

Cinematic Audio Source Separation Using Visual Cues

Anonymous CVPR submission

Abstract

001 *Cinematic Audio Source Separation (CASS) aims to de-*
002 *compose mixed film audio into speech, music, and sound*
003 *effects, enabling applications like dubbing and remaster-*
004 *ing. Existing CASS approaches are audio-only, overlook-*
005 *ing the inherent audio-visual nature of films, where sounds*
006 *often align with visual cues. We present the first frame-*
007 *work for audio-visual CASS (AV-CASS), leveraging visual*
008 *context to enhance separation quality. Our method formu-*
009 *lates CASS as a conditional generative modeling problem*
010 *using conditional flow matching, enabling multimodal au-*
011 *dio source separation. To address the lack of cinematic*
012 *datasets with isolated sound tracks, we introduce a train-*
013 *ing data synthesis pipeline that pairs in-the-wild audio and*
014 *video streams (e.g., facial videos for speech, scene videos*
015 *for effects) and design a dedicated visual encoder for this*
016 *dual-stream setup. Trained entirely on synthetic data, our*
017 *model generalizes effectively to real-world cinematic con-*
018 *tent and achieves strong performance on synthetic, real-*
019 *world, and audio-only CASS benchmarks.*

020 1. Introduction

021 Cinematic audio is composed of layered sound elements
022 such as speech, music, and sound effects, which collectively
023 enrich storytelling and immersion. The goal of Cinematic
024 Audio Source Separation (CASS) is to separate a mixed
025 movie audio into these three distinct tracks (Fig. 1). CASS
026 enables a wide range of applications, including multilingual
027 dubbing, film remastering, audio editing, and accessibility
028 enhancement. As video streaming platforms continue to
029 grow, the need for CASS becomes increasingly important
030 to enable automatic tools that precisely control individual
031 sound components.

032 While related audio separation tasks, such as speech sep-
033 aration [8, 24, 62] and music demixing [14, 15, 44, 52],
034 have seen significant progress, CASS remains underex-
035 plored. The introduction of the Divide and Remaster (DnR)
036 dataset [48] formalized CASS as a three-way separation
037 problem and initiated new research [57, 59, 60]. How-

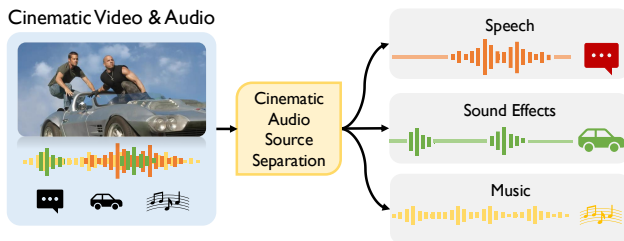


Figure 1. **Illustration of the Cinematic Audio Source Separation (CASS) task.** The audio stream from a movie is separated into distinct tracks: speech, sound effects, and music.

ever, existing methods are purely audio-based and overlook
a defining characteristic of film: its inherently audio-visual
nature.

In cinematic content, sound is often tightly coupled
with visual events. Speech typically co-occurs with lip
movements, and sound effects frequently align with ob-
ject interactions or visual actions. Prior work in audio-
visual learning has shown that using visual context, such
as facial motion, body movement, and scene composition,
significantly improves audio separation and enhance-
ment [1, 7, 16, 20, 21, 55]. Yet, to our knowledge, visual
information has not been utilized in the context of CASS.

The primary motivation of this paper is to achieve cine-
matic audio source separation using visual cues. Accord-
ingly, we first formulate CASS as a conditional genera-
tive modeling problem. We adopt conditional flow match-
ing [39] as our generative framework, which has demon-
strated strong performance in both image [17, 27] and au-
dio [22, 31, 35, 58] generation tasks. Our model gener-
ates clean speech, sound effects, and music tracks from the
mixture, conditioned on both audio and visual inputs.

Second, despite the strong potential of visual cues, ap-
plying them to CASS remains challenging. A key obstacle
is the lack of publicly available audio-visual datasets with
clean source tracks, which are difficult – if not impossible –
to obtain from real films. This raises an important question:
Can we leverage individually available in-the-wild audio-
visual data to train an effective audio-visual CASS model?
To address this, we propose a novel training data synthesis

067	strategy. In the absence of paired datasets such as films and	
068	their clean sound source tracks, we construct a pipeline that	
069	synthetically pairs individually available audio and video	
070	sources, <i>e.g.</i> , using facial videos for speech and scene video	
071	clips for sound effects. This results in a dual-stream video	
072	setup that provides controllable, source-specific visual su-	
073	per- vision, for which we also design a visual feature extrac-	
074	tor to leverage this setup. Importantly, although our model	
075	is trained using this synthetic dual-video configuration, we	
076	demonstrate that it generalizes effectively to real-world cin-	
077	ematic content without architectural changes, highlighting	
078	the practicality and robustness of our training approach.	
079	In summary, our contributions are:	
080	• We introduce the first framework for audio-visual cine-	
081	matic audio source separation (AV-CASS), incorporating	
082	visual cues to separate speech, sound effects, and music	
083	in film audio.	
084	• We formulate CASS as a generative task using condi-	
085	tional flow matching, enabling flexible and principled	
086	modeling of multimodal audio decomposition.	
087	• We propose a training data synthesis strategy, along with	
088	a dedicated visual encoder, to enable training without the	
089	need for original source-separated film data.	
090	• We achieve strong performance on synthetic data, real-	
091	world cinematic examples, and standard audio-only	
092	CASS benchmarks.	
093	2. Related Work	
094	2.1. Cinematic Audio Source Separation	
095	Cinematic Audio Source Separation (CASS) was formal-	
096	ized by [48, 49], introducing the Divide and Remaster	
097	(DnR) dataset for the separation of speech, sound effects,	
098	and music in film audio. BandIt [59] applied Band-split	
099	RNNs [42] to this task, achieving improved performance	
100	over earlier methods. More recently, DnRv3 [60] ex-	
101	panded the dataset with multilingual content and introduced	
102	mixing strategies aligned with industrial audio production	
103	pipelines. DnR-nonverbal [23] further expands the dataset	
104	by incorporating nonverbal sounds, such as laughter and	
105	screams. While these efforts have advanced audio-based	
106	CASS, they remain limited to audio-only learning. To our	
107	knowledge, no prior work has explored audio-visual ap-	
108	proaches for CASS, despite the inherently multimodal na-	
109	ture of cinematic content. This gap largely stems from	
110	the difficulty of acquiring datasets with both isolated audio	
111	tracks and temporally aligned video. In contrast, we pro-	
112	pose the first audio-visual CASS framework and introduce	
113	a training data synthesis pipeline that constructs audio-visual	
114	training samples from in-the-wild video sources, without re-	
115	quiring ground-truth source audio with film datasets. Our	
116	approach enables effective audio-visual learning and gener-	
117	alizes well to real-world cinematic content.	
	2.2. Audio-Visual Source Separation	118
	Incorporating visual information has proven highly effec-	119
	tive in sound source separation tasks. Prior work has shown	120
	that visual cues such as lip movements align strongly with	121
	spoken content [10, 43], while facial features provide cross-	122
	modal biometric information that enhances speech separa-	123
	tion [1, 3, 12, 16, 46, 47]. Beyond speech, visual informa-	124
	tion from instrument motion or class-level appearance cues	125
	has also been leveraged to improve music separation [5, 18,	126
	19, 66]. More generally, recent works have demonstrated	127
	that visual context benefits generic sound separation across	128
	a wide range of categories [7, 9, 21, 28, 29, 54–56, 63]. De-	129
	spite these advances, applying audio-visual learning to the	130
	complex, multi-source nature of cinematic audio remains	131
	largely unexplored. In this work, we introduce audio-visual	132
	learning to Cinematic Audio Source Separation (CASS) for	133
	the first time. Unlike prior AVSS approaches that condition	134
	on a single visual cue (<i>e.g.</i> , facial motion [10, 43] or ob-	135
	ject appearance [21, 29, 55]), our framework integrates two	136
	complementary visual streams, facial and scene context, de-	137
	derived from the same video, enabling individual-source train-	138
	ing while remaining applicable to real-world single-video	139
	inputs.	140
	2.3. Flow Matching Models	141
	Flow matching [17, 39] has recently gained attention as	142
	an efficient alternative to diffusion models, offering faster	143
	inference by following shorter and more direct generation	144
	trajectories. Recent works have applied flow matching to	145
	audio synthesis, separation, and enhancement [22, 29, 31,	146
	35, 45, 58, 64], demonstrating its potential to produce high-	147
	quality, natural-sounding outputs. Previous non-generative,	148
	masking-based separation models often introduce artifacts	149
	(<i>e.g.</i> , spectral holes) as noted in [64], rendering the out-	150
	put unsuitable for downstream tasks like audio editing. We	151
	therefore adopt a generative flow-matching model for the	152
	CASS task, which effectively resolves this issue. To our	153
	knowledge, this is the first visually conditioned generative	154
	flow-matching approach to CASS, designed to yield percep-	155
	tually natural and artifact-free separated audio suitable for	156
	cinematic production.	157
	3. Methodology	158
	We propose a framework for audio-visual cinematic audio	159
	source separation (AV-CASS). As shown in Fig. 2, it con-	160
	sists of a Vision Extractor that generates a fused repre-	161
	sentation c^V from input videos to condition the separation	162
	model. Next, a flow-based generative model performs	163
	source separation by mapping Gaussian noise to three dis-	164
	tinct audio components: speech, sound effects, and mu-	165
	sic. Finally, to enable training without source-separated film	166
	data, we introduce a training data synthesis pipeline that	167

