# Motif Mining: Finding and Summarizing Remixed Image Content

William Theisen[1], Daniel Gonzalez Cedre[1], Zachariah Carmichael[1], Daniel Moreira[1], Tim Weninger[1], and Walter Scheirer[1]

University of Notre Dame, Notre Dame IN 46556, USA `cse@nd.edu`
http://cse.nd.edu

**Abstract.** On the internet, images are no longer static; they have become dynamic content. Thanks to the availability of smartphones with cameras and easy-to-use editing software, images can be remixed (*i.e.*, redacted, edited, and recombined with other content) on-the-fly and with a world-wide audience that can repeat the process. From digital art to memes, the evolution of images through time is now an important topic of study for digital humanists, social scientists, and media forensics specialists. However, because typical data sets in computer vision are composed of static content, the development of automated algorithms to analyze remixed content has been limited. In this paper, we introduce the idea of *Motif Mining* — the process of finding and summarizing remixed image content in large collections of unlabeled and unsorted data. In this paper, this idea is formalized and a reference implementation is introduced. Experiments are conducted on three meme-style data sets, including a newly collected set associated with the information war in the Russo-Ukrainian conflict. The proposed motif mining approach is able to identify related remixed content that, when compared to similar approaches, more closely aligns with the preferences and expectations of human observers.

**Keywords:** Image Retrieval, Image Clustering, Remixed Image Content, Digital Humanities, Computational Social Science, Media Forensics

## 1 Introduction

As the number of images posted online has grown, it has become impossible to summarize trends in this information by hand. Moreover, even though some computer vision algorithms have been proposed for this problem [32,37], the need for labelled ground truth presents a significant hurdle. The cost, in both money and time, is far too high. This is particularly evident when analyzing social trends, which often move so quickly that well-labelled data becomes obsolete by the time it is prepared. Moreover, the prevalence of remixed image content like memes, whereby an image is edited multiple times to remove information or incorporate other content, has raised questions about how related images should be associated. In this paper we describe the automated process of discovering trends in a large collection of remixed images as *Motif Mining*. Several different communities are interested in this concept, including digital humanists studying new

**Fig. 1.** Given a large unsorted and unlabelled data set from social media, the motif mining strategy groups related content — especially that which is remixed — together in an unsupervised manner. Above are two motifs mined from a new data set of images related to the Russo-Ukrainian conflict collected from the Telegram platform [28].

participatory art movements, computational social scientists studying the role of visual communication in conflicts, and media forensics specialists attempting to detect disinformation [25,33,15].

There are several challenges that must be overcome to achieve robust and accurate motif mining. To-date, the concept has been applied informally in the literature [36,3,4,29], leaving questions about optimization strategies that can be applied to the problem, as well as the structure of the output. With respect to a viable algorithm, no image feature exists that works with both globally similar images and images that are similar only in small local regions — what is commonly observed in remixed content. Additionally there has yet to be a large study of how the different combinations of image features and clustering algorithms affect the human perception of mined motifs. As the purpose of motif mining is to aid human observers, this is an important question to answer.

To address these challenges, this paper makes the following contributions:

1. A formal description of the problem of motif mining, providing a roadmap on how to more-easily discover salient trends in large, unsorted image corpora.
2. A solution to this problem in the form of an end-to-end pipeline.[1]
3. An image feature strategy for this problem combining both local and global information to aid in human-preferred clustering.
4. A new data set of over half of a million posts and remixed and static images collected from Telegram over the past six years, including the beginning of the 2022 invasion of Ukraine by the Russian Federation.
5. An empirical study of the proposed pipeline, including comparisons to related approaches on three data sets containing remixed and static image content.

**Related Work.** Within computer vision, motif mining is most closely related to content-based image retrieval (CBIR) [5], with differences in their respective inputs and purposes. CBIR takes as input (i) an image of interest (namely, a query) and (ii) a corpus of potentially related images (namely, the gallery), and aims to retrieve the images from the gallery that are similar to the query, sorting them from the most to the least similar, according to a well-defined similarity criteria. Depending upon the CBIR system user's intent, the similarity criteria

---

[1] This system will be open-sourced pending publication of this work.

may range from retrieving images that are semantically similar (*e.g.*, images that depict the same type of objects as the query), to retrieving images that are near-duplicates (*e.g.*, images that are minor variations of the query, thus sharing portions of pixels that come from the same imaging pipeline).

Regardless of the intent, CBIR solutions focus on reducing the semantic gap between the values of the image pixels and the user's objective. To overcome this gap, typical CBIR solutions operate in a multi-level image representation approach. At the lower level, either local or global features are extracted from the pixel values. While local features describe portions of the image that depict interesting visual phenomena (such as corners, edges, blobs, etc.), global features individually describe the entire image content. As expected, while many local features are usually extracted to represent a single image, only one to a few global features are used to represent the same content.

Independent of being local or global, features are either handcrafted (*i.e.*, carefully engineered by a CBIR expert who targeted a particular set of visual phenomena), or learned (*i.e.*, they are the outcome of a data-driven machine learning solution). Popular examples of handcrafted local features include SIFT [13] and SURF [2], while more recent learned local feature approaches include LIFT [35], DELF [18], and LISRD [21]. They are most commonly used for performing image registration, when different pictures of the same object or scene taken from distinct standpoints are stitched together in a post-processing operation. Handcrafted global features, in turn, comprise the concatenation of patch-wise LBP [19] and PHASH [14], only to name a few. Their common use case is also the retrieval of near-duplicate images, or different captures of the same object. More recently, the use of intermediary convolutional layers (before the fully connected layers) of popular architectures of neural networks as global image descriptors has become possible, including features from VGG [26], ResNet [9], and MobileNet [10] (hereafter MOBILE). These features are useful for the retrieval of semantically similar images because they leverage the learned content classification ability of their respective networks.

One level up, CBIR solutions aim to index the low-level features to reduce their inverted file index (IVF) storage space through feature compression and by speeding up the retrieval of feature-wise k-nearest neighbors. The standard solution in feature indexing is based on the retrieval of approximate nearest neighbor (ANN) features for each one the query's features, supported by optimized product quantization (OPQ [7]) of all the features. FAISS [11] is a popular library that implements different indexing strategies, including OPQ. Finally, at the highest level, once a set of features from the gallery images is retrieved for all the query's features, voting schemes based on IVF are used to find the gallery images that are the most similar to the query. Leveraging the voting count, gallery images can be sorted from most to least similar, constituting the desired output: a ranked list of images similar to the query.

In contrast to CBIR, motif mining takes as input a large corpus of images of interest only; there is no query to take as a reference. Moreover, rather than returning a ranked list of similar images, the purpose of motif mining is to find

different motifs, *i.e.*, groups (or clusters) of images whose similarity of interest is not known at execution time. This aspect is something to be discovered, as part of the problem formulation, being sometimes semantic (in the case of conceptually similar images) and sometimes based on pixel values (in the case of images sharing templates, such as memes containing stock character macros [25]), or even both (such as the example motif depicted on the left-hand side of Fig. 1).

Despite their differences in both input and purpose, we show that motif mining borrows from CBIR the combined use of local and global low-level features (because there is no clear definition of whether semantically similar images or near-duplicates are desired), as well as the best feature indexing strategies. In a similar fashion, recent work such as the algorithm described by Niu et al. [17] have focused on merging features of different modalities, such as visual and textual content. We instead combine only visual features. At the upper level, in turn, image clustering is performed, similar to the methods proposed for the specific case of image-based memes by Zannettou et al. [36], Beskow et al. [3], Dubey et al. [4], and Theisen et al. [29]. The novelty of the present work includes on a new formalization of the problem with reference to human-acceptable output to guide a new algorithmic design. Further, our approach is not constrained to just remixed content like memes, and it works as a method to identify visually similar static images in an unsupervised manner as well.

## 2    Formalization of Motif Mining

Traditionally, computer vision problems have been framed as an optimization over some metric calculated in reference to ground truth data for a task. Although this provides high-quality baselines for comparison, there is often little effort expended on demonstrating that higher metric scores actually result in more useful output for human observers for tasks like image retrieval. As an alternative, the methods and procedures of visual psychophysics from psychology have been recently recommended as a way to use human behavioral responses to evaluate algorithms [23,24]. Taking insight from these prior works, we can formalize motif mining.

**The Motif Mining Problem.** The purpose of motif mining is to allow human observers to quickly gain insights about visual trends in a large collection of unsorted and unlabeled data. The most common method for finding multiple trends in a given data set is clustering. However, as the purpose of motif mining is to aid people, the clusters must be optimized around some feedback mechanism. We differentiate between clusters produced with no human feedback and human-preferred clusters, and we structure our experiments around this idea.

For an example of a useful cluster, the right-hand side of Fig. 1 shows a number of different airplanes, several of which are fighter jets. This example is drawn from the current Russo-Ukrainian conflict. An increase in the number of militaristic images being posted online might prefigure action in a conflict [33] and could also potentially leak useful and/or damaging intelligence to third parties. The Ukrainian government recently addressed this concern specif-

ically, with "Ukraine's defense minister, Oleksii Reznikov, [...] calling on viewers to share images of Russia's assault" and "a local Telegram channel urged its 400,000 subscribers to 'carefully film' and share video of passing Russian troops so Ukrainian fighters could hunt them down" [8]. These examples were taken from 851 motifs mined from a subset of 16,433 images from the Ukrainian data set collected from Telegram [28]. The ideal number of clusters and the distribution of the images across them is best formalized as a clustering optimization problem with the task accuracy being derived from human feedback.
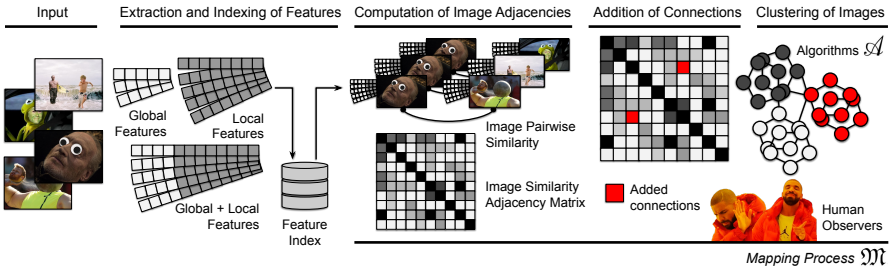
**Optimization for Human Observers.** Given a large, unlabelled corpus of images, our goal is to automatically discover trends in this data set by classifying those images that can be thought of as being "conceptually similar" or "derived from the same picture" in some intuitive sense. Because our classification task is both inherently intuitive and difficult to formally specify, and because the quantity of data far exceeds any human annotators' ability to manually label, we develop an unsupervised system for clustering these images together and verify them *a posteriori* by humans.

