

Guaranteed Diversity and Optimality in Cost Function Network Based Computational Protein Design Methods[†]

Manon Ruffini¹, Jelena Vucinic², Simon de Givry¹, George Katsirelos³, Sophie Barbe², Thomas Schiex^{1,*}

¹ Université Fédérale de Toulouse, ANITI, INRAE, UR 875, Toulouse, France

² TBI, Université de Toulouse, CNRS, INRAE, INSA, ANITI, Toulouse, France

³ MIA-Paris - Mathématiques et Informatique Appliquées, INRAE, Paris, France

* Correspondence: thomas.schiex@inrae.fr

[†] This paper is an extended version of our paper published in the Proceedings of the 2019 IEEE 31st International Conference on Tools with Artificial Intelligence, Portland, OR, USA, 4-6 Nov. 2019.

Abstract: Proteins are the main active molecules of Life. While natural proteins play many roles, as enzymes or antibodies for example, there is a need to go beyond the repertoire of natural proteins to produce engineered proteins that precisely meet application requirements, in terms of function, stability, activity or other protein capacities. Computational Protein Design aims at designing new proteins from first principles, using full-atom molecular models. However, the size and complexity of proteins require approximations to make them amenable to energetic optimization queries. These approximations make the design process less reliable and a provable optimal solution may fail. In practice, expensive libraries of solutions are therefore generated and tested. In this paper, we explore the idea of generating libraries of provably diverse low energy solutions by extending Cost Function Network algorithms with dedicated automaton-based diversity constraints on a large set of realistic full protein redesign problems. We observe that it is possible to generate provably diverse libraries in reasonable time and that the produced libraries do enhance the Native Sequence Recovery, a traditional measure of design methods reliability.

Keywords: Computational Protein Design; Graphical Models; Automata; Cost Function Networks; Structural Biology; Diversity.

Citation: Ruffini M; Vucinic J; de Givry S; Katsirelos G; Barbe S.; Schiex T Guaranteed Diversity and Optimality for Computational Protein Design. *Algorithms* 2021, 1, 0. <https://doi.org/>

Received:
Accepted:
Published:

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Copyright: © 2021 by the authors. Submitted to *Algorithms* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Proteins are complex molecules that govern much of how cells work, in humans, plants, and microbes. They are made of a succession of simple molecules called α -amino acids. All α -amino acids share a common constant linear body and a variable side-chain. The side-chain defines the nature of the amino acid. There are 20 natural amino acid types, each having a distinct side-chain offering specific physico-chemical properties. In a protein, the successive amino acids are connected one to the other by peptidic bonds, defining a long linear polymer called the protein backbone. In solution, most proteins fold into a 3D shape, determined by the physico-chemical properties of their amino acid side-chains. Because of their large variety of functions, and their potentials for applications in medicine, environment, biofuels, green chemistry, etc, new protein sequences are sought, that present desired new or enhanced properties and functions. As function is closely related to three-dimensional (3D) structure [1], Computational Protein Design (CPD) methods aim at finding a sequence that folds into a target 3D structure, that corresponds to the desired properties and functions. A general formulation of this problem being highly intractable, simplifying assumptions have been made (see Figure 1): the target protein structure (or backbone) is often assumed to be rigid, the continuous space of flexibility of amino acids side-chains is represented as a discrete set of conformations called rotamers and the atomic forces that control the protein stability are represented as a decomposable energy function, defined as the sum

36 of terms involving at most two bodies (amino acids). The problem of design is then
37 reduced to a purely discrete optimization problem: given a rigid backbone, one must
38 find a combination of discrete side-chain natures and conformations (rotamers) that
39 minimizes the energy. The resulting sequence and associated side-chain conformations
40 define the Global Minimum Energy Conformation (GMEC) for the target backbone.
41 A rotamer library for all 20 natural amino acids containing typically few hundreds of
42 conformations, the discrete search space becomes very quickly challenging to explore
43 and the problem has been shown to be NP-hard [2] (decision NP-complete). It has
44 been naturally approached by stochastic optimization techniques such as Monte Carlo
45 simulated annealing [3], as in the commonly used Rosetta software [4]. Such stochastic
46 methods offer only asymptotic optimality guarantees. Another possible approach is
47 to use provable optimization techniques, that instead offer finite-time deterministic
48 guarantees. In the last decade, Constraint programming-based algorithms for solving
49 the Weighted Constraint Satisfaction Problem (WCSP) on Cost Function Networks
50 (CFN) have been proposed to tackle CPD instances [5,6]. These provable methods have
51 shown unprecedented efficiency at optimizing decomposable force fields on genuine
52 protein design instances [6], leading to successfully characterized new proteins [7]. Cost
53 Function Networks are one example of a larger family of mathematical models that aim
54 at representing and analyzing decomposable functions, called Graphical Models [8,9].

55 Even if provable methods definitely remove the possibility of failed optimization,
56 they cannot fight the simplifying assumptions that appear in the CPD problem formula-
57 tion. First, the optimized pairwise decomposed energetic criterion only approximates
58 the actual molecule energy. Then, the rigid backbone and discrete-chain conformations
59 ignore the actual continuous protein flexibility [10]. Ultimately, even with a perfect
60 energy function, an alternative backbone structure may well exist that gives the GMEC
61 sequence an even better energy. This usually requires expensive post-hoc filtering based
62 on structure prediction (forward folding [11]). Therefore, even with provable methods, a
63 library of highly relevant mutants is usually produced for further experimental testing,
64 with the hope that the probability of identifying a functional protein will be increased.
65 Provable Branch and Bound-based WCSP algorithms have the native ability of enumer-
66 ating solutions within a threshold of the optimal solution. Empirically, one can observe
67 that the set of sequences that lie within this threshold grows very quickly in size with
68 the energy threshold, but is mostly composed of sequences that are very similar to the
69 optimal GMEC sequence. Ideally, a design library should be a set of low energy but
70 also diverse solutions. With yeast-display capacity to simultaneously express and test
71 thousands of proteins, libraries of diversified designs become increasingly important.
72 The hope is that sequence diversity will improve the likelihood that a protein endowed
73 of desired function is found. In this paper, we are therefore interested in algorithmic
74 methods that can provide such a set of guaranteed diverse low energy solutions and
75 then to empirically check if enforcing diversity in a library while optimizing energy does
76 improve the protein design process.

77 Because of their important applications, protein sequences can be subject to patents.
78 Ideally, a newly designed sequence should satisfy a *minimum* Hamming distance con-
79 straint to existing patented sequences. Specific design targets may require to escape
80 known patterns such as, *e.g.*, antigenic sub-sequences that would be recognized by
81 the Major Histocompatibility Complexes [12]. This again raises the need to produce
82 sequences satisfying minimum distance requirement to given sequences. Finally, CPD
83 input structures often come from existing, experimentally resolved natural (or native)
84 proteins. In this case, a native sequence exists, that has usually acquired desirable prop-
85 erties following the billions of years of natural evolution and optimization it has been
86 through. In many cases, to avoid disrupting the native protein properties (*e.g.* catalytic
87 capacities), the protein designer may want to bound the maximum number of mutations
88 introduced in the design. This raises the need to produce sequences satisfying *maximum*
89 distance requirement to given sequences.

90 In this paper, given an initial rigid backbone, we consider the problem of producing
91 a set of diverse low energy sequences that also provably satisfy a set of minimum and
92 maximum distance requirements w.r.t. given sequences. We observe that, beyond struc-
93 tural biology and bioinformatics, this problem of producing a set of diverse solutions
94 of a Graphical Model has been considered by many authors, either on discrete Boolean
95 Graphical Models such as Constraint Networks (used in Constraint Programming), or
96 on stochastic graphical models such as Markov Random Fields. While this shows that
97 the interest for the problem of diverse solutions generation goes well beyond Compu-
98 tational Protein Design, we observe that these approaches either offer no guarantee,
99 or are limited to specific tractable sub-classes of functions, such as sub-modular func-
100 tions [13]. Our approach instead relies on the reduction of distance requirements to
101 discrete automaton-based constraints that can be decomposed and compressed into
102 three-bodies (ternary) or two-bodies (binary) terms, using suitable dual and hidden
103 encoding [14,15]. These constraints can then be processed natively by existing WCSP
104 algorithms. While our approach is general and generally applicable to the production of
105 libraries of solutions of arbitrary discrete graphical models, its design is motivated by
106 Computational Protein Design. We therefore empirically evaluate this approach for the
107 generation of a library of diverse sequences on a set of protein design problems. We first
108 observe that this approach is capable of producing provably diverse sets of solutions on
109 Computational Protein Design problems of realistic sizes in reasonable time. Going back
110 to our initial aim, we also observe that sufficiently diverse libraries do offer better Native
111 Sequence Recovery rates (NSR), a usual metric for protein design methods evaluation
112 that measures how well it is able to reproduce Nature's optimization.

113 2. Computational Protein Design

114 A CPD instance is first composed of an input target 3D structure, defined by the
115 Cartesian coordinates of all the atoms in the protein backbone. The target protein struc-
116 ture can come from an existing protein backbone, that was determined experimentally
117 on an existing protein; or from a model that can be derived from existing 3D structures
118 of similar proteins; or from a completely new structure, as it is done in *de novo* design.
119 Once a backbone has been chosen, the design space must be fixed. One may choose to do
120 a full redesign, where the amino acids at all positions of the protein can be changed, or
121 redesign only a subset of all positions, focusing on positions that are key for the targeted
122 function. Overall, each position in the protein sequence will be set by the designer
123 as either *fixed*, *flexible*, or *mutable*. If the position is *fixed*, the side-chain is fixed and
124 rigid: the amino acid type and orientation are determined in the input target structure.
125 If the position is *flexible*, the residue type is fixed to the amino acid type of the input
126 structure, but the side-chain might adopt several conformations in space. If the position
127 is *mutable*, all or a restricted set of amino acid types are allowed at the position, along
128 with different conformations of their side-chain. Because of the supposed rigidity of the
129 backbone, the sequence-conformation search space is therefore characterized by two
130 decision levels: the sequence space, which corresponds to all the possible sequences \mathbf{s}
131 enabled by the mutable positions, and the conformation space, which must be searched
132 to identify the best side-chain conformation at each flexible or mutable position. The
133 possible conformations, or rotamers, for each amino acid are indexed in *rotamer libraries*,
134 such as the Dunbrack [16] or the Penultimate libraries [17]. Each library gathers a finite
135 set of conformations, capturing a representative subset of all frequently adopted con-
136 formations in experimentally determined structures. In the Rosetta design software [4],
137 that relies on the Dunbrack library, a fully mutable position will be typically associated
138 with around 400 possible rotamers. Designing a 10-residue peptide actually requires the
139 exploration of $400^{10} \approx 10^{26}$ conformations.

Given a backbone structure and a rotamer library, the CPD problem seeks a stable and functional sequence-conformation. The protein functionality is assumed to result from its conformation and its stability is captured by an energy function E , that allows

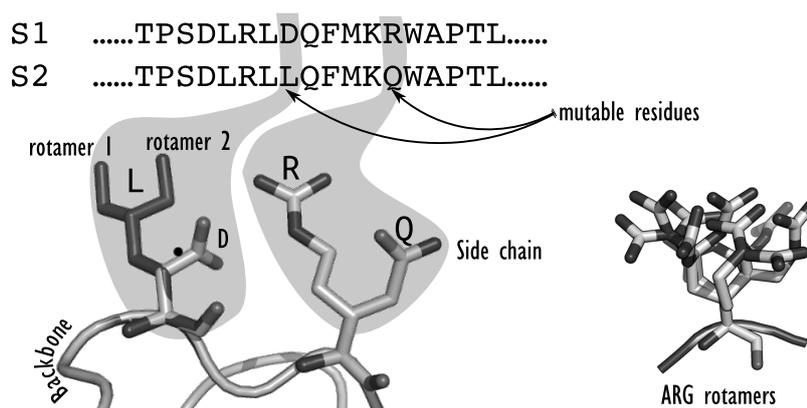


Figure 1. An example of two protein sequences (top) where two mutable amino acids have been redesigned. At the first position, the amino acid D (an aspartic acid) has been changed to a L (leucine), in a specific conformation (orientation). At the second position, the arginine R, with its very long and flexible side chain has been changed to a glutamine Q. The figure on the right illustrates the potential flexibility of the long arginine side chain, showing a sample of several possible superimposed conformations, representing a fraction of all possible conformations for an arginine side chain in existing rotamer libraries.

to compute the energy of any sequence-conformations on the target backbone. The task at hand is the minimization of this energy function. The optimal sequence is the best possible choice for the target rigid backbone. To model the energy, score functions are used. They can be physics-based, as the energetic force fields AMBER [18] and CHARMM [19]. They capture various atomic interactions including bond and torsion angle potentials, van der Waals potentials, electrostatic interactions, hydrogen bonds forces and entropic solvent effects. Score functions may also be enriched by “knowledge-based” energy terms, that result from the statistical analysis of known protein structures. For instance, Rosetta ref2015 and beta_nov_16 score functions [4,20] also integrate rotamer log-probabilities of apparition in natural structures, as provided in the Dunbrack library, in a specific energy term. To be optimized, the energy function should be easy to compute while remaining as accurate as possible, so as to predict relevant sequences. To try to meet these requirements, *additive pairwise decomposable* approximations of the energy have been chosen for protein design approaches [6,21]. The decomposable energy E of a sequence-conformation $\mathbf{r} = (r_1, \dots, r_n)$ where r_i is the rotamer used at the position i in the protein sequence can be written as:

$$E(\mathbf{r}) = E_{\emptyset} + \sum_{1 \leq i \leq n} E_i(r_i) + \sum_{1 \leq i < j \leq n} E_{ij}(r_i, r_j)$$

140 The term E_{\emptyset} is a constant that captures interactions within the rigid backbone. For
 141 $1 \leq i \leq n$, the unary (or one body) terms E_i capture the interactions between the rotamer
 142 r_i at position i and the backbone, as well as interactions internal to the rotamer r_i . For
 143 $1 \leq i < j \leq n$, the binary terms E_{ij} capture the interactions between rotamers r_i and
 144 r_j at positions i and j respectively. These energy terms only vary with the rotamers,
 145 thanks to the rigid backbone assumption. Protein design dedicated software, such as
 146 OSPREY [22] or Rosetta [4], compute all the constant, unary, and binary energy terms, for
 147 each rotamer and combination of rotamers, for each position and pair of positions. While
 148 this requires quadratic time in the protein length in the worst-case, distance-cutoffs
 149 make these computations essentially linear in this length. Once computed, these values
 150 are stored in energy matrices. During the exploration of the sequence-conformation
 151 space, conformation energies can be efficiently computed by summing the relevant
 152 energy terms fetched from the energy matrix. CPD methods aim at finding the optimum

153 conformation, called the *global minimum energy conformation* (GMEC). Despite all these
154 simplifications, this problem remains decision NP-complete [23].

