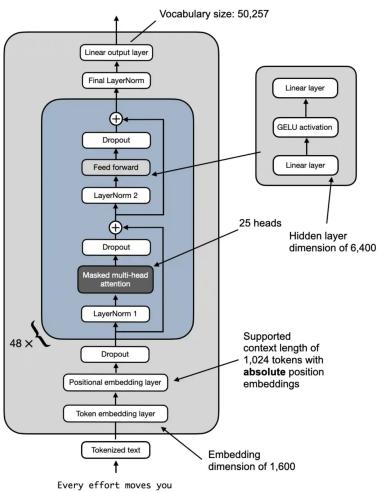
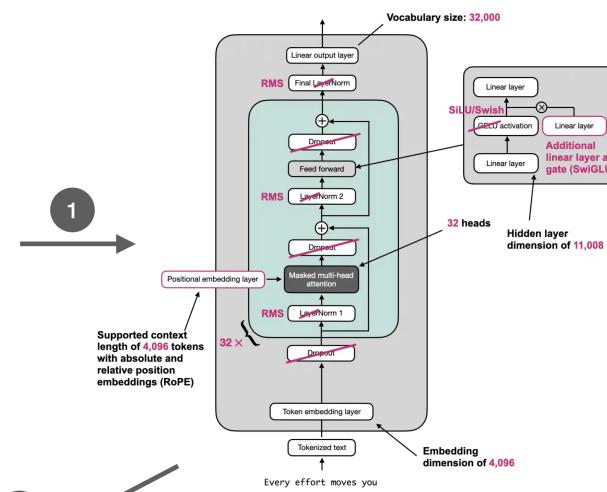


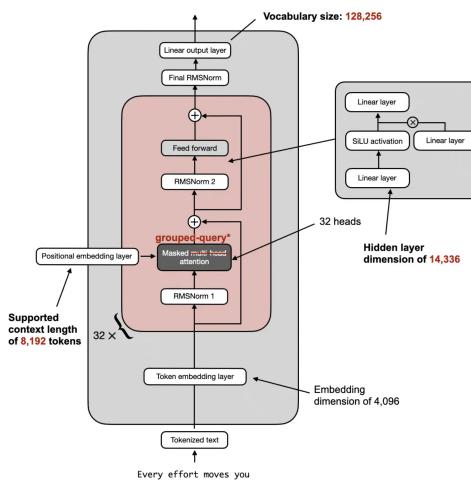
## GPT-2 XL 1.5B



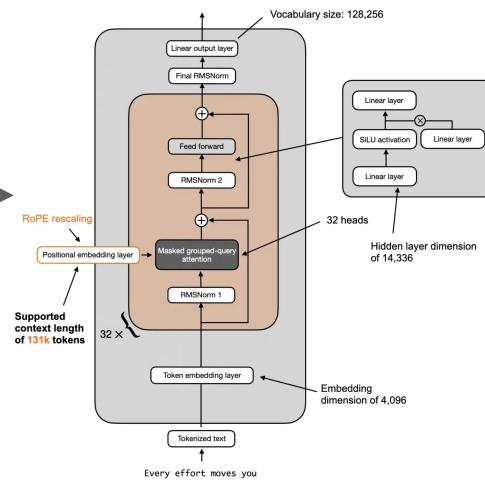
## Llama 2 7B



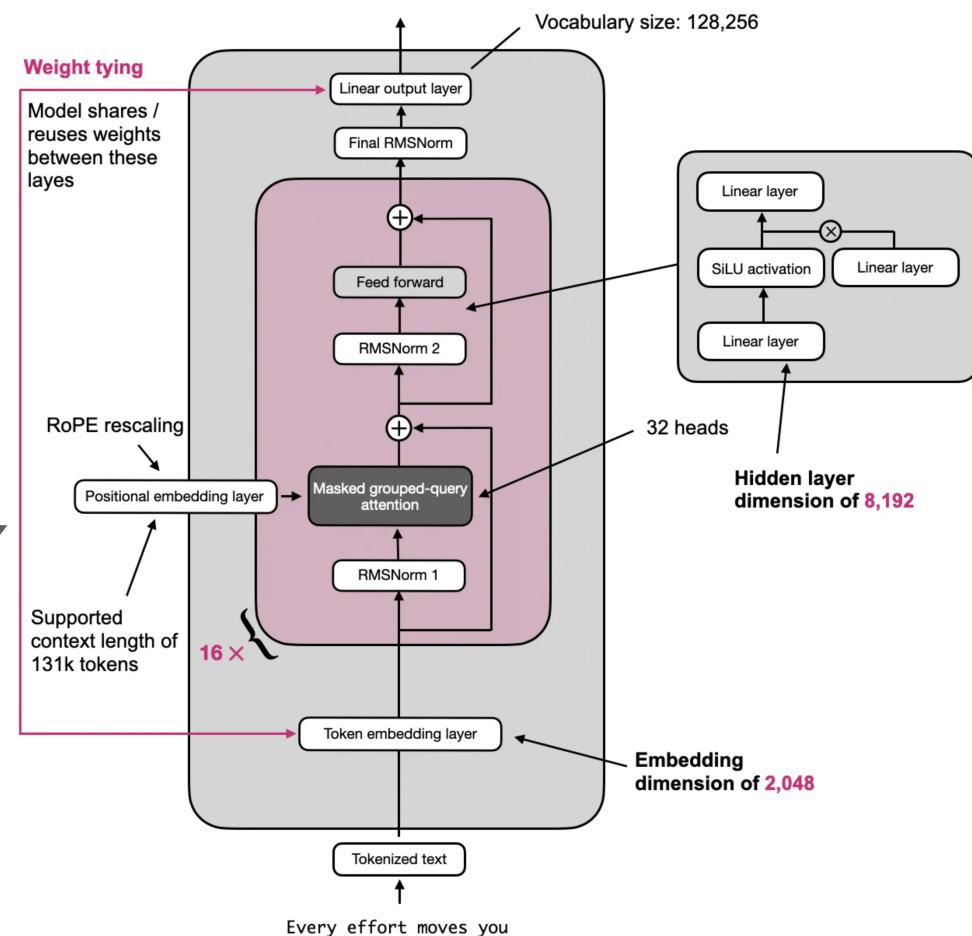
## Llama 3 8B



## Llama 3.1 8B



## Llama 3.2 1B



\* The larger Llama 2 34B and 70B also used grouped-query attention