# Real-Time Robust Video Object Detection System Against Physical-World Adversarial Attacks

Husheng Han,Xing Hu *IEEE,* Kaidi Xu,Pucheng Dang,Ying Wang, Yongwei Zhao,Zidong Du *IEEE,* Qi Guo *IEEE,* Yanzhi Yang,Tianshi Chen

**Abstract**—DNN-based video object detection (VOD) powers autonomous driving and video surveillance industries with rising importance and promising opportunities. However, adversarial patch attack yields huge concern in live vision tasks because of its practicality, feasibility, and powerful attack effectiveness. This work proposes *Themis*, a software/hardware system to defend against adversarial patches for real-time robust video object detection. We observe that adversarial patches exhibit extremely localized superficial feature importance in a small region with non-robust predictions, and thus propose the adversarial region detection algorithm for adversarial effect elimination. *Themis* also proposes a systematic design to efficiently support the algorithm by eliminating redundant computations and memory traffics. Experimental results show that the proposed methodology can effectively recover the system from the adversarial attack with negligible hardware overhead.

**Index Terms**—Deep learning security;Real-time and embedded systems;Adversarial patch attack;Video object detection.

✦

## 1 INTRODUCTION

POWERED by deep neural network (DNN) techniques, video recognition achieves tremendous success and starts to boost existing industries, such as autonomous driving, surveillance systems, drones, and robots. For example, autonomous driving based on video recognition, whose market is predicted to leap to $77 billion (25% of the whole automotive market) by 2035 [1], has attracted the attention of giants including Tesla and Waymo [2], [3].

Despite the promising opportunities and rising importance of DNN-powered video recognition, the vulnerability of DNNs emerges as an important problem to video recognition tasks, especially in life-critical scenarios. DNN techniques have shown to be vulnerable to adversarial attacks. For example, wearing the

*The preliminary version is published in NeurIPS 2021.*

- *Husheng Han, Pucheng Dang and Yongwei Zhao are with the SKL of Processors, Institute of Computing Technology, CAS, Beijing 100190, China, the University of Chinese Academy of Sciences, Beijing, 100049, China, and also with Cambricon Technologies, Beijing, China.*
  *E-mail:{hanhusheng20z,dangpucheng20g,zhaoyongwei}@ict.ac.cn*
- *Xing Hu and Qi Guo are with the SKL of Processors, Institute of Computing Technology, CAS, Beijing, 100190, China.*
  *E-mail:{huxing,guoqi}@ict.ac.cn*
- *Kaidi Xu is with Department of Computer Science, College of Computing & Informatics, Drexel University. Philadelphia, 19104, USA.*
  *E-mail:kx46@drexel.edu*
- *Zidong Du is with the SKL of Processors, Institute of Computing Technology, CAS, Beijing, 100190, China, and also with Cambricon Technologies, Beijing, China. E-mail:duzidong@ict.ac.cn*
- *Tianshi Chen is with Cambricon Technologies, Beijing, China.*
  *E-mail:chentianshi@ict.ac.cn*
- *Ying Wang is with SKL of Computer Architecture, Institute of Computing Technology, CAS, Beijing 100190, China. Email:wangying2009@ict.ac.cn*
- *Yanzhi Yang is with Department of Electrical and Computer Engineering, Northeastern University, Boston, 02115, USA*
  *Email:yanz.wang@northeastern.edu*

*(Corresponding author: Xing Hu.)*

T-shirt with adversarial patch printing on it, which effectively fools DNN-based person detectors in physical environments even under diverse scenarios like people walking, sitting, and running [4], [5]. Such attacks are malicious in the surveillance and autonomous vehicle application scenarios, which evade the video detectors in the physical world and incur life-or-death problems. Therefore, robust video recognition that defends against such adversarial attacks and eliminates the adversarial effects is urgent and important.

For live vision scenarios, the defensive methodology should meet the following two requirements: 1) Effectively recover (more than detection only) the system from the adversarial attacks considering video recognition is usually adopted in real-time decision-making scenarios. 2) The proposed defensive methodology should introduce lightweight performance overhead to achieve the goal of real-time object detection. Effectively recovering the VOD system from adversarial attack in real-time is a highly challenging task. Existing pioneering studies for robust image classification fail to meet these two requirements: They either detect abnormal inputs only without recovery [6], [7], [8], [9], or introduce too much overhead that cannot be born in real-time VOD systems [10], [11]. MRD introduces extremely large overhead (costs about 1446s for one ImageNet-class image), which is not feasible in real-time scenario [10]. Patchguard++ relies on analyzing deep feature of every single image and cannot utilize the temporal and spatial redundancy of adjacent frames in video data [12].

To this end, focusing on the important live vision scenario that is widely adopted in autonomous driving and surveillance systems, this work proposes the real-time robust video object detection system, *Themis*, to defend the adversarial patch attacks that practically introduce damaging consequences in video recognition tasks. We propose localized important superficial feature (LISF) based defensive methodology, which not only effectively identifies the adversarial regions but is also able to leverage temporal

and spatial redundancy of video data for efficient robust object detection.

*Themis Algorithm*[1]: We draw the key observation that adversarial inputs induce VOD to be overshadowed by the localized but inductive superficial features, and the effect of adversarial patches can be eliminated facilely without aggravating the prediction accuracy of benign images when moving out LISFs. Hence, we propose LISF-based detection and recovery method based on prediction stability testing by moving out LISFs. Results show that Themis algorithm effectively locates the adversarial regions and eliminate adversarial effects.

*Themis Architecture:* Although *Themis* algorithm is effective at detecting the adversarial regions in input data, it also introduces challenges due to the additional multiple rounds of inference during the prediction stability testing by occluding LISFs. To reduce performance overhead to support real-time robust video object detection, we propose the *Themis* architecture to eliminate both the inter-frame and intra-frame redundant computations. 1) Inter-frame: *Themis* leverages the spatial and temporal redundancy of video data to eliminate unnecessary computations in non-key frames. Specifically, for key frames of video data, the complete defensive algorithm is performed to locate adversarial region or features. For non-key frames, by leveraging the temporal and spatial locality in video data, Themis approximately predicts the adversarial region location or features in non-key frames by estimating the motion movement of adversarial regions and features (i.e. optical flow in this work). 2) Intra-frame: In observing the redundant computations of benign features during LISF-based detection and recovery, we propose the efficient hardware design for benign feature computation reuse and eliminate redundant unnecessary computation and memory traffic. Themis architecture can be easily integrated with existing DNN accelerators to friendly support state-of-the-art performance-oriented or accuracy-oriented video recognition methodologies In summary, this work has the following contributions:

- We propose the LISF-based methodology to accurately identify the adversarial region locations in input data. The detection methodology can effectively work under adaptive attack when an adversary has white-box knowledge of defensive approaches.
- We propose the defensive framework that significantly reduces defending overhead in non-key frames by leveraging temporal and spatial locality in video data, which can be friendly integrated with the state-of-the-art video object detection frameworks.
- We propose lightweight hardware customization to efficiently support the defensive framework and fully exploit the computation reuse of benign features, which can be easily adapted to existing deep learning accelerators.
- We make an extensive experimental evaluation and the results show that *Themis* system can effectively and efficiently defend adversarial attacks in real-time. Compared to the system without defensive mechanisms, *Themis* improves the average mAP of real-time object detection from 0.03 to 0.66, with 1.08% of hardware overhead.

1. Themis, (Greek: "Order") in Greek religion, personification of justice, a **blindfolded** goddess holding a pair of scales.

## 2 BACKGROUND

### 2.1 Video Object Detection Preliminary

Video object detection, recognizing instances of visual objects (e.g, humans, cars, animals) and their locations in digital videos, is a fundamental important computer vision task. It forms the basis of many other computer vision tasks, such as instance segmentation [13], object tracking [14], and image captioning [15], etc.

**Object Detection in a Single Image:** Image object detection is the foundation of detecting objects in videos. Image object detector solves the following two subtasks: 1) Predicting how many objects are in the images. 2) Classifying these objects and estimating their locations with bounding boxes. The most important evaluation metric for object detector prediction accuracy is mAP (mean Average Precision) which considers both precision and recall rate. Recently, video object detection capability has been largely boosted by deep learning techniques with milestones of CNN-based image detectors, such as RCNN, YOLO [16], SSD [17], RetinaNet [18], etc.

**Object Detection in Video Data:** Video data consists of many time-sequential images. Hence, video object detection can be achieved by performing image object detection in every image (referred to as accuracy-oriented (AO) schema). For better computing efficiency, prior video object detection methodologies are proposed to leverage temporal redundancy across frames in time-sensitive applications [19], [20], [21] (referred to as performance-oriented (PO) schema). The key design concept of PO schema is to undergo the precise computation of key frames, while approximately computing non-key frames based on key frame features and the motion, trajectory, or optic flow information.

