# **MoE & recap**

Fine-tuning, RL, LoRA, and all that jazz

# Overview

| Train a base model | Instructions fine-tuning | Reinforcement Learning |
|---|---|---|



Hoffmann et al. (NeurIPS 2022)

Jacobs et al. (Neural Computation 1991)

Shazeer et al. (ICLR 2017)

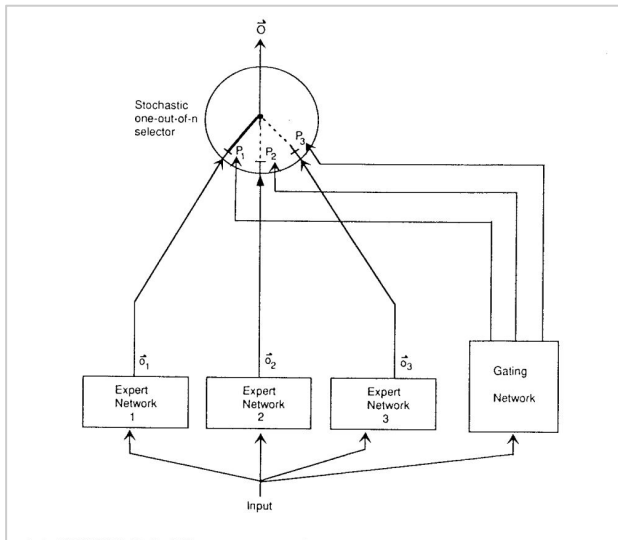Fedus et al. (JMLR 2022)

# MoE



Figure 1: A system of expert and gating networks. Each expert is a feed-forward network and all experts receive the same input and have the same number of outputs. The gating network is also feedforward, and typically receives the same input as the expert networks. It has normalized outputs $p_j = \exp(x_j) / \sum_i \exp(x_i)$, where $x_j$ is the total weighted input received by output unit $j$ of the gating network. The selector acts like a multiple input, single output stochastic switch; the probability that the switch will select the output from expert $j$ is $p_j$.



Figure 2: Illustration of a Switch Transformer encoder block. We replace the dense feed forward network (FFN) layer present in the Transformer with a sparse Switch FFN layer (light blue). The layer operates independently on the tokens in the sequence. We diagram two tokens ($x_1$ = "More" and $x_2$ = "Parameters" below) being routed (solid lines) across four FFN experts, where the router independently routes each token. The switch FFN layer returns the output of the selected FFN multiplied by the router gate value (dotted-line).

# RLHF



Figure 2: Diagram of our human feedback, reward model training, and policy training procedure.
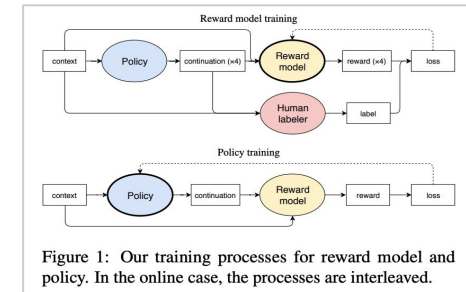
2020

2022

Figure 1: Our training processes for reward model and policy. In the online case, the processes are interleaved.

2019