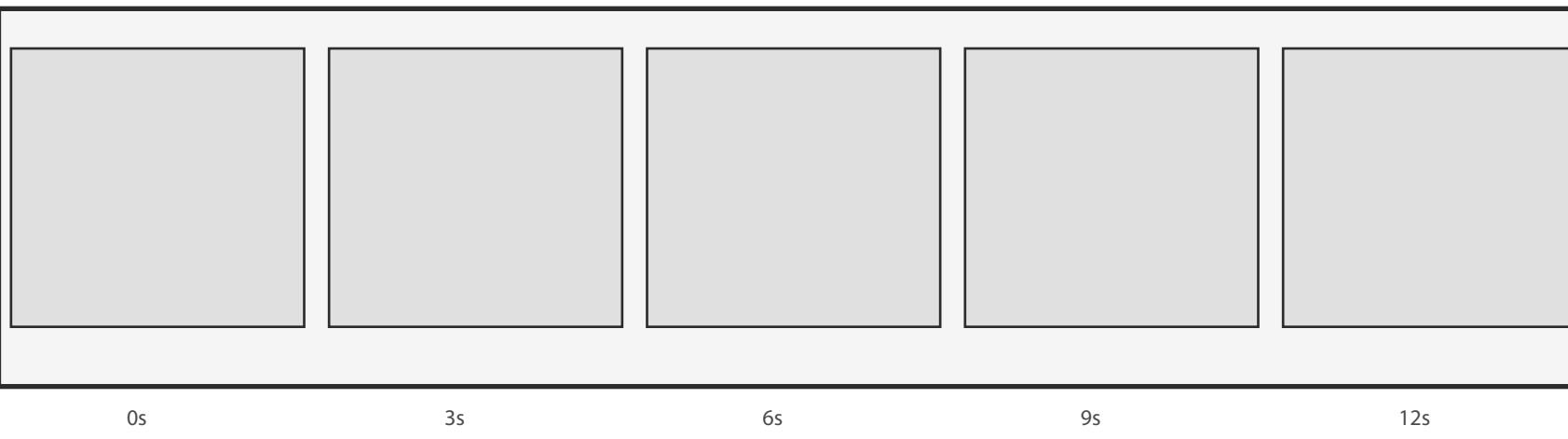


HERBench: High Evidence Requirement Benchmark

(1) Referring & Tracking

[AGAR] Appearance-Grounded Attribute Recognition

How does the person's pace appear throughout their movement in the video?

