# E-ConvNeXt: A Lightweight and Efficient ConvNeXt Variant with Cross-Stage Partial Connections

Fang Wang[a,1], Huitao Li[a], Wenhan Chao[b], Zheng Zhuo[a,*], Yiran Ji[a], Chang Peng[a], Yupeng Sun[a] and jiyu zhang[a]

[a]*College of Information Engineering, Beijing Institute of Petrochemical Technology, Beijing, 102617, People's Republic of China*
[b]*School of Computer Science and Engineering, Beihang University, Beijing, 100191, People's Republic of China*

## ARTICLE INFO

## ABSTRACT

Many high-performance networks were not designed with lightweight application scenarios in mind from the outset, which has greatly restricted their scope of application. This paper takes ConvNeXt as the research object and significantly reduces the parameter scale and network complexity of ConvNeXt by integrating the Cross Stage Partial Connections mechanism and a series of optimized designs. The new network is named E-ConvNeXt, which can maintain high accuracy performance under different complexity configurations. The three core innovations of E-ConvNeXt are : (1) integrating the Cross Stage Partial Network (CSPNet) with ConvNeXt and adjusting the network structure, which reduces the model's network complexity by up to 80%; (2) Optimizing the Stem and Block structures to enhance the model's feature expression capability and operational efficiency; (3) Replacing Layer Scale with channel attention. Experimental validation on ImageNet classification demonstrates E-ConvNeXt's superior accuracy-efficiency balance: E-ConvNeXt-mini reaches 78.3% Top-1 accuracy at 0.9GFLOPs. E-ConvNeXt-small reaches 81.9% Top-1 accuracy at 3.1GFLOPs. Transfer learning tests on object detection tasks further confirm its generalization capability.

## 1. Introduction

Neural networks[1] have developed rapidly in the field of computer vision. The 2010s can be called the era of convolutional neural networks: AlexNet [2] pioneered the history of CNN-dominated computer vision; ResNet[3], with its wide adaptability to various downstream tasks, opened the chapter of image classification networks as backbone networks; and EfficientNet[4] promoted neural networks toward a multi-variant development trend through compound scaling strategies. Entering the 2020s, the introduction of Transformer[5] has made the field of computer vision more open and diverse. Vision Transformer (ViT)[6] and Convolutional Neural Networks (ConvNets) learn from each other, giving birth to a series of new networks, whose development trends focus on higher accuracy, faster speed, and lower computational complexity.

ConvNeXt[7] emerges as a representative of such hybrid networks that bridge the gap between traditional ConvNet and Transformers. As a key innovation in the 2020s, it inherits the strengths of ResNet's convolutional foundation while integrating core design ideas from Vision Transformer, aiming to combine the efficiency of pure convolutional architectures with the expressive power of Transformer-based models. ConvNeXt, based on ResNet, draws on the design concepts of Vision Transformer. It enlarges the convolution kernel to increase the effective receptive field, introduces depthwise convolution[8] to control network parameters and network complexity, and simplifies the network architecture.

It retains the efficient feature extraction capability of pure convolutional networks while absorbing the architectural advantages of Transformer. ConvNeXt outperforms Swin-Transformer[9] in image classification task[10] and also demonstrates strong feature extraction capability in downstream tasks[11][12].

However, despite the excellent performance of ConvNeXt, it has shortcomings in lightweight adaptation[8][13][14]. Unlike EfficientNet-Lite[4] and FasterNet-tiny[15], which are specifically designed for lightweight requirements, ConvNeXt lacks small-parameter variants. As a result, it is challenging to find suitable ConvNeXt variant that meet low-computation requirements in resource-constrained scenarios, such as mobile and embedded devices.

ConvNeXt-tiny is the smallest variant of ConvNeXt, but the parameters of ConvNeXt-tiny are 28M, the network complexity of ConvNeXt-tiny is 4.5GFLOP. Compared with other lightweight networks, ConvNeXt still has a huge gap in terms of parameter count and model complexity, which greatly limits its application scenarios.

To address this critical gap, E-ConvNeXt is developed as a lightweight derivative of ConvNeXt. Explicitly engineered to mitigate ConvNeXt's inadequacies in lightweight support, E-ConvNeXt achieves significant reductions in network complexity and parameter through structural refinements and parameter recalibration, while preserving competitive performance metrics. This design enables seamless adaptation to resource-constrained scenarios, thereby extending the applicability of the ConvNeXt framework across a broader spectrum of practical use cases. The main contributions of E-ConvNeXt are as follows:

1. This paper introduce CSPNet[16] into ConvNeXt and incorporate a transition hyperparameter into the

---

*Corresponding author

✉ fangwang@bipt.edu.cn (F. Wang); 2024520241@bipt.edu.cn (H. Li); chaowenhan@buaa.edu.cn (W. Chao); 0020230022@bipt.edu.cn (Z. Zhuo)

ORCID(s): 0000-0001-6170-0463 (F. Wang); 0009-0005-4750-6273 (H. Li); 0000-0001-8028-5285 (W. Chao); 0000-0001-8028-5285 (Z. Zhuo)

CSP Stage, thereby significantly reducing the network complexity of ConvNeXt.

2. We improve the model's operational efficiency and accuracy through a series of optimizations, specifically including: adjusting the Stem structure; replacing Layer Normalization with Batch Normalization; and using the channel attention mechanism to replace Layer Scale.

3. Based on the aforementioned improvements, this paper proposes a new network named E-ConvNeXt. We conducted extensive experiments on ImageNet-1K to verify the effectiveness of E-ConvNeXt. Additionally, we transferred E-ConvNeXt to object detection tasks to validate its generalization capability.

## 2. Related Work

### 2.1. Lightweight Neural Network

With the surging demand for mobile and edge computing scenarios, lightweight networks[8][4] have gradually become an important research direction in the field of computer vision. Traditional ConvNets fail to meet the lightweight requirements in terms of parameters and network complexity.For example, the network complexity of ResNet-18 is 1.8GFLOPs, it's parameters are 11M. MobileNet[8] replaces traditional convolution with Depthwise Separable Convolution, which significantly reduces parameters and network complexity of ConvNets. MobileNetV1_x0_25 network complexity is 0.07GFLOPs , it parameters are 0.460M. Early lightweight neural networks[13][14][17][18] only focus on lightweight application scenarios and achieve good results in such scenarios, but they can't meet the requirements in high-performance scenarios.

EfficientNet[4] can cater to multiple application scenarios by controlling the network's depth, width, and input image size. EfficientNet has multiple variants. The network complexity of EfficientNet-B0 is 0.7GFLOPs, that of EfficientNetB7 is 72.3GFLOPs. EfficientNet has suitable variants to cope with lightweight application scenarios and high-performance application scenarios. Many ConvNets[19][15] draw inspiration from EfficientNet and set up multiple variants to cope with different application scenarios.

However, many network do not consider lightweight application scenarios, and how to apply them to lightweight scenarios remains a relatively under-explored research direction.This paper takes ConvNeXt as the research object and combines CSPNet with ConvNeXt, significantly reducing the number of model parameters and network complexity while ensuring a certain level of accuracy.

### 2.2. Cross Stage Partial Network (CSPNet)

Cross Stage Partial Networks (CSPNet) respects the variability of the gradients by integrating feature maps from the beginning and the end of a network stage. CSPNet is combined with many networks[20][21][3], it can reduce computations by 20% with equivalent or even superior accuracy on the ImageNet dataset.

