

Contents

List of Figures	iii
List of Tables	iv
List of Equations	v
List of Abbreviations	vi
1 Introduction	1
2 Materials & Methods	2
3 Residue contacts predicted by evolutionary covariance extend the application of <i>ab initio</i> molecular replacement to larger and more challenging protein folds	3
4 Evaluation of ROSETTA distance-restraint energy functions on contact-guided <i>ab initio</i> structure prediction	4
4.1 Introduction	5
4.2 Materials & Methods	6
4.2.1 Target selection	6
4.2.2 Covariance-based contact prediction	6
4.2.3 Contact pair to ROSETTA distance restraint formatting	7
4.2.4 <i>Ab initio</i> structure prediction	9
4.2.5 Molecular Replacement	10
4.3 Results	10
4.3.1 Direct comparison of three contact metapredictors	10
4.3.2 Protein structure prediction with two ROSETTA energy functions .	15
4.3.3 Impact of metapredictors and energy functions on unconventional Molecular Replacement	22
4.4 Discussion	34
5 Alternative <i>ab initio</i> structure prediction algorithms for AMPLE	36
6 Decoy subselection using contact information to enhance MR search model creation	37
6.1 Introduction	38
6.2 Materials & Methods	38
6.2.1 Target selection	38
6.2.2 Computation of range-specific satisfaction scores	39
6.2.3 Decoy subselection	39
6.2.4 Molecular Replacement	40
6.3 Results	40

6.3.1	Contact pair satisfaction correlates with decoy quality	40
6.3.2	Long-range contact satisfaction metric to filter decoy sets	43
6.3.3	AMPLE's cluster-and-truncate approach with filtered decoy sets . .	46
6.3.4	Decoy subselection extends AMPLE's performance	50
7	Protein fragments as search models in Molecular Replacement	52
7.1	Introduction	53
7.2	Materials & Methods	54
7.2.1	Target selection	54
7.2.2	Fragment picking using FLIB	54
7.2.3	Molecular Replacement in MRBUMP	56
7.2.4	Assessment of FLIB fragments	56
7.3	Results	57
7.3.1	Precision of FLIB input data	57
7.3.2	FLIB fragment picking	61
7.3.3	FLIB fragment selection for Molecular Replacement	65
7.3.4	Molecular Replacement using FLIB fragments	71
7.4	Discussion	76
8	Conclusion	79
A	Appendix	80
	Bibliography	84

List of Figures

4.1	Schematic comparison of ROSETTA energy functions	7
4.2	Precision analysis of three metapredictors	11
4.3	Sequence coverage and contact precision analysis	12
4.4	Contact singleton analysis for three metapredictors	12
4.5	Comparison of contact precision for three metapredictors	14
4.6	Metapredictor contact pair similarity analysis	15
4.7	Median TM-score comparison between ROSETTA energy functions	17
4.8	Top TM-score comparison between ROSETTA energy functions	18
4.9	TM-score distribution by fold category and ROSETTA energy function	19
4.10	Median TM-score analysis by fold category and ROSETTA energy function	20
4.11	Influence of target chain length and restraint precision on median TM-score	21
4.12	Structure solution count from AMPLE-derived search models	23
4.13	Relationship between median TM-score and search model size of AMPLE ensembles	24
4.14	Relationship between median TM-score and decoy count in SPICKER cluster	26
4.15	SPICKER cluster sizes grouped by restraint condition	26
4.18	RIO score analysis of successful targets	29
4.19	Example of successfully placed AMPLE search model	30
4.20	Example of successfully placed AMPLE search model	32
4.21	Example of successfully placed AMPLE search model	33
6.1	Linear regression model between decoy TM-score and contact satisfaction	42
6.2	Top-1 decoy TM-score and contact satisfaction analysis	43
6.3	TM-score comparison pre- and post-decoy subselection	45
6.4	Effect of decoy subselection on SPICKER clusters	48
6.5	Effect of decoy subselection on THESEUS variance	49
6.6	Molecular Replacement summary of decoy-subselected AMPLE ensembles	51
7.1	PSIPRED schema for FLIB targets	58
7.2	Contact map comparison for FLIB targets	60
7.3	SPIDER2 torsion angle prediction analysis of FLIB targets	61
7.4	FLIB fragment library comparison	62
7.5	Coverage and precision of Flib fragment libraries	64
7.6	Spearman rank-order correlation coefficient analysis of FLIB fragments	66
7.7	Correlation analysis for final FLIB Molecular Replacement (MR) fragments	67
7.8	Distribution of contact precision for FLIB fragments	68
7.9	Fragment search models derived from FLIB	70
7.10	MR structure solutions by FLIB target	71
7.11	MR structure solutions by FLIB library	72
7.12	MR structure solutions by input parameters	73
7.13	Relationship between fragment chain length and normalised RIO scores.	74
7.14	Example of FLIB fragment to MR solution	75

List of Tables

4.1	Summary of AMPLE keyword arguments for FADE and SIGMOID ROSETTA energy functions.	9
6.1	Correlation analysis between decoy TM-score and contact satisfaction	41
7.1	Contact prediction summary for FLIB targets	59
7.2	FLIB fragment characterics across four protein targets	62
A.1	Summary of the ORIGINAL dataset.	81
A.2	Summary of the PREDICTORS dataset.	82
A.3	Summary of the TRANSMEMBRANE dataset.	83

List of Equations

List of Abbreviations

ACL	Average Chain Length
CC	Correlation Coefficient
CMO	Contact Map Overlap
eLLG	expected Log-Likelihood Gain
FP	False Positive
KDE	Kernel Density Estimate
MAE	Mean Absolute Error
MR	Molecular Replacement
MSA	Multiple Sequence Alignment
NMR	Nuclear Magnetic Resonance
PDB	Protein Data Bank
RIO	Residue-Independent Overlap
RMSD	Root-Mean-Square Deviation
TM-score	Template-Modelling score
TP	True Positive

Chapter 1

Introduction

Chapter 2

Materials & Methods

Chapter 3

Residue contacts predicted by evolutionary covariance extend the application of *ab initio* molecular replacement to larger and more challenging protein folds

Chapter 4

Evaluation of ROSETTA distance-restraint energy functions on contact-guided *ab initio* structure prediction

4.1 Introduction

The extended tractability of the AMPLE program for globular protein targets through the use of residue-residue contact information to restrain *ab initio* structure prediction has been highlighted in chapter XYZ. However, that study only focused on PCONSC2 as a metapredictor without considering alternatives, and thus served only as a proof-of-principle work for applications of contact information in unconventional MR.

Besides the individual contact prediction algorithms employed by the PCONSC2 protocol, numerous metapredictors have been developed exploiting different combinations of starting alignments and individual contact predictors to identify the strongest correlating pairs for optimal contact prediction [1–7]. Furthermore, each of those protocols typically includes its own post-prediction algorithms to find a consensus amongst individual predictions and/or further identify patterns characteristic for residue pairings between secondary structure elements in a protein fold. Thus, depending on the overall protocol, the resulting predictions may differ significantly despite the same underlying algorithms to generate starting alignments and to predict residue contact pairs.

Furthermore, the precision of contact predictions used as distance restraints in *ab initio* structure prediction improves the accuracy of the folding process significantly. However, a diversity of structure prediction protocols, whether fragment-based or not, have been applied and each with a unique integration of contact information as distance restraints [3, 8–12]. Such divergence results in three major problems: (1) researchers cannot directly compare results, and thus have to test each protocol against their own with every newly published approach; (2) novice users might find it difficult to make appropriate decisions given the diversity of algorithms and lack of comparative studies; and (3) users only interested in the information encoded in predicted contact pairs are at risk of picking the most readily available approach over the most accurate for their problem.

Thus, the work presented in this chapter was aimed at extensively comparing state-of-the-art contact- and structure-prediction protocols with a focus on the use of such decoys for AMPLE users.

4.2 Materials & Methods

4.2.1 Target selection

This study was conducted using 18 out of 27 targets from the PREDICTORS dataset (??). The nine targets with alignment depths of < 100 in the Pfam Multiple Sequence Alignment (MSA) were excluded (Table A.2).

4.2.2 Covariance-based contact prediction

Residue contacts for each target sequence were predicted using three different metapredictors, namely METAPSICOV [3], GREMLIN [1], and PCONSC2 [2]. Online servers for METAPSICOV (<http://bioinf.cs.ucl.ac.uk/METAPSICOV>) and GREMLIN (<http://gremlin.bakerlab.org>) were used to predict two sets of contact pairs. The choice of online servers over local installations was justified to directly imitate most AMPLE users. Both servers were used with default settings.

The GREMLIN web server returns the raw contact prediction files as well as pre-formatted ROSETTA distance restraints. The raw contact prediction files were downloaded to allow different contact selection thresholds as well as local conversion into ROSETTA restraints files. The METAPSICOV web server returned two contact prediction files, one after Stage 1 and another after Stage 2 post-prediction processing. In this study, contact predictions after Stage 1 (referred to as METAPSICOV from here onwards) were chosen. The PCONSC2 contact prediction set was obtained using a local installation of PCONSC2 due to downtime of the web server at the time of this study. Additionally to the three main contact predictions outlined above, a set of BBCONTACTS restraints was obtained for protein targets containing β -strands (see ??).

The sequence-database versions of all three metapredictors, whether on- or offline, were identical to those used in Chapter 3.

4.2.3 Contact pair to ROSETTA distance restraint formatting

Contact restraints for *ab initio* protein structure prediction were generated by selecting the top-ranking contact pairs from each prediction and reformatting them into a ROSETTA-readable format. The number of top-ranking contact pairs varied according to the two energy functions used (FADE cutoff: L ; SIGMOID cutoff: $3L/2$; where L corresponds to the number of residues in the protein chain). Both energy functions are sigmoidal functions and introduced into the ROSETTA folding protocol in the same fashion.

Neither energy function enforces a specified distance between restrained atoms but reward those that meet it. The two energy functions (Fig. 4.1) differ in that the FADE function does not only have an upper but also a lower bound. Based on previous findings [2, 9], the FADE function was set to acknowledge a formed restraint if the participating C β atoms (C α in case of Gly) were within 9Å. In comparison, the SIGMOID function was defined with amino acid specific distances for C β atoms (C α in case of Gly) to recognise the different sizes of each amino acid [1, 11].

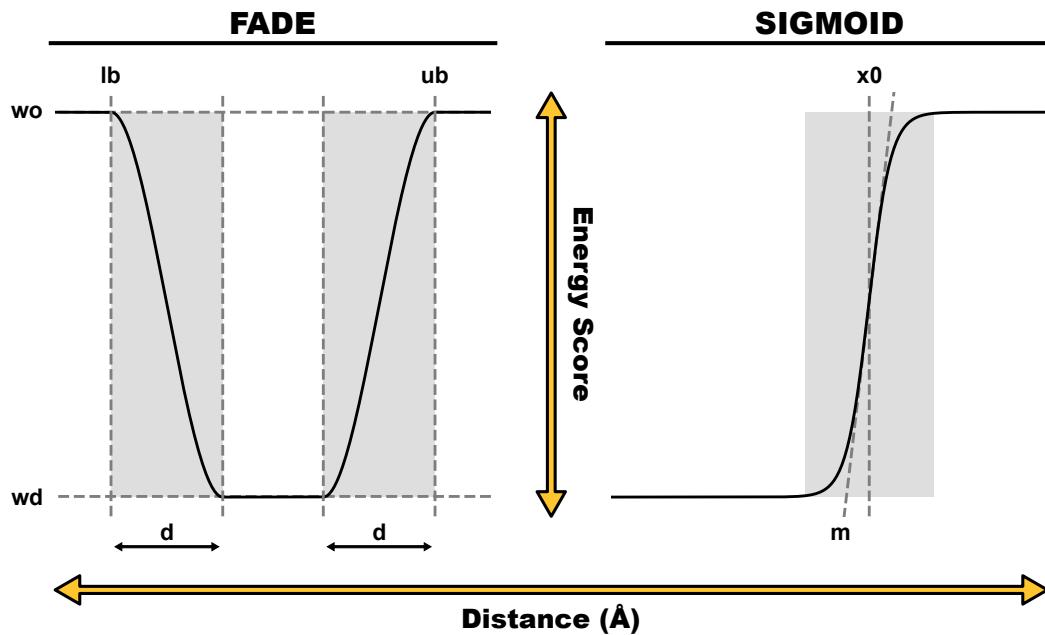


Figure 4.1: ROSETTA energy function comparison. Abbreviations corresponds to input parameters.

To explore the effects of the varying energy function definitions, we created six lists of contact restraints for each α -helical target and nine lists for each β -structure containing

one. The top-ranking contact pairs per prediction were converted using the PCONSFOLD definition of the FADE function [9], the GREMLIN definition of the SIGMOID function [11], and additionally the PCONSC2 BBCONTACTS definition of the FADE function for β -structure containing targets (see Chapter 3).

The conversion was handled in AMPLE and invoked with the keywords outlined in Table 4.1. The `-restraints_factor` keyword defines the factor used to select contact pairs based on the target chain length, i.e. a factor of 1.5 would correspond to $3L/2$ contact pairs. The `-distance_to_neighbour` keyword defines the minimum distance in sequence space between contact pair participating residues, which were set to 5 residues for the FADE function [9] and 3 for the SIGMOID function [11]. Additionally, all distance restraints were given an additional weight when introduced via the SIGMOID energy function to balance its energy term with all remaining terms in the ROSETTA scoring function (Sergey Ovchinnikov, personal communication). This was achieved by using the `-restraints_weight` keyword and weights of 1.0 and 3.0 for the FADE and SIGMOID energy functions.

The addition of BBCONTACTS to existing sets of contacts was achieved with the FADE function in an identical manner as described in Chapter 3. In comparison, the SCALARWEIGHTED term in the GREMLIN implementation of the SIGMOID energy function [11] was multiplied by the number of occurrences of each contact pair in the combined map.

Table 4.1: Summary of AMPLE keyword arguments for FADE and SIGMOID ROSETTA energy functions.

Energy Function	AMPLE keywords
FADE	-contact_file <FILENAME>
	-contact_format <FORMAT>
	-energy_function FADE
	-restraints_factor 1.0
	-distance_to_neighbour 5
(BBCONTACTS)	-restraints_weight 1.0
	-contact_file <FILENAME>
	-contact_format <FORMAT>
	-energy_function FADE
	-restraints_factor 1.0
SIGMOID	-distance_to_neighbour 5
	-restraints_weight 1.0
	-contact_file <FILENAME>
	-contact_format <FORMAT>
	-energy_function SIGMOID
(BBCONTACTS)	-restraints_factor 1.5
	-distance_to_neighbour 3
	-restraints_weight 3.0
	-contact_file <FILENAME>
	-contact_format <FORMAT>
SIGMOID	-energy_function SIGMOID_bbcontacts
	-restraints_factor 1.5
	-distance_to_neighbour 3
	-restraints_weight 3.0

4.2.4 *Ab initio* structure prediction

Six or nine individual lists of contact restraints generated for each target were used in separate ROSETTA *ab initio* protein structure prediction runs. Additionally, protein

structures were predicted without any contact restraints to acquire a control set of decoys. Homologous fragments were excluded during fragment library generation to imitate the folding process of a target with unknown fold. Fragment libraries were generated once per target and used throughout. In total, 1,000 *ab initio* decoys were generated per run using ROSETTAs default settings [13] and one of the seven contact conditions described previously. In total, 162 sets of models were generated across 18 protein targets.

4.2.5 Molecular Replacement

Besides considering model quality, one key interest of this study was the assessment of the model sets created in the previous step as *ab initio* MR search model templates. To reduce the enormous computational cost linked to trialling 162 sets of models, 108 sets were chosen from the following conditions: simple Rosetta, PCONSC2 prediction and FADE function, GREMLIN prediction and SIGMOID function, METAPSICOV prediction and FADE function, and where applicable, PCONSC2 BBCONTACTS, GREMLIN BBCONTACTS and METAPSICOV STAGE 1 BBCONTACTS predictions and FADE function. Overall, this resulted in four MR runs for the six α -helical targets, seven runs for the six all- β , and seven runs for the six mixed α - β targets. The resulting 108 model sets were trialled in AMPLE v1.1.0. Structure solution success was assessed as described in ??.

4.3 Results

4.3.1 Direct comparison of three contact metapredictors

In this study, a direct comparison between three metapredictors — GREMLIN, METAPSICOV and PCONSC2 — was carried out. Residue-residue contact pairs were predicted for 18 protein target sequences with a range of chain lengths and numbers of effective sequences in their Pfam MSAs.

METAPSICOV is the most precise contact predictor across the protein target dataset in this study (Fig. 4.2). The difference between the three metapredictors is most evident in the highest-scoring contact pairs ($L/10$). The median precision values for METAPSICOV

and PCONSC2 contact predictions are above 50% up to L contact pairs. GREMLIN, in comparison, predicts contacts with a median precision score at least 20% worse than that of METAPSICOV and 15% worse than PCONSC2. However, at $3L/2$ contact pairs the median precision scores are much more similar across the three different metapredictors: METAPSICOV and PCONSC2 are near identical, and GREMLIN is at most 12% worse compared to the other two. Inspecting the mean precision scores over a continuous range of selection cutoff values illustrates further the difference between METAPSICOV, PCONSC2 and GREMLIN (Fig 4.3). The former two similarly high precision scores compared to the average precision scores for GREMLIN, which are 0.2 precision score units lower. Added to the difference in precision scores is the difference in sequence coverage (Fig. 4.3). Although producing the on-average worst contact predictions out of the three metapredictors used in this study, GREMLIN contact predictions have the highest sequence coverage. However, an analysis of singleton contact pairs, usually with high degrees of false positives, revealed a positive correlation ($\rho_{Pearson} = 0.47; p < 0.001$) between the fraction of singleton contact pairs and sequence coverage and hints to a weak negative correlation ($\rho_{Pearson} = -0.27; p < 0.05$) between the fraction of singleton contact pairs and contact precision (Fig. 4.4).

