

CCF-CV走进高校系列报告会
西安电子科技大学, 2016.12.28

Deep Fusion: A Neural Network Architecture Design Pattern

Jingdong Wang
Lead Researcher
Microsoft Research, Beijing, China

Deep learning in the past 10 years

- Reducing the dimensionality of data with neural networks, Science, 2006
 - Fast learning algorithms for Restricted Boltzmann machine

Do not work as expected

- ImageNet Classification with deep convolutional neural networks, NIPS, 2012
 - Dramatic performance improvement
 - ImageNet, GPU

Win almost in all the applications

Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such "autoencoder" networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

Dimensionality reduction facilitates the classification, visualization, communication, and storage of high-dimensional data. A simple and widely used method is principal components analysis (PCA), which finds the directions of greatest variance in the data set and represents each data point by its coordinates along each of these directions. We describe a nonlinear generalization of PCA that uses an adaptive, multilayer "encoder" network

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

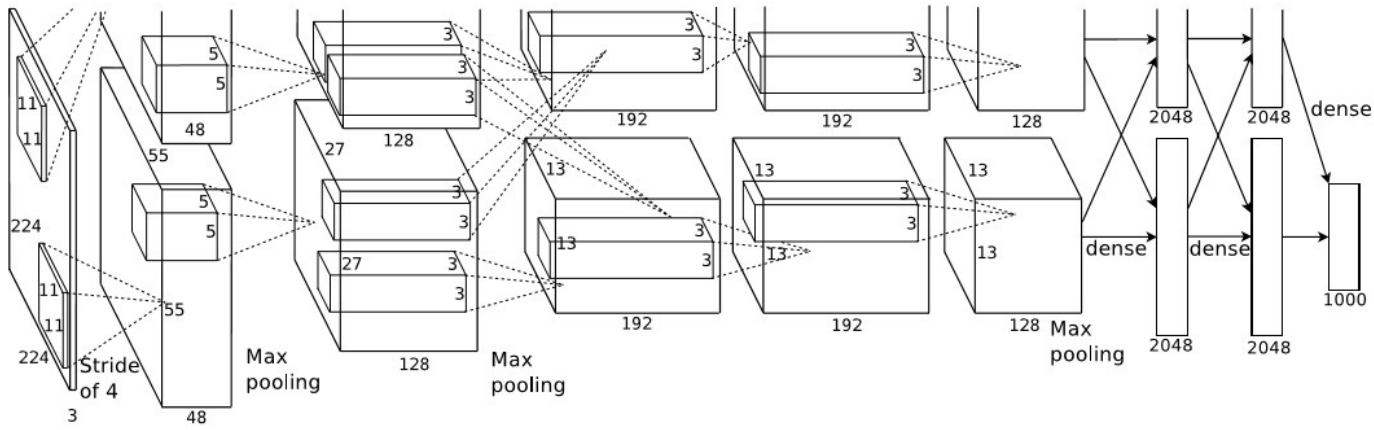
Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

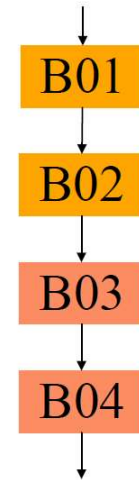
Abstract

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet ILSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called "dropout" that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

AlexNet, 2012



7 layers



AlexNet
2012

GoogleNet
2014

ResNet
2015

FractalNet
2016

DFN-Merge
and Run
2016

VGGNet
2014

Highway
2015

Deeply-Fused Net
2016

DenseNet
2016

.....

Dropout: A Simple Way to Prevent Neural Networks from Overfitting

Nitish Srivastava

Geoffrey Hinton

Alex Krizhevsky

Ilya Sutskever

Ruslan Salakhutdinov

Department of Computer Science

University of Toronto

10 Kings College Road, Rm 3302

Toronto, Ontario, M5S 3G4, Canada.

NITISH@CS.TORONTO.EDU

HINTON@CS.TORONTO.EDU

KRIZ@CS.TORONTO.EDU

ILYA@CS.TORONTO.EDU

RSALAKHU@CS.TORONTO.EDU

Editor: Yoshua Bengio

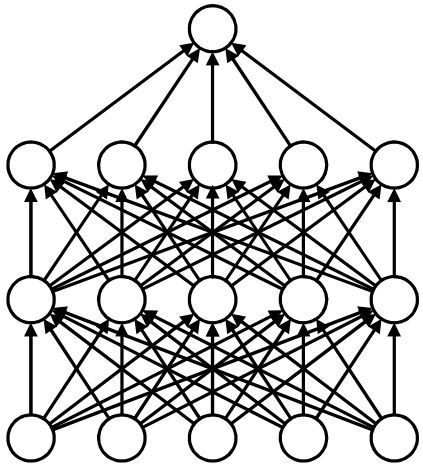
Dropout \approx Ensembling

Abstract

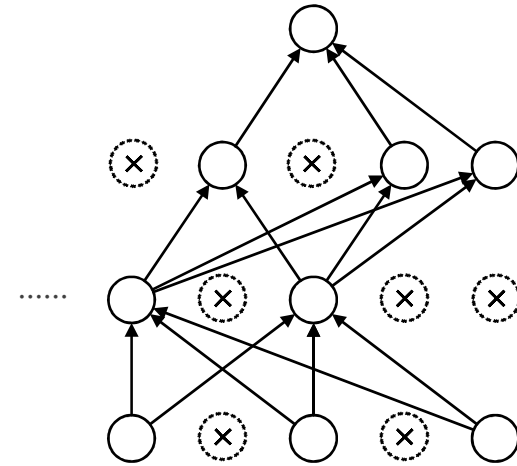
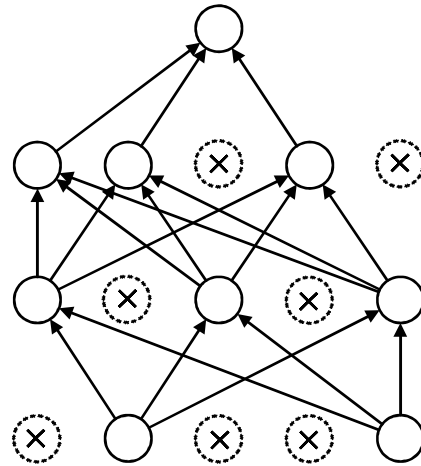
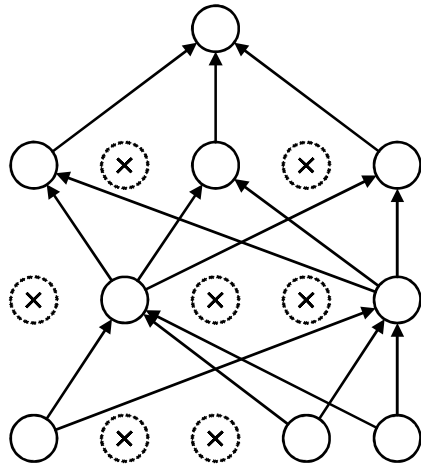
Deep neural nets with a large number of parameters are very powerful machine learning systems. However, overfitting is a serious problem in such networks. Large networks are also slow to use, making it difficult to deal with overfitting by combining the predictions of many different large neural nets at test time. Dropout is a technique for addressing this problem. The key idea is to randomly drop units (along with their connections) from the neural network during training. This prevents units from co-adapting too much. During training, dropout samples from an exponential number of different “thinned” networks. At test time, it is easy to approximate the effect of averaging the predictions of all these thinned networks by simply using a single unthinned network that has smaller weights. This significantly reduces overfitting and gives major improvements over other regularization methods. We show that dropout improves the performance of neural networks on supervised learning tasks in vision, speech recognition, document classification and computational biology, obtaining state-of-the-art results on many benchmark data sets.

Keywords: neural networks, regularization, model combination, deep learning

Dropout



Predict with an approximate ensemble

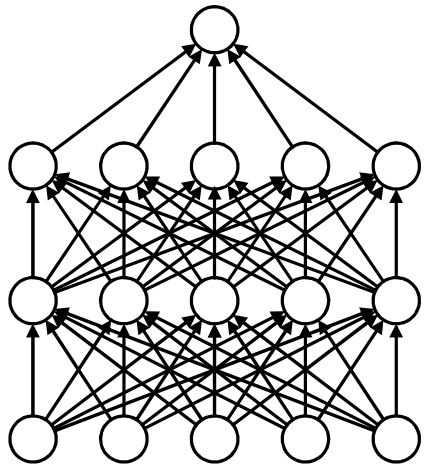


Train thinned nets generated from dropout: Each iteration trains a thinned net

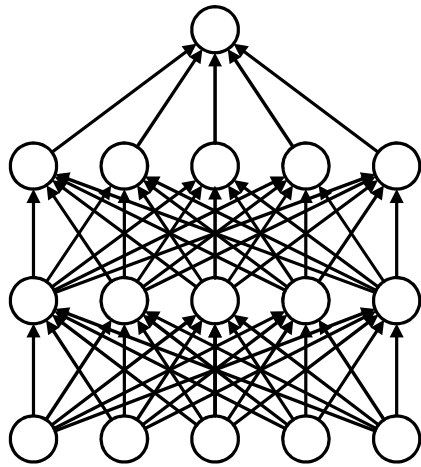
Train exponential number of thinned networks with the same training data

Predict from one network like **ensembling** an exponential number of thinned networks

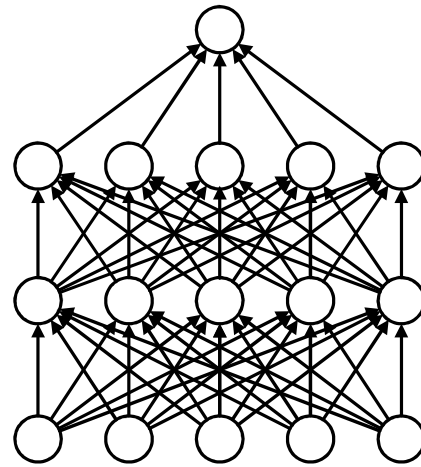
DisturbLabel



Predict with
an approximate ensemble

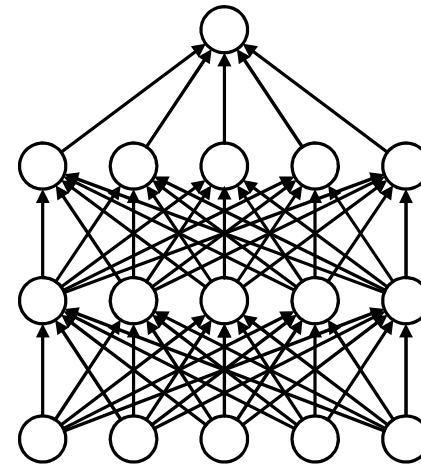


$\{(\mathbf{x}_n, \tilde{\mathbf{y}}_n^1)\}$



$\{(\mathbf{x}_n, \tilde{\mathbf{y}}_n^2)\}$

.....



$\{(\mathbf{x}_n, \tilde{\mathbf{y}}_n^t)\}$

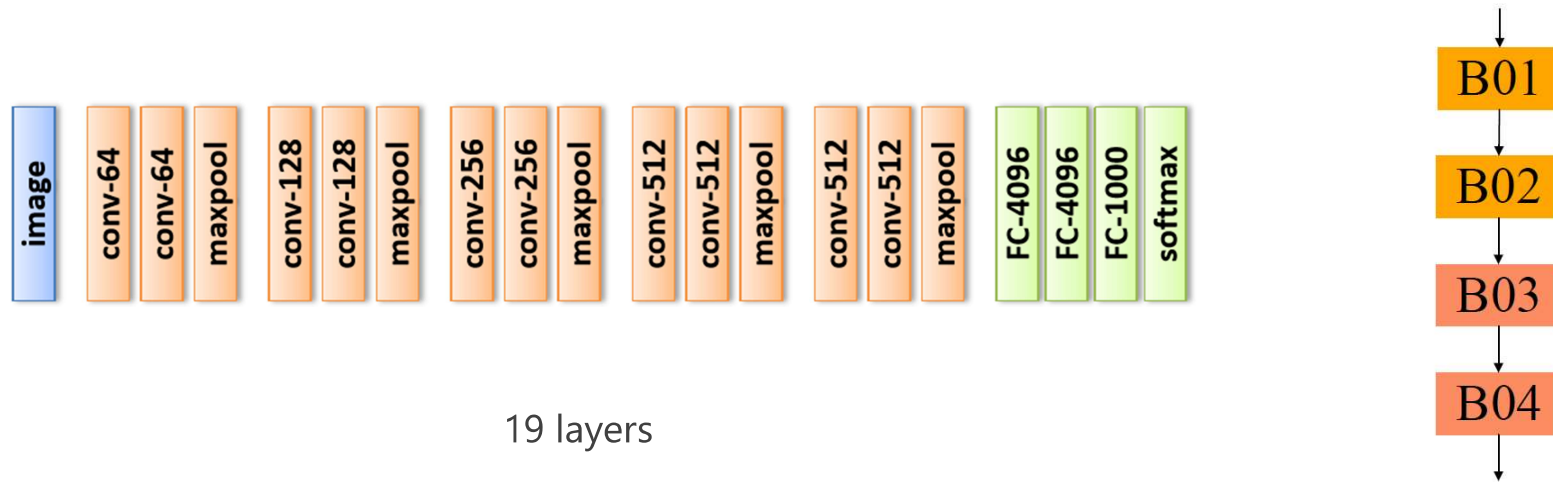
Train the network: with a different training set in each iteration

Train the network of same structure with exponential number of different training data

Predict from one network like **ensembling** an exponential number of networks of same structure

[Lingxi Xie](#), Jingdong Wang, [Zhen Wei](#), [Meng Wang](#), [Qi Tian](#): DisturbLabel: Regularizing CNN on the Loss Layer. CVPR (2016)

VGGNet, 2014



AlexNet
2012

GoogleNet
2014

ResNet
2015

FractalNet
2016

DFN-Merge
and Run
2016

VGGNet
2014

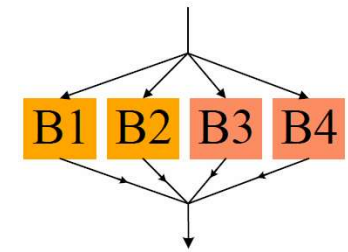
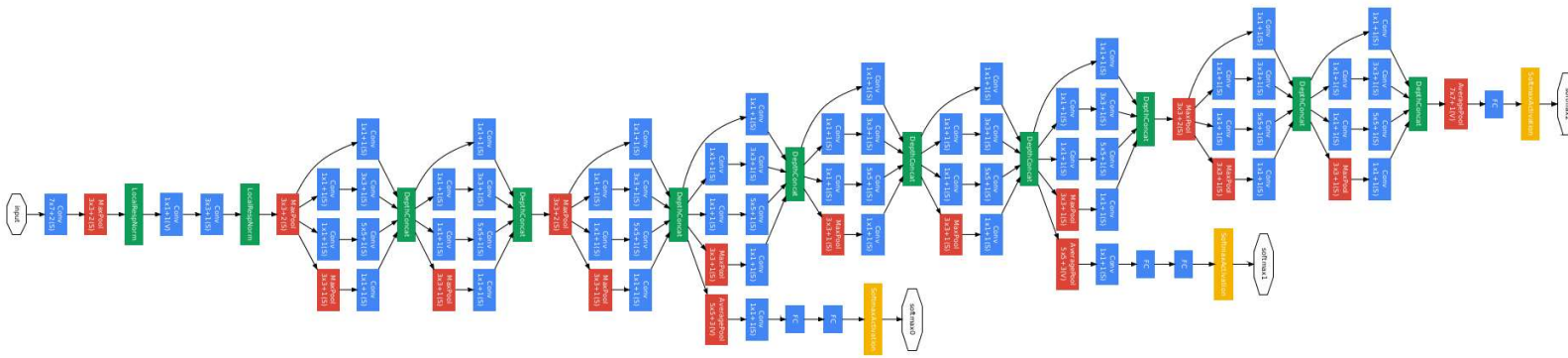
Highway
2015

Deeply-Fused Net
2016

DenseNet
2016

.....

GoogLeNet, 2014



22 layers, Inception: multi-branch, concatenation fusion

AlexNet
2012

GoogleNet
2014

ResNet
2015

FractalNet
2016

DFN-Merge
and Run
2016

VGGNet
2014

Highway
2015

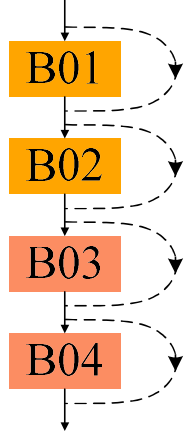
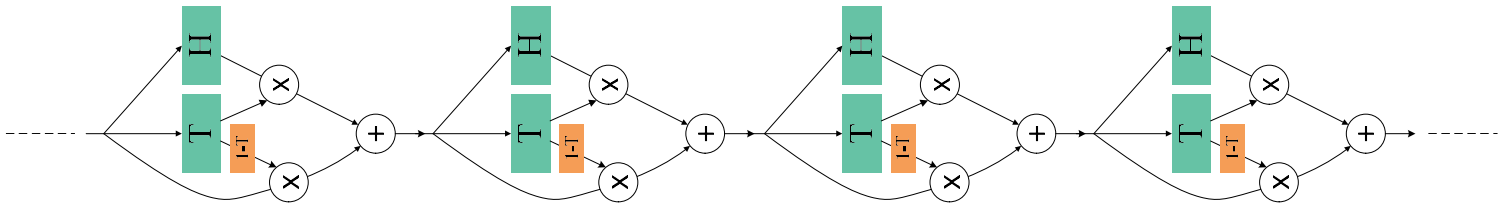
Deeply-Fused Net
2016

DenseNet
2016

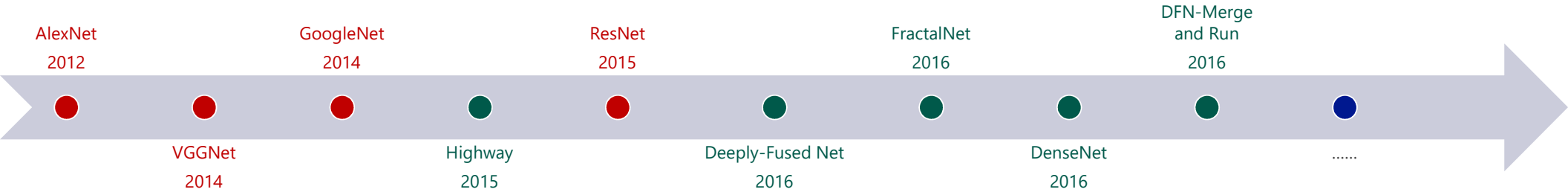
.....



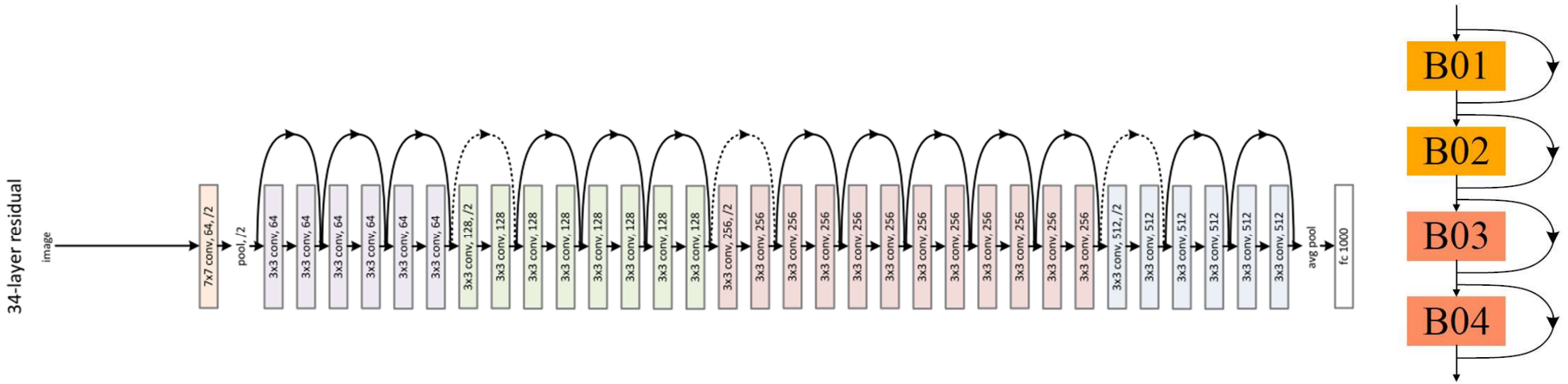
Highway, 2015



100+ layers, skip connection **eases** the **training** of a very deep network



ResNet, 2015



AlexNet
2012

GoogleNet
2014

ResNet
2015

FractalNet
2016

DFN-Merge
and Run
2016

VGGNet
2014

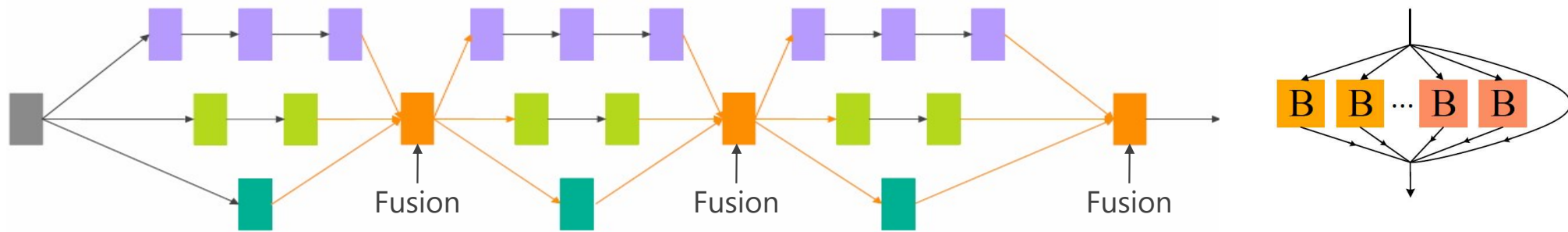
Highway
2015

Deeply-Fused Net
2016

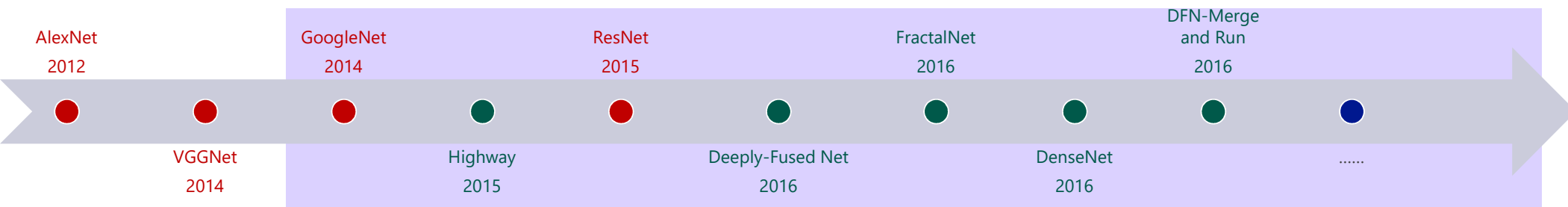
DenseNet
2016

.....

