# NTSEBench: Cognitive Reasoning Benchmark for Vision Language Models

Pranshu Pandya[1†], Vatsal Gupta[1†], Agney S. Talwarr[1],
Tushar Kataria[2], Dan Roth[3], Vivek Gupta[4]*

†Equal Contribution, *Corresponding Author

[1]IIT Guwahati, [2]University of Utah, [3]University of Pennsylvania, [4]Arizona State University

Indian Institute of Technology, Guwahati

Arizona State University

NAACL 2025

Penn UNIVERSITY of PENNSYLVANIA

NLP

## 1. NTSEBench

Benchmark to evaluate the **cognitive reasoning** capabilities of SOTA **LLMs** and **VLMs**.



## 2. Dataset Highlights

- Questions are sampled from the National Talent Search Examination (NTSE), India.
- Mental Ability Test questions: Tests **reasoning rather than rote learning**
- 2728 MCQs with **4642 images**.
- Question covering **26 distinct categories across 8 cognitive dimensions**.
- Contains both **multimodal** (text-images) and **text-only** questions.

| Text Only | | Vision + Text | |
|---|---|---|---|
| Categories | # Samples | Categories | # Samples |
| Series | 256 | Non-Verbal Series | 95 |
| Alphabet Test | 94 | Missing Character | 127 |
| Odd one out | 170 | Embedded Figure | 96 |
| Analogy | 151 | Non-Verbal odd one out | 70 |
| Coding-Decoding | 149 | Non-Verbal Analogy | 100 |
| Number and Ranking | 139 | Paper Folding & Cutting | 96 |
| Blood Relation | 126 | Incomplete Figure | 94 |
| Mathematical Operations | 99 | Figure Partition | 71 |
| Puzzle Test | 95 | Cube and Dice | 89 |
| Syllogisms | 44 | Dot problem | 23 |
| Statement & Conclusions | 104 | Direction Sense | 96 |
| Data Sufficiency | 90 | Time and Clock | 51 |
| | | Mirror, Water and Images | 92 |
| | | Venn diagrams | 111 |

Table: **NTSEBench categories count**

| Question | Options | Solutions | # Samples |
|---|---|---|---|
| ✗ | ✗ | ✗ | 1199 |
| ✗ | ✗ | ✓ | 381 |
| ✗ | ✓ | ✗ | 70 |
| ✗ | ✓ | ✓ | 18 |
| ✓ | ✗ | ✗ | 330 |
| ✓ | ✗ | ✓ | 126 |
| ✓ | ✓ | ✗ | 403 |
| ✓ | ✓ | ✓ | 201 |

Table: **Modality Variations Question Count**

| | |
|---|---|
| Pattern Recognition | Logical Deduction |
| Spatial Reasoning | Relational Reasoning |
| Quantitative Analysis | Classification |
| Contextual Interpretation | Verbal Reasoning |

