You May Speak Freely: Improving the Fine-Grained Visual Recognition Capabilities of MLLMs with Answer Extraction

Logan Lawrence

Oindrila Saha

Megan Wei

University of Massachusetts, Amherst

Chen Sun

Subhransu Maji

Brown University

Grant Van Horn



Problem

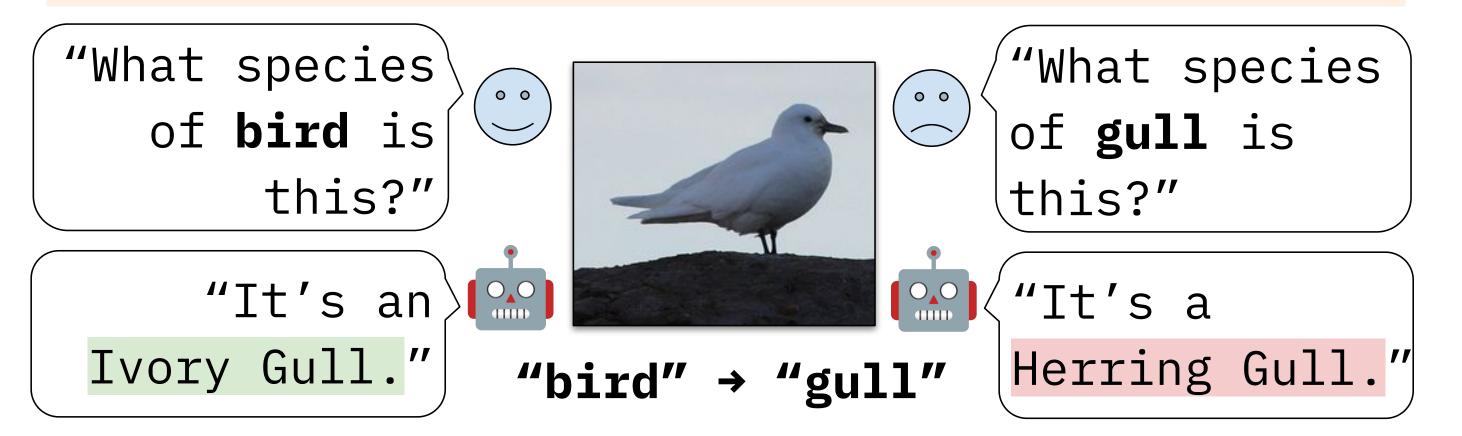
P1: FGVC choice counts are too big for classic MCQ evaluation.

Image: Coarse (# classes ≤ 5): Fine (# classes ≫ 200):

