

PCA-RECT: An Energy-efficient Object Detection Approach for Event Cameras

Bharath Ramesh*^[0000-0001-8230-3803], Andrés Ussa, Luca Della Vedova, Hong Yang, and Garrick Orchard

Temasek Laboratories, National University of Singapore, Singapore 117411
bharath.ramesh03@u.nus.edu
<http://sites.google.com/view/bharath-ramesh/>

Abstract. We present the first purely event-based, energy-efficient approach for object detection and categorization using an event camera. Compared to traditional frame-based cameras, choosing event cameras results in high temporal resolution (order of microseconds), low power consumption (few hundred mW) and wide dynamic range (120 dB) as attractive properties. However, event-based object recognition systems are far behind their frame-based counterparts in terms of accuracy. To this end, this paper presents an event-based feature extraction method devised by accumulating local activity across the image frame and then applying principal component analysis (PCA) to the normalized neighborhood region. Subsequently, we propose a backtracking-free k -d tree mechanism for efficient feature matching by taking advantage of the low-dimensionality of the feature representation. Additionally, the proposed k -d tree mechanism allows for feature selection to obtain a lower-dimensional dictionary representation when hardware resources are limited to implement dimensionality reduction. Consequently, the proposed system can be realized on a field-programmable gate array (FPGA) device leading to high performance over resource ratio. The proposed system is tested on real-world event-based datasets for object categorization, showing superior classification performance and relevance to state-of-the-art algorithms. Additionally, we verified the object detection method and real-time FPGA performance in lab settings under non-controlled illumination conditions with limited training data and ground truth annotations.

Keywords: Object recognition · Neuromorphic vision · Silicon retinas · Low-power FPGA · Object detection · Event cameras.

1 Introduction

Through these fruitful decades of computer vision research, we have taken huge strides in solving specific object recognition tasks, such as classification systems for automated assembly line inspection, hand-written character recognition in mail sorting machines, bill inspection in automated teller machines, to name a few. Despite these successful applications, generalizing object appearance, even

under moderately controlled sensing environments, for robust and practical solutions for industrial challenges like robot navigation and sense-making is a major challenge. This paper focuses on the industrially relevant problem of real-time, low-power object detection using an asynchronous event-based camera [2] with limited training data under unconstrained lighting conditions. Compared to traditional frame-based cameras, event cameras do not have a global shutter or a clock that determines its output. Instead, each pixel responds independently to temporal changes with a latency ranging from a low of tens of microseconds to a high of few milliseconds. This local sensing paradigm naturally results in a wider dynamic range (120 dB), as opposed to the usual 60 dB for frame-based cameras.

Most significantly, event cameras do not output pixel intensities, but only a spike output with a precise timestamp, also termed an event, that signifies a sufficient change in log-intensity of the pixel. As a result, event cameras require lower transmission bandwidth and consume only a few hundred mW vs. a few W by standard cameras [21]. In summary, event-based cameras offer a fundamentally different perspective to visual imaging while having a strong emphasis on low-latency and low-power algorithms [7,3,16,4].

Despite the notable advantages of event cameras, there still remains a significant performance gap between event camera algorithms and frame-based counterparts for various vision problems. This is partly due to a requirement of totally new event-by-event processing paradigms. However, the burgeoning interest in event-based classification/detection is focused on closing the gap using deep spiking neural networks [17,9], something that again entails dependence on powerful hardware like its frame-based counterpart. On the other hand, a succession of frames captured at a constant rate (say 30 Hz), regardless of the scene dynamics and ego-motion, works well with controlled scene condition and camera motion. Frame-based computer vision algorithms have benefited immensely from sophisticated methodologies that reduce the computational burden by selecting and processing only informative regions/keypoints within an image [12,5,23,29]. In addition, frame-based sensing has led to high hardware complexity, such as powerful GPU requirements for state-of-the-art object detection frameworks using deep neural networks [25,26].

In contrast to the above works, this paper introduces a simple, energy-efficient approach for object detection and categorization. Fig. 1 illustrates the local event-based feature extraction pipeline that is used for classification using a dictionary-based method. Accordingly, efficient feature matching with the dictionary is required, which is handled by a backtracking-free branch-and-bound k -d tree. This proposed system was ported to a field programmable gate array (FPGA) with certain critical design decisions, one of which demanded a virtual dimensionality reduction method based on the k -d tree, to accommodate very low-power computational needs.

