

# Solving K-MDPs

Jonathan Ferrer-Mestres  
Conservation Decisions  
Land and Water, CSIRO

Tom Dietterich  
Electrical Engineering and  
Computer Science  
Oregon State University

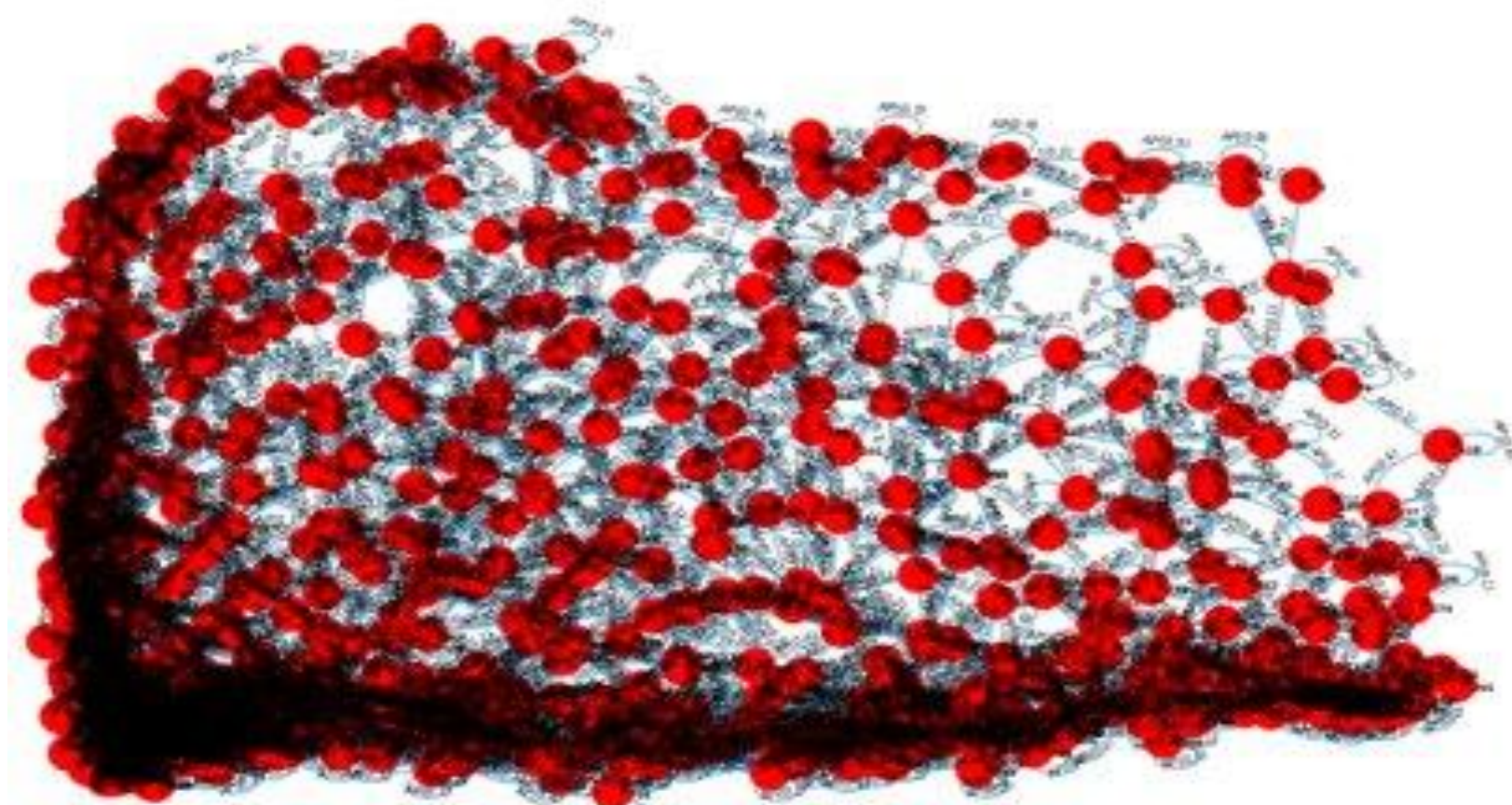
Olivier Buffet  
INRIA

Iadine Chades  
Conservation Decisions  
Land and Water, CSIRO

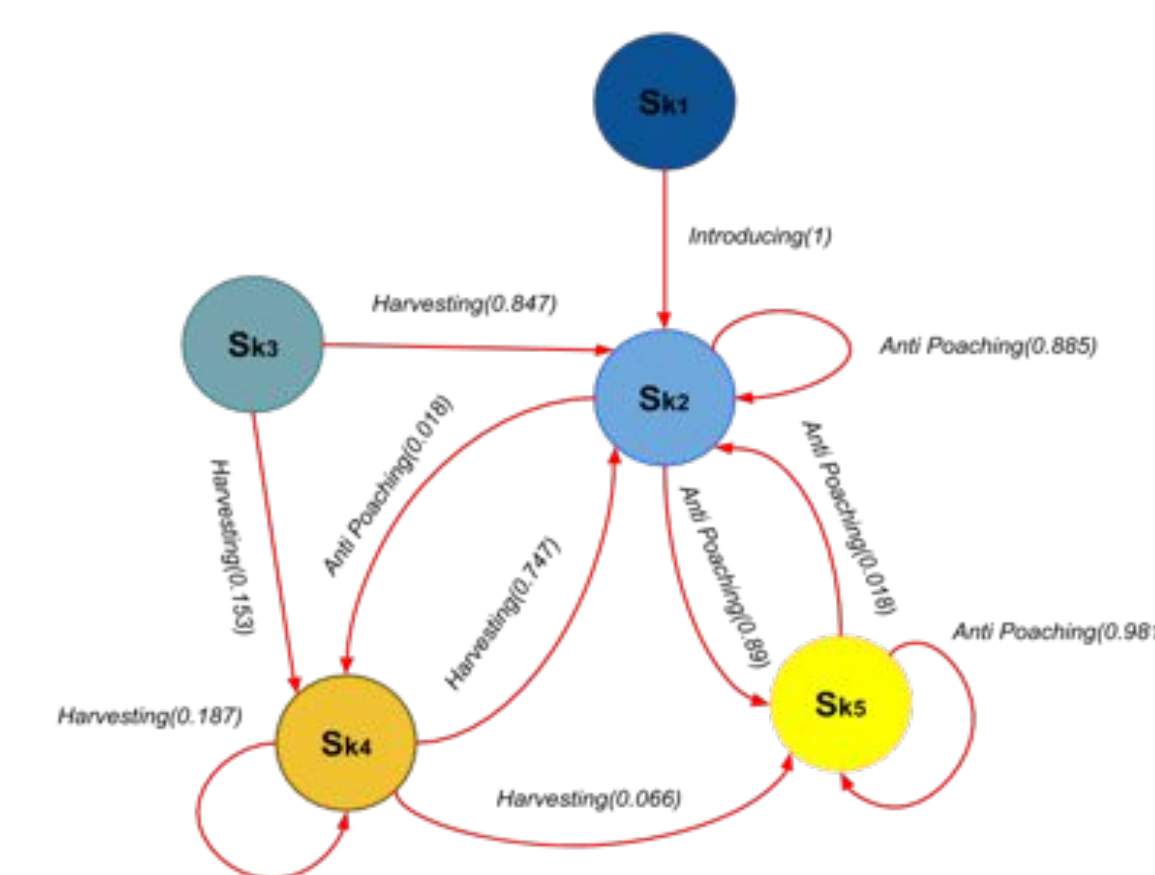
LAND AND WATER  
www.csiro.au

Markov decision processes (**MDPs**) are a convenient mathematical model for tackling sequential decision-making problems under uncertainty. MDPs have been applied to domains such as conservation of biodiversity and health. Policies computed for MDPs with thousands of states are in practice difficult to understand by humans. Developing **easy-to-interpret** solutions is crucial to increase uptake of MDP policies. We propose to increase the **interpretability** of MDPs by solving **K-MDPs**, i.e., given an original MDP and a number of states (**K**), generate a reduced state space MDP that **minimizes** the difference between the original and reduced optimal solutions.

