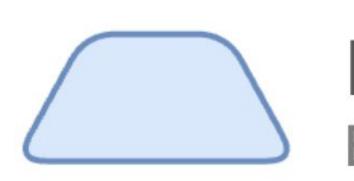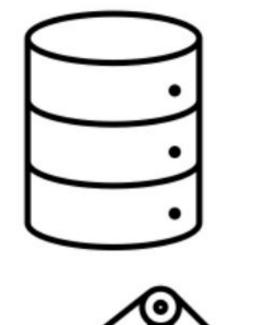# A Data-Centric Revisit of Pre-Trained Vision Models for Robot Learning

*by Xin Wen, Bingchen Zhao, Yilun Chen, Jiangmiao Pang, and Xiaojuan Qi*

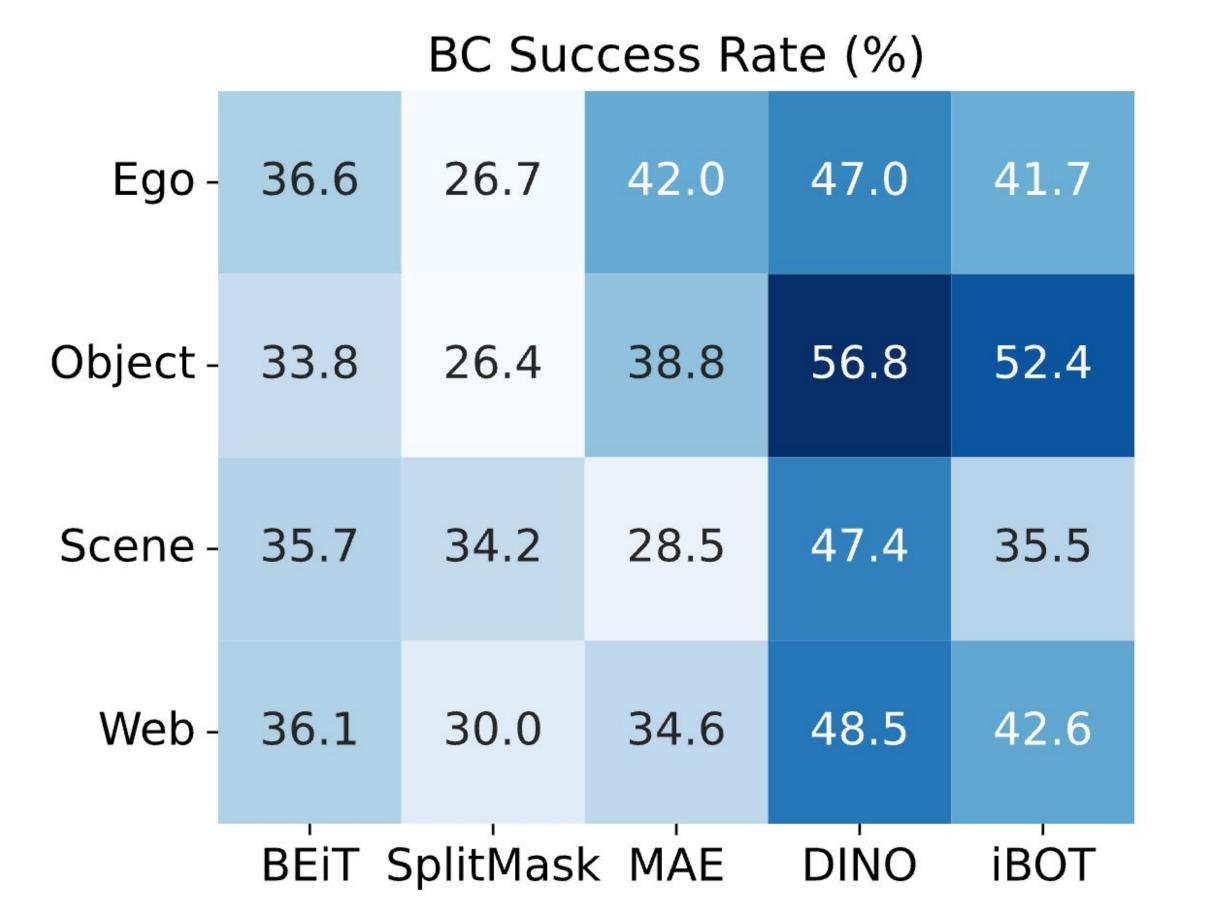## 1. How Pre-Training Data Affect Vision Models on Robot Tasks?

