# Audio-visual target speaker enhancement on multi-talker environment using event-driven cameras

Ander Arriandiaga[1], Giovanni Morrone[2], Luca Pasa[3], Leonardo Badino[3], and Chiara Bartolozzi[1]

[1]iCub Facility, Istituto Italiano di Tecnologia, Genova, Italy
[2]Department of Engineering "Enzo Ferrari", University of Modena and Reggio Emilia, Modena, Italy
[3]CTNSC, Istituto Italiano di Tecnologia, Ferrara, Italy

*Abstract*—We propose a method to address audio-visual target speaker enhancement in multi-talker environments using event-driven cameras. State of the art audio-visual speech separation methods shows that crucial information is the movement of the facial landmarks related to speech production. However, all approaches proposed so far work offline, using frame-based video input, making it difficult to process an audio-visual signal with low latency, for online applications. In order to overcome this limitation, we propose the use of event-driven cameras and exploit compression, high temporal resolution and low latency, for low cost and low latency motion feature extraction, going towards online embedded audio-visual speech processing. We use the event-driven optical flow estimation of the facial landmarks as input to a stacked Bidirectional LSTM trained to predict an Ideal Amplitude Mask that is then used to filter the noisy audio, to obtain the audio signal of the target speaker. The presented approach performs almost on par with the frame-based approach, with very low latency and computational cost.

**Index Terms**:speech separation, event-driven camera, optical-flow, LSTM, deep learning

## I. INTRODUCTION

The ability to disentangle and correctly recognise the speech of a single speaker among other speakers (the well known cocktail party effect [1]) is paramount for effective speech interaction in unconstrained environments. As such, it is an extremely useful feature for any artificial device that relies on speech interaction such as robots and mobile devices. To this aim, it is crucial to devise efficient speaker enhancement techniques that rely on small datasets and low power sensing and computation. Humans solve this problem using complementary and redundant strategies such as physical sound source separation (thanks to stereo sound acquisition [2]) and using cues from lips motion [3].

Artificial systems use single-channel audio signals as input to Long-Short Memory Networks (LSTM) [4]–[7] or dilated convolutional layers [8] for speaker-independent enhancement. However, the number of speakers has to be known in advance, as well as the correspondence between the target speaker and the output clean speech. An alternative is to give as

input to the model speaker dependent target features [9], [10], using an LSTM-based speaker encoder to produce speaker-discriminative embeddings. However, this solution needs a reference utterance of the speaker and an additional trainable Deep Neural Network (DNN), making the speech separation performance conditioned on the performance of the speaker encoder network, and computationally heavy.

Inspired on the findings that viewing the target speaker's face improves the listener ability to track the speech [3], methods that combine visual cues and speech processing achieved remarkably good results. They were based on residual networks (ResNet [11]) pre-trained on a word-level lip-reading task [12], [13], or based on a pre-trained face recognition model, in combination with 15 dilated convolutional layers [14]. Such architectures, however, are computationally heavy and require heterogeneous and large audio-visual datasets for training. An approach that allows to use smaller datasets (such as the GRID dataset [15]) is to rely on pre-trained models, with the use of images and corresponding optical flow as inputs to a pre-trained dual tower ResNet extracting visual features [16].

