

Article

Object Detection Using Improved Bi-Directional Feature Pyramid Network

Tran Ngoc Quang, Seunghyun Lee  and Byung Cheol Song * 

Department of Electronic Engineering, Inha University, Incheon 22212, Korea;
quangtrandn93@gmail.com (T.N.Q.); lsh910703@gmail.com (S.L.)

* Correspondence: bcsong@inha.ac.kr

Abstract: Conventional single-stage object detectors have been able to efficiently detect objects of various sizes using a feature pyramid network. However, because they adopt a too simple manner of aggregating feature maps, they cannot avoid performance degradation due to information loss. To solve this problem, this paper proposes a new framework for single-stage object detection. The proposed aggregation scheme introduces two independent modules to extract global and local information. First, the global information extractor is designed so that each feature vector can reflect the information of the entire image through a non-local neural network (NLNN). Next, the local information extractor aggregates each feature map more effectively through the improved bi-directional network. The proposed method can achieve better performance than the existing single-stage object detection methods by providing improved feature maps to the detection heads. For example, the proposed method shows 1.6% higher average precision (AP) than the efficient featurized image pyramid network (EFIPNet) for the MicroSoft Common Objects in COntext (MS COCO) dataset.

Keywords: object detection; non-local neural network; feature pyramid network



Citation: Quang, T.N.; Lee, S.; Song, B.C. Object Detection Using Improved Bi-Directional Feature Pyramid Network. *Electronics* **2021**, *10*, 746. <https://doi.org/10.3390/electronics10060746>

Academic Editor: Yoichi Hayashi

Received: 16 January 2021

Accepted: 16 March 2021

Published: 22 March 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

For a long time, object detection has been one of the major topics in the computer vision field, because it is highly utilized in various applications. Recent deep learning-based object detection methods have two mainstreams, i.e., two-stage [1,2] and single-stage approaches [3,4]. The two-stage model is an early methodology that combines a region proposal network (RPN) that estimates the location of an object and a network that classifies the estimated object. It has received attention because it provides much higher performance than existing detection techniques. However, the two-stage model inherently requires considerable computational cost, so a single-stage model was born. The single-stage model extracts feature maps at several stages of convolutional neural networks (CNNs) and concurrently estimates location and class through detection heads. Thus, it provides a more efficient structure than the two-stage model. Because the single-stage model is suitable for real edge devices, it has recently attracted great attention.

The latest single-stage object detection techniques are mainly based on the feature pyramid network (FPN) [4–11]. FPN contributed to improving the performance of existing single-stage methods by aggregating information of feature maps of different sizes in a top-down or bottom-up manner. However, previous FPN-based techniques still have a problem with information loss occurring during the aggregation process, because they adopt a somewhat naive architecture [6,9,11]. In detail, aggregated feature maps do not have enough multi-scale feature information, which degrades performance. Meanwhile, a few approaches to improving performance by delivering global feature information to FPN have been reported [7,8,10]. Since they used a shallow CNN to extract global information, they have a problem with the extracted information still having local characteristics. Therefore,

we claim that a more sophisticated network architecture is required to meet the purpose of each part of the proposed method in order to further improve the performance of FPN.

To solve the above-mentioned problems, this paper proposes advanced modules to improve detection performance. The proposed algorithm is composed of aggregation architecture and a global information extraction model (GIEM). First, the proposed aggregation architecture is designed to minimize information loss by properly fusing the aggregation modules of existing FPN-based techniques. Second, the proposed global information extraction model is constructed so that each point of the feature maps has global information through a non-local neural network (NLNN). Since the proposed method is more advanced than the conventional schemes, it can effectively provide information necessary for detection. As a result, the proposed detector can achieve high performance at a reasonable cost. The experimental results proved that the proposed method is superior to the existing techniques for PASCAL Visual Object Classes (VOC) and MS COCO datasets. For example, in the MS COCO dataset, the proposed method provides 1.6% higher AP than EFIPNet [10], a state-of-the-art (SOTA) technique.

The contribution points of the proposed method are as follows.

- This paper proposes an improved Bi-FPN that minimizes information loss by combining aggregation models of existing FPN-based single object detectors.
- GIEM introduces NLNN to understand the relation between the feature vector and the entire feature map, i.e., obtain global information.