Table: **Cognitive Dimensions in NTSEBench.**

## 3. Methodology

**Dataset Construction:**



## Modeling strategies



Figure: **Examples Showing Input to Different Proposed modelling strategies.**(A) Text Only Standard QA strategy(B) Standard VQA (C) Interleaved Strategy (D) Image Only.
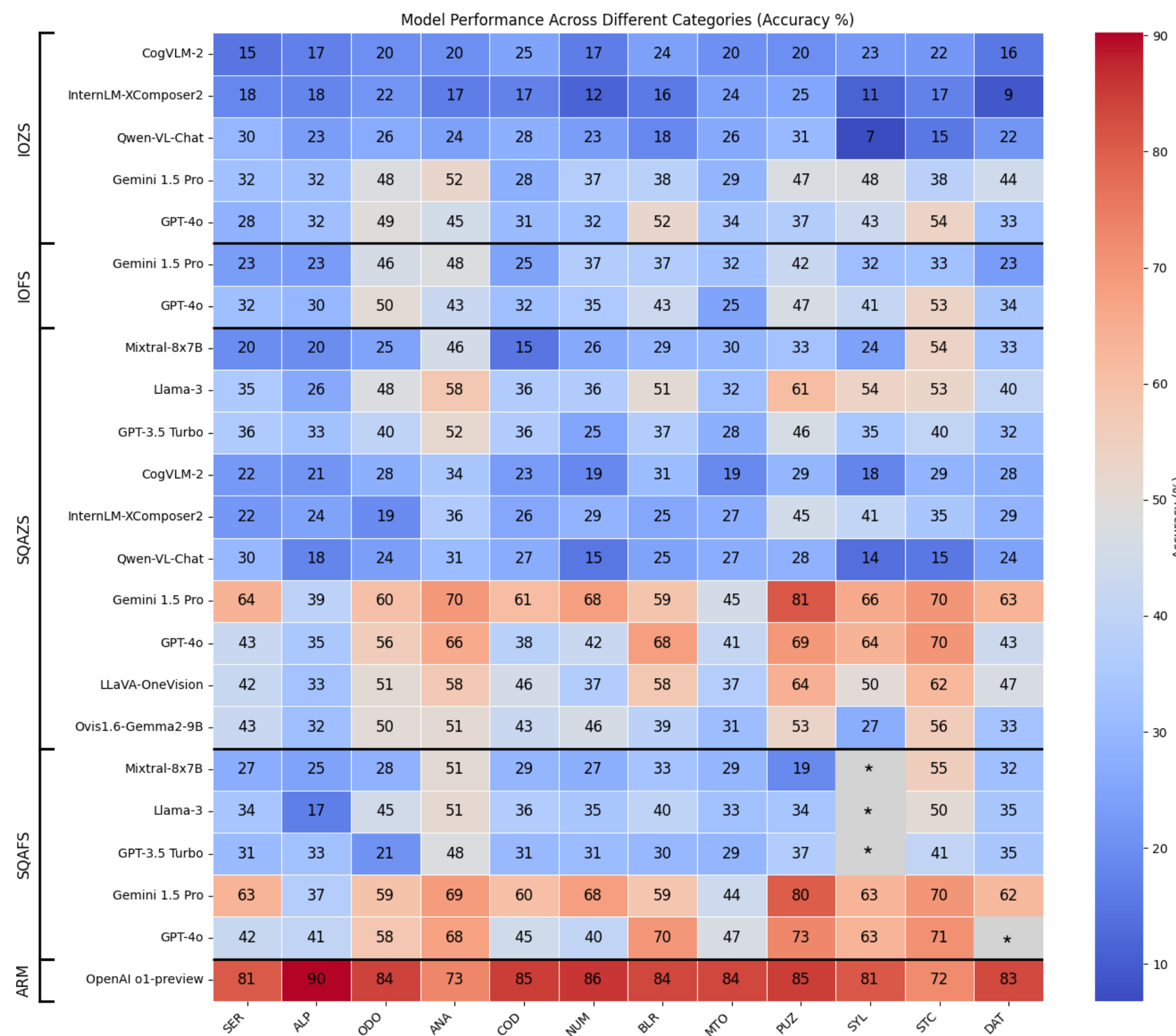
## 4. Results and Observations

