

Asynchronous Convolutional Networks for Object Detection in Neuromorphic Cameras

Marco Cannici Marco Ciccone Andrea Romanoni Matteo Matteucci
Politecnico di Milano, Italy

{marco.cannici, marco.ciccone, andrea.romanoni, matteo.matteucci}@polimi.it

Abstract

Event-based cameras, also known as neuromorphic cameras, are bioinspired sensors able to perceive changes in the scene at high frequency with low power consumption. Becoming available only very recently, a limited amount of work addresses object detection on these devices. In this paper we propose two neural networks architectures for object detection: YOLE, which integrates the events into surfaces and uses a frame-based model to process them, and fcYOLE, an asynchronous event-based fully convolutional network which uses a novel and general formalization of the convolutional and max pooling layers to exploit the sparsity of camera events. We evaluate the algorithm with different extensions of publicly available datasets, and on a novel synthetic dataset.

1. Introduction

Fundamental techniques underlying Computer Vision are based on the ability to extract meaningful features. To this extent, Convolutional Neural Networks (CNNs) rapidly became the first choice in many computer vision applications such as image classification [18, 45, 13, 48], object detection [42, 41, 25], semantic scene labeling [49, 37, 26], and they have been recently extended also to non-euclidean domains such as manifolds and graphs [16, 31]. In most of the cases the input of these networks are images.

In the meanwhile, neuromorphic cameras [43, 36, 3] are becoming more and more widespread. These devices are bio-inspired vision sensors that attempt to emulate the functioning of biological retinas. As opposed to conventional cameras, which generate frames at a constant frame rate, these sensors output data only when a brightness change is detected in the field of view. Whenever this happens, an event $\mathbf{e} = \langle x, y, ts, p \rangle$ is generated indicating the position (x, y) , the instant ts at which the change has been detected and its polarity $p \in \{1, -1\}$, *i.e.*, if the brightness change is positive or negative. The result is a sen-

sor able to produce a stream of asynchronous events that sparsely encodes changes with microseconds resolution and with minimum requirements in terms of power consumption and bandwidth. The growth in popularity of these type of sensors, and their advantages in terms of temporal resolution and reduced data redundancy, have led to fully exploit the advantages of event-based vision for a variety of applications, *e.g.*, object tracking [39, 29, 11], visual odometry [32, 40], and optical flow estimation [2, 24, 47].

Spiking Neural Networks (SNNs) [27], a processing model aiming to improve the biological realism of artificial neural networks, are one of the most popular neural model able to directly handle events. Despite their advantages in terms of speed and power consumption, however, training deep SNNs models on complex tasks is usually very difficult. To overcome the lack of scalable training procedures, recent works have focused on converting pre-trained deep networks to SNNs, achieving promising results even on complex tasks [14, 5, 8].

An alternative solution to deal with event-based cameras is to make use of frame integration procedures and conventional frame-based networks [35] which can instead rely on optimized training procedures. Recently, other alternatives to SNNs making use of hierarchical time surfaces [20] and memory cells [46] have also been introduced. Another solution, proposed in [33], suggests instead the use of LSTM cells to accumulate events and perform classification. An extension of this work making use of attention mechanisms has also been proposed in [4].

Although event-cameras are becoming increasingly popular, only very few datasets for object detection in event streams are available, and a limited number of object detection algorithms has been proposed [23, 6, 38].

In this paper we introduce a novel hybrid approach to extract features for object detection problems using neuromorphic cameras. The proposed framework allows the design of object detection networks able to sparsely compute features while still preserving the advantages of conventional neural networks. More importantly, networks implemented using the proposed procedure are asynchronous, meaning that

computation is only performed when a sequence of events arrive and only where previous results need to be recomputed.

In Section 3 the convolution and max-pooling operations are reformulated by adding an internal state, *i.e.*, a memory of the previous prediction, that allows us to sparsely recompute feature maps. An asynchronous fully-convolutional network for event-based object detection which exploits this formulation is finally described in Section 3.4.

2. Background

Leaky Surface. The basic component of the proposed architectures is a procedure able to accumulate events. Sparse events generated by the neuromorphic camera are integrated into a *leaky surface*, a structure that takes inspiration from the functioning of Spiking Neural Networks (SNNs) to maintain memory of past events. A similar mechanism has already been proposed in [7]. Every time an event with coordinates (x_e, y_e) and timestamp ts^t is received, the corresponding pixel of the surface is incremented of a fixed amount Δ_{incr} . At the same time, the whole surface is decremented by a quantity which depends on the time elapsed between the last received event and the previous one. The described procedure can be formalized by the following equations:

$$q_{x_s, y_s}^t = \max(p_{x_s, y_s}^{t-1} - \lambda \cdot \Delta_{ts}, 0) \quad (1)$$

$$p_{x_s, y_s}^t = \begin{cases} q_{x_s, y_s}^t + \Delta_{incr} & \text{if } (x_s, y_s)^t = (x_e, y_e)^t \\ q_{x_s, y_s}^t & \text{otherwise} \end{cases}, \quad (2)$$

where p_{x_s, y_s}^t is the value of the surface pixel in position (x_s, y_s) of the leaky surface and $\Delta_{ts} = ts^t - ts^{t-1}$. To improve readability in following equations, we name the quantity $(ts^t - ts^{t-1}) \cdot \lambda$ as Δ_{leak} . Notice that the effects of λ and Δ_{incr} are related: Δ_{incr} determines how much information is contained in each single event whereas λ defines the decay rate of activations. Given a certain choice of these parameters, similar results can be obtained by using, for instance, a higher increment Δ_{incr} and a higher temporal λ . For this reason, we fix $\Delta_{incr} = 1$ and we vary only λ based on the dataset to be processed. Pixel values are prevented from becoming negative by means of the max operation.

