

# SU-YOLO: Spiking Neural Network for Efficient Underwater Object Detection

Chenyang Li<sup>a,b,1</sup>, Wenxuan Liu<sup>c,b,1</sup>, Guoqiang Gong<sup>a,\*</sup>, Xiaobo Ding<sup>a</sup>,  
Xian Zhong<sup>b,d</sup>

<sup>a</sup>College of Computer and Information Technology, China Three Gorges University,  
Yichang, 443002, China

<sup>b</sup>Hubei Key Laboratory of Transportation Internet of Things, School of Computer Science  
and Artificial Intelligence, Wuhan University of Technology, Wuhan, 430070, China

<sup>c</sup>State Key Laboratory for Multimedia Information Processing, School of Computer  
Science, Peking University, Beijing, 100091, China

<sup>d</sup>State Key Laboratory of Maritime Technology and Safety, Wuhan University of  
Technology, Wuhan, 430063, China

---

## Abstract

Underwater object detection is critical for oceanic research and industrial safety inspections. However, the complex optical environment and the limited resources of underwater equipment pose significant challenges to achieving high accuracy and low power consumption. To address these issues, we propose Spiking Underwater YOLO (SU-YOLO), a Spiking Neural Network (SNN) model. Leveraging the lightweight and energy-efficient properties of SNNs, SU-YOLO incorporates a novel spike-based underwater image denoising method based solely on integer addition, which enhances the quality of feature maps with minimal computational overhead. In addition, we introduce Separated Batch Normalization (SeBN), a technique that normalizes feature maps independently across multiple time steps and is optimized for integration with residual structures to capture the temporal dynamics of SNNs more effectively. The redesigned spiking residual blocks integrate the Cross Stage Partial Network (CSPNet) with the YOLO architecture to mitigate

---

\*Corresponding author.

*Email addresses:* [chenyang@ctgu.edu.cn](mailto:chenyang@ctgu.edu.cn) (Chenyang Li), [lwxfight@163.com](mailto:lwxfight@163.com) (Wenxuan Liu), [guoqiang\\_gong@163.com](mailto:guoqiang_gong@163.com) (Guoqiang Gong), [84479265@qq.com](mailto:84479265@qq.com) (Xiaobo Ding), [zhongx@whut.edu.cn](mailto:zhongx@whut.edu.cn) (Xian Zhong)

<sup>1</sup>Contributed equally to this work.

spike degradation and enhance the model’s feature extraction capabilities. Experimental results on URPC2019 underwater dataset demonstrate that SU-YOLO achieves mAP of 78.8% with 6.97M parameters and an energy consumption of 2.98 mJ, surpassing mainstream SNN models in both detection accuracy and computational efficiency. These results underscore the potential of SNNs for engineering applications. The code is available in <https://github.com/lwxfight/snn-underwater>.

*Keywords:*

Spiking neural networks, Underwater object detection, Image denoising

---

## 1. Introduction

Underwater exploration has diverse applications, including aquatic life detection, facility safety inspections, and salvage operations. With advancements in remotely operated vehicles (ROVs) and autonomous underwater vehicles (AUVs), unmanned exploration methods are increasingly replacing traditional manual approaches, in which object detection plays a critical role. However, the complex underwater optical environment, characterized by low image brightness and high noise levels [1], renders underwater object detection significantly more challenging than general object detection. In addition, the size and power constraints of ROVs and AUVs necessitate lightweight detection models to minimize resource consumption, further complicating the task. These factors make underwater object detection a demanding research area in computer vision.

Several algorithms optimized for underwater object detection have been proposed, such as YOLOv9s-UI [2], which demonstrate accuracy improvements over baseline models. However, incorporating attention mechanisms in these models often increases computational complexity, rendering them unsuitable for resource-constrained devices like AUVs. Lightweight models [3, 4] have been introduced to address this issue by reducing parameter counts and computational demands. Unfortunately, these simplifications typically result in lower detection accuracy compared to state-of-the-art general object detection models.

Spiking neural networks (SNNs) offer a promising approach to achieving both low computation and high accuracy. As third-generation neural networks [5, 6, 7], SNNs are inspired by the neuronal communication mechanisms in the brain. Unlike conventional artificial neural networks (ANNs),

SNNs transmit information via discrete spikes rather than continuous values, which significantly reduces computational cost and energy consumption. Recent studies have explored applying SNNs to object detection, including Spiking-YOLO [8], EMS-YOLO [9], and SpikeYOLO [10]. However, these works primarily focus on general object detection and lack optimizations for underwater conditions, such as image preprocessing and noise suppression, thereby limiting their performance in underwater environments.

To address these limitations, we propose SU-YOLO, an SNN-based underwater object detection model built on the “you only look once” (YOLO) framework [11]. Given that YOLOv9 [12] currently stands as a stable and high-performance object detection framework, SU-YOLO adopts its glean network as the baseline. First, to accommodate the energy constraints of underwater devices, we simplify the network structure and convert it into a spiking architecture to achieve minimal power consumption. During this process, we found that the conventional ResNet [13] structure commonly used in SNNs leads to spike degradation even in shallow networks; consequently, we redesigned the lightweight SU-Blocks based on the CSPNet [14] structure as the fundamental unit. Second, recognizing that existing SNNs lack effective methods for image noise removal, we developed a novel spiking underwater image denoising module. This module processes feature maps rather than the original image, removing both positive and negative noise to improve detection accuracy. In addition, existing batch normalization (BN) methods either ignore the temporal dynamics of SNNs [15, 16], thereby losing potential information, or disregard the residual structure in networks [17, 18], leading to excessive neuron charging and activation. To overcome these issues, we propose an optimized normalization method.

To the best of our knowledge, the proposed SU-YOLO is the first SNN specifically designed for underwater object detection. As shown in Fig. 1, our approach significantly enhances SNN performance in underwater object detection, achieving an  $mAP_{0.5}$  of 78.8% on URPC2019 with 6.97M parameters and an energy consumption of 2.98 mJ, exceeding the performance of existing lightweight SNN and several ANN models.

The main contributions of this work are fourfold:

- We develop a lightweight SNN model based on the YOLO framework that supports direct training while effectively mitigating spike degradation. The model outperforms existing SNNs and several ANNs in underwater object detection, achieving optimal accuracy with only four time steps.

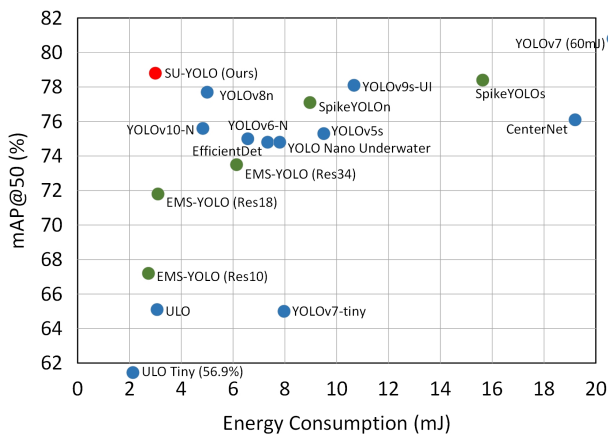


Figure 1: **Comparison of detection performance and energy consumption across various models on URPC2019.** The blue dots represent ANN models, while the green and red dots indicate SNN models.

- We design a novel spiking-compatible underwater image denoising method that integrates seamlessly into SNN architectures. This method requires only integer addition calculations, significantly improving feature map quality with minimal computational overhead.
- We propose an improved batch normalization method for SNNs, termed SeBN, which concatenates feature maps across time steps along the channel dimension for BN operations. SeBN is optimized for residual structures and outperforms classical BN and alternative methods, such as tdBN.
- Experiments on URPC2019 and UDD underwater image datasets validate the feasibility and superior performance of SNNs for underwater object detection.

## 2. Related Work

### 2.1. SNN Object Detection

Object detection is a fundamental problem in computer vision. Since the introduction of R-CNN [19], ANNs have dominated deep learning-based object detection methods. The development of Spiking-YOLO [8] marked the first successful application of SNNs to object detection. Spiking-YOLO



employs an ANN-SNN conversion method [20], representing activation values through spike firing rates by converting a pre-trained ANN into an SNN. However, this approach requires more than 1000 time steps to accurately represent the original values, which diminishes the computational efficiency of SNNs. Methods such as FSHNN [21] and Spike Calibration [22] mitigate this issue by reducing the required time steps to 300 through optimized network structures. More recently, Fast-SNN [23] introduces a signed IF neuron model and a layer-wise fine-tuning mechanism to address quantization and accumulation errors in ANN-SNN conversion, further compressing the time steps to 3. However, the use of negative spike computation in Fast-SNN compromises the binary nature of spikes, resulting in a non-standard spiking network structure.

An alternative to conversion is to directly train SNNs [24], though this approach poses challenges due to the non-differentiability of the spike firing function. Lee *et al.* [25] addressed this by treating the membrane potential as a differentiable parameter, which allows for gradient calculation and the development of deep, trainable SNNs. Subsequent surrogate gradient methods [26], using functions such as the sigmoid and arctangent, have further facilitated SNN training. Notable examples include the spatiotemporal backpropagation algorithm (STBP) [27] and the spike-based backpropagation algorithm [28].

Advancements in training methods have spurred research on SNN-based object detection. EMS-YOLO [9] employs a fully spiking ResNet [13], addressing issues in prior non-spike-based computations observed in MS-ResNet [29] and SEW-ResNet [30], and demonstrates robust performance in general object detection. Similarly, SpikeYOLO [10] introduces the I-LIF spiking neuron to convert non-spike-based computations into spike-based computations, thereby improving detection accuracy. However, these approaches are primarily designed for general object detection and perform suboptimally in underwater scenarios due to challenges such as blurry objects and severe noise.

## 2.2. Underwater Object Detection

Underwater object detection is particularly challenging because of the complex underwater environment and the performance limitations of underwater equipment. Several studies have attempted to overcome these challenges. For example, Song *et al.* [31] developed an algorithm based on Mask R-CNN [32] that employs multi-scale retinal enhancement and transfer learning to improve detection accuracy. These two-stage object detection methods, however, typically incur high computational costs.