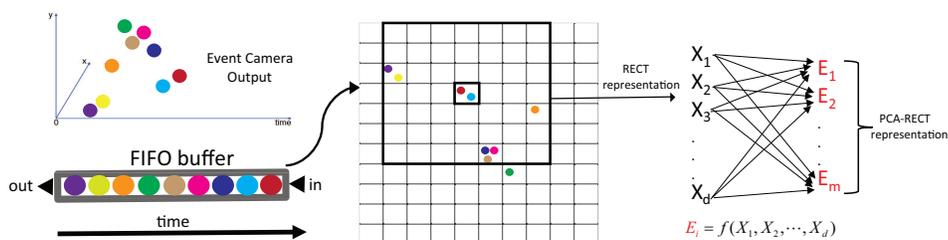


Fig. 1: PCA-RECT representation (best viewed on monitor). Useful events are sub-sampled and filtered after applying nearest-neighbor temporal filtering and refractory filtering, termed as rectangular event context transform (RECT). The sparser RECT event representation is updated dynamically using a first in, first out (FIFO) buffer. Subsequent feature extraction is carried out by applying principal component analysis (PCA) to project RECT onto a lower-dimensional subspace to obtain the final PCA-RECT feature representation

2 Event Cameras

For real-time experiments, we use the commercial event camera, the Dynamic and Active-pixel Vision Sensor (DAVIS) [2]. It has 240×180 resolution, 130 dB dynamic range and 3 microsecond latency. The DAVIS can concurrently output a stream of events and frame-based intensity read-outs using the same pixel array. An event consists of a pixel location (x, y) , a binary polarity value (p) for positive or negative change in log intensity and a timestamp in microseconds (t) . In this work, polarity of the events are not considered, and only the event stream of the DAVIS is used.

2.1 Related Work

Since event-based vision is relatively new, only a limited amount of work addresses object detection using these devices [11,10]. Liu *et al.* [11] focuses on combining a frame-based CNN detector to facilitate the event-based module. We argue that works using deep neural networks for event-based object detection may achieve good performance with lots of training data and computing power, but they go against the idea of low-latency, low-power event-based vision. In contrast, [10] presents a practical event-based approach to face detection by looking for pairs of blinking eyes. While [10] is applicable to human faces in the presence of activity, we develop a general purpose event-based, object detection method using a simple feature representation based on local event aggregation. Thus, this paper is similar in spirit to the recently spawned ideas of generating event-based descriptors, such as histogram of averaged time surfaces [28] and log-polar grids [24,22]. Moreover, the proposed object detection and categorization method was accommodated on FPGA to demonstrate energy-efficient low-power vision.

3 Method

We follow the event-based classification framework proposed in [22], with the following crucial changes: a new descriptor (PCA-RECT), a virtual dimensionality reduction technique using k -d trees (vPCA) and a simplified feature matching mechanism to account for hardware limitations. The framework [22] consists of four main stages: feature extraction, feature matching with a dictionary, dictionary representation followed by a linear classifier. Additionally, we incorporate an object detector in the framework as explained in the following subsections.

3.1 PCA-RECT

Each incoming event, $\mathbf{e}_i = (x_i, y_i, t_i, p_i)^T$ with pixel location x_i and y_i , timestamp t_i , polarity p_i , is encoded as a feature vector \mathbf{x}_i . To deal with hardware-level noise from the event camera, two main steps are used: (1) nearest neighbour filtering and (2) refractory filtering. We define a spatial Euclidean distance between events as,

$$D_{i,j} = \left\| \begin{pmatrix} x_i \\ y_i \end{pmatrix} - \begin{pmatrix} x_j \\ y_j \end{pmatrix} \right\|. \quad (1)$$

Using the above distance measure, for any event we can define a set of previous events within a spatial neighborhood, $N(\mathbf{e}_i, \gamma) = \{\mathbf{e}_j \mid j < i, D_{i,j} < \gamma\}$, where $\gamma = \sqrt{2}$ for an eight-connected pixel neighbourhood. When the time difference between the current event and the most recent neighboring event is less than a threshold, Θ_{noise} , the filter can be written as

$$F_{noise}(\mathbf{e}) = \{\mathbf{e}_i \mid N(\mathbf{e}_i, \sqrt{2}) \setminus N(\mathbf{e}_i, 0) \ni \mathbf{e}_j \mid t_i - t_j < \Theta_{noise}\}. \quad (2)$$

When the neighborhood is only the current pixel, $\gamma = 0$, the set of events getting through the refractory filter F_{ref} are those such that,

$$F_{ref}(\mathbf{e}) = \{\mathbf{e}_j \mid t_i - t_j > \Theta_{ref} \forall j \mid \mathbf{e}_j \in N(\mathbf{e}_j, 0)\}. \quad (3)$$

Cascading the filters, we can write the filtered incoming events as,

$$\{\hat{\mathbf{e}}\} = F_{noise}(F_{ref}(\mathbf{e})). \quad (4)$$

As shown in Fig. 1, the incoming events $\hat{\mathbf{e}}_i$ are first pushed into a FIFO buffer. The FIFO queue is then used to update an event-count matrix $C \in \mathbb{R}^{m \times n}$, where m and n denote the number of rows and columns of the event camera output.