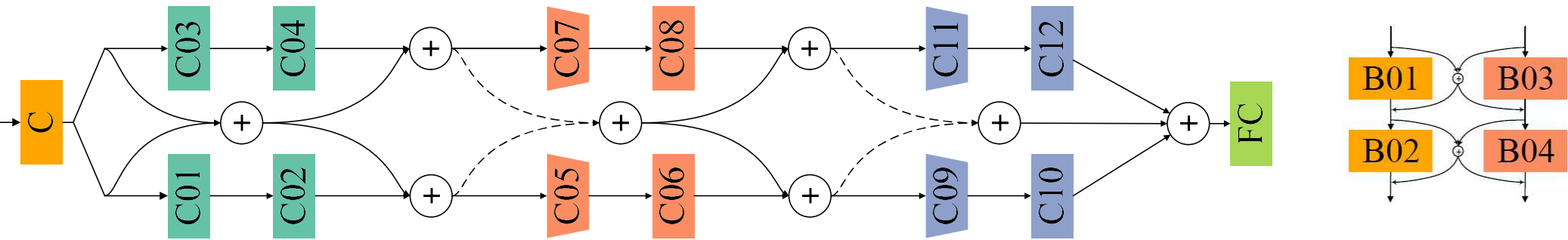
Deeply-fused net, 2016



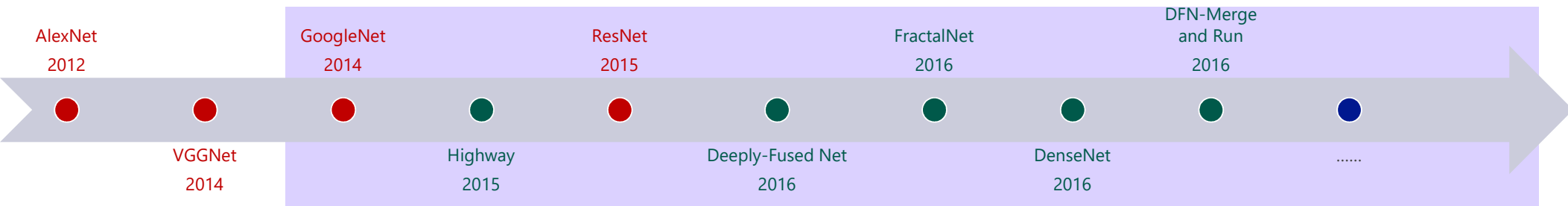
Multiple weight-shared paths: **Unifying** GoogLeNet, Highway, ResNet



Merge and run, 2016



Deep fusion by merge and run: **shallower could be better**

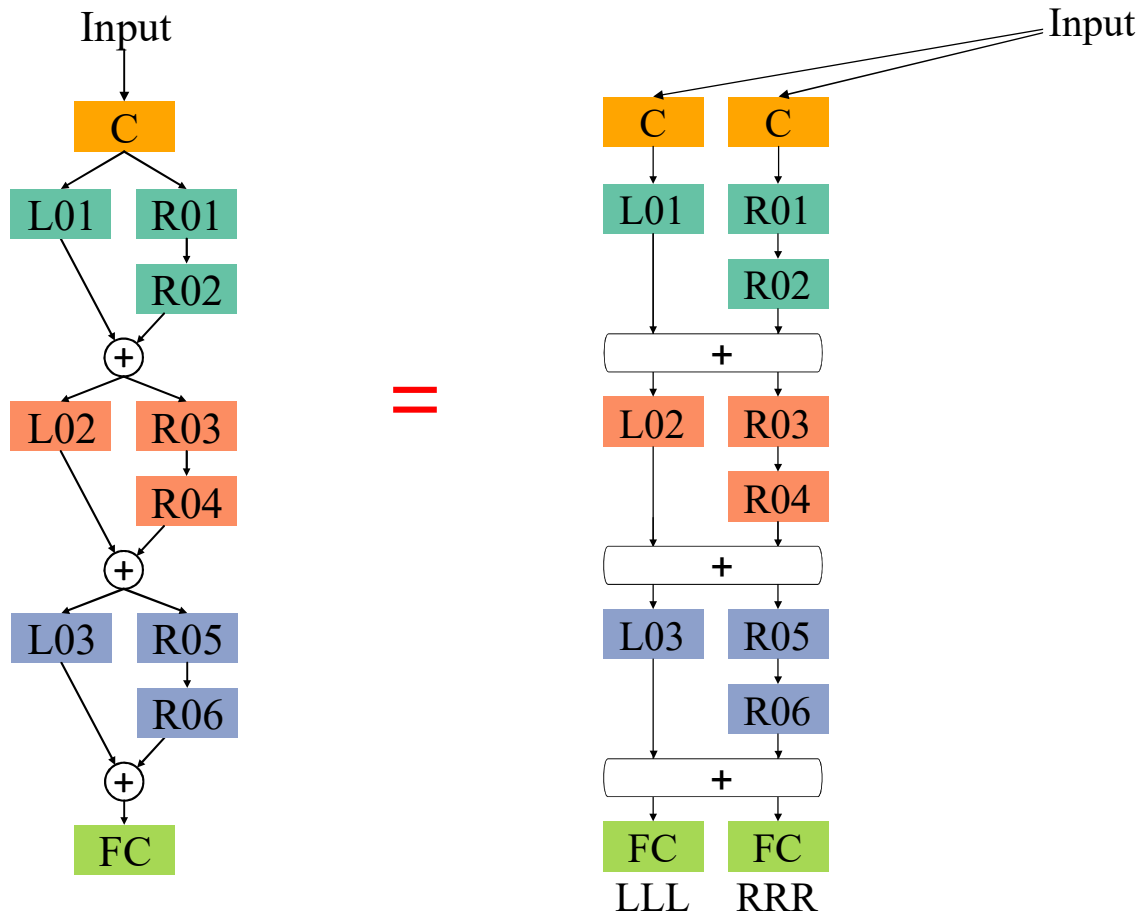


Deeply-Fused Nets

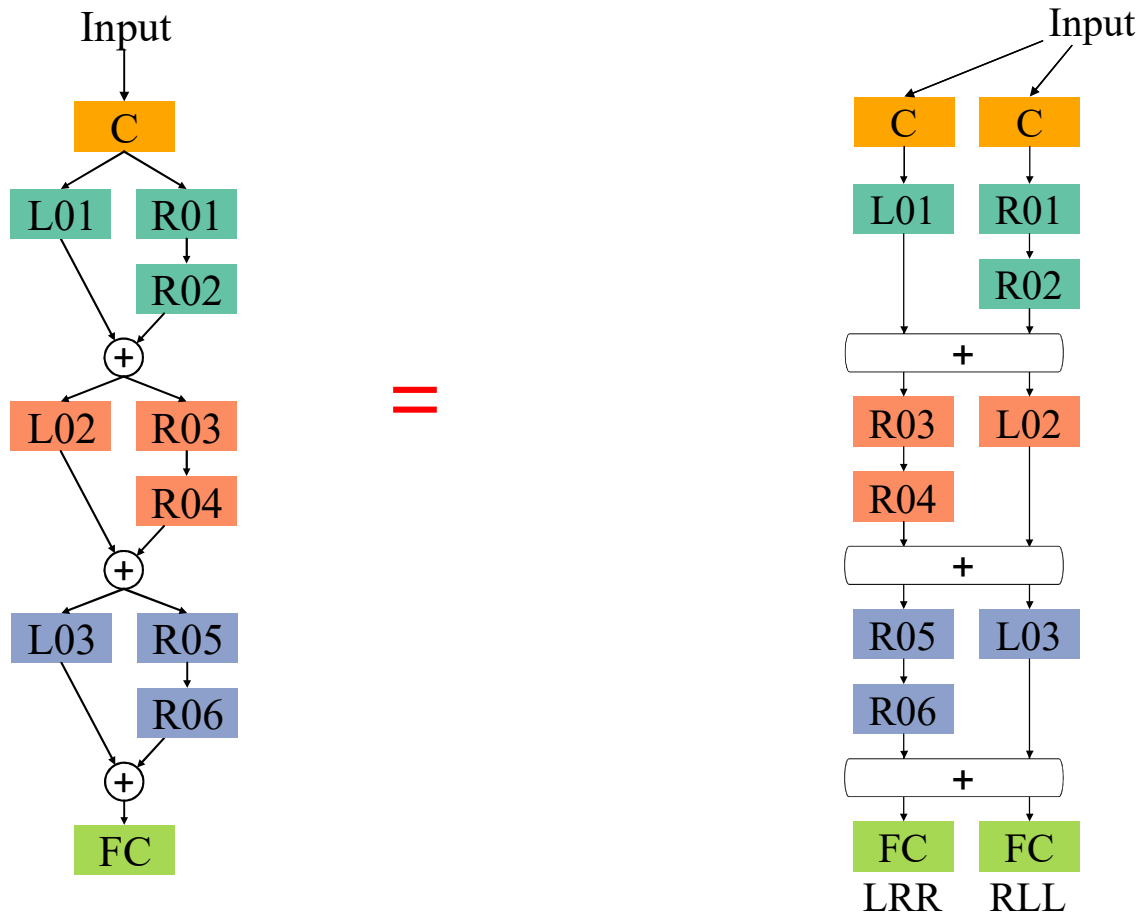
Jingdong Wang, [Zhen Wei](#), [Ting Zhang](#), [Wenjun Zeng](#): Deeply-Fused Nets. [CoRR abs/1605.07716](#) (2016)

A deeply-fused net can be formed from **many different** base networks

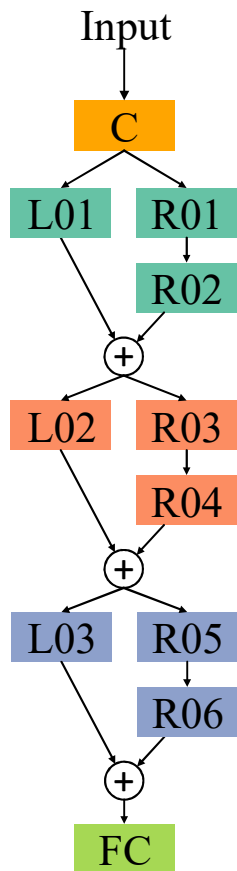
A deeply-fused net can be formed from **many different** base networks



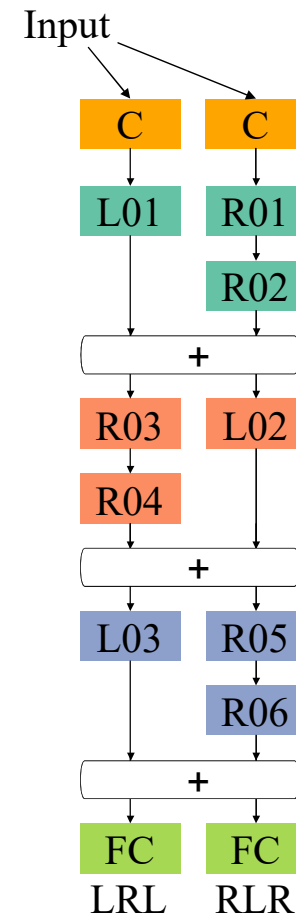
A deeply-fused net can be formed from **many different** base networks



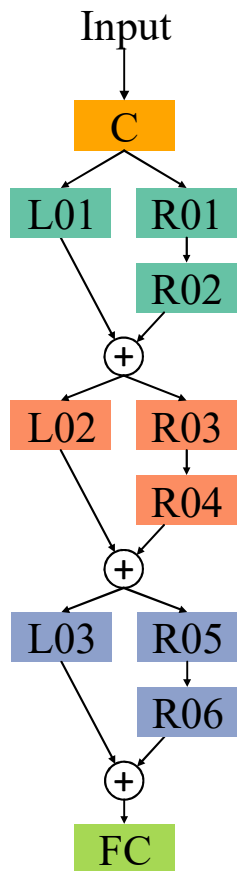
A deeply-fused net can be formed from **many different** base networks



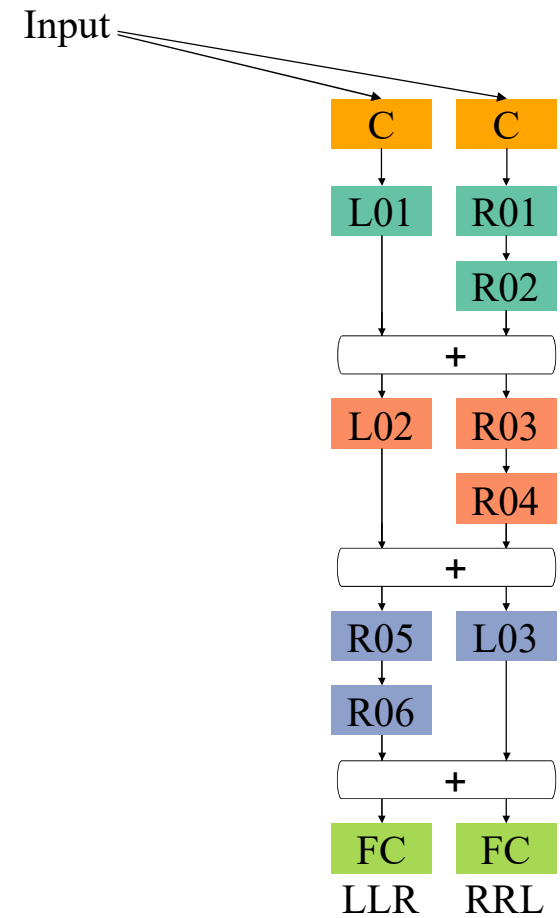
=



A deeply-fused net can be formed from **many different** base networks

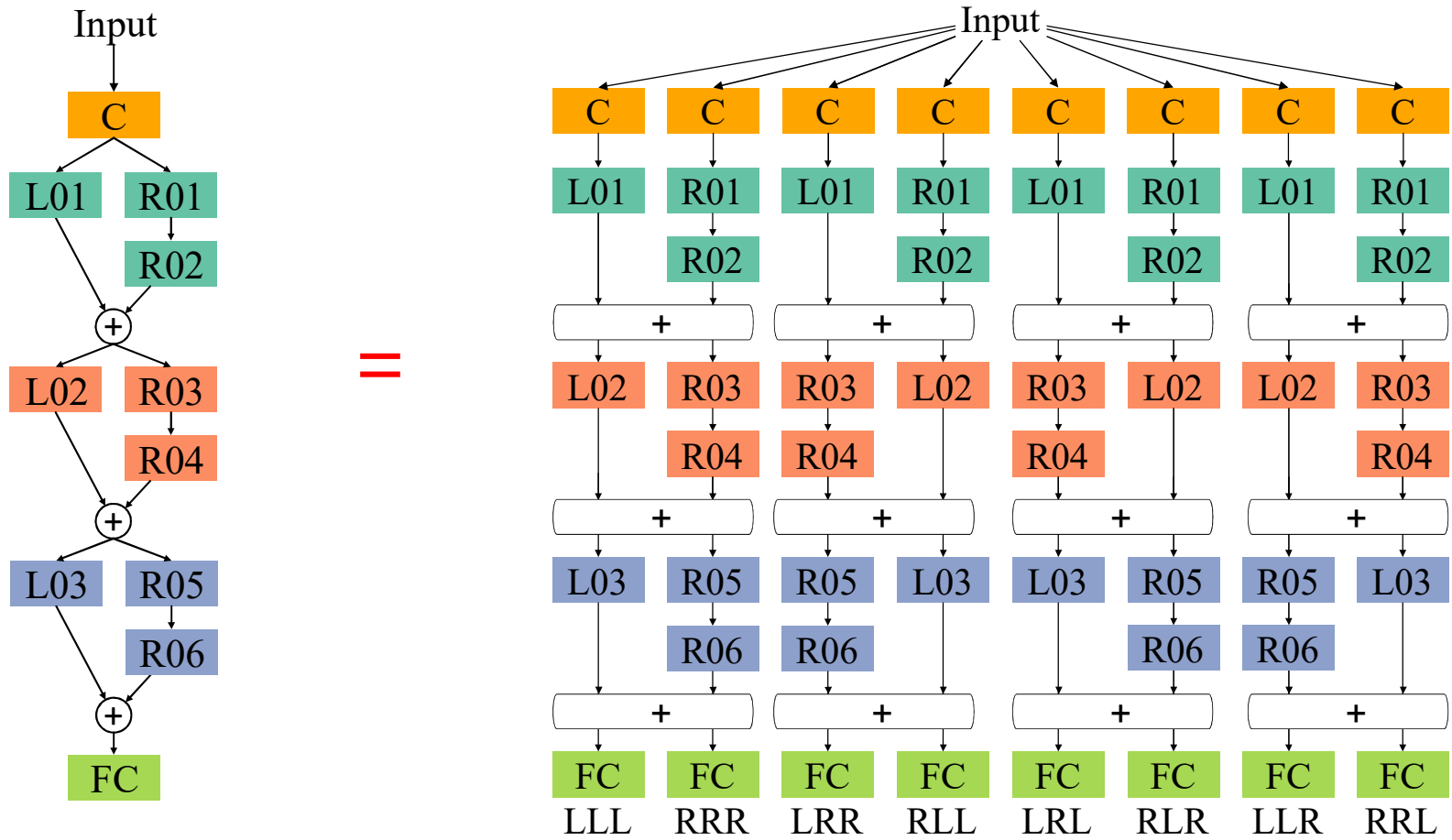


=

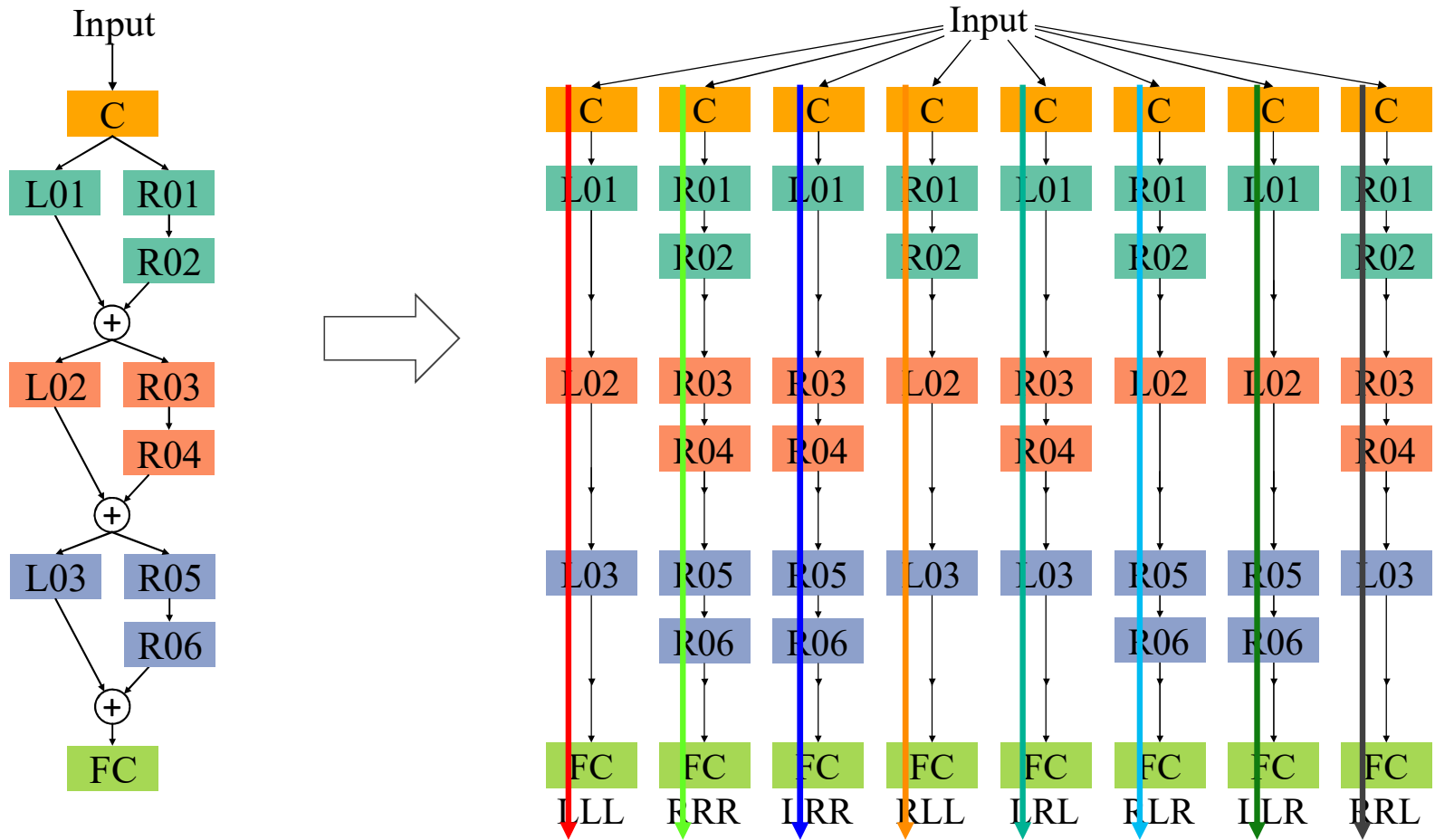


A deeply-fused net is a **multi-path multi-scale** network

A deeply-fused net is a **multi-path multi-scale** network

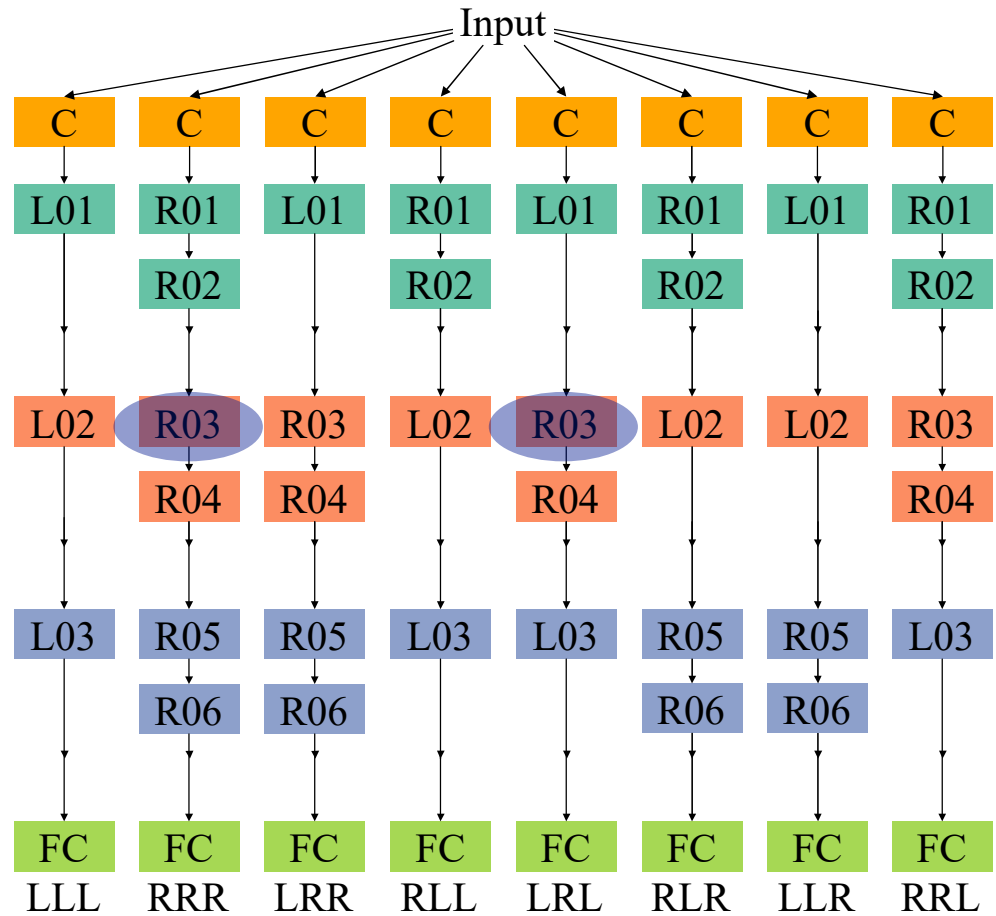
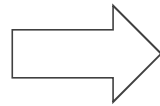
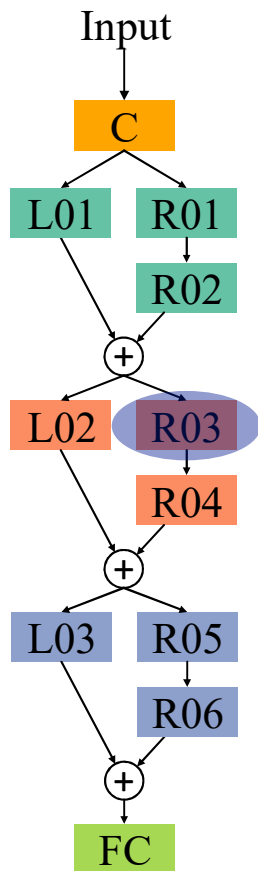