Here, our formal framework specifies our data set of images as a weighted graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, w : \mathcal{E} \to \mathbb{R}_+)$, where $\mathcal{V}$ is a vertex set for a graph $\mathcal{G}$ and $\mathcal{E}$ is its edge set. Here, $w : \mathcal{E} \to \mathbb{R}_+$ denotes a function that assigns positive, real-valued weights to each of the edges of $\mathcal{G}$. The vertices in this graph represent images from the data set and whose weighted edges represent the strength of the similarity between two adjacent images. Within this framework, the task becomes computing an unsupervised clustering of $\mathcal{V}$ such that it disagrees as little as possible with what human observers expect. We test this using the *intruder detection task* [31]. This means finding some subset $\mathcal{C}$ of the power set of the vertices $\mathcal{P}(\mathcal{V})$ such that, for a given pair of clusters $c_1, c_2 \in \mathcal{C}$, if a human were presented with $k$ images from $c_1$ and one intruder-image from $c_2$, the human would be able to pick the intruder. We define this formally as follows:

$$\min_{\mathcal{C} \in \mathfrak{P}(\mathcal{V})} \left( \sum_{c \neq \tilde{c} \in \mathcal{C}} \sum_{v_1, \ldots v_k \in c} \sum_{\tilde{v} \in \tilde{c}} (1 - H(v_1, \ldots v_k, \tilde{v})) \, \gamma(c, \tilde{c}) \right), \tag{1}$$

where (i) $k \in \mathbb{N}_+$, (ii) $\mathcal{G} = (\mathcal{V}, \mathcal{E}, w : \mathcal{E} \to \mathbb{R}_+)$ is the graph, (iii) $\mathfrak{P}(\mathcal{V})$ is the set of partitions of $\mathcal{V}$, (iv) $\mathcal{C}$ is a clustering on $\mathcal{G}$, (v) $v_1, \ldots v_k$ all belong to the same cluster $c$, (vi) $\tilde{v}$ belongs to a different cluster $\tilde{c}$, (vii) $H : \mathcal{V}^{k+1} \to \{0, 1\}$ returns 1 iff $\tilde{v}$ is correctly identified by a human, and (viii) $\gamma(\cdot, \cdot)$ is a normalizing factor.

In order to specify this graph $\mathcal{G}$, a mapping process $\mathfrak{M}$, with some corresponding parameters, is needed to map the image corpus onto a weighted graph. With this in mind, the problem can be further thought of as an optimization task like Eq. (1) for each weighted graph realizing the given data set. Thus, Eq. 1 will be minimized for a given clustering $\mathcal{C}$ of $\mathcal{G}$ only when human observers agree with the quality of the clustering. In order to find such a high-quality clustering, it is computationally infeasible to simply enumerate all the possible clusterings. Instead, a more principled approach would parameterize a clustering algorithm $\mathscr{A}$ suited for clustering weighted graphs and perform this optimization task over

**Fig. 2.** Motif mining pipeline. The process starts with a large set of images of interest as input and ends with the mined motifs, *i.e.*, clusters of the images that summarize remixed content. Four steps constitute the pipeline: extraction and indexing of low-level image features, computation of image adjacencies according to their pairwise similarities, computation of additional similarity connections, and clustering of the images. Human observers provide feedback on the quality of the motifs.

the set of $\mathscr{A}$'s parameters. This would, however, require an enormous number of human observers to check the quality of each clustering produced during this optimization.

In this paper, we employ a small variety of effective graph clustering algorithms and check their performance directly against the human observers for a few different realizations of the graph produced by $\mathfrak{M}$. Every time $\mathfrak{M}$ produces a weighted graph, we apply one of the clustering algorithms $\mathscr{A}_i$ to the graph and evaluate the quality of those clusters. This heuristic approach is a step in the direction of finding the true minimum alluded to in Eq. 1.

## 3    Implementation of Motif Mining

Fig. 2 summarizes our pipeline for motif mining. Each one of the steps depicted within it are detailed as follows.

**Extraction and Indexing of Features.** The first step towards motif mining is to determine the kind of features that will be used to generate the vectors in the index. This decision is informed primarily by the types of motifs an observer wishes to find in the data set. Global and local feature extraction methods will produce intuitively different results in the types of images returned by a given query and therefore the types of motifs mined. PHASH [12] features, for instance, will return only duplicate or near-duplicate images. MOBILE [10] and VGG [26] features will return images that are similar globally in a semantic sense; this often manifests in something akin to images that all contain airplanes, though those airplanes may be different shapes, sizes, positions, and styles. SURF [2] features lead to connections that may be visually distinct and share only some very small local information, such as a hand gesture or a logo on a flag.

It is the combination of these two types of features that creates compelling and robust connections. Quite frequently, SURF features will not return near duplicate images in the top of their query results as there are smaller, more subtle matches somewhere in the image. However, global visual similarity is

very apparent and useful to human observers. If, for example, military groups begin to post edited pictures of tanks or inflammatory extremist symbols, human observers will want to quickly identify these trends. To capture both of these cases, each image has a single global feature extracted in addition to its SURF features. The full set of global features is then subjected to Principal Component Analysis (PCA) [22] and each vector is reduced to 16 dimensions, to match the length of the smallest of the descriptors, namely PHASH. A further discussion about this parameter is provided in Section D.1 of the Appendix.

This global "tag" is then appended to each of the respective images' SURF features. Therefore, global context for any given image is incorporated into all of its individual local feature vectors. The effect of adding this global tag can be seen in both the airplane motif on the far right of Fig. 1, and the flag motif in the third row and third column of Fig. 3.

The local features come to the fore when the motif is a smaller object in an image, such as the Indonesian ballot boxes with "KPU" written on them as seen in column three, row two of Fig. 3.

To index the features, the proposed pipeline is implemented using the FAISS [11] library as a foundation. Available within FAISS, OPQ [7] allows for efficient mass-vector indexing and retrieval. Similar to the work previously done by Theisen et al. [29], an IVF is made with 256 centroids (an exploration of how the number of centroids affects the index and resulting graphs can be found in Section D.1 of the appendix). This index, once built, provides the function used to generate the graph that is then clustered. Readers familiar with the workings of OPQ may at this point wonder why the centroid clusters that are inherent to OPQ and already generated by FAISS are not just used for the end product, without the extra hassle of producing another graph on top of the index. This is discussed in a subsection of Section 4. Building this index allows us to quickly construct an approximate affinity matrix, and a subsequent graph, by leveraging the efficiency of OPQ.

**Computation of Image Adjacencies.** We need to perform image pairwise comparisons to compute image similarities prior to establishing the motif clusters. To compute similarities between images leveraging the index built in the previous step, we need to select a set of images as starting points to query their features from the index. Given IVF indices are built at the feature level, the indexed features have a many-to-one relationship with their respective source images. As a consequence, retrieved features need to be "mapped" back to the image from which they were extracted, since we are interested in image-level similarity. Considering the querying of the features of a selected starting-point image, each result $r \in \mathcal{R}$ is a tuple $(f, i, d)$, where $f$ is a feature corresponding to an image $i$, and $d$ is the distance between $f$ and the queried feature as computed by OPQ. If we focus on the subset $\mathcal{R}_i$ of the retrieved features that belong to image $i$, we can compute the similarity $s_i$ between that image and the selected starting-point one as follows:

$$s_i = \sum_{(f,i,d)\in\mathcal{R}_i} 1 - \tanh{(d)}. \tag{2}$$

We elect to use the nonlinear operator $\tanh(\cdot)$ as it is nicely bounded within the interval $[0, 1)$ for all non-negative $d$.

Intuitively, the nonlinear weighting rewards smaller distances and penalizes distances more harshly as they become larger.

Note that this similarity computation is a loop only for the local features (including the tagged ones). For the global features, since their relationship is 1:1 with the source image, we can simply take the single distance value returned by OPQ and compute the similarity score in a similar fashion to Eq. 2.

To realize a graph out of image similarity computations, we create an $N \times N$ adjacency matrix, whose each row/column is determined by one of the $N$ images in the data set.

We thus define entry $(i, j)$ of this matrix to be the similarity value computed through Eq. 2. The entries in this matrix will then correspond to weighted edges between vertices, which represent the images. Several strategies have been explored for selecting the starting-point images and filling in the scores in the matrix. Prior work [29] has simply taken a smaller subset of the images in the set and hoped that the resulting connections were diverse enough to form a representative graph. In this work, we instead continue selecting a random subset of isolated images in the graph (*i.e.*, images whose columns and rows within the adjacency matrix sum up to zero), until all the images are visited. This method is cheaper than querying all $N$ images while still eliminating any isolated vertices, unlike prior work. Note that this does not ensure one singular connected component, though some features (*e.g.*, SURF) do still lead to fully-connected graphs on smaller data sets.

**Addition of Image Connections.** The fact that our image similarity graph is not fully connected imposes an *a priori* constraint on any proposed vertex clusters since most clustering algorithms treat the presence of multiple components as a strong prior. So, if the connected components do not capture the images' underlying homophily (*i.e.*, if perceptually-similar images are distributed across different connected components), then this will act as a strong prior biasing even the best algorithms away from a high-quality clustering. To combat this effect, one could connect the disjoint components according to some principled (*i.e.*, data-driven) or heuristic schema. Here, we propose a random baseline connection schema, specified by an Erdős-Rényi [6] model over the connected components, and the *Best* and *Average* heuristic approaches.

The baseline Erdős-Rényi model takes a parameter $p$, which specifies the probability of adding an edge between any two components. The weights assigned to these new edges are proportional to the average weights of the edges in the two components being connected.

We find $p$ so that the expected number of new edges added to the graph is linear in the initial number of components. This avoids needlessly changing the density of the graph.

The *Best* and *Average* connection strategies work similarly to the Erdős-Rényi approach but with different strategies for determining when components get connected with each other. Given $N_C$ total components and a proposed pair

of components $C_i$ and $C_j$, these algorithms compare the components by extracting their vertices' associated feature vectors. The cosine similarity of these feature vectors then determines the "similarity" between the two components. The *Best* approach assigns a similarity to the pair $(C_i, C_j)$ based on the *most similar* pair of vertices found from $C_i$ and $C_j$. The *Average* approach, by contrast, assigns a similarity to $(C_i, C_j)$ based on the *average similarity* of their corresponding vertex pairs.

