

MULTIGAIN 2.0: MDP controller synthesis for multiple mean-payoff, LTL and steady-state constraints*

Severin Bals
severin.bals@tum.de

Technical University of Munich, Germany

Jan Křetínský
jan.kretinsky@tum.de

Technical University of Munich, Germany
Masaryk University Brno, Czech Republic

Alexandros Evangelidis ✉

alexandros.evangelidis@tum.de

Technical University of Munich, Germany

Jakob Waibel

jakob.waibel@tum.de

Technical University of Munich, Germany

ABSTRACT

We present MULTIGAIN 2.0, a major extension to the controller synthesis tool MULTIGAIN, built on top of the probabilistic model checker PRISM. This new version extends MULTIGAIN’s multi-objective capabilities, by allowing for the formal verification and synthesis of controllers for probabilistic systems with multi-dimensional long-run average reward structures, steady-state constraints, and linear temporal logic properties. Additionally, MULTIGAIN 2.0 can modify the underlying linear program to prevent unbounded-memory and other unintuitive solutions and visualizes Pareto curves, in the two- and three-dimensional cases, to facilitate trade-off analysis in multi-objective scenarios.

CCS CONCEPTS

• **Computing methodologies** → **Control methods**; • **Theory of computation** → **Logic and verification**; • **Mathematics of computing** → **Mathematical optimization**.

KEYWORDS

Markov decision process, quantitative verification, probabilistic model checking, controller synthesis

ACM Reference Format:

Severin Bals, Alexandros Evangelidis ✉, Jan Křetínský, and Jakob Waibel. 2024. MULTIGAIN 2.0: MDP controller synthesis for multiple mean-payoff, LTL and steady-state constraints. In *27th ACM International Conference on Hybrid Systems: Computation and Control (HSCC '24)*, May 14–16, 2024, Hong Kong, Hong Kong. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3641513.3650135>

*This research was supported by the German Research Foundation (DFG) project 427755713 GOPro and the MUNI Award in Science and Humanities (MUNI/1/1757/2021) of the Grant Agency of Masaryk University.

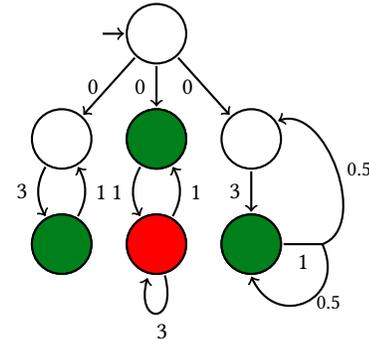
Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

HSCC '24, May 14–16, 2024, Hong Kong, Hong Kong

© 2024 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-0522-9/24/05.

<https://doi.org/10.1145/3641513.3650135>



- linear temporal logic (LTL)

$$G(\bullet \Rightarrow X(\neg \bullet U \bullet))$$

- steady-state constraints (SS)

$$\bullet \geq 0.6$$

- long-run average reward (LRA)

$$\lim_{n \rightarrow \infty} (\inf) \frac{1}{n} \sum_{i=1}^n r(A_i)$$

Figure 1: An MDP and its heterogeneous specification

1 INTRODUCTION

Markov decision processes (MDP), e.g., [14], are the basic model for decision making in uncertain environments. The policy synthesis problem is the problem of resolving the choices so that a given specification is satisfied. In verification, there are many types of properties considered; in this work, we focus on *infinite-horizon properties*. Firstly, *Linear Temporal Logic* (LTL) [13] is mainstream in verification [5]. It can express complex temporal relationships, abstracting from the concrete quantitative timing, e.g., *after every request, there is a grant* (not saying when exactly). Secondly, *Steady-State Policy Synthesis* (SS) [3] constrains the frequency with which states are visited, providing a more quantitative perspective. Recently, it has started receiving more attention also in AI planning [4, 11, 16]. Thirdly, rewards provide a classic framework for quantitative properties. In the setting of infinite horizon, a key role is played by the *long-run average reward* (LRA, a.k.a. mean payoff), e.g., [14], which constrains the reward gained on average per step.

Example 1.1. An example of an MDP with these specifications is shown in Fig. 1. There is a non-trivial choice at the beginning, deciding, intuitively, in which set of states we shall be circulating forever. Such a set is called a maximal end component (MEC). Further, there is another choice in the middle MEC and a probabilistic transition in the right one. The LTL formula in the example specifies that whenever a red state occurs, it is followed by non-red ones until a green one occurs. This can be satisfied in all the MECs of this example; hence they are called *accepting MECs*. The steady-state constraint determines that we stay in green states at least 60% of the time. Finally, the rewards are decorating the edges and then the average reward will be maximized on the “red” self-loop in the middle MEC. If all specifications are considered together, the reward is maximized in the right MEC only, because of the SS constraint. Δ

We consider MDP with the LTL+SS+LRA specifications *combining all these three types*, as introduced and theoretically solved in [11]. We build upon MULTIGAIN [6], a tool extending PRISM [10] with multi-dimensional long-run average reward. Our tool synthesizes a policy maximizing the LRA reward among all policies, ensuring the LTL specification (with the given probability) and adhering to the steady-state constraints.

