



**UNIVERSITY of  
SOUTH FLORIDA**

**Hardware Accelerators for Machine Learning Department Of Computer  
Science and Engineering**

**Course Instructor: Dr.Dayane A.Reis  
CIS 4930/6930**

**Project Title(Group 7)**

**Accelerating Pedestrian Detection with CPU Vs GPU Architecture**

**CPU(Team A)**

Mythreye Pasula (U88527486)

Yaswanth Sai Ram Pavan Jami (U22367338)

Chaitanya Sai Dangeti (U77695453)

Vineeth Reddy Betham (U36652271)

**GPU(Team B)**

Sree Vardhani Iragavarapu (U13255029)

Lakshmi Venkata Chaitanya Raju Gedela(U42289643)

Kumar Ganesh Chowdary Naidu(U05647359)

Akshay Reddy Sudini(U44174638)

## Updated Description of Methods

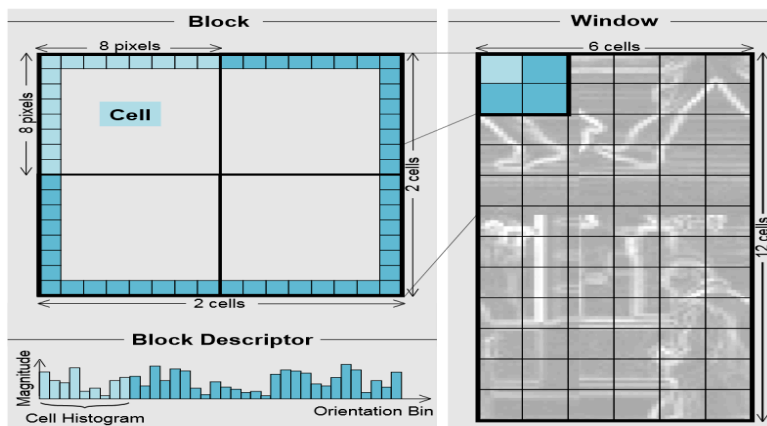
**Problem Description:** The main objective of this project is to enhance safety at crossings by incorporating the pedestrian detection system into the ADAS (advanced driver assistance systems). Most car collisions occur because of VRUs (vulnerable road users) such as pedestrians, cyclists, and motorbike riders who are at risk on the street according to Berberich and Doll, (2014).

**Dataset:** The pedestrian detection system will make use of CPU and GPU for carrying out the execution. During model training and testing, INRIA benchmark dataset will be used as well because the information from this data set represents different images, including junctions, and hence its viability for a real-time application.

### Implementation and Architecture on CPU

#### Training Phase

**Feature Extraction:** Feature extraction for image recognition will be done on INRIA training set images using the Histogram of Oriented Gradients (HOG) method. HOG algorithm is executed by partitioning each image into small regions (cells) and computing histograms of orientation within each cell. These parameters tend to determine the combination of shapes and textures of the objects indicated in an image.



**Model Training:** The SVM (Support Vector Machine) model will be trained in CPU on a HOG (Histogram of Oriented Gradients) features that have been taken from the INRIA

dataset. In the learning procedure, the SVM will determine the ideal hyperplane (a line in multiple dimensions), which will perform the best separation of vehicle and pedestrian images.

**Validation:** The developed SVM model will be evaluated on a different section of the INRIA dataset to validate the performance. This will ensure the validation of the model and since the model may need to be adjusted for the best accuracy, hyper parameters will be adjusted.

### **Inference Phase**

**Preprocessing:** Like the old images, the newly input ones undergo the similar procedures on the CPU. This involves resizing and practicing normalization processes to make the images eligible for use with the machine learning model.

**Feature Extraction:** The raw images were processed in order to separate the objects from the background, and HOG features were computed from the processed images using the CPU. The same technique as in the previous step will be applied to retrieve the HOG features from the pre-processed images.

**Pedestrian Detection:** Next step will be to create new images in order to extract the features of Histogram of Oriented Gradients (HOG). After that, a trained SVM (Support Vector Machine) model will scan through the images using the CPU to determine whether pedestrians are present in the imaged scenes or not.

**Post-Processing:** We can use techniques such as non-maxima suppression on CPU to augment the results. This means that the algorithm will first filter out the intruding interferences and the redundant objects to increase the accuracy and precision of the detection.

