

Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

Ze Liu^{1,2†*} Yutong Lin^{1,3†*} Yue Cao^{1*} Han Hu^{1*‡} Yixuan Wei^{1,4†}

Zheng Zhang¹ Stephen Lin¹ Baining Guo¹

¹Microsoft Research Asia ²University of Science and Technology of China

³Xian Jiaotong University ⁴Tsinghua University

{v-zeliul, v-yutlin, yuecao, hanhu, v-yixwe, zhez, stevelin, bainguo}@microsoft.com

Abstract

*This paper presents a new vision Transformer, called Swin Transformer, that capably serves as a general-purpose backbone for computer vision. Challenges in adapting Transformer from language to vision arise from differences between the two domains, such as large variations in the scale of visual entities and the high resolution of pixels in images compared to words in text. To address these differences, we propose a hierarchical Transformer whose representation is computed with **Shifted windows**. The shifted windowing scheme brings greater efficiency by limiting self-attention computation to non-overlapping local windows while also allowing for cross-window connection. This hierarchical architecture has the flexibility to model at various scales and has linear computational complexity with respect to image size. These qualities of Swin Trans-*

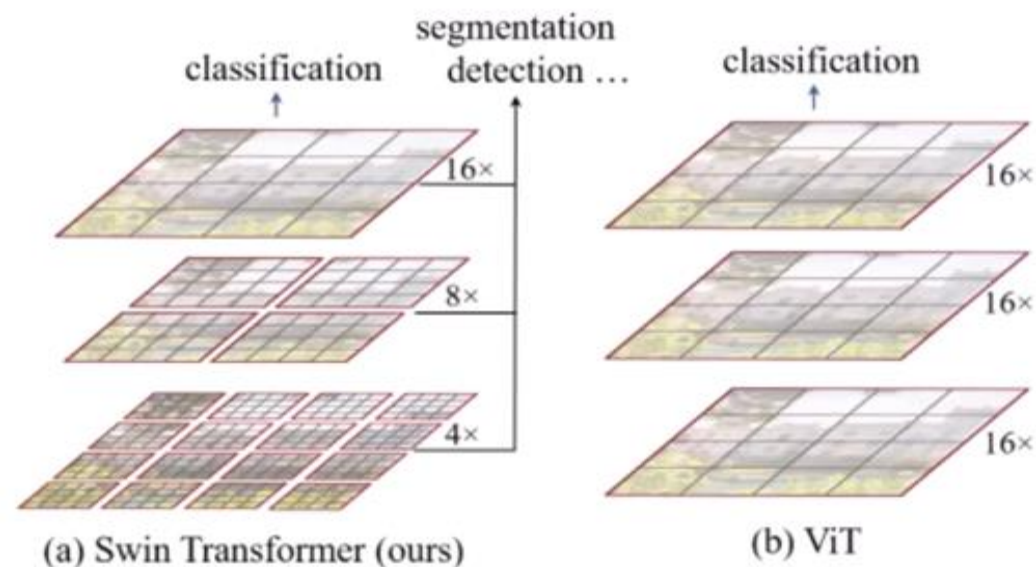
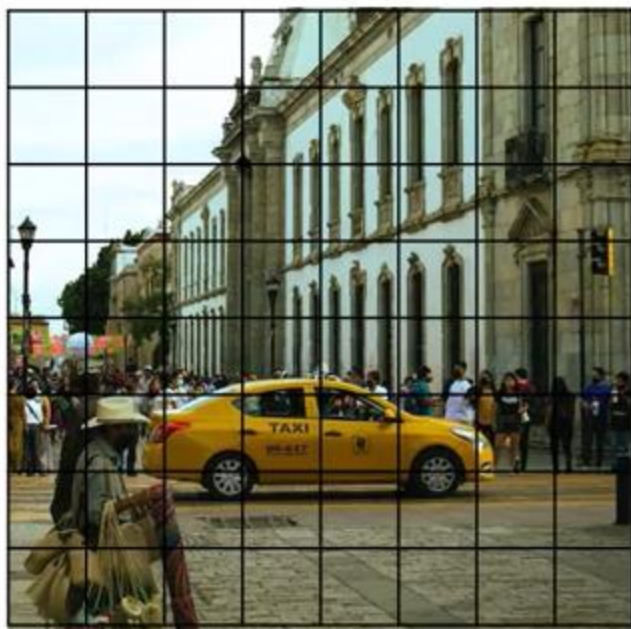
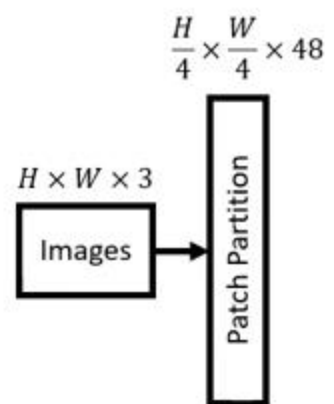
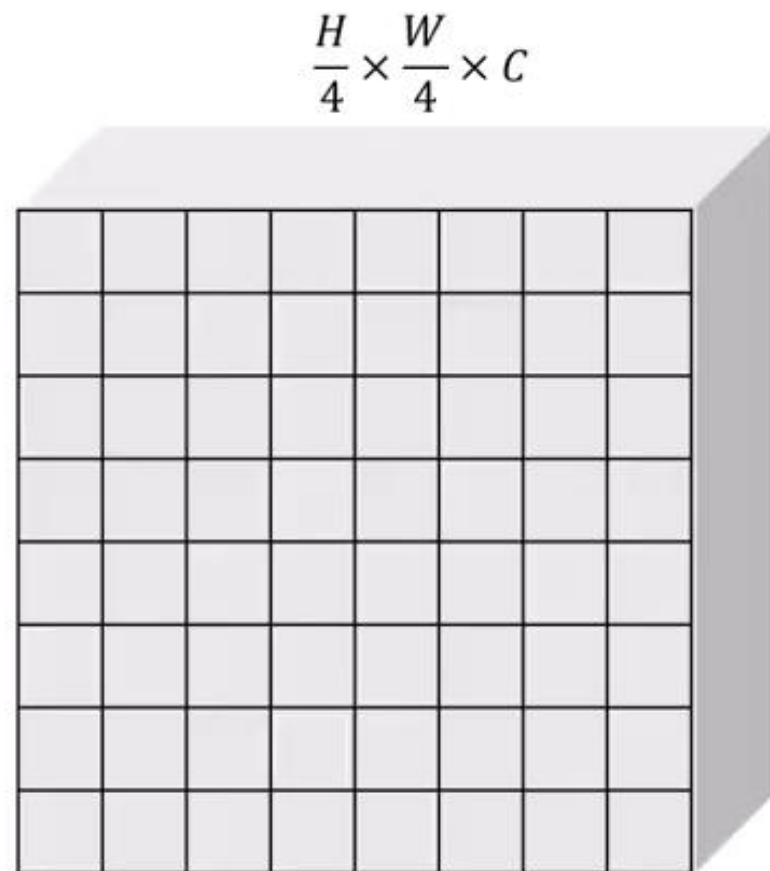
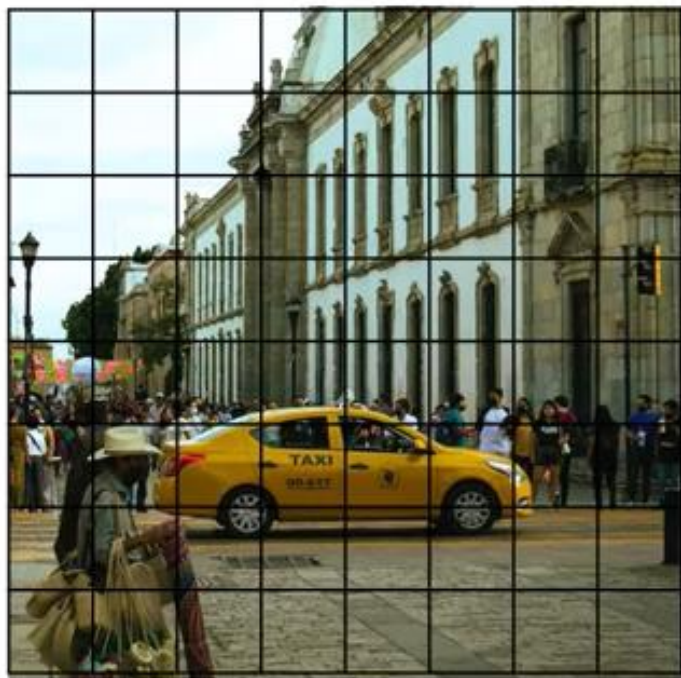
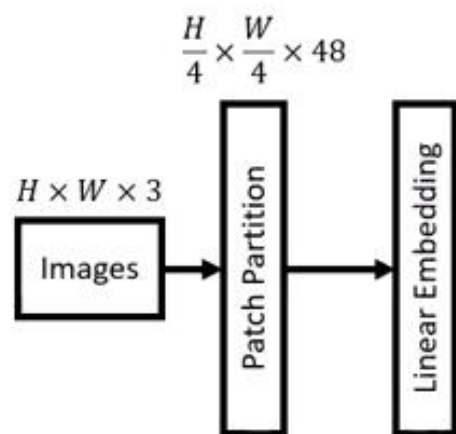
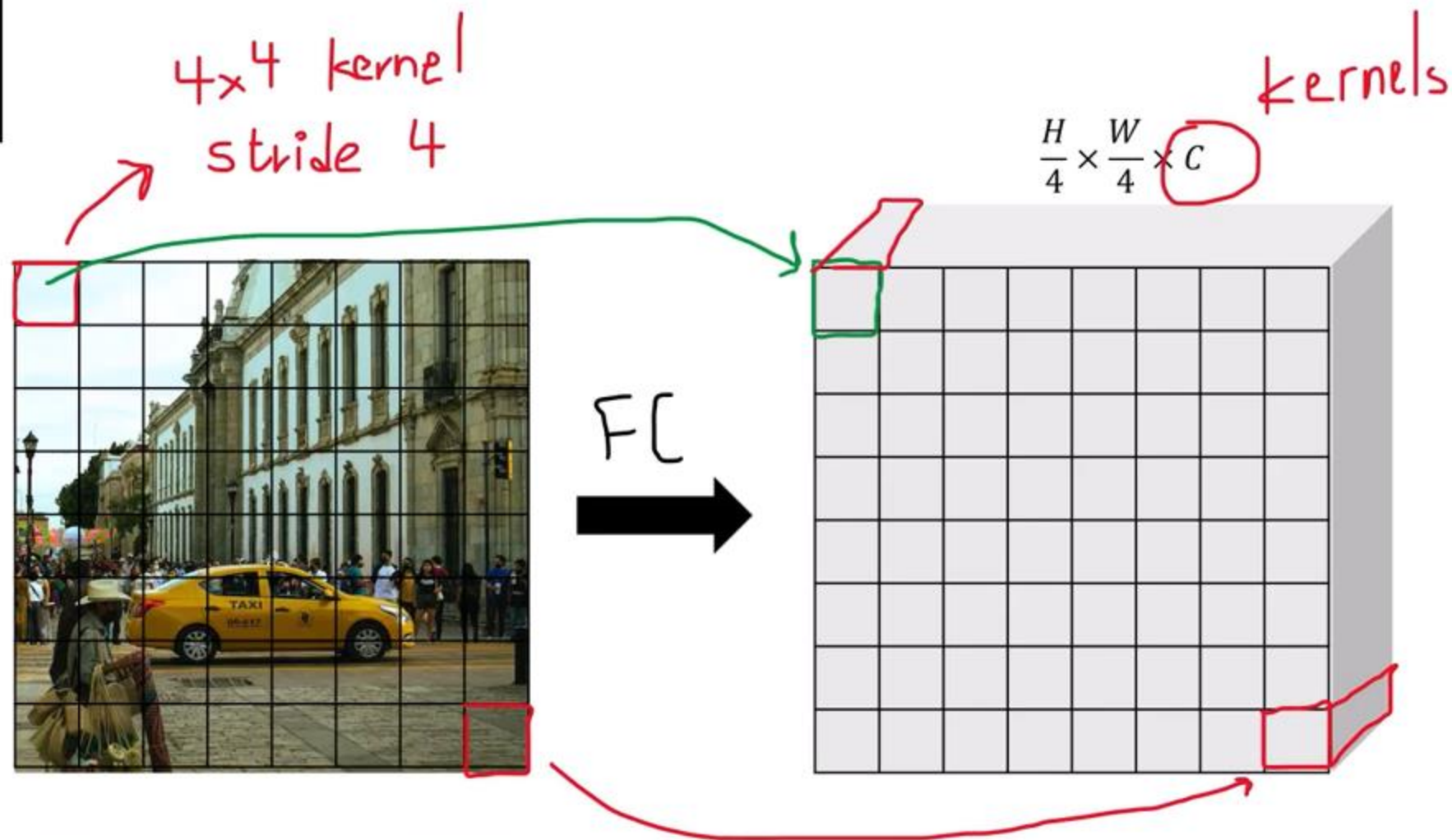
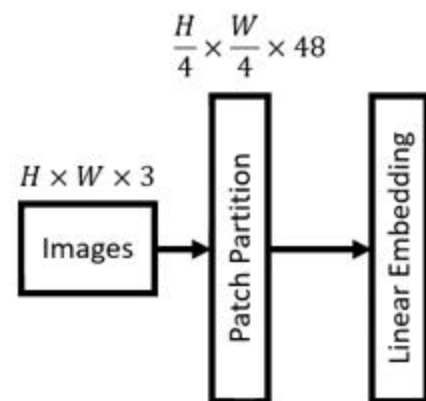
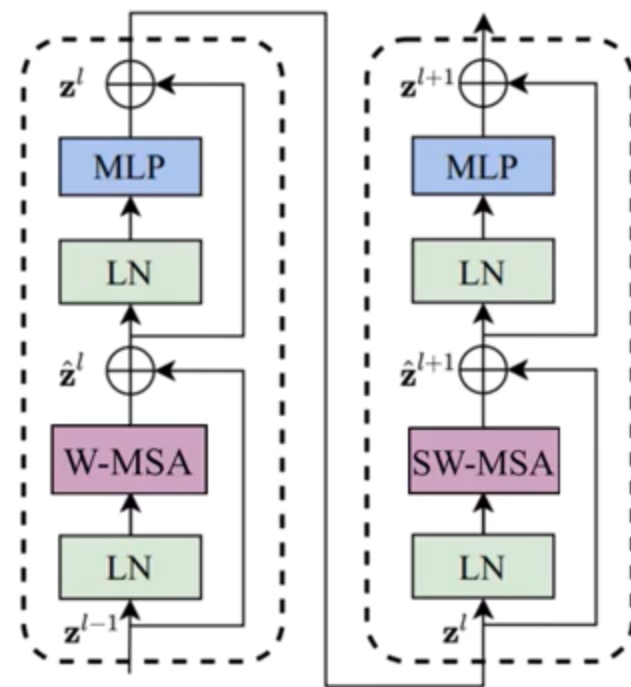
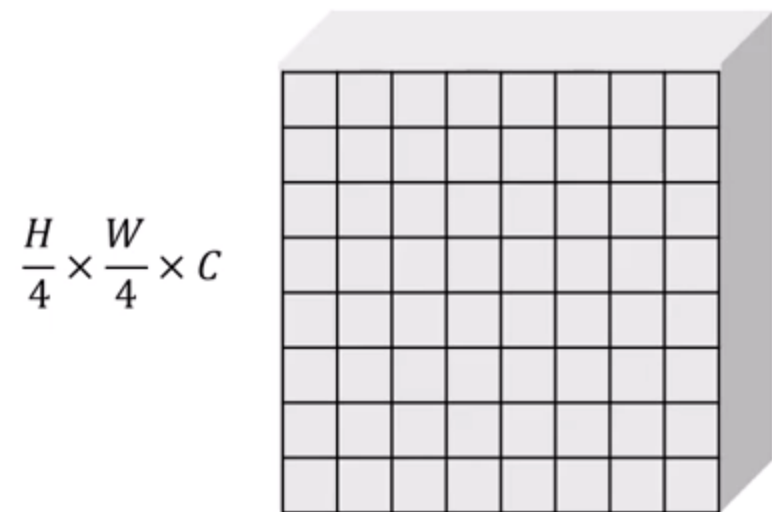
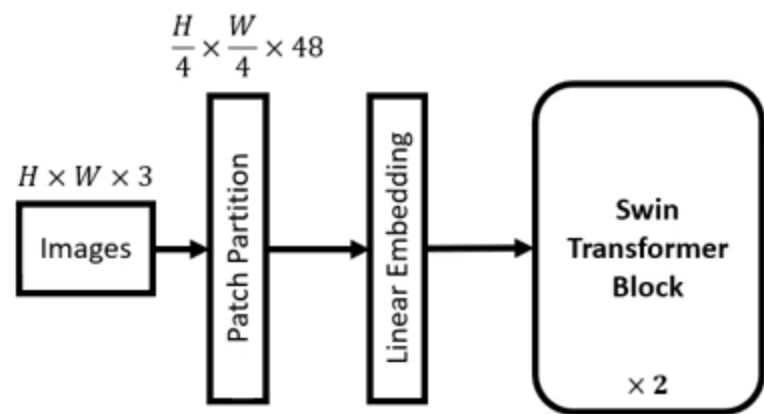


