

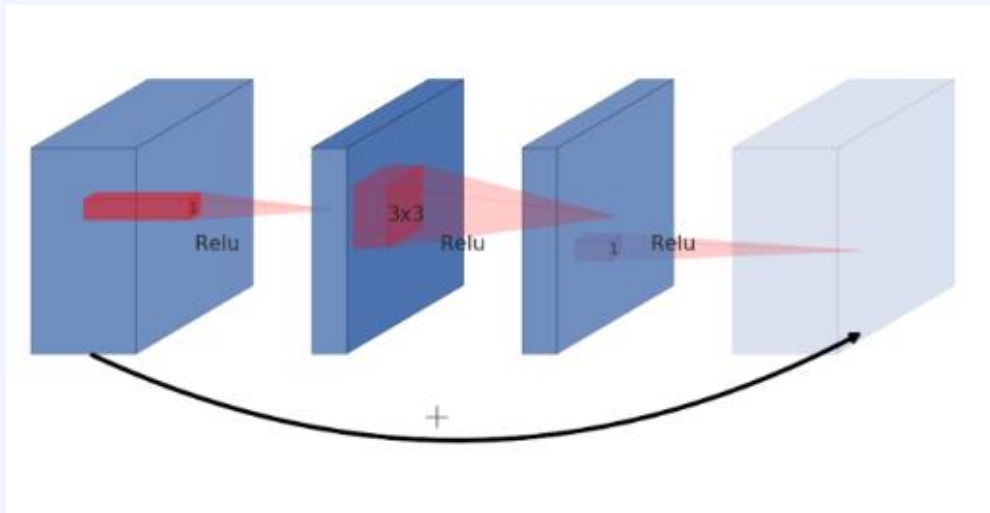
Ch2. CoAtNet

CoAtNet: Marrying Convolution and Attention for All Data Sizes

MB Convolution Block

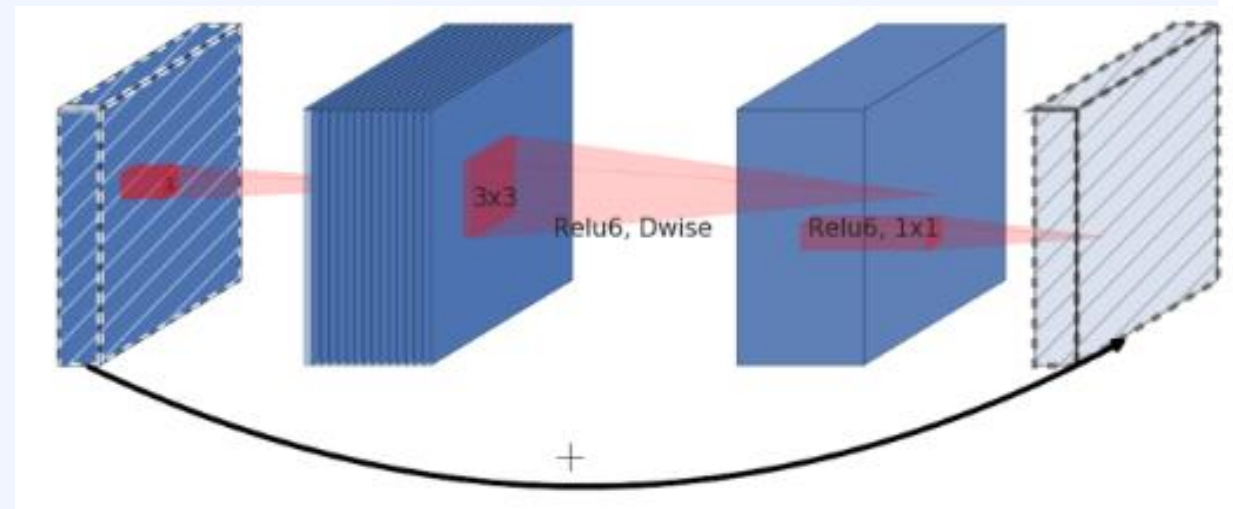
- Depthwise convolution을 사용하는 개선된 inverted residual bottleneck