A deeply-fused net is a **multi-path multi-scale** network

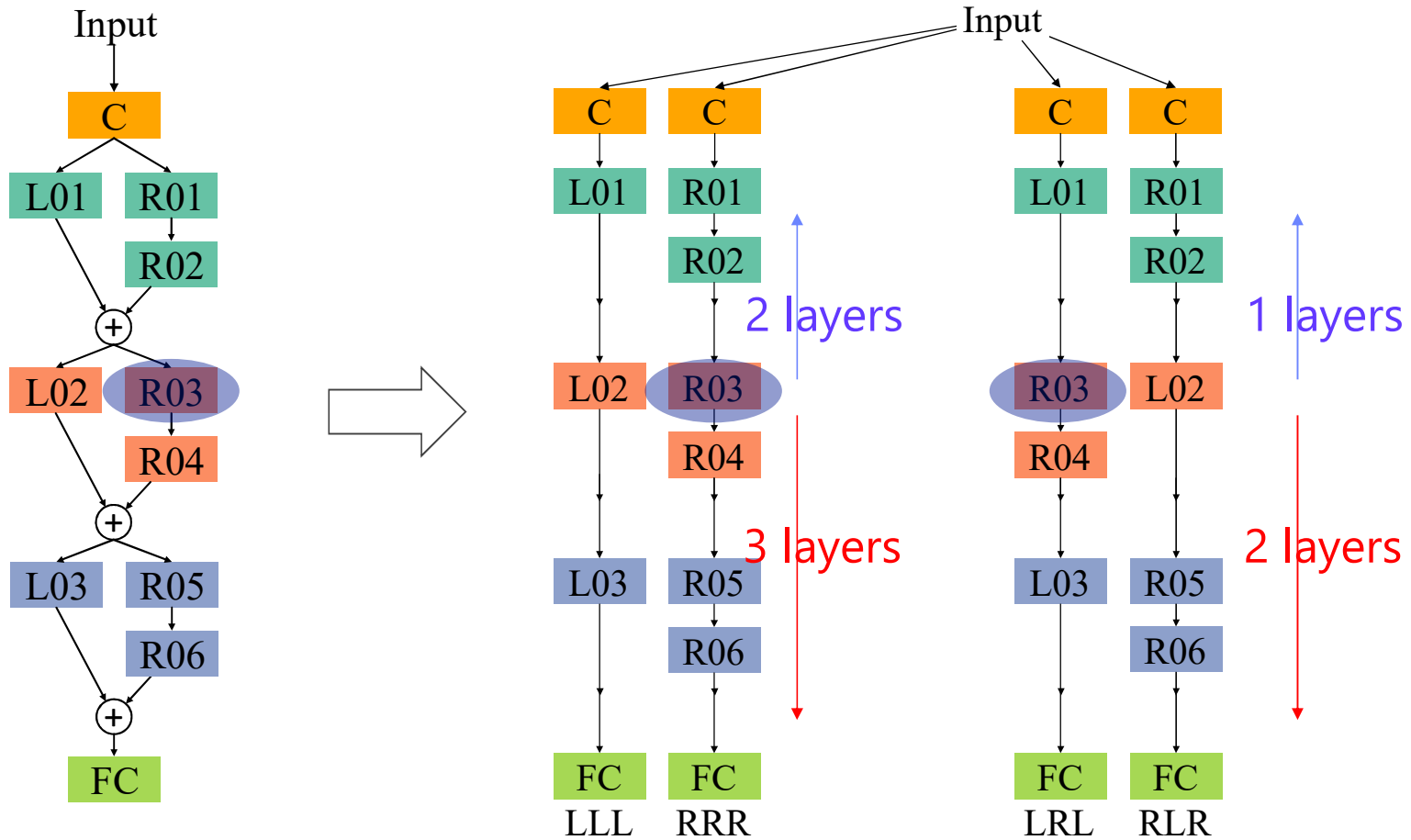


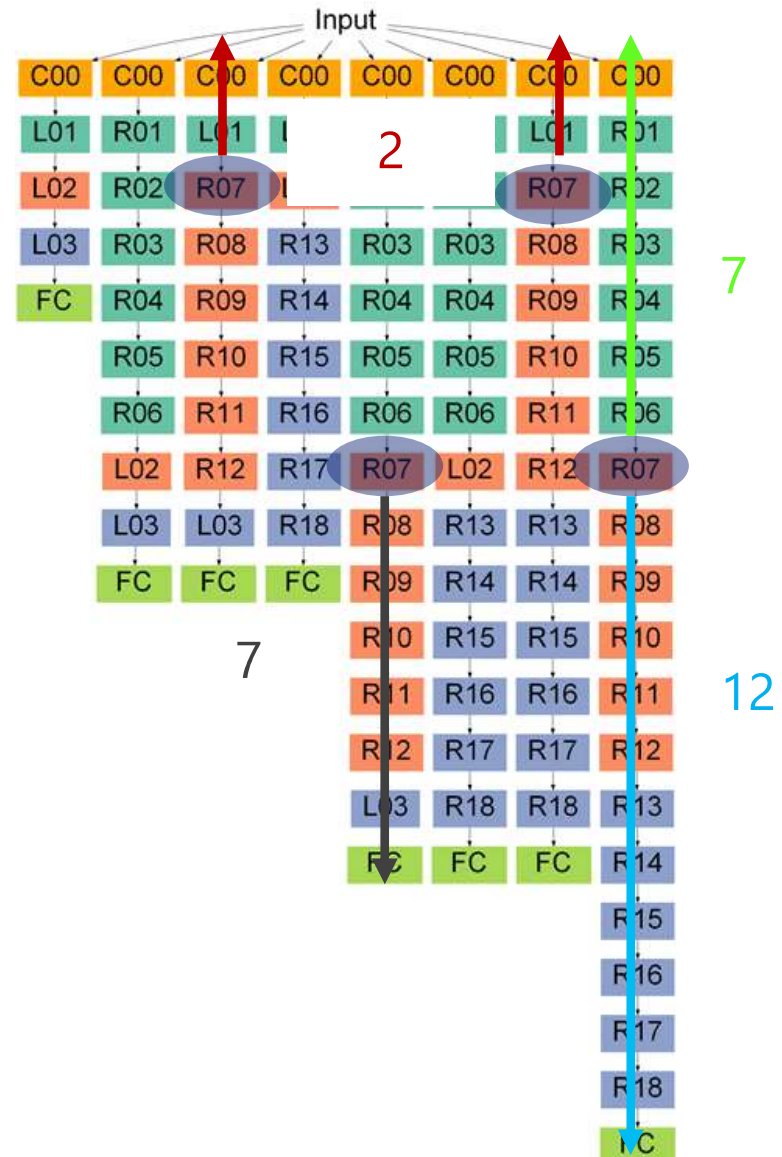
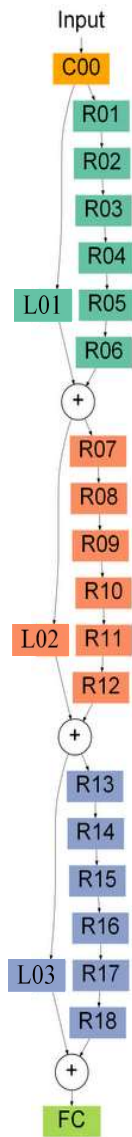
Each layer has an **express way** to input and output

Each layer has an **express way** to input and output

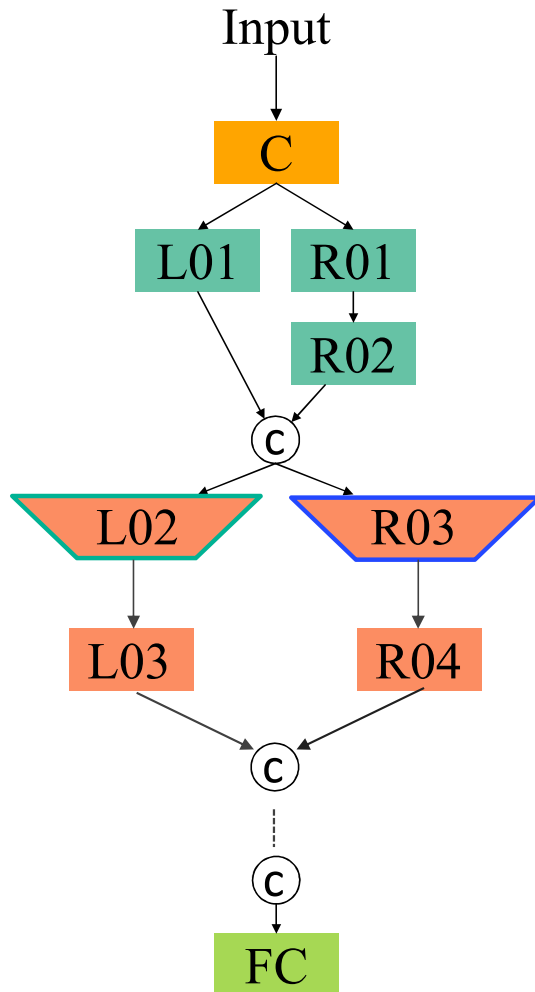


Each layer has an **express way** to input and output

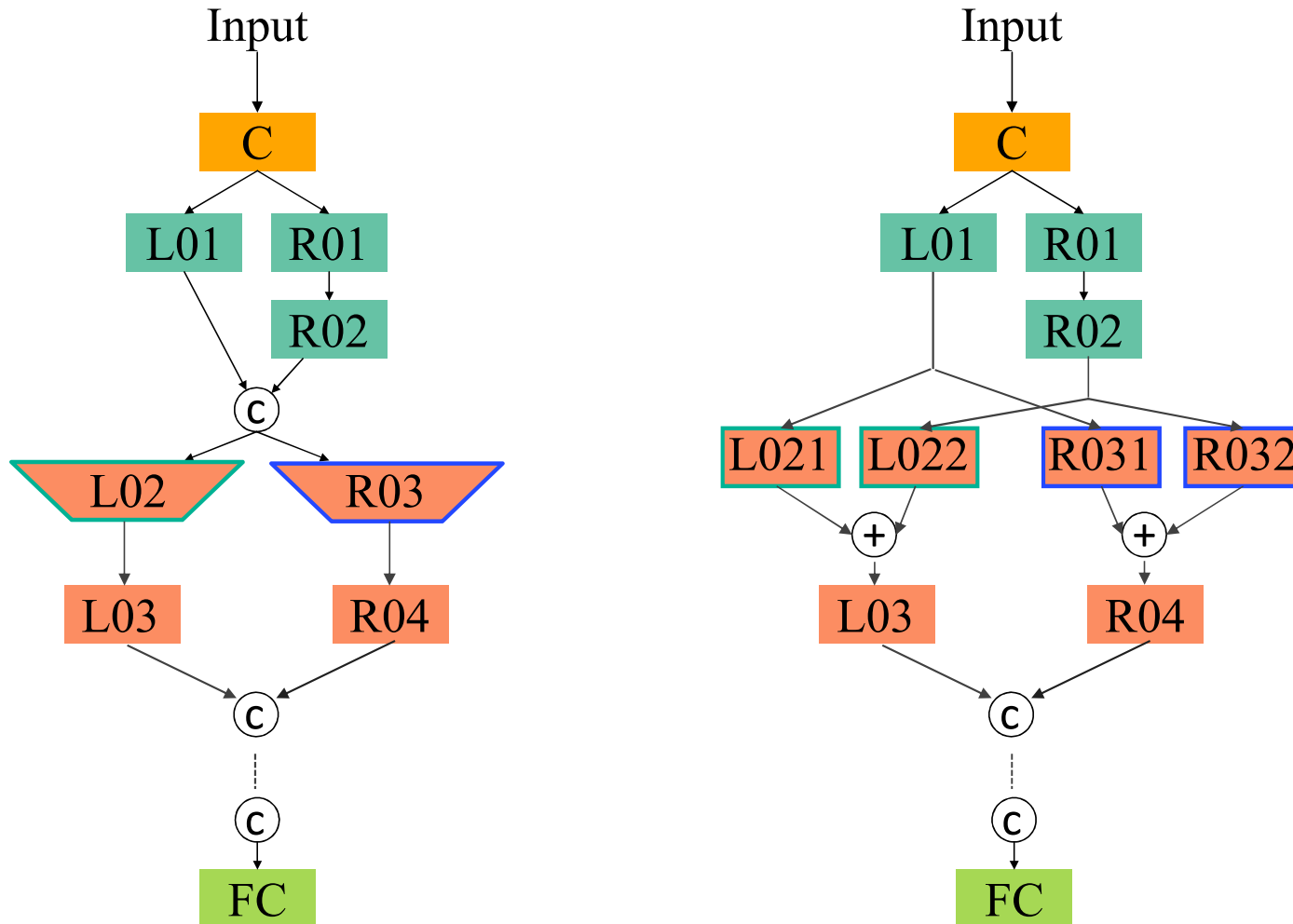




Concatenation in GoogLeNet inception is a sum fusion



Concatenation in GoogLeNet inception is a sum fusion

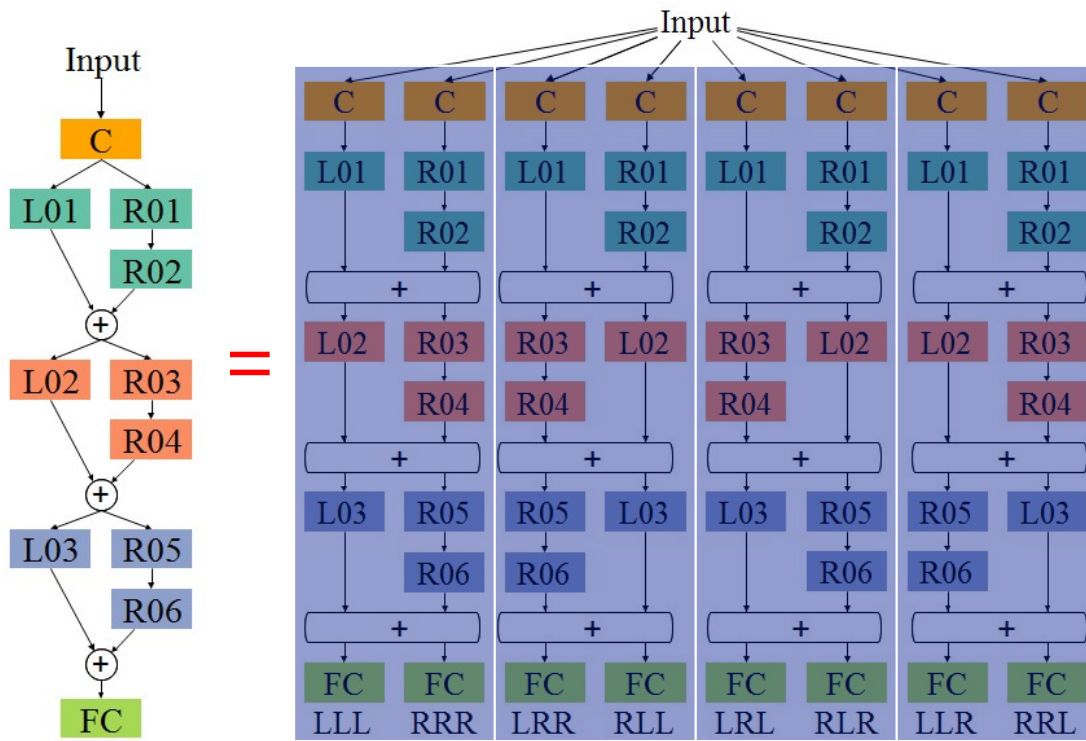


Deep fusion behaves like ensembling

[Liming Zhao](#), Jingdong Wang, [Xi Li](#), [Zhuowen Tu](#), [Wenjun Zeng](#): On the Connection of Deep Fusion to Ensembling. [CoRR abs/1611.07718](#) (2016)

The **architecture** of a deeply-fused net **resembles** an ensemble

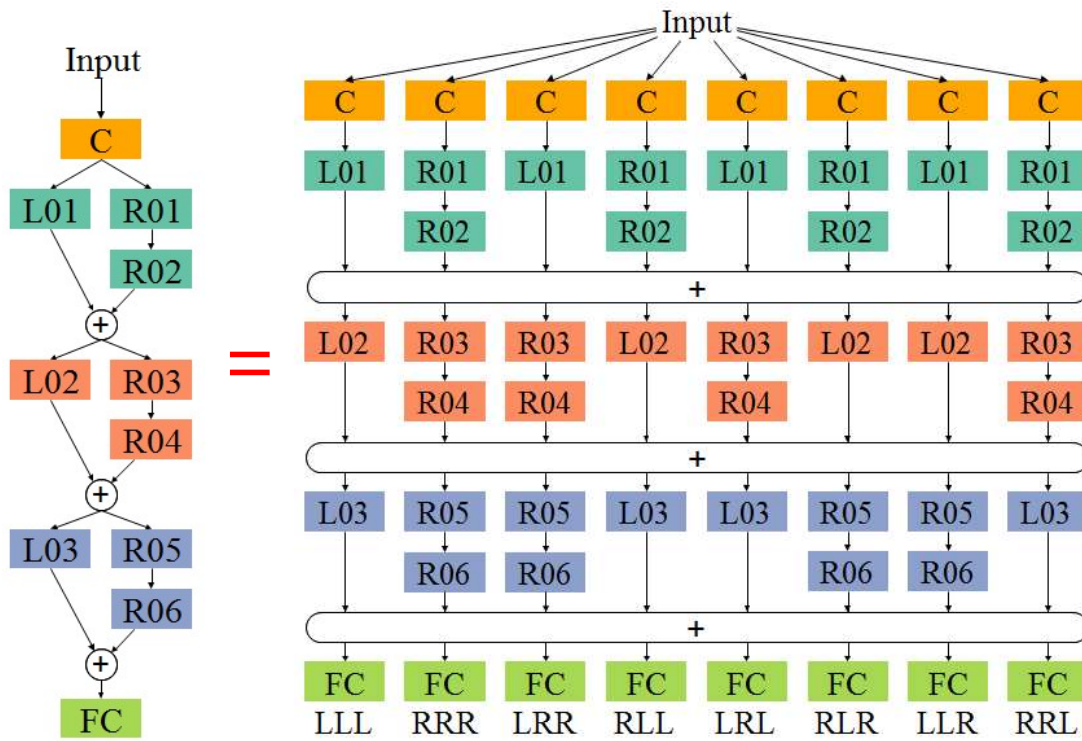
The **architecture** of a deeply-fused net **resembles** an ensemble



Deep fusion

Expanded view

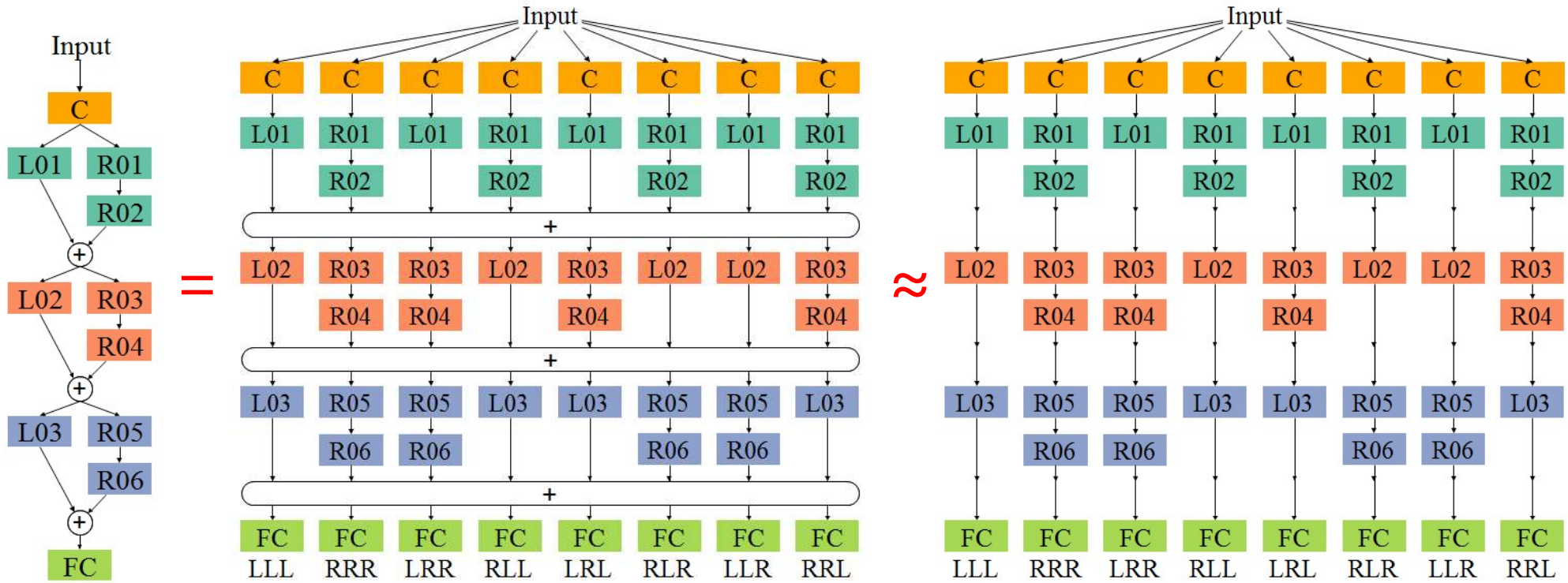
The **architecture** of a deeply-fused net **resembles** an ensemble



