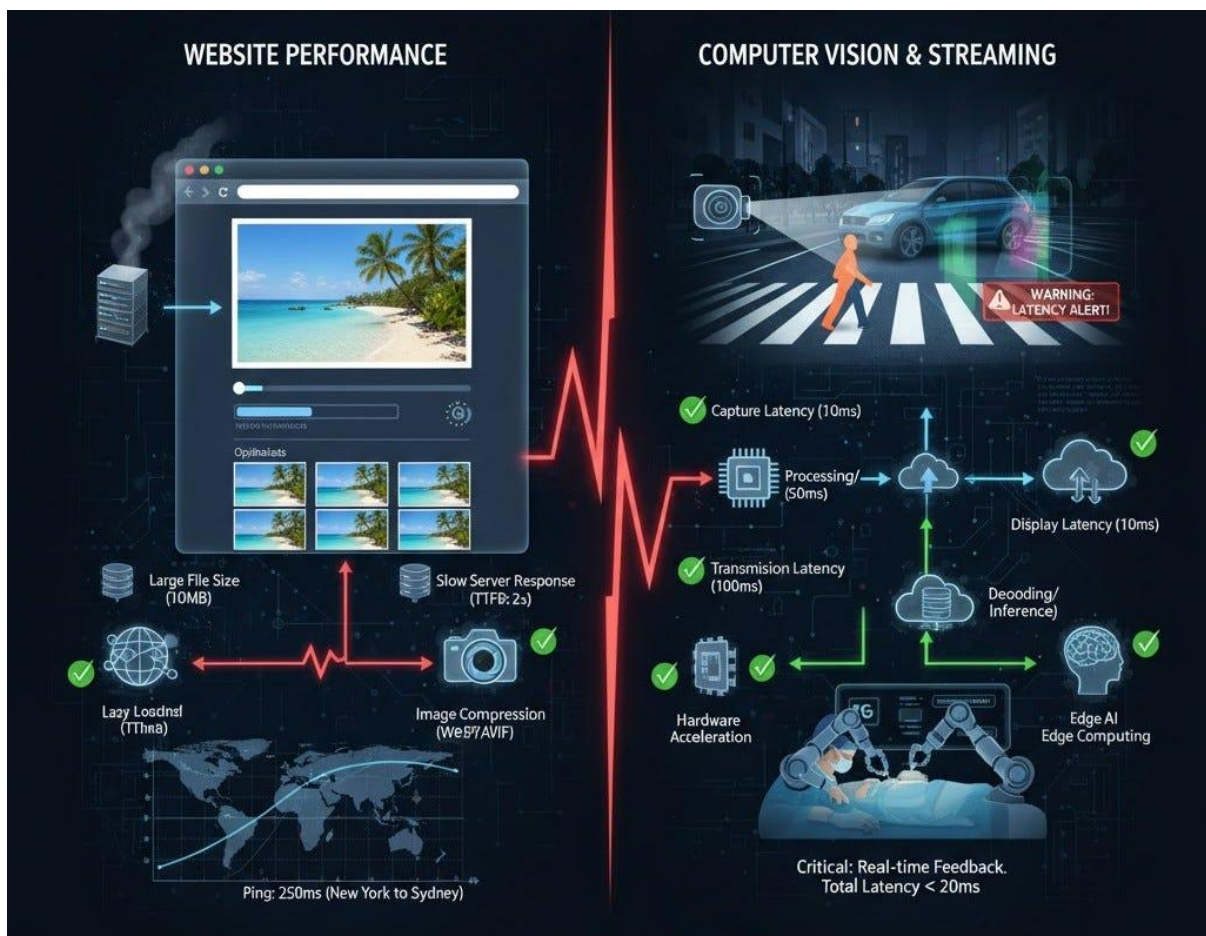


# Understanding Latency in Multi-Agent GenAI Systems

Why your AI Agents are slower than you think (And How to Fix It)

A practical guide to measuring and optimizing performance in production multi-agentic systems



