# SEMANTIC IMAGE FUSION

**P.R. Hill**
Visual Information Laboratory
University of Bristol, UK
paul.hill@bristol.ac.uk

**D.R. Bull**
Visual Information Laboratory
University of Bristol, UK
dave.bull@bristol.ac.uk

October 14, 2021

## ABSTRACT

Image fusion methods and metrics for their evaluation have conventionally used pixel based or low level features. However, for many applications the aim of image fusion is to effectively combine the semantic content of the input images. This paper proposes a novel system for the semantic combination of visual content using pre-trained CNN network architectures. Our proposed semantic fusion is initiated through the fusion of the top layer feature map outputs (for each input image) through gradient updating of the fused image input (so called image optimisation). Simple 'choose maximum' and 'local majority' filter based fusion rules are utilised for feature map fusion. This provides a simple method to combine layer outputs and thus a unique framework to fuse single channel and colour images within a decomposition pre-trained for classification and therefore aligned with semantic fusion. Furthermore, class activation mappings of each input image are used to combine semantic information at a higher level. The developed methods are able to give equivalent low level fusion performance to state of the art methods while providing a unique architecture to combine semantic information from multiple images.

***Keywords*** Image Fusion · CNN

## 1 Introduction

Image fusion is the combination of multiple images into a single image that aims to combine the most important visual information from all sources [1]. Image fusion has been motivated by the need to improve visual representations, visualisation, scene understanding and situational awareness in multi-sensor and multi-camera applications such as remote sensing [2], medicine [3] and surveillance [4].

Image fusion has been driven by device and sensor limitations. For example, not all important visual information can be captured by one type of sensor (e.g. IR, visible etc.) or within one single shooting setting (i.e. focus, angle etc.). Furthermore, complementary imaging modalities co-exist within domains such as remote sensing and medicine that contain very different and important visual information. The effective combination of all the visual information from all image sources is therefore the aim of image fusion. Such a combined image is often effective for subsequent tasks such as scene understanding and target recognition.

Image fusion has been a highly researched area over the last half century. During this time image fusion is performed at or decision-level, feature-level and pixel-level [5]. Simple signal processing based pixel-level image fusion has given excellent results in preceding decades and continues to be used due to its high-efficiency and lack of a need for training data [5]. Pixel-level fusion can be further classified into decomposition based methods or Sparse Representation (SR) based methods [5]. In decomposition based techniques, the input images are decomposed into transform domains using methods such as complex wavelets (DT-CWT) [1], the Discrete Wavelet Transform (DWT) [6] and the contourlet transform [7]. Within the transform domain, coefficients are combined using suitably defined fusion rules (such as weighted-averaging [8] and "choose maximum" [9]).

More recently, state-of-the-art image fusion techniques have focused on the use of network based methods [10, 11, 12, 13, 14, 15, 16, 17]. Due to the requirements of needing training data these methods are often domain focused with state-of-the-art results reported for IR/Visible fusion [12, 13], remote sensing (multispectral) [14] and multi-focus areas [15, 16]. Recent work has also been focused on universal image fusion methods that can effectively combine multiple sources within all of these domains [17].

There have only been a very small number of previously developed fusion methods that utilise semantic information for image fusion [18, 19]. However, these methods do not use the joint classification and class activation maps to semantically fuse the input sources as proposed in our work. These previous works also only use semantic information in a very limited way.

## 1.1 Contributions

The contributions and characteristics of our work are summarised as follows.

- An unsupervised fusion technique is proposed that can combine the outputs of pre-trained low level feature maps using choose maximum and majority filter based fusion rules (using image optimisation).

- An unsupervised semantic fusion method that uses class activation maps.

- An unsupervised fusion method that combines both direct semantic fusion using class activation maps together with the combination of low-level network layer outputs (i.e. combining the above two approaches).

## 2 Feature Map Fusion using Image Optimisation

Feature map based loss functions within pre-trained networks for image optimisation have been used extensively in Neural Style Transfer (NST) [20] and "Deep Dream" methods [21]. This field was initiated by the seminal work in NST by Gatys et al. [22]. Gatys' NST system applied previously developed texture synthesis methods [23, 24] to the combination of two images (a "content" image and a "style" image) through an image optimisation method utilising a pre-trained Convolutional Neural Network (CNN): VGG19 [25]. Although these methods have provided amazing visual results they have yet to be utilised for general image processing applications such as image fusion and denoising. NST using image optimised can be summarised through the generation of the output style transferred image ($I^*$) using the minimisation of a loss function:

$$I^* = \arg\min_I \mathcal{L}_{total}(I|I_c, I_s) \tag{1}$$

$$= \arg\min_I \alpha\mathcal{L}_c(I|I_c) + \beta\mathcal{L}_s(I|I_s) \tag{2}$$

where $I^*$ is the output image, $I_s$ is the "style" image and $I_c$ is the "content" image and $\alpha$ and $\beta$ are the loss weighting parameters. The content and style losses ($\mathcal{L}_c, \mathcal{L}_s$) compare the content and style representations between the style and content images ($I_s, I_c$).

