

ASAudio: A Survey of Advanced Spatial Audio Research

Anonymous ACL submission

Abstract

With the rapid development of spatial audio technologies today, applications in AR, VR and other scenarios have garnered extensive attention. Unlike traditional mono sound, spatial audio offers a more realistic and immersive auditory experience. Despite notable progress in the field, there remains a lack of comprehensive surveys that systematically organize and analyze these methods and their underlying technologies. In this paper, we provide a comprehensive overview of spatial audio and systematically review recent literature in the area. To address this, we chronologically outline existing work related to spatial audio and categorize these studies based on input-output representations, as well as generation and understanding tasks, thereby summarizing various research aspects of spatial audio. In addition, we review related datasets, evaluation metrics, and benchmarks, offering insights from both training and evaluation perspectives. Related materials are available at <https://github.com/ASAudio/ASAudio>.

1 Introduction

Spatial audio delivers an immersive, three-dimensional listening experience by simulating how sound propagates and is perceived in space, representing the culmination of audio’s evolution from mono to surround sound (Poeschl et al., 2013). Fueled by its adoption as a core feature in products from Apple, Google, and Meta, the technology has seen accelerated development and widespread application in film, gaming, and the emerging metaverse (Chen et al., 2025; Wuolli and Moreira Kares, 2023; Lee et al., 2023a; Broderick et al., 2018; Murphy and Neff, 2011), which in turn has sharpened the focus of academic research.

As illustrated in Fig. 1, the research landscape of spatial audio has undergone a significant evolution. Before 2021, efforts primarily center on understanding tasks like sound event localization

and detection (SELD) and source separation, dominated by foundational CNN-based models (Zhou et al., 2020; Gao and Grauman, 2019; Wu et al., 2021; Richard et al., 2021; Nguyen et al., 2022; Shimada et al., 2021) and limited by the scale of early datasets (Donley et al., 2021; Morgado et al., 2020). Since 2022, the field has entered a new phase of rapid, synergistic advancement in both understanding and generation, fueled by breakthroughs in generative models and the proliferation of multimodal datasets (Zheng et al., 2024; Zhang et al., 2025; Kim et al., 2025; Sun et al., 2024). This period sees the rise of powerful generation models like ImmerseDiffusion (Heydari et al., 2025) and DiffSAGe (Kushwaha et al., 2025), which drastically improves audio quality and realism. Crucially, the underlying technologies, such as attention mechanisms and large language models, also revolutionize understanding. This propels the task from traditional signal-level analysis toward higher-level semantic reasoning, as seen in advanced models for attention-based separation (Ye et al., 2024) and LLM-based spatial inference (Zheng et al., 2024).

To systematically review these advances in representation, understanding, generation, datasets, and evaluation protocols, this paper is organized as follows: Section 2 discusses input-output representations, Sections 3 and 4 analyze understanding and generation tasks, and Section 5 summarizes existing datasets and evaluation standards.

2 Representations of Spatial Audio

2.1 Inputs Representations

Input representations aim to capture semantic, acoustic, and spatial information. They are provided alone or in combination as mono audio, text, visual signals, or spatial coordinates. We provide a detailed explanation of input representation and their primary processing method in Figure 2.

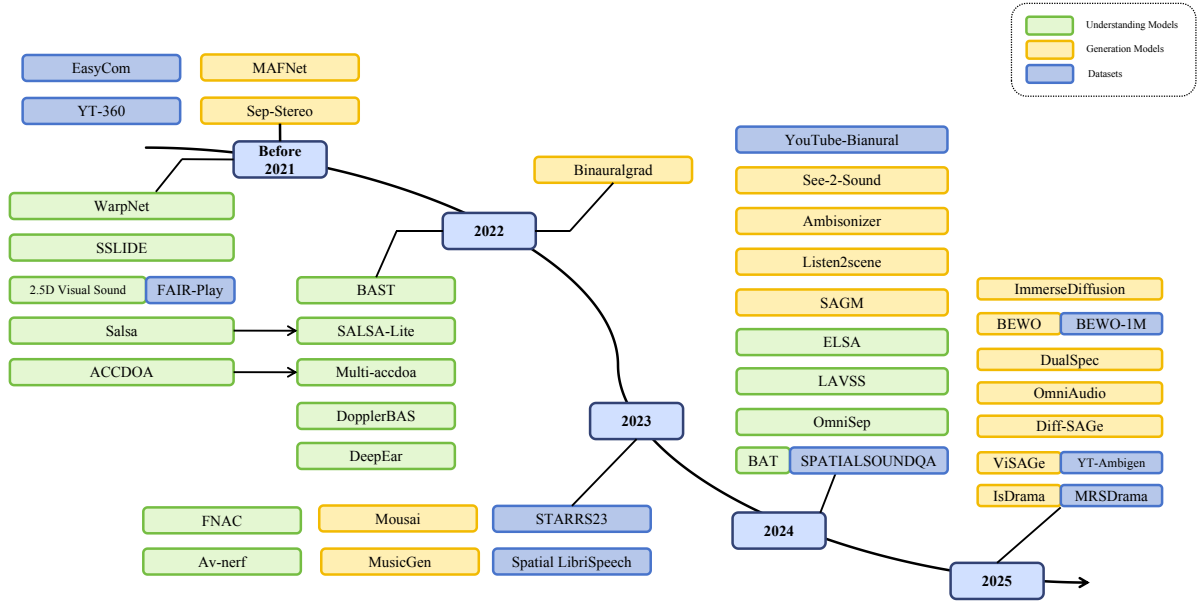


Figure 1: A timeline of recent spatial audio models & datasets in recent years. The timeline is established mainly according to the release date of the technical paper for each model. We mark the understanding models in green and the generation models in yellow, while datasets are marked in blue. Arrows indicate the evolution of models.

Natural Language Prompts Natural language prompts specify semantic content and spatial attributes in an intuitive way. They describe events for generation (Kreuk et al., 2022; Liu et al., 2023) or serve as queries in understanding tasks. For example, BAT (Zheng et al., 2024) uses a large language model to process question-answer pairs about sound event detection, direction estimation, and spatial reasoning, and it extracts spatial information from natural language.

Spatial Position Explicit spatial position data, such as Cartesian or spherical coordinates, provides direct guidance to place sources in generation tasks and serves as ground truth for localization models in understanding tasks. Some studies (Liu et al., 2022; Zhang et al., 2025) also include radial velocity and orientation. They simulate Doppler effects to enhance dynamic properties.

Visual Information Visual information (images or videos) strongly correlates with sound and provides valuable spatial and semantic context. It offers key cues for audio-visual source separation and localization (Zhao et al., 2018; Ye et al., 2024; Zhou et al., 2018) and for audio-visual acoustic matching (Chen et al., 2022). It also guides mono-to-spatial generation (Gan et al., 2019; Gao and Grauman, 2019) and video-to-spatial-audio genera-

tion (Liu et al., 2025a) tasks.

Monoaural Audio Mono audio serves as the base acoustic content in many generation tasks. It supplies core timbral and spectral cues. In two-stage systems, the mono stream is first processed and then “upmixed” into multichannel or binaural formats under the guidance of spatial inputs such as visuals or positions information.

2.2 Spatial Cues and Physical Modeling

A core aspect of spatial audio is the accurate modeling of sound propagation and perception in three-dimensional space, with two key concepts being the room impulse response (RIR) and the head-related transfer function (HRTF).

Room Impulse Response (RIR) The room impulse response (RIR) characterizes all acoustic paths from a source to a receiver, bridging virtual and real acoustics. As direct measurement is costly, research has focused on alternatives. Some methods estimate RIRs from visual inputs to avoid acoustic measurements (Kim et al., 2019; Ratnarajah et al., 2024; Majumder et al., 2022), while others use simulation tools to generate data for training and improving tasks like source separation (Roman et al., 2024; Ahn et al., 2023; Jeub et al., 2009; Vacher et al., 2014; Mittag et al.,

2017; Di Carlo et al., 2021; Grondin et al., 2020; Xu et al., 2021). To support complex applications, precise RIRs have been measured for specific scenarios like dense grids or dynamic sources (Koyama et al., 2021; Ratnarajah et al., 2022; Politis et al., 2020; McKenzie et al., 2021b,a), and perceptual evaluation often relies on measured binaural RIRs (BRIRs) to assess synthesis authenticity (Brinkmann et al., 2017).

Head-Related Transfer Function (HRTF)

Head-related transfer function (HRTF) is a subject-specific filter describing how an individual’s anatomy alters incoming sound, encoding the binaural and monaural cues essential for 3D perception. Because HRTFs are highly individualized, personalization is critical to avoid perceptual artifacts like in-head localization and front-back confusion. To this end, researchers have developed several methods. Some predict HRTFs from anthropometric features like ear shape using neural networks (Warnecke et al., 2022; Arbel et al., 2024; Zhao et al., 2022). Others select the best-matching HRTF from a database, guided by perception-aligned metrics (Lee et al., 2023b; Marggraf-Turley et al., 2024). The most mainstream approach, however, is spatial upsampling from sparse data, which uses deep models to interpolate a full HRTF from a few measurements. This includes using various deep architectures like CNNs and Transformers for reconstruction (Jiang et al., 2023; Ito et al., 2022; Hogg et al., 2024; Ma et al., 2023; Zhang et al., 2023), incorporating physical priors to improve performance (Chen et al., 2023; Thuillier et al., 2024), and leveraging neural fields to represent HRTFs as continuous functions (Zhang et al., 2023; Masuyama et al., 2024). Future work aims to fuse these methods and deploy them on consumer devices (Warnecke et al., 2022; Jiang et al., 2023).

2.3 Output Representations

Spatial audio is mainly represented in three formats. Channel-based formats (e.g., 5.1 or 7.1 surround) assign signals to predefined loudspeaker positions. Scene-based formats (e.g., higher-order Ambisonics (HOA)) represent the full three-dimensional sound field using spherical harmonic decomposition. Object-based formats, such as Dolby Atmos, treat each source as an independent object with positional metadata and render it dynamically at playback. We analyze three output paradigms and

discuss binaural rendering separately.

Channel-Based Audio Channel-based audio maps signals to predefined loudspeaker positions, such as stereo, 5.1, or 7.1. Spatial position is implied by level and time differences across channels. The psychoacoustic basis is summing localization. Amplitude panning follows the sine law:

$$\sin \theta_I = \frac{E_L - E_R}{E_L + E_R} \sin \theta_0. \quad (1)$$

This paradigm is widely used but depends on standardized layouts. It has a small “sweet spot” and limited flexibility and scalability.

Scene-Based Audio Scene-based audio aims to capture and physically reproduce the entire sound field within a region. Key methods include higher-order Ambisonics (HOA) and wave field synthesis (WFS). Ambisonics represents the 3D field by spherical harmonic decomposition:

$$P(\mathbf{x}, \omega) = \sum_{n=0}^N \sum_{m=-n}^n \alpha_n^m(\omega) j_n(kr) Y_n^m(\hat{\mathbf{x}}). \quad (2)$$

This paradigm produces a wide and stable listening area. However, it places high demands on the system, which limits adoption in the consumer market.

Object-Based Audio Object-based audio treats each source as an independent audio object that carries content and metadata, such as position and trajectory. The final mix is rendered in real time on the playback device. Dolby Atmos is a representative system. By decoupling content from the physical playback setup, this paradigm achieves strong scalability and interactivity and becomes a core of next-generation immersive media.

Binaural Audio Binaural audio is a key rendering method and the final form that delivers advanced spatial formats to the ears over headphones. It uses HRTFs to reconstruct the ear-canal pressure and thus tricks the brain into perceiving a 3D scene. Convincing experiences require dynamic head tracking and room acoustics (reverberation) modeling. These components reduce front-back confusion and promote externalization.

