

# Experiment report -2

Presented by: Amit Hayun, Bar Loupo

Lecturer: Dr. Marina Litvak

## Task definition

Following the conclusion of experiment number-1.

- We will try a different gradient descent optimization algorithm to explore the different performance of the model.[see experiment hyperparameter section]
- In order to thoroughly examine our model performance in a real-time environment, we built a prototype **Abuse Detection System [ADS]**. We bought an IP camera and built a Data Stream pipeline that takes as an input video frame capture by the IP camera, and as an output returns an np frame of size(149,224,224,5) . The output will use as input for the model in order to make a prediction on the frames [see ADS evaluation section]

The **model architecture, workspace environment, data preprocessing method, data augmentation techniques**, and **Data sets** are the same as in experiment number 1.

## Experiment hyperparameter

We chose to try the Adam optimization algorithm, with the command value for the parameters  $\beta_1$   $\beta_2$   $\epsilon$  as shown in **Table 1**

**Table 1:experiment hyperparameter**

<b>learning rate <math>\alpha</math></b>	<b>0.01</b>
<b>epsilon <math>\epsilon</math></b>	<b>1e-07</b>
<b>beta <math>\beta_1</math></b>	<b>0.9</b>
<b>beta <math>\beta_2</math></b>	<b>0.999</b>
<b>batch size</b>	<b>6</b>
<b>Number of workers</b>	<b>6</b>
<b>Number of epoch</b>	<b>30</b>
<b>GPU</b>	<b>1x Titan RTX</b>

## Result

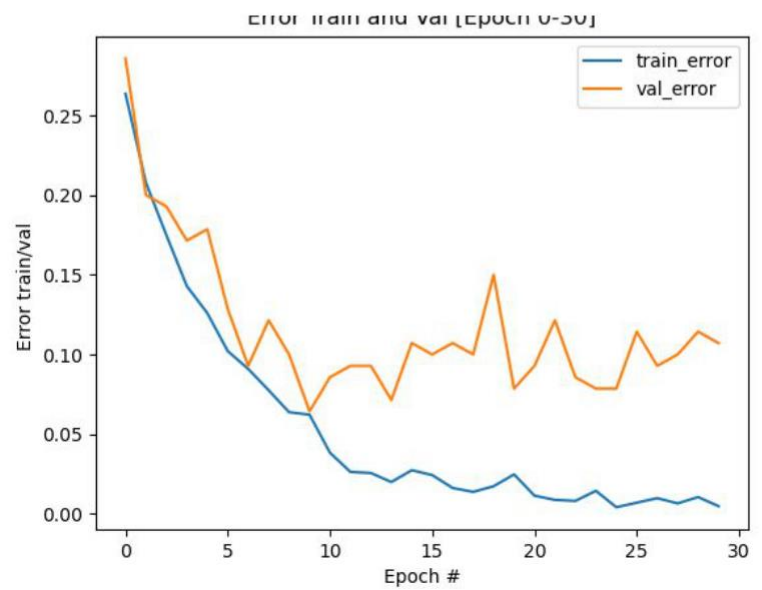
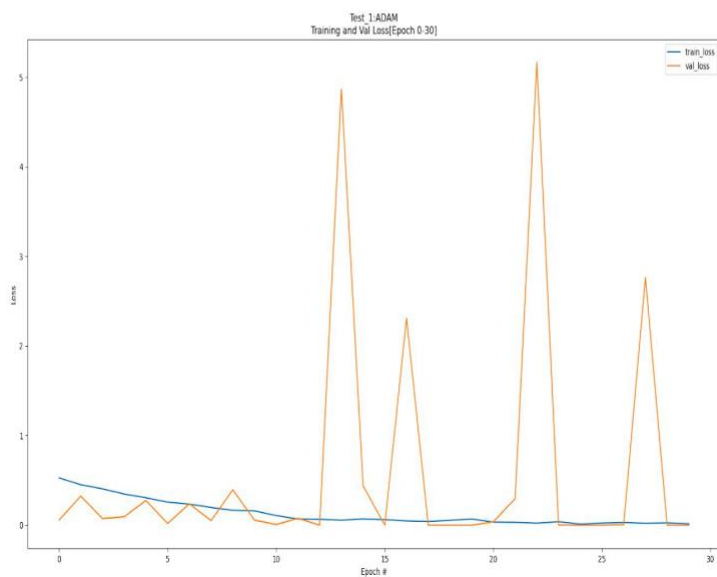
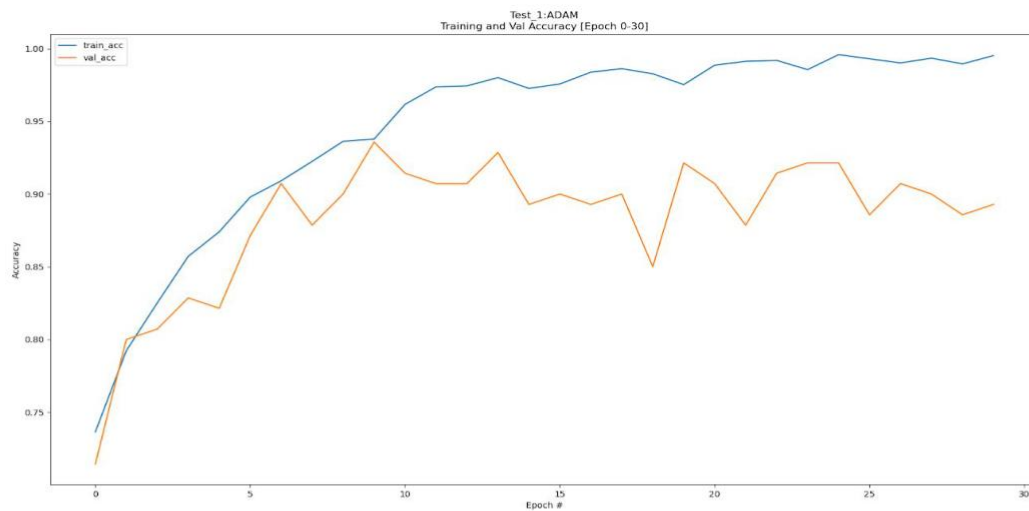
### Training result

