

Supplementary Materials for “MTN: Forensic Analysis of MP4 Video Files Using Graph Neural Networks”

Ziyue Xiang[†] Amit Kumar Singh Yadav[†] Paolo Bestagini[‡] Stefano Tubaro[‡] Edward J. Delp[†]

[†]Video and Image Processing Lab (VIPER), School of Electrical and Computer Engineering,
Purdue University, West Lafayette, Indiana, USA

[‡]Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, Milano, Italy

1. The Tree Restructuring Algorithm

The tree restructuring algorithm is shown in Alg. 1. The `TreeRestructure` procedure is used on an MP4 tree \mathcal{T} to generate the restructured tree \mathcal{T}_r .

2. String Processing

The string processing step consists of the following steps:

- Noise removal: non-printable and non-ASCII characters are removed from the input string
- String partition: the input string is broken down into parts separated by white space characters, where each part can be a textual word, URL, version string, ratio string, or numeric string
- URL decomposition: the input strings can also contain URLs which are decomposed into multiple parts; for example, `www.example.com/about/index.html` will be broken down into `www.example.com`, `about`, `index.html`
- Version string decomposition: version strings found in the input string are decomposed into multiple numeric parts; the numeric parts are surrounded by two special words `#VER_BEG` and `#VER_END` to emphasize that they are the members of a version string
- Ratio string conversion: ratio strings such as `1/48000` found in the input string are converted into their corresponding numeric values
- Numeric string conversion: string representations of floating point numbers and integers found in the input string are converted into their corresponding numeric values

3. MTN Model Details

3.1. Graph Pooling (GP) Readout

Denote the output of the Graph Pooling (GP) layer by E' . The first dimension of E' is variable. To produce a fixed-sized vector for classification, a readout mechanism [1] is defined as follows:

$$\text{ReadOut}(E') = \text{MaxP}(E') \parallel \text{MinP}(E'). \quad (1)$$

Here, `MaxP` denotes the max pooling operation that determines the maximum element along each column of E' ,

`MinP` denotes the min pooling operation, and `||` denotes concatenation. The readout operation always produces a $2C$ -dimensional vector, which can be used by a classifier. Here, C is the number of “channels” in the Graph Convolution (GC) layer.

In MTN, we use the hierarchical pooling architecture described in [1]. That is, there are multiple GC+GP layers in the Graph Neural Network (GNN). The readout operation is used for each GP layer in the GNN, and the results are summed together to generate a single vector for classification.

3.2. Model Size of the Node Embedding Network (NEN)

The relationship between the choice of M (embedding size) and the number of parameters in the Node Embedding Network (NEN) are shown in Tab. 1.

Table 1. The choice of M and the corresponding size of the NEN.

M	128	256	384
Num. Parameters	6.56M	15.39M	26.52M

3.3. Model Architecture

We implemented MTN using the `PyTorch` library¹. We used the GC and GP implementation from the `PyG` library².

We used `PyTorch`’s transformer encoder implementation (`TransformerEncoderLayer`). In the number encoder, we chose $\tau_{\max} = 10^{12}$ and $\tau_{\min} = -\tau_{\max}$ (see Sec. 2.2.1). The String Object List (SOL) encoder in the NEN contains 4 transformer encoder layers with $M = 384$ and $h = 4$. The Graph Attention Network (GATN) in the Graph Analysis Module (GAM) contains 6 transformer encoder layers with $M = 384$ and $h = 4$. The choice of M is discussed in Sec. 4.1. For GC, we used the `GCNConv` layer from the `PyG` library with $C = 128$. For GP, we used the `SAGPool` layer from the `PyG` library with $\alpha = 0.5$. The definition of C and α can be

¹<https://pytorch.org/docs/stable/index.html>

²<https://pytorch-geometric.readthedocs.io/en/latest/>

Algorithm 1: The tree restructuring procedure.

```

1 Procedure TreeRestructure( $\mathcal{T}$ ):
   Input: MP4 tree  $\mathcal{T}$ 
   Output: restructured MP4 tree  $\mathcal{T}_r$ 
2   Create a new empty tree  $\mathcal{T}_r$ ;
   // Make sure the new tree  $\mathcal{T}_r$  has the same topology as  $\mathcal{T}$ 
3   foreach node  $n$  in  $\mathcal{T}$  do Create a new node  $n'$  in  $\mathcal{T}_r$  with tag  $n$ ;
4   foreach node  $n$  in  $\mathcal{T}$  do
5     foreach child node  $c$  of  $n$  do
6       Find the corresponding nodes of  $n$  and  $c$  in  $\mathcal{T}_r$ , which are denoted by  $n'$  and  $c'$ ;
7       Connect  $n' \rightarrow c'$  in  $\mathcal{T}_r$ ;
8     end
9   end
   // Add data from  $\mathcal{T}$  to  $\mathcal{T}_r$ 
10  foreach node  $n$  in  $\mathcal{T}$  do
11    Find the corresponding node of  $n$  in  $\mathcal{T}_r$ , which is denoted by  $n'$ ;
12    foreach key-value pair  $(key, val)$  in  $n$  do RecursiveAdd( $(key, val), n'$ );
13  end
14  return  $\mathcal{T}_r$ 
15 Procedure RecursiveAdd( $obj, u'$ ):
   Input: an object to be added ( $obj$ ), which can be a string, a number, a key-value pair, a list, a dictionary; the node
   where the data is added ( $u'$ )
   // Recursively adds data into the restructured tree  $\mathcal{T}_r$ 
16  if  $obj$  is a number or a string then
17    Create a new node  $v'$  in  $\mathcal{T}_r$  with tag  $v$ ;
18    Set the data of  $v'$  as  $obj$ ;
19    Connect  $u' \rightarrow v'$  in  $\mathcal{T}_r$ ;
20  else if  $obj$  is a key-value pair  $(key, val)$  then
21    Create a new node  $k'$  in  $\mathcal{T}_r$  with tag  $k$ ;
22    Set the data of  $k'$  as  $key$ ;
23    Connect  $u' \rightarrow k'$  in  $\mathcal{T}_r$ ;
24    RecursiveAdd( $val, k'$ )
25  else if  $obj$  is a dictionary then
26    Create a new node  $w'$  in  $\mathcal{T}_r$  with tag  $n$ ;
27    Set the data of  $w'$  as #DICT;
28    Connect  $u' \rightarrow w'$  in  $\mathcal{T}_r$ ;
29    foreach  $(key, val)$  in  $obj$  do RecursiveAdd( $(key, val), w'$ );
30  else if  $obj$  is a list then
31    Create a new node  $w'$  in  $\mathcal{T}_r$  with tag  $l$ ;
32    Set the data of  $w'$  as #LIST;
33    Connect  $u' \rightarrow w'$  in  $\mathcal{T}_r$ ;
34  foreach item in  $obj$  do RecursiveAdd( $item, w'$ );

```

found in Sec. 2.3. The MTN has 3 GC+GP layers. Under this setting, the GC+GP layers contain 281k parameters. For classification, we used a three layer multilayer perceptron network with $2C, 2C, K$ neurons. Here, K is the number of classes in the video forensics task.

References

- [1] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention Graph Pooling. *Proceedings of 2019 International Conference*

on Machine Learning, pages 3734–3743, 2019. Long Beach, CA, USA. 1