

Spot the Difference by Object Detection

Junhui Wu^{1,2*} Yun Ye² Yu Chen² Zhi Weng²

¹ Tsinghua University ² JD Group

wujh10@outlook.com, {yeyun, chenyu6, wengzhi}@jd.com

Abstract

In this paper, we propose a simple yet effective solution to a change detection task that detects the difference between two images, which we call “spot the difference”. Our approach uses CNN-based object detection by stacking two aligned images as input and considering the differences between the two images as objects to detect. An early-merging architecture is used as the backbone network. Our method is accurate, fast and robust while using very cheap annotation. We verify the proposed method on the task of change detection between the digital design and its photographic image of a book. Compared to verification based methods, our object detection based method outperforms other methods by a large margin and gives extra information of location. We compress the network and achieve 24 times acceleration while keeping the accuracy. Besides, as we synthesize the training data for detection using weakly labeled images, our method does not need expensive bounding box annotation.

1. Introduction

This paper considers a problem of change detection between two book cover images, of which one is the digital design and the other is the photographic printed image. Book cover digital design may somehow change at the time of publication, but the digital design displayed on the Internet or stored in database system is still the original one. Book sellers want to make sure that a printed book cover conforms to its digital design in the database system, because differences between the two images may cause commercial dispute. In some cases, the book cover changes globally. In other cases, book cover changes locally but the change can not be ignored. For example, the first part and second part of one book series change locally, but the difference is very important. Some changed book cover photographic images and the original digital designs are shown in Figure 1.

Although changed book covers are rarely seen, it is nec-



Figure 1: Some examples of changed book cover data. Digital designs are on the top and photographic images are on the bottom. The differences are (a) text about version information, (b) local patterning, (c) book girdle, (d) the whole image respectively. The task is to spot these differences between digital design and photographic image.

essary for a book seller to check if any real printed book cover is exactly the same as the digital design, which is very labor intensive and error-prone if done by human. So we design a device and propose an computer vision algorithm to save labor costs and improve accuracy. The structure of the device and the pipeline of the application is shown in Figure 2. The book is put on a black conveyor belt, and a camera is put on the top to take picture of the book cover. At the same time, the bar code on the back of the book is scanned, and the digital design of the book cover is retrieved from the database. After alignment, our algorithm judges if the photographic image and the digital design are same or different.

Different from general verification task, which compares key features, this task values local features, as small changes in character or patterning mean the two images are different. For example, in face verification [22], faces with 10% or more areas occluded are considered the same, but a book cover with a black area occlusion is considered differ-

*This work is done when Junhui Wu was an intern at JD Group

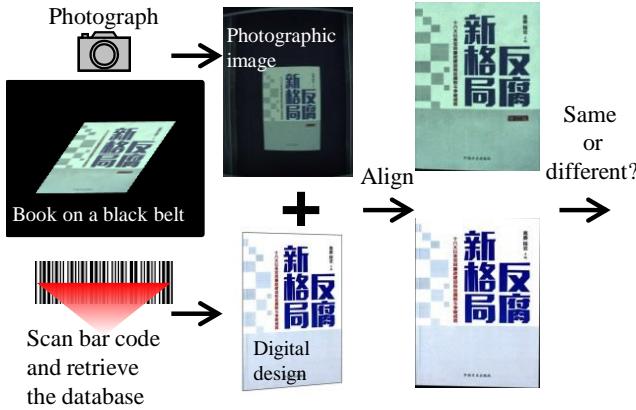


Figure 2: The pipeline of the application: The book is put on a black conveyor belt, and a camera on the top take picture of book cover. At the same time, the bar code at the back of the book is scanned, and the digital design of the book cover is retrieved form the database. After alignment, the algorithm judges if the photographic image and the digital design are same or different.

ent from its digital design. So the conventional verification methods are not suitable for this task.

The task is similar to the game “Spot the Difference”¹. However, it is more complicated than the game, because the photographic image is interfered with illumination, color, and other noises. In other words, the same parts between two images are not exactly looked the same, and the criteria of “difference” need to be represented robustly. Subtraction of pixels dose not work well.

To solve the spotting the difference problem, we propose a new method by CNN (Convolutional Neural Networks) based object detection. We stack the two RGB images as a 6-channel image, and consider the differences as objects to detect, then follow the pipeline of state-of-the-art object detection framework. In particular, the object detection framework is Faster R-CNN [17], as shown in Figure 3. We also propose an effective and efficient backbone network architecture based on object detection network by stacking two branches of input data after several convolutional layers, which we call two-branch architecture.