Deep fusion

Expanded view

The **architecture** of a deeply-fused net **resembles** an ensemble

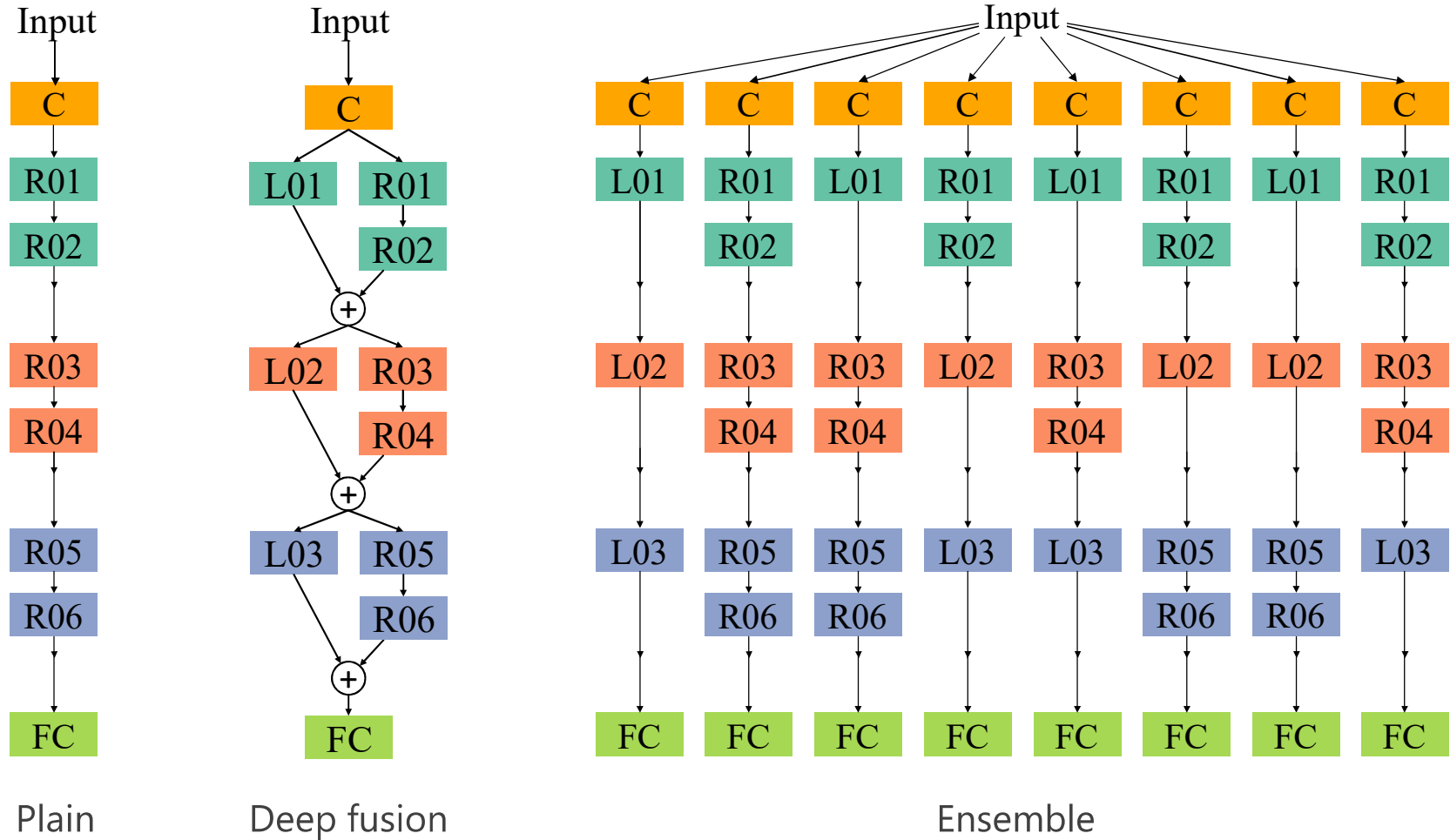


Deep fusion

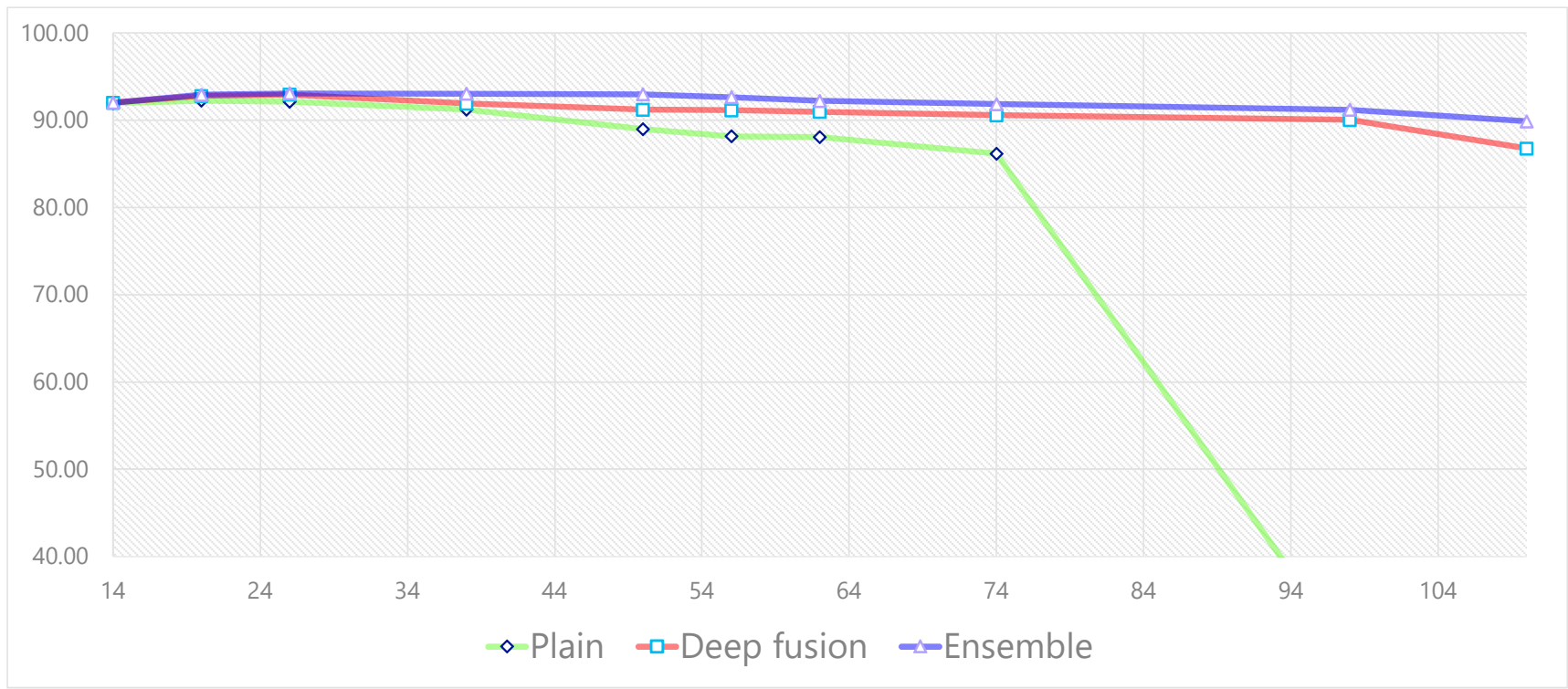
Expanded view

An ensemble

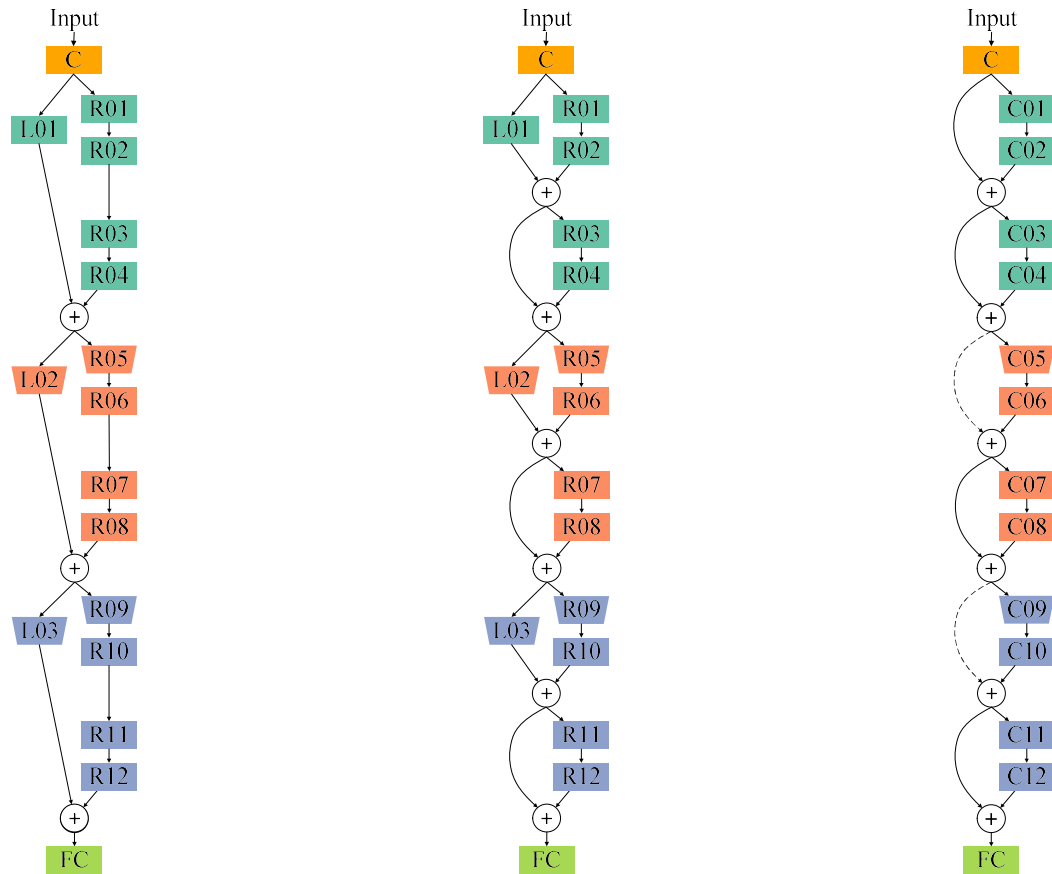
The **performances** of deeply-fused nets **resemble** ensembles



The **performances** of deeply-fused nets **resemble** ensembles



Larger ensemble sizes leads to better performances

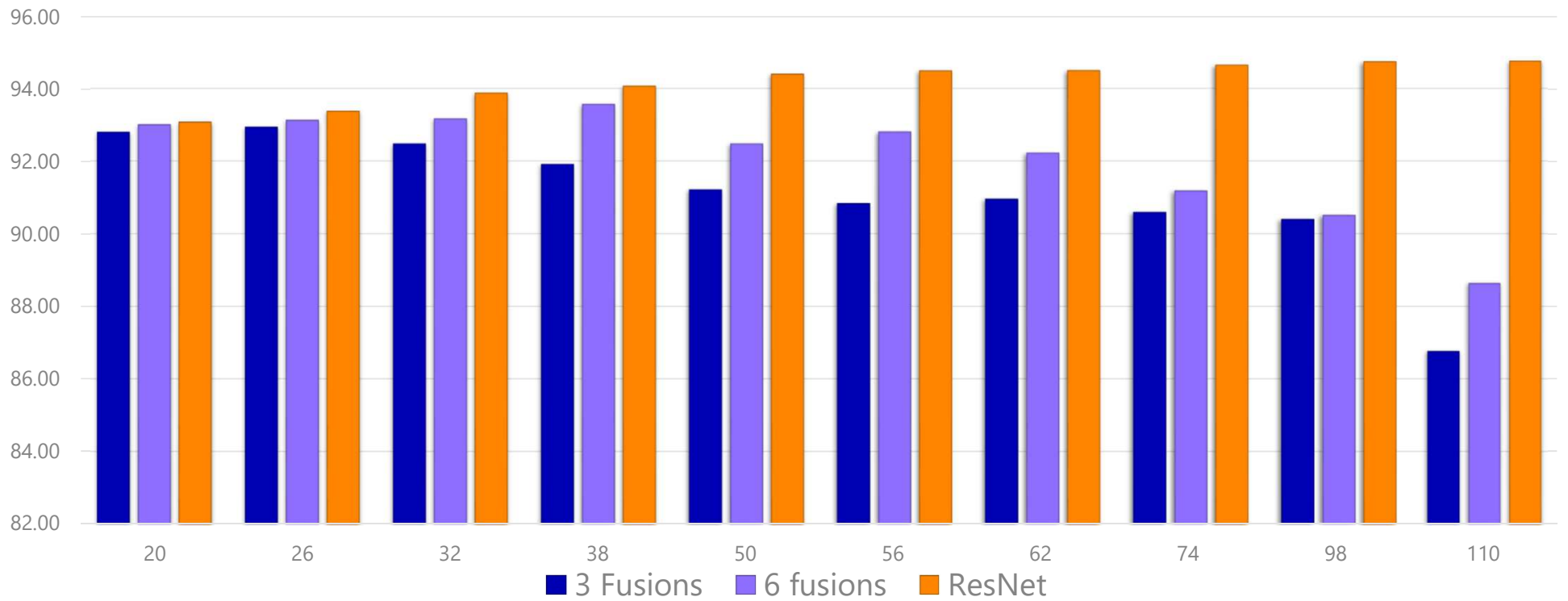


3 fusions

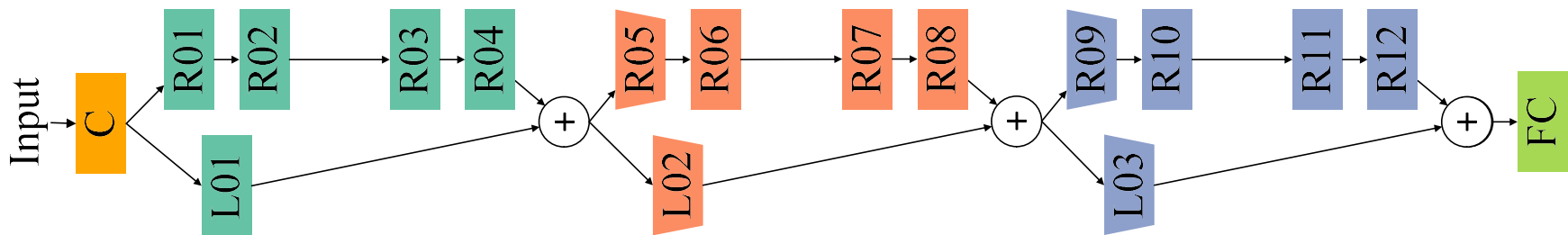
6 fusions

ResNet

Larger ensemble sizes leads to better performances

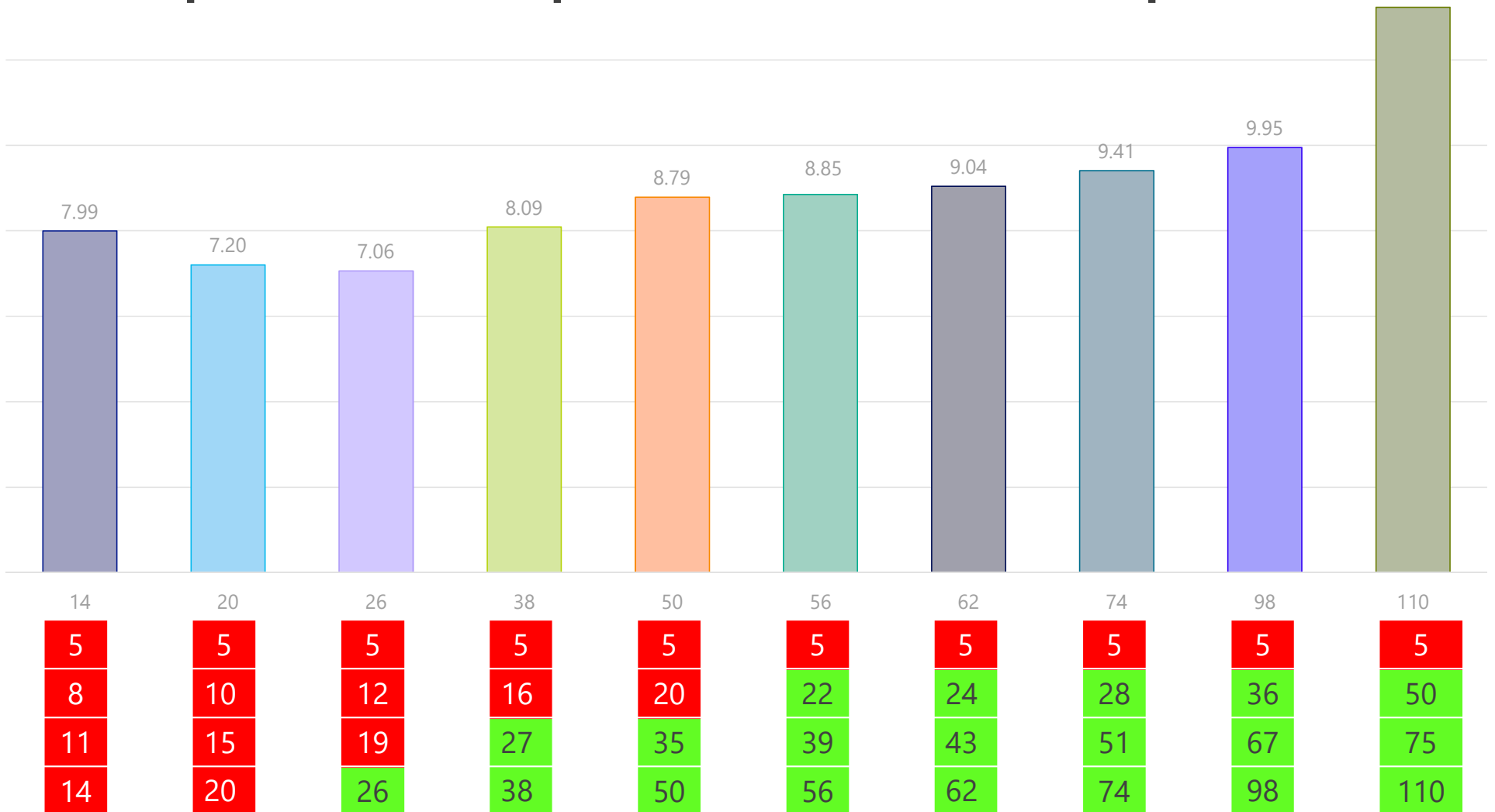


The **capabilities of component networks** affect the **performances**



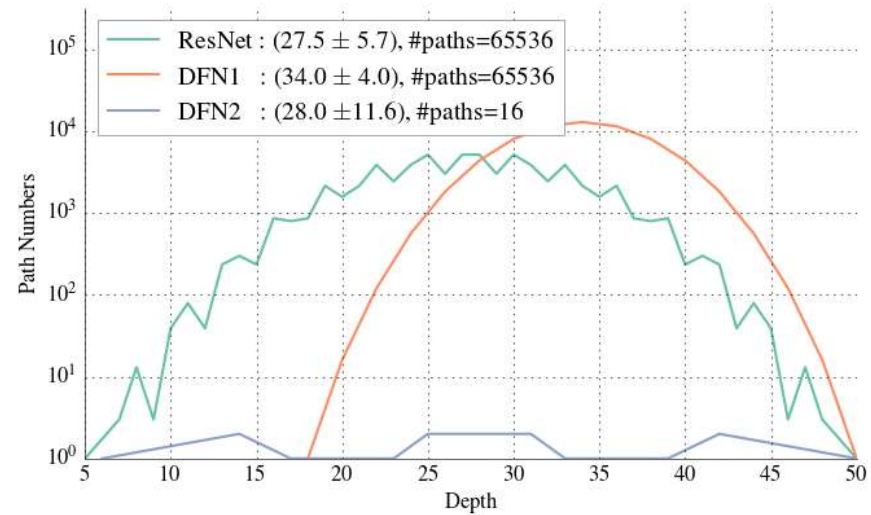
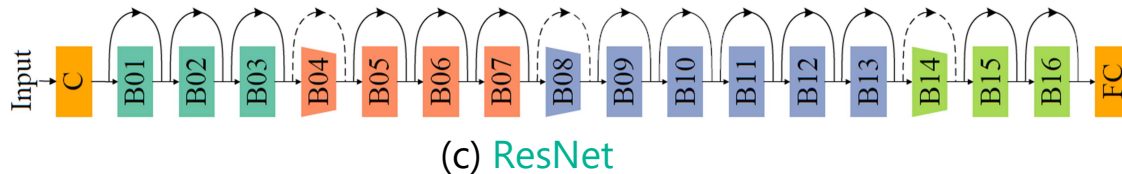
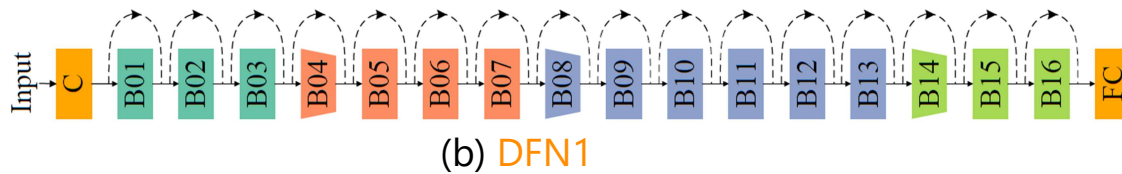
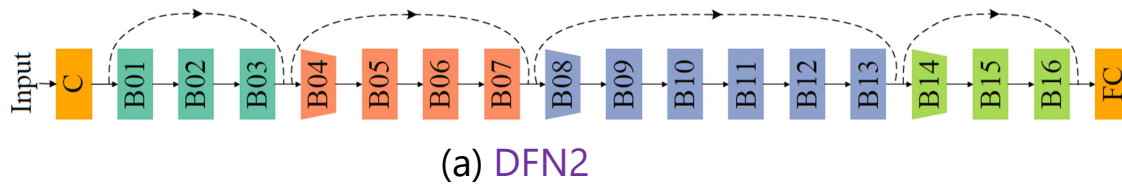
Deeply-fused nets: Different depths, 3 fusions

The capabilities of component networks affect the performances



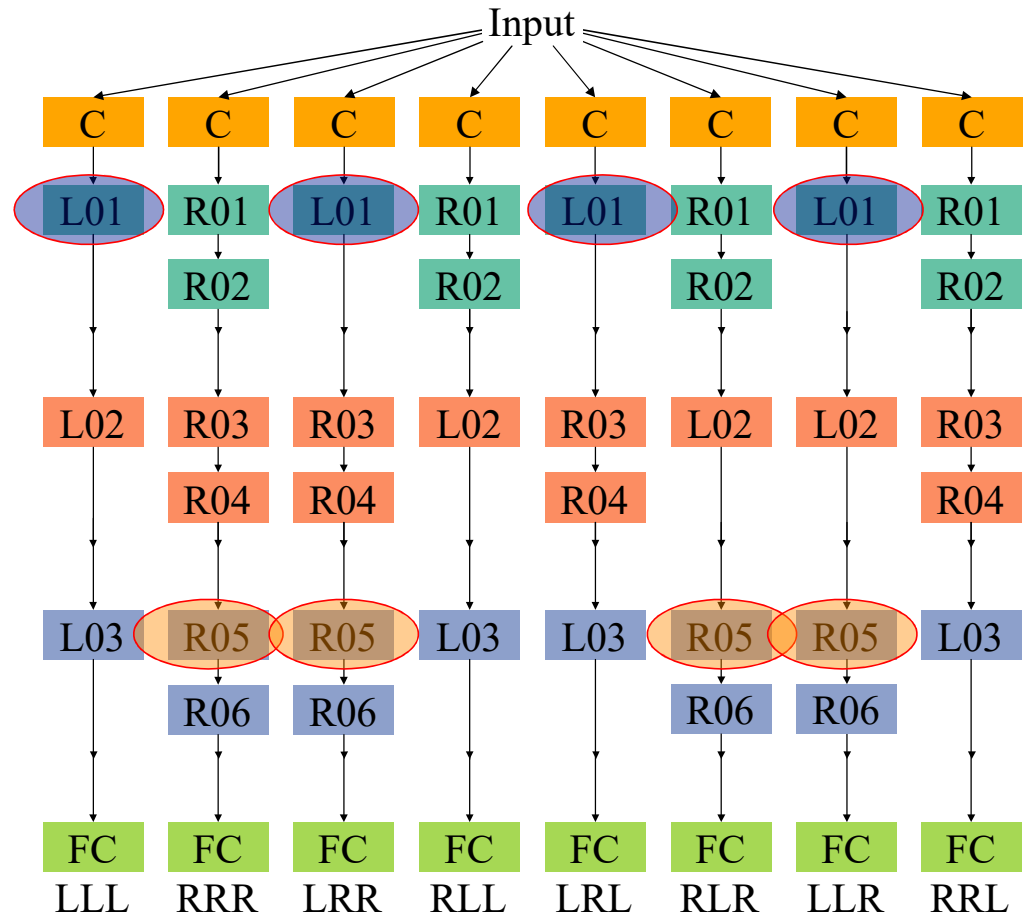
Empirical study on ImageNet

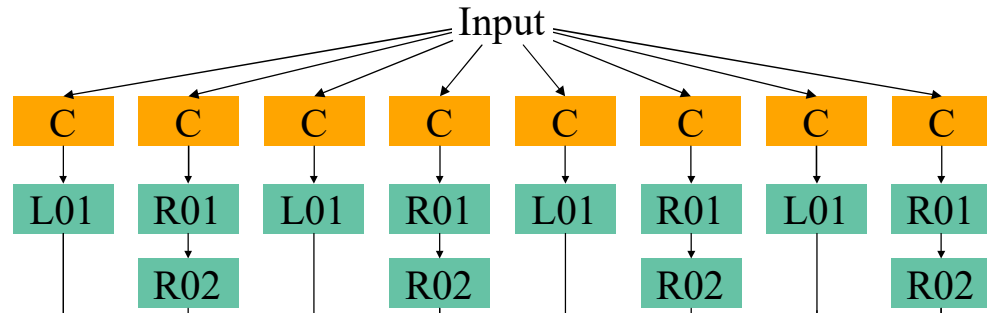
	ResNet	DFN1	DFN2
Top1 validation error	24.94	25.10	27.17
Top5 validation error	7.46	7.85	8.98



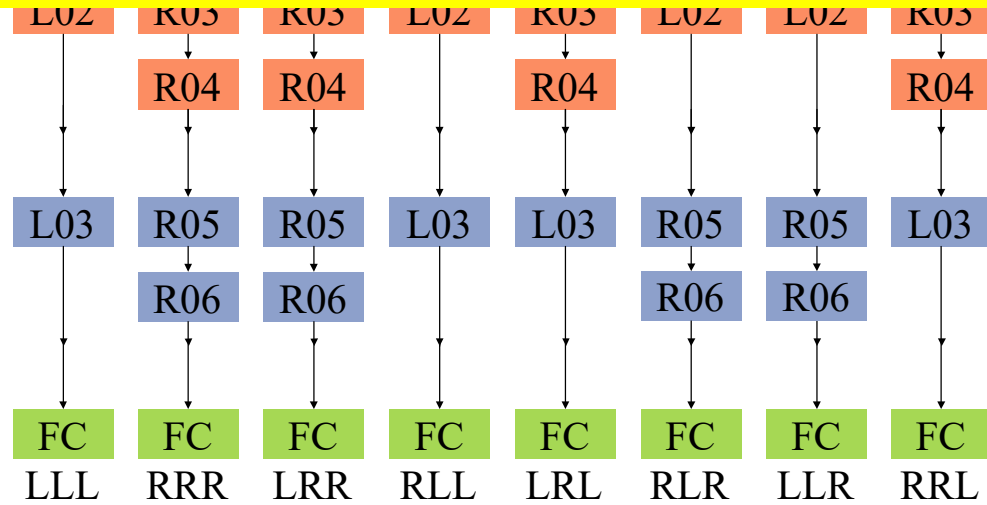
What the depth brings?

[Liming Zhao](#), Jingdong Wang, [Xi Li](#), [Zhuowen Tu](#), [Wenjun Zeng](#): On the Connection of Deep Fusion to Ensembling. [CoRR abs/1611.07718](#) (2016)

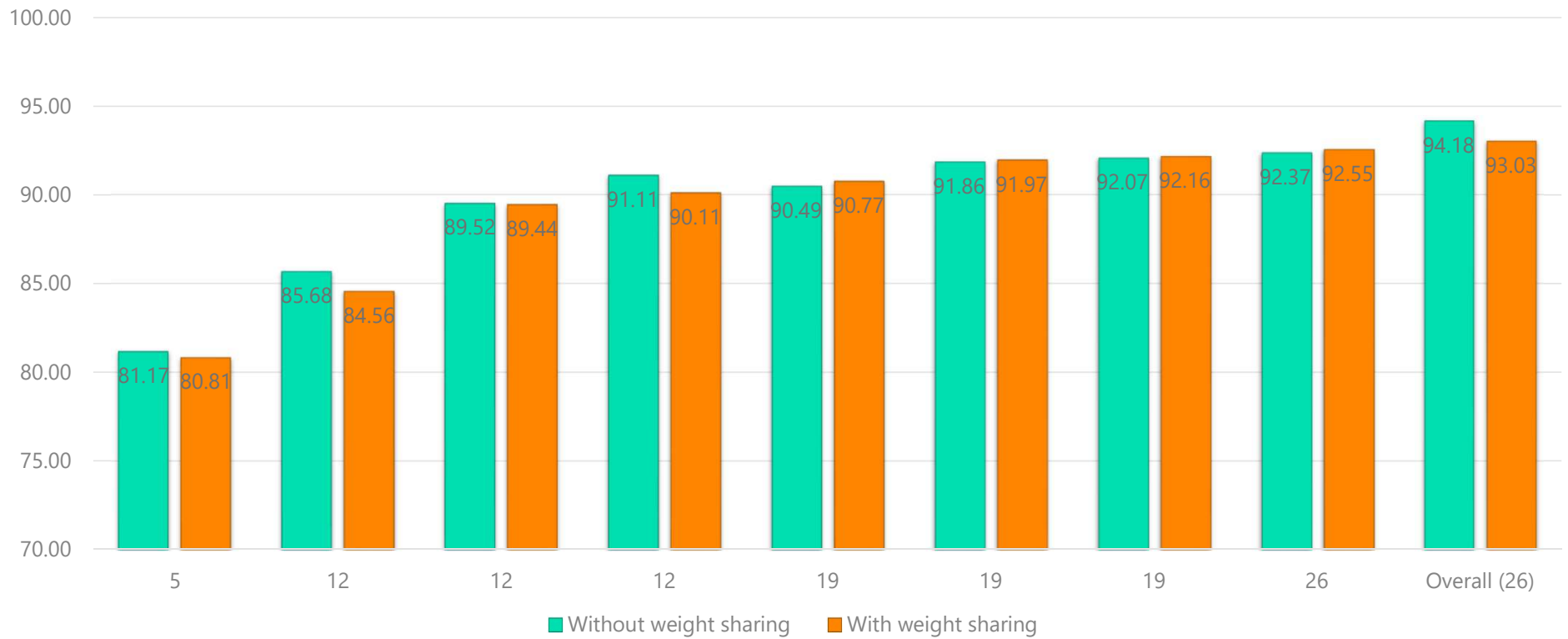




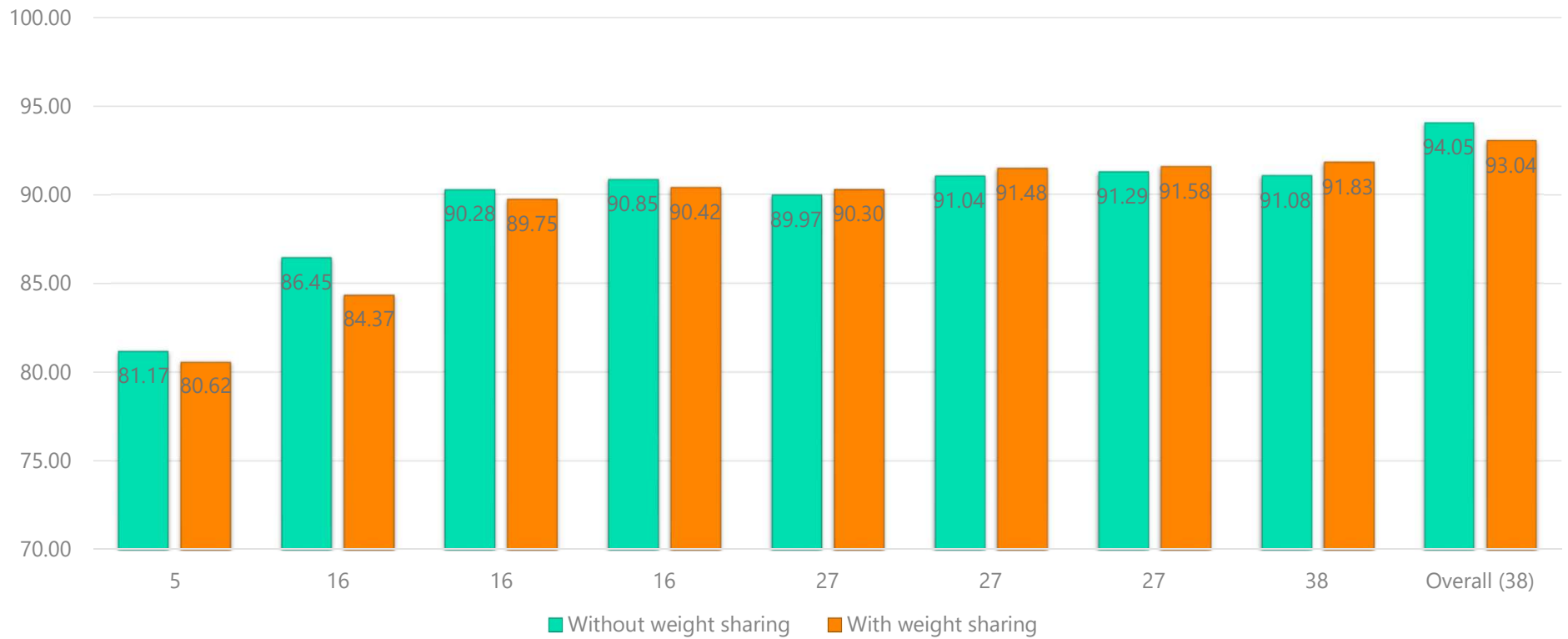
How sharing weight affect the performance of each component?