As an initial step, each of these two images are decomposed into the layer outputs of a VGG19 CNN network [25]. The two losses are calculated as functional comparisons between the layer outputs of the two images. Image optimisation is then achieved through gradient updates through the network using the gradient on the image w.r.t. the total loss $\mathcal{L}_{total}$. Since the original paper, many updates and optimisations to NST have been reported. Li and Wand [26] have proposed a Markov Random Field (MRF) based loss function that generates more plausible visual outputs. Computational optimised NST methods include Johnson et al. [27] and Ulyanov et al. [28]. These methods are similar to the Gatys' method in principle but are implemented using a single forward pass of a pre-trained network. Multiple-Style-Per-Model NST methods have included Dumoulin et al. [29], Li et al. [30] and Zhang and Dana [31]. Finally GANs [32], CycleGANs [33] and image transformers [34] have been recently used for NST. Although there have been many advances in this field, the Gatys method is still considered to be the gold standard by most researchers in terms of the quality of its results [20]. Therefore, although it is not computationally optimised or based on more complex transforms (GANs or Image Transformers), we have based our work on this type of image gradient update as it gives excellent results and is conceptually easy to understand and manipulate. Computational optimisation of our developed methods can be implemented as future work.

## 2.1 Fusion Rules and Loss Functions for Image Optimisation based Image Fusion

The image optimisation methods described above have almost exclusively been used for Neural Style Transfer. However, we propose such feature map image optimisation methods for image fusion. For image fusion, a variety of fusion rules and loss functions have been tested over an exhaustive set of layer subsets of various pre-trained CNN networks. It was found that taking losses across single layers gave the most effective results. Furthermore, it is recognised that as the layers get further from the input images the semantic information increases. This gives the possibility that fusion at such layers can combine more and more abstract semantic information. However, such high level layers have a reduced resolution and therefore increased spatial support of each feature in the feature map. It was found that fusing such layers generated unwanted artefacts such as banding. Layers with the same resolution as the input images were found to give the most effective results (e.g. the first two layers of the VGG19 CNN).

Although very sophisticated fusion rules and loss functions have been considered, it was found that a simple choose maximum fusion rule combined with a $l_2$ loss function in most cases gave the best results. The choice and utility of such a choose maximum fusion rule is motivated by its use within the wavelet transform domain [1, 6] i.e. large magnitude coefficients (or feature map outputs) correlates with perceptually important content.

Algorithm 1 illustrates the gradient update method used within the feature map based image fusion method. This algorithm is based on the key concept of updating the input image $I^*$, input to a network $N$ through a gradient update with respect to a loss function such as (6) (as implemented within the style transfer methods and the deep dream method). The loss function within the style transfer method is a weighted combination of the "content" and "style" losses defined through the comparison of the content and style image inputs. The loss function within the DeepDream method [21] is just based on the absolute magnitude of a chosen network output, layer or layer output. Our fusion loss function is given in (4). This is just the $L_2$ norm distance between the fused layer outputs (calculated just once outside the optimisation loops) and the iterated layer output feature map $F$. The actual image fusion optimisation process is defined in (3) being very similar to the NST equation (1) (but with the input images being the images to be fused $I_0$ and $I_1$ rather than the style and content images: $I_s$ and $I_c$).

### 2.1.1 Loss Functions

The fused image $I^*$ is generated using image minimisation of the fusion loss function $\mathcal{L}_f$:

$$I^* = \arg\min_I \mathcal{L}_f(I|I_0, I_1) \tag{3}$$

$$\mathcal{L}_f(I|I_0, I_1) = \sum_{l \in \{l_f\}} || \left( \Psi \left( (N^l(I_0), N^l(I_1)) - N^l(I) \right) ||_2^2 \tag{4}$$

where $N^l(I)$ is the network feature map output at layer $l$ when the network is input image $I$. $l_f$ is the set of layers to calculate the loss over and $\Psi$ is the fusion rule (defined below). Although an exhaustive combination of layers have been tested we have only utilised the top layer for our VGG19 CNN i.e. $l_f = \{'conv1\_1'\}$.

### 2.1.2 Fusion Rules

Our initial fusion rule $\Psi_0$ for combining the layer outputs is the choose maximum rule:

$$\Psi_0(F_0, F_1) = max\{F_x : x = 0...1\} \tag{5}$$

where $F_0$ and $F_1$ are the network outputs for the chosen layer (or set of layers) for input images 0 and 1 respectively. The maximum operator operates on a feature by feature basis in the selected output feature map. The output is therefore a tensor the same size as the two input tensors where the tensor elements are selected according to the maximum absolute magnitude of the two corresponding tensor elements.

