

Bridging different languages, countries, and cultures by Speech-to-speech Translation Research

Satoshi Nakamura

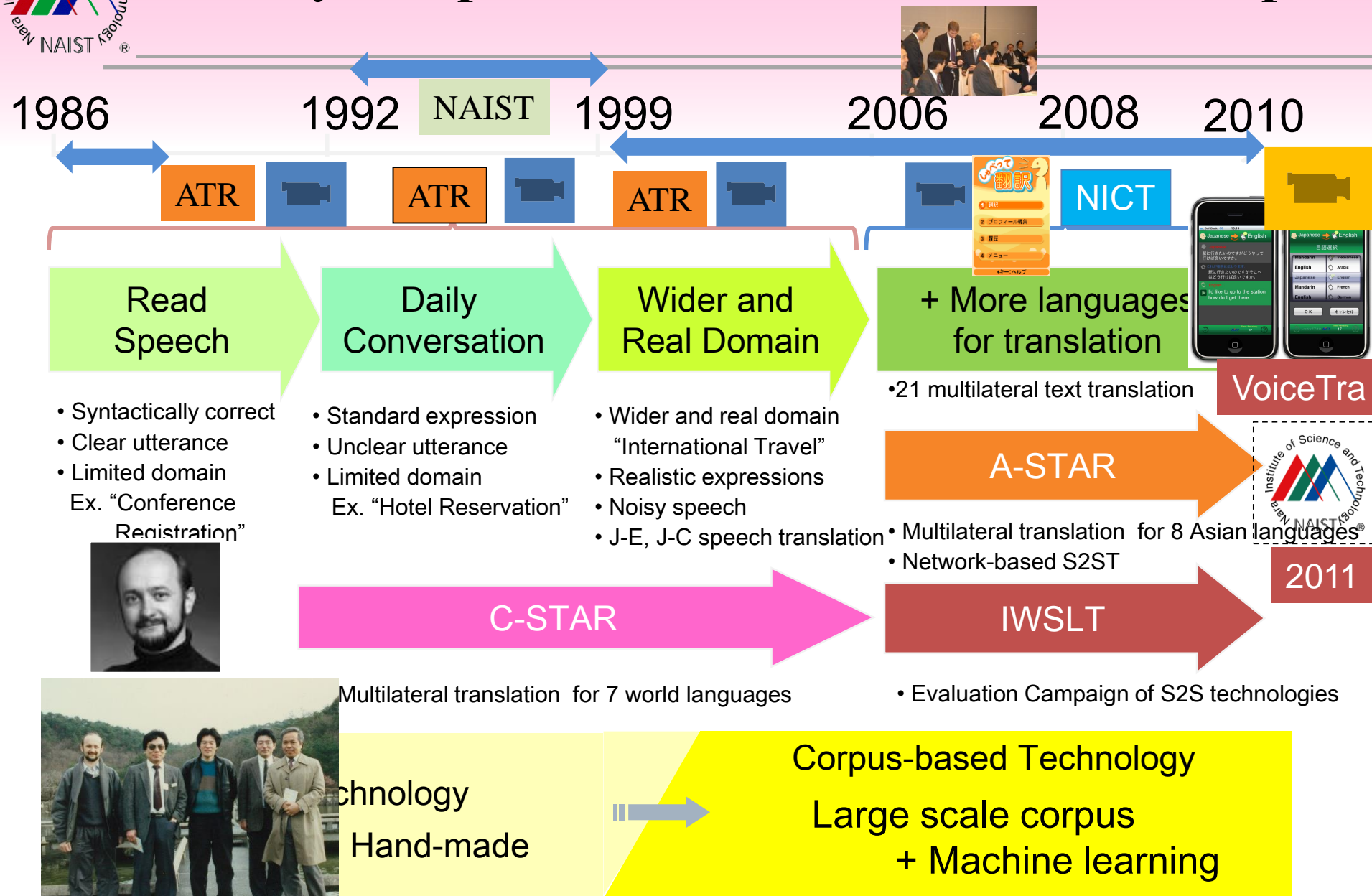
Director, Augmented Human Communication Laboratory,
Graduate School of Information Science,
Nara Institute of Science and Technology, Japan
(Formerly ATR and NICT)

Nakamura SRG

Augmented Human Communication (AHC) Laboratory

NAIST

History of Speech Translation Research in Japan



Speech Recognition

▶ 1986 at ATR

- HMM + n-gram : Collaboration with CMU (K. Lee, Sphinx)
- Spectrogram Reading : Collaboration with MIT (Victor Zue)
- Neural Network : Dr. Alex Waibel from CMU

▶ TDNN:

- Time Delay Neural Network
“Phoneme Recognition Using Time-Delay Neural Networks”, A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, K. Lang, IEEE Transactions on Acoustics, Speech and Signal Processing, March, 1989

→ **The IEEE 1990 Senior Best Paper award of the IEEE Acoustics, Speech and Signal Processing Society**

- MS-TDNN: Multi-scale Time Delay Neural Network
“Modularity and Scaling in Large Phonemic Neural Networks”, A. Waibel, H. Sawai, K. Shikano. IEEE Transactions on Acoustics, Speech and Signal Processing, December, December 1989

