# MEG: Multi-Evidence GNN for Multimodal Semantic Forensics

Ekraam Sabir, Ayush Jaiswal, Wael AbdAlmageed and Prem Natarajan

USC Information Sciences Institute

Marina Del Rey

California, USA

Email: {esabir, ajaiswal, wamageed, pnataraj}@isi.edu

*Abstract*—Fake news often involves semantic manipulations across modalities such as image, text, location etc and requires the development of multimodal semantic forensics for its detection. Recent research has centered the problem around images, calling it image repurposing – where a digitally unmanipulated image is semantically misrepresented by means of its accompanying multimodal metadata such as captions, location, etc. The image and metadata together comprise a multimedia package. The problem setup requires algorithms to perform multimodal semantic forensics to authenticate a query multimedia package using a reference dataset of potentially related packages as evidences. Existing methods are limited to using a single evidence (retrieved package), which ignores potential performance improvement from the use of multiple evidences. In this work, we introduce a novel graph neural network based model for multimodal semantic forensics, which effectively utilizes multiple retrieved packages as evidences and is scalable with the number of evidences. We compare the scalability and performance of our model against existing methods. Experimental results show that the proposed model outperforms existing state-of-the-art algorithms with an error reduction of up to $25\%$.

## I. Introduction

Credibility of news on social media has been very low recently because of fake news [1]. The severity of the problem is indicated by its ability to affect elections, manipulate public opinion and in general, spread potentially malicious misinformation [2][3]. Recognizing this problem, social media platforms, such as Twitter, have conducted studies and invested in research for understanding this phenomena [4]. The social or moral obligation to contain falsehood notwithstanding the economic and political ramifications of fake news increase the importance of developing methods to detect fake news. The term *fake news* colloquially refers to factually incorrect information. Technically speaking, *fake news* and alternative terms such as *hoax* and *rumor* usually refer to manipulated multimedia (e.g. text, images, video, etc.) used for spreading incorrect information. That includes image-centric manipulations where images are digitally edited and/or associated with altered metadata.

Image repurposing is a relatively new research problem that deals with multimodal semantic forensics. Jaiswal et al. [5] define image repurposing as the problem of semantically misrepresenting an image with falsified multimodal metadata such as captions, location, etc. The image and accompanying metadata together comprise a *multimedia package*. Figure



Fig. 1. An example of multimodal semantic manipulation or image repurposing, in a multimedia package. The image was originally taken in Japan and has been used to falsify location semantics.

1 shows an example of image repurposing presented as a multimedia package. The image was originally taken at a replica of the Statue of Liberty located at Odaiba Island in Japan, but has been repurposed to misrepresent its location. In the example, the supporting multimodal information is text and global positioning system (GPS) coordinates. The problem setup involves detecting whether a multimedia package has been semantically repurposed with the help of an external reference dataset, which is a knowledge base of packages that are assumed to contain unmanipulated information. In this setup, each package in the reference dataset is a potential evidence for verifying a query package. Sabir *et al.* [6] introduced the multimodal entity image repurposing (MEIR) dataset with challenging manipulations and presented a deep multimodal model (DMM) which achieved state-of-the-art performance on MEIR. However, a shortcoming of DMM is that it is capable of utilizing only one evidence from the reference dataset for repurposing detection. DMM does not scale to handle multiple retrieved packages and hence, does not leverage additional information for performance improvement.

In this paper, we propose MEG – a multi-evidence graph

neural network (GNN) model for multimodal semantic forensics. The GNN in our model inherently provides invariance to the order of packages retrieved from the reference dataset. This paper has the following contributions:

- A new GNN-based model for multimodal semantic forensics that achieves state-of-the-art performance
- A scalable model for assimilating arbitrary number of evidences for semantic repurposing detection

The remainder of this paper is organized as follows: Section II discusses related work. Section III describes the proposed model. Results of our evaluation are discussed in Section IV. Finally, Section V concludes the paper and provides directions for future work.

## II. RELATED WORK

While image repurposing itself is a type of fake news, most fake news forensics do not involve images or multimodal semantics. The problem is instead tackled solely from a natural language processing (NLP) perspective. There are two popular ways to tackle the problem: (1) classification of flat feature vectors [7][8] and (2) as an epidemic on a social graph [9][10]. Feature vector based approaches capture information about text content such as presence of swear words, urls or information about users such as number of followers. In [7] and [8] tweets are classified according to tweet content and metadata information for real time analysis. Jin *et al.*[10] investigate the use of a diffusion model to classify rumors. Wu *et al.*[9] model message propagation as a tree and use graph kernels to measure similarities for classifying tweets.

Digital image manipulation is a component of fake news, considering it can be used to misrepresent information. The connection between fake news and forged images is supported in [11], where Zampoglou *et al.* identify the need for checking images for journalism. They develop a web-based graphical user interface (GUI) for verifying images using existing algorithms. Digital image manipulation detection has been studied extensively [12][13][14]. Pixel level image manipulations fall into copy-move, image-splicing, resampling and retouching categories [14].