일반 block



- wide - narrow - wide 한 형태

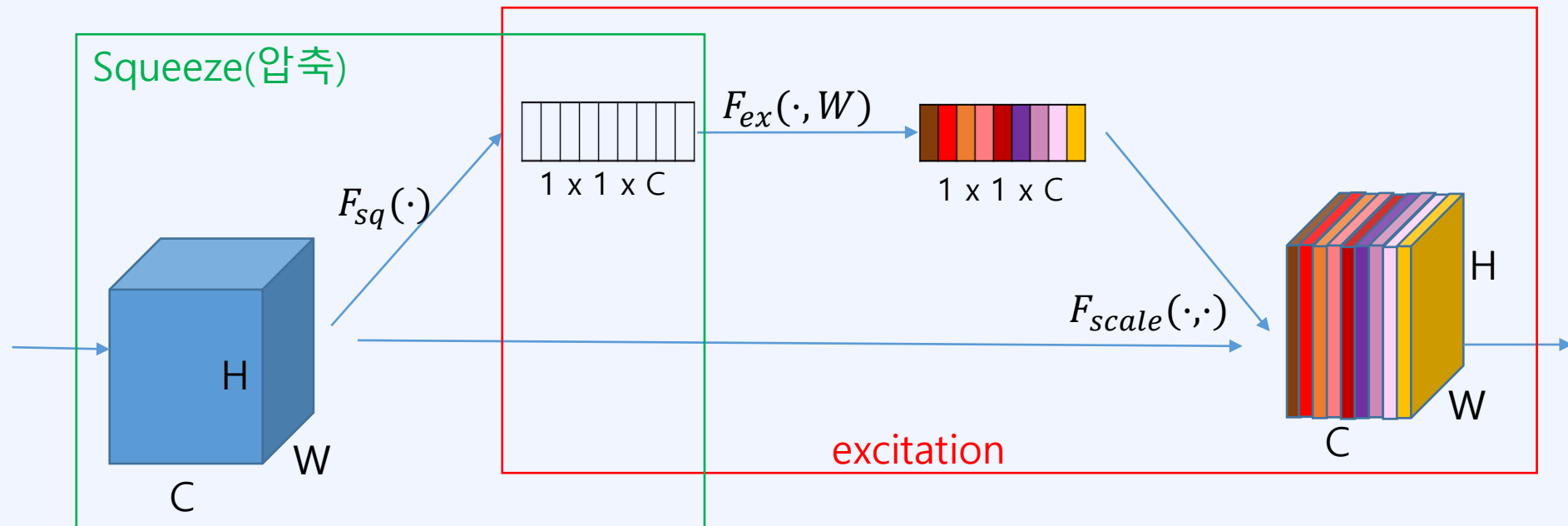
MB block



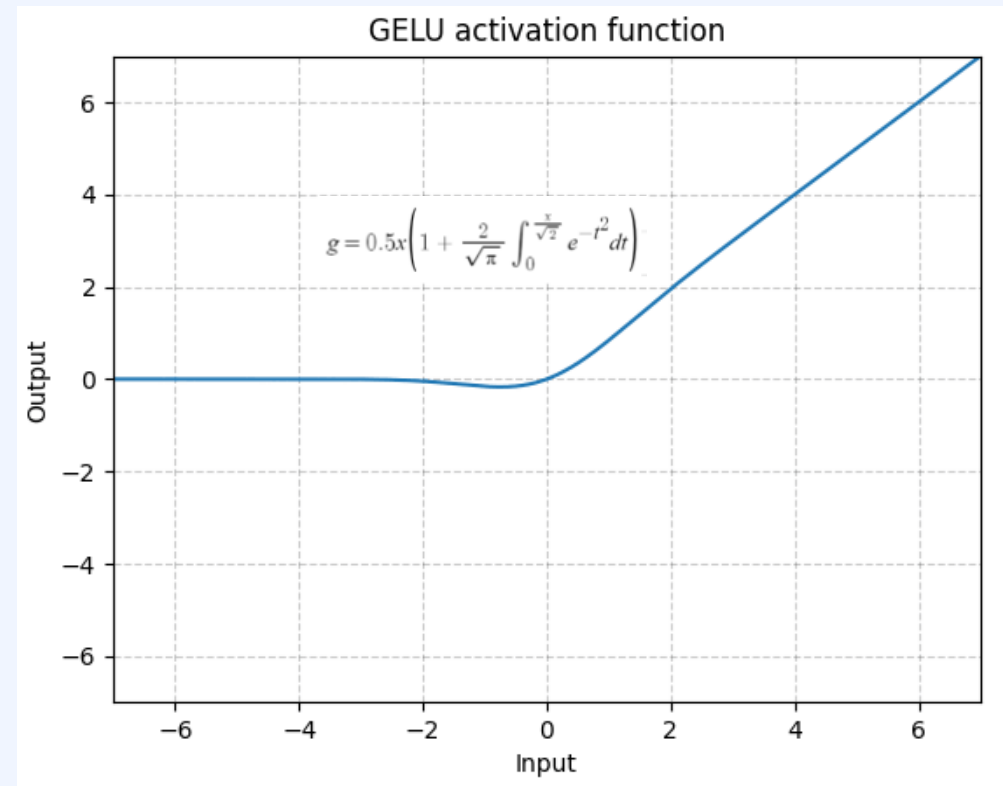
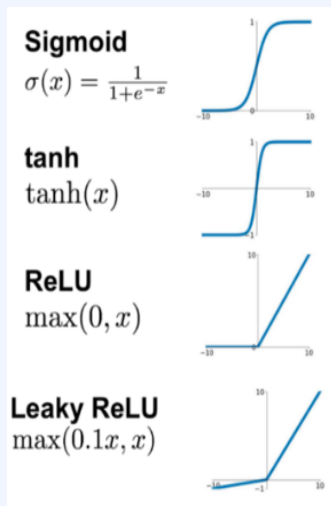
- narrow - wide - narrow한 형태

