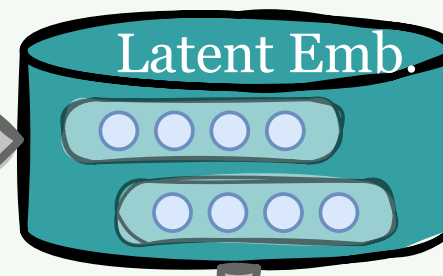
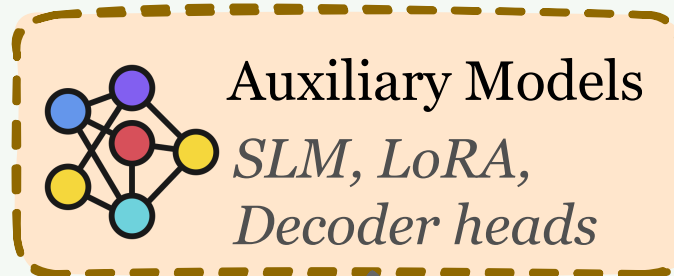


### (a) Generate

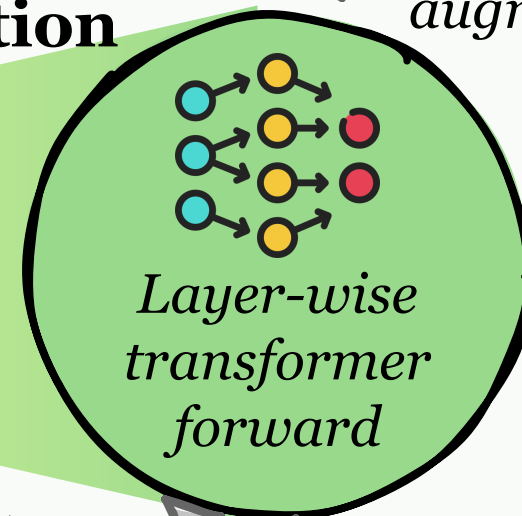


### LLM Internal Calculation



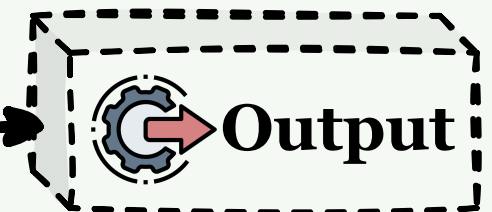
Inspection  
Closer

Text labels indicating the inspection and closer components of the internal calculation.



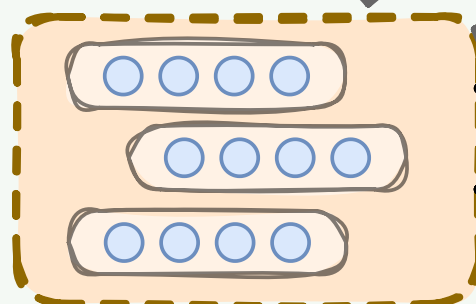
Interfere /  
augment

Text label indicating the interference or augmentation step.



KV cache/  
Intermediate  
Embeddings

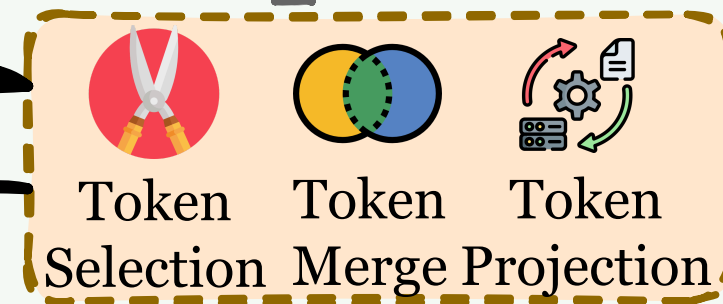
Text label indicating the storage of KV cache or intermediate embeddings.



### (b) Resue

augment

Text label indicating the augmentation step.



### (c) Transform