we search for the minimum gap between Test error and Val error

$$Test\ Error = 1 - Test\ Accuracy$$

$$Val\ Error = 1 - Val\ Accuracy$$

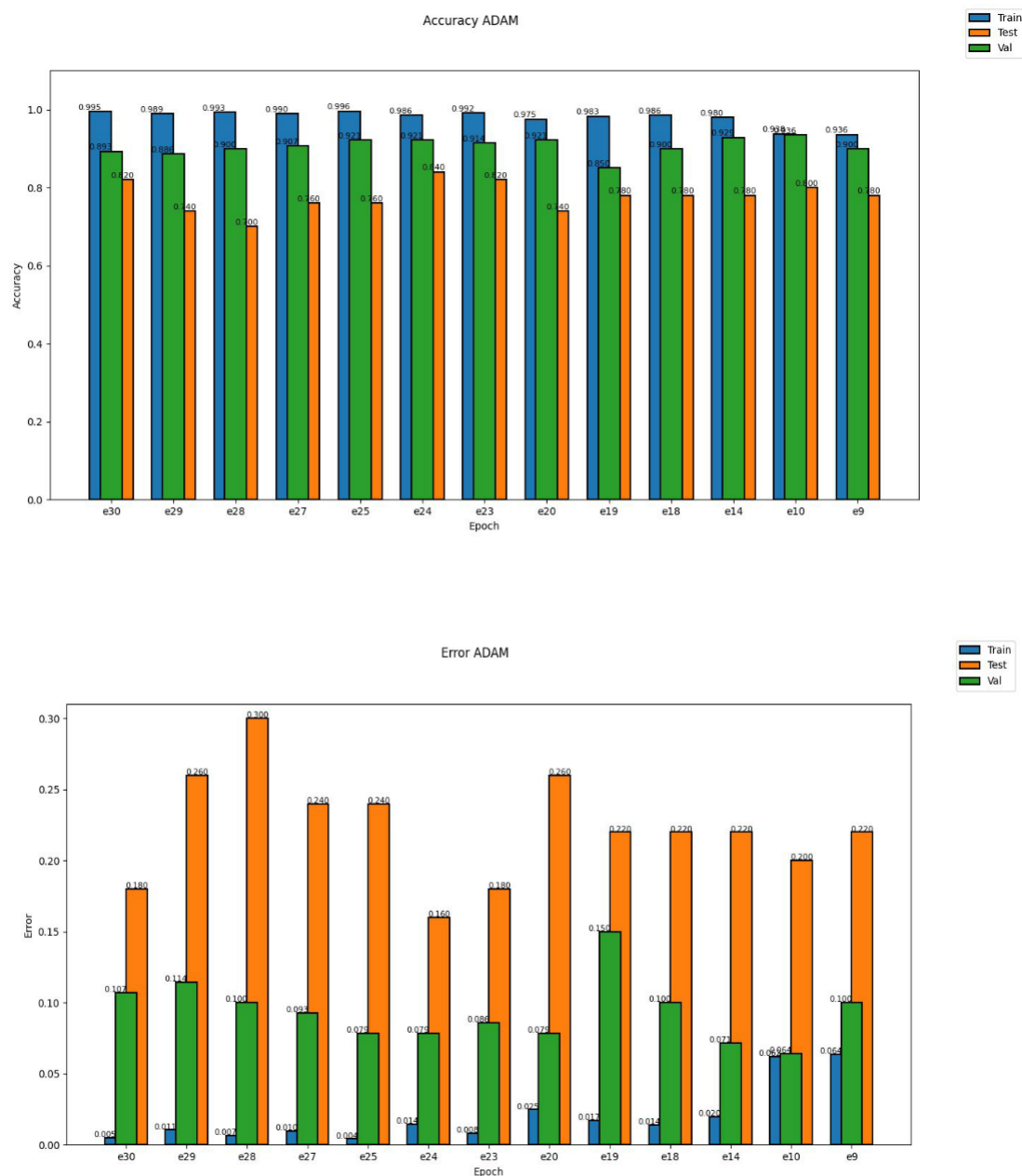
The following images show the results of training and Val accuracy, loss, error during the experiment



We selected the models with the lowest error rate and highest accuracy on the dev set, and found the following models to be the most promising for testing on our test set: models at epoch number [30,29,28,27,25,24,23,20,19,18,14,10,9].

After testing those models we will be analyzing the model performance on our test set using evaluation metric and then choose the best model by this measurement.

the following image shows the Training and Val result accuracy, error, on the chosen model



