
Learning Document Graphs with Attention for Image Manipulation Detection

Hailey James

Lendbuzz

Boston, MA 02110

hailey.james@lendbuzz.com

Otkrist Gupta

Lendbuzz

Boston, MA 02110

otkrist.gupta@lendbuzz.com

Dan Raviv

Lendbuzz

Boston, MA 02110

dan.raviv@lendbuzz.com

Abstract

Detecting manipulations in images is becoming increasingly important for combating misinformation and forgery. While recent advances in computer vision have lead to improved methods for detecting spliced images, most state-of-the-art methods fail when applied to images containing mostly text, such as images of documents. We propose a deep-learning method for detecting manipulations in images of documents which leverages the unique structured nature of these images in comparison with those of natural scenes. Specifically, we re-frame the classic image splice detection problem as a node classification problem, in which Optical Character Recognition (OCR) bounding boxes form nodes and edges are added according to a text-specific distance heuristic. We propose a system composed of a Variational Autoencoder (VAE)-based embedding algorithm and a graph neural network with attention, trained end-to-end for robust manipulation detection. Our proposed model outperforms both a state-of-the-art image splice detection method and a document-specific method.

1 Introduction

Identifying spliced images, or images in which content has been added during post-processing, has become an important area of research in computer vision. While spliced images of natural scenes can be used to propagate fake news, many real-world image splicing examples occur in images containing primarily text, or images of documents. Digital documents are commonly used to verify information in domains such as finance and administration, which creates a natural opportunity for forgery and manipulation. With a background of largely white space, images of digital documents pose a problem for most state-of-the-art image manipulation detection models, as they are traditionally trained to detect inconsistencies in texture and background features. In addition, the forged region of a document often looks very similar to the original image, as both often are simply black text on a white background (see Figure 2).

In this paper, we propose an image manipulation detection system designed for improved performance on digital documents. Unlike in natural scene images, forgery detection in digital documents allows us to use Optical Character Recognition (OCR) to narrow the search space of potential manipulations to the textual content of the image. As the number of text boxes on a given page varies and the existence of meaningful relationships between boxes can be inferred, the resulting OCR bounding boxes are natural candidates for learning with Graph Neural Networks (GNNs). Specifically, we can

construct a graph G from a single image of a document in which the OCR bounding boxes form nodes $\mathbf{v} = \{v_1, v_2, v_3, \dots, v_{|V|}\}$ where $|V|$ is the number of detected OCR text boxes on the page. The nodes are connected by edges using a document-specific distance heuristic, such that

$$e_{v_i, v_j} = \begin{cases} 1, & \text{if } d_{v_i, v_j} \leq T \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

where d_{v_i, v_j} is the distance between nodes v_i and v_j , T is a numeric threshold tuned such that connected edges exist within the same paragraph or textual region on a page. We can then re-frame the splice detection in digital documents as a node classification problem, in which text nodes are classified either spliced or genuine.

The first challenge for this approach is initializing the node embeddings to effectively create feature vectors $\mathbf{h} = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_{|V|}\}$ for effective splice manipulation detection. Most methods for learning text-based features from images (such as OCR) are designed to extract high-level features such as character values while ignoring low-level features such as text blur, spacing, or font properties. However, these are the exact properties which are most likely to be relevant for identifying a spliced box, while high-level features such as the character value or word are less likely to be discriminative. In the case of manipulation detection we are interested specifically in ignoring high-level features such as text content while extracting low-level character-based features.

To accomplish this task, we designed a feature extraction algorithm based on pre-training a Variational Autoencoder (VAE) on a set of input and output image pairs. In each image pair, the content varies while the low-level features (character weight, font, spacing, etc) are held constant. This results in a latent space embedding that is sensitive to manipulation artifacts while agnostic to the text content. The encoder from this feature extraction model forms the first part of the manipulation detection model, with weights pre-trained on the specially-designed images pairs.

The second challenge is to design a mechanism which leverages the context of a given block of text for improved classification. In the case of forgery detection, the relationships between text nodes are particularly important, as an anomalous node is most likely to be detected by comparing its features with those of its neighbors. We would expect to see subtle feature differences in spliced text, as it likely would have been spliced from another document or image with different character-level features.

Our solution involves adapting a graph attention layer [1], in which attention weights are trained to update the node embeddings according to the edges between nodes $\mathbf{h} = \{\vec{h}_1, \vec{h}_2, \dots, \vec{h}_{|V|}\} \rightarrow \mathbf{h}' = \{\vec{h}'_1, \vec{h}'_2, \dots, \vec{h}'_{|V|}\}$. This allows our model to appropriately update a node embedding by attending to its neighbors for binary classification (see Figure 1).

We evaluate our method on a public datasets of programmatically-generated document splice dataset. Our proposed model outperforms both a state-of-the-art image splice detection method and a document-specific method, both of which are generally unable to differentiate the spliced text from the genuine text.

2 Related Work

2.1 Image Splice Detection

Image splicing is a manipulation operation in which content is copied and pasted from a source image onto a destination image. Current methods for detecting splices in images can be broadly categorized as methods that search for local inconsistencies in the image background noise [2, 3, 4, 5, 6, 5, 7, 8] or compression discrepancies [9, 10, 11, 12], methods that leverage camera-based artifact inconsistencies [13, 14, 15, 16, 17, 18, 19, 20, 21], or deep-learning based methods [22, 23, 24, 25, 26, 27], with deep-learning based methods generally achieving the best performance. Several public datasets for benchmarking performance on image splice detection in natural scenes have been released [28, 29, 30, 31, 32]. In this paper, we compare our model with Self Consistency [24], a state-of-the-art image splice detection method with a publicly available trained model.

