

# Secure Multiparty Computation Meets Deep Learning

Yoshi234

2024-01-30

# Table of contents

<b>Preface</b>	<b>3</b>
Contents . . . . .	3
Resources . . . . .	3
<b>1 Object Detection Algorithms</b>	<b>4</b>
1.1 SSD: Single Shot MultiBox Detector . . . . .	4
1.2 A Comprehensive Review of YOLO (v1 to v8+) . . . . .	4
1.2.1 Applications of YOLO . . . . .	4
1.2.2 Evaluation Metrics . . . . .	5
1.2.3 Non-Maximum Suppression . . . . .	6
1.2.4 YOLO . . . . .	7
<b>2 Red Light Violation Detection</b>	<b>15</b>
2.1 Motivation . . . . .	15
2.2 V2I Algorithms for RLR Detection . . . . .	15
2.2.1 Red Light Violation Detection Algorithm . . . . .	15
2.3 Thao et.al on Traffic Violation Detection . . . . .	16
2.3.1 Problem Setting . . . . .	16
2.3.2 System Design and Solution Approach . . . . .	16
2.3.3 Primary Contributions . . . . .	18
2.4 Goma et.al on RLR Detection with SSD (Single Shot Detector) . . . . .	18
2.4.1 Methods and System Design . . . . .	18
<b>3 Traffic Flow Forecasting</b>	<b>19</b>
3.1 Attention Based Spatial-Temporal Graph Convolutional Networks for Traffic Flow Forecasting . . . . .	19
3.1.1 Core Contributions of ASTGCN . . . . .	19
<b>4 Summary</b>	<b>20</b>
<b>References</b>	<b>21</b>

# Preface

This is a Quarto book. To learn more about Quarto books visit <https://quarto.org/docs/books>.

## Contents

The quarto book contains an organized structure of notes on a variety of secure multiparty computation protocols, implementation details, and a discussion of v2x applications and relevant deep learning models. The authors hope that any readers find these resources useful.

## Resources

The [matcha editor](#) is used to construct some of the mathematical diagrams shown in this book. In order to export a matcha diagram, you need to enter the full-screen mode for the diagram, and click on the “export” drop-down which becomes available. This will allow you to save the math diagram as a png image. We will make liberal use of these throughout the book, especially for explaining complex security primitives such as beaver’s triples and homomorphic encryption.

# 1 Object Detection Algorithms

## 1.1 SSD: Single Shot MultiBox Detector

The single shot multibox detector is an algorithm presented by W. Liu (2016) for the purpose of taking a 300 x 300 input and generating bounding boxes on objects of interest within the image. The paper is linked [here](#)

## 1.2 A Comprehensive Review of YOLO (v1 to v8+)

J. Terven (2023) present review and analysis of the evolution of the Yolo algorithm, with a focus on the innovations and contributions made by each iteration, as well as the major changes in network architecture (and training tricks) which have been implemented over time.

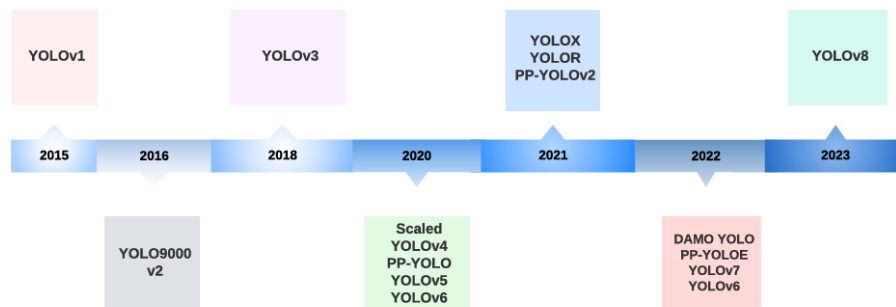


Figure 1: A timeline of YOLO versions.