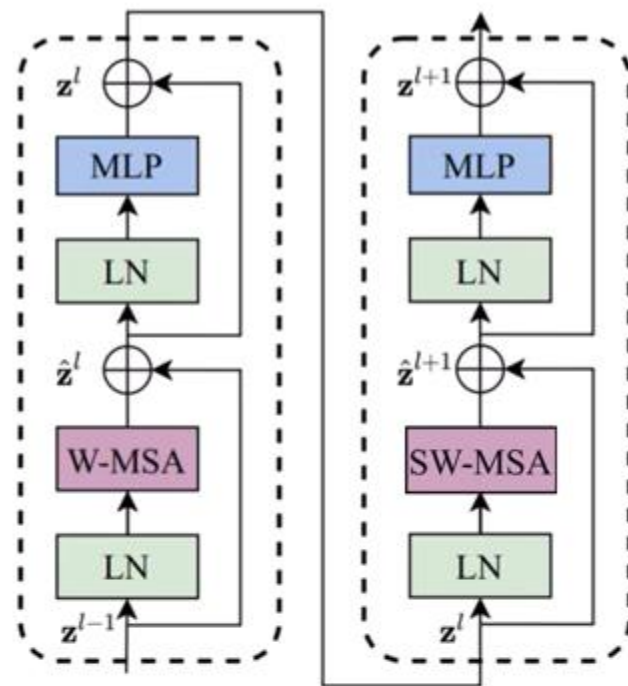
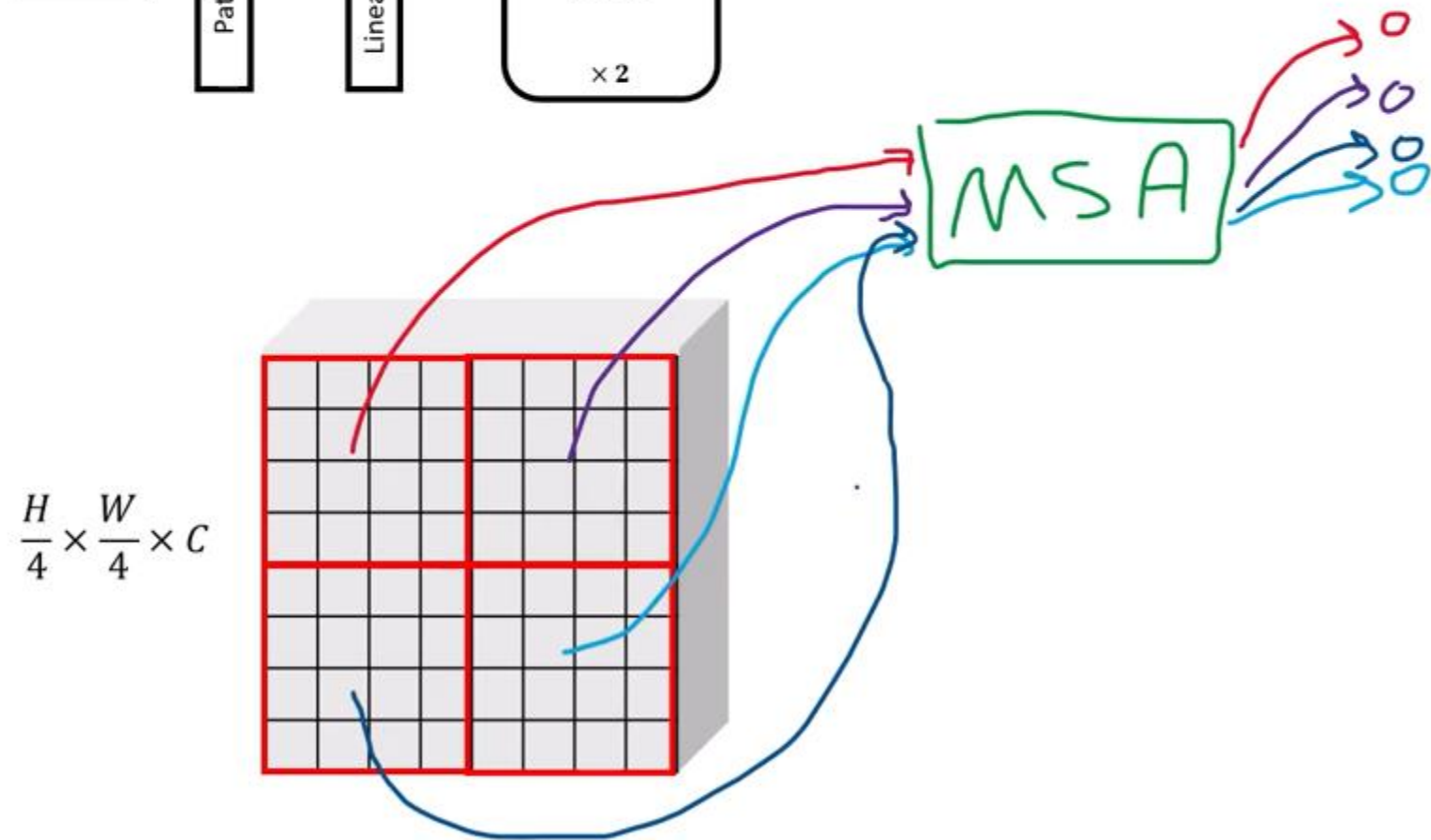
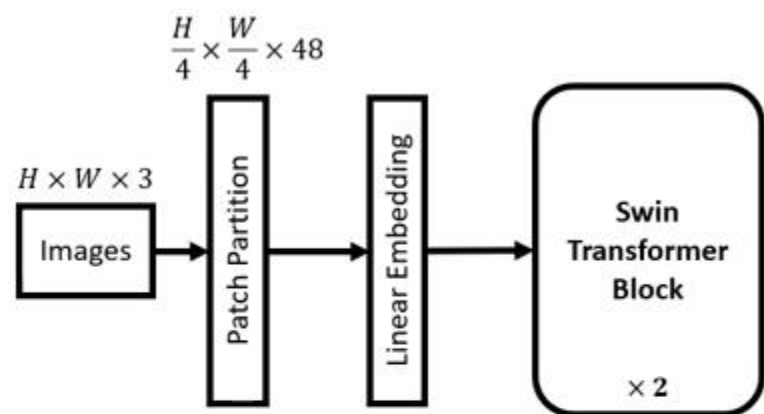
Figure 1. (a) The proposed Swin Transformer builds hierarchical feature maps by merging image patches (shown in gray) in deeper layers and has linear computation complexity to input image size due to computation of self-attention only within each local window (shown in red). It can thus serve as a general-purpose back-



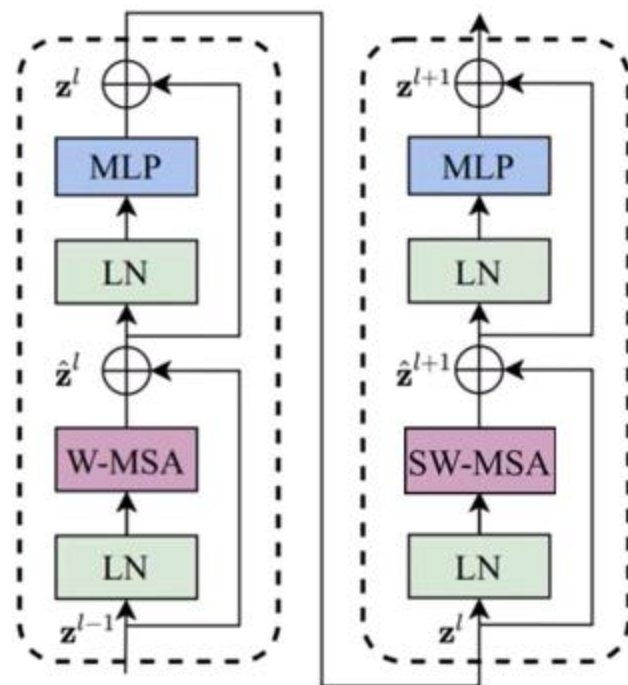
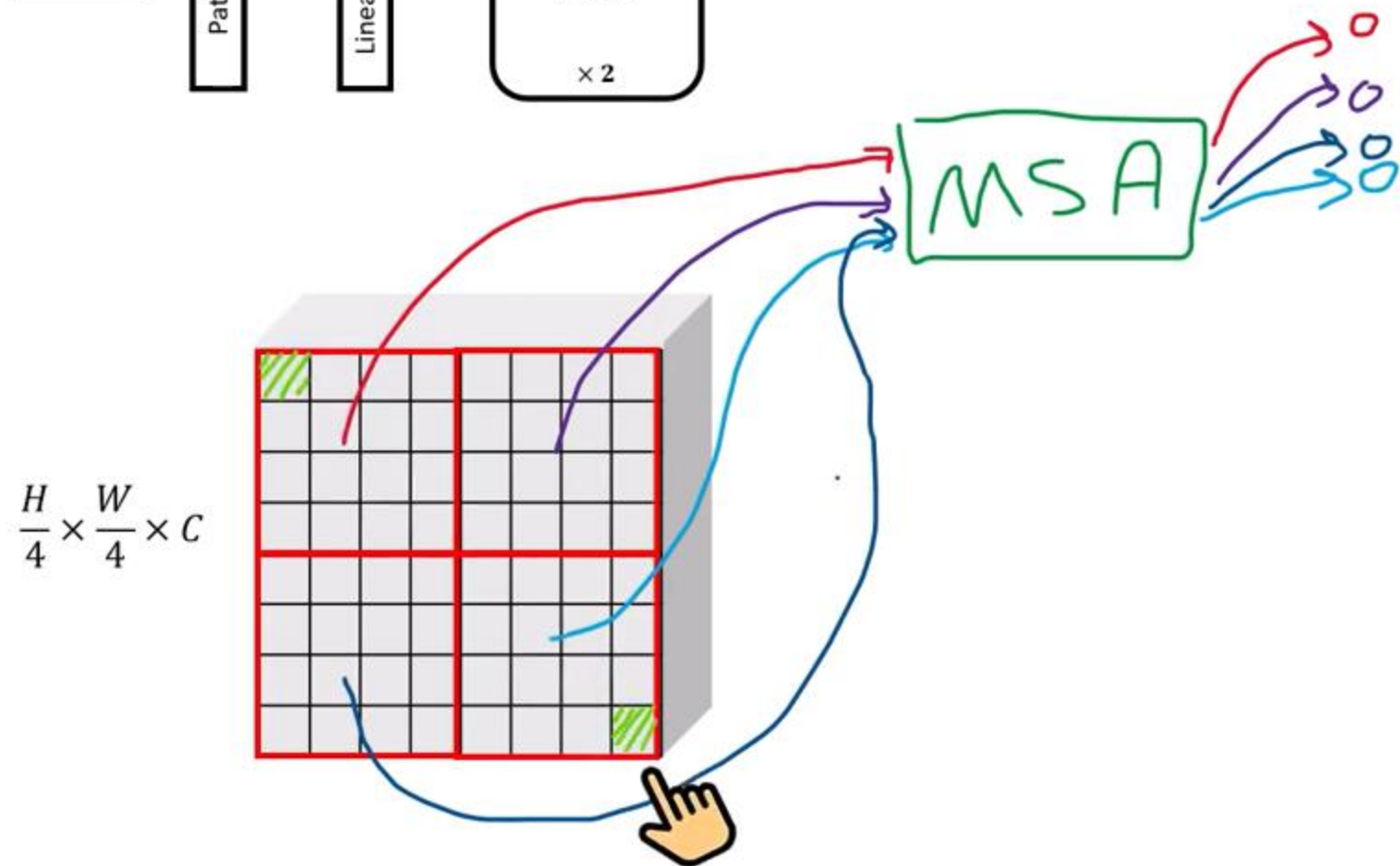
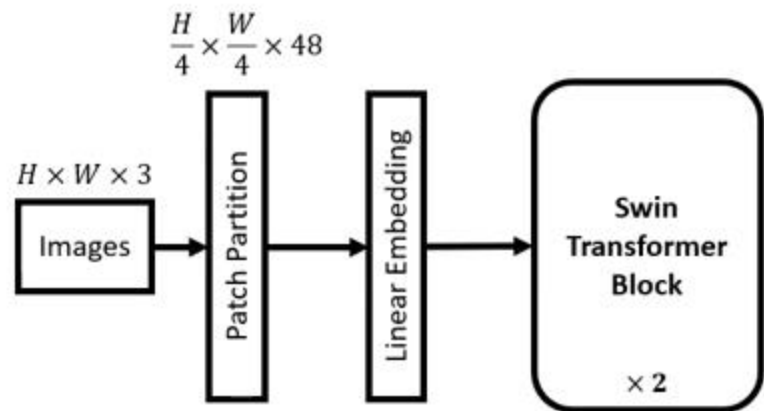




