# GASS: GENERALIZING AUDIO SOURCE SEPARATION WITH LARGE-SCALE DATA

*Jordi Pons\*    Xiaoyu Liu\*    Santiago Pascual    Joan Serrà*

Dolby Laboratories

## ABSTRACT

Universal source separation targets at separating the audio sources of an arbitrary mix, removing the constraint to operate on a specific domain like speech or music. Yet, the potential of universal source separation is limited because most existing works focus on mixes with predominantly sound events, and small training datasets also limit its potential for supervised learning. Here, we study a single general audio source separation (GASS) model trained to separate speech, music, and sound events in a supervised fashion with a large-scale dataset. We assess GASS models on a diverse set of tasks. Our strong in-distribution results show the feasibility of GASS models, and the competitive out-of-distribution performance in sound event and speech separation shows its generalization abilities. Yet, it is challenging for GASS models to generalize for separating out-of-distribution cinematic and music content. We also fine-tune GASS models on each dataset and consistently outperform the ones without pre-training. All fine-tuned models (except the music separation one) obtain state-of-the-art results in their respective benchmarks.

*Index Terms*— General audio source separation, deep learning.

## 1. INTRODUCTION

Audio source separation consists of isolating the sources present in an audio mix. Most previous works frame the problem as a source-specific task, as in speech source separation [1] (separating various speakers), or music source separation [2, 3] (separating vocals, bass, and drums). For such tasks, a source-specific model is trained on dedicated datasets tailored to the task at hand. In contrast to source-specific separation tasks, universal source separation was recently proposed [4, 5], which consists of building source-agnostic models that are not constrained to a specific domain (like music or speech), and targets at separating an unknown number of sources given an arbitrary mix. However, existing universal source separation works predominantly focus on separating mixes similar to field recordings (with mostly sound events like dog barking or alarms). Further, most supervised learning methods for this task rely on small training sets [4–8]. For instance, the commonly-used FUSS dataset contains only 23 hours of single-source recordings [5]. Considering the number of different sounds in the world, most audio sources might be under-represented in such small datasets. Hence, the potential of universal source separation is yet to be fully explored because (i) most previous works separate mixes with predominantly sound events instead of simultaneously separating a broader set of sources including speech, music, and sound events, and (ii) supervised universal source separation models have never been trained with large-scale data.

Here, we explore training a unified model with large-scale data to address general audio source separation holistically[1], with the goal of separating any source from a given mix, including speech, music, and/or sound events. First, we scale up our audio source separation dataset by collecting 15,499 hours of recordings including speech, music, and sound events. Note that our dataset contains 3 orders of magnitude more data than FUSS [5], the commonly-used dataset for supervised learning (Table 1). Next, to investigate the feasibility of general audio source separation[1], we train 3 state-of-the-art models with our large and diverse dataset. We are also interested in the generalization capabilities of the trained models. Hence, in addition to evaluating the models on different partitions of the same dataset (in-distribution), we also evaluate them on 4 standard downstream test sets, each one representing a different use case with different data and mixing pipelines (out-of-distribution). While in some cases the out-of-distribution results are competitive, in some others the separation results are not as satisfactory. Finally, we show that out-of-distribution performance can be improved by fine-tuning the pre-trained general audio source separation models on each task.

To our best knowledge, we offer the first study on supervised general audio source separation at scale without prior knowledge about the sources. Previous works [9, 10] also consider speech and music in supervised universal source separation, but they assume the availability of a target source embedding to identify and separate the desired source from a mix. Unsupervised approaches can also leverage large scale (noisy) data, but they tend to underperform supervised methods [11–13]. Previous research also looked at the generalization capability of speech separation models [14, 15], but here we study the general audio source separation problem with a much more diverse set of out-of-distribution downstream tasks. Finally, our work is also conceptually similar to fine-tuning problem-agnostic self-supervised models [16], since we fine-tune source-agnostic audio separation models on source-specific tasks.

