

# Solving Constrained Stochastic Shortest Path Problems with Scalarisation

Johannes Schmalz<sup>a,b,\*</sup> and Felipe Trevizan<sup>b</sup>

<sup>a</sup>Saarland University, Germany

<sup>b</sup>Australian National University, Australia

ORCID (Johannes Schmalz): <https://orcid.org/0000-0003-3389-2490>, ORCID (Felipe Trevizan):

<https://orcid.org/0000-0001-5095-7132>

## Abstract.

Constrained Stochastic Shortest Path Problems (CSSPs) model problems with probabilistic effects, where a primary cost is minimised subject to constraints over secondary costs, e.g., minimise time subject to monetary budget. Current heuristic search algorithms for CSSPs solve a sequence of increasingly larger CSSPs as linear programs until an optimal solution for the original CSSP is found. In this paper, we introduce a novel algorithm CARL, which solves a series of unconstrained Stochastic Shortest Path Problems (SSPs) with efficient heuristic search algorithms. These SSP subproblems are constructed with scalarisations that project the CSSP’s vector of primary and secondary costs onto a scalar cost. CARL finds a maximising scalarisation using an optimisation algorithm similar to the subgradient method which, together with the solution to its associated SSP, yields a set of policies that are combined into an optimal policy for the CSSP. Our experiments show that CARL solves 50% more problems than the state-of-the-art on existing benchmarks.

## 1 Introduction

Stochastic Shortest Path Problems (SSPs) model problems where the effects of actions can be probabilistic and Constrained SSPs (CSSPs) extend the model by allowing constraints over secondary costs, enforcing that the incurred secondary costs do not exceed a specified threshold over expectation. This can represent interesting real-world problems, e.g., planning an aeroplane’s route that minimises fuel usage while avoiding bad weather and satisfying certain timing requirements [14]. SSPs can be solved with the primal formulation, where we optimise over variables that represent the agent’s cost-to-go (a.k.a. value functions); or with the dual formulation over occupation measures, which represent the expected number of times an action is applied in each state. Algorithms over the primal representation tend to be faster and indeed, the state-of-the-art algorithms for solving SSPs optimally use the primal representation, e.g., iLAO\* [18], LRTDP [8], CG-iLAO\* [26]. In contrast, the only heuristic search algorithms for CSSPs use the dual representation: i-dual [31] and i<sup>2</sup>-dual [32]. In this paper, we present the first heuristic search algorithm that solves CSSPs optimally using the primal space, and show that it outperforms the existing algorithms.

Our algorithm CARL builds on existing work for finding deterministic policies for CSSPs (which are generally suboptimal) [22], with

the key difference that CARL finds optimal (potentially stochastic) policies. CARL works by running heuristic search in the primal space using scalarisation, where a scalarisation  $\lambda$  projects the CSSP’s cost vector onto a scalar cost function, thereby inducing an unconstrained SSP. CARL searches for an optimal scalarisation  $\lambda^*$  with optimisation techniques similar to the subgradient method [30] and solves SSPs induced by the encountered  $\lambda$  scalarisations as subproblems. Once CARL has found  $\lambda^*$ , it computes the optimal costs-to-go  $V^*$  associated with that scalarisation, and builds an optimal policy for the CSSP by combining all of  $V^*$ ’s greedy policies. CARL’s subproblems are solved efficiently with heuristic search algorithms for SSPs and its outer problem over scalarisations is blind; this is reversed in the dual methods (i-dual and i<sup>2</sup>-dual) where the construction of subproblems is guided with a heuristic, but the subproblems themselves are solved blindly by the LP solver. In our experiments, this tradeoff pays off and CARL can solve 1264 out of 1290 of our benchmark problems, whereas the state-of-the-art dual methods only solve 808. For some problems CARL is able to solve all instances while the state-of-the-art is not able to solve any, and over the problems where all algorithms solve all instances CARL offers an average speedup of  $10\times$  w.r.t. the strongest baseline.

## 2 Background

**Stochastic Shortest Path Problems (SSPs)** [4] are defined by the tuple  $\langle S, s_I, G, A, P, C \rangle$  where  $S$  is a finite set of states;  $s_I \in S$  is the initial state;  $G \subset S$  is a set of goal states that must be reached;  $A$  is a finite set of applicable actions, and we write  $A(s)$  to denote the applicable actions in state  $s$ ;  $P$  is the probability transition function where  $P(s'|s, a)$  is the probability of reaching  $s'$  after applying  $a$  to  $s$ ; and  $C : A \rightarrow \mathbb{R}_{>0}$  is a cost function where  $C(a)$  gives the cost of applying the action  $a$ .

**Policies** map states to actions and describe solutions to SSPs. Policies come in two flavours: **deterministic policies**  $\pi : S \rightarrow A$  and **stochastic policies**  $\pi : S \rightarrow \text{distr}(A)$ , which respectively map each state onto a single action, or onto a probability distribution over actions from which an action should be selected randomly;  $\pi(s, a)$  is the probability that  $\pi$  applies  $a$  in  $s$ , and we may write  $\pi(s) = a$  if  $\pi(s, a) = 1$ . The envelope of policy  $\pi$  from  $s$ , written  $S^{\pi, s}$ , denotes the set of states reachable by following  $\pi$  from  $s$ . The policy  $\pi$  is closed w.r.t.  $s$  if  $\pi$  is defined for all  $s' \in S^{\pi, s}$ , and  $\pi$  is proper w.r.t.  $s$  if following  $\pi$  from  $s$  guarantees that  $G$  is reached with proba-

---

\* Corresponding Author. Email: [schmalz@cs.uni-saarland.de](mailto:schmalz@cs.uni-saarland.de).

bility 1. If  $\pi$  is not closed or not proper it is called open or improper, respectively. If  $s$  is omitted in these expressions, then we assume  $s = s_I$ . We extend the notion of an envelope, and write  $\text{supp}(\pi)$  to denote the state-action pairs that might be encountered by  $\pi$ , i.e.,  $\text{supp}(\pi) = \{(s, a) : s \in S^\pi, \pi(s, a) > 0\}$ . We overload the  $C$  symbol and write  $C(\pi)$  to denote the expected cost incurred by following the proper policy  $\pi$  from  $s_I$ . We make two standard assumptions for SSPs: a proper policy exists from each state (reachability), and any improper policy incurs infinite cost. Then, a policy  $\pi$  is optimal if it is closed and minimises  $C(\pi)$ . Under these assumptions, SSPs always have an optimal deterministic policy [4].

**Primal algorithms for SSPs** work over value functions  $V : S \rightarrow \mathbb{R}_{\geq 0}$ , which represent the cost-to-go for each state, i.e.,  $V(s)$  indicates the expected cost that the agent must incur to reach the goal from  $s$ . For each  $V$ , its  $Q$ -values are  $Q(s, a) = C(a) + \sum_{s' \in S} P(s'|s, a)V(s')$ .  $V^*$  is the optimal value function and denotes the cheapest possible cost-to-go for each state.  $V^*$  is the unique solution of the Bellman equations:

$$V(s) = \min_{a \in A(s)} Q(s, a) \quad \forall s \in S \setminus G \text{ and } V(s) = 0 \quad \forall s \in G.$$

For  $V$ , its greedy policy  $\pi_V$  is defined by  $\pi_V(s) = \text{argmin}_{a \in A(s)} Q(s, a)$ , where the argmin operator breaks ties arbitrarily to yield a single policy.  $V^*$ 's greedy policy gives an optimal deterministic policy for the SSP. To compute  $V^*$ , we can start with some  $V$  and iteratively apply Bellman backups  $V(s) \leftarrow \min_{a \in A(s)} Q(s, a)$  over states  $s$ . Value Iteration (VI) [3] applies Bellman backups over the whole state space in each step, and converges to  $V^*$  in the limit as the number of steps increases.

**Heuristic search algorithms** use heuristic functions to focus on a promising subset of states and actions, and only apply backups there. For deterministic problems the canonical heuristic-search algorithm is  $A^*$  [20], and the state-of-the-art algorithms iLAO\* [18], LRTDP [8], and CG-iLAO\* [26] generalise  $A^*$  to the SSP setting. A heuristic function is a value function  $H : S \rightarrow \mathbb{R}_{\geq 0}$  that estimates  $V^*$ , so that the algorithm can avoid states with large  $H(s)$  and prefer states with low  $H(s)$ . We require heuristics to be admissible, i.e., to be lower bounds on the optimal value function  $H(s) \leq V^*(s) \quad \forall s \in S$ . With admissible heuristics the heuristic-search algorithms we consider can guarantee optimality. For VI, we measure the change in successive value functions with Bellman residual  $\text{RES}(s) = |V(s) - \min_{a \in A(s)} Q(s, a)|$ , and stop search once the changes in  $V$  are sufficiently small. In heuristic search, we stop when  $V$  is  $\epsilon$ -consistent [7], i.e., when for some greedy policy  $\pi_V$  we have  $\text{RES}(s) \leq \epsilon \quad \forall s \in S^{\pi_V}$ . This condition guarantees that  $V = V^*$  as  $\epsilon \rightarrow 0$ , and in practice, for sufficiently small  $\epsilon$ ,  $\epsilon$ -consistent value functions induce an optimal policy.

**CG-iLAO\*** [26] is a heuristic search algorithm over the primal space, which we explain because it provides crucial features to solve our subproblems efficiently. CG-iLAO\* is presented in alg. 1. It works over value function  $V$  and a partial SSP  $\hat{S}$ , which contains a subset of the original SSP's states and actions. Non-goal states in  $\hat{S}$  are called "expanded" if they have at least one applicable action in  $\hat{S}$ , and are called "fringes" if they have none. In each step, CG-iLAO\* applies Bellman backups over the greedy policy  $\hat{\pi}_V$ 's envelope  $\mathcal{E}$ , and expands  $\mathcal{E}$ 's fringes. Thus, CG-iLAO\* simultaneously works towards making  $\hat{\pi}_V$  closed by expanding fringes, and towards making  $V$   $\epsilon$ -consistent with the Bellman backups. Since  $\hat{S}$  need not contain all actions, CG-iLAO\* looks for expanded states  $s$  with actions  $a$  such that  $Q(s, a) < V(s)$ , because this suggests  $a$  can improve  $V(s)$  and should be added if it is missing.  $\Gamma$  is a set of all state-action pairs

---

**Algorithm 1: CG-iLAO\***

---

```

1 Function CG-iLAO* (SSP  $\mathbb{S}$ , heuristic  $H$ ,  $\epsilon \in \mathbb{R}_{>0}$ )
2    $\hat{S} \leftarrow$  partial SSP containing only  $s_I$ 
3    $V \leftarrow$  value function initialised by  $H$ 
4   repeat
5      $\hat{\pi}_V, \hat{\pi}_{\text{old}} \leftarrow$  greedy policy for  $V$ , restricted to  $\hat{S}$ 
6      $\mathcal{E} \leftarrow$  post-order DFS traversal of  $\hat{\pi}_V$  from  $s_I$ 
7      $\hat{S}, \mathcal{E}, \hat{\pi}_V \leftarrow$  expand fringes of  $\mathcal{E}$ 
8      $V, \text{RES}, \hat{\pi}_V, \Gamma \leftarrow$  apply Bellman backups over  $\mathcal{E}$ ,
       recording increases and decreases to  $V$  in  $\Gamma$ 
9      $V, \text{RES}, \Gamma, \hat{S} \leftarrow$  fix cases of  $V(s) > Q(s, a)$  in  $\Gamma$ 
10    until  $\text{fringes}(\mathcal{E}) = \emptyset$  and  $\hat{\pi}_{\text{old}} = \hat{\pi}_V$  and  $\text{RES} \leq \epsilon$ 
11    return  $V$ 

```

---

that potentially satisfy  $Q(s, a) < V(s)$ , i.e., it is a superset of improving actions.  $\Gamma$  is efficiently maintained by tracking changes to  $V$  and line 9 fixes any issues by setting  $V(s) \leftarrow \min\{V(s), Q(s, a)\}$  for all  $(s, a) \in \Gamma$  and adding  $a$  as required. Note that this may introduce more instances of  $V(s) > Q(s, a)$ , which are themselves recorded in  $\Gamma$  to be checked later.