2.4 Representation Discussion

Input representations We observe three axes that govern design choices: (i) **Abstraction vs Control precision**. Natural language and vision information are human-friendly and scalable for

high-level intent, but suffer from ambiguity and lower precision; spatial coordinates deliver exact, reproducible control but lack semantics and are tedious to author. (ii) **Semantics vs Geometry**. High-level intents require an interpretation layer (often an LLM or structured parsers) to map semantics to machine-executable spatial parameters; geometric inputs bypass this layer but reduce expressivity. (iii) **Content vs Spatialization**. Monaural audio supplies core acoustic content (timbre, pitch), while other modalities guide spatial rendering; a two-stage pipeline (content generation, then spatialization) yields modularity and controllability. We make a concise comparison and discussion in Appendix A.1, Table 1.

Output representations The three output forms differ in device dependence, scalability, listening freedom, and playback-side complexity. Channel-based formats have high device dependence but low playback complexity. Scene-based formats offer high listening freedom but place strict demands on the system. Object-based formats provide unmatched flexibility and scalability and act as a core driver of next-generation immersive media. These paradigms are not mutually exclusive, and each suits different applications best. A concise comparison is deferred to Appendix A.2, Table 2.

3 Understanding Approaches

Spatial audio understanding aims to analyze complex acoustic scenes by exploiting spatial cues. Core tasks include sound event localization and detection (SELD), spatial audio separation, and joint learning with visual and language modalities.

3.1 SELD Tasks

Sound event localization and detection (SELD) answers two questions at once: what sound occurs (sound event detection, SED) and where it comes from (direction of arrival estimation, DOAE). Traditional methods rely on signal processing while modern work increasingly adopts deep learning models on SELD tasks.

Deep learning achieves strong progress on SELD with diverse network architectures. Early work (May et al., 2010) models binaural cues (ITD/ILD) with Gaussian mixtures to estimate azimuth and lays the foundation for later studies. SELDnet (Adavanne et al., 2018a) uses a CRNN to process SED and DOA in parallel and becomes a key baseline. To further improve performance, researchers

explore alternative representations and mappings. For example, (Pavlidis et al., 2015) estimates the active intensity vector, while (Rana et al., 2019) builds an automated pipeline for Ambisonics estimation from audio–visual features. For binaural devices such as hearing aids, DeepEar (Yang and Zheng, 2022) designs a multi-sector network that localizes multiple sources. To handle unknown numbers of sources in the wild, (Kim et al., 2023) proposes a YOLO-inspired, event-driven localizer that is robust to concurrent events.

Jointly learning SED and DOA often degrades performance. Several strategies address this issue. (Cao et al., 2019) shows that two-stage training allows SED features to benefit DOAE. (Cao et al., 2021) introduces a track-wise output format, permutation-invariant training (PIT), and soft parameter sharing to avoid sacrificing subtask accuracy. (Shimada et al., 2021, 2022) proposes ACCDOA and its multi-target extension, which unify SELD as a single-target regression problem and remove the need to balance multi-task losses. SALSA (Nguyen et al., 2022) designs a joint time–frequency feature that maps signal energy and directional cues with high precision.

To fuse complementary strengths, (Yasuda et al., 2020) combines physics-based intensity vector (IV) estimation with DNN denoising and source separation to handle overlaps. With listener motion, (Krause et al., 2023) confirms the benefit of motion cues for localization, and (García-Barrios et al., 2022) analyzes how head rotations affect accuracy. In model design, self-supervised methods (Sun et al., 2023; Santos et al., 2024) and audio–visual learning (Gan et al., 2019; Tian et al., 2018) reduce dependence on large labeled sets. Recent architectures including CRNNs with SE modules (Naranjo-Alcazar et al., 2020), Transformers (Kuang et al., 2022), autoencoders (Huang et al., 2020; Wu et al., 2021), and VAEs (Bianco et al., 2020, 2021) capture time–frequency structure and support unsupervised or semi-supervised settings.

3.2 Spatial Audio Separation

Source separation aims to recover individual sources from a mixture. With binaural or multi-channel inputs, inter-channel spatial cues provide strong leverage, especially for challenging “cocktail party” scenarios.

Binaural Audio Separation Binaural separation uses ITD and ILD cues between the two ears to dis-

entangle overlapping sources. Early machine learning approaches, such as (Weiss et al., 2009), employ probabilistic models. To support human–robot interaction, (Deleforge and Horaud, 2012) proposes a generative model with active binaural hearing so that a robot performs robust separation and localization in cocktail-party conditions. To handle multi-speaker separation under reverberation, (Zhang and Wang, 2017) introduces a novel 2D ITD feature, while (Wang and Wang, 2018) tightly integrates spectral and spatial features in a deep framework. To preserve spatial cues that matter to downstream applications, (Han et al., 2020) proposes MIMO TasNet for real-time speech separation with binaural cue retention.

Audio–visual fusion is another major direction. The pioneering 2.5D Visual Sound (Gao and Grauman, 2019) adopts a mix-and-separate strategy, where visual cues guide binaural separation. To go beyond systems that only model acoustics and ignore spatial position, LAVSS (Ye et al., 2024) introduces audio–visual spatial source separation (AVSS). It encodes object locations explicitly to steer the separation process.

Multichannel Audio Separation Multichannel separation uses richer spatial information and array geometry to address the underdetermined case where sources outnumber channels. Traditional methods such as spatial clustering (Wang et al., 2018) cluster time–frequency bins with GMMs using inter-channel cues (ITD, ILD, etc.). Early DNN work (Nugraha et al., 2016) combines DNN-modeled spectra with a classical multichannel Gaussian model to exploit spatial structure. Recent unsupervised methods, such as (Zmolikova et al., 2021), adopt variational Bayes to unify spectral and spatial cues and achieve end-to-end spatial separation. In addition, (Wang et al., 2018) proposes an efficient algorithm that extends two-channel deep clustering to arbitrary microphone arrays. Similarly, (Morgado et al., 2018) converts mono audio to multichannel spatial audio via video analysis and implicitly separates and localizes unknown sources.

3.3 Joint Learning

To reach comprehensive scene understanding, spatial audio is increasingly learned together with other modalities, such as vision and natural language. The goal is to align and exploit the rich cues present across modalities.

Alignment Between Audio & Visual Information Aligning spatial audio with vision is key to cross-modal reasoning. (Morgado et al., 2020) propose audio–visual spatial alignment (AVSA) and use contrastive learning to capture correspondences between 360° videos and their spatial audio. (Yang et al., 2020) design a self-supervised task that asks the model to detect whether left–right audio channels are swapped. This task forces the model to learn spatial correspondence between audio modality and video modality.

Environment Information Understanding room acoustics is essential for realistic reproduction. (Liang et al., 2023) integrates propagation priors into NeRF to synthesize spatial audio consistent with novel views. (Luo et al., 2022) proposes neural acoustic fields (NAFs) that learn an implicit representation of sound propagation directly from impulse responses. Many studies (Savioja and Svensson, 2015; Ratnarajah et al., 2024; Bryan, 2020; ISO, 2009; Coldenhoff et al., 2024; Majumder et al., 2022; Srivastava et al., 2021) simulate or measure room impulse responses to analyze indoor acoustic parameters and capture geometry and material properties.

Visual Segmentation & Depth Estimation Depth and segmentation provide precise geometric supervision for spatial audio processing. (Liu et al., 2025b) integrates YOLOv8 (Varghese and Sambath, 2024) detection with Depth Anything to estimate depth. It then computes accurate 3D source positions and supplies key cues for downstream spatialization.

Natural Language Guided Natural language guidance is a new frontier for spatial audio understanding. Because existing audio foundation models usually lack spatial awareness, ELSA (Devnani et al., 2024) uses contrastive learning and spatial regression targets to align spatial audio with text for the first time. BAT (Zheng et al., 2024) builds a new dataset, SPATIALSOUNDQA, with spatial question–answer pairs and fine-tunes a large language model (LLaMA-2). It shows the strong potential of LLMs for spatial audio reasoning.

3.3.1 Future Work

Future work moves beyond simple acoustic-event perception toward higher-level cognitive scene analysis. It aims to develop unified models that reason about causality in complex acoustic environ-

ments. Spatial-audio understanding now undergoes a profound shift from perception to cognition. As a result, models not only process acoustic signals but also infer context, relations, and causality within the entire sound scene.

4 Spatial Audio Generation Methods

Spatial audio generation evolves from traditional digital signal processing to advanced deep learning methods. This progress is driven by rapid advances in generative models. This section reviews recent developments, covering both cascade models and end-to-end models. A summary of recent deep learning models is presented in the Appendix B and Table 3 with their input/output format and model framework.

4.1 Cascade Models

This part focuses on a core topic in spatial audio generation: upmixing monaural audio into binaural audio with three-dimensional spatial cues. The “mono-to-binaural” process builds immersive listening. It aims to reproduce the spatial cues that humans perceive and traces the technical path from structured physical models to deep, especially vision-guided, frameworks.

Traditional Methods Humans localize sound with binaural hearing. This mechanism involves an ITD, an ILD, and spectral changes described by HRTF. Early work such as (Brown and Duda, 1998) explicitly models wave propagation and diffraction with a simplified time-domain description. The model is interpretable and efficient. With deep learning, (Richard et al., 2021) introduces a neural rendering network that synthesizes binaural waveforms from a mono input and the listener position. The work shows the limits of a plain L_2 loss on raw waveforms.

Visually Guided Audio Spatialization A mono signal lacks spatial location information. Visually guided spatialization uses synchronized video to provide key context. The pioneering 2.5D Visual Sound framework (Gao and Grauman, 2019) employs a deep convolutional network to recover spatial cues and sets the basic paradigm. (Li et al., 2024c) adds object-level visual cues and designs a cyclic locate-and-upmix (CLUP) framework. It jointly learns visual source localization and binaural generation. To improve accuracy, researchers add 3D geometry. (Parida et al., 2022) stresses

depth maps and designs an encoder–decoder with hierarchical attention. (Garg et al., 2021) separates geometry cues with a multi-task network and learns geometry-aware features. Efficient cross-modal fusion becomes a focus. (Zhang and Shao, 2021) proposes the multi-attention fusion network (MAFNet). (Liu et al., 2024) adds a novel audio–visual matching loss. (Zheng et al., 2022) defines a “binaural ratio” linked to physical cues to improve interpretability. (Li et al., 2024b) introduces a GAN framework with shared visual guidance and proposes a new spatial metric.

Audio Quality Enhancement After solving localization, another line improves audio fidelity and physical realism. (Leng et al., 2022) first applies diffusion. It generates shared and ear-specific information in two stages. (Liu et al., 2022) adds a plug-and-play DopplerBAS module that uses radial velocity to handle Doppler effects. (Lee and Lee, 2023) proposes the Neural Fourier Shift (NFS) network, which renders in the Fourier domain and predicts early reflections, cutting computation.

Weakly-Supervised/Self-Supervised Paradigms To break data limits, researchers propose new learning paradigms. (Xu et al., 2021) creates PseudoBinaural. It uses physical priors to make pseudo labels from many mono videos. (Rachavarapu et al., 2021) uses source localization as a proxy task for weak supervision. Multi-task and self-supervised learning also help. Sep-Stereo (Zhou et al., 2020) adds visual-guided separation as a second task. (Lin and Wang, 2021) enforces left–right consistency. (Li et al., 2021) adds a channel-flip classification task for self-supervision.

4.2 End-to-End Models

End-to-end spatial audio generation no longer upmixes an existing mono track. It synthesizes a complete sound field from high-dimensional, multi-modal inputs such as silent video, natural language, or 3D geometry. The rise of diffusion models, large multimodal datasets, and cross-modal representation learning (e.g., CLIP) drives this paradigm. Early systems include the VQ-VAE framework in (Huang et al., 2022) and the surround-to-binaural network in (Yang et al., 2022).

Video-Driven Spatial Audio Generation The video-driven generation paradigm turns AI from a post-production tool into a creative engine. ViS-AGe (Kim et al., 2025) generates first-order Am-

bisonics (FOA) from silent video and surpasses cascade methods. With VR/AR, generating immersive audio for 360° videos becomes important. Omni-Audio (Liu et al., 2025a) tackles the 360V2SA task with a dual-branch design that uses panoramic and normal views. Other work (Rana et al., 2019; Liang et al., 2023) estimates 3D source positions from audio–visual cues and encodes them in panoramic sound.

