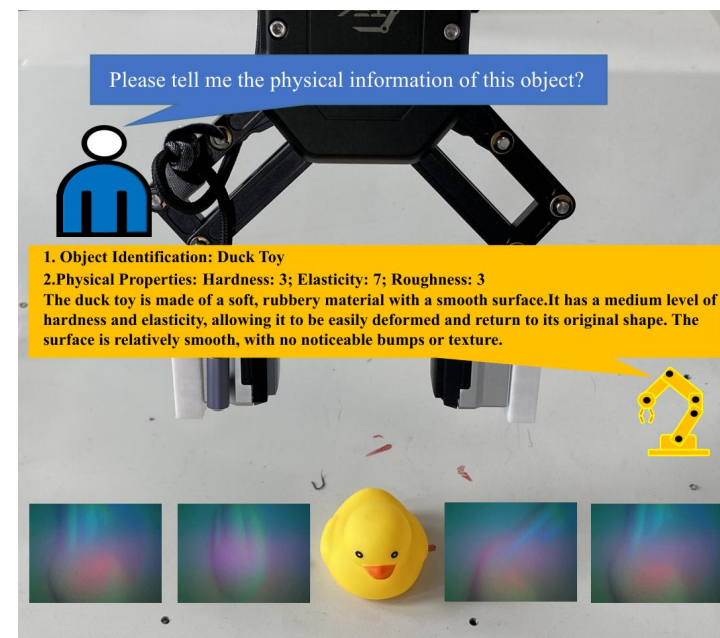
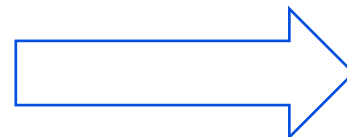


Octopi-X: Robotic Perception with a Large Tactile-Vision-Language Model for Physical Property Inference

Zexiang Guo^{1*}, Hengxiang Chen^{1*}, Xinheng Mai^{1*}, Qiusang Qiu¹, Gan Ma², Zhanat Kappassov³, Qiang Li^{1†}, Nutan Chen⁴



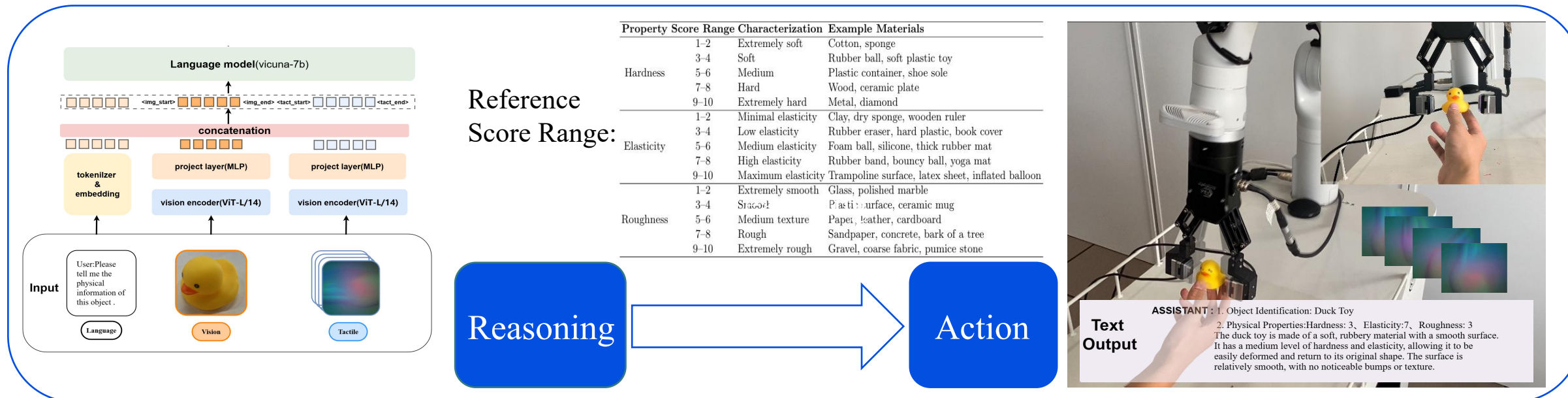
Project Homepage:



Hardness, elasticity and bumpiness ✓
 Only Tactile-Language reasoning ✗
 Pair-level reasoning ✗
 Zero-Shot Generalization Ability ✗

Hardness, elasticity and roughness ✓
 Tactile-Vision-Language reasoning ✓
 Instance-level reasoning ✓
 Zero-Shot Generalization Ability ✓

(i) Methodologies

(ii) Zero-Shot Generalization - Spearman Rank Correlation (ρ) with Ground Truth Measurements

35
unseen
objects

| Attribute | Method | Correlation Coefficient | P-value |
|------------|----------------------------|-------------------------|---------------|
| Hardness | Octopi-ViTAL | 0.501 | 0.005 |
| | Octopi-ViTAL (vision only) | 0.307 | 0.099 |
| | Octopi (fine-grained) | 0.307 | 0.099 |
| | Octopi (original) | 0.015 | 0.935 |
| Elasticity | Octopi-ViTAL | 0.530 | 0.003 |
| | Octopi-ViTAL (vision only) | 0.452 | 0.012 |
| | Octopi (fine-grained) | 0.053 | 0.781 |
| | Octopi (original) | -0.060 | 0.753 |
| Roughness | Octopi-ViTAL | 0.643 | 0.0001 |
| | Octopi-ViTAL (vision only) | 0.413 | 0.023 |
| | Octopi (fine-grained) | -0.010 | 0.959 |
| | Octopi (original) | 0.118 | 0.534 |

(iii) Conclusions

- i) Integrated vision, tactile, and language inputs → Improved **property inference**.
- ii) Our structured prompting → Achieved strong results on **hardness, elasticity, and roughness**.
- iv) Demonstrated **zero-shot generalization** across unseen objects.

Project Homepage:

