

A Multimodal Simultaneous Interpretation Prototype: Who Said What

Xiaolin Wang, Masao Utiyama and Eiichiro Sumita
{xiaolin.wang, mutiyama, eiichiro.sumita}@nict.go.jp

Advanced Translation Research and Development Promotion Center
National Institute of Information and Communications Technology, Japan

Conventional Way of Simultaneously Interpreting Video Streams

- Subtitle ← Translation



[Video 1](#)

So at some point I still hate my life so far, but I really have to decide to do what I want to do if I really have only one time.

Readable Way of Simultaneously Interpreting Video Streams

- Chat Log ← Speaker + Translation (Who Said What)



So at some point I still hate my life so far, but I really have to decide to do what I want to do if I really have only one time.

ago.

I thought I'd come and do what I liked, but after a while I thought I couldn't do it, so I gave up.

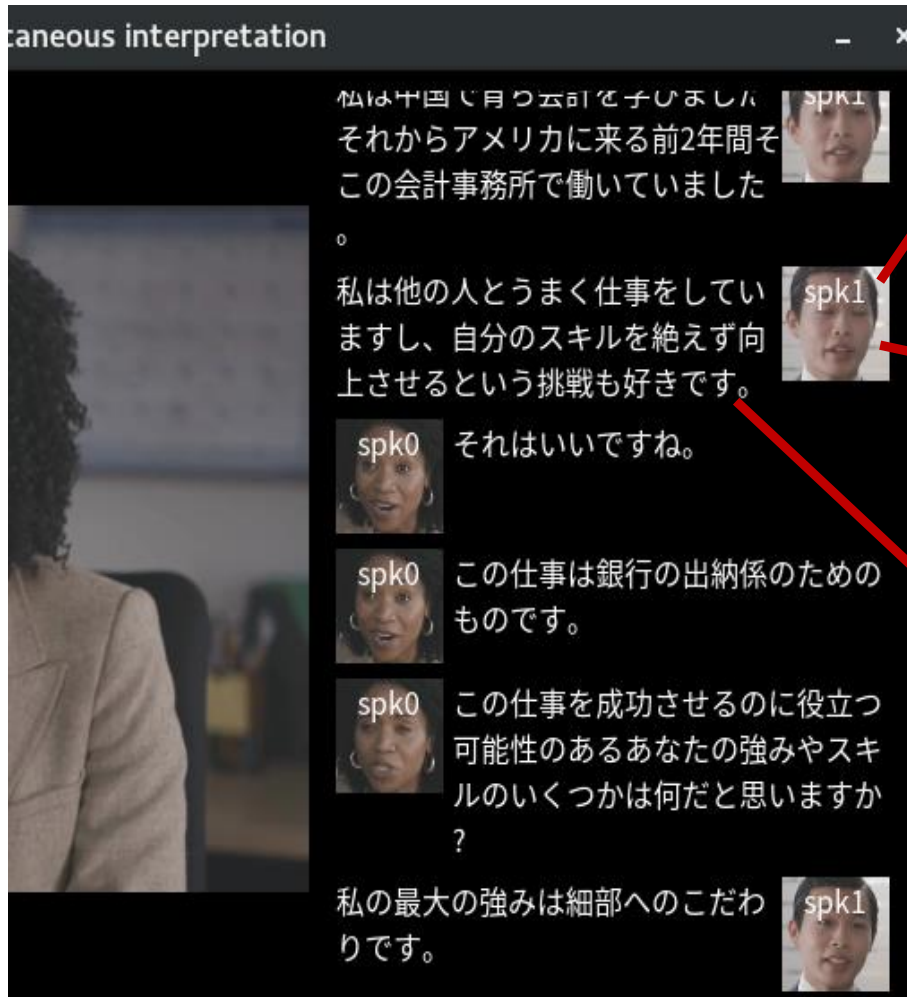
But after three years, I still want to do something in the world of cooking, so I quit my job and now I'm working at a restaurant, and there are quite a few people who want to open a restaurant after five years.

So at some point I still hate my life so far, but I really have to decide to do what I want to do if I really have only one time.