When ResNeXt-50 is transformed into CSPResNeXt-50, its network complexity decreases from 5.05 GFLOPs to 3.97 GFLOPs. When DenseNet-201[20] is transformed into CSPDenseNet-201, its network complexity drops from 4.35 GFLOPs to 3.05 GFLOPs. However, even after a 20% reduction in network complexity, these networks still have a huge gap compared with lightweight networks such as MobileNet and EfficientNet, and thus cannot meet the lightweight requirements.This paper significantly reduces the network complexity and parameters of model by adjusting the structure of CSPNet. Meanwhile, it can maintain good performance.

### 2.3. Channel Attention

Channel attention mechanisms aim to adaptively recalibrate feature channel weights by emphasizing discriminative channels and suppressing irrelevant ones, thereby enhancing model representational capacity. Pioneered by SENet[22], this paradigm introduced the squeeze-and-excitation (SE) block, which aggregates global spatial information via global average pooling (GAP) in the squeeze module and models cross-channel dependencies using multi-layer perceptrons (MLPs) in the excitation module.
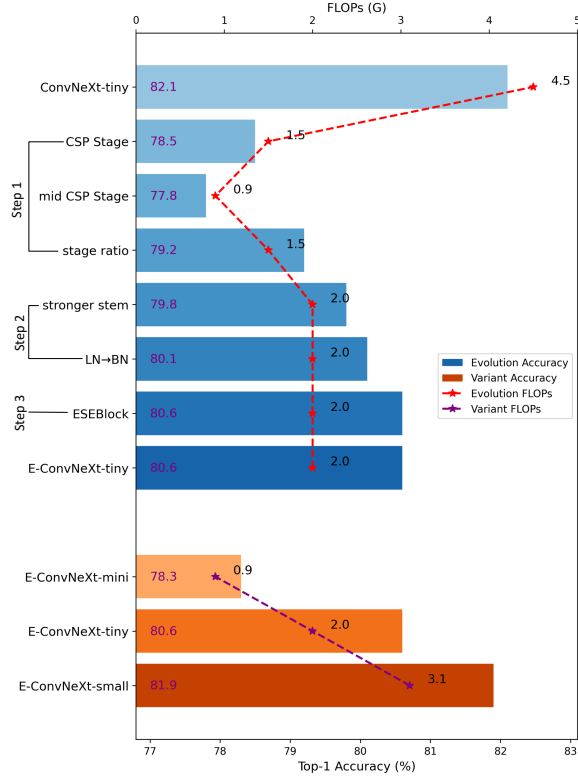
Subsequent works have refined this framework: ECANet [23] replaced MLPs with 1D convolutions to capture local cross-channel interactions without dimensionality reduction, reducing computational complexity while preserving performance. GSoP-Net [24] improved the squeeze module by incorporating global second-order pooling to model high-order statistical relationships between channels, enhancing global information capture. SRM[25] combined style pooling (mean and standard deviation) with lightweight channel-wise fully-connected layers, balancing efficiency and expressiveness. FcaNet[26] reinterpreted GAP through the lens of discrete cosine transforms (DCT), leveraging multi-spectral information to enrich global feature representation.

ConvNeXt uses Layer Scale to help stabilize network training. This paper will explore how to use attention mechanisms to replace Layer Scale to further enhance the model's feature extraction capability.

## 3. The proposed E-ConvNeXt

In this section, we provide a trajectory going from the original ConvNeXt to our E-ConvNeXt. Figure 1 shows the ideas and effects of each improvement step. To ensure fairness, we use ConvNeXt-tiny as the baseline and maintain the training strategy for all steps. All steps are implemented based on this trajectory and baseline setting, with specific improvements as follows:

1. Combine CSPNet with ConvNeXt: we combine CSP-Net with ConvNeXt and improve the CSP Module through the following three steps: 1) CSP Stage; 2) Introducing transition hyperparameter; 3) Stage ratio. As shown in Figure 1, this step significantly reduces network complexity. The specific details will be described in section 3.1;

**Figure 1: E-ConvNeXt Evolution Figure** shows the steps from ConvNeXt to E-ConvNeXt. The bar chart represents the Top-1 accuracy on ImageNet-1K, and the line chart shows the changes in model FLOPs.

2. Optimization of the ConvNeXt structure: 1) A stepped stem is proposed to alleviate the problem of information loss caused by continuous downsampling. 2) replacing Layer Normalization[27] (LN) with Batch Normalization[28] (BN): BN is used to replace LN tor further improve efficiency.

3. Introducing channel attention: We solved the problem that adding channel attention to ConvNeXt would cause training collapse, and added the Effective Squeeze-and-Excitation Block[29] (ESE Block) to the network.

As shown in Figure 1, step 1 can significantly reduce the network complexity, and step 2 and 3 enhance the model's efficiency and accuracy. Finally, we get E-ConvNeXt-tiny from ConvNeXt-tiny. Building on E-ConvNeXt-tiny, we further extend it into two variants: E-ConvNeXt-mini and E-ConvNeXt-small. Among them, E-ConvNeXt-mini network complexity is 0.9 GFLOPs and Top-1 accuracy is 78.3%; E-ConvNeXt-tiny network complexity is 2.0 GFLOPs and Top-1 accuracy is 80.6%; E-ConvNeXt-small network complexity is 3.1 FLOPs and Top-1 accuracy is 81.9%.

## 3.1. Combine CSPNet with ConvNeXt
### 3.1.1. CSP Stage

CSPNet reduces parameters and FLOPs while maintaining accuracy by adjusting the channel numbers in some

convolutions. As shown in Figure 2, CSPResNet splits the feature map into two feature maps. The channel numbers of each new map are half of the original ones. One feature map passes through the blocks and the other does not. Then the two parts are merging together. Compared to ResNet, CSPResNet only changes the input channel and output channel of the 1x1 convolution from 256 to 128.The formula for calculating the convolution FLOPs is

$$FLOPs = C_{in} \times H_{out} \times W_{out} \times K^2 \times C_{out}, \quad (1)$$

The FLOPs of ResNet Block are

$$((256 \times 56 \times 56 \times 1^2 \times 64) \approx 51,4M)+$$
$$((64 \times 56 \times 56 \times 3^2 \times 64) \approx 116M)+$$
$$((64 \times 56 \times 56 \times 1^2 \times 256) \approx 51.4M) \approx 218.8M.$$

The FLOPs of CSPResNet Block are

$$((128 \times 56 \times 56 \times 1^2 \times 64) \approx 25.7M)+$$
$$((64 \times 56 \times 56 \times 3^2 \times 64) \approx 116M)+$$
$$((64 \times 56 \times 56 \times 1^2 \times 128) \approx 25.7M) \approx 167.4M.$$

Compared with ResNet, CSPResNet reduced 10% FLOPs and maintained the same accuracy.