2. Related Work

2.1. Multi-Scale Feature Representation

As for the existing single-stage object detection techniques [4–6,9,11–15], it is most common to aggregate feature maps of various sizes. The lower the level in a feature pyramid, the larger the feature map, so semantic information tends to be weak, but geometric information tends to be strong. On the other hand, as the level increases, the opposite tendency exists. Thus, if feature maps of multiple levels are properly aggregated, detection performance can be improved, because diverse information can be utilized together. Note that iterative connections (IC) are mainly adopted as a way to aggregate feature maps [6,9,12]. However, IC aggregates feature maps in a simple weighted sum way. Thus, the geometric and semantic information of each feature map is not properly aggregated, resulting in information loss. In particular, the disadvantage is that the detection performance of small objects is not satisfactory.

2.2. Global Information Extraction Network

Recently, many methods of utilizing additional information to improve the performance of single-stage object detectors have been proposed. Among them, a few methods that employ the entire image information, that is, global information, are attracting attention because of their relatively high performance [7,8,10]. They first resize an input image to a small resolution version and then extract features through CNNs. Next, local and global information are fused by injecting the extracted global information into each level of the feature pyramid. The authors of [7,10] have a disadvantage in that they cannot extract feature maps of sufficiently various sizes inherently because global information is obtained only by shallow CNNs. To mitigate this problem, the authors of [8] adopted convolutional layers with three different dilations. However, since the above-mentioned methods simply extract feature maps through convolutional layers, the extracted feature maps are still quite dependent on local characteristics.

3. Methods

In order to improve the existing naïve aggregation framework that causes information loss, this paper proposes a new single-stage object detector. Figure 1 shows the overall structure. VGG-16 [16] was used as the backbone network. The local information extraction model (LIEM) extracted multi-scale object context from feature maps (conv4_3, fc7, conv8_2,

conv9_2) of the backbone network. At the same time, the global information extraction model (GIEM) extracted global information from down-sampled images. Finally, the aggregation network fused the feature maps extracted by LIEM and GIEM. As a result, the proposed method can improve classification and regression performance for objects of various sizes. From the following sections, each part of the proposed method is described in detail.