Our contribution can be summarized as follows:

- We extend MULTIGAIN to analyze an MDP with a heterogeneous LTL+SS+LRA specification for maximizing the long-run average reward under the LTL and steady-state constraints, as described in [11]. Additionally, we extend the specification and the algorithm to cater for further constraints. For example, satisfaction by policies that are deterministic (as in [17]), unichain policies remaining in a single MEC, or policies with a bound on the size of their memory.
- We extend the syntax of the PRISM language slightly to accommodate the richer queries. Further, we produce Pareto frontiers and display the two- and three-dimensional ones, to visualize the trade-offs.
- We conduct a series of experiments to demonstrate the scalability of the tool.

Related tools. To the best of our knowledge, there are no tools that can simultaneously handle multi-dimensional LRA reward computation, LTL, and steady-state specifications. There are, however, two tools that handle multi-dimensional LRA objectives: (i) the previous version of MULTIGAIN implements this functionality through linear programming which is also compatible with the solution offered in [11] and implemented here; and (ii) STORM [15], which implements the same functionality more efficiently through value iteration. Additionally, the Partial Exploration Tool (PET) [12] includes an implementation for LRA reward analysis, by focusing on partial exploration of the state space. However, it does not account for additional objectives such as LTL or steady-state specifications. Furthermore, the work of [17] presents a solution concept for finding deterministic *unichain* policies under LTL and steady-state constraints, however, it does not include any reward structures.

2 FUNCTIONALITY

The main functionality of our tool is to answer multi-objective LRA queries constrained by LTL and steady-state specifications

for MDPs and to synthesize a policy, if possible. We begin with an overview of the tool’s functionality, followed by a description of the various types of queries that are currently supported, including their syntax and semantics. Finally, we discuss additional functionalities that can be accessed via the command-line interface, and highlight key attributes of our tool. The tool and supporting files for the results in the next section are available from [2].

2.1 Workflow

Our tool functions according to the workflow depicted in Fig. 2. The input consists of an MDP defined in the standard PRISM language¹, and an infinite-horizon property. This property is specified using an extension of PRISM’s property specification language² that we developed. PRISM starts by constructing the MDP from the input file and translates the specified LTL property into a Deterministic Rabin Automaton (DRA). Then, it forms the product between the MDP and the DRA, also known as the *product MDP*, which is then passed as an input to the novel component of our tool. Here, an LP is constructed, following the methodology described in [11], and fed to an LP solver. Finally, after the LP is solved, MULTIGAIN 2.0 extracts the solution from the solver and, if required, synthesizes a policy.

2.2 Example

Grid world models have been used extensively for the performance evaluation of various MDP algorithms and tools in fields such as reinforcement learning [9], motion planning [7] and formal verification [17]. Here, we use two-dimensional grids of size $N \times N$, where an agent can traverse between the cells (or states) using one of the four actions, `left`, `down`, `up`, `right`, available in all states. Note that there are no actions to stay in a state i.e., no self loop actions. We show the grid world for $N = 3$ in Fig. 3a. Additionally, two cells are labeled `danger` and one cell is labeled `water_can` to indicate that they are on fire and the presence of a watering can, respectively. The initial state, labeled `home` is the cell at the top left corner. Also, we define a reward structure, denoted as “`extinguish`”, to assign a reward of 1 to the two fire states, intuitively encouraging the agent to repeatedly extinguish resurging flames.

An example query for our 3×3 grid world with all three types of properties is shown below.

```
multi(R{"extinguish"}max=? [S], P>=1 [(! "danger") U
↔ "water_can"], S>=0.25 ["home"])
```

Intuitively, it asks: “*What is the maximum expected long-run average value of reward structure “extinguish” (LRA), such that: (i) the agent does not visit any fire states before collecting the watering can first (LTL) and (ii) at least 25% of the time in the long run it stays “home” (SS)?*”

To compute a result for the query, as discussed in Section 2.1 and shown in Fig. 3b, the product MDP model is constructed, and the maximized LRA reward of 0.5 is returned. Intuitively, the product model consists of three copies of the grid, with a non-accepting MEC (gray grid on the left), which is reached when traversing to