Inspired by CSPResNet, we proposed E-ConvNeXt by combining CSPNet with ConvNeXt. To further reduce the parameters and FLOPs, we adjust the channel numbers for all convolutions in the block. The number of channels in the 7×7 depthwise convolution and the MLP layer is halved. As shown in Figure 2, the number of channels in E-ConvNeXt is reduced from 96 to 48 and from 384 to 192 compared to ConvNeXt. The formula for calculating depthwise convolution is

$$FLOPs = H_{out} \times W_{out} \times K^2 \times C_{out}, \quad (2)$$

The FLOPs of ConNeXt Block are

$$((56 \times 56 \times 7^2 \times 96) \approx 15M)+$$
$$((96 \times 56 \times 56 \times 1^2 \times 384) \approx 116M)+$$
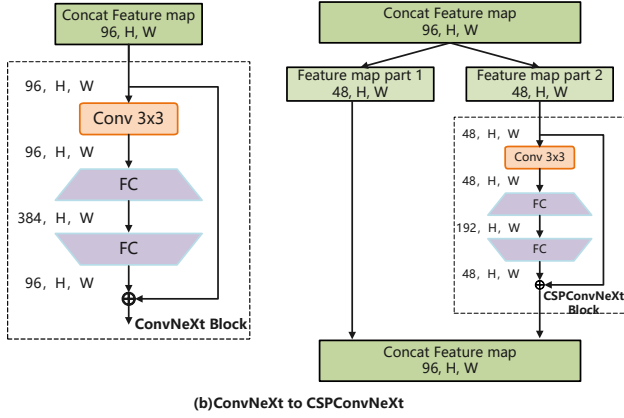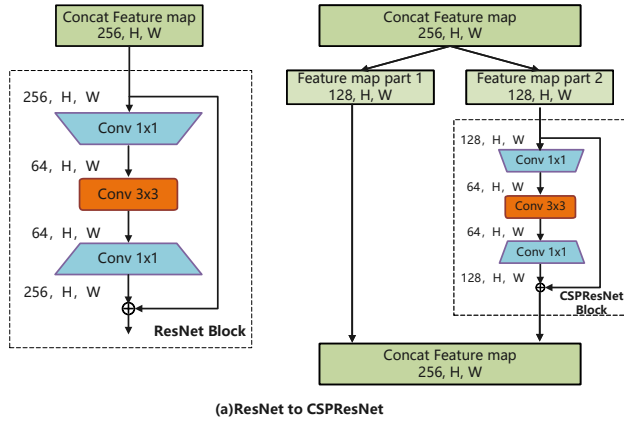$$((384 \times 56 \times 56 \times 1^2 \times 96) \approx 116M) \approx 257M.$$

The FLOPs of CSPConNeXt Block are

$$((56 \times 56 \times 7^2 \times 58) \approx 7.5M)+$$
$$((48 \times 56 \times 56 \times 1^2 \times 192) \approx 29.5M)+$$
$$((384 \times 56 \times 56 \times 1^2 \times 96) \approx 29.5M) \approx 66.5M.$$

Obviously, by introducing the CSP Module, the FLOPs of the original ConvNeXt is reduced by 60%.

### 3.1.2. Introducing transition hyperparameter

The structure of CSP Stage is shown in Figure 3 A. The CSP Stage uses splits the channels of the feature map into two parts via split operation. We use 1×1 convolution to

**(a)ResNet to CSPResNet**



**(b)ConvNeXt to CSPConvNeXt**

**Figure 2: CSPConvNets demonstration figure** shows how CSPNet transforms ResNet into CSPResNet, and how we transform ConvNeXt into CSPConvNeXt.



**Figure 3:** (A) is Original CSPConvNeXt CSP Stage. (B) is $ch_{\mathrm{mid}}$ CSPConvNeXt CSP Stage

complexity decreases, the model accuracy exhibits a corresponding decline. To achieve a balance between FLOPs and accuracy, we increase the FLOPs of the $ch_mid$ CSP-ConvNeXt to the same level as the Original CSPConvNeXt by adjusting the stage ratio. The adjusted stage ratio is as follows: input channels are [65, 128, 256, 512], output channels are [128, 256, 512, 1024], blocks are [3, 3, 9, 3].

The downsampling layers in ConvNeXt consist of four layers, namely one 4× downsampling layer followed by three 2× downsampling layers. To better adapt to the structure of CSP Module, the downsampling architecture of CSPConvNeXt is modified to five 2× downsampling layers, and each layer is composed of a 2x2 convolution with a stride of 2. The final structure of CSPConvNeXt is shown in Figure 4. Compared with the original CSPConvNeXt, it has the same FLOPs but the accuracy increases 0.7%.

### 3.2. Optimization of Network Structure Design
#### 3.2.1. Stepped Stem

As shown in Figure 5 (a), the stem layer of ConvNeXt employs a 4×4 convolution with a stride of 4. The CSPConvNeXt stem structure is modified by replacing the single 4×4 convolution with two successive 2×2 convolutions with a stride of 2, the specific structure is shown in Figure 5 (b). The first two downsampling operations in CSPConvNeXt will not only adjust the size of the feature map but also frequently modify the number of channels. This process will reduce the network accuracy.

Drawing inspiration from the designs of $ResNet_{\mathrm{vc}}$ (as shown in Figure 5 c) and ConvNeXt, the stem layer of E-ConvNeXt is modified by replacing the original 2×2 convolution with a stepped convolution combination. This combination can resolve the above issue. It consists of one 2×2 convolution and two 3×3 convolutions, with the number of channels changing from 3 to 32 and then to 64. Specifically, the 2×2 convolution retains the Patchify concept. The two 3×3 convolutions enable better feature extraction. The stepped channel adjustment further mitigates the loss

replace the Split operation, which can better split the feature map.

The feature map channel transformation process of the original CSPNet is

$$ch_{\mathrm{in}} \xrightarrow{\substack{\text{down-}\\\text{sampling}}} ch_{\mathrm{out}} \xrightarrow{\text{Split}} 2 \times (\frac{ch_{\mathrm{out}}}{2}) \xrightarrow{\text{concat}} ch_{\mathrm{out}}$$
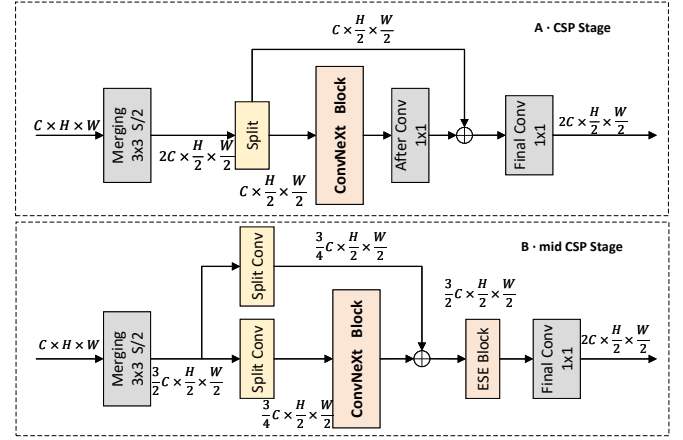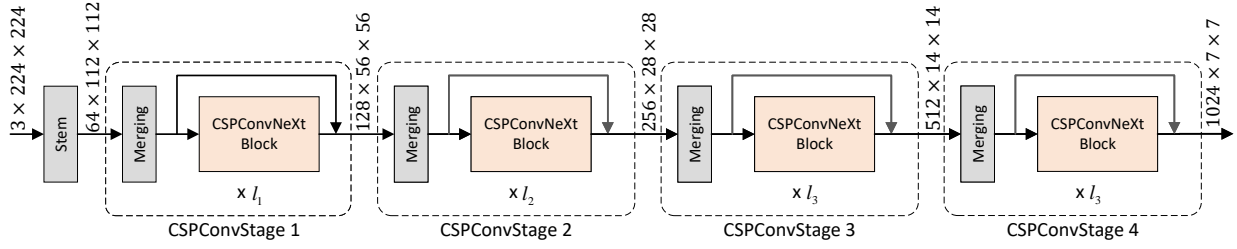
To further reduce the number of model parameters, we introduce an transition hyperparameter $ch_{\mathrm{mid}}$ . The formula for $ch_{\mathrm{mid}}$ is $ch_{\mathrm{mid}} = \frac{ch_{\mathrm{in}}+ch_{\mathrm{out}}}{2}$. The new feature map channel transformation process is

$$ch_{\mathrm{in}} \xrightarrow{\substack{\text{down-}\\\text{sampling}}} ch_{\mathrm{mid}} \xrightarrow{\substack{\text{1x1}\\\text{Conv}}} 2 \times \frac{ch_{\mathrm{mid}}}{2} \xrightarrow{\text{concat}} ch_{\mathrm{mid}} \xrightarrow{\substack{\text{1x1}\\\text{Conv}}} ch_{\mathrm{out}}$$

We call this version as $ch_{\mathrm{mid}}$ CSPConvNeXt. As shown in Figure 3 B, by introducing the $ch_{\mathrm{mid}}$ hyperparameter, it further reduces the FLOPs by 40% compared to the original CSPConvNeXt.