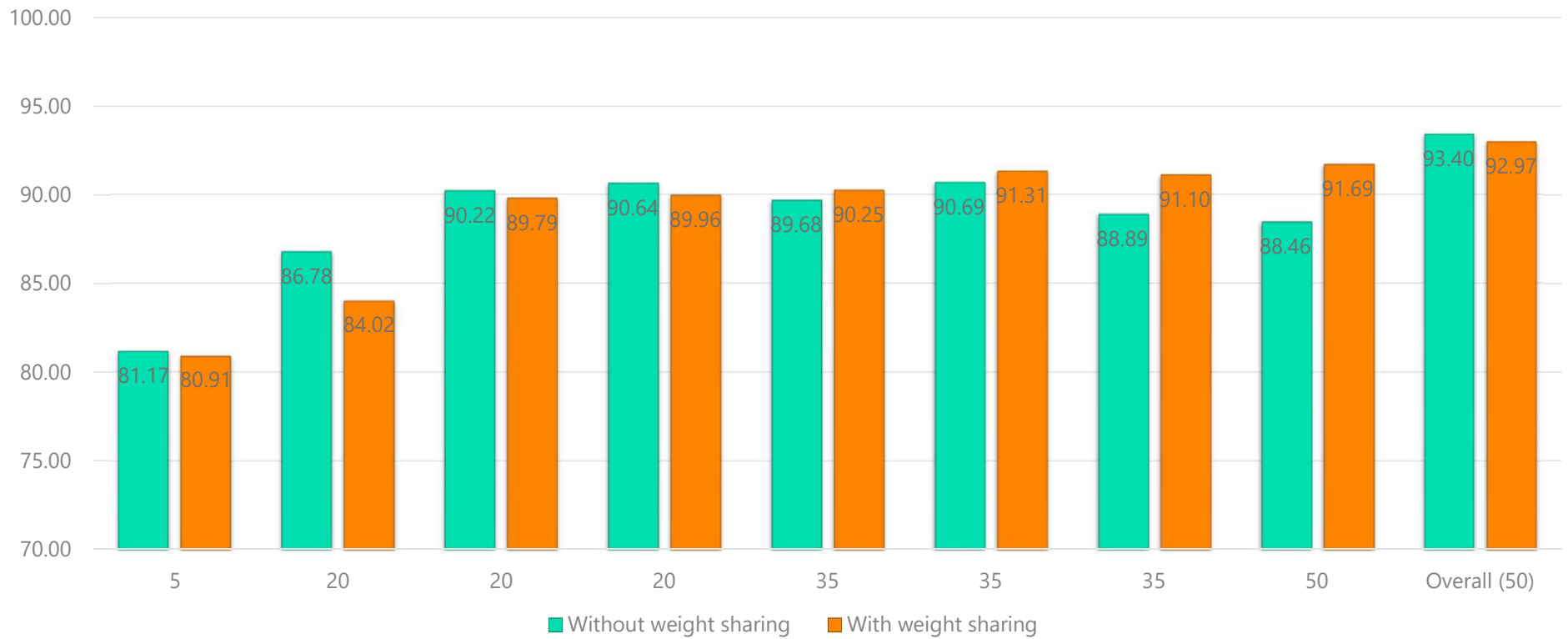
How sharing weight affect the performance of each component?



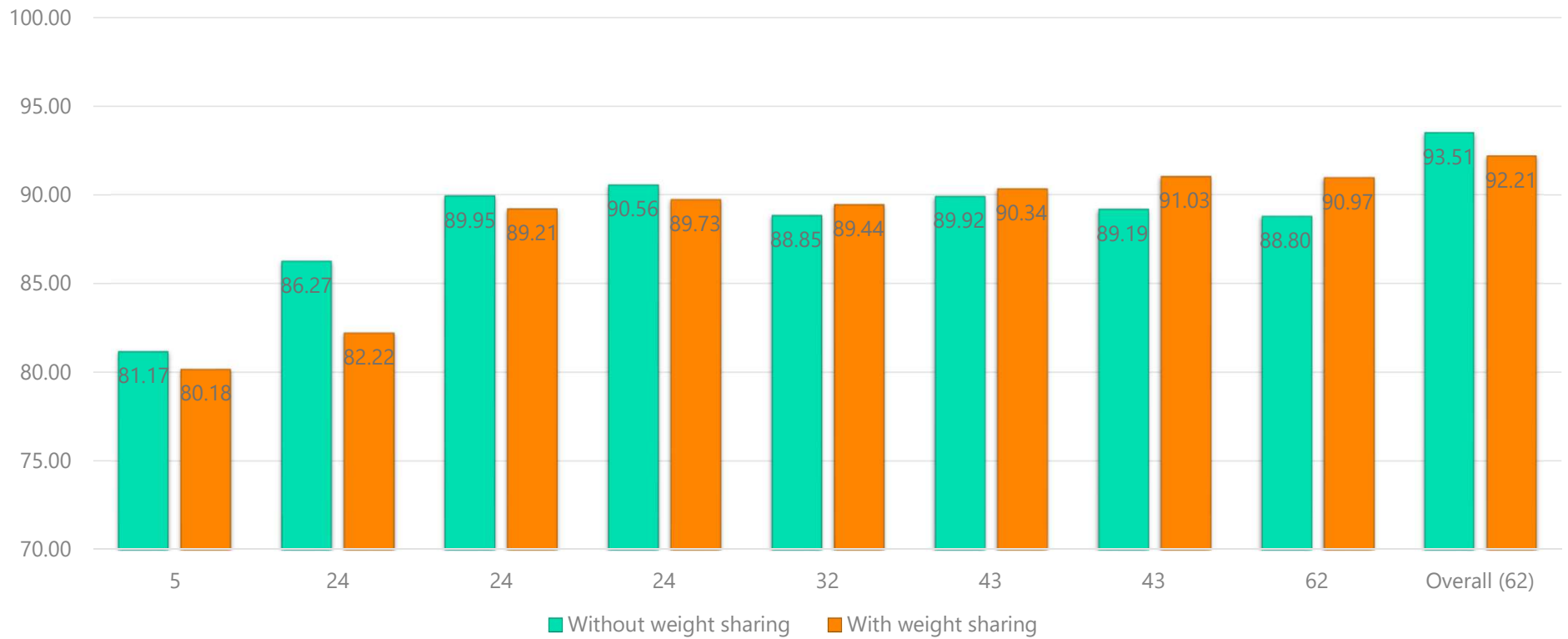
How sharing weight affect the performance of each component?



How sharing weight affect the performance of each component?



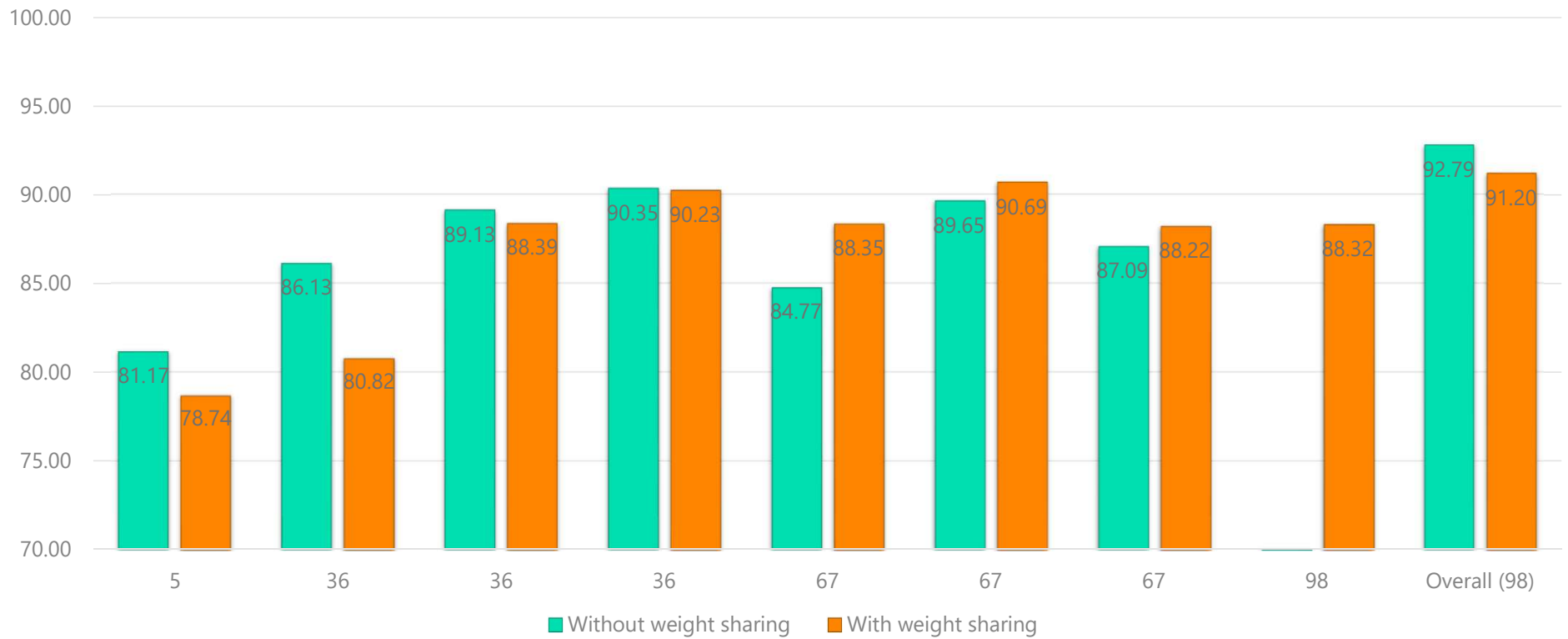
How sharing weight affect the performance of each component?



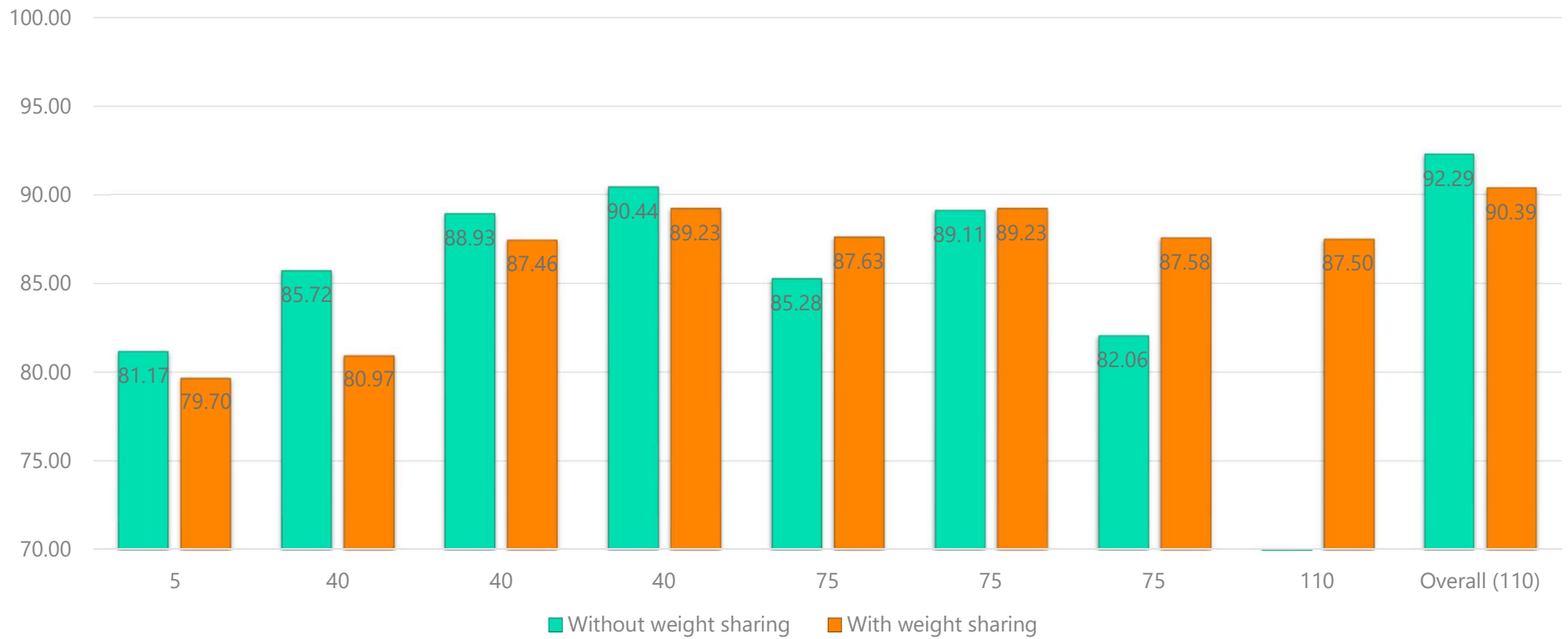
How sharing weight affect the performance of each component?

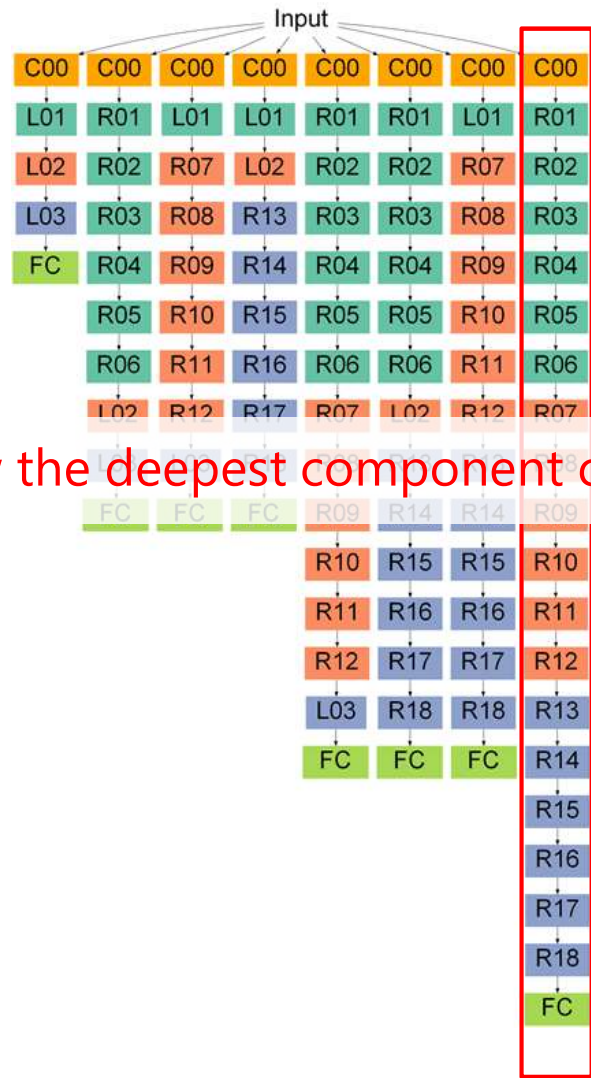


How sharing weight affect the performance of each component?

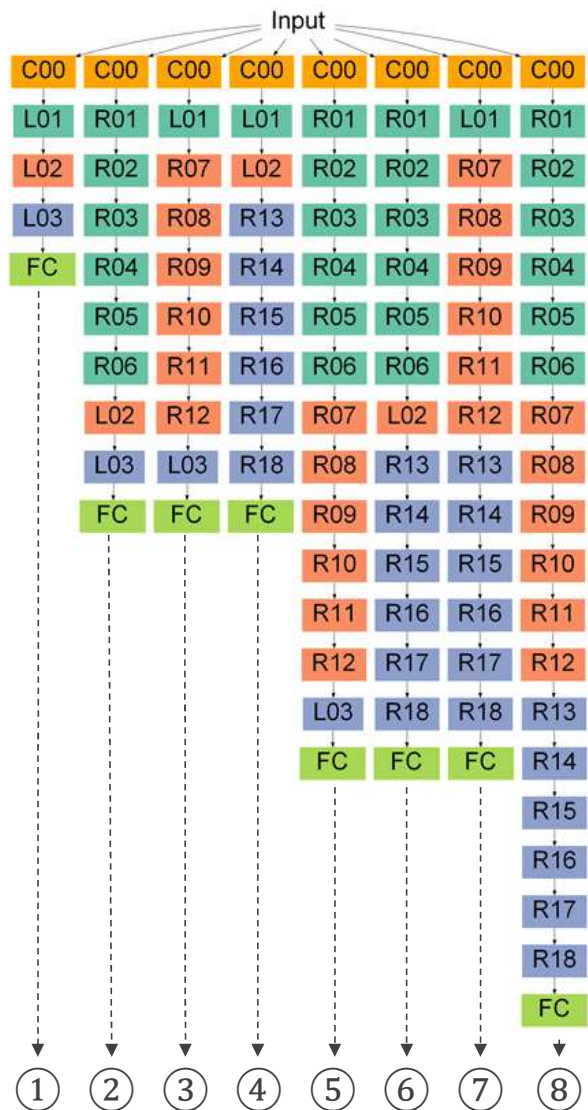


How sharing weight affect the performance of each component?

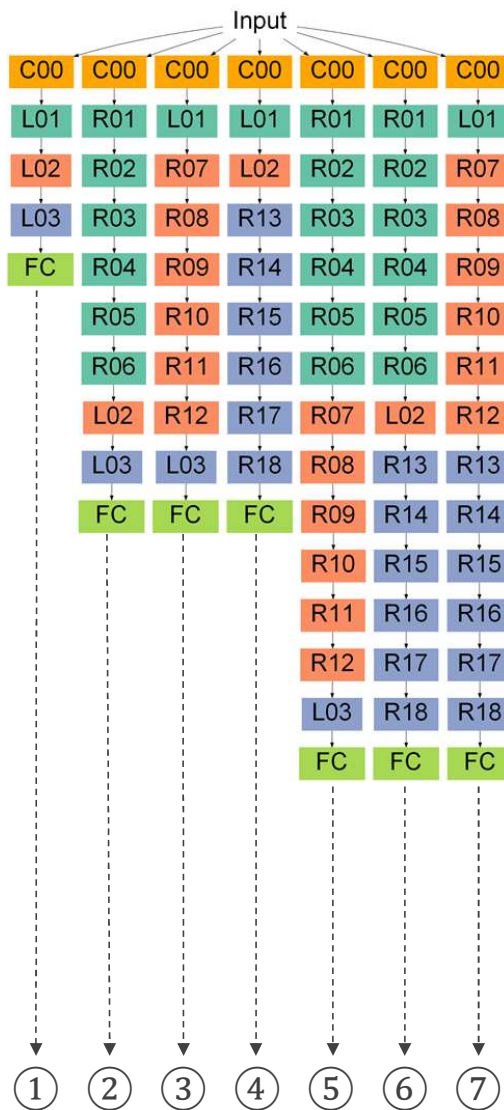




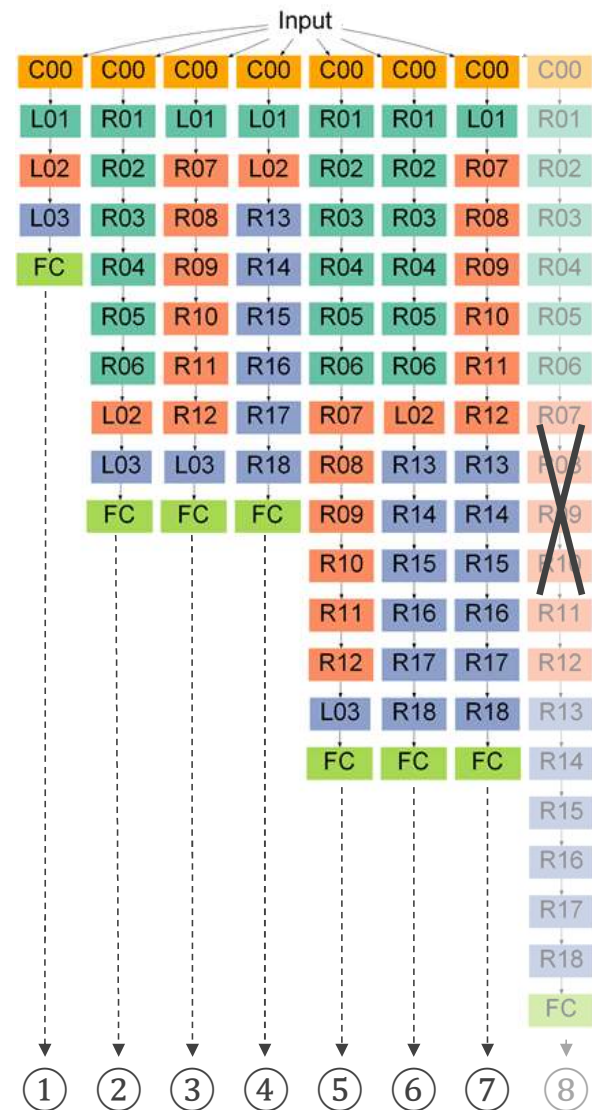
How the deepest component contributes?