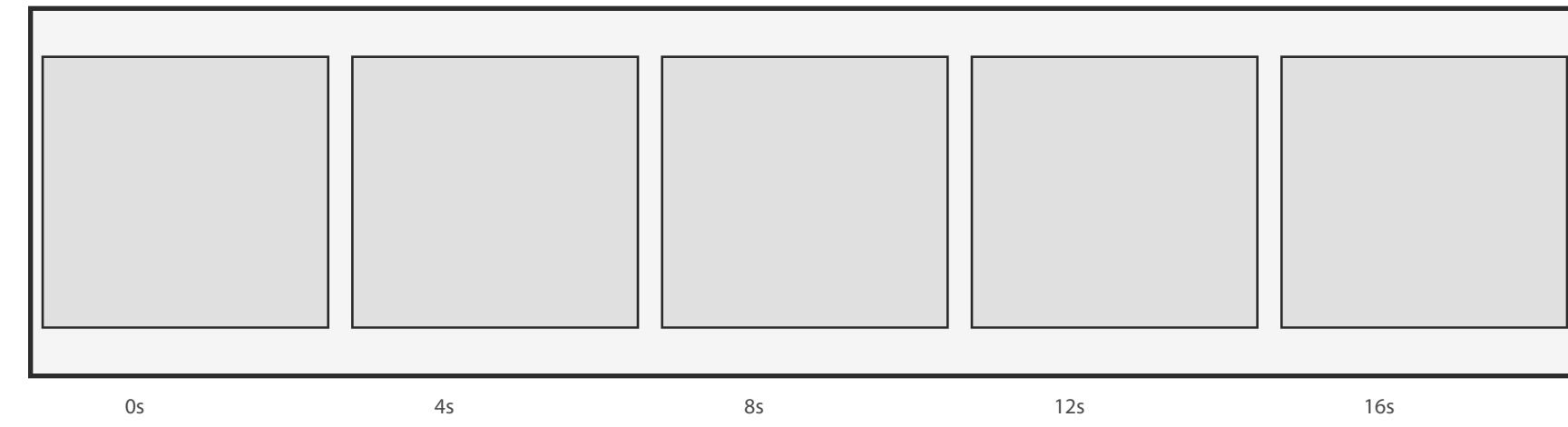


0s 3s 6s 9s 12s

Answer: They maintain a consistent pace without stopping

[AGBI] Appearance-Grounded Behavior Interactions

Who is accompanying the person as they move across the middle horizontal section?

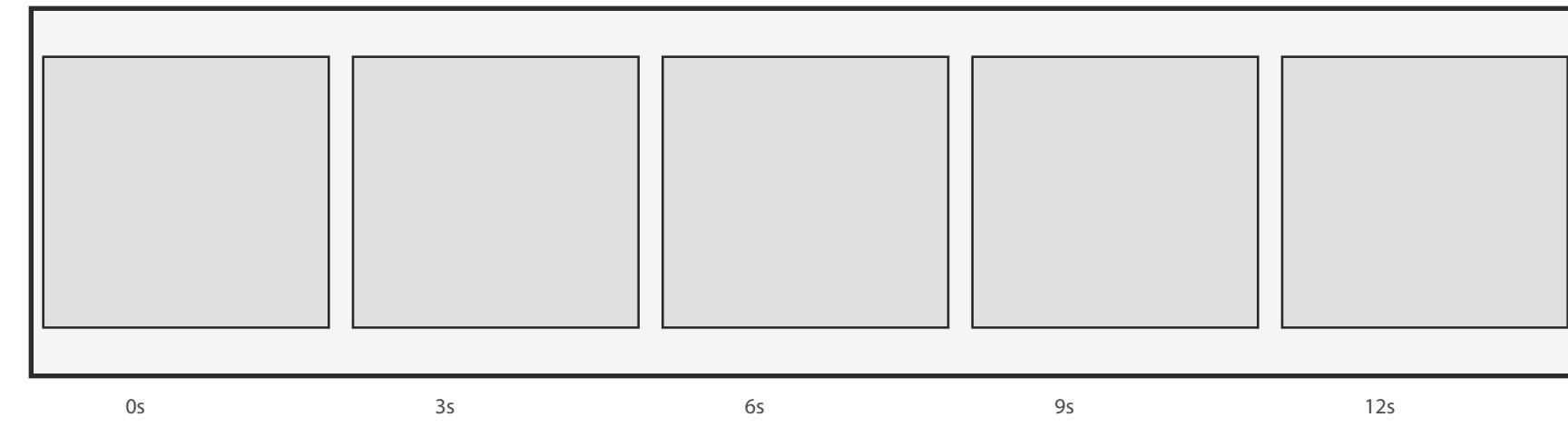


0s 4s 8s 12s 16s

Answer: A female wearing a light-colored coat and dark boots

[AGLT] Appearance-Grounded Localization Trajectory

In the video, how does the person move from the left to the right edge?