#### 3.1.3. Stage ratio

In both the original CSPConvNeXt and the $ch_{\mathrm{mid}}$ CSP-ConvNeXt, we have achieved a significant reduction in parameters and network complexity. However, as network

**Figure 4:** After adjusting the Stage ratio, CSPConvNeXt consists of one Stem and four CSPConvStages. Each Stage includes a downsampling layer and multiple CSPConvNeXt Blocks. The output of each Stage is twice the size of its input. The diagram shows the channel configuration of CSPConvNeXt-tiny.

of spatial features. The final stem structure is illustrated in Figure 5 d, which is called the Stepped Stem.

### 3.2.2. Layer Normalization → Batch Normalization

ConvNeXt employs LN instead of BN. There are two modes of LN in ConvNeXt. Within the blocks, LN adopts the channel last mode, which is faster but requires transposing the feature maps. On the other hand, LN uses the channel first mode, which does not need feature map transposition but its speed is relatively slow. However, CSP Module requires a large number of convolution operations to adjust the number of channels, and the use of LN will seriously affect the running speed of the model.

**Table 1**

Comparison of efficiency between LN and BN: *Downsample*1 employs a 2×2 convolution with a stride of 2, followed by LN. *Downsample*2 utilizes a 2×2 convolution with a stride of 2, followed by BN. The input feature map has a size of [8, 64, 56, 56]. Repeating 10,000 forward and backward propagation processes

| Name | DownSample-1 | DownSample-2 |
|---|---|---|
| Deploy-time | 16.41s | 29.65s |
| Percentage | 55.3% | 100% |

We tested LN and BN in the channel first mode. As shown in Table 1, we designed two downsampling configurations for comparison. The results show that the downsampling operation using BN runs significantly faster than that using LN. Therefore, all LN in non-block ConvLayers are replaced with BN, and all non-block ConvLayers are converted into "Convolution + BN + GELU".

As shown in Figure 6 (a), the LN in the ConvNeXt Block uses the channel last mode, which avoids inefficiency. However, this requires two transpose operations in the block layer to transpose the feature map, and also needs to use Fully Connected layers instead of 1x1 convolutions. Thus, we replaced LN with BN in the block layer, removed transpose operations, and used 1x1 convolutions instead of Fully Connected Layers. The modified block structure is shown in Figure 6 (b). The modified block structure does not

change the network's FLOPs or parameters, but increases the network speed by 20% and brings a slight improvement in accuracy.

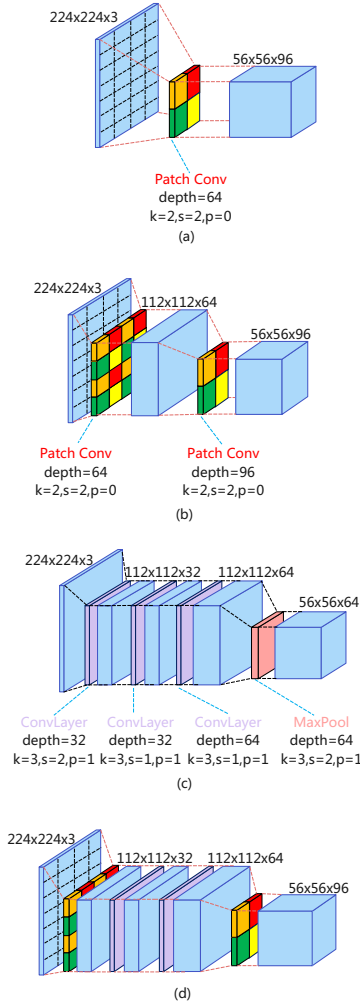### 3.3. Introducing Channel Attention

The original ConvNeXt architecture uses LayerScale[30] for feature scaling. In this paper, we replace this module with channel attention module, aiming to improve the accuracy of network. However, if channel attention is directly added to the end of the block, it will lead to training instability or model collapse.

To address the above issue, a normalization layer is introduced before the channel attention module. Referring the same method of introducing channel attention module in SENet and EfficientNet, the SE Block is added into ConvNeXt. The specific approach is shown in Figure 7. From the perspective of parameters and network complexity: the channel numbers in the SENet-style configuration is $c$, while that in the EfficientNet-style configuration reaches as high as $4c$. Although both result in a similar increase in FLOPs, the channel expansion in EfficientNet introduces a substantial number of additional parameters, significantly boosting the model's parameters. Consequently, the integration style of SENet is more aligned with the design goal of E-ConvNeXt in maintaining lightweight characteristics. We incorporate the ESE Block into the network as a channel attention module. After comparison, we integrate the ESE module into E-ConvNeXt.

### 3.4. The overall architecture of E-ConvNeXt

By combining CSPNet with ConvNeXt and improving the network architecture, E-ConvNeXt is proposed for balancing accuracy, FLOPs and parameters. Figure 8 shows the overall architecture. It has five hierarchical stages, each stage is preceded by a $2 \times 2$ convolution with stride 2 for spatial downsampling and channel expansion.

The first stage is the stem layer, which consists of a $2 \times 2$ convolution with stride 2, followed by two regular $3 \times 3$ convolution with stride 1. The number of channels in the Stem layer follows a stepped pattern. The following four stages are CSPConvStage. Each CSPConvStage contains a $2 \times 2$ convolution with stride 2 for downsampling, followed
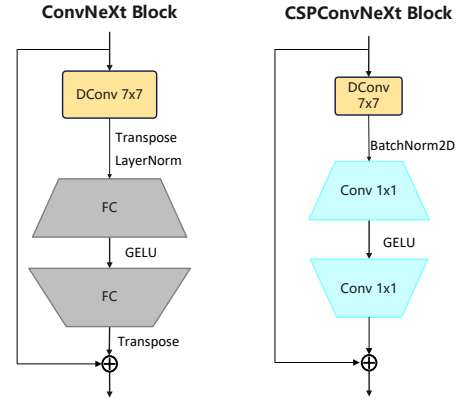
**Figure 6:** Comparison diagram of Block structures: The ConvNeXt Block uses fully connected layers to form the MLP layer, with Layer Normalization (LN) for normalization. The CSPConvNeXt Block employs 1x1 convolutions to construct the MLP layer, and adopts Batch Normalization (BN) for normalization.



