

# Sat2Sound: A Unified Framework for Zero-Shot Soundscape Mapping

Subash Khanal<sup>1</sup> Srikumar Sastry<sup>1</sup> Aayush Dhakal<sup>1</sup> Adeel Ahmad<sup>1,2</sup> Nathan Jacobs<sup>1</sup>

<sup>1</sup>Washington University in St. Louis <sup>2</sup>Taylor Geospatial Institute  
 {k.subash, s.sastry, a.dhakal, aadeel, jacobsn}@wustl.edu

## Abstract

We present Sat2Sound, a multimodal representation learning framework for soundscape mapping, designed to predict the distribution of sounds at any location on Earth. Existing methods for this task rely on satellite image and paired geotagged audio samples, which often fail to capture the diversity of sound sources at a given location. To address this limitation, we enhance existing datasets by leveraging a Vision-Language Model (VLM) to generate semantically rich soundscape descriptions for locations depicted in satellite images. Our approach incorporates contrastive learning across audio, audio captions, satellite images, and satellite image captions. We hypothesize that there is a fixed set of soundscape concepts shared across modalities. To this end, we learn a shared codebook of soundscape concepts and represent each sample as a weighted average of these concepts. Sat2Sound achieves state-of-the-art performance in cross-modal retrieval between satellite image and audio on two datasets: GeoSound and SoundingEarth. Additionally, building on Sat2Sound’s ability to retrieve detailed soundscape captions, we introduce a novel application: location-based soundscape synthesis, which enables immersive acoustic experiences. Our code and models will be publicly available.

## 1 Introduction

Imagine exploring our planet and listening to the sounds of a specific place, or creating a map that highlights locations that resemble your imagined soundscape. This is the essence of soundscape mapping—predicting and mapping the sound distribution of any location on Earth. Such a capability would provide immersive acoustic experiences that could be integrated into augmented reality or global mapping platforms, allowing us to hear the world as we explore it. For public health and urban planning, soundscape mapping can help monitor the acoustic ecosystem of an area and design sound-conscious interventions [4]. In addition, soundscape maps can be a valuable tool for real estate buyers and tourists, helping them find locations that match their acoustic preferences [29].

Recognizing the value of these capabilities, recent efforts have focused on developing frameworks for soundscape mapping [19, 20]. These frameworks represent locations using satellite images and learn a trimodal embedding space that links the satellite image, audio, and textual descriptions of the audio. Datasets used to train these frameworks contain geotagged audios collected from various crowdsourced platforms. However, these audio samples typically fail to capture the full diversity of sound sources at the recorded locations. In this work, we propose augmenting existing datasets by incorporating detailed textual descriptions of the soundscape, generated by querying a powerful vision-language model (VLM), LLaVA [26]. Soundscape descriptions produced by LLaVA capture a broader range of sound sources at each location, enhancing the semantic understanding of the soundscape associated with the satellite imagery.

Prior works for this task, such as GeoCLAP [19] and PSM [20], learn a contrastively trained trimodal embedding space, where each sample is represented by a global embedding. As a result, the

embedding space focuses primarily on global alignment between the modalities. However, in practice, even within a single satellite image, there may be multiple coexisting sound sources, and a particular sound category can arise from different parts of the scene. This motivates the need for local-level alignment across modalities. Recent works, such as FILIP [46] and FDT [5], have explored local alignment between modalities. FILIP adapts the contrastive objective by using token-wise maximum similarity between visual and textual tokens. FDT, on the other hand, uses a finite set of learnable tokens (a codebook) and represents samples for both modalities as a weighted aggregate of these tokens. For soundscape mapping, we hypothesize that the soundscape of any location is also a weighted aggregate of various soundscape concepts. Therefore, we adopt this codebook learning strategy into our framework to learn a shared codebook of soundscape concepts while encouraging local alignment between satellite images and their soundscape descriptions.

In our work, we focus on learning a discriminative multimodal embedding space for soundscape mapping. Nonetheless, this problem can also be approached from a generative perspective—synthesizing semantically meaningful soundscapes conditioned on satellite imagery. However, the inherently ambiguous relationship between satellite images and their corresponding audio makes training such generative models particularly challenging. In this regard, one possible training-free strategy is to first generate a detailed soundscape caption using a VLM, and then use that caption as input to a text-to-audio model [25, 12, 27, 17] to generate semantically meaningful sound for a given location. However, cascading two generative models in this way incurs high computational costs. Therefore, instead of relying on generation, we propose learning to retrieve VLM-generated image captions by treating them as an additional textual modality in Sat2Sound. The accurate retrieval of image captions and their pre-computed synthetic audio allows users to experience the sounds of any location with near-zero latency. The main contributions of our work are as follows.

- We propose a state-of-the-art soundscape mapping framework, Sat2Sound, which, in addition to learning from audio, textual descriptions of audio, and satellite image at the capture location, also learns from semantically rich synthetic soundscape descriptions.
- In Sat2Sound, we integrate a learnable codebook that represents a finite set of soundscape concepts shared across modalities, enhancing local alignment between image patches and soundscape concepts.
- By utilizing accurate image-to-text retrieval capability of Sat2Sound, we unify the strengths of image-to-text and text-to-audio generative models enabling novel application of location-conditioned soundscape synthesis.

## 2 Related Work

**Audio Visual Learning:** There is a strong semantic relationship between the acoustic and visual signals in a given audio-visual sample. Several studies [5, 16, 19, 20, 34, 48, 36, 37, 11] have leveraged this relationship to develop powerful audio-visual models. In the context of conditional audio generation, recent works [36, 42] have utilized existing foundational models to generate semantically meaningful audio from input images, while [50, 37] proposed models that generate images from audio. For soundscape mapping, recent works [19, 20] have focused on learning a shared embedding space between audio and satellite images. Building upon these efforts, our work aims to improve the alignment between satellite images and the diverse soundscape concepts.

**Contrastive Learning:** Contrastive learning is an effective strategy to learn shared embedding space between multiple modalities [31, 22, 47, 41, 20, 13]. The shared embedding space can be either deterministic or probabilistic. For example, [31] used contrastive learning to align large-scale image-text pairs, learning a deterministic image-text embedding space. Some recent works have proposed learning a probabilistic embedding space [8, 7] between modalities. While most of these methods focus on learning a single representation per sample to encourage global alignment between modalities, some recent works such as FILIP [46] and FDT [5] have adapted contrastive learning to encourage local alignment between image and text. MGA-CLAP [23] takes a similar approach to FDT and learns alignment between audio and text. Motivated by these approaches, we also train our framework by adapting a contrastive learning method [5], which fosters local alignment between satellite images and their soundscape descriptions.

**Discrete Representation Learning:** Discrete representation learning focuses on learning a discrete latent space (codebook) composed of a fixed number of concepts. These discrete latent concepts are

typically learned as intermediate weights within encoder-decoder frameworks, such as the Vector-Quantized Variational Autoencoder (VQ-VAE) [40]. VQ-VAE-style codebook learning has been applied to various conditional generation tasks including text-to-image generation [14], image-to-audio generation [18], and text-to-audio generation [45]. Beyond generative models, codebook learning has also been adopted in different cross-modal representation learning frameworks [24, 44, 23, 5]. Drawing inspiration from these works, we design our framework to learn a discrete set of soundscape concepts shared across our modalities: satellite images, text, and audio.

### 3 Method

This section describes our proposed framework: Sat2Sound, a multimodal representation learning framework for soundscape mapping.

Figure 1 provides an overview of Sat2Sound, which incorporates encoders for satellite images, audio, and text. Sat2Sound is trained on two types of text: first, the audio captions that capture the semantics of the geotagged audio in the dataset, and second, the textual descriptions of the soundscape at the location (image caption), generated by querying a powerful VLM, LLaVA [26]. More details on querying LLaVA can be found in Appendix A. Furthermore, we incorporate associated metadata for each sample (audio source, audio caption source, location, month, and time) and leverage the multiscale nature of satellite imagery, to enable metadata-aware, multi-scale soundscape mapping within our framework.