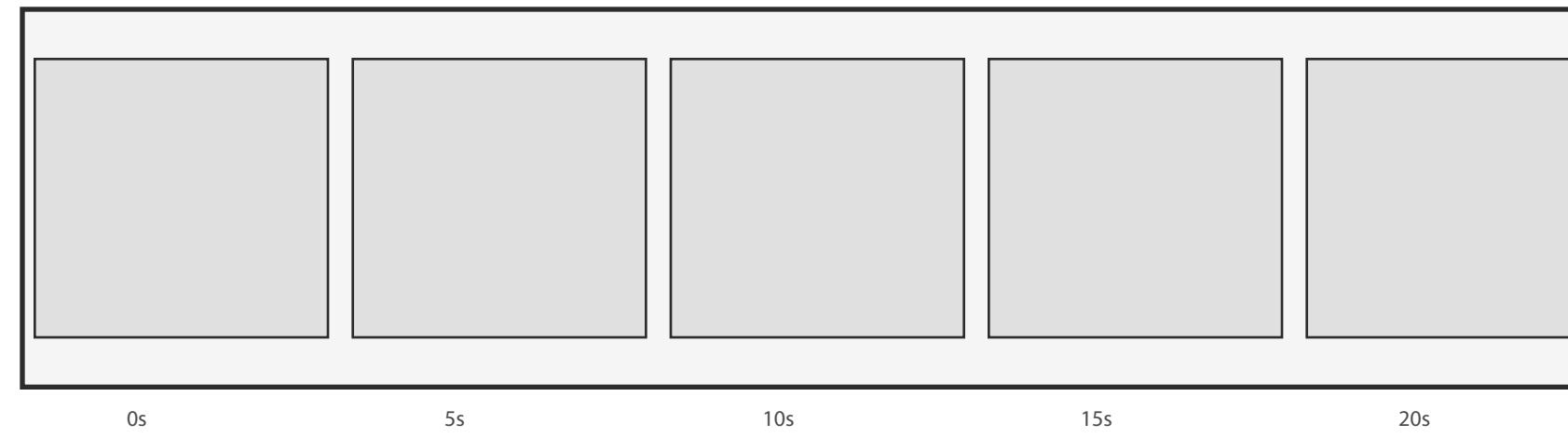
0s 3s 6s 9s 12s

Answer: Mostly straight without noticeable stops

(2) Multi-Entity Aggregation & Numeracy

[MEGL] Multi-Entities Grounding & Localization

Which people appeared in the video? (descriptions must match exactly)

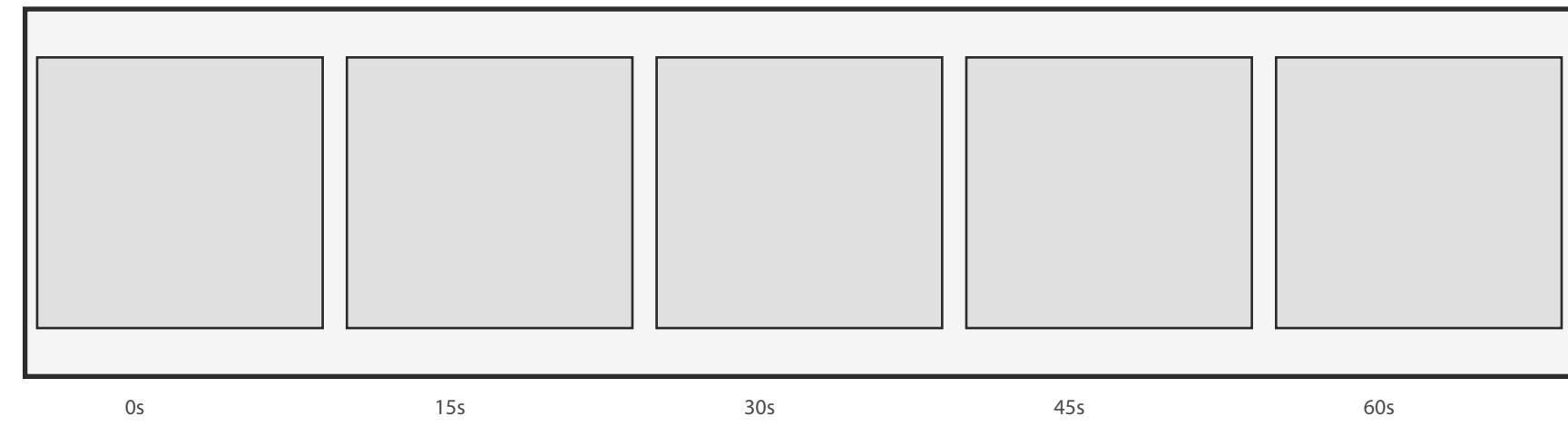


0s 5s 10s 15s 20s

Answer: All of them

[RLPC] Region-Localized People Counting

How many people entered the frame through the top edge?