To avoid costly data annotation, our training images are synthesized from weakly labeled image pairs without bounding boxes. The results show detection model trained on synthetic data performs well on both synthetic data and realistic data.

In the experiments, we show that the object detection based method achieves much better performance than the general verification methods. We also compress the network and observe that even with smaller and faster net-

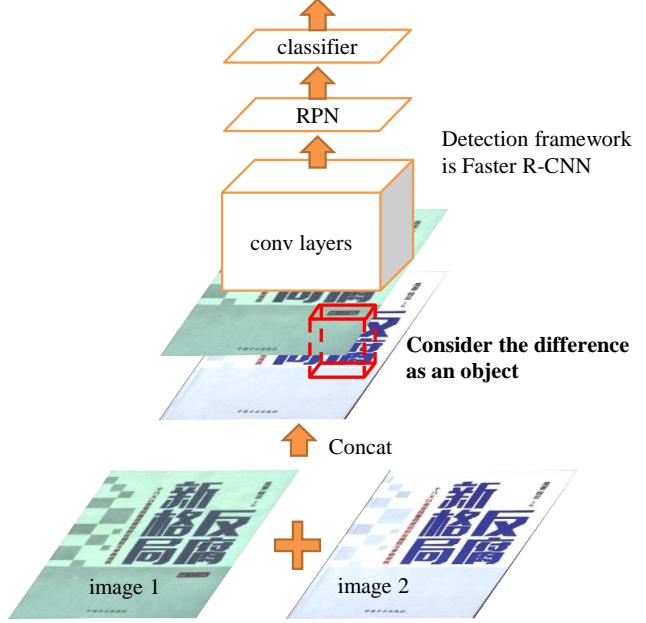


Figure 3: Spot the difference by object detection: two input images are concatenated along the channel dimension and the difference is considered as an object to detect. Faster R-CNN is adopted for object detection.

works, our method can still achieve good results.

The rest of the paper is organized as follows. Section 2 briefly reviews the related work. Then algorithm details are introduced in section 3. In section 4, we conduct various experiments. Section 6 concludes the paper.

2. Related Work

Many researches have been conducted in fields of image verification and change detection. However, most of these methods are not practical in this task. In this paper, we propose a new solution to this spotting the differences task, which is related to object detection. In the following, we introduce similar task solutions and explain why they are not practical in this task, including verification and change detection. We also introduce methods we used – object detection.

2.1. Verification

Verification problem has been developed relatively well in recent decades, including face verification [3, 23, 22] and signature verification [12]. The common outline of face verification consists of four stage: detect, align, represent, classify [23]. In represent stage, feature extractor is designed to extract key features like structure of eyes and nose, but ignore other element like glasses, scars and wrinkles. In this spotting the differences task, there is no such key

¹See wiki: https://en.wikipedia.org/wiki/Spot_the_difference

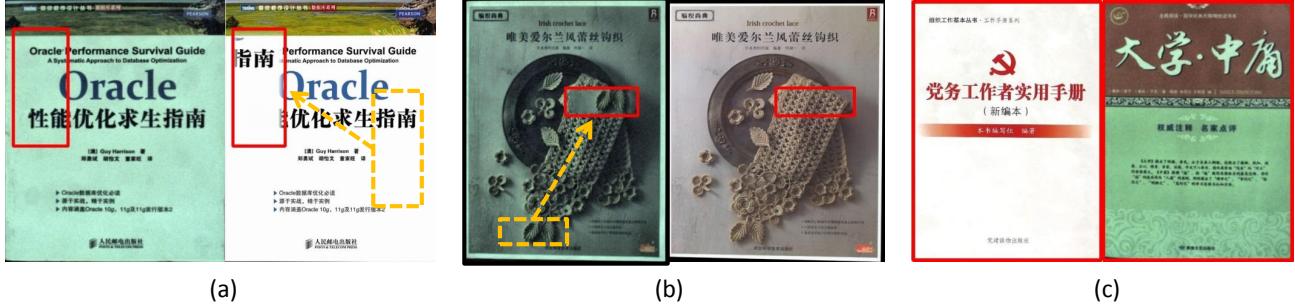


Figure 4: Three types of synthetic training data. For type (a), we sample a rectangle patch(in orange) from the digital design (in the right) and copy it to somewhere else(in red) as a difference from the photographic image (in the left). For type (b), we sample a rectangle patch from the photographic image and copy it to somewhere else as a difference from the digital design. For type (c), we sample miss-matched images as a training pair, of which the difference is the whole image.

features, for every character and patterning matters. For example, a book cover image with a version number or without should be different. So the flow of verification can not solve it.

2.2. Change Detection

Change Detection problem contains lots of applications like video surveillance, remote sensing, civil infrastructure, etc [14]. However, different applications and image sources have different processing steps. For change detection in video surveillance, image sequence is calculated and the relationship of images in time series is important [1, 26]. In remote sensing application [21, 25, 8], land-cover and land-use change information is considered, a core technique of state-of-the-art solution is doing segmentation and classification to images first, and then some comparison and analysis [8]. A solution for detecting changes of the three-dimensional structure of an outdoor scene is to estimate depth first and then compare the depth [18]. One of the document image change detection methods uses OCR method to recognize characters in images first [9]. Generally speaking, changes in different applications vary a lot and the solutions are not very practical in our task.

2.3. Object Detection

Numerous works have been proposed for object detection, and CNN based methods show leading results on detection benchmarks. From R-CNN [6], many CNN based methods are proposed to improve both accuracy and speed, including Fast-R-CNN [5], Faster-R-CNN [17], SSD [10], YOLO [15, 16]. CNN-based detection works for specific object like face [13] and text [29] also show great advantage. Different from hand-crafted feature extractor like SIFT [11] and HOG [4], features of a kind of object is learned from a deep neural network in these CNN based detection methods. Since the RGB pixel input can lead to features that describe a face or a text, it is reasonable that a

pair of RGB pixel input or 6 channel input leads to features that describe “difference”.

3. Algorithm

3.1. Image Alignment

The photographic image of book cover is taken on a black conveyor belt and the book cover is not at a fixed position of the image, thus we take a preprocessing pipeline to get a regular rectangle book cover, and then match it to the digital design. Following is the processing details.

First, the photographic image minus the pure black conveyor belt image background which still contains some texture. Then, we extract contour of book cover in the photographic image. We use Canny algorithm to extract edges and binary the image. After abandoning tiny pattern, we extract contour, approximate edges to a polygon, and finally extract a rectangle contour. Finally, we match and align the two images. We extract SIFT features from two images and match them, then filter the matches based on RANSAC and apply affine transformation to the digital design.

3.2. Training Data Synthetic

In the image collection process, it may cost a lot to collect enough image pairs that have local differences to train a detection network. In other words, realistic data doesn't have enough data that contain “objects”. On the one hand, changed book cover image pairs are in the minority, for in most cases, the book cover is not changed when printed. On the other hand, in this spotting the difference task, manually annotation makes neglect easily, as human eyes get tired soon.

To solve the lack of training data, we synthetic different book pairs from same ones, and the result shows the model trained on synthetic data performs well on realistic test set data. Some train set data is shown in Figure 4. Synthetic details are as follows.

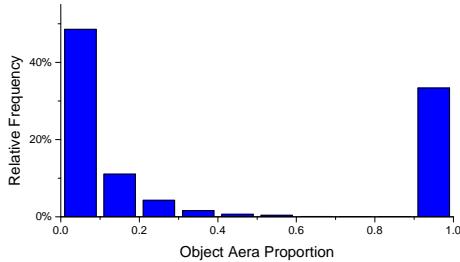


Figure 5: Object(difference) area distribution in synthetic training data. Small differences which are less than 10% of the image area occupies the most, and big differences as large as the whole image are the second.

About 5k same image pairs are collected. We synthesize different image pairs from these totally same pairs, of which two-thirds are locally different and one third are globally different.

To synthesize locally different pairs, we choose a random image from a pair, either the photographic image or digital design. Then we sample a random rectangle patch in the selected image and copy it somewhere else in the same image. Some disturbances like margin blur and small-range scaling is implemented on the copied patch to make it seems nature. Thus, the new patch in the image is a difference or an object we need to detect.

We make sure the copied patch is different from the where it stuck to by simply compare the color histogram between the two patches. Otherwise, there would be many incorrectly labeled objects in the training data. After implementing the comparison of color histogram, there are still some wrongly labeled objects which are not easily discovered by the view of eyes, but the amount is significant reduced. In this way, tiny amount of noise data has little effect on the training.

For globally different image pairs, we choose a photographic image, then randomly choose miss-matched digital design, and resize it the same as the photographic image. These two images form a pair, and the object to detect is labeled as large as the whole image.

In this way, We synthesize only one object in every image pair, so we have a lot of objects or foregrounds to train, but backgrounds are not augmented more than the original 5k same image pairs. In order to make the model learn about backgrounds better, we still need to collect more data in the future.

