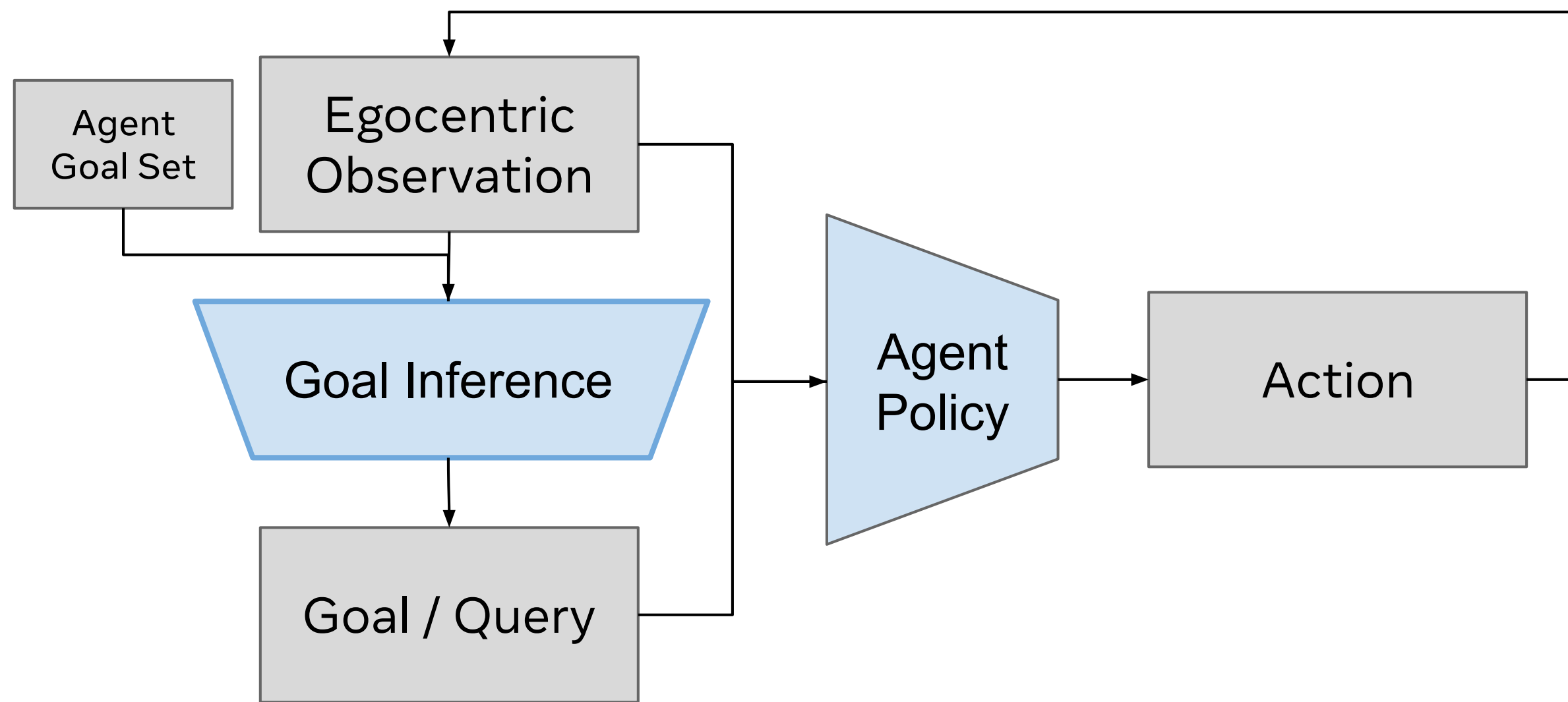
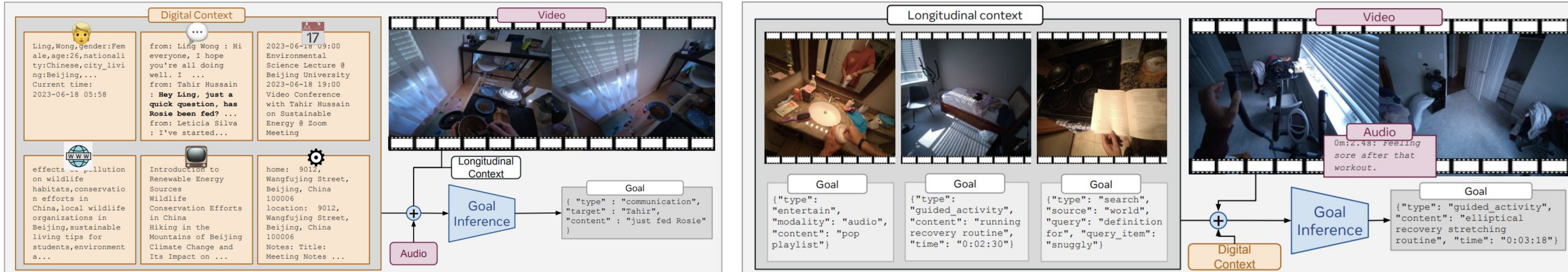


