

---

**torchiva**

***Release 0.0.1***

**Robin Scheibler, Kohei Saijo**

**Jan 30, 2023**



## **CONTENTS:**

<b>1</b>	<b>Quick Start</b>	<b>3</b>
<b>2</b>	<b>Separation using Pre-trained Model</b>	<b>5</b>
<b>3</b>	<b>Training</b>	<b>7</b>
<b>4</b>	<b>Authors</b>	<b>9</b>
<b>5</b>	<b>License</b>	<b>11</b>
<b>6</b>	<b>API</b>	<b>13</b>
6.1	T_ISS . . . . .	13
6.2	AuxIVA_IP . . . . .	14
6.3	AuxIVA_IP2 . . . . .	15
6.4	WPE . . . . .	16
6.5	MVDR beamformer . . . . .	16
6.6	MWF beamformer . . . . .	17
6.7	GEV beamformer . . . . .	18
6.8	FIVE . . . . .	18
<b>7</b>	<b>Indices and tables</b>	<b>21</b>
<b>Index</b>		<b>23</b>



---

A package for blind source separation and beamforming in [pytorch](#) .

- supports many BSS and beamforming methods
- supports memory efficient gradient computation for training neural source models
- supports batched computations
- can run on GPU via pytorch



---

**CHAPTER  
ONE**

---

**QUICK START**

The package can be installed via pip:

```
pip install torchiva
```



---

CHAPTER  
TWO

---

## SEPARATION USING PRE-TRAINED MODEL

We provide a pre-trained model in `trained_models/tiss`. You can easily try separation with the pre-trained model:

```
# Separation
python -m torchiva.separate INPUT OUTPUT
```

where INPUT is either a multichannel wav file or a folder containing multichannel wav files. If a folder, then all the files inside are separated. The output is saved to OUTPUT. The model stored in `trained_models/tiss` is automatically downloaded to `$HOME/.torchiva_models`. The path or url to the model can also be manually provided via the `--model` option. The model was trained on the [WSJ1-mix dataset](#) with the same configuration as `./examples/configs/tiss.json`.



---

## CHAPTER THREE

---

## TRAINING

We provide some simple training scripts. We support training of **T-ISS, MWF, MVDR, GEV**:

```
cd examples

# install some modules necessary for training
pip install -r requirements.txt

# training
python train.py PATH_TO_CONFIG PATH_TO_DATASET
```

Note that our example scripts assumes using WSJ1-mix dataset. If you want to use other datasets, please change the script in the part that loads audios.

Test your trained model with checkpoint from epoch 128:

```
# python ./test.py --dataset ../wsj1_6ch --n_fft 2048 --hop 512 --n_iter 40 --iss-
↪hparams checkpoints/tiss_delay1tap5_2ch/lightning_logs/version_0/hparams.yaml --epoch
↪128 --test
```

Export the trained model for later use:

```
python ./export_model.py ../trained_models/tiss checkpoints/tiss_delay1tap5_2ch/
↪lightning_logs/version_0 128 146 148 138 122 116 112 108 104 97
```

Run the example script using the exported model:

```
python ./example_dnn.py ../wsj1_6ch ../trained_models/tiss -m 2 -r 100
```



---

**CHAPTER  
FOUR**

---

**AUTHORS**

- Robin Scheibler
- Kohei Saijo



---

**CHAPTER  
FIVE**

---

**LICENSE**

2022 (c) Robin Scheibler, Kohei Saijo, LINE Corporation.

All of this code is released under [MIT License](#)



## 6.1 T\_ISS

```
class torchiva.T_ISS(n_iter=10, n_taps=0, n_delay=0, n_src=None, model=None, proj_back_mic=0,  
use_dmc=False, eps=None)
```

Joint dereverberation and separation with *time-decorrelation iterative source steering* (T-ISS)<sup>1</sup>.

Parameters can also be specified during a forward call. In this case, the forward argument is only used in that forward process and **does not rewrite class attributes**.

### Parameters

- **n\_iter** (*int, optional*) – The number of iterations. (default: 10)
- **n\_taps** (*int, optional*) – The length of the dereverberation filter. If set to 0, this method works as the normal AuxIVA with ISS update<sup>2</sup> (default: 0).
- **n\_delay** (*int, optional*) – The number of delay for dereverberation (default: 0).
- **n\_src** (*int, optional*) – The number of sources to be separated. When n\_src < n\_chan, a computationally cheaper variant (Over-T-ISS)<sup>3</sup> is used. If set to None, n\_src is set to n\_chan (default: None)
- **model** (*torch.nn.Module, optional*) – The model of source distribution. Mask estimation neural network can also be used. If None, spherical Laplace is used (default: None).
- **proj\_back\_mic** (*int, optional*) – The reference mic index to perform projection back. If set to None, projection back is not applied (default: 0).
- **use\_dmc** (*bool, optional*) – If set to True, memory efficient Demixing Matrix Checkpointing (DMC)<sup>4</sup> is used to compute the gradient. It reduces the memory cost to that of a single iteration when training neural source model (default: False).
- **eps** (*float, optional*) – A small constant to make divisions and the like numerically stable (default:None).

```
forward(n_iter=None, n_taps=None, n_delay=None, n_src=None, model=None, proj_back_mic=None,  
use_dmc=None, eps=None)
```

