Learning Spatial Features from Audio-Visual Correspondence in Egocentric Videos

Sagnik Majumder^{1,2} Ziad Al-Halah³ Kristen Grauman^{1,2}
¹UT Austin ²FAIR, Meta ³University of Utah

Abstract

We propose a self-supervised method for learning representations based on spatial audio-visual correspondences in egocentric videos. Our method uses a masked auto-encoding framework to synthesize masked binaural (multi-channel) audio through the synergy of audio and vision, thereby learning useful spatial relationships between the two modalities. We use our pretrained features to tackle two downstream video tasks requiring spatial understanding in social scenarios: active speaker detection and spatial audio denoising. Through extensive experiments, we show that our features are generic enough to improve over multiple state-of-theart baselines on both tasks on two challenging egocentric video datasets that offer binaural audio, EgoCom and Easy-Com. Project: http://vision.cs.utexas.edu/projects/ego_av_corr.

1. Introduction

Egocentric videos provide a first-person view of how we perceive and interact with our surroundings in daily life, and they are pushing a new frontier in multi-modal learning [10, 27, 35, 70]. A key aspect of ego-video is that it can provide a rich stream of first-person spatial audio lalongside visual frames when the audio is captured with multiple microphones [12, 56]. The coupling of such audio and vision provides strong spatial information about the sound sources (e.g. where the sound sources are, if they are in motion or not) in the context of the surrounding physical space (e.g. how big or small the room is, if there is a large wall nearby), as well as the camera wearer's attention in the scene revealed by how they move their head.

Such spatial cues are especially important for social settings of multiple people talking to each other, where it is valuable to be able to 1) focus on the voice(s) of interest from among various competing sounds and 2) understand

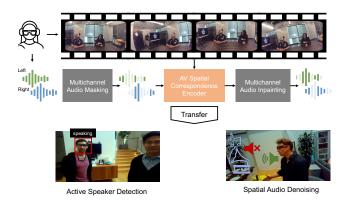


Figure 1. Given an egocentric video and binaural audio, we aim to learn spatial correspondences between vision and audio by solving the pretext task of inpainting segments of the binaural audio. The features benefit downstream social tasks where spatial localization is important: active speaker detection and audio denoising.

where people are directing their speech activity, for better comprehension and communication. In this way, future AR applications in conversational settings could allow a hearing-impaired person to determine who is speaking in order to redirect their attention, or enhance the received audio to make it more intelligible for any listener.

We argue that this creates the need for human-centric spatially-grounded understanding of audio-visual events. Can we learn representations from video that capture audio-visual events in the context of the persistent physical space of the environment and the human speakers in it? Such representations would be useful for answering questions like "who is speaking right now?" and "what would the voices sound like without the audio noise (distractor sounds)?" While the former requires inferring the source location for a voice in the scene, the latter requires understanding how the perceived audio is a function of the source locations, the listener, and the surrounding environment.

Despite the significance of these problems, today's models for audio-visual learning are not human-centric and they lack spatial grounding. On the one hand, current audio-visual representation learning methods exclusively tackle exocentric (third-person) video with monaural audio

¹Throughout we use the term *spatial audio* to refer to binaural audio, including the special case of two-channel *binaural* audio as perceived by two human ears. In contrast, single-channel *monaural* audio lacks spatial information.

[2, 24, 26, 33, 40, 57, 59]. That domain sidesteps challenges inherent to ego-video arising from the camera wearer's head motion and relatively limited field of view. On the other hand, the limited prior work exploring self-supervised objectives using multi-channel audio and video [19, 23, 51, 78] are also outside of the egocentric and social contexts (e.g., for sounds in empty homes [23], musical instruments [19], or single human speakers [78]), where the need for spatial understanding of multiple sound sources is limited.

We propose to learn audio-visual representations via spatial correspondence between an egocentric video and its binaural audio, for analyzing social (conversational) settings. In particular, we design a novel pretext task where the goal is to inpaint binaural (two-channel) audio using both video and audio. Given a social egocentric video clip with binaural audio, we mask segments of it and train a model based on masked autoencoding (MAE) [6, 16, 29, 34, 74] to predict the missing segments on the basis of the video and the unmasked segments in the audio. See Figure 1 (top). Additionally, we introduce a novel spatial audio masking strategy that facilitates learning strong audio-visual spatial correspondences while maintaining learning stability when vision alone is insufficient for the binauralization task. Once trained, our model's encoder provides spatial audio-visual features that can be used to address multiple downstream tasks, as we demonstrate using multiple different backbones and social egocentric video datasets.

In particular, motivated by the AR applications discussed above, we validate our feature learning method on two downstream social egocentric tasks that require strong audiovisual spatial reasoning: 1) active speaker detection: predicting which person in the field of view of an egocentric video is speaking, and 2) spatial audio denoising: separating audio noise (any sounds from non-speakers) from the input audio. See Figure 1 (bottom). We test the generality of our method by evaluating on two social egocentric video datasets, EgoCom [56] and EasyCom [12]—to our knowledge, the only two publicly available video datasets with binaural sound and social settings. On both, our method significantly outperforms multiple state-of-the-art task-specific and audio-visual spatial feature learning models.

2. Related Work

Audio-visual self-supervised pretraining Past work [2, 4, 40, 51, 55, 57–59] extensively studies the synergy of vision and audio for learning representations through self-supervision. They explore using both modalities to construct pretext tasks based on synthesis [58, 59], alignment [2, 4, 23, 40, 57], and masked auto-encoding (MAE) [24, 26, 33], and they focus on *semantic* downstream tasks like audio-visual event classification and retrieval. However, none of these methods are designed to extract spatial cues from video and multi-channel audio, nor do they analyze the social egocen-

tric setting. On the contrary, we propose self-supervised learning of spatial audio-visual features from egocentric video. Further, different from existing MAE-style models [24, 26, 33], we propose a specialized masking strategy that better learns spatial audio-visual cues. Our masking idea promotes the encoding of spatial and semantic information in the learned multimodal representation, thereby improving its effectiveness for transfer learning in downstream tasks that require nuanced reasoning about both *what* and *where* aspects, such as active speaker detection and spatial audio-visual denoising. This differs from previous methods [24, 26, 33], which mainly use a learning objective that emphasizes the encoding of semantic cues and tailor to tasks like multimodal event classification or retrieval.

Audio-visual spatial correspondence learning Learning the *spatial* alignment between video and audio is important for self-supervision [51, 68, 77, 78], audio generation [8, 19, 46, 50, 52, 64, 81], audio-visual embodied learning [7, 9, 44, 45] and 3D scene mapping [47, 63]. However, these methods are either restricted to exocentric settings [8, 19, 50, 51, 64, 68, 78], or else tackle egocentric settings [9, 44, 46, 47] in simulated 3D environments that lack realism and diversity, both in terms of the audio-visual content of the videos (no people are visible, no objects are moving) and their lack of continuous camera motion from a camera-wearer's physical movements. In contrast, we learn an audio-visual representation from real-world egocentric video.

More closely related to our work are Telling Left from Right (TLR) [78], 2.5D Visual Sounds (2.5D-VS) [19], and audio-visual stereo sound ranking (SSR) [68], all of which learn spatial audio-visual features, albeit for exocentric data only. TLR predicts whether the left and right binaural channels are swapped, and SSR ranks the similarity of video to different stereo audio samples through self-supervision both of which provide only coarse spatial information about the scene. 2.5D-VS learns to "lift" the mono input to binaural audio, which can be underconstrained from the singlechannel audio and video alone. We design a novel pretext task using audio-visual inpainting of binaural audio, which is both fine-grained (requiring to capture subtleties about the arrangement of speakers in the environment) and, through our novel masking strategy, exposes better multi-modal constraints that improve learning stability and performance. Our results in Sec. 4 show our model's advantages over all three prior methods [19, 68, 78].

Active speaker detection Active speaker detection (ASD) entails predicting the active speaker(s) from among all detected faces in a video, and is a special case of generic 2D sound localization [18, 31, 37, 49, 57, 79]. While early ASD methods rely on lip motion and facial gestures [15], recent methods employ ensemble networks [3] or 3D CNNs [39, 69, 73], relation context modules [80],

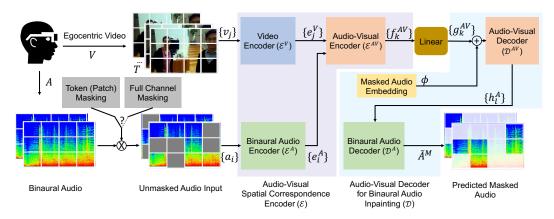


Figure 2. Our model learns the spatial correspondence between vision and binaural audio by inpainting masked tokens of the audio channels through the use of an audio-visual encoder-decoder model. We combine random token masking (which requires solving a more local binauralization task) with complete audio channel masking (which requires more global cues to synthesize unseen binaural segments). For downstream evaluation, we fuse the features from the audio-visual encoder with the backbones for downstream tasks, and finetune them.

attention [3, 73], or graph neural networks [41, 48]. Multichannel audio improves ASD [37], but requires privileged information of the speaker's pose for training. Recent work explores using supervised learning to infer not only who is talking, but also to whom a camera-wearer is listening [67]. In contrast, our goal is to learn spatial audio-visual features purely from in-the-wild egocentric video through self-supervision—features generic enough to benefit multiple ASD models, as we show for both TalkNet [73] and SPELL [48].

