## SELF-SUPERVISED AUDIO-VISUAL CO-SEGMENTATION

Andrew Rouditchenko\*1, Hang Zhao\*1, Chuang Gan², Josh McDermott1, Antonio Torralba1

# <sup>1</sup>MIT <sup>2</sup>MIT-IBM Watson AI Lab

{roudi, hangzhao, jhm, torralba}@mit.edu, ganchuang1990@gmail.com

#### **ABSTRACT**

Segmenting objects in images and separating sound sources in audio are challenging tasks, in part because traditional approaches require large amounts of labeled data. In this paper we develop a neural network model for visual object segmentation and sound source separation that learns from natural videos through self-supervision. The model is an extension of recently proposed work that maps image pixels to sounds [1]. Here, we introduce a learning approach to disentangle concepts in the neural networks, and assign semantic categories to network feature channels to enable independent image segmentation and sound source separation after audio-visual training on videos. Our evaluations show that the disentangled model outperforms several baselines in semantic segmentation and sound source separation.

*Index Terms*— audio-visual, co-segmentation, disentangled, self-supervised, source separation

#### 1. INTRODUCTION

Semantic segmentation of images [2, 3] and sound source separation in audio [4, 5, 6, 7] are two important and popular tasks in the computer vision and computational audition communities. Traditional approaches have relied on large, labeled datasets, but recent work has leveraged the natural correspondence between vision and sound to apply supervised learning without explicit labels. One approach is to use the signal or features from one modality to supervise the other. For example, [8] used visual features to supervise the learning of audio networks, and [9] used sound signals as supervision to train vision networks. Other models used sound and vision to jointly supervise each other in order to localize visual objects that make sound [10, 11, 12, 13], and to explore the relationship between unlabelled speech and visual input [14, 15]. More recently, [1, 16, 11] used audio-visual correspondence to separate sound sources. Another key direction is to design cross-modal representations that are robust and interpretable [17, 18]. Our contribution is to develop a model for audio-visual co-segmentation using videos.

In the Mix-and-Separate framework proposed in [1], neural networks are trained on videos through self-supervision to perform image segmentation and sound source separation. However, following training, the model could only be applied to videos with synchronized audio.

Here we seek to enable a system that can perform segmentation and separation tasks using test input containing only video frames or sound mixtures. We introduce a learning approach that disentangles concepts learned by neural networks, enabling independent inference of images and audio mixtures without needing to combine visual and auditory features. The proposed learning approach relies on an activation function schedule that uses the sigmoid activation function during the training stage and a softmax activation during the fine-tuning stage, producing sparse activations that could correspond to semantic categories in the input. Following learning, semantic categories are assigned to intermediate network feature channels using labels available in the training dataset. Given a video frame or sound excerpt, the category-to-feature-channel correspondence can be used to select a particular type of source or object for resynthesis or segmentation. The disentangling thus enables both independent inference and model interpretability because the network feature channels respond sparsely to semantic concepts.

We evaluate performance on image-only and audio-only tasks, which was not possible using the previous model. Furthermore, we substantially extend the scale of previous work [1] by training on a video dataset of naturally occurring audio-visual events with over 4000 videos [19]. The results show that we can achieve promising semantic segmentation and source source separation performance.

## 2. APPROACH

#### 2.1. Self-supervised Cross-modal Training

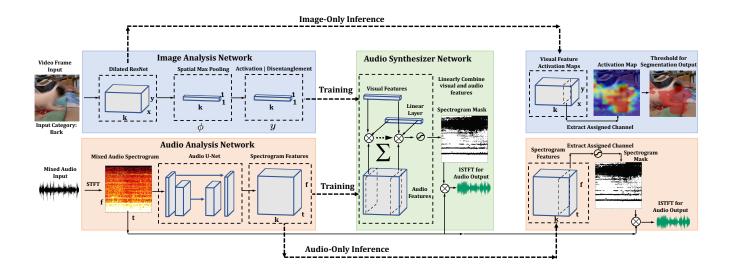
Our approach adopts the Mix-and-Separate framework used in [1], which first generates a synthetic sound separation training set by mixing the audio signals from two different videos, and then trains a neural network to separate the audio mixture conditioned on the visual input corresponding to one of the audio signals. Critically, although the neural network is trained in a supervised fashion, it does not require labeled data. Thus the training pipeline can be considered as self-supervised learning.