---

<sup>1</sup> T. Nakashima, R. Scheibler, M. Togami, and N. Ono, “Joint dereverberation and separation with iterative source steering”, ICASSP, 2021, <https://arxiv.org/pdf/2102.06322.pdf>.

<sup>2</sup> R. Scheibler, and N Ono, “Fast and stable blind source separation with rank-1 updates” ICASSP, 2021,

<sup>3</sup> R. Scheibler, W. Zhang, X. Chang, S. Watanabe, and Y. Qian, “End-to-End Multi-speaker ASR with Independent Vector Analysis”, arXiv preprint arXiv:2204.00218, 2022, <https://arxiv.org/pdf/2204.00218.pdf>.

<sup>4</sup> K. Saito, and R. Scheibler, “Independence-based Joint Speech Dereverberation and Separation with Neural Source Model”, arXiv preprint arXiv:2110.06545, 2022, <https://arxiv.org/pdf/2110.06545.pdf>.

**Parameters**

**X** (`torch.Tensor`) – The input mixture in STFT-domain, shape `(..., n_chan, n_freq, n_frames)`

**Returns**

**Y** – The separated and dereverberated signal in STFT-domain

**Return type**

`torch.Tensor`, shape `(..., n_src, n_freq, n_frames)`

---

**Note:**

This class can handle various BSS methods with ISS update rule depending on the specified arguments:

- IVA-ISS: `n_taps=0, n_delay=0, n_chan==n_src, model=LaplaceModel()` or `GaussModel()`
  - ILRMA-ISS: `n_taps=0, n_delay=0, n_chan==n_src, model=NMFModel()`
  - DNN-IVA-ISS: `n_taps=0, n_delay=0, n_chan==n_src, model="DNN"`
  - OverIVA-ISS: `n_taps=0, n_delay=0, n_chan < n_src`
  - ILRMA-T-ISS<sup>Page 13, 1</sup>: `n_taps>0, n_delay>0, n_chan==n_src, model=NMFModel()`
  - DNN-T-ISS<sup>Page 13, 4</sup>: `n_taps>0, n_delay>0, n_chan==n_src, model="DNN"`
  - Over-T-ISS<sup>Page 13, 3</sup>: `n_taps>0, n_delay>0, n_chan > n_src`
- 

**References**

## 6.2 AuxIVA\_IP

`class torchiva.AuxIVA_IP(n_iter=10, n_src=None, model=None, proj_back_mic=0, eps=None)`

Independent vector analysis (IVA) with iterative projection (IP) update<sup>5</sup>.

We do not support ILRMA-T with IP updates.

**Parameters**

- **n\_iter** (`int, optional`) – The number of iterations. (default: 10)
- **n\_src** (`int, optional`) – The number of sources to be separated. When `n_src < n_chan`, a computationally cheaper variant (OverIVA)<sup>6</sup> is used. If set to None, `n_src` is set to `n_chan` (default: None)
- **model** (`torch.nn.Module, optional`) – The model of source distribution. If None, spherical Laplace is used (default: None).
- **proj\_back\_mic** (`int, optional`) – The reference mic index to perform projection back. If set to None, projection back is not applied (default: 0).
- **eps** (`float, optional`) – A small constant to make divisions and the like numerically stable (default:None).

<sup>5</sup> N. Ono, “Stable and fast update rules for independent vector analysis based on auxiliary function technique”, WASSPA, 2011.

<sup>6</sup> R. Scheibler, and N Ono, “Independent vector analysis with more microphones than sources”, WASSPA, 2019, <https://arxiv.org/pdf/1905.07880.pdf>.

