

Transformer-based identification of stochastic information cascades in social networks using text and image similarity

Panagiotis Kasnesis†, Ryan Heartfield*, Xing Liang*, Lazaros Toumanidis†, Georgia Sakellari*, Charalampos Patrikakis†, George Loukas*

†*University of West Attica*, **University of Greenwich*

Abstract

Identifying the origin of information posted on social media and how this may have changed over time can be very helpful to users in determining whether they trust it or not. This currently requires disproportionate effort for the average social media user, who instead has to rely on fact-checkers or other intermediaries to identify information provenance for them. We show that it is possible to disintermediate this process by providing an automated mechanism for determining the information cascade where a post belongs. We employ a transformer-based language model as well as pre-trained ResNet50 model for image similarity, to decide whether two posts are sufficiently similar to belong to the same cascade. By using semantic similarity, as well as image in addition to text, we increase accuracy where there is no explicit diffusion of reshares. In a new dataset of 1,200 news items on Twitter, our approach is able to increase clustering performance above 7% and 4.5% for the validation and test sets respectively over the previous state of the art. Moreover, we employ a probabilistic subsampling mechanism, reducing significantly cascade creation time without affecting the performance of large-scale semantic text analysis and the quality of information cascade generation. We have implemented a prototype that offers this new functionality to the user and have deployed it in our own instance of social media platform Mastodon.

Keywords: Information Cascade, Semantic Textual Similarity, Image Similarity, Deep Learning

1. Introduction

When coming across a new piece of information posted on social media, users may wish to assess its trustworthiness. To do so, they either rely solely on their own knowledge and intuition or take considerable time to check where this information came from in the first place and whether it has been modified since first published. However, investigation on information provenance is not trivial and as such, many social media users will not have the time, motivation or knowledge to conduct it. Instead, they may rely on intermediaries, such as third-party fact-checkers or the social media platforms to do it for them. Even if we assume that these intermediaries are always correct and trustworthy themselves, by the time a false rumour has been fact-checked, it has already spread to a large part of the population. In fact, there is a trade-off between the number of people required to flag a post before it is forwarded for professional assessment versus the number of people exposed to it until it is assessed [1]. At the same time, misinformation travels faster than reliable information (one sixth of the time it took truth to reach 1500 people in [2]), and posts made by individuals or organisations who are experts in a particularly subject or topic (which is going viral) may not necessarily be visible to users due to author/post popularity (e.g., followers, likes, re-shares etc.) [3].

If users themselves were able to identify more easily the provenance of a post’s information at the point of accessing it, they would think twice before resharing it and this would naturally curb the spread of “infodemics”. Here, we take the first steps towards such a provision. Contrary to most existing research in this area, where information cascades are built in a deterministic manner based on explicit resharing (e.g., retweets on Twitter), our approach is stochastic, looking at the degree of similarity between different posts. The little prior work that exists in this area has used statistical word similarity, which however misses posts where the semantics may be the same even if the wording is not. In addition, the previous work has used only textual similarity, while the spreading of news or rumours on social media makes heavy use of images (the average number of reposts with images being estimated to be 11 times larger than those without images [4]). Here, we explore whether incorporating image similarity together with textual similarity can improve the identification of information cascades in social media.

Specifically, this paper introduces the following novel contributions to the body of machine learning techniques for addressing misinformation in social

38 media [5]:

- 39 • A method for monitoring implicit information diffusion and its resulting
40 information cascades over social networks
- 41 • A method for improving clustering performance by combining textual
42 and image similarity detection based on deep learning
- 43 • An efficient post subsampling method to increase the scalability of our
44 approach based on sentence embeddings
- 45 • A prototype tool implementing automated information cascade identi-
46 fication on an existing social media platform

47 2. Related Work

48 2.1. Identifying information cascades

49 Information diffusion has been studied since the beginning of the social
50 media phenomenon as part of the pattern and knowledge discovery dimen-
51 sion of Camacho et al.’s four dimensions of social media analysis [6]. Using
52 explanatory or predictive modelling, the aim is typically to derive latent in-
53 formation about users and communities of users [7]; why information has
54 been diffused in a particular way; where it will be diffused in the future [8]
55 and whether [9, 10] or how [11] it will “go viral” (for marketing [12], polit-
56 ical [13] or other reasons). In terms of provenance of information in social
57 media, most existing research has focused on explicit diffusion, as captured
58 for example through retweets on Twitter and shares on Facebook [14, 15].
59 This kind of provenance is deterministic, as the social media platform itself
60 guarantees the path the information travelled. However, after users come
61 across a post on social media, they may repeat its content without explicitly
62 resharing it word for word. The information is still spreading, yet this cannot
63 be captured by explicit diffusion models.

64 Having utilised post similarity between users’ own posts and their friends’
65 recent posts to reconstruct information cascades, Barbosa et al. [16] reported
66 that at least 11% of interactions are not captured by the explicit reply and
67 retweet/share mechanisms. Taxidou et al. [17] have also shown that limit-
68 ing to explicit resharing cannot capture accurately the influence that a post
69 has had. Instead, they proposed looking at implicit diffusion too, and in

70 their work they suggested reconstructing information cascades using statis-
71 tical word similarity based on TF-IDF (Term Frequency–Inverse Document
72 Frequency). Here, we adopt the same direction of implicit diffusion lead-
73 ing to stochastic information cascades, but we progress beyond statistical
74 similarity to semantic similarity, as different users may describe the same
75 information using very different wording. In addition, the same or very sim-
76 ilar images may be used to describe the same piece of news even if the text
77 appears different. In these cases, considering image similarity in conjunction
78 with semantic text similarity can add context that has not been previously
79 considered in identifying information cascades in social media.

80 *2.2. Transformers in text similarity tasks*

81 For Natural Language Processing (NLP) tasks, such as those gaining
82 increasing attention in social media for analysing information provenance and
83 credibility, Deep Learning (DL) models and in particular Recurrent Neural
84 Networks (RNNs) empowered with Long-Short Term Memory (LSTM), have
85 gained widespread popularity [18] because of their ability to capture the
86 semantics of the words and in consequence generalize over a range of contexts.
87 Recent works use baseline machine learning models such as Latent Dirichlet
88 Allocation (LDA) empowered with word semantics to improve clustering of
89 aspect terms according to their aspect category [19] and topic modeling [20].
90 Support Vector Machines (SVM) have also been used towards this direction
91 by being fed with two dense vectors to determine the degree of semantic
92 similarity between two input sentences. The first one utilizes word-to-word
93 similarity based on Word2Vec embeddings [21] and the latter is built using
94 the word-to-word similarity based on external sources of knowledge [22].

95 However, these DL and baseline NLP architectures have been observed
96 to lack the capability to support inductive transfer learning when it comes
97 to new NLP tasks, because fine-tuning pretrained word embeddings (e.g.
98 Word2Vec [21], Glove [23]) only target a model’s first layer and also be-
99 cause the main task model (e.g., the specific NLP task to be addressed) re-
100 quires training from scratch. In response to this limitation, Language Mod-
101 els (LM) have been proposed [24], which distinguish contextually between
102 similar words and phrases by incorporating the distribution over sequences
103 of words into model weights. Initially, LM architectures were found to lack
104 computational efficiency, since they preclude parallelization, making it a con-
105 straint when it comes to training big sequence lengths. However, recent work
106 based on Transformers-based network architectures [25] have revolutionized

107 NLP problems by replacing the RNNs with Multi-Head Self-Attention (see
108 Subsection 3.1). Transformers rely on an encoder-decoder architecture to
109 extract the meaning from word representations and their relationships, and
110 can be fine-tuned on a wide range of NLP tasks, such as question answering
111 and paraphrase identification, without substantial architecture modifications
112 [26].

113 Bidirectional Encoder Representations from Transformers (BERT) [26] is
114 a LM based on a transformer network [25] designed to pretrain deep bidirec-
115 tional representations from unlabeled text by jointly conditioning on both
116 left and right context in all layers. For pretraining, BERT relies on self-
117 supervised learning and, in particular, has two objectives: a) Masked Lan-
118 guage Modeling (MLM), and b) Next Sentence Prediction (NSP). In MLM,
119 a random sample of the tokens (15% of the input sentence) is removed and
120 replaced with the special token [MASK]. The objective of the model is to
121 predict the masked tokens using a cross-entropy loss function. Regarding
122 NSP, it is a binary classification task that aims at predicting whether two
123 sentences follow each other in the original text, thus negative examples are
124 artificially created by pairing sentences from different documents.

125 Robustly Optimized BERT Pretraining Approach (RoBERTa) [27] is an
126 optimized BERT successor with several modifications to improve the LM
127 pretraining: a) training the model longer, with bigger batches, over more
128 data; b) removing the NSP objective; c) training on longer sequences; and
129 (d) dynamically changing the masking pattern of the MLM. As a result,
130 RoBERTa has managed to surpass BERT’s performance on every NLU task
131 included in GLUE (General Language Understanding Evaluation) bench-
132 mark [28], including Paraphrase Identification (PI) and Semantic Textual
133 Similarity (STS) tasks.

134 Surprisingly, despite their generalizability in several tasks, BERT and
135 RoBERTa do not provide efficient sentence embeddings [29]. Averaging the
136 word embeddings of BERT provides worse latent sentence representations
137 than other models trained on this task, such as Universal Sentence Encoder
138 (USE) [30], a transformer-based network combined with a deep averaging
139 network [31] specifically trained to produce meaningful sentence embeddings.
140 To this end, Sentence-BERT (SBERT) and Sentence-RoBERTa (SRoBERTa)
141 models have been introduced in [29]. They are comprised of two identical
142 networks (e.g., BERT), where each one has a different sentence as input and
143 the objective is to decide whether the two sentences are semantically similar
144 by using cosine similarity as a distance metric, extracting useful embeddings

145 in this way.