# WAGIBench: Wearable Agent Goal Inference Benchmark for Assistive Wearable Agents

## Infer user goals without explicit queries!



## Identify modalities with sufficient context for proactive goal inference!



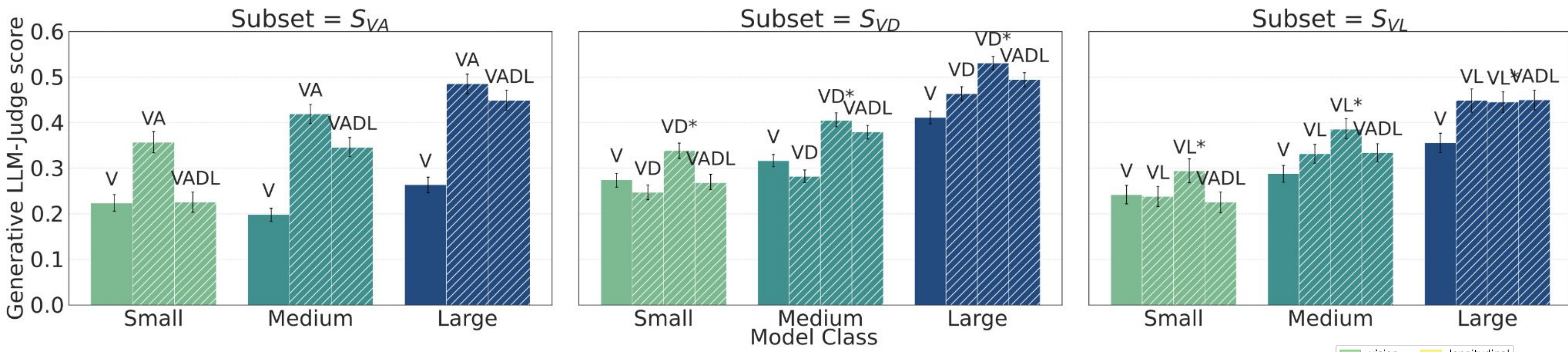
## Goal inference with audiovisual, digital (left) and longitudinal (right) context

- Introducing the first large-scale scripted proactive **Wearable Agent Goal Inference Benchmark: WAGIBench**.
- Scripts capture diverse real-life scenarios. User's rich **multimodal** context comprised of goal-relevant cues among distractors.
- 178 unique scenarios spanning 221 hours of video from ~300 participants wearing Aria glasses.