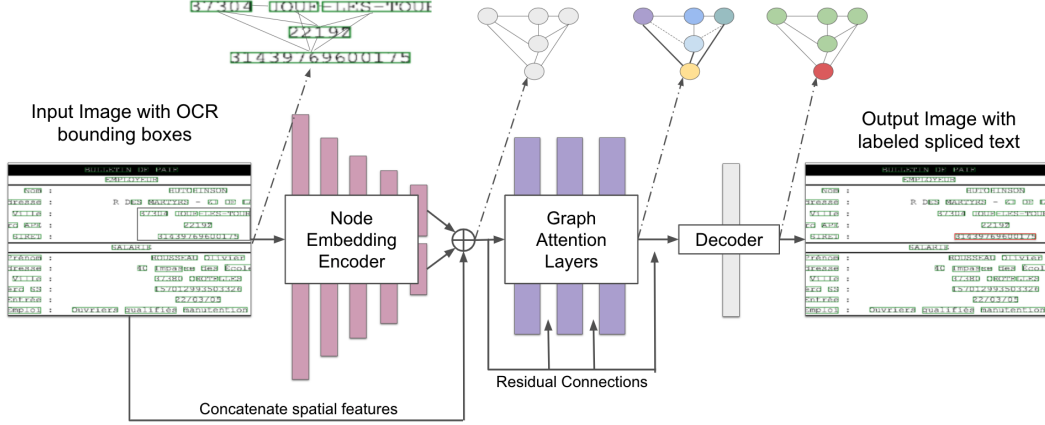


Figure 1: **Architecture overview.** Our proposed system is designed to label text on an image of a document as original to the document or as spliced (copied and pasted from another image). Left to right: First Optical Character Recognition (OCR) bounding boxes are obtained, and a graph structure is formed by treating the text boxes as nodes and adding edges accordingly to a document-based distance heuristic. Each graph is passed through the pre-trained encoder, embedding each node based on low-level features (Section 3). The graph’s spatial features are then concatenated with the encoder output, as represented by \oplus . Next, the node embeddings are updated by a trained graph attention layer, which implicitly weights edges between nodes to encode relational information. Finally, a decoder classifies each node as spliced or genuine, at which point they can be mapped back onto the input image.

a Control number	d Employee's social security number	a Control number	d Employee's social security number
c Employee's name, address, and ZIP code	b Employer identification number	c Employee's name, address, and ZIP code	b Employer identification number
ACME CORPORATION	94-6002123	UNIV OF CALIFORNIA - BERKELEY	94-6002123
BUSINESS SVCS.- PAYROLL #1104		BUSINESS SVCS.- PAYROLL #1104	
BERKELEY, CA 94720-1104		BERKELEY, CA 94720-1104	
e Employee's name, suffix	f Employee's address and ZIP code	e Employee's name, suffix	f Employee's address and ZIP code
BEAR, OSKI		BEAR, OSKI	
951 BEARS ROAD		951 BEARS ROAD	
BERKELEY, CA 94720		BERKELEY, CA 94720	

Figure 2: **Example of a spliced document image** Acme Corporation has replaced Univ of California Berkeley. Unlike in spliced images of natural scenes, there is usually little texture or background information from which to create meaningful features. Instead, we can leverage character-level features to classify blocks of text as original or spliced.

2.2 Document Manipulation Detection

Methods for detecting manipulations in documents often try to leverage document-specific features. Examples include methods that use intrinsic document elements, or portions of the document that can be matched against a template to gauge authenticity [33], font-based features [34, 35, 36], alignment or skew-based features [37]. For images of physical documents, this has included printer identification [38, 39, 40, 41]. [42] present findings from deploying manipulation techniques on real-world datasets.

2.3 Graph Neural Networks

Graph Neural Networks (GNNs), first introduced in 1997 [43], are intended to deal with arbitrarily structured data in order to generalize learning methods to the non-Euclidean domain [44]. Since then, GNNs have been successfully applied to applications ranging from LiDAR device information extraction [45, 46, 47] to scene graph generation [48, 49] to text classification [50, 51, 52]. Typically, Graph Neural Networks perform an iterative process of updating node states until equilibrium is reached. The resulting node states can be learned by consecutive trainable layers in order to classify