Image repurposing represents a component of fake news where images are not digitally manipulated, rather semantically misrepresented. It brings multimodal semantics into focus which has not been accounted for in previous approaches. Jaiswal *et al.*[5] introduced the problem of image repurposing. However the manipulations introduced in their dataset were unlikely to fool people and their reference dataset did not have directly linked evidences to query packages. They proposed a joint embedding model for images and text followed by outlier detection. This approach was unsupervised and worked on manipulations requiring common world knowledge. Sabir *et al.*[6] introduced the multimodal entity image repurposing (MEIR) dataset with more realistic manipulations. Their deep multitask (DMM) model retrieves one package from the reference dataset as evidence and uses it for repurposing/semantic manipulation detection. However, their model does not scale to multiple evidences because their multitask framework covers

only one retrieved package. Jaiswal *et al.* [15] introduced an adversarial model for image repurposing detection, where an active counterfieter helps train the detection network. Recently Budack *et al.* [16] attempted to detect manipulations with a web-based application that verifies cross-modal information.

Processing of sequential inputs have been well studied with recurrent neural networks [17]. However, sequential networks may not be optimal when processing unordered inputs as shown in [18]. Vinyals *et al.* [18] introduced a read-process-write (RPW) network for order invariant processing of set inputs. Their working mechanism is to learn a content-based representation that is invariant to input order. The read network involves learning representations for each input sample, the process network encodes all read network outputs into an order agnostic representation and the write network decodes the process network output. Application of order invariant methods is catching up, such as for object detection and multiclass image classification in [19].

Graph Neural Networks (GNNs) were introduced by Scarselli *et al.*[20]. A GNN receives graph structured data as input and processes it with neural networks. Li *et al.*[21] introduced a variant of GNNs called gated graph neural networks (GG-NNs) with gated recurrent units (GRU) [22] for updating information between nodes. Li *et al.*[21] demonstrated the use of GG-NNs on graph reachability problems and bAbI tasks [23]. They have also been used to learn from knowledge graphs for image classification [24] and for learning properties of chemical molecules [25]. There are variations of GNNs available in literature for different applications. In [21], the authors propose another architecture - gated graph sequence neural networks (GGS-NNs), which is a sequence of GG-NNs producing intermediate outputs. Graph convolutional networks (GCNs) [26][27] are inspired from convolutional neural networks (CNNs) and learn local receptive fields on graph structured input. In the context of multimodal semantic forensics, GNN provides a way to assimilate multiple evidences.

## III. MULTIMODAL SEMANTIC FORENSICS

As discussed in Section I, reference datasets in image repurposing problems often have multiple evidences for verification of the query package. However, previous methods for image repurposing detection cannot handle variable number of retrieved packages [5][6]. This shortcoming prevents these models from leveraging potentially multiple instances of related packages for performance improvement. As such, a driving motivation of our model design is to make it scalable to multiple retrieved packages. Additionally, as discussed in [6], it is possible for a modality to be missing from a package or the dataset itself. For example an image may be accompanied by text or location information or both. Under such circumstances, it is also important to keep the model flexible to handle an arbitrary number of modalities. In order to ensure this, our model processes each modality in a different branch which is architecturally the same. Effectively, each branch has the same architecture, but with different learned weights for each modality. Each branch has three major
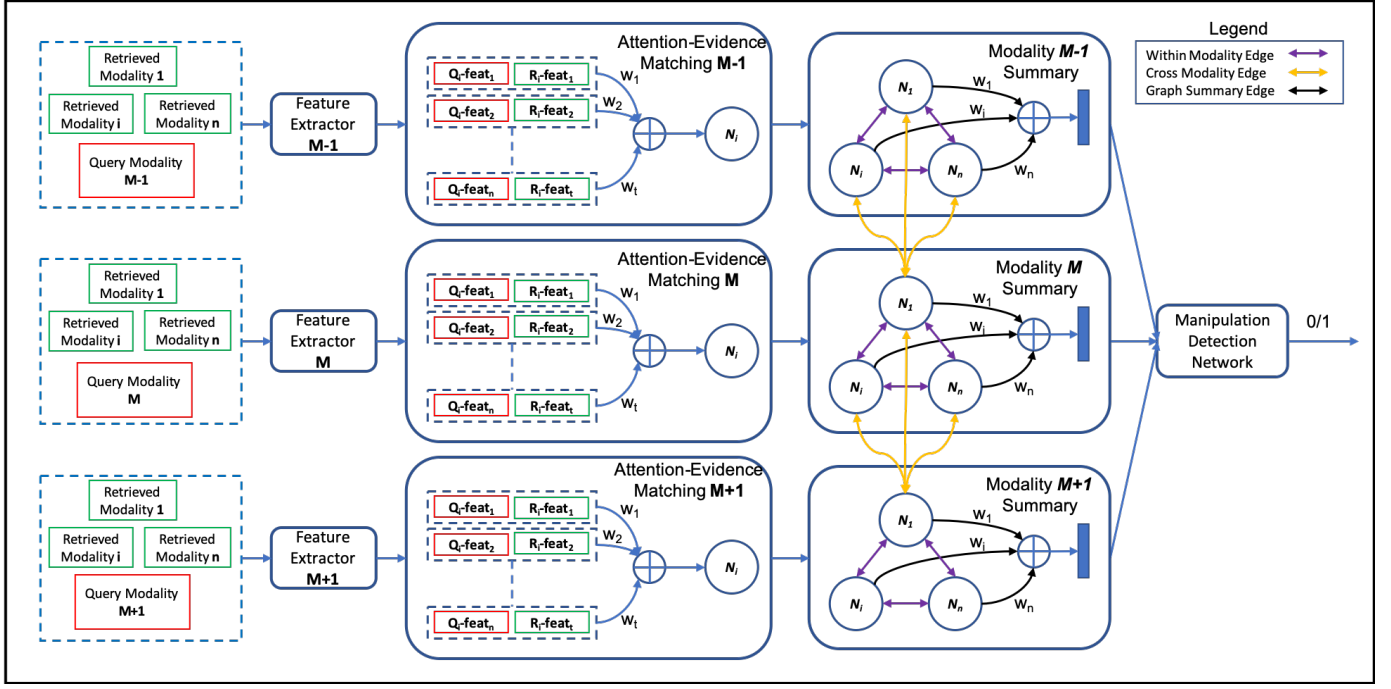
Fig. 2. Our Model diagram. Modalities from each retrieved package are organized together and are processed through a dedicated branch of the modality. The Evidence Matching layer matches query and retrieved features side-by-side and weights it to produce a node initialization. A graph neural network is used to summarize each modality for final manipulation detection. The crossmodal connections form a complete graph in implementation, but few connections are shown for simplicity.