If video features are extracted without using trainable methods, the neural networks are smaller and can be trained with smaller datasets without overfitting. Following this idea, [17] used face landmark movements as input visual features to a bidirectional LSTM that achieved good speaker-independent results on the GRID dataset. In this work, the use of landmark motion features rather than positional features turned out to be a key factor. Inspired by this finding, we propose to substitute the visual pipeline implemented with traditional frame-based sensors, face tracking and extraction of motion landmarks, with an equivalent pipeline, based on the use of a novel type of vision sensors – the event-driven cameras (EDC) – from which the extraction of motion is available at lower computational cost and latency. EDCs asynchronously measure the brightness change for each pixel, featuring a temporal resolution as high as 1 $\mu$s, extremely low latency and data compression (as only active pixels communicate data). With such an input, the audio-visual system can use the same temporal discretization of the auditory pipeline (10 ms), rather than
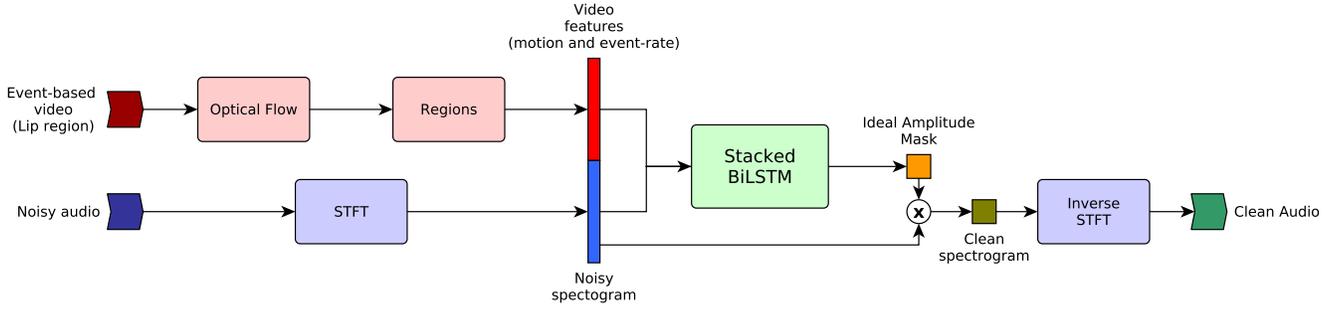
Fig. 1. Audio-Visual Speech Separation pipeline.

the one of the visual pipeline (30 fps is the standard frame-rate of traditional sensors). Event-driven vision sensors have been widely used with good results for object tracking [18], [19], detection [20] and segmentation [21], and for gesture recognition [22]. Recently, they have been applied in the context of speech processing: vision-only speech recognition (i.e., lip-reading) on GRID exploited EDCs as input to a Deep Neural Network architecture [23]; lip movements detected by an EDC were used to detect speech activity and enable an auditory-based voice activity detection [24], for embedded applications that require low computational cost. This work presents an audio-visual target speaker enhancement system on multi-talker environment using event-driven vision sensors that compute motion at lower latency and computational cost. Following [17], we propose a non trainable method to extract visual features combined with deep learning techniques. We use the GRID corpus in order to compare this approach with frame-based methods. To the best of our knowledge, this is the first work that presents an audio-visual target speaker enhancement system that uses event-driven cameras.

## II. METHODS

This work is based on [17] where audio and visual motion features were used as inputs of LSTM-based models to generate time-frequency masks. These masks are then applied to generate the clean spectrogram of the target speaker. Our contribution stands in the use of EDC for the acquisition of the visual signal and for the computation of the motion features of speech-production facial landmarks, based on the estimation of the normal optical flow from the events. Fig. 1 shows the block diagram of the whole system we propose for audio-visual target speaker enhancement: the computed visual motion features relative to the target speaker are concatenated with audio features to train a Recurrent Neural Network (RNN) that estimates a time-frequency mask that, multiplied by the noisy spectrogram, separates the clean speech produced by the target speaker.

### A. Event-driven motion features extraction

EDCs output asynchronous events whenever a pixel detects changes in log intensity larger than a threshold. Each event has an associated timestamp, $t$, pixel position, $<x, y>$, and

polarity (log intensity increase or decrease), $p$ [25]. The events are emitted with high temporal resolution and low latency, only when there is relative motion between the camera and the scene, increasing for fast moving objects and decreasing for less active scenarios, such as that of a target speaker talking in front of the camera. Fig. 2 shows the typical output in such scenario, where only the motion of the person and his/her mouth and eyes generate events, leading to a low amount of information to process and the possibility to have an always on front end visual acquisition and processing for audio-visual tasks. The different data structure and content from EDC require algorithms for the estimation of the optical flow, that can rely on the precise time of each event and the continuous observation of the events produced by contrast edges moving from one pixel to its neighbours. Even though the state of the art for EDC optical flow estimation is based on the use of deep learning [26], as the motion of lips and other facial landmarks are mostly perpendicular to the edge, we resorted to a temporally and computationally efficient algorithm for the estimation of the normal optical flow [27].