In either case, we then find a threshold $\theta$ so that the number of pairs $(C_i, C_j)$ with similarity scores above $\theta$ is proportional to $N_C$. Those pairs of components are then connected as follows: the *Best* adds edges between those vertices that had the most similar feature vectors; the *Average* approach randomly connects $k$-many pairs of vertices (by default, $k = 1$).

These new edges are weighted in proportion to the components' similarity.

**Clustering of Images.** A variety of unsupervised clustering techniques, well-suited to finding communities in weighted graphs, were tested. Louvain clustering, Markov clustering, and Spectral clustering. Spectral clustering was chosen so that results can be compared to the prior literature [29].

The *Louvain* method for community detection maximizes the modularity of the graph — a measure comparing the density within and across clusters — using a two-stage iterative optimization.

*Markov* clustering, on the other hand, is a random-walk based clustering algorithm that computes transition probabilities between the vertices of a weighted graph by modeling random walks over the graph as Markov chains.

Finally, *Spectral* clustering refers to a very popular approach to clustering data according to the eigenvalues of some similarity matrix computed over the data. In the context of clustering for graphs, the Spectral approach involves applying $k$-means clustering to the vertices of the graph using the $k$ largest eigenvalues of the graph's Laplacian matrix $L = D - A$ as features, where $A$ is the graph's (weighted) adjacency matrix, and $D$ is its diagonal.

These clustering algorithms provide an unsupervised way of exposing underlying trends in the data — the way in which remixed and static images would naturally be agglomerated by a human. These clusters are what human observers will judge during evaluation, so choosing a high-quality algorithm with sensible parameters should noticeably impact our performance. While it is not obvious which of these algorithms should be most similar to how a human would cluster the data, we can intuit that Louvain, with its modularity optimization, should maximize the separation between clusters and thus ensure as much difference between images from different clusters as is possible to infer from the graph. Markov clustering, which relies on the idea of distributing flow among the edges of the graph, might produce a clustering with comparatively "softer" boundaries. Finally, eigenvalue-based approaches like Spectral clustering, known to be related to probabilistic diffusion processes on graphs [16], rely on finding cuts in the graph based on its Laplacian's eigenvalues, and thus may intuitively fill a gap in-between Louvain and Markov clustering. Which of these, if any, most closely corresponds to human intuition can be revealed only through human evaluation.
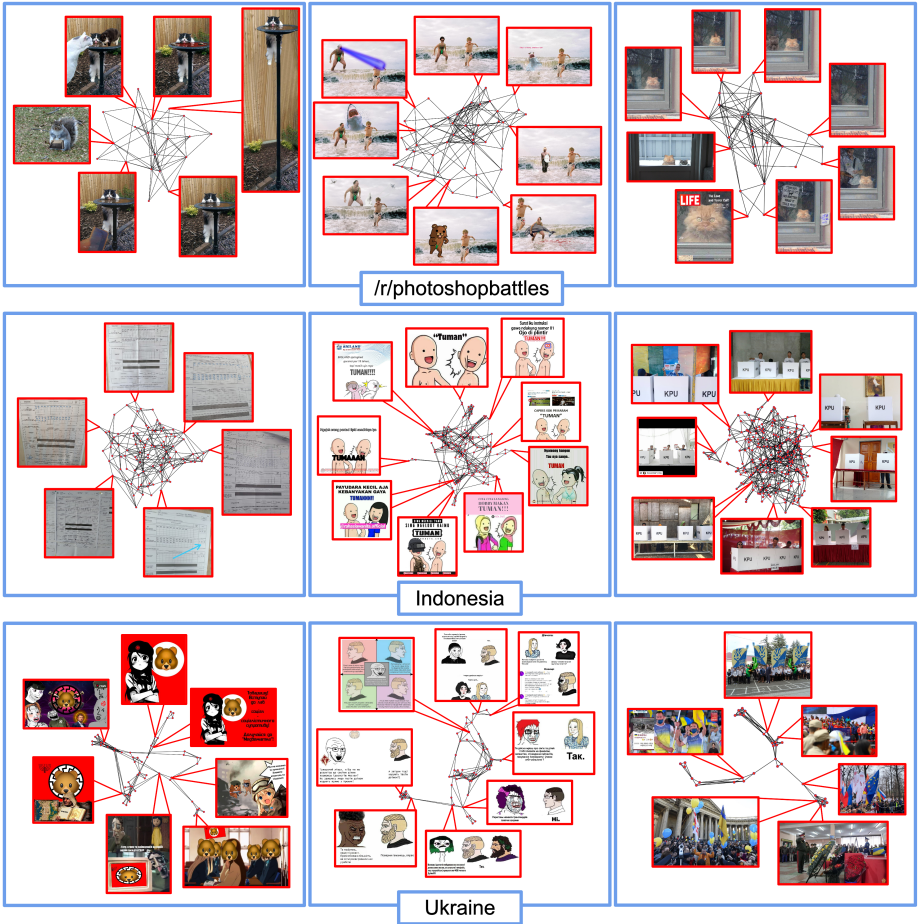
## 4    Experiments and Results

In total, 252 different configurations of the Motif Mining pipeline described in Sec. 3 were tested in order to identify the most effective ones. Figure 4 provides a summary of all of the these results, while Section B of the appendix presents them individually in a more detailed tabular form. Of primary interest is a particular combination's accuracy on the Imposter-Host test, which serves as a proxy as to whether or not the produced motifs are valuable to human observers. To describe the configurations we use the format *feature_type-connection_type-clustering_type* where feature type is one of: PHASH, MOBILE, VGG, SURF, SURF_PHASH, SURF_MOBILE, SURF_VGG. Connection type is one of: average (avg), best, Erdős-Rényi (er), or unconnected (reg). Finally the clustering method is one of Louvain, Markov, or Spectral.

**Data Sets.** The first benchmark data set used for experiments was the *Reddit Photoshop Battles* [15] data set, an image remix-specific set for content-based image retrieval and image clustering. It consists of 10,586 images taken from threads on Reddit's r/photoshopbattles subreddit. It proves to be a particularly challenging collection due to the diversity of the submitted images as seen in Fig. 3. This data set provides the 'purest' collection of remixed images, as the purpose of the subreddit is to take a piece of a 'donor' image and insert it into a variety of different sub-images. The second set consisted of 44,612 images (including memes) related to the 2019 Indonesian Presidential election (shortened to 'Indonesia'), as used in [29], to allow for an additional point of comparison to previous work.

A new data set of interest to the human rights community was also collected (referred to as 'Ukraine'). Scraped from Telegram and starting in the year 2016 and continuing through the beginning of March 2022, it focuses on content related to the Russo-Ukrainian conflict. Containing controversial images relating to nationalism, fascism, xenophobia, racism, homophobia, and the growth of militia and para-military groups, it provides an unparalleled look into the growth of online tensions surrounding the conflict. Comprised of 665,725 images and 721,441 posts, it is relevant to use for the testing of tools aimed at aiding human rights activists in the fight against online hate and disinformation. Out of this a subset of 16,433 images was selected as a test set for the purpose of these experiments.
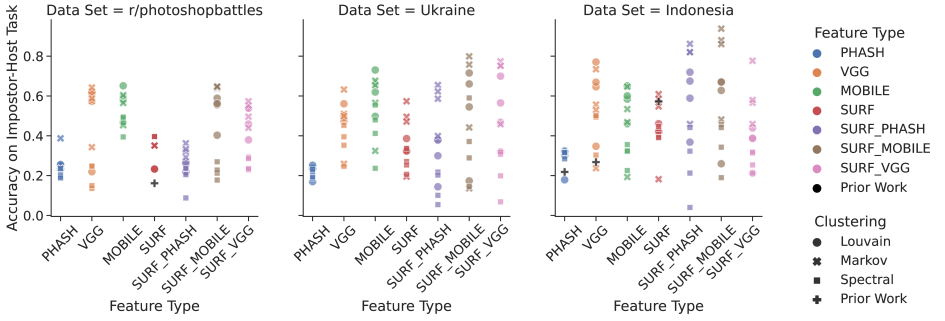
With the help of experts within Ukraine, a list of Telegram users was compiled. The channels associated with these users were then scraped, thus forming the data set. A list of these users may be found in Section F of the Appendix. Both the data set and the scraping tools will be released alongside the publication of this work. Each post consists of a JSON file containing, in addition to the post title and text, relevant meta-data such as the post date, view count at the time of scraping, image links associated with the post (Telegram allows more than one image per post), and the raw image files.

**Assessing Cluster Relatedness with an Imposter-Host Test.** As the goal of the pipeline is to create clusters for humans to review, testing the "accuracy" of the clusters requires human input. To this end, we use a version of the

**Fig. 3.** Nine motifs discovered from the Reddit Photoshop Battles dataset [15], 2019 Indonesian National Election dataset [29], and a newly collected dataset associated with the Russo-Ukrainian conflict. A variety of global motifs, local motifs, and remixed content can be seen across the nine examples. From cats to alleged voting fraud, the new pipeline can discover a diverse array of motifs in any data set.

*Imposter-Host* test [30] outlined in Theisen et al. [29], a standard way of evaluating classification tasks in a human-centric way. This test consisted of asking 50 Amazon Mechanical Turk [1] workers to find which image out of a set of five was the most different 25 times (with 5 of those 25 being control questions). Four of the images shown were from a single "host" cluster, and the fifth image — that the Turk workers are tasked with identifying — was taken from a randomly selected "imposter" cluster. Intuitively, the more related images in any given cluster are, the easier it is for the Turk workers to pick out the imposter image. Since selections are made from a set of five images, the baseline accuracy (computed by randomly picking one of the five images) is 20%. This was done

**Fig. 4.** The accuracy scores of the Imposter-Host test across the three data sets. Each of the three clustering methods is noted with a different shape. For further distinction with respect to the connection algorithm used we point to Section B of the Appendix.

for each of the 84 different possible combinations of feature types, connection types, and clustering methods and was done for all 3 data sets.

Reddit, the smallest data set, seems to imply that a global feature yields the best results in terms of observer accuracy per Figure 4. However as explained in Section 4.1, the global features are constrained by the number of centroids that the index is initialized with. For a smaller data set like Reddit, containing only 10,588 images, 256 clusters is enough to achieve high accuracy scores (this intuition is based on there being 186 reddit threads comprising the data set, which if we take as a proxy for classes implies that 186 motifs mined would be the perfect answer). However as the data sets grow larger the available number of clusters stays at 256. This is in contrast to the globally tagged features, which continue to grow in the number of clusters (though not strictly in proportion). In this work, increasing accuracy scores are seen amongst the global-local features as the number of images increases. Should the number of images continue to increase it seems intuitive to believe that the accuracy of the global features would begin to decrease as more and more visually diverse images will have to be fit to a maximum of 256 clusters, while the global-local features allow for more clusters as the data grows.

