# NAVER LABS

Europe

## ICLR
Vienna, 2024

hrough

enon

epfel    Christian Wolf

Guilla

(a) Siamese    (b) Siamese + m

**Goal**

**Observation**

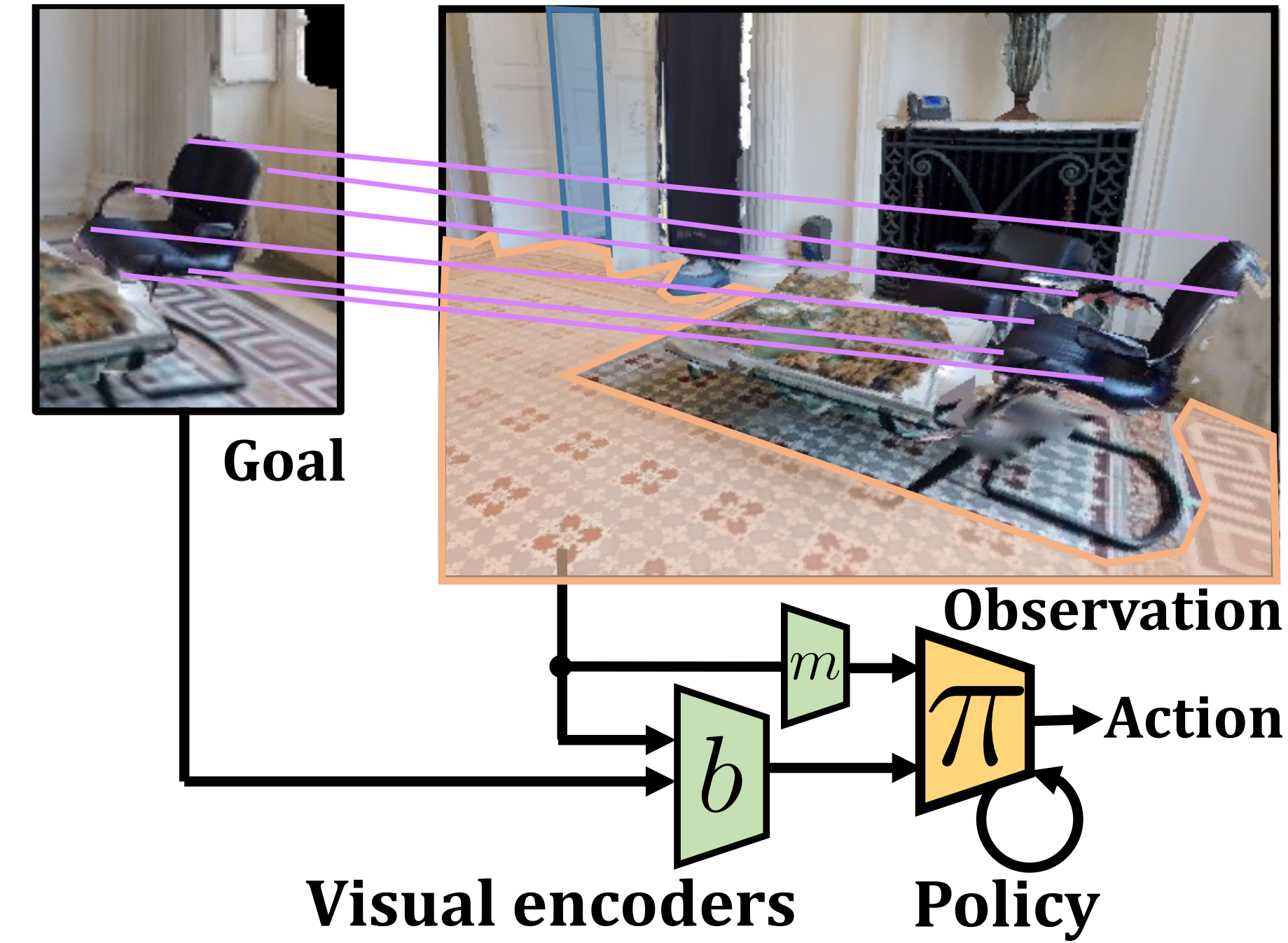**Action**

$m$

$b$

$\pi$

**Visual encoders**    **Policy**

## Motiva

**Task:** ImageNav,: navigation to a goal specified by an

**Requires:** (1) Nav skills: detection of navigable space, exi
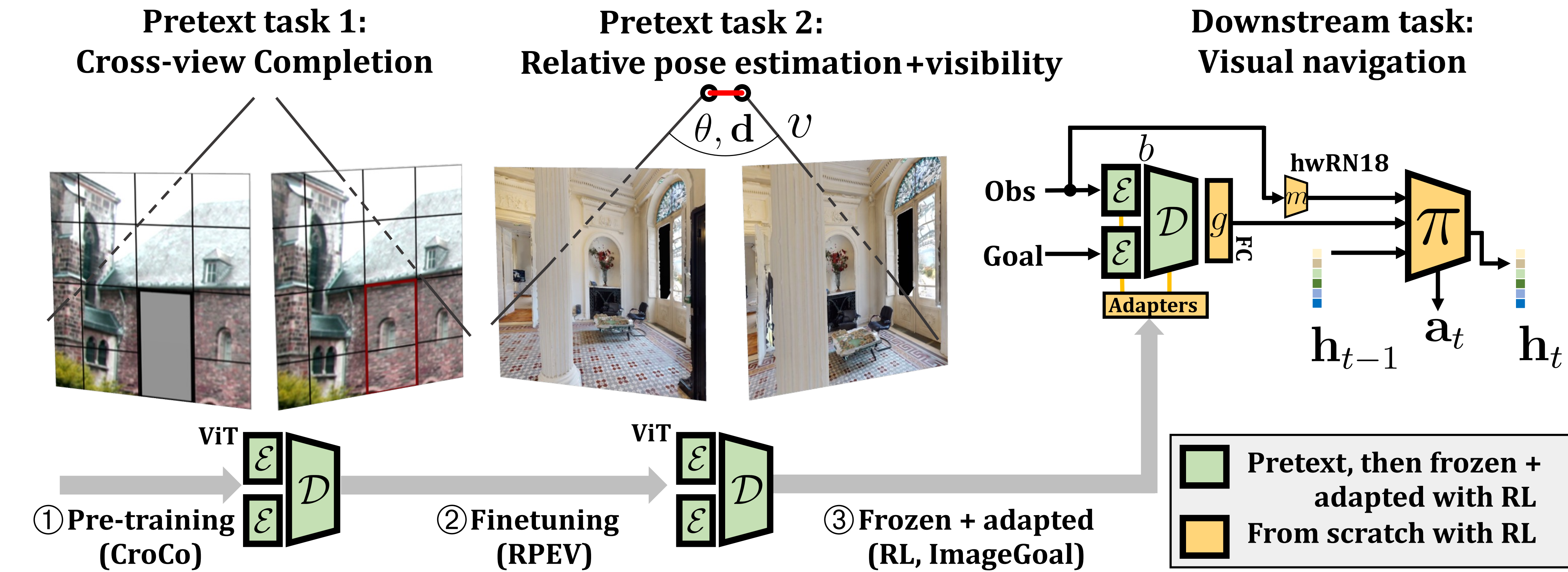
(2) Detection of relative pose wrt the goal

**Cur SOTA:** Train visual encoder + policy end-to-end, or out
goal detection to local feature based methods.

**Idea:** Pre-train a <u>binocular</u> ViT with cross-attention o
different losses:

- Cross-view competition as in Weinzaepfel e
NeurIPS 2022

- Extremely-wide baseline relative pose + visi

Use as visual encoder in an end-to-end policy tr
with RL. Encoder is frozen + adaptors.

### RPEV Results during Nav



### nd sizes

| | correct poses | | Vis-acc | Nav. perf. | |
|---|---|---|---|---|---|
| | 10° | 2m&20° | (%) | SR (%) | SPL (%) |