**Figure 5:** a is the Stem of ConvNeXt; b is the Stem of CSP-ConvNeXt after Stage ratio adjustment and the downsampling layer in the first CSPStage; c is the Stem layer of $ResNet_{vc}$; d is the Step Stem. The Step Stem is formed by integrating the characteristics of the Stem layers of ConvNeXt and $ResNet_{vc}$.
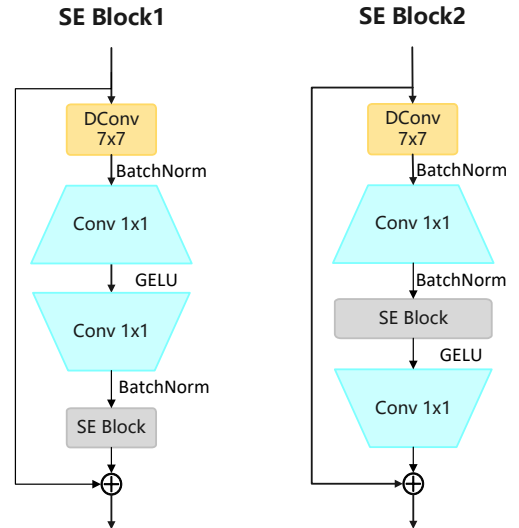
by splitting the feature map into two parts.ne feature map passes through the blocks and the other does not, and the feature maps are finally merging together. CSPConvStage introduces $ch_{mid}$ to further reduce the network's parameters and FLOPs. 1x1 conv is used to replace the Split operation.

Each CSPConvStage contains multiple Conv blocks. Each block consists of a 7x7 Dpethwise Conv followed by two 1x1 Conv layers. Together, they form an inverted residual block. We replace LN with BN, significantly improving the network effective. As for the activation layer, we keep use GELU. The final two layers are global average pooling and a fully connected layer, used for feature transformation and classification.

To meet a wide range of application requirements under different computational budgets, we present three variants of E-ConvNeXt. We call them CSPConvNeXt-mini, E-ConvNeXt-tiny, and E-ConvNeXt-Small respectively.They share a similar architecture but differ in depth and width.

**Figure 7:** SE Block1 adopts the integration method of SENet, with the number of channels in the SE Block being c. SE Block2 follows the integration method of EfficientNet, where the number of channels in the SE Block is 4c.

## 4. Experiments

This paper used two benchmark datasets: ImageNet1K[10] and ImageNet-100[1].

- ImageNet-1K: It is widely used in image classification tasks. It consists of 1,000 classes and contains approximately 1.28 million training images and 50,000 validation images. Each class represents a distinct category, ensuring the diversity and robustness of model evaluation.

[1]https://huggingface.co/datasets/ilee0022/ImageNet100

**Figure 8:** The overall architecture of E-ConvNeXt, illustrating the hierarchical stages, including the Stepped Stem, CSPConvStages, and the integration of key components such as ESE Block, Batch Normalization, and depthwise convolutions.

**Table 2**

**ImageNet-1k/100 (pre-)training settings.** The experimental setup of E-ConvNeXt on ImageNet-1k and ImageNet-100, for more detailed settings, please refer to the github repository.

| (pre-)training config | E-ConvNeXt -M/T/S ImageNet-1K $224^2$ | E-ConvNeXt -M/T/S ImageNet-100 $224^2$ |
|---|---|---|
| Weight init | trunc. normal (0.2) | trunc. normal (0.2) |
| Optimizer | AdamW[31] | AdamW[31] |
| Base learning rate | 1.25e-4 each 128 batch size | 1.25e-4 each 128 batch size |
| Weight decay | 0.05 | 0.05 |
| Optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ | $\beta_1, \beta_2 = 0.9, 0.999$ |
| Training epochs | 300 | 100 |
| Learning rate schedule | cosine decay | cosine decay |
| Warmup epochs | 20 | 5 |
| Warmup schedule | linear | linear |
| Layer-wise lr decay | None | None |
| Randaugment[32] | (9, 0.5) | (9, 0.5) |
| Mixup[33] | 0.8 | 0.8 |
| Cutmix[34] | 1.0 | 1.0 |
| Random erasing[35] | 0.25 | 0.25 |
| Label smoothing[36] | 0.1 | 0.1 |
| Head init scale | None | None |
| Gradient clip | None | None |

- ImageNet-100: It is a subset of ImageNet-1K, containing 100 classes. In in paper, it is primarily used to quickly verify the impact of different structures on model performance.

### 4.1. Settings

We provide the pretraining settings of E-ConvNeXts on ImageNet-1K and ImageNet-100 in Table 2. All E-ConvNeXt variants use the same settings, except that the base learning rate increases proportionally with the batch size. For example, when training E-ConvNeXt-tiny on 4 NVIDIA 4080 GPUs, the batch size is set to 256 and the base learning rate is set to 2.5e-4. In contrast, when training on 8 NVIDIA 4090 GPUs, the batch size is set to 1024 and the base learning rate is set to 1e-3.

### 4.2. Comparison of different CSP architectures

To compare the performance of various CSP structures, we train and validate each version of CSPConvNeXt on ImageNet-1K.

- original CSPConvNeXt: it combines CSPNet with ConvNeXt while preserving the channel ratio between convolution and MLP in the Block of ConvNeXt.

- $ch_{mid}$ CSPConvNeXt: it uses 1×1 convolution to replace the Split operation and introduces the transition hyperparameter $ch_{mid}$.

- final CSPConvNeXt: it is obtained by adjusting the downsampling layer and stage ratio of $ch_{mid}$ CSPConvNeXt.

As shown in the table 3, The original CSPConvNeXt can reduce FLOPs from 4.5G to 1.5G. The $ch_{mid}$ CSPConvNeXt can further reduce FLOPs to 0.9G. However, as FLOPs decrease, the accuracy of the models also shows a downward trend.

To further balance FLOPs and Accuracy, we adjusted the stage ratio of the $ch_{mid}$ CSPConvNeXt, resulting in the final CSPConvNeXt. The FLOPs of the final CSPConvNeXt are increased to the same level as the original CSPConvNeXt, but its accuracy is 0.7% higher. Therefore, we selected the final CSPConvNeXt as the base model for subsequent experiments.

**Table 3**
**Compare different CSP structures**

|  | Top-1 Acc | FLOPs |
|---|---|---|
| original CSPConvNeXt | 78.5 | 1.5G |
| $ch_{mid}$ CSPConvNeXt | 77.8 | 0.9G |
| final CSPConvNeXt | 79.2 | 1.5G |

**Table 4**
**Compare different Stem structures**

|  | TOP-1 ACC | FLOPs |
|---|---|---|
| CSPConvNeXt Stem | 79.2 | 1.5G |
| $ResNet_{vc}$ Stem | 79.7 | 2.0G |
| Stepped Stem | 79.8 | 2.0G |

**Table 5**
**Compare different Downsampling structure**

|  | Downsample LayerNorm | Downsample BatchNorm |
|---|---|---|
| Top-1 Acc | 79.8 | 79.9 |
| FLOPs | 2.0G | 2.0G |

### 4.3. Comparison of different Stem structures

To verify whether the Stepped Stem can alleviate the feature loss caused by continuous downsampling, we compare the performance with the original Stem of CSPConvNeXt and the $ResNet_{vc}$ Stem. The dataset used in this experiment is ImageNet-1k, with the image size fixed at 224x224.

As shown in Table 4, compared with the original Stem of CSPConvNeXt, the $ResNet_{vc}$ Stem improved accuracy by 0.5% while the FLOPs increased 0.5G . We further combined the Patchify Stem with the $ResNet_{vc}$ Stem to obtain the Stepped Stem. It further improved accuracy by 0.1% with the FLOPs unchanged. Finally, we selected the Stepped Stem as the Stem of the network.

