

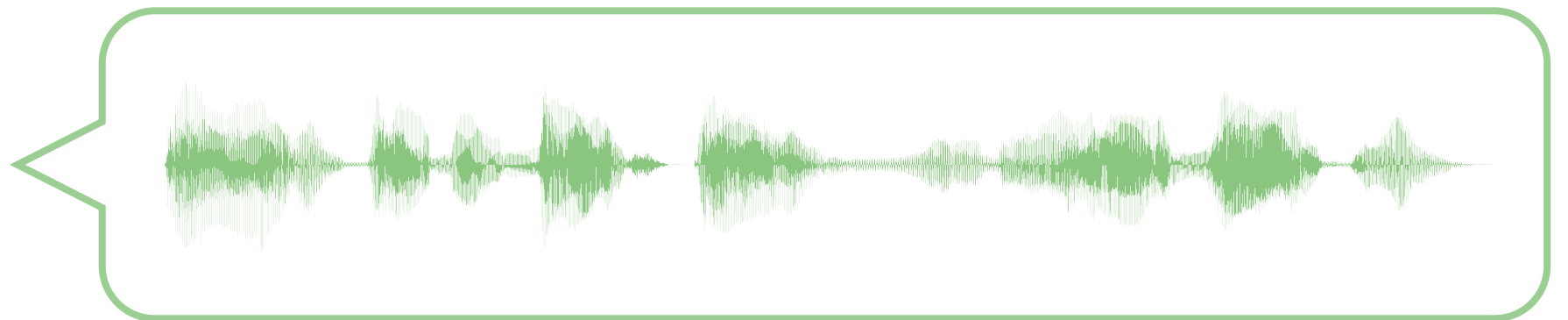
# **A Revisit to Feature Handling for High-quality Voice Conversion Based on Gaussian Mixture Model**