We count the object size distribution as shown in figure 5. Except for globally different pairs, we don't control the distribution of object area on purpose. We can see small objects which are less than 10% of the whole image occupy the majority. The detection model trained only on synthetic

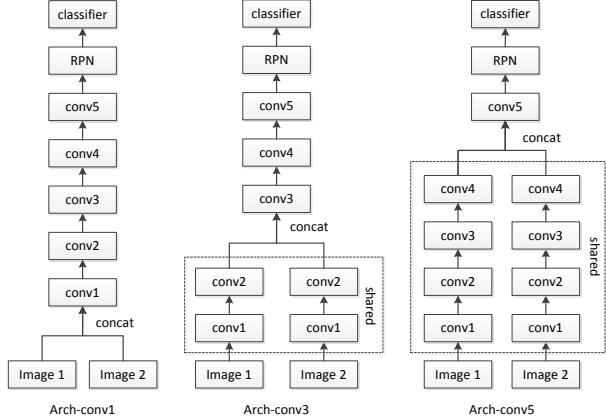


Figure 6: Three architectures of our method. Conv1 to conv5 represent layer blocks in ZF net. ReLU, LRN and Pooling layers are omitted here. In the three architectures, two network branches are concatenated at different layers. Left: the Early-merging architecture. Middle: two branches are concatenated before conv3. Right: two branches are concatenated before conv5.

data generalizes well not only in synthetic validation data, but also on realistic data with various sizes of objects or multiple objects, as shown in Figure 10. It is noted that the synthetic validation data is not participated in the training process.

3.3. Early-merging Architecture

The main idea of spotting the differences by object detection is to consider the two RGB 3-channel input images as one 6-channel image, and consider the differences as objects to detect. We use Faster R-CNN [17] as the detection pipeline, i.e., “CNN feature extraction + region proposal + classification”. Input layer is modified to accept two images and two branches of data are combined by a Concat layer. The idea is shown in Figure 3, and the Early-merging network is shown in Figure 6. The CNN feature extraction network is ZF net [28] which has five convolutional layers, and two image data layers are fed into a Concat layer immediately.

3.4. Two-branch Architecture

Considering the physical significance of two images, we design two-branches architecture. Two input images are fed into several Conv layers respectively, then are concatenated after several layers. Layers before the Concat layer share the same parameters. As shown in Figure 6 middle and right, conv1 to conv5 represent layer blocks in ZF net, as other layers like ReLU, LRN and Pooling is omitted. Explain in detail, in arch-conv3, channels concatenate after

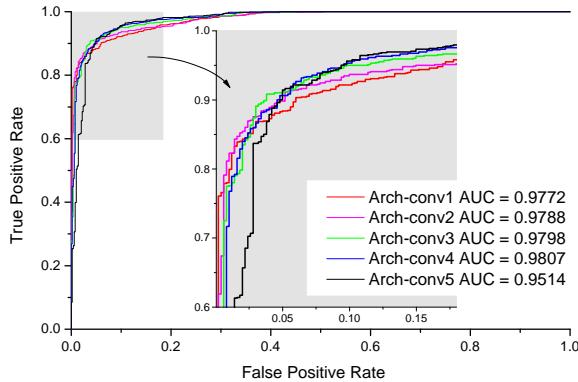


Figure 7: ROC Curves of five architectures. The label “Arch-conv1” represents two images are concatenated before conv1; “Arch-conv2” represents two branches are concatenated before conv2, and so on.

pool2 and before conv3 layer; in arch-conv5, channels concatenate after relu5 and before conv5.

4. Experiments

4.1. Dataset and evaluation protocol

A sample in the dataset consists of two images – a real book cover photographic image and its digital design, and is labeled simple “same” or “different”, without coordinate of bounding box. Realistic book cover image is photographed and then aligned with the digital design. Test set contains 1483 pairs of images, of which 982 pairs are same and 501 pairs different. In the different set, 332 pairs are globally different and 169 pairs are locally different.

For evaluation, we consider same image pairs as “positive” and different pairs as “negative”. The evaluation measure is *Receiving Operating Characteristic Curve* (ROC Curve). As in practical applications, different pairs are more sensitive, which means if you mistakenly put the different pair as the same, it may cause big trouble like customer complaints; on the contrary, if you mistakenly put the same pair as the different, it may only increase some labor cost. So in our evaluation, not only the *area under ROC Curve* (AUC) is considered, but also the slope of the curve.

