



STDP-based spiking deep convolutional neural networks for object recognition

Saeed Reza Kheradpisheh ^{a,b,*}, Mohammad Ganjtabesh ^a, Simon J. Thorpe ^b,
Timothée Masquelier ^b

^a Department of Computer Science, School of Mathematics, Statistics, and Computer Science, University of Tehran, Tehran, Iran

^b CERCO UMR 5549, CNRS –Université Toulouse 3, France

ARTICLE INFO

Article history:

Received 7 May 2017

Received in revised form 23 November 2017

Accepted 8 December 2017

Available online 23 December 2017

Keywords:

Spiking neural network

STDP

Deep learning

Object recognition

Temporal coding

ABSTRACT

Previous studies have shown that spike-timing-dependent plasticity (STDP) can be used in spiking neural networks (SNN) to extract visual features of low or intermediate complexity in an unsupervised manner. These studies, however, used relatively shallow architectures, and only one layer was trainable. Another line of research has demonstrated – using rate-based neural networks trained with back-propagation – that having many layers increases the recognition robustness, an approach known as deep learning. We thus designed a deep SNN, comprising several convolutional (trainable with STDP) and pooling layers. We used a temporal coding scheme where the most strongly activated neurons fire first, and less activated neurons fire later or not at all. The network was exposed to natural images. Thanks to STDP, neurons progressively learned features corresponding to prototypical patterns that were both salient and frequent. Only a few tens of examples per category were required and no label was needed. After learning, the complexity of the extracted features increased along the hierarchy, from edge detectors in the first layer to object prototypes in the last layer. Coding was very sparse, with only a few thousands spikes per image, and in some cases the object category could be reasonably well inferred from the activity of a single higher-order neuron. More generally, the activity of a few hundreds of such neurons contained robust category information, as demonstrated using a classifier on Caltech 101, ETH-80, and MNIST databases. We also demonstrate the superiority of STDP over other unsupervised techniques such as random crops (HMAX) or auto-encoders. Taken together, our results suggest that the combination of STDP with latency coding may be a key to understanding the way that the primate visual system learns, its remarkable processing speed and its low energy consumption. These mechanisms are also interesting for artificial vision systems, particularly for hardware solutions.

© 2017 Elsevier Ltd. All rights reserved.

1. Introduction

Primate's visual system solves the object recognition task through hierarchical processing along the ventral pathway of the visual cortex (DiCarlo, Zoccolan, & Rust, 2012). Through this hierarchy, the visual preference of neurons gradually increases from oriented bars in primary visual cortex (V1) to complex objects in inferotemporal cortex (IT), where neural activity provides a robust, invariant, and linearly-separable object representation (DiCarlo & Cox, 2007; DiCarlo et al., 2012). Despite the extensive feedback connections in the visual cortex, the first feed-forward wave of spikes in IT ($\sim 100\text{--}150$ ms post-stimulus presentation) appears to be sufficient for crude object recognition (Hung, Kreiman, Poggio,

& DiCarlo, 2005; Liu, Agam, Madsen, & Kreiman, 2009; Thorpe, Fize, Marlot, et al., 1996).

During the last decades, various computational models have been proposed to mimic this hierarchical feed-forward processing (Fukushima, 1980; LeCun & Bengio, 1998; Lee, Grosse, Ranganath, & Ng, 2009; Masquelier & Thorpe, 2007; Serre, Wolf, Bileschi, Riesenhuber, & Poggio, 2007). Despite the limited successes of the early models (Ghodrati, Farzamhd, Rajaei, Ebrahimpour, & Khaligh-Razavi, 2014; Pinto, Barhom, Cox, & DiCarlo, 2011), recent advances in deep convolutional neural networks (DCNN) led to high performing models (Krizhevsky, Sutskever, & Hinton, 2012; Simonyan & Zisserman, 2014; Zeiler & Fergus, 2014). Beyond the high precision, DCNNs can tolerate object variations as humans do (Kheradpisheh, Ghodrati, Ganjtabesh, & Masquelier, 2016a, b), use IT-like object representations (Cadieu et al., 2014; Khaligh-Razavi & Kriegeskorte, 2014), and match the spatio-temporal dynamics of the ventral visual pathway (Cichy, Khosla, Pantazis, Torralba, & Oliva, 2016).

* Corresponding author at: Department of Computer Science, School of Mathematics, Statistics, and Computer Science, University of Tehran, Tehran, Iran.

E-mail addresses: kheradpisheh@ut.ac.ir (S.R. Kheradpisheh), mgtabesh@ut.ac.ir (M. Ganjtabesh), simon.thorpe@cnrs.fr (S.J. Thorpe), timothee.masquelier@cnrs.fr (T. Masquelier).

Although the architecture of DCNNs is somehow inspired by the primate's visual system (LeCun, Bengio, & Hinton, 2015) (a hierarchy of computational layers with gradually increasing receptive fields), they totally neglect the actual neural processing and learning mechanisms in the cortex.

The computing units of DCNNs send floating-point values to each other which correspond to their activation level, while, biological neurons communicate to each other by sending electrical impulses (i.e., spikes). The amplitude and duration of all spikes are almost the same, so they are fully characterized by their emission time. Interestingly, mean spike rates are very low in the primate visual systems (perhaps only a few of hertz Shoham, O'Connor, & Segev, 2006). Hence, neurons appear to fire a spike only when they have to send an important message, and some information can be encoded in their spike times. Such spike-time coding leads to a fast and extremely energy-efficient neural computation in the brain (the whole human brain consumes only about 10–20 W of energy Maass, 2002).

The current top-performing DCNNs are trained with the supervised back-propagation algorithm which has no biological root. Although it works well in terms of accuracy, the convergence is rather slow because of the credit assignment problem (Rolls & Deco, 2002). Furthermore, given that DCNNs typically have millions of free parameters, millions of labeled examples are needed to avoid over-fitting. However, primates, especially humans, can learn from far fewer examples while most of the time no label is available. They may be able to do so thanks to spike-timing-dependent plasticity (STDP), an unsupervised learning mechanism which occurs in mammalian visual cortex (Huang et al., 2014; McMahon & Leopold, 2012; Meliza & Dan, 2006). According to STDP, synapses through which a presynaptic spike arrived before (respectively after) a postsynaptic one are reinforced (respectively depressed).

To date, various spiking neural networks (SNN) have been proposed to solve object recognition tasks. A group of these networks are actually the converted versions of traditional DCNNs (Cao, Chen, & Khosla, 2015; Diehl, Zarrella, Cassidy, Pedroni, & Neftci, 2016; Hunsberger & Eliasmith, 2015). The main idea is to replace each DCNN computing unit with a spiking neuron whose firing rate is correlated with the output of that unit. The aim of these networks is to reduce the energy consumption in DCNNs. However, the inevitable drawbacks of such spike-rate coding are the need for many spikes per image and the long processing time. Besides, the use of back-propagation learning algorithm and having both positive (excitatory) and negative (inhibitory) output synapses in a neuron are not biologically plausible. On the other hand, there are SNNs which are originally spiking networks and learn spike patterns. First group of these networks exploit learning methods such as auto-encoder (Burbank, 2015; Panda & Roy, 2016) and back-propagation (Bengio, Lee, Bornschein, & Lin, 2015) which are not biologically plausible. The second group consists of SNNs with bioinspired learning rules which have shallow architectures (Brader, Senn, & Fusi, 2007; Diehl & Cook, 2015; Habenschuss, Bill, & Nessler, 2012; Querlioz, Bichler, Dollfus, & Gamrat, 2013; Zhao, Ding, Chen, Linares-Barranco, & Tang, 2015) or only one trainable layer (Beyeler, Dutt, & Krichmar, 2013; Kheradpisheh, Ganjtabesh, & Masquelier, 2016; Masquelier & Thorpe, 2007).

In this paper we proposed a STDP-based spiking deep neural network (SDNN) with a spike-time neural coding. The network is comprised of a temporal-coding layer followed by a cascade of consecutive convolutional (feature extractor) and pooling layers. The first layer converts the input image into an asynchronous spike train, where the visual information is encoded in the temporal order of the spikes. Neurons in convolutional layers integrate input spikes, and emit a spike right after reaching their threshold. These

layers are equipped with STDP to learn visual features. Pooling layers provide translation invariance and also compact the visual information (Serre, Wolf et al., 2007). Through the network, visual features get larger and more complex, where neurons in the last convolutional layer learn and detect object prototypes. At the end, a classifier detects the category of the input image based on the activity of neurons in the last pooling layer with global receptive fields.

We evaluated the proposed SDNN on Caltech face/motorbike and ETH-80 datasets with large-scale images of various objects taken from different viewpoints. The proposed SDNN reached the accuracies of 99.1% on face/motorbike task and 82.8% on ETH-80, which indicates its capability to recognize several natural objects even under severe variations. Based on our knowledge, there is no other spiking deep network which can recognize large-scale natural objects. We also examined the proposed SDNN on the MNIST dataset which is a benchmark for spiking neural networks, and interestingly, it reached 98.4% recognition accuracy. In addition to the high performance, the proposed SDNN is highly energy-efficient and works with a few number of spikes per image, which makes it suitable for neuromorphic hardware implementation. Although current state-of-the-art DCNNs achieved stunning results on various recognition tasks, continued work on brain-inspired models could end up in strong intelligent systems in future, which can even help us to improve our understanding of the brain itself.

2. Proposed spiking deep neural network

A sample architecture of the proposed SDNN with three convolutional and three pooling layers is shown in Fig. 1. Note that the architectural properties (e.g., the number of layers and receptive field sizes) and learning parameters should be optimized for the desired recognition task.