155 3. CPD as a Weighted Constraint Satisfaction Problem

A Cost Function Network (CFN) \mathcal{C} is a mathematical model that aims at representing functions of many discrete variables that decompose as a sum of simple functions (with small arity or concise representation). It is a member of a larger family of mathematical models called graphical models [9], that all rely on multivariate function decomposability. A CFN is defined as a triple $\mathcal{C} = (\mathbf{X}, \mathbf{D}, \mathbf{C})$ where $\mathbf{X} = (X_1, \dots, X_n)$ is a set of variables, $\mathbf{D} = (\mathbf{D}_1, \dots, \mathbf{D}_n)$ is a set of finite domains, and \mathbf{C} is a set of cost functions. Each variable X_i takes its values in the domain \mathbf{D}_i . Each cost function $c_{\mathbf{S}} \in \mathbf{C}$ is a non negative integer function that depends on the variables in \mathbf{S} , called the scope of the function. Given a set of variables $\mathbf{S} \subset \mathbf{X}$, the set $\mathbf{D}_{\mathbf{S}} \prod_{X_i \in \mathbf{S}} \mathbf{D}_i$ denotes the Cartesian product of the domains of the variables in \mathbf{S} . For a tuple $\mathbf{t} \in \mathbf{Y}$, with $\mathbf{S} \subset \mathbf{Y} \subset \mathbf{X}$, the tuple $\mathbf{t}[\mathbf{S}]$ denotes the projection of \mathbf{t} over the variables of \mathbf{S} . A cost function $c_{\mathbf{S}} \in \mathbf{C}$ maps tuples of $\mathbf{D}_{\mathbf{S}}$ to integer costs in $\{0, \dots, \top\}$. In this paper, we assume, as is usual in most graphical models, that the default representation of a cost function $c_{\mathbf{S}}$ is a multidimensional cost table (or tensor) that contains the cost of every possible assignment of the variables in its scope. This representation requires space that grows exponentially in the cost function arity $|\mathbf{S}|$ which explains why arity is often assumed to be at most two. The joint function is defined as the bounded sum of all cost functions in \mathbf{C} :

$$\begin{aligned} C_{\mathcal{C}} : \mathbf{D}_{\mathbf{X}} &\longrightarrow \{0, \dots, \top\} \\ \mathbf{t} &\longmapsto \sum_{c_{\mathbf{S}} \in \mathbf{C}} c_{\mathbf{S}}(\mathbf{t}[\mathbf{S}]) \end{aligned}$$

where the bounded sum $+^{\top}$ is defined with $a +^{\top} b = \min(a + b, \top)$. The maximum cost $\top \in \mathbb{N} \cup \{+\infty\}$ is used for forbidden partial assignments and represents a sort of infinite or unbearable cost. Cost functions that take their values in $\{0, \top\}$ represent hard constraints. The Weighted Constraint Satisfaction Problem consists of finding the assignment of all the variables in \mathbf{X} with minimum cost:

$$\mathbf{t}^* = \min_{\mathbf{t} \in \mathbf{D}_{\mathbf{X}}} C_{\mathcal{C}}(\mathbf{t}) = \min_{\mathbf{t} \in \mathbf{D}_{\mathbf{X}}} \sum_{c_{\mathbf{S}} \in \mathbf{C}} c_{\mathbf{S}}(\mathbf{t}[\mathbf{S}])$$

156 Notice that when the maximum cost $\top = 1$, the cost function network becomes a
157 constraint network [24], where cost functions encode only constraints. Tuples that are
158 assigned cost 0 are valid, i.e., they satisfy the constraint, and tuples that are assigned
159 cost $1 = \top$ are forbidden. The Constraint Satisfaction Problem then consists of finding a
160 satisfying assignment, one that satisfies all the constraints.

Graphical models also encompass stochastic networks, such as discrete Markov Random Fields (MRF) and Bayesian Nets [25]. A discrete *Markov Random Field* is a graphical model $\mathcal{M} = (\mathbf{X}, \mathbf{D}, \Phi)$ where $\mathbf{X} = (X_1, \dots, X_n)$ is a set of random variables, $\mathbf{D} = (\mathbf{D}_1, \dots, \mathbf{D}_n)$ is a set of finite domains, and Φ is a set of potential functions. A potential function $\varphi_{\mathbf{S}}$ maps $\mathbf{D}_{\mathbf{S}}$ to $[0, +\infty]$. The *joint potential function* is defined as:

$$\begin{aligned} P = \Phi_{\mathcal{M}} : \mathbf{D}_{\mathbf{X}} &\longrightarrow [0, +\infty] \\ \mathbf{t} &\longmapsto \prod_{\varphi_{\mathbf{S}} \in \Phi} \varphi_{\mathbf{S}}(\mathbf{t}[\mathbf{S}]) \end{aligned}$$

161 Instead of the sum used to combine functions in Cost Function Networks, the product
162 is used here. The normalization of the potential function P by the partition function
163 $Z = \sum_{\mathbf{t} \in \mathbf{D}_{\mathbf{X}}} P(\mathbf{t})$ defines the probability function $p = \frac{1}{Z} P$ of the MRF. The Maximum A
164 Posteriori probability corresponds to the assignment with maximum probability (and
165 maximum potential) $\max_{\mathbf{t}} p(\mathbf{t})$.

166 MRF can also be expressed using additive energy functions $e_{\mathbf{S}} \in \mathbf{E}$, a logarithmic
167 transformation $e_{\mathbf{S}} = -\log \varphi_{\mathbf{S}}$ of the potential functions. The potential function is then
168 an exponential of the energy $P(\mathbf{t}) = \exp(-\sum_{e_{\mathbf{S}} \in \mathbf{E}} e_{\mathbf{S}})$. While potential functions are

169 multiplied together, energies simply add up. Therefore, Cost function networks are
 170 closely related to the energetic expression of Markov random fields. The main difference
 171 lies in the the fact that CFNs deal with non negative integers only, whereas MRF energies
 172 are real numbers. If $\top = +\infty$, a CFN can be transformed into an MRF through an
 173 exponential transformation, and given a precision factor, an MRF can be transformed
 174 into a CFN through a log transform. Zero potentials are mapped to cost \top and minimum
 175 energy means maximum probability.

176 Given a cost function network, the weighted constraint satisfaction problem can
 177 be answered by the exploration of a search tree in which nodes are CFNs induced
 178 by *conditioning*, i.e., domain reductions on the initial network. A *Branch-and-Bound*
 179 strategy is used to explore the search tree, that relies on a *lower bound* on the optimum
 180 assignment cost in the current sub-tree. If the lower bound is higher than the best
 181 joint cost found so far, it means that no better solution is to be found in the subtree,
 182 and it can be pruned. Each time a new solution is found, the maximum cost \top is
 183 updated to the corresponding cost, as we are only interested in finding solutions with
 184 lower cost. Ordering strategies are crucial and can lead to huge improvement in the
 185 empirical computation time: decisions should be made, that lead to low cost solutions
 186 (and decrease the maximum cost \top) and that enable early pruning.

187 The efficiency of the branch-and-bound strategy relies on strength of the lower
 188 bound on solution costs. In CFNs, since cost functions are non-negative, the empty-
 189 scoped cost function c_{\emptyset} provides a naive lower bound on the optimum cost. To efficiently
 190 compute tight lower bounds, local reasoning strategies are used, that aim at pushing
 191 as much cost as possible in the constant cost c_{\emptyset} , for better pruning. They are based on
 192 *equivalence preserving transformations* that perform cost transfers between cost functions,
 193 while maintaining the solutions' joint costs unchanged [26], i.e., the joint cost function is
 194 preserved (these operations are called reparameterizations in MRFs). Specific sequences
 195 of equivalence preserving transformations can be applied to a CFN to improve the lower
 196 bound c_{\emptyset} until a specific target property is reached on the CFN. These properties, called
 197 *local consistencies*, aim at creating zero costs, while improving c_{\emptyset} . These sequences of
 198 operations should converge to a fixpoint (closure). Various local consistency properties
 199 have been introduced, as node consistency, arc consistency, full directional arc consis-
 200 tency, existential directional arc consistency, or virtual arc consistency [26]. For binary
 201 CFNs, that involve functions of arity at most two, these properties can be enforced in
 202 polynomial time in the size of the network. Among these, Virtual arc consistency has
 203 been shown to solve the WCSP on networks composed of submodular functions [27].
 204 Note however that the complexity of local consistency enforcing remains exponential in
 205 the arity of the cost functions (as is the size of the corresponding tensors).

206 Sometimes, one may need to include specific functions with a large scope in the
 207 description of the joint function. Because of the exponential size of the tensor description,
 208 these *global cost functions* must be represented with dedicated concise descriptions and
 209 require dedicated algorithms for local consistency propagation. More formally, a global
 210 cost function, denoted $\text{GCF}(\mathbf{S}, \mathcal{A})$, is a family of cost functions, with scope \mathbf{S} and possible
 211 parameters \mathcal{A} . A global cost function is said to be *tractable* when its minimum can be
 212 computed in polynomial time.

The CFN formulation of computational protein design is straightforward [5,6,28,29]
 (see Figure 2): given a CPD instance with pairwise decomposable energy function
 $E = E_{\emptyset} + \sum_{1 \leq i \leq n} E_i + \sum_{1 \leq i < j \leq n} E_{ij}$, let $\mathcal{C} = (\mathbf{X}, \mathbf{D}, \mathbf{C})$ be a cost function network with
 variables $\mathbf{X} = (X_1, \dots, X_n)$, with one variable X_i for each flexible or mutable position i
 in the protein sequence, domains $\mathbf{D} = (\mathbf{D}_1, \dots, \mathbf{D}_n)$ where the domain \mathbf{D}_i of the variable
 X_i consists of the available rotamers at position i , i.e., the amino acid types and their
 side-chain conformations. The cost functions are the empty-scoped, unary and binary
 energy terms:

$$\mathbf{C} = \{ E_{\emptyset} \} \cup \{ E_i, 1 \leq i \leq n \} \cup \{ E_{ij}, 1 \leq i < j \leq n \}$$

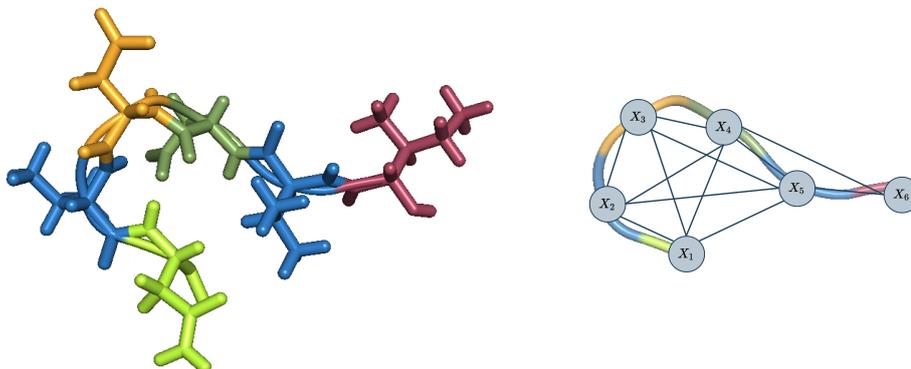


Figure 2. Input backbone and cost function network representation of a corresponding CPD instance with 6 mutable or flexible residues.

213 Energy terms, which are floating point numbers, are transformed into non-negative
 214 integer values by being shifted by a constant, multiplied by a large precision factor, and
 215 having their residual decimal truncated.

216 In this encoding, variables represent rotamers, combining information on the nature
 217 and the geometry (conformation) of the side-chain. In practice, it is often useful to
 218 add extra variables that represent the sequence information alone. The CFN \mathcal{C} can be
 219 transformed into $\mathcal{C}' = (\mathbf{X}', \mathbf{D}', \mathcal{C}')$, that embeds sequence variables, as follows:

220 **Variables** We add sequence variables to the network: $\mathbf{X}' = \mathbf{X}^{seq} \cup \mathbf{X}$, where $\mathbf{X}^{seq} =$
 221 $\{X_i^{seq} \mid X_i \in \mathbf{X}\}$. The value of X_i^{seq} represents the amino acid type of the rotamer
 222 value of X_i .

223 **Domains** $\mathbf{D} = \mathbf{D}^{seq} \cup \mathbf{D}$ where $\mathbf{D}^{seq} = \{\mathbf{D}_i^{seq} \mid \mathbf{D}_i \in \mathbf{D}\}$ where the domain \mathbf{D}_i^{seq} of X_i^{seq}
 224 is the set of available amino acid types at position i .

225 **Constraints** The new set of cost functions \mathcal{C}' is made of the initial functions \mathcal{C} ; and
 226 sequence constraints, that ensure that X_i^{seq} is the amino acid type of rotamer X_i .
 227 Such a function $c_{X_i, X_i^{seq}}$ just forbids (map to cost \top) pairs of values (r, a) where the
 228 amino acid identity of rotamer r does not match a . All other pairs are mapped to
 229 cost 0.

230 Such sequence variables depend functionally on the rotamer variables. They do not
 231 modify the search space and merely offer a facility to define properties on the protein
 232 sequence, if needed, as will be the case here.

233 4. Diversity and Optimality

234 In this section, we assume that we have a CFN $\mathcal{C} = (\mathbf{X}, \mathbf{D}, \mathcal{C})$ and we want to
 235 express diversity requirements on the values taken by variables in $\mathbf{S} \subset \mathbf{X}$. In the case of
 236 CPD, these variables will be the previously introduced *sequence* variables.

237 4.1. Measuring diversity

238 The task at hand is the production of a *set of diverse* and *low cost* solutions of \mathcal{C} . First,
 239 we need a measure of diversity between pairs of solutions, and among sets of solutions.