### Text based questions



Figure: **IO**: Image-Only, **SQA**: Standard QA, **ARM**: Advanced Reasoning Model, **ZS**: Zero-Shot, **FS**: Few-Shot
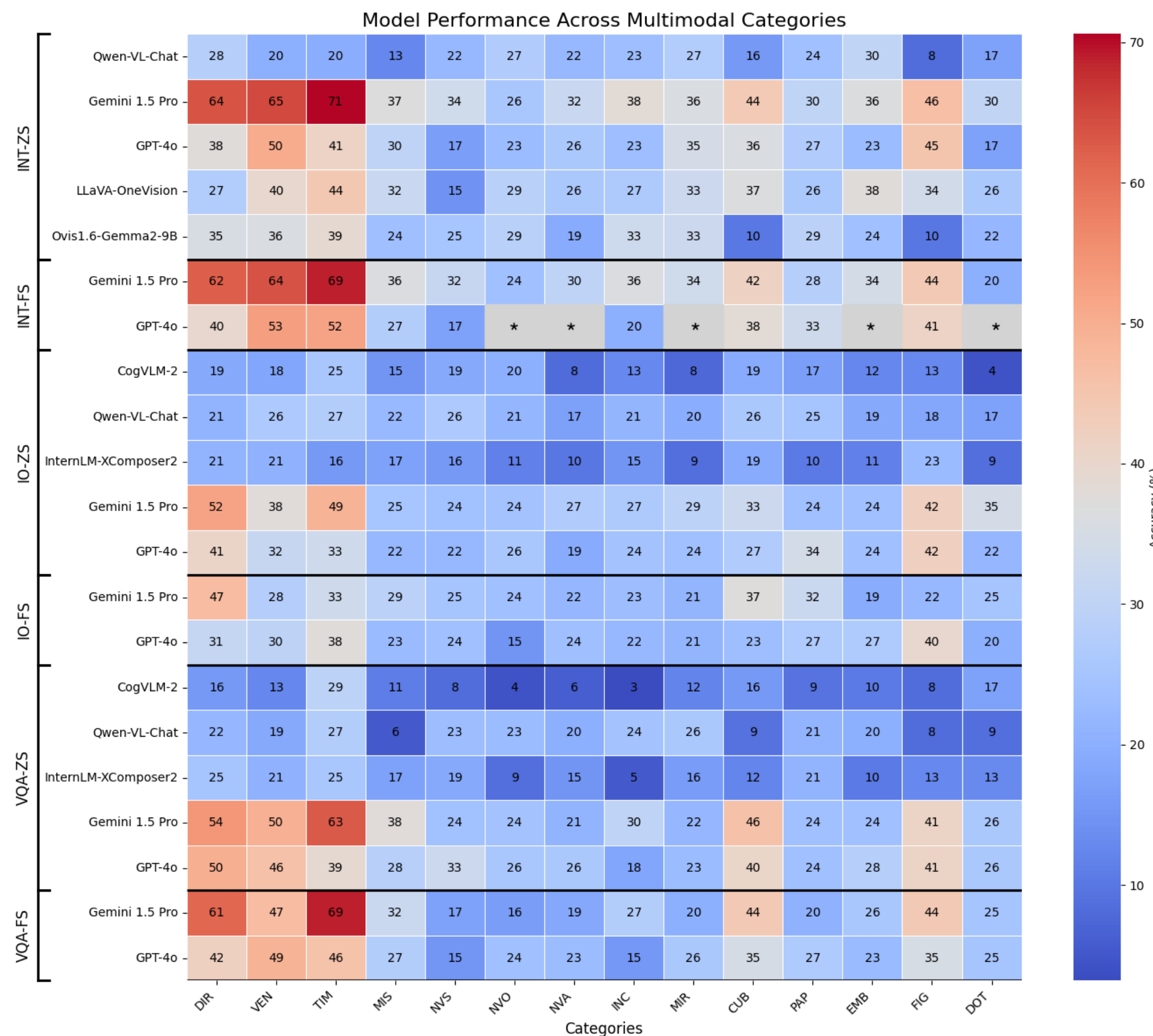
### Text+Image questions



Figure: **INT**: Interleaved, **IO**: Image-Only, **VQA**: Visual QA, **ZS**: Zero-Shot, **FS**: Few-Shot

