

IR-VIC: Unsupervised Discovery of Sub-goals for Transfer in RL

Nirbhay Modhe¹



Prithvijit Chattopadhyay¹



Mohit Sharma¹



Abhishek Das¹



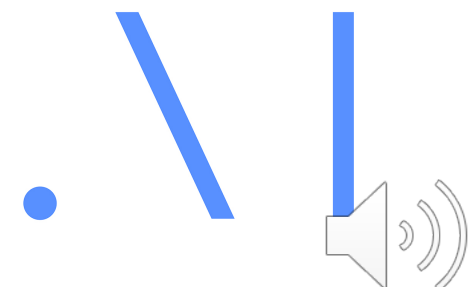
Devi Parikh^{1,2}



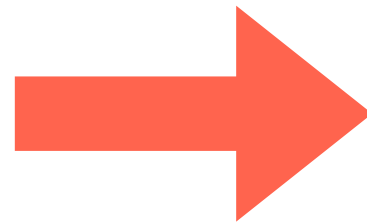
Dhruv Batra^{1,2}



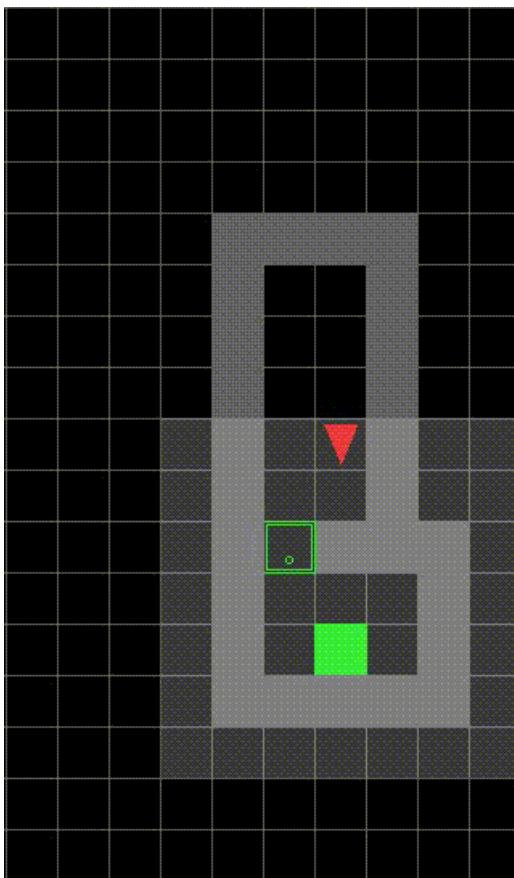
Ramakrishna Vedantam²



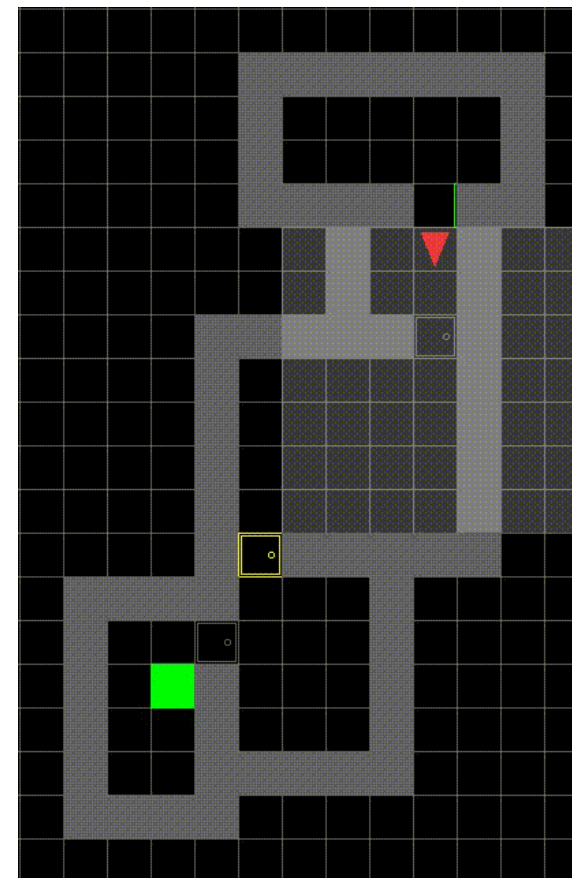
Exploration with Sparse Rewards



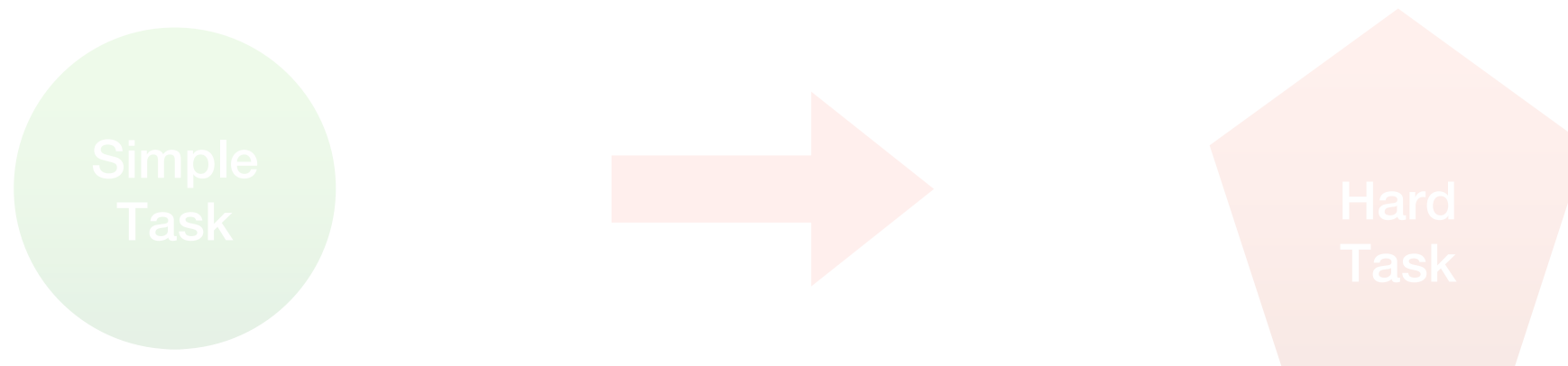
Transfer knowledge?



Sparse reward, +1
for reaching goal



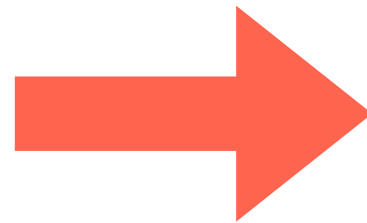
Transfer of Sub-goals



“...decomposing a complex problem into a set of simpler ones”
- McGovern & Barto, 2001^[1]



Transfer of Sub-goals

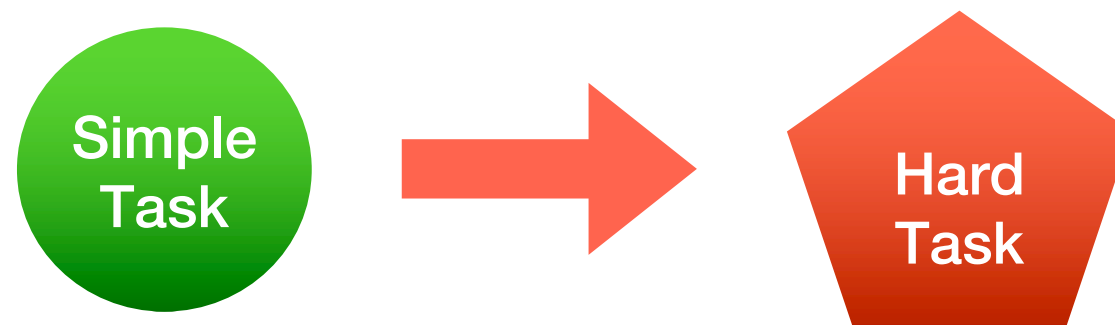


Learn sub-goal detector in
source environment

Transfer sub-goal detector
to target environment



Challenges of Transfer



Challenges for sub-goal transfer:



Detector should be **easy to train** in source environments



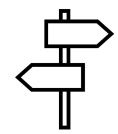
Transferable to **novel target environments**