240 The simplest diversity measure between pairs of solutions is the *Hamming distance*,
 241 defined hereafter. It counts the number of variables in \mathbf{S} that take different values in two
 242 solutions. In the CPD framework, sequence variables represent amino acid identities:
 243 the Hamming distance measures the number of introduced mutations (or substitutions),
 244 a very usual notion in protein engineering.

Definition 1. Given a set of variables $\mathbf{S} \subset \mathbf{X}$ and two assignments \mathbf{t} and $\mathbf{t}' \in \mathbf{D}_{\mathbf{S}}$, the Hamming distance between \mathbf{t} and \mathbf{t}' is defined as follows:

$$d_H(\mathbf{t}, \mathbf{t}') = \sum_{X_i \in \mathbf{S}} \mathbb{1}(\mathbf{t}[X_i] \neq \mathbf{t}'[X_i])$$

245 The Hamming distance can be generalized to take into account dissimilarity scores
246 between values. The resulting distance is a semi-metric, defined as a sum of variable-wise
247 dissimilarities, as follows:

Definition 2. Given a zero-diagonal symmetric positive matrix D , that defines value dissimilarities, and two assignments $\mathbf{t}, \mathbf{t}' \in \mathbf{D}_{\mathbf{S}}$, the weighted-Hamming distance between \mathbf{t} and \mathbf{t}' is defined as:

$$d_D(\mathbf{t}, \mathbf{t}') = \sum_{X_i \in \mathbf{S}} D(\mathbf{t}[X_i], \mathbf{t}'[X_i])$$

248 In computational biology, protein sequences are often compared using dedicated
249 similarity matrices, such as the BLOSUM62 matrix [30]. A protein similarity matrix S
250 can be transformed into a dissimilarity matrix D such that $D_{i,j} = (S_{i,i} + S_{j,j})/2 - S_{i,j}$.

251 **Definition 3.** Given a set \mathbf{Z} of solutions, we define:

- 252 • its average dissimilarity: $\bar{d}(\mathbf{Z}) = \frac{2}{|\mathbf{Z}|(|\mathbf{Z}| - 1)} \sum_{\mathbf{t} \neq \mathbf{t}' \in \mathbf{Z}} d(\mathbf{t}, \mathbf{t}')$
- 253 • its minimum dissimilarity: $\check{d}(\mathbf{Z}) = \min_{\mathbf{t} \neq \mathbf{t}' \in \mathbf{Z}} d(\mathbf{t}, \mathbf{t}')$

254 We are aiming at producing a library of solutions that is guaranteed to be diverse.
255 The *average dissimilarity* does not match this need: a set of solutions might have a satisfy-
256 ing average dissimilarity value, with several occurrences of the same assignment, and
257 one or a few very dissimilar ones. So, to guarantee diversity, the *minimum dissimilarity*
258 will be the diversity measure used throughout this paper.

259 So, producing a set of *diverse* solutions requires that all solution pairs have their
260 distance above a given threshold. This can be encoded in cost functions representing
261 constraints, taking their values in $\{0, \top\}$ only:

Definition 4. Given two sets of variables \mathbf{S}, \mathbf{S}' of the same cardinality, a dissimilarity matrix D and a diversity threshold δ , we define the global cost function:

$$\begin{aligned} \text{DIST}(\mathbf{S}, \mathbf{S}', D, \delta) : \mathbf{D}_{\mathbf{S}} \times \mathbf{D}_{\mathbf{S}'} &\longrightarrow \{0, \top\} \\ (\mathbf{t}, \mathbf{t}') &\longmapsto \begin{cases} 0 & \text{if } \text{sign}(\delta) \cdot d(\mathbf{t}, \mathbf{t}') \geq \delta \\ \top & \text{otherwise.} \end{cases} \end{aligned}$$

Allowing both positive and negative threshold δ allows the DIST cost function to express either minimum or maximum diversity constraints. When $\delta > 0$, the cost function expresses a minimum dissimilarity requirement between the assignments \mathbf{t} and \mathbf{t}' :

$$\text{DIST}(\mathbf{t}, \mathbf{t}', D, \delta) = 0 \Leftrightarrow d(\mathbf{t}, \mathbf{t}') \geq \delta$$

If $\delta < 0$, the cost function represents the fact that \mathbf{t} and \mathbf{t}' must be similar, with a dissimilarity lower than the absolute value of δ :

$$\text{DIST}(\mathbf{t}, \mathbf{t}', D, \delta) = 0 \Leftrightarrow -d(\mathbf{t}, \mathbf{t}') \geq \delta \Leftrightarrow d(\mathbf{t}, \mathbf{t}') \leq -\delta = |\delta|$$

262 If needed, both maximum and minimum requirements can be imposed using two
263 constraints.

264 4.2. Diversity given sequences of interest

265 In the CPD context, minimum and maximum distance requirements with known
266 sequences may be useful in practice in at least two situations.

- 267 • A native functional sequence \mathbf{s}_{nat} is known for the target backbone. The designer
268 wants that less than δ_{nat} mutations be introduced on some sensitive region of the
269 native protein, in order to avoid disrupting a crucial protein property.
- 270 • A patented sequence \mathbf{s}_{pat} exists for the same function, and sequences with more
271 than δ_{pat} mutations are required for the designed sequence to be usable without
272 requiring a license.

273 The distance here is the Hamming distance based on matrix H which equals 1
274 everywhere, except for its zero diagonal. Using sequence variables, the following
275 diversity constraint-encoding cost functions need to be added to the CFN model:

- 276 • $\text{DIST}(\mathbf{X}^{seq}, \mathbf{s}_{nat}, H, -\delta_{nat})$
- 277 • $\text{DIST}(\mathbf{X}^{seq}, \mathbf{s}_{pat}, H, \delta_{pat})$

278 4.3. Sets of diverse and good quality solutions

279 The problem of producing a set of diverse and good quality solutions, *i.e.*, such that
280 all pairs of solutions satisfy the diversity constraint, and the solutions have minimum
281 cost, can be expressed as follows:

282 **Definition 5** (DIVERSESET). *Given a dissimilarity matrix D , an integer M and a dissimilarity*
283 *threshold δ , the problem DIVERSESET(C, D, M, δ) consists in producing a set \mathbf{Z} of M solutions*
284 *of C such that:*

285 **Diversity** *For all $\mathbf{t} \neq \mathbf{t}' \in \mathbf{Z}$, $d(\mathbf{t}, \mathbf{t}') \geq \delta$, *i.e.*, $\text{DIST}(\mathbf{t}, \mathbf{t}', D, \delta) = 0$.*

286 **Quality** *The solutions have minimum cost, *i.e.* $\sum_{\mathbf{t} \in \mathbf{Z}} C_C(\mathbf{t})$ is minimum.*

287 For a CFN C with n variables, solving DIVERSESET requires to simultaneously
288 decide the value of nM variables. It can be solved by making M copies of C with variable
289 sets \mathbf{X}^1 to \mathbf{X}^M and adding $\frac{M \cdot (M - 1)}{2}$ constraints $\text{DIST}(X^i, X^j, D, \delta)$ for all $1 \leq i < j \leq M$.
290 If the upper bound \top is finite, all its occurrences must be replaced by $M \cdot (\top - 1) + 1$.
291 While very elegant, this approach yields a CFN instance where the number of variables
292 is multiplied by the number of wanted solutions. The WCSP and CPD problems being
293 NP-hard, one can expect that the resulting model will be quickly challenging to solve.
294 We empirically confirmed this on very tiny instances: we tested it on problems with 20
295 variables and maximum domain size bounded by six, asking for just for four 15-diverse
296 solutions. This elegant approach took more than 23 hours to produce 4 solutions. We
297 therefore decided to solve a relaxed version of DIVERSESET : an iterative algorithm
298 provides a greedy approximation of the problem that preserves most of the required
299 guarantees. Using this approach, the problem of producing four 15-diverse solutions of
300 the above tiny problem takes just 0.28 seconds.

301 **Definition 6** (DIVERSESEQ). *Given a dissimilarity matrix D , an integer M and a dissimilarity*
302 *threshold δ , the set of assignments \mathbf{Z} of DIVERSESEQ(C, D, M, δ) is built recursively:*

- 303 • *The first solution $\mathbf{Z}[1]$ is the optimum of C*
- 304 • *When solutions $\mathbf{Z}[1..(i - 1)]$ are computed, $\mathbf{Z}[i]$ is such that:*
305 *for all $1 \leq j < i$, $\text{DIST}(\mathbf{Z}[i], \mathbf{Z}[j], D, \delta) = 0$ and $\mathbf{Z}[i]$ has minimum cost.*
306 *That is, $\mathbf{Z}[i]$ is the minimum cost solution, among assignments that are at distance at least*
307 *δ from all the previously computed solutions.*

308 The set of solutions DIVERSESEQ requires to optimally assign n variables M times,
309 instead of the $n \cdot M$ variables. Given the NP-hardness of the WCSP, solving DIVERSESEQ

310 may be exponentially faster than DIVERSESET while still providing guarantees that dis-
 311 tance constraints are satisfied together with a weakened form of optimality, conditional
 312 on the solutions found in previous iterations. The solution set is still guaranteed to
 313 contain the GMEC (the first solution produced).

314 5. Relation with existing work

315 In the case of Boolean functions ($\top = 1$), the work of [31] considers the optimization
 316 of the solution set cardinality M or the diversity threshold δ using average or minimum
 317 dissimilarity. The authors prove that enforcing Arc Consistency on a constraint requiring
 318 sufficient average dissimilarity \bar{d} is polynomial but NP-complete for minimum dissim-
 319 ilarity \bar{d} . They evaluate an algorithm for incremental production of a set maximizing
 320 \bar{d} . The papers [32] and [33] later addressed the same problems using global constraints
 321 and knowledge compilation techniques. Being purely Boolean, these approaches cannot
 322 directly capture cost (or energy) which is crucial for CPD. More recently, [34] proposed a
 323 Constraint Optimization Problem approach to provide diverse high-quality solutions.
 324 Their approach however trades diversity for quality while diversity is a requirement in
 325 our case.

326 The idea of producing diverse solutions has also been explored in the more closely
 327 related area of discrete stochastic Graphical Models (such as Markov Random Fields).
 328 In the paper of Batra et al. [35], the Lagrangian relaxation of minimum dissimilarity
 329 constraints is shown to add only unary cost functions to the model. This approach
 330 can be adapted to cost function networks, but a non zero duality gap remains and
 331 ultimately, no guarantee can be offered. This work was extended in [36] using higher-
 332 order functions to *approximately* optimize a trade-off between diversity and quality. More
 333 recently, [37] addressed the DIVERSESET problem, but using optimization techniques
 334 that either provide no guarantee or are restricted to tractable variants of the WCSP
 335 problem, defined by submodular functions [13].

336 In the end, we observe that none of these approaches simultaneously provides
 337 guarantees on quality and diversity. Closest to our target, [38] considered the problem
 338 of incrementally producing the set of the best M δ -modes of the joint distribution $J_N(X)$.

339 **Definition 7** ([39]). *A solution \mathbf{t} is said to be a δ -mode iff there exists no better solution than \mathbf{t}*
 340 *in the Hamming ball of radius δ centered in \mathbf{t} (implying that \mathbf{t} is a local minimum).*

341 In [38–41], an exact dynamic programming algorithm, combined with an A* heuristic
 342 search and tree-decomposition was proposed to exactly solve this problem with the
 343 Hamming distance. This algorithm relies however on NP-complete lower bounds and is
 344 restricted to a fixed variable order, a restriction that is known to often drastically hamper
 345 solving efficiency. It however provides a diversity guarantee: indeed, a δ -mode will
 346 always be *strictly* more than δ away from another one and will be produced by greedily
 347 solving DIVERSESEQ.

348 **Theorem 1.** *Given a cost function network \mathcal{C} , a diversity threshold δ , and $D = H$ the Ham-*
 349 *ming dissimilarity matrix, for any δ -mode \mathbf{t} , there exists a value M' such that the solution of*
 350 *DIVERSESEQ($\mathcal{C}, H, M', \delta + 1$) contains \mathbf{t} .*

351 **Proof.** If a δ -mode \mathbf{t} is not in the solution of DIVERSESEQ($\mathcal{C}, H, M', \delta + 1$), this must be
 352 because it gets forbidden by a DIST constraint. Consider the iteration i which forbids
 353 \mathbf{t} for the first time: a solution with a cost lower than the cost of \mathbf{t} was produced (else \mathbf{t}
 354 would have been produced instead) but this solution is strictly less than $\delta + 1$ away from
 355 \mathbf{t} (since \mathbf{t} gets forbidden). But this contradicts the fact that \mathbf{t} is a δ -mode. \square

356 For a sufficiently large M , the sequence \mathbf{Z} of DIVERSESEQ($N, H, M, \delta + 1$) solutions
 357 will therefore contain all δ -modes and possibly some extra solutions. Interestingly, it is
 358 not difficult to separate modes from non-modes.

359 **Theorem 2.**

- 360 1. Any assignment \mathbf{t} of a CFN $\mathcal{C} = (\mathbf{X}, \mathbf{D}, \mathbf{C})$ is a δ -mode iff it is an optimal solution of the
 361 CFN $(\mathbf{X}, \mathbf{D}, \mathbf{C} \cup \{\text{DIST}(\mathbf{X}, \mathbf{t}, H, -\delta)\})$
 362 2. For bounded δ , this problem is in P.

363 **Proof.** 1. The function $\text{DIST}(\mathbf{X}, \mathbf{t}, H, -\delta)$ restricts \mathbf{X} to be within δ of \mathbf{t} . If \mathbf{t} is an optimal
 364 solution of $(\mathbf{X}, \mathbf{C} \cup \{\text{DIST}(\mathbf{X}, \mathbf{t}, H, -\delta)\})$ then there is no better assignment than \mathbf{t} in the
 365 δ -radius Hamming ball and \mathbf{t} is a δ -mode.

366 2. For bounded δ , a CFN with n variables and at most d values in each domain, there is
 367 $O((nd)^\delta)$ tuples within the Hamming ball, because from \mathbf{t} , we can pick any variable (n
 368 choices) and change its value (d choices), δ times. Therefore the problem of checking if \mathbf{t}
 369 is optimal is in P. \square

370 **6. Representing the diversity constraint**

371 The key to guaranteeing quality and diversity of solutions in a cost function net-
 372 work is the dissimilarity cost function DIST . Given a complete assignment \mathbf{t} , a dissimi-
 373 larity D and a threshold δ , we need to concisely encode the global diversity constraint
 374 $\text{DIST}(\mathbf{X}, \mathbf{t}, D, \delta)$.