**Constrained SSPs (CSSPs)** [1, 31] are extensions to SSPs defined by  $\mathbb{C} = \langle S, s_I, G, A, P, \mathbf{C}, \mathbf{u} \rangle$  with two changes from SSPs: (1) the vector cost function  $\mathbf{C} : A \rightarrow \mathbb{R}_{>0} \times \mathbb{R}_{\geq 0}^n$ , where  $C_0 : A \rightarrow \mathbb{R}_{>0}$  is the primary cost, and  $C_i : A \rightarrow \mathbb{R}_{\geq 0}$  for  $i \in \{1, \dots, n\}$  are the secondary costs; (2) we have a vector of upper bounds  $\mathbf{u} \in \mathbb{R}_{\geq 0}^n$ . Note that we use boldface to denote vectors.  $C_i(\pi)$  denotes the expected cost of  $\pi$  on the  $i^{\text{th}}$  cost; we assume reachability and that improper policies incur infinite primary cost. A policy  $\pi$  is feasible if it satisfies all its secondary-cost constraints, i.e.,  $C_i(\pi) \leq u_i \quad \forall i \in \{1, \dots, n\}$ , and  $\pi$  is optimal if it is closed, feasible, and minimises the primary cost  $C_0(\pi)$  w.r.t. the other feasible policies. In contrast to SSPs, CSSPs may have no optimal deterministic policies and only strictly stochastic ones [1, 31]. This difference occurs because there may be deterministic policies that are infeasible on their own, but can be mixed to satisfy the secondary-cost constraints over expectation.

There are no primal method for optimally solving CSSPs and all optimal CSSP algorithms rely on the dual formulation. The dual formulation optimises over the space of occupation measures  $\mathbf{x}$  where  $x_{s,a}$  represents the expected number of times that  $s$  is reached and then  $a$  applied. LP 1 is the occupation measure LP for CSSPs [31], where  $\llbracket \cdot \rrbracket$  is the Iverson bracket and we use the macros  $\text{out}(s) = \sum_{a \in A(s)} x_{s,a}$  and  $\text{in}(s) = \sum_{s' \in S, a' \in A(s')} x_{s',a'} P(s|s', a')$  for each  $s \in S$ . LP 1 can be interpreted as a network flow, where a unit of flow is injected into  $s_I$ , and must be routed through the actions so that all of it reaches the goals. An optimal solution  $\mathbf{x}$  for LP 1 induces the optimal policy  $\pi_{\mathbf{x}}$  with  $\pi_{\mathbf{x}}(s, a) = x_{s,a} / \text{out}(s) \quad \forall s \in S, a \in A(s)$ .

$$\min_{\mathbf{x}} \sum_{s \in S, a \in A(s)} x_{s,a} C_0(a) \text{ s.t. C1-C4} \quad (\text{LP 1})$$

$$\text{out}(s) - \text{in}(s) = \llbracket s = s_I \rrbracket \quad \forall s \in S \setminus G \quad (\text{C1})$$

$$\sum_{g \in G} \text{in}(g) = 1 \quad (\text{C2})$$

$$x_{s,a} \geq 0 \quad \forall s \in S, a \in A(s) \quad (\text{C3})$$

$$\sum_{s \in S, a \in A(s)} x_{s,a} C_i(a) \leq u_i \quad \forall i \in \{1, \dots, n\} \quad (\text{C4})$$

There is an important connection between deterministic and stochastic policies: any stochastic policy can be represented as a convex combination (a.k.a. mixture) of deterministic ones [14]. We write this decomposition of the stochastic policy  $\pi$  as  $\pi = \mu_0 \pi_0 + \dots +$

$\mu_k \pi_k$  where  $\mu_0, \dots, \mu_k \in \mathbb{R}_{>0}$ ,  $\mu_0 + \dots + \mu_k = 1$ , and  $\pi_0, \dots, \pi_k$  are  $\pi$ 's constituent deterministic policies. We give an example of a CSSP and its policies in appendix A.

### 3 Solving CSSPs with Scalarisation

In this section we explain how to solve CSSPs with scalarisation, and implement this framework in our novel algorithm CARL. First, we give a high-level overview of the algorithm, and then give details for its three main steps.

Underlying the scalarisation approach is LP 2, the primal LP for LP 1. Intuitively, the  $V_s$  variables represent the cost-to-go  $V(s)$ , so we will use  $V(s)$  and  $V_s$  interchangeably. In fact, if we construct LP 2 for an unconstrained SSP, then the  $\lambda_i$  and  $u_i$  terms disappear and the LP becomes the standard primal LP for solving SSPs which encodes the Bellman equations for  $V$ . If we fix  $\lambda$ , then LP 2 again encodes the Bellman equations for an SSP, but with the modified cost function  $C_\lambda(a) = [1 \ \lambda] \cdot C(a)$  (from C5) and the constant one-time terminal cost  $-\sum_{i=1}^n \lambda_i u_i$  (from the objective).  $\mathbb{S}(\lambda)$  denotes such an SSP parameterised by  $\lambda$ .

**Definition 1** (Scalarised SSP). *Given  $\mathbb{C}$  and a scalarisation  $\lambda \in \mathbb{R}_{\geq 0}^n$ ,  $\mathbb{S}(\lambda)$  is an SSP relaxation of  $\mathbb{C}$  with the cost function  $C_\lambda : A \rightarrow \mathbb{R}_{\geq 0}$  s.t.  $C_\lambda(a) = [1 \ \lambda] \cdot C(a)$  and the terminal cost  $-\sum_{i=1}^n \lambda_i u_i$ , which is incurred once when a goal is encountered, and is constant for  $\lambda$ . We write a proper policy's cost as  $C_\lambda(\pi) = [1 \ \lambda] \cdot C(\pi) - \sum_{i=1}^n \lambda_i u_i$ .*

We can interpret each  $\lambda$  as a scalarisation that projects the CSSP's vector cost function onto a scalar cost function. An optimal solution  $V^*, \lambda^*$  for LP 2 gives the solution to the scalarised SSP's Bellman equations that are maximal over all scalarisations. Thanks to the strong duality of LPs, the solution  $V^*, \lambda^*$  can be transformed into an optimal solution for LP 1 using complementary slackness, yielding an optimal stochastic policy for the CSSP. This transformation can be interpreted as extracting all the greedy deterministic policies for  $V^*$ , and combining them into an optimal policy.

$$\begin{aligned} \max_{V, \lambda} \quad & V_{s_I} - \sum_{i=1}^n \lambda_i u_i \quad \text{s.t. C5-C6} \quad (\text{LP 2}) \\ V_s \leq \quad & C_0(a) + \sum_{i=1}^n \lambda_i C_i(a) + \sum_{\substack{s' \in S \\ \forall s \in S \setminus G, a \in A(s)}} P(s'|s, a) V_{s'} \\ V_g = 0 \quad & \forall g \in G \quad (\text{C6}) \\ \lambda_i \geq 0 \quad & \forall i \in \{1, \dots, n\} \quad (\text{C7}) \end{aligned}$$

To avoid solving LP 2 directly, we separate the optimisation over  $V$  and  $\lambda$  into  $\max_\lambda L(\lambda)$ , where  $L(\lambda) = \max_V V_{s_I} - \sum_{i=1}^n \lambda_i u_i$ , i.e.,  $L(\lambda)$  is the optimal policy cost for  $\mathbb{S}(\lambda)$ . If we draw  $L(\lambda)$  for all values of  $\lambda$ , we get a surface that is piecewise linear concave [22]; we give some examples of this in appendix B. Thus, with an oracle that gives  $L(\lambda)$  and a subgradient for each  $\lambda$ , we can use any subgradient method to find  $\lambda^* = \max_\lambda L(\lambda)$ . Finally, we must find  $V^*$  that describes all optimal deterministic policies for  $\mathbb{S}(\lambda^*)$ , from which we extract an optimal stochastic policy for the CSSP.

Alg. 2 presents the high-level pseudocode of this scalarisation algorithm, named CARL. In the remainder of this section, we explain CARL's three steps: (1) finding  $\lambda^*$  (line 3), (2) finding  $V^*$  given  $\lambda^*$  (line 4), and (3) extracting the optimal stochastic policy (line 5). To conclude, we prove CARL's correctness, and give more detail about CARL's error terms.

---

#### Algorithm 2: Solve CSSP with Scalarisation

---

```

1 Function solve(CSSP  $\mathbb{C}$ , heuristic  $H : S \rightarrow \mathbb{R}_{\geq 0}^{n+1}$ )
2   //  $\lambda^* \in \mathbb{R}_{\geq 0}^n, V : S \rightarrow \mathbb{R}_{\geq 0}^{n+1}, V_\nabla^* : S \rightarrow \mathbb{R}_{\geq 0}$ 
3    $\lambda^*, V \leftarrow \text{find-}\lambda^*(\mathbb{C}, H)$ 
4    $V_\nabla^* \leftarrow \text{find-all-opts-for-}\mathbb{S}(\lambda^*, V, \mathbb{C}, H)$ 
5    $\pi^* \leftarrow \text{extract-opt-policy}(\lambda^*, V_\nabla^*, \mathbb{C})$ 
6   return  $\pi^*$ 

7 Function find- $\lambda^*(\text{CSSP } \mathbb{C}, H : S \rightarrow \mathbb{R}_{\geq 0}^{n+1})$ 
8   //  $\lambda \in \mathbb{R}_{\geq 0}^n, V : S \rightarrow \mathbb{R}_{\geq 0}^{n+1}, g \in \mathbb{R}^n$ 
9    $\lambda \leftarrow 0, V \leftarrow H$ 
10  repeat
11     $V \leftarrow \text{solve-}\mathbb{S}(\lambda, V, \mathbb{C})$ 
12    subgradient  $g \leftarrow [V_1(s_I) - u_1 \ \dots \ V_n(s_I) - u_n]$ 
13     $\lambda \leftarrow \text{move from } \lambda \text{ with } g$  (see section 3.1)
14  until  $\lambda$  has converged (see section 3.5)
15  return  $\lambda, V$ 

16 Function solve- $\mathbb{S}(\lambda, V, \mathbb{C})$ 
17    $V \leftarrow \text{solve } \mathbb{S}(\lambda)$  (see defn. 1) with SSP algorithm
18   modified to work on vector  $V$  (see section 3.1)
19   return  $V$ 

20 Function find-all-opts-for- $\mathbb{S}(\lambda, \mathbb{C}, H)$ 
21    $V_\nabla^* \leftarrow \text{solve-}\mathbb{S}(\lambda, V, \mathbb{C})$  with modification that satisfies
22   strong  $\epsilon$ -consistency (see eq. (1) in section 3.2)
23   return  $V_\nabla^*$  s.t.  $V_\nabla^*(s) = [1 \ \lambda^*] \cdot V_\nabla^*(s) \ \forall s \in S$ 

24 Function extract-opt-policy( $\lambda, V_\nabla^*, \mathbb{C}$ )
25    $x^* \leftarrow$  solution to SOL 1 with  $\lambda, V_\nabla^*, \mathbb{C}$ 
26   return  $\pi_{x^*}$ 

```

---

We point out that CARL follows the concepts of Lagrangian decomposition and uses the same framework for finding  $\lambda^*$  as the algorithm by Hong and Williams [22]. The novelty of CARL is that it produces optimal (potentially stochastic) policies and improves on previous techniques, which we discuss further in section 4.