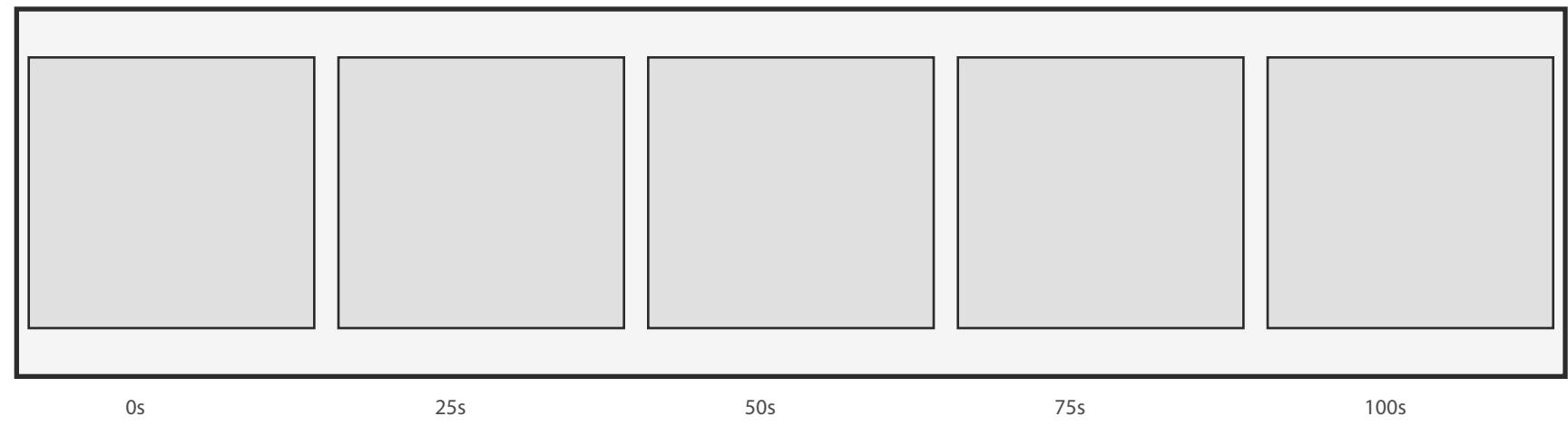


0s 15s 30s 45s 60s

Answer: 219-298

[AC] Action Counting

How many times does the action-object pair 'stir potatoes' occur?



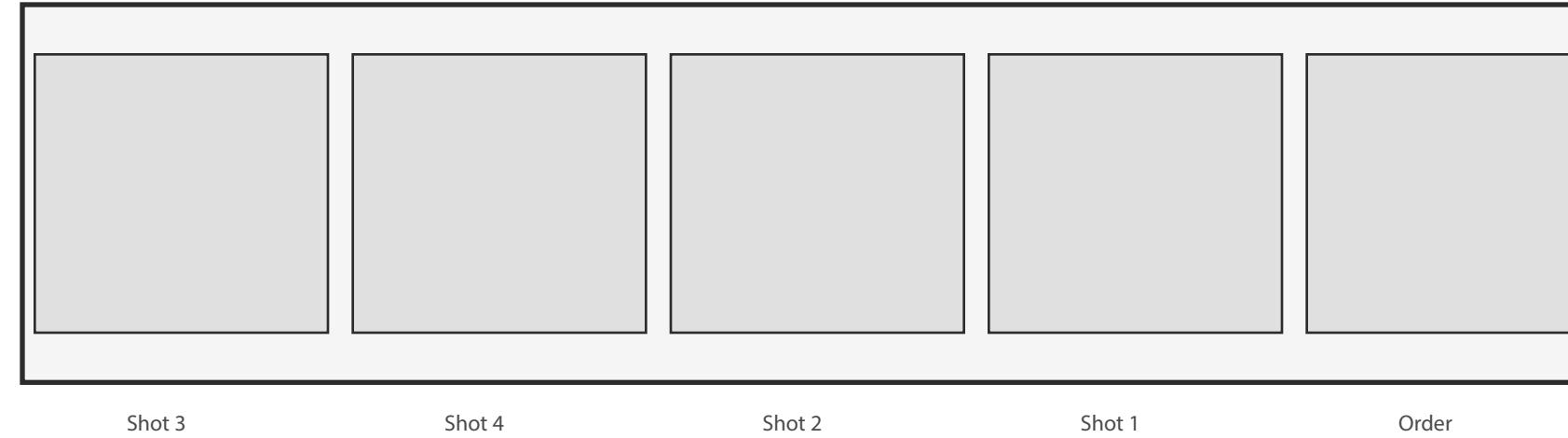
0s 25s 50s 75s 100s

Answer: 31

(3) Temporal Reasoning & Chronology

[TSO] Temporal Shot Ordering

What is the correct chronological order of these 4 shots?

