# Streaming

## Definition

Streaming in voice agents processes and delivers audio/text incrementally rather than waiting for complete results. Enables real-time interaction by starting playback before generation completes. Critical for maintaining conversational flow.

## Streaming ASR

Streaming speech recognition processes audio as it arrives, producing partial transcripts during speech. Provides 50-100ms latency versus 200-500ms for batch processing. Essential for responsive voice agents.

Partial results appear before user finishes speaking. Early words guide downstream processing. Final results include corrections and complete context. Applications must handle both partial and final transcripts.

Streaming ASR enables speculative processing. LLMs begin inference on partial transcripts, saving 200-500ms when predictions correct. Wrong predictions discard gracefully without impacting user experience.

## Streaming LLM

Language models stream tokens as generated rather than waiting for completion. Time-to-first-token determines perceived responsiveness. Under 500ms feels natural, over 1000ms feels broken.

Token-level streaming sends each word immediately. Minimizes latency but complicates TTS coordination. Synthesis must handle incomplete thoughts and potential rephrasing.

Sentence-level streaming buffers tokens until complete sentences form. Adds 200-500ms latency but produces coherent speech units. Balances responsiveness with speech quality. Most production systems use sentence-level streaming.

Chunk-based streaming groups tokens into semantic units (phrases, clauses). More flexible than full sentences, faster than waiting for periods. Requires natural language parsing to identify chunk boundaries.

## Streaming TTS

Text-to-speech streaming generates audio chunks as text arrives. First audio playback begins in 100-200ms. Maintains 1:1 generation ratio—one second of text produces one second of audio at real-time speed.

Streaming TTS requires careful buffer management. Audio playback cannot tolerate gaps. Systems maintain 150-250ms audio buffer to handle network jitter and generation variance.

Phoneme-level streaming generates audio in smallest units. Maximizes responsiveness but increases complexity. Word-level streaming balances simplicity with acceptable latency. Sentence-level streaming prioritizes quality over speed.

## Protocol Implementation

WebSocket provides full-duplex communication for streaming. Enables simultaneous send and receive. Lower overhead than HTTP polling. Standard for real-time web applications.

Server-sent events (SSE) offer simpler unidirectional streaming. Client receives updates, sends requests via HTTP. Easier to implement than WebSocket but less flexible. Suitable for streaming LLM responses to browser clients.

WebRTC streams audio with minimal latency using UDP transport. Handles packet loss gracefully through error correction. Industry standard for voice/video communication. More complex setup than WebSocket but superior audio quality.

gRPC streaming supports bidirectional communication with strong typing. Efficient binary protocol reduces bandwidth. Excellent for service-to-service streaming. Steeper learning curve than REST.

## Buffer Management

Client-side buffering prevents playback interruptions from network jitter. Too small (50ms) causes choppy audio. Too large (500ms+) increases latency. Adaptive buffering adjusts to conditions, typically maintaining 150-250ms.

Server-side buffering accumulates data before transmission. Batching reduces network overhead but increases latency. Target 50-100ms server buffers for voice applications. Larger buffers acceptable for non-real-time use cases.

Jitter buffers compensate for variable network delays. Reorder packets arriving out-of-sequence. Discard severely delayed packets rather than disrupting playback. Essential for maintaining smooth audio.

## Error Handling

Connection failures require graceful degradation. Buffer partial results to resume after reconnection. Notify users of connectivity issues. Implement exponential backoff for reconnection attempts.

Partial result handling manages incomplete data. Mark streaming results as provisional. Update with corrections when final results arrive. Applications must reconcile partial and final versions.

Timeout management prevents infinite waits. Set progressive timeouts—short for first response, longer for continuation. If streaming stalls, fall back to batch processing or cached responses.

## Backpressure

Backpressure occurs when producers generate data faster than consumers process. Audio playback has fixed speed—streaming faster than real-time creates unbounded buffering. Implement flow control to match production to consumption.

Token-bucket algorithms limit streaming rate. Producers wait when buffer fills. Prevents memory exhaustion from runaway generation. Maintains system stability under load.

Reactive streams provide standardized backpressure handling. Consumers signal readiness for more data. Producers respect consumption capacity. Prevents overwhelming downstream systems.

## Optimization Techniques

Predictive buffering anticipates user actions. Pre-buffer common responses during silence. Enables instant playback when predictions match. Wastes bandwidth when wrong but improves perceived responsiveness.

Parallel streaming processes multiple pipeline stages simultaneously. Begin TTS on early LLM tokens while generation continues. Start audio playback while synthesis completes. Reduces end-to-end latency 30-50%.

Adaptive quality adjusts streaming parameters based on conditions. Reduce audio quality on poor connections to maintain real-time performance. Increase quality on strong connections. Balances quality with reliability.

## Architecture Patterns

Pipeline architecture chains streaming stages sequentially. ASR streams to LLM, LLM streams to TTS, TTS streams to client. Simple to reason about but introduces cumulative latency. Each stage adds delay.

Fan-out architecture broadcasts streaming data to multiple consumers. Single ASR stream feeds multiple LLM instances. Enables A/B testing, parallel processing, and redundancy. Increases resource usage but improves reliability.

Hybrid architecture combines streaming and batch processing. Stream partial results for low latency. Process complete batches for accuracy. Reconcile differences asynchronously. Provides best of both approaches.

## Testing Strategies

Synthetic streams test pipeline behavior without external dependencies. Generate mock audio/text at controlled rates. Verify buffer management, error handling, and backpressure under various conditions.

Network simulation introduces latency, jitter, and packet loss. Tools like tc (traffic control) simulate poor conditions. Ensures graceful degradation. Identifies breaking points before production deployment.

Load testing validates streaming scalability. Simulate hundreds or thousands of concurrent streams. Monitor CPU, memory, network bandwidth. Identify bottlenecks and capacity limits.

## Monitoring and Metrics

Stream health metrics track connection stability. Monitor connection duration, disconnection rate, reconnection frequency. High disconnection rates indicate network issues or implementation problems.

Latency distribution reveals streaming performance. Track p50, p90, p99 for first-token time and chunk delivery. Percentiles show typical experience and worst-case scenarios. Target p99 under 1000ms for voice agents.

Buffer utilization indicates capacity headroom. Monitor buffer fill levels across pipeline stages. Consistently full buffers signal backpressure issues. Empty buffers suggest underutilization.

## Best Practices

Design for failure from the start. Implement timeouts, retries, and fallbacks. Test failure modes explicitly. Graceful degradation preserves user experience during issues.

Use appropriate protocols for use case. WebRTC for voice, WebSocket for text, gRPC for services. Don't force unsuitable protocols. Each has strengths and trade-offs.

Minimize serialization overhead. Binary formats (protobuf, msgpack) outperform JSON for high-frequency streaming. Profile hot paths. Optimize critical sections for throughput.

Implement circuit breakers for external dependencies. Prevent cascade failures when downstream services struggle. Fast-fail unhealthy connections. Preserve system stability.