With respect to the Imposter-Host task accuracy scores, state of the art results are achieved. As seen in Table 1 prior literature on the Reddit data set achieved only a 16.15% accuracy, worse than random chance, claiming that "due to the visual complexity of the data set, they [Turk Workers] were able to find connections that weren't intended to link images" [29]. The new pipeline achieves a highest accuracy score of 65.11% using the `mobile-best-louvain` combination. This represents an increase of nearly 50 percentage points (48.96). In addition to top results on the Reddit data set, double-digit improvements were also seen on the Indonesian data. Prior work found an accuracy of 57.25% on the data. The new pipeline achieves an accuracy of 93.81% with `surf_mobile-best-markov` (a difference of 36.56 percentage points).

|  | [29] (Top) | Ours (Top - Spectral) | Ours (Top) |
|---|---|---|---|
| Reddit - PHASH | N/A | 23.53% | **38.73%** |
| Reddit - VGG | N/A | 24.79% | **64.25%** |
| Reddit - SURF | 16.15% | **39.62%** | **39.62%** |
| Indonesia - PHASH | 21.83% | 31.81% | **32.53%** |
| Indonesia - VGG | 26.79% | 50.07% | **77.05%** |
| Indonesia - SURF | 57.25% | 44.66% | **60.94%** |

**Table 1.** The top accuracy scores on the Imposter-Host task compared against Theisen et al. [29]. Shown are the top scores from the proposed pipeline using the Spectral clustering method to allow for a fair comparison to how the clustered graph is generated in previous work, followed by the top score for any combination. This shows that, except for the Reddit - SURF combination, Spectral clustering maximize accuracy.

On the new Ukrainian/Russian data set, the accuracy scores are similarly high. A top score of 79.91% is seen using the `surf_mobile-er-markov` method. Fig. 4 further shows the tagged features' high accuracy as the data set grows.

There is, however, an important caveat to this result. Much as Theisen et al. [29] found that Spectral clustering produced a singular "mega-cluster" that contained the vast majority of the images and thus making the method undesirable, the Markov clustering method had a similar issue, inadvertently skewing the results towards the higher side. Instead of placing all of the odd-out images into a single massive cluster, it placed each into its own individual cluster. With the Imposter-Host test requiring at least 4 images in a host cluster, these individual image clusters were ignored. This left a number of clusters containing only near duplicates which would thus improve the Imposter-Host accuracy scores. Due to this quirk, unless observers were interested in small near-duplicate clustering, we would recommend using the Louvain clustering method in a real-world implementation of motif mining as it results in a much more even and usable distribution of images with still state-of-the-art accuracy scores on the Imposter-Host task (a plot showing this distribution may be seen in Section C.2 of the Appendix).

**Underlying Graph and Cluster Structures.** Although the new graph creation method ensures that there are no isolated vertices in the graph, it is not guaranteed to produce a single connected component. Instead, the graph contains several interesting patterns emerging from the type of feature used in the creation of the index, and from choices relating to the initialization of the index.

To explore this further an experiment was run in which each index was recreated on the Reddit data set with a different number of centroids (128, 256, 512, 1024). The global feature types always resulted in graphs with a number of components equal to the number of centroids the index was initialized with. The global features resulting in connected components equal to the number of centroids implies that all the querying step of the pipeline is doing is exposing the

underlying centroid-space that FAISS has already prepared and therefore could be done away with entirely (see Section D.1 in the Appendix for details).

Although SURF features by themselves resulted in a highly connected graph and global features resulted in simply mirroring the underlying cluster space FAISS had already computed, the "tagged" (*i.e.*, global + local) features resulted in a higher number of components than the number of centroids, implying that further sub-structures of similar images were discovered within the centroids. The tagged features producing more components percolates down the pipeline to the clustering step, where the clusters using these features give a better image-cluster spread.

**Qualitative Results.** During the processing of the Ukrainian data an initial test of the pipeline was run with a `surf_mobile-reg-louvain` configuration, which we recommend as the best option when running in the wild due to the combination of speed, accuracy, and image distribution. The results were interesting enough that we proceeded to run it on all three data sets producing the results seen in Fig. 1 and Fig. 3. In Fig. 3 the leftmost cluster in the top row shows five images of a cat drinking out of a birdbath and a spurious match on a squirrel with a briefcase. The middle motif is of a man chasing a child out of the ocean. Here we see much stronger examples of remixing, with laser beams, sharks, and a bear all being added. The final motif is of a cat sitting at a door. This grouping highlights the usefulness of the local matching as one of the images has a stark global contrast from all of the others but the local features allow the matching on the shared cat's head.

The second row illustrates three motifs from the Indonesian data set. Again we can see a strong cluster of remixed content in the middle. On the left can be seen images of voting tallies, which the Prabowo campaign used as alleged evidence of fraud in a failed attempt to contest the 2019 election [20]. The third cluster again demonstrates why local features are extremely useful in motif mining. Many different images of Indonesians at ballot boxes were present in this motif but only shared the locally similar 'KPU' logo on the boxes.

Finally the bottom row of Fig. 3 shows results from the Ukrainian dataset. A cartoon bear head used as a logo by one of the extremist meme channels is found in the left-most panel. Note that this imagery is often superimposed over a Sonnenrad, which is a co-opted Nazi rune [27]. In the middle are several variations of a Ukrainian version of the Yes-Chad meme [34]. On the right is a cluster of Ukrainian coloured flags. Of particular interest is the photo-shopped image top-center, showing school children bearing a number of flags.

## 5   Conclusions

The newly proposed pipeline, combined with a novel combination of image features, achieves state-of-the-art results on the motif mining problem. With increases reaching nearly 50 percentage points improvement over previous methods, it demonstrates a path forward for aiding human responses to emerging trends in online social media. In addition, a new data set has been collected

from Telegram to allow further bench-marking in this space. Its timely release will allow researchers to gain unparalleled views into the increasing tensions online between Ukrainian and Russian actors, mirroring the growing tensions happening on the ground.

# 6   Acknowledgements

# References

1. Amazon mechanical turk. https://www.mturk.com/ 11
2. Bay, H., Ess, A., Tuytelaars, T., Van Gool, L.: Speeded-up robust features (SURF). Elsevier Computer Vision and Image Understanding **110**(3), 346–359 (2008) 3, 6
3. Beskow, D., Kumar, S., Carley, K.M.: The evolution of political memes: Detecting and characterizing internet memes with multi-modal deep learning. Information Processing and Management **57**(2) (2020) 2, 4
4. Dubey, A., Moro, E., Cebrian, M., Rahwan, I.: Memesequencer: Sparse matching for embedding image macros. In: In Proceedings of WWW'18 (2018) 2, 4
5. Dubey, S.R.: A decade survey of content based image retrieval using deep learning. IEEE Transactions on Circuits and Systems for Video Technology (2021) 2
6. Erdős, P., Rényi, A., et al.: On the evolution of random graphs. Publ. Math. Inst. Hung. Acad. Sci **5**(1), 17–60 (1960) 8
7. Ge, T., He, K., Ke, Q., Sun, J.: Optimized product quantization for approximate nearest neighbor search. In: 2013 IEEE Conference on Computer Vision and Pattern Recognition. pp. 2946–2953 (2013). https://doi.org/10.1109/CVPR.2013.379 3, 7
8. Harwell, D., Lerman, R.: How ukrainians have used social media to humiliate the russians and rally the world. https://www.washingtonpost.com/technology/2022/03/01/social-media-ukraine-russia/ (2022) 5
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016) 3
10. Howard, A., Sandler, M., Chu, G., Chen, L., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., Le, Q.V., Adam, H.: Searching for mobilenetv3. In: In Proceedings of ICCV'19 (2019) 3, 6
11. Johnson, J., Douze, M., Jégou, H.: Billion-scale similarity search with GPUs. IEEE Transactions on Big Data **7**(3), 535–547 (2019) 3, 7
12. Klinger, E., Starkweather, D.: phash: The open source perceptual hash library. https://www.phash.org (2013) 6
13. Lowe, D.: Distinctive image features from scale-invariant keypoints. Springer International Journal of Computer Vision **60**(2), 91–110 (2004) 3
14. Monga, V., Evans, B.L.: Perceptual image hashing via feature points: performance evaluation and tradeoffs. IEEE Transactions on Image Processing **15**(11), 3452–3465 (2006) 3
15. Moreira, D., Bharati, A., Brogan, J., Pinto, A., Parowski, M., Bowyer, K., Flynn, P., Rocha, A., Scheirer, W.: Image provenance analysis at scale. IEEE Transactions on Image Processing **27**, 6109–6123 (08 2018). https://doi.org/10.1109/TIP.2018.2865674 2, 10, 11
16. Nadler, B., Lafon, S., Kevrekidis, I., Coifman, R.: Diffusion maps, spectral clustering and eigenfunctions of fokker-planck operators. In: Weiss, Y., Schölkopf, B., Platt, J. (eds.) Advances in Neural Information Processing Systems. vol. 18. MIT Press (2005), https://proceedings.neurips.cc/paper/2005/file/2a0f97f81755e2878b264adf39cba68e-Paper.pdf 9
17. Niu, Y., Lu, Z., Wen, J.R., Xiang, T., Chang, S.F.: Multi-modal multi-scale deep learning for large-scale image annotation. IEEE Transactions on Image Processing **28**(4), 1720–1731 (2018) 4
18. Noh, H., Araujo, A., Sim, J., Weyand, T., Han, B.: Large-scale image retrieval with attentive deep local features. In: IEEE International Conference on Computer Vision. pp. 3456–3465 (2017) 3