¹<https://www.prismmodelchecker.org/manual/ThePRISMLanguage/Introduction>

²<https://www.prismmodelchecker.org/manual/PropertySpecification/Introduction>

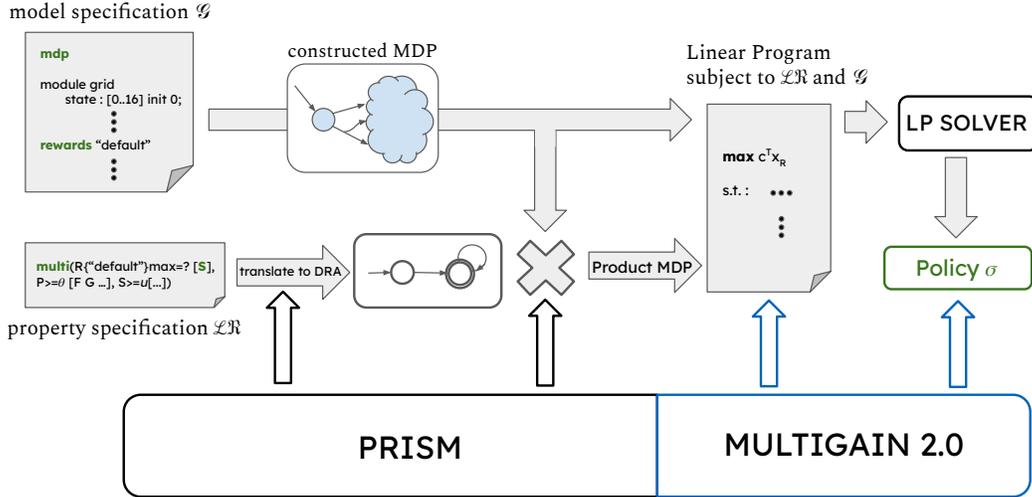
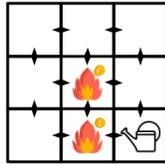
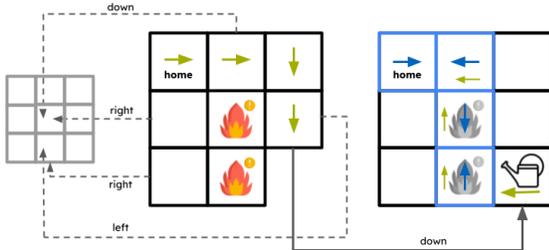


Figure 2: The workflow of MULTIGAIN 2.0 from input specifications to output policy.



(a) An example instance of the adjusted grid world model.



(b) The computed product model and policy. The transient part of the policy is described with green arrows, the recurrent part with blue arrows. The blue framed state have positive switch probability from transient to recurrent behavior.

Figure 3: Example application of MULTIGAIN 2.0 to an adjusted grid world model (a) and the corresponding solution (b).

a fire cell before visiting the watering can, and the unique accepting MEC (right), reached when visiting the watering can cell first. The transient, recurrent and switching behavior of the policy is indicated by arrows in the product model in Fig. 3b. After the transient part directly guides around the fires to the watering can, it continues to traverse to any state with positive frequency in the

long-run. At each such state the policy has a probability to switch to recurrent behavior, which suggests to loop both on the fire and the home cells. It may be noticed that while the four blue states have all positive occupation measure and are all reached by the policy, the recurrent behavior consists of two disconnected cycles here. Such unintuitive results can be avoided using further functionality of the tool.

2.3 Infinite-horizon properties

As described before, a multi-objective query for MULTIGAIN 2.0 consists of the following specifications:

- (1) **Long-run average:** Two types of LRA properties can be specified: (i) a *numerical* property, which seeks to determine the maximum LRA achievable, or (ii) a *Boolean* property that determines whether the LRA surpasses a certain threshold or not. In PRISM’s syntax, a *numerical* or a *Boolean* LRA property could be represented as $R\{\text{"rewardStruct"}\} \max=?[S]$ or $R\{\text{"rewardStruct"}\} \geq 0.5[S]$ respectively.
- (2) **LTL:** The tool only supports a single Boolean query for LTL specifications, since multiple LTL formulae can be conjoined to form one formula. An example could look like $P \geq 0.75 [G F \text{"stateLabel"}]$, which expresses that with probability ≥ 0.75 , states with the label `stateLabel` are reached infinitely often.
- (3) **Steady-state:** An SS property of type $S \leq 0.1[\text{"stateLabel"}]$ requires that the steady-state probability distribution of the states with label `stateLabel` is bounded from above by 0.1.

