

Multiple Parenting Relationships in Image Phylogeny: Tracking Down Forgeries and Their Creators Online*

Alberto Oliveira¹ and Anderson Rocha¹

¹ Institute of Computing – University of Campinas (UNICAMP)
Campinas – SP – Brazil

{alberto.oliveira, anderson.rocha}@ic.unicamp.br

Abstract. *Due to the large amount of images shared on the web, tracking the spread and evolution of their content have become an increasingly important problem. As an image might be a composition created through the combination of the semantic information existent in two or more source images, establishing a relationship between the sources and the composite is an ever-growing problem of interest. We name as Multiple Parenting Phylogeny the problem of identifying such relationships in a set containing near-duplicate subsets of source and composition images. To tackle this problem, this work presents a three-step solution: (1) separation of near-duplicate groups; (2) classification of the relations between the groups; and (3) identification of the images used to create the original composition. Furthermore, we extend upon this framework by introducing key improvements, such as better identification of when two images share content, and improved ways to compare this content. Evaluation of the proposed method is performed by means of quantitative metrics established for evaluating the accuracy in reconstructing phylogenies and finding multiple parenting relationships in the different datasets. Finally, we also analyze the results qualitatively, with images obtained from the web*

1. Introduction

Once a digital document is shared online, a common fate for it is to be copied, transformed, and re-shared, often with a completely different intention from the original. Multimedia files portraying the same semantic content, but diverging by minor image processing transformations, receive the name of *near duplicates*, and their detection and recognition (NDDR) has been greatly studied in the literature, often focused in the image domain [Zhao and Ngo 2009, Zhao et al. 2007].

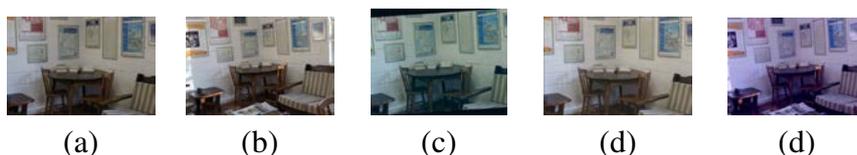


Figure 1. Example of a set of **Semantically Similar Images (SSIs)** and **Near-Duplicate Images (NDIs)**. (a) and (b), are original Images. (c) and (d) are NDIs obtained from (a) while (e) is an NDI obtained from (b).

Such works, however, overlooked any causal relationships between near duplicate images. A parenting relationship exists between *Near-Duplicate Images*(NDIs) *A* and *B*,

*We thank the support of UNICAMP, CAPES DéjàVu, CNPq, FAPESP, and the European Union.

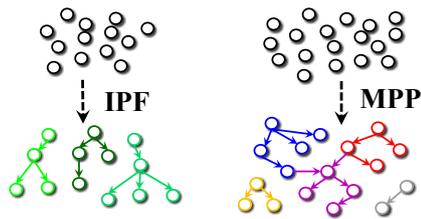


Figure 2. Image Phylogeny Forest (IPF) reconstructs a forest to represent the entire SSI set. Multiple Parenting Phylogeny (MPP) reconstructs multiple trees, finding, if existent, a relationship between a composition and the images used to create it.

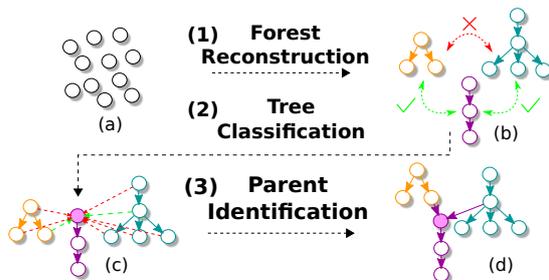


Figure 3. Multiple Parenting Phylogeny framework. (a) Set of images. (b) Searching for content relationships between IPTs in order to classify them. (c) The purple colored tree in the center, having a content relationship with the other two, is the composition tree. (d) We search for the host and alien parent of the composition root in the remaining IPTs.

when one was transformed generating the other, as Figure 1 illustrates. Recently, several works have been concerned with modeling such relationships in a set of NDIs, with the objective of better understanding how content evolves when shared and remixed on the internet [Kennedy and Chang 2008, De Rosa et al. 2010, Dias et al. 2010, Dias et al. 2012]. The work of Dias et al. [Dias et al. 2010, Dias et al. 2012] introduced the *Image Phylogeny Tree* (IPT), a directed graph representing all the parenting relationships between images of an NDI set. An asymmetric *dissimilarity* measure is employed to estimate the likelihood that an image A is parent of an image B . Once computed for all pairs of images, in both directions, the IPT is reconstructed by using a minimum spanning tree algorithm adapted to directed graphs. The root of the IPT is expected to be the original image that spawned the set.