Figure 1.1: A timeline of YOLO development

### 1.2.1 Applications of YOLO

Yolo has proven invaluable for a number of different applications

- autonomous vehicles
  - enables quick identification and tracking of objects like vehicles, pedestrians, bicycles and other obstacles
- action recognition
- video surveillance
- sports analysis
- human-computer interaction
- crop, disease, pest detection and classification
- face detection - biometrics, security, facial recognition
- cancer detection
- skin segmentation
- pill identification
- remote sensing
  - satellite and aerial imagery object detection / classification
  - land use mapping
  - urban planning
- security systems
- smart transportation systems
- robotics and drones

## 1.2.2 Evaluation Metrics

Average Precision (AP) and Mean Average Precision (mAP) are the most common metrics used in the object detection task. It measures average precision across all categories, providing a single value to compare different models

### 1.2.2.1 How AP works

- mAP is the average precision for accuracy of predictions across all classes of objects contained within an image
  - individual AP values are determined for each category separately.
- IOU (intersection over union)
  - measures the proportion of the predicted bounding box which overlaps which overlaps with the true bounding box

Different methods are used to compute AP when evaluating object detection methods on the COCO and VOC datasets (PASCAL-VOC)

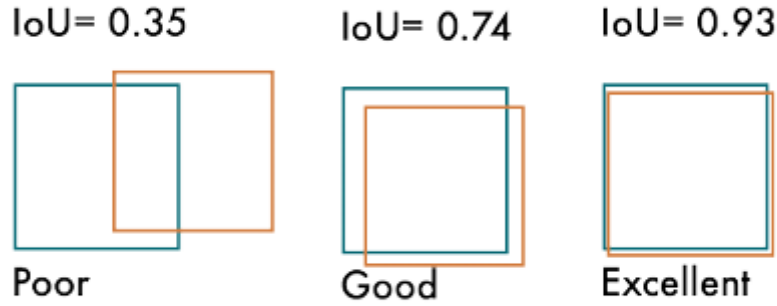


Figure 1.2: Intersection over union in practice

### 1.2.3 Non-Maximum Suppression

A post-processing technique - reduces number of overlapping boxes and improves detection quality. Object detectors typically generate multiple bounding boxes around the same object. Non-max suppression picks the best ones and gets rid of the others.

The algorithm for this is defined below:

---

**Algorithm 1** Non-Maximum Suppression Algorithm

---

**Require:** Set of predicted bounding boxes  $B$ , confidence scores  $S$ , IoU threshold  $\tau$ , confidence threshold  $T$   
**Ensure:** Set of filtered bounding boxes  $F$

```

1:  $F \leftarrow \emptyset$ 
2: Filter the boxes:  $B \leftarrow \{b \in B \mid S(b) \geq T\}$ 
3: Sort the boxes  $B$  by their confidence scores in descending order
4: while  $B \neq \emptyset$  do
5:   Select the box  $b$  with the highest confidence score
6:   Add  $b$  to the set of final boxes  $F$ :  $F \leftarrow F \cup \{b\}$ 
7:   Remove  $b$  from the set of boxes  $B$ :  $B \leftarrow B - \{b\}$ 
8:   for all remaining boxes  $r$  in  $B$  do
9:     Calculate the IoU between  $b$  and  $r$ :  $iou \leftarrow IoU(b, r)$ 
10:    if  $iou \geq \tau$  then
11:      Remove  $r$  from the set of boxes  $B$ :  $B \leftarrow B - \{r\}$ 
12:    end if
13:  end for
14: end while

```

---

Figure 1.3: Non-Max Suppression Alg

A useful visualization is also provided:



Figure 1.4: Non-Max Supression Vis

## 1.2.4 YOLO

The original authors of YOLO titled it as such for the reason that it only required a single pass on the image to accomplish the detection task. This is contrast to the other approaches used by Fast R-CNN and sliding window methods.

