



# **SelfRemaster: Self-Supervised Speech Restoration with Analysis-by-Synthesis Approach Using Channel Modeling**

**Takaaki Saeki**, Shinnosuke Takamichi, Tomohiko Nakamura,  
Naoko Tanji, and Hiroshi Saruwatari

The University of Tokyo, Japan

- Background
- Related Work
- Methods
- Experimental Evaluation
- Summary

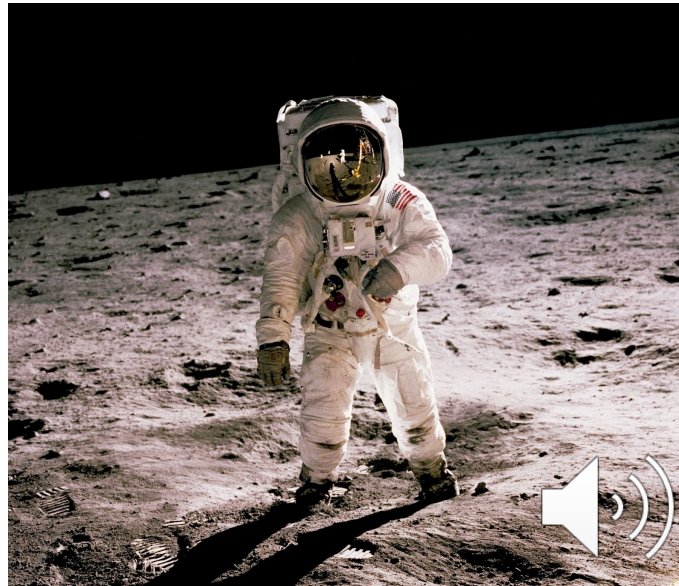
# Background: Speech Restoration

3/38

**Need to use and analyze existing degraded speech data.**

E.g.) Historical audio materials, telephone recordings, etc.

Containing low-resource languages, endangered cultures, etc.



[https://parade.com/.image/c\\_limit%2Ccs\\_srgb%2Cq\\_auto:good%2Cw\\_700/MTkwNTc5NTlyNDMyNTQyNTg4/1-19-martin-luther-king-ftr.webp](https://parade.com/.image/c_limit%2Ccs_srgb%2Cq_auto:good%2Cw_700/MTkwNTc5NTlyNDMyNTQyNTg4/1-19-martin-luther-king-ftr.webp)

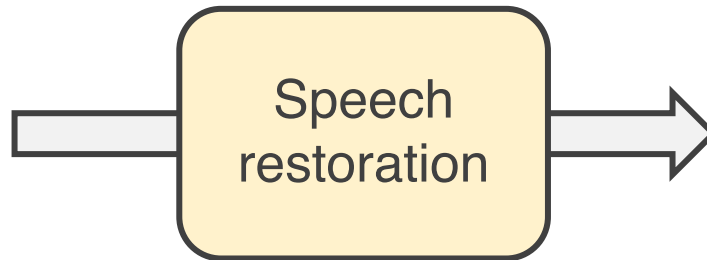
<https://i.insider.com/5d2faa4d7e76cc3f20437ff6?width=700>

<https://www.fluentu.com/blog/japanese/wp-content/uploads/sites/6/2021/09/classic-japanese-movies-5.jpg>

**Speech restoration:** generating clean speech from degraded speech.

Low-quality old recordings

High-quality restored audio



**Speech restoration:** generating clean speech from degraded speech.

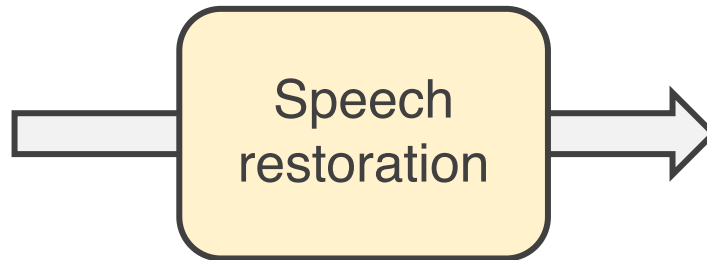
Speech restoration of real data is highly challenging.

Paired training data are not available.

Cannot use information on acoustics distortions (e.g., audio devices).

Low-quality old recordings

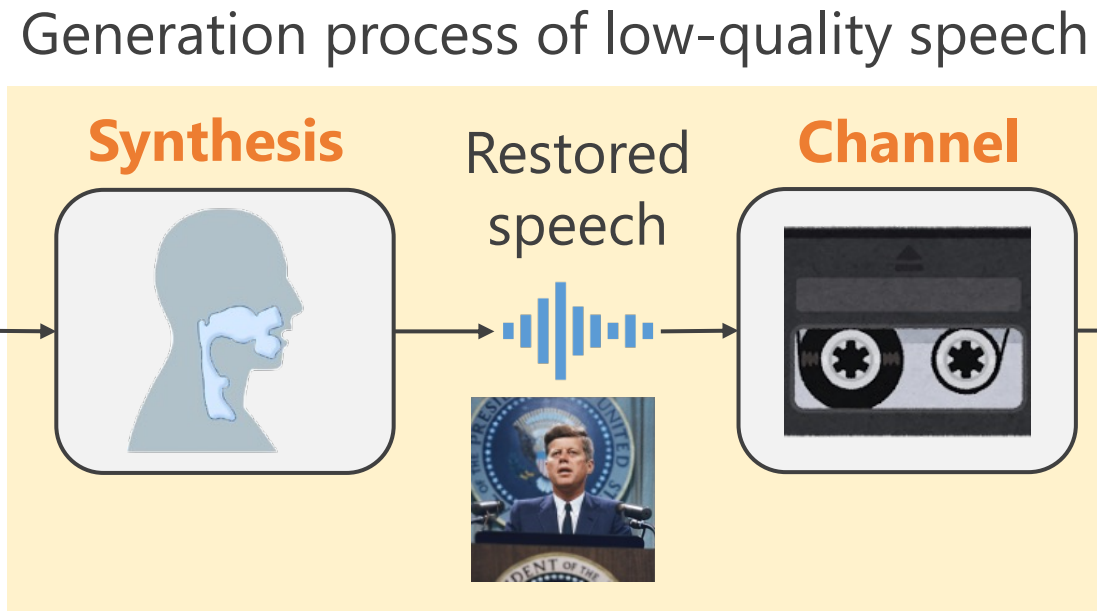
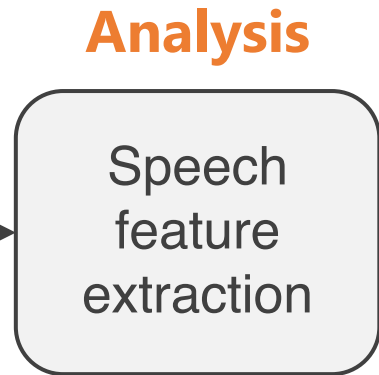
High-quality restored audio



Learning speech restoration model **without paired data**.

