

CNN Attention-based Networks

+ Attention, CNN Review

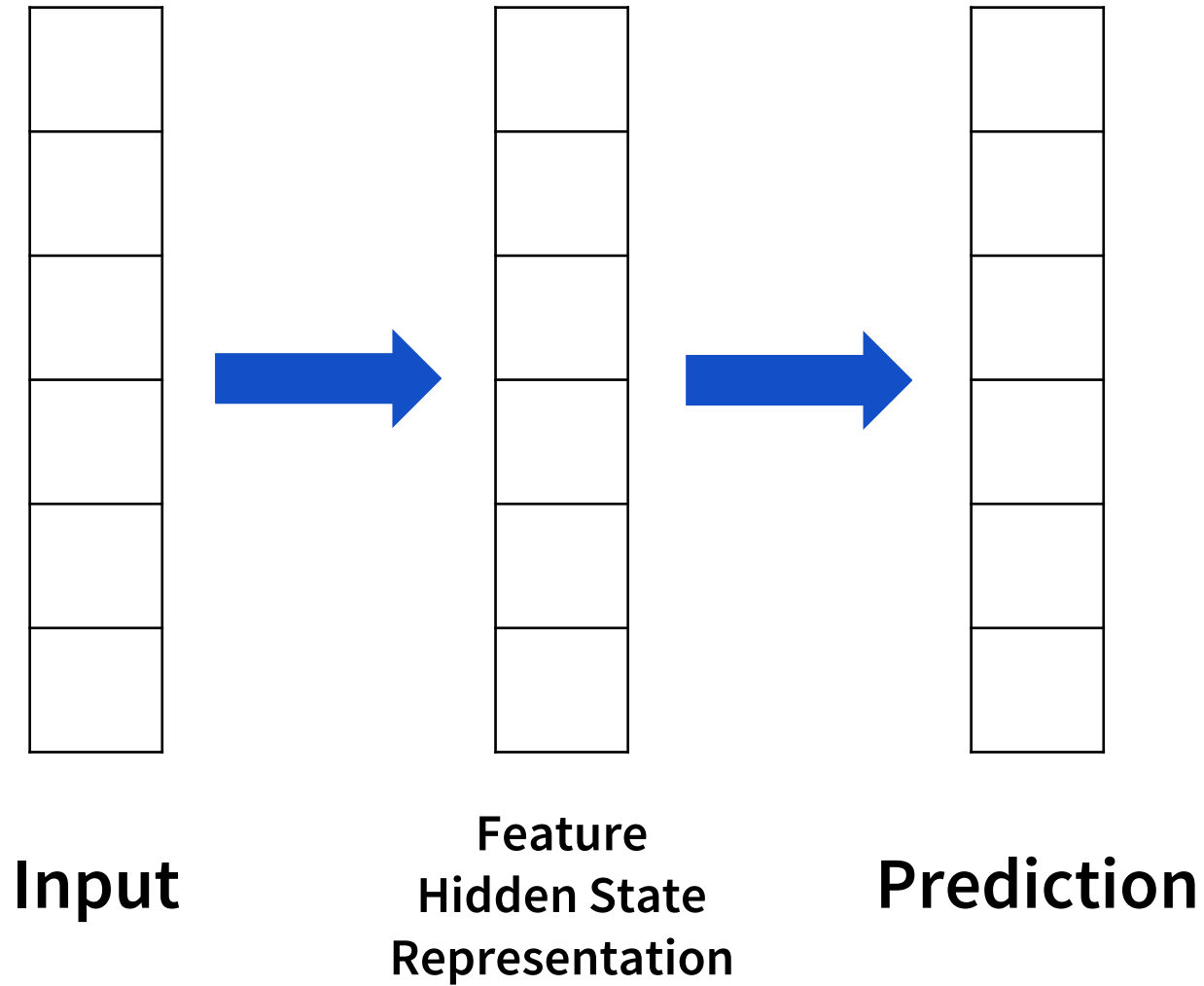
Tensorflow-KR **PR-163**, Taeoh Kim
MVPLAB, Yonsei Univ

Contents

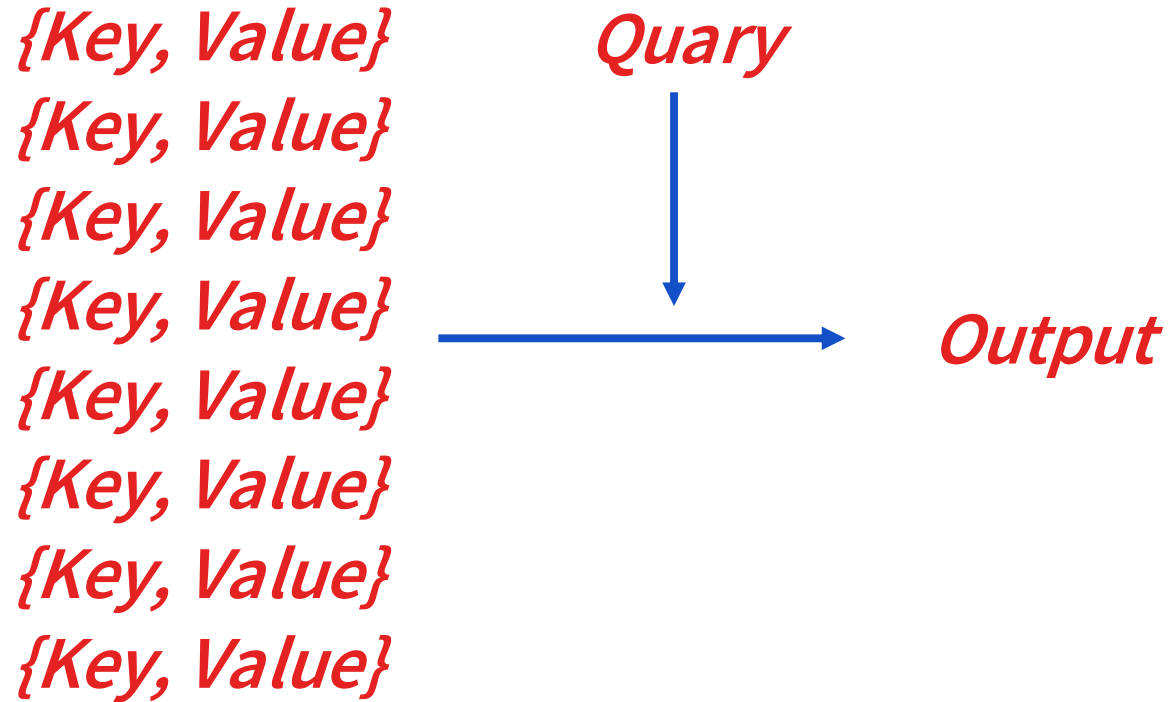
- *Attention, Self-Attention in NLP*
- *CNN-Review*
- *CNN Attention Networks for Recognition*
- *CNN Attention Networks for Other Vision Tasks*

Review 1: Attention

Neural Networks



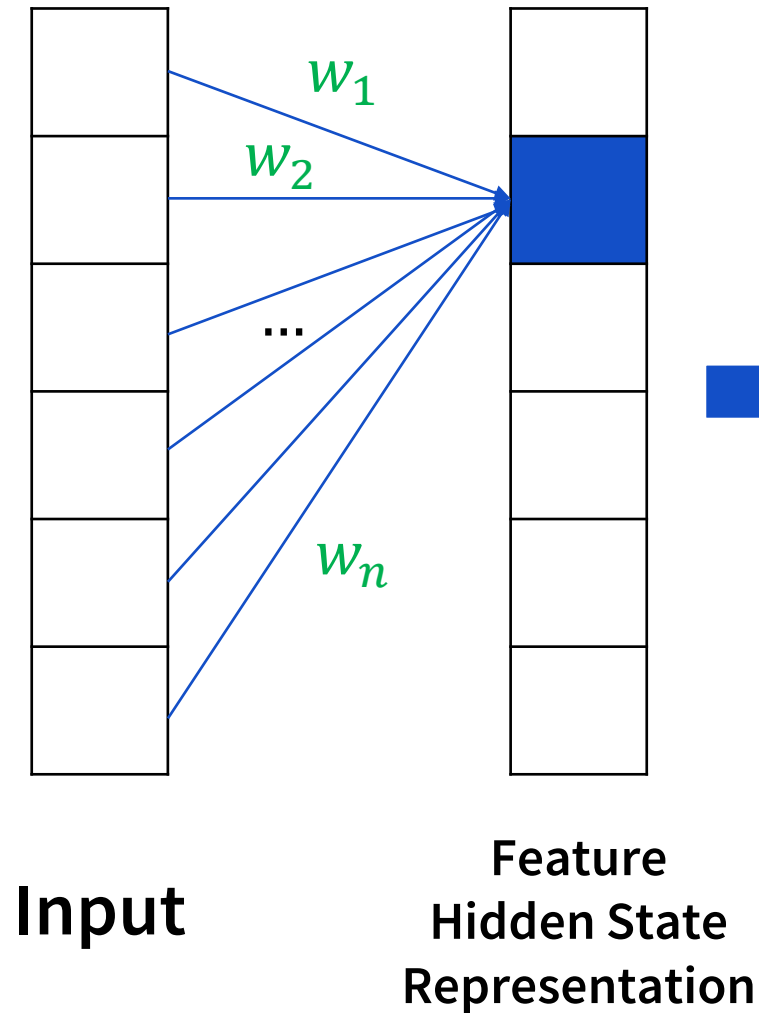
Attention



$$y = \sum_i w_i x_i$$

$$y = \sum_i f(\text{Query}, \text{Key}) \times \text{Value}$$

Fully Connected Neural Network

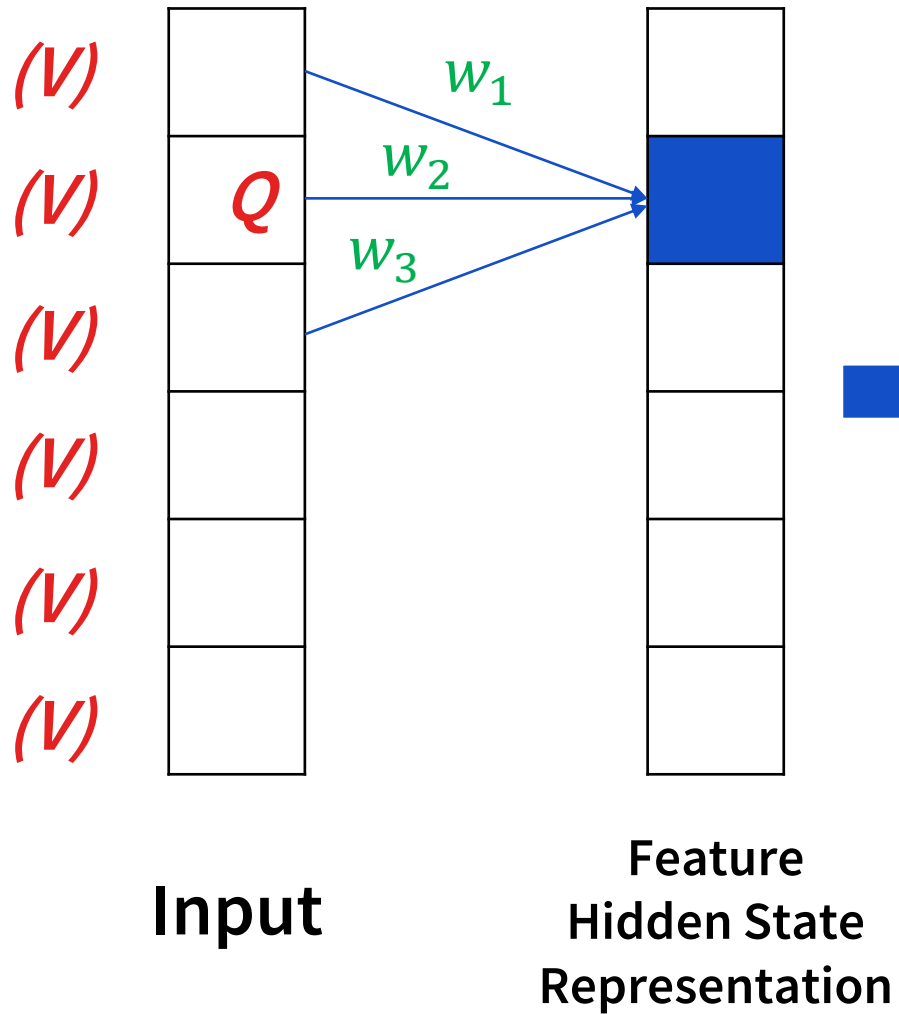


Fully Connected NN

*Represent **Blue**
using **Weighted** Sum of **Inputs***

without Constraint

Convolutional Neural Network

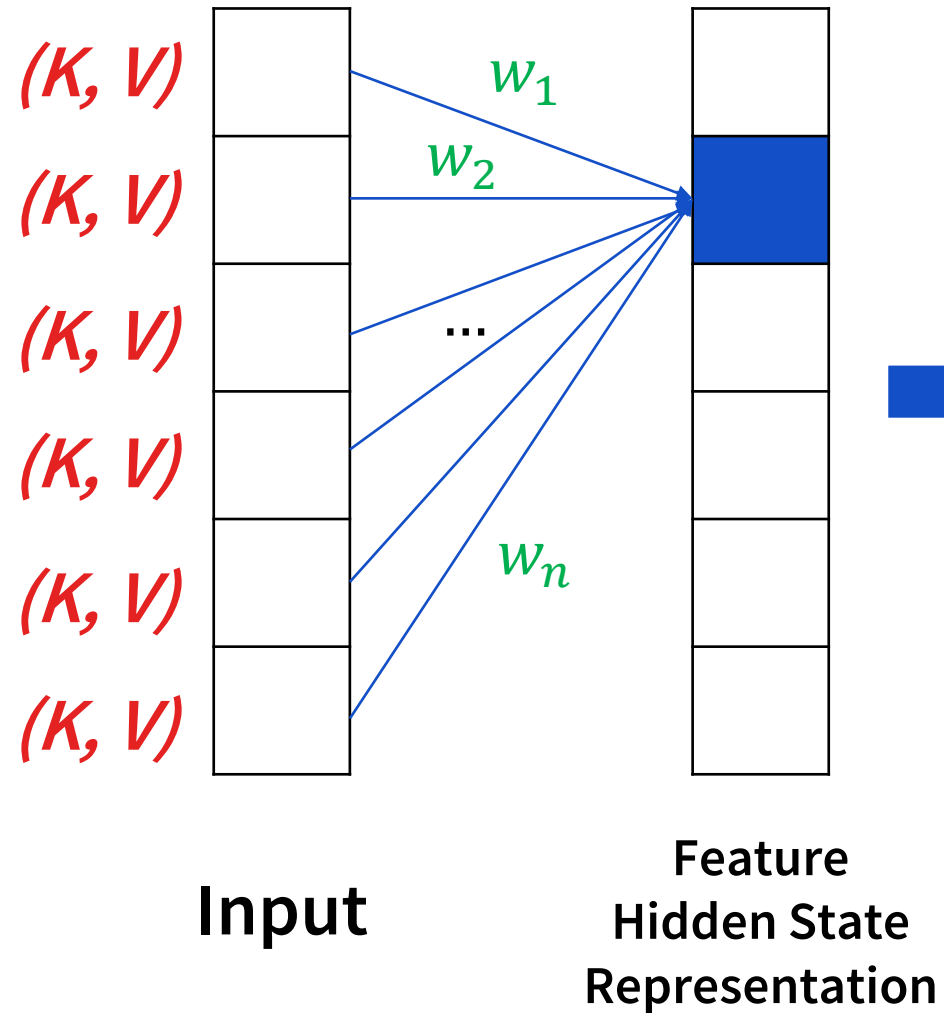


Convolutional NN

*Represent **Blue**
using **Weighted** Sum of **Inputs***

from Current Position

Attention



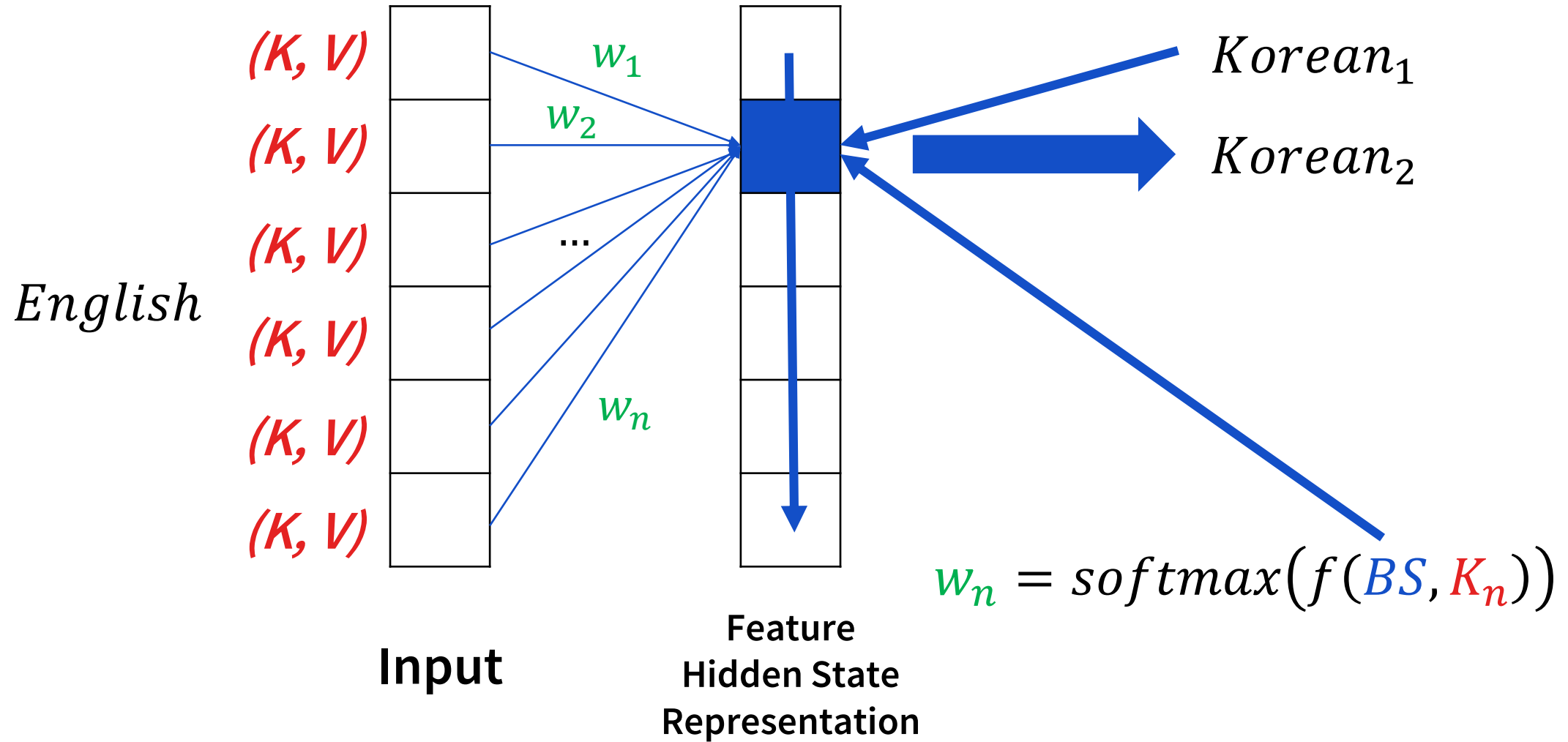
Attention

*Represent **Blue**
using **Weighted** Sum of **Inputs***

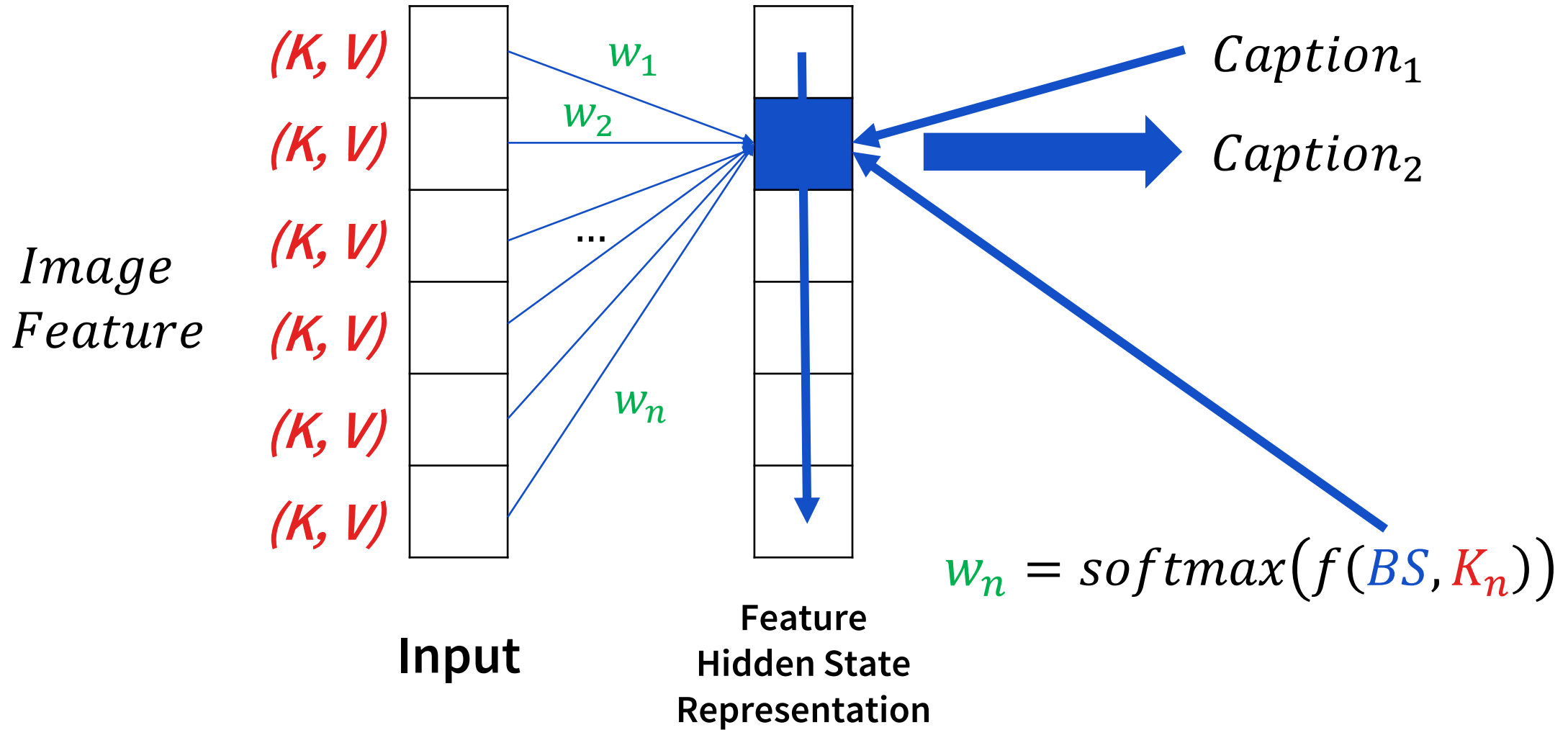
*from **Blue State** and **Inputs***

$$w_n = \text{softmax}(f(\text{BS}, K_n))$$

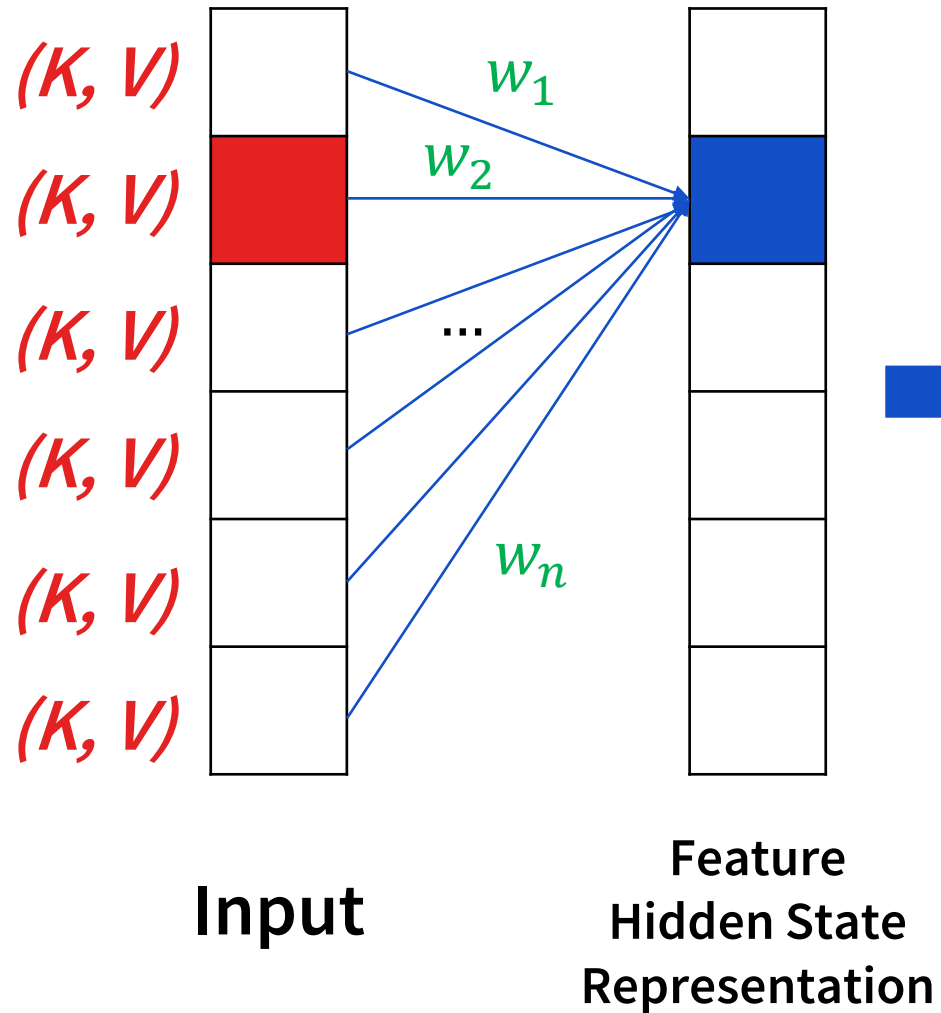
ex) Machine Translation (PR-055)