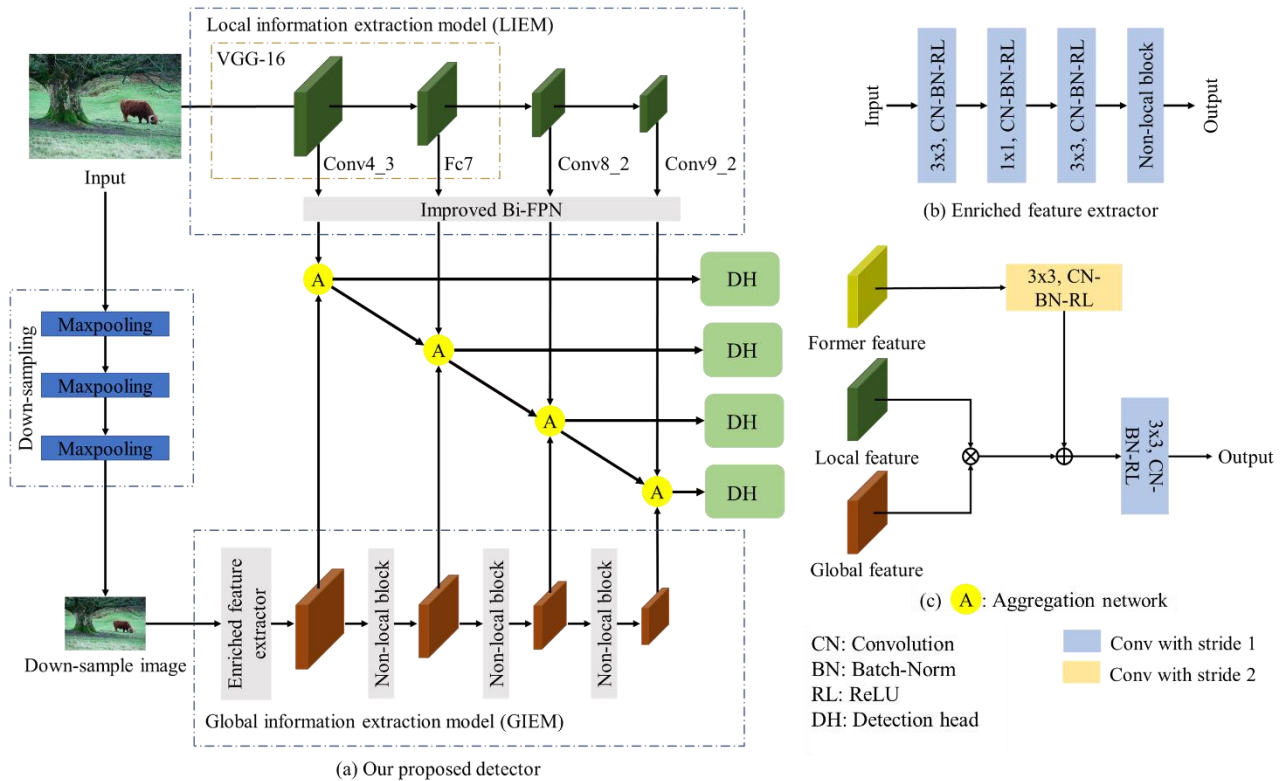


Figure 1. (a) The overall block diagram of the proposed method. Here, the detection head follows the method of [6]. (b) The enriched feature extractor used in a global information extraction model (GIEM). (c) The aggregation network that combines feature maps extracted by a local information extraction model (LIEM) and GIEM. Here, the first aggregation network uses only local features and global features.

3.1. Local Information Extraction Model (LIEM)

This section describes the LIEM part of the proposed method. A total of four levels of feature maps are generated by the backbone network VGG-16 and extra layers as in Figure 1. Then, the improved bi-directional FPN (IBi-FPN) extracts local information from these feature maps (see Figure 2). Next, feature maps are iteratively fused in top-down and bottom-up manners. When fusing feature maps of different levels, the following weighted sum is used.

$$WS(f_1, f_2) = Conv\left(\frac{w_1 * f_1 + w_2 * Resize(f_2)}{w_1 + w_2 + \epsilon}\right), \quad (1)$$

where f_1 and f_2 are two input feature maps, w_1 and w_2 are trained weights, and ϵ is set to 10^{-4} . $Conv$ is composed of a convolutional layer with 3×3 kernel and stride of 1, batch normalization, and ReLU. $Resize$ is a function to match the size of feature maps of two levels and usually uses bilinear interpolation or a max pooling layer. The fused feature maps have richer context and location information than the original feature maps. However, due to such a fusion, features necessary for each level are smoothened, so objects of a specific size may not be detected. To prevent this problem, the Maxout function (MO) extracts

important information from the original feature map and the fused feature map through the following equation.

$$Lc^{mx} = MO(f_1, f_2), \tag{2}$$

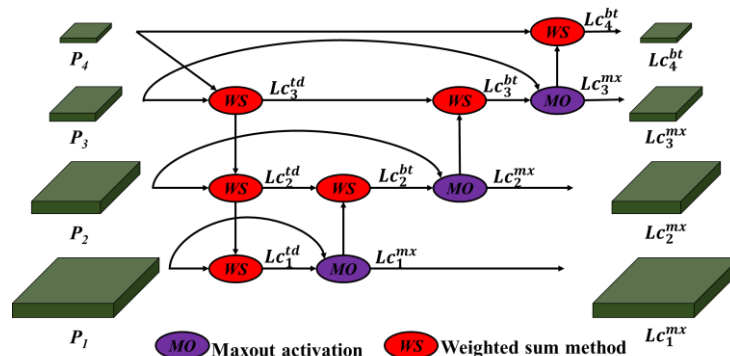


Figure 2. The improved bi-directional feature pyramid network (Bi-FPN).

Here, the smallest feature map is used as it is because it has little effect on fusion. Finally, feature maps extracted from LIEM are expressed as follows.

$$Lc = \{Lc_1^{mx}, Lc_2^{mx}, Lc_3^{mx}, Lc_4^{bt}\}, \tag{3}$$

where *mx* and *bt* indicate *Maxout* and bottom-up manner fusion, respectively. Therefore, the proposed LIEM not only allows feature maps to have richer context and location information but also allows each level to retain important information originally possessed by the Maxout function.

