

# Argument (Symbolic)

"All A are B. All B are C. Therefore, all A are C"

## (a) Formalisation

$\mathcal{P}_1 : \forall(s, m_1)$   
 $\mathcal{P}_2 : \forall(m_2, p)$   
 $\mathcal{P}_1 \wedge \mathcal{P}_2 \rightarrow \mathcal{C} : \forall(s, p)$   
 where,  $m_1 \equiv m_2$

## (b) Clean Input

$(s, m_1, m_2, p)$

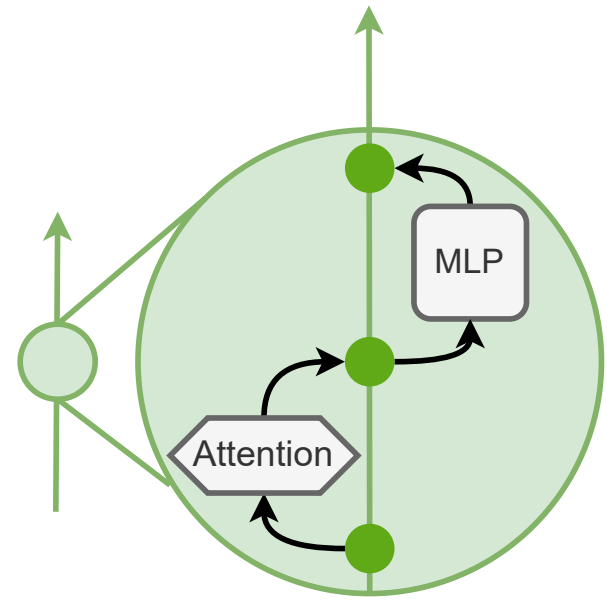
## Incomplete Syllogism

$x = [\mathcal{P}_1; \mathcal{P}_2; \mathcal{C} \setminus \{p\}]$

$m_2 \rightarrow m'_2$

$(s, m_1, m'_2, p)$

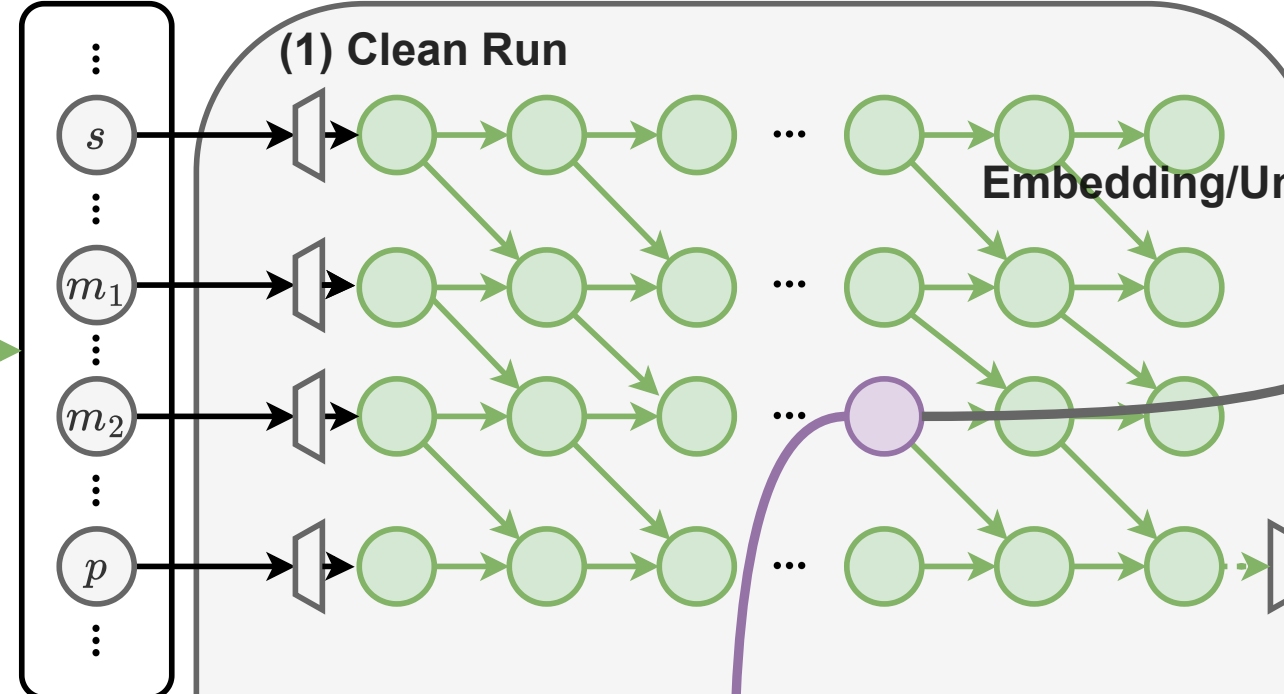
## (c) Corrupted Input



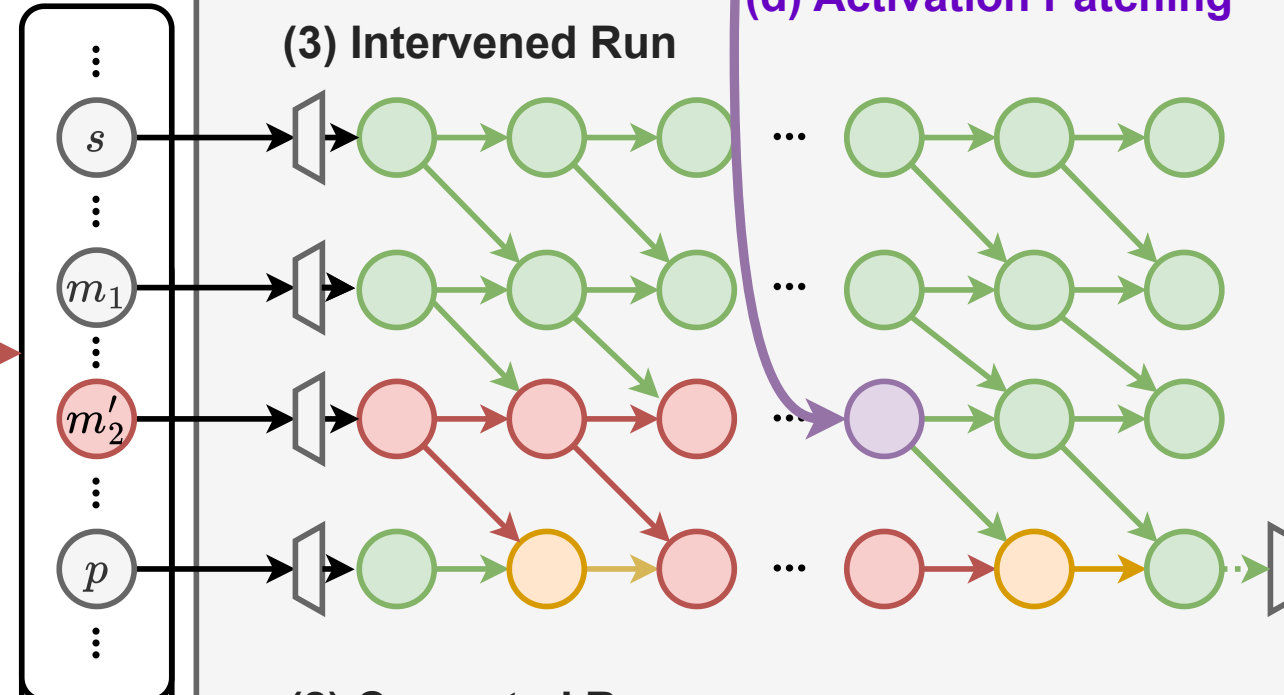
