# Ranking Micro-Influencers: a Novel Multi-Task Learning and Interpretable Framework

Adam Elwood
*lastminute.com*
Chiasso, Switzerland
adam.elwood@lastminute.com

Alberto Gasparin
*lastminute.com*
Chiasso, Switzerland
alberto.gasparin@lastminute.com

Alessandro Rozza
*lastminute.com*
Chiasso, Switzerland
alessandro.rozza@lastminute.com

*Abstract*—With the rise in use of social media to promote branded products, the demand for effective influencer marketing has increased. Brands are looking for improved ways to identify valuable influencers among a vast catalogue; this is even more challenging with "micro-influencers", which are more affordable than mainstream ones but difficult to discover. In this paper, we propose a novel multi-task learning framework to improve the state of the art in micro-influencer ranking based on multimedia content. Moreover, since the visual congruence between a brand and influencer has been shown to be good measure of compatibility, we provide an effective visual method for interpreting our models' decisions, which can also be used to inform brands' media strategies. We compare with the current state-of-the-art on a recently constructed public dataset and we show significant improvement both in terms of accuracy and model complexity. The techniques for ranking and interpretation presented in this work can be generalised to arbitrary multimedia ranking tasks that have datasets with a similar structure.

*Index Terms*—Influencer Marketing, Neural Networks, Multi-Task Learning, Multimodal, Interpretability

## I. INTRODUCTION

As people spend increasingly more time on the internet, brands are keen to make use of online marketing platforms [1]. When connected to search and social media, these platforms are able to personalise adverts, increasing the probability that consumers will interact with them [2]. This has been effective enough for many of the websites on the internet to be mainly funded through advertising revenues. However, as time has progressed, consumers have become desensitised to online advertisements, increasing their use of ad-blocking software and their suspicion of online adverts [3]–[5].

To find a solution to this issue, brands have started forming partnerships with "influencers", social media accounts that have a significant and dedicated following [6], [7]. Unlike celebrity sponsorship, influencer marketing can have the advantage of costing less and being able to target very specific demographics [4]. Influencers also typically have an intimate relationship with their audiences, making them a more trusted source of information and more effective advertisers. With the continuous rise in the use of social media platforms, this new form of marketing has had significant success [8].

Despite the promise of this new style of marketing, finding appropriate influencers to advertise their product is a major challenge for brands that want to engage in it [9]. In contrast to the relatively limited pool of online marketing platforms

and celebrities, there are many more individuals that can be defined as influencers [10]. Particularly enticing are "micro-influencers", those with between 5000 and 100 000 social media followers [11]. Due to their relative obscurity, their advertising services can be obtained cheaply, while being very effective in the case that their audience is a good match for a brand. At this point there are hundreds of thousands of micro-influencers, so the challenge is sorting through them and finding the most relevant candidates for a particular brand.

This problem has been made tractable with recent advances in machine learning for computer vision and natural language processing, along with the large quantity of data produced by the social media platforms. One promising technique for ranking micro-influencers based on such techniques was introduced by Gan *et al.* [12]. Influencers are ranked for brands, based on the similarity of the images and text in influencer and brand posts, which are pooled and embedded via a custom neural network. This methodology is backed up with further research into influencer marketing, which has helped to support the hypothesis that the visual congruence between a brand and influencer is a good measure of compatibility [13]. Other relevant literature is covered in Section II.

In this work we present significant improvements on the algorithm presented by Gan *et al.* [12] through the introduction of novel neural network architectures, described in Section III.

To develop these methods, we make use of a dataset of brands and influencers taken from Instagram, originally introduced in [12] and described further in Section IV-A. The details of the experiments run to demonstrate the efficacy of our methods can be found in Section IV.

We also introduce novel methods for the interpretability of influencer ranking models in Section V. This allows the techniques introduced to provide brands with tools for better understanding the impact of their media content, along with the potential for ranking micro-influencers at scale. Recent research in the interpretability of neural networks can be used to understand the exact components of influencer visual content that make them suitable for a certain brand. This allows brands to better tailor their media content to suit specific audiences.

Finally, we summarise our findings and suggest avenues for future work in Section VI.

The main contributions of this work are as follows:

- Significant improvements over the state-of-the-art of micro-influencer ranking across major metrics through the addition of multi-task learning to the technique presented in [12].
- Novel neural network architectures for ranking, based around a trainable inner product, improving performance with a significant reduction in the number of model parameters.
- The introduction of techniques for interpreting the ranking results to provide brands with tools for improving their media content.

## II. RELATED WORK

### A. Influencer marketing

As influencer marketing has grown in popularity, there has been significant interest from the traditional marketing research community [14]–[18]. This has led to it becoming an established part of the marketing strategy of many big brands. Due to the challenges of finding appropriate influencers detailed above, there has also been an increase in interest in using deep learning algorithms to identify influencers and analyse their interactions with their audience [13].

The most directly relevant work on influencer ranking with deep learning, and the inspiration for this work, comes from Gan *et al.* [12], which includes a summary of the relevant learning to rank techniques that have been applied here [19]–[22]. Additionally, Aleksandr *et al.* [23] presented a technical demonstration of an influencer discovery marketplace, SoMin, although did not provide any concrete methods. On top of this, Gelli *et al.* proposed a framework to predict the popularity of social images [24], along with a learn to rank technique designed to help brands find relevant media content [25].