SENet

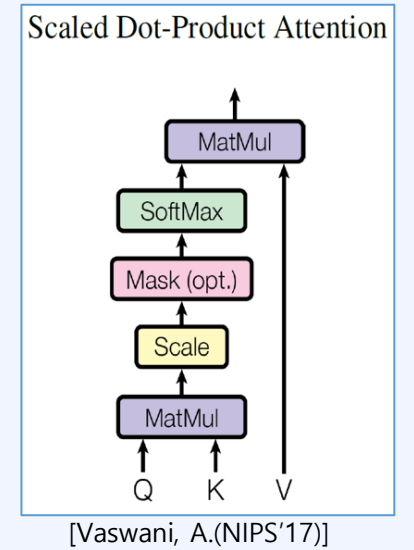
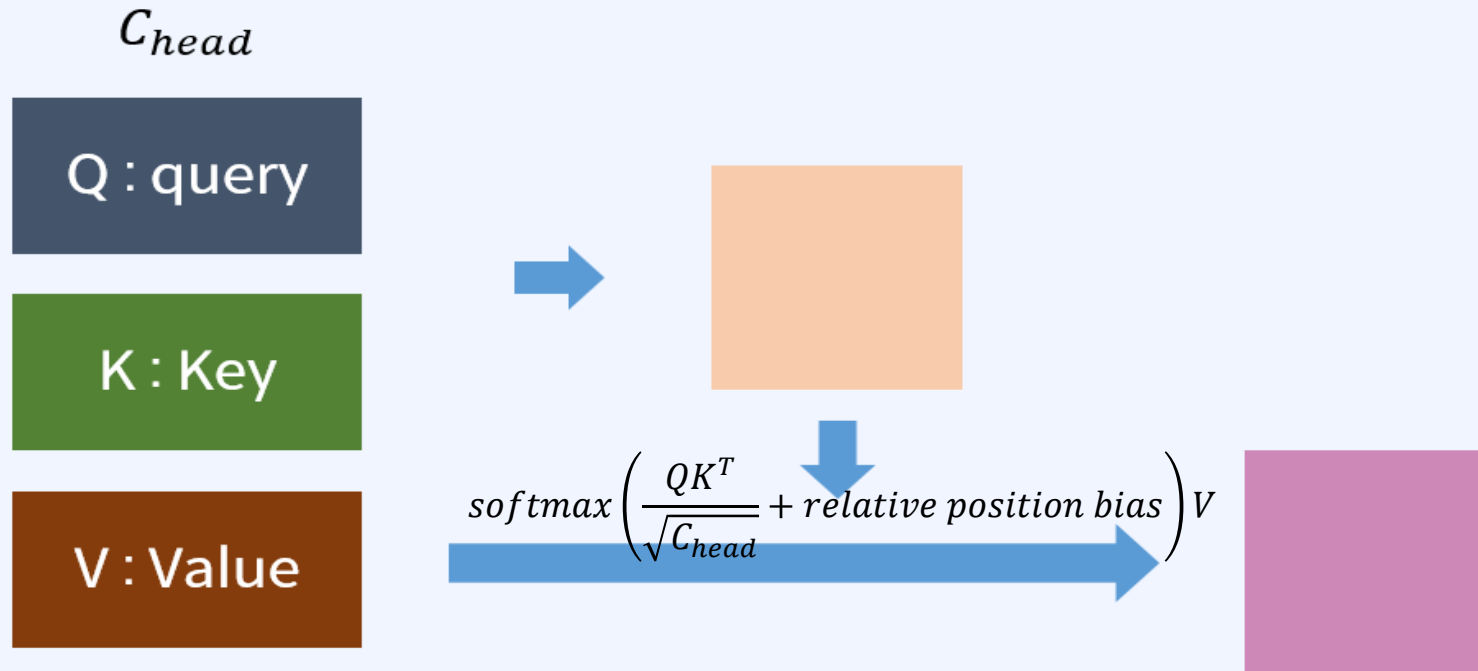
- Squeeze-and-excitation network
- 모델 성능 향상을 위해 기존의 네트워크에 붙여 사용
 - ResNet에서 테스트 시, 0.6% 정확도 향상



