

An Exact Algorithm for Side-Chain Placement in Protein Design

Stefan Canzar¹, Nora C. Toussaint², and Gunnar W. Klau¹

¹CWI, Life Sciences group, Science Park 123, 1098 XG Amsterdam, the Netherlands,
{stefan.canzar, gunnar.klau}@cw.nl

²University of Tübingen, Center for Bioinformatics, Applied Bioinformatics, Sand 14, 72076
Tübingen, Germany, toussaint@informatik.uni-tuebingen.de

Abstract

Computational protein design aims at constructing novel or improved functions on the structure of a given protein backbone and has important applications in the pharmaceutical and biotechnical industry. The underlying combinatorial side-chain placement problem consists of choosing a side-chain placement for each residue position such that the resulting overall energy is minimum. The choice of the side-chain then also determines the amino acid for this position. Many algorithms for this \mathcal{NP} -hard problem have been proposed in the context of homology modeling, which, however, reach their limits when faced with large protein design instances.

In this paper, we propose a new exact method for the side-chain placement problem that works well even for large instance sizes as they appear in protein design. Our main contribution is a dedicated branch-and-bound algorithm that combines tight upper and lower bounds resulting from a novel Lagrangian relaxation approach for side-chain placement. Our experimental results show that our method outperforms alternative state-of-the-art exact approaches and makes it possible to optimally solve large protein design instances routinely.

1 Introduction

Protein design aims at constructing novel or improved functions on the structure of a given protein backbone. Since proteins are key players in virtually all biological processes, the ability to design proteins is of great practical interest, *e.g.*, to the pharmaceutical and biotechnological industry. Experimental protein design methods, such as directed evolution [3], have been applied successfully. However, since experimental methods are time- and money-consuming, computational approaches are an attractive alternative.

Computational protein design is related to the side-chain placement (SCP) problem in protein homology modeling. Given the modeled backbone of a protein, the amino acid side-chains have to be placed on this backbone in the energetically most

favorable conformation. Two assumptions are commonly made: (i) side-chains adopt only statistically dominant low-energy side-chain conformations, the so-called rotamers [10], and (ii) the energy of a protein is the sum of intrinsic side-chain energies and pairwise interaction energies. These assumptions lead to the following discrete optimization problem: For each residue position choose a rotamer such that the total energy of the protein is minimum. This problem has been shown to be NP-hard [20] and inapproximable [5].

In protein design the candidate rotamers at each position do not only come from a single amino acid but from several potential amino acids, yielding very large problem instances. Previous *in silico* approaches to protein design differ in their choice of rotamer library, energy function, and optimization method. Utilization of a higher-resolution rotamer library and a more accurate energy function will improve the results. On the other hand, it will increase computation time and problem size. Regarding the optimization methods, computational protein design approaches generally employ computationally expensive heuristics such as the Monte Carlo method [7, 6, 21]. Other heuristics, which have been proposed for SCP in protein homology modeling, could also be applied [24, 25, 9, 28, 22]. However, Voigt *et al.* [23] have shown that these inexact algorithms become less accurate with increasing problem size. Thus, exact methods capable of solving large protein design instances are desirable. Several approaches to solving the SCP problem exactly have been proposed, including dead end elimination [8, 11, 19] (combined with systematic search [17] or residue reduction [26]), integer linear programming [1, 14], branch-and-bound [4, 24] and tree decomposition [27]. While most of these approaches work well for homology modeling, they reach their limits when applied to protein design.

In this paper, we propose a novel exact method for SCP that works well even for large instance sizes as they appear in protein design. After presenting the combinatorial problem formally in Section 2, we describe our new method in Section 3. Our main contribution is a dedicated branch-and-bound algorithm that combines tight upper and lower bounds resulting from a novel Lagrangian relaxation approach for SCP. In Section 4 we present and discuss our experimental results, in which we show that our method outperforms alternative state-of-the-art exact approaches and makes it possible to optimally solve large protein design instances routinely.

2 Combinatorial Problem Formulation and Notation

We study the following graph-theoretic formulation of the side-chain placement problem:

Problem 1 (SCP). *Given a k -partite graph $G = (V, E)$, $V = V_1 \cup \dots \cup V_k$, with node costs c_v , $v \in V$, and edge costs c_{uv} , $u, v \in E$, determine an assignment $a : \{1, \dots, k\} \rightarrow V$ with $a(i) \in V_i$, $1 \leq i \leq k$, such that the cost*

$$\sum_{i=1}^k c_{a(i)} + \sum_{i=1}^{k-1} \sum_{j=i+1}^k c_{a(i)a(j)}$$

of the induced graph is minimum.