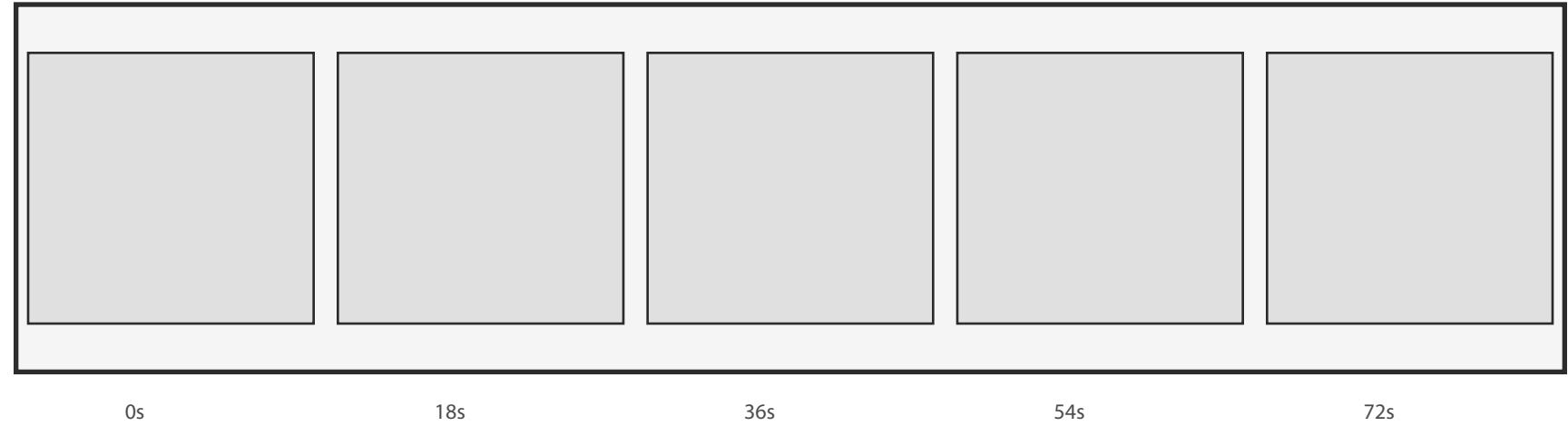


Shot 3 Shot 4 Shot 2 Shot 1 Order

Answer: 3>4>2>1

[ASII] Action Sequence Integrity & Identification

What is the correct temporal order of the 5 narrated events?