[Video 2](#)

User Interface : Chat Log



- **Speaker Tag**
 - Plain text of “*spk n*”
 - facilitate post-editing
- **Face Icon**
 - *Speakers' sentiments*
 - *Double check speaker tag*
- **Translation**
 - *Content of speech*

Post-Editing with Speaker Tags

- Plain transcript

wonderful
and what would you say are some of your weaknesses
one of my biggest weaknesses is asking for help when I need it
I 'd like to do better at that
I appreciate your honesty mister wang
what can you tell me about some of your goals over the next few years
my primary goal is to gain more work experience
so a position like this would help me meet that goal

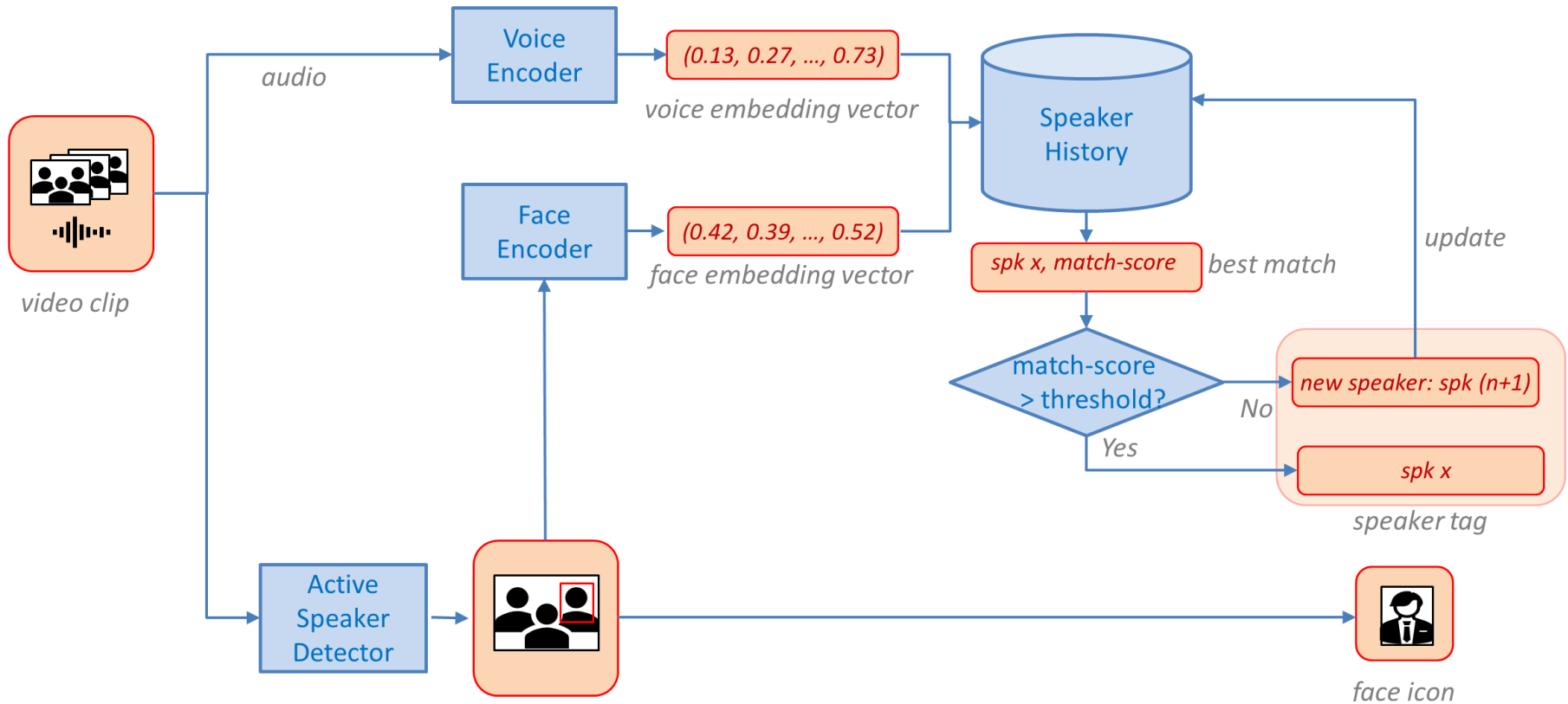
- Annotated transcript

<i>spk0</i>	wonderful
<i>spk0</i>	and what would you say are some of your weaknesses
<i>spk1</i>	one of my biggest weaknesses is asking for help when I need it
<i>spk1</i>	I 'd like to do better at that
<i>spk0</i>	I appreciate your honesty mister wang
<i>spk0</i>	what can you tell me about some of your goals over the next few years
<i>spk1</i>	my primary goal is to gain more work experience
<i>spk1</i>	so a position like this would help me meet that goal