### Key Observations:

- Proprietary models outperform open-source models.
- Interleaving text and images performs better than Standard VQA and Image Only.
- Multimodal reasoning is significantly harder.
- Human accuracy exceeds 80%, far surpassing the top proprietary model (62% text, 42% visual).
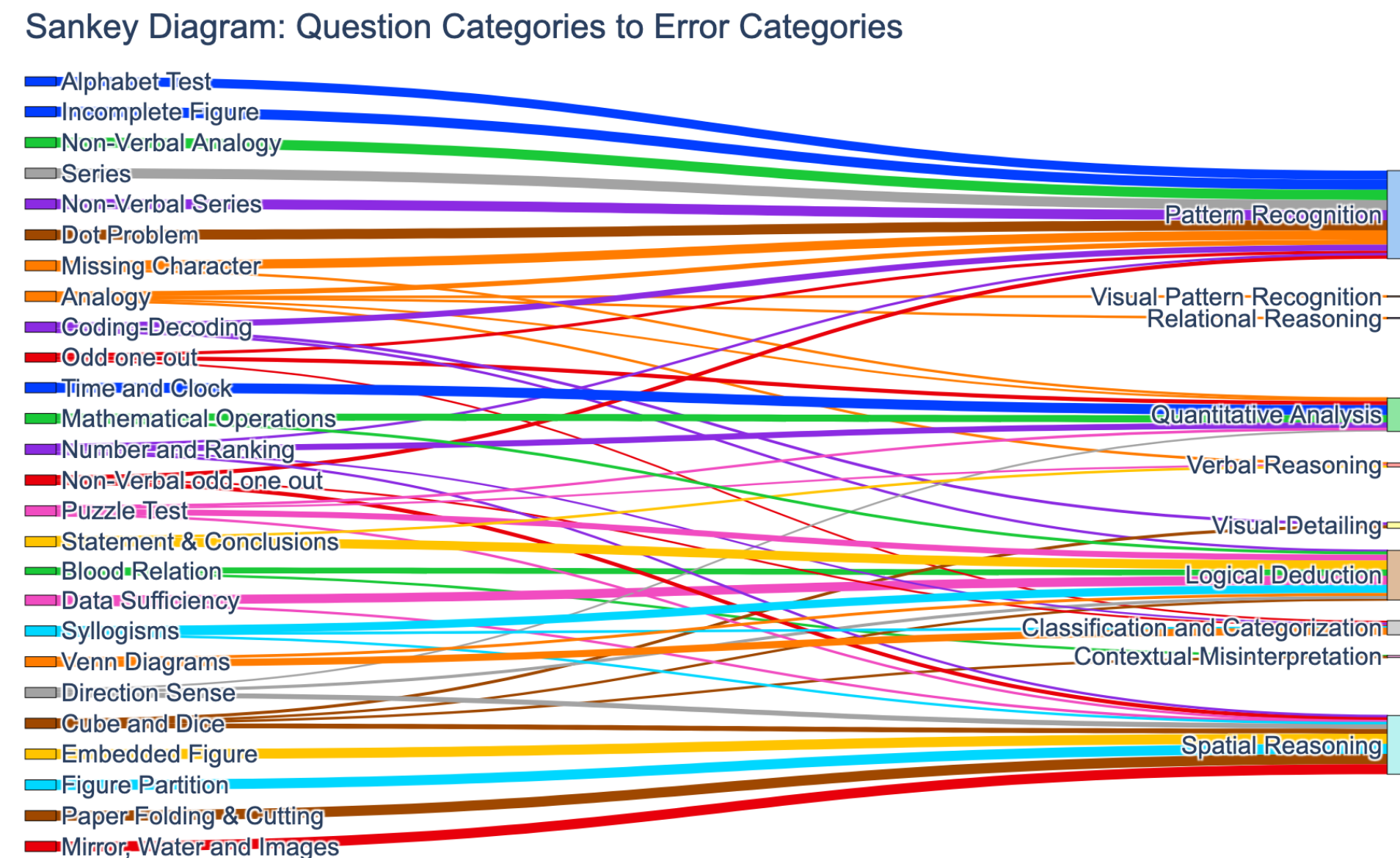
## Error Analysis:



Figure: **Overview of errors Gemini 1.5 Pro.**

- Challenges in *Pattern recognition, Spatial reasoning, and Logical deduction*.