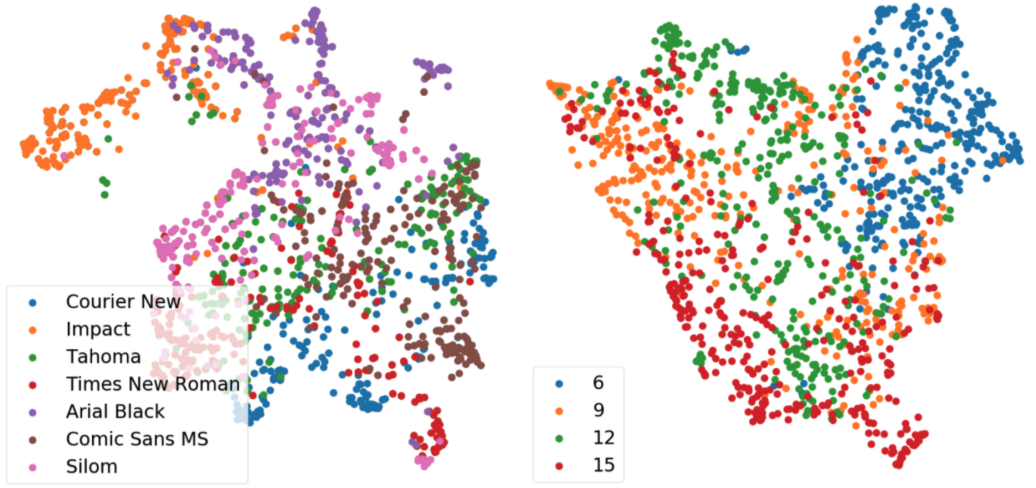


Figure 3: **Encoder embeddings capture low-level font features.** UMAP [53] clusters of randomly-generated text node embeddings labeled by font (left) and font size (right). The VAE is able to separate text iamges based on these low level features. Notably, the embeddings reveal clusters based on character size, despite each text box being scaled to a uniform height and cropped to a uniform width before being passed to the model.

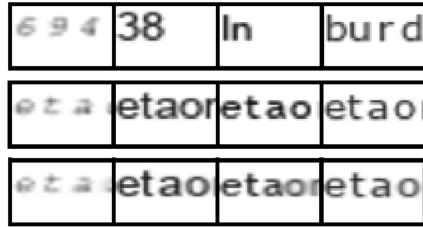


Figure 4: **Data used for pre-training the encoder.** VAE inputs (top), VAE outputs (middle), and VAE reconstructions (bottom). The inputs are constructed by drawing random strings of text (not OCR text images) from OCR results on the training images. The input images are constructed by assigning the random string specific properties, and the outputs are always the same string ("etaonishldcmufpgwybvkjqz") with matching properties to the input. This pairing enables the encoder to capture low-level manipulation artifacts while ignoring content-like character value or sequence patterns.

nodes, graphs, or edges. [1] proposed a Graph Attention Networks (GATs) for learning an attention mechanism on the relationships between nodes.

3 Node Embedding Initialization

One of the primary challenges in learning the document graph for forgery localization is initializing the node embeddings. The node embeddings must have uniform size and contain information that is likely to be useful for manipulation detection. While we can naïvely obtain OCR bounding boxes to represent each node, bounded regions of pixel values will be variable in size and obscure important information features. When text is spliced from one image to another, it is unlikely that features like the character size, font, spacing, and aspect ratio will be identical from the source to the destination. Thus while most learning methods using images of text seek to extract the high-level features such as character value or words (i.e. using OCR), we are interested in an embedding algorithm which can capture these low-level features. One option would be to try to directly learn the fonts, character spacing values, and font sizes and form a corresponding feature vector. However, this would require constructing a pre-defined set of fonts and other low-level features and would be unlikely to generalize well between document sets.

Instead, we train a generative model on image pairs in which the content (character values) change but low-level forgery related features are held constant (see Figures 4 and 3). Using a Variational Autoencoder (VAE) architecture constrains our model to form smooth latent space representations that are conducive to interpolation when unseen data is encountered. In each image pair, the input is formed by sampling a random string from the text in the training images (as obtained through OCR) and creating an image formed from this text and a random character spacing, size, font, and level of Gaussian blur. The output image is formed by taking a fixed string (the same string for every output) and creating an image with the same character spacing, size, font, and blur values as the input. We then take a random crop of 30x50 pixels of the input image and a fixed crop (from the top left corner) of the output image of the same dimensions. While any fixed string could be used as the output string, we chose to use the alphabetical characters ordered by frequency of use ("etaonisrhlcdmufpgwybvjkxqz") in the English language.

Our VAE architecture is composed of four fully-connected encoding layers (sizes 1500, 1204, 908, and 612) with ReLU activation and four fully-connected decoding layers (size 612, 908, 1204, and 1500) with ReLU activation, with a latent dimension of size twenty. We use a final Sigmoid activation function, Xavier initialization, batch size 32, learning rate 0.0001. The loss function is composed of the reconstruction loss and a regularization term, the Kullback-Leibler (KL) divergence between the means and variances encoded in latent space and that of a unit Gaussian. To train our model, we adapt the approximation proposed in [54]

$$\sum_j KL(q_j(z|x) \parallel \mathcal{N}(0, 1)) + \mathbf{L}(x, \hat{x})$$

where \mathbf{L} is the reconstruction error, j is the latent space dimension, KL is the KL divergence, $q_j(z|x)$ is the encoded distribution, and $\mathcal{N}(0, 1)$ is a unit Gaussian. To train our model, we adapt the approximation proposed in [54], or

$$\mathcal{L}(\theta; \hat{x}^{(i)}, x^{(i)}) \simeq \frac{1}{2} \sum_j (1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2) + BCE_{\theta}(\hat{x}^{(i)}, x^{(i)})$$

where \mathcal{L} is the loss given the parameters θ on a text image $x^{(i)}$ from the synthetic dataset, $\mu_j^{(i)}$ and $\sigma_j^{(i)}$ are the mean and variance of the distribution j on image i in the latent space, and $BCE_{\theta}(\hat{x}^{(i)}, x^{(i)})$ is the binary cross entropy loss between the output text image $\hat{x}^{(i)}$ from image pair i and the reconstructed image $\hat{x}^{(i)}$.