#### 3.1 Finding a Maximal Scalarisation

In this section, we describe *find- $\lambda^*$*  (line 7), which searches over  $\lambda$ s with a subgradient method (outer problem). This technique must compute subgradients, which we do by solving  $\mathbb{S}(\lambda)$  as subproblems

**Outer optimisation over  $\lambda$ .** Suppose that we have an oracle that returns  $L(\lambda)$  and a subgradient  $g$  for each  $\lambda$ . Recall that our optimisation space, i.e., the surface obtained by plotting  $L(\lambda)$  for all  $\lambda$ s, is piecewise linear concave. This function has “sharp kinks” and can not be differentiated there, which is why we need subgradients [30].<sup>1</sup> To find  $\max_\lambda L(\lambda)$ , we can use any optimisation method that accommodates concave and non-smooth search spaces. We follow Hong and Williams [22] and use an exact line search within coordinate search, which has been shown to be efficient for piecewise linear concave models with few constraints.

Coordinate search maximises  $L(\lambda)$  by focusing on one coordinate at a time. That is, it sequentially solves for each  $i \in \{1, \dots, n\}$  the subproblem  $\max_{\lambda_i} L(\lambda)$ , where all  $\lambda_j$  for  $j \neq i$  are fixed from the previous step. Thus, for each coordinate  $i$  it maximises  $L([\lambda_1 \ \dots \ \lambda_n])$  where  $\lambda_i$  is the only variable and all other terms are fixed. This subproblem is solved with an exact line search which exploits that the subproblem is a piecewise linear concave problem with a single variable. The approach is reminiscent of binary search:

<sup>1</sup> Concave functions technically have supergradients, but these are still called subgradients in the literature.

it starts with  $l = 0$  and  $u \in \mathbb{R}_{>0}$  which give lower and upper bounds on the optimal assignment to  $\lambda_i$ . Then, the subgradients are computed at  $\lambda$  with  $\lambda_i = l$  and  $\lambda_i = u$ , the intersection of these subgradients is computed as  $m$ , and either  $l$  or  $u$  is updated to  $m$ , depending on  $m$ 's subgradient. This process repeats until  $l$  and  $u$  converge or their subgradients have the same sign (see [22] for more details). We visualise how  $\lambda$  is updated with coordinate search in appendix B. Unfortunately, coordinate search is incomplete for non-smooth problems [29], as we exemplify in appendix C. Fortunately, these edge cases occur rarely, as we see in our experiments (section 5), and are detectable: if we extract a stochastic policy from  $V^\dagger, \lambda^\dagger$  whose primary cost is greater than  $L(\lambda^\dagger)$ , then we know that  $\lambda^\dagger \neq \lambda^*$ , because  $L(\lambda^*) = C_0(\pi^*)$  where  $\pi^*$  is extracted from  $V^*, \lambda^*$ . If we detect that coordinate search failed we can fall back on a complete method to find  $\lambda^*$ ; we use the subgradient method with projection to ensure  $\lambda \geq 0$  [30]. Line 13 represents a single update of  $\lambda$  using coordinate search or the subgradient method as a fall-back, as we have described.

**Oracle via  $\mathbb{S}(\lambda)$  Subproblem.** Recall that  $\mathbb{S}(\lambda)$  (defn. 1) is an SSP relaxation of the CSSP, and therefore solved optimally by deterministic policies.  $\mathbb{S}(\lambda)$ 's optimal deterministic policy  $\pi_\lambda^*$  gives  $L(\lambda)$  and a subgradient  $g \in \mathbb{R}^n$  in the following way:  $L(\lambda) = C_\lambda(\pi_\lambda^*)$  and  $g = [C_1(\pi_\lambda^*) - u_1 \dots C_n(\pi_\lambda^*) - u_n]$  [22]. These two steps correspond to lines 11 and 12, respectively. Thus, we can implement an oracle for  $\lambda$  by finding  $\pi_\lambda^*$  and then evaluating it w.r.t.  $C_\lambda$  and each secondary-cost function. To find  $\pi_\lambda^*$ , we can use any optimal SSP algorithm with any heuristic that is admissible for SSPs, e.g., h-max [6], LMcut [21], or ROC [32]. In the remainder of this section, we show how to implement *solve-S* (line 16) efficiently.

**Efficient Heuristic Search over  $\mathbb{S}(\lambda)$ s.** The  $\mathbb{S}(\lambda)$  subproblems are unconstrained SSPs and can be solved efficiently with existing heuristic-search algorithms such as CG-iLAO\* [26], but we can improve efficiency further by exploiting some additional properties. First, instead of using a scalar value function  $V : \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$ , we use a vector value function  $\mathbf{V} : \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}^{n+1}$  with entries for each cost function. Then, we replace the scalar Bellman backups with a vectorised variant, where a  $Q$ -value is defined as  $Q(s, a) = C(a) + \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s')$  and the Bellman backup becomes  $\mathbf{V}(s) \leftarrow Q(s, a_{\min})$  for some  $a_{\min}$  that minimises  $[1 \ \lambda] \cdot Q(s, a)$ . With some minor technicalities (see appendix D for more details), SSP algorithms can be modified to use this vector Bellman backup, and they still solve  $\mathbb{S}(\lambda)$  as an unconstrained SSP because the vector value function is projected onto  $V_\lambda(s) = [1 \ \lambda] \cdot \mathbf{V}(s)$ . However, we now work on  $\mathbf{V}$  directly with the advantage that we simultaneously find an optimal policy  $\pi^*$  for  $\mathbb{S}(\lambda)$  and evaluate the different costs of  $\pi^*$ , giving us  $C_i(\pi_\lambda^*)$  for each cost function  $i$ .

To ensure optimality, we require an admissible heuristic to solve  $\mathbb{S}(\lambda)$ , i.e.,  $H(s) \leq V_\lambda^*(s) \ \forall s \in \mathcal{S}$ ; in addition, we also require that the heuristic is a vector in  $\mathbb{R}_{\geq 0}^{n+1}$  so that it can be used with  $\mathbf{V}$ . To address this, we introduce  $\lambda$  heuristics. These compute some admissible scalar heuristic  $H(s)$ , then identify which actions are selected by  $H(s)$  and sum their vector costs to obtain a heuristic cost vector  $\mathbf{H}$  with  $[1 \ \lambda] \cdot \mathbf{H}(s) = H(s) \leq V_\lambda^*(s)$ . For example, if a scalar heuristic uses actions  $a_0$  and  $a_1$  to construct its estimate  $H(s) = [1 \ \lambda] \cdot C(a_0) + [1 \ \lambda] \cdot C(a_1)$ , then we assign  $\mathbf{H}(s) \leftarrow C(a_0) + C(a_1)$  which indeed satisfies  $[1 \ \lambda] \cdot \mathbf{H}(s) = H(s)$ . For most heuristics, we can use their internal data-structures to extract the selected actions, e.g., for LMcut [21] we consider the representative action from each cut and its weight, and for ROC [32] we consider actions' operator counts. Another way to obtain admissible vector

heuristics is with ideal-point (IP) heuristics [15]. Such heuristics are  $\mathbf{H} = \langle H_0, \dots, H_n \rangle$  where  $H_i$  are all admissible w.r.t. their cost function  $i$ . The advantage of  $\lambda$  heuristics is that they are computed for the specific  $\mathbb{S}(\lambda)$  and are therefore more informative, but with the trade-off that they have to be recomputed for each  $\mathbb{S}(\lambda)$ , whereas an IP heuristic is admissible for all  $\mathbb{S}(\lambda)$ s.

**Warm starting  $\mathbb{S}(\lambda)$ s.** Notice that all  $\mathbb{S}(\lambda)$ s are the same problem only differing in their cost functions. This makes  $\mathbf{V}$  from a previous  $\mathbb{S}(\lambda_{\text{old}})$  a good candidate to start solving a new  $\mathbb{S}(\lambda_{\text{new}})$ , a technique known as warm start in Operations Research. However, we have to be careful in reusing  $\mathbf{V}$ . Suppose  $\mathbf{V}$  is a solution to  $\mathbb{S}(\lambda_{\text{old}})$  and we have  $V_{\lambda_{\text{old}}}$  with  $V_{\lambda_{\text{old}}}(s) = [1 \ \lambda_{\text{old}}] \cdot \mathbf{V}(s) \ \forall s \in \mathcal{S}$  and  $V_{\lambda_{\text{new}}}$  with  $V_{\lambda_{\text{new}}}(s) = [1 \ \lambda_{\text{new}}] \cdot \mathbf{V}(s) \ \forall s \in \mathcal{S}$ . Even if  $V_{\lambda_{\text{old}}}$  is admissible for  $\mathbb{S}(\lambda_{\text{old}})$ , there is no guarantee that  $V_{\lambda_{\text{new}}}$  is admissible for  $\mathbb{S}(\lambda_{\text{new}})$ , which may break the optimality of heuristic search algorithms. Conveniently, CG-iLAO\* has a built-in mechanism for efficiently handling potentially inadmissible states by tracking states whose values have increased or decreased, i.e., we mark the states where  $[1 \ \lambda_{\text{new}}] \cdot \mathbf{V}(s)$  decreases or increases w.r.t.  $[1 \ \lambda_{\text{old}}] \cdot \mathbf{V}(s)$  in  $\Gamma$ . As explained in section 2, this ensures that CG-iLAO\* fixes any issues in the value function and finds the optimal policy for  $\mathbb{S}(\lambda_{\text{new}})$ . Reusing  $\mathbf{V}$  in this way speeds up CARL significantly because the optimal policy for one  $\mathbb{S}(\lambda)$  is often good for the next one, and in the best case, if the change in  $C_\lambda$  does not affect the optimal policy, then  $\mathbf{V}$  immediately solves the new  $\mathbb{S}(\lambda')$ .

### 3.2 Finding an Optimal Value Function

Once we have  $\lambda^*$ , we need to compute  $V^* : \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$  that describes all optimal deterministic policies for  $\mathbb{S}(\lambda^*)$ , which is done in the pseudocode by *find-all-opts-for-S* (line 19). Let  $\Pi(V)$  denote the set of tied-greedy policies for  $V$ , that is, the greedy policies of  $V$  obtained by breaking ties in all possible ways. Then, if  $V^*$  is the unique solution to the Bellman equations,  $\Pi(V^*)$  gives precisely all optimal deterministic policies for the SSP. For us,  $V^*$  is produced by an SSP algorithm that only guarantees  $\epsilon$ -consistency [7]. This means  $V^*$  need only have a small error within one policy's envelope, and any other policies' envelopes can have arbitrarily larger value functions, i.e.,  $V^*$  need only encode a single greedy policy. We give an example to illustrate this in appendix E. In order to obtain a value function that encodes all (approximately) optimal policies, without solving the Bellman equations exactly, we introduce strong  $\epsilon$ -consistency. This condition requires a value function to be  $\epsilon$ -consistent for all of its greedy policies. Formally,  $V$  is strongly  $\epsilon$ -consistent when

$$\forall s \in \bigcup_{\pi \in \Pi(V)} \mathcal{S}^\pi \quad \text{RES}(s) \leq \epsilon. \quad (1)$$

When a value function is strongly  $\epsilon$ -consistent, we write it as  $V_\epsilon^*$ . Note that this has not been defined for SSPs before, because there was no need for more than one policy.