## III. EXPERIMENTAL SETUP

### A. Dataset

We focussed our analysis on a challenging and common scenario, where the quantity of available data and resources are limited, using the GRID dataset [15]. The dataset consists of 3 seconds long audio and video recordings of 34 speakers pronouncing 1000 sentences in front of a frame-based camera and microphone. The camera records data at 25 frames per second, while the microphone data is recorded at 50kHz. We used a subset of the GRID corpus, consisting of 200 sentences from 33 speakers (one was discarded because the videos were not available). To test speech enhancement of a target speaker, the audio signals of two different speakers were mixed so that for each speaker there are 600 mixed-audio sample recordings. From the total amount of samples, samples from 25 speakers were for training, from 4 speakers for validation and from the last 4 speakers for testing the model. The videos were upscaled to 60 frames per second using video processing software to have more temporal information and avoid artefacts in the generation of events. The event-based
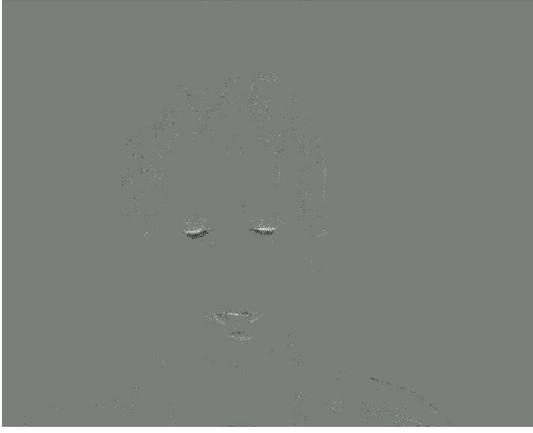
Fig. 2. Snapshot of a person talking in front of an EDC camera



Fig. 3. Optical flow representation when using regions of $10 \times 10$ pixels

data stream was generated by pointing the ATIS event-driven camera [25] ($240 \times 304$ pixels with 8mm lens) towards a high definition LED monitor while the upscaled videos were played. Due to the low quality of the original videos ($360 \times 288$ pixels resolution) and in order to preserve the details in lip movements, we cropped the mouth area over $100 \times 50$ pixels from the event stream.

### B. Model training

*1) Audio pre- and post-processing:* Following the state of the art in speech separation and enhancement, the audio original waveforms were pre-processed through Short Time Fourier Transform (STFT) applied over the over the audio waveforms resampled at 16kHz. STFT was applied using Fast Fourier Transform (FFT) size of 512, Hann window of length 25 ms, and hop length of 10 ms. The spectrogram $|x|^p$ of each input audio sample was obtained performing power-law compression of the STFT magnitude with p = 0.3. Finally, the data was normalized per-speaker with 0 mean and 1 standard deviation. To reconstruct the clean audio, on the post-processing stage, the inverse STFT to the estimated clean spectrogram was applied using the phase of the noisy input signal.

*2) Video pre-processing:* First, we compute the optical flow with the method explained in [27]. To align the visual and audio features, we generated frames from the optical flow stream every 10 ms. Over each frame. However, due to the nature of event-driven cameras, the number pixels that generate optical flow in each frame is different and therefore, the number of video features in each frame is different. To avoid this problem, we generate regions of same size across the $100 \times 50$ pixels.

For each region, we compute the mean of the $x$ component and $y$ component of optical flow and the event-rate, the number of events on each location at each frame. For example, with regions of $10 \times 10$ pixels we have a total of 50 regions and if we compute the event-rate and the mean of the $x$ and $y$ components, we have 150 video features. Fig. 3 shows an
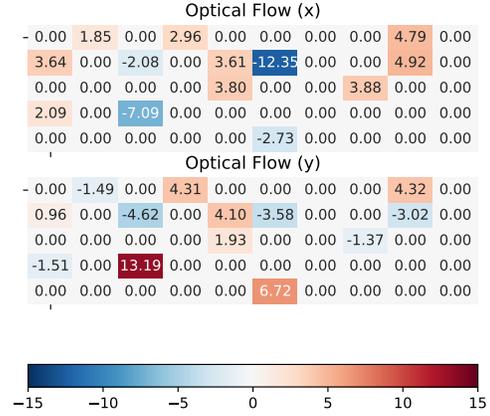
example of the $x$ component and $y$ component of the optical flow for a specific frame.