Here, each node set V_i corresponds to the candidate rotamers for the residue set at position i . Node costs model self energies of rotamers and edge costs model interaction energies between pairs of rotamers. A solution is given by selecting for each residue position i , $1 \leq i \leq k$, exactly one rotamer $a(i)$. Clearly, the choice of the rotamer determines also the amino acid at this position.

In the description of our algorithm we will also use a function $r : V \rightarrow \{1, \dots, k\}$ that denotes the residue position of a rotamer v , that is, $r(v) = i$ if and only if $v \in V_i$.

3 Lagrangian Relaxation Based Branch-and-Bound

We now present our novel approach to solve the SCP problem to provable optimality. Its core is the computation of sharp upper and lower bounds using a novel Lagrangian relaxation technique within a dedicated branch-and-bound approach.

3.1 Upper and Lower Bounds by Lagrangian Relaxation

Our relaxation builds on an integer linear programming (ILP) formulation for SCP that has been introduced by Althaus *et al.* [1] and extended by Kingsford *et al.* [14]:

$$\min \sum_{v \in V} c_v x_v + \sum_{uv \in E'} c_{uv} y_{uv} \quad (1)$$

$$\text{s.t. } \sum_{v \in V_i} x_v = 1 \quad 1 \leq i \leq k \quad (2)$$

$$\sum_{u \in V_i} y_{uv} = x_v \quad 1 \leq i \leq k, \text{ for all } v \in V_j \text{ with } j \in \mathcal{N}_i^+ \quad (3)$$

$$\sum_{\substack{u \in V_i \\ c_{uv} < 0}} y_{uv} \leq x_v \quad 1 \leq i \leq k, \text{ for all } v \in V_j \text{ with } j \notin \mathcal{N}_i^+ \quad (4)$$

$$x_v \in \{0, 1\} \quad \text{for all } v \in V \quad (5)$$

$$y_{uv} \in \{0, 1\} \quad \text{for all } uv \in E' \quad (6)$$

The formulation contains binary variables x_v , for nodes $v \in V$, and y_{uv} , for edges $uv \in E$, with the interpretation that a variable is 1 if the corresponding node or edge is part of the induced subgraph and 0 otherwise. Constraints (2) express that exactly one rotamer must be chosen per residue position. Constraints (3) and (4) link node and edge variables. When a rotamer v of a residue position j is chosen, i.e., $x_v = 1$, exactly one incident edge from each other residue position $i \neq j$ must be chosen as well for residue positions i that share positively weighted edges with residue position j , i.e., $j \in \mathcal{N}_i^+ := \{\ell \in \{1, \dots, k\} \mid \exists uv \in E, c_{uv} > 0, u \in V_i, v \in V_\ell\}$. If the two residue positions i and j are linked only by non-positive edges, i.e., $j \notin \mathcal{N}_i^+$, the relaxed constraints (4) apply: zero-weighted edges do not have to be forced to be in the solution and the corresponding variables can be removed from the ILP. Let $E' := E \setminus \{uv \in E \mid u \in V_i, v \in V_j, j \notin \mathcal{N}_i^+, c_{uv} = 0\}$ be the set of remaining edges.

The distinction between pairs of residues i, j with $j \in \mathcal{N}_i^+$ and pairs i, j with $j \notin \mathcal{N}_i^+$ in constraints (3) and (4) leads to a considerably smaller number of variables in practice and to a much better performance. Nevertheless, it is not crucial for the understanding of our approach and we thus drop this distinction in the remainder of this work and treat all pairs of residues as in constraint (3) for the sake of clarity of the presentation, that is, without removing any variables. In our implementation, however, we treat constraints (3) and (4) differently.

While previous work [1, 14] focuses on solving the linear programming (LP) relaxation of (1)–(6), we propose a Lagrangian relaxation approach. In the case of SCP this leads to a much more efficient algorithm, because we exploit structural knowledge of the SCP problem. The idea of Lagrangian relaxation is to relax constraints of an intractable problem, *e.g.*, the SCP ILP, such that the relaxed problem can be solved efficiently. The relaxed constraints are moved to the objective function, penalized by so-called Lagrangian multipliers. An optimal solution of the original problem, *i.e.*, an energy-minimum choice of candidate rotamers, is also a solution of the relaxed problem, and every optimal solution of the relaxed problem provides a lower bound on the optimal score of the original problem. The Lagrangian multipliers are adjusted iteratively such that the lower bound increases gradually. Also, after each iteration, we can evaluate the solution of the relaxed problem and thus obtain a new upper bound to the SCP problem. During the iterative process, the lowest upper and highest lower bound move closer and closer together. If they coincide, a provably optimal SCP has been found. Otherwise, we stop the process after a fixed number of iterations and use the bounds within the branch-and-bound framework.

