

On Policy Reuse: An Expressive Language for Representing and Executing Policies that Call Other Policies

Blai Bonet, Dominik Drexler, Hector Geffner

ICAPS, 2024

Hector Geffner
RWTH Aachen University
Aachen, Germany

Linköping University
Linköping, Sweden



Motivation

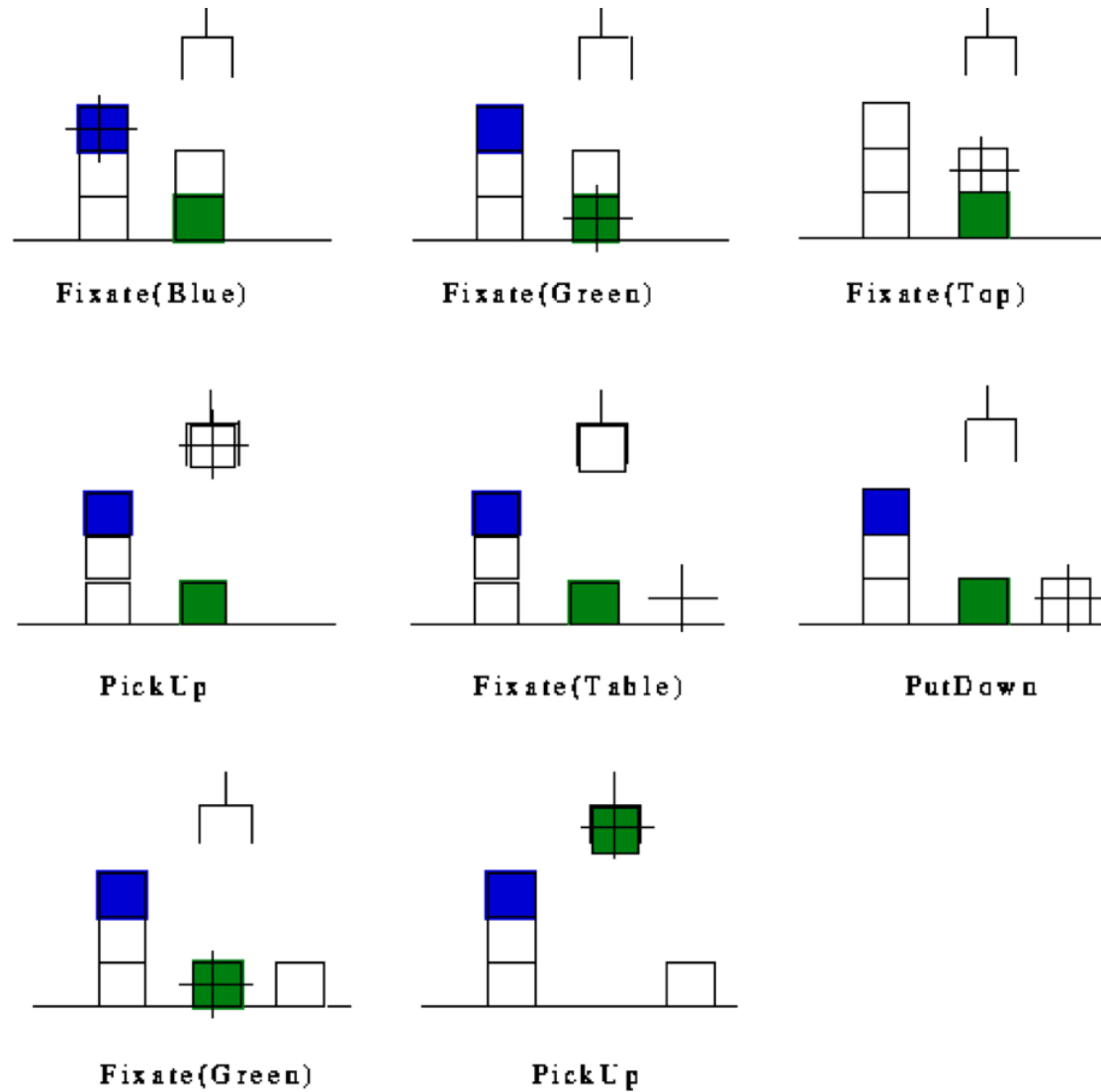
More expressive languages for **encoding** and **learning** general policies and sketches that support:

- **Reuse:** ability to call other policies **by passing parameters**
 - ▷ **Composition** and orchestration of subpolicies
 - ▷ **Bottom-up** construction of hierarchies, as opposed to top-down
 - ▷ **Answers:** “Can policy for $on(x, y)$ be reused to construct arbitrary towers?”
- **Indexicals:** ability to refer to objects **functionally**, not by name
 - ▷ **Features** for capturing general policies/sketches simplified
 - ▷ **Active perception:** what to observe and when
 - ▷ **Determine action to do** without considering other actions/transitions

Related Research Threads

- **Planning programs and inductive programming** [2, 11, 12, 5].
 - ▷ Dreamcoder: Growing generalizable knowledge with program learning. K. Ellis *et al.*; 2020
 - ▷ Generalized planning as heuristic search. J. Segovia, S. Jimenez, A. Jonsson, AIJ 2021
- **General policies** [9, 10, 6], [13], [15, 8, 14].
 - ▷ Learning generalized policies using concept languages. M. Martin, H. G., KR 2000
- **Deictic representations** [4, 1, 3, 7].
 - ▷ David Chapman. Penguins can make cake. AI Magazine, 1989
 - ▷ Deictic codes for the embodiment of cognition. D. Ballard *et al.*, BBS 1997
 - ▷ The thing that we tried didn't work very well: Deictic representation in RL, S. Finney *et al.*, UAI 2022.

Example: Pick up green blocks; Ballard *et al.* 1997



This Work

Extensions to the language of **general policies** and **sketches**:

- Indexical pointers to objects
- Memory states
- Ground actions
- Modules that call other modules (reuse)

Example: General Policy for $clear(x)$

- **Policy** π for class \mathcal{Q}_{clear} of problems with goal $clear(x)$ in Blocks:

$$\{\neg H, n > 0\} \mapsto \{H, n \downarrow\}$$

$$\{H, n > 0\} \mapsto \{\neg H\}$$

- **Features** $\Phi = \{H, n\}$: 'holding' and 'number of blocks above x '
- **Meaning:**
 - ▷ If $\neg H$ & $n > 0$, move to successor state where H holds and n **decreases**
 - ▷ If H & $n > 0$, move to successor state where $\neg H$ holds, n **doesn't change**
- **Shortcomings:**
 - ▷ Policy doesn't select actions directly; e.g. `pickup(A)`, if A top block above x
 - ▷ Feature n for 'number of blocks above x ', is "complex"

Example: New indexical policy for $clear(x)$

Concepts: used as features and to sample objects

- H_1 Boolean, whether block in τ_1 is being held
- $Table_1$: Boolean, whether block in τ_1 on table
- X : concept only contains given block x
- T_0 : concept that contains block on block in register τ_0 (if any)
- T_1 : concept that contains block on block in register τ_1 (if any)
- Initial memory state is always m_0 ; rule application change m_i

% Internal rules (update registers and internal memory; no state transitions involved)

$r_0 = m_0 \parallel \{X > 0\} \mapsto \{Load(X, \tau_0), T_0?\} \parallel m_1$ (Load x into register τ_0)

$r_1 = m_1 \parallel \{T_0 > 0\} \mapsto \{Load(T_0, \tau_1), T_1?\} \parallel m_2$ (Load block on x in τ_1 , if any)

$r_2 = m_2 \parallel \{T_1 > 0\} \mapsto \{Load(T_1, \tau_1), T_1?\} \parallel m_2$ (**Loop.** Load block on τ_1 in τ_1)