3.2. Global Information Extraction Model (GIEM)

This section describes GIEM to obtain global information. First, in order to lower computational complexity, an input image is down-sampled. In order to reduce the computational cost in down-sampling, we employ three max pooling layers rather than a single large max pooling layer. Next, the down-sampled feature map is transformed into an enriched feature map (EFM) by three convolutional layers and one non-local block. This process is expressed by the following equation.

$$Glb_1 = \varphi(I), \tag{4}$$

where φ indicates a combination of three convolutional layers, batch normalization, ReLU, and one non-local block. I denotes a down-sampled image. Figure 3 shows the non-local block used here. A non-local block (NLB) contributes to extracting global information of a given feature map by analyzing the relation between information in the feature map. Unlike a conventional NLB, the proposed NLB separates g and (Φ, θ) . This separation contributes to more focus on finding relations and improving features. The g function generates richer feature maps through a squeeze-and-excitation block. In addition, in order to obtain feature maps of the same size as the feature map with local information obtained beforehand, the size of the feature map is repeatedly reduced through a non-local block. This is expressed as the following equation.

$$Glb_i = \partial(Glb_k), \tag{5}$$

where k is an index with $k = 1, \dots, n-1$ th. Additionally, ∂ denotes one non-local block. The output of GIEM is defined by

$$Glb = \{Glb_1, Glb_2, Glb_3, Glb_4\}, \tag{6}$$

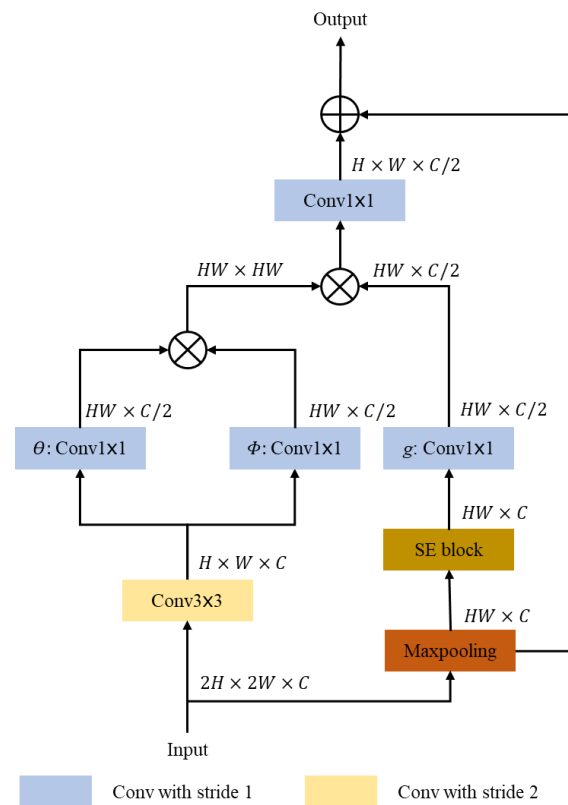


Figure 3. The architecture of a non-local block. A feature map has a feature dimension of $H \times W \times C$. H , W , and C indicate height, width, and channel, respectively. \oplus denotes element-wise addition and \otimes denotes matrix multiplication. SE denotes a Squeeze Excitation block.

The proposed GIEM allows pixels of each feature map to improve information through global relation through a non-local block. Therefore, the proposed GIEM provides richer information than the previous techniques.

3.3. Aggregation Network

This section describes how to aggregate feature maps obtained from the proposed LIEM and GIEM. The proposed aggregation network is constructed based on [8] (see Figure 1c). In this paper, the function for aggregating local and global features was determined heuristically, which is represented by

$$f_i = Conv_i^1(Lc_i \otimes Glb_i) \oplus Conv_i^2(f_{i-1}), \tag{7}$$

where i is an index with $i = 1, \dots, n$ and $Conv^s$ operation is composed of a convolution (kernel 3×3 , stride s), BatchNorm, and ReLU layers. Further, \otimes and \oplus indicate element-wise multiplication and addition, respectively. Therefore, the information obtained from the proposed LIEM and GIEM is aggregated in this way, which minimizes information loss and improves context and location information.