## Evaluation metrics

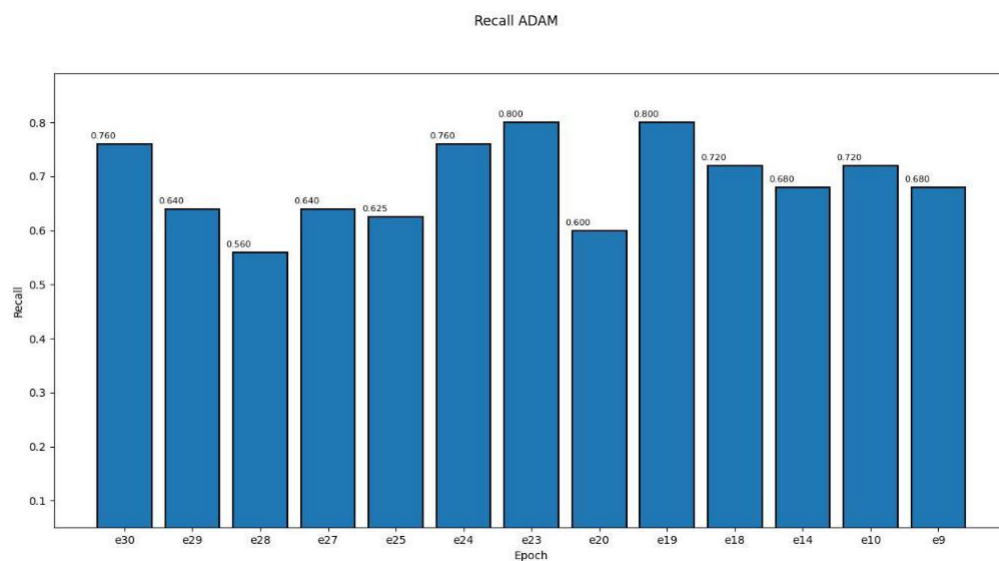
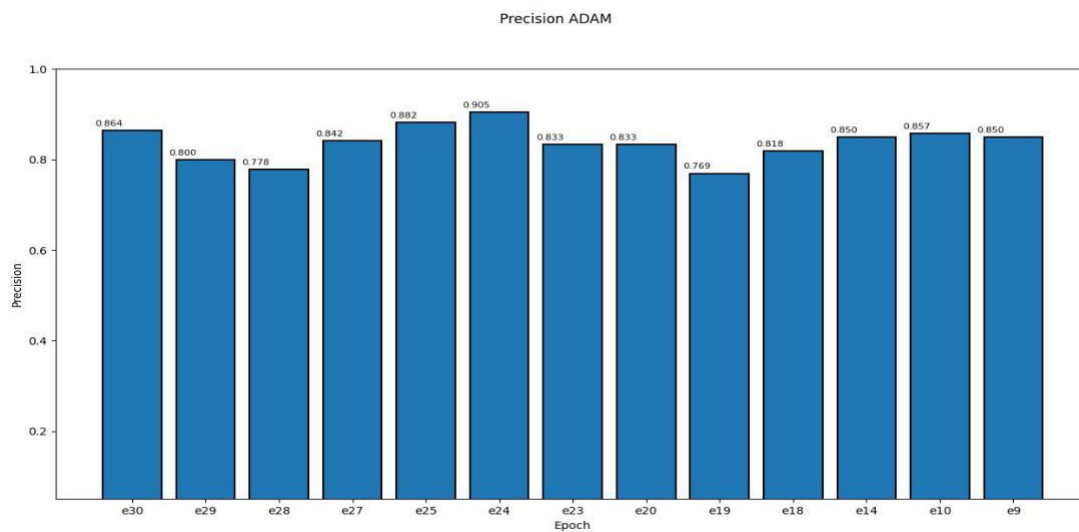
We used the same evaluation techniques as stated in the experiment number 1 report.

$$Recall = \frac{TP}{TP+FN}, P\ precision = \frac{TP}{TP+FP}, Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

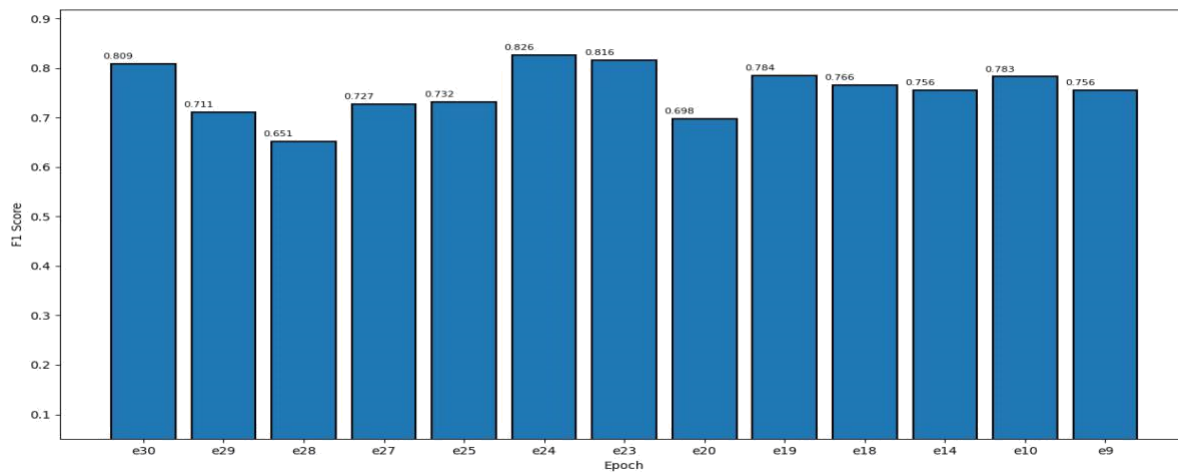
$$F1_{score} = \frac{2 \times Recall \times P\ precision}{Recall + P\ precision}$$