The output coordinates of the bounding box were detected using more straightforward regression techniques

### 1.2.4.1 YOLOv1

AP: 63.4%

YOLOv1 predicted all bounding boxes simultaneously by the following process:

1. divide image into  $S \times S$  grid
2. predict  $B$  bounding boxes of the same class and confidence for  $C$  different classes per grid element
3. each bounding box had five values:
  1.  $P_c$  - confidence score for the bounding box - how likely it contains an object and the accuracy of the box
  2.  $bx$  and  $by$  - coordinates of center of box relative to grid cell.
  3.  $bh$  and  $bw$  - height and width of box relative to full
4. output an  $S \times S \times (B \times 5 + C)$  tensor
5. (optional) NMS used to remove redundant bounding boxes

Here is an example of that output:

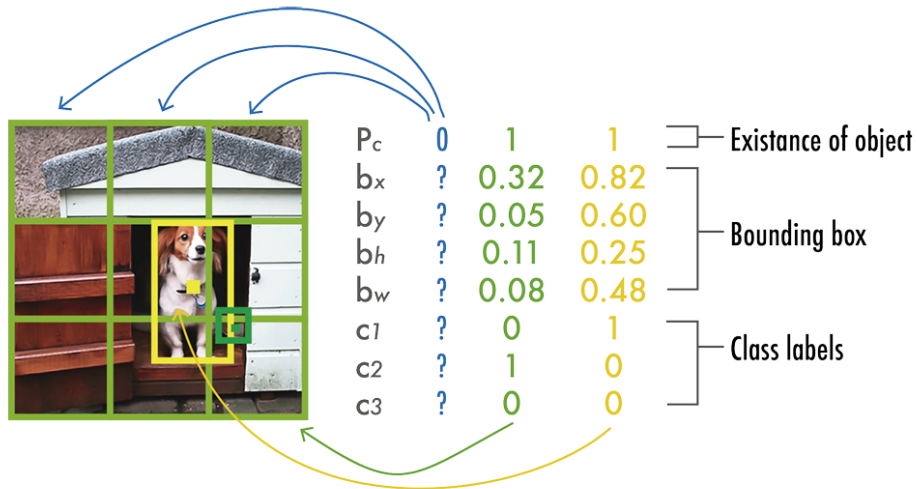


Figure 1.5: yolo output prediction

#### 1.2.4.1.1 v1 Architecture

##### Normal Architecture

- 24 conv layers
  - $1 \times 1$  conv layers are used - reduce number of feature maps and keep parameters lower
  - leaky rectified linear unit activations
- 2 fc layers
  - predict bounding box coordinates / probs
  - linear activation function for final layer

##### FastYOLO

- Used 9 conv layers instead of 24 for greater speed (at the cost of reduced accuracy)

#### 1.2.4.1.2 v1 Training

Basic training process:

1. pretrain first 20 layers at resolution  $224 \times 224$  with *ImageNet* dataset
2. add last four layers with randomly initialized weights - fine tune model with PASCAL VOC 2007 and PASCAL VOC 2012 at resolution  $448 \times 448$



	Type	Filters	Size/Stride	Output
	Conv	64	$7 \times 7 / 2$	$224 \times 224$
	Max Pool		$2 \times 2 / 2$	$112 \times 112$
	Conv	192	$3 \times 3 / 1$	$112 \times 112$
	Max Pool		$2 \times 2 / 2$	$56 \times 56$
1×	Conv	128	$1 \times 1 / 1$	$56 \times 56$
	Conv	256	$3 \times 3 / 1$	$56 \times 56$
	Conv	256	$1 \times 1 / 1$	$56 \times 56$
	Conv	512	$3 \times 3 / 1$	$56 \times 56$
	Max Pool		$2 \times 2 / 2$	$28 \times 28$
4×	Conv	256	$1 \times 1 / 1$	$28 \times 28$
	Conv	512	$3 \times 3 / 1$	$28 \times 28$
	Conv	512	$1 \times 1 / 1$	$28 \times 28$
	Conv	1024	$3 \times 3 / 1$	$28 \times 28$
	Max Pool		$2 \times 2 / 2$	$14 \times 14$
2×	Conv	512	$1 \times 1 / 1$	$14 \times 14$
	Conv	1024	$3 \times 3 / 1$	$14 \times 14$
	Conv	1024	$3 \times 3 / 1$	$14 \times 14$
	Conv	1024	$3 \times 3 / 2$	$7 \times 7$
	Conv	1024	$3 \times 3 / 1$	$7 \times 7$
	Conv	1024	$3 \times 3 / 1$	$7 \times 7$
	FC		4096	4096
	Dropout 0.5			4096
	FC		$7 \times 7 \times 30$	$7 \times 7 \times 30$