Spatial audio denoising Audio denoising, which requires separating a target sound from noise, has traditionally been studied with single-channel (non-spatial) audio, both in the audio-only setting [32, 71, 72, 76] and audio-visual settings [1, 14, 17, 20–22, 57, 60, 62]. Using spatial audio captured with multiple microphones [13, 54, 82] naturally makes the task simpler. Different from the above, we learn task-agnostic audio-visual spatial features. That is, our contribution is the feature learning idea (which benefits both denoising and ASD), rather than a novel denoising approach.

3. Learning spatial features from egocentric audio-visual correspondence

The spatial sound perceived in an egocentric setup is shaped by the environment in which it is emitted and its source location relative to the camera-wearer. Based on this knowledge, we hypothesize that trying to solve the pretext task of audio-visual inpainting of binaural audio—that is, synthesizing missing audio segments by extracting related visual cues about the scene and the source location—can lead to learning useful audio-visual spatial correspondences. To validate our hypothesis, we propose a novel feature-learning task.

Formally, we consider an egocentric video clip C = (V, A), where V and A refer to the visual and binaural

audio streams, respectively. See Fig. 2 left. The visual clip V comprises T frames, such that $V = \{V_1, \dots, V_T\}$. We generate a set of visual tokens \hat{V} by splitting V into P non-overlapping tubelets, such that $\hat{V} = \{\hat{V}_1, \dots, \hat{V}_P\}$, where \hat{V}_k denotes the k^{th} tubelet consisting of a contiguous sequence of patches spanning all T frames. We represent the binaural audio A as two Mel-spectrograms [34], $A = \{A^L, A^R\}$, where A^L and A^R are the spectrograms for the left and right channels, respectively. We create a set of audio tokens \hat{A} by splitting A into Q non-overlapping patches, such that $\hat{A} = \{\hat{A}_1, \dots, \hat{A}_Q\}$.

Next, we mask a portion of the audio tokens in \hat{A} and obtain complementary subsets of masked and unmasked tokens, \hat{A}^M and \hat{A}^U , respectively, where $\hat{A}^M = \{\ddot{A}_1, \ldots, \ddot{A}_S\}$, $\hat{A}^U = \{\bar{A}_1, \ldots, \bar{A}_{Q-S}\}$, and S is the number of masked tokens. Given $\{\hat{V}, \hat{A}^M, \hat{A}^U\}$, we aim to learn a self-supervised model \mathcal{F} comprising an encoder \mathcal{E} and decoder \mathcal{D} , such that $\mathcal{F} = \mathcal{D} \circ \mathcal{E}$ and $\mathcal{F}(\hat{V}, \hat{A}^U) = \tilde{A}^M$, where \tilde{A}^M is an estimate of the masked audio tokens in \hat{A}^M . By training on this pretext task, our encoder \mathcal{E} can learn rich audio-visual spatial correspondences that can be leveraged for multiple downstream tasks that require the synergy of vision and spatial audio, as we show in results.

In our method (see Fig. 2), \mathcal{E} (Sec. 3.1) is an audio-visual (AV) spatial correspondence encoder that learns an implicit representation of the spatial relationships between the visual and unmasked binaural audio tokens, while \mathcal{D} (Sec. 3.2) is an audio-visual decoder for binaural audio inpainting that uses this implicit representation to synthesize the masked audio tokens. We also devise a simple yet novel masking protocol (Sec. 3.3) for our inpainting task, which mixes masking random audio tokens with masking a full audio channel, and helps the model learn stronger audio-visual spatial associations that facilitate the downstream social tasks (Sec. 3.5). We train \mathcal{F} to minimize the prediction error in the masked

audio tokens (Sec. 3.4). Next, we describe our model design, audio masking protocol, training objective and network architecture, and downstream tasks.

3.1. Audio-visual spatial correspondence encoder

The audio-visual spatial correspondence encoder \mathcal{E} (Fig. 2 center) extracts features from the visual and unmasked audio tokens $\{\hat{V}, \hat{A}^U\}$. It begins by embedding the visual and audio tokens using separate transformer encoders [16] for individually capturing the spatio-temporal features in the two modalities. Next, it uses a shared transformer encoder to jointly encode the audio and visual features, and produces a multi-modal representation suitable for binaural audio inpainting. Next, we describe the individual encoders next.

Video and audio encoders. We first encode the visual tokens \hat{V} using a linear layer to generate visual features v, such that $v = \{v_1, \dots, v_P\}$. We encode the audio tokens \hat{A}^U with another linear layer to produce audio features a, such that $a = \{a_1, \dots, a_{Q-S}\}$, where S is the number of masked tokens out of a total of Q audio tokens (cf. Sec. 3). For each visual feature v_i , we add a sinusoidal positional embedding p_i^V [75] to it, where p_i^V captures cues about the 3D position of the j^{th} tubelet in the visual clip V. For an audio feature a_i , we add a sinusoidal positional embedding p_i^A and a learnable channel embedding $c \in \{c_L, c_R\}$ to it to convey information about the 2D location of the ith unmasked audio token in the spectrogram and also the audio channel to which it belongs. Next, we feed the transformed visual and audio features to separate transformer encoders, \mathcal{E}^V and \mathcal{E}^A , respectively, and obtain visual features $e^V = \{e_1^V, \dots, e_P^V\}$ and audio features $e^A = \{e_1^A, \dots, e_{Q-S}^A\}$.

Shared audio-visual encoder. Given the visual features e^V and audio features e^A , we concatenate them into e^{AV} , such that $e^{AV} = \left\{e_1^V, \ldots, e_P^V, e_1^A, \ldots, e_{Q-S}^A\right\}$, and re-add the sinusoidal positional embeddings p^V and p^A to the features of the respective modalities in e^{AV} . Furthermore, unlike existing audio-visual masked auto-encoders [24, 26, 33], we add the channel embeddings e^V to the audio features, and learnable modality embeddings e^V to the audio features, and learnable modalities. Next, a shared audio-visual transformer e^V encoder takes e^V as input and outputs audio-visual features e^V , which implicitly hold spatio-temporal information required for accurate inpainting of audio.

3.2. Audio-visual decoder for binaural audio inpainting

Our audio-visual decoder \mathcal{D} (Fig. 2 right) takes f^{AV} as input and attempts to synthesize the masked binaural audio tokens by leveraging spatio-temporal cues in f^{AV} . It first

projects f^{AV} to a lower-dimensional feature set g^{AV} . It then appends a learnable embedding for the masked audio tokens to g^{AV} and passes it through a shared audio-visual transformer decoder [29]. Next, it feeds the audio feature outputs of the shared decoder to another transformer decoder and uses its outputs to predict an estimate of the masked binaural audio tokens. The decoders are lightweight compared to the encoders, ensuring that the encoders are primarily responsible for driving the inpainting task and producing good audio-visual features for strong downstream performance. We next describe each component of $\mathcal D$ in detail.

Shared audio-visual decoder. We first create a lower-dimensional projection g^{AV} of the audio-visual encodings f^{AV} by passing it through a linear layer, and append a learnable embedding ϕ corresponding to each of the S masked audio tokens to g^{AV} . Next, we add the positional embeddings p^V and p^A , the audio channel embeddings c, and the modality embeddings m to g^{AV} , and feed it to a shallow transformer decoder \mathcal{D}^{AV} that outputs an audio-visual feature set h^{AV} . We then take the audio features h^A from h^{AV} and pass them to the audio decoder for further processing.

Audio decoder. The audio decoder \mathcal{D}^A re-adds the positional embeddings p^A and channel embeddings c to g^A , and feeds it to a transformer decoder, which outputs audio features d^A .

Prediction of masked audio tokens. Finally, we take the subset of all audio features d^A corresponding to the masked audio tokens \hat{A}^M , upsample them by passing through a linear layer, and reshape them to obtain an estimate \tilde{A}^M of the masked tokens \hat{A}^M , such that $\tilde{A}^M = \{\tilde{A}_1, \dots, \tilde{A}_S\}$.

3.3. Audio masking

Different from other masked auto-encoding counterparts [24, 26, 33], we design an audio masking protocol that is customized to help our model better extract spatial audio-visual cues during self-supervised pretraining. In particular, we mix the strategy of randomly masking a full audio channel with that of randomly masking audio tokens using a hyperparameter r during training, where r represents the probability with which we randomly drop a full audio channel and r is sampled from a uniform distribution U(0,1):

$$\operatorname{mask}(\hat{A}) = \begin{cases} \hat{A}^M = A^L \text{ or } \hat{A}^M = A^R & \text{if} \quad x \sim U(0,1) \leq r \\ \hat{A}^M \subseteq \{\hat{A}_1, \dots, \hat{A}_Q\} & \text{Otherwise} \end{cases}$$

On the one hand, *token masking* could lead to tokens from the same location in the two audio channels being present among the unmasked tokens, providing additional spatial cues to the model and resulting in a simpler, stabler optimization objective for the inpainting task. In addition, since token masking involves masking a short span in both frequency and

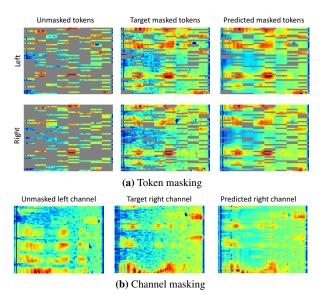


Figure 3. Masked targets and predictions shown alongside the unmasked inputs for (a) token masking and (b) channel masking. Our predictions accurately capture the global patterns in the target spectrograms, which depend on the scene's spatial properties.

time domains, the model can rely more on local audio-visual cues and tolerate the global noise in both the visual and audio streams due to a camera-wearer's motion. On the other hand, *channel masking* forces the model to solve a more challenging binauralization task solely on the basis of vision, which could help it learn even stronger spatial features. This encourages the model to reason about the camera motion at a more global scale (over the entire clip span). Towards achieving high performance on the downstream tasks, we aim to strike a fine balance between these two strategies and combine the benefits of reasoning at both temporal scales.

