

Learning to Compose Skills Himanshu Sahni, Saurabh Kumar, Farhan Tejani, Charles Isbell

Georgia College of Tech Computing

College of Computing, Georgia Institute of Technology contact: himanshu@gatech.edu

Abstract

We present a differentiable framework capable of learning a wide variety of compositions of simple policies, which we call skills. By recursively composing skills with themselves, hierarchies can be created that display complex behavior. Skill networks are trained to generate skill-state embeddings which are provided as inputs to a trainable composition function, which in turn outputs a policy for the overall task. Our experiments on an environment consisting of multiple collect and evade tasks show that this architecture is able to quickly build complex skills from simpler ones. The learned composition function displays transfer to unseen combinations of skills, allowing for zero-shot generalizations.

Background

Reinforcement Learning

- Maximize a notion of long term reward
- Typically formulated as an MDP <S, A, T, R, γ>

Policy Gradient

$$J(\theta) = \mathbb{E}_{p(S_{1:T};\theta)} \sum_{t=1}^{T} r_t = \mathbb{E}_{p(S_{1:T};\theta)}[R]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{p(S_{1:T};\theta)} \sum_{t=1}^{T} \nabla_{\theta} log[\pi(a_t|s_{1:t};\theta)] R_t$$

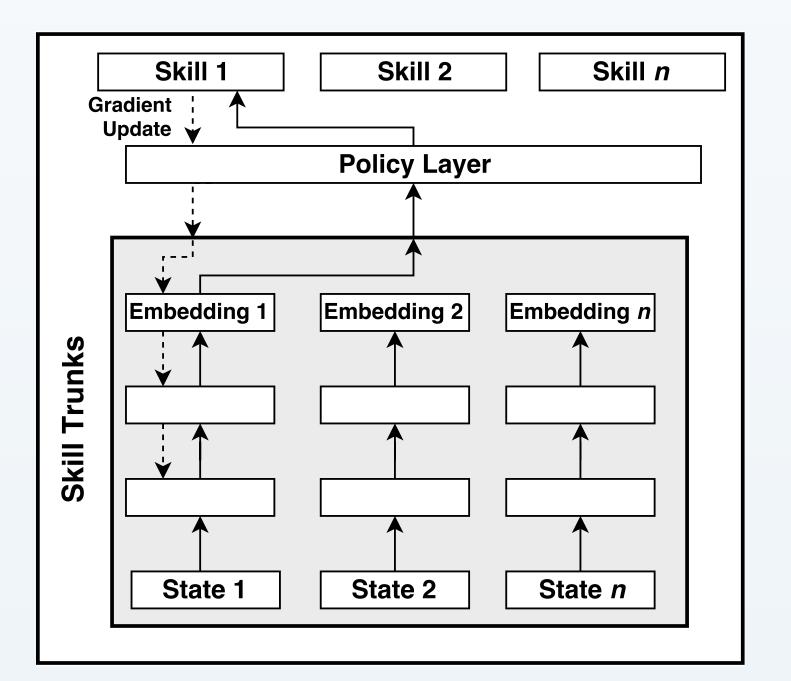
Linear Temporal Logic

- \mathcal{U} until
- \langle eventually
- · O- next
- Combined with basic logical connectives