Difficult to interpret



We propose to compute MDPs with K states: **K-MDPs**



Interpretable solution

## Methods

The aim of this research is to generate a **reduced state space** MDP that **minimizes** the difference between the original optimal MDP value function and the reduced optimal value function.

We solve **K-MDPs** using **state abstractions**. Formally, a **K-MDP** is a tuple  $\langle S_K, A, T_K, r_K, \gamma, \Phi \rangle$ , where:

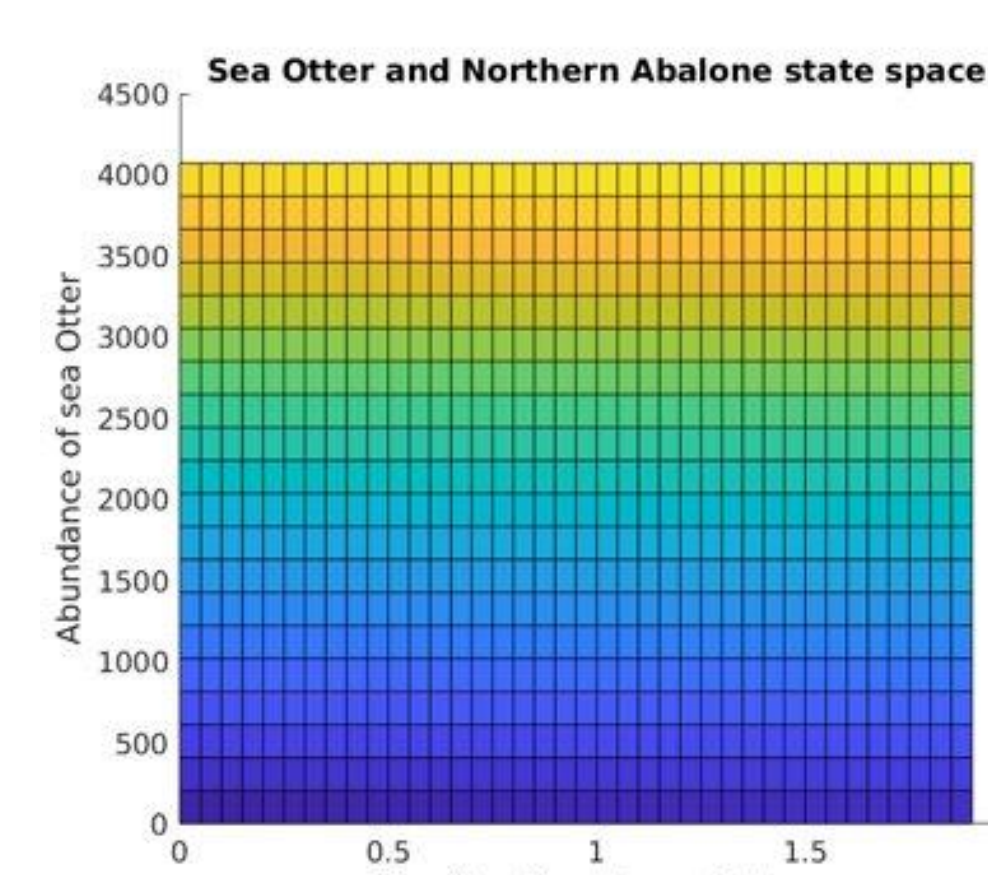
- $S_K$  is the abstract state space;
- $A$  is a set of actions;
- $T_K$  is the abstract **K-MDP** probability transition function;
- $r_K$  is the abstract reward function;
- $\gamma$  is the discount factor;
- $\Phi$  is a mapping function from  $S$  to  $S_K$ .