Other frame integration procedures, such as the one in [35], divide the time in predefined constant intervals. Frames are obtained by setting each pixel to a binary value (depending on the polarity) if at least an event has been received in each pixel within the integration interval. With this mechanism however, time resolution is lost and the same importance is given to each event, even if it represents noise. The adopted method, instead, performs continuous and incremental integration and is able to better handle noise.

Similar procedures capable of maintaining time resolution have also been proposed, such as those that make use of exponential decays [7, 19] to update surfaces, and those relying on histograms of events [28]. Recently, the concept of *time surface* has also been introduced in [20] where surfaces are obtained by associating each event with temporal features computed applying exponential kernels to the event neighborhood. Extensions of this procedure making use of memory cells [46] and event histograms [1] have also been proposed. Although these event representations better describe complex scene dynamics, we make use of a simpler formulation to derive a linear dependence between consecutive surfaces. This allows us to design the event-based framework discussed in Section 3 in which time decay is applied to every layer of the network.

Event-based Object Detection. We identified YOLO [41] as a good candidate model to tackle the object detection problem in event-based scenarios: it is fully-differentiable and produces predictions with small input-output delays. By means of a standard CNN and with a single forward pass, YOLO is able to simultaneously predict not only the class, but also the position and dimension of every object in the scene. We used the YOLO loss and the previous leaky surface procedure to train a baseline model which we called YOLE: "You Only Look at Events". The architecture is depicted in Figure 1. We use this model as a reference to highlight the strengths and weaknesses of the framework described in Section 3, which is the main contribution of this work. YOLE processes 128×128 surfaces, it predicts $B = 2$ bounding boxes for each region and classifies objects into C different categories.

Note that in this context, we use the term YOLO to refer only to the training procedure proposed by [41] and not to the specific network architecture. We used indeed a simpler structure for our models as explained in Section 4. Nevertheless, YOLE, *i.e.*, YOLO + leaky surface, does not exploit the sparse nature of events; to address this issue, in the next section, we propose a fully event-based asynchronous framework for convolutional networks.

3. Event-based Fully Convolutional Networks

Conventional CNNs for video analysis treat every frame independently and recompute all the feature maps entirely, even if consecutive frames differ from each other only in small portions. Beside being a significant waste of power and computations, this approach does not match the nature of event-based cameras.

To exploit the event-based nature of neuromorphic vision, we propose a modification of the forward pass of fully convolutional architectures. In the following the convolution and pooling operations are reformulated to produce the final prediction by recomputing only the features corresponding to regions affected by the events. Feature maps

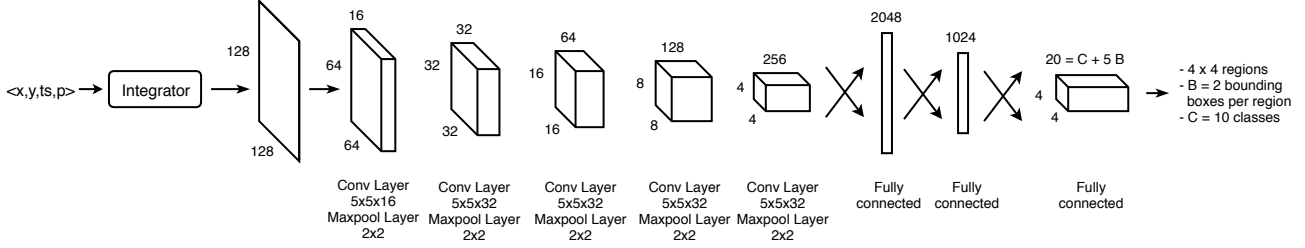


Figure 1. The YOLE detection network based on YOLO used to detect MNIST-DVS digits. The input surfaces are divided into a grid of 4×4 regions which predict 2 bounding boxes each.

maintain their state over time and are updated only as a consequence of incoming events. At the same time, the leaking mechanism that allows past information to be forgotten, acts independently on each layer of the CNN. This enables features computed in the past to fade away as their visual information starts to disappear in the input surface. The result is an asynchronous CNN able to perform computation only when requested and at different rates. The network can indeed be used to produce an output only when new events arrive, dynamically adapting to the timings of the input, or to produce results at regular rates by using the leaking mechanism to update layers in absence of new events.

The proposed framework has been developed to extend the YOLE detection network presented in Section 2. Nevertheless, this method can be applied to any convolutional architecture to perform asynchronous computation. A CNN trained to process frames reconstructed from streams of events can indeed be easily converted into an event-based CNN without any modification on its layers composition, and by using the same weights learned while observing frames, maintaining its output unchanged.

3.1. Leaky Surface Layer

The procedure used to compute the leaky surface described in Section 2 is embedded into an actual layer of the proposed framework. Furthermore, to allow subsequent layers to locate changes inside the surface, the following information are also forwarded to the next layer: (i) the list of incoming events. (ii) Δ_{leak} , which is sent to all the subsequent layers to update feature maps not affected by the events. (iii) the list of surface pixels which have been reset to 0 by the max operator in Equation (1).

3.2. Event-based Convolutional Layer (e-conv)

The *event-based convolutional* (e-conv) layer we propose uses events to determine where the input feature map has changed with respect to the previous time step and, therefore, which parts of its *internal state*, *i.e.*, the feature map computed at the previous time step, must be recomputed and which parts can be reused. We use a procedure similar to the one described in the previous section to let

features decay over time. However, due to the transformations applied by previous layers and the composition of their activation functions, Δ_{leak} may act differently in different parts of the feature map. For instance, the decrease of a pixel intensity value in the input surface may cause the value computed by a certain feature in a deeper layer to decrease, but it could also cause another feature of the same layer to increase. The update procedure, therefore, must also be able to accurately determine how a single bit of information is transformed by the network through all the previous layers, in any spatial location. We face this issue by storing an additional feature map, $F_{(n)}$, and by using a particular class of activation functions in the hidden layers of the network.

