# See It From My Perspective:
## How Language Affects Cultural Bias in Image Understanding

♛ Amith Ananthram, ♜ Elias Stengel-Eskin, ♜ Mohit Bansal, ♛ Kathleen McKeown

COLUMBIA UNIVERSITY
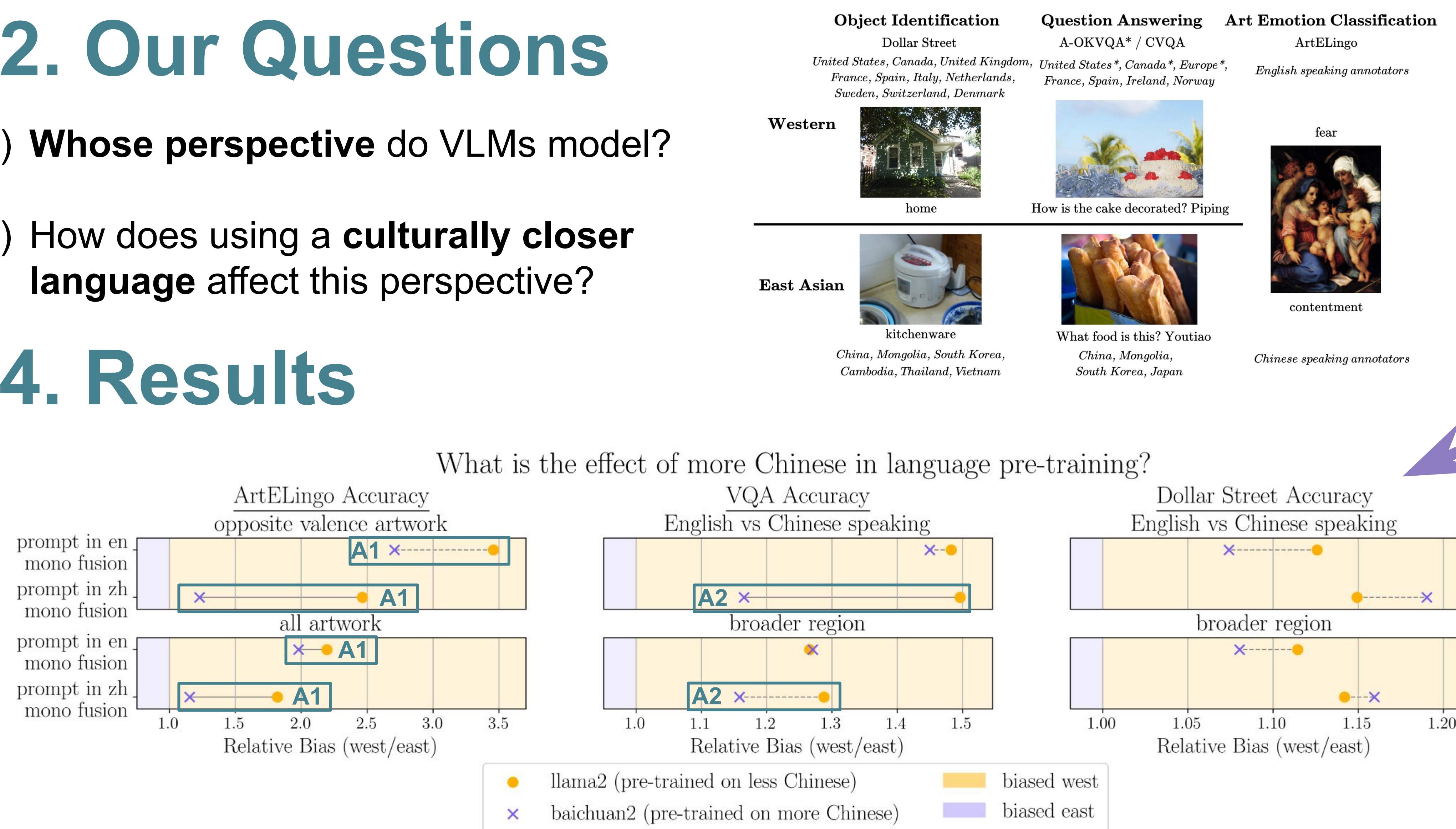THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

