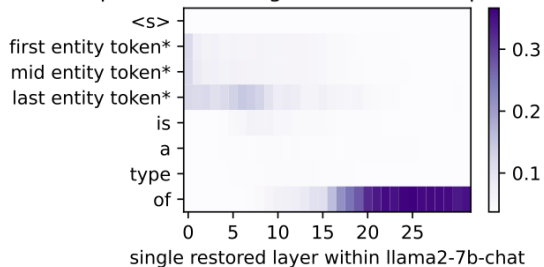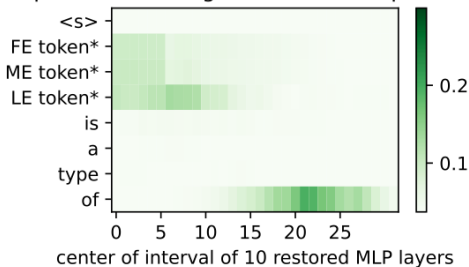**(A) Instance-Level Knowledge Locating**
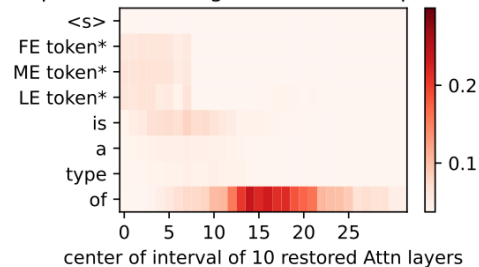
Impact of restoring state after corrupted input | Impact of restoring MLP after corrupted input | Impact of restoring Attn after corrupted input

single restored layer within llama2-7b-chat | center of interval of 10 restored MLP layers | center of interval of 10 restored Attn layers

**(B) Concept-Level Knowledge Locating**

Impact of restoring state after corrupted input | Impact of restoring MLP after corrupted input | Impact of restoring Attn after corrupted input
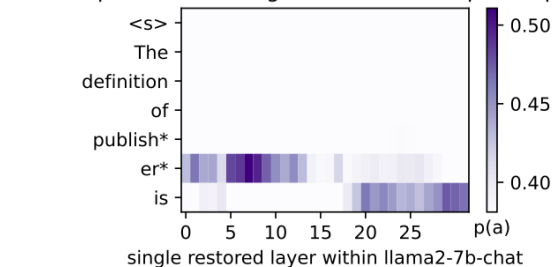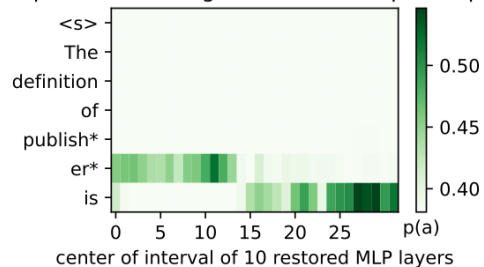
single restored layer within llama2-7b-chat | center of interval of 10 restored MLP layers | center of interval of 10 restored Attn layers
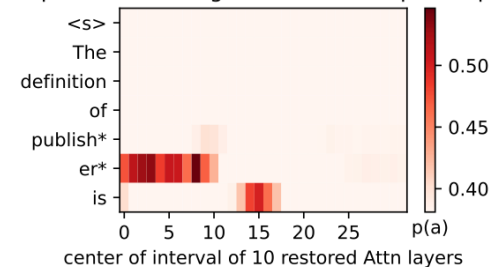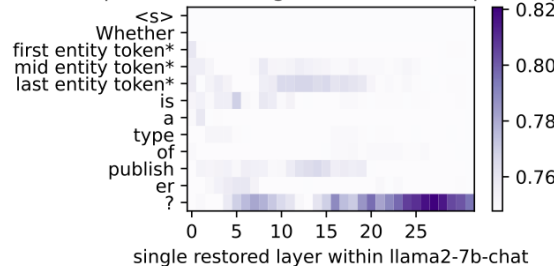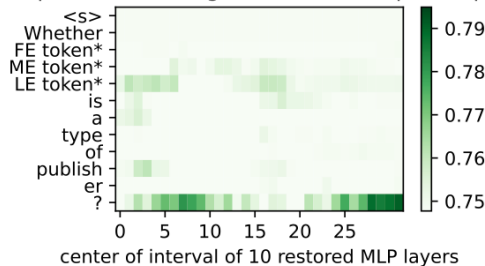
**(C) Combine Instance and Concept Level Locating**

Impact of restoring state after corrupted input | Impact of restoring MLP after corrupted input | Impact of restoring Attn after corrupted input

single restored layer within llama2-7b-chat | center of interval of 10 restored MLP layers | center of interval of 10 restored Attn layers