$TP$  = predict Abuse and actual result is Abuse  
 $FN$  = predict Not Abuse and actual result is Abuse  
 $TN$  = predict Not Abuse and actual result is Not Abuse  
 $FP$  = predict Abuse and actual result is Not Abuse

The following charts show the Precision, Recall, F1-score, confusion matrix result by the number of epoch and present all the evaluation metrics in one concluding table[Table -2]



F1 Score ADAM



Epoch e30

TP: 19	FP: 3
FN: 6	TN: 22

Epoch e29

TP: 16	FP: 4
FN: 9	TN: 21

Epoch e28

TP: 14	FP: 4
FN: 11	TN: 21

Epoch e27

TP: 16	FP: 3
FN: 9	TN: 22

Epoch e25

TP: 15	FP: 2
FN: 9	TN: 23

Epoch e24

TP: 19	FP: 2
FN: 6	TN: 23

Epoch e23

TP: 20	FP: 4
FN: 5	TN: 21

Epoch e20

TP: 15	FP: 3
FN: 10	TN: 22

Epoch e19

TP: 20	FP: 6
FN: 5	TN: 19

Epoch e18

TP: 18	FP: 4
FN: 7	TN: 21

Epoch e14

TP: 17	FP: 3
FN: 8	TN: 22

Epoch e10

TP: 18	FP: 3
FN: 7	TN: 22

Epoch e9

TP: 17	FP: 3
FN: 8	TN: 22

we summarize the result as shown in Table 2

**Table 2:result representation**

<b>mode at epoch</b>	<b>Train Accuracy</b>	<b>Val Accuracy</b>	<b>Test Accuracy</b>	<b>Recall</b>	<b>Precision</b>	<b>F1-Score</b>
<u>9</u>	<u>0.936</u>	<u>0.900</u>	<u>0.78</u>	<u>0.68</u>	<u>0.850</u>	<u>0.756</u>
<u>10</u>	<u>0.938</u>	<u>0.936</u>	<u>0.80</u>	<u>0.72</u>	<u>0.857</u>	<u>0.783</u>
<u>14</u>	<u>0.980</u>	<u>0.929</u>	<u>0.78</u>	<u>0.68</u>	<u>0.850</u>	<u>0.756</u>
<u>18</u>	<u>0.986</u>	<u>0.900</u>	<u>0.78</u>	<u>0.72</u>	<u>0.818</u>	<u>0.766</u>
<u>19</u>	<u>0.983</u>	<u>0.850</u>	<u>0.78</u>	<u>0.8</u>	<u>0.769</u>	<u>0.784</u>
<u>20</u>	<u>0.972</u>	<u>0.921</u>	<u>0.74</u>	<u>0.6</u>	<u>0.833</u>	<u>0.698</u>
<u>23</u>	<u>0.992</u>	<u>0.914</u>	<u>0.82</u>	<u>0.8</u>	<u>0.833</u>	<u>0.816</u>
<u>24</u>	<u>0.986</u>	<u>0.921</u>	<u>0.84</u>	<u>0.76</u>	<u>0.905</u>	<u>0.826</u>
<u>25</u>	<u>0.996</u>	<u>0.921</u>	<u>0.76</u>	<u>0.62</u>	<u>0.882</u>	<u>0.732</u>
<u>27</u>	<u>0.990</u>	<u>0.907</u>	<u>0.76</u>	<u>0.64</u>	<u>0.842</u>	<u>0.727</u>
<u>28</u>	<u>0.993</u>	<u>0.900</u>	<u>0.70</u>	<u>0.56</u>	<u>0.778</u>	<u>0.651</u>
<u>29</u>	<u>0.989</u>	<u>0.886</u>	<u>0.74</u>	<u>0.64</u>	<u>0.800</u>	<u>0.711</u>
<u>30</u>	<u>0.995</u>	<u>0.892</u>	<u>0.82</u>	<u>0.76</u>	<u>0.864</u>	<u>0.809</u>