The choice between each of the spatial positions within the network layer outputs is effectively a binary decision mask. This mask can be noisy and contain spatial inconsistencies. We therefore use a localised majority filter as an alternative fusion rule $\Psi_1$.

$$\Psi_1(F_0, F_1) = med_c\{|F_{0,i+r,c}| > |F_{1,i+r,c}| : r \in W\} \tag{6}$$

$F_{0,i,c}$ is the feature map value at vector spatial position $i$ and channel $c$ when when image $I_0$ is input into the network, $r$ is the spatial index vector of a sliding window ($3 \times 3$ in our case) and $med_c$ is the median value of the $3 \times 3$ boolean

map on a channel by channel basis. This median filter applied on a boolean map therefore gives a spatial majority filter on a channel by channel basis for the considered feature map.

---

**Algorithm 1:** Layer-Based Fusion using Gradient Update: We k = 100, e = 100, in our experiments.

---

**input** : Input Images: $I_0$, $I_1$
**output :** Fused Image: $I^*$
Initialise $I^*$: $I^* = mean(I_0, I_1)$
Fuse layer outputs: $G^l = \Psi\left((N^l(I_0), N^l(I_1))\right)$
**for** *e Epochs* **do**
    Update ADAM learning rate: $\lambda$;
    **for** *k Iterations* **do**
        Calculate layer feature loss: $\mathcal{L}_f(I^*|I_0, I_1) = \sum_{l \in \{l_f\}} || \left(G^l - N^l(I^*)\right) ||_2^2$
        Update input image: $I^* = I^* + \lambda \nabla_{I^*}\left(\mathcal{L}_f(I^*|I_0, I_1)\right)$
    **end**
**end**

---

## 3 Semantic Fusion Using Class Activation Mappings (CAMs)

The use of the more abstract layers within a CNN theoretically provides the ability to semantically fuse images using the image optimisation method described above. However, due to these feature maps having very limited spatial resolution such fusion results in significant edge based artefacts. We have therefore chosen to utilise Class Activation Mappings (CAMs) that give a spatial mapping of how relevant each pixel has been in generating a classification of a given class.

### 3.1 Class Activation Mappings

CAMs were first introduced by Zhou et al. [35] and generate an importance spatial mapping according to a class $c$ through the sum of output weights of an input image for the last layer of a CNN. This method was generalised to the use of backpropagated class gradients Grad-CAMs [36]. These methods have been further extended using methods such as Eigen-CAM [37], Grad-CAM using Vision-Transformers [38] etc. However, for our work the mappings from Grad-CAM gave the best results (due to their spatial consistency).

Given the definitions of a Grad-CAM mapping by Selvaraju et al. [36] a spatial mapping (the same resolution as the input image) can be represented as $M_c^l(x, y)$ where $c$ is the class under consideration, $l$ is the analysed layer (usually the spatial layer closest to the actual classification layer). We found that using the lowest resolution feature map, the mappings give the least spatial localisation. Using higher resolution feature map layers gives better spatial localisation but less accurately reflects the abstract class activation's as the lower (more abstract) layers. We therefore combine the lowest three spatial layers as follows:

$$P_c(x, y) = Norm\left(\prod_{l \in lset} M_c^l(x, y)\right) \tag{7}$$

where $Norm$ is the normalisation function:

$$Norm(z) = \frac{z - min(z)}{max(z) - min(z)} \tag{8}$$

As the $Norm$ function normalises the combination of the CAM mappings to the range [0,1] it can be considered as the probability $P_c$ of each pixel contributing to the classification of the top object recognised in each image (class $c_0$ for image $I_0$ and class $c_1$ for image $I_1$). For our utilised VGG19 CNN, $lset$ is the set of the last three coarsest resolution layers i.e. $lset = \{'conv3\_4',' conv4\_4',' conv5\_4'\}$.

These probabilities are termed $P_{c_0,I_0}(x, y)$ and $P_{c_0,I_1}(x, y)$ for input images $I_0$ and $I_1$ at spatial positions $(x, y)$. From these probabilities we need to generate a mixing ratio for the image fusion (also in the form of a probability termed $P_{M_0}(x, y)$: the probability that the output image should contain input image $I_0$). $P_{M_0}(x, y)$ is generated as an exclusive combination of $P_{c_0,I_0}(x, y)$ and $P_{c_1,I_1}(x, y)$ i.e.