**Optical Flow:** Optical flow is widely used to utilize the temporal redundancy in video data and to approximate the non-key frames.Optical flow describes the apparent motion of image objects in consecutive frames caused by the movement of objects or the camera. Specifically, as shown in Fig. 1, optical flow(c) is a 2D vector field where each vector is a displacement (dx, dy), showing the movement of pixels from first frame (a) to second (b). Once we obtain the optical flow, the pixels or feature map of non-key frames can be estimated by warping the pixels or feature map of their predecessor key frame with optical flow:

$$V(x+dx, y+dy, k+1) \triangleq V(x, y, k)$$

where $V$ can be the pixel values or feature activation and $\triangleq$ can be implemented with different interpolation methods. The optical flow should be resized to warp the feature maps with different size. Optical flow can be calculated based on pixel matching [22] or CNN [23], with the assumption that objects maintain the same intensity (brightness) between consecutive frames.

**AO and PO Video Object Detection Details:** Fig. 2 illustrates accuracy-oriented (AO) and performance-oriented (PO) video object detection frameworks. For AO framework, the image pixels of every frame are input to the full video object detector for precise computation. For PO framework, the DNN-based video object detector is divided into two parts: DNN-prefix for more generic feature extraction with much heavier computation and DNN-suffix for final prediction with low computation overhead [21], [24]. PO framework first calculates the optic flow between non-key frames and the key frame. Then, warps the key-frame features with the resized and scaled optic flow information to approximately compute the non-key-frame features. The predicted non-key-frame
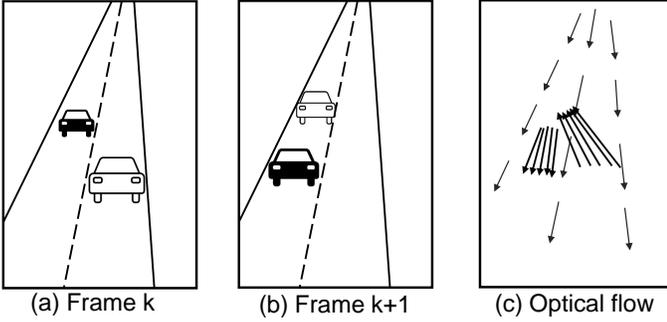
Fig. 1. Optical flow (2D vector field) describes the motion of objects in two frames.
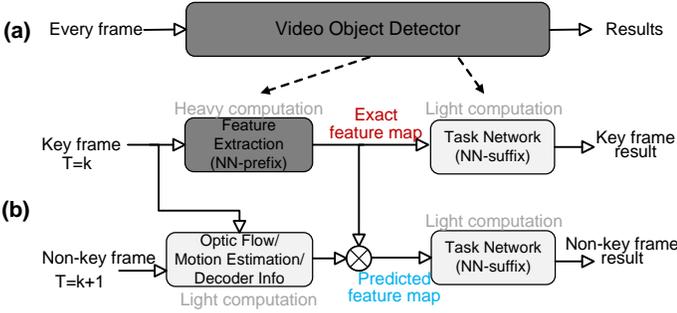


Fig. 2. Video recognition framework: (a) Accuracy-oriented (AO) framework precisely computes every frame; (b) Performance-oriented (PO) framework precisely computes key frames, while approximately computes non-key frame based on optical flow.

feature map is input to the DNN-suffix for the non-key frame result computation. Due to the large gap between the computation of DNN-prefix and DNN-suffix, the computation overhead of non-key frames is significantly reduced. These two VOD frameworks (image-based and feature-based) are adopted in different scenarios that are optimized for accuracy or performance. It is important to support effective and efficient defensive mechanisms in both these two cases. We propose *Themis* framework to achieve this goal.

## 2.2 Adversarial Attacks in Video Recognition

**Attack Formalization.** Adversarial patch attack can largely damage video recognition tasks by evading video object detectors. It manipulates the victim model to output malicious results by adding the patch perturbation in the object of video data. Formally, the goal of adversarial patch attack is to generate the adversarial patch, $\hat{P}$, to maximize the expectation of possibility for classifier $h$ to output targeted malicious label $y_p$ with all adversarial inputs $x_p$ derived from data set $X$.

$$\hat{P} = arg \max_p E_X[logPr(h(x_p) = y_p|x_p)] \qquad (1)$$

Patched image $x_p$ is generated by applying patch $p$ to $x$ in the input dataset $X$, which can be formalized as:

$$x_p = A(p,x), x \in X \qquad (2)$$

where $p$ is the adversarial patch, $x$ is the clean image, and A is the transformation function applying the adversarial patch on the clean image (environmental noises, resizing, rotations, and deformations).

The adversarial patch determines the prediction results with a very small region of pixels (adversarial region) for a relatively broad range of input images. By moving out the adversarial region, the adversarial effects are eliminated and the prediction results are

recovered. Hence, autonomously identifying the adversarial region in the video data is the key foundation of robust object detection.

**Patch Attacks vs. Example Attacks.** Compared to *adversarial example attacks* [25] that have been largely studied in image classification tasks, patch attacks have the following advantages: 1) *Better universality*, because the adversarial patch is independent of input images. Such a scene-independent feature enables physical-world attack without prior knowledge of the scene. Adversarial example attacks, on the other hand, generate perturbation noises highly dependent on the input images, which hinders their deployment in the physical world. 2) *Better robustness to environmental noises and geometric distortion*. For example, human being wearing the T-shirt printed with the adversarial patch can be ignored in the object detection systems in different environments and body gestures [4]. It is not only effective in a static figure input, but also in complex scenarios like walking, sitting, and running [4]. More studies prove that adversarial patches are robust to not only environmental noises but also geometric distortion.

Adversarial example attacks, however, are scene-dependent and transfer poorly in different inputs. In real cases, the adversary neither can obtain the attack scene in advance nor compute the adversarial examples for every frame in real-time. Therefore, adversarial example attack is not a practical attack model in physical environments for video recognition tasks. In this following, we do not consider the adversarial example attacks.

**Limitation of Existing Defenses.** Although some prior researches propose the adversarial example detection methodology [26], [27], the differences between adversarial patch and the adversarial example noises hinder these countermeasures' applicability in the patch attacks. Additionally, these countermeasures can only recognize whether the input is a benign image or not, but cannot recover the DNN system from the attack effect. In this study, we aim to rescue the video recognition tasks by not only detecting the malicious input but also eliminating the adversarial effect by moving out the malicious patches. More adversarial defense comparisons are introduced in Section 7.

## 3 ADVERSARIAL VIDEO PATCH CHARACTERIZATION

Intuitively, patched images rely on the extremely localized important neurons in the adversarial regions to deceive and induce the object detector to output the incorrect results, while benign images perform stable object detection without relying on extremely localized important neurons. Hence, we perform the LISF characterization and have the following observations:

1) *LISFs are good candidates for detecting adversarial regions in single frame.*

LISF Distribution: The important neurons in feature maps of superficial layers exhibit a localized pattern in patched images while scattered in benign images. The superficial important features refer to the neurons that contribute significantly to the feature map value in the superficial layers (in this paper, we use the first layer). Specifically, we take the superficial important neurons as the Top-K ones having the biggest value in the output feature map of the first layer. We visualize these important neurons of a patched image example in Fig. 3(a). Intuitively, the patched image exhibits an extremely localized region of important neurons. Further, we make statistical counting about the distribution of Top-200 neurons of 12K images randomly selected from FLIC dataset [28] based on two metrics: cluster distance and cluster number. We compare the

**(a)**



**(b)**

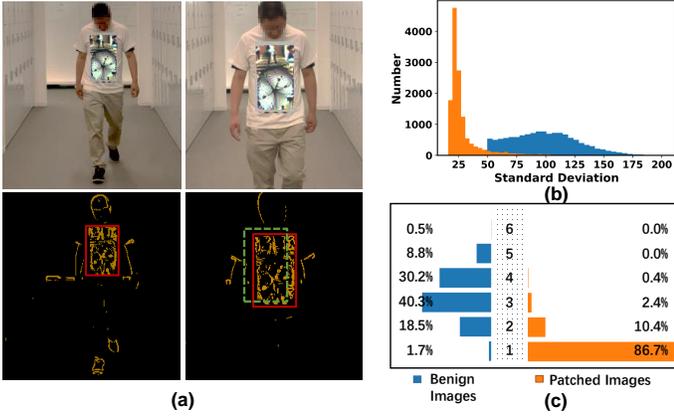| | | |
|---|---|---|
| 0.5% | 6 | 0.0% |
| 8.8% | 5 | 0.0% |
| 30.2% | 4 | 0.4% |
| 40.3% | 3 | 2.4% |
| 18.5% | 2 | 10.4% |
| 1.7% | 1 | 86.7% |

■ Benign Images   ■ Patched Images

**(c)**

Fig. 3. Distribution of important neurons. (a) A showcase of the important neuron distribution and temporal association of LISFs along consecutive frames. (b) The distance deviation of the important neurons in 12K images. (c) The average important neuron cluster number distribution in 12K images.



Fig. 4. Prediction stability of occluding LISFs: benign inputs vs. patched inputs.

standard deviation of the distance between the highlighted nodes to the central node in benign images and patched images. The patched images' distance deviation is much smaller than benign images, as shown in Fig. 3(b). We also cluster the Top-k nodes with the classic MeanShift clustering algorithm. As shown in Fig. 3(c), for patched images, about 86% of cases have only one cluster and 97% of cases have no more than two clusters. While for benign images, about 80% of the cases have more than three clusters. Both of these two metrics show the localized distribution of superficial important neurons in patched images.