**Hitoshi Suda, Gaku Kotani,  
Shinnosuke Takamichi, and Daisuke Saito  
(The University of Tokyo)**

# Voice Conversion

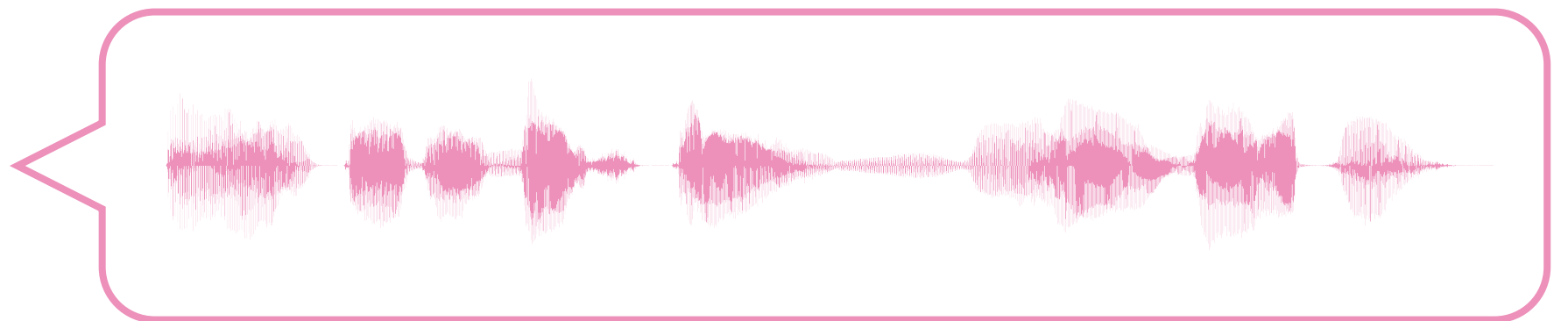
---

Source  
Speaker



Modifying personalities  
in the utterance

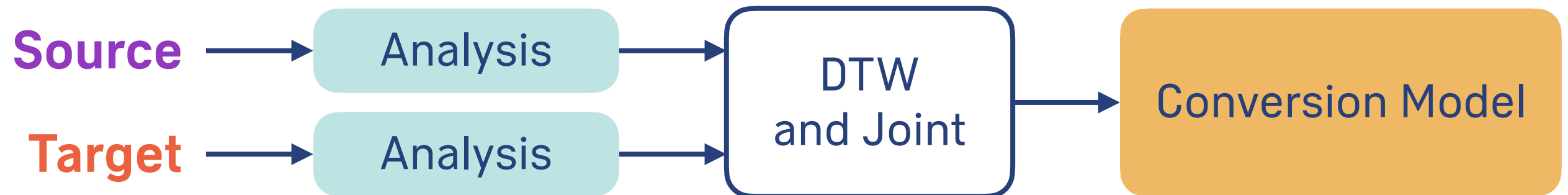
Target  
Speaker



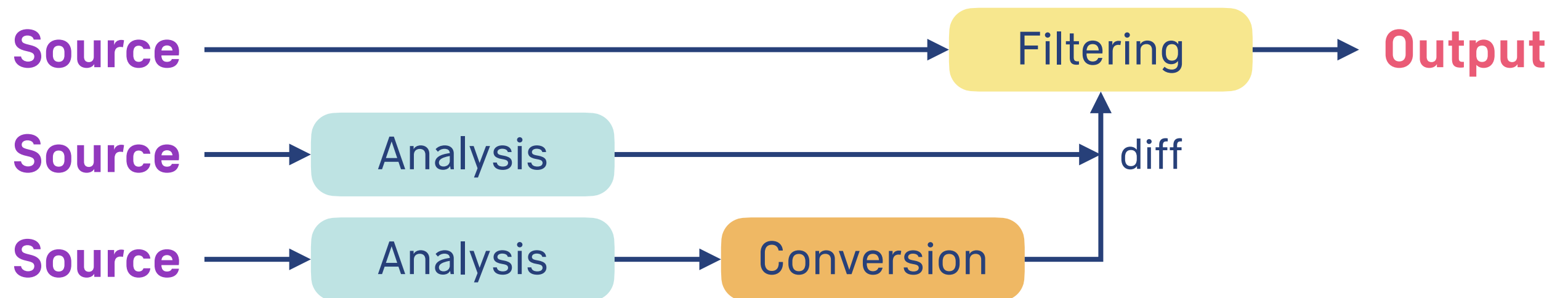
# Voice Conversion Framework

---

- Training



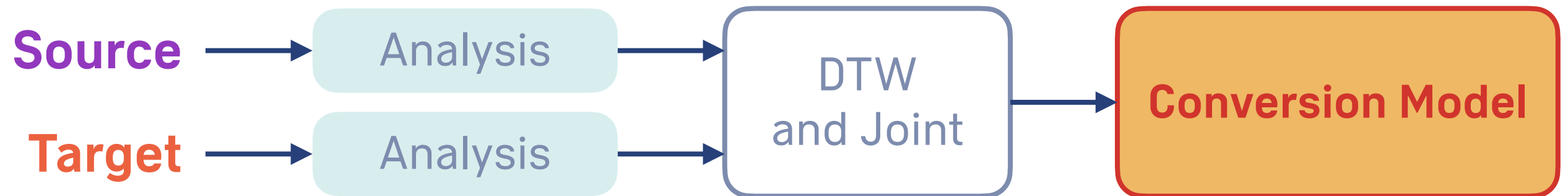
- Conversion



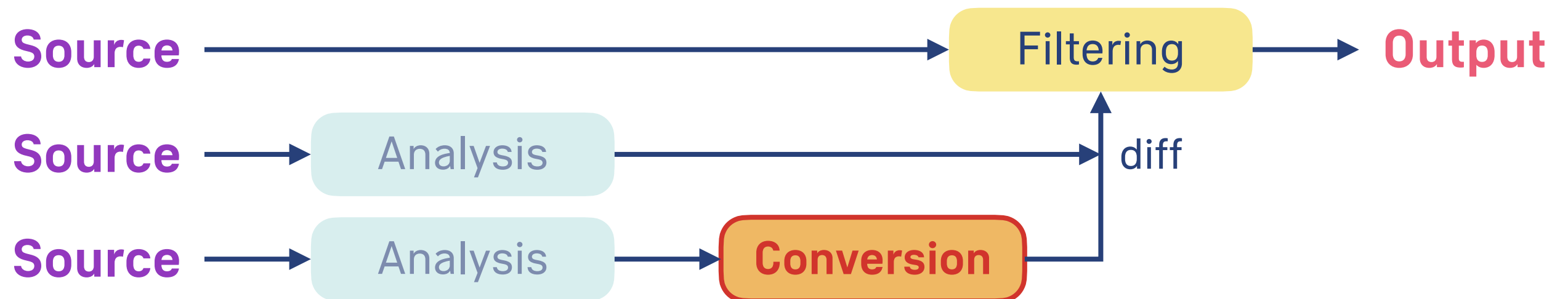
# Voice Conversion Framework

---

- Training



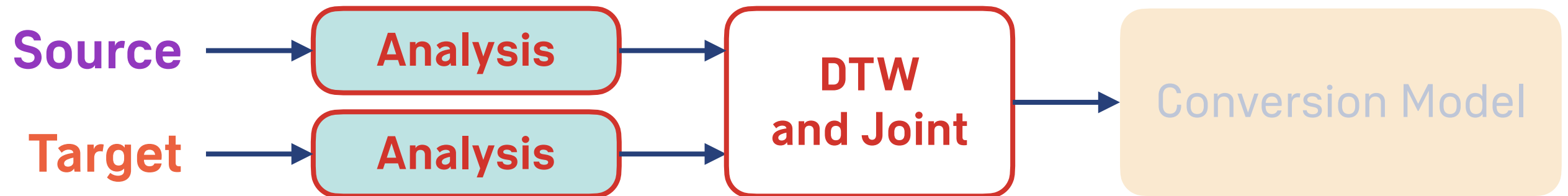
- Conversion



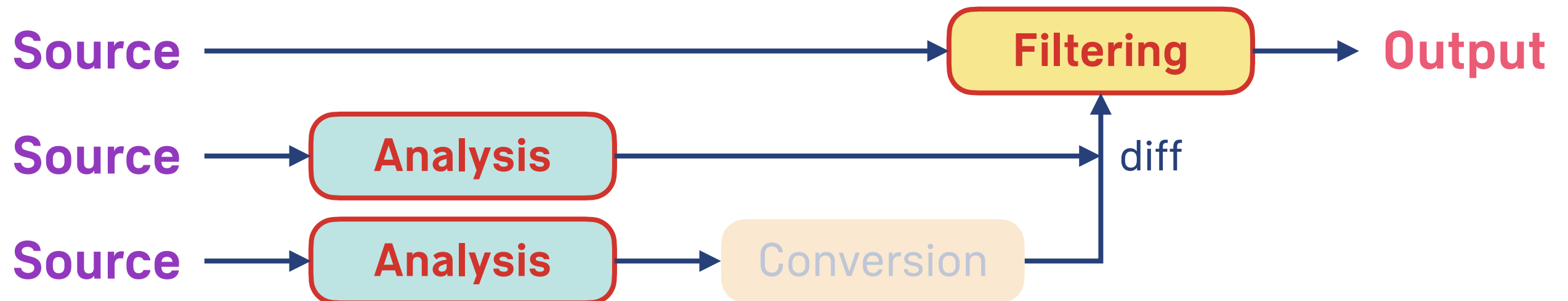
# Voice Conversion Framework

---

- Training



- Conversion



# Aim of This Study

---

## A Revisit to **Feature Handling** for High-quality Voice Conversion Based on Gaussian Mixture Model

- To **improve conversion quality** of existing voice conversion frameworks
- To experimentally reveal influences of feature handling
  - via subjective experiments

# Experimental Setups

---

- 50 sentences from ATR Japanese phonetically balanced sentence sets [Kurematsu+, 1990]
  - 40 for training, 10 for evaluation
- Sampling frequency: 22050 Hz
- Analysis and synthesis: WORLD [Morise+, 2016]
- Speakers: 2 males and 2 females
- Only intra-gender conversion / No  $F_0$  conversion
- 23 listeners answered questions in each preference test via our crowdsourcing system