components - (1) feature extraction, (2) evidence matching and (3) modality summary. They are preceded by package retrieval for evidences and followed by the manipulation detection layer. Figure 2 gives an overview of our model. We describe the motivation and design of each component of our model below.

### A. Package Retrieval

Verification of a query package requires additional information from a reliable reference dataset. Since packages retrieved from the reference dataset form the basis for authentication, the package retrieval system is an important component of the overall method. We use a package retrieval system similar to [6]. In [6], the authors score each modality of the query package against the corresponding modality of all packages in the reference dataset. A reference package with the top-1 combined score across all modalities is retrieved. We extend this method to retrieve $k$ packages with the highest scores.

### B. Feature Extraction

Learned feature extraction is an important component of deep learning models. Previous literature in image repurposing detection has used pretrained models to extract features from all modalities of a package. We follow a similar approach using convolutional neural network (CNN) based models, word2vec and global positioning system (GPS) coordinates for image, text and location feature extraction respectively. Specific details on models used for feature extraction are discussed in Section IV-A.

### C. Attention-based Evidence Matching

Semantic forensics can involve a subtle but specific change of detail which can be hard to detect at a glance. Additionally, the information (entity, location, etc.) involving both the manipulation and evidence in query and retrieved packages is unlikely to be previously seen. It is therefore prudent to develop a method that can compare a previously unseen instance of manipulation and evidence without memorizing it. This requirement is in contrast to classical computer vision models that reward memorization of training examples such as associating the word *dog* with a corresponding image in a standard classification task. We address the problem of dealing with previously unseen manipulations and evidences with an attention-based evidence matching module, shown in Figure 2. Evidence matching compares concatenated query and retrieved features with a soft attention mechanism for selecting important matches. For query and retrieved features $q$ and $r$ respectively, a concatenated feature vector $[q, r]$ is processed with 1D convolution network (CNN) for matching. A soft attention model on top of the concatenated representation, followed by a dense layer to compute a matched feature $feat$. This layer can be represented by Equation 1

$$feat = FC\Big(\sigma\big(Conv([q, r])\big) \odot ([q, r])\Big) \tag{1}$$

where $Conv$ is a 1D CNN, $\sigma$ is soft attention and $FC$ is a dense layer for dimensionality matching across modalities.