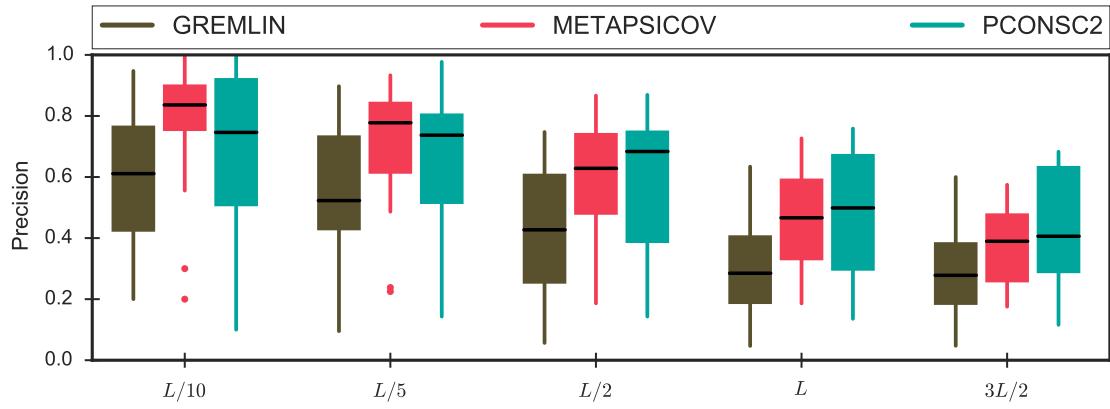


Figure 4.2: Precision spread for three metapredictors computed at five contact selection cutoff values relative to the target chain length (L).

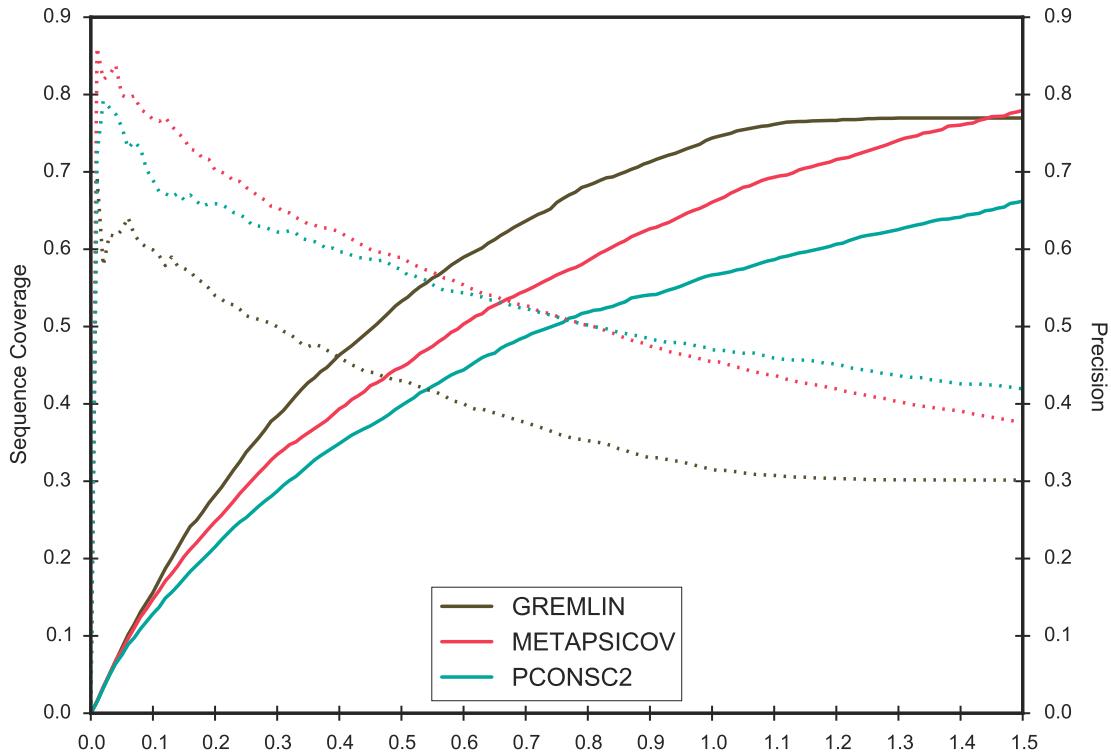


Figure 4.3: Average sequence coverage (line) and contact prediction precision scores (dashed) across a continuous range of contact selection cutoffs ranging from [0.0, 1.5] for all targets.

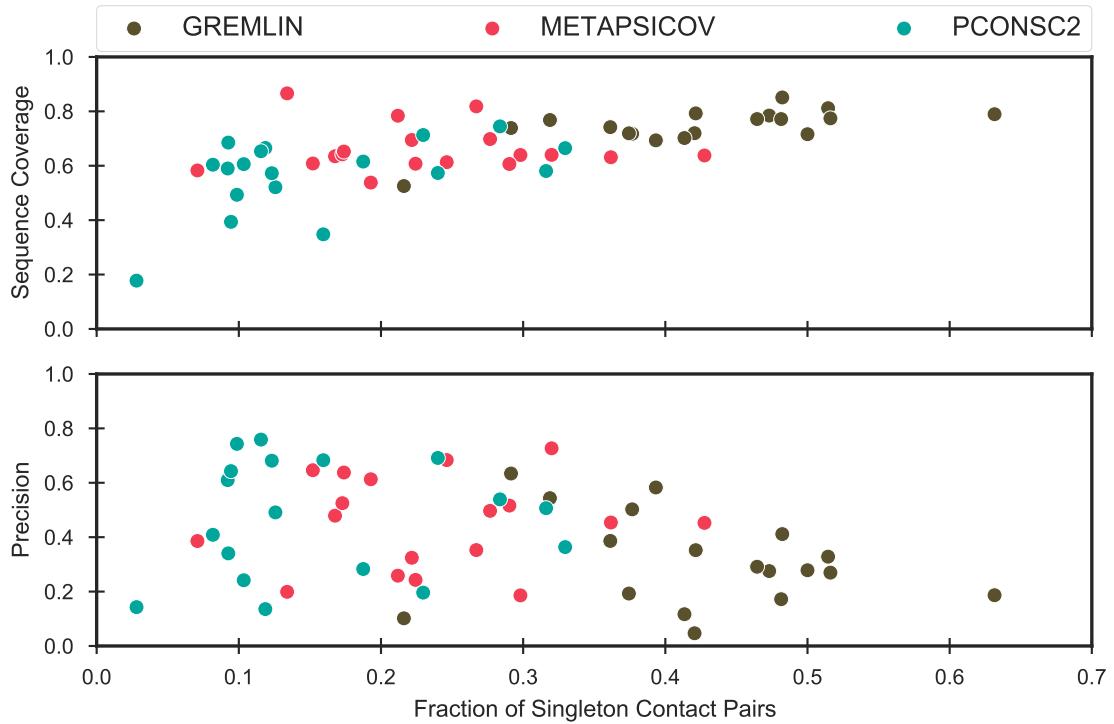


Figure 4.4: Contact singleton analysis compared against the precision of L contact pair lists for three metapredictors.

Given that the overall precision of contact pairs predicted by the three metapredictors differs, it is important to understand where the difference originates. To investigate this, a comparison of the precision values at different cutoff levels on a per-target basis was performed. For the majority of targets the precision scores are very similar across the three metapredictors (Fig. 4.5). However, the prediction precision of some targets differs significantly. For example, the METAPSICOV prediction for the human retinoic acid nuclear receptor HRAR (Protein Data Bank (PDB): 1fcy) contains high precision in its highest scoring (top- $L/10$) contact pairs (Fig. 4.5). In comparison, GREMLIN and PCONSC2 predictions for the same target contain less precise contact pairs ($\Delta\text{Precision}_{\text{METAPSICOV-GREMLIN}} L/10 = -0.522$; $\Delta\text{Precision}_{\text{METAPSICOV-PCONSC2}} L/10 = -0.435$). However, the addition of further contact pairs up to $3L/2$ results in near-identical precision across the three metapredictors for this target. A second example illustrating such a difference are the contact predictions for the human galectin-3 CRD sequence (PDB: 4lbj). In contrast to the previous example, the data shows high precision scores for the METAPSICOV and PCONSC2 predictions for this target, yet low precision for the top GREMLIN contact pairs ($\Delta\text{Precision}_{\text{METAPSICOV-GREMLIN}} L/10 = -0.231$; $\Delta\text{Precision}_{\text{METAPSICOV-PCONSC2}} L/10 = +0.077$).

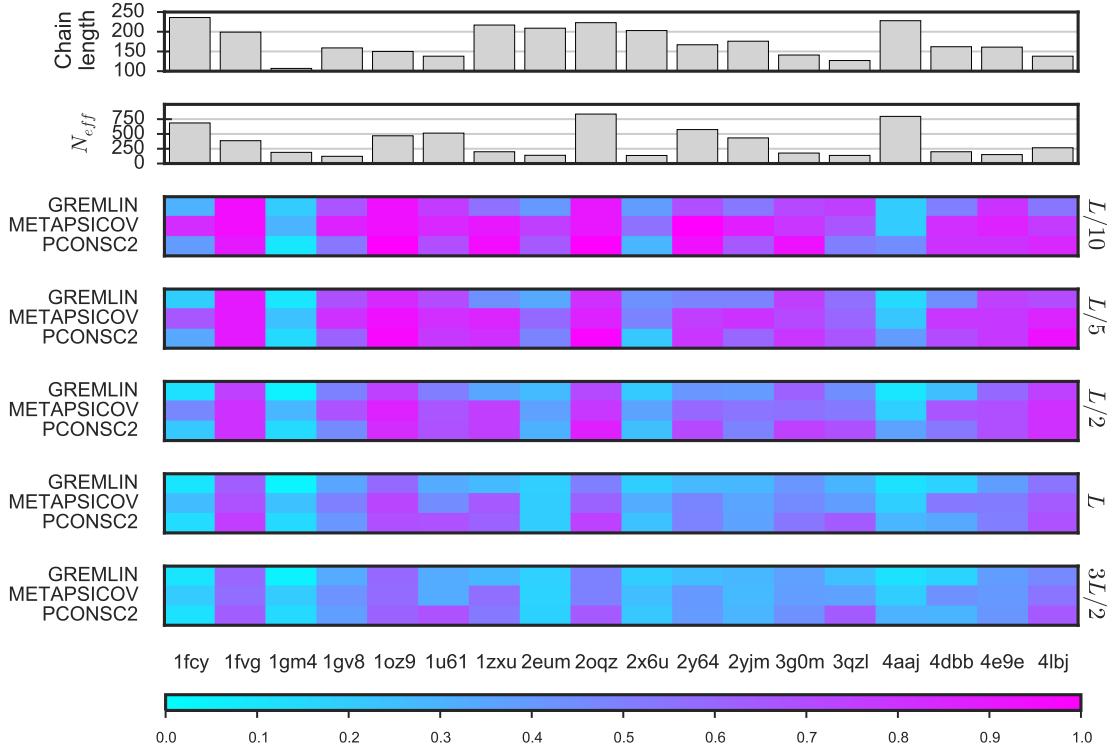


Figure 4.5: Contact prediction precision scores from three metapredictors for 18 targets at different contact pair selection thresholds. The Pfam alignment depth is given by means of number of effective sequences (N_{eff}). The color scale corresponds to the precision in $[0, 1]$.

The data presented in Fig. 4.5 also indicates that there is no direct link between chain length or N_{eff} and the precision of the resulting contact predictions. The N-(5'-phosphoribosyl)anthranilate isomerase sequence (PDB: 4aaaj) with a chain length of 228 residues and 750 effective sequences in its Pfam MSA yielded a mean precision at $L/10$ contact pairs of 0.283 (top- L : 0.195) across the three metapredictors. This strongly contrasts with the sequence of sortase B (PDB: 2oqz), which shows similar characteristics yet obtained mean precision at $L/10$ contact pairs of 0.938 (top- L : 0.622).

Although the contact predictions differ in precision, an interesting question rests with the similarity of the predicted contact pairs amongst the sets. Thus, the similarity of contact predictions across the three metapredictors is an important metric to evaluate the most appropriate algorithm for AMPLE users. Using the Jaccard similarity index to evaluate the direct overlap of contact pairs across sets of predictions, the data suggests very little similarity between the contact predictions of the three metapredictors for each target (Fig. 4.6). As with the differences in precision scores at higher cutoff thresholds,

the Jaccard index is also lower — indicating less overlap — at higher cutoff thresholds. However, it is worth noting that the Jaccard index only considers identical matches and does not consider the neighbourhood of a contact pair. Thus, the index does not highlight similar regions with contact pairs in both maps.

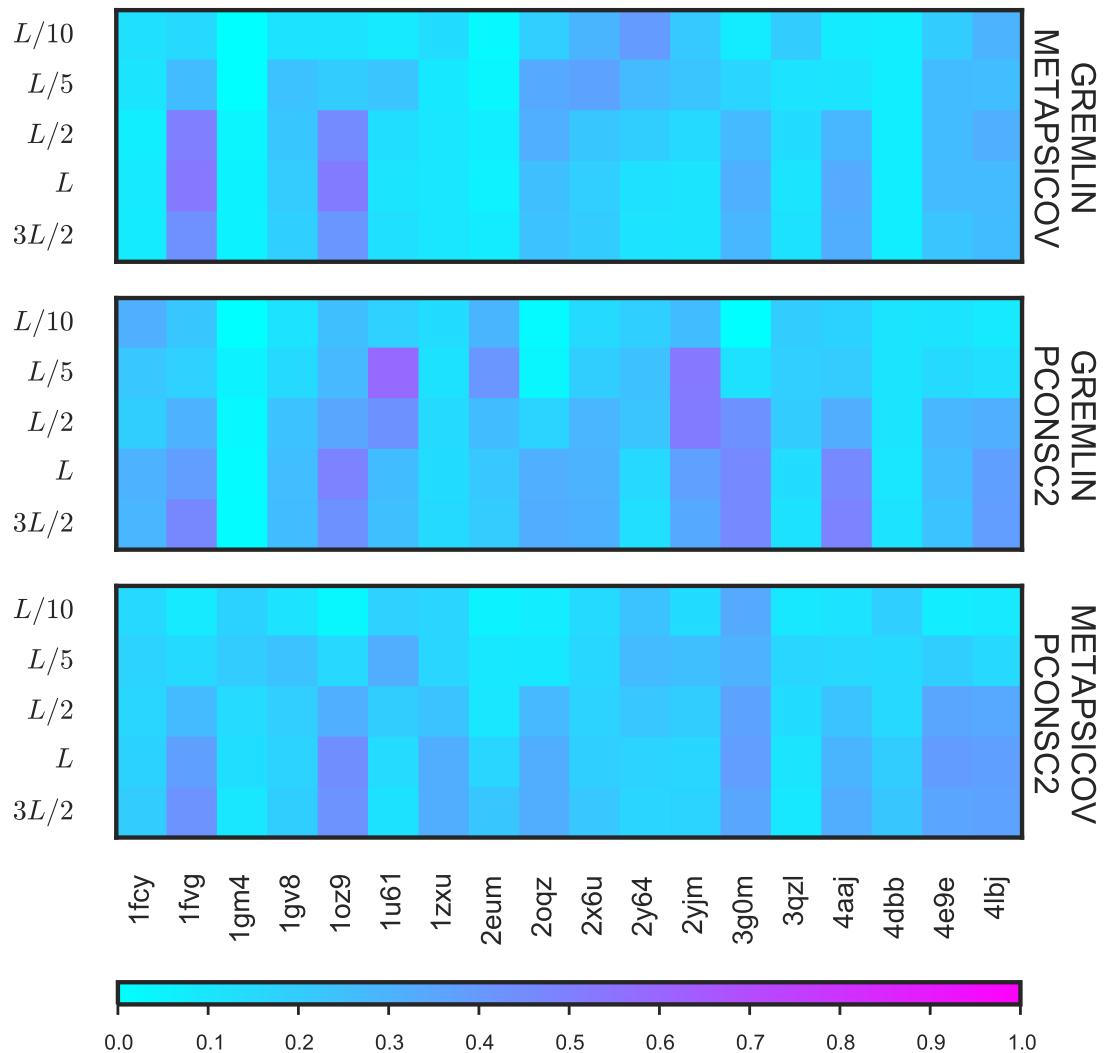


Figure 4.6: Jaccard similarity index illustrates a higher degree of overlap between metapredictor contact predictions with increasing numbers of contact pairs included in the calculation. The three panels show the different comparisons. The color scale corresponds to the Jaccard index in $[0, 1]$.

4.3.2 Protein structure prediction with two ROSETTA energy functions

The accuracy of the starting decoys is a major factor for an AMPLE run to succeed [14]. Thus, the quality of the decoys is of great essence to this study. Given the two different ROSETTA energy functions, FADE and SIGMOID, all contacts predicted were subjected

to individual *ab initio* structure prediction runs. Additionally, all contact predictions were enriched with BBCONTACTS for all β -containing targets in separate trials. A total of 234,000 individual decoys were generated in this study through all permutations of targets, contact predictions and ROSETTA energy function combinations.

Separating these individual decoys solely by the ROSETTA energy function (excluding unrestrained ROSETTA decoys) shows that the FADE energy function results in marginally more accurate decoys (median TM-score FADE: 0.3541; median TM-score SIGMOID: 0.2969). To further investigate which energy function is more suitable for the target dataset used in this study, the decoy sets were grouped by two additional characteristics: the fold of the target, and the source of distance restraints used. The results strongly suggest that the FADE energy function results in more accurate decoy sets (Fig. 4.7), outperforming the SIGMOID energy function by median TM-score in two-thirds of all decoys sets (FADE: 58; SIGMOID: 32). A split of the decoy sets into separate categories by fold and the addition of BBCONTACTS reveals that the SIGMOID energy function only yields similar results for all- β targets in combination with BBCONTACTS-supported distance restraints. Although the total count of decoy sets with higher accuracies between the two energy functions in this category are similar, the actual differences in TM-scores further supports the strength of the FADE energy function compared to the SIGMOID.