| | 7.5 | 98.9 | **94.0** | 82.0 | **59.6** |
| | 2.5 | 96.8 | 89.3 | **83.0** | 55.6 |
| | 2.7 | 93.5 | 81.6 | 79.6 | 52.1 |

## Training

**Pretext task 1:**
**Cross-view Completion**

**Pretext task 2:**
**Relative pose estimation+visibility**

$\theta, \mathbf{d}$  $v$

**Downstream**
**Visual naviga**



**Obs** — $\mathcal{E}$

**Goal** — $\mathcal{E}$ — $\mathcal{D}$ — $g$ — $b$ — hwRN18

$m$

**Adapters**

$\mathbf{h}_{t-1}$

**ViT** $\mathcal{E}$ $\mathcal{D}$ $\mathcal{E}$

①**Pre-training**
**(CroCo)**

**ViT** $\mathcal{E}$ $\mathcal{D}$ $\mathcal{E}$

②**Finetuning**
**(RPEV)**

③**Frozen + adapted**
**(RL, ImageGoal)**

Pretext, th
adapt

From scratch with RL



TN

$*<\tau, v>\tau)$, FN $(v^*>\tau, v<\tau)$

### Alignment architecture / loss



(a) Siamese    (b) Siamese + m

(c) CNN+TinyT    (d) DEBiT (Ours)

## Correspondence solutions emerge after pre-training (CroCo + RPEV)



Shown: cross-attention
of the last head average
over all heads. Patches
selected based on total
attention.