As shown in Fig. 1, the framework we use consists of three components: an image analysis network, an audio analysis network, and an audio synthesizer network. During training, we randomly select two videos to form a synthetic training example. The image analysis network extracts visual features on a video frame from one of the videos, then uses spatial max pooling to compress the features into a visual feature vector. The audio signals of the two videos are summed and the spectrogram is extracted. An audio analysis network then processes the mixture spectrogram into audio features, where each channel contains features of different components of the input sound mixture. Finally, an audio synthesizer network combines the visual and audio features to predict a spectrogram mask to generate the audio signal for the selected video.

#### 2.2. Disentangling Internal Representations

The models in [1] rely on synchronized video and audio as input, and thus can only perform joint audio-visual source separation, limiting their use in real applications where synchronized data are not available. Here we aim to use the image analysis network and audio analysis network independently after training, without needing the audio synthesizer network to combine the visual and audio features.

<sup>\*</sup> These authors contributed equally to this work.



**Fig. 1.** Joint audio-visual training and independent image and audio inference. After training on synthetic mixtures of videos and assigning the dataset categories to network feature channels, the image analysis network performs image-only segmentation and the audio analysis network performs audio-only source separation.

Specifically, we design a learning schedule to disentangle the learned internal representations before the audio synthesizer network combines audio and visual features. Disentanglement is a method to create interpretable representations that enhance functionality and that make feature channels more robust to changes in other units [20]. As shown in Fig. 1, the outputs of both the image and audio analysis networks have K channels, where K is larger than the number of dataset categories. Ideally, each channel would correspond to a separate concept and each input category would uniquely activate one channel. We attempt to achieve this with a learning schedule that causes the intermediate feature representations before audio-visual fusion to be sparse.

Our technique is inspired from [21], who studied the effects of annealing the temperature parameter in a softmax activation function in order to push output activations towards one-hot vectors. As the temperature parameter T in the softmax activation function changes from high to low, the shape of the output distribution changes from uniform to one-hot:

$$y_k = \frac{\exp(\frac{\phi_k}{T})}{\sum_{i=1}^n \exp(\frac{\phi_i}{T})},\tag{1}$$

where  $y_k$  is the value of the  $k_{th}$  visual feature channel after activation, T is the temperature, and  $\phi_i$  is the value of the  $i_{th}$  visual feature channel before activation. We used this idea by first training the model without imposing sparsity in the features, and then gradually changing the hyperparameters to encourage sparse and disentangled representations. The model is initially trained using a sigmoid activation on the visual feature vector  $\phi$ , which leads to diverse activations and helps with convergence to an initial solution. The sigmoid activation is then replaced with the softmax activation, and the temperature is gradually decreased, pushing the visual feature vector toward a one-hot vector, and causing the visual and audio feature representations to become sparse and disentangled. Note that we did not sample from the Gumbel distribution as described in [21], but instead incorporated the decaying temperature parameter in the softmax activation.

## 2.3. Category to Channel Assignment for Co-segmentation

After training the networks without labels, we then use the category labels in the dataset to match categories to network feature channels, so that a particular type of source or object can be selected for resynthesis or segmentation. We use the validation set to compute the visual feature vector for each video and make a normalized table of these activations, which represents the cost of assigning each dataset category to each network feature channel. We then use a matching algorithm for the linear sum assignment problem [23], which assigns each dataset category to a network feature channel. For example, the dataset category, "cars," could correspond to the first network channel, "male speech" to the second network channel, and so forth. We can measure the validity of the assignment via classification accuracy: we measure what percentage of the input video frames in the validation set activate their assigned channel most strongly. These results are reported in Sec. 4.

