

Static Visual Spatial Priors for DoA Estimation

Pawel Swietojanski*, *Member, IEEE*, and Ondrej Miksik*

Abstract—As we interact with the world, for example when we communicate with our colleagues in a large open space or meeting room, we continuously analyse the surrounding environment and, in particular, localise and recognise acoustic events. While we largely take such abilities for granted, they represent a challenging problem for current robots or smart voice assistants as they can be easily fooled by high degree of sound interference in acoustically complex environments. Preventing such failures when using solely audio data is challenging, if not impossible since the algorithms need to take into account wider context and often *understand* the scene on a *semantic level*. In this paper, we propose what to our knowledge is the first multi-modal *direction of arrival* (DoA) of sound, which uses *static visual spatial prior* providing an auxiliary information about the environment to suppress some of the false DoA detections. We validate our approach on a newly collected real-world dataset, and show that our approach consistently improves over classic DoA baselines.

Index Terms—Direction of Arrival, visual prior, voice assistants

I. INTRODUCTION

DIRECTION OF SOUND ARRIVAL (DoA) [1]–[3], or in general acoustic source localisation, played an important role in recent widespread adoption of voice assistants, in particular for devices that are more designed like robots with some degree of freedom in the environment. DoA is typically used for improving spatial scene understanding and as such forms basis for decision making, *e.g.* to take specific physical actions (rotate to the user, steer to an object of interest, *etc.*). Thus, it plays an important role in the overall user experience.

However, sound source localisation often becomes inherently ambiguous whenever the acoustic environment gets more complex. Consider, for instance, a single sound source and a device (microphone array) placed next to a glass wall; strong sound reflections from the wall often lead to unwanted interference that confuses DoA estimates. Rotating the robot to such location instead to the user can entirely break the user experience. But how can we *identify* the *true* sound source? Often, it is difficult or even impossible to disambiguate between the two using raw audio signal alone without any understanding of the wider context or having auxiliary prior knowledge about expected behaviour (*e.g.* estimate DoAs of people in the room, rather than DoAs of noises on the street).

Multi-sensory and multi-modal processing (sensors capture information from different origins) has been found to greatly improve performance in various machine learning perception tasks, in particular the ones with inherent ambiguity. Thus, it is viable to provide the DoA models with a spatial prior describing the vicinity of the device. With recent advances in computer vision, visual modality makes a good candidate

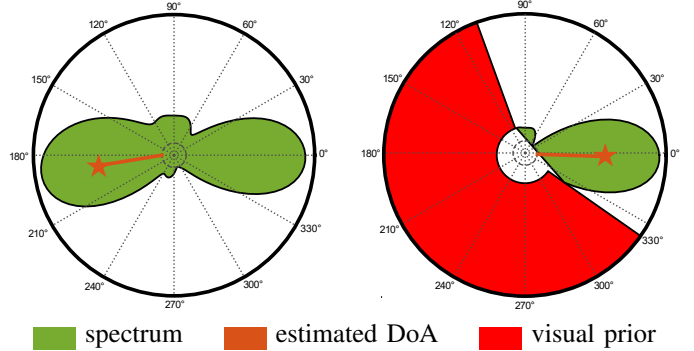


Fig. 1: (left) DoA is often confused in acoustically complex environments. (right) While identifying the true DoA using acoustic data alone may be challenging, injecting the algorithm with *static visual prior* provides an auxiliary information about the environment, that can be used to suppress false detections. Ground-truth prediction for this example is 0° .

for building such a prior since i) it can provide detailed information about the environment [4]–[6] and ii) majority of potential multimedia hardware are equipped with cameras, hence such algorithms come with no extra (hardware) cost.

In this paper, we propose what to our knowledge is the first multi-modal DoA which uses *static visual spatial priors* to suppress false DoA detections (*cf.* Fig. 1). Specifically, we assume a hardware platform equipped with a microphone array and a standard monocular camera mounted on a moving head which can rotate around its vertical axis. Using this platform, we capture 360° images which we use to predict *free space* and *obstacles*, *i.e.* to estimate plausible regions that can be occupied by people. To this end, we use semantic image segmentation as a proxy for sound source regions of interest. Then, we inject such static visual prior into classic audio-based DoA algorithms and show that it significantly reduces errors. We focus on injecting *static priors*, which is orthogonal to most multi-modal approaches that typically consider simultaneous and (often) synchronized data streams. The key difference is, that the static prior is estimated infrequently (*i.e.* calibration stage), when compared to the information throughput of the *primary* audio data stream (always on). This allows to use visual information with a constant compute cost, and in a broader set of situations – our approach is not sensitive to performance degradation in low-light conditions, nor assumes the users to be present within the camera field-of-view. Additionally, it does not require specialised hardware with active depth sensing (though in general could use it).

