



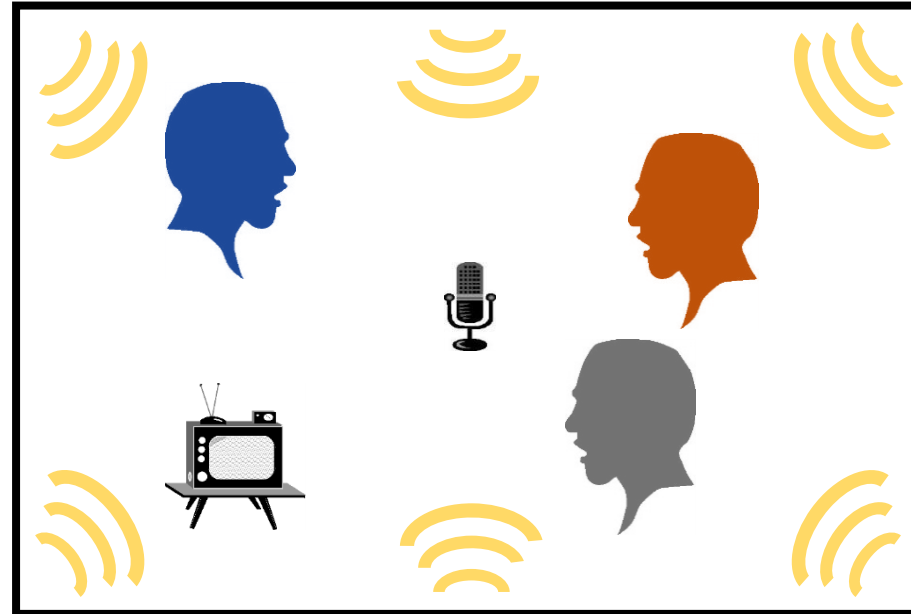
# ALL-NEURAL ONLINE SOURCE SEPARATION, COUNTING, AND DIARIZATION FOR MEETING ANALYSIS

**Thilo von Neumann**<sup>1,2</sup>, Keisuke Kinoshita<sup>1</sup>, Marc Delcroix<sup>1</sup>, Shoko Araki<sup>1</sup>,  
Tomohiro Nakatani<sup>1</sup>, Reinhold Haeb-Umbach<sup>2</sup>

<sup>1</sup> NTT Communication Science Laboratories, NTT Corporation, Kyoto, Japan

<sup>2</sup> Paderborn University, Department of Communications Engineering, Paderborn, Germany

# Tackling Problems in Meeting Scenarios



## Problems:

- multiple sources
- number of sources not known
- long recordings

⇒ source separation

⇒ source count estimation

⇒ blockwise/online processing

# NN-based Source Separation Methods



## Permutation Invariant Training (PIT)

[Kolbaek2017]

- purely neural network based
- model structure depends on number of sources to be estimated

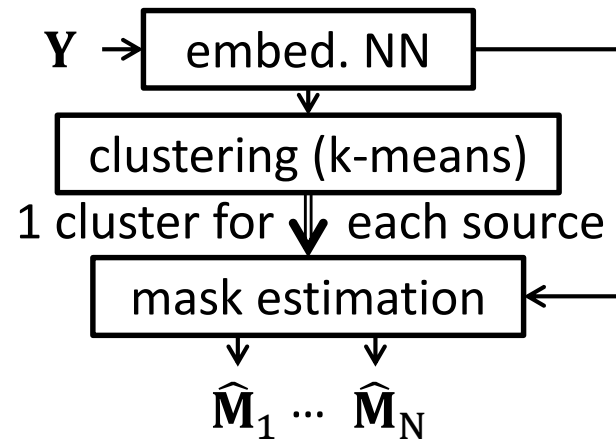


- ~ max number of sources must be known during training
- ~ source counting

## Deep Clustering (DC) / Deep Attractor Networks (DAN)

[Isik2016]/[Chen2017]

- 2 stages: embedding + clustering



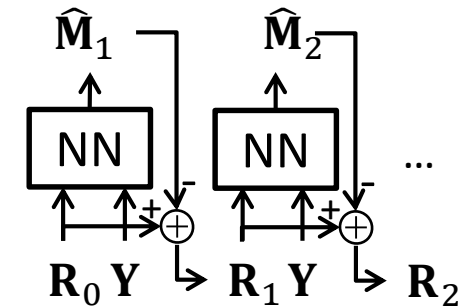
- ~ number of sources must be known or estimated separately
- ~ source counting

✗ block-online processing

## Recurrent Selective Attention Network (RSAN)

[Kinoshita2018]




- purely neural network based
- iterative source extraction



- ✓ source separation for arbitrary number of sources
- ✓ source counting

# Block-Online Processing






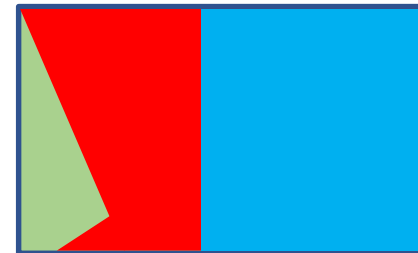
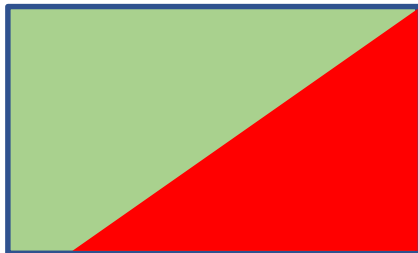
-  source 1
-  source 2
-  source 3



# Block-Online Processing

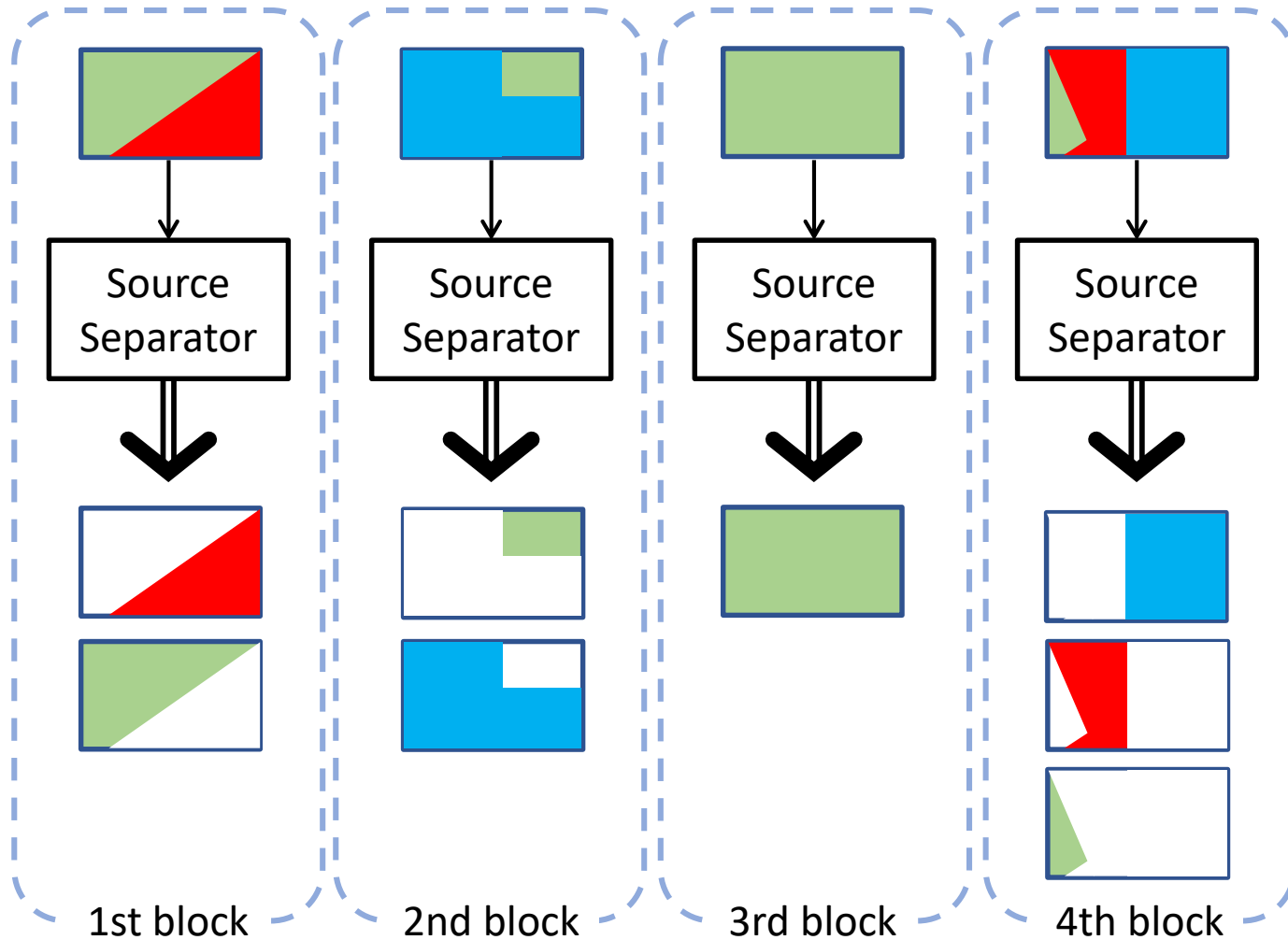


-  source 1
-  source 2
-  source 3



# Difficulties in Block-Online Processing

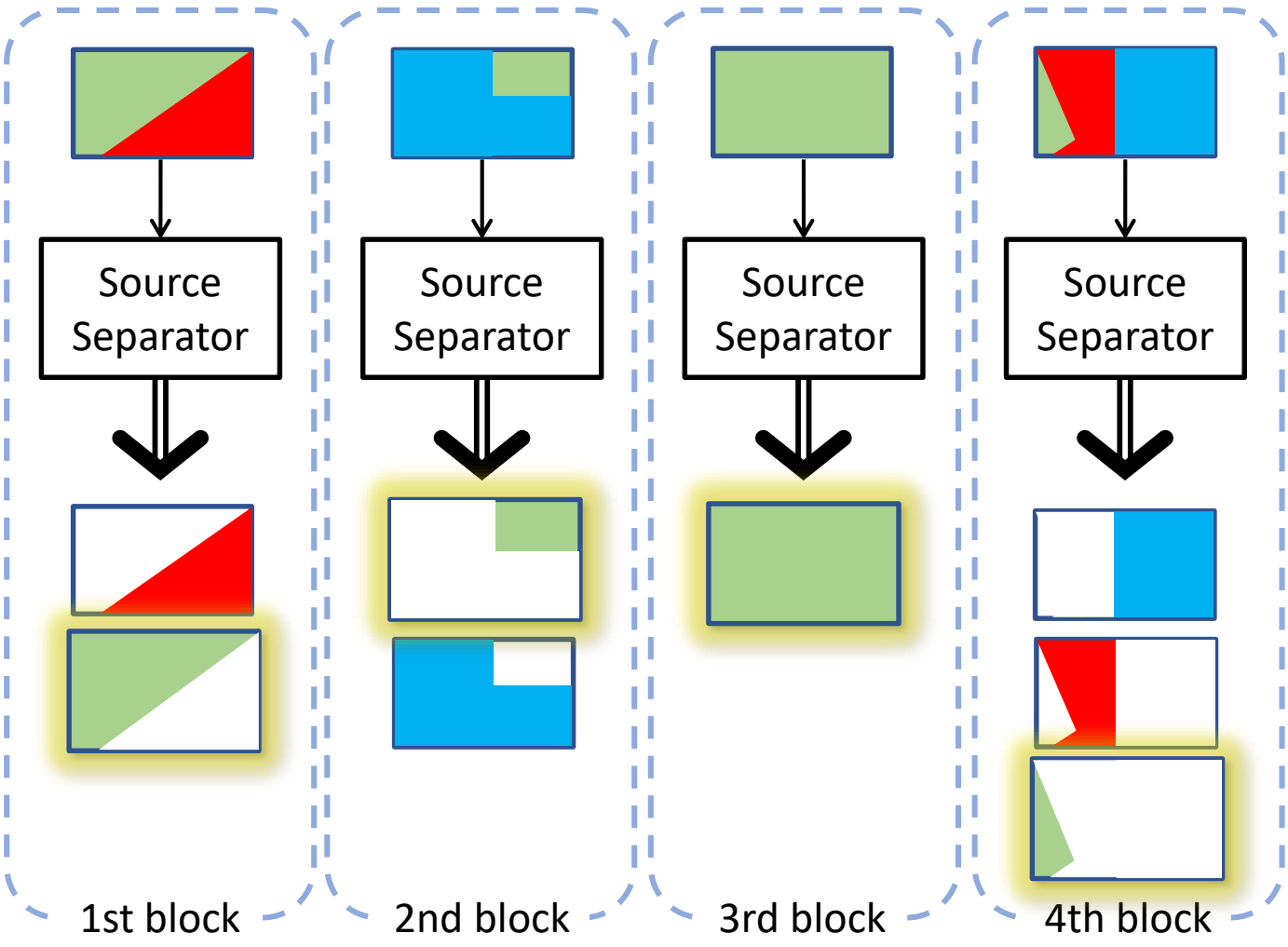
- src 1
- src 2
- src 3



# Difficulties in Block-Online Processing



- src 1
- src 2
- src 3

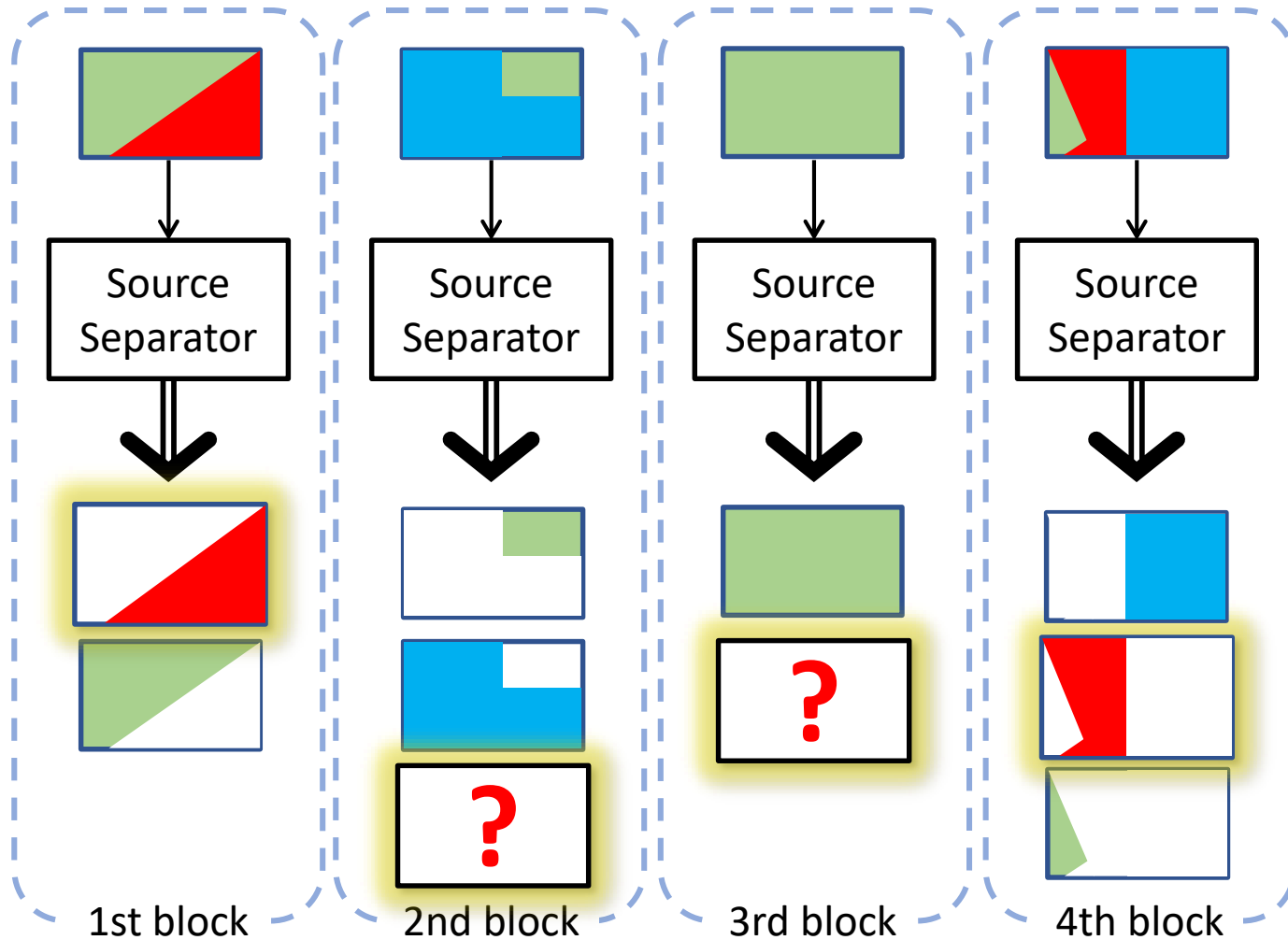


**Block Permutation Problem**

- The output **order** in each block is **unknown**

# Difficulties in Block-Online Processing

- src 1
- src 2
- src 3



## Silent Speakers

- **Notice** silent speakers
- **Remember** silent speakers over gaps

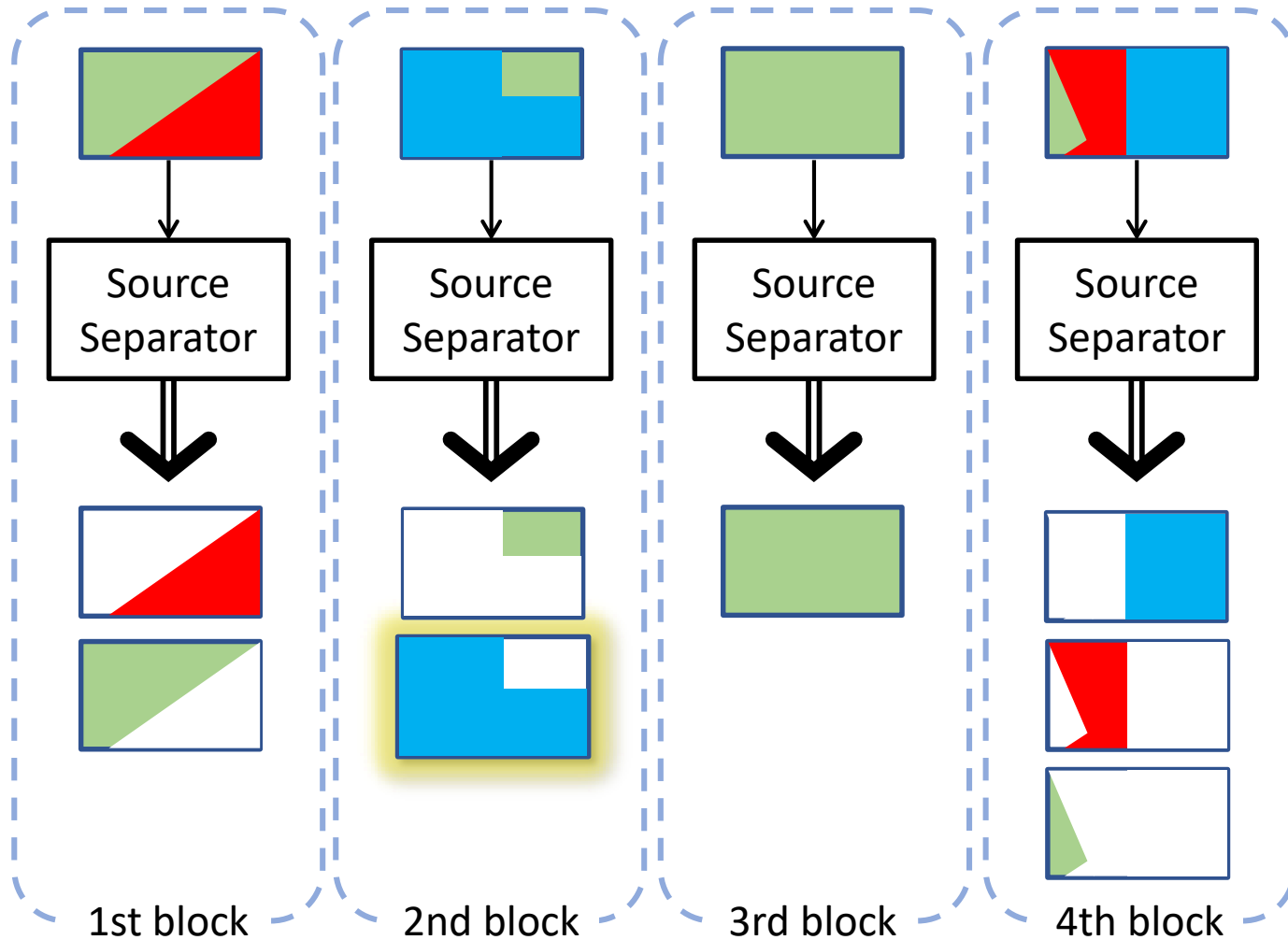
## Block Permutation Problem

- The output **order** in each block is **unknown**



# Difficulties in Block-Online Processing

- src 1
- src 2
- src 3



## New Speakers

- Detect **new speaker in each block**

## Silent Speakers

- **Notice** silent speakers
- **Remember** silent speakers over gaps

## Block Permutation Problem

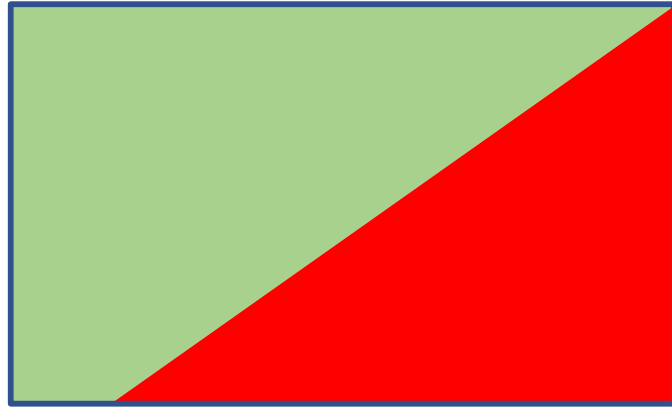
- The output **order** in each block is **unknown**

# Recurrent Selective Attention Network

A fully Neural Network based Source Separation and Source Number Counting  
Approach

K. Kinoshita, L. Drude, M. Delcroix and T. Nakatani. "Listening to Each Speaker One by One with Recurrent Selective Hearing Networks." *ICASSP* (2018).

# Recurrent Selective Attention Network (RSAN)



input spectrogram  
 $|Y|$

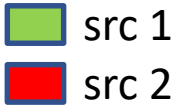


source 1

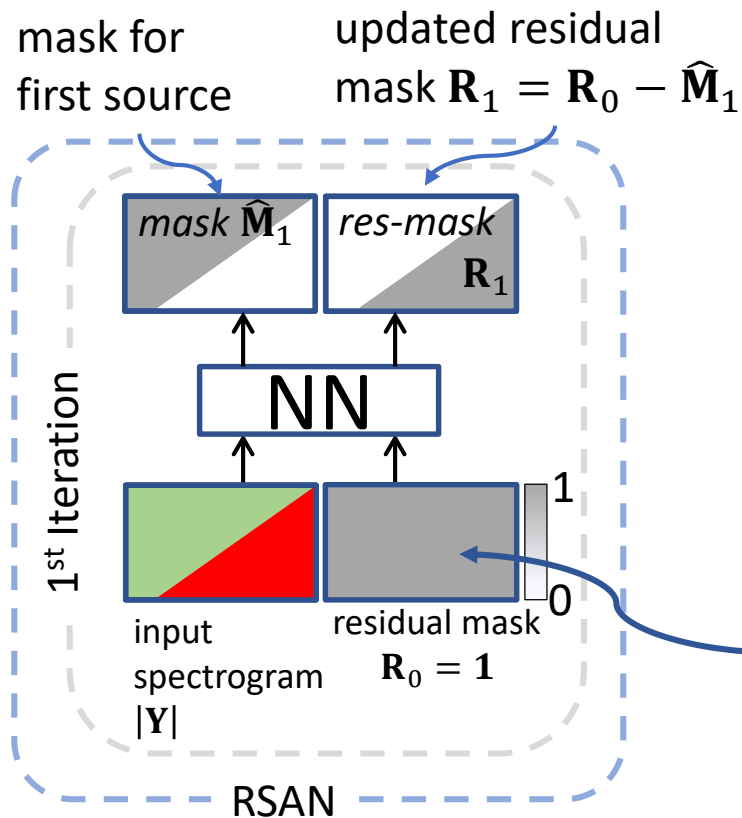


source 2

# Recurrent Selective Attention Network (RSAN)

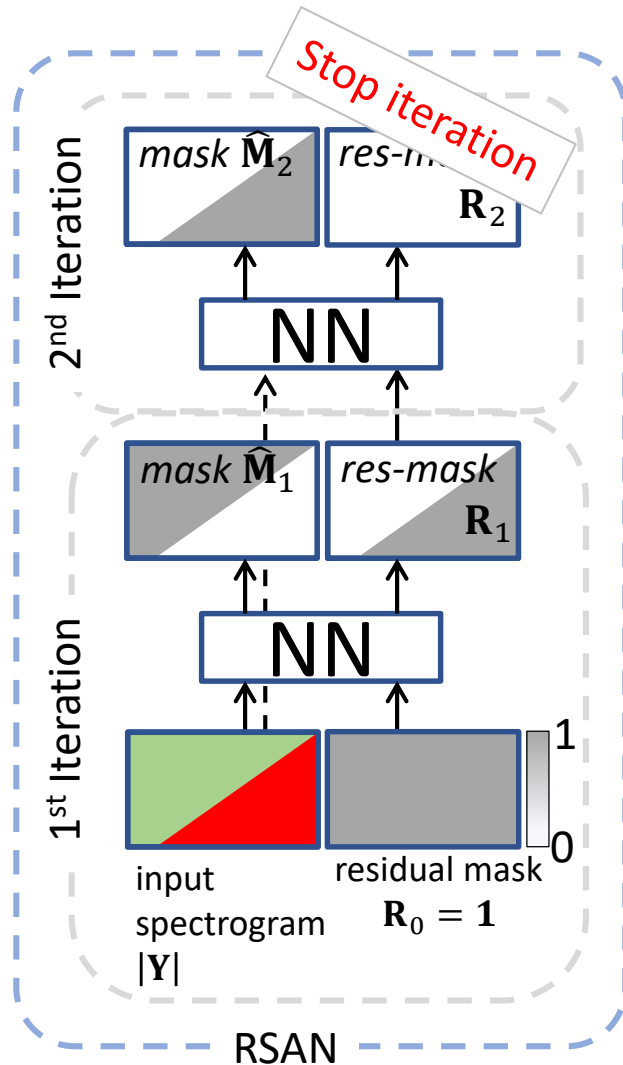
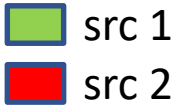


