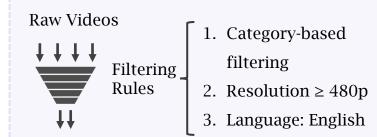
Automatic Processing

Step 1. Conversational Video Collection



14 K Videos Collected

Step 2. Basic Filtering

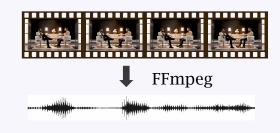


Step 3. Scene Cut

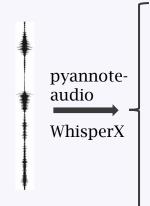


Basic Video Processing

Step 4. Audio Extraction

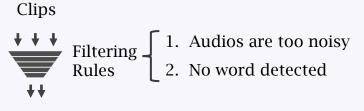


Step 5. Speech Recognition & Speaker Diarization



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

Step 6. Audio Filtering



Audio Extraction and Filtering

Step 7. Pose Estimation



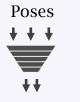








Step 8. Pose Filtering

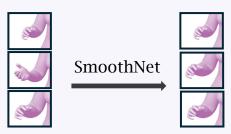


Filtering

Rules

- 1. Only two people
- 2. No speaker looking back
- 3. No missing joint of the upper body
- 4. Appropriately-sized speakers
- 5. Relatively noticeable gesture without change of the speaker positions

Step 9. Smoothing



Pose Estimation and Filtering

Manual Processing

Step 10. Video Filtering



X Nonconversational



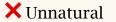
X Dual-person shots with external voice

Step 11. Pose Filtering





X Jittery...



X Unnatural

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		
	"SPEAK -	ER_00"
2. Manual Alignment		
	→ "SPEAKE.	R_00" : right ←