II. RELATED WORK

Direction of arrival. Estimating direction of arrival (DoA) of sound requires an access to a multi channel source of

* Equal contribution.

P. Swietojanski is with the School of Computer Science and Engineering, The University of New South Wales, Australia. (p.swietojanski@unsw.edu.au)

O. Miksik is with Emotech Labs, UK.

acoustic signal, typically captured by an M -element microphone array of (an ideally) known geometry. DoA can be then estimated directly by time-aligning the signals captured by pairs of microphones, using for example, Generic Cross Correlation with Phase Transform (GCC-PHAT) approach [7], [8]. Computing pair-wise DoAs in time delay domain, however, does not allow to fully utilise redundant information that results from combining several microphones in signal domain (and which is important in more challenging acoustic environments). Steered Response Power with PHAT (SRP-PHAT) [1] scans through candidate directions and seeks for the peaks in the value of GCC-PHAT averaged over all microphone pairs, this value is assumed as a desired DoA.

Signal sub-spaces methods estimate the so called *spatial covariance matrix* between multiple channels [2] and assume that the true and noise sources are independent and uncorrelated, thus occupy different subspaces. This has been generalised using maximum likelihood [9], [10], by exploiting statistical regularities *i.e. test of orthogonality of projected subspaces* (TOPS) [3] or weighting subspaces when deriving DoAs [11].

Robust estimation of DoA in reverberant environments remain an active research area, some directions include techniques for smoothing DoA trajectories [12]–[14] or the use of distributed sensors [15], [16]. The former approach is hard to apply in low-latency settings, while the latter is not always practical, though our method remains complementary to either. Recently, a trainable neural-net-based DoA was proposed [17].

Using audio-visual information to improve speaker localisation and tracking has also been studied, *e.g.* [18] or [19] used parallel audio-visual information for speaker tracking. It also remains an active research area in acoustic SLAM [20]. However, those approaches assume parallel data-streams and processing, whereas our work is concerned with independent asynchronous and non-real-time injection of visual data.

Visual free space prediction. Detecting *free space* from visual data has been widely explored in robotics. Free space is usually characterised by semantic classes corresponding rather to *stuff* than *objects* [21], hence this task is typically formulated as dense labelling (instead of using sparse representations such as bounding boxes) where the goal is to assign a (binary) label corresponding to free space or obstacles to each pixel. In robotics, semi-supervised methods used LIDARs or stereo-cameras to propose weak labels [22], [23]. Later, this was adapted to a single monocular camera [24]. Modern approaches go beyond binary labelling and rather segment the scene into semantically meaningful regions (*e.g.* floor, wall, TV, ...), typically using multi-class structured prediction frameworks [25]–[27]. The main advantages of multi-class segmentation is that i) it provides more information about the environment while still can be (indirectly) interpreted as binary (space / obstacle) labels and ii) the fact that multiple large-scale and annotated datasets are available [28]–[30].

In contrast to microphone arrays, cameras suffer from limited field-of-view (except for omnidirectional cameras). Thus, we need to process multiple images to “understand” the whole scene (*e.g.* a room), which unfortunately introduces two major difficulties: i) predictions from independently processed images are often inconsistent [31], [32] and ii) such data

has to be projected into a common representation (coordinate frame) shared with microphone array. While the former can be suppressed by explicit data association [33], [34], the latter is typically addressed by projecting the data on a common representation such as semantic maps [4], [5], [35], [36].

III. APPROACH

A. Incorporation of static visual priors into DoA

In this letter, we consider the DoA methods that rely on optimising some objective w.r.t. a set of candidate solutions. Let us denote a *discrete* set of potential angular directions φ of a circular array as $\mathcal{G} = \{\varphi_i \mid 0 \leq \varphi_i \leq 2\pi, i \in \mathbb{N}\}$, the expected DoA is then assumed to be at φ for which the cost function $f(\cdot)$ is maximised (or minimised). For example, for SRP-PHAT [1], one could search through all candidates in \mathcal{G} , and select the one with the tallest peak, *i.e.*

$$\varphi^* = \operatorname{argmax}_{\varphi \in \mathcal{G}} f(\varphi) \quad (1)$$

Similarly, we define static visual prior information as another set $\mathcal{P} = \{\varphi_i, Z_i \mid 0 \leq \varphi_i \leq 2\pi, Z_i \in \{0, 1\}, i \in \mathbb{N}\}$, in which Z maps a given angular direction φ_i into category of *free* ($Z = 1$) or *obstructed* ($Z = 0$) space. Then, one can alter the original search space \mathcal{G} with an additional information in \mathcal{P} , *i.e.* filter out directions that have been annotated as an obstructed space, and use \mathcal{G}' in (1) instead of \mathcal{G} :

$$\mathcal{G}' = \mathcal{G} \setminus \{ \cdot \} \cap \mathcal{P} \{ \cdot, Z = 1 \} \quad (2)$$

Note that (2) is independent of particular DoA model, thus the underlying characteristics behind back-end DoA estimator remain unchanged. Priors are estimated separately, hence can be easily swapped given the new \mathcal{P} is available (*e.g.* when device was moved to a different location). Also, visual prior linearly decreases computational cost for grid-based search algorithms, as one does not need to evaluate irrelevant directions when searching for the peaks. Finally, although we consider circular microphone arrays, our approach is applicable to arbitrary geometries assuming \mathcal{G} and \mathcal{P} were estimated correctly.