TDNN: Time Delay Neural Network

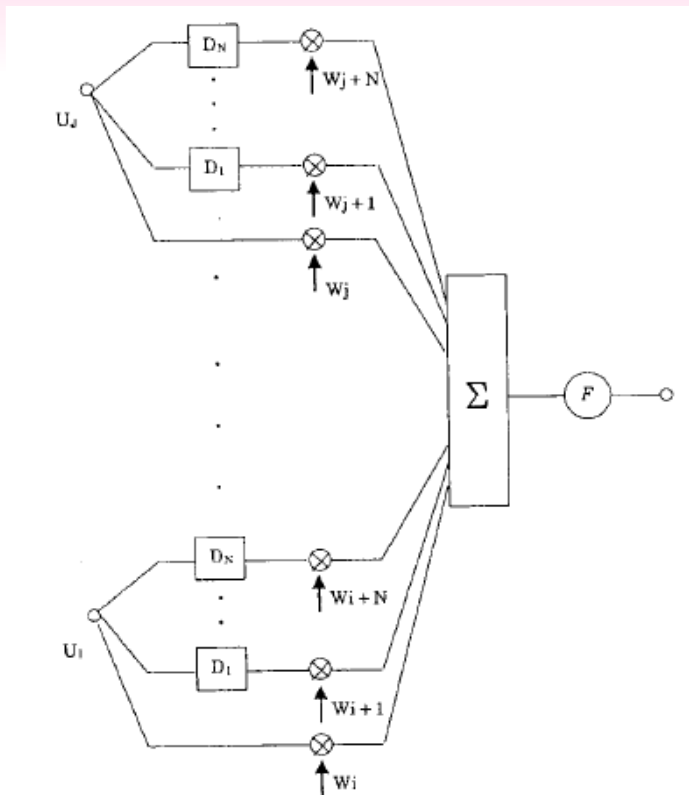


Fig. 1. A Time-Delay Neural Network (TDNN) unit.

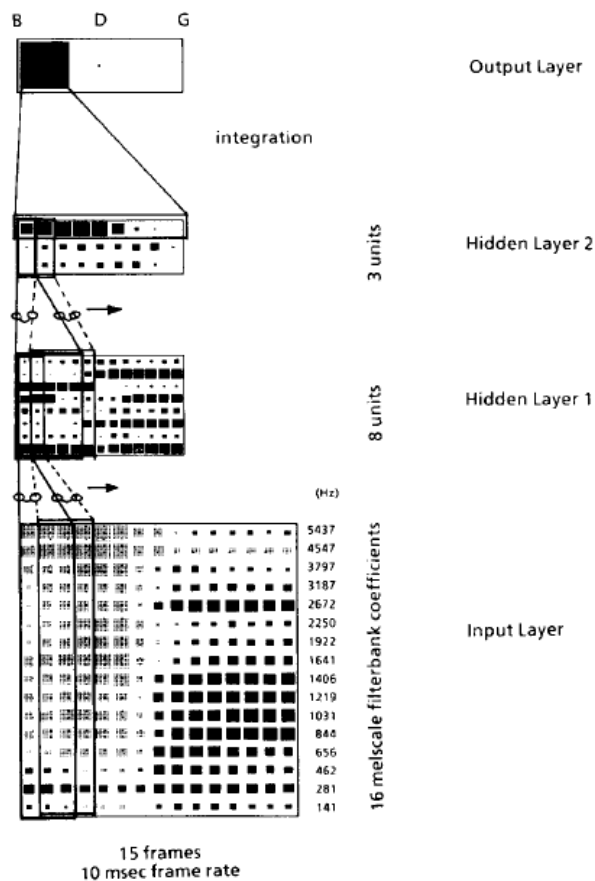
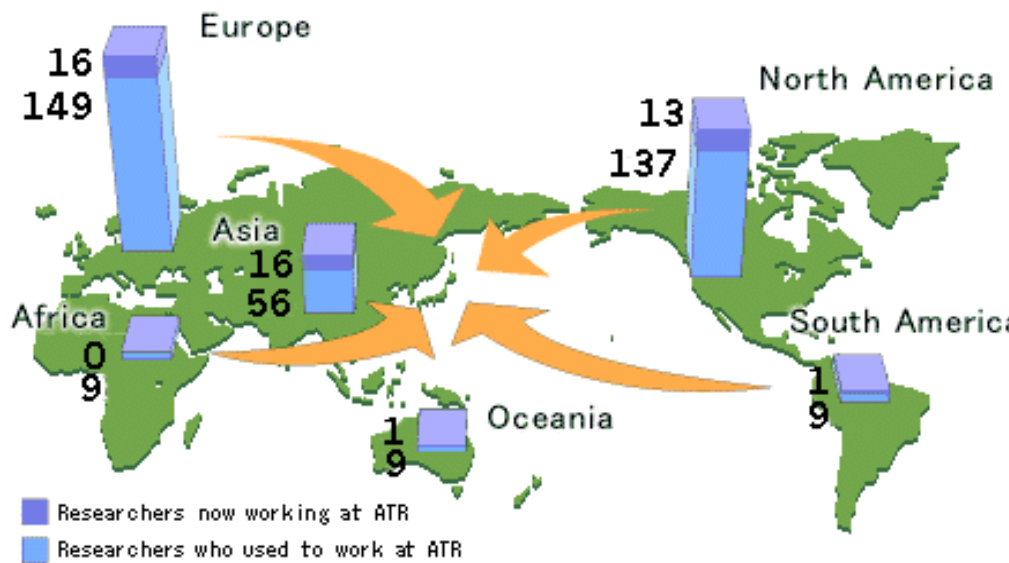
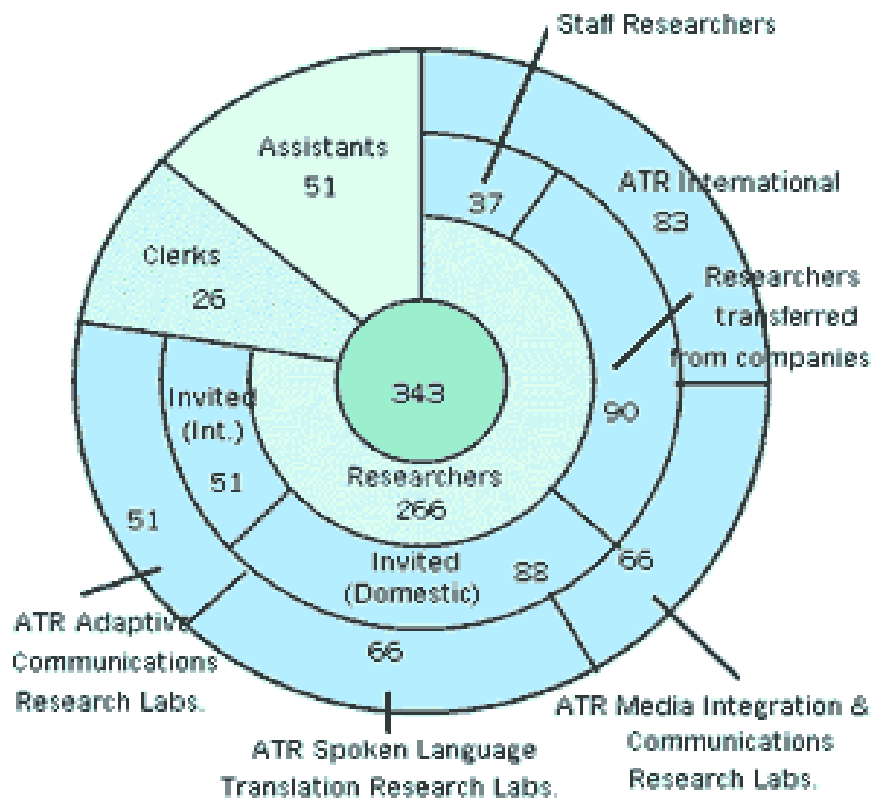