By - Rajesh Srivastava

## Table of Contents

Introduction .....	4
The Evolution: From Chatbots to Multi-Agent Systems .....	4
Part 1: The Foundation - Core Latency Metrics.....	5
1. The Basics Everyone Should Track.....	5
a) Average Latency (Mean Response Time) .....	5
b) Median Latency (P50) .....	5
c) Tail Latencies: Where Real Problems Hide .....	5
2. Operational Metrics That Reveal Bottlenecks .....	6
a) Request Queuing Delay .....	6
b) Processing/Service Time .....	6
c) Network Latency .....	6
d) Cold Start Latency .....	6
e) Timeout Rate & Jitter .....	6
Part 2: GenAI-Specific Latency Metrics .....	7
1. Pre-Processing: Before the Model Even Sees Your Request .....	7
a) Prompt Processing Latency .....	7
b) Embedding Latency.....	7
2. Retrieval (RAG Systems): The Hidden Time Sink .....	7
a) Retrieval Latency.....	7
3. Model Inference: Understanding Token-Level Performance .....	7
a) TTFT (Time to First Token) .....	7
b) Token Generation Rate.....	7
c) TTLX (Time to Last Token).....	8
4. Post-Processing: The Final Mile.....	8
a) Post-Processing Latency .....	8
b) Response Aggregation Latency .....	8
5. The Metrics That Matter Most for User Experience .....	8
a) User-Perceived Latency.....	8
b) Agent Coordination Latency.....	8
Part 3: The Reality of Multi-Agent Architectures .....	9
1. The Cast of Characters .....	9
a) Orchestrator/Controller Agent.....	9
b) Retriever Agent .....	9
c) Reasoning Agent.....	9
d) Generator Agent .....	9

e) Evaluator/Critic Agent .....	9
f) Memory/Knowledge Agent (optional).....	9
g) Tool/Action Agent (optional) .....	9
2. The Latency Formula .....	10
a) Sequential Operations (These Add Up).....	10
b) Parallel Operations (Take Maximum Time).....	10
c) Overhead (The Hidden Tax).....	10
d) Cache Adjustment.....	11
Part 4: The Latency Flow .....	12
Part 5: A Real-World Example .....	13
Part 6: How to Actually Measure This.....	15
OpenTelemetry: .....	15
Anatomy of a Good Trace Span.....	15
Part 7: Optimization Strategies That Actually Work.....	17
1. Parallelize Everything You Can (if possible) .....	17
2. Cache Aggressively (But intelligently) .....	17
3. Stream Everything.....	18
4. Optimize Context Windows.....	18
5. Batch Similar Requests.....	19
6. Eliminate Cold Starts .....	19
7. Co-locate Services .....	20
8. Right-Size Your Models.....	21
9. Make Non-Critical Operations Async .....	21
10. Monitor Tail Latencies Continuously.....	22
Part 8: What Matters Most .....	23
1. For User Experience .....	23
2. For System Performance .....	23
3. For Cost Efficiency .....	23
Final Thoughts.....	25

# Introduction

You have built a sophisticated multi-agent AI system. It's intelligent, it's powerful, and... it's slow.

Your users are waiting **8 seconds** for responses. Your demo looks impressive, but in production, it feels sluggish. You thought the LLM was the bottleneck, but when you checked, the model inference is only **40%** of the total latency.

## What is going on?

Welcome to the complex world of multi-agent system latency (and of course, distributed systems), where the coordination between components matters just as much as the performance of individual pieces.

This article breaks down everything you need to know about latency in modern GenAI architectures → what to measure, why it matters and most importantly, how to optimize it.

---

## The Evolution: From Chatbots to Multi-Agent Systems

Early GenAI applications were simple:

User sends prompt → Model generates response → Done

### Modern GenAI systems look more like this:

- User sends request
- Orchestrator decides which agents to invoke
- Retrieval agent searches vector databases
- Memory agent fetches conversation history
- Reasoning agent plans the response strategy
- Generator agent calls the LLM
- Evaluator agent validates output quality
- Response gets assembled and returned

**Each step adds latency.** And here is what most teams don't realize they don't know where the time is actually going until they measure it properly.

# Part 1: The Foundation- Core Latency Metrics

Before we dive into GenAI specific metrics, let's establish the baseline. These metrics apply to any distributed system, whether you are running a simple REST API or a complex AI pipeline.

