

Switch Transformer

This is a miniature [PyTorch](#) implementation of the paper [Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity](#). Our implementation only has a few million parameters and doesn't do model parallel distributed training. It does single GPU training, but we implement the concept of switching as described in the paper.

The Switch Transformer uses different parameters for each token by switching among parameters based on the token. Therefore, only a fraction of parameters are chosen for each token. So you can have more parameters but less computational cost.

The switching happens at the Position-wise Feedforward network (FFN) of each transformer block. Position-wise feedforward network consists of two sequentially fully connected layers. In switch transformer we have multiple FFNs (multiple experts), and we chose which one to use based on

```
39 import torch
40 from torch import nn
41
42 from labml_helpers.module import Module
43 from labml_nn.transformers.feed_forward import FeedForward
44 from labml_nn.transformers.mha import MultiHeadAttention
45 from labml_nn.utils import clone_module_list
```

a router. The output is a set of probabilities for picking a FFN, and we pick the one with the highest probability and only evaluate that. So essentially the computational cost is the same as having a single FFN. In our implementation this doesn't parallelize well when you have many or large FFNs since it's all happening on a single GPU. In a distributed setup you would have each FFN (each very large) on a different device.

The paper introduces another loss term to balance load among the experts (FFNs) and discusses dropping tokens when routing is not balanced.

Here's [the training code](#) and a notebook for training a switch transformer on Tiny Shakespeare dataset.



Routing among multiple FFNs

- `capacity_factor` is the capacity of each expert as a factor relative to ideally balanced load
- `drop_tokens` specifies whether to drop tokens if more tokens are routed to an expert than the capacity
- `is_scale_prob` specifies whether to multiply the input to the FFN by the routing probability
- `n_experts` is the number of experts
- `expert` is the expert layer, a [FFN module](#)

```
48 class SwitchFeedForward(Module):  
  
53     def __init__(self, *,  
54                 capacity_factor: float,  
55                 drop_tokens: bool,  
56                 is_scale_prob: bool,  
57                 n_experts: int,  
58                 expert: FeedForward,  
59                 d_model: int):
```

- `d_model` is the number of features in a token embedding
- `d_ff` is the number of features in the hidden layer of the FFN
- `dropout` is dropout probability in the FFN

```

70     super().__init__()
71
72     self.capacity_factor = capacity_factor
73     self.is_scale_prob = is_scale_prob
74     self.n_experts = n_experts
75     self.drop_tokens = drop_tokens

```

make copies of the FFNs

```

78     self.experts = clone_module_list(expert, n_experts)

```

Routing layer and softmax

```

80     self.switch = nn.Linear(d_model, n_experts)
81     self.softmax = nn.Softmax(dim=-1)

```

- `x` is the input to the switching module with shape `[seq_len, batch_size, d_model]`

```

83     def forward(self, x: torch.Tensor):

```

Capture the shape to change shapes later

```

89         seq_len, batch_size, d_model = x.shape

```

Flatten the sequence and batch dimensions

```

91         x = x.view(-1, d_model)

```

Get routing probabilities for each of the tokens.

```

97         route_prob = self.softmax(self.switch(x))

```

$$p_i(x) = \frac{e^{h(x)_i}}{\sum_j^N e^{h(x)_j}}$$

where N is the number of experts `n_experts` and $h(\cdot)$ is the linear transformation of token embeddings.

Get the maximum routing probabilities and the routes. We route to the expert with highest probability

```
101 route_prob_max, routes = torch.max(route_prob, dim=-1)
```

Get indexes of tokens going to each expert

```
104 indexes_list = [torch.eq(routes, i).nonzero(as_tuple=True)[0] for i in range(self.n_experts)]
```

Initialize an empty tensor to store outputs

```
107 final_output = x.new_zeros(x.shape)
```

Capacity of each expert.

```
113 capacity = int(self.capacity_factor * len(x) / self.n_experts)
```

expert capacity = $\frac{\text{tokens per batch}}{\text{number of experts}} \times \text{capacity}$



Number of tokens routed to each expert.

```
115 counts = x.new_tensor([len(indexes_list[i]) for i in range(self.n_experts)])
```

Initialize an empty list of dropped tokens

```
118 dropped = []
```

Only drop tokens if `drop_tokens` is `True`.

```
120 if self.drop_tokens:
```

Drop tokens in each of the experts

```
122     for i in range(self.n_experts):
```

Ignore if the expert is not over capacity

```
124         if len(indexes_list[i]) <= capacity:
125             continue
```