### 4.4. Comparison of LN and BN

To verify which normalization layer is more suitable for the downsampling layer of E-ConvNeXt, we replaced LN in the downsampling layer with BN. As shown in Table 5, the network using BN achieves higher accuracy compared to that using LN. Additionally, BN operates more efficiently than LN in the channel first mode. Therefore, BN is selected as the normalization layer for the downsampling layer of CSPConvNeXt.

To ensure network consistency, we replaced the LN with the BN in block. As shown in Table 6, compared to the Block using LN, the Block using BN achieved higher accuracy and 25% speedup without altering the network's computational complexity. This indicates that BN is more suitable for E-ConvNeXt.

**Table 6**
**Compare with different block**

| Evaluation Metrics | LN Block | BN Block |
|---|---|---|
| Top-1 Acc | 79.9 | 80.1 |
| FLOPs | 2.0G | 2.0G |
| FPS | 447 | 549 |

### 4.5. Comparison of Different Channel Attention Modules

We integrated various mainstream channel attention modules (namely SE, CBAM, ECA and ESE) into the baseline model and compared their performances. ImageNet-100 dataset is used in this experiment. The total number of training epochs is set to 60. The experimental results are shown in Table 7. Among all tested modules, ESE Block proformed the best. It achieved the highest accuracy and fastest inference time. Therefore, we add the ESE Block to E-ConvNeXt.

**Table 7**
**Copare with different channel attention**

|  | None | SE Moudle | ECA | CBAM | ESE Layer |
|---|---|---|---|---|---|
| Acc | 82.1 | 83.4 | 84.2 | 83.9 | 85.4 |
| Time | 1.82ms | 19.3ms | 19.3ms | 2.05ms | 1.89ms |

### 4.6. Comparison of state-of-the-arts

To evaluate the proformence of E-ConvNeXt, we conducted a series of experiments on the ImageNet-1k. We compared our model with mainstream models such as ConvNeXt, StarNet, FasterNet, and Swin-Transformer. To evaluate the performance of E-ConvNeXt under different network complexities, we compared E-ConvNeXt-mini, E-ConvNeXt-tiny, and E-ConvNeXt-small with other models of the same FLOPs, respectively. The experimental results are shown in Table 8, where the parameters, FLOPs, and ImageNet-1K Top1-ACC of other methods are all derived from relevant references.

E-ConvNeXt achieves performance comparable to state-of-the-art models. Specifically, E-ConvNeXt-mini, with 0.93G FLOPs and 78.3% accuracy, performed as well as the best model StarNet-S4. Compared with networks such as FasterNet-T2 and MobileViT-S, E-ConvNeXt-tiny achieved the highest accuracy under similar FLOPs. It is able to balance FLOPs and accuracy. Compared with other networks of the same FLOPs, E-ConvNeXt-small belong to the first tier in terms of accuracy. Meanwhile, E-ConvNeXt-small has lower FLOPs.
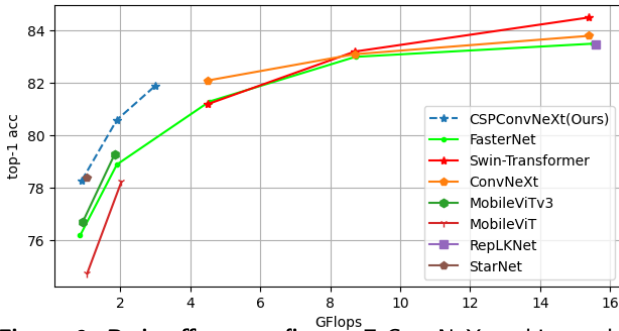
To further demonstrate the superiority of E-ConvNeXt, we plotted the data from Table 8 into trade-off curves as shown in Figure 9. As indicated in Figure 9, E-ConvNeXt establishes a new state-of-the-art in balancing accuracy and FLOPs among all the networks evaluated. The above results

**Table 8**

**Comparison on ImageNet-1k benchmark.** Models with similar FLOPs are grouped together. For each group, our E-ConvNeXt achieves the optimal balance between accuracy and FLOPs.

| model | #param. | FLOPs | IN-1K top-1 acc. |
|---|---|---|---|
| GhostNet ×1.3 [37] | 7.4M | 0.24G | 75.7 |
| ShuffleNetV2 ×2 [38] | 7.4M | 0.24G | 75.7 |
| MobileNetV2 ×1.4 [13] | 6.1M | 0.60G | 74.7 |
| MobileViT-XS [39] | 2.3M | 1.05G | 74.8 |
| EdgeNeXt-XS [40] | 2.3M | 0.54G | 75.0 |
| FasterNet-T1 [15] | 7.6M | 0.85G | 76.2 |
| StarNet-S4 [41] | 7.5M | 1.07G | 78.4 |
| E-ConvNeXt-Mini (Ours) | 7.6M | 0.93G | 78.3 |
| CycleMLP-B1 [42] | 15.2M | 2.10G | 79.1 |
| PoolFormer-S12 [43] | 11.9M | 1.82G | 77.2 |
| MobileViT-S [39] | 5.6M | 2.03G | 78.4 |
| EdgeNeXt-S [40] | 5.6M | 1.26G | 79.4 |
| FasterNet-T2 [15] | 15.0M | 1.91G | 78.9 |
| E-ConvNeXt-Tiny (Ours) | 13.2M | 2.04G | 80.6 |
| ResNet50 | 25.6M | 4.11G | 78.8 |
| CycleMLP-B2 [42] | 26.8M | 3.90G | 81.6 |
| PoolFormer-S24 [43] | 21.4M | 3.41G | 80.3 |
| PoolFormer-S36 [43] | 30.9M | 5.00G | 81.4 |
| PVT-Small [44] | 24.5M | 3.83G | 79.8 |
| PVT-Medium [44] | 44.2M | 6.69G | 81.2 |
| FasterNet-S [40] | 31.1M | 4.56G | 81.3 |
| Swin-T [9] | 28.3M | 4.51G | 81.3 |
| ConvNeXt-T [7] | 28.6M | 4.47G | 82.1 |
| E-ConvNeXt-Small (Ours) | 19.4M | 3.12G | 81.9 |



**Figure 9: Rade-off curves figure.** E-ConvNeXt achieves the highest efficiency in balancing accuracy and FLOPs.

fully demonstrate the application prospects of E-ConvNeXt in lightweight scenarios.

## 5. E-ConvNeXt on downstream tasks

We conducted extensive experiments to verify the performance of E-ConvNeXt in object detection. The experiments were carried out on the sonar image dataset and the DUO dataset[45]. In the experiments,we pretrained E-ConvNeXt on ImageNet-1K and used it as the backbone network. It is then equipped with the PP-YOLOE[46] detector and YOLOv10[47] detector.

### 5.1. Experiment Setup

To be fairness, PP-YOLOE and YOLOv10 use the same hyperparameter settings: 36 training epochs, AdamW optimizer, base learning rate 3.125e-4, batch size 8, CosineDecay scheduler, and LinearWarmup strategy. All experiments are conducted on a server equipped with two NVIDIA GeForce 1080Ti GPUs (11GB memory) and the Ubuntu 20.04 operating system.

The data augmentation methods for PP-YOLOE and YOLOv10 follow those in their original papers: PP-YOLOE uses Random Crop, Random Horizontal Flip, Color Distortion, and Multi-scale; YOLOv10 uses Mosaic, Mixup, and Copy-paste.