### B. Multi-task feature learning

In multi-task learning, deep neural networks are trained to be able to perform well simultaneously in multiple different, yet related, domains [26]–[28]. Provided there is a latent feature space that is relevant to all the tasks, multi-task learning can help networks to learn more robust features during the optimisation process. This often leads to improved performance in the main task, despite good performance in auxiliary tasks not being strictly necessary.

There are several different approaches to multi-task learning [27], but the most relevant to our use case is high-level feature learning. In early versions of this approach, the learning of different tasks is decoupled by learning the feature covariance for all the tasks, obtaining task specific hidden representations within a shallow network [29], [30]. In the deep learning setting, a common approach is for the different tasks to share several of the first hidden network layers, with only one or two dense layers used for task-specific parameters [31]–[34]. Due to its general purpose applicability, multi-task learning can be useful when training many different deep learning architectures, from image processing convolutional networks [35] to graph neural networks [36].

### C. Deep neural network interpretability

Deep neural networks are very powerful when applied to computer vision and natural language applications, but act as black-boxes by their nature. This makes it difficult to interpret why they make decisions relevant to the task they are trained to perform on. This has led to increased interest in finding ways to interpret why the networks make the decisions they do [37], [38]. Many of the techniques for interpretation are best applied to individual examples, for example providing insight into which parts of an image were most relevant to a convolutional neural network deciding on a certain classification. One of the most useful techniques for computer vision applications is Grad-CAM [39], which uses the gradients of a target neural network component flowing into the final convolutional layer to highlight the most important regions of the original image. This kind of interpretation is valuable in the influencer marketing use case, as it allows brands to understand the visual content that led to them being matched to certain influencers. Along with checking that appropriate visual content is being used, this information can also be valuable to help the brand design future marketing strategies.

## III. METHOD

### A. Problem formulation

Given a set of $n$ brands $\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2, \ldots, \mathbf{b}_n\}$, with $\mathbf{b}_j \in \mathbb{R}^{(d_v + d_t)}$ and a set of $m$ influencers $\mathcal{I} = \{\mathbf{i}_1, \mathbf{i}_2, \ldots, \mathbf{i}_m\}$ with $\mathbf{i}_k \in \mathbb{R}^{(d_v + d_t)}$, both coming with visual ($d_v$) and textual features ($d_t$), our goal is to recommend a given brand $\mathbf{b}$ with a list of candidate micro-influencers. We assume that for each brand $\mathbf{b}_*$ a list of associated influencers is available. In the following, we'll refer to these as positive examples for the given brand and describe the set as $\mathcal{I}_+(\mathbf{b}_*)$. In similar fashion, we'll refer to all non-associated micro-influencers, for a given brand $\mathbf{b}_*$, as negative examples, $\mathcal{I}_-(\mathbf{b}_*)$. Each brand and micro-influencer comes with both textual and visual information for each of its posts. We represent a brand $\mathbf{b} = [\mathbf{b}_t, \mathbf{b}_v]$, where $\mathbf{b}_t \in \mathbb{R}^{d_t}$ is a pooled summary of the textual features, obtained by a neural network embedding of word tokens via the spaCy library (based on Tok2Vec) [40], while $\mathbf{b}_v \in \mathbb{R}^{d_v}$ is a pooled summary of the visual features obtained from a pretrained VGG-16 [41] used as a feature extractor. More formally, given the last $N_p$ posts by a given brand, a pooling operation $p : \mathbb{R}^{N_p \times d} \mapsto \mathbb{R}^d$ is applied on the image (text) embeddings in order to obtain a unique account representation. In this work, the pooling operation is a simple average over the posts dimension. The preprocessing steps described above holds both for brand and micro-influencer accounts, and are deeply inspired by [12].

Given a brand and a set of micro-influencers $(\mathbf{b}, \{\mathbf{i}_1, \ldots, \mathbf{i}_k\})$ the objective is to provide a ranking (induced by a scoring function) of these micro-influencer w.r.t. the given brand. In the remaining sections we will describe two approaches to assign a score to each influencer, given a brand, and induce the previously described ranking.
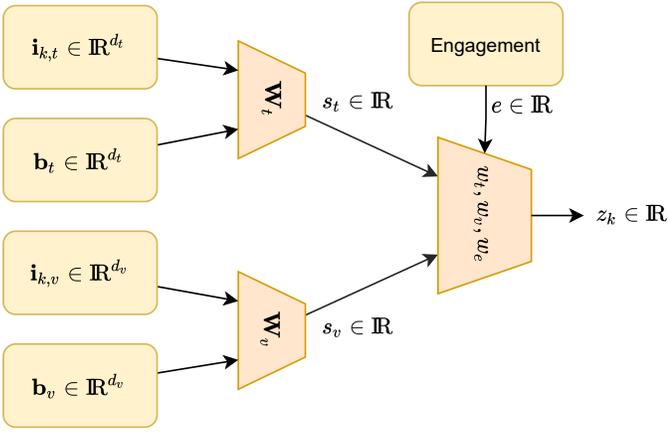
Fig. 1. WSim architecture with its parameters $\boldsymbol{\Theta} = [\mathbf{W}_t, \mathbf{W}_v, w_t, w_v, w_e]$. The input to the model is the pooled summary of the textual and visual features (obtained via VGG-16 and spaCy) for both brand and influencers.