$$C(x_i, y_i) = C(x_i, y_i) + 1. \quad (5)$$

Before pushing the latest event, the FIFO buffer of size s is popped to make space and simultaneously update the count matrix C ,

$$C(x_{i-s}, y_{i-s}) = C(x_{i-s}, y_{i-s}) - 1. \quad (6)$$

The event-count C is pooled to build local representations, which are further aggregated to obtain the RECT representation of each event. In particular, let A be a $p \times p$ square filter, the 2-D convolution is defined as,

$$R(j, k) = \sum_p \sum_q A(p, q)C(j - p + 1, k - q + 1) , \quad (7)$$

where p run over all values that lead to legal subscripts of $A(p, p)$ and $C(j - p + 1, k - p + 1)$. In this work, we consider a filter containing equal weights (commonly known as an averaging filter) for simplicity, while it is worth exploring Gaussian-type filters that can suppress noisy events. The resultant 2-D representation is termed as filtered matrix $R \in \mathbb{R}^{(m/p) \times (n/p)}$, where the filter dimensions are chosen to be give integer values for m/p and n/p or conversely C is zero-padded sufficiently. Subsequently, the RECT representation for \hat{e}_i is obtained as a patch \mathbf{u}_i of dimension d centered at $R(y/p, x/p)$. Subsequently, the filtered event-count patch is projected on-to a lower-dimensional subspace using principal component analysis (PCA) for eliminating noisy dimensions and improving classifier accuracy.

3.2 Feature Selection and Matching using K -d Trees

The PCA-RECT feature representation for each event is classified using a dictionary type method [22] that can handle the recognition of the desired object categories. However, exhaustive search is too costly for nearest neighbor matching with a dictionary, and approximate algorithms can be orders of magnitude faster than exact search, while almost achieving par accuracy.

In the vision community, k -d tree nearest-neighbor search is popular [27,14], as a means of searching for feature vectors in a large training database. Given n feature vectors $\mathbf{x}_i \in \mathbb{R}^{d'}$, the k -d tree construction algorithm recursively partitions the d' -dimensional Euclidean space into hyper-rectangles along the dimension of maximum variance. However, for high dimensional data, backtracking through the tree to find the optimal solution takes a lot of time.

This paper proposes a simple, backtracking-free branch-and-bound search for dictionary matching, taking advantage of the low-dimensionality of the PCA-RECT representation. The hypothesis is that, in general, the point recovered from the leaf node is a good approximation to the nearest neighbor in low-dimensional spaces, and performance degrades rapidly with increase in dimensionality, as inferred from the intermediate results in [1]. In other words, with $(\log_2 n) - 1$ scalar comparisons, nearest neighbor matching is accomplished without an explicit distance calculation. While the PCA-RECT representation is useful for software implementations, an extra PCA projection step can be computationally demanding on FPGA devices. To this end, we propose a virtual PCA-RECT representation based on the k -d tree, termed as vPCA-RECT.

vPCA-RECT A key insight is that only a fraction of the data dimensions are used to partition the k -d tree, especially when the dictionary size is only a few

times more than the feature dimension. Therefore, instead of using the PCA-RECT representation, an alternative dimensionality reduction scheme can be implemented by discarding the unused dimensions in the k -d tree structure. In other words, the RECT representation is first used to build a k -d tree that selects the important dimensions (projection π), which are then utilized for dictionary learning and classification. It is worth noting that exactly the same k -d tree will be obtained if the RECT data is first projected by π onto a subspace that is aligned with the coordinate axes. Since no actual projection takes place, we refer to this as a virtual projection – the irrelevant dimensions chosen by the k -d tree are discarded to obtain a lower-dimensional feature representation.

3.3 Event-based Object Categorization and Detection

The learning stage: Using either the PCA-RECT or vPCA-RECT event representation, the learning process corresponds to creating a set of K features denoted as $M = \{1, 2, \dots, K\}$ to form the dictionary. First, a simple sampling process is carried out such that, during training, a large pool of event representations of various categories and at random positions are extracted from a target set of events. In our setup, the dictionary features are learned from the sampled training set using clustering for all the objects jointly.

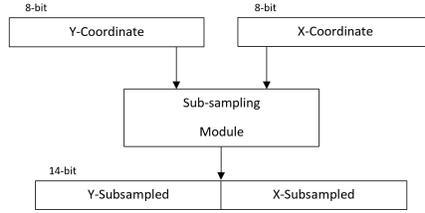
The learning stage for detection builds on top of the categorization module, in such a way that the learning process corresponds to selecting a subset of features from the dictionary for each object. In contrast to the learning phase of the categorization module, the detector features are selected from the whole training set in a supervised one-vs-all manner.

