

Incremental Policy Iteration with Guaranteed Escape from Local Optima in POMDP Planning

Marek Grzes¹
School of Computing, University of Kent, UK

Pascal Poupart²
David Cheriton School of Computer Science, University of Waterloo, Canada

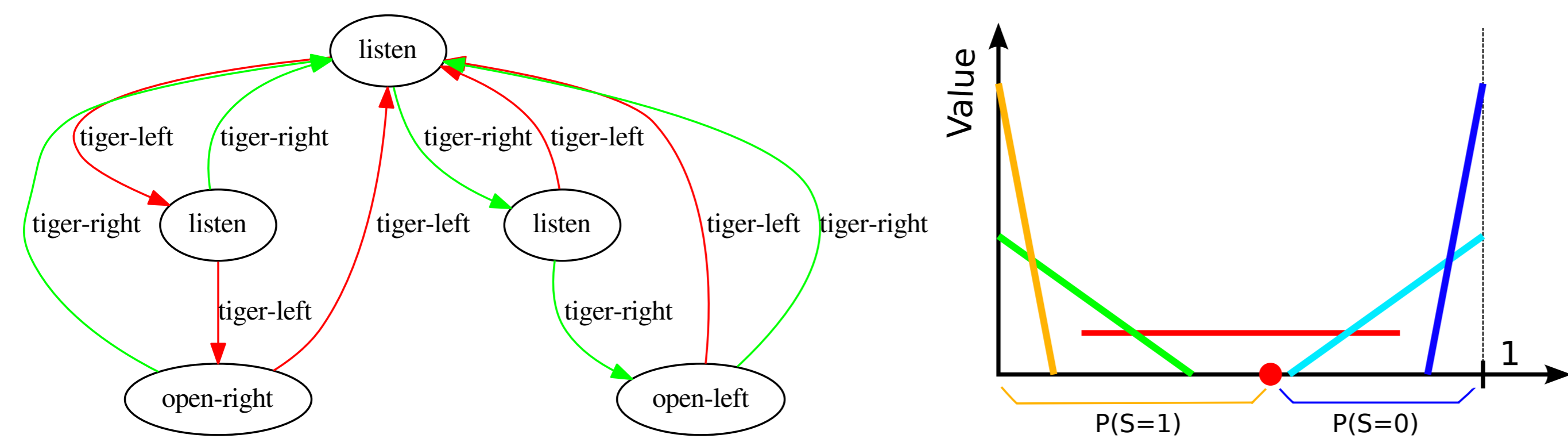
University of
Kent
UNIVERSITY OF
WATERLOO

AAMAS-15

Motivation and Contribution

- Finite-state controllers (FSCs) are the most energy efficient POMDP policies (see Grzes *et al.* "Energy Efficient Execution of POMDP Policies.") which shows their suitability for mobile applications
- Efficient and robust algorithms that compute small policies/controllers become desirable
- We investigate incremental methods that guarantee the escape from local optima
- We push the understanding and the performance of policy iteration for POMDPs to the point that for the first time they are competitive with the state-of-the-art point-based methods

Finite-state Controllers for POMDPs



Node Improvement in Bounded Policy Iteration (BPI)