B. Building visual-based static spatial prior

We draw inspiration from biological systems which use cognitive maps of the environment for decision making [37]–[39]. This is known in robotics as semantic maps and is typically built using *semantic* Simultaneous Localisation and Mapping (SLAM) [6]. Such maps exist in various forms ranging from sparse symbols, through semi-dense and/or layered representations to fully-dense metric 3D maps.

We opt for *layered panoramic* representation since it does not require specific sensors such as Kinect/stereo-cameras or data-driven depth estimation and provides reliable predictions even at large distances (tens of meters); hence it is considerably simpler and faster to build than dense metric 3D maps. At the same time, DoA typically does not estimate proximity of the sound source, therefore omitting *depth* from spatial prior is not overly restrictive. The fact we use only passive monocular camera is of paramount importance for practical applications as it is the most widely used imaging sensor

already commonly available on majority of potential hardware platforms such as Amazon Echo Spot/Show2 [40] or mobile phones. Given a set of input images $\mathcal{I} = \{I_1, I_2, \dots, I_n\}$, our goal is to output a spatial prior map expressed as continuous angular representation $\mathcal{P} = \{\varphi, Z\}$ where φ is an angle and $Z \in \{0, 1\}$ denotes obstacles and free space (cf. §III-A). We assume that all input images \mathcal{I} were collected using a sensor spinning around its vertical axis (cf. §IV) so they (approximately) share the same camera center and hence induce a homography [41]. In other words, we can estimate a homography matrix $\mathbf{H}_{i,j}$ representing a 1-to-1 mapping (warp) between any pair of sufficiently overlapping images i and j . We use a standard approach of [42] which takes a set of input images \mathcal{I} and outputs a set of corresponding pairwise homography matrices $\mathcal{H} = \{\mathbf{H}_{1,2}, \mathbf{H}_{2,3}, \dots, \mathbf{H}_{n,1}\}$, using sparse feature matching, robust RANSAC-based homography fitting and bundle adjustment. Rotation matrices $\mathbf{R}_{i,j} \in SO(3)$ can be extracted from homography $\mathbf{H}_{i,j}$ using *e.g.* [43]. To avoid time-consuming building of image dataset (computer vision models typically require tens thousands of labelled images [28]–[30]), we predict *free space* and *obstacles* indirectly, using semantic segmentation. This allows to use existing large-scale datasets, such as ADE20K [30] (consisting of 20210 training images labelled with per-pixel ground-truth). To this end, we learn a nonlinear function $f_\theta : I \rightarrow S$ mapping image $I \in \mathbb{R}^{w \times h \times 3}$ to output $S \in \mathbb{R}^{w \times h \times L}$. Here, each pixel of output S represents an L -dimensional scores vector corresponding to L semantic classes and w and h are image dimensions. The multi-class predictor f is implemented as a convolutional neural network (CNN) and parametrized by θ . Finally, the output scores are mapped into binary labels (free space / obstacles) and projected into spatial prior map \mathcal{P} using a corresponding homography matrix $\mathbf{H}_{i,j}$.

Implementation details. For panorama stitching, we adapted a public implementation of [42] from OpenCV to support semantic images. For semantic segmentation, we used the ADE20K Dataset [30] to train the DilatedNet model [44] which drops `pool4` and `pool5` from fully convolutional VGG-16 network, and replaces the following convolutions with dilated (atrous) convolutions, and bilinear upsampling layer at the end. Finally, to convert the labelling into the angular representation of *free space* and *obstacles* $\mathcal{P} = \{\varphi, Z\}$, we check whether the fraction of pixels with semantic classes typically found in free areas (floor, ceiling, desk, chair, ...) within the current camera frustum is above per-class thresholds set using cross-validation on the `Dev` fold (cf. Appendix A-B).