- When $P_{c_0,I_0}(x, y)$ and $P_{c_1,I_1}(x, y)$ are approximately equal (for all values between 0 and 1), $P_{M_0}(x, y)$ should give an equal mix of each image in the output (i.e. $P_{M_0}(x, y) \approx 0.5$).
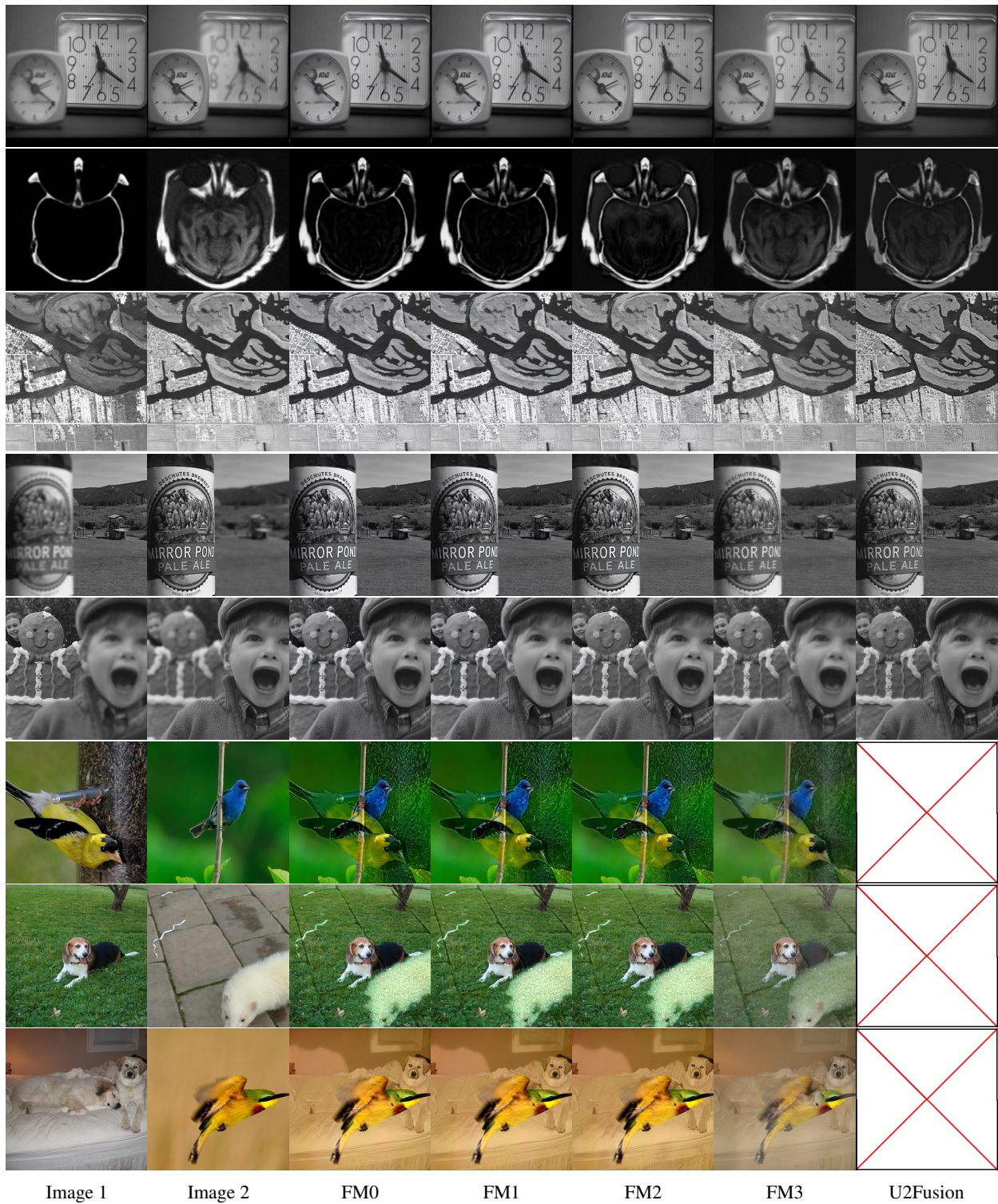
Figure 1: Fused image results. From top to bottom the fusion pairs are labelled "Clock", "Head", "Remote", "Bottle", "Gingerbread", "Goldfinch/Indigo_bunting", "Beagle/Ferret", "Great_Pyrenees/Bee-Eater". U2Fusion results are not available for the last three rows as it is not possible to be used for colour images.

- When $P_{c_0,I_0}(x,y)$ is small (i.e. near 0) and $P_{c_1,I_1}(x,y)$ is large (i.e. near 1) then $P_{M_0}(x,y)$ should approximate 0.
- When $P_{c_1,I_1}(x,y)$ is small (i.e. near 0) and $P_{c_0,I_0}(x,y)$ is large (i.e. near 1) then $P_{M_0}(x,y)$ should approximate 1.