To create the node embeddings during training on the spliced dataset, a word box is first scaled to 30 pixels in height and then cropped to 50 pixels in width. It is then passed through the encoder to get the two vectors of size twenty, representing the mean and standard deviations of the distributions in the latent space dimension. These vectors are concatenated with the original bounding box coordinates to incorporate the spatial information of each box for updating the node embeddings while learning the graph structure. This results in final node embeddings of size $20 + 20 + 4 = 44$ to be passed to the first graph attention layer.

4 Graph Learning and Classification

While in rare cases it might be possible to classify a block of text in an image based on its feature embeddings alone, in most cases a block of text is most likely to be classified by comparing it with the nearby text in the document. Because the amount of text varies from page to page, learning these relationships requires either a complex padding method or a network capable of handling input data of various sizes. In addition, text on a page is naturally structured into clusters of paragraphs, tables, headings, footers, etc. These two characteristics make splice detection in images of documents a natural candidates for a graph learning.

We construct the graph by adding edges between OCR bounding boxes given a document-based heuristic (see Equation 1). Because boxes can vary significantly in size, considering a midpoint to midpoint euclidean distance d is not practical for adding edges, as larger boxes would be biased towards having a lower degree than smaller boxes. Instead, we propose an edge addition formula

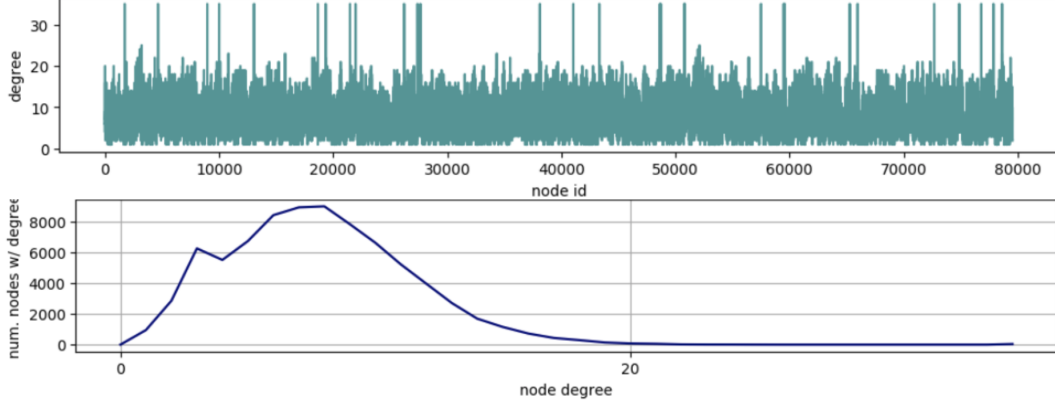


Figure 5: Graph statistics on a sample of 256 images from the Auto-Splice Dataset. Top: Degree by number of nodes with that degree on a sample set. Bottom: Number of nodes by degree. The median degree is 8, which intuitively makes sense as the number of nearby words that might be examined to determine if a block of text has been spliced.

which adds edges based on a function of vertical and horizontal gaps between the bounding boxes and the respective widths and heights of each box. Specifically, edges are added when the vertical gap g_v between the boxes is less than three times the minimum height h_{min} and the horizontal gap g_h is less than four times the minimum height h_{min} (see Equation 2). Additionally, the degree of all nodes is capped at 35. See Figure 5.

$$e_{v_i, v_j} = \begin{cases} 1, & \text{if } g_v \leq 3 \times h_{min} \cap g_h \leq 4 \times h_{min} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

The nature of text on a page means that, for the purpose of splice detection, not all edges can be considered equal. For example, a block of text that differs from the rest of the text on a horizontal line may be more likely to have been spliced than a block of text that differs from the text above or below. In order to learn these relationships and implicitly apply different weights to different edges, we developed a model based on graph attention layers, as proposed in [1] and [55]. The graph attention layer takes as input the set of node features, $\mathbf{h} = \{\vec{h}_1, \vec{h}_2 \dots \vec{h}_{|V|}\}$, $\vec{h}_i \in \mathbb{R}^{F_n}$, where $|V|$ is the number of nodes in the graph and F_n is the node embedding dimension at the n th graph attention layer. In the first graph attention layer, the input is the encoder latent space embedding concatenated with the scaled bounding box coordinates of the text nodes, or $2 \times \text{latent dim} + 4$ (see Section 3). The outputs of a graph attention layer are an updated set of node features, potentially with a new node embedding dimension, or $\mathbf{h}' = \{\vec{h}'_1, \vec{h}'_2 \dots \vec{h}'_N\}$, $\vec{h}'_i \in \mathbb{R}^{F_{n+1}}$.