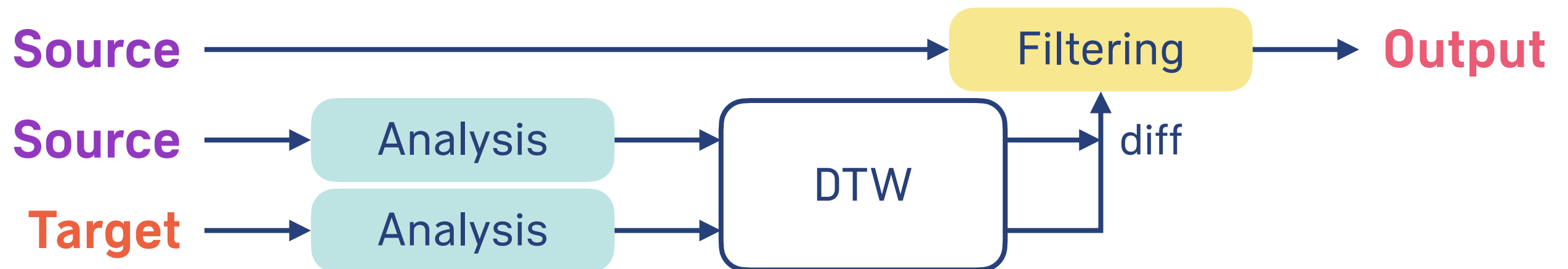
# 3 Experiments

---

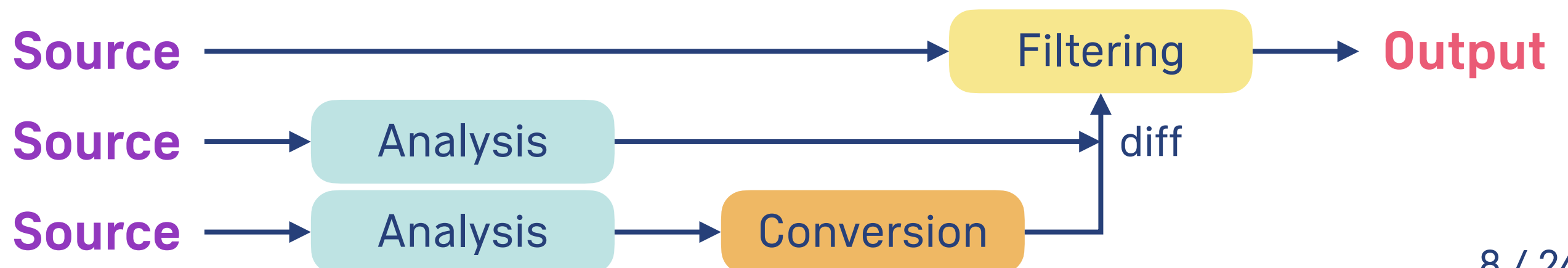
## 1. Analysis conditions



## 2. Conversion system without statistical mapping



## 3. Total conversion system





# Experiment 1: Analysis Conditions

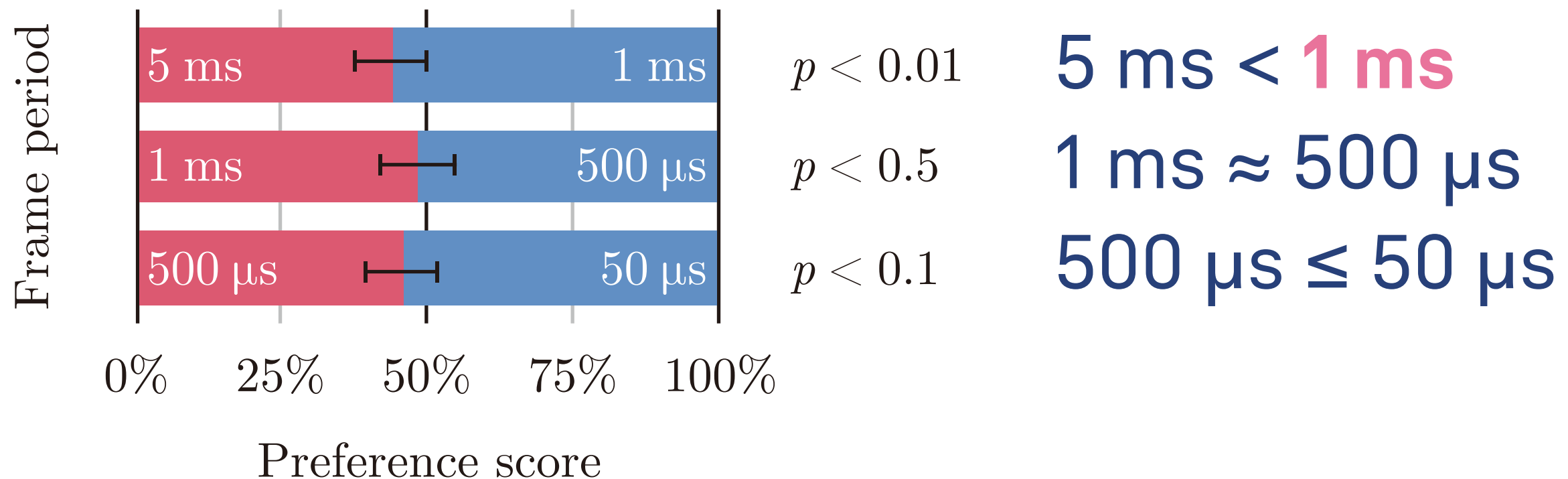
---



- To reveal the effects of conditions of analysis
  - Frame periods (or frame shifts)
    - How precisely the waveforms are analyzed in time domain
  - Order of mel-cepstral coefficients (mcep)
    - How precisely the spectral envelopes are represented

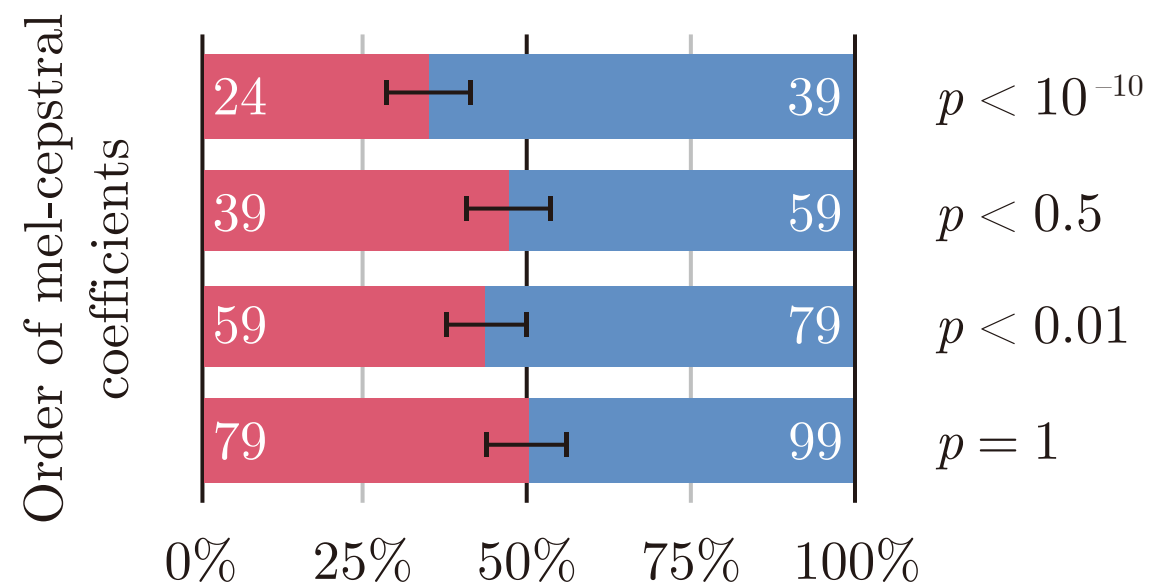
# Experiment 1: Analysis Conditions

- Frame periods:



# Experiment 1: Analysis Conditions

- Order of mcep:



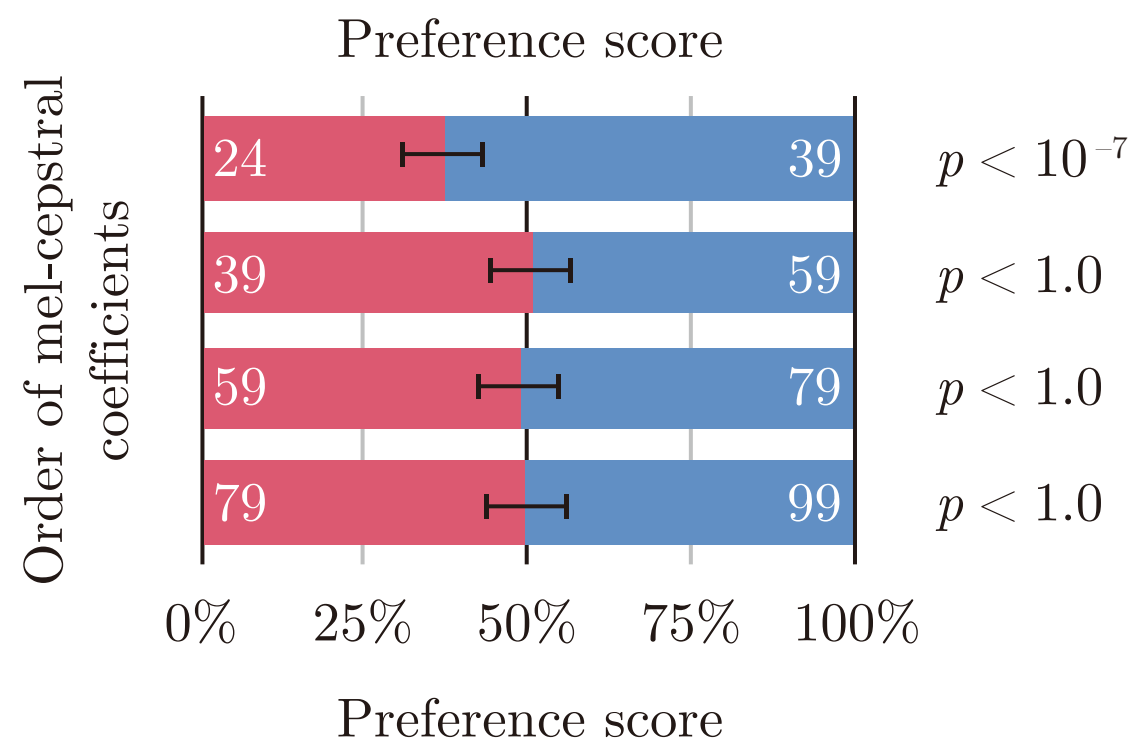
with 1 ms analysis

24 < 39

39  $\approx$  59

59 < 79

79  $\approx$  99



with 50  $\mu$ s analysis

24 < 39

39  $\approx$  59

59  $\approx$  79

79  $\approx$  99

# Experiment 1: Analysis Conditions

---

- Frame periods:

$$5 \text{ ms} < \mathbf{1 \text{ ms}} \approx 500 \text{ } \mu\text{s} \leq 50 \text{ } \mu\text{s}$$

- Order of mcep:

$$24 < 39 \approx 59 < \mathbf{79} \approx 99 \quad (\text{with } 1 \text{ ms frames})$$

$$24 < \mathbf{39} \approx 59 \approx 79 \approx 99 \quad (\text{with } 50 \text{ } \mu\text{s frames})$$