375 *6.1. Using automata*

376 Given a solution \mathbf{t} , a dissimilarity matrix D and a diversity threshold $\delta > 0$, the
 377 cost function $\text{DIST}(\cdot, \mathbf{t}, D, \delta)$ needs to be added to the cost function network. Note that
 378 the function may involve all the network variables: it is a *global cost function* and its
 379 representation as one huge, exponential size tensor is not possible.

380 To encode this function concisely, we exploit the fact that the set of authorized
 381 tuples defines a regular language, that can be encoded into a finite state automaton and
 382 then decomposed in ternary functions [42,43]. Here, we use a weighted automaton to
 383 define the weighted regular language of all tuples with their associated cost. A weighted
 384 automaton $\mathcal{A} = (\Sigma, \mathbf{Q}, \Delta, Q_0, \mathbf{F})$ encoding $\text{DIST}(\mathbf{X}, \mathbf{t}, D, \delta)$ can be defined as follows:

- 385 • The alphabet is the set of possible values, *i.e.*, the union of the variable domains
- 386 $\Sigma = \bigcup_{i=1}^n \mathbf{D}_i$
- The set of states \mathbf{Q} gathers $(\delta + 1) \cdot (n + 1)$ states denoted q_i^d :

$$\mathbf{Q} = \{q_i^d \mid 0 \leq i \leq n, 0 \leq d \leq \delta\}$$

387 that represent the fact that the first i values of \mathbf{X} have distance d to the first i values
 388 of \mathbf{t} . For $d = \delta$, automaton state q_i^δ represents the fact that the first i values of \mathbf{X}
 389 have distance $\geq \delta$ to the first i values of \mathbf{t} .

- In the initial state, no value of \mathbf{X} has been read, and the dissimilarity is 0:

$$Q_0 = q_0^0$$

- The assignment is accepted if it has dissimilarity from \mathbf{t} higher than the threshold δ ,
 hence the accepting state:

$$\mathbf{F} = \{q_n^\delta\}$$

- 390 • For every value r of X_i , the transition function $\Delta : \mathbf{Q} \times \Sigma \times \mathbf{Q}$ defines a 0-cost
 391 transition from q_i^d to $q_{i+1}^{\min(d+D(r, \mathbf{t}[i+1]), \delta)}$. All other transitions have infinite cost \top .

392 This weighted automaton contains $O(n \cdot (\delta + 1) \cdot d)$ finite cost transitions, were d is
 393 the maximum domain size. An assignment \mathbf{t}' of \mathbf{X} is accepted if and only if $d(\mathbf{t}', \mathbf{t}) \geq \delta$;
 394 and the automaton represents the DIST cost function. An example of a DIST encoding
 395 automaton is given in Figure 3.

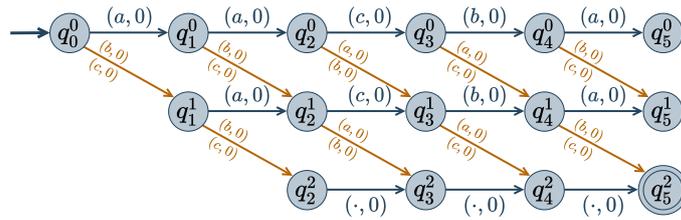


Figure 3. Weighted automaton representing $\text{DIST}(\mathbf{X}, \mathbf{t}, H, \delta)$ where \mathbf{X} is a set of 5 variables, with domains $\mathbf{D}_i = \{a, b, c\}$, $\mathbf{t} = aacba$, H represents the Hamming distance, and δ is set to 2. State q_i^d means that values $X_1 \dots X_i$ are such that $H(X_1 \dots X_i, t[X_1 \dots X_i]) = d$ (or $\geq \delta$ if $d = \delta$). A labeled arrow $q \xrightarrow{(v,w)} q'$ means $\Delta(q, v, q') = w$, i.e., there is a transition from q to q' with value v and weight w .

396 6.2. Exploiting automaton function decomposition

397 It is known that the WREGULAR cost function, encoding automaton \mathcal{A} , can be de-
 398 composed into a sequence of ternary cost functions [43]. The decomposition is achieved
 399 by adding $n + 1$ state variables Q_0, \dots, Q_n , and n ternary cost functions $w_{Q_i, X_{i+1}, Q_{i+1}}^A$
 400 such that $w_{Q_i, X_{i+1}, Q_{i+1}}^A(q_i, x_{i+1}, q_{i+1}) = c$ if and only if there exists a transition from q_i
 401 to q_{i+1} in \mathcal{A} labeled with value x_{i+1} and cost c . Variable Q_0 is restricted to the starting
 402 states and variable Q_n to the accepting states. Additional variables and ternary functions
 403 are represented in Figure 4. The resulting set of ternary functions is logically equivalent
 404 to the original global constraint.

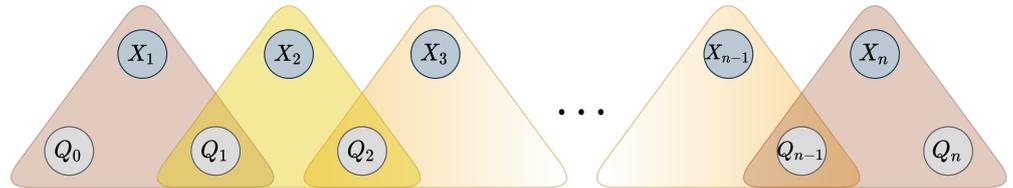


Figure 4. Hypergraph representation of the decomposition of a WREGULAR cost function with additional state variables Q_i and transition-encoding ternary functions.

405 The DIST function satisfies however several properties that can be exploited to
 406 further improve its space and time efficiency. One can first notice that the set of forbidden
 407 solutions does not depend on the order of the variables (the distance measure is the sum
 408 of independent variable-wise dissimilarities). Therefore, the order of the variables in the
 409 automaton can be chosen freely. We can use the DAC-ordering [43]. This order is known
 410 to preserve the strength of the lower bounds produced by CFN local consistencies when
 411 applied to the decomposition (instead of the initial full DIST function, with its potentially
 412 exponential size table).

413 Then, in the case of the DIST cost function, each state variable has $(\delta + 1) \cdot (n + 1)$
 414 values (the number of states in the automaton) and each ternary function cost table
 415 describes costs of $(\delta + 1)^2 \cdot (n + 1)^2 \cdot d$ tuples, where d is the domain size. To speed up
 416 the resolution through better soft local consistency enforcing, we exploit the properties
 417 of DIST and the dissimilarity matrix D .

418 6.3. Compressing the encoding

419 The encoding of a DIST cost function in a sequence of n ternary functions, described
 420 in cost tables of size $(\delta + 1)^2 \cdot (n + 1)^2 \cdot d$ can be reduced along several lines.

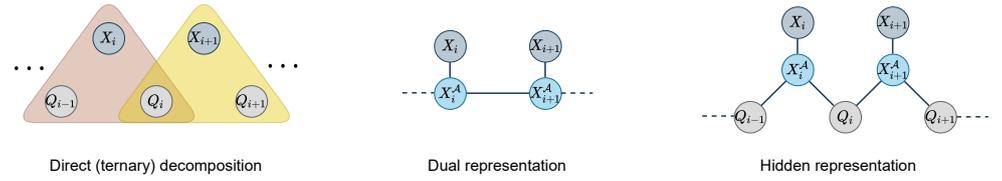


Figure 5. Representation of the ternary decomposition, and its dual and hidden representations.

First, for DIST, we know that states s_i^d can only be reached after i transitions, *i.e.*, the reachable states of variable Q_i are the states in the i -th column in the DIST automaton (see Figure 3). The domains of the variables Q_i can be reduced to the $\delta + 1$ states s_i^d :

$$\mathbf{D}_{Q_i} = \{q_i^d \mid 0 \leq d \leq \delta\}$$

421 Furthermore, our semi-metrics are defined by a non-decreasing sum of non-negative
 422 elements of D . Therefore, any state q_i^d can reach the accepting state q_n^δ if and only if the
 423 maximum dissimilarity md_i that can be achieved from variable i to variable n is larger
 424 than the remaining diversity to reach $\delta - d$. All such maximum dissimilarities md_i can
 425 be pre-computed in one pass over all variables in \mathbf{X} as follows:

- 426 • $md_n = 0$
- 427 • For $0 \leq i < n$, $md_i = md_{i+1} + \max_{v, v' \in \mathbf{D}_{i+1}} D(v, v')$

428 In the Hamming case, the distance can increase by 1 at most, *i.e.*, $\max_{v, v' \in \mathbf{D}_{i+1}} H(v, v') =$
 429 1, therefore $md_i = n - i$.

430 A symmetric argument holds for the starting state q_0^0 . These simplifications reduce
 431 ternary cost tables to $O((\delta + 1)^2 \cdot d)$.

432 For a given dissimilarity matrix D , let $\#D$ denote the *number of distinct values* that
 433 appear in D . If variables have domains of maximum size d and ignoring the useless 0
 434 matrix, we know that $2 \leq \#D \leq 1 + \frac{d \cdot (d+1)}{2}$. However, distance matrices are usually
 435 more structured. For example, the BLOSUM62 similarity matrix leads to $\#D = 12$ levels.

436 In the Hamming case, there are $\#H = 2$ dissimilarity levels. This means that a
 437 state q_i^d can only reach states q_{i+1}^d or q_{i+1}^{d+1} . This sparsity of the transition matrix can be
 438 exploited, provided it is made visible. This can be achieved using extended variants
 439 of the *dual* and *hidden* encoding of constraint networks [14,15]. These transformations,
 440 detailed hereafter, are known to preserve the set of solutions and their costs.

441 In constraint networks, the *dual* representation of a constraint network $\mathcal{X} =$
 442 $(\mathbf{X}, \mathbf{D}, \mathbf{C})$ is a new network $\mathcal{X}' = (\mathbf{X}', \mathbf{D}', \mathbf{C}')$ with:

- One variable X_S per constraint $c_S \in \mathbf{C}$:

$$\mathbf{X}' = \{X_S \mid c_S \in \mathbf{C}\}$$

- Domain \mathbf{D}_{X_S} of variable X_S is the set of tuples $\mathbf{t} \in \mathbf{D}_S$ that satisfy the constraint c_S :

$$\mathbf{D}' = \{\mathbf{D}_{X_S} \mid c_S \in \mathbf{C}\} \quad \mathbf{D}_{X_S} = \{\mathbf{t} \in c_S\}$$

- For each pair of constraints $c_S, c_{S'} \in \mathbf{C}$ with overlapping scopes $\mathbf{S} \cap \mathbf{S}' \neq \emptyset$, there is
 a constraint $c_{X_S, X_{S'}}$ that ensures that tuples assigned to X_S and $X_{S'}$ are compatible,
i.e., they have the same values on the overlapping variables:

$$\mathbf{C}' = \{c_{X_S, X_{S'}} \mid X_S, X_{S'} \in \mathbf{X}', \mathbf{S} \cap \mathbf{S}' \neq \emptyset\}$$

where

$$c_{X_S, X_{S'}} = \{(\mathbf{t}, \mathbf{t}') \in \mathbf{D}_{X_S} \times \mathbf{D}_{X_{S'}} \mid \mathbf{t}[\mathbf{S} \cap \mathbf{S}'] = \mathbf{t}'[\mathbf{S} \cap \mathbf{S}']\}$$

We apply this transformation to the reduced $w_{Q_i, X_{i+1}, Q_{i+1}}^A$ functions (see Figure 5). The dual variable of $w_{Q_i, X_{i+1}, Q_{i+1}}^A$ is a variable X_i^A that contains all pairs $(q, q') \in Q_i \times Q_{i+1}$ such that there is a transition from q to q' in \mathcal{A} . For the Hamming case, the variable X_i^A has at most $2\delta + 1$ values. It is connected to X_i by a pairwise function:

$$c_{X_i^A, X_i} : \begin{array}{l} \mathbf{D}_{X_i^A \times \mathbf{D}_i} \longrightarrow \{0, \dots, \top\} \\ ((q, q'), v) \longmapsto \Delta(q, v, q') \end{array}$$

443 where Δ is the weighted transition function of the automaton \mathcal{A} .

In this new dual representation, for every pair of consecutive dual variables X_{i-1}^A and X_i^A , we add a function on these two variables to ensure that the arriving state of X_{i-1}^A is the starting state of X_i^A :

$$c_{X_{i-1}^A, X_i^A} : ((q_{i-1}, q'_{i-1}), (q_i, q'_i)) \mapsto \begin{cases} 0 & \text{if } q'_{i-1} = q_i \\ \top & \text{otherwise.} \end{cases}$$

444 In the worst case, this function has size $O(\#D^2 \cdot \delta^2)$ ($O(\delta^2)$ in the Hamming case). Only
445 n extra variables are required.

446 The *hidden* representation of a constraint network $\mathcal{X} = (\mathbf{X}, \mathbf{D}, \mathbf{C})$ is a network
447 $\mathcal{X}'' = (\mathbf{X}'', \mathbf{D}'', \mathbf{C}'')$ with:

- All the variables in \mathbf{X} and the variables X_S from the dual network (and associated domains):

$$\mathbf{X}'' = \mathbf{X} \cup \mathbf{X}'$$

- For any dual variable X_S , and each $X_i \in \mathbf{S}$, the set of constraints \mathbf{C}'' contains a function involving X_i and X_S :

$$c_{X_i, X_S} : (v, \mathbf{t}) \in \mathbf{D}_i \times \mathbf{D}_{X_S} \mapsto \begin{cases} 0 & \text{if } \mathbf{t}[X_i] = v \\ \top & \text{otherwise.} \end{cases}$$

As before, this transformation is applied to the reduced $w_{Q_i, X_{i+1}, Q_{i+1}}^A$ functions only (see Figure 5). In this new hidden representation, we keep variables Q_i and create two pairwise functions involving each Q_i and respectively X_i^A and X_{i+1}^A :

$$c_{Q_i, X_{i-1}^A} : (q'', (q, q')) \mapsto \begin{cases} 0 & \text{if } q'' = q \\ \top & \text{otherwise.} \end{cases}$$

$$c_{Q_i, X_i^A} : (q'', (q, q')) \mapsto \begin{cases} 0 & \text{if } q'' = q' \\ \top & \text{otherwise.} \end{cases}$$

448 These functions ensure that the state value of Q_i is consistent with the arriving state of
449 the transition represented in X_{i-1}^A and the starting state of X_i^A . In the worst case, these
450 functions have size $O(\#D \cdot \delta^2)$ ($O(\delta^2)$ in the Hamming case).