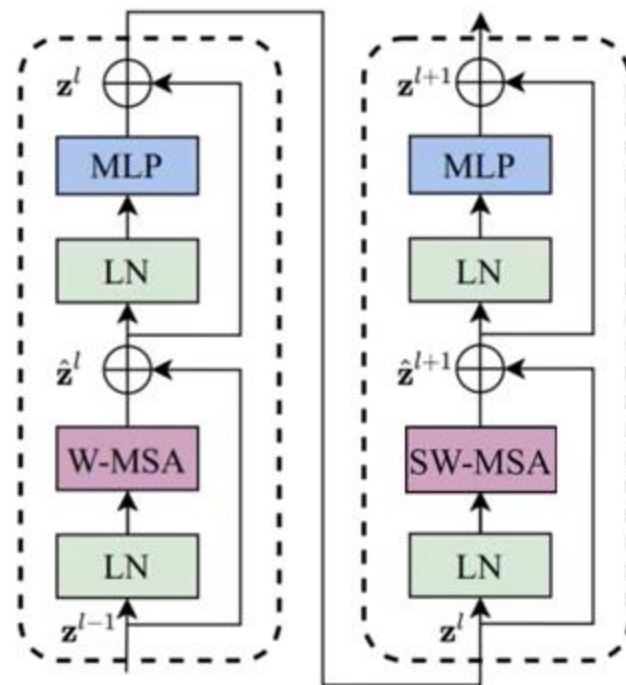
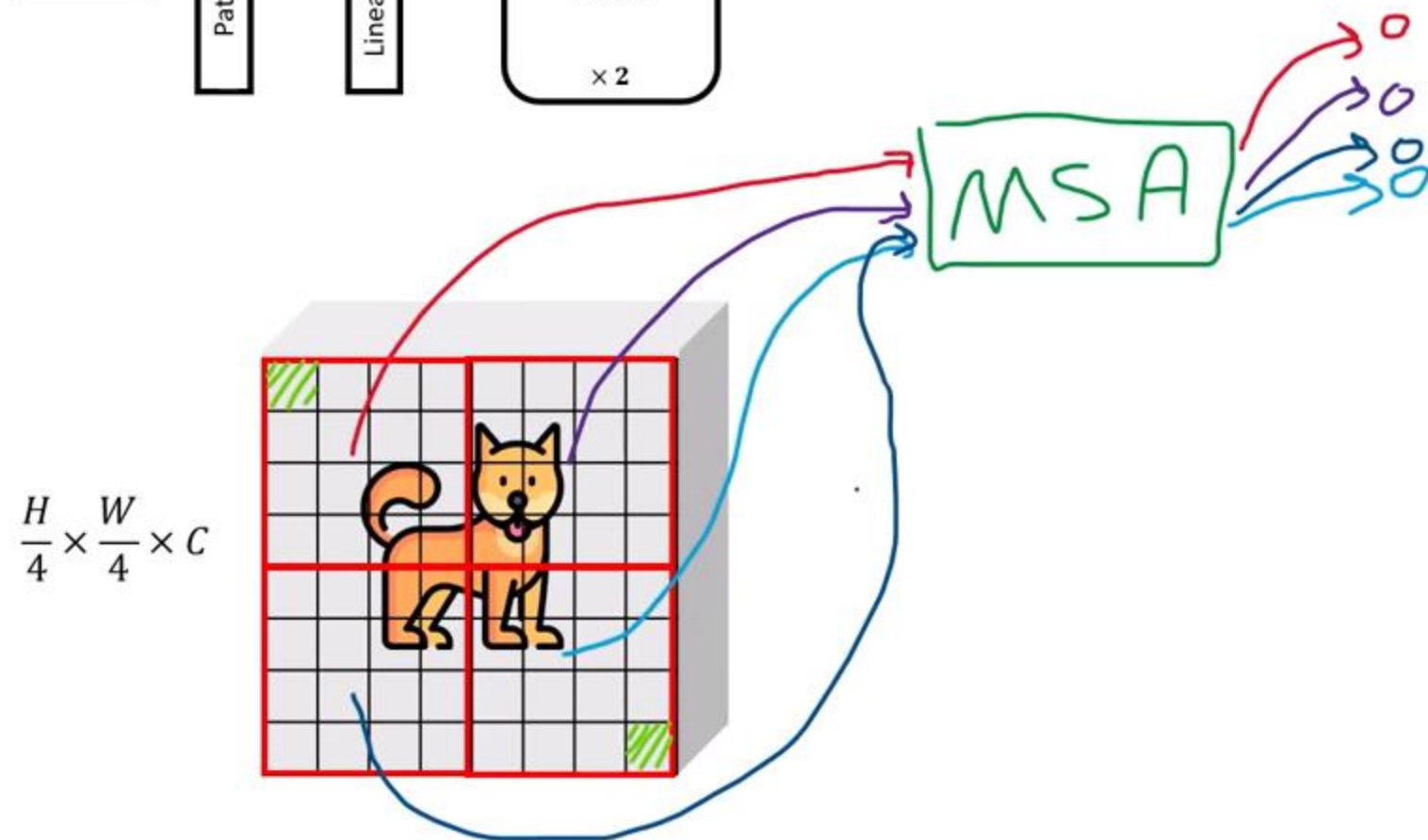
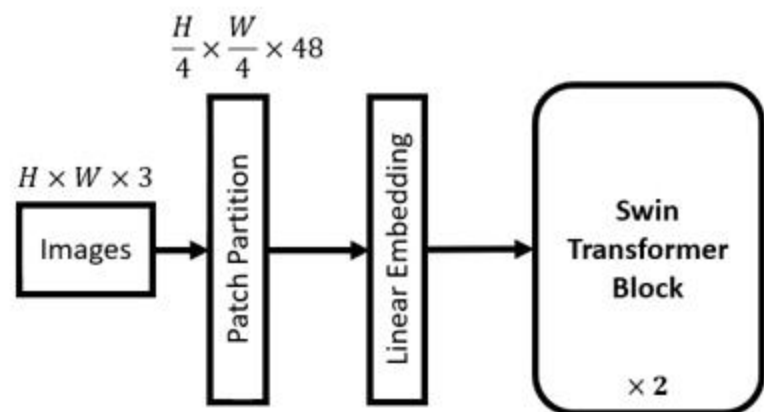




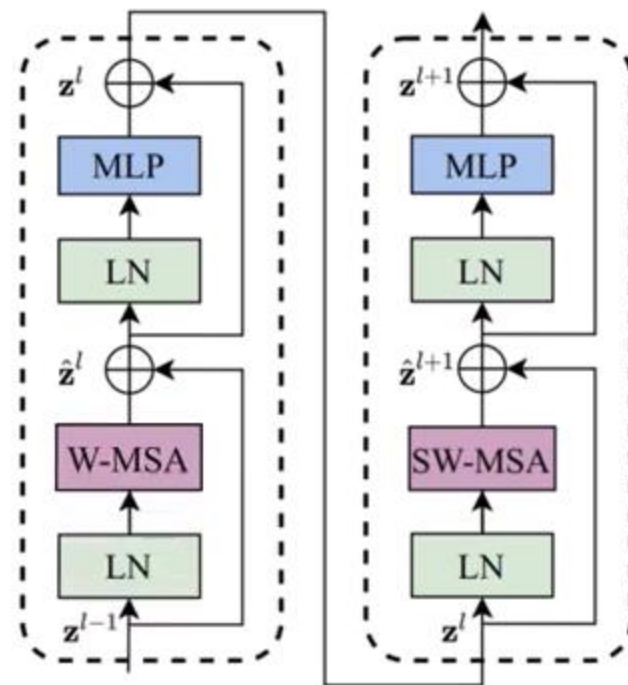
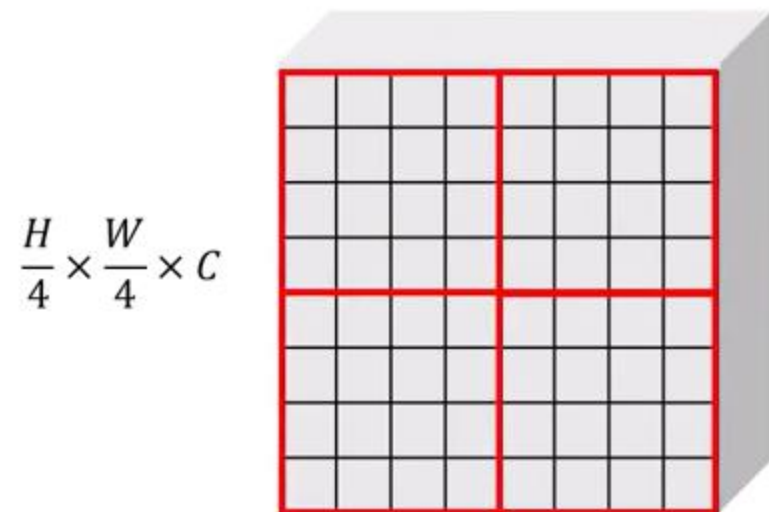
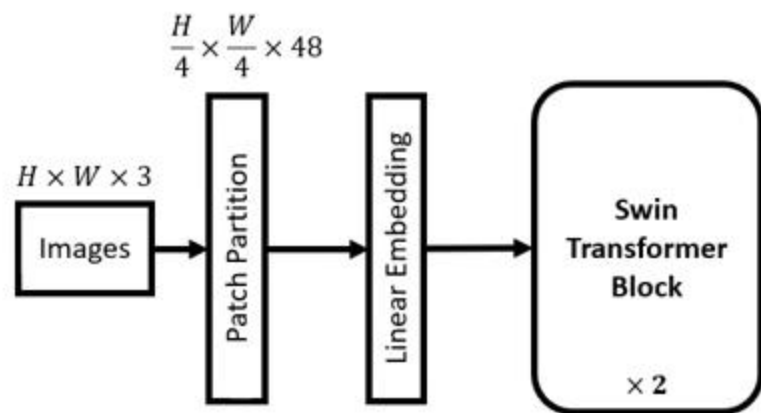
W-MSA



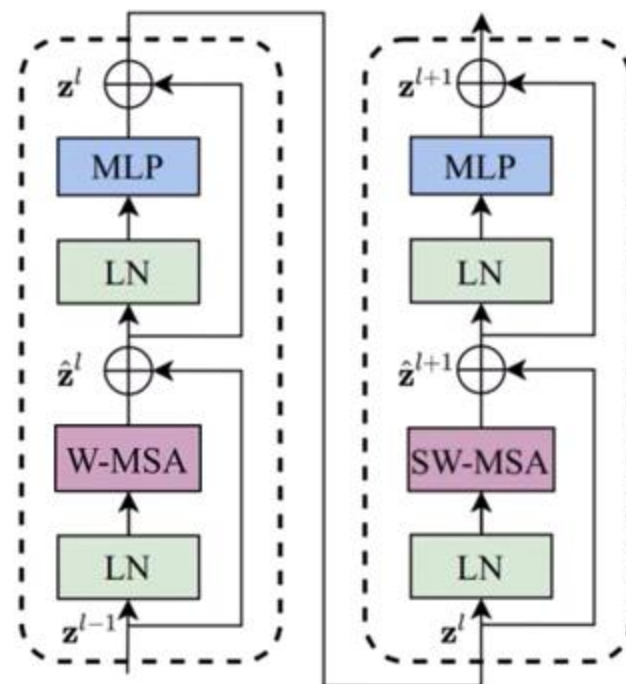
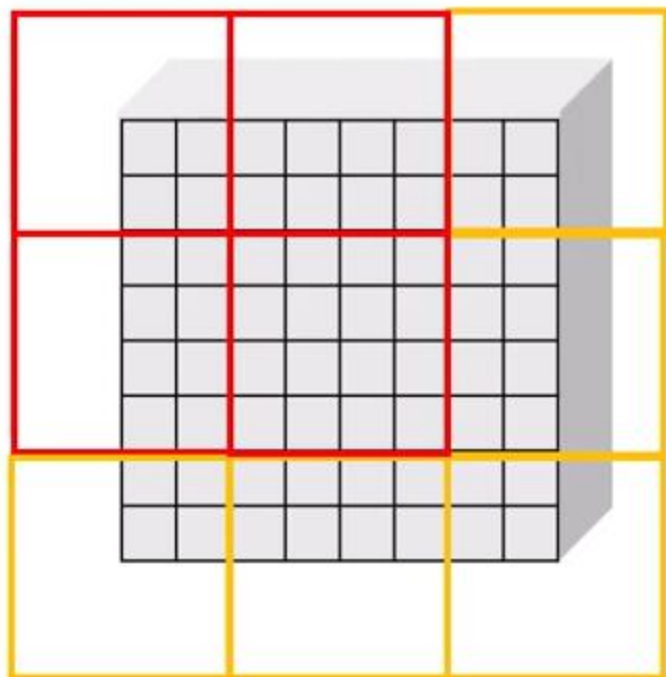
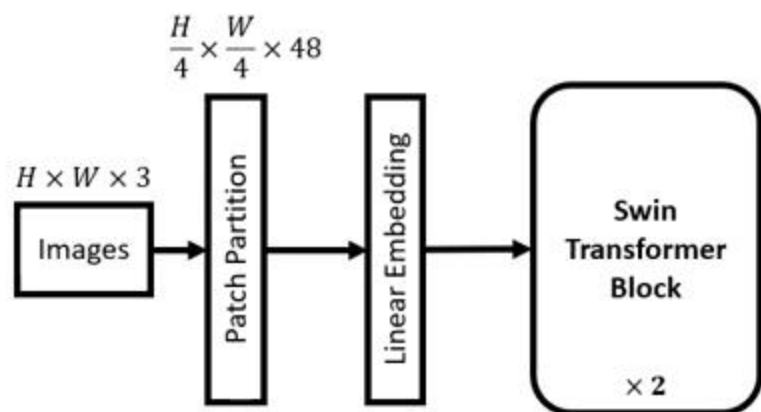
W-MSA