Prediction Stability by moving out LISFs: When occluding the localized superficial important features, patched images can be recovered mostly without affecting the prediction accuracy of benign images. As shown in Fig. 4, we compare the detection rate of benign objects and patched objects before and after moving out the localized important features. The detection rate of benign objects decreases to 97.4% slightly, while the detection rate of patched data increases from 17.9% to 72.8% significantly, which means that the prediction results of benign objects are much more stable than patched objects. For benign objects, the detection rate is not sensitive to the localized features. For the patched objects, patch effect would been effectively eliminated after moving the localized important superficial features.

2) *LISFs exhibit temporal association in video data, which enables us to leverage temporal redundancy to eliminate the recovery overhead in non-key frames.* Superficial feature computing is very closed to the input, hence the important superficial features exhibit the similar temporal association as the image frames in video data. The bottom two subfigures in Fig. 3 (a) show the LISFs marked within the solid red boxes in two consecutive frames. The LISF in the following frames can be predicted by warping the optical flow with the LISF in the previous key frame (predicted LISF marked in green box).

In summary, LISF-based methodology not only effectively detects adversarial region and recovers the object detection in key frames, but is also able to leverage temporal redundancy in videos and eliminate defensive overhead in non-key frames. Based on these two observations, we then propose the LISF-based robust video objector detection system.
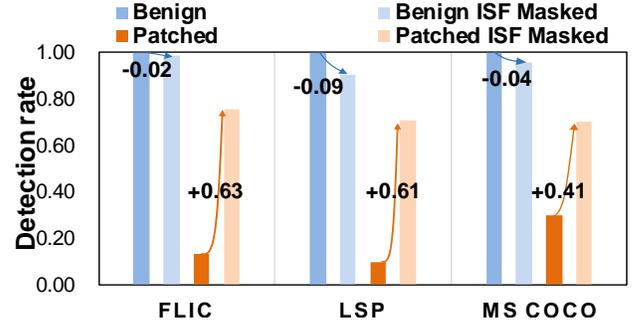
## 4 THEMIS SYSTEM ARCHITECTURE

### 4.1 Overview

The overall *Themis* system is shown in Fig. 5 with algorithm, framework, and hardware designs.

*Algorithms:* LISF-based methodology effectively targets the adversarial regions in input images and recovers the prediction results by moving out those regions.

*Framework:* Although the *Themis* algorithm already reduces the adversarial region searching space with LISF-based methodology, multiple inferences are introduced in the occluding testing stage and incur large overhead. To achieve the goal of real-time robust object detection, *Themis* proposes systematic design to efficiently remove redundant computations for better computing efficiency.

To further reduce the detection overhead, *Themis* proposes the inter-frame and intra-frame optimizations to reduce redundant computations. 1) Inter-frame optimization: By leveraging the spatial and temporal locality between frames, *Themis* only performs complete adversarial detection in key frames. For non-key frames, *Themis* framework either predicts the adversarial region locations based on regions detected in key frames (image-based warping) or reuse the clean features in key frames after eliminating the adversarial effects (feature-based warping). These two warp strategies can be easily integrated with existing AO and PO video object detection frameworks. 2) Intra-frame optimization: The algorithm introduces a large volume of redundant calculations of benign features during occluding prediction stage for the key frames. Hence, *Themis* scheduler reduces the computing overhead with computation reuse of benign features (Section 4.2). Additionally, *Themis* provides the interfaces of setting the following configurations for better feasibility: key frame proportion, the video object detector model segmentation strategy, and the configurable parameters in adversarial patch detection algorithms.

*Hardware:* Although *Themis* algorithm can be implemented in pure software, it is inefficient because of the following reasons: 1) Searching the heat-map of the input activation for the adversarial candidates is inefficient in the DNN accelerator due to the lack of the computing parallelism between the searching process and typical inference process; 2) During the voting stage, multiple patch candidate regions will be occluded to calculate the inference results. In this process, a significant volume of benign features is computed repeatedly and introducing large unnecessary performance overhead. To address these issues, we propose the *Themis* hardware architecture for efficient defense. The top-level block diagram of hardware architecture is shown in Fig. 5 (c). Apart from the typical DNN accelerators with PE array, scalar function unit (SFU), global buffer, and control logic, *Themis* is augmented with

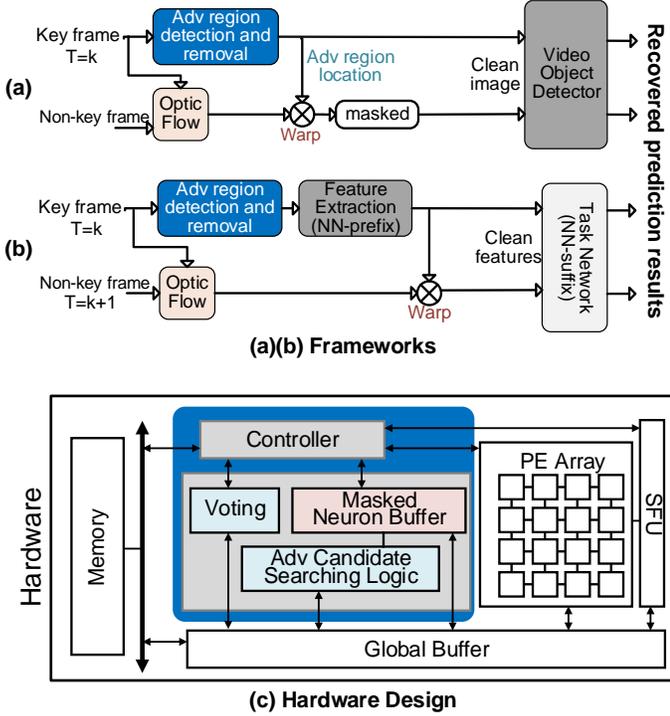**(a)(b) Frameworks**



**(c) Hardware Design**

Fig. 5. Themis Architecture Overview. Themis reduce the defending overhead in non-key frames for both AO framework (a) and PO framework (b); Hardware optimization for benign feature computation reuse (c).

the LISF searching logic to search the candidate regions according to the heat map of the first superficial layer feature map, Masked Neuron Buffer to optimize the data and computation reuse of benign features, and the voting logic to decide the final prediction results.

## 4.2 Framework

**Inter-frame Optimization based on Tempo-Spatial Redundancy.** The overview of *Themis* framework is shown in Fig. 5(a) and Fig. 5(b) with support for AO and PO video object detection frameworks. Under both AO and PO scenarios, the defensive mechanisms for key frames are the same: *Themis* algorithm detects the adversarial region of frames and masks them to eliminate the adversarial effect. For non-key frames, *Themis* proposes two warp strategies with optical flow information which are friendly integrated in existing AO and PO video object detection frameworks, as shown in Fig. 5(a) and Fig. 5(b).

1) *Image-based Warping in AO Frameworks:* For AO framework, we approximately estimate the patch location in non-key frames by warping the patch location in key-frame images with the optical flow. Specifically, every frame will be forwarded to the object detector for a complete inference. *Themis* only performs the adversarial patch detection in key frames and obtains the adversarial region location. Then warps the detected adversarial region in key-frames with the optical flow information to estimate the adversarial region in non-key frames. The corresponding area of non-key frames is masked and the masked non-key frames are forwarded to the object detector for inference. The rationality is that the adversarial regions in the input data also exhibit temporal and spatial redundancy in live vision. Moreover, compared to the object to be segmented, the adversarial region is much less sensitive

to the derivation brought by the inaccurate optical flow information (more validation in Section 6.1).

2) *Feature-based Warping in PO Frameworks:* For PO framework, we directly warp the clean features in key-frames with the resized and scaled optical flow to compute the non-key-frames features. Then, the predicted feature map is forwarded to the DNN-suffix for the final object detection results. Under both AO and PO cases, the defensive overhead of non-key frames is minimized.

**Intra-frame optimization for Benign Feature Computation Reuse.** With the obtained coordinates of the patch candidates, *Themis* then makes the prediction decisions by masking them from the original image individually. During this process, multiple inference rounds of inputs with different mask locations are introduced, which incurs a large execution overhead. Hence, we propose the scheduling methodology with computation reuse to alleviate the execution overhead and eliminate redundant computations. As shown in Fig. 7, the masked images and the original image share the most same pixels, with only the differences of masked regions. Thus, we compute the masked regions separately and splice the masked region back to the complete feature map for the final prediction results, to eliminate the most redundant computations. The detailed mask region computing is as follows: Given an image (224×224 pixels) with the masked regions (50×50 pixels). The pixels in the masked region are set as 0 to occlude their effect on the results. However, the features of the masked region through neural network layers are not simply set to 0, because the the existence of weight bias and the kernels that larger than 1. So we need to take the padding number into consideration to compute the features of the masked region. For example, when computing the C1 layer in Fig. 7 for masked images, the complete region (masked region + padded region) for recomputation is 53×53, with an additional purple part. The detailed computing process and dataflow optimization are in Section 4.4.

**Reconfigurable Design Knobs.** *Themis* framework provides the interfaces to configure the following design knobs that affect overall performance efficiency: the key frame proportion, the strategy of dividing the object detector to feature extraction DNN prefix and classification DNN suffix, the candidate numbers during adversarial patch detection.

*Key frame:* We set the key frame rate as 10% with fixed length. *Themis* framework can support the adaptive key frame strategies proposed by prior video object detection studies [21]. However, how to choose the key frame is out of the scope of this work.