The key idea of our Lagrangian relaxation approach is to define a total order, denoted by $<$, on the residue positions, to split the constraints that link node and edge variables into a left and a right part and then relax the right part of the constraints. W.l.o.g. and for ease of notation we assume that the residue positions have already been ordered, that is, residue position i denotes the i th residue position according to $<$.

First, we rewrite constraints (3) as left and right parts, that is, (3) becomes

$$\sum_{u \in V_i} y_{uv} = x_v \quad 1 \leq i \leq k-1, \text{ for all } v \in V_j \text{ with } j > i \quad (7)$$

$$\sum_{u \in V_i} y_{uv} = x_v \quad 2 \leq i \leq k, \text{ for all } v \in V_j \text{ with } j < i \quad (8)$$

We dualize constraints (8) for all non-neighboring residues, *i.e.*, constraints for which $i > j + 1$, with Lagrangian multipliers λ_v^i . To simplify notation, we further introduce $\lambda_v^i := 0$ for neighboring residues, *i.e.*, for which $r(v) + 1 = i$ holds. For fixed Lagrangian multipliers λ_v^i we obtain the following relaxation, which we denote by

(LR_λ) :

$$\min \sum_{j=1}^k \sum_{v \in V_j} \left(c_v + \sum_{i>j+1} \lambda_v^i \right) x_v + \sum_{\substack{uv \in E \\ r(u) < r(v)}} (c_{uv} - \lambda_u^{r(v)}) y_{uv} \quad (9)$$

$$\text{s.t. } \sum_{v \in V_i} x_v = 1 \quad 1 \leq i \leq k \quad (10)$$

$$\sum_{u \in V_i} y_{uv} = x_v \quad 1 \leq i \leq k-1, \text{ for all } v \in V_j \text{ with } j > i \quad (11)$$

$$\sum_{u \in V_i} y_{uv} = x_v \quad 2 \leq i \leq k \text{ for all } v \in V_{i-1} \quad (12)$$

$$x_v, y_{uv} \in \{0, 1\} \quad (13)$$

Note that a distinction between pairs of residues according to constraints (3) and (4) requires the Lagrangian multipliers λ_v^i associated with constraints (4) to be restricted in sign in order to guarantee that an optimal solution to (LR_λ) yields a lower bound on the optimal score of the SCP problem.

An integral variable assignment that satisfies constraints (10), (12), and constraints (11) for neighboring residues, *i.e.*, for $j = i + 1$, encodes a path p in the corresponding k -partite graph from a node in V_1 to a node in V_k that traverses exclusively edges between neighboring residues. The remaining constraints of (11) involve y -variables that do not appear in any other constraint and can thus be chosen independently of each other. In other words, we can determine the best possible contribution of a vertex v to the overall objective value, under the assumption that v lies on path p , by simply picking for every residue $i < r(v) - 1$ the edge of minimum weight between v and a node in V_i . More formally, we define the *profit* δ of a node v as

$$\delta(v) = (c_v + \sum_{i>r(v)+1} \lambda_v^i) + \sum_{i=1}^{r(v)-2} \min_{u \in V_i} (c_{uv} - \lambda_u^{r(v)}) ,$$

where the first term in brackets denotes the coefficient of variable x_v in the objective function. Then the score of a feasible solution to (LR_λ) that induces a path $p = (v_1, v_2, \dots, v_k)$ with $v_i \in V_i$ in graph G is

$$\sum_{i=1}^k \delta(v_i) + \sum_{i=1}^{k-1} c_{v_i v_{i+1}} .$$

Let graph G' be derived from G by removing all edges between non-neighboring residues, *i.e.*, edges $u'v'$ with $|r(u') - r(v')| > 1$, and by defining the weights of the remaining edges uv as $c_{uv} + \delta(v)$. Then, an optimal solution to (LR_λ) corresponds to a shortest path in G' from a node in V_1 to a node in V_k , see Figure 1.

Theorem 1. *An optimal solution to (LR_λ) can be computed in time $\mathcal{O}(|V|^2)$.*

Proof. The profits of all nodes can clearly be computed in time $\mathcal{O}(|V|^2)$. Graph G' is acyclic and thus a shortest path can be computed in time linear in the number of

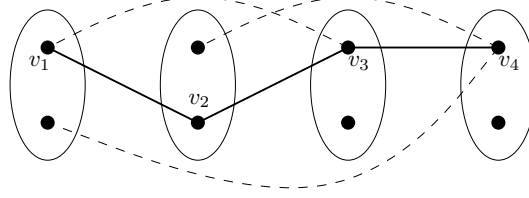


Figure 1: The structure of a feasible solution to (LR_λ) . The polygon drawn in solid line denotes the corresponding path $p = (v_1, v_2, v_3, v_4)$. Every node on path p has exactly one incident edge to every residue left of it, except to its direct neighbor, depicted by the dashed lines.