In contrast, single-stage detection algorithms, such as the YOLO series, feature simpler architectures and lower computational requirements, making them more suitable for resource-constrained underwater devices. Zhang *et al.* [33] proposed a lightweight underwater detection network based on YOLOv4 [34], which incorporates attention-based feature fusion to achieve accuracy comparable to the original model with fewer parameters. Further improvements were made by Zhang *et al.* [35] on YOLOv5 using newly designed blocks and MLLE image enhancement, surpassing the performance of YOLOv7 [36] and YOLOv8 [37]. TC-YOLO [38], an extension of YOLOv5s, integrates the Transformer [39] self-attention mechanism, significantly enhancing feature extraction for underwater objects. Despite these advances, the incorporation of attention mechanisms increases computational load, limiting their applicability in resource-constrained settings. Lightweight models such as ULO [4] reduce parameter counts and computational costs (3.94M parameters and 3.42 GFLOPs) but still achieve an mAP of only 65.13% on URPC, which lags behind mainstream object detection models. These limitations motivate our proposal of SU-YOLO, an SNN-based underwater object detection model.

### 2.3. Underwater Image Denoising

Image denoising is critical in underwater object detection. Conventional denoising methods, including mean filtering, Gaussian filtering, and non-local means (NLM), often blur noise and object boundaries, resulting in a loss of fine details. To address these shortcomings, Li *et al.* [40] proposed an adaptive color and contrast enhancement (ACCE) framework combined with denoising, which leverages frequency-domain decomposition and variational optimization to enhance visual quality. You *et al.* [41] employed a wavelet transform-based approach to denoise images prior to edge detection, using multi-resolution analysis to better preserve edge details. Although these methods provide improvements, they often struggle to remove complex noise without significantly blurring the image. Recent deep learning approaches have also been proposed. For example, Tian *et al.* [42] introduced the cross Transformer denoising network (CTNet) that combines Transformer and CNN for image denoising, significant advantages in images with complex noise. Ju *et al.* [43] proposed a patch-based denoising diffusion model to address image blurring and color distortion, which effectively overcomes underwater suspended particles and light absorption. However, these methods typically involve intensive floating-point computations or require running separate

models, making them less compatible with SNN architectures. Therefore, we have designed a dedicated lightweight underwater image denoising method tailored for SNNs.

#### 2.4. Normalization

BN [44] is widely used in ANNs to mitigate internal covariate shift by normalizing input statistics, thereby preventing gradient vanishing or explosion during training. However, for SNNs, which include an additional temporal dimension, conventional BN is less effective. Several improved normalization techniques have been proposed. NeuNorm [15] normalizes input features along the channel dimension before updating membrane potentials, yet it does not fully resolve the gradient vanishing problem. Temporal dimension BN (tdBN) [16] normalizes inputs across all time steps, thus accounting for both temporal and spatial dimensions. The SpikingJelly [45] framework also applies BN by flattening the time dimension into the batch dimension.

Although these methods are promising, our experiments indicate that independently normalizing each time step, rather than merging features across time steps, yields better performance. Approaches such as BNTT [17] and TEBN [18] follow this design but fail to consider local parallel structures, such as shortcut connections, or the need for batch-scale fusion in SNNs. To address these issues, we propose SeBN, an optimized normalization method, which is detailed in subsequent sections.

### 3. Proposed Method

#### 3.1. Spiking Underwater YOLO

Our model is based on the YOLO framework and comprises the SU-Block1, SU-Block2, SpikeSPP, Encoder, SpikeDenoiser, and DetectHead modules, as illustrated in Fig. 2. These components are carefully designed and optimized to enhance detection performance in underwater environments. During the forward pass, normalized RGB image data (scaled to the range  $[0,1]$ ) is fed into the network as floating-point values and then converted to spike outputs via integrate-and-fire (IF) neurons in the Encoder module. The SpikeDenoiser then performs image denoising and downsampling. Subsequently, the resulting feature map is processed by the backbone, which consists of three SU-Block1 residual blocks with increasing channel dimensions. Multi-scale features are extracted using the SpikeSPP module, followed by feature fusion in the Neck, and the DetectHead produces the final predictions.

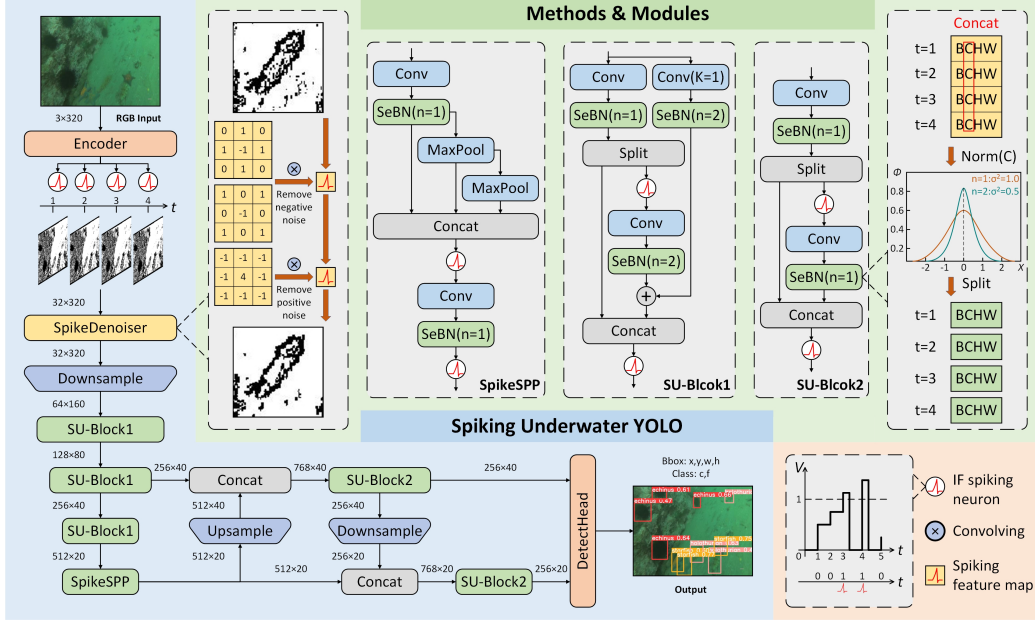


Figure 2: **Diagram of SU-YOLO and the structure of its modules.** In the main structural diagram, narrower graphics represent feature maps with reduced dimensions.

In SU-YOLO, all spiking neurons are implemented as IF neurons. Unlike Leaky IF (LIF) neurons, IF neurons do not exhibit membrane voltage decay over time. Their simpler mathematical formulation and lower computational cost reduce overall power consumption. The IF neuron model is described by:

$$V_i^l(t) = (1 - \Theta_i^l(t)) \left( V_i^l(t-1) + \sum_j (\Theta_j^{l-1}(t) W_{i,j}^l + B_{i,j}^l) \right), \quad (1)$$

where  $V_i^l(t)$  denotes the membrane voltage of the  $i$ -th neuron in layer  $l$  at time  $t$ , while  $W$  and  $B$  represent the weights and biases of the corresponding convolutional layer. The spike indicator  $\Theta$  is defined as:

$$\Theta_i^l(t) = H(V_i^l(t) - V_{th}), \quad (2)$$

with  $V_{th}$  as the threshold voltage and  $H(x)$  being the Heaviside step function:

$$H(x) = \begin{cases} 1, & \text{if } x \geq 0, \\ 0, & \text{if } x < 0. \end{cases} \quad (3)$$

When the accumulated membrane voltage exceeds the threshold, the neuron emits a spike and resets to zero potential, awaiting the next charge.

### 3.2. SpikeDenoiser

Noise severely degrades detection performance in underwater object detection. Existing SNN-based methods [8, 9, 10] target general object detection and lack specialized denoising techniques, rendering them vulnerable to noise interference in underwater environments. To overcome this limitation while preserving the fully spiking nature of the network, we propose a denoising method that operates directly on the feature maps. In SNNs, the binary nature of feature maps causes underwater salt-and-pepper noise to appear as isolated pixels with inverted binary values compared to their neighbors. To effectively extract and correct these isolated points, we employ three convolution filters with the following kernel definitions:

$$K_1 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & -1 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \quad (4)$$

$$K_2 = \begin{bmatrix} 1 & 0 & 1 \\ 0 & -1 & 0 \\ 1 & 0 & 1 \end{bmatrix}, \quad (5)$$

$$K_3 = \begin{bmatrix} -1 & -1 & -1 \\ -1 & 4 & -1 \\ -1 & -1 & -1 \end{bmatrix}. \quad (6)$$

The input feature map, encoded with spikes, is sequentially convolved with  $K_1$ ,  $K_2$ , and  $K_3$ . If the convolution result equals 4, the current pixel is inverted. This process is mathematically expressed as:

$$I_{\text{conv}}(x, y) = \sum_{i=-1}^1 \sum_{j=-1}^1 I_{\text{in}}(x+i, y+j) \times K_n(i+1, j+1), \quad (7)$$

$$I_{\text{out}}(x, y) = \begin{cases} 1 - I_{\text{in}}(x, y), & \text{if } I_{\text{conv}}(x, y) = 4, \\ I_{\text{in}}(x, y), & \text{if } I_{\text{conv}}(x, y) < 4, \end{cases} \quad (8)$$

where  $I_{\text{in}}(x, y)$  is the pixel value at coordinates  $(x, y)$ ,  $I_{\text{conv}}(x, y)$  is the convolution result, and  $I_{\text{out}}(x, y)$  is the output. Here,  $K_n$  ( $n = 1, 2, 3$ ) denotes the corresponding convolution kernel.

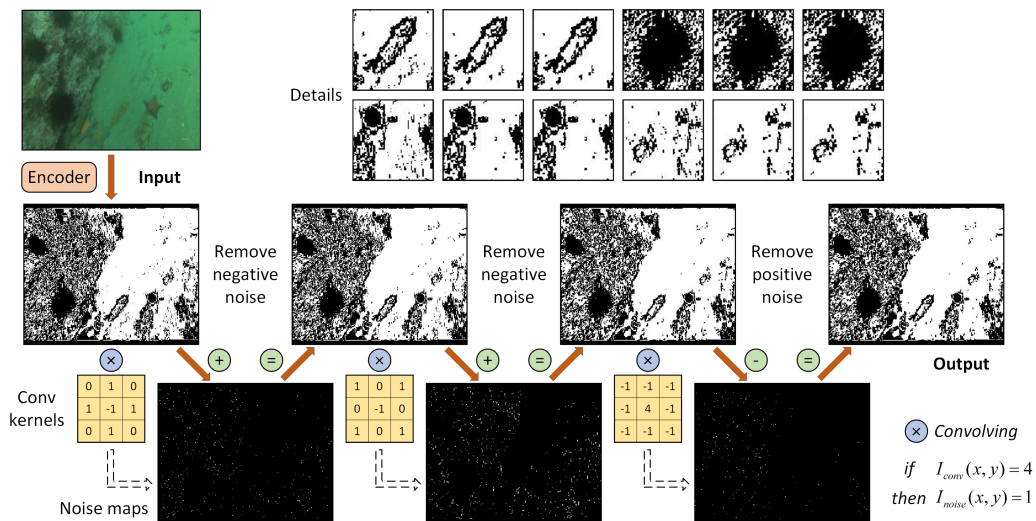


Figure 3: Schematic diagram illustrating the processing flow and effects of the SpikeDenoiser module.