## Proposed algorithms

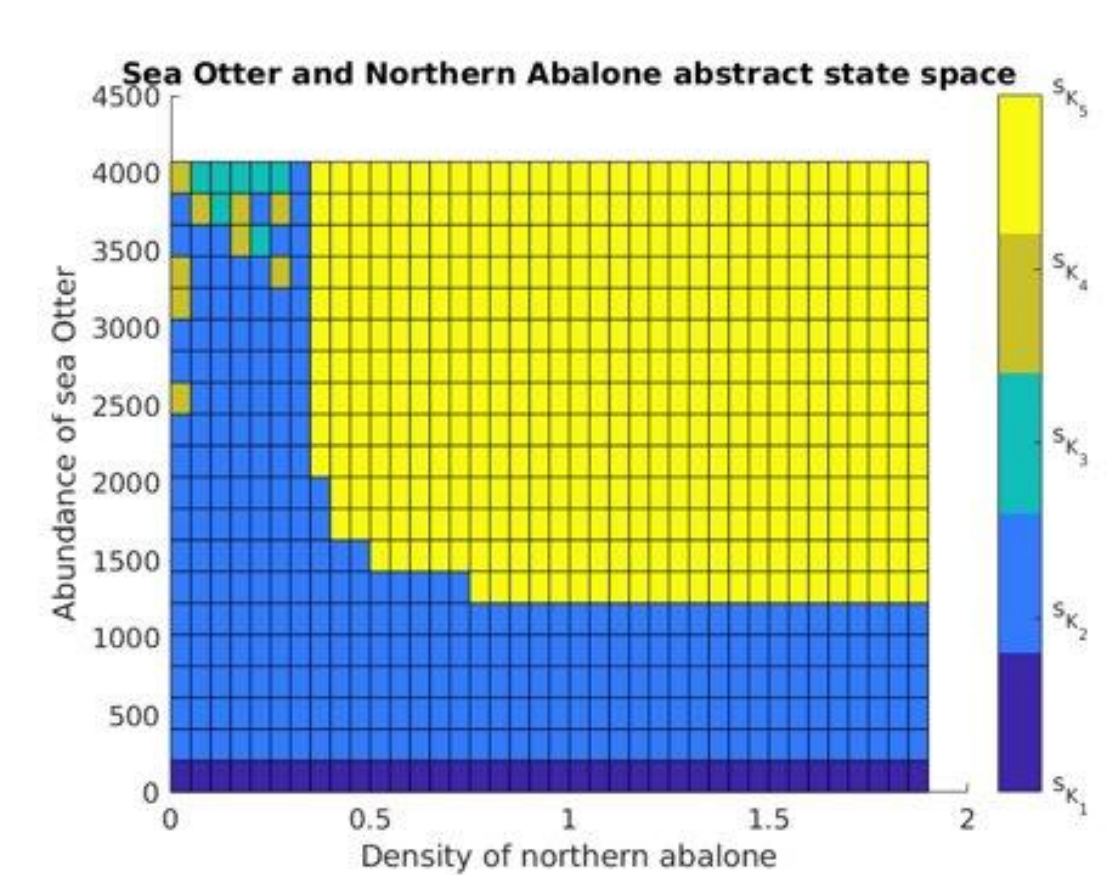
- Our algorithms group states into **bins**.
  - If two or more states are in the same bin, they can be **aggregated**.
  - **Challenge**: How to define the size of the bin if we want **K** states?
- We propose:
- Four algorithms based on **binary search** on  $\epsilon$  and  $d$ .
  - Finding best  $\epsilon$  and  $d$  that will guarantee the best performance.
  - One algorithm based on the **clustering technique** *k*-means++.