IV. EXPERIMENTAL RESULTS

Experimental Protocol. To the best of our knowledge, there is no publicly available dataset consisting of 360° audio-visual data annotated with ground-truth directions of arrivals. Therefore, we collected around 2 hours of acoustic data with the corresponding visual snapshots representing office and home environments. `Dev` set comprises around 6 minutes of natural speech (collected in 2 office rooms, 15 sound source locations and 4 different microphone/camera placements). For the test set we collected two variants - `Test-Clean` comprising

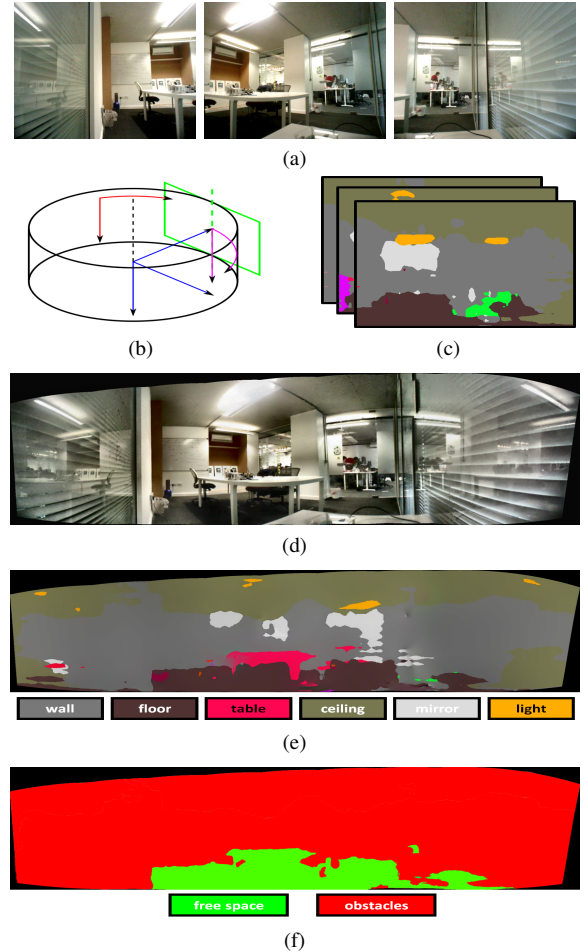


Fig. 2: Illustration of visual spatial prior construction. (a) we capture a set of images covering 360° around the device. These images are in parallel used to (b) estimate homography matrices and (c) to predict semantic segmentation. This is then combined to built *layered panoramic representation* (d) and (e), which is interpreted as *free space* and *obstacles* (f).

around 1h50min of re-recorded Librispeech data [45] as a primary sound source (2 rooms, 4 mic/camera placements, 5 sound-source and 2 noise-source locations). `Test-Noise` is a parallel variant (another 1h50min) collected in the identical conditions as `Test-Clean` but with present competing noise source (simulating a TV/radio) based on MUSAN data [46]. We also include a fully synthetic control benchmark referred to as `Syn`. Refer to the Appendix A-A for more details.

DoA predictions, unless stated otherwise, are computed on 256 ms long analysis windows. Signals are transformed to frequency domain with discrete Fourier transform, followed by an energy-based voice activity thresholding. DoA is calculated on frequency bins representing range from 500 Hz to 8 kHz. DoA detection is carried with 180 uniform bins, representing 2° resolution ($|\mathcal{G}| = 180$). We used SRP-PHAT [1], MUSIC [2] and TOPS [3] as back-end DoA algorithms. SRP-PHAT was found to give the best results, thus in the remainder all analyses are based on this method (full results for all three techniques are reported in the Appendix A-C). Experiments were carried out using an open source Pyromacoustics toolkit [47].

Results. Tab. I shows the main results for the three data

TABLE I: Average error rates and ± 5 deg bin accuracies (in square brackets) for the synthetic, dev and test sets obtained with SRP-PHAT algorithm with and without spatial priors.

Prior	Avg. Error (deg) [± 5 deg bin acc (%)]			
	Syn.	Dev	Test	
			Clean	Noise
None	3.2 [98.5]	25.1 [52.6]	46.2 [52.9]	73.7 [35.1]
Visual	N/A	14.5 [58.5]	31.1 [59.2]	52.7 [44.1]
Expert	1.6 [99.5]	11.6 [62.5]	25.3 [61.9]	39.3 [49.8]

folds, Syn., Dev and Test. We report both average errors, as well as average bin accuracies defined as a percentage of predictions falling into a bin of assumed width on either side of the ground truth DoA (± 5 degrees unless stated otherwise). DoAs with automatically extracted visual priors reduced average errors by an absolute 15.5° and increased bin accuracy by 16.2% relative (on Dev and Test sets, note CV priors are not available for Syn.). This effect was relatively stronger for Test-Noise condition, where automatically derived priors offered 25.9% rel. bin accuracy increase. Similar trends are observed with ground truth (expert) priors, where average degree errors were roughly halved in each of tested condition. Likewise, bin accuracies increased on average by 25.9% rel. and this effect was stronger in Test-Noise variant at 41.8% better vs. avg. 17.9% for Dev and Test-Clean. Note that those numbers concern difficult cases such as a device next to the wall, or in the corner. In case where device is further away from the walls, baseline errors are lower at 10% (*i.e.* as for M3 position in Fig. A1 in the Appendix A-A).