This is achieved by defining $P_{M_0}(x,y)$ as (illustrated in figure 2):

$$P_{M_0}(x,y) = 0.5(1 + P_{c_0,I_0}(x,y) - P_{c_1,I_1}(x,y)) \tag{9}$$
$$P_{M_1}(x,y) = 1 - P_{M_0}(x,y) \tag{10}$$



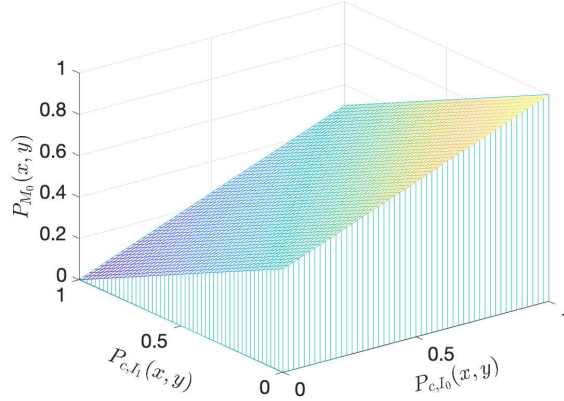Figure 2: Generating function (9) to obtain mixing probability $P_{M_0}(x,y)$ from $P_{c_0,I_0}(x,y)$ and $P_{c_1,I_1}(x,y)$

Fusion is therefore achieved using $P_{M_0}(x,y)$ and $P_{M_1}(x,y)$ as the mixing ratios of the two input images. The fused image $I^*$ is therefore defined as (dropping the $x,y$ indices).

$$I^* = P_{M_0}I_0 + P_{M_1}I_1 \tag{11}$$

## 4   Results

Table 1 shows objective results for all the methods (defined below) for the single channel input images. These have been compared to state of the art methods U2Fusion [17] and GMC [39]. The fusion metrics include: Wang ($Q_0$ [40]), Xydeas ($P_e$ [41]) and Piella ($Q$ [42]). These metrics can only be effectively used for single channel images; colour images are therefore not included in this table. This table shows that for the majority of the image pairs and metrics, the FM0 method (image optimisation based image fusion) gives the best performance. The GMC method gives the best results for all metrics for the head image pair. This is unsurprising as GMC was specifically designed to be used in this domain [39].

### 4.1   Method Definitions

**FM0,** Fusion Method 0: This method utilises image optimisation using (3) and (4) and fusion rule (choose maximum) $\Psi_0$ defined in (5).

**FM1,** Fusion Method 1: This method utilises image optimisation using (3) and (4) and fusion rule (majority filter) $\Psi_1$ defined in (6).

**FM2,** Fusion Method 2: This combines method FM0 and FM3 (i.e. feature maps $F_0$ and $F_1$ are multiplied by CAM probabilities $P_{c_0,I_0}$ and $P_{c_1,I_1}$ respectively).

**FM3,** Fusion Method 3: Fusion using CAMs utilising (7),(8),(9),(10) and (11).

### 4.2   Classification for Semantic Image Fusion

In order to generate the Class Activation Mappings (CAMs) the chosen network (VGG19) was used to get the top classification class $c$ for each of the input images i.e. class $c_0$ is the top classification class for input image $I_0$ and $c_1$

is the top classification class for input image $I_1$. These classifications don't always make sense for all the considered applications and domains (e.g. the "head" and "remote" images in Figure 1). However, table 2 shows the top classification classes for the image pairs shown in figure 3. Figure 3 also shows the CAM mappings for the top classes for each of the image pairs. This figure shows how the different semantic objects are highlighted by the class activation mappings and how this is utilised to combine and semantically fuse the input images. This figure shows how the utilisation of CAM mappings can distinguish between spatial regions that are in focus in multi-focus fusion pairs (see the clock CAMs in the top row).