*PO framework splitting:* Designing DNN prefix and suffix in PO frameworks is the trade-off between performance and detection accuracy. When the prefix is close to the classification layers, the feature map is too small and may introduce larger accuracy degradation during optical flow estimation. We test extensive datasets and set the splitting spot when the feature map size is smaller than 56×56.

*Adversarial region detection parameters:* In the adversarial patch detection stage, the more adversarial candidates, the more rounds of the additional inferences and larger performance overhead are. The number of the patch candidates is determined by the heat spot selected threshold ($\beta$ and $\theta$). We set $\beta = 0.75$ and $\theta = 0.85$.

## 4.3 Adversarial Detection and Recovery

In observing that important superficial features exhibit different spatial distribution characteristics and exert distinct influence on the

prediction results in benign and adversarial input data, we propose the LISF-based patch candidate searching methodology and then detect and eliminate the adversarial effect by occluding testing.

**LISF-based Patch Candidate Searching.** We perform the LISF searching in the output feature map of the first layer. The size of the searching window is the upper limit of the patch sizes and the searching stride is 1. When the number of important neurons in one searching window is larger than the threshold ($\theta$), it is recognized as an important window and marked as the patch candidate. When several important windows are overlapped, the central important window will be retained as patch candidate with all the other deleted. The detailed searching logic implementation is described in Section 4.4.

**Occluding Testing and Recovery.** We recover the prediction results with the following two steps:

1) *Masked image execution.* After obtaining the candidate locations of adversarial patches, we generate masked images by occluding the patch candidate locations individually from the original images. These masked images are taken into the victim model to produce the prediction decisions.

2) *Monopolist-occluded voting.* Themis performs the prediction decision analysis to detect patched images by examining whether there is a monopolist patch candidate that determines the prediction results. It is true that such a methodology will introduce the true-negatives. However, the results show that such cases are rare and *Themis* can achieve good detection effectiveness. The detailed patch candidate searching and voting mechanisms are introduced as follows. There are $k$ candidates: $P_1$, $P_2$, ..., $P_k$. The prediction result of the original image is $L_0$, and the corresponding prediction results of masked images: $L_1$, $L_2$, ..., $L_k$. We first detect the patch by checking whether there is a $L_i$ distinct from others, while all the other labels are the same: If positive, the $P_i$ is the monopolist that dominates the prediction results, which is the adversarial patch. Only when tearing it off the image, the classifier predicts the robust and benign label $L_i$.

If there is only one candidate, i.e. $k = 1$, we compare $L_1$ and $L_0$ to determine whether there is an adversarial patch. If they are different, $P_1$ is the adversarial patch and $L_1$ would be recovered label.

If there is no such particular label (either all the labels are the same, or several labels are different), no monopolist is detected. Then the image is recognized as a benign image. Themis then performs the majority voting to obtain the predicted label.

## 4.4 Hardware Design

The Themis architecture can be integrated into the typical DNN hardware accelerator to support real-time detection with small overhead.

**LISF searching logic.** The LISF searching logic outputs the coordinates of clustered important neurons, which infers the possible candidate locations in the input image. LISF searching consists of the following three steps:

1) Obtain the binary important (output) feature map of the first layer. Computing Top-K neurons in the feature map is time-consuming and introduces large hardware overhead. Therefore, we make the estimation about the adaptive threshold according to the maximum value of the feature map. All neurons with activation values larger than the threshold ($\beta * feature_{max}$) are selected as important neurons. The important neurons map is stored in the buffer of LISF searching logic.
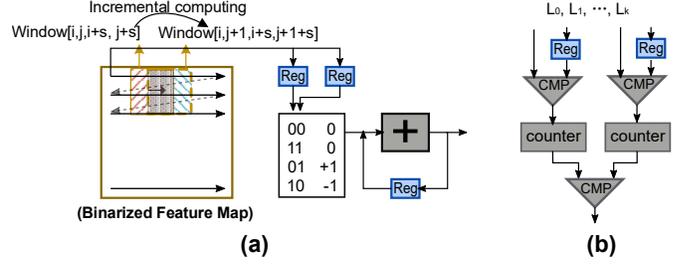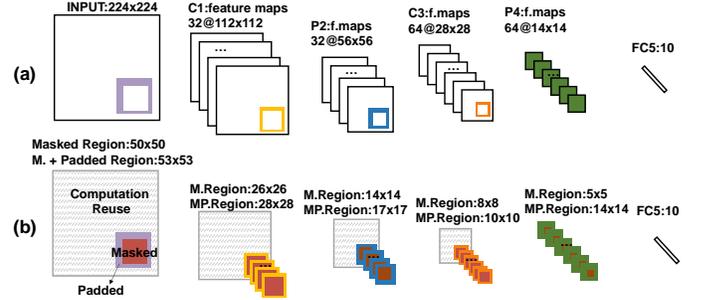


Fig. 6. Searching&Voting logic.



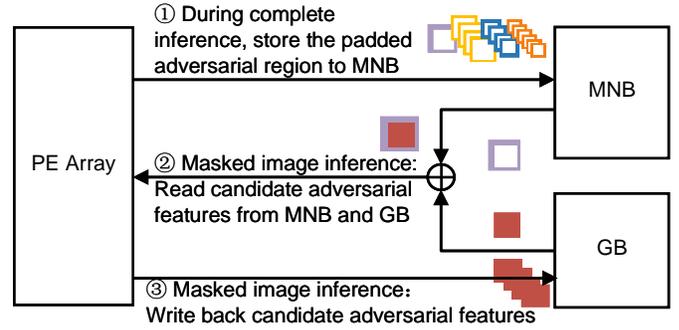Fig. 7. Computing flow for original images and the masked images.



Fig. 8. Data reuse flow

2) Identify the important windows. We slide the fixed-size window to make statistic counting about the important neuron numbers in one window. When the important neuron numbers occupy more than the threshold ($\theta$) of the total neurons, we mark this window as an important neuron window, which will be highlighted as the patch candidates. To reduce the hardware overhead, we make the incremental accumulation of the important neurons in one window. As shown in Fig. 6(a), with the number of important neurons in Window[i,j,i+s,j+s], to compute the Window[i,j+1,i+s,j+1+s], we simply subtract the important neurons in column j+1 and add the important neurons in column j+1+s between row i and row i+s.

3) Delete the overlapped important windows. When two sliding important windows have more than 30% of the area overlapped, we take them as one single patch candidate.

Noted, the LISF searching is not on the critical path of DNN inference, which is parallelized with the processing of the original images.

**Masked Neuron Buffer (MNB).** To support efficient computation reuse of benign features and optimize the data accesses of the adversarial features, the masked neuron buffer is proposed to buffer the padded bounding area for feature computation of candidate adversarial regions.

The data reuse flow of computing candidate adversarial regions

is as follows (Fig. 8): With the coordinates of potential adversarial region candidates, the padded bounding areas through all the neural network layers in Fig. 7(a) are determined. Through the typical inference, all the activation values of the padded bounding areas are stored in Masked Neuron Buffer. Then, during the inferences of masked images, all the masked regions are batched for computation. For every layer, the PE array reads the candidate adversarial features from global buffers and the padded bounding areas from masked neuron buffers, as the red box and purple box shown in Fig. 8②. The PE array combines these neurons by padding the red box with the purple box, takes it as input, and computes the candidate adversarial feature of next layer. After completing the computation of this layer, the adversarial features are stored back to global buffers. Because the padded bounding area is very small compared to the full activation map, we configure the mask neuron buffer with a size of 8KB for every PE array.

**Voting Logic.** The basic voting logic is as shown in Fig. 6(b). We use a compactor array to perform the pairwise comparison between the prediction labels of masked images ($L_0$, $L_1$, ..., $L_k$). If there is an orphan label $L_i$, this image is identified as a patched image and the recovered label is $L_i$. Otherwise, it is a benign image and Themis performs the majority voting to obtain the recovered label.

**Computing Flow.** The overall computing flow is as follows: taking in a new input image, the DNN accelerator first performs the normal inference procedure on the image. Meanwhile, the searching logic obtains the candidate masked regions in the input based on the localized superficial important features. During the inference procedure of the original images, the feature data of corresponding padded bounding box are stored in the masked neuron buffer. After the completion of the original image inference, masked image inference rounds start to get their prediction results. After that, taking in those prediction results, voting logic detects the adversarial patches and output the recovered results.

For one input image, searching and voting operations are only performed once, while the inference rounds are determined by the searching candidates. Among these stages, multiple inferences introduce the most performance cost. Voting is simple and performed within several cycles. Although the searching algorithm is time-consuming when offloaded to the CPU platform, its customized hardware largely boosts the performance and the overhead is less than 0.5% of the multiple inferences. The detailed experimental evaluation is illustrated in Fig. 15.

# 5 EXPERIMENTAL METHODOLOGY

In the following sections, we will evaluate the defensive effectiveness and architecture efficiency of *Themis*, which is complementary to enable the robust and real-time video object detection.

## 5.1 Validation on Algorithm Accuracy

We test the attack success rate and the defensive effectiveness in both single-frame object detection tasks and the video object recognition tasks. For the previous scenario, we focus on exploring and validating the adversarial patch detection capability of *Themis* algorithm on static images. In the latter scenario, we focus on exploring the defensive effectiveness on the non-key frames when *Themis* framework leverages the temporal and spatial information in video data.