### 5.2. Evaluation Metrics

In this experiment, mAP (mean Average Precision) is adopted as the core evaluation metric. mAP is one of the most authoritative metrics for object detection. It is calculated based on Average Precision (AP), specifically as the arithmetic mean of the average precision values across all classes, with the calculation formula as follows:

$$mAP = \frac{1}{N} \sum_{i=1}^{N} AP_i \qquad (3)$$

$N$ is the total number of classes in the dataset, and $AP_i$ denotes the average precision of the $i$-th class.

Average Precision (AP) is referring to the area under the Precision-Recall (P-R) curve for a single class. A higher AP value indicates better detection performance of the model on that class. Precision (P) and recall (R) are the core components of the P-R curve. The calculation formula of precision as follows:

$$P = \frac{TP}{TP + FP}$$

where $TP$ (True Positive) denotes correctly detected positive samples, and $FP$ (False Positive) denotes negative samples incorrectly detected as positive ones.

The calculation formula as follows:

$$R = \frac{TP}{TP + FN}$$

where $FN$ (False Negative) denotes actual positive samples that are not detected. Together, they form the vertical and horizontal axes of the P-R curve, serving as the core inputs for AP calculation.

### 5.3. Object detection method for underwater sonar images

In this section, we conduct extensive experiments to evaluate the effectiveness of E-ConvNeXt as the backbone network for PP-YOLOE and YOLOv10.

We use the forward sonar data from Ocean Space Environment Awareness (Orca) open-source project[2]. There are 5,000 images in total, of which 4,000 for training and 1,000 for test.Table 9 and 10 shows details about object categories and the number of sonar images for each category.

---

[2]https://github.com/violetweir/Sonor_dataset

**Table 9**
Sonor object information in training data

| Object category | Number of objects | Number of images |
|---|---|---|
| Ball | 1943 | 1941 |
| Circle cage | 386 | 383 |
| Cube | 1752 | 1749 |
| Cylinder | 402 | 401 |
| Human body | 684 | 683 |
| Metal bucket | 403 | 402 |
| Square cage | 655 | 655 |
| Tyre | 852 | 850 |
| ToTal | 7077 | 4000 |

**Table 10**
Sonor object information in training data

| Object category | Number of objects | Number of images |
|---|---|---|
| Ball | 595 | 595 |
| Circle cage | 86 | 86 |
| Cube | 424 | 423 |
| Cylinder | 39 | 39 |
| Human body | 379 | 379 |
| Metal bucket | 44 | 43 |
| Square cage | 169 | 169 |
| Tyre | 108 | 108 |
| ToTal | 1844 | 1000 |

### 5.3.1. Results

**Table 11**
Performance comparison of different backbones (including PaNet, ConvNeXt-T, and E-ConvNeXt-T) integrated with YOLOv10-L and YOLOv10-M on the sonar image dataset, evaluated by metrics: FLOPs, FPS, mAP, $mAP^m$ (medium-sized objects), and $mAP^l$ (large-sized objects).

| Backbone | FLOPs | FPS | mAP | mAP$^m$ | mAP$^l$ |
|---|---|---|---|---|---|
| YOLOv10-L PaNet | 120.3 | 22.4 | 46.1 | 46.0 | 41.5 |
| YOLOv10-M ConvNeXt-T | 94.3 | 23.1 | 50.1 | 50.2 | 32.8 |
| YOLOv10-M PaNet | 59.1 | 28.2 | 45.6 | 45.5 | 43.6 |
| YOLOv10-M E-ConvNeXt-T | 69.9 | 26.9 | 51.3 | 51.4 | 39.0 |

To compare the performance of networks in different FLOPs, we selected YOLOv10-L and YOLOv10-M. As shown in Table 11, although YOLOv10-L has 0.5% higher mAP than YOLOv10-M, YOLOv10-L's FLOPs are 103% higher than YOLOv10-M, and the large increase in FLOPs only brings a slight improvement in accuracy. When the backbone of YOLOv10-M is replaced with ConvNeXt-Tiny and E-ConvNeXt-Tiny, the accuracy of YOLOv10 improves compared to the original YOLOv10-L and YOLOv10-M.

Among the above models, E-ConvNeXt-tiny as the backbone of YOLOv10 achieves the highest mAP with the small FLOPs. The mAP of YOLOv10 is significantly improved after being equipped with E-ConvNeXt-tiny and ConvNeXt-tiny. This demonstrates the effectiveness of the ConvNeXt structure. Meanwhile, YOLOv10 with E-ConvNeXt-tiny achieves a higher mAP and lower FLOPs compared to YOLOv10 with ConvNeXt-tiny. This further illustrates the effectiveness of E-ConvNeXt as a backbone in the field of sonar images.

**Table 12**
Performance comparison of different backbones (including E-ConvNeXt-Mini, E-ConvNeXt-Tint) integrated with PPYOLOE-S and PPYOLOE-L on the sonar image dataset, evaluated by metrics: FLOPs, FPS, mAP, $mAP^m$ (medium-sized objects), and $mAP^l$ (large-sized objects).

| Backbone | FLOPs | FPS | mAP | mAP$^m$ | mAP$^l$ |
|---|---|---|---|---|---|
| PPYOLOE-S | 17.4 | 46.0 | 42.6 | 42.0 | 28.1 |
| PPYOLOE-S E-ConvNeXt-Mini | 20.56 | 43.2 | 50.6 | 50.8 | 44.8 |
| PPYOLOE-L | 110.7 | 33.8 | 49.0 | 48.8 | 35.5 |
| PPYOLOE-L E-ConvNeXt-Tiny | 82.6 | 36.3 | 51.3 | 51.3 | 44.4 |

To compare the performance of networks in different FLOPs, we selected two versions of PP-YOLOE: PP-YOLOE-S and PPYOLO-L. As shown in Table 12, the mAP steadily increases as the network complexity from small to large. By comparing the performances in terms of $mAP^l$ and $mAP^m$, we find that PP-YOLOE behaves poorly in recognizing large target objects

When using E-ConvNeXt-Mini as the backbone for PP-YOLOE-S, the mAP reached 50.6 and $mAP^L$ was 44.8. When using CSPConvNeXt-Tiny as the backbone for PP-YOLOE-L, the mAP reached 51.3 and $mAP^L$ was 44.4. Compared with the original PP-YOLOE, PP-YOLOE with E-ConvNeXt backbone shows significant improvements in sonar image object detection, especially fot large target objects.

## 5.4. Object detection for underwater optical images

In Section 5.3, compared with PP-YOLOE, YOLOv10 with E-ConvNeXt achieves better performance. Therefore, we use YOLOv10 with E-ConvNeXt in this experiment.

Detecting Underwater Objects (DUO)[3] dataset contains 15632 optical images. It is a typical object detection dataset which has been widely used. As shown in Table 13 and Table 14, DUO includes 74903 annotated objects across four common categories. Therefore, we select DUO as this section dataset.