edges in G' , *i.e.*, $\mathcal{O}(\sum_{i=1}^{k-1} |V_i| \cdot |V_{i+1}|)$. Note that a topological sorting of the vertices is implicitly given by the k -partition. \square

We apply a standard *subgradient optimization* technique [12] to find those Lagrangian multipliers λ_v^i that yield the largest lower bound to our relaxation. This iterative adaption of the Lagrangian multipliers only requires the profits of a small fraction of the vertices to be recomputed from scratch in each iteration. In practice, the shortest path computation for given profits by a simple dynamic programming scheme dominates the overall running time needed to resolve the Lagrangian relaxation (LR_λ) for modified multipliers λ_v^i . In other words, the running time will be linear in the number of edges in G' rather than G .

In order to improve the practical performance of our approach we sort the residues by increasing number of rotamers. This ordering results in a minimum total number of dualized constraints.

3.2 Branch-and-Bound

We embed our Lagrangian bounding scheme into a branch-and-bound framework to obtain an energy-minimum rotamer assignment. The general idea is to divide the overall SCP problem into easier subproblems by fixing rotamers at individual residue positions and to solve the resulting problems recursively. To avoid a complete enumeration of all possible rotamer assignments, we employ our Lagrangian bounds to prune large parts of the enumeration tree. In particular, let S_k be a subproblem in which certain residue positions have been assigned a rotamer, and let \hat{x} be an arbitrary solution to the original SCP problem. If the lower bound \underline{z}^k on the minimal energy assignment for subproblem S_k is larger than $z(\hat{x})$, *i.e.*, the total energy of rotamer assignment \hat{x} , then no optimal solution to the original problem can be obtained from S_k and we can prune the subtree rooted at S_k from the search space.

Branching scheme. Starting from an SCP problem instance S , a straightforward branching rule is to impose the constraint $x_{ij} = 1$ in the left child node of the search tree, and $x_{ij} = 0$ in the right, *i.e.*, fixing and forbidding a rotamer j at residue position i , respectively. However, since every rotamer j is contained in a constraint (2) for a

residue position i , this leads to an unbalanced search tree, because the right child node leaves $k - 1$ possible rotamers for residue position i , whereas the left child leaves only one possibility. While partitioning the set of rotamers of a given residue position into two roughly equally sized sets avoids this imbalance, we experienced, however, a significantly smaller number of nodes in the tree with the following, alternative branching scheme. Instead of creating two subproblems, we create one for each rotamer of a selected residue position i , by fixing rotamer j_k in subproblem k . Only for very large design instances, we first partitioned the sets of possible rotamers that were larger than some threshold p , into two smaller sets. The effectiveness of this scheme is mainly based on two properties. First, when fixing a rotamer j for a residue position i , we reduce the problem instance by incorporating the interaction energy between rotamer j and any other rotamer j' of all residue positions $i' \neq i$ into the self energy of j' . In contrast, in our relaxation only the profits δ of rotamers of residue positions $i' > i$ would take into account a further subdivision of the set of rotamers of residue position i . The rotamer assignment to residue positions $i' < i$ would still rely on a correct choice of the incoming edges for the remaining rotamers of residue position i , which could only be accomplished by iteratively adapting the Lagrangian penalties. Second, a large enough set of child nodes will give our depth-first search traversal of the branch-and-bound tree the freedom to pick a promising node first.

Choosing a constraint. The question remains which position constraint (2) to choose for a branching step. We adopt the idea of *strong branching* [2]. The rough idea is to estimate the progress, *i.e.*, increase in lower bound, for the residue positions before actually branching on one of them. This is done by successively fixing each rotamer of a given residue position and solving the resulting Lagrangian subproblem. Based on the progress of the single rotamers, we compute an overall score of the residue position, see below, and pick the one with the highest score as the next residue position to branch on. Since the computation time per node of this procedure would be enormous, we try to estimate the locally best residue position by simplifying this search in three different aspects.

- First, we restrict the evaluation to the most promising residue positions. More precisely, we order the residue positions in increasing order of their maximum (primal) fractional value of its rotamers. We recover the primal solutions by taking the convex combinations of the last k solutions x^t produced in the course of the subgradient optimization. Then, at a node of depth l of the branch-and-bound tree, we consider the first $\gamma(l)$ percent of the sorted residue positions. Note that the concept of a residue position with minimal maximal fractional value of its rotamer variables can be considered as a generalization of the *most infeasible branching* rule for binary variables, to constraints of the form (2).
- Second, to estimate the increase of the objective function when fixing a rotamer, only a few subgradient iterations are performed, along with an aggressive multiplier adjustment.