Shuffle indexes before dropping	127	<code>indexes_list[i] = indexes_list[i][torch.randperm(len(indexes_list[i]))]</code>
Collect the tokens over capacity as dropped tokens	129	<code>dropped.append(indexes_list[i][capacity:])</code>
Keep only the tokens upto the capacity of the expert	131	<code>indexes_list[i] = indexes_list[i][:capacity]</code>
Get outputs of the expert FFNs	134	<code>expert_output = [self.experts[i](x[indexes_list[i], :]) for i in range(self.n_experts)]</code>
Assign to final output	137 138	<code>for i in range(self.n_experts): final_output[indexes_list[i], :] = expert_output[i]</code>
Pass through the dropped tokens	141 142 143 144 145	<code>if dropped: dropped = torch.cat(dropped) final_output[dropped, :] = x[dropped, :] if self.is_scale_prob:</code>
Multiply by the expert outputs by the probabilities $y = p_i(x)E_i(x)$	147 148	<code>final_output = final_output * route_prob_max.view(-1, 1) else:</code>
Don't scale the values but multiply by $\frac{p}{\hat{p}} = 1$ so that the gradients flow (this is something we experimented with).	151	<code>final_output = final_output * (route_prob_max / route_prob_max.detach()).view(-1, 1)</code>
Change the shape of the final output back to <code>[seq_len, batch_size, d_model]</code>	154	<code>final_output = final_output.view(seq_len, batch_size, d_model)</code>
Return	165	<code>return final_output, counts, route_prob.sum(0), len(dropped), route_prob_max</code>

- the final output
- number of tokens routed to each expert
- sum of probabilities for each expert
- number of tokens dropped.
- routing probabilities of the selected experts

These are used for the load balancing loss and logging

Switch Transformer Block

This is the same as normal transformer block with handling extra outputs of switch feedforward module.

- `d_model` is the token embedding size
- `attn` is the attention module
- `feed_forward` is the feed forward module (which is the switching module in this case)
- `dropout_prob` is the probability of dropping out after self attention and FFN

```
168 class SwitchTransformerLayer(Module):
```

```
176     def __init__(self, *,
177                 d_model: int,
178                 attn: MultiHeadAttention,
179                 feed_forward: SwitchFeedForward,
180                 dropout_prob: float):
```

```
187         super().__init__()
188         self.size = d_model
189         self.attn = attn
190         self.feed_forward = feed_forward
```

```
191         self.dropout = nn.Dropout(dropout_prob)
192         self.norm_self_attn = nn.LayerNorm([d_model])
193         self.norm_ff = nn.LayerNorm([d_model])
```

```
195     def forward(self, *,
196                 x: torch.Tensor,
197                 mask: torch.Tensor):
```

Normalize the vectors before doing self attention

```
199         z = self.norm_self_attn(x)
```

Run through self attention, i.e. keys and values are from self

```
201         self_attn = self.attn(query=z, key=z, value=z, mask=mask)
```

Add the self attention results

```
203         x = x + self.dropout(self_attn)
```

Normalize for feed-forward

```
206         z = self.norm_ff(x)
```

Pass through the switching feed-forward network

```
208         ff, counts, route_prob, n_dropped, route_prob_max = self.feed_forward(z)
```

Add the feed-forward results back

```
210         x = x + self.dropout(ff)
211
212         return x, counts, route_prob, n_dropped, route_prob_max
```

Switch Transformer

```
215 class SwitchTransformer(Module):
```

```
220     def __init__(self, layer: SwitchTransformerLayer, n_layers: int):
221         super().__init__()
```

Make copies of the transformer layer

```
223         self.layers = clone_module_list(layer, n_layers)
```

Final normalization layer

225 `self.norm = nn.LayerNorm([layer.size])`

227 `def forward(self, x: torch.Tensor, mask: torch.Tensor):`

Run through each transformer layer

229 `counts, route_prob, n_dropped, route_prob_max = [], [], [], []`
230 `for layer in self.layers:`
231 `x, f, p, n_d, p_max = layer(x=x, mask=mask)`
232 `counts.append(f)`
233 `route_prob.append(p)`
234 `n_dropped.append(n_d)`
235 `route_prob_max.append(p_max)`

Finally, normalize the vectors

237 `x = self.norm(x)`

239 `return x, torch.stack(counts), torch.stack(route_prob), n_dropped, torch.stack`
 `(route_prob_max)`