# CoWs on Pasture:
## Baselines and Benchmarks for Language-Driven Zero-Shot Object Navigation

Samir Yitzhak Gadre, Mitchell Wortsman,
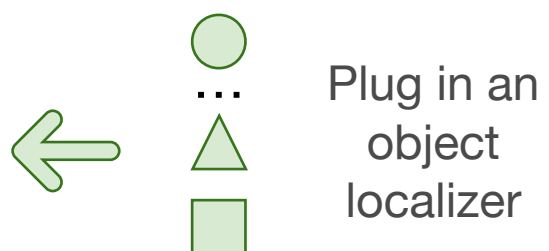Gabriel Ilharco, Ludwig Schmidt, Shuran Song

COLUMBIA UNIVERSITY IN THE CITY OF NEW YORK

UNIVERSITY of WASHINGTON

## We turn CLIP models into zero-shot object navigators without additional training

Paper:

Code:

**1 Task**

Give a natural language description of an object, with potentially many attributes, the task is to find the object.
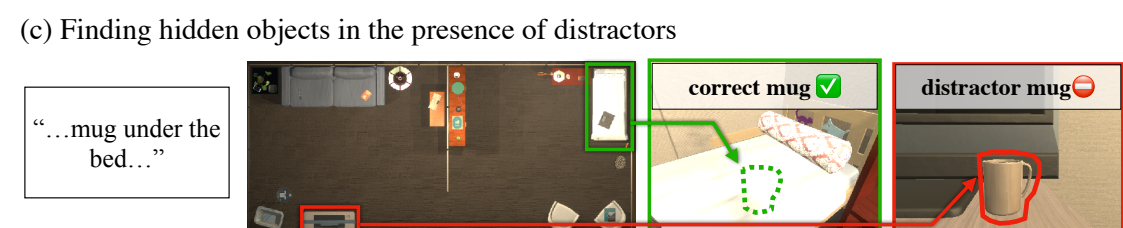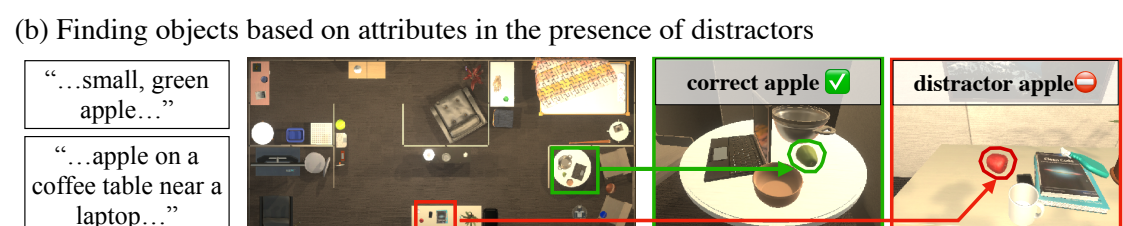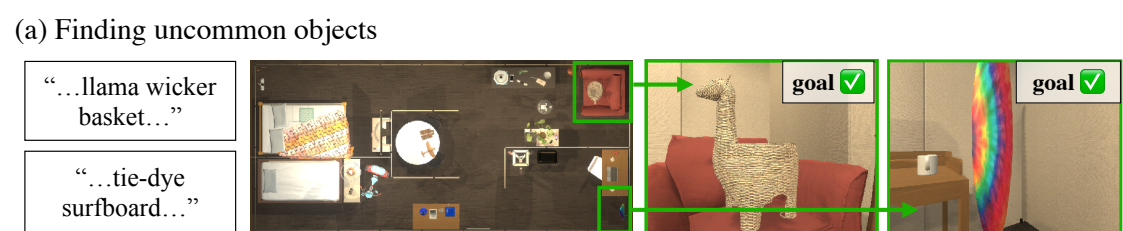
**2 CoW Baselines**

```
If object is in view:
    move to it
else:
    explore
```

··· Plug in an object localizer

··· Plug in a policy

**3 Pasture Benchmark**

Sample tasks



| Sample tasks | Top-down visualization | Egocentric Observations |
|---|---|---|