| Dataset | Method | $Q_0$ | $P_e$ | $Q$ |
|---|---|---|---|---|
| Clock | FM0 | 0.8294 | **0.6308** | **0.9580** |
| | FM1 | 0.8293 | 0.6262 | 0.9572 |
| | FM2 | **0.8296** | 0.6199 | 0.9544 |
| | FM3 | 0.8293 | 0.5881 | 0.9495 |
| | U2Fusion | 0.8267 | 0.5409 | 0.8970 |
| | GMC | 0.8286 | 0.5208 | 0.9325 |
| Head | FM0 | 0.8035 | 0.3425 | 0.3045 |
| | FM1 | 0.8034 | 0.3411 | 0.3049 |
| | FM2 | 0.8045 | 0.5539 | 0.6614 |
| | FM3 | 0.8066 | 0.4753 | 0.7117 |
| | U2Fusion | 0.8072 | 0.2732 | 0.7056 |
| | GMC | **0.8216** | **0.6660** | **0.8002** |
| Remote | FM0 | 0.8100 | **0.6490** | **0.8640** |
| | FM1 | 0.8103 | 0.6394 | 0.8604 |
| | FM2 | 0.8106 | 0.6364 | 0.8597 |
| | FM3 | 0.8121 | 0.5827 | 0.8544 |
| | U2Fusion | 0.8097 | 0.5458 | 0.8487 |
| | GMC | **0.8156** | 0.6210 | 0.7708 |
| Bottle | FM0 | 0.8196 | **0.7112** | **0.9491** |
| | FM1 | 0.8191 | 0.6993 | 0.9470 |
| | FM2 | **0.8199** | 0.6960 | 0.9476 |
| | FM3 | 0.8164 | 0.4688 | 0.8819 |
| | U2Fusion | 0.8168 | 0.6111 | 0.9286 |
| | GMC | 0.8166 | 0.4830 | 0.8725 |
| Gingerbread | FM0 | 0.8217 | **0.6681** | **0.9522** |
| | FM1 | 0.8215 | 0.6633 | 0.9517 |
| | FM2 | **0.8230** | 0.6497 | 0.9489 |
| | FM3 | 0.8208 | 0.5218 | 0.9190 |
| | U2Fusion | 0.8204 | 0.5947 | 0.9430 |
| | GMC | 0.8205 | 0.5518 | 0.9033 |

Table 1: The objective results of different methods (these comparisons are only possible for single channel image pairs). Metrics: Wang ($Q_0$ [40]), Xydeas ($P_e$ [41]), Piella ($Q$ [42]). FM0-3 defined in section 4.1. Comparison techniques: U2Fusion [17] and GMC [39].

| **Fusion Pair** | $c_0$, $\mathbf{P}(c_0)$ | $c_1$, $\mathbf{P}(c_1)$ | $c_{FM3}$, $\mathbf{P}(c_{FM3})$ |
|---|---|---|---|
| Clocks (Top row figure 3) | Analog_Clock, 0.524 | Analog_Clock, 0.722 | Analog_Clock, 0.663 |
| Beagle/Ferret (Second row figure 3) | Beagle, 0.970 | Black-Footed_Ferret, 0.449 | Beagle, 0.927 |
| Great_Pyrenees/Bee-Eater (Third row figure 3) | Great_Pyrenees, 0.858 | Bee-Eater, 0.991 | Great_Pyrenees, 0.656 |
| Goldfinch/Indigo_bunting (Fourth row figure 3) | Goldfinch, 0.9999 | Indigo_Bunting, 0.999 | Goldfinch, 0.780 |

Table 2: VGG19 top classification results (class and probability) for input and FM3 fused images (from figure 3).

## 5   Conclusion

This paper has defined four novel fusion methods that utilise: image optimisation; choose maximum and majority filter fusion rules; and class activation mappings for semantic fusion. The image optimisation method was able to

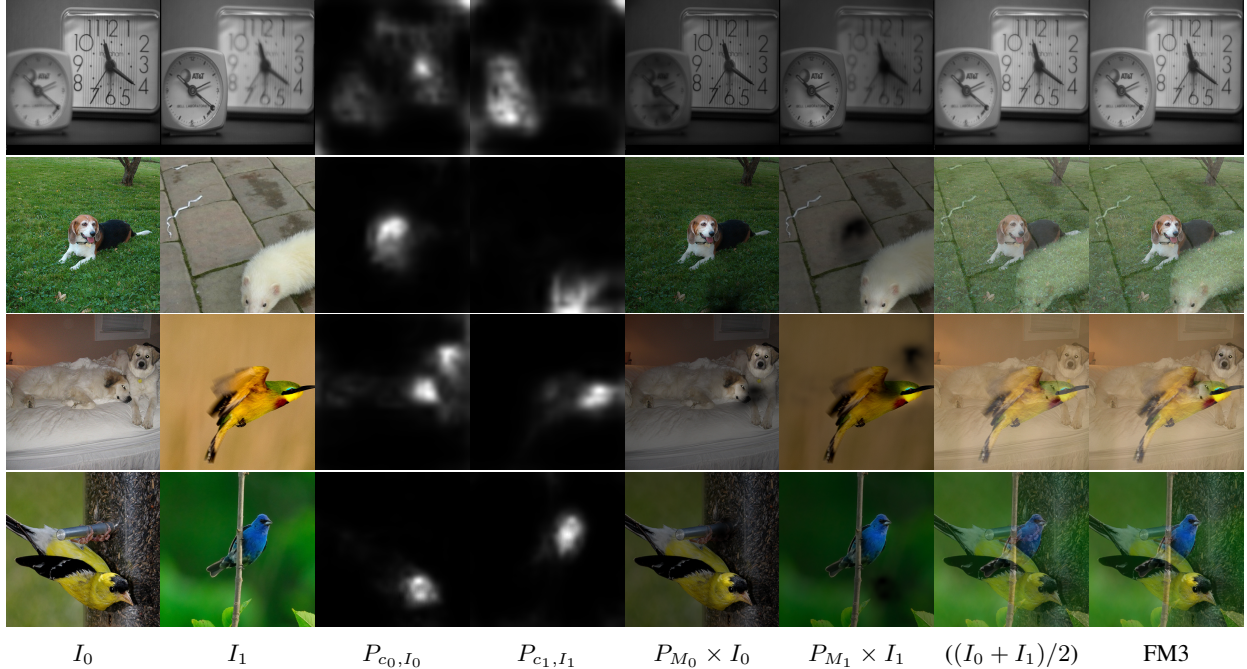|  $I_0$  |  $I_1$  |  $P_{c_0,I_0}$  |  $P_{c_1,I_1}$  |  $P_{M_0} \times I_0$  |  $P_{M_1} \times I_1$  |  $((I_0 + I_1)/2)$  |  FM3  |