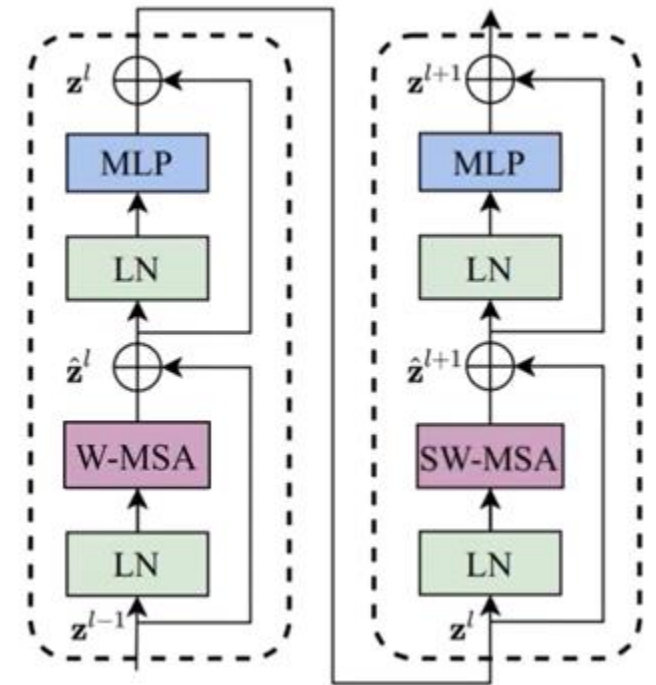
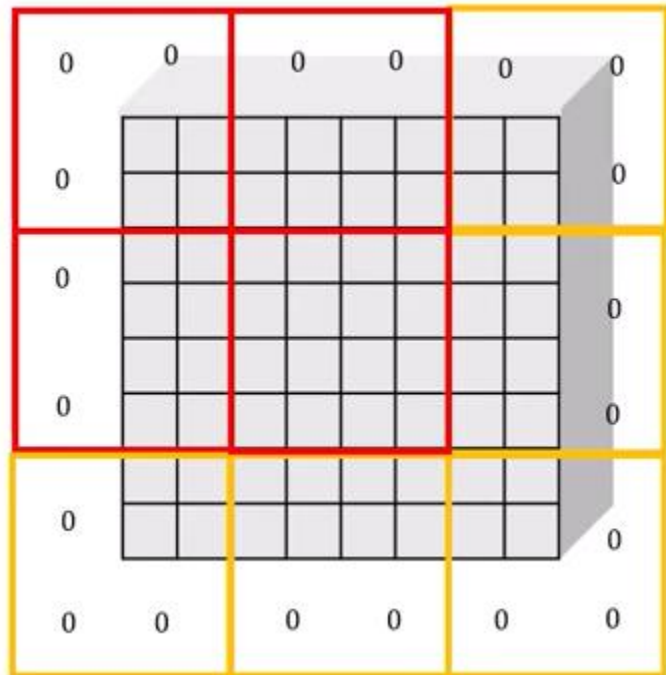
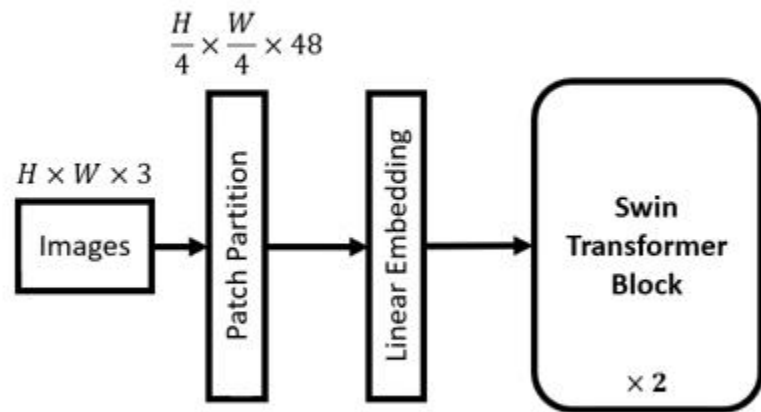
W-MSA



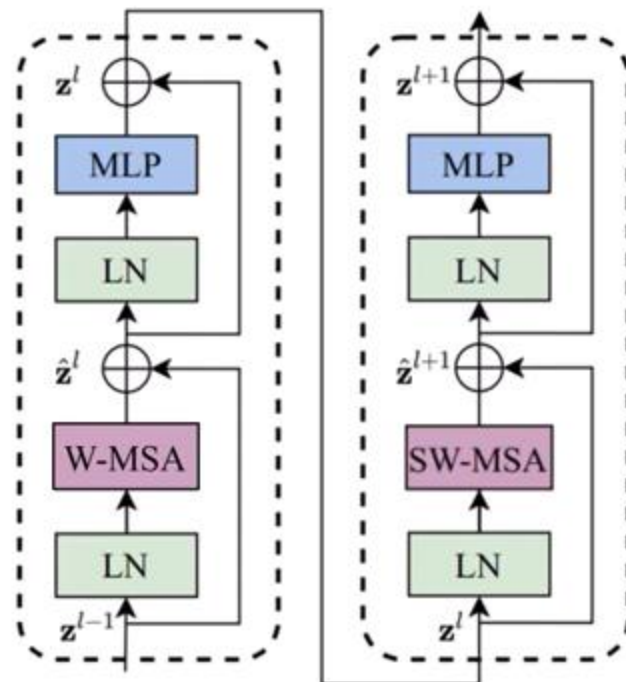
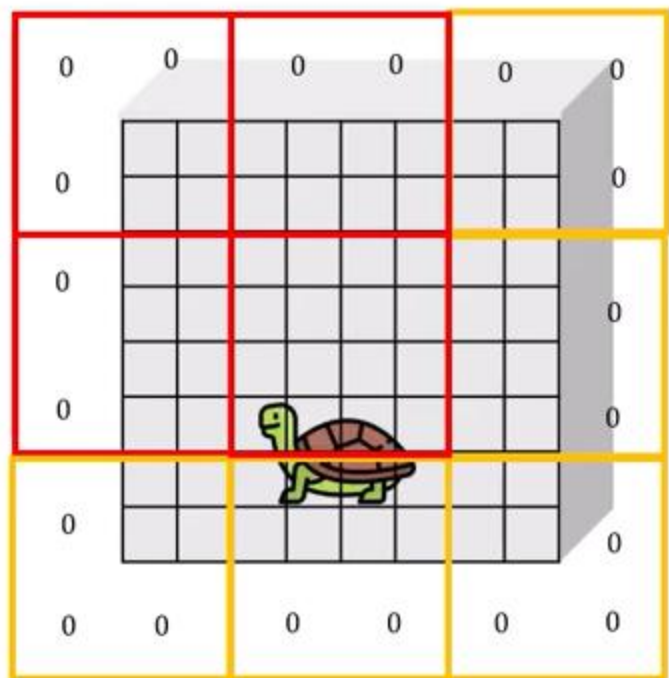
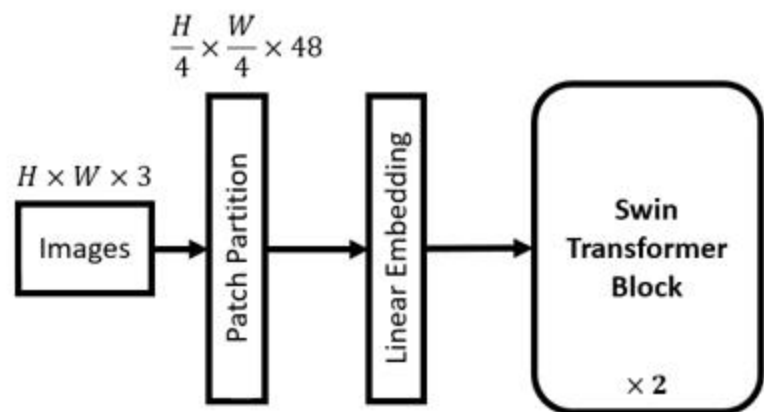
SW-MSA



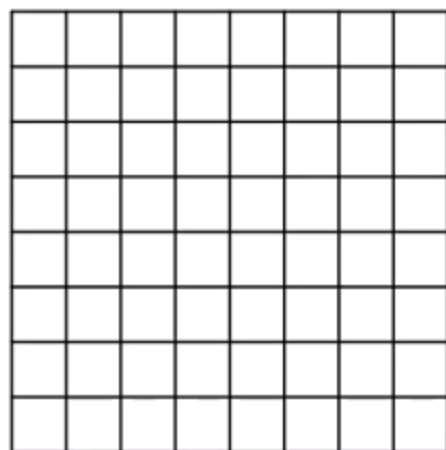
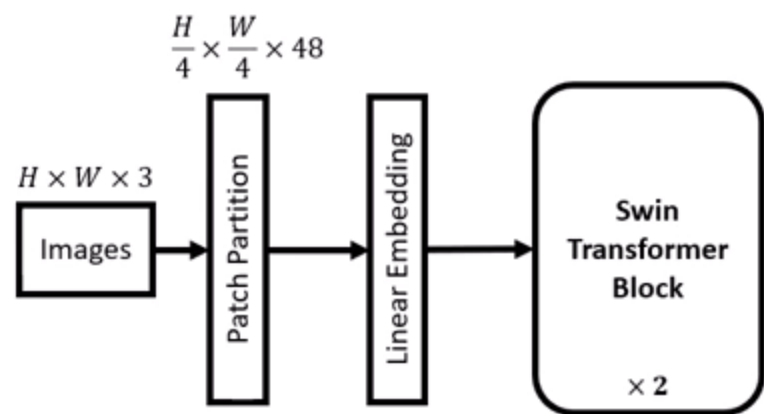
SW-MSA



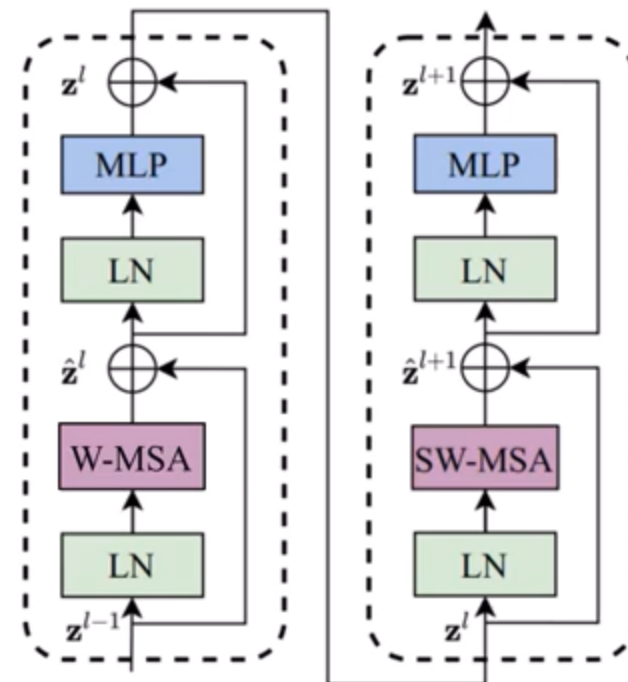
SW-MSA



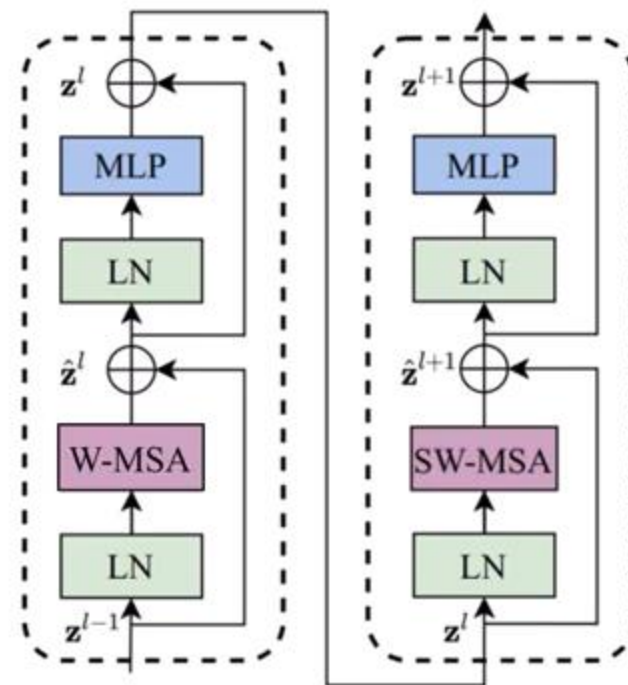
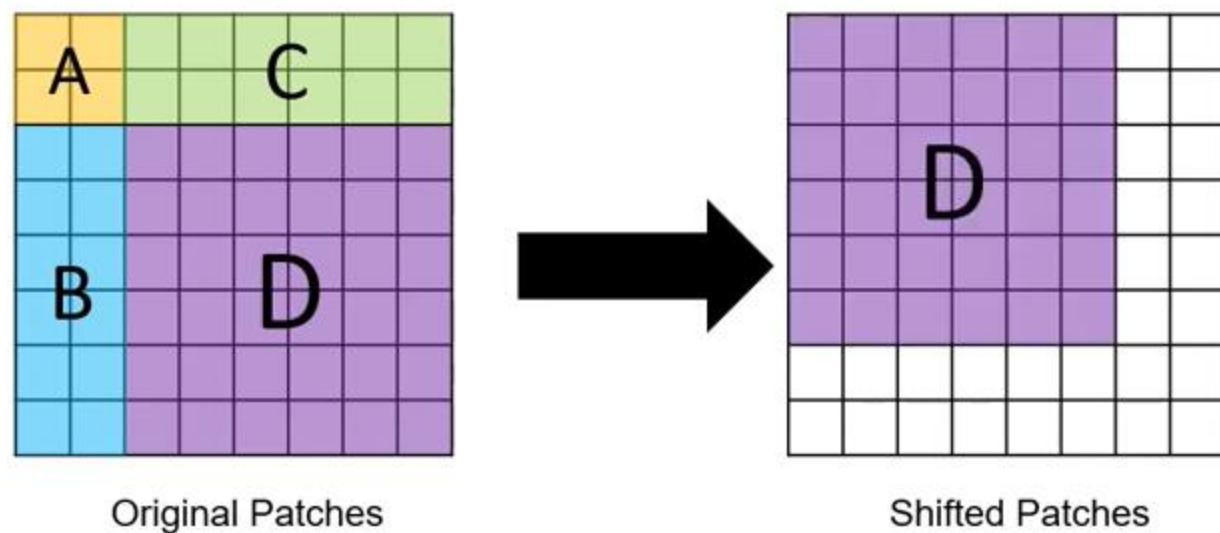
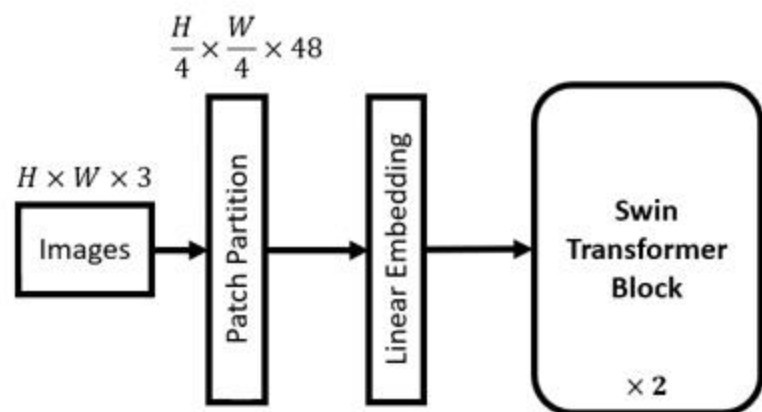
SW-MSA