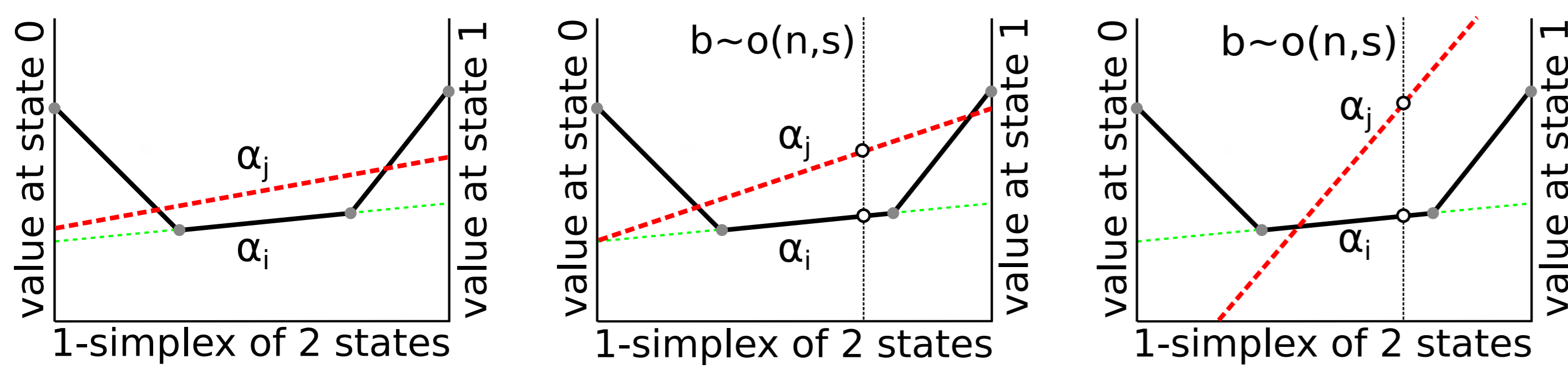


Figure 1: Alpha vectors α_j for improved nodes.

The Need to Escape Local Optima

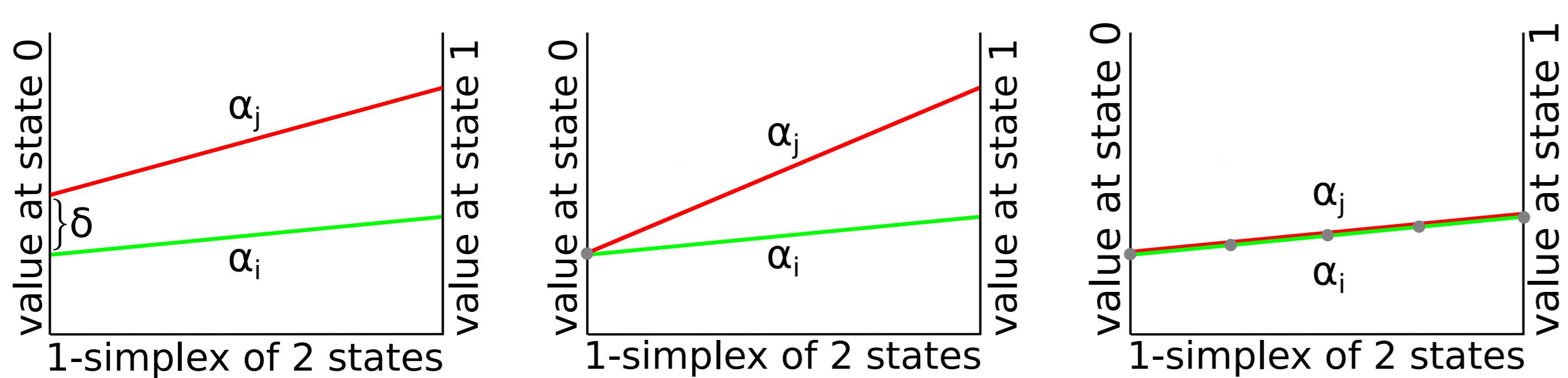


Figure 2: Possible node improvements computed in BPI; α_i corresponds to an old node and α_j to a new node.

Find a New Node with Maximal Improvement Over Entire Belief Simplex

The naïve way to compute the exact DP update for POMDPs is to enumerate all possible alpha vectors. We can compute the set $\Gamma^{a,o}$ of vectors $V_n^{a,o}(s)$ for each $\langle a, o \rangle$ pair by applying a DP update:

$$\Gamma^{a,o} \leftarrow V_n^{a,o}(s) = \frac{R^a(s)}{|O|} + \gamma \sum_{s' \in S} P(o|s', a) P(s'|a, s) V_n^\pi(s'), \forall_n$$