- NN estimates a **mask  $\hat{M}_i$**  for one **source** (which it can **choose by itself**)
- Residual Mask defines region where to look for sources



“Residual”- or “Attention”-mask to guide where to extract sources

# Recurrent Selective Attention Network (RSAN)



- NN estimates a **mask  $\hat{M}_i$**  for one **source** (which it can **choose by itself**)
- Residual Mask defines region where to look for sources
- **Iteratively** repeated until no source is left
  - based on thresholding on residual mask
- Count iterations for **source count estimation**

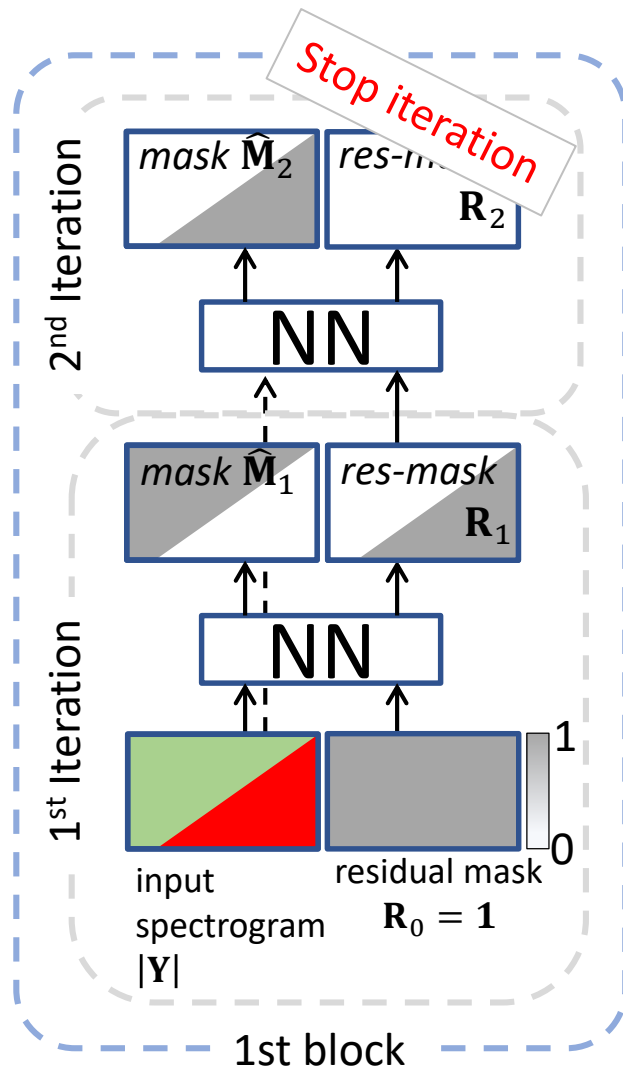
# Proposed Method

How to do block-online processing with RSAN

# Proposed Method



- src 1
- src 2
- src 3

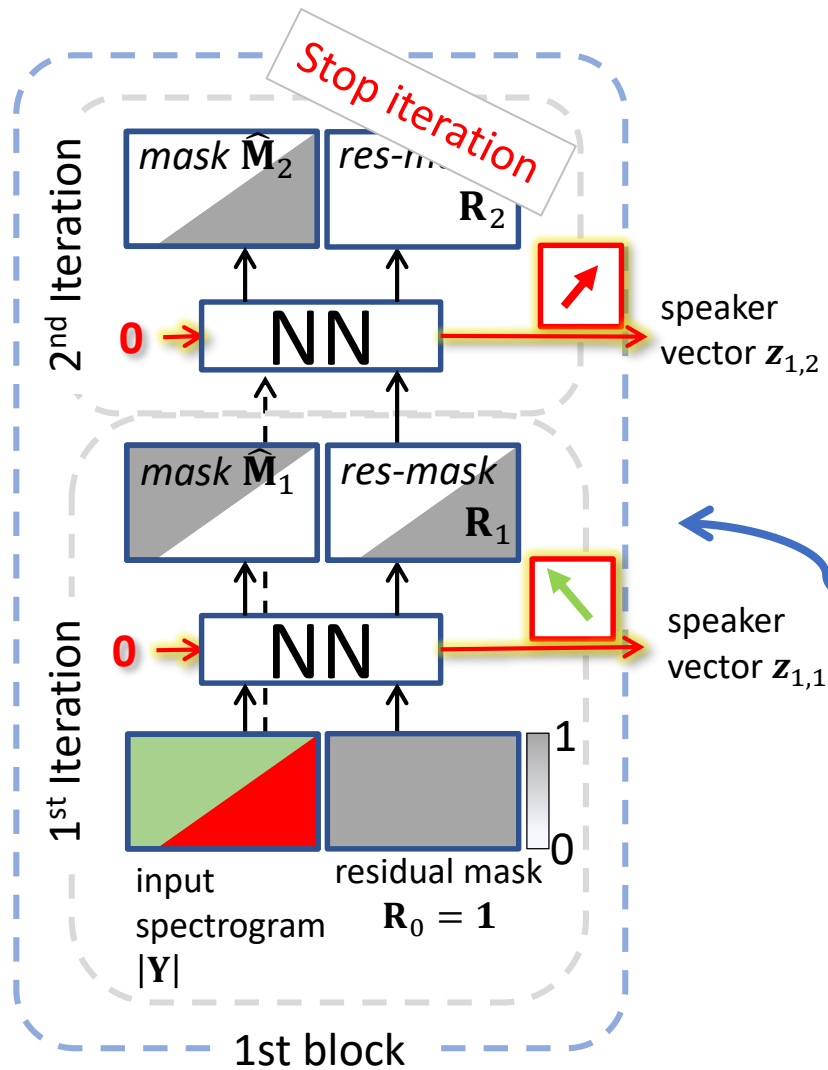


Recurrent Selective Attention Network (RSAN)

# Proposed Method



- src 1
- src 2
- src 3



Recurrent Selective Attention Network (RSAN)

