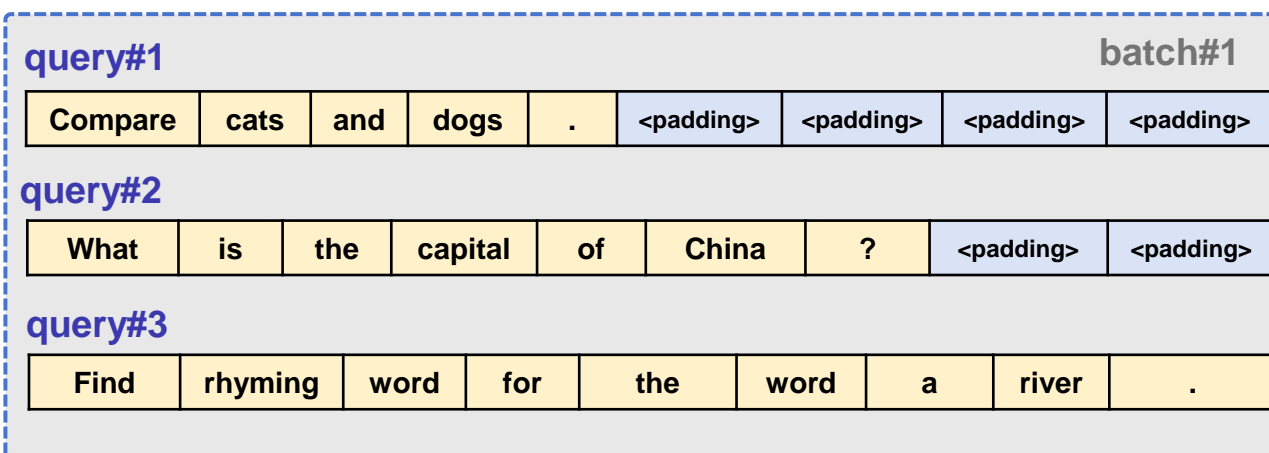


Default

KV Cache (batching part)



KV Cache (inference part)



UELMLM