We propose to evaluate the balanced matches Y_+^k to each dictionary feature f_k from the target events against the matches Y_-^k for all the other events to the respective feature. Mathematically, the ratio

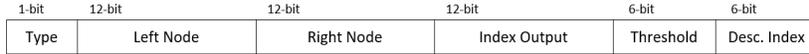
$$D(k) = \frac{\beta_+^k Y_+^k}{\beta_-^k Y_-^k}, \text{ where } \beta_+^k = \frac{|Y_+^k|}{\sum_{k=1}^K |Y_+^k|}, \text{ and } \beta_-^k = \frac{|Y_-^k|}{\sum_{k=1}^K |Y_-^k|}, \quad (8)$$

is to be maximized. The balancing component β_+^k denotes the percentage of target events matched to the dictionary feature f_k . Similarly, β_-^k denotes the percentage of non-target events matched to the dictionary feature f_k . Thus, choosing the detector features with the D-largest ratios completes the learning phase.

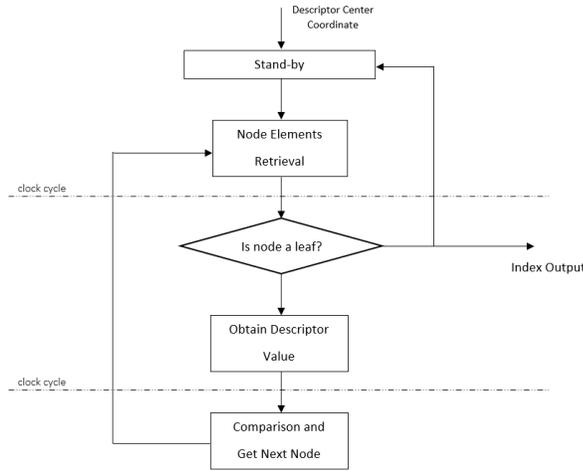
The classification/detection stage: At runtime, the event representations are propagated through the k -d tree. On the one hand, the distribution of the dictionary features are then extracted and further passed to a simple linear classifier (we experimented with both linear SVM and Kernel Methods). On the other hand, the event representations propagated through the k -d tree are matched with the detector features. Those matched events are used to update a location map for the target object and the region with the highest activation is considered to be the final detection result.



(a) Sub-sampling module



(b) A k -d tree node in hardware



(c) Recursive logic-driven k -d tree implemented in hardware

Fig. 2: FPGA implementation details

4 FPGA Implementation

4.1 Categorization Pipeline

In order to showcase energy-efficient event-based object recognition, the FPGA implementation of the algorithm is designed as a series of four independent hardware units: event sub-sampling, vPCA-RECT generation, a recursive k -d tree and a SVM classifier output on an event-by-event basis, each of which has an independent block design. Generally, these hardware counterparts are not a direct application of the algorithm presented in the earlier section, i.e., certain design decisions were taken for this task, among them, to desist the use of an extra PCA projection along the pipeline.

The sub-sampling block receives the filtered event locations as input values x and y , each 8-bit in size, which are used to update the zero-padded count matrix $C \in \mathbb{R}^{m \times n}$ (Eq. (5) and Eq. (6)). The sub-sampling behavior can be achieved in hardware through a combinatorial module that performs the division by shifting the inputs by one bit, and subsequently adding p and q to that value to obtain the sub-sampled representation (Eq. (7)). This results in two 7-bit values which are then concatenated to output a single memory address (Fig. 2(a)).

The next block uses the cell-count matrix $R \in \mathbb{R}^{(m/p) \times (n/q)}$, created by a block of distributed RAM of depth $((m/p) \times (n/q))$ and $\log(s)$ -bits width, corresponding to the FIFO buffer size s , initialized to zero for generating the vPCA-RECT representation. To generate a descriptor with respect to the last event received would add a considerable overhead, since each element of the descriptor would have to be read sequentially from the block RAM while being stored by the next module. Instead, the address corresponding to the center of the descriptor is provided, i.e. the input address of the count matrix is passed over to the k -d tree module. This allows to trigger the k -d tree in one clock cycle once the count matrix is updated and later read the descriptor values based on this single coordinate. However, a new issue arises, the count matrix then can not be modified while the k -d tree exploration is being performed. Hence a buffering element is added between the sub-sampling and count matrix modules that will only provide the next address once there is a valid output from the tree.

The k -d tree nodes are represented in a 49-bit number stored in a previously initialized single port ROM of depth equal to the number of nodes. This number is conformed by the elements of a node: type, left node, right node, index output, split value and split dimension; these are concatenated and their width is shown in Fig. 2(b).

The k -d tree module follows a three steps cycle (Fig. 2(c)). The split dimension of a k -d tree node provides the address that needs to be read from the cell-count matrix block RAM to get the relevant descriptor value. Next, the descriptor value is compared to the previously stored split value from the node, taking a path down the tree, left or right, depending on the boolean condition. The corresponding node to get is then retrieved from the respective left or right address element acquired in the retrieval step. This cycle repeats until the node type belongs to a leaf, then the leaf node output is made available for the classifier module. It is worth mentioning that in the software implementation of this algorithm, once the descriptor is formed, it is then normalized before being passed to the k -d tree. A normalization step in hardware would add a big overhead to the pipeline, disturbing its throughput, and it was removed from the FPGA implementation after verifying that the overall performance was not affected harshly.