- Finally, we choose our scoring function of the residue positions in such a way, that they can be computed quickly while still giving a good estimate on the overall progress in the dual (lower) bound. Let \mathcal{Q} be the subproblem corresponding to the current node of the tree and let problem \mathcal{Q}_i^j be obtained from \mathcal{Q} by fixing rotamer j in residue r_i . Then the score of a residue position r_i is given by

$$\xi(r_i) = \min_{j \in r_i} \Delta_i^j,$$

where $\Delta_i^j = z(\mathcal{Q}_i^j) - z(\mathcal{Q})$. To determine the score of a given residue position i , we test the rotamers in decreasing order of their fractional values. The goal is to evaluate rotamers j with small progress Δ_i^j first, since the subgradient optimization for the remaining rotamers can be aborted as soon as their increase in objective function exceeds the smallest progress seen so far. Also notice that sorting the residue positions as described above is beneficial for the computation of the residue scores. Whenever we encounter a rotamer with a progress Δ_i^j that is smaller than the smallest progress determined for a previous residue position $i' \neq i$, we do not have to consider the remaining rotamers of i' , since we are interested in the residue position with maximal score.

Choosing a node. A primal feasible solution that gives a good upper bound on the minimum total energy is necessary to prune the enumeration tree significantly. Therefore we follow a depth-first search (DFS) strategy to descend in the branch-and-bound tree as quickly as possible, increasing the chances of finding a new and hopefully better feasible solution. Furthermore, the Lagrangian subproblems corresponding to a node and to one of its immediate descendants differ only in one residue position. Therefore, a subproblem can be resolved faster when starting from the multiplier vector λ determined in the immediate parent node. On the downside, once being in a wrong branch, one may spend a long time in this subtree before getting back on a path leading to an improved solution. We thus combine the advantages of DFS and a best-node first strategy. We fix the rotamers of a given residue position in increasing order of their dual (lower) bounds. Following this approach led to a considerably smaller number of nodes evaluated in the tree, and we were able to find the optimal solution much faster in most of the cases.

4 Experimental Results

We have implemented our combinatorial algorithm for finding an optimal solution to the Lagrangian relaxation (LR_λ) in C++ using the LEDA and BALL libraries [18, 13]. We iteratively improve the obtained lower bounds on the optimal solution of the SCP problem by applying a standard subgradient approach. In each iteration, we derive from the Lagrangian solution a feasible solution to the original problem and thus an upper bound on the optimal score by evaluating the subgraph induced by the nodes lying on the shortest path from V_l to V_k , see Section 3.1. We exploit the upper and lower bounds in a branch-and-bound manner to prune large parts of the search space and to derive a provably optimal solution to the SCP problem.

In order to determine an initial upper bound for the branch-and-bound framework, we employ a simple local search procedure: Given an initial configuration in which each residue position is assigned the rotamer with the lowest self energy, residue positions are selected randomly and optimized, *i.e.*, the respective position is assigned the rotamer yielding the best energy within the current conformation. This minimization proceeds until the energy could not be improved several times in a row or a maximum of 100 iterations is reached.

The only state-of-the-art exact method for SCP that can cope with protein design instances is the ILP based method proposed by Kingsford *et al.* [14]. Available software packages for DEE or treewidth based approaches such as R3 [26] or TreePack [27] do not allow several candidate amino acids at each position and are thus not applicable to protein design instances. Furthermore, our experiments show that even small protein design instances already have treewidths of 10 to 20 as compared to 3 to 4 for most homology modeling instances [27]. Since the complexity of the TreePack algorithm grows exponentially in the treewidth, a reasonable performance on protein design instances is not to be expected. For DEE-based methods, a similar argument holds because reduced protein design instances are still too large to be processed in reasonable time by residue unification or other enumeration techniques. We therefore compare our Lagrangian based approach only to an implementation that solves the ILP proposed by Kingsford *et al.* [14] using CPLEX 12.2¹ with Concert Technology.

In our experiments, we used two different benchmark sets. The first set consists of protein design energy files provided by Kingsford *et al.* [14]. It comprises 25 proteins with 11 to 124 flexible residue positions. Surface residues are fixed. At each core position up to six different amino acids are allowed. The employed energy function comprises statistical potentials and van der Waals interactions. We omit the experimental results on the simpler homology modeling instances, since almost all of these instances can be solved in a fraction of a second by both our Lagrangian relaxation approach and the CPLEX based method. The second set of protein design instances was taken from Yanover *et al.* [28]. This set comprises 97 proteins with 40 to 180 amino acids. All residue positions are flexible and at each position all 20 amino acids are allowed yielding very large problem instances. Here, the more realistic Rosetta energy function [16] was used to determine self and interaction energies.