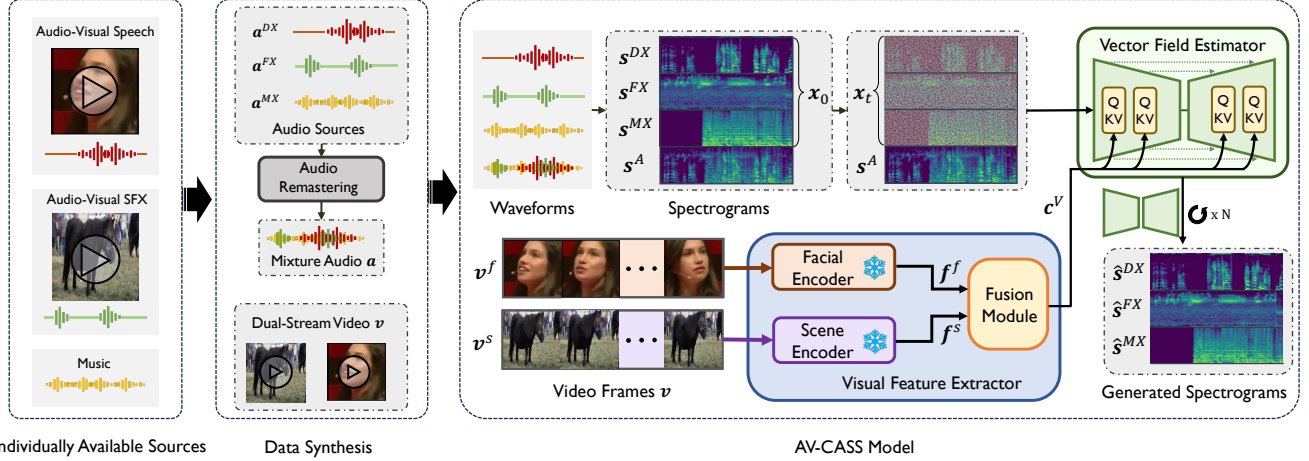


Figure 2. **Architecture of AV-CASS.** The fusion module integrates visual features from the facial and scene encoders into c^V , which serves as a conditioning input along with a mixture audio s^A for the vector field estimator u_θ .

168 leverages individually available in-the-wild audio and video
169 sources.

170 3.1. Problem Setup

171 Given a cinematic audio mixture \mathbf{a} composed of three additive
172 sources, speech \mathbf{a}^{DX} , sound effects \mathbf{a}^{FX} , and music
173 \mathbf{a}^{MX} , the goal of CASS is to recover each source track
174 from the mixture: $\mathbf{a} = \mathbf{a}^{DX} + \mathbf{a}^{FX} + \mathbf{a}^{MX}$. We formulate
175 this task as a conditional generation problem, where the model
176 is conditioned not only on the mixture audio but also on
177 visual input \mathbf{v} from the corresponding video. The video
178 input contains two components: facial frames \mathbf{v}^f and scene
179 frames \mathbf{v}^s , which provide complementary cues for speech
180 and sound effects, respectively.

181 3.2. Input Data for Training

182 To enable training without the need for source-separated
183 cinematic data, we synthesize a training data by combining
184 individual video and audio segments from diverse sources.
185 For each training sample, we: (1) select speech clips from
186 an audio-visual speech dataset to create the speech track
187 \mathbf{a}^{DX} and the corresponding video as face stream \mathbf{v}^f ; (2)
188 select sound effect clips from an audio-visual sound dataset
189 to create the sound effects track \mathbf{a}^{FX} and the corresponding
190 video as scene stream \mathbf{v}^s ; and (3) select background
191 music clips from a music dataset to create the music track
192 \mathbf{a}^{MX} . These source tracks are mixed to create the input
193 mixture \mathbf{a} , while the video consists of two parallel streams,
194 $\mathbf{v} = \{\mathbf{v}^f, \mathbf{v}^s\}$: a face stream and a scene stream that represent
195 different semantic aspects of the sound. This strategy
196 offers a diverse, controllable, and scalable source of training
197 pairs with ground-truth supervision for all components.
198 Details of the training data synthesis process are provided in
199 Sec. 4. All audio waveforms are converted into their corresponding
200 spectrograms where s^A denotes the mixture spec-

201 trogram, while s^{DX} , s^{FX} , and s^{MX} represent the individual
202 source spectrograms. The video streams \mathbf{v}^f and \mathbf{v}^s
203 are represented as sequences of frames.

204 3.3. AV-CASS Model

205 Our audio-visual cinematic source separation model (AV-
206 CASS) consists of two main components: (1) visual feature
207 extraction and fusion, and (2) conditional flow matching for
208 source generation.