2.4 Syntax and semantics

As described in [11], there are several types of queries one can formulate by combining the properties discussed above. Moreover, as previously discussed, we extended PRISM’s syntax to allow for

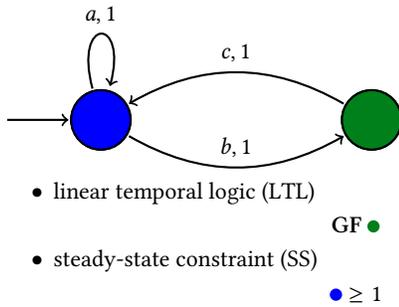


Figure 4: An example MDP with LTL and steady-state constraints. Only policies visiting the accepting state less and less frequently satisfy both specifications, requiring unbounded memory for the information on a current run’s history.

new types of queries using the following notation:

keyword ($[N,]prop_1, prop_2, \dots, prop_k$)

where the keyword can be one among `multi`, `mlessmulti`, `detmulti`, or `unichain` and each $prop_i$ is an LRA, LTL, or SS property. Note that there can be only one LTL property in the syntax of a query. We now explain the semantics of the four different keywords.

2.4.1 multi. The semantics of the `multi` keyword is similar to its meaning in PRISM and MULTIGAIN, i.e., computing a policy that satisfies the conjunction of all the individual properties. Here, the type of result obtained depends on the number of numerical LRA properties. If there are no numerical properties, the query only consists of Boolean LRA, LTL, and SS properties and the result is either *true* or *false*, depending on whether all of them can be satisfied or not. In the case of a single numerical LRA property, the tool returns the maximum (or minimum) achievable value for that LRA reward while satisfying all of the other properties. When more than one numerical property is given, the tool approximates the corresponding Pareto curve while satisfying all of the other properties. An example of a `multi` query is shown below:

```
multi(R{"reward1"}max=? [S], R{"reward2"}max=? [S],
  ↪ R{"reward3"}>=0.5 [S], P>=0.75[G F]
  ↪ "stateLabel1", S>=0.5 ["stateLabel2"],
  ↪ S<=0.5["stateLabel2"])
```

2.4.2 mlessmulti. The `mlessmulti` keyword, used in the previous version of our tool, solves a different problem in its current implementation. In general, a policy computed by the LP (i.e., a `multi` query) might visit the accepting states less and less often to satisfy the remaining constraints, thus requiring unbounded memory as seen in Fig. 4. Relaxing the LRA and SS specifications by an arbitrary factor $\delta > 0$ lifts this restriction [11], such that a finite-memory policy exists for the model.

The implementation of a `mlessmulti` query addresses this problem from a different angle. It allows the user to specify an additional integer N , signifying the maximum number of steps on (the long-run) average before an accepting state is revisited. The tool subsequently computes and outputs the resulting minimal factor δ to uniformly relax all steady-state specifications and long-run

average rewards, with regards to this fixed accepting frequency. Hence, exporting the strategy yields a finite-memory policy, more specifically a 2-memory policy [11] consisting of a memoryless transient policy, which switches to a memoryless recurrent policy. Note that due to the modified objective function it is not possible to define numerical LRA properties in a `mlessmulti` query.

An example of a `mlessmulti` query is shown below:

```
mlessmulti(1000, R>=0.5 ["sLabel"], P>=1[G F]
  ↪ "tLabel", S>=1 ["sLabel"])
```

2.4.3 detmulti. Depending on the underlying model and property specification, the policy computed by the `multi` keyword may exhibit two significant characteristics. Firstly, it is typically randomizing, and secondly, the policy may require an infinite amount of memory to remember the current history. To address both of these issues, the `detmulti` keyword implements the approach by [17], which is based on a mixed-integer linear program. The resulting policy, which is defined over the original MDP rather than the product, is both deterministic and finite-memory. The `detmulti` queries may contain a single LTL property and arbitrarily many steady-state specifications. The result, other than an exportable policy, is either the optimal LRA reward or a boolean value indicating whether a solution was found or not. An example of a `detmulti` query is shown below:

```
detmulti(P>=0.75 [(! "stateLabel1") U "stateLabel2"],
  ↪ S>=0.75 ["stateLabel3"])
```

2.4.4 unichain. We introduce the keyword `unichain`, which computes a *unichain* solution for the `multi` query, i.e., the recurrent behavior of the policy resides only in a single MEC and thus can be turned into a “single” behavior happening with probability 1. Formally, a policy is called *unichain* if the induced Markov chain has only one recurrent class and all the other states are transient. This is computed by exploring each MEC (or accepting MEC if an LTL specification is present) individually. The implementation concept follows the idea presented in [11, Section 6]. If no numerical LRA properties are specified, the tool explores the MECs until a *unichain* solution is found and outputs the corresponding boolean value. For a single numerical LRA property, our tool searches for the *unichain* solution maximizing (or minimizing) the reward structure and outputs the corresponding reward. Multiple numerical rewards are not allowed for this keyword, as this would result in comparing multiple Pareto curves. An example of a `unichain` query is shown below:

```
unichain((R{"reward1"}max=? [S], R{"reward2"}>=0.5
  ↪ [S], P>=0.75[G F "stateLabel1"], S>=0.5
  ↪ ["stateLabel2"], S<=0.5["stateLabel2"])
```