Fig. 2. The architecture of the TDNN.

A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, K. Lang, "Phoneme Recognition Using Time-Delay Neural Networks", IEEE Transactions on Acoustics, Speech and Signal Processing, March, 1989

Y. LeCun and Y. Bengio, "Convolutional Networks for Images, Speech, and Time-series," in The Handbook of Brain Theory and Neural Networks. MIT Press, 1995

ATR as of 2002



14 M Euro/

Year for Each Lab.

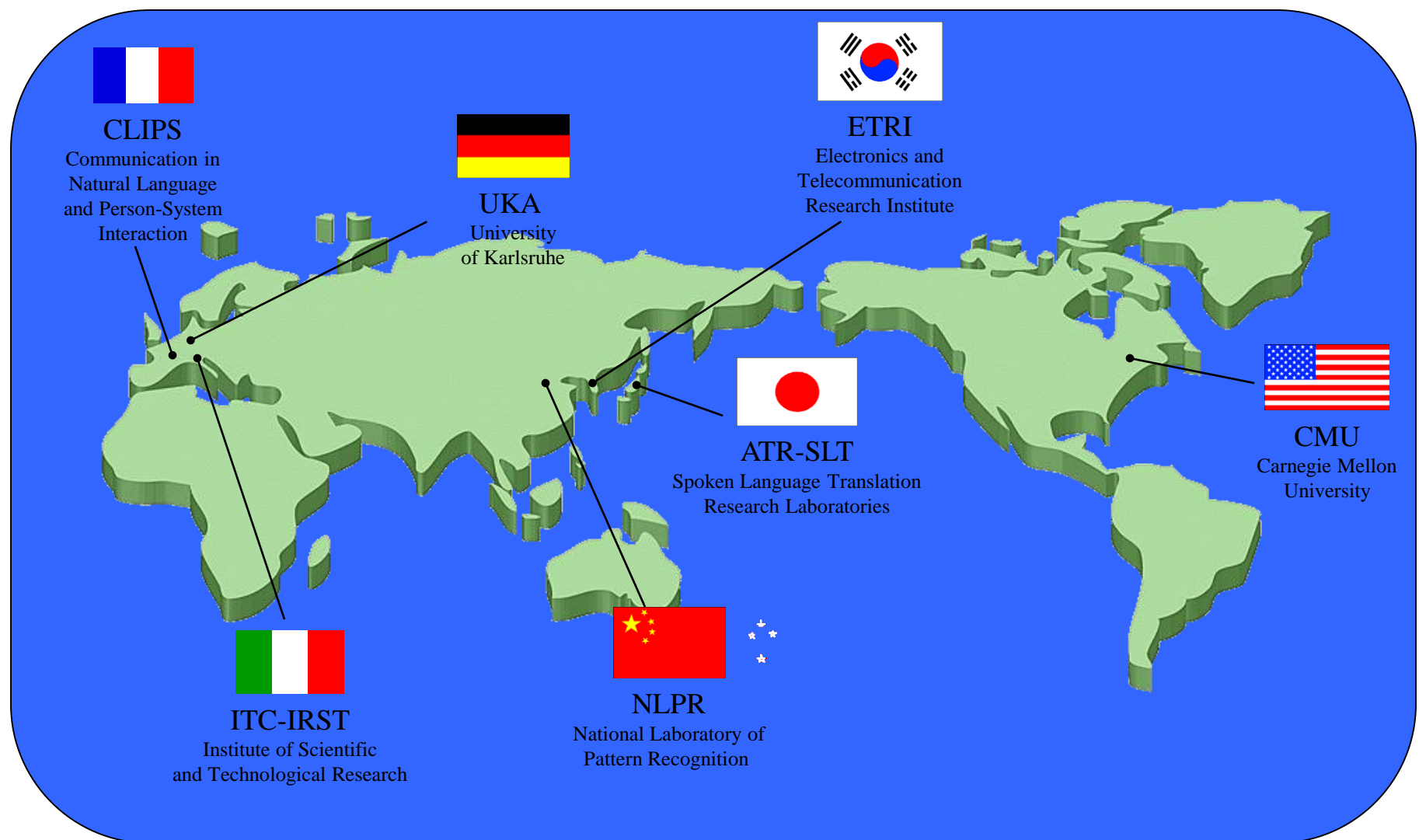
[Researchers from around the world at ATR]
18% of researchers are non-Japanese.

ATR bridged many international researchers !

