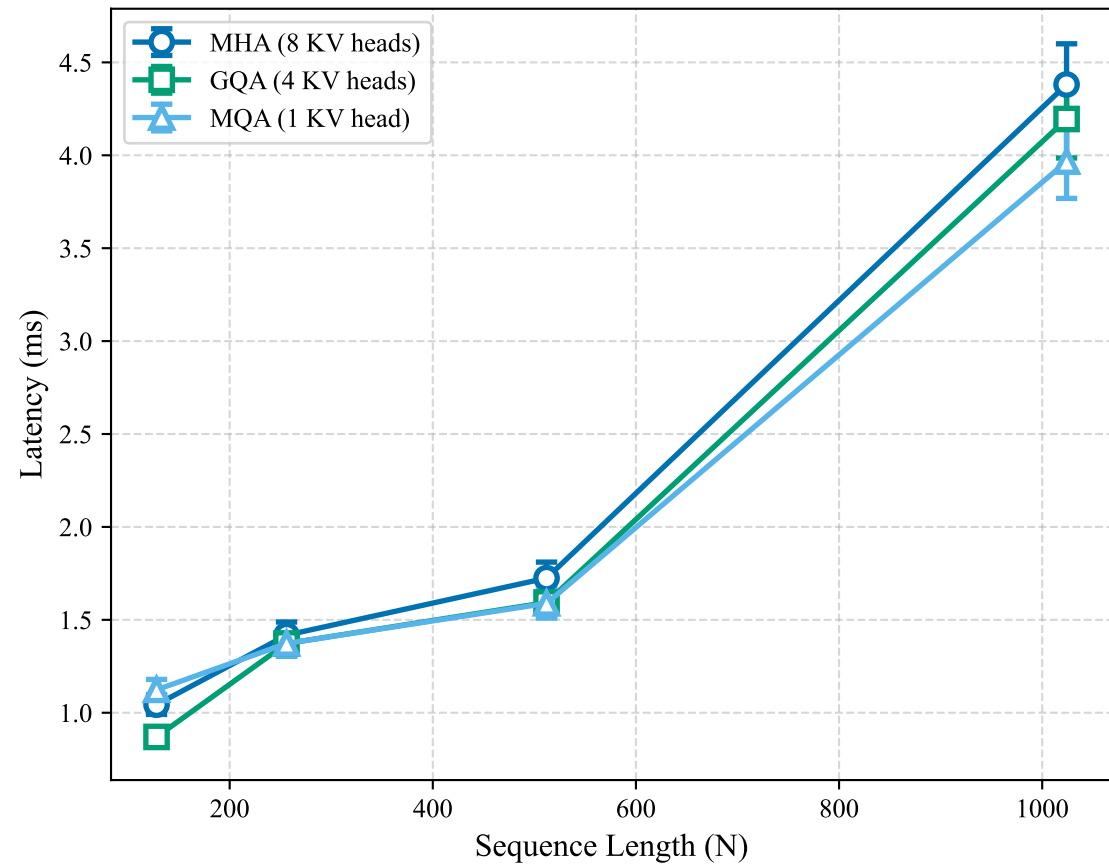**(a) Dense Attention Variants**

**(b) KV Cache Memory Requirements**