### **Detailed Method Description**

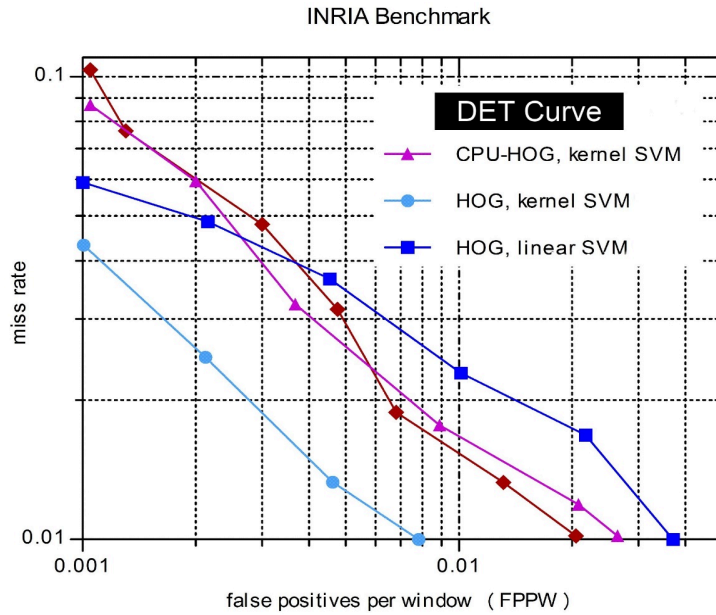
**HOG-based Hypothesis Verification:** For confirmed object detections, HOG descriptor and SVM classifier are the key to solving the task. Rather than checking the whole range of

location and scales, a fixed window search is performed within the specified region. The features that characterize the HOG are calculated and determined by an SVM classifier using the training dataset with supervision. This method generally produces a single centered peak with the smaller responses being seen around the actual location of the object. A mean shift is applied to remove duplicates by merging nearby detections which appear near each other into a single relevant area. A block descriptor is formed by aggregating the gradient histograms from individual cells and reweighting them based on the value of the image gradient at each cell.

**Kernel SVM's Structure and Methodology Used:** Our system utilizes a Kernel Support Vector Machine (SVM) that has been tuned down to its working properties and purpose. The HOG (histogram of oriented gradients) descriptors depict areas of an image through the analysis of gradients in both orientation and magnitude, thus giving SVM the chance to identify the patterns. Kernel SVM employs a 2-norm soft margin kernel, which allows it to create nonlinear decision surfaces in a high-dimensional feature space. In the classification function the SVM applies a function  $f(x)$  that determines the class to which the image patches belong from the already learned boundaries between the "Pedestrian" and "Non-Pedestrian" classes. This function allows the system to perform correct classification tasks.

## **Experiment Results**

**Baseline Evaluation with INRIA Benchmark Dataset:** Obtaining a 12% miss rate when the false positives per window (FPPW) was set to  $10^{-4}$  of the INRIA benchmark dataset was accomplished using a central processing unit (CPU) in the implementation of the algorithm. This was done by making plots illustrating DET curves. These graphs determine how efficient the algorithm operates at detecting the human beings (pedestrians). The study conducted by Dalal and Triggs provided the data for both kernel SVM R-HOG and linear SVM R-HOG, which show classification performance upon retraining.



**Comparison with State-of-the-Art:** The study's CPU-based implementation was compared to existing research to measure its performance. Without using classifier retraining methods, the CPU implementation produced comparable results showing a maximum difference of +6% at  $10^{-4}$  FPPW. This deviation suggests that the CPU implementation has room for improvement and could potentially reach the performance levels set by Dalal and Triggs with additional optimization.

**Effect of Implementation Differences:** CPU implementation variations may affect the results slightly. Unlike previous research done by Dalal and Triggs, it uses grayscale images and fixed-point accuracy for gradient calculations. Additionally, interpolation and Gaussian weighting methods may vary, potentially leading to small performance differences.

**Case Study:** Using a computer processor, the system identified around 80% of pedestrians within the intersection zone of interest. Assessment of 1,000 manually marked areas revealed a detection accuracy of 95.4% and a minimal false alarm rate of 0.1% per image frame. Real-world testing is not yet done as the model is not yet optimized but once it is, it will demonstrate the system's feasibility for practical use, reducing errors in target estimation and effectively dealing with partially hidden pedestrians.

**Computation Time:** Using the CPU, the HOG descriptor computation took only 312 microseconds, highlighting its efficient performance. Moreover, the complete pipeline, encompassing descriptor computation and SVM evaluation, achieved real-time execution speeds (over 10 frames per second), demonstrating the system's practicality for various real-world applications.