## D. Modality Summary

The retrieved packages represent a *bag-of-packages* without specific order. We use a graph neural network (GNN) for each modality that considers all possible comparisons between the query and retrieved packages. The graph network makes the overall system (1) flexible enough to scale to an arbitrary number of retrievals and (2) invariant to the order of retrieved packages. Each node in the graph network is updated with respect to its adjacent nodes allowing simultaneous updates. The graph is then summarized into one modality-level representation. Each node contributes directly to the final graph summary. This is different from recurrent networks where the latent embedding is updated in sequence, making them order-dependent. A node $v$ in a GNN is represented by the hidden state $h_v$. A forward pass through a GNN is divided into propagation and output steps. The propagation step updates nodes along edges in the graph for $T$ timesteps. It can be thought of as a gated recurrence along paths in the graph, similar to long short-term memory network (LSTM) recurrence. The output step produces a graph level vector representation by combining hidden states of nodes with an attention mechanism. The model is summarized by Equations 2-7

$$h_v^1 = [x_v, 0]^T \tag{2}$$

$$a_v^{(t)} = A_v^T [h_1^{(t-1)} ... h_N^{(t-1)}]^T + b \tag{3}$$

$$z_v^t = \sigma(W^z a_v^{(t)} + U^z h_v^{(t-1)}) \tag{4}$$

$$r_v^t = \sigma(W^r a_v^{(t)} + U^r h_v^{(t-1)}) \tag{5}$$

$$\widetilde{h_v^{(t)}} = tanh(W a_v^{(t)} + U(r_v^t \odot h_v^{(t-1)})) \tag{6}$$

$$h_v^{(t)} = (1 - z_v^t) \odot h_v^{(t-1)} + z_v^t \odot \widetilde{h_v^{(t)}} \tag{7}$$

The hidden state is initialized with an initial representation $x_v$ according to the application and padded with 0 to match dimensions if needed. $A$ is the adjacency matrix and $a_v$ is the summation of adjacent node embeddings based on edge type. Equations 4-7 represent updates using a GRU. The graph neural network effectively summarizes the potential for manipulation in a learned representation for each modality. We use a complete graph (adjacency matrix of ones except along the diagonal) with one timestep of propagation. This allows simultaneous update of all nodes throughout the graph in an order agnostic manner. The final graph output $G_m$ for modality $m$ of our model is a weighted average of activation all $N$ nodes as shown in Equation 8:

$$G_m = \sum_{v=1}^{N} \left( h_v^1 \odot Att(h_v^1) \right) \tag{8}$$

The weights are estimated using a neural network $Att$. Since the model is set up for variable number of inputs, the scale of adjacent node embeddings $a_v$ may fluctuate by an order of magnitude. To control for the variation, we modify Equation 3 by scaling it down by the number of adjacent nodes. For our fully connected graph setup, modality $m$ with $N$ nodes has $N - 1$ adjacent nodes, as shown in 3.

$$a_v^{(t)} = \frac{A_v^T [h_1^{(t-1)} ... h_N^{(t-1)}]^T + b}{N - 1 + \epsilon} \tag{9}$$

This summarizes each modality into a single graph output, with nodes of the same modality. However, it has been shown that cross-modal learning helps with multimodal tasks [28]. To incorporate cross-modal learning into our model we add cross-modal graph connections. The adjacency matrix is expanded to include nodes from adjacent modalities. We validate the performance of cross-modal connections later in Section IV-B. For $m$ modalities, each with $N_m$ nodes, a general update to Equation 9 for arbitray nodes in adjacent modalities is shown in Equation 10.

$$a_v^{(t)} = \frac{A_v^T [h_1^{(t-1)} ... h_N^{(t-1)}]^T + b}{N_i - 1 + \sum_{j=1, j \neq i}^{m} N_j + \epsilon} \tag{10}$$

However, considering that each modality generates an equal number of nodes and fully-connected cross-modal edges are used, Equation 10 is simplified to Equation 11.

$$a_v^{(t)} = \frac{A_v^T [h_1^{(t-1)} ... h_N^{(t-1)}]^T + b}{m * N - 1 + \epsilon} \tag{11}$$

The $\epsilon$ term takes care of zero adjacency for each node in a graph. Finally, the output representation for each modality is combined as described next.

## E. Manipulation Detection

A feed-forward network on top of concatenated modality summary outputs is used for the final manipulation detection. This layer combines all branches of modalities into a single binary prediction.

## F. Implementation Details

Our model was implemented in Keras and trained with ADAM optimizer with an initial learning rate of 0.001. All parameters had default values, unless otherwise mentioned. All edge layers and feedforward layers have ReLU activation function. We trained all our models with a batch size of 32 and subsampled models within each epoch for selecting the best model.

## IV. EVALUATION

This section describes benchmarks in Section IV-A and both quantitative and qualitative results in Section IV-B.

## A. Benchmark Datasets

We perform experimental evaluation on MEIR [6], which is the most challenging dataset for image repurposing detection. We also evaluate on *Google Landmarks* [29] and *Painter by Numbers* [30] datasets which were originally released for different tasks, but can be adapted for semantic forensics. The adapted splits for Google Landmarks and Painter by Numbers used in [15] have repeated locations and painters

in training and testing. Keeping in line with the idea in [6] that manipulations are unseen in training and test, we adapt a new split with mutually exclusive manipulations in training and test set.