Figure 1.6: yolo v1 architecture

Loss functions:

- scaling factors
  - $\lambda_{coord} = 5$  - gives more weight to boxes with objects
  - $\lambda_{noobj} = 0.5$  - reduces importance of boxes with no object
- localization loss:
  - first two terms
  - computes error in predicted bounding box locations  $(x, y)$  and  $(w, h)$
  - only penalizes boxes with objects in them
- confidence loss:
  - confidence error when object is detected (third term)
  - confidence error when no object is in box (fourth term)
- classification loss:
  - squared error of class conditional probabilities for each class if an object appears in the cell

#### 1.2.4.2 YOLOv2 (YOLO 9000)

AP: 78.6%

##### Improvements / Changes

1. Batch normalization - included on all convolutional layers
2. Higher resolution classifier - pretrained model (224 x 224) and then fine-tuned with ImageNet at a higher resolution (448 x 448) for ten epochs
3. fully convolutional - remove dense layers and use fully conv architecture
4. use anchor boxes to predict bounding boxes
  1. anchor box - box with predefined shapes for prototypical objects
  2. defined for each grid cell
  3. system predicts coordinates and class for every anchor box
5. Dimension clusters - pick good anchor boxes using k-means clustering on the training bounding boxes - improves accuracy of bounding boxes
6. Direct Location Prediction
7. Finer-grained features
  1. removed pooling layer - get feature 13 x 13 feature map for 416 x 416 images
  2. passthrough layer - 26 x 26 x 512 feature map -> stack adjacent features into different channels