## 1. The Basics Everyone Should Track

### a) Average Latency (Mean Response Time)

The arithmetic mean of all your response times. It's useful for trends, but it lies to you. A few slow outliers can make a fast system look mediocre.

*Example: 90 requests at 200ms + 10 requests at 5000ms = 680ms average. But 90% of users experienced 200ms.*

### b) Median Latency (P50)

The middle value, literally the 50th percentile. Half your users are faster, half are slower. This is your **typical** user experience.

*This is often 2-5x better than your average, which tells you outliers are skewing your averages.*

### c) Tail Latencies: Where Real Problems Hide

This is where things get interesting and where most systems fail.

- **P90 (90th Percentile):** 10% of your users experience this latency or worse
- **P95 (95th Percentile):** 5% of your users experience this latency or worse
- **P99 (99th Percentile):** 1% of your users experience this latency or worse
- **P99.9 (99.9th Percentile):** 0.1% of your users (1 in 1000) experience this latency or worse

**Here is why tail latencies matter more than averages:** At scale, your P99 is somebody's consistent experience.

If you have 1 million daily users and your P99 is 10 seconds, that means 10,000 users per day are waiting 10+ seconds. That is not an edge case, that's a significant user segment having a terrible experience.

## 2. Operational Metrics That Reveal Bottlenecks

### a) Request Queuing Delay

Time requests spend waiting before processing even begins. High queuing means you are resource-constrained, your agents can't keep up with demand.

*Common culprit: Not enough worker threads, containers, or GPU instances.*

### b) Processing/Service Time

Actual execution duration, excluding queue time. This tells you if your code is efficient or if you are doing unnecessary work.

### c) Network Latency

Round-trip time between client and server. For remote API calls (OpenAI, Anthropic, Gemini, Cohere etc), this can be 50-200ms before you even start processing.

### d) Cold Start Latency

The warm-up delay when a function or container starts for the first time. In serverless environments, this can add 500ms to 5 seconds on the first request.

*Pro tip: Keep instances warm with periodic pings, or use provisioned concurrency.*

### e) Timeout Rate & Jitter

Timeout rate shows how often requests exceed your threshold. Jitter measures variability, high jitter means inconsistent performance, which users hate more than consistently slow responses.

---

## Part 2: GenAI-Specific Latency Metrics

Now let's layer in the GenAI specific metrics that traditional monitoring doesn't capture.

### 1. Pre-Processing: Before the Model Even Sees Your Request

#### a) Prompt Processing Latency

Time to prepare the prompt → template rendering, variable substitution, context assembly. This is often overlooked but can add 50-200ms.

#### b) Embedding Latency

Converting text to vectors for semantic search. Typically 50-200ms per operation, but it compounds when you are embedding user queries, documents, and retrieved chunks.

### 2. Retrieval (RAG Systems): The Hidden Time Sink

#### a) Retrieval Latency

Time spent searching vector databases or making retrieval API calls.

- **Pinecone/Weaviate:** 50-200ms (purpose-built for speed)
- **Elasticsearch:** 100-500ms (more flexible, slightly slower)
- **Custom solutions:** Highly variable (often 200-1000ms without optimisation)

*This is where caching delivers massive wins (more on that later)*

### 3. Model Inference: Understanding Token-Level Performance

This is what everyone focuses on.

#### a) TTFT (Time to First Token)

The delay before the model outputs its first token. This is what users feel as **response delay**.

#### b) Token Generation Rate

How fast tokens are produced during streaming, measured in tokens per second.

### c) TTLX (Time to Last Token)

Total generation time for the complete response.

- **Short responses (100 tokens):** 2-4 seconds
- **Medium responses (500 tokens):** 5-12 seconds
- **Long responses (1000+ tokens):** 15-30+ seconds

## 4. Post-Processing: The Final Mile

### a) Post-Processing Latency

Time for output validation, formatting, ranking, or evaluation. Often 50-200ms, but can balloon to 500ms+ if you are running complex validation logic.

### b) Response Aggregation Latency

Time to combine outputs from multiple agents. In sequential architectures, this happens last. In parallel architectures, you are waiting for the slowest agent.

## 5. The Metrics That Matter Most for User Experience

### a) User-Perceived Latency

How long users wait before seeing the first visible token. This is your most important metric for conversational AI.