C-star partners (2002)

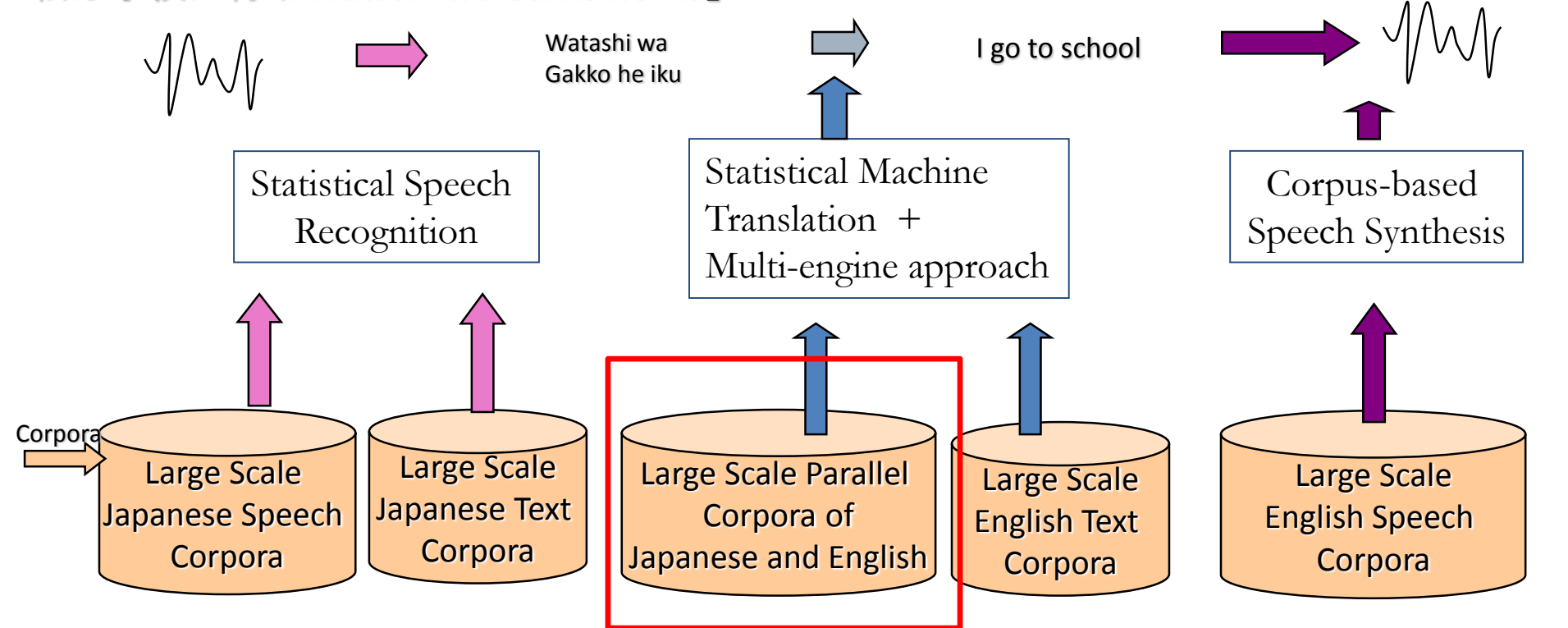
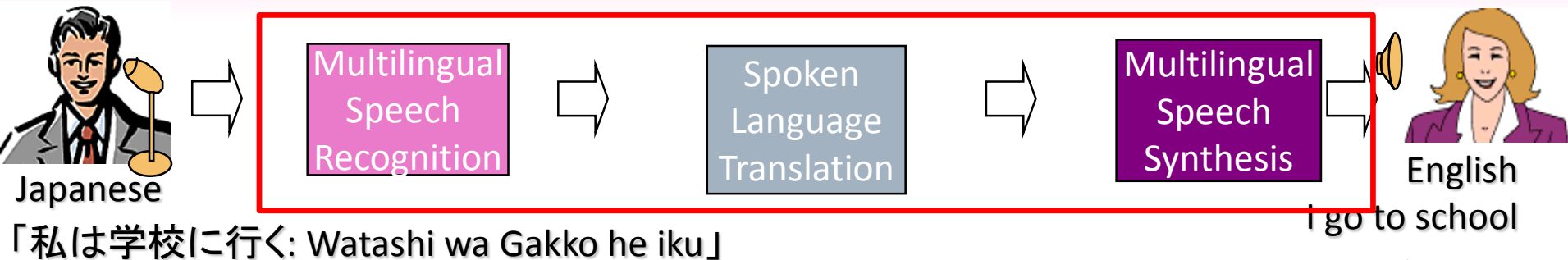
(Consortium for Speech Translation)

C-STAR bridged many international research institutes !



Speech-to-speech translation system

Speech-to-speech translation bridged speech and NLP researchers!



ATR BTEC Corpus

Basic

12.2% (7)

- greet someone
- ask a question
- state one's purpose
- ...

Trouble

12.1% (20)

- luggage
- emergency
- medicine
- assistance
- ...

Shopping

10.0% (13)

- buy something
- gather information
- price
- wrapping
- ...

Move

8.4% (8)

- transportation
- buy a ticket
- rental car
- trouble
- ...

Stay

8.2% (11)

- make/change a reservation
- check-in
- trouble
- ...

| | |
|----------------------|-----------|
| Sightseeing | 7.7% (11) |
| Restaurant | 7.3% (11) |
| Communication | 6.4% (6) |
| Airport | 5.5% (14) |
| Business | 5.3% (26) |
| Contact | 4.0% (6) |
| Airplane | 3.6% (11) |
| Homestay | 2.3% (11) |

| | |
|-----------------------|-----------|
| Study Overseas | 1.6% (14) |
| Drink | 1.3% (4) |
| Exchange | 1.2% (5) |
| Snack | 1.2% (4) |
| Beauty | 0.8% (5) |
| Go Home | 0.6% (4) |
| Research | 0.1% (12) |