ex) Image Captioning



Self-Attention



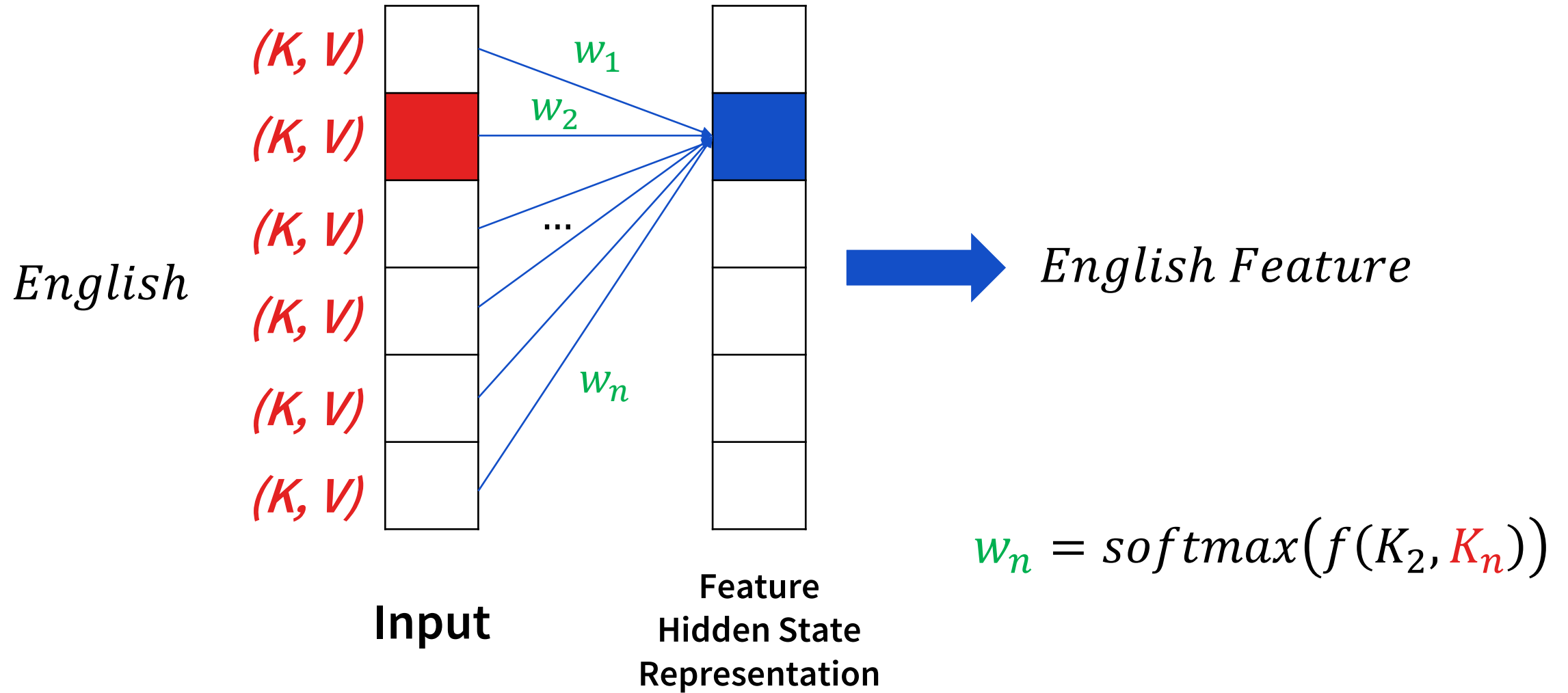
Self-Attention

*Represent **Blue**
using **Weighted Sum of Inputs***

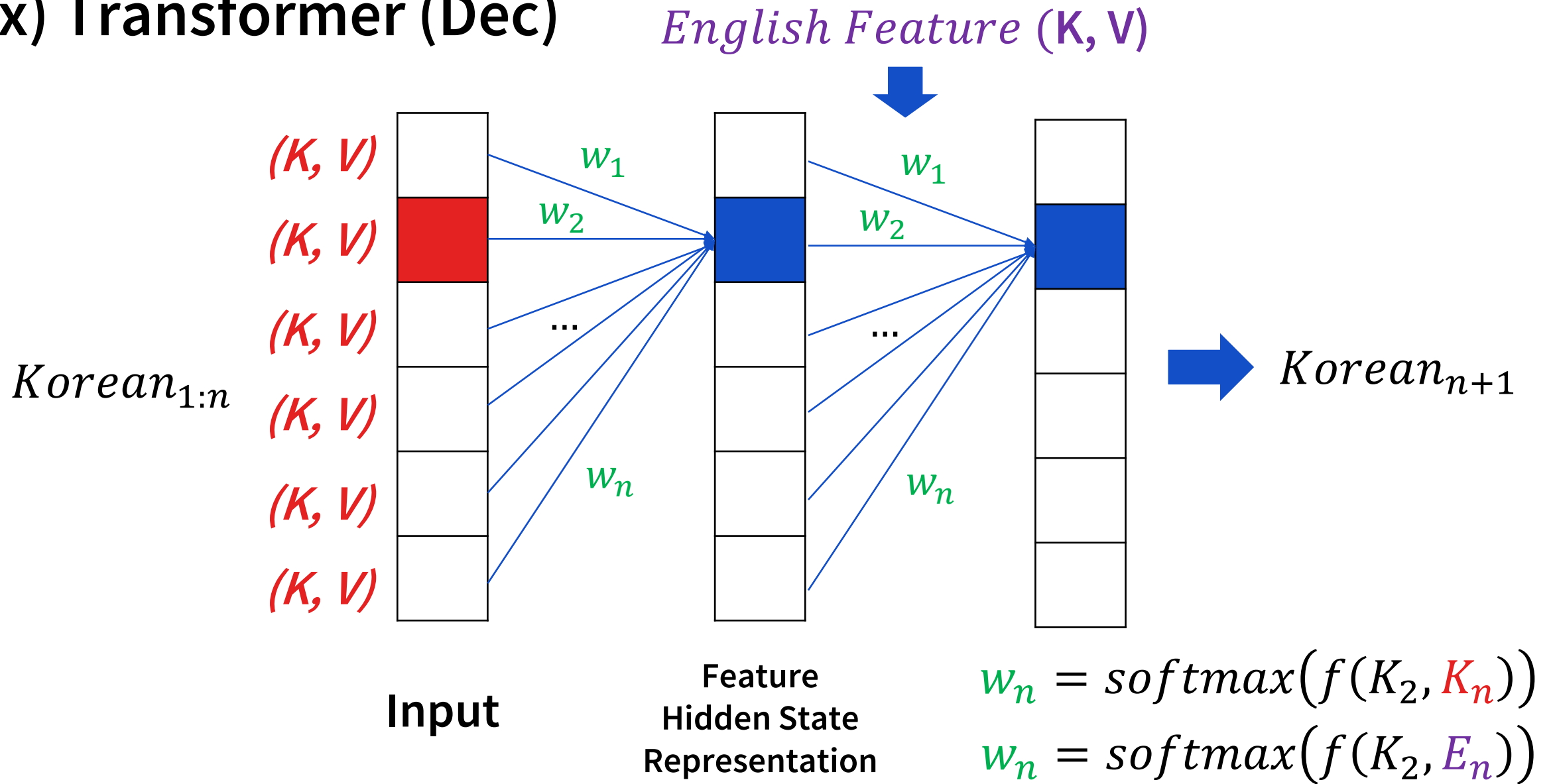
*from **Input and Inputs***

$$w_n = \text{softmax}(f(K_2, K_n))$$

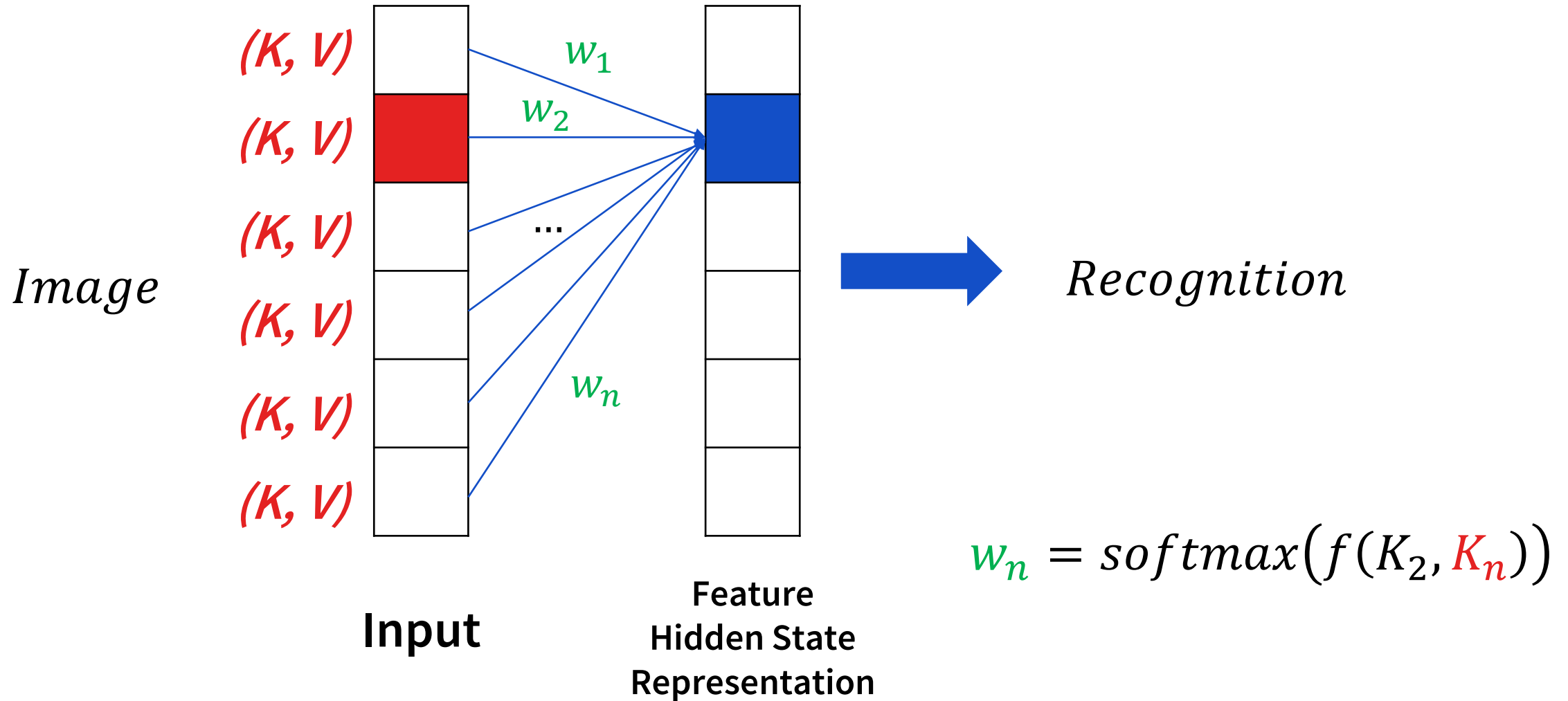
ex) Transformer (Enc) (PR-049, PR-161)



ex) Transformer (Dec)

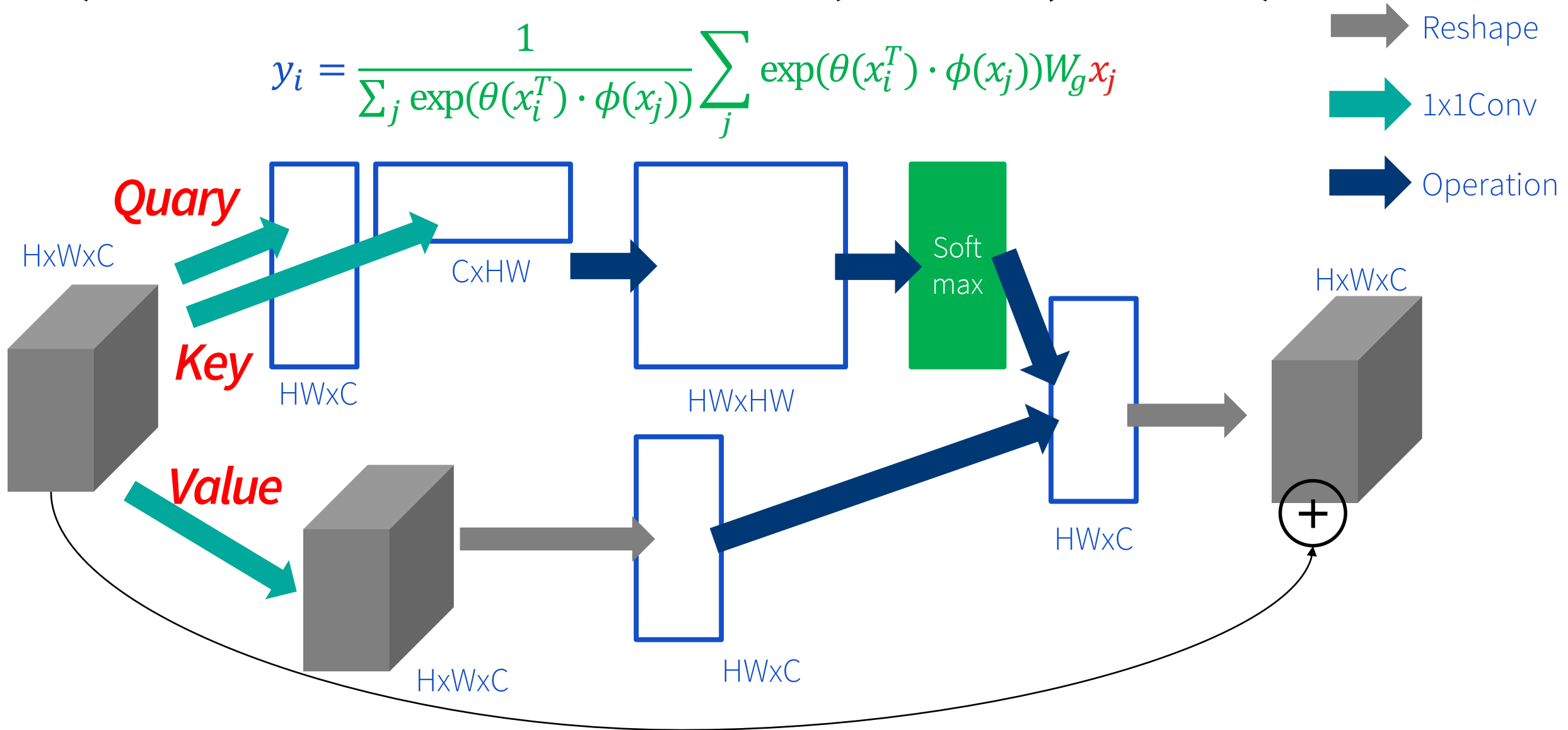


CNN Self-Attention for Image = Representation

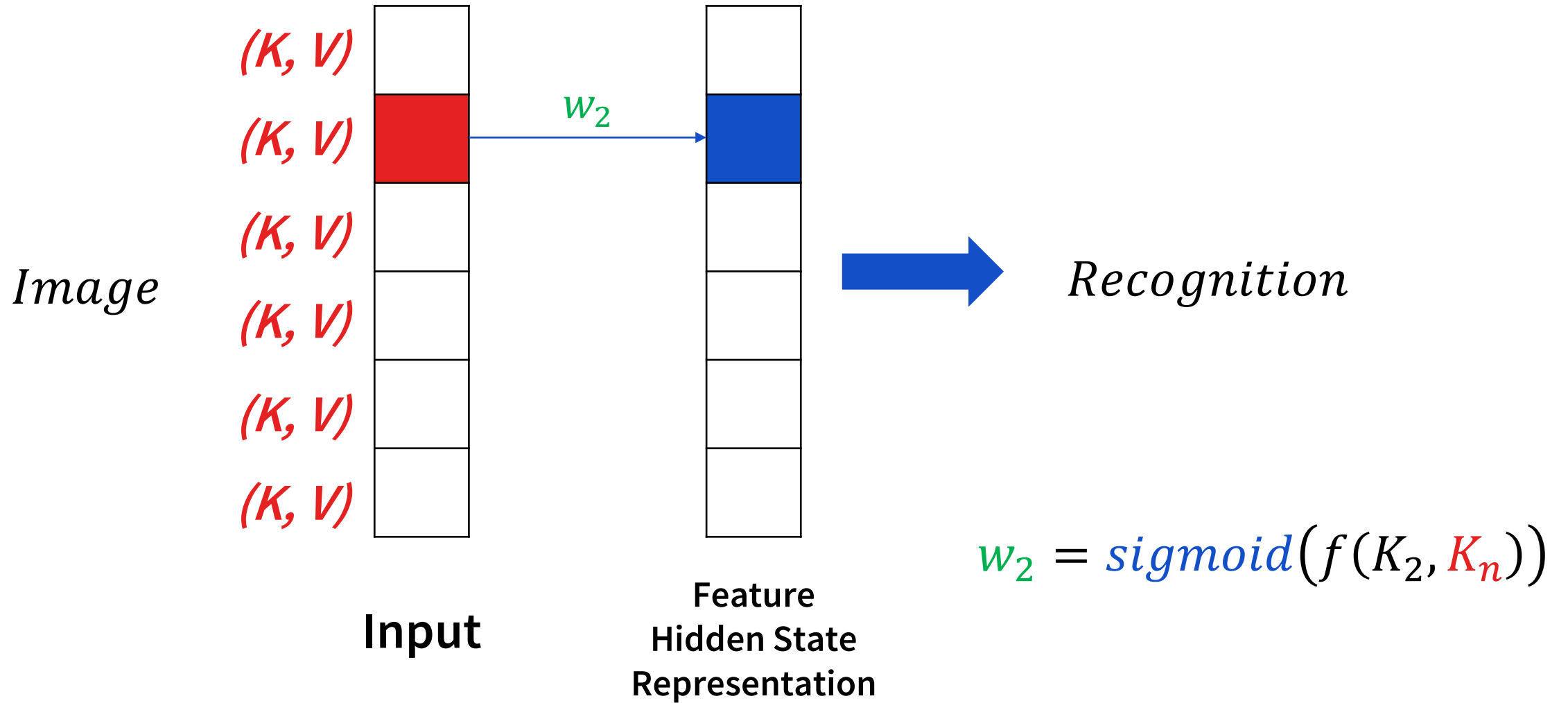


ex) Non-local Neural Networks (CVPR18, PR-083)

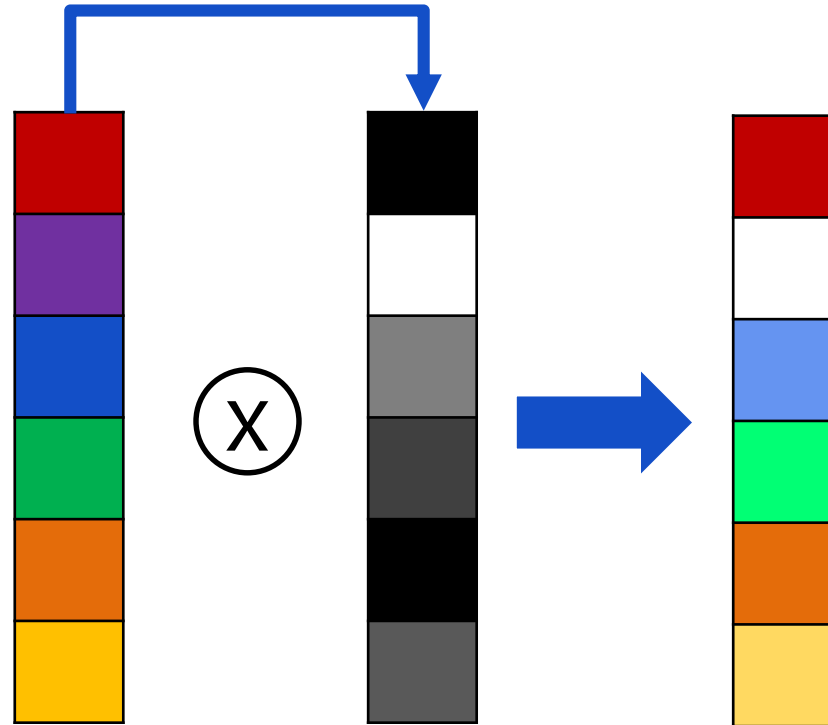
$$y_i = \frac{1}{\sum_j \exp(\theta(x_i^T) \cdot \phi(x_j))} \sum_j \exp(\theta(x_i^T) \cdot \phi(x_j)) W_g x_j$$



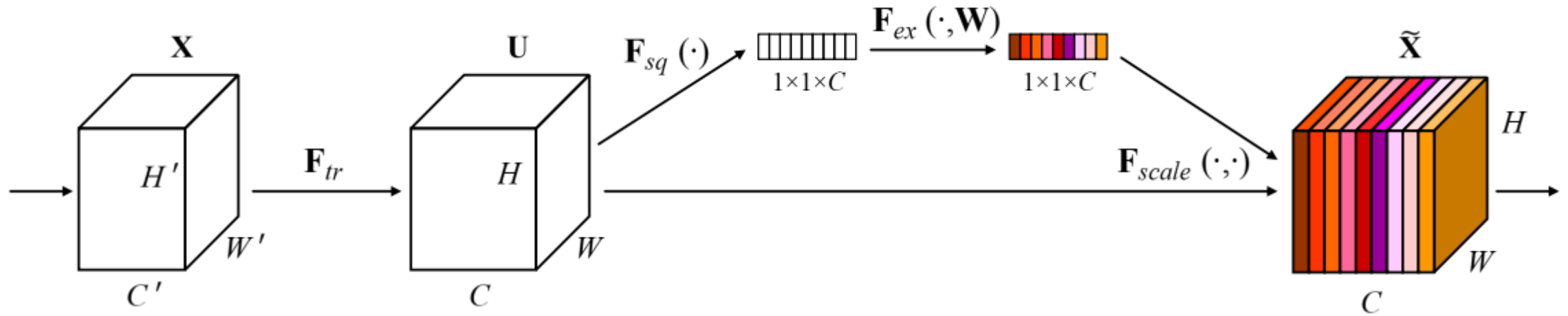
CNN Simplified-Attention for Image = Recalibration



CNN Simplified-Attention for Image = Recalibration



ex) Squeeze-and-Excitation Networks (CVPR18)

