

Investigating Abstraction Capabilities of the o3 Model Using Textual and Visual Modalities

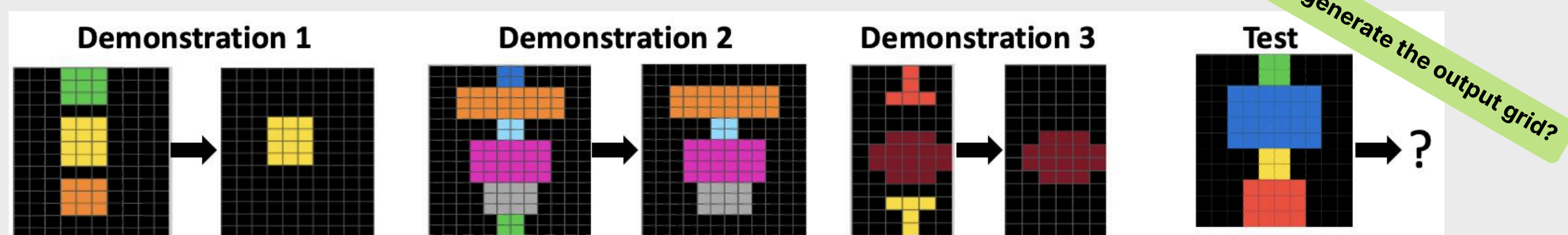
Claas Beger¹, Shuhao Fu¹, Ryan Yi¹, Arseny Moskvichev², Melanie Mitchell¹

¹Santa Fe Institute, ²Advanced Micro Devices

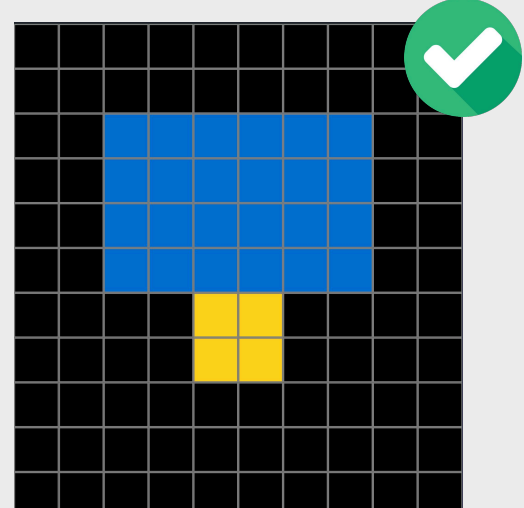
Motivation

OpenAI's o3-preview model outperformed humans on the **Abstraction and Reasoning Corpus**. But does it solve tasks using humanlike abstractions, or less generalizable "shortcuts"?

We investigated this using the **ConceptARC benchmark**, in which simple tasks focus on a single concept (e.g., "top vs bottom").



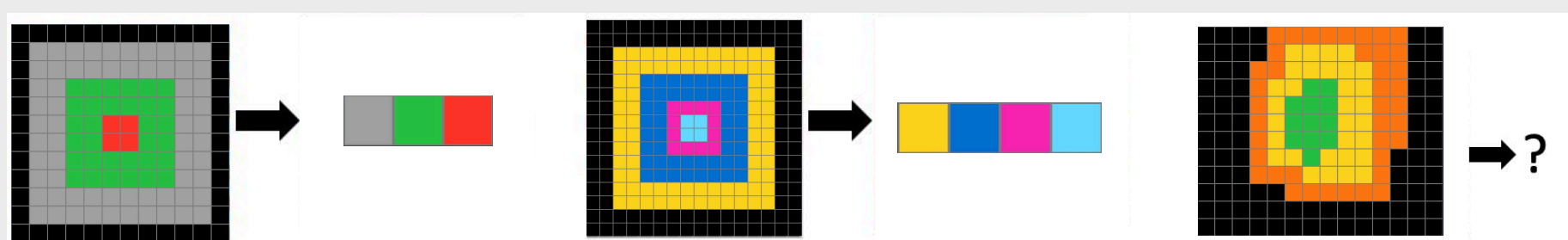
o3 generates the correct grid... ...but it does not use the intended abstraction "top/bottom" and fails on demonstrations 1 and 3!