Improve performance in target environment's task

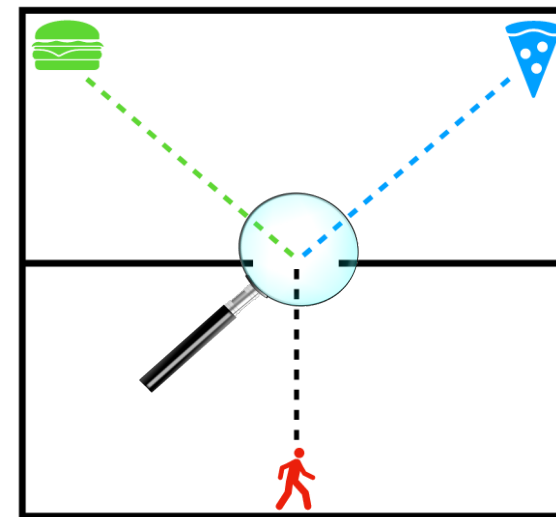


Information-based Sub-goals



Peaks in **Relevant Goal Information (RGI)**

Goal-informed actions

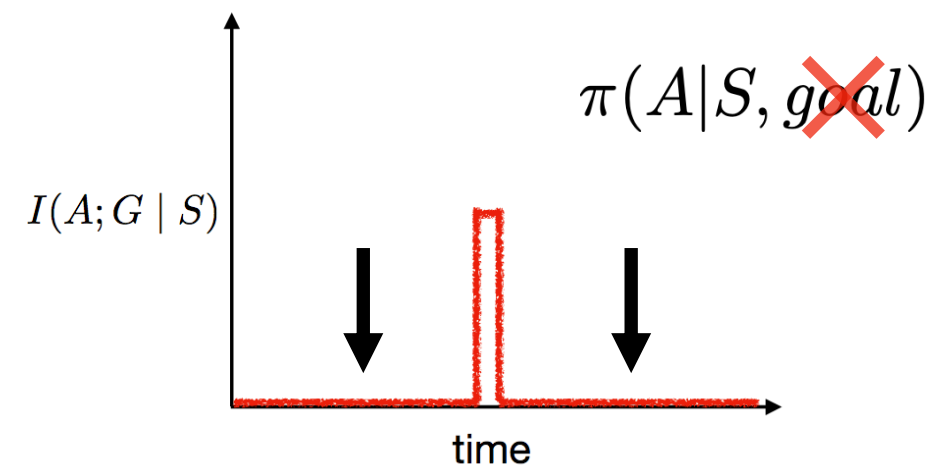


$$G \sim (\text{burger}, \text{pizza})$$



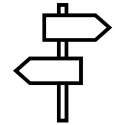
Default Behavior (**Low RGI**)

Goal-independent actions

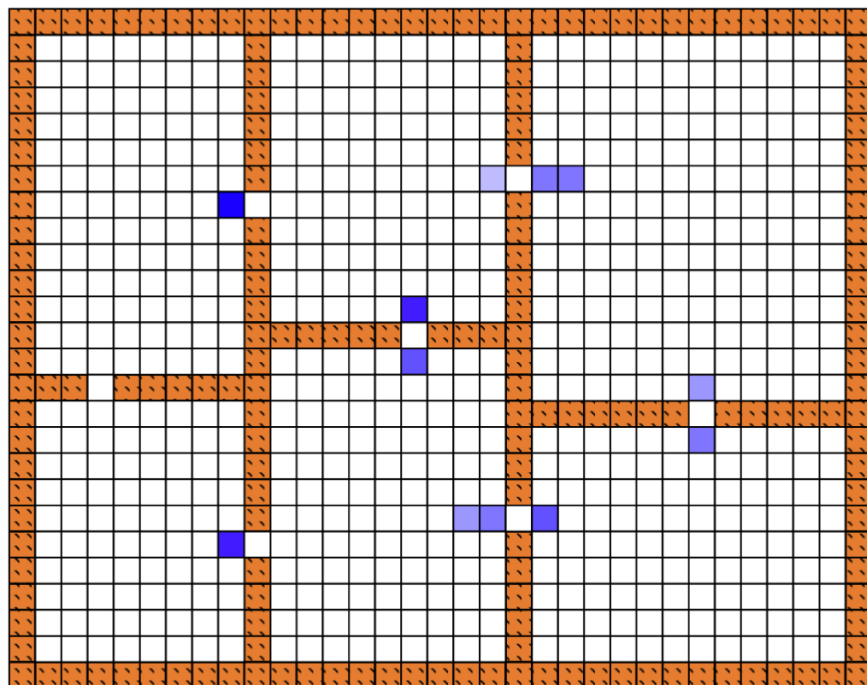




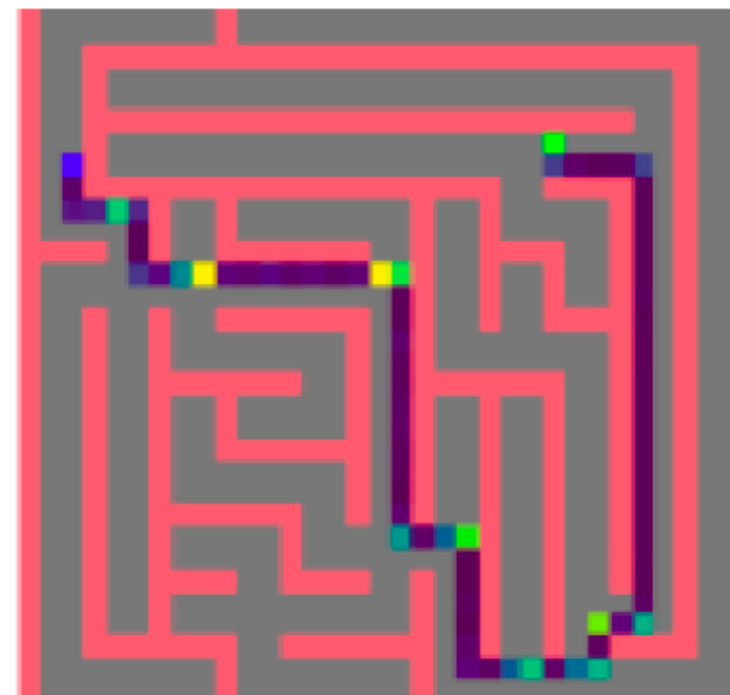
Sub-goal Identification



Sub-goals discovered with goal-driven (extrinsic) rewards



Dijk & Polani "Grounding Subgoals in Information Transitions." 2011

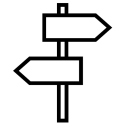


Goyal et. al. "InfoBot: Transfer and Exploration via the Information Bottleneck." 2019



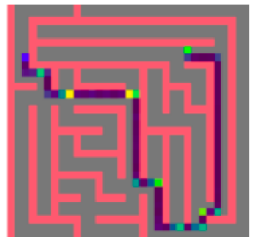
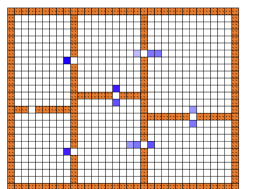
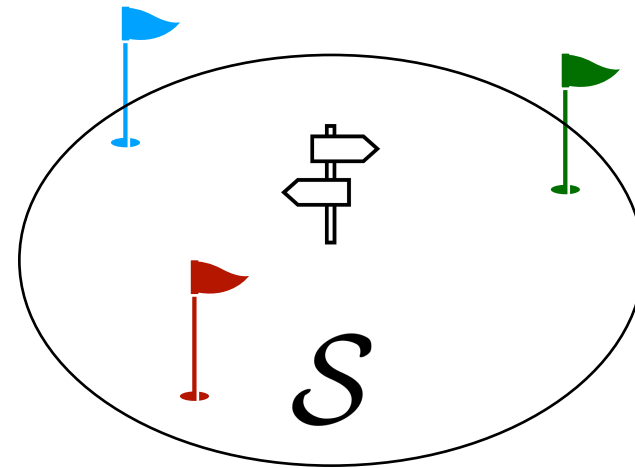


