| | house | burns | [SEP] |
|---|---|---|---|
| **Embedding cosine similarities** | | | |
| [CLS] | 0.61 | 0.63 | 0.4 |
| fire | 0.62 | 0.81 | 0.41 |
| starts | 0.59 | 0.73 | 0.41 |
| [SEP] | 0.78 | 0.51 | 0.38 |

| | house | burns | [SEP] |
|---|---|---|---|
| **Attention scores** | | | |
| [CLS] | 0.0086 | 0.029 | 0 |
| fire | 0.005 | 0.078 | 0 |
| starts | 0.004 | 0.87 | 0 |
| [SEP] | 0 | 0 | 0 |

**Causal strength: 0.73**

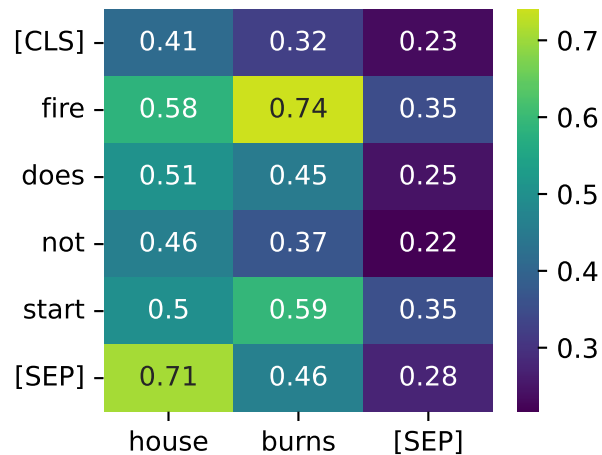| | house | burns | [SEP] |
|---|---|---|---|
| [CLS] | 0.0053 | 0.019 | 0 |
| fire | 0.0031 | 0.063 | 0 |
| starts | 0.0024 | 0.64 | 0 |
| [SEP] | 0 | 0 | 0 |

| | house | burns | [SEP] |
|---|---|---|---|
| **Embedding cosine similarities** | | | |
| [CLS] | 0.41 | 0.32 | 0.23 |
| fire | 0.58 | 0.74 | 0.35 |
| does | 0.51 | 0.45 | 0.25 |
| not | 0.46 | 0.37 | 0.22 |
| start | 0.5 | 0.59 | 0.35 |
| [SEP] | 0.71 | 0.46 | 0.28 |

| | house | burns | [SEP] |
|---|---|---|---|
| **Attention scores** | | | |
| [CLS] | 0.4 | 0.59 | 0 |
| fire | 0 | 0 | 0 |
| does | 0 | 0.0015 | 0 |
| not | 0 | 0.0024 | 0 |
| start | 0 | 0.0019 | 0 |
| [SEP] | 0 | 0 | 0 |

**Causal strength: 0.36**

| | house | burns | [SEP] |
|---|---|---|---|
| [CLS] | 0.16 | 0.19 | 0 |
| fire | 0 | 0 | 0 |
| does | 0 | 0.00066 | 0 |
| not | 0 | 0.00087 | 0 |
| start | 0 | 0.0011 | 0 |
| [SEP] | 0 | 0 | 0 |