4.2. Object Detection Method

4.2.1 Training details

We adopt Faster R-CNN [17] as the object detection framework, with ZF net [28] as CNN feature extraction layers. Two images are inputted, and then fed into a Concat layer immediately or after several layers (Figure 6). We initialize Conv layers before Concat layer by pretrained model from

ImageNet, and Conv layers after the Concat layer with the “xavier” method [7]. It is noted the experiment results show it doesn’t help to initialize layers after Concat layer by pre-trained model from ImageNet. In both training and testing, we use single-scale method, where the image’s shorter side is $s = 600$ pixels. The horizontally flipping is implemented as training data augmentation.

The network is trained end-to-end as described in [17] as approximate joint training. The hyper-parameters for training Faster R-CNN are almost the same in [17]. We set global learning rate as 0.001 at the beginning. After 10 epochs, we lower the learning rate to 0.0001 to train for 4 epochs.

In evaluation, as the testing set images are labeled binary “same” or “different”, we consider the highest object score as the distance between the two images. In other words, the larger object score means higher possibility the two images are different, lower score means the two images are same. Then we consider same image pairs as “positive” and draw ROC Curve to analysis as described in section 4.1.

4.2.2 Amount of synthetic training data

There are about 5k same image pairs in total. We synthesize training data as described in 3.2. We experiment on different amount of synthetic training data to explore if more synthetic training data is better. The result is shown in Figure 8 (a). The tag means synthetic magnification(TODO:check) to the original 5k image pairs, as “1/4” means the amount of synthetic training data is one fourth of the 5k image pairs, “2” means twice as the 5k image pairs, and so on.

The ROC Curve shows that there is a tendency that more synthetic training data leads to better test performance, but when the magnification larger than 1, the increase becomes slow, even become worse when magnification larger than 6. The result shows it helps a little to synthetic training data more than the original 5k same image pairs. The reason is that we can synthesize infinite objects if we want, but the amount of background can not be expanded beyond these 5k pairs. In order to make the model learn about backgrounds better, we still need to collect more data in the future.

4.2.3 Two-branch Architecture

We experiment different architectures that two branches are concatenated after several layers. More than shown in Figure 6, we experiment five architectures, trying all the Concat layer locations before five Convolution Layers. The amount of synthetic training data is twice the original 5k.

The performance of the five architectures is in Figure 7. The AUCs of five architectures are fairly close, among which Arch-conv4 has a small advantage. But if we look into the slope of the curve, we find the Arch-conv1(Early-mearing architecture) close to the y axis most, which means

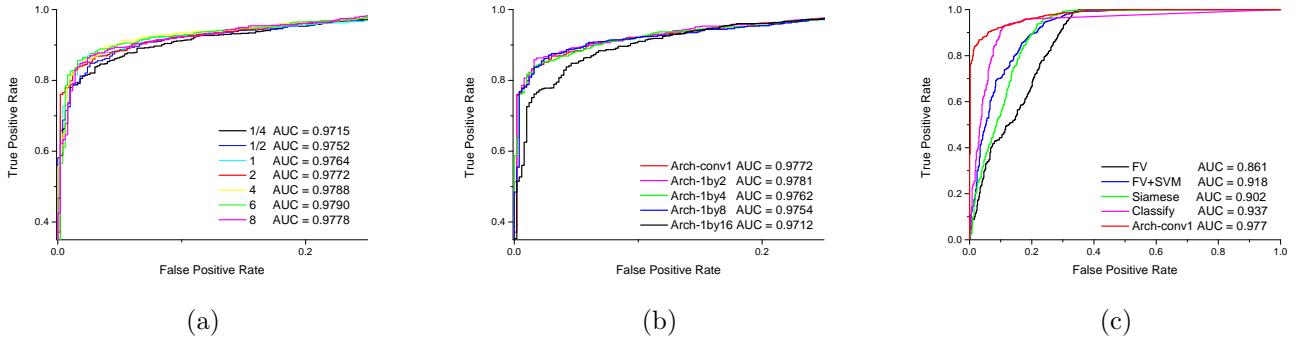


Figure 8: ROC Curves of three experiments. (a) Amount of synthetic training data. (b) Simplified networks. (c) Object detection method (Arch-conv1) compared with other methods.