451 The dedicated dual and hidden representations require the description of $O(\delta \cdot d +$
452 $\#D^2 \cdot \delta^2)$ and $O(\delta \cdot d + \#D \cdot \delta^2)$ tuples respectively (it is $O(\delta \cdot d + \delta^2)$ in the Hamming
453 case), instead of the $O(d \cdot \delta^2)$ tuples in $w_{Q_i, X_{i+1}, Q_{i+1}}^A$.

454 7. Greedy DiverseSeq

455 The task at hand is the resolution of $\text{DIVERSESET}(\mathcal{C}, D, M, \delta)$, *i.e.*, the generation
456 of a set of M solutions with minimum cost, that satisfy a minimum pairwise diversity
457 constraint. The exact computation being too expensive, we are tackling a greedy com-
458 putation of $\text{DIVERSESEQ}(\mathcal{C}, D, M, \delta)$, a set of diverse good solutions that approximates
459 DIVERSESET . The DIVERSESEQ computation is iterative:

- 460 1. The CFN \mathcal{C} is solved using branch-and-bound while maintaining soft local consis-
461 tencies [26].

- 462 2. If a solution \mathbf{t} is found, it is added to the ongoing solution sequence \mathbf{Z} .
 463 3. If M solutions have been produced, the algorithm stops.
 464 4. Otherwise, the cost function $\text{DIST}(\mathbf{X}, \mathbf{t}, D, \delta)$ is added to the previously solved
 465 problem.
 466 5. We loop and solve the problem again (Step 1)
 467 At Step 2, if no solution exists, the sequence of solutions \mathbf{Z} can provably not be extended
 468 to length M and the problem has no solution (but a shorter sequence has been produced).

Algorithm 1: Incremental production of $\text{DIVERSESEQ}(\mathcal{C}, D, M, \delta)$

```

1 Procedure Solve ( $\mathcal{C}, lb, ub$ )
2   Compute optimum solution  $\mathbf{t}^*$  of  $\mathcal{C}$  with  $lb \leq \text{Cost}(\mathbf{t}^*) < ub$ ;
3   if  $\mathbf{t}^*$  exists then
4     | return ( $\mathbf{t}^*, true$ )
5   else
6     | return ( $\emptyset, false$ );

7 Procedure IncrementalSearch ( $\mathcal{C}, lb, ub, \Delta^h, \mathbf{Z}, M, D, \delta$ )
8    $\mathbf{t}^*, solved \leftarrow \text{Solve}(\mathcal{C}, lb, \text{Cost}(\mathbf{Z}[i-1] + \Delta^h))$ ;
9   if not solved then
10    |  $\mathbf{t}^*, solved \leftarrow \text{Solve}(\mathcal{C}, lb, ub)$ ; ▷ Upper bound prediction failed
11  if not solved then
12    | return  $\mathbf{Z}$ ; ▷  $\mathbf{Z}$  can not be extended to length  $M$ 
13   $\mathbf{Z} \leftarrow \mathbf{Z} \cup \{\mathbf{t}^*\}$ ;
14  if  $|\mathbf{Z}| = M$  then
15    | return  $\mathbf{Z}$ ; ▷ Enough solutions have been produced
16  Add  $\text{DIST}(\mathbf{X}, \mathbf{t}^*, D, \delta)$  to  $\mathcal{C}$ ;
17  Propagate and save local consistencies in  $\mathcal{C}$ ;
18  Update  $\Delta^h$ ; ▷ Using Cost ( $\mathbf{t}^*$ )
19   $lb \leftarrow \text{Cost}(\mathbf{t}^*)$ ;
20   $\text{DIVERSESEQ}(\mathcal{C}, lb, ub, \Delta^h, \mathbf{Z}, M, D, \delta)$ ;

21 Procedure DiverseSeq ( $\mathcal{C}, M, D, \delta$ )
22  | return  $\text{IncrementalSearch}(\mathcal{C}, 0, \top, 0, \emptyset, M, D, \delta)$ ;

```

469 This basic schema has been improved in three different ways (see Algorithm 1):

470 **Incrementality** Since the problems solved are increasingly constrained, all the equivalence preserving transformations and pruning that have been applied to enforce local consistencies at iteration $i-1$ are still valid in the following iterations. Instead of restarting from a problem $\mathcal{C} = (\mathbf{X}, \mathbf{D}, \mathcal{C} \cup \bigcup_{1 \leq j < i} \{\text{DIST}(\mathbf{X}, \mathbf{Z}[j], D, \delta)\})$, we reuse the problem solved at iteration $i-1$ after it has been made locally consistent, add the $\text{DIST}(\mathbf{X}, \mathbf{Z}[i-1], D, \delta)$ constraint and reinforce local consistencies. Similarly to incremental SAT solvers, adaptive variable ordering heuristics that have been trained at iteration $i-1$ are reused at iteration i .

478 **Lower bound** Since the problems solved are increasingly constrained, we know that the optimal cost oc^i obtained at iteration i cannot have a lower cost than the optimum cost oc^{i-1} reported at iteration $i-1$. When large plateaus are present in the energy landscape, this allows stopping the search as soon as a solution of cost oc^{i-1} is reached, avoiding a useless repeated proof of optimality.

483 **Upper bound prediction** Even if there are no plateaus in the energy landscape, there may be large regions with similar variations in energy. In this case, the difference in energy between oc^{i-1} and oc^i will remain similar for several iterations. Let $\Delta_i^h = \max_{\max(2, i-h) \leq j < i} (oc^j - oc^{j-1})$ be the maximum variation observed in the last h iterations (we used $h = 5$). At iteration i , we can first solve the problem

PDB ID	n	d	PDB ID	n	d	PDB ID	n	d
1aho	56	378	3i8z	50	354	1ten	81	392
2fjz	53	324	2cg7	82	380	1ucs	60	342
1b9w	78	386	3rdy	65	396	2bwf	69	347
2gkt	45	357	2erw	47	446	2evb	68	323
1f94	53	386	3vdj	67	391	2o37	60	386
2pne	77	401	2fht	64	346	2o9s	48	327
1hyp	66	385	1bxy	52	384	3f04	87	356
2pst	61	357	1ctf	68	349	3fym	70	348
1uln	66	367	1czp	76	373	3gqs	67	344
1uoy	56	337	1fqt	85	377	3gva	87	348
2ca7	44	348	1guu	47	350	3i2z	67	360
1yzm	46	294	1t8k	68	361			

Table 1. List of protein structures used in our benchmark set, for full redesign: pdb identifier, domain length n (number of variables in the resulting CFN) and maximum domain size d .

488 with a temporary upper bound $k' = \min(k, oc_{i-1} + 2 \cdot \Delta_i^h)$ that should preserve a
 489 solution. If $k' < k$, this will lead to increased determinism, additional pruning, and
 490 possibly exponential savings. Otherwise, if no solution is found, the problem is
 491 solved again with the original upper bound k . We call this *predictive bounding*.

492 Each of these three improvements has the capacity to offer exponential time savings,
 493 and all are used in the experiments presented in the next sections.

494 8. Results

495 We implemented the iterative approach described above in its direct (ternary)
 496 decomposition, hidden and dual representations in the CFN open-source solver `toulbar2`
 497 and experimented with it on various CFNs representing real Bayesian Networks [44].
 498 All three decompositions offered comparable efficiency but empirically, as expected,
 499 the dual encoding was almost systematically more efficient. It is now the default for
 500 diversity encoding in `toulbar2`. All `toulbar2` preprocessing algorithms dedicated to exact
 501 optimization that do not preserve suboptimal solutions were deactivated at the root node
 502 (variable elimination, dead-end elimination, variable merging). We chose to enforce
 503 strong virtual arc consistency (flag `-A` in `toulbar2`). The computational cost of VAC,
 504 although polynomial, is high, but amortized over the M resolutions. During tree search,
 505 the default existential directional arc consistency (EDAC) was used. All experiments
 506 were performed on one core of a Xeon Gold 6140 CPU at 2.30 GHz. Wall-clock times
 507 could be further reduced using a parallel implementation of the underlying Hybrid
 508 Best-First search engine [45], currently under development in `toulbar2`.

509 Following our main motivation for protein design, we extracted two sets of prepared
 510 protein backbones for full redesigns from the benchmark sets built by [46] and [47] with
 511 the aim of checking if, as expected, diverse libraries can improve the overall design
 512 process. In the benchmark of monomeric proteins of less than 100 residues, with an X-ray
 513 resolved structure below 2 Å, with no missing or nonstandard residues and no ligand
 514 from [46], we selected the 20 proteins that had required the least CPU-time to solve, as
 515 indicated in the Excel sheet provided in the supplementary information of paper [46].
 516 The harder instances from [47] correspond to proteins with diverse secondary structure
 517 compositions and fold classes. We selected the 17 instances that required less than 24
 518 hours of CPU-time for the full redesigned GMEC to be computed by `toulbar2`. These
 519 instances are listed in Table 1.

520 Full redesign was performed on each protein structure, and CFN instances were
 521 generated using the Dunbrack library [16] and Rosetta `ref2015` score function [48].
 522 Alternate rotamer libraries and score functions can be used if required as the algorithms
 523 presented here are not specialized for Rosetta (and not even for CPD, see [44]). The

524 resulting networks have from 44 to 87 rotamer variables, and maximum domain sizes
525 range from 294 to 446 rotamers. The number of variables is doubled after sequence
526 variables are added.

527 Predictive bounding contribution

528 $M = 10$ solutions with diversity threshold $\delta = 10$ for each problem from [46] were
529 generated, with and without predictive bounding. The worst CPU-time spent on the
530 resolution without predictive bounding was 32 minutes. It was reduced to 17 minutes
531 with predictive bounding. The average computation time was 201s per problem. This
532 shows that predictive bounding provides a simple and efficient boost and that real
533 CPD instances can be solved in a reasonable time, even when relatively large diversity
534 requirements are used.

535 Diversity improves prediction quality

536 For all instances, sets of $M = 10$ solutions were generated with diversity threshold
537 δ ranging from 1 to 15. For $\delta = 1$, the set of solutions produced is just the set of the 10
538 best (minimum energy) sequences.

539 These CPD problems use real protein backbones, determined experimentally. A
540 native sequence exists for these backbones, therefore it is possible to measure the im-
541 provements diversity brings in terms of recovering native sequences, known to be folded
542 and functional. Two measures are often used to assess computational protein design
543 methods. The Native Sequence Recovery (NSR) is the percentage of amino acid residues
544 in the designed protein which are identical to the amino acid residues in the native
545 sequence that folds on the target backbone. The NSR can be enhanced by taking into
546 account *similarity scores* between the amino acid types. Such scores are provided by
547 similarity matrices, like BLOSUM62 [30]. The *Native Sequence Similarity Recovery* (NSSR)
548 is the fraction of positions where the designed and native sequences have a positive
549 similarity score. NSR and NSSR measure how much the redesigned protein resembles
550 the natural in terms of sequence. While often used, these measures have their own
551 limitations: while protein design targets maximal stability, natural protein only require
552 sufficient marginal stability. In the end, they therefore provide a useful but imperfect
553 proxy for computation protein design evaluation: a perfect (100%) recovery would not
554 necessarily indicate the best algorithm (also because the approximate energy function
555 plays a major role here).

556 If solution diversity helps, the maximum NSR/NSSR over the 10 sequences should
557 improve when δ is large compared to when $\delta = 1$, as long as the costs remain close to the
558 optimum. A solution cost too far from the optimum, which could be generated because
559 of a diversity threshold set too high, would mean a poor quality of the solution. Even
560 with $\delta = 15$, the maximum difference in energy we observed with the global minimum
561 energy never exceeded 4.3 kcal/mol (with an average of 2.1 kcal/mol).

562 For each protein, and each diversity threshold $\delta = 2 \dots 15$, we compared the best
563 NSR (resp. NSSR) that could be obtained with δ to the best obtained with diversity
564 threshold 1, *i.e.*, the simple enumeration of the 10 best sequences. Results are plotted
565 in Figure 6 (resp. Figure 7). While somewhat noisy, they show a clear general increase
566 in the best NSR (resp. NSSR) when diverse sequences are produced, even with small δ .
567 To validate the statistical significance of the diversity improvement of sequence quality,
568 p -values for a unilateral Wilcoxon signed-rank test comparing the sample of best NSR
569 (resp. NSSR) for each $\delta = 2 \dots 15$ with $\delta = 1$ were computed. They are shown in Table 2
570 and confirm the improvement brought by increasingly diverse libraries.