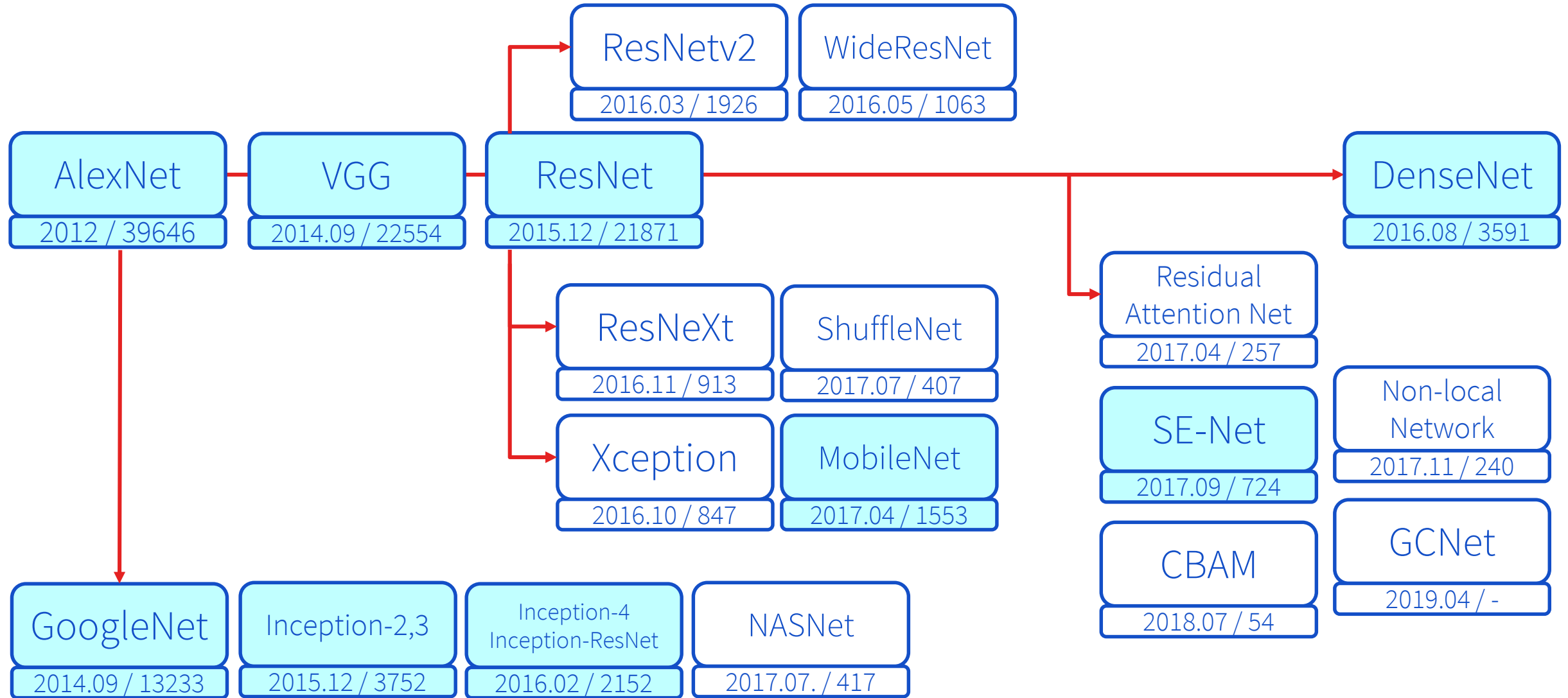


Summary

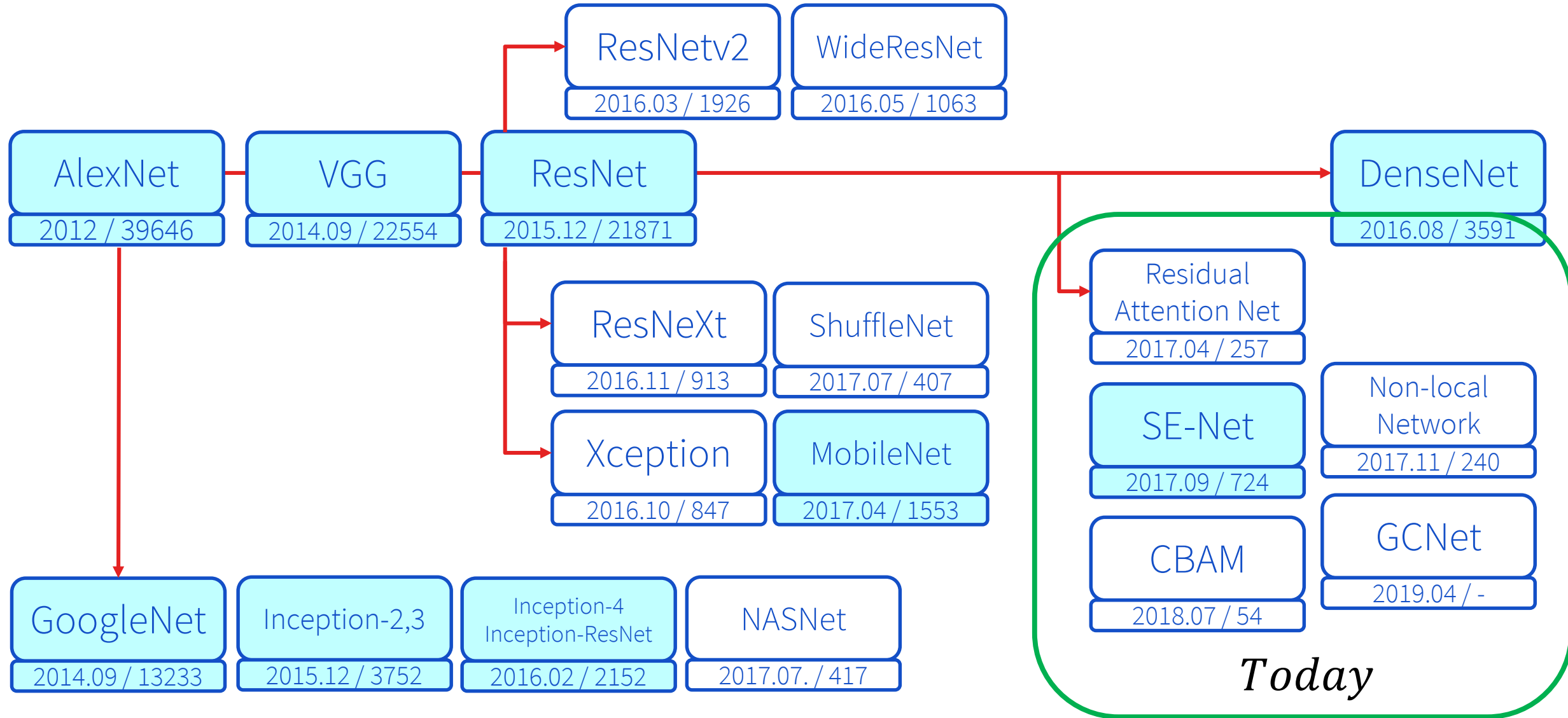
Attention	Quary	Structure	Objective	Examples
Attention	Current States	Recurrent	Representation	NMT, Captioning, VQA
Self-Attention	Input Itself	Feed-Forward	Representation	Transformer Non-local NN
	Input Itself	Feed-Forward	Recalibration	SE-Net, RAN, CBAM

Review 2: CNN Networks

CNN Review



CNN Review



Plain Networks

ILSVRC12

ILSVRC14

AlexNet

VGG

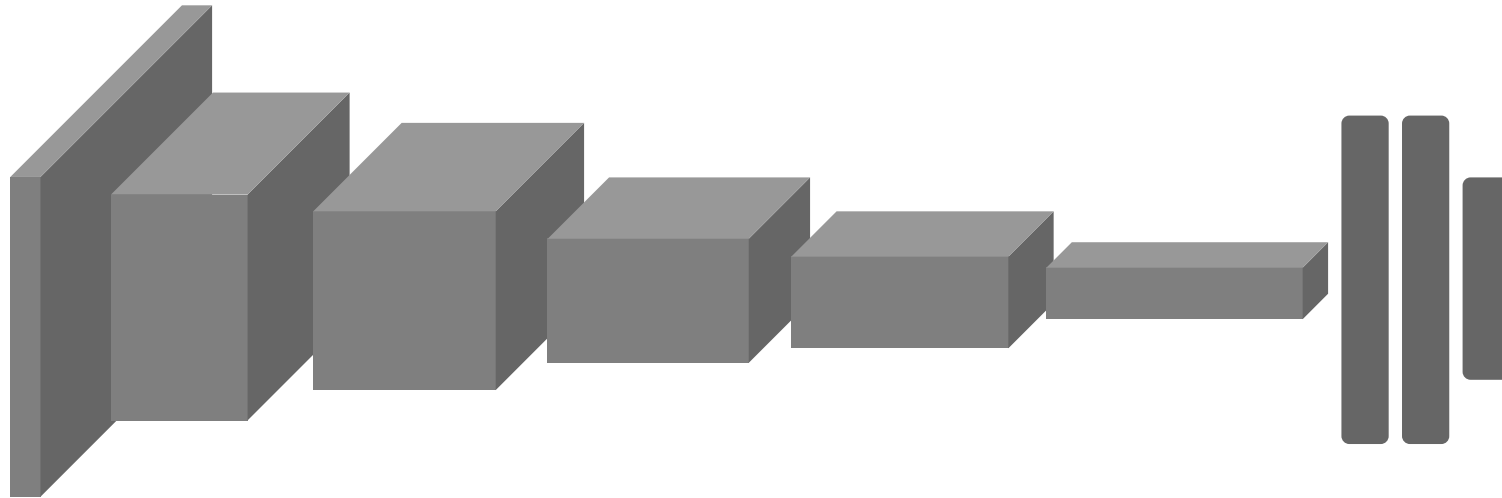
ResNet

2012 / 39646

2014.09 / 22554

2015.12 / 21871

- *Plain Networks using Max-Pooling*
- *Low Performance / Large Parameters, Operations*

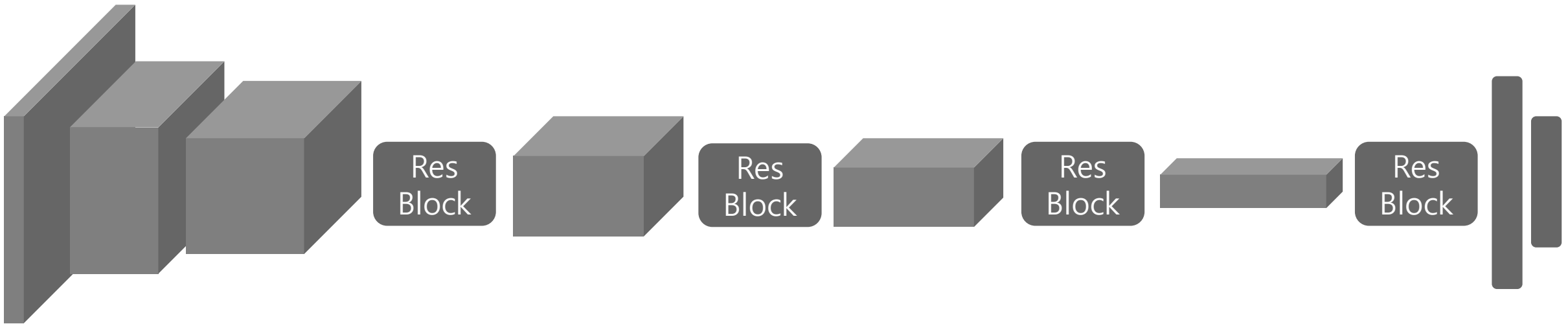


ResNet

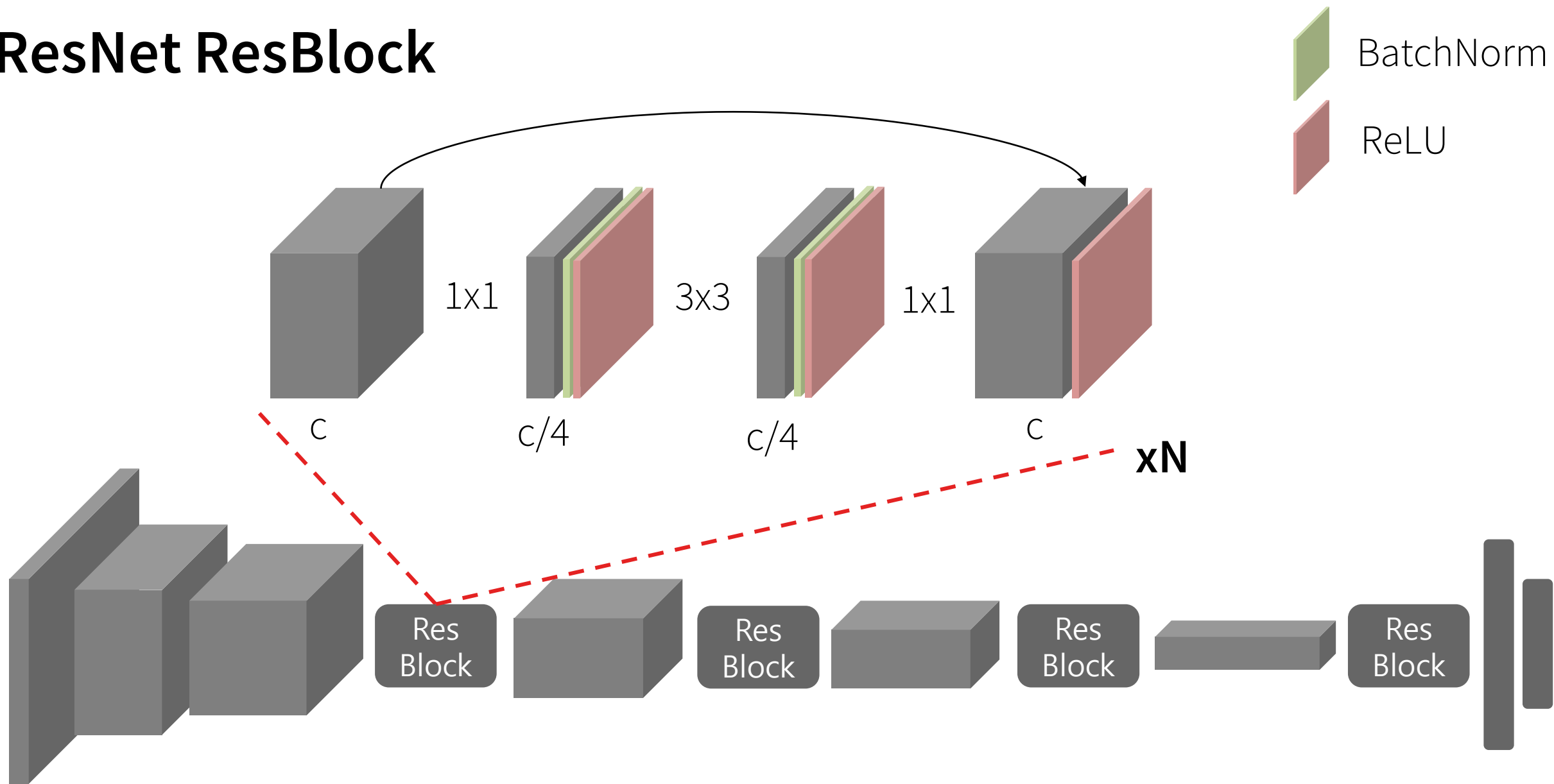
ILSVRC15

AlexNet	VGG	ResNet
2012 / 39646	2014.09 / 22554	2015.12 / 21871

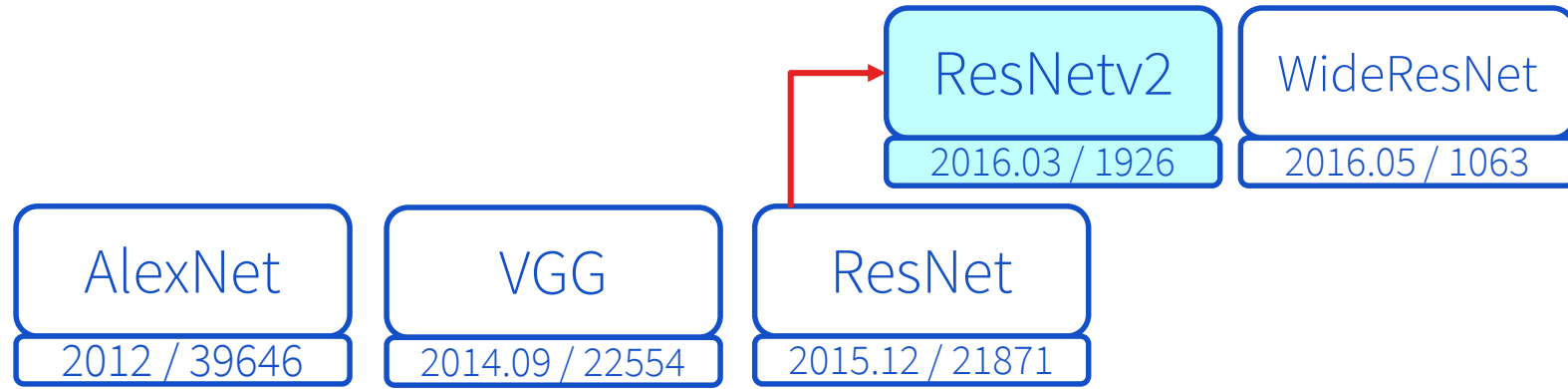
- *Deeper Networks using Skip-Connection*