**Table 13**
Object information in training data

| Object category | Number of objects | Number of images |
|---|---|---|
| holothurian | 6808 | 3001 |
| echinus | 42955 | 5872 |
| scallop | 1707 | 487 |
| starfish | 12528 | 4078 |
| total | 63998 | 13438 |

**Table 14**
Object information in test data

| Object category | Number of objects | Number of images |
|---|---|---|
| holothurian | 1079 | 490 |
| echinus | 7201 | 967 |
| scallop | 217 | 78 |
| starfish | 2020 | 659 |
| total | 10517 | 2194 |

**Table 15**
Performance comparison of different backbones (PaNet and E-ConvNeXt-T) with YOLOv10-L and YOLOv10-M on the COCO underwater dataset, including metrics of FLOPs, FPS, mAP, $mAP^m$ (medium object mAP), and $mAP^l$ (large object mAP).

| Backbone | FLOPs | FPS | mAP | $mAP^m$ | $mAP^l$ |
|---|---|---|---|---|---|
| YOLOv10-L PaNet | 120.3 | 22.4 | 56.4 | 40.7 | 62.0 |
| YOLOv10-L E-ConvNeXt-T | 78.0 | 25.4 | 61.2 | 49.7 | 70.1 |
| YOLOv10-M PaNet | 59.1 | 28.2 | 56.2 | 41.9 | 60.7 |
| YOLov10-M E-ConvNeXt-T | 36.7 | 32.8 | 60.1 | 56.0 | 67.0 |

#### 5.4.1. Results

As shown in Table 15, we selected two versions of YOLOv10: YOLOv10-L and YOLOv10-M. The FLOPs of large version is twice that of medium version, but the mAP does not differ significantly. This is mainly because the original YOLOv10 typically requires hundreds of training epochs. To be fairness, all models were trained with 36 epochs in this experiment.

Using E-ConvNeXt as the backbone, YOLOv10 achieves better performance on underwater optical image detection task. Specifically, the mAP of the YOLOv10-L with E-ConvNeXt-tiny is 61.2, which is 5.3 higher than that of the original YOLOv10-L. The mAP of the YOLOv10-M with E-ConvNeXt-tiny also shows an improvement compared with the original YOLOv10-M. This result indicates that E-ConvNeXt has significant advantages in feature extraction tasks for optical images, and it is particularly suitable for

---

³https://github.com/chongweiliu/DUO

lightweight application scenarios such as object detection tasks in complex optical environments.

## 6. Conclusion

This paper presents the E-ConvNeXt framework, aiming to balance performance and efficiency, and address the limitations of existing ConvNeXt-based models in lightweight application scenarios. The proposed integration of CSPNet with ConvNeXt, along with optimized Stem and Block structures and the introduction of the ESE Block, significantly reduces network complexity while enhancing feature expression capability and operational efficiency . The training and testing strategies further improve performance and efficiency. Extensive experiments on the ImageNet-1K dataset, as well as transfer learning tests on object detection tasks, validate the effectiveness, efficiency, and transferability of this method in lightweight scenarios and downstream tasks. Our work encourages further exploration across tasks, with codes available at https://github.com/violetweir/E-ConvNeXt

**Limitations and Future Work.** E-ConvNeXt, despite its advantages in balancing accuracy and efficiency, still has certain limitations. For instance, its performance in more extreme lightweight scenarios (e.g., ultra-low power consumption embedded devices) needs further verification and improvement. In the future, we will focus on optimizing the model structure to further reduce parameters and computational complexity while maintaining accuracy, to better adapt to extreme lightweight application scenarios.

## References

[1] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Back-propagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.

[2] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.

[3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[4] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.

[5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[7] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022.

[8] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam.

Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[9] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.

[10] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.

[11] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127(3):302–321, 2019.

[12] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[13] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.

[14] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1314–1324, 2019.

[15] Jierun Chen, Shiu-hong Kao, Hao He, Weipeng Zhuo, Song Wen, Chul-Ho Lee, and S-H Gary Chan. Run, don't walk: chasing higher flops for faster neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12021–12031, 2023.

[16] Chien-Yao Wang, Hong-Yuan Mark Liao, Yueh-Hua Wu, Ping-Yang Chen, Jun-Wei Hsieh, and I-Hau Yeh. Cspnet: A new backbone that can enhance learning capability of cnn. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 390–391, 2020.

[17] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shuflenet: An extremely efficient convolutional neural network for mobile devices, 2017.

[18] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shuflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6848–6856, 2018.

[19] Mingxing Tan and Quoc Le. Efficientnetv2: Smaller models and faster training. In *International conference on machine learning*, pages 10096–10106. PMLR, 2021.

[20] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.

[21] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017.

[22] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.

[23] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11534–11542, 2020.

[24] Zilin Gao, Jiangtao Xie, Qilong Wang, and Peihua Li. Global second-order pooling convolutional networks. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 3024–3033, 2019.

[25] HyunJae Lee, Hyo-Eun Kim, and Hyeonseob Nam. Srm: A style-based recalibration module for convolutional neural networks. In *Proceedings of the IEEE/CVF International conference on computer vision*, pages 1854–1862, 2019.

[26] Zequn Qin, Pengyi Zhang, Fei Wu, and Xi Li. Fcanet: Frequency channel attention networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 783–792, 2021.

[27] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.

[28] Sergey Ioffe. Batch renormalization: Towards reducing minibatch dependence in batch-normalized models. *Advances in neural information processing systems*, 30, 2017.

[29] Youngwan Lee and Jongyoul Park. Centermask: Real-time anchor-free instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13906–13915, 2020.

[30] Régis Vaillant, Christophe Monrocq, and Yann Le Cun. Original approach for the localisation of objects in images. *IEE Proceedings-Vision, Image and Signal Processing*, 141(4):245–250, 1994.

[31] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[32] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020.

[33] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.

[34] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019.

[35] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 13001–13008, 2020.

[36] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.

[37] Kai Han, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. Ghostnet: More features from cheap operations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1580–1589, 2020.

[38] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shuflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European conference on computer vision (ECCV)*, pages 116–131, 2018.

[39] Sachin Mehta and Mohammad Rastegari. Mobilevit: light-weight, general-purpose, and mobile-friendly vision transformer. *arXiv preprint arXiv:2110.02178*, 2021.

[40] Muhammad Maaz, Abdelrahman Shaker, Hisham Cholakkal, Salman Khan, Syed Waqas Zamir, Rao Muhammad Anwer, and Fahad Shahbaz Khan. Edgenext: efficiently amalgamated cnn-transformer architecture for mobile vision applications. In *European conference on computer vision*, pages 3–20. Springer, 2022.

[41] Xu Ma, Xiyang Dai, Yue Bai, Yizhou Wang, and Yun Fu. Rewrite the stars. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5694–5703, 2024.

[42] S Chen, E Xie, C Ge, R Chen, D Liang, and P Luo. Cyclemlp: A mlp-like architecture for dense prediction. arxiv 2021. *arXiv preprint arXiv:2107.10224*.

[43] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what

you need for vision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10819–10829, 2022.

[44] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 568–578, 2021.

[45] Chongwei Liu, Haojie Li, Shuchang Wang, Ming Zhu, Dong Wang, Xin Fan, and Zhihui Wang. A Dataset And Benchmark Of Underwater Object Detection For Robot Picking. *arXiv e-prints*, page arXiv:2106.05681, June 2021.

[46] Shangliang Xu, Xinxin Wang, Wenyu Lv, Qinyao Chang, Cheng Cui, Kaipeng Deng, Guanzhong Wang, Qingqing Dang, Shengyu Wei, Yuning Du, et al. Pp-yoloe: An evolved version of yolo. *arXiv preprint arXiv:2203.16250*, 2022.

[47] Ao Wang, Hui Chen, Lihao Liu, Kai Chen, Zijia Lin, Jungong Han, et al. Yolov10: Real-time end-to-end object detection. *Advances in Neural Information Processing Systems*, 37:107984–108011, 2024.