The assignment of input categories to network feature channels allows independent image and audio processing without needing the audio-synthesizer network to combine the features. In principle, one could select any activated channel and resynthesize its source signal or segmentation. Here, for evaluation purposes, we simply use the channel corresponding to the video's label in the dataset.

For object segmentation, the last spatial max pooling layer of the image analysis network is removed to preserve activation feature maps instead of a visual feature vector. Given an input video frame, the activation map in the channel assigned to the video's category is selected, upsampled to the input size, and thresholded to obtain a predicted segmentation.

Given an audio mixture, the audio analysis network outputs spectrogram features in K channels. The channels assigned to the two source video categories are selected, and used as a spectrogram mask for the respective source. Each spectrogram mask is then applied to the mixture spectrogram in order to separate the corresponding sound source from the mixture.

| Model Name               | Learning Schedule |               |            |              |             |      | Classification |
|--------------------------|-------------------|---------------|------------|--------------|-------------|------|----------------|
|                          | Softmax Epochs    | Initial Temp. | Decay Rate | Decay Epochs | Final Temp. |      |                |
| Baseline-Sigmoid Only    | -                 | -             | -          | -            | -           | 0.38 | 6.30%          |
| Baseline-Softmax Only    | 25                | 1.0           | 0.3        | 10, 20       | 0.090       | 0.99 | 38.7%          |
| Sigmoid & Softmax A      | 20                | 10.0          | 0.5        | 4, 8, 12, 16 | 0.625       | 0.93 | 18.1%          |
| Sigmoid & Softmax B      | 20                | 1.5           | 0.75       | 4, 8, 12, 16 | 0.475       | 0.97 | 37.1%          |
| Sigmoid & Softmax C      | 25                | 1.0           | 0.3        | 4, 8         | 0.090       | 0.99 | 40.3%          |
| Sigmoid & Softmax D      | 25                | 1.0           | 0.3        | 3, 6, 9, 12  | 0.008       | 0.99 | 24.0%          |
| Sigmoid & Softmax E      | 25                | 1.0           | 0.5        | 5, 10, 15    | 0.125       | 0.99 | 45.9%          |
| ResNet-18 Features & SVM | -                 | -             | -          | -            | -           | _    | 68.4%          |

**Table 1**. Classification performance and activation sparsity for the proposed model with different learning schedules and baselines. Decay Epochs indicates the epochs at which the temperature was decayed.

| Model Name                   | Sound S | eparation | Seman. Seg. |
|------------------------------|---------|-----------|-------------|
|                              | SDR     | SIR       | IoU         |
| Baseline-Sigmoid Only        | 0.865   | 6.04      | 0.204       |
| Baseline-Softmax Only        | 0.172   | 3.37      | 0.207       |
| Sigmoid & Softmax A          | -0.536  | 4.52      | 0.112       |
| Sigmoid & Softmax B          | 0.341   | 6.23      | 0.152       |
| Sigmoid & Softmax C          | 0.716   | 6.21      | 0.232       |
| Sigmoid & Softmax D          | -1.88   | 2.82      | 0.205       |
| Sigmoid & Softmax E          | 1.03    | 6.37      | 0.225       |
| Nonnegative Matrix Fact. [5] | 0.196   | 3.94      | -           |
| Class Activation Maps [22]   | -       | -         | 0.190       |

**Table 2**. Quantitative sound separation and semantic segmentation performance.

## 3. EXPERIMENTAL SETUP

#### 3.1. Models

The video analysis network is a dilated variation of the ResNet-18 model [24]. The dilated convolutions preserve larger visual feature activation maps, which are used after training for image segmentation. For an input video frame with size  $224 \times 224$ , it outputs K output activation maps of size  $14 \times 14$ . Spatial max pooling is then applied to compress the visual features into a visual feature vector with K channels.