146 *2.3. Image in information diffusion tasks*

147 In addition to text, information diffusion in social media has also been
148 studied in relation to images, for predicting the future popularity of a given
149 piece of information [32, 33] or the proliferation of misinformation [34]. For
150 example, Jin et al. [4] have found that images used in disinformation can have
151 distinctive distribution patterns both visually and statistically. McParlane
152 et al. [35] have focused on image popularity prediction by considering visual
153 appearance, content and context. Relevant to our work is Cheng et al.’s
154 work [33] which used image matching to identify copies of the same image
155 and place them into corresponding cascades, but without considering text
156 similarity.

157 More recently, pretrained deep learning models such as VGG16, VGG19,
158 ResNet50, InceptionV3, Xception, InceptionResNetV2 are increasingly adopted
159 to retrieve high level image features [32][36][37][38]. In [36], pre-trained model
160 InceptionResNet V2 was used to derive useful information from photos for
161 popularity prediction in social media. VGG19 was adopted in [37] to extract
162 deep features in addition to extracting basic features including texture and
163 colour of images. Galli et al. [32] have used VGG16 to take sentiment into
164 consideration for social media popularity prediction.

165 In this paper, we propose the use of two deep learning architectures to ex-
166 tract both visual and textual information and fuse them together afterwards
167 to evaluate how similar two posts are. In particular, we collected posts from
168 Twitter to monitor how information spreads in social media by identify-
169 ing diffusion of the posts containing the same or similar content (i.e., text
170 and/or images). This could benefit not only misinformation detection but
171 also various pattern recognition applications such as information retrieval,
172 classification, clustering and change detection.

173 **3. Discovering Probabilistic Information Cascades**

174 In social media, implicit information diffusion processes [17] between
175 posts can manifest over varied conditions based on their content, such as
176 whether a post contains text, an image, video, URL, or any combination
177 of these. If different posts have sufficient similarity between these respective
178 content features, they can be considered the same or slightly different versions
179 of the same information. Here, we focus on discovering information cascades

180 taking into consideration both text and image content similarity. Below, we
181 provide a detailed overview of the algorithms and models used, with exam-
182 ples demonstrating the objectives of these methods for reliably linking posts
183 within associated information cascades, followed by an overview of their in-
184 tegration into a systematic information cascade discovery pipeline.

185 3.1. Text similarity

186 Text similarity deals with determining how similar two pieces of text are.
187 It is considered to be a Natural Language Understanding (NLU) problem
188 that, unlike NLP, deals with machine reading comprehension. Therefore, the
189 objective of text similarity is to identify whether two or more pieces of text
190 represent the same information, albeit with varied use of language, and as
191 such, a trained Artificial Intelligence (AI) model should be able to process
192 natural language in a way that is flexible and not exclusive to a single task,
193 genre or dataset. Typically, in the field of NLP and NLU, this is considered
194 to be an AI-hard problem[28].

195 To develop our text similarity evaluation for information cascade discov-
196 ery, we have chosen RoBERTa_{LARGE} model for the text similarity and feature
197 extraction tasks. RoBERTa follows an encoder-decoder network architecture.
198 The encoder part is composed of a stack of $N = 12$ identical layers, where
199 each of them has two sub-layers connected in a residual manner and followed
200 by layer normalization. The first sub-layer is a Multi-Head Self-Attention
201 mechanism, and the second is a fully connected feed-forward neural network.
202 Residual blocks introduce skip connections are employed around each of the
203 two sub-layers and finally produce embedding outputs of dimension $d_{model} =$
204 1024. The decoder is composed of a stack of $N = 12$ identical layers, but
205 includes a further sub-layer (three in total) to perform Multi-Head Attention
206 over the output of the encoder. Like the encoder, residual connections are
207 used for merging their outputs, followed by layer normalization [25, 27].

208 The efficiency of transformers is mostly based on the Multi-Head Self-
209 Attention mechanism, which defines which parts of a sentence are highly
210 related with each other. In practice, this mechanism makes use of a set of
211 queries Q applied to a set of keys K and provides the most relevant values
212 V . The Self-Attention is given by:

$$A = softmax\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (1)$$

213 where d is the dimensionality of the key vectors used as a scaling factor.
214 The Multi-Head Self-Attention enables the model to attend to several and
215 different representation subspaces at different positions by concatenating the
216 outputs of the heads.

$$MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)W^o \quad (2)$$

217 where h denotes the number of heads, which in the RoBERTa_{LARGE} case are
218 equal to 16. W^o represents the weights of the dense layer that follows the
219 Multi-Head Self-Attention.

220 An advantage of RoBERTa against BERT for text-based information cas-
221 cade identification is its pretrained architecture which benefits from a more
222 diverse range of datasets (larger corpus). For example, its training corpus
223 includes the CommonCrawl News dataset¹ which contains 63 million English
224 news articles and has a larger vocabulary size (50 thousand units) compared
225 to BERT’s (30 thousand units).

226 We used the SRoBERTa large² model pretrained on two NLI datasets,
227 SNLI [39] and MultiNLI [40]. SNLI consists of 570,000 sentence pairs anno-
228 tated with the labels entailment, contradiction, and neutral, while MultiNLI
229 is a collection of 433,000 crowdsourced sentence pairs, containing the same la-
230 bels but covering a range of genres of spoken and written text. SRoBERTa_{LARGE}
231 was trained using a batch size of 16, Adam optimizer with learning rate 2e-5,
232 and a linear learning rate warmup over 10% of the training data [29].

233 The model was retrained and evaluated afterwards on the Semantic Tex-
234 tual Similarity Benchmark (STS-B) dataset [41] reaching a score of 86.39 in
235 Spearman’s rank correlation; it is a collection of sentence pairs, comprised
236 by 7,000 training and 1,400 test samples, drawn from news headlines, video
237 and image captions, and NLI data. The pairs are human-annotated with a
238 similarity score from 1 (lowest) to 5 (highest), while the task is to predict
239 these scores. A model’s performance on this task is evaluated using Pearson
240 and Spearman correlation coefficients, while it should be noted that it is a
241 regression task.

242 In our approach, we exploit the retrained on STS-B model to extract
243 useful text embeddings from the input posts. The extracted embeddings are
244 represented by an array consisting of 1024 float numbers. After acquiring

¹<http://commoncrawl.org/2016/10/newsdataset-available>

²<https://github.com/UKPLab/sentence-transformers>

Table 1: Examples included in the STS-B train set

Examples	Normalized STS score
1: A man is smoking. 2: A man is skating.	0.10
1: Three men are playing chess. 2: Two men are playing chess.	0.52
1: A man is playing the cello. 2: A man seated is playing the cello.	0.85

245 the embeddings of an input posts we apply cosine similarity to identify the
 246 N most similar existing posts and pass them to STS service.

247 In relation to textual similarity and paraphrase identification, we used
 248 two alternative approaches to train our STS model. The RoBERTa_{LARGE} was
 249 trained separately on the Microsoft Research Paraphrase Corpus (MRPC) for
 250 Paraphrase Identification (PI) and on STS-B for STS. The MRPC dataset
 251 [42] is a corpus of sentence pairs (3,700 training and 1,700 test samples)
 252 included in online news sources, annotated by humans to define whether
 253 the sentences in the pair are semantically equivalent; it is imbalanced (68%
 254 positive, 32% negative pairs). Unlike STS-B, MRPC is a dataset handled by
 255 classification algorithms.

256 We trained the model for both datasets using a batch size of 8 and Adam
 257 optimizer with learning rate 1e-5 for 5 epochs, achieving results almost iden-
 258 tical with those reported in [27]. The evaluation of these two models is
 259 presented in Section 4. It should be noted that while the RoBERTa_{LARGE}
 260 MRPC model produces outputs from 0 to 1, the RoBERTa_{LARGE} STS-B
 261 model produces outputs from 1 to 5. So, during the decision process, they
 262 are normalized by dividing by 5. Table 1 presents some examples from the
 263 STS-B training set.

264 3.2. Image similarity

265 Due to context, such as date and occasion, the conditions for assessing
 266 image similarity in information diffusion tend to be stricter than text sim-
 267 ilarity. For example, consider two separate images of a politician taken in
 268 direct point-of-view, standing at the exact same lectern, in the exact same
 269 room, holding a government press conference on television on different days.
 270 In both images, the politician is wearing a suit, in one image blue, and in
 271 the other black. In this case, the images are likely to yield high similarity

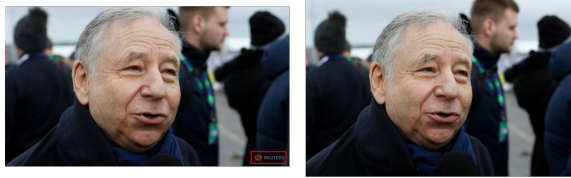
272 with respect to their content, but they should be considered different images
273 and representative of different information contextually. On the contrary,
274 considering the two images of the same nature, where the politician wear-
275 ing the black suit, on the same day, with a news broadcasting logo overlaid
276 on the bottom right of the image, and the other with no news channel logo
277 visible, should be considered the same image and representative of the same
278 information contextually.

279 However, relying on similarity analysis of images alone for reliable infor-
280 mation cascade discovery is naturally prone to false positives, because images
281 related to branding and advertisements (e.g., the “breaking news” image or
282 a company’s logo) are often reused. This may cause the erroneous creation
283 of information cascades between them when there is no real connection be-
284 tween them other than the reuse of a generic image. To address this, it
285 makes sense to combine image similarity with text similarity and deriving a
286 combined similarity metric.