## 2. METHODOLOGY

### 2.1. Creating a Large-scale Source Separation Dataset

We collect recordings from public and licensed datasets to scale up general audio source separation with $\approx$ 1.9 M recordings of speech, music, and sound events. We mix recordings $r_k$ at various gains $g_k$:

$$m = \sum_{k=1}^{K} s_k = \sum_{k=1}^{K} g_k r_k,$$

where we normalize $r_k$'s amplitudes to 1 before mixing, and $K$ is the number of resulting sources $s_k$ present in the mix. Note that $K$ is assumed to be unknown during training/inference and, following common practice [5], we set $K \in \{1, 2, 3, 4\}$. Also, defining what constitutes a source is a significant challenge. We find that the

---

definition of "any recording with one source" might be impractical. For instance, considering separating two speakers talking in a cafeteria, it may be unnecessary to separate every individual sound in the background like the cutlery and the crowd noise. Similarly, in a mix with background music, it may not be desirable to separate out each instrument. In our view, incorporating low-volume, non-dominant background sounds as a single, combined source to be separated together could enhance the realism of the resulting mixes. Hence, to build our dataset, we rely on the following definition: "any recording with one source, except for low-volume background events that can contain one or more sources". We distinguish between foreground and background sources by simply applying higher gains to foreground sources. Table 1 presents those gains $g_k$, together with the number of collected recordings $r_k$ and their source types:

- **Speech foreground** is a multilingual collection of public and licensed clean speech recordings, each with 1 speaker. A large portion of the recordings we use are public: AVSpeech [17], VCTK [18], DAPS [19], and TIMIT [20].

- **Sound event foreground and background** are a combination of public and licensed datasets. The largest public dataset we use is (most of the content in) Freesound. Extensive listening finds that shorter Freesound recordings tend to be single-source, and longer ones tend to contain multiple sources. Hence, Freesound recordings shorter than 8 sec are used as foreground, and longer ones as background. We also use other background datasets, including: WHAM! [21] and DEMAND [22].

- **Music foreground and background** are a combination of public and licensed datasets. Public single-source datasets include: Slakh [23], ENST-drums [24], VocalSet [25], QMUL singing database [26], MUSIC [27], and EGFxSet [28]. Hence, foreground music mostly contains vocals, bass, drums, guitar, and keys, but also includes synthesizers, percussion, and classical instruments. Background music includes licensed music mixes.

Note that our collection is significantly larger than FUSS [5], the most common benchmark for universal source separation. After collecting our data, we define a set of rules to create the artificial mixes. These rules can be summarized into the following 3 upstream tasks:

- **Speech separation**. These mixes always contain at least 1 speech foreground source. Other sources are sampled from the following sets: speech foreground, sound events foreground/background, and music background to create mixes for speech denoising and speech source separation (from 1 to 4 speakers) use cases.

- **Sound event separation**. Sources are sampled from sound events foreground/background and music background to create mixes similar to previous universal source separation works [4, 5].

- **Music separation**. Sources are sampled from music foreground and sound events background to create mixes for music denoising and music source separation (from 1 to 4 sources) use cases.

Hence, to generate training data we randomly select: an upstream task (speech, sound event, or music separation with a probability of 0.25, 0.25, and 0.5, respectively), the number of sources $K$ (uniformly from 1 to 4), the recordings $r_k$ (which fragments and when they start in the mix), and the gains $g_k$ (sampled from a Beta distribution $Beta(2, 1)$ within the ranges in Table 1). We then downmix all the data to mono, zero-pad or truncate each sample to 8 sec, and resample them to 48 kHz. Note that our large-scale dataset covers various sampling rates and bandwidths, all resampled to 48 kHz, since we observe in preliminary experiments that models trained on this dataset perform competently at various (lower) sampling rates.

**Table 1**: Our large-scale general audio source separation dataset.

| Source type | $g_k$ (dB) | Single-source | # Recordings |
|---|---|---|---|
| Speech foreground | $[-10, 0]$ | ✓ | 759,397 |
| Sound event foreground | $[-10, 0]$ | ✓ | 314,652 |
| Sound event background | $[-20, -10]$ | ✗ | 398,360 |
| Music foreground | $[-3, 0]$ | ✓ | 75,639 |
| Music background | $[-20, -10]$ | ✗ | 379,565 |
| All dataset | | | 15,499 hours |
| FUSS [5] | | | 23 hours |

### 2.2. Models and Upstream Training

**TDANet-Wav** (10.8 M parameters). TDANet [1] is a state-of-the-art waveform-based speech source separation model based on an encoder-separator-decoder architecture. We adopt the official implementation[2] and increase the encoder dimension to 1024 and proportionally double the dimension of the separator layers.