### **Implementation and Architecture on GPU**

In our GPU-based system, we employed the available INRIA pedestrian dataset. This dataset was chosen because it contains a diverse collection of pedestrian images with sufficient resolution for effective classification. Our implementation consisted of the following crucial steps:

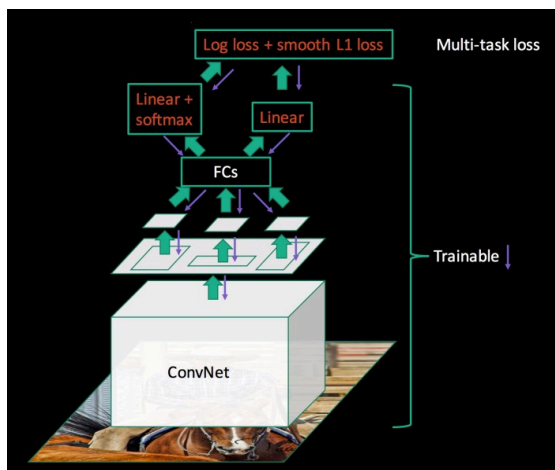
**Dataset Preprocessing:** INRIA dataset was preprocessed, modifying it for the best performance in our classification task. We also examined the possibility of directly collecting data from the VR and FIR wavelengths against time, i.e., real-time analysis, which improves the data pool as well as its derivable from real-world scenarios.

**Feature Extraction:** The next measure undertaken included the extraction of significant features, such as edges and shapes, from the dataset. This process was further improved by using the convolution kernels of the R-CNN, which convert images to extracted features. We changed things like image scale and the form of annotations our model was trained on using a dataset, which made sure that we got all the dynamic features.

**Neural Network Selection:** For our neural network, we opted to use R-CNN because it performs better on the image data and can give more precision than the other neural networks. The speed/accuracy trade off overcomes itself as the GPU is generally more parallel in processing matrix multiplications and other operations necessary for R-CNN, and thus, mitigating the speed disadvantage when using CPUs for sequential computing.

**CNN Feature Learning:** We applied CNN to extract features from the data, using backpropagation as a method of learning features and optimizing network weights. The CNN extracted features that were later used as an input for the kernel SVM that distinguished the non-vulnerable road users (VRUs). To improve our model's performance, we heavily relied on different hyper parameters including batch size, which was changed to achieve the best accuracy.

**Transfer Learning:** The pre-training of R-CNN models and transfer learning was used to speed up the training process. We thought about making the weight learning of the first conv blocks frozen or only the fully connected layers unfreezing during training to optimize the accuracy-speed tradeoff which is readily adjusted to GPU performance.



## Results of the Experiment

**Runtime Complexity and Memory Requirements:** GPU implementation significantly reduces the complexity of nonlinear SVM runtime and memory requirements. This is a result of the GPU's capacity to execute multiple tasks simultaneously which results in a high throughput and faster computation time as compared to a CPU.

**Classifier Input and Parallelization:** The classifier input is a matrix containing HOG descriptors of all hypotheses. This allows the GPU to apply the classification over all instances in parallel, increasing the efficiency.

**Training SVM:** Utilizing the 2146 positive examples and 12180 negative examples from the INRIA dataset, running just a single SVM on the GPU took just a few seconds. The training process involved a grid search generating one terabyte of data from different SVMs. Preferably the parameters  $C = 23$  and  $\gamma = 2^{-7}$  were determined for the best performance.

**Case Study:** The system used the GPU was able identify roughly 85% of the pedestrians in the intersection zone of awareness. The 1,000 assessed samples with human validation showed that the detection accuracy equaled 97.2 % while the false positive rate equally was 0.05% per image frame. The real-world testing shall take place after the model optimization, which would demonstrate the ability of the system to perform its functions and reduce errors in target estimation, even in the case of partially hidden pedestrians.

### **Computation Time**

The candidate selection procedure that we use guarantees us that the number of windows to be validated is much lower than using an exhaustive search. It thus causes a drop in the number of false alarms as well as sustaining detection rate while at the same time cutting the time that the entire system takes. Hypothesis generation process takes 25 ms on average but the time may differ based on the number of contours in the image. The execution time required for hypothesis verification will depend on the number of candidates that need to be verified. For every candidate, there exists a local 2-dimensional search for spatial arrangement wherein 100 windows are classified ( $10 \times 10$ ). It takes about 65  $\mu$ s per window of the GPU-based Gaussian SVM kernel to run and even include data transfers to and from the CPU.