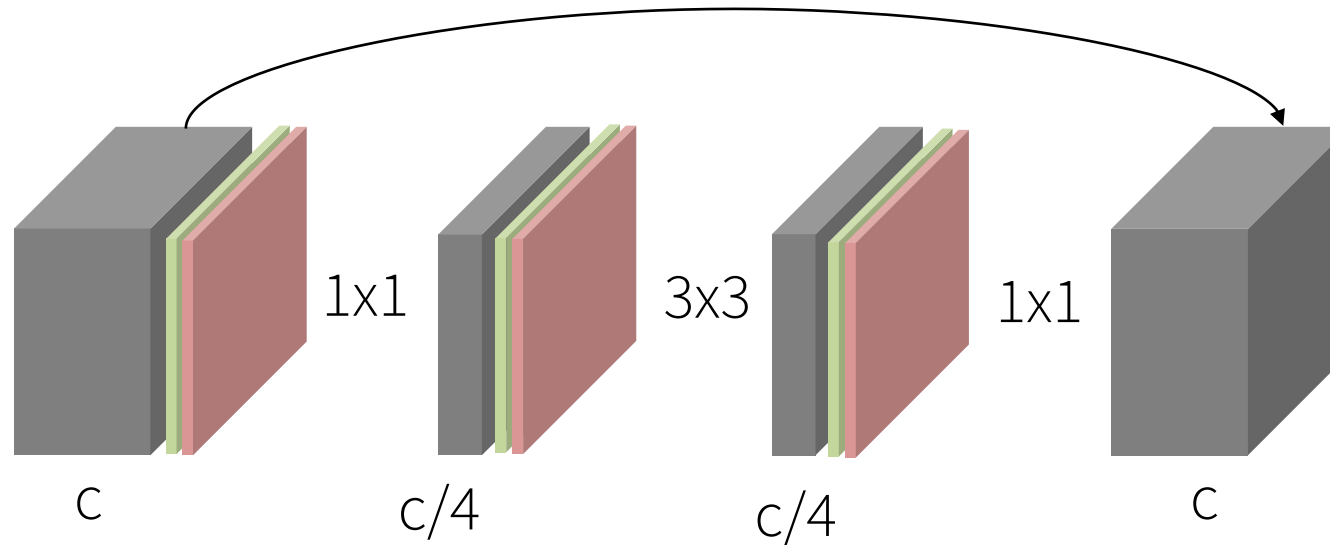
ResNet ResBlock



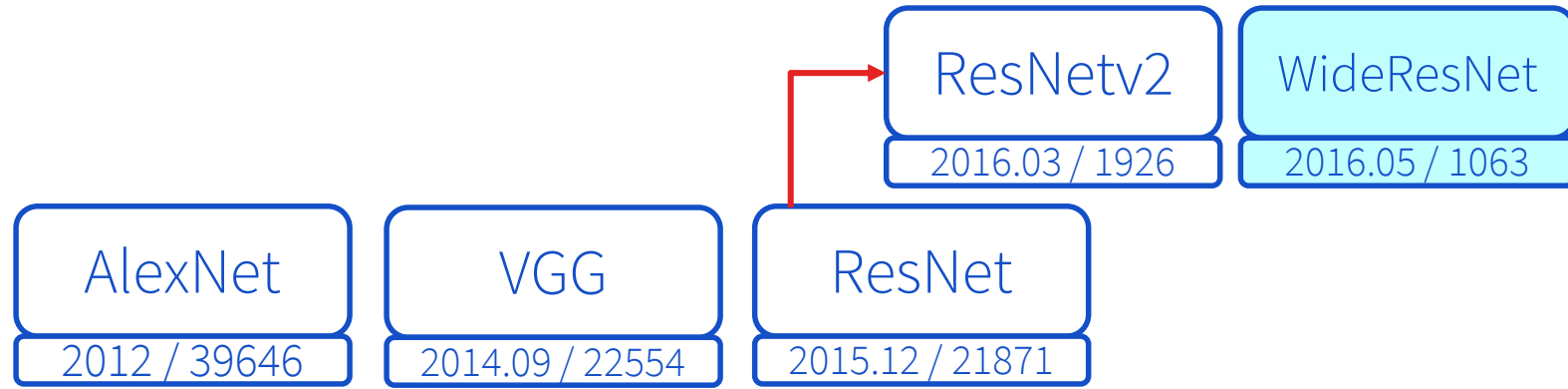
ResNet Variants



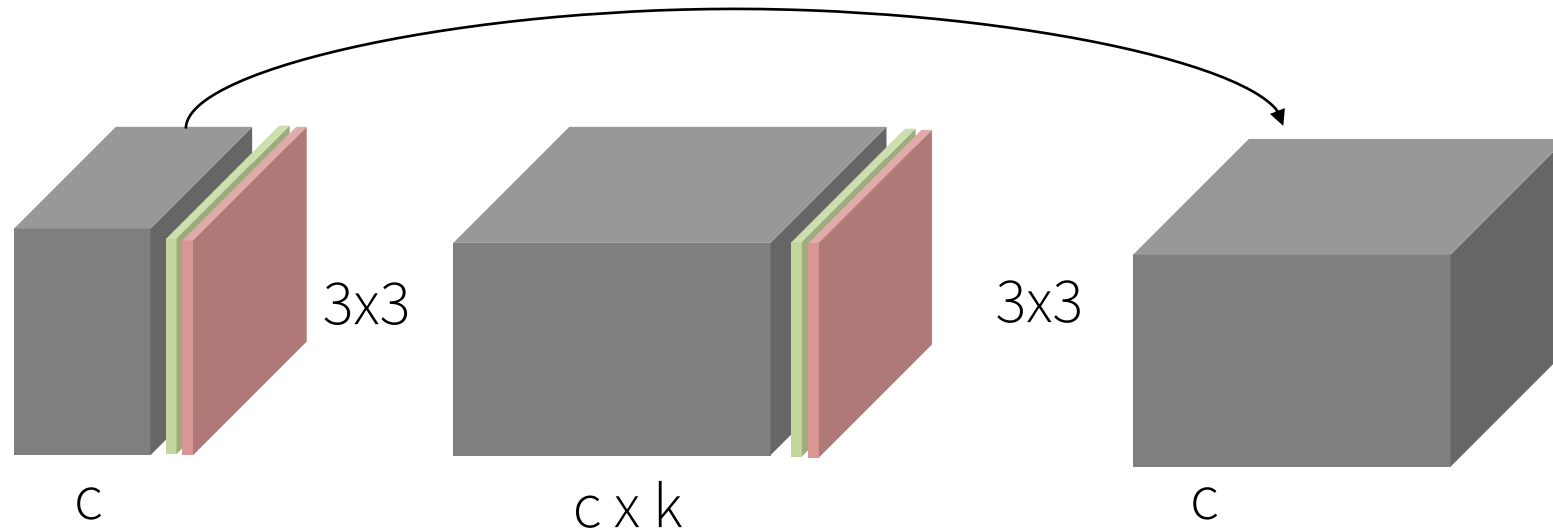
- *Pre-activation ResNet*



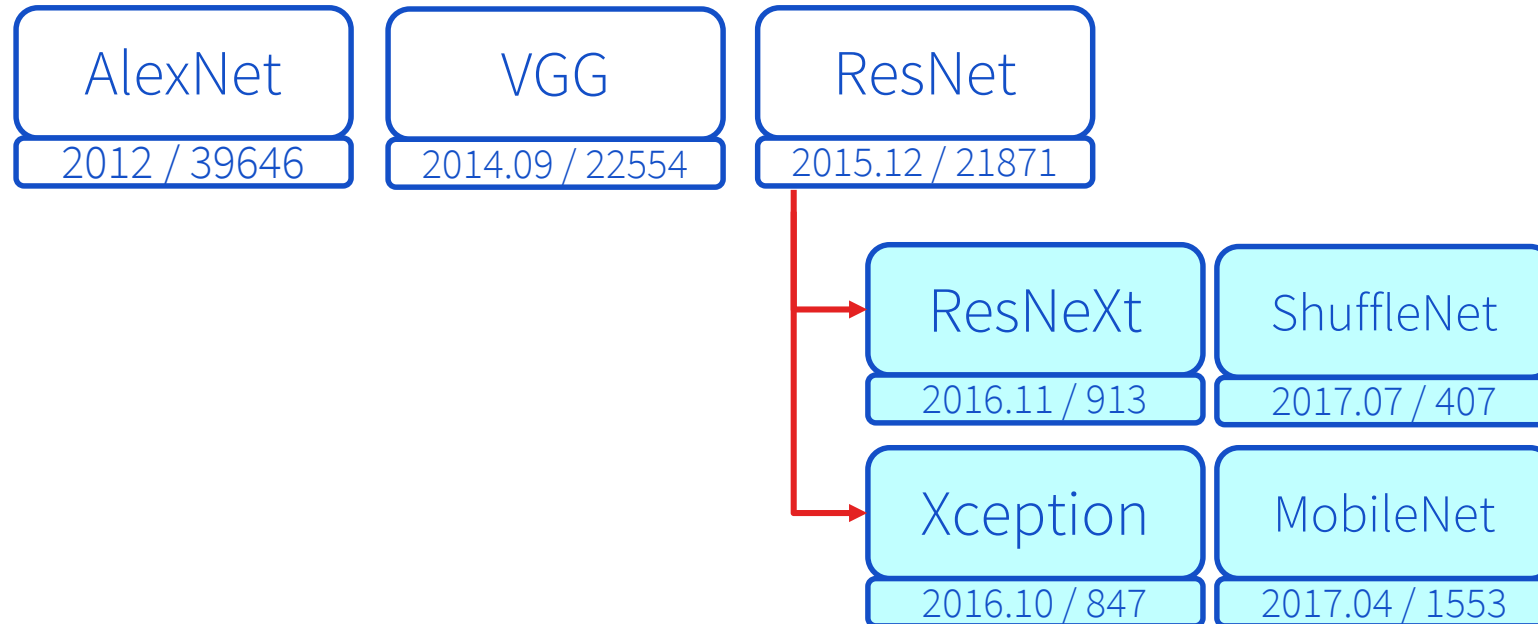
ResNet Variants



- *Wider Channel ResNet*



ResNet with Cardinality

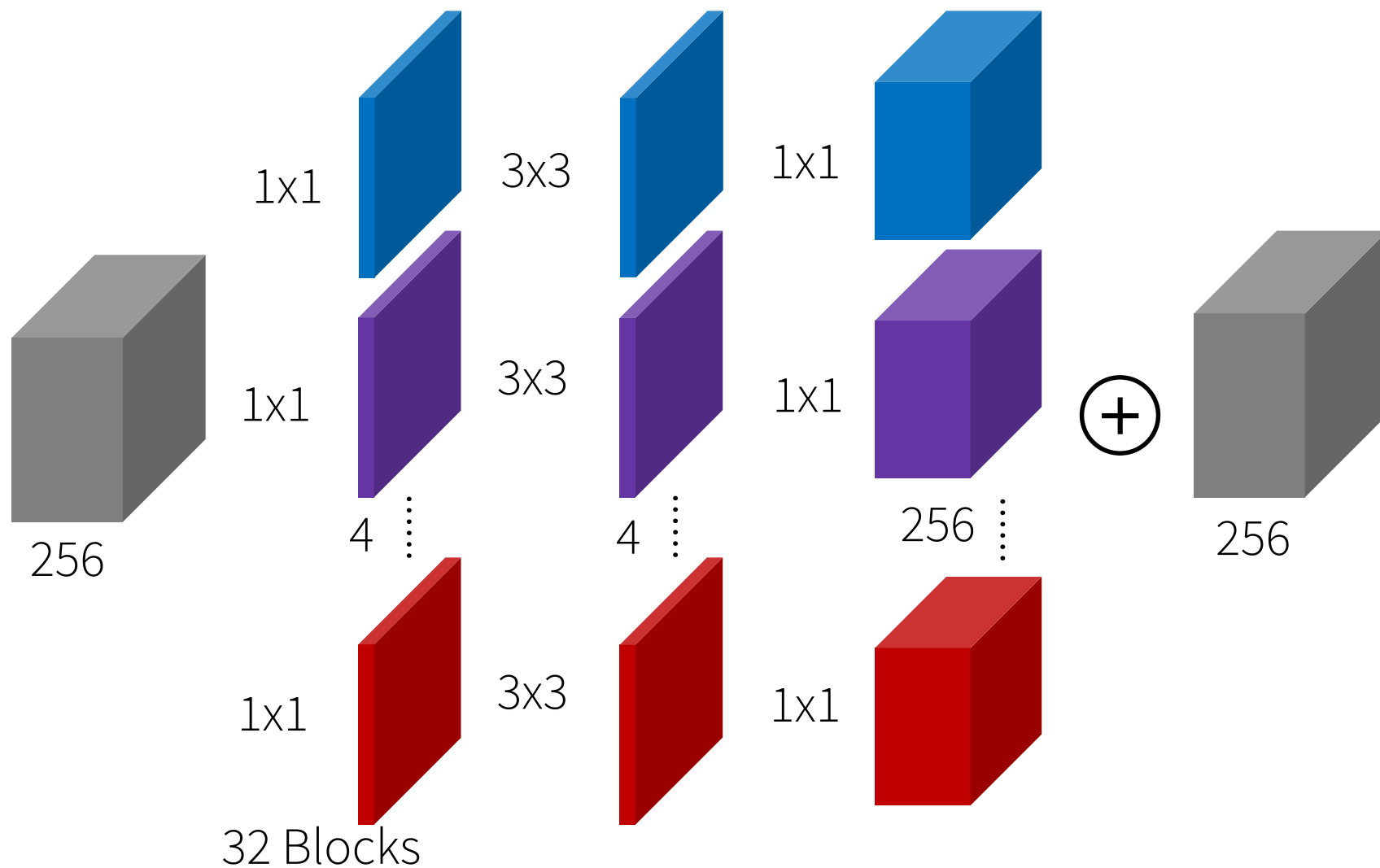


- *Cardinality*
= *Group Conv*

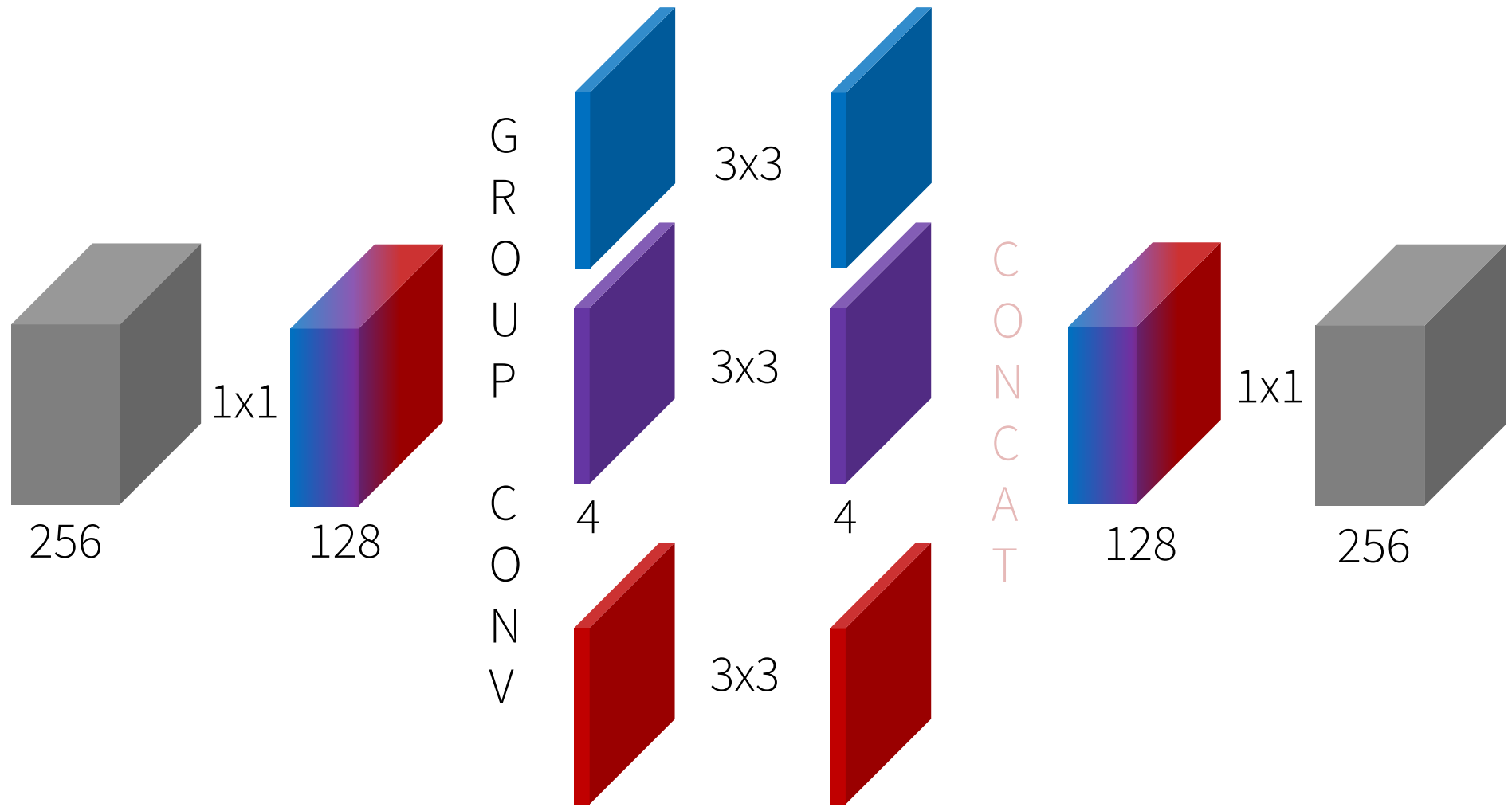
- *Modify Convolution Operators*

- *PR-034: Xception*
- *PR-044: MobileNet*
- *PR-054: SuffleNet / ResNeXt*

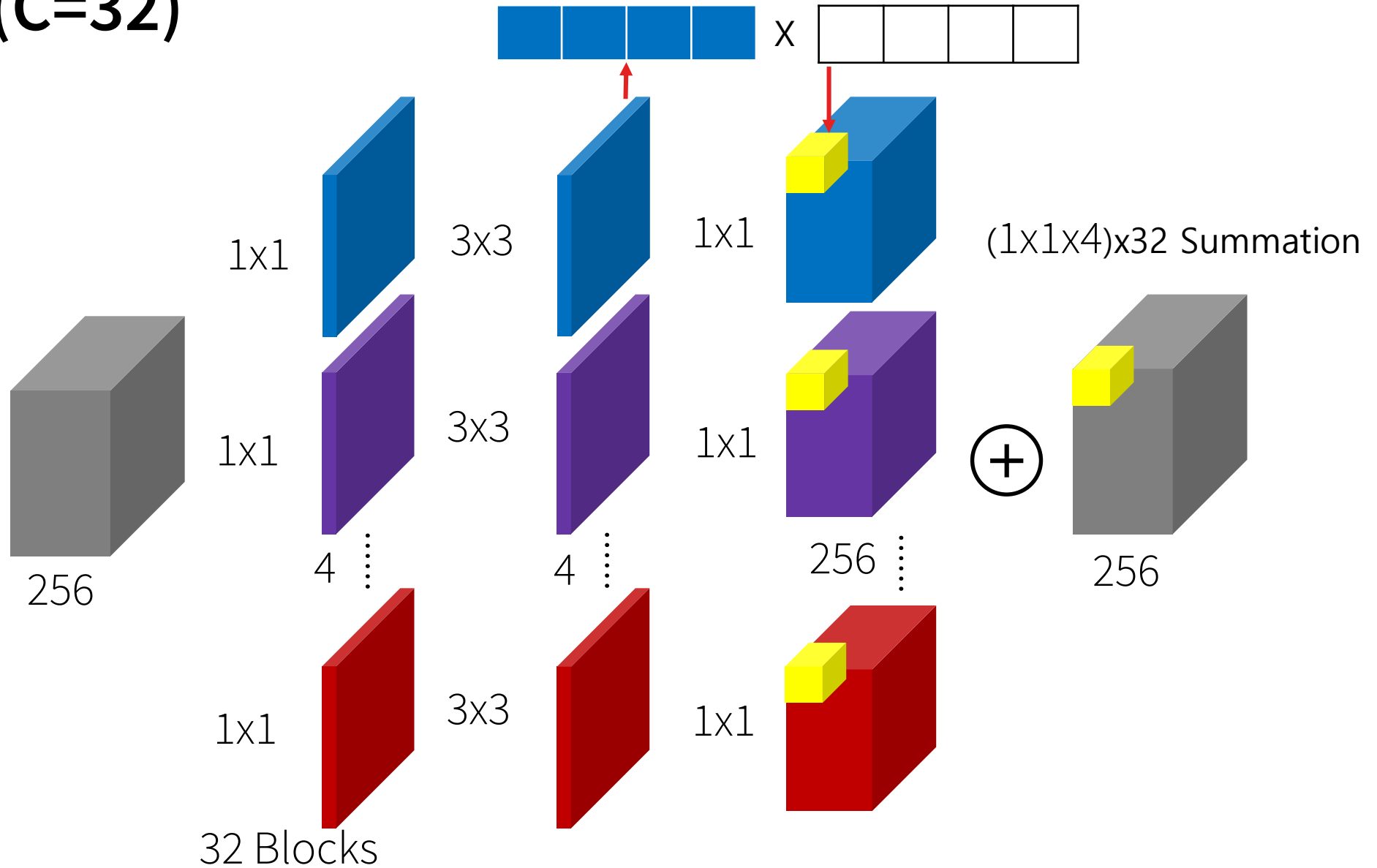
ResNeXt (C=32)



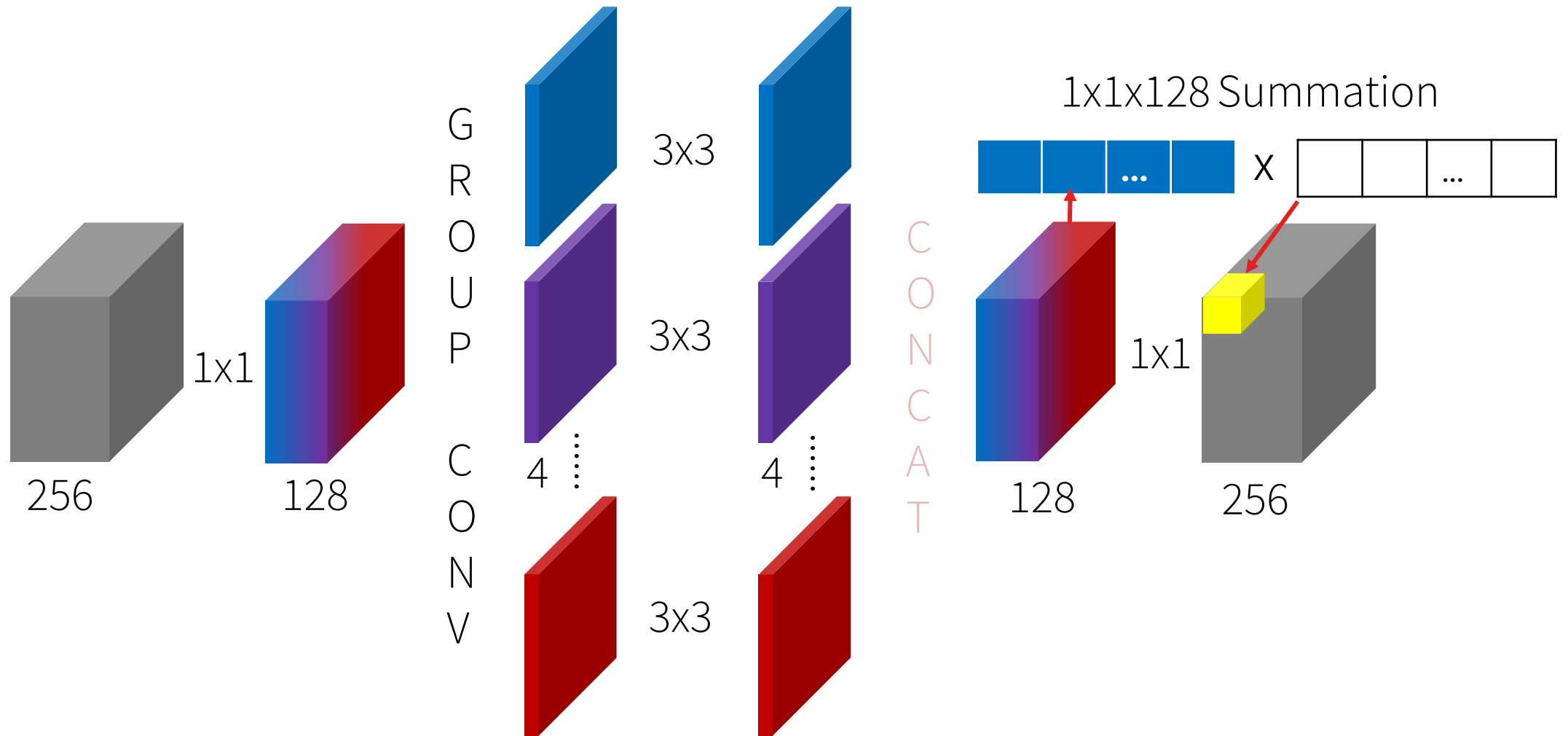
ResNeXt (C=32) with Group Conv.



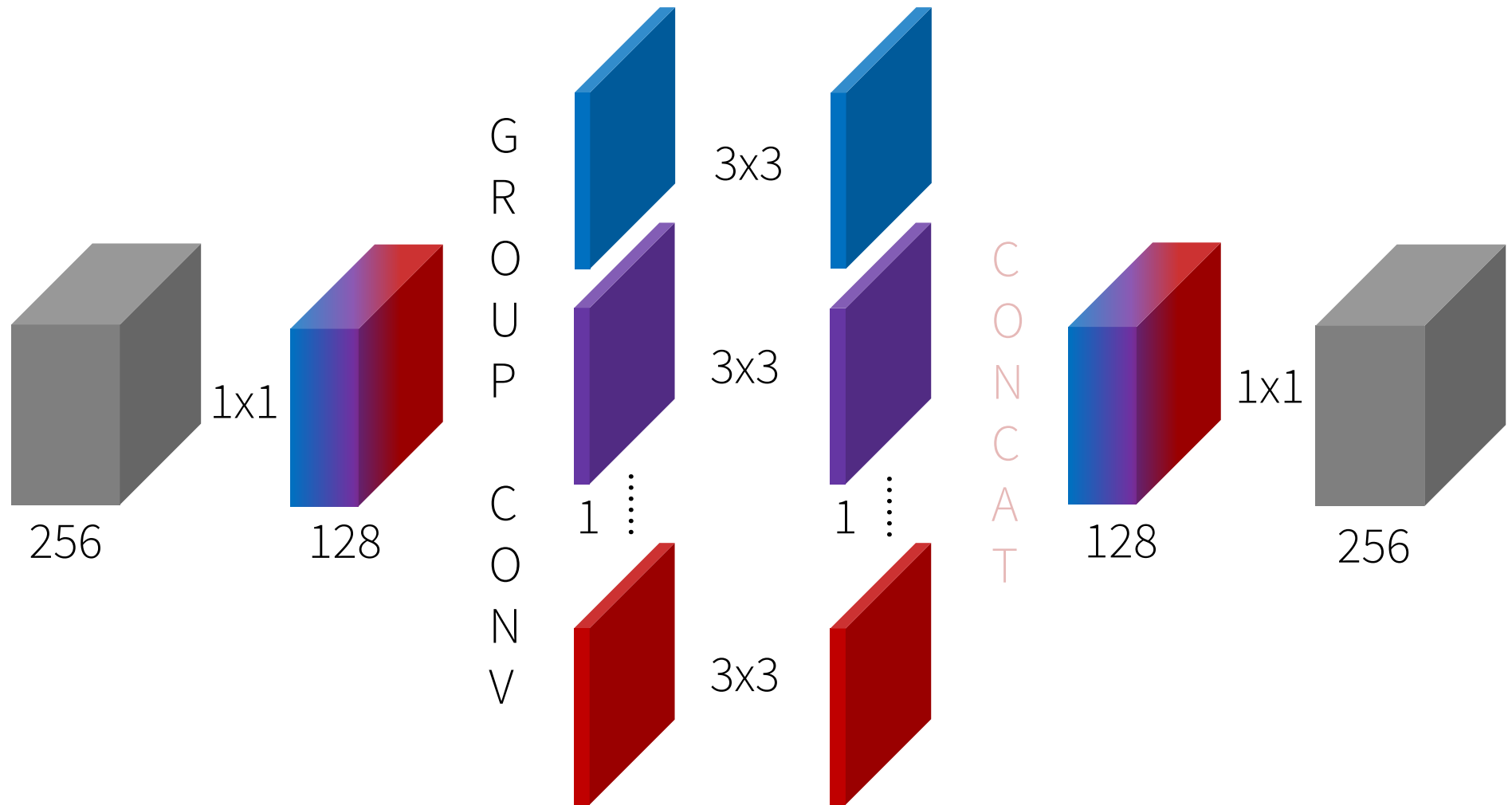
ResNeXt (C=32)



ResNeXt (C=32) with Group Conv.



Xception (G=Channel)



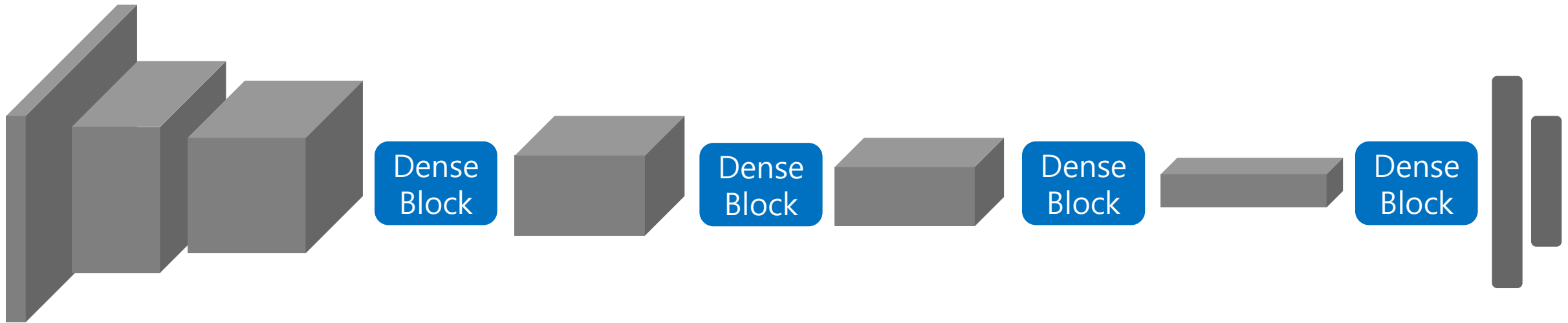
MobileNet / ShuffleNet

- *MobileNet: Lightweight Xception*
 - *Xception: $1 \times 1 \rightarrow 3 \times 3$ Depthwise*
 - *MobileNet: 3×3 Depthwise $\rightarrow 1 \times 1$*
- *ShuffleNet: Lightweight ResNeXt + Channel Shuffle*

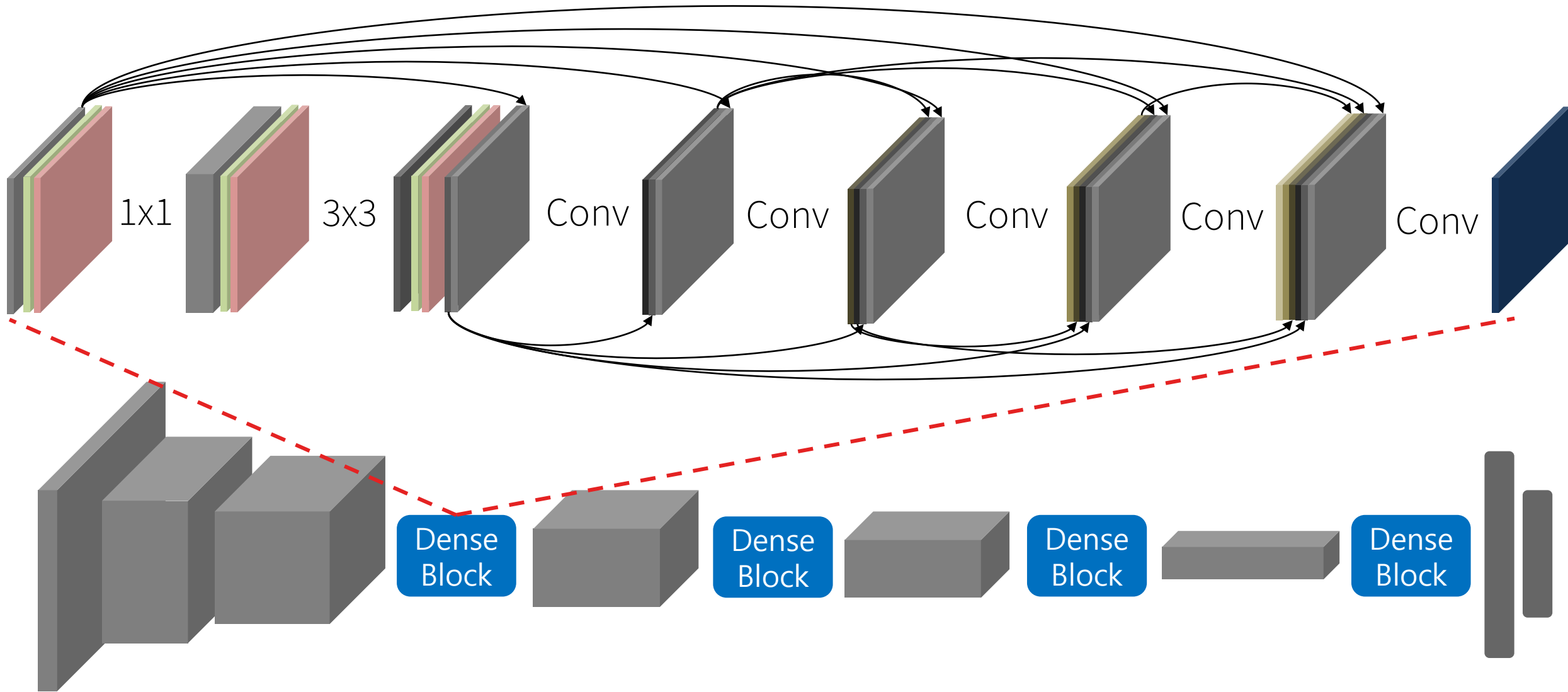
DenseNet



- *DenseNet: Concat Previous Layers: PR-028*



DenseBlock



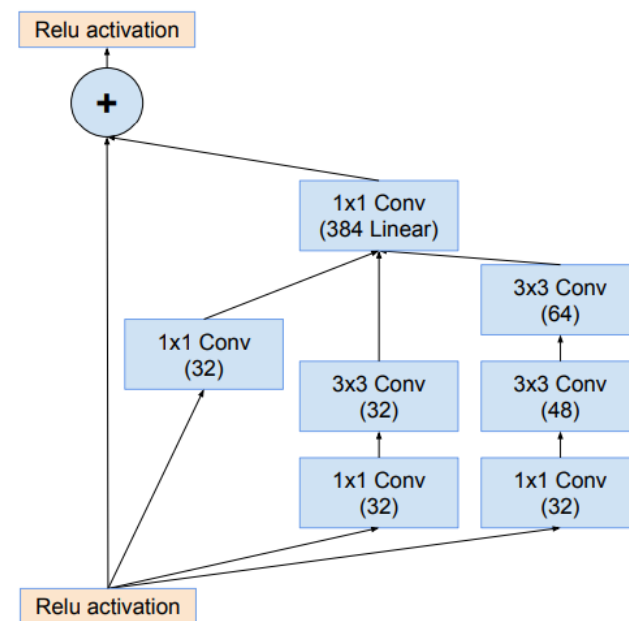
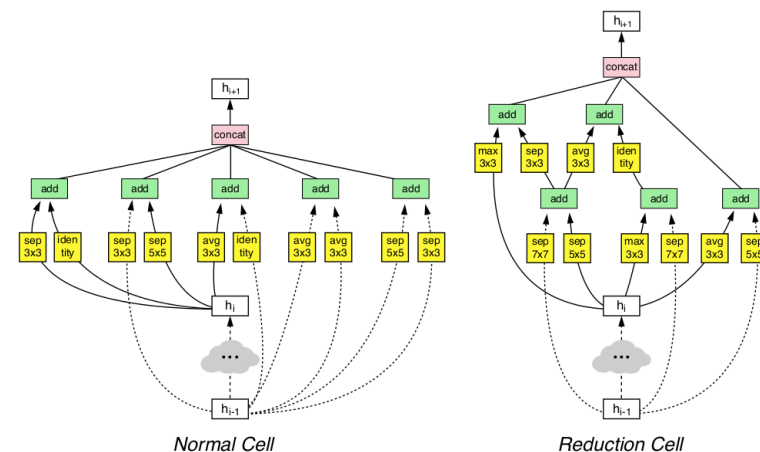
Inception / NASNet