An attention weight matrix, $\mathbf{W} \in \mathbb{R}^{F_n \times F_{n+1}}$, performs self-attention on each of the nodes by sharing the attention mechanism a for coefficient computation. While the weight matrix is applied over all nodes in the graph, structural information is preserved by masking the attention to include only first-order neighbors of a node i . Attention is normalized across all edges by using the softmax function. In our proposed model, the attention mechanism is a single-layer neural network with Leaky ReLU activation.

$$\alpha_{ij} = \frac{\exp(a(\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_j))}{\sum_{z \in N_x} \exp(a(\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_k))}$$

The normalized attention coefficients serve to form a linear combination of input features, and form the final output node embeddings for the layer. We found that using multi-head attention and *ELU* activation improved training stability, consistent with results in [56] and [1]. The output features can thus be represented by

$$\vec{h}'_i = \parallel_{k=1}^K \text{ELU}(\sum_{j \in N_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j)$$

APPENDIX A

Table A.1: 1980 Census Population by age groups

Age groups	U.S. population		
	Proportion (total)	Proportion (20+ years)	Total
Under 1 year	0.0156		3,533,692
1 - 2 smaller	0.0287		6,493,373
9 - 5 further	0.0419		9,483,880
6 - 11 years	0.0920		20,834,439
12 - 19 years	0.1418		32,113,079
20 - 29 years	0.1803	0.2650	40,839,623
30 - 39 years	0.1392	0.2046	31,526,222
40 - 49 years	0.1005	0. Alaska	22,759,163

Figure 6: An example image from the spliced image dataset with three splices ("smaller", "further" and "Alaska." In order to classify each block of text as spliced or genuine, the model is trained to recognize the subtle variations in font or character width in comparison to its surrounding text.

where K represents the independent attention heads, and \parallel denotes the concatenation of the K outputs, or averaging in the case of the final graph attention layer.

The entire model architecture is composed of four fully-connected encoding layers with 1500, 1204, 908 and 612 nodes respectively (Section 3). The encoder is followed by a three graph attention layers of sizes 36, 29 and 22 and with 4, 4, and 6 attention heads respectively. Residual connection are added between graph attention layer. Finally the updated embeddings are classified by a decoder with a single fully connected layer of size 22. We use simple binary cross-entropy loss for training the model end to end. See Figure 1 for a visualization of the model pipeline.

5 Auto-Splice Dataset

Forgery examples with images of text often include private information, which has resulted in a lack of publicly available datasets for comparing results of manipulation detection models. As a result, we have created and release for public use a dataset of spliced text in images of non-private documents. The original source of the documents is PDF files from the ICDAR 2013 Table Competition [57].¹ The original PDFs contain a wide variety of text types, ranging from paragraphs in articles to numeric tables. This property helps ensure a wide variety of character-level features both within and between documents. We carefully designed an automatic splicing process in order to create splices that are realistic and difficult to identify.

First, PDF-level word boxes are extracted from the PDF document directly. Importantly, the boxes are extracted directly from the PDFs and *not* using OCR, which would otherwise give our model an advantage by being able to bound each potential spliced region with certainty. Each page is then converted to a PNG image to form a set of unmanipulated document images. To calculate the number of splices that will occur in the dataset, we multiply the number of PDF word boxes in the dataset by 0.05, such that in the dataset overall, 5% of the PDF boxes will be spliced. To perform each splice, a source image is chosen at random, from which a source box is drawn. A destination box and destination image are similarly drawn at random, and the source box is scaled and cropped in order to fit the dimensions of the destination box and pasted. This process is repeated until the specified number of splices have been performed, with the constraint that every document retains at least 85% of its original text boxes. Splicing does not occur between instances of the same document, and there is no constraint that enforces boxes to be spliced from a document with different character-level features (font, font size, etc). See Figure 6 for an example from the resulting dataset.

From each of the 200 total PDF pages, we create 50 identical images to serve as an initial images. The required number of splicing operations is performed based on the percentage of splices – we created datasets of both 1% and 5% spliced text boxes, as is reported in Table 1. The resulting dataset has 10,000 unique images of size 1700x2200, complete with ground truth splice bounding boxes.

¹The images in this dataset are images from PDFs from academic works. The PDFs include articles from sociology journals, some of which discuss violent content and may be upsetting to certain readers.

6 Experiments

To train our model, we use Tesseract OCR [58] by way of pytesseract [59] (Apache Software License) to extract word-level bounding boxes. From these bounding boxes, we construct graphs from the page using the distance heuristic described in Equation 2. For computational optimization purposes, we construct a single large disconnected graph for each dataset (training, validation and test). During training and inference, a graph ID mask is used to split the dataset into batches, such that the batch size corresponds to the number of distinct document graphs in the batch, or the number of connected components in the overall dataset graph. Before concatenating the spatial information (see Figure 1) we standardize the data by fitting a scaler on the training data and using it to transform the training, validation, and test data. We created our model using [60] (BSD license) and make our code available for community experimentation. For the graph attention layer, we started from the implementation provided by [61] (MIT License) and made the necessary modifications.

We experimented with a variety of architectures and ranges for the hyperparameters, including an encoder with up to 16 fully-connected layers, up to 5 graph attention layers, up to 4 fully-connected decoding layers, and the presence or absence of residual connections. We note that our architecture exploration was far from exhaustive and that further experimentation in architecture design could lead to improved performance. Based on our limited experimentation facilitated by Tune [62] [63], we decided on the architecture as described above, with 3,689,615 total trainable parameters. Due to computational constraints and the relatively large amount of data, we simply tuned the hyperparameters on the entire validation set rather than performing cross validation. To conserve resources, we stopped experiments if a new maximum validation F1 score had not been reached for 100 consecutive epochs, and spread the experiments over four NVIDIA GeForce 29C RTX P8 GPUs, with each experiment trained on a single GPU. The hyperparameters of the final model are batch size of 64, dropout of 0.4, learning rate 0.0005, and Adam optimization with weight decay 0.0005. The top ten resulting architectures and hyperparameters can be found in the supplementary materials.

6.1 Comparison Models

We compare our proposed model with two state-of-the-art splice detection models, *Intrin. Feats* [64] and *Self Consistency* [24]. [64] is a model especially designed for detecting manipulations in images of documents. Similar to our proposed model, it starts from OCR bounding boxes obtained through Tesseract OCR. Following OCR extraction, a feature vector for each text bounding box is constructed using the character size, principal inertia axis, horizontal alignment, and Hu moments [65]. One of the main contributions of the method is the use of character-level features that are intrinsic to the document, in comparison with other methods which rely on the presence of extrinsic feature watermarks or a master template. [24] is a deep learning based method designed to search for inconsistencies in the low-level pixel information for the purpose of splice detection. It has achieved state-of-the-art performance on several challenging splice datasets of images of natural scenes, such as *Columbia* [30], *Carvalho et al.* [29] and *Realistic Tampering* [66]. We made use of the publicly available code and pre-trained model for evaluation released by the authors to perform the comparison on the spliced document image datasets.

6.2 Results

While object detection methods often report a mean Average Precision (mAP) score, this would primarily measure the quality of the OCR engine rather than splice detection model performance. This is because the inputs to our model and [64] are OCR bounding boxes. Instead, to evaluate our proposed model’s ability to detect splices in digital documents, we measure both the F1 score with respect to labeling boxes as manipulated or genuine. Because our model takes inputs already partitioned by OCR bounding boxes, it is straightforward to calculate the F1 score given the predicted labels and ground truth labels.

In order to calculate an F1 score for comparison with [24], we calculate a score for each OCR bounding box, which is simply the mean of all the pixels in the model output within a given bounding box. Each box is then assigned a final label, which is determined by a threshold score, and an F1 score is calculated. To determine the threshold, we simply iterated through a list of linearly spaced potential threshold values and selected the one which gave the maximum F1 score.

Table 1: The results of our model compared with two other splice detection models on our dataset of digital document splices. The state-of-the-art image manipulation detection models are in effect unable to identify splices in documents, despite being far more computationally expensive than our proposed model. Each model was run on a single NVIDIA GeForce 29C RTX P8 GPU, which is reflected in the mean seconds per image.

	Auto-Splice (5% Manip)	Auto-Splice (1% Manip.)	Auto-Splice
	Test F1		Mean Seconds per Image
Self Cons.	0.0171	0.0013	104.40 s
Intrin. Doc.	0.0559	0.0348	17.64 s
Proposed Model	0.9038	0.6528	8.01 s

Our model outperforms both of the state-of-the-art splice detection models, despite being far less computationally expensive (See Table 1). For inference on a sample image of the same size and amount of text as in the training data, our proposed model takes 2.8500 seconds to extract the OCR boxes with Tesseract, 5.1470 seconds to construct the document graph, and 0.0087 seconds to pass the input through the trained model.

7 Conclusion

We present a new method and dataset for the purpose of detecting splice manipulations in images of documents. In contrast with other splice manipulation detection methods, our method leverages the structured nature of text in document images, re-framing the problem as a node classification task. To initialize useful embeddings for manipulation detection, we train a Variational Auto-encoder (VAE) on text-image pairs which differ in content but are consistent in character-level features. Next, we construct a graph from the OCR bounding boxes in an image, adding edges according to a custom document-based distance heuristic. Finally we combine the node embedding encoder with a graph attention layer and a final decoding layer for robust splice detection. Our model compares favorably to both comparison models.

Future work could consider new graph construction methods, including methods which are trainable alongside the rest of the network. One example could include adding edges between blocks of text that are the same throughout the document. Additionally, we anticipate significant improvement could come from employing a more robust OCR method to improve inputs to our system.

This work is limited to splicing manipulations in which text is spliced between two images containing different text. While we anticipate this approach could be adapted for additional manipulation types (e.g. copy-move), it is not intended to generalize to splicing operations which cannot be detected by OCR (such as non-text splicing) or to splicing between images with identical text.

Methods for detecting manipulations in documents can reduce the spread of misinformation and prevent financial fraud and scams. Developing automated methods can help reduce the burden on human reviewers who are tasked with authenticating digital documents. As with any machine learning model, caution should be exercised during planning and deployment. In particular, all stakeholders should be involved in determining acceptable false positive rates, designing recourse policies, and ensuring adequate model monitoring.

References

- [1] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- [2] Xin Wang, Bo Xuan, and Si-long Peng. Digital image forgery detection based on the consistency of defocus blur. In *2008 International Conference on Intelligent Information Hiding and Multimedia Signal Processing*, pages 192–195. IEEE, 2008.
- [3] Bo Liu and Chi-Man Pun. Splicing forgery exposure in digital image by detecting noise discrepancies. *International Journal of Computer and Communication Engineering*, 4(1):33, 2015.