GeLU(Gaussian Error Linear units)



2D Relative Attention



Relative position bias

- Relative coordinates

Ex) Feature map size = 3

Ex) Feature Map size = 5

x axis			Δx								
1	2	3	1	2	3	4	5	6	7	8	9
1	2	3	0	0	0	-1	-1	-1	-2	-2	-2
4	5	6	0	0	0	-1	-1	-1	-2	-2	-2
7	8	9	0	0	0	-1	-1	-1	-2	-2	-2

1	0	0	0	-1	-1	-1	-2	-2	-2
2	0	0	0	-1	-1	-1	-2	-2	-2
3	0	0	0	-1	-1	-1	-2	-2	-2
4	1	1	1	0	0	0	-1	-1	-1
5	1	1	1	0	0	0	-1	-1	-1
6	1	1	1	0	0	0	-1	-1	-1
7	2	2	2	1	1	1	0	0	0
8	2	2	2	1	1	1	0	0	0
9	2	2	2	1	1	1	0	0	0

y axis			Δy								
1	2	3	1	2	3	4	5	6	7	8	9
4	5	6									
7	8	9									

1	0	-1	-2	0	-1	-2	0	-1	-2
2	1	0	-1	1	0	-1	1	0	-1
3	2	1	0	2	1	0	2	1	0
4	0	-1	-2	0	-1	-2	0	-1	-2
5	1	0	-1	1	0	-1	1	0	-1
6	2	1	0	2	1	0	2	1	0
7	0	-1	-2	0	-1	-2	0	-1	-2
8	1	0	-1	1	0	-1	1	0	-1
9	2	1	0	2	1	0	2	1	0

Relative position bias

- Relative Position index

Δx

	1	2	3	4	5	6	7	8	9
1	2	2	2	1	1	1	0	0	0
2	2	2	2	1	1	1	0	0	0
3	2	2	2	1	1	1	0	0	0
4	3	3	3	2	2	2	1	1	1
5	3	3	3	2	2	2	1	1	1
6	3	3	3	2	2	2	1	1	1
7	4	4	4	3	3	3	2	2	2
8	4	4	4	3	3	3	2	2	2
9	4	4	4	3	3	3	2	2	2

Δy

	1	2	3	4	5	6	7	8	9
1	2	1	0	2	1	0	2	1	0
2	3	2	1	3	2	1	3	2	1
3	4	3	2	4	3	2	4	3	2
4	2	1	0	2	1	0	2	1	0
5	3	2	1	3	2	1	3	2	1
6	4	3	2	4	3	2	4	3	2
7	2	1	0	2	1	0	2	1	0
8	3	2	1	3	2	1	3	2	1
9	4	3	2	4	3	2	4	3	2