Fig. 3 offers more insight into other operating points on Dev. In particular, Fig. 3 (top) shows average errors and bin accuracies for analysis windows ranging from 32–512 ms. As expected, longer windows (more snapshots) offer lower errors and higher bin accuracy, though at the cost of latency (could be an issue for some applications like mapping acoustic events to spatial locations). Fig. 3 (left) shows how bin accuracies varies for different bin widths - this is of interest as even coarse prediction (*i.e.* ± 30 deg bin) is still acceptable, as it may put the sound source within camera’s field-of-view, which can be then used to fine-tune DoA [18]–[20]. In either scenario, trends are consistent and both expert and visual priors significantly outperform the baseline. Fig. 3 (right) shows how errors varies for different prior widths (*i.e.* corner, wall). Interestingly, tighter priors offer larger relative improvements and this effect increases for shorter analysis windows (*cf.* results in the Appendix A-C). Finally, Fig. 4 illustrates an example sequence of DoA prediction under all three settings.

V. DISCUSSION AND CONCLUSIONS

While our approach has demonstrated promising results, there are many possible extensions. For instance, one could improve static priors by exploiting richer semantic information, *i.e.* detect acoustically active noise sources such as TV or speakers and further constrain *free* space to regions that are more likely to be occupied by people vs. devices (for talker localisation). Similarly, such priors could be used to discover and set nulls for known sources of acoustic interference in the generic sidelobe canceller class of beamformers [48].

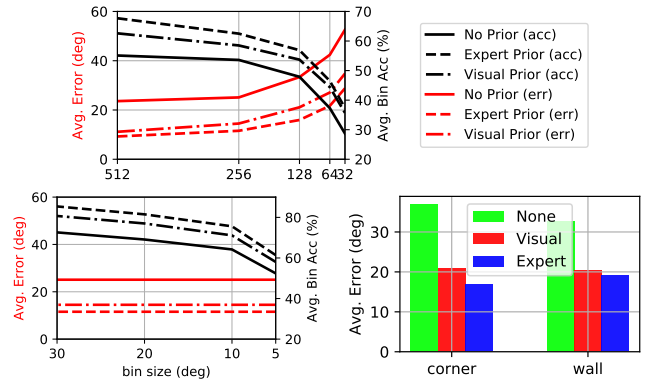


Fig. 3: Effects of priors on avg. errors (red) and bin acc. (black) on Dev data as a function of (top) window sizes (left) bin sizes and (right) prior widths (*i.e.* corner, wall).

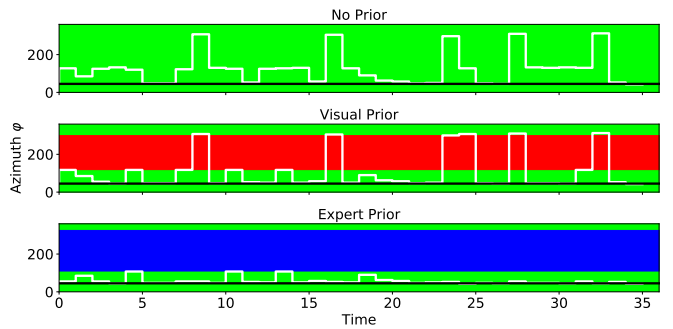


Fig. 4: Visual priors for free space (green) and obstacles (red, blue) define allowed search regions for DoA algorithms. Solid black line is the ground-truth (at 45°), white lines denote predicted DoAs in analysis windows. (Top) Baseline DoA predictions without spatial priors. (Middle) Estimated DoAs with automatically derived spatial prior using computer vision (obstacles between 120° and 300°). (Bottom) The expert derived ground-truth prior (obstacles between 110° and 320°). This scene corresponds to the one depicted in the Fig. 2 (d) in which the device “observes” the room from the corner.

Our layered panoramic representation is only an example of spatial prior map construction, however, it offers a number of interesting features as it is i) computationally and memory efficient, ii) continuous (*i.e.* no angular quantization, temporal consistency) and iii) scalable (arbitrary number of classes, depth, materials, surface normals, *etc.*). Experiments proved this is an efficient approach (*cf.* §IV), however, a variety of alternative approaches exists; for instance, one could learn a CNN to predict free space directly, instead of semantic segmentation (given appropriately annotated data). Another option might be using dense metric 3D representation, if more advanced sensors are available.

We have proposed the first multi-modal DoA, which uses static visual spatial prior to reduce potential false detections. We have validated our approach on a newly collected real-world dataset, and showed that our approach consistently improves over a wide range of DoA baselines using the ground-truth prior as obtained by an expert. Finally, we have demonstrated a real-world performance of our approach using a simple method for deriving spatial prior automatically.