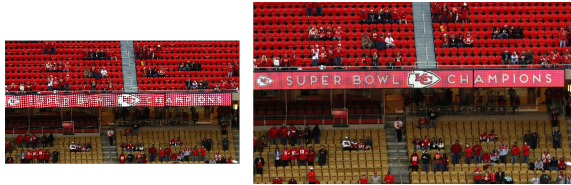
287 To illustrate the requirements of strict image similarity in information
288 cascade generation, in Figure 1 we provide an example of three pairs of
289 social media post images from our TNCD dataset, which are related to the
290 same piece of information. In *Exhibit A*, we can observe that the exact same
291 image has been used between two posts, with a small news logo overlaid
292 on the bottom right the first image, and with the images being a different
293 resolution. In *Exhibit B*, the same image has been used, with the first image
294 being a lower quality, and smaller resolution than the second. Finally, in
295 *Exhibit C*, it is clear these are different images but are related to the same
296 sportsperson, at the same event.

297 To evaluate image similarity in the context of information cascade discov-
298 ery, we adopted existing approaches in image embeddings and metric learn-
299 ing. In image embedding, a robust and discriminative descriptor is learned to
300 represent each image as a compact feature embedding. Typical descriptors
301 include SIFT [43], LBP [44], ORB [45], HOG [46] and Convolutional Neural
302 Network (CNN) embedding’s [47]. In this work we employ feature descrip-
303 tors generated by an existing CNN which employs unsupervised learning to
304 extract latent features, implemented in a Keras pretrained model, as the
305 base for our image feature embeddings generation. For the purposes of com-
306 parison, we have used two CNNs: *ResNet50* [48] (Figure 2) and the Visual
307 Geometry Group (VGG) submission to the ImageNet Challenge [49].

308 The image embeddings are extracted by the deep CNN network, which has
309 multiple layer (M) and n_m neurons in the m^{th} layer ($m= 1,2, \dots M$). For a given



(a) **Exhibit A (0.962 similarity)**: The same image, with different resolution and news broadcaster logo on bottom right of left image



(b) **Exhibit B (0.945 similarity)**: The same image with different resolution and varying image quality



(c) **Exhibit C (0.689 similarity)**: Different images of same sportsman at the same event.

Figure 1: Comparison of image similarity based on strict cascade link requirements

310 image, E_m is the output of the m layer, where $E_m = \sigma(W_m x + b^m)$: W_m is the
 311 projection matrix to be learnt in the m^{th} layer and b^m bias vector; σ is the
 312 non-linear activation function. In each of the CNN networks, a parametric
 313 non-linear function f : image $\rightarrow E_m$ projects an image of D dimensions into a
 314 sub-space of N dimensions in the m^{th} layer. In this sub-space similar images
 315 would be closer to each other and dissimilar images to be further apart.

316 Residual Networks (ResNets) introduce skip connections to skip blocks
 317 of convolutional layers, forming a residual block [48]. These stacked residual
 318 blocks greatly improve training efficiency and largely resolve the vanishing
 319 gradient problem present in deep networks. With a top five accuracy of
 320 93.29%, ResNet50 model won the ImageNet challenge [49] (or ILSVRC),
 321 which is an annual competition using a subset of ImageNet [50] (a large
 322 visual database designed for use in visual object recognition of over 15 million
 323 labelled high-resolution images belonging to roughly 22,000 categories) and

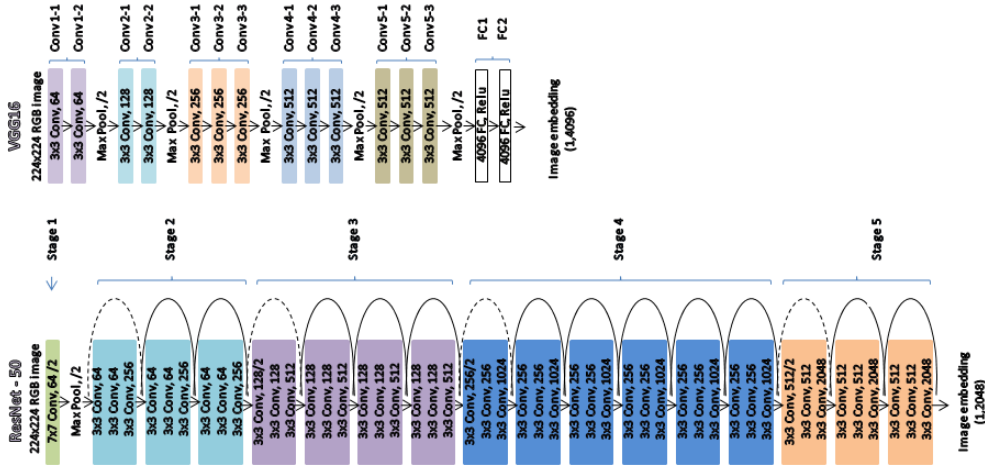


Figure 2: VGG16 and ResNet50 image embeddings

324 is designed to foster the development and benchmarking of state-of-the-art
 325 algorithms. ResNet50 learns a 2048 N dimensional embeddings of an image
 326 from the last layer of stage five (see Figure 2). In contrast, VGG16 has 13
 327 convolutional and 3 Fully Connected (FC) layers, and was employed to learn
 328 a 4096 N dimensional embeddings of an image from FC2 layer. See Figure 2
 329 for a process comparison between ResNet50 and VGG16 image embedding.

330 In metric based learning, a distance metric is utilised to learn from CNN-
 331 embeddings in an latent space to effectively measure the similarity of images.
 332 Considerable efforts have been made to define intuitive image distances in
 333 information retrieval [51, 52, 53, 54], including Cosine similarity, which mea-
 334 sures the similarity between two vectors of an inner product space. It is
 335 measured by the cosine of the angle between two vectors and determines
 336 whether they are pointing in roughly the same direction. It is often used to
 337 measure image similarity as well as document similarity in text analysis (as
 338 in Section 3.1).

339 For each pair of images (I_i, I_j) with image embeddings (E_{m_i}, E_{m_j}) , image
 340 similarity is computed by cosine similarity on image embedding features

341 based on Eq. 3:

$$ImageSimilarity_{(I_i, I_j)} = \frac{\sum_{n=1}^N E_{(mi,n)} * E_{(mj,n)}}{\sqrt{\sum_{n=1}^N E_{(mi,n)}^2} * \sqrt{\sum_{n=1}^N E_{(mj,n)}^2}} \quad (3)$$

342 3.2.1. Information cascade pipeline

343 To identify information cascades in a manner which is practical for real-
344 world deployment, we have developed a pipeline for iterative evaluation of
345 social media posts as they are shared online. When a new post is published,
346 we immediately assess its similarity against all existing posts published up to
347 that point. This is possible by employing an efficient subsampling technique
348 using cosine similarity analysis, which we describe in step 2 of the *Informa-*
349 *tion Cascade Pipeline* below. To demonstrate the utility of the subsampling
350 process, in Figure 7 (see section 4.3) we illustrate how the pipeline cosine sub-
351 sampling latency, combined with RoBERTa_{LARGE} STS latency (i.e., using a
352 fixed subsample of i posts based on the highest cosine scores, as described in
353 step 3 of the pipeline), is capable of processing millions of posts in under 5
354 s using our single computer testbed configuration. The *Information Cascade*
355 *Pipeline* can therefore support information cascade discovery in webscale on-
356 line social media platforms.

357 The *Information Cascade Pipeline* implements the following steps: 1)
358 Extract Feature Embeddings, 2) Subsample Candidate Posts, 3) Semantic
359 Text Similarity, 4) Post link threshold algorithm. In Figure 3, the *Informa-*
360 *tion Cascade Pipeline* is illustrated visually with notation for each processing
361 steps’ algorithmic inputs and outputs (See Table 2 for notations).

362 In step 1, after a new post p is published to a social media platform and
363 stored in the platform database, its text content p_t and image content p_m
364 are extracted to generate post sentence embeddings $(p_{t,f})$ and image sen-
365 tence embeddings $(p_{m,f})$, using our RoBERTa_{LARGE} and ResNet50 models,
366 respectively. Note that the Information Cascade Pipeline is only activated
367 for newly published posts if the post contains at least three words, with or
368 without an image. Where this condition is met, extracted feature embed-
369 dings are stored in a post database alongside existing original post content
370 for future post similarity analysis (i.e., when new posts are published). In
371 step 2, the set of all existing post feature vectors E is queried from the post
372 database and a pairwise comparison of the newly published post text and
373 image feature embeddings (p_t, p_m) is made with each of the existing posts’
374 in E $(e_{t,f}, e_{m,f})$. For each pairwise comparison, for both text $e_{t,f}$ and image