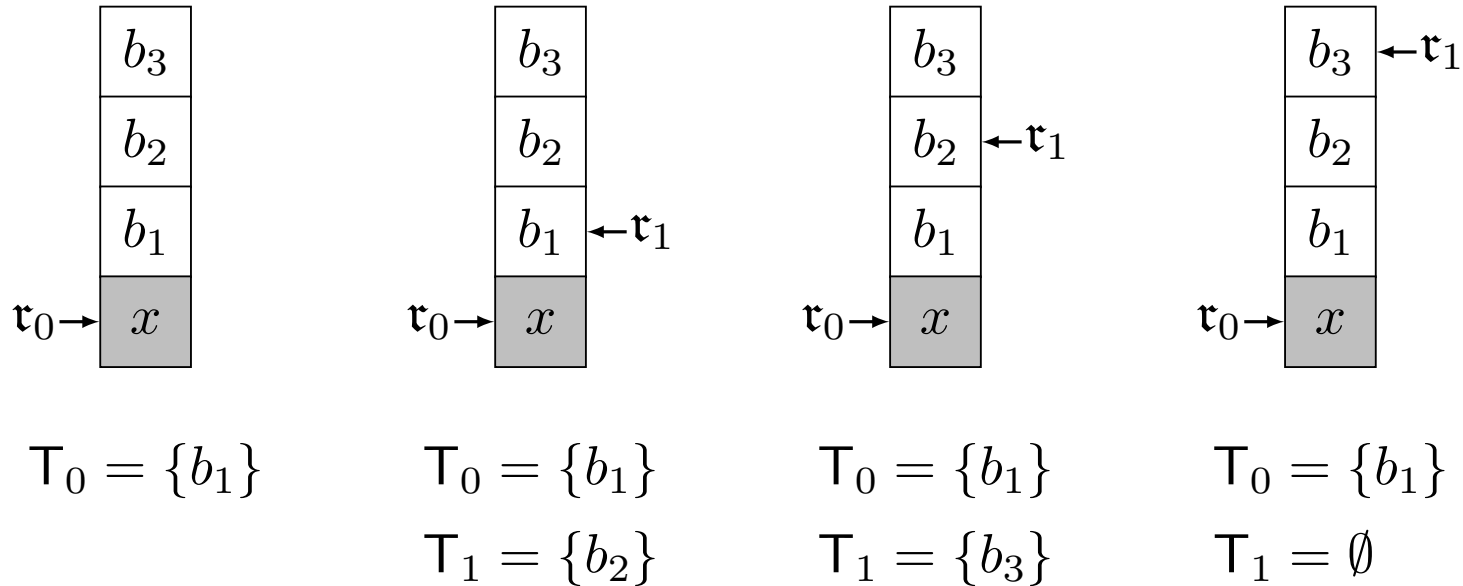
$r_3 = m_2 \parallel \{T_1 = 0\} \mapsto \{\} \parallel m_3$ until no such blocks, then

% External rules (state transitions involved)

$r_4 = m_3 \parallel \{\neg H_1\} \mapsto \{H_1\} \parallel m_3$ Unstack τ_1)

$r_5 = m_3 \parallel \{H_1\} \mapsto \{Table_1, \neg H_1\} \parallel m_1$ (Put block being held on table, and **loop**)

Example: Execution of new indexical policy for $clear(x)$



- Initially, load x in register r_0 ; equivalently, **mark** x with r_0
- Put r_1 **mark** on block that is on the one marked with r_0
- Move r_1 **mark** to block that is on the one marked with r_1
- Until block with r_1 **mark** is clear and can be picked up directly

Extended Sketch/Policy Language

- **Concepts** C (unary predicates) used explicitly as **Boolean features**, $C > 0$, **numerical features** C_{\downarrow} , and for **sampling objects**
- **Registers** r_i can be “loaded” with objects sampled from concepts; $Load(C, r_i)$; registers are concepts too.
- **Memory states** m_i control flow along with Boolean conditions; e.g.,
 $m_1 \parallel \{C\} \mapsto \{E\} \parallel m_2$
- Rules with load effects or empty effects deemed as **internal rules**; others as **external rules**
- Memory states of internal rules and external rules different
- See paper for **formal syntax** and **semantics**