3.4. Training objective and network architecture

We train our model to minimize the error in predicting the masked audio tokens. In particular, we compute the mean-squared error \mathcal{L} averaged over all masked audio tokens:

$$\mathcal{L} = \frac{1}{S} \sum_{i=1...S} ||\ddot{A}_i - \tilde{a}_i||_2^2.$$
 (1)

We visualize our predicted audio tokens in Fig. 3 for the cases of token (Fig. 3a) and channel (Fig 3b) masking. Our model is able to accurately predict the masked targets and capture the global patterns in the spectrograms, which are often determined by the spatial audio-visual cues captured from of the scene (visual input not shown in Fig 3 for brevity), thereby further emphasizing our model's ability to learn useful spatial features.

Our uni-modal encoders, \mathcal{E}^A and \mathcal{E}^V , have 8 layers, while the audio-visual encoder \mathcal{E}^{AV} has 6 layers. All encoders have 12 attention heads and use 768-dimensional hidden

embeddings. The audio-visual decoder \mathcal{D}^{AV} and audioonly decoder \mathcal{D}^A have 1 and 3 layers, respectively. Both decoders have 6 attention heads and use 384-dimensional hidden embeddings. To pretrain our model, we set the relative frequency of dropping an audio channel in our masking protocol for training to r=20% based on disjoint validation data (see Supp.). We train our model for 200 epochs using the AdamW [43] optimizer with a weight decay of 10^{-5} , and a learning rate scheduler that reaches a peak learning rate of 2×10^{-4} over 10 warmup epochs, and then decays it through half-cycle cosine annealing [42]. For data agumentation, we perform random flipping of audio channels and video clips along the frame width. For downstream evaluation, we fuse the features from the audio-visual encoder with the backbones for downstream tasks, and finetune them. See Supp. for further details on architecture and training.

3.5. Downstream tasks

We explore two downstream tasks with our pretrained features: active speaker detection and spatial audio denoising. Active speaker detection (ASD) involves matching an audio clip with an appropriate face track from the corresponding video clip, i.e., answering "which person is speaking now?". While current SOTA methods [48, 73] rely on semantic similarities between monaural audio and vision to solve this task, we hypothesize that leveraging spatial audio can additionally reveal the sound source location in the video. In spatial audio denoising, the goal is to separate the target audio from distractors. In particular, we aim to remove the audio from sources extraneous to the conversation—out-of-view sounds from other parts of the scene. We detail the backbone models for each in the next section.

4. Experiments

We validate our learned representations on two downstream tasks and two datasets, and we compare with prior models for spatial audio-visual feature learning [19, 24, 68, 78], as well as various baselines and ablations.

Datasets. We train and evaluate our model on two challenging egocentric datasets containing video and binaural audio of people having conversations: 1) EgoCom [56], and 2) EasyCom [12], detailed in Supp. To our knowledge, these are the only two publicly available datasets offering binaural audio with conversations in video, whether exocentric or egocentric. In particular, Ego4D [27] and EPIC [10] do not comprise social scenarios and are not applicable. Whereas EasyCom primarily has participants sitting around a table and talking, EgoCom has videos of participants moving around a room, turning their face and body, standing up, etc. They test our method's robustness in diverse scenarios of varying difficulty. See Supp. for more details.

4.1. Active speaker detection

We first evaluate on active speaker detection (ASD).

Backbone models. We consider two SOTA ASD models as the backbones for leveraging our pretrained representations: 1) TalkNet [73], and 2) SPELL [48]. TalkNet encodes a face track and an audio clip using attention for learning intra- and inter-modal semantic and temporal features. Next, it fuses these features and performs binary classification to predict if the face in the track is active. SPELL extracts audiovisual features for each face in a clip using ResNets [3], and learns long and short-term semantic relations among them using a graph neural network. Finally, it performs binary classification of these features for predicting active speakers.

Pretrained feature fusion. To fuse our pretrained features with the ASD backbones, we use a transformer decoder that cross-attends to the feature outputs of the shared encoder \mathcal{E}^{AV} using sinsuoidal embedding as queries, with each embedding representing a clip frame index. Next, we append the decoder outputs to the cross-attention outputs in TalkNet, or the audio-visual encoder outputs in SPELL, frame by frame. In essence, while the original audio-visual encoders leverage *semantic* correlations between vision and audio, our features provide strong complementary *spatial* cues.

Baselines. For both TalkNet and SPELL, we compare against multiple baselines comprising both the unmodified backbone and improved versions of it, in addition to some naive methods:

- All-active: a naive model that predicts that all visible faces are always active
- All-inactive: a naive model that predicts that all visible faces are always inactive
- Random: a naive model that emits a random ASD confidence score for every visible speaker
- Backbone w/o audio: a vision-only version of the backbone with no audio input
- Backbone: the originally-proposed backbone that processes only faces and monaural audio
- Backbone-binaural: an improved backbone using binaural audio instead of monaural, alongside positional encodings for the faces, indicative of their relative position and depth, for better matching the face to the audio
- Backbone-binaural w/ scene: a further improvement over the backbone, where we also provide the scene images (uncropped video frames) to backbone-binaural
- Backbone w/ TLR [78]: fuses features from the SOTA Telling Left from Right (TLR) [78], which learns audiovisual spatial correspondences by predicting the spatial alignment between vision and binaural audio
- Backbone w/ 2.5D-VS [19]: fuses features from the SOTA audio-visual binauralization model, 2.5D Visual Sounds (2.5D-VS) [19]

	TalkN	let [73]	SPEL	L [48]
Model	EgoCom	EasyCom	EgoCom	EasyCom
No pretraining				
All-active	32.9	30.1	32.9	30.1
All-inactive	32.9	30.1	32.9	30.1
Random	30.8	28.0	30.8	28.0
B w/o audio	41.5	50.1	60.4	63.2
В	52.8	45.7	60.9	59.0
B-binaural	60.0	59.6	63.1	60.3
B-binaural w/ scene	60.8	66.9	61.2	61.4
Alternate pretraining method	s			
B w/ TLR [78]	47.9	59.3	61.3	61.7
B w/ 2.5D-VS [19]	57.7	63.7	61.2	59.7
B w/ 2.5D-VS [19]++	63.4	68.3	65.1	64.5
B w/ SSR [68]++	61.2	70.6	61.2	67.4
B w/ AV-MAE [24]	61.1	61.3	64.4	65.2
Ours	63.9	71.8	65.6	70.2
Ours w/o pretrain	62.7	62.9		
Ours w/ pretrain monaural	61.0	69.4	63.9	69.0

Table 1. Mean average precision (%) for active speaker detection with TalkNet [73] and SPELL [48] backbones on both datasets. Higher is better. 'B' refers to backbone. Note that SPELL requires storing pretrained features in the graph nodes; therefore it does not allow training from scratch.

- Backbone w/ 2.5D-VS [19]++: fuses features from 2.5D-VS with a transformer architecture
- Backbone w/ SSR [68]++: fuses features from the SOTA self-supervised audio-visual stereo sound ranking (SSR) [68] model with a transformer architecture
- Backbone w/ AV-MAE [24]: fuses features from the SOTA modality-inpainting AV-MAE [24] model

For all alternate feature-learning methods [19, 24, 68, 78], we pretrain them on our datasets and use our feature fusion method. Thus, any advantages in performance of our approach over these SOTA representations will be attributable to our modeling ideas. Importantly, the 2.5D-VS [19]++, SSR [68]++, and AV-MAE [24] features all rely on transformers and have similar model capacity as ours (see Supp. for a detailed analysis on model capacity). We use the standard **mean average precision** (mAP) metric.

Results. Table 1 (top) reports our ASD results on both val and test splits. The three naive baselines perform poorly on both EgoCom [56] and EasyCom [12], emphasizing the difficulty of the task. For both TalkNet [73] and SPELL [48], the unchanged backbone model generally performs better than the model without audio, showing that both vision and audio are required. Upgrading from monaural to binaural audio further boosts performance, as the model can now leverage both spatial and semantic information. Additionally using scene features lets the backbone explicitly match the scene area around the inferred source location with the face, and further improves ASD, especially for EgoCom, where the background scene changes more often.