**Simulating the generation process** of recorded audio.

Low-quality  
speech



Reconstructed  
speech



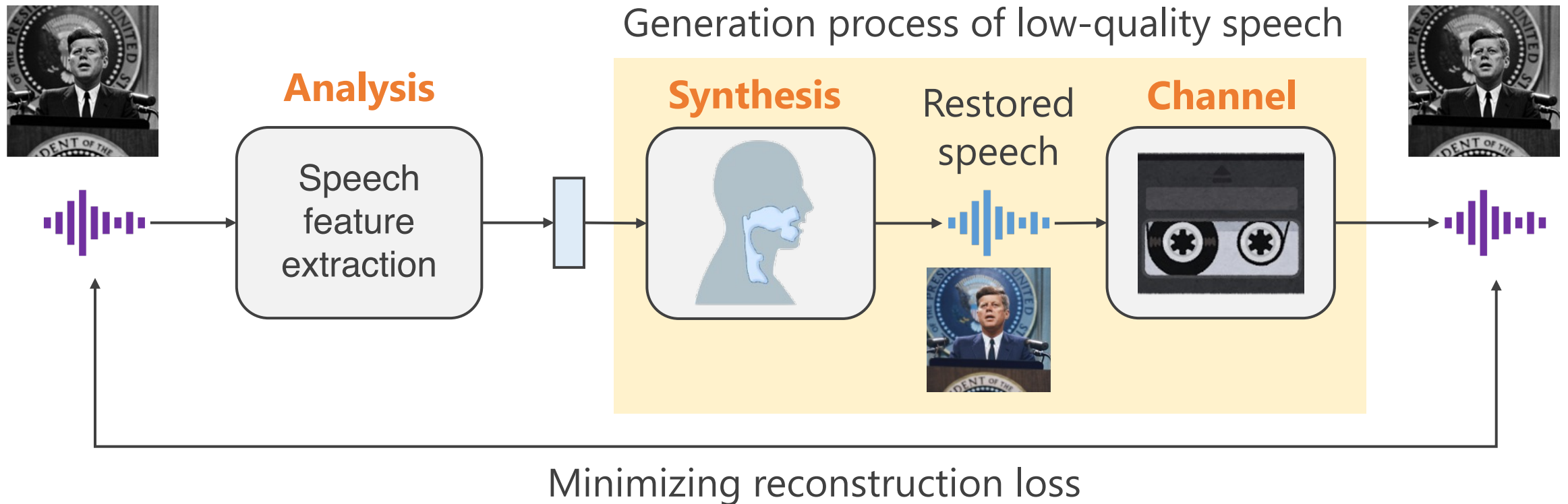
Minimizing reconstruction loss

Learning speech restoration model **without paired data**.  
**Simulating the generation process** of recorded audio.

Generation process of low-quality speech



Learning speech restoration model **without paired data**.  
Consisting of **analysis**, **synthesis**, and **channel** modules.



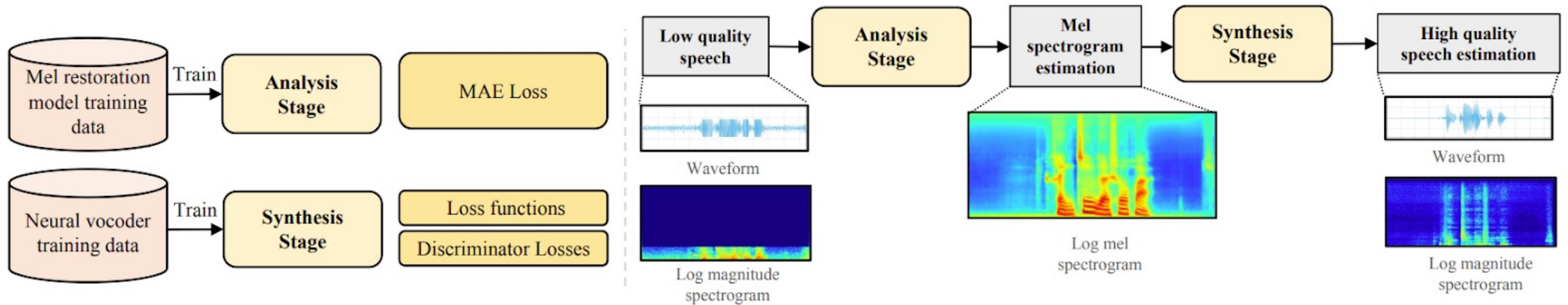


- Background
- Related Work
- Methods
- Experimental Evaluation
- Summary

**Supervised learning** for speech restoration [Liu+, 2021]

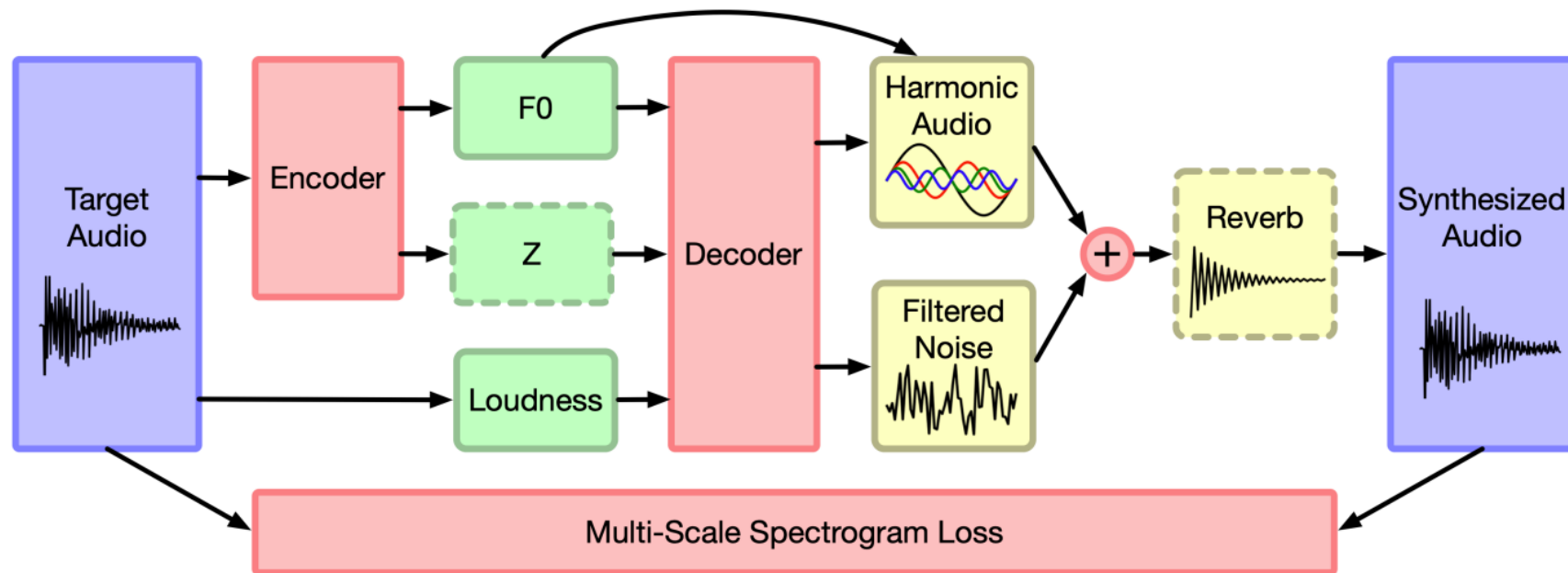
Training **analysis** and **synthesis** modules separately.

Creating artificial **paired** training data.



Our approach uses **real data** based on self-supervised learning.

