

# Value Representation in Large Factored State Spaces



**Wendelin Böhmer**  
 <wendelin@ni.tu-berlin.de>  
 Neural Information Processing Group, Technische Universität Berlin

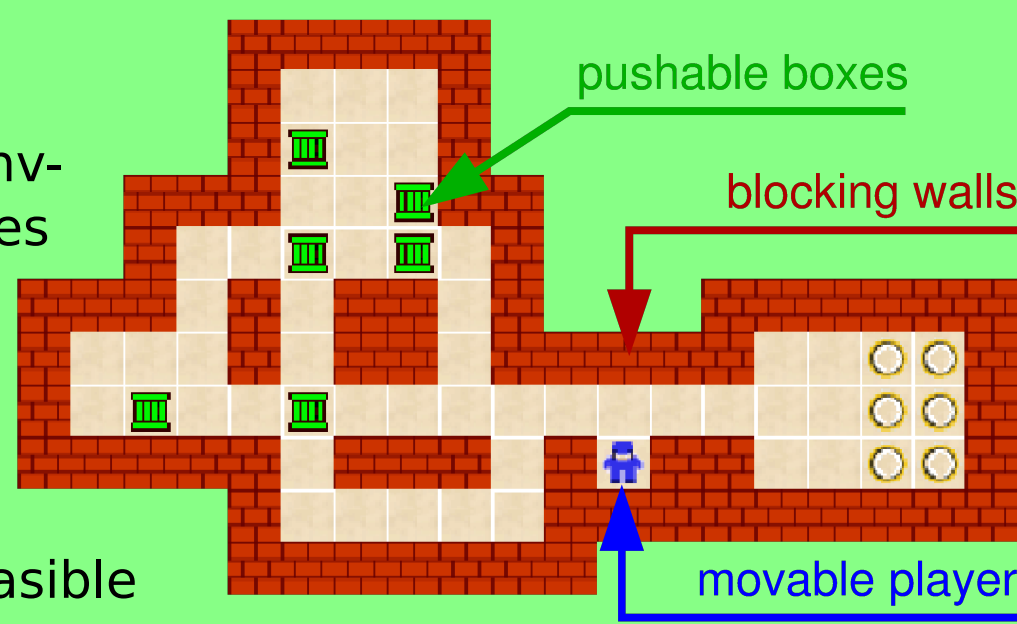
**Klaus Obermayer**  
 <oby@ni.tu-berlin.de>  
 Neural Information Processing Group, Technische Universität Berlin



## Exploiting Factorization in Large Metric State Spaces

### Overview:

- Autonomous agents in dynamic environments, e.g. with pushable boxes
- Stationary planning induces high dimensional state spaces ( $\mathbb{O}(2^n)$ )
- Example: puzzle game *Sokoban*



**Problem A:** sampling all states infeasible

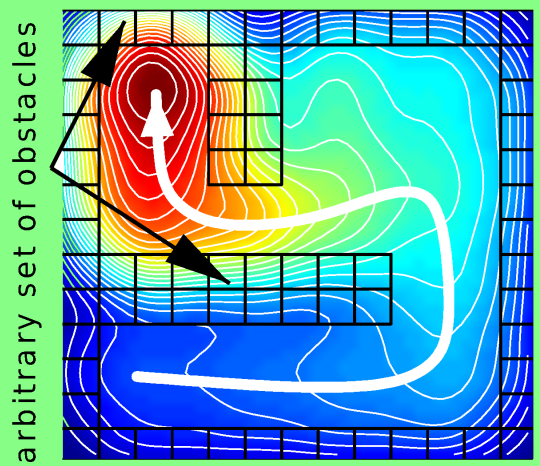
**Problem B:** no transfer to new environments

### Factored structure of induced state spaces:

- Often state variables describe "objects", e.g. detected from vision
- Reward functions often factorize: reward depends only on few variables
- Factored MDP can be exploited, but transition models rarely factorize

### Divide & Conquer models:

- Relational description of transition rules
- Metric transition effect models
- Decomposable factorization
- Adapts to new situations



