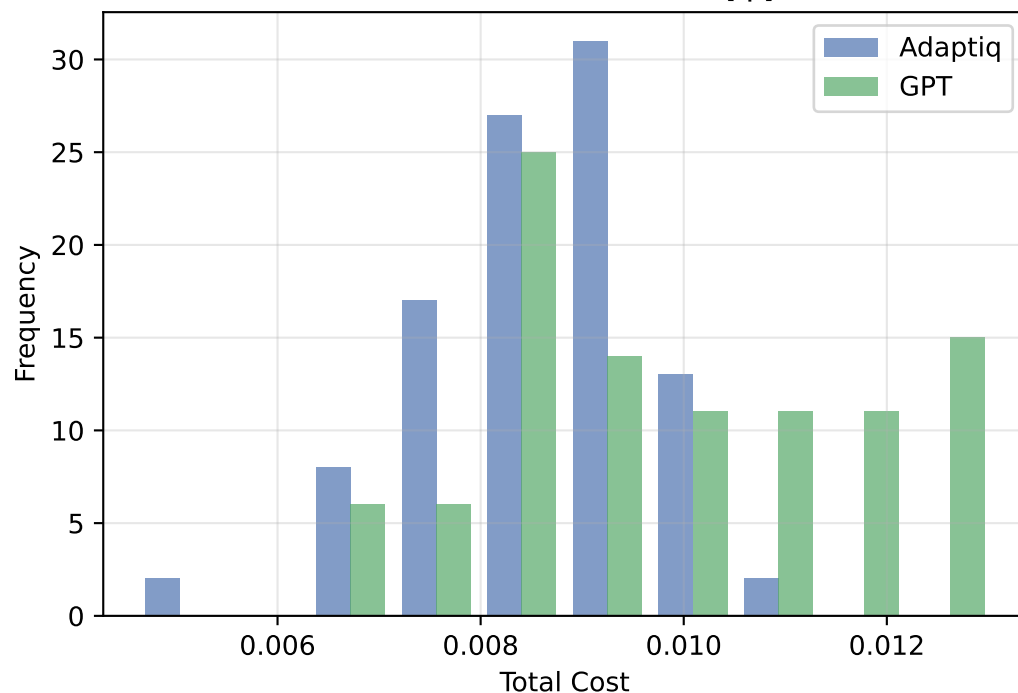
## Latency Distribution (s)



## Token Usage Distribution



## Cost Distribution (\$)



## Statistical Tests

Mann-Whitney U Test Results:

latency:  $p=0.0000$  \*\*\*

tokens:  $p=0.3656$  ns

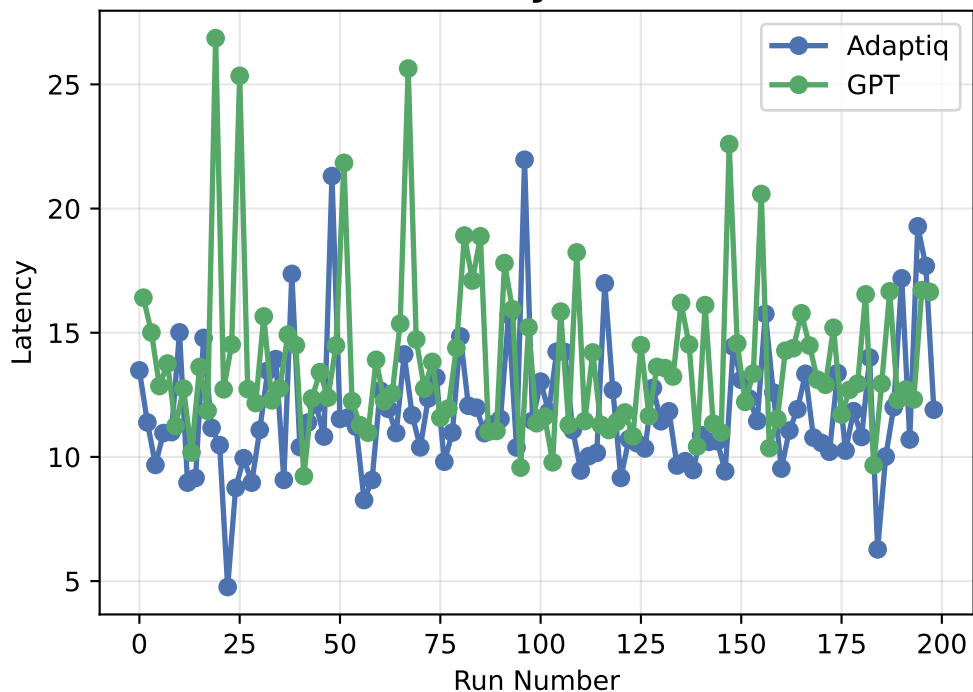
total\_cost:  $p=0.0000$  \*\*\*