At runtime in a software implementation, the classification is performed by a linear combination of the weights and a feature vector created by the k -d tree after a buffer time of S events. To achieve this in a hardware implementation, the depth of the feature vector would have to be transversed while performing several multiplications which would require a considerable amount of multiplier

Algorithm 1 Event-based FPGA Object Detection

Input: Filtered event stream $\{\hat{\mathbf{e}}\}$, detector landmarks l , number of events S **Output:** Mean object location (x_{obj}, y_{obj})

```

1: Initialize detector count  $D(y, x) = 0_{m,n}$ , detector cut-off  $threshold = 0$ 
2: for  $t = 1 : S$  do
3:   For each incoming event  $\hat{\mathbf{e}}_t = (x_t, y_t, t, p_t, \mathbf{x}_t^T)^T$ 
4:   For  $\mathbf{x}_t$  get leaf node index  $l_t$  using  $k$ -d tree
5:   if  $l_t \in l$  then
6:      $D(y_t, x_t) = D(y_t, x_t) + 1$ 
7:     if  $D(y_t, x_t) > threshold$  then
8:        $threshold = threshold + 1$ 
9:       Reset detector mean calculation FIFO
10:    end if
11:    if  $D(y_t, x_t) = threshold$  then
12:      Push  $x_t, y_t$  into the mean calculation FIFO
13:    end if
14:  end if
15: end for
16: Output the mean of the coordinates in the FIFO as  $(x_{obj}, y_{obj})$ 

```

elements from the FPGA, and would affect the speed of the module. Thus, it was desired to avoid this solution and the following was proposed.

The elements of the linear combination mentioned would be acquired as readily available and would be added to an overall sum vector of length equal to the number of classes to classify, hence performing the dot product operation as one addition per event. Then, after S events, a resulting vector is formed, which is equal to the result of the same linear combination first mentioned in the software implementation. Thus, the final module to perform the classification receives the output index from the k -d tree and adds its corresponding classifier parameter to a sum vector of length equal to the number of classes. In parallel, this index value is stored in another FIFO element. When the queue is full, the oldest value would be passed to the module to be subtracted from the sum. This allows to have a classification output at any point in time, corresponding to the last S events.

4.2 Detection Pipeline

Parallel to the modules performing the classification pipeline, the aim of the detection process is to find the coordinates corresponding to “landmarks” with the highest activation after S events, and then find the most probable location for the object. Again, the algorithm was divided into multiple coherent hardware modules that would produce the same results as the original software version. The designed blocks are: landmarks detector, detection heat map and mean calculation.

First, the dictionary features corresponding to the landmarks that were calculated offline are loaded into a binary memory block. This module receives as

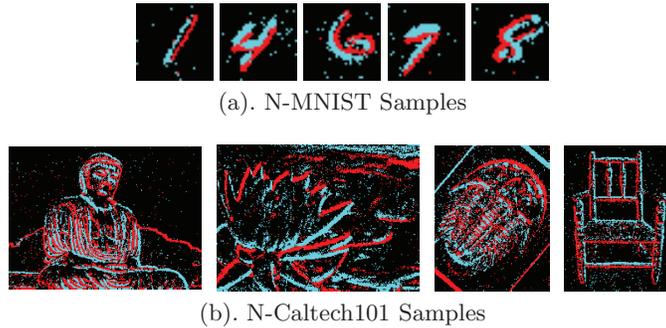


Fig. 3: Samples from the Event-based Benchmark Datasets

input the dictionary feature index provided by the k -d tree for the current event. If the feature is found as one of the landmarks, the respective event coordinates x and y are passed as a concatenated address to the next module in the pipeline. Next, a stage corresponding to the heat map is utilized. This module holds a matrix represented as a block RAM of depth $m \times n$, since the coordinates are not sub-sampled and have the ranges $1 \leq x \leq m$ and $1 \leq y \leq n$. For each new input address, its value in memory is incremented.

Since the aim of the detection algorithm is to calculate the average of the coordinates with the highest activation, it would be inefficient to find these event addresses after S events. Therefore, the coordinates with the highest count are stored in a FIFO element while the counting is performed. At the end, this will contain all the x and y coordinates needed for the average calculation. Once the classification flag is triggered, all the coordinates stored in the previous step (which belong to the highest activation) are acquired for calculating the total activation (the divisor). Subsequently, it will calculate the sum of the respective x and y values, and pass these as dividends to hardware dividers that will provide the final coordinates of the detected object. Alg. 1 summarizes the above object detection hardware pipeline clearly.

5 Experiments and Discussion

5.1 Event-based Object Categorization

This section compares the proposed object categorization system to state-of-the-art event-based works and thus software implementation is used with double numeric precision. We validated our approach on two benchmark datasets [19], namely the N-MNIST and N-Caltech101, that have become de-facto standards for testing event-based categorization algorithms. Fig. 3 shows some representative samples from N-MNIST and N-Caltech101.