Figure 3: Fused image results. From top to bottom the fusion pairs are labelled "Clock", "Beagle/Ferret", "Great_Pyrenees/Bee-Eater" and "Goldfinch/Indigo_bunting". This image shows the input image, CAM mappings, the CAM based mixtures from each image, the averaged image output and the fused image output (FM3).

achieve state of the art image fusion results measured using conventional image fusion metrics in multiple domains. This method is also directly extendable to colour images (not historically a major focus of image fusion).

Furthermore the CAM based methods can, for the first time, directly exploit semantic information from the top classified class in each of the fused images to generate true semantic level image fusion. It is conjectured that combining regions that are important to the classification of a set of images will most effectively combine the semantic information from the input images. Images combined in this semantically meaningful way are shown to retain the important semantic information from both images.

This method would easily extend to multiple images and the fusion of the top-5 classes in each image.

## References

[1] P. R. Hill, C. N. Canagarajah, and D. R. Bull, "Image fusion using complex wavelets." in *BMVC*, 2002, pp. 1–10.

[2] H. Ghassemian, "A review of remote sensing image fusion methods," *Information Fusion*, vol. 32, pp. 75–89, 2016.

[3] A. P. James and B. V. Dasarathy, "Medical image fusion: A survey of the state of the art," *Information fusion*, vol. 19, pp. 4–19, 2014.

[4] N. Paramanandham and K. Rajendiran, "Multi sensor image fusion for surveillance applications using hybrid image fusion algorithm," *Multimedia Tools and Applications*, vol. 77, no. 10, pp. 12 405–12 436, 2018.

[5] S. Li, X. Kang, L. Fang, J. Hu, and H. Yin, "Pixel-level image fusion: A survey of the state of the art," *information Fusion*, vol. 33, pp. 100–112, 2017.

[6] A. Ben Hamza, Y. He, H. Krim, and A. Willsky, "A multiscale approach to pixel-level image fusion," *Integrated Computer-Aided Engineering*, vol. 12, no. 2, pp. 135–146, 2005.

[7] S. Yang, M. Wang, L. Jiao, R. Wu, and Z. Wang, "Image fusion based on a new contourlet packet," *Information Fusion*, vol. 11, no. 2, pp. 78–84, 2010.

[8] J. H. Jang, Y. Bae, and J. B. Ra, "Contrast-enhanced fusion of multisensor images using subband-decomposed multiscale retinex," *IEEE Transactions on Image Processing*, vol. 21, no. 8, pp. 3479–3490, 2012.

[9] J. Hu and S. Li, "The multiscale directional bilateral filter and its application to multisensor image fusion," *Information Fusion*, vol. 13, no. 3, pp. 196–206, 2012.

[10] F. D. Javan, F. Samadzadegan, S. Mehravar, A. Toosi, R. Khatami, and A. Stein, "A review of image fusion techniques for pan-sharpening of high-resolution satellite imagery," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 171, pp. 101–117, 2021.

[11] H. Hermessi, O. Mourali, and E. Zagrouba, "Multimodal medical image fusion review: Theoretical background and recent advances," *Signal Processing*, p. 108036, 2021.

[12] H. Li, X.-J. Wu, and T. Durrani, "Nestfuse: An infrared and visible image fusion architecture based on nest connection and spatial/channel attention models," *IEEE Transactions on Instrumentation and Measurement*, vol. 69, no. 12, pp. 9645–9656, 2020.

[13] H. Li, X.-J. Wu, and J. Kittler, "Infrared and visible image fusion using a deep learning framework," in *2018 24th international conference on pattern recognition (ICPR)*.    IEEE, 2018, pp. 2705–2710.

[14] Z. Shao and J. Cai, "Remote sensing image fusion with deep convolutional neural network," *IEEE journal of selected topics in applied earth observations and remote sensing*, vol. 11, no. 5, pp. 1656–1669, 2018.