19. Ojala, T., Pietikainen, M., Maenpaa, T.: Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. IEEE Transactions on Pattern Analysis and Machine Intelligence **24**(7), 971–987 (2002) 3
20. Paddock, R.C.: Indonesia court rejects presidential candidate's voting fraud claims. https://www.nytimes.com/2019/06/27/world/asia/indonesia-widodo-prabowo-election-fraud.html (2019) 14
21. Pautrat, R., Larsson, V., Oswald, M.R., Pollefeys, M.: Online invariance selection for local feature descriptors. In: Springer European Conference on Computer Vision. pp. 707–724 (2020) 3
22. Pearson, K.: Liii. on lines and planes of closest fit to systems of points in space. The London, Edinburgh, and Dublin philosophical magazine and journal of science **2**(11), 559–572 (1901) 7
23. RichardWebster, B., Anthony, S., Scheirer, W.: Psyphy: A psychophysics driven evaluation framework for visual recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence **41**(9), 2280–2286 (2018) 4
24. RichardWebster, B., Kwon, S.Y., Clarizio, C., Anthony, S.E., Scheirer, W.J.: Visual psychophysics for making face recognition algorithms more explainable. In: Proceedings of the European Conference on Computer Vision (ECCV) (September 2018) 4
25. Shifman, L.: Memes in Digital Culture. The MIT Press (2013) 2, 4
26. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: International Conference on Learning Representations (2015) 3, 6
27. Sonnenrad. https://www.adl.org/education/references/hate-symbols/sonnenrad 14
28. Telegram FZ LLC and Telegram Messenger Inc.: Telegram. https://telegram.org/ 2, 5
29. Theisen, W., Brogan, J., Thomas, P.B., Moreira, D., Phoa, P., Weninger, T., Scheirer, W.: Automatic discovery of political meme genres with diverse appearances. Fifteenth International AAAI Conference on Web and Social Media **15**, 714–726 (2021) 2, 4, 7, 8, 9, 10, 11, 12, 13
30. Weninger, T., Bisk, Y., Han, J.: Document-topic hierarchies from document graphs. In: In Proceedings of CIKM'12 (2012) 11
31. Weninger, T., Bisk, Y., Han, J.: Document-topic hierarchies from document graphs. In: Proceedings of the 21st ACM international conference on Information and knowledge management. pp. 635–644 (2012) 5
32. Yang, X., Deng, C., Zheng, F., Yan, J., Liu, W.: Deep spectral clustering using dual autoencoder network. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4066–4075 (2019) 1
33. Yankoski, M., Theisen, W., Verdeja, E., Scheirer, W.J.: Artificial intelligence for peace: An early warning system for mass violence. Towards an International Political Economy of Artificial Intelligence pp. 147–175 (2021) 2, 4
34. Yes chad. https://knowyourmeme.com/memes/yes-chad 14
35. Yi, K.M., Trulls, E., Lepetit, V., Fua, P.: LIFT: Learned invariant feature transform. In: Springer European conference on computer vision. pp. 467–483 (2016) 3
36. Zanneettou, S., Caulfield, T., Blackburn, J., Cristofaro, E.D., Sirivianos, M., Stringhini, G., Suarez-Tangil, G.: On the origins of memes by means of fringe web communities. In: In Proceedings of IMC'18 (2018) 2, 4

37. Zhao, J., Lu, D., Ma, K., Zhang, Y., Zheng, Y.: Deep image clustering with category-style representation. In: Springer European Conference on Computer Vision. pp. 54–70 (2020) 1

# Appendices

## A    Runtimes

The various run times for the pipeline vary widely depending on what feature type is used for extraction. PHASH and SURF features were the quickest due to their ease of parallelization and in SURF's case, its ability to run on a GPU. MOBILE features are noticeably slower but still much faster than VGG features, which took more than twice as long as MOBILE features on the Indonesia data set. It is for this reason primarily the we recommend against using VGG features. PHASH and SURF features, while fast, achieved low scores on both the Reddit and Ukraine data set in their individual forms, and slightly higher in the "tagged" combination of the two. Surprisingly the SURF_PHASH score on the Indonesia data set was quite high and comparable to the top scores. It's unclear whether this was a fluke, due to some quirk in the data set, or due to the increasing size of the data set. More work needs to be done to explore this but if speed was of the absolute essence it would be worth trying this feature combination to explore a sufficiently large data set. If speed does not matter as much we recommend a variety of the MOBILE features. It is important to keep in mind that the MOBILE features by themselves will be stuck to a number of clusters equal to the number of centroids the OPQ index is initialized with and thus we prefer the SURF_MOBILE combination which allows for more clusters and thus achieving a better image/cluster ratio.

Adding images to the index is extremely quick and should not be a serious consideration when exploring motif mining. On the other hand, the graph creation, connection, and clustering has serious run time implications. Local feature querying is significantly slower than the global features due to the voting required to map back to the images from the features retrieved from the index. On top of this, the graph connection process is costly. Do note that the runtimes for this portion of the pipeline include all three connection methods and in practice only one would need to be used. Even with the connection methods sped up using dynamic programming the BEST and AVG connection methods still averaged seven hours approximately for the Indonesia data set. We don't believe this is worth the CPU time due to no noticeable increase in the accuracy scores on the resulting graphs. Human observers seem not to notice whether or not a graph has been connected prior to clustering. The clustering run times include all three methods on all four graphs and therefore, in an implementation in which one were to run only a single combinations, are of no serious concern.

An important note is that these run-times are a sample size of 1 and therefore should only be used as rough guidelines to how long one might expect the pipeline to run, but times might vary depending on the hardware and other activities on the machine. These run times were collected on a machine with an Intel Xeon E5-2620 v3 (12) @ 3.200GHz (CPU), 256 GB of RAM and a Titan X and Titan Z (GPUs).

| Feature Extraction (Total Runtime) | Reddit (10586) | Ukraine (16433) | Indonesia (44612) |
|---|---|---|---|
| PHASH (CPU, PE=6) | 00:01:02 (00:01:17) | 00:00:17 (00:00:36) | 00:01:59 (00:02:47) |
| MOBILE (CPU, PE=1) | 00:41:10 (00:41:38) | 00:39:51 (00:40:27) | 02:04:54 (02:06:27) |
| VGG (GPU, PE=1) | 01:41:31 (01:41:54) | 02:25:53 (02:26:27) | 06:58:22 (06:59:48) |
| SURF (GPU, PE=6) | 00:05:08 (00:20:27) | 00:06:20 (00:25:14) | 00:14:05 (01:02:55) |
| SURF_PHASH (GPU/CPU, PE=6) | 00:06:30 (00:23:36) | 00:06:31 (00:29:12) | 00:14:59 (01:15:10) |
| SURF_MOBILE (GPU/CPU PE=1) | 01:14:49 (1:39:40) | 00:56:30 (01:19:26) | 02:46:13 (03:47:54) |
| SURF_VGG (GPU, PE=1) | 02:17:33 (02:30:37) | 02:35:49 (02:58:35) | 07:33:56 (08:35:24) |

**Table 2.** CPU, GPU indicates which device the feature extraction was performed on. PE gives the number of parallel processes used during the feature extraction. Due to its low overhead, PHASH is trivial to parallelize which decreases the time needed to extract features. Times are expressed in the "hours:minutes:seconds" format.

| Index Add | Reddit (10586) | Ukraine (16433) | Indonesia (44612) |
|---|---|---|---|
| PHASH | 00:00:02 | 00:00:02 | 00:00:02 |
| MOBILE | 00:00:02 | 00:00:03 | 00:00:03 |
| VGG | 00:00:03 | 00:00:02 | 00:00:02 |
| SURF | 00:00:31 | 00:00:42 | 00:02:46 |
| SURF_PHASH | 00:00:31 | 00:00:43 | 00:01:55 |
| SURF_MOBILE | 00:00:32 | 00:00:44 | 00:02:43 |
| SURF_VGG | 00:00:25 | 00:00:43 | 00:02:08 |

**Table 3.** The time spent to add all feature vectors to the index, for each feature type. Times are expressed in the "hours:minutes:seconds" format.

## B   Imposter-Host Accuracy Tables

Below are the full tabular results for the Imposter-Host test accuracy scores. The scores marked as N/A were invalid due to there be a number of clusters equal to the number of images in the data set and therefore no purpose in running the task. There is no apparent pattern in which graph connection method observers preferred and for that reason we mostly recommend against using them, for runtime purposes. However, if time is of no concern a number of top scores were produced using the BEST connection method and could be tried. While Markov clustering produced the highest scores we recommend the Louvain method due to the healthier spread of images amongst the clusters.

| Graph/Cluster Creation | Reddit (10586) | Ukraine (16433) | Indonesia (44612) |
|---|---|---|---|
| PHASH | 00:00:04/00:12:05/00:02:45 | 00:00:04/00:14:29/00:02:35 | 00:00:13/02:14:17/00:16:17 |
| MOBILE | 00:00:05/00:17:24/00:03:47 | 00:00:05/00:25:55/00:03:23 | 00:00:11/02:30:59/00:06:08 |
| VGG | 00:00:04/00:12:08/00:03:13 | 00:00:05/00:28:13/00:03:45 | 00:00:09/02:22:41/00:13:45 |
| SURF | 00:43:46/00:00:28/00:07:29 | 00:49:41/00:04:05/00:08:34 | 04:51:25/00:39:42/00:21:42 |
| SURF_PHASH | 00:32:10/00:17:18/00:03:28 | 00:59:48/01:44:55/00:05:30 | 07:11:42/14:02:41/00:13:53 |
| SURF_MOBILE | 00:43:46/00:13:07/00:03:59 | 01:08:49/02:33:48/00:05:38 | 07:57:32/14:53:36/00:15:54 |
| SURF_VGG | 00:30:07/00:09:58/00:01:53 | 01:04:06/01:50:55/00:09:00 | 05:43:41/18:56:54/00:13:42 |

**Table 4.** Times spent to create the clusters and mine the motifs, for each feature type. Times are expressed in the "hours:minutes:seconds" format.