### B. Bilinear Similarity

The first model we propose is a straightforward generalization of the cosine similarity model. In particular, given a pair of brand and micro-influencer ($\mathbf{b} \in \mathbb{R}^{(d_v+d_t)}, \mathbf{i}_k \in \mathbb{R}^{(d_v+d_t)}$) we compute two similarity scores for the visual and textual features separately:

$$s_t = \mathbf{b}_t^T \mathbf{W}_t \mathbf{i}_{k,t} \qquad (1)$$
$$s_v = \mathbf{b}_v^T \mathbf{W}_v \mathbf{i}_{k,v} \qquad (2)$$

with $\mathbf{W}_t \in \mathbb{R}^{d_t \times d_t}, \mathbf{W}_v \in \mathbb{R}^{d_v \times d_v}$ being diagonal matrices whose weights are the free-parameters of the model. It is worth noting that this bilinear similarity model generalizes the dot product, allowing one to compute a similarity score which is based on the minimization (maximization) of a principled objective. The text and visual scores are then combined with a precomputed engagement score $e$ that takes into account the popularity of the posts of a given brand (micro-influencer) as in [12]. A convex combination of the three scores is considered to compute the final score $z_k$ of a micro-influencer $\mathbf{i}_k$ given a brand $\mathbf{b}$:

$$z_k = w_t s_t + w_v s_v + w_e e \quad \text{s.t.} \ w_t + w_v + w_e = 1 \qquad (3)$$

The model's parameter $\boldsymbol{\Theta} = [\mathbf{W}_t, \mathbf{W}_v, w_t, w_v, w_e]$ are learned by minimizing the ranking objective described in Section III-D; for further details on the optimization strategy we refer the reader to Section IV-B. In the remaining of the paper we will refer to this model as **WSim** and its architecture is depicted in Figure 1.

### C. Multi-task Learning

The previously presented model builds on top of the feature representations of general-purpose pretrained models such as VGG-16 and the language model provided by spaCy. In order to further exploit the information available within the dataset and to enhance the learned representation for both influencers and brands, we first feed the pooled text embeddings to a neural network $f_\theta$ and the pooled visual embeddings to another

neural network $g_\phi$. Note that the same networks are used for brand and influencers to learn a common latent space. This parameter sharing is highlighted in the model's architecture depicted in Figure 2. In principle, $f_\theta$ and $g_\phi$ have no predefined functional form and all kind of models can be used. In this work, we employed two fully connected neural networks with two hidden layers each. More details on the architecture used in our experiments are provided in Section IV-B.

The new embeddings for each influencer and brand are obtained by multiplying the visual and text representation obtained through $f_\theta$ and $g_\phi$, which are constrained to have the same output dimension, $d_r$:

$$\mathbf{e}_b = f_\theta(\mathbf{b}) \odot g_\phi(\mathbf{b}) \qquad (4)$$
$$\mathbf{e}_{i,k} = f_\theta(\mathbf{i}_k) \odot g_\phi(\mathbf{i}_k) \qquad (5)$$

where $\odot$ is the Hadamard product. These new embeddings are then fed through a bilinear similarity layer, with a learnable parameter matrix $\mathbf{W}_r \in \mathbb{R}^{d_r \times d_r}$, and a score $z(k = 1, 2, \dots)$ is computed for each brand and micro-influencer pair as follows:

$$z_k = (1 - w_e)(\mathbf{e}_b^T \mathbf{W}_r \mathbf{e}_{i,k}) + w_e e \qquad (6)$$

where $w_e$ is a trainable weight. Similarly to what has been described in the previous section, a ranking loss is used and will be referred as $\mathcal{L}_{main}$ hereafter. Along with the standard ranking task, an auxiliary task has been defined in order to allow the network to learn even better representations. In particular, we introduce a classification loss $\mathcal{L}_{ce}$ (which corresponds to the standard cross-entropy loss) to predict, given $\mathbf{e}_b$ or $\mathbf{e}_{i,k}$ to which macro-category a brand or influencer belong through a shallow neural network $h_\psi$. The overall training loss can be formalized as:

$$\mathcal{L}(\cdot) = \mathcal{L}_{main}(\cdot) + \lambda \mathcal{L}_{ce,brand}(\cdot) + \gamma \mathcal{L}_{ce,infl}(\cdot) \qquad (7)$$

with $\lambda$ and $\gamma$ two hyperparameters that controls the intensity of the auxiliary tasks for the brand and the micro-influencer being evaluated. The macro-categories are detailed in Section IV-A and are one of the key feature available in the dataset used for the experiments. In the remaining of this work we'll refer to this model as **WSim-MT** due to the obvious similarities between it and the model described in the previous section. Indeed, they both heavily rely on a learnable bilinear similarity within the model and the final score for the given micro-influencer and brand is computed in almost the same way (see Equation 3 and 6).

### D. Ranking Loss

Learning to rank is a well known problem and several loss functions that address it have been defined over the years [19]. In this work, given a brand and a pool of $K$ micro-influencers we compute the top one probability of each candidate as in [21]:

$$\mathbf{p} = \text{softmax}(\mathbf{z}); \qquad \mathbf{z} \in \mathbb{R}^K \qquad (8)$$
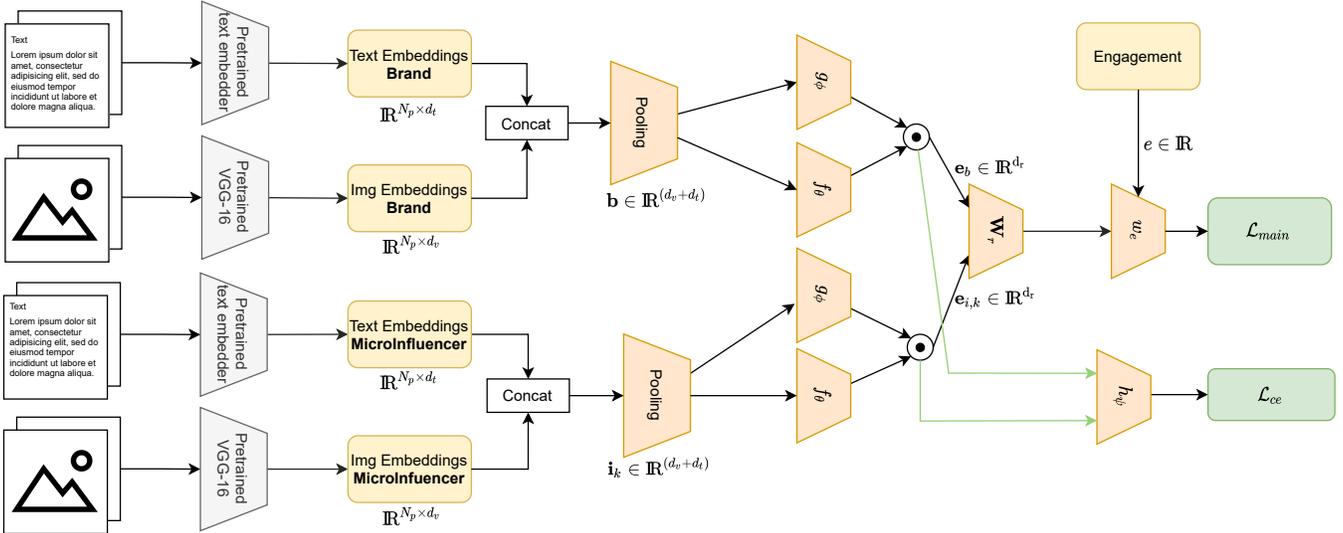
Fig. 2. Overall architecture for the multi-task model **WSim-MT**.

where $\mathbf{z}$ is a vector whose elements are obtained according to Equation 3 or 6 depending on the model being used. Then the model is trained with the cross entropy loss:

$$\mathcal{L}_{main}(\mathbf{y}, \mathbf{p}) = -\sum_{i}^{P} y_i log(p_i) \qquad (9)$$

where $p_i$ is the i-th element of vector $\mathbf{p}$ obtained as per Equation 8 and $y_i = \frac{1}{|\mathcal{I}_+(\mathbf{b})|}$ if micro-influencer $i$ is associated with brand $\mathbf{b}$ or 0 otherwise. In order to exploit the rank information of partial sequences we employed a simplified variant of the approach describe in [12]. In particular, given a brand $\mathbf{b}$, $P$ pools of $K$ candidate micro-influencers are created, such that in each pool a varying number of positive micro-influencers is always present. In particular, the number of positive influencer in a pool can range between 1 and $K$. Negative examples within the pool are randomly sampled without replacement from the dataset. In order to use all the valuable information provided within the dataset, the number of pools $P$ scales with the amount of positive examples for each given brand. For example, if for a given brand there are 5 positive examples, and the pool size is $k = 3$, then 15 pools will be created: 5 pools containing only one positive instance, 5 pools containing 2 positive instances and 5 pools with all positive instances. In this way, we guarantee that all the positive examples are used within the training procedure while keeping the computation tractable. This would not happen if $P$ would contain all the possible permutations of the partial sequences of positive examples, as the number of pools for a given brand $\mathbf{b}$ would be: $\sum_{k=1}^{K} \binom{|\mathcal{I}_+(\mathbf{b})|}{k}$ which can quickly become intractable. When employing the "partial-sequence mode" instead of the full "listwise" one, the definition of $y_i$ change slightly. The normalizing factor is not the number of positive examples associated with a brand anymore, instead it

is the number of positive examples associated with a brand within a given pool and it can range between 1 and $K$.

## IV. EXPERIMENTS

### A. Dataset

Open Source datasets for micro-influencer recommendation are hard to find and expensive to create, so there are a limited set of baseline tasks to choose from. Recently, the authors of [12] shared the first brand, micro-influencer multimedia dataset to spark research in the field. We therefore evaluate our methods on it in this work.

The dataset targets brand accounts on Instagram belonging to one of 12 macro categories: Airline, Auto, Clothing, Drink, Electronics, Entertainment, Food, Jewellery, Makeup, Non-Profit, Shoes and Services. 30 brands are selected for each category for a total of 360 brands. A micro-influencer is assumed to be associated with a brand when it is cited in one of the last 1000 posts of the brand. Moreover, to be eligible, a micro-influencer should have a number of followers ranging between 5000 and 100000. On average, each brand is associated with 11 micro-influencers. For each brand and micro-influencer in the dataset a mix of profile and posts-related information is available. In particular, the images, text, number of likes/comments of their last 50 posts is available ($N_p = 50$).

### B. Experimental Setup

We performed 5-fold cross validation to remain consistent with the 80-20 split performed in [12] and ensure a fair comparison. We enforce that all brand categories are available in each train-test split. The evaluation is done by scoring a brand against all the candidate micro-influencers. The obtained scores are then sorted in descending order to obtain the brand-specific micro-influencer ranking.