Dias et al. [Dias et al. 2013] later expanded their work for a more generic set of *Semantically Similar Images* (SSIs). Two images are SSIs if both depict the same scene, but were not necessarily generated from the same source image. An example of the occurrence of this scenario is when two pictures are taken from the same scene, but using different cameras, as Figure 1 depicts. Costa et al. [Costa et al. 2014] further introduced several new IPF reconstruction approaches.

The aforementioned works, however, considered that a parenting relationship could only exist between two NDIs. A common scenario on the web is the creation of new images through the combination of existing ones. Splicings, montages, and mosaics are possible ways to create *compositions* by using already existing content. The parenting relationship existent between compositions and the source images used to create them gave rise to a new problem, which we refer to as *Multiple Parenting Phylogeny* (MPP). Given a set of images of varied content, the objective of MPP is to (1) identify the different phylogenies existent and (2) find the multiple parenting relationships, if any. Figure 2 contrasts the objectives of MPP and IPF. Applications of MPP range from forensic and copyright enforcement, by providing proof that an image is an forgery or created from protected content, as well as the tracking of how viral contents evolve and are created on the internet.

In this work, we tackle the MPP problem, considering *splicing* compositions, created by inserting an object extracted from an *alien* image to the background provided by a

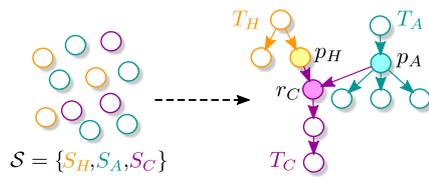


Figure 4. A multiple parenting scenario for splicing compositions. The studied set of images comprises three subsets of host, alien and composition NDIs.

host image. Given a set containing alien, host, and composition NDIs, we aim at identifying the IPTs of each of those three subsets, as well as the multiple parenting relationship existent between the subsets. For that, we introduce a three step framework: (1) Forest Reconstruction, (2) Tree Classification, and (3) Parent Identification. Figure 3 illustrates the proposed framework. In addition to the proposed framework, another contribution of this work is the creation of three datasets, comprising of compositions and their source images: We have also laid out an evaluation plan for this scenario, proposing several metrics to quantitatively assess the quality of the MPP approach. The following sections further detail this work’s contributions.

2. Multiple Parenting Phylogeny Framework

Within the Multiple Parenting Phylogeny problem, we focus our interest on *splicing*-type compositions, created by combining the object of an *alien* image with the background of a *host* image. Our test scenario is one in which three groups of NDIs exist, containing host (S_H), alien (S_A) and composition (S_C) NDIs. Each set forms an IPT, T_H , T_A , and T_C , respectively, and the root of T_C , or r_C , was created by combining a host parent p_H from T_H with an alien parent p_A from T_A . This scenario is illustrated by Figure 4. Our objectives within MPP are twofold: (1) reconstruct T_H , T_A , and T_C , and (2) identify r_C , p_A , and p_H . For this, we proposed the framework depicted by Figure 3. Next, we detail each step of the framework.

2.1. Forest Reconstruction

Before finding the MPP relationships, it is important to identify and separate which images belong to the host, alien and composition set of NDIs. Dias et al. [Dias et al. 2013] proposed the *Automatic Oriented Kruskal* (AOK) algorithm for Image Phylogeny Forests arguing that the dissimilarity between SSIs from different sources was significantly larger than the dissimilarity from NDIs. This allowed them to compute an adaptive threshold from the dissimilarity values, in order to identify the different IPTs of the forest. Costa et al. [Costa et al. 2014] later proposed the *Extended Automatic Optimum Branching* (E-AOB), employing the Optimum Branching algorithm in a similar fashion.

We compared both algorithms as NDI group separation tools in the MPP scenario. The rationale is that, if both algorithms work well in the SSI scenario, whereby images have strongly related content, it should also work well in the MPP scenario, whereby images have either strongly related (compositions and hosts), weakly related (compositions and aliens) or unrelated (hosts and aliens) content. Moreover, employing an IPF algorithm allows us to integrate MPP to the already existent Image Phylogeny frameworks. Our results have shown that, in fact, both algorithms perform really well in the MPP scenario, with a particular highlight to the performance in root identification, which was around 80% for both. Although E-AOB outperformed AOK in the controlled datasets, we observed that in the web-scenario AOK performed better, thus pointing out that the best choice for algorithm is largely dependent on the type of problem faced.

2.2. Tree Classification

Once NDIs are separated into trees, we need to identify which is the host, alien, and composition trees. To do this, we took advantage of the content relationships existent between images of different trees. We define two images as having a content relationship when some of their content is shared. A composition and a host share the background, a composition and an alien share the spliced object, while a host and an alien have no content relationship at all. Taking advantage of this information, we used randomly selected images from each IPT, and using the content relationships between each pair of them, find out which tree is most likely to contain composition images.