Let's consider the first layer of a CNN which processes surfaces obtained using the procedure described in the previous section and which computes the convolution of a set of filters W with bias b and activation function $g(\cdot)$. The computation performed on each receptive field is:

$$y_{ij(1)}^t = g \left(\sum_h \sum_k x_{h+i,k+j}^t W_{hk(1)} + b_{(1)} \right) = g(\tilde{y}_{ij(1)}^t), \quad (3)$$

where h, k select a pixel $x_{h+i,k+j}^t$ in the receptive field of the output feature (i, j) and its corresponding value in the kernel W , whereas the subscript (1) indicates the hidden-layer of the network (in this case the first after the leaky surface layer).

When a new event arrives, the leaky surface layer decreases all the pixels by Δ_{leak} , *i.e.*, a pixel not directly affected by the event becomes: $x_{hk}^{t+1} = x_{hk}^t - \Delta_{leak}^{t+1}$, with $\Delta_{leak}^{t+1} > 0$. At time $t + 1$ Equation (3) becomes:

$$\begin{aligned} y_{ij(1)}^{t+1} &= g \left(\sum_h \sum_k x_{h+i,k+j}^{t+1} W_{hk(1)} + b_{(1)} \right) \\ &= g \left(\sum_h \sum_k (x_{h+i,k+j}^t - \Delta_{leak}^{t+1}) W_{hk(1)} + b_{(1)} \right) \\ &= g \left(\tilde{y}_{ij(1)}^t - \Delta_{leak}^{t+1} \sum_h \sum_k W_{hk(1)} \right). \end{aligned} \quad (4)$$

If $g(\cdot)$ is (i) a piecewise linear activation function $g(x) = \{\alpha_i \cdot x \text{ if } x \in D_i\}$, as ReLU or Leaky ReLU, and we assume that (ii) the updated value does not change which linear segment of the activation function the output falls onto and, in this first approximation, (iii) the leaky surface layer does not restrict pixels using $\max(\cdot, 0)$, Equation 4 can be rewritten as it follows:

$$y_{ij(1)}^{t+1} = y_{ij(1)}^t - \Delta_{leak}^{t+1} \alpha_{ij(1)} \sum_h \sum_k W_{hk(1)}, \quad (5)$$

where $\alpha_{ij(1)}$ is the coefficient applied by the piecewise function $g(\cdot)$ which depends on the feature value at position (i, j) . When the previous assumption is not satisfied, the feature is recomputed as its receptive field was affected by an event (*i.e.*, applying the filter W locally to x^{t+1}).

Consider now a second convolutional layer attached to the first one:

$$\begin{aligned} y_{ij(2)}^{t+1} &= g \left(\sum_{h,k} y_{i+h,j+k(1)}^{t+1} W_{hk(2)} + b(2) \right) \\ &= g \left(\sum_{h,k} \left(y_{i+h,j+k(1)}^t - \Delta_{leak}^{t+1} \alpha_{i+h,j+k(1)} \sum_{h',k'} W_{h'k'(1)} \right) W_{hk(2)} + b(2) \right) \\ &= y_{ij(2)}^t - \Delta_{leak}^{t+1} \alpha_{ij(2)} \sum_{h,k} \left(\alpha_{i+h,j+k(1)} \sum_{h',k'} W_{h'k'(1)} \right) W_{hk(2)} \\ &= y_{ij(2)}^t - \Delta_{leak}^{t+1} \alpha_{ij(2)} \sum_{h,k} F_{i+h,j+k(1)}^{t+1} W_{hk(2)} = y_{ij(2)}^t - \Delta_{leak}^{t+1} F_{ij(2)}^{t+1}. \end{aligned} \quad (6)$$

The previous equation can be extended by induction as it follows:

$$y_{ij(n)}^{t+1} = y_{ij(n)}^t - \Delta_{leak}^{t+1} F_{ij(n)}^{t+1},$$

$$\text{with } F_{ij(n)}^{t+1} = \alpha_{ij(n)} \sum_h \sum_k F_{i+h,j+k(n-1)}^{t+1} W_{hk(n)} \text{ if } n > 1, \quad (7)$$

where $F_{ij(n)}$ expresses how visual inputs are transformed by the network in every receptive field (i, j) , *i.e.*, the composition of the previous layers activation functions.

Given this formulation, the max operator applied by the leaky surface layer can be interpreted as a ReLU, and Equation (5) becomes:

$$y_{ij(1)}^{t+1} = y_{ij(1)}^t - \Delta_{leak}^{t+1} \alpha_{ij(1)} \sum_h \sum_k F_{i+h,j+k(0)}^{t+1} W_{hk(1)}, \quad (8)$$

where the value $F_{i+h,j+k(0)}$ is 0 if the pixel $x_{i+h,j+k} \leq 0$ and 1 otherwise.

Notice that $F_{ij(n)}$ needs to be updated only when the corresponding feature changes enough to make the activation function use a different coefficient α , *e.g.*, from 0 to 1 in case of ReLU. In this case $F_{ij(n)}$ is updated locally in correspondence of the change by using the update matrix of the previous layer and by applying Equation 7 only for the features whose activation function has changed. Events are

used to communicate the change to subsequent layers so that their update matrix can also be updated accordingly.

The internal state of the e-conv layer, therefore, comprises the feature maps $y_{(n)}^{t-1}$ and the update values $F_{(n)}^{t-1}$ computed at the previous time step. The initial values of the internal state are computed making full inference on a blank surface; this is the only time the network needs to be executed entirely. As a new sequence of events arrives the following operations are performed (see Figure 3(a)):

- i. Update $F_{(n)}^{t-1}$ locally on the coordinates specified by the list of incoming events (Eq. (7)). Note that we do not distinguish between actual events and those generated by the use of a different slope in the linear activation function.
- ii. Update the feature map $y_{(n)}$ with Eq. (7) in the locations which are not affected by any event and generate an output event where the activation function coefficient has changed.
- iii. Recompute $y_{(n)}$ by applying W locally in correspondence of the incoming events and output which receptive field has been affected.
- iv. Forward the feature map and the events generated in the current step to the next layer.

3.3. Event-based Max Pooling Layer (e-max-pool)

The location of the maximum value in each receptive field of a max-pooling layer is likely to remain the same over time. An event-based pooling layer, hence, can exploit this property to avoid recomputing every time the position of maximum values.

The internal state of an event-based max-pooling (e-max-pool) layer can be described by a *positional matrix* $I_{(n)}^t$, which has the shape of the output feature map produced by the layer, and which stores, for every receptive field, the position of its maximum value. Every time a sequence of events arrives, the internal state $I_{(n)}^t$ is sparsely updated by recomputing the position of the maximum values in every receptive field affected by an incoming event. The internal state is then used both to build the output feature map and to produce the *update matrix* $F_{(n)}^t$ by fetching the previous layer on the locations provided by the indices $I_{ij(n)}^t$. For each e-max-pool layer, the indices of the receptive fields where the maximum value changes are communicated to the subsequent layers so that the internal states can be updated accordingly. This mechanism is depicted in Figure 3(b).

Notice that the leaking mechanism acts differently in distinct regions of the input space. Features inside the same receptive field can indeed decrease over time with different speeds as their update values $F_{ij(n)}^t$ could be different. Therefore, even if no event has been detected inside a region, the position of its maximum value might change.

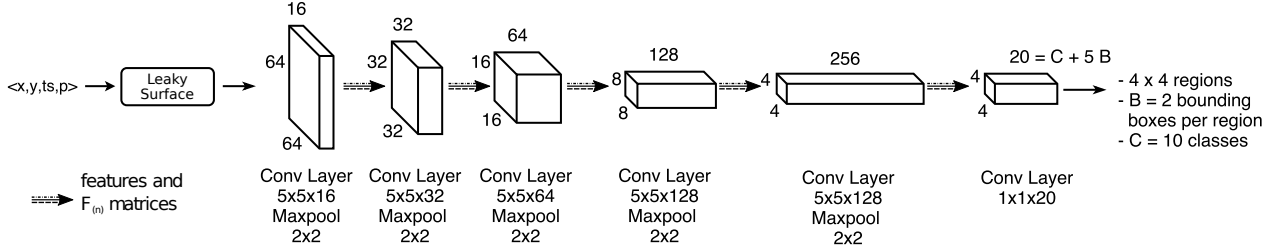


Figure 2. fcYOLO: a fully-convolutional detection network based on YOLE. The last layer is used to map the feature vectors into a set of 20 values which define the parameters of the predicted bounding boxes.

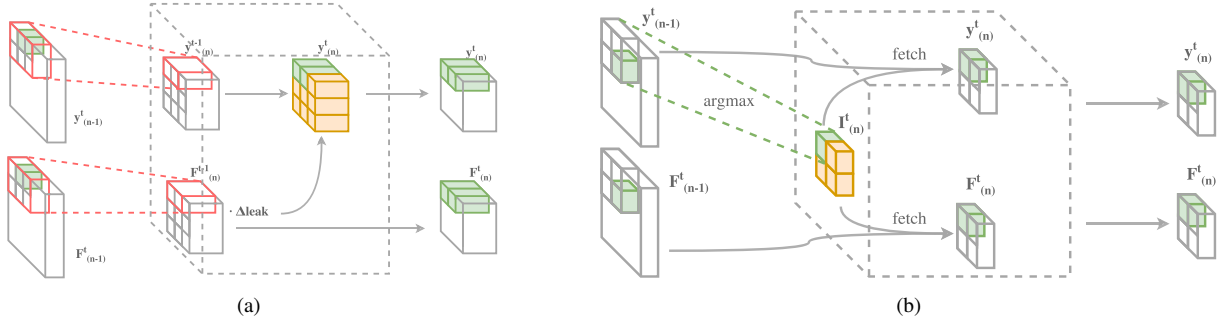


Figure 3. The structure of the e-conv (a) and e-max-pooling layers (b). The internal states and the update matrices are recomputed locally only where events are received (green cells) whereas the remaining regions (depicted in yellow) are obtained reusing the previous state.

However, if an input feature M has the minimum update rate $F_{M(n-1)}$ among features in its receptive field R and it also corresponds to the maximum value in R , the corresponding output feature will decrease slower than all the others in R and its value will remain the maximum. In this case, its index $I_{(n)R}^t$ does not need to be recomputed until a new event arrives in R . We check if the maximum has to be recomputed for each receptive field affected by incoming events and also in all positions where the previous condition does not hold.

3.4. Event FCN for Object Detection (fcYOLO)

To fully exploit the event-based layers presented so far, the YOLE model described in Section 2 needs to be converted into a fully convolutional object detection network replacing all its layers with their event-based versions (see Figure 3). Moreover, fully-connected layers are replaced with 1×1 e-conv layers which map features extracted by the previous layers into a precise set of values defining the bounding boxes parameters predicted by the network. Training was first performed on a network composed of standard layers; the learned weights were then used with e-conv and e-max-pool layers during inference.

This architecture divides the 128×128 field of view into a grid of 4×4 regions that predicts 2 bounding boxes each and classify the detected objects into C different classes. The last 1×1 e-conv layer is used to decrease the dimensionality of the feature vectors and to map them into the right set of parameters, regardless of their position in the

field of view.

Moreover, this architecture can be used to process surfaces of different sizes without the need to re-train or re-design it. The subnetworks processing 32×32 regions, in fact, being defined by the same set of parameters, can be stacked together to process even larger surfaces.

4. Experiments

4.1. Datasets

Only few event-based object recognition datasets are publicly available in the literature. The most popular ones are: N-MNIST [34], MNIST-DVS [44], CIFAR10-DVS [22], N-Caltech101 [34] and POKER-DVS [44]. These datasets are obtained from the original MNIST [21], CIFAR-10 [17] and Caltech101 [10] datasets by recording the original images with an event camera while moving the camera itself or the images of the datasets. We performed experiments on N-Caltech101 and on modified versions of N-MNIST and MNIST-DVS for object detection, *i.e.*, *Shifted N-MNIST* and *Shifted MNIST-DVS*, and on an extended version of POKER-DVS, namely *OD-Poker-DVS*. Moreover we also perform experiments on a synthetic dataset, named *Blackboard MNIST*, showing digits written on a blackboard. A detailed description of these datasets is provided in the supplementary materials.

Shifted N-MNIST The N-MNIST [34] dataset is a conversion of the popular MNIST [21] image dataset for computer vision. We enhanced this collection by building a

slightly more complex set of recordings. Each sample is indeed composed of two N-MNIST samples placed in two random non-overlapping locations of a bigger 124×124 field of view. Each digit was also preprocessed by extracting its bounding box which was then moved, along with the events, in its new position of the bigger field of view. The final dataset is composed of 60,000 training and 10,000 testing samples.

Shifted MNIST-DVS We used a similar procedure to obtain Shifted MNIST-DVS recordings. We first extracted bounding boxes with the same procedure used in Shifted N-MNIST and then placed them in a 128×128 field of view. We mixed MNIST-DVS *scale4*, *scale8* and *scale16* samples within the same recording obtaining a dataset composed of 30,000 samples.

OD-Poker-DVS The Poker-DVS dataset is a small collection of neuromorphic samples showing poker card pips obtained by extracting 31×31 symbols from three deck recordings. We used the tracking algorithm provided with the dataset to track pips and enhance the original uncut deck recordings with their bounding boxes. We finally divided these recordings into a set of shorter examples obtaining a collection composed of 218 training and 74 testing samples.

Blackboard MNIST We used the DAVIS simulator released by [32] to build our own set of synthetic recordings. The resulting dataset consists of a number of samples showing digits written on a blackboard in random positions and with different scales. We preprocessed a subset of images from the original MNIST dataset by removing their background and by making them look as if they were written with a chalk. Sets of digits were then placed on the image of a blackboard and the simulator was finally run to obtain event-based recordings and the bounding boxes of every digit visible within the camera field of view. The resulting dataset is the union of three simulations featuring increasing level of variability in terms of camera trajectories and digit dimensions. The overall dataset is composed of 2750 training and 250 testing samples.

N-Caltech101 The N-Caltech101 [34] collection is the only publicly available event-based dataset providing bounding boxes annotations. We split the dataset into 80% training and 20% testing samples using a stratified split. Since no ground truth bounding boxes are available for the *background* class, we decided not to use this additional category in our experiments.

4.2. Experiments Setup

Event-based datasets, especially those based on MNIST, are generally simpler than the image-based ones used to train the original YOLO architecture. Therefore, we designed the MNIST object detection networks taking inspiration from the simpler LeNet [21] model composed of 6 conv-pool layers for feature extraction. Both YOLE and

fcYOLO share the same structure up to the last regression/classification layers.

For what concerns the N-Caltech101 dataset, we used a slightly different architecture inspired by the structure of the VGG16 model [45]. The network is composed by only one layer for each group of convolutional layers, as we noticed that a simpler architecture achieved better results. Moreover, the dimensions of the last fully-connected layers have been adjusted such that the surface is divided into a grid of 5×7 regions predicting $B = 2$ bounding boxes each. As in the original YOLO architecture we used Leaky ReLU for the activation functions of hidden layers and a linear activation for the last one.

In all the experiments the first 4 convolutional layers have been initialized with kernels obtained from a recognition network pretrained to classify target objects, while the remaining layers using the procedure proposed in [12]. All networks were trained optimizing the multi-objective loss proposed by [41] using Adam [15], learning rate 10^{-4} , $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. The batch size was chosen depending on the dataset: 10 for Shifted N-MNIST, 40 for Shifted MNIST-DVS and N-Caltech101, 25 for Blackboard MNIST and 35 for Poker-DVS with the aim of filling the memory of the GPU optimally. Early-stopping was applied to prevent overfitting using validation sets with the same size of the test set.

4.3. Results and Discussion

Detection performance of YOLE. The YOLE network achieves good detection results both in terms of mean average precision (mAP) [9] and accuracy, which in this case is computed by matching every ground truth bounding box with the predicted one having the highest intersection over union (IOU), in most of the datasets. The results we obtained are summarized in Table 3.

We used the Shifted N-MNIST dataset also to analyze how detection performance changes when the network is used to process scenes composed of a variable number of objects, as reported in Table 4. We denote as $v1$ the results obtained using scenes composed of a single digit and with $v2$ those obtained with scenes containing two digits in random locations of the field of view. We further tested the robustness of the proposed model by adding some challenging noise. We added non-target objects ($v2fr$) in the form of five 8×8 fragments, taken from random N-MNIST digits using a procedure similar to the one used to build the *Cluttered Translated MNIST* dataset [30], and 200 additional random events per frame ($v2fr+ns$).

In case of multiple objects the algorithm is still able to detect all of them, while, as expected, performance drops both in terms of accuracy and mean average precision when dealing with noisy data. Nevertheless, we achieved very good detection performance on the Shifted MNIST-DVS,

Table 1. YOLE Top-20 average precisions on N-Caltech101. Full table provided in the supplemental material.

	Motorbikes	airplanes	Faces easy	metronome	laptop	dollar bill	umbrella	watch	minaret	grand piano	menorah	inline skate	saxophone	stapler	windsor chair	rooster	yin yang	Leopards	trilobite	garfield
AP	97.8	95.8	94.7	88.3	88.1	86.5	85.9	84.2	81.3	81.3	80.7	75.1	68.4	68.1	65.2	64.5	63.3	62.9	62.5	62.3
N_{train}	480	480	261	20	49	32	45	145	46	61	53	19	24	27	34	31	36	120	52	22

Table 2. fcYOLE Top-20 average precisions on N-Caltech101. Full table provided in the supplemental material.

	Motorbikes	airplanes	Faces easy	watch	dollar bill	car side	grand piano	menorah	metronome	umbrella	yin yang	saxophone	minaret	soccer ball	Leopards	dragonfly	stop sign	windsor chair	accordion	buddha
AP	97.5	96.8	92.2	75.7	74.4	70.3	69.5	67.7	63.4	61.0	60.4	59.7	59.5	57.3	57.2	55.6	55.1	52.3	48.3	46.5
N_{train}	480	480	261	145	32	75	61	53	20	45	36	24	46	40	120	42	40	34	33	51

Table 3. Performance comparison between YOLE and fcYOLE.

S-MNIST-DVS				Blackboard MNIST			
fcYOLE		YOLE		fcYOLE		YOLE	
acc	mAP	acc	mAP	acc	mAP	acc	mAP
94.0	87.4	96.1	92.0	88.5	84.7	90.4	87.4
OD-Poker-DVS				N-Caltech101			
fcYOLE		YOLE		fcYOLE		YOLE	
acc	mAP	acc	mAP	acc	mAP	acc	mAP
79.10	78.69	87.3	82.2	57.1	26.9	64.9	39.8

Table 4. YOLE performance on S-N-MNIST variants.

	S-N-MNIST				
	v1	v2	v2*	v2fr	v2fr+ns
accuracy	94.9	91.7	94.7	88.6	85.5
mAP	91.3	87.9	90.5	81.5	77.4

Blackboard MNIST and Poker-DVS datasets which represent a more realistic scenario in terms of noise. All of these experiments were performed using the set of hyperparameters suggested by the original work from [41]. However, a different choice of these parameters, namely $\lambda_{coord} = 25.0$ and $\lambda_{noobj} = 0.25$, worked better for us increasing both the accuracy and mean average precision scores (v2*).

The dataset in which the proposed model did not achieve noticeable results is N-Caltech101. This is mainly explained by the increased difficulty of the task and by the fact that the number of samples in each class is not evenly balanced. The network, indeed, usually achieves good results when the number of training samples is high such as with *Airplanes*, *Motorbikes* and *Faces_easy*, and in cases in which samples are very similar, e.g., *inline_skate* (see Table 1 and supplementary material). As the number of training samples decreases and the sample variability within the class increases, however, the performance of the model becomes worse, behavior which explains the poor aggregate scores we report in Table 3.

Detection performance of fcYOLE. With this fully-convolutional variant of the network we registered a slight decrease in performance w.r.t. the results we obtained using YOLE, as reported in Table 3 and Table 2. This gap in per-

formance is mainly due to the fact that each region in fcYOLE generates its predictions by only looking at the visual information contained in its portion of the field of view. Indeed, if an object is only partially contained inside a region the network has to guess the object dimensions and class by only looking at a restricted region of the surface. It should be stressed, however, that the difference in performance between the two architectures does not come from the use of the proposed event layers, whose output are the same as the conventional ones, but rather from the reduced expressive power caused by the absence of fully-connected layers in fcYOLE. Indeed, not removing them would have allowed us to obtain the same performance of YOLE, but with the drawback of being able to exploit event-based layers only up to the first FC-layer, which has not been formalized yet in an event-based form. Removing the last fully-connected layers allowed us to design a detection network made of only event-based layers and which uses also a significantly lower number of parameters. In the supplementary materials we provide a video showing a comparison between YOLE and fcYOLE predictions.

To identify the advantages and weaknesses of the proposed event-based framework in terms of inference time we compared our detection networks on two datasets, Shifted N-MNIST and Blackboard MNIST. We group events into batches of 10ms and average timings on 1000 runs. In the first dataset the event-based approach achieved a 2x speedup (22.6ms per batch), whereas in the second one it performed slightly slower (43.2ms per batch) w.r.t. a network making use of conventional layers (34.6ms per batch). The second benchmark is indeed challenging for our framework since changes are not localized in restricted regions. Our current implementation is not optimized to handle noisy scenes efficiently. Indeed, additional experiments showed that asynchronous CNNs are able to provide a faster prediction only up to a 80% of event sparsity (where with sparsity we mean the percentage of changed pixels in the reconstructed image). Further investigations are out of the scope of this paper and will be addressed in future works.