2.5 Interface

The tool is used via a command line interface, which requires the user to specify two files as input arguments, containing the model and the queries. The approximated Pareto curve can be exported to a file by using the flag `--exportpareto`. Furthermore, the tool includes a Python script that enables the visualization of Pareto frontiers with two or three dimensions. In Fig. 5, we show example plots of two- and three-dimensional Pareto frontiers produced by

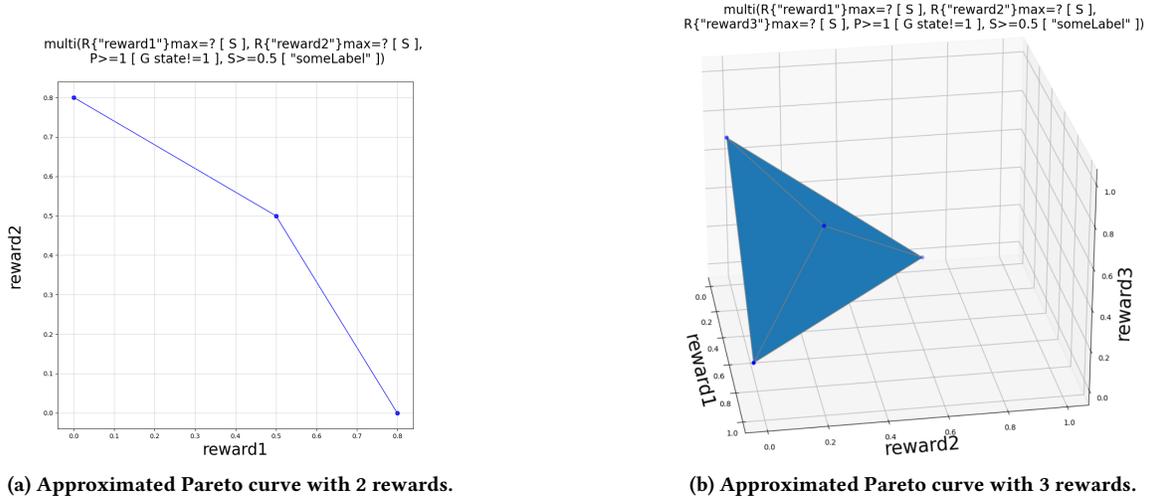


Figure 5: Example plots of approximated Pareto curves with 2 and 3 dimensions.

the tool. Moreover, for all queries except Pareto approximation, we provide the option to export the computed policy, which may have an unbounded memory, to a file in various formats.

2.6 Implementation characteristics

In this section, we report on the quality of MULTIGAIN 2.0 by highlighting some of its key characteristics.

Extensibility. The underlying LP solver implements a general interface and can thus be easily switched for every run. This implementation allows the simple extension and addition of further LP solvers. Currently, the tool supports the use of `lp_solve` [1] and `Gurobi` [8]. After solving the LP, the tool extracts the solution from the solver and, if required, synthesizes a policy.

For approximating Pareto curves a new generic class has been implemented which takes as input a weight function, mapping weights to reward structures. This class lifts the Pareto curve approximation from MULTIGAIN 2.0 so that other PRISM-based tools could utilize it. Furthermore, Pareto curves of any dimension can be approximated, contrary to the two-dimensional limit of the previous version of the tool. Since the tool is implemented in the unifying approach of the PRISM pipeline, it can be extended at a variety of entry points, as seen in Fig. 2. For example, new deterministic automata could be implemented alongside the translation of the LTL and building the product model, without changing the tool's core functionality.

3 EXPERIMENTAL EVALUATION

In this section, we assess the performance of our tool in terms of its ability to solve the types of queries described in Section 2.4. We conducted multiple experiments to evaluate the performance of our tool. We first discuss the experimental setup, followed by the technical details regarding our experiments, and then we give a detailed overview of our experimental results in Section 3.1.

Experimental setup. Our evaluation consists of three parts: (i) the evaluation of the full property suite (LRA, SS, and LTL properties) using a grid world model; (ii) a scalability analysis of the tool

Table 1: Average running time (in seconds) over 20 randomly grid world labeled instances.

LTL	LRA $R_{\max=?}^{\text{rew}_c} [S]$	Average running time for each grid				
		4×4	16×16	32×32	64×64	128×128
$G(\neg b) \wedge (GFa)$	×	0.121	0.231	0.466	1.296	26.104
$(GFa) \vee (FGb)$	×	0.029	0.091	0.210	0.581	2.498
$(Fa)Ub$	×	0.020	0.074	0.238	0.852	4.042
$(Fa) \wedge (Fb) \wedge (Fc)$	×	0.042	0.147	0.499	2.396	21.678
$G(\neg b) \wedge (GFa)$	✓	0.025	0.092	0.623	10.12	128.197
$(GFa) \vee (FGb)$	✓	0.021	0.081	0.555	5.245	106.772
$(Fa)Ub$	✓	0.013	0.105	0.668	11.053	158.448
$(Fa) \wedge (Fb) \wedge (Fc)$	✓	0.017	0.298	3.232	81.389	883.446

regarding its performance with an increasing number of steady-state constraints; and (iii) an evaluation of how different LP solvers impact the tool's efficiency, including both runtime performance and memory usage, in the context of handling queries.