| Reddit | Louvain | Markov | Spectral |
|---|---|---|---|
| PHASH | 23.46% - AVG | 38.73% - AVG | 19.52% - AVG |
|  | 25.18% - BEST | 23.63% - BEST | 20.34% - BEST |
|  | 24.87% - ER | N/A - ER | 23.53% - ER |
|  | 25.57% - REG | N/A - REG | 18.83% - REG |
| MOBILE | 46.62% - AVG | 45.31% - AVG | 39.44% - AVG |
|  | 65.11% - BEST | 59.55% - BEST | 48.88% - BEST |
|  | 58.96% - ER | 60.43% - ER | 49.39% - ER |
|  | 57.86% - REG | 56.49% - REG | 46.45% - REG |
| VGG | 21.93% - AVG | 34.28% - AVG | 24.61% - AVG |
|  | 61.13% - BEST | 57.48% - BEST | 13.59% - BEST |
|  | 57.27% - ER | 64.25% - ER | 24.79% - ER |
|  | 62.00% - REG | 58.76% - REG | 14.95% - REG |
| SURF | 23.29% | 35.08% | 39.62% |
| SURF_PHASH | 21.49% - AVG | 29.22% - AVG | 20.49% - AVG |
|  | 21.96% - BEST | 32.09% - BEST | 23.41% - BEST |
|  | 25.68% - ER | 36.23% - ER | 23.74% - ER |
|  | 26.65% - REG | 33.26% - REG | 08.77% - REG |
| SURF_MOBILE | 40.28% - AVG | 63.83% - AVG | 21.26% - AVG |
|  | 58.88% - BEST | 64.32% - BEST | 22.79% - BEST |
|  | 55.63% - ER | 64.96% - ER | 17.77% - ER |
|  | 56.22% - REG | 64.67% - REG | 26.94% - REG |
| SURF_VGG | 37.94% - AVG | 44.01% - AVG | 22.92% - AVG |
|  | 54.01% - BEST | 55.12% - BEST | 29.30% - BEST |
|  | 45.87% - ER | 49.62% - ER | 28.53% - ER |
|  | 53.90% - REG | 57.35% - REG | 23.49% - REG |

| Indonesia | Louvain | Markov | Spectral |
|---|---|---|---|
| PHASH | 32.53% - AVG<br>17.90% - BEST<br>32.02% - ER<br>30.12% - REG | N/A - AVG<br>N/A - BEST<br>N/A - ER<br>N/A - REG | 31.07% - AVG<br>31.81% - BEST<br>31.61% - ER<br>28.45% - REG |
| MOBILE | 46.04% - AVG<br>58.43% - BEST<br>60.06% - ER<br>65.11% - REG | 19.30% - AVG<br>64.71% - BEST<br>53.42% - ER<br>46.85% - REG | 22.55% - AVG<br>32.10% - BEST<br>32.39% - ER<br>35.55% - REG |
| VGG | 34.73% - AVG<br>64.61% - BEST<br>77.05% - ER<br>66.92% - REG | 23.72% - AVG<br>55.69% - BEST<br>52.95% - ER<br>73.46% - REG | 30.31% - AVG<br>50.03% - BEST<br>49.44% - ER<br>50.07% - REG |
| SURF | 42.67% - AVG<br>42.36% - BEST<br>46.08% - ER<br>41.71% - REG | 60.94% - AVG<br>58.58% - BEST<br>54.73% - ER<br>18.18% - REG | 40.70% - AVG<br>39.38% - BEST<br>44.66% - ER<br>39.08% - REG |
| SURF_PHASH | 36.81% - AVG<br>71.95% - BEST<br>58.89% - ER<br>67.48% - REG | 45.93% - AVG<br>81.91% - BEST<br>86.19% - ER<br>82.02% - REG | 32.39% - AVG<br>44.16% - BEST<br>03.99% - ER<br>21.26% - REG |
| SURF_MOBILE | 25.96% - AVG<br>66.99% - BEST<br>62.78% - ER<br>67.01% - REG | 48.19% - AVG<br>93.81% - BEST<br>88.02% - ER<br>86.05% - REG | 45.76% - AVG<br>32.26% - BEST<br>44.13% - ER<br>18.98% - REG |
| SURF_VGG | 21.19% - AVG<br>38.49% - BEST<br>38.75% - ER<br>43.92% - REG | 45.94% - AVG<br>56.61% - BEST<br>57.95% - ER<br>77.68% - REG | 31.28% - AVG<br>21.23% - BEST<br>25.36% - ER<br>31.76% - REG |

| Ukraine | Louvain | Markov | Spectral |
|---|---|---|---|
| PHASH | 16.98% - AVG<br>23.47% - BEST<br>23.60% - ER<br>25.19% - REG | N/A - AVG<br>N/A - BEST<br>N/A - ER<br>N/A - REG | 18.73% - AVG<br>19.36% - BEST<br>21.93% - ER<br>23.22% - REG |
| MOBILE | 49.76% - AVG<br>66.43% - BEST<br>61.99% - ER<br>73.04% - REG | 32.43% - AVG<br>67.68% - BEST<br>65.37% - ER<br>56.56% - REG | 23.66% - AVG<br>55.47% - BEST<br>47.97% - ER<br>41.21% - REG |
| VGG | 50.10% - AVG<br>56.08% - BEST<br>47.45% - ER<br>48.96% - REG | 25.91% - AVG<br>63.24% - BEST<br>51.01% - ER<br>49.96% - REG | 24.68% - AVG<br>39.51% - BEST<br>45.26% - ER<br>35.25% - REG |
| SURF | 32.86% - AVG<br>32.86% - BEST<br>32.29% - ER<br>38.61% - REG | 19.51% - AVG<br>57.36% - BEST<br>49.56% - ER<br>47.15% - REG | 33.68% - AVG<br>27.11% - BEST<br>20.64% - ER<br>25.35% - REG |
| SURF_PHASH | 14.39% - AVG<br>38.35% - BEST<br>29.97% - ER<br>37.83% - REG | 39.89% - AVG<br>65.60% - BEST<br>62.26% - ER<br>58.62% - REG | 21.78% - AVG<br>20.28% - BEST<br>05.43% - ER<br>10.09% - REG |
| SURF_MOBILE | 17.35% - AVG<br>71.55% - BEST<br>54.51% - ER<br>66.05% - REG | 44.18% - AVG<br>13.61% - BEST<br>79.91% - ER<br>75.77% - REG | 28.89% - AVG<br>37.12% - BEST<br>59.02% - ER<br>14.52% - REG |
| SURF_VGG | 31.91% - AVG<br>56.52% - BEST<br>46.76% - ER<br>69.95% - REG | 45.82% - AVG<br>77.42% - BEST<br>75.35% - ER<br>75.15% - REG | 30.74% - AVG<br>06.81% - BEST<br>19.88% - ER<br>31.81% - REG |

# C   Cluster Structures.

## C.1   Cluster Statistics.

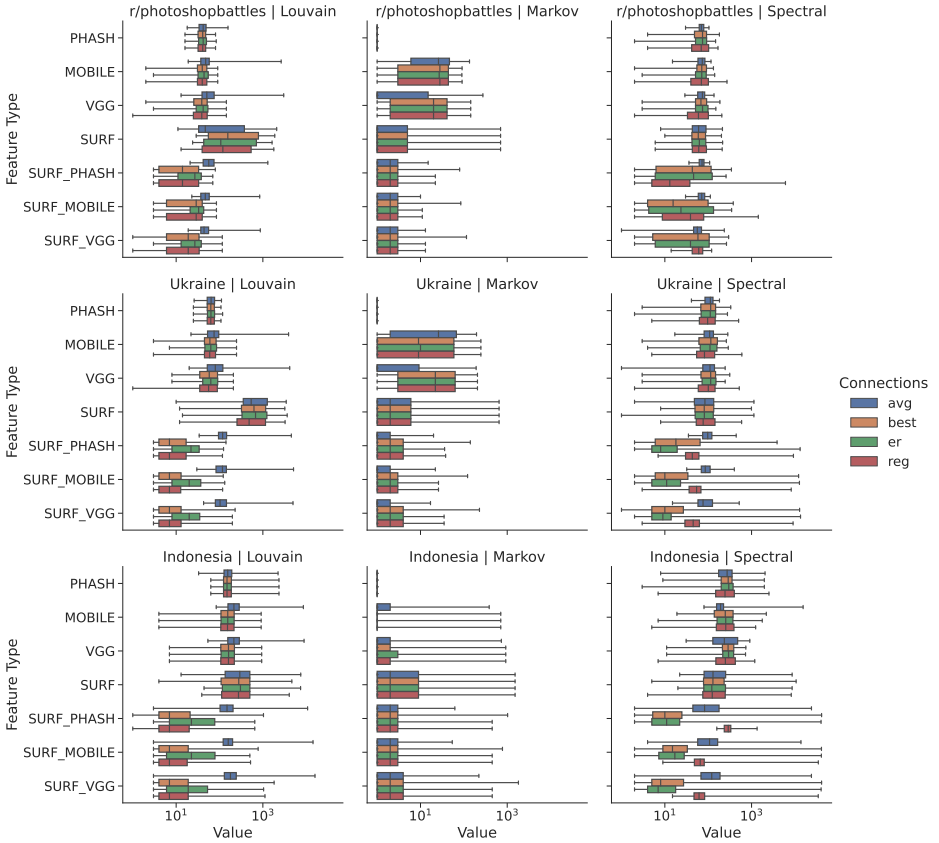| Reddit | Louvain | Markov | Spectral |
|---|---|---|---|
| **PHASH** | 244 - AVG | <u>10586 - AVG</u> | 150 - AVG |
| | 257 - BEST | 10586 - BEST | 150 - BEST |
| | 244 - ER | 10586 - ER | <u>150 - ER</u> |
| | <u>256 - REG</u> | 10586 - REG | 150 - REG |
| **MOBILE** | 164 - AVG | 355 - AVG | 150 - AVG |
| | <u>257 - BEST</u> | 394 - BEST | 150 - BEST |
| | 238 - ER | <u>393 - ER</u> | <u>150 - ER</u> |
| | 128 - REG | 161 - REG | 150 - REG |
| **VGG** | 127 - AVG | 809 - AVG | 150 - AVG |
| | 255 - BEST | 425 - BEST | 150 - BEST |
| | 236 - ER | <u>424 - ER</u> | <u>150 - ER</u> |
| | <u>256 - REG</u> | 425 - REG | 150 - REG |
| **SURF** | 28 - AVG | 760 - AVG | 150 - AVG |
| | 28 - BEST | 760 - BEST | 150 - BEST |
| | 28 - ER | 760 - ER | 150 - ER |
| | 28 - REG | 760 - REG | 150 - REG |
| **SURF_PHASH** | 158 - AVG | 4827 - AVG | 150 - AVG |
| | 537 - BEST | 4263 - BEST | 150 - BEST |
| | 408 - ER | <u>4261 - ER</u> | <u>150 - ER</u> |
| | <u>535 - REG</u> | 4260 - REG | 150 - REG |
| **SURF_MOBILE** | 203 - AVG | 5059 - AVG | 150 - AVG |
| | <u>397 - BEST</u> | 4705 - BEST | 150 - BEST |
| | 326 - ER | <u>4707 - ER</u> | 150 - ER |
| | 396 - REG | 4704 - REG | <u>150 - REG</u> |
| **SURF_VGG** | 173 - AVG | 4388 - AVG | 150 - AVG |
| | <u>394 - BEST</u> | 3905 - BEST | <u>150 - BEST</u> |
| | 319 - ER | 3904 - ER | 150 - ER |
| | 391 - REG | <u>3905 - REG</u> | 150 - REG |