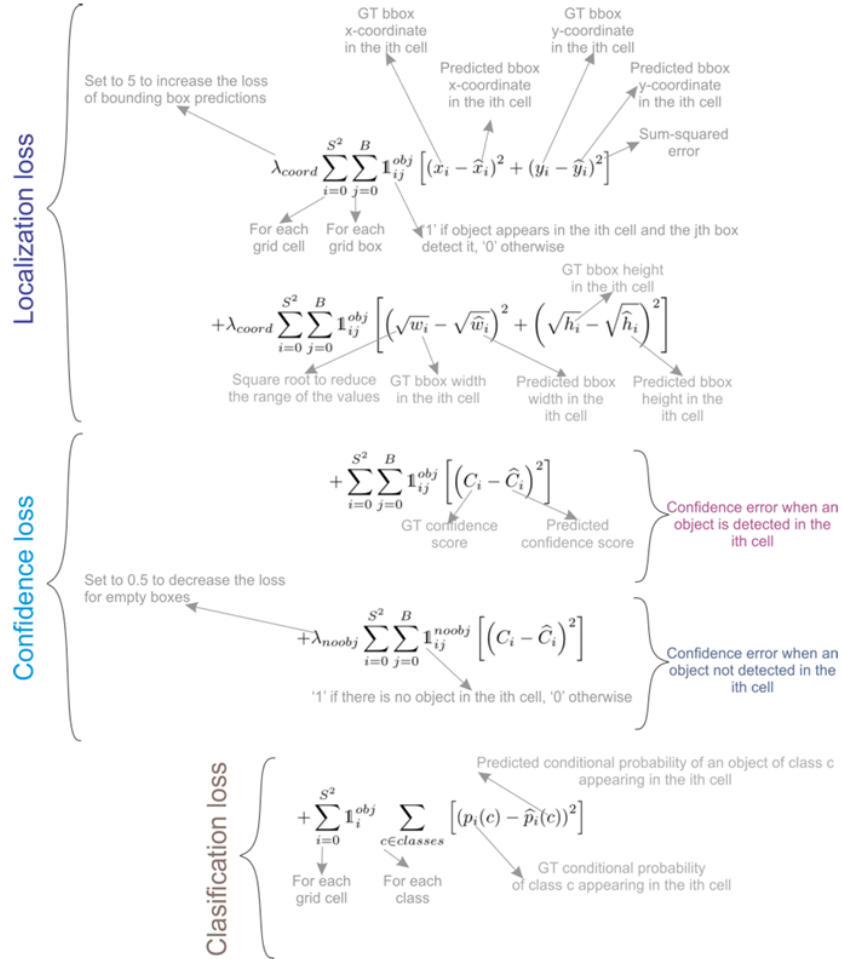


Figure 1.7: yolo v1 loss function

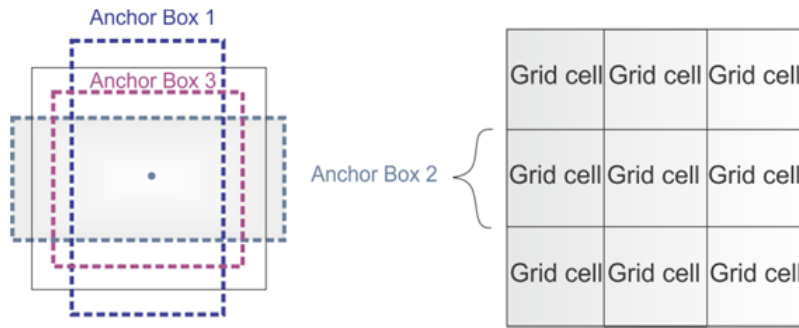


Figure 6: Anchor boxes. YOLOv2 defines multiple anchor boxes for each grid cell.

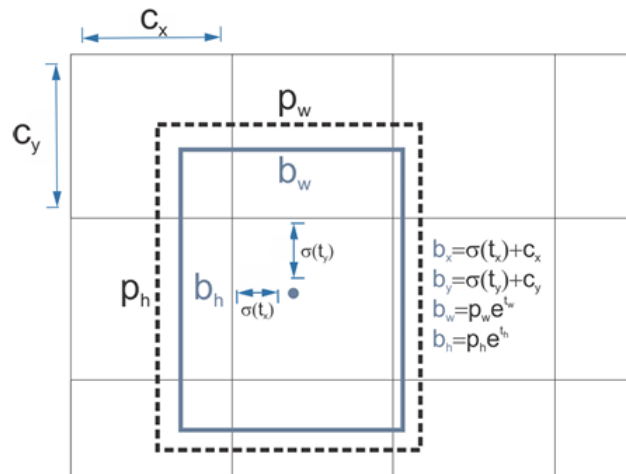


Figure 1.8: yolo v2 anchor boxes

8. Multi-scale training - train on different input sizes to make model robust to different input types

#### **1.2.4.2.1 v2 Architecture**

- backbone architecture -> Darknet-19
  - 19 conv layers
  - 5 max pool layers
    - \* non-linear operation - uses OT to perform efficiently
  - use  $1 \times 1$  conv between  $3 \times 3$  to reduce parameters
  - batch normalization to help convergence
  - object classification head (replaces last 4 conv layers of YOLOv1)
    - \* 1 conv layer (1000 filters)
    - \* GAP layer
    - \* Softmax classifier