The audio analysis network is a modified U-Net [25] architecture. It has 7 down-convolution layers and 7 up-convolution layers with skip connections in between. For an input audio spectrogram with size  $256\times256,$  it outputs K spectrogram feature maps of size  $256\times256.$ 

The audio synthesizer network is a linear layer that is applied to combine the audio and visual features into a spectrogram mask that is multiplied element-wise with the input spectrogram. The inverse STFT is applied to the predicted magnitude spectrogram with the phase of the input spectrogram to recover the waveform. The network outputs could be either binary or floating point masks, and we chose to use binary masks with a per pixel cross entropy loss.

#### 3.2. Dataset

To train our models on a diverse set of audio-visual events, we used the Audio-Visual Event (AVE) dataset, containing 4143 videos covering 28 event categories [19]. The dataset spans categories such as cars, musical instruments, and speech, thus offering a collection that is wider in scope than other audio-visual datasets, such those limited to speech or instruments. The dataset is divided into the following splits: training (3339 videos), validation (402 videos), and test (402 videos).

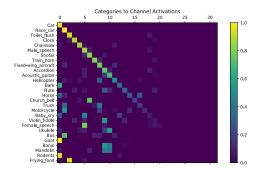


Fig. 2. Channel activations correspond to categories for the best performing model.

The dataset was preprocessed to extract video sections containing audio-visual events, in which the sound source is visible and the sound it produces is audible. The videos were cropped into 6-second clips, and the video frames were downsampled to 2 frames per second. The audio signals were resampled to 11kHz and converted to spectrograms via the Short-Time Fourier Transform (STFT). The STFT used a window size of 1022 samples and a hop length of 256 samples, which resulted in  $512 \times 256$  time-frequency spectrograms. For efficient model training, these spectrograms were further resampled on a on a log-frequency scale to obtain  $256 \times 256$  spectrograms, which is similar to applying mel-frequency spacing.

## 3.3. Activation Learning Schedule and Sparsity

The learning schedule to produce sparse visual feature activations was implemented with two stages: training and fine-tuning. The training stage used a fixed sigmoid activation function and the fine-tuning stage used a softmax activation function with custom schedules for the temperature parameter. The custom schedules varied the initial temperature, the number of epochs for fine-tuning, the decay rate, and decay epochs, which proved to be important. Besides the decaying temperature, the learning rates were also reduced by a factor of five in the fine-tuning stage. We used a measure of sparsity from computational neuroscience [26] to evaluate the sparsity of our model activations:

$$S(\mathbf{x}) = \frac{1 - (\frac{\mathbf{x} \cdot \mathbf{u}}{\|\mathbf{x}\| \|\mathbf{u}\|})^2}{1 - 1/K},$$
 (2)

where  $\mathbf{x}$  is the channel activation,  $\mathbf{u}$  is the uniform distribution, and K is the number of channels in the model. We measure the sparsity of the visual feature vector after activation to quantify the extent of disentanglement.

#### 4. EXPERIMENTAL RESULTS

#### 4.1. Disentanglement and Classification

In Table 1, we show the performance of the proposed model ("Sigmoid & Softmax") with several different hyperparameter settings, as well as of the baseline models. To measure the extent of disentanglement, we evaluated the visual feature vector sparsity and classification performance on the AVE validation set. A random search over the hyperparameters was conducted to find the best performing models. We also tested different numbers of channels, K, and found 32 to work well and train efficiently. There are 28 categories in the AVE dataset, such that 32 channels is enough to match each category with a channel and to have extra channels for content that is not accounted for by the categories, such as silence or noise.

The proposed model with the best hyperparameter setting, "Sigmoid & Softmax E," achieved a classification accuracy of 45.9%, significantly higher than the baseline variants of the model. To compare this result with a supervised baseline, we also trained a linear SVM on features from a ResNet-18, pre-trained on ImageNet. Although this supervised baseline achieves a higher classification accuracy of 68.9%, its features result from label supervision, potentially enabling fine-grained distinctions not possible using self-supervised learning. Moreover, our model has a much smaller feature vector, to enable the selection of discrete sources.