Model	SI-SDRi \uparrow	$STFT\downarrow$
No pretraining		
B w/o vision	1.61	7.36
В	1.46	7.27
B w/ ImageNet features	1.48	6.95
Alternate pretraining methods		
B w/ TLR [78]	1.41	7.79
B w/ 2.5D-VS [19]	1.67	7.34
B w/ 2.5D-VS [19]++	2.11	6.60
B w/ SSR [68] ++	2.04	6.70
B w/ AV-MAE [24]	2.07	6.62
Ours	2.20	6.51
B w/o pretrain	1.90	7.25
B w/ pretrain monaural audio	2.00	6.75

Table 2. Audio denoising with U-Net [78] backbone (referred to as 'B' in table) for 0 dB noise (maximum). See Supp. for varying noise levels. All STFT distance measures use base 10^{-3} .

Among alternate feature learning methods, 2.5D-VS [19]++, SSR [68]++ and AV-MAE [24] consistently outperform TLR [78] and 2.5D-VS [18], and also the basic and enhanced backbones, showing that self-attention and higher model capacity enhance feature quality. Besides, 2.5D-VS outperforms TLR, and 2.5D-VS++ and AV-MAE generally outperform SSR++, indicating that objectives that promote reasoning directly at the spectrogram level improve results.

Our model substantially outperforms all baselines—including the SOTA AV representation learning methods—for both backbones (TalkNet and SPELL) on both datasets. This shows that our method learns stronger spatial features that are both backbone- and dataset-agnostic. In contrast, methods developed for exocentric settings with more stationary cameras (such as TLR and 2.5D-VS) rely more on the global visual context and seem to struggle in our setting, where the camera moves frequently and the sound source leaves the field of view. Finally, our improvement over 2.5D-VS++, SSR++ and AV-MAE, which use similar encoders as ours, disentangles the benefits of our masking strategy and model design from those of the model capacity.

Model analysis. Table 1 (bottom) shows ablations of our method. Upon training for ASD from scratch, we see sharp drop in performance, showing that our advantage is not solely due to our model design, but also our self-supervised pretraining stage. Performance also declines upon pretraining with monaural audio instead of binaural, showing that our model goes beyond learning semantic features and successfully captures spatial features useful for ASD.

4.2. Spatial audio denoising

Next we evaluate spatial audio denoising on EgoCom.² To instantiate this task, we add the binaural audio of a target clip with the downscaled binaural audio from another randomly chosen clip, where the downscaling factor depends on the desired noise level. The goal is to extract the target sound from the mixture. We evaluate three noise levels, expressed using the signal-to-noise (SNR) ratio: 1) 0 dB, 2) 2.5 dB, and 3) 5 dB. The different noise levels test our model's robustness to varying levels of task difficulty—the lower the SNR value, the higher the noise and difficulty.

Backbone model. We adopt the commonly used U-Net [65] for audio-visual source separation [19, 78] as the backbone, which produces a binaural ratio mask for the target audio (see Supp. for details). We multiply the predicted ratio mask with the mixed magnitude spectrogram to get the predicted magnitude spectrogram, then convert it to a waveform using inverse short-time Fourier transform with the mixed audio phase.

Pretrained feature fusion. To use our features for denoising, we reshape the visual features and unmasked audio features produced by our audio-visual encoder \mathcal{E}^{AV} to form multi-channel 2D maps, where the features align with their corresponding tokens vis-a-vis the raster order. Next, we pass the feature maps to separate convolutional layers, concatenate the outputs channel-wise, and use them to replace the visual features at the U-Net [78] bottleneck.

Baselines. We compare against the following baselines and existing methods:

- U-Net w/o vision: an audio-only blind denoising model
- U-Net: the original backbone without any alterations
- U-Net w/ ImageNet: pretrains the visual encoder on ImageNet [11]
- U-Net w/ TLR [78]: fuses features from TLR [78]
- U-Net w/ 2.5D-VS [19]: fuses pretrained features from 2.5D-VS [19]
- U-Net w/ 2.5D-VS [19]++: fuses features from the transformer-based version of 2.5D-VS
- U-Net w/ SSR [68]++: fuses features from the transformer-based version of SSR [68]
- U-Net w/ AV-MAE [68]: fuses features from the modality inpainting AV-MAE [24] model

Evaluation metric. For evaluating our denoising quality, we use standard metrics: 1) **STFT** distance (the L2 error

²For EasyCom, the task setup is ill-posed for all models because mixing audio from a different EasyCom clip as noise leads to spatially overlapping sound sources, since all clips in the dataset are recorded at the same physical location (people sitting around the same table in the same room).

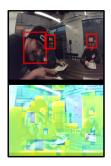






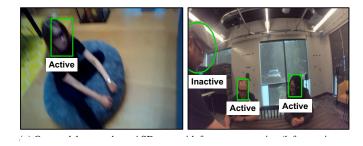
Figure 4. Heat maps showing the image areas our model's AV encoder attends to, placed alongside the images. Brighter yellow means higher attention score. Our model attends to image regions (*e.g.* faces of speakers, sound-reflecting flat regions like floor and table, etc.) that strongly determine the spatial properties of the audio, including direct sources of sound (marked in red).

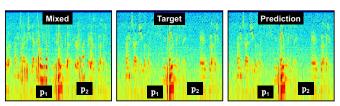
between the predicted and ground-truth spectrograms) expressed using base 10^{-3} and **2**) **SI-SDRi**: the improvement in SI-SDR [66], a scale-invariant estimate of the level distortion in the audio, over using the mixed audio as the prediction.

Results. Table 2 (top) shows spatial audio denoising results on the challenging EgoCom dataset with 0 dB, the most difficult noise level. See Supp. for similar results with the other noise levels. The unchanged U-Net backbone lowers the STFT distance compared to the version that lacks vision, showing that like ASD, vision is crucial for better denoising. Using pretrained features of 2.5D-VS [19] (++), SSR [68]++ or AV-MAE [24] further improves the performance, showing that learning spatial audio-visual features aids denoising.

Our method outperforms all baselines ($p \leq 0.05$) across both metrics. While the improvement over the baselines that do not use self-supervised pretraining emphasizes the utility of learning spatial audio-visual relationships through self-supervision, the gain over other pretraining methods underlines the strengths of our self-supervised method design—which are consistently realized for both ASD and denoising. Further, our performance margins are larger for higher noise levels (0 and 2.5 dB), indicating that our features play a bigger role in the more difficult denoising settings.

Model analysis. In Table 2 (bottom), we ablate our pretraining method. Similar to ASD, training from scratch on the denoising task hurts performance. This disentangles the impact of our pretext task design from the model architecture and shows that our pretraining stage helps the backbone with learning better audio-visual features, leading to superior denoising quality. Furthermore, pretraining with monaural audio also degrades performance, re-emphasizing that our method is not restricted to learning semantic features—in contrast to prior work [24, 26, 33].





(b) Our model denoises accurately—note the noise patches in mixed audio above points p_1 and p_2 , which are successfully removed in our prediction.

Figure 5. Success cases for ASD (a) and denoising (b)

See Supp. for additional analysis of the effect of alternate masking choices, multi-level positional embeddings, tasks-specific backbones, and our finetuning strategy on performance.

4.3. Qualitative analysis

In Fig. 4, we analyze the visual attention maps of our shared encoder \mathcal{E}^{AV} . Note that the regions of high attention are not only limited to the direct sound sources (*e.g.*, regions in and around faces of active speakers across examples), but also include large sound-reflecting objects (*e.g.*, the flat surface of the table on the left; the walls on the left and in the middle; the floor on the right, etc.) that determine how sound spatializes through early reflections, late reverberations, etc. Interestingly, our model also attends to multiple people if they are speaking at the same time (see left), thereby facilitating the detection of multiple active speakers. See Supp. for additional visualizations showing how, depending on the scene's spatial layout, our model uses one audio channel more than the other to attend to important image locations.

In Fig. 5, we qualitatively show our model's success cases. On ASD (Fig. 5a), our model can tackle drastic camera movements, multiple active speakers, and partially visible faces. On denoising (Fig. 5b), our model is able to remove interferences from distractor sources, and make predictions that closely match the ground truth in spectrogram structure.

We also observe some limitations. Our model's performance on ASD declines when there are drastic movements of the camera wearer, or there is a high overlap in speech from different conversation participants. On denoising, our model struggles when the noisy audio is semantically and acoustically similar to the target, or when it cannot extract spatial cues due to occlusions or out-of-view speakers. Refer

to our Supp. video for both success and failure cases.

5. Conclusion

We introduce a novel self-supervised approach for learning audio-visual representations in social egocentric videos via spatial correspondence between the video and its binaural audio. Through extensive evaluation, we show that our learned features are strong and generic enough to improve over multiple backbone methods on multiple downstream tasks. In future work, we will explore how the learned spatial audio-visual cues may reveal the social attention between speakers.

Acknowledgements: UT Austin is supported in part by NSF CCRI and the IFML NSF AI Institute. KG is paid as a research scientist by Meta, and SM is a visiting researcher at Meta.