## Project Goal A

### Bayesian Learning of Relational Models for Transitions and Reward from Data

#### Work Package A1: lean models

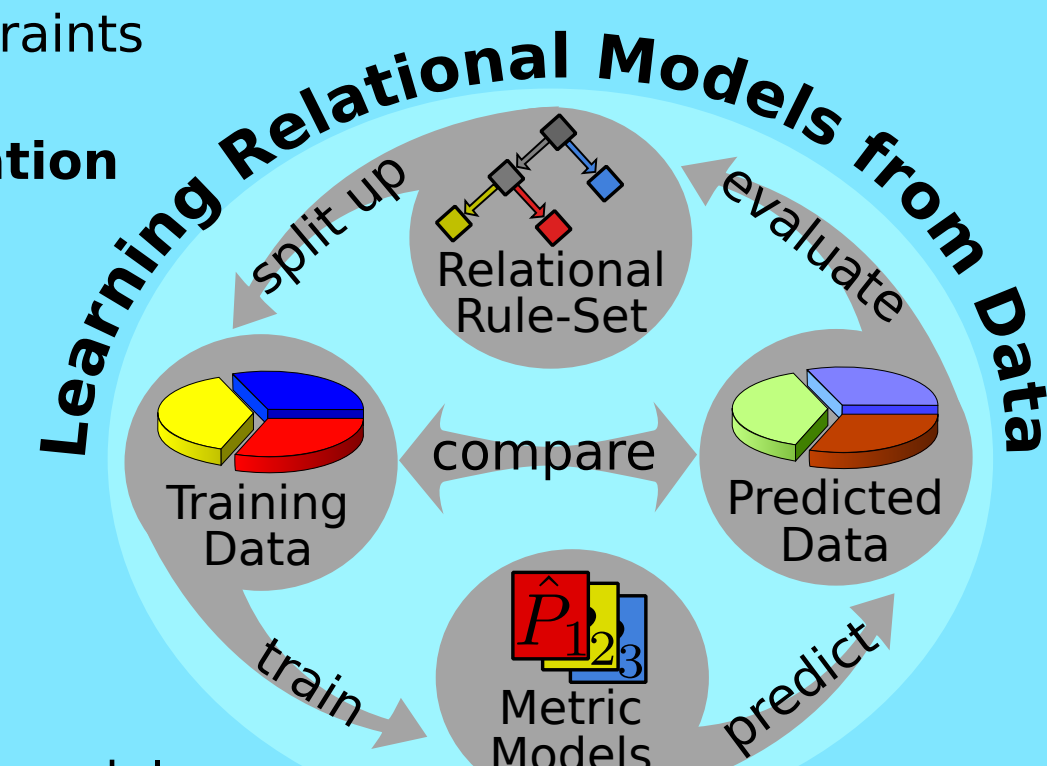
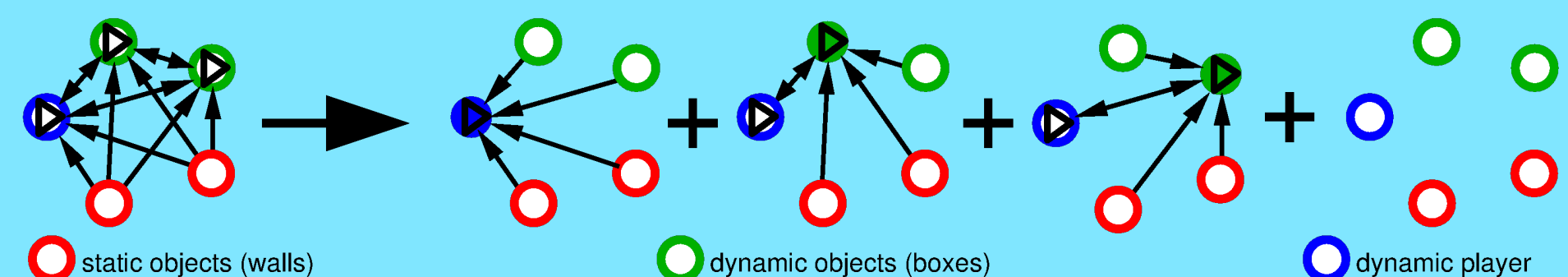
- Learning relational models is greedy
- Bayesian distribution intractable in practice
- Sparse coding allows tractable approximations
- Learn rule-sets that split training data
- Include structural (factored) constraints

#### Work Package A2: active exploitation

- Training data is finite but may be extended by active exploration
- Find the most informative states, i.e. where likely models contradict
- Explore those states actively, i.e. scientific hypothesis testing

Exploit factored structure by decomposition in multiple metric transition models:

full dependency graph = "player moves unrestricted" or "player pushes box A" or "player pushes box B" or "player cannot move"



## Representation with "Linear Factored Functions"

$$\text{LFF: } f(\vec{x}) = \vec{a}^\top \vec{\varphi}(\vec{x}) = \vec{a}^\top \left[ \prod_{\alpha} \vec{\varphi}^{\alpha}(x_{\alpha}) \right] = \vec{a}^\top \left[ \prod_{\alpha} \mathbf{B}^{\alpha} \vec{\phi}^{\alpha}(x_{\alpha}) \right]$$

