



Query Input

Query & Question

Q: "What did the man holding the microphone do after he saw the hand gesture of the lady with cap?"

Options

- A. Shakes her hand
- B. Signal to the rest to stop
- C. Smile
- D. Walk away
- E. Stand at the side

↓ Input ↓

1. Decomposer

For question "What did the man..lady with cap?", it is determined to be a **single-intent query**. The query content does not need to be decomposed and is "What did the man..lady with cap?"



Video Input



↓ Input ↓



2. Multi-source Agent



Temporal Agent

Recorded in order...the lady made a gesture...the man responded...the band stopped; the memory bank...action was a key link to...

Answer: (B)



Web-based Agent

When searching ... "gesture interaction" ... accompanied by friendly responses, and "smiling" is a common...

Answer: (C)



Video-language Agent

Located time segments 00:18-00:20 (lady's gesture)... 00:23-00:24 ... the man... conforming to the function of "signaling to stop"...

Answer: (B)



4. Final Answer

Answer

- A. Shakes her hand
- B. Signal to the rest to stop
- C. Smile
- D. Walk away
- E. Stand at the side

Reason

lady gestured, man acted, band stopped... chain shows man signaled others to stop.

↑ Output ↑



3. Decision Agent

Synthesizing answers... consistency above threshold... two point to the same result. Web-Based... deviant, but overall consistency meets requirements. Polished answer....