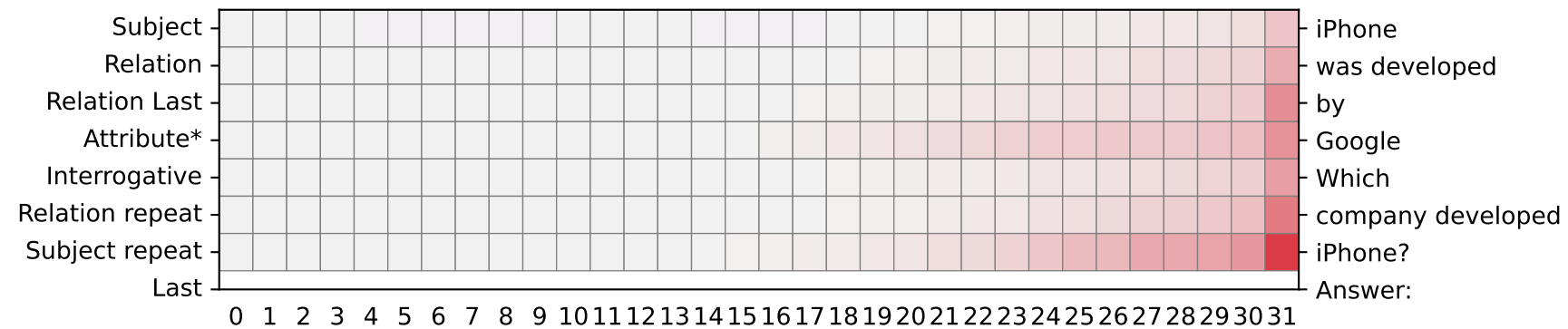


Layer



Logit of Factual



Layer



Logit of Counterfactual