Ensemble-8

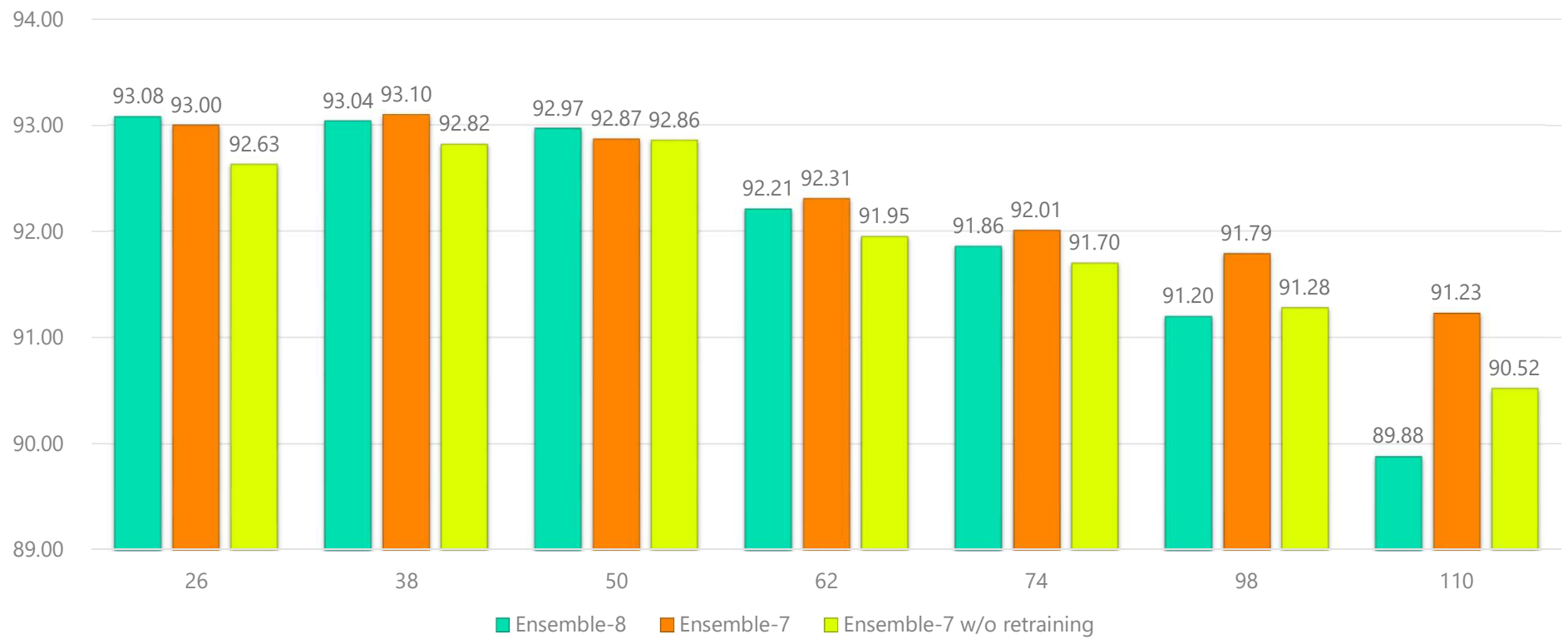


Ensemble-7

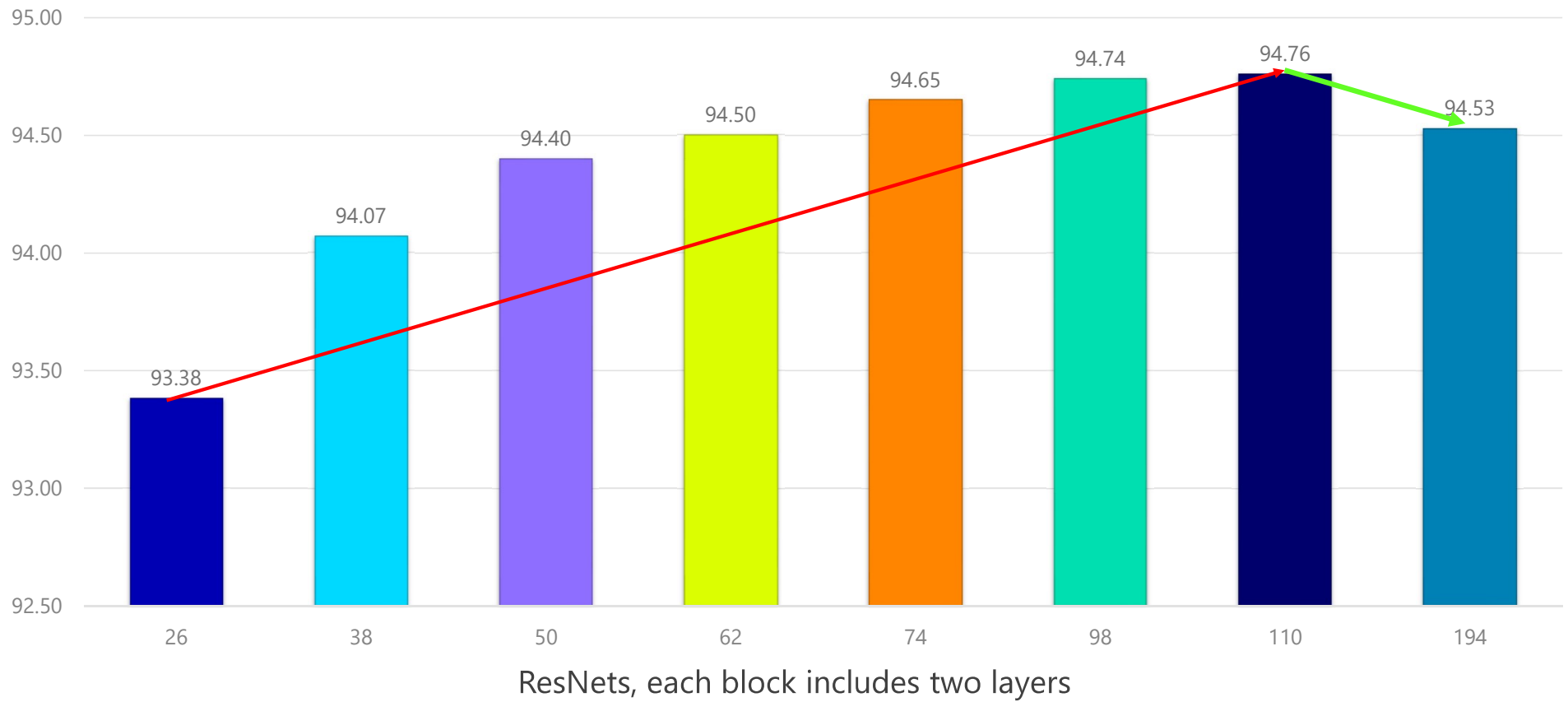


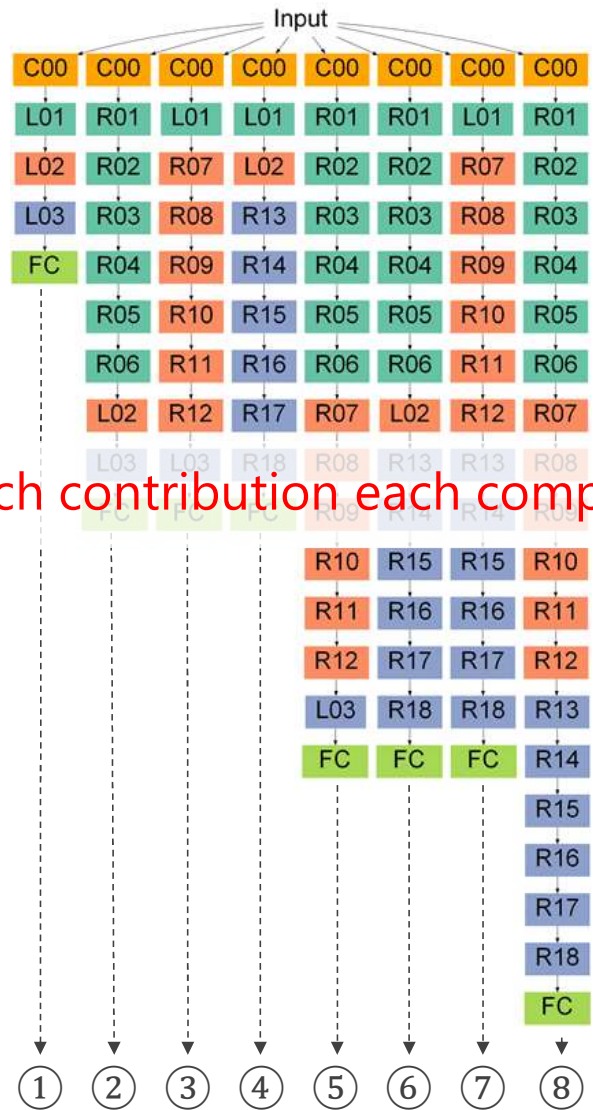
Ensemble-7 w/o retraining

The deepest component **harms** the training

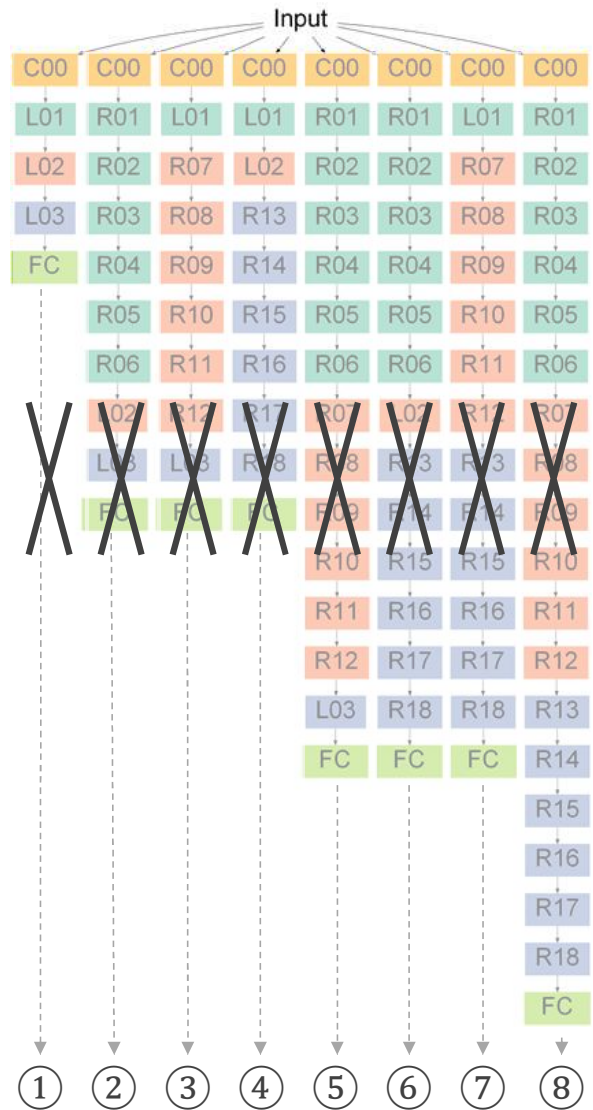


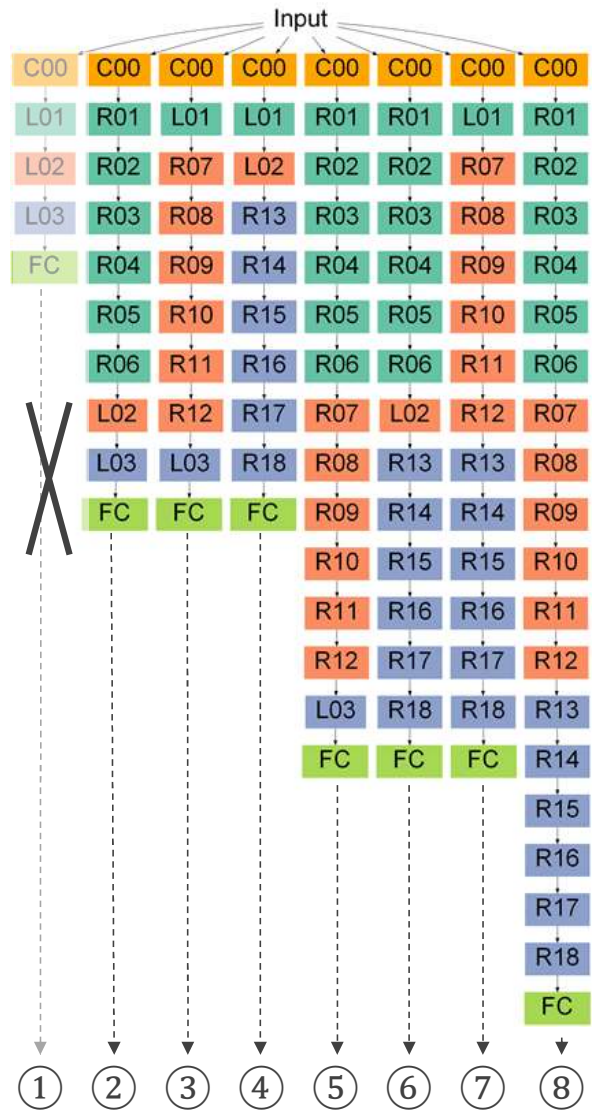
Greater depths leads to **larger** ensemble sizes, but **weaker** components

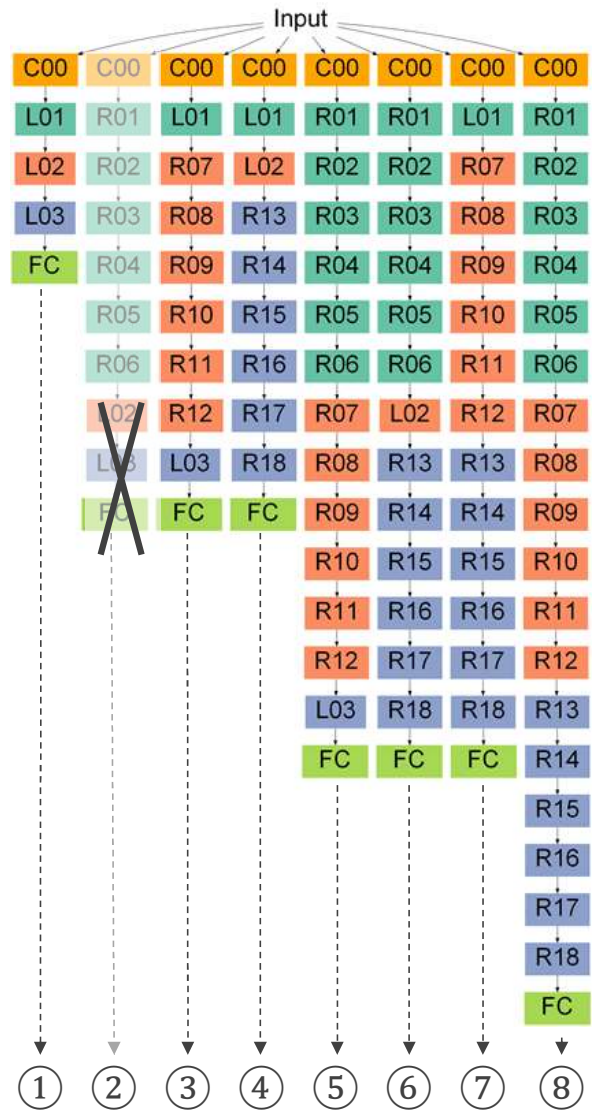


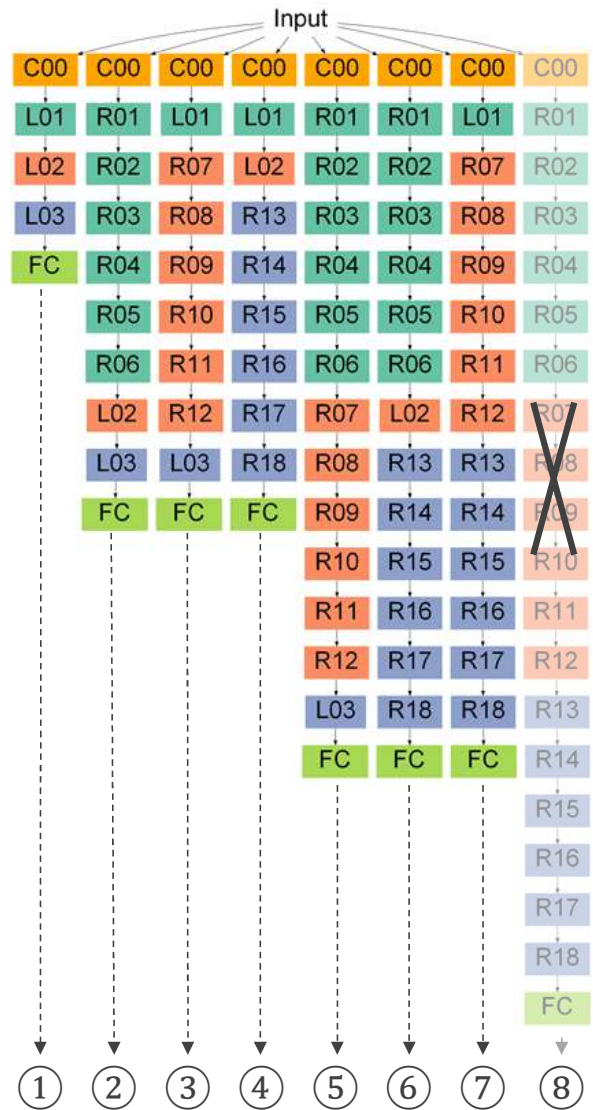


How much contribution each component makes?

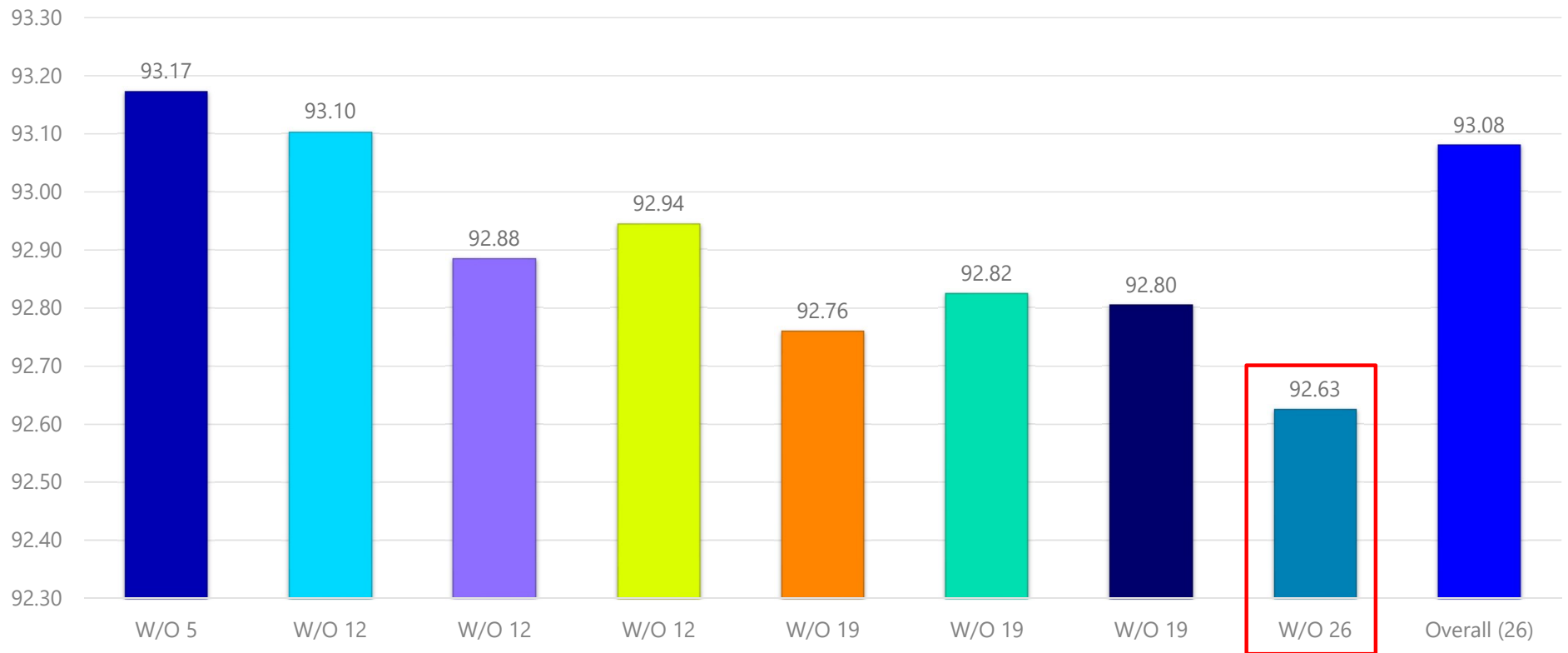




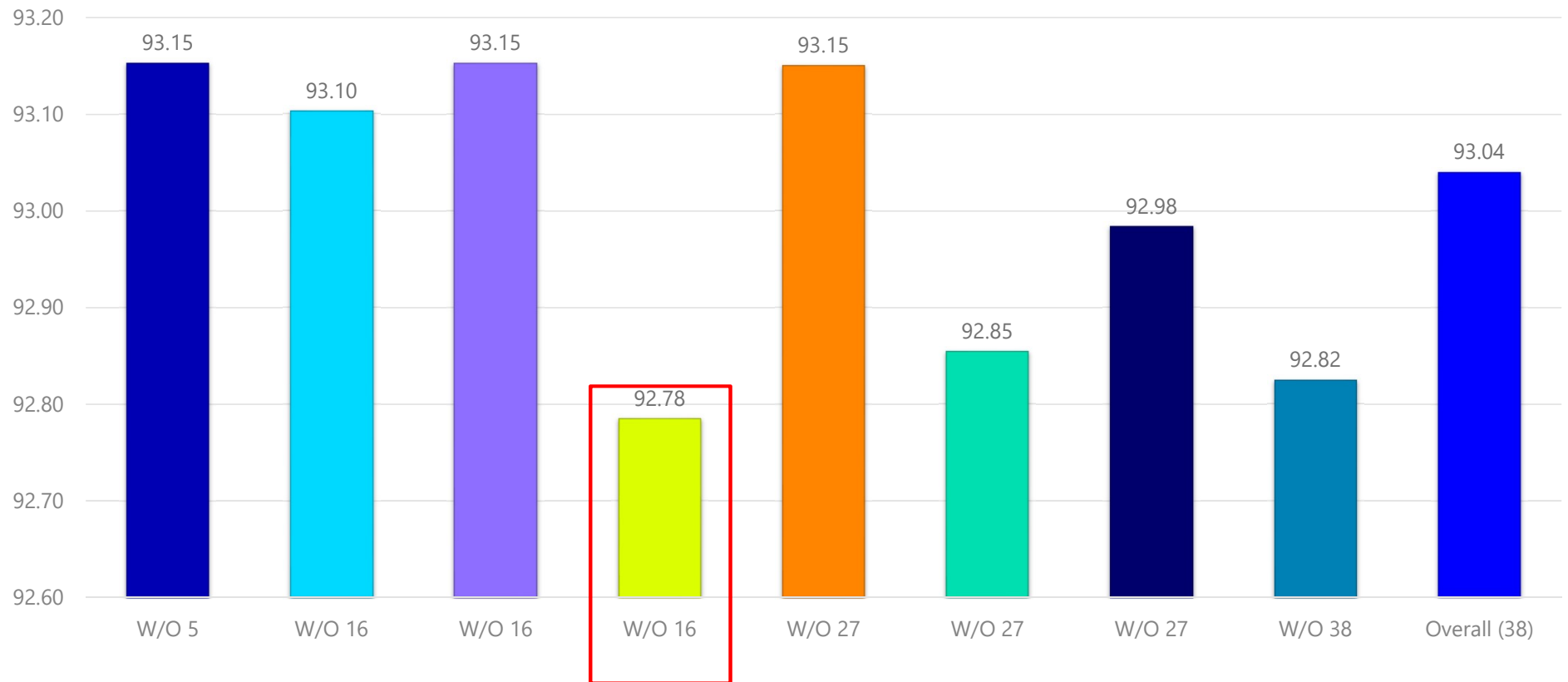




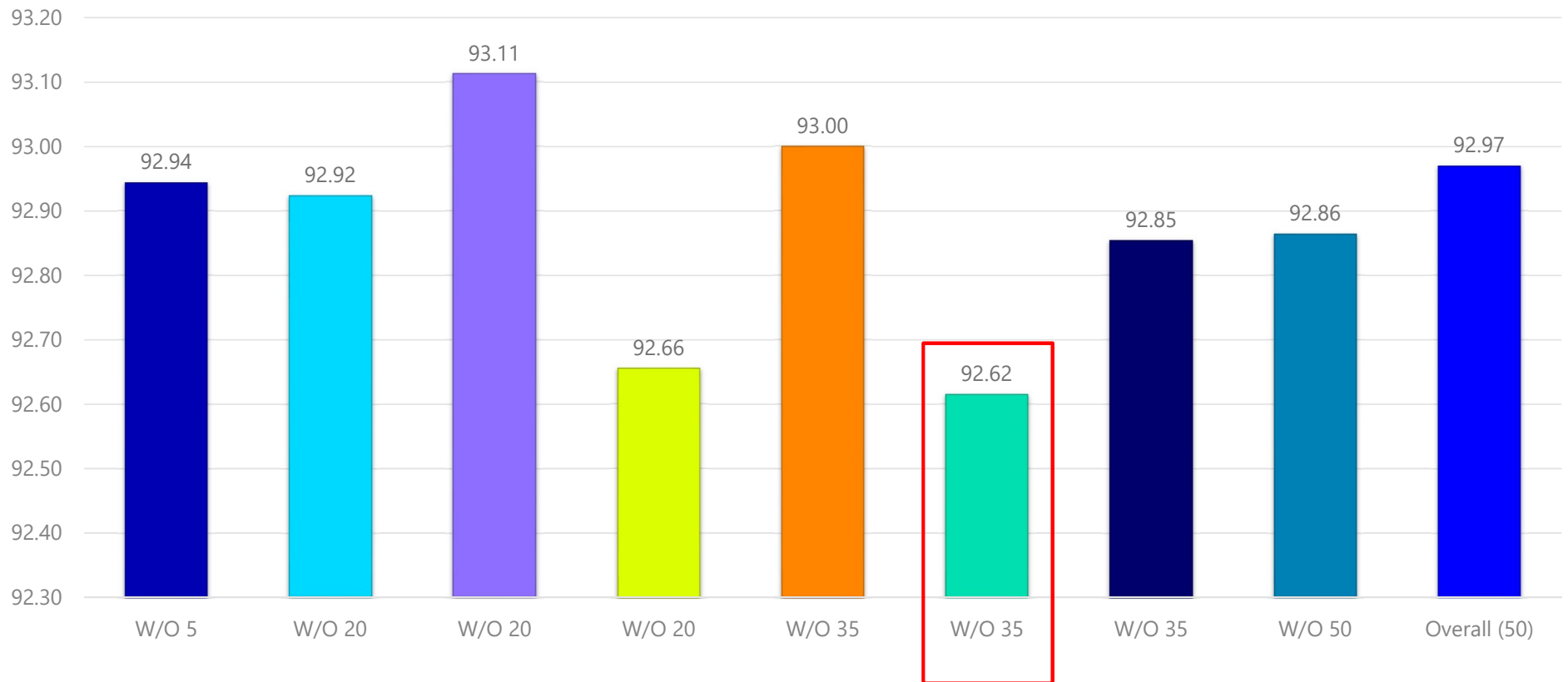
How much contribution each component makes?