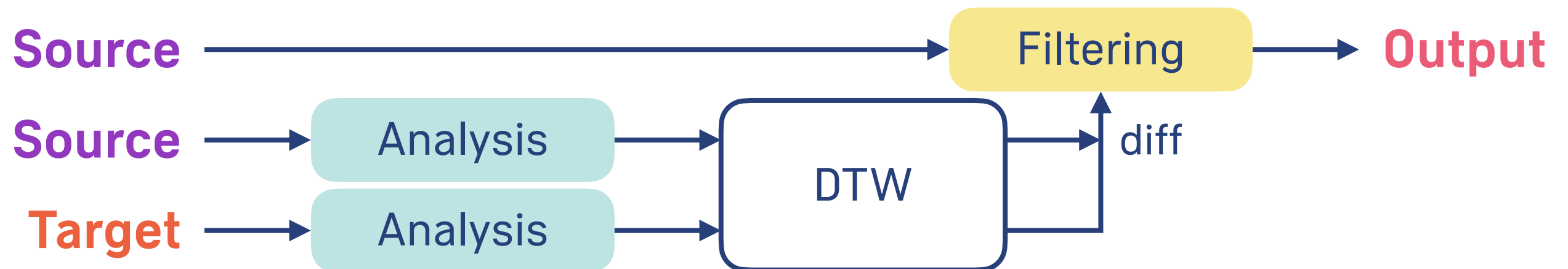
# 3 Experiments

---

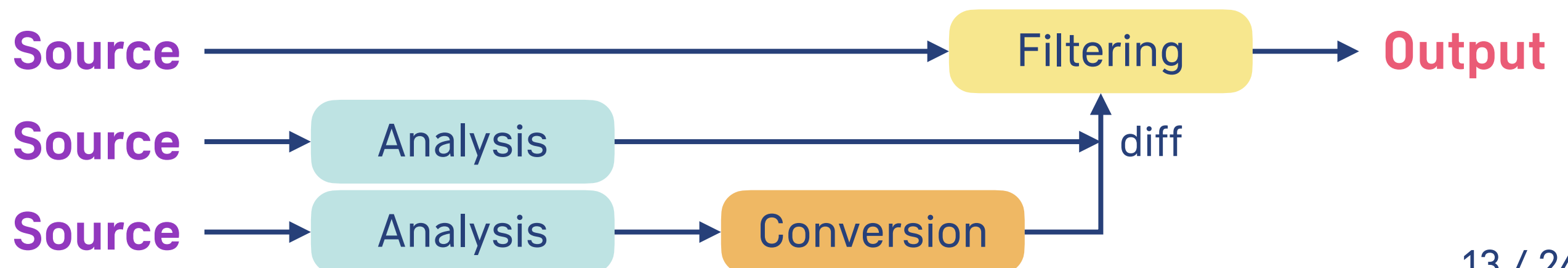
## 1. Analysis conditions



## 2. Conversion system without statistical mapping

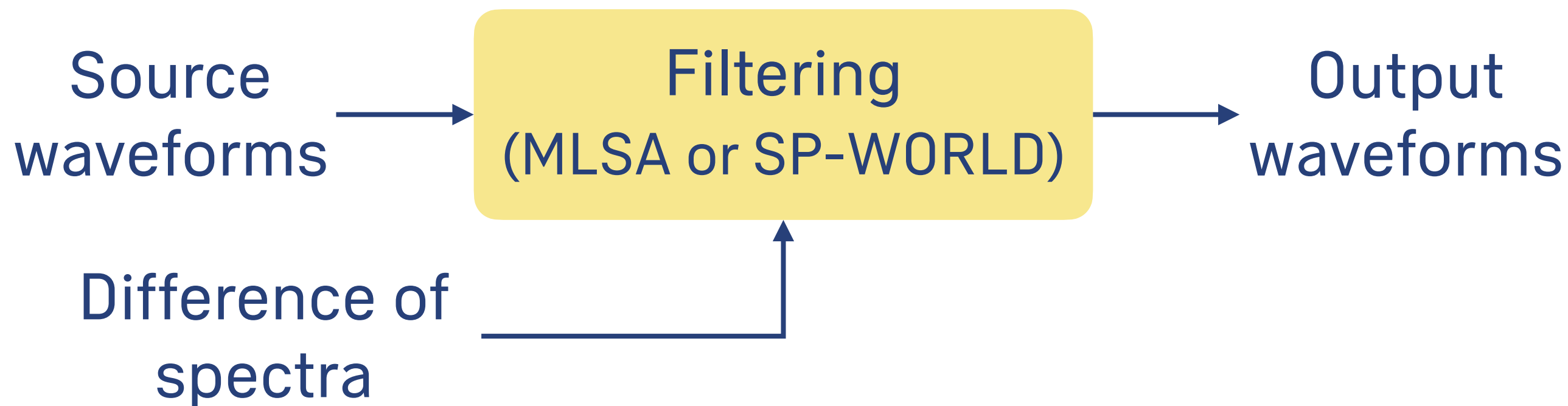


## 3. Total conversion system



# Differential-spectrum Compensation

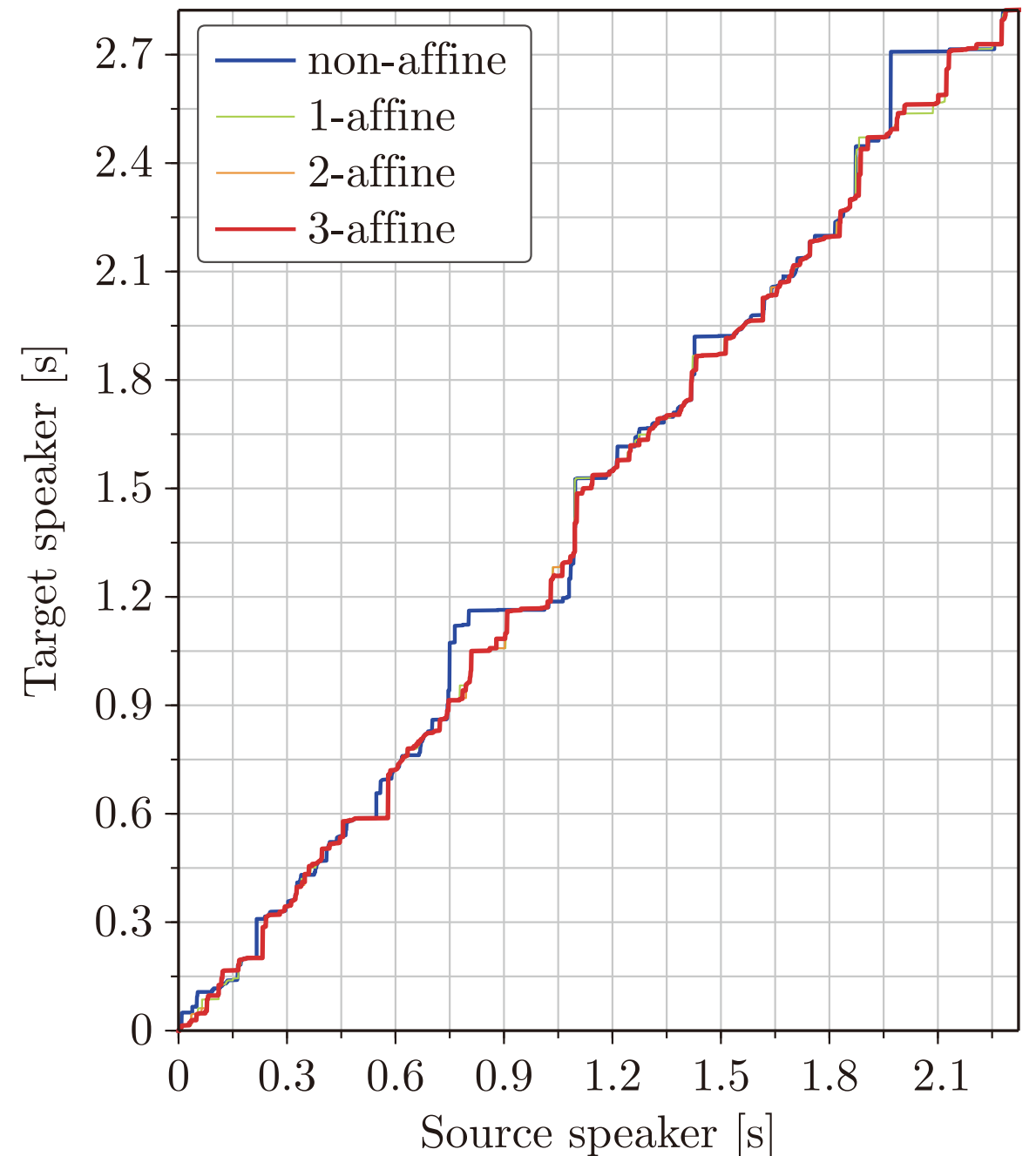
(diffspec) [Kobayashi+, 2014]



- Famous implementation: Mel log spectrum approximation (MLSA) Filtering [Imai+, 1983]
- We introduce a diffspec method “**SP-WORLD**” inspired by WORLD vocoder

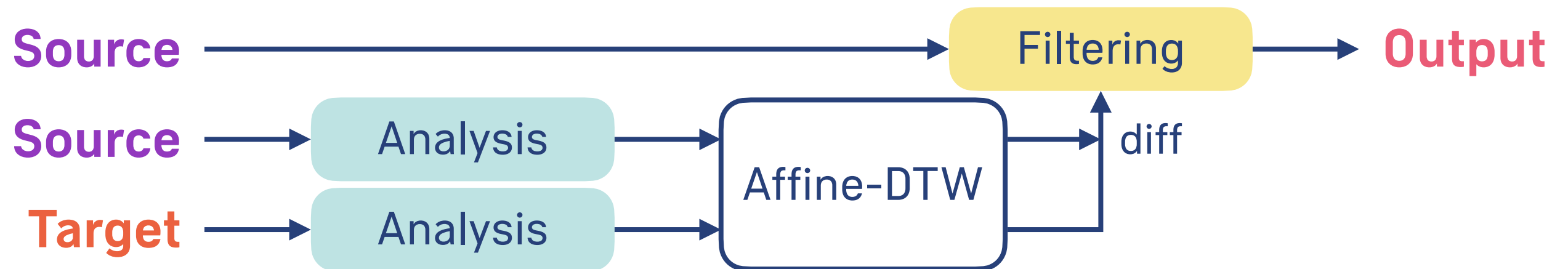
# Dynamic Time Warping (DTW)