**Table 5.** The number of clusters produced from each of the 52 combinations on the Reddit data set. The number that correlates with the combination that achieved the top accuracy score on the Imposter-Host task is underlined.

| Indonesia | Louvain | Markov | Spectral |
|---|---|---|---|
| PHASH | <u>256 - AVG</u><br>257 - BEST<br>256 - ER<br>256 - REG | 44612 - AVG<br>44612 - BEST<br>44612 - ER<br>44612 - REG | 150 - AVG<br><u>150 - BEST</u><br>144 - ER<br>147 - REG |
| MOBILE | 159 - AVG<br>256 - BEST<br>254 - ER<br><u>256 - REG</u> | 2264 - AVG<br><u>3187 - BEST</u><br>3187 - ER<br>3186 - REG | 150 - AVG<br>150 - BEST<br>150 - ER<br><u>150 - REG</u> |
| VGG | 154 - AVG<br>257 - BEST<br><u>254 - ER</u><br>256 - REG | 3157 - AVG<br>1609 - BEST<br>1590 - ER<br><u>1607 - REG</u> | 150 - AVG<br>148 - BEST<br>144 - ER<br><u>149 - REG</u> |
| SURF | 69 - AVG<br>72 - BEST<br><u>68 - ER</u><br>73 - REG | <u>3103 - AVG</u><br>3136 - BEST<br>3103 - ER<br>3103 - REG | 150 - AVG<br>150 - BEST<br><u>150 - ER</u><br>150 - REG |
| SURF_PHASH | 197 - AVG<br><u>1456 - BEST</u><br>846 - ER<br>1609 - REG | 16197 - AVG<br>13659 - BEST<br><u>13620 - ER</u><br>13648 - REG | 150 - AVG<br><u>146 - BEST</u><br>147 - ER<br>150 - REG |
| SURF_MOBILE | 183 - AVG<br>1597 - BEST<br>846 - ER<br><u>1609 - REG</u> | 17531 - AVG<br><u>14670 - BEST</u><br>14639 - ER<br>14668 - REG | 150 - AVG<br><u>149 - BEST</u><br>146 - ER<br>150 - REG |
| SURF_VGG | 154 - AVG<br>2000 - BEST<br>1150 - ER<br><u>2008 - REG</u> | 15280 - AVG<br>11712 - BEST<br>11687 - ER<br><u>11703 - REG</u> | 150 - AVG<br>149 - BEST<br>150 - ER<br><u>150 - REG</u> |

**Table 6.** The number of clusters produced from each of the 52 combinations on the Indonesia data set.

| Ukraine | Louvain | Markov | Spectral |
|---|---|---|---|
| PHASH | 252 - AVG | 16433 - AVG | 150 - AVG |
| | 257 - BEST | 16433 - BEST | 150 - BEST |
| | 252 - ER | 16433 - ER | 148 - ER |
| | 256 - REG | 16433 - REG | 147 - REG |
| MOBILE | 162 - AVG | 416 - AVG | 150 - AVG |
| | 256 - BEST | 511 - BEST | 149 - BEST |
| | 257 - ER | 501 - ER | 148 - ER |
| | 256 - REG | 510 - REG | 150 - REG |
| VGG | 138 - AVG | 1068 - AVG | 150 - AVG |
| | 252 - BEST | 437 - BEST | 150 - BEST |
| | 238 - ER | 437 - ER | 150 - ER |
| | 256 - REG | 436 - REG | 150 - REG |
| SURF | 17 - AVG | 2169 - AVG | 149 - AVG |
| | 18 - BEST | 2169 - BEST | 148 - BEST |
| | 17 - ER | 2169 - ER | 150 - ER |
| | 21 - REG | 2169 - REG | 148 - REG |
| SURF_PHASH | 94 - AVG | 8398 - AVG | 150 - AVG |
| | 1203 - BEST | 5453 - BEST | 148 - BEST |
| | 694 - ER | 5449 - ER | 148 - ER |
| | 1202 - REG | 5451 - REG | 150 - REG |
| SURF_MOBILE | 97 - AVG | 9084 - AVG | 150 - AVG |
| | 1282 - BEST | 5730 - BEST | 149 - BEST |
| | 658 - ER | 5722 - ER | 147 - ER |
| | 1286 - REG | 5727 - REG | 150 - REG |
| SURF_VGG | 98 - AVG | 8351 - AVG | 150 - AVG |
| | 1191 - BEST | 5523 - BEST | 146 - BEST |
| | 645 - ER | 5511 - ER | 148 - ER |
| | 1189 - REG | 5520 - REG | 150 - REG |

**Table 7.** The number of clusters produced from each of the 52 combinations on the Ukraine data set.

**Fig. 5.** Box plots of the distribution of cluster sizes for each data set and each combination of feature type, clustering algorithm, and connection type. Note that the $x$-axis has a logarithmic scale.

### C.2    Cluster Image Distributions.

As the motif mining pipeline is intended to aid and be consumed by human observers, we believe the distribution of images amongst the clusters is of the utmost importance. Figure 5 shows box and whisker plots for all of the possible combinations. While the Markov clustering algorithm delivers the highest accuracy scores on the imposter host test, it is important to realize that the majority of the clusters are of size 1, or in other words useless to reviewers. The highest realized accuracy score was SURF_MOBILE-BEST-MARKOV on the Indonesian data set. However, the second quartile for the image distribution was at 2 images per cluster and the third quartile is only 3 images per cluster. Out of these clusters only 63.38% were of a size larger than 1, and only 20.59% contained more than 3 images (I.E. valid for the Imposter-Host task). From

Figure 5 we can see that this trend holds for almost all possible combinations when Markov clustering is used. It is for this reason that we recommend Louvain clustering used with the globally tagged features. In contrast to the Markov statistics, SURF_MOBILE-BEST-LOUVAIN, on the Indonesian data set, has a second quartile at 7 images and the third quartile is 19 images. Additionally 100% of the clusters have more than 1 image per cluster and 78.46% have more than 3 images. Per Figure 5 this trend holds similar for all combinations and on all three data sets.

If one were to look at just Figure 5 they might come to the conclusion that Spectral clustering achieves a similar distribution to Louvain clustering and may wonder why the authors recommend Louvain clustering over Spectral clustering. It is for this reason that the whiskers are important. The maximum cluster size for SURF_MOBILE-BEST-SPECTRAL is 40,909 images. The data set contains 44,612 images. With 40,909 images in a single cluster this means that 91.69% of the images are essentially unsorted. We consider this case more or less useless to human reviewers, in much the same way as Markov clustering putting thousands of images into their own clusters. It is for these reasons that we believe Louvain clustering is the best of the three methods tested for motif mining.

# D   Graph Structures.

| Components, Edges | Reddit | Ukr | Indo |
|---|---|---|---|
| PHASH | 256C, 38068E | 256C, 58537E | 256C, 908593E |
| MOBILE | 256C, 38523E | 256C, 58565E | 256C, 202405E |
| VGG | 256C, 41440E | 256C, 63952E | 256C, 193731E |
| SURF | 1C, 161253E | 1C, 197858E | 14C, 475000E |
| SURF_PHASH | 412C, 24877E | 935C, 35938E | 1085C, 209128E |
| SURF_MOBILE | 336C, 21728E | 1112C, 32887E | 1237C, 203859E |
| SURF_VGG | 324C, 18837E | 984C, 33389E | 1372C, 213988E |

**Table 8.** The number of components and edges the generated graph contained for each feature type for each data set. Of particular interest is each global feature resulting in 256 components (due to the number of FAISS centroids), SURF features producing 1, 1, and 14 components (due to their locality and diversity of query results), and the married features resulting in a relatively high number of components implying the discovery of 'sub-structures' of similar images within the already calculated FAISS centroids.

## D.1   Centroid and Tag Number Experiments.

| Components, Edges | 128 Centroids | 256 Centroids | 512 Centroids | 1024 Centroids |
|---|---|---|---|---|
| PHASH | 128C, 50406E | 256C, 38068E | 512C, 25491E | 1024C, 19200E |
| MOBILE | 128C, 59138E | 256C, 38523E | 512C, 26390E | 1024C, 19937E |
| VGG | 128C, 68027E | 256C, 41440E | 512C, 27070E | 1024C, 19652E |
| SURF | 1C, 158000E | 1C, 161253E | 1C, 159750E | 1C, 157600E |
| SURF_PHASH | 233C, 24753E | 412C, 24877E | 733C, 24086E | 1257C, 22406E |
| SURF_MOBILE | 205C, 22353E | 336C, 21728E | 599C, 21466E | 1116C, 20432E |
| SURF_VGG | 200C, 19667E | 324C, 18837E | 588C, 19008E | 1083C, 18039E |

**Table 9.** The number of components and edges the resulting graphs had when the index was created with 128, 256, 512, and 1024 centroids. This shows that regardless of the number of centroids chosen all the global features accomplish is exposing the pre-existing centroid space from the OPQ index.
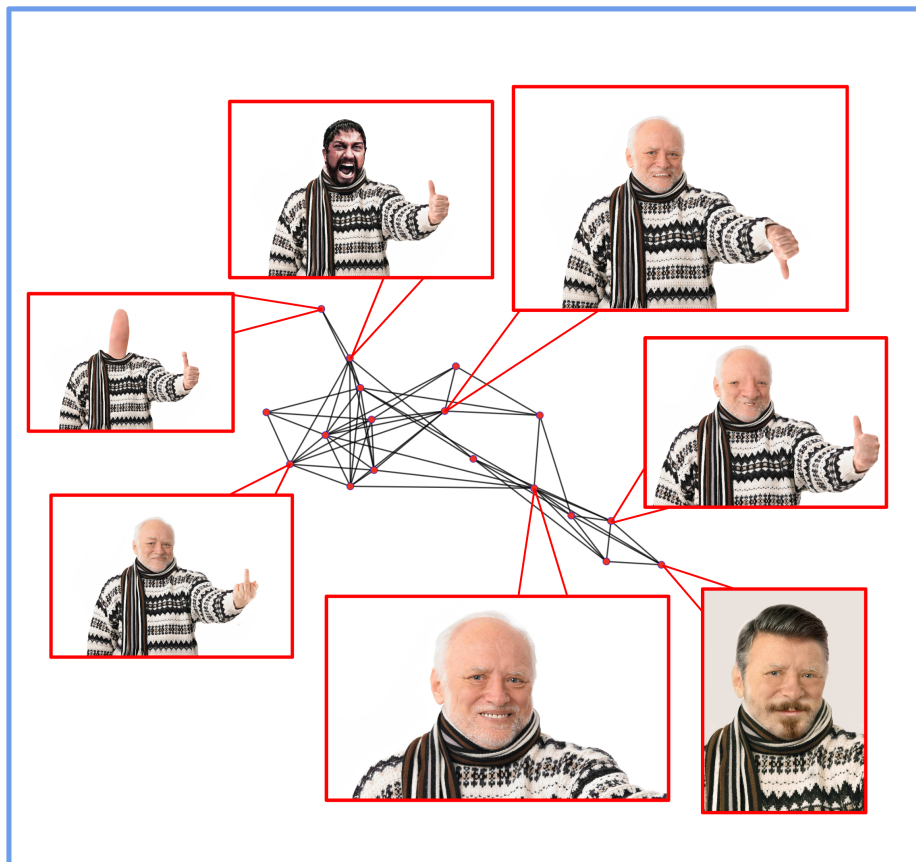
| Components, Edges | 8 Length Tag | 16 Length Tag | 32 Length Tag | 64 Length Tag |
|---|---|---|---|---|
| SURF_PHASH | 259C, 27575E | 412C, 24877E | N/A | N/A |
| SURF_MOBILE | 362C, 20789E | 336C, 21728E | 633C, 18447E | 577C, 18706E |
| SURF_VGG | 340C, 18777E | 324C, 18837E | 558C, 16040E | 533C, 16189E |