---

```
forward(X, n_iter=None, n_src=None, model=None, proj_back_mic=None, eps=None)
```

**Parameters**

**X** (`torch.Tensor`) – The input mixture in STFT-domain, shape `(..., n_chan, n_freq, n_frames)`

**Returns**

**Y** – The separated signal in STFT-domain

**Return type**

`torch.Tensor`, shape `(..., n_src, n_freq, n_frames)`

**Note:**

This class can handle two BSS methods with IP update rule depending on the specified arguments:

- AuxIVA-IP: `n_chan==n_src`, `model=LaplaceModel()` or `GaussModel()`
  - ILRMA-IP: `n_chan==n_src`, `model=NMFModel()`
  - OverIVA\_IP<sup>6</sup>: `n_taps=0`, `n_delay=0`, `n_chan==n_src`, `model=NMFModel()`
- 

**References**

### 6.3 AuxIVA\_IP2

```
class torchiva.AuxIVA_IP2(n_iter=10, model=None, proj_back_mic=0, eps=None)
```

Blind source separation based on independent vector analysis with alternating updates of the mixing vectors<sup>7</sup>

**Parameters**

- **n\_iter** (`int`, *optional*) – The number of iterations (default: 10).
- **model** (`torch.nn.Module`, *optional*) – The model of source distribution. If None, spherical Laplace is used (default: None).
- **proj\_back\_mic** (`int`, *optional*) – The reference mic index to perform projection back. If set to None, projection back is not applied (default: 0).
- **eps** (`float`, *optional*) – A small constant to make divisions and the like numerically stable (default:None).

```
forward(X, n_iter=None, model=None, proj_back_mic=None, eps=None)
```

**Parameters**

**X** (`torch.Tensor`) – The input mixture in STFT-domain, shape `(..., n_chan, n_freq, n_frames)`

**Returns**

**X** – The separated signal in STFT-domain.

**Return type**

`torch.Tensor`, shape `(..., n_chan, n_freq, n_frames)`

---

<sup>7</sup> N. Ono, “Fast stereo independent vector analysis and its implementation on mobile phone”, IWAENC, 2012.

## References

## 6.4 WPE

```
class torchiva.WPE(n_iter=3, n_delay=3, n_taps=5, model=None, eps=1e-05)
```

Weighted prediction error (WPE)<sup>9</sup>.

### Parameters

- **n\_iter** (*int, optional*) – The number of iterations. (default: 3)
- **n\_taps** (*int, optional*) – The length of the dereverberation filter (default: 5).
- **n\_delay** (*int, optional*) – The number of delay for dereverberation (default: 3).
- **model** (*torch.nn.Module, optional*) – The model of source distribution. If None, time-varying Gaussian is used. (default: None).
- **eps** (*float, optional*) – A small constant to make divisions and the like numerically stable (default: 1e-5).

### Returns

**Y** – The dereverberated signal in STFT-domain.

### Return type

`torch.Tensor, shape (..., n_src, n_freq, n_frames)`

## References

## 6.5 MVDR beamformer

```
class torchiva.MVDRBeamformer(mask_model, ref_mic=0, eps=1e-05, mvdr_type='rtf', n_power_iter=None)
```

Implementation of MVDR beamformer. This class is basically assumes DNN-based beamforming. also supports the case of estimating three masks

### Parameters

- **mask\_model** (*torch.nn.Module*) – A function that is given one spectrogram and returns 2 or 3 masks of the same size as the input. When 3 masks (1 for target and the rest 2 for noise) are estimated, they are utilized as in<sup>10</sup>
- **ref\_mic** (*int, optional*) – Reference channel (default: 0)
- **eps** (*float, optional*) – A small constant to make divisions and the like numerically stable (default: 1e-5).
- **mvdr\_type** (*str, optional*) – The way to obtain the MVDR weight. If set to `rtf`, relative transfer function is computed to obtain MVDR. If set to ‘scm’, MVDR weight is obtained directly with spatial covariance matrices<sup>11</sup> (default: `rtf`).

<sup>9</sup> T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi, and B. H. Juang, “Speech dereverberation based on variance-normalized delayed linear prediction”, IEEE Trans. on Audio, Speech, and Lang. Process., 2010.

<sup>10</sup> C. Boeddeker et al., “Convulsive Transfer Function Invariant SDR training criteria for Multi-Channel Reverberant Speech Separation”, ICASSP, 2021.

<sup>11</sup> Mehrez Souden, Jacob Benesty, and Sofiene Affes, “On optimal frequency-domain multichannel linear filtering for noise reduction”, IEEE Trans. on audio, speech, and lang. process., 2009.

- **n\_power\_iter** (*int, optional*) – Use the power iteration method to compute the relative transfer function instead of the full generalized eigenvalue decomposition (GEVD). The number of iteration desired should be provided. If set to `None`, the full GEVD is used (default: `None`).

**forward**(*X, mask\_model=None, ref\_mic=None, eps=None, mvdr\_type=None, n\_power\_iter=None*)

#### Parameters

**X** (*torch.Tensor*) – The input mixture in STFT-domain, shape `(..., n_chan, n_freq, n_frames)`

#### Returns

**Y** – The separated signal in STFT-domain

#### Return type

`torch.Tensor, shape (..., n_src, n_freq, n_frames)`

## References

## 6.6 MWF beamformer

**class** `torchiva.MWFBeamformer(mask_model, ref_mic=0, eps=1e-05, time_invariant=True)`

Implementation of MWF beamformer described in<sup>12</sup>. This class is basically assumes DNN-based beamforming.

#### Parameters

- **mask\_model** (*torch.nn.Module*) – A function that is given one spectrogram and returns 2 masks of the same size as the input.
- **ref\_mic** (*int, optional*) – Reference channel (default: `0`)
- **eps** (*float, optional*) – A small constant to make divisions and the like numerically stable (default:`1e-5`).
- **time\_invariant** (*bool, optional*) – If set to `True`, this flag indicates that we want to use the time-invariant version of MWF. If set to `False`, the time-varying MWF is used instead (default: `True`).

**forward**(*X, mask\_model=None, ref\_mic=None, eps=None, time\_invariant=None*)

#### Parameters

**X** (*torch.Tensor*) – The input mixture in STFT-domain, shape `(..., n_chan, n_freq, n_frames)`

#### Returns

**Y** – The separated signal in STFT-domain

#### Return type

`torch.Tensor, shape (..., n_src, n_freq, n_frames)`

<sup>12</sup> Y. Masuyama et al., “Consistency-aware multi-channel speech enhancement using deep neural networks”, ICASSP, 2020.

## References

## 6.7 GEV beamformer

```
class torchiva.GEVBeamformer(mask_model, ref_mic=0, eps=1e-05)
```

Implementation of GEV beamformer. This class is basically assumes DNN-based beamforming.

### Parameters

- **mask\_model** (`torch.nn.Module`) – A function that is given one spectrogram and returns 2 masks of the same size as the input.
- **ref\_mic** (`int, optional`) – Reference channel (default: 0)
- **eps** (`float, optional`) – A small constant to make divisions and the like numerically stable (default:1e-5).

```
forward(X, mask_model=None, ref_mic=None, eps=None)
```

### Parameters

**X** (`torch.Tensor`) – The input mixture in STFT-domain, shape `(..., n_chan, n_freq, n_frames)`

### Returns

**Y** – The separated signal in STFT-domain

### Return type

`torch.Tensor, shape (..., n_src, n_freq, n_frames)`

## 6.8 FIVE

```
class torchiva.FIVE(n_iter=10, model=None, proj_back_mic=0, eps=None, n_power_iter=None)
```

Fast independent vector extraction (FIVE)<sup>8</sup>. FIVE extracts one source from the input signal.

### Parameters

- **n\_iter** (`int, optional`) – The number of iterations (default: 10).
- **model** (`torch.nn.Module, optional`) – The model of source distribution (default: LaplaceModel).
- **proj\_back\_mic** (`int, optional`) – The reference mic index to perform projection back. If set to None, projection back is not applied (default: 0).
- **eps** (`float, optional`) – A small constant to make divisions and the like numerically stable (default: None).
- **n\_power\_iter** (`int, optional`) – The number of power iterations. If set to None, eigen-vector decomposition is used instead. (default: None)

```
forward(X, n_iter=None, model=None, proj_back_mic=None, eps=None)
```

### Parameters

**X** (`torch.Tensor`) – The input mixture in STFT-domain, shape `(..., n_chan, n_freq, n_frames)`

<sup>8</sup> R. Scheibler, and N Ono, “Fast independent vector extraction by iterative SINR maximization”, ICASSP, 2020, <https://arxiv.org/pdf/1910.10654.pdf>.

**Returns**

**Y** – The extracted *one* signal in STFT-domain.

**Return type**

`torch.Tensor, shape (..., n_freq, n_frames)`

**References**



---

CHAPTER  
**SEVEN**

---

## **INDICES AND TABLES**

- genindex
- modindex
- search



# INDEX

## A

`AuxIVA_IP` (*class in torchiva*), 14  
`AuxIVA_IP2` (*class in torchiva*), 15

## F

`FIVE` (*class in torchiva*), 18  
`forward()` (*torchiva.AuxIVA\_IP method*), 14  
`forward()` (*torchiva.AuxIVA\_IP2 method*), 15  
`forward()` (*torchiva.FIVE method*), 18  
`forward()` (*torchiva.GEVBeamformer method*), 18  
`forward()` (*torchiva.MVDRBeamformer method*), 17  
`forward()` (*torchiva.MWFBeamformer method*), 17  
`forward()` (*torchiva.T\_ISS method*), 13

## G

`GEVBeamformer` (*class in torchiva*), 18

## M

`module`  
    `torchiva`, 11  
`MVDRBeamformer` (*class in torchiva*), 16  
`MWFBeamformer` (*class in torchiva*), 17

## T

`T_ISS` (*class in torchiva*), 13  
`torchiva`  
    `module`, 11

## W

`WPE` (*class in torchiva*), 16