BTEC

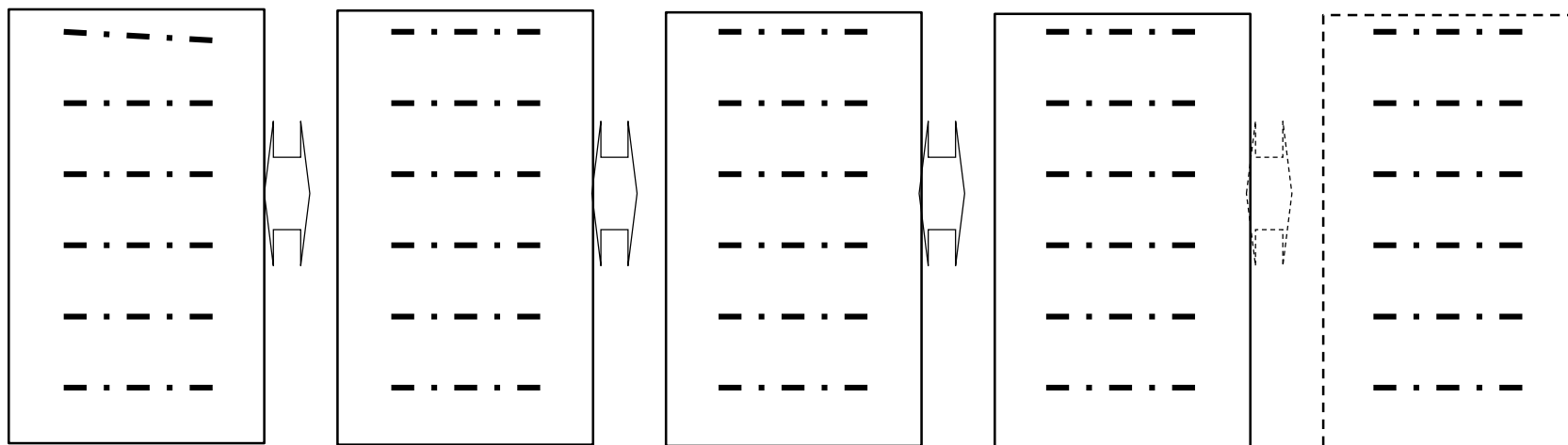
Japanese

English

Chinese

Korean

New lang.



Parallel sentences

Parallel sentences bridged many languages!

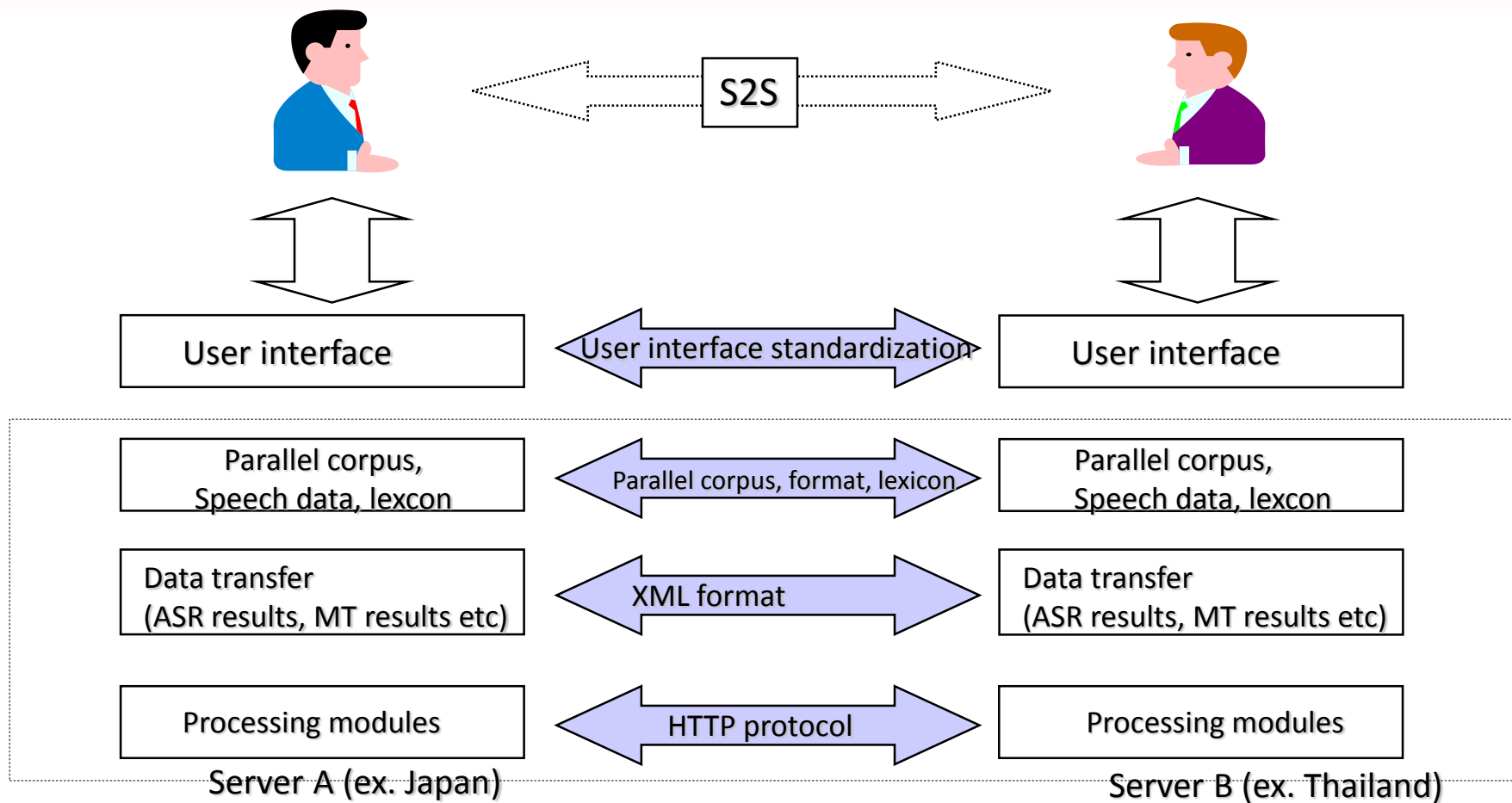


Status of Shared Corpus (uploaded to CSTAR server)