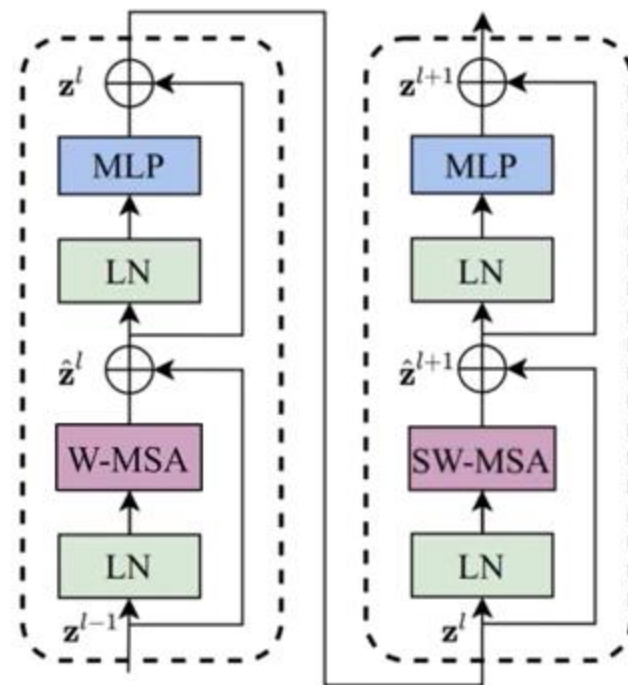
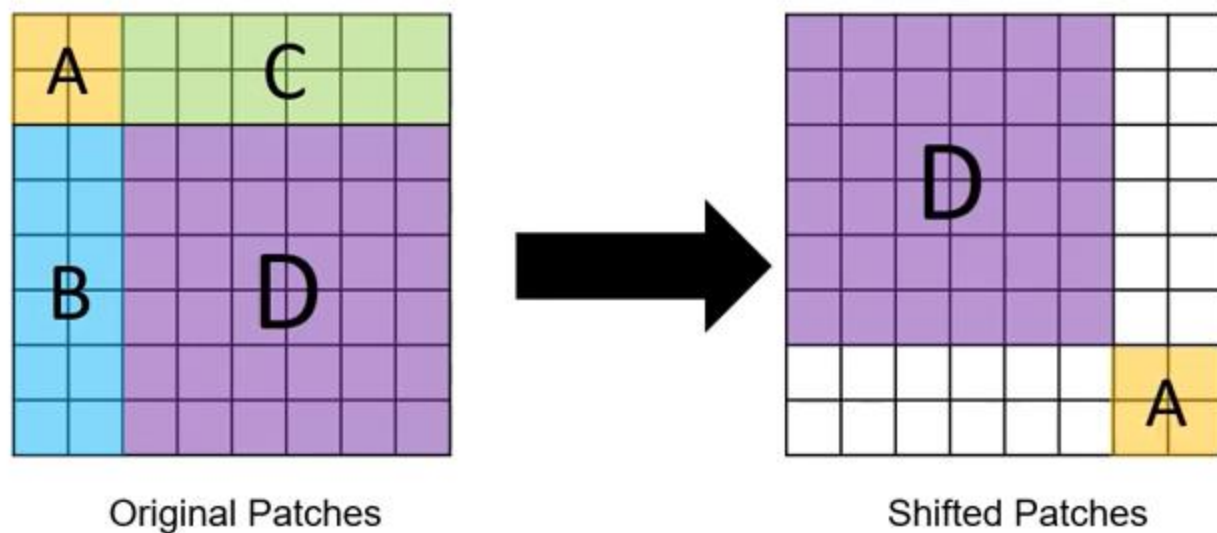
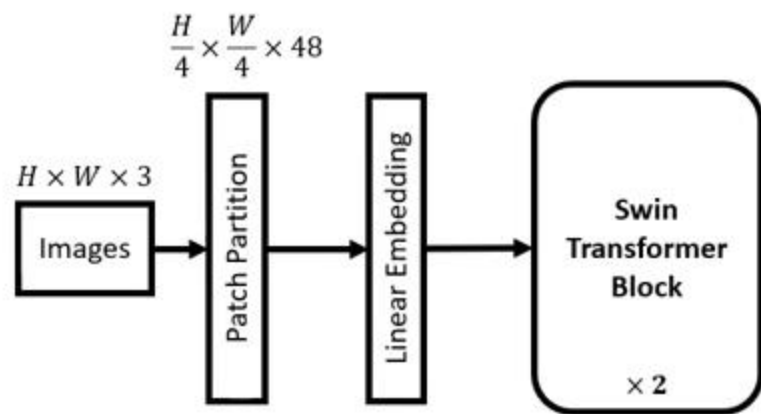
Original Patches



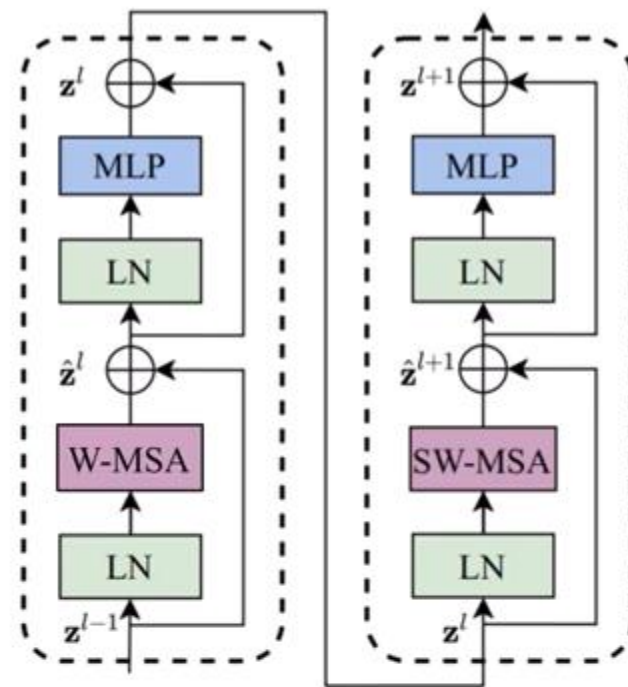
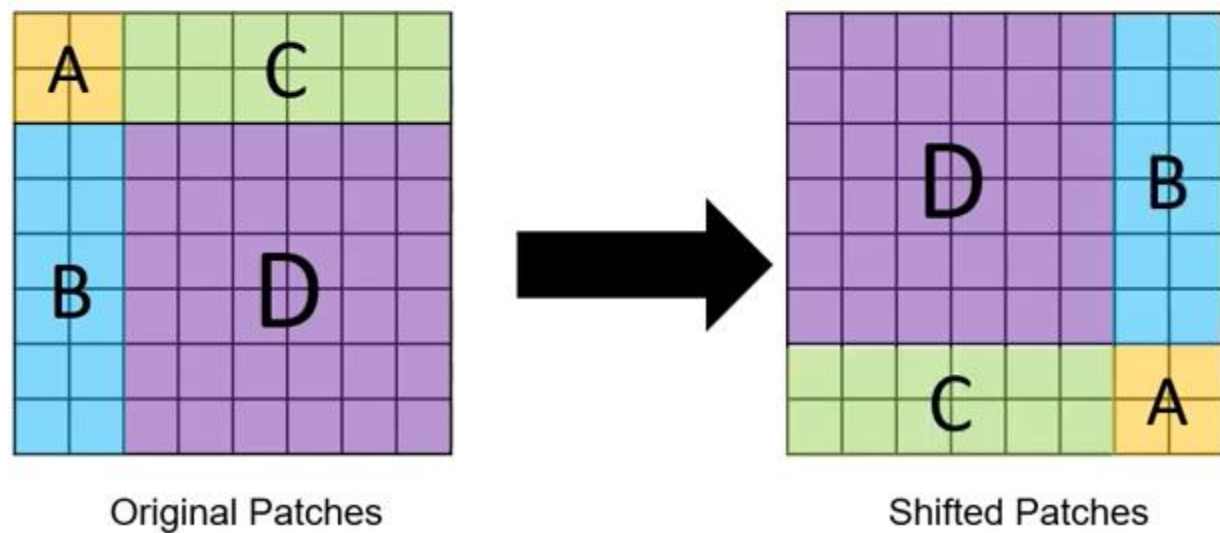
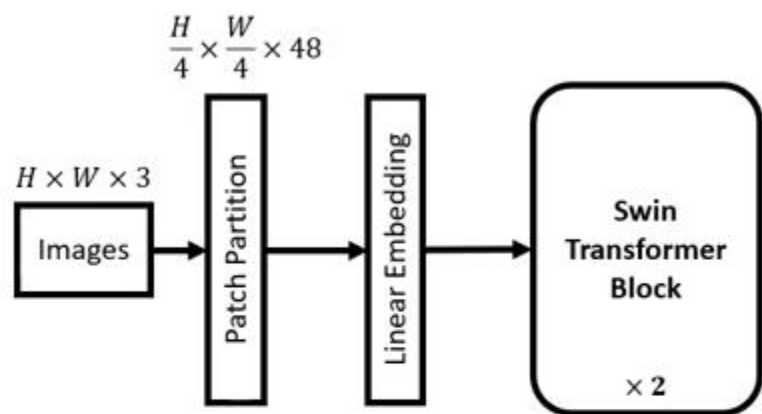
SW-MSA



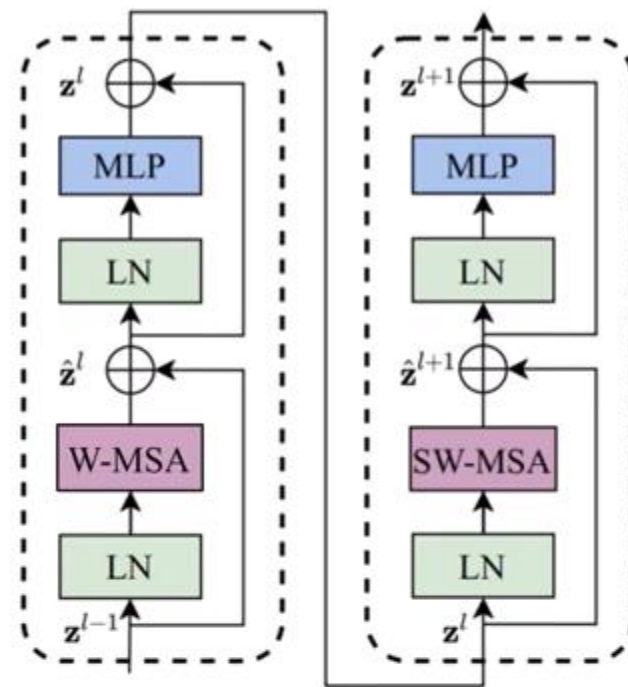
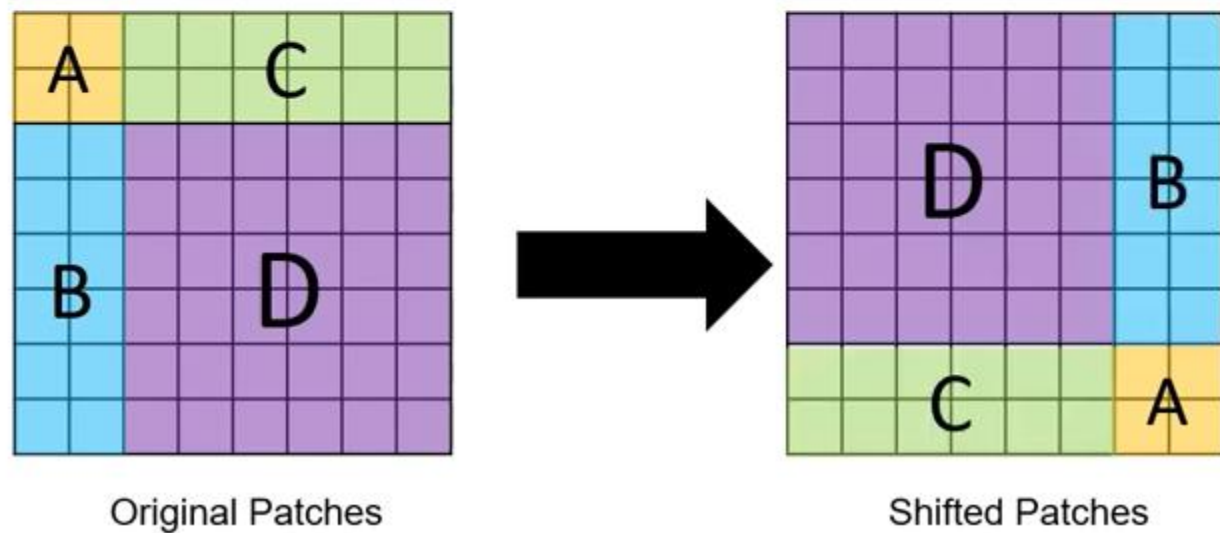
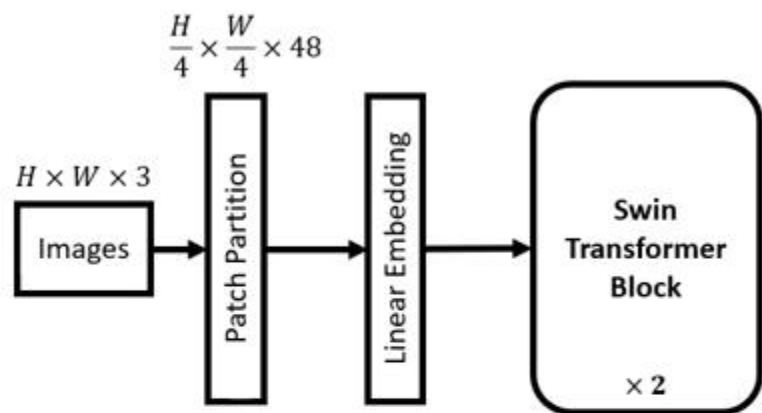
SW-MSA