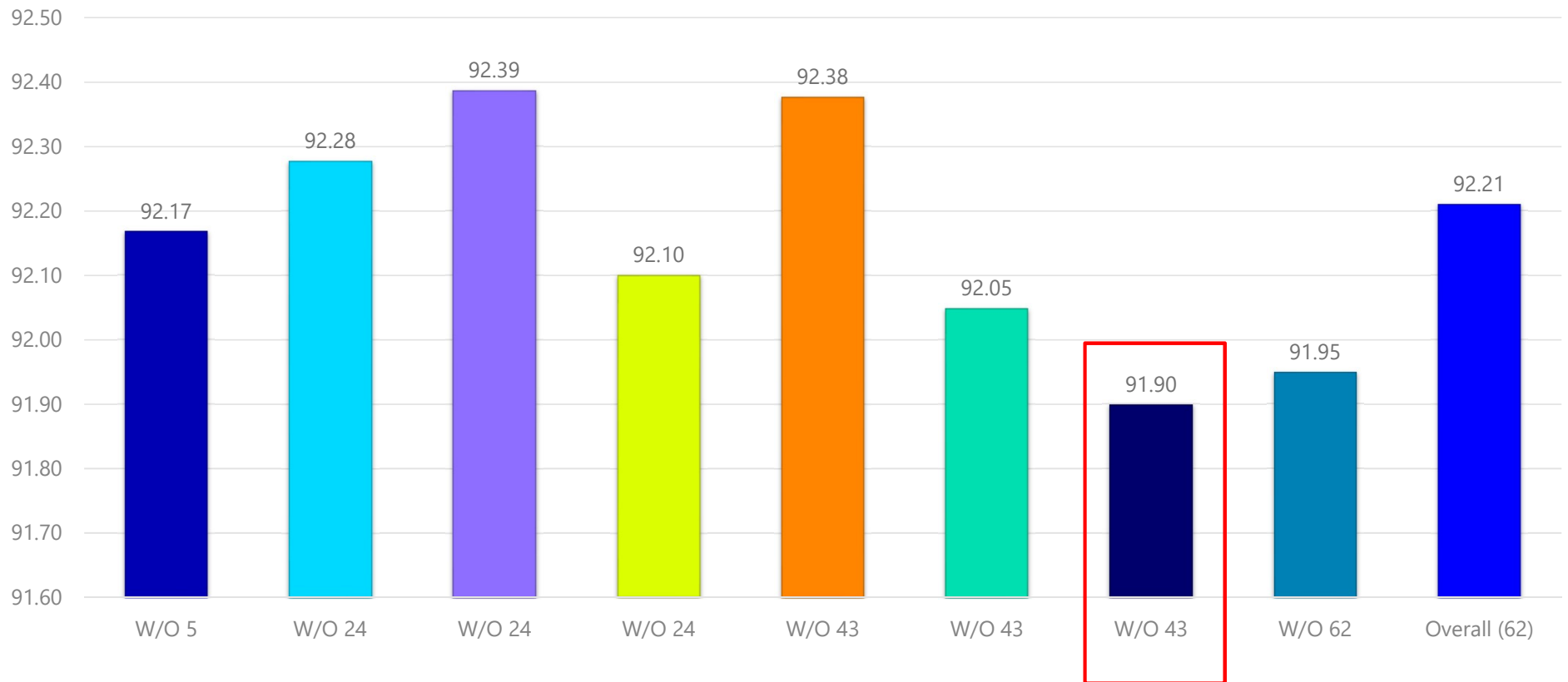
How much contribution each component makes?



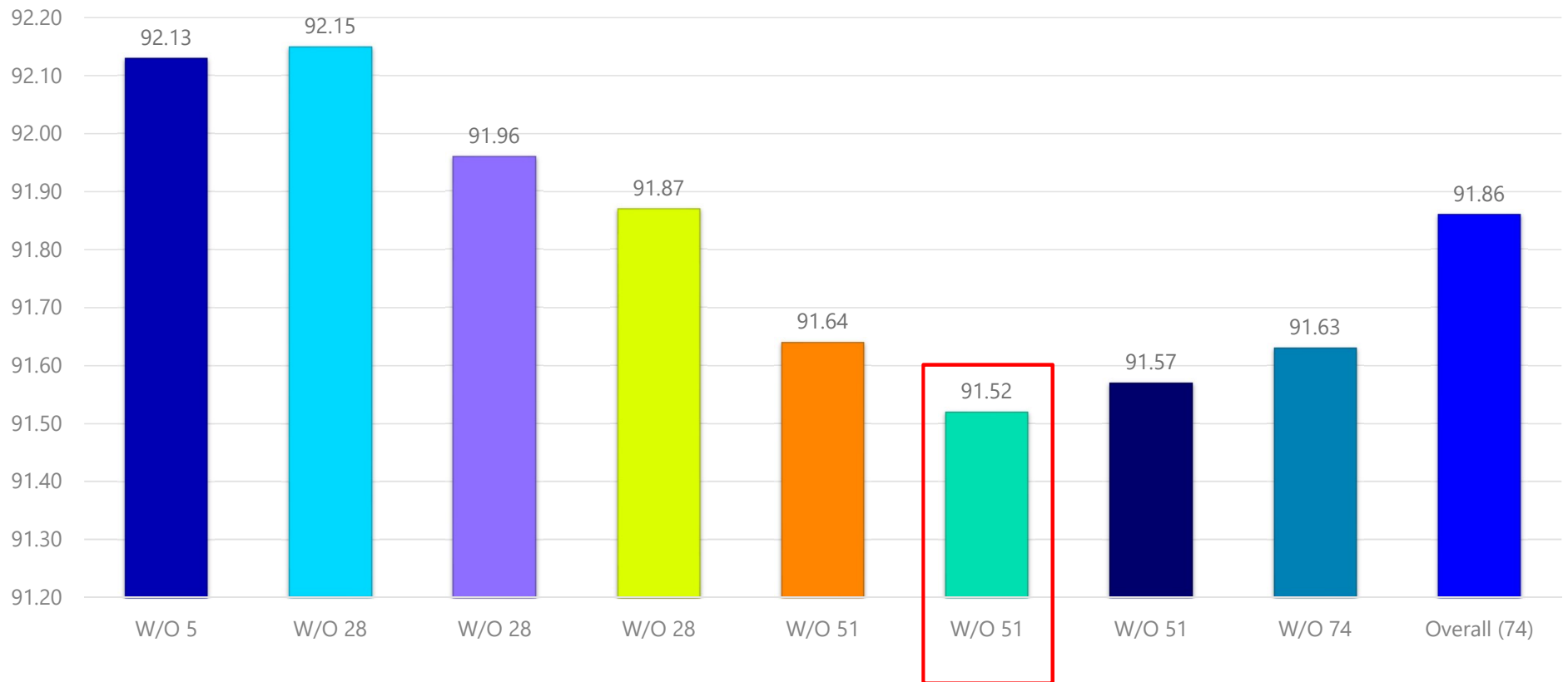
How much contribution each component makes?



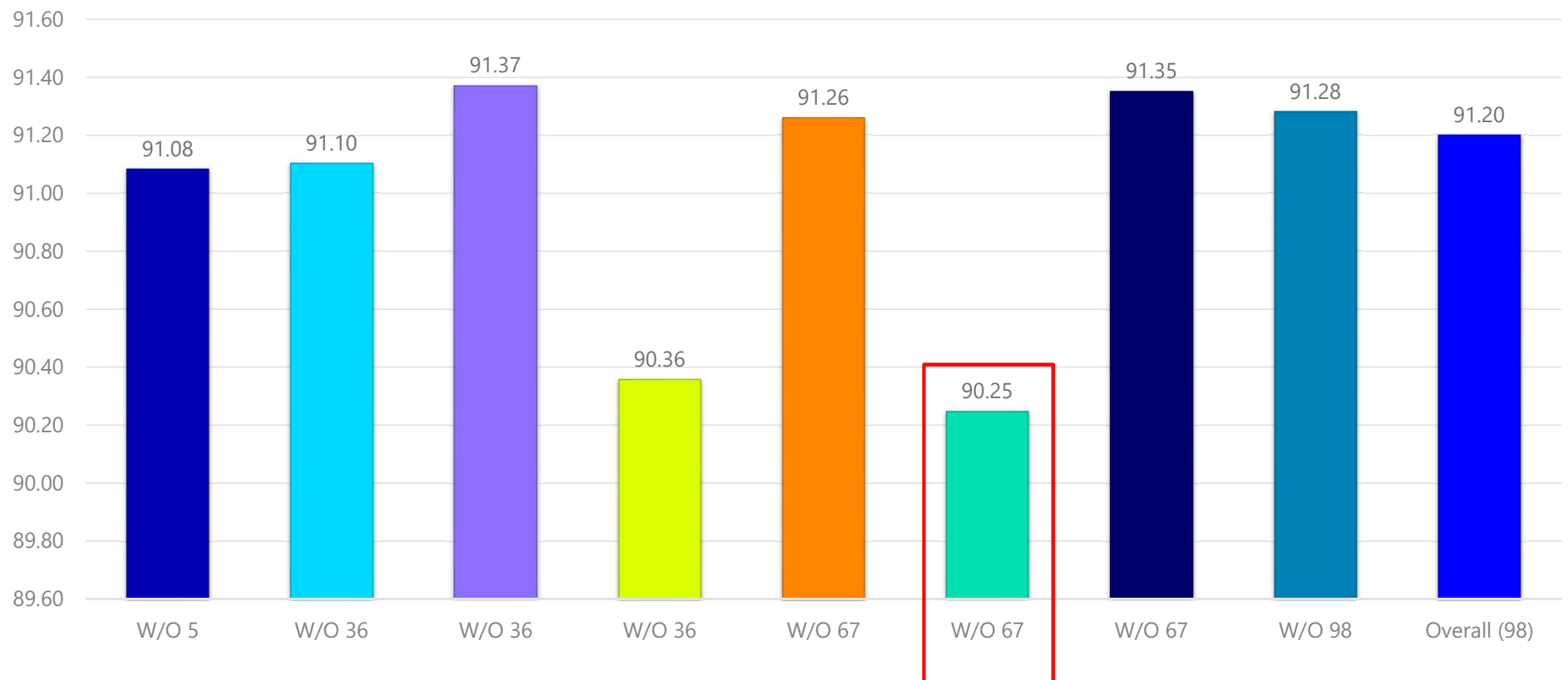
How much contribution each component makes?



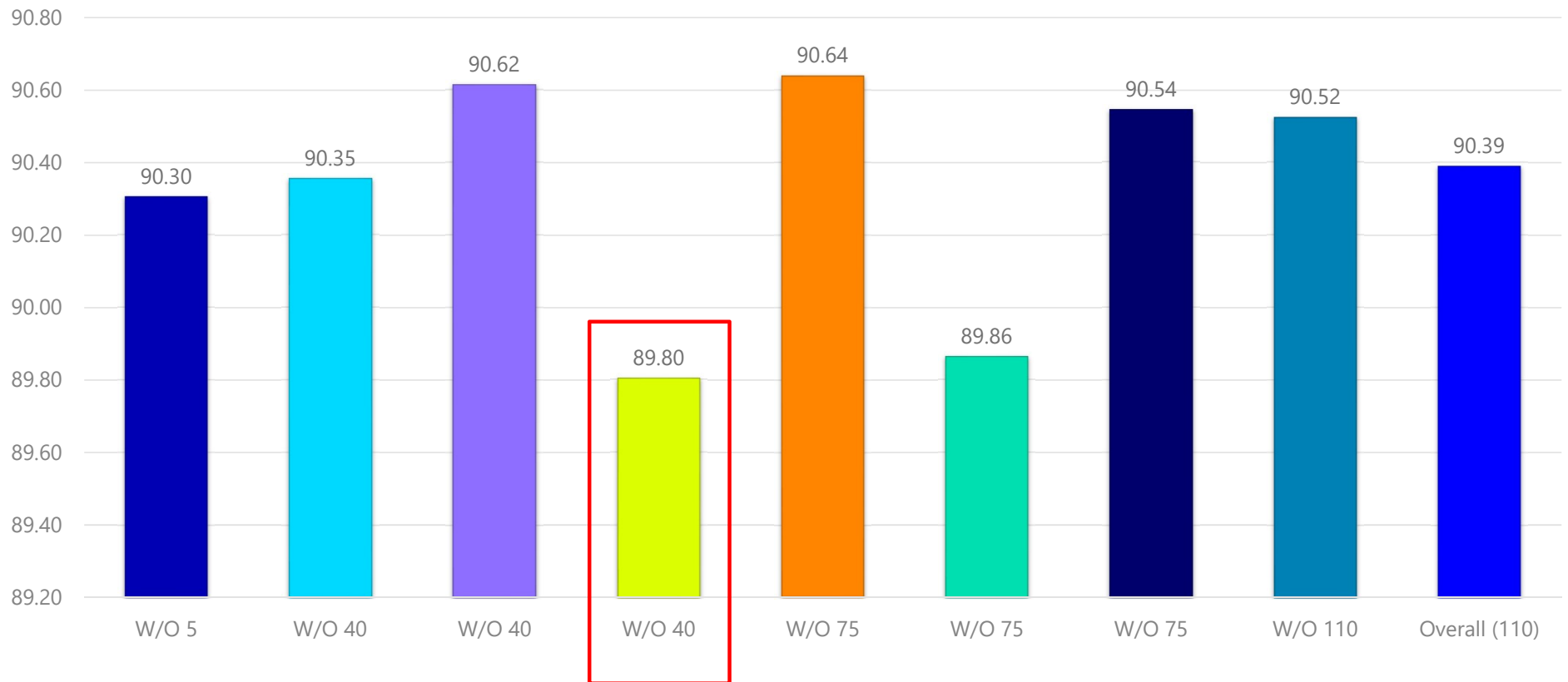
How much contribution each component makes?



How much contribution each component makes?



How much contribution each component makes?



Our approach: Merge-and-Run

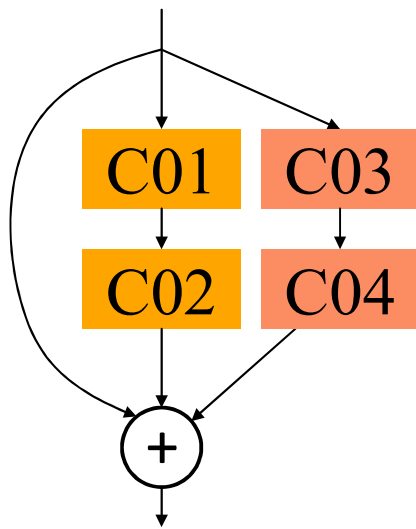
[Liming Zhao](#), Jingdong Wang, [Xi Li](#), [Zhuowen Tu](#), [Wenjun Zeng](#): On the Connection of Deep Fusion to Ensembling. [CoRR abs/1611.07718](#) (2016)

Design rules:

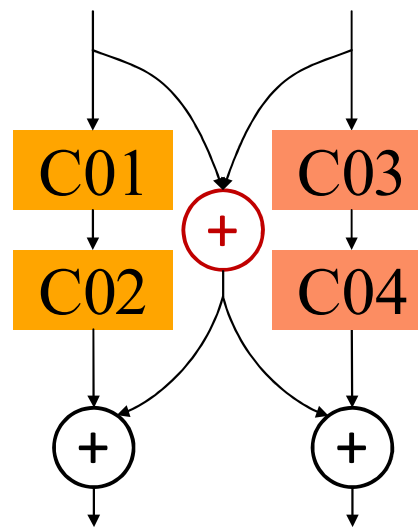
(1) Large ensemble size; (2) Avoid too deep networks

Design rules:

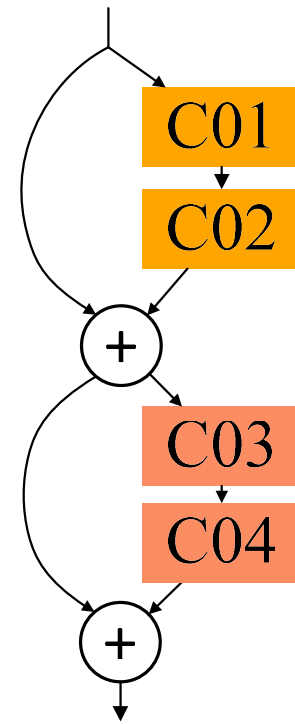
(1) Large ensemble size; (2) Avoid too deep networks



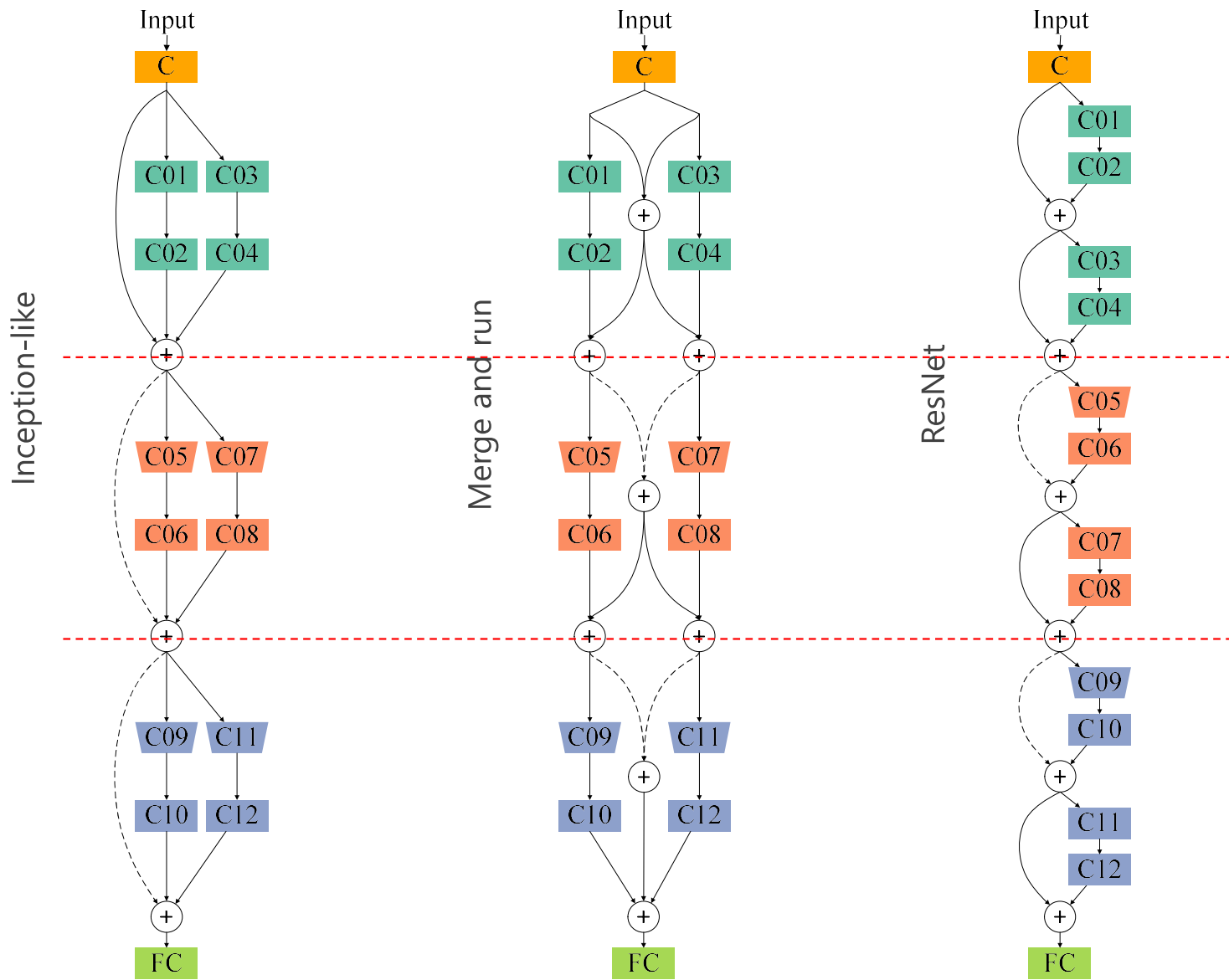
Inception-like
3 paths

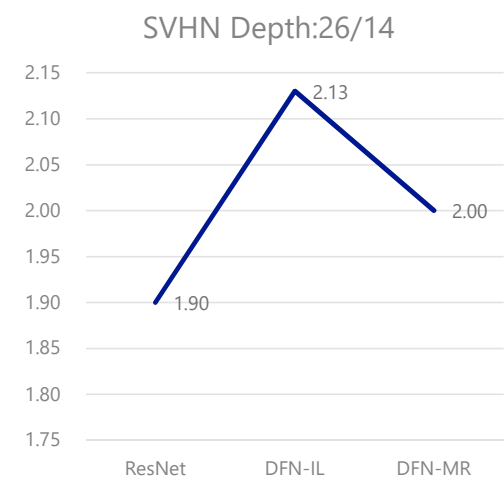
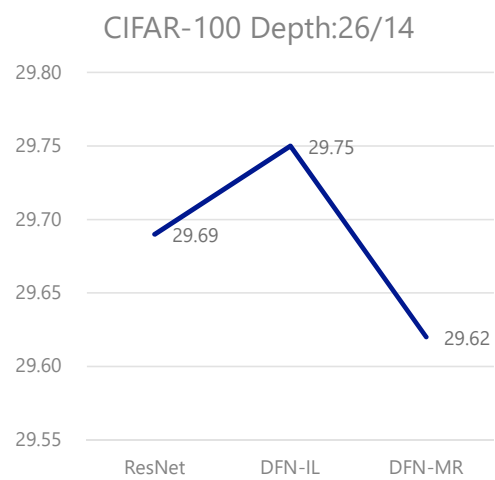
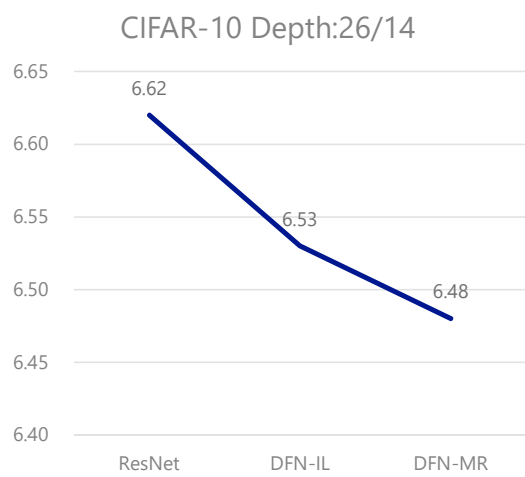


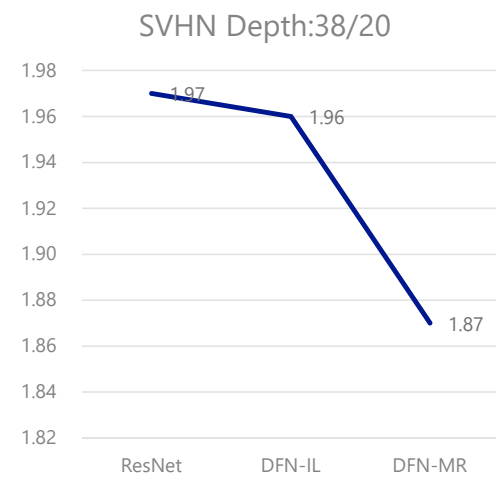
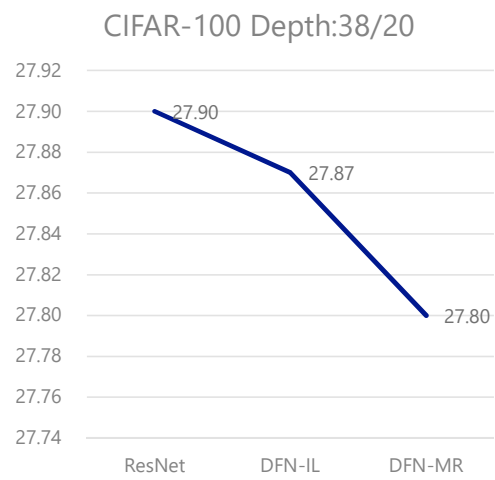
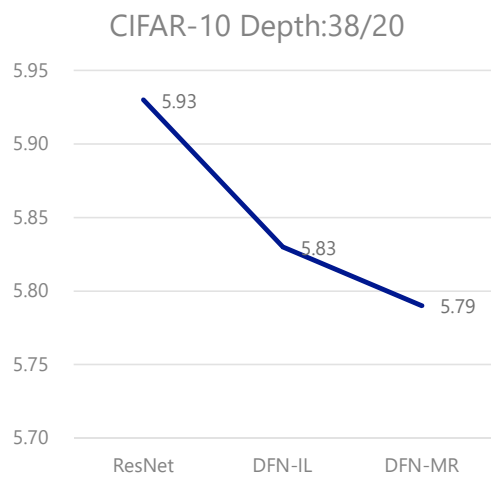
Merge and run
6 paths