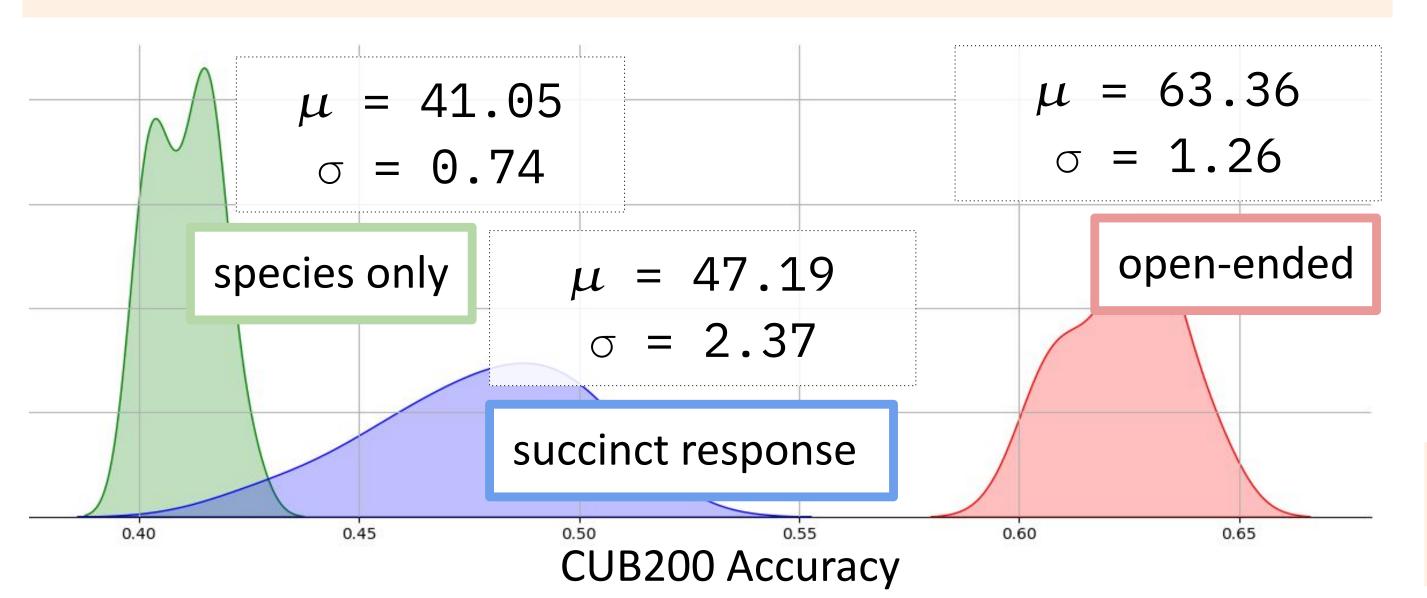


- (a) "Bird"
- (a)"Gray Catbird"
- (b) "Cat"
- (b)"N. Mockingbird"
- (c) "Dog"
- (d) "Airplane"
- (eb) "Ivory Gull"

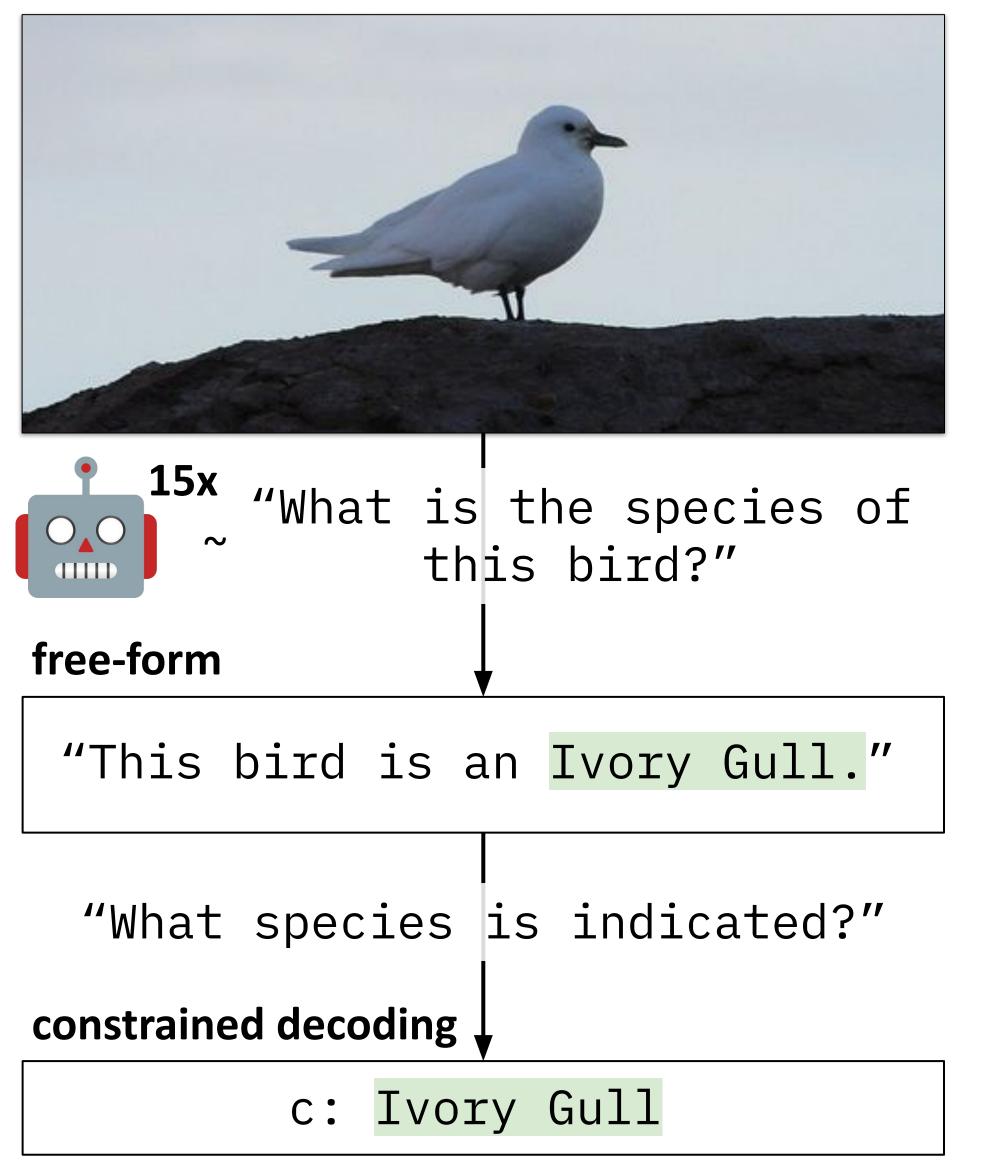
P2: MLLM outputs are not robust to small changes.