[15] H. Tang, B. Xiao, W. Li, and G. Wang, "Pixel convolutional neural network for multi-focus image fusion," *Information Sciences*, vol. 433, pp. 125–141, 2018.

[16] M. Amin-Naji, A. Aghagolzadeh, and M. Ezoji, "Ensemble of cnn for multi-focus image fusion," *Information fusion*, vol. 51, pp. 201–214, 2019.

[17] H. Xu, J. Ma, J. Jiang, X. Guo, and H. Ling, "U2fusion: A unified unsupervised image fusion network," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[18] R. Fernandez-Beltran, J. M. Haut, M. E. Paoletti, J. Plaza, A. Plaza, and F. Pla, "Remote sensing image fusion using hierarchical multimodal probabilistic latent semantic analysis," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 11, no. 12, pp. 4982–4993, 2018.

[19] F. Fan, Y. Huang, L. Wang, X. Xiong, Z. Jiang, Z. Zhang, and J. Zhan, "A semantic-based medical image fusion approach," *arXiv preprint arXiv:1906.00225*, 2019.

[20] Y. Jing, Y. Yang, Z. Feng, J. Ye, Y. Yu, and M. Song, "Neural style transfer: A review," *IEEE transactions on visualization and computer graphics*, vol. 26, no. 11, pp. 3365–3385, 2019.

[21] A. Mordvintsev, C. Olah, and M. Tyka, "Deepdream-a code example for visualizing neural networks," *Google Research*, vol. 2, no. 5, 2015.

[22] L. A. Gatys, A. S. Ecker, and M. Bethge, "A neural algorithm of artistic style," *arXiv preprint arXiv:1508.06576*, 2015.

[23] ——, "Texture synthesis and the controlled generation of natural stimuli using convolutional neural networks," in *Bernstein Conference 2015*, 2015, pp. 219–219.

[24] J. Portilla and E. P. Simoncelli, "A parametric texture model based on joint statistics of complex wavelet coefficients," *International journal of computer vision*, vol. 40, no. 1, pp. 49–70, 2000.

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

[26] C. Li and M. Wand, "Combining markov random fields and convolutional neural networks for image synthesis," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2479–2486.

[27] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European conference on computer vision*.    Springer, 2016, pp. 694–711.

[28] D. Ulyanov, V. Lebedev, A. Vedaldi, and V. S. Lempitsky, "Texture networks: Feed-forward synthesis of textures and stylized images." in *ICML*, vol. 1, no. 2, 2016, p. 4.

[29] V. Dumoulin, J. Shlens, and M. Kudlur, "A learned representation for artistic style," *arXiv preprint arXiv:1610.07629*, 2016.

[30] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang, "Diversified texture synthesis with feed-forward networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3920–3928.

[31] H. Zhang and K. Dana, "Multi-style generative network for real-time transfer," in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2018, pp. 0–0.

[32] L. Zhang, Y. Ji, X. Lin, and C. Liu, "Style transfer for anime sketches with enhanced residual u-net and auxiliary classifier gan," in *2017 4th IAPR Asian Conference on Pattern Recognition (ACPR)*.   IEEE, 2017, pp. 506–511.

[33] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.

[34] Y. Deng, F. Tang, X. Pan, W. Dong, C. Xu *et al.*, "Stytrˆ 2: Unbiased image style transfer with transformers," *arXiv preprint arXiv:2105.14576*, 2021.

[35] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2921–2929.

[36] 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*, 2017, pp. 618–626.

[37] M. B. Muhammad and M. Yeasin, "Eigen-cam: Class activation map using principal components," in *2020 International Joint Conference on Neural Networks (IJCNN)*.   IEEE, 2020, pp. 1–7.

[38] A. Aldahdooh, W. Hamidouche, and O. Deforges, "Reveal of vision transformers robustness against adversarial attacks," *arXiv preprint arXiv:2106.03734*, 2021.

[39] N. Anantrasirichai, R. Zheng, I. Selesnick, and A. Achim, "Image fusion via sparse regularization with non-convex penalties," *Pattern Recognition Letters*, vol. 131, pp. 355–360, 2020.

[40] Q. Wang, Y. Shen, and J. Jin, "Performance evaluation of image fusion techniques," *Image fusion: algorithms and applications*, vol. 19, pp. 469–492, 2008.

[41] C. Xydeas and V. Petrovic, "Objective image fusion performance measure," *Electronics letters*, vol. 36, no. 4, pp. 308–309, 2000.

[42] G. Piella and H. Heijmans, "A new quality metric for image fusion," in *Proceedings 2003 International Conference on Image Processing (Cat. No. 03CH37429)*, vol. 3.   IEEE, 2003, pp. III–173.