*3) Recurrent Neural Network:* The RNN model consists of 5 stacked Bidirectional Long Short-Term Memory (BiLSTM) layers, with 250 neurons in each layer. The inputs of the model are the audio and visual features concatenated. The output of the network are an Ideal Amplitude Mask (IAM) and the loss function $J_{mr}$:

$$\mathbf{p}_t[f] = \frac{\mathbf{s}_t[f]}{\mathbf{y}_t[f]} \tag{1}$$

$$J_{mr} = \sum_{t=1}^{T} \sum_{f=1}^{d} (\hat{\mathbf{p}}_t[f] \cdot \mathbf{y}_t[f] - \mathbf{s}_t[f])^2 \tag{2}$$

Where $p_t[f]$ is IAM, $s_t[f]$ is the target clean spectrogram, $y_t[f]$ is the noisy spectrogram and $\hat{p}_t[f]$ is the estimated IAM at each frequency bin $f \in [1, ..., d]$.

We train the model using the Adam optimizer and 20% of dropout to avoid overfitting. Each model is trained up to 500 epochs and early stopping is applied on the validation set to stop the training process.

### IV. RESULTS

To measure the performance of each model, we use the well known source-to-distortion ratio (SDR) and PESQ [28], to quantify the separation of the target speech from the concurrent speech and the quality of cleaned speech (i.e. the speech enhancement measure), respectively.

Table I shows the results from three different models. To train the first model we used 150 video features as input (concatenated with the audio features). These 150 features correspond to the $x$ and $y$ components of the optical-flow and the event-rate (the number of events) for each of the 50 regions ($10 \times 10$ pixels each region). The results are quite good, with higher than 7.0 SDR and on pare with the frame-based approach on PESQ performance. This shows that,

|                                    | SDR  | PESQ |
|------------------------------------|------|------|
| Noisy signal                       | 0.21 | 1.94 |
| Frame-based approach [17]          | **7.37** | **2.65** |
| Event-based approach (150 features)| 7.03 | **2.65** |
| Event-based approach (400 features)| 6.58 | 2.59 |
| Event-based approach (LSTM)        | 3.79 | 2.22 |

TABLE I

GRID DATASET RESULTS.

although the original GRID dataset is frame-based, the event-based approach shows similar performance as the frame-based approach, despite the low quality and noisy events of the dataset used, that was not recorded live with the subjects, but obtained by recording a movie played back on a high resolution monitor.

In the next experiment we decrease the size of each region to $5 \times 5$ pixels in order to have more localized video features. However, the number of input features increases enormously. That is why we only used $x$ and $y$ components of the optical-flow not including the event rate like in the previous case(400 visual input features in total). Although the results are good (6.58 SDR and 2.59 PESQ), they are not close to those achieved with 150 input features. Besides, the training and inference time increases using 400 visual input features.

One of the drawbacks of BiLSTMs used in the previous experiments is that they need to pass all the features forward and backward before giving a prediction. That means that BiLSTM has higher latency than unidirectional RNN architectures. That is why we carried out one final experiment using LSTM instead of BiLSTM. However, the results show that the performance of deep LSTM is far from that yielded by the models with BiLSTM.

Finally, we compare the computation time of both approaches. Computing the face landmarks movements on frame-based approach for each video file (each video is 3
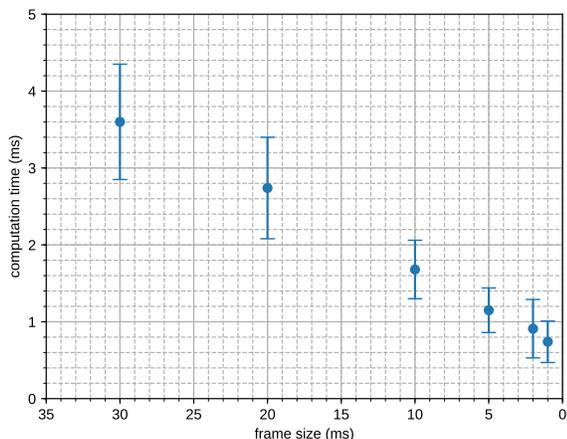
seconds long) the mean computation time using Dlib [29] is 2.980 seconds with 0.825 standard deviation (SD). On the other hand, for a event-based approach the mean is 1.126 seconds with 0.212 SD. The computation time of the event-based approach is almost three times less than the frame-based approach. The computation time of the event-based approach is divided as follows: 0.679 seconds computing the optical-flow and 0.447 seconds mapping optical-flow to the regions. In Figure 4 the computation time of the optical-flow for different frame sizes is shown e.g. for a frame size of 10 ms we accumulate events every 10 ms to generate a frame and we compute the optical for the events in each frame. It can be seen that for all the cases the computation time of the optical-flow is lower than the frame size and there is enough time for mapping the optical-flow to the regions i.e. it is possible to extract the visual features before the next frame arrives without leaks. The hardware used to measure the computation time is and Intel® Core™ i7-7500U CPU @ 2.70GHz x 4.

