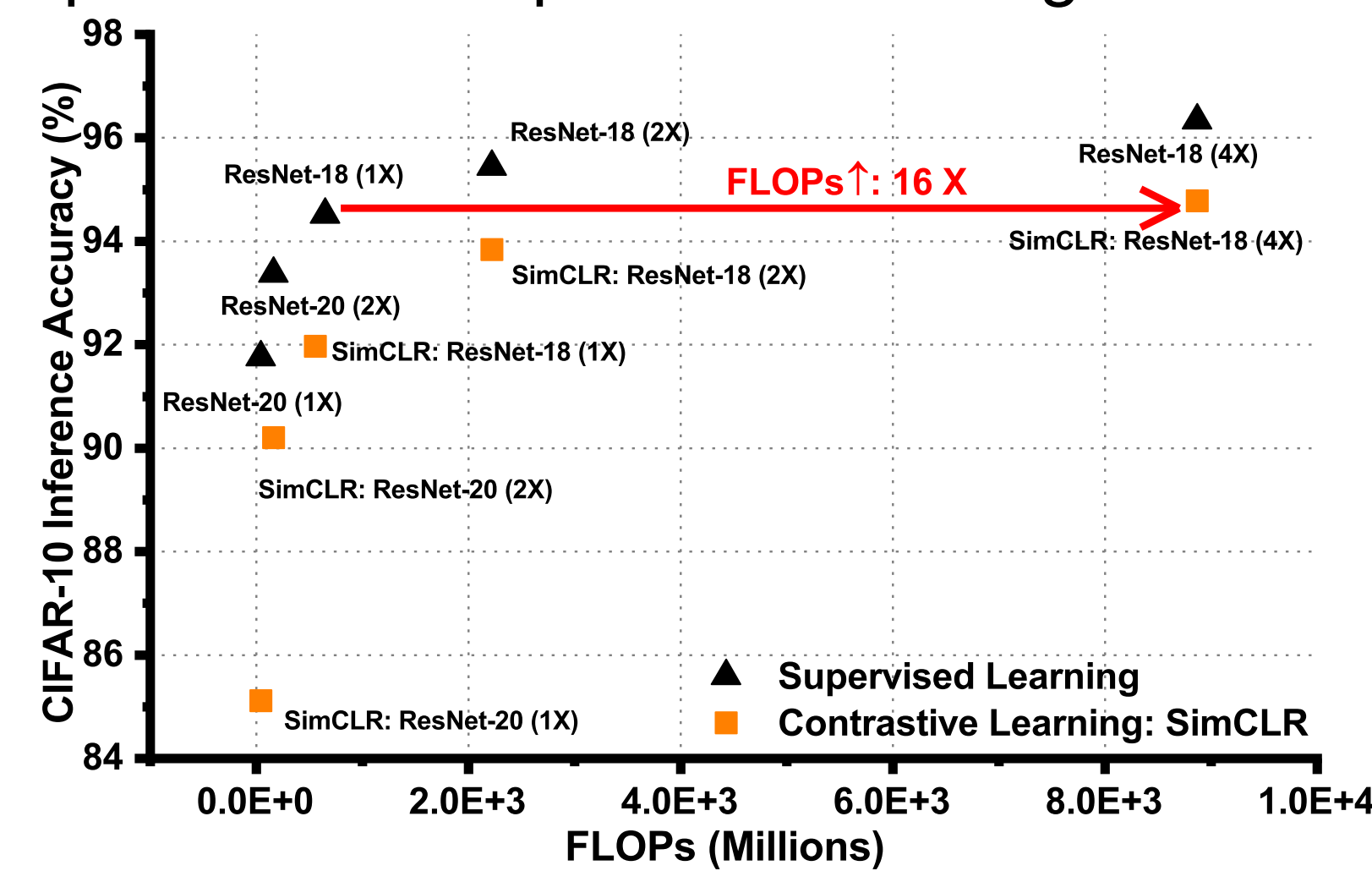


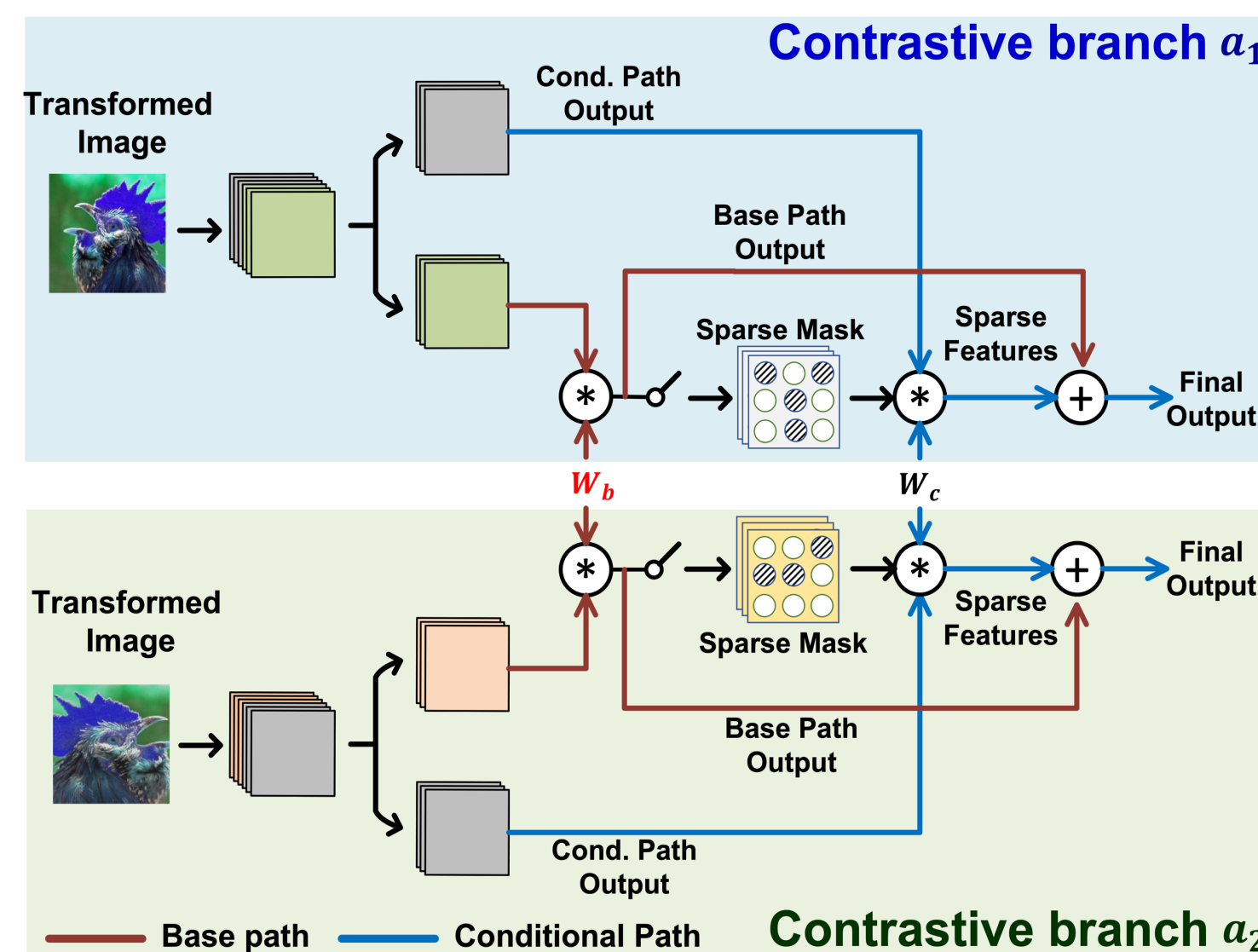
## Introduction

- Recent contrastive learning-based self-supervised learning works require wide and deep models (e.g., 4X wider) to achieve comparable performance as supervised training works (1X).
- The extraordinary computation cost necessitates efficient computation reduction techniques for self-supervised learning.