209 3.3.1. Visual Feature Encoding and Fusion

210 We extract visual features from two video streams: facial
211 frames \mathbf{v}^f and scene frames \mathbf{v}^s , using separate encoders
212 suited to their semantic roles. The facial encoder, adopted
213 from AVDiffuSS [37], is designed for lip-synced speech
214 videos. The scene encoder, based on the CAVP model [41],
215 captures temporally and semantically aligned features of
216 sounding objects and events. Both encoders are frozen
217 during training. The outputs are feature sequences
218 $\mathbf{f}^f \in \mathbb{R}^{T_f \times D_f}$ and $\mathbf{f}^s \in \mathbb{R}^{T_s \times D_s}$, where T_f and T_s
219 denote the number of frames, and D_f , D_s are feature
220 dimensions. To fuse the representations, we project each
221 stream into a shared feature space of dimension C using
222 separate MLPs:

$$222 \mathbf{f}^f \rightarrow \mathbb{R}^{T_f \times C}, \quad \mathbf{f}^s \rightarrow \mathbb{R}^{T_s \times C}. \quad (1)$$

223 The resulting features are concatenated along the temporal
224 axis and passed through a final fusion MLP to obtain the
225 visual condition vector:

$$226 \mathbf{c}^V \in \mathbb{R}^{(T_f + T_s) \times C'}. \quad (2)$$

227 This visual condition is used to guide the audio generation
228 process via cross-attention in the U-Net backbone.

229 3.3.2. Flow Matching for Multisource Separation

230 We adopt conditional flow matching [39] to model the
231 conditional joint distribution of source spectrograms $p_1(\mathbf{x}) :=$

232 $p(\mathbf{s}^{DX}, \mathbf{s}^{FX}, \mathbf{s}^{MX})$ given the mixture spectrograms \mathbf{s}^A
 233 and visual conditioning vector \mathbf{c}^V . Conditional flow match-
 234 ing defines a conditional mapping between the Gaussian
 235 noise distribution $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and the target joint distri-
 236 bution of source spectrograms $\mathbf{x}_1 \sim p_1(\mathbf{x})$. This mapping
 237 defines a time-varying probability density governed by an
 238 ordinary differential equation:

$$239 \quad d\mathbf{x}_t = \mathbf{u}_\theta(\mathbf{x}_t, t|\mathbf{c})dt, \quad (3)$$

240 where \mathbf{u}_θ is a vector field estimator representing the gradi-
 241 ent of the probability density w.r.t. time t at point \mathbf{x}_t ; \mathbf{x}_t
 242 is a point in the probability density space at time t ; and \mathbf{c}
 243 is conditioning variable includes mixture audio \mathbf{s}^A and its
 244 visual context \mathbf{c}^V .

245 To construct this mapping, we train \mathbf{u}_θ to approximate
 246 a reference vector field \mathbf{u}_t , which constructs a probabilis-
 247 tic path between the noise distribution $p_0(\mathbf{x})$ and the target
 248 distribution $p_1(\mathbf{x})$, conditioned on \mathbf{c} . The vector field for
 249 a noise-data pair $(\mathbf{x}_0, \mathbf{x}_1)$ is defined as $\mathbf{u}_t = \mathbf{x}_1 - \mathbf{x}_0$. In
 250 [39], given a target distribution sample \mathbf{x}_1 , the data point \mathbf{x}_t
 251 for timestep t on the path is defined as:

$$252 \quad \mathbf{x}_t = (1-t)\mathbf{x}_0 + t\mathbf{x}_1, \quad (4)$$

253 where $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $t \in [0, 1]$. In practice, \mathbf{u}_θ is
 254 trained to approximate \mathbf{u}_t by minimizing the following loss:

$$255 \quad \mathcal{L} = \mathbb{E}_{t, \pi_1(\mathbf{x}_1), \pi_0(\mathbf{x}_0)} \|\mathbf{u}_\theta(\mathbf{x}_t, t|\mathbf{c}) - (\mathbf{x}_1 - \mathbf{x}_0)\|_2^2. \quad (5)$$

256 Inspired by [17], we sample the timestep t from a logit-
 257 normal distribution, as this has been shown to enhance gen-
 258 eration quality by placing more emphasis on the interme-
 259 diate timesteps during training. In practice, we first sample
 260 a random variable z from a standard Gaussian distribution
 261 $z \sim \mathcal{N}(0, 1)$, then map it with a logistic function as follows:

$$262 \quad t = \frac{1}{1 + e^{-z}}, z \sim \mathcal{N}(0, 1). \quad (6)$$

263 For the vector field estimator \mathbf{u}_θ , we adopt a CNN-based
 264 U-Net architecture and three sources are concatenated along
 265 the channel dimension to form the input, as commonly done
 266 in diffusion-based image generation models [26, 51].

267 3.3.3. Inference

268 At inference, the trained conditional flow matching model
 269 generates separated source spectrograms from a mixture
 270 audio and its associated visual context. As outlined ear-
 271 lier, our training setup uses two video streams for a single
 272 audio mixture. Although this differs from the real-world
 273 one-video-one-audio setting, our approach remains effec-
 274 tive for real-world cinematic content. Since different visual
 275 regions (e.g., faces, environments, background elements)
 276 correspond to distinct sound sources, as shown in Fig. 3, we

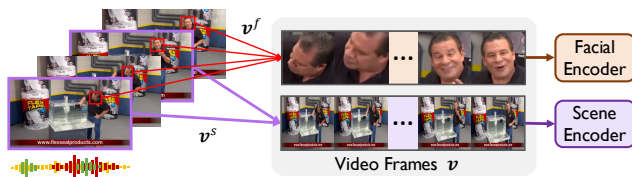


Figure 3. **Extraction of dual-stream visual inputs from a real-world cinematic video during inference.** Since no architectural changes are required, the AV-CASS model can be used with real-world cinematic videos for inference.

277 extract the facial regions as one stream for speech cues (pro-
 278 cessed by the Facial Encoder) and the full scene frames as
 279 another for sound effects (processed by the Scene Encoder).
 280 This design enables our model to process real-world sam-
 281 ples without architectural changes. The model takes as input
 282 the mixture spectrogram \mathbf{s}^A and fused visual condition
 283 \mathbf{c}^V and outputs the separated components: speech, sound
 284 effects, and music.

285 To formulate inference under the conditional flow match-
 286 ing framework, we follow the sampling procedure defined
 287 in [39]. Specifically, we initialize the sample \mathbf{x}_0 as Gaus-
 288 sian noise, $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and iteratively update it along the
 289 vector field predicted by vector field estimator \mathbf{u}_θ , which is
 290 conditioned on both the mixture audio and visual context.
 291 We apply the forward Euler method to integrate the flow:

$$292 \quad \mathbf{x}_{t+\eta} = \mathbf{x}_t + \eta \mathbf{u}_\theta(\mathbf{x}_t, t | \mathbf{s}^A, \mathbf{c}^V), \quad (7)$$

293 where $\eta = 1/N$ is the step size and N is the total number of
 294 sampling steps. The time variable t progresses from 0 to 1,
 295 and at each step, the vector field guides the sample closer to
 296 the data distribution of clean source spectrograms. After N
 297 steps, the final output \mathbf{x}_1 represents a concatenated output
 298 of three spectrograms, each corresponding to a separated
 299 source: $\{\hat{\mathbf{s}}^{DX}, \hat{\mathbf{s}}^{FX}, \hat{\mathbf{s}}^{MX}\}$. These spectrograms are con-
 300 verted back to the waveform domain using inverse STFT.

301 4. Training Data Construction Pipeline

302 Training an audio-visual CASS model requires synchron-
 303 ized film video and isolated source audio tracks, which
 304 are rarely available for real-world cinematic content. To
 305 address this, we design a pipeline that synthetically pairs
 306 independently sourced audio and video from existing
 307 datasets to create realistic training pairs resembling cine-
 308 matic data. This pipeline produces multimodal training data
 309 that mimics cinematic audio-visual patterns while preserv-
 310 ing ground-truth alignment across all audio stems and video
 311 streams.