We trained all our models in an end-to-end fashion using the Adam optimiser [42] for 200 epochs with Early Stopping (where 20% of the samples in the training set were held out and used for validation purposes) and a learning rate of 0.001. The visual and textual embeddings' size obtained from the pretrained models are 25088 and 300 respectively. In **WSim-MT** $f_\theta$ and $g_\phi$ are fully connected neural networks with two hidden layers composed by (300, 512) and (4096, 512) hidden units respectively. We used ReLU activation and a droupout rate of 0.5, the intensity regularizer $\lambda$ and $\gamma$ described in Equation 7 for the multi-task model are both set to 0.5. Also, $h_\psi$ is an affine transformation followed by a softmax activation.

## C. Results

In this section we will compare the methods introduced in Section III-B with several baselines, some of which represent the state-of-the-art in brand/micro-influencers recommendations. The models considered for comparison are:

- **Random**: given a brand, the associated list of micro-influencers are randomly sampled among the available ones.
- **SimCos**: concatenate the images and text embeddings to obtain a unique vector representation for brand and micro-influencers. Once the concatenated vectors are obtained, each micro-influencer is scored according to its cosine similarity w.r.t the brand representation.
- **MIR(k)**: The method proposed in [12]. Differently from our methods, no trainable similarity is used between micro-influencer and brand; moreover, the engagement score is combined with the brand/micro-influencer similarity in a static, non trainable manner. Finally, multi-task learning was not considered in this previous work.

To remain consistent with previous work we employed four main metrics to evaluate our models. *AUC*, that measures the probability of a positive example being ranked higher then a negative one. *Recall @ k*, with $k \in \{10, 50\}$ as a classic recommender system metric; it considers the fraction of positive examples, among all the available one for a given brand, that are included in the top-k ranked items. Finally, *MedR* considers the median position of the highest ranked positive example for a given brand in the test set.

Table I summarizes the results obtained by our methods compared against the previous ones. For **MIR(k)**, we report both the results presented in the original paper (the ones with no uncertainty) and the results obtained via our implementation. Some discrepancy is expected as our implementation has been tested on multiple different test datasets.

We notice that the cosine similarity model, despite it's simplicity, already provides a strong baseline. It is therefore reasonable to expect **WSim**, a straightforward generalization of the previous model, to perform better. Indeed, by introducing minimal complexity into the model, we can move from a task-agnostic model to one which can specialize for the task at hand by optimizing a principled objective. This improvement is clear in Table I, where we can also observe some kind
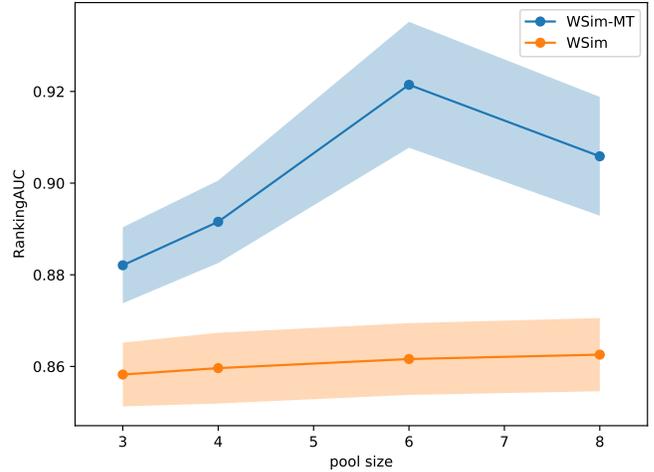


Fig. 3. RankingAUC as a function of the pol size for both WSim and Wsim-MT.

of insensitivity to the size of the pool of candidates micro-influencers used during training.

**WSim-MT** increases the complexity of the model, but gives a large improvement in all the considered metrics, in particular we improved by 8.2% in AUC, 42.9% in Recall @ 10 and 28% in Recall @50 compared to the previous state-of-the-art. We note that in **WSim** the use of a larger pool size has no significant impact on performance. It can be seen in Figure 3 that the algorithm is almost insensitive to this kind of change. The same does not hold for **WSim-MT**, where using a pool size of $k = 6$ provide the best results overall. We also note that continuously increasing the pool size at some point has a negative impact on performance. We attribute this behaviour to the sampling mechanism we adopted. Indeed, in order to keep the computation tractable, the number of samples that we use for training scales linearly with the pool size. This is not enough for larger pool sizes, as the model would likely need many more samples to learn an optimal ranking in this case.

## V. INTERPRETATION

Up to this point, we have introduced influencer ranking algorithms that perform well on overall metrics. However, as these techniques rely heavily on deep neural networks, it can be difficult for a human to understand why the posts of a particular influencer lead to its ranking position for a certain brand. Being able to better interpret why an influencer is matched well with a brand can be of interest to both brands and influencers. Along with providing verification that the ranking algorithm is functioning as expected, this information can help influencers or brands tailor their content to better match with the content of a social media account that has an audience that they wish to target.

As our dataset is based on Instagram posts, in which images are predominant, we focus on interpreting the visual content for individual posts that could be of interest to a particular brand or influencer. To understand which components of an