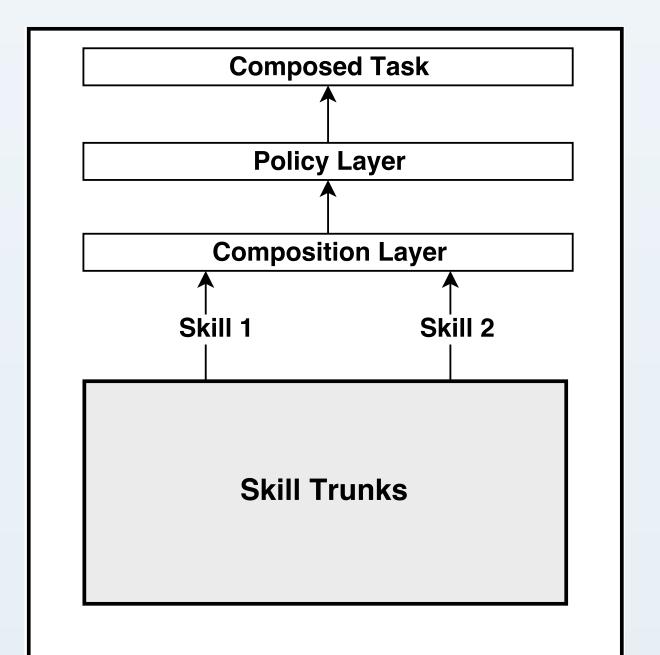
Can specify a wide variety of RL tasks

ComposeNet



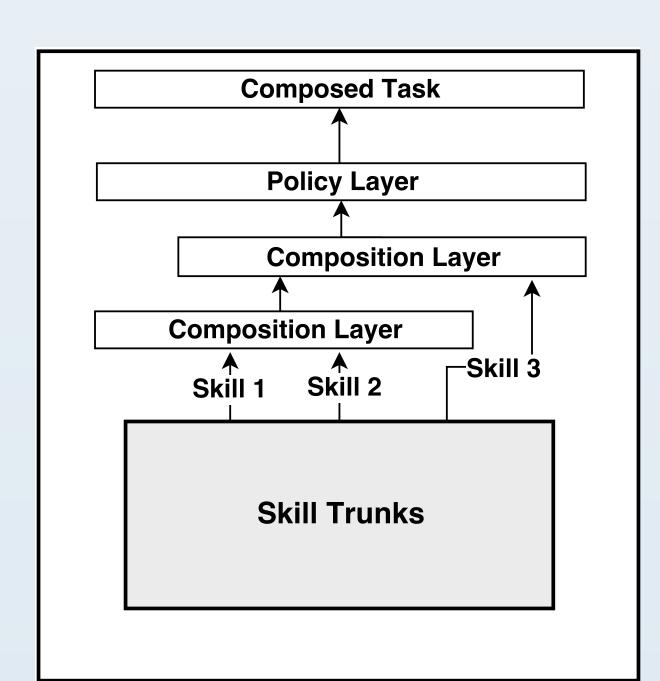
Skill trunks are trained with common policy layer

This forces topmost layer of skill trunks to embed skill-state information.



Trunk and policy layer weights are frozen.

Skill-state embeddings are composed and fed into policy layer.