**TDANet-STFT** (7.4 M parameters). We modify TDANet-Wav such that the encoder/decoder are replaced by STFT/iSTFT, and reuse the phase of the mixture for the iSTFT. The separator then outputs a mask over the STFT domain, not over a latent space as in TDANet-Wav. We use 32 and 8 ms frame length and stride, respectively. The bottleneck size is 384 and the separator layers follow the recommended ratio of feature maps with respect to the bottleneck size [1].

**BSRNN** (21.8 M parameters). Band-Split RNN is a powerful model for music source separation [3] and speech enhancement [29], also based on an encoder-separator-decoder architecture. Its encoder splits complex-valued STFT bins into bands and projects each band to a latent. We create 43 bands for our 48 kHz model, 2 more bands on top of the setup proposed for separating vocals from music at 44.1 kHz [3]. The separator consists of 12 interleaved band-level and sequence-level blocks with bidirectional LSTMs. The decoder undoes the band splitting and predicts complex-valued STFT masks. We adopt an available open-source implementation[3].

**IRM** (oracle). We compute the Ideal Ratio Mask (IRM) as an oracle upper bound using the magnitude STFT of the ground truth sources.

**Upstream training**. All models are trained on the upstream large-scale dataset for 10 M steps using the Adam optimizer with a batch size of 10 and a cyclical learning rate between $10^{-7}$ and $10^{-4}$ spanning 400 k steps per cycle. All models predict 4 sources $\hat{s}_k$ given a mix $m$. When there are fewer targets during training ($K<4$), the extra targets are set to zeros. Permutation invariant training [30] (PIT) aligns the predictions with the targets, and we minimize the logarithmic-MSE loss with a threshold $\tau$ set to $-30$ dB [5]:

$$\mathcal{L}(s_k, \hat{s}_k) = \begin{cases} 10 \log_{10} \left( \|\hat{s_k}\|^2 + \tau \|m\|^2 \right) & \text{if } s_k = 0, \\ 10 \log_{10} \left( \|s_k - \hat{s}_k\|^2 + \tau \|s_k\|^2 \right) & \text{otherwise.} \end{cases}$$

### 2.3. Evaluation Framework

**Upstream (in-distribution) evaluation**. For each upstream task (speech, sound event, and music separation), we set aside 3,000 mixes made of unseen recordings, which are sampled and mixed based on the same pipeline used for upstream training.

**Downstream (out-of-distribution) evaluation**. We study the generalization capabilities of our models with out-of-distribution datasets. We consider the following 4 downstream tasks:

---
[2] https://github.com/JusperLee/TDANet
[3] https://github.com/sungwon23/BSRNN

**Table 2**: Upstream (in-distribution) results for speech, sound event, and music separation. SI-SDR column: SI-SDRs/SI-SDRi (dB). US/ES/OS: source count rate (%).

| Task | Model | SI-SDR ↑ | US ↓ | ES ↑ | OS ↓ |
|---|---|---|---|---|---|
| Speech | TDANet-Wav | 53.5/**14.3** | **6.9** | **87.8** | 5.3 |
| | TDANet-STFT | 80.6/13.8 | 14.1 | 83.3 | **2.6** |
| | BSRNN | 44.3/12.8 | 13.8 | 80.1 | 6.1 |
| | IRM | 85.7/19.3 | 0 | 100 | 0 |
| Sound events | TDANet-Wav | 49.1/20.1 | 14.2 | 79.6 | 6.2 |
| | TDANet-STFT | 71.8/**22.1** | 17.8 | 78.0 | **4.2** |
| | BSRNN | 49.9/20.3 | **12.6** | **81.6** | 5.8 |
| | IRM | 78.0/28.3 | 0 | 100 | 0 |
| Music | TDANet-Wav | 52.6/14.6 | 5.9 | 90.9 | 3.2 |
| | TDANet-STFT | 80.8/14.6 | 9.1 | 89.1 | **1.8** |
| | BSRNN | 46.2/**18.2** | **3.9** | **93.2** | 2.9 |
| | IRM | 88.8/17.8 | 0 | 100 | 0 |

- **FUSS** is a universal source separation dataset with 1 to 4 sources, with mixes at 16 kHz similar to field recordings [5] (mostly sound events). We select the standard reverberated FUSS version for our downstream evaluation. Since FUSS is a subset of FSD50K [31], we exclude FSD50K from our upstream dataset.