| | training corpus | | evaluation corpus | |
|-----------------------|----------------------------------|--------------------|-------------------|------------------------------|
| | <i>sentences</i> | <i>paraphrases</i> | <i>sentences</i> | <i>multiple references</i> |
| Japanese (ATR) | 162k | — | 506 sen | ≤ 16 [5x3 sen] |
| English (ATR) | | — | | ≤ 16 [5x3 sen] |
| Korean (ETRI) | | 309k | | ≤ 2 [1x2 sen] |
| Chinese (NLPR) | | — | | ≤ 3 [1x3 sen] |
| Italian (IRST) | 48k | 7k | | ≤ 6 [1x3 sen] |
| Spanish (CMU) | 6k + α? | 2k | | ≤ 10 [2x5 sen] |
| German (UKA) | β ? | — | — | — |
| French (CLIPS) | γ ? | — | — | — |

IWSLT bridged many international researchers!

Standardization at S2ST



- ◆ Activity start for standardization of Network-based S2ST at ITU-T SG16
- ◆ Session period : 2009- 2010
- ◆ NICT is the editor for S2ST standardization at ITU-T SG16, WP2 Q21/22

| Document | Title | Scope |
|----------|---|--|
| F.745 | Functional Requirements for Network-based S2ST | - Definition of Network-based S2ST - Functions and service requirements of network-based S2ST |
| H.625 | Architectural Requirements for Network-based S2ST | - Requirements of S2ST architecture - Definition of interface for Network-based S2ST |

- ◆ Not only language conversion but also the potentially added module like sign language are taken into account :
S2ST -> Modality conversion

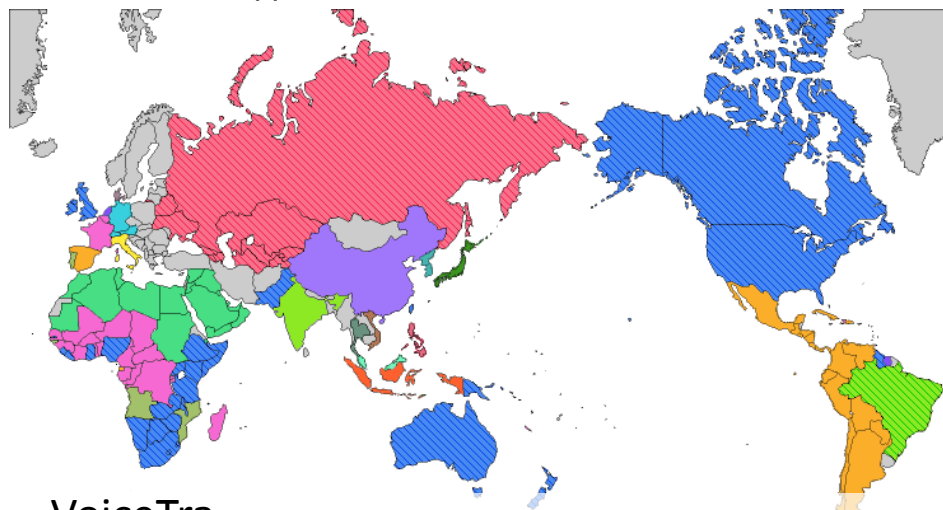
ITU-T standards bridge more languages and services!

VoiceTra, TexTra on iPhone (2010)

- A new speech translation software “VoiceTra” is available for iPhone at AppStore (released on July 29, 2010).
- **21** languages (including Ja, En, Ch, Ko) are covered, while **6** languages (including Ja, En, Ch) are **voice** enabled.
- So far approx. **800,000 download and 10 million accesses** have been achieved after the software release. (2012)



* Text-translation application, TexTra is released at the same time.