## conclusions

As we can see in Table 3, the best performance was at epoch 24, and the model at epoch 24 got an **accuracy of 84%**.

The confusion matrix at epoch 24:

**Table-3**

TP 19 76%	FP 2 8%
FN 6 24%	TN 23 92%

The Adam optimization algorithm has shown better performance on our test set than the previous SGD optimization algorithm.

We can see an improvement of 5% in F1-Score and an 8% improvement in recall.

mode at epoch	Train Accuracy	Val Accuracy	Test Accuracy	Recall	Precision	F1-Score
<u>SGD</u> epoch 19	<u>0.89</u>	<u>0.829</u>	<u>0.8</u>	<u>0.68</u>	<u>0.89</u>	<u>0.77</u>
<u>ADAM</u> epoch 24	<u>0.986</u>	<u>0.921</u>	<u>0.84</u>	<u>0.76</u>	<u>0.905</u>	<u>0.826</u>

Now, after we found the model with the best result, In order to gain deeper insights on the model performance in a real-time environment, the next step is to evaluate the model on our ADS prototype.

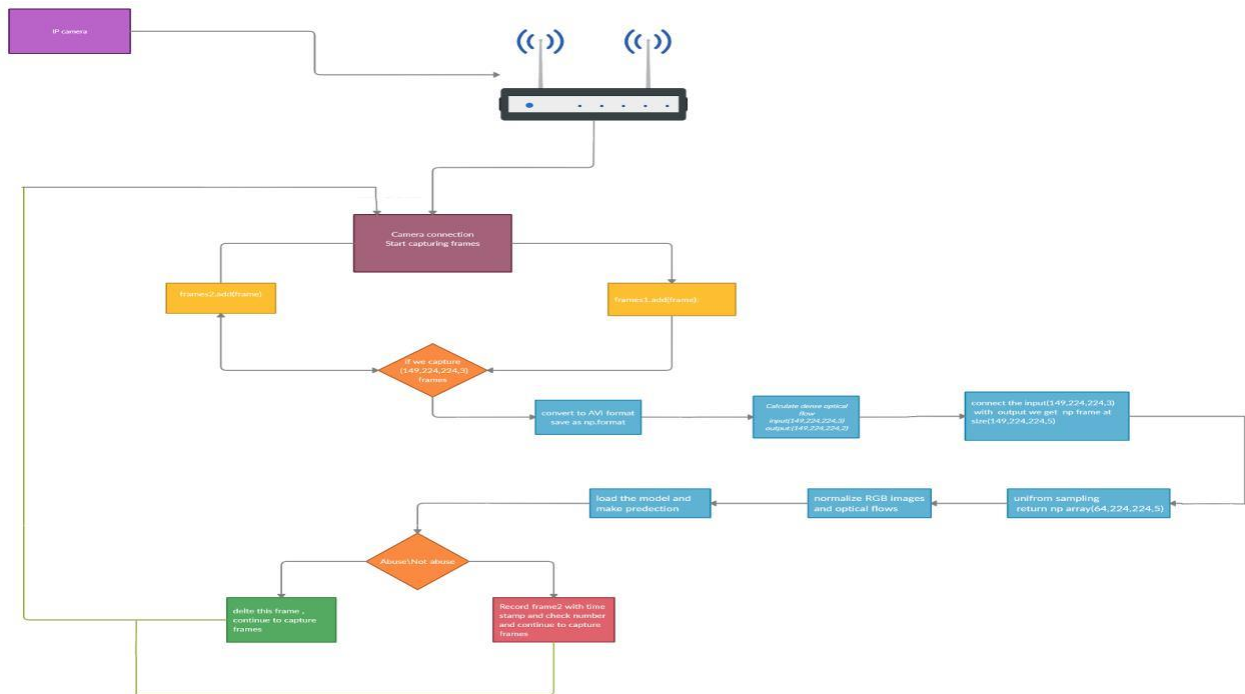


## ADS evaluation

- **ADS structure :**

frame1 = the frame that we will use to make a prediction - Abuse\Not abuse

frame2 = the original frame format captured by the IP camera at the same time.



We built a data pipeline [Blue color] that executes the following every time when capturing 149 frames during broadcasting.

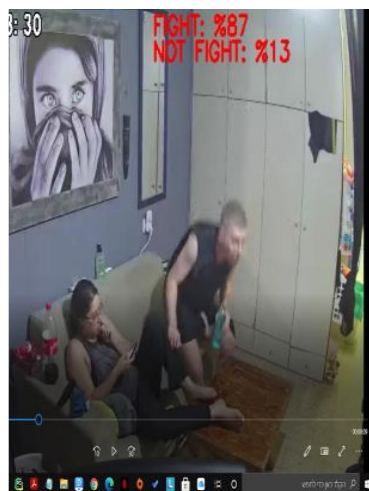
- convert all 149 frames to AVI format
- reshape frame size to (224,224)
- *Calculate dense optical flow input(149,224,224,3) output:(149,224,224,2)*
- connect the input(149,224,224,3) with output to get np frame of size(149,224,224,5)
- *uniform sampling return np array(64,224,224,5)*
- *normalize RGB images and optical flows*
- *load the model and make a prediction of Abuse /Not abuse on the np array(64,224,224,5)*

If the model prediction is Abuse then we will save frame2 as an Abuse event, else we delete frame1 and frame2 and go back to capturing 149 frames.

In order to stimulate an adult abuse environment, we tested the following scenarios

- fast violence moment in ranging 5 seconds each.
- direct violence moment containing punching and kicking
- hidden violence moment when only the back of the attacker is seen