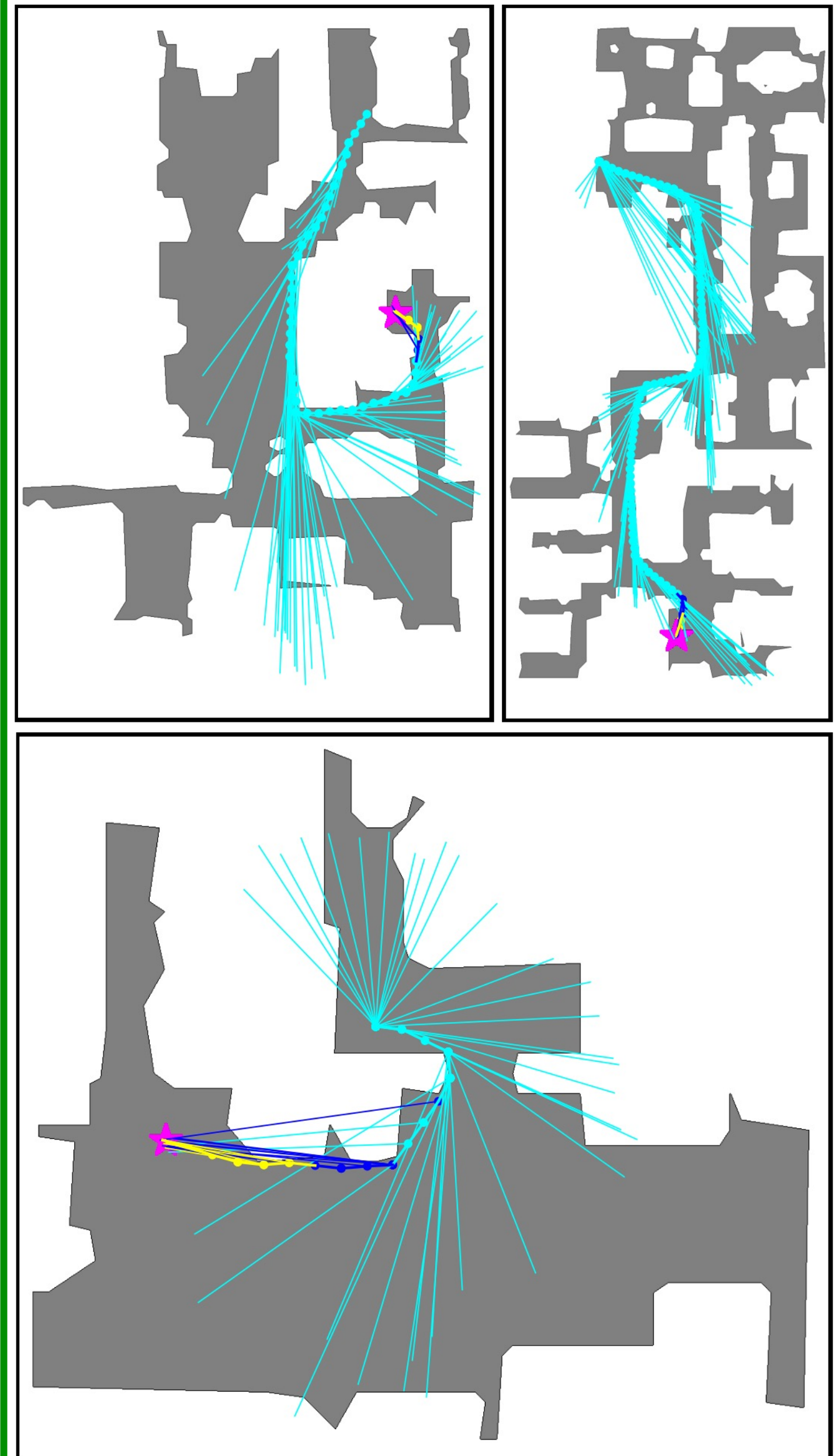| | #parms | SR | SPL |
|---|---|---|---|
| Ours (DEBiT-L), no adapters | 200M | 82.0 | 59.6 |
| Ours (DEBiT-L) + adapters | 200M | **94.0** | **71.7** | Frozen + adapted |