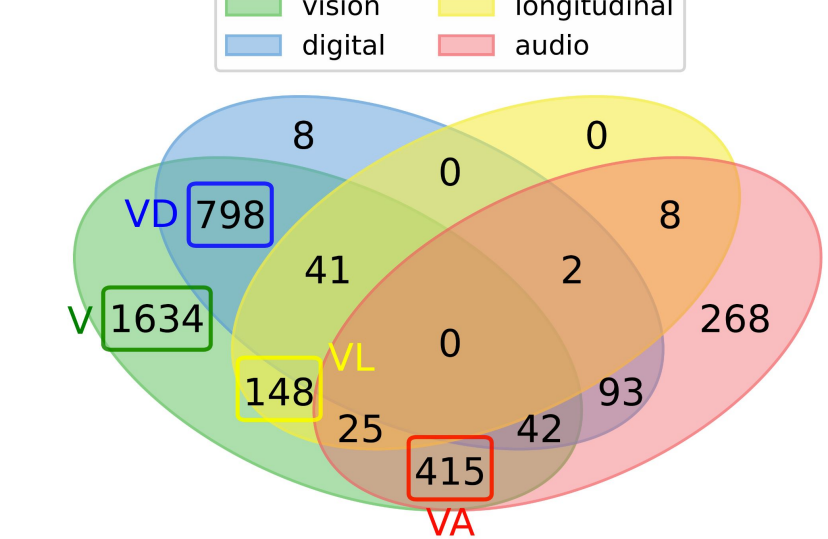
### Comparison of dataset statistics with prior work

Paper	Dataset	Videos	Questions	Ground Truth	Task	Modalities
MM-Ego	Ego4D	629	7,026	LLM (narrations)	Agent Policy	
EgoLife	EgoLife	6	6,000	LLM (captions)	Agent Policy	× T
PARSE-Ego4D	Ego4D	10,133	19,255	LLM (narrations)	Goal Inference	or
Ours	Ours	3,477	3,477	Scripted	Goal Inference	,  ,  × T

## Results – Modality ablation on generative goal inference across model sizes

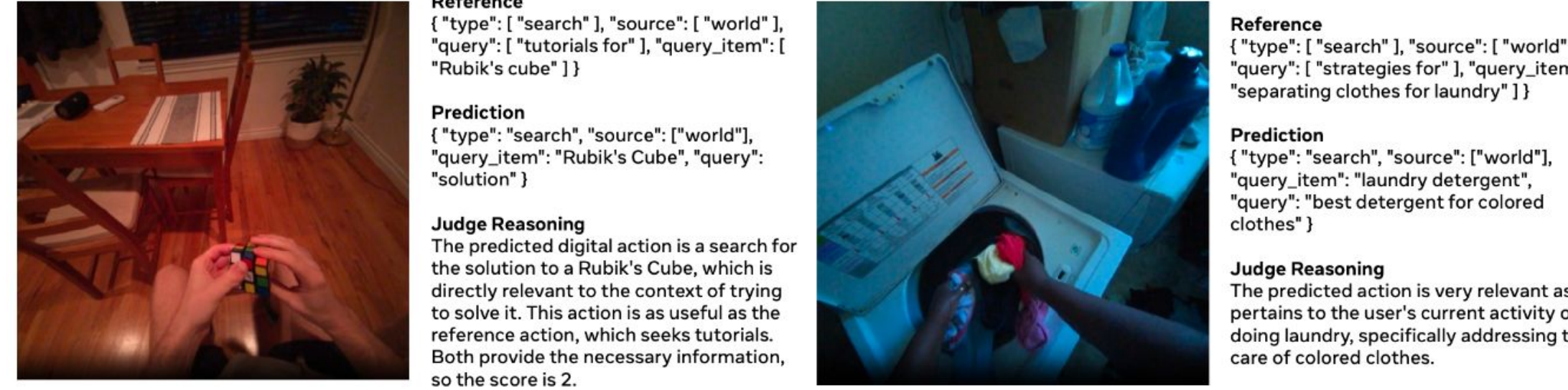


- Performance improves with model size in our scaling law experiments.
- Multi-modal context (e.g. Vision+Audio) strengthens performance
- Large models suffer less interference from mixed modalities, better disentangle relevant features.