**Formula:** *Network latency + Queue time + Orchestration + TTFT*

### b) Agent Coordination Latency

Time spent in inter-agent communication, planning, or tool execution. In distributed multi-agent systems, this can secretly dominate your latency budget.

- **Same region/VPC:** 1-5ms per hop
- **Cross-region:** 50-200ms per hop
- **Cross-cloud provider:** 100-300ms per hop

## Part 3: The Reality of Multi-Agent Architectures

Here is where theory meets practice. A production multi-agent system typically includes:

### 1. The Cast of Characters

#### **a) Orchestrator/Controller Agent**

The traffic cop. Manages task flow, decides which agents to invoke, and aggregates results. Adds 50-200ms of coordination overhead.

#### **b) Retriever Agent**

Fetches relevant context from vector databases or APIs. Can take 100-300ms, or 5ms if you get a cache hit.

#### **c) Reasoning Agent**

Interprets user intent, plans subtasks, makes routing decisions. Typically 50-150ms of planning time.

#### **d) Generator Agent**

Calls the LLM to produce outputs. This is your 2-20 second time sink.

#### **e) Evaluator/Critic Agent**

Validates, scores, or refines responses. Adds 100-500ms but improves quality.

#### **f) Memory/Knowledge Agent (optional)**

Maintains context between interactions. Fast if cached (5-20ms), slower if fetching from persistent storage (100-400ms).

#### **g) Tool/Action Agent (optional)**

Executes API calls, database updates, or external actions. Highly variable, anywhere from 100ms to several seconds depending on the external service.

## 2. The Latency Formula

**Total E2E Latency** = Sequential Operations + max(Parallel Operations) + Overhead + Cache Adjustment

Let me break this down:

### a) Sequential Operations (These Add Up)

- Network latency (50-200ms)
- API gateway routing (10-50ms)
- Queue wait time (0-500ms, load-dependent)
- Orchestrator decision logic (50-200ms)
- Prompt assembly (20-100ms)
- LLM prompt encoding (50-300ms, scales with context size)
- TTFT (100ms-2s)
- Token generation (2-20s)
- Post-processing (50-200ms)
- Response formatting (10-50ms)

### b) Parallel Operations (Take Maximum Time)

When agents run concurrently, you only wait for the slowest one:

```
max(  
  Retrieval Agent → 250ms,  
  Memory Agent → 180ms,  
  Embedding Service → 100ms  
) = 250ms total
```

→ **Not 530ms** (if they ran sequentially).

### c) Overhead (The Hidden Tax)

- Inter-agent communication (10-80ms per hop)
- Serialization/deserialization (5-20ms per operation)

- Logging & metrics (5-30ms)
- Error handling & retries (0-2000ms if triggered)