Relative position bias

- Relative Position index

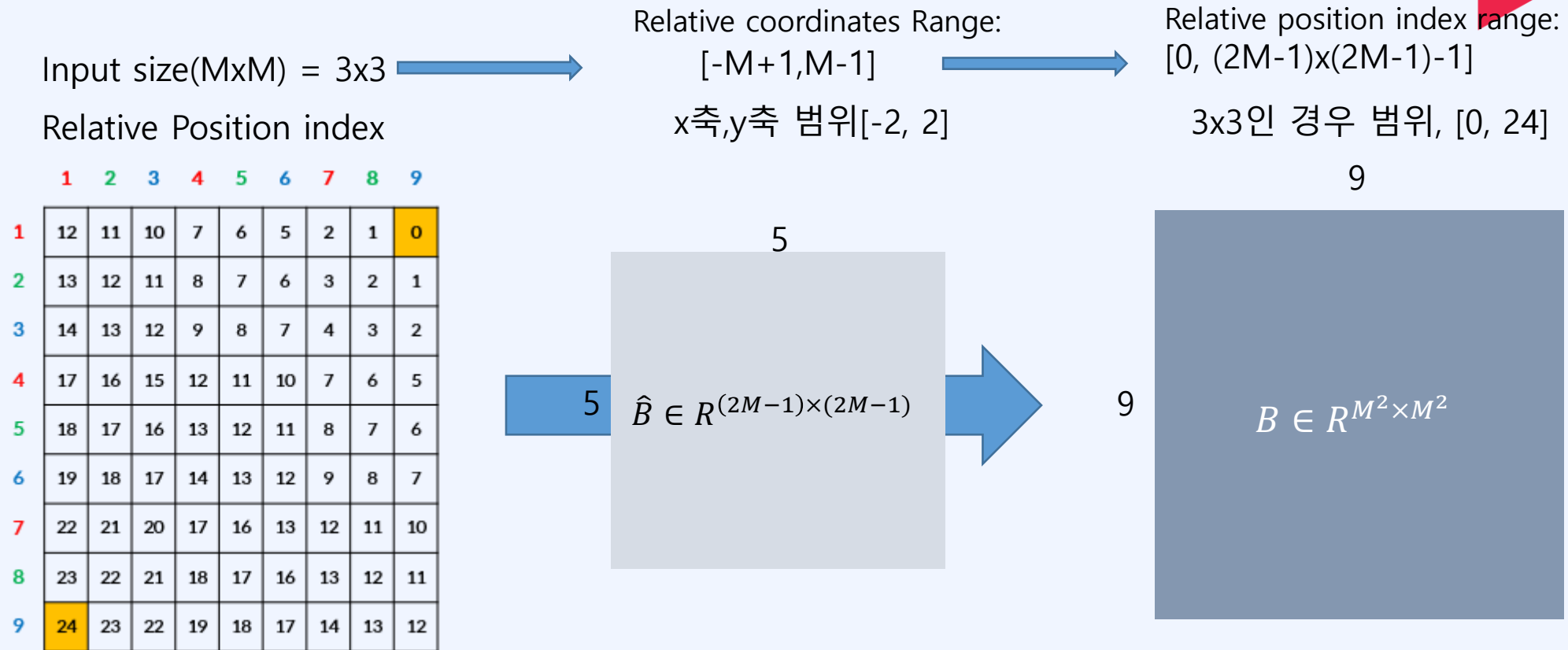
$$\Delta x = \Delta x * (2 * input_{width} - 1)$$

	1	2	3	4	5	6	7	8	9
1	10	10	10	5	5	5	0	0	0
2	10	10	10	5	5	5	0	0	0
3	10	10	10	5	5	5	0	0	0
4	15	15	15	10	10	10	5	5	5
5	15	15	15	10	10	10	5	5	5
6	15	15	15	10	10	10	5	5	5
7	20	20	20	15	15	15	10	10	10
8	20	20	20	15	15	15	10	10	10
9	20	20	20	15	15	15	10	10	10

$$\text{Relative Position matrix} = \Delta x + \Delta y$$

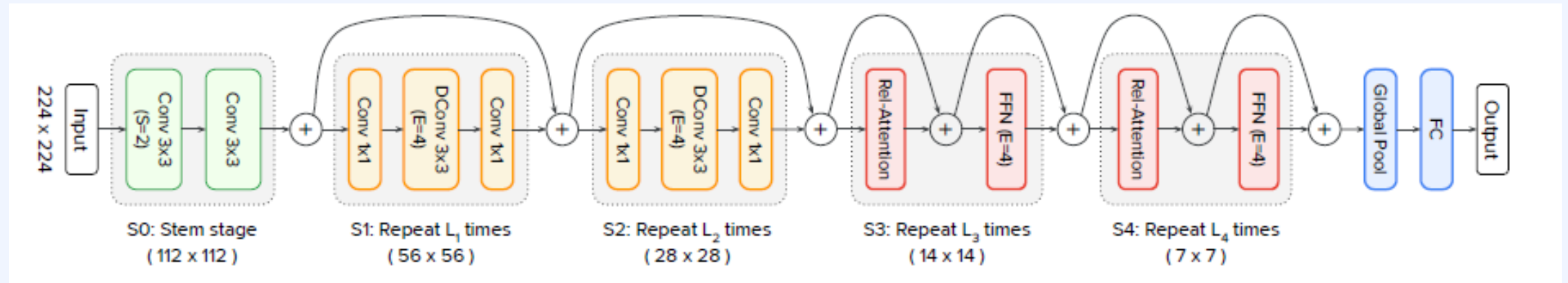
	1	2	3	4	5	6	7	8	9
1	12	11	10	7	6	5	2	1	0
2	13	12	11	8	7	6	3	2	1
3	14	13	12	9	8	7	4	3	2
4	17	16	15	12	11	10	7	6	5
5	18	17	16	13	12	11	8	7	6
6	19	18	17	14	13	12	9	8	7
7	22	21	20	17	16	13	12	11	10
8	23	22	21	18	17	16	13	12	11
9	24	23	22	19	18	17	14	13	12