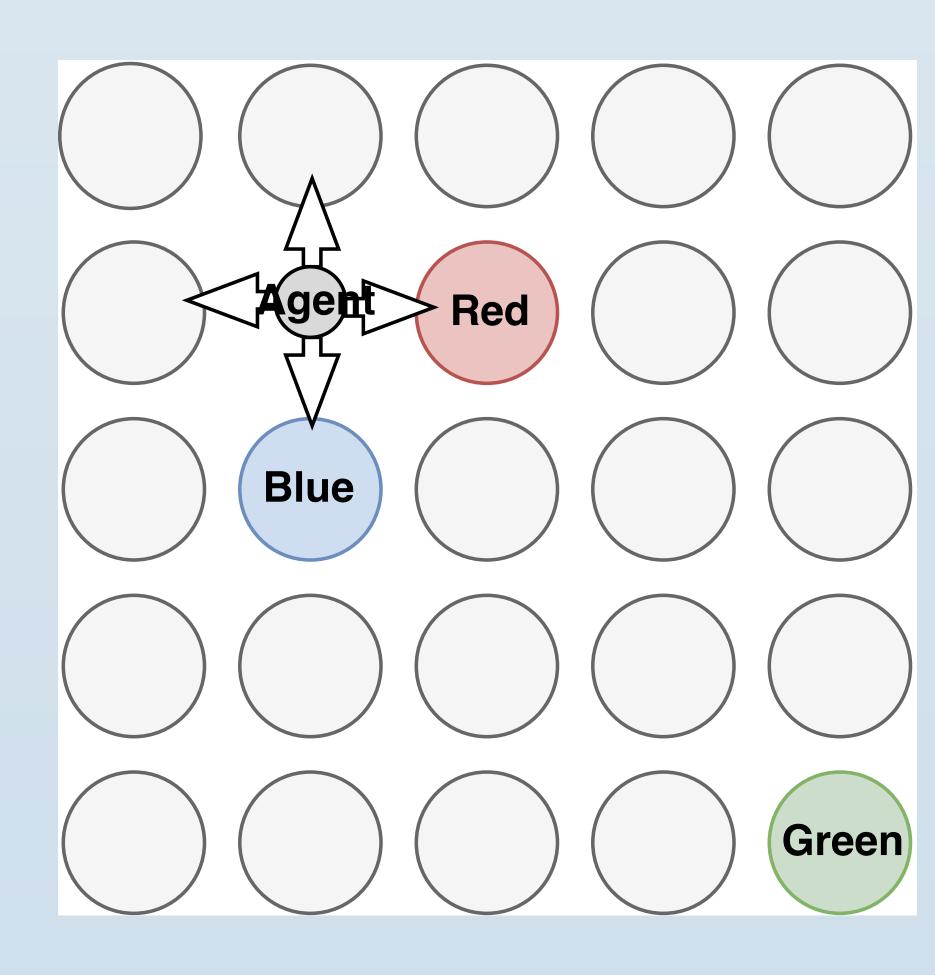


Hierarchies of compositions can be created post-hoc.

$$C: \langle \iota_e^{(1)}(S), \iota_e^{(2)}(S), \dots \iota_e^{(n)}(S) \rangle \to \iota_c(S)$$

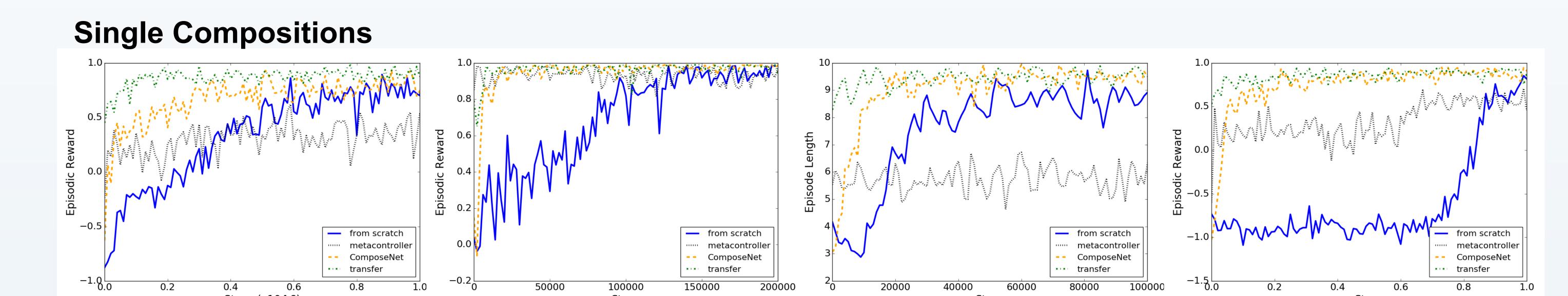
$$\pi_k: \iota_e^{(k)}(S) \to p(A)$$

Environment



Agent must collect target objects and evade enemies.