## 1. Background

- our knowledge and our beliefs are informed by culture (Goldstein, 1957)
- culture affects *how* we see things [color grouping (Chiao & Harada, 2008); attentional focus (Nisbett, 2001)]
- LLMs exhibit a Western worldview in knowledge & beliefs (Xu, 2024)
- VLMs inherit knowledge (Tsimpoukelli, 2021) and multilinguality from their LLMs
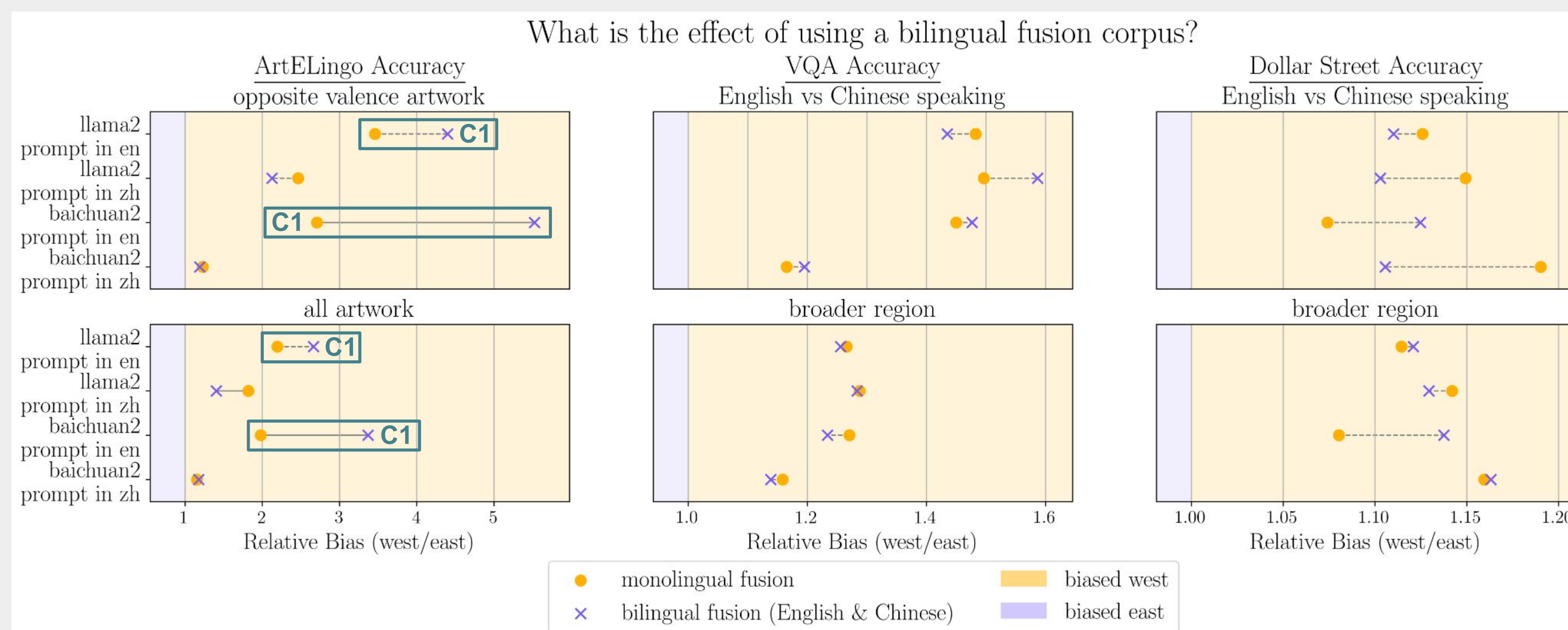
## 2. Our Questions

1) **Whose perspective** do VLMs model?

2) How does using a **culturally closer language** affect this perspective?

## 3. Approach



**Step 1: Bias Characterization**

Do off-the-shelf ($OTS_i$) VLMs have a Western bias?

$$\text{bias}_{OTS_i} = \frac{w_i}{e_i}$$

**Step 2: Bias Sourcing**

How does language affect Western bias in VLMs?

C: Fusion Corpus (mono vs bi)
en: a dog on a couch.
zh: 描述一下这个图像。

Llama2-Chat or Baichuan2-Chat — A: Base LLM

B: Prompt Language
en: Describe this image.
zh: 描述一下这个图像。

train our own **mLLaVA variants**

## 4. Results



What is the effect of more Chinese in language pre-training?