As illustrated in Fig. 3, the method primarily fills missing pixels (Negative Noise) and removes isolated pixels (Positive Noise). Kernels  $K_1$  and  $K_2$  identify and fill pixels where the center is 0 and the surrounding pixels are 1, while  $K_3$  targets and removes isolated pixels where the center is 1 and the surrounding pixels are 0. Experimental results indicate that an aggressive approach for filling Negative Noise, combined with a conservative strategy for removing Positive Noise, yields optimal denoising performance. Notably, the SpikeDenoiser operates only during inference and incurs minimal computational cost, less than 0.1 GSOPs and 0.01 mJ for a  $320 \times 320$  image, while improving mAP by 3.5% on UDD. This method maintains the fully spiking nature of the network since both the input and output feature maps remain binary, and the operations rely solely on integer additions, reducing energy consumption by up to  $9\times$  and  $37\times$  compared to floating-point addition and multiplication, respectively [46]. Its compatibility with various SNN architectures makes it a versatile solution for underwater image denoising.

### 3.3. SU-Block

SU-Block1 and SU-Block2 serve as the fundamental residual blocks in our model. Underwater object detection typically requires shallow, computationally efficient networks due to device performance constraints. Although

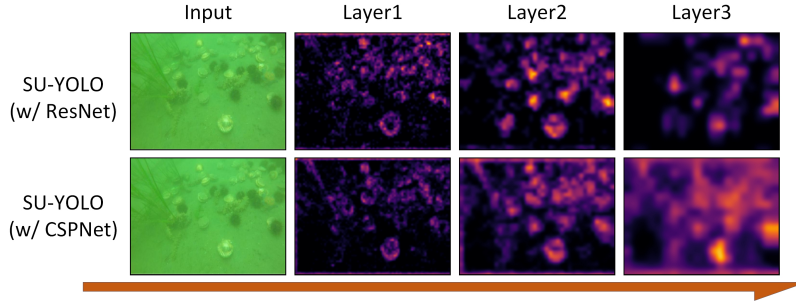


Figure 4: **Feature maps extracted from backbone layers when using ResNet or CSPNet in SNN.**

ResNet is commonly employed in existing SNNs [9, 29, 30], the inherent low quality of underwater images limits the feature extraction capabilities of shallow ResNet architectures, leading to reduced detection accuracy. In contrast, CSPNet (Cross Stage Partial Network) provides enhanced feature representation even under challenging imaging conditions and with shallower network depths, thereby maintaining high accuracy while reducing computational costs. Moreover, we observed that a simple ResNet structure suffers from spike degradation in SNNs (see Fig. 4), as the firing rate decreases with increasing network depth, which can result in the loss of critical object features. CSPNet, however, effectively mitigates this degradation by maintaining robust neuron activity throughout the backbone. Therefore, we designed the SU-Block based on the CSPNet architecture.

SU-Block1, employed in the backbone for feature extraction, downsampling, and channel expansion, integrates CSPNet and shortcut connections (see Fig. 2). It is defined as:

$$X_1^l, X_2^l = \text{Split}(\text{Conv}(X_{\text{out}}^{l-1})), \quad (9)$$

$$X_{\text{out}}^l = \text{IF}(\text{Concat}(X_1^l, \text{Conv}(\text{IF}(X_2^l)) + \text{Conv}(X_{\text{out}}^{l-1}))), \quad (10)$$

where  $X_{\text{out}}^l$  denotes the output of the  $l$ -th layer,  $\text{IF}(\cdot)$  represents the IF neuron activation,  $\text{Conv}(\cdot)$  denotes a convolutional layer,  $\text{Split}(\cdot)$  splits the input equally along the channel dimension, and  $\text{Concat}(\cdot)$  concatenates inputs along the channel dimension.

Following a  $3 \times 3$  convolution, the feature map is split into two groups. One group remains unchanged while the other undergoes an additional  $3 \times 3$  convolution for further feature extraction. The two groups are then concatenated, activated, and output. To compensate for the absence of a residual



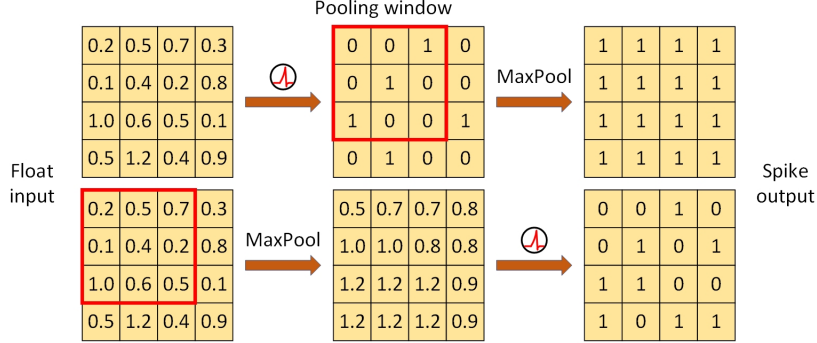


Figure 5: **Impact of pooling and activation sequence on output.**

path in the second group, a shortcut connection is introduced to preserve gradient flow. SU-Block2, used in the Neck for feature fusion, follows a similar design but without downsampling or additional residual connections, thereby enhancing parameter efficiency. It is defined as:

$$X_1^l, X_2^l = \text{Split}(\text{Conv}(X_{\text{out}}^{l-1})), \quad (11)$$

$$X_{\text{out}}^l = \text{IF}(\text{Concat}(X_1^l, \text{Conv}(\text{IF}(X_2^l)))). \quad (12)$$

### 3.4. SpikeSPP

The SpikeSPP module is a simplified version of spatial pyramid pooling (SPP) [47], as depicted in Fig. 5. Its operations are defined by:

$$X_1^l = \text{Conv}(X_{\text{out}}^{l-1}), \quad (13)$$

$$X_2^l = \text{MaxPool}(X_1^l), \quad (14)$$

$$X_3^l = \text{MaxPool}(X_2^l), \quad (15)$$

$$X_{\text{out}}^l = \text{IF}(\text{Conv}(\text{IF}(\text{Concat}(X_1^l, X_2^l, X_3^l)))). \quad (16)$$

where  $\text{MaxPool}(\cdot)$  represents max-pooling, and  $\text{Concat}(\cdot)$  concatenates along the channel dimension. Although SPP provides multi-level pooling for extracting multi-scale features, pooling after activation can lead to significant information loss in SNNs due to the binary nature of spikes. In underwater object detection, where salt-and-pepper noise is prevalent, pooling prior to activation better preserves information. Experimental results further indicate that removing the third max-pooling layer reduces computational load without sacrificing performance. The resulting features are then passed to the IF neurons, which generate spikes with appropriate firing rates.



### 3.5. Encoder and DetectHead

The Encoder module converts floating-point input images into spikes. It applies preliminary feature extraction and normalization using conventional convolution and BN. The resulting feature map is duplicated  $T$  times and passed through spike neurons, yielding binary spike outputs that complete the image encoding process.

The DetectHead adopts the decoupled head design from YOLOv9 [12], modified to handle spike-based inputs and outputs. Separate convolutional layers compute class probabilities and bounding box coordinates. Multi-level feature fusion is achieved by accepting inputs from two SU-Block2 modules, thereby capturing richer features. In the final spiking neuron of the DetectHead, the membrane potential threshold is set to infinity, allowing it to accumulate all floating-point values over each time step, which ensures more accurate predictions. Moreover, the parameter  $n$  in SeBN is set to the number of time steps  $T$ , which reduces the range of single-step outputs and stabilizes the accumulated membrane potential, thereby improving gradient stability. The final network output is obtained by reading the membrane potential of these neurons.

### 3.6. Separated Batch Normalization

In SNNs, each spiking neuron charges and fires independently at each time step, meaning that presynaptic inputs may exhibit different distribution characteristics over time. Normalizing feature maps across all time steps [16] can obscure these differences by forcing them into a single distribution. By assigning distinct normalization parameters for each time step, the network preserves these temporal variations, thereby enhancing its information capacity and expressive power. Additionally, residual network structures require careful adjustment of data distributions prior to summation [17] to prevent neuron overcharging and excessive spike firing. To address these issues, we propose Separated Batch Normalization (SeBN), which normalizes feature maps independently across time steps while optimizing for integration with residual structures. SeBN is defined as:

$$y_{t,k} = \gamma_{t,k} \left( \frac{x_{t,k} - \mu_{t,k}}{\sqrt{n(\sigma_{t,k}^2 + \varepsilon)}} \right) + \beta_{t,k}, \quad (17)$$

where  $x_{t,k}$  is the feature map of the  $k$ -th channel at time step  $t$ ,  $n$  is a positive integer representing the number of SeBN modules connected after the current