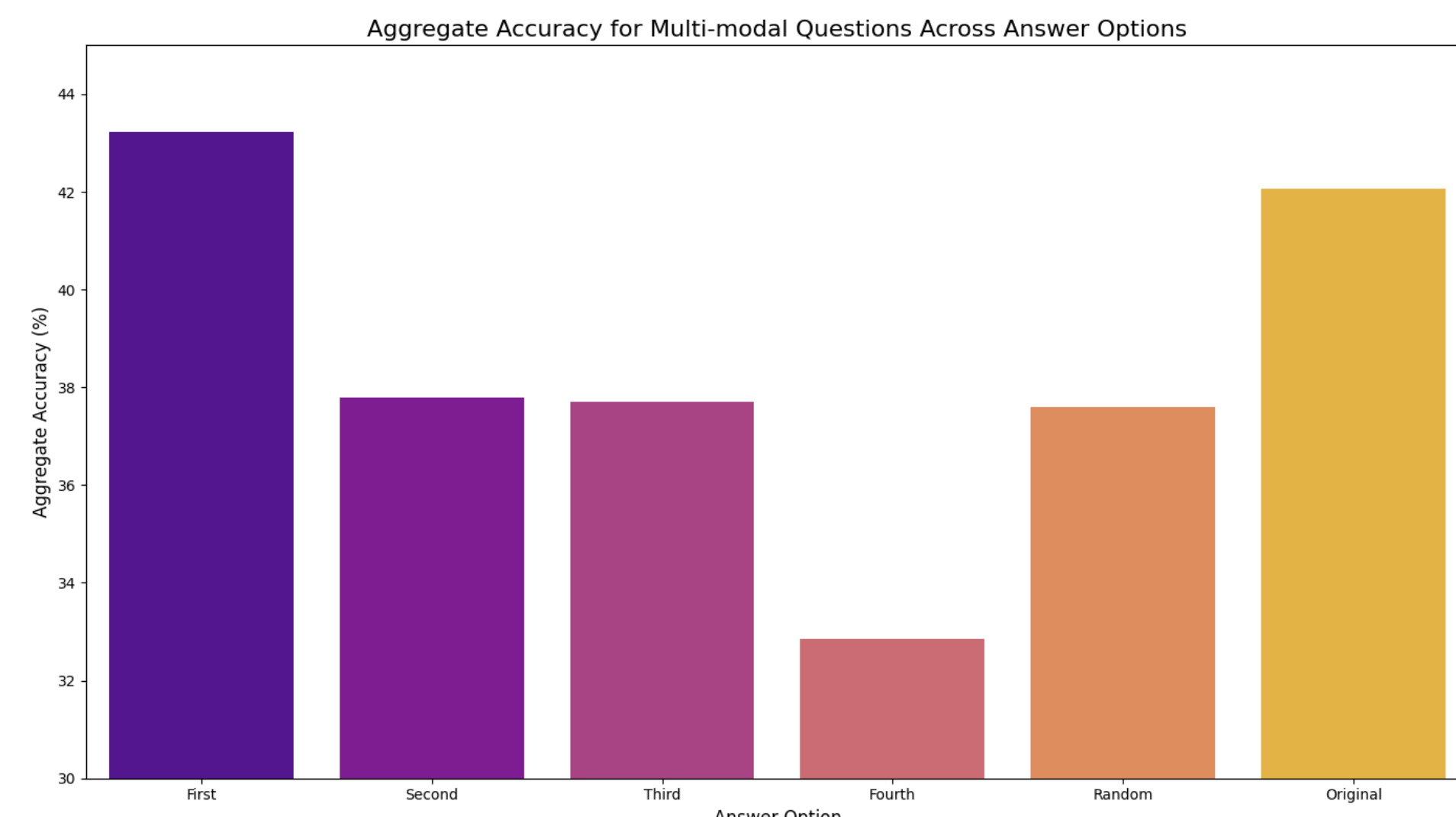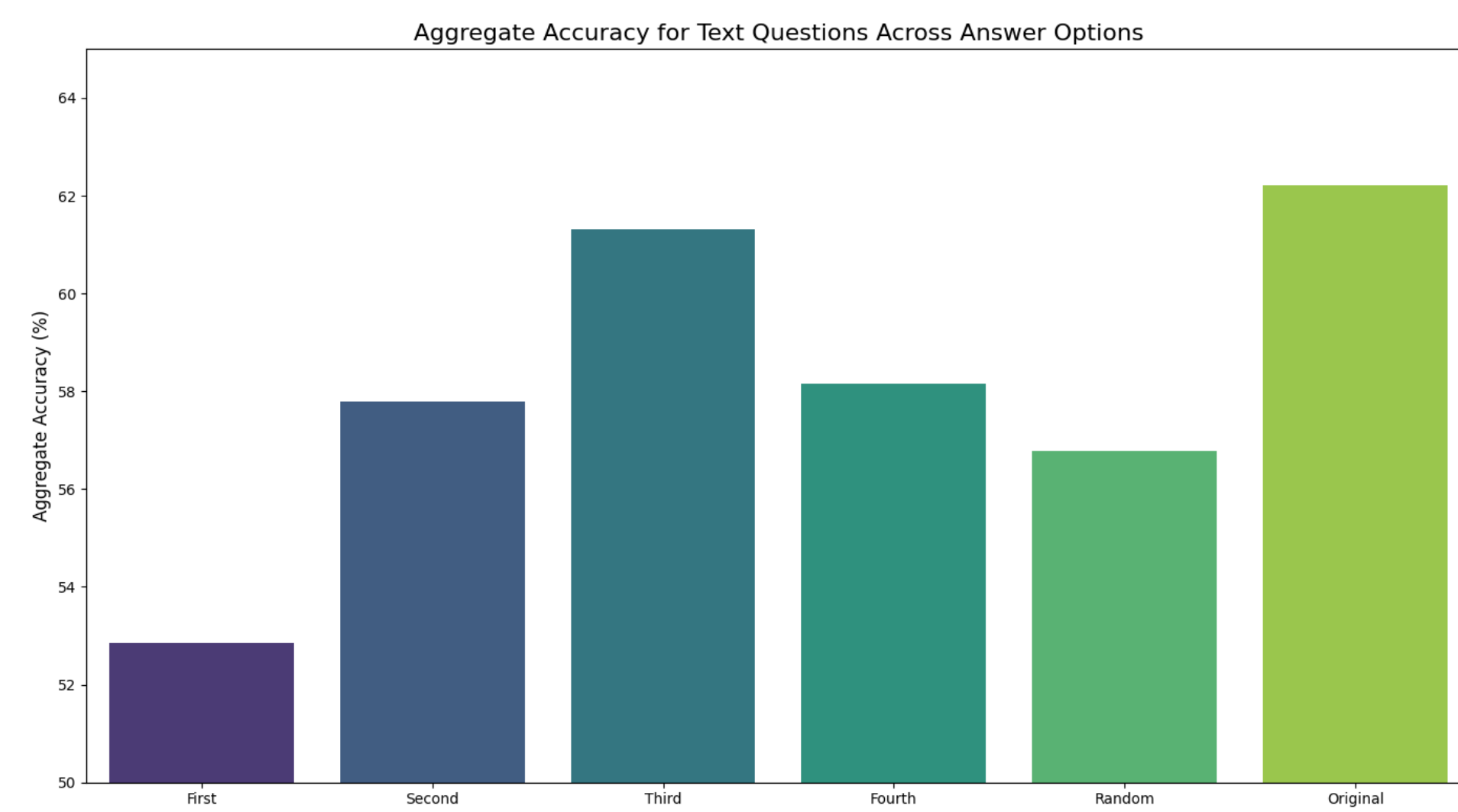### Option Shuffling on Gemini 1.5 Pro



Figure: Text based



Figure: Text+Image

- Option bias ablation shows answer placement impacts model performance.

## 5. Future Directions

- **Data Augmentation:** Expand the dataset with new tasks and perturbations.
- **Advanced Architectures:** Explore generative VLMs for complex tasks.
- **Fine-tuning:** Assess fine-tuned performance on multi-modal reasoning.

## Website and Contact Info

**Website:** ntsebench.github.io

**Author Emails:**
{p.pandya,g.vatsal,t.agney}@alumni.iitg.ac.in ,
tkataria@cs.utah.edu, danroth@seas.upenn.edu,
vgupt140@asu.edu