

□ Streaming (Real-time Token Streaming) – 100 Questions & Answers

Concepts & Basics

1.

Q: What is streaming in AI models?

A: Streaming is the process of sending model-generated tokens incrementally in real time instead of waiting for the full response.

2.

Q: Why is streaming useful in conversational AI?

A: It provides faster perceived response time, making interactions feel more natural.

3.

Q: How does token streaming differ from batch responses?

A: Streaming sends tokens one-by-one as generated, while batch responses send the entire output at once.

4.

Q: What is a “token” in streaming?

A: A token is the smallest unit of text (word or sub-word) that the model outputs sequentially.

5.

Q: Which part of an LLM enables streaming?

A: The autoregressive decoding mechanism that generates one token at a time.

6.

Q: Can streaming reduce latency?

A: Yes, because tokens are delivered progressively instead of waiting for the entire generation.

7.

Q: Is streaming synchronous or asynchronous?

A: Typically asynchronous, using callbacks or event-based handlers.

8.

Q: What is the main advantage of streaming in user experience?

A: Users can start reading the response immediately, improving interactivity.

9.

Q: What is “partial completion” in streaming?

A: Receiving and displaying tokens before the entire response is finished.

10.

Q: What is “real-time feedback” in streaming?

A: Immediate delivery of generated tokens, enabling dynamic updates for the user.

Implementation

11.

Q: How is streaming implemented in Python APIs?

A: By enabling a streaming flag and attaching a callback/event handler for tokens.

12.

Q: Which Python libraries support AI streaming?

A: OpenAI SDK, LangChain, and Async frameworks (like asyncio).

13.

Q: What is a stream_handler?

A: A function or object that processes tokens as they arrive.

14.

Q: Which method handles streamed output in OpenAI Python SDK?

A: `stream = client.chat.completions.create(..., stream=True)`

15.

Q: Can multiple callbacks be used in streaming?

A: Yes, for logging, UI updates, and custom processing simultaneously.

16.

Q: What is a buffer in streaming?

A: A temporary storage for received tokens before final assembly.

17.

Q: How do you display streaming text in a CLI app?

A: Print tokens as they arrive instead of storing until completion.

18.

Q: How do you handle streaming in a web app?

A: Use WebSockets or Server-Sent Events (SSE) for token updates.

19.

Q: What Python feature is most used for streaming?

A: Asynchronous generators (async for).

20.

Q: Why is yield useful in streaming implementations?

A: It allows progressive token output in generator functions.

Technical Behavior

21.

Q: Does streaming affect model quality?

A: No, only delivery changes, not model output.

22.

Q: Can streaming be paused?

A: Yes, by halting the event loop or stream consumer.

23.

Q: What happens if a stream is interrupted?

A: Partial results are available, but the full response may be incomplete.

24.

Q: Can you restart a stream from the middle?

A: No, you must re-request completion from the model.

25.

Q: Is streaming deterministic?

A: The generation process is the same; streaming just reveals it earlier.

26.

Q: Does streaming consume fewer tokens?

A: No, token usage remains the same.

27.

Q: Can streaming handle multiple responses in parallel?

A: Yes, with async workers or threads.

28.

Q: What is the streaming format of OpenAI completions?

A: A series of delta updates containing new tokens.

29.

Q: How do you detect the end of a stream?

A: Via a special event/flag like `finish_reason`.

30.

Q: What is backpressure in streaming?

A: The slowdown when consumer processing lags behind token generation.

Performance & Optimization

31.

Q: Does streaming reduce total latency?

A: Not overall, but improves perceived latency.

32.

Q: How can you optimize token streaming?

A: Use efficient callbacks and avoid heavy processing per token.

33.

Q: Can streaming overload the client?

A: Yes, if tokens are produced faster than consumed.

34.

Q: What is batching tokens in streaming?

A: Sending tokens in small groups instead of one-by-one.

35.

Q: Why would you batch tokens instead of single?

A: To reduce overhead and improve rendering speed.

36.

Q: Does streaming affect memory usage?

A: Slightly lower, since text is processed progressively.

37.

Q: Can you measure token speed in streaming?

A: Yes, by calculating tokens per second (TPS).

38.

Q: Why is async I/O important in streaming?

A: It prevents blocking while waiting for new tokens.

39.

Q: What role does concurrency play in streaming?

A: Enables handling multiple streams or users simultaneously.

40.

Q: How do you ensure smooth UI updates in streaming?

A: Buffer small chunks and throttle rendering.

Use Cases

41.

Q: Why is streaming good for chatbots?

A: It mimics real human typing.

42.

Q: How is streaming used in transcription apps?

A: For real-time speech-to-text updates.

43.

Q: How does streaming improve coding assistants?

A: Developers see code generated progressively.

44.

Q: What role does streaming play in live demos?

A: It creates immediate feedback for audiences.

45.

Q: Why is streaming critical in low attention apps?

A: Keeps users engaged with partial output.

46.

Q: Is streaming used in real-time translation?