Having $V_n^{a,o}(s)$ for each $\langle a, o \rangle$ pair, we can formulate the following optimisation problem:

$$\begin{aligned} \max: & \sum_{a,n',o,s} w(s) P(n', a|o) V_{n'}^{a,o}(s) - \beta \\ \text{s.t.} & \sum_s w(s) = 1; \sum_{n',a} P(n', a|o) = 1; \\ & \forall_{a,o_1,o_2} \sum_{n_1} P(n_1, a|o_1) = \sum_{n_2} P(n_2, a|o_2) \\ & \forall_n \beta \geq \sum_s w(s) V_n^\pi(s); \\ & \forall_s w(s) \in \mathbf{R}; \forall_{n',a,o} P(n', a|o) \in [0, 1] \end{aligned}$$

Figure 3: A quadratic optimization program to search for a new node that provides maximal improvement at the entire belief simplex; the belief w (witness belief) is the belief at which the improvement happens. Decision variables are the witness belief w , the current value β at belief w , and node parameters $P(n', a|o)$ which, when interpreted as probabilities, correspond to $P(n'|o)P(a)$.

Optimal Solution to the Escape Problem

Theorem

There always exists an optimal solution to the quadratic problem shown in Fig. 3 that is integral $\forall_{n',a,o} P(n', a|o)$, i.e., there exists an optimal solution that corresponds to a deterministic node.

Our Algorithm

$$\begin{aligned} \max: & \sum_{a,n',o,s} y(s, n', a, o) V_{n'}^{a,o}(s) - \beta \\ \text{s.t.} & \sum_s w(s) = 1; \sum_{n',a} P(n', a|o) = 1; \\ & \forall_{a,o_1,o_2} \sum_{n_1} P(n_1, a|o_1) = \sum_{n_2} P(n_2, a|o_2) \\ & \forall_n \beta \geq \sum_s w(s) V_n^\pi(s); \\ & \forall_s w(s) \in \mathbf{R}; \forall_{n',a,o} P(n', a|o) \in \{0, 1\}; \\ & \forall_{s,a,o,n} 0 \leq y(s, a, o, n') \leq P(n', a|o); \\ & \forall_{s,a,o,n} w(s) + P(n', a|o) - 1 \leq y(s, a, o, n') \leq w(s); \end{aligned}$$

■ Real variables ■ Integer variables

Figure 4: The McCormick transformation of the quadratic program in Fig. 3.

- Thanks to the above theorem, McCormick relaxation finds an optimal, deterministic node
- MILP is intractable, but we don't need the optimal solution
- Even a linear relaxation of our MILP can be sufficient (see the paper for interesting properties)

Practical Implementation with Fast Heuristics

Algorithm 1: IPI(-LP): Incremental Policy Iteration for POMDPs.

```

Data: POMDP
Result: FSC for POMDP
1 FSC.N ← {n1}; /* the first node */
2 FSC.N ← FSC.N ∪ {n2}; /* the second node */
3 while impr = true do
4   Policy evaluation
5   for n ∈ FSC.N do
6     impr ← IMPROVENODE(FSC, n); /* DP or LP */
7   if ¬impr then /* escape is required */
8     impr ← ONPOLICYLH(FSC)
9     if ¬impr then
10      impr ← BESTOF(OFFPOLICYLH, SPLIT, CORNER)
11      if ¬impr then
12        impr ← MILPESCAPE(FSC)
13 Prune dominated nodes
    
```