- llama2 (pre-trained on less Chinese)
- baichuan2 (pre-trained on more Chinese)
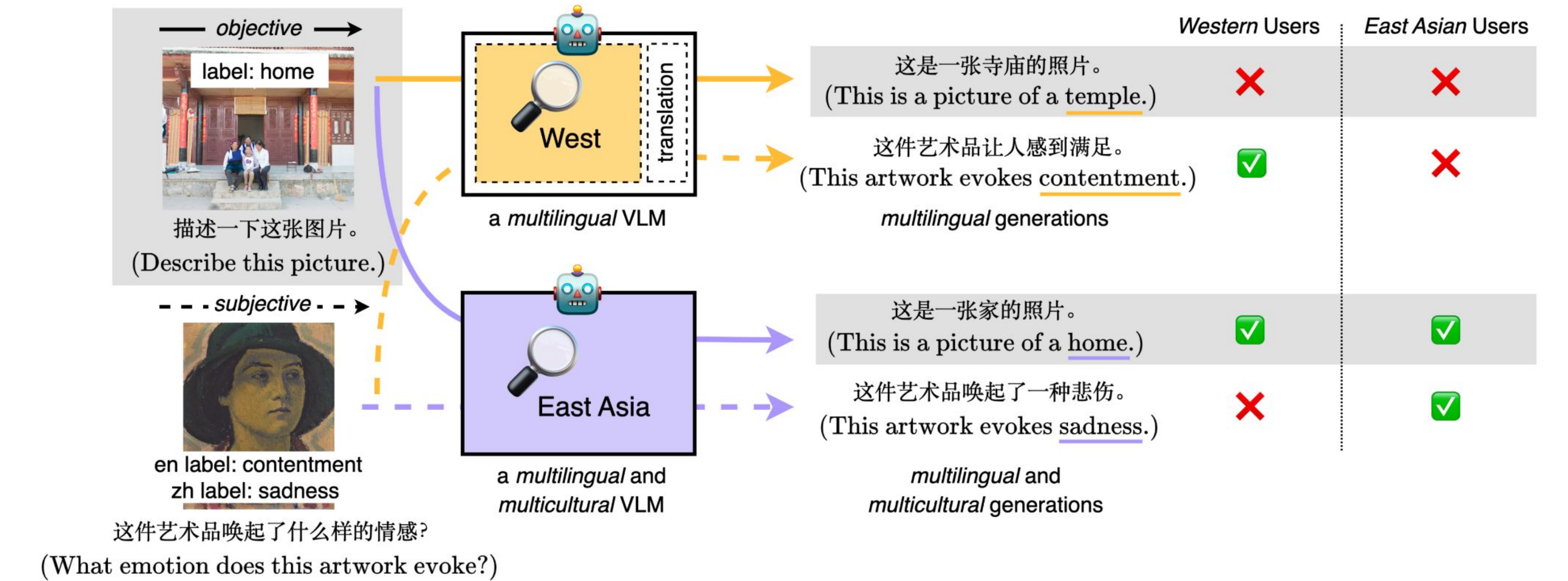- biased west
- biased east

A1. On *subjective* tasks, it reduces bias when prompting in both *English* and *Chinese*.
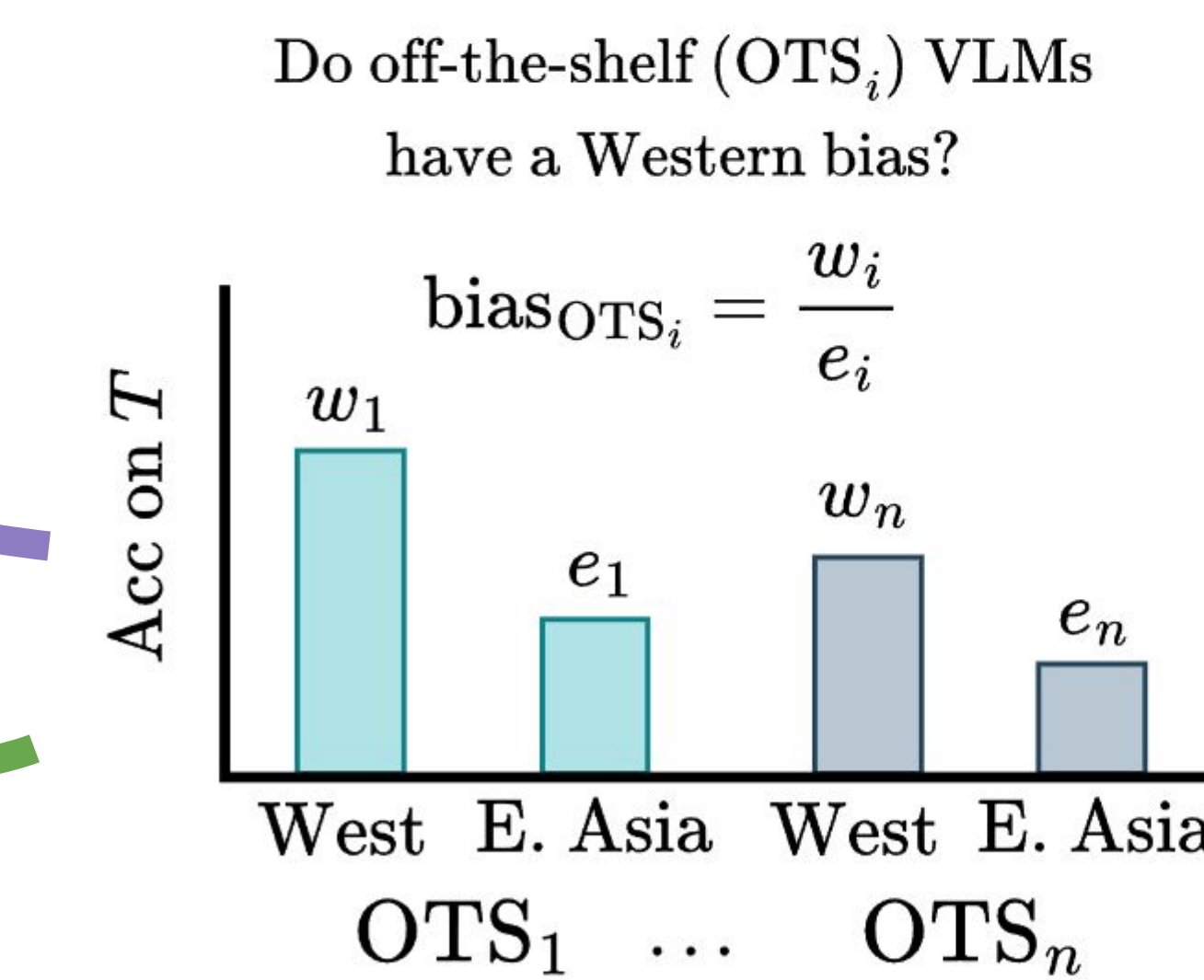A2. On *objective* tasks, it reduces bias when prompting in *Chinese*.