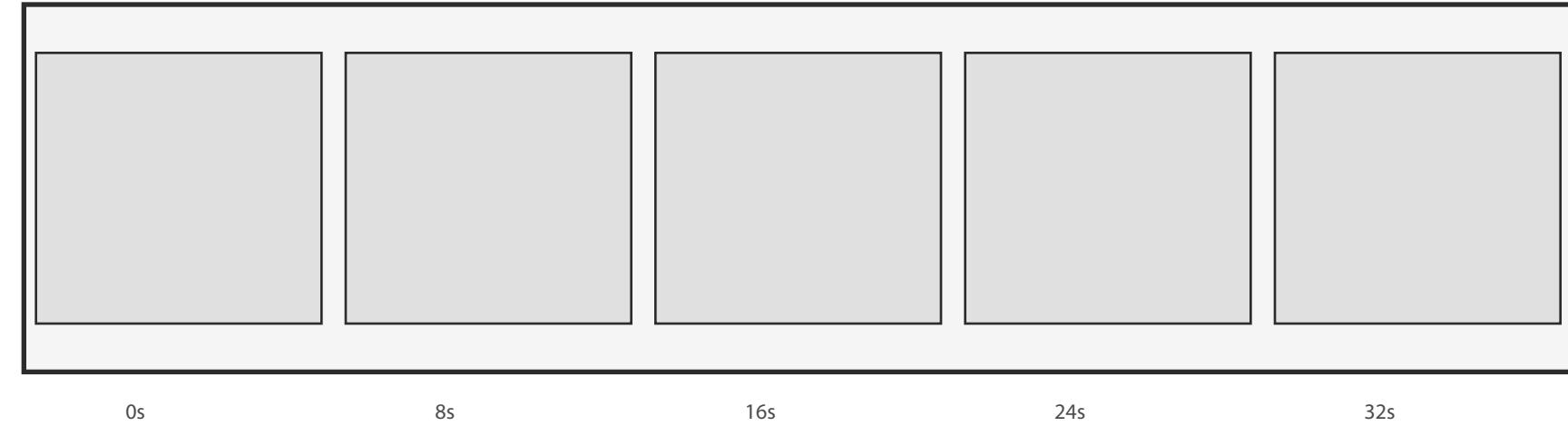


0s 18s 36s 54s 72s

Answer: Close door -> Spoon rice -> Pick kettle -> Pick sriracha -> Pick soy sauce

[MPDR] Multi-Person Duration Reasoning

Who entered and exited the frame first?



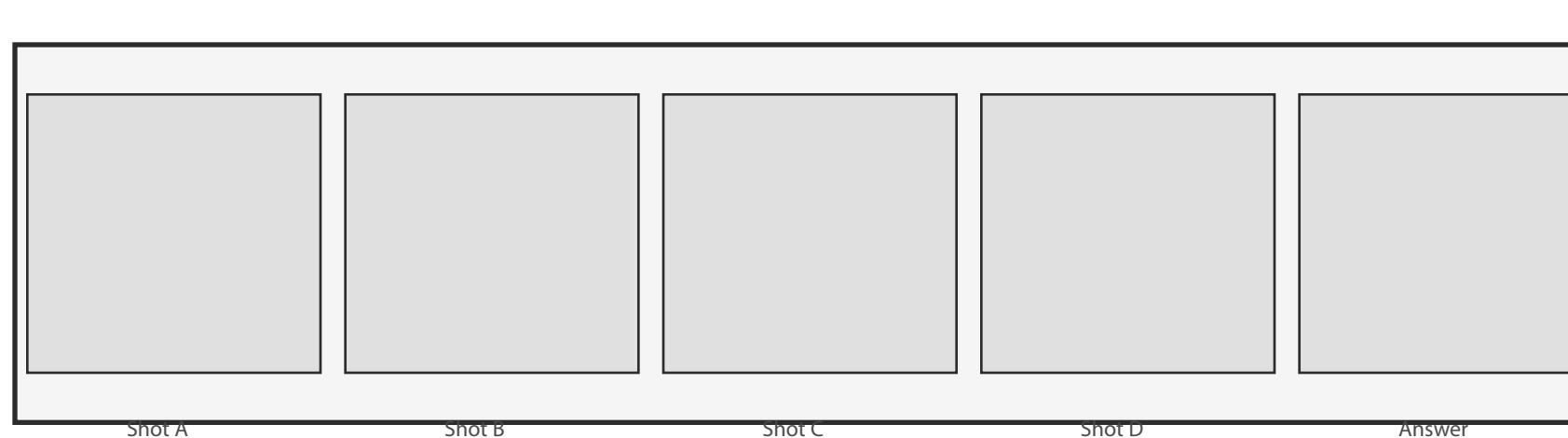
0s 8s 16s 24s 32s

Answer: Person 2 entered first and exited first

(4) Global Consistency & Verification

[SVA] Scene Verification Arrangement

From the correctly described shots, which appears last in the video?