**Attack Methodologies:** For single-frame testing scenario, we adopt the digital-synthesized attack methodology that randomly

TABLE 1
Hardware Platform Configurations

| Mem | Global Buffer/Array | RegisterFile/PE | Bandwidth | MNB/Array |
|---|---|---|---|---|
| | 64KB | 256B | 26GB/s | 8KB |
| Comp | # of Arrays | # of PEs/Array | DataType | |
| | 2×2 | 24×24 | INT8 | |

**MNB** denotes **M**asked **N**euron **B**uffer.

attaches the digital adversarial patches onto the bounding box regions of the objects in MS COCO, FLIC, LSP datasets and random locations of images in ImageNet dataset with random rotated angles. The patch size is scaled with the area of bounding boxes ranging from 33x33 to 130x130 pixels. Since we focus on the defensive effectiveness of *Themis*, we only perform attacks on the objects or images that can be correctly identified by the detector. For video data testing scenario, we adopt physical attack videos released in the state-of-the-art attack methodology [4] (Adv T-shirt), which significantly damages the functionality of YOLOv2 object detector. The patch size is variable and the adversarial patch has been significantly deformed during the human movement.

**Optical Flow Methodologies:** *Themis* framework is compatible with different optical flow methodologies. We use both the CV-based and DNN-based optic flow methodologies: DIS [22] and SpyNet [23], to validate the defensive effectiveness and architecture efficiency. The SpyNet is customized with scaled input image sizes and reduced pyramid levels.

**Evaluation Metrics:** We adopt the commonly-used metrics, the detection rate and the standard mean average precision (mAP) score, to measure the accuracy of video object detection. In terms of performance efficiency, the frame per second (fps) is used to indicate how fast the framework process the video data.

## 5.2 Validation on Architecture Efficiency

**NN Accelerator Hardware Implementation.** Our methodology can be generalized adopted in diverse neural network accelerators. In this work, we evaluate our methodology based on the classic Eyeriss accelerator. Specifically, we augment the Eyeriss accelerator with the patch detector consisting of the LISF searching, masked neuron buffer, and voting logic. To prove the generality and its low overhead, we test both the server and edge accelerators with different configurations, as shown in Table 1.

We implement the *Themis* hardware design in Verilog RTL. To obtain the area and power, we synthesis and place&route the RTL code with Synopsys toolchains under TSMC 28nm technology. We use CACTI 7 and DESTINY to model the DRAM memory and on-chip SRAM buffers. Due to the unbearable long duration of silicon simulation, we also tailor an open-sourced simulator, nn-dataflow to support *Themis* for the total execution latency simulation. The simulator also reports the exact memory traces and module activities, which are then used to calculate dynamic energy consumption.

Besides, we deploy *Themis* system on an off-the-shelf FPGA, Zynq UltraScale+ MPSoC ZCU104 board, to show the real performance. Specifically, we implement the proposed architecture components like Adv Candidate Search Logic and integrate them with a DPU (a soft IP for NN inference) using Vitis-ai framework.

# 6 EXPERIMENTAL RESULTS

We first show the defensive effectiveness of *Themis* in terms of single frame and video scenarios (Section 6.1). Then we show the
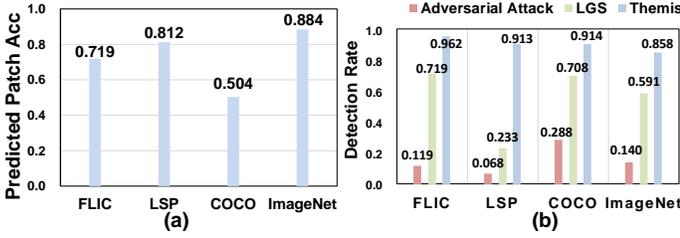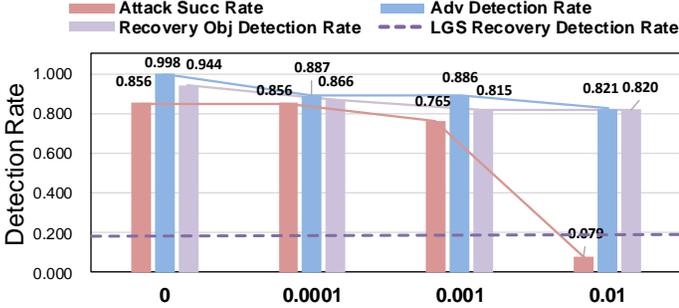
Fig. 9. Detection Effectiveness.



Fig. 10. Defensive Effectiveness under Adaptive Attack.The x-axis refers to the penalty coefficient $\alpha$.

performance and energy efficiency of *Themis* (Section 6.2), which introduces negligible overhead to real-time video object detection. Finally, we show that *Themis* adds negligible area overhead to the baseline DNN accelerator(Section 6.2.4).

## 6.1 Defensive Effectiveness Evaluation

We validate the defensive effectiveness of Themis under the following two scenarios: single-frame data and the video data with sequential frames.

### 6.1.1 Single-frame Defensive Effectiveness

For single-frame testing, we use FLIC [28], LSP [29], MS COCO [30] and ImageNet [31] datasets that are commonly used in object detection domain. We first evaluate the adversarial patch detection accuracy based on the metric of overlapped area proportion in Fig. 9(a). The predicted patch area and the actual patch area have an average of 72%, 81%, 50% and 88% overlapped region compared to the actual patch area size. The results show that LISF-based searching methodology can accurately identify the location of adversarial patches. We then show the object detection rate before and after *Themis* defense in Fig. 9(b) compared with Local Gradient Smoothing (LGS) [32] method. LGS locates the patch using local gradient of image pixels with the basic assumption that patch pixels are not smoothing. The detailed attack methodology is as follows: the adversary randomly attaches the adversarial patch in the person bounding box of the images in the datasets, so that the object detectors are evaded to ignore the persons. With the adversarial patch attack, the object detection rate is 11.9%, 6.8%, 22.8%, 14.0% for FLIC, LSP, MS COCO and ImageNet. Compared to LGS that improves the detection rates to 71.9%, 23.3%, 70.8%, 59.1%. , our *Themis* defensive mechanisms works better and improves them to 96.2%, 91.3% , 91.4% and 85.8%. The results show that *Themis* can eliminate the adversarial patch effect effectively.

**Defensive Effectiveness Against Adaptive Attacks.** In the further step, we evaluate the defensive effectiveness under the

strong adaptive attack [33] that the adversary gets known the full knowledge of the defensive strategies. To steer clear of the detection of *Themis*, the adversary trains the adversarial patch with Equation(3) that a penalty loss for the superficial activation value of patch is considered compared to Equation(1). In this way, the adversary aims to build the adversarial patch with good poisoning effects, but also trying to escape from the adversary candidate searching. $\alpha$ is the parameter to control the scale of the penalty loss in superficial activation value.

$$loss = -logPr(h(x_p) = y_p) + \alpha * \sum(w_1 * patch) \qquad (3)$$

We perform the adaptive attacks on ImageNet dataset and the results are shown in Fig. 10. Adversary detection rate refers to the rate that *Themis* correctly identifies the adversarial region in the adversarial inputs or the benign input with no adversarial region. Recovery object detection rate refers to the rate that *Themis* correctly identifies the object in the image. The results show that, it is indeed that both the adversary detection rate and the recovery object detection rate decrease to 82.1% and 82.0% respectively, when $\alpha$ increases from 0 to 0.01. However, compared to the gentle slope of adversary detection rate and the recovery object detection rate, the attack success rate decreases much more drastically. When $\alpha$ is set to 0.01, the adversarial attack success rate is dropped to 7.9%. As a comparison, after adding total variation loss, the detection rate of LGS drops signficantly from 59.1% to 19.8% while attack success rate maintains high. These results indicate that the adversary cannot maintain the two goals of high attack success rate and good stealthiness simultaneously. *Themis* can work effectively even under adaptive attack.

### 6.1.2 Video Frame Defensive Effectiveness

For the video frames, we adopt the adversarial video benchmarks in the state-of-the-art adversarial attack study [4], where the people wearing the adversarial T-shirt moving in indoor and outdoor scenarios and perform the practical attacks in the physical environment. The video object detector is based on YOLOv2. We evaluate the object detection rate and mAP under the following scenarios:

- *AO-ND*: AO framework with no defensive mechanisms.
- *AO-Full*: Defensive AO framework that examines every frames.
- *AO-Dis*: Defensive AO framework that completely examines key frames, but CV-based methodology to predict the adversarial patch locations in non-key frame.
- *AO-Spynet*: Defensive AO framework that completely examines key frames, but DNN-based optical flow information to predict the adversarial patch locations in non-key frame.
- *PO-ND*: PO framework with no defensive mechanisms.
- *PO-Spynet*: Defensive PO framework with DNN-based optical flow (spynet) .