**Finding  $V_\epsilon^*$ .** Existing SSP algorithms are designed to find  $\epsilon$ -consistent solutions, so we must modify them to find strongly  $\epsilon$ -consistent solutions. This is done by updating the termination condition, and ensuring that Bellman backups are applied for all greedy policies for  $V$ . We make this modification concrete for the version of CG-iLAO\* from section 3.1, but it can be done for other algorithms too. Recall that this version of CG-iLAO\* works over the vector value function  $\mathbf{V} : \mathcal{S} \rightarrow \mathbb{R}^{n+1}$ , so we are looking for  $V_\epsilon^*$  that induces  $V_\epsilon^*$  with  $V_\epsilon^*(s) = [1 \ \lambda^*] \cdot V_\epsilon^*(s)$ . Instead of tracking the candidate policy  $\hat{\pi}_V$ , we track the union of

all tied-greedy policies (up to  $\epsilon$ ) with  $\text{tied}(s) = \{a \in A(s) : Q_\lambda(s, a) \leq \min_{a' \in A(s)} Q_\lambda(s, a') + \epsilon\}$ . Then, CG-iLAO\*'s policy envelope  $\mathcal{E}$  over all tied-greedy policies is found by running the depth-first search from  $s_I$  over  $\text{tied}$ . Additionally, any actions with  $Q_\lambda(s, a) \leq V_\lambda(s) + \epsilon$  must be added to  $\mathbb{S}$  in line 9 so that new tied-greedy actions are caught. Also, we replace the termination condition  $\hat{\pi}_{\text{old}} = \hat{\pi}_V$ , now requiring that no new actions have been added to any  $\text{tied}(s)$ . We only care about new actions and do not track removals, because they only shrink the envelope and do not affect  $V$ . These changes ensure that CG-iLAO\* outputs a strongly  $\epsilon$ -consistent solution which captures all  $\epsilon$ -consistent policies. To prove this, we use similar arguments to Schmalz and Trevizan [26] to obtain that  $\forall s \in \mathcal{E} \text{ RES}(s) \leq \epsilon$ . But  $\mathcal{E}$ , obtained by DFS over the tied-greedy actions (up to  $\epsilon$ ), captures all tied-greedy policies on  $V$  by construction, so  $V$  must be strongly  $\epsilon$ -consistent.

### 3.3 Extracting Stochastic Policies

Suppose we have  $\lambda^*$  and  $V_\psi^*$  which encodes all optimal deterministic policies for  $\mathbb{S}(\lambda^*)$ . The final step of CARL is to extract an optimal (potentially stochastic) policy, which is done by *extract-opt-policy* (line 22). To implement this, recall that  $V_\psi^*$  and  $\lambda^*$  give the primary cost of all (potentially stochastic) optimal policies for the CSSP with  $C_0(\pi^*) = V_\psi^*(s_I) - \sum_{i=1}^n \lambda_i^* u_i$ . Furthermore, we can extract  $\pi^*$ 's envelope and support from the envelopes and supports of  $V_\psi^*$ 's greedy policies. We now show how to use these insights to extract  $\pi^*$ , and thereby solve the CSSP optimally.

Given a policy  $\pi$ , we can evaluate its primary cost  $C_0(\pi)$  with LP 1 by restricting the LP's variables to the support of  $\pi$ , adding the constraints  $x_{s,a}/\text{out}(s) = \pi(s, a)$  for all  $s, a \in \text{supp}(\pi)$ , and removing the objective function. The resulting System of Linear Equations and Inequalities (SOL) finds  $C_0(\pi)$  as  $\sum_{s \in \mathbb{S}, a \in A(s)} x_{s,a} C_0(a)$ . Our case is the converse since we have  $\pi^*$ 's cost  $C_0(\pi^*)$  but do not have the distributions  $\pi(s, \cdot)$ . So, instead of adding the constraints  $x_{s,a}/\text{out}(s) = \pi(s, a)$ , we add the constraint  $\sum_{s \in \mathbb{S}, a \in A(s)} x_{s,a} C_0(a) = C_0(\pi^*) = V_\psi^*(s_I) - \sum_{i=1}^n \lambda_i^* u_i$ ; and we approximate  $\text{supp}(\pi^*)$  with  $\text{supp}(V_\psi^*)$ , which is defined as  $\bigcup_{\pi \in \Pi(V_\psi^*)} \text{supp}(\pi)$ , noting that  $\text{supp}(V_\psi^*) \supseteq \text{supp}(\pi^*)$  as we explain later. This yields SOL 1 without C9 and C10 which we also explain later. Thus, a solution  $\mathbf{x}$  for SOL 1 (with or without C9 and C10) induces a feasible policy  $\pi_{\mathbf{x}}$  with the same cost as an optimal policy by construction. SOL 1's only decision is how to distribute the actions in  $\text{supp}(V_\psi^*)$  so that the resulting policy has the required costs; in other words, it computes how to combine  $\mathbb{S}(\lambda^*)$ 's deterministic policies into an optimal feasible policy for  $\mathbb{C}$ .

find  $\mathbf{x}$  s.t. C1–C3 over  $\text{supp}(V_\psi^*)$  and C8–C10 (SOL 1)

$$\sum_{s, a \in \text{supp}(V_\psi^*)} x_{s,a} C_0(a) = V_\psi^*(s_I) - \sum_{i=1}^n \lambda_i^* u_i \quad (\text{C8})$$

$$\sum_{s, a \in \text{supp}(V_\psi^*)} x_{s,a} C_i(a) \leq u_i \quad \forall i \in \{1, \dots, n\} \text{ s.t. } \lambda_i^* = 0 \quad (\text{C9})$$

$$\sum_{s, a \in \text{supp}(V_\psi^*)} x_{s,a} C_i(a) = u_i \quad \forall i \in \{1, \dots, n\} \text{ s.t. } \lambda_i^* > 0 \quad (\text{C10})$$

Formally, SOL 1 comes from complementary slackness, a technique that lets us transform an optimal solution of an LP into an optimal solution for its dual [5]. Suppose  $V_\psi^*, \lambda^*$  is an optimal basic feasible solution for LP 2 that is non-degenerate; intuitively, a basic feasible solution corresponds to a deterministic policy (rather than a stochastic one) and non-degeneracy is a technical requirement that the solution does not have too many variables set to zero. With such

$V_\psi^*, \lambda^*$ , we can recover a solution for LP 2's dual (LP 1) with complementary slackness by constructing a system of linear equations according to the following rules, and then solving it:

$$V_\psi^*(s) > 0 \implies \text{C1 for } s \text{ is tight, i.e., } \text{out}(s) - \text{in}(s) = \llbracket s = s_I \rrbracket$$

$$\lambda_i > 0 \implies \text{C4 for } i \text{ is tight, i.e., } \sum_{s \in \mathbb{S}, a \in A(s)} x_{s,a} C_i(a) = u_i$$

$$\text{C5 for } s, a \text{ is loose, i.e., } V_\psi^*(s) < Q_\psi^*(s, a) \implies x_{s,a} = 0.$$

These rules describe the constraints of SOL 1. In particular, note that these rules tighten the secondary-cost constraints to an equality  $\sum_{s \in \mathbb{S}, a \in A(s)} x_{s,a} C_i(a) = u_i$  whenever  $\lambda_i^* > 0$ , yielding the constraints C9 and C10. Unfortunately, our  $V_\psi^*, \lambda^*$  is often degenerate, so this system of linear equations is underspecified, and we must include additional constraints from the original problem. In particular, we have to reintroduce the secondary cost constraints C4 where  $\lambda_i = 0$ , which is what we do in C9. We also include C8 to ensure the resulting policy is indeed optimal.

We have been claiming that  $\text{supp}(\pi^*) \subseteq \text{supp}(V_\psi^*)$  and that SOL 1 is constructing  $\pi^*$  from the optimal deterministic policies of  $\mathbb{S}(\lambda^*)$ . This follows from the last complementary slackness rule: occupation measures  $x_{s,a}$  are only allowed to be non-zero if  $V_\psi^*(s) = Q_\psi^*(s, a)$ , i.e., if  $a$  is a tied-greedy action and thereby  $\pi(s) = a$  for some  $\pi \in \Pi(V_\psi^*)$ .

### 3.4 Correctness of CARL

In this section we show that CARL always finds an optimal solution. First, *find- $\lambda^*$* ( $\mathbb{C}, \mathbf{H}$ ) (alg. 2) is guaranteed to find  $\lambda^*$ . This is because  $\max_\lambda L(\lambda)$  is a piecewise linear concave problem, solving  $\mathbb{S}(\lambda)$  yields subgradients [22], and we use a complete subgradient method that is guaranteed to converge to  $\lambda^*$  (if coordinate search fails, it falls back on the subgradient method, see section 3.1). It remains to show that CARL finds  $V^*$  for  $\mathbb{S}(\lambda^*)$  such that the solution encodes an optimal policy.

**Lemma 1.** *If  $\pi^*$  is an optimal policy for the CSSP, then its constituent deterministic policies must be optimal for  $\mathbb{S}(\lambda^*)$ .*

*Proof.* We first show that  $\pi^*$  is optimal for  $\mathbb{S}(\lambda^*)$ , and then conclude that its constituent deterministic policies must also be optimal for  $\mathbb{S}(\lambda^*)$ . Policy  $\pi^*$ 's cost in  $\mathbb{S}(\lambda^*)$  can be rewritten as  $C_{\lambda^*}(\pi^*) = C_0(\pi^*) + \sum_{i=1}^n \lambda_i^* (C_i(\pi^*) - u_i)$ . But  $\pi^*$  is feasible, so  $C_i(\pi^*) - u_i \leq 0$  for all  $i \in \{1, \dots, n\}$ , and therefore  $C_{\lambda^*}(\pi^*) \leq C_0(\pi^*)$ . On the other hand, we know that  $C_{\lambda^*}(\pi) \geq C_0(\pi^*)$  for any policy  $\pi$ , due to the semantic of  $\lambda^*$  within strong duality. Thus,  $C_{\lambda^*}(\pi^*) = C_0(\pi^*)$ , which is the minimal cost for  $\mathbb{S}(\lambda^*)$ , and  $\pi^*$  is optimal for  $\mathbb{S}(\lambda^*)$ . We know that  $\pi^*$  can be decomposed into the convex combination of deterministic policies  $\mu_0 \pi_0 + \dots + \mu_k \pi_k$ , and  $C_{\lambda^*}(\pi_i) = C_{\lambda^*}(\pi^*)$  for each  $i \in \{0, \dots, k\}$ , because otherwise  $\pi^*$  can not be optimal for  $\mathbb{S}(\lambda^*)$ .  $\square$

**Theorem 1.** *CARL's find-all-opts-for- $\mathbb{S}(\lambda^*, \mathbb{C}, \mathbf{H})$  finds  $V_\psi^*$  so that SOL 1 with  $\lambda^*$  and  $V_\psi^*$  constructed by  $V_\psi^*(s) = \lambda^* \cdot V_\psi^*(s)$  produces an optimal stochastic policy.*

*Proof.* Suppose  $V_\psi^*$  encodes all the optimal deterministic policies for  $\mathbb{S}(\lambda^*)$ . If an optimal policy  $\pi^*$  exists for the CSSP, then its constituent deterministic policies must be optimal for  $\mathbb{S}(\lambda^*)$  by lem. 1; therefore its constituent policies are a subset of  $\Pi(V_\psi^*)$ . It follows that SOL 1 contains the whole support of such policy  $\pi^*$ . Since  $C_0(\pi^*) = V_\psi^*(s_I) - \sum_{i=1}^n \lambda_i^* u_i$  due to strong duality, it must be the case that SOL 1 induces an optimal policy, thanks to its connection

to LP 1. We have argued that the version of CG-iLAO\* presented in section 3.2 finds all policies for which  $V$  can be  $\epsilon$ -consistent, which gives approximately optimal policies for small  $\epsilon$ , and precisely the optimal policies as  $\epsilon \rightarrow 0$ . We give more detail about CARL’s error terms in section 3.5.  $\square$

### 3.5 Error Terms

In this section, we describe the error terms that are present in CARL and explain how to bound their effects on CARL’s policies. We have already discussed that CG-iLAO\*, which solves CARL’s  $\mathbb{S}(\lambda)$  subproblems, uses strong  $\epsilon$ -consistency with  $\epsilon \in \mathbb{R}_{>0}$  as its stopping condition, extending regular  $\epsilon$ -consistency for SSPs (see section 3.2). The error associated with  $\epsilon$ -consistency is well-understood and accepted in the SSP literature: for a fixed  $\epsilon$  it is possible to construct a pathological problem where an  $\epsilon$ -consistent policy is arbitrarily worse than the optimal policy, but for any SSP, as  $\epsilon \rightarrow 0$ ,  $\epsilon$ -consistent policies become optimal [7, 25, 17]. We inherit the same properties in strong  $\epsilon$ -consistency, so our version of CG-iLAO\* produces  $V'_\psi \approx V^*_\psi$  where the error disappears as  $\epsilon \rightarrow 0$ . Another error term appears in *find- $\lambda^*$* , which we have omitted until now for simplicity. Recall that to find  $\lambda^*$  we use coordinate search, and if it fails we fall back on the subgradient method (see section 3.1). Both require a tolerance term  $\eta \in \mathbb{R}_{>0}$ . In practice, coordinate search stops when it iterates over all coordinates and no single coordinate was able to improve  $L(\lambda)$  by more than  $\eta$ . The subgradient method stops when its step size (which is monotonically decreasing over the algorithm’s iterations) is smaller than  $\eta$ . In both cases, this yields an approximate solution  $\lambda' \approx \lambda^*$ , where the error disappears as  $\eta \rightarrow 0$ .