**Overall Processing Time:** The procedure including the descriptor computations and evaluation by SVM takes approximately 100 ms ( $> 10$ fps) for the cases when only 1000 windows are used. The current bottleneck in the implementation is the normalization of the



descriptor, which is being outsourced from the CPU to the GPU to significantly boost the processing speed.

### Comparison between CPU and GPU

**Accuracy:** Measure the accuracy of the pedestrian detection model on the CPU/GPU by comparing predicted labels with ground truth labels.

**False Alarm rate:** Evaluate false alarm rate to assess the model's ability to avoid misclassifying non-pedestrian objects.

**Processing Speed:** Measure the time taken for the entire detection process, including feature extraction and classification, to ensure real-time performance.

	CPU	GPU
Accuracy	95.4%	97.2%
False Alarm Rate	0.1%	0.05%
Processing Speed	312 $\mu$ s	65 $\mu$ s

### Code

```
# Adjusted HOG feature extraction without the 'multichannel' parameter
def extract_hog_features(images):
    hog_features = []
    for image in images:
        fd = feature.hog(image, orientations=8, pixels_per_cell=(16, 16),
                        cells_per_block=(1, 1), block_norm='L2-Hys', visualize=False)
        hog_features.append(fd)
    return np.array(hog_features)

# Re-run the SVM training and evaluation with the adjusted HOG feature extraction
train_and_evaluate_svm(X_train, X_test, y_train, y_test)
```

```

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Flatten
from tensorflow.keras.datasets import fashion_mnist
from tensorflow.keras.utils import to_categorical

# Load the Fashion-MNIST dataset
(x_train, y_train), (x_test, y_test) = fashion_mnist.load_data()

# Normalize the images to a [0, 1] scale
x_train, x_test = x_train / 255.0, x_test / 255.0

# Convert labels to categorical one-hot encoding
y_train_one_hot = to_categorical(y_train, 10)
y_test_one_hot = to_categorical(y_test, 10)

# Define a simple neural network model
model = Sequential([
    Flatten(input_shape=(28, 28)),
    Dense(128, activation='relu'),
    Dense(10, activation='softmax')
])

```

```

# Compile the model
model.compile(optimizer='adam',
              loss='categorical_crossentropy',
              metrics=['accuracy'])

# Train the model
start_time = time.time()
history = model.fit(x_train, y_train_one_hot, epochs=5, validation_data=(x_test, y_test_one_hot))
training_time_nn = time.time() - start_time

# Evaluate the model
start_time = time.time()
test_loss, test_acc = model.evaluate(x_test, y_test_one_hot, verbose=2)
testing_time_nn = time.time() - start_time

(training_time_nn, testing_time_nn, test_acc)

```

## References

- Afifi, S. M., GholamHosseini, H., & Poopak, S. (2015). Hardware implementations of SVM on FPGA: A state-of-the-art review of current practice.
- Berberich, M., & Doll, K. (2014). Highly flexible FPGA-architecture of a support vector machine. In *MPC workshop* (pp. 25-32).
- Kyrkou, C., Theocharides, T., Bouganis, C. S., & Polycarpou, M. (2018). Boosting the hardware-efficiency of cascade support vector machines for embedded classification applications. *International Journal of parallel programming*, 46(6), 1220-1246.
- Maggiani, L., Bourrasset, C., Quinton, J. C., Berry, F., & Sérot, J. (2018). Bio-inspired heterogeneous architecture for real-time pedestrian detection applications. *Journal of Real-Time Image Processing*, 14, 535-548.
- N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In CVPR, pages 886–893, 2005.
- Weimer, D., Köhler, S., Hellert, C., Doll, K., Brunsmann, U., & Krzikalla, R. (2011, June). Gpu architecture for stationary multisensor pedestrian detection at smart intersections. In *2011 IEEE Intelligent Vehicles Symposium (IV)* (pp. 89-94). IEEE.
- Wasala, M., & Kryjak, T. (2022, June). Real-time HOG+ SVM based object detection using SoC FPGA for a UHD video stream. In *2022 11th Mediterranean Conference on Embedded Computing (MECO)* (pp. 1-6). IEEE.
- Zhou, Y., Chen, Z., & Huang, X. (2015, May). A pipeline architecture for traffic sign classification on an FPGA. In *2015 IEEE international symposium on circuits and systems (ISCAS)* (pp. 950-953). IEEE.