## Contrastive Gating

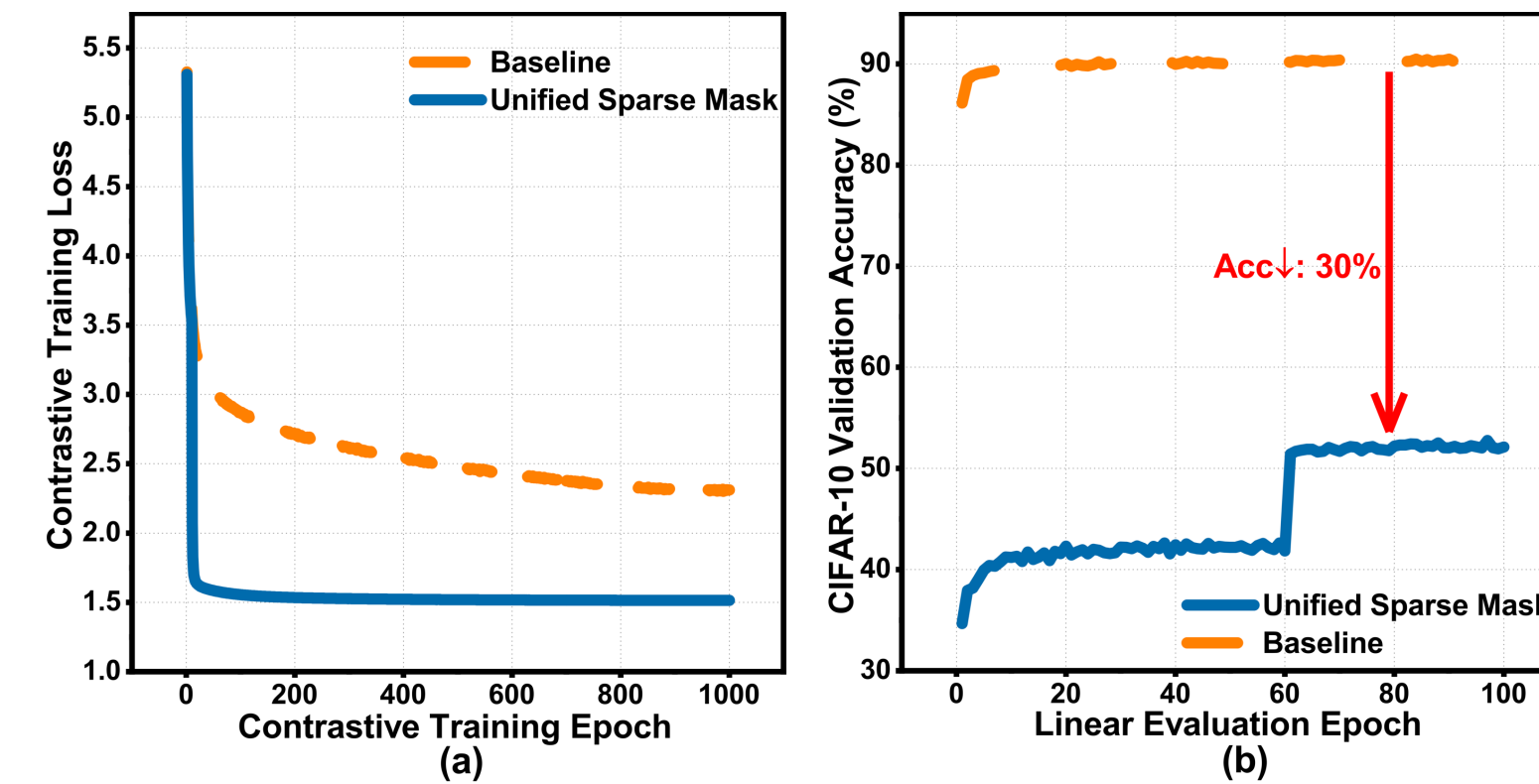
- Learning the sparse features in both contrastive branches during the unsupervised learning process.



- For each layer of the encoder, the input feature maps and weights  $(X, W)$  are divided into base  $(X_b, W_b)$  and conditional  $(X_c, W_c)$  paths.
- The saliency of the base path output determines the computation skipping decision  $M_c \in \{0,1\}$  of the conditional path.