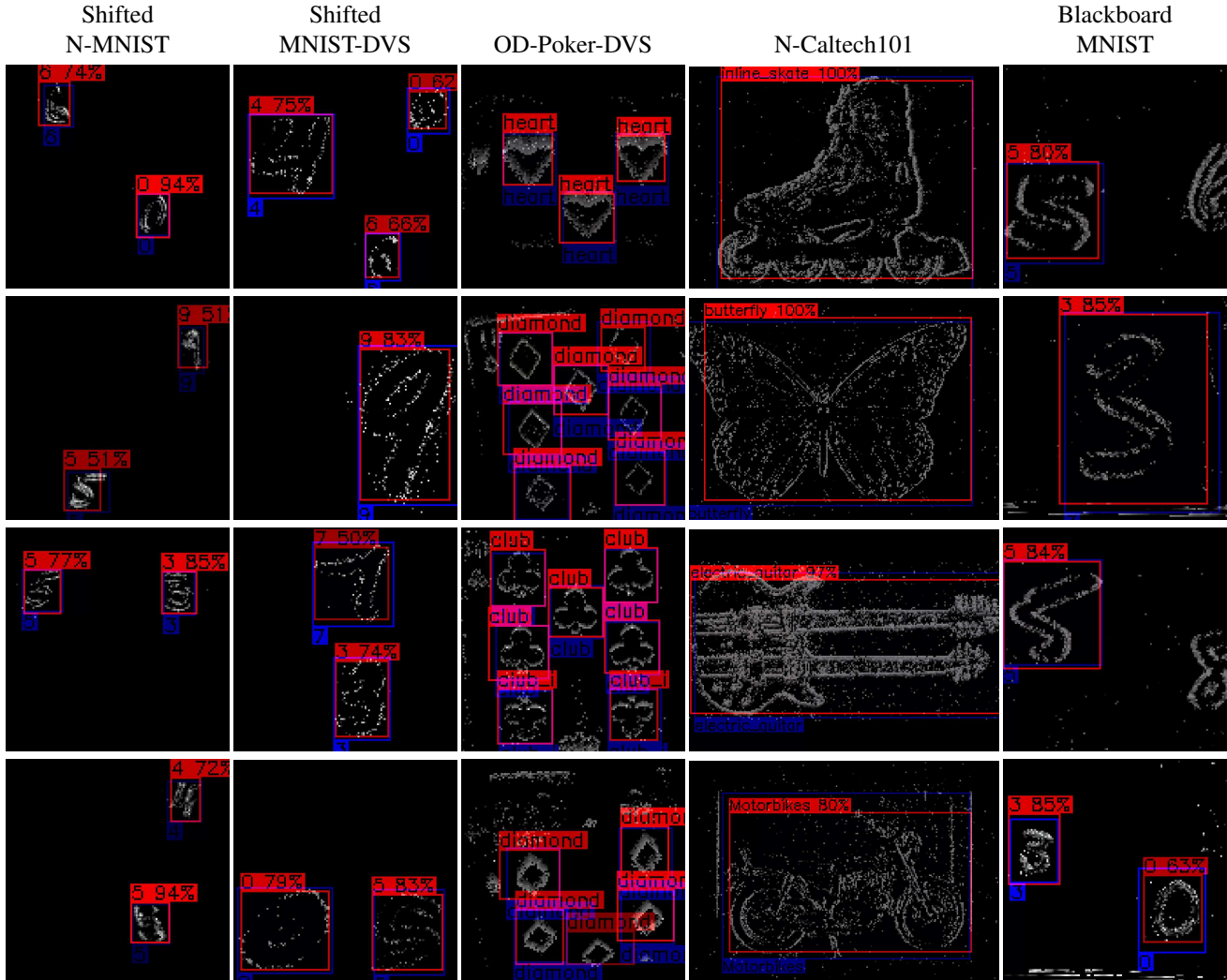


Figure 4. Examples of YOLE predictions.

5. Conclusions

We proposed two different methods, based on the YOLO architecture, to accomplish object detection in event-based cameras. The first one, namely YOLE, integrates events into a unique leaky surface. Conversely, fcYOLE relies on a more general extension of the convolutional and max pooling layers to directly deal with events and exploit their sparsity by preventing the whole network to be reprocessed. The obtained asynchronous detector dynamically adapts to the timing of the events stream by producing results only as a consequence of incoming events and by maintaining its internal state, without performing any additional computation, when no events arrive. This novel event-based framework can be used in every fully-convolutional architecture to make it usable with event-cameras, even conventional CNN for classification, although in this paper it has been applied to object detection networks.

We analyzed the timing performance of this formaliza-

tion obtaining promising results. We are planning to extend our framework to automatically detect at runtime when the use of event-based layers speeds up computation (*i.e.*, changes affect few regions of the surface) or a complete re-computation of the feature maps is more beneficial in order to exploit the benefits of both approaches. Nevertheless, we believe that a ad-hoc hardware implementation, would allow to better exploit the advantages of the proposed method enabling a fair timing comparison with SNNs, which are usually implemented in hardware.

Acknowledgements We would like to thank Prophesee for helpful discussions on YOLE. The research leading to these results has received funding from project TEINVEIN: TECnologie INnovative per i VEicoli Intelligenti, CUP (Codice Unico Progetto - Unique Project Code): E96D17000110009 - Call “Accordi per la Ricerca e l’Innovazione”, cofunded by POR FESR 2014-2020 (Programma Operativo Regionale, Fondo Europeo di Sviluppo Regionale Regional Operational Programme, European Regional Development Fund).