References

- Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. The conversation: Deep audio-visual speech enhancement. arXiv preprint arXiv:1804.04121, 2018.
- [2] Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised learning of audio-visual objects from video. In Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16, pages 208–224. Springer, 2020.
- [3] Juan Leon Alcazar, Fabian Caba Heilbron, Long Mai, Federico Perazzi, Joon-Young Lee, Pablo Arbeláez, and Bernard Ghanem. Active speakers in context. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 12462–12471, 2020. 2, 3, 6
- [4] Relja Arandjelović and Andrew Zisserman. Look, listen and learn. 2017 IEEE International Conference on Computer Vision (ICCV), pages 609–617, 2017. 2
- [5] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. arXiv preprint arXiv:1607.06450, 2016.
- [6] Alan Baade, Puyuan Peng, and David F. Harwath. Mae-ast: Masked autoencoding audio spectrogram transformer. In *Interspeech*, 2022. 2
- [7] Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicenc Amengual Gari, Ziad Al-Halah, Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman. Soundspaces: Audiovisual navigation in 3d environments. In Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16, pages 17–36. Springer, 2020. 2
- [8] Changan Chen, Ruohan Gao, Paul Calamia, and Kristen Grauman. Visual acoustic matching. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 18858–18868, 2022. 2
- [9] Changan Chen, Carl Schissler, Sanchit Garg, Philip Kobernik, Alexander Clegg, Paul Calamia, Dhruv Batra, Philip Robinson, and Kristen Grauman. Soundspaces 2.0: A simulation platform for visual-acoustic learning. Advances in Neural Information Processing Systems, 35:8896–8911, 2022. 2

- [10] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision (IJCV)*, 130:33–55, 2022. 1, 5
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition, pages 248–255, 2009.
- [12] Jacob Donley, Vladimir Tourbabin, Jung-Suk Lee, Mark Broyles, Hao Jiang, Jie Shen, Maja Pantic, Vamsi Krishna Ithapu, and Ravish Mehra. Easycom: An augmented reality dataset to support algorithms for easy communication in noisy environments. 2021. 1, 2, 5, 6, 13, 15, 16, 17
- [13] Ngoc QK Duong, Emmanuel Vincent, and Rémi Gribonval. Under-determined reverberant audio source separation using a full-rank spatial covariance model. *IEEE Transactions on Audio, Speech, and Language Processing*, 18(7):1830–1840, 2010. 3
- [14] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation. arXiv preprint arXiv:1804.03619, 2018. 3
- [15] Mark Everingham, Josef Sivic, and Andrew Zisserman. Hello! my name is... buffy" – automatic naming of characters in tv video. In *British Machine Vision Conference*, 2006. 2
- [16] Christoph Feichtenhofer, Yanghao Li, Kaiming He, et al. Masked autoencoders as spatiotemporal learners. *Advances in neural information processing systems*, 35:35946–35958, 2022. 2, 4, 16
- [17] John W Fisher III, Trevor Darrell, William Freeman, and Paul Viola. Learning joint statistical models for audio-visual fusion and segregation. In *Advances in Neural Information Processing Systems*. MIT Press, 2000. 3
- [18] Chuang Gan, Hang Zhao, Peihao Chen, David Cox, and Antonio Torralba. Self-supervised moving vehicle tracking with stereo sound. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7053–7062, 2019. 2, 7
- [19] Ruohan Gao and Kristen Grauman. 2.5D visual sound. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 324–333, 2019. 2, 5, 6, 7, 8, 13, 14, 16, 17
- [20] Ruohan Gao and Kristen Grauman. Co-separating sounds of visual objects. In *Proceedings of the IEEE/CVF International* Conference on Computer Vision, pages 3879–3888, 2019. 3
- [21] Ruohan Gao and Kristen Grauman. Visualvoice: Audiovisual speech separation with cross-modal consistency. *arXiv* preprint arXiv:2101.03149, 2021.
- [22] Ruohan Gao, Rogerio Feris, and Kristen Grauman. Learning to separate object sounds by watching unlabeled video. In Proceedings of the European Conference on Computer Vision (ECCV), pages 35–53, 2018. 3

- [23] Ruohan Gao, Changan Chen, Ziad Al-Halah, Carl Schissler, and Kristen Grauman. Visualechoes: Spatial image representation learning through echolocation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*, pages 658–676. Springer, 2020. 2
- [24] Mariana-Iuliana Georgescu, Eduardo Fonseca, Radu Tudor Ionescu, Mario Lucic, Cordelia Schmid, and Anurag Arnab. Audiovisual masked autoencoders. arXiv preprint arXiv:2212.05922, 2022. 2, 4, 5, 6, 7, 8, 14, 16
- [25] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Inter*national Conference on Artificial Intelligence and Statistics, 2010. 16
- [26] Yuan Gong, Andrew Rouditchenko, Alexander H. Liu, David Harwath, Leonid Karlinsky, Hilde Kuehne, and James R. Glass. Contrastive audio-visual masked autoencoder. In *The Eleventh International Conference on Learning Representations*, 2023. 2, 4, 8
- [27] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 18995–19012, 2022. 1, 5
- [28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 16
- [29] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 16000– 16009, 2022. 2, 4
- [30] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv: Learning*, 2016. 16
- [31] John Hershey and Javier Movellan. Audio vision: Using audio-visual synchrony to locate sounds. In *Advances in Neural Information Processing Systems*. MIT Press, 1999. 2
- [32] P. Huang, M. Kim, M. Hasegawa-Johnson, and P. Smaragdis. Deep learning for monaural speech separation. In 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1562–1566, 2014. 3
- [33] Po-Yao Huang, Vasu Sharma, Hu Xu, Chaitanya Ryali, Haoqi Fan, Yanghao Li, Shang-Wen Li, Gargi Ghosh, Jitendra Malik, and Christoph Feichtenhofer. Mavil: Masked audio-video learners. *arXiv preprint arXiv:2212.08071*, 2022. 2, 4, 8
- [34] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. *Advances in Neural Information Processing Systems*, 35:28708–28720, 2022. 2, 3, 15
- [35] Jaesung Huh, Jacob Chalk, Evangelos Kazakos, Dima Damen, and Andrew Zisserman. Epic-sounds: A large-scale dataset of actions that sound. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023. 1

- [36] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 448–456, Lille, France, 2015. PMLR. 16
- [37] Hao Jiang, Calvin Murdock, and Vamsi Krishna Ithapu. Egocentric deep multi-channel audio-visual active speaker localization. In *Proceedings of the IEEE/CVF Conference on Com*puter Vision and Pattern Recognition, pages 10544–10552, 2022. 2, 3
- [38] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. 16, 17
- [39] Okan Köpüklü, Maja Taseska, and Gerhard Rigoll. How to design a three-stage architecture for audio-visual active speaker detection in the wild. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1193–1203, 2021.
- [40] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. Advances in Neural Information Processing Systems, 31, 2018.
- [41] Juan Le'on-Alc'azar, Fabian Caba Heilbron, Ali K. Thabet, and Bernard Ghanem. Maas: Multi-modal assignation for active speaker detection. 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pages 265–274, 2021.
- [42] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017. 5
- [43] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 5
- [44] Sagnik Majumder and Kristen Grauman. Active audio-visual separation of dynamic sound sources. In Computer Vision– ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIX, pages 551– 569. Springer, 2022. 2
- [45] Sagnik Majumder, Ziad Al-Halah, and Kristen Grauman. Move2hear: Active audio-visual source separation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 275–285, 2021.
- [46] Sagnik Majumder, Changan Chen, Ziad Al-Halah, and Kristen Grauman. Few-shot audio-visual learning of environment acoustics. In Advances in Neural Information Processing Systems, 2022.
- [47] Sagnik Majumder, Hao Jiang, Pierre Moulon, Ethan Henderson, Paul Calamia, Kristen Grauman, and Vamsi Krishna Ithapu. Chat2map: Efficient scene mapping from multi-ego conversations. arXiv preprint arXiv:2301.02184, 2023. 2
- [48] Kyle Min, Sourya Roy, Subarna Tripathi, Tanaya Guha, and Somdeb Majumdar. Learning long-term spatial-temporal graphs for active speaker detection. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*, pages 371–387. Springer, 2022. 3, 5, 6, 14, 15, 16, 17, 18

- [49] Shentong Mo and Pedro Morgado. Localizing visual sounds the easy way. In Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII, pages 218–234. Springer, 2022. 2
- [50] Pedro Morgado, Nuno Nvasconcelos, Timothy Langlois, and Oliver Wang. Self-supervised generation of spatial audio for 360°video. In Advances in Neural Information Processing Systems. Curran Associates, Inc., 2018. 2
- [51] Pedro Morgado, Yi Li, and Nuno Nvasconcelos. Selfsupervised generation of spatial audio for 360° video. Advances in Neural Information Processing Systems, 33:4733– 4744, 2020. 2
- [52] Giovanni Morrone, Daniel Michelsanti, Zheng-Hua Tan, and Jesper Jensen. Audio-visual speech inpainting with deep learning. In *ICASSP*, 2021. 2
- [53] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML* 2010, pages 807–814, 2010. 16
- [54] Kazuhiro Nakadai, Ken-ichi Hidai, Hiroshi G Okuno, and Hiroaki Kitano. Real-time speaker localization and speech separation by audio-visual integration. In *Proceedings 2002 IEEE International Conference on Robotics and Automation* (Cat. No. 02CH37292), pages 1043–1049. IEEE, 2002. 3
- [55] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and A. Ng. Multimodal deep learning. In International Conference on Machine Learning, 2011. 2
- [56] Curtis Northcutt, Shengxin Zha, Steven Lovegrove, and Richard Newcombe. Egocom: A multi-person multi-modal egocentric communications dataset. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2020. 1, 2, 5, 6, 13, 14, 15, 16, 17
- [57] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *Proceedings of the European Conference on Computer Vision* (ECCV), pages 631–648, 2018. 2, 3
- [58] Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H Adelson, and William T Freeman. Visually indicated sounds. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2405–2413, 2016. 2
- [59] Andrew Owens, Jiajun Wu, Josh H. McDermott, William T. Freeman, and Antonio Torralba. Ambient sound provides supervision for visual learning. In *European Conference on Computer Vision*, 2016.
- [60] S. Parekh, S. Essid, A. Ozerov, N. Q. K. Duong, P. Pérez, and G. Richard. Motion informed audio source separation. In 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 6–10, 2017. 3
- [61] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely. The kaldi speech recognition toolkit. In *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding*. IEEE Signal Processing Society, 2011. IEEE Catalog No.: CFP11SRW-USB. 16
- [62] Jie Pu, Yannis Panagakis, Stavros Petridis, and Maja Pantic. Audio-visual object localization and separation using low-rank and sparsity. In 2017 IEEE International Conference