What is the effect of using a bilingual fusion corpus?

- monolingual fusion
- bilingual fusion (English & Chinese)
- biased west
- biased east

C1. On *subjective* tasks, it ties a language to its speakers' perspective (esp. in English).



What is the effect of prompting in Chinese?

- prompt in English
- prompt in Chinese
- biased west
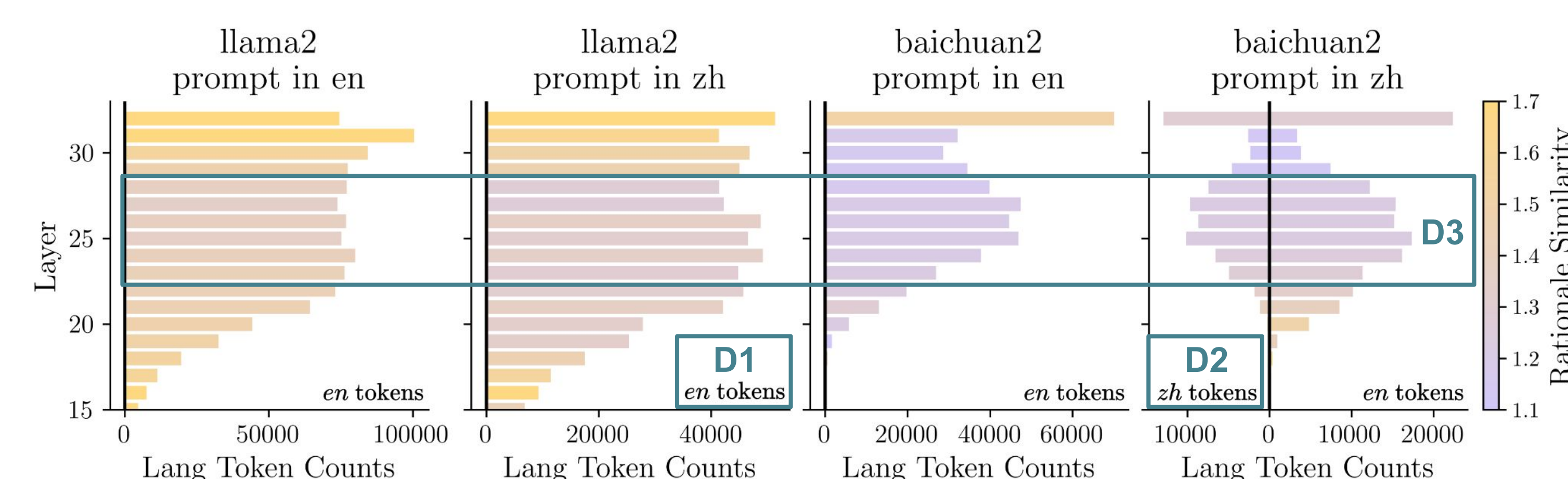- biased east

B1. On *subjective* tasks, it reduces bias *even if* little Chinese was seen during pre-training.
B2. However, it is more effective when Chinese was *common during pre-training*.



Mechanistic Exploration of Hidden States

D1. In Llama2, decode to English.

D2. In Baichuan2, decode to English and Chinese.

D3. In Baichuan2, *more similar* to rationales from Chinese language annotators than in Llama2.

## Takeaway:

A prompt language *can* reduce cultural bias in VLMs but the <u>text-only pre-training language mix</u> matters more; MT / bilingual fusion are insufficient proxies.