Modules: Reusing Policies

- Policies and sketches wrapped into *modules*
- **Modules may call other modules and do recursion passing parameters**
- Execution model uses a **stack** and caller/callee protocol, as in prog. languages
- Orchestration of collections $\{\text{mod}_0, \text{mod}_1, \text{mod}_2, \dots\}$ of modules
- Additional external rules in modules:
 - ▷ **Call rules:** $m \parallel C \mapsto \text{mod}(C_1, C_2, \dots, C_k) \parallel m'$ where C is condition, and m and m' are memory states, to call mod with C_1, \dots, C_k as arguments
 - ▷ **Do rules:** $m \parallel C \mapsto \text{act}(C_1, C_2, \dots, C_k) \parallel m'$ to apply a **ground action** $\text{act}(o_1, o_2, \dots, o_k)$ with objects $o_i \in C_i$, for $i = 1, 2, \dots, k$

Example: Modules for $on(x, y)$

Module $On(X, Y)$:

- $r_0 = m_0 \parallel \{\neg On\} \mapsto \text{Clear}(X) \parallel m_1$ (Call **Clear** with argument X)
- $r_1 = m_1 \parallel \{\} \mapsto \text{Clear}(Y) \parallel m_2$ (Call **Clear** with argument Y)
- $r_2 = m_2 \parallel \{\neg H_X\} \mapsto \{H_X\} \parallel m_3$ (Pick block x , either **unstack** or **pickup**)
- $r_3 = m_3 \parallel \{H_X\} \mapsto \text{stack}(X, Y) \parallel m_3$ (Apply **stack** to put x on y)

Module $Clear(X)$:

- $r_0 = m_0 \parallel \{X > 0\} \mapsto \{Load(X, \tau_0), T_0?\} \parallel m_1$ (Load x in register τ_0)
- $r_1 = m_1 \parallel \{T_0 > 0\} \mapsto \{Load(T_0, \tau_1), T_1?\} \parallel m_2$ (Load block above x in τ_1 , if any)
- $r_2 = m_2 \parallel \{T_1 > 0\} \mapsto \{Load(T_1, \tau_1), T_1?\} \parallel m_2$ (Loop. Load block above τ_1 in τ_1)
- $r_3 = m_2 \parallel \{T_1 = 0\} \mapsto \{\} \parallel m_3$ (Go to external rules)
- $r_4 = m_3 \parallel \{\neg H\} \mapsto \text{unstack}(\tau_1, B) \parallel m_3$ (Apply **unstack** to pick τ_1)
- $r_5 = m_3 \parallel \{H\} \mapsto \text{putdown}(\tau_1) \parallel m_1$ (Apply **putdown** to put τ_1 on table)

Example: Building One Tower with Module $\text{Tower}(\mathbf{O}, \mathbf{X})$

- Objective is to build tower $\bigwedge_{i=1}^k \text{on}(x_i, x_{i-1}) \wedge \text{ontable}(x_0)$
- Role argument $\mathbf{O} = \{(x_i, x_{i-1}) \mid i = 1, \dots, k\}$
- \mathbf{X} is concept for *lowest* block in tower that is *misplaced*
- \mathbf{M} is concept for block to be placed on τ_0 according to \mathbf{O} (if any)
- \mathbf{W} is concept for block below τ_0 according to \mathbf{O} (if any)

Module $\text{Tower}(\mathbf{O}, \mathbf{X})$:

$r_0 = m_0 \parallel \{X > 0\} \mapsto \{\text{Load}(X, \tau_0), \mathbf{M}?, \mathbf{W}?\} \parallel m_1$ (Load X into register τ_0)

$r_1 = m_1 \parallel \{W = 0\} \mapsto \text{On-Table}(\tau_0) \parallel m_2$ (**On-Table** to put X on table)

$r_2 = m_1 \parallel \{W > 0\} \mapsto \text{On}(\tau_0, W) \parallel m_2$ (**On**(τ_0, W) to well-place τ_0)

$r_3 = m_2 \parallel \{M > 0\} \mapsto \text{Tower}(\mathbf{O}, \mathbf{M}) \parallel m_3$ (Continue building tower from \mathbf{M})