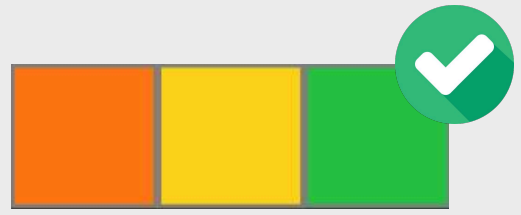
o3
"Delete every color region that touches any edge of the grid, leaving only regions entirely surrounded by background."

Human
"Copy the grid from the example then remove the topmost and bottom-most shapes. The middle shapes remain unchanged."

How about abstractions that correctly explain the input demonstrations?



o3 generates the correct grid... ...but it does not meet the intended rule, despite solving demonstrations (correct-unintended)!



o3
"List all non-background colours in descending order of frequency"

Human
"The grid should be 1 block (height) x width (...). Fill the grid from left to right with the colours, starting from the one used on the outside (...) and finishing with the colour on the inside."

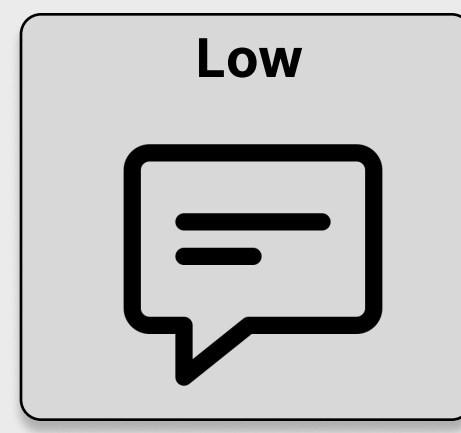
Methodology

Experiments

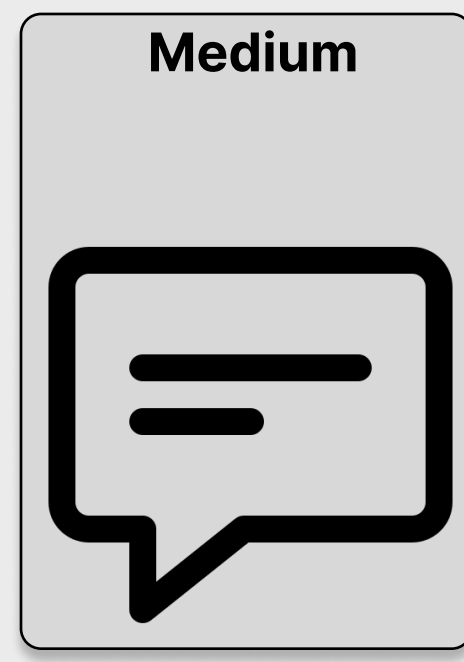
We test o3 on ConceptARC, systematically varying:



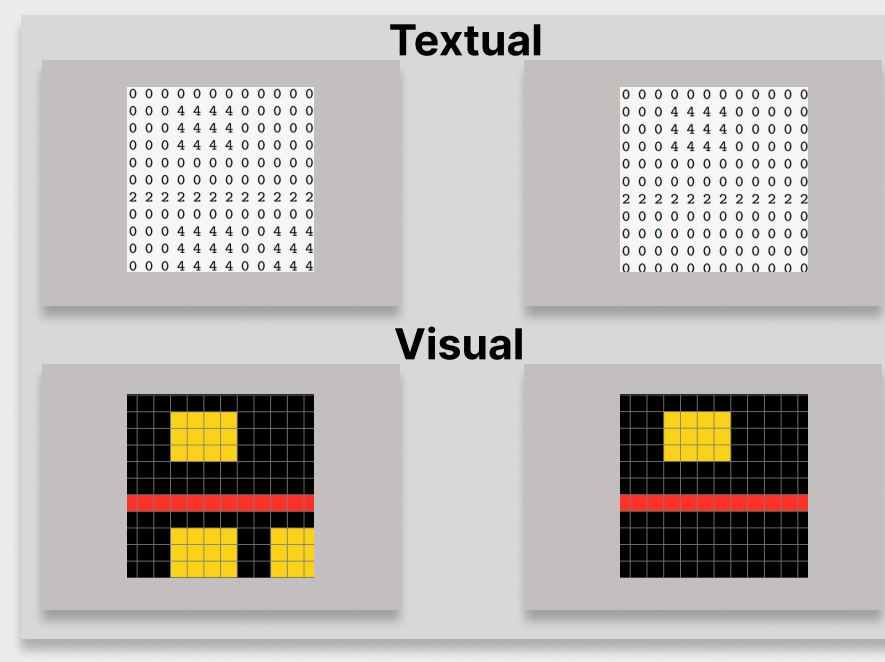
Enabling Python Code
Tool Usage



Reasoning Token Budget



Medium



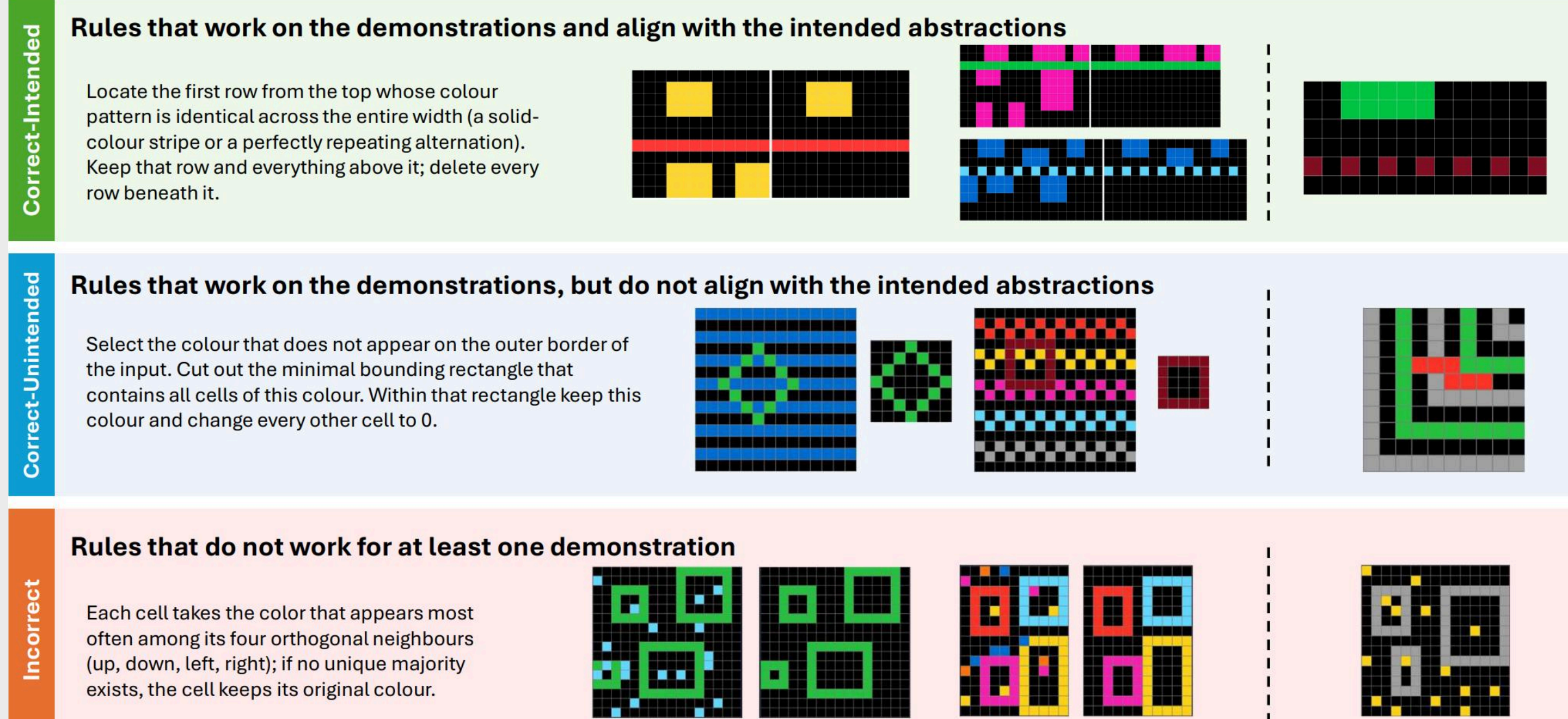
Input Modality

Output

We prompt o3 to generate both **output grid** and **natural-language rule** for each task

Evaluation

Output-Grid Accuracy: measured via exact match with ground truth
Rule Correctness:



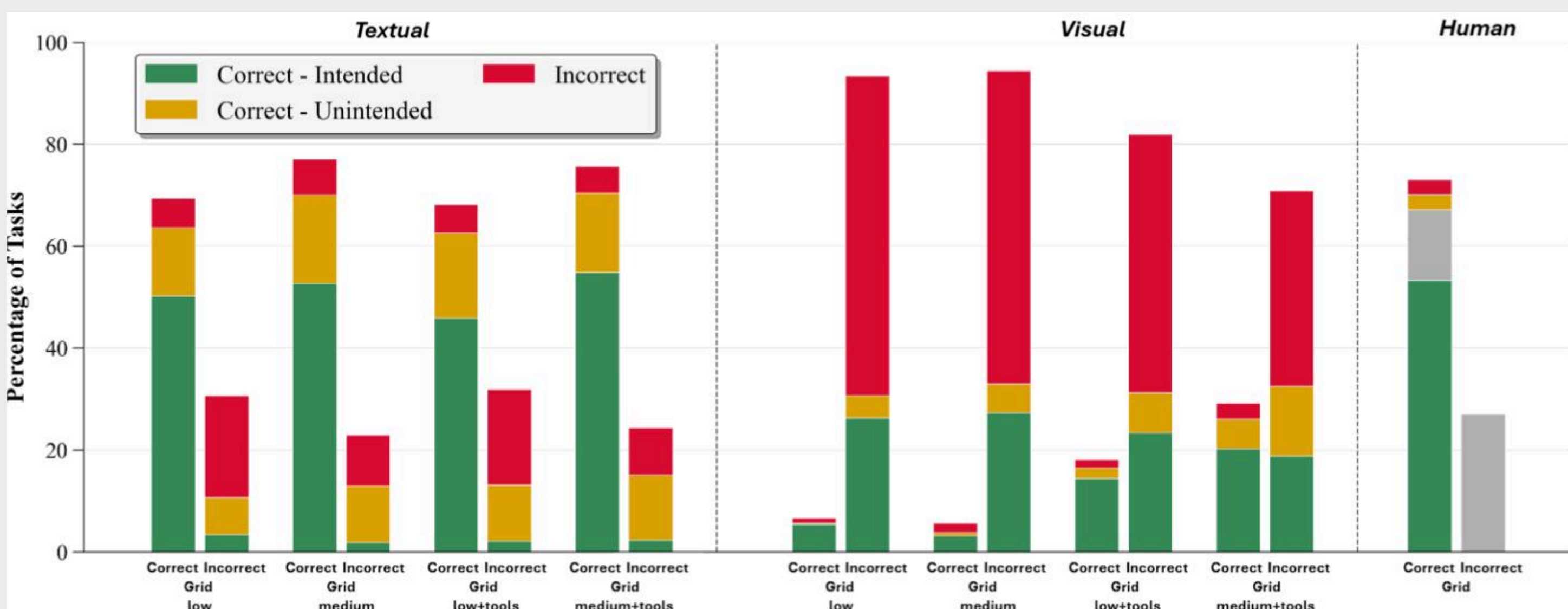
Results

Output-Grid Accuracy

Model	Modality	Low effort	Medium effort	Low effort + tools	Medium effort + tools
o3	Textual	68.3%	77.1%	67.9%	75.6%
	Visual	6.7%	5.6%	18.1%	29.2%

Textual modality improves with reasoning effort.
Visual modality improves with Python tool use.