VoiceTra



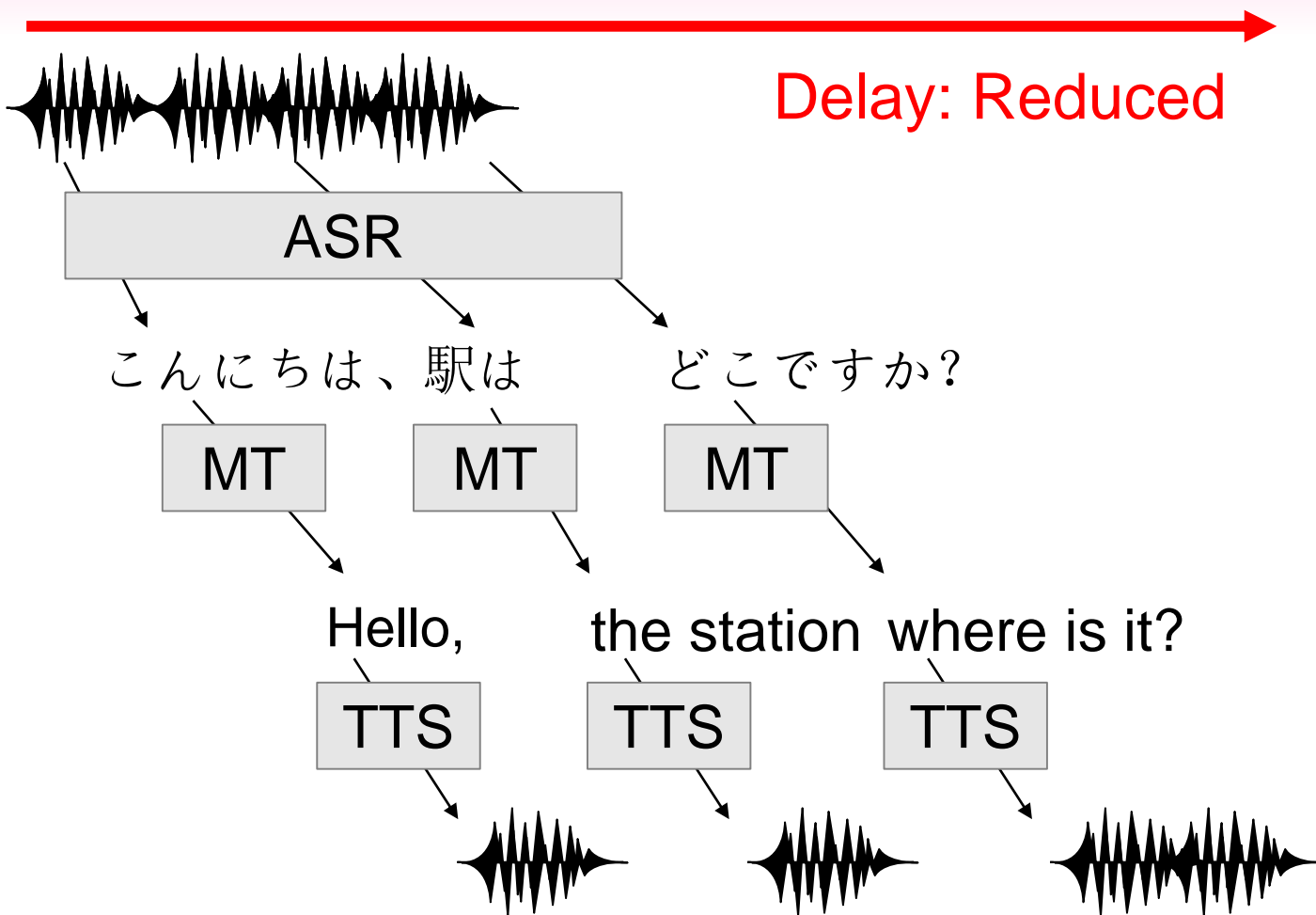
Japanese, English, Mandarin, Taiwanese
Mandarin, German, French, Dutch, Danish,
Italian, Spanish, Portuguese, Brazilian
Portuguese, Russian, Arabic, Hindi, Indonesian,
Malay, Thai, Tagalog, Vietnamese, Korean
 ※ Language in red can be input/output in voices.
 ※ There is no text input support for Hindi or Vietnamese.



Global Communication Project in Japan

- ▶ 2020 Olympic and Paralympic Game in Tokyo
- ▶ Increasing incoming tourists
 - 20 million international tourists to Japan in 2015
- ▶ Language support service in Japan
 - Initiated Ministry of Information and Communication, Internal Affairs of Japan, MIC
 - National Institute of Information and communications technology, NICT
 - Global Communication Project Consortium, GCP
 - Research and Development Group (Nakamura, NAIST)
 - Technology Transfer Group (Usami, KDDI)
 - Industry Consortium
 - Panasonic organizes contract project funded by MIC.
 - NTT, Fujitsu, NEC,
- ▶ Target: Speech-to-speech translation service for
 - Shopping, sightseeing, living, troubles, disaster management
 - Stores, hotels, sightseeing spots, hospital

Simultaneous Incremental Speech Translation



But, this is not easy!

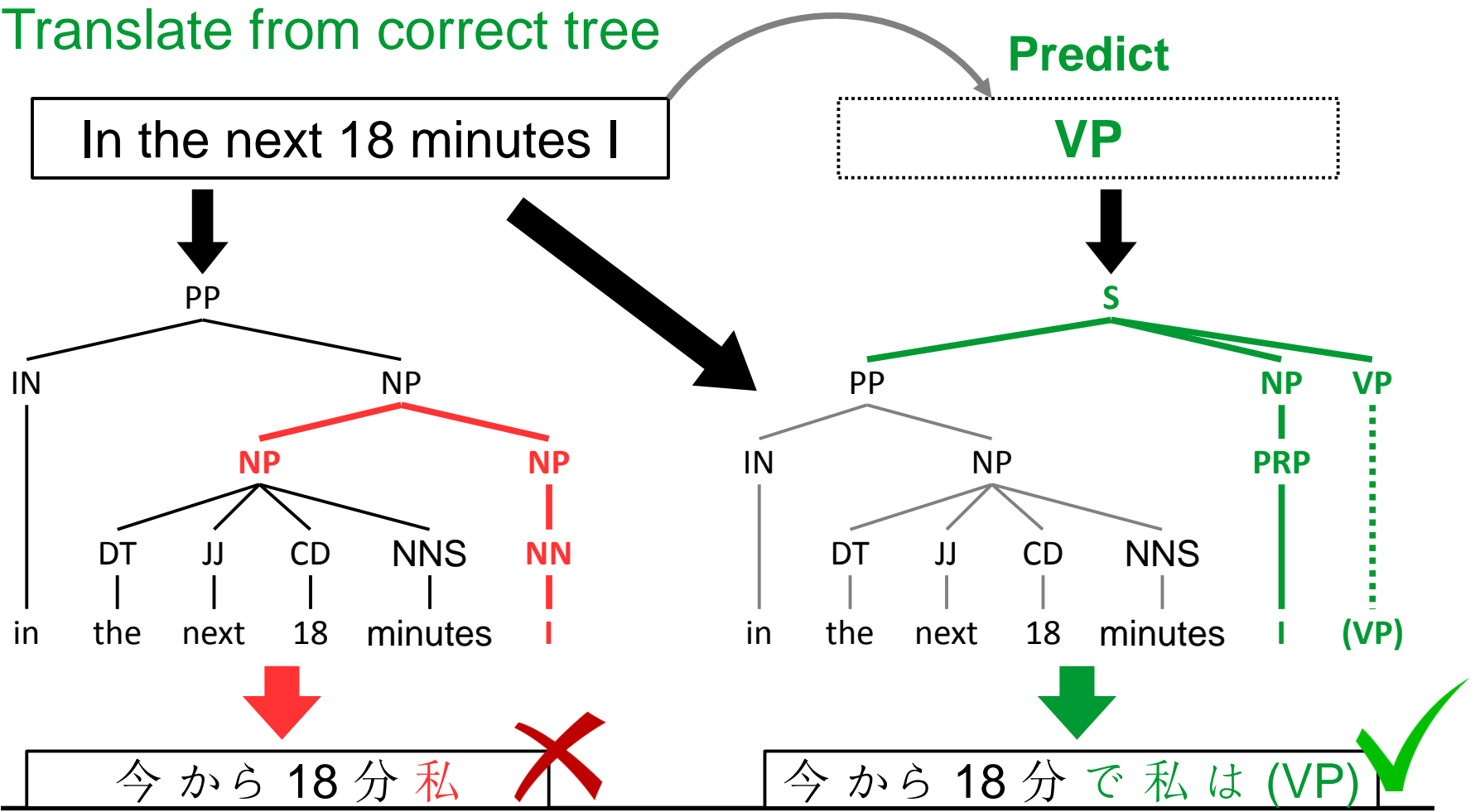
Can We Do the Same in Speech Translation Systems?