\*  $p<0.05$ , \*\*  $p<0.01$ , \*\*\*  $p<0.001$

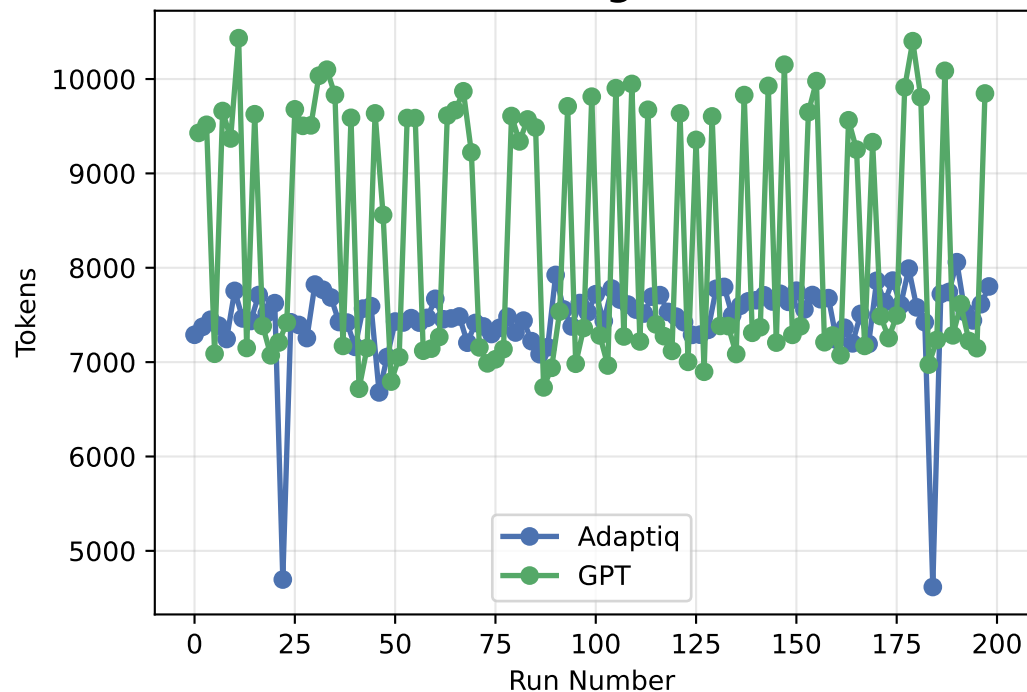
ns = not significant

# Performance Trends Over