375 $e_{m,f}$ feature vectors, a cosine similarity score is generated with the results
 376 $s_{t,\alpha}$, $s_{m,\alpha}$ added to cosine similarity sets S_t and S_m , respectively. Next, a sub-
 377 set of text T_t and image T_m samples is selected from each cosine similarity
 378 set S_t , S_m , based on the highest respective cosine score, for example, where
 379 $T_t = S_{t,a_i} \cup \{max(S_t \setminus S_{t_i})\}$. In our experiment, for text, we have selected
 380 $i = 8$ as the upper limit of existing posts to forward to semantic text sim-
 381 ilarity analysis, for images as our aim to find the most similar image in all
 382 existing posts, we have used $i = 1$. In step 3, for each $s_t \in T_t$, we compute
 383 the semantic text similarity (STS) score $s_{t,\beta}$ (using our STS-b fine-tuned
 384 RoBERTa_{LARGE} model) for all eight existing post text feature embeddings
 385 in T_t , adding these to the set $T_{t,\beta}$, forwarding the computed STS scores for
 386 cascade link threshold analysis. Step 4 represents the final processing step
 387 where the sets of subsampled STS scores $T_{t,\alpha}$ and image cosine similarity
 388 scores $S_{m,\alpha}$ are assessed by the post link threshold algorithm which evaluates
 389 whether the text and image similarity scores satisfy a predefined threshold
 390 for creating a cascade link. Here, θ_t , θ_m represent the link threshold for se-
 391 mantic text similarity and image cosine similarity, respectively. Based on
 392 our experiments, we have derived optimal θ for text and image similarity
 393 using a gridsearch during the RoBERTa_{LARGE} and ResNet50 fine-tuning
 394 process. In step 4, the algorithm also checks if the existing subsampled post
 395 has a cascade ID $s_c \neq 0$, or not $s_c = 0$ (i.e., where 0 refers to the default
 396 cascade ID for singleton posts that have no cascade association). If the sub-
 397 sampled post’s text and image similarity with the new post is equal to or
 398 above the required similarity threshold the subsampled post’s is checked to
 399 see if it has an existing cascade ID assigned to it. If the the subsampled post
 400 has a cascade ID, the newly published post p_c , linking the newly published
 401 post to corresponding information cascade. Otherwise, a new cascade ID is
 402 created for both the new and subsampled post by selecting the next highest
 403 cascade number in the existing set of cascade IDs C queried from the post
 404 database, where $p_c = 1 + max_{c \in C}$. If no comparison threshold is satisfied, the
 405 newly published post is considered a singleton post and is assigned the de-
 406 fault cascade ID $s_c = 0$. Note that in the case of STS and cosine score ties for
 407 the new post across multiple subsampled posts, time is used as a tiebreaker
 408 to ensure a single link is created for a post in any given information cascade.

409 In Figure 4 an example of the *Information Cascade Pipeline* output is
 410 shown for an identified cascade in our TNCD dataset. Here, the pipeline
 411 shows that it has linked primarily via semantic text similarity, where $\theta_t = 0.5$,
 412 as derived from the gridsearch optimisation, and $s_{t,\beta} \geq \theta_t$). Note that, should

Table 2: List of symbols for Information Cascade Pipeline

Variable	Definition
p_t	Raw text from post
p_m	Raw image from post
$p_{t,f}$	Extracted SRoBERTa _{LARGE} text feature embeddings
$p_{m,f}$	Extracted ResNet50 image feature embeddings
E	Set of existing post text & feature embeddings ($e_{t,f}, e_{m,f}$) & cascade IDs (e_c)
$e_{t,f}$	Text feature embeddings for post $e \in E$
$e_{m,f}$	Image feature embeddings for post $e \in E$
S_t	Set of all text feature embeddings cosine scores $\forall e \in E$
S_m	Set of all image feature embeddings cosine scores $\forall e \in E$
$s_{t,\alpha}$	Text cosine similarity for post $s_t \in S_t$
$s_{m,\alpha}$	Image cosine similarity for post $s_m \in S_m$
T_t	Set of top i cosine scores for text $s_{t,\alpha}$ in set S_t
T_m	Set of top i cosine scores for text $s_{m,\alpha}$ in set S_m
$s_{t,\beta}$	Semantic textual similarity (STS) score for subsampled post $s_t \in S_t$
$T_{t,\beta}$	Subset of text semantic text similarity(STS) scores
s_c	Cascade ID for subsampled post $s \in S$
C	Set of all existing Cascade IDs
θ_t	Link threshold for text similarity
θ_m	Link threshold for image similarity
p_c	Assigned cascade ID for new post p

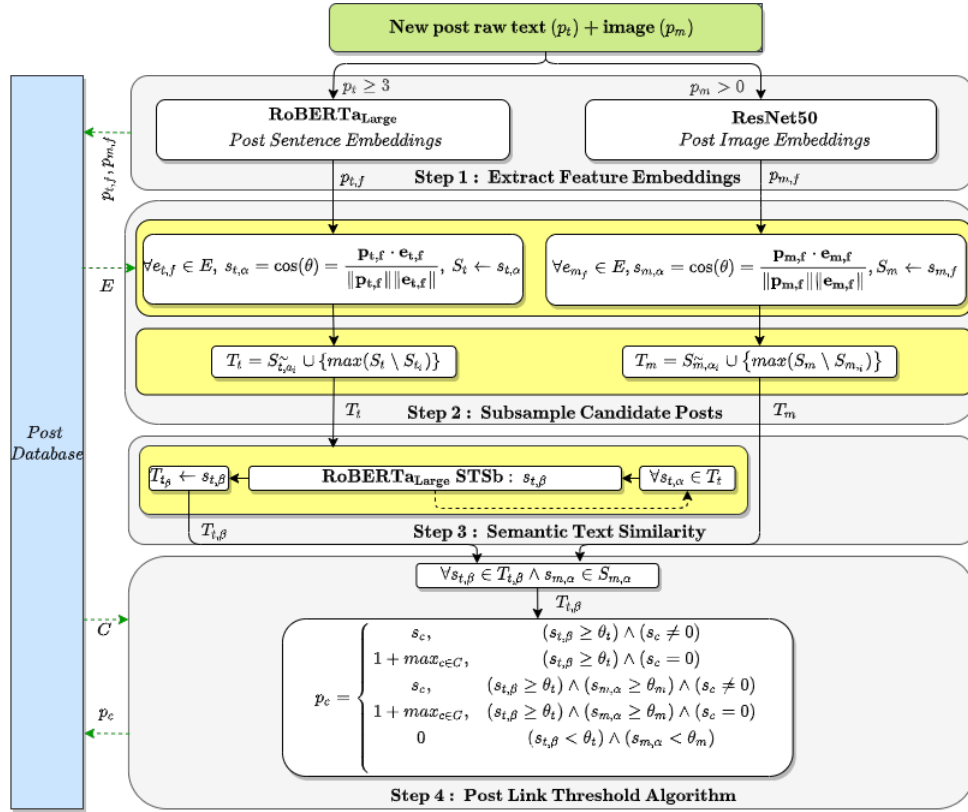


Figure 3: Information Cascade Pipeline

413 $s_{t,\beta} < \theta_t$ (0.5 in this case) for the fourth post in the case, the cascade pipeline
 414 would still have correctly linked the fifth post in the cascade, based on its
 415 image similarity cosine score.

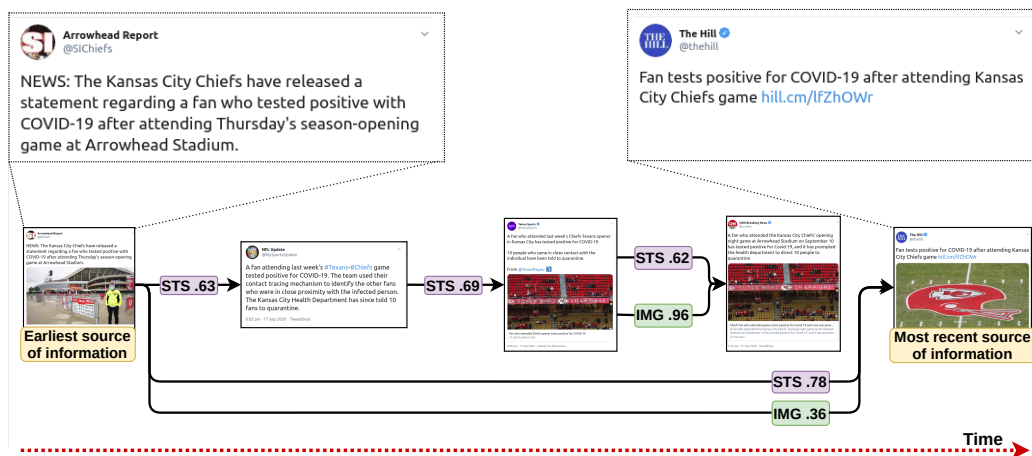


Figure 4: Example of Information Cascade Pipeline output for a identified cascade in the TNCD dataset

416 4. Experimental Analysis and Validation

417 4.1. Experiment methodology and testbed

418 For the experimental analysis of the Information Cascade Pipeline, we
 419 have pre-trained multiple models for text and image similarity, where each
 420 set of models was validated on publicly available datasets optimised for their
 421 respective inference tasks. The experiments were executed on a single com-
 422 puter workstation equipped with a NVIDIA GTX 1080 Ti GPU featuring
 423 11gigabytes RAM, 3584 CUDA cores and a bandwidth of 484GB/s. We
 424 used the Python numpy library for matrix multiplication, Re library for text
 425 preprocessing (i.e., regular expression operations), emoji³ library to convert
 426 emojis into text and Transformers⁴ and Simple Transformers⁵ frameworks
 427 for retraining and evaluating the RoBERTa model. In the case of TF-IDF,

³<https://github.com/carpedm20/emoji>

⁴<https://github.com/huggingface/transformers>

⁵<https://github.com/ThilinaRajapakse/simpletransformers>

Table 3: Details on TNCD

Parameters	Validation set	Test set
no. of posts	600	600
no. of posts in cascade	306	281
no. of cascades	57	60
no. of posts with images	579	599
min no. of posts in a cascade	2	2
max no. of posts in a cascade	12	13

428 we used the NLTK library⁶ to remove English stop words and scikit-learn⁷ to
 429 compute the features. To accelerate the tensor multiplications, we used the
 430 CUDA Toolkit with cuDNN, which is the NVIDIA GPU-accelerated library
 431 for deep neural networks.