571 The improvements are more clearly visible when one compares the absolute im-
572 provement in NSR (and NSSR) obtained when one compares the best sequences pro-
573 duced inside a library using a guaranteed diversity of 15 versus just 1, as illustrated in
574 Figure 8. On most backbones (X-axis), the increased diversity yields a clear improvement
575 in the NSR (Y-axis), with the largest absolute improvement exceeding 15% in NSR. For a

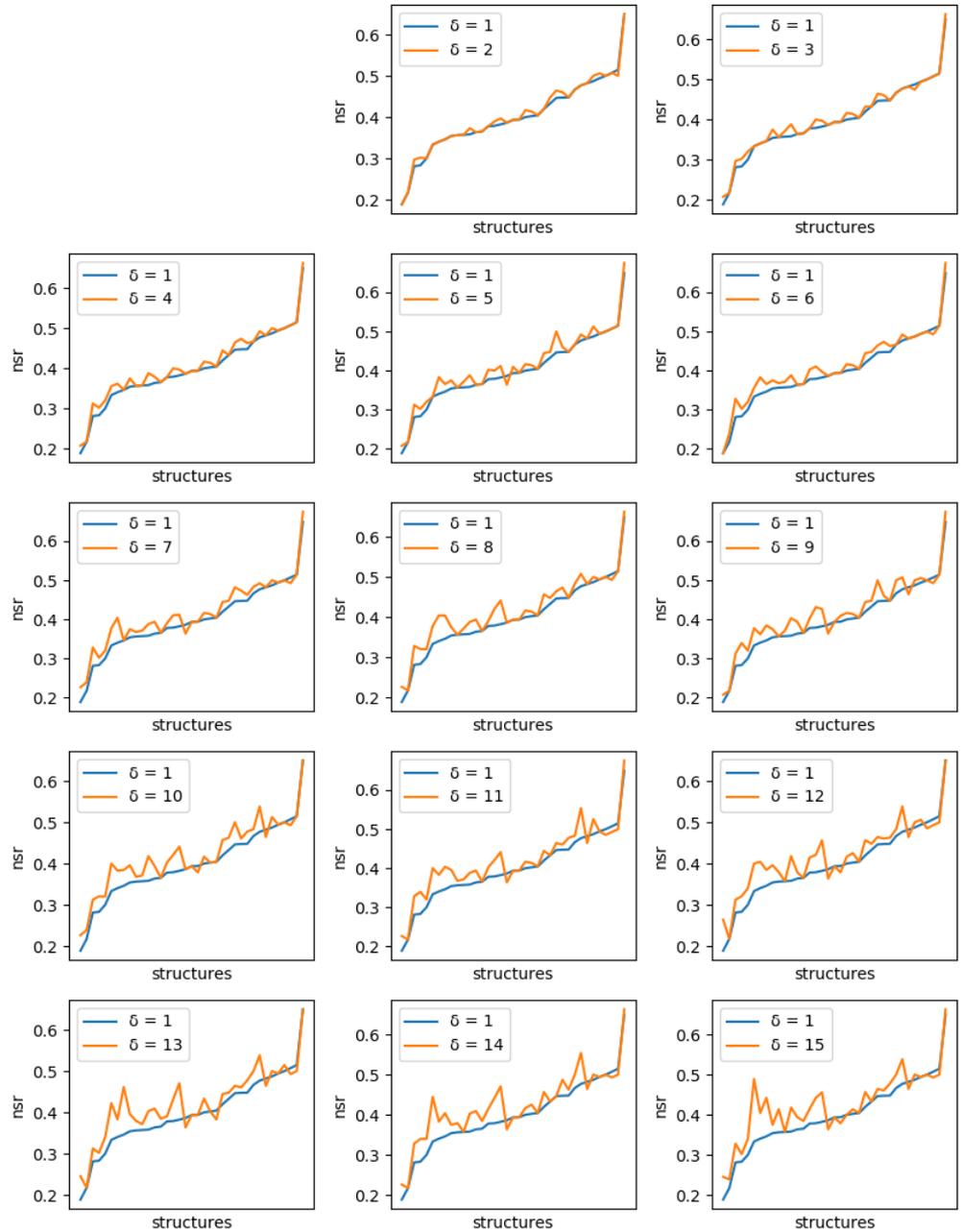


Figure 6. Comparison of the best NSR value obtained with ten 1-diverse sequences ($\delta = 1$, blue curve) with the best NSR value obtained with libraries of ten sequences of increased diversity. Each plot corresponds to a specific additional value of δ ($\delta = 2$ to 15, golden curve). Plots are ordered lexicographically from top-left to bottom right, with increasing values of diversity (δ). In each plot, the X-axis ranges over all tested backbones, sorted in increasing order of NSR value for the 1-diverse case and the Y-axis gives the corresponding NSR value. As the diversity requirement increases, the NSR value indicated by the golden curve increases also visibly.

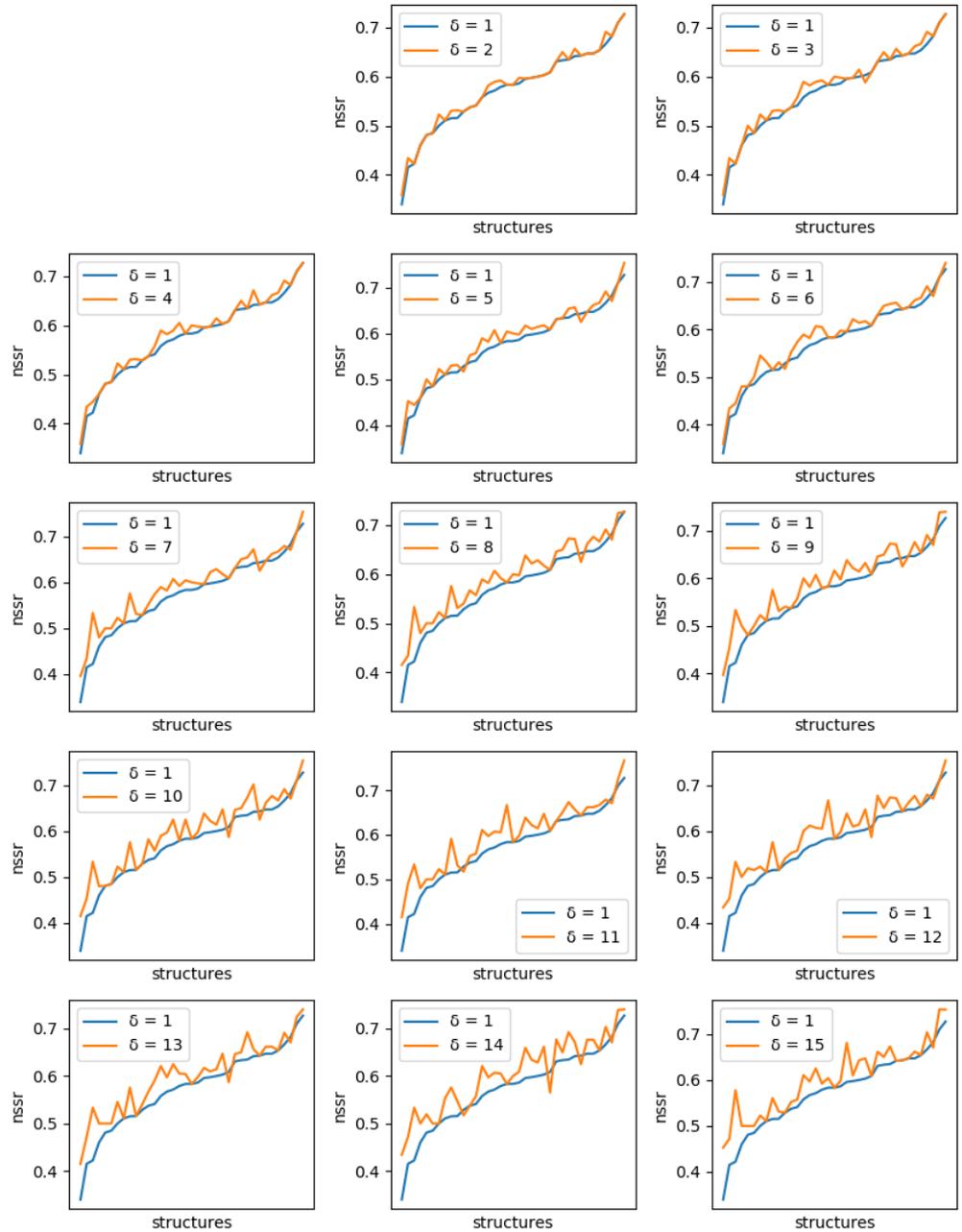


Figure 7. Comparison of the best NSSR value obtained with ten 1-diverse sequences ($\delta = 1$, blue curve) with the best NSSR value obtained with libraries of ten sequences of increased diversity. Each plot corresponds to a specific additional value of δ ($\delta = 2$ to 15, golden curve). Plots are ordered lexicographically from top-left to bottom right, with increasing values of diversity (δ). In each plot, the X-axis ranges over all tested backbones, sorted in increasing order of NSSR value for the 1-diverse case and the Y-axis gives the corresponding NSSR value. As the diversity requirement δ increases, the NSSR value indicated by the golden curve increases also visibly.

δ	Exact resolution		Subopt. resolution	
	NSR	NSSR	NSR	NSSR
2	2.88e-03	1.10e-03	4.11e-01	1.35e-01
3	3.87e-04	1.01e-04	1.14e-01	2.60e-02
4	4.42e-05	6.58e-05	4.48e-03	1.06e-03
5	8.11e-05	1.54e-05	1.98e-03	2.15e-03
6	1.51e-05	4.39e-06	7.47e-05	4.49e-05
7	1.88e-05	4.23e-06	8.86e-06	3.50e-05
8	1.27e-05	1.49e-06	8.19e-06	1.77e-05
9	2.76e-05	5.97e-06	2.07e-05	2.80e-06
10	1.14e-05	1.18e-05	1.78e-05	3.06e-05
11	4.27e-05	5.81e-07	2.32e-05	1.73e-05
12	6.63e-05	2.26e-06	1.75e-05	1.18e-05
13	4.43e-05	2.52e-06	2.48e-06	6.15e-06
14	2.29e-05	5.76e-06	5.26e-06	3.89e-07
15	3.92e-05	1.58e-06	2.68e-05	4.86e-05

Table 2. p -values for a unilateral Wilcoxon signed rank test comparing the sample of best NSR (resp. NSSR) for each $\delta = 2 \dots 15$ with $\delta = 1$, for optimal and suboptimal (3 kcal/mol allowed energy gap to real optimum) resolution.

576 small fraction of backbones, there is absolutely no improvement. These are backbones for
 577 which the GMEC actually provides the best NSR in both the 1-diverse and the 15-diverse
 578 cases. Then an even smaller fraction of backbones show an actual decrease in NSR: one
 579 close to GMEC solution did better than any of the 15-diverse sequences. The degradation
 580 here is very limited and likely corresponds to substitutions to similar amino acids. This
 581 is confirmed by the NSSR curve that takes into account amino acid properties through
 582 the underlying BLOSUM62 matrix used. Here, only one case show a degradation in
 583 NSSR.

584 Diversity also has the advantage that it is more likely to provide improved se-
 585 quences when the predicted 1-diverse sequences are poor. Indeed, with a initial sequence
 586 with NSR equal to r , the introduction of a random mutation will move us away from the
 587 native in $r\%$ of cases (we mutate a correct amino acid) and otherwise (we mutate a wrong
 588 amino acid) leave us with a wrong amino acid again (in 18/19 cases, leaving the NSR
 589 unchanged) or get us closer to the native sequence with 1/19 probability. On average,
 590 a mutation should therefore decrease the number of correct positions with probability
 591 $(r - \frac{1-r}{19})$, which increases with r and is positive as soon as the sequence has NSR higher
 592 than 5% ($\frac{1}{20}$). Our results confirm this trend, as shown in the NSR figure on the left of
 593 Figure 8: among the ten backbones with the highest improvement in NSR, nine had a
 594 1-diverse NSR below the average 1-diverse NSR. Conversely, only 50% of the ten less
 595 improved backbones had a below-average 1-diverse NSR. While these improvements
 596 underline the approximate nature of the energy function, showing that it is often worth
 597 to explore the sequence space beyond the GMEC, they also confirm that energy, on
 598 average, does guide us towards native sequences: instead of degrading NSR as soon as
 599 $r > \frac{1}{20}$, energy optimization pushes the introduced mutations to improve NSR in most
 600 cases, even with 15 introduced mutations and initial NSRs ranging from 20 to 60%.

601 Each dissimilarity constraint adds n extra variables to the network (with the dual
 602 representation). These variable domain sizes increase with the diversity threshold δ and
 603 contribute to the construction of increasingly large CFN instances that need to be solved.
 604 Computation times are plotted in Figure 9. As expected, a threshold increase leads to an
 605 exponential increase in the computation time. The small set of points on the top right
 606 corner of the plot correspond to the protein structure 3FYM. This protein, with 70 amino
 607 acids, is not the largest in our benchmark set, a reminder that size is not necessarily the

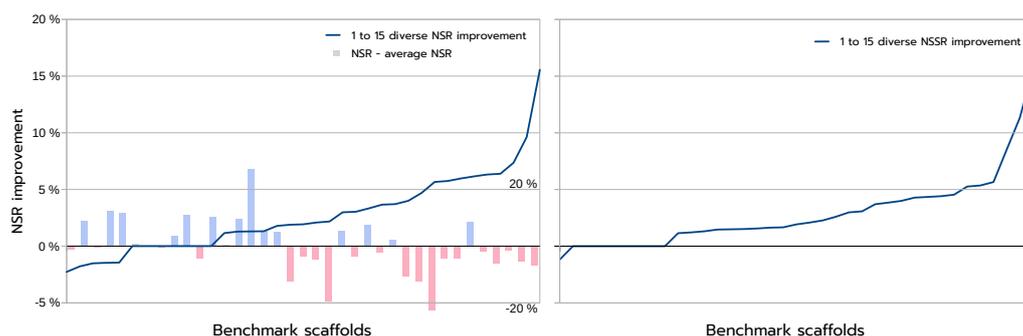


Figure 8. The blue curves above give the absolute change in NSR (Y axis, left figure) and NSSR (Y-axis, right figure) between the best 15-diverse and the best 1-diverse sequences found for each backbone. Backbones (on the X-axis) are ordered in increasing order of the corresponding measure. In the left figure, the bar-plot shows the difference between each backbone 1-diverse NSR and average 1-diverse NSR over all backbones. The corresponding NSR change scale appears on the the right with $\pm 20\%$ labels. Red bars indicate a below average 1-diverse NSR while blue bars indicate an above average 1-diverse NSR. The most improved NSRs, on the right of the left figure, mostly appear on weak (red, below average) 1-diverse NSRs.

608 best predictor of empirical computational cost in NP-hard problem solving. On this
 609 backbone, for high diversity thresholds, the 24h computation time limit was reached
 610 and less than 10 sequences were produced.

611 Given that the optimized function is an approximation of the real intractable energy,
 612 solving the problem to optimality might seem exaggerated. The requirement for opti-
 613 mality that we have used in previous experiments can be trivially relaxed to a relative or
 614 absolute approximation guarantee using artificially tightened pruning rules as originally
 615 proposed in [49] in the Weighted-A* algorithm. This pruning mechanism is already
 616 implemented in the toulbar2 solver (using the `-rgap` and `-agap` flags respectively).

617 For diversity threshold $\delta = 2 \dots 15$, we generated sets of 10 suboptimal diverse
 618 sequences that satisfy the diversity constraints, but allowing for a 3 kcal/mol energy
 619 gap to the optimum. Our optimizations are still provable, but the optimality guarantee
 620 is now reduced to a bounded sub-optimality guarantee. Empirically, the maximum
 621 energy degradation we observed with the global minimum energy over the 10 diverse
 622 sequences produced never exceeded 5.65 kcal/mol (with an average energy difference
 623 of 3.86 kcal/mol). This is only slightly more than the 4.3 kcal/mol worse degradation
 624 (average 2.1 kcal/mol) of the resolution, when exact optimum are used.