In a preprocessing phase, we apply established rules [11, 26] to decrease the size of the problem instances while preserving optimality properties. Tables 1 and 2 show the running times of our Lagrangian relaxation branch-and-bound approach and the CPLEX based method using default settings on the resulting instances on a compute cluster with two 2.26 GHz Intel Quad Core processors with 24 GB of RAM on each node, running 64 bit Linux. We applied a time limit of 12 hours and a memory limit of 16 GB. Computations exceeding one of these limits were aborted.

The first three columns of the table give the characteristics of the instances, *i.e.*, their PDB identifier, the number of residues and the total number of rotamers. The following two columns give characteristics of the branch-and-bound proce-

¹<http://www.cplex.com>

Name	Instance		Lagrangian B&B			CPLEX	S
	#res	#rot	N	H	time/s	time/s	
laac	105	1523	2	1	1.73	3.40	2.0
laho	64	981	1	0	0.01	0.02	2.0
1b9o	123	2056	3	1	2.09	2.56	1.2
1c5e	95	1108	1	0	0.12	0.25	2.1
1c9o	66	1130	2	1	0.33	1.96	5.9
1cc7	72	1396	1	0	0.28	0.59	2.1
1cex	197	2556	9	2	13.37	33.25	2.5
1cku	85	1093	1	0	0.03	0.09	3.0
1ctj	89	1021	1	0	0.07	0.27	3.9
1cz9	139	2332	1	0	3.7	18.10	4.9
1czp	98	1170	1	0	0.54	4.32	8.0
1d4t	104	1636	1	0	0.37	2.36	6.4
1igd	61	926	1	0	0.01	0.02	2.0
1mfm	153	2134	25	5	21.89	145.63	6.7
1plc	99	1156	2	1	1.50	6.08	4.1
1qj4	256	4080	313	10	8,424.56	31,636.40	3.8
1qq4	198	2045	16	4	32.56	38.89	1.2
1qtn	152	2516	1	0	1.50	3.22	2.1
1qu9	126	1817	2	1	0.31	0.66	2.1
1rcf	169	2396	2	1	4.76	12.85	2.7
1vfy	67	939	1	0	0.01	0.01	1.0
2pth	193	3077	66	6	322.28	518.51	1.6
3lzt	129	2074	7	2	3.20	10.64	3.3
5p21	166	2874	52	4	106.09	115.01	1.1
7rsa	124	1958	1	0	0.78	3.31	4.2

Table 1: Running times of our Lagrangian relaxation branch-and-bound approach and the CPLEX based method on the design instances from [14]. We further give the number of residues (#res) and the total number of rotamers (#rot) of the instance, the number of nodes (N) and height (H) of the branch-and-bound tree as well as the speedup S (ratio of running times).

duration: Columns N and H give the total number of evaluated nodes and the height of the branch-and-bound tree, respectively. The remaining columns give the running times in seconds of our Lagrangian based approach and the CPLEX based method as well as the ratio of running times S. Note that we include the time spent in the local search heuristic for the Lagrangian branch-and-bound approach.

On the first dataset, our Lagrangian based approach outperforms the state-of-the-art CPLEX based method on all 25 instances. The small number of nodes evaluated in the course of the branch-and-bound procedure indicates sharp lower and upper bounds derived from the Lagrangian solutions. On the more challenging second dataset, our method could solve 52 of the 97 instances within 12 h, whereas the CPLEX based method could only finish 12 instances within the time and memory limits.