Parameter Settings. The time thresholds for the nearest neighbour filter and the refractory filter are nominally set as $\Theta_{noise} = 5$ ms and $\Theta_{ref} = 1$ ms re-

Table 1: Comparison of classification accuracy on event-based datasets (%).

	N-MNIST	N-Caltech101
H-First	71.20	5.40
HOTS	80.80	21.0
Gabor-SNN	83.70	19.60
HATS	99.10	64.20
vPCA-RECT (this work)	98.72	70.25
PCA-RECT (this work)	98.95	72.30
Phased LSTM	97.30	-
Deep SNN	98.70	-

spectively, as suggested in [20]. We used a FIFO buffer size of 5000 events for dynamically updating the count matrix as and when events are received. Subsequently, the RECT representation with a 2 by 2 averaging filter without zero padding at the boundaries is used to obtain a 9×9 feature vector for all event locations. We also experimented with other feature vector dimensions using a 3×3 , 5×5 , 7×7 sampling region and found that increasing the context improved the performance slightly. For obtaining the PCA-RECT representation, the number of PCs can be chosen automatically by retaining the PCs that hold 95% eigenenergy of the training data, which is typically about 60 in our case. For testing on the benchmark datasets, a dictionary size of 3000 was universally used with a k -d tree with backtracking to find precise feature matches.

Results on the Benchmark Datasets. The results on the N-MNIST and N-Caltech101 datasets are given in Table 1. As it is common practice, we report the results in terms of classification accuracy. The baselines methods considered were HATS [28], HOTS [8], HFirst [18] and Spiking Neural Networks (SNN) [9,15] (Gabor-SNN reported in [28]).

On the widely reported N-MNIST dataset, our method is as good as the best performing HATS method. Moreover, other SNN methods are also in the same ballpark, which is due to the simple texture-less digit event streams giving distinct features for most methods. Therefore, it is a good benchmark as long as a proposed method performs in the high 90's. A test on the challenging NCaltech-101 dataset will pave way for testing the effectiveness close to a real-world scenario.

Our method has the highest classification rate ever reported for an event-based classification method on the challenging N-Caltech101 dataset. The unpublished HATS work is the only comparable method in terms of accuracy, while the other learning mechanisms fail to reach good performance. Fig. 4 shows the performance of the PCA-RECT representation as the number of PCs are varied. It is worth noticing that just retaining five dimensions can give better performance compared to available works.

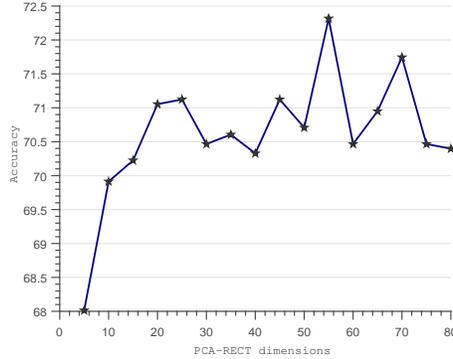


Fig. 4: Number of principal component dimensions vs. accuracy on N-Caltech101



Fig. 5: Dataset samples: Landing platform, UAV and Thumper (BG: Empty Floor).

5.2 Event-based Object Detection

Dataset. The datasets described in the previous section are good for only evaluating the categorization module. In addition, as the benchmark datasets were generated by displaying images on a monitor with limited and predefined motion of the camera, they do not generalize well to real-world situations. To overcome these limitations, we created a new dataset by directly recording objects in lab environment with a freely moving event-based sensor. The in-house dataset, called as Neuromorphic Single Object Dataset (N-SOD), contains three objects with samples of varying length in time (up to 20 s). The three objects to be recognized are a thumper 6-wheel ground robot, an unmanned aerial vehicle, a landing platform along with a background class (Fig. 5).

Results on N-SOD. For testing on the N-SOD dataset, we divide the dataset into training and testing, with 80% temporal sequence samples per class for training and the remaining for testing. Using the training features, a dictionary is generated. Since the temporal sequences are of different length, for a fixed number of events, say every 10^5 events, a dictionary representation is extracted and a linear SVM classifier is trained. Similarly for testing, for every 10^5 events, the dictionary representation is classified using the SVM.

Based on the above setup, an accuracy of 97.14% was obtained (Tab. 2) with a dictionary size of 950, which resulted in a k -d tree with 10 layers. We also experimented with lower dictionary sizes such as 150, 300, 450, etc., and the

Table 2: Confusion Matrix (%) for the best result on the in-house N-SOD dataset.

	Background	LP	Thumper	UAV
Background	95.4128	0.3058	3.3639	0.9174
LandingPlatform	0	99.2268	0.5155	0.2577
Thumper	0	1.9257	96.9739	1.1004
UAV	0	0	3.1884	96.8116

performance drop was insignificant ($> 96\%$). On the other hand, using a k -d tree with backtracking, descriptor normalization, etc., achieved close to 100% accuracy on offline high-performance PCs, which of course does not meet low-power and real-time requirements. In summary, the proposed vPCA-RECT method with a backtracking-free k -d tree implementation mildly compromises on accuracy to handle object detection and categorization using an event camera in real-time.