Text and Multimodal Conditioned Generation

Controlling spatial audio with natural language is a cutting-edge direction. Diffusion models drive this change. (Heydari et al., 2025) uses a latent diffusion model to produce 3D immersive soundscapes from text. It supports descriptive and parametric control. (Sun et al., 2024) notes that plain text embeddings blur spatial cues. It proposes Spatial-Sonic, which adds a spatial encoder and an azimuth-elevation matrix for explicit guidance. Architectural innovation then improves controllability. DualSpec (Zhao et al., 2025) introduces a pretrained separator and a channel-shift loss to enhance spatialization. Other studies, such as (Kushwaha et al., 2025; Zang et al., 2024), generate FOA from class labels and positions or directly from text. The trend extends to complex dialog and music. IS-Drama (Zhang et al., 2025) accepts scripts, video, and pose and produces multi-speaker spatial dialog with dramatic prosody. MusicGen (Copet et al., 2023), Moûsai (Schneider et al., 2023), and (Evans et al., 2024b) generate high-quality stereo music from text input.

Environmental Acoustic Modeling For higher realism and interactivity, research splits into two philosophies: holistic and compositional. Environmental acoustic modeling represents the holistic view. (Ratnarajah and Manocha, 2024) renders sound for a 3D scene with a graph neural network that encodes material and geometry. (Kim et al., 2019) estimates room geometry and acoustics from 360° images to synthesize scene-aware audio. Modular and zero-shot generation illustrates the compositional view. SEE-2-Sound (Dagli et al., 2024) breaks the visual-to-audio task into region recognition, 3D localization, mono generation, and spatialization. The modular design lets the system produce matching spatial audio for novel visual content and shows strong generalization.

4.3 Future Work

Despite significant progress in both cascade and end-to-end spatial audio generation models, several challenges remain. First, improving computational efficiency and real-time performance while maintaining quality is crucial. Second, achieving high-quality spatial audio generation without high-quality annotated data is another important research direction. Finally, deeper integration of spatial audio with other modalities is a key trend for future research.

5 Dataset and Evaluation of Spatial Audio

5.1 Datasets

Spatial audio data exists in a variety of formats, each reflecting different characteristics and tailored to specific tasks. This section provides an in-depth analysis of existing spatial audio datasets, illustrating the diverse methods of data collection and processing, and explaining how these elements contribute to the understanding of spatial audio. Sources including real-world recordings, physics-based simulations, and web-crawled material are shown in Appendix C and Table 4.

5.1.1 Multi-Channel Audio Datasets

Multi-channel datasets are crucial for developing far-field speech interaction and scene analysis systems. Early corpora like REVERB Challenge (Kinoshita et al., 2016), DIRHA (Ravanelli et al., 2015), and Sweet-Home (Vacher et al., 2014) focus on speech enhancement and ASR in reverberant home environments. To support more precise spatial hearing research, datasets such as Voice-Home (Bertin et al., 2016), SECL-UMons (Brousmiche et al., 2020), and AVRI (Qian et al., 2022) provide detailed geometric annotations for localization and speaker tracking. Recent efforts capture dynamic and complex scenes, including pedestrian environments in the Wearable SELD dataset (Nagatomo et al., 2022) and diverse indoor/outdoor settings in the high-channel-count RealMAN dataset (Yang et al., 2024).

5.1.2 First-Order Ambisonics Datasets

First-Order Ambisonics (FOA) is a standard format for tasks requiring 3D acoustic information, with datasets collected via crawling, simulation, and real-world recording. Crawled datasets like YT-ALL (Morgado et al., 2018) and YT-360 (Morgado et al., 2020) provide large-scale, in-the-wild data

for pre-training, while YT-AMBIGEN (Kim et al., 2025) improves alignment by filtering for camera metadata. Simulated datasets, including the TUT Sound Events series (Adavanne et al., 2018a) and DCASE2021 Task 3 (Politis et al., 2021), offer controlled benchmarks for SELD, whereas Spatial LibriSpeech (Sarabia et al., 2023) and SonicSet (Li et al., 2024a) spatialize large existing corpora. Scarce but highly realistic recorded datasets like REC-STREET (Morgado et al., 2018) and the STARSS series (Politis et al., 2022; Shimada et al., 2023) provide invaluable data for outdoor scenes and high-resolution SELD benchmarks.

5.1.3 Binaural Datasets

Binaural audio offers a perceptually plausible format for headphone-based immersion by directly mimicking human hearing. Real-world recordings capture naturalistic scenes, from musical performances in FAIR-Play (Gao and Grauman, 2019) to challenging noisy conversations in EasyCom (Donley et al., 2021) and head-tracked dialogues in the dataset by Richard et al. (Richard et al., 2021). Simulated datasets like SimBinaural (Garg et al., 2023) enable large-scale, controllable data generation, while hybrid approaches like YouTube-Binaural (Garg et al., 2023) convert existing surround audio to a pseudo-binaural format. Recent efforts integrate richer multimodal and semantic information, with BEWO-1M (Sun et al., 2024) enabling text-guided generation and MRSDrama (Zhang et al., 2025) providing a unique corpus of expressive spatial speech for narrative tasks.

5.2 Evaluation Metrics

5.2.1 Evaluation Metrics for Understanding

SELD Evaluation covers SED and DOA estimation. SED uses segment-based F-score and error rate (ER) (Mesaros et al., 2016). DOA uses two frame-wise metrics: DOA error, which measures the angular deviation between estimates and references, and frame recall, which measures the fraction of frames with the correct number of detected sources (Adavanne et al., 2018b). DOA error averages the assignment cost between reference DOAs DOA_R^t and estimated DOAs DOA_E^t based on the Hungarian algorithm.

Spatial Audio Separation Separation quality is measured with `mir_eval` metrics such as signal-to-distortion ratio (SDR) and signal-to-interference ratio (SIR) (Ye et al., 2024).

Joint Learning For audio–visual tasks, evaluation often uses binary classification metrics, such as audio–visual correspondence (AVC-Bin) and audio–visual spatial alignment (AVSA-Bin) (Morgado et al., 2020). Downstream tasks, such as semantic segmentation, use pixel accuracy and mean Intersection over Union (mIoU).

5.2.2 Evaluation Metrics for Generation

Monaural-to-Binaural Audio Generation Fidelity is evaluated with objective measures in the time domain (Wave L_2), spectral domain (Amplitude L_2 , Phase L_2 , multi-resolution STFT loss), and perceptual scores (PESQ, MOS) (Leng et al., 2022; Liu et al., 2022). The multi-resolution STFT loss (MRSTFT) combines spectral convergence \mathcal{L}_{SC} and log-magnitude loss \mathcal{L}_{mag} .

End-to-End Binaural Audio Generation Evaluation focuses on key spatial cues. Objective metrics include mean absolute error (MAE) of interaural phase difference (IPD) and interaural level difference (ILD) (Zhang et al., 2025). Perceptual evaluation often measures cosine similarity between angle/distance embeddings from a pretrained model (e.g., SPATIAL-AST (Zheng et al., 2024)) and those from generated audio.

End-to-End FOA Generation Evaluation combines spatial accuracy, codec quality, and perceptual plausibility (Heydari et al., 2025). Spatial accuracy reports errors of azimuth (θ), elevation (ϕ), and distance (d), which are derived from the intensity vector of FOA channels. The overall spatial-angle error $\Delta_{\text{Spatial-Angle}}$ is also reported (Van Brummelen, 2012). Codec quality uses STFT and Mel distances. Plausibility uses Fréchet Audio Distance (FAD) and KL divergence. The CLAP score measures consistency between text prompts and generated audio.

Detailed formulas are presented in Appendix D.

6 Conclusion

This paper presents a comprehensive survey of the rapidly advancing spatial audio field covering foundational spatial audio input and output representations; the core research paradigms of understanding and generation; and the landscape of datasets and evaluation metrics. We hope this survey serves as a valuable resource for researchers, further guiding future work and fostering innovation in immersive audio technology.

Limitations

While this survey provides a broad overview of the algorithmic and data-centric aspects of spatial audio, its scope has certain limitations, leaving several important areas underexplored.

First, our review is heavily centered on software, models, and datasets, with only a cursory treatment of the specialized hardware that underpins the entire spatial audio pipeline. We do not offer a detailed analysis of different microphone array geometries (e.g., spherical, tetrahedral), the design of dedicated audio processors (DSPs) for real-time rendering, or the technologies behind head-tracking sensors (e.g., IMUs) and their integration into consumer devices. A deeper dive into these hardware components would be necessary for a complete picture of the field's engineering challenges.

Second, while we touch upon perceptual concepts like HRTF personalization and evaluation metrics like MOS, the survey does not delve deeply into the fundamentals of psychoacoustics and human spatial hearing. A dedicated discussion on the perceptual mechanisms that enable sound localization and immersion would provide crucial context for the engineering solutions presented. Similarly, our section on evaluation metrics focuses extensively on objective, formula-based measures but does not detail the methodologies of subjective listening tests (e.g., MUSHRA, A/B testing), which remain the gold standard for assessing the perceptual quality of spatial audio systems.

Ethical Considerations

The increasing sophistication and accessibility of spatial audio technologies also raise important ethical considerations that the research community must address. The deployment of multi-microphone arrays in personal and public spaces for high-fidelity spatial audio capture creates significant privacy risks. Such systems could be used for covert surveillance, capturing and localizing conversations without consent. Clear guidelines and robust privacy-preserving mechanisms are needed to prevent misuse.

References

Sharath Adavanne, Archontis Politis, Joonas Nikunen, and Tuomas Virtanen. 2018a. Sound event localization and detection of overlapping sources using convolutional recurrent neural networks. *IEEE Journal of Selected Topics in Signal Processing*, 13(1):34–48.

Sharath Adavanne, Archontis Politis, and Tuomas Virtanen. 2018b. Direction of arrival estimation for multiple sound sources using convolutional recurrent neural network. In *2018 26th European Signal Processing Conference (EUSIPCO)*, pages 1462–1466. IEEE.

Byeongjoo Ahn, Karren Yang, Brian Hamilton, Jonathan Sheaffer, Anurag Ranjan, Miguel Sarabia, Oncel Tuzel, and Jen-Hao Rick Chang. 2023. Novel-view acoustic synthesis from 3d reconstructed rooms. *arXiv preprint arXiv:2310.15130*.

Lior Arbel, Ishwarya Ananthabhotla, Zamir Ben-Hur, David Lou Alon, and Boaz Rafaely. 2024. On hrtf notch frequency prediction using anthropometric features and neural networks. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 816–820. IEEE.

Nancy Bertin, Ewen Camberlein, Emmanuel Vincent, Romain Lebarbenchon, Stéphane Peillon, Éric Lamandé, Sunit Sivasankaran, Frédéric Bimbot, Irina Illina, Ariane Tom, and 1 others. 2016. A french corpus for distant-microphone speech processing in real homes. In *Interspeech 2016*.

Michael J Bianco, Sharon Gannot, Efren Fernandez-Grande, and Peter Gerstoft. 2021. Semi-supervised source localization in reverberant environments with deep generative modeling. *IEEE Access*, 9:84956–84970.

Michael J Bianco, Sharon Gannot, and Peter Gerstoft. 2020. Semi-supervised source localization with deep generative modeling. In *2020 IEEE 30th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6. IEEE.

Fabian Brinkmann, Alexander Lindau, and Stefan Weinzierl. 2017. On the authenticity of individual dynamic binaural synthesis. *The Journal of the Acoustical Society of America*, 142(4):1784–1795.

James Broderick, Jim Duggan, and Sam Redfern. 2018. The importance of spatial audio in modern games and virtual environments. In *2018 IEEE games, entertainment, media conference (GEM)*, pages 1–9. IEEE.

Mathilde Brousmiche, Jean Rouat, and Stéphane Dupont. 2020. Secl-umons database for sound event classification and localization. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 756–760. IEEE.

C Phillip Brown and Richard O Duda. 1998. A structural model for binaural sound synthesis. *IEEE transactions on speech and audio processing*, 6(5):476–488.

Nicholas J Bryan. 2020. Impulse response data augmentation and deep neural networks for blind room acoustic parameter estimation. In *ICASSP 2020-2020*

820	<i>IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 1–5. IEEE.	874
821		875
822	Yin Cao, Turab Iqbal, Qiuqiang Kong, Fengyan An, Wenwu Wang, and Mark D Plumbley. 2021. An improved event-independent network for polyphonic sound event localization and detection. In <i>ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 885–889. IEEE.	876
823		877
824		878
825		879
826		880
827		881
828		882
829	Yin Cao, Qiuqiang Kong, Turab Iqbal, Fengyan An, Wenwu Wang, and Mark D. Plumbley. 2019. Polyphonic sound event detection and localization using a two-stage strategy . <i>CoRR</i> , abs/1905.00268.	883
830		884
831		885
832		886
833	Changan Chen, Ruohan Gao, Paul Calamia, and Kristen Grauman. 2022. Visual acoustic matching. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 18858–18868.	887
834		888
835		889
836		890
837	Guodong Chen, Sizhe Wang, Jacob Chakareski, Dimitrios Koutsonikolas, and Mallesham Dasari. 2025. Spatial video streaming on apple vision pro xr headset. In <i>Proceedings of the 26th International Workshop on Mobile Computing Systems and Applications</i> , pages 115–120.	891
838		892
839		893
840		894
841		895
842		896
843	Xingyu Chen, Fei Ma, Yile Zhang, Amy Bastine, and Prasanga N Samarasinghe. 2023. Head-related transfer function interpolation with a spherical cnn. <i>arXiv preprint arXiv:2309.08290</i> .	897
844		898
845		899
846		900
847	Jozef Coldenhoff, Andrew Harper, Paul Kendrick, Tijana Stojkovic, and Milos Cernak. 2024. Multi-channel mosra: Mean opinion score and room acoustics estimation using simulated data and a teacher model. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 381–385. IEEE.	901
848		902
849		903
850		904
851		905
852		906
853		907
854	Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. 2023. Simple and controllable music generation. <i>Advances in Neural Information Processing Systems</i> , 36:47704–47720.	908
855		909
856		910
857		911
858		912
859	Rishit Dagli, Shivesh Prakash, Robert Wu, and Houman Khosravani. 2024. See-2-sound: Zero-shot spatial environment-to-spatial sound. <i>arXiv preprint arXiv:2406.06612</i> .	913
860		914
861		915
862		916
863	Antoine Deleforge and Radu Horaud. 2012. The cocktail party robot: Sound source separation and localisation with an active binaural head. In <i>Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction</i> , pages 431–438.	917
864		918
865		919
866		920
867		921
868	Bhavika Devnani, Skyler Seto, Zakaria Aldeneh, Alessandro Toso, Elena Menyaylenko, Barry-John Theobald, Jonathan Sheaffer, and Miguel Sarabia. 2024. Learning spatially-aware language and audio embeddings. <i>Advances in Neural Information Processing Systems</i> , 37:33505–33537.	922
869		923
870		924
871		925
872		926
873		927
		928
		929
	Diego Di Carlo, Pinchas Tandeyitnik, Cédric Foy, Antoine Deleforge, Nancy Bertin, and Sharon Gannot. 2021. dechorate: a calibrated room impulse response database for echo-aware signal processing. <i>arXiv preprint arXiv:2104.13168</i> .	
	Jacob Donley, Vladimir Tourbabin, Jung-Suk Lee, Mark Broyles, Hao Jiang, Jie Shen, Maja Pantic, Vamsi Krishna Ithapu, and Ravish Mehra. 2021. Easycom: An augmented reality dataset to support algorithms for easy communication in noisy environments. <i>arXiv preprint arXiv:2107.04174</i> .	
	Zach Evans, CJ Carr, Josiah Taylor, Scott H Hawley, and Jordi Pons. 2024a. Fast timing-conditioned latent audio diffusion. In <i>Forty-first International Conference on Machine Learning</i> .	
	Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. 2024b. Long-form music generation with latent diffusion. <i>arXiv preprint arXiv:2404.10301</i> .	
	George Forman and Martin Scholz. 2010. Apples-to-apples in cross-validation studies: pitfalls in classifier performance measurement. <i>Acm Sigkdd Explorations Newsletter</i> , 12(1):49–57.	
	Chuang Gan, Hang Zhao, Peihao Chen, David Cox, and Antonio Torralba. 2019. Self-supervised moving vehicle tracking with stereo sound. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 7053–7062.	
	Ruohan Gao and Kristen Grauman. 2019. 2.5 d visual sound. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 324–333.	
	Guillermo García-Barrios, Daniel Aleksander Krause, Archontis Politis, Annamaria Mesaros, Juana M Gutiérrez-Arriola, and Rubén Fraile. 2022. Binaural source localization using deep learning and head rotation information. In <i>2022 30th European Signal Processing Conference (EUSIPCO)</i> , pages 36–40. IEEE.	
	Rishabh Garg, Ruohan Gao, and Kristen Grauman. 2021. Geometry-aware multi-task learning for binaural audio generation from video. <i>arXiv preprint arXiv:2111.10882</i> .	
	Rishabh Garg, Ruohan Gao, and Kristen Grauman. 2023. Visually-guided audio spatialization in video with geometry-aware multi-task learning. <i>International Journal of Computer Vision</i> , 131(10):2723–2737.	
	François Grondin, Jean-Samuel Lauzon, Simon Michaud, Mirco Ravanelli, and François Michaud. 2020. Bird: Big impulse response dataset. <i>arXiv preprint arXiv:2010.09930</i> .	
	Cong Han, Yi Luo, and Nima Mesgarani. 2020. Real-time binaural speech separation with preserved spatial cues. In <i>ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 6404–6408. IEEE.	

1044	Kai Li, Wendi Sang, Chang Zeng, Runxuan Yang, Guo Chen, and Xiaolin Hu. 2024a. Sonicsim: A customizable simulation platform for speech processing in moving sound source scenarios. <i>arXiv preprint arXiv:2410.01481</i> .	1099
1045		1100
1046		1101
1047		1102
1048		1103
1049	Sijia Li, Shiguang Liu, and Dinesh Manocha. 2021. Binaural audio generation via multi-task learning. <i>ACM Transactions on Graphics (TOG)</i> , 40(6):1–13.	1104
1050		1105
1051		1106
1052	Zhaojian Li, Bin Zhao, and Yuan Yuan. 2024b. Cross-modal generative model for visual-guided binaural stereo generation. <i>Knowledge-Based Systems</i> , 296:111814.	1107
1053		
1054		
1055		
1056	Zhaojian Li, Bin Zhao, and Yuan Yuan. 2024c. Cyclic learning for binaural audio generation and localization. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 26669–26678.	1108
1057		1109
1058		1110
1059		1111
1060		1112
1061	Susan Liang, Chao Huang, Yapeng Tian, Anurag Kumar, and Chenliang Xu. 2023. Av-nerf: Learning neural fields for real-world audio-visual scene synthesis. <i>Advances in Neural Information Processing Systems</i> , 36:37472–37490.	1113
1062		
1063		
1064		
1065		
1066	Yan-Bo Lin and Yu-Chiang Frank Wang. 2021. Exploiting audio-visual consistency with partial supervision for spatial audio generation. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pages 2056–2063.	1114
1067		1115
1068		1116
1069		1117
1070		1118
1071	Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. 2023. Audioldm: Text-to-audio generation with latent diffusion models. <i>arXiv preprint arXiv:2301.12503</i> .	1119
1072		1120
1073		
1074		
1075		
1076	Huadai Liu, Tianyi Luo, Qikai Jiang, Kaicheng Luo, Peiwen Sun, Jialei Wan, Rongjie Huang, Qian Chen, Wen Wang, Xiangtai Li, and 1 others. 2025a. Omniaudio: Generating spatial audio from 360-degree video. <i>arXiv preprint arXiv:2504.14906</i> .	1121
1077		1122
1078		1123
1079		1124
1080		1125
1081	Jinglin Liu, Zhenhui Ye, Qian Chen, Siqi Zheng, Wen Wang, Qinglin Zhang, and Zhou Zhao. 2022. Dopplerbas: Binaural audio synthesis addressing doppler effect. <i>arXiv preprint arXiv:2212.07000</i> .	1126
1082		1127
1083		1128
1084		1129
1085	Miao Liu, Jing Wang, Xinyuan Qian, and Xiang Xie. 2024. Visually guided binaural audio generation with cross-modal consistency. In <i>ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 7980–7984. IEEE.	1130
1086		
1087		
1088		
1089		
1090		
1091	Xiaojing Liu, Ogulcan Gurelli, Yan Wang, and Joshua Reiss. 2025b. Visual-based spatial audio generation system for multi-speaker environments. <i>arXiv preprint arXiv:2502.07538</i> .	1131
1092		1132
1093		1133
1094		1134
1095	Andrew Luo, Yilun Du, Michael Tarr, Josh Tenenbaum, Antonio Torralba, and Chuang Gan. 2022. Learning neural acoustic fields. <i>Advances in Neural Information Processing Systems</i> , 35:3165–3177.	1135
1096		1136
1097		1137
1098		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152