Both  $\epsilon$  and  $\eta$  can affect CARL’s solution, and we explain how this can be mitigated. Assume for now that  $\epsilon$  is sufficiently small. If CARL finds  $\lambda' \approx \lambda^*$  and an associated policy  $\pi$ , we have

$$L(\lambda') \leq C_0(\pi^*) \leq C_0(\pi)$$

where  $\pi^*$  is an optimal policy for the CSSP. This gives us the optimality gap  $C_0(\pi) - L(\lambda')$ , which upper bounds how far  $\pi$ ’s primary cost deviates from  $\pi^*$ ’s. If the gap is too large, we can reduce  $\eta$ . Now, we return our attention to  $\epsilon$ . The value of  $\epsilon$  affects the subgradient search oracle’s value of  $L(\lambda)$  and subgradient  $g$  at each  $\lambda$ , which in turn can significantly affect the found policy’s quality and breaks the optimality gap for  $\eta$ . Recall that this is standard for algorithms using  $\epsilon$ -consistency, and they only guarantee an optimal policy as  $\epsilon \rightarrow 0$ . For additional guarantees, there are “outer loop” methods that run the algorithm and select  $\epsilon$  automatically in order to give the stronger guarantee of  $\epsilon$ -optimality [17, 19]. CARL is in a similar position, and guarantees an optimal policy as  $\epsilon \rightarrow 0$ , but requires an outer loop method to select  $\epsilon$  automatically for stronger guarantees. The outer loop methods for ensuring  $\epsilon$ -optimality apply to the strongly  $\epsilon$ -consistent subproblems with minimal changes, which lets us recover a correct optimality gap. Thus, it is possible to modify CARL to produce solutions with stronger guarantees. In practice, values of  $\epsilon = \eta = 0.0001$  are usually “good enough.”

## 4 Related Work

In the Operations Research (OR) community, it is common to consider the Lagrangian dual of a constrained problem, and then use some version of the subgradient method to find a solution. This framework has been used to find plans for constrained deterministic shortest path problems [16], stochastic policies for Constrained Partially Observable Markov Decision Processes (POMDPs) [23], and

deterministic policies for CSSPs [22]. Our presentation differs from theirs, but LP 2 is functionally identical to the Lagrangian dual of LP 1, and CARL fits within this OR framework. They have the same  $\lambda$  as our scalarisation term, but they call it a Lagrangian multiplier.

Hong and Williams [22]’s algorithm is the most similar to CARL. Both solve CSSPs and start by finding the CSSPs’  $\lambda^*$ . The difference is that CARL looks for optimal, potentially stochastic, policies whereas they focus on deterministic ones. Our method for finding  $\lambda^*$  is based on theirs but we make substantial algorithmic improvements. They solve  $\mathbb{S}(\lambda)$ s with iLAO\* and reuse  $V$  as a warm start (see section 3.1); however, they need to run VI on iLAO\*’s expanded states to make  $V$  admissible for the next step. CARL uses CG-iLAO\*, allowing it to focus only on the changes to  $\lambda$  as opposed to VI’s blind search. As a result, CARL finds  $\lambda^*$  up to  $4.7\times$  faster than their method with a mean speedup of  $2\times$ .

Once  $\lambda^*$  has been found, the similarities between our algorithms end. Strong duality does not hold for deterministic policies, thus  $\lambda^*$  does not directly produce an optimal deterministic policy for the original CSSP. To overcome this, they use an expensive second stage that enumerates deterministic policies until the optimal one is found. CARL does not require this stage because strong duality holds for stochastic policies. Instead, we must ensure that  $V^*$  encodes all optimal policies for  $\mathbb{S}(\lambda^*)$ , which can be done cheaply with dynamic programming, and then we use SOL 1 to extract the optimal policy for the CSSP. Both these steps have no analogue in their algorithm. Thus, CARL is guaranteed to find an optimal policy, stochastic or not, for CSSPs with relatively little computational effort after finding  $\lambda^*$ , whereas Hong and Williams [22]’s method can only find an optimal *deterministic* policy, which can cost more than the optimal stochastic policy, and requires an expensive second stage to do so.

The algorithm by Lee et al. [23] shares with CARL that they get deterministic policies from their subproblems, which must be combined into a stochastic one. However, their setting and approach is quite different: they use Monte-Carlo sampling to approximate their  $Q$ -values, and their policies map histories onto action distributions, obtained by solving an LP for each history. We compute  $V^*$  and its  $Q$ -values with a modified SSP search, and then solve SOL 1 once.

Multi-Objective (MO) planning is also relevant. There are various models of probabilistic MO problems, e.g., MOMDPs with queries from model checking [13] and MOSSPs from planning [11]. These models are similar to CSSPs in that they have multiple cost functions, but the difference is that MO problems are solved by finding all undominated trade-offs between cost functions, e.g., consider  $\pi_1, \pi_2, \pi_3$  with  $C(\pi_1) = [1 \ 3]$ ,  $C(\pi_2) = [3 \ 1]$ ,  $C(\pi_3) = [3 \ 2]$ ;  $\pi_1$  and  $\pi_2$  are undominated and give different trade-offs of  $C_0$  and  $C_1$ , and  $\pi_3$  is dominated by  $\pi_2$ . The set of undominated policies is called the Pareto front. In contrast, CSSPs are solved by a single policy that minimises  $C_0$  while satisfying the secondary-cost constraints. Consequently, MOSSPs are more general than CSSPs, and their algorithms can be applied to solve CSSPs, e.g., [13, 11], but these are much less efficient because they need to track an entire set of policies in the Pareto front, rather than searching for a single policy.

Scalarisation is a standard technique in MO planning [12], and has been applied to MOMDPs with queries [13]. Their algorithm’s similarity to CARL is that both project the cost vector onto scalars to induce single-objective SSPs (or MDPs) as their subproblems. However, the way that the scalarisations are explored and the purpose is completely different: they use different scalarisations to find different trade-offs, and CARL uses scalarisation in the context of Lagrangian optimisation where scalarisations  $\lambda$  are also called Lagrangian multipliers. Importantly, CARL’s scalarisation is a single-

objective technique, and it progressively finds a better scalarisation vector, whereas Forejt et al. [13] accumulate a set of non-dominated scalarisation vectors. Another difference is that Forejt et al. [13] use Value Iteration to solve the subproblems, and we use heuristic search, which increases CARL’s efficiency but presents technical challenges as addressed in section 3.1 and section 3.2. Heuristic search has been shown to dominate non-guided search approaches from model checking, e.g., [2]. We point out that our method for efficiently solving subproblems can be applied to the algorithm from Forejt et al. [13], establishing another contribution of this work.

## 5 Experiments

We compare our novel algorithm CARL with i-dual [31] and  $i^2$ -dual [32], the state-of-the-art algorithms for finding optimal stochastic policies for CSSPs. For CARL, we consider the heuristics LMcut [21] and ROC [32] as  $\lambda$  heuristics (section 3.1), called  $\lambda$ -LMcut and  $\lambda$ -ROC. In appendix F we also test LMcut and ROC as ideal-point (IP) heuristics (section 3.1). For i-dual, we use ideal-point LMcut (IP-LMcut) and C-ROC [32].  $i^2$ -dual uses C-POM [32], which is built into the algorithm. In our test problems, we allow CSSPs to have dead ends, but transform them into CSSPs without dead ends with the finite-penalty transformation, which adds give-up actions to each state that lead to a goal with a large cost penalty. We require a user-specified dead-end penalty for each cost function [31]. For problems with dead ends before the finite-penalty transformation, we augment ROC and  $i^2$ -dual’s built-in C-POM with h-max as a dead-end detector [32]. We run experiments over the following benchmark domains:

**Search and Rescue (SAR) [31].** A drone must rescue one of multiple survivors on an  $n \times n$  grid as quickly as possible ( $n \in \mathbb{N}$ ). One survivor’s position is known and has distance  $d \in \mathbb{N}$  from the drone’s initial position. Some other locations, selected according to a density  $r \in [0, 1]$ , have a predetermined probability of having another survivor. There is a single secondary-cost constraint over fuel usage.

**Elevators (Elev) [31].** In a 10-floor building,  $w \in \mathbb{N}$  people have called the elevator and are waiting, and  $h \in \mathbb{N}$  hidden people have not pressed the button yet, and do so probabilistically in each timestep. The  $e \in \mathbb{N}$  elevators must be routed to deliver all people to their destinations in the fewest timesteps. Each person has a constraint over the time spent waiting, and the time spent travelling in an elevator, resulting in  $2(w + h)$  secondary-cost constraints.

**Exploding Blocks World (ExBW) [31].** Based on Exploding Blocks World for SSPs [9], we must rearrange  $N \in \mathbb{N}$  blocks into a given arrangement with the minimal number of moves. This is complicated by the blocks being rigged with an explosive that can explode, destroying the block or table immediately underneath it. In the CSSP variant, the table can be repaired with a large primary cost, but destroyed blocks can not be repaired. There is a single secondary-cost constraint, limiting the expected number of destroyed blocks to  $c \in \mathbb{R}_{\geq 0}$ . We consider specific starting and goal arrangements of blocks from the 2008 International Probabilistic Planning Competition [10], and identify specific problems by the parameter  $id$ .

**PARC Printer (PARC) [32].** Based on the IPC domain, a modular printer must print four pages with various requirements. Three of the printer’s components are unreliable, where a page can get jammed with probability 0.1, simultaneously ruining the page and rendering the component unusable. There are two secondary-cost constraints:  $f \in [0, 1]$  constrains how many components are allowed to fail and

$u \in \mathbb{N}$  limits how often a particular component is allowed to be used ( $u = \infty$  deactivates this constraint).

**Triangle Tireworld (CTW) [27].** In the Triangle Tireworld for SSPs we must drive across a triangular network of cities, where  $n \in \mathbb{N}$  specifies the size of the network, and  $d \in \mathbb{N}$  gives the starting distance from the goal [24, 26]. Between any two cities the car gets a flat tyre with probability 0.5; with a spare tyre the flat can be replaced and the journey can continue; if no spare is loaded, the car is stuck and no actions are available. The car only fits one spare tyre at a time, and spare tyres can only be acquired in select cities. In the CSSP variant, tyres can be purchased in these cities for one of  $c \in \mathbb{N}$  currencies, and the problem has secondary-cost constraints over the  $c$ -many currencies.

We ran each planner and heuristic on 30 instances of each problem on a cluster with Intel Xeon 3.2 GHz CPUs limited to one CPU core with 30 mins and 4GB RAM. For SAR and Elev, a problem instance is a randomly generated problem for the specified parameters. For the other domains, the parameters specify a unique problem and an instance couples this problem with a random seed that is passed to the algorithm. We use CPLEX Version 22.1.1 to solve LPs and SOLs. The source code and benchmarks are available at [28].

The results of our experiments are presented in tab. 1. We give each algorithm’s coverage (the number of instances solved within its time and memory limits), and over the instances where the algorithm converged we give its mean runtime (and 95% C.I.). We highlight the fastest algorithms, where an algorithm is considered the fastest for a particular problem if it has the maximum coverage and then its mean time plus 95% C.I. is less than the next best algorithm’s mean time minus 95% C.I. We now identify the best performing algorithm per domain.