Name	Instance		Lagrangian B&B			CPLEX	S
	#res	#rot	N	H	time/s	time/s	
lbrf	44	3524	9	4	293.97	469.87	1.6
lbr7	25	1048	1	0	0.54	5.77	10.7
ld3b	66	5732	1	0	530.37	9,577.68	18.1
len2	59	2689	1	0	19.41	39.94	2.1
lezg	58	1653	2	1	185.11	441.23	2.4
lg6x	51	3190	1	0	23.96	160.64	6.7
lgcq	65	5442	4	2	903.82	5,270.08	9.8
li07	52	3186	4	1	187.45	166.20	0.9
lkth	49	3330	18	4	798.57	642.42	0.8
lrb9	43	3307	7	2	127.93	9,535.72	74.5
lsem	54	4348	192	8	5,020.55	6,470.37	1.3
lvfy	58	3951	16	2	2,540.86	†	n/a
4rxn	45	3636	1	0	220.33	3,034.57	13.8
la8o	62	4510	6	2	1,418.71	†	n/a
lb67	66	5543	27	4	3,822.09	†	n/a
lbbz	52	3935	7	2	1,329.26	†	n/a
lb4	60	5289	12	3	1,875.32	†	n/a
lc75	63	4323	25	2	7,175.69	†	n/a
lcc8	69	6515	26	2	16,508.10	†	n/a
ld3b	66	5732	1	0	530.37	†	n/a
lfr3	61	5100	22	4	6,997.76	†	n/a
lgut	62	4945	22	2	7,745.17	†	n/a
lhg7	65	5047	5	2	987.51	†	n/a
li27	69	5934	39	6	4,070.20	†	n/a
ligd	60	5207	18	4	3,163.14	†	n/a
ligq	53	4582	18	5	4,294.16	†	n/a
liqz	75	5412	15	2	2,137.58	†	n/a
lj75	55	4861	14	4	5,704.83	†	n/a
ljo8	54	4680	41	4	1,830.23	†	n/a
lkq1	58	5244	7	1	3,990.65	†	n/a
ll9l	69	5518	4	2	1,514.50	†	n/a
lidd	71	6383	8	1	4,582.84	†	n/a
lljo	69	6428	14	3	5,624.10	†	n/a
lmhn	53	4454	3	2	570.15	†	n/a
lnkd	56	4148	9	4	1,119.56	†	n/a
loai	56	4330	73	3	20,021.10	†	n/a
lplc	92	7955	34	4	24,308.50	†	n/a
lpwt	58	4876	2	1	886.94	†	n/a
lr69	60	4926	49	2	27,862.50	†	n/a
lwap	65	5551	14	3	5,267.68	†	n/a
2igd	59	5262	13	3	6,332.26	†	n/a
lc4q	65	5598	843	10	30,789	†	n/a
lc9o	60	5305	97	8	11,627.20	†	n/a
lctj	84	6232	265	10	7,083.65	†	n/a
ldj7	69	5571	22	4	12,581.90	†	n/a
le0b	58	4715	147	6	30,840.10	†	n/a
lerv	101	9150	76	8	31,539.70	†	n/a
lfk5	81	5714	28	5	14,625.30	†	n/a
lg2b	59	4926	129	8	6,880.45	†	n/a
lmgq	71	6250	4	1	2,613.06	†	n/a
lvie	56	4803	3	111	866.54	†	n/a

Table 2: Running times of our Lagrangian relaxation branch-and-bound approach and the CPLEX based method on the design instances from [28]. Instances which cannot be solved by both approaches within a time limit of 12 hours are omitted. The † sign indicates that a computation exceeded either the time limit (12 h) or the memory limit (16 GB).

5 Conclusions and Outlook

We have constructed a Lagrangian relaxation of the Kingsford ILP formulation of the SCP problem that allowed us to obtain strong bounds by solving a modified shortest path problem on the underlying k -partite graph. By utilizing these bounds within a branch-and-bound framework we achieved running times that outperform a state-of-the-art exact method that uses the professional mathematical programming solver CPLEX. Our implementation of the Lagrangian branch-and-bound approach as well as the data sets used in this paper are freely available as the package SCP of the planet lisa software library [15].

Future work on exact side-chain placement should explore possible connections to the recently introduced method by Sontag *et al.* [22], which is based on belief propagation. This heuristic algorithm is currently the best non-exact method and finds optimal solutions astonishingly often. In our opinion, underlying ideas from the area of belief propagation may be useful also in a truly exact method.

The mathematical model of the SCP problem as studied in this work appears in a wide range of applications including image understanding, error correcting codes, and frequency assignment in telecommunications. We believe that our approach can be applied successfully in these areas, too.

Acknowledgements. Much of the work has been carried out while NCT visited CWI on a CWI internship. The authors thank Inken Wohlers and Oliver Kohlbacher for useful discussions. Computational experiments were sponsored by the NCF for the use of supercomputer facilities, with financial support from NWO.

References

- [1] E. Althaus, O. Kohlbacher, H.-P. Lenhof, and P. Müller. A combinatorial approach to protein docking with flexible side chains. *J Comput Biol*, 9(4):597–612, 2002.
- [2] D. Applegate, R. Bixby, V. Chvátal, and W. Cook. Finding cuts in the TSP. Technical Report 95-05, DIMACS, 1995.
- [3] J. D. Bloom, M. M. Meyer, P. Meinhold, C. R. Otey, D. MacMillan, and F. H. Arnold. Evolving strategies for enzyme engineering. *Curr Opin Struct Biol*, 15(4):447–452, 2005.
- [4] A. A. Canutescu, A. A. Shelenkov, and R. L. Dunbrack. A graph-theory algorithm for rapid protein side-chain prediction. *Protein Sci*, 12(9):2001–2014, Sep 2003.
- [5] B. Chazelle, C. Kingsford, and M. Singh. A semidefinite programming approach to side chain positioning with new rounding strategies. *INFORMS J. on Computing*, 16(4):380–392, 2004.