312 4.1. Individually Available Sources

313 We leverage two large-scale audio-visual datasets and one
 314 music dataset, each corresponding to one of the three target
 315 stems in CASS:

316 **Speech (DX):** We use LRS3 [2], a lip-synchronized video
 317 dataset with high-quality speech that may reflect cinematic
 318 dialogue through natural prosody and visual expressiveness.
 319 **Sound Effects (FX):** For ambient, object or action-driven
 320 sounds, we use VGGSound [6], an audio-visual dataset
 321 containing everyday events and objects. Unlike prior
 322 works [48, 60] using audio-only FSD50K, VGGSound pro-
 323 vides aligned video for learning visually grounded effects.
 324 **Music (MX):** Since background music is typically not vi-
 325 sually grounded, we follow standard practice and use the
 326 FMA dataset [13], which contains a wide variety of high-
 327 quality music.

328 4.2. Audio Preprocessing and Stream Synthesis

329 To ensure each sample contains a single, uncontaminated
 330 source track, we filter VGGSound and FMA with the
 331 SMAD model [30], removing all segments containing
 332 speech and music. After filtering, we obtain a total of 152K
 333 DX clips, ~ 62 K FX clips, and ~ 49 K MX samples.

334 Following the protocol in DnRv3 [60], we synthesize
 335 a cinematic audio by using the DX, FX, and MX tracks.
 336 For each track, we randomly sample short clips, concate-
 337 nate them with overlapping transitions, and apply loudness
 338 normalization to meet cinematic mastering standards. The
 339 resulting tracks are then mixed by addition:

$$340 \mathbf{a}^A = \mathbf{a}^{DX} + \mathbf{a}^{FX} + \mathbf{a}^{MX}. \quad (8)$$

341 All audio is converted to mono and resampled at 16kHz.
 342 Statistics of the resulting training data are in *Appendix C*.

343 4.3. Visual Stream Synthesis

344 We extract the video clips corresponding to the DX and FX
 345 streams. Using the timestamps of each selected audio seg-
 346 ment, we retrieve the facial video aligned with DX (from
 347 LRS3) and the scene video aligned with FX (from VG-
 348 GSound). The MX stream has no associated visual input.

349 Each training sample contains a mixed audio stream with
 350 ground-truth sources and two video streams, *i.e.*, facial and
 351 scene, reflecting how speech and sound effects are visually
 352 grounded in real films. Although the dataset is synthetically
 353 constructed, it enables precise supervision for multimodal
 354 learning and aligns well with cinematic audio-visual con-
 355 ventions. As shown in the Sec. 5, our model trained on this
 356 data generalizes effectively to real-world movies, validating
 357 our training data pipeline.

358 5. Experiments

359 5.1. Experimental Setup

360 **Baselines.** We compare our method with existing CASS
 361 models, including MRX [48] and BandIt [59], as well
 362 as musical instrument separation models such as Hybrid

Demucs [14], HT Demucs [52], and MSDM [44]. Be-
 yond audio-only baselines, we also include the audio-visual
 sound separation model DAVIS-Flow [29] which is the cur-
 rent state-of-the-art model for audio-visual sound separa-
 tion, to assess the contribution of visual conditioning and
 to highlight the fundamental differences between CASS
 and generic audio-visual separation tasks. All models are
 trained from scratch on the same dataset, using their orig-
 inal training configurations and appropriately modified to
 operate under the CASS setting for fair comparison.

Metrics. We use Fréchet Audio Distance (FAD) [32] and
 Kullback-Leibler divergence (KL) from AudioLDM [40]
 to measure distributional similarity between generated and
 real audio. We also report Perceptual Evaluation of Speech
 Quality (PESQ) [50] for speech and Scale-Invariant Signal-
 to-Distortion Ratio improvement (SI-SDRi) [36] in dB,
 following prior works [44, 48]. For FAD, KL, and SI-
 SDRi, we report averages across all three sources. In ad-
 dition, we introduce a new metric, *Wrong Placement Ratio*
 (*WPR*), to estimate the proportion of residual or misplaced
 components from other stems. *WPR* is computed using
 PANNs [34], a pretrained sound event detection model, and
 reflects stem-level separation quality without ground-truth
 references; lower values indicate better isolation. Details
 on metric calculations is in the *Appendix E*.

Training and Implementation Details. Before training
 the full audio-visual model, we apply audio-denoiser warm-
 up to stabilize optimization. Specifically, we first train the
 model using only the audio component of the synthetic
 audio-visual dataset described in Sec. 4. In practice, it im-
 proves training stability, accelerates convergence, and pre-
 vents the model from prematurely overfitting to visual cues
 in the early stages of training.

After warm-up, we train the full audio-visual model on
 the same synthetic dataset. Visual cues from both the fa-
 cial and scene streams are introduced gradually using zero-
 initialized convolution layers, following ControlNet [65]
 strategy, while the video encoders remain frozen. This de-
 sign preserves the stabilized audio representation while al-
 lowing the model to gradually integrate visual information.

We use the Adam optimizer [33] with $\beta_1 = 0.9$, $\beta_2 =$
 0.999 , a fixed learning rate of 10^{-4} , and no weight decay.
 Full audio-visual training runs for 600k steps with a batch
 size of 8 across four RTX 4090 GPUs. We use 128 sampling
 steps during evaluation. Additional implementation details
 and pseudo code are provided in the *Appendix A*.

5.2. Main Results

5.2.1. Evaluation on real-world samples

Subjective evaluation. Generalization to real-world movie
 samples is critical for CASS models. To evaluate our model,
 we randomly select clips from the Condensed Movies
 dataset [4], manually verifying that each contains all three

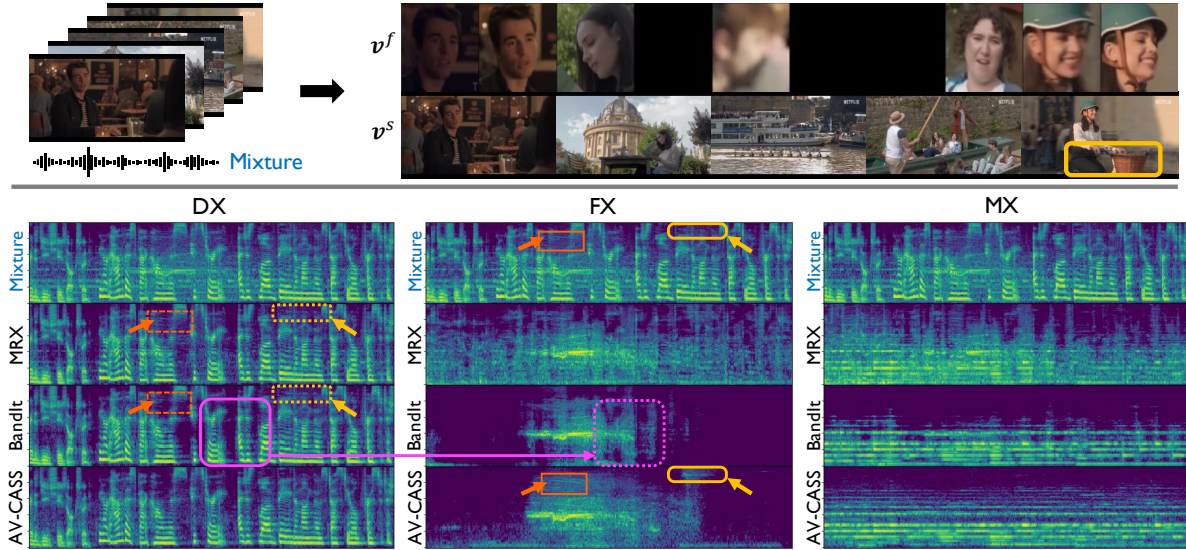


Figure 4. **Comparison of MRX, BandIt, and AV-CASS on a real-world movie sample.** Input video frames v^f and v^s are shown at the top, with the input audio spectrogram s^A placed for each stem. Yellow boxes highlight the bicycle bell, red boxes indicate cheering, and dotted boxes show elements misplaced in non-target stems. The dotted pink box in BandIt’s FX shows speech artifacts. Best viewed when zoomed in. This sample can also be viewed in the supplementary video.