- on Acoustics, Speech and Signal Processing (ICASSP), pages 2901–2905, IEEE, 2017. 3
- [63] Senthil Purushwalkam, Sebastia Vicenc Amengual Gari, Vamsi Krishna Ithapu, Carl Schissler, Philip Robinson, Abhinav Gupta, and Kristen Grauman. Audio-visual floorplan reconstruction. In *Proceedings of the IEEE/CVF Interna*tional Conference on Computer Vision, pages 1183–1192, 2021. 2
- [64] Kranthi Kumar Rachavarapu, Aakanksha, Vignesh Sundaresha, and A. N. Rajagopalan. Localize to binauralize: Audio spatialization from visual sound source localization. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 1930–1939, 2021. 2
- [65] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. Unet: Convolutional networks for biomedical image segmentation. In Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18, pages 234–241. Springer, 2015. 7
- [66] Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R. Hershey. Sdr – half-baked or well done? ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 626–630, 2018. 8
- [67] Fiona Ryan, Hao Jiang, Abhinav Shukla, James M Rehg, and Vamsi Krishna Ithapu. Egocentric auditory attention localization in conversations. In *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, pages 14663–14674, 2023. 3
- [68] Tomoya Sato, Yusuke Sugano, and Yoichi Sato. Self-supervised learning for audio-visual relationships of videos with stereo sounds. *IEEE Access*, 10:94273–94284, 2022. 2, 5, 6, 7, 8, 14, 16
- [69] Rahul Sharma, Krishna Somandepalli, and Shrikanth S. Narayanan. Crossmodal learning for audio-visual speech event localization. ArXiv, abs/2003.04358, 2020. 2
- [70] Gunnar A. Sigurdsson, Abhinav Gupta, Cordelia Schmid, Ali Farhadi, and Karteek Alahari. Charades-ego: A largescale dataset of paired third and first person videos. *CoRR*, abs/1804.09626, 2018.
- [71] Paris Smaragdis, Bhiksha Raj, and Madhusudana Shashanka. Supervised and semi-supervised separation of sounds from single-channel mixtures. In *Independent Component Analysis* and Signal Separation, pages 414–421, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg. 3
- [72] Martin Spiertz and Volker Gnann. Source-filter based clustering for monaural blind source separation. In in Proceedings of International Conference on Digital Audio Effects DAFx'09, 2009. 3
- [73] Ruijie Tao, Zexu Pan, Rohan Kumar Das, Xinyuan Qian, Mike Zheng Shou, and Haizhou Li. Is someone speaking? exploring long-term temporal features for audio-visual active speaker detection. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 3927–3935, 2021. 2, 3, 5, 6, 14, 15, 16, 17, 18
- [74] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Video-MAE: Masked autoencoders are data-efficient learners for

- self-supervised video pre-training. In Advances in Neural Information Processing Systems, 2022. 2
- [75] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017. 4, 13
- [76] Tuomas Virtanen. Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria. *IEEE Trans. On Audio, Speech and Lang. Processing*, 2007. 3
- [77] Shanshan Wang, Archontis Politis, Annamaria Mesaros, and Tuomas Virtanen. Self-supervised learning of audio representations from audio-visual data using spatial alignment. *IEEE Journal of Selected Topics in Signal Processing*, 16(6): 1467–1479, 2022. 2
- [78] Karren Yang, Bryan Russell, and Justin Salamon. Telling left from right: Learning spatial correspondence of sight and sound. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9932–9941, 2020. 2, 5, 6, 7, 14, 16, 17, 18
- [79] Karren Yang, Michael Firman, Eric Brachmann, and Clément Godard. Camera pose estimation and localization with active audio sensing. In Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII, pages 271–291. Springer, 2022. 2
- [80] Yuanhang Zhang, Susan Liang, Shuang Yang, Xiao Liu, Zhongqin Wu, Shiguang Shan, and Xilin Chen. Unicon: Unified context network for robust active speaker detection. In Proceedings of the 29th ACM International Conference on Multimedia, pages 3964–3972, 2021. 2
- [81] Hang Zhou, Ziwei Liu, Xudong Xu, Ping Luo, and Xiaogang Wang. Vision-infused deep audio inpainting. In *ICCV*, 2019.
- [82] Özgür Yılmaz and Scott Rickard. Blind separation of speech mixtures via time-frequency masking. In IEEE TRANSAC-TIONS ON SIGNAL PROCESSING (2002) SUBMITTED, 2004. 3

6. Supplementary Material

In this supplementary material, we provide additional details about:

- Video (with audio) for qualitative illustration of our pretext task and qualitative evaluation of our model on the downstream tasks (Sec. 6.1), as noted in Sec. 4.3 in main.
- Spatial audio denoising results with varying noise levels (Sec. 6.2), as mentioned in Sec. ?? in main
- Evaluation of the impact of the channel masking probability r (from Sec. 3.3 and 3.4 in main) in our audio masking protocol (Sec. 6.3)
- Analysis of the effect of alternate audio masking choices (Sec. 6.4), as referenced in Sec. 4.2 in main
- Study of the effect of multi-level positional embeddings (Sec. 7), as noted in Sec. 4.2 in main
- Analysis of the effect of task-specific backbones (Sec. 8), as mentioned in Sec. 4.2 in main
- Comparison of model capacity and computational cost among different pretraining methods (Sec. 8.1), as noted in Sec. 4.2 in main
- Qualitative analysis of the visual attention maps for left and right audio channels separately (Sec. 8.2), as referenced in Sec. 4.3 in main
- Analysis of the impact of our finetuning strategy (Sec. 8.3), as noted after Sec. 4.2 in main
- Evaluation of the impact of our model parameter initialization on the downstream performance (Sec. 8.4)
- Additional dataset details (Sec. 8.5), as mentioned in Sec. 4 in main
- Additional model architecture and hyperparameter details for both self-supervised pretraining and downstream training (Sec. 8.6), as referenced in Sec. 3.4 in main

6.1. Supplementary video

The supplementary video provides a qualitative illustration of our pretraining task for learning spatial features from audio-visual correspondence in egocentric videos, and our proposed approach. Moreover, we provide video samples from the both EgoCom [56] and EasyCom [12] datasets to illustrate the unique challenges posed by the egocentric videos. Additionally, we demonstrate our model's prediction quality for both active speaker detection and spatial audio denoising, and analyze common failure models for our model on both tasks. The video is available on http://vision.cs.utexas.edu/projects/ego_av_corr.

6.2. Denoising with varying noise levels

In table 2 in main, we evaluated denoising with 0 dB noise. Here, we analyze the effect of varying the noise level. Table 3 reports the results with 2.5 dB and 5 dB noise. We observe general similarity in performance trends across all noise levels. Whereas our model outperforms the baselines in the high-noise settings (0 and 2.5 dB), using 2.5D-VS [19]++ improves the separation quality for 5 dB, underlining that our features are especially important for tackling the more challenging high-noise settings.

6.3. Channel masking probability r

Here, we analyze the effect of the channel masking probability r in our audio masking protocol (Sec. 3.3 in main) on the downstream

task performance. Table 4 reports the active speaker detection (ASD) results on the more challenging EgoCom [56] dataset, and table 5 reports the denoising results for different noise levels. We notice that the performance on both ASD and denoising, especially at the higher noise levels, declines upon increasing or decreasing the value of r from our choice of 20 % based on the downstream validation performance (Sec. 3.4 in main), which helps our model achieve a fine balance between the two complementary strategies of masking a complete channel and randomly masking audio tokens. Whereas randomly masking a channel of the binaural audio entails solving the more under-constrained and consequently complex binauralization task, thereby helping our model learn stronger spatial associations between vision and audio, randomly masking audio tokens helps with improving training stability.

6.4. Alternate audio masking choices

Here, we evaluate alternate masking choices, namely time, frequency, and time-frequency masking, in place of randomly masking audio patches as part of our proposed masking strategy, in table 6 and 7. Our model outperforms the versions with these alternatives, showing that random patch masking when combined with channel dropping enables learning more useful features in our setup. This happens possibly because in random patch masking, dropping a full frequency band or time segment is highly improbable thereby allowing our model to extract useful information from the unmasked regions of all frequency bands and time segments of the audio spectrograms.