AlexNet	VGG	ResNet
2012 / 39646	2014.09 / 22554	2015.12 / 21871

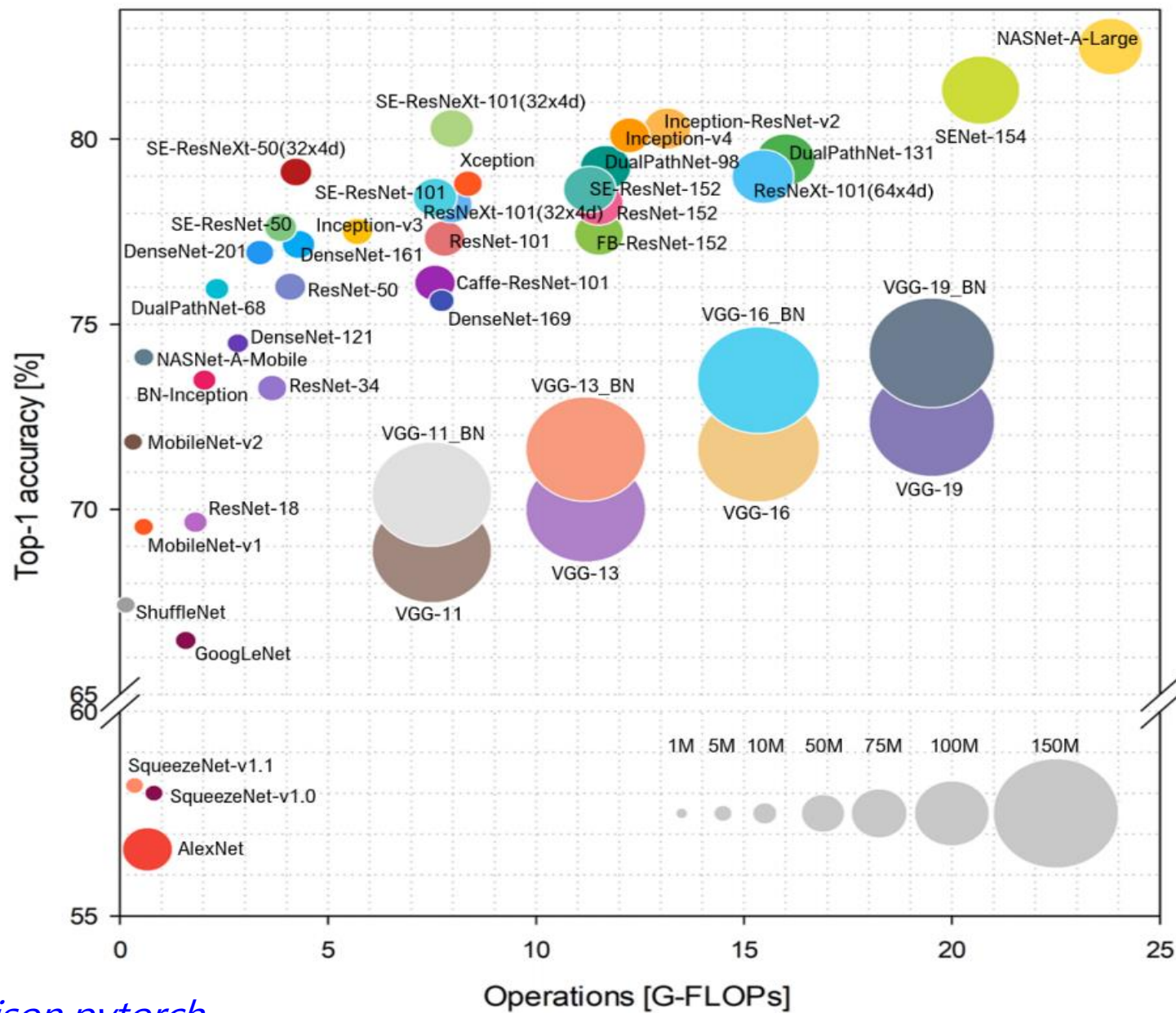
- *Engineered Networks*
- *PR-034: Inception*
- *PR-069: NASNet*

ILSVRC14

GoogleNet	Inception-2,3	Inception-4 Inception-ResNet	NASNet
2014.09 / 13233	2015.12 / 3752	2016.02 / 2152	2017.07. / 417



CNN Performances



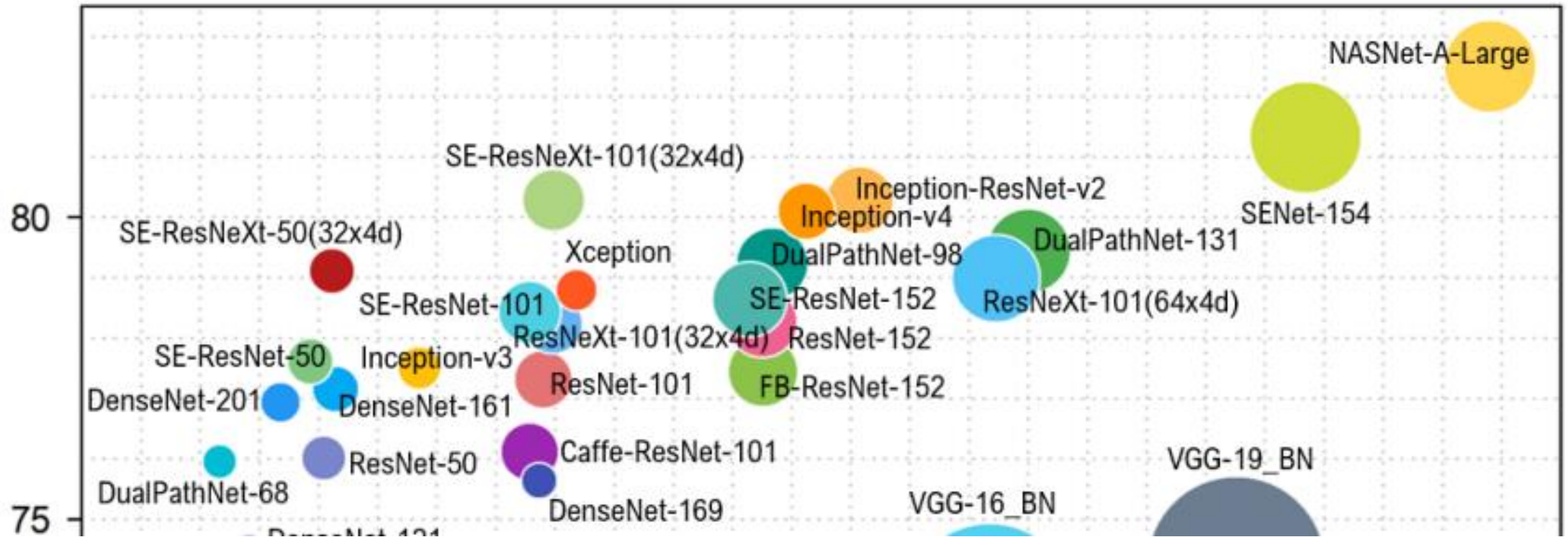
CNN Review

Category	Networks	Pros	Cons
Plain	AlexNet, VGG	Simple Good Transfer	Low Performance
ResNet	ResNet	Simple	
Cardinality	ResNeXt/Xception MobileNet/ShuffleNet	Cost Efficient + Performance	Group Conv
DenseNet	DenseNet	Cost Efficient + Performance	Memory I/O
Engineering	Inception NASNet	SoTA	Complex

CNN Review

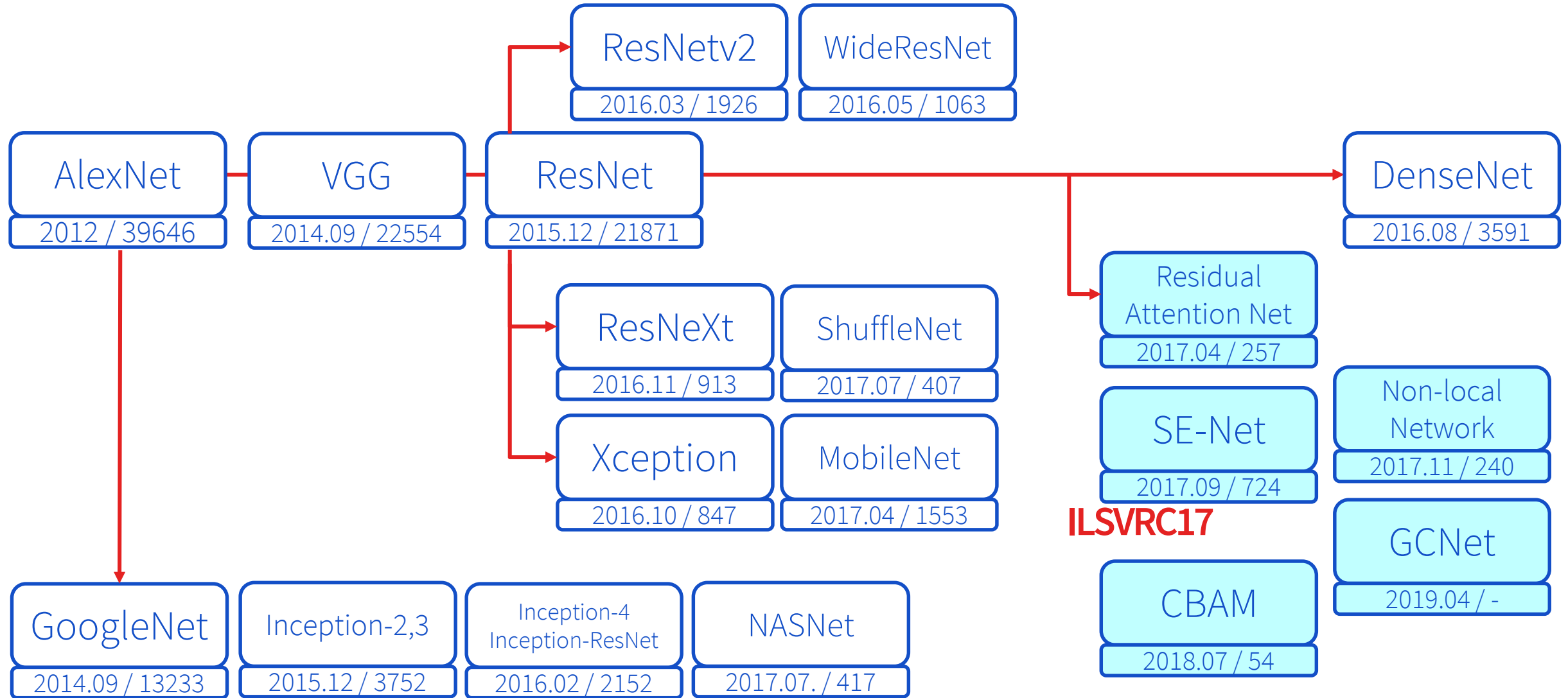
Category	Networks	Pros	Cons
Plain	AlexNet, VGG	Simple Good Transfer	Low Performance
ResNet	ResNet	Simple	
Cardinality	ResNeXt/Xception MobileNet/ShuffleNet	Cost Efficient + Performance	Group Conv
DenseNet	DenseNet	Cost Efficient + Performance	Memory I/O
Engineering	Inception NASNet	SoTA	Complex
Attention Module	SENet, CBAM, GCNet	Simple + Performance	

CNN Review



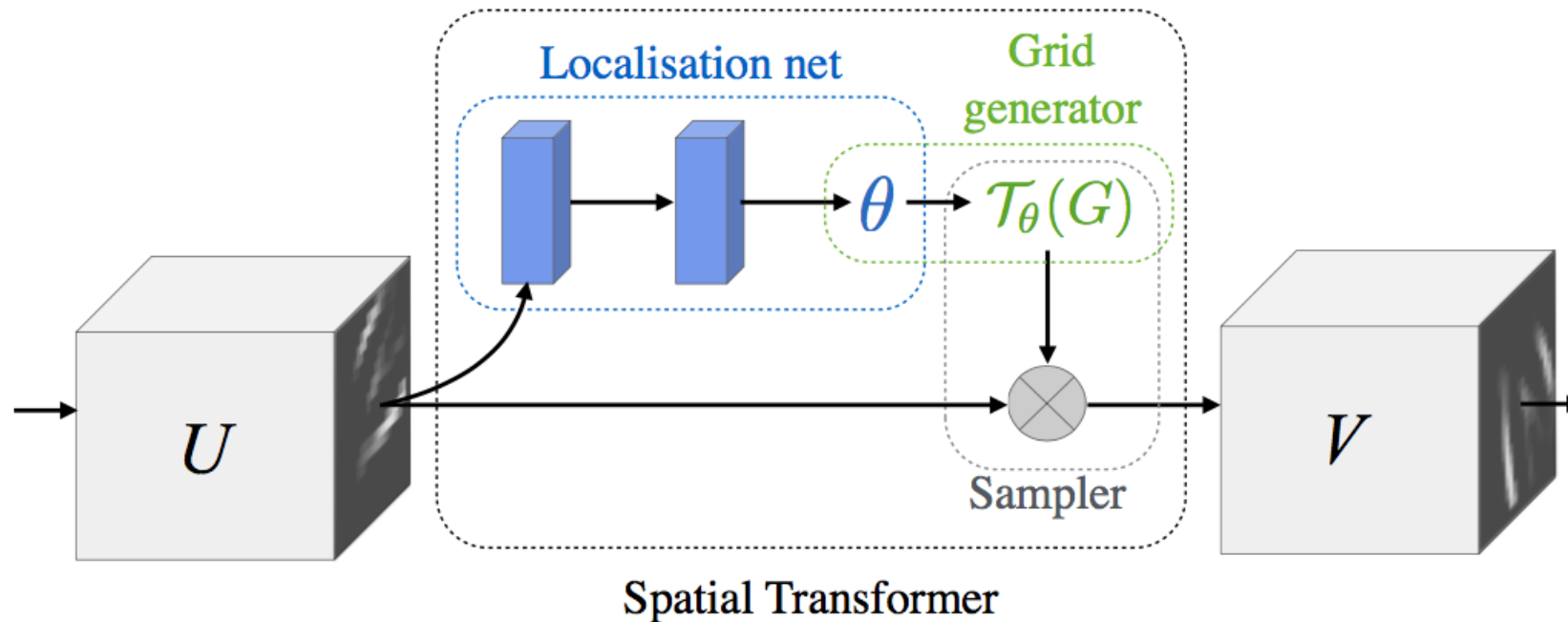
CNN x Attention

CNN Attention-Networks

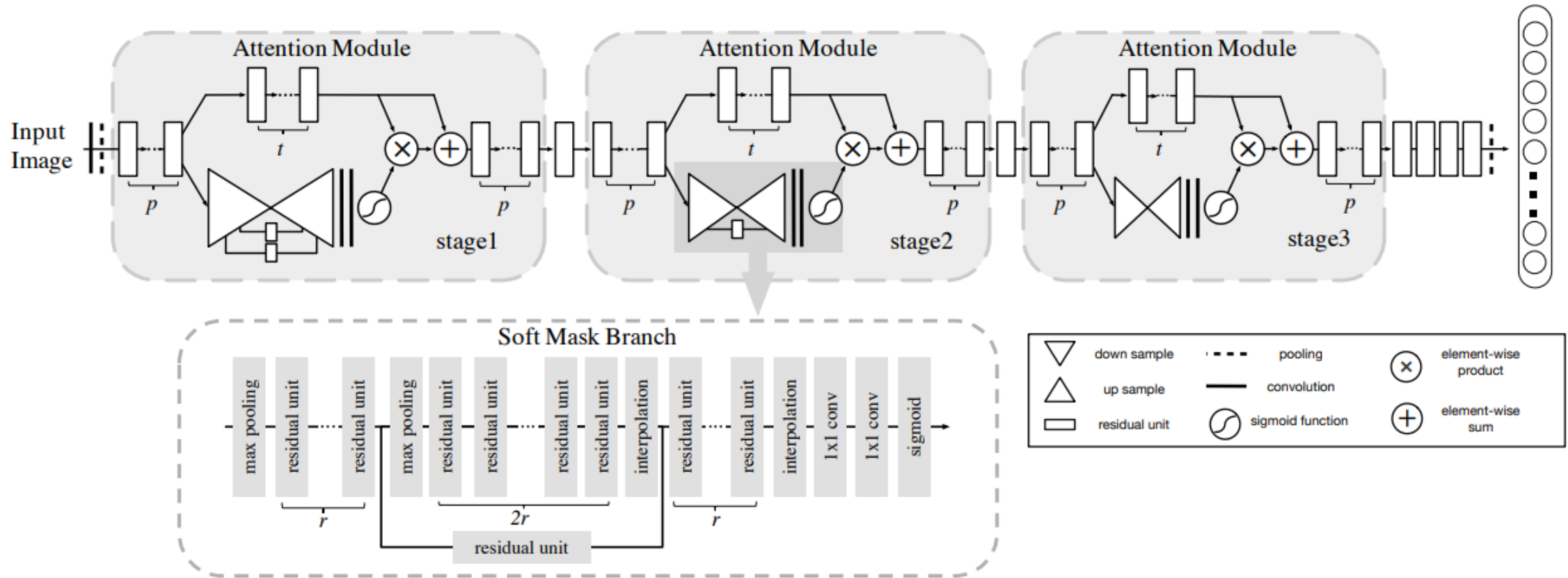


Spatial Transformer Networks (NIPS15, [PR-011](#))

- Recalibration (with Transform)*

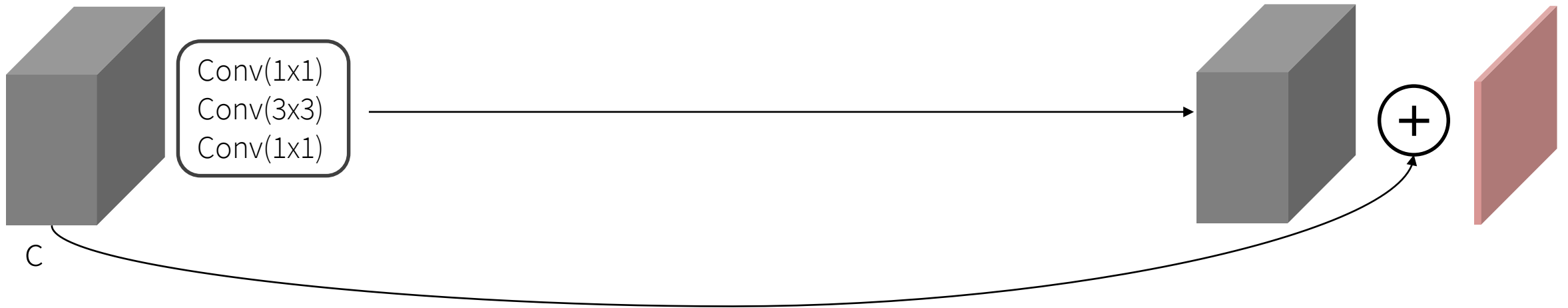


Residual Attention Network (CVPR17)



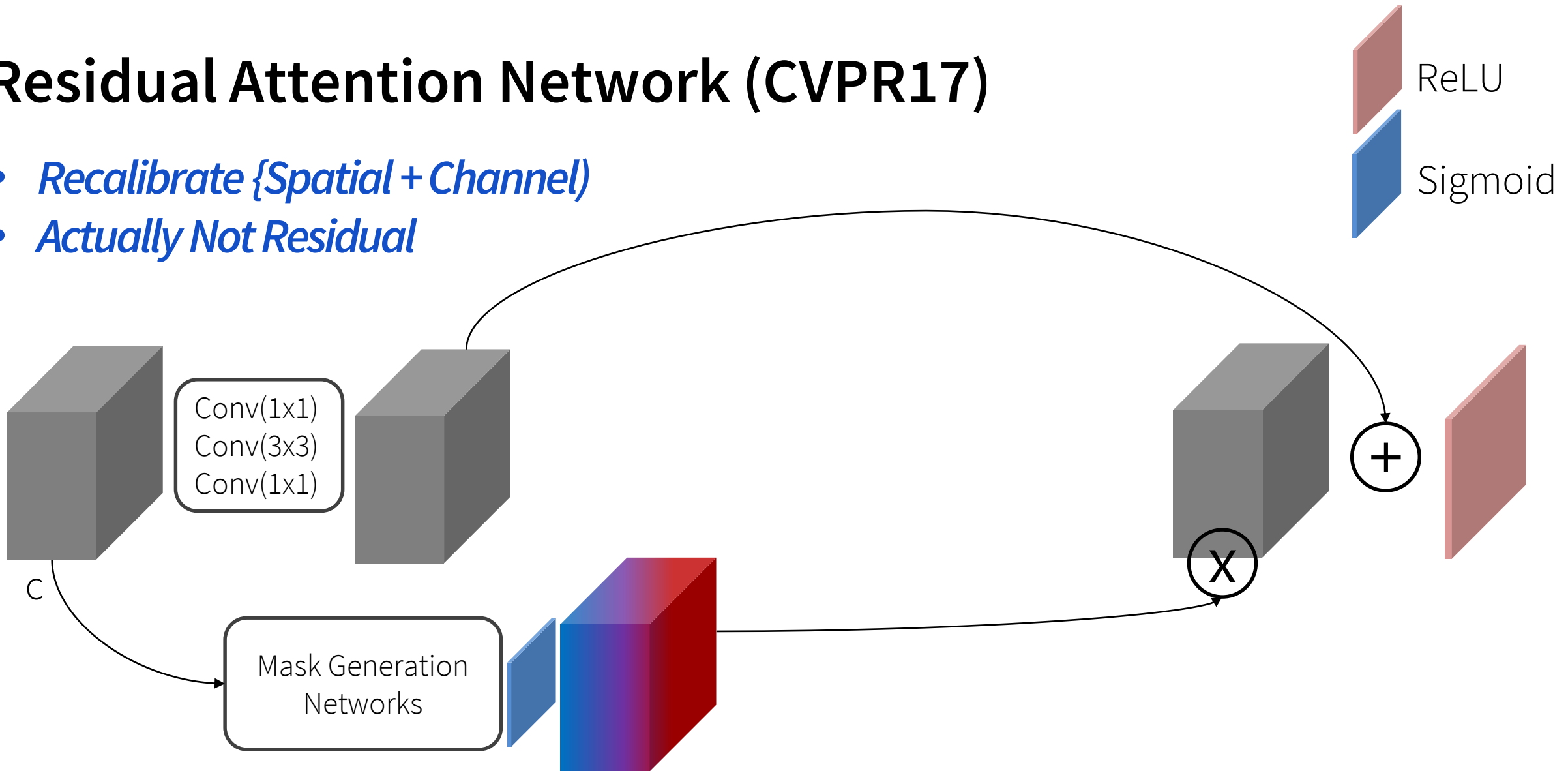
Residual Attention Network (CVPR17)