REFERENCES

- [1] J. H. DiBiase, *A high-accuracy, low-latency technique for talker localization in reverberant environments using microphone arrays*, Ph.D. thesis, Brown University, Providence, RI, 2000.
- [2] R. Schmidt, “Multiple emitter location and signal parameter estimation,” *T-AP*, 1986.
- [3] Y. Yeo-Sun, L. M. Kaplan, and J. H. McClellan, “Tops: new doa estimator for wideband signals,” *T-SP*, 2006.
- [4] J. P. C. Valentin, S. Sengupta, J. Warrell, A. Shahrokni, and P. H. S. Torr, “Mesh based semantic modelling for indoor and outdoor scenes,” in *CVPR*, 2013.
- [5] V. Vineet, O. Miksik, M. Lidegaard, M. Nießner, S. Golodetz, V. A. Prisacariu, O. Köhler, D. W. Murray, S. Izadi, P. Perez, and P. H. S. Torr, “Incremental dense semantic stereo fusion for large-scale semantic scene reconstruction,” in *ICRA*, 2015.
- [6] A. J. Davison, “Futuremapping: The computational structure of spatial AI systems,” *CoRR*, vol. abs/1803.11288, 2018.
- [7] C. Knapp and G. Carter, “The generalized correlation method for estimation of time delay,” *T-ASSP*, 1976.
- [8] M. S. Brandstein and H. F. Silverman, “A robust method for speech signal time-delay estimation in reverberant rooms,” in *ICASSP*, 1997.
- [9] P. Stoica and K. C. Sharman, “Maximum likelihood methods for direction-of-arrival estimation,” *IEEE Trans. on Acoustics, Speech, and Signal Processing*, vol. 38, no. 7, pp. 1132–1143, July 1990.
- [10] M. Pesavento and A. B. Gershman, “Maximum-likelihood direction-of-arrival estimation in the presence of unknown nonuniform noise,” *IEEE Trans. on Signal Processing*, vol. 49, no. 7, pp. 1310–1324, July 2001.
- [11] E. D. di Claudio and R. Parisi, “Waves: weighted average of signal subspaces for robust wideband direction finding,” *T-SP*, 2001.
- [12] X. Anguera, C. Wooters, and J. Hernando, “Acoustic beamforming for speaker diarization of meetings,” *IEEE Trans. on Audio, Speech, and Language Processing*, vol. 15, no. 7, pp. 2011–2021, September 2007.
- [13] B. Rafaely and D. Kolossa, “Speaker localization in reverberant rooms based on direct path dominance test statistics,” in *Proc. IEEE ICASSP*, 2017, pp. 6120–6124.
- [14] C. Evers, B. Rafaely, and P. A. Naylor, “Speaker tracking in reverberant environments using multiple directions of arrival,” in *2017 Hands-free Speech Communications and Microphone Arrays (HSCMA)*, March 2017, pp. 91–95.
- [15] A. Canciani, F. Antonacci, A. Sarti, and S. Tubaro, “Acoustic source localization with distributed asynchronous microphone networks,” *IEEE Trans. on Audio, Speech, and Language Processing*, vol. 21, no. 2, pp. 439–443, Feb 2013.
- [16] C. Evers, E. A. P. Habets, S. Gannot, and P. A. Naylor, “Doa reliability for distributed acoustic tracking,” *IEEE Signal Processing Letters*, vol. 25, no. 9, pp. 1320–1324, Sep. 2018.
- [17] S. Adavanne, A. Politis, and T. Virtanen, “Direction of arrival estimation for multiple sound sources using convolutional recurrent neural network,” in *Proc. IEEE EUSIPCO*, IEEE, 2018, pp. 1462–1466.
- [18] N. Strobel, S. Spors, and R. Rabenstein, “Joint audio-video signal processing for object localization and tracking,” in *Microphone Arrays: Signal Processing Techniques and Applications*, 2001.
- [19] T. Gehrig, K. Nickel, H. K. Ekenel, U. Klee, and J. McDonough, “Kalman filters for audio-video source localization,” in *Workshop on Applications of Signal Processing to Audio and Acoustics*, 2005.
- [20] I. D. Gebru, C. Evers, P. A. Naylor, and R. Horaud, “Audio-visual tracking by density approximation in a sequential bayesian filtering framework,” in *2017 Hands-free Speech Communications and Microphone Arrays (HSCMA)*, March 2017, pp. 71–75.