Example: Building Many Towers

- Argument O is **role** that contains the pairs describing the towers to build
- L is concept for *lowest misplaced blocks* according to O

Module `Blocks(O)`:

$r_0 = m_0 \parallel \{L > 0\} \mapsto \{Load(L, \tau_0)\} \parallel m_1$ (Load X into register τ_0)

$r_1 = m_1 \parallel \{\} \mapsto Tower(O, \tau_0) \parallel m_0$ (Build tower on τ_0)

Summary. Future

Language extensions for **encoding** and **learning** general policies and sketches:

- **Reuse** and **bottom up** composition of policies
- Don't **learn policies from scratch**; reuse those learned
- **Indexicals** (registers) simplify features, determine actions to do, active perception
- **Interpreter** available, but **not learning yet**
- **Limitations.** Language is:
 - ▷ “too much”: hard to learn and verify, too many alternative encodings
 - ▷ “too little”: flexibility lacking for handling negative interactions
- One step; others to follow

References

- [1] Philip E. Agre and David Chapman. What are plans for? *Robotics and Autonomous Systems*, 6:17–34, 1990.
- [2] Javier Segovia Aguas, Sergio Jiménez Celorrio, and Anders Jonsson. Generalized planning with procedural domain control knowledge. In *Proc. ICAPS*, pages 285–293, 2016.
- [3] Dana H. Ballard, Mary M. Hayhoe, Polly K. Pook, and Rajesh P. N. Rao. Deictic codes for the embodiment of cognition. *Behavioral and Brain Sciences*, 20:723–742, 1996.
- [4] David Chapman. Penguins can make cake. *AI magazine*, 10(4):45–45, 1989.
- [5] Kevin Ellis, Lionel Wong, Maxwell Nye, Mathias Sable-Meyer, Luc Cary, Lore Anaya Pozo, Luke Hewitt, Armando Solar-Lezama, and Joshua B Tenenbaum. Dreamcoder: growing generalizable, interpretable knowledge with wake–sleep bayesian program learning. *Phil. Trans. R. Soc. A*, 381:20220050, 2023.
- [6] Alan Fern, Sungwook Yoon, and Robert Givan. Approximate policy iteration with a policy language bias: Solving relational markov decision processes. *Journal of Artificial Intelligence Research*, 25:75–118, 2006.
- [7] Sarah Finney, Natalia Gardiol, Leslie Pack Kaelbling, and Tim Oates. The thing that we tried didn’t work very well : Deictic representation in reinforcement learning. *CoRR*, abs/1301.0567, 2013.
- [8] Sankalp Garg, Aniket Bajpai, and Mausam. Generalized neural policies for relational mdps. In *Proc. ICML*, 2020.
- [9] Roni Khardon. Learning action strategies for planning domains. *Artificial Intelligence*, 113:125–148, 1999.
- [10] Mario Martín and Hector Geffner. Learning generalized policies from planning examples using concept languages. *Applied Intelligence*, 20(1):9–19, 2004.
- [11] Javier Segovia-Aguas, Sergio Jiménez, and Anders Jonsson. Computing programs for generalized planning using a classical planner. *Artificial Intelligence*, 272:52–85, 2019.
- [12] Javier Segovia-Aguas, Sergio Jiménez, and Anders Jonsson. Generalized planning as heuristic search. In *ICAPS*, pages 569–577, 2021.
- [13] Siddharth Srivastava, Neil Immerman, and Shlomo Zilberstein. A new representation and associated algorithms for generalized planning. *Artificial Intelligence*, 175(2):393–401, 2011.

- [14] Simon Ståhlberg, Blai Bonet, and Hector Geffner. Learning general policies with policy gradient methods. In *Proc. KR*, pages 647–657, 2023.
- [15] Sam Toyer, Sylvie Thiébaux, Felipe Trevizan, and Lexing Xie. Asnets: Deep learning for generalised planning. *Journal of Artificial Intelligence Research*, 68:1–68, 2020.