A: Yes, to output translated text as input arrives.

47.

Q: What industry benefits most from streaming?

A: Customer support systems for quick response feel.

48.

Q: How can streaming aid accessibility?

A: Screen readers can read tokens as they appear.

49.

Q: How does streaming affect decision-making apps?

A: Provides early hints before full results are ready.

50.

Q: Why is streaming popular in collaborative editing tools?

A: It enables live co-creation.

Error Handling

51.

Q: What happens if connection drops mid-stream?

A: Tokens already received are safe, but rest are lost.

52.

Q: How do you retry a failed stream?

A: Restart the request with the same input.

53.

Q: What is a timeout in streaming?

A: A limit on waiting time for new tokens.

54.

Q: Can tokens arrive out of order?

A: No, they always arrive sequentially.

55.

Q: How do you detect streaming errors?

A: Via exception handling or error callbacks.

56.

Q: What is graceful termination in streaming?

A: Closing the stream cleanly without data corruption.

57.

Q: How do you debug a streaming issue?

A: Log each token and check network stability.

58.

Q: Can partial results be finalized after error?

A: Yes, but they may be incomplete.

59.

Q: What is a fallback strategy if streaming fails?

A: Request a full non-streamed completion.

60.

Q: How do you handle network congestion in streaming?

A: Use buffering and throttling mechanisms.

Advanced Concepts

61.

Q: Can you stream structured outputs?

A: Yes, but partial JSON must be carefully parsed.

62.

Q: How do you validate partial JSON in streaming?

A: Use incremental parsers.

63.

Q: What is speculative decoding in streaming?

A: Predicting upcoming tokens to reduce latency.

64.

Q: Can you stream images or multimodal outputs?

A: Currently limited, but research is ongoing.

65.

Q: What is token-level delay?

A: The time gap between consecutive tokens.

66.

Q: Can streaming be combined with speech synthesis?

A: Yes, to speak tokens as they arrive.

67.

Q: How does streaming integrate with tracing?

A: Logs token-by-token generation.

68.

Q: Can multiple agents stream collaboratively?

A: Yes, each agent can stream its part.

69.

Q: What is token prefetching?

A: Generating tokens ahead of need to reduce wait.

70.

Q: What is streaming checkpointing?
A: Saving partial progress for recovery.

Comparison

71.

Q: Streaming vs Batch: Which is faster overall?
A: Same total time, but streaming feels faster.

72.

Q: Which is easier for error handling: streaming or batch?
A: Batch, because the result is complete.

73.

Q: Which is better for interactivity?
A: Streaming.

74.

Q: Which consumes more network requests?
A: Streaming, due to multiple small chunks.

75.

Q: Which uses more CPU: streaming or batch?
A: Streaming may use more due to frequent callbacks.

76.

Q: Which is preferred for offline generation?
A: Batch.

77.

Q: Which provides better user experience?
A: Streaming.

78.

Q: Which mode is more reliable?

A: Batch, since it avoids mid-stream failures.

79.

Q: Which mode allows immediate partial use?

A: Streaming.

80.

Q: Which works better with low bandwidth?

A: Batch, to reduce packet overhead.

Practical

81.

Q: Can you stop a stream midway?

A: Yes, by aborting the request.

82.

Q: How do you collect the final response from stream?

A: Concatenate tokens progressively.

83.

Q: How do you detect user interruptions in streaming?

A: Listen for cancel/abort signals.

84.

Q: How do you measure latency in streaming?

A: Track first-token time and last-token time.

85.

Q: Can you store streamed tokens in a DB?

A: Yes, append as they arrive.

86.

Q: Can streaming be logged for debugging?

A: Yes, record tokens with timestamps.

87.

Q: How do you avoid flicker in UI with streaming?

A: Use buffered display updates.

88.

Q: Can you merge multiple streams into one?

A: Yes, with multiplexing.

89.

Q: How do you secure streaming?

A: Use encrypted channels (HTTPS/WSS).

90.

Q: Can you replay a stream later?

A: Yes, if tokens were logged.

Edge Cases

91.

Q: What if the first token is delayed?

A: The user perceives lag; prefetching can help.

92.

Q: What if the model generates nothing?

A: The stream closes with empty output.

93.

Q: What happens if streaming is too fast for UI?

A: UI may lag unless throttled.

94.

Q: What if tokens contain formatting symbols?

A: Render progressively, then adjust formatting at end.

95.

Q: Can streaming output overlap?

A: No, tokens are sequential.

96.

Q: Can the model change earlier tokens during streaming?

A: No, once emitted, tokens are final.

97.

Q: What happens if multiple users stream simultaneously?

A: Each gets their own channel/session.

98.

Q: What if stream output is cut mid-word?

A: Word completes in following tokens.

99.

Q: Can you combine streaming with retries?

A: Yes, restart from last checkpoint.

100.

Q: What is the ultimate goal of streaming in AI?

A: To create responsive, natural, and interactive user experiences.

