**Input Plot**

Value of leading corporate brands in the
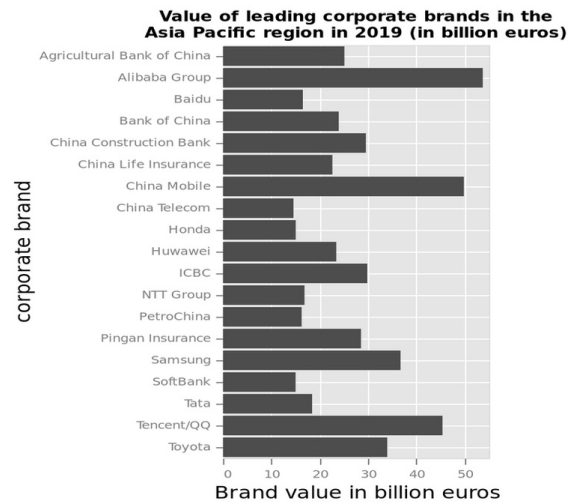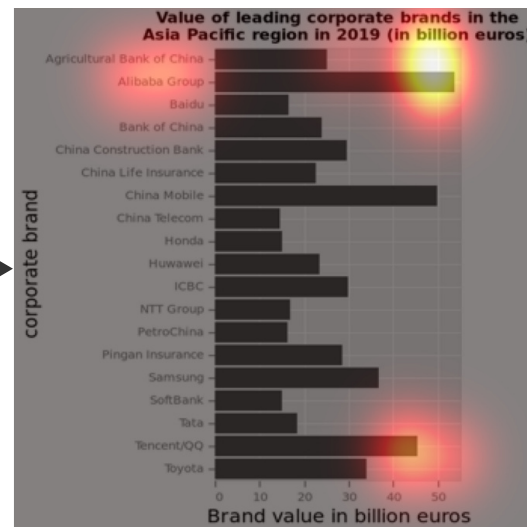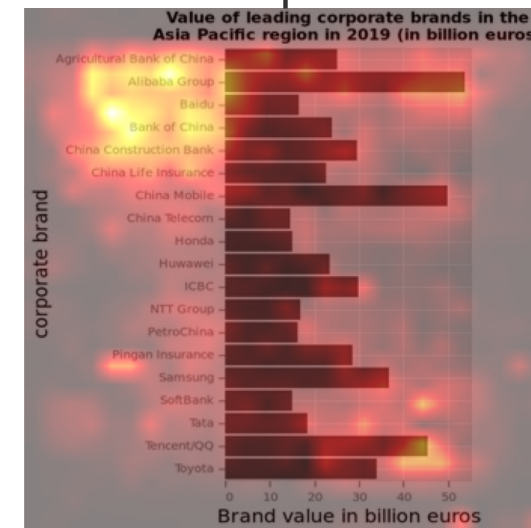Asia Pacific region in 2019 (in billion euros)

**Human Annotation and Recording Eye-gaze Data**

**Human Eye-gaze Attention Matrix (G)**

$$\mathcal{L}_{\text{W-MSE}} = \frac{1}{N} \sum_{i=1}^{N} w_i \cdot (G_i - A_i)^2$$

**Question and Answers Generated Based on the Chart Summaries From VisText**

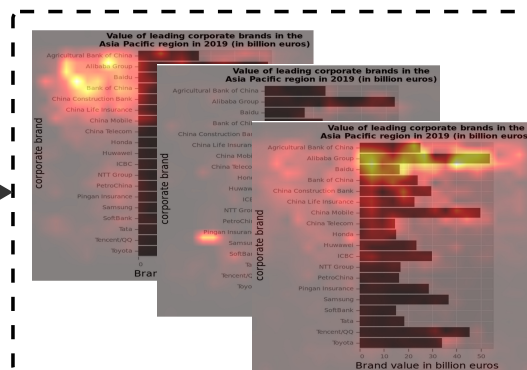Was Alibaba Group the most valuable brand at over €50 billion?

Answer: True

**LARGE VISION LANGUAGE MODEL**

**Attention Between Text and Vision Tokens Averaged over All Heads and Text Tokens**

**Average Over the First M Layers**

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{LM}} + \mathcal{L}_{\text{W-MSE}}$$

**Model Output Tokens:**
True </s>

$$\mathcal{L}_{\text{LM}} = -\sum_{t=1}^{T} \log P(y_t \mid x_{\leq t})$$