#### **d) Cache Adjustment**

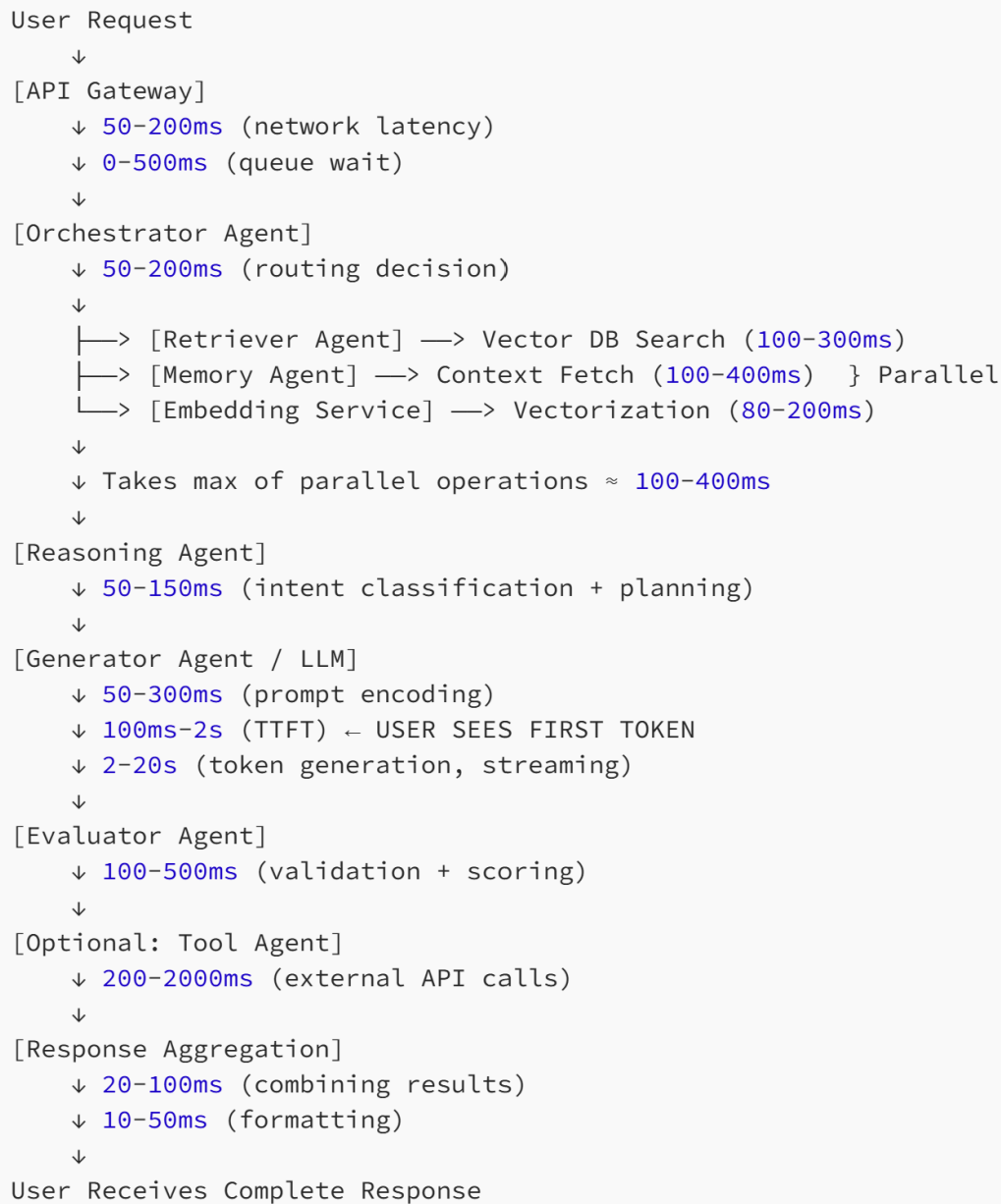
- **Cache hit:** 90-99% latency reduction on cached components
- **Cache miss:** Full latency stack

This is why cache hit rate is one of your most important KPIs.

---

## Part 4: The Latency Flow

Here is how a request actually flows through your system:



**The insight:** Every arrow represents measurable latency. It's the entire orchestration.

## Part 5: A Real-World Example

Let's make this concrete. Imagine a customer support chatbot that uses RAG to answer policy questions.

**Scenario:** User asks "What's your refund policy?"

### Before Optimization:

Step	Operation	Latency
1	Network + API Gateway	80ms
2	Queue Wait	150ms
3	Orchestrator Routing	120ms
4a	Vector Search (parallel)	250ms
4b	Memory Lookup (parallel)	180ms
4c	Query Embedding (parallel)	100ms
→	<b>Parallel Maximum</b>	<b>250ms</b>
5	Prompt Assembly	60ms
6	LLM Prompt Encoding	200ms
7	TTFT	600ms
8	Token Generation (200 tokens @ 30/sec)	6,700ms
9	Response Validation	150ms
10	JSON Formatting	30ms
11	Inter-agent Overhead	80ms
<b>TOTAL</b>		<b>~8.5 seconds</b>