Technical details. All experiments were performed on a desktop computer with 16 GB of RAM and an Intel i7-8550U CPU @ 1.80GHz, running Ubuntu 22.04.3 LTS.

For the grid world model, the average running time over 20 runs was recorded, as a countermeasure to the high variance of individual running times. All results are rounded to three decimal places.

3.1 Results

LRA+LTL+SS queries. In Table 1, we present the results for various multi queries that involve the combination of all three types of properties. These queries are categorized into two groups: those containing an LRA property, denoted by a ✓ symbol, and those that do not, represented by a × symbol. Following the experiments in [17], the states were randomly divided into four equally-sized subsets, and each subset was labeled with an atomic proposition from the set $AP = a, b, c, d$ before a run of the tool. Additionally, we

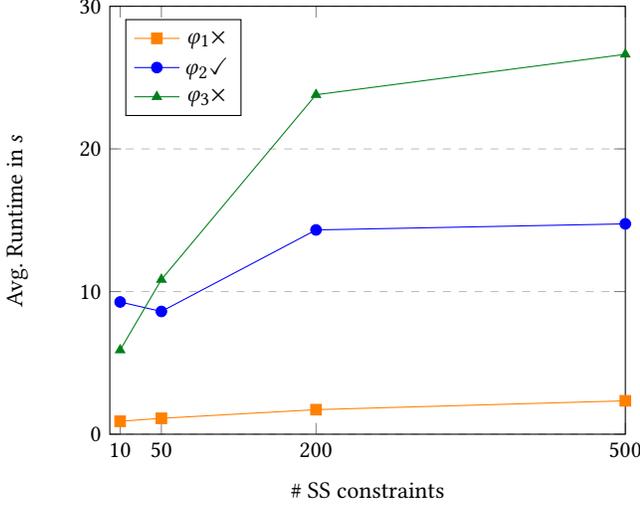


Figure 6: LP solver running times of multi queries on 64×64 grid world models with the LTL formulae φ_1 , φ_2 and φ_3 from above, based on the number of steady-state constraints specified.

created a reward structure, rew_c , that assigns a reward of 1 to each state labeled with c . In each run, we required the LTL formula to be fulfilled with a probability threshold of $\theta = 0.5$ and steady-state constraints of $S \geq 0.01[\"d\"]$, $S \leq 0.5[\"d\"]$.

The upper half of Table 1 highlights the efficient performance of our tool when no LRA property is specified. Even for the largest grid world MDP model, the longest running time is still only a few seconds. Also, when an LRA property is specified (lower half of Table 1), for the majority of the cases in the 64×64 grid, it remains quite efficient (< 11 seconds); however, scalability issues start emerging for the 128×128 grid. The LRA property used here maximizes the average reward of rew_c , denoted as $R_{\max=?}^{rew_c}[S]$. For larger grid sizes, such as the 128×128 case, the runtimes are reasonable except for the last LTL property $(Fa) \wedge (Fb) \wedge (Fc)$. This trend aligns with the findings in [17] and can be attributed to the random generation procedure of the grid world instances, which may have a bias toward certain types of models.

Scaling the number of steady-state constraints. In the next set of experiments, we systematically assess the computational overhead as a function of the number of steady-state constraints introduced per specification. We consider the following three different LTL formulae: $\varphi_1 = (GFa) \vee (FGb)$, $\varphi_2 = (Fa)Ub$, and $\varphi_3 = Fa \wedge Fb \wedge Fc$, and we employ instances of a 64×64 grid world for our experimental setup. To instantiate non-trivial steady-state constraints, for each experimental iteration, we stochastically select a subset of states, denoted \mathcal{S} , that are labeled with d . Every state $s \in \mathcal{S}$ is then attributed a distinct label, denoted l_s . Formally, for a given label l_s , its corresponding constraint is such that: $U(l_s) = 0.5$, and $\sum_{s \in \mathcal{S}} L(l_s) \leq 0.25$, where U and L denote an upper and lower bound, respectively. For each of these designated labels l_s , we append a steady-state constraint to the property specification, ensuring that the cumulative lower bounds do not exceed a threshold,

Table 2: Average LP solver runtimes (in seconds) over 20 runs of respective grid world instances, recorded using Gurobi and lp_solve. The faster runtime of each problem is marked in green.