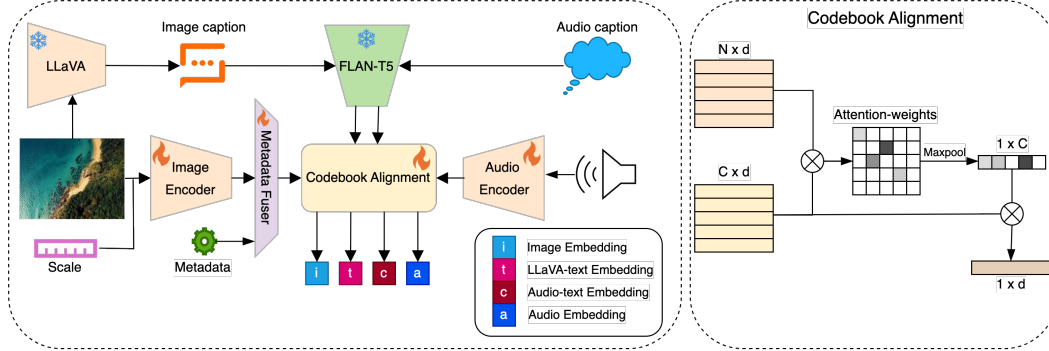


Figure 1: Sat2Sound framework learns a shared multimodal embedding space between satellite images, audio, audio captions, and image captions. Modality-specific encoders generate token embeddings for each modality, which are aligned into a shared codebook through an attention-score-based concept aggregation process.

#### 3.1 Encoding Modalities

Each sample ( $X$ ) used in our training consists of geotagged audio  $X^a$ , its corresponding audio caption  $X^c$ , a satellite image at a scale  $s$ , taken at the audio-recording location,  $X_s^i$ , and the associated image caption  $X^t$ . Modality-specific encoders,  $E_{audio}$ ,  $E_{text}$ , and  $E_{image}$ , are used to obtain patch/token-level representations, each projected into a  $d$ -dimensional embedding space:  $h^a \in \mathbb{R}^{N^a \times d}$ ,  $h^c \in \mathbb{R}^{N^c \times d}$ ,  $h^{i,s} \in \mathbb{R}^{N^i \times d}$ , and  $h^{t,s} \in \mathbb{R}^{N^{t,s} \times d}$  for audio, audio caption, image at scale  $s$ , and image caption, respectively. Here,  $N^a$  represents the number of frames in the audio feature,  $N^c$  is the number of tokens in the audio caption,  $N^i$  is the number of patches in the satellite image, and  $N^{t,s}$  is the number of tokens in the image caption.

$$h^m = E_m(X^m), \quad (1)$$

where  $m \in \{\text{audio, audio caption, image, image caption}\}$  and  $E_m$  is the modality-specific encoder.

To learn a metadata-aware embedding space, we adopt an early-fusion strategy, where we combine 5 metadata components (geolocation, month, hour, audio source, and audio caption source) with the patch embeddings for the satellite image, as obtained from Equation 1. Specifically, each

metadata component is first embedded into  $d$ -dimensional representations using shallow linear layers, and these representations are concatenated with the image patch embeddings, along the feature dimension. The concatenated input is then passed through a transformer-based module to obtain a metadata-conditioned satellite image representation.

$$h^{i'} = E_{meta}(h^i, \text{metadata}) \quad (2)$$

where  $E_{meta}$  represents the metadata fusion module, and  $h^{i'} \in \mathbb{R}^{(N^i+5) \times d}$  is the resulting metadata-conditioned satellite image patch embeddings.

### 3.2 Codebook Alignment

Once the modality-specific encoders compute the patch/token embeddings for each modality, the next step is to project them into a shared embedding space. To achieve this, we adopt a discrete representation learning strategy [5, 23], which learns a shared codebook,  $C \in \mathbb{R}^{M \times d}$ , representing  $M$  soundscape concepts shared across the modalities.

To illustrate the process, let us assume that our current modality of interest is satellite image ( $i$ ). For a given codebook token in  $C$ , the relevance scores for the input image patch embeddings are computed using an inner product operation, followed by selecting the maximum value across all patches ( $p_j^i, j \in 1, \dots, N^i$ ):

$$r_m^i = \max_j \langle p_j^i, C_m \rangle, \quad (3)$$

where  $r_m^i$  represents the maximum relevance score between visual patches and the  $m$ -th codebook token and  $C_m \in \mathbb{R}^{1 \times d}$  is the  $m$ -th codebook token. These relevance scores are then normalized using the Softmax function to obtain the attention weights:

$$w_m^i = \frac{e^{r_m^i}}{\sum_{m=1}^M e^{r_m^i}}. \quad (4)$$

Moreover, following [5], we also apply a Sparsemax function [28] to these attention weights to obtain sparser weights, which reduce noise and improve interpretability for grounding.

Using the same process, normalized attention weights for other modalities—audio caption ( $w_m^c$ ), image caption ( $w_m^t$ ), and audio ( $w_m^a$ )—are computed using the token/frame embeddings from their respective encoders and the shared codebook  $C$ . These attention weights enable each modality to dynamically attend to relevant codebook tokens, facilitating cross-modal alignment in the shared embedding space.

Finally, the pooled embeddings for each modality ( $r$ ) are obtained as a weighted sum of all the codebook concepts, as follows:

$$f^r = \sum_{m=1}^M w_m^r \cdot C_m, \quad r \in \{i, a, c, t\}, \quad (5)$$

where  $f^i, f^a, f^c$ , and  $f^t$  are the codebook-aligned embeddings for the image ( $i$ ), audio ( $a$ ), audio caption ( $c$ ), and image caption ( $t$ ) respectively.

### 3.3 Multimodal Contrastive Learning

Finally, the codebook-aligned embeddings obtained from Equation 5 are used in our multimodal contrastive learning framework. For a pair of modalities ( $u, v$ ), we use the InfoNCE loss [30, 31] which is defined as follows:

$$\mathcal{L}_{u,v} = -\frac{1}{2B} \left( \sum_{n=1}^B \log \frac{\exp((f_n^u \cdot f_n^v)/\tau_{uv})}{\sum_{s=1}^B \exp((f_s^u \cdot f_s^v)/\tau_{uv})} + \sum_{n=1}^B \log \frac{\exp((f_n^v \cdot f_n^u)/\tau_{uv})}{\sum_{s=1}^B \exp((f_s^v \cdot f_s^u)/\tau_{uv})} \right) \quad (6)$$

where  $\tau_{uv}$  is the learnable temperature parameter and  $B$  is the batch size during training.

Given the one-to-many nature of the matching between satellite images and audio, there can be pseudo-positives within a batch. Pseudo-positives are samples labeled as negatives but are semantically

similar enough to be considered positives, identified based on their similarity score in the latent space relative to the ground truth. Therefore, following [7, 20], we also include the contribution of these samples in the contrastive loss computation:

$$\mathcal{L}_{u,v}^\dagger = \mathcal{L}_{u,v} + \alpha \cdot \mathcal{L}_{u,v}^{\text{pseudo}} \quad (7)$$

where  $\mathcal{L}_{u,v}^{\text{pseudo}}$  is the contrastive loss computed while including the pseudo-positive matches in a batch as positive matches, and  $\alpha$  is the weight controlling its contribution. Using Equation 7, we compute the total loss for the following four modality pairs: image and audio ( $\mathcal{L}_{i,a}^\dagger$ ), image and audio caption ( $\mathcal{L}_{i,c}^\dagger$ ), audio and audio caption ( $\mathcal{L}_{a,c}^\dagger$ ), and image and image caption ( $\mathcal{L}_{i,t}^\dagger$ ). Finally, our trimodal contrastive learning objective is formulated as:

$$\mathcal{L}_{\text{tri}} = (\mathcal{L}_{i,a}^\dagger + \mathcal{L}_{i,c}^\dagger + \mathcal{L}_{a,c}^\dagger)/3 \quad (8)$$

Additionally, we acknowledge that the audio-to-image retrieval task can be modified to a composed audio-to-image retrieval task. In this setting, an audio query is also accompanied by its caption. To explicitly embrace this scenario during training, we create a composed audio embedding by combining the audio embedding with the audio-caption embedding:  $f^{a+c} = f^a + f^c$ . The contrastive loss is then computed between this composed audio embedding and the image embedding ( $\mathcal{L}_{i,a+c}^\dagger$ ).

Finally, we design our overall objective function as follows:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{tri}} + \mathcal{L}_{i,a+c}^\dagger + \mathcal{L}_{i,t}^\dagger \quad (9)$$

## 4 Evaluation

We evaluate Sat2Sound on two datasets: GeoSound and SoundingEarth. Experimental details, including information on the datasets, input processing, soundscape caption generation using LLaVA, encoders for each modality, and training hyperparameters, are provided in Appendix A.

**Baselines:** We compare our performance with prior works on zero-shot soundscape mapping: GeoCLAP [19] and PSM. [20]. For comparison of results, we present the performance of these methods as reported in the paper for the existing SOTA, PSM.

**Evaluation Metrics:** Following the baselines, we evaluate Sat2Sound for the task of cross-modal retrieval between satellite images and audios using Recall@10% (R10) and Median-Rank (MR) as our evaluation metrics. We use the image-to-audio retrieval (I2A-R10) performance on the validation set to select the best checkpoint for evaluation.