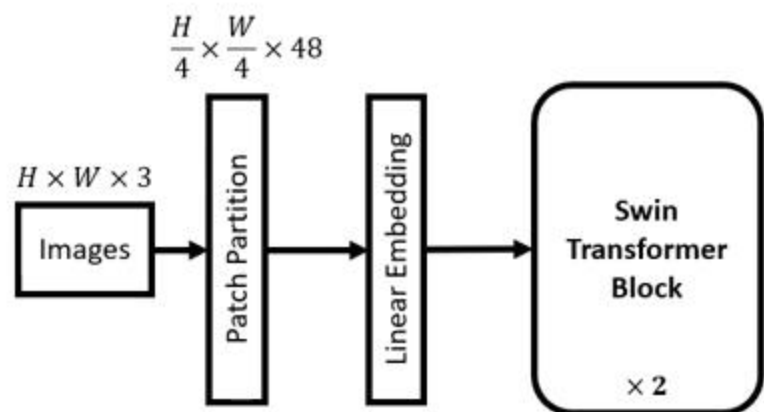
SW-MSA



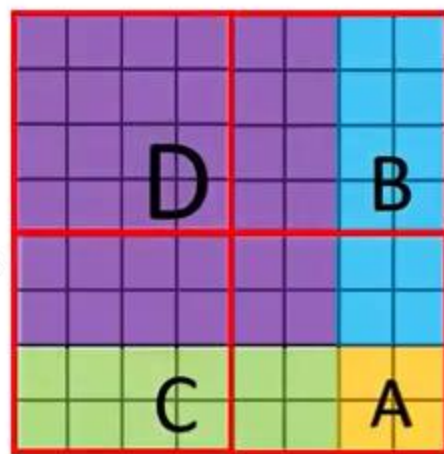
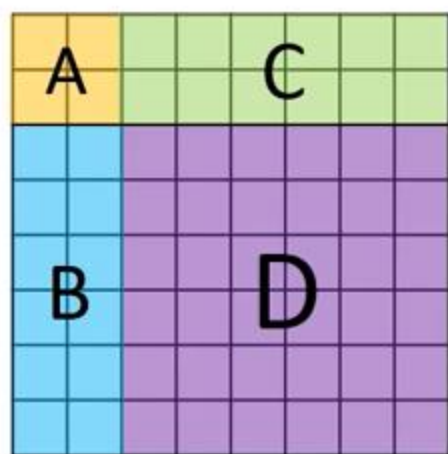
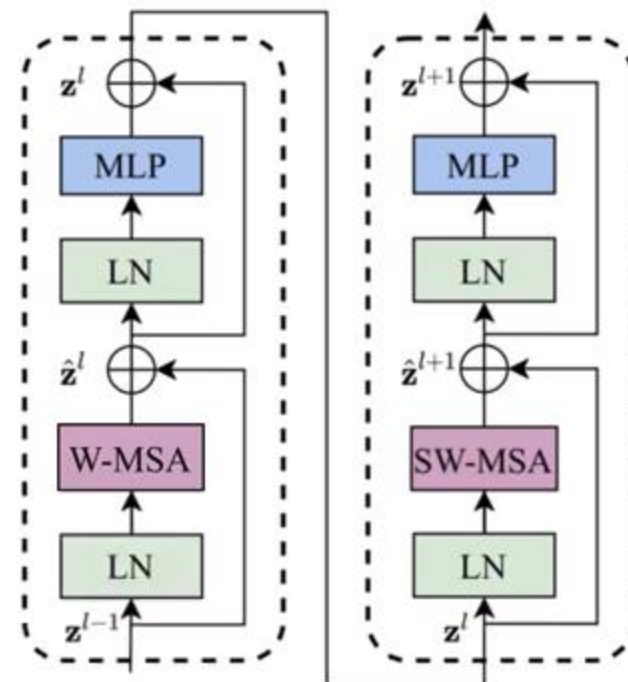
SW-MSA



SW-MSA

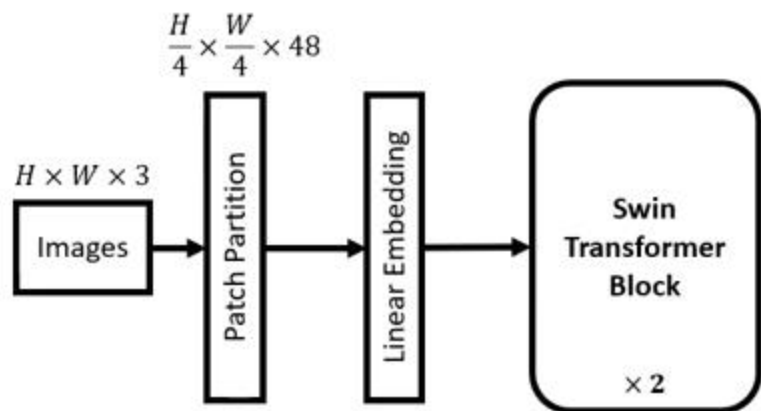


SW-MSA

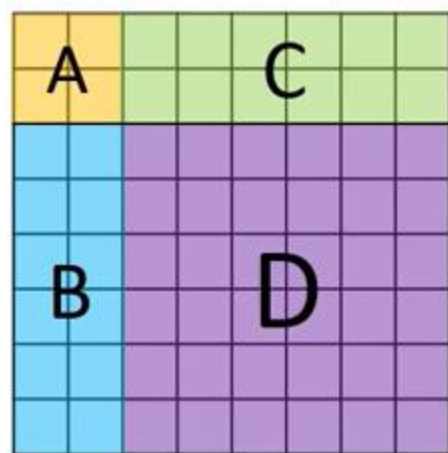
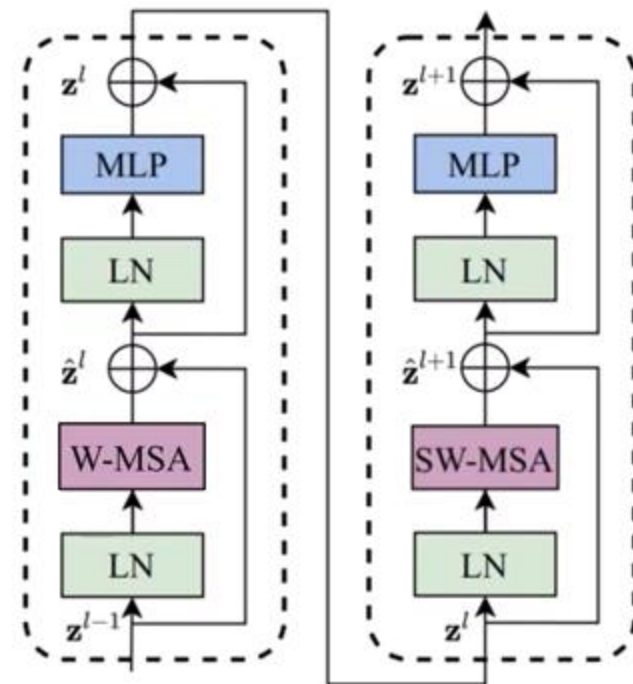


Original Patches

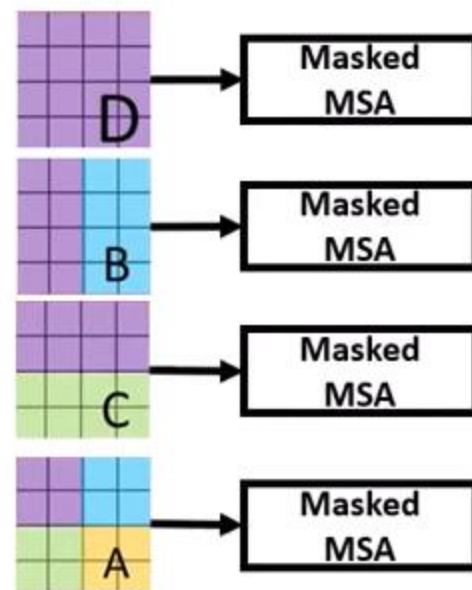
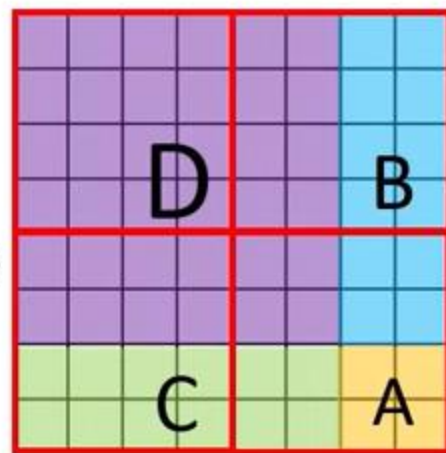
Shifted Patches



SW-MSA

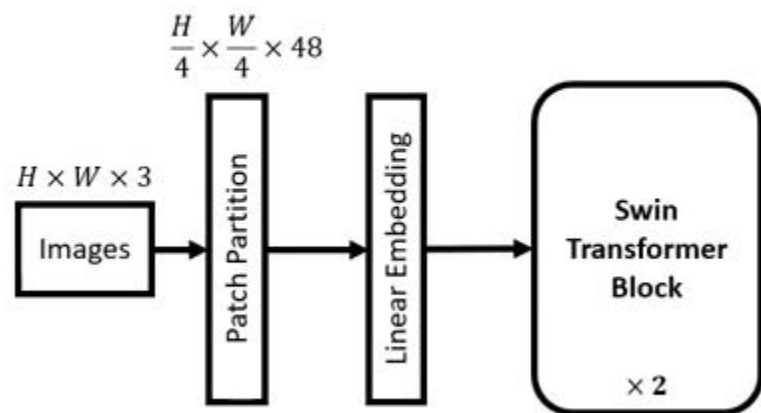


Cyclic Shift

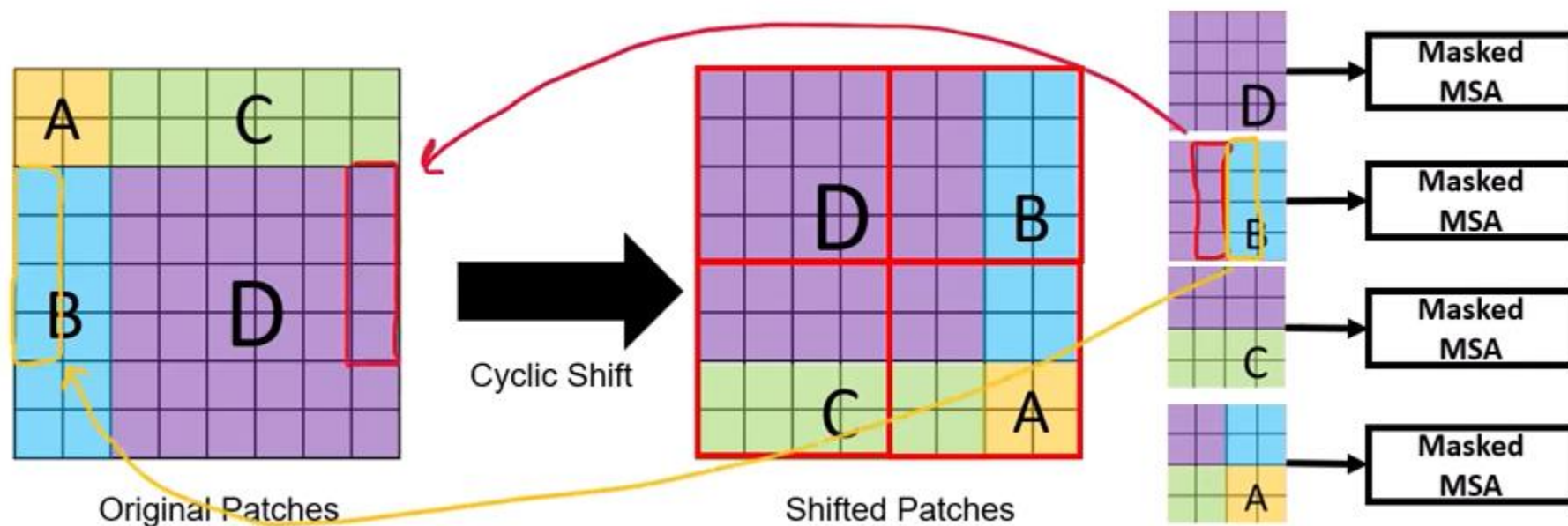
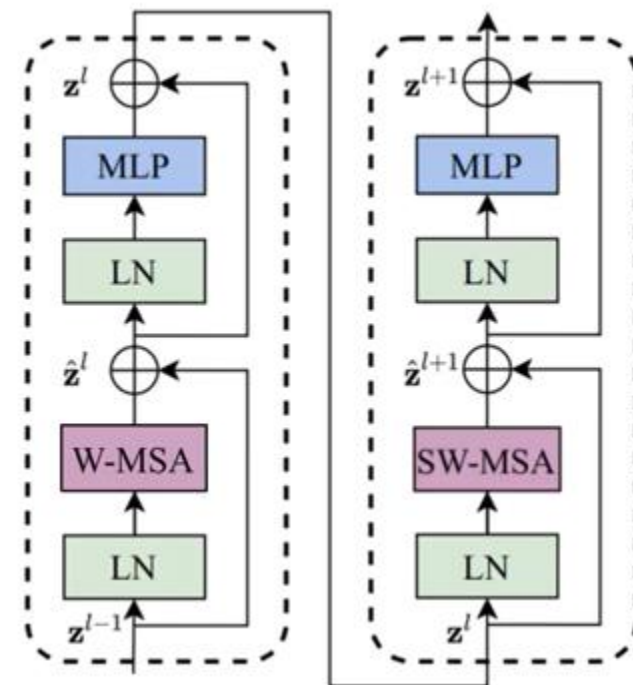