432 4.2. TNCD dataset

433 To evaluate the performance of our approach we collected 1,200 news
 434 items posted on Twitter. We call this the Twitter News Cascade Dataset
 435 (TNCD). It contains posts (text and images) retrieved from sportsmen,
 436 politicians and news channels accounts, most from September 2020. We
 437 used the tweepy⁸ library to access the Twitter API. The posts are human-
 438 annotated regarding whether they belong to a particular information dif-
 439 fusion cascade or not. Table 3 presents some of the characteristics of the
 440 created dataset. It is equally split into validation and test set, with each set
 441 containing 600 posts. This was done in order to tune the values of θ_t and θ_m
 442 (see next subsection). It should be noted that all posts contain text but not
 443 all contain images.

444 4.3. Performance evaluation

445 To assess the effectiveness of our *Information Cascade Pipeline* and demon-
 446 strate the usefulness of its hybrid text and image similarity detection model
 447 ensemble (using RoBERTa_{large} for semantic text similarity, finetuned on the
 448 STS-b dataset), we have conducted a comparative analysis of the TNCD

⁶<https://www.nltk.org/>

⁷https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html

⁸<https://www.tweepy.org/>

449 dataset across four different algorithms that could be applied in step 3 of In-
 450 formation Cascade Pipeline. Namely, different pipeline configurations for se-
 451 mantic text analysis which leverage 1) a standard pretrained SRoBERTa_{LARGE}
 452 text similarity model (pretrained on the SNLI and MRPC datasets), 2) a
 453 pretrained RoBERTa_{LARGE} text similarity model fine-tuned for paraphrase
 454 identification classification tasks using the MRPC dataset, 3) a TF-IDF fea-
 455 ture extraction model using cosine similarity (based on work in presented in
 456 [17], and 4) pretrained RoBERTa_{LARGE} text similarity model fined-tuned on
 457 the STS-B dataset. All of the above were evaluated, also by combining them
 458 with the ResNet50 image similarity model, (as well as with VGG16 combined
 459 with pretrained RoBERTa_{LARGE} on the STS-B, for comparison) as part of
 460 the hybrid text and image cascade generation process. Each pipeline config-
 461 uration was evaluated using the “Post Link Threshold Heuristic” defined in
 462 the Information Cascade Monitoring pipeline architecture (see Figure 3).

463 For evaluating each pipeline configuration’s performance, we have selected
 464 the Fowlkes-Mallows index (FMI) [55], which is typically used to determine
 465 the degree of similarity between clusters of data points obtained via a clus-
 466 tering algorithm. Common evaluation metrics such as accuracy and F1-score
 467 used in classification are not applicable to clustering algorithms, or machine
 468 learning approaches which assign a group-based identity to data points, since
 469 their performance evaluation is not as simple as counting the number of false
 470 positives and false negatives, or the precision and recall. This is due to the
 471 fact that the evaluation metric should not consider the exact values of the
 472 cluster labels but rather check whether a cluster is comprised of similar data
 473 according to a set of ground truth labels. The FMI metric provides a suitable
 474 metric for measuring the performance of information cascade generation ac-
 475 cording to the confusion matrix analysis used in our experiment training and
 476 testing results (e.g., True Positive (TP) - post correctly linked to a cascade,
 477 True Negative (TN) - post correctly not added to a cascade, False Positive
 478 (FP) - post incorrectly added to a cascade, False Negative (FN) - post in-
 479 correctly not added to a cascade). This is because information cascades can
 480 be naturally grouped as clusters of interrelated data points. The FMI score
 481 itself is represented in a range from 0 to 1, where the higher the value the
 482 more similar the datapoints within a given information cascade:

$$FMI = \frac{TP}{\sqrt{(TP + FP)(TP + FN)}} \quad (4)$$

483 where TP depicts the true positives, i.e. the number of pairs of posts that
484 belong to the same cascade in both the ground truth labels and the predicted
485 ones), FP the false positives, i.e. the number of pairs of posts that belong to
486 the same cascade in the true labels but not in the predicted labels, and FN
487 the false positives, i.e. the number of pairs of posts that belong in the same
488 cascade in the predicted labels and not in the true labels.

489 During the preprocessing phase, for the case of the transformer-based
490 approaches we removed usernames (e.g., USER) and URLs, while the in-
491 cluded emojis were “deemojified” into text (e.g., :smile). On the other hand,
492 for the TF-IDF approach we removed also the English stop words from the
493 posts’ texts and punctuation before computing the TF-IDF features. Af-
494 terwards, to optimise the selection of text and image similarity threshold
495 parameters θ_t and θ_m in the “Post Link Threshold Heuristic”, we perform
496 a grid-search of their parameters. In Figure 5, each heatmap illustrates
497 the FMI score achieved for different text and image similarity cascade link
498 thresholds parameters across each grid-search iteration. We have excluded
499 the RoBERTa_{LARGE} MRPC model from the best θ_t search since it is trained
500 on a binary classification task, and as a result this threshold is already de-
501 fined to 0.5. The best θ_t was 0.25 for TF-IDF, 0.5 for RoBERTa_{LARGE} STS-B
502 and 0.6 for USE and SRoBERTa_{LARGE}. For all approaches, the optimal θ_m
503 was 0.9.

504 The evaluation of each pipeline semantic text similarity configuration (Ta-
505 ble 4) shows that RoBERTa_{LARGE} fine-tuned on MRPC achieves the lowest
506 performance. This was expected as the MRPC dataset focuses on para-
507 phrase classification rather than semantic text similarity. By comparison,
508 SRoBERTa_{LARGE} pretrained on the NLI and STS-B datasets model achieves
509 a higher FMI score (validation +11.5% and test +9.66%), while the USE
510 model reached even higher FMI scores, 84.06% and 84.03% for the valida-
511 tion and test set respectively. For RoBERTa_{LARGE} fine-tuned on STS-B,
512 the model outperforms RoBERTa_{LARGE} MRPC by a FMI score of over 18%.
513 This improvement in performance is reasonable given the problem defini-
514 tion of information cascade monitoring focuses on the semantic similarity
515 between text (STS-B), and fine-tuning the model further on this dataset op-
516 timises its attention task towards semantic text similarity tasks. Moreover,
517 the RoBERTa_{LARGE} STS-B surpasses by over 7% the performance of TF-
518 IDF-based approach presented in [17] on the validation set and over 4.5%
519 for the case of the test set. Finally, our proposed text (RoBERTa_{LARGE}
520 STS-b) and image (ResNet50) ensemble detection model obtained the high-

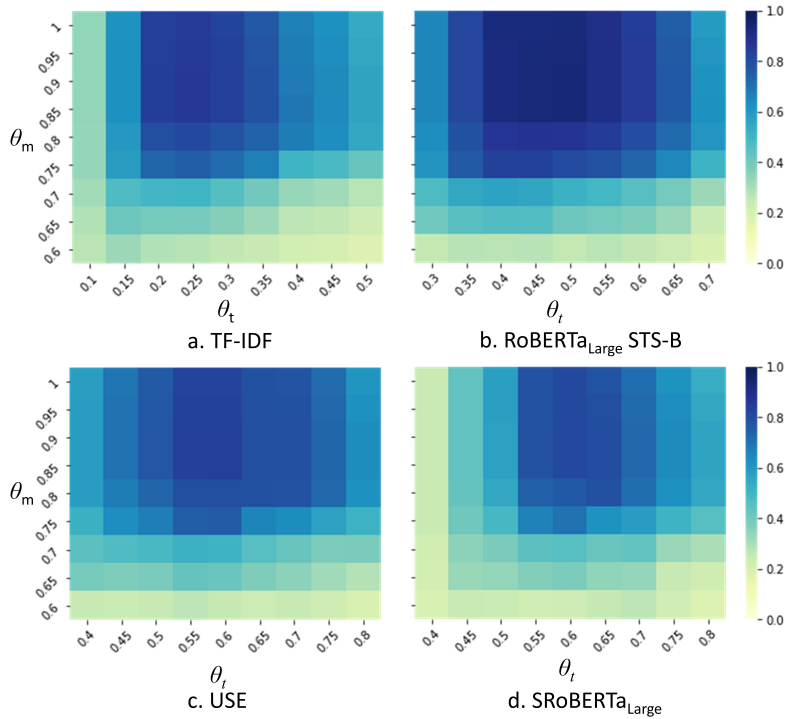


Figure 5: Heatmaps representing the influence of the θ_t , θ_m values to the obtained FMI on the: a. TF-IDF, b. RoBERTa_{LARGE} STS-B, c. USE, and d. SRoBERTa_{LARGE} on the validation set

521 est FMI score and incidentally provided the most accurate configuration for
 522 information cascade monitoring. We observe that including image similar-
 523 ity in the information cascade monitoring process has led to a meaning-
 524 ful performance benefit for all model configurations we have tested and for
 525 RoBERTa_{LARGE} STS-b (validation: from 92.07% to 93.40%, test: 91.55% to
 526 92.02%), which was the best performing model. Furthermore, we examined
 527 also the use of VGG16 embeddings obtaining almost identical scores with
 528 those of ResNet50 (validation: 93.40%, test: 91.96%); however, extracting
 529 embeddings in VGG16 is more computationally expensive (VGG16 has ap-
 530 proximately five times the number of model parameters defined in ResNet50),
 531 which results in significantly increased execution latency (i.e., for 1,000 itera-
 532 tions the inference time per image is 0.117 ms for the VGG16 while only 0.052
 533 ms for the ResNet50). Moreover, it is worth mentioning that by following a
 534 greedy approach (i.e., excluding sentence embedding-based subsampling) we

Table 4: Information cascade discovery performance

Model Integration	Validation FMI	Test FMI
TF-IDF[17] (text)	84.40%	86.80%
TF-IDF[17] (text) + ResNet50 (image)	86.00%	87.50%
USE (text)	82.10%	84.23%
USE (text) + ResNet50 (image)	84.06%	84.93%
RoBERTa _{LARGE} MRPC (text)	66.41%	66.37%
RoBERTa _{LARGE} MRPC (text) + ResNet50 (image)	69.92%	73.34%
SRoBERTa _{LARGE} (text)	81.42%	81.15%
SRoBERTa _{LARGE} (text) + ResNet50 (image)	81.42%	83.00%
RoBERTa _{LARGE} STS-B (text)	92.07%	91.55%
RoBERTa _{LARGE} STS-B (text) + ResNet50 (image)	93.40%	92.02%

535 obtained the same cascade FMI scores for both VGG16 and ResNet50 image
536 similarity models when used in conjunction with our STS-B model. There-
537 fore, for execution latency performance reasons alone, we selected ResNet50
538 as the image similarity deep learning architecture in our Information Cascade
539 Pipeline.