**User experience:** Stares at loading spinner for 8.5 seconds. Probably thinks something broke.

## After Optimization:

Step	Operation	Latency	Improvement
1	Same-region deployment	50ms	✓ 30ms saved
2	Autoscaling enabled	20ms	✓ 130ms saved
3	Cached routing	80ms	✓ 40ms saved
4a	<b>Vector search cached</b>	5ms	✓ 245ms saved
4b	<b>Memory cached</b>	5ms	✓ 175ms saved
4c	<b>Embedding cached</b>	5ms	✓ 95ms saved
→	<b>Parallel Maximum</b>	<b>5ms</b>	✓ Huge win
5	Smaller context	40ms	✓ 20ms saved
6	Optimized encoding	150ms	✓ 50ms saved
7	Faster model (quantized)	400ms	✓ 200ms saved
8	<b>Streaming @ 50 tokens/sec</b>	4,000ms	✓ 2.7s saved
9	Async validation	50ms	✓ 100ms saved
10	Minimal formatting	20ms	✓ 10ms saved
11	Co-located services	30ms	✓ 50ms saved
<b>TOTAL</b>		<b>~4.8 seconds</b>	<b>43% faster</b>
<b>User Perceived</b>	<i>Time to first token</i>	<b>~725ms</b>	<b>Feels 11x faster</b>

**The magic of streaming:** Even though total generation takes ~5 seconds, users see results in under 1 second. This is transformative for user experience.

---

## Part 6: How to Actually Measure This

Theory is great. Implementation is better. Here is how to instrument your system.

### OpenTelemetry:

Use **OpenTelemetry (OTEL)** or a similar tool to add distributed tracing to every component. This gives you end-to-end visibility across your entire stack.

#### What OTEL gives you:

- Trace context that flows through every service
- Parent-child relationships between operations
- Precise timing for each component
- Ability to correlate requests across services
- Integration with Grafana, Datadog, Honeycomb, and other observability platforms

### Anatomy of a Good Trace Span

Each operation in your pipeline should create a span with:

**1. Component Name:** Clear, consistent naming: RetrieverAgent, VectorDB.search, LLM.generate

**2. Timestamps:** Start and end times with millisecond precision

**3. Duration:** Calculated automatically from timestamps

#### **4. Attributes (Tags)**

- request\_id: For correlating across services
- user\_id: For per-user analysis
- model\_name: Which model was used
- cache\_hit: True/false for cache analysis
- token\_count: Input and output tokens
- error: If something failed

**5. Parent-Child Relationships:** Shows which operations are nested within others, revealing the actual call flow

**Example Trace Structure**



This visualization shows you exactly where time is spent.

Metric	What It Tells You	Alert Threshold
TTFT P95	User experience quality	> 2 seconds
Pipeline P99	Worst-case scenarios	> 15 seconds
Cache Hit Rate	Optimization effectiveness	< 60%
Vector DB P90	Retrieval performance	> 500ms
Token Generation Rate	Model efficiency	< 20 tokens/sec
Error Rate	System reliability	> 1%
Timeout Rate	Resource adequacy	> 0.5%
Queue Depth	Load handling	> 100 requests

# Part 7: Optimization Strategies That Actually Work

Now for the good part, how to make your system faster.

## 1. Parallelize Everything You Can (if possible)

a) **The Problem:** Running agents sequentially multiplies latency.

```
# ❌ Bad: Sequential execution (530ms total)
retrieval_result = await retriever.execute() # 250ms
memory_result = await memory.execute() # 180ms
embed_result = await embedder.execute() # 100ms

# ✅ Good: Parallel execution (250ms total)
results = await asyncio.gather(
    retriever.execute(), # 250ms
    memory.execute(), # 180ms (runs simultaneously)
    embedder.execute() # 100ms (runs simultaneously)
)
# Total time = max(250, 180, 100) = 250ms
```

b) **Impact:** 40-70% reduction in coordination overhead.

## 2. Cache Aggressively (But intelligently)

a) **The Problem:** Repetitive expensive operations that don't change often.

b) **What to cache:**

- **Embeddings** (text → vector conversions): 90% hit rate possible
- **Vector search results** for common queries: 60-80% hit rate
- **LLM responses** for deterministic queries: 40-70% hit rate
- **Intermediate agent outputs:** Variable, but worth it

c) **Cache invalidation strategies:**

- **Time-based (TTL):** 5-60 minutes for dynamic content
- **Event-based:** Invalidate when underlying data changes
- **LRU (Least Recently Used):** For bounded cache sizes

d) **Real-world example:**

```
# ❌ Without cache: 250ms every time
results = await vector_db.search(query_embedding)

# ✅ With cache: 5ms on hit, 250ms on miss
cache_key = hash(query_embedding)
results = cache.get(cache_key)
if results is None:
    results = await vector_db.search(query_embedding)
```

e) **Impact:** 90-99% latency reduction on cache hits. At 70% hit rate, that's a 63-69% average latency reduction.