References

- [1] L. Y. Alex Zihao Zhu. Ev-flownet: Self-supervised optical flow estimation for event-based cameras. *Robotics: Science and Systems*, Jan 2018. 2
- [2] P. Bardow, A. J. Davison, and S. Leutenegger. Simultaneous optical flow and intensity estimation from an event camera. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 884–892, 2016. 1
- [3] R. Berner, C. Brandli, M. Yang, S.-C. Liu, and T. Delbruck. A 240×180 10mw 12us latency sparse-output vision sensor for mobile applications. pages C186–C187, 01 2013. 1
- [4] M. Cannici, M. Ciccone, A. Romanoni, and M. Matteucci. Attention mechanisms for object recognition with event-based cameras. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1127–1136, Jan 2019. 1
- [5] Y. Cao, Y. Chen, and D. Khosla. Spiking deep convolutional neural networks for energy-efficient object recognition. *International Journal of Computer Vision*, 113(1):54–66, 2015. 1
- [6] N. F. Y. Chen. Pseudo-labels for Supervised Learning on Dynamic Vision Sensor Data, Applied to Object Detection under Ego-motion. *arXiv*, Sep 2017. 1
- [7] G. K. Cohen. *Event-Based Feature Detection, Recognition and Classification*. PhD thesis, Université Pierre et Marie Curie - Paris VI, Sep 2016. 2
- [8] P. U. Diehl, D. Neil, J. Binas, M. Cook, S. Liu, and M. Pfeiffer. Fast-classifying, high-accuracy spiking deep networks through weight and threshold balancing. In *2015 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, July 2015. 1
- [9] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The Pascal Visual Object Classes (VOC) Challenge. *Int. J. Comput. Vision*, 88(2):303–338, Jun 2010. 6
- [10] L. Fei-Fei, R. Fergus, and P. Perona. One-shot learning of object categories. *IEEE Trans. Pattern Anal. Mach. Intell.*, 28(4):594–611, Apr 2006. 5
- [11] D. Gehrig, H. Rebecq, G. Gallego, and D. Scaramuzza. Asynchronous, photometric feature tracking using events and frames. In *Eur. Conf. Comput. Vis.(ECCV)*, 2018. 1
- [12] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. *PMLR*, pages 249–256, Mar 2010. 6
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1
- [14] S. Kim, S. Park, B. Na, and S. Yoon. Spiking-yolo: Spiking neural network for real-time object detection. *arXiv preprint arXiv:1903.06530*, 2019. 1
- [15] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *arXiv*, Dec 2014. 6
- [16] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017. 1
- [17] A. Krizhevsky. Learning multiple layers of features from tiny images. 04 2009. 5
- [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems* 25, pages 1097–1105. Curran Associates, Inc., 2012. 1
- [19] X. Lagorce, G. Orchard, F. Galluppi, B. E. Shi, and R. B. Benosman. HOTS: A Hierarchy of Event-Based Time-Surfaces for Pattern Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(7), Jul 2016. 2
- [20] X. Lagorce, G. Orchard, F. Galluppi, B. E. Shi, and R. B. Benosman. Hots: a hierarchy of event-based time-surfaces for pattern recognition. *IEEE transactions on pattern analysis and machine intelligence*, 39(7):1346–1359, 2017. 1, 2
- [21] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, Nov 1998. 5, 6
- [22] H. Li, H. Liu, X. Ji, G. Li, and L. Shi. CIFAR10-DVS: An Event-Stream Dataset for Object Classification. *Front. Neurosci.*, 11:309, May 2017. 5
- [23] J. Li, F. Shi, W. Liu, D. Zou, Q. Wang, H. Lee, P.-K. Park, and H. E. Ryu. Adaptive temporal pooling for object detection using dynamic vision sensor. *British Machine Vision Conference (BMVC)*, 2017. 1
- [24] M. Liu and T. Delbruck. Adaptive time-slice block-matching optical flow algorithm for dynamic vision sensors. Technical report, 2018. 1
- [25] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. 1
- [26] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 1
- [27] W. Maass. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 10(9):1659–1671, 1997. 1
- [28] A. I. Maqueda, A. Loquercio, G. Gallego, N. Garca, and D. Scaramuzza. Event-based vision meets deep learning on steering prediction for self-driving cars. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2
- [29] A. Mitrokhin, C. Fermuller, C. Parameshwara, and Y. Aloimonos. Event-based moving object detection and tracking. *arXiv preprint arXiv:1803.04523*, 2018. 1
- [30] V. Mnih, N. Heess, A. Graves, et al. Recurrent models of visual attention. In *Advances in neural information processing systems*, pages 2204–2212, 2014. 6
- [31] F. Monti, D. Boscaini, J. Masci, E. Rodola, J. Svoboda, and M. M. Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5115–5124, 2017. 1

- [32] E. Mueggler, H. Rebecq, G. Gallego, T. Delbruck, and D. Scaramuzza. The Event-Camera Dataset and Simulator: Event-based Data for Pose Estimation, Visual Odometry, and SLAM. *arXiv*, Oct 2016. 1, 6
- [33] D. Neil, M. Pfeiffer, and S.-C. Liu. Phased lstm: Accelerating recurrent network training for long or event-based sequences. In *Advances in Neural Information Processing Systems*, pages 3882–3890, 2016. 1
- [34] G. Orchard, A. Jayawant, G. K. Cohen, and N. Thakor. Converting Static Image Datasets to Spiking Neuromorphic Datasets Using Saccades. *Front. Neurosci.*, 9, Nov 2015. 5, 6
- [35] J. A. Pérez-Carrasco, B. Zhao, C. Serrano, B. Acha, T. Serrano-Gotarredona, S. Chen, and B. Linares-Barranco. Mapping from frame-driven to frame-free event-driven vision systems by low-rate rate coding and coincidence processing—application to feedforward ConvNets. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(11):2706–2719, Nov 2013. 1, 2
- [36] C. Posch, 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. *IEEE J. Solid-State Circuits*, 46(1):259–275, Jan 2011. 1
- [37] A. Raj, D. Maturana, and S. Scherer. Multi-scale convolutional architecture for semantic segmentation. page 14, 01 2015. 1
- [38] B. Ramesh, H. Yang, G. Orchard, N. A. L. Thi, and C. Xiang. DART: Distribution Aware Retinal Transform for Event-based Cameras. *arXiv*, Oct 2017. 1
- [39] B. Ramesh, S. Zhang, Z. W. Lee, Z. Gao, G. Orchard, and C. Xiang. Long-term object tracking with a moving event camera. 2018. 1
- [40] H. Rebecq, T. Horstschaefer, G. Gallego, and D. Scaramuzza. Evo: A geometric approach to event-based 6-dof parallel tracking and mapping in real time. *IEEE Robotics and Automation Letters*, 2(2):593–600, 2017. 1
- [41] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. 1, 2, 6, 7
- [42] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 1
- [43] T. Serrano-Gotarredona and B. Linares-Barranco. A 128×128 1.5 % Contrast Sensitivity 0.9 % FPN $3 \mu\text{s}$ Latency 4 mW Asynchronous Frame-Free Dynamic Vision Sensor Using Transimpedance Preamplifiers. *IEEE J. Solid-State Circuits*, 48(3):827–838, Mar 2013. 1
- [44] T. Serrano-Gotarredona and B. Linares-Barranco. Poker-DVS and MNIST-DVS. Their History, How They Were Made, and Other Details. *Front. Neurosci.*, 9, Dec 2015. 5
- [45] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 1, 6
- [46] A. Sironi, M. Brambilla, N. Bourdis, X. Lagorce, and R. Benosman. Hats: Histograms of averaged time surfaces for robust event-based object classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1731–1740, 2018. 1, 2
- [47] T. Stoffregen and L. Kleeman. Simultaneous optical flow and segmentation (sofas) using dynamic vision sensor. *arXiv preprint arXiv:1805.12326*, 2018. 1
- [48] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*, volume 4, page 12, 2017. 1
- [49] F. Yu and V. Koltun. Multi-Scale Context Aggregation by Dilated Convolutions. In *International Conference on Learning Representations (ICLR)*, 2016. 1