The first layer of the network uses Difference of Gaussians (DoG) filters to detect contrasts in the input image. It encodes the strength of these contrasts in the latencies of its output spikes (the higher the contrast, the shorter the latency). Neurons in convolutional layers detect more complex features by integrating input spikes from the previous layer which detects simpler visual features. Convolutional neurons emit a spike as soon as they detect their preferred visual feature which depends on their input synaptic weights. Through the learning, neurons that fire earlier perform the STDP and prevent the others from firing via a winner-take-all mechanism. In this way, more salient and frequent features tend to be learned by the network. Pooling layers provide translation invariance using maximum operation, and also help the network to compress the flow of visual data. Neurons in pooling layers propagate the first spike received from neighboring neurons in the previous layer which are selective to the same feature. Convolutional and pooling layers are arranged in a consecutive order. Receptive fields gradually increase through the network and neurons in higher layers become selective to complex objects or object parts.

It should be noted that the internal potentials of all neurons are reset to zero before processing the next image. Also, learning only happens in convolutional layers and it is done layer by layer. Since the calculations of each neuron is independent of other adjacent neurons, to speed-up the computations, each of the convolution, pooling, and STDP operations are performed in parallel on GPU.

2.1. DoG and temporal coding

The important role of the first stage in SNNs is to encode the input signal into discrete spike events in the temporal domain. This temporal coding determines the content and the amount of information carried by each spike, which deeply affects the neural

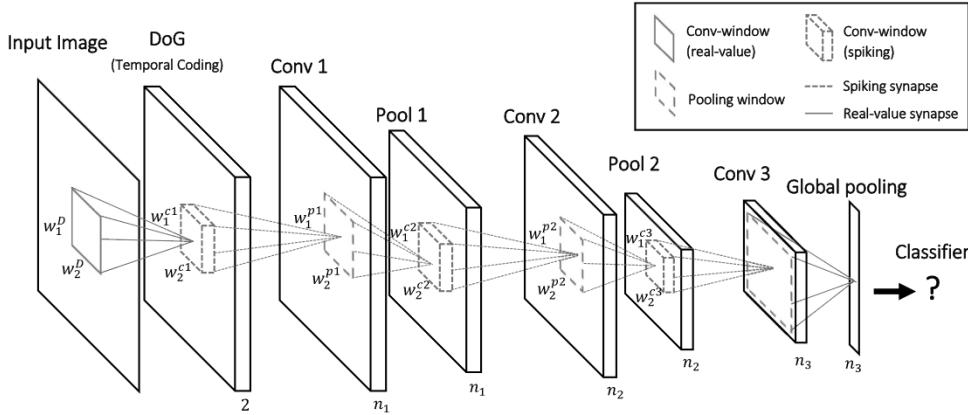


Fig. 1. A sample architecture of the proposed SDNN with three convolutional and three pooling layers. The first layer applies ON- and OFF-center DoG filters of size $w_1^D \times w_2^D$ on the input image and encode the image contrasts in the timing of the output spikes. The i th convolutional layer, Conv i , learns combinations of features extracted in the previous layer. The i th pooling layer, Pool i , provides translation invariance for features extracted in the previous layer and compress the visual information using a local maximum operation. Finally the classifier detects the object category based on the feature values computed by the global pooling layer. The window size of the i th convolutional and pooling layers are indicated by $w_{1,2}^{ci}$ and $w_{1,2}^{pi}$, respectively. The number of the neuronal maps of the i th convolutional and pooling layer is also indicated by n_i below each layer.

computations in the network. Hence, using efficient coding scheme in SNNs can lead to fast and accurate responses. Various temporal coding schemes can be used in visual processing (see Ref. Thorpe, Delorme, & Van Rullen, 2001). Among them, rank-order coding is shown to be efficient for rapid processing (even possibly in retinal ganglion cells) (Portelli et al., 2016; Van Rullen & Thorpe, 2001).

Cells in the first layer of the network apply a DoG filter over their receptive fields to detect positive or negative contrasts in the input image. DoG well approximates the center-surround properties of the ganglion cells of the retina. When presented with an image, these DoG cells detect the contrasts and emit a spike; the more strongly a cell is activated (higher contrast), the earlier it fires. In other word, the order of the spikes depends on the order of the contrasts. This rank-order coding is shown to be efficient for obtaining V1 like edge detectors (Delorme, Perrinet, & Thorpe, 2001) as well as complex visual features (Kheradpisheh, Ganjtabesh et al., 2016; Masquelier & Thorpe, 2007) in higher cortical areas.

DoG cells are retinotopically organized in two ON-center and OFF-center maps which are respectively sensitive to positive and negative contrasts. A DoG cell is allowed to fire if its activation is above a certain threshold. Note that this scheme guarantees that at most one of the two cells (positive or negative) can fire in each location. As mentioned above, the firing time of a DoG cell is inversely proportional to its activation value. In other words, if the output of the DoG filter at a certain location is r , the firing time of the corresponding cell is $t = 1/r$. For efficient GPU-based parallel computing, the input spikes are grouped into equal-size sequential packets. At each time step, spikes of one packet are propagated simultaneously. In this way, a packet of spikes with near ranks (carrying similar visual information) is propagated in parallel, while, the next spike packet will be processed in the next time step.

2.2. Convolutional layers

A convolutional layer contains several neuronal maps. Each neuron is selective to a visual feature determined by its input synaptic weights. Neurons in a specific map detect the same visual feature but at different locations. To this end, synaptic weights of neurons belonging to the same map should always be the same (i.e., weight sharing). Within a map, neurons are retinotopically arranged. Each neuron receives input spikes from the neurons located in a determined window in all neuronal maps of the

previous layer. Hence, a visual feature in a convolutional layer is a combination of several simpler feature extracted in the previous layer. Note that the input windows of two adjacent neurons are highly overlapped. Hence, the network can detect the appearance of the visual features in any location.

Neurons in all convolutional layers are non-leaky integrate-and-fire neurons, which gather input spikes from presynaptic neurons and emit a spike when their internal potentials reach a prespecified threshold. Each presynaptic spike increases the neuron's potential by its synaptic weight. At each time step, the internal potential of the i th neuron is updated as follows:

$$V_i(t) = V_i(t-1) + \sum_j W_{j,i} S_j(t-1), \quad (1)$$

where $V_i(t)$ is the internal potential of the i th convolutional neuron at time step t , $W_{j,i}$ is the synaptic weight between the j th presynaptic neuron and the i th convolutional neuron, and S_j is the spike train of the j th presynaptic neuron ($S_j(t-1) = 1$ if the neuron has fired at time $t-1$, and $S_j(t-1) = 0$ otherwise). If V_i exceeds its threshold, V_{thr} , then the neuron emits a spike and V_i is reset:

$$V_i(t) = 0 \text{ and } S_i(t) = 1, \quad \text{if } V_i(t) \geq V_{thr}. \quad (2)$$

Also, there is a lateral inhibition mechanism in all convolutional layers. When a neuron fires, in a specific location, it inhibits other neurons in that location belonging to other neuronal maps (i.e., resets their potentials to zero) and does not allow them to fire until the next image is shown. In addition, neurons are not allowed to fire more than once. These together provide a sparse but highly informative coding, because, there can be at most one spike at each location which indicates the existence of a particular visual feature in that location.

2.3. Local pooling layers

Pooling layers help the network to gain invariance by doing a nonlinear max pooling operation over a set of neighboring neurons with the same preferred feature. Some evidence suggests that such a max operation occurs in complex cells in visual cortex (Serre, Wolf et al., 2007). Thanks to the rank-order coding used in the proposed network, the maximum operation of pooling layers simply consists of propagating the first spike emitted by the afferents (Rousselet, Thorpe, & Fabre-Thorpe, 2003).

A neuron in a neuronal map of a pooling layer performs the maximum operation over a window in the corresponding neuronal map of the previous layer. Pooling neurons are integrate-and-fire neurons whose input synaptic weights and threshold are all set to one. Hence, the first input spike activates them and leads to an output spike. Regarding to the rank-order coding, each pooling neuron is allowed to fire at most once. It should be noted that no learning occurs in pooling layers.

Another important role of pooling layers is to compress the visual information. Regarding to the maximum operation performed in pooling layers, adjacent neurons with overlapped inputs would carry redundant information (each spike is sent to many neighboring pooling neurons). Hence, in the proposed network, the overlap between the input windows of two adjacent pooling neurons (belonging to the same map) is set to be very small. It helps to compress the visual information by eliminating the redundancies, and also, to reduce the size of subsequent layers.

2.4. STDP-based learning

As mentioned above, learning occurs only in convolutional layers which should learn to detect visual features by combining simpler features extracted in the previous layer. The learning is done layer by layer, i.e., the learning in a convolutional layer starts when the learning in the previous convolutional layer is finalized. When a new image is presented, neurons of the convolutional layer compete with each other and those which fire earlier trigger STDP and learn the input pattern.

A simplified version of STDP (Masquelier & Thorpe, 2007) is used:

$$\Delta w_{ij} = \begin{cases} a^+ w_{ij}(1 - w_{ij}), & \text{if } t_j - t_i \leq 0, \\ a^- w_{ij}(1 - w_{ij}), & \text{if } t_j - t_i > 0, \end{cases} \quad (3)$$

where i and j respectively refer to the index of post- and presynaptic neurons, t_i and t_j are the corresponding spike times, Δw_{ij} is the synaptic weight modification, and a^+ and a^- are two parameters specifying the learning rate. Note that the exact time difference between two spikes does not affect the weight change, but only its sign is considered. Also, it is assumed that if a presynaptic neuron does not fire before the postsynaptic one, it will fire later. These simplifications are equivalent to assuming that the intensity-latency conversion of DoG cells compresses the whole spike wave in a relatively short time interval (say, 20–30 ms), so that all presynaptic spikes necessarily fall close to the postsynaptic spike time, and the time lags are negligible. The multiplicative term $w_{ij}(1 - w_{ij})$ ensures the weights remain in the range [0,1] and thus maintains all synapses in an excitatory mode in adding to implementing soft-bound effect.