Num	Type	Filters	Size/Stride	Output
1	Conv/BN	32	$3 \times 3 / 1$	$416 \times 416 \times 32$
2	Max Pool		$2 \times 2 / 2$	$208 \times 208 \times 32$
3	Conv/BN	64	$3 \times 3 / 1$	$208 \times 208 \times 64$
4	Max Pool		$2 \times 2 / 2$	$104 \times 104 \times 64$
5	Conv/BN	128	$3 \times 3 / 1$	$104 \times 104 \times 128$
6	Conv/BN	64	$1 \times 1 / 1$	$104 \times 104 \times 64$
7	Conv/BN	128	$3 \times 3 / 1$	$104 \times 104 \times 128$
8	Max Pool		$2 \times 2 / 2$	$52 \times 52 \times 128$
9	Conv/BN	256	$3 \times 3 / 1$	$52 \times 52 \times 256$
10	Conv/BN	128	$1 \times 1 / 1$	$52 \times 52 \times 128$
11	Conv/BN	256	$3 \times 3 / 1$	$52 \times 52 \times 256$
12	Max Pool		$2 \times 2 / 2$	$52 \times 52 \times 256$
13	Conv/BN	512	$3 \times 3 / 1$	$26 \times 26 \times 512$
14	Conv/BN	256	$1 \times 1 / 1$	$26 \times 26 \times 256$
15	Conv/BN	512	$3 \times 3 / 1$	$26 \times 26 \times 512$
16	Conv/BN	256	$1 \times 1 / 1$	$26 \times 26 \times 256$
17	Conv/BN	512	$3 \times 3 / 1$	$26 \times 26 \times 512$
18	Max Pool		$2 \times 2 / 2$	$13 \times 13 \times 512$
19	Conv/BN	1024	$3 \times 3 / 1$	$13 \times 13 \times 1024$
20	Conv/BN	512	$1 \times 1 / 1$	$13 \times 13 \times 512$
21	Conv/BN	1024	$3 \times 3 / 1$	$13 \times 13 \times 1024$
22	Conv/BN	512	$1 \times 1 / 1$	$13 \times 13 \times 512$
23	Conv/BN	1024	$3 \times 3 / 1$	$13 \times 13 \times 1024$
24	Conv/BN	1024	$3 \times 3 / 1$	$13 \times 13 \times 1024$
25	Conv/BN	1024	$3 \times 3 / 1$	$13 \times 13 \times 1024$
26	Reorg layer 17			$13 \times 13 \times 2048$
27	Concat 25 and 26			$13 \times 13 \times 3072$
28	Conv/BN	1024	$3 \times 3 / 1$	$13 \times 13 \times 1024$
29	Conv	125	$1 \times 1 / 1$	$13 \times 13 \times 125$

Figure 1.9: yolo v2 architecture

## 2 Red Light Violation Detection

### 2.1 Motivation

Secure Red Light Violation detection is an important application of secure machine learning protocols. Oftentimes, these systems will require the use of image segmentation and object recognition protocols. The images used for this task expose often-times sensitive data about individual users.

In the practical setting of interest, this means exposing license-plate information, associated vehicles, and location information available from the images themselves.

### 2.2 V2I Algorithms for RLR Detection

The authors of this paper have constructed V2I mechanisms for red light running (RLR) detection, wrong way entry (WWE), and an array of other import tasks in the context of V2X. See the citation Dokur and Katkoori (2022)

#### 2.2.1 Red Light Violation Detection Algorithm

The proposed system utilizes the following logic to detect whether a car will violate a red light. A car which is approaching an intersection is connected to road-side units (RSUs) which are installed at traffic lights in an intersection.

Each light is said to be located at points  $B(x_2, y_2, z_2)$ ,  $C(x_3, y_3, z_3)$ ,  $D(x_4, y_4, z_4)$  and  $E(x_5, y_5, z_5)$  respectively.

Unlike image-based systems, this system assumes V2I communication between the traffic lights and the vehicle in question. This means that the traffic state does not need to be determined by an image classifier. Rather, we already have this information by default.

## 2.3 Thao et.al on Traffic Violation Detection

The paper proposed by L. Thao (2022) introduces a mechanism for detecting red light violations automatically. There paper is titled: *Automatic Traffic Red-Light Violation Detection Using AI*

### 2.3.1 Problem Setting

The reason AI technologies (image classification and detection) systems are better suited than standard sensor technologies is that they are able to operate more consistently, even when the number of vehicles in the setting increases dramatically.