Sub-goal Identification



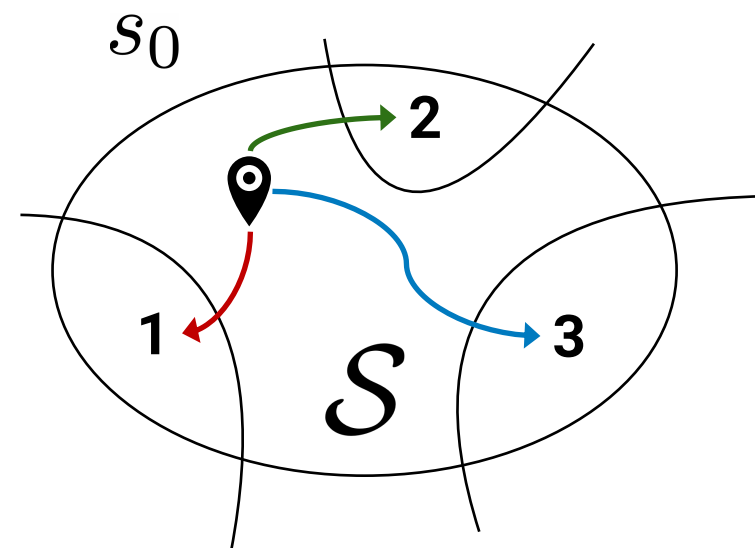
Sub-goals discovered with **goal-driven (extrinsic)** rewards

$$\text{signpost} = f(\text{task}, \text{environment})$$



Sub-goals discovered with **intrinsic** rewards

$$\text{signpost} = f(\text{environment})$$



self-supervised goals





Sub-goal Identification

$$\text{Icon} = f(\text{task}, \text{environment})$$

- Require **task specification** in source environment
- May not generalize to **different tasks** in similar environments

$$\text{Icon} = f(\text{environment})$$

- Unsupervised, **task-independent** objective
- Exploits **environment structure** alone, generalizing better to similar environments



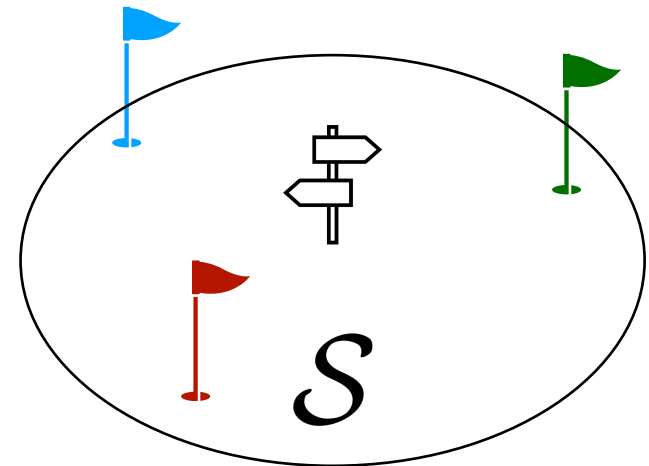


Unsupervised Sub-goal Discovery

Supervised goal:  External Reward  $I(\text{actions; goals})$

$$\max_{\pi_{\theta}} [r - \beta I(A; G)]$$

Maximize task reward with penalty for using goal information

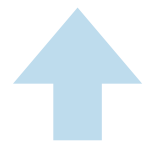




Unsupervised Sub-goal Discovery

IR-VIC: Sub-goal discovery without external tasks

Supervised goal:



External Reward



$I(\text{actions; goals})$

Unsupervised goal:



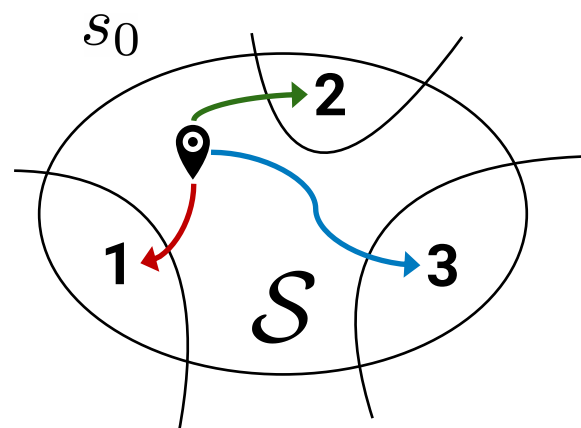
Intrinsic Control



$I(\text{actions; options})$

Learn intrinsic options

Look at option sparingly



Maximize intrinsic reward
with penalty for using option
information

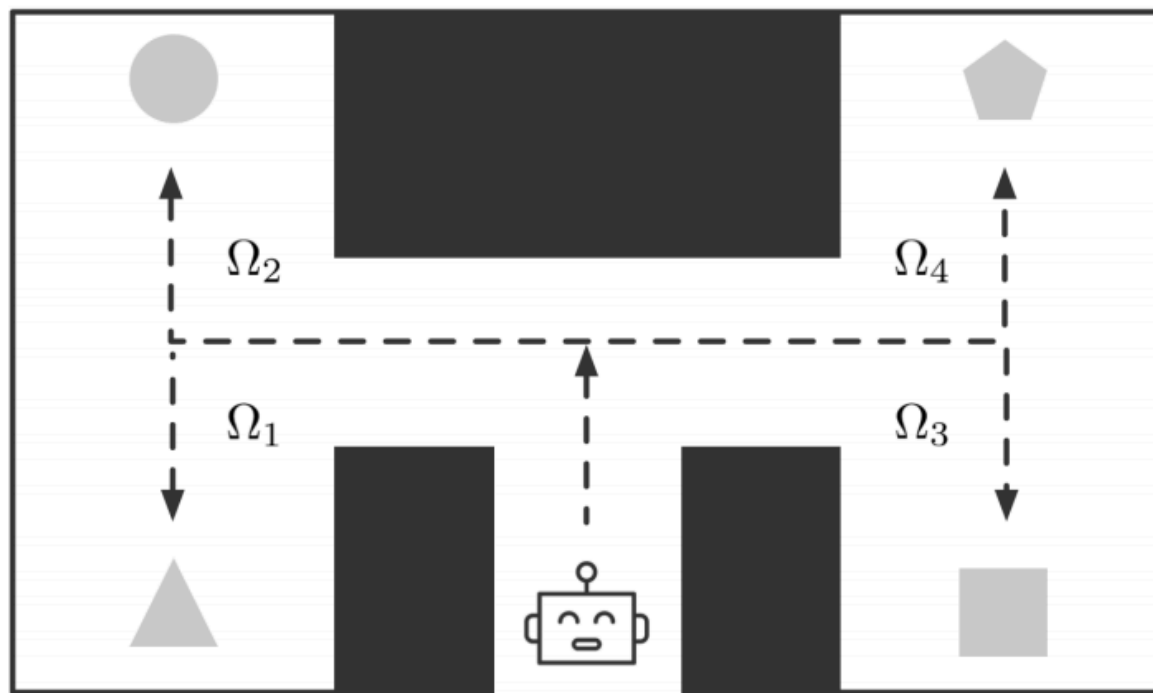


Learning Intrinsic Options

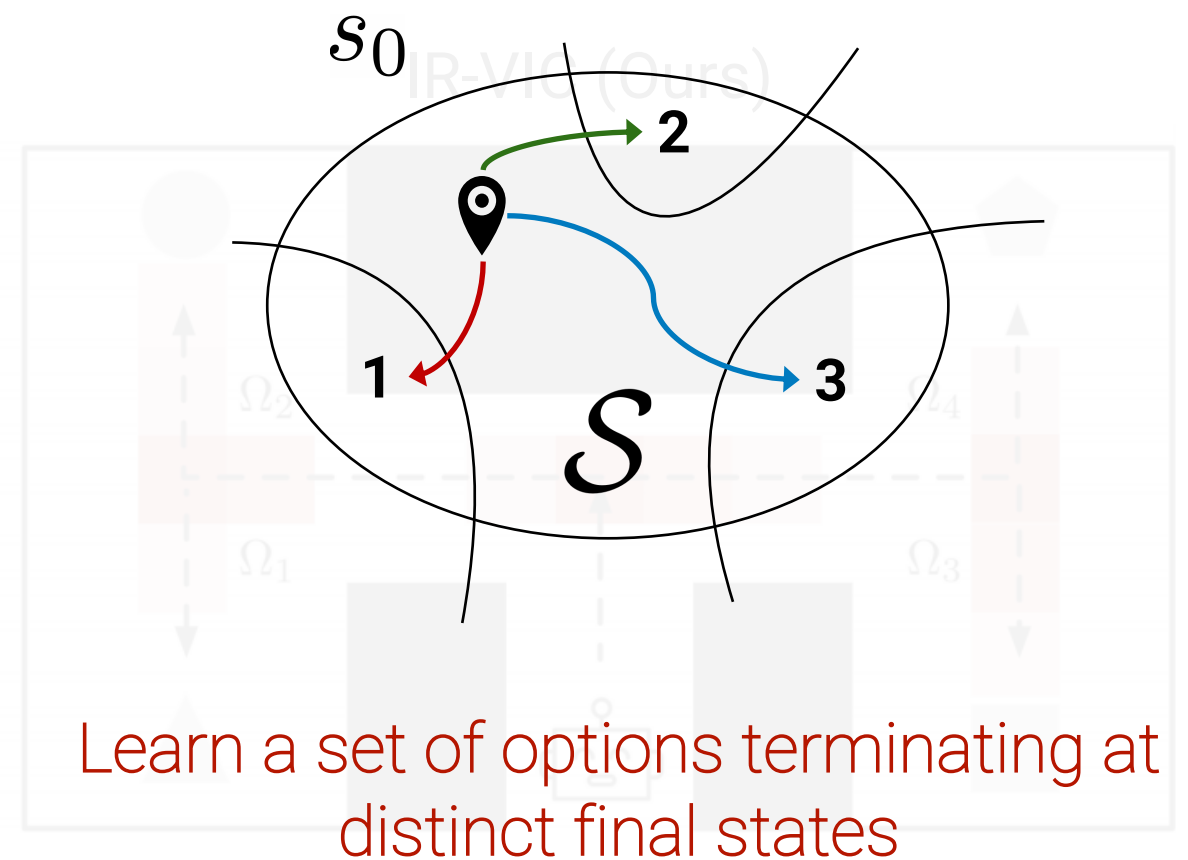
Unsupervised goal:  **Intrinsic Control**  $I(\text{actions}; \text{options})$

Learn intrinsic options

VIC: Variational Intrinsic Control^[1]
(Gregor et al., 2016)



Ω : option

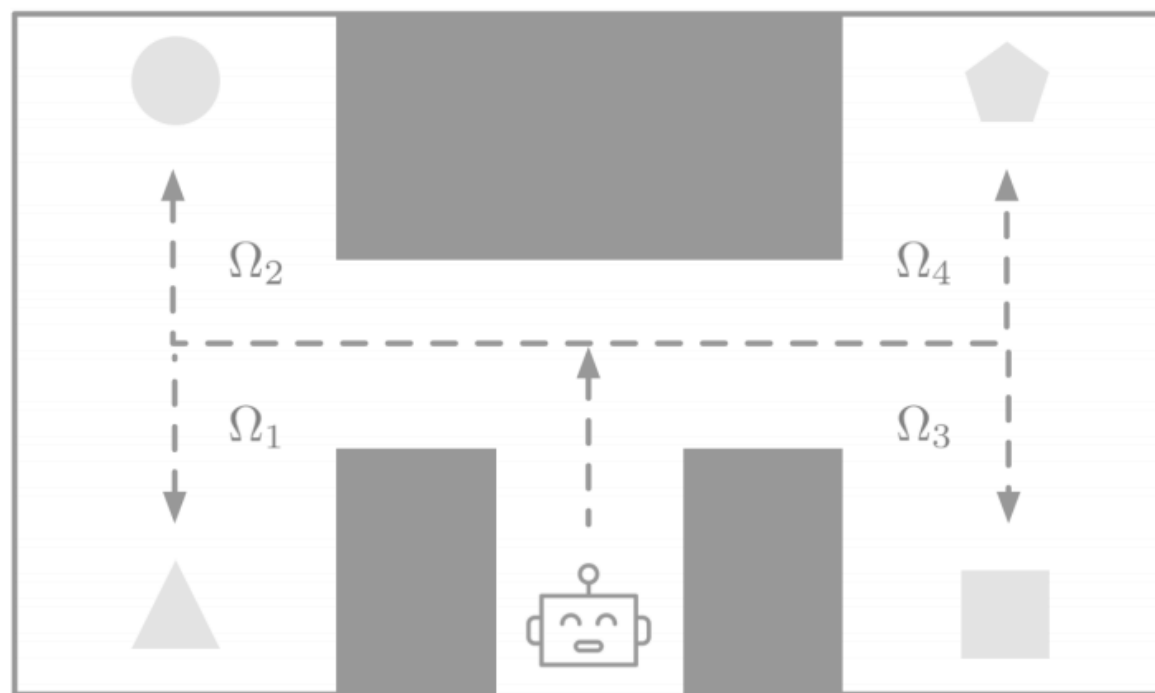


Information Minimization

Unsupervised goal:  Intrinsic Control  $I(\text{actions}; \text{options})$

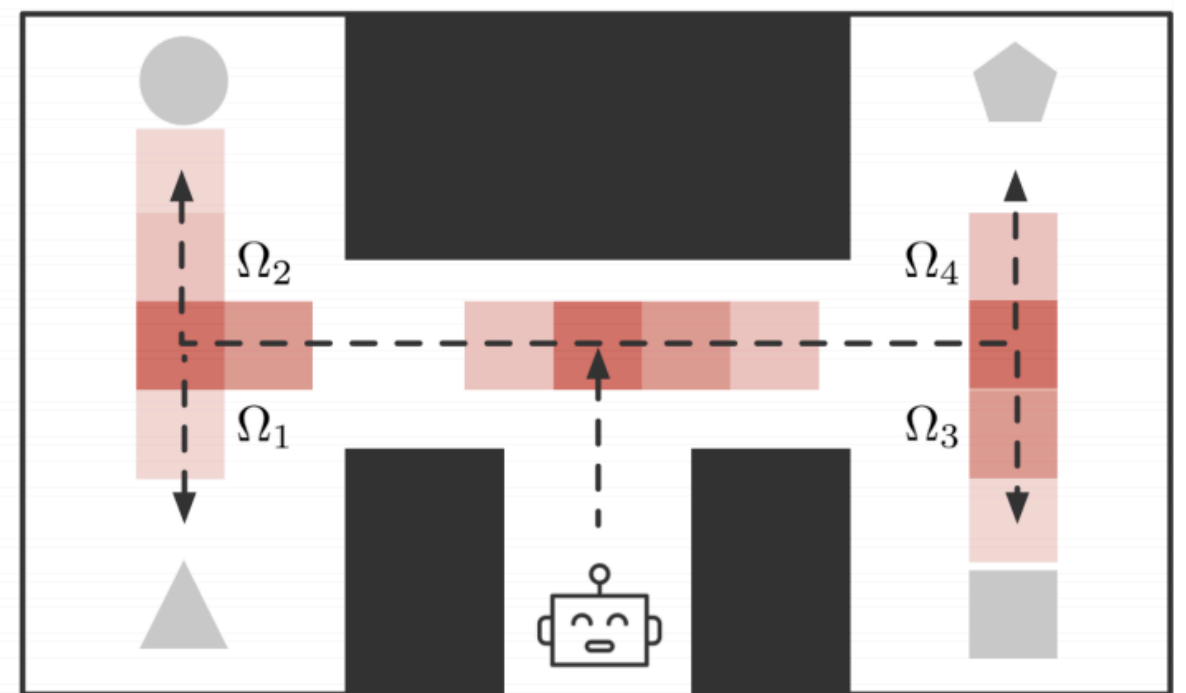
Enforce "default behavior" / option-independent actions 

VIC: Variational Intrinsic Control^[1]
(Gregor et al., 2016)



Ω : option

IR-VIC (Ours)