A qualitative evaluation of the performance of the best model is shown in Fig. 2, which shows how the visual feature channels activate for different input categories. Generally speaking, each visual input category only activates one or a few channels. Some channels respond to semantically related categories, indicating that the misclassifications mostly arise due to relatively fine-grained confusions.

The hyperparameters in the softmax fine-tuning stage proved to be important to achieve disentanglement. The softmax activation is necessary for the activations to become sparse, as shown by the low sparsity measurement from the "Sigmoid Only" model. The best schedule turned out to be a gradual decay from an initial temperature of 1 to about 0.01, as indicated by model "Sigmoid & Softmax E." The results show that ending with a temperature too high or too low can lead to suboptimal performance.

## 4.2. Source Separation

In the previous version of the model [1], source separation was only possible given synchronized audio-visual input because the network's representations of audio and video were entangled. By contrast, the current model can perform audio-only tasks following training because the sparse activations lead network feature channels to tend to correspond to semantic categories. We conducted audio-only sound source separation on the AVE test set, and show quantitative results in Table. 2 and qualitative results in Fig. 3. The Signal to Distortion Ratio (SDR) and the Signal to Interference Ratio (SIR) are two commonly used sound source separation metrics [27], and were calculated using the mir-eval library [28]. We include a baseline approach of Nonnegative Matrix Factorization [5]. The model which achieved the highest classification accuracy, "Sigmoid & Softmax E," also achieved the highest SDR and SIR. Qualitatively, the model succeeds in separating the sound from different sources to a large extent, which is visible in the source spectrogram recovery.

## 4.3. Semantic Segmentation

The current model can now perform vision-only tasks following training, without the fusion of visual and audio features as in the

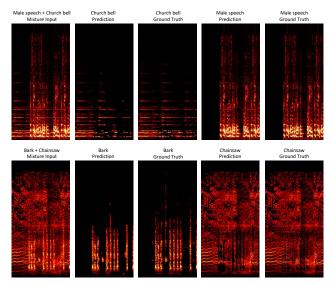


Fig. 3. Source separation results from the audio analysis network.



Fig. 4. Object segmentation results from the image analysis network.

previous version of the model [1]. To quantitatively evaluate the segmentation results, we labelled the middle video frame in each video from the AVE test set with polygons corresponding to the objects making sounds in the videos. The quantitative results, measured by Intersection over Union (IoU), are shown in Table 2 and qualitative results are shown Fig. 4. We include a baseline approach of Class Activation Mapping [22], which is a weakly supervised method used for object localization. The best semantic segmentation performance was achieved by "Sigmoid & Softmax C," but the version with the highest classification performance, "Sigmoid & Softmax E," performed nearly as well. As evident in Fig. 4, the boundaries of the predicted masks were often imperfect. This could result from the low resolution of the activation maps and the weak supervision used during training.

#### 5. CONCLUSION

We developed a self-supervised audio-visual co-segmentation approach to segment visual objects and separate sound sources. The approach relied on training networks for source separation through self-supervision on a large dataset of videos. We propose a method for learning disentangled feature representations and an assignment of dataset categories to network feature channels that enables independent image segmentation and sound source separation. Experimental results on the AVE dataset show that our approach achieves promising results on semantic segmentation and source separation.