To discover if two images share content, we start by matching local features between them. Shared content means that parts from both images are the same, and thus many local matches between the pair should conform to the same rigid transformation. Therefore, it is possible to cluster matches based on their spatial transformation, and if a big enough cluster is observed, we assume that it is likely that the matched images share content. Additionally, by examining the area covered by the clustered matches, it is possible to reason that the shared content is the whole background, or only a small object. By registering this information when looking for shared content between images from different IPTs, we can use it after identifying which is the composition IPT to classify the remaining IPTs into host or alien.

When analyzing the complete test cases (no unrelated nodes and no missing nodes), the classification scheme has yielded very good results, with composition IPT classification accuracy of over 85% for the hardest, professionally made splicings, and over 90% classification accuracy for both host and alien IPTs for all datasets. Those results, depend strongly on a good IPF reconstruction.

2.3. Parent Identification

After classifying the IPTs, the last step of the framework consists in identifying which nodes participated in the composition process. The composition root (r_C) was created by combining one of the host nodes (p_H) with one of the alien nodes (p_A). Finding r_C comes directly from correctly reconstructing the IPF and classifying the composition IPT, as r_C is the root of that tree. For all datasets, in the complete scenario, the r_C identification accuracy ranged from 73% to 81%, considering all datasets.

Once we identify r_C , we compare it against all the candidate nodes from the host and alien IPTs in order to find the two nodes most likely to be p_H and p_A . For this, we developed the *Local Dissimilarity*, an approach to compare images considering only their shared content. First, we detect the shared content between a pair of candidate images, by clustering their local feature matches. Then, only the region inside the convex hull of the clustered matching features is used for comparison. The remainder comparison is similar to that of the original dissimilarity procedure, including registration, color adjustment and compression adjustment. The host and alien image with smaller local dissimilarity to the r_C are pointed as p_H and p_A . In the hardest, professionally made dataset, our p_H and p_A was at least 72% and 53%, respectively. It is significantly harder to identify p_A , as the region it covers is much smaller and harder to locate.

3. Datasets

For validation we built two datasets containing compositions and the source images used to create them. The first was an *Automatic Splicing Dataset*, which is further divided into two different pasting processes: *Direct Pasting* and *Poisson Blending*. In the former, we paste the alien object onto the host image as is, with no image processing techniques employed. In the later, the pasting process uses Poisson Blending to better blend the pasted object and the background. This dataset uses images from popular image retrieval and image segmentation datasets, and employs the same procedure from to create phylogeny trees. Once an IPT is generated, in the same fashion employed by Dias et al. [Dias et al. 2010, Dias et al. 2012], for the alien and host images, two nodes were randomly picked to generate the composition.

We also created a *Professional Splicing Dataset*, for which we hired a professional artist to create compositions intended to fool human viewers. Additionally, the host and alien images are all high resolution (obtained through image sharing websites).

Once the compositions were created, we generated three sets of test cases. The first comprises complete IPTs of host, alien, and composition images. The second, is the expanded set, adding unrelated images to the complete set. The last, is the missing nodes scenario, with removal of randomly chosen nodes, as well as nodes involved in the composition process (compositon root, alien parent, and host parent). Both expanded and missing nodes scenario were only created for the Professional Splicing Dataset.

4. Web Scenario

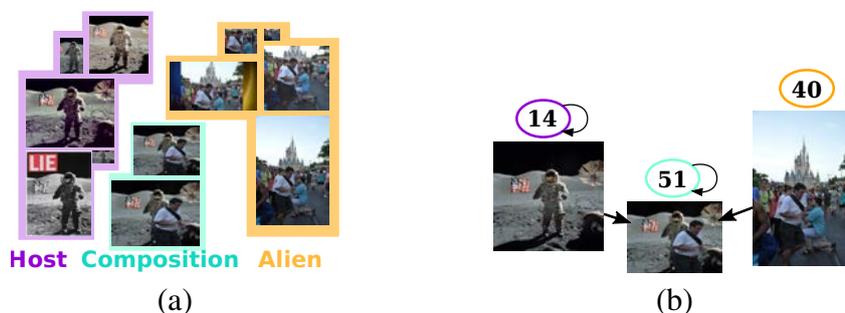


Figure 5. (a) Web scenario example images. (b) Multiple parents on web scenario. The alien parent is scaled down.