540 To validate the performance of our heuristic algorithm which integrates
541 the combination of text and image similarity for cascade link selection, we
542 have performed an experimental comparison with related research by Sakaki
543 et al., who proposed in [56] an alternative formula for combining text (lin-
544 ear SVM classifier over Bag of Words) and image (Scale-invariant feature
545 transform with SVM) similarity models:

$$Score_{combined} = Score_{text} \times a + Score_{image} \times (1 - a) \quad (5)$$

546 where a is set as a ratio of the text score and an image score to combine two
547 scores appropriately. The authors used a equal to 0.244. However, for the
548 case of our dataset we found out that the best a is 0.95 and the $Score_{combined}$
549 term should be above or equal to 0.45 in order for a post to be included
550 in a cascade. Table 5 presents the obtained results using ResNet50 and
551 RoBERTa_{LARGE} for image and text similarity respectively. Experimental
552 results with our Twitter dataset reported that our heuristic algorithm out-
553 performs the method proposed by Sakaki et al., which reported a validation
554 FMI score of 92.64% and test score of 91.85%, compared to 93.40 and 92.02
555 for our approach, respectively. At the time of writing and to the best of our
556 knowledge, there has been no study other than Sakaki et al.’s exploring the

Table 5: Performance of comparison of text and image integration heuristic algorithms

Integration Algorithm	Validation FMI	Test FMI
Sakaki et al. (2014) [56]	92.64%	91.85%
Our Method*	93.40%	92.02%

*See Figure 3 - Step 4

557 integration between post text and image similarity modelling in social media
 558 information cascade or diffusion analysis.

559 Figure 6 shows a tree-based representation of the information cascades
 560 identified. Here, black links represent TP connections in the cascade, while
 561 the red links represent FP connections in the cascade. As shown, the largest
 562 information cascade presented in our TNCD dataset is correctly identified
 563 to consist of 13 posts. Delving deeper into the predicted FP links (Table 6),
 564 we can observe that some can easily be confused as similar even by human
 565 annotators. The first example presented in Table 6 presents two posts that
 566 talk about the political relationship between the U.S and Iran, with the first
 567 mentioning that the U.N. sanctions against Iran have been restored, while
 568 the second one that they will be reimposed. The posts included in the second
 569 example pair refer both to fatal car accidents, and the street number included
 570 in the first post equals the age of the driver in the second post.

Table 6: False positive examples on the TNCD test set

Examples	STS score
1: The Trump administration has declared that all U.N. sanctions against Iran have been restored, a move most of the rest of the world rejects as illegal. 2: U.S. says U.N. sanctions on Iran to be reimposed Saturday. What does that mean?	0.5386
1: The driver who died heading eastbound in a pickup truck on State Road 40 when the driver of a sport utility vehicle entered a curve and veered into the eastbound lane. 2: The man in his 40s was fatally injured and pronounced dead at the scene	0.5177

571 In line with the previous state of the art [17], we evaluate the performance
 572 of the *Information Cascade Pipeline* with respect to its computation latency

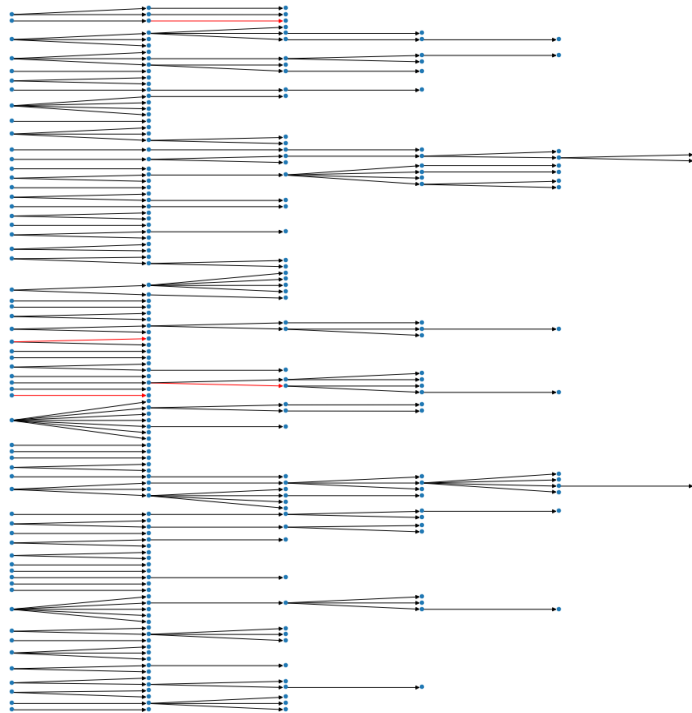


Figure 6: Created cascades in the TNCD test dataset (red links represent the false positive links)

573 when processing newly published posts on a social media platform. As our
 574 objective is to integrate the *Information Cascade Pipeline* into a real tool
 575 for supporting the assessment of information trustworthiness in social media
 576 (Section 5), our analysis takes into consideration the latency for information
 577 cascade analysis of each new post published. Therefore, here, processing
 578 latency represents the total processing time required to assess information
 579 cascade association for every new post. In Figure 7, we first compare the
 580 processing latency of a new post with all existing posts E , for up to 10,000
 581 posts across three methods: 1) bruteforce (greedy) pairwise STS processing
 582 with no subsampling, 2) hierarchical clustering subsampling ([57]) + STS
 583 subset (subset $i = 8$), 3) cosine similarity subsampling mechanism + STS
 584 subset (subset $i = 8$), and 4) TF-IDF estimation followed by cosine similar-
 585 ity. In subfigure 7a, we observe an expected high linear increase in process-
 586 ing latency as the number of stored posts for brutefore comparison increases
 587 (approximately 29 minutes for 10,000 posts), whereas for clustering, cosine

588 comparison (which includes a fixed STS subsample of eight posts) and TF-
589 IDF, the processing latency is orders of magnitude lower and relatively stable
590 as the number of stored posts increases. Subfigure 7b shows that hierarchical
591 clustering also follows a relatively linear processing delay compared to cosine
592 subsampling, albeit with significantly reduced processing time compared to
593 STS bruteforce (approximately 30 s for 10,000 posts). In subfigure 7c, co-
594 sine subsampling takes approximately 4 s to process 1,000,000 posts. The
595 results demonstrate that our *Information Cascade Pipeline* cosine similarity
596 subsampling with a fixed-size STS subset, can support web scale analysis pro-
597 viding lower estimation time than the previous TF-IDF approach [17] above
598 10,000 examples. This is due to the fact that the estimation of TF-IDF index,
599 similarly to that of the cosine similarity, increases as the number of stored
600 posts increase, while the RoBERTa-based STS estimation is applied only to
601 8 posts, and is therefore constant. Subfigure 7c displays, also, the computa-
602 tional expense of including the image processor in our pipeline. Similarly to
603 text similarity, finding similar images is based on applying cosine similarity
604 over ResNet’s embeddings so it is highly dependent on the number of stored
605 posts. As result, the computational cost of including image similarity as
606 well is almost twice as high when compared to ours text-based similarity ap-
607 proach, however, it is still reasonable; it is few milliseconds higher ($\approx 160\text{ms}$)
608 than using only the TF-IDF based text analysis, while having a much higher
609 information cascade discovery performance ($\approx 7\%$). Moreover, it is worth
610 highlighting that the estimation of TF-IDF requires updating the already
611 estimated and stored TF-IDF in the database TF-IDF. By comparison, our
612 method storage of text embeddings is static and does not require continuous
613 updates. Note that this functional behaviour is not reflected in the plots,
614 which display only the estimation times and not the transactions with the
615 database.

616 5. Prototype implementation of the Information Cascade Pipeline 617 mechanism

618 In this section, we provide an overview of our prototype *Information*
619 *Cascade Pipeline* implementation on a private instance of the decentralised
620 social media platform *Mastodon* created for the *EUNOMIA*⁹ project. Figure

⁹<https://eunomia.social>

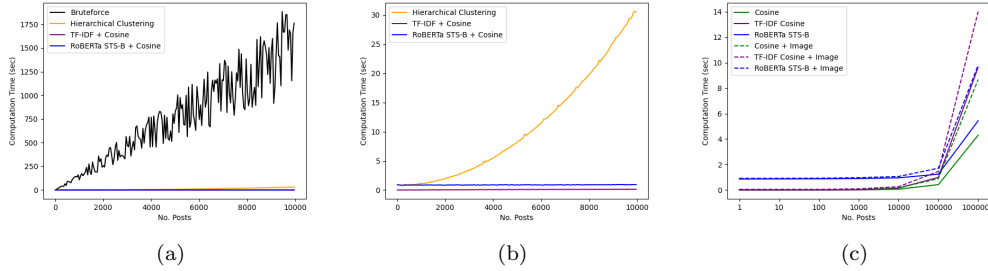


Figure 7: Performance evaluation of Information Cascade Pipeline computation time: a) text similarity only (including brute force); b) text similarity (no brute force); c) text Vs. text+image (no brute force)

621 8 is a high-level illustration of the *Information Cascade Pipeline* integration
 622 within the platform. Specifically, the information cascade monitoring proto-
 623 type is an independent module which interfaces with *EUNOMIA*'s private
 624 Mastodon API to access posts' information, whilst receiving new published
 625 post content via the *EUNOMIA* services orchestrator. Our prototype im-
 626 plements a post analysis component which communicates with the internal
 627 post database, text and image similarity components. Here, the *Information*
 628 *Cascade Pipeline* described in Figure 3 is activated when a published post
 629 meets a predefined minimum word length for cascade processing.

