

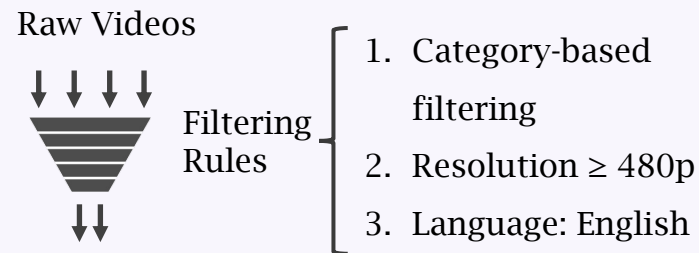
Automatic Processing

Step 1. Conversational Video Collection



14 K Videos Collected

Step 2. Basic Filtering

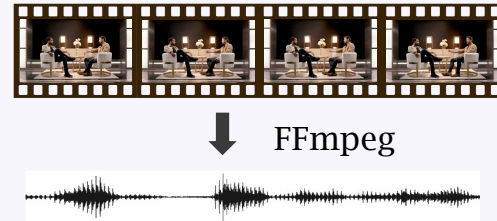


Step 3. Scene Cut

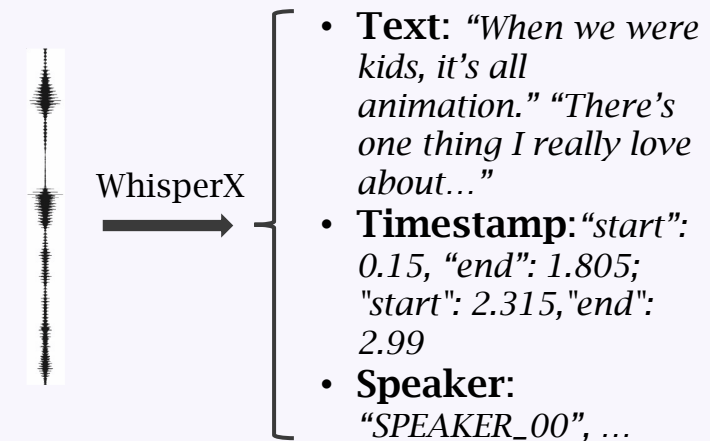


Basic Video Processing

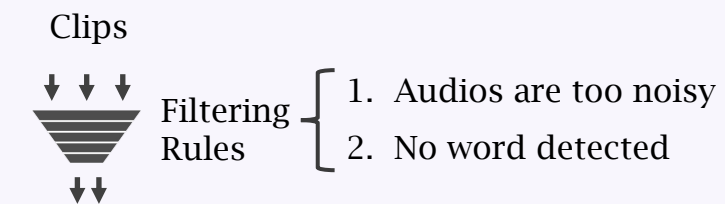
Step 4. Audio Extraction



Step 5. Speech Recognition & Speaker Diarization

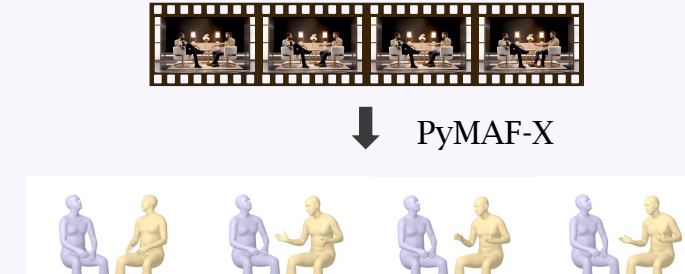


Step 6. Audio Filtering

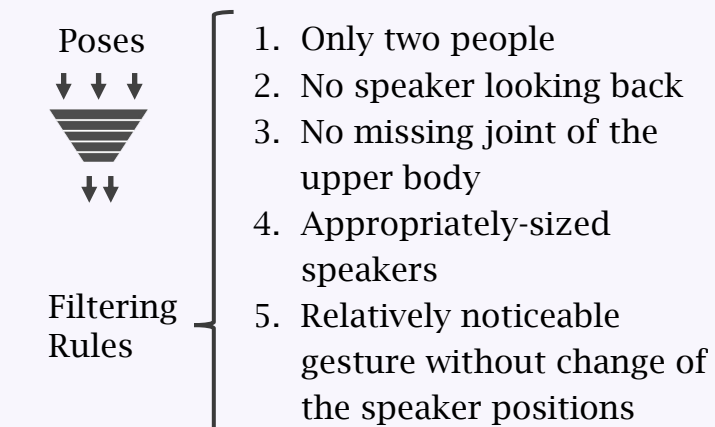


Audio Extraction and Filtering

Step 7. Pose Estimation



Step 8. Pose Filtering



Step 9. Smoothing



Pose Estimation and Filtering

Manual Processing

Step 10. Video Filtering



✗ Non-conversational

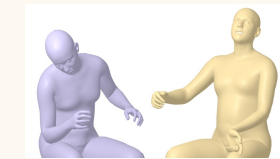


✗ Dual-person shots with external voice

Step 11. Pose Filtering



✗ Unnatural



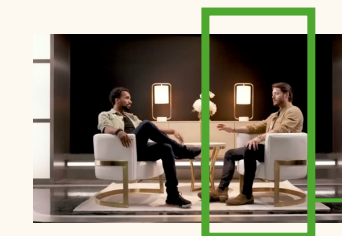
✗ Unnatural

✗ Jittery...

Step 12. Audio&Speaker Alignment

Text	Timestamps	Speaker
"When we were kids, it's all animation."	"start": 0.15, "end": 1.805	"SPEAKER_00"
"There's one thing I really love about..."	"start": 2.315, "end": 2.99	"SPEAKER_01"

👉 1. Manual Correction
"SPEAKER_00"



👉 2. Manual Alignment
"SPEAKER_00" : right