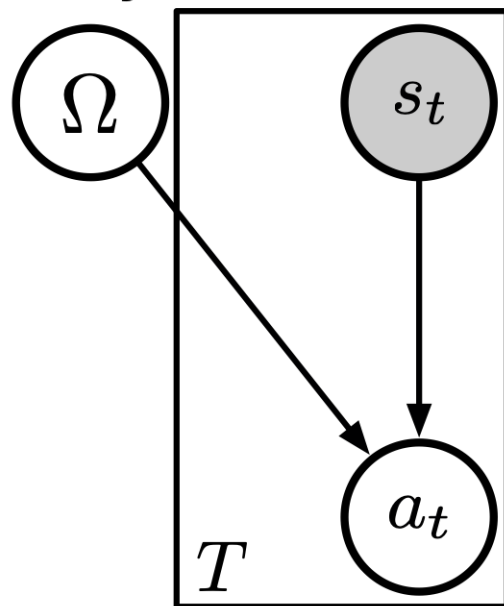


Information Minimization

Unsupervised goal:  Intrinsic Control  $I(\text{actions}; \text{options})$

Enforce "default behavior" / option-independent actions 

Policy Parameterization



policy
 $\pi(a_t | s_t, \Omega)$

$$I(a_t; \Omega | s_t) \xrightarrow{\text{pink arrow}} I(z_t; \Omega | s_t)$$

Ideal: Minimize action-goal information



Information Minimization

Unsupervised goal:



Intrinsic Control

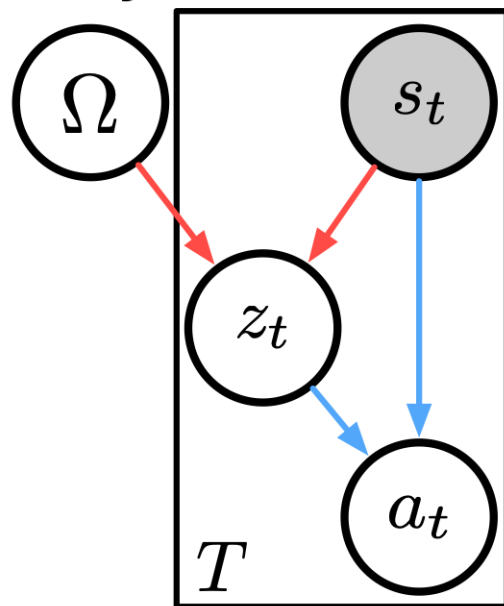


$I(\text{actions}; \text{options})$

Enforce "default behavior" / option-independent actions



Policy Parameterization



policy

$$\pi(a_t | s_t, \Omega)$$

encoder

$$p_{enc}(z_t | s_t, \Omega)$$

decoder

$$p_{dec}(a_t | s_t, z_t)$$

$$I(a_t; \Omega | s_t) \leq I(z_t; \Omega | s_t) \quad \text{pink downward arrow}$$

Practical: Minimize upper bound

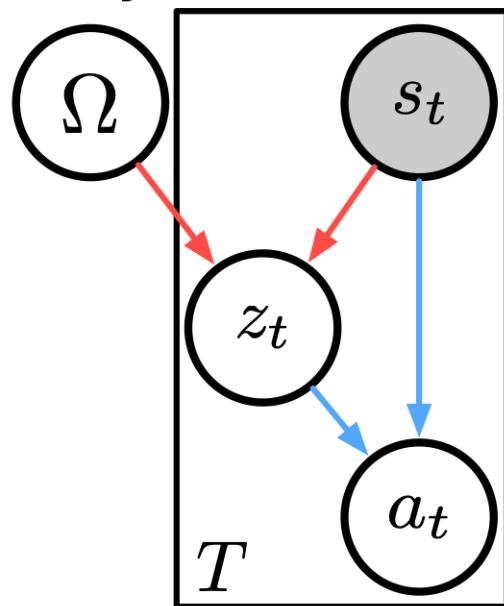


Information Minimization

Unsupervised goal:  Intrinsic Control  $I(\text{actions}; \text{options})$

Enforce "default behavior" / option-independent actions 

Policy Parameterization



policy
 $\pi(a_t | s_t, \Omega)$
 encoder
 $p_{enc}(z_t | s_t, \Omega)$
 decoder
 $p_{dec}(a_t | s_t, z_t)$

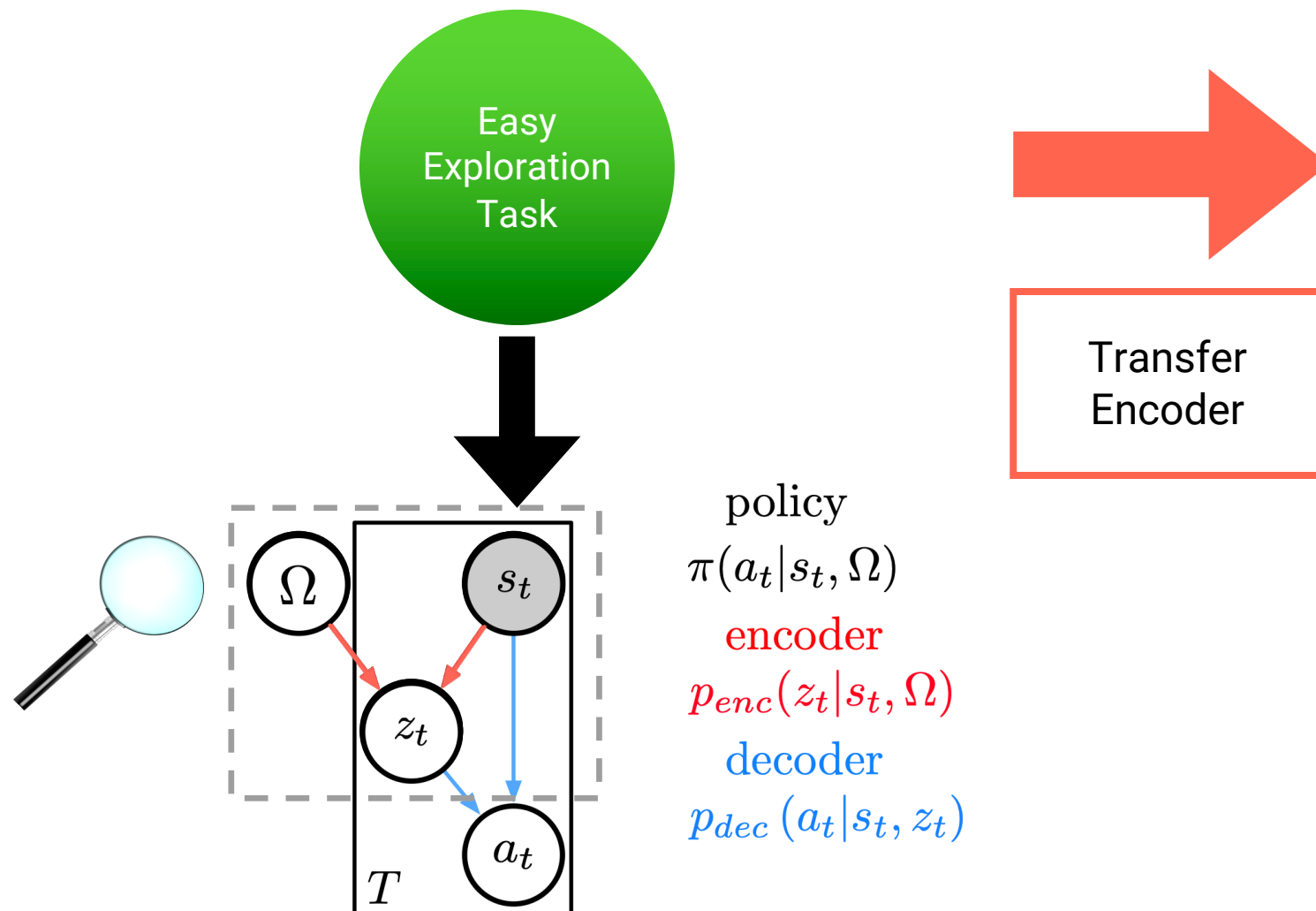
$$I(a_t; \Omega | s_t) \leq I(z_t; \Omega | s_t) \quad \text{↓}$$