Method	MRX [48]	BandIt [59]	AV-CASS (Ours)
MOS (\uparrow)	2.55 ± 0.10	3.78 ± 0.10	4.13 ± 0.09

Table 1. **MOS results of CASS models on real-world samples.** The scores are computed based on 95% confidence intervals.

Method	DX (\downarrow)	FX (\downarrow)	MX (\downarrow)
MRX [48]	2.63	33.72	6.01
BandIt [59]	<u>0.60</u>	22.41	<u>0.35</u>
DAVIS-Flow [29]	5.88	14.58	35.94
AV-CASS (Ours)	0.46	<u>19.81</u>	0.32

Table 2. **Wrong Placement Ratio (WPR [%]) on real-world samples.** Lower is better.

415 target tracks. As this process requires extensive effort, we
 416 collected 30 samples. Since ground-truth separation data is
 417 unavailable, we conduct a subjective evaluation using mean
 418 opinion scores (MOS) through a human study, detailed in
 419 the *Appendix D*. Some of the movie samples used for MOS
 420 are available in the supplementary video. We strongly en-
 421 courage readers to view real-world samples there, along
 422 with comparisons to existing CASS methods. The numeri-
 423 cal results are provided in Tab. 1. Participants compared
 424 and evaluated the separated outputs for each target stem us-
 425 ing two criteria: clarity of separation and completeness of
 426 target reconstruction. These criteria allows raters to evalu-
 427 ate both how well the model suppresses non-target sources
 428 and how fully it preserves the target source without loss of
 429 content or quality. Each separation was rated on a 5-point
 430 Likert scale, with 1 meaning ‘‘Poor’’ and 5 meaning ‘‘Excel-
 431 lent’’. As shown in Tab. 1, our model receives higher scores,
 432 reflecting user preference and demonstrating the naturalness
 433 and sound quality of the outputs. Overall, this result indi-
 434 cates generalization to real-world movie clips.

435 **Objective evaluation.** To assess AV-CASS on real-world
 436 movies, we report quantitative WPR across DX, FX, and
 437 MX stems. As shown in Tab. 2, our model achieves the
 438 lowest WPR for both DX and MX, and delivers competitive
 439 performance on FX, indicating strong separation fidelity
 440 for dialogue and music. While DAVIS-Flow [29] attains

441 a lower WPR on FX, this is expected since it is specifically
 442 designed for generic object-centric sound separation; how-
 443 ever, it performs poorly on the other tracks. Overall, the
 444 results demonstrate that AV-CASS achieves robust cross-track
 445 isolation in complex, in-the-wild cinematic audio, support-
 446 ing the effectiveness of our formulation and synthetic train-
 447 ing pipeline.

448 **Qualitative result.** We also provide a qualitative visual
 449 analysis (Fig. 4) illustrating residual cross-track compo-
 450 nents in existing methods, whereas our approach yields
 451 cleaner separated tracks. Yellow boxes in v^s and the spec-
 452 trograms highlight the bicycle bell, and red boxes indicate
 453 cheering. Dotted boxes show elements incorrectly placed
 454 in non-target stems. Other methods retain these effects in
 455 the speech track, while ours correctly assigns them to FX.
 456 The dotted pink box in BandIt’s FX shows speech artifacts,
 457 which AV-CASS avoids. This sample and more are shown
 458 in supplementary videos.

5.2.2. Evaluation on AVDnR

459 While real-world movie audio offers valuable qualitative
 460 insights, it does not provide clean ground-truth stems for
 461 quantitative evaluation. To enable a more controlled and
 462 comprehensive assessment, we construct a fully supervised
 463

Method	A-V	FAD (↓)	KL (↓)	SI-SDRi (↑)	PESQ (↑)	WPR (↓)
<i>Predictive Model</i>						
Hybrid Demucs [14]	✗	2.05	1.03	13.57	2.16	5.24
HT Demucs [52]	✗	2.08	1.06	13.41	2.06	9.23
MRX [48]	✗	3.47	1.67	10.60	1.89	14.91
BandIt [59]	✗	2.15	1.14	14.40	2.15	4.65
<i>Generative Model</i>						
MSDM [44]	✗	2.90	2.90	11.63	2.12	5.65
DAVIS-Flow [29]	✓	5.94	1.64	9.25	1.96	12.14
AV-CASS (Ours)	✓	0.84	0.93	12.32	2.26	1.84
Ours (Audio-only)	✗	<u>1.63</u>	1.15	12.23	2.08	<u>2.01</u>

Table 3. Results on AVDnR dataset (objective scores). All models are trained on our training data. A-V indicates audio-visual.

Method	H-Demucs	HT Demucs	MRX	BandIt	MSDM	Ours
MOS (↑)	<u>3.14 ± 0.15</u>	3.01 ± 0.14	1.90 ± 0.13	3.12 ± 0.14	2.79 ± 0.14	3.90 ± 0.13

Table 4. MOS results on AVDnR dataset.

audio-visual test set, AVDnR, using the same data synthesis pipeline described in Sec. 4. We strictly partition all source clips into disjoint training and testing splits to avoid overlap. The final AVDnR benchmark contains 1K audio-visual samples, each 60 seconds long, providing a reliable testbed with complete ground-truth sources.

Objective evaluation. We train all methods on our training set and report results on AVDnR in Tab. 3. Methods are grouped into *predictive* and *generative* categories. As commonly observed [61, 64], predictive models which optimized with reconstruction or SNR-based losses tend to achieve higher SI-SDRi, reflecting their stronger alignment with SNR-based metrics. However, they often produce over-smoothed estimates that limit perceptual fidelity [38, 53]. In contrast, generative models focus on producing realistic samples and therefore excel on perceptual metrics. Within this context, AV-CASS achieves the best FAD, KL, PESQ, and WPR scores across all methods, indicating superior perceptual quality and cleaner cross-track separation.

Notably, AV-CASS outperforms the audio-visual baseline DAVIS-Flow [29], even though both methods use visual information. This suggests that performance in AV-CASS depends not only on using vision, but on the way of providing visual cues. By incorporating facial and scene streams in a dual-stream setup and adopting a multi-source formulation, AV-CASS receives source-specific visual context that is not available in the single-target design of DAVIS-Flow, leading to more consistent and reliable track disentanglement. Detailed per-track metrics are in the Appendix F.