- Alignment of features
- Sensitive to difference of individuality
- We introduce “Affine-DTW”
  - Iteration of alignment and coarse conversion



# Experiment 2: Ideal Conversion

---

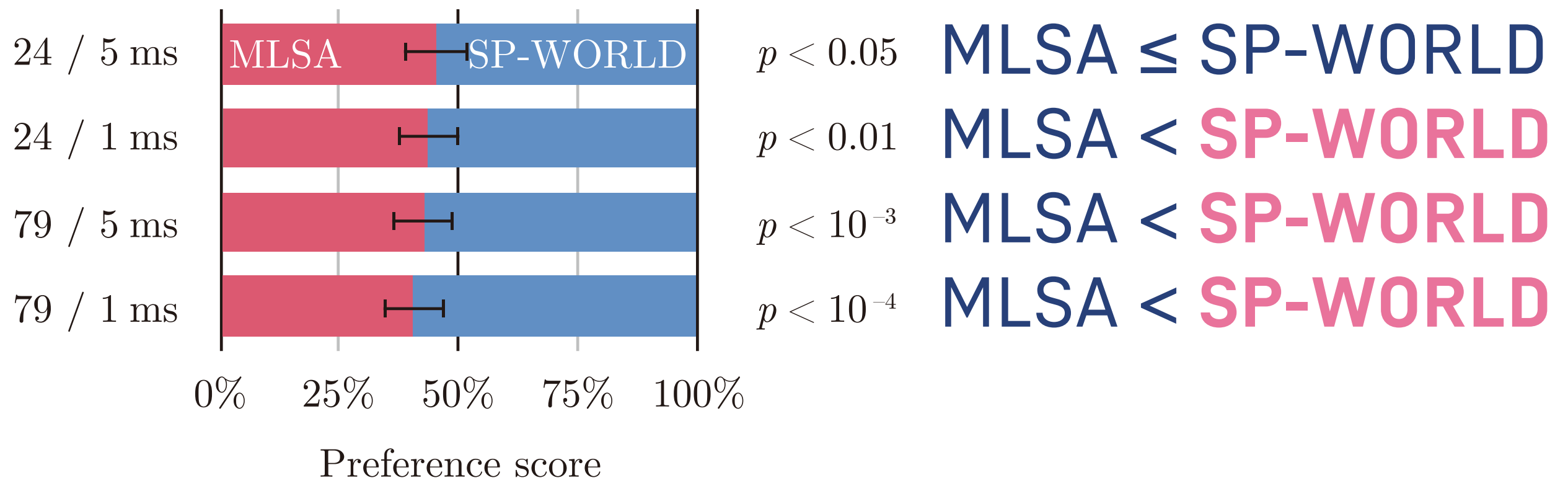


- To reveal the effects on quality of conversion of the conditions except mapping models
  - Diffspec method: MLSA or SP-WORLD
  - Analysis conditions
    - Frame period and order of mcep



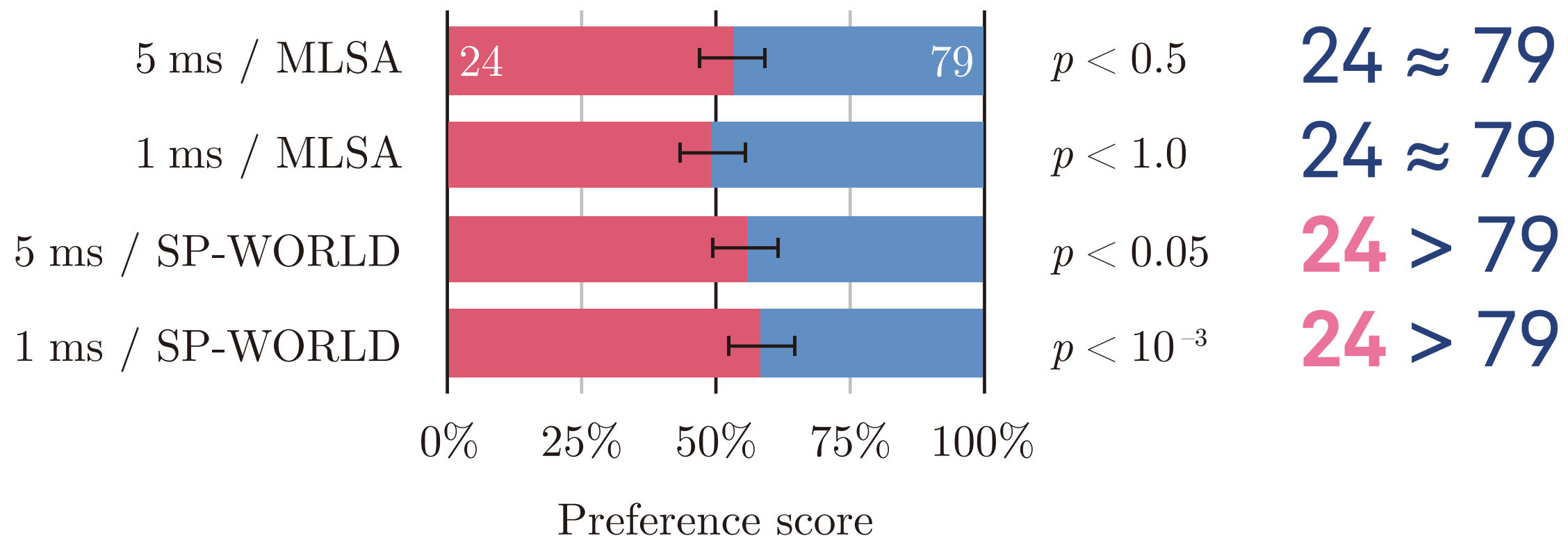
# Experiment 2: Ideal Conversion

- Diffspec method:



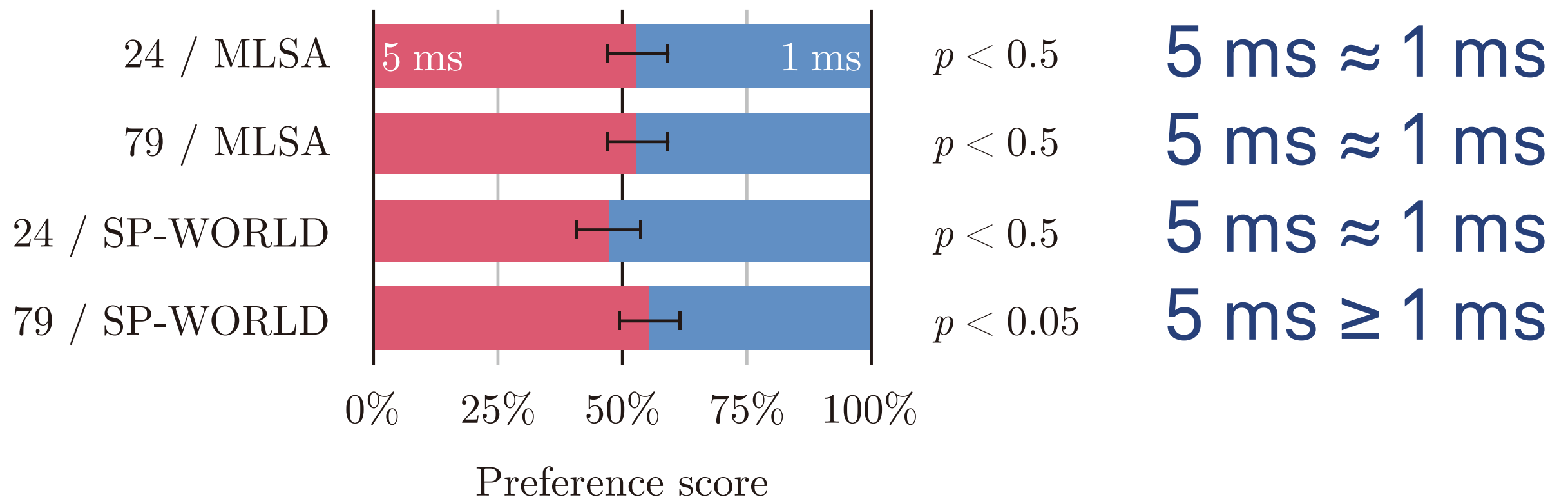
# Experiment 2: Ideal Conversion

- Order of mcep:



# Experiment 2: Ideal Conversion

- Frame periods:



# Experiment 2: Ideal Conversion

---

- Diffspec method: MLSA < **SP-WORLD**
- Order of mcep: **24** > 79 (with SP-WORLD)  
24  $\approx$  79 (with MLSA)
- Frame periods:  
5 ms  $\geq$  1 ms (with SP-WORLD / 79-order)  
5 ms  $\approx$  1 ms (with other conditions)

