
Mixture of Experts

Motivation

- Solving Complex Problems
- Enhancing Performance
- Scalability
- Diverse Applications

Motivation

- Solving Complex Problems
- Enhancing Performance
- Scalability
- Diverse Applications

Motivation

Solving Complex Problems

- By utilizing multiple experts, MoE can effectively handle complex and diverse problems with various aspects.
- Each expert can focus on a specific part of the problem, enhancing overall performance.

Motivation

Enhancing Performance

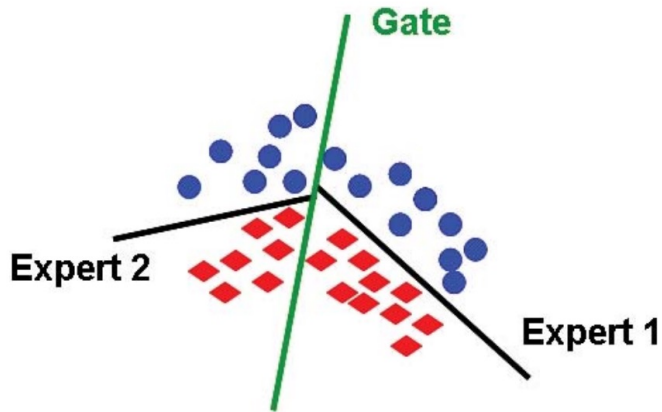
- MoE uses a gating mechanism to select the most suitable expert for each specific case, ensuring that each part of the problem is tackled by the best expert.
- Minimizes resource usage by activating only the necessary experts, saving time and memory.

Motivation

Scalability

- The system can be expanded by adding new experts without changing the existing structure.
- MoE allows for the distribution of work among experts, enhancing the ability to process large datasets and improve scalability.

MoE



- Simplified classification example for ME
- The blue circles and the red diamonds belong to classes 1 and 2, respectively, and they present a nonlinear classification example

MoE

- Let $D = X, Y$ denote the data where $X = x^{(n)}_{n=1}^N$ is the input, $Y = y^{(n)}_{n=1}^N$ is the target, and N is the number of training points. Also, let $\theta = \{\theta_g, \theta_e\}$ denote the set of all parameters where θ_e is the set of expert parameters
- Given an input vector x and a target vector y , the total probability of observing y can be written in terms of the experts, as

$$\begin{aligned} P(y|x, \theta) &= \sum_{i=1}^I P(y, i|x, \theta) \\ &= \sum_{i=1}^I P(i|x, \theta_g) P(y|i, x, \theta_e) \\ &= \sum_{i=1}^I g_i(x, \theta_g) P(y|i, x, \theta_e) \end{aligned}$$

MoE

$$\begin{aligned}P(y|x, \theta) &= \sum_{i=1}^I P(y, i|x, \theta) \\&= \sum_{i=1}^I P(i|x, \theta_g) P(y|i, x, \theta_e) \\&= \sum_{i=1}^I g_i(x, \theta_g) P(y|i, x, \theta_e)\end{aligned}$$

- Where I is the number of experts, the function $g_i(x, \theta_g) = P(i|x, \theta_g)$, represents the gate's rating, i.e., the probability of the i th expert given x , and $P(y|i, x, \theta_e)$ is the probability of the i th expert generating y given x . The latter will be denoted by $P_i(y)$ from now on.

MoE

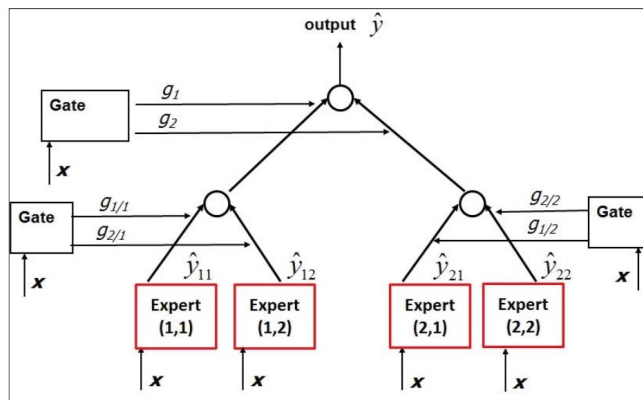
- During the training of ME, the gate and experts get decoupled, so the architecture attains a modular structure. Using this property

$$P(y|x, \theta) = \sum_{i=1}^I g_i(x, \theta_{g_i}) \sum_{j=1}^{J_i} g_{(j|i)}(x, \theta_{g_{(i|j)}}) P_{ij}(y, \theta_e)$$

Where I is the number of nodes connected to the gate at the top layer, and J_i is the number of nodes connected to the i th lower-level gating network, g_i is the output of the gate in the top layer, $g_{(j|i)}$ is the output of the j th gate connected to the i th gate of the top layer, and θ_{g_i} and $\theta_{g_{(i|j)}}$ are their parameters, respectively.

MoE

$$P(y|x, \theta) = \sum_{i=1}^I g_i(x, \theta_{g_i}) \sum_{j=1}^{J_i} g_{(j|i)}(x, \theta_{g(i|j)}) P_{ij}(y, \theta_e)$$



MoE

For both classification and regression, the gate is defined by the softmax function

$$g_i(\mathbf{x}, \mathbf{v}) = \frac{\exp(\beta_i(\mathbf{x}, \mathbf{v}))}{\sum_{j=1}^I \exp(\beta_j(\mathbf{x}, \mathbf{v}))}$$

$$\beta_i(\mathbf{x}, \mathbf{v}) = \mathbf{v}_i^T [\mathbf{x}, 1].$$

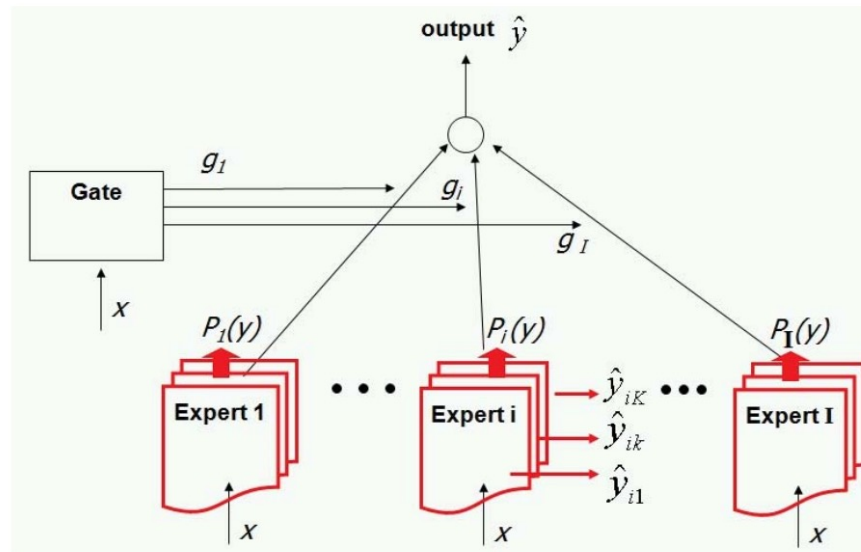
MoE Classification Model

K is the number of classes,
 $\{\mathbf{w}_{ik}\}_{k=1}^K$ is parameter of expert i th

$$\hat{y}_{ik} = \frac{\exp(\mathbf{w}_{ik}^T [\mathbf{x}, 1])}{\sum_{r=1}^K \exp(\mathbf{w}_{ir}^T [\mathbf{x}, 1])}$$

$$P_i(\mathbf{y}) = \prod_k \hat{y}_{ik}^{y_k}$$

$$\hat{y}_k = \sum_i g_i(\mathbf{x}, \mathbf{v}) \hat{y}_{ik}$$

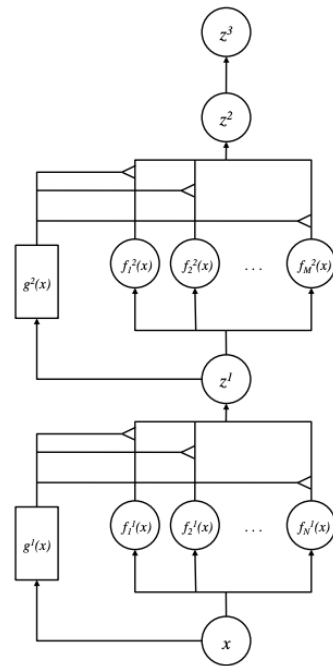
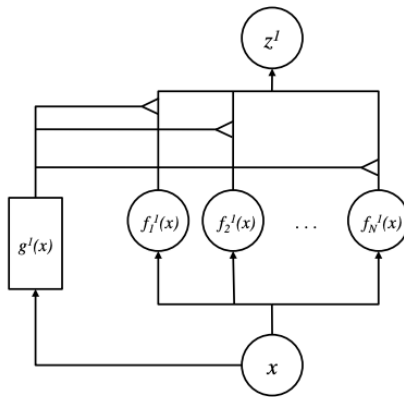


Deep Mixture of Experts

$$z^1 = \sum_{i=1}^N g_i^1(x) f_i^1(x)$$

$$z^2 = \sum_{j=1}^M g_j^2(z^1) f_j^2(z^1)$$

$$F(x) = z^3 = \text{softmax}(f^3(z^2))$$



Deep Mixture of Experts

Note:

The experts at each layer that perform best for the first few examples end up overpowering the remaining experts. This happens because the first examples increase the gating weights of these experts, which in turn causes them to be selected with high gating weights more frequently. This causes them to train more, and their gating weights to increase again, ad infinitum.

Let $G_i^l(t) = \sum_{t'=1}^t g_i^l(x_{t'})$ be the running total assignment to expert i of layer l at step t , and let $\overline{G^l(t)} = \frac{1}{N} \sum_{i=1}^N G_i^l(t)$ be their mean (here, $x_{t'}$ is the training example at step t'). Then for each expert i , we set $g_i^l(x_t) = 0$ if $G_i^l(t) - \overline{G^l(t)} > m$ for a margin threshold m .

THE SPARSELY GATED MIXTURE OF EXPERTS LAYER

$G(x)$ and $E_i(x)$ the output of the gating network and the output of the i -th expert network for a given input x .

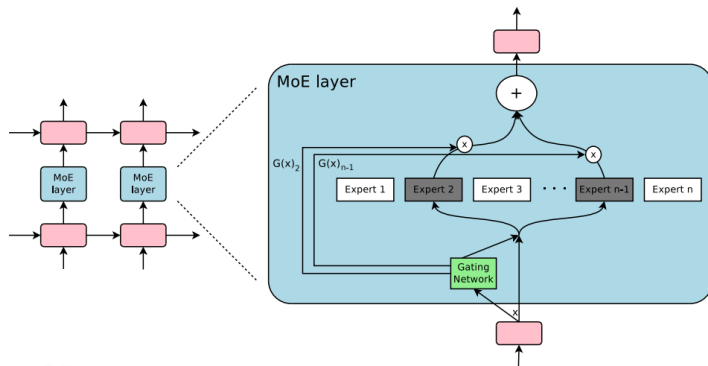
$$y = \sum_{i=1}^n G(x)_i E_i(x)$$

$$G_{\sigma}(x) = \text{Softmax}(x \cdot W_g)$$

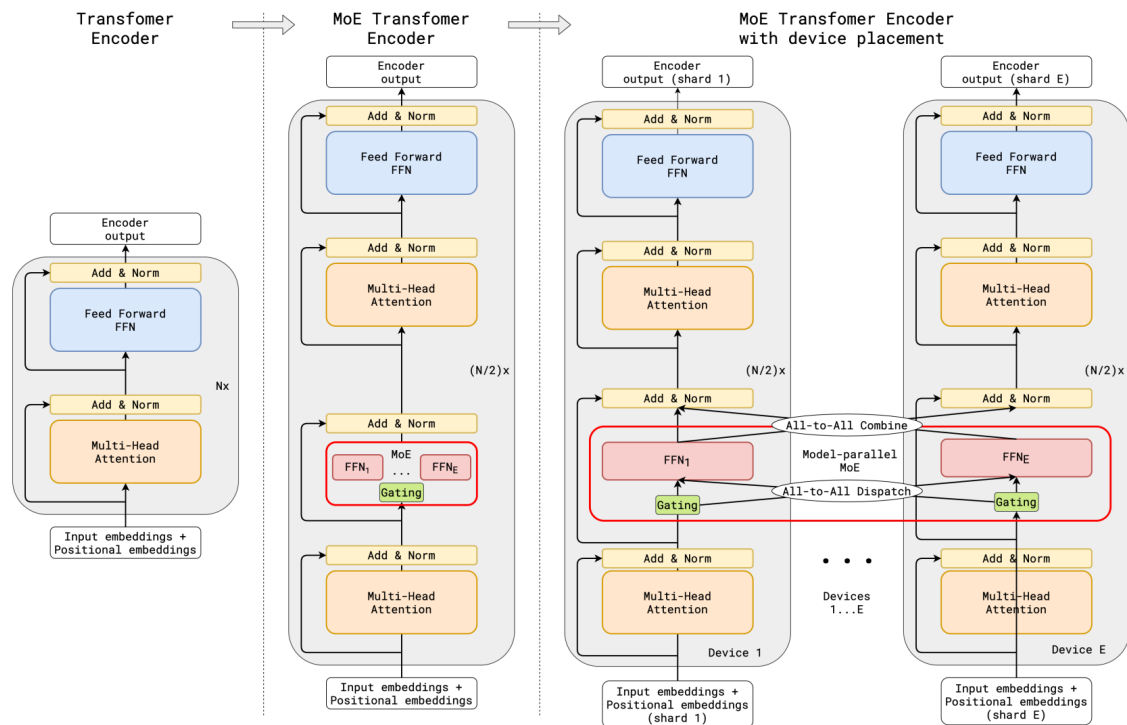
$$G(x) = \text{Softmax}(\text{KeepTopK}(H(x), k))$$

$$H(x)_i = (x \cdot W_g)_i + \text{StandardNormal}() \cdot \text{Softplus}((x \cdot W_{\text{noise}})_i)$$

$$\text{KeepTopK}(v, k)_i = \begin{cases} v_i & \text{if } v_i \text{ is in the top } k \text{ elements of } v. \\ -\infty & \text{otherwise.} \end{cases}$$



GShards



GShards

Algorithm 1: Group-level top-2 gating with auxiliary loss

Data: x_S , a group of tokens of size S

Data: C , Expert capacity allocated to this group

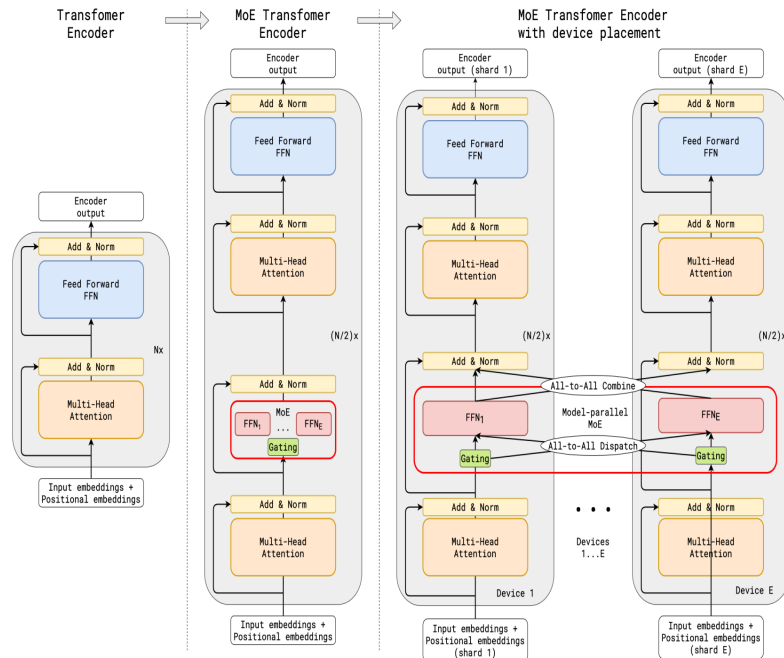
Result: $\mathcal{G}_{S,E}$, group combine weights

Result: ℓ_{aux} , group auxiliary loss

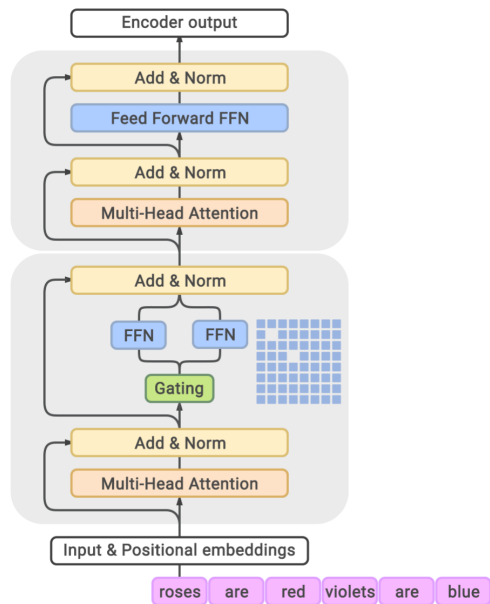
```

(1)  $c_E \leftarrow 0$                                  $\triangleright$  gating decisions per expert
(2)  $g_{S,E} \leftarrow \text{softmax}(wg \cdot x_S)$          $\triangleright$  gates per token per expert,  $wg$  are trainable weights
(3)  $m_E \leftarrow \frac{1}{S} \sum_{s=1}^S g_{s,E}$              $\triangleright$  mean gates per expert
(4) for  $s \leftarrow 1$  to  $S$  do
(5)    $g1, e1, g2, e2 = \text{top}_2(g_{s,E})$            $\triangleright$  top-2 gates and expert indices
(6)    $g1 \leftarrow g1 / (g1 + g2)$                  $\triangleright$  normalized  $g1$ 
(7)    $c \leftarrow c_{e1}$                            $\triangleright$  position in  $e1$  expert buffer
(8)   if  $c_{e1} < C$  then
(9)      $\mathcal{G}_{s,e1} \leftarrow g1$                      $\triangleright e1$  expert combine weight for  $x_s$ 
(10)  end
(11)   $c_{e1} \leftarrow c + 1$                        $\triangleright$  incrementing  $e1$  expert decisions count
(12) end
(13)  $\ell_{aux} = \frac{1}{E} \sum_{e=1}^E \frac{c_e}{S} \cdot m_e$ 
(14) for  $s \leftarrow 1$  to  $S$  do
(15)    $g1, e1, g2, e2 = \text{top}_2(g_{s,E})$            $\triangleright$  top-2 gates and expert indices
(16)    $g2 \leftarrow g2 / (g1 + g2)$                  $\triangleright$  normalized  $g2$ 
(17)    $rnd \leftarrow \text{uniform}(0, 1)$                $\triangleright$  dispatch to second-best expert with probability  $\propto 2 \cdot g2$ 
(18)    $c \leftarrow c_{e2}$                            $\triangleright$  position in  $e2$  expert buffer
(19)   if  $c < C \wedge 2 \cdot g2 > rnd$  then
(20)      $\mathcal{G}_{s,e2} \leftarrow g2$                      $\triangleright e2$  expert combine weight for  $x_s$ 
(21)  end
(22)   $c_{e2} \leftarrow c + 1$ 
(23) end

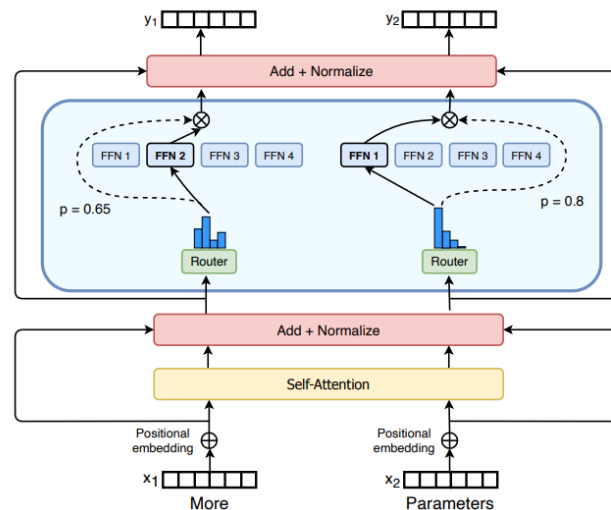
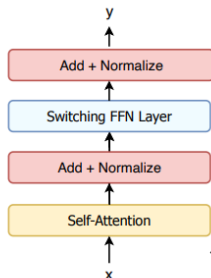
```



GLaM and Switch Transformers



GLaM



Switch transformers