4. Experiments

4.1. Datasets

PASCAL VOC [17] is a dataset with 20 classes. We constructed a training set composed of 16,000 images from the trainval of VOC2007 and VOC2012 and trained the network with the dataset. Additionally, it was evaluated with the VOC2007 test set consisting of 5000 images. The threshold of intersection over union (IOU) for calculating mean average precision (mAP) was set to 0.5.

MS COCO [18] is a dataset with 80 classes. A network was trained with a trainval dataset composed of 120,000 images and was evaluated with a test set composed of 40,000 images. For the evaluation with the MS COCO dataset, well-known pycocotools were used, and the average precision (AP) was calculated by averaging the results for the IOU thresholds in the range of 0.5 to 0.95.

4.2. Implementation Details

The backbone network VGG-16 was pre-trained on ImageNet [19]. For the MS COCO dataset, the initial learning rate was set to 2×10^{-3} and decreased by 10% at 90 and 120 epochs. For the PASCAL VOC dataset, the learning rate was set to 2×10^{-3} and gradually decreased to 10^{-4} and 2×10^{-5} at 150 and 200 epochs, respectively. We deployed the same settings as the standard SSD [3], where we used the same loss function, scales, and aspect ratios of the default's boxes and the data augmentation method. Further, the weight decay was set to 0.0005, the momentum was set to 0.9, and the batch size was set to 32. We considered two input sizes, i.e., 300×300 and 320×320 .

4.3. Evaluation Results

PASCAL VOC 2007: We compared the proposed method with various single- and two-stage methods. Table 1 shows the comparison result on the PASCAL VOC 2007 test set. Input image size and performance figures of other techniques are cited as they are in the paper. In the case of the 300×300 input size and VGG-16 backbone network, the proposed method provides a significant performance gain of about 2% over SSD. However, the proposed method provides a little bit lower performance than SOTA models such as EFIPNet [10] and deep feature pyramid reconfiguration (DFPR) [20]. On the other hand, compared with different backbone networks, the proposed method provides a gain of 1.02% in terms of mAP over gated feature reuse-deeply supervised object detectors (GFR-DSOD) [9]. Although the proposed method effectively utilizes the information of feature maps to improve the performance of SSD, its performance is slightly less than the SOTA technique's. Despite the slightly insufficient performance, we have obtained insight that if the components of the proposed method are fused with other algorithms, they can contribute to improving the performance of the algorithms. Detailed analysis is covered in the ablation study. Further, in the case of the 320×320 input size, the proposed object detector achieved mAP of 80.37, which is a significant gain of 2.02% over SSD. However, the proposed scheme shows 0.23% and 1.00% lower performance than dynamic anchor feature selection (DAFS) [21] and enriched feature guided refinement network (EFGRNet), respectively. In general, the performance of object detection techniques is highly dependent on datasets. In Table 1, the performance of the proposed method is inferior to that of EFIPNet, but in Table 2 regarding MS COCO, the proposed method is superior. The processing speed of the algorithm normally relies on the computer hardware. In fact, the computing power used by EFIPNet is stronger than that of ours. This leads to a difference in speed. The authors urge the reviewer to consider this tendency.

Table 1. Comparison of the proposed method and the conventional models on the PASCAL VOC 2007 dataset.

Model	Backbone	Input Size	mAP	FPS
Two-Stage Detectors:				
Faster RCNN [2]	ResNet101	1000 × 600	76.4	-
R-FCN [22]	ResNet101	1000 × 600	80.5	-
Single-Stage Detectors:				
CHFANet [23]	VGG16	300 × 300	79.9	-
FERNet [24]	VGG16	300 × 300	80.2	-
EFIPNet [10]	VGG16	300 × 300	80.4	111
SFDet [25]	VGG16	300 × 300	78.8	-
RFBNet [26]	VGG16	300 × 300	78.6	-
DES [27]	VGG16	300 × 300	79.7	76.8
DFPR [20]	VGG16	300 × 300	80	-
GFR-DSOD [9]	DSOD	300 × 300	78.9	-
SSD [3]	VGG16	300 × 300	77.93	64.23
IBi-FPN (Ours)	VGG16	300 × 300	79.92	45.68
EFGRNet [7]	VGG16	320 × 320	81.37	44.4
DAFS [21]	VGG16	320 × 320	80.6	-
RefineDet [28]	VGG16	320 × 320	80	40.3
DSSD [29]	ResNet101	321 × 321	78.6	9.5
GFR-DSOD [9]	DSOD	320 × 320	79.2	-
RON [30]	VGG16	320 × 320	76.6	-
SSD [3]	VGG16	320 × 320	78.35	-
IBi-FPN (Ours)	VGG16	320 × 320	80.37	44

Table 2. Comparison of the proposed method and the conventional models on the MS COCO dataset.