*a) MEIR::* It is a multimodal dataset comprising images, text and location modalities. Manipulations are present in text and location modalities and comprise three types of entity manipulations — person, location and organization. Manipulations are also coherent within a package i.e. a location manipulation within text will result in corroborating manipulations in GPS coordinates. The dataset comprises 82,156 packages in reference dataset and 57,940 packages split between training, test and validation sets. It should be noted that all packages in training, test and validation conform to different events. This helps in evaluating the generalizability of models to unseen semantic manipulations.

*b) Google Landmarks::* We use Google Landmarks for further evaluating location manipulation, since locations can easily confuse people, especially if the landmark in the photo is not well recognized by the person. This is one of the manipulations present in MEIR, but is mixed with all other manipulations. This dataset is available as a part of a Kaggle competition[1]. The modified task for semantic forensics on this dataset is to identify if the landmark associated with a query package is correct. The complete dataset is extremely large with over 1.2 million images and 14,951 different landmarks. We prune the dataset, keeping landmarks with at least 3 images and at most 50 images. This leaves us with 152,074 images split into 78,573 images for reference and 73,501 images for train, test and validation which is further split in a 70-10-20 ratio. We create believable semantic manipulations by swapping similar images. Images are determined to be similar using a kd-tree search. During test, we ensure that landmarks in training, test and validation form a disjoint set. This ensures that the model is robust in identifying unseen landmarks. Image features are generated using NetVLAD [31], followed by principal component analysis (PCA) and $l_2$ normalization as used in [15]. Landmark IDs are encoded using 50 dimensional random embeddings. We also measure our retrieval accuracy using mean average precision (MAP). With this scheme of manipulation and feature embedding, a cosine similarity based retriever as described in [6], but generalized to 5 package retrieval achieves 0.81 mean average precision.

*c) Painter by Numbers::* We evaluate the proposed system on painting forgeries, which is an old problem with counterfeits being created for paintings by famous artists. This is a high stakes problem with art experts being called in to validate paintings. We create a painting repurposing dataset from the *Painters by Numbers* dataset. This dataset was also released as a part of a Kaggle challenge[2]. We restructure the dataset for semantic forensics, where the identity of the artist for a given painting is potentitally manipulated. After ensuring

TABLE I
ABLATION EXPERIMENTS FOR VERIFYING SCALABILITY. WE REPLACE THE MODALITY SUMMARY (GNN) COMPONENT OF MEG WITH OTHER MODELS. ALL VARIANTS ARE TRAINED ON TWO AND TESTED ON FIVE PACKAGES. AUC SCORES ARE REPORTED.

| Ablation Model | Train on 2 Packages | Test on 5 Packages | Relative Drop |
|---|---|---|---|
| **MEG (Ours)** | **0.91** | **0.91** | **0%** |
| MEG - GNN + GRU | 0.90 | 0.88 | 20% |
| MEG - GNN + LSTM | 0.90 | 0.85 | 50% |
| MEG - GNN + RPW | 0.89 | 0.89 | 0% |
| MEG - scaling | 0.90 | 0.85 | 50% |

that each artist has at least three images, the dataset is split into 36,669 reference images, and 36,164 images for train, test and validation. There are 1000 different artists in the dataset. To create manipulations, we use a kd-tree for finding similar paintings and swap the artists. There is no overlap between artists in training, test and validation, to ensure generalization. Image features are extracted using the winner's model from the competition[3], similar to [15]. For painter IDs, we generate 50 dimensional random embeddings. Again, using the retrieval scheme in [6], for top-5 packages, we achieve 0.72 mean average precision.

*B. Evaluation Results*

We use accuracy, area under the receiver operating characteristic curve (AUC), and $F_1$-clean and $F_1$-tampered ($F_1$ scores for unmanipulated and manipulated class respectively) scores as evaluation metrics. Previous works have used these metrics [6][5]. We perform ablation experiments to test the scalability and order invariance of our model. A summary of the results is discussed. We also evaluate our model on benchmark datasets and discuss the quantitative and qualitative results.

*a) Scalability::* A contribution of our model is the ability to handle variable number of related packages. The modality summary module of our model is responsible for providing scalability. Keeping this in mind, we perform two categories of ablation experiments as shown in Table I: (1) replacing the modality summary network with standard recurrent networks (GRU and LSTM) and the read-process-write (RPW) network from [18] (2) removing scaling modifications we made to the graph network in Section III-D. For this set of experiments we train our model for up to two packages and test on 5 packages. A drop in performance at 5 packages indicates that the model does not scale. We train all models with a minimum of two packages to avoid an empty adjacency matrix for GNN. The results clearly support the scalability of the proposed model.

