# Opportunities and Challenges for Domain-Independent Planning with Deep Reinforcement Learning

Forest Agostinelli
University of South Carolina

# Students

## Ph.D. Students



Rojina Panta

Vedant Khandelwal

Misagh Soltani

Cale Workman
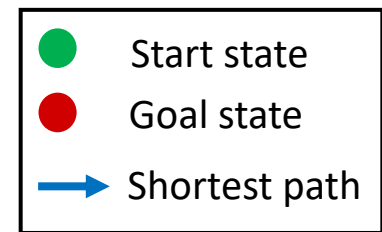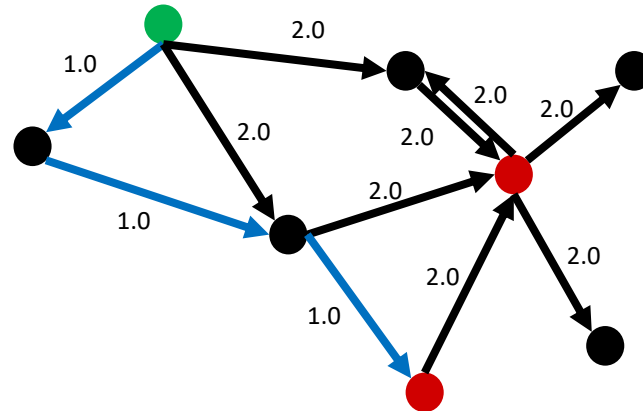
## B.S. Students



Christian Geils

William Edwards

# Outline

- Background
- Generalizing over states
- Generalizing over goals
- Generalizing over domains
- Towards obtaining approximately admissible heuristic functions
- Generalizing to domains with unknown transition functions

# Pathfinding

- The objective of pathfinding is to find a sequence of actions that forms a path between a given start state and a given goal
  - A goal is a set of states
  - Preference for minimum cost paths
- A pathfinding problem can be represented as a weighted directed graph where nodes represent states, edges represent actions that transition between states, and edge weights represent transition costs
  - The cost of a path is the sum of transition costs

# Pathfinding Domains

- Pathfinding problems can be found throughout mathematics, computing, and the natural sciences
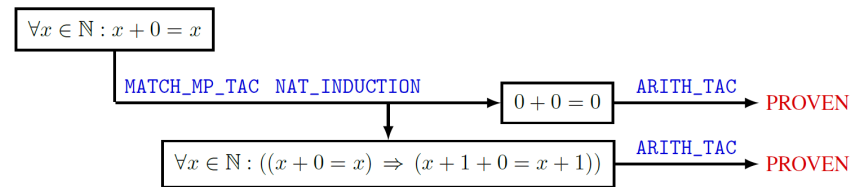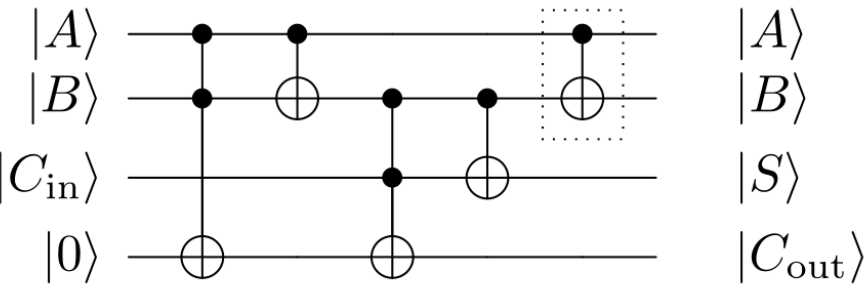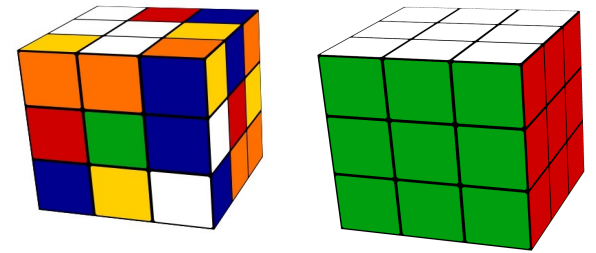  - Puzzle solving, chemical synthesis, quantum circuit synthesis, theorem proving, program synthesis, robotics
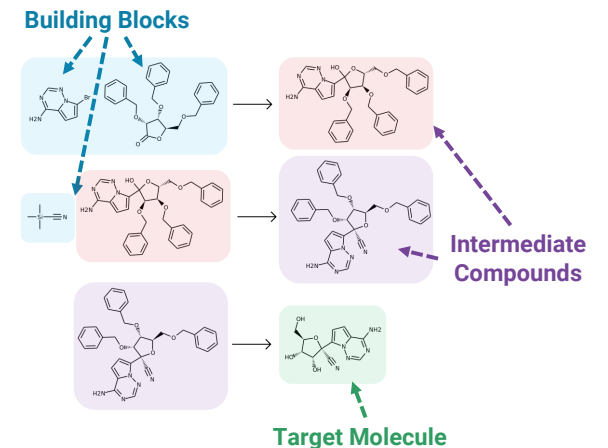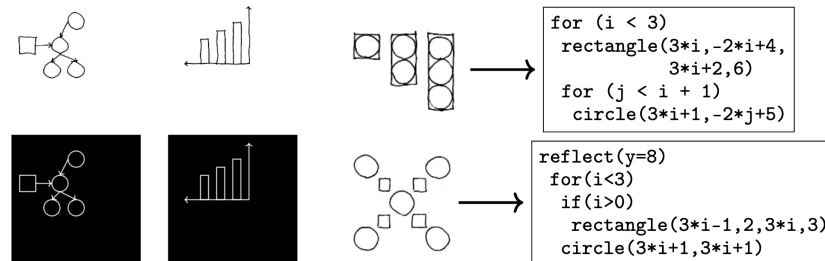


Figure 1: Formally proving $\forall x \in \mathbb{N} : x + 0 = x$.

# Pathfinding Domain Definition

- The entire state space graph cannot be given to a pathfinding problem solver because the number of states in a pathfinding problem can be very large.
  - Rubik's cube: $\sim 10^{19}$
  - 48-puzzle: $\sim 10^{62}$
  - Organic chemistry: $\sim 10^{60}$ (exact number unknown)
- Assumptions on what is given
  - Action space
  - State transition function
  - Transition cost function
  - Goal specification language
  - Goal test function
- Objective: Create a domain independent algorithm
  - Input: Pathfinding domain definition, start state, goal specification
  - Output: Path to a goal state

# Outline

- Background
- Generalizing over states
- Generalizing over goals
- Generalizing over domains
- Towards obtaining approximately admissible heuristic functions
- Generalizing to domains with unknown transition functions

# Learned Heuristic Functions

- Heuristic function maps a state to an estimate of the cost of a shortest path from that state, also known as the cost-to-go

# Value Iteration

- Value iteration is a dynamic programming algorithm and is a foundational algorithm in reinforcement learning

- In the context of pathfinding, value iteration is an algorithm for computing the cost-to-go of finding a shortest path for each state in the state space

- Tabular value iteration loops over all states and applies the following update until convergence ($h$ stops changing)
  - $h(s) = \min_a(c^a(s) + h(T(s,a)))$
  - Guaranteed to converge to $h^*$ in the tabular setting

- $s$: state

- $a$: action

- $T$: state transition function

- $c^a$: transition cost function

# Value Iteration: Visualization

- Actions: up, down, left, right

- Transition costs
  - 1 if square is blank
  - 10 if square has a rock
  - 50 if square has a plant

- Goal: shovel

- Updates propagate outwards from the goal

# Approximate Value Iteration

- As the state space grows, tabular value iteration becomes infeasible
- Approximate value iteration uses an approximation architecture to approximate the value iteration update
- When using a deep neural network as the approximation architecture, we refer to this as deep approximate value iteration (DAVI)
- The update is approximated using the following loss function
    - $L(\theta) = \left( \min_{a}(c^a(s) + h_{\theta^-}(T(s, a))) - h_\theta(s) \right)^2$
    - Target is set to zero if $s$ is a terminal state
- $s$: state
- $a$: action
- $T$: state transition function
- $c^a$: transition cost function
- $\theta$: parameters
- $\theta^-$: parameters for target network
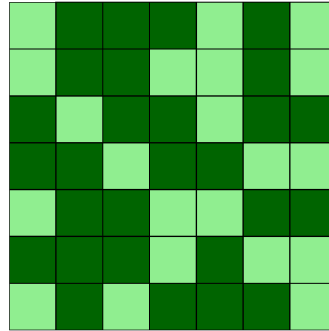    - Is periodically updated to $\theta$ throughout training

# Application to Puzzle Solving



Rubik's cube



24 puzzle



Lights Out (7×7)



Sokoban

1. Rubik's Cube
2. 15-puzzle
3. 24-puzzle
4. 35-puzzle
5. 48-puzzle
6. Lights Out
7. Sokoban

Largest state space is $3.0 \times 10^{62}$ (48-puzzle)

- Prioritized sweeping: Generate training data by taking moves in reverse from the goal



Goal

# Training

- Deep neural network
  - Input layer -> Two fully connected layers -> Four residual blocks -> Linear output layer
  - Same type of architecture used for all puzzles
    - 24-puzzle has two more residual blocks
- Training
  - Batch size of 5,000
  - ~1,000,000 training iterations
  - Parameters for target network updated when loss goes below some target threshold
    - Future work updates based on greedy policy performance

# Greedy Policy Performance

- Behave greedily with respect to the heuristic function

- $\pi(s) = \operatorname*{argmin}_a (c^a(s) + h_\theta(T(s,a)))$

- Does not solve all states

# Integration with A* Search

- Learned heuristic function can be used as a heuristic in A* search
- A* Search
  - Maintains a search tree where nodes are states and edges are actions
  - Initialized with a start node representing the start state
  - Expands nodes according to the priority
    - $f(n) = g(n) + h(n.s)$
    - $f(n)$: cost
    - $g(n)$: path cost (cost to get from start node to $n$)
    - $h(n.s)$: heuristic (estimated cost-to-go from $n.s$ to a closest goal state)
  - Terminates when a node associated with a goal state is selected for expansion
- Weighted A* Search
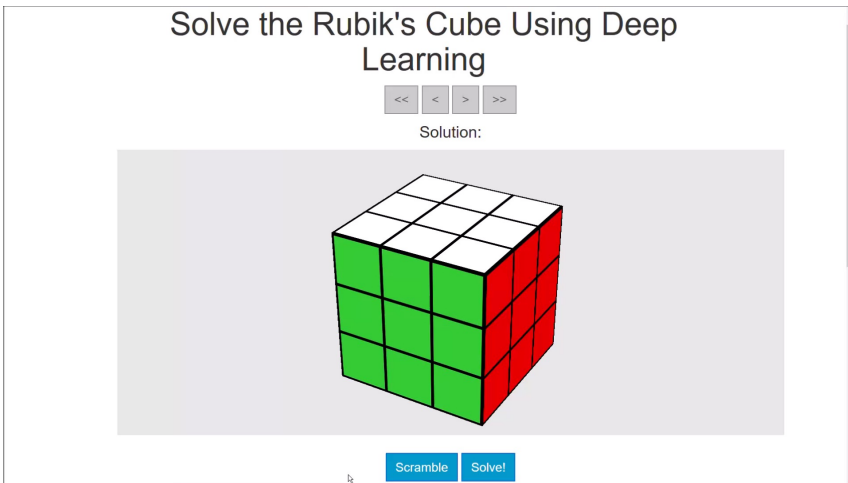  - Decreasing the weight on the path cost may result in expanding fewer nodes while possibly increasing the length of paths found
  - $f(n) = \lambda * g(n) + h(n.s)$
- Batch weighted A* Search
  - Can take advantage of parallelism provided by GPUs by expanding multiple nodes at a time

Agostinelli, Forest, et al. "Obtaining approximately admissible heuristic functions through deep reinforcement learning and A* search." *ICAPS PRL Workshop*. 2021.
Li, Tianhua, et al. "Optimal search with neural networks: Challenges and approaches." *Proceedings of the International Symposium on Combinatorial Search*. Vol. 15. No. 1. 2022.

# DeepCubeA: Results

- When applied to seven different puzzles, it was able to solve all test instances and found a shortest path in the majority of verifiable cases

- http://deepcube.igb.uci.edu/

| Puzzle | Solution Length | Percent Optimal | Time (seconds) |
|---|---|---|---|
| Rubik's Cube | 21.50 | 60.3% | 24.22 |
| 15-puzzle | 52.03 | 99.4% | 10.28 |
| 24-puzzle | 89.49 | 96.98% | 19.33 |
| 35-puzzle | 124.64 | N/A | 28.45 |
| 48-puzzle | 253.35 | N/A | 74.46 |
| Lights Out | 24.26 | 100.0% | 3.27 |
| Sokoban | 32.88 | N/A | 2.35 |

Agostinelli, Forest, et al. "Solving the Rubik's cube with deep reinforcement learning and search." *Nature Machine Intelligence* 1.8 (2019): 356-363.
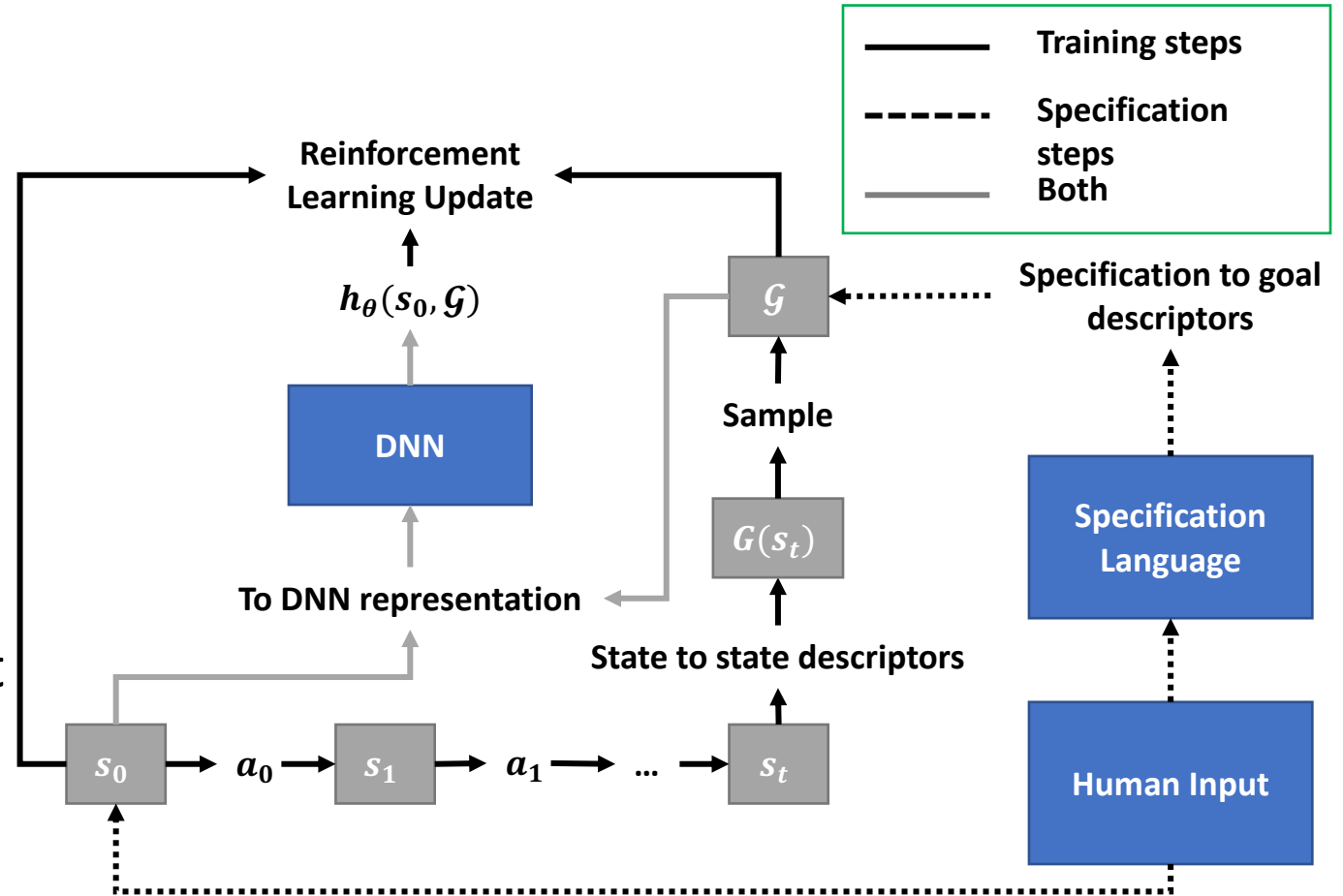
# Limitations

- The goal is pre-determined
  - Specifying a new goal requires re-training the DNN
- The domain is pre-determined
  - A change in the state transition function requires re-training the DNN
- Heuristic functions are not as amenable to analysis as domain-independent heuristics derived from PDDL
  - No admissibility guarantees

Muppasani, Bharath, Vishal Pallagani, Biplav Srivastava, and Forest Agostinelli. "Specifying Goals to Deep Neural Networks with Answer Set Programming." *ICAPS HSDIP Workshop 2024*

# Outline

- Background
- Generalizing over states
- Generalizing over goals
- Generalizing over domains
- Towards obtaining approximately admissible heuristic functions
- Generalizing to domains with unknown transition functions

- In the previous work, the goal is predetermined

- We build on hindsight experience replay to generalize over sets of goal states

- In our work
  - State descriptors: assignments of values to variables
  - Specification language: Answer set programming (ASP)
  - ASP will be used to describe goals at a high-level using formal logic and an answer set solver will be used to find assignments that represent a subset of the goal
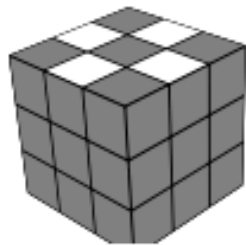
# State Representation

- In a given pathfinding domain, there are $V$ variables
  - A variable, $x_i$, can be assigned a single value from its (variable) domain, $D(x_i)$
- An assignment is an assignment is a set of assignments of values to variables $\{x_i = v_i\}$
  - All $v_i \in D(v_i)$
  - If $x_i$ is not in the assignment then it is unassigned
- An assignment is a complete assignment iff all variables have been assigned values
- A state is a complete assignment
- For example, for the Rubik's cube, variables are stickers and values are their colors
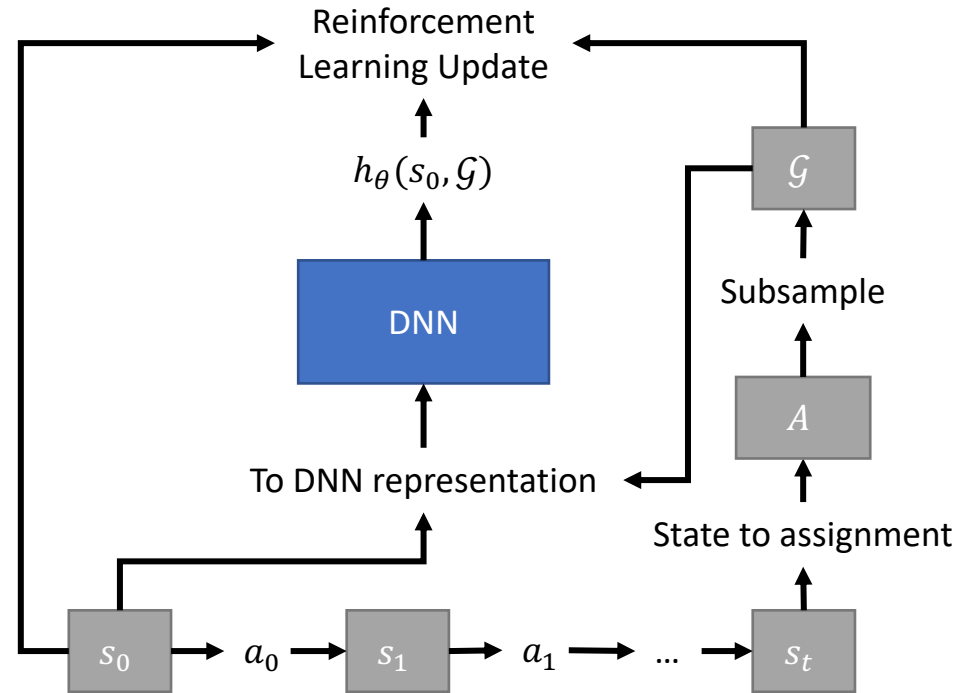
# Goal Representation

- An assignment is a partial assignment iff at least one variable has not been assigned a value

- A goal is a complete or partial assignment

- An assignment, $A$, represents a set of states, $\mathcal{S}_A$
  - A complete assignment always represents a set of states of size 1

- A state, $s$, is in $\mathcal{S}_A$ iff $A \subseteq s$
  - In other words, all assignments in $A$ are present in $s$
  - An empty assignment represents the set of all possible states

- For example, a visualization of an assignment for the "white cross" pattern for the Rubik's cube and a state that is in the set of states represented by this assignment

# Training

- Generate a start state
- Take a random walk whose length is somewhere between 0 and T
  - Future work could use artificial curiosity
- Convert the end state to its representation as an assignment
- Subsample to obtain a goal
- Convert this representation into one suitable for the DNN
  - One-hot representation
  - Graph
  - Etc.
- RL Update

  - $L(\theta) = \left( \min_{a}(c^a(s) + h_{\theta-}(T(s,a)), \mathcal{G}) - h_\theta(s,\mathcal{G}) \right)^2$

# Experiments

- ASP will be used to find assignments; therefore, we compare our method (DeepCubeA$_g$) to other methods capable of finding paths to goals that can be represented as complete or partial assignments

- 500-1,000 test start and goal pairs

- 200 second time limit to solve test states

- **DeepCubeA**
  - Predefined goal

- **Fast Downward Planner**
  - Can automatically construct heuristics given a formal definition of the domain (including the transition function) in the planning domain definition language (PDDL)
  - Goal count heuristic, fast forward heuristic, causal graph heuristic
  - A* search

- **PDBs**
  - Divides into subproblems and enumerates all possible combinations of the subproblem to create heuristic
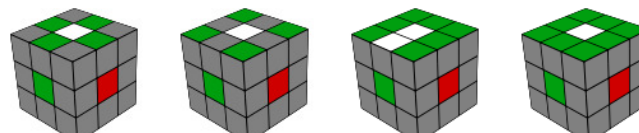  - Predefined goal
  - IDA* search

# Performance

- Canon: Canonical goal states

- Rand: Random assignment selected as goal
  - Can be as small as the empty assignment
  - Methods that require a pre-definied goal cannot be applied to this scenario without considerable overhead

- PDBs+: Also includes group theory knowledge

- DeepCubeA$_g$ consistently outperforms fastdownard in terms of percentage of states solved

| Puzzle | Solver | Path Cost | % Solved | % Opt | Nodes | Secs | Nodes/Sec |
|---|---|---|---|---|---|---|---|
| RC (Canon) | PDBs$^+$ | **20.67** | **100.00%** | **100.0%** | **2.05E+06** | **2.20** | **1.79E+06** |
| | DeepCubeA | 21.50 | **100.00%** | 60.3% | 6.62E+06 | 24.22 | 2.90E+05 |
| | DeepCubeA$_g$ | 22.03 | **100.00%** | 35.00% | 2.44E+06 | 41.99 | 5.67E+04 |
| | FastDown (GC) | - | 0.00% | 0.0% | - | - | - |
| | FastDown (FF) | - | 0.00% | 0.0% | - | - | - |
| | FastDown (CG) | - | 0.00% | 0.0% | - | - | - |
| RC (Rand) | DeepCubeA$_g$ | 15.22 | **99.40%** | - | 1.91E+06 | 32.24 | 5.19E+04 |
| | FastDown (GC) | 7.18 | 32.80% | - | 2.67E+06 | 13.79 | 1.41E+05 |
| | FastDown (FF) | 6.49 | 31.20% | - | 4.87E+05 | 13.83 | 2.93E+04 |
| | FastDown (CG) | 7.85 | 33.80% | - | 1.12E+06 | 11.62 | 5.81E+04 |
| 15-P (Canon) | PDBs | **52.02** | **100.00%** | **100.0%** | **3.22E+04** | **0.002** | **1.45E+07** |
| | DeepCubeA | 52.03 | **100.00%** | 99.4% | 3.85E+06 | 10.28 | 3.93E+05 |
| | DeepCubeA$_g$ | **52.02** | **100.00%** | **100.0%** | 1.81E+05 | 2.61 | 6.94E+04 |
| | FastDown (GC) | 36.75 | 0.80% | 0.80% | 9.05E+07 | 102.11 | 8.66E+05 |
| | FastDown (FF) | 52.75 | 80.80% | 24.80% | 2.92E+06 | 42.11 | 6.93E+04 |
| | FastDown (CG) | 41.95 | 4.40% | 1.20% | 2.00E+07 | 80.58 | 2.47E+05 |
| 15-P (Rand) | DeepCubeA$_g$ | **33.98** | **100.00%** | - | **1.11E+05** | **1.60** | **6.16E+04** |
| | FastDown (GC) | 14.92 | 38.00% | - | 1.61E+07 | 18.77 | 5.46E+05 |
| | FastDown (FF) | 32.66 | 89.20% | - | 1.24E+06 | 17.39 | 5.65E+04 |
| | FastDown (CG) | 20.45 | 51.20% | - | 3.90E+06 | 21.41 | 1.20E+05 |
| 24-P (Canon) | PDBs | **89.41** | **100.00%** | **100.00%** | 8.19E+10 | 4239.54 | **1.91E+07** |
| | DeepCubeA | 89.49 | **100.00%** | 96.98% | 6.44E+06 | 19.33 | 3.34E+05 |
| | DeepCubeA$_g$ | 90.47 | **100.00%** | 55.24% | **3.38E+05** | **5.22** | 6.48E+04 |
| | FastDown (GC) | - | 0.00% | 0.00% | - | - | - |
| | FastDown (FF) | 81.00 | 1.01% | 0.40% | 2.68E+06 | 89.84 | 2.91E+04 |
| | FastDown (CG) | - | 0.00% | 0.00% | - | - | - |
| 24-P (Rand) | DeepCubeA$_g$ | 66.28 | **99.60%** | - | 3.10E+05 | 4.91 | 6.16E+04 |
| | FastDown (GC) | 9.86 | 10.00% | - | 9.54E+06 | 11.88 | 4.27E+05 |
| | FastDown (FF) | 26.35 | 26.00% | - | 5.99E+05 | 19.57 | 2.41E+04 |
| | FastDown (CG) | 13.75 | 12.60% | - | 1.42E+06 | 14.42 | 6.85E+04 |
| Sokoban | DeepCubeA | 32.88 | **100.00%** | - | **5.01E+03** | 2.71 | 1.84E+03 |
| | DeepCubeA$_g$ | **32.02** | **100.00%** | - | 1.80E+04 | 0.95 | 1.79E+04 |
| | FastDown (GC) | 31.94 | 99.80% | - | 3.17E+06 | 5.93 | 5.85E+05 |
| | FastDown (FF) | 33.15 | **100.00%** | - | 2.92E+04 | **0.32** | **7.49E+04** |
| | FastDown (CG) | 33.12 | **100.00%** | - | 4.43E+04 | 0.51 | 7.25E+04 |

# ASP Specifications

- We build on this using answer set programming to describe goals with first-order logic and use answer set solvers to solve for assignments that make these goals true

- For the Rubik's cube
  - Define basic background knowledge
    - Colors, faces, cubelets
    - Constraints: Cannot have two stickers of the same color on the same cubelet, cannot have two stickers from the same cubelet on opposite faces
  - Given basic background knowledge, specifications often only require a few lines of code
    - `face_same(F) :- face_col(F, FCol), #count{Cbl : onface(Cbl, FCol, F)}=9.`
    - `canon_solved :- #count{F : face_same(F)}=6.`
  - Our specifications contain combinations of common patterns
    - Note: the training procedure is unaware of what the specification will be at test time



(a) Cross    (b) X    (c) Cup    (d) Spot

# Reaching Goals



- If our specification behaves monotonically, then all candidate states are goal states
  - Therefore, we can randomly sample assignments from $\Pi$ until we find one that we can reach

- Some of these assignments may represent the empty set

- The answer set solver (we use clingo) used is agnostic to the cost of a shortest path

# Results

| Goal | Path Cost | % Solved | # Models | Model Time | Search Time |
|------|-----------|----------|----------|------------|-------------|
| Rubik's Cube (Canon) | 24.41 | 100% | 1 | 0.?? | ?.?? |
| Rubik's Cube (Cross6) | 13.11 | 100% | 1 | 0.? | |
| Rubik's Cube (Cup4) | 24.33 | 100% | 42.? | 3 | |
| Rubik's Cube (CupSpot) | 17.99 | 100% | 27.68 | 38.66 | 241.08 |
| Rubik's Cube (Checkers) | 23.85 | 100% | 1 | 0.49 | 4.2 |
| Sokoban (Immov) | 35.15 | 100% | 6.37 | 6.83 | 16.16 |
| ?ban | 3.77 | 88% | 1.89 | 0.58 | 6.08 |
| ?ban | 4.42 | 77% | 1.26 | 0.38 | 4.09 |

(a) Example 1    (b) Example 2

**Cross6**

**Cup4**

**CupSpot**

**Checkers**

Agostinelli, Forest, Rojina Panta, and Vedant Khandelwal. "Specifying Goals to Deep Neural Networks with Answer Set Programming." *ICAPS 2024*

# Results



All boxes are immoveable

| Goal | Path Cost | % Solved | # Models | Model Time | Search Time |
|------|-----------|----------|----------|------------|-------------|
| Rubik's Cube (Canon) | 24.41 | 100% | 1 | 0.37 | 4.39 |
| Rubik's Cube (Cross6) | 13.11 | 100% | 1 | 0.41 | 2.14 |
| Rubik's Cube (Cup4) | 24.33 | 100% | 42.5 | 34.65 | 374.11 |
| Rubik's Cube (CupSpot) | 17.99 | 100% | 27.68 | 38.66 | 241.08 |
| Rubik's Cube (Checkers) | 23.85 | 100% | 1 | 0.49 | 4.2 |
| Sokoban (Immov) | 35.15 | 100% | 6.37 | 6.83 | 16.16 |
| Sokoban (BoxBox) | 33.77 | 88% | 1.89 | 0.58 | 6.08 |
| Sokoban (AgentInBox) | 34.42 | 77% | 1.26 | 0.38 | 4.09 |



A box of boxes

Boxes at the four corners of the agent

Agostinelli, Forest, Rojina Panta, and Vedant Khandelwal. "Specifying Goals to Deep Neural Networks with Answer Set Programming." *ICAPS 2024*

# Handling Non-Monotonicity

- If negation as failure is used in a program, $\Pi$, then $\Pi$ can exhibit non-monotonic behavior
    - A logic program is non-monotonic if some atoms that were previously derived can be retracted by adding new knowledge
    - Therefore, we can have a state that is a candidate state but not a goal state
- For example, a white cross with no yellow stickers on the white face
    - The assignment for this specification is just a white cross
    - However, there can be a state that is a specialization of this assignment, but has yellow on the white face
- To address this, we use a conflict-driven approach that specializes assignments based on why a state is not a goal state

# Handling Non-Monotonicity

$\Pi$: Answer set program

$\mathcal{S}_\Pi$: set of states represented by program

$\mathcal{S}_A$: set of states represented by assignment

| Goal | SpecOp | Cost | %Solve | #Itr | #Assign | %reach | %not goal | $\frac{\text{Secs}}{\text{Spec}}$ | $\frac{\text{Secs}}{\text{Path}}$ | Secs |
|---|---|---|---|---|---|---|---|---|---|---|
| RC:∀diffCtrW | - | 11.54 | 70 | **3.34** | 33.43 | 7.68 | **0** | 12.77 | 7.5 | 564.94 |
| RC:¬∃sameCtrW | Rand | 1.67 | 99 | 7.2 | | 87.84 | | **0.06** | | 95.46 |
| | Conflict | **1.26** | **100** | 5.43 | | **99.34** | | **0.06** | | **5.98** |
| 24p:r0SumEven | - | 24.55 | **100** | 9.24 | | **100** | | **0.2** | | 42.52 |
| 24p:¬r0SumOdd | Rand | 3.16 | **100** | 4.27 | | **100** | | **0.2** | | 6.64 |
| | Conflict | **2.51** | **100** | **4.06** | **31.6** | **100** | 22.13 | 0.21 | 0.04 | **6.58** |
| 24p:∀rSumEven | - | 83.71 | **100** | 9.19 | 91.9 | 50.41 | **0** | 0.88 | 1.77 | 250.18 |
| 24p:¬∃rSumOdd | Rand | 17.07 | **100** | 10.23 | 92.05 | 99.98 | 85.51 | **0.1** | **0.08** | 21.72 |
| | Conflict | **12.87** | **100** | **8.66** | **77.1** | **100** | 79.72 | 0.11 | **0.08** | **17.08** |

**All stickers on the white face are different than the center sticker**



Start     Mono: path cost 12     Non-mono: path cost 1

**All rows sum to an even number**



Start     Mono: path cost 93     Non-mono: path cost 4

Agostinelli, Forest. "A Conflict-Driven Approach for Reaching Goals Specified with Negation as Failure." *ICAPS 2024 HAXP Workshop*

# Outline

- Background
- Generalizing over states
- Generalizing over goals
- Generalizing over domains
- Towards obtaining approximately admissible heuristic functions
- Generalizing to domains with unknown transition functions

# Example



Start state        Goal state

- If using only canonical actions, the cost-to-go is 16

- If including diagonal actions, the cost-to-go is 2

- To differentiate between these two scenarios, information about the domain must also be given to the heuristic function

# Training

- For each example, randomly sample a domain

- For that domain, randomly sample a state

- RL Update

  - $L(\theta) = \left( \min_a (c^a(s) + h_{\theta-}(T(s,a),D)) - h_\theta(s,D) \right)^2$

    - $D$: Domain

# Preliminary Experiments

- For the 15-puzzle, generate different domains by sampling a subset of {U, D, L, R, UL, UR, DL, DR} actions for each tile position
  - 8 actions for each of the 16 positions, max $2^{8*16} \approx 3.4 \times 10^{38}$ domains
  - Ensure all sampled domains are reversible, for simplicity
- Represent the domain as a one-hot vector of which actions are allowed in each position
- Compare heuristic performance with true cost-to-go for random states from domains
  - True cost-to-go computed with merge-and-shrink heuristic
- Compare when training a heuristic function across domains without domain information
- Compare heuristic function with DeepCubeA trained for a fixed domain

# Results

- Adding action information significantly improves performance

- Performs slightly worse when compared to DeepCubeA trained on that specific domain
  - However, unlike DeepCubeA, it does not need to be re-trained for that domain



(a) C: P vs GT

(b) D: P vs GT

(c) C+D: P vs GT

(d) C: DCA vs GT

(e) D: DCA vs GT

(f) C+D: DCA vs GT

(a) Without Action Info

(b) With Action Info

Khandelwal, Vedant, Amit Sheth, Forest Agostinelli. "Towards Learning Foundation Models for Heuristic Functions to Solve Pathfinding Problems." *arxiv, 2024*

# Results

- Repeat training for 8-puzzle and 24-puzzle
- Proposed approach compares favorably to the fast downward planner with the fast forward heuristic
- Is slightly worse than DeepCubeA, which must be re-trained for each domain
- Future work could build on work by Felipe Trevizan and Sylvie Thiebaux on using graph neural networks to encode PDDL domains

| Domain | Solver | Len | Opt | Nodes | Secs | Nodes/Sec | Solved |
|---|---|---|---|---|---|---|---|
| 8 Puzzle (C) | DeepCubeA | 18.38 | 100% | 3.59E+04 | 0.69 | 5.2E+04 | 100% |
| 8 Puzzle (C) | Proposed | 18.38 | 100% | 7.17E+04 | 1.76 | 4.07E+04 | 100% |
| 8 Puzzle (C) | FD(FF) | 18.8 | 81% | 5.56E+02 | 0.11 | 4.7E+03 | 100% |
| 8 Puzzle (D) | DeepCubeA | 1.44 | 100% | 1.95E+01 | 0.01 | 2.92E+03 | 100% |
| 8 Puzzle (D) | Proposed | 1.44 | 100% | 4.05E+01 | 0.01 | 4.92E+03 | 100% |
| 8 Puzzle (D) | FD(FF) | 1.44 | 100% | 2.45E+00 | 0.2 | 1.23E+01 | 100% |
| 8 Puzzle (C+D) | DeepCubeA | 11.84 | 100% | 6.2E+04 | 1.18 | 5.26E+04 | 100% |
| 8 Puzzle (C+D) | Proposed | 11.84 | 100% | 6.23E+04 | 1.56 | 3.97E+04 | 100% |
| 8 Puzzle (C+D) | FD(FF) | 12.9 | 54.2% | 8.68E+01 | 0.13 | 6.59E+02 | 100% |
| 15 Puzzle (C) | DeepCubeA | 52.03 | 99.4% | 1.82E+05 | 4.31 | 4.21E+04 | 100% |
| 15 Puzzle (C) | Proposed | 52.18 | 93.76% | 3.62E+05 | 10.39 | 3.49E+04 | 100% |
| 15 Puzzle (C) | FD(FF) | 52.75 | 24.80 | 2.92E+06 | 42.11 | 6.93E+04 | 80.80% |
| 15 Puzzle (D) | DeepCubeA | 10.8 | 100% | 8.2E+02 | 0.03 | 2.43E+04 | 100% |
| 15 Puzzle (D) | Proposed | 10.81 | 99.8% | 1.64E+03 | 0.05 | 3.01E+04 | 100% |
| 15 Puzzle (D) | FD(FF) | 10.86 | 96.8% | 4.18E+01 | 0.21 | 1.96E+02 | 100% |
| 15 Puzzle (C+D) | DeepCubeA | 25.66 | 100% | 1.78E+05 | 3.74 | 4.78E+04 | 100% |
| 15 Puzzle (C+D) | Proposed | 25.67 | 99.8% | 1.78E+05 | 4.72 | 3.78E+04 | 100% |
| 15 Puzzle (C+D) | FD(FF) | 29.32 | 13.4% | 8.4E+03 | 1.17 | 3.56E+03 | 100% |
| 24 Puzzle (C) | DeepCubeA | 89.48 | 96.98% | 3.34E+05 | 8.05 | 4.15E+04 | 100% |
| 24 Puzzle (C) | Proposed | 92.80 | 22.03% | 7.6E+05 | 24.06 | 3.16E+04 | 100% |
| 24 Puzzle (C) | FD(FF) | 81.00 | 0.40 | 2.68E+06 | 89.84 | 2.91E+04 | 1.01% |
| 24 Puzzle (D) | DeepCubeA | 14.9 | 100% | 2.55E+04 | 0.47 | 5.46E+04 | 100% |
| 24 Puzzle (D) | Proposed | 14.92 | 99.8% | 5.1E+04 | 1.35 | 3.78E+04 | 100% |
| 24 Puzzle (D) | FD(FF) | 15.16 | 89.2% | 2.64E+02 | 0.12 | 2.05E+03 | 100% |
| 24 Puzzle (C+D) | DeepCubeA | 31.33 | 100% | 2.27E+05 | 4.83 | 4.69E+04 | 100% |
| 24 Puzzle (C+D) | Proposed | 31.34 | 99.6% | 2.27E+05 | 6.78 | 3.34E+04 | 100% |
| 24 Puzzle (C+D) | FD(FF) | 36.81 | 13.8% | 1.7E+04 | 5.35 | 1.77E+03 | 99.4% |

Khandelwal, Vedant, Amit Sheth, Forest Agostinelli. "Towards Learning Foundation Models for Heuristic Functions to Solve Pathfinding Problems." *arxiv, 2024*
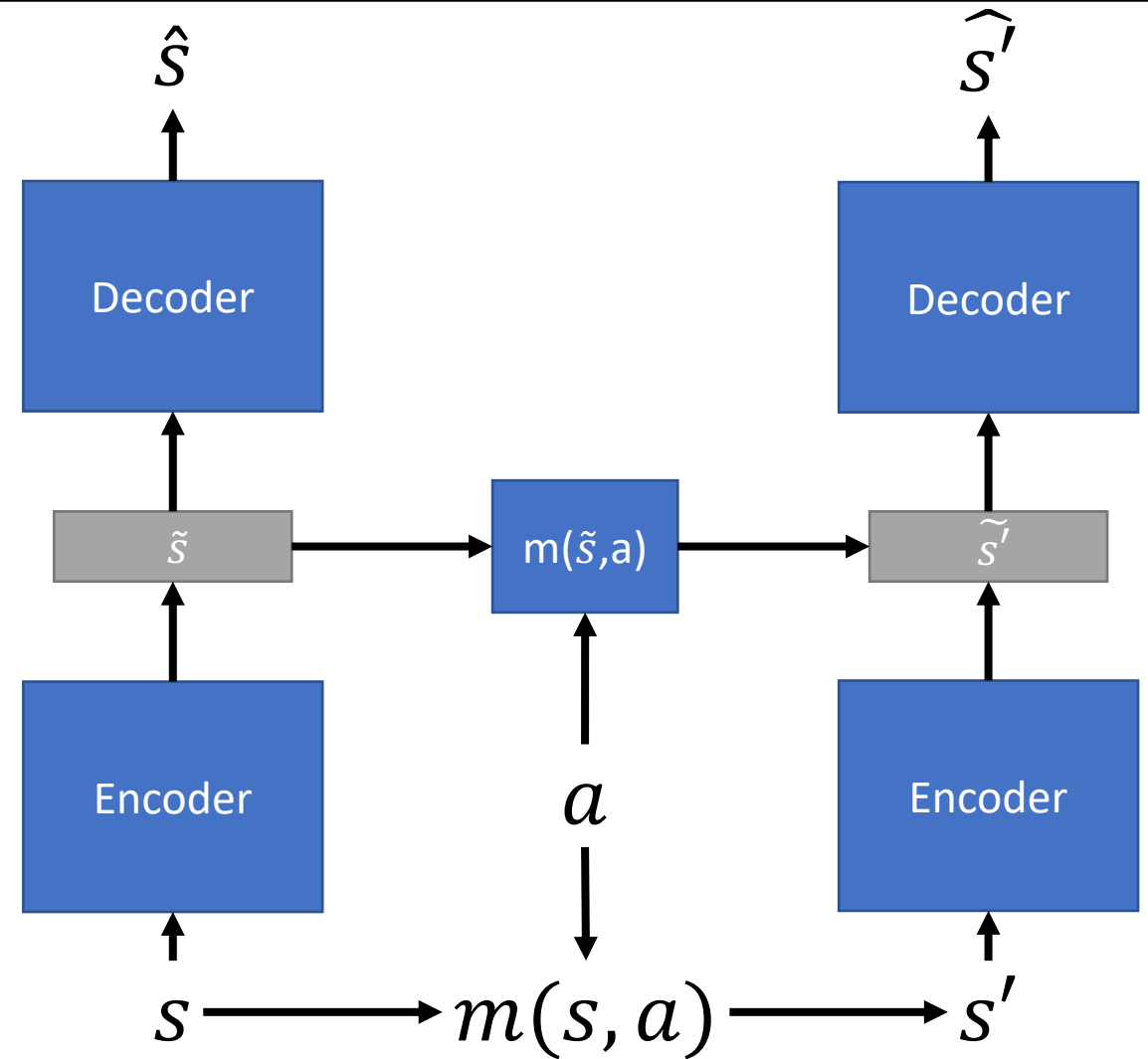Chen, Dillon Z., Sylvie Thiébaux, and Felipe Trevizan. "Learning Domain-Independent Heuristics for Grounded and Lifted Planning." *AAAI*. Vol. 38. No. 18. 2024.

# Outline

- Background
- Generalizing over states
- Generalizing over goals
- Generalizing over domains
- Towards obtaining approximately admissible heuristic functions
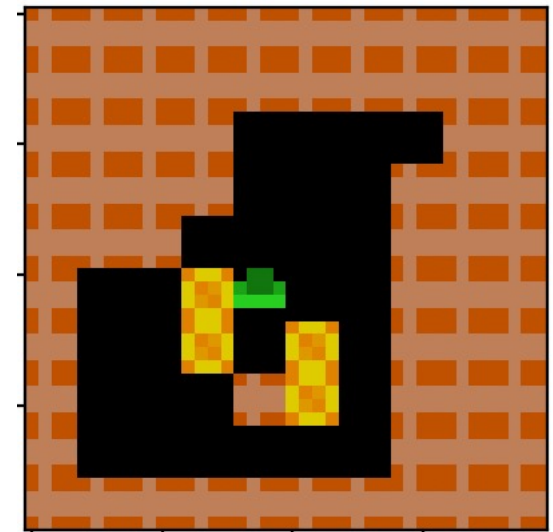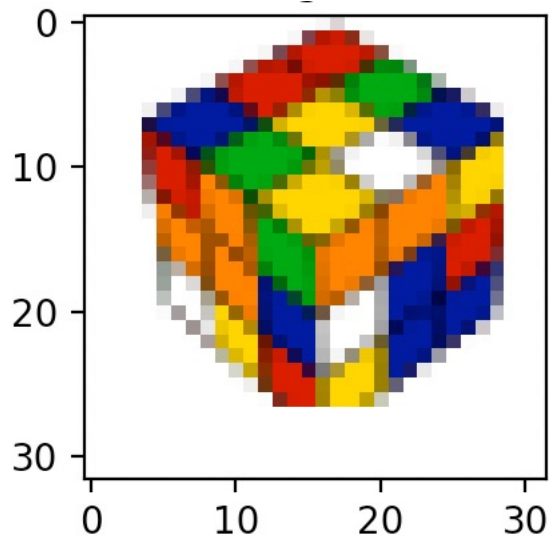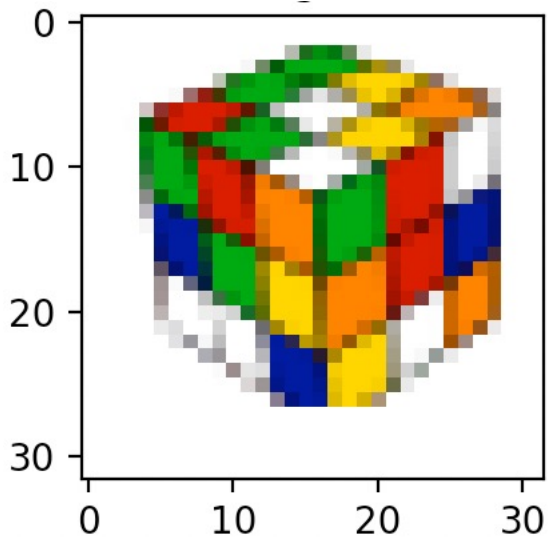- Generalizing to domains with unknown transition functions

# Approach

- When performing A* search with an admissible heuristic function, every node popped from OPEN is a lower bound on the cost-to-go

- We perform A* search with an admissible heuristic on a representative set of states for a given domain to get their lower bounds
  - Start with a trivially admissible heuristic function of always zero

- We can then use these lower bounds to correct an inadmissible heuristic function based on its maximum overestimation

- We can then repeat this process with the adjusted inadmissible heuristic function to get an improved estimation of the lower bound

- The larger the representative set of states, the more accurate the adjustment process will be

**Algorithm 1:** Approximately Admissible Conversion

**Input:**
- $h$: Inadmissible heuristic function
- $\mathcal{X}$: Representative set
- $\eta$: target increment for $h^a$

**Output:**
- $h'$: Converted approximately admissible heuristic function

$h^a(x) \leftarrow 0, \forall x \in \mathcal{X}$
$is\_solved(x) \leftarrow False, \forall x \in \mathcal{X}$
**while** $\exists x \in \mathcal{X} \mid is\_solved(x) == False$ **do**
$\quad h' \leftarrow adjust(h^a, h, \mathcal{X})$
$\quad$ **for** $x \in \mathcal{X}$ **do**
$\quad\quad h^a(x), is\_solved(x) \leftarrow A^*(x, h', h^a(x) + \eta)$

$h' \leftarrow adjust(h^a, h, \mathcal{X})$
**Return** $h'$

# Results

- For the 15-puzzle
  - Before adjustment: 71.37% overestimation,  max overestimation: 8.28
  - After adjustment: 0.0019% overestimation,  max overestimation: 0.62

- The larger the representative set, the better the adjustment



(a) Before approximately admissible conversion

(b) After approximately admissible conversion



(a) Max overestimation

(b) Percent inadmissible

(c) Average heuristic value

Agostinelli, Forest, et al. "Obtaining approximately admissible heuristic functions through deep reinforcement learning and A* search." *ICAPS PRL Workshop 2021.*

# Outline

- Background
- Generalizing over states
- Generalizing over goals
- Generalizing over domains
- Towards obtaining approximately admissible heuristic functions
- Generalizing to domains with unknown transition functions

# Learning Discrete World Models

- Addressing previous shortcomings
  - Small errors in prediction can be corrected by simply rounding
  - Can reidentify states by comparing two vectors
- Encoder
  - Maps the state to a discrete representation
  - To allow training with gradient descent, use a straight through estimator
- Decoder
  - Maps the discrete representation to the state
  - Ensures the discrete representation is meaningful
- Environment model
  - Maps discrete states and actions to next discrete state

# Experiments

- Rubik's cube
  - Two 32x32 RGB images showing both sides of the cube
- Sokoban
  - One 40x40 RGB image
- Generate offline dataset of 300,000 episodes of 30 random steps, each

# Discrete vs Continuous Model Performance

- The continuous model eventually accumulates error for the Rubik's cube



(a) Rubik's Cube                     (b) Sokoban

# Discrete vs Continuous Model Performance
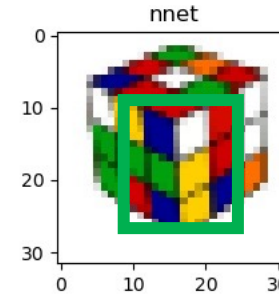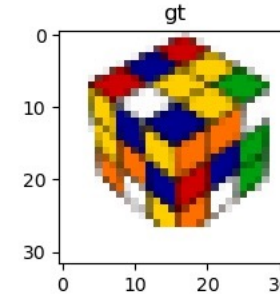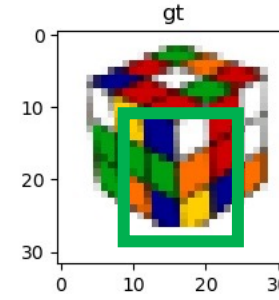
# Heuristic Learning and Search with Discrete Model

- DeepCubeAI – DeepCubeA + "Imagination"
  - Learn discrete world model with offline data
  - Use offline data and the learned world model to generate training data
  - Heuristic learning: Q-learning with hindsight experience replay
    - Generalize over goal states
  - Heuristic search: Q* search
    - Helps when model uses computationally expensive DNN

| Domain | Solver | Len | Opt | Nodes | Secs | Nodes/Sec | Solved |
|---|---|---|---|---|---|---|---|
| RC | PDBs$^+$ | 20.67 | 100.0% | 2.05E+06 | 2.20 | 1.79E+06 | 100% |
| | DeepCubeA | 21.50 | 60.3% | 6.62E+06 | 24.22 | 2.90E+05 | 100% |
| | Greedy (ours) | - | 0% | - | - | - | 0% |
| | DeepCubeAI (ours) | 22.85 | 19.5% | 2.00E+05 | 6.21 | 3.22E+04 | 100% |
| RC$_{rev}$ | Greedy (ours) | - | 0% | - | - | - | 0% |
| | DeepCubeAI (ours) | 22.81 | 21.92% | 2.00E+05 | 6.30 | 3.18+04 | 99.9% |
| Sokoban | LevinTS | 39.80 | - | 6.60E+03 | - | - | 100% |
| | LevinTS (*) | 39.50 | - | 5.03E+03 | - | - | 100% |
| | LAMA | 51.60 | - | 3.15E+03 | - | - | 100% |
| | DeepCubeA | 32.88 | - | 1.05E+03 | 2.35 | 5.60E+01 | 100% |
| | Greedy (ours) | 29.55 | - | - | 1.68 | - | 41.9% |
| | DeepCubeAI (ours) | 33.12 | - | 3.30E+03 | 2.62 | 1.38E+03 | 100% |

# Questions?

- Papers
  - Agostinelli, Forest, et al. "Solving the Rubik's cube with deep reinforcement learning and search." Nature Machine Intelligence 1.8 (2019): 356-363.
  - Agostinelli, Forest, Rojina Panta, and Vedant Khandelwal. "Specifying Goals to Deep Neural Networks with Answer Set Programming." *ICAPS 2024*
  - Agostinelli, Forest and Soltani, Misagh "Learning Discrete World Models for Heuristic Search." *Reinforcement Learning Conference 2024*
  - Khandelwal, Vedant, Amit Sheth, Forest Agostinelli. "Towards Learning Foundation Models for Heuristic Functions to Solve Pathfinding Problems." *arxiv, 2024*
  - Agostinelli, Forest, et al. "Obtaining approximately admissible heuristic functions through deep reinforcement learning and A* search." *ICAPS PRL Workshop 2021.*
  - Agostinelli, Forest. "A Conflict-Driven Approach for Reaching Goals Specified with Negation as Failure." *ICAPS 2024 HAXP Workshop*

- Code
  - Many of these algorithms are publicly available on GitHub
  - https://github.com/forestagostinelli/deepxube

Email: foresta@cse.sc.edu

Website: https://cse.sc.edu/~foresta/