1153	David Murphy and Flaithrí Neff. 2011. Spatial sound	Archontis Politis, Kazuki Shimada, Parthasaarathy	1209
1154	for computer games and virtual reality. In <i>Game</i>	Sudarsanam, Sharath Adavanne, Daniel Krause,	1210
1155	<i>sound technology and player interaction: Concepts</i>	Yuichiro Koyama, Naoya Takahashi, Shusuke Taka-	1211
1156	<i>and developments</i> , pages 287–312. IGI Global Scien-	hashi, Yuki Mitsufuji, and Tuomas Virtanen. 2022.	1212
1157	tific Publishing.	Starss22: A dataset of spatial recordings of real	1213
		scenes with spatiotemporal annotations of sound	1214
1158	Kento Nagatomo, Masahiro Yasuda, Kohei Yatabe,	events. <i>arXiv preprint arXiv:2206.01948</i> .	1215
1159	Shoichiro Saito, and Yasuhiro Oikawa. 2022. Wear-		
1160	able seld dataset: Dataset for sound event localization	Xinyuan Qian, Zhengdong Wang, Jiadong Wang, Guo-	1216
1161	and detection using wearable devices around head.	hui Guan, and Haizhou Li. 2022. Audio-visual	1217
1162	In <i>ICASSP 2022-2022 IEEE International Confer-</i>	cross-attention network for robotic speaker tracking.	1218
1163	<i>ence on Acoustics, Speech and Signal Processing</i>	<i>IEEE/ACM Transactions on Audio, Speech, and Lan-</i>	1219
1164	<i>(ICASSP)</i> , pages 156–160. IEEE.	<i>guage Processing</i> , 31:550–562.	1220
1165	Javier Naranjo-Alcazar, Sergi Perez-Castanos, Jose	Kranthi Kumar Rachavarapu, Vignesh Sundaresha,	1221
1166	Ferrandis, Pedro Zuccarello, and Maximo Cobos.	AN Rajagopalan, and 1 others. 2021. Localize to	1222
1167	2020. Sound event localization and detection using	binauralize: Audio spatialization from visual sound	1223
1168	squeeze-excitation residual cnns. <i>arXiv preprint</i>	source localization. In <i>Proceedings of the IEEE/CVF</i>	1224
1169	<i>arXiv:2006.14436</i> .	<i>International Conference on Computer Vision</i> , pages	1225
		1930–1939.	1226
1170	Thi Ngoc Tho Nguyen, Karn N Watcharasupat,	Aakanksha Rana, Cagri Ozcinar, and Aljosa Smolic.	1227
1171	Ngoc Khanh Nguyen, Douglas L Jones, and Woon-	2019. Towards generating ambisonics using audio-	1228
1172	Seng Gan. 2022. Salsa: Spatial cue-augmented log-	visual cue for virtual reality. In <i>ICASSP 2019-2019</i>	1229
1173	spectrogram features for polyphonic sound event lo-	<i>IEEE International Conference on Acoustics, Speech</i>	1230
1174	calization and detection. <i>IEEE/ACM Transactions on</i>	<i>and Signal Processing (ICASSP)</i> , pages 2012–2016.	1231
1175	<i>Audio, Speech, and Language Processing</i> , 30:1749–	IEEE.	1232
1176	1762.		
		Anton Ratnarajah, Sreyan Ghosh, Sonal Kumar, Purva	1233
1177	Aditya Arie Nugraha, Antoine Liutkus, and Emmanuel	Chiniya, and Dinesh Manocha. 2024. Av-rir: Audio-	1234
1178	Vincent. 2016. Multichannel audio source separa-	visual room impulse response estimation. In <i>Pro-</i>	1235
1179	tion with deep neural networks. <i>IEEE/ACM Trans-</i>	<i>ceedings of the IEEE/CVF Conference on Computer</i>	1236
1180	<i>actions on Audio, Speech, and Language Processing</i> ,	<i>Vision and Pattern Recognition</i> , pages 27164–27175.	1237
1181	24(9):1652–1664.		
		Anton Ratnarajah and Dinesh Manocha. 2024. Lis-	1238
1182	Kranti Kumar Parida, Siddharth Srivastava, and Gaurav	ten2scene: Interactive material-aware binaural sound	1239
1183	Sharma. 2022. Beyond mono to binaural: Generating	propagation for reconstructed 3d scenes. In <i>2024</i>	1240
1184	binaural audio from mono audio with depth and cross	<i>IEEE Conference Virtual Reality and 3D User Inter-</i>	1241
1185	modal attention. In <i>Proceedings of the IEEE/CVF</i>	<i>faces (VR)</i> , pages 254–264. IEEE.	1242
1186	<i>winter conference on applications of computer vision</i> ,		
1187	pages 3347–3356.	Anton Ratnarajah, Zhenyu Tang, Rohith Aralikatti, and	1243
		Dinesh Manocha. 2022. Mesh2ir: Neural acoustic	1244
1188	Despoina Pavlidi, Symeon Delikaris-Manias, Ville	impulse response generator for complex 3d scenes.	1245
1189	Pulkki, and Athanasios Mouchtaris. 2015. 3d lo-	In <i>Proceedings of the 30th ACM International Con-</i>	1246
1190	calization of multiple sound sources with intensity	<i>ference on Multimedia</i> , pages 924–933.	1247
1191	vector estimates in single source zones. In <i>2015 23rd</i>		
1192	<i>European Signal Processing Conference (EUSIPCO)</i> ,	Mirco Ravanelli, Luca Cristoforetti, Roberto Gretter,	1248
1193	pages 1556–1560. IEEE.	Marco Pellin, Alessandro Sosi, and Maurizio Omol-	1249
		ogo. 2015. The dirha-english corpus and related tasks	1250
1194	Sandra Poeschl, Konstantin Wall, and Nicola Doering.	for distant-speech recognition in domestic environ-	1251
1195	2013. Integration of spatial sound in immersive vir-	ments. In <i>2015 IEEE Workshop on Automatic Speech</i>	1252
1196	tual environments an experimental study on effects	<i>Recognition and Understanding (ASRU)</i> , pages 275–	1253
1197	of spatial sound on presence. In <i>2013 IEEE Virtual</i>	282. IEEE.	1254
1198	<i>Reality (VR)</i> , pages 129–130. IEEE.		
		Alexander Richard, Dejan Markovic, Israel D Gebru,	1255
1199	Archontis Politis, Sharath Adavanne, Daniel Krause,	Steven Krenn, Gladstone Alexander Butler, Fernando	1256
1200	Antoine Deleforge, Prerak Srivastava, and Tuomas	Torre, and Yaser Sheikh. 2021. Neural synthesis of	1257
1201	Virtanen. 2021. A dataset of dynamic reverberant	binaural speech from mono audio. In <i>International</i>	1258
1202	sound scenes with directional interferers for sound	<i>Conference on Learning Representations</i> .	1259
1203	event localization and detection. <i>arXiv preprint</i>		
1204	<i>arXiv:2106.06999</i> .	Iran R Roman, Christopher Ick, Sivan Ding, Adrian S	1260
		Roman, Brian McFee, and Juan P Bello. 2024. Spa-	1261
1205	Archontis Politis, Sharath Adavanne, and Tuomas Vir-	tial scaper: a library to simulate and augment sound-	1262
1206	tanen. 2020. A dataset of reverberant spatial sound	scapes for sound event localization and detection in	1263
1207	scenes with moving sources for sound event localiza-	realistic rooms. In <i>ICASSP 2024-2024 IEEE Interna-</i>	1264
1208	tion and detection. <i>arXiv preprint arXiv:2006.01919</i> .	<i>tional Conference on Acoustics, Speech and Signal</i>	1265
		<i>Processing (ICASSP)</i> , pages 1221–1225. IEEE.	1266

1267	Orlem Santos, Karen Rosero, Bruno Masiero, and	aware contrastive learning. In <i>Proceedings of the</i>	1324
1268	Roberto de Alencar Lotufo. 2024. w2v-seld: A	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>	1325
1269	sound event localization and detection framework	<i>tern Recognition (CVPR)</i> , pages 6420–6429.	1326
1270	for self-supervised spatial audio pre-training. <i>IEEE</i>		
1271	<i>Access</i> .		
1272	Miguel Sarabia, Elena Menyaylenko, Alessandro Toso,	Christian Templin, Yanda Zhu, and Hao Wang. 2025.	1327
1273	Skyler Seto, Zakaria Aldeneh, Shadi Pirhosseinloo,	Sonicmotion: Dynamic spatial audio soundscapes	1328
1274	Luca Zappella, Barry-John Theobald, Nicholas Apos-	with latent diffusion models. <i>arXiv preprint</i>	1329
1275	toloff, and Jonathan Sheaffer. 2023. Spatial lib-	<i>arXiv:2507.07318</i> .	1330
1276	rispeech: An augmented dataset for spatial audio		
1277	learning. <i>arXiv preprint arXiv:2308.09514</i> .	Etienne Thuillier, Craig T Jin, and Vesa Välimäki. 2024.	1331
1278	Lauri Savioja and U Peter Svensson. 2015. Overview	Hrtf interpolation using a spherical neural process	1332
1279	of geometrical room acoustic modeling techniques.	meta-learner. <i>IEEE/ACM Transactions on Audio,</i>	1333
1280	<i>The Journal of the Acoustical Society of America</i> ,	<i>Speech, and Language Processing</i> , 32:1790–1802.	1334
1281	138(2):708–730.		
1282	Flavio Schneider, Ojasv Kamal, Zhijing Jin, and Bern-	Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and	1335
1283	hard Schölkopf. 2023. Mo [^] usai: Text-to-music	Chenliang Xu. 2018. Audio-visual event localization	1336
1284	generation with long-context latent diffusion. <i>arXiv</i>	in unconstrained videos. In <i>Proceedings of the Euro-</i>	1337
1285	<i>preprint arXiv:2301.11757</i> .	<i>pean conference on computer vision (ECCV)</i> , pages	1338
1286	Kazuki Shimada, Yuichiro Koyama, Naoya Takahashi,	247–263.	1339
1287	Shusuke Takahashi, and Yuki Mitsufuji. 2021. Acc-		
1288	doa: Activity-coupled cartesian direction of arrival	Michel Vacher, Benjamin Lecouteux, Pedro Chahuara,	1340
1289	representation for sound event localization and de-	François Portet, Brigitte Meillon, and Nicolas Bonne-	1341
1290	tection. In <i>ICASSP 2021-2021 IEEE international</i>	fond. 2014. The sweet-home speech and multimodal	1342
1291	<i>conference on acoustics, speech and signal process-</i>	corpus for home automation interaction. In <i>The 9th</i>	1343
1292	<i>ing (ICASSP)</i> , pages 915–919. IEEE.	<i>edition of the Language Resources and Evaluation</i>	1344
1293	Kazuki Shimada, Yuichiro Koyama, Shusuke Takahashi,	<i>Conference (LREC)</i> , pages 4499–4506.	1345
1294	Naoya Takahashi, Emiru Tsunoo, and Yuki Mitsufuji.		
1295	2022. Multi-acdoa: Localizing and detecting	Glen Van Brummelen. 2012. <i>Heavenly mathematics:</i>	1346
1296	overlapping sounds from the same class with auxil-	<i>The forgotten art of spherical trigonometry</i> . Prince-	1347
1297	iary duplicating permutation invariant training. In	ton University Press.	1348
1298	<i>ICASSP 2022-2022 IEEE international conference</i>		
1299	<i>on acoustics, speech and signal processing (ICASSP)</i> ,	Rejin Varghese and M Sambath. 2024. Yolov8: A	1349
1300	pages 316–320. IEEE.	novel object detection algorithm with enhanced per-	1350
1301	Kazuki Shimada, Archontis Politis, Parthasaarathy Su-	formance and robustness. In <i>2024 International Con-</i>	1351
1302	darsanam, Daniel A Krause, Kengo Uchida, Sharath	<i>ference on Advances in Data Engineering and Intelli-</i>	1352
1303	Adavanne, Aapo Hakala, Yuichiro Koyama, Naoya	<i>gent Computing Systems (ADICS)</i> , pages 1–6. IEEE.	1353
1304	Takahashi, Shusuke Takahashi, and 1 others. 2023.		
1305	Starss23: An audio-visual dataset of spatial record-	Zhong-Qiu Wang, Jonathan Le Roux, and John R Her-	1354
1306	ings of real scenes with spatiotemporal annotations	shey. 2018. Multi-channel deep clustering: Discrim-	1355
1307	of sound events. <i>Advances in neural information</i>	inative spectral and spatial embeddings for speaker-	1356
1308	<i>processing systems</i> , 36:72931–72957.	independent speech separation. In <i>2018 IEEE Inter-</i>	1357
1309	Prerak Srivastava, Antoine Deleforge, and Emmanuel	<i>national conference on acoustics, speech and signal</i>	1358
1310	Vincent. 2021. Blind room parameter estimation us-	<i>processing (ICASSP)</i> , pages 1–5. IEEE.	1359
1311	ing multiple multichannel speech recordings. In <i>2021</i>		
1312	<i>IEEE Workshop on Applications of Signal Processing</i>	Zhong-Qiu Wang and DeLiang Wang. 2018. Combin-	1360
1313	<i>to Audio and Acoustics (WASPAA)</i> , pages 226–230.	ing spectral and spatial features for deep learning	1361
1314	IEEE.	based blind speaker separation. <i>IEEE/ACM Trans-</i>	1362
1315	Peiwen Sun, Sitong Cheng, Xiangtai Li, Zhen Ye,	<i>actions on audio, speech, and language processing</i> ,	1363
1316	Huadai Liu, Honggang Zhang, Wei Xue, and Yike	27(2):457–468.	1364
1317	Guo. 2024. Both ears wide open: Towards language-		
1318	driven spatial audio generation. <i>arXiv preprint</i>	Michaela Warnecke, Sharon Jamison, Sebastian Pre-	1365
1319	<i>arXiv:2410.10676</i> .	pelita, Paul Calamia, and Vamsi Krishna Ithapu. 2022.	1366
1320	Weixuan Sun, Jiayi Zhang, Jianyuan Wang, Zheyuan	Hrtf personalization based on ear morphology. In	1367
1321	Liu, Yiran Zhong, Tianpeng Feng, Yandong Guo,	<i>Audio Engineering Society Conference: 2022 AES</i>	1368
1322	Yanhao Zhang, and Nick Barnes. 2023. Learning	<i>International Conference on Audio for Virtual and</i>	1369
1323	audio-visual source localization via false negative	<i>Augmented Reality</i> . Audio Engineering Society.	1370
		Ron J Weiss, Michael I Mandel, and Daniel PW Ellis.	1371
		2009. Source separation based on binaural cues and	1372
		source model constraints.	1373
		Yifan Wu, Roshan Ayyalasomayajula, Michael J Bianco,	1374
		Dinesh Bharadia, and Peter Gerstoft. 2021. Sslide:	1375
		Sound source localization for indoors based on deep	1376
		learning. In <i>ICASSP 2021-2021 IEEE International</i>	1377
		<i>Conference on Acoustics, Speech and Signal Process-</i>	1378
		<i>ing (ICASSP)</i> , pages 4680–4684. IEEE.	1379