Note that choosing large values for the learning parameters (i.e., a^+ and a^-) will decrease the learning memory, therefore, neurons would learn the last presented images and unlearn previously seen images. Also, choosing tiny values would slow down the learning process. At the beginning of the learning, when synaptic weights are random, neurons are not yet selective to any specific pattern and respond to many different patterns, therefore, the probability for a synapse to get depressed is higher than being potentiated. Hence, by setting a^- to be greater than a^+ , synaptic weights gradually decay insofar as neurons cannot reach their threshold to fire anymore. Therefore, a^+ is better to be greater than a^- , however, by setting a^+ to be much greater than a^- , neurons will tend to learn more than one pattern and respond to all of them. All in all, it is better to choose a^+ and a^- not too big and not too small, and it is better to set a^+ a bit greater than a^- .

During the learning of a convolutional layer, neurons in the same map, detecting the same feature in different locations, integrate input spikes and compete with each other to do the STDP.

The first neuron which reaches the threshold and fires, if any, is the winner (global intra-map competition). The winner triggers the STDP and updates its synaptic weights. As mentioned before, neurons in different locations of the same map have the same input synaptic weights (i.e., weight sharing) to be selective to the same feature. Hence, the winner neuron prevents other neurons in its own map to do STDP and duplicates its updated synaptic weights into them. Also, there is a local inter-map competition for STDP. When a neuron is allowed to do the STDP, it prevents the neurons in other maps within a small neighborhood around its location from doing STDP. This competition is crucial to encourage neurons of different maps to learn different features.

Because of the discretized time variable in the proposed model, it is probable that some competitor neurons fire at the same time step. One possible scenario is to pick one randomly and allow it to do STDP. But a better alternative is to pick the one which has the highest potential indicating higher similarity between its learned feature and input pattern.

Synaptic weights of convolutional neurons initiate with random values drawn from a normal distribution with the mean of $\mu = 0.8$ and STD of $\sigma = 0.05$. Note that by choosing a small μ , neurons would not reach their threshold to fire and will not learn anything. Also, by choosing a large σ , some initial synaptic weights will be smaller (larger) than others and have less (more) contribution to the neuron activity, and regarding the STDP rule, they have a higher tendency to converge to zero (one). In other words, dependency on the initial weights will be higher for a large σ .

As the learning of a specific layer progresses, its neurons gradually converge to different visual features which are frequent in the input images. As mentioned before, learning in the subsequent convolutional layer starts whenever the learning in the current convolutional layer is finalized. Here we measure the learning convergence of the l th convolutional layer as

$$C_l = \sum_f \sum_i w_{f,i}(1 - w_{f,i})/n_w \quad (4)$$

where, $w_{f,i}$ is the i th synaptic weight of the f th feature and n_w is the total number of synaptic weights (independent of the features) in that layer. C_l tends to zero if each of the synaptic weights converge toward zero or one. Therefore, we stop the learning of the l th convolutional layer, whenever C_l was sufficiently close to zero (i.e. $C_l < 0.01$).

2.5. Global pooling and classification

The global pooling layer is only used in the classification phase. Neurons of the last layer perform a global max pooling over their corresponding neuronal maps in the last convolutional layer. Such a pooling operation provides a global translation invariance for prototypical features extracted in the last convolutional layer. Hence, there is only one output value for each feature, which indicates the presence of that feature in the input image. The output of the global pooling layer over the training images is used to train a linear SVM classifier. In the testing phase, the test object image is processed by the network and the output of the global pooling layer is fed to the classifier to determine its category.

To compute the output of the global pooling layer, first, the threshold of neurons in the last convolutional layer were set to be infinite, and then, their final potentials (after propagating the whole spike train generated by the input image) were measured. These final potentials can be seen as the number of early spikes in common between the current input and the stored prototypes in the last convolutional layer. Finally, the global pooling neurons compute the maximum potential at their corresponding neuronal maps, as their output value.

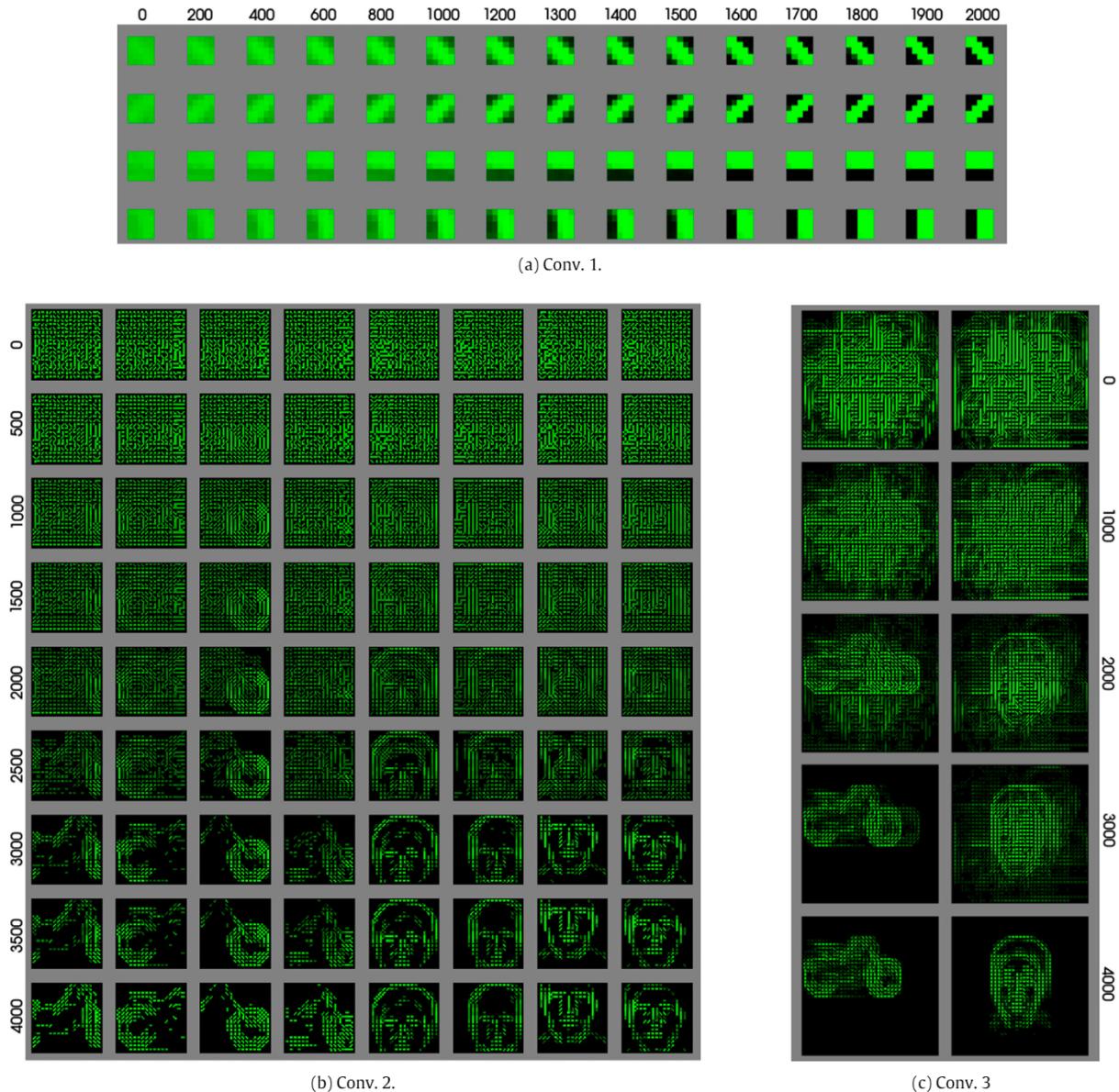


Fig. 2. The synaptic changes of some neuronal maps in different layers through the learning with the Caltech face/motorbike dataset. (A) The first convolutional layer becomes selective to oriented edges. (B) The second convolutional layer converges to object parts. (C) The third convolutional layer learns the object prototype and respond to whole objects.

3. Results

3.1. Caltech face/motorbike dataset

We evaluated our SDNN on the face and motorbike categories of the Caltech 101 dataset available at <http://www.vision.caltech.edu> (see Fig. 4 for sample pictures). The training set contains 200 randomly selected images per category, and remaining images constitute the test set. The test images are not seen during the learning phase but used afterward to evaluate the performance on novel images. This standard cross-validation procedure allows measuring the system's ability to generalize, as opposed to learning the specific training examples. All images were converted to grayscale values and rescaled to be 160 pixels in height (preserving the aspect ratio). In all the experiments, we used linear SVM classifiers with penalty parameter $C = 1.0$ (optimized by a grid search in the range of $(0, 10]$).

Here, we used a network similar to Fig. 1, with three convolutional layers each of which followed by a pooling layer. For the

first layer, only ON-center DoG filters of size 7×7 and standard deviations of 1 and 2 pixels are used. The first, second and third convolutional layers consists of 4, 20, and 10 neuronal maps with conv-window sizes of 5×5 , $16 \times 16 \times 4$, and $5 \times 5 \times 20$ and firing thresholds of 10, 60, and 2, respectively. The pooling window sizes of the first and second pooling layers are 7×7 and 2×2 with the strides of 6 and 2, correspondingly. The third pooling layer performs a global max pooling operation. The learning rates of all convolutional layers are set to $a^+ = 0.004$ and $a^- = 0.003$. In addition, each image is processed for 30 time steps.