625 We compared these samples of suboptimal sequences to the set of 10 exact best
 626 sequences. Results for NSR and NSSR are shown in Figures 10 and 11 respectively, and
 627 corresponding p -values are displayed in Table 2 (unilateral Wilcoxon signed rank test).
 628 With dissimilarity threshold $\delta \geq 6$, it is clear that the set of diverse suboptimal sequences
 629 have better quality than the 10 best enumerated sequences. Moreover, as shown in
 630 Figure 12, for harder instances, suboptimal diverse resolution becomes faster than exact
 631 enumeration.

632 So, when predicting a library of sequences, if the instance is hard, it seems empiri-
 633 cally wise to generate suboptimal diverse sequences, instead of enumerating minimum
 634 energy sequences, without diversity. Doing so, there is a higher chance of predicting
 635 better sequences “in practice” faster.

636 9. Conclusions

637 Producing a library of diverse solutions is a very usual requirement when an
 638 approximate or learned model is used for optimal decoding. In this paper, we show that
 639 with an incremental provable CFN approach that directly tackles a series of decision NP-
 640 complete problems, using diversity constraints represented as weighted automata that
 641 are densely encoded in a dedicated dual encoding together with predictive bounding,

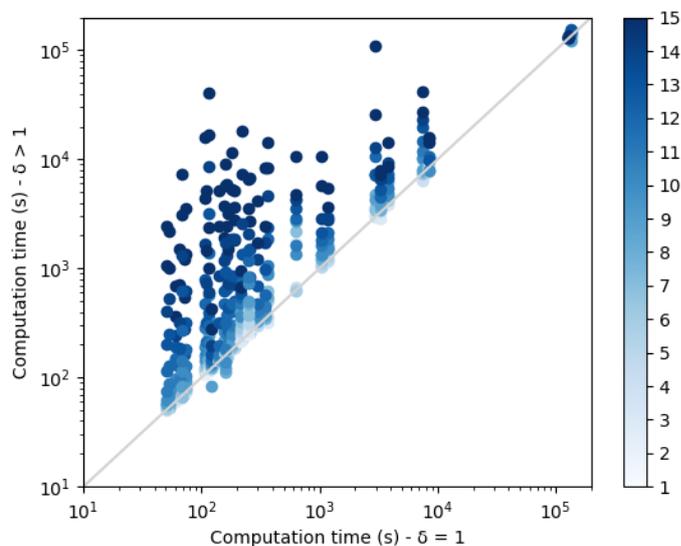


Figure 9. Comparison of the computation times of sequence sets without diversity $\delta = 1$, with sequence sets with diversity $\delta > 1$. The color scale on the right indicates the corresponding value of δ .

642 it is possible to produce sequences of solutions that satisfy guarantees on diversity
 643 on realistic full redesign Computational Protein Design instances. This guarantee is
 644 obtained on dense problems, with non-permuted-submodular functions while also
 645 guaranteeing that each new solution produced is the best given the previously identified
 646 solutions.

647 We also showed that the stream of diverse solutions that our algorithm produces
 648 can be filtered and each solution efficiently identified as being a δ -mode or not. δ -mode
 649 represent local minima, each defining its own basin in the protein sequence energy
 650 landscape. Beyond their direct use for design, the guaranteed diversity provided by our
 651 algorithm could also be of interest to perform more systematic analyses of such energy
 652 landscapes.

653 On real protein design problems, we observe that small and large diversity require-
 654 ments do improve the quality of sequence libraries when native proteins are fully re-
 655 designed. Moreover, large diversity requirements on suboptimal sequences also improve
 656 the quality of sequence libraries, compared to a simple enumeration of the minimum
 657 energy sequences. In the context of optimizing an approximate or learned function,
 658 the requirement for an optimal cost solution may be considered as exaggerated. Our
 659 guaranteed suboptimal resolution is useful, given that even computationally expensive
 660 approaches with asymptotic convergence results such as Simulated Annealing may fail
 661 with unbounded error [46].

662 Two directions could extend this work. Beyond the language of permuted sub-
 663 modular functions, the other important tractable class of CFN is the class of CFN with
 664 a graph of bounded tree-width. This parameter is exploited in several protein pack-
 665 ing [50] and design [51] algorithms and is also exploited in dedicated branch and bound
 666 algorithms, also implemented in *toulbar2* [45,52]. These tree-search algorithms are
 667 able to trade space for time and are empirically usable on problems with a tree-width
 668 that is too large to make pure dynamic programming applicable, mostly because of
 669 its space-complexity (in $O(d^w)$ where d is the domain size and w is the width of the
 670 tree-decomposition used). On such problems, it would be desirable to show that the
 671 decomposed ternary or binary functions we use for encoding DIST can be arranged in
 672 such a way that tree-width can be preserved or more likely, not exaggeratedly increased.

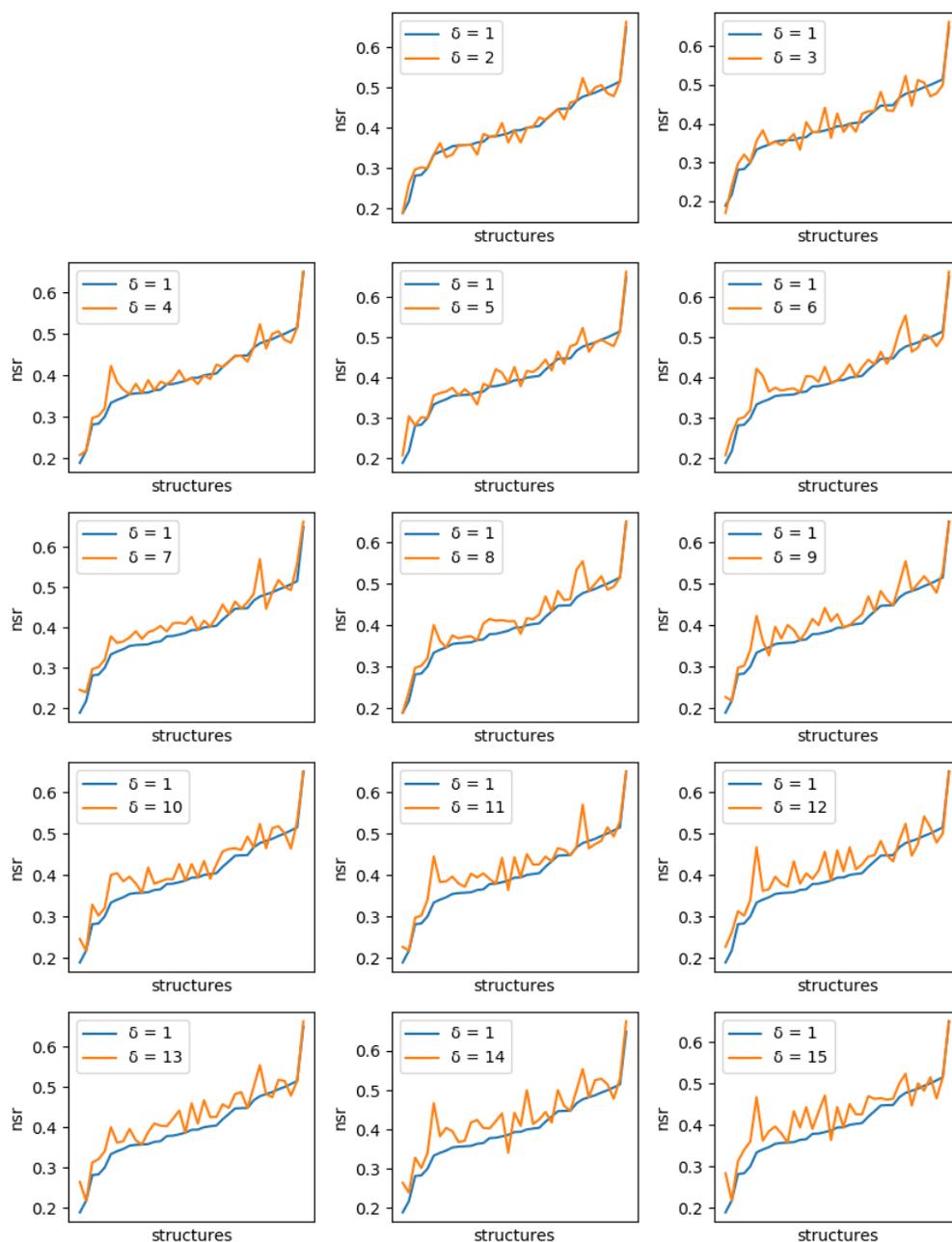


Figure 10. Comparison of the best NSR value obtained with ten 1-diverse sequences ($\delta = 1$, blue curve) with the best NSR value obtained with libraries of ten sequences of increased diversity all predicted with an allowed gap top optimal energy of 3 kcal/mol. Each plot corresponds to a specific additional value of δ ($\delta = 2$ to 15, golden curve). Plots are lexicographically ordered from top-left to bottom right, with increasing values of diversity (δ). In each plot, the X-axis ranges over all tested backbones, sorted in increasing order of NSR value for the 1-diverse case and the Y-axis gives the corresponding NSR value. As the diversity requirement δ increases, the NSR value indicated by the golden curve increases also visibly.

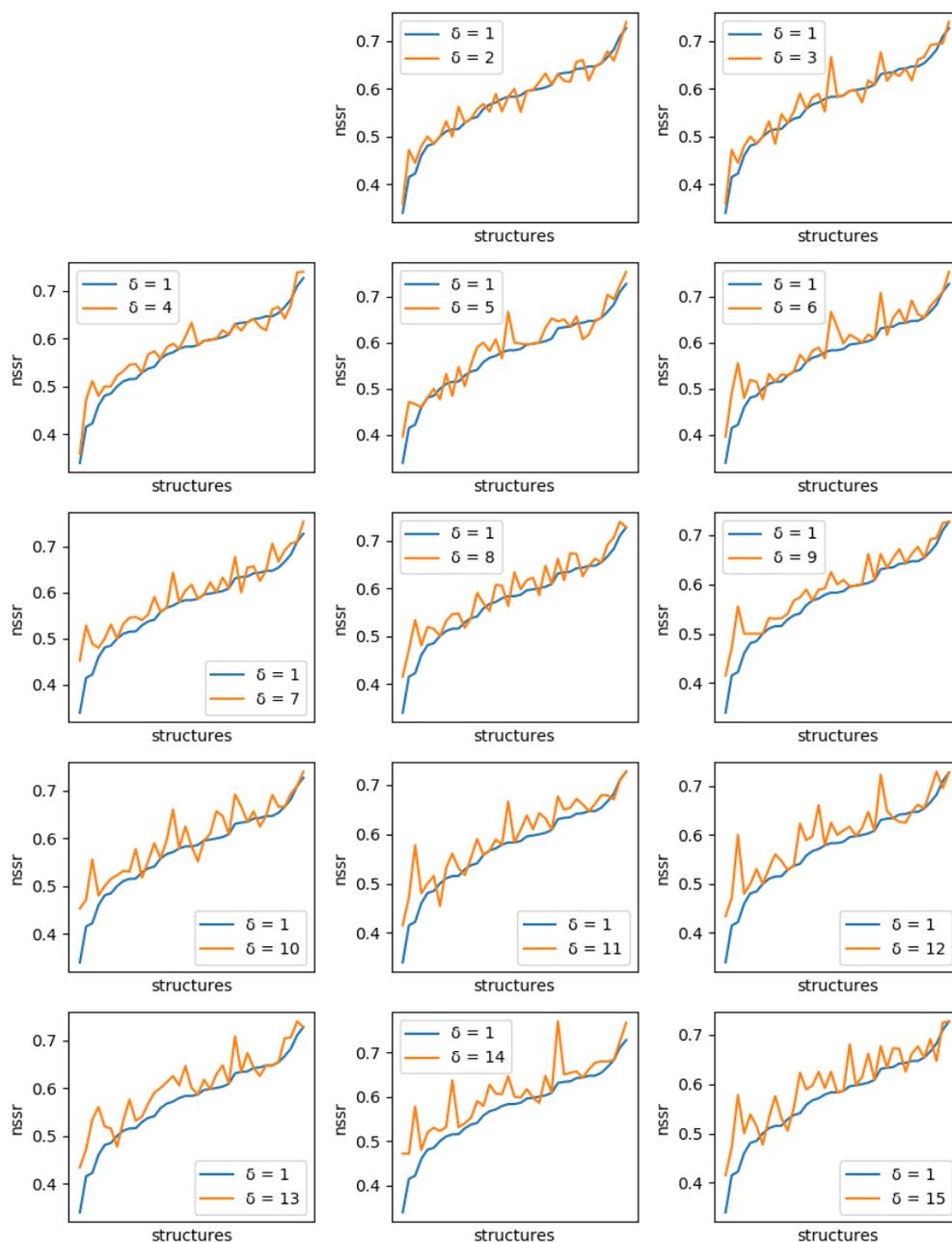


Figure 11. Comparison of the best NSSR value obtained with ten 1-diverse sequences ($\delta = 1$, blue curve) with the best NSSR value obtained with libraries of ten sequences of increased diversity all predicted with an allowed gap top optimal energy of 3 kcal/mol. Each plot corresponds to a specific additional value of δ ($\delta = 2$ to 15, golden curve). Plots are lexicographically ordered from top-left to bottom right, with increasing values of diversity (δ). In each plot, the X-axis ranges over all tested backbones, sorted in increasing order of NSSR value for the 1-diverse case and the Y-axis gives the corresponding NSSR value. As the diversity requirement δ increases, the NSSR value indicated by the golden curve increases also visibly.