Algorithm	Function	Predicate	Value Loss	Trans.
$\phi_{Q_2}$ <i>K-MDP-ILP new</i>	$\phi_{Q_2}$ [Abel et al., 2016]	$\max_a  Q^*(i, a) - Q^*(j, a)  \leq \epsilon$	$\max_{s \in S} V^{\pi^*}(s) - V_{\phi_{Q_2}}^{\pi_K}(s) \leq \frac{2\epsilon R_{max}}{(1-\gamma)^2}$	No
$\phi_{Q_2}$ <i>K-MDP new</i>	$\phi_{Q_2}$ [Abel et al., 2018]	$\forall a \left  \frac{Q^*(i, a)}{d} \right  = \left\lfloor \frac{Q^*(j, a)}{d} \right\rfloor$	$\max_{s \in S} V^{\pi^*}(s) - V_{\phi_{Q_2}}^{\pi_K}(s) \leq \frac{2d R_{max}}{(1-\gamma)^2}$	Yes
$\phi_{a_d^*}$ <i>K-MDP new</i>	$\phi_{a_d^*}$ new	$a_i^* = a_j^* \wedge \left\lfloor \frac{V^*(i)}{d} \right\rfloor = \left\lfloor \frac{V^*(j)}{d} \right\rfloor$	$\max_{s \in S} V^{\pi^*}(s) - V_{\phi_{a_d^*}}^{\pi_K}(s) \leq \frac{2d R_{max}}{(1-\gamma)^2}$	Yes
$\phi_{Q_2}$ Greedy <i>K-MDP new</i>	$\phi_{Q_2}$ [Abel et al., 2016]	$\max_a  Q^*(i, a) - Q^*(j, a)  \leq \epsilon$	$\max_{s \in S} V^{\pi^*}(s) - V_{\phi_{Q_2}}^{\pi_K}(s) \leq \frac{2\epsilon R_{max}}{(1-\gamma)^2}$	No
<i>k</i> -means++ <i>K-MDP new</i>	- new	-	-	Yes

## Experiments



Chades et al., 2012  
MDP model with 819 states



Ferrer Mestres et al., 2020  
*K*-MDP model with 5 states

An optimal solution for a **K-MDP** can be applied to the original MDP using the mapping function  $\Phi$ :

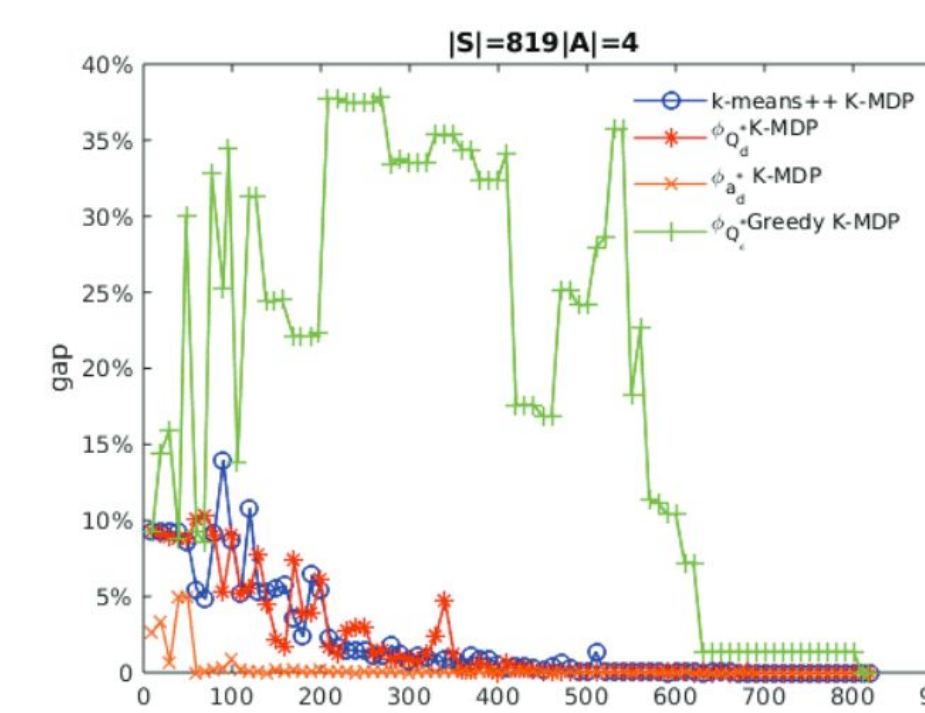
$$V_{\phi}^{\pi_K}(s) = E \left( \sum_{t=0}^{t=H} \gamma^t r(s_t, \pi_K^*(\phi(s_t))) \mid s_0 = s \right).$$