Fig. 2 shows the preferred visual features of some neuronal maps in the first, second and third convolutional layers through the learning process. To visualize the visual feature learned by a neuron, a backward reconstruction technique is used. Indeed, the visual features in the current layer can be reconstructed as the weighted combinations of the visual features in the previous layer. This backward process continues until the first layer, whose preferred visual features are computed by DoG functions. As shown

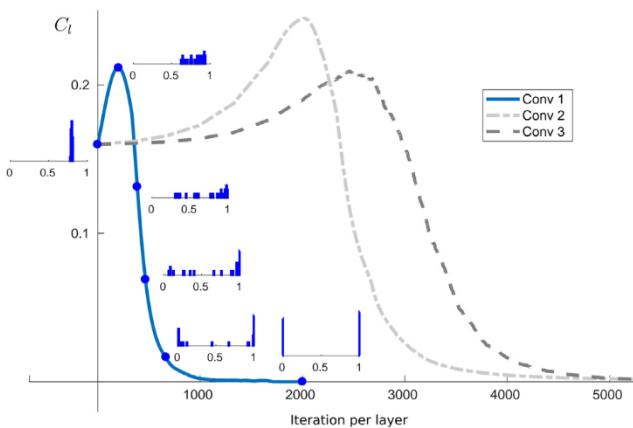


Fig. 3. Each curve shows the variation of the convergence index through the learning of a convolutional layer. The weight histogram of the first convolutional layer at some critical points during the learning is shown next to its convergence curve.

in Fig. 2A, interestingly, each of the four neuronal maps of the first convolutional layer converges to one of the four orientations: $\pi/4$, $\pi/2$, $3\pi/4$, and π . This shows how efficiently the association of the proposed temporal coding in DoG cells and unsupervised learning method (the STDP and learning competition) led to highly diverse edge detectors which can represent the input image with edges in different orientations. These edge detectors are similar to the simple cells in primary visual cortex (i.e., V1 area) (Delorme et al., 2001).

Fig. 2B shows the learning progress for the neuronal maps of the second convolutional layer. As mentioned, the first convolutional layer detects edges with different orientations all over the image, and due to the used temporal coding, neurons corresponding to edges with higher contrasts (i.e., salient edges) will fire earlier. On the other hand, STDP naturally tends to learn those combination of edges that are consistently repeating in the training images (i.e., common features between the target objects). Besides, the learning competition tends to prevent the neuronal maps from learning similar visual features. Consequently, neurons in the second convolutional layer learn the most salient, common, and diverse visual features of the target objects, and do not learn the backgrounds that drastically change between images. As seen in Fig. 2B, each of the maps gradually learns a different visual feature (combination of oriented edges) representing a face or motorbike feature.

The learning progress for two neuronal maps of the third convolutional layer is shown in Fig. 2C. As seen, one of them gradually becomes selective to a complete motorbike prototype as a combination of motorbike features such as back wheel, middle body, handle, and front wheel detected in the second layer. Also, the other map learns a whole face prototype as a combination of facial features. Indeed, the third convolutional layer learns the whole object prototypes using intermediate complexity features detected in the previous layer. Neurons in the second layer compete with each other and send spikes toward the third layer as they detect their preferred visual features. Since, different combinations of these features are detected for each object category, neuronal maps of the third layer will learn different prototypes of different categories. Therefore, the STDP and the learning competition mechanism direct neuronal maps of the third convolutional layer to learn highly category specific prototypes.

As mentioned in the previous section, learning at the l th convolutional layer stops whenever the learning convergence index C_l tends to zero. Here, in Fig. 3, we presented the evolution of C_l

during the learning of each convolutional layer. Also, we provided the weight histogram of the first convolutional layer at some critical points during the learning process. Since the initial synaptic weights are drawn from a normal distribution with $\mu = 0.8$ and $\sigma = 0.05$, at the beginning of the learning, the convergence index of the l th layer starts from $C_l \simeq 0.16$. Since features are not formed at the early iterations, neurons respond to almost every pattern, therefore many of the synapses are depressed most of the times and gradually move toward 0.5, and C_l peaks. As a consequence, neurons' activity decreases, and they start responding to a few of the patterns and not to others. From then, synapses contributed in these patterns are repeatedly potentiated and others are depressed. Due to the nature of the employed soft-bound STDP, learning is faster when the weights are around 0.5 and C_l rapidly decreases. Finally, at the end of the learning, as features are formed and synaptic weights converge to zero or one, C_l tends to zero.

As seen in Fig. 3, the weight changes in lower layers are faster than in higher layers. This is mainly due to the stronger competition among the feature maps in higher layers. For instance, in the first layer, there are only four feature maps with small receptive fields. If a neuron from one feature map loses the competition at one location, another neuron can still win the competition at some other location and do the STDP. While, in the third layer, there are more feature maps with larger receptive fields. If a neuron wins the competition, it prevents neurons of other feature maps from doing the STDP in a much wider area. Hence, for each image, only a few of the feature maps update their weights. Therefore, the competition is stronger in higher layers and learning takes a much longer time. Beside the competition factor, neurons in higher layers become selective to more complex and less frequent features so they do not reach their threshold as often as neurons in lower layers (e.g., all images contain edges).

Fig. 4 shows the accumulated spiking activity of the DoG and the following three convolutional layers over all time steps, for two face and motorbike sample images. For each layer, the preferred features of some neuronal maps with color coded borders are demonstrated on top, and their corresponding spiking activity are shown in panels below them. Each colored point inside a panel indicates the neuronal map of the neuron which has fired in that location at a time step. As seen, neurons of the DoG layer detect image contrasts, and edge detectors in the first convolutional layer detect the orientation of edges. Neurons in the second convolutional layer, which are selective to intermediate complexity features, detect their preferred visual feature by combining input spikes from edge detector cells in the first layer. Finally, the coincidence of these features activates neurons in the third convolutional layer which are selective to object prototypes. As seen, when a face (motorbike) image is presented, neurons in the face (motorbike) maps fire. To better illustrate the learning progress of all the layers as well as their spiking activity in the temporal domain, we prepared a short video (see Video 1).

As mentioned in the previous section, the output of the global pooling layer is used by a linear SVM classifier to specify the object category of the input images. We trained the proposed SDNN on training images and evaluated it over the test images, where the model reached the categorization accuracy of $99.1 \pm 0.2\%$. It shows how the object prototypes, learned in the highest layer, can well represent the object categories. Furthermore, we also calculated the single neuron accuracy. In more details, we separately computed the recognition accuracy of each neuron in the global pooling layer. Surprisingly, some single neurons reached an accuracy of 93%, and the mean accuracy was 89.8%. Hence, it can be said that single neurons in the highest layer are highly class specific, and different neurons carry complementary information which altogether provide robust object representations.

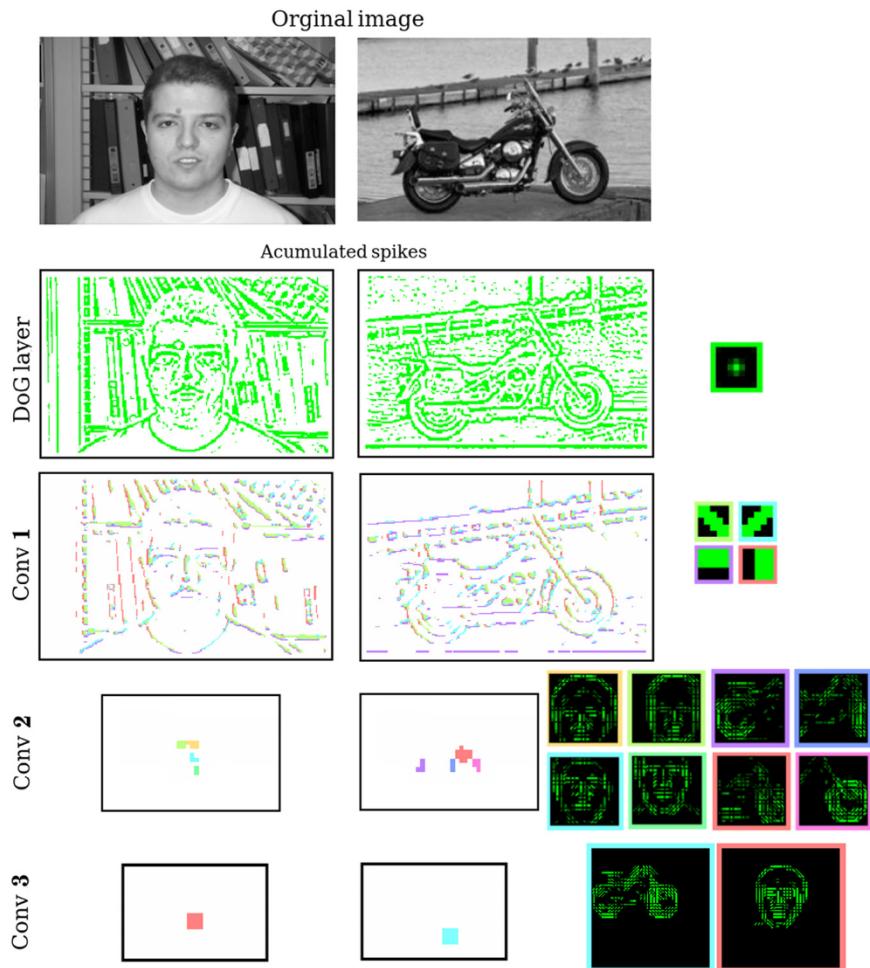


Fig. 4. The spiking activity of the convolutional layers with the face and motorbike images. The preferred features of neuronal maps in each convolutional layer are shown on the right. Each feature is coded by a specific color border. The spiking activity of the convolutional layers, accumulated over all the time steps, is shown in the corresponding panels. Each point in a panel indicates that a neuron in that location has fired at a time step, and the color of the point indicates the preferred feature of the activated neuron. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 1
Recognition accuracies of the proposed SDNN with random features in different convolutional layers.

Conv layers with random features	Non	3rd	2nd & 3rd	1st & 2nd & 3rd
Accuracy (%)	99.1	80.2	67.8	66.3

Table 2
Recognition accuracies of the proposed SDNN for different amounts of noise.

Noise level	$\alpha = 0\%$	$\alpha = 5\%$	$\alpha = 10\%$	$\alpha = 20\%$	$\alpha = 30\%$	$\alpha = 40\%$	$\alpha = 50\%$
Accuracy (%)	99.1	95.4	91.6	84.3	63.7	57.6	54.2

To better demonstrate the role of the STDP learning rule in the proposed SDNN, we used random features in different convolutional layers and assessed the final accuracy. To this end, we first trained all the three convolutional layers of the network (using STDP) until they all converged (synaptic weights had a bimodal distribution with 0 and 1 as centers). Then, for a certain convolutional layer, we counted the number of active synapses (i.e., close to one) of each of its feature maps. Corresponding to each learned feature in the first convolutional layer, we generated a random feature with the same number of active synapses. For the second and third convolutional layers, the number of active synapses in the random features was doubled. Note that with fewer active synapses, neurons in the second and third convolutional layers could not reach their threshold, and with more active synapses, many of the neurons tend to fire together. Anyways, we evaluated the

network using these random features. **Table 1** presents the SDNN accuracy when we had random features in the third, or second and third, or in all the three convolutional layers. As seen, by replacing the learned features of the lower layers with random ones, the accuracy decreases more. Especially, using random features in the second layer, the accuracy severely drops. This shows the critical role of intermediate complexity features for object categorization.

We also evaluated how the proposed SDNN is robust to noise. To this end, we added a white noise to the neurons' threshold during both training and testing phases. Indeed, for each image, we added to the threshold of each neuron, V_{thr} , a random value drawn from a uniform distribution in range $\pm \alpha\%$ of V_{thr} . We evaluated the proposed DCNN for different amount of noise (from $\alpha = 5\%$ to 50%) and the results are provided in **Table 2**. Up to the 20% of noise, the accuracy is still reasonable, but by increasing the noise level, the

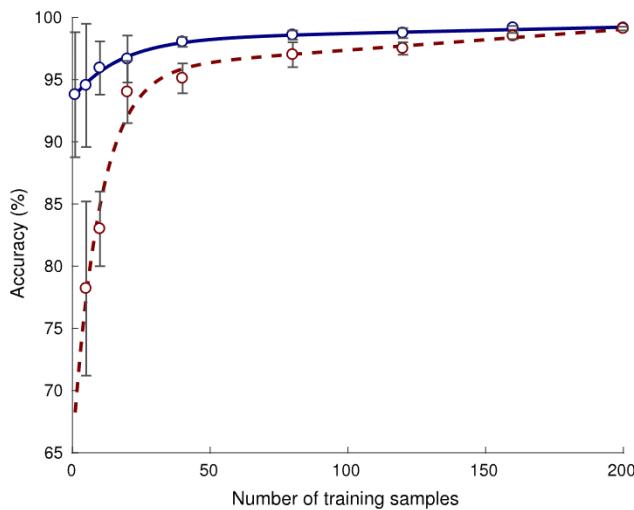


Fig. 5. Recognition accuracies (mean \pm std) of the proposed SDNN for different number of training images per category used in the STDP-based feature learning and/or training the classifier. The red curve presents the model's accuracy when the different number of images are used to train the network (by unsupervised STDP) and the classifier as well. The blue curve shows the model's accuracy when the layers of the network are trained using STDP with 400 images (200 from each category), and the classifier is trained with different number of labeled images per category.

accuracy dramatically drops and reaches to the chance level in case of 50% of noise. In other words, the network is more or less able to tolerate the instability caused by the noise below 20%, but when it goes further, the neurons' behavior drastically changes during the learning and STDP cannot extract informative features from the input images.

In another experiment, we changed the number of training samples and calculated the recognition accuracy of the proposed SDNN. For instance, we trained the network and the classifier with 5 samples from each category, and then evaluated the final system with the test samples. As shown by the red curve in Fig. 5, with 5 images per category the model reached the accuracy of 78.2%, and only 40 images from each category are sufficient to

reach 95.1% recognition accuracy. Although having more training samples leads to higher accuracies, the proposed SDNN can extract diagnostic features and reach reasonable accuracies even using a few tens of training images. Due to the unsupervised nature of STDP, the proposed SDNN does not suffer much from the overfitting challenge caused by small training set size in supervised learning algorithms such as back-propagation. In a real world, the number of labeled samples is very low. Therefore, learning in humans or other primates is mainly unsupervised. Here, in another experiment, we trained the network using STDP over 200 images per category, then we used different portions of these images as labeled samples (from 1 image to 200 images per category) to train the classifier. It lets us see whether the visual feature that the network has learned from the unlabeled images (using unsupervised STDP) is sufficient for solving the categorization task with only a few labeled training samples. As shown by the blue curve in Fig. 5, with only one sample per category the model could reach the average accuracy of 93.8%. It suggests that the unsupervised STDP can provide a rich feature space well explaining the object space and reducing the need for labeled examples.

3.2. ETH-80 dataset

The ETH-80 dataset contains eight different object categories: apple, car, toy cow, cup, toy dog, toy horse, pear, and tomato (10 instances per category). Each object is photographed from 41 view-points with different view angles and different tilts. Some examples of objects in this dataset are shown in Fig. 6. ETH-80 is a good benchmark to show how the proposed SDNN can handle multi-object categorization tasks with high inter-instance variability, and how it can tolerate huge view-point variations. In all the experiments, we used linear SVM classifiers with penalty parameter $C = 1.2$ (optimized by a grid search in the range of $(0, 10]$).

Five randomly chosen instances of each object category are selected for the training set used in the learning phase. The remaining instances constitute the testing set, and are not seen during the learning phase. All the object images were converted to grayscale values. To evaluate the proposed SDNN on ETH-80, we used a network architecturally similar to the one used for Caltech face/motorbike dataset. The other parameters are also similar,

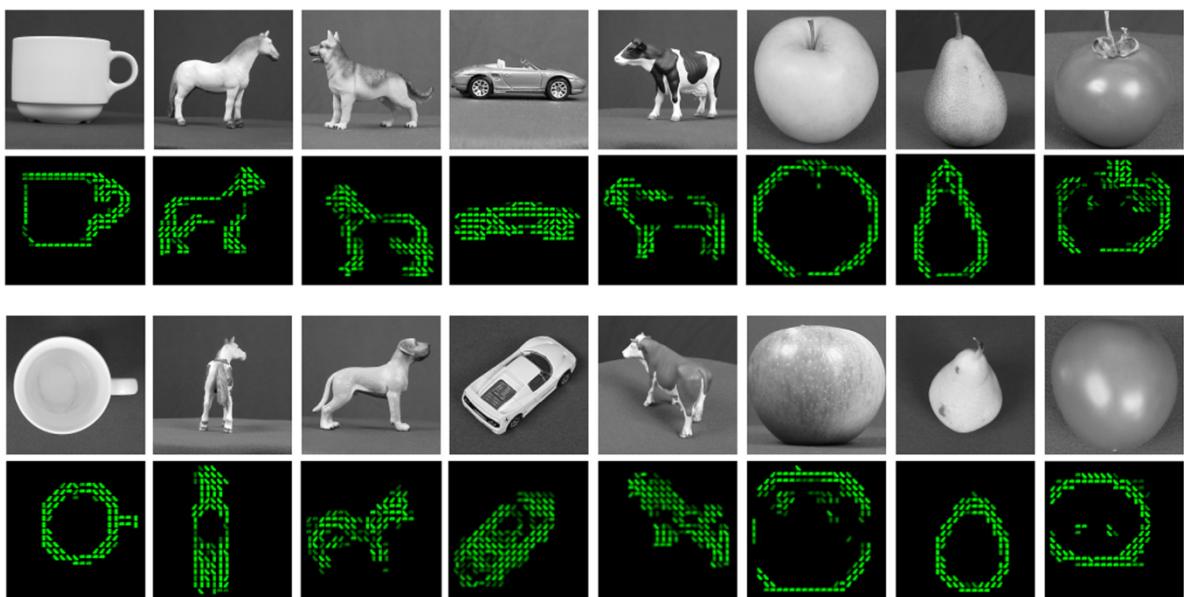


Fig. 6. Some sample images of different object categories of ETH-80 in different viewpoints. For each image, the preferred feature of an activated neuron in the third convolutional layer is shown in below.

Table 3

Recognition accuracies of the proposed SDNN and some other methods over the ETH-80 dataset. Note that all the models are trained on 5 object instances of each category and tested on other 5 instances.

Method	Accuracy (%)
HMAX (Kheradpisheh, Ganjtabesh et al., 2016)	69.0
Convolutional SNN (Kheradpisheh, Ganjtabesh et al., 2016)	81.1
Pre-trained AlexNet	79.5
Fine-tuned AlexNet	94.2
Supervised DCNN	81.9
Unsupervised DCA	80.7
Proposed SDNN	82.8

except for the number of neuronal maps in the second and third convolutional and pooling layers. Here we used 400 neuronal maps in each of these layers.

Similar to the Caltech dataset, neuronal maps of the first convolutional layer converged to the four oriented edges. Neurons in the second and third convolutional layers also became selective to intermediate features and object prototypes, respectively. Fig. 6 shows sample images from the ETH-80 dataset and the preferred features of some neuronal maps in the third convolutional layer which are activated for those images. As seen, neurons in the highest layer respond to different views of different objects, and altogether, provide an invariant object representation. Thus, the network learns 2D and view-dependent prototypes of each object category to achieve 3D representations.

As mentioned before, we evaluated the proposed SDNN over the test instances of each object category which are not shown to the network during the training. The recognition accuracy of the proposed SDNN along with some other models on the ETH-80 dataset are presented in Table 3. HMAX (Serre, Oliva, & Poggio, 2007) is one of the classic computational models of the object recognition process in visual cortex. It has 5000 features and uses a linear SVM as the classifier (see Kheradpisheh, Ganjtabesh et al., 2016 for more details). Convolutional SNN (Kheradpisheh, Ganjtabesh et al., 2016) is an extension of the Masquelier and Thorpe (2007) which had one trainable layer with 1200 visual features and used linear SVM for classification. AlexNet is the first DCNN that significantly improved recognition accuracy on Imagenet Dataset. Here, we compared our proposed SDNN with both Imagenet pre-trained and ETH-80 fine-tuned versions of AlexNet. To obtain the accuracy of the pre-trained AlexNet, images were shown to the model and feature vectors of the last layer were used to train and test a linear SVM classifier. Also, to fine-tune the Imagenet pre-trained AlexNet on ETH-80, its decision layer was replaced by an eight-neuron decision layer and trained by the stochastic gradient descent learning algorithm.

Regarding the fact that the pre-trained AlexNet has not seen ETH-80 images and fine-tuned AlexNet has already trained by millions of images from Imagenet dataset, the comparison may not be fair enough. Therefore, we compared the proposed SDNN to a supervised DCNN with the same structure but having two extra dense and decision layers on top with 70 and 8 neurons, respectively. The DCNN was trained on the gray-scaled images of ETH-80 using the backpropagation algorithm. The ReLU and softmax activation functions were employed for the intermediate and decision layers, respectively. We used a cross-entropy loss function with L2 kernel regularization and L1 activity regularization terms. We also performed a 50% dropout regularization on the dense layer. The hyper-parameters such as learning rates, momentums, and regularization factors were optimized using a grid search. Also, we used an *early stopping strategy* to prevent the network from overfitting. Eventually, the supervised DCNN reached the average accuracy of 81.9% (see Table 3). Note that we tried to evaluate the DCNN with a dense layer of more than 70 neurons, but all

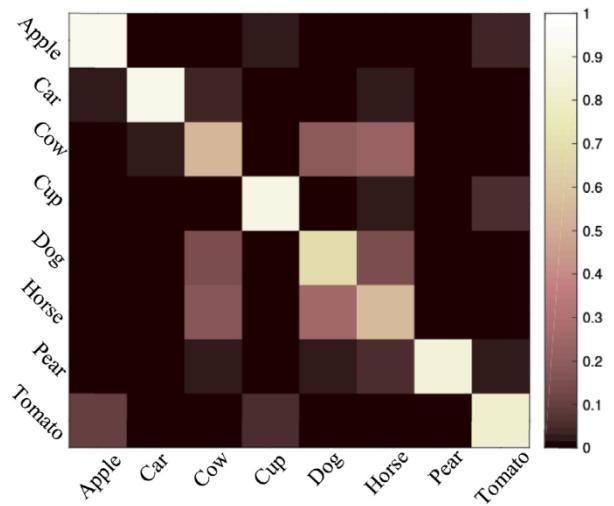


Fig. 7. The confusion matrix of the proposed SDNN over the ETH-80 dataset.

the time the network got quickly overfitted on the training data with no accuracy improvement over the test set. It seems that the supervised DCNN suffers from the lack of sufficient training data.

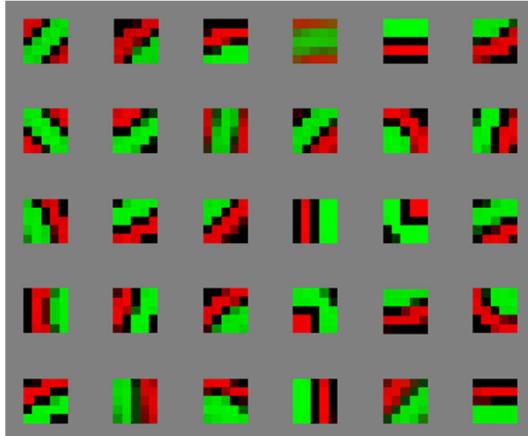
In addition, we compared the proposed SDNN to a deep convolutional autoencoder (DCA) which is one of the best unsupervised learning algorithms in machine learning. We developed a DCA with an encoder network having the same architecture as our SDNN followed by a decoder network with reversed architecture. The ReLU activation function was used in the convolutional layers of both encoder and decoder networks. We used the cross-entropy loss function and stochastic gradient descent learning algorithm to train the DCA. The learning parameters (i.e., learning rate and momentum) were optimized through a grid search. When the learning had converged, we eliminated the decoder part and used the encoder's output representations to train and test a linear SVM classifier. Note that DCA was trained on gray-scaled images of ETH-80. The DCA reached the average accuracy of 80.7% (see Table 3) and the proposed SDNN could outperform the DCA network with the same structure.

We also evaluated the proposed SDNN that has two convolutional layers. Indeed, we removed the third convolutional layer, applied the global pooling over the second convolutional layer, and trained a new SVM classifier over its output. The model's accuracy dropped by 5% and reached to 77.4%. Although the features in the second layer represent object parts with intermediate complexity, it is the combination of these features in the last layer which culminates the object representations and makes the classification easier. In other words, the similarity between the smaller parts of objects of different categories but with similar shapes (e.g., apple and potato) is higher than the whole objects, hence, the visual features in the higher layers can provide better object representations.

In a subsequent analysis, we computed the confusion matrix, to see which categories are mostly confused with each other. Fig. 7 illustrates the confusion matrix of the proposed SDNN over the ETH-80 dataset. As seen, most of the errors are due to the misclassification of dogs, horses, and cows. We checked whether these errors belong to the some specific viewpoints or not. We found out that the errors are uniformly distributed between different viewpoints. Some misclassified samples are shown in Fig. 8. Overall, it can be concluded that these categorization errors are due to the overall shape similarity between these object categories.

The other important aspect of the proposed SDNN is the computational efficiency of the network. For each ETH-80 image, on average, about 9100 spikes are emitted in all the layers, i.e., about 0.02

Images:								
Category:	Car	Apple	Cow	Cup	Cat	Horse	Pear	Tomato
Misclassified as:	Horse	Tomato	Dog	Apple	Horse	Cow	Tomato	Apple

Fig. 8. Some misclassified samples of the ETH-80 dataset by the proposed SDNN.**Fig. 9.** The Gabor-like features learned by the neuronal maps of the first convolutional layer from the MNIST images. The red and green colors receptively indicate the strength of input synapses from ON- and OFF-center DoG cells. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

spike per neuron per image. Note that the number of inhibitory events is equal to the number of spikes. These together points to the fact that the proposed SDNN can recognize objects with high precision but low computational cost. This efficiency is caused by the association of the proposed temporal coding and STDP learning rule which led to a sparse but informative visual coding.

3.3. MNIST dataset

MNIST (LeCun, Bottou, Bengio, & Haffner, 1998) is a benchmark dataset for SNNs which has been widely used (Diehl & Cook, 2015; Diehl et al., 2015; Hussain, Liu, & Basu, 2014; O'Connor, Neil, Liu, Delbrück, & Pfeiffer, 2013; Querlioz et al., 2013; Zhao et al., 2015). We also evaluated our SDNN on the MNIST dataset which contains 60,000 training and 10,000 test handwritten single-digit images. Each image is of size 28×28 pixels and contains one of the digits 0–9. For the first layer, ON- and OFF-center DoG filters with standard deviations of 1 and 2 pixels are used. The first and second convolutional layers respectively consist of 30 and 100 neuronal maps with 5×5 convolution-window and firing thresholds of 15 and 10. The pooling-window of the first pooling layer was of size 2×2 with the stride of 2. The second pooling layer performs a global max operation. Note that the learning rates of all convolutional layers were set to $a^+ = 0.004$ and $a^- = 0.003$. In all the experiments, we used linear SVM classifiers with penalty parameter $C = 2.4$ (optimized by a grid search in the range of $(0, 10]$).

Fig. 9 shows the preferred features of some neuronal maps in the first convolutional layer. The green and red colors correspond to ON- and OFF-center DoG filters. Interestingly, this layer converged to Gabor-like edge detectors with different orientations, phase and polarity. These edge features are combined in the next layer and provide easily separable digit representation.

Recognition performance of the proposed method and some recent SNNs on the MNIST dataset are provided in Table 4. As seen,

the proposed SDNN outperforms unsupervised SNNs by reaching 98.4% recognition accuracy. Besides, the accuracy of the proposed SDNN is close to the 99.1% accuracy of the totally supervised rate-based SDNN (Diehl et al., 2015) which is indeed the converted version of a traditional DCNN trained by back-propagation.