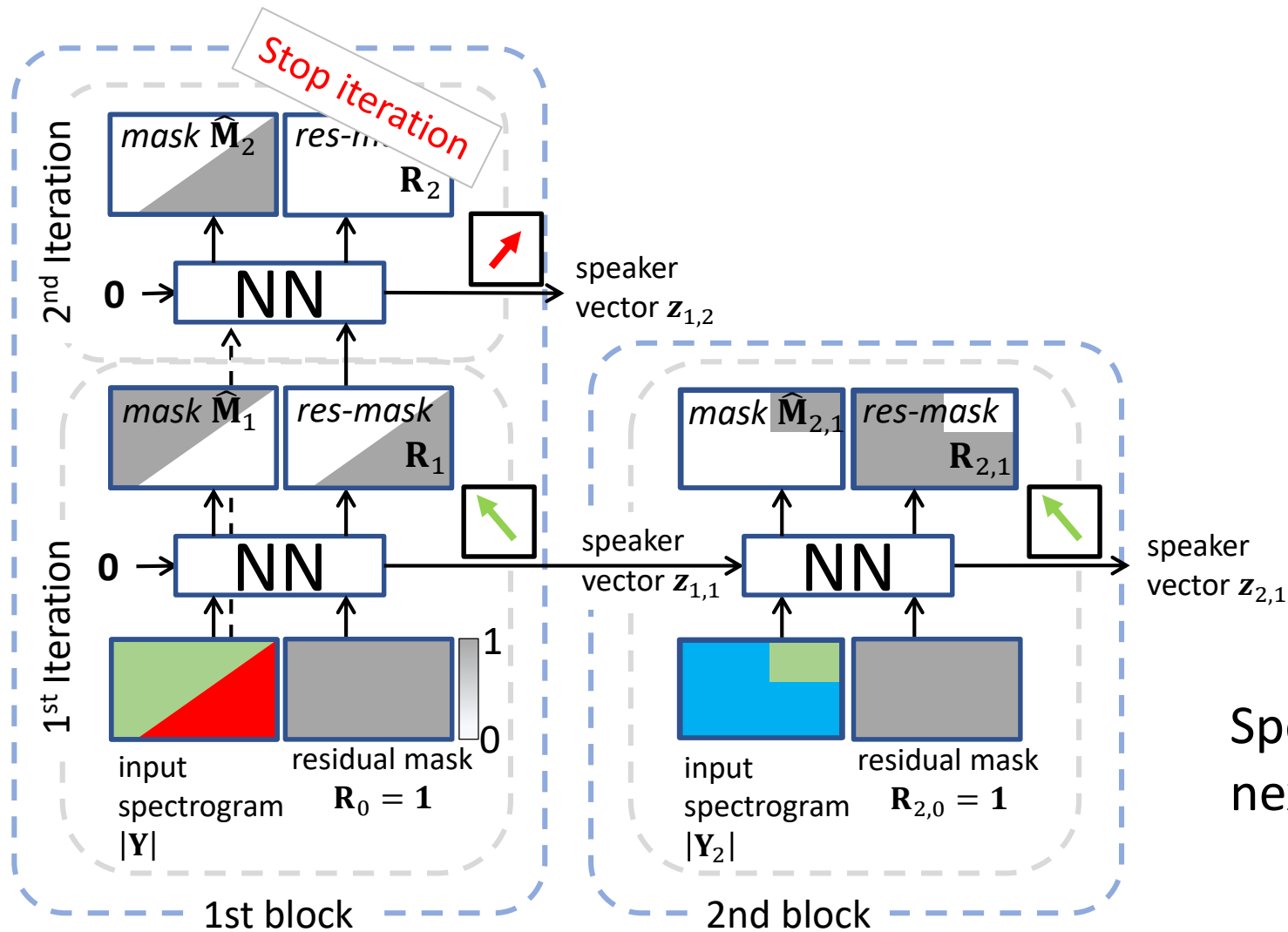
+ speaker embedding input & output



# Proposed Method



- src 1
- src 2
- src 3

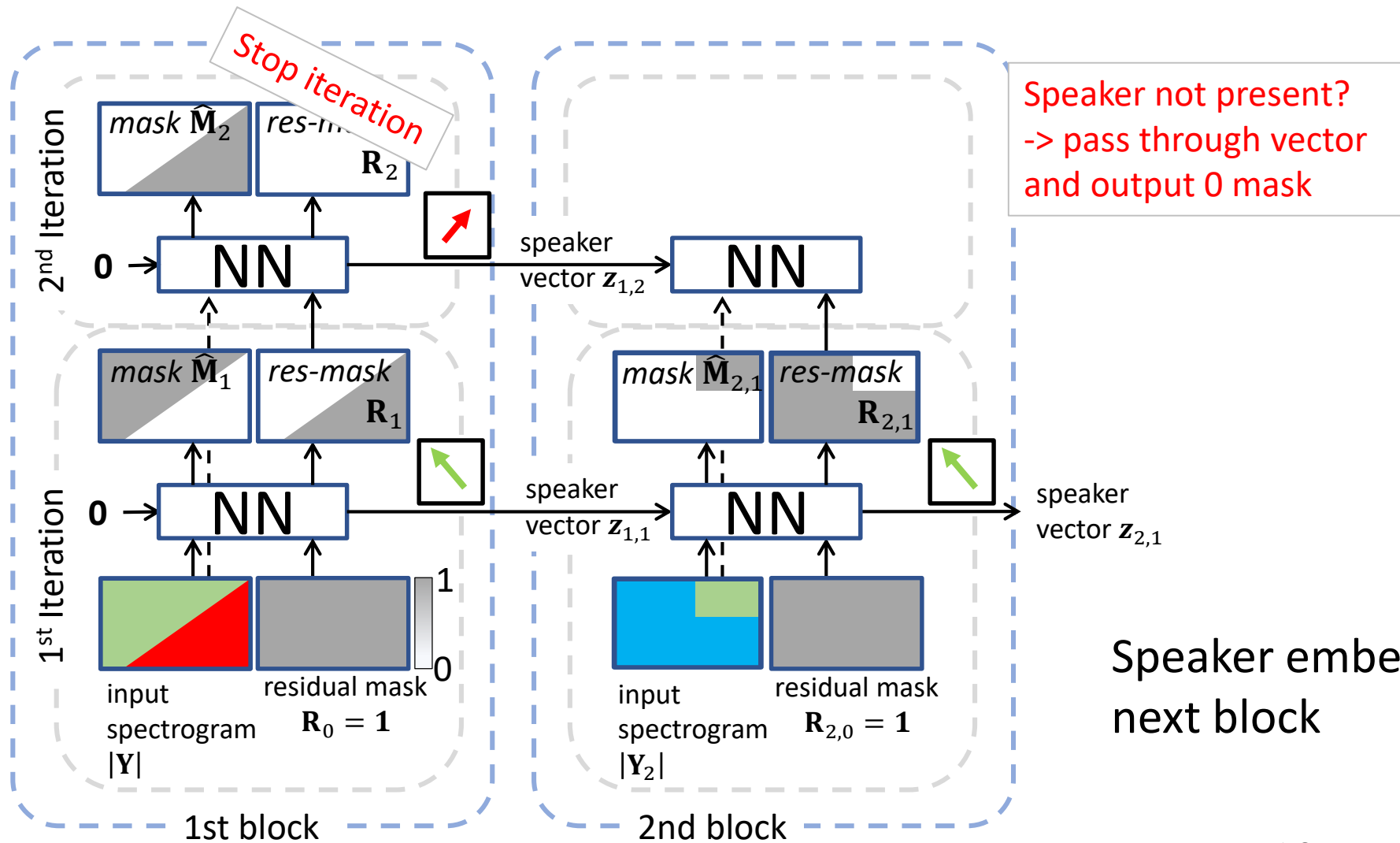


Speaker embedding is passed to next block

# Proposed Method



- src 1
- src 2
- src 3

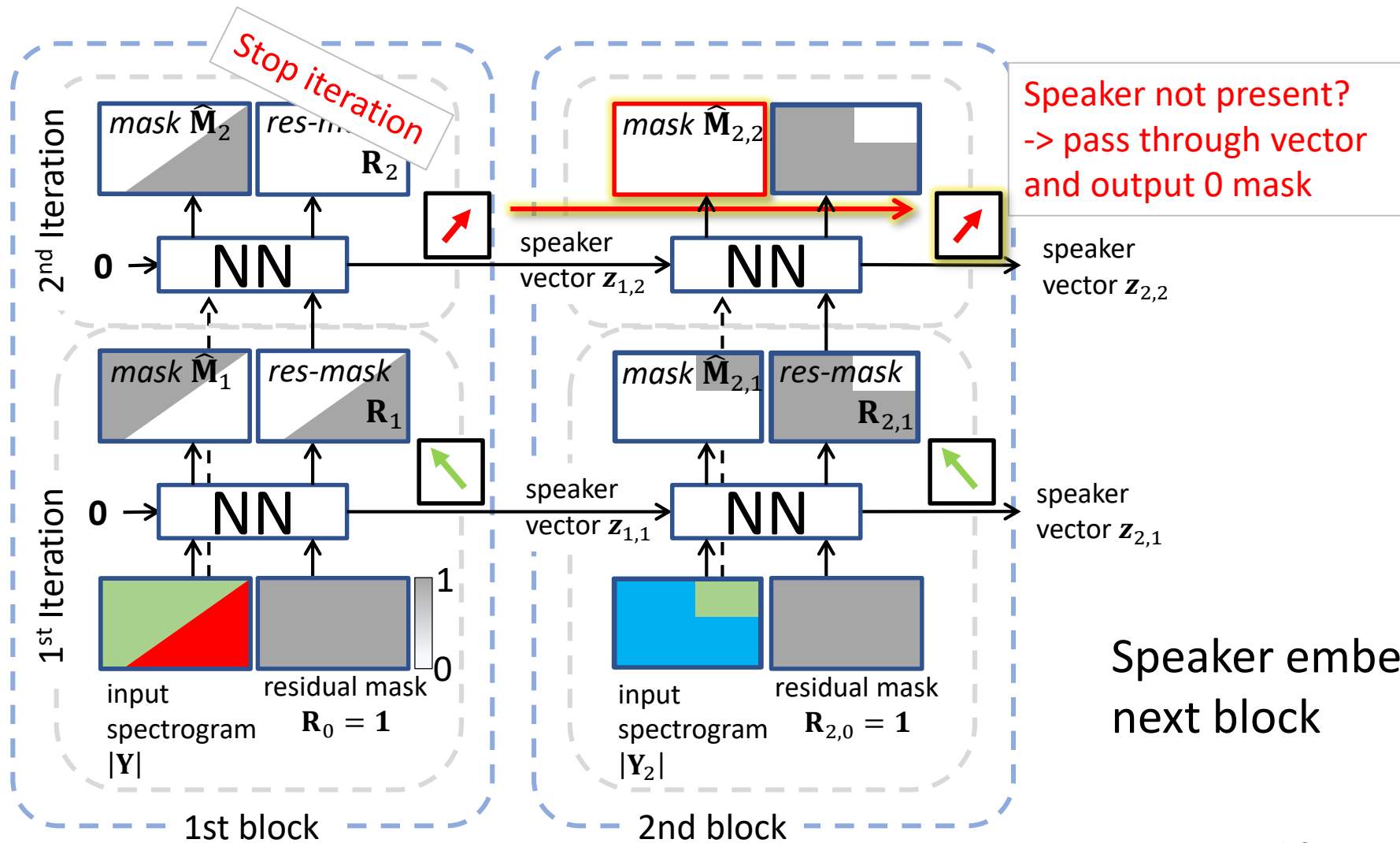


Speaker embedding is passed to next block

# Proposed Method

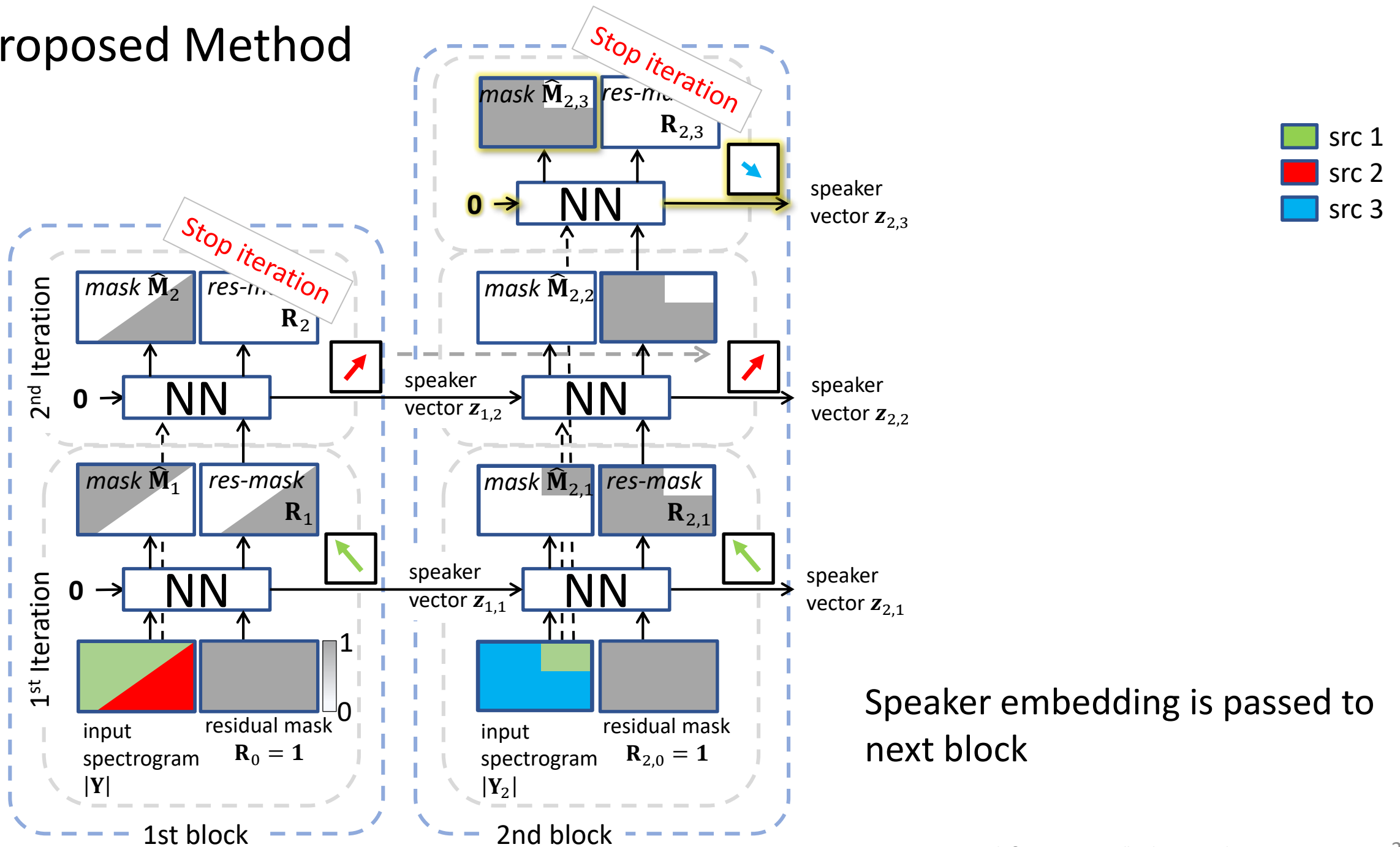