### 3. Stream Everything

a) **The Problem:** Users wait for complete generation before seeing anything.

b) **The Solution:** Stream tokens as they are generated.

```
# Batch: User waits 7 seconds
response = await llm.generate(prompt) # 7000ms
return response

# ✅ Streaming: User sees results in 600ms
async for token in llm.stream(prompt):
    yield token # First token appears at ~600ms (TTFT)
# Subsequent tokens appear every 30-50ms
```

c) **Impact:** Perceived latency reduced by 60-90%. Users feel the system is 6-10x faster, even though total generation time is similar.

d) **Pro tip:** Combine streaming with thinking indicators for even better UX:

```
"Searching knowledge base..." → 200ms
"Found 3 relevant articles..." → 400ms
[First token appears] → 600ms
[Streaming continues] → 7000ms total
```

### 4. Optimize Context Windows

a) **The Problem:** Long prompts increase encoding time and TTFT.

Every 1000 tokens of context adds:

- 50-150ms encoding time
- 100-300ms to TTFT
- Higher API costs

b) **Solutions:**

- **Summarize retrieved documents:**

```
# ❌ Include full documents (5000 tokens)
context = "\n\n".join([doc.full_text for doc in docs])

# ✅ Use summaries or relevant excerpts (500 tokens)
context = "\n\n".join([doc.summary or doc.text[:500] for doc in docs])
```

- **Implement sliding window for long conversations:**

```
# Keep last N messages + system prompt
recent_history = conversation_history[-10:] # Last 10 messages
```

- **Use prompt compression:**

```
# Libraries like LongLLMLingua can compress prompts by 50-70%
# while maintaining quality
compressed = compressor.compress(long_prompt)
```

c) **Impact:** 100-500ms saved per 1000 tokens reduced. Plus cost savings.

## 5. Batch Similar Requests

a) **The Problem:** Processing requests one at a time underutilizes GPU/compute.

b) **The Solution:** Batch similar requests together.

```
# Process multiple similar queries in one batch
batch = collect_requests(timeout=100ms, max_size=32)
results = await llm.generate_batch(batch)
```

c) **The Tradeoff:**

- **Throughput:** 3-10x improvement
- **Per-request latency:** 20–200ms additional queuing delay

d) **When to use batching:**

- Background processing jobs
- High-traffic systems with many similar queries
- Don't use in real-time user-facing applications where latency matters more than throughput

e) **Impact:** Massive cost reduction and throughput improvement. Use carefully for user-facing features.

## 6. Eliminate Cold Starts

a) **The Problem:** First request takes 500ms-5s longer due to initialization.

b) **Solutions:**

- **Keep model instances loaded:**

```
# Pre-load model at startup, not on first request
model = load_model("gpt-5") # Do this once at startup
```

```
# Not this:
```

```
def generate(prompt):
    model = load_model("gpt-5") # ✗ Adds 2-5s per request
    return model.generate(prompt)
```

- **Use warm pools in serverless:**

```
# AWS Lambda provisioned concurrency
ProvisionedConcurrency: 5 # Keep 5 instances always warm
```

- **Implement connection pooling:**

```
# Reuse database connections
db_pool = ConnectionPool(min_size=5, max_size=20)

# Instead of:
def query():
    conn = create_connection() # ✗ 50-200ms penalty
    result = conn.query(...)
    conn.close()
```

c) **Impact:** Removes 500ms-5s initialization delays. Especially critical for serverless deployments.

## 7. Co-locate Services

a) **The Problem:** Cross-region or cross-cloud network latency.

b) **Latency by location:**

- Same VPC/region: <5ms
- Cross-region (same cloud): 50-200ms
- Cross-cloud (AWS→GCP): 100-300ms
- On-prem → Cloud: 50-500ms

c) **Solution:** Deploy related services together.

d) **Architecture principles:**

- Keep vector DB and retrieval agent in same region
- Deploy orchestrator close to model endpoints
- Use regional CDN for global users
- Consider multi-region deployment for global scale

e) **Impact:** 50-200ms saved per inter-service call. In a system with 5-10 hops, that's 250ms-2s total.