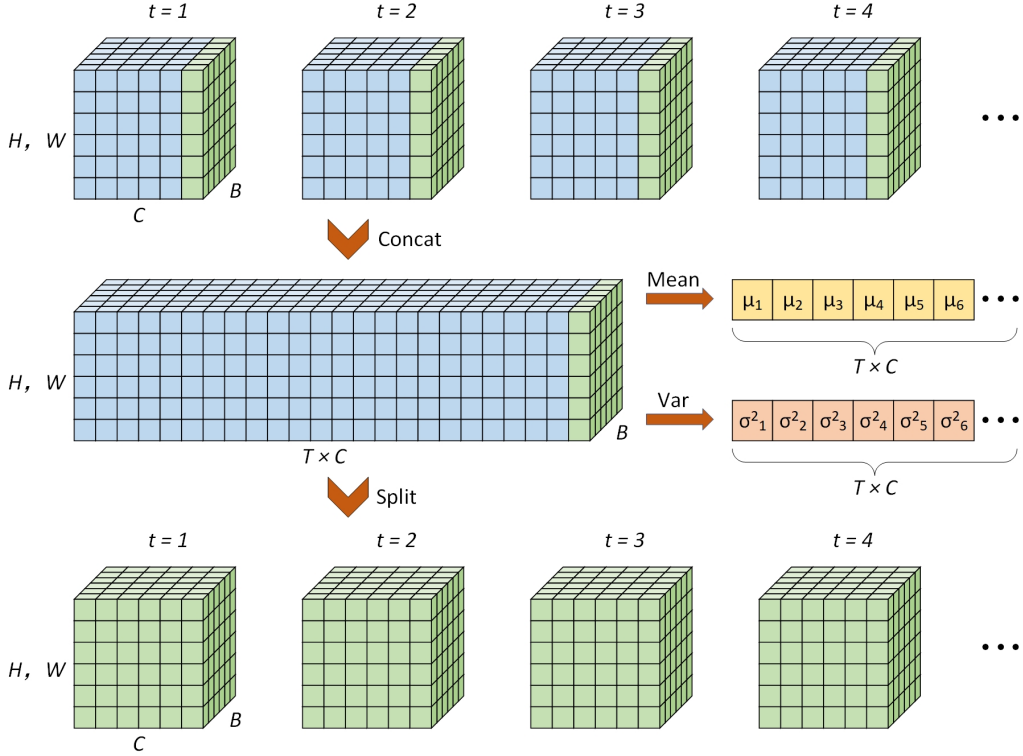


Figure 6: **Calculation and fusion methods of SeBN.** In the SeBN computation diagram, cubes represent the feature maps at each time step, where  $B$  denotes the batch size,  $C$  indicates the number of channels,  $H$  and  $W$  are the height and width of the feature maps, and  $T$  represents the total number of time steps. The symbols  $\mu_k$  and  $\sigma_k^2$  denote the mean and variance of the  $k$ -th feature map.

module, and  $\varepsilon$  is a small constant to prevent division by zero. The parameters  $\gamma_{t,k}$  and  $\beta_{t,k}$  are trainable, and  $y_{t,k}$  is the normalized output. The mean  $\mu_{t,k}$  and variance  $\sigma_{t,k}^2$  are computed as:

$$\mu_{t,k} = \frac{1}{N} \sum_{i=1}^N (x_{t,k})_i, \quad (18)$$

$$\sigma_{t,k}^2 = \frac{1}{N} \sum_{i=1}^N ((x_{t,k})_i - \mu_{t,k})^2, \quad (19)$$

with  $N$  representing the batch size.

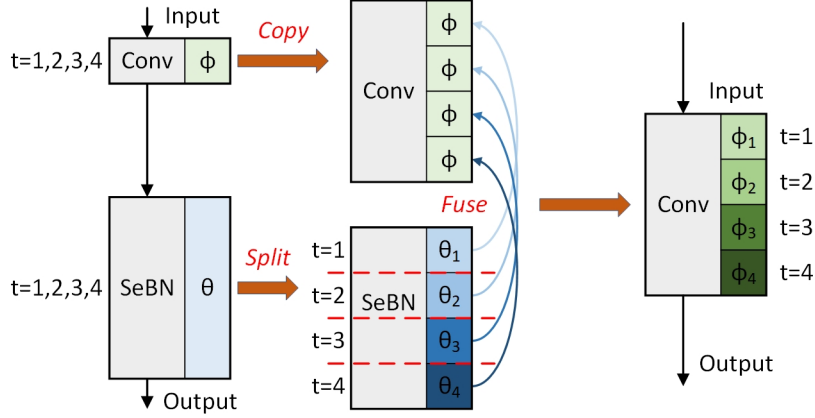


Figure 7: **Process of batch-scale fusion of SeBN.** Here,  $\theta_t$  represents all parameters of the SeBN layer at time step  $t$ ,  $\phi$  represents parameters shared by the convolutional layer across all time steps, and  $\phi_t$  denotes parameters of the convolutional layer at time step  $t$ .

As shown in Fig. 6, feature maps from all  $T$  time steps are concatenated along the channel dimension and normalized using SeBN. The normalized feature maps are then evenly split and redistributed across the  $T$  time steps. In this normalization, the factor  $n$  adjusts the distribution of the input data from  $N(0, 1)$  to  $N(0, \frac{1}{n})$ , ensuring that the sum in residual structures maintains a distribution of  $N(0, 1)$  and preventing excessive neuron activation. During training, running means and variances for each batch across  $T \times C$  are updated via a moving average to compute  $\mu$  and  $\sigma_{\text{run},t,k}^2$ , which are then used during inference.

To maintain the fully spiking nature of the SNN during inference, the normalization layer is removed, and a batch-scale fusion method is introduced for SeBN. As shown in Fig. 7, the weights and biases  $\phi$  shared across all time steps in the original convolutional layer are duplicated for each time step. The parameters for each time step,  $\theta_t$ , are fused with  $\phi$  to generate  $T$  new sets of weights and biases, denoted by  $\phi_t$ . This fusion is performed as follows:

$$W'_{t,k} = \gamma_{t,k} \left( \frac{W_k}{\sqrt{n(\sigma_{\text{run},t,k}^2 + \varepsilon)}} \right), \quad (20)$$

$$B'_{t,k} = \gamma_{t,k} \left( \frac{B_k - \mu_{\text{run},t,k}}{\sqrt{n(\sigma_{\text{run},t,k}^2 + \varepsilon)}} \right) + \beta_{t,k}, \quad (21)$$

where  $W_k$  and  $B_k$  are the weights and biases for the  $k$ -th output channel in the original convolutional layer, and  $\mu_{\text{run},t,k}$  and  $\sigma_{\text{run},t,k}^2$  are the running mean and variance computed using an exponential moving average (EMA) over the training set. The fused weights and biases,  $W'_{t,k}$  and  $B'_{t,k}$ , are then used for each time step during inference, eliminating the need for the SeBN layer and thereby reducing computational load and accelerating inference.

In the fused convolutional layer, each time step utilizes its own dedicated parameters. Although this increases the model’s file size, the computational load and inference time remain unchanged. Given that storage is a relatively inexpensive resource, this trade-off is acceptable considering the significant accuracy improvements achieved by SeBN.

## 4. Experimental Results

### 4.1. Datasets and Implementation Details

#### 4.1.1. Datasets

To evaluate our model’s performance in underwater object detection, we conduct experiments on URPC2019 and UDD datasets. URPC2019 dataset [48], released during the 2019 Underwater Robot Picking Contest, contains 4,707 labeled underwater images spanning four object classes: scallops, starfish, echinus, and holothurian. It is partitioned into a training set with 3,767 images, a validation set with 695 images, and a test set with 245 images. UDD dataset [49] comprises 2,227 labeled underwater images featuring three classes (scallops, holothurian, and echinus), and is divided into a training set of 1,827 images and a validation set of 400 images.

For experiments involving SeBN, we also utilize PASCAL VOC 2012 dataset [50], which consists of 11,540 labeled images across 20 object classes. This dataset is split into 5,717 training images and 5,823 validation images with object detection annotations.

#### 4.1.2. Implementation Details

All experiments are performed on a single NVIDIA RTX 2080Ti (22GB) GPU, paired with an AMD Ryzen 5600 CPU and 32GB of RAM, running Ubuntu 22.04 LTS with Python 3.8 and PyTorch 2.0.0. For training, all neurons in our model are implemented as IF neurons using the SpikingJelly [45] framework with a reset potential  $V_{\text{rst}} = 0$  and a threshold potential  $V_{\text{th}} = 1.0$ . We employ the SGD optimizer with a learning rate of 0.001, a batch size of 16, and an input image size of  $320 \times 320$ . The model is trained for 300 epochs.

#### 4.1.3. Computational Cost

In conventional ANNs, computational cost is measured in FLOPs, which represent the number of floating-point operations (additions or multiplications). However, since SNNs communicate via discrete spikes, calculating FLOPs is less straightforward. Instead, we use synaptic operations (SOPs) as the metric for computational cost. SOPs account for the convolution operations in each layer and the spike emissions from preceding layers, while also considering the floating-point additions involved in neuron charging. The formulas for FLOPs and SOPs are defined as follows:

$$\text{FLOPs} = \sum_{i=1}^N (C_{\text{in},i} \times k_{h,i} \times k_{w,i} \times C_{\text{out},i} \times H_{\text{out},i} \times W_{\text{out},i}), \quad (22)$$

$$\begin{aligned} \text{SOPs} = \sum_{t=1}^T \sum_{i=1}^N & \left( (f_{i-1,t} \times C_{\text{in},i} \times k_{h,i} \times k_{w,i} + 1) \right. \\ & \left. \times C_{\text{out},i} \times H_{\text{out},i} \times W_{\text{out},i} \right), \end{aligned} \quad (23)$$

where  $T$  is the total number of time steps,  $N$  is the total number of network layers,  $f_{i-1,t}$  denotes the firing rate of neurons in the  $(i-1)$ -th layer at time  $t$ ,  $C_{\text{in},i}$  and  $C_{\text{out},i}$  are the numbers of input and output channels of the  $i$ -th convolutional layer,  $k_{w,i}$  and  $k_{h,i}$  are the kernel dimensions, and  $W_{\text{out},i}$  and  $H_{\text{out},i}$  are the output feature map dimensions.

Note that while FLOPs account for both multiplications and additions, SOPs only account for additions, thereby leading to lower energy consumption for the same operation count. For fair comparison, we also estimate energy consumption. Given that in convolutional neural networks the ratio of multiplications to additions is approximately 1:1, we approximate the number of multiply-accumulate operations (MACs) as FLOPs divided by 2. According to estimates in [46], on a typical 45nm chip, a MAC operation consumes about 4.6 pJ and an addition consumes about 0.9 pJ.