We formulate the problem of finding the best reduced state space as a **gap minimization** problem:

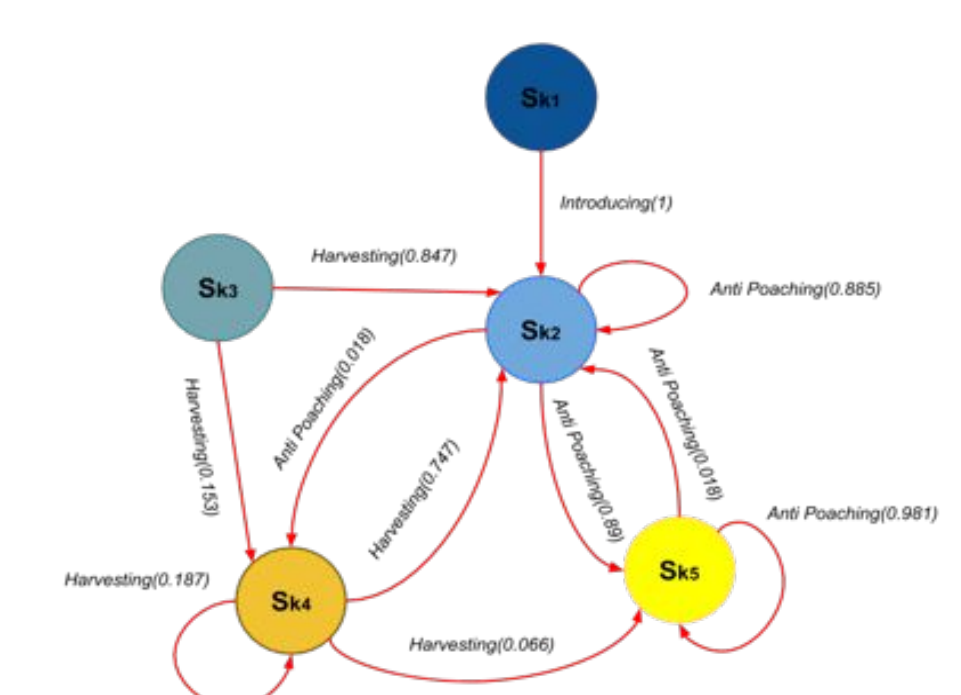
$$gap^* = \min_{S_K \in P(S), |S_K| \leq K} \max_{s \in S} [V^{\pi^*}(s) - V_{\phi}^{\pi_K}(s)]$$

## Discussion

Our approach aims at increasing **uptake** of MDPs in human operated systems by providing easier to **interpret** models and solutions. We have proposed to solve **K-MDPs** using state abstraction functions and a clustering technique. Future work will test the **interpretability** of proposed **K-MDP** solutions and will address visualization and interpretability challenges in problems with a large number of state variables.



Performance of *K*-MDP algorithms.



Reduced policy graph with  $K=5$  abstract states.