- [4] Bo Liu, Chi-Man Pun, and Xiao-Chen Yuan. Digital image forgery detection using jpeg features and local noise discrepancies. *The Scientific World Journal*, 2014, 2014.
- [5] Chi-Man Pun, Bo Liu, and Xiao-Chen Yuan. Multi-scale noise estimation for image splicing forgery detection. *Journal of visual communication and image representation*, 38:195–206, 2016.
- [6] Thibaut Julliard, Vincent Nozick, and Hugues Talbot. Image noise and digital image forensics. In *International Workshop on Digital Watermarking*, pages 3–17. Springer, 2015.
- [7] Miroslav Goljan, Jessica Fridrich, and Rémi Cogramne. Rich model for steganalysis of color images. In *2014 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 185–190. IEEE, 2014.
- [8] Babak Mahdian and Stanislav Saic. Using noise inconsistencies for blind image forensics. *Image and Vision Computing*, 27(10):1497–1503, 2009.
- [9] Shuiming Ye, Qibin Sun, and Ee-Chien Chang. Detecting digital image forgeries by measuring inconsistencies of blocking artifact. In *2007 IEEE International Conference on Multimedia and Expo*, pages 12–15. Ieee, 2007.
- [10] Yuting Su, Jing Zhang, and Jie Liu. Exposing digital video forgery by detecting motion-compensated edge artifact. In *2009 International Conference on Computational Intelligence and Software Engineering*, pages 1–4. IEEE, 2009.
- [11] Wu-Chih Hu and Wei-Hao Chen. Effective forgery detection using dct+ svd-based watermarking for region of interest in key frames of vision-based surveillance. *International Journal of Computational Science and Engineering*, 8(4):297–305, 2013.
- [12] Ashima Gupta, Nisheeth Saxena, and SK Vasistha. Detecting copy move forgery using dct. *International Journal of Scientific and Research Publications*, 3(5):1, 2013.
- [13] Pasquale Ferrara, Tiziano Bianchi, Alessia De Rosa, and Alessandro Piva. Image forgery localization via fine-grained analysis of cfa artifacts. *IEEE Transactions on Information Forensics and Security*, 7(5):1566–1577, 2012.
- [14] Harpreet Kaur, Jyoti Saxena, and Sukhjinder Singh. Simulative comparison of copy-move forgery detection methods for digital images. *International Journal of Electronics, Electrical and Computational System*, 4:62–66, 2015.
- [15] Bo Liu and Chi Man Pun. Hsv based image forgery detection for copy-move attack. In *Applied Mechanics and Materials*, volume 556, pages 2825–2828. Trans Tech Publ, 2014.
- [16] Feng Zeng, Wei Wang, Min Tang, and Zhanghua Cao. Exposing blurred image forgeries through blind image restoration. In *2015 10th International Conference on P2P, Parallel, Grid, Cloud and Internet Computing (3PGCIC)*, pages 466–469, 2015.
- [17] Wu-Chih Hu, Wei-Hao Chen, Deng-Yuan Huang, and Ching-Yu Yang. Effective image forgery detection of tampered foreground or background image based on image watermarking and alpha mattes. *Multimedia Tools and Applications*, 75(6):3495–3516, 2016.
- [18] Aniket Roy, Rahul Dixit, Ruchira Naskar, and Rajat Subhra Chakraborty. Copy-move forgery detection with similar but genuine objects. In *Digital Image Forensics*, pages 65–77. Springer, 2020.
- [19] Irene Amerini, Lamberto Ballan, Roberto Caldelli, Alberto Del Bimbo, Luca Del Tongo, and Giuseppe Serra. Copy-move forgery detection and localization by means of robust clustering with j-linkage. *Signal Processing: Image Communication*, 28(6):659–669, 2013.
- [20] Longyin Wen, Honggang Qi, and Siwei Lyu. Contrast enhancement estimation for digital image forensics. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 14(2):1–21, 2018.

- [21] I-Cheng Chang, J Cloud Yu, and Chih-Chuan Chang. A forgery detection algorithm for exemplar-based inpainting images using multi-region relation. *Image and Vision Computing*, 31(1):57–71, 2013.
- [22] Jawadul H Bappy, Amit K Roy-Chowdhury, Jason Bunk, Lakshmanan Nataraj, and BS Manjunath. Exploiting spatial structure for localizing manipulated image regions. In *Proceedings of the IEEE international conference on computer vision*, pages 4970–4979, 2017.
- [23] Belhassen Bayar and Matthew C Stamm. A deep learning approach to universal image manipulation detection using a new convolutional layer. In *Proceedings of the 4th ACM Workshop on Information Hiding and Multimedia Security*, pages 5–10, 2016.
- [24] Minyoung Huh, Andrew Liu, Andrew Owens, and Alexei A Efros. Fighting fake news: Image splice detection via learned self-consistency. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 101–117, 2018.
- [25] Ronald Salloum, Yuzhuo Ren, and C-C Jay Kuo. Image splicing localization using a multi-task fully convolutional network (mfcn). *Journal of Visual Communication and Image Representation*, 51:201–209, 2018.
- [26] Peng Zhou, Xintong Han, Vlad I Morariu, and Larry S Davis. Learning rich features for image manipulation detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1053–1061, 2018.
- [27] Yue Wu, Wael AbdAlmageed, and Premkumar Natarajan. Mantra-net: Manipulation tracing network for detection and localization of image forgeries with anomalous features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9543–9552, 2019.
- [28] Jing Dong, Wei Wang, and Tieniu Tan. Casia image tampering detection evaluation database. In *2013 IEEE China Summit and International Conference on Signal and Information Processing*, pages 422–426. IEEE, 2013.
- [29] Tiago José De Carvalho, Christian Riess, Elli Angelopoulou, Helio Pedrini, and Anderson de Rezende Rocha. Exposing digital image forgeries by illumination color classification. *IEEE Transactions on Information Forensics and Security*, 8(7):1182–1194, 2013.
- [30] Tian-Tsong Ng, Shih-Fu Chang, and Q Sun. A data set of authentic and spliced image blocks. *Columbia University, ADVENT Technical Report*, pages 203–2004, 2004.
- [31] Paweł Korus and Jiwu Huang. Multi-scale analysis strategies in prnu-based tampering localization. *IEEE Transactions on Information Forensics and Security*, 12(4):809–824, 2016.
- [32] Vladimir V Kniaz, Vladimir Knyaz, and Fabio Remondino. The point where reality meets fantasy: Mixed adversarial generators for image splice detection. 2019.
- [33] Amr Gamal Hamed Ahmed and Faisal Shafait. Forgery detection based on intrinsic document contents. In *2014 11th IAPR International Workshop on Document Analysis Systems*, pages 252–256. IEEE, 2014.
- [34] Romain Bertrand, Oriol Ramos Terrades, Petra Gomez-Kramer, Patrick Franco, and Jean-Marc Ogier. A conditional random field model for font forgery detection. In *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, pages 576–580. IEEE, 2015.
- [35] Francisco Cruz, Nicolas Sidere, Mickael Coustaty, Vincent Poulain D’Andecy, and Jean-Marc Ogier. Local binary patterns for document forgery detection. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 1, pages 1223–1228. IEEE, 2017.
- [36] Yulia S Chernyshova, Mikhail A Aliev, Ekaterina S Gushchanskaia, and Alexander V Sheshkus. Optical font recognition in smartphone-captured images and its applicability for id forgery detection. In *Eleventh International Conference on Machine Vision (ICMV 2018)*, volume 11041, page 110411J. International Society for Optics and Photonics, 2019.