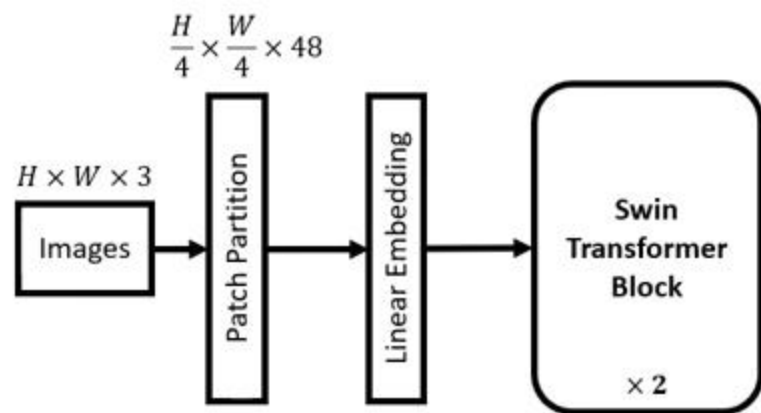
Original Patches

Shifted Patches

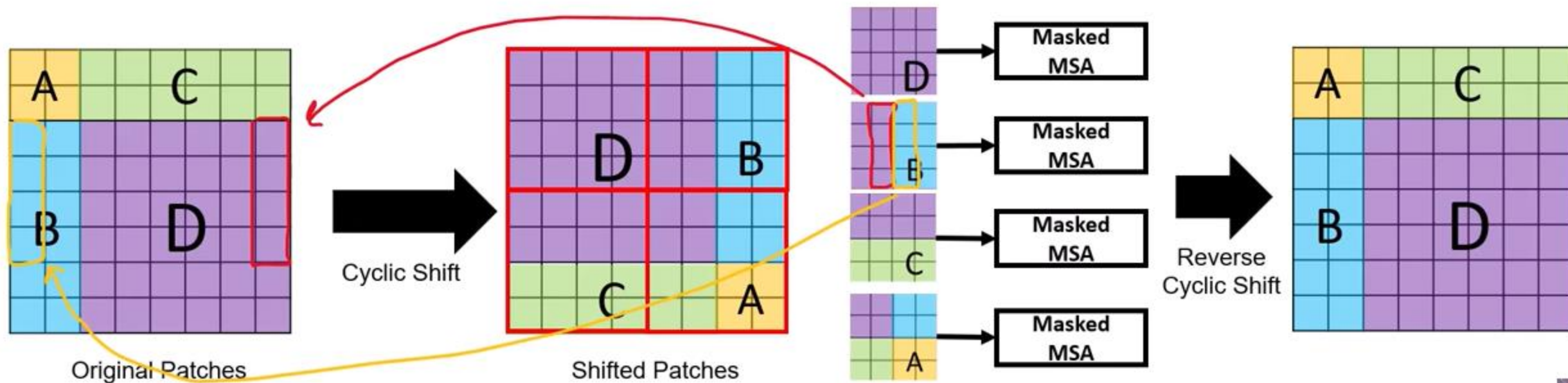
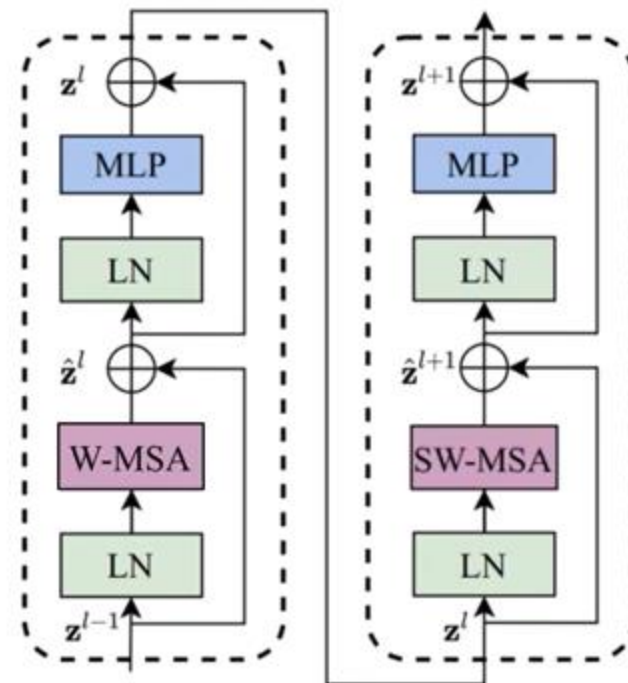


SW-MSA





SW-MSA



```

222     if self.shift_size > 0:
223         # calculate attention mask for SW-MSA
224         H, W = self.input_resolution
225         img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226         h_slices = (slice(0, -self.window_size),
227                     slice(-self.window_size, -self.shift_size),
228                     slice(-self.shift_size, None))
229         w_slices = (slice(0, -self.window_size),
230                     slice(-self.window_size, -self.shift_size),
231                     slice(-self.shift_size, None))
232         cnt = 0
233         for h in h_slices:
234             for w in w_slices:
235                 img_mask[:, h, w, :] = cnt
236                 cnt += 1
237
238         mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239         mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240         attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241         attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242     else:
243         attn_mask = None

```

```

222     if self.shift_size > 0:
223         # calculate attention mask for SW-MSA
224         H, W = self.input_resolution
225         img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226         h_slices = (slice(0, -self.window_size), [0: -4]
227                     slice(-self.window_size, -self.shift_size), [-4: -2]
228                     slice(-self.shift_size, None))[-2: ]
229         w_slices = (slice(0, -self.window_size), [0: -4]
230                     slice(-self.window_size, -self.shift_size), [-4: -2]
231                     slice(-self.shift_size, None))[-2: ]
232         cnt = 0
233         for h in h_slices:
234             for w in w_slices:
235                 img_mask[:, h, w, :] = cnt
236                 cnt += 1
237
238         mask_windows = window_partition(img_mask, self.window_size) # nW, window_size, window_size, 1
239         mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240         attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241         attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242     else:
243         attn_mask = None

```

0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2


```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

0	0	0	0	1	1	0	0
0	0	0	0	1	1	0	0
0	0	0	0	1	1	0	0
0	0	0	0	1	1	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nW, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
3	3	3	3	0	0	0	0
3	3	3	3	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222     if self.shift_size > 0:
223         # calculate attention mask for SW-MSA
224         H, W = self.input_resolution
225         img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226         h_slices = (slice(0, -self.window_size), [0: -4]
227                     slice(-self.window_size, -self.shift_size), [-4: -2]
228                     slice(-self.shift_size, None))[-2: ]
229         w_slices = (slice(0, -self.window_size), [0: -4]
230                     slice(-self.window_size, -self.shift_size), [-4: -2]
231                     slice(-self.shift_size, None))[-2: ]
232         cnt = 0
233         for h in h_slices:
234             for w in w_slices:
235                 img_mask[:, h, w, :] = cnt
236                 cnt += 1
237
238         mask_windows = window_partition(img_mask, self.window_size) # nW, window_size, window_size, 1
239         mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240         attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241         attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242     else:
243         attn_mask = None

```

0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
3	3	3	3	4	4	0	0
3	3	3	3	4	4	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nW, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
3	3	3	3	4	4	5	5
3	3	3	3	4	4	5	5
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2


```

222     if self.shift_size > 0:
223         # calculate attention mask for SW-MSA
224         H, W = self.input_resolution
225         img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226         h_slices = (slice(0, -self.window_size), [0: -4]
227                     slice(-self.window_size, -self.shift_size), [-4: -2]
228                     slice(-self.shift_size, None)) [-2: ]
229         w_slices = (slice(0, -self.window_size), [0: -4]
230                     slice(-self.window_size, -self.shift_size), [-4: -2]
231                     slice(-self.shift_size, None)) [-2: ]
232         cnt = 0
233         for h in h_slices:
234             for w in w_slices:
235                 img_mask[:, h, w, :] = cnt
236                 cnt += 1
237
238         mask_windows = window_partition(img_mask, self.window_size) # nW, window_size, window_size, 1
239         mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240         attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241         attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242     else:
243         attn_mask = None

```

0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
0	0	0	0	1	1	2	2
3	3	3	3	4	4	5	5
3	3	3	3	4	4	5	5
6	6	6	6	0	0	0	0
6	6	6	6	0	0	0	0

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None))[-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None))[-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nW, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

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

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

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

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

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

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```



4 4 4

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

D

B

C

A

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2


```

222     if self.shift_size > 0:
223         # calculate attention mask for SW-MSA
224         H, W = self.input_resolution
225         img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226         h_slices = (slice(0, -self.window_size), [0: -4]
227                     slice(-self.window_size, -self.shift_size), [-4: -2]
228                     slice(-self.shift_size, None)) [-2: ]
229         w_slices = (slice(0, -self.window_size), [0: -4]
230                     slice(-self.window_size, -self.shift_size), [-4: -2]
231                     slice(-self.shift_size, None)) [-2: ]
232         cnt = 0
233         for h in h_slices:
234             for w in w_slices:
235                 img_mask[:, h, w, :] = cnt
236                 cnt += 1
237
238         mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239         mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240         attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241         attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242     else:
243         attn_mask = None

```



4 4 4

→ [4, 16]

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

D

B

C

A

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None)) [-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None)) [-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```

Handwritten annotations on the code:

- A hand icon points to line 240.
- Arrows point from the following lists to line 239:
 - $[4, 1, 16]$ (with a red '4' above it)
 - $[4, 16, 1]$ (with a green '4' above it)
 - $[4, 16]$ (with a green '4' above it)

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

Labels on the table:

- D** is in the top-left 4x4 purple region.
- B** is in the top-right 4x4 blue region.
- C** is in the bottom-left 4x4 green region.
- A** is in the bottom-right 4x4 yellow region.

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2


```

222 if self.shift_size > 0:
223     # calculate attention mask for SW-MSA
224     H, W = self.input_resolution
225     img_mask = torch.zeros((1, H, W, 1)) # 1 H W 1
226     h_slices = (slice(0, -self.window_size), [0: -4]
227                 slice(-self.window_size, -self.shift_size), [-4: -2]
228                 slice(-self.shift_size, None))[-2: ]
229     w_slices = (slice(0, -self.window_size), [0: -4]
230                 slice(-self.window_size, -self.shift_size), [-4: -2]
231                 slice(-self.shift_size, None))[-2: ]
232     cnt = 0
233     for h in h_slices:
234         for w in w_slices:
235             img_mask[:, h, w, :] = cnt
236             cnt += 1
237
238     mask_windows = window_partition(img_mask, self.window_size) # nw, window_size, window_size, 1
239     mask_windows = mask_windows.view(-1, self.window_size * self.window_size)
240     attn_mask = mask_windows.unsqueeze(1) - mask_windows.unsqueeze(2)
241     attn_mask = attn_mask.masked_fill(attn_mask != 0, float(-100.0)).masked_fill(attn_mask == 0, float(0.0))
242 else:
243     attn_mask = None

```



$[4, 16, 1]$
 $[4, 16]$
 $[4, 16]$
 Different
 Same

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

why?

Assumptions:

$H = W = 8$

window_size = 4

shift_size = 2

$$Attention(Q, K, V) = Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

+attn_mask

3. Apply the softmax function to each input:

$$\text{softmax}(-100) = \frac{e^{-100}}{1} \approx 3.72 \times 10^{-44}$$

$$\text{softmax}(0) = \frac{e^0}{1} = 1$$

4. Normalize the outputs so that they sum to 1:

$$\text{softmax}(-100) \approx \frac{3.72 \times 10^{-44}}{1 + 3 \times 3.72 \times 10^{-44}} = \frac{3.72 \times 10^{-44}}{1} \approx 3.72 \times 10^{-44}$$

$$\text{softmax}(0) \approx \frac{1}{1 + 3 \times 3.72 \times 10^{-44}} = \frac{1}{1} = 1$$

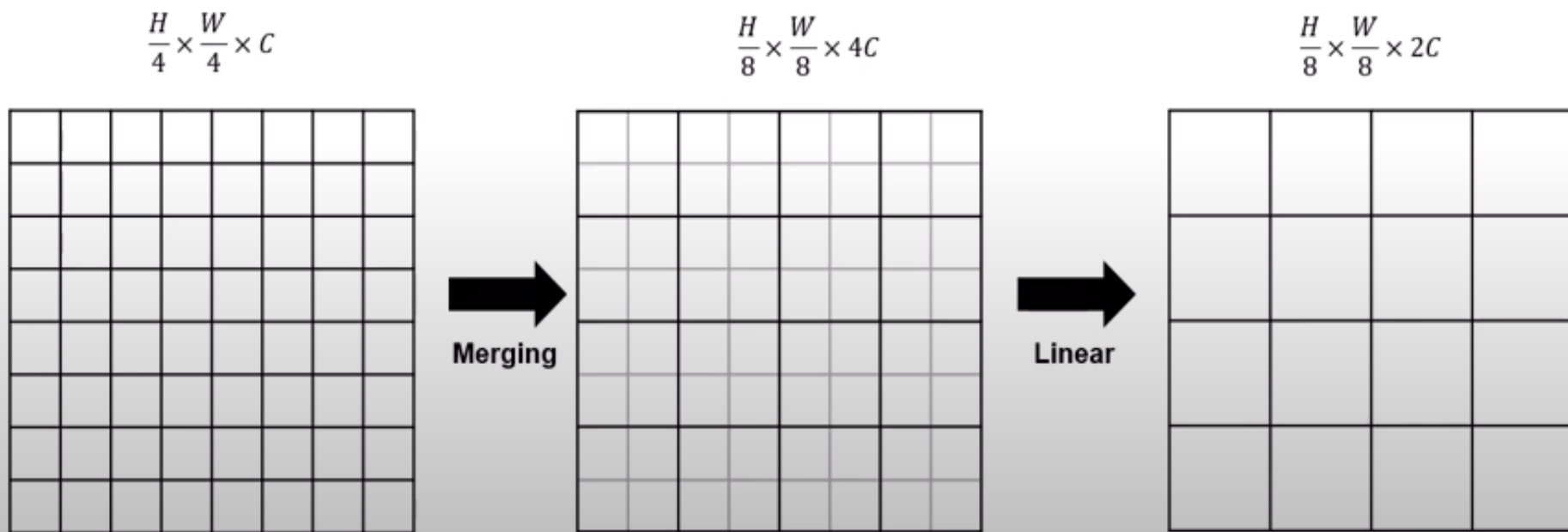
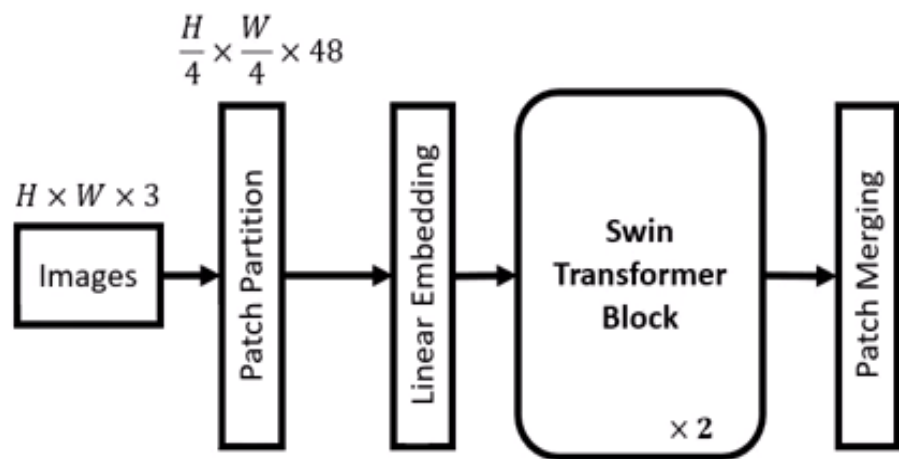
Since the terms 3.72×10^{-44} are exceedingly small compared to 1, the output for the softmax will effectively be:

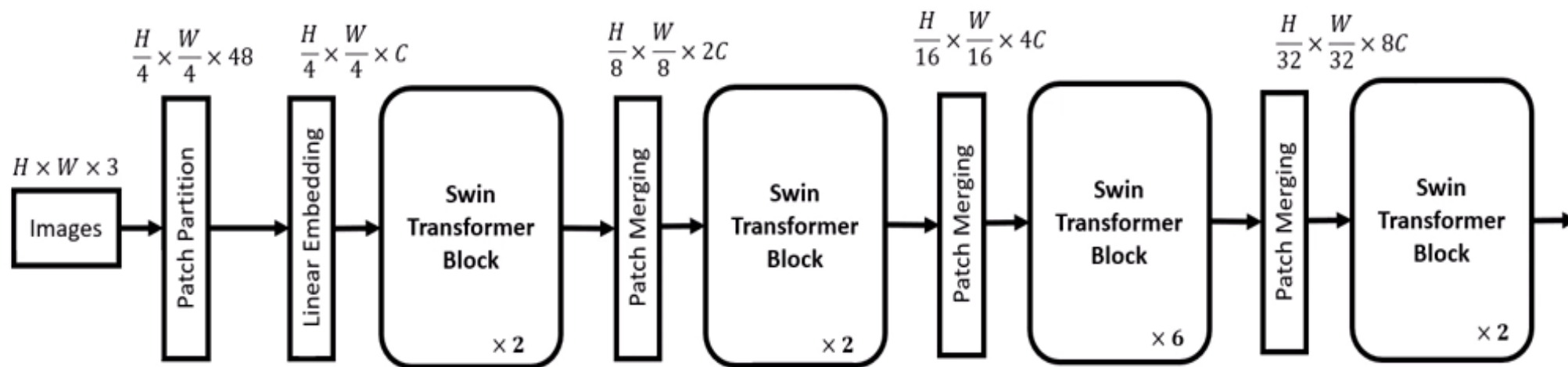
$$\text{softmax}([-100, -100, -100, 0]) = [3.72 \times 10^{-44}, 3.72 \times 10^{-44}, 3.72 \times 10^{-44}, 1]$$

When normalized, the values will sum to 1:

$$\text{softmax}([-100, -100, -100, 0]) \approx [0, 0, 0, 1]$$

This indicates that the probability is almost entirely concentrated on the last input, which is 0, as it is significantly larger than the others.





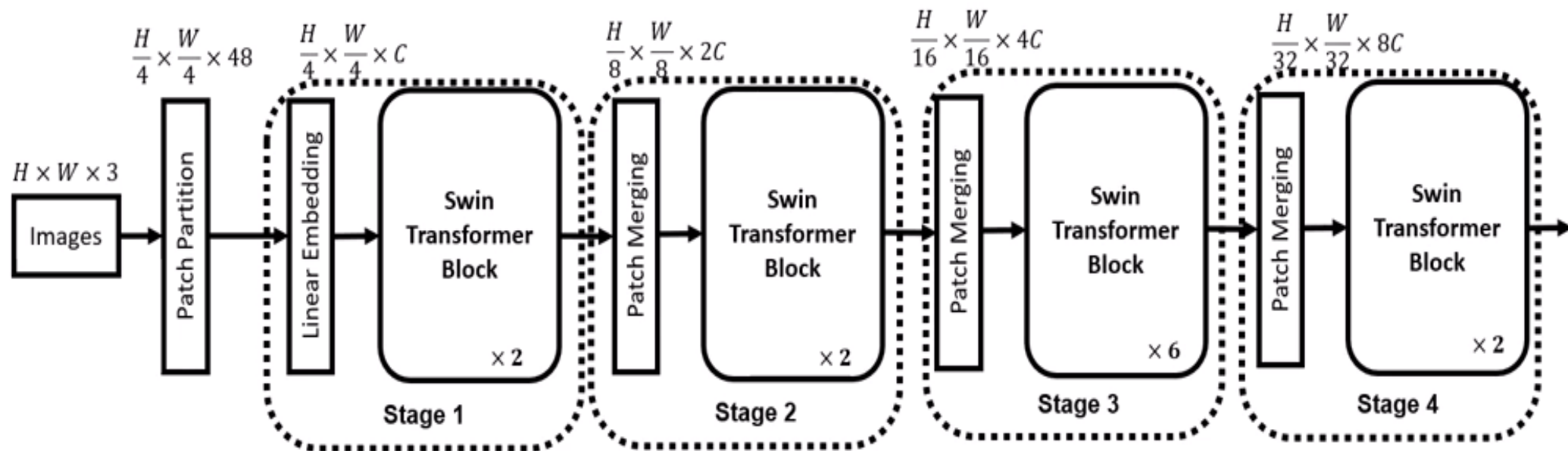


Image classification: Output of the last stage is passed through a linear layer.

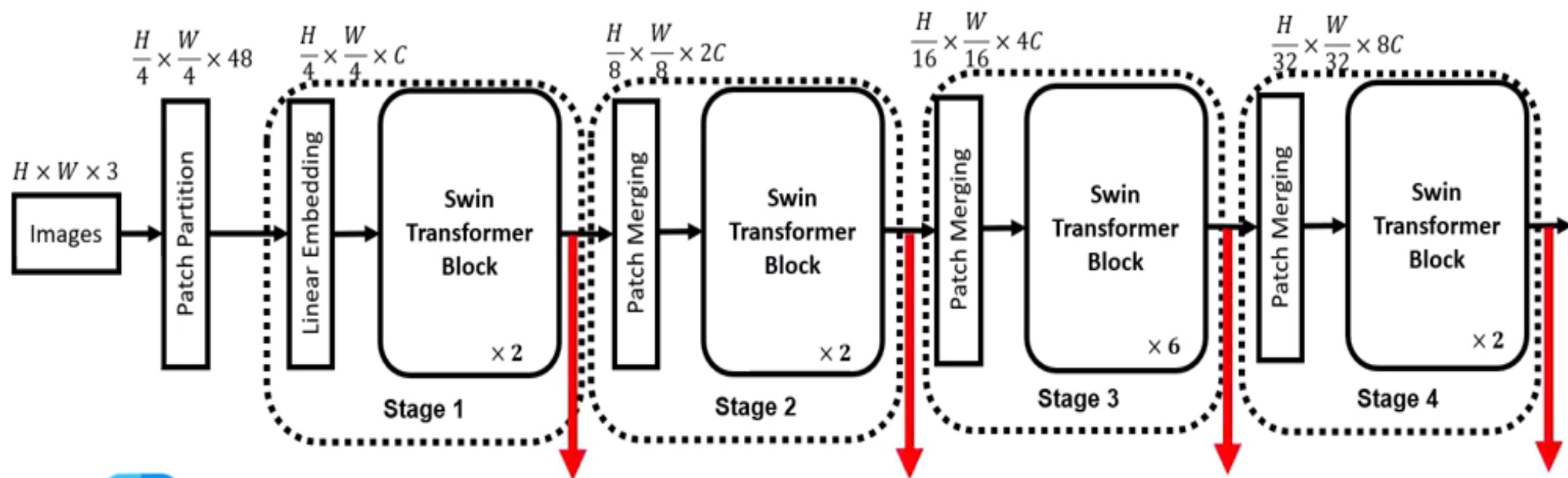


Image classification: Output of the last stage is passed through a linear layer.



Object detection and Image segmentation: Output of all the stages are used as features.

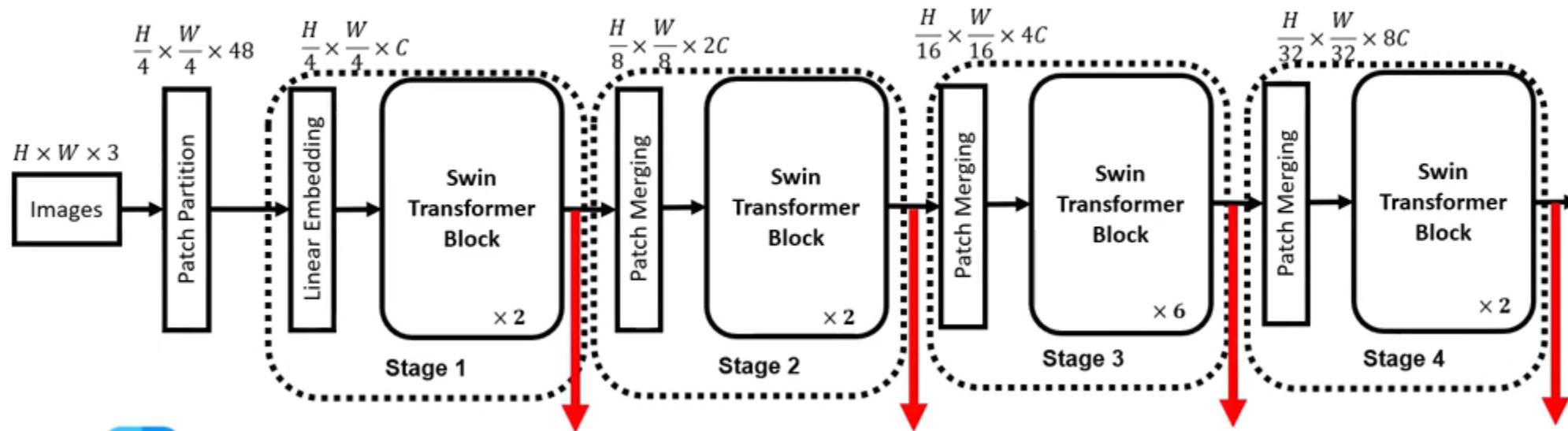
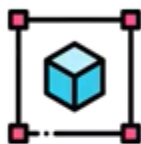


Image classification: Output of the last stage is passed through a linear layer.



Object detection and Image segmentation: Output of all the stages are used as features.

- Swin-T: $C = 96$, layer numbers = $\{2, 2, 6, 2\}$
- Swin-S: $C = 96$, layer numbers = $\{2, 2, 18, 2\}$
- Swin-B: $C = 128$, layer numbers = $\{2, 2, 18, 2\}$
- Swin-L: $C = 192$, layer numbers = $\{2, 2, 18, 2\}$

(a) Regular ImageNet-1K trained models

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
RegNetY-4G [44]	224 ²	21M	4.0G	1156.7	80.0
RegNetY-8G [44]	224 ²	39M	8.0G	591.6	81.7
RegNetY-16G [44]	224 ²	84M	16.0G	334.7	82.9
ViT-B/16 [19]	384 ²	86M	55.4G	85.9	77.9
ViT-L/16 [19]	384 ²	307M	190.7G	27.3	76.5
DeiT-S [57]	224 ²	22M	4.6G	940.4	79.8
DeiT-B [57]	224 ²	86M	17.5G	292.3	81.8
DeiT-B [57]	384 ²	86M	55.4G	85.9	83.1
Swin-T	224 ²	29M	4.5G	755.2	81.3
Swin-S	224 ²	50M	8.7G	436.9	83.0
Swin-B	224 ²	88M	15.4G	278.1	83.5
Swin-B	384 ²	88M	47.0G	84.7	84.5

(b) ImageNet-22K pre-trained models

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
R-101x3 [34]	384 ²	388M	204.6G	-	84.4
R-152x4 [34]	480 ²	937M	840.5G	-	85.4
ViT-B/16 [19]	384 ²	86M	55.4G	85.9	84.0
ViT-L/16 [19]	384 ²	307M	190.7G	27.3	85.2
Swin-B	224 ²	88M	15.4G	278.1	85.2
Swin-B	384 ²	88M	47.0G	84.7	86.4
Swin-L	384 ²	197M	103.9G	42.1	87.3


	ImageNet		COCO		ADE20k
	top-1	top-5	AP ^{box}	AP ^{mask}	mIoU
w/o shifting	80.2	95.1	47.7	41.5	43.3
shifted windows	81.3	95.6	50.5	43.7	46.1
no pos.	80.1	94.9	49.2	42.6	43.8
abs. pos. 	80.5	95.2	49.0	42.4	43.2
abs.+rel. pos.	81.3	95.6	50.2	43.4	44.0
rel. pos. w/o app.	79.3	94.7	48.2	41.9	44.1
rel. pos.	81.3	95.6	50.5	43.7	46.1

Table 4. Ablation study on the *shifted windows* approach and different position embedding methods on three benchmarks, using the Swin-T architecture. w/o shifting: all self-attention modules adopt regular window partitioning, without *shifting*; abs. pos.: absolute position embedding term of ViT; rel. pos.: the default settings with an additional relative position bias term (see Eq. (4)); app.: the first scaled dot-product term in Eq. (4).