Subjective evaluation. Since objective scores do not always capture perceived audio quality, especially for generative models that may introduce realistic details beyond the reference, we additionally conduct a user study on AVDnR. As shown in Tab. 4, AV-CASS achieves the highest MOS among all methods. This perceptual preference aligns with our strong FAD, KL, and PESQ performance (Tab. 3), indicating that AV-CASS produces outputs that listeners consis-

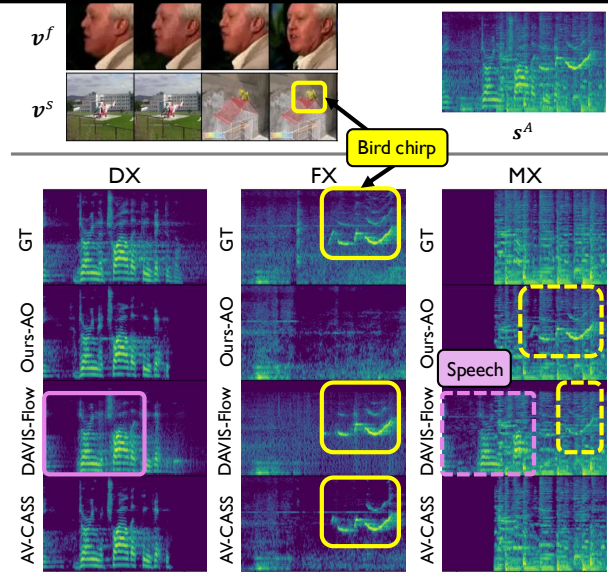


Figure 5. Comparison of our audio-only model (Ours-AO), DAVIS-Flow [29], and our audio-visual model (AV-CASS) on a clip from the AVDnR testset. The input video frames and the GT audio spectrograms are shown at the top. Yellow boxes highlight the bird chirping present in v^s and the FX tracks. Dotted boxes indicate misplaced segments. Better viewed when zoomed in.

tently find clearer and more natural. Together, these results confirm that the perceptual advantages of our model extend beyond objective measures and translate directly into improved user experience.

Qualitative results. Fig. 5 compares separated spectrograms from our audio-only model (Ours-AO), DAVIS-Flow [29], and our audio-visual model (AV-CASS). As the yellow boxes indicate, both audio-visual models correctly separate bird chirping sound into FX, guided by the bird visible in v^s . In contrast, Ours-AO incorrectly places the chirping sound in the MX track. This shows that visual cues can guide correct source separation, resolving ambiguity about where the sound came from. While DAVIS-Flow separates FX correctly, it fails to isolate MX clearly because its design requires a corresponding visual input for every target source. This comparison highlights that AV-CASS effectively exploits visual information. Taken together with the results on real-world samples, the gap between AV-CASS and DAVIS-Flow further confirms that the CASS task differs from generic audio-visual sound separation and requires specialized architectures tailored to its multi-source separation demands, such as ours. Additional examples are shown in Appendix H.

5.2.3. Extendability towards audio-only separation

As discussed earlier, the CASS task is traditionally defined in the audio-only domain, with standard benchmarks such as DnRv2 [49] and DnRv3 [60]. Although our model is designed for audio-visual separation, it can also operate in

	FAD (\downarrow)	KL (\downarrow)	SI-SDRi (\uparrow)	PESQ (\uparrow)	WPR (\downarrow)	
DnRv2 [48]	Hybrid Demucs [14]	3.66	1.47	9.19	<u>2.03</u>	4.77
	HT Demucs [52]	3.72	<u>1.37</u>	8.59	1.95	9.14
	MRX [48]	5.01	1.77	7.48	1.78	16.56
	BandIt [59]	<u>2.75</u>	1.55	7.68	1.97	<u>3.67</u>
	MSDM [44]	6.25	1.54	<u>9.09</u>	2.07	5.33
	Ours (Audio-only)	1.95	1.33	8.10	1.93	2.13
DnRv3 [60]	Hybrid Demucs [14]	<u>3.12</u>	1.68	10.62	1.89	3.78
	HT Demucs [52]	3.17	1.59	<u>9.92</u>	1.83	8.49
	MRX [48]	5.02	2.26	8.94	1.65	16.64
	BandIt [59]	4.79	1.98	9.14	<u>1.86</u>	<u>3.53</u>
	MSDM [44]	5.51	1.75	9.19	1.65	4.49
	Ours (Audio-only)	2.62	<u>1.66</u>	9.36	<u>1.86</u>	1.91

Table 5. **Audio-only CASS results.** Metrics are averaged across three sources, except PESQ, which is evaluated only on the speech source. All models are trained on our dataset.

an audio-only setting, *e.g.*, when video frames are unavailable, by removing the visual encoder and cross-attention blocks in the U-Net. To evaluate this configuration, we train the audio-only variant on our dataset (as in previous experiments) and compare it with other methods on standard audio-only benchmarks and AVDnR. As shown in Tab. 3 and Tab. 5, our method is competitive with models specialized for audio-only CASS while showing clear superiority in FAD, indicating more natural and realistic outputs. Most importantly, our method achieves better WPR performance, indicating cleaner separation with less contamination. These results highlight the flexibility of our approach to support single-modality setups, though its primary scope remains solving CASS from an audio-visual perspective.

5.3. Ablation Study

Analysis on visual streams. We study the impact of visual stream by ablating facial and scene video streams. As shown in Tab. 6, adding either stream improves performance over the audio-only baseline, while using both yields the best results across all metrics, highlighting their complementary benefits and justifying our design choice for AV-CASS. We further analyze how each type of visual input affects the misplacement rates in Tab. 7. DX and FX columns clearly show that the corresponding visual inputs minimize misplaced segments: the lowest DX WPR is achieved with the facial stream, and the lowest FX WPR with the scene stream. By utilizing both visual cues, our model achieves the most balanced WPR performance across all stems and thus improves the overall perceptual performance. This fine-grained misplacement analysis provides additional evidence that visual cues not only improve signal quality metrics, but also substantially reduce semantic cross-contamination between tracks, a crucial property for practical cinematic audio applications. Taken together with the results in Tab. 6, these findings confirm that combining both visual streams offers the most consistent and reliable separation, yielding strong performance across metrics while minimizing residual cross-track contamination.

Method	v^f	v^s	FAD (\downarrow)	KL (\downarrow)	SI-SDRi (\uparrow)	PESQ (\uparrow)
Audio-only	\times	\times	1.63	1.15	12.23	2.08
+ Facial stream	\checkmark	\times	0.91	1.00	12.13	2.21
+ Scene stream	\times	\checkmark	0.87	1.00	12.27	2.24
+ Both (Ours)	\checkmark	\checkmark	0.84	0.93	12.32	2.26

Table 6. **Ablation results on visual streams.**

Method	v^f	v^s	DX	FX	MX
Audio-only	\times	\times	0.0265	4.6507	1.3443
+ Facial stream	\checkmark	\times	0.0164	4.3310	1.2028
+ Scene stream	\times	\checkmark	0.0289	4.0939	1.2282
+ Both (Ours)	\checkmark	\checkmark	<u>0.0255</u>	<u>4.2854</u>	<u>1.2077</u>

Table 7. **Ablation results on visual streams with Wrong Placement Ratio (WPR [%]) for each stem.** Lower is better.

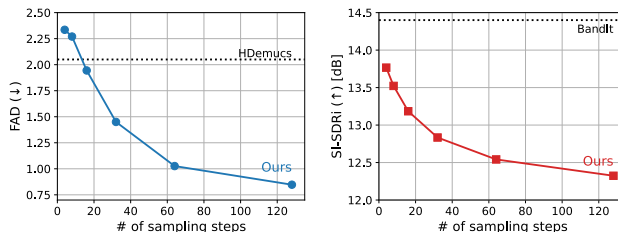


Figure 6. **Impact of the number of sampling steps.** Dotted lines denote Hybrid Demucs (second-best FAD) on the left and BandIt (highest SI-SDRi) on the right.