Four problems:

- **Segmentation:** When do we start translating?
- **Prediction:** Can we predict things that haven't been said?
- **Rewording:** Can we reword sentences to be conducive to simultaneous translation?
- **Evaluation:** How do we decide which results are better?

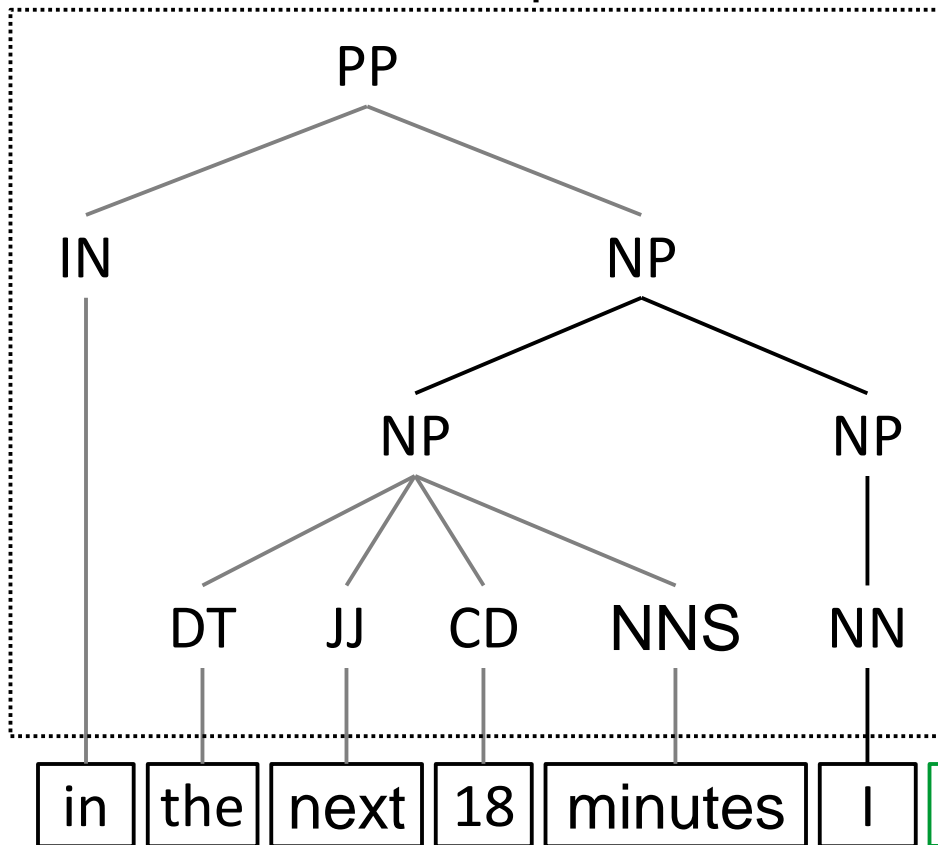
through Prediction of Unseen Syntactic Constituents [Oda+ ACL15]

- Predict unseen syntax constituents
- Translate from correct tree



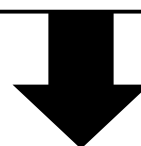
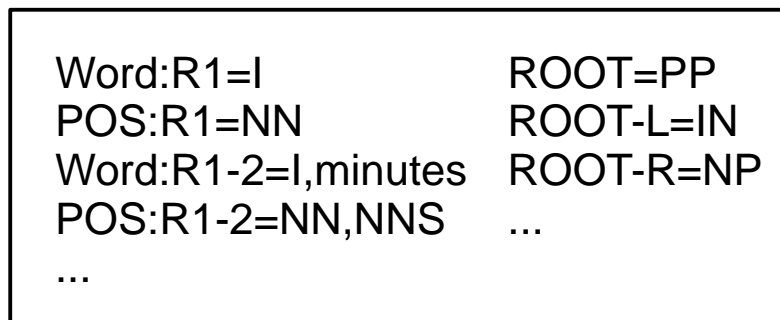
Syntax Prediction Process

1. *Parse* the input as-is



Input translation unit

2. *Extract* features



3. *Predict* the next tag (linear SVM)

VP ... 0.65
 NP ... 0.28
 nil ... 0.04
 ...

4. *Append* to sequence



5. *Repeat* until nil

Summary and Future Directions

▶ Speech-to-speech translation research

– bridged

- Different languages, international researchers, researchers in different fields, different services, international research institutes

– InterACT and C-Star contributed to the S2ST history.

▶ New Research Direction of Speech Translation at NAIST

– Simultaneous incremental speech-to-speech translation

– Emotion, para-linguistics, face, and gesture translation

– Neural MT modeling

▶ Future Work

– Para-linguistics and discourse structure

– Context and situation

– Background and specific domain knowledge

– Semantic structure and semantic analysis

– Communicability through the speech translation

NAIST joined interACT (2012)