#### 4.1.4. Time Step

In SNNs, time steps serve as discrete units of temporal processing, during which neuronal states (*e.g.*, membrane potential and spike generation) and synaptic dynamics are updated. As illustrated in Fig. 8, more time steps result in increased spike emissions, thereby enhancing the network’s expressiveness. However, this also increases computational overhead and reduces inference

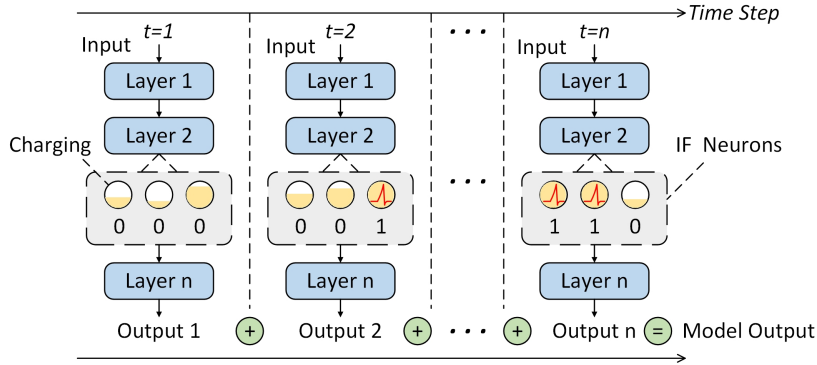


Figure 8: **Running process of SU-YOLO varies with time steps.**

speed. Consequently, for equivalent detection accuracy, methods that require fewer time steps are more advantageous in terms of energy efficiency and real-time performance.

#### 4.2. Comparison Experiments

We compared our model with existing methods on URPC2019 and UDD, with results summarized in Table 1. Our directly trained SNN outperforms models obtained via ANN-SNN conversion while using fewer time steps. Compared to existing directly trained SNNs, our model achieves higher mAP with fewer parameters. Moreover, its performance is comparable to that of lightweight ANNs.

EMS-YOLO [9] currently serves as a benchmark for SNN-based object detection. On URPC2019, the lightest version of EMS-YOLO, based on ResNet10, achieves an  $\text{mAP}_{0.5}$  of 0.672 with an energy consumption of 2.73 mJ, while a larger version using ResNet34 attains an  $\text{mAP}_{0.5}$  of 0.735. SpikeY-OLO [10] further improves performance, reaching an  $\text{mAP}_{0.5}$  of 0.784 with energy consumption of 8.64 mJ when using one time step and four virtual time steps. Our method surpasses these results, achieving an  $\text{mAP}_{0.5}$  of 0.788 with only 6.97M parameters and a computational load of 2.90 GSOPs, while consuming merely 2.98 mJ, significantly outperforming other SNN approaches in underwater object detection.

We also compared our model with several ANNs. Although our mAP is slightly lower than that of advanced models like YOLOv7 [36], our energy consumption is only 5% of YOLOv7’s. When compared to lightweight ANNs with similar energy profiles, our model demonstrates superior perfor-

Model	Params (M)	GFLOPs/ GSOPs	Time (Step)	mAP <sub>0.5</sub>		mAP <sub>0.5:0.95</sub>		Energy (mJ)
				URPC	UDD	URPC	UDD	
ANN								
CenterNet (Res50) [51]	23.51	8.35	-	0.761	-	-	-	19.21
EfficientDet-D3 [52]	12.00	3.19	-	0.748	-	-	-	7.34
YOLOv5s [53]	7.20	4.13	-	0.753	-	0.420	-	9.50
YOLOv6-N [54]	4.70	2.85	-	0.750	0.537	0.410	0.228	6.56
YOLOv7-tiny [36]	6.23	3.46	-	0.650	0.461	0.311	0.196	7.96
YOLOv7 [36]	36.49	25.88	-	0.808	0.533	0.434	0.236	59.52
YOLOv8n [37]	3.20	2.18	-	0.777	0.576	0.436	0.267	5.01
YOLOv10-N [55]	2.71	2.10	-	0.756	0.542	0.425	0.251	4.83
ULO Tiny [4]	3.26	0.93	-	0.569	-	-	-	2.14
ULO [4]	3.76	1.34	-	0.651	-	-	-	3.08
YOLO Nano Underwater [3]	4.33	3.39	-	0.748	-	-	-	7.80
YOLOv9s-UI [2]	4.11	4.63	-	0.781	-	0.487	-	10.65
ANN-SNN								
Spiking-YOLO [8] *	8.68	-	3500	0.537	0.389	0.224	0.138	-
Fast-SNN (Dark19) [23]	20.80	-	15	0.604	0.454	0.250	0.163	-
SNN								
EMS-YOLO (Res10) [9]	6.20	0.45 + 1.88	4	0.672	0.506	0.301	0.194	<b>2.73</b>
EMS-YOLO (Res18) [9]	9.34	0.45 + 2.28	4	0.718	0.529	0.349	0.211	3.09
EMS-YOLO (Res34) [9]	14.40	0.45 + 5.66	4	0.735	0.557	0.369	0.224	6.13
SpikeYOLOn [10]	12.64	0.08 + 11.05	4 × 1	0.735	0.499	0.392	0.220	10.13
SpikeYOLOn [10]	12.64	0.08 + 6.20	1 × 4	0.771	0.532	0.429	0.242	5.76
SpikeYOLOs [10]	22.05	0.11 + 16.61	4 × 1	0.761	0.529	0.418	0.244	15.20
SpikeYOLOs [10]	22.05	0.11 + 9.32	1 × 4	0.784	0.557	<b>0.441</b>	0.252	8.64
SU-YOLO (Ours)	6.97	0.16 + 2.90	4	<b>0.788</b>	<b>0.582</b>	0.429	<b>0.266</b>	2.98

Table 1: **Experimental results on URPC2019 and UDD.** mAP<sub>0.5</sub> denotes mean average precision at an IoU threshold of 0.5, and mAP<sub>0.5:0.95</sub> spans IoU thresholds from 0.5 to 0.95. For SNNs, computational load includes FLOPs (from encoding) and SOPs. SpikeYOLO’s Time is calculated as time steps × virtual time steps. mAP values marked with \* are from the ANN before conversion. The ANN-SNN conversion at  $T = 3500$  is assumed lossless.

mance. Additional experiments on UDD confirm that our model consistently outperforms alternative methods.

### 4.3. Ablation Experiments

A series of ablation experiments were performed using the same parameter settings as the comparative experiments to verify the effectiveness of our proposed modules.

BL	SD	SS	SB2	SB1	Params (M)	mAP <sub>0.5</sub>		mAP <sub>0.5:0.95</sub>	
						URPC	UDD	URPC	UDD
●	○	○	○	○	10.11	0.743	0.509	0.396	0.222
●	●	○	○	○	10.11	0.752	0.523	0.398	0.232
●	○	●	○	○	9.98	0.755	0.525	0.403	0.229
●	○	●	●	○	8.96	0.760	0.545	0.401	0.248
●	○	●	●	●	6.97	0.781	0.567	0.422	0.257
●	●	●	●	●	6.97	<b>0.788</b>	<b>0.582</b>	<b>0.429</b>	<b>0.266</b>

Table 2: **Ablation experiments on URPC2019 and UDD for each module in SU-YOLO.** BL, SD, SS, SB2, and SB1 represent Baseline, SpikeDenoiser, SpikeSPP, SU-Block2, and SU-Block1, respectively.

Method	mAP <sub>0.5</sub>		mAP <sub>0.5:0.95</sub>		Energy (mJ)
	URPC	UDD	URPC	UDD	
-	0.781	0.567	0.422	0.257	2.9700
Mean Filter [56]	0.780	0.557	0.414	0.248	<b>2.9732</b>
Gaussian Filter [57]	0.783	0.573	0.415	0.261	2.9821
SpikeDenoiser	<b>0.788</b>	<b>0.582</b>	<b>0.429</b>	<b>0.266</b>	2.9751

Table 3: **Performance and energy consumption when using different denoising methods in SU-YOLO.**

#### 4.3.1. Effectiveness of Modules in SU-YOLO

We first evaluated the impact of our improved modules by replacing SU-Block1 and SU-Block2 with the basic residual block RepNCSPPELAN4 from YOLOv9, substituting SPPELAN for SpikeSPP, and replacing all activation functions with spiking neurons to maintain a fully spiking network. The overall network architecture, including the number and size of feature maps, was kept consistent across all configurations to serve as a baseline. As shown in Table 2, each module incrementally reduced the parameter count and improved accuracy. With all improvements applied, mAP<sub>0.5</sub> increased by 4.5% and mAP<sub>0.5:0.95</sub> by 3.3%. Furthermore, introducing the SpikeDenoiser module independently resulted in a 0.9% mAP<sub>0.5</sub> gain over the baseline.

#### 4.3.2. Effectiveness of SpikeDenoiser

To evaluate computational cost, we compared conventional image denoising methods with our proposed SpikeDenoiser in Table 3. For fairness, both the Mean Filter [56] and Gaussian Filter [57] were applied to the raw images before network input to preserve the network’s binary nature. Experimental results indicate that the Mean Filter degrades mAP, likely due to its strong blurring effect, while the Gaussian Filter slightly improves accuracy by better



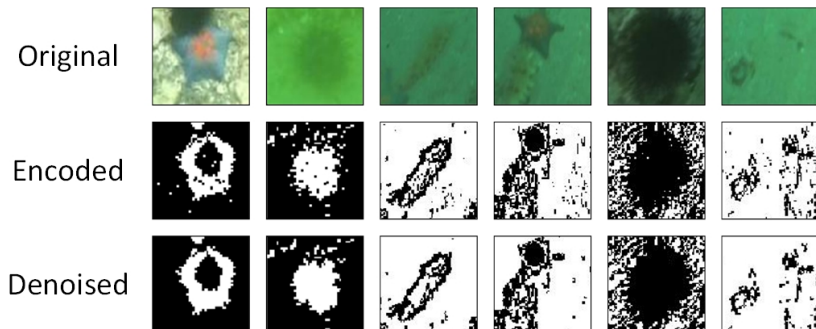


Figure 9: Denoising effect of the SpikeDenoiser.

Model	Method	URPC2019		PASCAL VOC 2012	
		mAP <sub>0.5</sub>	mAP <sub>0.5:0.95</sub>	mAP <sub>0.5</sub>	mAP <sub>0.5:0.95</sub>
EMS-YOLO (Res18)	BN [44]	0.698	0.334	0.514	0.310
	tdBN [16]	0.713	0.350	0.522	0.324
	TEBN [18]	0.705	0.335	0.519	0.316
	SeBN	<b>0.726</b>	<b>0.362</b>	<b>0.527</b>	<b>0.324</b>
SU-YOLO (Ours)	BN [44]	0.735	0.387	0.438	0.266
	tdBN [16]	0.776	0.423	0.536	0.346
	TEBN [18]	0.776	0.427	0.532	0.347
	SeBN	<b>0.788</b>	<b>0.429</b>	<b>0.537</b>	<b>0.350</b>

Table 4: Detection results for BN, tdBN, TEBN, and SeBN on URPC2019 and Pascal VOC 2012.

preserving edge information. In contrast, SpikeDenoiser yields significant mAP improvements on both datasets with minimal additional energy consumption.

Visualization results in Fig. 9 further demonstrate that SpikeDenoiser effectively removes noise from feature maps, delivering clearer inputs to the backbone network.

#### 4.3.3. Effectiveness of SeBN

We compared our novel SeBN with conventional BN methods, as summarized in Table 4. Models using standard BN perform relatively poorly, while tdBN and TEBN yield moderate improvements. Notably, in the EMS-YOLO model, where residual blocks are heavily utilized, TEBN underperforms tdBN, likely due to its inadequate handling of element-wise summation in residual structures. In contrast, in SU-YOLO, where fewer such residual blocks are present, TEBN and tdBN exhibit similar performance; replacing both with SeBN further improves detection accuracy. We also evaluated SeBN

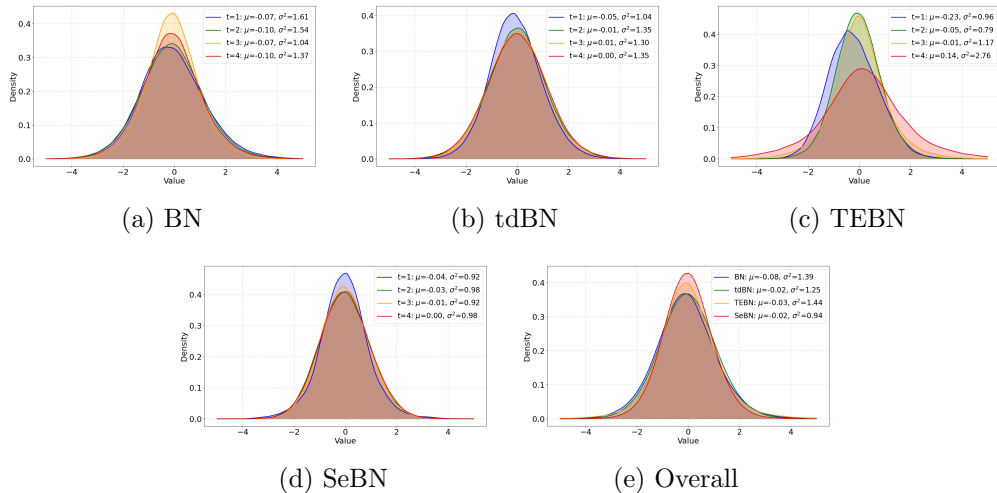


Figure 10: **Data distribution after a element summation residual block with different BN methods on URPC2019.**  $\mu$  and  $\sigma^2$  are the mean and variance of data.

on EMS-YOLO and on PASCAL VOC 2012 for general object detection, where SeBN consistently achieved the highest mAP. These results indicate that SeBN has broad applicability and generalizes well to other datasets and SNN architectures.

To illustrate the differences between SeBN and other normalization methods, we visualized the data distributions in Fig. 10. The distributions were derived from the outputs after the final residual block in layer 3 of SU-YOLO. The results reveal significant instability in the BN and TEBN distributions, with marked fluctuations in mean and variance across time steps. In particular, TEBN exhibits a variance as high as 2.76 at  $t = 1$ , likely due to unoptimized residual addition that amplifies the output range. Although tdbN shows smaller fluctuations, its variance significantly deviates from the desired standard normal distribution. In contrast, SeBN maintains a stable distribution that closely approximates  $N(0, 1)$  at each time step. An overall comparison, averaged across time steps, further confirms the superior stability and performance of SeBN.

#### 4.3.4. Residual Block Replacement

Residual blocks are the core of feature extraction and significantly impact detection accuracy. To assess the performance of our designed SU-Block, we experimented with different residual blocks while keeping the overall network

Residual Blocks	Params (M)	mAP <sub>0.5</sub>		mAP <sub>0.5:0.95</sub>		Firing Rate		
		URPC	UDD	URPC	UDD	Layer1	Layer2	Layer3
ANN-Res18	9.20	0.754	0.530	0.406	0.244	-	-	-
EMS-Res18	9.23	0.722	0.506	0.370	0.234	0.152	0.119	0.115
SU-Block	<b>6.97</b>	<b>0.788</b>	<b>0.582</b>	<b>0.429</b>	<b>0.266</b>	0.127	0.154	0.188

Table 5: Experimental results of replacing different residual blocks in SU-YOLO.

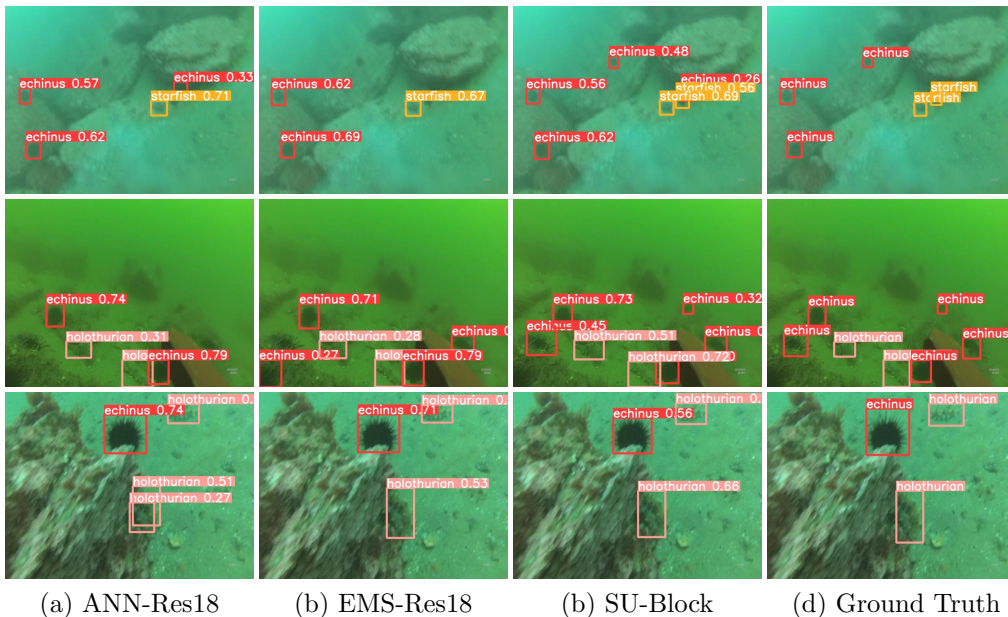


Figure 11: Detection results using various residual blocks on URPC2019. Ground Truth shows the original dataset labels.

structure constant. We calculated the average spike firing rate of various residual blocks on URPC2019. As shown in Table 5, SU-Blocks employing the CSPNet structure achieved firing rates of 12.7%, 15.4%, and 18.8%, which are higher than the firing rates of EMS-Blocks using the ResNet structure (15.2%, 11.9%, and 11.5%). This observation aligns with our visualization analysis: SU-Blocks effectively mitigate spike degradation, releasing more spikes and thereby enhancing feature extraction. Moreover, Table 5 demonstrates that SU-YOLO, using SU-Blocks, achieves an outstanding mAP<sub>0.5</sub> of 78.8%, significantly higher than that obtained using EMS-Blocks or the residual blocks employed in the ANN version.

Fig. 11 presents example detection results, highlighting that the SU-Block

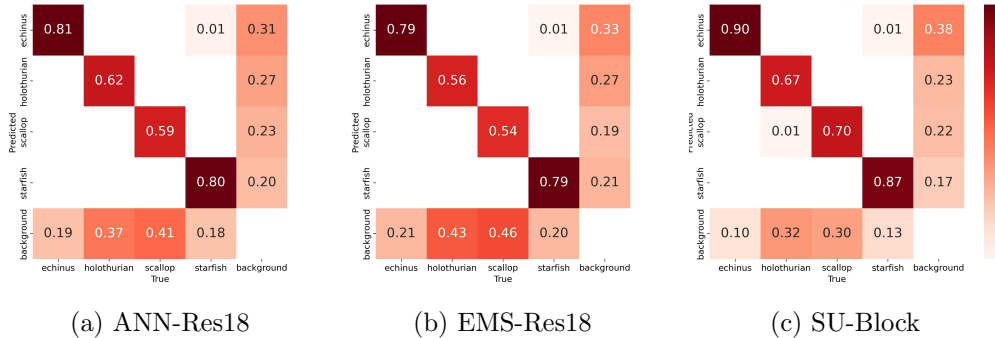


Figure 12: **Confusion matrices for models employing different residual blocks on URPC2019.** The horizontal axis represents actual objects, and the vertical axis represents predicted objects of the model.

model detects more objects. These results suggest that SU-Block can match or even surpass the performance of ANN residual blocks of the same scale in underwater object detection. Combined with the low-power consumption characteristics of SNNs, a fully SNN structure may offer advantages over hybrid ANN-SNN architectures.

#### 4.3.5. Confusion Matrix and Precision-Recall Curve

The confusion matrix in Fig. 12 illustrates the model’s classification performance. For starfish and echinus, all three models achieve high accuracy; however, holothurian and scallop detections exhibit lower accuracy, likely because their lighter colors blend with the water background, making them harder to distinguish. Notably, the SU-Block model demonstrates the highest accuracy across all object classes, particularly outperforming ANN-Res18 and EMS-Res18 by 12% and 16%, respectively, in scallop detection. This indicates that SU-Block has superior feature extraction capabilities, enabling it to detect less conspicuous objects.

The Precision-Recall curves in Fig. 13 further highlight the superior performance of the SU-Block model, especially in scallop detection, with an overall  $mAP_{0.5}$  of 78.8%, compared to 75.4% for ANN-Res18 and 72.2% for EMS-Block.