Methods	Input Size	Backbone	FPS	AP (0.5:0.95)	AP (0.5)	AP (0.75)	APs (0.5:0.95)	APm (0.5:0.95)	API (0.5:0.95)
FAENet [31]	300 × 300	VGG-16	-	28.3	47.9	29.7	10.5	30.9	41.9
SSD [3]	300 × 300	VGG-16	50	25.1	43.1	25.8	6.6	25.9	47.6
LSNet [8]	300 × 300	VGG-16	76.9	32	51.5	33.8	12.6	34.9	47
EFIPNet [10]	300 × 300	VGG-16	71.4	30	48.8	31.7	10.9	32.8	46.3
IBi-FPN (Ours)	300 × 300	VGG-16	45.43	31.6	50.3	33	10.7	36.9	47.7
EFGRNet [7]	320 × 320	VGG-16	47.62	33.2	53.4	35.4	13.4	37.1	47.9
DAFS [21]	320 × 320	VGG-16	46	31.2	50.8	33.4	10.8	34	47.1
PASSD [13]	320 × 320	VGG-16	40	31.4	51.6	33.6	12.0	35.1	45.8
M2Det [5]	320 × 320	VGG-16	33.4	33.5	52.4	35.6	14.4	37.6	47.6
IBi-FPN (Ours)	320 × 320	VGG-16	42.01	32	50.9	33.7	12	37	47.1

Figure 4 shows an example of detecting an object in an image of PASCAL VOC 2007 using the proposed method and SSD, respectively. We find that the proposed method shows better detection performance than the SSD by making good use of local and global information extracted from the SSD. In particular, the proposed method has qualitatively excellent performance by detecting objects that even the SSD cannot detect.



Figure 4. An example of object detection for PASCAL VOC dataset. We can see that the proposed method (a) can detect more objects than SSD (b).

MS COCO: Table 2 compares the proposed method with the latest single-stage models on MS COCO dataset. In the case of the 300×300 input size, the proposed object detector provides a gain of 6.5% over SSD. In the case of the 320×320 input size, the proposed object provides a gain of 0.8% over DAFS [29]. The proposed scheme shows 1.2% and 1.3% lower performance than EFGRNet [7] and M2Det, respectively. However, note that the proposed method has a fast speed of 42.01 FPS for 320×320 .

4.4. Ablation Study and Discussion

Through the following ablation study, each part of the proposed method is described in detail. To verify the performance of each module of the proposed method, Table 3 shows the results of the ablation study on the PASCAL VOC2007 dataset. Firstly, we compared the proposed LIEM with standard SSD. The proposed LIEM provided a gain of 1.49% in terms of mAP. Thus, we can find that the Maxout for fusion in the proposed method effectively fuses more features than the weighted sum of the existing BiFPN. In the case of GIEM, conventional NLNN and the proposed non-local block were compared. Both GIEMs effectively extracted global information and improved detection performance. In particular, the proposed non-local block shows 0.21% higher performance than NLNN, which proves that the proposed technique is effective. Therefore, we can say that every module of the proposed framework has a meaningful effect on the performance of object detection.

Table 3. The effect of each model of improved bi-directional FPN (IBi-FPN) on performance for the PASCAL VOC 2007 dataset.

Standard SSD	Methods				mAP
	BiFPN	LIEM	GIEM		
			NLNN	Split Branch	
✓	-	-	-	-	78.35
-	✓	-	-	-	79.53
-	-	✓	-	-	79.84
-	-	✓	✓	-	80.16
-	-	✓	✓	✓	80.37

It is true that the performance of the proposed method is somewhat inferior to that of the SOTA technique. In other words, since the proposed method requires a certain amount of additional computation, it shows a slightly lower FPS than the SOTA scheme with similar performance. Nevertheless, it is worth noting that the proposed method gives significant intuition to the object detection field. That is to say, a meaningful way to improve the existing naïve approach was suggested. Therefore, if a light-weighting

technique is applied to the proposed method, a computationally efficient object detector will be developed.