- Inner products & marginalization can be computed analytically as well as point-wise products (increases number of bases)
- Partial derivatives factorize in related LFF
- No analytical nonlinearities like max or inverse

### Developed greedy LFF algorithms:

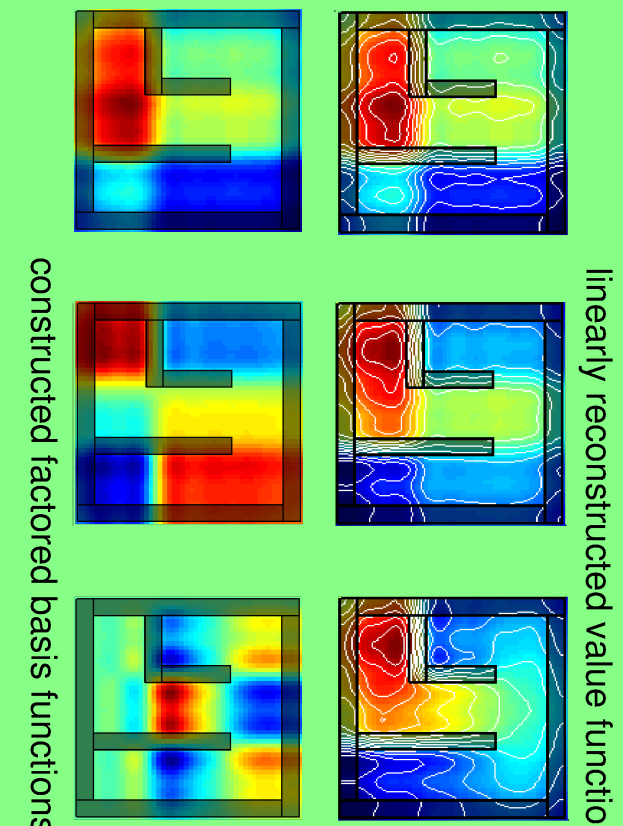
- Compression, e.g. after multiplication
- Density estimation from data  $\{\vec{x}_t\}_{t=1}^n \sim \xi$
- Regression from labels  $\{y(\vec{x}_t) + \delta_t\}_{t=1}^n$
- Sparse regression that adjusts uncertainty prior

$$\inf_{f, \vec{\mu}, \eta} \|y - f\|_{\xi}^2 + \frac{1}{\eta} D_{\text{KL}}(\vec{\mu} \| \frac{1}{d})$$

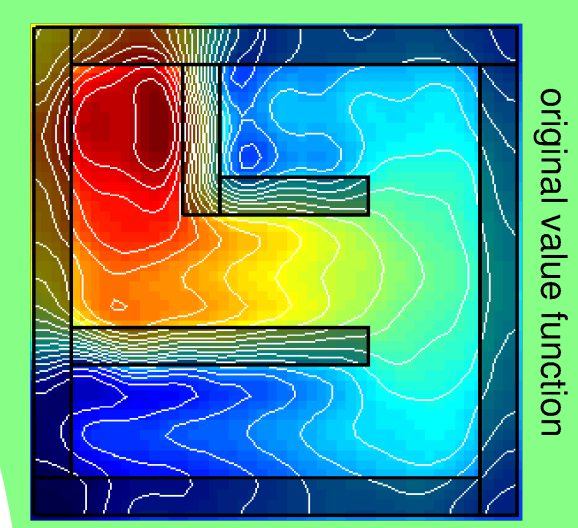
$$\text{s.t. } \|\varphi_i^{\alpha}\|_{\nu^{\alpha}} = 1, \vec{x} \leftarrow \vec{x} + \vec{e}(\vec{x})$$

$$\epsilon_{\alpha}(\vec{x}) \sim \mathcal{N}(0, \mu_{\alpha} \frac{d\nu}{d\xi}(\vec{x}))$$

$$\mu_{\alpha} \geq 0, \sum_{\alpha} \mu_{\alpha} = 1$$



- Current research focus:**
- Increase stability under multiplication with priors
  - Learn factorizing transition operators
  - Nonlinear policy improvement by symbolic reasoning



construct  
factored  
basis

relational  
policy iteration  
with metric  
evaluation

provide  
metric policies  
for action  
symbols

ground  
relational  
rules

$$\text{Q-Value Function}$$

$$Q^{\pi}(\vec{x}, \vec{a}) = \sum_i a_i \psi_i(\vec{z}) = \sum_i a_i \prod_{\alpha} \psi_i^{\alpha}(z_{\alpha})$$