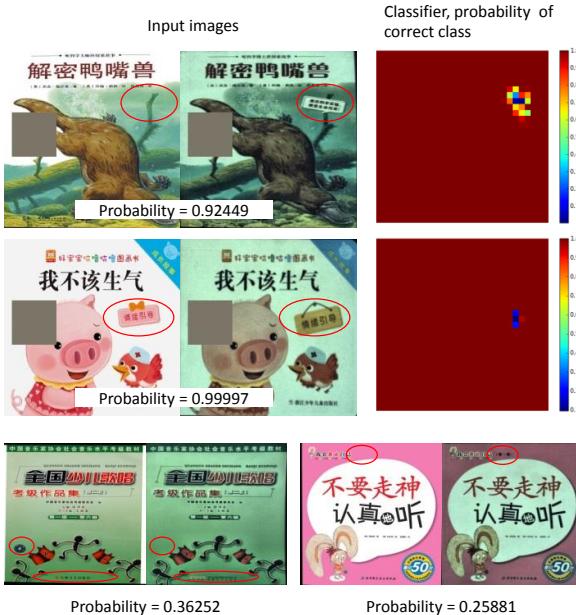


Figure 9: 6-Channel classification network shows attention to the feature of “differnce” and predicts relative large differences. The differences are marked with red ellipses. The visualization is a map of probability as a function of the position of the gray square [28]. However, small differences are still ignored by 6-channel network as shown in the bottom.

in lower false positive rate, Arch-conv1 can reach higher true positive rate. In other words, if the application prefers to rarely miss different image pair, Arch-conv1 may be the best one.

Besides, Arch-conv1 has the least parameters and least computational cost, for layers with same Convolution Layer name in all the architectures have same channels, except for

Table 1: Params and MAC of architectures

Architecture	Params	MAC
Arch-conv1	3.74M	17.9G
Arch-conv2	4.35M	23.82G
Arch-conv3	5.24M	26.04G
Arch-conv4	6.57M	29.35G
Arch-conv5	7.45M	31.56G
Arch-1by2	941.86k	5.55G
Arch-1by4	239k	1.92G
Arch-1by8	61.52k	748.22M
Arch-1by16	16.16k	320.93M

the Concat layer doubles the channel.

Table 1 shows the parameters and theoretical computational cost, which is given as the number of adds and multiplications(MAC) of five architectures. As the architectures are exactly the same from conv5 to the end, we only count the number of parameters and MAC from beginning to conv5. The size of input image is assumed as 1000*600.

This result is consistent with [27], of which experiments show good flexibility when relations between two images are learned from the beginning convolutional layer.

4.2.4 Model Simplification

As we consider the difference between two images as an object, it is an object with much simpler features, compared to objects like human, cat, vehicle, etc. In this situation, it is possible that simpler network with fewer parameters still has a good performance. We tried some simplified networks.

Based on Arch-conv1, we simply reduce channel of layers from beginning to conv5, rpn_conv, fc6 and fc7 to half (Arch-1by2) or quarter (Arch-1by4), and so on. Training details are the same as above. The ROC Curves (Figure 8 (b)) show that these simplified networks perform quite close

to the Arch-conv1, until we reduce the channel to one-sixteenth there appears a relatively large decline.

The parameters and MAC of these simplified architectures are shown in Table 1. It is noted that the layers after conv5 are not counted. Compared to Arch-conv1, the Arhc-1by8 has about 24 times fewer MAC, but still achieves close performance on test set data.

4.3. Comparison with other methods

The evaluation compares the object detection method with four other baselines. The methods are described in details below and we analysis at the end.

4.3.1 Fisher Vector

Before deep learning is widely used in image recognition, one of common methods is to encode dense features of an image to a single feature vector, and the vector represents the image. Fisher Vector(FV) encoding has been used in image recognition tasks and had a good performance [19, 20]. Following a classic recognition example in [24], we extract dense SIFT feature and encode it to a feature vector by FV, then compare Euclidean distance between two feature vectors of two images.

4.3.2 FV + SVM

After extract feature vectors of the two images by FV encoding, which is the same as the above, we carry out a SVM as diagonal metric learning method. We use a conventional linear SVM, and the feature fed in SVM is the vector of squared difference between the two compared FVs. As for train set data, we use all the 5k same image pairs we have collected and the same amount of different image pairs. Different images pairs includes about 1k realistic different image pairs, and the rest is random selected in all the synthetic training data described in 3.2.

4.3.3 Siamese Architecture

The Siamese architecture is widely used in object verification tasks, such as face verification [3], signature verification [2]. We use ZF net as parameter-shared convolutional network. 4096-dimension vector is extracted and fed into the loss function of *ContrastiveLoss*². The train set data is the same as the FV + SVM described above.