- **Libri2Mix** is a common benchmark for speech source separation, with recordings at 16 kHz containing 2 clean speech sources [15]. All LibriSpeech [32] data is excluded from our upstream dataset.

- **DnR** dataset targets at separating cinematic mixes at 44.1 kHz into speech, music, and sound effects [33]. Again, all involved datasets in DnR are excluded from our upstream dataset. Also note that DnR is a particular out-of-distribution case because it violates our source definition. We expect our models to separate each speaker, musical sources, and sound effect sources unless the music and sound effects are low-volume background events. However, DnR separates a mix into 3 combined stems: speech (with all speakers), music (with all musical sources), and sound effects (all together).

- **MUSDB** is a music source separation dataset at 44.1 kHz with 4 sources: vocals, bass, drum, and 'other' [34]. Yet, note that our models are trained to separate more musical sources, including vocals, bass, drums, keys, guitar, synthesizers, and classical instruments. Further, the 'other' stem in MUSDB also violates our source definition, since such sources come grouped in one stem. We exclude both MUSDB and MedleyDB from our upstream data.

Although DnR (all stems) and MUSDB ('other' stem) violate our source definition, we are still interested in those to study fine-tuning a pre-trained (upstream) general model on a separation task defined differently. We conduct 3 evaluations for each downstream task:

- **No-tuning**. The pre-trained upstream models are assessed without any modification. This setup can also be seen as a zero-shot source separation case, where the models are pre-trained on a large dataset and then evaluated on new datasets without any adaptation.

- **Fine-tuning**. The pre-trained upstream models are fine-tuned on the new downstream task alone with PIT. This setup studies the upstream model as a general model that can be pre-trained on a large dataset and then fine-tuned on a new use case. When there are fewer training targets ($K<4$), the extra targets are set to zeros.

- **From-scratch**. The models are trained from-scratch on each downstream task. This setup studies the performance of the models when they are not pre-trained on a large dataset.

Note, however, that the downstream datasets have different sampling rates. To unify our evaluation framework, we upsample the mixes

**Table 3**: Downstream (out-of-distribution) results on FUSS. SI-SDR column: SI-SDRs/SI-SDRi (dB). US/ES/OS: source count rate (%).

| Evaluation | Model | SI-SDR ↑ | US ↓ | ES ↑ | OS ↓ |
|---|---|---|---|---|---|
| No-tuning | TDANet-Wav | 32.7/15.1 | 39.3 | 54.7 | **6.0** |
| | TDANet-STFT | 30.0/**16.4** | 38.6 | 55.0 | 6.4 |
| | BSRNN | 30.5/16.0 | **36.6** | **57.0** | 6.4 |
| Fine-tuning | TDANet-Wav | 33.2/17.7 | **11.8** | 77.5 | 10.7 |
| | TDANet-STFT | 34.0/18.1 | 16.5 | 73.1 | 10.4 |
| | BSRNN | 33.7/**18.6** | 14.0 | **78.5** | **7.5** |
| From-scratch | TDANet-Wav | 33.0/13.7 | 22.2 | 65.2 | 12.5 |
| | TDANet-STFT | 33.1/14.4 | 20.6 | 67.7 | **11.7** |
| | BSRNN | 32.4/**14.4** | **13.7** | **70.6** | 15.7 |
| SOTA Oracle | Postolache et al. [8] | 35.3/13.8 | 23.6 | 63.9 | 12.5 |
| | IRM | 39.9/25.3 | 0 | 100 | 0 |

and targets to 48 kHz. In that way, we can compute the loss against the upsampled targets when fine-tuning and training from-scratch. To compute metrics with the original ground truth, we downsample the predicted sources back to the original sampling rates. In preliminary experiments, we observe that models trained from-scratch and evaluated in this way yield similar results as those obtained by models trained on the original datasets without resampling.

**Evaluation metrics**. We use the standard metrics for each task:

- **SI-SDR** (dB) in DnR. We use scale-invariant signal-to-distortion ratio [35] (SI-SDR) to measure the quality of the separations.