#### 4.3.6. Sequence of Pooling and Activation

The sequence of max-pooling and activation operations significantly affects model performance in SNN-based object detection. We conducted experiments

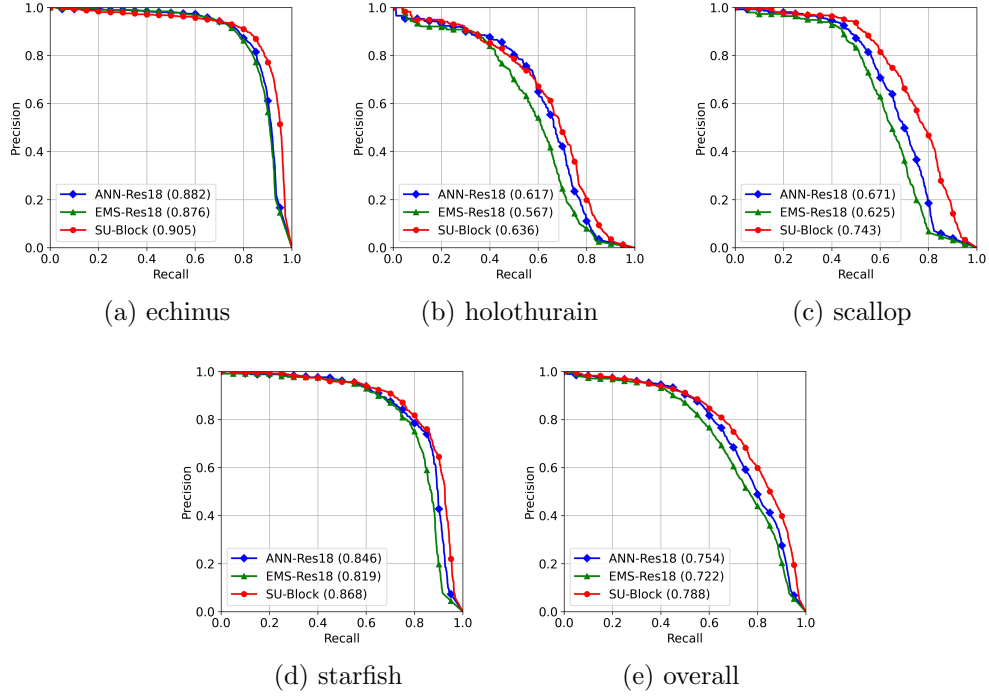


Figure 13: Precision-Recall curves for models employing different residual blocks on URPC2019. The values in the legends represent  $mAP_{0.5}$ .

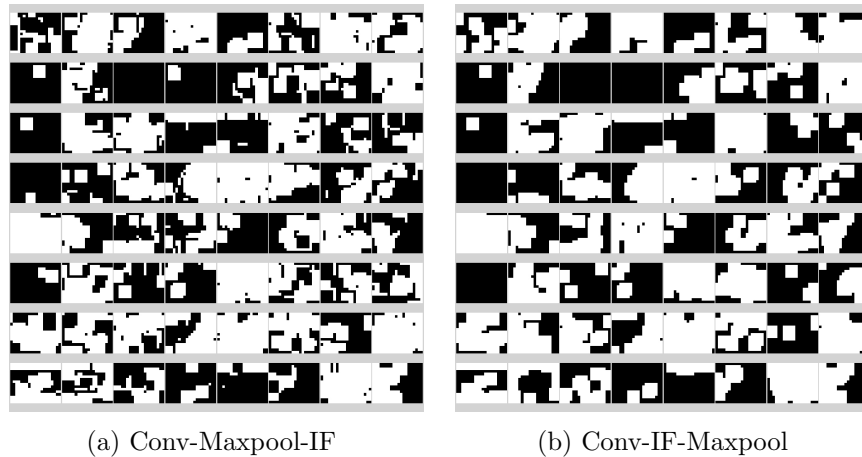


Figure 14: Pooled spike feature maps at time step  $T = 1$ , where black indicates a value of 0 and white indicates a value of 1.

Sequence	mAP <sub>0.5</sub>		mAP <sub>0.5:0.95</sub>	
	URPC	UDD	URPC	UDD
Conv-IF-MaxPool	0.775	0.548	0.416	0.253
Conv-MaxPool-IF	<b>0.788</b>	<b>0.582</b>	<b>0.429</b>	<b>0.266</b>

Table 6: **Effects of two different sequences on detection accuracy in the SpikeSPP module of SU-YOLO.**

Time Step $T$		1	2	3	4	5	6
URPC	mAP <sub>0.5</sub>	0.700	0.757	0.776	<b>0.788</b>	0.787	0.742
	mAP <sub>0.5:0.95</sub>	0.358	0.401	0.416	<b>0.429</b>	0.428	0.406
UDD	mAP <sub>0.5</sub>	0.488	0.532	0.537	<b>0.582</b>	0.578	0.567
	mAP <sub>0.5:0.95</sub>	0.219	0.238	0.250	0.266	<b>0.277</b>	0.257

Table 7: **mAP of the model as the time step  $T$  varies from 1 to 6.**

comparing two scenarios. As shown in Fig. 14, when pooling is applied before activation in the SpikeSPP, the resulting feature maps preserve more detail.

This observation is corroborated by the quantitative results in Table 6: placing the pooling layer before the IF neurons improves mAP<sub>0.5</sub> by 1.3% on URPC2019 and by 3.4% on UDD, thereby confirming our analysis.

#### 4.3.7. Length of Time Steps

The choice of time steps in SNNs directly impacts both performance and computational load. As presented in Table 7, SU-YOLO achieves optimal accuracy when the number of time steps,  $T$ , is set to 4. Increasing  $T$  beyond 4 does not further improve performance and may slightly reduce it. Therefore, we adopt 4 time steps to achieve the best balance between detection accuracy and power consumption.

## 5. Conclusion

In this study, we proposed SU-YOLO, a lightweight SNN model for underwater object detection designed for resource-constrained platforms. To address the challenge of noise in underwater images, we developed a spike-based image denoising method that proved highly effective in our experiments. Additionally, we introduced SeBN, a novel batch normalization method tailored to the unique characteristics of SNNs, which outperformed existing normalization techniques. The carefully designed modules within SU-YOLO significantly enhance detection capabilities in underwater environments, and

experimental results demonstrate that our model achieves state-of-the-art performance among current SNN approaches for underwater object detection.

However, our conclusions regarding energy efficiency remain theoretical, as they are based on estimates rather than actual hardware measurements. In future work, we aim to further enhance SU-YOLO's performance and implement it on physical hardware systems to achieve more reliable detection results. This advancement will bolster the practical applications of SNNs in engineering and contribute to the progress of underwater object detection technology. Furthermore, we plan to explore the generalization of the SU-YOLO architecture for broader object detection tasks by extending its applicability to additional datasets and applications.

### Acknowledgments

This work was supported in part by the National Natural Science Foundation of China under Grants 62271361, and the Hubei Provincial Key Research and Development Program under Grant 2024BAB039.

### References

- [1] Y. Ju, L. Li, X. Zhong, Y. Rao, Y. Liu, J. Dong, A. C. Kot, Underwater surface normal reconstruction via cross-grained photometric stereo transformer, *IEEE J. Ocean. Eng.* 50 (1) (2025) 192–203.
- [2] W. Pan, J. Chen, B. J. Lv, L. Peng, Optimization and application of improved yolov9s-ui for underwater object detection, *Appl. Sci.* (2024).
- [3] L. Wang, X. Ye, H. Xing, Z. Wang, P. Li, Yolo nano underwater: A fast and compact object detector for embedded device, *Glob. Oceans* (2020) 1–4.
- [4] L. Wang, X. Ye, S. Wang, P. Li, Ulo: An underwater light-weight object detector for edge computing, *Mach.* (2022).
- [5] W. Maass, Networks of spiking neurons: The third generation of neural network models, *Neural Networks* 10 (9) (1997) 1659–1671.
- [6] X. Zhong, S. Hu, W. Liu, W. Huang, J. Ding, Z. Yu, T. Huang, Towards low-latency event-based visual recognition with hybrid step-wise distillation spiking neural networks, in: *Proc. ACM Int. Conf. Multimedia*, 2024, pp. 9828–9836.

- [7] H. You, X. Zhong, W. Liu, Q. Wei, W. Huang, Z. Yu, T. Huang, Converting artificial neural networks to ultralow-latency spiking neural networks for action recognition, *IEEE Trans. Cogn. Dev. Syst.* 16 (4) (2024) 1533–1545.
- [8] S. J. Kim, S. Park, B. Na, S. Yoon, Spiking-yolo: Spiking neural network for energy-efficient object detection, in: *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 11270–11277.
- [9] Q. Su, Y. Chou, Y. Hu, J. Li, S. Mei, Z. Zhang, G. Li, Deep directly-trained spiking neural networks for object detection, in: *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2023, pp. 6532–6542.
- [10] X. Luo, M. Yao, Y. Chou, B. Xu, G. Li, Integer-valued training and spike-driven inference spiking neural network for high-performance and energy-efficient object detection, *arXiv:2407.20708* (2024).
- [11] J. Redmon, S. K. Divvala, R. B. Girshick, A. Farhadi, You only look once: Unified, real-time object detection, in: *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2016, pp. 779–788.
- [12] C. Wang, I. Yeh, H. M. Liao, Yolov9: Learning what you want to learn using programmable gradient information, *arXiv:2402.13616* (2024).
- [13] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [14] C. Wang, H. M. Liao, Y. Wu, P. Chen, J. Hsieh, I. Yeh, Cspnet: A new backbone that can enhance learning capability of CNN, in: *Proc. IEEE/CVF Int. Conf. Comput. Vis. Workshops*, 2020, pp. 1571–1580.
- [15] Y. Wu, L. Deng, G. Li, J. Zhu, Y. Xie, L. Shi, Direct training for spiking neural networks: Faster, larger, better, in: *Proc. AAAI Conf. Artif. Intell.*, 2019, pp. 1311–1318.
- [16] H. Zheng, Y. Wu, L. Deng, Y. Hu, G. Li, Going deeper with directly-trained larger spiking neural networks, in: *Proc. AAAI Conf. Artif. Intell.*, 2021, pp. 11062–11070.