$R_{\max=?}^{rew_c}[S]$ (LRA)	LTL	Solver	Average running time per grid size			
			4×4	16×16	32×32	64×64
\times	φ_1	Gurobi	0.002	0.029	0.062	0.264
		lp_solve	0.001	0.072	0.33	1.165
	φ_2	Gurobi	0.007	0.018	0.09	0.405
		lp_solve	0.001	0.069	1.502	1.617
	φ_3	Gurobi	0.002	0.027	0.233	2.006
		lp_solve	0.001	0.097	1.369	15.224
\checkmark	φ_1	Gurobi	0.006	0.038	0.342	4.566
		lp_solve	0.002	0.059	0.206	1.856
	φ_2	Gurobi	0.005	0.097	0.284	10.022
		lp_solve	0.002	0.071	0.412	3.842
	φ_3	Gurobi	0.003	0.014	2.311	79.566
		lp_solve	0.001	0.595	3.043	160.526

thereby reducing the likelihood of encountering infeasible scenarios.

Our experiments span configurations with 10, 50, 200, and 500 steady-state constraints, evaluated against the three distinct LTL formulae. To account for the variance in individual running times, as discussed in previous experiments, the recorded running times were averaged over 20 runs. We note that for φ_1 , the solver’s average runtime increases as more steady-state constraints are appended. Specifically, starting from an average runtime of ≈ 1 second with 10 constraints, it rose to 2.34 seconds with 500 constraints. On the other hand, φ_2 presented an interesting pattern as the average runtime of ≈ 9 seconds for 10 and 50 constraints, were fairly similar. However, a significant increase was observed as we introduced 200 constraints, reaching ≈ 14 seconds, and this growth seemed to stabilize by the time we integrated 500 constraints. Finally, φ_3 seemed to be the most computationally demanding, starting at ≈ 6 seconds with 10 constraints and reaching ≈ 27 seconds at 500 constraints.

This experiment demonstrates that adding steady-state constraints does not have a significant impact on overall runtime, and is therefore not a restriction on the user.

LP solver comparison. In this set of experiments, we evaluate how the performance of the queries is affected by the choice of the underlying LP solver. We consider the three LTL formulae used in the previous experiments. As stated in Section 2.6, MULTI-GAIN 2.0 supports the publicly available solver lp_solve as well as the well-known commercial state-of-the-art solver Gurobi. In Table 2 we present the average runtime over 20 runs of Gurobi and lp_solve on various grid world instances, with and without LRA maximization, marked by the \checkmark and \times symbols, respectively. The fastest runtime for each problem is marked in green.

Our results show a trend where the Gurobi significantly outperforms lp_solve as the size of the grid increases, especially on the 64×64 grid. However, even in this case, there are instances in which lp_solve performs reasonably well. For example, with the LTL formula φ_1 , when LRA reward maximization is considered, lp_solve’s average runtime of 1.856 seconds outperforms Gurobi’s of ≈ 4.5 seconds. Moreover, for the φ_3 formula under the same grid configuration and with LRA maximization included, Gurobi is $\approx 50\%$ faster compared to lp_solve, whereas in the case

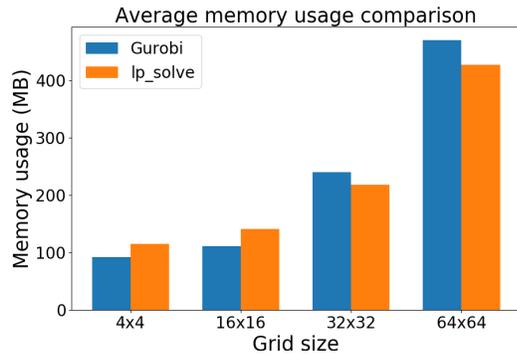


Figure 7: Average memory usage (in megabytes) over 20 runs of respective grid world instances, recorded using Gurobi and lp_solve.

without LRA maximization, Gurobi demonstrates $\approx 87\%$ reduction in runtime over lp_solve.

On the other hand, in smaller grid sizes, lp_solve becomes more competitive and in some cases it even outperforms Gurobi. For instance, when evaluating the φ_1 formula without LRA maximization on a 4×4 grid, lp_solve achieves a runtime of 0.001 seconds compared to Gurobi's 0.002 seconds, denoting a 50% more efficiency in handling smaller grid sizes.

Apart from the runtime performance of both solvers, we also compared the quality of the solution on selected experiments. Both solvers exhibited no noticeable issues in terms of finding a solution and solution quality.

In Fig. 7 we report on the average memory usage of our tool when using the Gurobi and lp_solve solvers, across different grid sizes. As expected, there is an increasing trend in terms of the tool's memory usage as the grid size increases, indicating greater memory requirements for solving larger problems. We note, that for the majority of the cases, the tool consumes more memory when using lp_solve compared to Gurobi, something which can be observed in smaller grid sizes. However, across every grid size the difference in memory usage between the two solvers is not significant, indicating that both solvers are viable alternatives from a memory consumption perspective.