- moments where 2 people are far apart, but 1 person performs a sudden action [such as putting a cup of hot coffee on the table] that will cause a big optical change.

The following pictures show the video that the model classified as Abuse



## **ADS conclusions**

The model will classify correctly violent events when accruing, and it only takes 4 seconds to make a prediction on np frame of size (149,224,224,5)

- The model classifies **video clips as violent** even when one person makes sudden action and there is no proximity between the two people in the video, which actually makes a lot of sense, since the model is looking for changes in the optical spacing between 2 neighboring frames.  
Therefore any fast action, regardless of the distance between two people, will cause a significant change in the optical current.  
This is a very big problem for our system, because most of the nurses/caregivers in a nursing home perform additional activities on a daily basis, for example folding laundry, changing bedding, etc.
- The model has difficulty analyzing frames that contain a TV/computer screen. The sharp movement of the frames on the TV causes a high result in the optical current thus misleading the model and causing it to classify the event as abuse. Moreover, because in most nursing homes / private homes the camera is usually aimed at the living room and/or bedrooms that almost always have a TV, it can be misleading for the model.

## **Possible solution to discuss with academic mentor**

In a broad look at the sub-tasks of computer vision, our project is a sub-task of action recognition. We will therefore look at other sub-tasks in computer vision to find solutions such as **Object Detection, Face Recognition, Visual Relationship Detection ect..**

We thought about using the YOLO algorithm to detect objects. We can identify the position of the TV and by using coordinates we can blacken the area of the TV and thus solve the problem of TV/computer screens. In addition by using the YOLO algorithm, we can identify people in frames, and by using their position coordinates we can calculate the distance between 2 people so that we will run our model only when the distance between 2 people reaches a certain threshold.