## Visualizations

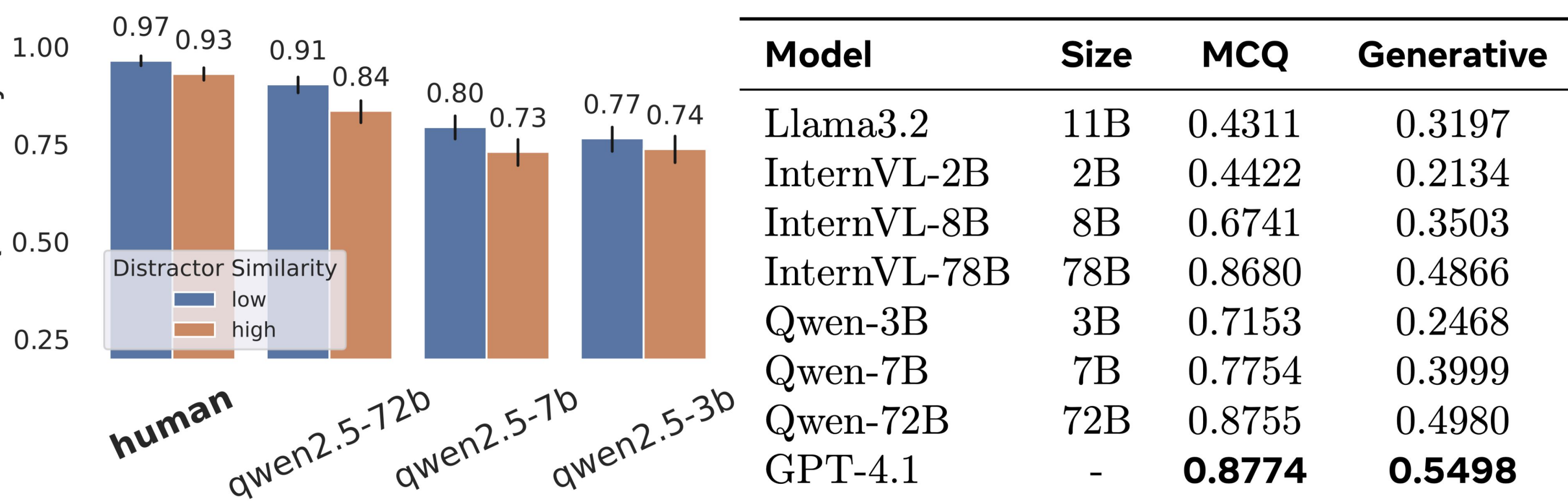
### Goal inference examples with only Vision context



### Goal inference examples with Digital contexts

Reference Goal	V Modality	D* Modality
<pre>{   "type": "search",   "source": "world",   "query": "how to play",   "query_item": "Clue Rivals Edition" }</pre>		<pre>{   "current_datetime": "2023-10-31 09:09",   "calendar": [     {       "title": "Playing Clue Rivals Edition with Kaito",       "start_datetime": "2023-11-01 19:00",       "end_datetime": "2023-11-01 21:30",       "location": "Game Center, Shinjuku-ku, Tokyo",       "attendees": [         "Kazuki Nakamura",         "Kaito Matsumoto"       ]     }   ] }</pre>
<b>Cues</b> ("Calendar": "The user sees a calendar event 'playing Clue Rivals Edition with friend'")  <b>Comment</b> Too many distractor objects in vision (for V) and large digital state for VD. Only VD* has the exact digital context to predict meaningful goals.	<b>Prediction with VD</b> <pre>{   "type": "entertain",   "modality": "audio",   "vision": "vision",   "content": "watching movies or shows" }</pre>	<b>Prediction with VD*</b> <pre>{   "type": "entertain",   "modality": "audio",   "vision": "vision",   "content": "board game strategy videos" }</pre>
Judge Score: 0.0	Judge Score: 0.0	Judge Score: 1.0