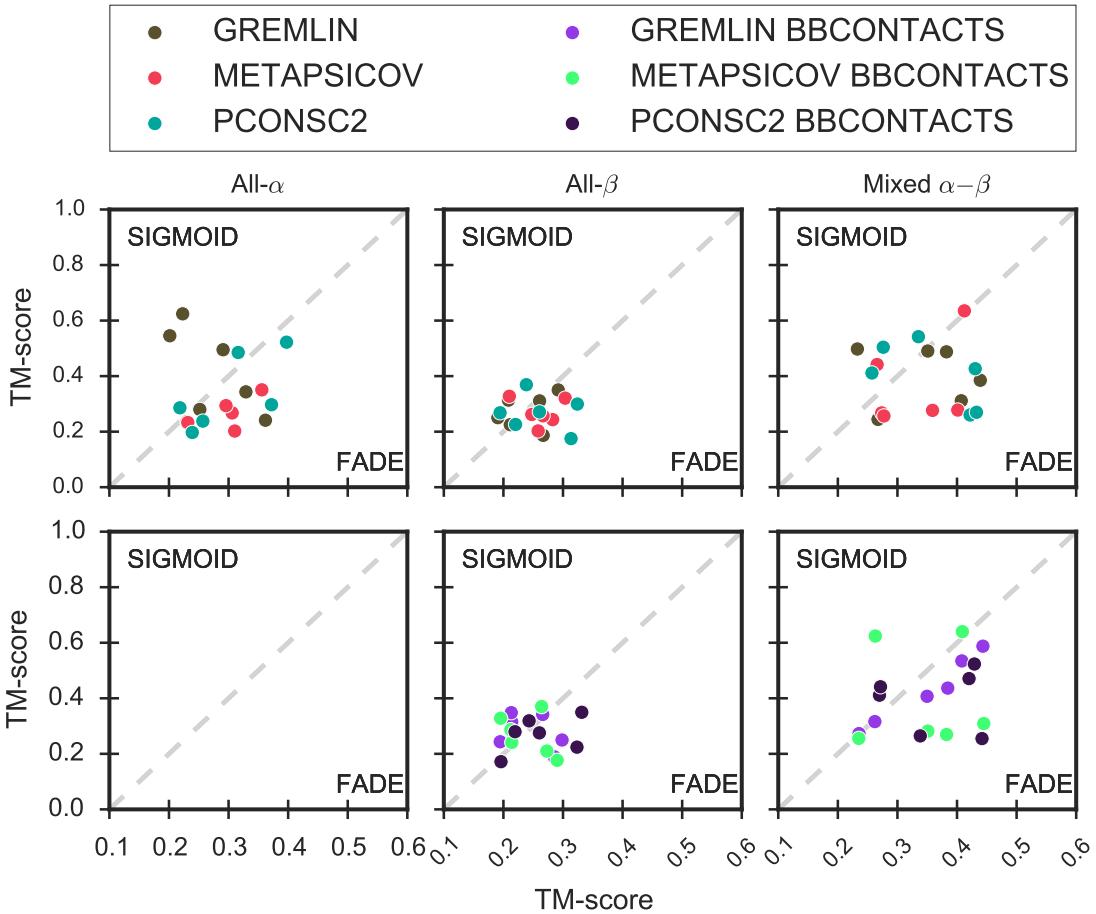


Figure 4.7: Median TM-score comparison of FADE and SIGMOID ROSETTA energy functions differentiated by fold and the addition of BBCONTACTS restraints.

Besides the structure prediction accuracy of each set of decoys, the single, most accurate decoy is also of great interest. If one energy function consistently predicts single decoys more accurately, it might be appropriate to reconsider the structure identification routine (i.e. clustering) in AMPLE for search model preparation. However, a similar difference to that of the decoy quality of entire sets is observed for the top-1 decoy in each set (Fig. 4.8). The FADE energy function outperforms the SIGMOID function for the majority of target-contact prediction permutations (FADE: 51; SIGMOID: 39). However, the GREMLIN distance restraints in combination with the SIGMOID energy function produce better top-1 decoys than GREMLIN restraints with the FADE energy function. This suggests that GREMLIN restraints and the SIGMOID energy function were tailored to complement each other with the ultimate goal of predicting single decoys to high accuracy over entire sets of decoys. Additionally, the spread of decoy quality differences between the two energy functions widens when only looking at the best decoy in each predicted set

(Δ Median Template-Modelling score (TM-score)_{ALL}: $\min = 0.002$, $\max = 0.429$; Δ Median TM-score_{TOP}: $\min = 0.002$, $\max = 0.456$).

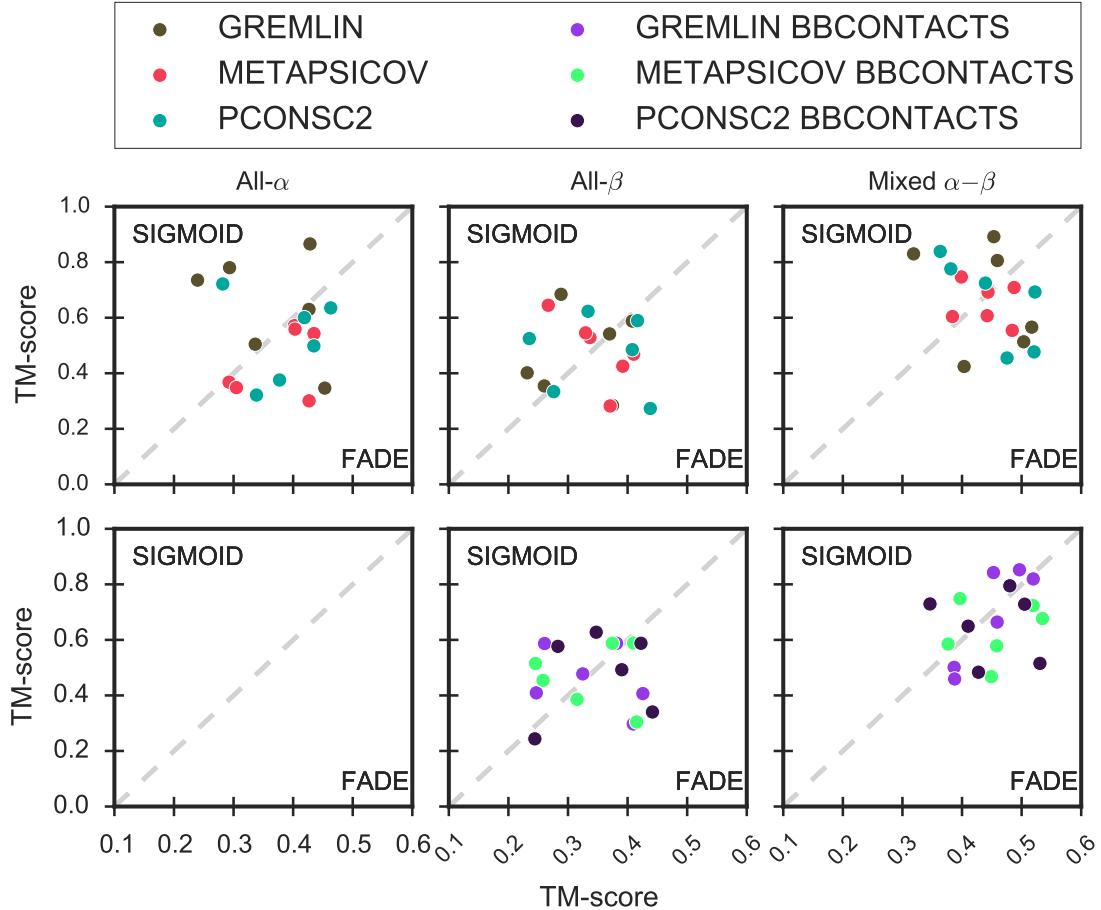


Figure 4.8: Top TM-score comparison of FADE and SIGMOID ROSETTA energy functions differentiated by fold and the addition of BBCONTACTS restraints.

A Kernel Density Estimate (KDE) of TM-scores using each predicted decoy was generated with the TM-scores of individual decoys separated only by fold class and ROSETTA energy function (Fig. 4.9). This density estimate further supports the results presented above: the FADE energy function generates more accurate decoys. However, a very important detail is highlighted by the estimates. Distinct regions with high density are visible in the estimates of the TM-scores of individual decoys for all- α and mixed α - β targets (Fig. 4.9). The bimodal distribution of decoy TM-scores from both energy functions strongly suggests that predicted structures are either native-like or not (based on the TM-score threshold of ≤ 0.5). However, the number of correctly predicted decoys versus incorrectly predicted decoys is in favour of the latter. The decoy sets of all- β targets do not show such distinct regions of high density for decoys with TM-scores < 0.5 units in any

of its density estimates (Fig. 4.9). The generally poor decoy quality of decoys predicted without any distance restraint information (ROSETTA) highlights the benefit of contact predictions to *ab initio* protein structure prediction.

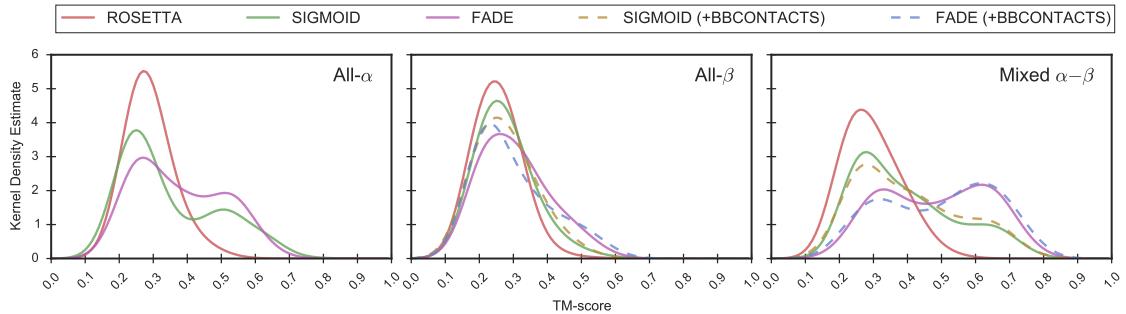


Figure 4.9: TM-score density estimate of all decoys in each respective fold class separating by ROSETTA energy function (SIGMOID or FADE) and no contact information used (ROSETTA). Dashed lines indicate decoys which were predicted with the addition of BBCONTACTS.

A further important aspect of this study is to explore the benefits of adding BBCONTACTS restraints to the structure prediction of β -containing targets. Although previous results ?? in combination with those presented above outline overall improvements in decoy quality, it is essential to understand which targets benefit from this treatment. Figure 4.10a highlights the effects of adding BBCONTACTS restraints to the structure prediction strategies employed here. In summary, the addition of BBCONTACTS restraints hardly affects the decoy quality of most targets under the various contact prediction and energy function combinations. Nevertheless, three target, contact prediction and energy function combinations yielded TM-score improvements of at least 0.1 TM-score units compared to the same condition without the addition of BBCONTACTS restraints. In contrast, the addition of BBCONTACTS restraints did not lower the median TM-score by more than 0.1 units for any target (Fig. 4.10b).

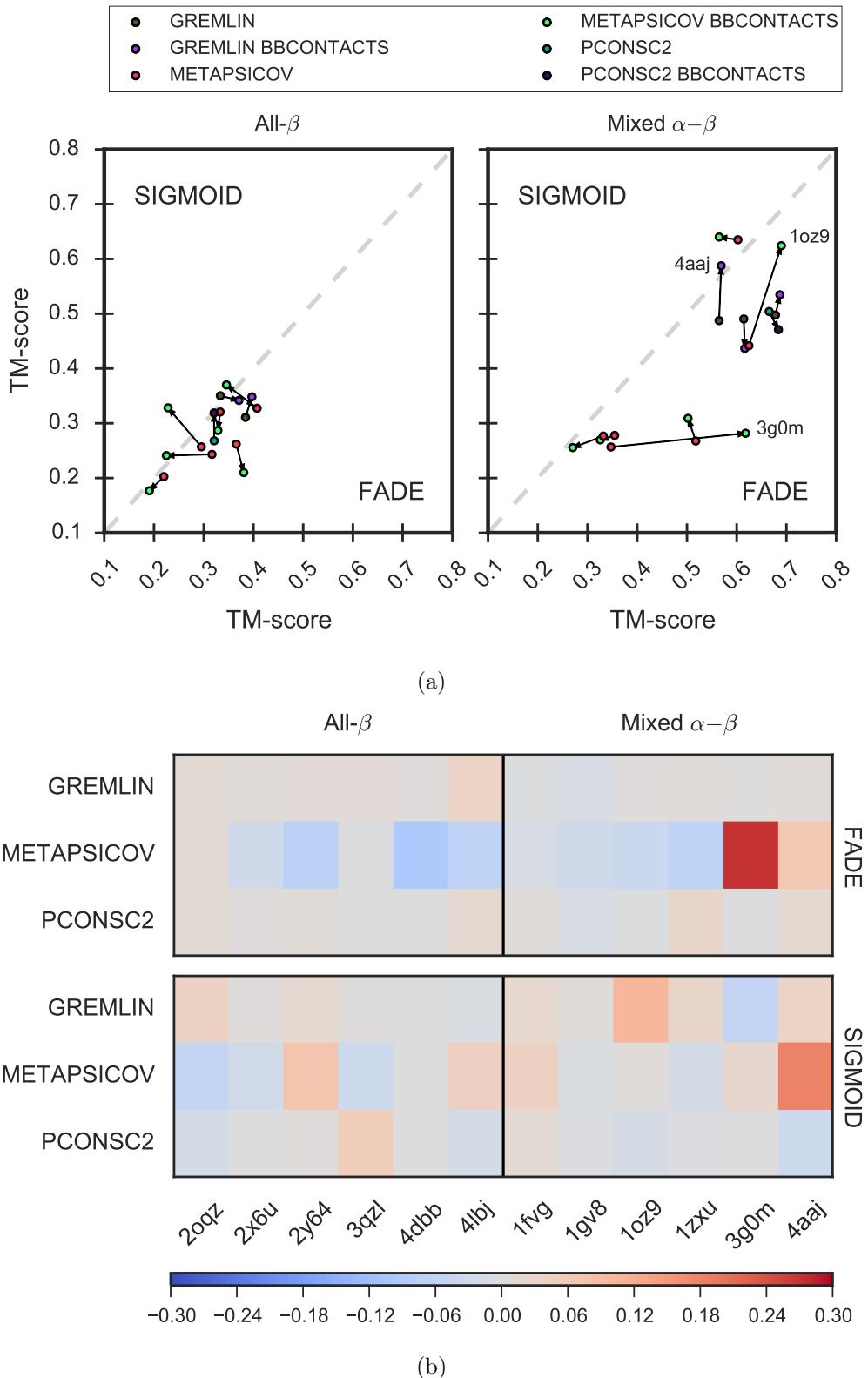


Figure 4.10: Median TM-score comparison of FADE and SIGMOID ROSETTA energy functions differentiated by fold (excl. all- α). (a) Arrows indicate the effect on decoy quality through the addition of BBCONTACTS restraints. Targets with a distance < 0.03 TM-score units between normal and BBCONTACTS-added conditions were excluded from the scatter plots. (b) Effect on decoy quality through the addition of BBCONTACTS restraints highlighted by heatmap difference. The color scale corresponds to the difference in median TM-score between normal and BBCONTACTS-added contact maps.

Two further aspects in understanding the differences in effects of the FADE and SIGMOID ROSETTA energy functions on decoy quality are the target chain length and restraints precision. The former appears to affect the final decoy quality of all 1,000 decoys insignificantly (Fig. 4.11). However, the restraint precision results in some differences between the two ROSETTA energy functions (Fig. 4.11). The FADE energy function (L restraints) generally appears to be less sensitive to restraint lists with higher false positive contact pairs. In contrast, the SIGMOID function ($3L/2$ restraints) produces less accurate decoys than the FADE function with more accurate restraints. Most strikingly, the FADE energy function generated decoys with a median TM-score of 0.678 for the N-(5'-phosphoribosyl)anthranilate isomerase domain (PDB: 4aaJ) compared to the SIGMOID function with a median TM-score of 0.498. Nevertheless, both energy functions appear to broadly follow a positive linear trend, i.e. better restraint precision results in more accurate decoys.

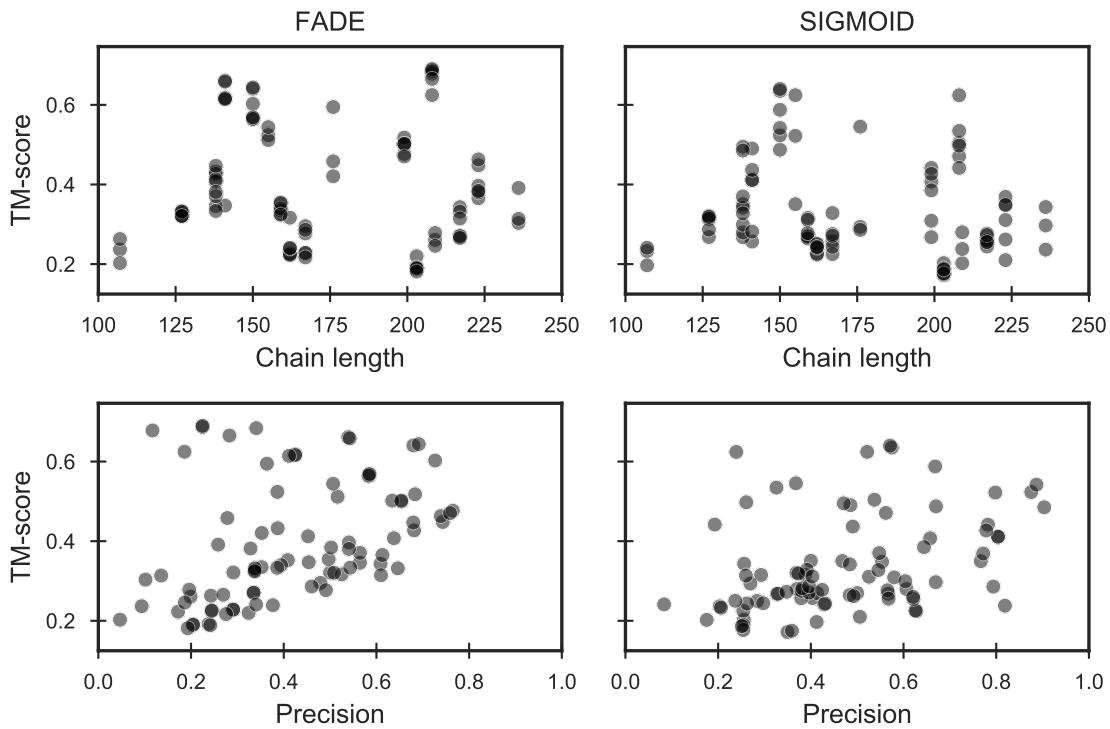


Figure 4.11: Effects of target chain length and restraint precision on the median TM-score for FADE and SIGMOID ROSETTA energy functions. Each scatter point represents a 1,000-decoy set.