- src 1
- src 2
- src 3

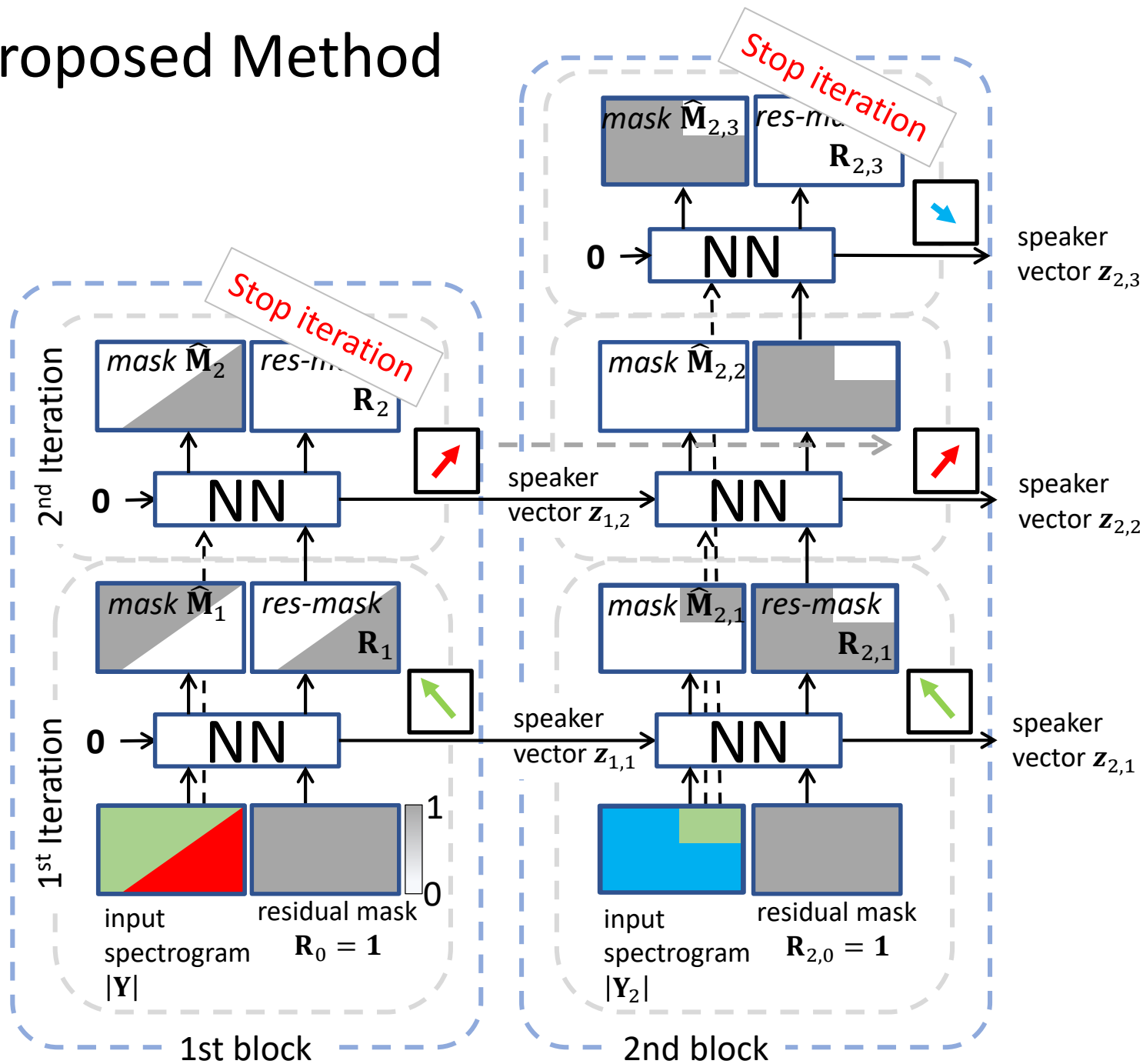


Speaker embedding is passed to next block

# Proposed Method



# Proposed Method



**New Speakers**  
Can estimate new speaker vector in each block

**Silent Speakers**  
Notices silent speakers and estimates 0 mask

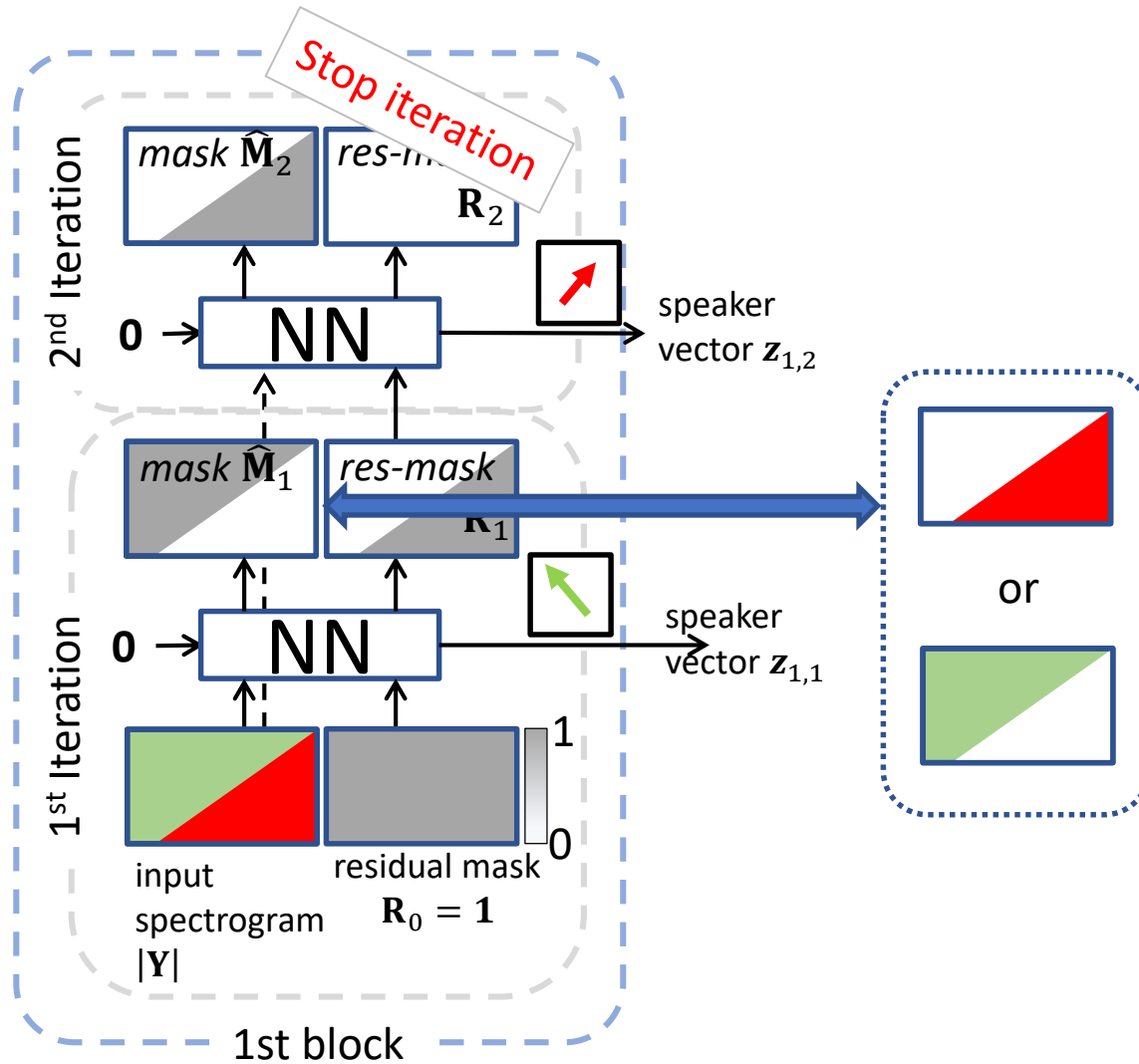
**Block Permutation Problem**  
Estimated signals always in the same order

