


```

K_cmp.append(cur_K)
V_cmp.append(cur_V)

K_cmp = torch.cat(K_cmp, dim = 2).transpose(1,2)
V_cmp = torch.cat(V_cmp, dim = 2).transpose(1,2)
print(K_cmp.shape) # torch.Size([1, 4, 16]) # 32->4
print(V_cmp.shape) # torch.Size([1, 4, 16]) # 32->4

```

2. tokens

```

1. idx_slc_start = idx * d
   idx_slc_end = idx * d + l
   K_slc = torch.randn(batch_size, t, d * select_top_k, dim)
   V_slc = torch.randn(batch_size, t, d * select_top_k, dim)
   for i in range(batch_size):
       for j in range(t):
           for k in range(select_top_k):
               K_slc[i, j, k * d : k * d + l, :] = K[i, idx_slc_start[i, j, k ]
: idx_slc_end[i, j, k ] , :]
               V_slc[i, j, k * d : k * d + l, :] = V[i, idx_slc_start[i, j, k ]
: idx_slc_end[i, j, k ] , :]
   print(K_slc.shape) # bs, seq_len, select_kv, dim, 1,32,16,16, 1t
select_kv
   print(V_slc.shape) # bs, seq_len, select_kv, dim 1,32,16,16, 1t
select_kv

```

3. tokens NSA context q kv KV

```

1. # built sliding window attention
   def get_window_mask(seq_len, window):
       mask = torch.ones(seq_len, seq_len)
       mask = torch.tril(mask)
       win_mask = torch.ones(seq_len - window, seq_len - window)
       win_mask = 1.0 - torch.tril(win_mask)
       mask>window:, :seq_len - window] = win_mask
       return mask
   print(get_window_mask(7, 3)) # test
   window_mask = get_window_mask(t, 8)

```

4. [1, 32, 16]

```
1. o_list = [o_cmp, o_slc, o_win]
   o_star = torch.zeros(batch_size, t, dim)
   for i in range(3):
       o_star += gate[:, :, i].unsqueeze(2) * o_list[i]
   print(o_star.shape)
```

????

[[image.png](NSA ████████) by deepseek/jxLimage.png](NSA ████████) by
deepseek/jxLimage.png)

Revision #6

Created 2025-03-08 08:42:58 UTC by Colin

Updated 2026-04-29 07:34:30 UTC by Colin