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.