$$\text{Value Function}$$

$$V^{\pi}(\vec{x}) = \sum_i w_i \varphi_i(\vec{x}) = \sum_i w_i \prod_{\alpha} \varphi_i^{\alpha}(x_{\alpha})$$

$$\text{Metric Action Policies}$$

$$\hat{\Pi}_k[\psi](\vec{x}) = \int \pi_k(d\vec{a}) \psi(\vec{x}, \vec{a}) = \psi^{\pi}(\vec{x}) \cdot \hat{\Pi}_k[\psi^{\pi}](\vec{x})$$

$$\text{Divide & Conquer Policies}$$

$$\hat{\Pi}_{\pi}[\psi] = \hat{\Pi}_{\tau}[\psi] + \sum_k p_k (\hat{\Pi}_k[\psi] - \hat{\Pi}_{\tau}[\psi])$$

$$\text{Grounded Relational Conditions}$$

$$p_k(\vec{x}) = \prod_{\alpha} p_k^{\alpha}(\vec{x}_{\alpha}, x_{\alpha})$$

**Relational Rules**

"collision predicate"  $c()$

- Sokoban Rule 1:  $\forall (k \in B \cup W). \neg c(x, k, a)$
- Sokoban Rule 2:  $\exists (b \in B). [c(x, b, a) \wedge \forall (k \in B \cup W \setminus b). \neg c(b, k, a)]$

**State and Action Symbols**

## Project Goal B

### Utilize Metric Information for Relational Policy Improvement, and to Adapt Symbols

#### Work Package B1: metric evaluation

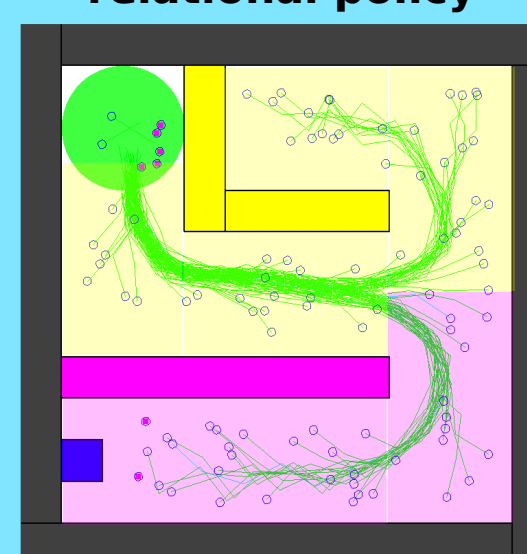
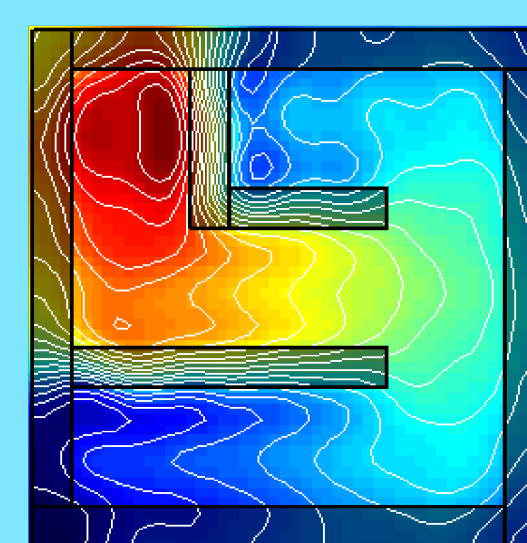
- Merge relational planning and metric values
- Relational RL yields action sequences, no policies
- Evaluate candidate AS with their metric value

#### Work Package B2: decision tree policies

- Grounded relational policies are decision trees
- Construct decision trees that maximize Q-value
- Greedy predicate selection vs. Bayesian methods

#### Work Package B3: adjust relational symbols

- Relational actions  $a_k$  have metric policies  $\pi_k$
- Most actions  $a_k$  depend on few variables
- Interprete Q-values as rewards for  $\pi_k$ , i.e. adjust  $\pi_k$  with standard RL in reduced space
- New actions: frequent local optimizations



Example relational policy on the left:

- Conditions select yellow and magenta areas
- Action  $a_m$ : keep magenta wall to your left
- Action  $a_y$ : keep yellow wall to your right



**Linking Metric and Symbolic Levels in Autonomous Reinforcement Learning**

**DFG** Schwerpunktprogramm 1527  
 Autonomous Learning