#### 6. REFERENCES

- [1] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba, "The sound of pixels," in *The European Conference on Computer Vision (ECCV)*, September 2018.
- [2] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba, "Scene parsing through ADE20K dataset," in *Proc. CVPR*, 2017.
- [3] Jonathan Long, Evan Shelhamer, and Trevor Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [4] Albert S Bregman, Auditory scene analysis: The perceptual organization of sound, MIT press, 1994.
- [5] Tuomas Virtanen, "Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria," *IEEE transactions on audio, speech, and language processing*, vol. 15, no. 3, pp. 1066–1074, 2007.
- [6] Josh H McDermott, "The cocktail party problem," *Current Biology*, vol. 19, no. 22, pp. R1024–R1027, 2009.
- [7] John R Hershey, Zhuo Chen, Jonathan Le Roux, and Shinji Watanabe, "Deep clustering: Discriminative embeddings for segmentation and separation," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 31–35.
- [8] Yusuf Aytar, Carl Vondrick, and Antonio Torralba, "Soundnet: Learning sound representations from unlabeled video," in Advances in Neural Information Processing Systems, 2016, pp. 892–900.
- [9] Andrew Owens, Jiajun Wu, Josh H McDermott, William T Freeman, and Antonio Torralba, "Ambient sound provides supervision for visual learning," in *European Conference on Computer Vision*. Springer, 2016, pp. 801–816.
- [10] Relja Arandjelovic and Andrew Zisserman, "Look, listen and learn," in 2017 IEEE International Conference on Computer Vision (ICCV). IEEE, 2017, pp. 609–617.
- [11] Andrew Owens and Alexei A Efros, "Audio-visual scene analysis with self-supervised multisensory features," *arXiv preprint arXiv:1804.03641*, 2018.
- [12] Jie Pu, Yannis Panagakis, Stavros Petridis, and Maja Pantic, "Audio-visual object localization and separation using low-rank and sparsity," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 2901–2905.
- [13] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon, "Learning to localize sound source in visual scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4358–4366.
- [14] Herman Kamper, Gregory Shakhnarovich, and Karen Livescu, "Semantic speech retrieval with a visually grounded model of untranscribed speech," *IEEE/ACM Transactions on Audio*, *Speech and Language Processing (TASLP)*, vol. 27, no. 1, pp. 89–98, 2019.

- [15] David Harwath, Adrià Recasens, Dídac Surís, Galen Chuang, Antonio Torralba, and James Glass, "Jointly discovering visual objects and spoken words from raw sensory input," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 649–665.
- [16] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T Freeman, and Michael Rubinstein, "Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation," arXiv preprint arXiv:1804.03619, 2018.
- [17] Grzegorz Chrupała, Lieke Gelderloos, and Afra Alishahi, "Representations of language in a model of visually grounded speech signal," arXiv preprint arXiv:1702.01991, 2017.
- [18] Kenneth Leidal, David Harwath, and James Glass, "Learning modality-invariant representations for speech and images," in 2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU). IEEE, 2017, pp. 424–429.
- [19] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu, "Audio-visual event localization in unconstrained videos," *ECCV*, 2018.
- [20] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework," in *Interna*tional Conference on Learning Representations, 2017.
- [21] Eric Jang, Shixiang Gu, and Ben Poole, "Categorical reparameterization with gumbel-softmax," *arXiv preprint arXiv:1611.01144*, 2016.
- [22] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba, "Learning deep features for discriminative localization," in *Computer Vision and Pattern Recognition*, 2016
- [23] Harold W Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [26] William E Vinje and Jack L Gallant, "Sparse coding and decorrelation in primary visual cortex during natural vision," *Science*, vol. 287, no. 5456, pp. 1273–1276, 2000.
- [27] Emmanuel Vincent, Rémi Gribonval, and Cédric Févotte, "Performance measurement in blind audio source separation," *IEEE transactions on audio, speech, and language processing*, vol. 14, no. 4, pp. 1462–1469, 2006.
- [28] Colin Raffel, Brian McFee, Eric J Humphrey, Justin Salamon, Oriol Nieto, Dawen Liang, Daniel PW Ellis, and C Colin Raffel, "mir\_eval: A transparent implementation of common mir

metrics," in In Proceedings of the 15th International Society for Music Information Retrieval Conference, ISMIR. Citeseer, 2014.