Differential digital signal processing (DDSP) autoencoder [Engel+, 2021]  
Learning disentangled audio features in a **self-supervised** manner



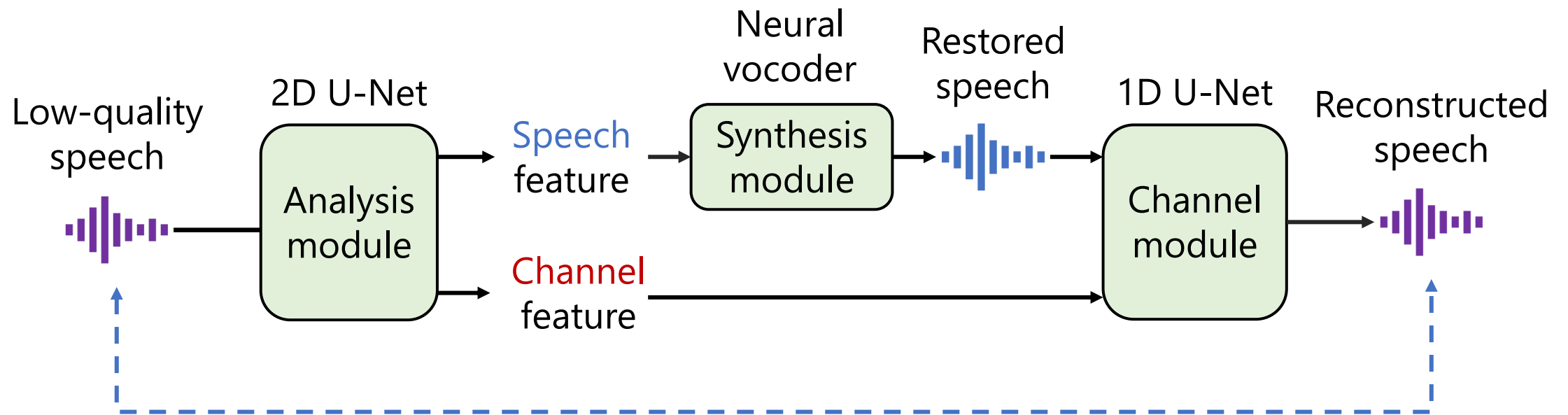
Our work focuses on **restoration of degraded speech**.

- Background
- Related Work
- Methods**
- Experimental Evaluation
- Summary

Analysis, Synthesis, Channel modules are all composed of neural networks.

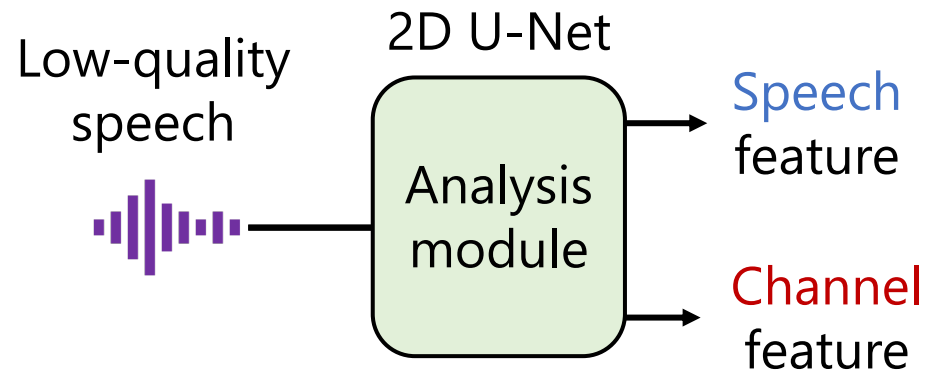
**Analysis** and **Channel** modules: 2D and 1D U-Net models

**Synthesis** module: HiFi-GAN [Kong+, 2020]

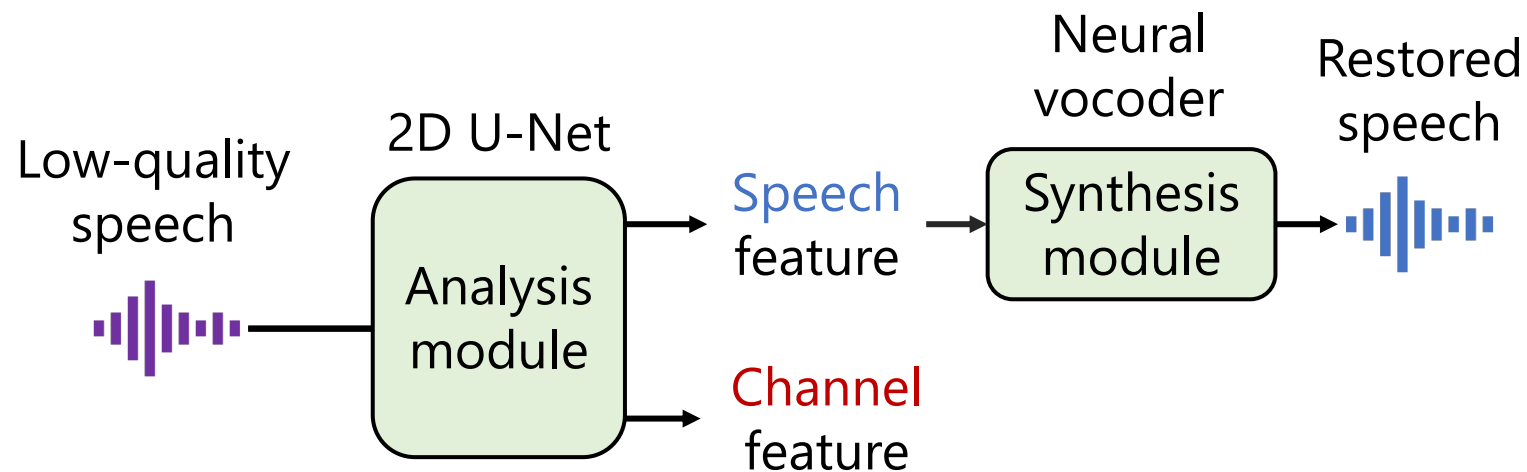


Multi-scale spectral loss  $\mathcal{L}_{\text{recons}} = \sum_i \{ \|\mathbf{s}_i - \hat{\mathbf{s}}_i\|_1 + \alpha \|\log \mathbf{s}_i - \log \hat{\mathbf{s}}_i\|_1 \}$

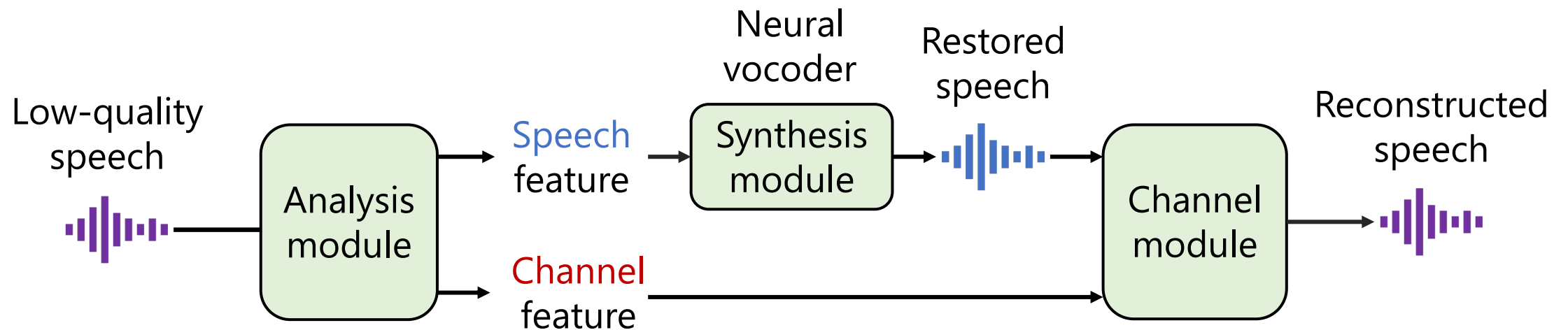
**Analysis** module estimates **speech** features and **channel** features.