**Composed Image-Audio Retrieval:** The existing SOTA method, PSM, reports results for cross-modal retrieval between satellite images and audio when the audio caption embedding is added to both the audio and image query embeddings during retrieval. In contrast, we introduce a more realistic and fair setting, where the audio caption embedding is added only to the audio query. This setup better reflects practical scenarios, where off-the-shelf audio captioning models [10, 6] can be used to potentially enhance audio queries for retrieval.

**Image-Text Retrieval:** Apart from image-audio retrieval, we also evaluate Sat2Sound on image-text retrieval. For this evaluation, the LLaVA-generated soundscape caption serve as our text to be retrieved for an image. Unlike the often noisy or missing text annotations for audios in the dataset, using LLaVA ensures a semantically rich soundscape description for each location. To compare the performance of Sat2Sound for image-to-text retrieval, we create a strong baseline that is similar to Sat2Sound but trained only with image and image-caption pairs, without any metadata. We refer to this baseline as Sat2Text in our paper. Sat2Text was trained using  $\mathcal{L}_{i,t}^\dagger$  (Equation 7) as the training objective. Evaluation of image-text retrieval is conducted using Image-Text Recall@10% (R10) and Median Rank (MR). Additionally, to assess the similarity between the ground-truth image caption and the top-1 retrieved caption, we compute standard machine translation metrics: METEOR, BLEU, and F1 BERTScore (BERT-F1).

Table 1: Cross-modal retrieval performance comparison of Sat2Sound across different datasets.

Method	Dataset	Metadata	Image-to-Audio		Audio-to-Image	
			R@10	MR	R@10	MR
GeoCLAP	GeoSound-Bing	$\times$	0.399	1500	0.403	1464
PSM			0.423	1401	0.428	1344
Ours			<b>0.534</b>	<b>872</b>	<b>0.535</b>	<b>850</b>
PSM		$\checkmark$	0.828	261	0.829	248
Ours			<b>0.871</b>	<b>168</b>	<b>0.875</b>	<b>164</b>
GeoCLAP	GeoSound-Sentinel	$\times$	0.459	1179	0.465	1141
PSM			0.474	1101	0.485	1061
Ours			<b>0.549</b>	<b>802</b>	<b>0.556</b>	<b>778</b>
PSM		$\checkmark$	0.802	294	0.804	283
Ours			<b>0.868</b>	<b>191</b>	<b>0.872</b>	<b>183</b>
GeoCLAP	SoundingEarth	$\times$	0.454	667	0.449	694
PSM			0.514	547	0.518	543
Ours			<b>0.570</b>	<b>438</b>	<b>0.562</b>	<b>463</b>
PSM		$\checkmark$	0.563	454	0.569	447
Ours			<b>0.626</b>	<b>358</b>	<b>0.621</b>	<b>372</b>

Table 2: Image-Text retrieval results for different frameworks on GeoSound dataset with Bing.

Method	Metadata	scale	I2T-R10	I2T-MR	T2I-R10	T2I-MR	METEOR	BLEU	BERT-F1
Sat2Text		1	0.904	164	0.912	138	0.697	0.524	0.926
Sat2Sound		1	0.881	183	0.900	166	0.682	0.497	0.921
Sat2Sound	$\checkmark$	1	0.908	160	0.914	136	0.688	0.520	0.925
Sat2Text		3	0.929	117	0.927	110	0.695	0.519	0.922
Sat2Sound		3	0.918	134	0.921	120	0.675	0.491	0.917
Sat2Sound	$\checkmark$	3	0.940	106	0.938	97	0.689	0.513	0.920
Sat2Text		5	0.910	132	0.915	120	0.651	0.480	0.916
Sat2Sound		5	0.895	153	0.903	138	0.635	0.457	0.911
Sat2Sound	$\checkmark$	5	0.920	127	0.926	114	0.650	0.468	0.914

## 5 Results and Discussion

### 5.1 Cross Modal Retrieval: Image-Audio

Table 1 shows the results for cross-modal retrieval between satellite image and audio, using models trained in two different settings: one with associated metadata and another without. The GeoSound dataset contains satellite images from three satellite image scales (1, 3, and 5), but the results presented in this section correspond to scale 1. Additionally, we also include results for composed retrieval settings. Results for other scales and composed retrieval settings can be found in Tables 12, 13 and 14 of the appendix. As the results in these tables demonstrate, Sat2Sound achieves state-of-the-art performance in the task of satellite image-to-audio cross-modal retrieval.

Table 1 demonstrates that when Sat2Sound is trained on the GeoSound dataset with Bing imagery without metadata, I2A-R10 performance improves from 0.423 (the performance of the existing SOTA, PSM) to 0.534. In a similar setting with metadata, performance increases from 0.828 to 0.871. A similar improvement is observed for the GeoSound dataset with Sentinel imagery: I2A-R10 improves from 0.474 to 0.549 without metadata, and from 0.802 to 0.868 with metadata. Similarly, results for the SoundingEarth dataset show I2A-R10 improving from 0.514 to 0.570 without metadata, and from 0.563 to 0.626 with metadata. Similar to gains in image to audio retrieval, a noticeable improvement in audio-to-image retrieval is observed for Sat2Sound across both datasets and settings. Additionally, as reported in Tables 12, 13 and 14, consistent performance-gain is achieved by Sat2Sound across all scales of satellite images and different settings for cross-modal retrieval.

The results of Sat2Sound demonstrate its effectiveness in multi-scale satellite image-to-sound retrieval, showing significant improvements when trained and evaluated with metadata. To assess the contribution of each metadata component, we conduct ablation studies by evaluating Sat2Sound with individual or subsets of metadata components. The results, shown in Tables 6 and 7, reveal that the most impactful component is the audio source, consistent with the findings in PSM. This is further supported by the results from the SoundingEarth dataset, where the performance difference

between the model trained with and without metadata is small, as all samples come from the same source (Aporee:Radio [2]) for this dataset. In contrast, the GeoSound dataset includes samples from four distinct sources—Freesound [1], iNaturalist [3], Aporee:Radio [2], and Flickr [38]—each contributing unique sound types, such as nature sounds from iNaturalist and human activity sounds from Flickr. Therefore, for models trained on GeoSound dataset, the *audio source* metadata plays a critical role, enabling the model to learn an embedding space that accounts for the biases of different sound data platforms. During inference, user can select the expected audio source for a given location, enabling improved retrieval and metadata-conditioned soundscape maps.

## 5.2 Cross Modal Retrieval: Image-Text

In this section, we discuss the results of Sat2Sound on cross-modal retrieval between satellite image and image caption. Table 2 reports the performance of models trained on the GeoSound dataset with Bing imagery. As shown in the table, Sat2Sound trained without any metadata achieves an I2T-R10 of 0.898, when averaged across three scales of satellite imagery, with average performance slightly improving to 0.923 when trained with metadata. This performance is similar to the baseline (Sat2Text), which has an average I2T-R10 of 0.914. In addition to retrieval performance, the average METEOR score between the ground truth image caption and the top-1 retrieved caption is 0.681 for the baseline, 0.664 for Sat2Sound trained without metadata, and 0.676 for Sat2Sound trained with metadata indicating high similarity between the retrieved and ground truth caption. Some examples of the top-1 image captions retrieved by our model are provided in Figure 4.

From these results, we observe that Sat2Sound accurately retrieves semantically relevant image captions. Unlike the results for cross-modal retrieval between satellite images and audio, as seen in Table 1, the performance of Sat2Sound does not differ much with or without metadata. This is expected, as LLaVA was queried to generate a caption for the satellite image, providing a detailed soundscape description of the location, independent of the metadata associated with the sample. Moreover, the performance of Sat2Sound which is trained with all modalities, is similar to that of the baseline, which was specifically trained for image-to-text retrieval only. Finally, the high retrieval and METEOR, BLEU and BERT-F1 scores achieved by Sat2Sound encourage us to use the top-1 text retrieved by our model as input for any text-to-audio generation model, such as TangoFlux [17], to synthesize a semantically rich sound for any location on Earth.

## 5.3 Soundscape Maps

Utilizing the multimodal embedding space of Sat2Sound we can create large-scale soundscape maps for any region. As illustrated in Figure 2 (A), first, the satellite images (Bing or Sentinel) covering the geographic region of interest are downloaded. Then, for the desired scale and metadata settings, image embeddings are computed for each image. Finally, cosine similarity scores between the query embedding and all of the pre-computed image embeddings are used to create soundscape maps.