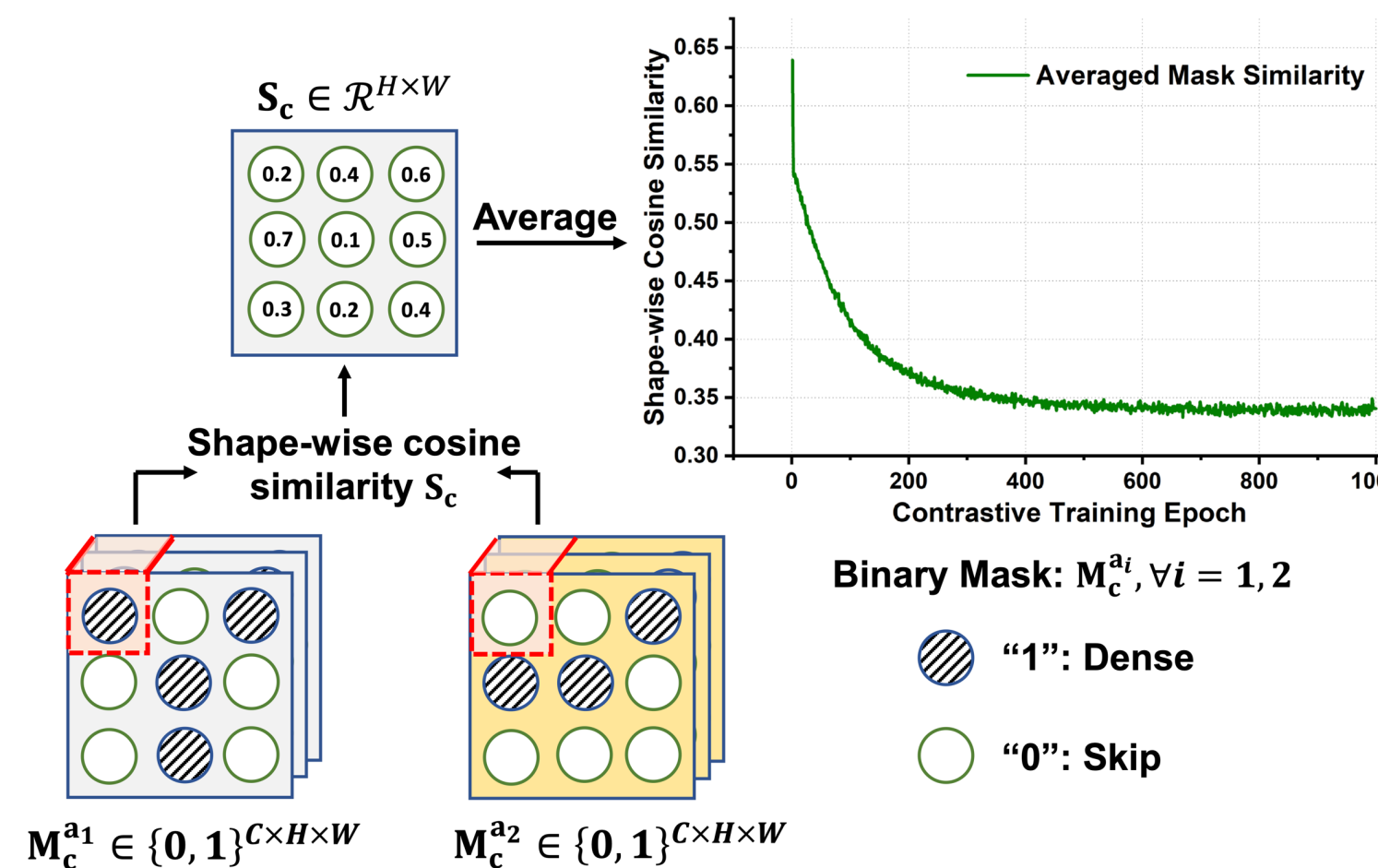
## SSL-based Sparse Feature Learning

- Non-transferability of dynamic sparse feature masks



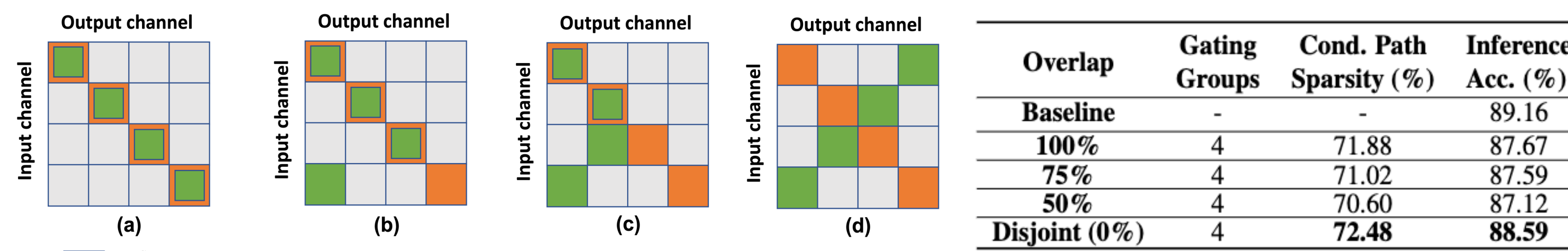
- Broadcasting the identical sparse feature mask  $M_c$  to both SimCLR [1] contrastive paths results in: (a) reduced contrastive training loss, and (b) defective generalizability with unsuccessful supervised linear evaluation.