**Synthesis** module synthesizes **restored speech** from speech features.



**Channel** module adds **channel** features to restored speech.

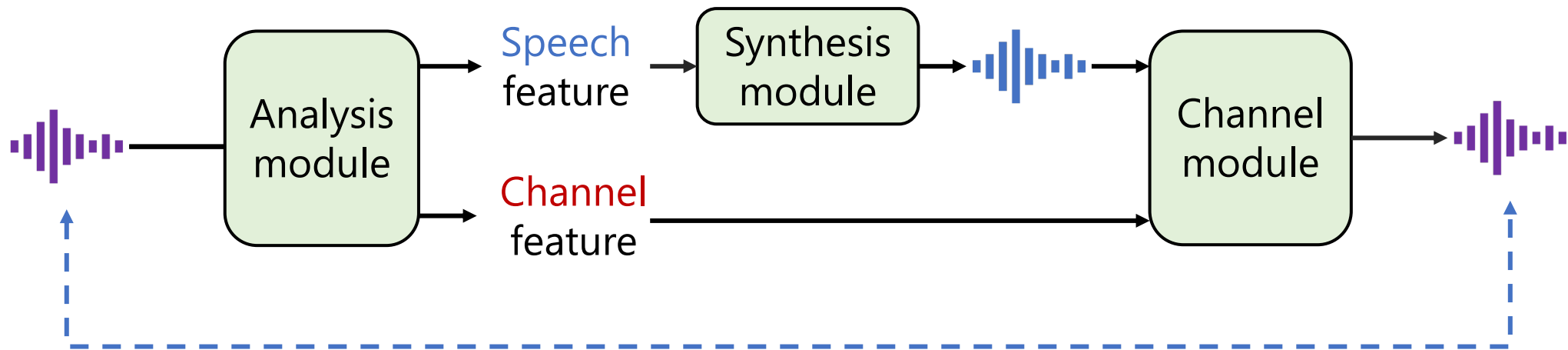




**Analysis** module estimates **speech** features and **channel** features.

**Synthesis** module synthesizes **restored speech** from speech features.

**Channel** module adds **channel** features to restored speech.

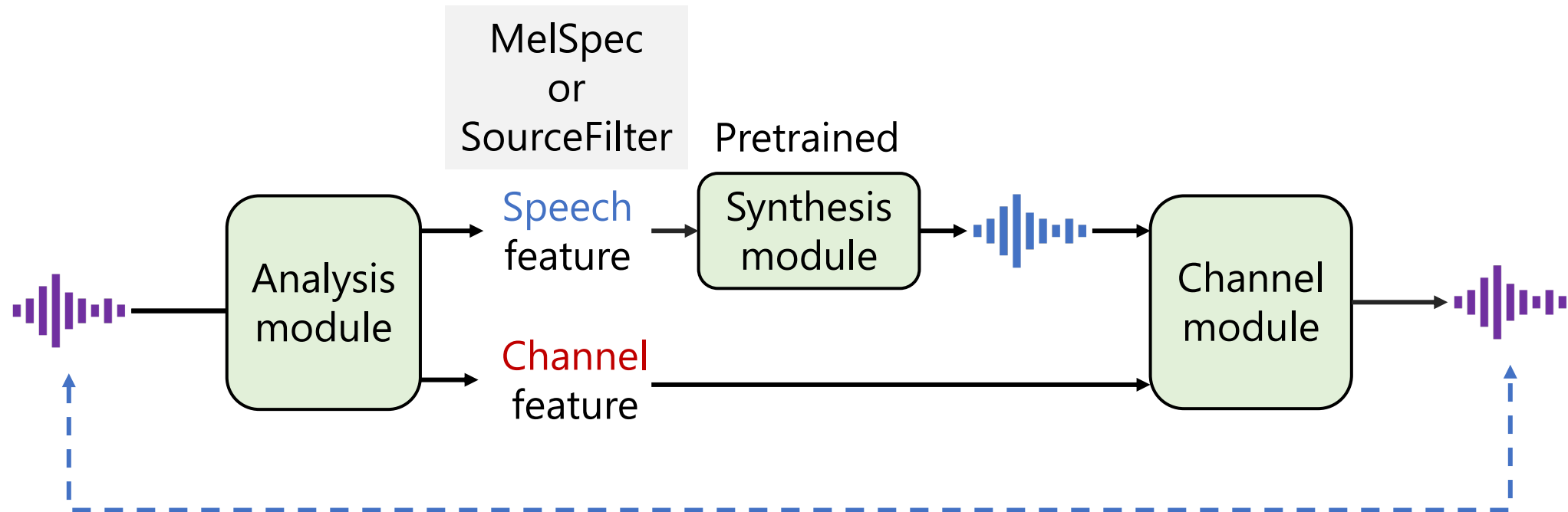


Multi-scale spectral loss  $\mathcal{L}_{\text{recons}} = \sum_i \{ \|\mathbf{s}_i - \hat{\mathbf{s}}_i\|_1 + \alpha \|\log \mathbf{s}_i - \log \hat{\mathbf{s}}_i\|_1 \}$

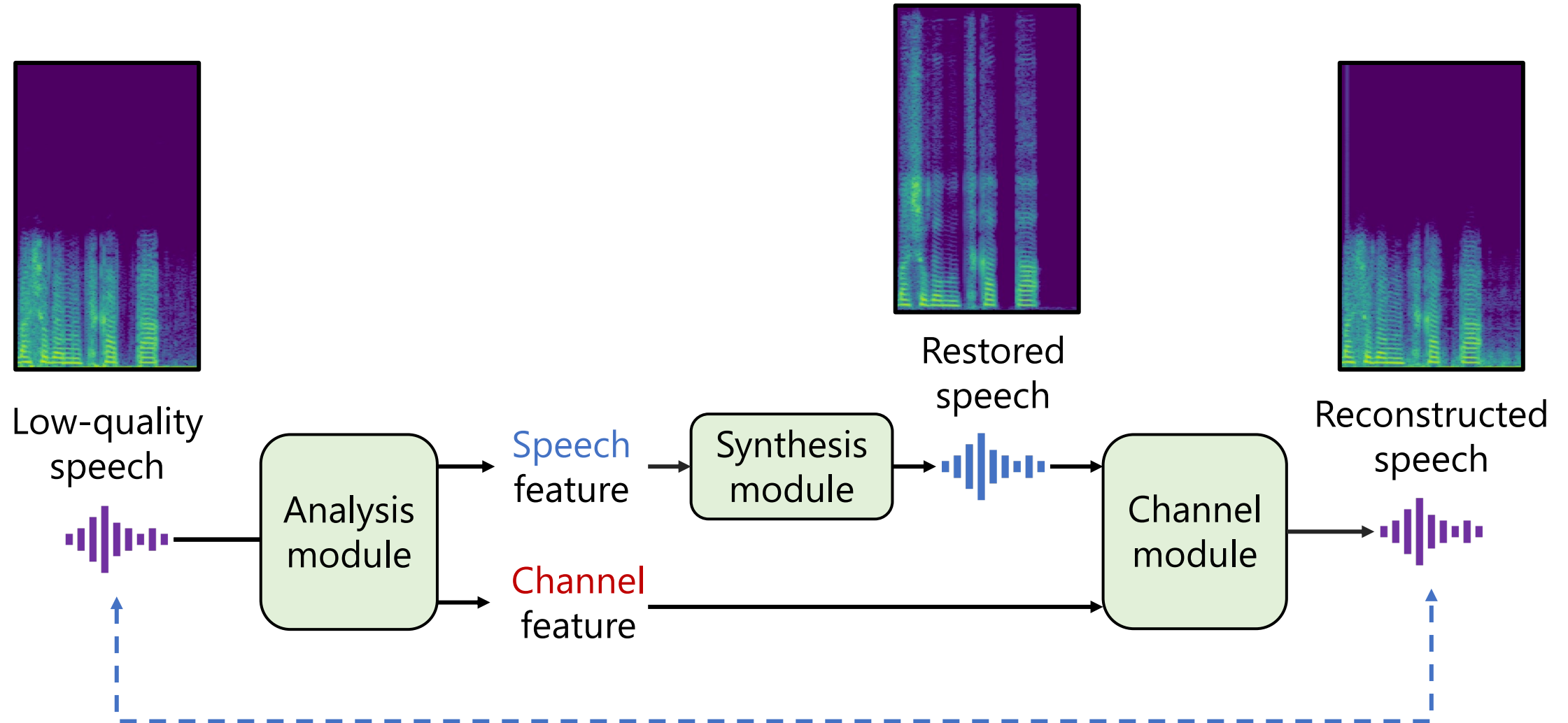
**MelSpec:** Using mel spectrogram to train synthesis module

**SourceFilter:** Using mel cepstrum + F0 to train synthesis module

Only pretraining synthesis module and freezing it.