Using our best-performing model trained with Bing imagery, Figure 2 (C) was created with Bing images at scale 1 (i.e.,  $300 \times 300$  px with 0.6m GSD at zoom level 18), covering the USA, based on two textual prompts. The first prompt describes the soundscape of a forest: *sounds of birds chirping, leaves rustling, and the occasional rustling of branches*, while the second prompt describes: *sounds of farm machines*. As shown in Figure 2 (C), the regions highlighted by these queries correspond to the relevant areas marked as forest and cropland, respectively, in a reference land-cover map. Following the same procedure, regional-scale soundscape maps were created, as shown in Figure 2 (B), for different textual queries and audio samples obtained from Freesound. In region *a*, with the textual query: *sounds of trains passing by on the tracks*, the area with a railway track is highly activated. In region *b* of Figure 2 (B), for the audio query *car horn*, higher activation is observed outside the dense park. For the audio query *rooster crowing*, activation appear outside the dense urban region. Similarly, in region *c*, urban areas are activated for the audio query *construction drilling*, while non-urban areas are activated for the textual query: *sounds of animals on a farm*. These qualitative results suggest the ability of Sat2Sound to create semantically meaningful soundscape maps.

## 5.4 Codebook-guided Local Alignment

As discussed earlier, Sat2Sound learns a codebook of soundscape concepts shared across modalities. Once trained, this codebook enables local alignment between image patches and their corresponding

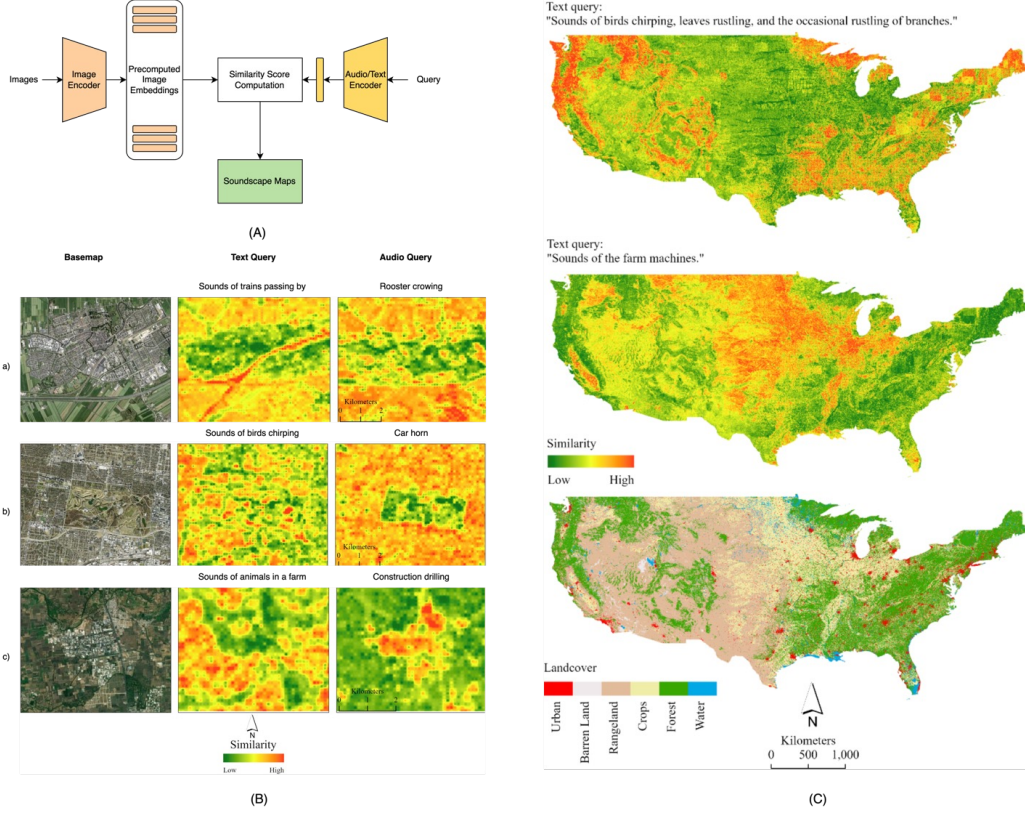


Figure 2: (A) Soundscape mapping framework using Sat2Sound’s encoders. (B) City-scale soundscape maps using different queries for cities in a) Netherlands; b) USA; c) India. (C) Country-scale soundscape maps created for queries over the USA with a reference land cover map for comparison.

soundscape concepts. Figure 3 demonstrates Sat2Sound’s ability to align words/phrases in text with patches in the image. Specifically, given token-level embeddings for the soundscape textual query ( $h^t \in \mathbb{R}^{N^t \times d}$ ), we compute the inner product between  $h^t$  and the learned codebook ( $C \in \mathbb{R}^{M \times d}$ ), resulting in attention scores between words and codebook concepts. Now for the target word, we first select the concept with the highest attention score and use the index of the selected concept to get attention scores for all image patches corresponding to the grounded concept. For phrases, attention scores grounded to each words are averaged across the image patches.

## 6 Location-based Soundscape Synthesis

In this section, we propose two training-free alternatives for location-based soundscape synthesis, where location is represented by satellite imagery.

**Cascaded Generative Framework:** This framework is built by cascading two types of generative models: image-to-text and text-to-audio. Specifically, we first query a powerful instruction-tuned VLM, LLaVA, to generate a detailed soundscape caption for a given image, and then use that caption to generate semantically diverse sound using the recent SOTA text-to-audio generative model, TangoFlux [17]. Assuming the availability of generative models, this approach is training-free and offers the flexibility to choose from the best generative models available at the time. However, considering the high computational cost and latency of this pipeline, it may not be suitable for real-time systems such as web services or augmented reality applications.

**Retrieval-based Generative Framework:** As reported in Table 2, Sat2Sound can accurately retrieve a soundscape caption for a given satellite image. Examples of the top-1 retrieved image captions by Sat2Sound can be found in Figure 4, where we observe a strong alignment between the semantics of the captions and the corresponding satellite images. Leveraging this capability, we can use the



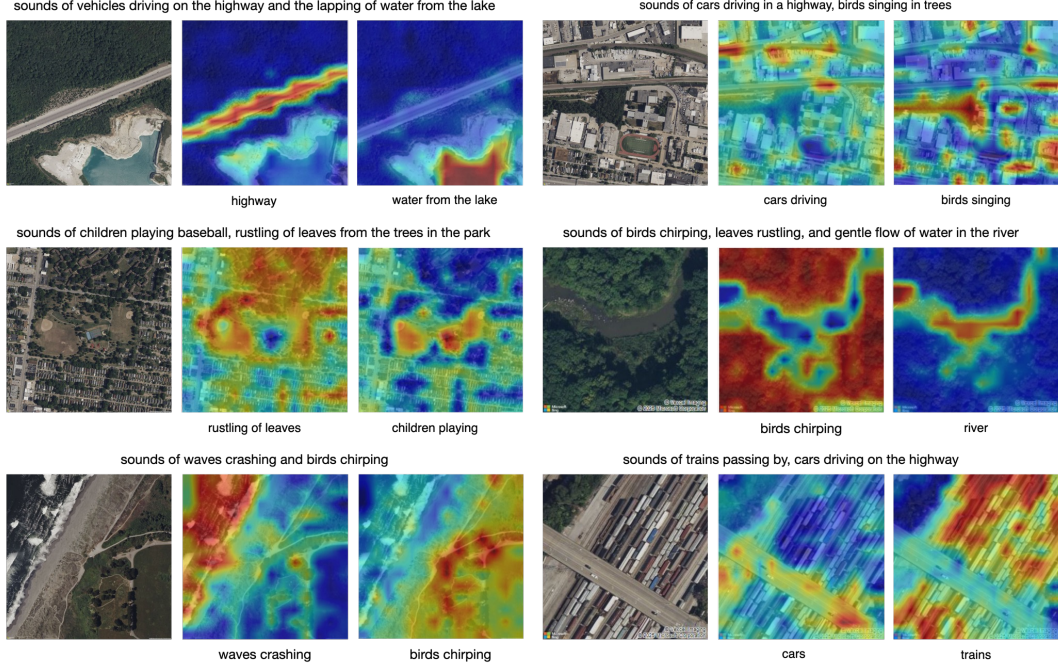


Figure 3: Alignment between patches in a single image and soundscape concepts in textual query.

retrieved caption to generate audio through a text-to-audio generative model. In practice, we can curate a gallery of diverse locations worldwide, pre-compute LLaVA-generated soundscape captions, and their corresponding TangoFlux-generated synthetic audio. During inference, this approach eliminates the computational cost of generation, making satellite image-to-sound generation purely retrieval-based, while still benefiting from generative models’ ability to produce semantically rich soundscapes for any location. We present a web demo of our retrieval-based generative framework in our supplemental video.

**Evaluation of Frameworks:** To assess and compare the semantic quality of synthetic audio generated by two frameworks (Cascaded Generative and Retrieval-based), we conducted a perceptual study. In our study, 16 participants rated 20 locations on a scale from 1 to 5, evaluating how likely the generated sound is to be heard at each location shown in the satellite image. More details of the study can be found in Appendix A. The average human ratings for synthetic audios are provided as *score* in Table 3. As shown in Table 3, our retrieval-based generative approach achieves ratings similar to a fully generative framework, while being significantly more efficient. With Sat2Sound’s 130M-parameter satellite image encoder, the retrieval-based approach operates with minimal computational cost (TFLOPS of 0.14, 0.5s latency on CPU). In contrast, the cascaded generative approach, with 7.57B parameters, 49.03 TFLOPS has latency of 102.5s on CPU, and about 4s on an NVIDIA H100 GPU. Finally, it is worth noting that, for the retrieval-based framework, as text-to-audio models improve, we can easily swap in a new gallery generated by the best available text-to-audio model, ensuring flexibility for future advancements.

Table 3: Comparison of location-based soundscape synthesis methods.

Approach	score	#params	TFLOPS
Generative	$3.77 \pm 0.51$	7.57B	49.03
Retrieval	$3.52 \pm 0.48$	130M	0.14

## 7 Conclusion

We introduced Sat2Sound, a multimodal representation learning framework that integrates audio, audio descriptions, satellite images, and soundscape descriptions from images. Sat2Sound advances the state-of-the-art in soundscape mapping. Our approach explicitly learns from multiple sound sources present at any location and aligns them to image regions through a learned codebook. Furthermore, we showed Sat2Sound’s applicability for location-based soundscape synthesis. We hope that our quantitative and qualitative results support our recommendation to use Sat2Sound as a unified framework for soundscape mapping.

## References

- [1] Freesound, <https://freesound.org/>.
- [2] Radio aporee: Maps - sounds of the world, <https://aporee.org>.
- [3] inaturalist, <https://www.inaturalist.org>.
- [4] Francesco Aletta, Tin Oberman, Andrew Mitchell, Mercede Erfanian, Matteo Lionello, Magdalena Kachlicka, and Jian Kang. Associations between soundscape experience and self-reported wellbeing in open public urban spaces: a field study. *The Lancet*, 394:S17, 2019.
- [5] Yuxiao Chen, Jianbo Yuan, Yu Tian, Shijie Geng, Xinyu Li, Ding Zhou, Dimitris N Metaxas, and Hongxia Yang. Revisiting multimodal representation in contrastive learning: from patch and token embeddings to finite discrete tokens. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15095–15104, 2023.
- [6] Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.
- [7] Sanghyuk Chun. Improved probabilistic image-text representations. In *The Twelfth International Conference on Learning Representations*, 2024.
- [8] Sanghyuk Chun, Seong Joon Oh, Rafael Sampaio De Rezende, Yannis Kalantidis, and Diane Larlus. Probabilistic embeddings for cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8415–8424, 2021.
- [9] Yezhen Cong, Samar Khanna, Chenlin Meng, Patrick Liu, Erik Rozi, Yutong He, Marshall Burke, David B. Lobell, and Stefano Ermon. SatMAE: Pre-training transformers for temporal and multi-spectral satellite imagery. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=WBhqzpF6KYH>.
- [10] Soham Deshmukh, Benjamin Elizalde, Rita Singh, and Huaming Wang. Pengi: An audio language model for audio tasks. *Advances in Neural Information Processing Systems*, 36: 18090–18108, 2023.
- [11] Junyu Gao, Hao Yang, Maoguo Gong, and Xuelong Li. Audio–visual representation learning for anomaly events detection in crowds. *Neurocomputing*, 582:127489, 2024.
- [12] Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction-tuned llm and latent diffusion model. *arXiv preprint arXiv:2304.13731*, 2023.
- [13] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190, 2023.
- [14] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10696–10706, 2022.
- [15] Claudia Hauff. A study on the accuracy of flickr’s geotag data. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*, pages 1037–1040, 2013.
- [16] Di Hu, Xuhong Li, Lichao Mou, Pu Jin, Dong Chen, Liping Jing, Xiaoxiang Zhu, and Dejing Dou. Cross-task transfer for geotagged audiovisual aerial scene recognition. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIV 16*, pages 68–84. Springer, 2020.

- [17] Chia-Yu Hung, Navonil Majumder, Zhifeng Kong, Ambuj Mehrish, Rafael Valle, Bryan Catanzaro, and Soujanya Poria. Tangoflux: Super fast and faithful text to audio generation with flow matching and clap-ranked preference optimization. *arXiv preprint arXiv:2412.21037*, 2024.
- [18] Vladimir Iashin and Esa Rahtu. Taming visually guided sound generation. In *The 32st British Machine Vision Virtual Conference*. BMVA Press, 2021.
- [19] Subash Khanal, Srikumar Sastry, Aayush Dhakal, and Nathan Jacobs. Learning tri-modal embeddings for zero-shot soundscape mapping. In *British Machine Vision Conference (BMVC)*, November 2023.
- [20] Subash Khanal, Xing Eric, Srikumar Sastry, Aayush Dhakal, Xiong Zhexiao, Adeel Ahmad, and Nathan Jacobs. Psm: Learning probabilistic embeddings for multi-scale zero-shot soundscape mapping. In *ACM Multimedia*, November 2024.
- [21] Konstantin Klemmer, Esther Rolf, Caleb Robinson, Lester Mackey, and Marc Rußwurm. Satclip: Global, general-purpose location embeddings with satellite imagery. *arXiv preprint arXiv:2311.17179*, 2023.
- [22] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021.
- [23] Yiming Li, Zhifang Guo, Xiangdong Wang, and Hong Liu. Advancing multi-grained alignment for contrastive language-audio pre-training. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 7356–7365, 2024.
- [24] Alexander H Liu, SouYoung Jin, Cheng-I Jeff Lai, Andrew Rouditchenko, Aude Oliva, and James Glass. Cross-modal discrete representation learning. *arXiv preprint arXiv:2106.05438*, 2021.
- [25] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*, 2023.
- [26] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [27] Navonil Majumder, Chia-Yu Hung, Deepanway Ghosal, Wei-Ning Hsu, Rada Mihalcea, and Soujanya Poria. Tango 2: Aligning diffusion-based text-to-audio generations through direct preference optimization. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 564–572, 2024.
- [28] Andre Martins and Ramon Astudillo. From softmax to sparsemax: A sparse model of attention and multi-label classification. In *International conference on machine learning*, pages 1614–1623. PMLR, 2016.
- [29] Pierluigi Morano, Francesco Tajani, Felicia Di Liddo, and Michele Darò. Economic evaluation of the indoor environmental quality of buildings: The noise pollution effects on housing prices in the city of bari (italy). *Buildings*, 11(5):213, 2021.
- [30] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [32] Colorado J Reed, Ritwik Gupta, Shufan Li, Sarah Brockman, Christopher Funk, Brian Clipp, Kurt Keutzer, Salvatore Candido, Matt Uyttendaele, and Trevor Darrell. Scale-mae: A scale-aware masked autoencoder for multiscale geospatial representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4088–4099, 2023.

- [33] Adam Roberts, Hyung Won Chung, Gaurav Mishra, Anselm Levskaya, James Bradbury, Daniel Andor, Sharan Narang, Brian Lester, Colin Gaffney, Afroz Mohiuddin, et al. Scaling up models and data with t5x and seqio. *Journal of Machine Learning Research*, 24(377):1–8, 2023.
- [34] Tawfiq Salem, Menghua Zhai, Scott Workman, and Nathan Jacobs. A multimodal approach to mapping soundscapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 2524–2527, 2018.
- [35] Srikumar Sastry, Subash Khanal, Aayush Dhakal, Adeel Ahmad, and Nathan Jacobs. Taxabind: A unified embedding space for ecological applications. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1765–1774. IEEE, 2025.
- [36] Roy Sheffer and Yossi Adi. I hear your true colors: Image guided audio generation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- [37] Kim Sung-Bin, Arda Senocak, Hyunwoo Ha, Andrew Owens, and Tae-Hyun Oh. Sound to visual scene generation by audio-to-visual latent alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6430–6440, June 2023.
- [38] Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. YFCC100M: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016. URL <http://cacm.acm.org/magazines/2016/2/197425-yfcc100m/fulltext>.
- [39] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*, 2025.
- [40] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- [41] Vicente Vivanco Cepeda, Gaurav Kumar Nayak, and Mubarak Shah. Geoclip: Clip-inspired alignment between locations and images for effective worldwide geo-localization. *Advances in Neural Information Processing Systems*, 36, 2024.
- [42] Heng Wang, Jianbo Ma, Santiago Pascual, Richard Cartwright, and Weidong Cai. V2a-mapper: A lightweight solution for vision-to-audio generation by connecting foundation models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 15492–15501, 2024.
- [43] Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.
- [44] Yan Xia, Hai Huang, Jieming Zhu, and Zhou Zhao. Achieving cross modal generalization with multimodal unified representation. *Advances in Neural Information Processing Systems*, 36, 2024.
- [45] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. Diffsound: Discrete diffusion model for text-to-sound generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:1720–1733, 2023.
- [46] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. *arXiv preprint arXiv:2111.07783*, 2021.
- [47] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research*, 2022.

- [48] Donghuo Zeng, Jianming Wu, Gen Hattori, Rong Xu, and Yi Yu. Learning explicit and implicit dual common subspaces for audio-visual cross-modal retrieval. *ACM Transactions on Multimedia Computing, Communications and Applications*, 19(2s):1–23, 2023.
- [49] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023.
- [50] Lin Zhang, Shentong Mo, Yijing Zhang, and Pedro Morgado. Audio-synchronized visual animation. In *European Conference on Computer Vision*, pages 1–18. Springer, 2025.

## A Experimental Details

**Datasets:** We experiment with two datasets: *GeoSound* and *SoundingEarth*. *GeoSound* contains 294019/5000/9931 train/validation/test samples and uses both 0.6m GSD *Bing* image tiles ( $1500 \times 1500$ ) and 10m GSD *Sentinel-2* image tiles ( $1280 \times 1280$ ). *SoundingEarth* with 0.2m GSD *Google Earth* satellite image tiles of size ( $1024 \times 1024$ ) contains 41469/3242/5801 train/validation/test samples.

**Input Processing:** We process our three input modalities: audio, text, and image as follows:

*Audio:* We convert all of our input audios to mono-audio, randomly sample a 10-second segment, and resample to 32000 Hz. The audio then undergoes Short-Time Fourier Transform (STFT) with window size of 1024 and hop length of 320, followed by conversion to mel-spectrogram with 64 bands for frequencies from 50 Hz to 14000 Hz.

*Text:* Both audio captions and image captions are tokenized using the `google/flan-t5-large` tokenizer.

*Image:* For the *GeoSound* dataset, we center-crop satellite images based on a scale (uniformly sampled from  $\{1, 3, 5\}$ ) multiplied by source-specific tile sizes (256px for Sentinel, 300px for Bing). For *SoundingEarth*, we only center-crop 256px. In both cases, we resize the cropped images to 224x224px and apply color jitter and normalization.

**Metadata:** Following PSM [20], Sat2Sound is also trained with metadata (geolocation, month, hour, audio source, and audio caption source) in addition to satellite imagery and associated audio and text. For the *GeoSound* dataset used in our work, the geotagged audios were collected from four different sources: Freesound, Aporee, iNaturalist, and Flickr. The audio caption can either come from the user-uploaded textual description or be generated using recent SOTA audio-to-text generation models such as Pengi [10] or Qwen-Audio [6], with the caption selection based on the caption’s CLAP score [43] with the ground-truth audio. Each metadata component is embedded into a 1024-dimensional vector and fused using Sat2Sound’s transformer-based metadata fusion module. To prevent overfitting, we apply a dropout rate of 0.5, independently dropping each metadata component during training.

**Image Captions:** For the cropped satellite images at each scale, we generate detailed soundscape captions using LLaVA [26], a powerful open-source Vision Language Model. Specifically, we query `llava-hf/llava-1.5-7b-hf` on HuggingFace using the following prompt:

*“What types of sounds can we expect to hear from the location captured by this aerial view image? Describe in up to two sentences.”*

**Encoders:** Following [20], we fine-tune the pre-trained checkpoint for the SATMAE-encoder [9] to encode satellite imagery while updating its positional embeddings with scale-aware GSDPE [32] to encode the scale of the satellite image. For audio, we fine-tune the pre-trained audio encoder of MGA-CLAP [23], which generates frame-level audio embeddings. The textual modality is processed using a frozen FLAN-T5 [33] model, which extracts token embeddings from texts for each sample.

**Hyper-parameters:** We set the embedding dimension of Sat2Sound ( $d$ ) to 1024 and the number of concepts in the codebook ( $M$ ) to 16000. We train our model using the AdamW optimizer with cosine-annealing with warm restarts as the learning rate scheduler with following parameters: learning rate of  $5e-5$ , weight decay of 0.2, and betas of (0.9, 0.98). We set the pseudo-positives contribution to loss ( $\alpha$ ) to 0.1. We train Sat2Sound for 20 epochs with train-batch size ( $B$ ) of 128.

**Compute Infrastructure:** All experiments were conducted on an NVIDIA H100 80GB GPU, using 16 workers to enable faster data loading. We employed full-precision training throughout.

**Human Study:** In this study, 16 participants were shown a Bing satellite image at scale 1 for 20 locations on Earth. These 20 locations were selected by clustering SatCLIP’s [21] geolocation embeddings of all the samples in our gallery, with the centroid of each cluster serving as a test location. Each satellite image was paired with two 10-second synthetic audios generated using TangoFlux[17] with parameters: `steps` = 50 and `guidance` = 4.5. One audio was generated using the top-1 retrieved image caption by Sat2Sound, and the second using the directly generated LLaVA caption passed to TangoFlux.

Table 4: Ablation of different loss components. Evaluated for cross modal image-to-audio retrieval on *GeoSound* with *Bing* imagery at scale 1.

trimodal	L(a+c)	L(i,t)	I2A-R@10	I2A-MdR	A2I-R@10	A2I-MdR
✓			0.866	192	0.871	179
✓		✓	0.852	206	0.852	203
✓	✓		0.873	182	0.876	169
✓	✓	✓	0.871	168	0.875	164

Table 5: Ablation of different loss components. Evaluated for cross model composite audio-to-image retrieval on *GeoSound* with *Bing* imagery at scale 1.

trimodal	L(a+c)	L(i,t)	I2A-R@10	I2A-MdR	A2I-R@10	A2I-MdR
✓			0.935	88	0.949	75
✓		✓	0.922	105	0.937	94
✓	✓		0.960	71	0.962	62
✓	✓	✓	0.955	70	0.958	64

## B Ablation Studies

### B.1 Loss Ablation

We conduct an ablation study on different components of the loss to assess their impact on the overall training objective (Equation 9). We observe that the addition of the composite audio-based loss ( $\mathcal{L}_{i,a+c}^\dagger$ ) slightly improves the performance of the standard audio-image cross-modal retrieval as observed in Table 4 and noticeably improves for composed audio-image cross modal retrieval as observed in Table 5. Furthermore, the inclusion of an additional image-text loss ( $\mathcal{L}_{i,t}^\dagger$ ) does not degrade performance and provides the benefit of accurate retrieval of text as reflected in Table 2, Figure 4 and our supplemental video.

### B.2 Metadata Ablation

This experiment is intended to identify the impact of different metadata components in the cross-modal retrieval performance of Sat2Sound. As observed in Table 6, the most contributing metadata is *audio source*. This is consistent with results in prior work, PSM [20].

In addition to the independent metadata ablation presented in Table 6, we also conduct a more detailed ablation study on metadata combinations, as shown in Table 7. These results demonstrate how our model’s performance progressively improves with the addition of different metadata components and provide further insights into the contribution of metadata components. Notably, with Bing imagery, performance improves from 0.787 to 0.866 when the easily available components—time and geolocation—are added. In contrast, the performance gain from incorporating the *audio caption source* is not significant, suggesting that any audio captioning model or user-written text can be used, as long as a reasonable audio caption is provided to query our model.

### B.3 Codebook Size Ablation:

We conduct ablation of codebook size of our framework. As seen in Tables 8 and 9, performance remains fairly consistent across different codebook sizes. We speculate that the sparsification operation[28] of the attention weights (Equation 4), encourages our framework to select only the relevant concepts, making the framework independent of the codebook size.

Our choice of codebook-based learning is motivated by the intuition that a fixed set of soundscape concepts can be shared across modalities. This approach offers a more interpretable way to align local features between the satellite image and corresponding soundscape elements. As shown in Figure 3, this local alignment serves as a valuable byproduct from an interpretability perspective, further discussed in Section E.

Table 6: Metadata ablation to evaluate Sat2Sound models trained on GeoSound dataset with satellite imagery at scale 1.

Imagery	latlong	month	time	a-source	c-source	I2A-R10	I2A-MR	A2I-R10	A2I-MR
Sentinel	✓					0.603	613	0.627	566
Sentinel		✓				0.585	682	0.604	613
Sentinel			✓			0.641	537	0.669	477
Sentinel				✓		<b>0.805</b>	<b>318</b>	<b>0.808</b>	<b>306</b>
Sentinel					✓	0.562	763	0.590	672
Bing	✓					0.627	547	0.644	499
Bing		✓				0.577	697	0.594	638
Bing			✓			0.629	564	0.651	521
Bing				✓		<b>0.787</b>	<b>325</b>	<b>0.793</b>	<b>320</b>
Bing					✓	0.551	785	0.566	742

Table 7: Composite metadata ablation to evaluate Sat2Sound models trained on GeoSound dataset with satellite imagery at scale 1.

Imagery	a-source	time	latlong	month	c-source	I2A-R10	I2A-MR	A2I-R10	A2I-MR
Sentinel	✓					0.805	318	0.808	306
Sentinel	✓	✓				0.819	295	0.820	278
Sentinel	✓	✓	✓	✓		0.850	239	0.851	227
Sentinel	✓	✓	✓	✓	✓	0.862	200	0.865	192
Sentinel	✓	✓	✓	✓	✓	<b>0.868</b>	<b>191</b>	<b>0.872</b>	<b>183</b>
Bing	✓					0.787	325	0.793	320
Bing	✓	✓				0.807	289	0.811	283
Bing	✓	✓	✓	✓		0.854	209	0.857	202
Bing	✓	✓	✓	✓	✓	0.866	175	0.871	169
Bing	✓	✓	✓	✓	✓	<b>0.871</b>	<b>168</b>	<b>0.875</b>	<b>164</b>

## C Simpler Baselines

In this section, we compare the performance of Sat2Sound with existing off-the-shelf multimodal embedding spaces. As shown in the image-text cross-modal retrieval results (Table 10), existing pre-trained image-text models underperform compared to Sat2Sound. We attribute this to the mismatch between the soundscape descriptions generated by LLaVA from satellite images and the textual data these models were originally trained on. In contrast, Sat2Sound is explicitly trained on these captions, giving it a clear advantage and resulting in significantly better performance than the compared vision-language baselines. A similar trend is observed in Table 11 for image-audio cross-modal retrieval. These findings underscore the limitations of existing state-of-the-art multimodal embedding spaces for soundscape mapping, highlighting the need for a specialized framework like Sat2Sound tailored to this task.

## D Multi-scale Cross-Modal Retrieval

Sat2Sound is trained on multi-scale satellite imagery for the GeoSound dataset. The results presented in the main paper are for satellite imagery at scale 1. In this section, we present results for two additional scales: 3 and 5, using both Sentinel and Bing imagery from the GeoSound dataset. Additionally, for both datasets (GeoSound and SoundingEarth), we provide results for composed retrieval settings where the audio caption embedding is added either only to the audio query embedding (indicated as *audio* in the tables) or to both the audio query and image query embeddings (indicated as *query* in the tables), as done in PSM [20]. As observed in Tables 12, 13, and 14, Sat2Sound outperforms existing baselines by a noticeable margin in almost all of the settings.

## E Analyzing codebook concepts

The codebook learned by Sat2Sound can be used to generate fine-grained soundscape maps for regions covered by a single satellite image, as illustrated in Figure 3. In this section, we qualitatively explore what the codebook has learned. Specifically, for our gallery of image captions, we first obtain



Table 8: Codebook ablation for Image-Text retrieval on GeoSound dataset with Bing imagery (scale=1) and corresponding image captions.

Codebook Size	I2T-R10	I2T-MR	T2I-R10	T2I-MR
4000	0.905	167	0.915	145
8000	0.899	163	0.917	146
16000	0.908	160	0.914	136
32000	0.902	165	0.914	148

Table 9: Codebook ablation for Image-Audio retrieval on GeoSound dataset with Bing imagery (scale=1) and corresponding audio.

Codebook Size	I2A-R10	I2A-MR	A2I-R10	A2I-MR
4000	0.868	167	0.870	161
8000	0.875	171	0.876	163
16000	0.871	168	0.875	164
32000	0.874	164	0.879	160

the corresponding codebook attention weights (Equation 4) and group together samples that share a similar set of highly activated codebook concepts. For a subset of these groups, we randomly sample examples to examine the behavior and semantic meaning captured by different sets of codebook concepts. Representative samples are visualized in Figure 5.

## F Broader Impact and Limitations

Sat2Sound has the potential to enhance commercial applications such as augmented reality and geospatial navigation systems by providing immersive auditory experience. It also offers significant value to public health and urban planning stakeholders, enabling more informed decisions related to acoustic ecosystem monitoring, and the design of sound-conscious urban environments.

Although our framework is trained and evaluated on GeoSound, the largest available geospatially diverse dataset of geotagged audios from four sources, it still suffers from an uneven global distribution of samples, limiting its ability to fully capture soundscapes in underrepresented regions. We hope our work encourages more citizen science initiatives to collect more geotagged audios from diverse locations, enabling further improvements of soundscape mapping frameworks. Moreover, the precision of geolocation information collected from geotagged audio metadata can vary [15], leading to potential mismatches between the context of the exact audio recording site and that visible in overhead images. Additionally, our framework learns soundscape representations primarily from an overhead perspective, and incorporating ground-level view could provide valuable complementary information for more precise representations. We leave this as an avenue for future work.

Table 10: Image-Text retrieval comparison with additional baselines. Results on GeoSound with Bing Imagery (scale=1).

Model	I2T-R10	I2T-MR	T2I-R10	T2I-MR
CLIP[31]	0.528	1420	0.461	1999
SigLIP[49]	0.340	3307	0.368	2652
SigLIP2[39]	0.449	2641	0.397	2400
Ours(w/o meta)	0.881	183	0.900	166
Ours(w meta)	<b>0.908</b>	<b>160</b>	<b>0.914</b>	<b>136</b>

Table 11: Image-Audio retrieval comparison with additional baselines. Results on GeoSound with Sentinel Imagery (scale=1).

Model	I2A-R10	I2A-MR	A2I-R10	A2I-MR
ImageBind[13]	0.214	3675	0.231	3541
TaxaBind[35]	0.235	3448	0.250	3400
Ours(w/o meta)	0.549	802	0.556	778
Ours(w meta)	<b>0.868</b>	<b>191</b>	<b>0.872</b>	<b>183</b>

Table 12: Image-Audio retrieval results for SoundingEarth with different composed audio-image settings.

Method	Composed	I2A-R10	I2A-MR	A2I-R10	A2I-MR
<i>Without Metadata</i>					
GeoCLAP	query	0.523	533	0.470	641
PSM	query	0.687	234	0.560	451
Ours	query	<b>0.847</b>	<b>94</b>	<b>0.564</b>	<b>448</b>
GeoCLAP	audio	0.478	624	0.470	641
PSM	audio	0.558	462	0.560	451
Ours	audio	<b>0.567</b>	<b>443</b>	<b>0.564</b>	<b>448</b>
<i>With Metadata</i>					
PSM	query	0.690	264	0.608	371
Ours	query	<b>0.855</b>	<b>91</b>	<b>0.862</b>	<b>129</b>
PSM	audio	0.606	380	0.608	371
Ours	audio	<b>0.855</b>	<b>127</b>	<b>0.862</b>	<b>129</b>

Table 13: Image-Audio retrieval results for GeoSound with Bing imagery at different scales.

Scale	Method	Composed	I2A-R10	I2A-MR	A2I-R10	A2I-MR
<i>Without Metadata</i>						
1	GeoCLAP	query	0.577	712	0.468	1141
	PSM	query	0.754	204	0.510	952
	Ours	query	<b>0.903</b>	<b>82</b>	<b>0.540</b>	<b>836</b>
	GeoCLAP	audio	0.464	1159	0.468	1141
	PSM	audio	0.503	980	0.510	952
	Ours	audio	<b>0.535</b>	<b>864</b>	<b>0.540</b>	<b>836</b>
3	GeoCLAP	none	0.408	1441	0.420	1389
	PSM	none	0.440	1302	0.443	1266
	Ours	none	<b>0.560</b>	<b>777</b>	<b>0.561</b>	<b>779</b>
	GeoCLAP	query	0.577	707	0.483	1056
	PSM	query	0.753	207	0.529	880
	Ours	query	<b>0.908</b>	<b>79</b>	<b>0.567</b>	<b>737</b>
5	GeoCLAP	audio	0.477	1092	0.483	1056
	PSM	audio	0.523	891	0.529	880
	Ours	audio	<b>0.564</b>	<b>751</b>	<b>0.567</b>	<b>737</b>
	GeoCLAP	none	0.409	1428	0.421	1373
	PSM	none	0.440	1302	0.448	1279
	Ours	none	<b>0.564</b>	<b>760</b>	<b>0.559</b>	<b>770</b>
5	GeoCLAP	query	0.581	698	0.489	1036
	PSM	query	0.753	209	0.532	863
	Ours	query	<b>0.910</b>	<b>78</b>	<b>0.567</b>	<b>748</b>
	GeoCLAP	audio	0.482	1071	0.489	1036
	PSM	audio	0.528	881	0.532	863
	Ours	audio	<b>0.554</b>	<b>764</b>	<b>0.567</b>	<b>748</b>
<i>With Metadata</i>						
1	PSM	query	0.901	113	0.943	100
	Ours	query	<b>0.970</b>	<b>33</b>	<b>0.958</b>	<b>64</b>
	PSM	audio	0.935	115	0.943	100
	Ours	audio	<b>0.955</b>	<b>70</b>	<b>0.958</b>	<b>64</b>
3	PSM	none	0.827	266	0.832	250
	Ours	none	<b>0.874</b>	<b>163</b>	<b>0.879</b>	<b>159</b>
	PSM	query	0.900	114	0.945	102
	Ours	query	<b>0.972</b>	<b>32</b>	<b>0.960</b>	<b>62</b>
5	PSM	audio	0.936	118	0.945	102
	Ours	audio	<b>0.957</b>	<b>66</b>	<b>0.960</b>	<b>62</b>
	PSM	none	0.821	281	0.826	261
	Ours	none	<b>0.877</b>	<b>167</b>	<b>0.882</b>	<b>167</b>
5	PSM	query	0.896	115	0.941	107
	Ours	query	<b>0.972</b>	<b>32</b>	<b>0.963</b>	<b>64</b>
	PSM	audio	0.929	124	0.941	107
	Ours	audio	<b>0.959</b>	<b>68</b>	<b>0.963</b>	<b>64</b>

Table 14: Image-Audio retrieval results for GeoSound with Sentinel imagery at different scales.

Scale	Method	Composed	I2A-R10	I2A-MR	A2I-R10	A2I-MR
<i>Without Metadata</i>						
1	GeoCLAP	query	0.546	827	0.553	804
	PSM	query	0.803	153	<b>0.595</b>	<b>664</b>
	Ours	query	<b>0.909</b>	<b>79</b>	0.566	748
	GeoCLAP	audio	0.542	809	0.553	802
	PSM	audio	<b>0.586</b>	<b>701</b>	<b>0.595</b>	<b>664</b>
	Ours	audio	0.555	765	0.566	748
3	GeoCLAP	none	0.454	1200	0.456	1197
	PSM	none	0.479	1086	0.487	1042
	Ours	none	<b>0.559</b>	<b>776</b>	<b>0.561</b>	<b>763</b>
	GeoCLAP	query	0.542	840	0.555	790
	PSM	query	0.799	159	<b>0.604</b>	<b>657</b>
	Ours	query	<b>0.910</b>	<b>81</b>	0.577	729
5	GeoCLAP	audio	0.548	812	0.555	790
	PSM	audio	<b>0.594</b>	<b>676</b>	<b>0.604</b>	<b>657</b>
	Ours	audio	0.561	757	0.577	729
	GeoCLAP	none	0.458	1194	0.457	1184
	PSM	none	0.459	1172	0.465	1138
	Ours	none	<b>0.545</b>	<b>804</b>	<b>0.560</b>	<b>774</b>
5	GeoCLAP	query	0.542	835	0.554	791
	PSM	query	0.796	158	<b>0.584</b>	<b>711</b>
	Ours	query	<b>0.909</b>	<b>82</b>	0.566	751
	GeoCLAP	audio	0.550	812	0.554	791
	PSM	audio	<b>0.579</b>	<b>720</b>	<b>0.584</b>	<b>711</b>
	Ours	audio	0.553	784	0.566	751
<i>With Metadata</i>						
1	PSM	query	0.872	142	0.940	104
	Ours	query	<b>0.972</b>	<b>35</b>	<b>0.959</b>	<b>70</b>
	PSM	audio	0.931	123	0.940	104
	Ours	audio	<b>0.956</b>	<b>78</b>	<b>0.959</b>	<b>70</b>
3	PSM	none	0.795	306	0.800	290
	Ours	none	<b>0.857</b>	<b>208</b>	<b>0.858</b>	<b>199</b>
	PSM	query	0.870	150	0.940	104
	Ours	query	<b>0.970</b>	<b>37</b>	<b>0.955</b>	<b>74</b>
5	PSM	audio	0.929	126	0.940	104
	Ours	audio	<b>0.949</b>	<b>83</b>	<b>0.955</b>	<b>74</b>
	PSM	none	0.794	316	0.794	299
	Ours	none	<b>0.846</b>	<b>220</b>	<b>0.851</b>	<b>216</b>
5	PSM	query	0.868	156	0.935	109
	Ours	query	<b>0.969</b>	<b>37</b>	<b>0.954</b>	<b>80</b>
	PSM	audio	0.926	131	0.935	109
	Ours	audio	<b>0.948</b>	<b>88</b>	<b>0.954</b>	<b>80</b>


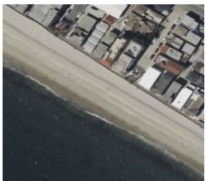
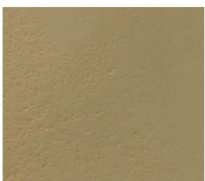

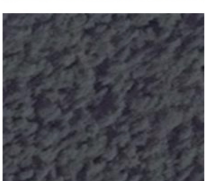
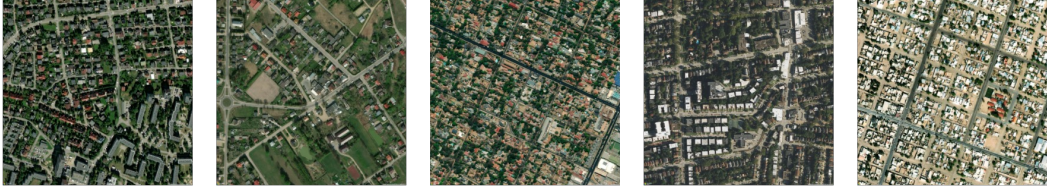
Image	Retrieved Captions
	From the location captured in the aerial view image, we can expect to hear the sounds of people walking, talking, and cheering, as well as the sounds of the game being played on the field.
	From the location captured in the aerial view image, we can expect to hear the sounds of waves crashing on the beach, as well as the occasional noise from the nearby houses and vehicles.
	From the location captured in the aerial view image, we can expect to hear the sounds of the wind blowing across the sandy surface, as well as the occasional rustling of the sand.
	From the aerial view image, we can expect to hear the sounds of cars driving on the street, pedestrians walking on the sidewalks, and the occasional noise from the surrounding buildings.
	From the location captured in this aerial view image, we can expect to hear the sounds of the dense vegetation, such as leaves rustling, branches creaking, and the occasional chirping of birds.

Figure 4: Examples of Top-1 retrieved LLaVA captions for a Bing image by Sat2Sound from our gallery which is the *test*-set of the GeoSound dataset.

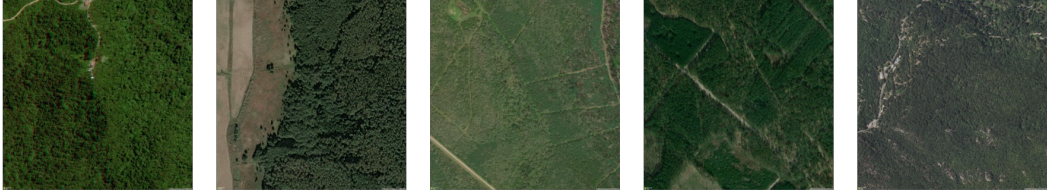
(a) top codebook ids: {13264, 1489, 7856}



(b) top codebook ids: {14500, 13264, 7224}



(c) top codebook ids: {13264, 8407, 7224}



(d) top codebook ids: {13264, 6454, 7224}



Figure 5: Some example groups from the GeoSound test set are shown, where each group shares a common set of highly activated codebook concepts, reflecting similar soundscapes of specific geographic areas. The samples in (a) correspond to residential soundscapes, (b) reflect the soundscape of open fields, (c) represent forested area soundscapes, and (d) capture the soundscape of landscapes with nearby water bodies.