**SAR, ExBW, and CTW.** CARL dominates in these domains: either CARL( $\lambda$ -ROC) or CARL( $\lambda$ -LMcut) is the best performing planner and, with exception of ExBW (1,5,.1) and (3,6,.91), the second best performing planning is also CARL. Moreover, CARL can obtain substantially larger coverage with the extreme case of ExBW (6,8,.3) in which CARL using both heuristics obtained full coverage while all other planners had zero coverage. Comparing only entries with full coverage in these domains, CARL has  $10\times$ ,  $19\times$ , and  $42\times$  speedup on average w.r.t. i-dual (C-ROC), i-dual (IP-LMcut) and  $i^2$ -dual, respectively. Between CARL( $\lambda$ -ROC) and CARL( $\lambda$ -LMcut), there is no clear winner on SAR; CARL( $\lambda$ -LMcut) dominates on ExBW; and CARL( $\lambda$ -ROC) dominates on CTW.

**Elev.** With one elevator i-dual dominates, and with two elevators CARL dominates. With one elevator, CARL fails to solve some instances for one of two reasons: (1) CARL’s coordinate search fails and falls back on the subgradient method, which is significantly slower than coordinate search; (2)  $\mathbb{S}(\lambda)$ s are too difficult to solve in some cases, which suggests the heuristic does not give adequate guidance. Regarding (1), we note that these Elev instances are the only cases in our benchmarks where coordinate search fails and CARL is forced to fall back on the subgradient method to ensure completeness (see section 3.1). Here, the subgradient method is significantly slower than coordinate search and failed to converge in time. However, we reiterate that falling back to the subgradient method ensures that CARL is complete, and given enough time and memory CARL would solve the instances. With two elevators, the problem structure changes so that coordinate search succeeds and (1) does not occur, and any instances where CARL fails to find a solution is because of (2), that is, the coordinate search’s subproblems take too long to





## A CSSP Example

Here, we present an example of a CSSP.

**Getting to Work.** The agent needs to make its way from home to work. It can run, use a taxi, or walk to the train station where it can try to take the train. The train is cancelled with 50% probability. Each action has a cost vector  $[t \ p \ e]$  in terms of time ( $t$ ), price ( $p$ ), and personal effort ( $e$ ). The agent’s task is to get to work in minimal time s.t. price  $\leq 15$  and effort  $\leq 10$  over expectation. The CSSP is shown in fig. 1. This problem has three proper deterministic policies:

- $\pi_{\text{run}} = \{s_I \mapsto \text{run}\}$ ,
- $\pi_{\text{taxi}} = \{s_I \mapsto \text{taxi}\}$ ,
- $\pi_{\text{train}} = \{s_I \mapsto \text{walk}, s_1 \mapsto \text{train}, s_2 \mapsto \text{train}\}$ .

This problem has a unique optimal policy  $\pi^*$  which is stochastic with  $\pi^*(s_I, \text{run}) = 0.5$  and  $\pi^*(s_I, \text{taxi}) = 0.5$  where

- $C_0(\pi^*) = 1$  (time),
- $C_1(\pi^*) = 15 \leq 15$  (price),
- $C_2(\pi^*) = 10 \leq 10$  (effort).

Note that  $\pi^*$  can be expressed as  $\pi^* = 0.5\pi_{\text{run}} + 0.5\pi_{\text{taxi}}$ .

## B Visualisation of $L(\lambda)$

In this section, we visualise the surface of  $L(\lambda)$  for two different CSSPs.

**Getting to Work.** First, we consider the “Getting to Work” example from appendix A. Fig. 2 shows this problem’s  $L(\lambda)$  as  $\lambda$  varies. Each linear surface of the plot corresponds to one of the deterministic policies; we have labelled the linear surfaces associated with  $\pi_{\text{run}}$  and  $\pi_{\text{train}}$  with “run” and “train” respectively, and  $\pi_{\text{taxi}}$  is obscured. A linear surface indicates that the corresponding policy is optimal for the  $\mathbb{S}(\lambda)$  at the relevant values of  $\lambda$ ; edges and vertices indicate that multiple deterministic policies are optimal at the  $\mathbb{S}(\lambda)$ . In particular, note that at  $\lambda^* = 0$  we have an edge between  $\pi_{\text{run}}$  and  $\pi_{\text{taxi}}$ , and  $\pi_{\text{train}}$  does not intersect there; this indicates that the optimal (in this case stochastic) policy consists of  $\pi_{\text{run}}$  and  $\pi_{\text{taxi}}$ , which we know to be true. This example is not interesting in terms of coordinate search, because  $\lambda^* = 0$  is found in the first step.

**Interesting CSSP for Coordinate Search.** Consider the CSSP in fig. 3, with  $u_1 = u_2 = 15$ . Its optimal policy  $\pi^*$  is stochastic with

- $\pi^*(s_I, a_2) = 1$
- $\pi^*(s_1, a_4) = \frac{1}{4}$     $\pi^*(s_1, a_5) = \frac{3}{4}$

and it has  $C_0(\pi^*) = 4$  and  $C_1(\pi^*) = C_2(\pi^*) = 15$ . Fig. 4 shows this CSSP’s  $L(\lambda)$  as  $\lambda$  varies. As before, each linear surface corresponds to a deterministic policy that is optimal for the relevant  $\mathbb{S}(\lambda)$ . We visualise the values of  $\lambda$  considered by coordinate search with the red lines across the surface, i.e., the line starts at the initial value of  $\lambda = 0$ , then moves to  $\lambda = [0.025 \ 0]$ , then  $\lambda = [0.025 \ 0.025]$ , and so on, eventually converging to  $\lambda^* = [0.2 \ 0.2]$ . Recall that coordinate search only changes a single coordinate per step, which gives the red lines’ “stair-case” shape.

## C Pathological Example for Coordinate Search

There are non-smooth problems where coordinate search fails. For example, consider the CSSP in fig. 5 with two secondary-cost constraints  $u_1 = u_2 = 1$ , i.e., a policy  $\pi$  is feasible if  $C_1(\pi) \leq 1$  and  $C_2(\pi) \leq 1$ . The CSSP has three actions that lead deterministically from  $s_I$  to the goal with costs

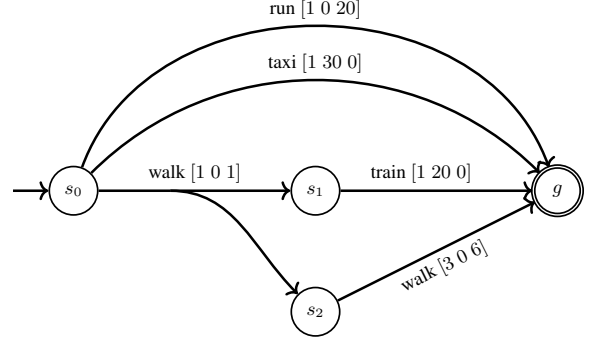


Figure 1. The CSSP associated with the “Getting to Work” example.

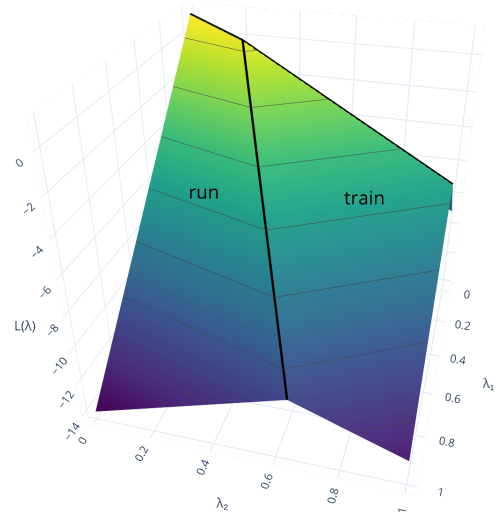


Figure 2.  $L(\lambda)$  for the “Getting To Work” example.

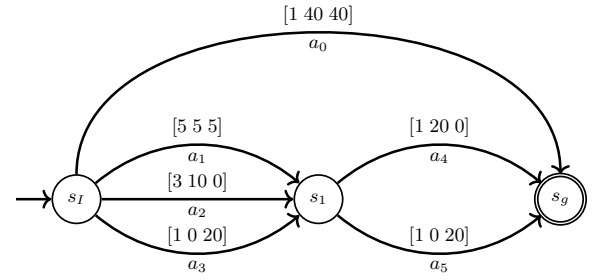
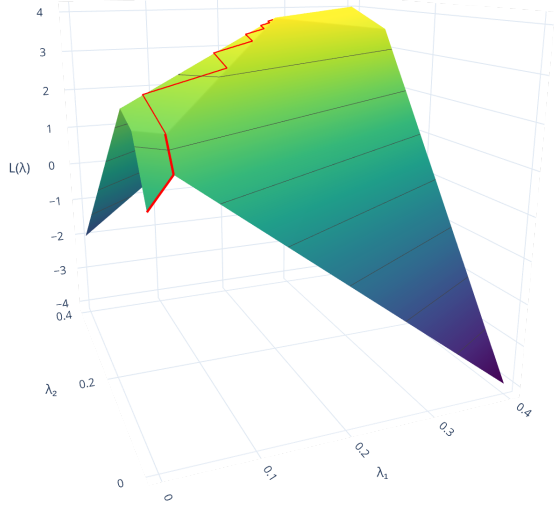


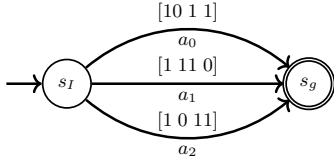
Figure 3. A CSSP that is interesting for coordinate search. It has  $u_1 = u_2 = 15$ .

- $C(a_0) = [10 \ 1 \ 1]$ ,
- $C(a_1) = [1 \ 11 \ 0]$ , and
- $C(a_2) = [1 \ 0 \ 11]$ ,

inducing the three deterministic policies  $\pi_0, \pi_1, \pi_2$  with  $\pi_i(s_I) = a_i$  for  $i \in \{0, 1, 2\}$ . Taking into account the terminal costs, the policy costs are



**Figure 4.**  $L(\lambda)$  for the “Interesting CSSP for Coordinate Search” example.



**Figure 5.** A pathological CSSP where coordinate search fails to find  $\lambda^*$ . It has  $u_1 = u_2 = 1$ .

- $C_\lambda(\pi_0) = 10 + \lambda_1(1 - u_1) + \lambda_2(1 - u_2) = 10$ ,
- $C_\lambda(\pi_1) = 1 + \lambda_1(11 - u_1) + \lambda_2(0 - u_2) = 1 + 10\lambda_1 - \lambda_2$ , and
- $C_\lambda(\pi_2) = 1 + \lambda_1(0 - u_1) + \lambda_2(11 - u_2) = 1 - \lambda_1 + 10\lambda_2$ .

$L(\lambda)$  selects the cheapest one of these, i.e.,  $L(\lambda) = \min\{10, 1 + 10\lambda_1 - \lambda_2, 1 - \lambda_1 + 10\lambda_2\}$ . This surface is shown in fig. 6; note there are three linear surfaces corresponding to the three deterministic policies (the flat one is  $\pi_0$ , and the sloped ones are  $\pi_1$  and  $\pi_2$ ). If we consider  $\lambda = 0$  as a starting point then  $L([0 \ 0]) = 1$  and we can see that there is no way to get a larger  $L(\lambda)$  by moving in a single coordinate. Algebraically,

$$L([x \ 0]) = \min\{10, 1 + 10x, 1 - x\} = 1 - x < 1$$

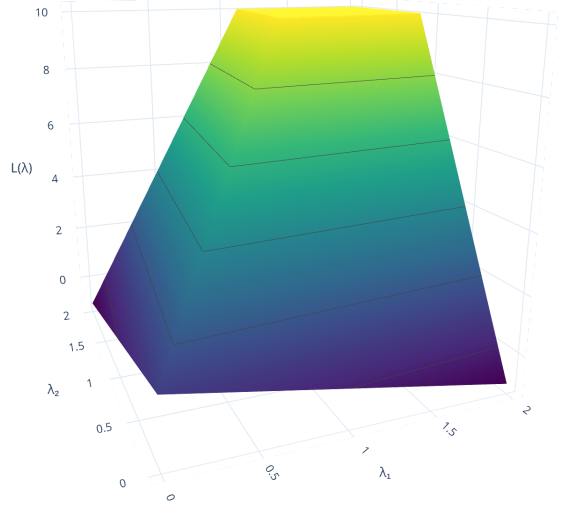
and symmetrically

$$L([0 \ x]) = \min\{10, 1 - x, 1 + 10x\} = 1 - x < 1$$

for all  $x \in \mathbb{R}_{>0}$ . Thus, coordinate search gets stuck at  $\lambda = [0 \ 0]$  with  $L(\lambda) = 1$  and fails to progress to an optimal solution such as  $\lambda^* = [2 \ 2]$  with  $L(\lambda^*) = 10$ .