P3: These changes result in large performance differences.



Method



Datasets

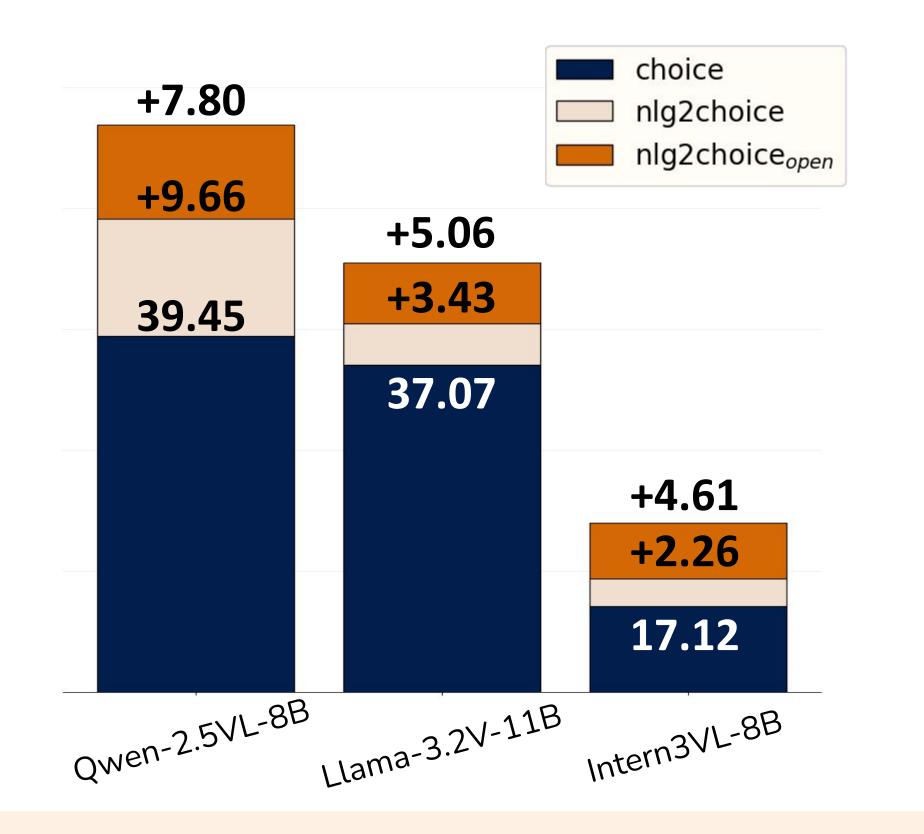
Benchmarked on **7 FGVC datasets:** CUB200, Flowers102,
FGVC Aircrafts, Stanford Cars,
Food101, NABirds, iNat21-Birds



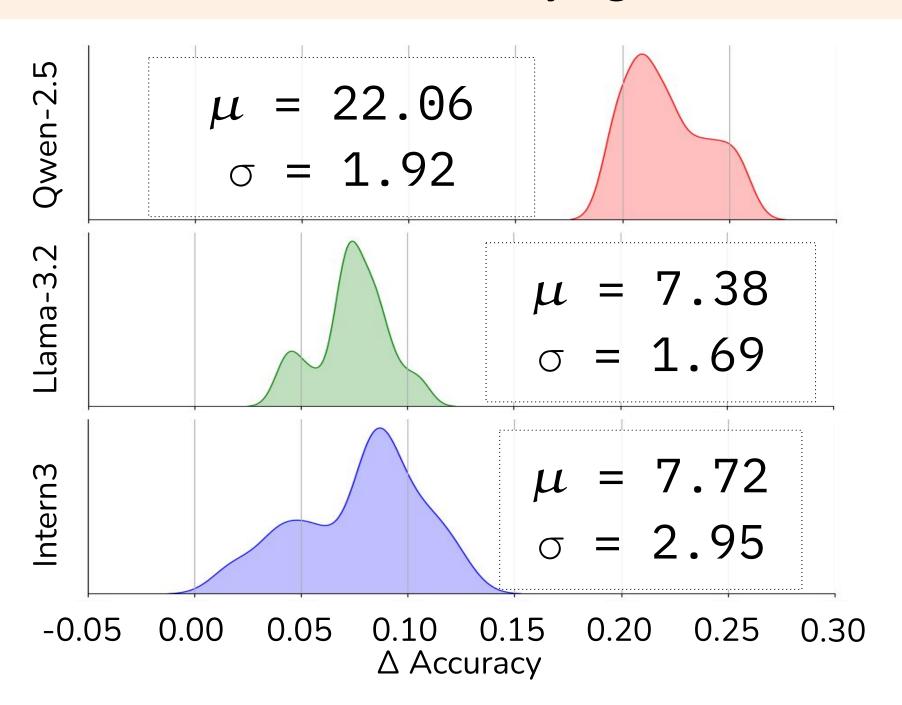
See project repo for our small labeled answer extraction dataset for fine-grained species!