Relative position bias



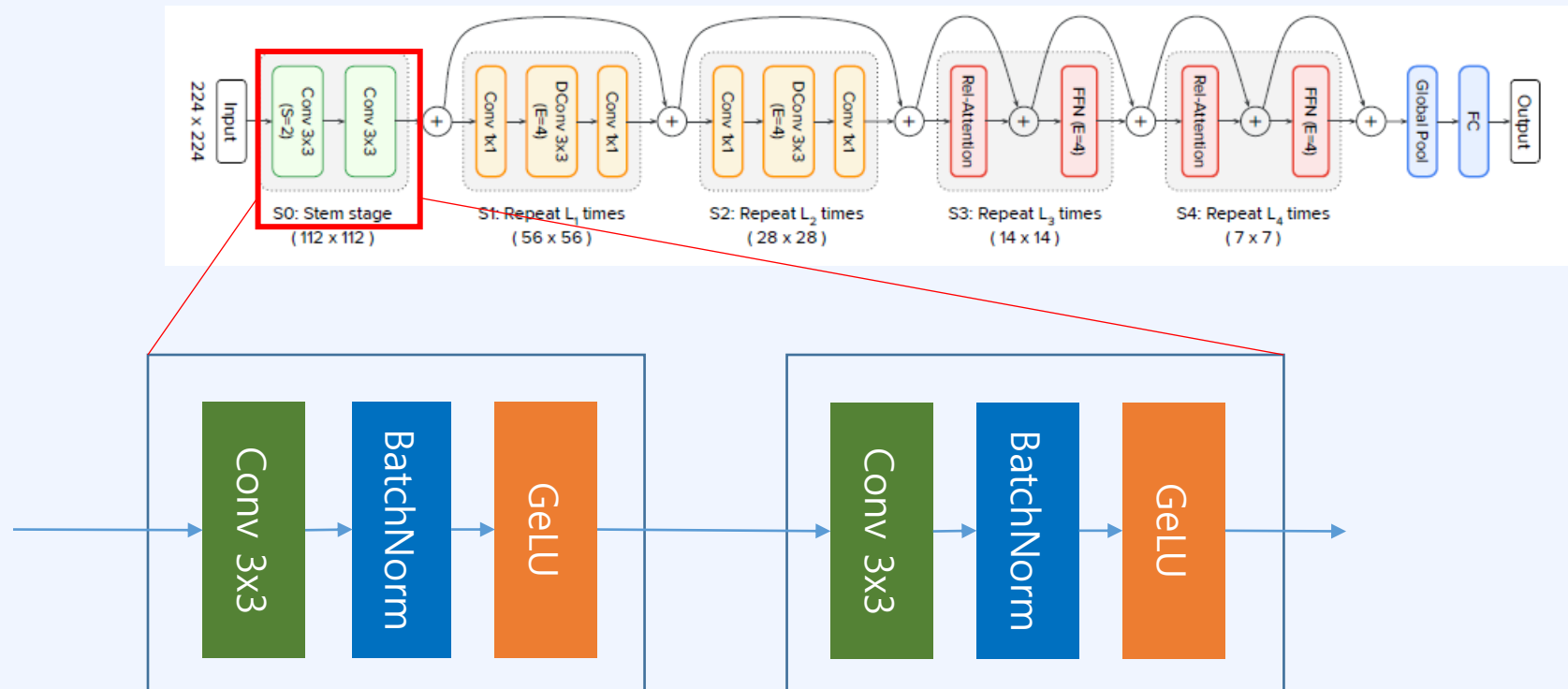
CoAtNet Architecture

CoAtNet Architecture

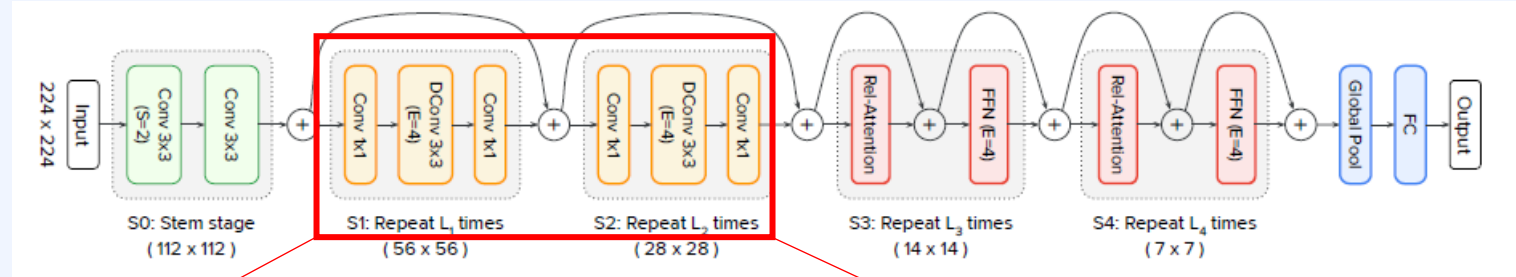