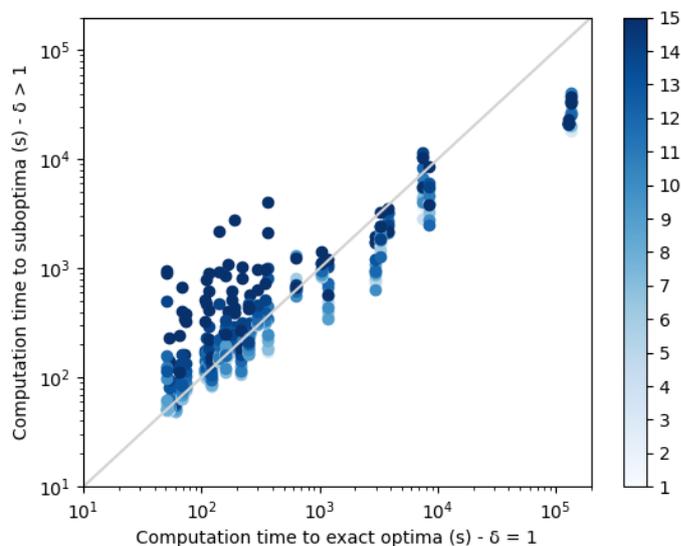


Figure 12. Comparison of the computation times of sequence sets without diversity $\delta = 1$, with suboptimal sequence sets with diversity $\delta > 1$. An energy gap of 3 kcal/mol is allowed to actual optimum.

673 This would enable the efficient production of diverse solutions for otherwise unsolved
674 structured instances.

675 Another direction would be to identify a formulation of the DIST (and possibly
676 DIV_{\min}) constraints that would provide better pruning or avoid the introduction of extra
677 variables that often disturb dynamic variable ordering heuristics. One possibility would
678 be to encode these using linear (knapsack) constraints for which dedicated propagators
679 would need to be developed.

680 **Author Contributions:** Conceptualization, S.B. and T.S.; methodology, M.R., G.K. and T.S.; soft-
681 ware, M.R., J.V.; S.d.G and T.S.; validation, M.R.; formal analysis, M.R.; resources, T.S. and S.B.;
682 data curation, M.R., J.V., T.S. and S.B.; writing—original draft preparation, M.R.; writing—review
683 and editing, M.R., T.S. and S.B.; supervision, T.S. and S.B.; funding acquisition, T.S. All authors
684 have read and agreed to the published version of the manuscript.

685 **Funding:** This research was funded by the French “Agence Nationale de la Recherche” through
686 grants ANR-16-CE40-0028 and ANR-19-PI3A-0004.

687 **Data Availability Statement:** No data specific to this paper. The benchmarks we use are available
688 on existing repositories and the code of *toulbar2*, that includes our contributions, is available on
689 GitHub at <https://github.com/toulbar2/toulbar2> under an OSI MIT license.

690 **Acknowledgments:** We thank the GenoToul (Toulouse, France) Bioinformatics platform and the
691 CALMIP HPC platform for their computational support.

692 **Conflicts of Interest:** The authors declare no conflict of interest.

693 Abbreviations

694 The following abbreviations are used in this manuscript:

695	CFN	Cost Function Network
	CPD	Computational Protein Design
	CSP	Constraint Satisfaction Problem
696	MRF	Markov Random Field
	NSR	Native Sequence Recovery
	NSSR	Native Sequence Similarity Recovery
	WCSP	Weighted Constraint Satisfaction Problem

References

1. Anfinsen, C.B. Principles that govern the folding of protein chains. *Science* **1973**, *181*, 223–230.
2. Pierce, N.A.; Winfree, E. Protein design is NP-hard. *Protein engineering* **2002**, *15*, 779–782.
3. Van Laarhoven, P.J.; Aarts, E.H. Simulated annealing. In *Simulated annealing: Theory and applications*; Springer, 1987; pp. 7–15.
4. Leaver-Fay, A.; Tyka, M.; Lewis, S.M.; Lange, O.F.; Thompson, J.; Jacak, R.; Kaufman, K.W.; Renfrew, P.D.; Smith, C.A.; Sheffler, W.; others. ROSETTA3: an object-oriented software suite for the simulation and design of macromolecules. In *Methods in enzymology*; Elsevier, 2011; Vol. 487, pp. 545–574.
5. Traoré, S.; Allouche, D.; André, I.; De Givry, S.; Katsirelos, G.; Schiex, T.; Barbe, S. A new framework for computational protein design through cost function network optimization. *Bioinformatics* **2013**, *29*, 2129–2136.
6. Allouche, D.; André, I.; Barbe, S.; Davies, J.; de Givry, S.; Katsirelos, G.; O’Sullivan, B.; Prestwich, S.; Schiex, T.; Traoré, S. Computational protein design as an optimization problem. *Artificial Intelligence* **2014**, *212*, 59–79.
7. Noguchi, H.; Addy, C.; Simoncini, D.; Wouters, S.; Mylemans, B.; Van Meervelt, L.; Schiex, T.; Zhang, K.Y.; Tame, J.R.; Voet, A.R. Computational design of symmetrical eight-bladed β -propeller proteins. *IUCr* **2019**, *6*, 46–55.
8. Schiex, T.; Fargier, H.; Verfaillie, G.; others. Valued constraint satisfaction problems: Hard and easy problems. *IJCAI (1)* **1995**, *95*, 631–639.
9. Cooper, M.; de Givry, S.; Schiex, T. Graphical models: queries, complexity, algorithms. *Leibniz International Proceedings in Informatics* **2020**, *154*, 4–1.
10. Bouchiba, Y.; Cortés, J.; Schiex, T.; Barbe, S. Molecular flexibility in computational protein design: an algorithmic perspective. *Protein Engineering, Design and Selection* **2021**, *34*.
11. Marcos, E.; Silva, D.A. Essentials of de novo protein design: Methods and applications. *Wiley Interdisciplinary Reviews: Computational Molecular Science* **2018**, *8*, e1374.
12. King, C.; Garza, E.N.; Mazor, R.; Linehan, J.L.; Pastan, I.; Pepper, M.; Baker, D. Removing T-cell epitopes with computational protein design. *Proceedings of the National Academy of Sciences* **2014**, *111*, 8577–8582.
13. Kirillov, A.; Shlezinger, D.; Vetrov, D.P.; Rother, C.; Savchynskyy, B. M-Best-Diverse Labelings for Submodular Energies and Beyond. NIPS, 2015, pp. 613–621.
14. Bacchus, F.; Van Beek, P. On the conversion between non-binary and binary constraint satisfaction problems. AAI/IAAI, 1998, pp. 310–318.
15. Larrosa, J.; Dechter, R. On the dual representation of non-binary semiring-based CSPs. CP’2000 workshop on soft constraints, 2000.
16. Shapovalov, M.V.; Dunbrack Jr, R.L. A smoothed backbone-dependent rotamer library for proteins derived from adaptive kernel density estimates and regressions. *Structure* **2011**, *19*, 844–858.
17. Lovell, S.C.; Word, J.M.; Richardson, J.S.; Richardson, D.C. The penultimate rotamer library. *Proteins: Structure, Function, and Bioinformatics* **2000**, *40*, 389–408.
18. Case, D.A.; Cheatham III, T.E.; Darden, T.; Gohlke, H.; Luo, R.; Merz Jr, K.M.; Onufriev, A.; Simmerling, C.; Wang, B.; Woods, R.J. The Amber biomolecular simulation programs. *Journal of computational chemistry* **2005**, *26*, 1668–1688.
19. Brooks, B.R.; Brooks III, C.L.; Mackerell Jr, A.D.; Nilsson, L.; Petrella, R.J.; Roux, B.; Won, Y.; Archontis, G.; Bartels, C.; Boresch, S.; others. CHARMM: the biomolecular simulation program. *Journal of computational chemistry* **2009**, *30*, 1545–1614.
20. Alford, R.F.; Leaver-Fay, A.; Jeliaskov, J.R.; O’Meara, M.J.; DiMaio, F.P.; Park, H.; Shapovalov, M.V.; Renfrew, P.D.; Mulligan, V.K.; Kappel, K.; others. The Rosetta all-atom energy function for macromolecular modeling and design. *Journal of chemical theory and computation* **2017**, *13*, 3031–3048.
21. Samish, I. *Computational protein design*; Springer, 2017.
22. Gainza, P.; Roberts, K.E.; Georgiev, I.; Lilien, R.H.; Keedy, D.A.; Chen, C.Y.; Reza, F.; Anderson, A.C.; Richardson, D.C.; Richardson, J.S.; others. OSPREY: protein design with ensembles, flexibility, and provable algorithms. In *Methods in enzymology*; Elsevier, 2013; Vol. 523, pp. 87–107.
23. Pierce, N.A.; Spriet, J.A.; Desmet, J.; Mayo, S.L. Conformational splitting: A more powerful criterion for dead-end elimination. *Journal of computational chemistry* **2000**, *21*, 999–1009.
24. Rossi, F.; van Beek, P.; Walsh, T., Eds. *Handbook of Constraint Programming*; Elsevier, 2006.
25. Koller, D.; Friedman, N. *Probabilistic graphical models: principles and techniques*; MIT press, 2009.
26. Cooper, M.C.; De Givry, S.; Sánchez, M.; Schiex, T.; Zytnecki, M.; Werner, T. Soft arc consistency revisited. *Artificial Intelligence* **2010**, *174*, 449–478.
27. Cooper, M.C.; De Givry, S.; Sánchez-Fibla, M.; Schiex, T.; Zytnecki, M. Virtual Arc Consistency for Weighted CSP. AAI, 2008, pp. 253–258.
28. Traoré, S.; Roberts, K.E.; Allouche, D.; Donald, B.R.; André, I.; Schiex, T.; Barbe, S. Fast search algorithms for computational protein design. *Journal of computational chemistry* **2016**, *37*, 1048–1058.
29. Traoré, S.; Allouche, D.; André, I.; Schiex, T.; Barbe, S. Deterministic Search Methods for Computational Protein Design. In *Computational Protein Design*; Springer, 2017; pp. 107–123.
30. Henikoff, S.; Henikoff, J.G. Amino acid substitution matrices from protein blocks. *Proceedings of the National Academy of Sciences* **1992**, *89*, 10915–10919.

31. Hebrard, E.; Hnich, B.; O'Sullivan, B.; Walsh, T. Finding diverse and similar solutions in constraint programming. *Proceedings of AAAI 2005*, 2005, Vol. 5, pp. 372–377.
32. Hebrard, E.; O'Sullivan, B.; Walsh, T. Distance Constraints in Constraint Satisfaction. *IJCAI*, 2007, Vol. 2007, pp. 106–111.
33. Hadžić, T.; Holland, A.; O'Sullivan, B. Reasoning about optimal collections of solutions. *International Conference on Principles and Practice of Constraint Programming*. Springer, 2009, pp. 409–423.
34. Petit, T.; Trapp, A.C. Finding diverse solutions of high quality to constraint optimization problems. *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.
35. Batra, D.; Yadollahpour, P.; Guzman-Rivera, A.; Shakhnarovich, G. Diverse M-best solutions in Markov Random Fields. *European Conference on Computer Vision*. Springer, 2012, pp. 1–16.
36. Prasad, A.; Jegelka, S.; Batra, D. Submodular meets structured: Finding diverse subsets in exponentially-large structured item sets. *Advances in Neural Information Processing Systems*, 2014, pp. 2645–2653.
37. Kirillov, A.; Savchynskyy, B.; Schlesinger, D.; Vetrov, D.; Rother, C. Inferring M-best diverse labelings in a single one. *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1814–1822.
38. Chen, C.; Kolmogorov, V.; Zhu, Y.; Metaxas, D.; Lampert, C. Computing the M most probable modes of a graphical model. *Proc. of Artificial Intelligence and Statistics*, 2013, pp. 161–169.
39. Chen, C.; Yuan, C.; Ye, Z.; Chen, C. Solving M-Modes in Loopy Graphs Using Tree Decompositions. *Proc. of the International Conference on Probabilistic Graphical Models*, 2018, pp. 145–156.
40. Chen, C.; Liu, H.; Metaxas, D.; Zhao, T. Mode estimation for high dimensional discrete tree graphical models. *Proceedings of Advances in neural information processing systems*, 2014, pp. 1323–1331.
41. Chen, C.; Yuan, C.; Chen, C. Solving M-Modes Using Heuristic Search. *Proc. of IJCAI'16*, 2016, pp. 3584–3590.
42. Pesant, G. A regular language membership constraint for finite sequences of variables. *International conference on principles and practice of constraint programming*. Springer, 2004, pp. 482–495.
43. Allouche, D.; Bessiere, C.; Boizumault, P.; De Givry, S.; Gutierrez, P.; Lee, J.H.; Leung, K.L.; Loudni, S.; Métivier, J.P.; Schiex, T.; others. Tractability-preserving transformations of global cost functions. *Artificial Intelligence* **2016**, *238*, 166–189.
44. Ruffini, M.; Vucinic, J.; de Givry, S.; Katsirelos, G.; Barbe, S.; Schiex, T. Guaranteed Diversity & Quality for the Weighted CSP. *2019 IEEE 31st International Conference on Tools with Artificial Intelligence (ICTAI)*. IEEE, 2019, pp. 18–25.
45. Allouche, D.; De Givry, S.; Katsirelos, G.; Schiex, T.; Zytnecki, M. Anytime hybrid best-first search with tree decomposition for weighted CSP. *International Conference on Principles and Practice of Constraint Programming*. Springer, 2015, pp. 12–29.
46. Simoncini, D.; Allouche, D.; de Givry, S.; Delmas, C.; Barbe, S.; Schiex, T. Guaranteed discrete energy optimization on large protein design problems. *Journal of chemical theory and computation* **2015**, *11*, 5980–5989.
47. Ollikainen, N.; Kortemme, T. Computational protein design quantifies structural constraints on amino acid covariation. *PLoS Comput Biol* **2013**, *9*, e1003313.
48. Park, H.; Bradley, P.; Greisen Jr, P.; Liu, Y.; Mulligan, V.K.; Kim, D.E.; Baker, D.; DiMaio, F. Simultaneous optimization of biomolecular energy functions on features from small molecules and macromolecules. *Journal of chemical theory and computation* **2016**, *12*, 6201–6212.
49. Pohl, I. Heuristic search viewed as path finding in a graph. *Artificial intelligence* **1970**, *1*, 193–204.
50. Xu, J.; Berger, B. Fast and accurate algorithms for protein side-chain packing. *Journal of the ACM (JACM)* **2006**, *53*, 533–557.
51. Jou, J.D.; Jain, S.; Georgiev, I.S.; Donald, B.R. BWM*: A novel, provable, ensemble-based dynamic programming algorithm for sparse approximations of computational protein design. *Journal of Computational Biology* **2016**, *23*, 413–424.
52. De Givry, S.; Schiex, T.; Verfaillie, G. Exploiting tree decomposition and soft local consistency in weighted CSP. *AAAI*, 2006, Vol. 6, pp. 1–6.