Results



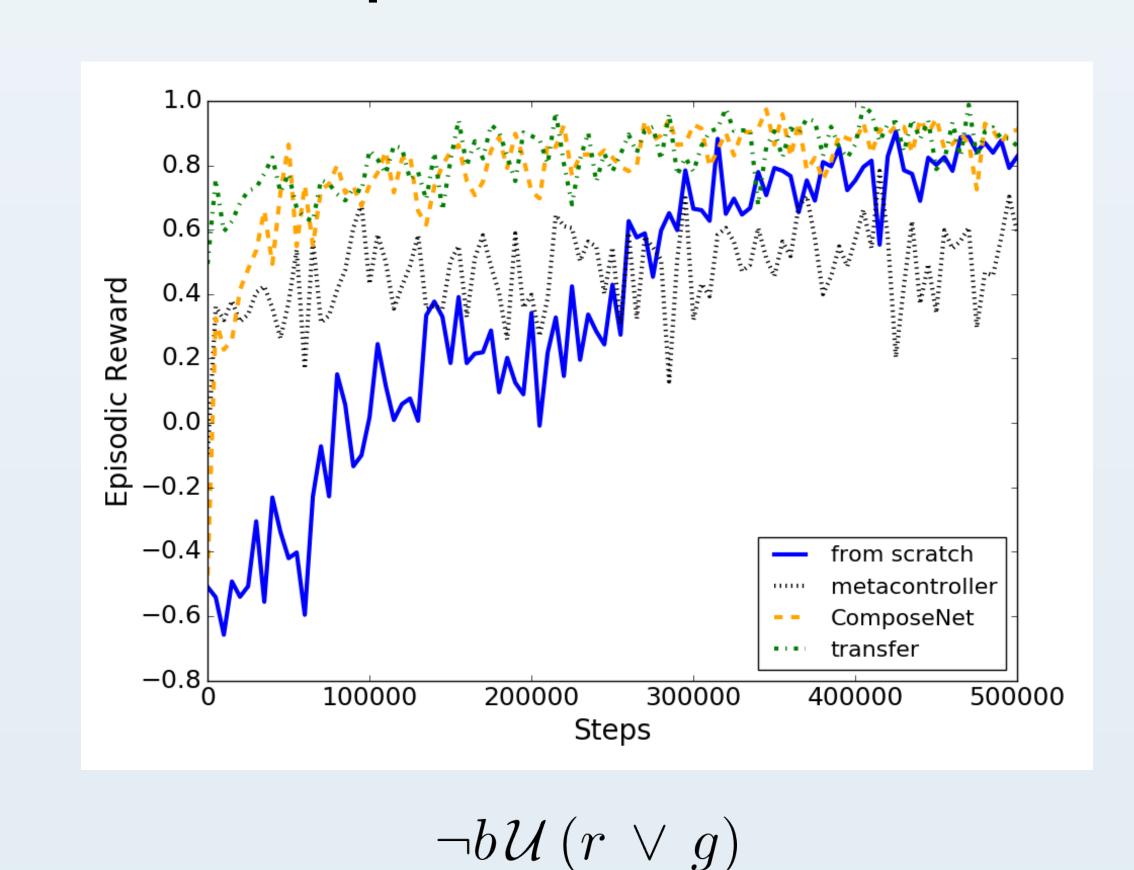
 $\Diamond r \vee \Diamond b$

zero-shot: 0.79

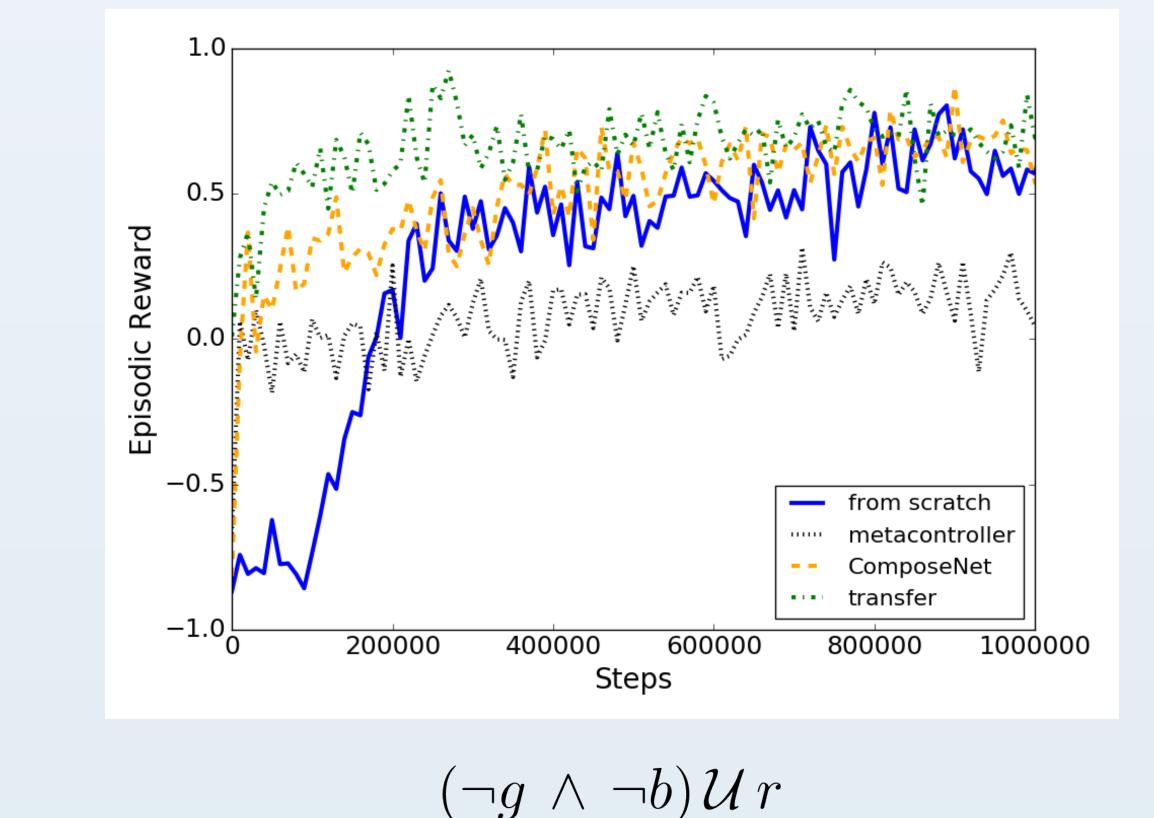
Hierarchical Compositions

 $\neg g \mathcal{U} b$

zero-shot: 0.45



zero-shot: 0.49



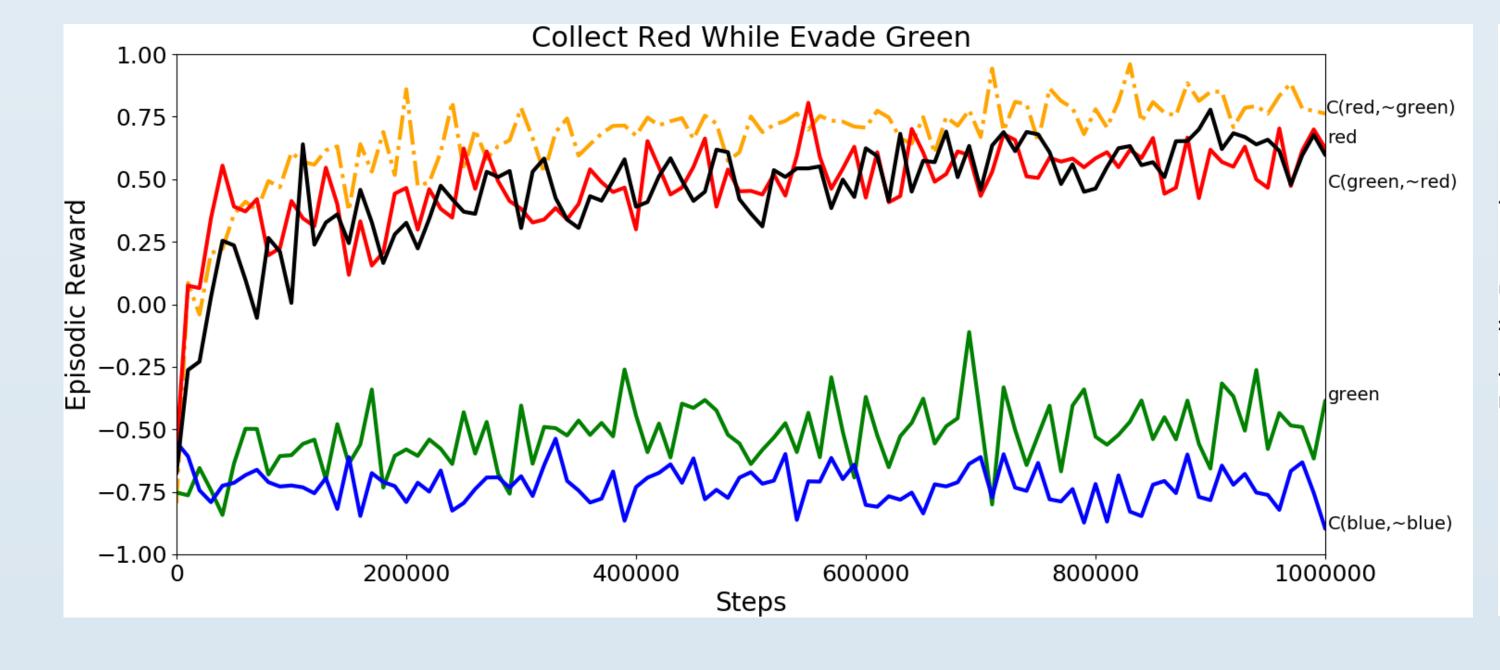