4.3.4 6-Channel Classification

In [27], to compare similarity of a pair of images, the input layer is a 2-channel image which is stacked by two gray-scale images, then followed a series of convolutional lay-

ers, and finally a fully connected linear decision layer. Similar to [27], we concatenate two 3-channel images as a 6-channel image and feed it to a ZF classification network.

4.3.5 Analysis

Figure 8 (c) shows the ROC Curves of different methods. The curve of 6-channel classification method is tagged as “Classify”. As we can see, our method has a comprehensive advantage.

Verification methods including FV, FV + SVM and Siamese can find out globally different image pairs precisely, but can not distinguish locally different image pairs. As we describe in Section 2.1, in verification solution outline, features are extracted to represent the image and then fed into a classifier. In Fisher Vector method, hand-crafted features are extracted and encoded, and Euclidean distances between two feature vectors represents the non-similarity between two images. Distance representation is learned in FV + SVM. In Siamese method, not only distance representation but also image feature is learned by training. These methods consider features among an image globally, and a small locally pattern’s feature is extracted indifferently among the whole image. So the comparison of a global feature vector can not distinguish image pairs that have slightly local differences.

6-Channel classification shows relatively large advantage. The relation between two images are leaned from the beginning, and we can see the classification network shows attention to the feature “difference” as shown in Figure 9. However, small but important differences are still ignored by the 6-channel network.

On the contrary, our method using object detection method to consider local information. The feature of “difference” is learned via CNN, and regions of possibly local differences are proposed by Region Proposal Network in Faster R-CNN. So in our object detection method, features of locally differences will not fade away in a whole image.

5. Discussion

Figure 10 shows some examples of qualitative results from Arch-conv1. We can see, the differences including characters and patterning can be detected precisely. Besides, the noise is ignored, including illumination, color, noisy points and slight stains. This is because these kinds of noises are treated as background in the training data. Correspondingly, the real differences, or objects, are treated as foreground. In the training processing, features of foreground are learned, as well as background.

However, there are also some false or missed detection. Figure 11 shows some false detection examples. Actually, in visual point of view, there are some differences between two images, but in application they should be divided into

²<http://caffe.berkeleyvision.org/tutorial/layers/contrastiveloss.html>



Figure 10: Some qualitative results. We can see the detection model trained on synthetic data performs well on differences of character, patterning, book girdle, and global design, no matter single or multiple objects in a pair. The samples on first row are synthetic validation data and others are realistic data.



Figure 11: Some false detection examples. In the left pair, some characters are printed in reflecting material; In the middle, a piece of characters moves a small distance; In the right, an icon is not printed but pasted on the book, which has not exactly the same pose or location as in the digital design. Technically, there are some differences in visual view, but in business application these differences are not important.

the noise category. In the left, some characters are printed in reflecting material, so the photographic image looked quite

different from the digital design. In the middle, a piece of patterning moves a small distance. Technically, this is



Figure 12: Some missed detection examples. In the left pair, the difference is too unapparent. In the middle, the reason is that the black margin is trained as background. In the right, two images are not aligned well.

a kind of difference, but is not very important in selling books. Similarly, in the right, an icon is not printed but pasted on the book, which has not exactly the same pose or location as in the digital design.

In these false detection cases, reflecting printing, small movements and stuck icons should be some kind of background, similar to illumination or color noise, but actually they are closer to differences to detect. At present, these cases account for a tiny proportion in the whole dataset. If they cause big problem in the future, which means there would be enough data of these cases, we have an improvement idea, that is, to treat the reflecting printing, small movement and stuck icon as other kinds of objects. In this way, there would be four classes of objects to detect, including real difference.

Some image pairs do have differences, but are not detected, as shown in Figure 12. In the left, the difference is too unapparent that even human eye can't find it easily. In the middle, the photographic image is somehow not completely and missed part is filled with black pixel. This case is not detected, because a lot image have black margins and the black margin is trained as background. In the right, the difference between two images is quite obvious but has a low detection score. The reason is that the two images are not aligned well, and this kind of data is fairly few in training set.

Besides verifying and detecting the changes between book cover digital design and its photographic image, our proposed method can generalize to similar applications, such as document change, places change, object change and so on.