**PVM: 5 pre-training methods**
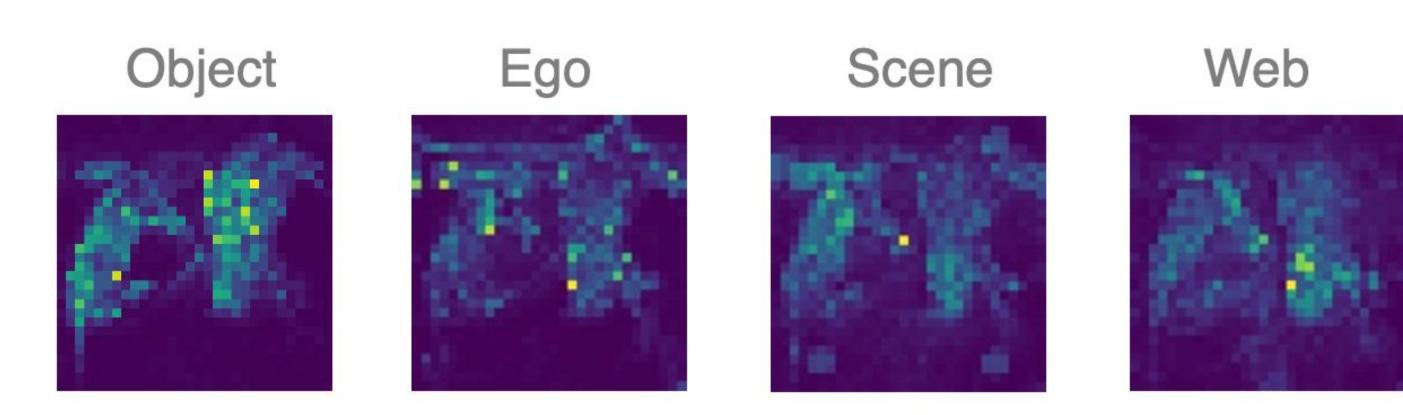BEiT, SplitMask, MAE, DINO, and iBOT

**Data: 4 pre-training datasets**
Ego-Centric, (Single-)Object-Centric, Scene-Centric, and Web-Crawled data

**Eval: 13 behavior cloning tasks**
Franka Kitchen and Meta-World



Pre-Training Data: Object-centric, Scene-centric, Web-crawled, Ego-centric

### BC Success Rate (%)

| | BEiT | SplitMask | MAE | DINO | iBOT |
|---|---|---|---|---|---|
| Ego | 36.6 | 26.7 | 42.0 | 47.0 | 41.7 |
| Object | 33.8 | 26.4 | 38.8 | 56.8 | 52.4 |
| Scene | 35.7 | 34.2 | 28.5 | 47.4 | 35.5 |
| Web | 36.1 | 30.0 | 34.6 | 48.5 | 42.6 |

Avg performance on manipulation tasks.
**DINO/iBOT rival other methods a lot!**

**Downstream Robot Learning Tasks**
Franka Kitchen, Meta-World, ObjectNav, ImageNav
Open Slide Cabinet, Pick Bin, Goal: dining table, Goal