 $\Diamond(r \land \Diamond g)$

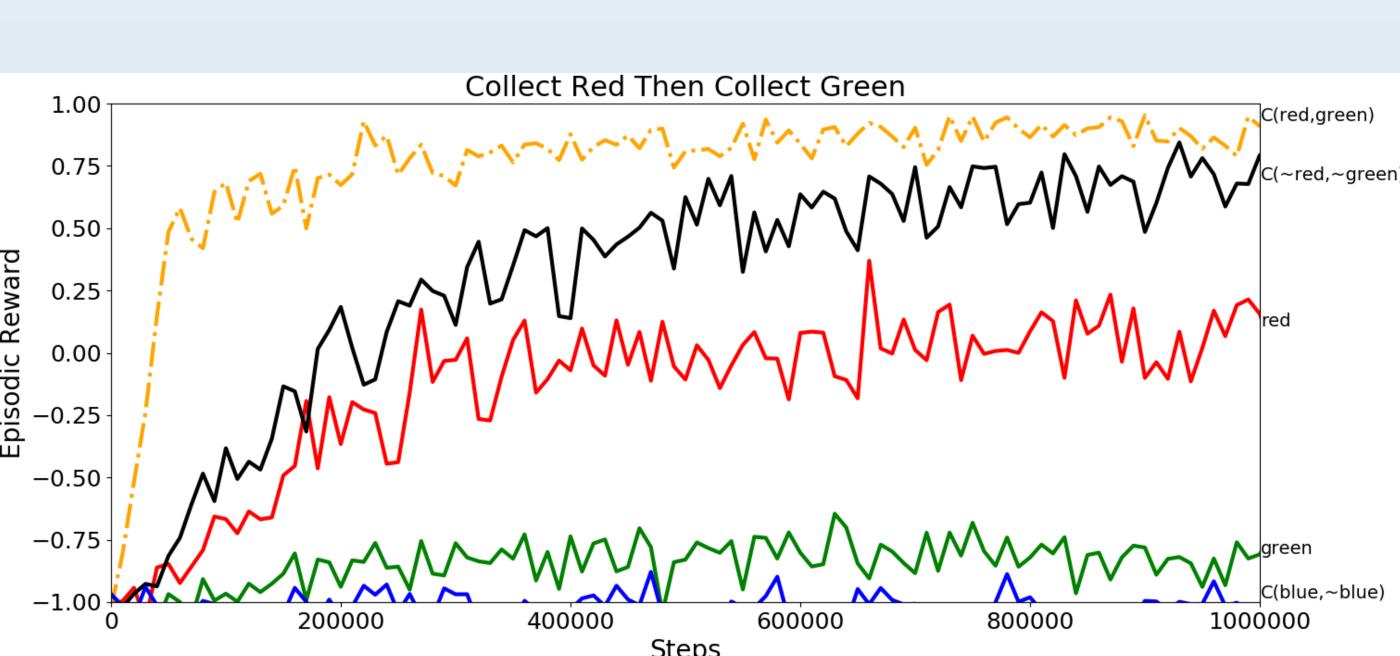
zero-shot: 0.53

 $\Box \neg r \land \Box \neg q$

zero-shot: 8.28

Ablations





zero-shot: 0.01

Conclusion and Future Work

- Individual skills can be successfully composed with the ComposeNet architecture.
- Learned faster than either of the baselines.
- Some zero-shot generalization, can be fine tuned to optimality.
- Hierarchies can be created by recursive compositions.
- Key to this are skill-state embeddings and a trainable composition function.
- Stay tuned for more compositions and environments!