### 2.3.2 System Design and Solution Approach

Separate the task into three parts:

1. vehicle violation detection
2. red signal change monitoring
3. vehicle recognition

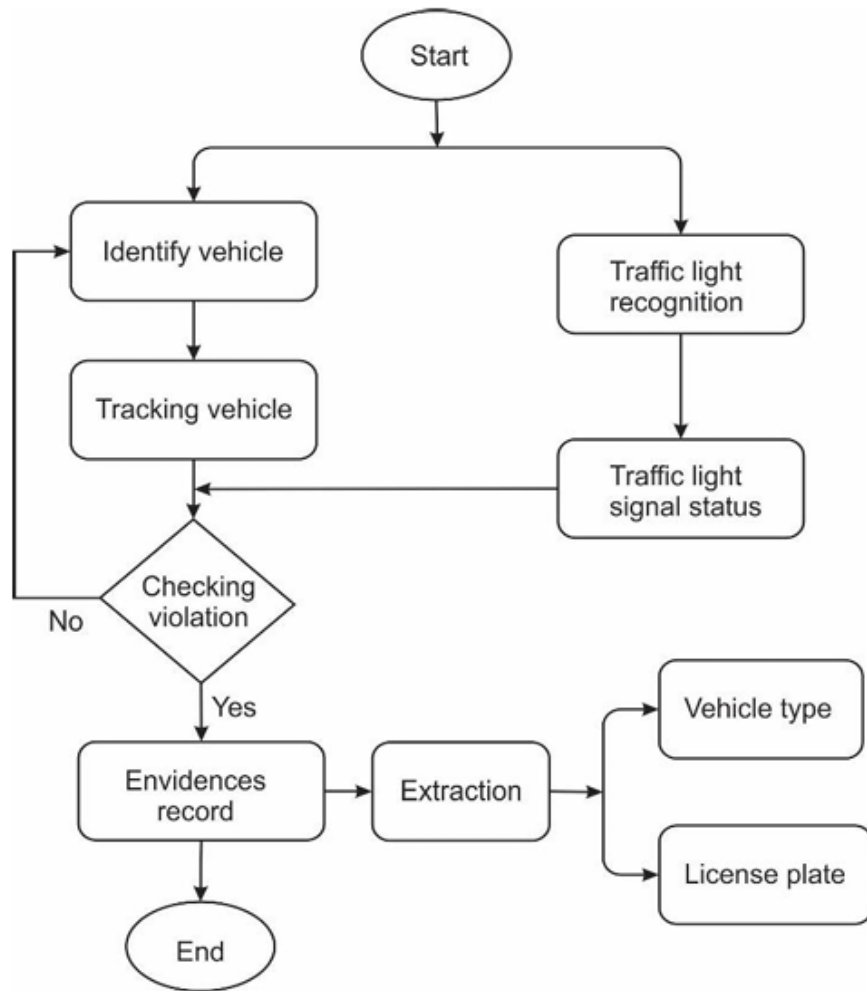
#### Vehicle Violation Detection

The YOLOv5s pretrained model (COCO dataset) is used for detecting violating vehicles. After detecting violation, following frames are used to try and determine the license plate (identify vehicle).

Below, a picture of the overall system flow is presented:

- **vehicle tracking** - performed every 5 frames
  - if IOU (intersection over union) of bounding box is close to one from a previous frame, then the car is assumed to be the same one from that frame.
- **violation line detection**
  - image processing is used to determine traffic lines
  - boundary lines are drawn onto frames captured by the camera later
- **traffic state detection**
  - color filters and image processing used to detect changes in the state of the traffic light