- [6] G. Dantas, C. Corrent, S. Reichow, J. Havranek, Z. Eletr, N. Isern, B. Kuhlman, G. Varani, E. Merritt, and D. Baker. High-resolution structural and thermodynamic analysis of extreme stabilization of human procarboxypeptidase by computational protein design. *Journal of Molecular Biology*, 366(4):1209–1221, 2007.
- [7] G. Dantas, B. Kuhlman, D. Callender, M. Wong, and D. Baker. A large scale test of computational protein design: folding and stability of nine completely redesigned globular proteins. *J Mol Biol*, 332(2):449–460, 2003.
- [8] J. Desmet, M. D. Maeyer, B. Hazes, and I. Lasters. The dead-end elimination theorem and its use in protein side-chain positioning. *Nature*, 356(6369):539–542, 1992.
- [9] J. Desmet, J. Spriet, and I. Lasters. Fast and accurate side-chain topology and energy refinement (FASTER) as a new method for protein structure optimization. *Proteins*, 48(1):31–43, 2002.
- [10] R. L. Dunbrack. Rotamer libraries in the 21st century. *Curr Opin Struct Biol*, 12(4):431–440, 2002.
- [11] R. F. Goldstein. Efficient rotamer elimination applied to protein side-chains and related spin glasses. *Biophys J*, 66(5):1335–1340, 1994.
- [12] M. Held and R. Karp. The traveling salesman problem and minimum spanning trees: part II. *Mathematical Programming*, 1:6–25, 1971.
- [13] A. Hildebrandt, A. K. Dehof, A. Rurainski, A. Bertsch, M. Schumann, N. C. Toussaint, A. Moll, D. Stöckel, S. Nickels, S. C. Mueller, H.-P. Lenhof, and O. Kohlbacher. BALL – biochemical algorithms library 1.3. *BMC Bioinformatics*, 11:531, 2010.
- [14] C. L. Kingsford, B. Chazelle, and M. Singh. Solving and analyzing side-chain positioning problems using linear and integer programming. *Bioinformatics*, 21(7):1028–1036, 2005.
- [15] G. W. Klau et al. planet lisa software library. <http://planet-lisa.net>.
- [16] B. Kuhlman and D. Baker. Native protein sequences are close to optimal for their structures. *Proc Natl Acad Sci U S A*, 97(19):10383–10388, 2000.
- [17] A. R. Leach and A. P. Lemon. Exploring the conformational space of protein side chains using dead-end elimination and the A* algorithm. *Proteins*, 33(2):227–239, 1998.
- [18] K. Mehlhorn and S. Näher. *The LEDA Platform of Combinatorial and Geometric Computing*. Cambridge University Press, Cambridge, 1999.
- [19] N. A. Pierce, J. A. Spriet, J. Desmet, and S. L. Mayo. Conformational splitting: A more powerful criterion for dead-end elimination. *Journal of Computational Chemistry*, 21:999–1009, 2000.

- [20] N. A. Pierce and E. Winfree. Protein design is NP-hard. *Protein Eng*, 15(10):779–782, 2002.
- [21] P. S. Shah, G. K. Hom, S. A. Ross, J. K. Lassila, K. A. Crowhurst, and S. L. Mayo. Full-sequence computational design and solution structure of a thermostable protein variant. *J Mol Biol*, 372(1):1–6, 2007.
- [22] D. Sontag, T. Meltzer, A. Globerson, T. Jaakkola, and Y. Weiss. Tightening LP relaxations for MAP using message passing. In D. McAllester and P. Myllymak, editors, *Conference on Uncertainty in Artificial Intelligence*. AUA Press, Corvallis, Oregon, 2008.
- [23] C. Voigt, D. Gordon, and S. Mayo. Trading accuracy for speed: a quantitative comparison of search algorithms in protein sequence design. *Journal of Molecular Biology*, 299(3):789–803, 2000.
- [24] L. Wernisch, S. Hery, and S. Wodak. Automatic protein design with all atom force-fields by exact and heuristic optimization. *Journal of Molecular Biology*, 301(3):713–736, 2000.
- [25] Z. Xiang and B. Honig. Extending the accuracy limits of prediction for side-chain conformations. *J Mol Biol*, 311(2):421–430, 2001.
- [26] W. Xie and N. V. Sahinidis. Residue-rotamer-reduction algorithm for the protein side-chain conformation problem. *Bioinformatics*, 22(2):188–194, 2006.
- [27] J. Xu and B. Berger. Fast and accurate algorithms for protein side-chain packing. *J. ACM*, 53(4):533–557, 2006.
- [28] C. Yanover, T. Meltzer, and Y. Weiss. Linear programming relaxations and belief propagation—an empirical study. *J Mach Learn Res*, 7:1887–1907, 2006.