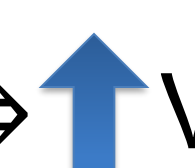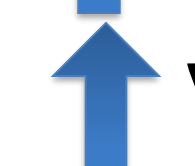# Unveiling Visual Perception in Language Models:
# An Attention Head Analysis Approach

Jing Bi, Junjia Guo, Yunlong Tang, Lianggong Bruce Wen, Zhang Liu, Chenliang Xu

## Introduction

➤ **Certain attention heads is specialized in visual perception**

➤ **Visual attention is measurable via attention distributions**

⬆ Image tokens ⇒ ⬆ Visual heads

⬆ Visual heads ≈ better performance

< 7B model ➕ > 2k token = ⬆ performance ✅

13B model ➕ > 2k token = More data ⚠️

➤ **Attention head can reveal the model-dataset behavior**

➤ **Quantifying attention head behavior can reduce the need for extensive benchmark testing**

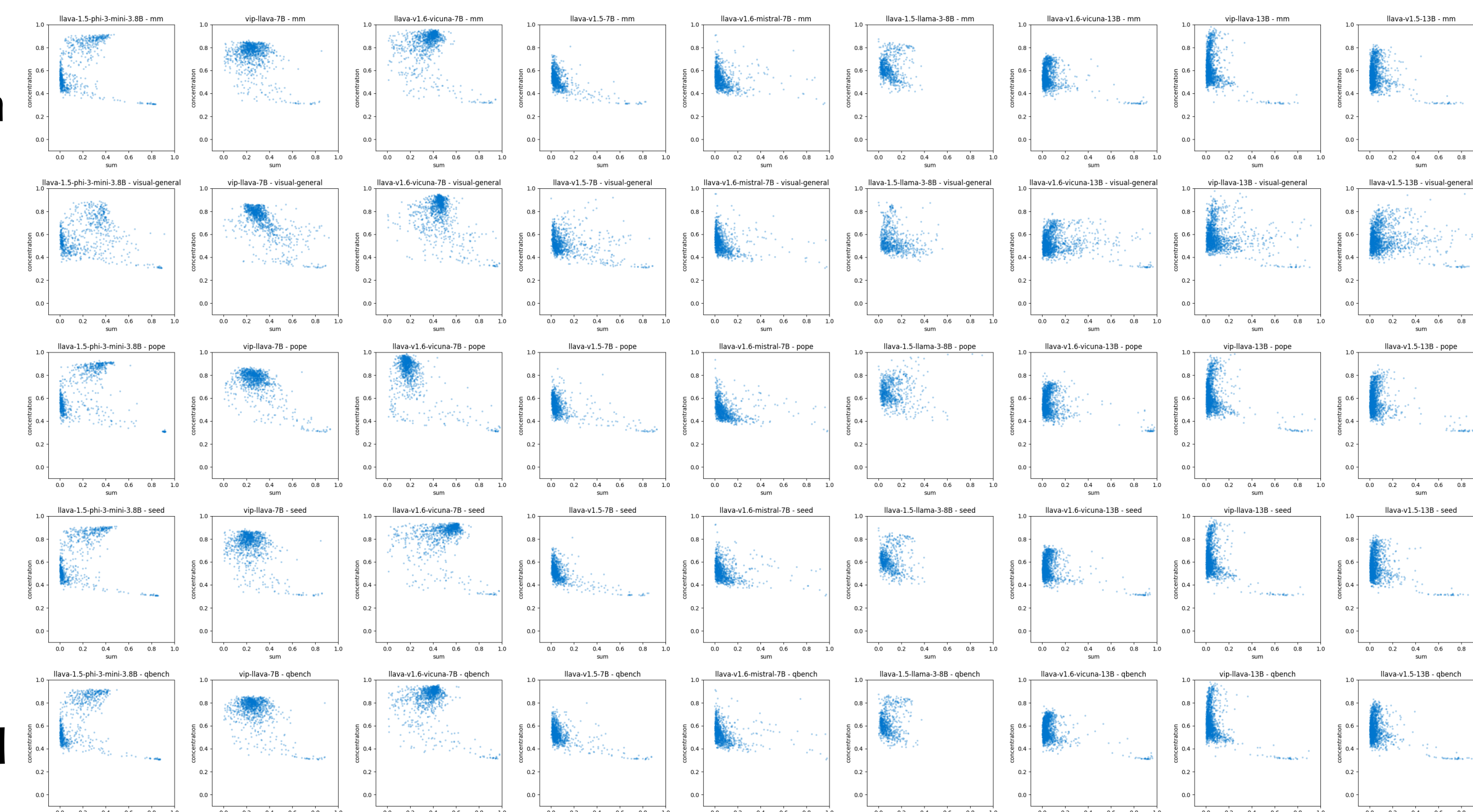| General Question | Obj Question | Super Question |
|---|---|---|
| How many of these are there? | How many **people** are pictured? | How many of these **beings** are there? |



| Model | LLM Family | Layer-Head | Resolution | Training Strategy | Visual Tokens |
|---|---|---|---|---|---|
| vip-phi-3-3.8B | Phi-3 | 24 × 32 | 336 × 336 | frozen vision encoder | 576 |
| 1.6-mistral-7B | Mistral-v0.2 | 32 × 32 | Dynamic Res | full model trainable | 576 × 1 ∼ 4 |
| vip-llama-3-8B | Llama-3 | 24 × 32 | 336 × 336 | frozen vision encoder | 576 |
| 1.5-7B | Vicuna-v1.5 | 32 × 32 | 336 × 336 | frozen vision encoder | 576 |
| 1.6-vicuna-7B | | 32 × 32 | Dynamic Res | full model trainable | 576 × 1 ∼ 4 |
| vip-7B | | 32 × 32 | 336 × 336 | frozen vision encoder | 576 |
| 1.5-13B | | 40 × 40 | 336 × 336 | frozen vision encoder | 576 |
| 1.6-vicuna-13B | | 40 × 40 | Dynamic Res | full model trainable | 576 × 1 ∼ 4 |
| vip-13B | | 40 × 40 | 336 × 336 | frozen vision encoder | 576 |