

Modality-invariant Visual Odometry for Indoor Navigation

Marius Memmel, Amir Zamir



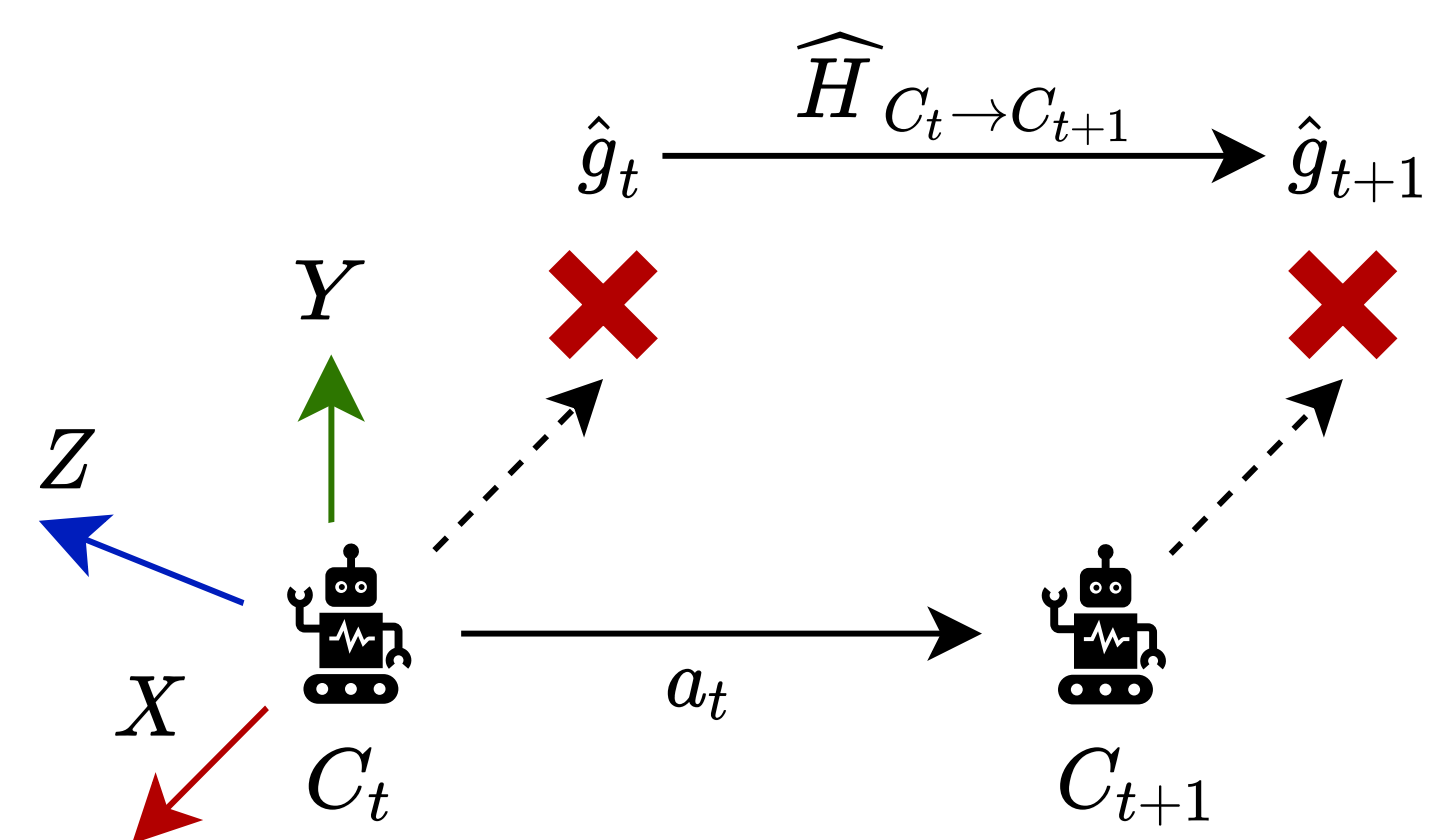
Code, paper, visuals

Problem

- *GPS+compass* are **not reliable for localization** in indoor environments [1]
- Visual Odometry (VO) is a sufficient alternative, but models can be **dominated by a single modality** in a multi-modal (RGBD) setting
- **Erroneous sensors** causes **catastrophic failure** of VO models