## V. CONCLUSIONS AND FUTURE WORK

This work presents a RNN for target speaker audio extraction on multi-talker environment using event-driven camera for visual motion feature estimation. The system is trained on the GRID dataset. We show that, although this approach does not outperform the frame based approach in terms of quality of speech enhancement, the performance is almost on pair to the frame-based approach. We believe that improvements on the quality will be obtained when using data recorded directly with the EDC, as it will improve the spatial resolution and signal to noise ratio of the dataset. At pair quality, the event-driven approach offers advantages in terms of computational cost and latency, that are critical for online, embedded applications. The proposed method required the pre-processing of the frame-based dataset, using upscaling to 60 fps, for recording the event-driven dataset. However, this operation is only required once and won't be required in an online system where the visual signal is directly recorded by means of an EDC. The same linear interpolation operation needs to be always performed on the frame-based implementation to align video with audio features. The computation of the visual features depends on the scene, but is as low as 4ms, leading to a very low latency system implementation. To further reduce the latency of the output clean speech, we substituted the BiLSTM, that requires the passing of all the features forward and backward, with an LSTM, however, the quality of the processed speech signal is far from being comparable to that of BiLSTM. Further work needs to address this problem.

To the best of our knowledge, this is the first work that uses event-driven cameras to address the target speaker extraction task. We showed that the method used to compute the optical flow to extract visual features is more efficient than the frame-based method used in [17] and is better suited for embedded applications. Finally, this work shows that the $x$ and $y$ components of the optical flow from the lip region can be useful video features for target speaker audio extraction.



Fig. 4. Computation time of optical-flow to extract visual features (Intel® Core™ i7-7500U CPU @ 2.70GHz x 4)

REFERENCES

[1] E. C. Cherry, "Some experiments on the recognition of speech, with one and with two ears," *Journal of the Acoustical Society of America*, no. 25, pp. 975–979, 1953.

[2] A. W. Bronkhorst, "The cocktail-party problem revisited: early processing and selection of multi-talker speech," *Atten Percept Psychophys*, vol. 77, no. 5, pp. 1465–1487, 2015.

[3] E. Zion Golumbic, G. B. Cogan, C. E. Schroeder, and D. Poeppel, "Visual input enhances selective speech envelope tracking in auditory cortex at a "cocktail party"," *Journal of Neuroscience*, vol. 33, no. 4, pp. 1417–1426, 2013. [Online]. Available: https://www.jneurosci.org/content/33/4/1417

[4] Y. Isik, J. L. Roux, Z. Chen, S. Watanabe, and J. R. Hershey, "Single-channel multi-speaker separation using deep clustering," in *Interspeech*, 2016.

[5] M. Kolbaek, D. Yu, Z.-H. Tan, and J. Jensen, "Multitalker speech separation with utterance-level permutation invariant training of deep recurrent neural networks," *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, vol. 25, no. 10, pp. 1901–1913, Oct. 2017.

[6] C. Han, Y. Luo, and N. Mesgarani, "Online deep attractor network for real-time single-channel speech separation," in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2019, pp. 361–365.

[7] Y. Luo and N. Mesgarani, "Tasnet: Time-domain audio separation network for real-time, single-channel speech separation," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2018, pp. 696–700.

[8] ——, "Conv-tasnet: Surpassing ideal time–frequency magnitude masking for speech separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 8, pp. 1256–1266, Aug 2019.

[9] K. Zmolikova, M. Delcroix, K. Kinoshita, T. Higuchi, A. Ogawa, and T. Nakatani, "Speaker-aware neural network based beamformer for speaker extraction in speech mixtures." in *Interspeech*, 2017, pp. 2655–2659.

