

OneSparse: A Unified Framework for Sparse Activation Layers in Vision Models

Supplementary Material

A. Visual Analysis of Functional Specialization

We visualize routing probabilities in deep ConvNeXt stages to test whether the Unified Router allocates broad context to Memory Units and fine semantic refinement to Expert Units. Fig. A1 shows a clear separation between the two branches, consistent with the intended adaptive processing strategy.

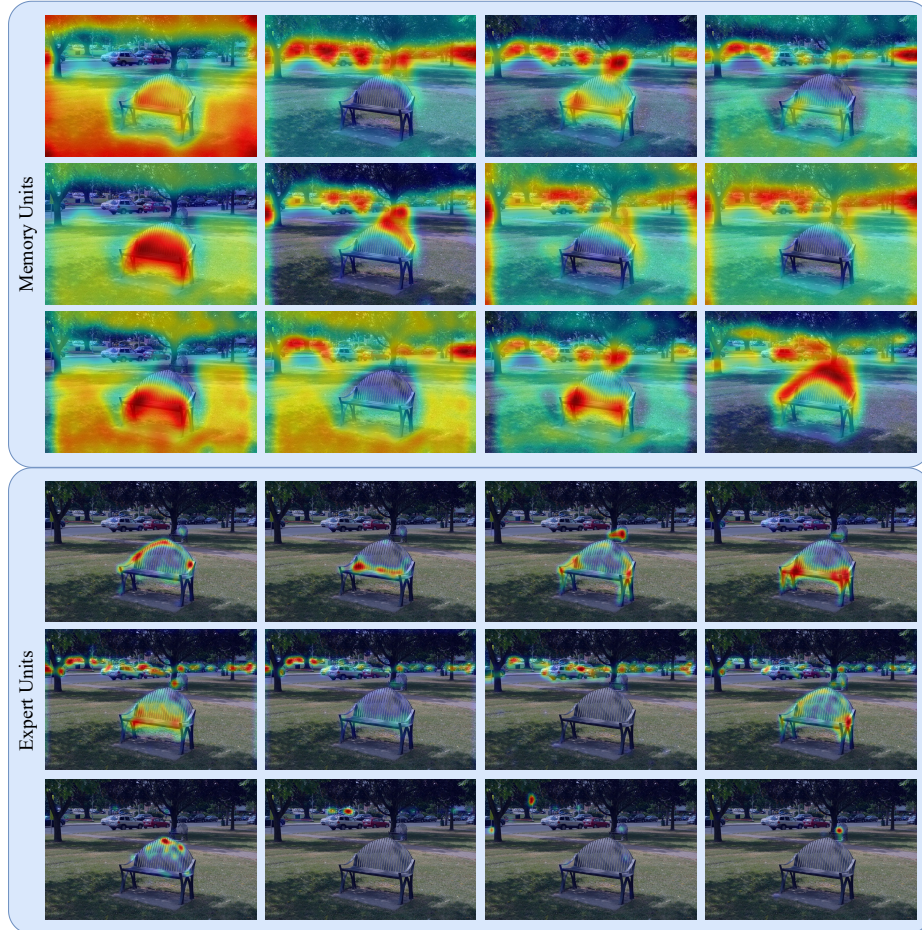


Figure A1. **Visualizing Functional Specialization. Top (Memory Units):** Diffuse, image-wide activations indicate coarse-grained global modeling. **Bottom (Expert Units):** Concentrated responses on salient objects and boundaries indicate fine-grained refinement.

As shown in Fig. A1, memory units exhibit broad image-wide attention, whereas expert units concentrate on salient objects and boundaries. This result supports the view that the Nexus Layer separates coarse context modeling from fine-grained refinement without explicit supervision.

B. Implementation of Nexus Layer

We provide PyTorch-style pseudo-code for the Nexus Layer, separating the core routing procedure from the processing sub-modules.