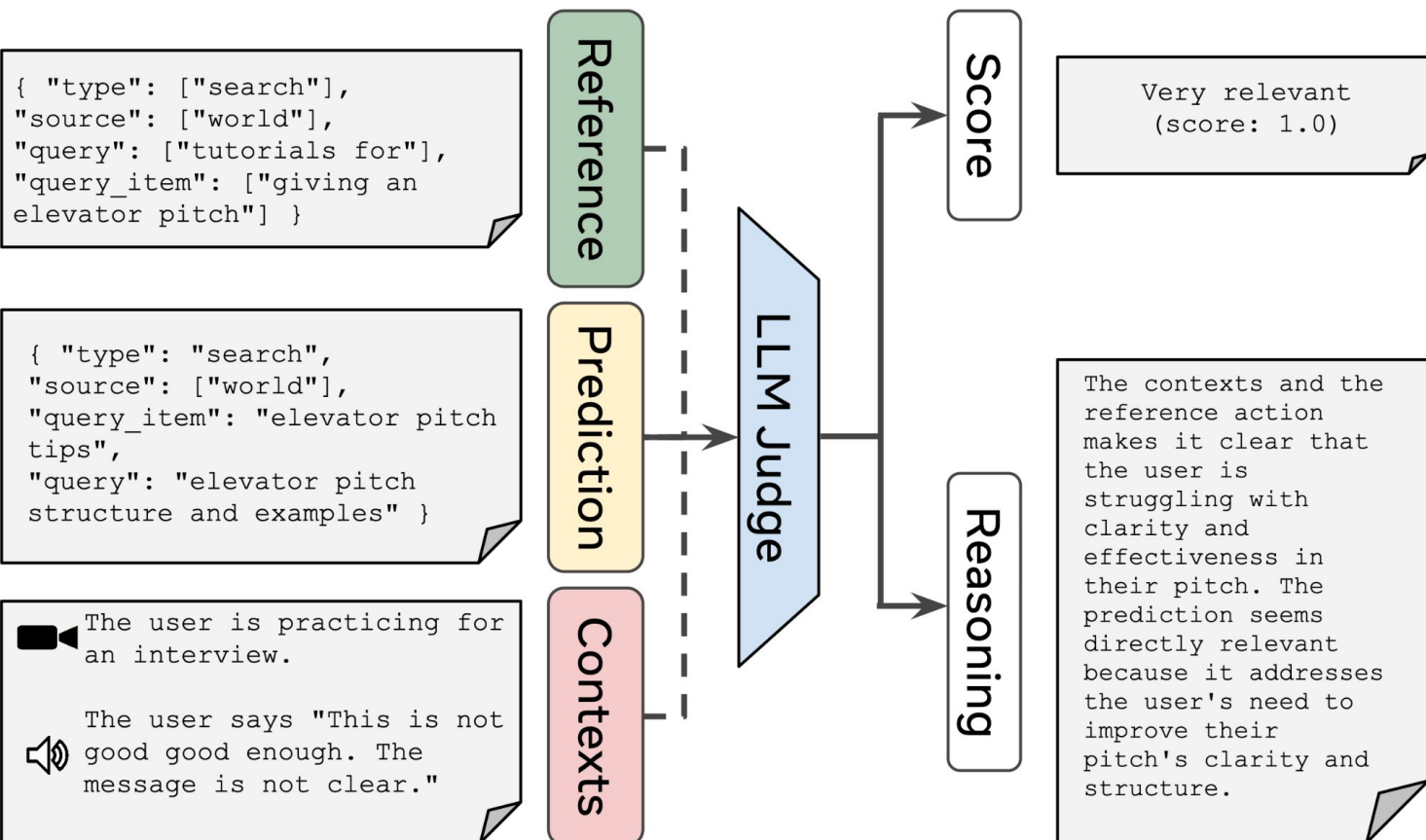
## Results – Multi-Choice Questions (MCQ) and Generative



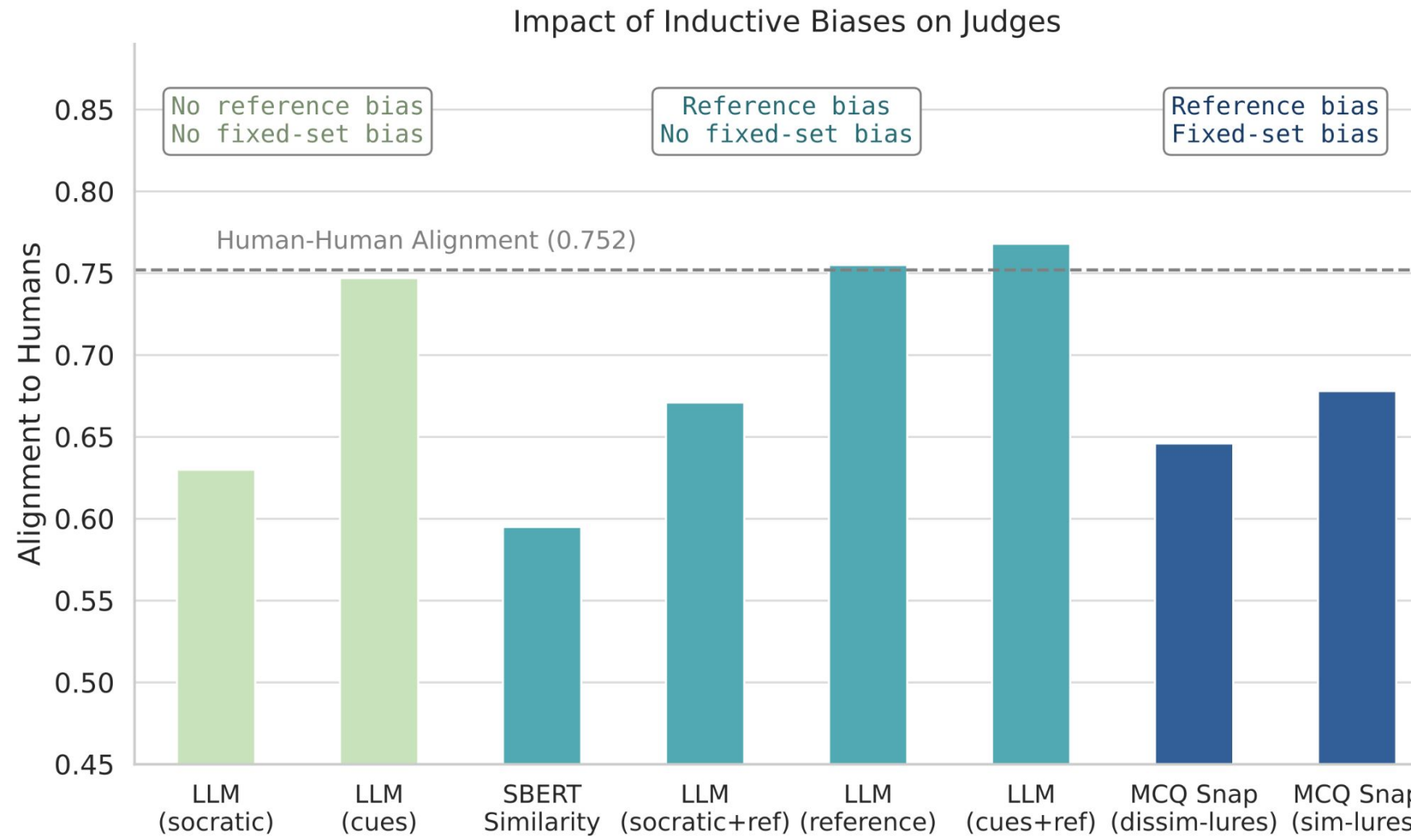
- Humans set an upper bound on goal inference in the MCQ setting
- Large VLMs trail humans yet still perform strongly on MCQ
- Even the most competent VLMs (GPT-4.1) scored only ~55% in the "generative" setting, implying significant room for improvement.

## Results – Meta Evaluation of LLM-Judges

### LLM-as-Judge for Generative Evaluation



### Alignment between Human raters and Judges



- LLM-Judges with access to the reference goal best align with human raters.
- The Judge model parameterized with both reference and script cues best aligns with human judgment (76.8%).

### Visualization of LLM-Judge ratings