Residual Stream

## GPT-2 medium

### (1) Clean Run



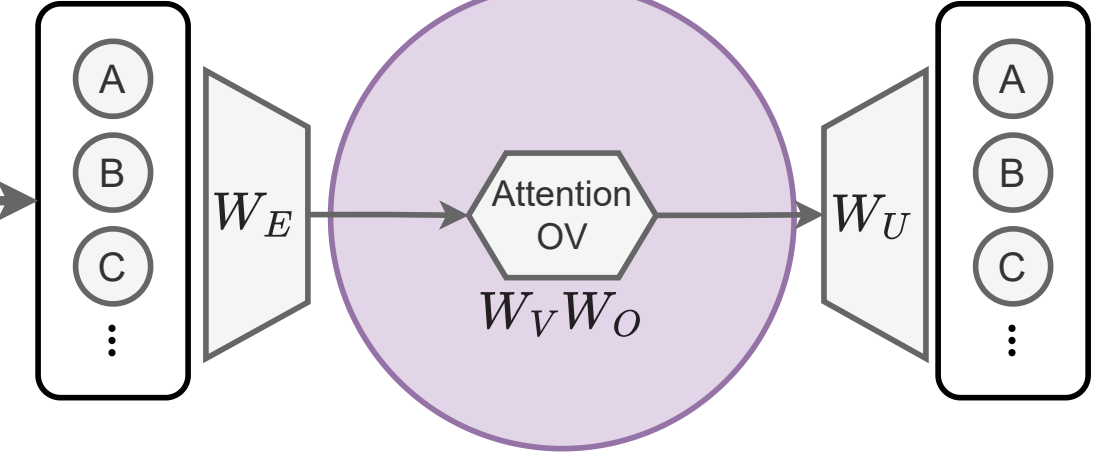
### (3) Intervened Run



### (2) Corrupted Run



## Embedding/Unembedding Projection



## (e) OV circuit Logit Lens

## Base Logit Difference

$\delta_+(p, m)$

## (d) Activation Patching

## (4) Quantification

$$\mathcal{S}_x = \frac{\delta_p(p, m) - \delta_-(p, m)}{\delta_+(p, m) - \delta_-(p, m)}$$

## Patched Logit Difference

$\delta_p(p, m)$

## Corrupted Logit Difference

$\delta_-(p, m)$

## Circuit Discovery

Symbolic Circuit  
(Schematic Logic)

Non-symbolic Circuit  
(World Knowledge,  
Beilef Bias)