**Evaluation Protocol: Behavior cloning with attentive probing on frozen PVMs.**
Report success rate on multiple trials.

## 2. Objectness Matters
But is hard to obtain on non-object-centric (NOC) data
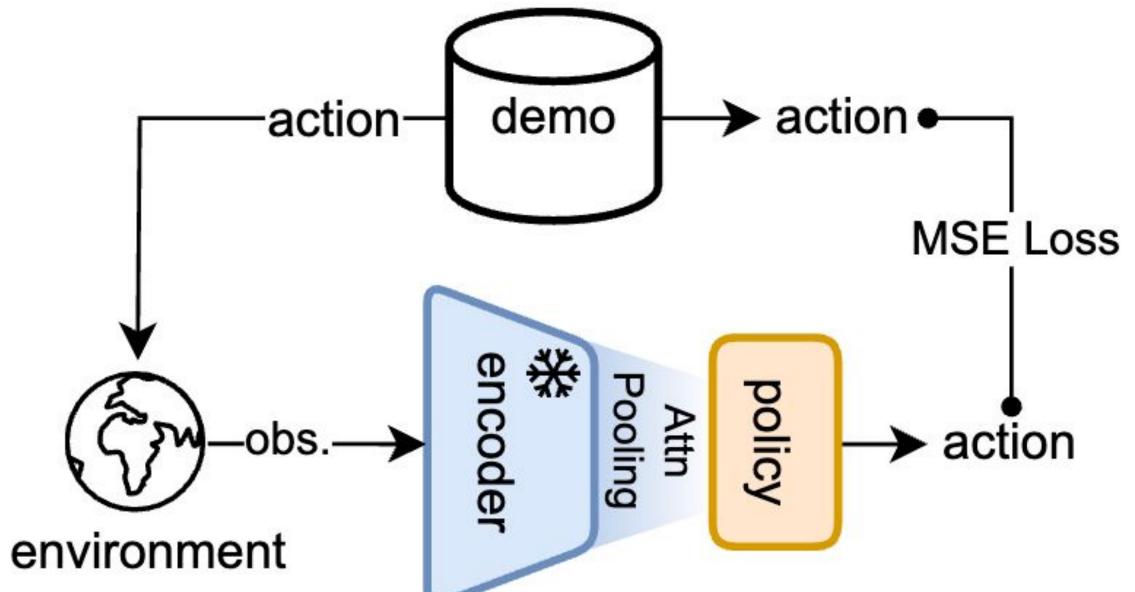


**Top: Attention Masks of DINO**
Poor objectness on non-object-centric data.
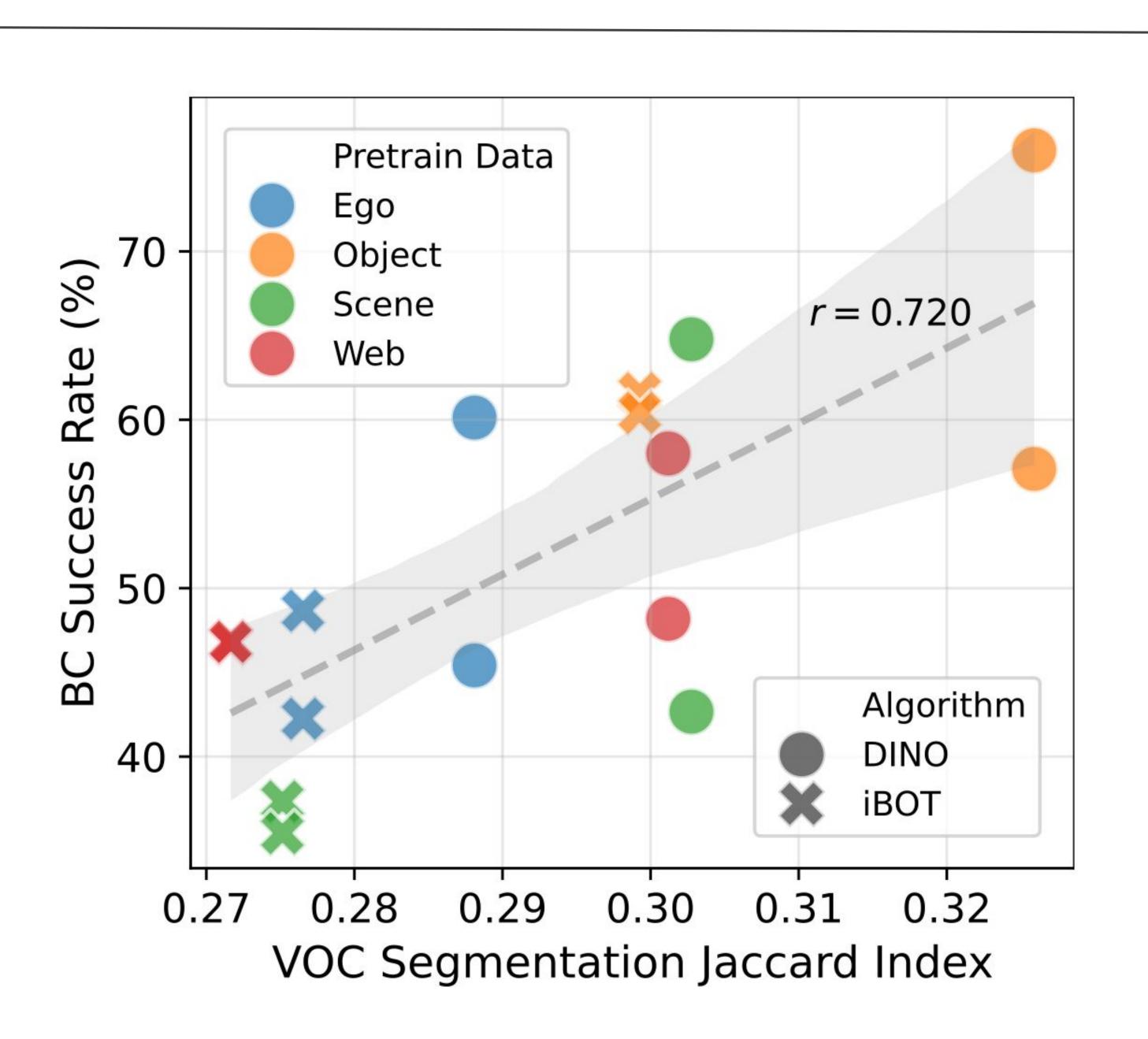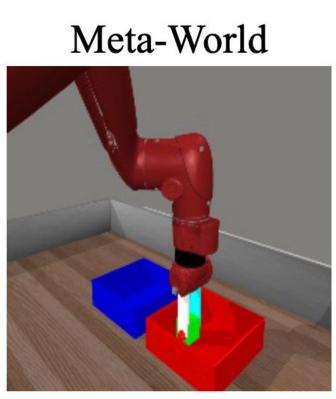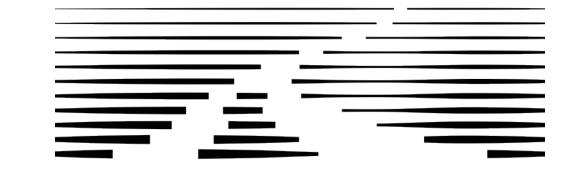**Right:** corr. (objectness v.s. performance)