*b) Order Invariance::* A result of the GNN based modality summary module in our model is invariance to input ordering. We perform ablation experiment by replacing the modality summary module with standard recurrent networks (GRU and LSTM) and an existing order agnostic model - read-process-write (RPW) network from [18]. It has been reported that recurrent networks suffer from order dependence

| Ablation Model | Before | After | Relative Drop |
|---|---|---|---|
| MEG (Ours) | **0.92** | **0.92** | **0.0%** |
| MEG - GNN + GRU | 0.91 | 0.83 | 88.8% |
| MEG - GNN + LSTM | 0.91 | 0.87 | 44.4% |
| MEG - GNN + RPW | 0.89 | 0.89 | 0.0% |

| Ablation Model | Order Invariance | Scalability | Score |
|---|---|---|---|
| MEG (Ours) | ✓ | ✓ | **0.92** |
| MEG - GNN + GRU | | | 0.91 |
| MEG - GNN + LSTM | | | 0.91 |
| MEG - GNN + RPW | ✓ | ✓ | 0.89 |
| MEG - scaling | ✓ | | 0.90 |

issues, resulting in performance drop for input order changes between training and test [18]. We train our model for 5 packages and test by reversing the training order. A drop in performance indicates that the model variation is not model invariant. Results are presented in Table II. It is evident that LSTM or GRU based variations of our model are not order invariant. As expected, replacing GNN with RPW [18] in the modality summary layer maintains order invariance, but leads to a performance drop.

*c) Ablation Summary::* We summarize the ablation results in Table III. Three properties are considered: scalability, order invariance and detection performance. Our method satisfies all properties while maintaining best performance for all comparisons.

*d) Performance::* We compare performance against previous methods from [6] - namely the deep multimodal model (DMM) and the semantic retrieval system (SRS). DMM is a deep learning based model which verifies a query package using top-1 retrieved package. SRS is a non-learning method which computes the Jacardian index on packages retrieved by individual modalities. It's performance is known to scale with the correctness of retrievals. Our model improves upon state-of-the-art performance across all three datasets as shown in Table IV.

*e) Analysis::* Examples from Painters by Numbers and Google Landmarks datasets are shown in Figure 3 and 4 respectively. True positive and false negative examples from MEIR are shown in Figure 5 and 6 respectively. It is visible from the results that image repurposing performance is dependent on package retrieval performance. To further test this hypothesis, we compare the average number of correct packages retrieved between successful (true positive and true negative) and unsuccessful (false positive and false negative) classifications. The results in Table V show a consistently
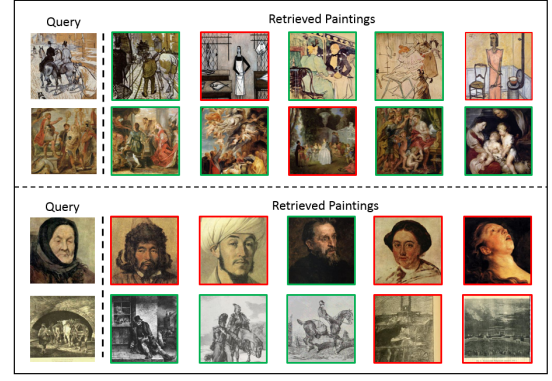


Fig. 3. The top two rows contain true positive examples and the bottom two rows contain false positive samples. Across both cases it is noticeable that the repurposing/manipulation is believable. In the bottom row, the correct retrievals are visually different from the query, leading to false alarms. Green and red borders indicate correct and incorrect retrievals respectively.
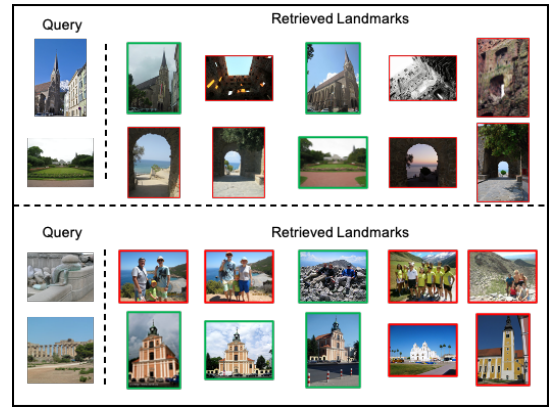


Fig. 4. The top two and bottom two rows contain true positive and false positive samples respectively. In the bottom row, the correct retrievals look significantly different, leading to a false alarm. Green and red borders indicate correct and incorrect retrievals respectively.

better retrieval for all correctly classified packages.

## V. CONCLUSION AND FUTURE WORK

Image repurposing detection is an important but emerging research area for multimodal semantic forensics and fake news detection. We presented a multi-evidence GNN model (MEG) for multimodal semantic forensics that improves upon previous state-of-the-art across three benchmark datasets. Our scaling modifications over a standard GNN make the proposed model scalable to multiple retrieved packages. Our model is order invariant compared to standard recurrent architectures.

Besides the improvements to image repurposing detection in this paper, there are still unexplored problems remaining. MEG does not localize the exact manipulation. While successful manipulation detection can alert users to semantic manipulations, successful localization can help users reason about manipulations. Another possible area to explore is real-time multimodal semantics i.e. using the web instead of a reference dataset. These directions are left for future work.