Visual Odometry for Indoor Navigation [1,2]

- Estimate transformation between agent's coordinate systems



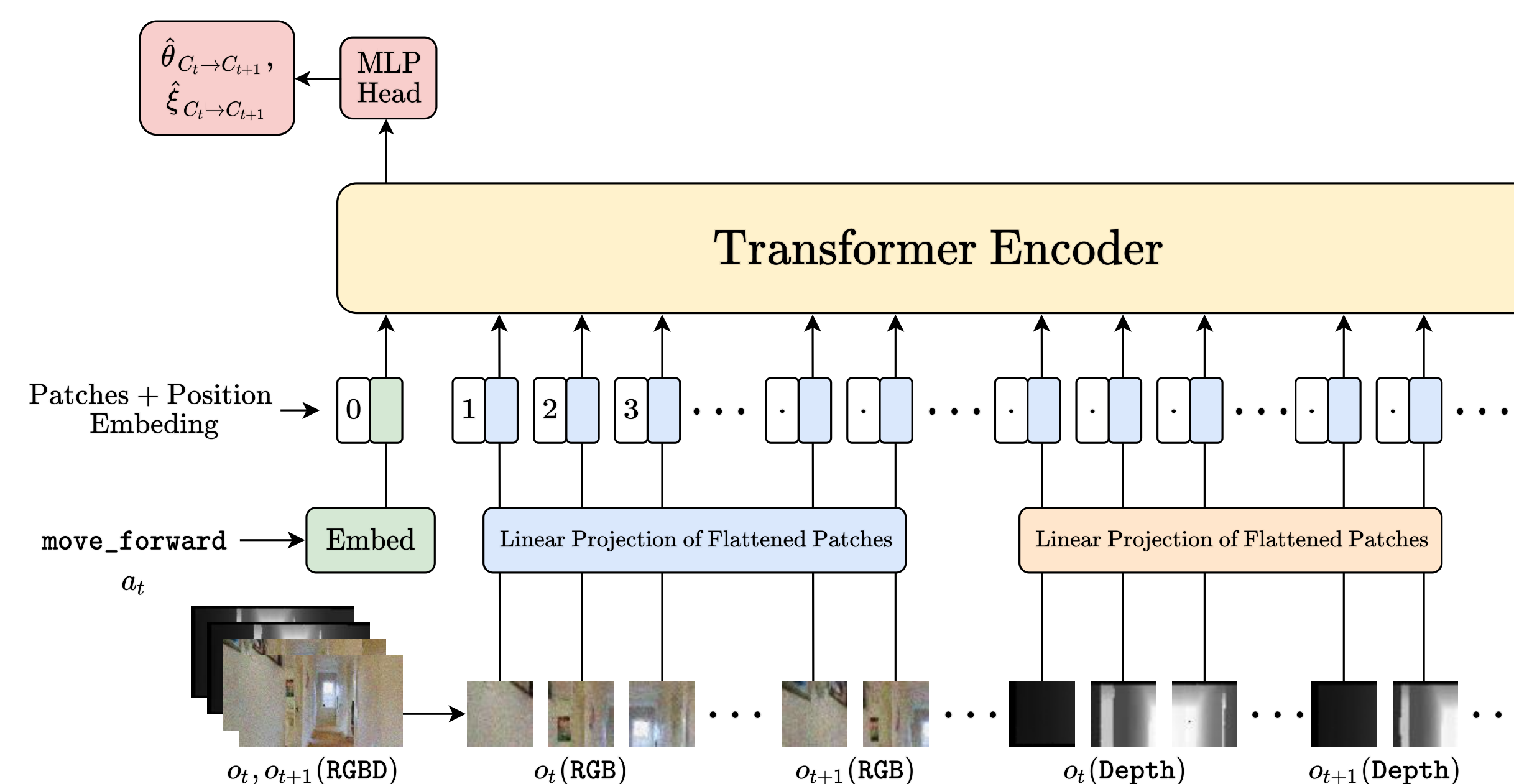
$$\hat{H}_{C_t \rightarrow C_{t+1}} = \begin{bmatrix} \hat{R}_{C_t \rightarrow C_{t+1}} & \hat{\xi}_{C_t \rightarrow C_{t+1}} \\ 0 & 1 \end{bmatrix}$$