| | AUC | Recall@10 | Recall@50 | MedR | N.Params |
|---|---|---|---|---|---|
| Random. | $0.493 \pm 0.003$ | $0.007 \pm 0.002$ | $0.041 \pm 0.004$ | $71.4 \pm 5.3$ | 0 |
| SimCos | $0.755 \pm 0.006$ | $0.156 \pm 0.007$ | $0.368 \pm 0.009$ | $2.5 \pm 0.5$ | 0 |
| MIR(k=4) | $0.821 \pm 0.008$ | $0.082 \pm 0.006$ | $0.307 \pm 0.014$ | $8.7 \pm 1.3$ | $\sim$105M |
| MIR(from [12]) | 0.849 | 0.135 | 0.428 | 6 | - |
| **WSim(k=4)** | $0.858 \pm 0.007$ | $0.175 \pm 0.012$ | $0.467 \pm 0.020$ | $3.0 \pm 0.6$ | $\sim$**25K** |
| **WSim(k=6)** | $0.861 \pm 0.008$ | $0.178 \pm 0.013$ | $0.469 \pm 0.023$ | $2.9 \pm 0.5$ | $\sim$**25K** |
| **WSim-MT(k=4)** | $0.890 \pm 0.009$ | $0.163 \pm 0.004$ | $0.479 \pm 0.018$ | $2.2 \pm 0.4$ | $\sim$105M |
| **WSim-MT(k=6)** | $\mathbf{0.920 \pm 0.014}$ | $\mathbf{0.193 \pm 0.011}$ | $\mathbf{0.548 \pm 0.037}$ | $\mathbf{2.0 \pm 0.4}$ | $\sim$105M |

image are important, we introduce a technique inspired by the Grad-CAM algorithm [39], which operates on image classification networks, highlighting the parts of the input image that are most relevant to the output classification. Grad-CAM is designed to work on any convolutional neural network by determining an importance score across all the dimensions in the final convolutional layer of the network. These importances can be turned into a heatmap that can be superimposed over the input image by summing over the non-spatial dimensions of the feature space. This heatmap then highlights the most relevant areas for the final classification. In the original algorithm the importance scores are determined by the size of the gradients of the top predicted class with respect to the activations in the last convolutional layer.

We adapt these ideas to our use case by introducing a new measure of importance based on the architecture of the **WSim** model introduced in Section III-B. In this model, input images are passed through VGG-16 and activations of the final convolutional layer are used as a feature representation, $\mathbf{X}_{\text{vgg}} \in \mathbb{R}^{s_1 \times s_2 \times f_N}$, where $s_1$ and $s_2$ are the two spatial dimensions of the image and $f_N$ is the dimension of convolutional filter channels in the final layer. As described earlier, the visual image dimension, $d_v$ is the unrolled values of $\mathbf{X}_{\text{vgg}}$, such that:

$$d_v \equiv s_1 \times s_2 \times f_N \qquad (10)$$

These representations are pooled across posts for a brand or influencer and a similarity score between the two obtained through a trained bilinear similarity, which is parameterised by a diagonal matrix $\mathbf{W}_v \in \mathbb{R}^{d_v \times d_v}$. Within this matrix there are $d_v$ parameters, which are one-to-one with the parameters in $\mathbf{X}_{\text{vgg}}$, with a magnitude that determines how relevant particular features are when determining the similarity of two images. They therefore provide a natural measure of importance, $\mathbf{I}_v \in \mathbb{R}^{d_v}$, for each component of the visual features $\mathbf{X}_{\text{vgg}}$:

$$\mathbf{I}_v \equiv \text{diag}(\mathbf{W}_v) \qquad (11)$$

We can now use $\mathbf{I}_v$ to replace the gradients used in Grad-CAM, while allowing a heatmap for a particular image, $\mathbf{H} \in \mathbb{R}^{s_1 \times s_2}$,

to be built in the same way:

$$\mathbf{H} = \sum_i^{f_N} \mathbf{I}_v \odot \mathbf{X}_{\text{vgg}} \qquad (12)$$

where the unrolled components of $\mathbf{X}_{\text{vgg}}$ are one-to-one with the components of $\mathbf{I}_v$ and the non-spatial dimensions are summed over. This heatmap can be used to visualise the components of an image that are most relevant to the **WSim** model when calculating the similarity of any two accounts. This process applied to the images from two different posts from fashion and food influencers can be seen in Fig. 4. In the food focused post, we can see that the model focuses on the pizza and ignores the background. Other posts would be scored as being similar to this post if they had similarly predominant food based visual components. The same can be said for the person and their clothes in the fashion focused post. This helps to give us confidence that our model is performing as expected, paying most attention to the parts of the post a human would have deemed most important.

The importance score $\mathbf{I}_v$ just highlights the general purpose features that are important for calculating the similarity of images. It can be made more specific to a particular brand or influencer by multiplying the bilinear similarity matrix by the pooled feature vector of the brand, $\mathbf{b}_v$, or influencer, $\mathbf{i}_v$, to get a new importance vector, $\mathbf{I}_v^b \in \mathbb{R}^{d_v}$ in the case of a brand:

$$\mathbf{I}_v^b \equiv \mathbf{W}_v \mathbf{b}_v$$

This importance score is again one-to-one with the final convolutional layer of VGG-16 and can be used to derive a heatmap that highlights the components of an image most relevant to the brand or influencer in question, directly replacing $\mathbf{I}_v$ with $\mathbf{I}_v^b$ in Equation 12. This can be seen in Fig. 5, where the importance heatmap of an image containing a woman and a car has been calculated with $\mathbf{b}_v$ taken from either a fashion brand, or a car brand. In the case of the fashion brand, this importance measure suggests that the model pays most attention to the woman when calculating similarity. For the car brand the model pays most attention to the car. A similar comparison is carried out in Fig. 6 for a food and jewellery brand. In all these cases it is clear that the model is paying attention to

Fig. 4. Heatmaps of general image features the influencer ranking network (WSim) gives most importance to when comparing image similarities. Shown for a post from a fashion influencer (top) and a food influencer (bottom), where red is more important