5. Conclusions

We propose a new single-stage object detection framework which contains LIEM, GIEM, and an aggregation network. LIEM based on a bi-directional model extracts high context and location information, and GIEM captures rich global features using non-local blocks. Additionally, the aggregation network generates informative feature maps by aggregating local and global features. Experimental results show that the proposed object detector achieves the desirable AP performance of 31.6% on MS COCO dataset, which is better than the existing single- and two-stage approaches. For example, the proposed method shows 1.6% higher AP than EFIPNet, which is a significant step forward in the object detection task. If further performance improvement is made in the future, it will be possible to apply the proposed method to the industry. Moreover, the processing speed of our object detector is 45.43 FPS, which is suitable for the real-time application of object detection. However, it is true that the processing speed of the proposed method still does not reach SOTA. Therefore, further improving the computational efficiency is our future research task.

Author Contributions: Conceptualization, T.N.Q. and S.L.; methodology, T.N.Q. and S.L.; software, T.N.Q.; validation, T.N.Q. and S.L.; investigation, T.N.Q. and S.L.; resources, T.N.Q. and S.L.; data curation, T.N.Q. and S.L.; writing—original draft preparation, T.N.Q. and S.L.; writing—review and editing, T.N.Q., S.L., and B.C.S. All authors have read and agreed to the published version of the manuscript.

Funding: This work was supported by the research grant of Inha University.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems*; Neural Information Processing Systems Foundation Inc.: San Diego, CA, USA, 2015; pp. 91–99.
2. He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Las Vegas, NV, USA, 27–30 June 2016; pp. 770–778.
3. Liu, W.; Anguelov, D.; Erhan, D.; Szegedy, C.; Reed, S.; Fu, C.Y.; Berg, A.C. Ssd: Single shot multibox detector. In *European Conference on Computer Vision*; Springer: Cham, Switzerland, 2016; pp. 21–37.
4. Lin, T.Y.; Dollár, P.; Girshick, R.; He, K.; Hariharan, B.; Belongie, S. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, HI, USA, 21–26 July 2017; pp. 2117–2125.
5. Zhao, Q.; Sheng, T.; Wang, Y.; Tang, Z.; Chen, Y.; Cai, L.; Ling, H. M2det: A single-shot object detector based on multi-level feature pyramid network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Honolulu, HI, USA, 27 January–1 February 2019; Volume 33, pp. 9259–9266.
6. Tan, M.; Pang, R.; Le, Q.V. Efficientdet: Scalable and efficient object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Seattle, WA, USA, 13–19 June 2020; pp. 10781–10790.
7. Nie, J.; Anwer, R.M.; Cholakkal, H.; Khan, F.S.; Pang, Y.; Shao, L. Enriched feature guided refinement network for object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, Seoul, Korea, 27 October–2 November 2019; pp. 9537–9546.
8. Wang, T.; Anwer, R.M.; Cholakkal, H.; Khan, F.S.; Pang, Y.; Shao, L. Learning rich features at high-speed for single-shot object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, Seoul, Korea, 27 October–2 November 2019; pp. 1971–1980.
9. Shen, Z.; Shi, H.; Yu, J.; Phan, H.; Feris, R.; Cao, L.; Liu, D.; Wang, X.; Huang, T.; Savvides, M. Improving object detection from scratch via gated feature reuse. *arXiv* **2017**, arXiv:1712.00886.
10. Pang, Y.; Wang, T.; Anwer, R.M.; Khan, F.S.; Shao, L. Efficient featurized image pyramid network for single shot detector. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Long Beach, CA, USA, 16–20 June 2019; pp. 7336–7344.
11. Pang, J.; Chen, K.; Shi, J.; Feng, H.; Ouyang, W.; Lin, D. Libra r-cnn: Towards balanced learning for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Long Beach, CA, USA, 16–20 June 2019; pp. 821–830.