**Detection Rate:** Fig. 11 shows the object detection rate under different scenarios. With adversarial attacks, both AO and PO object detectors have significant low object detection rates of 4.4% and 4.8%, which indicates that adversarial attacks can essentially damage the integrity and functionality of object detectors even in the physical environments. With *Themis* defensive algorithm, the detection rate is significantly improved above 93.8% for different defensive strategies. Compared to AO-Full that examines every frame, other approximate defensive methods achieve relatively good object detection rates within a gap of less than 5%.
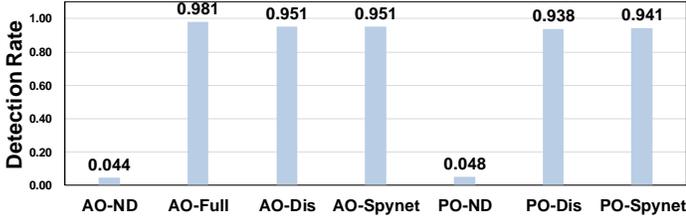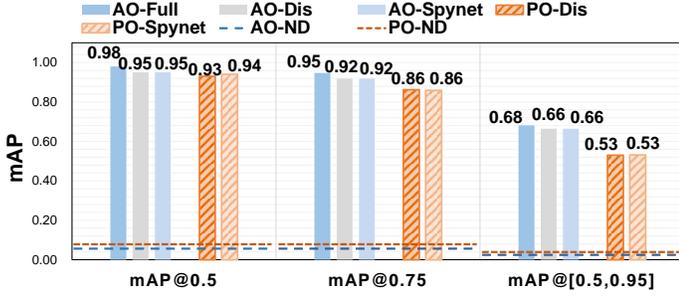
Fig. 11. Detection rate in video-frames.



Fig. 12. Defensive effectiveness in video-frames .



(a) Performance comparison



(b) Energy comparison

Fig. 13. Performance and Energy comparison among different defensive strategies.

Detection rate is only a coarse-grained metric that intuitively indicates the detection recall rate of object detectors. In the further step, we evaluate the mAP that considers both the prediction accuracy, recall rate, and the predicted bounding box accuracy in the following.

**mAP:** Fig. 12 shows the mAP results under IoU = 0.5 (mAP@0.5), IoU = 0.75 (mAP@0.75), and average mAP value where IoU ranges from 0.5 to 0.95 with the step of 0.05. IoU (Intersections over Union) is the metric to determine whether it is an accurate prediction of the bounding box, which is calculated as the rate of dividing the area of overlap by area of union. A larger IoU indicates a more strict criterion of mAP prediction accuracy. From the plot, specifically, we have the following observations:

1) Consistently, it is observed that adversarial attacks can effectively fool the object detector to ignore the human being with mAP as low as 0.03, 0.05 of AO-ND and PO-ND. With the guard of *Themis*, the functionality of object detector is recovered and the average mAP is improved to the range of (0.53, 0.68) under different defensive mechanisms.
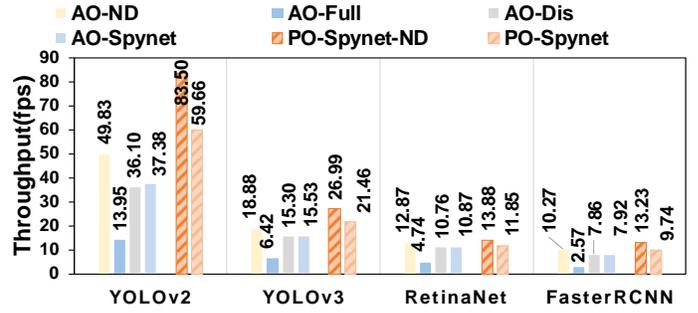
2) When IoU is low (IoU=0.5), the mAP of defensive PO frameworks is equally good to the defensive approaches that examine every frame (AO-Full). When IoU is high, PO series may introduce the reduction of mAP due to the shift and deviation of the optical flow information. Specifically, the gap between mAP@0.75 and the average mAP of PO-Spynet and AO-Full is 0.09 and 0.15.

3) Adversarial patch location prediction is less sensitive to the deviation of optical flow. Although optical flow information also introduces the deviation between the predicted and actual adversarial regions, such deviation does not markedly hurt the defensive effectiveness of *Themis* (less than 0.03), because of the prediction instability of adversarial regions, as analysis in Section 3.
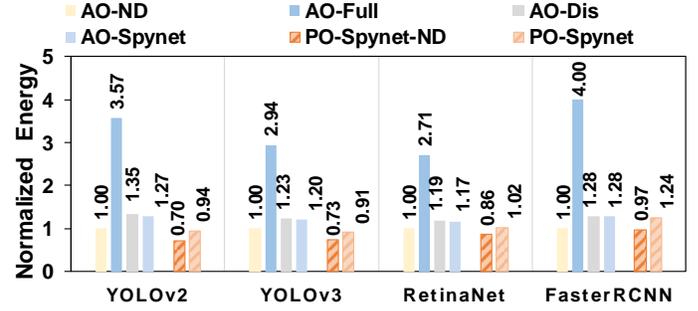
## 6.2 Architecture Efficiency Evaluation

### 6.2.1 Overall Performance and Energy Comparison

We also evaluate the performance and energy of AO and PO frameworks with different defensive strategies. Fig. 13(a) shows the object detection throughput for four commonly-used video

detectors: YOLOv2, YOLOv3, RetinaNet, and FasterRCNN. From the plot, we have the following conclusions:

1) PO frameworks significantly boost the performance when their model architectures can be divided into heavy NN-prefix and light NN-suffix.

2) Performing adversarial detection in every frame incurs heavy overhead. Compared to AO-ND, AO-Full incurs 3.3x execution latency, which remarkably reduces the fps in all four object detectors.

3) Eliminating the unnecessary recomputing for the adversarial region locations in non-key frames improves the performance and reduces the performance gap between defensive approaches with original approaches, while maintains the defensive effectiveness. Specifically, AO-SpyNet introduce 25.0%, 17.7%, 15.5%, 22.9% overhead compared to AO-ND. PO-SpyNet introduces 28.6%, 20.4%, 14.6%, and 26.4% overhead compared to PO-SpyNet-ND.

In summary, *Themis* can effectively defend against adversarial attacks in video tasks, while still maintain the throughput of about 36 fps and 59 fps for real-time object detection in AO and PO frameworks.

### 6.2.2 Scheduling Optimization

We first evaluate the effectiveness of scheduling optimization for computing masked images in four typical datasets: MS COCO, FLIC, LSP and T-shirt [4]. Fig. 14 shows the latency reduction ratio with benign feature computation reuse for different object detection models. Patches on MS COCO, FLIC, LSP, and T-shirt account for 2.81%, 1.47%, 3.60%, and 2.01% of the whole image pixels on average, respectively. We have the following observations:

1) *Overall, Themis reduces the latency for masked images effectively with scheduling optimization.* The average latency reduction ratio for MS COCO, FLIC, LSP, and T-shirt are 51.93%
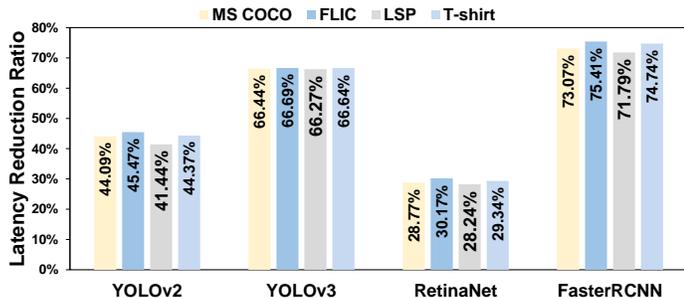
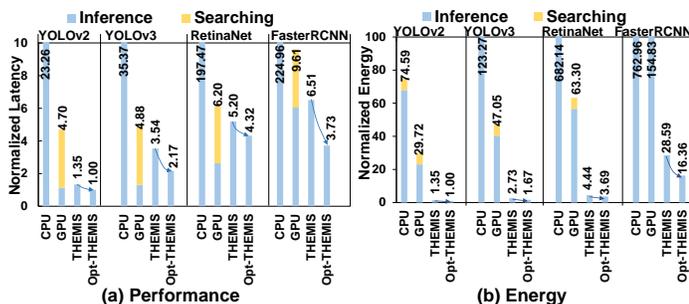Fig. 14. Latency reduction with benign feature computation reuse.



Fig. 15. Performance and Energy comparison with different architectures.

~53.77%, which indicates that the scheduling optimization in *Themis* is applicable to various scenarios.

2) *The scheduling optimization strategy in Themis is more efficient in object detectors with a shallower depth.* The average latency reduction ratio in FasterRCNN is 73.75%, which is significantly higher than 29.13% in RetinaNet. The intrinsic reason is that the computation reuse ratios in the first several layers are higher than the latter. Specifically, for the first layer in YOLOv3, the computation reuse ratio reaches 98.38%.

*Themis* can efficiently reduce latency of the detection process for masked images with scheduling optimization. With an average latency reduction of 53.31%, which significantly reduces the overhead incurred by the defense scheme.