Sequential Runs

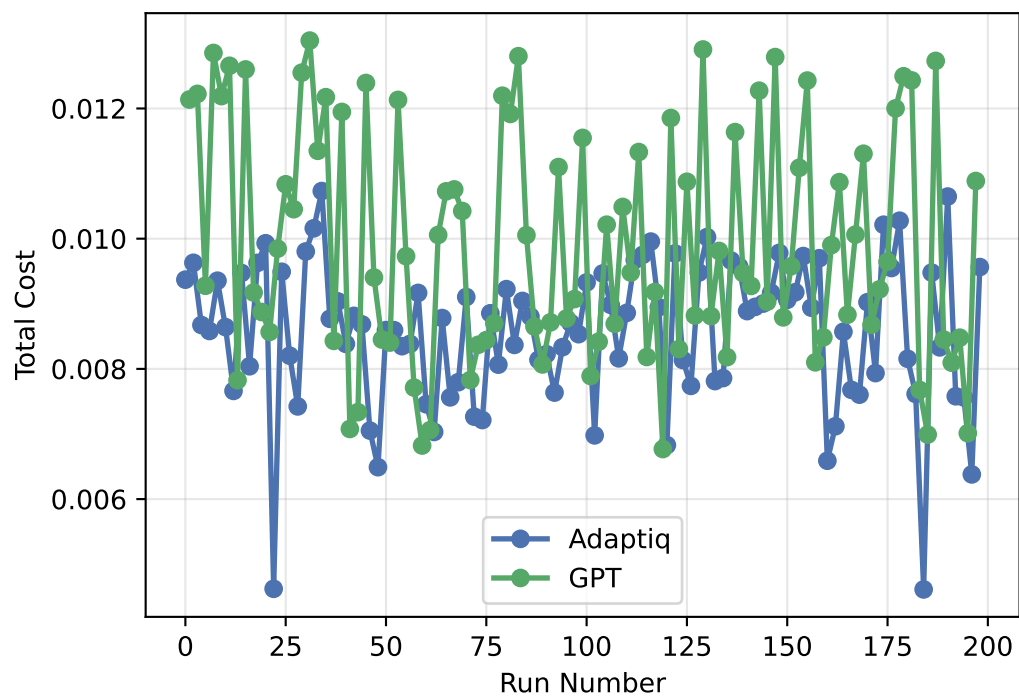
## Latency Trend



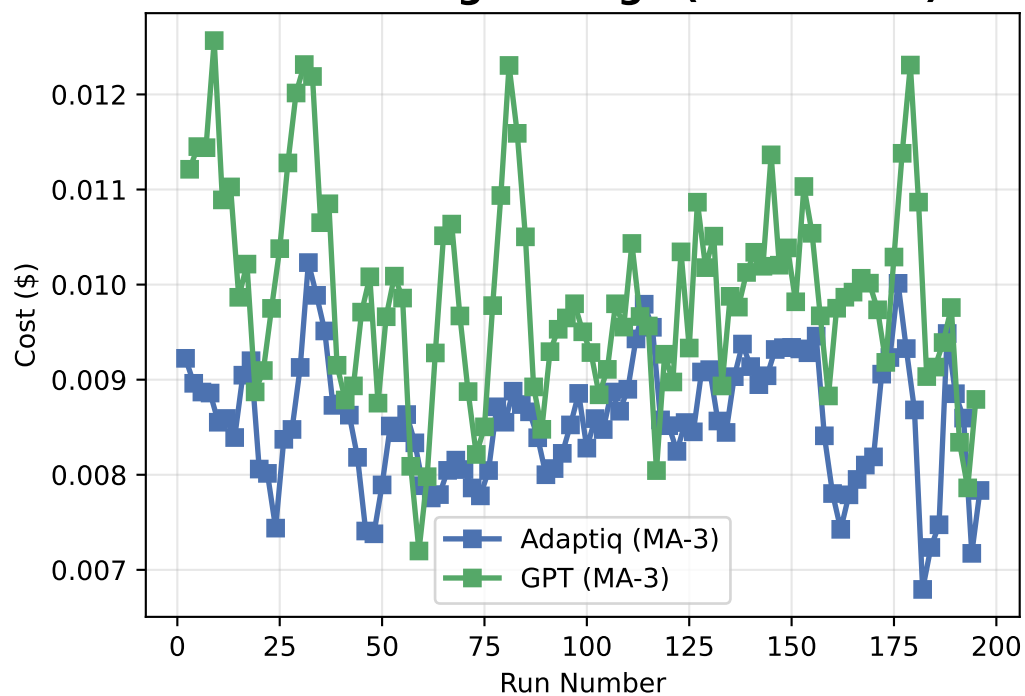
## Token Usage Trend



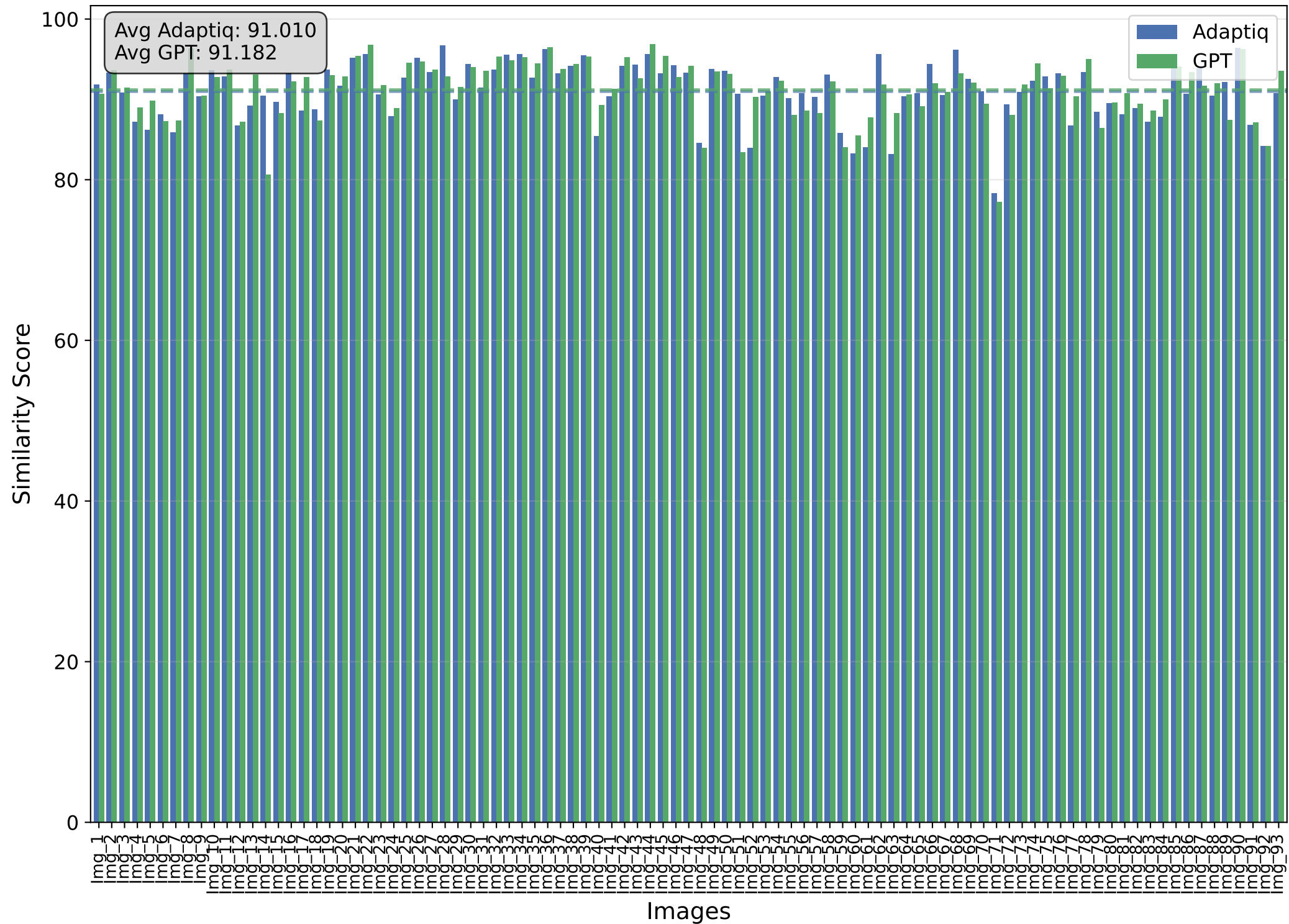
## Cost Trend



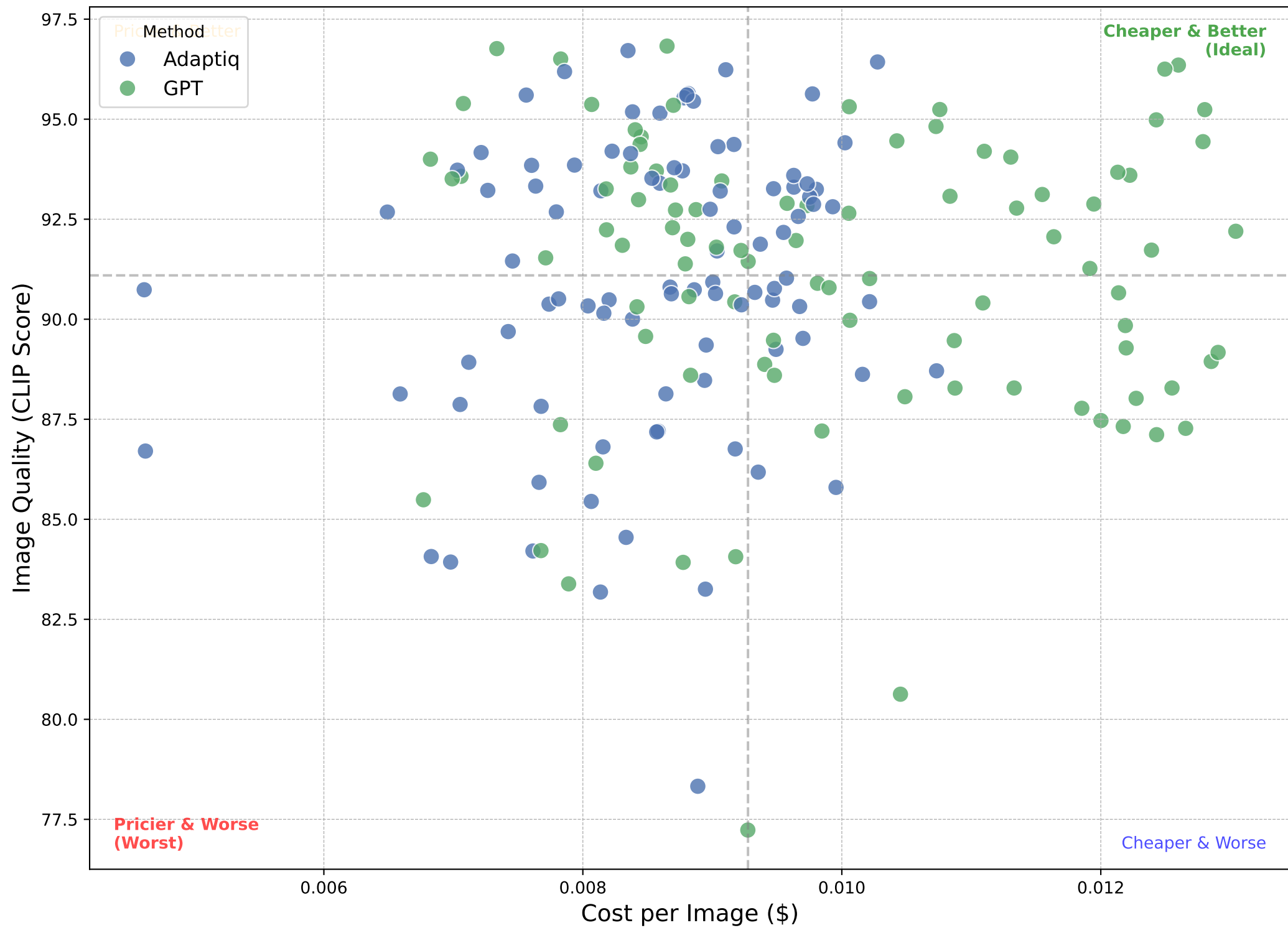
## Cost Moving Average (Window=3)



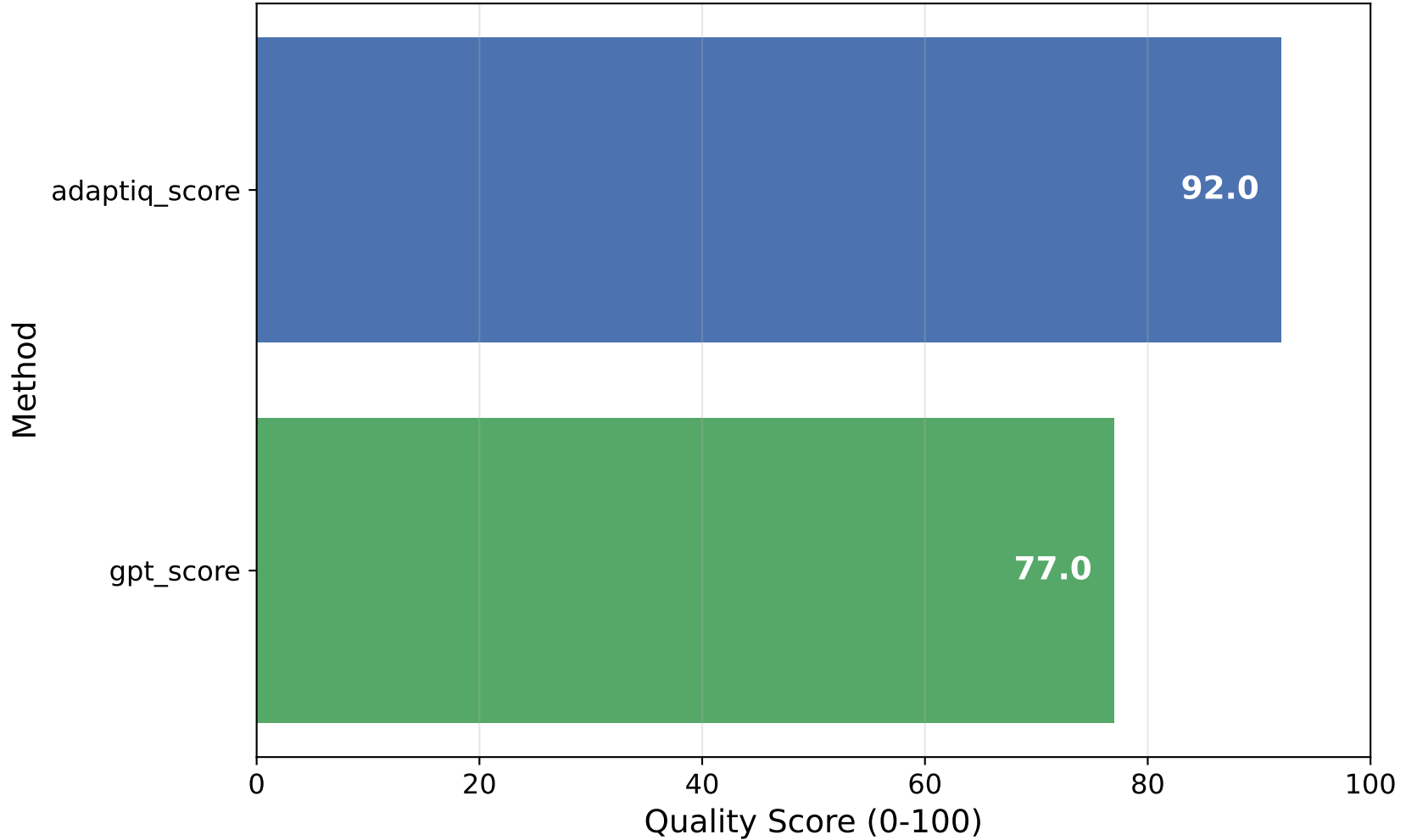
# CLIP Similarity Scores (All 93 Images)



# Image Quality vs. Generation Cost



# Prompt Engineering Quality Scores



**Summary:**  
Adaptiq outperforms GPT method across all criteria due to its granular, modular workflow, explicit error handling, and rigid structural design. It is more precise, reliable, easier to debug, and more effective at consistently producing optimal, safe prompts. While GPT's approach covers the core requirements, it trades strictness for flexibility, which can reduce effectiveness and traceability in demanding or high-stakes workflows.