# Hierarchical Assignment of Behaviours by Self-Organizing

### W. Moerman<sup>1</sup> B. Bakker<sup>2</sup> M. Wiering<sup>3</sup>

<sup>1</sup>M.Sc. Cognitive Artificial Intelligence Utrecht University

<sup>2</sup>Intelligent Autonomous Systems Group University of Amsterdam

> <sup>3</sup>Intelligent Systems Group Utrecht University

Neural Information Processing Systems 2007 Workshop Hierarchical Organization of Behaviour: Computational, Psychological and Neural Perspectives

(B)

**NIPS 2007** 

1/18

### Summary of Ideas

We propose:

- Shifting the design burden:
  - from task decomposition
  - to a suitable abstract representation of the state space

• • • • • • • • • • • • •

## Summary of Ideas

We propose:

- Shifting the design burden:
  - from task decomposition
  - to a suitable abstract representation of the state space

• Using self-organization to figure out which behaviours are needed

- starting with uncommitted policies
- learning (parts of) its hierarchical structure
- no designed or fixed pre/post conditions

• • • • • • • • • • • •

## Summary of Ideas

We propose:

- Shifting the design burden:
  - from task decomposition
  - to a suitable abstract representation of the state space

• Using self-organization to figure out which behaviours are needed

- starting with uncommitted policies
- learning (parts of) its hierarchical structure
- no designed or fixed pre/post conditions

### The algorithm is called HABS:

Hierachical Assignment of Behaviours by Self-organization

• • • • • • • • • • • •

### Outline

### Introduction

- Hierarchical Reinforcement Learning
- Abstractions

### Our Algorithm (HABS)

- High Level Policy and Subpolicies
- Self-Organizing Behaviours
- Relating HABS to Other Work

## 3 Experiments

- Setup
- Results

Using behaviours (temporally extended/high level actions, ...) allows



extended actions

• • • • • • • • • • • •

Using behaviours (temporally extended/high level actions, ...) allows

- "Divide and Conquer": decompose into smaller (easier) subtasks
  - task decomposition enables re-use of (sub)policies



Moerman, Bakker, Wiering (UU, UvA)

• • • • • • • • • • • •

Using behaviours (temporally extended/high level actions, ...) allows

- "Divide and Conquer": decompose into smaller (easier) subtasks
   task decomposition enables re-use of (sub)policies
- "Dæmon of Dimensionality": smaller state spaces on all levels



Using behaviours (temporally extended/high level actions, ...) allows

- "Divide and Conquer": decompose into smaller (easier) subtasks
   task decomposition enables re-use of (sub)policies
- "Dæmon of Dimensionality": smaller state spaces on all levels
- Different specific state abstractions for different (sub)policies (just wait until tomorow)



Using behaviours (temporally extended/high level actions, ...) allows

- "Divide and Conquer": decompose into smaller (easier) subtasks
   task decomposition enables re-use of (sub)policies
- "Dæmon of Dimensionality": smaller state spaces on all levels
- Different specific state abstractions for different (sub)policies (just wait until tomorow)
- Faster exploration



- Reinforcement Learning Description is random walk
- but behaviours do something consistent (hopefully)
  - they move agent non-randomly through state space



- Reinforcement Learning Description is random walk
- but behaviours do something consistent (hopefully)
  - they move agent non-randomly through state space



less random choices

A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

- Reinforcement Learning Description is random walk
- but behaviours do something consistent (hopefully)
  - they move agent non-randomly through state space



 less random choices
 more distance covered in random walk

- Reinforcement Learning Description is random walk
- but behaviours do something consistent (hopefully)
  - they move agent non-randomly through state space



- less random choices
- more distance covered in random walk
- faster exploration

A (1) > A (1) > A

- Reinforcement Learning Description is random walk
- but behaviours do something consistent (hopefully)
  - they move agent non-randomly through state space



- less random choices
- more distance covered in random walk
- faster exploration



Drunken Mans walk



Random Walks on street level

# **Abstract State Space**



A suitable Abstract State Space has these properties:

- it has an underlying "geometric" structure:
  - not constrained to "spatial" geometry
  - consistent mapping: points close together in state space should be near each other in abstract state space, and vice versa
- Abstract State Space significantly smaller than State Space

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

### **Behaviour Space**

- Behaviour Space: space of all possible transition vectors in the Abstract State Space
  - note: Abstract State Space treated as continuous
  - intuition: think of Behaviour Space as a sphere



## **Behaviour Space**

- Behaviour Space: space of all possible transition vectors in the Abstract State Space
  - note: Abstract State Space treated as continuous
  - intuition: think of Behaviour Space as a sphere



- Actually occuring transitions between abstract states need to be distributed non-uniformly in the Behaviour Space
  - behaviours (transitions) in abstract state space are vectors in Behaviour Space
  - can be characterized by a limited number of vectors if clustered





HABS (Hierachical Assignment of Behaviours by Self-organization)

- One (high level) Policy<sub>HL</sub> and a limited set of subpolicies
  - uses abstract state space

### HABS (Hierachical Assignment of Behaviours by Self-organization)

- One (high level) Policy<sub>HL</sub> and a limited set of subpolicies
   uses abstract state space
- Policy<sub>HL</sub> has subpolicies as its (extended) actions:



- 4 ∃ →

### HABS (Hierachical Assignment of Behaviours by Self-organization)

- One (high level) Policy<sub>HL</sub> and a limited set of subpolicies
  - uses abstract state space
- Policy<sub>HL</sub> has subpolicies as its (extended) actions:



subpolicies self-organize to cover required behaviours

### HABS (Hierachical Assignment of Behaviours by Self-organization)

- One (high level) Policy<sub>HL</sub> and a limited set of subpolicies
  - uses abstract state space
- Policy<sub>HL</sub> has subpolicies as its (extended) actions:



-∢ ∃ ▶

- subpolicies self-organize to cover required behaviours
- rewards received during a behaviour are accumulated and used for High Level Policy<sub>HL</sub> reward

#### HABS (Hierachical Assignment of Behaviours by Self-organization)

- One (high level) Policy<sub>HI</sub> and a limited set of subpolicies
  - uses abstract state space
- Policy<sub>H</sub> has subpolicies as its (extended) actions:



- subpolicies self-organize to cover required behaviours
- rewards received during a behaviour are accumulated and used for High Level Policy<sub>HL</sub> reward
- subpolicy rewarding is independent of overall task



- 4 ∃ →

### HABS (Hierachical Assignment of Behaviours by Self-organization)

- One (high level) Policy<sub>HL</sub> and a limited set of subpolicies
  - uses abstract state space
- Policy<sub>HL</sub> has subpolicies as its (extended) actions:



(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

- subpolicies self-organize to cover required behaviours
- rewards received during a behaviour are accumulated and used for High Level Policy<sub>HL</sub> reward
- subpolicy rewarding is independent of overall task

### Standard RL techniques (online, off policy) like Q-learning •

 $Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot Q(s_t, a_t) + \alpha \cdot (\mathit{reward}_t + \gamma \cdot \max_a Q(s_{t+1}, a))$ 

• Assume: actually occuring behaviours are clustered together

- use clustering algorithm
- assign subpolicy to cluster center (specialize)

• • • • • • • • • • • •

Assume: actually occuring behaviours are clustered together

- use clustering algorithm
- assign subpolicy to cluster center (specialize)
- Exploration: agent will stumble upon new abstract states
  - transitions to new states are behaviours

Assume: actually occuring behaviours are clustered together

- use clustering algorithm
- assign subpolicy to cluster center (specialize)
- Exploration: agent will stumble upon new abstract states
  - transitions to new states are behaviours
  - abstract states are nearby

Assume: actually occuring behaviours are clustered together

- use clustering algorithm
- assign subpolicy to cluster center (specialize)
- Exploration: agent will stumble upon new abstract states
  - transitions to new states are behaviours
  - abstract states are nearby
  - random walks on small distances are disproportionately better

Assume: actually occuring behaviours are clustered together

- use clustering algorithm
- assign subpolicy to cluster center (specialize)
- Exploration: agent will stumble upon new abstract states
  - transitions to new states are behaviours
  - abstract states are nearby
  - random walks on small distances are disproportionately better
  - will happen long before overall task is completed

• • • • • • • • • • • •

Assume: actually occuring behaviours are clustered together

- use clustering algorithm
- assign subpolicy to cluster center (specialize)
- Exploration: agent will stumble upon new abstract states
  - transitions to new states are behaviours
  - abstract states are nearby
  - random walks on small distances are disproportionately better
  - will happen long before overall task is completed

Agent can discover meaningful behaviours long before it has a chance to solve overall problem (due to Random Walks properties)

# **Clustering and Rewarding Subpolicies**

- Subpolicy terminates  $\iff$  new abstract state reached or timeout
- On subpolicy termination: compare actually executed behaviour to cluster center (characteristic behaviour) of terminated subpolicy
  - if closest match: move cluster center towards experience



# **Clustering and Rewarding Subpolicies**

- Subpolicy terminates  $\iff$  new abstract state reached or timeout
- On subpolicy termination: compare actually executed behaviour to cluster center (characteristic behaviour) of terminated subpolicy
  - if closest match: move cluster center towards experience



### Always train subpolicy using Reinforcement Learning

$\mathit{reward}_{\mathit{sub}} = \Big\{$	( 0	not terminated
	1	terminated: closest match
	κ <sub>r</sub>	terminated: another cluster center is closer
	κ <sub>f</sub>	terminated: timeout (failed to reach anything)
	•	

### repeat

// run HL-Policy

Policy<sub>HL</sub> selects SubPolicy SUB<sub>i</sub>;

// HL-action

update *Policy<sub>HL</sub>* with *reward<sub>HL</sub>*; until task solved or timeout<sub>HL</sub> // for executing SUB i

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Moerman, Bakker, Wiering (UU, UvA) Self-Organizing Hierarchical Behaviours

### repeat

```
Policy<sub>HL</sub> selects SubPolicy SUB<sub>i</sub>;

repeat // execute SUB i

SUB<sub>i</sub> selects and executes a primitive action;

if new abstract state then BREAK; // behaviour=>terminate

else update SUB<sub>i</sub> with 0; // sparse reward

until timeout<sub>SUB</sub>
```

update *Policy<sub>HL</sub>* with *reward<sub>HL</sub>*; until task solved or timeout<sub>HL</sub>

```
repeat
    Policy_{HL} selects SubPolicy SUB_i;
    repeat
        SUB<sub>i</sub> selects and executes a primitive action :
        if new abstract state then BREAK :
        else update SUB<sub>i</sub> with 0;
    until timeoutsur
    if timeout<sub>SUB</sub> then punish SUB<sub>i</sub> :
                                        // no new abs. state
    else
                                   // compare EXECuted with clusters
        if EXEC ∈ CLUSTER<sub>SUB</sub> then
            reward SUB_i:
                                                                        // match
            move CLUSTER<sub>SUB</sub> towards EXEC;
                                                                        // match
        else punish SUB<sub>i</sub>;
                                                                   // no match
    update Policy_{HI} with reward<sub>HI</sub>;
until task solved or timeout<sub>HI</sub>
```

**NIPS 2007** 

11/18

```
repeat
    reward_{HI} = 0:
                                                // for accumulating rewards
    Policy_{HL} selects SubPolicy SUB_i;
    repeat
         SUB<sub>i</sub> selects and executes a primitive action :
         reward_{HI} \leftarrow reward_{HI} + receivedReward;
                                                                      // accumulate
        if new abstract state then BREAK :
        else update SUB<sub>i</sub> with 0;
    until timeoutsub
    if timeout<sub>SUB</sub> then punish SUB<sub>i</sub> ;
    else
        if EXEC ∈ CLUSTER<sub>SUB</sub> then
             reward SUB_i:
             move CLUSTER<sub>SUB</sub> towards EXEC;
        else punish SUB<sub>i</sub>;
    update Policy_{HI} with reward<sub>HI</sub>;
until task solved or timeout<sub>HI</sub>
```

**NIPS 2007** 

11/18

```
repeat
    reward_{HI} = 0:
    Policy_{HL} selects SubPolicy SUB_i;
    repeat
         SUB<sub>i</sub> selects and executes a primitive action :
         reward_{HI} \leftarrow reward_{HI} + receivedReward;
         if new abstract state then BREAK :
         else update SUB<sub>i</sub> with 0;
    until timeoutsub
    if timeout<sub>SUB</sub> then punish SUB<sub>i</sub>;
    else
         if EXEC ∈ CLUSTER<sub>SUB</sub> then
             reward SUB_i:
              move CLUSTER<sub>SUB</sub> towards EXEC;
         else punish SUB<sub>i</sub>;
    update Policy_{HI} with reward<sub>HI</sub>;
until task solved or timeout<sub>HI</sub>
```

## The Difference: Shifting the Burden

HABS differs from other hierarchical RL approaches:

- no focus on defining a task decomposition (MAXQ, HEXQ, HAM)
  - no need to define start and stop conditions
- starts with uncommitted subpolicies that self-organize

HABS shifts design burden from task decomposition to defining a suitable abstract representation

Like many hierarchical approaches, HABS depends on a certain structure ("geometry") in the State Space

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

## Experiment Description ("Cleaner")

### • Gridworld environment:

- actions: North, East, South, West, Pickup<sub>object</sub>, Drop<sub>object</sub>
- walls, drop areas and portable objects (max. 1 per cell)
- reward for each object dropped at drop area

• • • • • • • • • • • • •

# Experiment Description ("Cleaner")

### • Gridworld environment:

- actions: North, East, South, West, Pickup<sub>object</sub>, Drop<sub>object</sub>
- walls, drop areas and portable objects (max. 1 per cell)
- reward for each object dropped at drop area

### • "Cleaner":

- many objects
- agent can carry 10
- high level policy: multilayer neural network
- subpolicies: multilayer neural network



• • • • • • • • • • • •

# Experiment Description ("Cleaner")

### • Gridworld environment:

- actions: North, East, South, West, Pickup<sub>object</sub>, Drop<sub>object</sub>
- walls, drop areas and portable objects (max. 1 per cell)
- reward for each object dropped at drop area

### • "Cleaner":

- many objects
- agent can carry 10
- high level policy: multilayer neural network
- subpolicies: multilayer neural network •



• • • • • • • • • • • • •

### • Task has spatial and non-spatial aspects!

# State Space and Abstract State Spaces

### State Space: (subpolicies)

 Agent has a simulated "radar" → observes object/wall/drop areas (~ 100 inputs)



# State Space and Abstract State Spaces

State Space: (subpolicies)

 Agent has a simulated "radar" ● observes object/wall/drop areas (~ 100 inputs)

Abstract State Space: (high level policy)

- "cleaner" task uses < area<sub>agent</sub>, cargo > to determine subpolicy termination
  - no info about other objects!

< □ > < □ > < □ > < □ >

# State Space and Abstract State Spaces

State Space: (subpolicies)

 Agent has a simulated "radar" ● observes object/wall/drop areas (~ 100 inputs) Ø.

Abstract State Space: (high level policy)

- "cleaner" task uses < area<sub>agent</sub>, cargo > to determine subpolicy termination
  - no info about other objects!
  - objects moving/(dis)appearing without the agent is no behaviour: agent behaviour only defined by its position and cargo

• • • • • • • • • • • •

# State Space and Abstract State Spaces

State Space: (subpolicies)

 Agent has a simulated "radar" ● observes object/wall/drop areas (~ 100 inputs) Ð.

Abstract State Space: (high level policy)

- "cleaner" task uses < area<sub>agent</sub>, cargo > to determine subpolicy termination
  - no info about other objects!
  - objects moving/(dis)appearing without the agent is no behaviour: agent behaviour only defined by its position and cargo
  - but objects are important for behaviour selection:
     "high level radar" observations (wider and coarser than low level)

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

# State Space and Abstract State Spaces

State Space: (subpolicies)

 Agent has a simulated "radar" ● observes object/wall/drop areas (~ 100 inputs)

Abstract State Space: (high level policy)

- "cleaner" task uses < area<sub>agent</sub>, cargo > to determine subpolicy termination
  - no info about other objects!
  - objects moving/(dis)appearing without the agent is no behaviour: agent behaviour only defined by its position and cargo
  - but objects are important for behaviour selection:
     "high level radar" observations (wider and coarser than low level)
- Better to say: "radar" observation is Abstract State and only a subset of the Abstract States is used for subpolicy termination

#### Experiments

Results

## **Cleaner: Results**



### Boxplots:

### - HABS

- "flat" learner

 $m{lpha} = \{ {\color{red} 0.02}, {\color{black} 0.01} \}$ 

HABS: 5 hidden (Policy<sub>HL</sub> ) 2 hidden (subpol.)

"flat": 15 hidden

### • HABS is much faster and only slightly suboptimal

- "flat" has wide variance in convergence value
- "flat" has far wider variance in convergence time

Results

### **HABS** Demonstration

- HABS near convergence
- shows suboptimality: agent ignores objects in area I (pink, center)

Moerman, Bakker, Wiering (UU, UvA) Self-Organizing Hierarchical Behaviours

NIPS 2007 16 / 18

# **Conclusions / Future Work**

### o conclusions about HABS:

- possible to shift design burden from task decomposition to state space abstraction
- learns conditions for behaviours by self-organizing
- can start with unspecified behaviours (no fixed pre/post conditions)

### • future work

- try HABS with 3 or more layers
- behaviour representation: limited to vectors?

• • • • • • • • • • • •



Moerman, Bakker, Wiering (UU, UvA)

Self-Organizing Hierarchical Behaviours

NIPS 2007 18 / 18

(日)

# Action Space and Primitive Actions

Action Space: space of all possible transition vectors in the state space

• primitive actions are vectors in the *Action Space* (a subset)



• primitive actions are not distributed evenly but clustered together

- primitive actions are not distributed evenly but clustered together
- only one primitive action for North instead of many North<sub>1</sub>, North<sub>2</sub>, North<sub>3</sub>,... for going north from state 1 to 352, from state 2 to 369, from state 3 to 792345, ...
- relative (vectors), not absolute (North = "in state A goto B")

back to behaviour space

# **Behaviours Mirror Primitive Actions**

<b>Primitive Actions</b>	Behaviours
vectors in action space	vectors in behaviour space
relative to state	relative to abstract state
clustered	hopefully clustered
1 time step	1 high level time step
action successful	reach new abstract state
action fails	timeout

back to behaviour space

イロト イヨト イヨト イヨト

## Training Subpolicies, How?

- a subpolicy starts with no knowledge (i.e. randomly initated)
  - what is its desired or characteristic behaviour?
- train on pairs of abstract states?
  - designer needs to specify pre/post conditions
- rewards independent of the overall task (Policy<sub>HL</sub>)
  - ▶ behaviour "A ⇒ goal" same as "B ⇒ C"
  - ▶ blue behaviour has high Q<sub>HL</sub>-value in A but low Q<sub>HL</sub>-value in B
  - red behaviour has high Q<sub>HL</sub>-value in B



### not dependence on Policy<sub>HL</sub> on fixed pre/post conditions!

▶ back to Policy<sub>HL</sub>

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

## **Reinforcement Learning**

An agent:

- observes a state
- executes an action
- receives a reward

Based on this information, it needs to learn what actions to select in what situations – using an RL algorithm:

- future rewards need to be discounted
- stored in tabular form or function approximator

back to intro

A (1) > A (2) > A

# Q-Learning and Advantage Learning

Q-Learning:

$$\begin{array}{rcl} Q(s_t, a_t) & \leftarrow & (1 - \alpha) \cdot Q(s_t, a_t) & + \\ & & \alpha \cdot (\mathit{reward}_{t+1} + \gamma \cdot \max_a Q(s_{t+1}, a)) \end{array}$$

Advantage Learning (Baird):

$$\begin{array}{lcl} \mathcal{A}(s_t, a_t) & \leftarrow & (1 - \alpha) \cdot \mathcal{A}(s_t, a_t) & + \\ & & & \alpha \cdot \left( \begin{array}{c} \max_{a_t} \mathcal{A}(s_t, a_t) & + \\ & & \frac{\mathsf{reward}_{t+1} + \gamma \max_a \mathcal{A}(s_{t+1}, a) - \max_{a_t} \mathcal{A}(s_t, a_t)}{k} \end{array} \right) \end{array}$$

• • • • • • • • • • • •

**NIPS 2007** 

23/18

where  $\alpha$  is the learning rate, and *k* the scaling factor ( $0 < k \le 1$ ). With k = 1 this equation reduces to Q-Learning

back to intro back to Proposed Algorithm back to Subpolicies

### Formula for Clustering

- used euclidean distance for determining closest cluster center
- if subpolicy was the winner, move cluster center:

 $char_{t+\Delta t} = (1 - \alpha) \cdot char_t + \alpha \cdot act_{t \to t+\Delta t}$ 

where  $act_{t \to t+\Delta t}$  is the actually executed behaviour,  $\Delta t$  the time it took to execute the subpolicy, and *char<sub>t</sub>* the characteristic behaviour vector for a subpolicy at time *t* 

- otherwise: do nothing with the clustering
  - subpolicy executed a behaviour that is better matched by another subpolicy

back to clustering

(ロ) (同) (E) (E) (E) (O)

## Details on Self Organizing

- When agent reaches new abstract state it experiences a behaviour
  - center of closest cluster moved towards the newly experienced behaviour
  - this is the clustering part
- characteristic behaviour for a subpolicy is represented by cluster center
  - reward subpolicy when actually executed behaviour closest to its own cluster center
  - forced "outward" by punishing "staying in the same abstract state"
  - the self-organizing part

• • • • • • • • • • • •

# Formulae for Rewarding Behaviours 1

At every time step during the execution of a subpolicy, the normal RL (Q-Learning, Advantage Learning, ...) update is applied:

$$Q_L(s_t, a_t) \leftarrow (1 - \alpha) \cdot Q_L(s_t, a_t) + \alpha \cdot (r_{t+1} + \gamma \cdot \max_a Q_L(s_{t+1}, a))$$

 $Q_{L}$  is the Q-Value for a (low level) subpolicy,  $\alpha$  and  $\gamma$  are the learning rate and discount as usual

### Sensors

- agent has radar-like sensorgrid (8 areas per ring)
  - details nearby, course observations far away
  - trade off between detail and amount of sensor data



observation: vector of area densities

- $\blacktriangleright \langle \frac{1}{28}, \frac{2}{24}, \frac{3}{28}, \frac{4}{24}, \dots, \frac{3}{8}, 0, 0, \frac{1}{6}, \dots, \frac{2}{3}, 0, \frac{1}{3}, 1, \dots, 1, 1, 0, 1, \dots \rangle$
- for walls / objects / drop areas, so 3 × 32 inputs

extra values coding for position, cargo, etc were added

back to state space

# Neural Nets for Subpolicies 1

### Linear neural network:



Multilayer Neural network:



Both are "feed forward" networks (no recurrence)

back to experiment description

Appendix

## Neural Nets for Subpolicies 2



- each subpolicy has 6 neural nets, one (shown on prev. slide) for each action a<sub>i</sub>, giving Q(s, a<sub>i</sub>) in state s
  - allows for different inputs for different actions (not tried here)
- advantage of seperate networks:
  - no interference between actions
  - less hidden neurons needed for each network
  - faster backpropagation (because only one action is updated)

back to experiment description

• • • • • • • • • • • •

## **Comparing with Flat Learners**

A "flat" reinforcement learning agent was used for comparison

- first experiment (maze): tabular
  - each observation is unique (because position is included)
  - more efficient storage in table (only store position)
- second experiment (cleaner): neural network
  - comparable to what the high level policy used
  - also tried with more hidden neurons

used best performance settings for the "flat" learner (took lots and lots of time to find)

## **Experiment: Maze**

### gridworld environment:

- actions: North, East, South, West, Pickup<sub>object</sub>, Drop<sub>object</sub>
- walls, drop areas and portable objects (max. 1 per cell)
- reward for each object dropped at drop area



- first experiment (Maze):
  - big  $(39 \times 36 \approx 1.4 \cdot 10^3 \text{ cells})$
  - only one object
  - high level policy tabular
  - subpolicies: linear neural network D

back to Cleaner Description

📔 🕨 back to Cleaner Resu

Appendix

# First Experiment (Maze): Results

- HABS compared with "flat" learner (best performance, tabular)
- for HABS, only coarse search was done, but no extensive fine tuning!



(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

- HABS is order of magnitude faster but sub optimal
- memory usage:
  - ▶ flat learner needs to store  $10^7$  Q-values ( $\approx 100$  megabyte)
  - ► HABS needs  $5 \cdot 10^3 \times numberOfSubpolicies$  ( $\approx 1$  megabyte)
  - Neural network storage is neglectable

back to Cleaner Description back to Cleaner Results