### 6.2.3 Comparison with different Architectures

We compare the performance and energy of *Themis* with CPU (Intel(R) Core(TM) i7-4770K at 3.50GHz) and GPU (NVIDIA TITAN V) platforms. Fig. 15 shows the normalized latency and energy for CPU, GPU and *Themis* architectures running four video detectors: YOLOv2, YOLOv3, RetinaNet, and FasterRCNN with MS COCO dataset. *Themis* hardware components can be integrated with any DNN architectures and we choose Eyeriss-like architecture as the basic DNN accelerator. Specifically, 'Themis' refers to the cases of adopting candidate searching hardware components and 'Themis-Opt' refers to the *Themis* architecture with computation reuse optimization. For each case, we break down the system into inference procedure and searching procedure (CPU data bars are too large in this figure and are rescaled with the value marking on it). We draw the following conclusions: 1) *Themis* is much more efficient compared to CPU and GPU platforms. Compared with CPU, *Themis* achieves speedup from 16.3x to 60.3x. Compared with Titan V, *Themis* reduces energy from 9.47x to 29.7x. 2) Candidate searching consumes non-negligible overhead in GPU platform. The customized searching logic achieves 745.7x speedup. 3) With the computation reuse optimization, Themis-Opt achieves speedup from 1.20x to 1.75x.

### 6.2.4 Hardware Area Overhead

The baseline Eyeriss-like DNN accelerator has an area of 12.60 $mm^2$. On top of the basic DNN accelerator design, *Themis* only introduces a total area overhead of 0.136 $mm^2$, which incurs only 1.08% of hardware overhead. Specifically, the LFI searching logic occupies 0.0746% of the hardware overhead, the masked neuron buffer occupies 1.005%, and the rest is attributed to the voting logic.

### 6.2.5 Implementation on FPGA

The experimental FPGA device (Xilinx ZCU104 evaluation board) consists two ARM processor cores and 16nm FinFET+ programmable logic. The ARM processors are used to run Linux system and control the operation flow, while the programmable logic is reconfigured to accelerate the user applications. We use the Vitis-ai framework to deploy *Themis* with a customized DPU IP. Compared to the officially provided Yolov2 demo that runs at 24 FPS, *Themis* achieves 22 FPS at AO defense mode and 35 FPS at PO defense mode.

## 7 RELATED WORK

### 7.1 Real-time Video Object Detection

Video object detection is the fundamental important computer vision task and grows rapidly with the increasing demands in autonomous driving and video surveillance [34]. Extended from image to video domain, DNN techniques largely boost video recognition capability. YOLO series [16], SSD [17], RCNN series [35], are proposed in rapid succession. Compared to single image object detection, video object detection has new attributes of existing both spatial and temporal correlations within consecutive frames. Previous studies leverage such attributes to improve the performance of video recognition tasks [19], [20], [21]. In addition to the optimized video recognition frameworks that make use of temporal-spatial information for computing efficiency, some recent studies enhance feature maps with tempo-spatial information to improve detection accuracy that originally degraded by motion blur, rare poses, video defocus, etc [34], [36]. Besides, LSTM-based [37], attention-based [38], tracking-based [39] video object detectors are also proposed. *Themis* algorithm can support such video recognition frameworks well, because it is able to directly identify the adversarial region of input and supports feature aggregation or propagation operations during leveraging the tempo-spatial video information.

### 7.2 Defense Against Adversarial Attack

Adversarial patch attack is one of the most practical DNN attack models that can effectively damage object detectors even in physical environments [4]. Envisioning its importance, prior studies contribute to the defense techniques against the adversarial patch attack by 1) building certified robust neural networks that resist the attack effect, such as IBP, de-randomized smoothing, etc [12], [40], [41]; However, all of these certified studies achieve distinctly low certified accuracy for large scale datasets such as ImageNet. A certified accuracy of 20.5%-36.3% from these methods can hardly be applied in the real world systems. 2) adversarial training to obtain robust model again patch attack [42], in which online re-training process introduces unacceptable cost. 3) performing patch detection based on empirical observations to eliminate the attack effect, such as digital watermarking (DW) [43], local gradient

smoothing (LGS) [32]. However, these defenses proved to be invalidated when confronting with the strong adaptive attacker that has the white-box knowledge of the defense [40]. Note that since locality is the intrinsic property of patch attack, our LISF-based detection algorithm prevents the adversary from maintaining both strong attack effect and stealthiness, which promises the effectiveness of *Themis* even under the adaptive attack.

The most related recovery methodology are MRD [10] and ObjectSeeker [11]. MRD [10] requires a large amount of iterative inference passes to obtain the prediction map by masking every small region sliding across the entire original input images. It is extremely time-consuming and costs $1446s$ for detection of one single image with sizes of $224{\times}224$ on an Eyeriss-scale accelerator. Similarly, ObjectSeeker [11] proposes patch-agnostic masking for certified objection detection that needs more than 100 inferences for different bands of the input image. Therefore, both MRD and ObjectSeeker are unacceptable for real-time detection due to their additional large cost.

In summary, existing countermeasures are unable to perform online defense for video recognition tasks. This work proposes a high-efficient and effective detection and recovery system to defend the adversarial attacks that practically introduce damaging consequences in video scenarios.

## 7.3 Feature Importance Analysis

Neuron importance has been widely used for abnormal input detection [44], [45]in previous studies. However, the metric (superficial feature importance) in our methodology is distinct from previous work. Previous studies focus more on the neurons that contribute significantly to the inference output (deep feature importance). We envision that deep feature importance is not a good candidate from the following two aspects: 1) both the benign images and the adversarial images have the deep feature importance. Its discrimination in the benign images and adversarial images is not straightforward, so that it requires more complex computation to identify the benign and adversarial images. 2) calculating such deep feature importance is time-consuming, which demands the gradient information and the complete backward propagation process. We propose the superficial input feature importance as the metric for discrimination analysis based on the intuition that in order to efficiently manage the output prediction results with a very small region of the input data, the adversarial patch must incur large activation from the first place instead of the accumulation of the deep feature extraction.

## 8 CONCLUSION

*Themis* efficiently and accurately recovers the DNN systems from the adversarial attacks with both algorithmic framework and the architectural support. At the algorithmic level, *Themis* prevents the classifier from being overshadowed by the trivial but extremely biased parts by tearing the patch off the original images. At the architectural level, *Themis* not only proposes efficient searching and voting logic, but also proposes the scheduling methodology to accelerate the masked image execution by eliminating the redundant computations and memory traffics. The results show that the proposed methodology can effectively recover the VOD system from the adversarial effect in real-time.

## REFERENCES

[1] Global Management Consulting. Autonomous vehicle adoption study, 2016.

[2] Electrek. Elon musk clarifies tesla's plan for level 5 fully autonomous driving: 2 years away from sleeping in the car, 2017.

[3] Waymo. Introducing waymo's suite of custom-build, self-driving hardware, 2017.

[4] Kaidi Xu, Gaoyuan Zhang, Sijia Liu, Quanfu Fan, Mengshu Sun, Hongge Chen, Pin-Yu Chen, Yanzhi Wang, and Xue Lin. Adversarial t-shirt! evading person detectors in a physical world. In *ECCV*, pages 665–681, 2020.

[5] Zuxuan Wu, Ser-Nam Lim, Larry S Davis, and Tom Goldstein. Making an invisibility cloak: Real world adversarial attacks on object detectors. In *European Conference on Computer Vision*, pages 1–17. Springer, 2020.

[6] Chong Xiang and Prateek Mittal. Patchguard++: Efficient provable attack detection against adversarial patches. *arXiv preprint arXiv:2104.12609*, 2021.

[7] Ahmed Abusnaina, Yuhang Wu, Sunpreet Arora, Yizhen Wang, Fei Wang, Hao Yang, and David Mohaisen. Adversarial example detection using latent neighborhood graph. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7687–7696, 2021.

[8] Jinyu Tian, Jiantao Zhou, Yuanman Li, and Jia Duan. Detecting adversarial examples from sensitivity inconsistency of spatial-transform domain. *arXiv preprint arXiv:2103.04302*, 2021.

[9] Ning Wang, Yimin Chen, Yang Xiao, Yang Hu, Wenjing Lou, and Thomas Hou. Manda: On adversarial example detection for network intrusion detection system. *IEEE Transactions on Dependable and Secure Computing*, 2022.

[10] Michael McCoyd, Won Park, Steven Chen, Neil Shah, Ryan Roggenkemper, Minjune Hwang, Jason Xinyu Liu, and David Wagner. Minority reports defense: Defending against adversarial patches. In *International Conference on Applied Cryptography and Network Security*, pages 564–582. Springer, 2020.

[11] Chong Xiang, Alexander Valtchanov, Saeed Mahloujifar, and Prateek Mittal. Objectseeker: Certifiably robust object detection against patch hiding attacks via patch-agnostic masking. *arXiv preprint arXiv:2202.01811*, 2022.

[12] Chong Xiang, Arjun Nitin Bhagoji, Vikash Sehwag, and Prateek Mittal. Patchguard: A provably robust defense against adversarial patches via small receptive fields and masking. 2021.

[13] Pedro O Pinheiro, Tsung-Yi Lin, Ronan Collobert, and Piotr Dollár. Learning to refine object segments. In *European conference on computer vision*, pages 75–91. Springer, 2016.

[14] Wongun Choi and Silvio Savarese. A unified framework for multi-target tracking and collective activity recognition. In *European Conference on Computer Vision*, pages 215–230. Springer, 2012.

[15] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Neural baby talk. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7219–7228, 2018.

[16] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7263–7271, 2017.

[17] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.