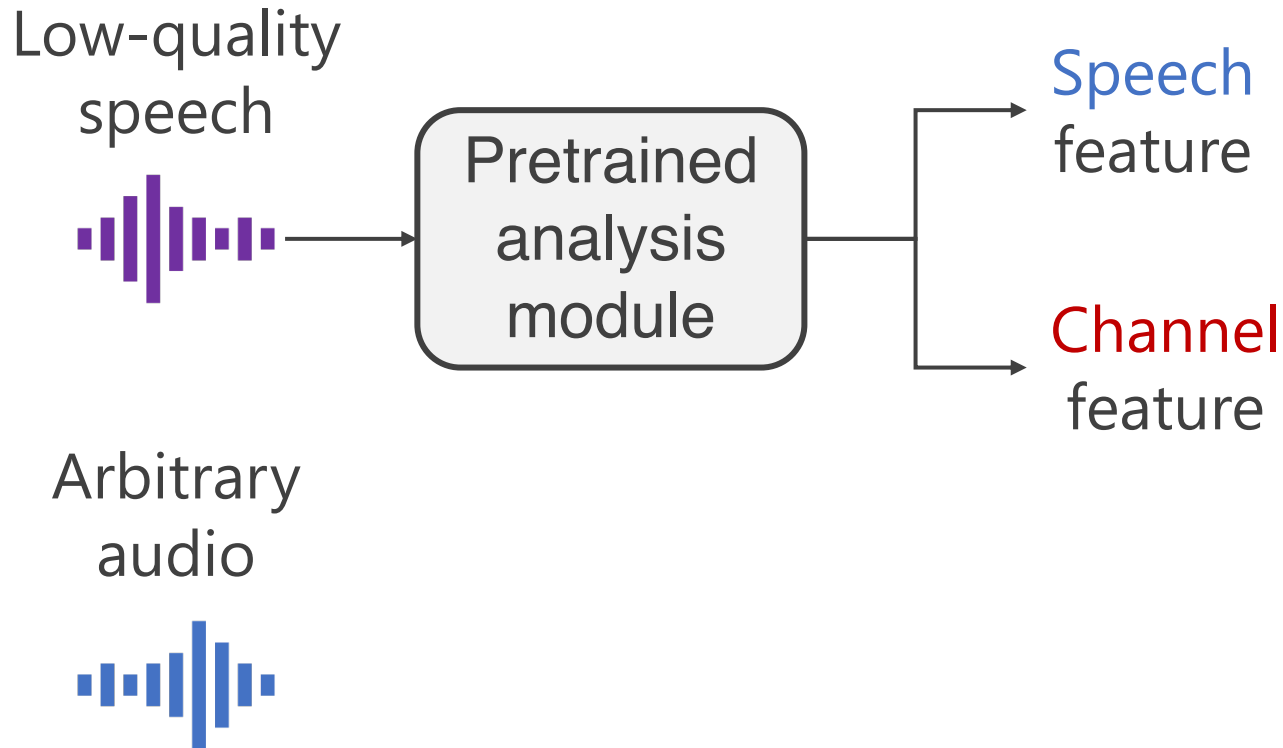


Multi-scale spectral loss  $\mathcal{L}_{\text{recons}} = \sum_i \{ \|\mathbf{s}_i - \hat{\mathbf{s}}_i\|_1 + \alpha \|\log \mathbf{s}_i - \log \hat{\mathbf{s}}_i\|_1 \}$