# Proposed Method: Loss



- src 1
- src 2
- src 3

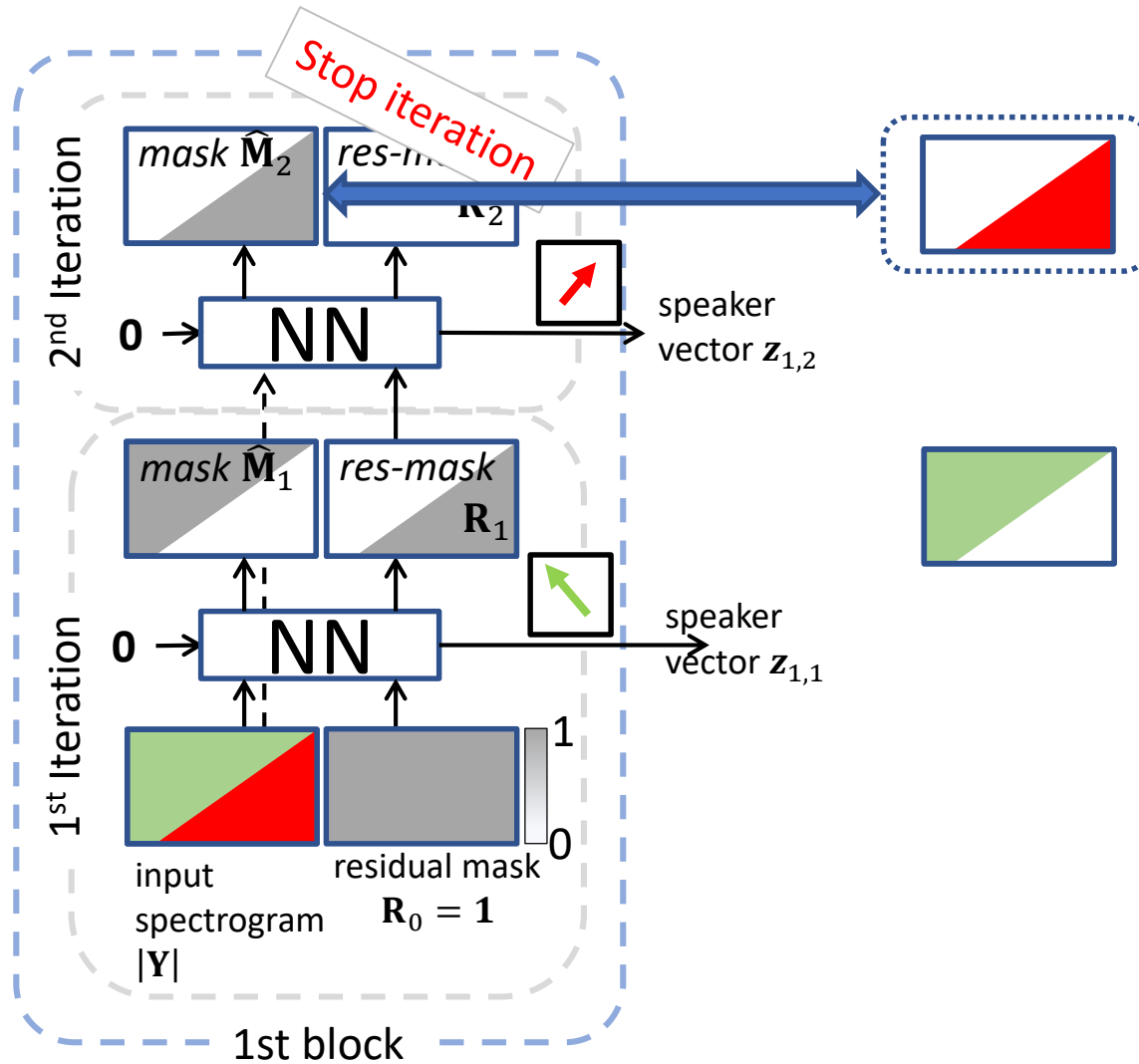
## permutation invariant loss



# Proposed Method: Loss

-  src 1
-  src 2
-  src 3

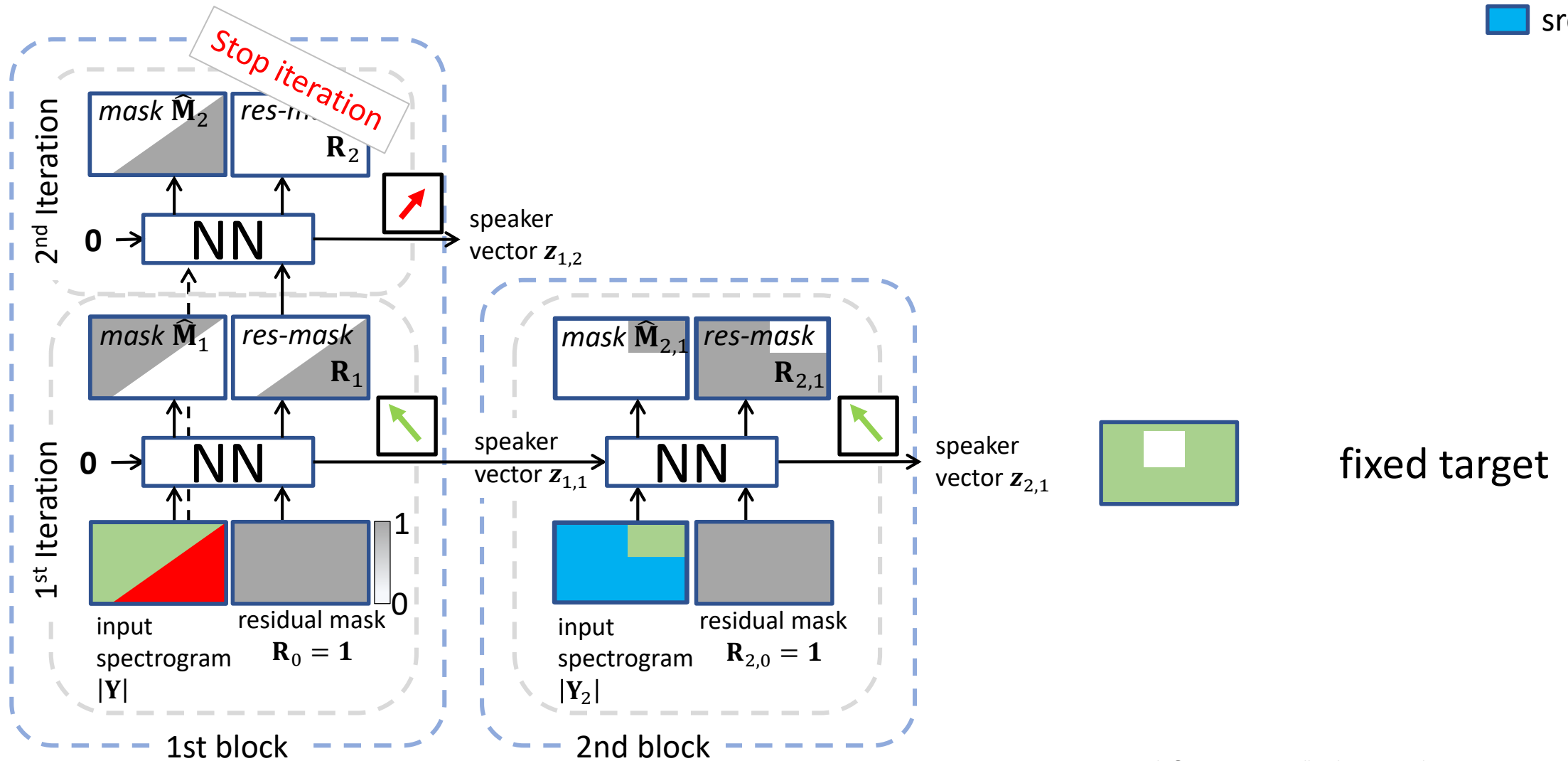
permutation invariant loss



# Proposed Method: Loss



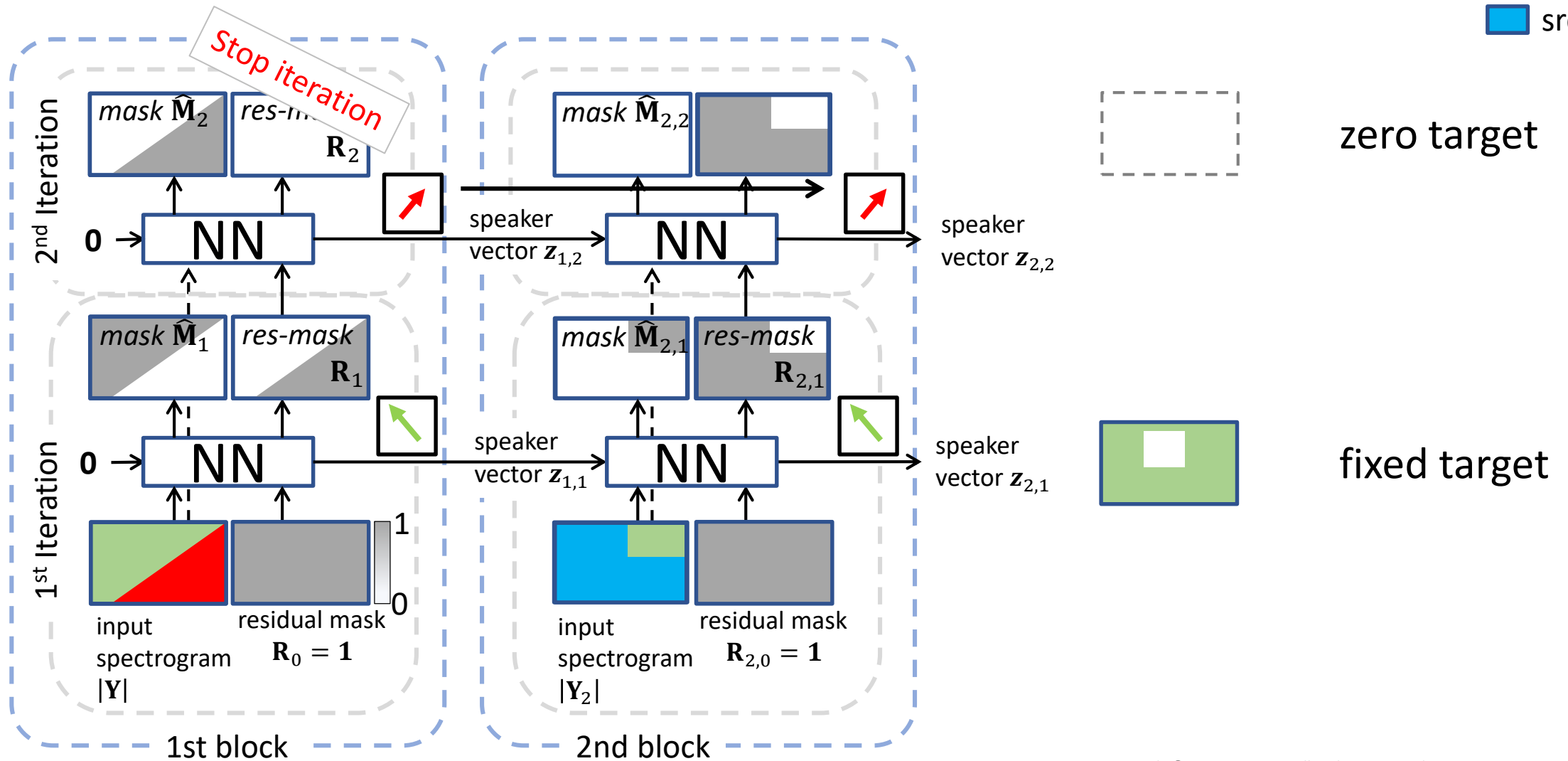
- src 1
- src 2
- src 3

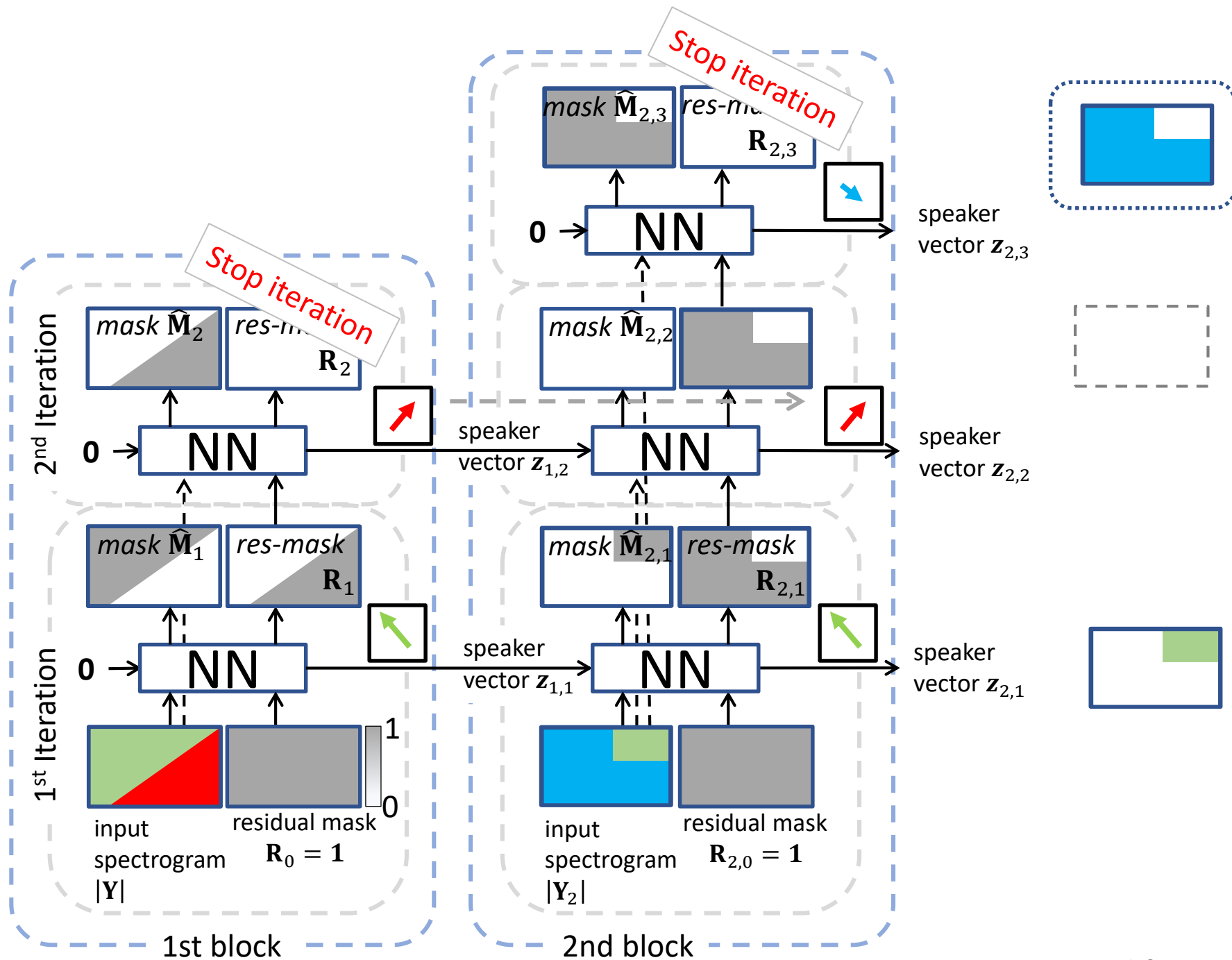




# Proposed Method: Loss

- src 1
- src 2
- src 3





permutation  
invariant loss on  
remaining  
target(s)

zero target

fixed target

# Proposed Method - Experiments

- Artificially mixed from wsj
- 1 or 2 speakers
- Blocks of 2.5s
- Random activity pattern
- $\geq 1$  speaker active in first block