Impact of sampling step. Since we formulate CASS as a generative task, an important factor is the number of sampling steps N , which directly impacts performance. Fig. 6 shows FAD and SI-SDRi results on AVDnR as we vary $N \in \{4, 8, 16, 32, 64, 128\}$. FAD consistently improves with more steps, indicating better perceptual quality. Also, our model with only 32 steps already surpasses the second-best method, Hybrid Demucs [14], while also improving SI-SDRi. For our experiments, we use 128 steps to maximize perceptual quality, though N can be adjusted depending on the target metric.

6. Conclusion

In this work, we present the first audio-visual framework for cinematic audio source separation (CASS). By shifting from predictive to conditional generative modeling and leveraging multimodal cues, our method delivers high-quality, perceptually realistic separation. We develop a synthetic training pipeline that pairs in-the-wild audio and video, enabling training without cinematic datasets containing clean separated tracks. Trained solely on synthetic data, AV-CASS generalizes seamlessly to real cinematic content. Extensive experiments demonstrate its effectiveness on synthetic benchmark, real-world movie samples, and even standard audio-only CASS benchmarks. These results highlight the potential of an audio-visual perspective for building scalable, generalizable CASS models.

593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

References

- [1] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. The conversation: Deep audio-visual speech enhancement. In *Proc. Interspeech*, 2018. 1, 2
- [2] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. LRS3-TED: a large-scale dataset for visual speech recognition, 2018. 5
- [3] Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised learning of audio-visual objects from video. In *ECCV*, 2020. 2
- [4] Max Bain, Arsha Nagrani, Andrew Brown, and Andrew Zisserman. Condensed movies: Story based retrieval with contextual embeddings. In *ACCV*, 2020. 5, 2
- [5] Moitreyia Chatterjee, Jonathan Le Roux, Narendra Ahuja, and Anoop Cherian. Visual scene graphs for audio source separation. In *ICCV*, 2021. 2
- [6] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP*, 2020. 5
- [7] Jiaben Chen, Renrui Zhang, Dongze Lian, Jiaqi Yang, Ziyao Zeng, and Jianbo Shi. iquery: Instruments as queries for audio-visual sound separation. In *CVPR*, 2023. 1, 2
- [8] Zhuo Chen, Yi Luo, and Nima Mesgarani. Deep attractor network for single-microphone speaker separation. In *ICASSP*, 2017. 1
- [9] Xize Cheng, Siqi Zheng, Zehan Wang, Minghui Fang, Ziang Zhang, Rongjie Huang, Shengpeng Ji, Jialong Zuo, Tao Jin, and Zhou Zhao. Omniseq: Unified omni-modality sound separation with query-mixup. In *ICLR*, 2025. 2
- [10] Joon Son Chung and Andrew Zisserman. Lip reading in the wild. In *ACCV*, 2017. 2
- [11] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. Voxceleb2: Deep speaker recognition. In *Proc. Interspeech*, 2018. 1
- [12] Soo-Whan Chung, Soyeon Choe, Joon Son Chung, and Hong-Goo Kang. Facefilter: Audio-visual speech separation using still images. In *Proc. Interspeech*, 2020. 2
- [13] Michaël Defferrard, Kirell Benzi, Pierre Vandergheynst, and Xavier Bresson. FMA: A dataset for music analysis. In *18th International Society for Music Information Retrieval Conference (ISMIR)*, 2017. 5
- [14] Alexandre Défossez. Hybrid spectrogram and waveform source separation. In *Proceedings of the ISMIR 2021 Workshop on Music Source Separation*, 2021. 1, 5, 7, 8, 4
- [15] Alexandre Défossez, Nicolas Usunier, Léon Bottou, and Francis Bach. Music source separation in the waveform domain, 2019. 1
- [16] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T. Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: a speaker-independent audio-visual model for speech separation. In *Proc. ACM SIGGRAPH*, 2018. 1, 2
- [17] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In *Proc. ICML*, 2024. 1, 2, 4
- [18] Chuang Gan, Deng Huang, Hang Zhao, Joshua B Tenenbaum, and Antonio Torralba. Music gesture for visual sound separation. In *CVPR*, 2020. 2
- [19] Ruohan Gao and Kristen Grauman. Co-separating sounds of visual objects. In *ICCV*, 2019. 2
- [20] Ruohan Gao and Kristen Grauman. Visualvoice: Audio-visual speech separation with cross-modal consistency. In *CVPR*, 2021. 1
- [21] Ruohan Gao, Rogerio Feris, and Kristen Grauman. Learning to separate object sounds by watching unlabeled video. In *ECCV*, 2018. 1, 2
- [22] Yiwei Guo, Chenpeng Du, Ziyang Ma, Xie Chen, and Kai Yu. Voiceflow: Efficient text-to-speech with rectified flow matching. In *ICASSP*, 2024. 1, 2
- [23] Takuya Hasumi and Yusuke Fujita. Dnr-nonverbal: Cinematic audio source separation dataset containing non-verbal sounds. In *Proc. Interspeech*, 2025. 2
- [24] John R. Hershey, Zhuo Chen, Jonathan Le Roux, and Shinji Watanabe. Deep clustering: Discriminative embeddings for segmentation and separation. In *ICASSP*, 2016. 1
- [25] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 3
- [26] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 4
- [27] Vincent Tao Hu, Wei Zhang, Meng Tang, Pascal Mettes, Deli Zhao, and Cees Snoek. Latent space editing in transformer-based flow matching. In *AAAI*, 2024. 1
- [28] Chao Huang, Susan Liang, Yapeng Tian, Anurag Kumar, and Chenliang Xu. DAVIS: High-quality audio-visual separation with generative diffusion models. In *ACCV*, 2024. 2
- [29] Chao Huang, Susan Liang, Yapeng Tian, Anurag Kumar, and Chenliang Xu. High-quality sound separation across diverse categories via visually-guided generative modeling. *arXiv preprint arXiv:2509.22063*, 2025. 2, 5, 6, 7, 1, 4
- [30] Yun-Ning Hung, Chih-Wei Wu, Iroko Orife, Aaron Hipple, William Wolcott, and Alexander Lerch. A large tv dataset for speech and music activity detection. *EURASIP Journal on Audio, Speech, and Music Processing*, 2022. 5
- [31] Chaeyoung Jung, Suyeon Lee, Ji-Hoon Kim, and Joon Son Chung. FlowAVSE: Efficient audio-visual speech enhancement with conditional flow matching. In *Proc. Interspeech*, 2024. 1, 2
- [32] Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. Fréchet audio distance: A metric for evaluating music enhancement algorithms. In *Proc. Interspeech*, 2019. 5, 3
- [33] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5
- [34] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D Plumbley. Panns: Large-scale pretrained audio neural networks for audio pattern recognition. *IEEE/ACM Trans. on Audio, Speech, and Language Processing*, 2020. 5, 3, 4