4.3.3 Impact of metapredictors and energy functions on unconventional Molecular Replacement

The results obtained from the decoy quality comparison outlined above highlighted differences between the FADE and SIGMOID ROSETTA energy functions. This difference is more pronounced for some targets and less so for others. Thus, the next step in this study was to analyse the consequences of these differences for unconventional MR using the automated pipeline AMPLE.

Overall, the decoys restrained with GREMLIN distance restraints via the SIGMOID energy function throughout the structure prediction process yielded six out of 18 possible structure solutions (Fig. 4.12). This result was the highest of all trialled conditions and only resulted in one more structure solution compared to unrestrained ROSETTA decoys. All remaining conditions resulted in fewer structure solutions than those from ROSETTA decoys. Furthermore, the conditions METAPSICOV (FADE function), METAPSICOV BBCONTACTS (FADE function) and PCONSC2 BBCONTACTS (FADE function) yielded no more than half of the structure solutions achieved by GREMLIN (SIGMOID function). The remaining two conditions — PCONSC2 (FADE function) and GREMLIN BBCONTACTS (FADE function) — resulted in four out of 18 structure solutions. The addition of BBCONTACTS did not improve decoy quality enough to increase the chances of structure solution success; however, the structure of the bovine peptide methionine sulfoxide reductase (PDB: 1fg) was only solved with the GREMLIN BBCONTACTS (FADE function) decoys further supporting the small but important value of BBCONTACTS restraint addition to separately determined contact predictions.

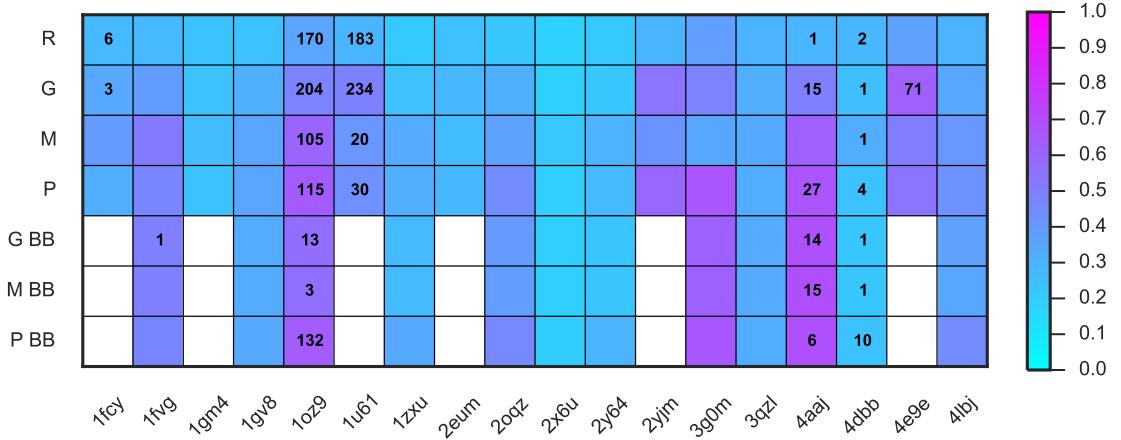


Figure 4.12: Structure solution count for AMPLE search models generated from decoys with varying contact prediction and ROSETTA energy function conditions: unrestrained ROSETTA (R); GREMLIN (G; SIGMOID function); METAPSICOV (M; FADE function); PCONSC2 (P; FADE function); GREMLIN BBCONTACTS (G BB; FADE function); METAPSICOV BBCONTACTS (M BB; FADE function); PCONSC2 BBCONTACTS (P BB; FADE function). The color scale of each square indicates the median TM-score of all 1,000 starting decoys.

The number of structure solutions obtained from the decoy sets subjected to the AMPLE pipeline are somewhat surprising given that ROSETTA decoys result in the second-most structure solutions. These results suggest that the current implementation cannot exploit the true value of more accurate decoy sets. This hypothesis is further supported when considering the decoy set quality and the number of structure solutions (Fig. 4.12). For example, PCONSC2 (FADE function) decoys predicted for the hypothetical protein AQ_1354 (PDB: 1oz9) yield high accuracy, and thus would generally be considered highly desirable starting structures for the AMPLE protocol; nevertheless, the AMPLE protocol was unable to exploit such highly accurate decoys for successful structure solutions of other targets, e.g. cysteine desulfurization protein SufE (PDB: 3g0m; median TM-score PCONSC2 BBCONTACTS (FADE function)=0.661). In comparison, the median TM-scores for all successful ROSETTA decoy sets do not exceed 0.355 TM-score units.

Naturally, one would expect the best decoys to result in the most accurate ensemble search models, which in turn yield the highest number of structure solutions per target. However, here we demonstrate that the most accurate decoys do not guarantee structure solution, and in contrast some poorly predicted decoy sets achieve structure solution. Thus, it is essential to investigate the stage in AMPLE’s cluster-and-truncate approach at

which the higher decoy quality results in less suitable ensemble search models for MR.

The data generated as part of this study reveals a positive correlation ($\rho_{Spearman} = 0.78$; $p < 0.001$) between the decoy quality and the number of resulting AMPLE ensemble search models (Fig. 4.13). The plotted data alongside a line of best fit further illustrate that small differences in decoy quality in the lower TM-score regions increases the total number of generated ensemble search models dramatically. However, once the threshold of 0.5 TM-score units [15] is surpassed the number of generated ensemble search models plateaus at around 350-400 ensemble search models, approaching the maximum number of search models generatable by AMPLE. Furthermore, the data suggests that sets containing fewer than 100 ensemble search models do not lead to structure solution, although this result needs to be considered with care given the difficulty of predicting which search model will lead to structure solution.

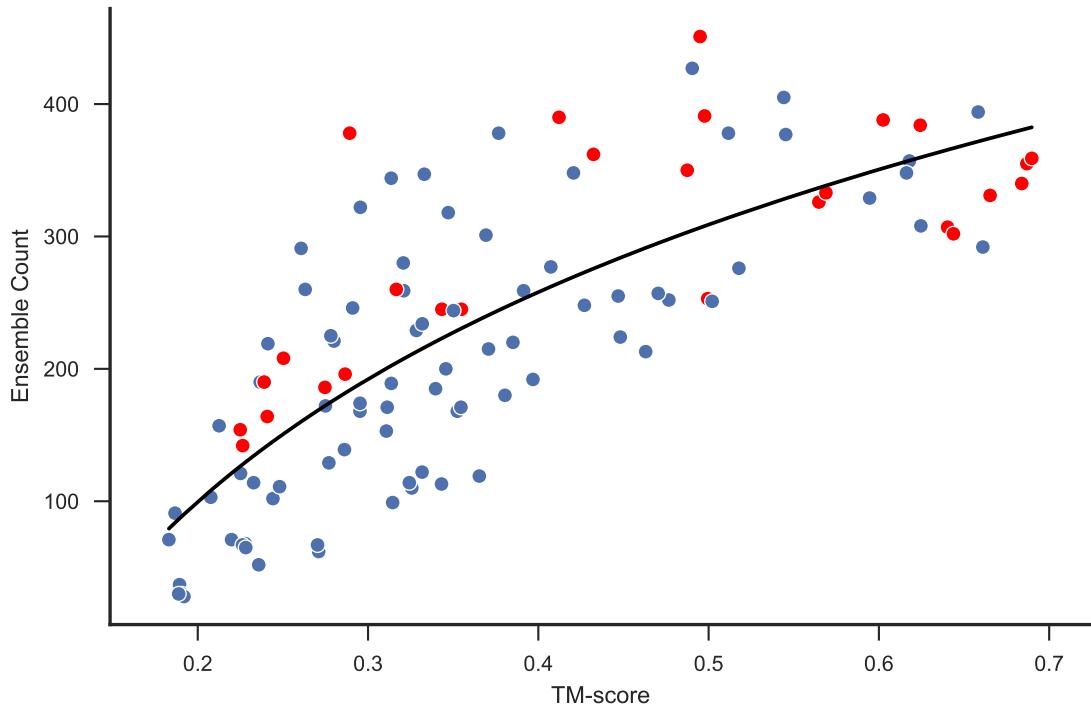


Figure 4.13: Comparison of median TM-score (per 1,000 decoys) against the resulting AMPLE ensemble search model count. The equation of the line of best fit is defined by $y = 228.50 * \ln(20.96 * x) - 227.95$. Red dots indicate successful ensemble sets.

Besides looking at the relationship between entire decoy sets and the resulting structure solutions on a per-target or per-condition basis, it is important to also consider individual ensemble search models, their origins and their properties in relation to MR metrics.

Previous findings highlighted the relationship between the number of decoys in the first cluster and the quality of the decoys it contains (see Chapter 3). Here, we further support these findings given the positive relationship between the median TM-scores and the corresponding size of the largest SPICKER cluster (Fig. 4.14). An analysis of the cluster sizes demonstrates the downstream benefits of increased decoy quality through contact restraints in the folding process (Fig. 4.15). The sizes of the first three clusters generated from most contact-restraint decoy sets greatly surpass their equivalent cluster sizes for unrestrained ROSETTA decoys. Given that cluster sizes correlate with decoy quality, the findings in this study also support that the mean C α Root-Mean-Square Deviation (RMSD) — as calculated by THESEUS for cluster truncation — is directly related to better decoy quality via the larger number of decoys in each cluster (Fig. 4.16a). The same mean C α RMSD is also related to the number of ensemble search models generated after subclustering (Fig. 4.16b), which hints towards a direct relationship between increased quality of 1,000 decoys per set and the total number of ensemble search models generated. Interestingly, GREMLIN decoys show similar C α RMSD per cluster compared to unrestrained ROSETTA decoys (Fig. 4.17), unlike all other contact restraint guided structure predictions. However, it is worth noting that almost no distinction can be made amongst the remaining contact restraint treatments albeit some differences in cluster size distributions exist (Fig. 4.15).

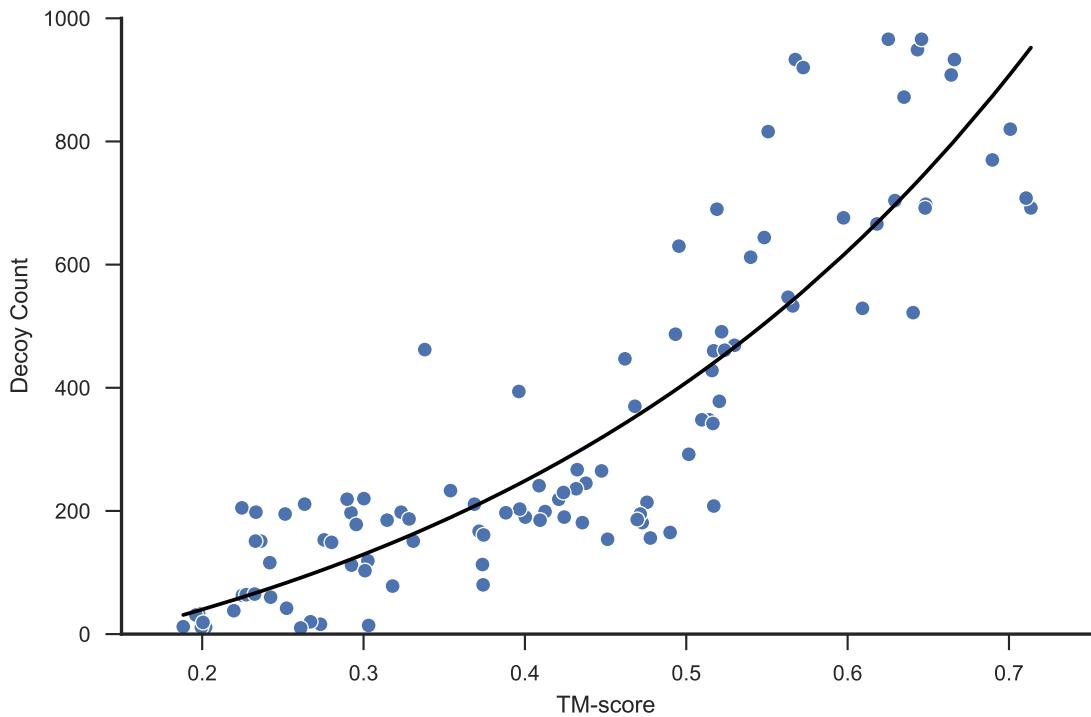


Figure 4.14: Relationship between cluster median TM-score and the number of cluster decoys. Blue line represents line of best fit with equation $y = 148.85 * \exp(2.90 * x) - 225.76$.

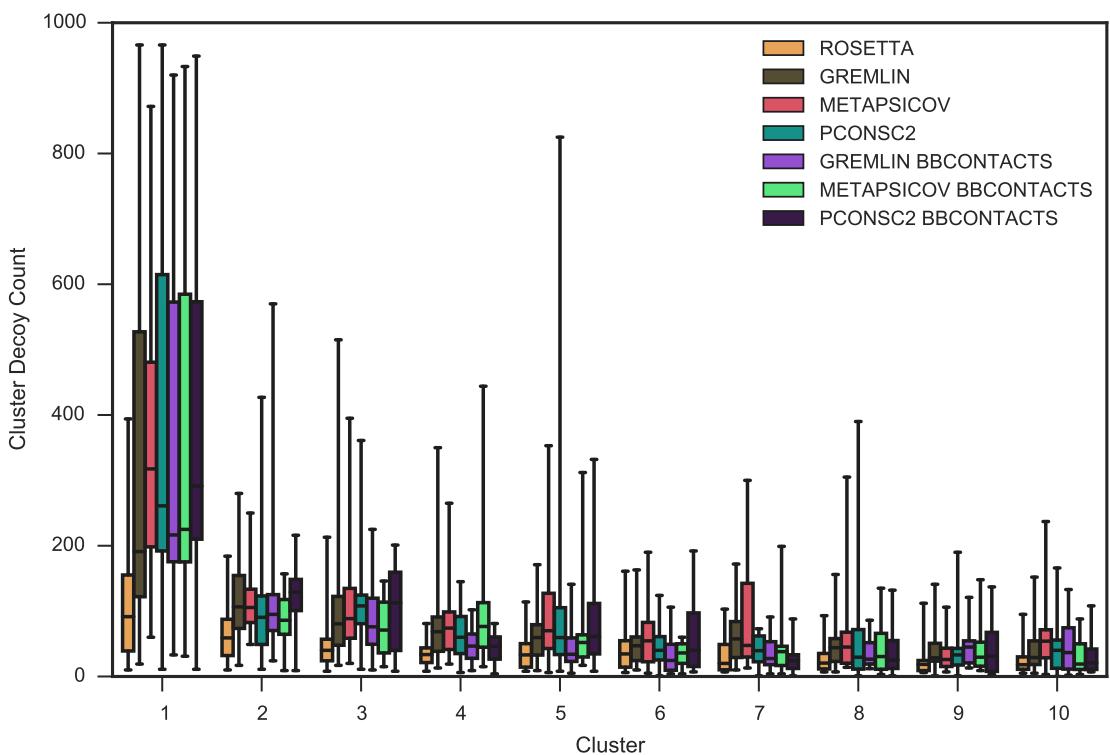


Figure 4.15: SPICKER cluster sizes of each target grouped the restraint condition used during the structure prediction protocol. Whiskers span the range from the minimum to maximum counts.

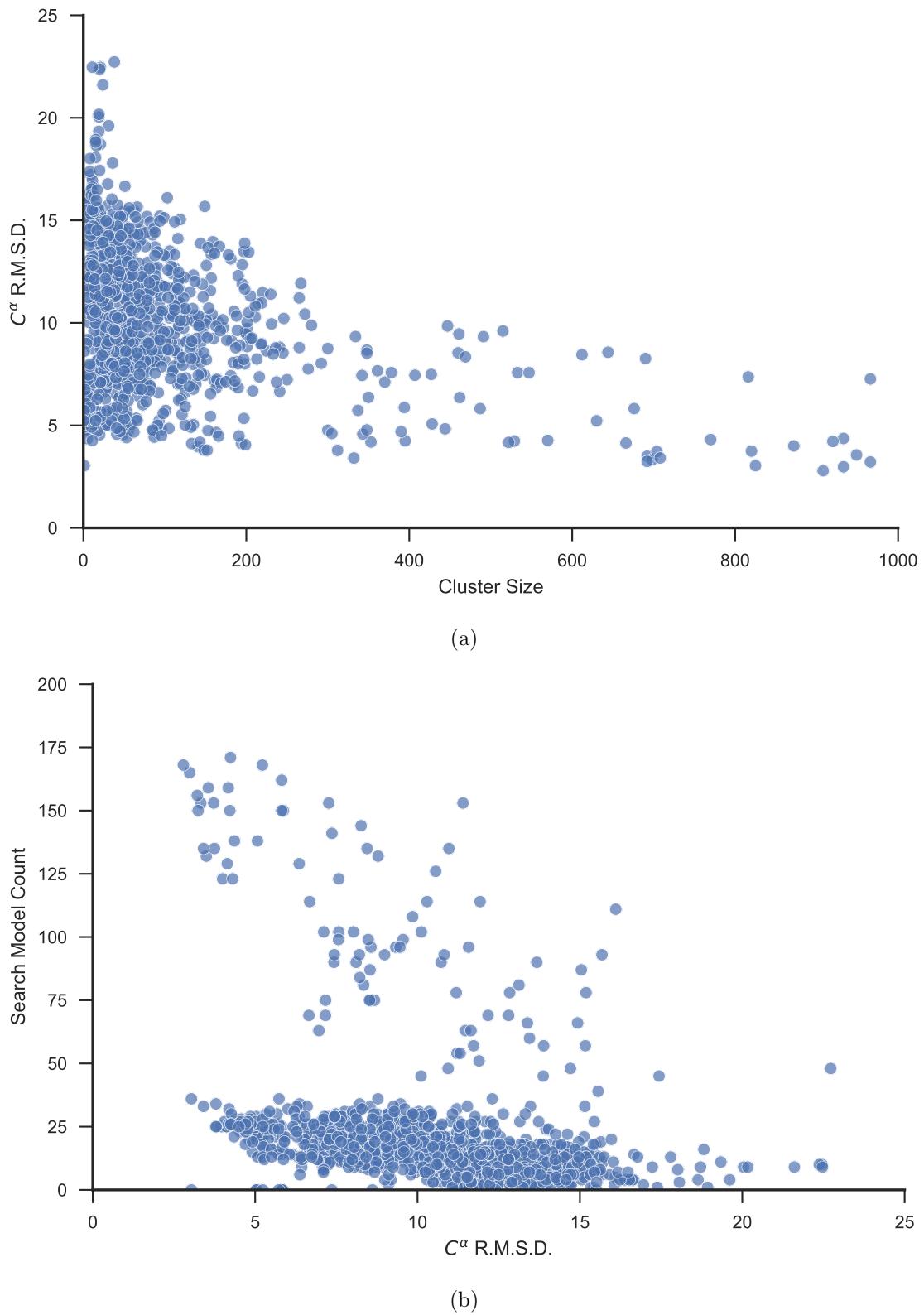


Figure 4.16: (a) Number of decoys per SPICKER cluster plotted against the mean C^α -atom RMSD for all decoys in each cluster. (b) Mean C^α -atom RMSD for decoys per cluster plotted against the number of search models derived from the cluster.