Stage 0

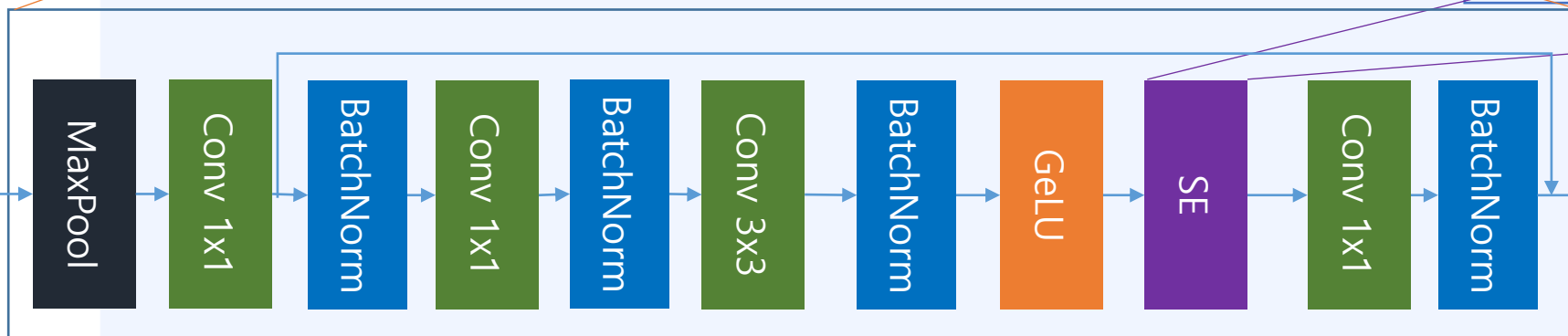
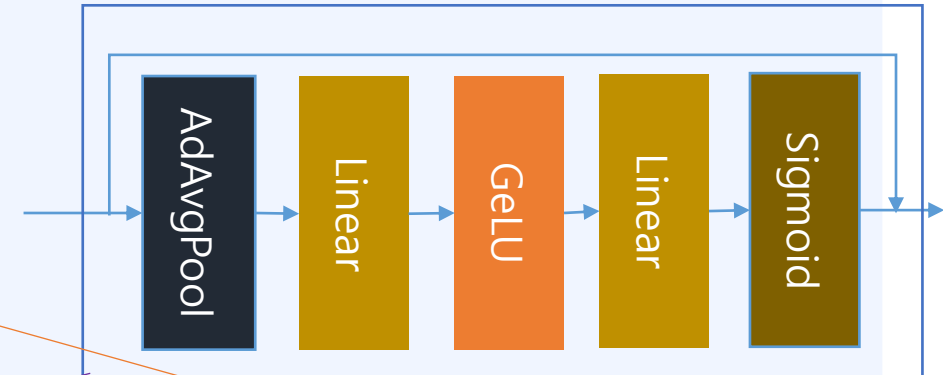
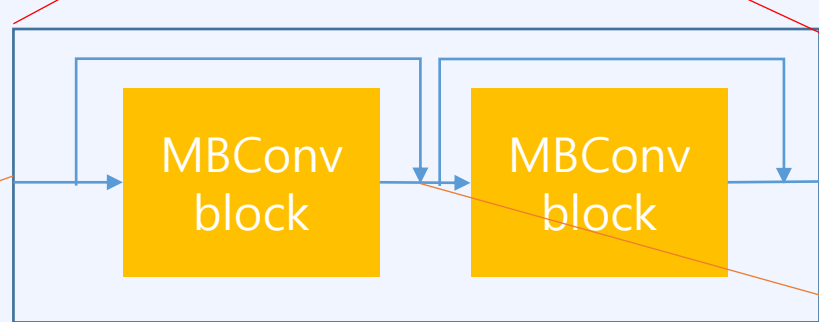
- 2 Layer Convolution stem