706	[35] Matthew Le, Apoorv Vyas, Bowen Shi, Brian Karrer, Leda Sari, Rashel Moritz, Mary Williamson, Vimal Manohar, Yossi Adi, Jay Mahadeokar, et al. Voicebox: Text-guided multilingual universal speech generation at scale. In <i>NeurIPS</i> , 2024. 1, 2		
707			
708			
709			
710			
711	[36] Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R Hershey. SDR-half-baked or well done? In <i>ICASSP</i> , 2019. 5, 3		
712			
713			
714	[37] Suyeon Lee, Chaeyoung Jung, Youngjoon Jang, Jaehun Kim, and Joon Son Chung. Seeing through the conversation: Audio-visual speech separation based on diffusion model. In <i>ICASSP</i> , 2024. 3, 1		
715			
716			
717			
718	[38] Jean-Marie Lemerrier, Julius Richter, Simon Welker, and Timo Gerkmann. StoRM: A diffusion-based stochastic regeneration model for speech enhancement and dereverberation. <i>IEEE/ACM Trans. on Audio, Speech, and Language Processing</i> , 2023. 7		
719			
720			
721			
722			
723	[39] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In <i>ICLR</i> , 2023. 1, 2, 3, 4		
724			
725			
726	[40] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. In <i>Proc. ICML</i> , 2023. 5, 3		
727			
728			
729			
730	[41] Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. Diff-foley: Synchronized video-to-audio synthesis with latent diffusion models. In <i>NeurIPS</i> , 2023. 3, 1		
731			
732			
733	[42] Yi Luo and Jianwei Yu. Music source separation with band-split rnn. <i>IEEE/ACM Trans. on Audio, Speech, and Language Processing</i> , 2023. 2		
734			
735			
736	[43] Pingchuan Ma, Stavros Petridis, and Maja Pantic. Visual speech recognition for multiple languages in the wild. <i>Nature Machine Intelligence</i> , 2022. 2		
737			
738			
739	[44] Giorgio Mariani, Irene Tallini, Emilian Postolache, Michele Mancusi, Luca Cosmo, and Emanuele Rodolà. Multi-source diffusion models for simultaneous music generation and separation. In <i>ICLR</i> , 2024. 1, 5, 7, 8, 4		
740			
741			
742			
743	[45] Shivam Mehta, Ruibo Tu, Jonas Beskow, Éva Székely, and Gustav Eje Henter. Matcha-TTS: A fast tts architecture with conditional flow matching. In <i>ICASSP</i> , 2024. 2		
744			
745			
746	[46] Arsha Nagrani, Samuel Albanie, and Andrew Zisserman. Seeing voices and hearing faces: Cross-modal biometric matching. In <i>CVPR</i> , 2018. 2		
747			
748			
749	[47] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In <i>ECCV</i> , 2018. 2		
750			
751			
752	[48] Darius Petermann, Gordon Wichern, Zhong-Qiu Wang, and Jonathan Le Roux. The cocktail fork problem: Three-stem audio separation for real-world soundtracks. In <i>ICASSP</i> , 2022. 1, 2, 5, 6, 7, 8, 4		
753			
754			
755			
756	[49] Darius Petermann, Gordon Wichern, Aswin Shanmugam Subramanian, Zhong-Qiu Wang, and Jonathan Le Roux. Tackling the cocktail fork problem for separation and transcription of real-world soundtracks. <i>IEEE/ACM Trans. on Audio, Speech, and Language Processing</i> , 2023. 2, 7		
757			
758			
759			
760			
761	[50] Antony W Rix, John G Beerends, Michael P Hollier, and Andries P Hekstra. Perceptual evaluation of speech quality	(PESQ)-a new method for speech quality assessment of tele-	763
762		phone networks and codecs. In <i>ICASSP</i> , 2001. 5, 4	764
		[51] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In <i>CVPR</i> , 2022. 4	765
			766
			767
		[52] Simon Rouard, Francisco Massa, and Alexandre Défossez. Hybrid transformers for music source separation. In <i>ICASSP</i> , 2023. 1, 5, 7, 8, 4	768
			769
			770
		[53] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. <i>IEEE TPAMI</i> , 2022. 7	771
			772
			773
		[54] Akira Takahashi, Shusuke Takahashi, and Yuki Mitsufuji. MMAudioSep: Taming video-to-audio generative model towards video/text-queried sound separation, 2025. 2	774
			775
			776
		[55] Efthymios Tzinis, Scott Wisdom, Aren Jansen, Shawn Hershey, Tal Remez, Daniel PW Ellis, and John R Hershey. Into the wild with audioscope: Unsupervised audio-visual separation of on-screen sounds. In <i>ICLR</i> , 2021. 1, 2	777
			778
			779
			780
		[56] Efthymios Tzinis, Scott Wisdom, Tal Remez, and John R Hershey. Audioscopev2: Audio-visual attention architectures for calibrated open-domain on-screen sound separation. In <i>ECCV</i> , 2022. 2	781
			782
			783
			784
		[57] Stefan Uhlich, Giorgio Fabbro, Masato Hirano, Shusuke Takahashi, Gordon Wichern, Jonathan Le Roux, Dipam Chakraborty, Sharada Mohanty, Kai Li, Yi Luo, Jianwei Yu, Rongzhi Gu, Roman Solovyev, Alexander Stempkovskiy, Tatiana Habruseva, Mikhail Sukhovei, and Yuki Mitsufuji. The sound demixing challenge 2023 – cinematic demixing track. <i>Trans. of the International Society for Music Information Retrieval</i> , 2024. 1	785
			786
			787
			788
			789
			790
			791
			792
		[58] Yongqi Wang, Wenxiang Guo, Rongjie Huang, Jiawei Huang, Zehan Wang, Fuming You, Ruiqi Li, and Zhou Zhao. Frieren: Efficient video-to-audio generation network with rectified flow matching. In <i>NeurIPS</i> , 2024. 1, 2	793
			794
			795
			796
		[59] Karn N. Watcharasupat, Chih-Wei Wu, Yiwei Ding, Iro Orife, Aaron J. Hipple, Phillip A. Williams, Scott Kramer, Alexander Lerch, and William Wolcott. A generalized band-split neural network for cinematic audio source separation. <i>IEEE Open Journal of Signal Processing</i> , 2023. 1, 2, 5, 6, 7, 8, 4	797
			798
			799
			800
			801
			802
		[60] Karn N. Watcharasupat, Chih-Wei Wu, and Iro Orife. Remastering divide and remaster: A cinematic audio source separation dataset with multilingual support. In <i>2024 IEEE 5th International Symposium on the Internet of Sounds</i> , 2024. 1, 2, 5, 7, 8	803
			804
			805
			806
			807
		[61] Yutong Wen, Ke Chen, Prem Seetharaman, Oriol Nieto, Jiaqi Su, Rithesh Kumar, Minje Kim, Paris Smaragdakis, Zeyu Jin, and Justin Salamon. PromptSep: Generative audio separation via multimodal prompting, 2025. 7	808
			809
			810
			811
		[62] Dong Yu, Morten Kolbæk, Zheng-Hua Tan, and Jesper Jensen. Permutation invariant training of deep models for speaker-independent multi-talker speech separation. In <i>ICASSP</i> , 2017. 1	812
			813
			814
			815
		[63] Yinfeng Yu and Shiyu Sun. Dgfnet: End-to-end audio-visual source separation based on dynamic gating fusion. In <i>Proc. ACM ICMR</i> , 2025. 2	816
			817
			818

- 819 [64] Yi Yuan, Xubo Liu, Haohe Liu, Mark D Plumbley, and
820 Wenwu Wang. FlowSep: Language-queried sound separa-
821 tion with rectified flow matching. In *ICASSP, 2025*. 2, 7
- 822 [65] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding
823 conditional control to text-to-image diffusion models. In
824 *ICCV, 2023*. 5
- 825 [66] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Von-
826 drick, Josh McDermott, and Antonio Torralba. The sound of
827 pixels. In *ECCV, 2018*. 2