7. Multi-level positional embeddings

Here, we evaluate the impact of our multi-level positional embeddings by comparing our model with the ablation where positional embeddings are used only at the input level. See table 8 for results on ASD and table 9 for results on denoising with 0 dB noise. Our model improves over the ablation on both tasks, showing that using multi-level positional embeddings is crucial for remembering the spatial layout of the tokens at different stages in the model.

8. Task-specific backbones

Here, we study the impact of using task-specific backbones on our model performance by evaluating two baselines, with the same architecture but without task-specific backbones (Ours w/o B)—one is learned from scratch and another is pretrained. See table 10 for results on ASD and table 11 for results on denoising with 0 dB noise. Our pretraining scheme leads to better performance than a from-scratch model even w/o B (table 10 and table 11 top), and we get the best results when we combine our features with task-specific backbones. This shows that while our audio-visual features provide important spatial cues to downstream models, they are not intended to replace the face-specific features used in ASD or the mixed audio features used in denoising.

8.1. Pretraining model capacity and computational cost

Here, we report the model capacity (parameter count) and GFLOPs of all pretraining methods in Table 12. Note that both the parameter count and GFLOPs of all transformer [75]-based methods

	2.5 d	lB	5 dl	В
Model	SI-SDRi \uparrow	$STFT\downarrow$	SI-SDRi ↑	$STFT\downarrow$
No pretraining				
U-Net w/o vision	1.91	5.32	2.02	3.04
U-Net	2.04	4.72	2.05	2.85
U-Net w/ ImageNet features	2.04	4.66	2.24	2.74
Alternate pretraining methods			ı	
U-Net w/ TLR [78] features	1.70	5.40	2.00	2.77
U-Net w/ 2.5D-VS [19] features	1.81	4.81	2.22	2.62
U-Net w/ 2.5D-VS [19]++ features	2.65	4.31	2.79	2.48
U-Net w/ SSR [68]++ features	2.25	4.63	2.21	2.80
U-Net w/ AV-MAE [24] features	2.46	4.60	2.14	2.93
Ours	2.72	4.22	2.46	2.70
Ours w/o pretraining	2.30	4.54	2.15	2.83
Ours w/ pretraining using monaural audio	2.58	4.38	2.31	2.81

Table 3. Audio denoising with U-Net [78] backbone for varying noise levels. All STFT distance measures use base 10^{-3} .

	TalkN	et [73]	SPEL	L [48]
r(%)	Val	Test	Val	Test
0	67.9	62.9		65.3
20 (Ours)	68.7	63.9	68.4	65.6
50	64.5	63.1	67.5	65.2
80	64.4	61.8	64.7	60.1
100	67.9	63.4	66.1	65.1

Table 4. Effect of r on the mean average precision (%) of our model for active speaker detection with two different backbones (TalkNet and SPELL).

(2.5DVS [19]++, AV-MAE [24] and SSR [68]++) are comparable to those of our model³, re-emphasizing that our improvements in performance on the downstream tasks are solely attributable to our better model design.

8.2. Visual attention maps per audio channel

In Fig. 6, we show the attention maps of our model (similar to Fig. 4 in main) separately for the left and right channels, on the more challenging EgoCom [56] dataset. We notice that our model uses the left channel to focus more on areas to the left of scene image and vice-versa, indicating that our model can reason about the spatial properties of the scene using both audio and vision.

Further, to better portray the larger trend, we measure the percentage of cases in our test set where the left audio channel attends more towards patches on the left side of the scene image, and the right channel attends more towards patches on the right. This measure comes out to be 62.4% for the left channel, and 57.2% for the right channel, showing that our model uses the left channel to focus more on areas to the left of the scene image and vice-versa, across the whole test set.

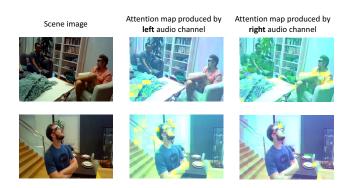


Figure 6. Heat maps for left and right audio channels, similar to Fig. 4 in main. Interestingly, our model uses the left channel to focus more on areas to the left of scene image and vice-versa.

8.3. Finetuning strategy

We mask audio tokens during finetuning primarily to reduce the computation overhead. Since our pretraining also involves token masking, our model can learn strong audio-visual features even when all audio tokens are not available during finetuning. However, to quantatively evaluate the effect of our finetuning strategy, we also finetune our model with all audio tokens and report the results in table 13 for ASD and table 14 for denoising. We don't see a significant change in performance when using all the tokens compared to masking when finetuning, but using all tokens is 1.7 times slower on average.

8.4. Model parameter initialization

To evaluate the effect of random parameter initialization on our model, we train our model on both tasks and datasets with 3 different random seeds. Across all runs, our standard errors are less than 0.01 on all metrics, showing that our model is robust to different random parameter initializations, and the improvements in performance are significantly larger than these small variations from randomness.

³The parameter count and GFLOPs of AV-MAE are a bit lower owing to its modality-inpainting architecture design, where the modality being inpainted is dropped from the input, leading to a slightly smaller model.

	0 dl		2.5 d		5 dI	
r(%)	SI-SDRi ↑	STFT ↓	SI-SDRi ↑	STFT ↓	SI-SDRi ↑	STFT ↓
0	2.17	6.60	2.57	4.38	2.85	2.35
20 (Ours)	2.20	6.51	2.72	4.22	2.46	2.70
50	1.92	7.19	2.30	4.60	2.09	2.80
80	1.82	7.55	1.98	5.05	1.68	3.30
100	2.11	6.60	2.65	4.31	2.79	2.48

Table 5. Effect of r on our model performance for audio denoising.

	TalkNet [73]		SPELL [48]]			
	Ego	Com	Easy	Com	Ego	Com	Easy	Com
Model	Val	Test	Val	Test	Val	Test	Val	Test
Ours w/ time masking							68.5	
Ours w/ frequency masking	64.1	62.1	59.2	70.9	67.6	63.2	68.7	69.4
Ours w/ time-frequency masking	65.4	63.1	56.3	63.3	67.5	65.1	68.6	69.1
Ours	68.7	63.9	60.5	71.8	68.4	65.6	68.9	70.2

Table 6. ASD with our model when pretrained with other audio masking choices [34].

	0 dB		2.5 d	lB	5 dB	
Model	SI-SDRi \uparrow	$STFT\downarrow$	SI-SDRi ↑	$STFT\downarrow$	SI-SDRi ↑	$STFT\downarrow$
Ours w/ time masking	1.82	7.41	1.98	4.88	2.07	2.80
Ours w/ frequency masking	2.05	7.04	2.33	4.85	2.25	2.79
Ours w/ time-frequency masking	1.91	7.12	2.14	5.15	1.81	3.13
Ours	2.20	6.51	2.72	4.22	2.46	2.70

Table 7. Denoising with our model when pretrained with other audio masking choices [34]. All STFT distance measures use base 10^{-3} .

	TalkNet		SPELL	
Model	EgoCom	EasyCom	EgoCom	EasyCom
Ours w/o multi-level PEs	59.2	70.2	60.4	65.6
Ours	63.9	71.8	65.6	70.2

Table 8. Effect of our multi-level positional embeddings on ASD.

Model	SI-SDRi ↑	STFT $(\times 10^{-3}) \downarrow$
Ours w/o multi-level PEs	1.30	7.88
Ours	2.20	6.51

Table 9. Effect of our multi-level positional embeddings on denoising with 0dB noise.

8.5. Dataset details

As discussed in main (Sec. 4), we use two public datasets containing egocentric videos with binaural audio, EgoCom [56] and EasyCom [12], for our experiments. For EgoCom, we follow the authors and split the data into train/val/test comprising 30.3/2.4/5.8 hours of data. For EasyCom, we randomly generate train/val/test splits with 4.5/0.38/0.39 hours of data, such that there is no overlap in conversation participants between any two splits. Next, we ex-

Model	Ego TalkNet		Easy TalkNet	
Ours w/o B (from-scratch)	61.1		62.0	
Ours w/o B (pretrained)	63.1		65	5.7
Ours	63.9	65.6	71.8	70.2

Table 10. Effect of task-specific backbones (denoted using 'B') on ASD.

Model	SI-SDRi ↑	STFT ($\times 10^{-3}$) \downarrow
Ours w/o B (from-scratch)	1.02	8.99
Ours w/o B (pretrained)	2.05	7.12
Ours	2.20	6.51

Table 11. Effect of task-specific backbones (denoted using 'B') on denoising with 0dB noise.

tract 1 second long clips from both datasets, where the video and binaural audio are sampled at 5 frames per second (fps) and 16 kHz, respectively. The frame resolution is 240×352 for EgoCom, and 198×352 for EasyCom. Furthermore, we choose audio channel 5 and 6 (corresponding to the in-ear microphones) as our binaural

	Model	parameter #	G	FLOPs
Model	ASD	Denoising	ASD	Denoising
2.5D-VS [19]	61.2	18.2	79.5	33.2
TLR [78]	57.5	18.1	75.9	33.7
2.5D-VS [19]++	180.9	75.3	174.0	90.2
AV-MAE [24]	178.6	74.1	171.6	87.5
SSR [68]++	180.9	75.3	174.0	90.2
Ours	180.9	75.3	174.0	90.2

Table 12. Model parameter count (in millions) and GFLOPs of different pretraining methods.

audio channels for EasyCom.

