

Mean drop of SplicedLLMLoss

4  
2  
0

0

2

4

6

8

10

Layer with intervention

Dataset kind=openwebtext-32-512

Vector multiplier