## 8. Right-Size Your Models

a) **The Problem:** Using Large LLM for everything when smaller would work.

b) **Implementation:**

```
def route_to_model(task_complexity):
    if task_complexity == "simple":
        return SmallModel() # 3-10x faster
    elif task_complexity == "moderate":
        return MediumModel()
    else:
        return LargeModel()
```

c) **Impact:** 3-10x faster inference for routine operations. Massive cost savings.

## 9. Make Non-Critical Operations Async

a) **The Problem:** Blocking on logging, analytics, or optional validation.

```
# ❌ Bad: Blocks user response
response = generate_response(query)
await log_to_analytics(response) # 100ms
await update_metrics(response) # 80ms
await send_to_data_lake(response) # 200ms
return response # User waits extra 380ms

# ✅ Good: Fire and forget
response = generate_response(query)

# Run in background
asyncio.create_task(log_to_analytics(response))
asyncio.create_task(update_metrics(response))
asyncio.create_task(send_to_data_lake(response))
return response # No additional wait
```

b) **What to make async:**

- Logging to external systems
- Analytics and metrics collection
- Non-critical validation
- Audit trail updates
- Email/notification sending

c) **Impact:** 50-500ms saved per request. Users never notice these operations.

## 10. Monitor Tail Latencies Continuously

a) **The Problem:** Averages hide the worst user experiences.

b) **The Solution:** Track and alert on P95, P99, and P99.9.

```
# Set up alerts
if p99_latency > 10_000: # 10 seconds
    alert("P99 latency exceeds threshold")

if p95_latency > 5_000: # 5 seconds
    warning("P95 latency degrading")
```

c) **What tail latencies reveal:**

- Resource contention (P95 spikes → insufficient capacity)
- Outlier requests (P99 spikes → specific edge cases)
- System instability (P99.9 spikes → infrastructure issues)

d) **Impact:** Catches performance degradation before it affects most users. Prevents small issues from becoming major incidents.

---

## Part 8: What Matters Most

Here are the priorities:

### 1. For User Experience

a) **TTFT (Time to First Token):** Keep under 1 second

- This is what users feel as responsiveness
- More important than total latency

b) **Streaming:** Always stream when possible

- Reduces perceived latency by 60–90%
- Makes 5-second responses feel like 1-second responses

c) **Consistency:** Low jitter matters

- Users prefer consistently fast over occasionally blazing
- A system that is always 2s beats one that is 1s sometimes and 5s other times

### 2. For System Performance

a) **Pipeline P95 latency:** Keep under 5 seconds

- 95% of requests should complete in 5s or less
- This is your typically good target

b) **Cache hit rate:** Target 70%+ for common queries

- Single biggest latency reducer
- Monitor and optimize continuously

c) **Tail latencies (P99):** Keep under 10 seconds

- Your worst 1% shouldn't wait more than 10s
- This protects your most frustrated users

### 3. For Cost Efficiency

a) **Parallel execution:** Minimize sequential operations

- Cut coordination overhead by 40-70%
- Especially critical in multi-agent systems

b) **Caching strategy:** Avoid redundant expensive operations

- Can reduce costs by 60-90% on cache hits
- ROI is immediate

c) **Right-sized models:** Use smallest model that meets quality bar

- 3-10x cost savings
  - Often 3-10x latency improvement too
-

## Final Thoughts

In multi-agent GenAI systems, latency is about the entire orchestra, not just the soloist.

Your LLM might be fast, but if your orchestrator takes 500ms, your vector search takes 300ms, your validation takes 200ms, and your network hops add 150ms, you have added 1.15 seconds before the model even starts generating.

The difference between a prototype and a production system:

- **Measuring every component** with distributed tracing
- **Optimizing the critical path** through parallelization and caching
- **Prioritizing user-perceived latency** through streaming and TTFT reduction
- **Monitoring tail latencies** to catch problems before they become widespread
- **Making architecture decisions** based on data, not intuition

Low latency builds trust. When your agents respond quickly and consistently, users engage more deeply. When they don't, even the most intelligent system feels broken.

The key insight is making the entire pipeline efficient, measurable, and resilient under real-world conditions.

---