BC Success Rate (%) vs VOC Segmentation Jaccard Index, $r = 0.720$
Pretrain Data: Ego, Object, Scene, Web
Algorithm: DINO, iBOT

## 3. Object-Centric Learning on Non-Object-Centric Data

### 1) Intuition: modeling objectness explicitly

Backbone: Tokens → Attention
Objectness of DINO emerges internally, which relies on (Single-)OC data bias.

Prototypes (can group tokens to objects): Tokens → Attention
Our method (SlotMIM) receives explicit objectness supervision externally.

### 2) Observation: iBOT (DINO on patches) **discovers objects via prototypes even when trained on NOC data, despite semantically misaligned.**



(a) **Clustering assignment of patch tokens.** Each patch is assigned to its nearest-neighbor prototype, with different colors indicating different prototypes.

(b) **Top-5 segments retrieved by the prototypes (by column).** A segment consists of patches assigned to the same prototype within an image. Each column shows the top-5 segments with the highest cosine similarity to the prototype corresponding to the column.

### 3) Method: DINO loss on image patches within and cross views, and contrastive learning between (grouped) object-level features.



## 4. Experiments: SoTA on Manip., Nav., Det., and Seg.

Prev. SoTA vs SlotMIM:
- Franka Kitchen: DINO 76.0, Ours 86.0
- Meta-World: DINO 80.0, Ours 84.2
- ImageNav: VC-1 67.9, Ours 69.8
- ObjectNav: VC-1 55.4, Ours 62.0
- COCO Det: iBOT 51.2, Ours 52.5
- ADE20K Seg: iBOT 50.3, Ours 51.4

| Benchmark Suite | RGB | Proprio. | Physics | Action | Goal | Learning |
|---|---|---|---|---|---|---|
| Franka Kitchen [26] | ✓ | ✓ | ✗ | Continu. | – | IL |
| Meta-World [76] | ✓ | ✓ | ✗ | Continu. | – | IL |
| ObjectNav [6] | ✓ | ✗ | ✓ | Discrete | Class | IL |
| ImageNav [81] | ✓ | ✗ | ✓ | Discrete | Image | RL |

Tasks: control and perception
Prev. SoTAs: DINO, iBOT, VC-1
Detailed comparisons available in the paper.

### Ablation Study

| | mask | $\mathcal{L}_{patch}^{cross}$ | $\mathcal{L}_{patch}^{within}$ | $\mathcal{L}_{slot}$ | $k$-NN | ADE | Jacc | $\overline{K}$ |
|---|---|---|---|---|---|---|---|---|---|
| 1 | ✗ | ✓ | ✗ | ✗ | 45.1 | 47.4 | 42.5 | 8.3 |
| 2 | ✓ | ✓ | ✗ | ✗ | 44.9 | 48.6 | 42.3 | 10.3 |
| 3 | ✓ | ✗ | ✓ | ✗ | 27.7 | 45.7 | 39.3 | 20.7 |
| 4 | ✓ | ✗ | ✓ | ✓ | 45.3 | 47.5 | 42.9 | 8.4 |
| 5 | ✓ | ✓ | ✓ | ✓ | **46.2** | **49.1** | **43.9** | 9.4 |

### Qualitative Results



Object, Ego, Scene, Web

SlotMIM learns objectiveness adaptively.

### Scaling Curves (Ctrl. & Percept.) highlighting best (model, data)



Franka Kitchen (Control), Meta-World (Control), ADE20K Sem. Seg. (Perception)
Method: MAE, DINO, iBOT, SlotMIM, V-Cond
Data: Ego, Object, Scene

### Scaling Curves (ImageNet L.P.)



OC Data (Linear Prob.), NOC Data (Linear Prob.)

**Better scaling on NOC data** means 1) less dependence on data curation, 2) better scalability, and 3) better data efficiency.

Future investigations on scalability could be valuable.