# PyTorch vs Triton Inference Server
## OpenCLIP ViT-B-32 on RTX A4000
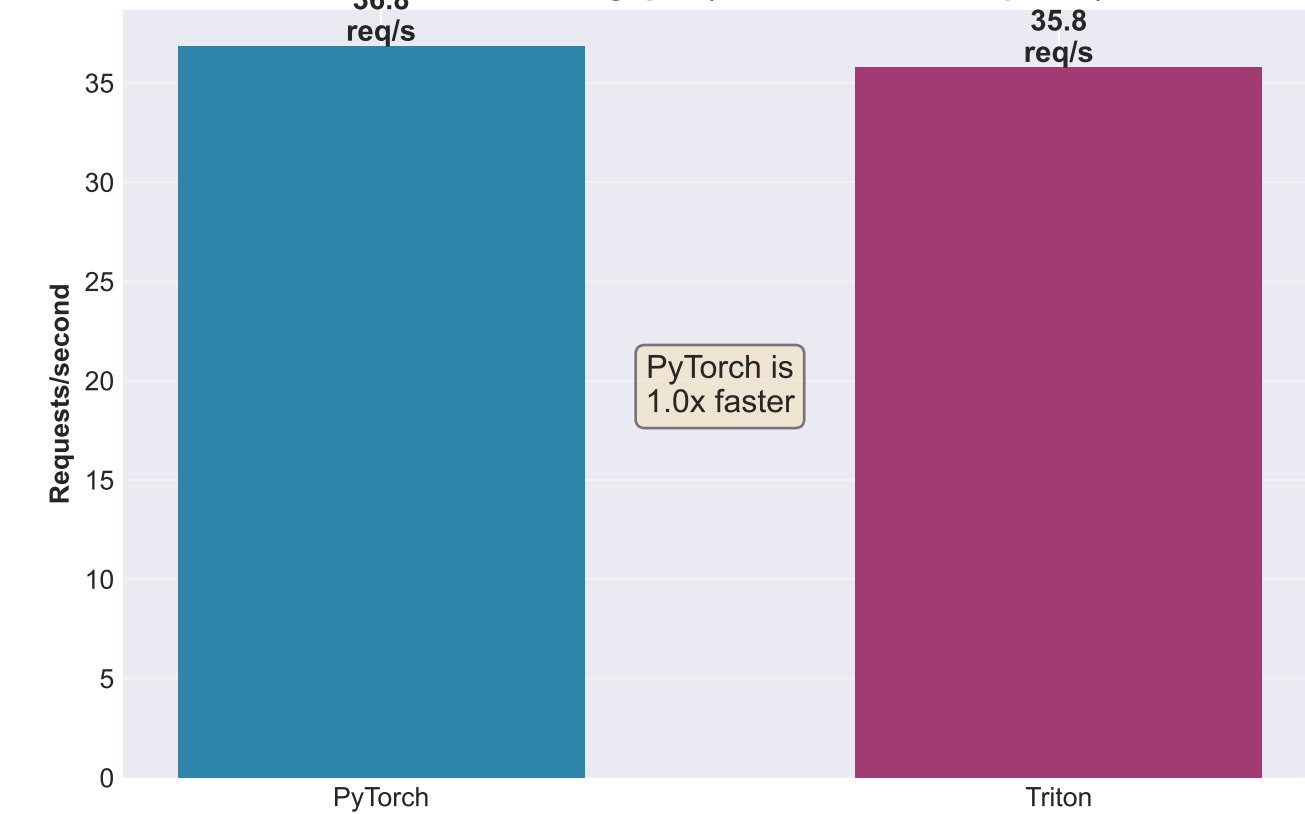
### Single-Image Latency



| | PyTorch | Triton |
|---|---|---|
| Mean | 189 | 319 |
| p50 | 167 | 229 |
| p95 | 250 | 534 |
| p99 | 470 | 2106 |

Latency (ms)

### Batch Processing Throughput

| | PyTorch | Triton |
|---|---|---|
| Batch 4 | 11.1 | 5.1 |
| Batch 8 | 17.9 | 5.0 |
| Batch 16 | 25.6 | 4.5 |
| Batch 32 | 29.9 | 5.7 |

Images/second

### Concurrent Throughput (16 workers, 200 requests)

PyTorch: 36.8 req/s
Triton: 35.8 req/s

PyTorch is 1.0x faster

Requests/second

### Cost to Process 10,000 Images
(@ $0.20/hour GPU)

PyTorch: $0.015 (4.5 min)
Triton: $0.016 (4.7 min)

Save 3% with PyTorch

Cost (USD)