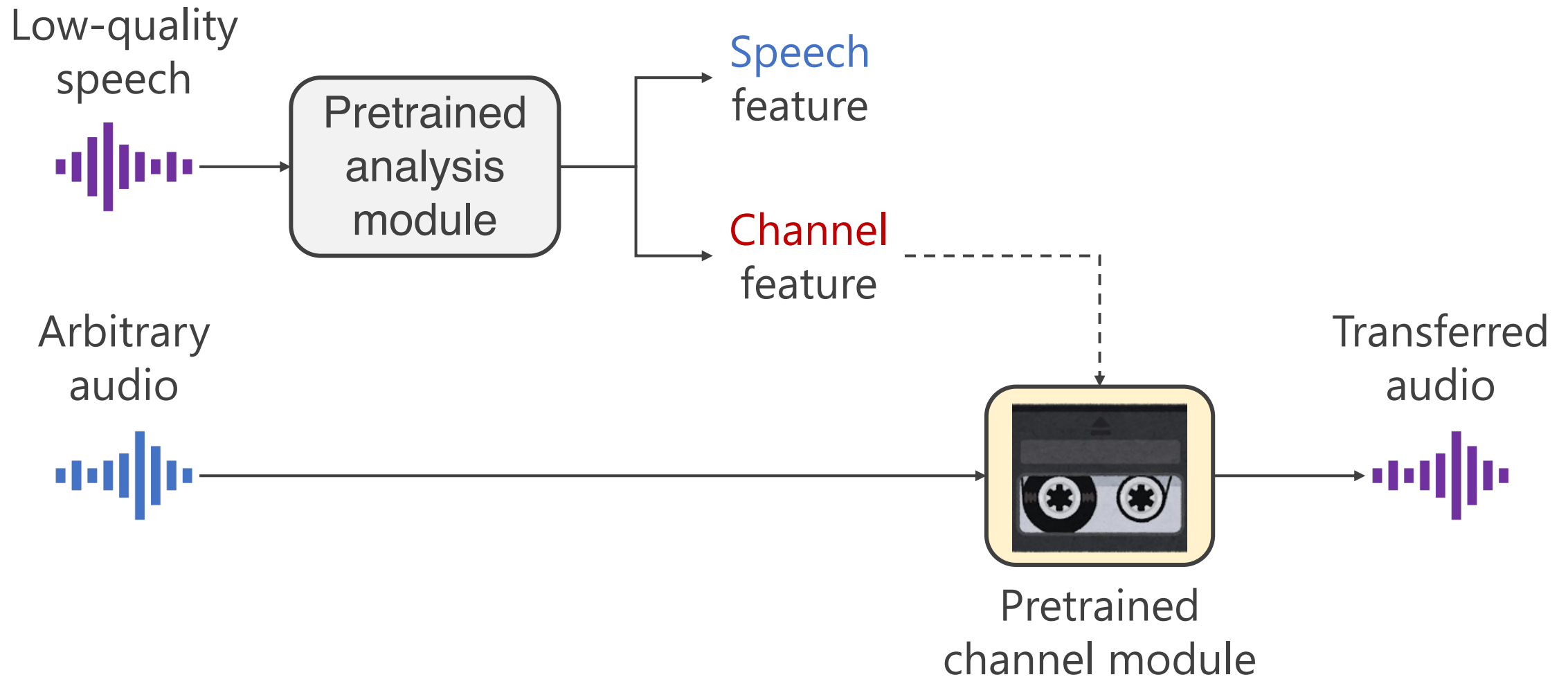


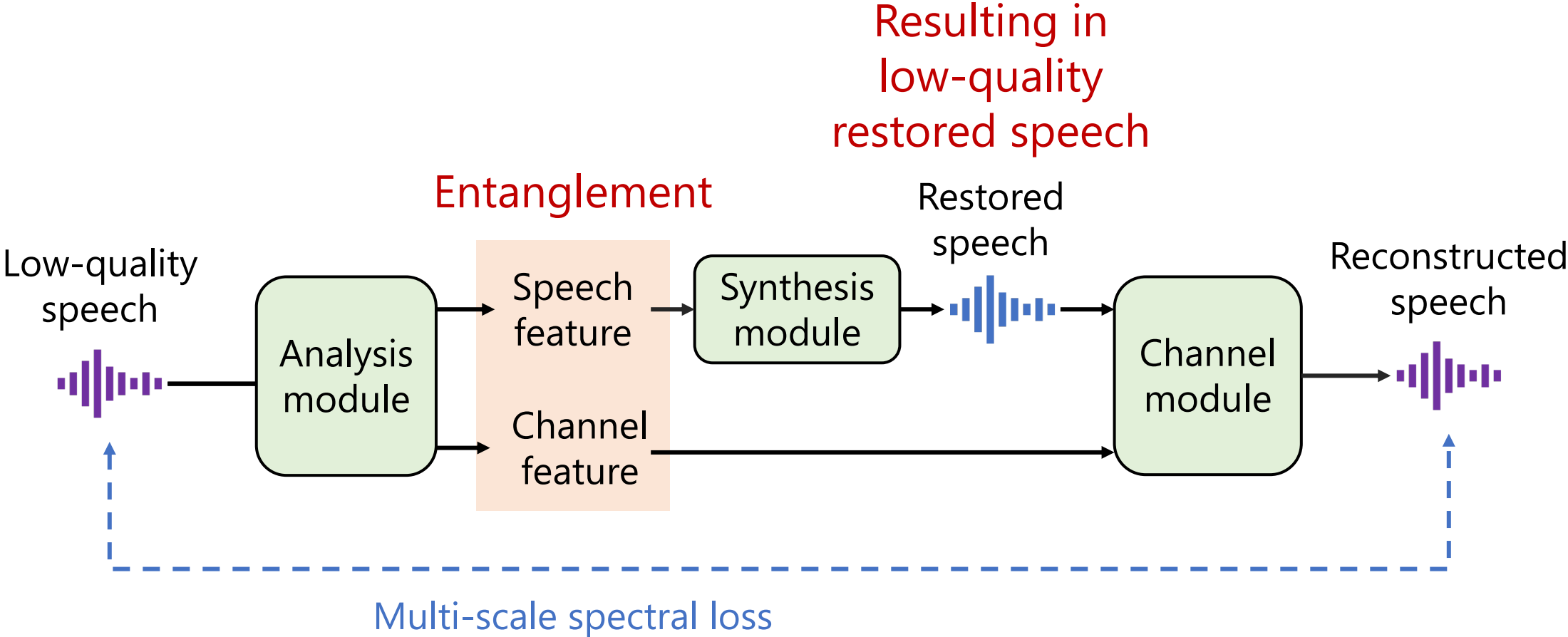
Multi-scale spectral loss  $\mathcal{L}_{\text{recons}} = \sum_i \{ \|\mathbf{s}_i - \hat{\mathbf{s}}_i\|_1 + \alpha \|\log \mathbf{s}_i - \log \hat{\mathbf{s}}_i\|_1 \}$

Proposed method **works as Audio effector** to extract and add channel features.

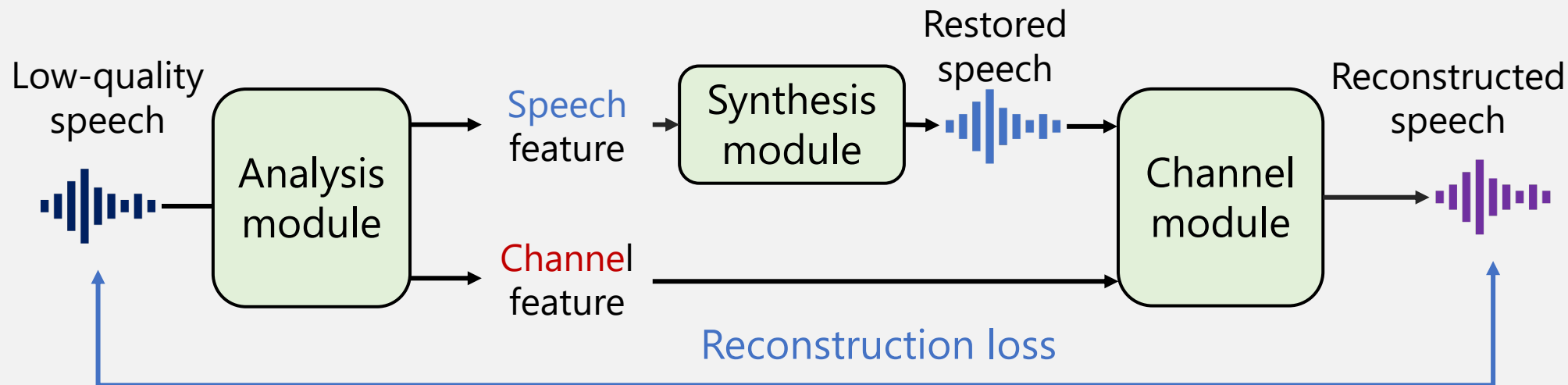


Proposed method **works as Audio effector** to extract and add channel features.



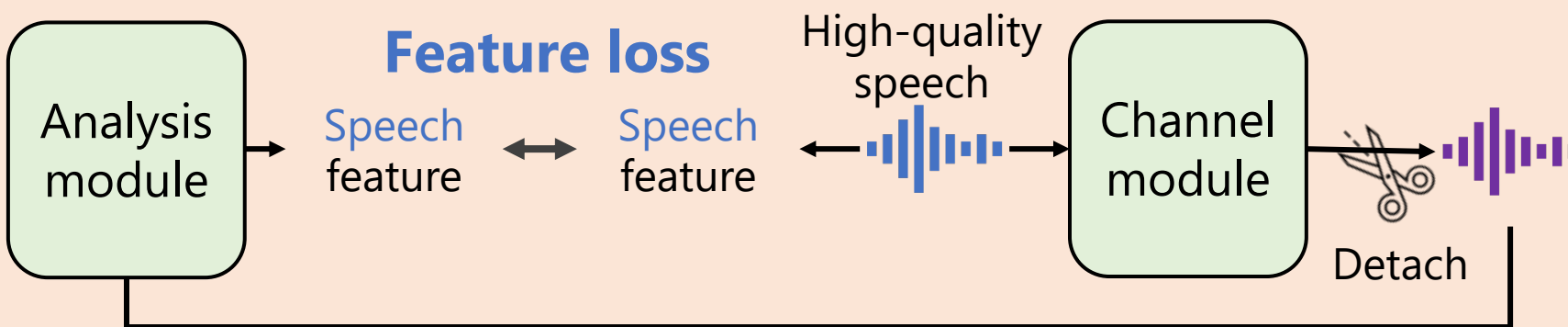


## Self-supervised learning with only degraded speech



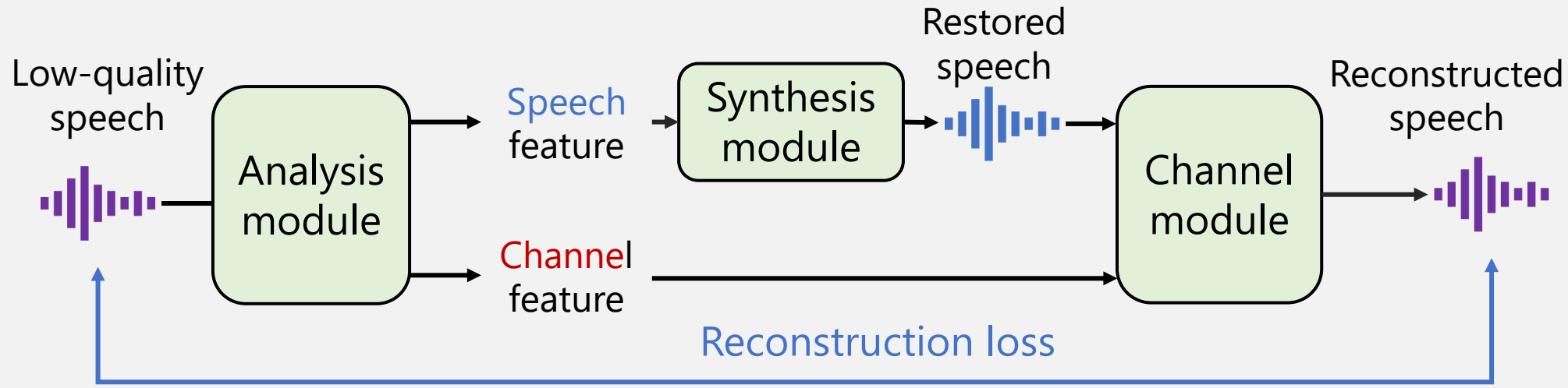
$\tau_{\text{forward}}$   
**Learning channel module** to adapt to input speech