- *Original ResNet (BottleNeck Block)*



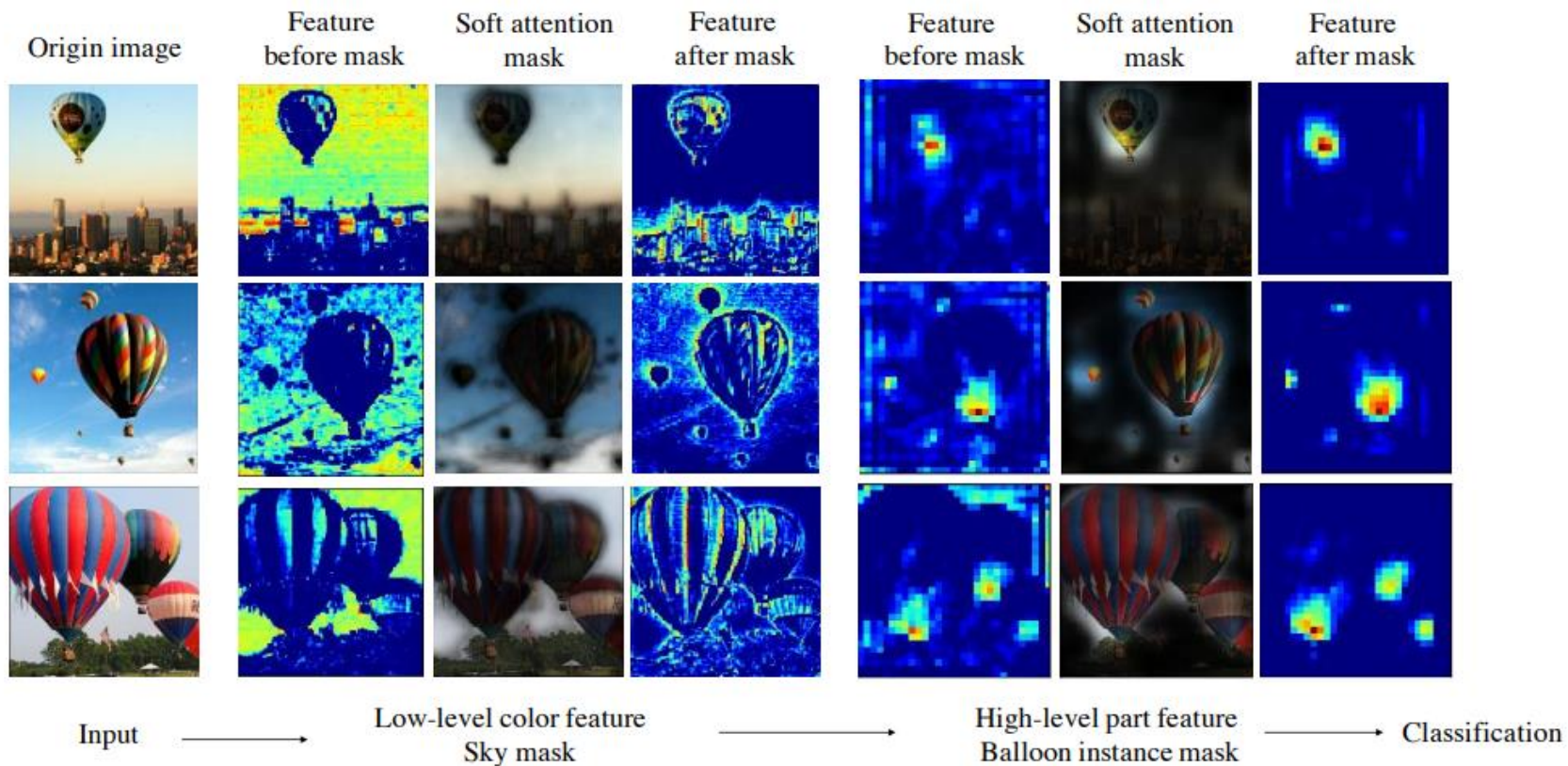
Residual Attention Network (CVPR17)

- *Recalibrate {Spatial + Channel}*
- *Actually Not Residual*



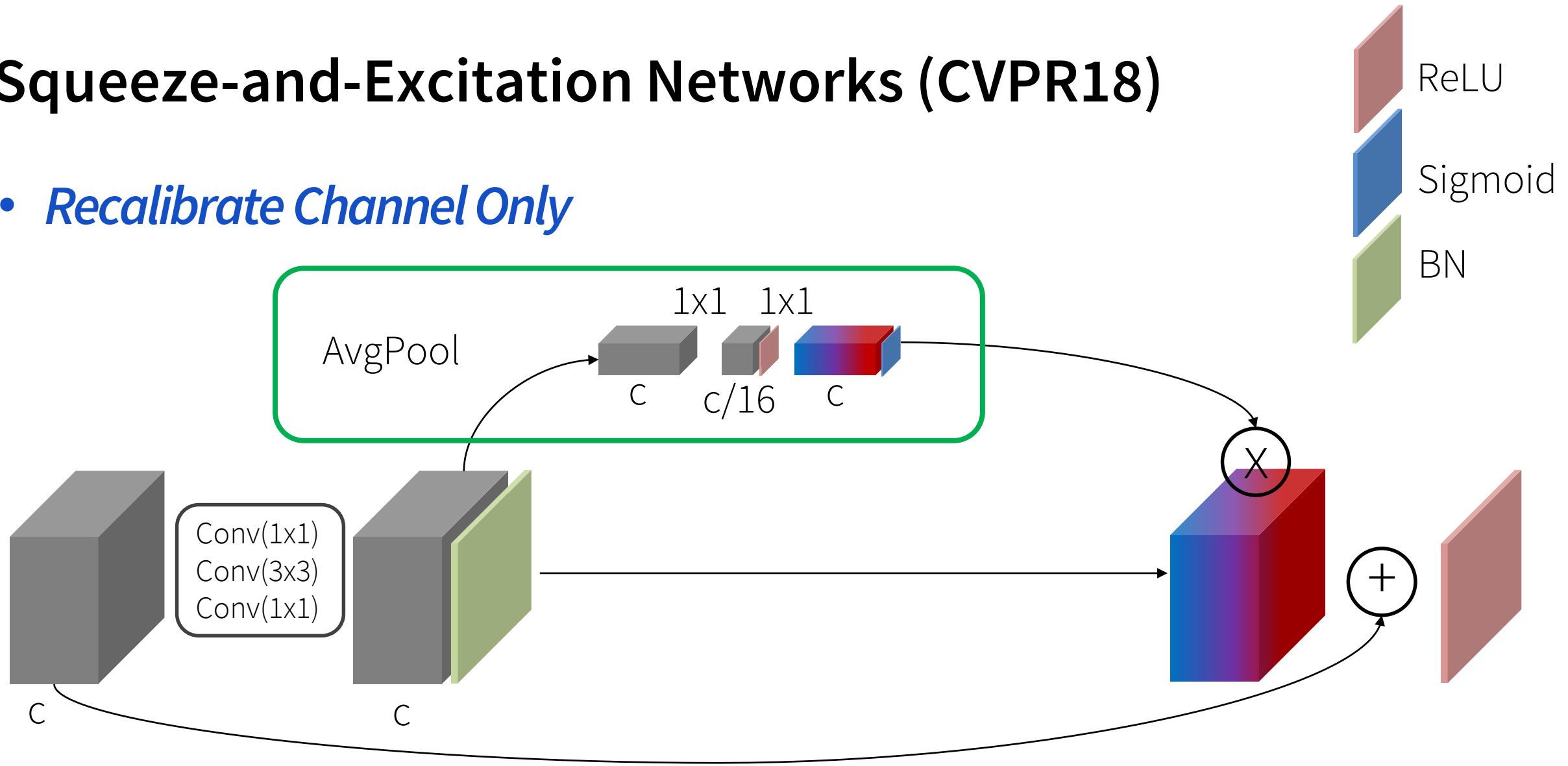
Residual Attention Network (CVPR17)

- Results: Interpretable Features*



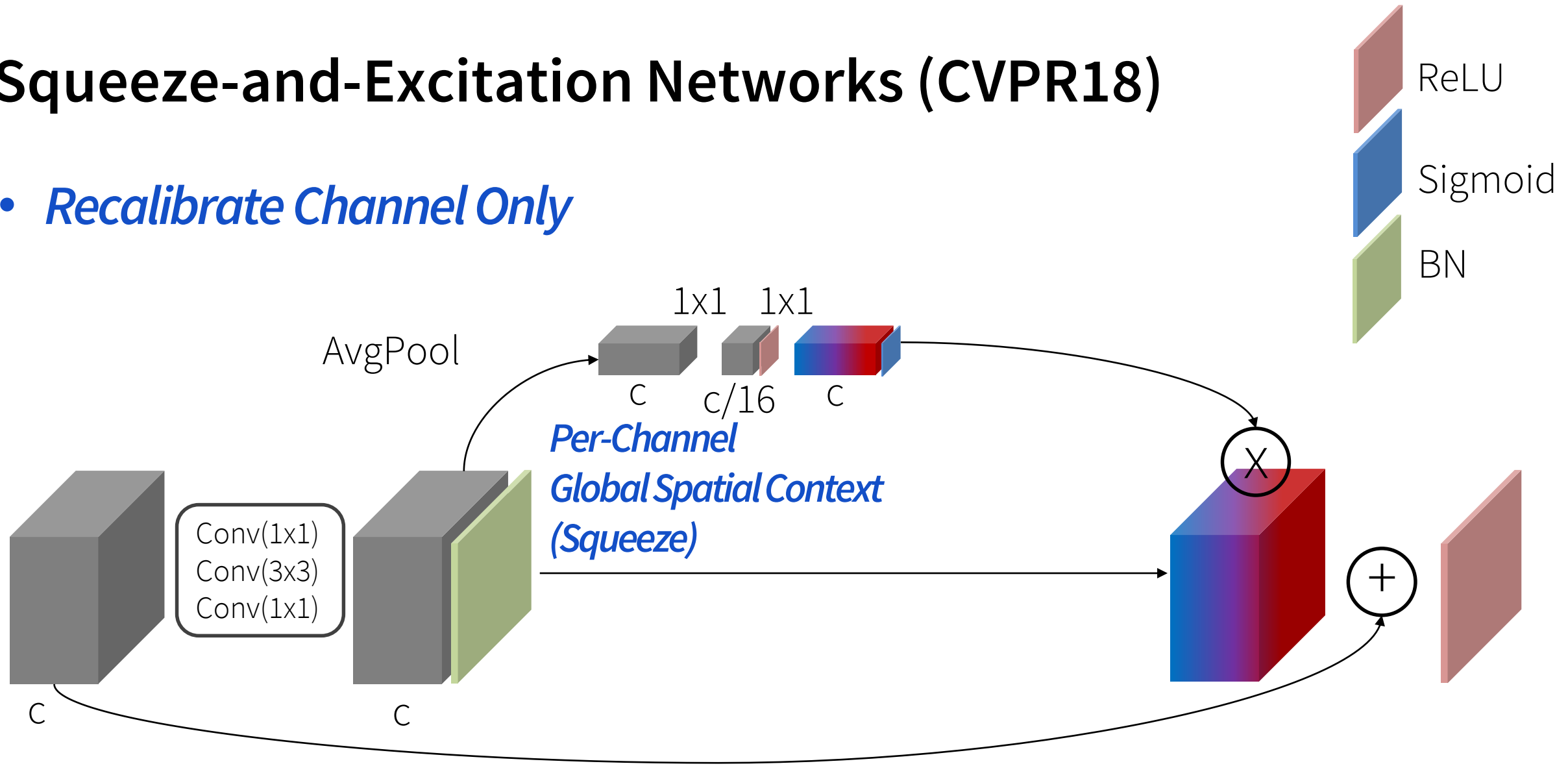
Squeeze-and-Excitation Networks (CVPR18)

- Recalibrate Channel Only*



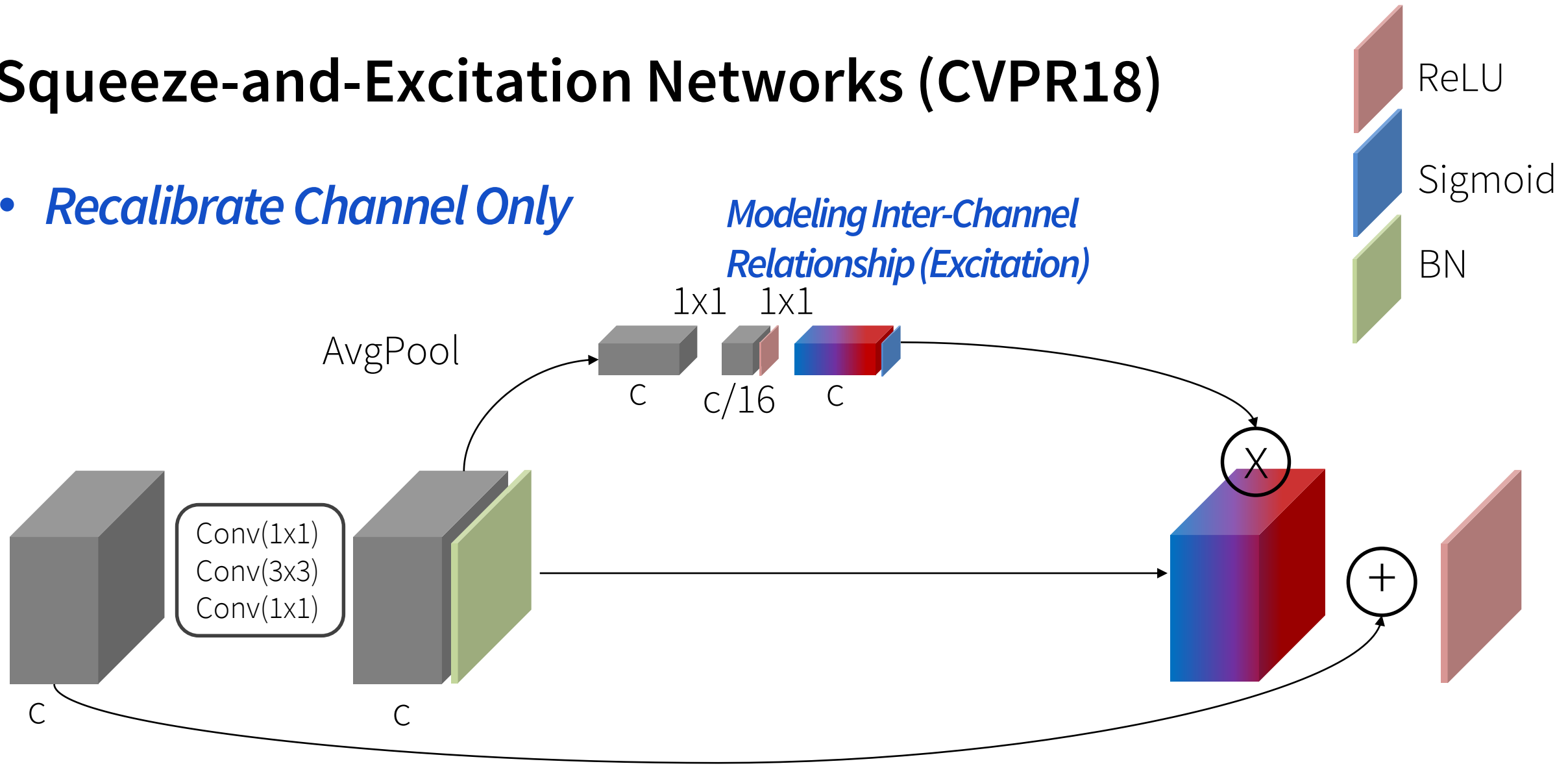
Squeeze-and-Excitation Networks (CVPR18)

- Recalibrate Channel Only*



Squeeze-and-Excitation Networks (CVPR18)

- Recalibrate Channel Only*



Squeeze-and-Excitation Networks (CVPR18)

	original		re-implementation			SENet		
	top-1 err.	top-5 err.	top-1 err.	top-5 err.	GFLOPs	top-1 err.	top-5 err.	GFLOPs
ResNet-50 [13]	24.7	7.8	24.80	7.48	3.86	23.29 _(1.51)	6.62 _(0.86)	3.87
ResNet-101 [13]	23.6	7.1	23.17	6.52	7.58	22.38 _(0.79)	6.07 _(0.45)	7.60
ResNet-152 [13]	23.0	6.7	22.42	6.34	11.30	21.57 _(0.85)	5.73 _(0.61)	11.32
ResNeXt-50 [19]	22.2	-	22.11	5.90	4.24	21.10 _(1.01)	5.49 _(0.41)	4.25
ResNeXt-101 [19]	21.2	5.6	21.18	5.57	7.99	20.70 _(0.48)	5.01 _(0.56)	8.00
VGG-16 [11]	-	-	27.02	8.81	15.47	25.22 _(1.80)	7.70 _(1.11)	15.48
BN-Inception [6]	25.2	7.82	25.38	7.89	2.03	24.23 _(1.15)	7.14 _(0.75)	2.04
Inception-ResNet-v2 [21]	19.9 [†]	4.9 [†]	20.37	5.21	11.75	19.80 _(0.57)	4.79 _(0.42)	11.76

Squeeze-and-Excitation Networks (CVPR18)

Ratio r	top-1 err.	top-5 err.	Params
2	22.29	6.00	45.7M
4	22.25	6.09	35.7M
8	22.26	5.99	30.7M
16	22.28	6.03	28.1M
32	22.72	6.20	26.9M
original	23.30	6.55	25.6M

Excitation	top-1 err.	top-5 err.
ReLU	23.47	6.98
Tanh	23.00	6.38
Sigmoid	22.28	6.03

Squeeze	top-1 err.	top-5 err.
Max	22.57	6.09
Avg	22.28	6.03

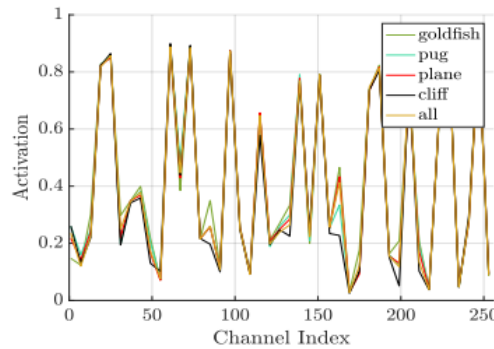
Design	top-1 err.	top-5 err.
SE	22.28	6.03
SE-PRE	22.23	6.00
SE-POST	22.78	6.35
SE-Identity	22.20	6.15

Stage	top-1 err.	top-5 err.	GFLOPs	Params
ResNet-50	23.30	6.55	3.86	25.6M
SE_Stage_2	23.03	6.48	3.86	25.6M
SE_Stage_3	23.04	6.32	3.86	25.7M
SE_Stage_4	22.68	6.22	3.86	26.4M
SE_All	22.28	6.03	3.87	28.1M

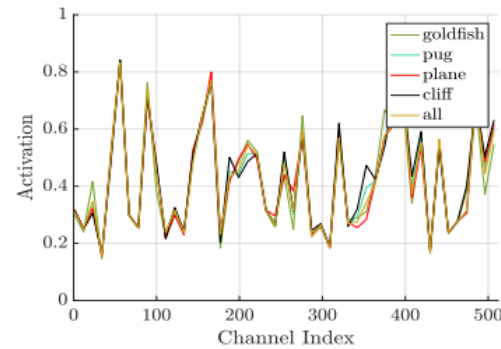
Design	top-1 err.	top-5 err.	GFLOPs	Params
SE	22.28	6.03	3.87	28.1M
SE_3×3	22.48	6.02	3.86	25.8M

Squeeze-and-Excitation Networks (CVPR18)

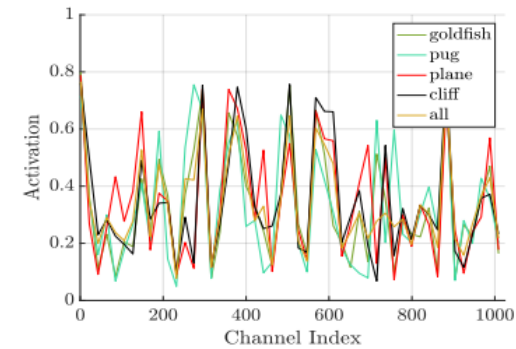
- Channel Recalibration Stats.*



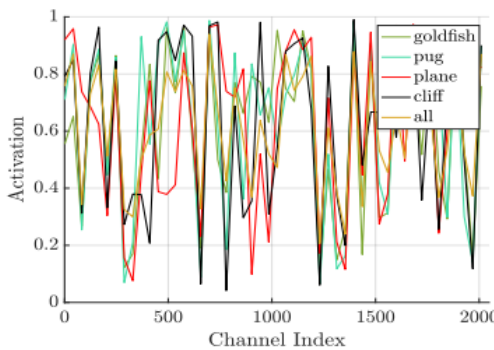
(a) SE_2_3



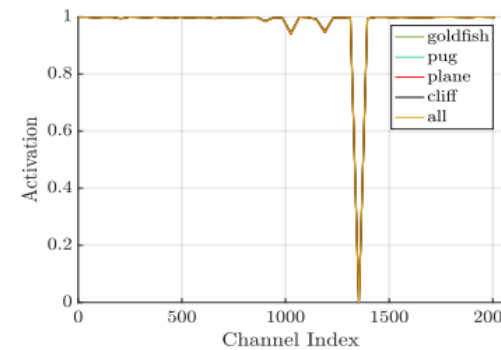
(b) SE_3_4



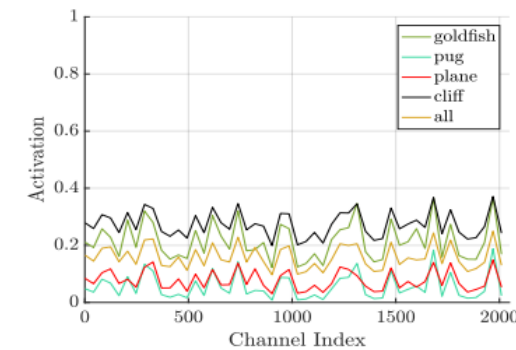
(c) SE_4_6



(d) SE_5_1

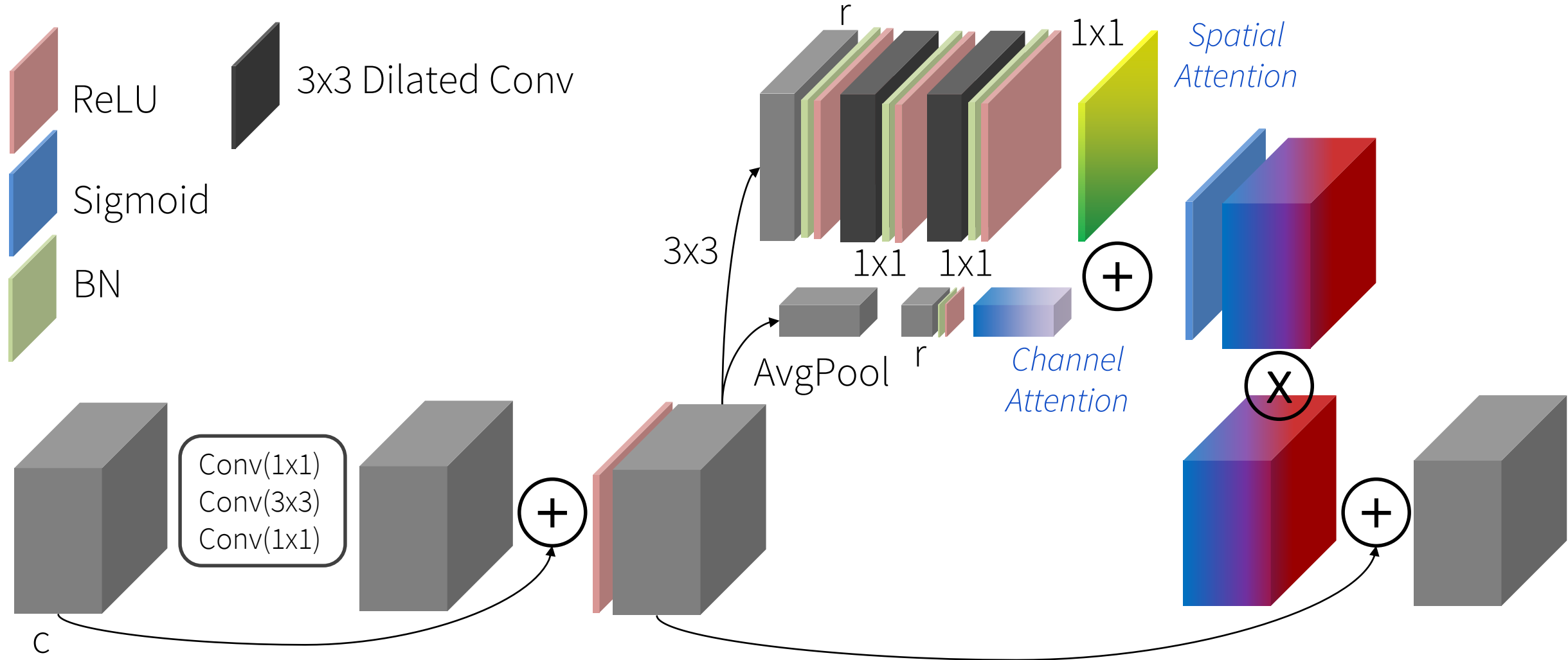


(e) SE_5_2

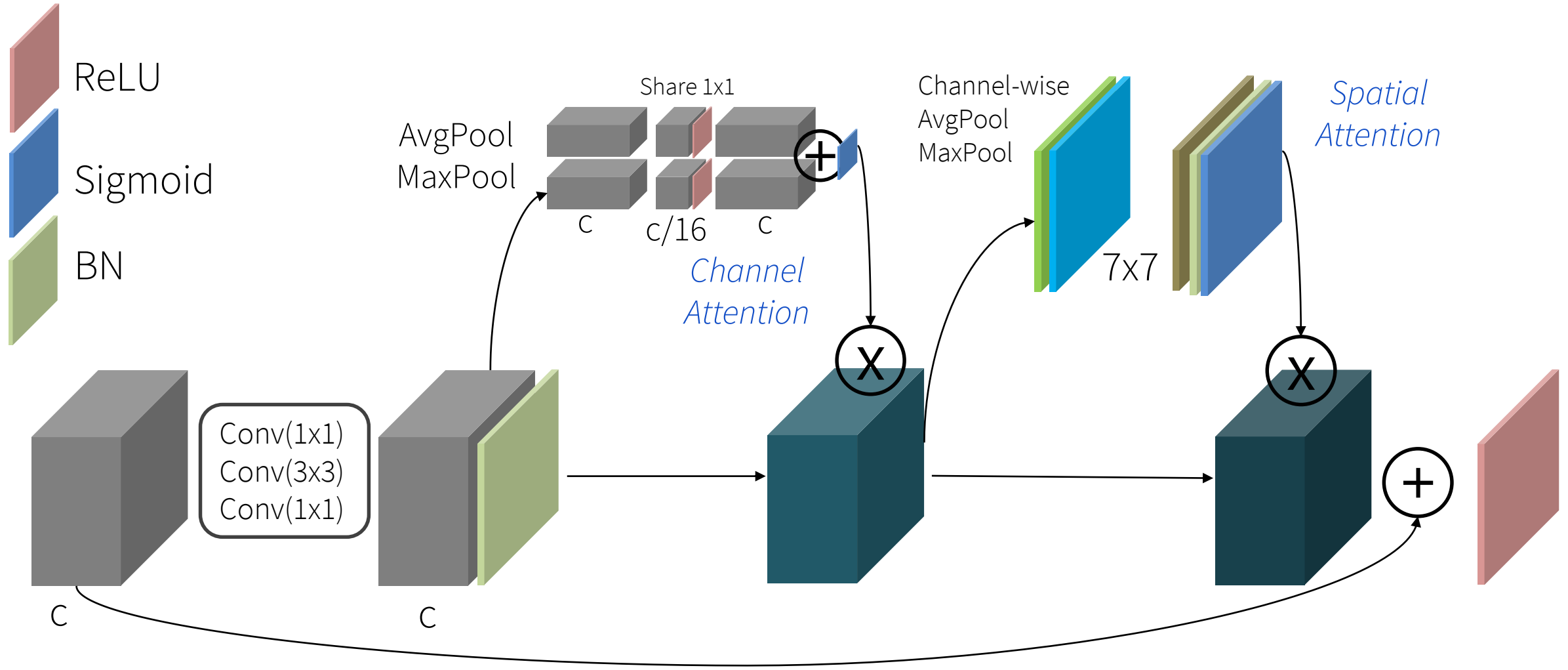


(f) SE_5_3

Bottleneck Attention Networks (BMVC18)