Implementations

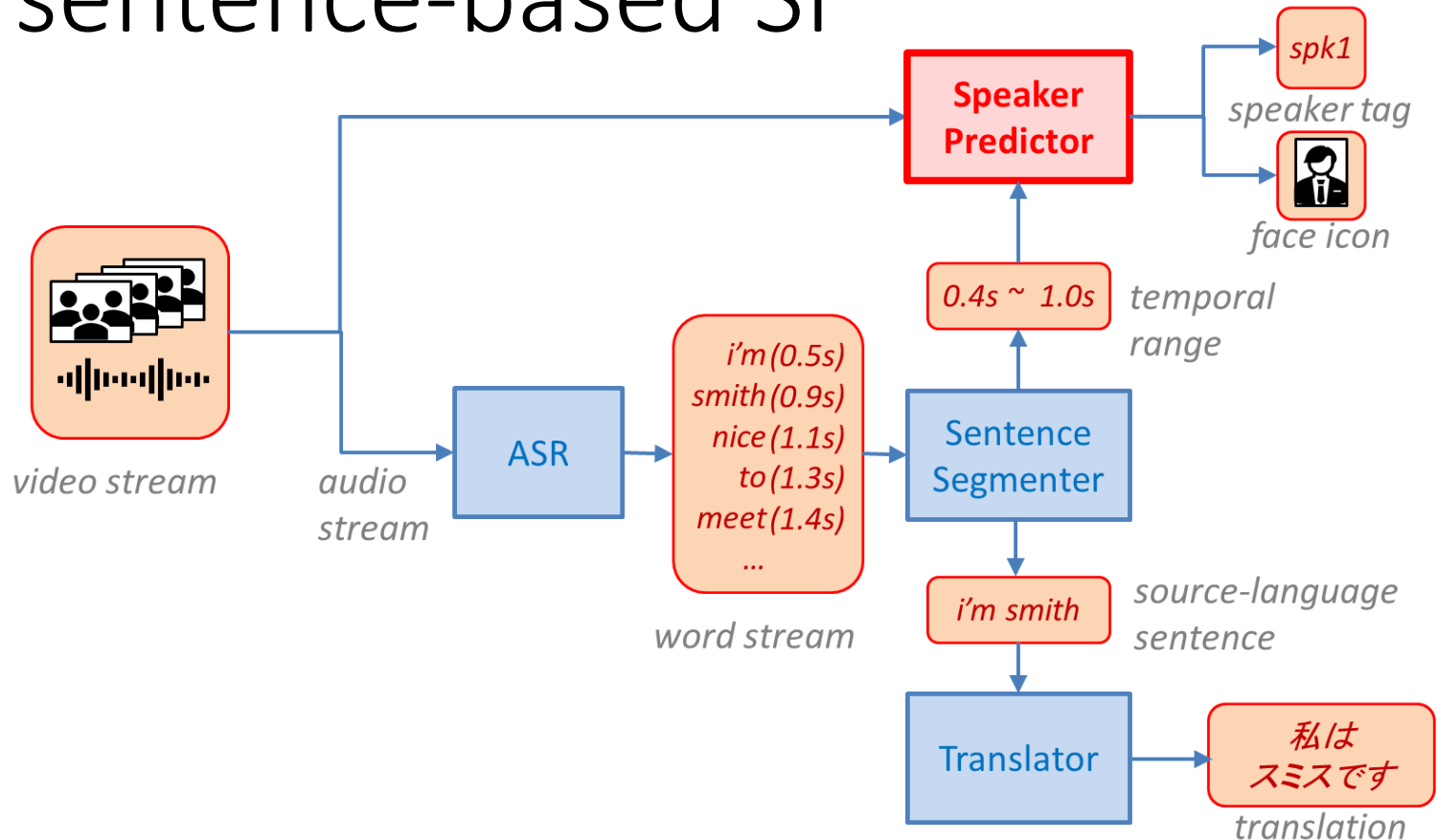
- Challenges
 - Recognize speakers from video streams
 - Maintain low latency for interpretation
- Solution:
 - Multimodal speaker recognition
 - Integrate speaker recognition with sentence-based SI

Multimodal speaker recognition



- Voice embedding
- Face embedding
- Active speaker detector: find the faces of speakers

Integrate speaker recognition with sentence-based SI



- Launch speaker predictor on video clips (video stream + temporal range)
 - Assemble large batches for GPU efficiency
- Parallelize translator with speaker predictor
 - Low latency

Weakness & Future

- Weakness
 - Highly dependent on textual-based sentence segmentor
 - Performs poorly on video with difficult speeches
 - Poor sentence segmentation
 - Poor translation quality
- Future: robust multimodal SI
 - Fuse audio and visual input with
 - Sentence segmenter
 - Translator
 - ASR

Thank you! (Q&A)