4 CONCLUSION

We presented MULTIGAIN 2.0, an MDP controller synthesis tool for multiple long-run average reward structures subject to LTL and steady-state constraints. Apart from the normal combination of these different objectives, it is also able to solve the δ -satisfaction problem, which relaxes the objectives by a small factor δ . We also implemented a new method sketched in [11] providing *unichain* solutions, and a method described in [17] for *deterministic* solutions. Our tool can export the Pareto curve and the policy, and it can also visualize two and three-dimensional Pareto curves.

This tool can be further extended to work for *omega-regular* objectives. Another useful direction for future work would be to explore combinations with other types of properties, such as non-linear, *finite-horizon* or *discounted*, rewards.

ACKNOWLEDGMENTS

The authors would like to thank Ismail R. Alkhouri for the comprehensive guidance on their approach in [17] and provision of example models, as well as Ayse Aybüke Ulusarslan for their initial help with the project. Further, the authors would like to thank the original developers of MULTIGAIN for allowing us to use the tool's name for our extension.

REFERENCES

- [1] [n. d.]. <https://sourceforge.net/projects/lpsolve/>
- [2] [n. d.]. Supporting material. <http://www.multigain.github.io/multigain2/>.
- [3] S. Akshay, Nathalie Bertrand, Serge Haddad, and Loïc Hélouët. 2013. The Steady-State Control Problem for Markov Decision Processes. In *QEST (Lecture Notes in Computer Science, Vol. 8054)*. Springer, 290–304.
- [4] George K. Atia, Andre Beckus, Ismail Alkhouri, and Alvaro Velasquez. 2020. Steady-State Policy Synthesis in Multichain Markov Decision Processes. In *IJCAI*. ijcai.org, 4069–4075.
- [5] Christel Baier and Joost-Pieter Katoen. 2008. *Principles of Model Checking*. MIT Press.
- [6] Tomáš Brázdil, Krishnendu Chatterjee, Vojtěch Forejt, and Antonín Kučera. 2015. MultiGain: A Controller Synthesis Tool for MDPs with Multiple Mean-Payoff Objectives. In *Tools and Algorithms for the Construction and Analysis of Systems*, Christel Baier and Cesare Tinelli (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 181–187.
- [7] Meng Guo and Michael M. Zavlanos. 2018. Probabilistic Motion Planning Under Temporal Tasks and Soft Constraints. *IEEE Trans. Automat. Control* 63, 12 (2018), 4051–4066. <https://doi.org/10.1109/TAC.2018.2799561>
- [8] Gurobi Optimization, LLC. 2023. Gurobi Optimizer Reference Manual. <https://www.gurobi.com>
- [9] L. P. Kaelbling, M. L. Littman, and A. W. Moore. 1996. Reinforcement learning: A survey. *Journal of Artificial Intelligence Research* 4 (1996), 237–285.
- [10] Marta Z. Kwiatkowska, Gethin Norman, and David Parker. 2011. PRISM 4.0: Verification of Probabilistic Real-Time Systems. In *CAV (LNCS, Vol. 6806)*. 585–591. https://doi.org/10.1007/978-3-642-22110-1_47
- [11] Jan Křetínský. 2021. LTL-Constrained Steady-State Policy Synthesis. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, Zhi-Hua Zhou (Ed.). International Joint Conferences on Artificial Intelligence Organization, 4104–4111. <https://doi.org/10.24963/ijcai.2021/565> Main Track.
- [12] Tobias Meggendorfer. 2022. PET – A Partial Exploration Tool For Probabilistic Verification. In *Automated Technology for Verification and Analysis: 20th International Symposium, ATVA 2022, Virtual Event, October 25–28, 2022, Proceedings*. Springer-Verlag, Berlin, Heidelberg, 320–326. https://doi.org/10.1007/978-3-031-19992-9_20
- [13] Amir Pnueli. 1977. The Temporal Logic of Programs. In *FOCS*. 46–57.
- [14] Martin L. Puterman. 1994. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley.
- [15] Tim Quatmann and Joost-Pieter Katoen. 2021. Multi-objective Optimization of Long-run Average and Total Rewards. In *TACAS (1) (Lecture Notes in Computer Science, Vol. 12651)*. Springer, 230–249.
- [16] Alvaro Velasquez. 2019. Steady-State Policy Synthesis for Verifiable Control. In *IJCAI*. ijcai.org, 5653–5661.
- [17] Alvaro Velasquez, Ismail Alkhouri, Andre Beckus, Ashutosh Trivedi, and George Atia. 2022. Controller Synthesis for Omega-Regular and Steady-State Specifications. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (Virtual Event, New Zealand) (AAMAS '22)*. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1310–1318.