The important advantage of the proposed SDNN is the use of much fewer spikes. Our SDNN uses only about 600 spikes for each MNIST images in total for all the layers, while the supervised rate-based SDNN uses thousands of spikes per layer (Diehl et al., 2015). Also, because of using rate-based neural coding in such networks, they need to process images for hundreds of time steps, while our network processes the MNIST images in 30 time steps only. Notably, whenever a neuron in our network fires it inhibits other neurons at the same position, therefore, the total number of inhibitory events per each MNIST image is equal to the number of spikes(i.e., 600). As stated in Section 2, the proposed SDNN uses a temporal code which encodes the information of the input image in the spike times, and each neuron in all layers, is allowed to fire at most once. This temporal code associated with unsupervised STDP rule leads to a fast, accurate, and efficient processing.

4. Discussion

Recent supervised DCNNs have reached high accuracies on the most challenging object recognition datasets such as Imagenet. Architecture of these networks are largely inspired by the deep hierarchical processing in the visual cortex. For instance, DCNNs use retinotopically arranged neurons with restricted receptive fields, and the receptive field size and feature complexity of the neurons gradually increase through the layers. However, the learning and neural processing mechanisms applied in DCNNs are inconsistent with the visual cortex, where neurons communicate using spikes and learn the input spike patterns in a mainly unsupervised manner. Employing such mechanisms in DCNNs can improve their energy consumption and decrease their need for an expensive supervised learning with millions of labeled images.

A popular approach in previous researches is to convert pre-trained supervised DCNNs into equivalent spiking network. To simulate the floating-point calculations in DCNNs, they have to use the firing rate as the neural code, which in result increases the number of required spikes and the processing time. For instance, the converted version of a simple two-layer DCNN for the MNIST dataset with images of 28×28 pixels requires thousands of spikes and hundreds of time steps per image. On the other hand, there are some SNNs (Bengio et al., 2015; Burbank, 2015; Panda & Roy, 2016) which are originally spiking network but employ firing rate coding and biologically implausible learning rules such as autoencoders and back-propagation.

Here, we proposed a STDP-based SDNN with a spike-time coding. Each neuron was allowed to fire at most once, where its spike-time indicates the significance of its visual input. Therefore, neurons that fire earlier are carrying more salient visual information, and hence, they were allowed to do the STDP and learn the input patterns. As the learning progresses, each layer converged to a set of diverse but informative features, and the feature complexity gradually increases through the layers from simple edge features to object prototypes. The proposed SDNN was evaluated

Table 4

Recognition accuracies of the proposed SDNN and some other SNNs over the MNIST dataset.

Architecture	Neuralcoding	Learning-type	Learning-rule	Accuracy (%)
Dendritic neurons (Hussain et al., 2014)	Rate-based	Supervised	Morphology learning	90.3
Convolutional SNN (Zhao et al., 2015)	Spike-based	Supervised	Tempotron rule	91.3
Two layer network (Querlioz et al., 2013)	Spike-based	Unsupervised	STDP	93.5
Spiking RBM (O'Connor et al., 2013)	Rate-based	Supervised	Contrastive divergence	94.1
Two layer network (Diehl & Cook, 2015)	Spike-based	Unsupervised	STDP	95.0
Convolutional SNN (Diehl et al., 2015)	Rate-based	Supervised	Back-propagation	99.1
Proposed SDNN	Spike-based	Unsupervised	STDP	98.4

on several image datasets and reached high recognition accuracies. This shows how the proposed temporal coding and learning mechanism (STDP and learning competition) lead to discriminative object representations.

The proposed SDNN has several advantages to its counterparts. First, our proposed SDNN is the first spiking neural network with more than one learnable layer which can process large-scale natural object images. Second, due to the use of an efficient temporal coding, which encodes the visual information in the time of the first spikes, it can process the input images with a low number of spikes and in a few processing time steps. Third, the proposed SDNN exploits the bio-inspired and totally unsupervised STDP learning rule which can learn the diagnostic object features and neglect the irrelevant backgrounds.

We compared the proposed SDNN to several other networks including unsupervised methods such as HMAX and convolutional autoencoder network, and supervised methods such as DCNNs. The proposed SDNN could outperform the unsupervised methods which shows its advantages in extracting more informative features from training images. Also, it was better than the supervised deep network which largely suffered from overfitting and lack of sufficient training data. Although, the state-of-the-art supervised DCNNs have stunning performance on large datasets like Imagenet, but contrary to the proposed SDNN, they fall in trouble with small and medium size datasets.

Our SDNN could be efficiently implemented in parallel hardware (e.g., FPGA Yousefzadeh, Serrano-Gotarredona, & Linares-Barranco, 2015) using address event representation (AER) (Sivilotti, 1991) protocol. With AER, spike events are represented by the addresses of sending and receiving neurons, and time is represented by the asynchronous occurrence of spike events. Since these hardware are much faster than biological hardware, simulations could run several order of magnitude faster than real time (Serrano-Gotarredona, Masquelier, Prodromakis, Indiveri, & Linares-Barranco, 2013). The primate visual system extracts the rough content of an image in about 100 ms (Hung et al., 2005; Kirchner & Thorpe, 2006; Liu et al., 2009; Thorpe et al., 1996). We thus speculate that some dedicated hardware will be able to do the same in the order of a millisecond or less.

Also, the proposed SDNN can be modified to use spiking retinal models (Martínez-Cañada, Morillas, Pino, Ros, & Pelayo, 2016; Wohrer & Kornprobst, 2009) as the input layer. These models mimic the spatiotemporal filtering of the retinal ganglion cells with center/surround receptive fields. Alternatively, we could use neuromorphic asynchronous event-based cameras such as dynamic vision sensor (DVS), which generate output events when they capture transients in the scene (Lichtsteiner, Posch, & Delbrück, 2007). Finally, due to the DoG filtering in the input layer of the proposed SDNN, some visual information such as texture and color are lost. Hence, future studies should focus on encoding these additional pieces of information in the input layer.

Biological evidences indicate that in addition to the unsupervised learning mechanisms (e.g., STDP), there are also dopamine-based reinforcement learning strategies in the brain (Pignatelli & Bonci, 2015). Besides, although it is still unclear how supervised learning is implemented in biological neural networks, it seems

that for some tasks (e.g., motor control and sensory inputs prediction) the brain must constantly learn temporal dynamics based on error feedback (Doya, 2000). Employing such reinforcement and supervised learning strategies could improve the proposed SDNN in different aspects which are inevitable with unsupervised learning methods. Particularly, they can help to reduce the number of required features and to extract optimized task-dependent features.

5. Supporting information

Video 1. We prepared a video available at <https://youtu.be/u32Xnz2hDkE> showing the learning progress and neural activity over the Caltech face and motorbike task. Here we presented the face and motorbike training examples, propagated the corresponding spike waves, and applied the STDP rule. The input image is presented at the top-left corner of the screen. The output spikes of the input layer (i.e., DoG layer) at each time step is presented in the top-middle panel, and the accumulation of these spikes is shown in the top-right panel. For each of the subsequent convolutional layers, the preferred features, the output spikes at each time step, and the accumulation of the output spikes are presented in the corresponding panels. Note that 4, 8, and 2 features from the first, second and third convolutional layers are selected and shown, respectively. As mentioned, the learning occurs layer by layer, thus, the label of the layer which is currently doing the learning is specified by the red color. As seen, the first layer learns to detect edges, the second layer learns intermediate features, and finally the third layer learns face and motorbike prototype features.

Acknowledgments

This research received funding from the European Research Council under the European Unions 7th Framework Program (FP/2007–2013)/ ERC Grant Agreement no. 323711 (M4 project). The authors thank the NVIDIA Academic Programs team for donating a GPU hardware.

References

- Bengio, Y., Lee, D.-H., Bornschein, J., & Lin, Z. (2015). Towards biologically plausible deep learning, *arXiv:1502.04156*.
- Beyeler, M., Dutt, N. D., & Krichmar, J. L. (2013). Categorization and decision-making in a neurobiologically plausible spiking network using a stdp-like learning rule. *Neural Networks*, 48, 109–124.
- Brader, J. M., Senn, W., & Fusi, S. (2007). Learning real-world stimuli in a neural network with spike-driven synaptic dynamics. *Neural Computation*, 19(11), 2881–2912.
- Burbank, K. S. (2015). Mirrored stdp implements autoencoder learning in a network of spiking neurons. *PLoS Computational Biology*, 11(12), e1004566.
- Cadieu, C. F., Hong, H., Yamins, D. L., Pinto, N., Arditia, D., Solomon, E. A., et al. (2014). Deep neural networks rival the representation of primate it cortex for core visual object recognition. *PLoS Computational Biology*, 10(12), e1003963.
- Cao, Y., Chen, Y., & Khosla, D. (2015). Spiking deep convolutional neural networks for energy-efficient object recognition. *International Journal of Computer Vision*, 113(1), 54–66.
- Cichy, R. M., Khosla, A., Pantazis, D., Torralba, A., & Oliva, A. (2016). Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical correspondence. *Scientific Reports*, 6.