We report the precision and recall of the detection results by ascertaining if the mean position of the detected result is within the ground truth bounding box. We obtained: (a) *Precision* - $(498/727) = 0.685$: The percentage of the detections belonging to the object that overlap with the groundtruth (b) *Recall* - $(498/729) = 0.683$: The percentage of correct detections that are retrieved by the system. The number of “landmarks” were set to 20 in the above experiments while similar results were obtained for values such as five and ten.

Comparison to CNN. In order to compare to state-of-the-art deep neural networks, we recorded a similar dataset to N-SOD using a frame-based camera and transfer learning via AlexNet classified the object images. With an equivalent train/test split compared to N-SOD, perfect performance can be achieved on the clearly captured test images, however, when we recorded a frame-based dataset under fast motion conditions (motion blur), an accuracy of only 79.20% was obtained. It was clear that the black UAV frame when blurred looks like the black-stripped background and creates much confusion. This confirms the disadvantage of using frame-based cameras to handle unconstrained camera motion. Note that fast camera motion leads to only an increase in data-rate for event-based cameras and has no effect on the output. In fact, recordings of N-SOD have significant amount of such fast motions.

FPGA Performance. The hardware implementation and performance of the Xilinx Zynq-7020 FPGA running at 100 MHz was evaluated by direct comparison with the results of the algorithm’s software version in MATLAB. The time taken for a single event to be classified for the worst possible k -d tree path was 550 nanoseconds. The Zynq was interfaced to a down-looking DAVIS camera, on-board an unmanned aerial vehicle flying under unconstrained lighting

Table 3: Hardware utilization report for the FPGA running the proposed modules.

	Utilization	Available	Utilization %
LUT	18238	53200	34.28
LUTRAM	12124	17400	69.68
FF	2065	106400	1.94
BRAM	48	140	34.29
DSP	4	220	1.82
IO	102	200	51.00

scenarios. We recommend viewing our submitted video¹ that clearly shows the classification/detection process better than still images.

A summary of utilization of hardware elements can be seen in Tab. 3. Due to our low hardware requirements, the dynamic on-chip power increased only by 0.37W while the base FPGA power consumption was in the order of 3W. As a rough comparison, FPGA-based recognition systems like [13,6,30], which present solutions running at equal or lower clock frequencies, consume more power than our implementation.

6 Conclusions

We have demonstrated object detection and categorization in an energy-efficient manner using event cameras, where the only information that is important for the tasks is how edges move, and the event camera naturally outputs it. The proposed PCA-RECT feature takes advantage of this sparsity to generate a low-dimensional representation. The low-dimensional representation is further exploited for feature matching using a k -d tree approach, capable of obtaining the best performance on the challenging Neuromorphic Caltech-101 dataset compared to state-of-the-art works. Most importantly, real-time FPGA implementation was achieved with several careful design considerations, such as a backtracking-free k -d tree for dictionary matching, a virtual PCA-RECT representation obtained by analyzing the k -d tree partitioning of the feature space, etc. To the best of our knowledge, this is the first work implementing a generic object recognition framework for event cameras on an FPGA device, verified in a lab demo setting under unconstrained motion and lighting setup, thereby demonstrating a high performance over resource ratio.

References

1. Beis, J.S., Lowe, D.G.: Shape indexing using approximate nearest-neighbour search in high-dimensional spaces. In: Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR). IEEE Computer Society (1997)