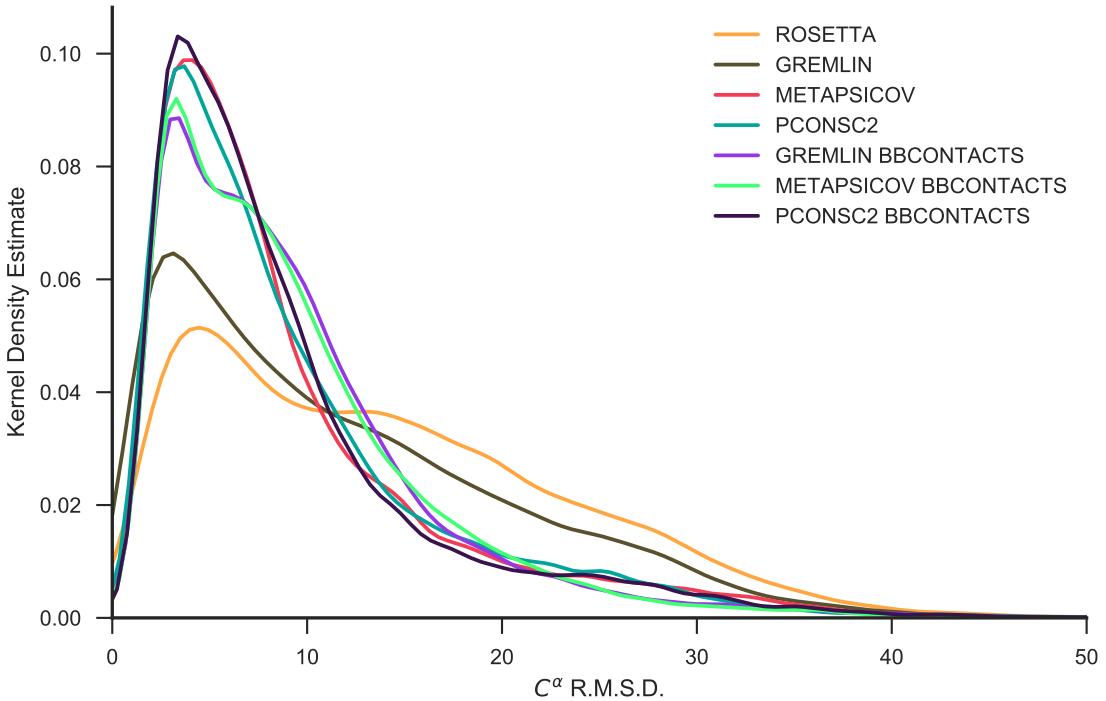


Figure 4.17: Kernel density estimate of $C\alpha$ interatomic RMSD for SPICKER clusters.

The structure solution through pipelines like AMPLE and other unconventional MR software [16, 17] can result from the placement of generated (ensemble) search models either in- or out-of-sequence register. The RIO metric [18] can reliably assess the register placement, and thus was used to analyse the MR placements of all search models of the seven targets with structure solutions from one or more decoy sets. The RIO scores for the hypothetical protein AQ_1354 (PDB: 1oz9) strongly support the high quality decoys used as input across all seven contact conditions (Fig. 4.18). Most search models are placed in-register and hardly any search models with out-of-register RIO scores failed either. In contrast, the search models of N-(5-phosphoribosyl)anthranilate isomerase (PDB: 4aaaj) — derived from high quality decoys in most conditions — shows a low percentage of AMPLE search models with RIO scores leading to structure solution (Fig. 4.18). Furthermore, the RIO scores normalized by the target chain length indicate that search models, independent of MR structure solution, were relatively small only exceeding 20% of the total target sequence in a few cases.

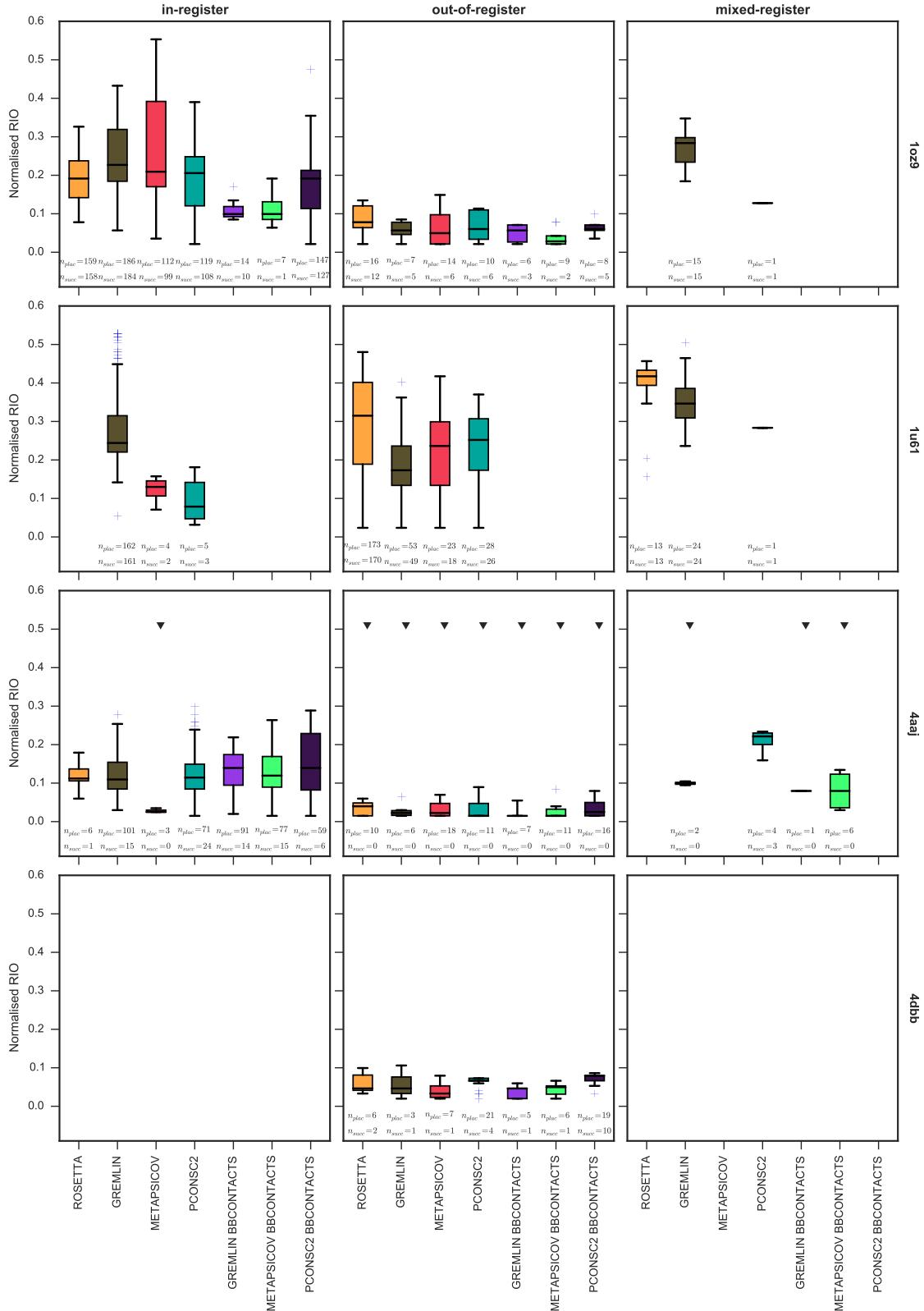


Figure 4.18: Normalised RIO score analysis of four successful targets in the MR dataset. Black triangles indicate AMPLEx search model sets without a structure solution.

One interesting target in this set with respect to the sequence register of the AMPLEx

search models leading to structure solution is putative ribonuclease III (PDB: 1u61). Although decoys from all contact conditions readily solved this target with at least 20 or more AMPLE search models, one interesting aspect arises from the RIO register analysis. Only GREMLIN decoys are primarily placed in-register (Fig. 4.18). AMPLE search models derived from the other three contact conditions, and in particular those from ROSETTA decoys, are primarily placed out-of-register with sequence coverage values of roughly 25%. In fact, a close analysis of the diversity of AMPLE search models highlights the accuracy of GREMLIN search models which represent a closely-matched substructure of the target protein (Fig. 4.19).

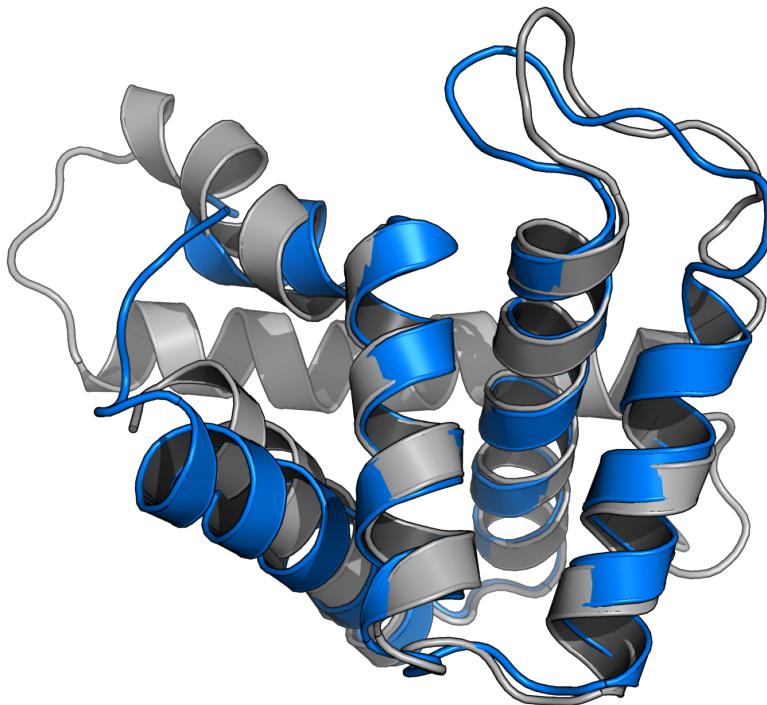


Figure 4.19: Successful search model (blue cartoon) post-PHASER placement superposed with the native structure (gray cartoon) for putative ribonuclease III (PDB: 1u61).

Compared to all other targets with structure solutions in at least one condition, the PTB domain of Mint1 (PDB: 4dbb) produced interesting yet somewhat surprising results. None of the search models, independent of their decoy source, achieved correct placement with any residue being in register. All structure solutions were obtained from out-of-register search model placements (Fig. 4.18). A visual inspection of all successful search models revealed that structure solutions were exclusively obtained with ide-

alised fragments. ROSETTA, GREMLIN and METAPSICOV decoys resulted in one or more single-helix ensemble search models that led to structure solution (Fig. 4.20). More interestingly though, PCONSC2, GREMLIN BBCONTACTS, METAPSICOV BBCONTACTS and PCONSC2 BBCONTACTS decoys yielded one or more two-strand β -sheets which, after successful MR, yielded fully built structures (Fig. 4.20).

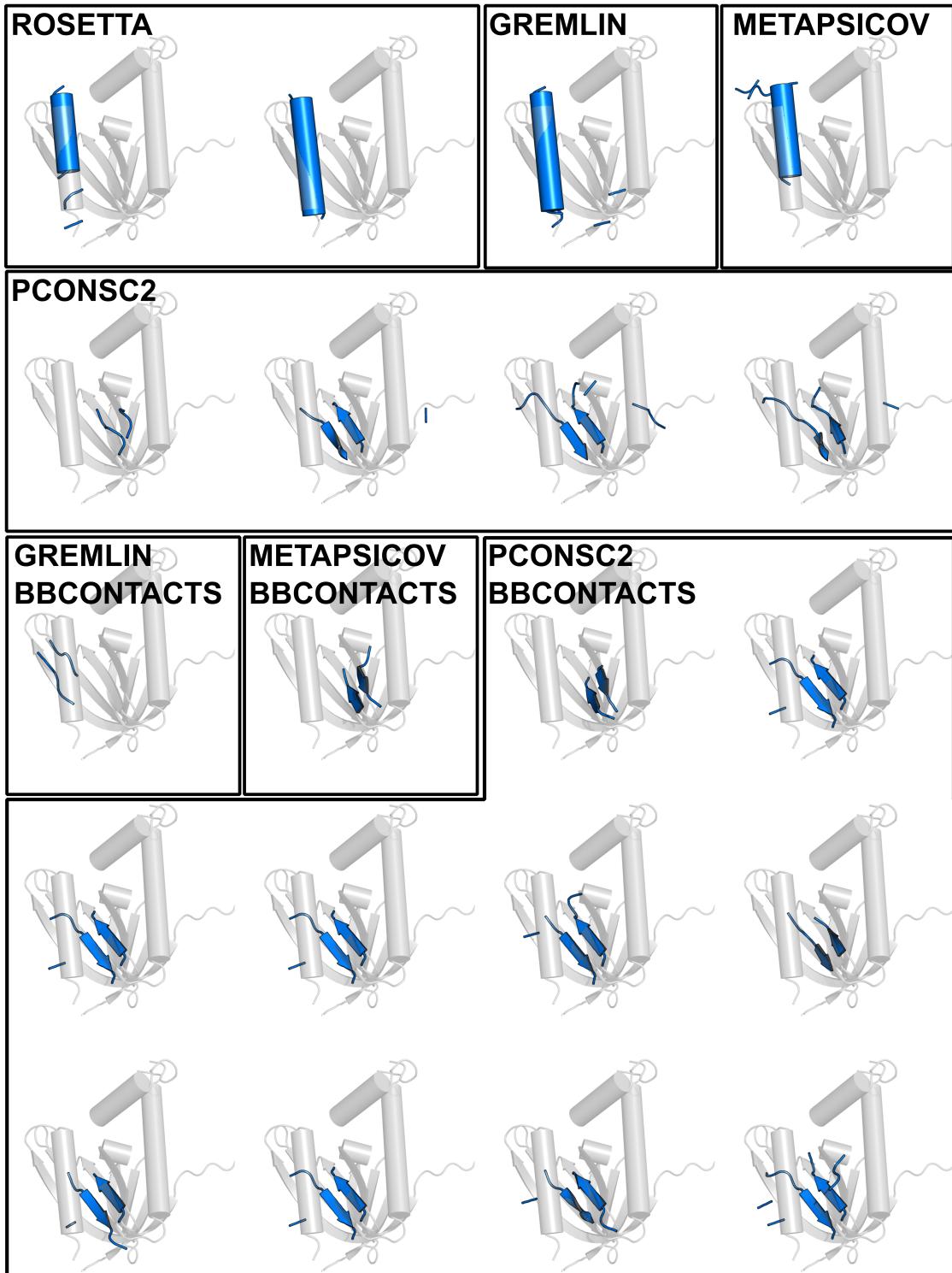


Figure 4.20: Successful search models post-PHASER placement (blue) superposed to the reference crystal structure (grey) for PTB domain of Mint1 (PDB: 4dbb).

Lastly, three targets were solved with one or two decoy sets alone. The structures of the retinoic acid nuclear receptor HRAR (PDB: 1fcy) and the peptide methionine

sulfoxide reductase (PDB: 1fgv) were only solved with a handful of AMPLE search models. Often singleton solutions like these are achieved through AMPLE's cluster-and-truncate procedure producing a single, idealised helix as search model. Here, we confirm such findings for target 1fcy, whereby single out-of-register helices derived from ROSETTA and GREMLIN decoys achieved structure solutions. However, the singleton search model derived from the GREMLIN BBCONTACTS decoys for the peptide methionine sulfoxide reductase (PDB: 1fgv) was placed in-register. A closer inspection of this AMPLE ensemble search model highlights a great success of the approach of adding BBCONTACTS distance restraints to separately predicted contact maps. In this instance, the successful AMPLE ensemble search model has 77% of its 49 residues placed in-register. More importantly, the search model is made up of two β -strands packing against each other, which was supported by BBCONTACTS predictions (Fig. 4.21). The last case, glycosylase domain of MBD4 (PDB: 4e9e), solved solely with GREMLIN decoys yielding 71 structure solutions. All successful AMPLE search models derived from the GREMLIN decoys were placed in-register.

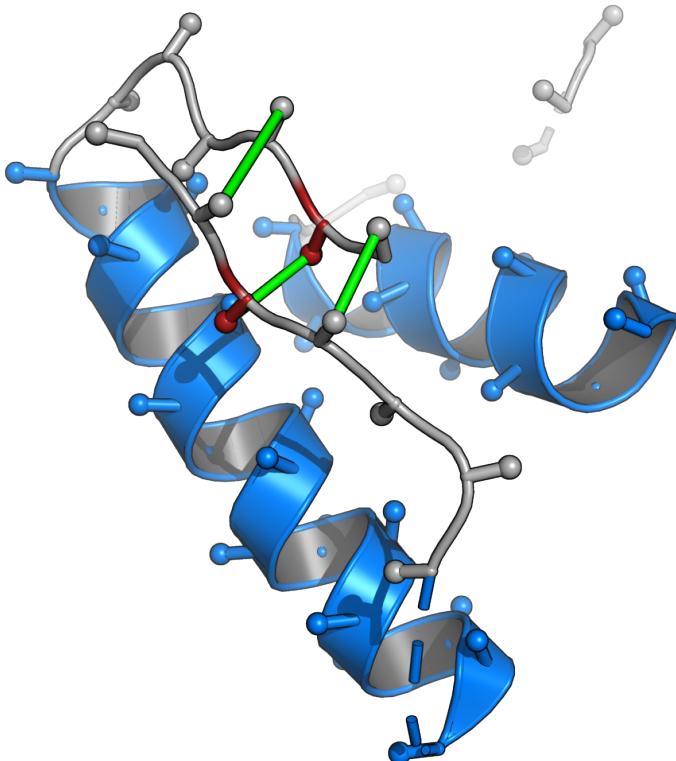


Figure 4.21: Successful search model post-PHASER placement for peptide methionine sulfoxide reductase (PDB: 1fgv). BBCONTACTS distance restraints are represented as green lines, α -helices in blue and β -strands in red. Secondary structure assignment calculated with STRIDE [19].