Algorithm A1 implements *Unified Routing*, while Algorithm A2 defines the expert and memory branches.

```

1 class NexusLayer(nn.Module):
2     def __init__(self, dim, num_experts, num_memory_units, slots_per_unit):
3         super().__init__()
4         self.num_experts = num_experts
5         self.num_memory = num_memory_units
6         self.total_units = num_experts + num_memory_units
7
8         # =====
9         # Unified Router Parameters
10        # =====
11        # Learnable queries for all processing slots.
12        # Experts and memory groups are treated uniformly as units.
13        # Q: (Total_Units, Slots, Dim)
14        self.router_queries = nn.Parameter(torch.randn(self.total_units, slots_per_unit, dim))
15        # Temperatures for the dispatch and combine softmaxes
16        self.router_temp = nn.Parameter(torch.tensor([1.0, 1.0]))
17
18        # =====
19        # Adaptive Processing Units
20        # =====
21        # Branch 1: Dynamic experts
22        self.experts = MultiExpertBlock(dim, num_experts)
23        # Branch 2: Static memory
24        self.memory = MemoryRetrievalBlock(dim, num_memory_units)
25
26    def forward(self, x):
27        """
28        Input: x (Batch, N_tokens, Dim)
29        Output: y (Batch, N_tokens, Dim)
30        """
31        # 1. Calculate Affinity Scores (Eq. 10)
32        # Project input tokens to the router's query space
33        # affinity: (Batch, N, Total_Units, Slots)
34        affinity = einsum("b n d, e s d -> b n e s", x, self.router_queries)
35
36        # 2. Generate Routing Tensors
37        # Dispatch tensor D: softmax over tokens
38        D = torch.softmax(affinity / self.router_temp[0], dim=1)
39
40        # Combine tensor B: softmax over slots
41        combine_logits = affinity / self.router_temp[1]
42        B = torch.softmax(combine_logits.flatten(start_dim=-2), dim=-1).view_as(affinity)
43
44        # 3. Dispatch Stage (Eq. 11 Inner Sum)
45        # Aggregate tokens into slots
46        z_in = einsum("b n e s, b n d -> b e s d", D, x)
47
48        # 4. Adaptive Processing (Heterogeneous Execution)
49        z_exp_in = z_in[:, :self.num_experts, ...]
50        z_mem_in = z_in[:, self.num_experts:, ...]
51
52        # Branch 1: Expert units
53        z_exp_out = self.experts(z_exp_in)
54        # Branch 2: Memory units
55        z_mem_out = self.memory(z_mem_in)
56
57        # Concatenate outputs
58        z_out = torch.cat([z_exp_out, z_mem_out], dim=1)
59
60        # 5. Combine Stage (Eq. 11 Outer Sum)
61        # Reconstruct final token representations
62        y = einsum("b n e s, b e s d -> b n d", B, z_out)
63
64        return y

```

Algorithm A1: Core implementation of the Nexus Layer for unified dispatch and combine.

```

1 class MultiExpertBlock(nn.Module):
2     """Vectorized implementation of multiple FFN experts."""
3     def __init__(self, dim, num_experts, mlp_ratio=4):
4         super().__init__()
5         # Flexible control over expert capacity
6         hidden_dim = int(dim * mlp_ratio)
7
8         # Batched weights for parallel execution of all experts
9         self.w1 = nn.Parameter(torch.randn(num_experts, dim, hidden_dim))
10        self.w2 = nn.Parameter(torch.randn(num_experts, hidden_dim, dim))
11
12    def forward(self, x):
13        # x: (Batch, Num_Experts, Slots, Dim)
14        # Standard FFN applied in parallel across the 'Num_Experts' dimension
15        x = einsum("b e s d, e d h -> b e s h", x, self.w1)
16        x = F.relu(x)
17        x = einsum("b e s h, e h d -> b e s d", x, self.w2)
18        return x
19
20 class MemoryRetrievalBlock(nn.Module):
21     """Memory processing via efficient Key-Value retrieval."""
22     def __init__(self, dim, num_units, num_keys=1024, topk=4):
23         super().__init__()
24         self.topk = topk
25         self.num_keys = num_keys
26
27         # Keys are used for computing similarity (Content addressing)
28         self.mem_keys = nn.Parameter(torch.randn(num_units, num_keys, dim))
29
30         # Values are stored in an EmbeddingBag for highly optimized weighted sum retrieval
31         # This avoids explicit gather + multiply + sum operations
32         self.value_bag = nn.EmbeddingBag(num_units * num_keys, dim, mode='sum')
33
34    def forward(self, x):
35        # x: (Batch, Num_Mem_Units, Slots, Dim)
36
37        # 1. Query the internal memory bank
38        # scores: (Batch, Num_Mem_Units, Slots, K_keys)
39        scores = einsum("b e s d, e k d -> b e s k", x, self.mem_keys)
40
41        # 2. Retrieve top-k indices and weights
42        topk_weights, topk_indices = torch.topk(scores, k=self.topk, dim=-1)
43        topk_weights = torch.softmax(topk_weights, dim=-1)
44
45        # Offset indices to align with the flattened EmbeddingBag
46        # offset: (Num_Mem_Units,) -> (1, Num_Mem_Units, 1, 1)
47        offset = torch.arange(x.shape[1], device=x.device) * self.num_keys
48        topk_indices = topk_indices + offset.view(1, -1, 1, 1)
49
50        # 3. Efficient Aggregation using EmbeddingBag
51        # Flatten for bag lookup: (Batch * Num_Units * Slots, TopK)
52        flat_indices = topk_indices.view(-1, self.topk)
53        flat_weights = topk_weights.view(-1, self.topk)
54
55        # The EmbeddingBag handles the weighted sum efficiently
56        out = self.value_bag(flat_indices, per_sample_weights=flat_weights)
57
58        # Reshape back to original 4D tensor
59        return out.view(x.shape)

```

Algorithm A2: Implementation of the processing sub-modules. The Expert Block performs dynamic FFN computation, while the Memory Block performs efficient key-value retrieval.