Convolutional Block Attention Networks (ECCV18)



BAM / CBAM Results

Architecture	Parameters	GFLOPs	Top-1(%)	Top-5(%)
ResNet18 [15]	11.69M	1.81	29.60	10.55
ResNet18 [15] + BAM	11.71M _(+0.02)	1.82 _(+0.01)	28.88	10.01
ResNet50 [15]	25.56M	3.86	24.56	7.50
ResNet50 [15] + BAM	25.92M _(+0.36)	3.94 _(+0.08)	24.02	7.18
ResNet101 [15]	44.55M	7.57	23.38	6.88
ResNet101 [15] + BAM	44.91M _(+0.36)	7.65 _(+0.08)	22.44	6.29
WideResNet18 [47] (widen=1.5)	25.88M	3.87	26.85	8.88
WideResNet18 [47] (widen=1.5) + BAM	25.93M _(+0.05)	3.88 _(+0.01)	26.67	8.69
WideResNet18 [47] (widen=2.0)	45.62M	6.70	25.63	8.20
WideResNet18 [47] (widen=2.0) + BAM	45.71M _(+0.09)	6.72 _(+0.02)	25.00	7.81
ResNeXt50 [43] (32x4d)	25.03M	3.77	22.85	6.48
ResNeXt50 [43] (32x4d) + BAM	25.39M _(+0.36)	3.85 _(+0.08)	22.56	6.40
MobileNet[18]	4.23M	0.569	31.39	11.51
MobileNet[18] + BAM	4.32M _(+0.09)	0.589 _(+0.02)	30.58	10.90
MobileNet[18] $\alpha = 0.7$	2.30M	0.283	34.86	13.69
MobileNet[18] $\alpha = 0.7$ + BAM	2.34M _(+0.04)	0.292 _(+0.009)	33.09	12.69
MobileNet[18] $\rho = 192/224$	4.23M	0.439	32.89	12.33
MobileNet[18] $\rho = 192/224$ + BAM	4.32M _(+0.09)	0.456 _(+0.017)	31.56	11.60
SqueezeNet v1.1 [22]	1.24M	0.290	43.09	20.48
SqueezeNet v1.1 [22] + BAM	1.26M _(+0.02)	0.304 _(+0.014)	41.83	19.58

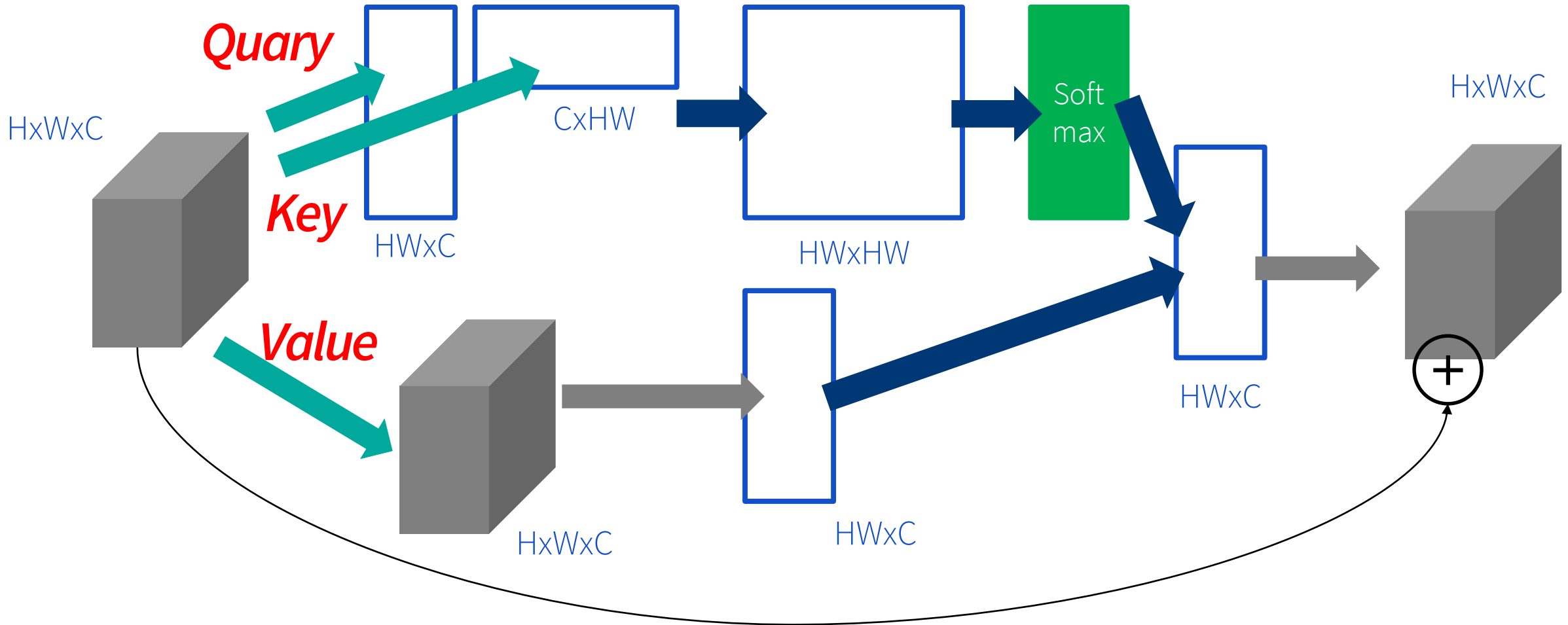
Architecture	Param.	GFLOPs	Top-1 Error (%)	Top-5 Error (%)
ResNet18 [5]	11.69M	1.814	29.60	10.55
ResNet18 [5] + SE [28]	11.78M	1.814	29.41	10.22
ResNet18 [5] + CBAM	11.78M	1.815	29.27	10.09
ResNet34 [5]	21.80M	3.664	26.69	8.60
ResNet34 [5] + SE [28]	21.96M	3.664	26.13	8.35
ResNet34 [5] + CBAM	21.96M	3.665	25.99	8.24
ResNet50 [5]	25.56M	3.858	24.56	7.50
ResNet50 [5] + SE [28]	28.09M	3.860	23.14	6.70
ResNet50 [5] + CBAM	28.09M	3.864	22.66	6.31
ResNet101 [5]	44.55M	7.570	23.38	6.88
ResNet101 [5] + SE [28]	49.33M	7.575	22.35	6.19
ResNet101 [5] + CBAM	49.33M	7.581	21.51	5.69
WideResNet18 [6] (widen=1.5)	25.88M	3.866	26.85	8.88
WideResNet18 [6] (widen=1.5) + SE [28]	26.07M	3.867	26.21	8.47
WideResNet18 [6] (widen=1.5) + CBAM	26.08M	3.868	26.10	8.43
WideResNet18 [6] (widen=2.0)	45.62M	6.696	25.63	8.20
WideResNet18 [6] (widen=2.0) + SE [28]	45.97M	6.696	24.93	7.65
WideResNet18 [6] (widen=2.0) + CBAM	45.97M	6.697	24.84	7.63
ResNeXt50 [7] (32x4d)	25.03M	3.768	22.85	6.48
ResNeXt50 [7] (32x4d) + SE [28]	27.56M	3.771	21.91	6.04
ResNeXt50 [7] (32x4d) + CBAM	27.56M	3.774	21.92	5.91
ResNeXt101 [7] (32x4d)	44.18M	7.508	21.54	5.75
ResNeXt101 [7] (32x4d) + SE [28]	48.96M	7.512	21.17	5.66
ResNeXt101 [7] (32x4d) + CBAM	48.96M	7.519	21.07	5.59

RAN / SE / BAM / CBAM Comparison

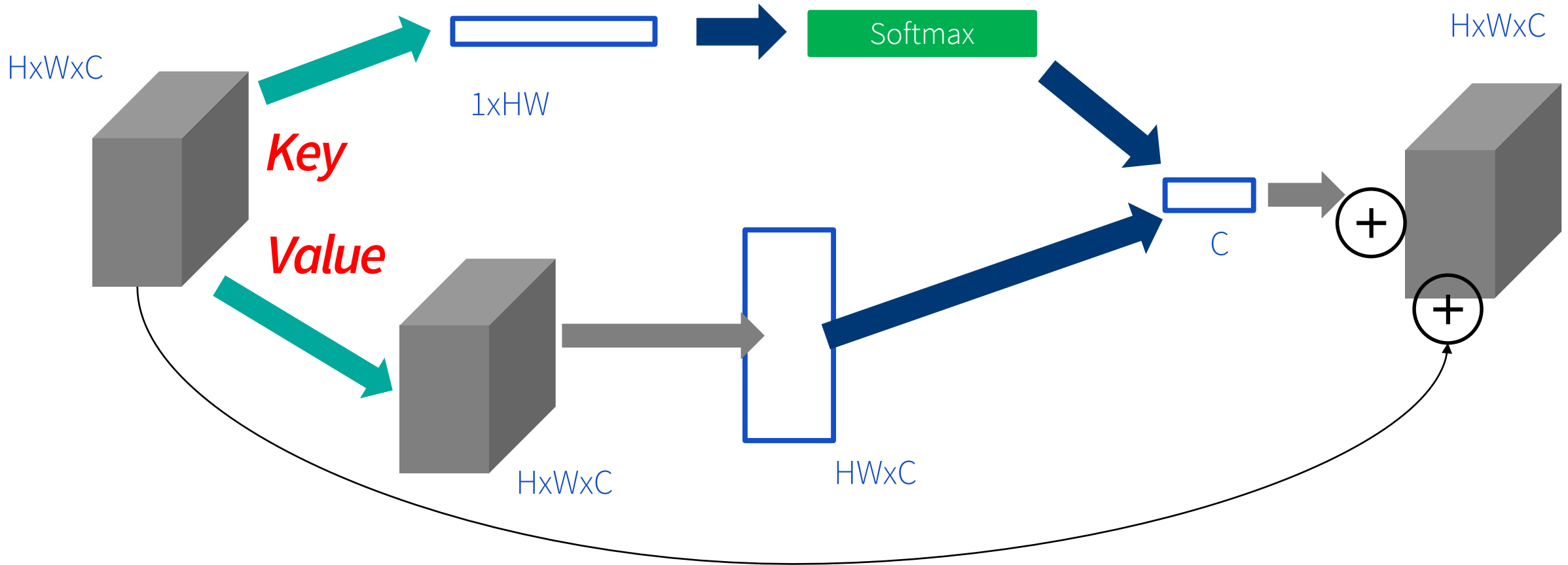
Network	Module Position	Attention
RAN (CVPR17)	Modified	ChannelxSpatial 3D
SE (CVPR18)	In the ResBlock	Channel
BAM (BMVC18)	Before the Stride=2 ResBlock	Channel, Spatial Parallel
CBAM (ECCV18)	In the ResBlock	Channel, Spatial Sequential

Non-local Networks

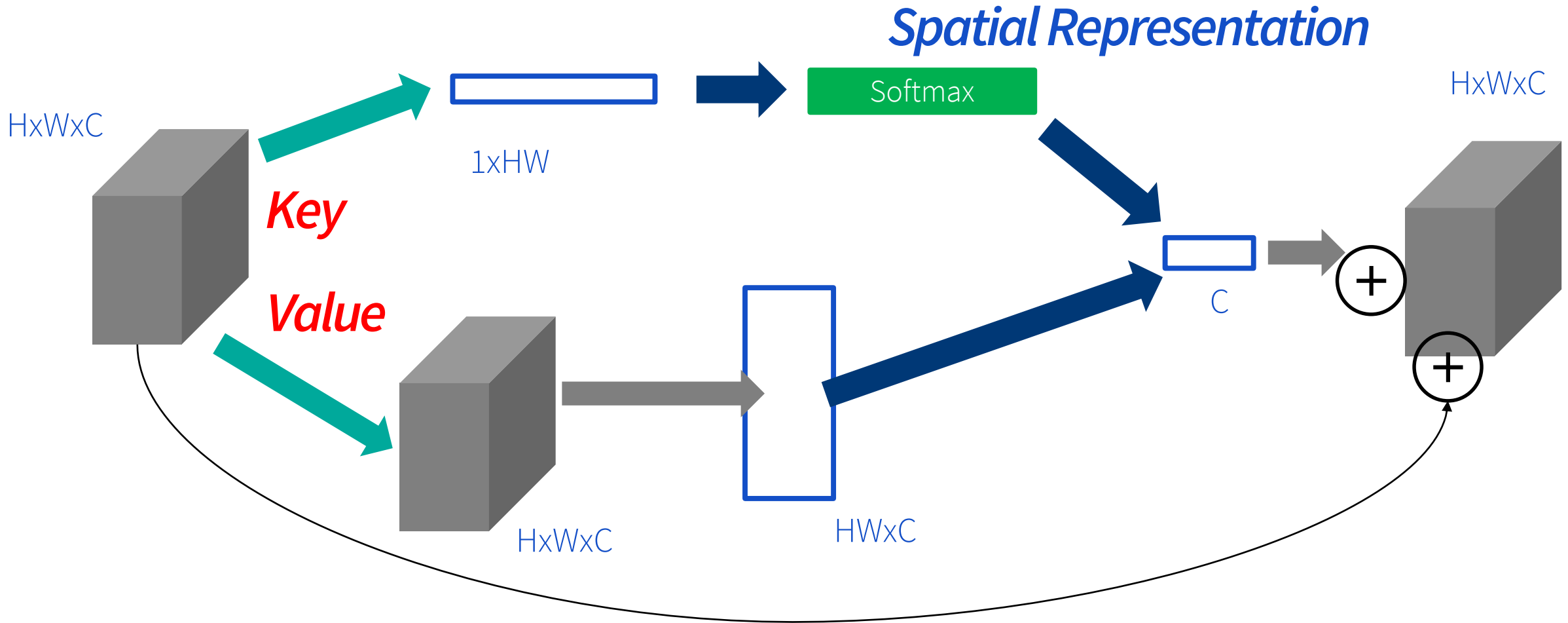
- *Represent Spatial-Only*



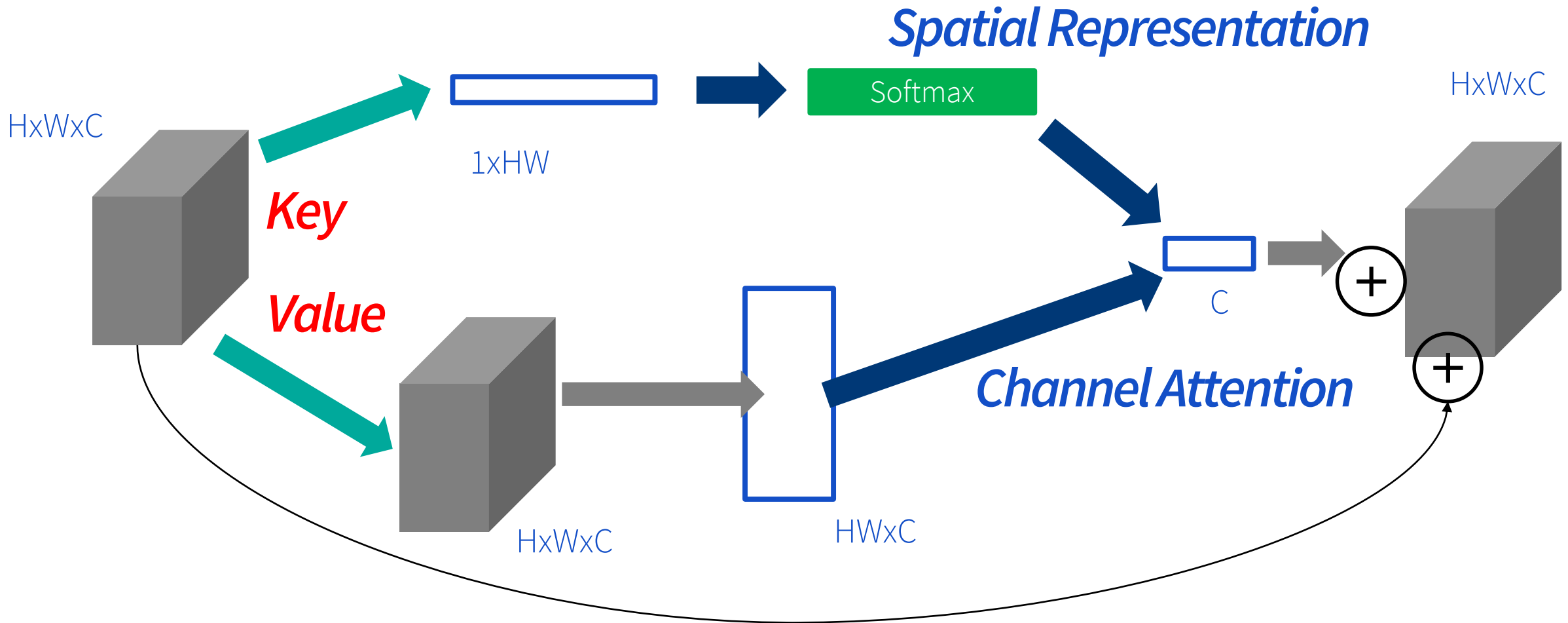
Global-Context Attention Networks



Global-Context Attention Networks



Global-Context Attention Networks

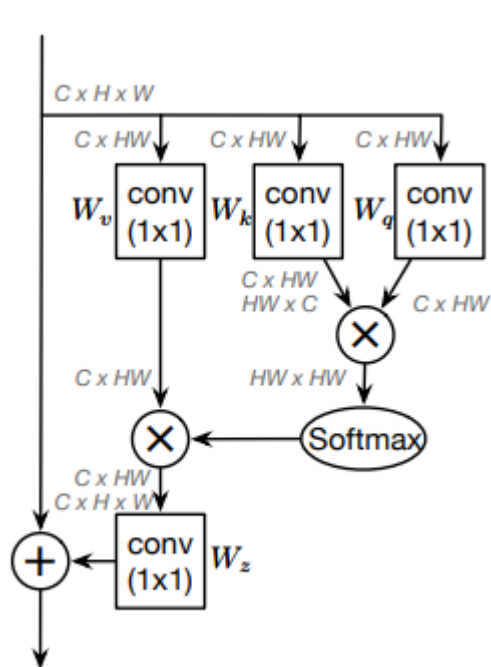


Global-Context Attention Networks



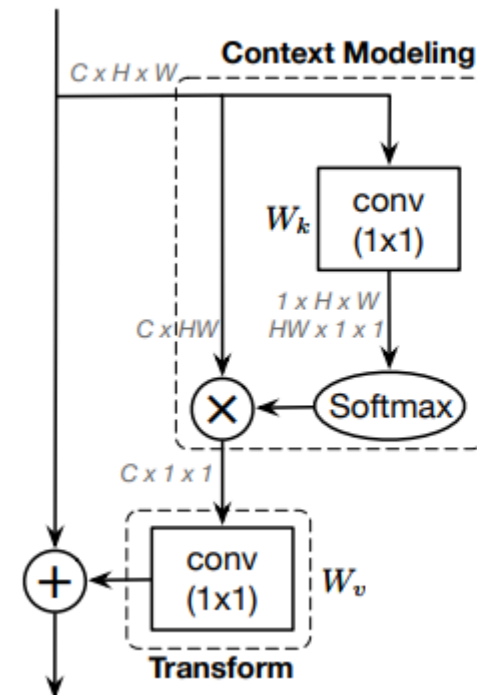
Query Independent Representation → Recalibration

Global-Context Attention Networks



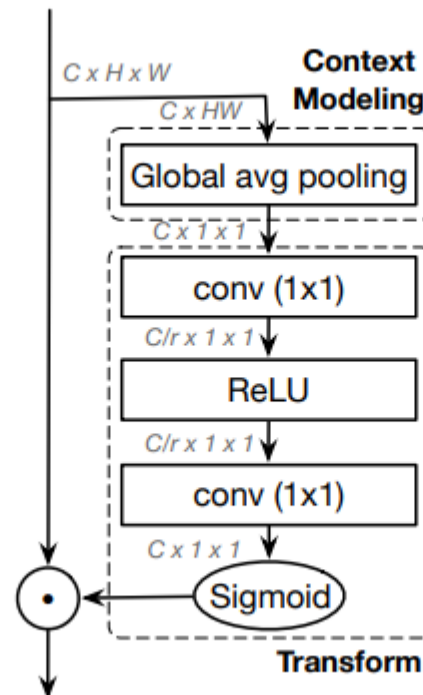
NLNet

Spatial Weighted Sum
Per Pixel
(HxW)



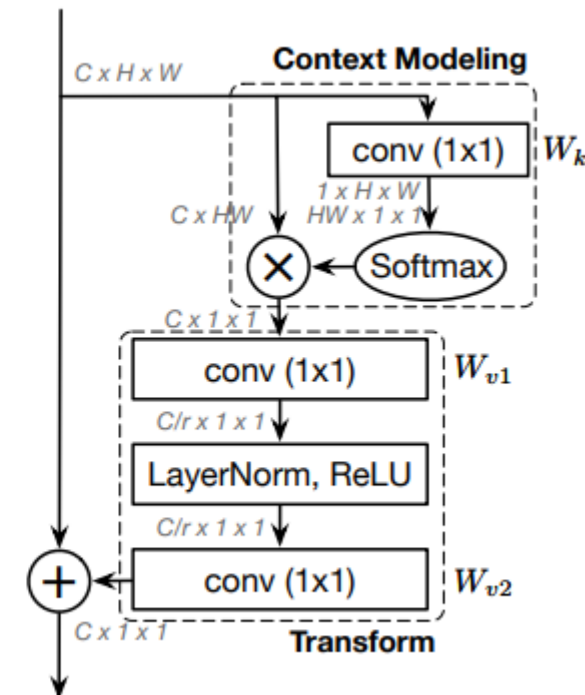
Simplified NLNet

Spatial Weighted Sum
Shared
(Scalar)



SENet

Spatial Aggregation
(Global Avg Pool)
→ for Channel Recalibration



GCNet

Spatial Weighted Sum
→ for Channel Recalibration

Global-Context Attention Networks



Query Independent Representation → Recalibration

Non-local Networks Meet Squeeze-Excitation Networks and Beyond (GCNet)

Global-Context Attention Networks

(a) test on validation set								
backbone		AP ^{bbbox}	AP ^{bbbox} ₅₀	AP ^{bbbox} ₇₅	AP ^{mask}	AP ^{mask} ₅₀	AP ^{mask} ₇₅	FLOPS
R50	baseline	37.2	59.0	40.1	33.8	55.4	35.9	279.4G
	+GC r16	39.4	61.6	42.4	35.7	58.4	37.6	279.6G
	+GC r4	39.9	62.2	42.9	36.2	58.7	38.3	279.6G
R101	baseline	39.8	61.3	42.9	36.0	57.9	38.3	354.0G
	+GC r16	41.1	63.6	45.0	37.4	60.1	39.6	354.3G
	+GC r4	41.7	63.7	45.5	37.6	60.5	39.8	354.3G
X101	baseline	41.2	63.0	45.1	37.3	59.7	39.9	357.9G
	+GC r16	42.4	64.6	46.5	38.0	60.9	40.5	358.2G
	+GC r4	42.9	65.2	47.0	38.5	61.8	40.9	358.2G
X101 +Cascade	baseline	44.7	63.0	48.5	38.3	59.9	41.3	536.9G
	+GC r16	45.9	64.8	50.0	39.3	61.8	42.1	537.2G
	+GC r4	46.5	65.4	50.7	39.7	62.5	42.7	537.3G
X101+DCN +Cascade	baseline	47.1	66.1	51.3	40.4	63.1	43.7	547.5G
	+GC r16	47.9	66.9	52.2	40.9	63.7	44.1	547.8G
	+GC r4	47.9	66.9	51.9	40.8	64.0	44.0	547.8G
(b) test on test-dev set								
X101 +Cascade	baseline	45.0	63.7	49.1	38.7	60.8	41.8	536.9G
	+GC r16	46.5	65.7	50.7	40.0	62.9	43.1	537.2G
	+GC r4	46.6	65.9	50.7	40.1	62.9	43.3	537.3G
X101+DCN +Cascade	baseline	47.7	66.7	52.0	41.0	63.9	44.3	547.5G
	+GC r16	48.3	67.5	52.7	41.5	64.6	45.0	547.8G
	+GC r4	48.4	67.6	52.7	41.5	64.6	45.0	547.8G

(a) Block Design				
	Top-1 Acc	Top-5 Acc	#params(M)	FLOPs(G)
baseline	76.88	93.16	25.56	3.86
+1NL	77.20	93.51	27.66	4.11
+1SNL	77.28	93.60	26.61	3.86
+1GC	77.34	93.52	25.69	3.86
+all GC	77.70	93.66	28.08	3.87
(b) Pooling and fusion				
	Top-1 Acc	Top-5 Acc	#params(M)	FLOPs(G)
baseline	76.88	93.16	25.56	3.86
avg+scale (SENet)	77.26	93.55	28.07	3.87
avg+add	77.40	93.60	28.07	3.87
att+scale	77.34	93.48	28.08	3.87
att+add	77.70	93.66	28.08	3.87

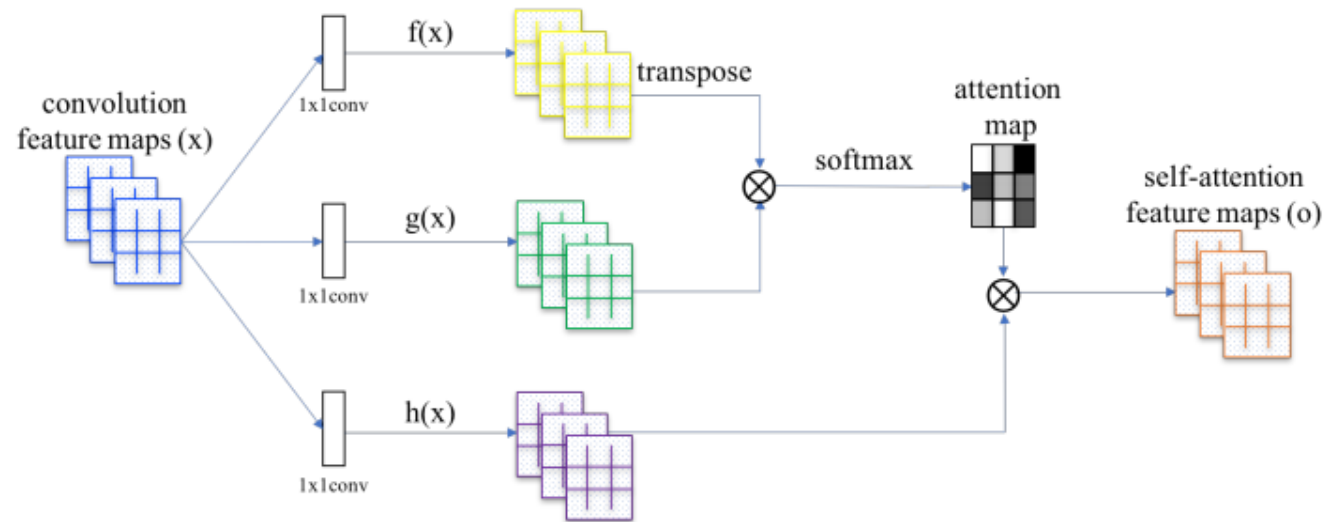
Table 4: **Ablation study** of GCNet with ResNet-50 on **im-
age classification** on ImageNet validation set.