## Backward training with arbitrary high-quality speech



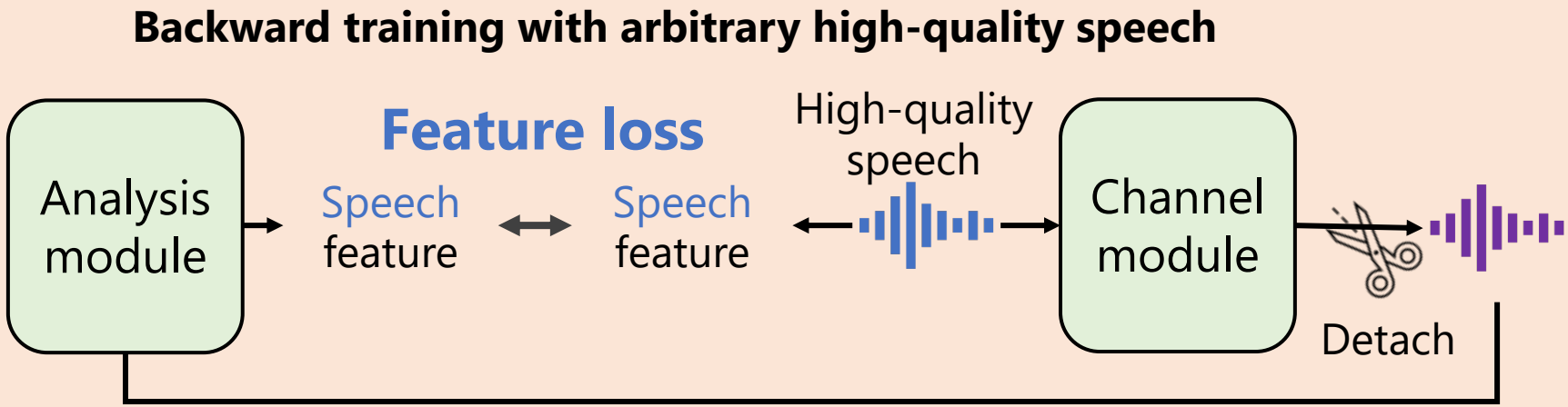
$\tau_{\text{backward}}$   
**Learning analysis module** to output clean speech features

## Self-supervised learning with only degraded speech



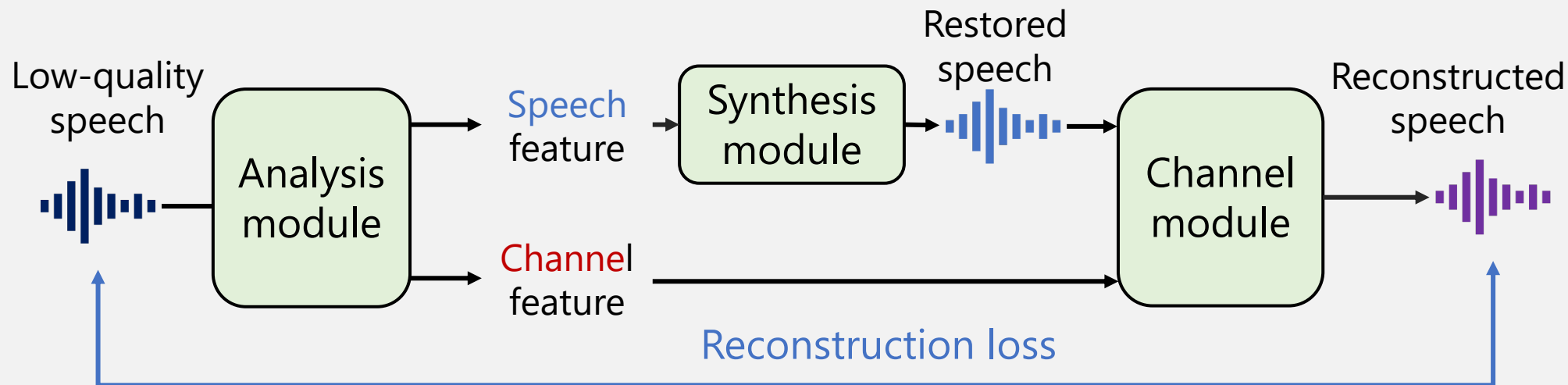
$\tau_{\text{forward}}$   
**Learning channel module** to adapt to input speech





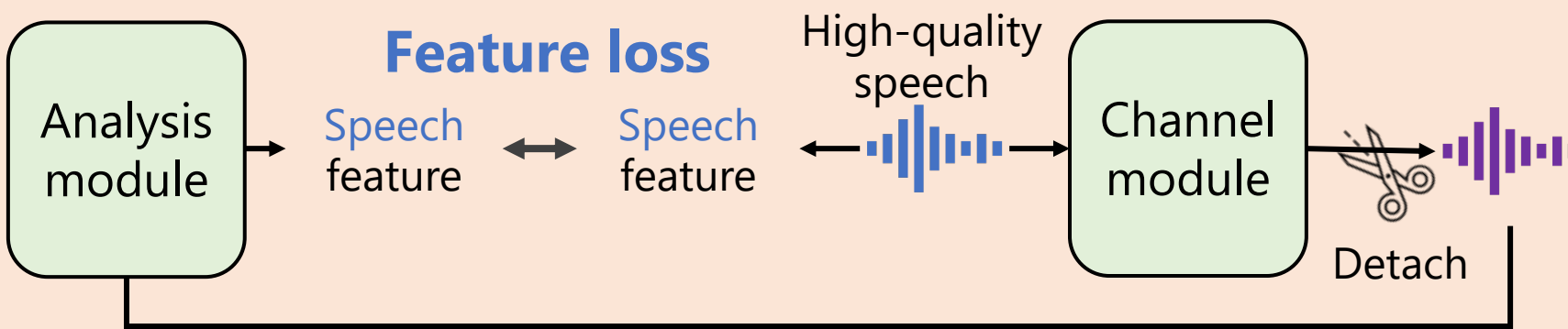
$\tau_{\text{backward}}$   
**Learning analysis module**  
to output clean speech features