$$\text{with } \hat{R}_{C_t \rightarrow C_{t+1}} = \begin{bmatrix} \cos(\hat{\beta}_{C_t \rightarrow C_{t+1}}) & -\sin(\hat{\beta}_{C_t \rightarrow C_{t+1}}) \\ \sin(\hat{\beta}_{C_t \rightarrow C_{t+1}}) & \cos(\hat{\beta}_{C_t \rightarrow C_{t+1}}) \end{bmatrix} \in SO(2)$$

$$\langle \hat{\beta}_{C_t \rightarrow C_{t+1}}, \hat{\xi}_{C_t \rightarrow C_{t+1}} \rangle = f_\phi(\mathbf{o}_t, \mathbf{o}_{t+1})$$

- Predict transformation parameters from observations $\mathbf{o}_t, \mathbf{o}_{t+1}$ using f_ϕ
- When f_ϕ is ConvNet, input size is fixed
- ▶ Dropping any modality leads to catastrophic failure!

Solution: Visual Odometry Vision Transformers (VOT)



- Vision Transformers are **flexible w.r.t. input modalities** even during test time due to the independence of their weights on the sequence length [3]
- Multi-modal **pre-training** (e.g. MultiMAE [4])
- Attention mechanism to **identify importance** of regions and modalities
- Pass **action prior** as separate token