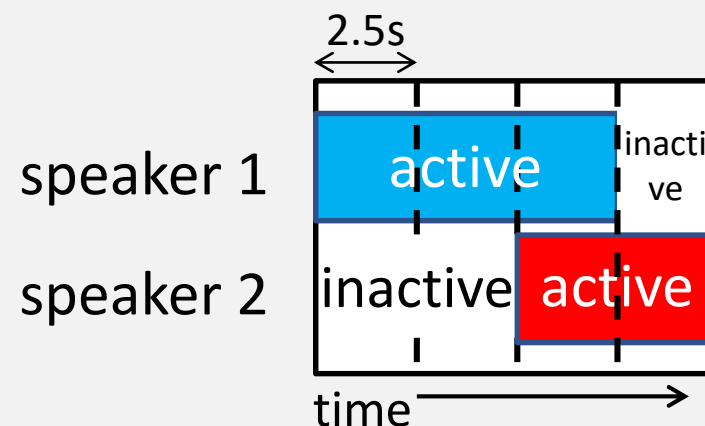
## Training Data

- 4 blocks / 10s

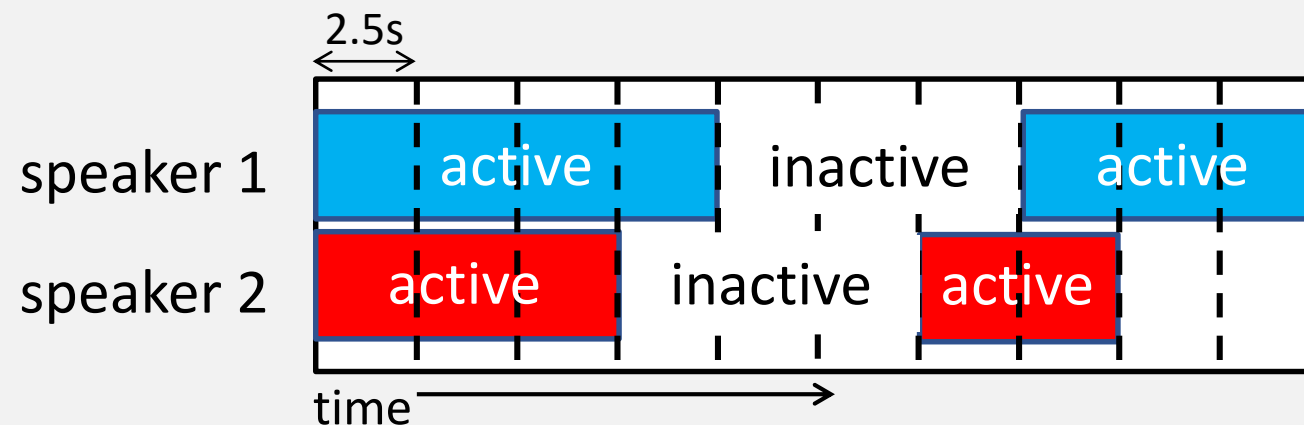
## Test Data

- 12 blocks / 30s
- ~5 hours
- Considerably longer than training data

## Example Train Mixture



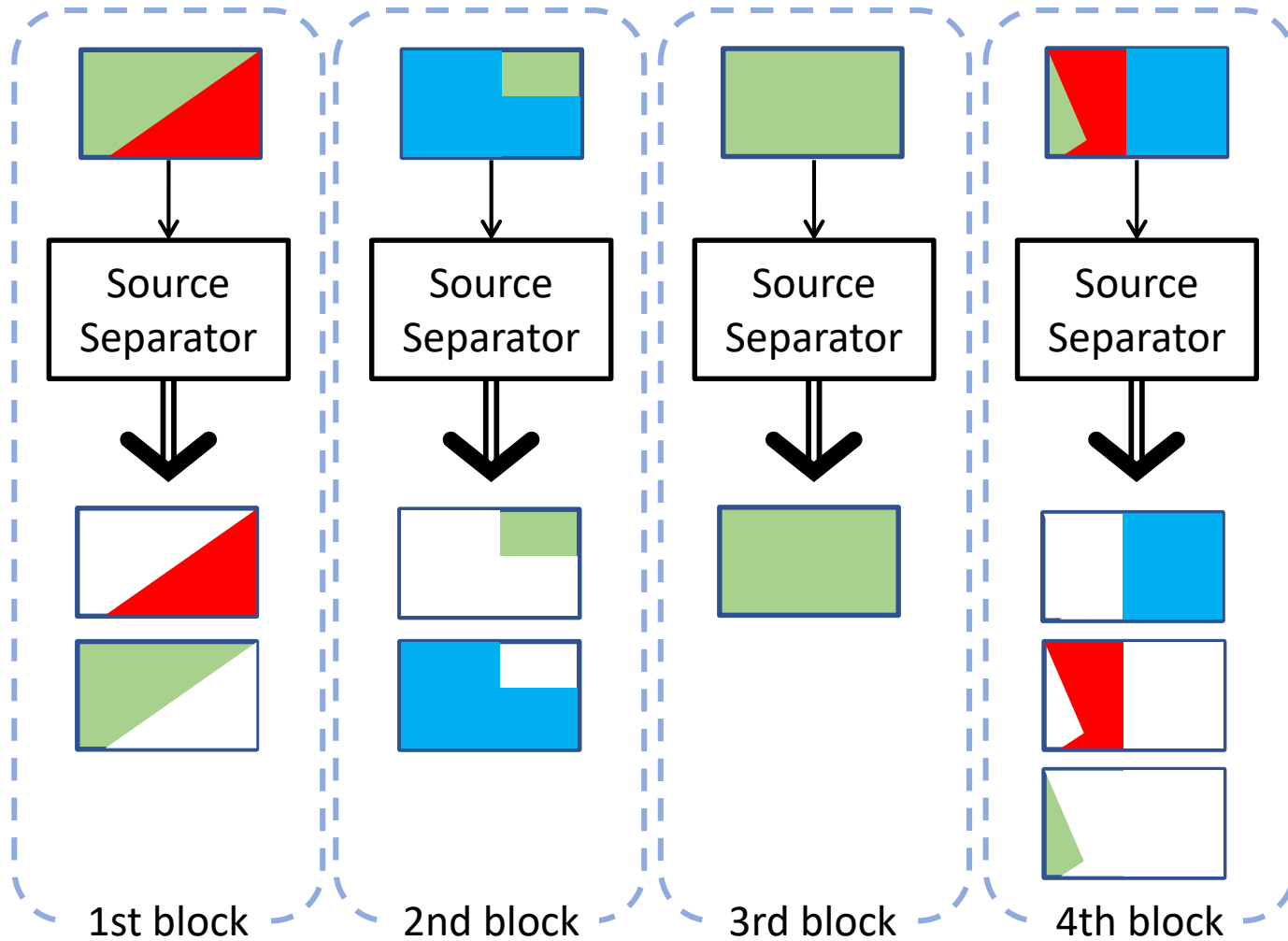
## Example Test Mixture



# Conventional 2-stage Method: Clustering of Speaker Characteristics



- src 1
- src 2
- src 3

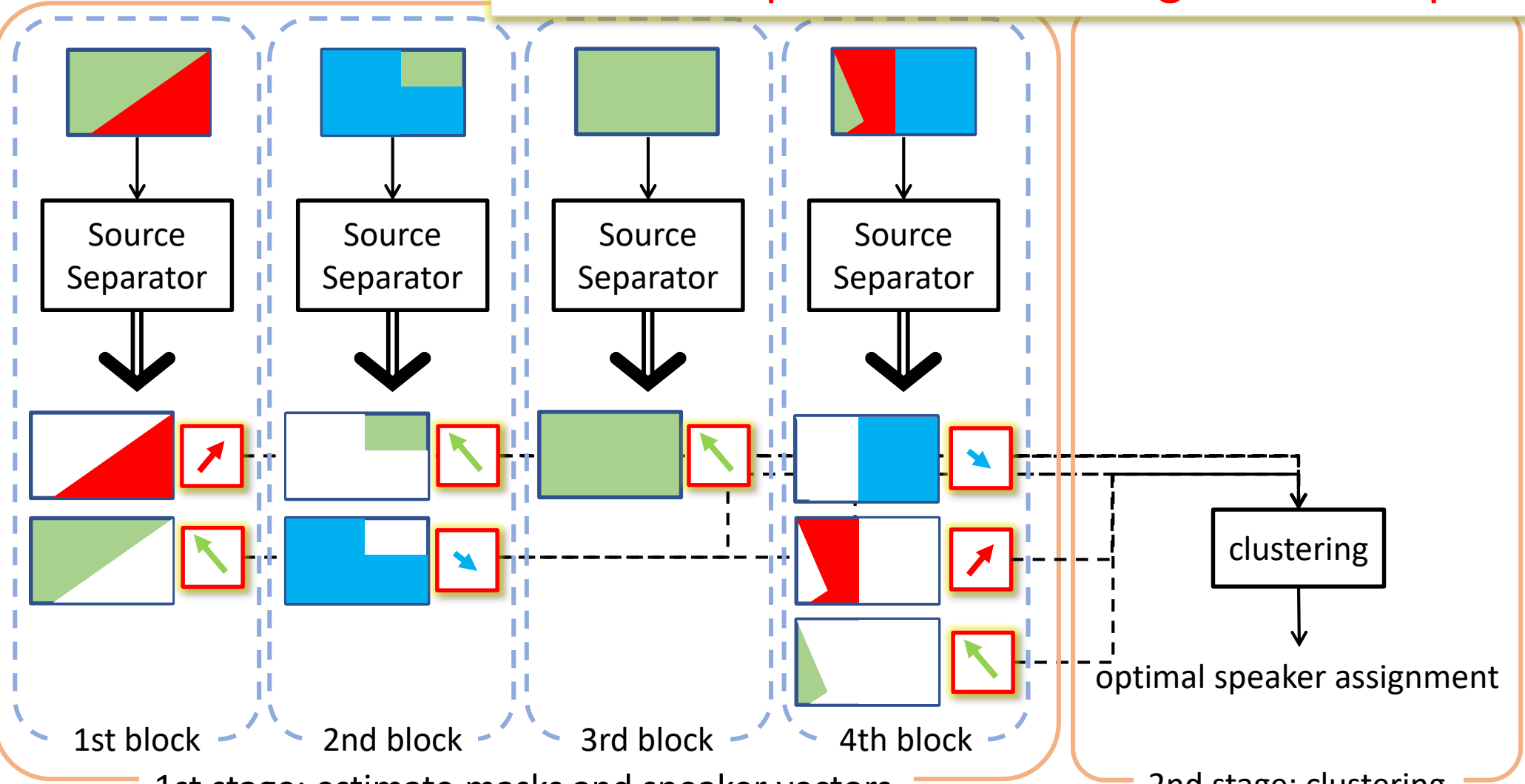


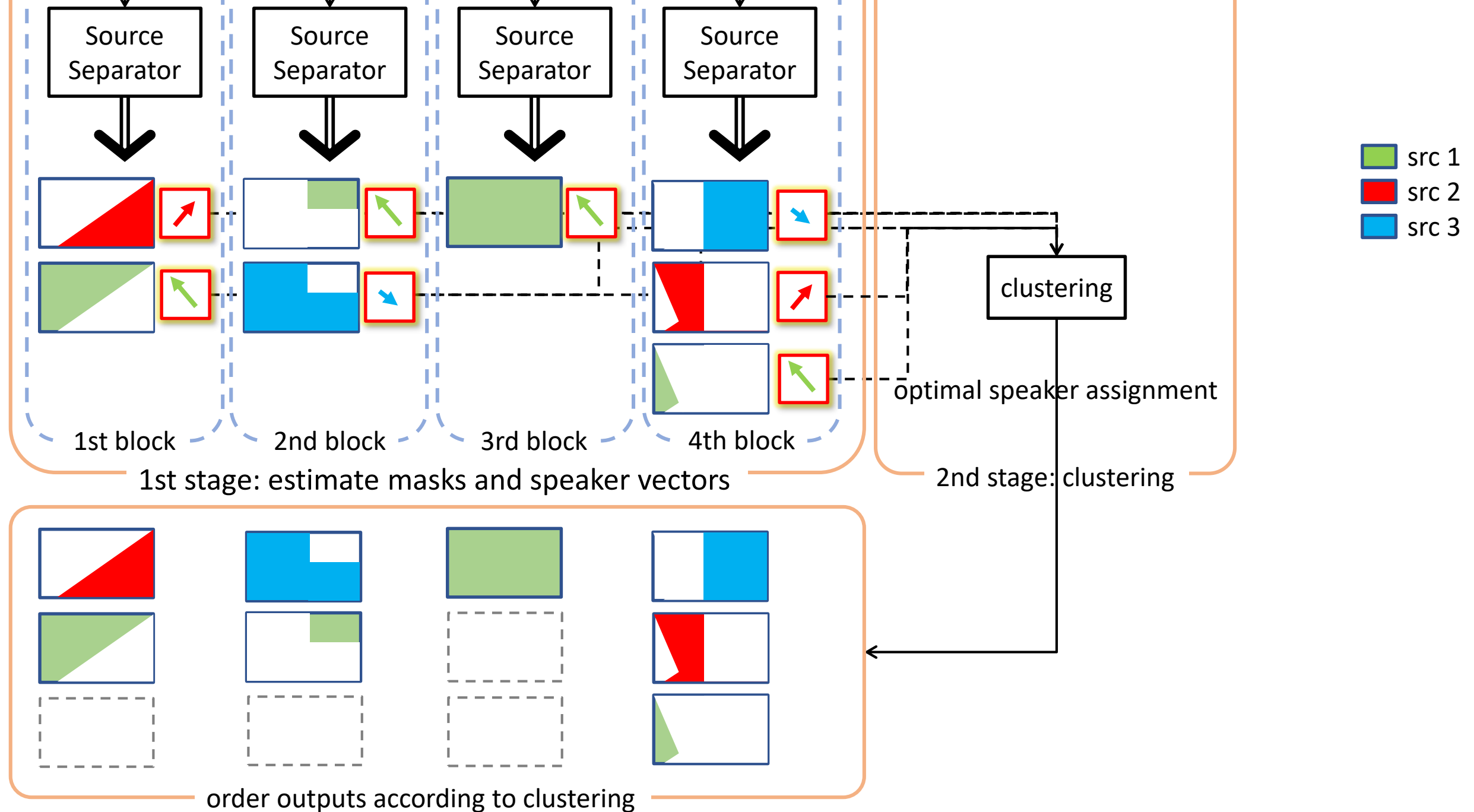
# Conventional 2-stage Method: Clustering of Speaker Characteristics



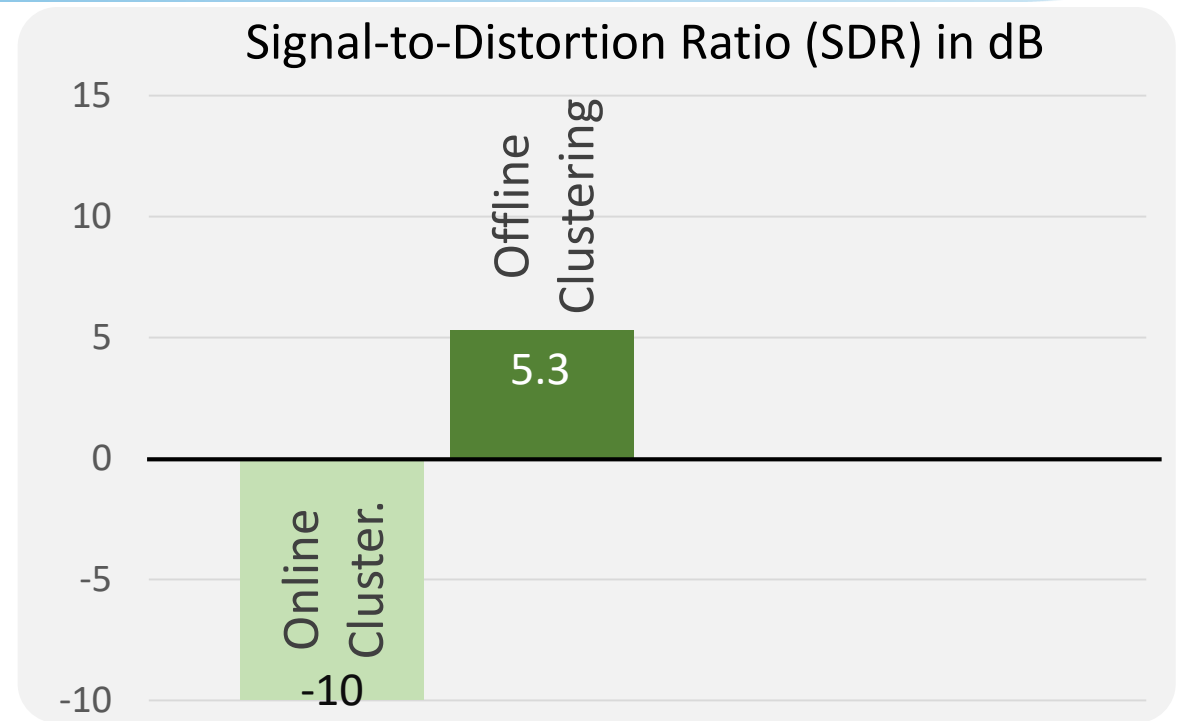
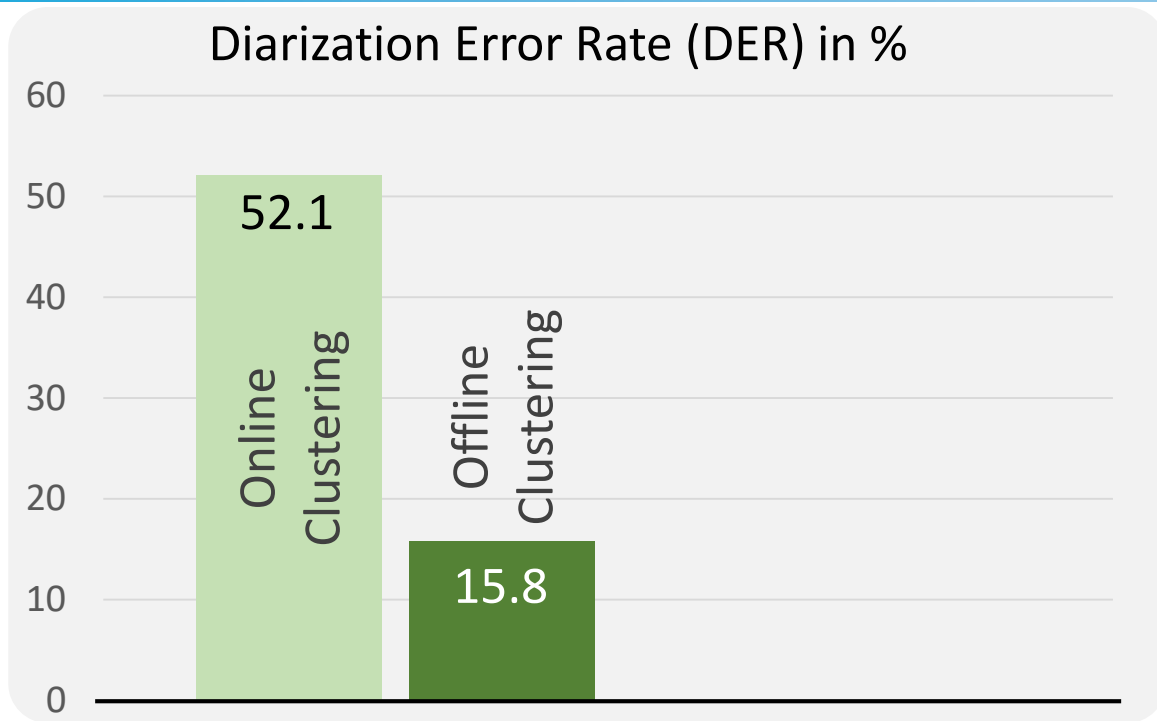
additional speaker embedding vector output