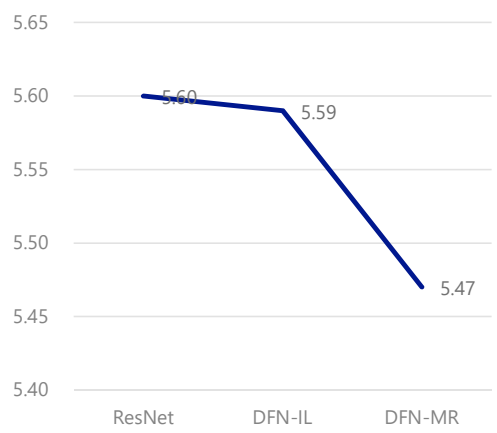
ResNet
4 paths



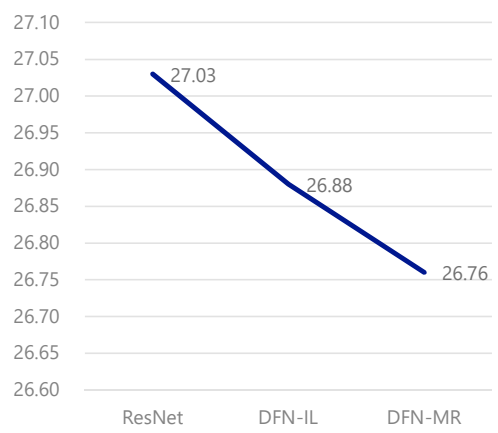




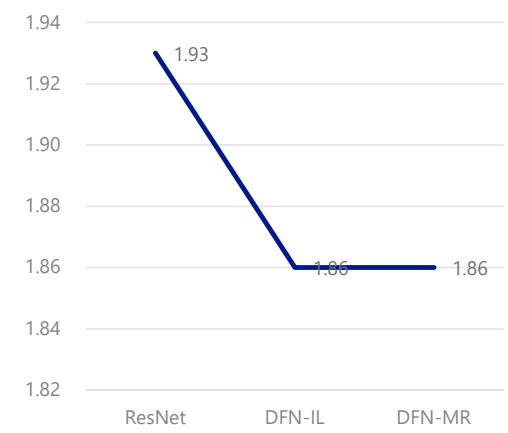
CIFAR-10 Depth:50/26



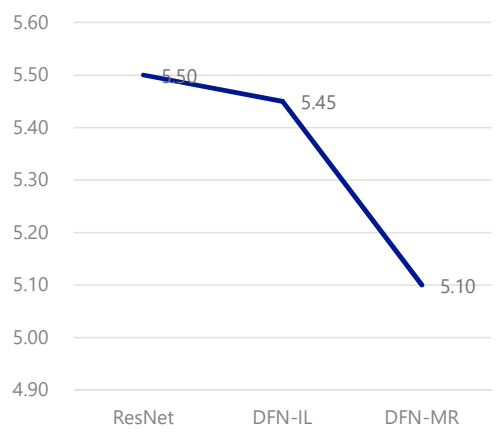
CIFAR-100 Depth:50/26



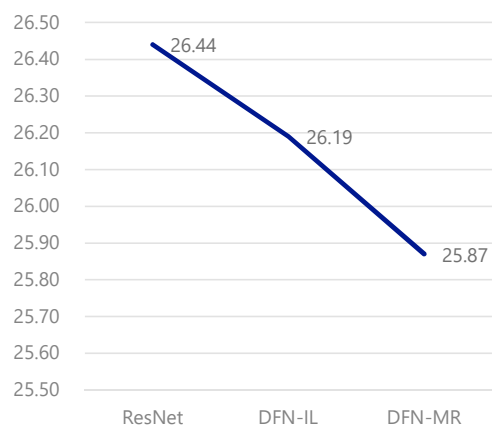
SVHN Depth:50/26



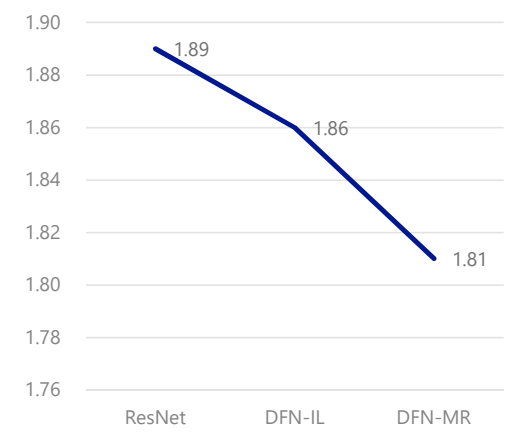
CIFAR-10 Depth:62/32



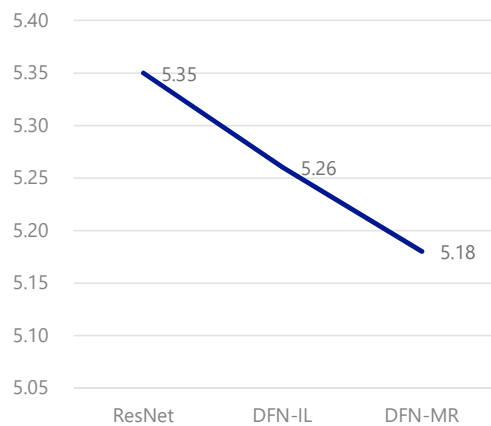
CIFAR-100 Depth:62/32



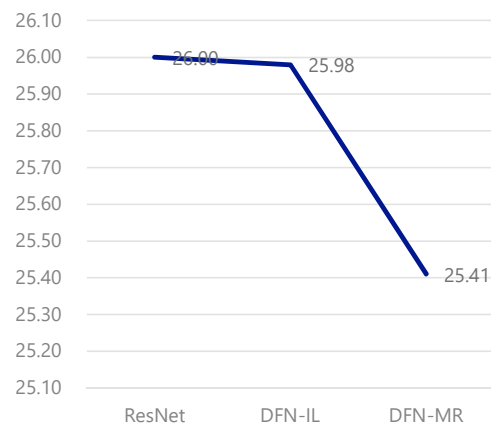
SVHN Depth:62/32



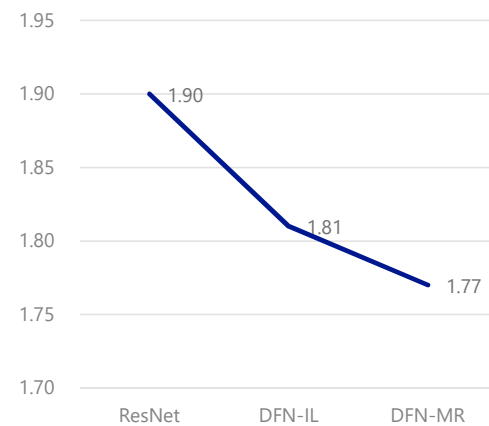
CIFAR-10 Depth:74/38

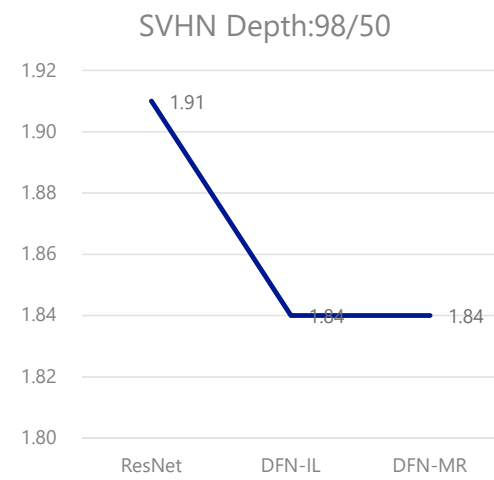
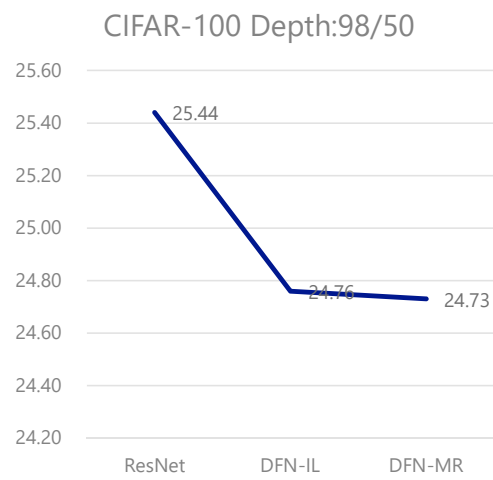
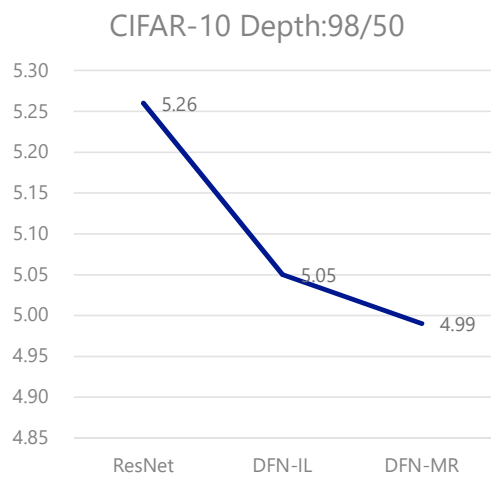


CIFAR-100 Depth:74/38

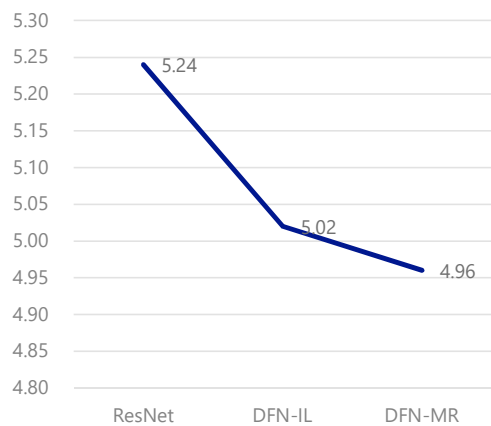


SVHN Depth:74/38

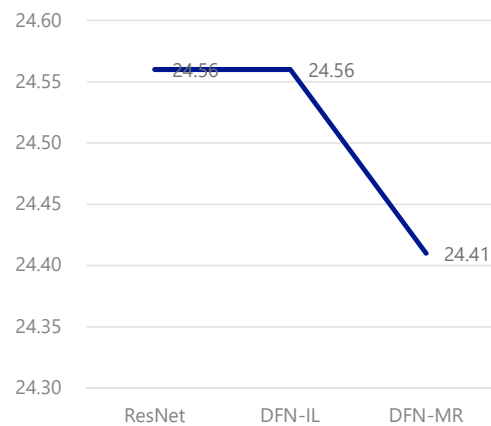




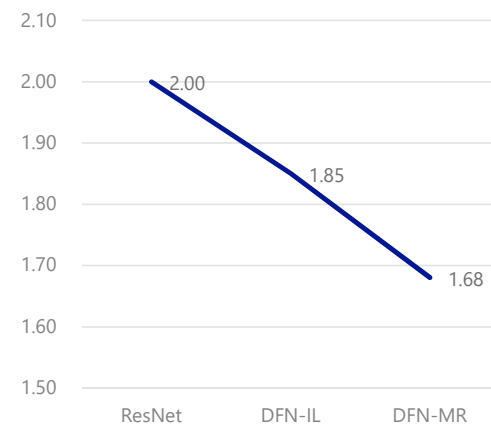
CIFAR-10 Depth:110/56



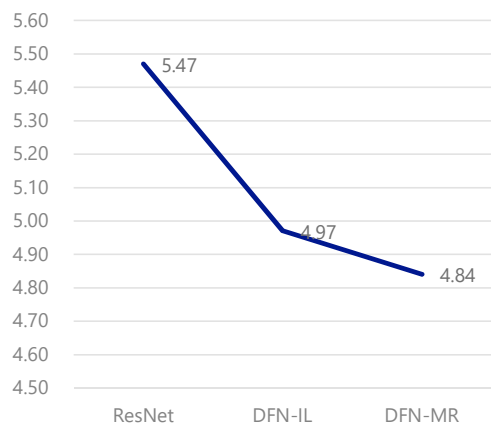
CIFAR-100 Depth:110/56



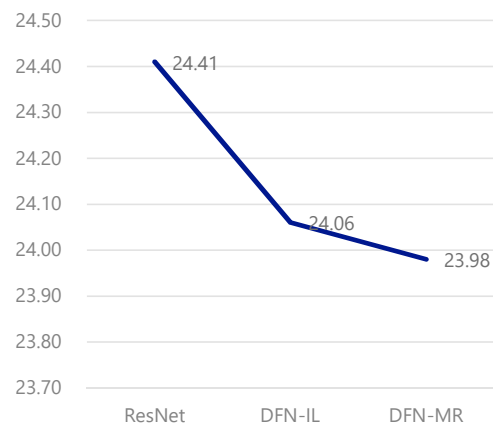
SVHN Depth:110/56



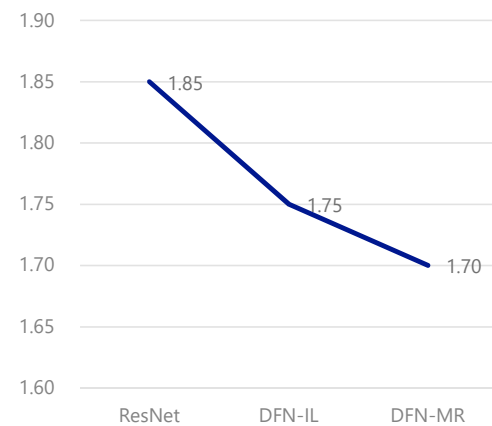
CIFAR-10 Depth:194/98



CIFAR-100 Depth:194/98



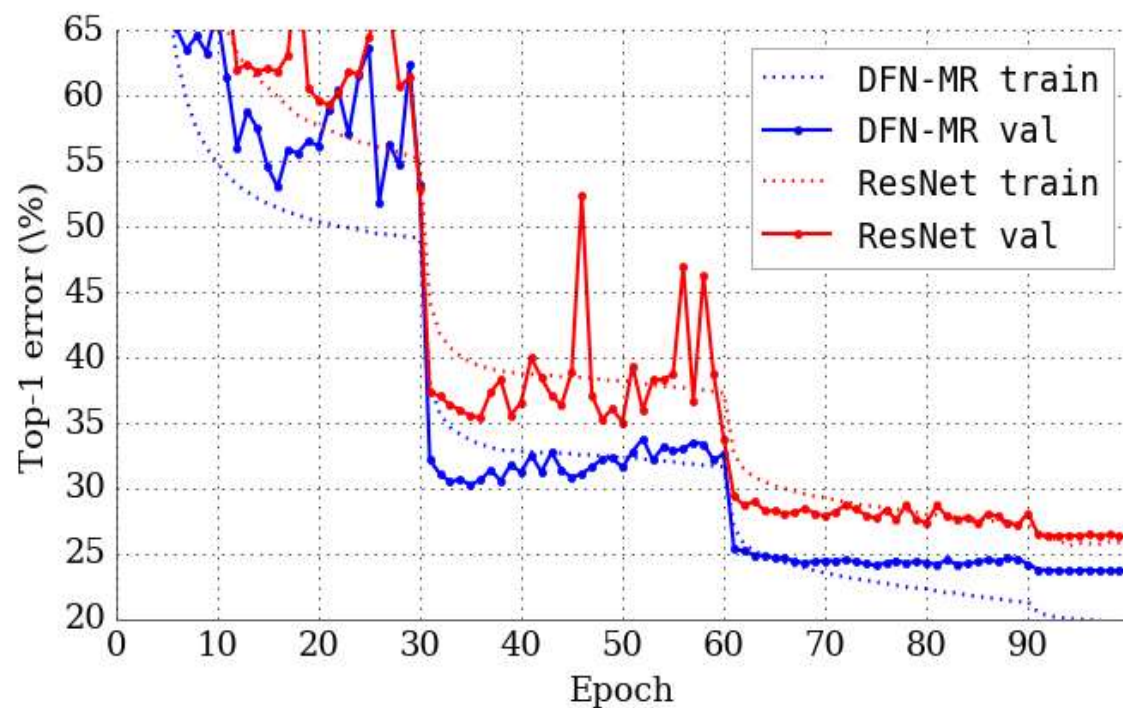
SVHN Depth:194/98



Comparison with state-of-the-arts

Method	Depth	#Params.	Cifar-10	Cifar-100	SVHN
DSN	-	-	7.97	34.57	1.92
FractalNet with DO/DP	21	38.6M	5.22	23.30	2.01
	21	38.6M	4.60	23.73	1.87
ResNet	110	1.7M	6.41	27.22	2.01
Multi ResNet	200	10.2M	4.35	20.42	-
Wide ResNet	16	11.0M	4.81	22.07	-
	28	36.5M	4.17	20.50	-
DenseNet	40	1.0M	5.24	24.42	1.79
	100	27.2M	3.74	19.25	1.59
DFN (ours)	56	1.7M	4.94	24.46	1.66
DFN (ours)	32	14.9M	3.94	19.25	1.51
DFN (ours)	50	24.8M	3.57	19.00	1.55

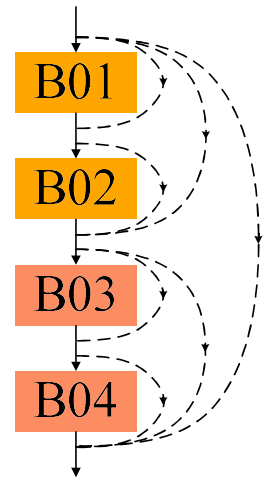
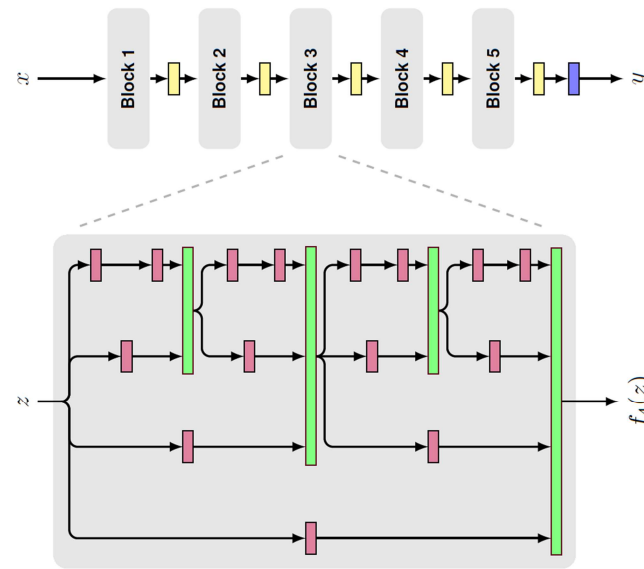
ImageNet Results



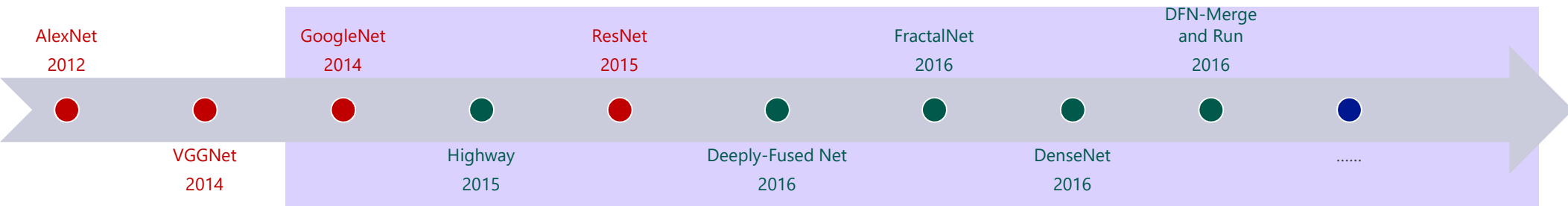
	ResNet-101 (He's)	ResNet-101	DFN
		44.5M	43.3M
Top1 validation error	23.60	26.41	23.66
Top5 validation error	7.10	8.50	6.81
Top-1 training error	17.00	25.75	19.72
Top-5 training error	-	8.12	6.59

Other network structures

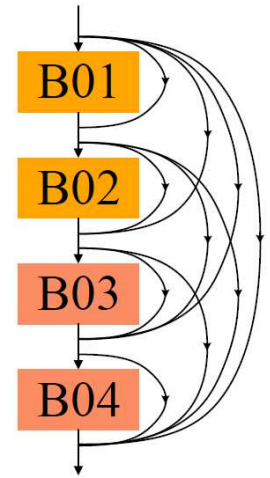
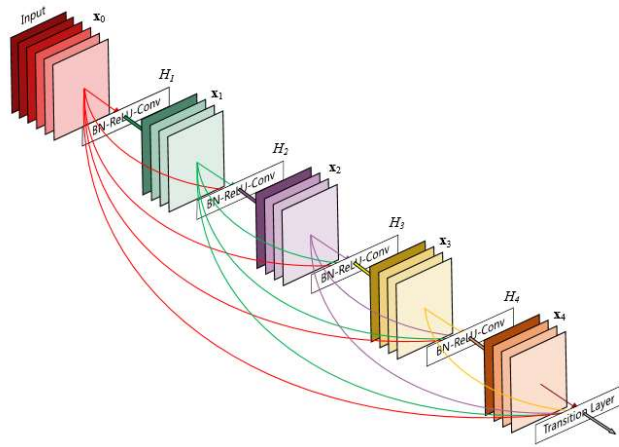
FractalNet, 2016



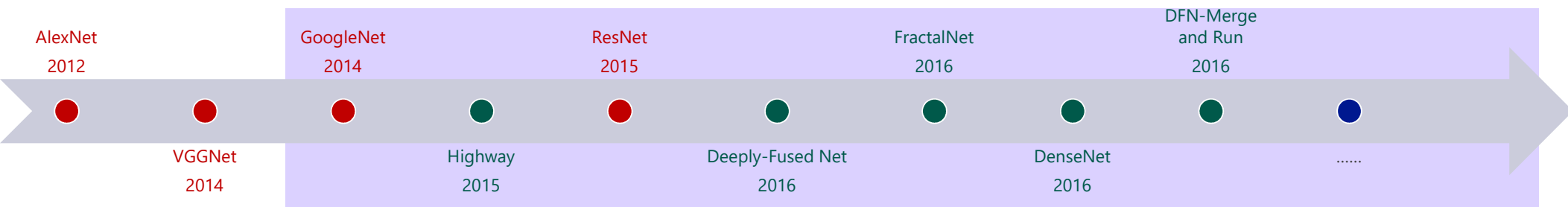
Multi-path, without identity connection, **deep fusion in deep fusion**



DenseNet, 2016



Dense connection: Express way to **all layers**, deeper is better
Deep graph fusion



Summary

Deep fusion

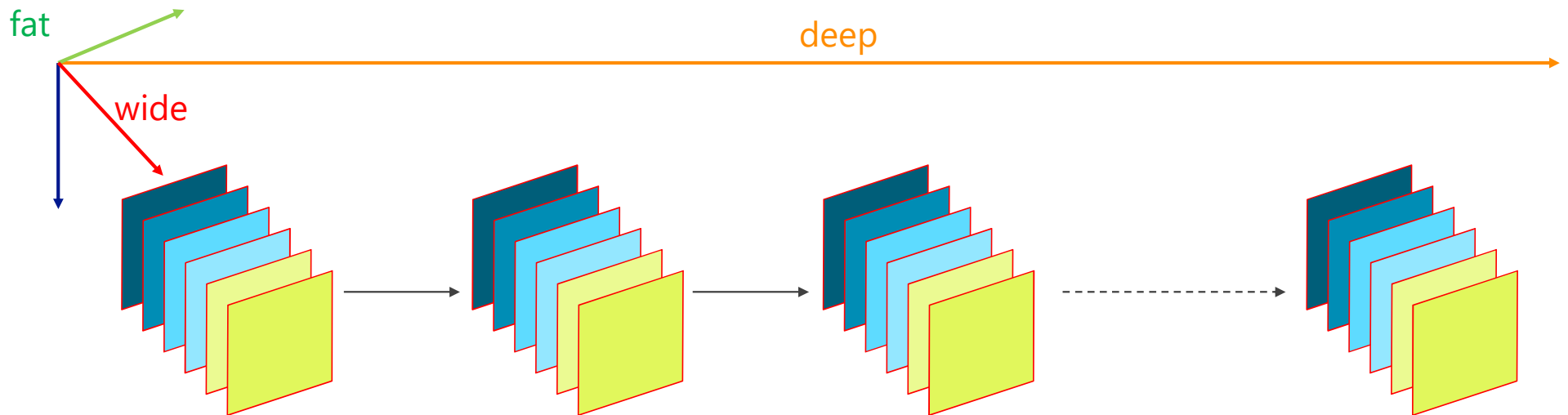
Multiple paths

Long and short

Express way between layers

Weight sharing

Ultra-deep is misleading!
What's next? **Wider?** **Fatter?**



Thanks!
Q&A