Besides the controlled scenarios mentioned before, we also validated our approach on a test case obtained from the web. This test case contained a total of 55 images, divided in 25 host NDIs, 25 alien NDIs, and 5 composition NDIs. Figure 5 illustrates some examples. Correct separation of trees is one of the biggest hurdles for MPP, as incorrectly joining the host and composition IPTs, or separating some of the IPTs in many smaller IPTs, might make it difficult to later classify the IPFs. In this regard, contrary to the controlled scenario, AOK performed better at IPF reconstruction than E-AOB, as it ended up finding a smaller number of IPTs. However, both algorithms separated the IPTs very well, with no images from different groups in the same IPT.

The classification of IPTs also proceeded correctly for the composition IPT, showing the robustness of the method. When classifying the remainder IPTs, only a single of them (composed of a single node) was mis-classified as alien, when it was a host. This

image in particular was heavily modified, thus presenting a challenge to classification. Finally, the images identified as p_H and p_A are shown by Figure 5. Although we have no way to confirm those are in fact the correct host and alien parent, they are reasonable choices, as no content present in the composition image is missing in them, which would indicate an incorrect result.

5. Conclusions and Future Research Directions

In this work, we have dealt with the Multiple Parenting Phylogeny problem, which rises when images are created by combining the content of other images. We proposed a 3-step framework for this scenario, and extensively tested it on controlled scenarios of varied difficulty, plus in one uncontrolled, web scenario. Our results show that we can, not only detect montages and compositions online, but also pinpoint the images used to create them. This is a prime step toward empowering forgery detection systems with capabilities of tracking the forgery creators other than just detecting such forgeries.

However, there are still some untied knots that need to be properly taken care of, which includes alternative forms for separating images from different NDI groups, as well as the design of better IPF reconstruction algorithms, crucial for a correct identification of the relationships between IPTs. Finally, it is also important to refine the ways that we detect and extract regions that contain similar or equal content between different images, as well as judging when shared content is not present.

6. Publications and Relevant Production

1. de Oliveira, A. A., Ferrara, P., De Rosa, A., Piva, A., Barni, M., Goldenstein, S., Dias, Z., and Rocha, A. (2015). Multiple parenting phylogeny relationships in digital images. *IEEE Transactions on Information Forensics and Security (TIFS)*, 11(2), 328-343. (**Impact Factor**: 2.408 Journal Citation Reports).
2. de Oliveira, A. A., Ferrara, P., De Rosa, A., Piva, A., Barni, M., Goldenstein, S., Dias, Z., and Rocha, A. (2014). Multiple parenting identification in image phylogeny. *IEEE Intl. Conference on In Image Processing (ICIP)*, (pp. 5347-5351). IEEE.
3. Costa, Filipe, de Oliveira, A. A., Ferrara, P., Goldenstein, S., Dias, Z., and Rocha, A. New Dissimilarity Measures for Image Phylogeny Reconstruction. *Elsevier Journal of Visual Communication and Image (JVCI)*, submitted, 2015. (**Impact Factor**: 1.218 Journal Citation Reports).

References

- Costa, F., Oikawa, M., Dias, Z., Goldenstein, S., and Rezende de Rocha, A. (2014). Image phylogeny forests reconstruction. *IEEE Transactions on Information Forensics and Security (TIFS)*, 9(10):1533–1546.
- De Rosa, A., Uccheddu, F., Costanzo, A., Piva, A., and Barni, M. (2010). Exploring image dependencies: a new challenge in image forensics. In *SPIE Media Forensics and Security II*, page 75410.
- Dias, Z., Rocha, A., and Goldenstein, S. (2010). First steps toward image phylogeny. In *IEEE Intl. Workshop on Information Forensics and Security (WIFS)*, pages 1–6.
- Dias, Z., Rocha, A., and Goldenstein, S. (2012). Image phylogeny by minimal spanning trees. *IEEE Transactions on Information Forensics and Security (TIFS)*, 7(2):774–788.
- Dias, Z., Rocha, A., and Goldenstein, S. (2013). Toward image phylogeny forests: Automatically recovering semantically similar image relationships. *Elsevier Forensic Science Intl.*, 231(1-3):178–189.
- Kennedy, L. and Chang, S.-F. (2008). Internet image archaeology: Automatically tracing the manipulation history of photographs on the web. In *ACM Intl. Conference on Multimedia (ACM-MM)*, pages 349–358.
- Zhao, W.-L. and Ngo, C.-W. (2009). Scale-rotation invariant pattern entropy for keypoint-based near-duplicate detection. *IEEE Transactions on Image Processing (TIP)*, 18(2):412–423.
- Zhao, W.-L., Ngo, C.-W., Tan, H.-K., and Wu, X. (2007). Near-duplicate keyframe identification with interest point matching and pattern learning. *IEEE Transactions on Multimedia (MM)*, 9(5):1037–1048.