Summary

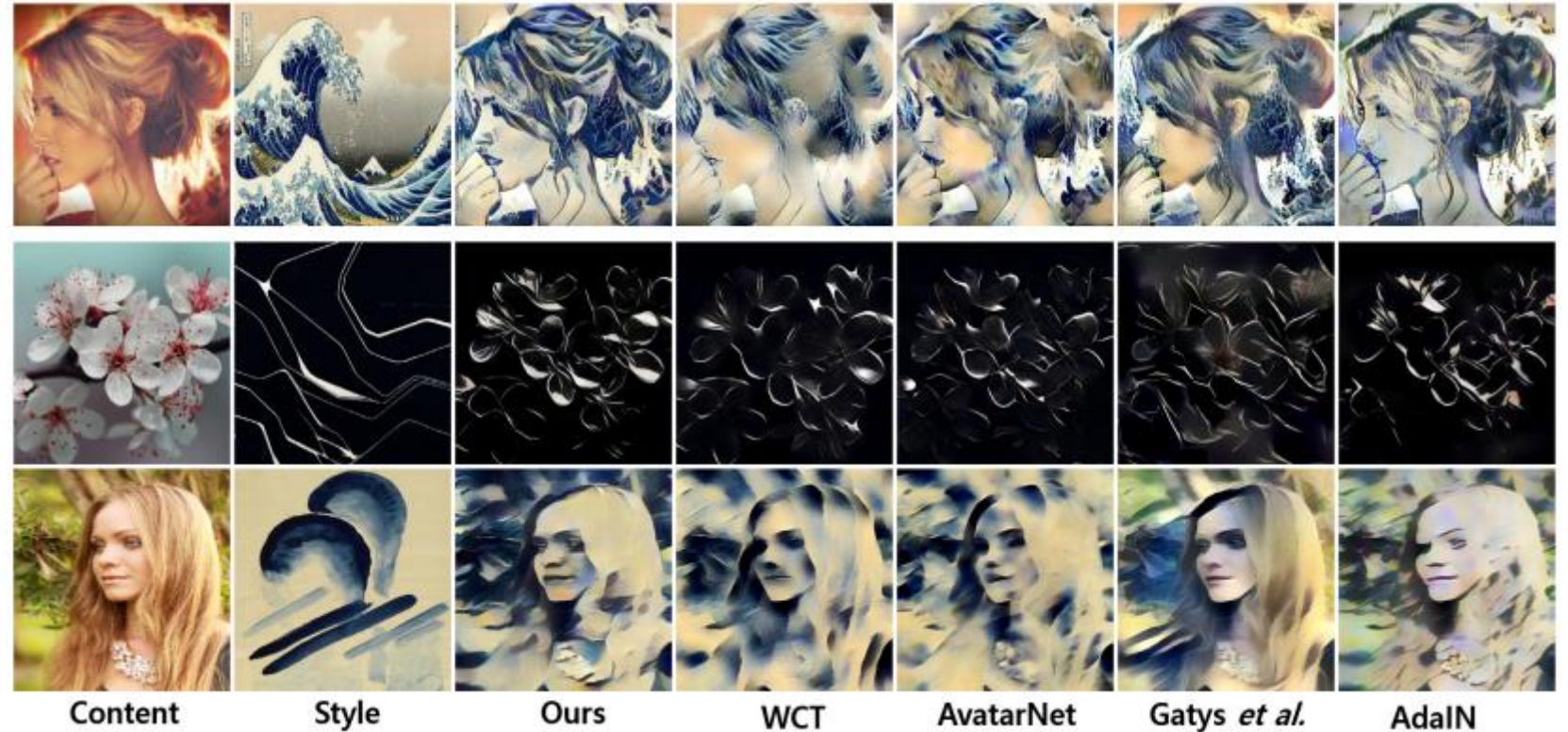
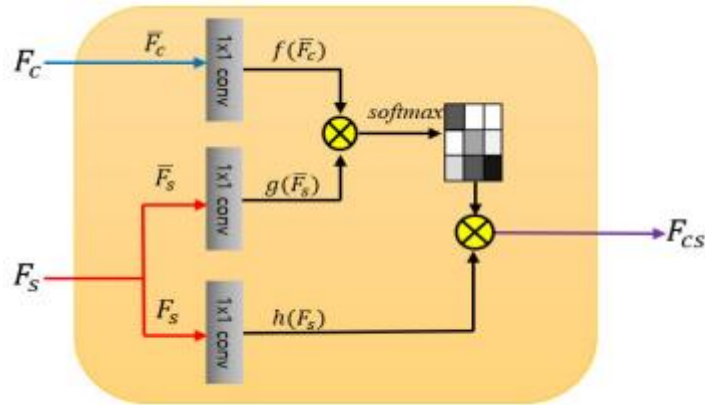
Network	Attention	Spatial Modeling
RAN (CVPR17)	ChannelxSpatial 3D	Network
SE (CVPR18)	Channel	Avg Pool
BAM (BMVC18)	Channel, Spatial Parallel	Avg Pool
CBAM (ECCV18)	Channel, Spatial Sequential	Avg Pool + Max Pool
NLNet (CVPR18)	Spatial (Representation)	Non-local Representation
GCNet (Preprint19)	Channel	Non-local Representation

CNN x Attention: Other Vision Tasks

Self-Attention GAN

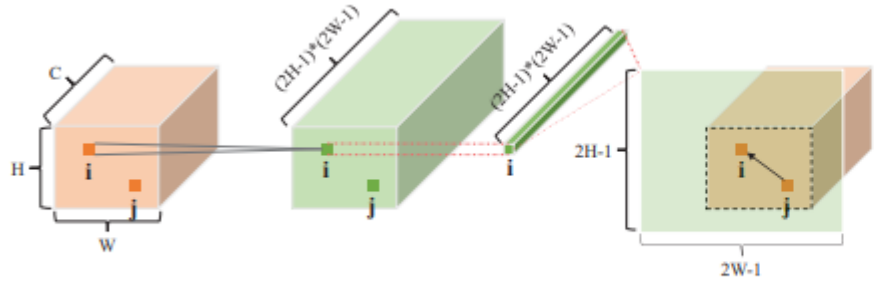


Style Transfer (CVPR19)

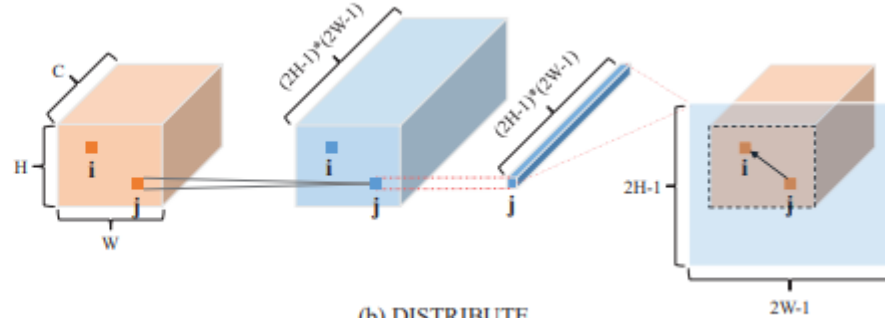


Arbitrary Style Transfer with Style-Attentional Networks

PSANet (ECCV18) / Context Encoding (CVPR18) / OCNet (2018)

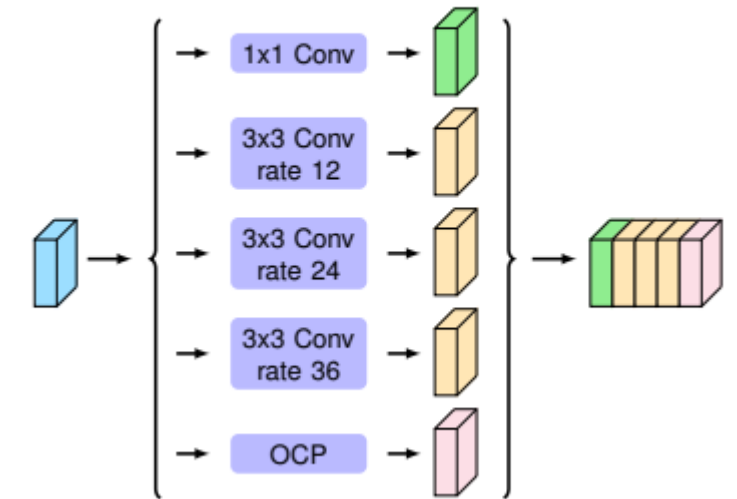
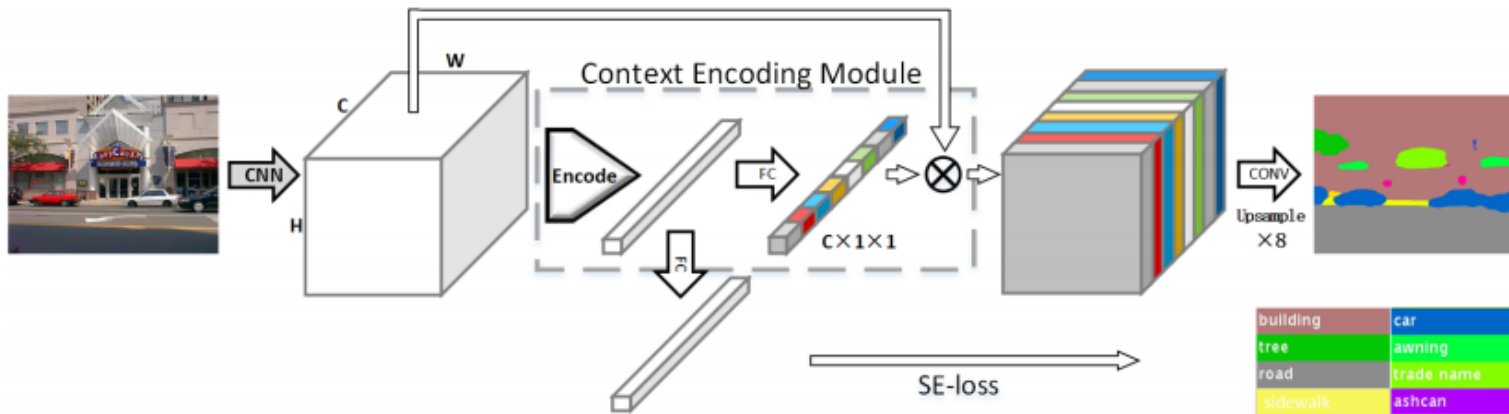


(a) COLLECT

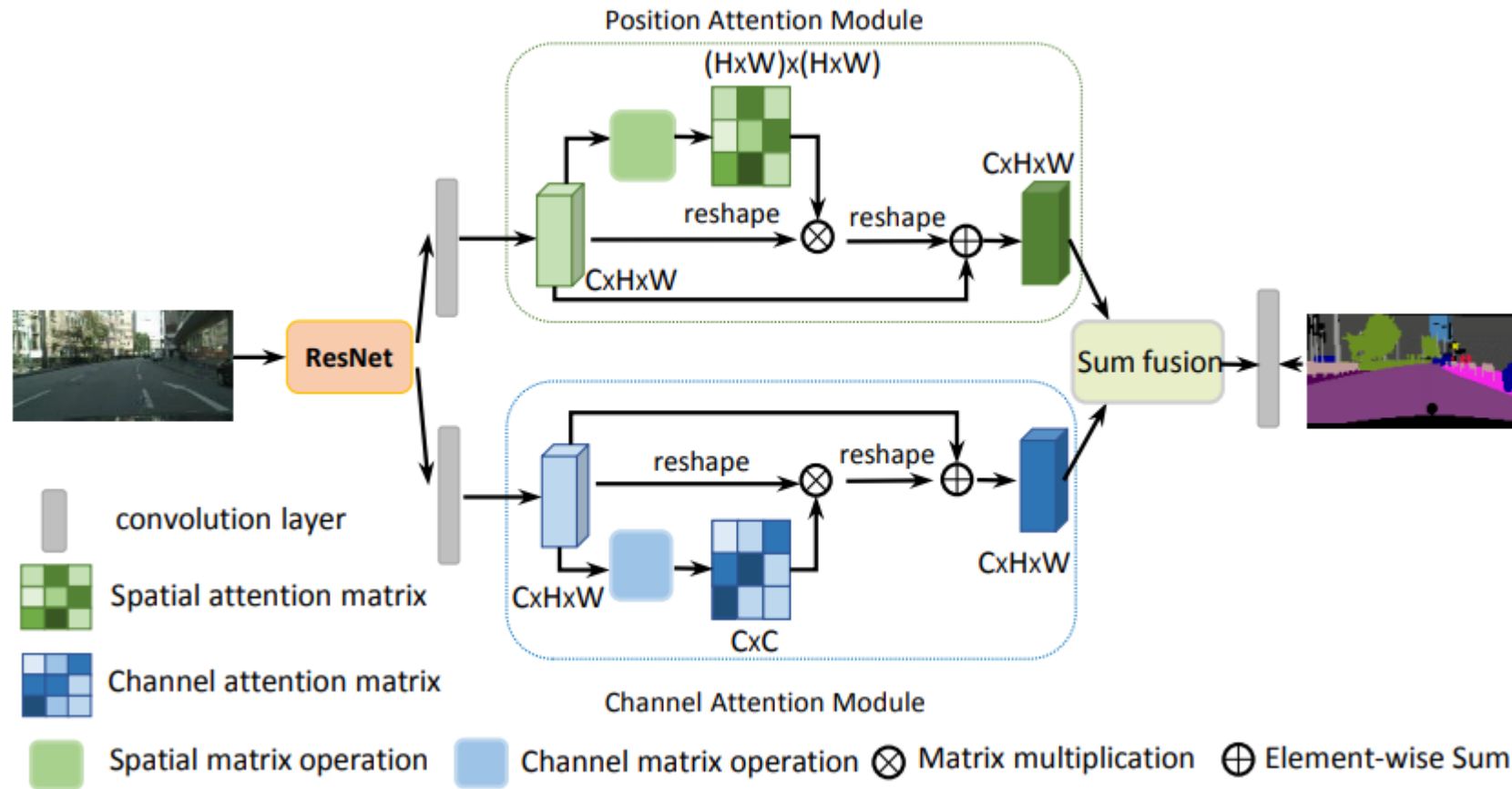


(b) DISTRIBUTE

(d) ASP-OC



Dual Attention Network (CVPR19)



Criss-Cross Non-local Attention Networks (2019)

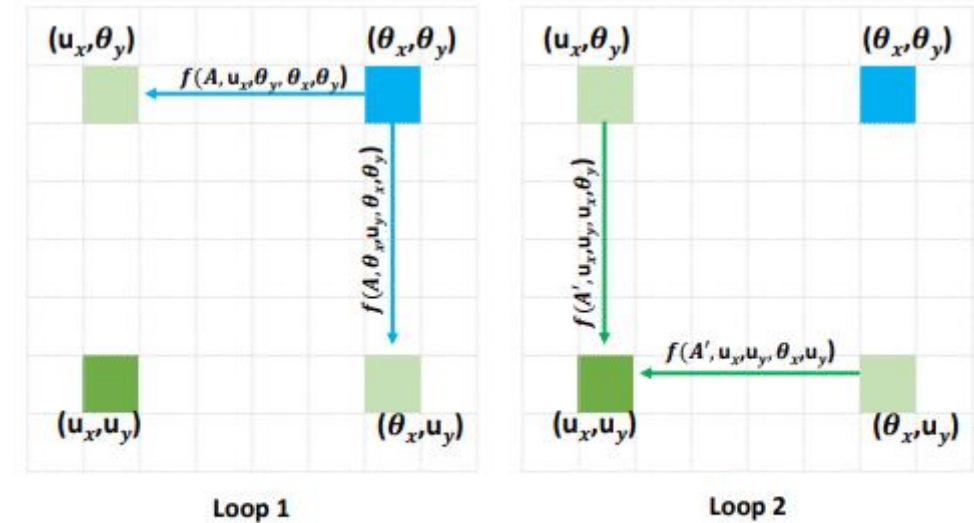
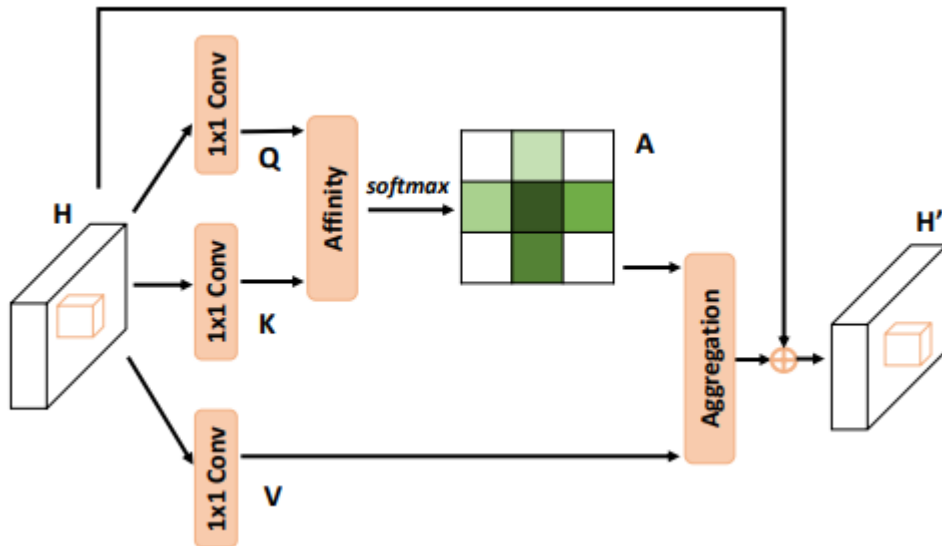


Figure 4. An example of information propagation when the loop number is 2.

Semantic Segmentation

Network	Performance (Cityscape) mIoU	Structure
DenseASPP (CVPR18)	80.6	DenseNet
PSANet (ECCV18)	80.1	Spatial Attention
Context Encoding (CVPR18)	-	Channel Attention
CCNet (Arxiv19)	81.4	Fast NL-Net
DANet (CVPR19)	81.5	NL-Net (Spatial + Channel)
OCNet (Arxiv18)	81.7	NL-Net + PSP

Non-local in SISR

NL-Means Method:
Buades (2005)

p, q neighbors define
a vector distance;

$$\| \vec{V}_p - \vec{V}_q \|^2$$

Filter with this:

No spatial term!



$$NLMF[I]_p = \frac{1}{W_p} \sum_{q \in S} \cancel{G_{\sigma_s}(\| \vec{p} - \vec{q} \|)} G_{\sigma_r}(\| \vec{V}_p - \vec{V}_q \|^2) I_q$$

Single Image Super-Resolution

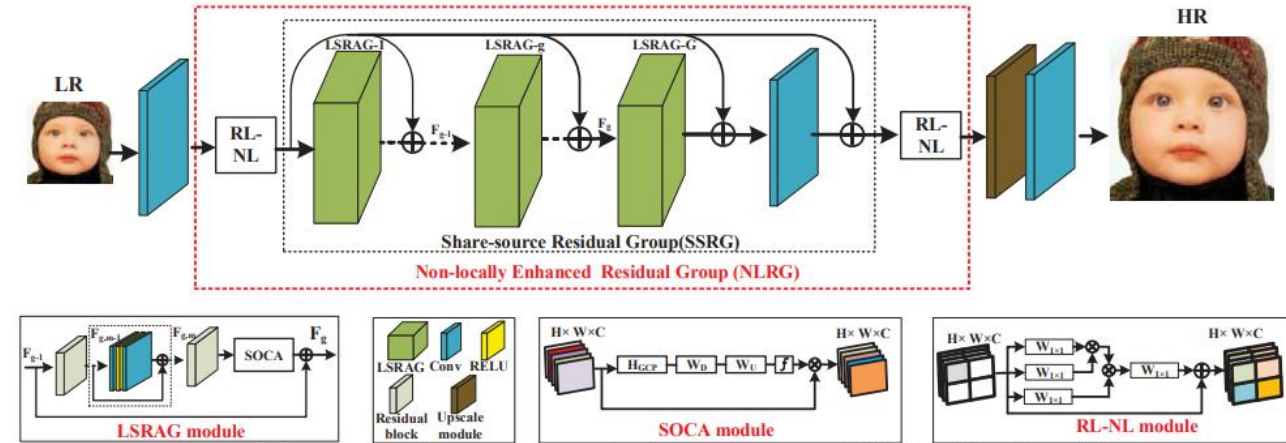
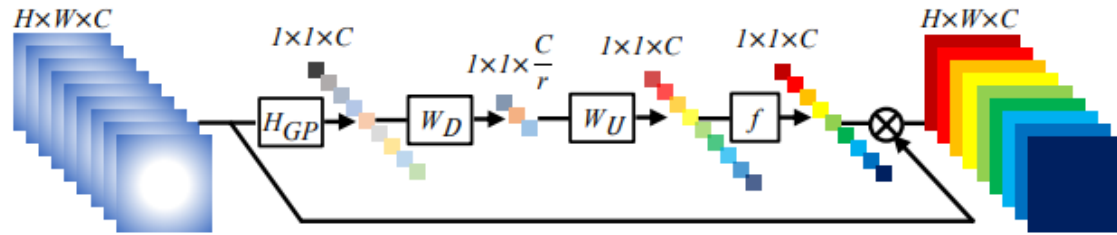


Figure 2. Framework of the proposed second-order attention network (SAN) and its sub-modules.

Network	Performance (set5, PSNR)	Structure
RDN (CVPR18)	38.24 / 32.47 (x2, x4)	DenseNet
RNRN (ICLR19)	38.17 / 32.49	NL-Net
RCAN (ECCV18)	38.27 / 32.63	Channel Attention
SAN (CVPR19)	38.31 / 32.64	Channel Attention + NL-Net

Conclusion

- *Attention (Recurrent) vs Self-Attention (Feed-Forward)*
- *Representation vs Recalibration*
- *Channel Attention: Simple*
- *Spatial Attention: Global Information*

Summary: CNN Architectures

- Many popular architectures available in model zoos
- ResNet and SENet currently good defaults to use
- Networks have gotten increasingly deep over time
- Many other aspects of network architectures are also continuously being investigated and improved
- Even more recent trend towards meta-learning
- Next time: Recurrent neural networks

Thank You
Q&A?