Results

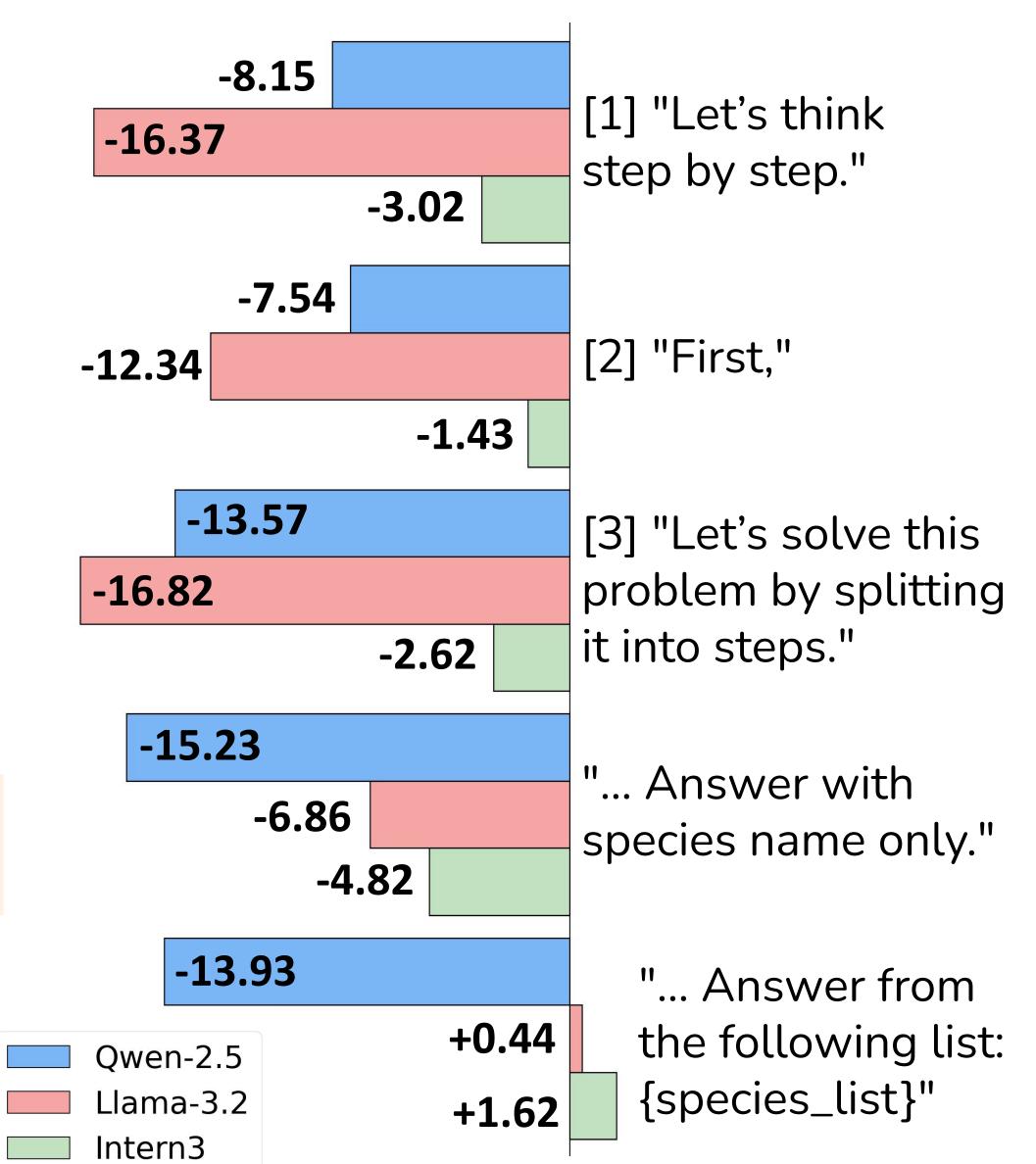




R2: With writing as source of randomness, results are **statistically significant**.



R3: Using CoT-inducing prompts or more specific instructions **degrades performance.**



References

- [1] Kojima et al. Large Language Models are Zero-Shot Reasoners.
- [2] Ahn et al. Do as I Can, Not as I Say: Grounding Language in Robotic Affordances.
- [3] Reynolds and McDonell, et al. Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm.