

# 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>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

# Summary of Ideas

We propose:

- A new **Hierarchical Reinforcement Learning** algorithm:
  - ▶ Shifting the focus to **state abstraction**
  - ▶ Using **self-organization** to learn behaviours

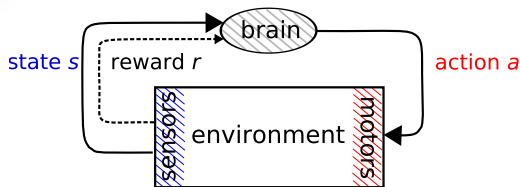
The proposed algorithm is called HABS:  
**Hierarchical Assignment of Behaviours by Self-organization**

# Outline

- 1 Introduction
  - Reinforcement Learning
  - Hierarchical Reinforcement Learning
- 2 Our Algorithm (HABS)
  - Focus on State Abstraction
  - Searching For Suitable Abstractions
  - Self-Organizing the Behaviours
- 3 Experiments
  - Setup
  - Results

# Agent–Environment–Interaction

- **Feedback** between agent and environment:



- **Agent:**

- ▶ **brain** (table or function approximator)
- ▶ **sensors**
- ▶ **motor controls**

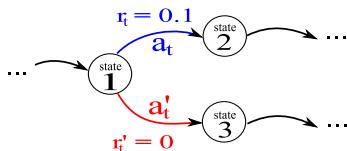
- **Environment:**

- ▶ **rewards**
- ▶ **states** (sensed by sensors)
- ▶ **actions** (executed by motor controls)

- Agent wants to **maximize the total reward** it receives

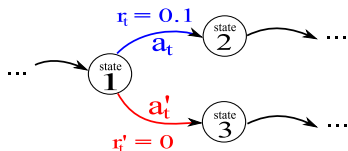
# Reinforcement Learning is “Trial-and-Error”

- Suppose we define “goodness” values based on the rewards
  - ▶ *being in* a state:  $V(\text{state}_2) = 0.1$  and  $V(\text{state}_3) = 0$
  - ▶ *or: going to* a state:  $Q(\text{state}_1, a_t) = 0.1$  and  $Q(\text{state}_1, a'_t) = 0$



# Reinforcement Learning is “Trial-and-Error”

- Suppose we define “goodness” values based on the rewards
  - ▶ *being in* a state:  $V(\text{state}_2) = 0.1$  and  $V(\text{state}_3) = 0$
  - ▶ or: *going to* a state:  $Q(\text{state}_1, a_t) = 0.1$  and  $Q(\text{state}_1, a_{t'}) = 0$



- Every time the agent visits a state, **update** the values:
  - ▶ move  $Q$  or  $V$  towards *reward*:

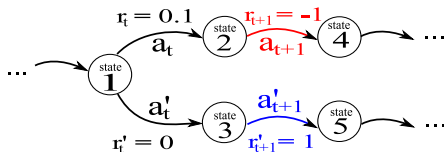
$$V(\text{state}) \leftarrow (1 - \alpha) \cdot V(\text{state}) + \alpha \cdot \text{reward}$$

or

$$Q(\text{state}, \text{action}) \leftarrow (1 - \alpha) \cdot Q(\text{state}, \text{action}) + \alpha \cdot \text{reward}$$

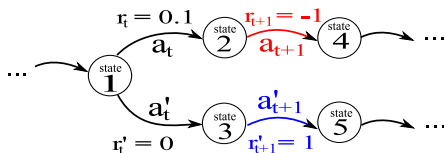
# Deceptive Rewards — Discounting the Future

- What if the future rewards are “deceptive”?
  - ▶ a **good** action followed by a **very bad** action
  - ▶ **future** needs to be taken into account



# Deceptive Rewards — Discounting the Future

- What if the future rewards are “deceptive”?
  - ▶ a good action followed by a very bad action
  - ▶ future needs to be taken into account

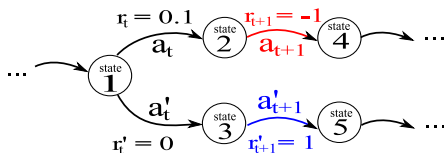


- Some measurement for future is needed: **FUTUREGOODNESS**

$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot Q(s_t, a_t) + \alpha (r_{t+1} + \text{FUTUREGOODNESS})$$

# Deceptive Rewards — Discounting the Future

- What if the future rewards are “deceptive”?
  - ▶ a good action followed by a very bad action
  - ▶ future needs to be taken into account

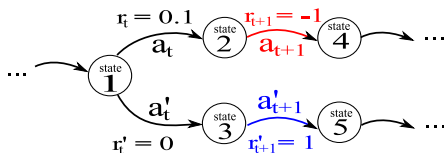


- Some measurement for future is needed: FUTUREGOODNESS
  - ▶ SARSA:  $\text{FUTUREGOODNESS} = \gamma \cdot Q(s_{t+1}, a_{t+1})$

$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot Q(s_t, a_t) + \alpha (r_{t+1} + \gamma \cdot Q(s_{t+1}, a_{t+1}))$$

# Deceptive Rewards — Discounting the Future

- What if the future rewards are “deceptive”?
  - ▶ a good action followed by a very bad action
  - ▶ future needs to be taken into account



- Some measurement for future is needed: **FUTUREGOODNESS**
  - ▶ SARSA:  $\text{FUTUREGOODNESS} = \gamma \cdot Q(s_{t+1}, a_{t+1})$
  - ▶ Q-Learning:  $\text{FUTUREGOODNESS} = \gamma \cdot \max_a Q(s_{t+1}, a)$

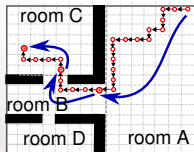
$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot Q(s_t, a_t) + \alpha (r_{t+1} + \gamma \cdot \max_a Q(s_{t+1}, a))$$



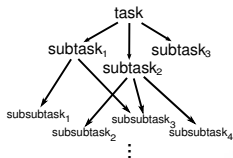
# Why use hierarchies at all?

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

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



extended actions

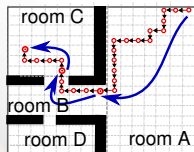


decomposition and re-use

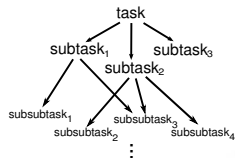
# Why use hierarchies at all?

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

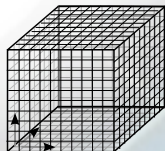
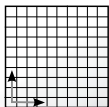
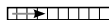
- “**Divide and Conquer**”: decompose into smaller (easier) subtasks
  - ▶ task decomposition enables **re-use of (sub)policies**
- Battling the “**Dæmon of Dimensionality**”
  - ▶ smaller state spaces on all levels



extended actions



decomposition and re-use

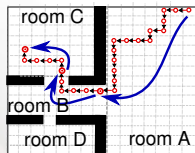


curse of dimensionality

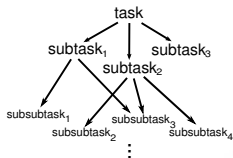
# Why use hierarchies at all?

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

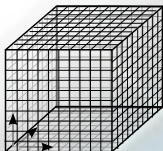
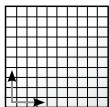
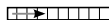
- “**Divide and Conquer**”: decompose into smaller (easier) subtasks
  - ▶ task decomposition enables **re-use of (sub)policies**
- Battling the “**Dæmon of Dimensionality**”
  - ▶ smaller state spaces on all levels
- **Faster exploration** >>



extended actions



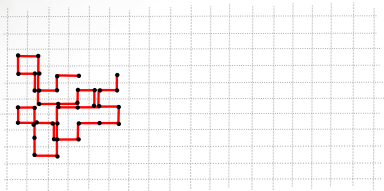
decomposition and re-use



curse of dimensionality

# Hierarchies Make Exploration Faster

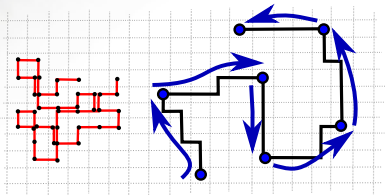
- Reinforcement Learning ▶ exploration is random walk
  - ▶ a new random action every time step



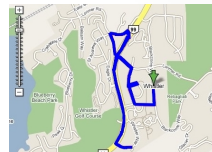
Drunken Mans Walk

# Hierarchies Make Exploration Faster

- Reinforcement Learning ► exploration is random walk
  - a new random action every time step



Drunken Mans Walk

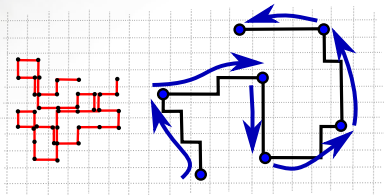


Random Walk on street level

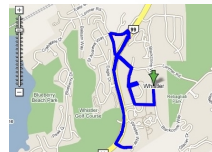
- But behaviours do something consistent (hopefully)
  - they move agent non-randomly through state space
  - Less random choices (only when selecting behaviour)
  - more distance: faster exploration

# Hierarchies Make Exploration Faster

- Reinforcement Learning ► exploration is random walk
  - a new random action every time step



Drunken Mans Walk



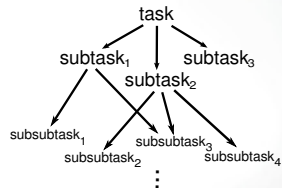
Random Walk on street level

- But behaviours do something consistent (hopefully)
  - they move agent non-randomly through state space
  - Less random choices (only when selecting behaviour)
  - more distance: faster exploration

Agent can discover meaningful behaviours long before it has a chance to solve overall problem

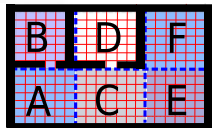
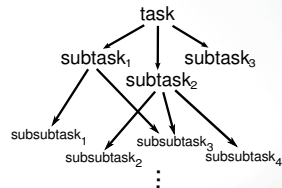
# Shifting Focus

- Many Hierarchical RL-approaches **focus on Task Decomposition**
  - ▶ designer **decomposes actions** into smaller and smaller actions
  - ▶ what if the task is **complicated** or the designer doesn't see the solution?
  - ▶ MAX-Q, HAM, HEX-Q, ...



# Shifting Focus

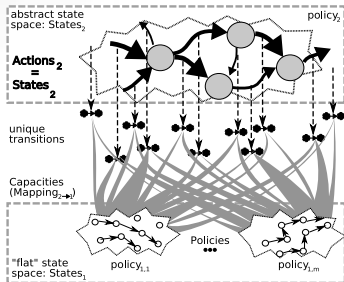
- Many Hierarchical RL-approaches **focus on Task Decomposition**
  - ▶ designer **decomposes actions** into smaller and smaller actions
  - ▶ what if the task is **complicated** or the designer doesn't see the solution?
  - ▶ MAX-Q, HAM, HEX-Q, ...



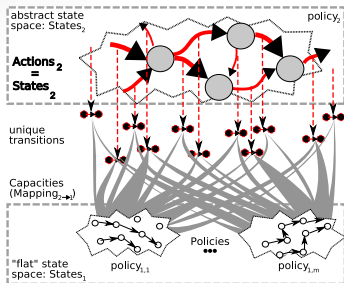
- **Focus on State Space Abstraction**

- ▶ complementary to task decomposition
- ▶ provide suitable **state abstraction**, and **let the RL-agent figure it out!**
- ▶ done in FEUDAL LEARNING and HASSLE

# Abstract States = High level Actions? (Goal Directed)

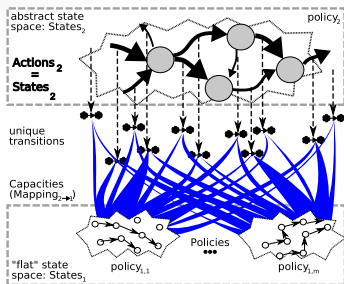


# Abstract States = High level Actions? (Goal Directed)



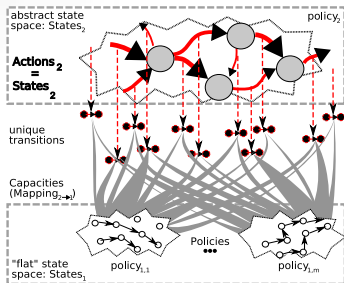
- Each transition (action) between high level states is **unique**

# Abstract States = High level Actions? (Goal Directed)



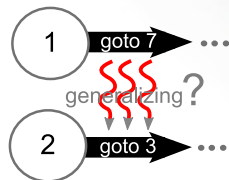
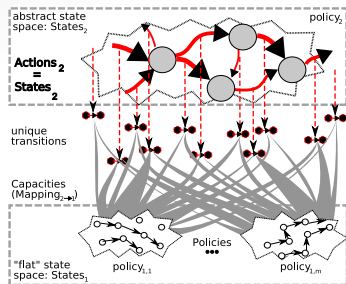
- Each transition (action) between high level states is **unique**
- Need to **map these to small set** of subpolicies

# Abstract States = High level Actions? (Goal Directed)



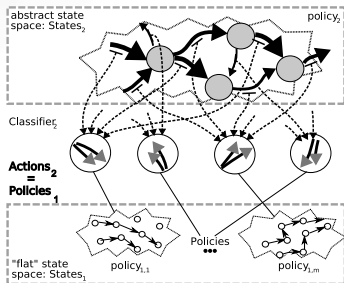
- Each transition (action) between high level states is **unique**
- Need to **map these to small set** of subpolicies
- Problem: High level states are **unique**
  - ▶  $||\text{high level actions}|| = ||\text{high level states}|| \sim \text{"explosion"}$

# Abstract States = High level Actions? (Goal Directed)

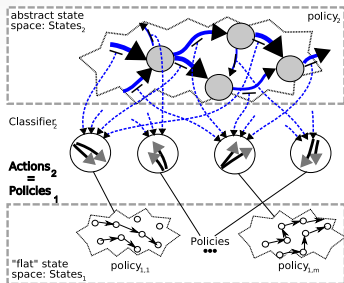


- Each transition (action) between high level states is **unique**
- Need to **map these to small set** of subpolicies
- Problem: High level states are **unique**
  - ▶  $||\text{high level actions}|| = ||\text{high level states}|| \sim \text{“explosion”}$
  - ▶ **no generalization or neural networks on high level!**  
(no structure in state descriptions)

# Abstract States $\neq$ High level Actions!

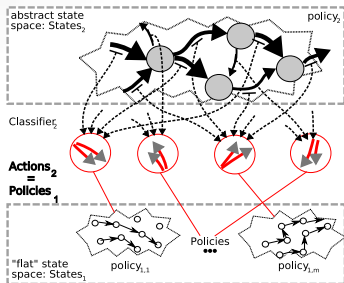


# Abstract States $\neq$ High level Actions!



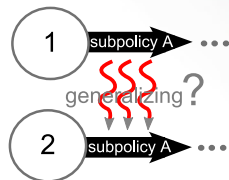
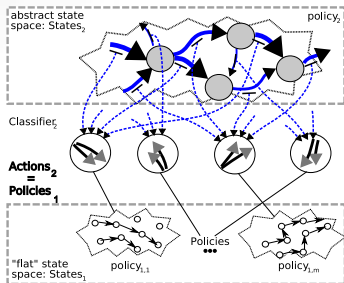
- Reverse order: first **many similar transitions** are classified

# Abstract States $\neq$ High level Actions!



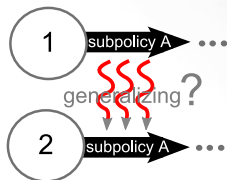
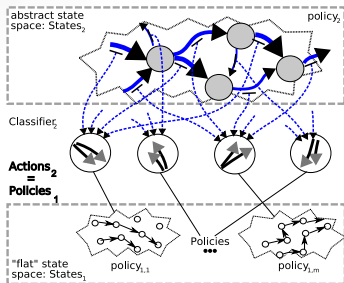
- Reverse order: first **many similar transitions** are classified
- Then each cluster **corresponds** with one subpolicy

# Abstract States $\neq$ High level Actions!



- Reverse order: first **many similar transitions** are classified
- Then each cluster **corresponds** with one subpolicy
- Use subpolicy as a **high level action** (still no generalizing)

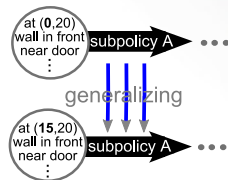
# Abstract States $\neq$ High level Actions!



- Reverse order: first **many similar transitions** are classified
- Then each cluster **corresponds** with one subpolicy
- Use subpolicy as a high level action (still no generalizing)
- **Problem: how do we find this classification?**
  - ▶ *given a priori* by designer
  - ▶ **discovered** by the agent during learning

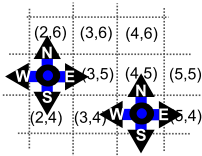
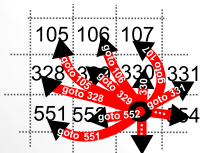
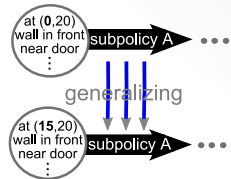
# Classification Requires Structure

- Abstract State Space needs structure similar to “flat” state space
  - ▶ **no structure, no heuristics** (e.g. hierarchies)!
  - ▶ structure makes **generalization** possible



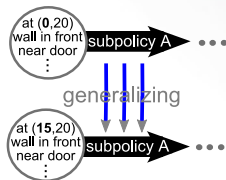
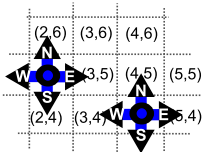
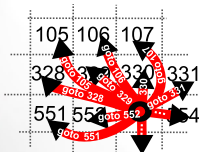
# Classification Requires Structure

- Abstract State Space needs **structure** similar to “flat” state space
  - ▶ no structure, no heuristics (e.g. hierarchies)!
  - ▶ structure makes generalization possible
- Primitive actions are not **goal directed** but **relative** to a state
  - ▶ **1000 cell** grid world with **only 4 actions!**



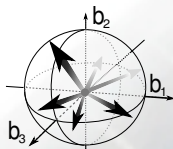
# Classification Requires Structure

- Abstract State Space needs **structure** similar to “flat” state space
  - ▶ no structure, no heuristics (e.g. hierarchies)!
  - ▶ structure makes generalization possible
- Primitive actions are not **goal directed** but **relative** to a state
  - ▶ 1000 cell grid world with only 4 actions!



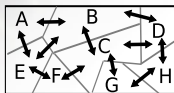
- **Action Space:**

- ▶ “space of all possible difference vectors in the state space”
- ▶ primitive actions **all map to few vectors** in the Action Space!

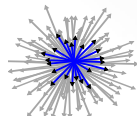


# Smart State Abstractions

- Suppose we would create this Abstract State Space ...



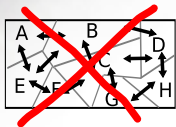
- ▶ ... yielding this Behaviour Space
- ▶ **no clustering, no natural groups!**



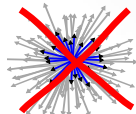
Note: “Behaviour Space” is short for “Abstract Action Space”

# Smart State Abstractions

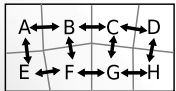
- Suppose we would create this Abstract State Space ...



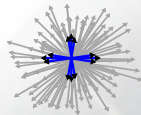
- ▶ ... yielding this Behaviour Space
- ▶ **no clustering, no natural groups!**



- Or this ...

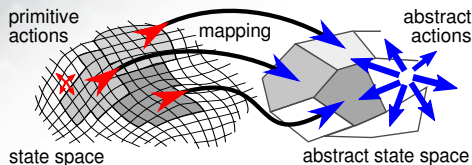


- ▶ ... yielding this Behaviour Space
- ▶ **groups suitable for classification!**
- ▶ like the primitive actions



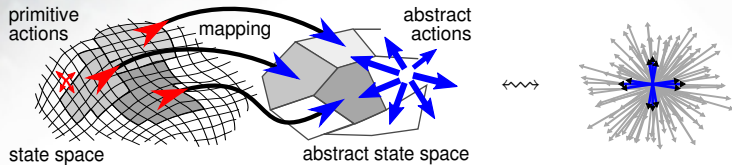
Note: "Behaviour Space" is short for "Abstract Action Space"

# Abstract State Space Properties



- 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
  - ▶ transitions in abstract state space correspond to meaningful behaviours in “flat” state space
- Abstract State Space **significantly smaller** than State Space

# Abstract State Space Properties



- 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
  - ▶ transitions in abstract state space correspond to meaningful behaviours in “flat” state space
- Abstract State Space **significantly smaller** than State Space
- **Actually occurring transitions** need to be distributed **non-uniformly** in Behaviour Space

# Proposed Algorithm

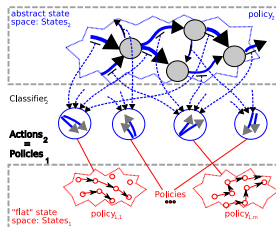
## HABS (Hierarchical Assignment of Behaviours by Self-organization)

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

# Proposed Algorithm

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

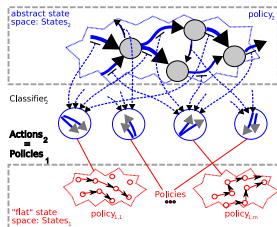
- One (high level)  $\text{Policy}_{\text{HL}}$  and a limited set of subpolicies
  - uses abstract state space
- $\text{Policy}_{\text{HL}}$  has subpolicies as its (extended) actions:
  - the classification and subpolicies are self-organizing
  - rewards received during a behaviour are accumulated and used for High Level  $\text{Policy}_{\text{HL}}$  reward
  - subpolicy rewarding is independent of overall task  $\ggg$



# Proposed Algorithm

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

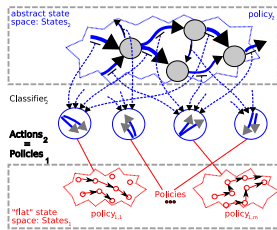
- One (high level)  $\text{Policy}_{\text{HL}}$  and a limited set of subpolicies
  - uses abstract state space
- $\text{Policy}_{\text{HL}}$  has subpolicies as its (extended) actions:
  - the classification and subpolicies are self-organizing
  - rewards received during a behaviour are accumulated and used for High Level  $\text{Policy}_{\text{HL}}$  reward
  - subpolicy rewarding is independent of overall task  $\ggg$
- Standard RL techniques** (online, off policy) like **Q-learning**



# Proposed Algorithm

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

- One (high level)  $\text{Policy}_{\text{HL}}$  and a limited set of subpolicies
  - uses abstract state space
- $\text{Policy}_{\text{HL}}$  has subpolicies as its (extended) actions:
  - the classification and subpolicies are self-organizing
  - rewards received during a behaviour are accumulated and used for High Level  $\text{Policy}_{\text{HL}}$  reward
  - subpolicy rewarding is independent of overall task  $\ggg$
- Standard RL techniques (online, off policy) like Q-learning
- First solve subtasks (subpolicies), then overall task ( $\text{Policy}_{\text{HL}}$ )



# Training Subpolicies, How?

- a subpolicy starts with **no knowledge** (i.e. randomly initiated)
  - ▶ what is its **desired or characteristic behaviour?**
- **train on pairs of abstract states?**
  - ▶ designer needs to **specify pre/post conditions**

# Training Subpolicies, How?

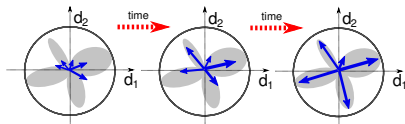
- a subpolicy starts with no knowledge (i.e. randomly initiated)
  - ▶ what is its desired or characteristic behaviour?
- train on pairs of abstract states?
  - ▶ designer needs to specify pre/post conditions
- rewards **independent** of the overall task ( $\text{Policy}_{HL}$ )
  - ▶ behaviour “**A**  $\Rightarrow$  **goal**” same as “**B**  $\Rightarrow$  **C**”
  - ▶ **blue behaviour** has **high**  $Q_{HL}$ -value in **A** but **low**  $Q_{HL}$ -value in **B**
  - ▶ **red behaviour** has **high**  $Q_{HL}$ -value in **B**



no dependence on  $\text{Policy}_{HL}$  on fixed pre/post conditions!

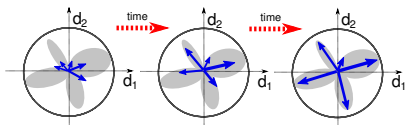
# Clustering and Rewarding Subpolicies

- 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

$$reward_{sub} = \begin{cases} 0 & \text{not terminated} \\ 1 & \text{terminated: closest match} \\ \kappa_r & \text{terminated: another cluster center is closer} \\ \kappa_f & \text{terminated: timeout (failed to reach anything)} \end{cases}$$



# HABS in Pseudo Code

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>

# HABS in Pseudo Code

repeat

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

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 *timeout<sub>SUB</sub>*

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<sub>i</sub>* ; // match

move *CLUSTER<sub>SUB</sub>* towards *EXEC* ; // match

else punish *SUB<sub>i</sub>* ; // no match

update *Policy<sub>HL</sub>* with *reward<sub>HL</sub>* ;

until task solved or *timeout<sub>HL</sub>*

# HABS in Pseudo Code

```

repeat
   $reward_{HL} = 0$  ; // for accumulating rewards
   $Policy_{HL}$  selects  $SubPolicy\ SUB_i$  ;
  repeat
     $SUB_i$  selects and executes a primitive action ;
     $reward_{HL} \leftarrow reward_{HL} + receivedReward$  ; // accumulate
    if new abstract state then BREAK ;
    else update  $SUB_i$  with 0 ;
  until  $timeout_{SUB}$ 
  if  $timeout_{SUB}$  then punish  $SUB_i$  ;
  else
    if  $EXEC \in CLUSTER_{SUB}$  then
      reward  $SUB_i$  ;
      move  $CLUSTER_{SUB}$  towards  $EXEC$  ;
    else punish  $SUB_i$  ;
  update  $Policy_{HL}$  with  $reward_{HL}$  ;
until task solved or  $timeout_{HL}$ 

```

# HABS in Pseudo Code

```

repeat
   $reward_{HL} = 0$  ;
   $Policy_{HL}$  selects SubPolicy  $SUB_i$  ;
  repeat
     $SUB_i$  selects and executes a primitive action ;
     $reward_{HL} \leftarrow reward_{HL} + receivedReward$  ;
    if new abstract state then BREAK ;
    else update  $SUB_i$  with 0 ;
  until  $timeout_{SUB}$ 
  if  $timeout_{SUB}$  then punish  $SUB_i$  ;
  else
    if  $EXEC \in CLUSTER_{SUB}$  then
      reward  $SUB_i$  ;
      move  $CLUSTER_{SUB}$  towards  $EXEC$  ;
    else punish  $SUB_i$  ;
  update  $Policy_{HL}$  with  $reward_{HL}$  ;
until task solved or  $timeout_{HL}$ 

```

# 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)
- ▶ rewarded for delivered objects (at drop areas)


Task has spatial **and non-spatial** aspects!

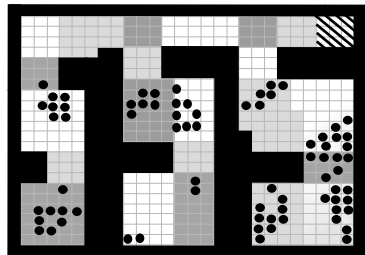
# Experiment Description (“Cleaner”)

- **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)
- ▶ rewarded for delivered objects (at drop areas)

- **“Cleaner”:**

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

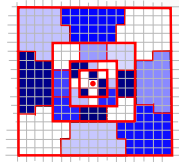


Task has spatial **and non-spatial** aspects!

# State Space and Abstract State Spaces


**State Space:** (subpolicies)

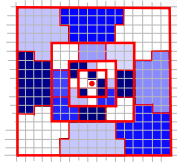
- Agent has a **simulated “radar”**  observes object/wall/drop areas ( $\sim 100$  inputs)



# State Space and Abstract State Spaces

**State Space:** (subpolicies)

- Agent has a simulated “radar”  observes object/wall/drop areas ( $\sim 100$  inputs)




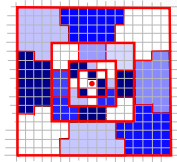
**Abstract State Space:** (Policy<sub>HL</sub>)

- “cleaner” task **only** uses  $\langle \text{area}_{\text{agent}}, \text{cargo} \rangle$  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 ( $\sim 100$  inputs)



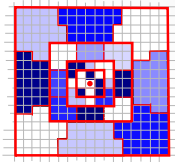
## Abstract State Space: (Policy<sub>HL</sub>)

- “cleaner” task **only** uses  $\langle \mathit{area}_{\mathit{agent}}, \mathit{cargo} \rangle$  to determine subpolicy termination
  - ▶ **no info about other objects!**
  - ▶ **objects moving/(dis)appearing on themselves** 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)

# State Space and Abstract State Spaces

**State Space:** (subpolicies)

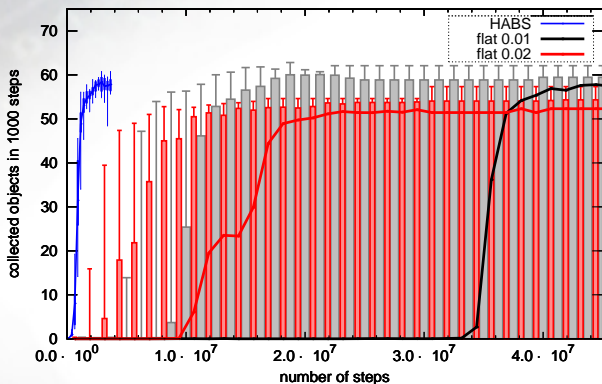
- Agent has a simulated “radar”  observes object/wall/drop areas ( $\sim 100$  inputs)



**Abstract State Space:** (Policy<sub>HL</sub>)

- “cleaner” task **only** uses  $\langle \text{area}_{\text{agent}}, \text{cargo} \rangle$  to determine subpolicy termination
  - ▶ no info about other objects!
  - ▶ objects moving/(dis)appearing on themselves 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** – shows that termination criteria can differ from abstract states

# Cleaner: Results



## Boxplots:

- **HABS**

- “flat” learner

$\alpha = \{0.02, 0.01\}$

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  $\updownarrow$

- ▶ “flat” has far wider variance in convergence time  $\longleftrightarrow$

# Demonstration

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

# Conclusions / Future Work

- Conclusions about HABS:
  - ▶ shifting design burden from task decomposition to state space abstraction is possible
  - ▶ learns conditions for behaviours by self-organizing
  - ▶ starts with unspecified behaviours (no fixed pre/post conditions)

# Conclusions / Future Work

## ● Conclusions about HABS:

- ▶ shifting design burden from task decomposition to state space abstraction is possible
- ▶ learns conditions for behaviours by **self-organizing**
- ▶ starts with **unspecified behaviours** (no fixed pre/post conditions)

## ● Future work

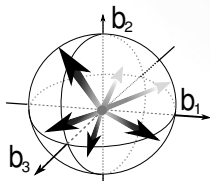
- ▶ try HABS with **3 or more layers** (possible, in theory)
- ▶ behaviour **representation**: limited to vectors?
- ▶ **automatically detect relevant features** in abstract state space which determine behaviour (instead of selecting them *a priori*)
- ▶ test **continuous termination criteria** (“distance traveled”, instead of “reached new abstract state”)



# 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
  - ▶ only one primitive action for *North* instead of many  $North_1, North_2, North_3, \dots$  for going north from *state 1* to 352, from *state 2* to 369, from *state 3* to 792345,  $\dots$
- **relative** (vectors), not absolute (*North* = “in state *A* goto *B*”)



▶ back to behaviour space

# Behaviours Mirror Primitive Actions

## Primitive Actions

## Behaviours

vectors in action space  
 relative to state  
 clustered  
 1 time step  
 action successful  
 action fails

vectors in behaviour space  
 relative to abstract state  
 hopefully clustered  
 1 high level time step  
 reach new abstract state  
 timeout

▶ back to behaviour space

# Training Subpolicies, How?

- a subpolicy starts with **no knowledge** (i.e. randomly initiated)
  - ▶ 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 ( $\text{Policy}_{HL}$ )
  - ▶ behaviour "**A**  $\Rightarrow$  **goal**" same as "**B**  $\Rightarrow$  **C**"
  - ▶ **blue behaviour** has **high**  $Q_{HL}$ -value in **A** but **low**  $Q_{HL}$ -value in **B**
  - ▶ **red behaviour** has **high**  $Q_{HL}$ -value in **B**



no dependence on  $\text{Policy}_{HL}$  on fixed pre/post conditions!

▶ back to  $\text{Policy}_{HL}$

# 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](#)

# Q-Learning and Advantage Learning

## Q-Learning:

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

## Advantage Learning (*Baird*):

$$A(s_t, a_t) \leftarrow (1 - \alpha) \cdot A(s_t, a_t) + \alpha \cdot \left( \frac{\max_{a_t} A(s_t, a_t) + \text{reward}_{t+1} + \gamma \max_a A(s_{t+1}, a) - \max_{a_t} A(s_t, a_t)}{k} \right)$$

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

# Formula for Clustering

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

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

where  $act_{t \rightarrow t+\Delta t}$  is the actually executed behaviour,  $\Delta t$  the time it took to execute the subpolicy, and  $char_t$  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

# 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

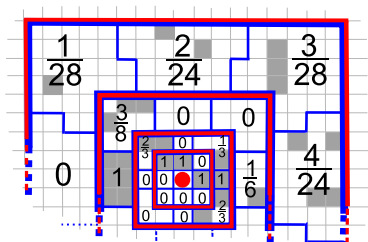
At every time step during the execution of a subpolicy, the normal RL (Q-Learning, Advantage Learning, ...) update is applied:

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

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

# Sensors

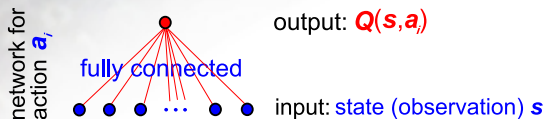
- agent has radar-like sensorgrid (8 areas per ring)
  - ▶ details nearby, coarse observations far away
  - ▶ trade off between detail and amount of sensor data



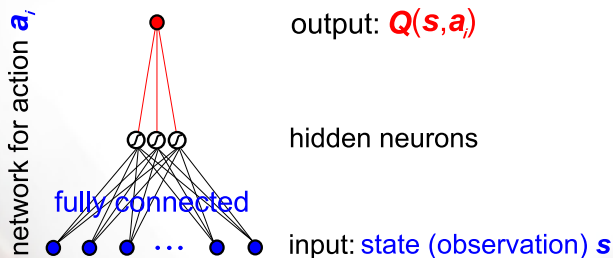
- observation: vector of area densities
  - ▶  $\langle \frac{1}{28}, \frac{2}{24}, \frac{3}{28}, \frac{4}{24}, \dots, \frac{3}{8}, 0, 0, \frac{1}{6}, \dots, \frac{2}{3}, 0, \frac{1}{3}, 1, \dots, 1, 1, 0, 1, \dots \rangle$
  - ▶ for walls / objects / drop areas, so  $3 \times 32$  inputs
- extra values coding for position, cargo, etc were added

# Neural Nets for Subpolicies 1

Linear neural network:

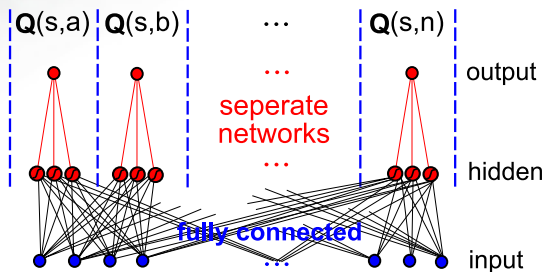


Multilayer Neural network:



Both are “feed forward” networks (no recurrence)

# Neural Nets for Subpolicies 2



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

# 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**):

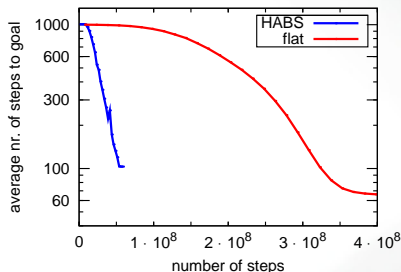
- ▶ big ( $39 \times 36 \approx 1.4 \cdot 10^3$  cells)
- ▶ only **one** object
- ▶ high level policy **tabular**
- ▶ subpolicies: **linear neural network** ▶

▶ back to Cleaner Description

▶ back to Cleaner Results

# First Experiment (Maze): Results

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



- 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  $5 \cdot 10^3 \times \text{numberOfSubpolicies}$  ( $\approx 1$  megabyte)
  - ▶ Neural network storage is neglectable

▶ back to Cleaner Description

▶ back to Cleaner Results