Average indirect effect of a single hidden vector | Average indirect effect of a run of 10 MLP lookups | Average indirect effect of a run of 10 Attn modules

Legend:
- First subject token
- Middle subject tokens
- Last subject token
- First subsequent token
- Further tokens
- Last token

Y-axis: Average indirect effect on p(o)
X-axis: Layer number in Llama-3.1-8B