[1] *Perf. from orig. papers;* [2] *Mono-view ablation of baseline in Table III of (Mezghani et al., 2022);*
[3] *Retrained in mono-view settings, see Table 1 of (Al-Halah et al., 2022)*
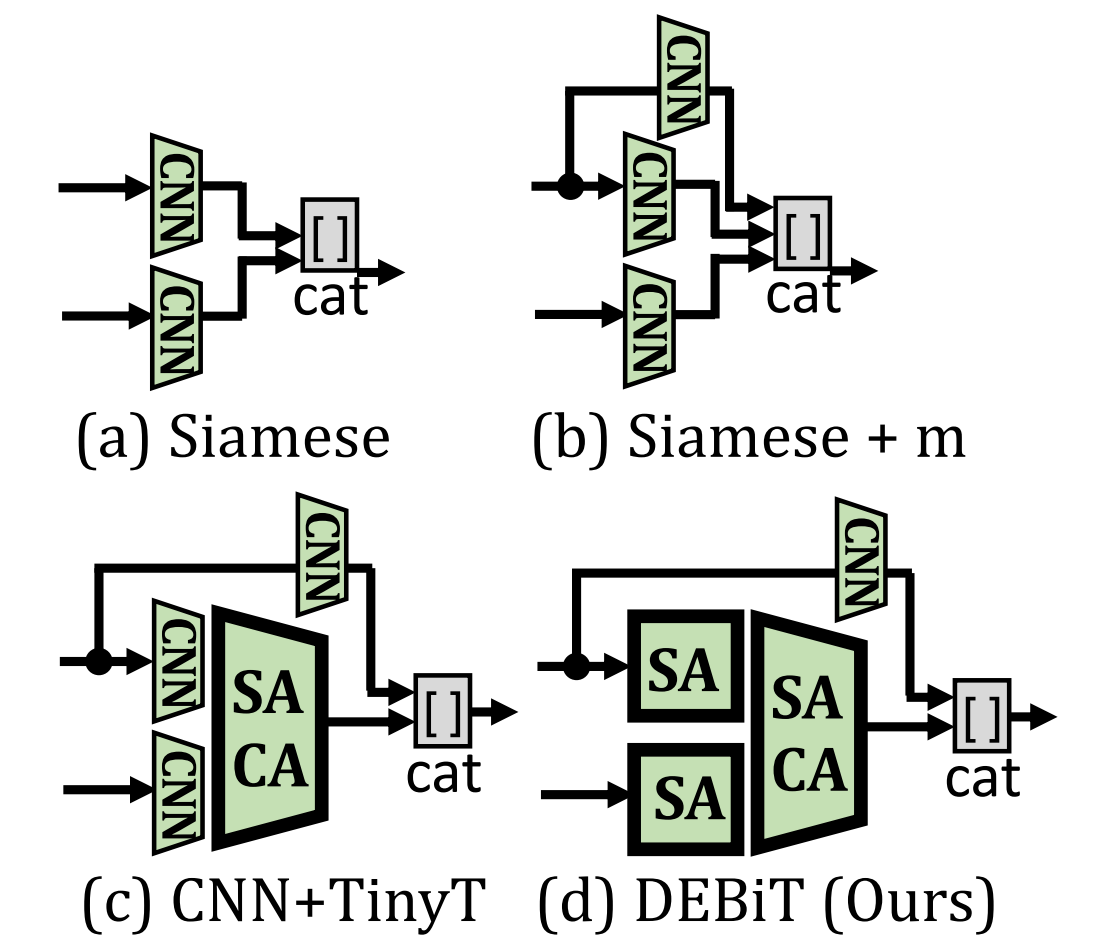
## Comparison w. SOTA: Instance Image

| Method | #steps | — SR (%) — | | — SPL (%) — | |
|---|---|---|---|---|---|
| | | max | avg | max | avg |
| (Krantz et al., 2022) | 3500M | 5.5 | n/a | 2.3 | n/a |
| (Krantz et al., 2023) | n/a | 56.1 | n/a | 23.3 | n/a |
| Ours(DEBiT-L)+adapters | 200M | **61.1** | 59.3 | 33.5 | 32.4 |

**Task:** camera
may differ bt

**Method:** pre
finetune poli
Instance-Ima



(a) Siamese    (b) Siamese + m

(c) CNN+TinyT    (d) DEBiT (Ours)

| Visual encoder | Pre-train | #parms | SR | SPL |
|---|---|---|---|---|
| (a) Siamese hwRN18* | No | 4.1M | 10.1 | 9.6 |
| (b) Siamese hwRN18*+$m$ | RPEV | 8.3M | 8.0 | 7.7 |
| (c) hwRN18+Cross-Att+$m$ | No | 10M | 7.4 | 4.7 |
| (c) hwRN18+Cross-Att+$m$ | RPEV | 10M | 7.4 | 7.2 |
| (d) DEBiT-B (Ours), no adapters | No | 60M | 6.8 | 4.0 |
| (d) DEBiT-B (Ours), no adapters | CroCo+RPEV | 60M | **83.0** | 55.6 |

* *Baseline in (Mezghani et al., 2022), inspired by (Zhu et al., 2017)*