- **SI-SDRs** (dB) in FUSS and upstream. For mixes with one source, we compute SI-SDRs = SI-SDR($s_k, \hat{s}_k$) = SI-SDR($m, \hat{s}_k$) [5], since with one-source mixes the goal is to bypass the mix. The 's' sub-index stands for single-source.

- **SI-SDRi** (dB) in FUSS, Libri2Mix, and upstream. For mixes with 2 to 4 sources, we report SI-SDRi = SI-SDR($s_k, \hat{s}_k$) − SI-SDR($s_k, m$) [5, 8]. The 'i' sub-index stands for improvement. To account for inactive sources, estimate-target pairs that have silent target sources are discarded.

- **US, ES, OS** (%) in FUSS and upstream. Note that our models implicitly count the number of sources to separate. To evaluate source counting, we compute the proportion of the samples for which the number of nonzero predictions are fewer than (under-separation, US), equal to (equal-separation, ES), or more than (over-separation, OS) the number of nonzero targets [5]. A prediction is considered nonzero if its average energy is above $-20$ dB relative to the softest nonzero target source [5].

- **SDR** (dB) in MUSDB. Defined in [36], is the per-source median across the median SDR over all 1 second chunks in each song.

## 3. RESULTS

Separations produced by our models are available on our website[4].

### 3.1. Upstream (In-distribution) Evaluation

Table 2 lists the results for the 3 in-distribution tasks, showing that it is possible, with a single model, to perform general audio source separation (including speech, sound events, and music) without prior knowledge about the source types and the number of sources (up to 4). Comparing with the IRM, we see that the models are competitive. Interestingly, each model stands out at a different task: TDANet-Wav for speech separation, TDANet-STFT for sound event

---

[4]http://www.jordipons.me/apps/GASS

**Table 4**: Downstream (out-of-distribution) SI-SDRi (dB, ↑) results on Libri2Mix: speech source separation of 2 speakers.

|  | No-tuning | Fine-tuning | From-scratch |
|---|---|---|---|
| TDANet-Wav | **11.4** | **17.9** | **17.5** |
| TDANet-STFT | 9.5 | 13.3 | 12.7 |
| BSRNN | 8.7 | 16.0 | 15.2 |
| Li et al. [1] (SOTA) | - | - | 17.4 |
| IRM (Oracle) | 13.3 | - | - |
| TDANet-Wav-FUSS | −6.6 | - | - |

**Table 5**: Downstream (out-of-distribution) SI-SDR (dB, ↑) results on DnR for speech (S), music (M), and sound effects (FX).

|  | No-tuning | | | Fine-tuning | | | From-scratch | | |
|---|---|---|---|---|---|---|---|---|---|
|  | S | M | FX | S | M | FX | S | M | FX |
| TDANet-Wav | **8.1** | **0.6** | **−0.7** | **14.8** | 6.0 | 7.7 | **14.4** | 5.6 | 7.1 |
| TDANet-STFT | 7.7 | −1.9 | −1.4 | 13.1 | 5.4 | 7.0 | 12.9 | 4.8 | 6.5 |
| BSRNN | 7.9 | 0.3 | −1.5 | 14.4 | **6.5** | **7.9** | 14.0 | **6.0** | **7.4** |
| Unprocessed mixes | 1.0 | −6.8 | −5.0 | - | - | - | - | - | - |
| Petermann et al. [33] | - | - | - | - | - | - | 12.3 | 4.2 | 5.7 |
| IRM (Oracle) | 15.6 | 8.5 | 10.7 | - | - | - | - | - | - |

separation, and BSRNN for music separation. BSRNN outperforms IRM for music separation, showing the advantage of operating on the complex STFT for this task. Also, the relatively high equal-separation rates (ES) show that the models are often able to count/separate the sources correctly. Among the miscounting cases, models tend to under-separate (US). Finally, the high SI-SDRs values show that the models are able to bypass single-source inputs.

### 3.2. Downstream (Out-of-distribution) Evaluation

**FUSS** (Table 3). First, the no-tuning SI-SDRi results consistently outperform those of the models trained from-scratch. This reflects that, for FUSS, the upstream models are capable to generalize. Yet, the under-separation rates are higher for the no-tuning models. We hypothesize that allowing low-volume multi-source backgrounds in the upstream tasks could cause under-separation in FUSS. After fine-tuning, however, we improve both the source counting accuracy and SI-SDRi, denoting how transferable to FUSS the upstream models are. Note that the fine-tuned models are also significantly better than the state-of-the-art.

**Libri2Mix** (Table 4). We observe that the best no-tuning model approaches the IRM result, denoting that the upstream model can generalize to the Libri2Mix task. Yet, there is still a gap when compared to the models trained from-scratch, which is not surprising considering the more general upstream task we address with the same model capacity. Fine-tuning always outperforms training from-scratch, but the improvements are much smaller if compared to those obtained by fine-tuning on FUSS (Table 3). This shows that the no-tuning performance is indicative of how transferable the models are between tasks. Nonetheless, the fine-tuned TDANet-Wav obtains state-of-the-art results. We also evaluate a TDANet-Wav trained on FUSS (TDANet-Wav-FUSS) data to study the capacity of FUSS as an upstream dataset. Its failure denotes the limitations of current supervised universal source separation to separate an arbitrary mix.

**DnR** (Table 5). The upstream models with no-tuning do not perform competently on this downstream task that violates our source definition, since DnR aims at '3-group' separation but our models are trained to separate each source. For this reason, we observe high over-separation rates (≈ 95% of the time, the 4th output contains

**Table 6**: Downstream (out-of-distribution) SDR (dB, ↑) results on MUSDB for vocals (V), bass (B), drums (D), and other (O).

| Evaluation | Model | V | B | D | O | Avg |
|---|---|---|---|---|---|---|
| No-tuning | TDANet-Wav | **1.2** | 0.4 | 2.2 | 0.2 | 1.0 |
|  | TDANet-STFT | 0.7 | **1.1** | **3.2** | **0.3** | **1.3** |
|  | BSRNN | 0.0 | −0.2 | 0.0 | 0.0 | 0.0 |
| No-tuning | TDANet-Wav-M | 3.4 | 1.1 | 4.9 | 0.1 | 2.4 |
| Fine-tuning | TDANet-Wav | 7.0 | **9.6** | **9.8** | 4.7 | **7.8** |
|  | TDANet-STFT | 6.8 | 5.9 | 6.5 | 4.4 | 5.9 |
|  | BSRNN | **8.6** | 7.7 | 8.3 | **5.3** | 7.5 |
| From-scratch | TDANet-Wav | 6.5 | **9.6** | **9.6** | 4.7 | **7.6** |
|  | TDANet-STFT | 6.7 | 6.1 | 6.4 | 4.1 | 5.8 |
|  | BSRNN | **8.2** | 7.4 | 8.4 | **5.2** | 7.3 |
| From-scratch | BSRNN w/ PIT | 8.3 | 7.2 | 8.4 | 5.1 | 7.2 |
| SOTA | Luo & Yu [3] | 10.0 | 7.2 | 9.0 | 6.7 | 8.2 |
| Oracle | IRM | 9.4 | 7.1 | 8.5 | 7.9 | 8.2 |

non-negligible predictions). When comparing no-tuning results with the unprocessed mixes, one notes that the no-tuning models are able to perform some degree of separation, but are much worse than the ones trained from-scratch. However, the fine-tuned models perform better than the ones trained from-scratch, indicating the transferability of the upstream models to a differently-defined task. Note that we outperform Petermann et al. [33] (the best published result on DnR).

**MUSDB** (Table 6). On MUSDB, all no-tuning models perform poorly. First, we hypothesize that PIT may cause this problem, since we are not aware of prior works using PIT for music source separation. Hence, we compare two BSRNN models trained from-scratch with and without PIT[5], to find out that their results are comparable and PIT is not the problem. Next, we train a TDANet-Wav with only upstream musical mixes to study the generalization capabilities of this music-specific model (TDANet-Wav-M). Yet, despite improving upon the general model, this model still performs much worse than the from-scratch TDANet-Wav. This fact, combined with the good in-distribution performance of our general models (Table 2), suggests that we have a mismatch between the upstream musical mixes and MUSDB mixes. Overall, these observations suggest future investigations, including collecting more music foreground data (note in Table 1 that we only collected 75,639 music foreground recordings), increasing the model capacity, and probing interference between different upstream tasks.

## 4. CONCLUSION

We studied general audio source separation models trained in a supervised fashion with large-scale data. To study their generalization capabilities, we evaluated both in- and out-of-distribution performance. The in-distribution results show that the models are able to separate an unknown number of sources from a variate set of mixes that include speech, music, and sound events. Among the out-of-distribution results, the no-tuning models achieved competitive performance for sound event and speech separation, but we also noted that our models had challenges for generalizing to separate cinematic and music mixes. Moreover, with fine-tuning consistently outperforming from-scratch, we show how transferable the upstream models are to a diverse set of downstream tasks, even when there is a mismatch between the source definitions of the upstream and downstream tasks. All fine-tuned models (except the music separation one) obtain state-of-the-art results in their respective benchmarks.

---

[5]Our BSRNNs did not achieve the results as in Luo & Yu [3] because we used a much smaller model and a single banding structure for all sources.

# 5. REFERENCES

[1] K. Li, R. Yang, and X. Hu, "An Efficient Encoder-Decoder Architecture with Top-down Attention for Speech Separation," in *ICLR*, 2023.

[2] S. Rouard, F. Massa, and A. Défossez, "Hybrid Transformers for Music Source Separation," in *ICASSP*, 2023, pp. 1–5.

[3] Y. Luo and J. Yu, "Music Source Separation With Band-Split RNN," *IEEE/ACM TASLP*, 2023.

[4] I. Kavalerov, S. Wisdom, H. Erdogan, B. Patton, K. Wilson, J. Le Roux, and J. R. Hershey, "Universal Sound Separation," in *WASPAA*, 2019, pp. 175–179.

[5] S. Wisdom, H. Erdogan, D. P. W. Ellis, R. Serizel, N. Turpault, E. Fonseca, J. Salamon, P. Seetharaman, and J. R. Hershey, "What's All the Fuss About Free Universal Sound Separation Data?," in *ICASSP*, 2021, pp. 186–190.

[6] E. Tzinis, Z. Wang, X. Jiang, and P. Smaragdis, "Compute and Memory Efficient Universal Sound Source Separation," *Journal of Signal Processing Systems*, pp. 245–259, 2022.

[7] E. Tzinis, S. Wisdom, J. R. Hershey, A. Jansen, and D. P. W. Ellis, "Improving Universal Sound Separation Using Sound Classification," in *ICASSP*, 2020, pp. 96–100.

[8] E. Postolache, J. Pons, S. Pascual, and J. Serrà, "Adversarial Permutation Invariant Training for Universal Sound Separation," in *ICASSP*, 2023, pp. 1–5.

[9] Q. Kong, K. Chen, H. Liu, X. Du, T. Berg-Kirkpatrick, S. Dubnov, and M. D. Plumbley, "Universal Source Separation with Weakly Labelled Data," *arXiv:2305.07447*, 2023.

[10] X. Liu, Q. Kong, Y. Zhao, H. Liu, Y. Yuan, Y. Liu, R. Xia, Y. Wang, M. D. Plumbley, and W. Wang, "Separate Anything You Describe," *arXiv:2308.05037*, 2023.

[11] S. Wisdom, A. Jansen, R. J. Weiss, H. Erdogan, and J. R. Hershey, "Sparse, Efficient, and Semantic Mixture Invariant Training: Taming In-the-wild Unsupervised Sound Separation," in *WASPAA*, 2021, pp. 51–55.

[12] H. Dong, N. Takahashi, Y. Mitsufuji, J. McAuley, and T. Berg-Kirkpatrick, "CLIPSep: Learning Text-queried Sound Separation with Noisy Unlabeled Videos," *ICLR*, 2023.

[13] S. Wisdom, E. Tzinis, H. Erdogan, R. Weiss, K. Wilson, and J. Hershey, "Unsupervised Sound Separation Using Mixture Invariant Training," *NeurIPS*, vol. 33, pp. 3846–3857, 2020.

[14] B. Kadıoğlu, M. Horgan, X. Liu, J. Pons, D. Darcy, and V. Kumar, "An Empirical Study of Conv-TasNet," in *ICASSP*, 2020, pp. 7264–7268.

[15] J. Cosentino, M. Pariente, S. Cornell, A. Deleforge, and E. Vincent, "LibriMix: An Open-Source Dataset for Generalizable Speech Separation," *arXiv:2005.11262*, 2020.

[16] W. Hsu, B. Bolte, Y. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units," *IEEE/ACM TASLP*, vol. 29, pp. 3451–3460, 2021.

[17] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein, "Looking to Listen at the Cocktail Party: A Speaker-Independent Audio-Visual Model for Speech Separation," in *SIGGRAPH*, 2018.

[18] C. Veaux, J. Yamagishi, and K. MacDonald, "CSTR VCTK Corpus: English Multi-speaker Corpus for CSTR Voice Cloning Toolkit," 2017.

[19] G. J. Mysore, "Can We Automatically Transform Speech Recorded on Common Consumer Devices in Real-World Environments into Professional Production Quality Speech?—A Dataset, Insights, and Challenges," *IEEE Signal Processing Letters*, vol. 22, no. 8, pp. 1006–1010, 2014.

[20] J. S. Garofolo et al., "TIMIT Acoustic Phonetic Continuous Speech Corpus," *Linguistic Data Consortium*, 1993.

[21] G. Wichern, J. Antognini, M. Flynn, L. R. Zhu, E. McQuinn, D. Crow, E. Manilow, and J. Le Roux, "WHAM!: Extending Speech Separation to Noisy Environments," in *INTERSPEECH*, 2019, pp. 1368–1372.

[22] J. Thiemann, N. Ito, and E. Vincent, "The Diverse Environments Multi-Channel Acoustic Noise Database (DEMAND): A Database of Multichannel Environmental Noise Recordings," in *Meetings on Acoustics*, 2013.

[23] E. Manilow, G. Wichern, P. Seetharaman, and J. Le Roux, "Cutting Music Source Separation Some Slakh: A Dataset to Study the Impact of Training Data Quality and Quantity," in *WASPAA*, 2019, pp. 45–49.

[24] O. Gillet and G. Richard, "ENST-Drums: An Extensive Audio-Visual Database for Drum Signals Processing," in *ISMIR*, 2006.

[25] J. Wilkins, P. Seetharaman, A. Wahl, and B. Pardo, "VocalSet: A Singing Voice Dataset," in *ISMIR*, 2018, pp. 468–474.

[26] D. A. Black, M. Li, and M. Tian, "Automatic Identification of Emotional Cues in Chinese Opera Singing," in *ICMPC*, 2014.

[27] H. Zhao, C. Gan, A. Rouditchenko, C. Vondrick, J. McDermott, and A. Torralba, "The Sound of Pixels," in *ECCV*, 2018, pp. 570–586.

[28] H. Pedroza, G. Meza, and I. R. Roman, "EGFxSet: Electric Guitar Tones Processed Through Real Effects of Distortion, Modulation, Delay and Reverb," in *ISMIR*, 2022.

[29] J. Yu, H. Chen, Y. Luo, R. Gu, and C. Weng, "High Fidelity Speech Enhancement with Band-split RNN," in *INTERSPEECH*, 2023, pp. 2483–2487.

[30] M. Kolbæk, D. Yu, Z. Tan, and J. Jensen, "Multitalker Speech Separation with Utterance-level Permutation Invariant Training of Deep Recurrent Neural Networks," *IEEE/ACM TASLP*, vol. 25, no. 10, pp. 1901–1913, 2017.

[31] E. Fonseca, X. Favory, J. Pons, F. Font, and X. Serra, "FSD50K: An Open Dataset of Human-Labeled Sound Events," *IEEE/ACM TASLP*, vol. 30, pp. 829–852, 2021.

[32] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR Corpus Based on Public Domain Audio Books," in *ICASSP*, 2015, pp. 5206–5210.

[33] D. Petermann, G. Wichern, Z. Wang, and J. Le Roux, "The Cocktail Fork Problem: Three-Stem Audio Separation for Real-world Soundtracks," in *ICASSP*, 2022, pp. 526–530.

[34] Z. Rafii, A. Liutkus, F. Stöter, S. I. Mimilakis, and R. Bittner, "MUSDB18-HQ - An Uncompressed Version of MUSDB18," 2019.

[35] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR–Half-Baked or Well Done?," in *ICASSP*, 2019, pp. 626–630.

[36] F. Stöter, A. Liutkus, and N. Ito, "The 2018 Signal Separation Evaluation Campaign," in *LVA/ICA*, 2018, pp. 293–305.