Fig. 5. Heatmaps of brand specific image features the influencer ranking network (WSim) gives most importance to. On top is the original image, below is the importance for a fashion brand (bottom left) and a car brand (bottom right), where red is more important.
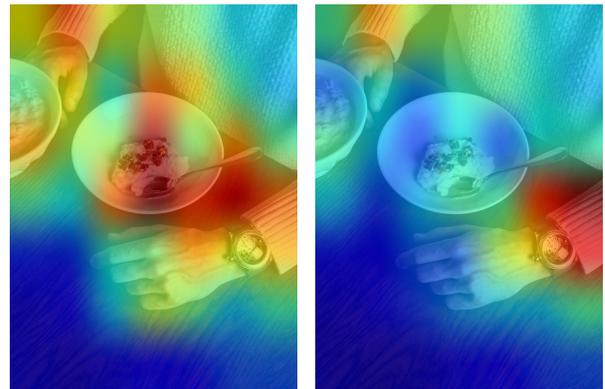


Fig. 6. Heatmaps of brand specific image features the influencer ranking network (WSim) gives most importance to. On top is the original, below is the importance for a food brand (bottom left) and a jewellery brand (bottom right), where red is more important.

the part of the image one would naively have thought most relevant to the brand in question.

These techniques for interpreting the **WSim** model help us to confirm the model is behaving as intuitively expected. However, they can also be used by brands or influencers that wish to analyse their media content. When deciding on images for new posts, for example, these techniques can help to highlight which components of new images are most similar to previous posts made by the brand. This allows a brand media manager to understand and direct their visual content. Such techniques could also be used to help brands reach new audiences, by choosing images most similar to the posts made by an influencer with an audience a brand would like to target.

## VI. CONCLUSIONS AND FUTURE WORK

We have improved on the state-of-the-art in micro-influencer ranking, building on the work originally presented in [12]. This was achieved by introducing two novel deep neural network architectures based on a trainable inner product (**WSim**), which drastically reduces the number of trainable parameters required, or the introduction of multi-task learning (**WSim-MT**), which performs best overall.

On top of this, inspired by Grad-CAM [39], we make use of the architecture of the **WSim** model to introduce a layer of interpretability of the ranking decisions made by our models. This has helped us validate that the ranking is taking into

account the features we would expect and provides tools for brands and influencers to improve their media content in the future.

It is worth noting that this work focuses on micro-influencers, but the tools developed, particularly the interpretability layer, can also be useful for influencers with large followings. These techniques can also be generalised to arbitrary multimedia ranking tasks that have datasets with a similar structure.

To expand on this work further, it would be interesting to test the performance on other social media platforms in which text is more relevant. If it works well, the techniques for interpretability could be expanded to text, allowing brands to isolate key words and phrases that help target particular audiences. Additionally, it could be interesting to investigate if the neural network architectures presented could be used for predicting the engagement of individual posts. In the case that this is successful, the interpretability layer could help brands isolate the components of images that are most important for producing engaging content.

## REFERENCES

[1] N. Zietek, "Influencer marketing: the characteristics and components of fashion influencer marketing," 2016.

[2] G. Appel, L. Grewal, R. Hadi, and A. T. Stephen, "The future of social media in marketing," *Journal of the Academy of Marketing Science*, vol. 48, no. 1, pp. 79–95, 2020.

[3] K.-Y. Goh, C.-S. Heng, and Z. Lin, "Social media brand community and consumer behavior: Quantifying the relative impact of user-and marketer-generated content," *Information Systems Research*, vol. 24, no. 1, pp. 88–107, 2013.

[4] D. Vrontis, A. Makrides, M. Christofi, and A. Thrassou, "Social media influencer marketing: A systematic review, integrative framework and future research agenda," *International Journal of Consumer Studies*, vol. n/a, no. n/a.

[5] C.-H. Cho and U. o. T. a. A. i. a. as, "Why do people avoid advertising on the internet?," *Journal of advertising*, vol. 33, no. 4, pp. 89–97, 2004.

[6] P. L. Breves, N. Liebers, M. Abt, and A. Kunze, "The perceived fit between instagram influencers and the endorsed brand: How influencer–brand fit affects source credibility and persuasive effectiveness," *Journal of Advertising Research*, vol. 59, no. 4, pp. 440–454, 2019.

[7] C. Lou and S. Yuan, "Influencer marketing: how message value and credibility affect consumer trust of branded content on social media," *Journal of Interactive Advertising*, vol. 19, no. 1, pp. 58–73, 2019.

[8] C. Abidin, "Communicative intimacies: Influencers and perceived interconnectedness," *Ada*, vol. 8, pp. 1–16, 2015.

[9] S. Woods, "# sponsored: The emergence of influencer marketing," 2016.

[10] A. P. Schouten, L. Janssen, and M. Verspaget, "Celebrity vs. influencer endorsements in advertising: the role of identification, credibility, and product-endorser fit," *International journal of advertising*, vol. 39, no. 2, pp. 258–281, 2020.

[11] B. Wissman, "Micro-influencers: The marketing force of the future," *Erişim adresi: https://www. forbes. com/sites/barrettwissman/2018/03/02/micro-influencers-the-marketing-force-of-the-future*, 2018.

[12] T. Gan, S. Wang, M. Liu, X. Song, Y. Yao, and L. Nie, "Seeking micro-influencers for brand promotion," in *Proceedings of the 27th ACM International Conference on Multimedia*, MM '19, (New York, NY, USA), p. 1933–1941, Association for Computing Machinery, 2019.