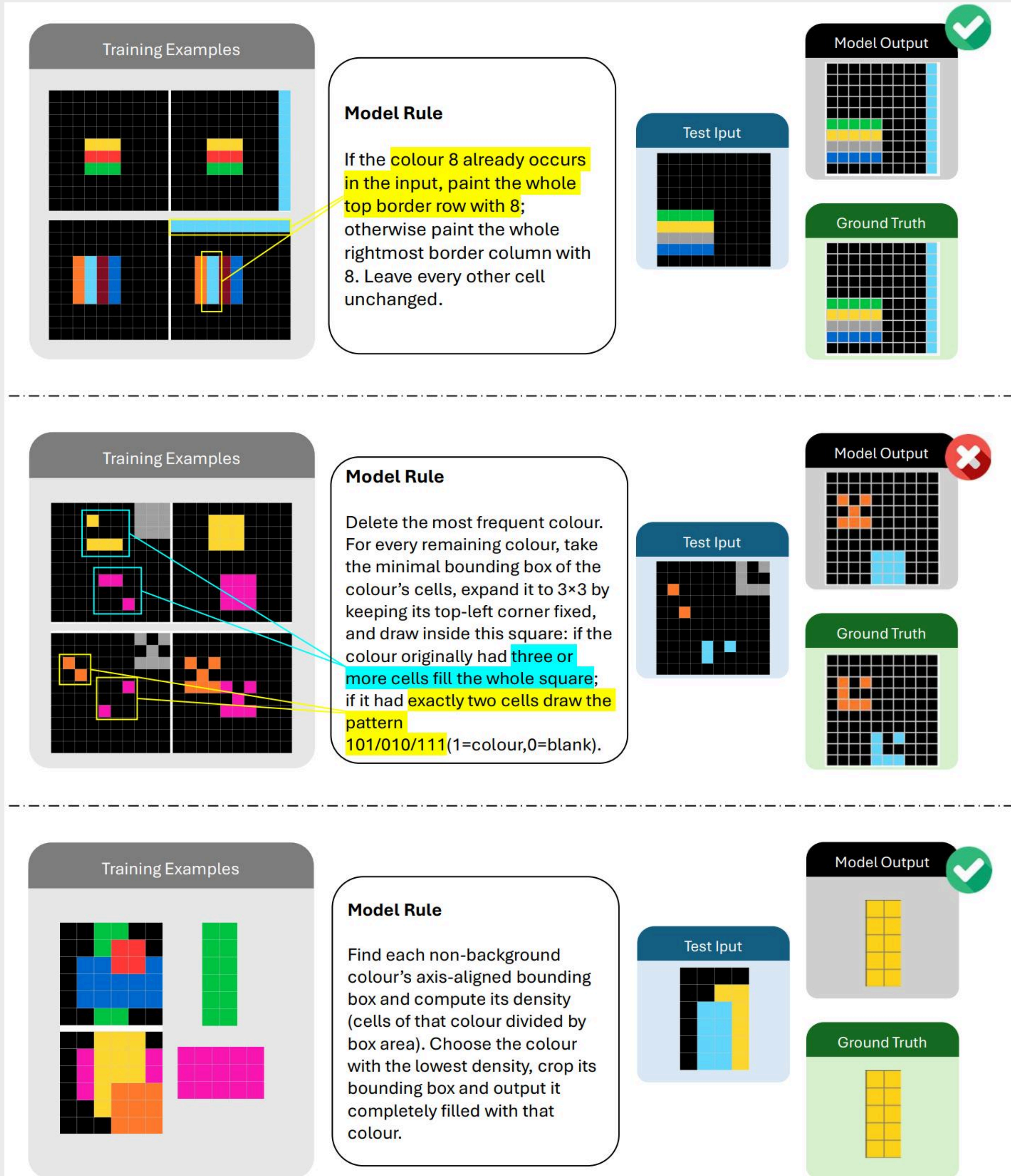
Rule Correctness



For each experimental setting the two bars show the percent of correct and incorrect grid outputs over the 480 ConceptARC tasks. Each bar shows the fraction of tasks for which the rule is *correct-intended*, *correct-unintended*, and *incorrect*. Gray areas in the human-result bars represent rules that we could not classify due to missing data or nonsensical responses.

Textual modality: Accuracy alone overestimates o3's abstract reasoning ability
Visual modality: Accuracy alone underestimates o3's abstract reasoning ability

Examples of o3's "Shortcuts"



Conclusion

Main Takeaway

Evaluating abstract reasoning tasks **solely based on accuracy can be misleading**

- Regardless of modality, **o3 generates correct-intended rules at a lower rate than humans** (around 70% in models vs 90% in Humans)
- o3 frequently produce **shortcuts or heuristics in place of humanlike abstractions**

Textual Modality

While o3 achieved high output-grid accuracy on ConceptARC, it often did so via non-humanlike reasoning. Thus, output accuracy alone **overestimates o3's abstract reasoning abilities** in this modality.

Visual Modality

While o3's output-grid accuracy was very low, it **still produced correct-intended rules fairly often**. Thus, output accuracy alone **underestimates o3's abstract reasoning abilities** in this modality.

Effects of Experimental Variations

- Increasing reasoning effort** primarily has a positive effect in the **textual** modality
- Allowing Python-tool usage** mainly helps in **visual** modality