- [17] Y. Kim, P. Panda, Revisiting batch normalization for training low-latency deep spiking neural networks from scratch, *Front. Neurosci.* 15 (2020).
- [18] C. Duan, J. Ding, S. Chen, Z. Yu, T. Huang, Temporal effective batch normalization in spiking neural networks, in: *Adv. Neural Inf. Process. Syst.*, 2022.
- [19] S. Ren, K. He, R. B. Girshick, J. Sun, Faster R-CNN: towards real-time object detection with region proposal networks, *IEEE Trans. Pattern Anal. Mach. Intell.* 39 (6) (2017) 1137–1149.
- [20] B. Rueckauer, I. Lungu, Y. Hu, M. Pfeiffer, S. Liu, Conversion of continuous-valued deep networks to efficient event-driven networks for image classification, *Front. Neurosci.* 11 (2017).
- [21] B. Chakraborty, X. She, S. Mukhopadhyay, A fully spiking hybrid neural network for energy-efficient object detection, *IEEE Trans. Image Process.* 30 (2021) 9014–9029.
- [22] Y. Li, X. He, Y. Dong, Q. Kong, Y. Zeng, Spike calibration: Fast and accurate conversion of spiking neural network for object detection and segmentation, *arXiv:2207.02702* (2022).
- [23] Y. Hu, Q. Zheng, X. Jiang, G. Pan, Fast-snn: Fast spiking neural network by converting quantized ANN, *IEEE Trans. Pattern Anal. Mach. Intell.* 45 (12) (2023) 14546–14562.
- [24] A. Tavanaei, M. Ghodrati, S. R. Kheradpisheh, T. Masquelier, A. Maida, Deep learning in spiking neural networks, *Neural Networks* 111 (2019) 47–63.
- [25] J. Lee, T. Delbrück, M. Pfeiffer, Training deep spiking neural networks using backpropagation, *Front. Neurosci.* 10 (2016).
- [26] E. O. Neftci, H. Mostafa, F. Zenke, Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks, *IEEE Signal Process. Mag.* 36 (6) (2019) 51–63.
- [27] Y. Wu, L. Deng, G. Li, J. Zhu, L. Shi, Spatio-temporal backpropagation for training high-performance spiking neural networks, *Front. Neurosci.* 12 (2017).

- [28] C. Lee, S. S. Sarwar, K. Roy, Enabling spike-based backpropagation in state-of-the-art deep neural network architectures, arXiv:1903.06379 (2019).
- [29] Y. Hu, Y. Wu, L. Deng, G. Li, Advancing residual learning towards powerful deep spiking neural networks, *IEEE Trans. Neural Networks Learn. Syst.* (2021).
- [30] W. Fang, Z. Yu, Y. Chen, T. Huang, T. Masquelier, Y. Tian, Deep residual learning in spiking neural networks, in: *Adv. Neural Inf. Process. Syst.*, 2021, pp. 21056–21069.
- [31] S. Song, J. Zhu, Research on underwater biological object recognition based on mask r-cnn and transfer learning, *Comput. Appl. Res* 37 (2020) 386–391.
- [32] K. He, G. Gkioxari, P. Dollár, R. B. Girshick, Mask R-CNN, *IEEE Trans. Pattern Anal. Mach. Intell.* 42 (2) (2020) 386–397.
- [33] M. Zhang, S. Xu, W. Song, Q. He, Q. Wei, Lightweight underwater object detection based on YOLO v4 and multi-scale attentional feature fusion, *Remote. Sens.* 13 (22) (2021) 4706.
- [34] A. Bochkovskiy, C. Wang, H. M. Liao, Yolov4: Optimal speed and accuracy of object detection, arXiv:2004.10934 (2020).
- [35] J. Zhang, J. Zhang, K. Zhou, Y. Zhang, H. Chen, X. Yan, An improved yolov5-based underwater object-detection framework, *Sensors* 23 (7) (2023) 3693.
- [36] C. Wang, A. Bochkovskiy, H. M. Liao, Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2023, pp. 7464–7475.
- [37] R. Varghese, M. Sambath, Yolov8: A novel object detection algorithm with enhanced performance and robustness, *Proc. Int. Conf. Adv. Data Eng. Intell. Comput. Syst.* (2024) 1–6.
- [38] K. Liu, L. Peng, S. Tang, Underwater object detection using TC-YOLO with attention mechanisms, *Sensors* 23 (5) (2023) 2567.

- [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: *Adv. Neural Inf. Process. Syst.*, 2017, pp. 5998–6008.
- [40] X. Li, G. Hou, K. Li, Enhancing underwater image via adaptive color and contrast enhancement, and denoising, *Eng. Appl. Artif. Intell.* 111 (2021) 104759.
- [41] N. You, L. Han, D. Zhu, W. Song, Research on image denoising in edge detection based on wavelet transform, *Appl. Sci.* (2023).
- [42] C. Tian, M. Zheng, W. Zuo, S. Zhang, Y. Zhang, C. Lin, A cross transformer for image denoising, *Inf. Fusion* 102 (2024) 102043.
- [43] Y. Ju, J. Xiao, C. Zhang, H. Xie, A. Luo, H. Zhou, J. Dong, A. C. Kot, Towards marine snow removal with fusing fourier information, *Inf. Fusion* 117 (2025) 102810.
- [44] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, in: *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 448–456.
- [45] W. Fang, Y. Chen, J. Ding, Z. Yu, T. Masquelier, D. Chen, L. Huang, H. Zhou, G. Li, Y. Tian, Spikingjelly: An open-source machine learning infrastructure platform for spike-based intelligence, *Sci. Adv.* 9 (2023).
- [46] M. Horowitz, 1.1 computing’s energy problem (and what we can do about it), in: *Proc. IEEE Int. Solid-State Circuits Conf.*, 2014, pp. 10–14.
- [47] K. He, X. Zhang, S. Ren, J. Sun, Spatial pyramid pooling in deep convolutional networks for visual recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 37 (9) (2015) 1904–1916.
- [48] underwater fish, Urpc2019 dataset, <https://universe.roboflow.com/underwater-fish-f6cri/urpc2019-nrbk1> (2023).
- [49] C. Liu, Z. Wang, S. Wang, T. Tang, Y. Tao, C. Yang, H. Li, X. Liu, X. Fan, A new dataset, poisson GAN and aquanet for underwater object grabbing, *IEEE Trans. Circuits Syst. Video Technol.* 32 (5) (2022) 2831–2844.

- [50] M. Everingham, S. M. A. Eslami, L. V. Gool, C. K. I. Williams, J. M. Winn, A. Zisserman, The pascal visual object classes challenge: A retrospective, *Int. J. Comput. Vis.* 111 (1) (2015) 98–136.
- [51] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, Q. Tian, Centernet: Keypoint triplets for object detection, in: *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 6568–6577.
- [52] M. Tan, R. Pang, Q. V. Le, Efficientdet: Scalable and efficient object detection, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 10778–10787.
- [53] G. J. et. al., ultralytics/yolov5: v6.0 - YOLOv5n 'Nano' models, Roboflow integration, TensorFlow export, OpenCV DNN support (2021). [doi: 10.5281/zenodo.5563715](https://doi.org/10.5281/zenodo.5563715).
- [54] C. Li, L. Li, Y. Geng, H. Jiang, M. Cheng, B. Zhang, Z. Ke, X. Xu, X. Chu, Yolov6 v3.0: A full-scale reloading, *ArXiv:2301.05586* (2023).
- [55] A. Wang, H. Chen, L. Liu, K. Chen, Z. Lin, J. Han, G. Ding, Yolov10: Real-time end-to-end object detection, *arXiv:2405.14458* (2024).
- [56] C. Shao, P. Kaur, R. Kumar, An improved adaptive weighted mean filtering approach for metallographic image processing\*, *J. Intell. Syst.* 30 (1) (2021) 470–478.
- [57] Y. Yu, N. Yang, C. Yang, N. Tashi, Memristor bridge-based low pass filter for image processing, *J. Syst. Eng. Electron.* (2019).

## Author Biography



**Chenyang Li** is a M.S. student in computer technology at China Three Gorges University. He received his B.S. degree in water conservancy and hydropower engineering from China Three Gorges University in 2021. His research interests include computer vision and neuromorphic computing.



**Wenxuan Liu** received her Ph.D. from Wuhan University of Technology in 2024. Now she is a postdoc at Peking University. Her research interests include neuromorphic computing, spike vision, action recognition, and multimedia content analysis.



**Guoqiang Gong** received his B.S. from Lanzhou University in 1998 and Ph.D. from Tongji University in 2010. He is currently a Professor at China Three Gorges University. His research interests include intelligent signal processing and computer vision.



**Xiaobo Ding** received his B.S. from Beijing Electronic Science and Technology Institute in 1995 and Ph.D. from Harbin Engineering University in 2015. He is currently a Associate Professor at China Three Gorges University. His research interests include IoT technology and system virtualization.



**Xian Zhong** received his B.S. degree in Computer Science from Wuhan University, China, in 2007, and his Ph.D. degree in Computer Science from Huazhong University of Science and Technology, China, in 2013. He is currently a Professor in the School of Computer Science and Artificial Intelligence at Wuhan University of Technology, China. His research interests include multimodal information processing, intelligent transportation systems, and neuromorphic computing.

## Figure Captions

Fig. 1: **Comparison of detection performance and energy consumption across various models on URPC2019.** The blue dots represent ANN models, while the green and red dots indicate SNN models.

Fig. 2: **Diagram of SU-YOLO and the structure of its modules.** In the main structural diagram, narrower graphics represent feature maps with reduced dimensions.

Fig. 3: **Schematic diagram illustrating the processing flow and effects of the SpikeDenoiser module.**

Fig. 4: **Feature maps extracted from backbone layers when using ResNet or CSPNet in SNN.**

Fig. 5: **Impact of pooling and activation sequence on output.**

Fig. 6: **Calculation and fusion methods of SeBN.** In the SeBN computation diagram, cubes represent the feature maps at each time step, where  $B$  denotes the batch size,  $C$  indicates the number of channels,  $H$  and  $W$  are the height and width of the feature maps, and  $T$  represents the total number of time steps. The symbols  $\mu_k$  and  $\sigma_k^2$  denote the mean and variance of the  $k$ -th feature map.

Fig. 7: **Process of batch-scale fusion of SeBN.** Here,  $\theta_t$  represents all parameters of the SeBN layer at time step  $t$ ,  $\phi$  represents parameters shared by the convolutional layer across all time steps, and  $\phi_t$  denotes parameters of the convolutional layer at time step  $t$ .

Fig. 8: **Running process of SU-YOLO varies with time steps.**

Fig. 9: **Denoising effect of the SpikeDenoiser.**

Fig. 10: **Data distribution after a element summation residual block with different BN methods on URPC2019.**  $\mu$  and  $\sigma^2$  are the mean and variance of data.

Fig. 11: **Detection results using various residual blocks on URPC2019.** Ground Truth shows the original dataset labels.

Fig. 12: **Confusion matrices for models employing different residual blocks on URPC2019.** The horizontal axis represents actual objects, and the vertical axis represents predicted objects of the model.

Fig. 13: **Precision-Recall curves for models employing different residual blocks on URPC2019.** The values in the legends represent  $mAP_{0.5}$ .

Fig. 14: **Pooled spike feature maps at time step  $T = 1$ , where black indicates a value of 0 and white indicates a value of 1.**

# Graphical Abstract