TABLE IV
PERFORMANCE OF OUR PROPOSED MODEL (MEG) AGAINST EXISTING METHODS FROM [6] ACROSS ALL THREE BENCHMARK DATASETS.

| Metric | MEIR | | | Painter by Numbers | | | Google Landmarks | | |
|---|---|---|---|---|---|---|---|---|---|
| | SRS | DMM | MEG | SRS | DMM | MEG | SRS | DMM | MEG |
| $F_1$-clean | 0.51 | 0.80 | **0.84** | 0.70 | 0.59 | **0.83** | 0.82 | **0.87** | **0.87** |
| $F_1$-tampered | 0.66 | 0.80 | **0.84** | 0.80 | 0.67 | **0.79** | 0.86 | **0.87** | **0.87** |
| Accuracy | 0.60 | 0.80 | **0.84** | 0.76 | 0.63 | **0.82** | 0.84 | **0.88** | 0.87 |
| AUC | 0.67 | 0.88 | **0.92** | 0.77 | 0.74 | **0.86** | 0.93 | 0.93 | **0.94** |



Fig. 5. The two rows show true positive samples of our model. The first and second package have location and organization manipulation respectively. Green and red borders indicate correct and incorrect retrievals respectively. Metadata highlighted in red in query package is manipulation.

TABLE V
WE MEASURE THE AVERAGE NUMBER OF CORRECTLY RETRIEVED PACKAGES OUT OF TOP-5 RETRIEVALS FOR CORRECTLY CLASSIFIED (TRUE POSITIVE AND TRUE NEGATIVE) AND MISCLASSIFIED (FALSE POSITIVE AND FALSE NEGATIVE) QUERY PACKAGES. PACKAGE RETRIEVAL ACCURACY POSITIVELY AFFECTS FINAL MODEL PERFORMANCE.

| Dataset | Classification Category | |
|---|---|---|
| | TP+TN | FP+FN |
| MEIR | **3.05** | 2.60 |
| Painter by Numbers | **3.25** | 1.00 |
| Google Landmarks | **3.05** | 2.15 |

## REFERENCES

[1] M. Schmierbach and A. Oeldorf-Hirsch, "A Little Bird Told Me, So I Didn't Believe It: Twitter, Credibility, and Issue Perceptions," *Communication Quarterly*, vol. 60, no. 3, pp. 317–337, Jul. 2012. [Online]. Available: https://doi.org/10.1080/01463373.2012.688723

[2] D. Spohr, "Fake news and ideological polarization: Filter bubbles and selective exposure on social media," *Business Information Review*, vol. 34, no. 3, pp. 150–160, Sep. 2017. [Online]. Available: https://doi.org/10.1177/0266382117722446

[3] H. Allcott and M. Gentzkow, "Social Media and Fake News in the 2016 Election," *Journal of Economic Perspectives*, vol. 31, no. 2, pp. 211–236, May 2017. [Online]. Available: https://www.aeaweb.org/articles?id=10.1257/jep.31.2.211

[4] S. Vosoughi, D. Roy, and S. Aral, "The spread of true and false news online," *Science*, vol. 359, no. 6380, pp. 1146–1151, 2018.

[5] A. Jaiswal, E. Sabir, W. AbdAlmageed, and P. Natarajan, "Multimedia Semantic Integrity Assessment Using Joint Embedding Of Images And Text," in *Proceedings of the 2017 ACM on Multimedia Conference*, ser. MM '17. New York, NY, USA: ACM, 2017, pp. 1465–1471. [Online]. Available: http://doi.acm.org/10.1145/3123266.3123385

[6] E. Sabir, W. AbdAlmageed, Y. Wu, and P. Natarajan, "Deep Multimodal Image-Repurposing Detection," in *Proceedings of the 26th ACM International Conference on Multimedia*, ser. MM '18. New York, NY, USA: ACM, 2018, pp. 1337–1345. [Online]. Available: http://doi.acm.org/10.1145/3240508.3240707

[7] A. Gupta, P. Kumaraguru, C. Castillo, and P. Meier, "TweetCred: Real-Time Credibility Assessment of Content on Twitter," in *Social Informatics*, ser. Lecture Notes in Computer Science. Springer, Cham, Nov. 2014, pp. 228–243. [Online]. Available: https://link.springer.com/chapter/10.1007/978-3-319-13734-6_16

[8] X. Liu, A. Nourbakhsh, Q. Li, R. Fang, and S. Shah, "Real-time Rumor Debunking on Twitter," in *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, ser. CIKM '15. New York, NY, USA: ACM, 2015, pp. 1867–1870. [Online]. Available: http://doi.acm.org/10.1145/2806416.2806651

[9] K. Wu, S. Yang, and K. Q. Zhu, "False rumors detection on Sina Weibo by propagation structures," in *2015 IEEE 31st International Conference on Data Engineering*, Apr. 2015, pp. 651–662.