4.4 Discussion

This study was designed to explore the state-of-the-art metapredictor pipelines for residue-residue contact prediction. The main focus of this work was to distinguish differences in three key parts: raw contact predictions, their use in *ab initio* structure prediction and finally the effects on unconventional MR using AMPLE.

Key findings in this study revealed METAPSICOV and PCONSC2 metapredictors to yield the most precise contact predictions regardless of target fold or size. These results are in line with previous findings, which independently confirmed METAPSICOV contact predictions to yield the highest precision across numerous prediction algorithms [20, 21]. However, work in this study cannot confirm their findings, which demonstrate more precise contact predictions for all- β and mixed α - β protein targets compared to all- α ones. Several reasons might give insights into this discrepancy: (1) a much smaller sample size was trialled in this study (Wuyun et al. [20]: 680; De Oliveira et al. [21]: 3500); (2) the targets were chosen to deliberately sample various alignment depths including relatively low Neff (< 200) values; (3) only final contact predictions were analysed as part of this work, thus benefiting from post-prediction consensus finding and contact map processing through unsupervised machine-learning algorithms.

Furthermore, we demonstrated in this study that two similar ROSETTA energy functions yield different structure prediction results. The FADE function on average achieves more accurate structure predictions compared to the SIGMOID one. This result seems striking at first; however, a closer inspection of each of the energy function parameters gives possible insights into the reasons for the different outcomes. The FADE energy function defines both a maximum and minimum distance. The FADE energy function also does not consider amino acid-specific distances while the SIGMOID function does [1]. Furthermore, a custom weight factor is added for SIGMOID restraints to balance the restraint term in the overall energy term of each decoy (Sergey Ovchinnikov, personal communication). Thus, small changes in each of those definitions could have significant effects on the final structure prediction. Unfortunately, it is out of the scope of this study to explore all variations, and thus results aid primarily as guide for future work and AMPLE users. This study highlighted again the benefits of adding BBCONTACTS predictions to

existing contact maps to further restrain β -rich regions during structure prediction.

Lastly, part of the comparison carried out in this study was aimed specifically at macromolecular crystallographers and, in particular, AMPLE users. Beyond the proof-of-principle study described in Chapter 3, this work further illustrates how important additional restraint information can be to increase the chances of unconventional MR success. However, this work also highlighted limitations in the AMPLE routine whereby decoys that were restrained by residue-residue contacts achieved much higher decoy quality compared to unrestrained ROSETTA decoys, yet solved fewer targets. The idea that restrained decoys might benefit from a different kind of processing was further supported by the most successful decoy sets, which were obtained with GREMLIN contact predictions. Given that GREMLIN and ROSETTA decoys achieved similar decoy qualities for a large set, their structure solutions were identical for all of ROSETTAs successful solutions. GREMLIN decoys outperformed ROSETTA decoys solely on the basis that it acquired highly accurate decoys for one further target, and thus achieved the most structure solutions in this study.

Therefore, further work is required to identify the optimal strategy for decoy sets with high structural similarities to the native fold. Such work could focus on the recent idea of selecting decoys based on their long-range contact precision [21, 22] to specifically eliminate the worst decoys, and thus enhance a more fine-grained clustering approach in SPICKER. Alternatively, truncation could be guided by alternative means, such as the importance of each residue in the predicted contact map. Ultimately, it is key to improve the AMPLE protocol to exploit the much higher decoy quality to enhance the users chance of success.

Chapter 5

Alternative *ab initio* structure prediction algorithms for AMPLE

Chapter 6

Decoy subselection using contact information to enhance MR search model creation

6.1 Introduction

Work presented in previous chapters highlighted the much improved *ab initio* decoy quality achievable by restraining the conformational search space with residue-residue contact information. Furthermore, the data also highlighted that this improvement extends AMPLE’s tractability of achieving structure solution for more challenging targets. However, the data also indicated that AMPLE’s current protocol is not tailored towards decoy sets with overall much higher accuracy. Decoy sets with correctly predicted folds (average TM-score > 0.5 per 1,000 decoys) did not generate any or many ensemble search models leading to MR structure solution. It also became apparent that certain decoy sets contained few high-quality decoys that were lost in the process of clustering, since none of the top-10 SPICKER clusters contained that fold.

Beyond the limitations observed in AMPLE, *ab initio* decoy similarity in exceptional cases approaches a near-identical fold ($\text{RMSD} < 1.5\text{\AA}$) to the crystallised one. Although challenging by current means to identify these decoys, it is of great interest since these decoys might be sufficient by themselves as MR search models. Contact information, which was used to restrain the folding protocol, might provide enough information to drive such filtering. Indeed, De Oliveira et al. [21] found that long-range residue-residue contact pair satisfaction correlates well with decoy quality. Additionally, Adhikari and Cheng [12] use long-range contact satisfaction routinely in CONFOLD2 to exclude the worst decoys amongst the set predicted ones.

Thus, this chapter focuses on exploring alternative strategies of decoy selection in AMPLE, and if contact information can be used beyond the distance-restraint application in *ab initio* protein structure prediction.

6.2 Materials & Methods

6.2.1 Target selection

The dataset for this study consisted of 113 ROSETTA decoy sets generated throughout the works outlined in previous chapters. The 113 decoy sets covered all targets in

the ORIGINAL (Table A.1), PREDICTORS (Table A.2) and TRANSMEMBRANE (Table A.3) datasets. Top- L (> 5 residues sequence separation) CCMPRED [23], PCONSC2 [2], METAPSICOV STAGE 1 [3] and MEMBRAIN [24] contact pairs were used in combination with the *FADE* energy function to restrain the *ab initio* structure prediction process.

6.2.2 Computation of range-specific satisfaction scores

The satisfaction of short- (> 6 residues sequence separation), medium- (> 12 residues sequence separation) and long-range contact pairs (> 23 residues sequence separation; see ??) were computed for each decoy in each set. Hereby, the contact pairs of the original set of contact pairs used to restrain the *ab initio* structure prediction protocol were extracted, matched against the contact pairs extracted from individual decoys and the contact pair range-specific satisfaction score evaluated.

6.2.3 Decoy subselection

Each set of decoys was then ranked in descending order by their long-range contact pair satisfaction scores and the n decoys with the lowest scores removed from each set. The number of decoys to remove n were selected using a number of different strategies:

- *NONE*: leave the original set unchanged
- *LINEAR*: remove the worst 500 decoys
- *CUTOFF*: remove all decoys with a score of < 0.287
- *SCALED*: remove all decoys with a scaled score of < 0.5 , where the scaled score is score divided by set average

The fixed definition in the *CUTOFF* strategy was determined by De Oliveira et al. [21]. The scaled score used by the *SCALED* strategy was computed by dividing each decoy's long-range contact pair satisfaction by the set's average.

6.2.4 Molecular Replacement

To evaluate the benefits of such subselection to MR in AMPLE, a subset of 35 decoy sets (spanning 35 unique targets) were processed as described above and subjected to AMPLE v1.2.0 and CCP4 v7.0.28. Default options were chosen with few exceptions: decoys in all 10 clusters were used, subcluster radii thresholds were set to 1 and 3 Å, and side-chain treatments were set to `polyala` only. This change in protocol was shown to be advantageous in most cases by Jens Thomas (PhD Thesis), and thus trialled in this context.

To allow comparability of these results to previous AMPLE runs, an additional condition was added, namely *NONE_classic*. The decoy set from the *NONE* strategy was hereby subjected to the AMPLE protocol with default settings.

Each MR run was assessed using the criteria defined in ??.

6.3 Results

This chapter focuses on identifying further uses of predicted residue-residue contact pairs in unconventional MR. In particular, the exclusion of *ab initio* decoys by their contact satisfaction scores is under investigation. A total of 113 decoy datasets were used to identify potential means of identifying the best or worst decoys. Furthermore, three strategies were trialled alongside two standards to test the viability of excluding the worst decoys in ensemble search model preparation in AMPLE.

6.3.1 Contact pair satisfaction correlates with decoy quality

Kosciolek and Jones [25] previously identified a correlation between the TM-score of a decoy and its fraction of satisfied contact pairs. Albeit the striking positive correlations (short-range: $\rho = 0.50$; medium-range: $\rho = 0.57$; long-range: $\rho = 0.87$) for top-1 decoys, the study by Kosciolek and Jones [25] was limited to 10 representative targets with a maximum chain length of 158 residues. Furthermore, FRAGFOLD [26] was used for *ab initio* protein structure prediction, a method with inferior performance to ROSETTA [13] when using

the decoys in unconventional MR (see Chapter 5). Thus, the more diverse set of decoys generated in this study might be more representative in determining a correlation between decoy TM-scores and contact pair satisfaction.

A correlation analysis with 35 ROSETTA decoy sets representing 35 globular and transmembrane targets shows a positive linear correlations between a decoy's TM-score and short-, medium- and long-range contact satisfaction (Table 6.1). Furthermore, separating the correlation analysis of all targets by fold classification reveals that all- α , mixed α - β and transmembrane protein targets show the strongest positive correlations for long-range contact satisfaction (Table 6.1). All- β and mixed α - β decoy sets show the strongest correlations for short- and medium-range contact satisfaction, whereby the former shows a stronger positive correlation between the decoy's TM-score and its medium-range contact satisfaction than its long-range contact satisfaction (medium-range: $\rho = 0.54$; long-range: $\rho = 0.50$) (Table 6.1). Notably, the decoys of transmembrane protein targets show no correlation between TM-score and short-range contact satisfaction ($\rho = 0.08$; Table 6.1).

Table 6.1: Pearson's Correlation Coefficient (CC) analysis between a ROSETTA decoy's TM-score and short-, medium- and long-range contact satisfaction. Probability values for all p coefficients is < 0.01 .

Target class	Pearson's CC		
	Short-range	Medium-range	Long-range
all	0.11	0.18	0.64
all- α	0.30	0.44	0.69
all- β	0.40	0.54	0.50
mixed α - β	0.42	0.55	0.69
transmembrane	0.08	0.48	0.70

Following on from the correlation analysis, a linear regression model was fitted to individual subsets of the data used for the correlation analysis to see if a decoy's TM-score could be predicted from its contact satisfaction score. However, weak coefficients of determination indicate that only some cases show models with reasonably good fits to the data (Fig. 6.1). Nevertheless, all models further support the positive linear correlations between a decoy's TM-score and its range-dependent contact satisfaction. Interestingly, the strongest and best fits of the linear regression model to its corresponding data is for

long-range contact pairs, where the linear regression models are also near identical between the different fold categories (Fig. 6.1).

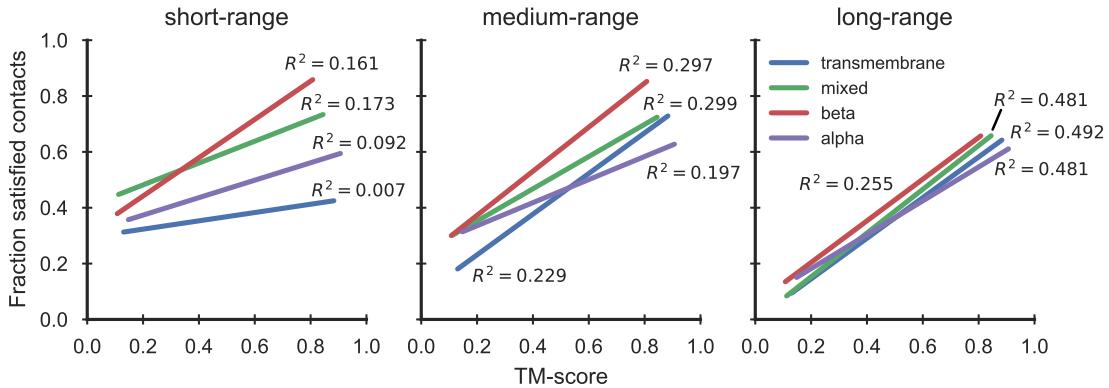


Figure 6.1: Linear regression model fitted to decoy TM-scores and corresponding fractions of satisfied, range-dependent contacts. Targets were further separated by fold classification. Coefficients of determination (R^2 -values) added alongside each regression model.

An analysis of the correlation between the TM-score and long-range contact satisfaction of individual decoy sets further highlights the potential to subselect decoy sets by their long-range contact satisfaction. Thirty decoy sets show statistically significant positive correlations between decoy TM-scores and their long-range contact satisfaction (ρ -values in range of 0.09 to 0.97 with p -value < 0.01). A singleton ROSETTA decoy set, derived for the Glycolipid transfer protein with PDB ID 2eum and restrained with METAPSICOV STAGE 1 contact data, shows a weak negative correlation ($\rho = -0.10$, $p < 0.01$). The remaining four decoy sets, derived for targets with PDB IDs 1chd, 1gm4, 2x6u and 3ouf and restrained with METAPSICOV STAGE 1 contact data except for 2x6u (PCONSC2), show no statistically significant correlation between the TM-score and long-range contact satisfaction of the decoy sets.

A further subdivide of the previously presented data by metapredictor highlights that no predictor outperforms the others. Decoy sets from all metapredictors result in decoy sets with stronger and weaker correlations. Similarly, target chain length and fold do not show overall stronger or weaker correlations.

So far, all analyses focused on entire sets of decoys (1,000 decoys per set); however, it is often desirable to know if we could better estimate the accuracy of the best decoy by some measure. Kosciolak and Jones [25] demonstrated strong positive correlations for

short-, medium- and long-range contact satisfaction with a decoy’s corresponding TM-score. In this work, these findings are confirmed albeit the strength of the correlation for long-range contact satisfaction is much weaker than observed previously (short-range: no correlation; medium-range: $\rho = 0.52$; long-range: $\rho = 0.69$) (Fig. 6.2). The weak positive correlation for short-range contact satisfaction is statistically non-significant, and thus cannot be validated.

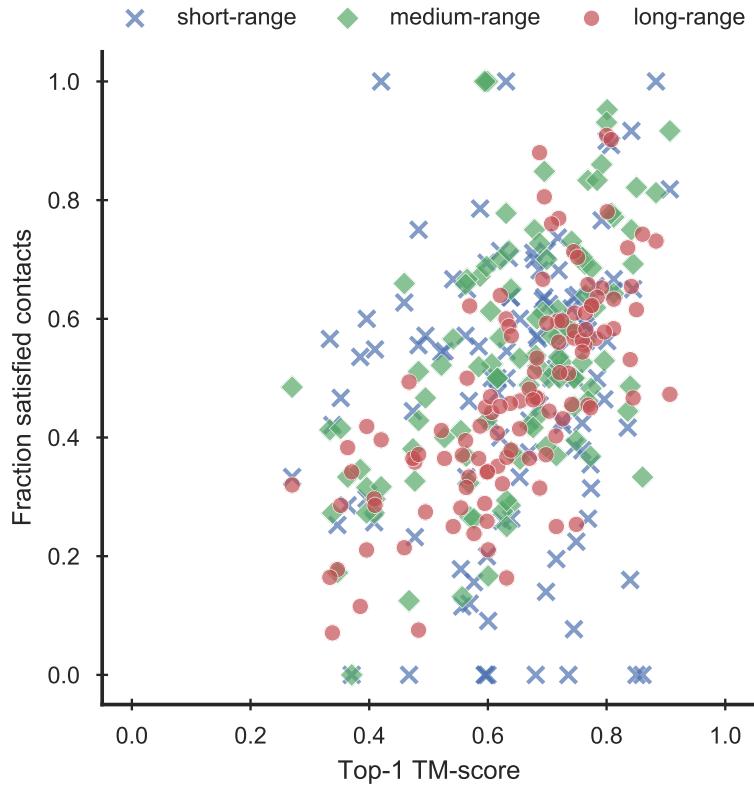


Figure 6.2: Analysis of the relationship between TM-score and contact satisfaction for the top-1 decoy (as ranked by TM-score) in each decoy set.

6.3.2 Long-range contact satisfaction metric to filter decoy sets

In the previous section, the data highlighted that decoy quality correlates positively with contact satisfaction. In particular, a strong positive correlation between long-range contact satisfaction and decoy quality could be established for almost all decoy sets in this study. A key ambition in this work is to determine, if this correlation could be used to include or exclude decoys prior to submission to the AMPLE processing pipeline to enrich the quality and number of ensemble search models trialled in MR.

The difference in mean TM-score of each decoy set before and after applying a subselection strategy (see Section 6.2.3) is shown in Fig. 6.3. Estimating a decoy’s quality by short-range contact satisfaction results in marginal mean TM-score changes of decoy sets ($\Delta_{CUTOFF} = -0.003$; $\Delta_{LINEAR} = 0.008$; $\Delta_{CUTOFF} = 0.001$). In comparison, medium- ($\Delta_{CUTOFF} = 0.005$; $\Delta_{LINEAR} = 0.015$; $\Delta_{CUTOFF} = 0.002$) and long-range ($\Delta_{CUTOFF} = 0.025$; $\Delta_{LINEAR} = 0.032$; $\Delta_{CUTOFF} = 0.005$) contact satisfaction are better estimates to improve the mean TM-score values of each decoy set. Notably, per-decoy long-range contact satisfaction provides the best estimate for identifying and excluding the least accurate decoys independent of the subselection strategy.