- [21] E. H. Adelson, “On seeing stuff: the perception of materials by humans and machines,” 2001.
- [22] H. Dahlkamp, A. Kaehler, D. Stavens, S. Thrun, and G. R. Bradski, “Self-supervised monocular road detection in desert terrain,” in *Robotics: Science and Systems*, 2006.
- [23] R. Hadsell, P. Sermanet, J. Ben, A. Erkan, M. Scoffier, K. Kavukcuoglu, U. Muller, and Y. LeCun, “Learning long-range vision for autonomous off-road driving,” *IJRR*, 2009.
- [24] O. Miksik, P. Petyovsky, L. Zalud, and P. Jura, “Robust detection of shady and highlighted roads for monocular camera based navigation of ugv,” in *ICRA*, 2011.
- [25] L. Ladicky, C. Russell, P. Kohli, and P. H. Torr, “Associative hierarchical random fields,” *T-PAMI*, 2013.
- [26] LC Chen, A. Schwing, A. Yuille, and R. Urtasun, “Learning deep structured models,” in *ICML*, 2015.
- [27] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. S. Torr, “Conditional random fields as recurrent neural networks,” in *ICCV*, 2015.
- [28] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, “The cityscapes dataset for semantic urban scene understanding,” in *CVPR*, 2016.
- [29] T.-Y. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft COCO: common objects in context,” *CoRR*, vol. abs/1405.0312, 2014.
- [30] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, “Scene parsing through ade20k dataset,” in *CVPR*, 2017.
- [31] S. Sengupta, P. Sturgess, L. Ladicky, and P. H. S. Torr, “Automatic dense visual semantic mapping from street-level imagery,” in *IROS*, 2012.
- [32] V. Badrinarayanan et al., “Segnet: A deep convolutional encoder-decoder architecture for image segmentation,” *T-PAMI*, 2017.
- [33] O. Miksik, D. Munoz, J. A. Bagnell, and M. Hebert, “Efficient temporal consistency for streaming video scene analysis,” in *ICRA*, 2013.
- [34] A. Kundu, V. Vineet, and V. Koltun, “Feature space optimization for semantic video segmentation,” in *CVPR*, 2016.
- [35] S. Sengupta, E. Greveson, A. Shahrokni, and P. H. S. Torr, “Urban 3d semantic modelling using stereo vision,” *ICRA*, 2013.
- [36] A. Hermans, G. Floros, and B. Leibe, “Dense 3D Semantic Mapping of Indoor Scenes from RGB-D Images,” in *ICRA*, 2014.
- [37] B. M. Lake, T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman, “Building machines that learn and think like people,” *CoRR*, vol. abs/1604.00289, 2016.
- [38] R. J. Dolan and P. Dayan, “Goals and habits in the brain,” *Neuron*, 2013.
- [39] N. D. Daw, Y. Niv, and P. Dayan, “Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control,” *Nature Neuroscience*, 2005.
- [40] “Amazon echo spot/show2,” <https://www.amazon.com/All-new-Echo-Show-2nd-Gen/dp/B077SXWSRP>, Accessed: 2018-10-23.
- [41] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, 2004.
- [42] M. Brown and D. G. Lowe, “Automatic panoramic image stitching using invariant features,” *IJCV*, 2007.
- [43] E. Malis and M. Vargas, “Deeper understanding of the homography decomposition for vision-based control,” Research Report RR-6303, INRIA, 2007.
- [44] F. Yu and V. Koltun, “Multi-scale context aggregation by dilated convolutions,” in *ICLR*, 2016.
- [45] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: An asr corpus based on public domain audio books,” in *In Proc. IEEE ICASSP*, April 2015, pp. 5206–5210.
- [46] D. Snyder, G. Chen, and D. Povey, “MUSAN: A Music, Speech, and Noise Corpus,” 2015, arXiv:1510.08484v1.
- [47] R. Scheibler, E. Bezzam, and I. Dokmanić, “Pyroomacoustics: A python package for audio room simulation and array processing algorithms,” in *ICASSP*, 2018.
- [48] H. L. Van Trees, *Optimum array processing*, John Wiley & Sons, 2002.