## Self-supervised learning with only degraded speech



$\tau_{\text{forward}}$   
**Learning channel module** to adapt to input speech

## Backward training with arbitrary high-quality speech



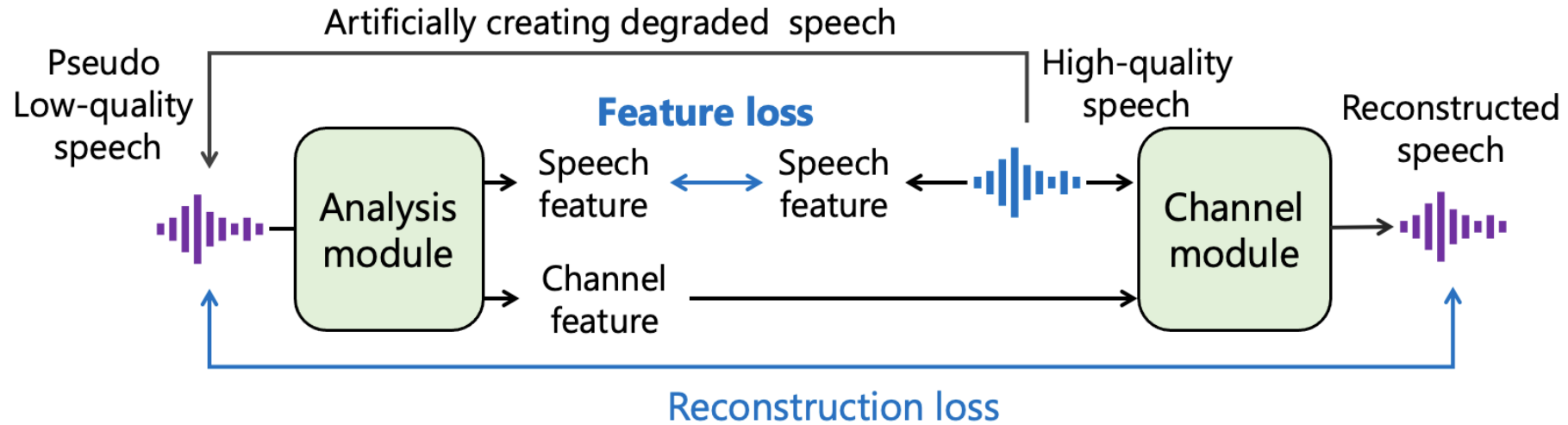
$\tau_{\text{backward}}$   
**Learning analysis module** to output clean speech features

We cannot get so much data from a single historical audio material.

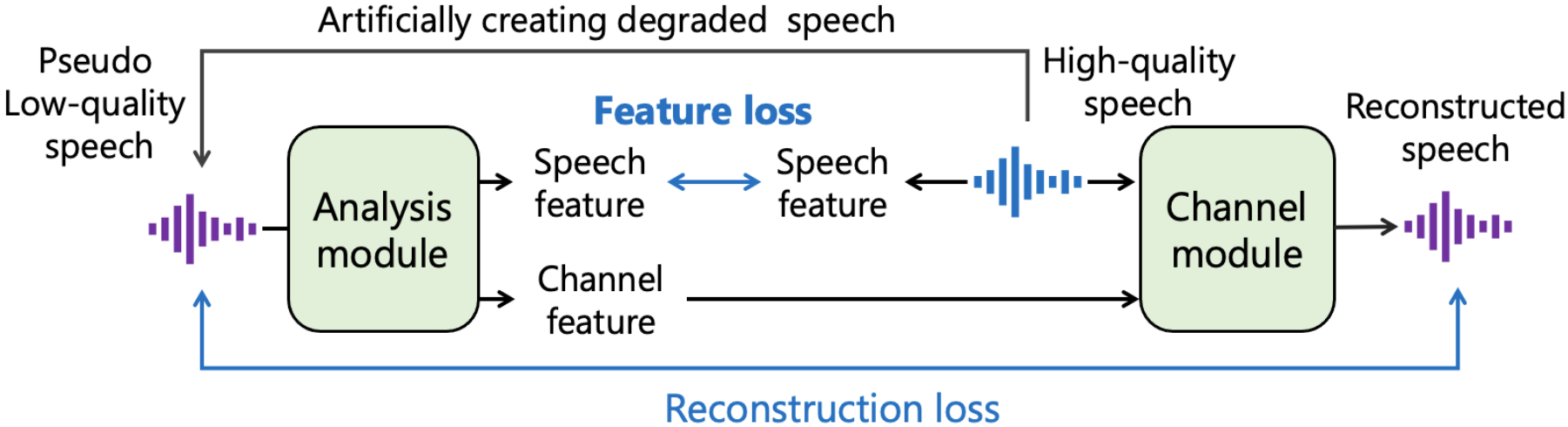
Hard to learn modules with low-resource data (< 1 hour).

➤ Introducing supervised pretraining to tackle data scarcity.

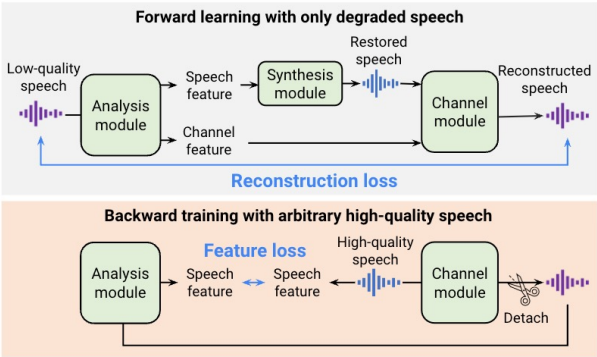
## Supervised pretraining with pseudo low-quality speech data



## Supervised pretraining with pseudo low-quality speech data



## Self-supervised learning with real data



- Background
- Related Work
- Methods
- Experimental Evaluation
- Summary

Compared our method and previous supervised method [Liu+, 2021].

1) **Simulated datasets** based on high-quality speech corpus [Takamichi+, 2021]

Applied four types of distortions to **6-hour** single-speaker data

- a) Band-limited
- b) Clipped
- c) Quantized & Resampled
- d) Overdrive

Compared our method and previous supervised method [Liu+, 2021].

1) **Simulated datasets** based on high-quality speech corpus [Takamichi+, 2021]

Applied four types of distortions to **6-hour** single-speaker data





















- a) Band-limited
- b) Clipped
- c) Quantized & Resampled
- d) Overdrive

2) **Real historical audio material** recorded on an analog tape recorder

Around **20 minutes'** multi-speaker data recorded in 1960s – 1970s







**Mean opinion score (MOS) test** of speech quality with 40 listeners in each case

	(a) Band-limited		(b) Clipped		(c) Quantized & Resampled		(d) Overdrive	
	MOS	Sample	MOS	Sample	MOS	Sample	MOS	Sample
Ground-truth	4.51		4.58		4.67		4.65	
Input	2.38		2.45		1.73		1.54	
<b>Supervised</b> [Liu+, 2021]	3.74		3.01		2.80		2.00	
<b>Proposed</b> (MelSpec)	<b>4.20</b>		<b>3.49</b>		<b>3.27</b>		<b>2.68</b>	
<b>Proposed</b> (SourceFilter)	3.46		2.49		2.66		2.58	

Proposed method achieved significantly higher MOS than previous supervised method.





Evaluated proposed method with **real historical audio** (around 20 min).

	MOS	Sample
Input	2.98	
<b>Supervised</b> [Liu+, 2021]	2.80	
Proposed (MelSpec)	2.96	
Proposed + pretraining (MelSpec)	<b>3.06</b>	

Evaluated proposed method with **real historical audio** (around 20 min).

Statistical significance  
(p-value < 0.05)  
in side-by-side test

Effectiveness  
for real data

	MOS	Sample
Input	2.98	
<b>Supervised</b> [Liu+, 2021]	2.80	
Proposed (MelSpec)	2.96	
Proposed + pretraining (MelSpec)	<b>3.06</b>	





Still needs  
pretraining  
for real data

Similarity MOS (SMOS) test for **similarity of audio characteristics**.

- **Source**: Original high-quality audio samples
- **Mean spec. diff**: Applying differential spectrum to original audio samples
- **Proposed**: Performing audio effect transfer with our method

Similarity MOS (SMOS) test for **similarity of audio characteristics**.

- **Source**: Original high-quality audio samples
- **Mean spec. diff**: Applying differential spectrum to original audio samples
- **Proposed**: Performing audio effect transfer with our method

	Simulated (Quantized & Resampled)	Real
Target	3.98 	2.99
Source	1.16 	1.30
<b>Mean spec. diff</b>	1.68 	-
<b>Proposed</b>	<b>3.44</b> 	<b>2.12</b>

**Self-supervised speech restoration** **without paired data**

Confirmed effectiveness with **real data** but need more data.

**Code**



**More audio samples**