8.6. Model architecture and training details

In addition to the provided details in Sec. 3.4, 4.1 and 4.2 in main, we provide here extra model architecture and training details for both pretraining and finetuning on downstream tasks, for reproducibility. We perform all training using 8 NVIDIA Tesla V100 SXM2 GPUs. We will release all code and data.

8.6.1 Pretraining

We described our model architecture and pretraining details in Sec. 3.4 in main. Here, we provide additional details about our model's input preparation, architecture, parameter initialization, and training.

Input preparation. We sample the video clips at their original resolution, normalize them using the per-color means and standard deviations computed using ImageNet [28], and generate a total of 330 and 286 visual tokens for EgoCom and EasyCom, respectively, by splitting the clips into non-overalapping tubelets containing a sequence of 5 patches, where each patch is 16×16 in size (L193-5 in main). We represent the binaural audio as two-channel Kaldicompliant [61] spectrograms with 98 temporal windows and 128 Mel-frequency bins, which we compute by using the binaural audio normalized to [-1,1], a window length of 25 ms and a hop length of 10 ms. We normalize the spectrograms by computing the mean and standard deviation of the Mel-spectrograms generated from all audio clips in each dataset. We next generate 392 audio tokens per spectrogram channel by splitting it into non-overlapping patches of size 2×16 .

Architecture. All hidden layers in each transformer block [16] emit features that are four times as long as the embedding size for the block. We always use LayerNorm [5] after every output of a transformer block unless it's a direct input to another transformer block.

Parameter initialization. We use Xavier [25] uniform initialization for all network parameters. For the LayerNorm [5] layers, we initialize their weights to 1 and biases to 0. We use a truncated normal distribution with a standard deviation of 0.02 and a sampling range of [-2,2] to initialize the learnable modality and channel embedding tokens, and initialize the mask tokens from a normal distribution with a standard deviation of 0.02.

Training. We set the batch size to 104 and weight decay to 10^{-5} during pretraining.

8.6.2 Active speaker detection

In Sec. 4.1 in main, we outlined our feature fusion method for active speaker detection (ASD). Here, we provide additional architectural details for feature fusion, and also describe our finetuning process.

Pretrained feature fusion. Fig. 7 and 8 show our feature fusion methods for TalkNet [73] and SPELL [48] ASD backbones, respectively. The single-layer transformer decoder (Sec. 4.1 in main), which we use for fusing our pretrained features with the backbones (Sec. 4.1 in main), generates 128 and 512 dimensional embeddings for TalkNet and SPELL, respectively. Since SPELL doesn't train any audio-visual features when training its graph neural network (GNN), we first pretrain the the transformer decoder for SPELL by using it with the TalkNet backbone. Towards that goal, we feed the decoder features to a single linear layer that maps the 512 dimensional features to 128 dimensional features, and is followed by GELU [30] activations and LayerNorm [5], before fusing the 128 dimensional features with the TalkNet backbone. After pretraining, we append the 512 dimensional outputs of the decoder with the outputs of the two-stream audio-visual encoder (L405-8 in main) for training the GNN in SPELL.

Training. For TalkNet, we train using Adam [38] for 25 epochs optimizer with an initial learning rate (LR) of 10^{-4} for the backbone and 10^{-5} for the pretrained components, both of which we decay using a step LR scheduler by a factor of 0.95 after every epoch. We set the batch size to 400.

For SPELL, we first train the two-stream audio-visual encoder for feature extraction for 100 epochs using the cross entropy loss and Adam [38] with an initial learning rate of 5×10^{-4} , which we decay by 0.1 after every 40 epochs. We set the batch size to 320. For training the GNN of SPELL, we train for 70 epochs by using a batch size of 320 again and learning rate of 10^{-3} , while setting all other hyperparameters per the original paper.

8.6.3 Spatial audio denoising

Backbone architecture. Following [78], our U-Net backbone for spatial audio denoising (Sec. 4.2 in main) is an audio-visual model comprising an audio encoder, a visual encoder, and a decoder for predicting an estimate of the target audio. The audio encoder takes the log magnitude spectrogram of the mixed binaural audio as input, and uses a stack of 5 convolutional (conv.) layers to produce a multi-channel 2D audio feature map. Each conv. layer has a kernel size of 4, padding of 1, and stride of 2, and is followed by leaky ReLU [53] activations with a slope of 0.2 for negative inputs, and batch normalization [36]. The conv. layers have 64, 128, 256, 512 and 512 output channels, respectively. The visual encoder has a ResNet-18 [28] architecture that outputs a multi-channel 2D visual feature map without feeding it to the average pooling or any subsequent layers. We push the ResNet outputs through another conv. layer to match its height and width with the audio features. The conv layer has a kernel size of (1, 4), a padding of (0, 0) for EgoCom [56] and (1, 0) for EasyCom [12], and 128 output channels.

			SPELL [48]		
	EgoCom	EasyCom	EgoCom EasyC	Com	
Model	Val Test	Val Test	Val Test Val	Гest	
Ours w/ finetuning all audio tokens					
Ours	68.7 63.9	60.5 71.8	68.4 65.6 68.9 7	/0.2	

Table 13. ASD with our model when all tokens are used in downstream training.

	0 dB		2.5 dB		5 dB	
Model	SI-SDRi \uparrow	$STFT\downarrow$	SI-SDRi ↑	$STFT\downarrow$	SI-SDRi ↑	$STFT\downarrow$
Ours w/ finetuning all audio tokens	2.18	6.47	2.50	4.49	2.58	2.52
Ours	2.20	6.51	2.72	4.22	2.46	2.70

Table 14. Denoising with our model when all audio tokens are used in downstream training. All STFT distance measures use base 10^{-3} .

Further, we remove the last feature column from the output of the conv. layer for all channels for EasyCom. We concatenate the perframe features along the channel dimension and generate the visual features. We then concatenate the visual features with the audio features channel-wise, and feed the concatenated features to the audio decoder, which predicts an estimate of the ratio mask [19, 78] for the target audio magnitude spectrogram. The audio decoder first uses a stack of 5 transpose convolutional (conv.) layers, which are connected to the corresponding encoder layers through skip connections. The transpose conv. layers have a kernel size of 4, stride of 2, and a padding of (1, 1), except for the last layer, which has a padding of (2, 1). The transpose conv. layers have 1152, 1024, 512, 256 and 128 output channels, respectively. Next, the audio decoder feeds the output of the transpose conv. layers to a conv. layer with 2 input and output channels, and a kernel size of (2, 1) to emit the predicted ratio mask.

Input preparation. To transform the audio waveforms into magnitude spectrograms, we first normalize them to [-1, 1] and then compute the short-time Fourier transform with a window length of 128, hop length of 64, and 512 frequency bins.

Pretrained feature fusion. Fig. 9 shows our feature fusion method for spatial audio denoising. We reshape the visual features from the outputs of our audio-visual encoder \mathcal{E}^{AV} to form multichannel 2D visual feature maps (Sec. 4.2 in main), such that the 2D raster order of the features matches that of the tubelets in the video clip, and feed the reshaped features to a convolutional (conv.) layer with a kernel size of (3, 4), stride of (2, 3), padding of (1, 2) and (2, 2) for EgoCom [56] and EasyCom [12], respectively, and 128 input and 784 output channels. We similarly reshape the audio features, and feed them to another conv. layer with a kernel size of (1, 7), padding of 0, stride of (1, 6), and 128 input and 256 output channels. Both conv. layers are followed by leaky ReLU activations with a slope of 0.2 for the negative values, and batch normalization. Next, we concatenate the visual and audio features along the channel dimension, and further concatenate them with the audio encoder outputs channel-wise (Sec. 4.2 in main).

Training. We train using Adam [38] for 200 epochs optimizer with an learning rate (LR) of 5×10^{-4} . We set the batch size to 80

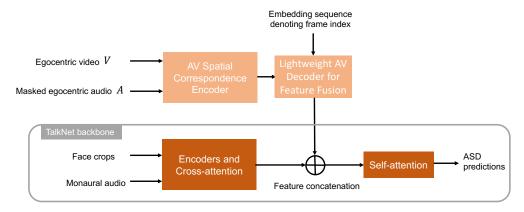


Figure 7. Method to fuse our pretrained features with TalkNet [73] for ASD.

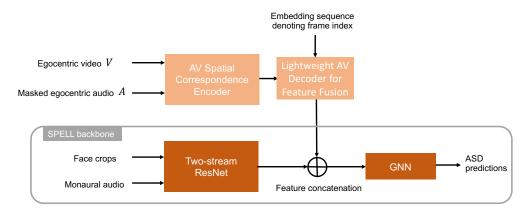


Figure 8. Method to fuse our pretrained features with SPELL [48] for ASD.

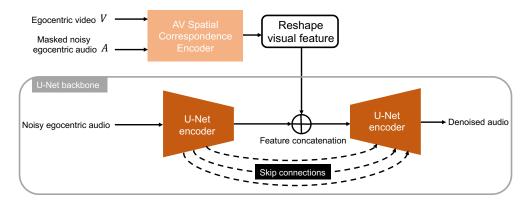


Figure 9. Method to fuse our pretrained features with U-Net [78] for spatial audio denoising.