- Disjoint sparse contrastive features



- Given the unanimous data transformation and identical base path, contrastive training encourages the encoders to highlight different contrastive features.

- Unbiased contrastive grouping



- Evenly activating the disjoint channels among the different contrastive paths will enhance the sparse feature learning during contrastive training.

## Contrastive Dual Gating (CDG)

- During the forward pass of the contrastive training, CDG generates pruning masks  $M_c^{a1}$  and  $M_c^{a2}$  for both contrastive branches:

$$M_c^{a_i} = \sigma(\text{normalize}(X_b^{a_i} * W_b^{a_i}) - \tau)$$

- The contrastive branches are selected along the diagonal and inverse-diagonal of the channel groups.  $\tau$  learns the gating decision during training.

## Structured Contrastive Dual Gating (SCDG)

- During the forward pass of the contrastive training, SCDG generates structured pruning masks  $M_c^{a1}$  and  $M_c^{a2}$  based on the averaged saliency:

$$M_{sc}^{a_i} = \sigma(\text{normalize}(\text{AvgPool}(X_b^{a_i} * W_b^{a_i})) - \tau)$$

## Experimental Results

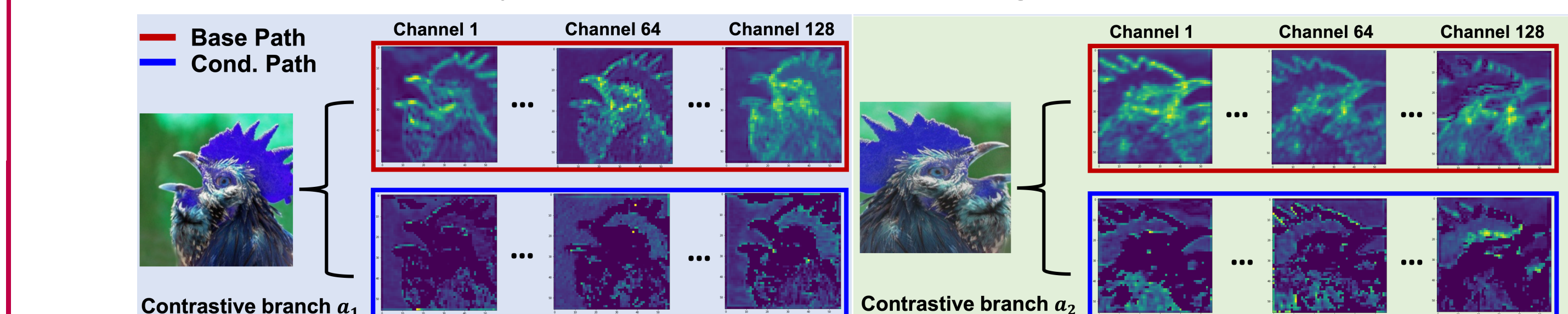
- The widely-used mini-neural network-based auxiliary saliency predictors (e.g., FBS [2], DGC [3]) are difficult to train from scratch, resulting in degraded inference accuracy.

Method	# of Gating Groups	Linear Eval. Inference Accuracy (%)	Fine-tuning Inference Accuracy (%)	FLOPs Reduction
This work (CDG-SimCLR)	4	88.84	90.74	2.12x
FBS-SimCLR	-	86.91	88.89	2.00x
DGC-SimCLR	4	73.10	81.77	2.11x
CGNet-SimCLR	4	87.40	89.26	2.09x

- Structured-CDG (SCDG) results with spatial feature group size =  $8 \times 1 \times 1$

Model	# of Gating Groups	Dataset	Conditional Path Sparsity (%)	Inference Acc. (%)	Top-1 Acc. Drop	FLOPs Reduction	Index Reduction
ResNet-18 (1x)	4	CIFAR-10	71.64	90.37	-0.89	2.16x	8x
		CIFAR-100	66.24	65.94	-1.84	1.98x	8x
		ImageNet-100	45.52	76.63	-2.24	1.53x	8x

- Similar accuracy as CDG with 8x sparse index reduction.
- Base paths preserves the details with dense convolution, while the sparse conditional path only keeps the important edges.



(CDG with other datasets and frameworks are presented in main paper.)