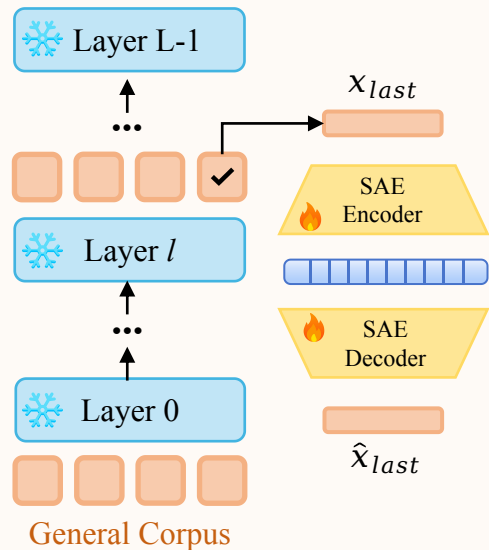


### Stage 1: Sequence-level SAE Pretraining



### Stage 2: SARM Reward Modeling