[13] Y. A. Argyris, Z. Wang, Y. Kim, and Z. Yin, "The effects of visual congruence on increasing consumers' brand engagement: An empirical investigation of influencer marketing on instagram using deep-learning algorithms for automatic image classification," *Computers in Human Behavior*, vol. 112, p. 106443, 2020.

[14] N. Zykun, Y. Zoska, A. Bessarab, V. Voronova, Y. Kyiashko, and D. Fayvishenko, "Branding as a social communication technology for managing consumer behavior," *International Journal of Management (IJM)*, vol. 11, no. 6, 2020.

[15] C. Campbell and J. R. Farrell, "More than meets the eye: The functional components underlying influencer marketing," *Business Horizons*, vol. 63, no. 4, pp. 469–479, 2020.

[16] M. De Veirman, V. Cauberghe, and L. Hudders, "Marketing through instagram influencers: the impact of number of followers and product divergence on brand attitude," *International journal of advertising*, vol. 36, no. 5, pp. 798–828, 2017.

[17] L. Hudders, S. De Jans, and M. De Veirman, "The commercialization of social media stars: a literature review and conceptual framework on the strategic use of social media influencers," *International Journal of Advertising*, pp. 1–49, 2020.

[18] C. Hughes, V. Swaminathan, and G. Brooks, "Driving brand engagement through online social influencers: An empirical investigation of sponsored blogging campaigns," *Journal of Marketing*, vol. 83, no. 5, pp. 78–96, 2019.

[19] T.-Y. Liu, "Learning to rank for information retrieval," 2011.

[20] H. Li, "A short introduction to learning to rank," *IEICE TRANSACTIONS on Information and Systems*, vol. 94, no. 10, pp. 1854–1862, 2011.

[21] Z. Cao, T. Qin, T.-Y. Liu, M.-F. Tsai, and H. Li, "Learning to rank: from pairwise approach to listwise approach," in *Proceedings of the 24th international conference on Machine learning*, pp. 129–136, 2007.

[22] T. Luo, D. Wang, R. Liu, and Y. Pan, "Stochastic top-k listnet," *arXiv preprint arXiv:1511.00271*, 2015.

[23] A. Farseev, K. Lepikhin, H. Schwartz, E. K. Ang, and K. Powar, "Somin. ai: Social multimedia influencer discovery marketplace," in *Proceedings of the 26th ACM international conference on Multimedia*, pp. 1234–1236, 2018.

[24] F. Gelli, T. Uricchio, M. Bertini, A. Del Bimbo, and S.-F. Chang, "Image popularity prediction in social media using sentiment and context features," in *Proceedings of the 23rd ACM international conference on Multimedia*, pp. 907–910, 2015.

[25] F. Gelli, T. Uricchio, X. He, A. Del Bimbo, and T.-S. Chua, "Beyond the product: Discovering image posts for brands in social media," in *Proceedings of the 26th ACM international conference on Multimedia*, pp. 465–473, 2018.

[26] R. Caruana, "Multitask," in *Learning to Learn*, pp. 95–133, Springer, Boston, MA., 1998.

[27] Y. Zhang and Q. Yang, "A survey on multi-task learning," *arXiv preprint arXiv:1707.08114*, 2017.

[28] R. Caruana, "Multitask learning," *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.

[29] A. Evgeniou and M. Pontil, "Multi-task feature learning," *Advances in neural information processing systems*, vol. 19, p. 41, 2007.

[30] A. Argyriou, T. Evgeniou, and M. Pontil, "Convex multi-task feature learning," *Machine learning*, vol. 73, no. 3, pp. 243–272, 2008.

[31] Z. Zhang, P. Luo, C. C. Loy, and X. Tang, "Facial landmark detection by deep multi-task learning," in *European conference on computer vision*, pp. 94–108, Springer, 2014.

[32] N. Mrkšić, D. O. Séaghdha, B. Thomson, M. Gašić, P.-H. Su, D. Vandyke, T.-H. Wen, and S. Young, "Multi-domain dialog state tracking using recurrent neural networks," *arXiv preprint arXiv:1506.07190*, 2015.

[33] S. Li, Z.-Q. Liu, and A. B. Chan, "Heterogeneous multi-task learning for human pose estimation with deep convolutional neural network," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pp. 482–489, 2014.

[34] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert, "Cross-stitch networks for multi-task learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3994–4003, 2016.

[35] D. Tellez, D. Höppener, C. Verhoef, D. Grünhagen, P. Nierop, M. Drozdzal, J. Laak, and F. Ciompi, "Extending unsupervised neural image compression with supervised multitask learning," in *Medical Imaging with Deep Learning*, pp. 770–783, PMLR, 2020.

[36] F. Manessi and A. Rozza, "Graph-based neural network models with multiple self-supervised auxiliary tasks," *arXiv preprint arXiv:2011.07267*, 2020.

[37] F. Fan, J. Xiong, and G. Wang, "On interpretability of artificial neural networks," *Preprint at https://arxiv. org/abs/2001.02522*, 2020.

[38] G. Montavon, W. Samek, and K.-R. Müller, "Methods for interpreting and understanding deep neural networks," *Digital Signal Processing*, vol. 73, pp. 1–15, 2018.

[39] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.

[40] M. Honnibal, I. Montani, S. Van Landeghem, and A. Boyd, "spaCy: Industrial-strength Natural Language Processing in Python," 2020.

[41] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[42] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings* (Y. Bengio and Y. LeCun, eds.), 2015.