**Table 10.** How the length of the global tag affects the number of components and edges in the resulting graph. The fact that PHASH features have a length of 16 was the primary driver of that length being used. One can see however that increasing the tag almost doubles the number of components between 16 and 32. If the goal is a larger number of discrete clusters this might be a worthwhile change.

# E    Extra Figures and Data

## E.1    Meme Clusters and Examples



**Fig. 6.** An example of a Reddit cluster demonstrating the kind of visual remixing done in the data set.

**Fig. 7.** A cluster containing remixes of a stack of tree frogs. This cluster shows the usefulness of the global tag, as all the images look very similar globally and their matching can benefit from the composite feature type.

**Fig. 8.** An example of the local features being used to create a cluster on the Indonesian data set. Each image, while globally very different, contains at least part of a map of Indonesia. The local features are able to find the shared map portions in each of the images and cluster them together.

**Fig. 9.** A collection of presidential campaign ads from the Indonesian election in 2019. The same base image is used throughout but is remixed in various contexts. This kind of campaign ad remixing was common in the data set for both candidates.

**Fig. 10.** A cluster of Ukrainian memes. While initially appearing different all of the memes share the same four panel structure and a subgroup of them share the same genre on top of which various topics are remixed.

**Fig. 11.** An example of a cluster which human observers may consider interesting but we would consider a failure from an algorithmic perspective. While human observers might be interested in exploring the online meme space on Telegram, visually the images in this cluster do not have much in common. While the implemented motif mining pipeline is good, it is far from perfect and not every cluster contains a recognizable motif from a computer vision stand-point.

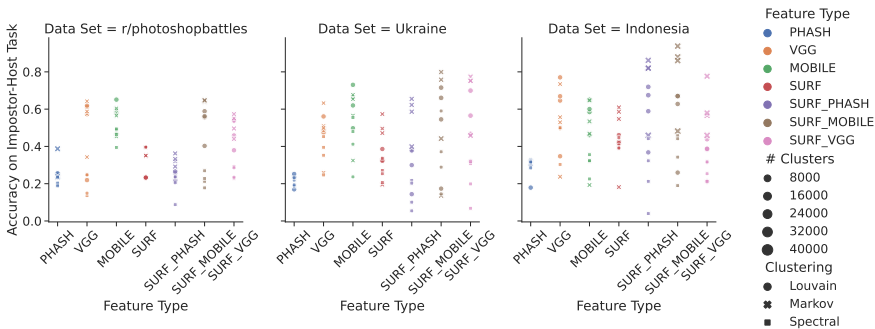## E.2    Connection Type Plots



**Fig. 12.** The accuracy scores of the Imposter-Host test across the three data sets for each connection type. Each of the three clustering methods is noted with a different shape. The size of each marker is proportional to the number of clusters.



**Fig. 13.** The accuracy scores of the Imposter-Host test across the three data sets for each connection type. Each of the three clustering methods is noted with a different shape. The size of each marker is proportional to the number of components in the graph.

**Fig. 14.** The accuracy scores of the Imposter-Host test across the three data sets for each connection type. Each of the three clustering methods is noted with a different shape. The size of each marker is proportional to the ratio of the number of components in the graph to the number of images in the dataset.
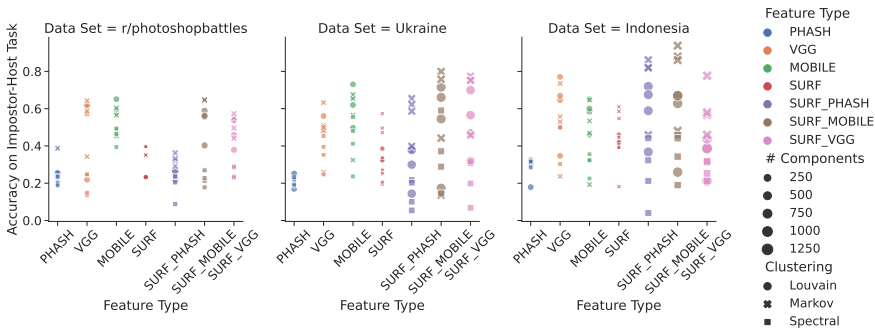


**Fig. 15.** The accuracy scores of the Imposter-Host test across the three data sets for each connection type. Each of the three clustering methods is noted with a different shape. The size of each marker is proportional to the number of edges in the graph.
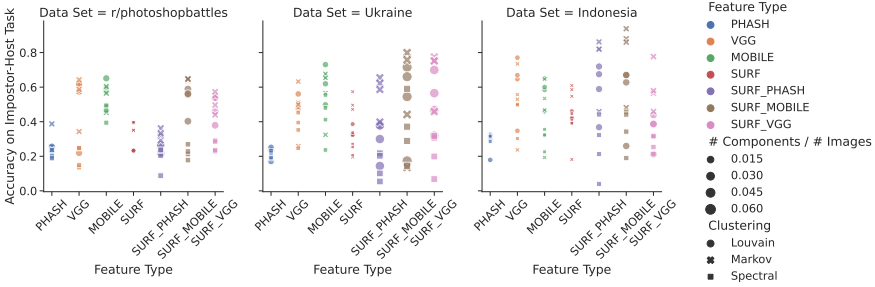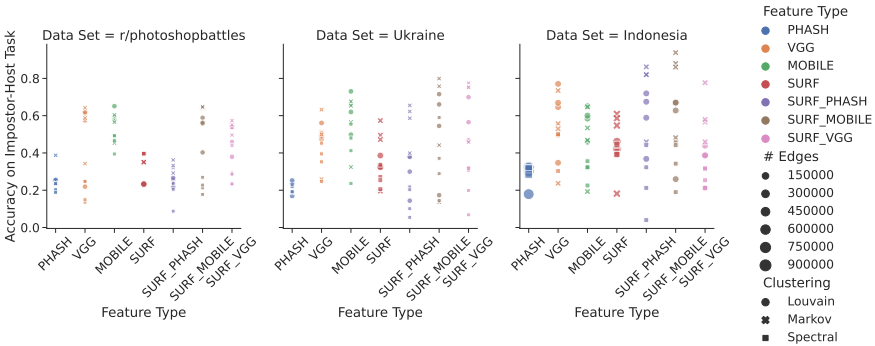
### E.3   Feature Type Plots



**Fig. 16.** The accuracy scores of the Imposter-Host test across the three data sets for each feature type. Each of the three clustering methods is noted with a different shape. The size of each marker is proportional to the number of clusters.



**Fig. 17.** The accuracy scores of the Imposter-Host test across the three data sets for each feature type. Each of the three clustering methods is noted with a different shape. The size of each marker is proportional to the number of components in the graph.

**Fig. 18.** The accuracy scores of the Imposter-Host test across the three data sets for each feature type. Each of the three clustering methods is noted with a different shape. The size of each marker is proportional to the ratio of the number of components in the graph to the number of images in the dataset.



**Fig. 19.** The accuracy scores of the Imposter-Host test across the three data sets for each feature type. Each of the three clustering methods is noted with a different shape. The size of each marker is proportional to the number of edges in the graph.

# F     Telegram Users

Medvezhatko1488, sashakots, russ_orientalist, white_powder2020, karpatska_sich, NSDviz, dadzibao, olifand_rolands, ASupersharij, BerezaJuice, dark_k, joker_ukr, kryuchoktv, legitimniy, notesdetective, rezident_ua, smolii_ukraine, tayni_deputata, thanksrinat, nationalcorps, nedotorkani, ivkolive, ze_konets, ukrnastup, dubinskypro, ruheight, AleksandrSemchenko, botsmanua, borodatayaba, gistapa, kachuratut, poliakovanton, BeregTime, MaksymZhorin, tradition_and_order, KlymenkoTime, sorosata, tsibulya_ua, Ten_NaPleten, donbasscase, lugansk_inside, sorok40russia, ze_landia, zv_kyiv, moh_zdoh, wargonzo, apleonkov, PiB88, format_W, gribvictoria, maksnazar, sheptoon, dobkinmm, UlejUA, spletnicca, razvedinfo, rus_demiurge, LastBP, zlobniaukr, mig41, catars_is, ukrain1an_news, korchynskiy, ua_stalker, project_solaris, liberaxy, orthodox_news, sooproon_bestiary, tasty_flashbacks, fascio_memes, intolerant_historian, Ironvoter, mem_lozha, knpu_division, kekistandivision, EternalMuscovites, nt_orthodox, intolerant_journalist, AD_i_OR, nazbolukr, odindrugqoom, DeepStateUA, ukrnastup, ep867, legion_of_kuchma, NFafaf, History_Q, vidardivision, avantguardia, ulpra, KARAS_EVGEN, GrantDetector, privatnamemarnya, OstanniyCapitalist, afemina, totalopir, intermariumnc, intolerant_warfighter, ukrmemesmineproblemes, evil_ukraine, national_resistance_ua, propala_gramota, postbased, ukrainianintolerant, korchynskiy, Ukrainianintolerantrezerv, RightLit, selo_divisionS, mayonez_sorosa, ubd_ua, national_corp_kyiv, centuriaua