¹ <https://youtu.be/h3SgXa47Kjc>

2. Brandli, C., Berner, R., Yang, M., Liu, S.C., Delbruck, T.: A 240 x 180 130 db 3 μ s latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits* **49**(10), 2333–2341 (Oct 2014)
3. Conradt, J., Berner, R., Cook, M., Delbruck, T.: An embedded AER dynamic vision sensor for low-latency pole balancing. In: *IEEE International Conference on Computer Vision Workshop*. pp. 780–785 (Sept 2009)
4. Delbruck, T., Lang, M.: Robotic goalie with 3 ms reaction time at 4% CPU load using event-based dynamic vision sensor. *Frontiers in Neuroscience* **7**, 223 (2013)
5. Galleguillos, C., Rabinovich, A., Belongie, S.: Object categorization using co-occurrence, location and appearance. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1–8. *IEEE Computer Society* (2008)
6. Hikawa, H., Kaida, K.: Novel fpga implementation of hand sign recognition system with som-hebb classifier. *IEEE Transactions on Circuits and Systems for Video Technology* **25**(1), 153–166 (2015)
7. Kueng, B., Mueggler, E., Gallego, G., Scaramuzza, D.: Low-latency visual odometry using event-based feature tracks. In: *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. pp. 16–23 (Oct 2016)
8. Lagorce, X., Orchard, G., Gallupi, F., Shi, B.E., Benosman, R.: Hots: A hierarchy of event-based time-surfaces for pattern recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **1**(99), 1–1 (2016)
9. Lee, J.H., Delbruck, T., Pfeiffer, M.: Training deep spiking neural networks using backpropagation. *Frontiers in Neuroscience* **10**, 508 (2016)
10. Lenz, G., Ieng, S., Benosman, R.: Event-based dynamic face detection and tracking based on activity. *CoRR* **abs/1803.10106** (2018), <http://arxiv.org/abs/1803.10106>
11. Liu, H., Moeys, D.P., Das, G., Neil, D., Liu, S.C., Delbruck, T.: Combined frame- and event-based detection and tracking. In: *IEEE International Symposium on Circuits and Systems (ISCAS)*. pp. 2511–2514 (May 2016)
12. Lowe, D.: Distinctive image features from scale-invariant keypoints. *International journal of computer vision* **60**(2), 91–110 (2004)
13. Mousoulitis, P.G., Panayiotou, K.L., Tsardoulis, E.G., Petrou, L.P., Symeonidis, A.L.: Expanding a robot’s life: Low power object recognition via fpga-based dcnn deployment. In: *7th International Conference on Modern Circuits and Systems Technologies (MOCAS)*. pp. 1–4 (2018)
14. Muja, M., Lowe, D.G.: Fast approximate nearest neighbors with automatic algorithm configuration. In: *VISAPP International Conference on Computer Vision Theory and Applications*. pp. 331–340 (2009)
15. Neil, D., Pfeiffer, M., Liu, S.C.: Phased LSTM: Accelerating Recurrent Network Training for Long or Event-based Sequences. In: *Neural Information Processing Systems*. pp. 3889–3897. *NIPS’16, Curran Associates Inc.* (2016)
16. Ni, Z., Bolopion, A., Agnus, J., Benosman, R., Regnier, S.: Asynchronous event-based visual shape tracking for stable haptic feedback in microrobotics. *IEEE Transactions on Robotics* **28**(5), 1081–1089 (2012)
17. O’Connor, P., Neil, D., Liu, S.C., Delbruck, T., Pfeiffer, M.: Real-time classification and sensor fusion with a spiking deep belief network. *Frontiers in Neuroscience* **7** (2013)
18. Orchard, G., Meyer, C., Etienne-Cummings, R., Posch, C., Thakor, N., Benosman, R.: HFirst: A Temporal Approach to Object Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **37**(10), 2028–2040 (2015)

19. Orchard, G., Jayawant, A., Cohen, G.K., Thakor, N.: Converting static image datasets to spiking neuromorphic datasets using saccades. *Frontiers in Neuroscience* **9**, 437 (2015)
20. Padala, V., Basu, A., Orchard, G.: A noise filtering algorithm for event-based asynchronous change detection image sensors on truenorth and its implementation on truenorth. *Frontiers in Neuroscience* **12**, 118 (2018)
21. Posch, C., Serrano-Gotarredona, T., Linares-Barranco, B., Delbruck, T.: Retinomorphic event-based vision sensors: Bioinspired cameras with spiking output. *Proceedings of the IEEE* **102**(10), 1470–1484 (Oct 2014)
22. Ramesh, B., Thi, N.A.L., Orchard, G., Xiang, C.: Spike context: A neuromorphic descriptor for pattern recognition. In: *IEEE Biomedical Circuits and Systems Conference (BioCAS)*. pp. 1–4 (Oct 2017)
23. Ramesh, B., Jian, N.L.Z., Chen, L., Xiang, C., Gao, Z.: Scalable scene understanding via saliency consensus. *Soft Computing* (Nov 2017)
24. Ramesh, B., Yang, H., Orchard, G., Thi, N.A.L., Xiang, C.: DART: distribution aware retinal transform for event-based cameras. *CoRR* **abs/1710.10800** (2017)
25. Redmon, J., Farhadi, A.: Yolov3: An incremental improvement. *CoRR* **abs/1804.02767** (2018), <http://arxiv.org/abs/1804.02767>
26. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **39**(6), 1137–1149 (2017)
27. Silpa-Anan, C., Hartley, R.: Optimised kd-trees for fast image descriptor matching. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1–8 (2008)
28. Sironi, A., Brambilla, M., Bourdis, N., Lagorce, X., Benosman, R.: HATS: histograms of averaged time surfaces for robust event-based object classification. *CoRR* **abs/1803.07913** (2018), <http://arxiv.org/abs/1803.07913>
29. Vikram, T.N., Tscherepanow, M., Wrede, B.: A saliency map based on sampling an image into random rectangular regions of interest. *Pattern Recognition* **45**(9), 3114 – 3124 (2012)
30. Zhai, X., Bensaali, F., McDonald-Maier, K.: Automatic number plate recognition on fpga. In: *International Conference on Electronics, Circuits, and Systems (ICECS)*. pp. 325–328. IEEE (2013)