- [37] Joost van Beusekom, Faisal Shafait, and Thomas M. Breuel. Text-line examination for document forgery detection. *International Journal on Document Analysis and Recognition (IJDAR)*, 16(2):189–207, June 2013.
- [38] Joost van Beusekom, Faisal Shafait, and Thomas M Breuel. Automatic authentication of color laser print-outs using machine identification codes. *Pattern Analysis and Applications*, 16(4):663–678, 2013.
- [39] Shize Shang, Nasir Memon, and Xiangwei Kong. Detecting documents forged by printing and copying. *EURASIP Journal on Advances in Signal Processing*, 2014(1):140, 2014.
- [40] Shize Shang, Xiangwei Kong, and X. You. Document forgery detection using distortion mutation of geometric parameters in characters. *J. Electronic Imaging*, 2015.
- [41] Surbhi Gupta and Munish Kumar. Forensic document examination system using boosting and bagging methodologies. *Soft Computing*, 24(7):5409–5426, April 2020.
- [42] Joost Van Beusekom, Armin Stahl, and Faisal Shafait. Lessons learned from automatic forgery detection in over 100,000 invoices. In *Computational Forensics*, pages 130–142. Springer, 2012.
- [43] Alessandro Sperduti and Antonina Starita. Supervised neural networks for the classification of structures. *IEEE Transactions on Neural Networks*, 8(3):714–735, 1997.
- [44] Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.
- [45] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019.
- [46] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4558–4567, 2018.
- [47] Gusi Te, Wei Hu, Amin Zheng, and Zongming Guo. Rgcnn: Regularized graph cnn for point cloud segmentation. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 746–754, 2018.
- [48] Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. Graph r-cnn for scene graph generation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 670–685, 2018.
- [49] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5410–5419, 2017.
- [50] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- [51] William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. *arXiv preprint arXiv:1706.02216*, 2017.
- [52] Diego Marcheggiani and Ivan Titov. Encoding sentences with graph convolutional networks for semantic role labeling. *arXiv preprint arXiv:1703.04826*, 2017.
- [53] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction, 2018. cite arxiv:1802.03426Comment: Reference implementation available at <http://github.com/lmcinnes/umap>.
- [54] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.

- [55] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2015.
- [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [57] Tami Hassan. Icdar 2013 table competition dataset, 2013.
- [58] Ray Smith. An overview of the tesseract ocr engine. In *Ninth international conference on document analysis and recognition (ICDAR 2007)*, volume 2, pages 629–633. IEEE, 2007.
- [59] Matthias Lee. pytesseract, Dec 2020.
- [60] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [61] Aleksa Gordić. pytorch-gat. <https://github.com/gordicaleksa/pytorch-GAT>, 2020.
- [62] Richard Liaw, Eric Liang, Robert Nishihara, Philipp Moritz, Joseph E Gonzalez, and Ion Stoica. Tune: A research platform for distributed model selection and training. *arXiv preprint arXiv:1807.05118*, 2018.
- [63] Liam Li, Kevin Jamieson, Afshin Rostamizadeh, Ekaterina Gonina, Moritz Hardt, Ben Recht, and Ameet Talwalkar. Massively parallel hyperparameter tuning, 2019.
- [64] Romain Bertrand, Petra Gomez-Kramer, Oriol Ramos Terrades, Patrick Franco, and Jean-Marc Ogier. A system based on intrinsic features for fraudulent document detection. In *2013 12th International conference on document analysis and recognition*, pages 106–110. IEEE, 2013.
- [65] Ming-Kuei Hu. Visual pattern recognition by moment invariants. *IRE transactions on information theory*, 8(2):179–187, 1962.
- [66] Paweł Korus and Jiwu Huang. Evaluation of random field models in multi-modal unsupervised tampering localization. In *2016 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6, 2016.