[10] Q. Wang, H. Muckenhirn, K. Wilson, P. Sridhar, Z. Wu, J. Hershey, R. A. Saurous, R. J. Weiss, Y. Jia, and I. L. Moreno, "Voicefilter: Targeted voice separation by speaker-conditioned spectrogram masking," 2018.

[11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 770–778.

[12] T. Afouras, J. S. Chung, and A. Zisserman, "The conversation: Deep audio-visual speech enhancement," in *INTERSPEECH*, 2018.

[13] J. Wu, Y. Xu, S.-X. Zhang, L.-W. Chen, M. Yu, L. Xie, and D. Yu, "Time domain audio visual speech separation," 2019.

[14] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein, "Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation," *ACM Trans. Graph.*, vol. 37, no. 4, pp. 112:1–112:11, Jul. 2018. [Online]. Available: http://doi.acm.org/10.1145/3197517.3201357

[15] M. Cooke, J. Barker, S. Cunningham, and X. Shao, "An audio-visual corpus for speech perception and automatic speech recognition," *The Journal of the Acoustical Society of America*, vol. 120, no. 5, pp. 2421–2424, November 2006. [Online]. Available: http://link.aip.org/link/?JAS/120/2421/1

[16] A. Gabbay, A. Ephrat, T. Halperin, and S. Peleg, "Seeing through noise: Visually driven speaker separation and enhancement," 2017, pp. 3051–3055.

[17] G. Morrone, S. Bergamaschi, L. Pasa, L. Fadiga, V. Tikhanoff, and L. Badino, "Face landmark-based speaker-independent audio-visual speech enhancement in multi-talker environments," in *2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6900–6904.

[18] A. Glover and C. Bartolozzi, "Robust visual tracking with a freely-moving event camera," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sep. 2017, pp. 3769–3776.

[19] L. A. Camuñas-Mesa, T. Serrano-Gotarredona, S. Ieng, R. Benosman, and B. Linares-Barranco, "Event-driven stereo visual tracking algorithm to solve object occlusion," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 9, pp. 4223–4237, Sep. 2018.

[20] M. Iacono, S. Weber, A. Glover, and C. Bartolozzi, "Towards event-driven object detection with off-the-shelf deep learning," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct 2018, pp. 1–9.

[21] T. Stoffregen, G. Gallego, T. Drummond, L. Kleeman, and D. Scaramuzza, "Event-based motion segmentation by motion compensation," 2019.

[22] J. Maro, G. Lenz, C. Reeves, and R. Benosman, "Event-based visual gesture recognition with background suppression running on a smartphone," in *2019 14th IEEE International Conference on Automatic Face Gesture Recognition (FG 2019)*, May 2019, pp. 1–1.

[23] X. Li, D. Neil, T. Delbruck, and S. Liu, "Lip reading deep network exploiting multi-modal spiking visual and auditory sensors," in *2019 IEEE International Symposium on Circuits and Systems (ISCAS)*, May 2019, pp. 1–5.

[24] A. Savran, R. Tavarone, B. Higy, L. Badino, and C. Bartolozzi, "Energy and computation efficient audio-visual voice activity detection driven by event-cameras," in *2018 13th IEEE International Conference on Automatic Face Gesture Recognition (FG 2018)*, May 2018, pp. 333–340.

[25] C. Posh, D. Matolin, and R. Wohlgenannt, "A QVGA 143 dB dynamic range frame-free PWM image sensor with lossless pixel-level video compression and time-domain CDS," in *IEEE Journal of Solid-State Circuits*, vol. 46, Jan. 2011, pp. 259–275.

[26] A. Z. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, "Unsupervised event-based learning of optical flow, depth, and egomotion," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 989–997.

[27] R. Benosman, C. Clercq, X. Lagorce, S. Ieng, and C. Bartolozzi, "Event-based visual flow," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 25, no. 2, pp. 407–417, Feb 2014.

[28] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, "Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs," in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.01CH37221)*, vol. 2, May 2001, pp. 749–752 vol.2.

[29] D. E. King, "Dlib-ml: A machine learning toolkit," *Journal of Machine Learning Research*, vol. 10, no. 60, pp. 1755–1758, 2009. [Online]. Available: http://jmlr.org/papers/v10/king09a.html