Stage1,2

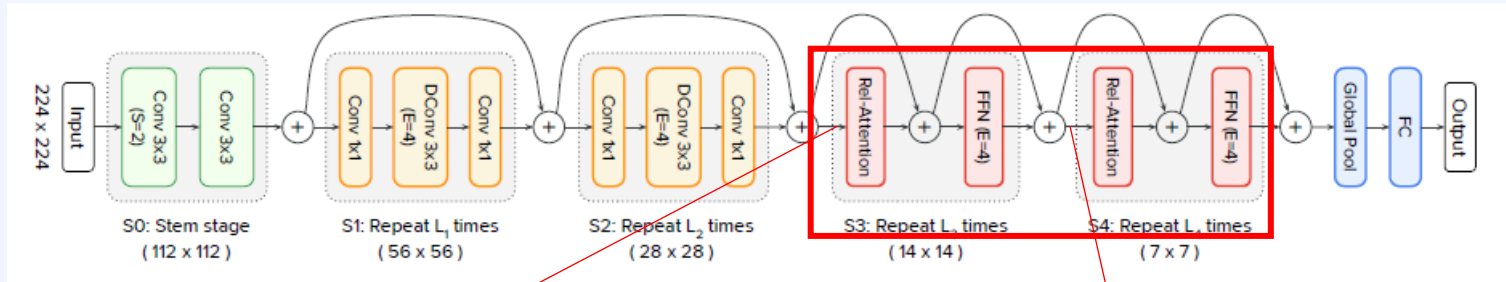


SENet

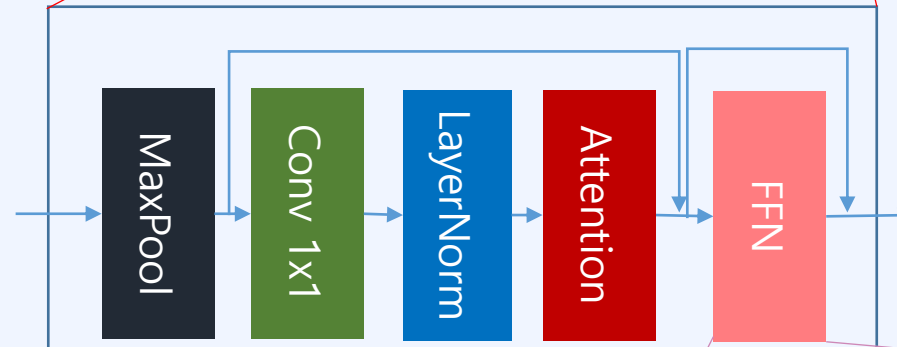


MBConv block

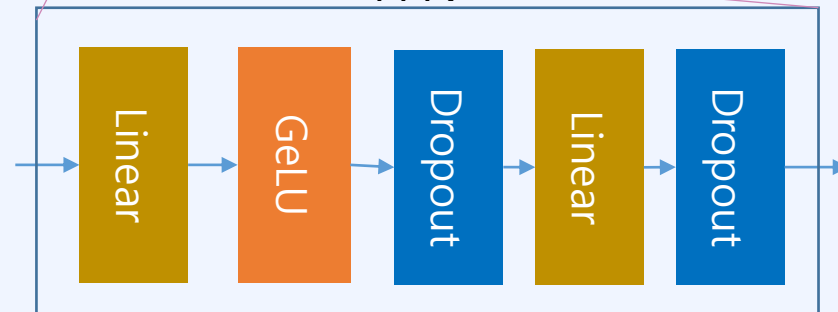
Stage 3,4



Transformer block



FFN



Summary