## D Bellman Backups over Vector Value Function

Instead of using a scalar value function  $V : S \rightarrow \mathbb{R}_{\geq 0}$ , we use a vector value function  $V : S \rightarrow \mathbb{R}_{\geq 0}^{n+1}$ , with entries for each cost function. This modification still solves the  $\mathbb{S}(\lambda)$  SSP by projecting



**Figure 6.** Surface of  $L(\lambda)$  for the pathological CSSP where coordinate search fails to find  $\lambda^*$ .

the vector value function onto  $V_\lambda(s) = [1 \ \lambda] \cdot V(s)$ , but we now apply Bellman backups on  $V$  directly. We write  $Q(s, a) = C(a) + \sum_{s' \in S} P(s'|s, a)V(s')$ , then the  $\lambda$  modified Bellman backups are given as follows:

**Definition 2** ( $\lambda$  Bellman backup). *Given the vector value function  $V : S \rightarrow \mathbb{R}_{\geq 0}^{n+1}$  and a scalarisation  $\lambda \in \mathbb{R}_{\geq 0}^n$ , a vector Bellman backup applied to state  $s$  sets*

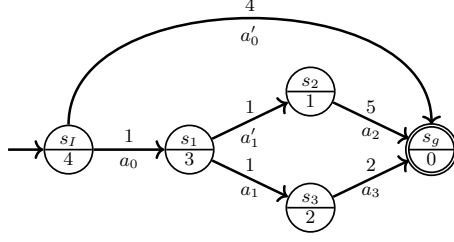
$$V(s) \leftarrow Q(s, a_{\min}) \quad \text{for} \quad a_{\min} \in \operatorname{argmin}_{a \in A(s)} ([1 \ \lambda] \cdot Q(s, a)).$$

By applying  $\lambda$  Bellman backups, we can solve  $\mathbb{S}(\lambda)$  and evaluate its policy for each cost function  $C_i(\pi_\lambda^*)$  at the same time. This idea was used by Hong and Williams [22].

The  $\lambda$  Bellman backups introduce two technical issues. First, to ensure that each  $C_i(\pi_\lambda^*)$  is evaluated accurately, we must redefine the Bellman residual so that each value function converges, otherwise the backups may converge to an optimal vector value function for the  $\mathbb{S}(\lambda)$  without the scalar value functions converging. We can do this by considering  $\max_{i=0}^n |V_i(s) - Q(s, a_{\min})|$  for  $s$ 's  $a_{\min}$ . The second technical issue follows from the redefined Bellman residual, because it is possible to get stuck in cycles where actions with similar scalarised  $Q$ -values change individual value functions significantly, so that the new residual never goes to 0. To address this, we need to use tie-breaking rules, e.g., if two actions have costs within  $\epsilon$  of each other, iterate through the value function indices and pick the action that has a smaller  $Q$ -value first.

## E Strong $\epsilon$ -consistency vs $\epsilon$ -consistency

For unconstrained SSPs, if a value function  $V$  is  $\epsilon$ -consistent [7], then it is guaranteed to induce a deterministic optimal policy (as  $\epsilon \rightarrow 0$ ); however, it may not induce all optimal deterministic policies. To produce all optimal deterministic policies, we require strong



**Figure 7.** An SSP with a value function that is  $\epsilon$ -consistent but not strongly  $\epsilon$ -consistent, and fails to describe all optimal policies.

$\epsilon$ -consistency. We give an example where  $\epsilon$ -consistency fails to capture all optimal deterministic policies: consider the SSP in fig. 7, with the value function  $V$  specified in the bottom of each node. The SSP has two optimal policies:

- $\pi_0^* = \{s_I \mapsto a'_0\}$
- $\pi_1^* = \{s_I \mapsto a_0, s_1 \mapsto a_1, s_3 \mapsto a_3\}$ .

$V$  has two greedy policies: the optimal  $\pi_0^*$  and a suboptimal policy

- $\pi' = \{s_I \mapsto a_0, s_1 \mapsto a'_1, s_2 \mapsto a_2\}$ .

$V$  is  $\epsilon$ -consistent w.r.t. the greedy policy  $\pi_0^*$ , and  $V$  is admissible; so,  $V$  could be produced by an optimal SSP algorithm such as CG-iLAO\* [26] (which is what we use to solve our  $\mathbb{S}(\lambda)$  SSPs). Importantly,  $V$  is not  $\epsilon$ -consistent w.r.t. the other greedy policy  $\pi'$  because  $\text{RES}(s_2) = |1 - 5| > \epsilon$  (assume small  $\epsilon$  such as 0.0001), and therefore  $V$  is not strongly  $\epsilon$ -consistent. This has the effect that  $V$  fails to capture the second optimal policy  $\pi_1^*$ , i.e.,  $\pi_1^*$  is not a greedy policy w.r.t.  $V$ .

## F More Results

We present an extended table of results in tab. 2. It has the same entries as tab. 1, and additionally includes CARL with ideal-point heuristics and C-ROC. Recall that C-ROC is not admissible for CARL’s  $\mathbb{S}(\lambda)$ s, so it provides no guarantee of optimality nor convergence. Nevertheless, in our experiments, wherever CARL(C-ROC) has a coverage entry, it found the optimal solution.

## References

- [1] E. Altman. Constrained Markov Decision Processes. Technical report, INRIA, 1995. URL <https://inria.hal.science/inria-00074109>.
- [2] P. Baumgartner, S. Thiébaux, and F. Trevizan. Heuristic Search Planning With Multi-Objective Probabilistic LTL Constraints. In *Proc. of Int. Conf. on Principles of Knowledge Representation and Reasoning (KR)*, pages 415–424, 2018. URL <https://cdn.aaai.org/ocs/18039/18039-78663-1-PB.pdf>.
- [3] R. Bellman. *Dynamic programming*. Princeton University Press, 1957.
- [4] D. Bertsekas and J. Tsitsiklis. An analysis of stochastic shortest path problems. *Mathematics of Operations Research*, pages 580–595, 1991. doi: 10.1287/moor.16.3.580.
- [5] D. Bertsimas and J. Tsitsiklis. *Introduction to Linear Optimization*. Athena Scientific, 1997.
- [6] B. Bonet and H. Geffner. Planning as heuristic search: New results. In *Recent Advances in AI Planning*, pages 360–372, 2000. doi: 10.1007/10720246\_28.
- [7] B. Bonet and H. Geffner. Faster heuristic search algorithms for planning with uncertainty and full feedback. In *Proc. of Int. Conf. on Automated Planning and Scheduling (ICAPS)*, pages 1233–1238, 2003. URL <https://www.dtic.upf.edu/~hgeffner/html/reports/hdp.pdf>.
- [8] B. Bonet and H. Geffner. Labeled RTDP: Improving the Convergence of Real-Time Dynamic Programming. In *Proc. of Int. Conf. on Automated Planning and Scheduling (ICAPS)*, pages 12–21, 2003. URL <https://dl.acm.org/doi/abs/10.5555/3036969.3036972>.
- [9] B. Bonet and R. Givan. 5th International Planning Competition: Non-deterministic track. call for participation, 2006. URL [https://www.researchgate.net/publication/228724778\\_5th\\_International\\_Planning\\_Competition\\_Non-deterministic\\_track\\_call\\_for\\_participation](https://www.researchgate.net/publication/228724778_5th_International_Planning_Competition_Non-deterministic_track_call_for_participation).
- [10] D. Bryce and O. Buffet. International Planning Competition Uncertainty Part: Benchmarks and Results, 2008. URL <https://ipc08.icaps-conference.org/probabilistic/wiki/images/0/03/Results.pdf>.
- [11] D. Chen, F. Trevizan, and S. Thiébaux. Heuristic Search for Multi-Objective Probabilistic Planning. In *Proc. of AAAI Conf. on Artificial Intelligence*, pages 11945–11954, 2023. doi: 10.1609/aaai.v37i10.26409.
- [12] M. Ehrgott. *Multicriteria Optimization*. Springer London, Limited, 2006. doi: 10.1007/3-540-27659-9.
- [13] V. Forejt, M. Kwiatkowska, and D. Parker. Pareto Curves for Probabilistic Model Checking. In *Automated Technology for Verification and Analysis*, pages 317–332, 2012. doi: 10.1007/978-3-642-33386-6\_25.
- [14] F. Geißer, G. Pováda, F. Trevizan, M. Bondouy, F. Teichteil-Königsbuch, and S. Thiébaux. Optimal and heuristic approaches for constrained flight planning under weather uncertainty. *Proc. of Int. Conf. on Automated Planning and Scheduling (ICAPS)*, pages 384–393, 2020. doi: 10.1609/icaps.v30i1.6684.
- [15] F. Geißer, P. Haslum, S. Thiébaux, and F. Trevizan. Admissible heuristics for multi-objective planning. *Proc. of Int. Conf. on Automated Planning and Scheduling (ICAPS)*, pages 100–109, 2022. doi: 10.1609/icaps.v32i1.19790.
- [16] G. Y. Handler and I. Zang. A dual algorithm for the constrained shortest path problem. *Networks*, pages 293–309, 1980. doi: 10.1002/net.3230100403.
- [17] E. Hansen and I. Abdoulahi. Efficient bounds in heuristic search algorithms for stochastic shortest path problems. *Proc. of AAAI Conf. on Artificial Intelligence*, pages 3283–3290, 2015. doi: 10.1609/aaai.v29i1.9643.
- [18] E. Hansen and S. Zilberstein. LAO\*: A heuristic search algorithm that finds solutions with loops. *Artificial Intelligence*, pages 35–62, 2001. doi: 10.1016/S0004-3702(01)00106-0.
- [19] E. A. Hansen. Error bounds for stochastic shortest path problems. *Mathematical Methods of Operations Research*, pages 1–27, 2017. doi: 10.1007/s00186-017-0581-5.
- [20] P. Hart, N. Nilsson, and B. Raphael. A formal basis for the heuristic determination of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics*, pages 100–107, 1968. doi: 10.1109/tssc.1968.300136.
- [21] M. Helmert and C. Domshlak. Landmarks, Critical Paths and Abstractions: What’s the Difference Anyway? *Proc. of Int. Conf. on Automated Planning and Scheduling (ICAPS)*, pages 162–169, 2009. doi: 10.1609/icaps.v19i1.13370.
- [22] S. Hong and B. C. Williams. An anytime algorithm for constrained stochastic shortest path problems with deterministic policies. *Artificial Intelligence*, page 103846, 2023. doi: 10.1016/j.artint.2022.103846.
- [23] J. Lee, G.-H. Kim, P. Poupart, and K.-E. Kim. Monte-Carlo Tree Search for Constrained POMDPs. In *Advances in Neural Information Processing Systems (NIPS)*, 2018. URL [https://proceedings.neurips.cc/paper\\_files/paper/2018/file/54c3d58c5efcf59ddeb7486b7061ea5a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/54c3d58c5efcf59ddeb7486b7061ea5a-Paper.pdf).
- [24] I. Little and S. Thiébaux. Probabilistic planning vs. replanning. In *ICAPS Workshop on IPC: Past, Present and Future*, 2007. URL <http://users.cecs.anu.edu.au/~iaai/icaps07.pdf>.
- [25] Mausam and A. Kolobov. *Planning with Markov Decision Processes: An AI Perspective*. Springer International Publishing, 2012. doi: 10.1007/978-3-031-01559-5.
- [26] J. Schmalz and F. Trevizan. Efficient constraint generation for stochastic shortest path problems. *Proc. of AAAI Conf. on Artificial Intelligence*, pages 20247–20255, 2024. doi: 10.1609/aaai.v38i18.30005.
- [27] J. Schmalz and F. Trevizan. Finding optimal deterministic policies for constrained stochastic shortest path problems. In *Proc. of European Conf. on Artificial Intelligence (ECAI)*, pages 4148–4156, 2024. doi: 10.3233/FAIA240986.
- [28] J. Schmalz and F. Trevizan. Code and Benchmarks for ECAI 2025 paper “Solving Constrained Stochastic Shortest Path Problems with Scalarisation”, 2025. doi: 10.5281/zenodo.16871056.
- [29] H.-J. M. Shi, S. Tu, Y. Xu, and W. Yin. A primer on coordinate descent algorithms, 2017.
- [30] N. Z. Shor. *Minimization Methods for Non-Differentiable Functions*. Springer Berlin Heidelberg, 1985. doi: 10.1007/978-3-642-82118-9.
- [31] F. Trevizan, S. Thiébaux, P. Santana, and B. Williams. Heuristic Search in Dual Space for Constrained Stochastic Shortest Path Problems. In *Proc. of Int. Conf. on Automated Planning and Scheduling (ICAPS)*, pages 326–334, 2016. doi: 10.1609/icaps.v26i1.13768.
- [32] F. Trevizan, S. Thiébaux, and P. Haslum. Occupation Measure Heuris-