12. Li, Y.; Pang, Y.; Shen, J.; Cao, J.; Shao, L. NETNet: Neighbor Erasing and Transferring Network for Better Single Shot Object Detection. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Seattle, WA, USA, 14–19 June 2020; pp. 13349–13358.
13. Jang, H.D.; Woo, S.; Benz, P.; Park, J.; Kweon, I.S. Propose-and-attend single shot detector. In Proceedings of the IEEE Winter Conference on Applications of Computer Vision, Snowmass Village, CO, USA, 1–5 March 2020; pp. 815–824.
14. Ghiasi, G.; Lin, T.Y.; Le, Q.V. Nas-fpn: Learning scalable feature pyramid architecture for object detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Long Beach, CA, USA, 16–20 June 2019; pp. 7036–7045.
15. Liu, S.; Qi, L.; Qin, H.; Shi, J.; Jia, J. Path aggregation network for instance segmentation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–22 June 2018; pp. 8759–8768.
16. Simonyan, K.; Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv* **2014**, arXiv:1409.1556.
17. Everingham, M.; van Gool, L.; Williams, C.K.; Winn, J.; Zisserman, A. The pascal visual object classes (voc) challenge. *Int. J. Comput. Vis.* **2010**, *88*, 303–338. [[CrossRef](#)]
18. Lin, T.Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; Zitnick, C.L. Microsoft Coco: Common Objects in Context. In *European Conference on Computer Vision*; Springer: Cham, Switzerland, 2014; pp. 740–755.
19. Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; et al. ImageNet Large Scale Visual Recognition Challenge. *Int. J. Comput. Vis.* **2015**, *115*, 211–252. [[CrossRef](#)]
20. Kong, T.; Sun, F.; Tan, C.; Liu, H.; Huang, W. Deep feature pyramid reconfiguration for object detection. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, 8–14 September 2018; pp. 169–185.
21. Li, S.; Yang, L.; Huang, J.; Hua, X.S.; Zhang, L. Dynamic anchor feature selection for single-shot object detection. In Proceedings of the IEEE International Conference on Computer Vision, Seoul, Korea, 27 October–2 November 2019; pp. 6609–6618.
22. Dai, J.; Li, Y.; He, K.; Sun, J. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in Neural Information Processing Systems*; Neural Information Processing Systems Foundation Inc.: San Diego, CA, USA, 2016; pp. 379–387.
23. Xu, X.; Luo, X.; Ma, L. Context-Aware Hierarchical Feature Attention Network for Multi-Scale Object Detection. In Proceedings of the 2020 IEEE International Conference on Image Processing (ICIP), Abu Dhabi, United Arab Emirates, 25–28 October 2020; pp. 2011–2015.
24. Fan, B.; Chen, W.; Cong, Y.; Tian, J. Dual Refinement Underwater Object Detection Network. *Constr. Side Channel Anal. Secur. Des.* **2020**, 275–291. [[CrossRef](#)]
25. Antioquia, A.M.C.; Tan, D.S.; Azcarraga, A.; Hua, K.L. Single-Fusion Detector: Towards Faster Multi-Scale Object Detection. In Proceedings of the 2019 IEEE International Conference on Image Processing (ICIP), Taipei, Taiwan, 22–25 September 2019; pp. 76–80.
26. Liu, S.; Huang, D. Receptive field block net for accurate and fast object detection. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, 8–14 September 2018; pp. 385–400.
27. Zhang, Z.; Qiao, S.; Xie, C.; Shen, W.; Wang, B.; Yuille, A.L. Single-shot object detection with enriched semantics. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–22 June 2018; pp. 5813–5821.
28. Zhang, S.; Wen, L.; Bian, X.; Lei, Z.; Li, S.Z. Single-shot refinement neural network for object detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–22 June 2018; pp. 4203–4212.
29. Fu, C.Y.; Liu, W.; Ranga, A.; Tyagi, A.; Berg, A.C. Dssd: Deconvolutional single shot detector. *arXiv* **2017**, arXiv:1701.06659.
30. Kong, T.; Sun, F.; Yao, A.; Liu, H.; Lu, M.; Chen, Y. Ron: Reverse connection with objectness prior networks for object detection. In Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 5936–5944.
31. Li, W.; Liu, G. A Single-Shot Object Detector with Feature Aggregation and Enhancement. In Proceedings of the 2019 IEEE International Conference on Image Processing (ICIP), Taipei, Taiwan, 22–25 September 2019; pp. 3910–3914.