- Delorme, A., Perrinet, L., & Thorpe, S. J. (2001). Networks of integrate-and-fire neurons using rank order coding b: Spike timing dependent plasticity and emergence of orientation selectivity. *Neurocomputing*, 38, 539–545.
- DiCarlo, J. J., & Cox, D. D. (2007). Untangling invariant object recognition. *Trends in Cognitive Sciences*, 11(8), 333–341.
- DiCarlo, J. J., Zoccolan, D., & Rust, N. C. (2012). How does the brain solve visual object recognition? *Neuron*, 73(3), 415–434.
- Diehl, P. U., & Cook, M. (2015). Unsupervised learning of digit recognition using spike-timing-dependent plasticity. *Frontiers in Computational Neuroscience*, 9, 99.
- Diehl, P. U., Neil, D., Binas, J., Cook, M., Liu, S.-C., & Pfeiffer, M. (2015). Fast-classifying, high-accuracy spiking deep networks through weight and threshold balancing. In *IEEE International joint conference on neural networks, IJCNN*, Killarney, Ireland (pp. 1–8).
- Diehl, P. U., Zarrella, G., Cassidy, A., Pedroni, B. U., & Neftci, E. Conversion of artificial recurrent neural networks to spiking neural networks for low-power neuromorphic hardware. In *IEEE International conference on rebooting computing*, San Diego, California, USA (pp. 1–8).
- Doya, K. (2000). Complementary roles of basal ganglia and cerebellum in learning and motor control. *Current Opinion in Neurobiology*, 10(6), 732–739.
- Fukushima, K. (1980). Neocognitron : a self organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36(4), 193–202.
- Ghodrati, M., Farzamhadi, A., Rajaei, K., Ebrahimpour, R., & Khaligh-Razavi, S.-M. (2014). Feedforward object-vision models only tolerate small image variations compared to human. *Frontiers in Computational Neuroscience*, 8(74), 1–17.
- Habenschuss, S., Bill, J., & Nessler, B. (2012). Homeostatic plasticity in Bayesian spiking networks as expectation maximization with posterior constraints. In *Advances in neural information processing systems*, Lake Tahoe, Nevada, USA (pp. 773–781).
- Huang, S., Rozas, C., Treviño, M., Contreras, J., Yang, S., Song, L., et al. (2014). Associative Hebbian synaptic plasticity in primate visual cortex. *The Journal of Neuroscience*, 34(22), 7575–7579.
- Hung, C. P., Kreiman, G., Poggio, T., & DiCarlo, J. J. (2005). Fast readout of object identity from macaque inferior temporal cortex. *Science*, 310(5749), 863–866.
- Hunsberger, E., & Eliasmith, C. (2015). Spiking deep networks with LIF neurons, *arXiv:1510.08829*.
- Hussain, S., Liu, S.-C., & Basu, A. (2014). Improved margin multi-class classification using dendritic neurons with morphological learning. In *IEEE International symposium on circuits and systems, ISCAS*, Melbourne, VIC, Australia (pp. 2640–2643).
- Khaligh-Razavi, S.-M., & Kriegeskorte, N. (2014). Deep supervised, but not unsupervised, models may explain cortical representation. *PLoS Computational Biology*, 10(11), e1003915.
- Kheradpisheh, S. R., Ganjtabesh, M., & Masquelier, T. (2016). Bio-inspired unsupervised learning of visual features leads to robust invariant object recognition. *Neurocomputing*, 205, 382–392.
- Kheradpisheh, S. R., Ghodrati, M., Ganjtabesh, M., & Masquelier, T. (2016a). Deep networks resemble human feed-forward vision in invariant object recognition. *Scientific Reports*, 6, 32672.
- Kheradpisheh, S. R., Ghodrati, M., Ganjtabesh, M., & Masquelier, T. (2016b). Humans and deep networks largely agree on which kinds of variation make object recognition harder. *Frontiers in Computational Neuroscience*, 10, 92.
- Kirchner, H., & Thorpe, S. J. (2006). Ultra-rapid object detection with saccadic eye movements: Visual processing speed revisited. *Vision Reseach*, 46(11), 1762–1776.
- Krizhevsky, A., Sutskever, I., & Hinton, G. (2012). ImageNet classification with deep convolutional neural networks. In *Neural information processing systems, NIPS*, Lake Tahoe, Nevada, USA (pp. 1–9).
- LeCun, Y., & Bengio, Y. (1998). Convolutional networks for images, speech, and time series. In *The handbook of brain theory and neural networks* (pp. 255–258). Cambridge, MA: MIT Press.
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444.
- LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278–2324.
- Lee, H., Grosse, R., Ranganath, R., & Ng, A. Y. (2009). Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. (ICML), (pp. 1–8). New York, New York, USA: ACM Press.
- Lichtsteiner, P., Posch, C., & Delbrück, T. (2007). An 128x128 120dB 15us-latency temporal contrast vision sensor. *IEEE Journal of Solid State Circuits*, 43(2), 566–576.
- Liu, H., Agam, Y., Madsen, J. R., & Kreiman, G. (2009). Timing, timing, timing: fast decoding of object information from intracranial field potentials in human visual cortex. *Neuron*, 62(2), 281–290.
- Maass, W. (2002). Computing with spikes. *Special Issue on Foundations of Information Processing of TELEMATIK*, 8(1), 32–36.
- Martínez-Cañada, P., Morillas, C., Pino, B., Ros, E., & Pelayo, F. (2016). A computational framework for realistic Retina modeling. *International Journal of Neural Systems*, 26(7), 1650030.
- Masquelier, T., & Thorpe, S. J. (2007). Unsupervised learning of visual features through spike timing dependent plasticity. *PLoS Computational Biology*, 3(2), e31.
- McMahon, D. B., & Leopold, D. A. (2012). Stimulus timing-dependent plasticity in high-level vision. *Current Biology*, 22(4), 332–337.
- Meliza, C. D., & Dan, Y. (2006). Receptive-field modification in rat visual cortex induced by paired visual stimulation and single-cell spiking. *Neuron*, 49(2), 183–189.
- O'Connor, P., Neil, D., Liu, S.-C., Delbrück, T., & Pfeiffer, M. (2013). Real-time classification and sensor fusion with a spiking deep belief network. *Frontiers in Neuroscience*, 7, 178.
- Panda, P., & Roy, K. (2016). Unsupervised regenerative learning of hierarchical features in spiking deep networks for object recognition. In *IEEE International joint conference on neural networks, IJCNN*, Vancouver, Canada (pp. 1–8).
- Pignatelli, M., & Bonci, A. (2015). Role of dopamine neurons in reward and aversion: a synaptic plasticity perspective. *Neuron*, 86(5), 1145–1157.
- Pinto, N., Barhom, Y., Cox, D. D., & DiCarlo, J. J. (2011). Comparing state-of-the-art visual features on invariant object recognition tasks. In *IEEE workshop on applications of computer vision, WACV*, Kona, Hawaii, USA (pp. 463–470).
- Portelli, G., Barrett, J. M., Hilgen, G., Masquelier, T., Maccione, A., Di Marco, S., et al. (2016). Rank order coding: a retinal information decoding strategy revealed by large-scale multielectrode array retinal recordings. *Eneuro*, 3(3), ENEURO-0134.
- Querlioz, D., Bichler, O., Dollfus, P., & Gamrat, C. (2013). Immunity to device variations in a spiking neural network with memristive nanodevices. *IEEE Transactions on Nanotechnology*, 12(3), 288–295.
- Rolls, E. T., & Deco, G. (2002). *Computational neuroscience of vision*. Oxford, UK: Oxford university press.
- Rousselet, G. A., Thorpe, S. J., & Fabre-Thorpe, M. (2003). Taking the max from neuronal responses. *Trends in Cognitive Sciences*, 7(3), 99–102.
- Serrano-Gotarredona, T., Masquelier, T., Prodromakis, T., Indiveri, G., & Linares-Barranco, B. (2013). STDP and STDP variations with memristors for spiking neuromorphic learning systems. *Frontiers in Neuroscience*, 7(FEBRUARY), 2.
- Serre, T., Oliva, A., & Poggio, T. (2007). A feedforward architecture accounts for rapid categorization. *Proceedings of the National Academy of Sciences*, 104(15), 6424–6429.
- Serre, T., Wolf, L., Bileschi, S., Riesenhuber, M., & Poggio, T. (2007). Robust object recognition with cortex-like mechanisms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(3), 411–426.
- Shoham, S., O'Connor, D. H., & Segev, R. (2006). How silent is the brain: is there a dark matter problem in neuroscience? *Journal of Comparative Physiology A*, 192(8), 777–784.
- Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition, *arXiv:1409.1556*.
- Sivilotti, M. (1991). *Wiring considerations in analog VLSI systems with application to field-programmable networks*. Pasadena, CA: Comput. Sci. Div., California Inst. Technol., (Ph.D. thesis).
- Thorpe, S., Delorme, A., & Van Rullen, R. (2001). Spike-based strategies for rapid processing. *Neural Networks*, 14(6), 715–725.
- Thorpe, S., Fize, D., Marlot, C., et al. (1996). Speed of processing in the human visual system. *Nature*, 381(6582), 520–522.
- Van Rullen, R., & Thorpe, S. J. (2001). Rate coding versus temporal order coding: what the retinal ganglion cells tell the visual cortex. *Neural Computation*, 13(6), 1255–1283.
- Wohrer, A., & Kornprobst, P. (2009). Virtual Retina: a biological retina model and simulator, with contrast gain control. *Journal of Computational Neuroscience*, 26(2), 219–249.
- Yousefzadeh, A., Serrano-Gotarredona, T., & Linares-Barranco, B. (2015). Fast Pipeline 128x128 pixel spiking convolution core for event-driven vision processing in FPGAs. In *Proceedings of 1st International conference on event-based control, communication and signal processing, EBCCSP 2015*. <http://dx.doi.org/10.1109/EBCCSP.2015.7300698>.
- Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *European conference on computer vision, ECCV*, Zurich, Switzerland (pp. 818–833).
- Zhao, B., Ding, R., Chen, S., Linares-Barranco, B., & Tang, H. (2015). Feedforward categorization on aer motion events using cortex-like features in a spiking neural network. *IEEE Transactions on Neural Networks and Learning Systems*, 26(9), 1963–1978.