[18] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[19] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. Deep feature flow for video recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2349–2358, 2017.

[20] Zhuoran Song, Feiyang Wu, Xueyuan Liu, Jing Ke, Naifeng Jing, and Xiaoyao Liang. Vr-dann: Real-time video recognition via decoder-assisted neural network acceleration. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 698–710. IEEE, 2020.

[21] Mark Buckler, Philip Bedoukian, Suren Jayasuriya, and Adrian Sampson. Eva$^2$: Exploiting temporal redundancy in live computer vision. In *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*, pages 533–546. IEEE, 2018.

[22] Till Kroeger, Radu Timofte, Dengxin Dai, and Luc Van Gool. Fast optical flow using dense inverse search. In *European Conference on Computer Vision*, pages 471–488. Springer, 2016.

[23] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4161–4170, 2017.

[24] Yu Feng, Paul Whatmough, and Yuhao Zhu. Asv: Accelerated stereo vision system. In *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*, pages 643–656, 2019.

[25] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *ICLR*, 2015.

[26] Bita Darvish Rouhani, Mohammad Samragh, Mojan Javaheripi, Tara Javidi, and Farinaz Koushanfar. Deepfense: Online accelerated defense against adversarial deep learning. In *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pages 1–8. IEEE, 2018.

[27] Xingbin Wang, Rui Hou, Boyan Zhao, Fengkai Yuan, Jun Zhang, Dan Meng, and Xuehai Qian. Dnnguard: An elastic heterogeneous dnn accelerator architecture against adversarial attacks. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 19–34, 2020.

[28] Ben Sapp and Ben Taskar. Modec: Multimodal decomposable models for human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3674–3681, 2013.

[29] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *Proceedings of the British Machine Vision Conference*, 2010. doi:10.5244/C.24.12.

[30] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[31] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.

[32] Muzammal Naseer, Salman Khan, and Fatih Porikli. Local gradients smoothing: Defense against localized adversarial attacks. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1300–1307, 2019.

[33] Florian Tramer, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1633–1645. Curran Associates, Inc., 2020.

[34] Haidi Zhu, Haoran Wei, Baoqing Li, Xiaobing Yuan, and Nasser Kehtarnavaz. A review of video object detection: Datasets, metrics and methods. *Applied Sciences*, 10(21):7834, 2020.

[35] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(6):1137–1149, 2016.

[36] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen Wei. Flow-guided feature aggregation for video object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 408–417, 2017.

[37] Mason Liu, Menglong Zhu, Marie White, Yinxiao Li, and Dmitry Kalenichenko. Looking fast and slow: Memory-guided mobile video object detection. *arXiv preprint arXiv:1903.10172*, 2019.

[38] Yihong Chen, Yue Cao, Han Hu, and Liwei Wang. Memory enhanced global-local aggregation for video object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10337–10346, 2020.

[39] Wenfei Yang, Bin Liu, Weihai Li, and Nenghai Yu. Tracking assisted faster video object detection. In *2019 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1750–1755. IEEE, 2019.

[40] Ping-Yeh Chiang, Renkun Ni, Ahmed Abdelkader, Chen Zhu, Christoph Studor, and Tom Goldstein. Certified defenses for adversarial patches. In *8th International Conference on Learning Representations (ICLR)*, 2020.

[41] Alexander Levine and Soheil Feizi. (De)randomized smoothing for certifiable defense against patch attacks. In *Conference on Neural Information Processing Systems, (NeurIPS)*, 2020.

[42] Sukrut Rao, David Stutz, and Bernt Schiele. Adversarial training against location-optimized adversarial patches. In *European Conference on Computer Vision*, pages 429–448. Springer, 2020.

[43] Jamie Hayes. On visible adversarial perturbations & digital watermarking. In *2018 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPR Workshops)*, pages 1597–1604, 2018.

[44] Yiming Gan, Yuxian Qiu, Jingwen Leng, Minyi Guo, and Yuhao Zhu. Ptolemy: Architecture support for robust deep learning. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 241–255. IEEE, 2020.

[45] Zirui Xu, Fuxun Yu, and Xiang Chen. Lance: A comprehensive and lightweight cnn defense methodology against physical adversarial attacks on emb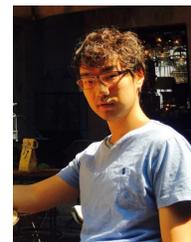edded multimedia applications. In *2020 25th Asia and South Pacific Design Automation Conference (ASP-DAC)*, pages 470–475. IEEE, 2020.

**Husheng Han** received the B.S degree from Tsinghua University, Beijing, China in 2020. Currently he is working toward the PhD degree in the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China, and the University of Chinese Academy of Science, Beijing, China. His current research interests include machine learning security and domain-specific hardware architectures.
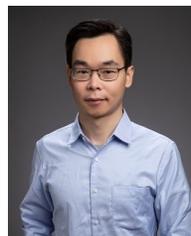
**Xing Hu** received the B.S. degree from Huazhong University of Science and Technology, Wuhan, China, and Ph.D. degree from University of Chinese Academy of Sciences, Beijing, China, in 2009 and 2014, respectively. She is currently an associate professor of State Key Laboratory of Processors, Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS), Beijing, China. Her current research interests include domain-specific hardware architectures and deep learning system.

**Kaidi Xu** received the B.S. and M.S. degrees from Sichuan University in 2015 and the Department of Computer Science at University of Florida in 2017 respectively and Ph.D. degree at Northeastern University in 2021. He is currently an Assistant Professor in Department of Computer Science at Drexel University, Philadelphia, USA. His primary research interest is the robustness of machine learning, including physical adversarial attacks, rigorous robustness verification and certified defenses.

**Pucheng Dang** received the B.S degree from Harbin Institute Of Technology, Harbin, China in 2019. Currently he is working toward the PhD degree in the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China, and the University of Chinese Academy of Science, Beijing, China. His current research interests include machine learning security and Computer Vision.

**Ying Wang** received the B.S and M.S. degree from School of Astronautics, Harbin Institute of Technology, Heilongjiang, China in 2007 and 2009 respectively and the PhD degree from the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China in 2014. He is currently an Associate Professor in State Key Laboratory of Computer Architecture at Institute of Computing Technology, Chinese Academy of Sciences. His research interests primarily focus on the area of reliable computer architecture and VLSI design, with an emphasis on memory systems, energy-efficient accelerators, and approximate/error-tolerant computing.

13

**Yongwei Zhao** recieved the bachelor's degree in computer science and technology from the Huazhong University of Science and Technology, Wuhan, China, in 2015 and the PhD degree from the institute of Computing Technology, Chinese Academy of Sciences, Beijing, China, in 2020. He is currently an assistant professor at the Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS), Beijing, China. His research interests primarily focus on the area of accelerator architecture.

**Zidong Du** (Member, IEEE) received his B.E. degree from the Department of Electronic Engineering, Tsinghua University, Beijing, China, in 2011, and the Ph.D. degree from the Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS), Beijing, in 2016. He is currently an associate professor at Intelligent Processor Research Center, Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS). His research interests mainly focus on novel architecture for artificial intelligence, including deep learning processors, inexact/approximate computing, neural network architecture, neuromorphic architecture. He has published over 20 top-tier computer architecture research papers, including ASPLOS, MICRO, ISCA, TC, TOCS, TCAD. For his innovative works on deep learning processors, he won the best paper award of ASPLOS'14, Distinguished Doctoral Dissertation Award of CAS (40/10000), Distinguished Doctoral Dissertation Award of China Computer Federation (10 per year).

**Qi Guo** (Member, IEEE) received the B.E. degree in computer science from Tongji University, Shanghai, China, in 2007, and the Ph.D. degree from the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China, in 2012.,From 2012 to 2014, he was a Staff Researcher at IBM Research, Beijing. From 2014 to 2015, he was a Postdoctoral Researcher with Carnegie Mellon University, Pittsburgh, PA, USA. He is currently a Professor with the Institute of Computing Technology, Chinese Academy of Sciences. His research interests include computer architecture, programming models, and machine learning.

**Yanzhi Wang** received his B.S. Degree in Electronic Engineering from Tsinghua University, Beijing, China, in 2009 and the Ph.D. Degree in Computer Engineering from University of Southern California (USC) in 2014, under the supervision of Prof. Massoud Pedram. He is currently an Associate Professor and Faculty Fellow in the Department of Electrical and Computer Engineering, and Khoury College of Computer Science (Affiliated) at Northeastern University. His research interests include real-time and energy-efficient deep learning and artificial intelligence systems, model compression of deep neural networks (DNNs), neuromorphic computing and non-von Neumann computing paradigms.

**Tianshi Chen** received the bachelor's degree in mathematics from the Special Class for the Gifted Yong, University of Science and Technology of China (USTC), Hefei, China, in 2005, and the PhD degree in computer science from the Department of Computer Science and Technology, University of Science and Technology of China, Hefei, China, in 2010. He received the China Computer Federation Distinguished Doctoral Dissertation Award, in 2011 and the Chinese Academy of Sciences Distinguished Doctoral Dissertation Award, in 2011 for his PhD work on computational complexity analysis of evolutionary algorithms. He is currently a professor with the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China. He is also serving as the CEO of a startup called Cambricon Technologies Corporation Limited, whose commercial processor products are named "Cambricon".