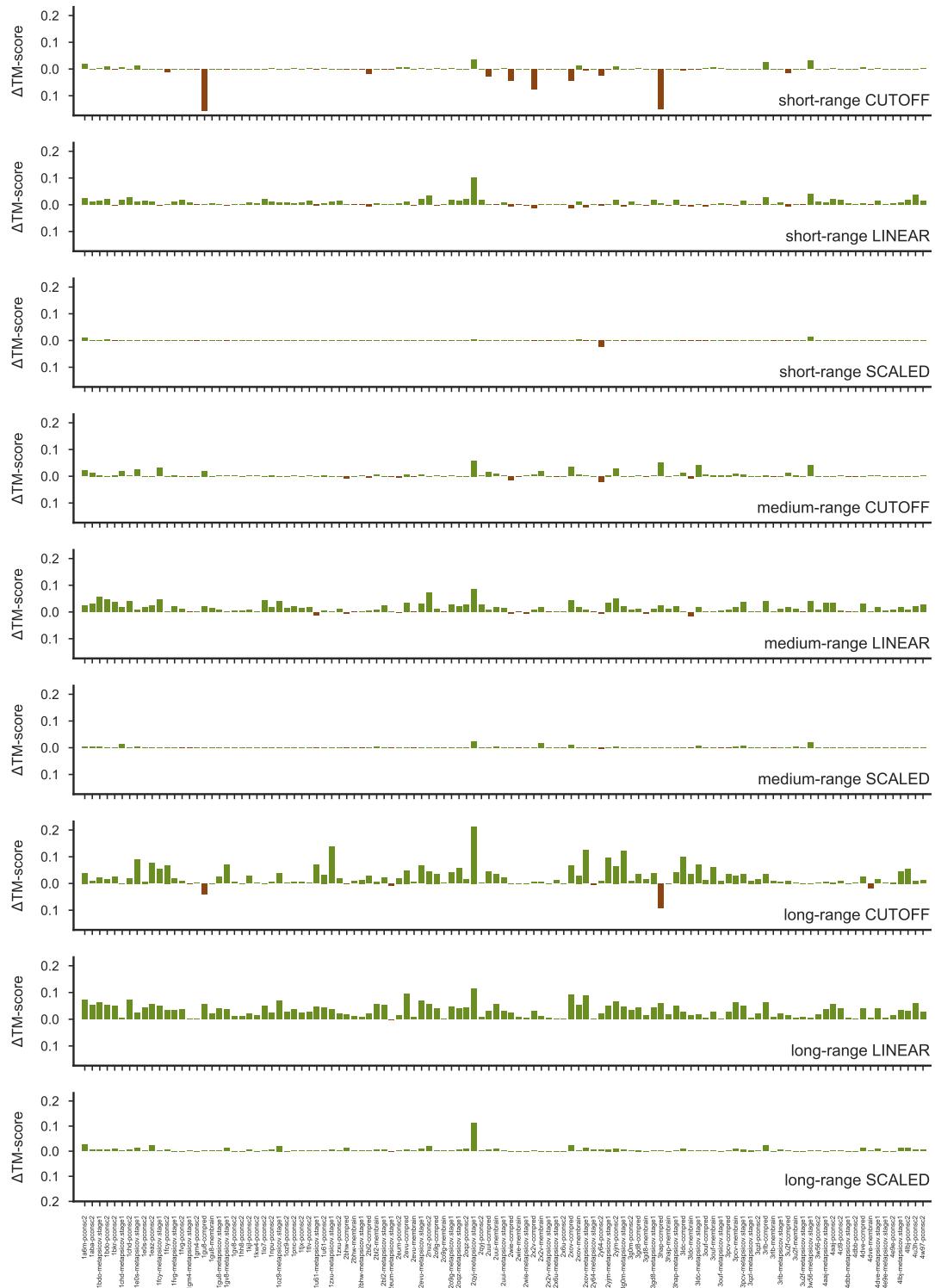


Figure 6.3: Differences in mean TM-score for decoy sets pre- and post-decoy subselection. Each subselection strategy is stated in each subplot along with the contact range used to establish decoy inclusion in the final set.

A comparison of Δ TM-scores between subselection strategies shows that the *LINEAR* strategy results in the largest improvements across all contact range categories. Nevertheless,

less, such great changes are expected compared to other strategies since the *LINEAR* strategy removes on average the most decoys from each set (*CUTOFF*=280; *LINEAR*=500; *SCALED*=66). However, the sample-dependent strategies (*CUTOFF* and *SCALED*) may remove a much greater number of decoys from a set if the corresponding satisfaction scores fall below a certain threshold.

In certain cases, some subselection strategies greatly changed the overall size and quality of the starting decoy set. The METAPSICOV STAGE 1 decoy set of the ankyrin sequence (PDB ID: 2qyj) shows overall quality improvements from 0.006 (short-range *SCALED*; $n_{models} = 958$) to 0.213 (long-range *CUTOFF*; $n_{models} = 218$). The CCMPRED decoy set of sensory rhodopsin II sequence (PDB ID: 1gu8) shows overall changes from -0.155 (short-range *CUTOFF*; $n_{models} = 2$) to 0.06 (long-range *LINEAR*; $n_{models} = 500$).

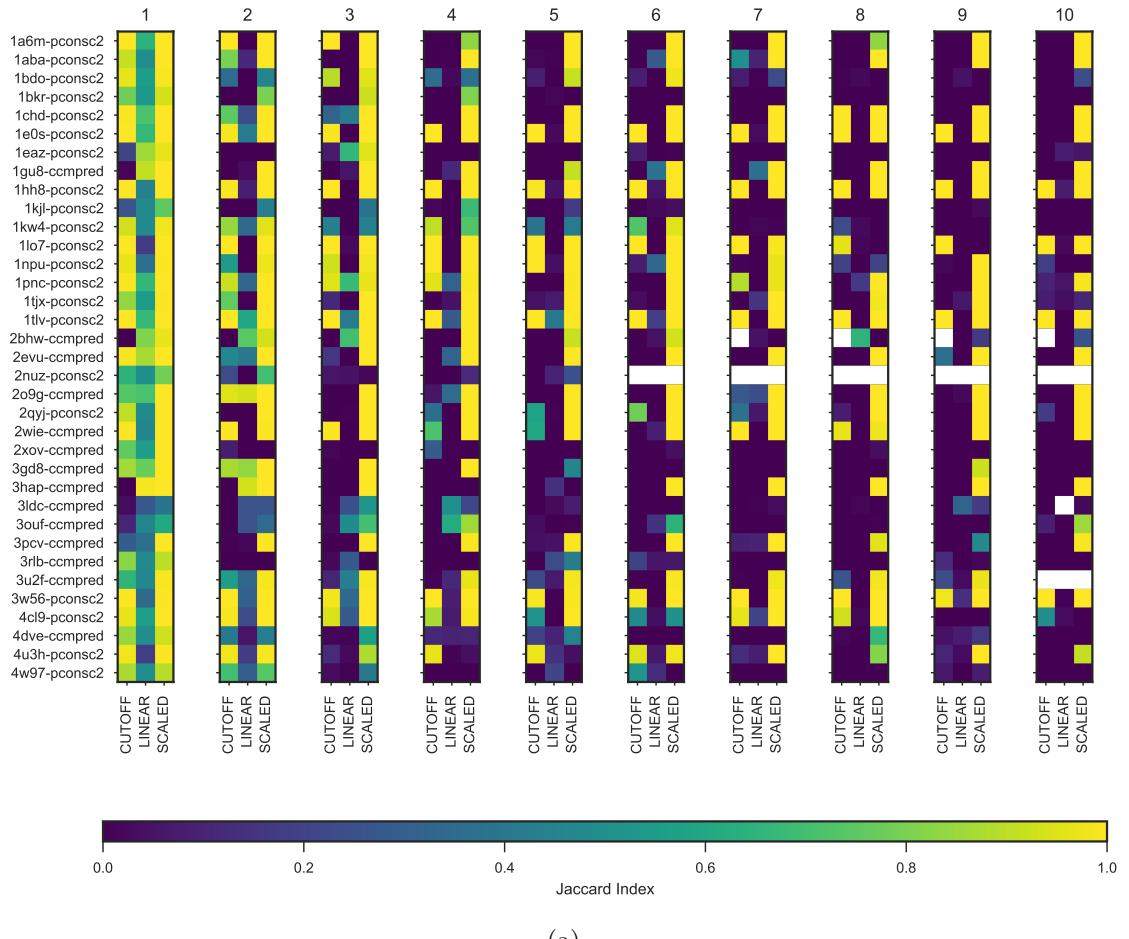
Overall, the optimal strategy to select or exclude decoys from a starting set of structures appears to be long-range contact satisfaction given the improved similarity metric of the set compared to the reference crystal structure.

6.3.3 AMPLE’s cluster-and-truncate approach with filtered decoy sets

The original dataset was limited to a sub-sample of 35 decoy sets spanning 35 unique targets (21 globular and 14 transmembrane targets) for MR trials. The contact prediction algorithms generating the restraints for the *ab initio* structure predictions were PCONSC2 (globular targets) and CCMPRED (transmembrane targets). Each decoy set was subjected to the AMPLE pipeline with certain decoys removed according to one of five subselection strategies.

The initial step in the AMPLE pipeline is the clustering of decoys. A comparison of SPICKER clusters between the *NONE* default strategy and the *CUTOFF*, *LINEAR* and *SCALED* subselection strategies highlights an important difference. Larger clusters — those ranked higher — show higher similarity between a subselection strategy and the default (Fig. 6.4a). The top SPICKER cluster shows high similarities between the *NONE* strategy and all other subselection ones, whereby it has to be noted that the *LINEAR*

strategy contains only 50% of the starting decoys and thus can at most show a Jaccard index of 0.5. With increasing cluster index, the overall similarity degrades and most of the decoys in cluster 10 are non-identical between each subselection strategy and the default. Furthermore, a similar analysis to compare the overall quality of each cluster compared to the target structure revealed less difference between the default and each subselection strategy for higher-order SPICKER clusters (Fig. 6.4b). With decreasing SPICKER cluster index, the difference in median TM-scores starts to alternate without any particular pattern. Thus, pre-selecting decoys prior to AMPLE’s cluster-and-truncate approach most certainly preserves the top cluster for the *CUTOFF* and *SCALED* subselection strategies, whereby lower clusters show more deviation from the default.



(a)

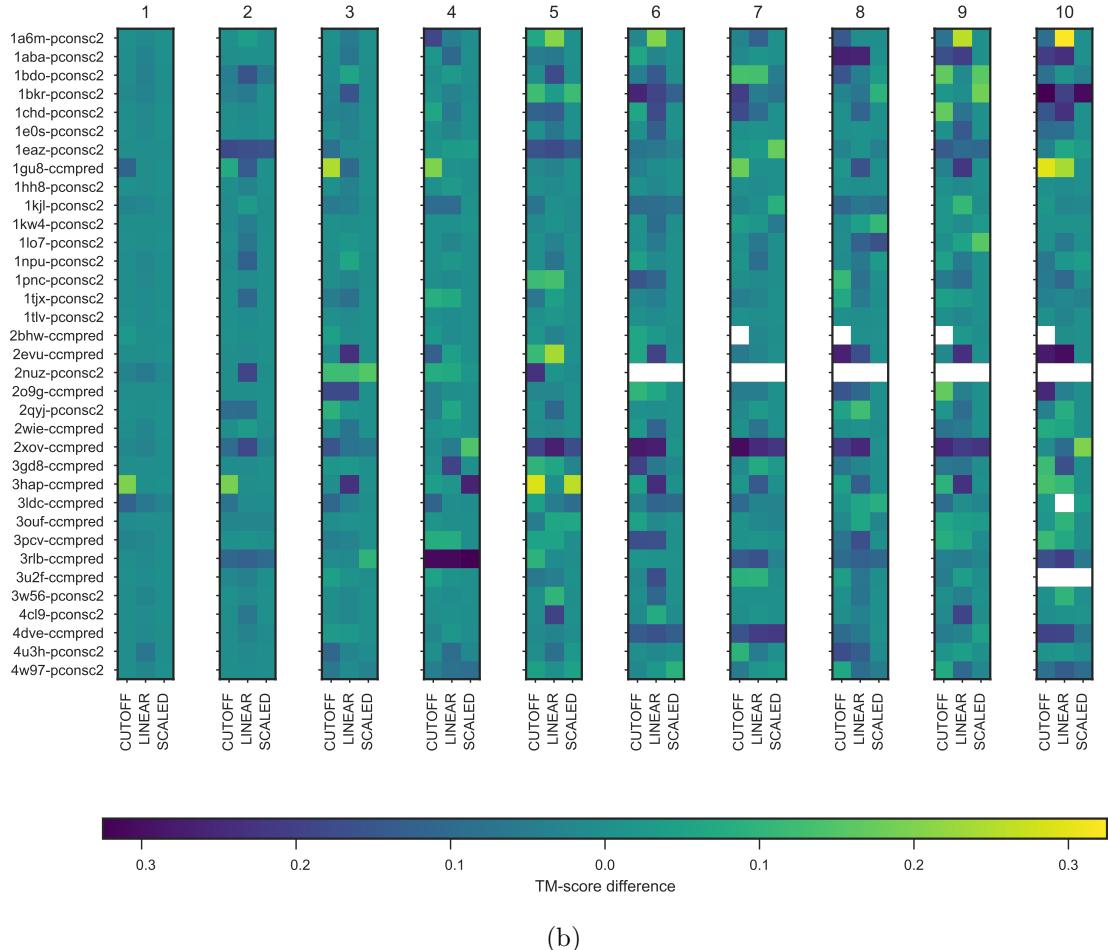


Figure 6.4: Effect of decoy subselection on SPICKER clusters. Effect illustrated by (a) the Jaccard Index and (b) median TM-score difference. Values were calculated for clusters resulting from the full starting set of decoys and the *CUTOFF*, *LINEAR* and *SCALED* subselection strategies.

The mean of the inter-decoy variance computed by THESEUS — used in AMPLE to guide truncation — is reduced in lower clusters compared to the *NONE* default strategy (Fig. 6.5). The clusters of decoys based on the galectin-3 domain (PDB ID: 1kjl) sequence show overall the highest reduction in mean inter-decoy variance up to -15\AA compared to the default strategy. Similarly, clusters 4 and 8 of the K^+ -channel protein domain (PDB ID: 3ouf) show reductions in mean inter-decoy variance of up to -20\AA . In general and also true for the aforementioned examples, clusters starting from *CUTOFF*-subselected decoys show mean inter-decoy variance reductions, followed by *LINEAR* and then *SCALED*-subselected decoys sets.

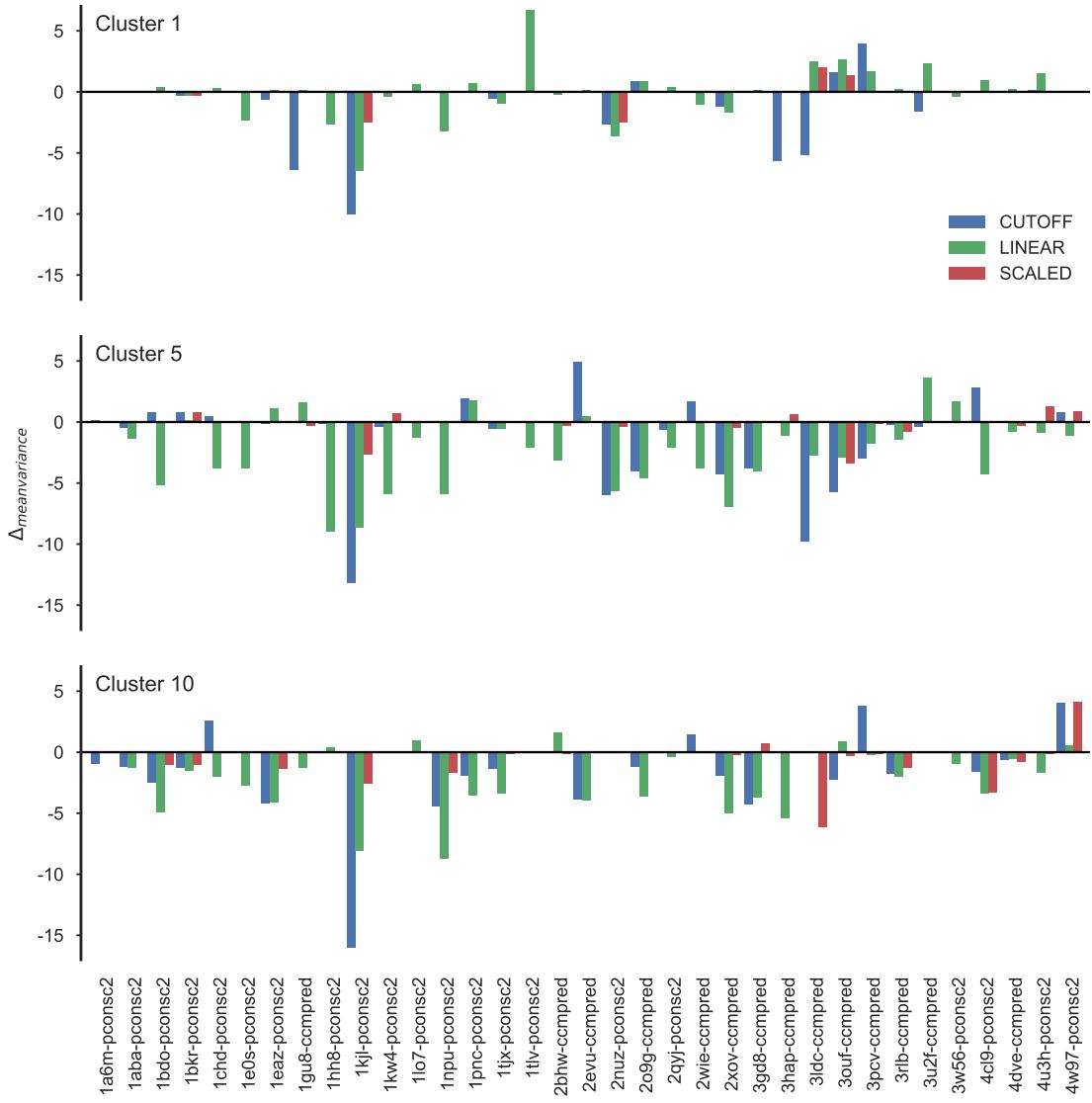


Figure 6.5: Effect of decoy subselection on mean inter-decoy THESEUS variance. Difference in mean variance calculated between the default and the three decoy subselection strategies *CUTOFF*, *LINEAR* and *SCALED*. Data for clusters 1, 5 and 10 shown as examples.