1380	Lauri Wuolio and Elina Moreira Kares. 2023. On the	1434
1381	potential of spatial audio in enhancing virtual user	1435
1382	experiences.	1436
1383	Xudong Xu, Hang Zhou, Ziwei Liu, Bo Dai, Xiaogang	1437
1384	Wang, and Dahua Lin. 2021. Visually informed bin-	1438
1385	aural audio generation without binaural audios. In	
1386	<i>Proceedings of the IEEE/CVF Conference on Com-</i>	
1387	<i>puter Vision and Pattern Recognition</i> , pages 15485–	
1388	15494.	
1389	Bing Yang, Changsheng Quan, Yabo Wang, Pengyu	
1390	Wang, Yujie Yang, Ying Fang, Nian Shao, Hui Bu,	
1391	Xin Xu, and Xiaofei Li. 2024. Realman: A real-	
1392	recorded and annotated microphone array dataset for	
1393	dynamic speech enhancement and localization. <i>Ad-</i>	
1394	<i>vances in Neural Information Processing Systems</i> ,	
1395	37:105997–106019.	
1396	Haici Yang, Sanna Wager, Spencer Russell, Mike Luo,	
1397	Minje Kim, and Wontak Kim. 2022. Upmixing via	
1398	style transfer: a variational autoencoder for disen-	
1399	tangling spatial images and musical content. In <i>ICASSP</i>	
1400	<i>2022-2022 IEEE International Conference on Acous-</i>	
1401	<i>tics, Speech and Signal Processing (ICASSP)</i> , pages	
1402	426–430. IEEE.	
1403	Karren Yang, Bryan Russell, and Justin Salamon. 2020.	
1404	Telling left from right: Learning spatial correspon-	
1405	dence of sight and sound. In <i>Proceedings of the</i>	
1406	<i>IEEE/CVF conference on computer vision and pat-</i>	
1407	<i>tern recognition</i> , pages 9932–9941.	
1408	Qiang Yang and Yuanqing Zheng. 2022. Deeppear:	
1409	Sound localization with binaural microphones. <i>IEEE</i>	
1410	<i>Transactions on Mobile Computing</i> , 23(1):359–375.	
1411	Masahiro Yasuda, Yuma Koizumi, Shoichiro Saito,	
1412	Hisashi Uematsu, and Keisuke Imoto. 2020. Sound	
1413	event localization based on sound intensity vector re-	
1414	fined by dnn-based denoising and source separation.	
1415	In <i>ICASSP 2020-2020 IEEE International Confer-</i>	
1416	<i>ence on Acoustics, Speech and Signal Processing</i>	
1417	<i>(ICASSP)</i> , pages 651–655. IEEE.	
1418	Yuxin Ye, Wenming Yang, and Yapeng Tian. 2024.	
1419	Lavss: Location-guided audio-visual spatial audio	
1420	separation. In <i>Proceedings of the IEEE/CVF Win-</i>	
1421	<i>ter Conference on Applications of Computer Vision</i> ,	
1422	pages 5508–5519.	
1423	Yongyi Zang, Yifan Wang, and Minglun Lee. 2024. Am-	
1424	bisonizer: Neural upmixing as spherical harmonics	
1425	generation. <i>arXiv preprint arXiv:2405.13428</i> .	
1426	Wen Zhang and Jie Shao. 2021. Multi-attention audio-	
1427	visual fusion network for audio spatialization. In	
1428	<i>Proceedings of the 2021 International Conference on</i>	
1429	<i>Multimedia Retrieval</i> , pages 394–401.	
1430	Xueliang Zhang and DeLiang Wang. 2017. Deep learn-	
1431	ing based binaural speech separation in reverberant	
1432	environments. <i>IEEE/ACM transactions on audio,</i>	
1433	<i>speech, and language processing</i> , 25(5):1075–1084.	
	You Zhang, Yuxiang Wang, and Zhiyao Duan. 2023.	1434
	Hrtf field: Unifying measured hrtf magnitude repre-	1435
	sentation with neural fields. In <i>ICASSP 2023-2023</i>	1436
	<i>IEEE International Conference on Acoustics, Speech</i>	1437
	<i>and Signal Processing (ICASSP)</i> , pages 1–5. IEEE.	1438
	Yu Zhang, Wenxiang Guo, Changhao Pan, Zhiyuan Zhu,	1439
	Tao Jin, and Zhou Zhao. 2025. Isdrama: Immersive	1440
	spatial drama generation through multimodal prompt-	1441
	ing. <i>arXiv preprint arXiv:2504.20630</i> .	1442
	Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl	1443
	Vondrick, Josh McDermott, and Antonio Torralba.	1444
	2018. The sound of pixels. In <i>Proceedings of the</i>	1445
	<i>European conference on computer vision (ECCV)</i> ,	1446
	pages 570–586.	1447
	Lei Zhao, Sizhou Chen, Linfeng Feng, Xiao-Lei Zhang,	1448
	and Xuelong Li. 2025. Dualspec: Text-to-spatial-	1449
	audio generation via dual-spectrogram guided diffu-	1450
	sion model. <i>arXiv preprint arXiv:2502.18952</i> .	1451
	Manlin Zhao, Zhichao Sheng, and Yong Fang. 2022.	1452
	Magnitude modeling of personalized hrtf based on	1453
	ear images and anthropometric measurements. <i>Ap-</i>	1454
	<i>plied Sciences</i> , 12(16):8155.	1455
	Tao Zheng, Sunny Verma, and Wei Liu. 2022. Inter-	1456
	pretable binaural ratio for visually guided binaural	1457
	audio generation. In <i>2022 International Joint Confer-</i>	1458
	<i>ence on Neural Networks (IJCNN)</i> , pages 1–8. IEEE.	1459
	Zhisheng Zheng, Puyuan Peng, Ziyang Ma, Xie Chen,	1460
	Eunsol Choi, and David Harwath. 2024. Bat: Learn-	1461
	ing to reason about spatial sounds with large language	1462
	models. <i>arXiv preprint arXiv:2402.01591</i> .	1463
	Hang Zhou, Xudong Xu, Dahua Lin, Xiaogang Wang,	1464
	and Ziwei Liu. 2020. Sep-stereo: Visually guided	1465
	stereophonic audio generation by associating source	1466
	separation. In <i>Computer Vision–ECCV 2020: 16th</i>	1467
	<i>European Conference, Glasgow, UK, August 23–</i>	1468
	<i>28, 2020, Proceedings, Part XII 16</i> , pages 52–69.	1469
	Springer.	1470
	Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and	1471
	Tamara L Berg. 2018. Visual to sound: Generating	1472
	natural sound for videos in the wild. In <i>Proceed-</i>	1473
	<i>ings of the IEEE conference on computer vision and</i>	1474
	<i>pattern recognition</i> , pages 3550–3558.	1475
	Katerina Zmolikova, Marc Delcroix, Lukáš Burget, To-	1476
	mohiro Nakatani, and Jan Honza Černocký. 2021.	1477
	Integration of variational autoencoder and spatial	1478
	clustering for adaptive multi-channel neural speech	1479
	separation. In <i>2021 IEEE Spoken Language Technol-</i>	1480
	<i>ogy Workshop (SLT)</i> , pages 889–896. IEEE.	1481

Appendix

A Extended Representation Discussions

A.1 Input Representations

Figure 2 shows the input modalities for spatial audio tasks, natural language, spatial position, visual information, and monaural audio, each offers a unique perspective for the system to perceive, interpret, or generate soundscapes. While they can be used independently, their true potential is often realized through synergistic multimodal combinations. The choice of input representation is not merely a technical decision but a fundamental architectural one that dictates the system’s capabilities, complexity, and the nature of its interaction with the user or environment. This section will comparatively analyze these input paradigms, examining their intrinsic properties, task suitability, and the emerging trends in their combined application.

As shown in Table 1, these input representations exhibit a core trade-off between the level of abstraction and control precision. Natural language and visual information reside at the highest level of abstraction. They are intuitive for humans and well-suited for high-level scene description or content querying. However, this intuitiveness introduces challenges of lower control precision and semantic ambiguity, necessitating complex models to bridge the gap between semantics and machine-processable signals.

Conversely, spatial position coordinates offer the highest control precision, making them ideal for defining precise source trajectories or serving as ground truth for evaluation. However, they lack semantic context, and manually specifying complex scenes is a tedious process. Monaural audio plays a unique role. Positioned at a low level of abstraction, it does not directly provide spatial control. Instead, it serves as the foundational acoustic content for generation tasks, providing core acoustic features such as timbre and pitch. It acts as raw material that other modalities spatialize.

Therefore, the selection of an input representation is fundamentally a trade-off between the intuitive, abstract control preferred by humans and the precise, geometric data required by machines, a choice contingent on the specific requirements of the task.

Abstract intent vs. geometric precision A fundamental trade-off exists among the different input representations: the opposition between the level

of abstraction in control and its precision. Natural language and visual information represent the pinnacle of abstract, human-centric control. Natural language provides an intuitive way to specify semantic content (e.g., "a bird is chirping") and relational spatial attributes ("on the left"). Similarly, visual information from images or videos offers rich spatial and semantic context. These inputs describe what exists in a scene and how its components are related, which aligns closely with human perception.

However, this intuitiveness comes at the cost of reduced precision. The system must infer precise physical parameters from abstract descriptions. The BAT model (Zheng et al., 2024) exemplifies this challenge, utilizing a large language model to interpret complex natural language queries regarding "sound event detection, direction and distance estimation, and spatial reasoning". This highlights a critical point: high-level abstract inputs require a sophisticated, AI-based interpretation layer to translate human intent into machine-executable instructions.

In contrast, spatial position data provides the highest degree of precision. Cartesian or spherical coordinates offer direct and unambiguous guidance for placing sound sources. This makes it indispensable for tasks requiring absolute accuracy, such as providing ground truth for training and evaluating sound localization models, or simulating precise physical phenomena like the Doppler effect by incorporating velocity vectors. The inherent trade-off is that this representation lacks semantic context and is non-intuitive and tedious for manually specifying complex acoustic scenes.

Monaural audio as the acoustic substrate Unlike other inputs that primarily define where a sound is, monaural audio defines what the sound itself is. It constitutes the "foundational acoustic content" for many spatial audio tasks, providing core acoustic characteristics such as timbre of a specific instrument or the phonetic features of speech. Therefore, monaural audio plays a unique role in the ecosystem of input representations.

Many advanced generative systems follow a two-stage principle: first, a source model (such as AudioGen (Creuk et al., 2022) or AudioLDM (Liu et al., 2023)) generates a monaural audio stream; then, this stream is spatialized or upmixed into a multichannel or binaural format under the guidance of other input modalities, such as visual or posi-

Attribute	Natural Language	Spatial Position	Visual Information	Monaural Audio
Primary Info	Semantic, relational, implicit spatial	Explicit spatial, dynamic	Semantic, spatial, dynamic	Acoustic (timbre, pitch, content)
Control Precision	Low	Very high	High	N/A
Abstraction Level	High	Low	High	Low
Interpretability	Indirect	Direct	Indirect	Indirect
Key Challenges	Ambiguity; semantic-signal gap	No semantics; tedious authoring	Ambiguity; occlusion; compute cost	Lack of spatial cues

Table 1: Comparative analysis of spatial audio *input* representations.