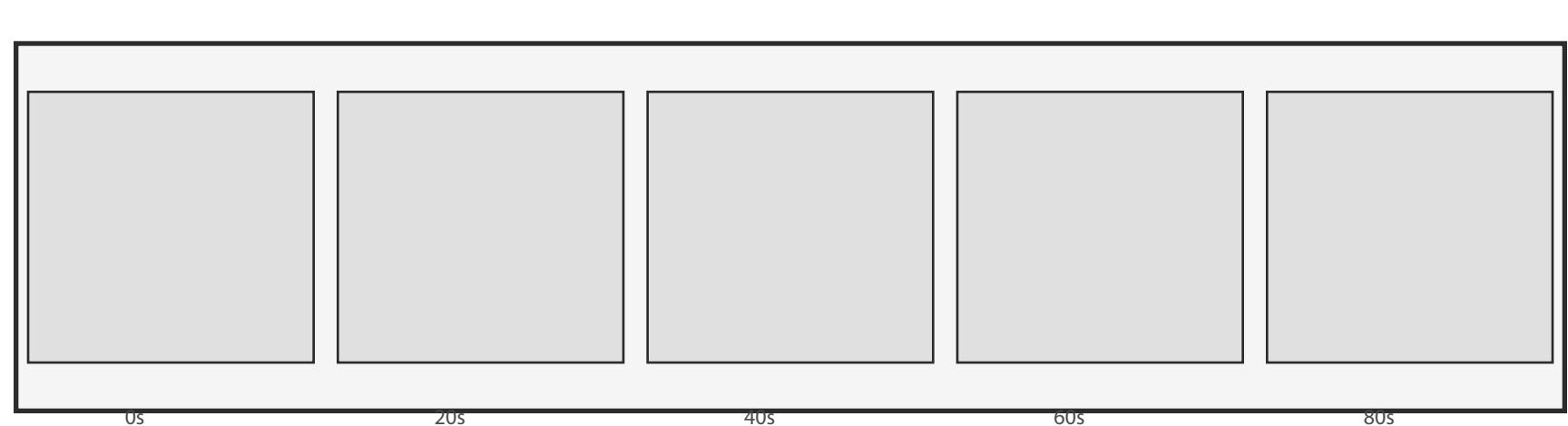


Shot A Shot B Shot C Shot D Answer

Answer: Person with dark hair at window overlooking stormy sea

[FAM] False Action Memory

Which of the following actions did NOT occur in the video?

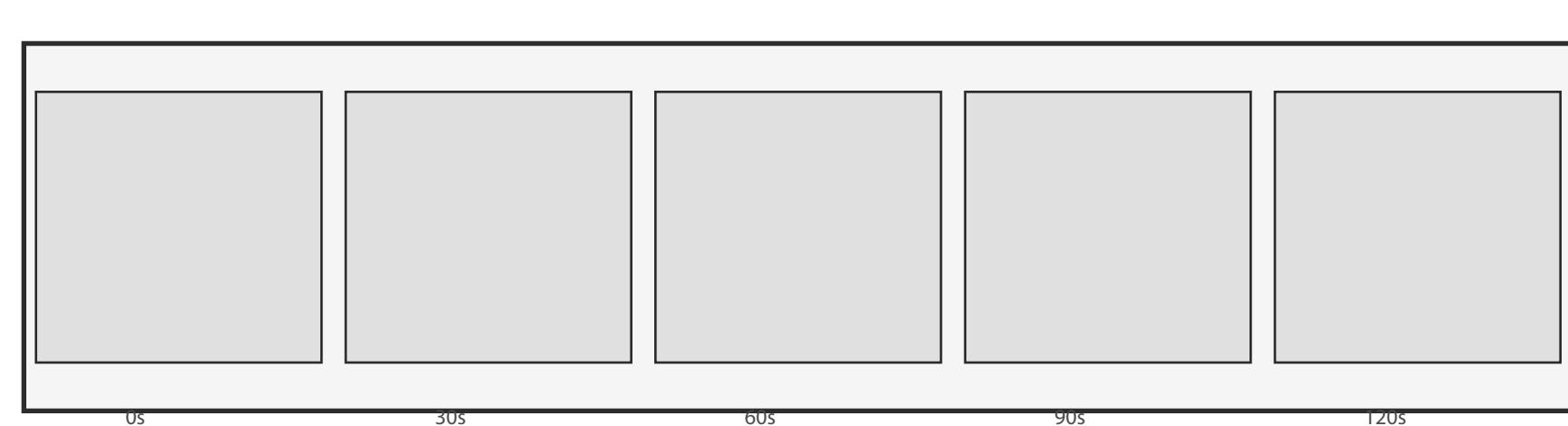


0s 20s 40s 60s 80s

Answer: Stir chopped onion

[FOM] False Object Memory

Which object did the camera wearer NOT interact with?



0s 30s 60s 90s 120s

Answer: Sweet potato