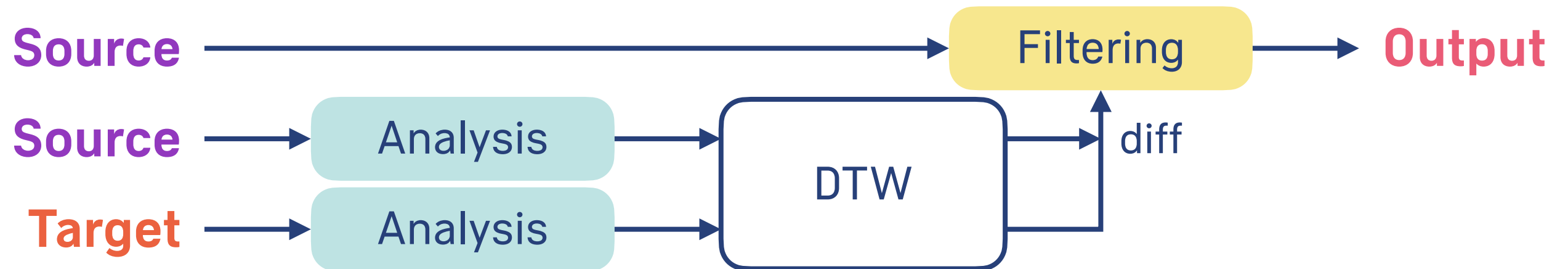
# 3 Experiments

---

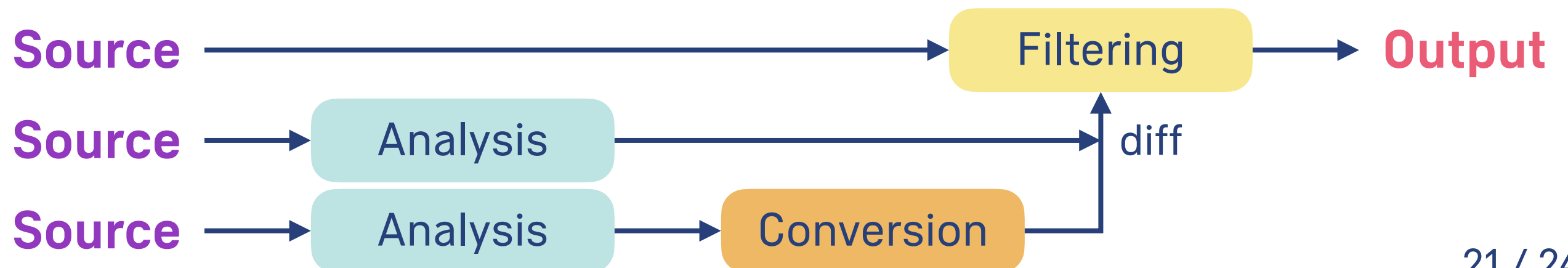
## 1. Analysis conditions



## 2. Conversion system without statistical mapping

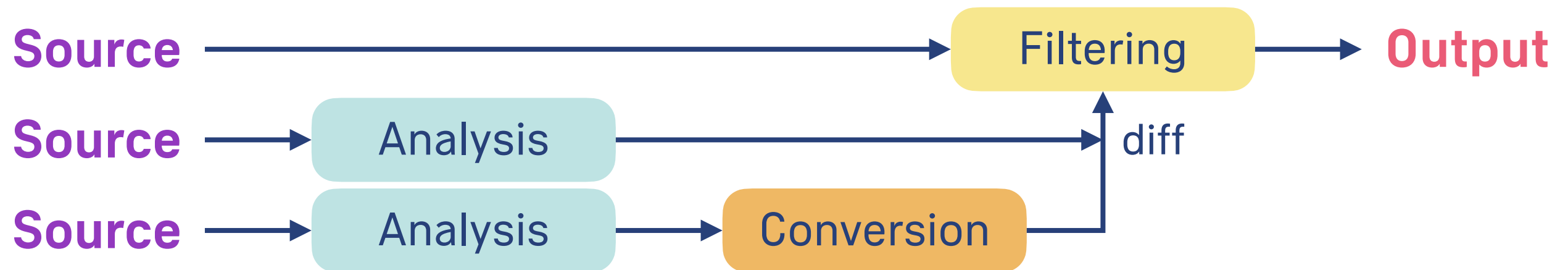


## 3. Total conversion system



# Experiment 3: Statistical Conversion

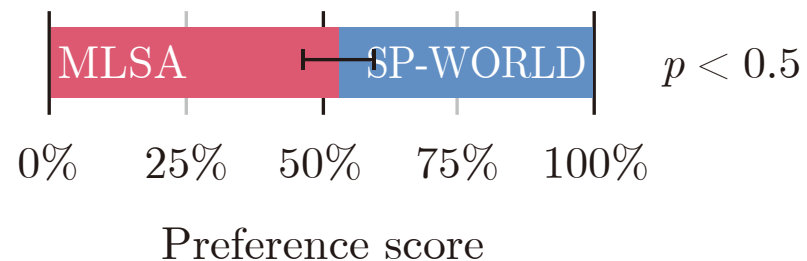
---



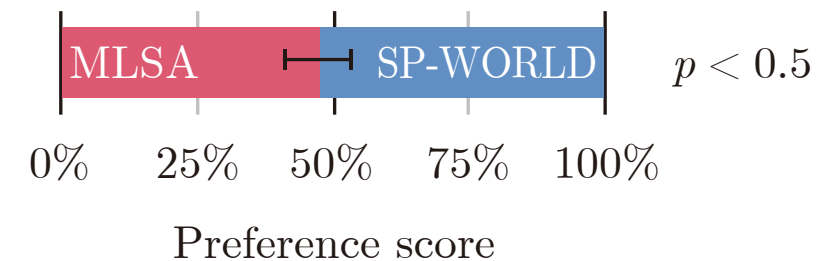
- To reveal the influences of below components
  - Diffspec method: MLSA or SP-WORLD
  - Sequence features [Toda+, 2007]
    - Dynamic features and global variances
- 1 ms period / 24-order of mcep

# Experiment 3: Statistical Conversion

- Diffspec method:  $\text{MLSA} \approx \text{SP-WORLD}$



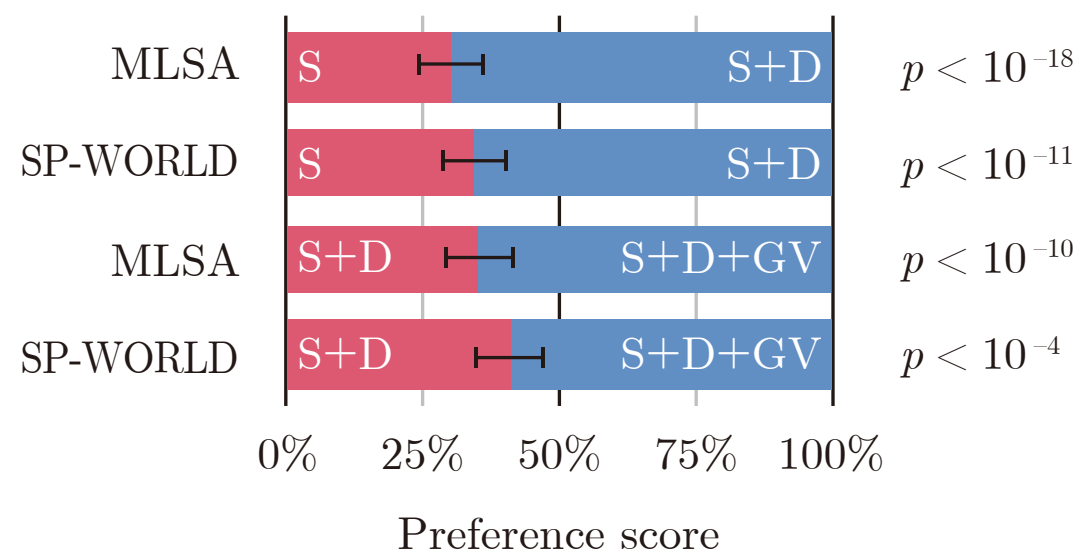
speaker similarity



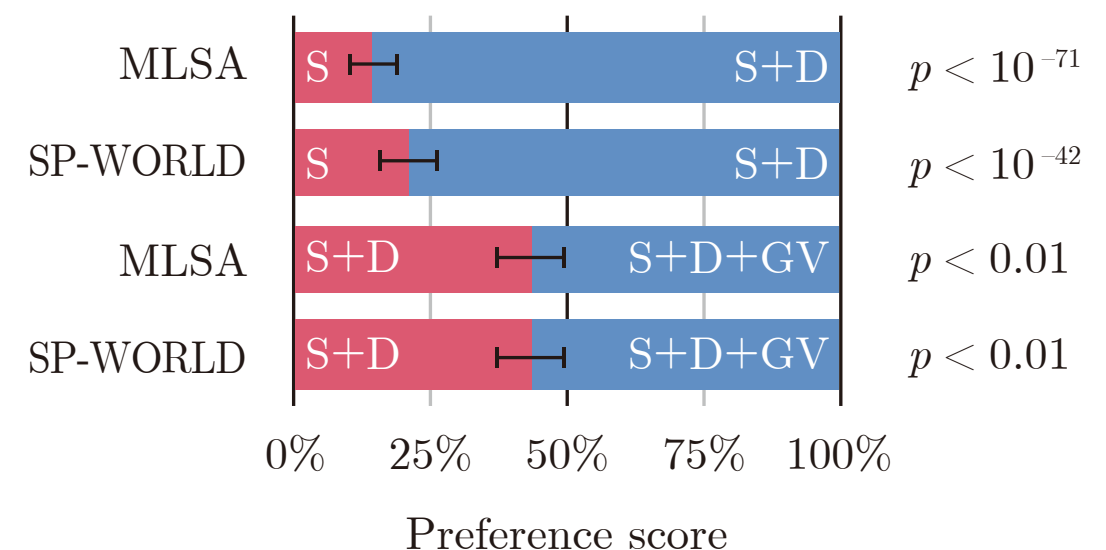
naturalness

- Sequence features:

static < static+dynamic < **static+dynamic+GV**



speaker similarity



naturalness

# Conclusion

---

- In GMM-based statistical voice conversion,
  - Dynamic features and GV: definitely effective
  - SP-WORLD: comparable to MLSA
    - Superior in ideal conversion
  - Higher order of features: **not** always effective
- High time-resolution analysis is effective in analysis-synthesis
  - Potential of effectiveness also in conversion

## Future Works

- $F_0$  conversion
- Break the 1 ms barrier in WORLD analysis
- Other mapping models such as neural networks