APPENDIX A

A. Data Collection

We have built an experimental platform consisting of a fixed base and a moving head equipped with an RGB wide-angle camera, two 4-microphone arrays with circular geometries¹, and a DC motor allowing for continuous 360° rotation around its vertical axis. The reference position of the moving head (with camera) was calibrated with respect to microphone array to allow mapping from visual to audio data. Using this platform, we have collected realistic acoustic and visual data necessary to construct and test injection of static visual priors.

We split this data into development set *Dev* comprising natural speech with no competing sound sources, but in rooms characterised by high reverberation, with T_{60} being approximately 600ms and 250ms., and test sets referred to as *Test-Clean* and *Test-Noise*. Locations of sound and noise sources and mic/camera positions are shown in Fig. A1. Synthetic data *Syn.* is sampled from normal distribution. To simulate different impeding angles (6 in total) the source channel is shifted w.r.t. itself $M - 1$ times using delay filter banks corresponding to the 1st mic array geometry. SNR augmented data is obtained by adding Gaussian noise at desired SNR levels separately to each of the shifted channels.

B. Mapping semantic segmentation to binary labels

To convert multi-class semantic segmentation to binary labels, we learnt per-class thresholds on the *Dev* fold, for the following classes {*floor*, *desk*, *table*, *chair*, *tv*} and assigned full images as *free space* if the number of pixels with these classes is above threshold (recall the images are sampled with 10° angular resolution). Of course, these are not the only available classes in the ADE20K Dataset, however, in practice worked well on the *Dev* fold and generalised to the test sets, as is demonstrated by presented results. To convert recognised free space to angular representation, we simply read out images, which were assigned to *free space* label as the images are captured relatively densely, however, one could in practice decompose the homography matrix, as described in §III and achieve much finer resolution if desired.

C. More results

Table AI shows results for SRP-PHAT, MUSIC and TOPS on *Dev* set. All share the same automatically extracted visual and expert priors as well as identical pre-processing pipelines (the number of DFT points was optimised for each method).

Table AII offers similar set of results to Table I, but for 128ms analysis window. The overall findings are consistent, but here priors injection offer larger relative gains across all sets (this trend increases as analysis window gets smaller).

Finally, Fig. A2 shows impact of expert priors on synthetic dataset. We test to what extent priors help under different SNR regimes Fig. A2 (a), analysis window and bins sizes (b) and (c) as well as prior widths (d), *i.e.* 270° approx. corresponds to a corner while 180° to a wall case. Note, all errors on *Syn.* set are considerably lower when compared to realistic folds, thus gains are primarily visible in more

challenging operating conditions (low SNR, short windows), however, overall findings are in-line with *Dev* and *Test* sets.

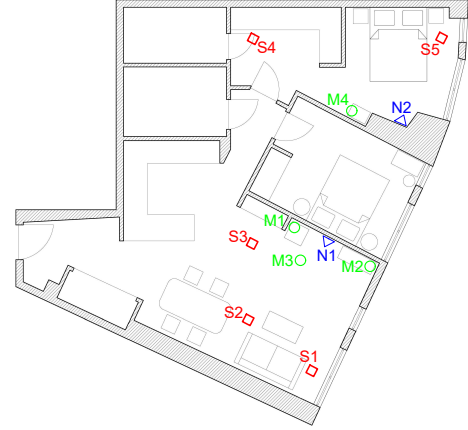


Fig. A1: Test set collection environment with annotated positions of speakers (S1-S5), microphones (M1-M4) and noise sources (N1-N2). Not to scale.

TABLE AI: Results on *Dev* set obtained with three DOA algorithms and considered priors for 128ms analysis windows.

Prior	Average Error (deg.) [± 5 deg. bin acc. (%)]		
	SRP-PHAT [1]	MUSIC [2]	TOPS [3]
None	33.4 [47.9]	56.9 [28.4]	38.5 [33.8]
Visual	21.1 [54.2]	33.7 [33.8]	18.2 [41.8]
Expert	15.9 [56.9]	28.0 [34.1]	13.8 [42.3]

TABLE AII: Results for 128ms long analysis window and SRP-PHAT algorithm with and without spatial priors.

Prior	Avg. Err (deg) [± 5 deg bin acc (%)]			
	Syn.	Dev	Test	
			Clean	Noise
None	6.3 [90.1]	33.4 [47.9]	55.3 [44.0]	77.8 [30.7]
Visual	N/A	21.1 [54.2]	36.4 [52.1]	54.7 [40.6]
Expert	1.8 [93.4]	15.9 [56.9]	28.4 [55.8]	41.4 [45.9]

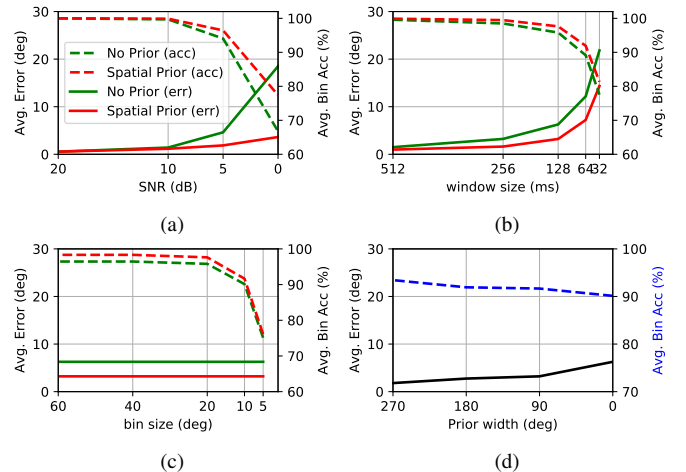


Fig. A2: Effects of prior on avg. errors (solid) and bin acc. (dashed) on *Syn.* data as a function of (a) SNRs (b) window sizes (c) bin sizes and (d) prior widths (*i.e.* corner, wall). Window 128ms, average over scores of all SNR levels.

¹The 1st array (used to collect *Dev*) had 28.3mm radius, the 2nd array (used to collect *Test*) was out-of-the-shelf XMOS VocalFusion XVF3100.