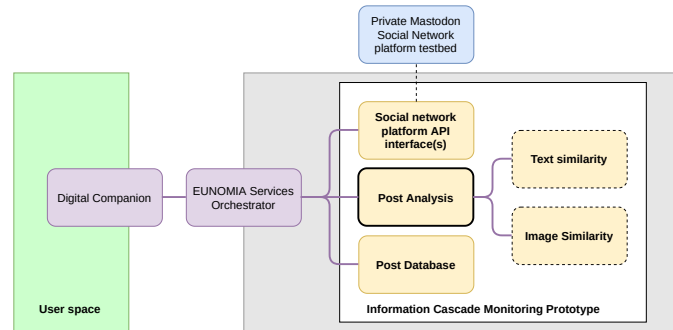


Figure 8: High level overview of the information cascade monitoring prototype within the the EUNOMIA system architecture

630 Figure 9 shows a screenshot of the prototype information cascade user
 631 interface, presented to the user as a side panel that is accessed via the "Show
 632 other similar posts" link shown under each post that belongs to a cascade.
 633 The *Information Cascade Pipeline* has identified an information cascade and

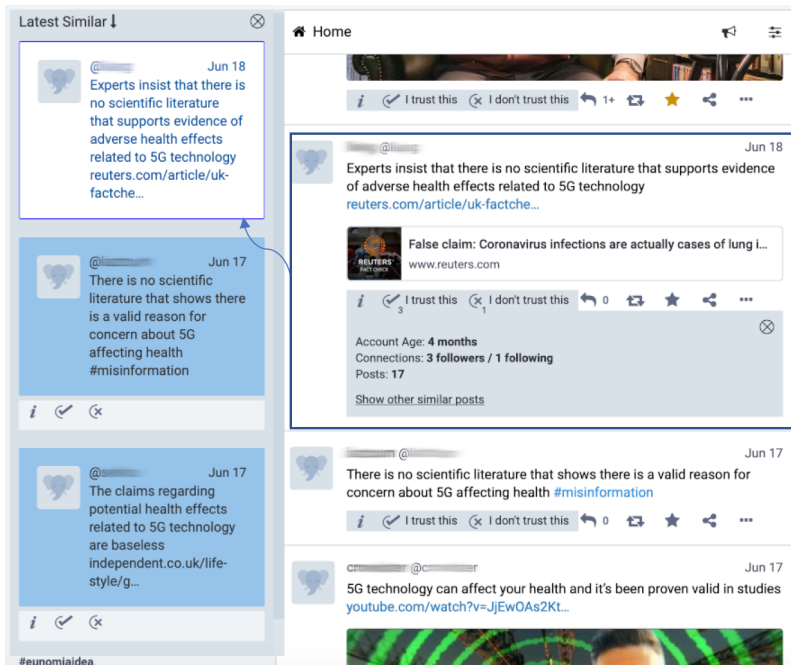


Figure 9: Screenshot of the information cascade as visualised to the EUNOMIA user

634 has ordered it chronologically, highlighting to the user the earliest and most
 635 recent posts in the cascade.

636 6. Conclusions

637 Identifying cases in social media where information has spread or been
 638 replicated by users, without them explicitly resharing it, is a complex task.
 639 Intelligent mechanisms capable of autonomously monitoring the implicit dif-
 640 fusion of information on social media can help analyse the true virality and
 641 spread of information as it propagates in real-time. Importantly, such mech-
 642 anisms can help a user identify the provenance of information and how it
 643 may have changed over time. Here, we progressed beyond the state of the
 644 art in this direction by applying semantic as opposed to statistical similarity,
 645 as well as by incorporating also image similarity. This involved employing a
 646 transformer-based model and a deep Convolutional Neural Network for tex-
 647 tual and image similarity respectively. In addition, our post subsampling
 648 approach was able to make our method applicable to real-world online social
 649 networks. We implemented and deployed our prototype in our own instance

650 of the decentralized social media platform Mastodon. While we have found
651 the prototype to already be practical, it is not able to re-evaluate the mem-
652 bership of posts in existing cascades. In particular, the similarity of orphan
653 posts (not yet included in a cascade) should be re-estimated after a certain
654 time. This would decrease false negatives, but would need to be performed
655 in a manner that is scalable for a real-world social media platform. Also,
656 larger transformer-based architectures could be exploited to increase the per-
657 formance of semantic textual similarity and information fusion mechanisms
658 extracting relational embeddings from text and image pairs, and in this way
659 enable an end-to-end approach. We consider these as interesting directions
660 for future research.

661 **Acknowledgment**

662 The work presented in this paper has been supported through the Eu-
663 ropean Commission’s H2020 Innovation Action programme under project
664 EUNOMIA (Grant agreement No. 825171).

665 **References**

- 666 [1] J. Kim, B. Tabibian, A. Oh, B. Schölkopf, M. Gomez-Rodriguez, Lever-
667 aging the crowd to detect and reduce the spread of fake news and misin-
668 formation, in: Proceedings of Eleventh ACM International Conference
669 on Web Search and Data Mining, ACM, 2018, pp. 324–332.
- 670 [2] S. Vosoughi, D. Roy, S. Aral, The spread of true and false news online,
671 *Science* 359 (2018) 1146–1151.
- 672 [3] T. Porat, P. Garaizar, M. Ferrero, H. Jones, M. Ashworth, M. A. Vadillo,
673 Content and source analysis of popular tweets following a recent case
674 of diphtheria in Spain, *European journal of public health* 29 (2019)
675 117–122.
- 676 [4] Z. Jin, J. Cao, Y. Zhang, J. Zhou, Q. Tian, Novel visual and statistical
677 image features for microblogs news verification, *IEEE Transactions on*
678 *Multimedia* 19 (2017) 598–608.
- 679 [5] M. Choraś, K. Demestichas, A. Giełczyk, Á. Herrero, P. Ksieniewicz,
680 K. Remoundou, D. Urda, M. Woźniak, *Advanced machine learning*

- 681 techniques for fake news (online disinformation) detection: A systematic
682 mapping study, *Applied Soft Computing* 101 (2020).
- 683 [6] D. Camacho, A. Panizo-LLedot, G. Bello-Orgaz, A. Gonzalez-Pardo,
684 E. Cambria, The four dimensions of social network analysis: An
685 overview of research methods, applications, and software tools, *Informa-
686 tion Fusion* 63 (2020) 88–120.
- 687 [7] S. Cavallari, E. Cambria, H. Cai, K. C.-C. Chang, V. W. Zheng, Embed-
688 ding both finite and infinite communities on graphs [application notes],
689 *IEEE Computational Intelligence Magazine* 14 (2019) 39–50.
- 690 [8] M. Li, X. Wang, K. Gao, S. Zhang, A survey on information diffusion
691 in online social networks: Models and methods, *Information* 8 (2017)
692 118.
- 693 [9] K. Subbian, B. A. Prakash, L. Adamic, Detecting large reshare cascades
694 in social networks, in: *Proceedings of 26th International Conference on
695 World Wide Web*, 2017, pp. 597–605.
- 696 [10] W. Xie, F. Zhu, J. Xiao, J. Wang, Social network monitoring for bursty
697 cascade detection, *ACM Transactions on Knowledge Discovery from
698 Data (TKDD)* 12 (2018) 1–24.
- 699 [11] S. S. Singh, A. Kumar, K. Singh, B. Biswas, Lapso-im: A learning-
700 based influence maximization approach for social networks, *Applied
701 Soft Computing* 82 (2019) 105554.
- 702 [12] C. Ashley, T. Tuten, Creative strategies in social media marketing: An
703 exploratory study of branded social content and consumer engagement,
704 *Psychology & Marketing* 32 (2015) 15–27.
- 705 [13] A. Panizo-LLedot, J. Torregrosa, G. Bello-Orgaz, J. Thorburn, D. Ca-
706 macho, Describing alt-right communities and their discourse on twitter
707 during the 2018 us mid-term elections, in: *Proceedings of International
708 conference on complex networks and their applications*, Springer, 2019,
709 pp. 427–439.
- 710 [14] I. Taxidou, P. M. Fischer, Online analysis of information diffusion in
711 twitter, in: *Proceedings of 23rd International Conference on World Wide
712 Web*, 2014, pp. 1313–1318.

- 713 [15] I. Taxidou, T. De Nies, R. Verborgh, P. M. Fischer, E. Mannens,
714 R. Van de Walle, Modeling information diffusion in social media as
715 provenance with w3c prov, in: Proceedings of 24th International Con-
716 ference on World Wide Web, 2015, pp. 819–824.
- 717 [16] S. Barbosa, R. M. Cesar-Jr, D. Cosley, Using text similarity to detect
718 social interactions not captured by formal reply mechanisms, in: Pro-
719 ceedings of 11th International Conference on e-Science, IEEE, 2015, pp.
720 36–46.
- 721 [17] I. Taxidou, S. Lieber, P. M. Fischer, T. De Nies, R. Verborgh, Web-scale
722 provenance reconstruction of implicit information diffusion on social me-
723 dia, *Distributed and Parallel Databases* 36 (2018) 47–79.
- 724 [18] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey,
725 M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah,
726 M. Johnson, X. Liu, L. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa,
727 K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa,
728 A. Rudnick, O. Vinyals, G. S. Corrado, M. Hughes, J. Dean, Google’s
729 neural machine translation system: Bridging the gap between human
730 and machine translation, *ArXiv abs/1609.08144* (2016).
- 731 [19] S. Poria, I. Chaturvedi, E. Cambria, F. Bisio, Sentic lda: Improv-
732 ing on lda with semantic similarity for aspect-based sentiment analy-
733 sis, 2016 International Joint Conference on Neural Networks (IJCNN)
734 (2016) 4465–4473.
- 735 [20] E. Ekinici, S. I. Omurca, NET-LDA: a novel topic modeling method
736 based on semantic document similarity, *Turkish Journal of Electrical*
737 *Engineering and Computer Sciences* 28 (2020) 2244–2260.
- 738 [21] T. Mikolov, K. Chen, G. S. Corrado, J. Dean, Efficient estimation of
739 word representations in vector space, *CoRR abs/1301.3781* (2013).
- 740 [22] H. T. Nguyen, P. H. Duong, E. Cambria, Learning short-text seman-
741 tic similarity with word embeddings and external knowledge sources,
742 *Knowl. Based Syst.* 182 (2019).
- 743 [23] J. Pennington, R. Socher, C. D. Manning, Glove: Global vectors for
744 word representation, in: Proceedings of Conference on Empirical Meth-
745 ods in Natural Language Processing, 2014, pp. 1532–1543.

- 746 [24] A. M. Dai, Q. V. Le, Semi-supervised sequence learning, in: Proceedings
747 of Advances in Neural Information Processing Systems 28, 2015.
- 748 [25] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez,
749 L. Kaiser, I. Polosukhin, Attention is all you need, in: Proceedings of
750 Annual Conference on Neural Information Processing Systems, 2017.
- 751 [26] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of
752 deep bidirectional transformers for language understanding, in: Pro-
753 ceedings of Annual Conference of the North American Chapter of the
754 Association for Computational Linguistics, 2019.
- 755 [27] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis,
756 L. Zettlemoyer, V. Stoyanov, Roberta: A robustly optimized bert
757 pretraining approach, in: Proceedings of International Conference on
758 Learning Representations, 2019.
- 759 [28] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, S. R. Bowman, Glue:
760 A multi-task benchmark and analysis platform for natural language un-
761 derstanding, in: Proceedings of EMNLP Workshop BlackboxNLP: An-
762 analyzing and Interpreting Neural Networks for NLP, 2018, pp. 353–355.
- 763 [29] N. Reimers, I. Gurevych, Sentence-bert: Sentence embeddings using
764 siamese bert-networks, in: Proceedings of Conference on Empirical
765 Methods in Natural Language Processing and the 9th International Joint
766 Conference on Natural Language Processing, 2019.
- 767 [30] D. M. Cer, Y. Yang, S. yi Kong, N. Hua, N. Limtiaco, R. S.
768 John, N. Constant, M. Guajardo-Cespedes, S. Yuan, C. Tar, Y.-H.
769 Sung, B. Strope, R. Kurzweil, Universal sentence encoder, ArXiv
770 abs/1803.11175 (2018).
- 771 [31] M. Iyyer, V. Manjunatha, J. L. Boyd-Graber, H. Daumé, Deep un-
772 ordered composition rivals syntactic methods for text classification, in:
773 Proceedings of ACL, 2015.
- 774 [32] F. Gelli, T. Uricchio, M. Bertini, A. D. Bimbo, Image popularity predic-
775 tion in social media using sentiment and context features, in: Proceed-
776 ings of 23rd ACM international conference, IEEE, 2015, pp. 907–910.

- 777 [33] J. Cheng, L. A. Adamic, P. A. Dow, J. Kleinberg, J. Leskovec, Can
778 cascades be predicted?, in: Proceedings of 23rd international conference
779 on World wide web, 2014, pp. 925–936.
- 780 [34] P. Qi, J. Cao, T. Yang, J. Guo, J. Li, Exploiting multi-domain visual
781 information for fake news detection, in: Proceedings of IEEE Interna-
782 tional Conference on Data Mining, 2019, pp. 518–527.
- 783 [35] P. J. McParlane, Y. Moshfeghi, J. M. Jose, Nobody comes here anymore,
784 it’s too crowded; predicting image popularity on flickr, in: Proceedings
785 of International Conference on Multimedia Retrieval, 2014.
- 786 [36] M. Meghawat, S. Yadav, D. Mahata, Y. Yin, R. R. Shah, R. Zimmer-
787 mann, A multimodal approach to predict social media popularity, in:
788 Proceedings of IEEE Conference on Multimedia Information Processing
789 and Retrieval, 2018.
- 790 [37] J. Lv, W. Liu, M. Zhang, H. Gong, B. Wu, H. Ma, Multi-feature fusion
791 for predicting social media popularity, in: Proceedings of 25th ACM
792 international conference on Multimedia, 2017.
- 793 [38] A. Canziani, A. Paszke, E. Culurciello, An analysis of deep neural
794 network models for practical applications, arXiv preprint (2017).
- 795 [39] S. R. Bowman, G. Angeli, C. Potts, C. D. Manning, A large anno-
796 tated corpus for learning natural language inference, in: Proceedings
797 of Conference on Empirical Methods in Natural Language Processing,
798 Association for Computational Linguistics, Lisbon, Portugal, 2015, pp.
799 632–642. doi:10.18653/v1/D15-1075.
- 800 [40] A. Williams, N. Nangia, S. Bowman, A broad-coverage challenge cor-
801 pus for sentence understanding through inference, in: Proceedings of
802 Conference of the North American Chapter of the Association for Com-
803 putational Linguistics: Human Language Technologies, Association for
804 Computational Linguistics, New Orleans, Louisiana, 2018, pp. 1112–
805 1122. doi:10.18653/v1/N18-1101.
- 806 [41] D. M. Cer, M. T. Diab, E. Agirre, I. Lopez-Gazpio, L. Specia, Semeval-
807 2017 task 1: Semantic textual similarity multilingual and crosslingual
808 focused evaluation, ArXiv abs/1708.00055 (2017).

- 809 [42] W. B. Dolan, C. Brockett, Automatically constructing a corpus of sen-
810 tential paraphrases, in: Proceedings of Third International Workshop
811 on Paraphrasing, 2005.
- 812 [43] D. G. Lowe, Distinctive image features from scale-invariant keypoints,
813 International journal of computer vision 60 (2004) 91–110.
- 814 [44] T. Ahonen, A. Hadid, M. Pietikainen, Face description with local binary
815 patterns: Application to face recognition, IEEE transactions on pattern
816 analysis and machine intelligence 28 (2006) 2037–2041.
- 817 [45] E. Karami, S. Prasad, M. Shehata, Image matching using sift, surf,
818 brief and orb: Performance comparison for distorted images, in: Pro-
819 ceedings of Newfoundland Electrical and Computer Engineering Confer-
820 ence, 2015.
- 821 [46] N. Dalal, B. Triggs, Histograms of oriented gradients for human detec-
822 tion, in: Proceedings of IEEE Computer Society conference on computer
823 vision and pattern recognition, IEEE, 2005, pp. 886–893.
- 824 [47] S. Appalaraju, V. Chaoji, Image similarity using deep cnn and curricu-
825 lum learning, arXiv preprint abs/1709.08761 (2017).
- 826 [48] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recog-
827 nition, in: Proceedings of Computer Vision and Pattern Recognition,
828 2016, pp. 770–778.
- 829 [49] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-
830 scale image recognition, in: Proceedings of International Conference on
831 Learning Representations, 2015.
- 832 [50] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with
833 deep convolutional neural networks, Advances in neural information
834 processing systems 25 (2012) 1097–1105.
- 835 [51] G. Palubinskas, Mystery behind similarity measures MSE and SSIM,
836 in: Proceedings of IEEE International Conference on Image Processing,
837 IEEE, 2014, pp. 575–579.
- 838 [52] L. Wang, Y. Zhang, J. Feng, On the euclidean distance of images, IEEE
839 Transactions on Pattern Analysis and Machine Intelligence 27 (2005)
840 1334–1339.

- 841 [53] D. Huttenlocher, G. Klanderman, W. Rucklidge, Comparing images
842 using the hausdorff distance, *IEEE Transactions on Pattern Analysis*
843 *and Machine Intelligence* 15 (1993) 850–863.
- 844 [54] D. Sejal, T. Ganeshsingh, K. R. Venugopal, S. S. Iyengar, L. M. Patnaik,
845 Acsir: Anova cosine similarity image recommendation in vertical search,
846 *International Journal of Multimedia Information Retrieval* 6 (2017) 143–
847 154.
- 848 [55] E. B. Fowlkes, C. L. Mallows, A method for comparing two hierarchical
849 clusterings, *Journal of the American statistical association* 78 (1983)
850 553–569.
- 851 [56] S. Sakaki, Y. Miura, X. Ma, K. Hattori, T. Ohkuma, Twitter user
852 gender inference using combined analysis of text and image processing,
853 in: *Proceedings of VL@COLING*, 2014.
- 854 [57] P. Kasnesis, R. Heartfield, L. Toumanidis, X. Liang, G. Loukas, C. Pa-
855 trikakis, A prototype deep learning paraphrase identification service for
856 discovering information cascades in social networks, in: *Proceedings*
857 *of IEEE International Conference on Multimedia and Expo Workshops*,
858 2020.