# PyTorch vs Triton Inference Server
## OpenCLIP ViT-B-32 on RTX 3070
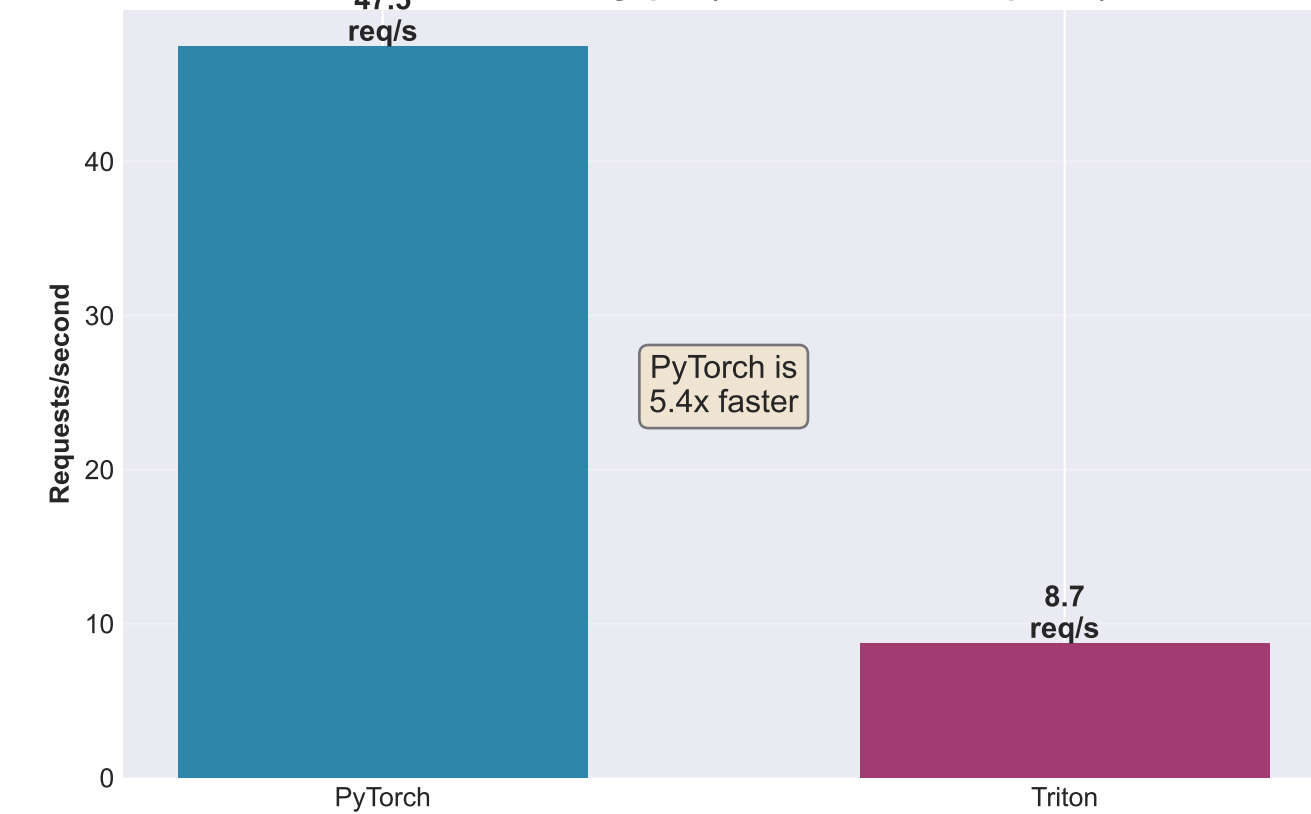
### Single-Image Latency

| Metric | PyTorch | Triton |
|---|---|---|
| Mean | 206 | 724 |
| p50 | 162 | 564 |
| p95 | 359 | 1874 |
| p99 | 892 | 2031 |

### Batch Processing Throughput

| Batch | PyTorch | Triton |
|---|---|---|
| Batch 4 | 12.5 | 2.3 |
| Batch 8 | 15.0 | 3.1 |
| Batch 16 | 29.2 | 3.0 |
| Batch 32 | 43.9 | 3.3 |

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

PyTorch: 47.5 req/s
Triton: 8.7 req/s

PyTorch is 5.4x faster

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

PyTorch: $0.012 (3.5 min)
Triton: $0.064 (19.1 min)

Save 82% with PyTorch