Results

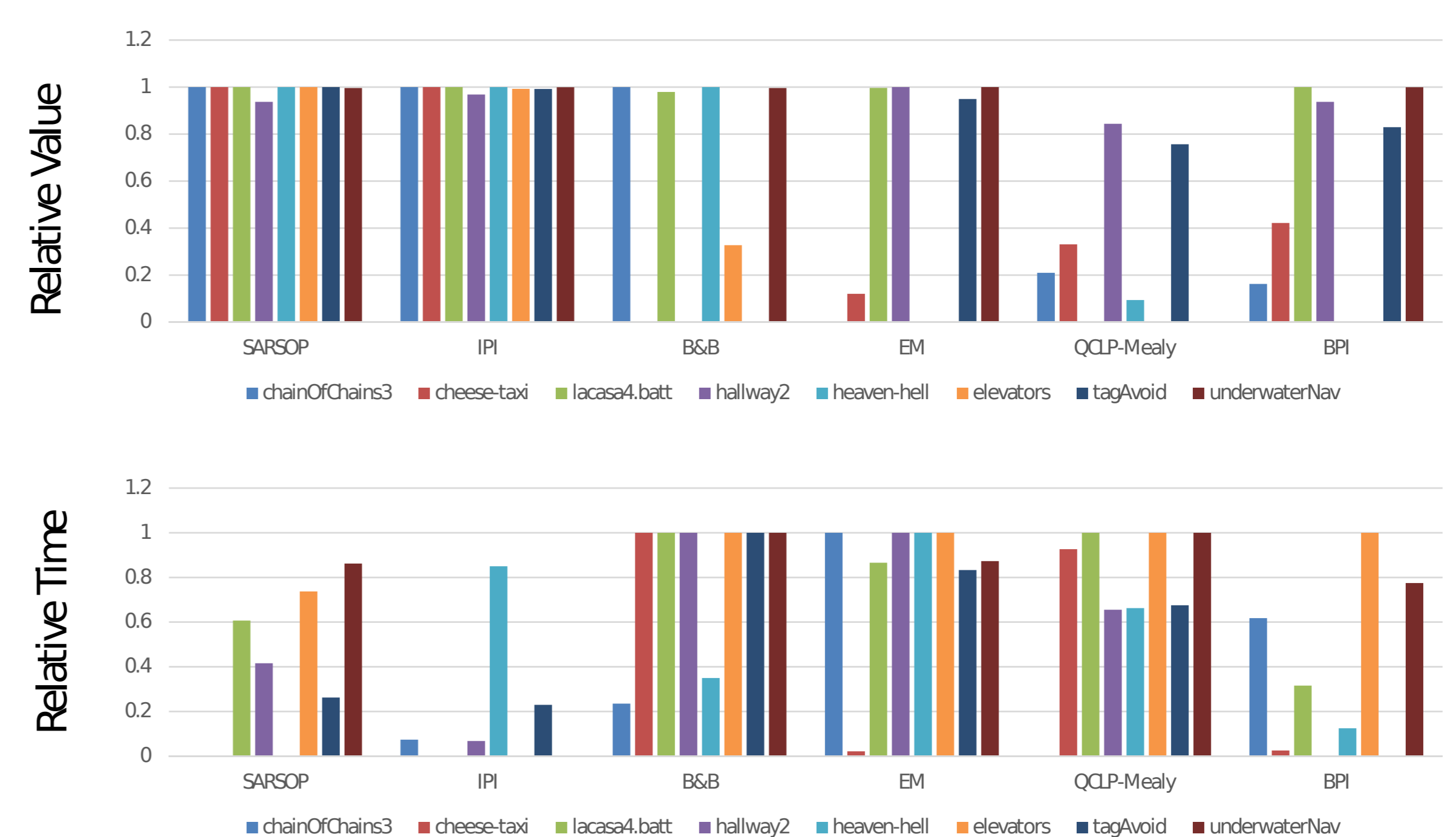


Figure 5: Relative values: normalised values so that the SARSOP upper bound is 1 and the worst value achieved by any algorithm is 0. Relative time: normalised time where the longest time is 1 and the shortest time is 0.

Conclusion

- A new view on principled methods for policy iteration in POMDPs
- A new efficient method for improving individual nodes
- An intuitive explanation of local optima and challenges in escaping it
- A guaranteed method for escape that facilitates fast, anytime execution
- The best node for escape is deterministic
- Heuristic methods analysed (with new connections identified) and used in a practical and well-justified manner

Acknowledgements

This research was sponsored by NSERC and MITACS.