(a) Finding uncommon objects

"…llama wicker basket…"     goal ✅   goal ✅

"…tie-dye surfboard…"

(b) Finding objects based on attributes in the presence of distractors

"…small, green apple…"      correct apple ✅   distractor apple ⊖

"…apple on a coffee table near a laptop…"

(c) Finding hidden objects in the presence of distractors

"…mug under the bed…"       correct mug ✅   distractor mug ⊖

Uncommon object navigation targets



"whiteboard saying CVPR"   "red and blue tricycle"
"tie-dye surfboard"        "white electric guitar"
"llama wicker basket"      "espresso machine"
"green plastic crate"      "wooden toy airplane"
"rice cooker"              "gingerbread house"
"maté gourd"               "graphics card"

## CLIP on Wheels (CoW)

**Exploration strategies**

Where to search next?

$I_t$ → GRU → $a_t$

Learning-based        Frontier-based

**Language driven localization strategies**

Am I looking at the `object`?

Gradient-based   Detector-based   Reference-based

a `object` in the center?

**Explore** or **Exploit**?

Take a CLIP model, put it **on Wheels**, CoW!

**4** CoW baselines struggle to take advantage of spatial and appearance object attributes. Distractor objects hurt performance less when finding hidden objects.
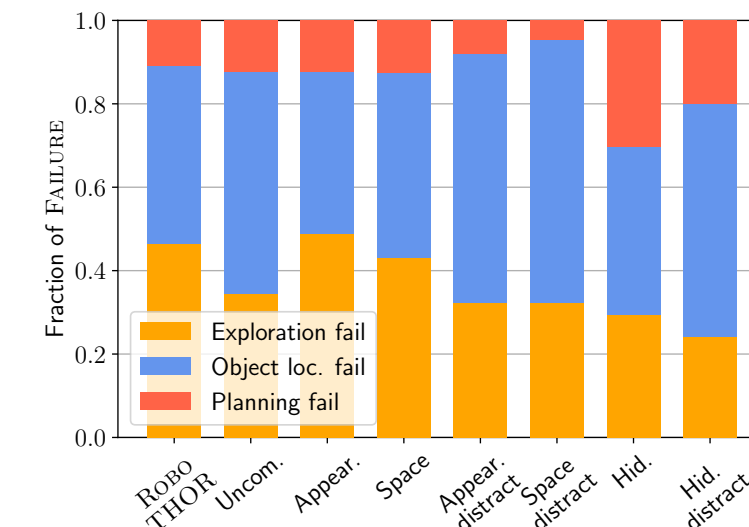
(a) Attribute object navigation

**5 LLM Priors**

Incorporating object and scene priors helps.

| | CoW breeds | | | PASTURE Uncom. | ROBOTHOR |
|---|---|---|---|---|---|
| ID | Loc. | Arch. | Obj. Prior | SPL | SPL |
| ▲ | OWL | B/32 | None | 20.5 | 32.8 | 16.8 | 26.7 |
| ▲ | OWL | B/32 | GPT-3.5 | 22.2 | 36.9 | 17.0 | 27.5 |

**6 Failure Analysis**

CoWs may improve with better perception and exploration strategies.



Exploration fail
Object loc. fail
Planning fail

**7 Comparison to Prior Art**

CoWs beat baselines that train for millions of steps.

| | CoW breeds | | HABITAT (MP3D) | | ROBOTHOR (subset) | | ROBOTHOR (full) | | Nav. training steps |
|---|---|---|---|---|---|---|---|---|---|
| ID | Loc. | Arch. | SPL | SR | SPL | SR | SPL | SR | |
| ▲ | CLIP-Grad. | B/32 | **4.9** | 9.2 | 15.0 | 23.7 | 9.7 | 15.2 | **0** |
| ▲ | OWL | B/32 | 3.7 | 7.4 | **20.8** | **32.5** | **16.9** | **26.7** | **0** |
| ▲ | EmbCLIP-ZSON [38] | – | – | – | – | 8.1 | – | 14.0* | 60M |
| | SemanticNav-ZSON [46] | 4.8 | **15.3** | – | – | – | – | – | 500M |