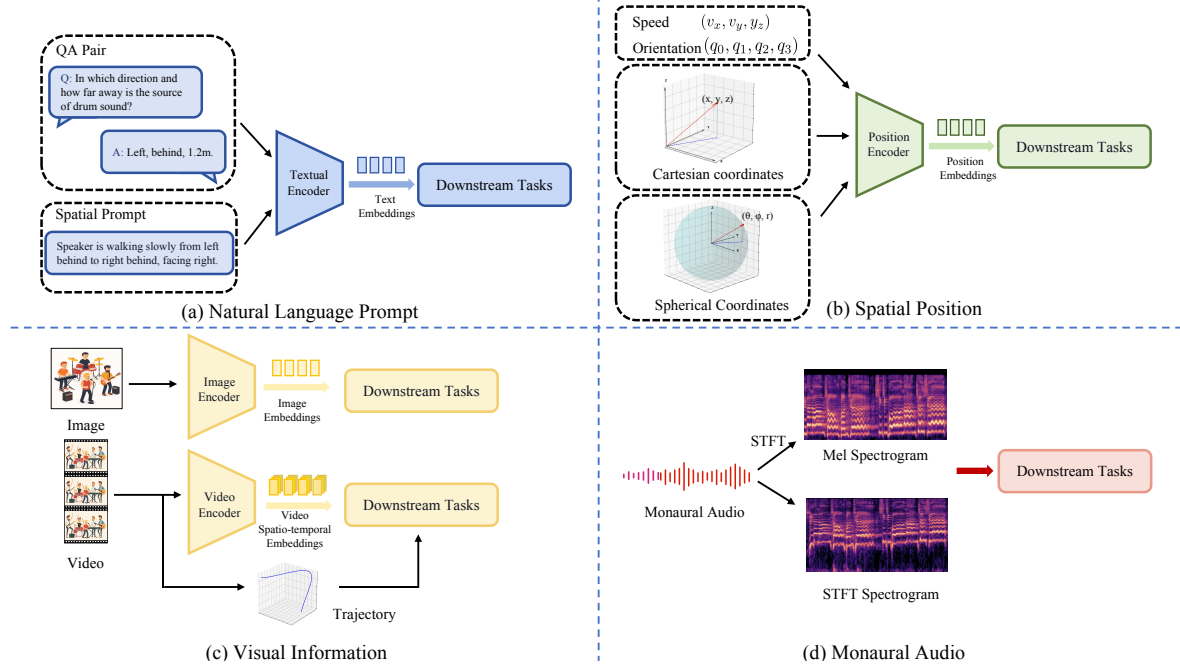


Figure 2: An overview of input representations of spatial audio and their primary processing methods.

tional data. This architecture clearly separates the problem of content generation from that of spatial rendering, enabling modular and flexible system design. Consequently, monaural audio is not an alternative option parallel to other input forms, but rather the fundamental substrate upon which they act.

Multimodal synergy The most powerful spatial audio systems are increasingly moving towards multimodality, creating comprehensive control schemes by combining the strengths of different input types to overcome the limitations of any single modality. The synergy between vision and audio is particularly potent. In audio-visual source separation tasks, the visual presence of an object (e.g., a speaking person) provides a strong, albeit implicit, cue for isolating its corresponding sound from a noisy mixture. In generation tasks, visual in-

formation can guide the spatialization process; for example, a U-Net architecture can take a monaural input and, guided by a video, render a spatially correct binaural or stereo output. The audio-visual matching task is considered crucial, highlighting the deeply learned correspondences between these modalities.

Similarly, adding explicit spatial position data (such as source orientation and velocity) to a monaural audio stream allows for the simulation of highly realistic dynamic effects, like the Doppler shift, elevating realism to a level unattainable with static spatialization.

A.2 Output Representations

Table 2 presents a comparative analysis of the three primary spatial audio output representations. Each paradigm possesses unique advantages and limitations, making it suitable for different application

Attribute	Channel-Based	Scene-Based	Object-Based
Freedom of Listening Position	Limited	High	Moderate
Playback System Dependency	Very high	High	Low
Scalability	Low	Moderate	Excellent
Playback-End Complexity	Low	High	Moderate
Common Formats	Stereo; 5.1/7.1 surround	Ambisonics; wave-field synthesis (WFS)	Dolby Atmos; DTS:X; MPEG-H 3D Audio

Table 2: Comparative analysis of spatial audio *output* representations.

scenarios and user requirements.

Playback system dependency and scalability are key to understanding the evolution of these three paradigms. Channel-based formats exhibit very high system dependency but poor scalability. This is because their audio mix is baked-in for a specific, standardized loudspeaker layout (e.g., 5.1 surround sound). Any playback system that deviates from this layout will degrade the intended spatial effect. In contrast, object-based formats feature low dependency and excellent scalability. They achieve this by decoupling the audio content from its metadata, which allows the playback device to render the audio in real-time according to its own arbitrary loudspeaker configuration. Consequently, a single master file can be adapted to any system. Scene-based formats occupy a middle ground. Their high dependency stems from the requirement for numerous loudspeakers and complex processing systems to physically reconstruct the sound field. Their moderate scalability is demonstrated by the ability to improve performance by increasing the system order (e.g., Higher-Order Ambisonics), though this significantly increases system cost and complexity.

Freedom of listening position and playback-end complexity are directly related to user experience and implementation cost. Channel-based formats confine the listener to a narrow sweet spot, but their playback-end complexity is low, requiring only simple channel-to-loudspeaker mapping. Scene-based formats offer high freedom, allowing listeners to move freely within a designated area. However, this comes at the cost of very high playback-end complexity, which involves real-time decoding and substantial signal processing. Object-based formats provide moderate freedom of movement (depending on the rendering system). Their moderate to high playback-end complexity arises from the need for a real-time rendering engine to process metadata and dynamically generate the mix.

Overall, these three paradigms are not mutu-

ally exclusive; rather, each has its optimal application domain. Channel-based technology retains its place in traditional media due to its simplicity and broad compatibility. Scene-based techniques offer irreplaceable advantages in applications requiring high physical fidelity and large-scale public experiences. Meanwhile, object-based technology, with its unparalleled flexibility and interactivity, has become the core driver for next-generation immersive media, such as VR/AR, gaming, and streaming. Understanding their fundamental differences is crucial for selecting and implementing the most appropriate spatial audio solution.

B Generation Details

This section provides a detailed description of the input and output formats for the generative models summarized in Table 3. These formats represent the diverse ways in which spatial audio systems are controlled and the types of immersive experiences they can produce.

Spatial audio generation has evolved from two-stage upmixing approaches to fully end-to-end synthesis, driven by increasingly powerful deep learning architectures. Early and still prevalent methods, often based on CNNs like U-Net, focus on spatializing existing audio. These models typically take a monaural audio track and visual information from an image or video as input, and output a corresponding binaural or multi-channel audio signal, as seen in pioneering works like 2.5D Visual-Sound (Gao and Grauman, 2019). More recent research has shifted towards direct, end-to-end synthesis from more abstract or multimodal inputs. Diffusion and flow-matching models are at the forefront of this trend, capable of generating high-fidelity FOA or binaural audio directly from text prompts, images, class labels, and explicit spatial positions (e.g., ImmerseDiffusion (Heydari et al., 2025), SonicMotion (Templin et al., 2025), OmniAudio (Liu et al., 2025a)). Transformer-based models excel at

integrating complex, heterogeneous data streams; for instance, ViSAGE (Kim et al., 2025) generates FOA audio from video combined with camera position metadata, while ISDrama (Zhang et al., 2025) synthesizes expressive binaural speech from a rich mix of video, audio, text, and positional data. Other architectures serve specialized functions: VAEs are often used to learn disentangled latent representations for flexible spatial manipulation or to generate intermediate outputs like impulse responses (IRs) from 360° images (Kim et al., 2019), while GANs can incorporate detailed geometric data like 3D meshes to generate physically accurate binaural IRs, as demonstrated by Listen2Scene (Ratnarajah and Manocha, 2024).

C Dataset Details

Spatial audio data exists in a variety of formats, each reflecting different characteristics and tailored to specific tasks. Due to variations in recording equipment and application scenarios, spatial audio data comes in multiple formats, often accompanied by annotations and auxiliary data from other modalities. Moreover, because recording spatial audio is typically costly and resource-intensive, many existing approaches resort to using simulation systems to generate synthetic data from current monaural audio datasets. Some datasets also include real-world spatial audio crawled from the YouTube platform. The following subsections focus on the acquisition and processing methods—both recorded and simulated—associated with various spatial audio formats, including multi-channel audio, First-Order Ambisonics, and binaural audio. A summary of commonly used datasets is presented in the Table 4.

D Evaluation Metrics Details

D.1 Evaluation Metrics for Spatial Audio Understanding

SELD. The SELD task is evaluated using separate metrics for Sound Event Detection (SED) and Direction-of-Arrival (DOA) estimation. For SED, the one-second segment F-score and Error Rate (ER) are commonly used (Mesaros et al., 2016).

For DOA estimation, two frame-wise metrics are frequently employed (Adavanne et al., 2018b): **DOA Error** and **Frame Recall**. Let T be the total number of time frames. Denote by DOA_R^t the set of reference DOAs at frame t and by DOA_E^t the

set of estimated DOAs. Define

$$D_R^t = |\text{DOA}_R^t|, \quad D_E^t = |\text{DOA}_E^t|. \quad (3)$$

The **DOA Error** is defined as

$$\frac{1}{\sum_{t=1}^T D_E^t} \sum_{t=1}^T \text{Hungarian}(\text{DOA}_R^t, \text{DOA}_E^t), \quad (3)$$

where $\text{Hungarian}(\cdot, \cdot)$ denotes the optimal assignment cost computed by the Hungarian algorithm, using as the pairwise cost the central angle between a reference DOA (ϕ_R, λ_R) and an estimated DOA (ϕ_E, λ_E) :

$$\sigma = \arccos(\sin \lambda_E \sin \lambda_R + \cos \lambda_E \cos \lambda_R \cos|\phi_R - \phi_E|). \quad (4)$$

Here $\phi \in [-\pi, \pi]$ is the azimuth and $\lambda \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ is the elevation.

To account for frames where the number of estimated DOAs does not match the number of reference DOAs, the **Frame Recall** is defined as

$$\text{Frame Recall} = \frac{1}{T} \sum_{t=1}^T \mathbf{1}(D_R^t = D_E^t), \quad (5)$$

where $\mathbf{1}(\cdot)$ is the indicator function, equal to 1 if its argument is true and 0 otherwise.

An ideal SELD method achieves an error rate of zero, an F-score of 1 (100%), a DOA Error of 0°, and a Frame Recall of 1 (100%). To compare submitted methods, each method is ranked individually for all four metrics, and final positions are determined by the cumulative minimum of these ranks.

The four cross-validation folds are treated as a single experiment: metrics are computed only after training and testing all folds. Intermediate measures (insertions, deletions, substitutions) are aggregated across folds before calculating the final metrics, rather than averaging per fold (Forman and Scholz, 2010).

Spatial Audio Separation. Metrics to measure the quality of separation, usually adopt the widely-used mir eval library metrics: Signal-to-Distortion Ratio (SDR) measures both interference and artifacts, Signal-to-Interference-Ratio (SIR) measures interference. Higher values indicate a better degree of separation (Ye et al., 2024).

Model	Input Format	Output Format	Framework
Xu et al. 2021	Mono; Image	Binaural	Diffusion-based
Binauralgrad (Leng et al., 2022)	Mono	Binaural	
Moûsai (Schneider et al., 2023)	Text	Binaural	
See-2-Sound (Dagli et al., 2024)	Image; (Text)	<i>Multi</i>	
Evans et al. 2024b	Text; (Audio; Duration)	Binaural	
DualSpec (Zhao et al., 2025)	Text	Binaural	
ImmerseDiffusion (Heydari et al., 2025)	Text; (Position)	FOA	
SonicMotion (Templin et al., 2025)	Text; Position	FOA	
Huang et al. 2022	Mono; Position	Binaural	Latent
Yang et al. 2022	Binaural/ <i>Multi</i>	<i>Multi</i>	
Lee and Lee 2023	Mono; Position; Orientation	Binaural	Transformer-based
MusicGen (Copet et al., 2023)	Text	Mono/Binaural	
Ambisonizer (Zang et al., 2024)	Mono/Binaural	FOA	
ViSAGe (Kim et al., 2025)	Video; Camera Position	FOA	
ISDrama (Zhang et al., 2025)	Video; Audio; Text; Position	Binaural	
OmniAudio (Liu et al., 2025a)	360° Video	FOA	Flow Matching
Diff-SAGe (Kushwaha et al., 2025)	Class Label; Position	FOA	
Listen2Scene (Ratnarajah and Manocha, 2024)	3D Mesh; (Source & Listener Position)	Binaural IRs	GANs
SAGM (Li et al., 2024b)	Mono; Video	Binaural	

Table 3: Comparison of current spatial audio generative models. FOA denotes first-order ambisonics; *Multi* denotes multi-channel audio. Inputs/outputs in parentheses are optional. CNN-based models are omitted.

Joint Learning. In joint learning, they typically employ two binary-classification-based evaluation metrics (Morgado et al., 2020). **AVC-Bin** (Audio–Visual Correspondence) determines whether an audio–video clip pair originates from the same video instance. **AVSA-Bin** (Audio–Visual Spatial Alignment) assesses the spatial consistency between the audio and visual streams.

For semantic segmentation, the model’s dense-prediction capability is evaluated using **pixel accuracy and mean Intersection-over-Union (mean IoU)**. Additionally, **clip-level accuracy** is employed for action recognition.

D.2 Evaluation Metrics for Spatial Audio Generation

Monoaural-to-Binaural Audio Generation. To comprehensively assess the fidelity of the synthesized binaural signal \hat{x} concerning the reference binaural recording x , previous works (Leng et al., 2022; Liu et al., 2022) on **monoaural-to-binaural audio generation** adopt both objective and subjective criteria. Except for the perceptual measures, PESQ and MOS, all metrics are lower-is-better. Notation is unified as follows: $n \in \{1, \dots, T\}$ indicates time-domain samples; $c \in \{L, R\}$ indexes the two output channels; $k \in \{1, \dots, K\}$ and $m \in \{1, \dots, M\}$ denote STFT frequency

and frame indices; $\text{STFT}\{\cdot\}$ yields the complex time–frequency representation.

For **Wave L₂**, The time-domain mean-squared error (MSE) captures sample-by-sample deviations:

$$\mathcal{L}_{L_2}^{\text{wave}} = \frac{1}{T} \sum_{n=1}^T \sum_{c \in \{L, R\}} (\hat{x}_c[n] - x_c[n])^2. \quad (6)$$

Although it provides a well-behaved gradient and is easy to implement, it ignores the non-uniform frequency sensitivity of human hearing.

For **Amplitude L₂**, after converting both signals to their magnitude spectra,

$$\begin{aligned} X(k, m) &= |\text{STFT}\{x\}(k, m)|, \\ \hat{X}(k, m) &= |\text{STFT}\{\hat{x}\}(k, m)|. \end{aligned} \quad (7)$$

The energy envelope mismatch is quantified as

$$\mathcal{L}_{L_2}^{\text{amp}} = \frac{1}{KM} \sum_{k=1}^K \sum_{m=1}^M (\hat{X}(k, m) - X(k, m))^2. \quad (8)$$

For **Phase L₂**, spatial cues rely strongly on interaural phase differences. To prevent phase-wrap

Dataset	Format	Collect	Hours	Type	Labels
Sweet-Home (Vacher et al., 2014)	Multi	Recorded	47.3	Speech	Text
Voice-Home (Bertin et al., 2016)	Multi	Recorded	2.5	Speech	Text, Geomtrtric
YT-ALL (Morgado et al., 2018)	FOA	Crawled	113	Audio	Video, Text
REC-STEEET (Morgado et al., 2018)	FOA	Recorded	3.5	Audio	Video
FAIR-Play (Gao and Grauman, 2019)	Binaural	Recorded	5.2	Audio	Video
SECL-UMons (Brousmiche et al., 2020)	Multi	Recorded	5	Audio	Text, Geometric
YT-360 (Morgado et al., 2020)	FOA	Crawled	246	Audio	Video
EasyCom (Donley et al., 2021)	Binaural	Recorded	5	Speech	Geometric, Text
Binaural(Richard et al., 2021)	Binaural	Recorded	2	Speech	Geometric
SimBinaural (Garg et al., 2023)	Binaural	Simulated	116	Audio	Video, Geometric
YouTube-Binaural (Garg et al., 2023)	Binaural	Crawled	27	Audio	Video
Spatial LibriSpeech (Sarabia et al., 2023)	FOA	Simulated	650	Speech	Text, Geometric
STARSS23 (Shimada et al., 2023)	FOA	Recorded	7.5	Audio	Video, Geometric
YT-Ambigen (Kim et al., 2025)	FOA	Crawled	142	Audio	Video
BEWO-1M (Sun et al., 2024)	Binaural	Simulated	2.8k	Audio	Text/Image, Geometric
MRSDrama (Zhang et al., 2025)	Binaural	Recorded	98	Speech	Text, Video, Geometric

Table 4: Comparison of current spatial audio datasets. FOA means first-order ambisonics, while Multi denotes multi-channel audio.

artefacts, we minimize the wrapped phase distance:

$$\mathcal{L}_{L_2}^{\text{phase}} = \frac{1}{KM} \sum_{k=1}^K \sum_{m=1}^M \left(\text{wrap}(\angle \hat{X}(k, m) - \angle X(k, m)) \right)^2, \quad (9)$$

where $\text{wrap}(\theta) \in [-\pi, \pi]$.

To align perceptual quality with spectral accuracy, we average three complementary losses over a bank of M STFT configurations $\{ \cdot^{(i)} \}_{i=1}^M$ as **Multi-Resolution STFT Loss (MRSTFT)**:

$$\mathcal{L}_{\text{SC}}^{(i)} = \frac{\| |X^{(i)}| - |\hat{X}^{(i)}| \|_F}{\| |X^{(i)}| \|_F},$$

$$\mathcal{L}_{\text{mag}}^{(i)} = \frac{1}{N^{(i)}} \| |X^{(i)}| - |\hat{X}^{(i)}| \|_1,$$

$$\mathcal{L}_{\text{log}}^{(i)} = \frac{1}{N^{(i)}} \| \log(|X^{(i)}| + \varepsilon) - \log(|\hat{X}^{(i)}| + \varepsilon) \|_1, \quad (10)$$

$$\mathcal{L}_{\text{MRSTFT}} = \frac{1}{M} \sum_{i=1}^M (\mathcal{L}_{\text{SC}}^{(i)} + \lambda_{\text{mag}} \mathcal{L}_{\text{mag}}^{(i)} + \lambda_{\text{log}} \mathcal{L}_{\text{log}}^{(i)}). \quad (11)$$

This compound objective balances global spectral convergence with fine-grained magnitude fidelity across multiple time–frequency resolutions.

For **Perceptual Evaluation of Speech Quality (PESQ)**, the ITU-T P.862 standard maps symmetric (d_{sym}) and asymmetric (d_{asym}) perceptual distortions onto a MOS-like scale:

$$\text{PESQ} = 4.5 - 0.1 d_{\text{sym}} - 0.0309 d_{\text{asym}}, \quad (12)$$

yielding scores in $[-0.5, 4.5]$. Higher values denote closer perceptual similarity.

Finally, subjective quality **Mean Opinion Score (MOS)** is obtained by averaging listener ratings over a five-point Likert scale:

$$\text{MOS} = \frac{1}{N} \sum_{i=1}^N s_i, \quad (13)$$

where s_i is the score from the i -th participant. MOS serves as the definitive benchmark to which all objective metrics are ultimately calibrated.

Wave/Amplitude/Phase L_2 losses provide gradient-friendly objectives that capture complementary signal aspects. MRSTFT augments them with multi-resolution spectral consistency. PESQ offers a single-ended perceptual estimate that correlates well with telecommunication speech quality, and MOS delivers the gold-standard human judgment. Together, this metric suite affords a balanced evaluation of both technical accuracy and perceptual realism in mono-to-stereo binaural conversion.

End-to-End Binaural Audio Generation. Evaluation metrics are highly varied for this task. In the case of binaural spatial audio, metrics can be computed based on interaural time difference (ITD), interaural level difference (ILD), and embeddings from a pretrained spatial-audio-understanding model (Zheng et al., 2024) to calculate specific performance indicators (Zhang et al., 2025). For the objective evaluation of IPD and ILD, they first convert the time-domain signal $x(n)$ into the frequency-domain signal $X(t, f)$ using the

short-time Fourier transform (STFT):

$$X_i(t, f) = \sum_{n=0}^{N-1} x_i(n) \cdot w(t - n) \cdot e^{-j2\pi fn}, i \in \{1, 2\} \quad (14)$$

where $w(t - n)$ is a window function, N is the window length, and i indicates the channel of the binaural audio. Next, they calculate the mel-spectrogram, IPD, and ILD based on the frequency-domain signals $X_i(t, f)$. The mel-spectrogram for each channel is calculated as:

$$S_i(t, m) = \log(|X_i(t, f)|^2 \times \text{melW}), \quad (15)$$

where melW is an M -bin mel filter bank. **IPD** is derived from the phase spectrograms of the left and right channels:

$$\text{IPD}(t, f) = \angle \frac{X_2(t, f)}{X_1(t, f)}. \quad (16)$$

Then, **ILD** is extracted from the loudness spectrum of the left and right channels:

$$\text{ILD}(t, f) = 20 \log_{10} \left(\frac{|X_2(t, f)| + \varepsilon}{|X_1(t, f)| + \varepsilon} \right), \varepsilon = 1e^{-10}. \quad (17)$$

They calculate Mean Absolute Error (MAE) metrics based on the IPD and ILD extracted from the ground truth (GT) and the predicted speech. Since the IPD here is in radians and the ILD uses \log_{10} , the resulting values are quite small, especially after averaging the MAE over the time dimension. So, they multiply by 100 to make the results more intuitive.

Additionally, they analyze angular and distance metrics using SPATIAL-AST (Zheng et al., 2024). SPATIAL-AST encodes **angle and distance** embedding for binaural audio. They compute and average the cosine similarity for each 1-second segment based on the GT and predicted audio.

End-to-End FOA Generation. Current methods usually assess spatial localization accuracy by measuring azimuth error, elevation error, distance error, and spatial-angle difference (Heydari et al., 2025). Codec quality is evaluated via STFT and Mel distances between original and reconstructed FOA audio on the test set, using AuraLoss with default settings (Evans et al., 2024a,b). Plausibility of generated clips is quantified by the **Fréchet Audio Distance (FAD)** between generated and reference embeddings, and by **KL divergence** computed with a pretrained ELSA model. The **CLAP**

score, the cosine similarity between spatial text embeddings and corresponding audio embeddings, is also reported. For the parametric model, KL divergence and CLAP are computed using spatial captions from the test set, despite training on non-spatial captions and parameters.

To measure spatial accuracy, they compare ground-truth and estimated **azimuth** θ , **elevation** ϕ , and **distance** d . Intensity vectors I_x, I_y, I_z are obtained by multiplying the omnidirectional channel W with the directional channels X, Y, Z :

$$I_x = W \cdot X, \quad I_y = W \cdot Y, \quad I_z = W \cdot Z \quad (18)$$

$$\theta = \tan^{-1} \frac{I_y}{I_x}, \quad \phi = \tan^{-1} \frac{I_z}{\sqrt{I_x^2 + I_y^2}}, \quad (19)$$

$$d = \sqrt{I_x^2 + I_y^2 + I_z^2} \quad (20)$$

They report the L1 norm of the differences for azimuth, elevation, and distance. For azimuth, they use the circular difference:

$$\text{L1}_\theta = ||(|\theta - \hat{\theta}|, 2\pi - |\theta - \hat{\theta}|)||_1 \quad (21)$$

Spatial-angle error $\Delta_{\text{Spatial-Angle}}$ is defined as (Van Brummelen, 2012):

$$a = \sin^2\left(\frac{\Delta_\phi}{2}\right) + \cos(\phi) \cos(\hat{\phi}) \sin^2\left(\frac{\Delta_\theta}{2}\right) \quad (22)$$

$$\Delta_{\text{Spatial-Angle}} = 2 \arctan 2(\sqrt{a}, \sqrt{1-a}) \quad (23)$$

Here, Δ_ϕ and Δ_θ denote the linear and circular differences for elevation and azimuth, respectively.

E Licenses and Availability

We respect the original licenses of all referenced artifacts and do not redistribute them. This survey does not create new deployable systems or redistribute data. Any consultation of third-party artifacts is limited to research/read-only use and complies with their intended-use statements and access conditions.