Sub-goals are peaks in MI



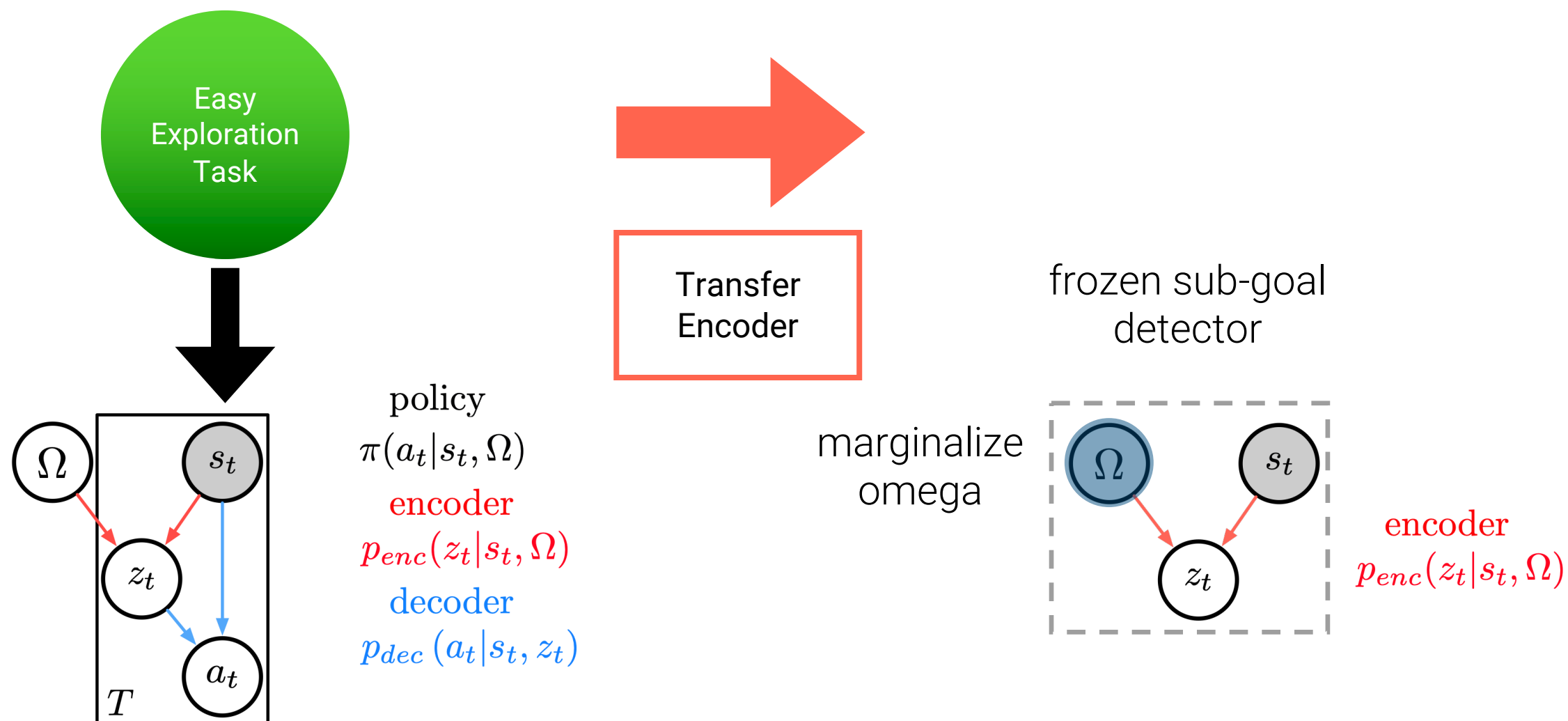
Sub-goal Transfer

Learn sub-goal detector
for simple task



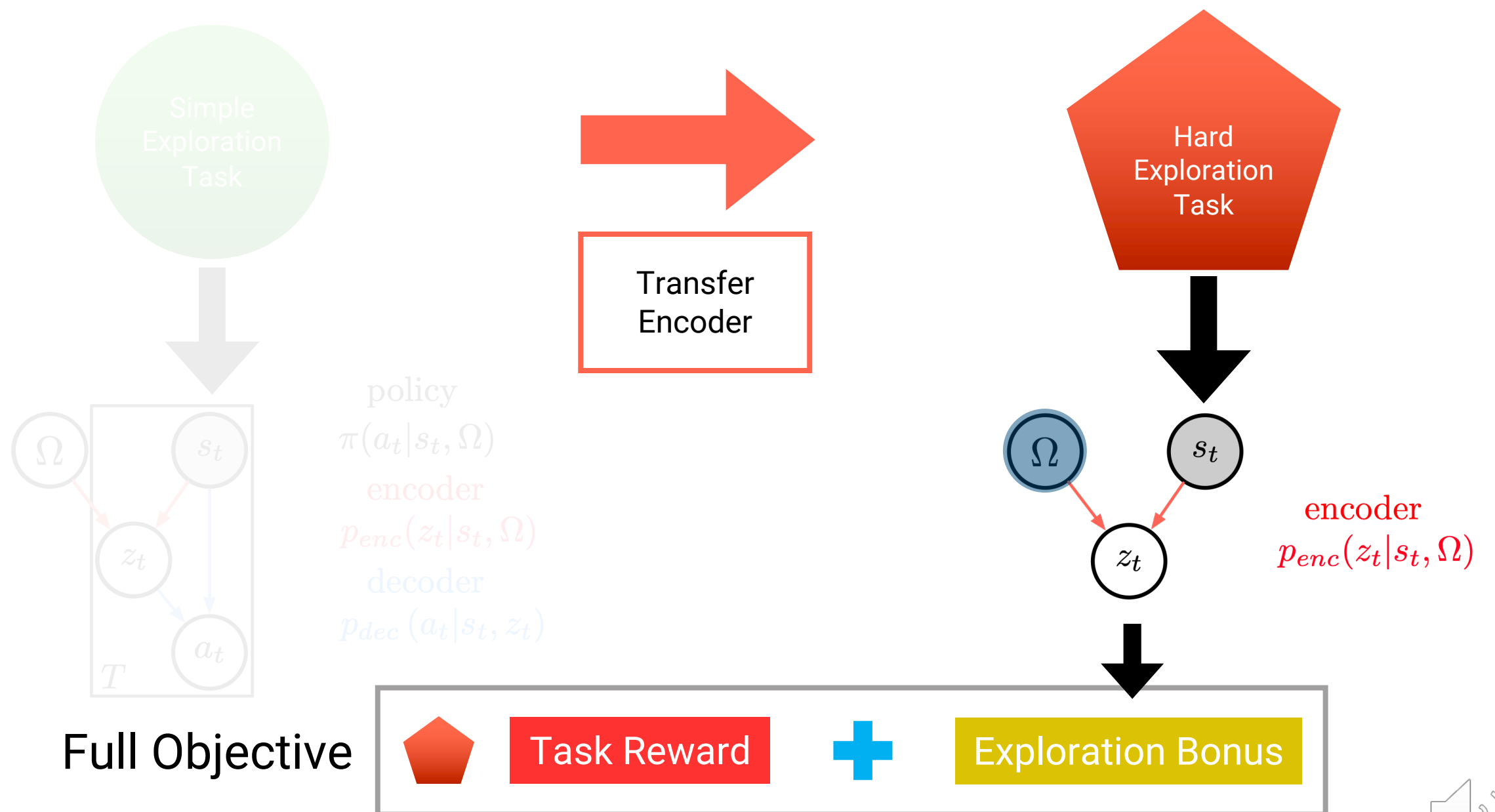
● → ◆ Sub-goal Transfer

Learn sub-goal detector
for simple task



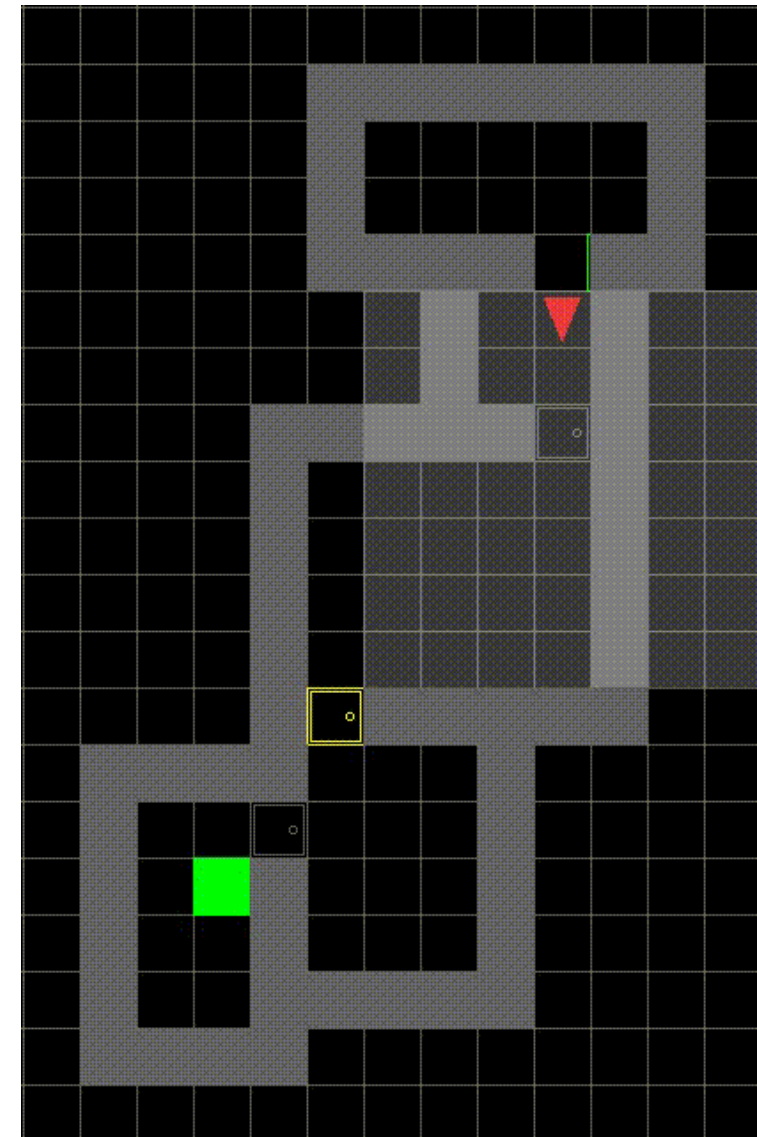
Sub-goal Transfer

Learn sub-goal detector
for simple task



Experimental Evaluation

- Set of environments: Gym-minigrid
 - $s_t \rightarrow N \times N$ Image
 - $G_t \rightarrow$ Vector to Goal
 - $\mathcal{A} = \{\text{fwd, left, right, toggle}\}$
- External Task: Point-goal navigation (green square)
- External Reward: +1 (decaying) on goal reached



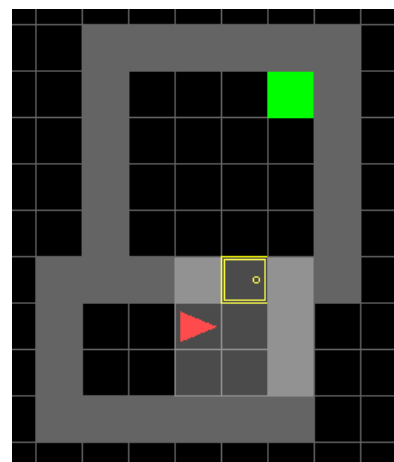