		$\lambda$ -ROC			IP-ROC			CARL C-ROC			$\lambda$ -LMcut			IP-LMcut			i-dual			i <sup>2</sup> -dual		
SAR (n, d, r)	(4, 3, 0.75)	<b>30</b>	<b>5.6± 1.4</b>	30	11.5± 2.3	30	12.1± 2.4	<b>30</b>	<b>4.3± 1.0</b>	30	9.2± 1.8	30	104.1± 51.0	30	45.8± 21.1	30	69.2± 56.4					
	(4, 4, 0.75)	<b>30</b>	<b>41.9± 5.3</b>	30	75.5± 7.6	30	79.5± 8.3	<b>30</b>	<b>34.1± 4.2</b>	30	65.3± 7.6	5	1345.7±534.1	9	979.8±296.6	8	495.5±417.8					
	(5, 3, 0.50)	<b>30</b>	<b>4.9± 0.9</b>	30	9.0± 1.6	30	9.3± 1.7	<b>30</b>	<b>4.8± 1.1</b>	30	10.0± 2.4	30	25.1± 8.5	30	13.1± 4.6	30	11.1± 4.0					
	(5, 3, 0.75)	<b>30</b>	<b>12.3± 2.5</b>	30	26.2± 5.8	30	27.9± 6.0	<b>30</b>	<b>13.3± 2.9</b>	30	33.1± 9.0	30	253.8±110.6	30	113.3± 46.0	30	168.0±105.9					
	(5, 4, 0.50)	<b>30</b>	<b>30.8± 6.8</b>	<b>30</b>	<b>52.3± 9.5</b>	<b>30</b>	<b>54.7± 9.8</b>	<b>30</b>	<b>36.2± 7.9</b>	<b>30</b>	<b>70.5± 12.9</b>	17	695.0±268.5	24	550.6±166.1	27	637.8±213.1					
	(5, 4, 0.75)	<b>30</b>	<b>126.6± 45.2</b>	<b>30</b>	<b>270.0± 80.7</b>	<b>30</b>	<b>282.9± 83.9</b>	<b>30</b>	<b>142.5± 49.3</b>	<b>30</b>	<b>353.9±105.7</b>	6	930.9±369.3	6	499.2±226.4	8	349.1±288.1					
Elev (e,w,h)	(1, 1, 1)	30	1.6± 0.6	30	2.0± 0.4	30	2.1± 0.3	30	1.2± 0.3	30	1.4± 0.3	30	1.5± 0.1	<b>30</b>	<b>0.8± 0.1</b>	30	6.2± 0.6					
	(1, 1, 2)	29	50.9± 26.8	29	42.7± 12.6	29	42.2± 11.9	29	35.0± 12.7	29	44.4± 14.1	30	25.3± 1.0	<b>30</b>	<b>21.1± 2.6</b>	30	280.0± 46.5					
	(1, 2, 1)	25	14.2± 5.6	25	17.7± 3.7	25	17.7± 3.7	25	13.1± 4.2	25	15.8± 4.3	30	13.9± 0.9	<b>30</b>	<b>9.8± 1.4</b>	30	174.2± 38.4					
	(1, 2, 2)	22	313.3±150.2	22	399.8± 99.9	22	407.9±102.0	22	336.9±115.8	22	460.1±118.5	<b>30</b>	<b>340.1± 35.7</b>	<b>30</b>	<b>392.5± 72.1</b>	1	989.2					
	(2, 1, 1)	<b>30</b>	<b>104.5± 23.2</b>	<b>30</b>	<b>131.7± 26.2</b>	<b>30</b>	<b>126.9± 25.1</b>	<b>30</b>	<b>121.8± 26.7</b>	<b>30</b>	<b>149.7± 31.5</b>	23	670.6±173.7	26	576.1±181.4	1	1799.7					
	(2, 1, 2)	<b>20</b>	<b>1321.5±161.1</b>	4	1281.0±114.2	5	1309.1±230.8	3	1344.2±341.5	0		1	1799.6	1	1799.1	0						
(2, 2, 1)	<b>28</b>	<b>569.2±131.7</b>	<b>28</b>	<b>830.8±160.8</b>	<b>28</b>	<b>796.4±158.7</b>	23	890.5±164.2	18	1157.4±180.4	5	1608.9±229.1	5	1444.4±352.4	0							
ExBW (id,N,c)	(01, 5, 0.1)	30	8.6± 0.1	30	8.5± 0.1	30	8.5± 0.1	30	1.6	<b>30</b>	<b>1.2</b>	30	97.2± 1.0	30	7.9± 0.1	30	626.6± 18.8					
	(02, 5, 0.07)	30	40.4± 0.3	30	36.0± 0.2	30	35.3± 0.3	30	18.7± 0.1	<b>30</b>	<b>15.2± 0.1</b>	0		30	732.9± 30.6	0						
	(03, 6, 0.91)	30	173.5± 5.1	30	138.0± 1.4	30	135.7± 1.2	30	7.5± 0.1	<b>30</b>	<b>3.3</b>	0		30	10.9± 0.2	0						
	(04, 6, 0.16)	30	161.2± 1.4	30	142.9± 1.0	30	142.1± 1.7	30	63.7± 0.5	<b>30</b>	<b>45.6± 0.3</b>	0		27	1462.2± 79.6	0						
	(05, 7, 0.01)	30	30.7± 0.2	30	28.1± 0.2	30	27.2± 0.3	30	25.8± 0.2	<b>30</b>	<b>12.8± 0.1</b>	30	82.7± 1.5	30	38.5± 0.6	30	451.0± 14.5					
	(06, 8, 0.3)	30	300.4± 10.9	30	276.0± 3.1	30	276.7± 4.1	30	91.6± 0.6	<b>30</b>	<b>48.2± 0.4</b>	0		0		0						
PARC (f,u)	(0.0, 1)	30	311.2± 13.6	30	183.7± 6.9	<b>30</b>	<b>1.5± 1.4</b>	30	270.9± 8.3	30	252.5± 8.7	30	126.2± 9.1	30	103.1± 2.4	30	42.8± 3.1					
	(0.0, ∞)	30	100.7± 0.8	30	52.1± 0.9	<b>30</b>	<b>0.1</b>	30	141.3± 1.0	30	64.7± 0.5	30	47.1± 2.5	30	42.1± 1.8	30	29.5± 1.8					
	(0.2, 1)	30	600.7± 3.6	30	234.8± 3.1	13	184.2± 3.5	30	863.4± 9.7	30	337.1± 3.8	0		0		<b>30</b>	<b>65.1± 11.1</b>					
	(0.2, ∞)	30	121.7± 1.9	30	91.0± 1.3	30	56.9± 0.6	30	185.7± 1.2	30	79.2± 0.6	30	1537.7± 72.8	0		<b>30</b>	<b>41.7± 4.9</b>					
	(0.4, 1)	30	882.8± 31.7	<b>30</b>	<b>247.5± 3.8</b>	<b>30</b>	<b>241.8± 3.2</b>	30	564.1± 12.8	30	308.2± 4.0	0		0		3	801.5±536.2					
	(0.4, ∞)	30	130.9± 2.5	30	98.0± 1.8	30	94.0± 1.7	30	188.1± 1.3	<b>30</b>	<b>82.0± 1.2</b>	0		0		3	665.9±308.3					
CTW (n,d,c)	(4, 4, 2)	<b>30</b>	<b>1.6± 0.1</b>	30	2.7± 0.1	30	3.2± 0.1	<b>30</b>	<b>1.8± 0.1</b>	<b>30</b>	<b>1.8± 0.1</b>	30	4.5± 0.1	30	24.2± 0.7	30	24.3± 0.8					
	(4, 4, 4)	<b>30</b>	<b>2.4</b>	30	4.6± 0.1	30	5.4± 0.1	30	3.2	30	3.2	30	8.2± 0.5	30	43.0± 1.0	30	36.1± 2.0					
	(4, 4, 6)	<b>30</b>	<b>3.1</b>	30	6.6± 0.1	30	7.4± 0.1	30	4.6± 0.1	30	4.7± 0.1	30	10.5± 0.4	30	56.0± 2.4	30	42.5± 1.5					
	(4, 5, 2)	<b>30</b>	<b>3.9± 0.2</b>	30	6.3± 0.3	30	7.4± 0.3	<b>30</b>	<b>4.2± 0.3</b>	<b>30</b>	<b>4.2± 0.2</b>	30	18.7± 0.8	30	119.1± 5.5	30	115.3± 6.3					
	(4, 5, 4)	<b>30</b>	<b>5.9± 0.1</b>	30	11.3± 0.3	30	13.4± 0.6	30	7.5± 0.1	30	7.6± 0.2	30	32.4± 1.0	30	239.5± 9.0	30	157.0± 5.3					
	(4, 5, 6)	<b>30</b>	<b>7.3± 0.3</b>	30	15.9± 0.4	30	18.6± 0.5	30	10.3± 0.4	30	11.1± 0.5	30	41.7± 1.0	30	310.6± 14.7	30	194.5± 7.1					
(4, 6, 2)	<b>30</b>	<b>9.3± 0.5</b>	30	14.9± 0.6	30	17.8± 0.8	<b>30</b>	<b>10.1± 0.6</b>	<b>30</b>	<b>9.8± 0.6</b>	30	90.9± 4.7	30	840.2± 47.1	30	937.7± 45.7						
(4, 6, 4)	<b>30</b>	<b>13.9± 0.3</b>	30	27.1± 0.8	30	31.3± 0.9	30	17.8± 0.3	30	18.5± 0.3	30	187.0± 5.9	13	1656.6± 76.8	30	1281.6± 64.0						
(4, 6, 6)	<b>30</b>	<b>18.2± 0.4</b>	30	37.7± 0.9	30	44.7± 0.9	30	24.9± 0.7	30	26.5± 0.6	30	225.1± 10.0	4	1424.3± 88.9	30	1238.9± 81.6						

**Table 2.** For the benchmark problems, we show each planner and heuristic’s coverage (out of 30) and over the converged runs the mean runtime (secs) with 95% C.I. The fastest planner and heuristic per problem are in bold.

tics for Probabilistic Planning. In *Proc. of Int. Conf. on Automated Planning and Scheduling (ICAPS)*, pages 306–315, 2017. doi: 10.1609/icaps.v27i1.13840.