[10] F. Jin, E. Dougherty, P. Saraf, Y. Cao, and N. Ramakrishnan, "Epidemiological Modeling of News and Rumors on Twitter," in *Proceedings of the 7th Workshop on Social Network Mining and Analysis*, ser. SNAKDD '13. New York, NY, USA: ACM, 2013, pp. 8:1–8:9. [Online]. Available: http://doi.acm.org/10.1145/2501025.2501027

Fig. 6. The two rows show false negative samples of our model. The first and second package have location and organization manipulation respectively. Green and red borders indicate correct and incorrect retrievals respectively. Metadata highlighted in red in query package is manipulation.

[11] M. Zampoglou, S. Papadopoulos, Y. Kompatsiaris, R. Bouwmeester, and J. Spangenberg, "Web and Social Media Image Forensics for News Professionals." in *SMN@ ICWSM*, 2016. [Online]. Available: http://www.aaai.org/ocs/index.php/ICWSM/ICWSM16/paper/download/13206/12860

[12] H. Farid, "Image forgery detection," *IEEE Signal Processing Magazine*, vol. 26, no. 2, pp. 16–25, Mar. 2009.

[13] Y. Wu, W. Abd-Almageed, and P. Natarajan, "Deep Matching and Validation Network: An End-to-End Solution to Constrained Image Splicing Localization and Detection," in *Proceedings of the 2017 ACM on Multimedia Conference*, ser. MM '17. New York, NY, USA: ACM, 2017, pp. 1480–1502. [Online]. Available: http://doi.acm.org/10.1145/3123266.3123411

[14] M. A. Qureshi and M. Deriche, "A bibliography of pixel-based blind image forgery detection techniques," *Signal Processing: Image Communication*, vol. 39, pp. 46–74, 2015. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S0923596515001393

[15] A. Jaiswal, Y. Wu, W. AbdAlmageed, I. Masi, and P. Natarajan, "Aird: Adversarial learning framework for image repurposing detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 330–11 339.

[16] E. Müller-Budack, J. Theiner, S. Diering, M. Idahl, and R. Ewerth, "Multimodal analytics for real-world news using measures of cross-modal entity consistency," in *Proceedings of the 2020 International Conference on Multimedia Retrieval*, 2020, pp. 16–25.

[17] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to Sequence Learning with Neural Networks," in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 3104–3112. [Online]. Available: http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf

[18] O. Vinyals, S. Bengio, and M. Kudlur, "Order Matters: Sequence to sequence for sets," *arXiv:1511.06391 [cs, stat]*, Nov. 2015. [Online]. Available: http://arxiv.org/abs/1511.06391

[19] S. H. Rezatofighi, V. K. B. G, A. Milan, E. Abbasnejad, A. Dick, and I. Reid, "DeepSetNet: Predicting Sets with Deep Neural Networks," in *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct. 2017, pp. 5257–5266.

[20] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE Transactions on Neural Networks*, vol. 20, no. 1, pp. 61–80, 2009.

[21] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel, "Gated Graph Sequence Neural Networks," in *International Conference on Learning Representations*, 2016. [Online]. Available: http://arxiv.org/abs/1511.05493

[22] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.

[23] J. Weston, A. Bordes, S. Chopra, A. M. Rush, B. van Merriënboer, A. Joulin, and T. Mikolov, "Towards ai-complete question answering: A set of prerequisite toy tasks," *arXiv preprint arXiv:1502.05698*, 2015.

[24] K. Marino, R. Salakhutdinov, and A. Gupta, "The More You Know: Using Knowledge Graphs for Image Classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2673–2681.

[25] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural Message Passing for Quantum Chemistry," in *International Conference on Machine Learning*, 2017, pp. 1263–1272.

[26] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams, "Convolutional networks on graphs for learning molecular fingerprints," in *Advances in neural information processing systems*, 2015, pp. 2224–2232.

[27] M. Niepert, M. Ahmed, and K. Kutzkov, "Learning Convolutional Neural Networks for Graphs," in *International Conference on Machine Learning*, 2016, pp. 2014–2023.

[28] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, "Multimodal deep learning," in *Proceedings of the 28th international conference on machine learning (ICML-11)*, 2011, pp. 689–696.

[29] H. Noh, A. Araujo, J. Sim, T. Weyand, and B. Han, "Large-scale image retrieval with attentive deep local features," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3456–3465.

[30] "Painter by Numbers." [Online]. Available: https://kaggle.com/c/painter-by-numbers

[31] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN Architecture for Weakly Supervised Place Recognition," 2016, pp. 5297–5307. [Online]. Available: https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Arandjelovic_NetVLAD_CNN_Architecture_CVPR_2016_paper.html