N: # of rooms
S: max room size



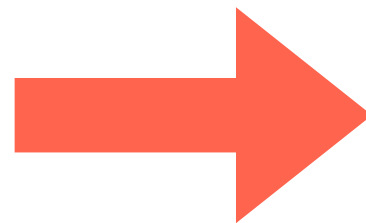
Easy Transfer

Transfer experiments from Goyal et al.^[1]

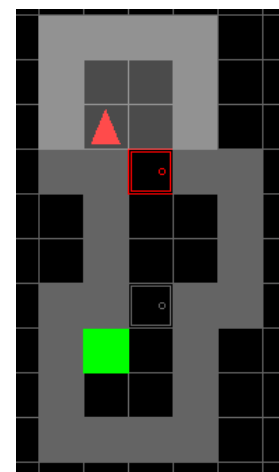
N: # of rooms
S: max room size



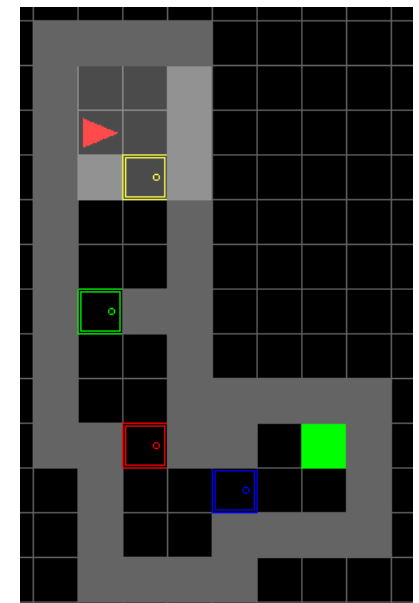
N2 S6



Increasing
of rooms



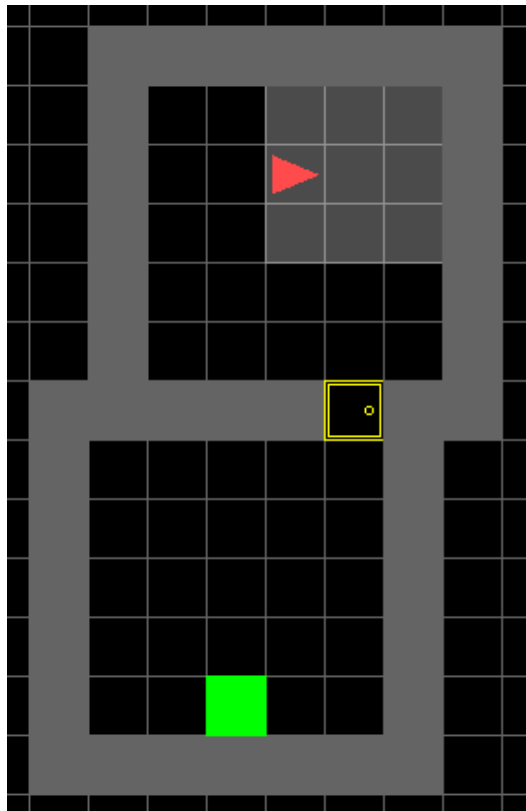
N3 S4



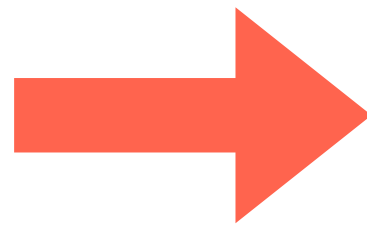
N3 S5



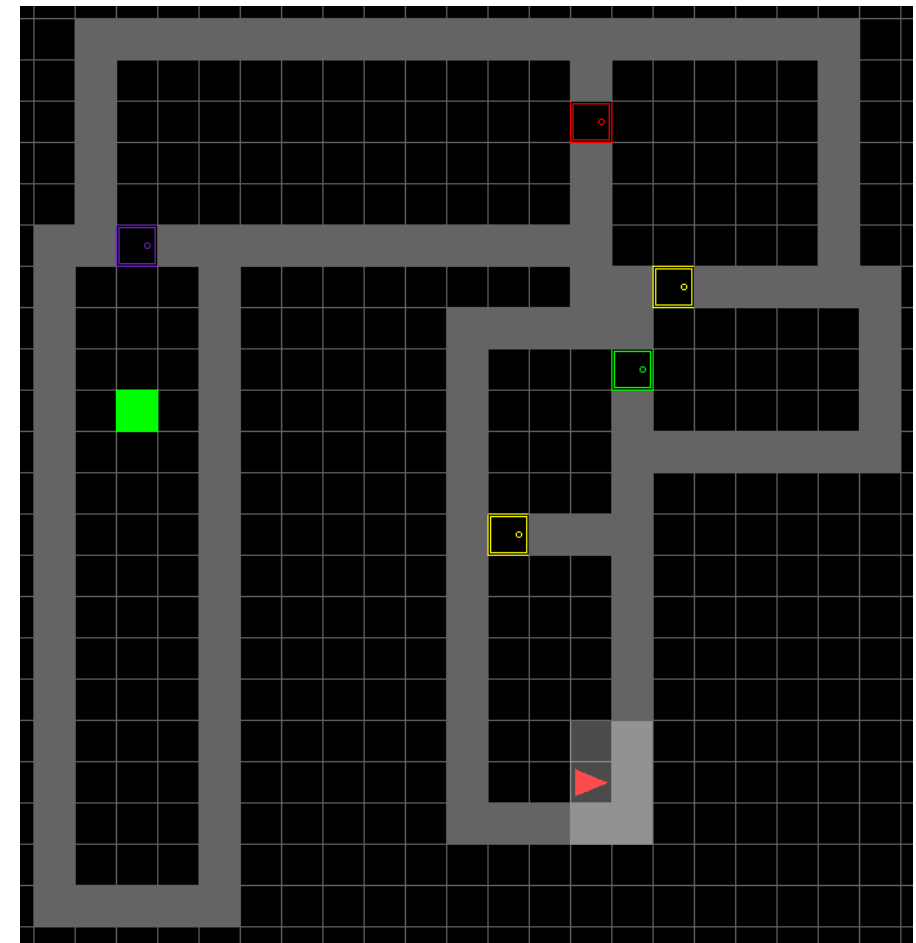
Challenging Transfer



N2 S10



Increasing
of rooms
and max size



N6 S25



Success Evaluation

Success: % of times goal reached over 512 different environments
S.E.M.: Standard error of mean over 10 random seeds