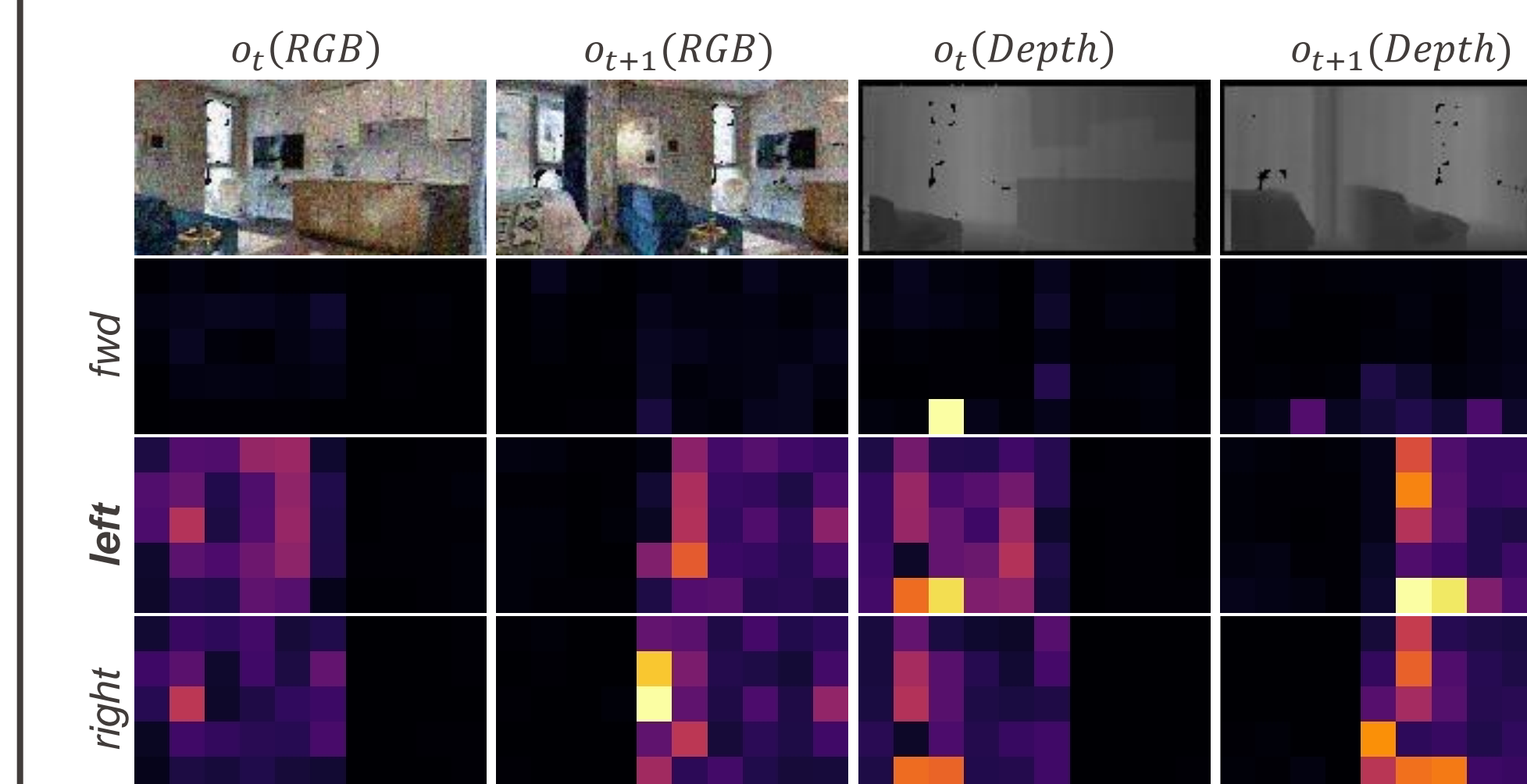
PointNav Results

- dropping modalities (drp) lets ConvNet approaches converge to a *blind* VO model
- VOT maintains sufficient localization capabilities, c.f., *SoftSPL* (SSPL) [1]

method	obs	drp	$S \uparrow$	SPL \uparrow	SSPL \uparrow
<i>blind</i>	-	-	0.00	0.00	5.40
<i>oracle</i>	-	-	97.89	74.80	73.10
[17]	RGB-D	-	64.50	48.90	65.40
[17]	RGB-D	RGB	0.00	0.00	5.40
[17]	RGB-D	Depth	0.00	0.00	5.40
VOT (ours)	RGB	-	59.30	45.40	66.70
VOT (ours)	Depth	-	93.30	71.70	72.00
VOT (ours)	RGB-D	-	88.20	67.90	71.30
VOT (ours)	RGB-D	RGB	75.90	58.50	69.90
VOT (ours)	RGB-D	Depth	26.10	20.00	58.70

Attention Maps

(VOT trained on RGBD and pre-trained with a MultiMAE)



- Action taken: *left*, Injected actions: *fwd*, *left*, *right*
- Embed *fwd* causes the attention to focus on the center
- Embed *left* and *right* move attention towards regions of the image that would be consistent across time steps $t, t+1$ in case of rotation

Limitations & Future Work

- ▶ Full robustness may require explicit measures (randomly drop modalities during training)
- ▶ Dropping modalities causes inconsistencies in the localization (fine-tune the policy to adapt)
- ▶ Only two modalities (extend by semantic segmentation)

Takeaways

- ▶ VOT can function when modalities are missing
- ▶ VOT implicitly learns some invariance to modalities
- ▶ Action prior primes model to attend relevant regions
- ▶ Pre-training reduces data requirements to 25%

PointNav Paths

- ▶ VOT helps agents reach the goal 1) - 3)
- when dropping modalities, agent still makes progress towards the goal 4) 5)
- however, those agents get stuck in narrow passages
- ▶ don't reach the goal in T time steps

