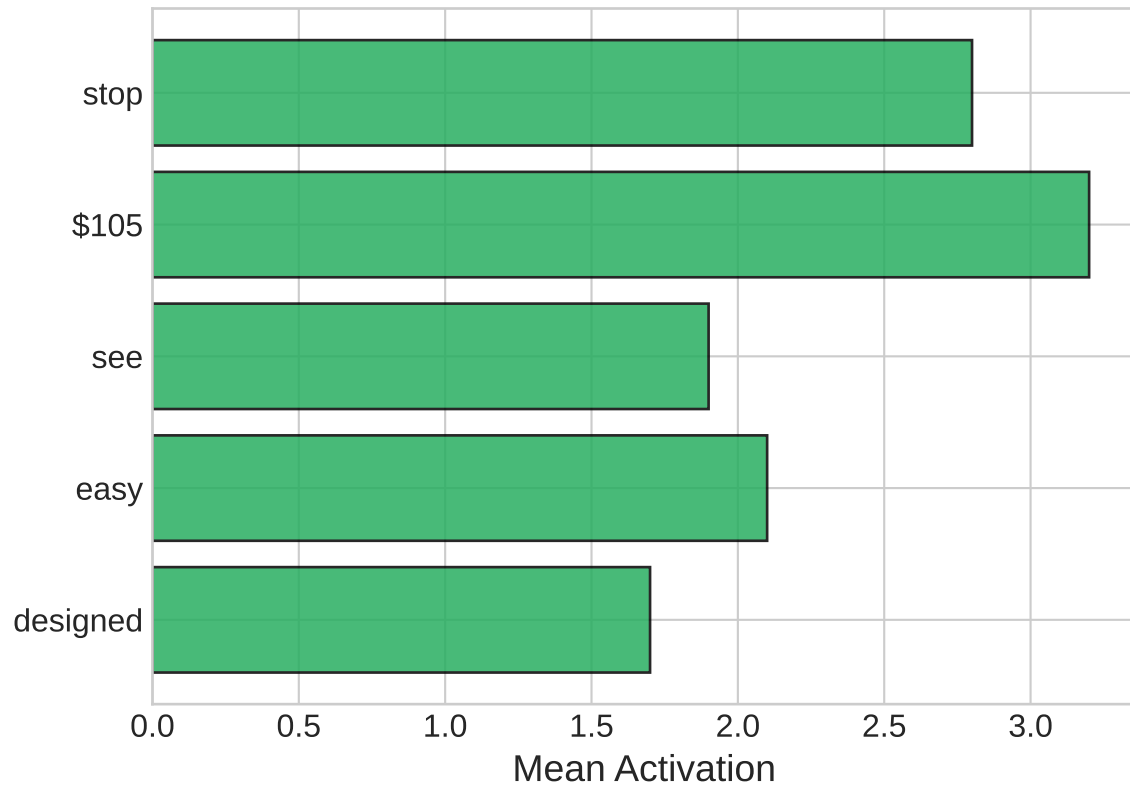


## Word-Feature Associations in LLM Responses

### Words Associated with Safe Features



### Words Associated with Risky Features