- src 1
- src 2
- src 3





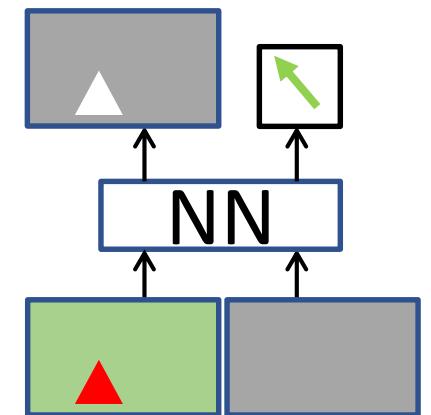
# Evaluation on 12-block (30s) Mixtures



**RSAN with Speaker Embedding + clustering**

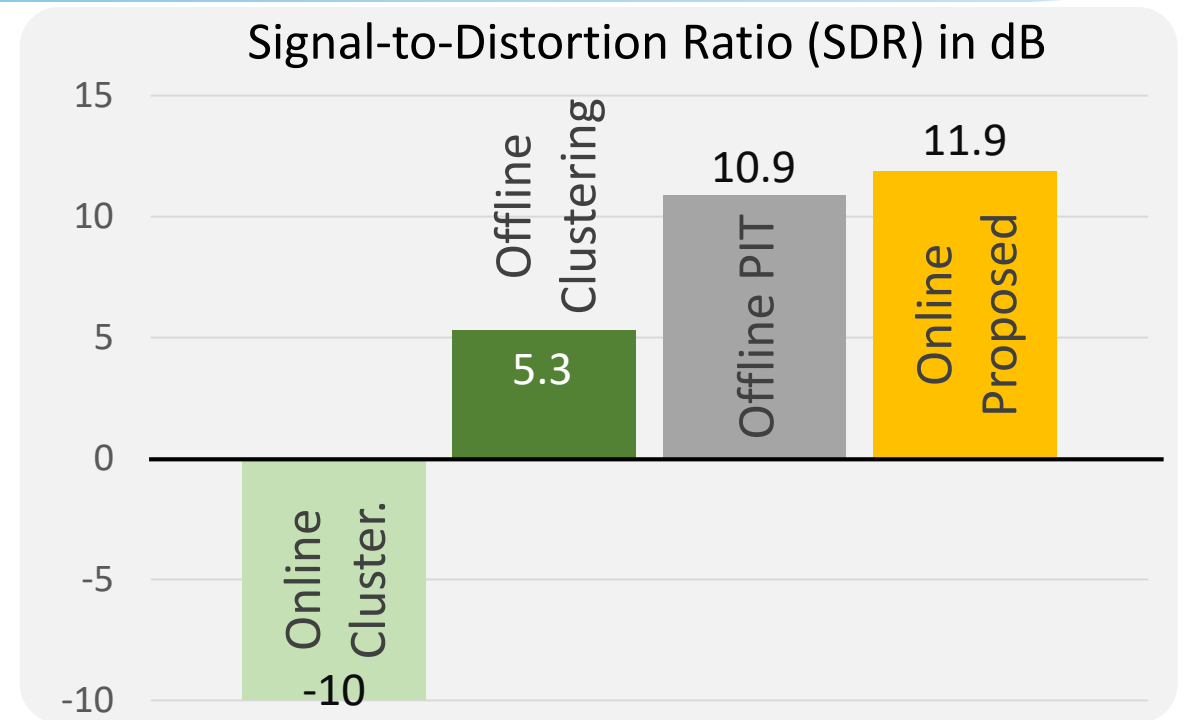
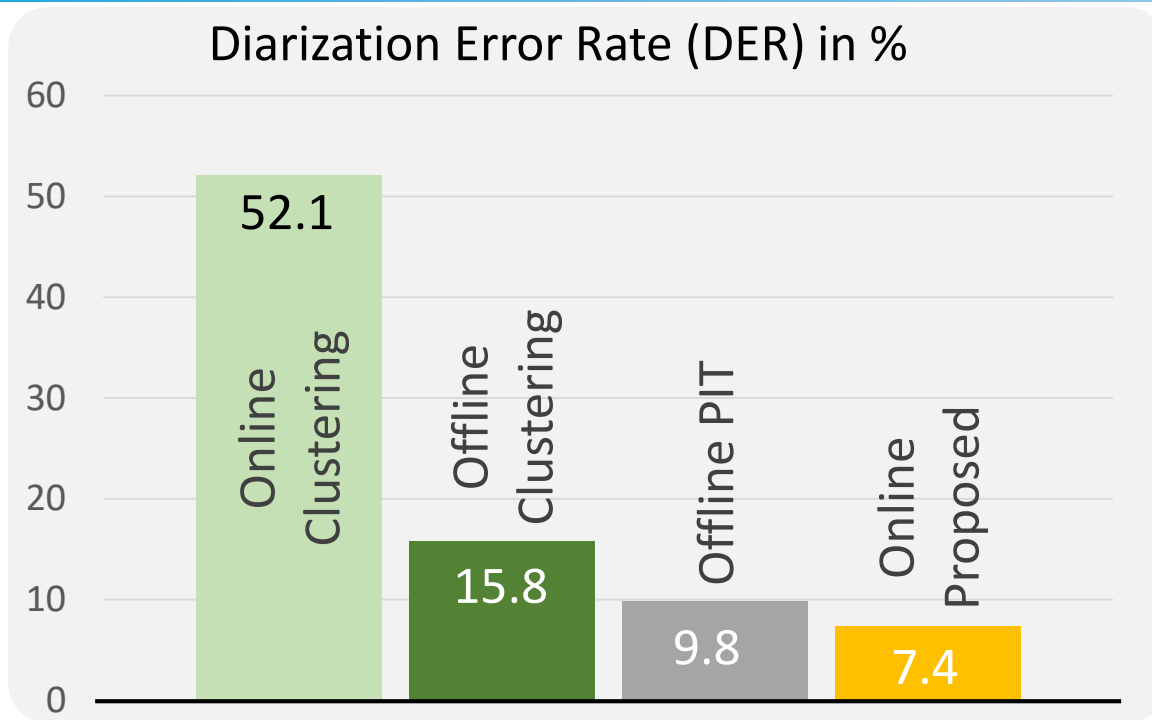
**Online Clustering:** leader-follower clustering

**Offline Clustering:** hierarchical clustering given the oracle #speakers



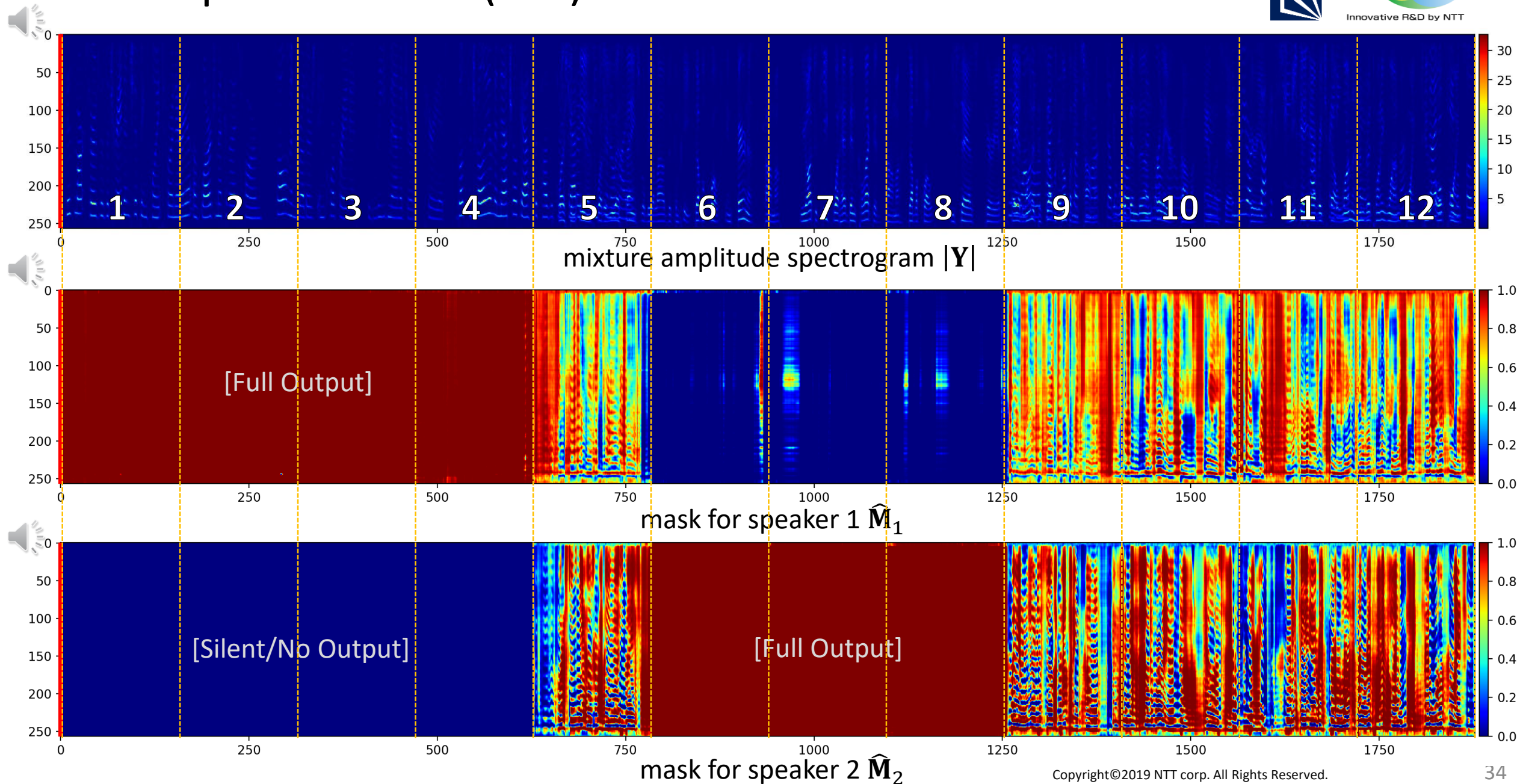


# Evaluation on 12-block (30s) Mixtures



- **Proposed Method generalizes well to an unseen number of blocks**
  - Trained on 4 blocks, evaluated on 12 blocks
- **Proposed Method outperforms the other approaches**

# Example: 12-block (30s) Mixture





## Problems in Meeting Scenarios:

- ⇒ source separation
- ⇒ source count estimation
  
- ⇒ blockwise/online processing

## Solved:

- ✓ **RSAN**
  - ✓ iterative source extraction
  - ✓ count of iterations
- ✓ **Proposed all-neural block-online method based on RSAN**
  - ✓ generalizes well to unseen number of blocks



# Thank you for your attention!