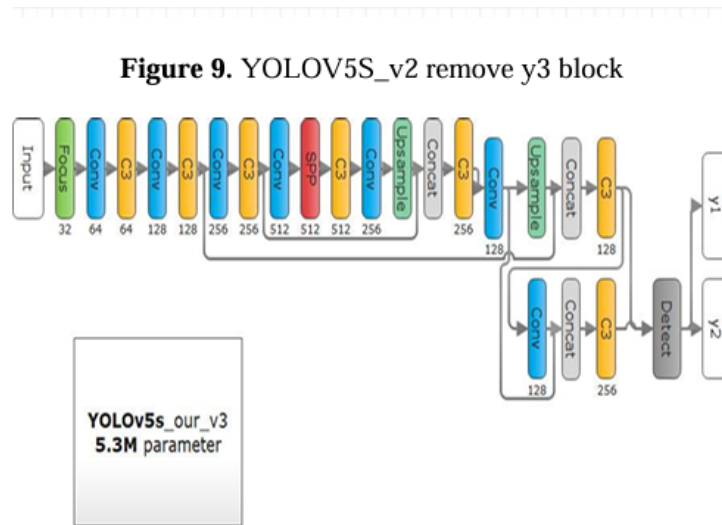


**Figure 2.** Algorithm of violation detection

Figure 2.1: system flow

### 2.3.3 Primary Contributions

1. Implementation of modified *YOLOv5s* model
  1. used parameter changes from original model
  2. achieved following accuracy results:
    1. 82% - vehicle identification
    2. 90% - traffic signal status change
    3. 86% - violation detection
2. Best Performing Architecture given below (v3 / v4)



**Figure 10. YOLOV5S\_v3 reduce ½ filter in Conv, remove y3 block**

Figure 2.2: modified Yolo architecture

## 2.4 Goma et.al on RLR Detection with SSD (Single Shot Detector)

Work by J. Goma (2020) demonstrates the ability to detect red-light-running and over-speeding with a high level of accuracy. Specifically, they achieve **100%** accuracy on red light running detection and **92.1%** accuracy for over-speeding violations. They accomplish these results using a CNN applied to an SSD (single deep neural network)

### 2.4.1 Methods and System Design

## 3 Traffic Flow Forecasting

### 3.1 Attention Based Spatial-Temporal Graph Convolutional Networks for Traffic Flow Forecasting

In this paper, S. Guo (2019) proposed a method for traffic flow prediction which utilized an **attention based spatial-temporal graph convolutional network**. They aimed to model several time-dependencies: (1) recent, (2) daily-periodic, and (3) weekly-periodic dependencies. The *attention mechanism* captures spatial-temporal patterns in the traffic data and the *spatial-temporal convolution* is used to capture spatial patterns while *standard convolutions* describe temporal features.

#### 3.1.1 Core Contributions of ASTGCN

Difficulties of the traffic forecasting problem

1. it is difficult to handle unstable and nonlinear data
2. prediction performance of models require extensive feature engineering
  1. domain expertise is necessary
3. cnn - spatial feature extraction from grid-based data, gcn - describe spatial correlation of grid based data
  1. fails to simultaneously model spatial temporal features and dynamic correlations of traffic data

Addressing these issues:

1. develop a *spatial-temporal attention mechanism*
  1. learns dynamic spatial-temporal correlations of traffic data
  2. temporal attention is applied to capture dynamic temporal correlations for different times
2. Design of *spatial-temporal convolution module*
  1. has graph convolution for modeling graph structure
  2. has convolution in temporal dimension (kind of like 3-d convolution)

## 4 Summary

In summary, this book has no content whatsoever.

## References

- Dokur, O., and S. Katkoori. 2022. “Vehicle-to-Infrastructure Based Algorithms for Traffic Light Detection, Red Light Violation, and Wrong-Way Entry Applications.” *IEEE International Symposium on Smart Electronic Systems*.
- J. Goma, R. Bautista, M. Eviota. 2020. “Detecting Red-Light Runners (RLR) and Speeding Violation Through Video Capture.” *IEEE 7th International Conference on Industrial Engineering and Applications*.
- J. Terven, D. Cordova-Esparaza. 2023. “A Comprehensive Review of YOLO: From YOLOv1 to YOLOv8 and Beyond.” *ACM Computing Surveys*. <https://arxiv.org/pdf/2304.00501v1.pdf>.
- L. Thao, N. Anh, D. Cuong. 2022. “Automatic Traffic Red-Light Violation Detection Using AI.” *International Information and Engineering Technology Association*.
- S. Guo, N. Feng, Y. Lin. 2019. “Attention Based Spatial-Temporal Graph Convolutional Networks for Traffic Flow Forecasting.” *The Thirty-Third AAAI Conference on Artificial Intelligence*.
- W. Liu, D. Erhan, D. Anguelov. 2016. “SSD: Single Shot MultiBox Detector.” *ECCV*.