6. Conclusion

We have presented an object detection method for the task of detecting changes between two images. Our method is fast, accurate and robust while using very cheap annotation. The method stacks two 3-channel images as one 6-channel image and treats the differences as objects. Faster

R-CNN is used to detect the differences. We design the Early-merging architecture that merges two branches of image data at the beginning of the network. Compared with other architectures, Early-merging performs better and faster by sharing convolution computation between the two branches. Besides, our method is robust. After model compression, our method can run 24 times faster while keeping the accuracy. To avoid expensive annotation, all the detection training data are synthesized from the weakly labeled image pairs without bounding boxes. Even so, the trained model can spot the difference effectively on both synthetic data and realistic data. We verify the proposed method on a real task that detects the difference between the digital design and photographic image of a book. Compared with the verification based methods, including Fisher Vector coding, Siamese, and 6-channel classification, our detection based method achieves the best performance.

References

- [1] S. Bianco, G. Ciocca, and R. Schettini. How far can you get by combining change detection algorithms? *arXiv preprint arXiv:1505.02921*, 2015. [3](#)
- [2] J. Bromley, J. W. Bentz, L. Bottou, I. Guyon, Y. LeCun, C. Moore, E. Säckinger, and R. Shah. Signature verification using a "siamese" time delay neural network. *IJPRAI*, 7(4):669–688, 1993. [7](#)
- [3] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pages 539–546. IEEE, 2005. [2, 7](#)
- [4] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pages 886–893. IEEE, 2005. [3](#)
- [5] R. Girshick. Fast r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1440–1448, 2015. [3](#)
- [6] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic

- segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014. 3
- [7] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Aistats*, volume 9, pages 249–256, 2010. 5
- [8] M. Hussain, D. Chen, A. Cheng, H. Wei, and D. Stanley. Change detection from remotely sensed images: From pixel-based to object-based approaches. *ISPRS Journal of Photogrammetry and Remote Sensing*, 80:91–106, 2013. 3
- [9] R. Jain and D. Doermann. Visualdiff: Document image verification and change detection. In *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*, pages 40–44. IEEE, 2013. 3
- [10] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *European Conference on Computer Vision*, pages 21–37. Springer, 2016. 3
- [11] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004. 3
- [12] R. Plamondon and G. Lorette. Automatic signature verification and writer identificationthe state of the art. *Pattern recognition*, 22(2):107–131, 1989. 2
- [13] H. Qin, J. Yan, X. Li, and X. Hu. Joint training of cascaded cnn for face detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3456–3465, 2016. 3
- [14] R. J. Radke, S. Andra, O. Al-Kofahi, and B. Roysam. Image change detection algorithms: a systematic survey. *IEEE transactions on image processing*, 14(3):294–307, 2005. 3
- [15] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 779–788, 2016. 3
- [16] J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. *arXiv preprint arXiv:1612.08242*, 2016. 3
- [17] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 2, 3, 4, 5
- [18] K. Sakurada, T. Okatani, and K. Deguchi. Detecting changes in 3d structure of a scene from multi-view images captured by a vehicle-mounted camera. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 137–144, 2013. 3
- [19] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek. Image classification with the fisher vector: Theory and practice. *International journal of computer vision*, 105(3):222–245, 2013. 7
- [20] K. Simonyan, O. M. Parkhi, A. Vedaldi, and A. Zisserman. Fisher vector faces in the wild. In *BMVC*, volume 2, page 4, 2013. 7
- [21] C. Song, C. E. Woodcock, K. C. Seto, M. P. Lenney, and S. A. Macomber. Classification and change detection using landsat tm data: When and how to correct atmospheric effects? *Remote sensing of Environment*, 75(2):230–244, 2001. 3
- [22] Y. Sun, X. Wang, and X. Tang. Deeply learned face representations are sparse, selective, and robust. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2892–2900, 2015. 1, 2
- [23] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014. 2
- [24] A. Vedaldi and B. Fulkerson. Vlfeat: An open and portable library of computer vision algorithms (2008), 2012. 7
- [25] V. Walter. Object-based classification of remote sensing data for change detection. *ISPRS Journal of photogrammetry and remote sensing*, 58(3):225–238, 2004. 3
- [26] Y. Wang, P.-M. Jodoin, F. Porikli, J. Konrad, Y. Benetech, and P. Ishwar. Cdnet 2014: an expanded change detection benchmark dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 387–394, 2014. 3
- [27] S. Zagoruyko and N. Komodakis. Learning to compare image patches via convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4353–4361, 2015. 6, 7
- [28] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014. 4, 5, 6
- [29] Y. Zhu, C. Yao, and X. Bai. Scene text detection and recognition: Recent advances and future trends. *Frontiers of Computer Science*, 10(1):19–36, 2016. 3