Method	MR-N3S4	MR-N5S4	MR-N6S25
$p_\phi(Z_t S_t, \Omega)$ pretrained on	MR-N2S6	MR-N2S6	MR-N2S10
InfoBot [Goyal <i>et al.</i> , 2019]	90%	85%	N / A
InfoBot (Our Implementation)	99.9% \pm 0.1%	79.1% \pm 11.6%	90.9% \pm 1.2%
Count-based Baseline	99.7% \pm 0.1%	99.7% \pm 0.1%	86.8% \pm 2.2%
DIAYN	99.7% \pm 0.1%	95.4% \pm 4.1%	0.1% \pm 0.1%
Random Network	99.9% \pm 0.1%	98.8% \pm 0.7%	79.5% \pm 5.2%
Heuristic Baseline	N / A	N / A	85.9% \pm 3.0%
Ours ($\beta = 10^{-2}$)	99.3% \pm 0.3%	99.4% \pm 0.2%	92.9% \pm 1.2%

Success % \pm s.e.m.

Easy Transfer

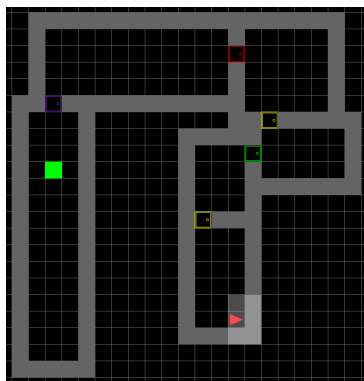
Challenging Transfer



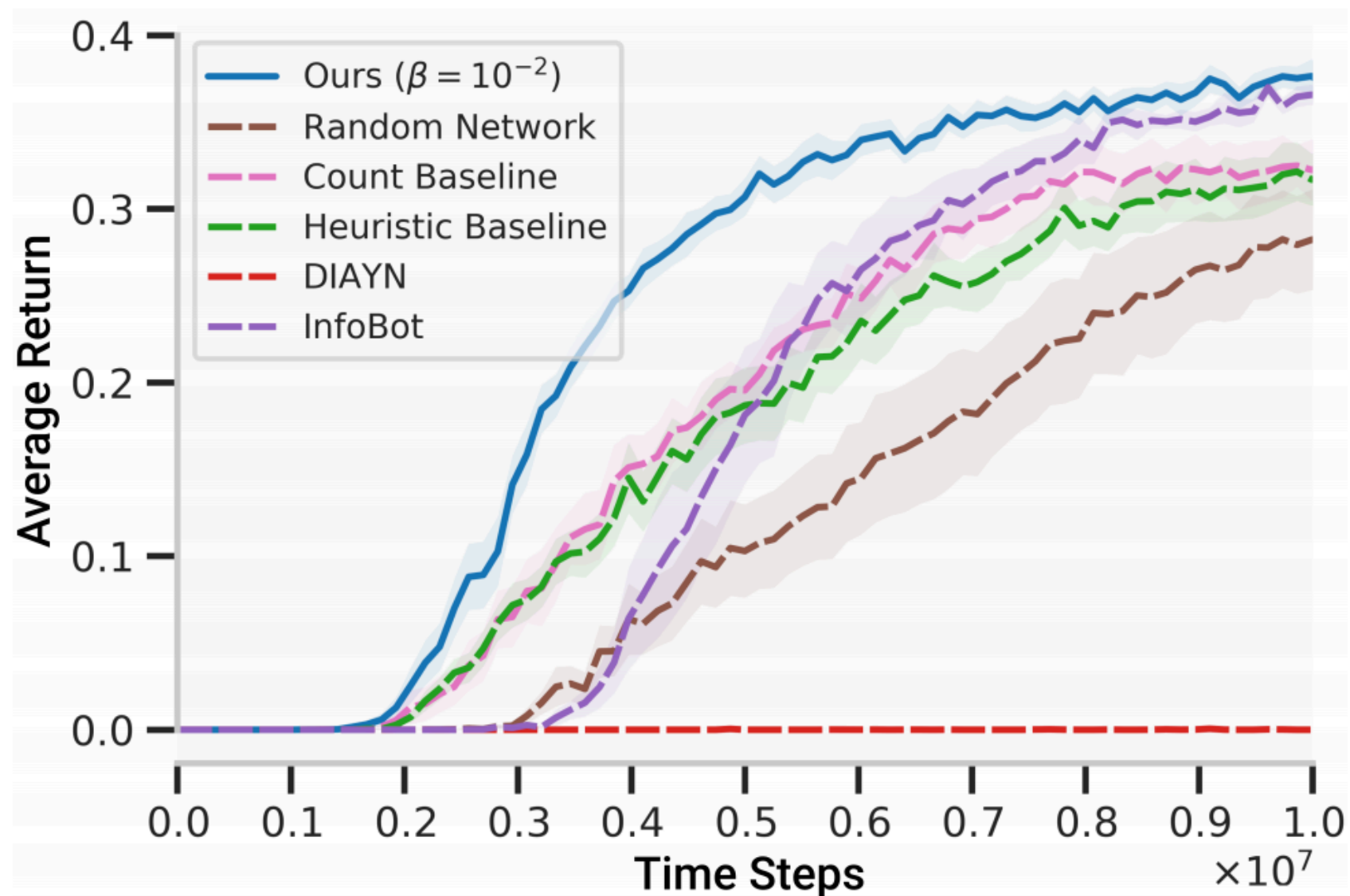
Average Return Evaluation



Sample Efficiency



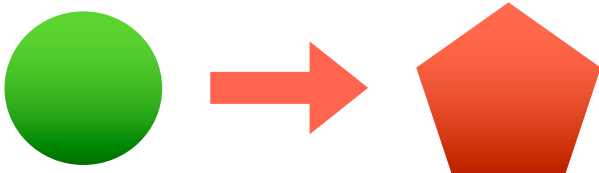
N = 6 rooms
S = 25 max room size

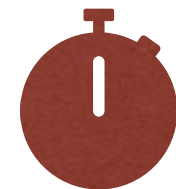
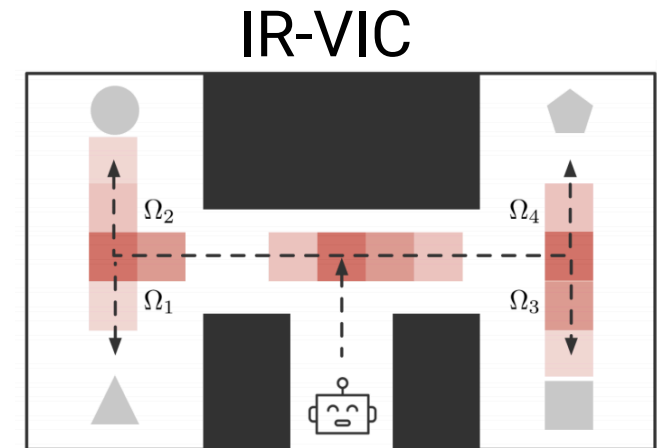


Challenging Transfer



Summary

- **Unsupervised objective** for sub-goal discovery
- **Transferable** sub-goals 
- **Better exploration** and **sample efficiency** in hard exploration tasks



Visit our poster or watch our 20 minute video for more details!

Code (coming soon): github.com/nirbhayjm/irvic