A comparison of intermediate stages in AMPLE pipeline resulting from differently subselected decoy sets is very difficult. Each strategy results in different starting sets, which result in different clusters. Since AMPLE's objective truncation procedure is based on the inter-decoy variance, it might be greatly affected by differing clusters. Nevertheless, structure solution is more likely when AMPLE generates more ensemble search models because a greater number of search models relates to greater inter-cluster decoy similarity and trialling a greater number should provide a higher chance of success. A count of generated AMPLE ensemble search models reveals that the *SCALED* strategy generates

the most search models ($n = 7,611$), which is roughly 300 more than the default *NONE* strategy ($n = 7,340$). The *CUTOFF* subselection strategy generates the least ensemble search models ($n = 7,237$), whilst the *LINEAR* strategy's count ($n = 7,401$) is very similar to the *NONE* one.

Further inspection of the number of AMPLE-generated ensemble search models by target reveal near identical numbers between the *NONE*, *LINEAR* and *SCALED* strategies (Fig. 6.6). In fact, only few outliers for each of those methods distinguish them from the others. The *CUTOFF* strategy shows greater deviation from the other three, especially for certain targets with differences up to approximately 200 ensemble search models (Fig. 6.6). If we compare all these strategies to the previous default processing in AMPLE (*NONE_classic*; further details in Section 6.2.4), we can see that the number of search models is greatly reduced (Fig. 6.6). A comparison of the previous default (*NONE_classic*) with the new one (*NONE*) shows on average 144 fewer ensemble search models, whilst sampling a larger range of folds through all 10 clusters.

6.3.4 Decoy subselection extends AMPLE's performance

The final step in this study is the assessment of AMPLE-generate ensemble search models in MR. In particular, the comparison of different decoy subselection strategies is of great interest, since it might allow us to extend AMPLE's performance beyond that described in previous chapters.

A comparison of the total number of targets solved by each subselection strategy shows that the *CUTOFF*-subselected decoys lead to most structure solutions (14 out of 35) whilst resulting in the fewest search model count (Fig. 6.6). Although slightly less successful, the *LINEAR* and *SCALED* subselection strategies lead to approximately 6% more solved targets than the current AMPLE default *NONE* strategy. The *LINEAR* and *SCALED* strategies are on par with AMPLE's previous default, the *NONE_classic* strategy (Fig. 6.6).

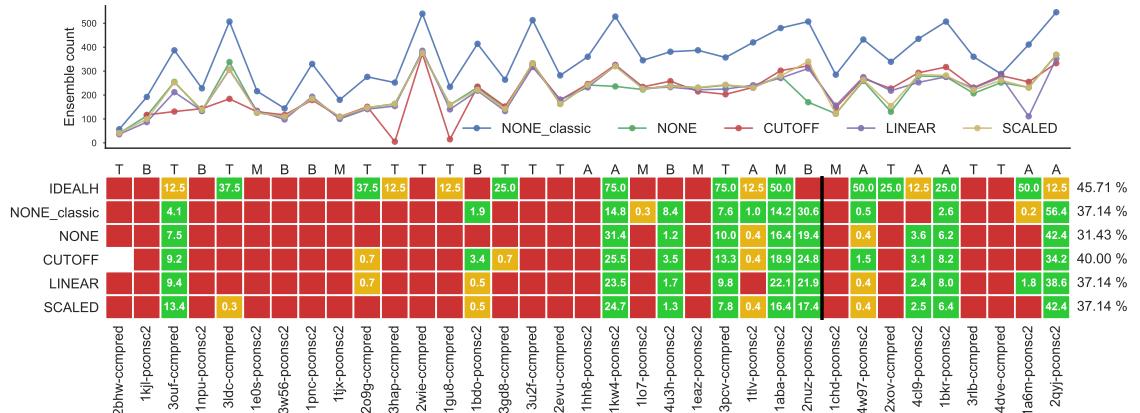


Figure 6.6: Molecular Replacement summary of decoy-subselected AMPLE ensembles. AMPLE-generated ensemble counts illustrated at the top with Molecular Replacement results in grid below: red cell equates to no solution; orange to a singleton solution; and green to multiple solutions. The number in the orange and green cells indicates the percentage of ensemble search models leading to structure solutions. One letter code above each column indicates the target fold: “T” for transmembrane; “A” for all- α ; “B” for all- β ; “M” for mixed α - β). Percentage alongside each row indicates the number of targets with solutions by each strategy. Targets are sorted from left to right with increasing median TM-score of the starting decoy set. The black line highlights the TM-score threshold of 0.5 for considering a fold native-like [15]. The subselection strategy *IDEALH* refers AMPLE’s ideal helix library.

Chapter 7

Protein fragments as search models in Molecular Replacement

7.1 Introduction

Ab initio structure prediction algorithms typically start with a coarse grained search of conformational space through the assembly of previously picked structural fragments. As such, the accuracy of structure prediction is heavily dependent on the similarity of fragments to the target fold for each position [27]. Thus, the necessary structural information for accurate structure prediction must be encoded in the fragment library for a given target sequence. This approach allows the modelling of new protein folds by considering them as assemblies of already known building blocks, such as super-secondary structure motifs [28]. Furthermore, fragments similar to those typically selected for *ab initio* structure prediction were successfully used in other areas of structural biology including Nuclear Magnetic Resonance (NMR) [29, 30] and X-ray crystallography [31] studies to elucidate unknown protein folds. Despite their modest success, almost all attempts neglected target-specific information generally available to structural biologists obtainable through bioinformatics software. This information includes the primary sequence of the target, torsion angle predictions, predicted solvent accessibility or co-evolution information. In theory, all additional information should improve the generation of such fragment libraries by aiding the selection process or cross-validating the identified fragments.

Over the last decade, efforts have been made to improve the precision of structural fragment libraries used in *ab initio* structure prediction [27, 32–38]. Various different algorithms have been developed to generate static and dynamic fragment libraries. Static fragment libraries are those pre-computed and generally consist of common super-secondary structure motifs. In comparison, dynamic fragment libraries consist of fragments of variable lengths acknowledging the fragment-dependent optimal length. Most commonly used in *ab initio* structure prediction are dynamic algorithms, such as FLIB [38], NNmake [27] or HHfrag [35]. Dynamic-library producing algorithms differ in their definition of ideal fragment lengths, the default number of fragments used per position and the way in which fragments are extracted. However, these algorithms typically share the same additional sequence-based information used to aid the selection of target fragments, which usually includes sequence similarity, three-state secondary structure prediction and torsion angle prediction.

Given that fragment libraries selected to perform *ab initio* structure prediction can contain high quality fragments or super-secondary structure motifs, those fragments must sometimes be suitable as MR search models. Correct identification of true positives should allow for dynamic fragment selection to achieve MR structure solution without the overhead of *ab initio* structure prediction. Furthermore, dynamic algorithms could pick fragments of varying lengths, possibly matching co-evolution data or other externally obtainable restraints to validate fragments prior to any MR attempt. As such, the work in this chapter focuses on exploring this idea using FLIB [38], a dynamic fragment picking algorithm considering co-evolution data to verify fragments during the picking process.

7.2 Materials & Methods

7.2.1 Target selection

Four targets were manually selected for this study. The crystallographic data needed a resolution of around 1.5Å with a single molecule in the asymmetric unit. The target chain length needed to be below 150 residues, and the fold of the protein structure to be either mixed α - β or all- β . A further target selection criterion was the availability of precise contact information for fragment selection.

The PDB identifiers of the selected targets are: 1aba, 1lo7, 1u06, and 5nfc. The former two are described in Table A.1. Target 1u06 is a recently published structure of α -spectrin SH3 domain (PDB ID: 1kjl in Table A.1) with a resolution of 1.49Å. Target 5nfc is a recently published structure of Galectin-3 (PDB ID: 1kjl in Table A.1) with a resolution of 1.59Å. This resulted in a dataset with similar attributes for each target: crystallographic data resolution of 1.5Å with a single molecule in the asymmetric unit, and the target chain length of < 150 residues. Each fold class, mixed α - β and all- β , contained two targets.

7.2.2 Fragment picking using FLIB

FLIB [38] requires four inputs: the predicted secondary structure, predicted torsion angles, residue-residue contact pair data and a copy of the PDB. The secondary structure for each

target was predicted using PSIPRED v4.0 [39] with default parameters. The torsion angles were predicted using SPIDER2 [40] with default parameters, and residue-residue contact pairs using METAPSICOV v1.04 [3] with default parameters. HHBLITS v2.0.16 [41] with database version `uniprot20_2016_02` was used by METAPSICOV to generate the MSA for contact prediction of each target sequence. BLASTP v2.2.31+ [42, 43] was used by PSIPRED with the UNIPROT database version `uniref90-2016_06`. The local copy of the PDB for fragment picking was downloaded on August 11, 2016.

Two modifications were made to the default FLIB v1.01 (<https://github.com/sauloho/FLIB-Coevo>, commit `abade3b`) protocol. The first focuses on exclusion of fragments with > 90% helical content (assigned by DSSP [19]). If fragments with > 90% helical content are allowed and residues are predicted to be part of an α -helix, fragment libraries tend to be overpopulated for these positions with short helices. This would generate fragment libraries similar to ideal helix libraries, which is not the purpose of this work. The second modification was to allow fragments with RMSD > 10.0 \AA to the reference structure to be considered. This modification to the FLIB algorithm was implemented for development purposes by the authors to validate the performance of the algorithm. However, to allow for the automatic calculation of RMSD value of each fragment without deliberately excluding less-similar fragments this modification was lifted.

Two-hundred fragments were picked per target sequence position. Top- L or $L/2$ contact pairs were selected from both METAPSICOV STAGE 1 and STAGE 2 predictions with a minimum sequence separation of either 6 or 12 residues. Helical fragments were either included or excluded. The fragment length ranged from either 6 or 12 (dependent on minimum sequence separation) to 63 residues. In all instances the `-coevo_only` flag was set to exclude fragments with starting residues undefined by any contact pair in the set¹. Overall, this generated 16 fragment libraries per target.

Each fragment library was then filtered to remove homologs of the target to be solved. BLASTP [42, 43] and HHPRED [44] searches were conducted to identify homologous PDB entries. The BLASTP search was performed identically to Oliveira et al. [38] using an E-

¹The `-coevo_only` flag was intended to select only fragments that satisfied at least one contact pair. This intended behaviour was not part of the source code throughout this study, and only detected post-analysis. The issue was reported to the developers and has since been fixed in the FLIB source code (commit "`b3eb01d`").

value cutoff of 0.05. The HHPRED search parameters were identical to the MPI-Toolkit [45] webserver version (<https://toolkit.tuebingen.mpg.de/>). Fragments derived from PDB entries identified by BLASTP and HHPRED (probability score of ≥ 20.0) were excluded from the fragment libraries.

All per-target fragments were then binned by their peptide lengths. Subsequently, they were ranked by FLIB scores and RMSD values, and the best fragment from each length-dependent bin selected. Partially redundant fragments of the same template structure consisting of the same region with varying flanking residues were kept, if they were ranked top for each fragment length group. Finally, the coordinates of the fragment backbone atoms were extracted to create poly-alanine search models.

Note, the FLIB score refers in this chapter to the predicted torsion angle score for a given fragment, which FLIB uses in its default routine to rank fragments with lower scores being more favourable [38].

7.2.3 Molecular Replacement in MRBUMP

The previously extracted fragments were subjected to the MR pipeline MRBUMP [46]. This uses PHASER [47] for MR, REFMAC5 [48] for refinement and SHELXE [49] for density modification and main-chain tracing. MRBUMP default parameters were used with exception of the PHASER RMSD estimate. Each fragment was subjected to MRBUMP using PHASER RMSD values of 0.1, 0.6 and 1.0 \AA .

7.2.4 Assessment of FLIB fragments

Fragment torsion angles — predicted by SPIDER2 [40] — were assessed using the Mean Absolute Error (MAE), which evaluates the average absolute difference between the predicted and experimentally determined angles [40]. To account for the periodicity of an angle, the smaller value of the absolute difference d_i and $360 - d_i$ was used. The coverage of a fragment library was assessed by the proportion of residues present in at least one fragment in the library. The precision of a fragment library was defined by the fraction of True Positive (TP) fragments. All fragments with an RMSD of $< 1.5\text{\AA}$ were considered

TP else False Positive (FP). The equation used to calculate the precision score is ??.

The RMSD value, as calculated by FLIB [38], was computed between the aligned residues of the corresponding crystal structure and the fragment. The number of satisfied contact pairs in each fragment was calculated by scoring the number of TP contact pairs by using a contact's residue indexes according to sequence alignment provided by FLIB. MR success for each search model was solely assessed by SHELXE scores, whereby a CC score of ≥ 25.0 combined with an Average Chain Length (ACL) score of ≥ 10.0 was required.

7.3 Results

In this study, the main objective was to determine if peptide fragments derived from protein structures in the PDB could be reliably selected and trialed in MR to achieve structure solutions. The fragment picking algorithm FLIB [38] was used to pick fragments given its novel approach of validating selected fragments against a set of predicted residue-residue contacts.

7.3.1 Precision of FLIB input data

The FLIB algorithm requires two sets of input data — the predicted secondary structure and per-residue torsion angles — for each target sequence alongside an optional third source of information in form of co-evolution data. The first part of the analysis in this study focuses on these data given that the FLIB fragment picking heavily relies on the individual features in the selection and scoring of each individual fragment [38]. Poor data at this stage could lead to poor fragments that would be unsuitable for MR trials given that high accuracy, i.e. a low RMSD value between the search model and target, is required.

The secondary structure prediction highlighted high precision between each target's prediction and the DSSP-assigned [19] secondary structure of the target reference structure (Fig. 7.1). The three targets with PDB identifiers 1aba, 1lo7 and 1u06 have secondary structure predictions with a precision of $> 89\%$. The fourth target, 5nfc, shows comparatively poor precision of 50.7% over all residues in the PSIPRED prediction and the DSSP

assignment using the reference crystal structure. However, 11 out of 13 secondary structure features are correctly predicted, suggesting successful fragment picking is possible.

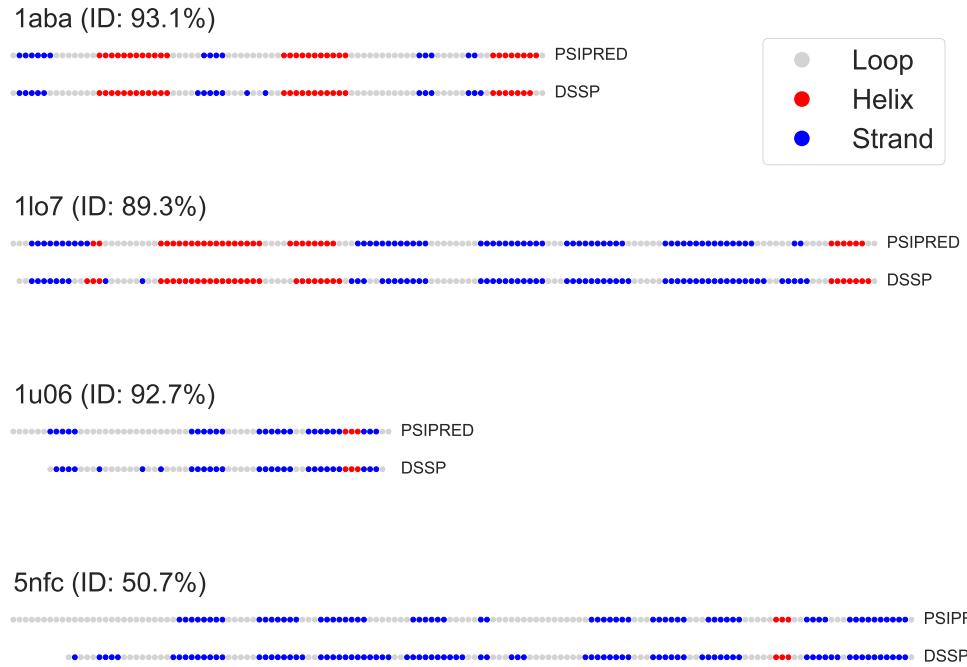


Figure 7.1: Schematic comparison of PSIPRED [39] secondary structure prediction and DSSP [19] assignment. Percentage identity is provided next to each identifier. The identity was computing using the Hamming distance over all positions present in the target sequence and reference structure.

The contact prediction data for METAPSICOV STAGE 1 and STAGE 2 predictions demonstrate the high precision scores achievable by this algorithm (Table 7.1). In this study, the top contact pairs at cutoffs L and $L/2$ were provided to the FLIB algorithm. All targets have precision scores for both sets of predictions at both cutoff levels of > 0.6 (Table 7.1). A comparison of the sets of contact pairs shows that only every third (for $L/2$ contacts) or every other (for L contacts) contact pair is shared between both METAPSICOV STAGE predictions highlighting the importance of trialling both when selecting FLIB fragments (Jaccard index in Table 7.1).

Table 7.1: Precision scores for METAPSICOV [3] STAGE 1 and STAGE 2 contact predictions. Jaccard index calculated for the same L -dependent selection of contact pairs between METAPSICOV STAGE 1 and STAGE 2 predictions.

Target	$L/2$ contact pairs			L contact pairs		
	Prec _{STAGE 1}	Prec _{STAGE 2}	Jaccard	Prec _{STAGE 1}	Prec _{STAGE 2}	Jaccard
1aba	0.884	0.884	0.303	0.713	0.759	0.513
1lo7	0.857	0.957	0.308	0.738	0.837	0.446
1u06	0.839	0.806	0.378	0.710	0.787	0.459
5nfc	0.822	0.836	0.327	0.619	0.762	0.434

Given the two METAPSICOV contact prediction files, both show localised clusters of contact pairs characteristic for secondary structure features (Fig. 7.2). These clusters are more populated with contact pairs in METAPSICOV STAGE 2 predictions. This behaviour is to-be-expected given that the second stage in METAPSICOV screens the first to remove singleton contact pairs whilst enriching the already existing clusters [3]. Besides the visual analysis, a cluster determination study on each of those contact maps further confirmed a higher singleton frequency in METAPSICOV STAGE 1 predictions. The latter contain on average 9% more singleton contact pairs, and thus a higher degree of noise.

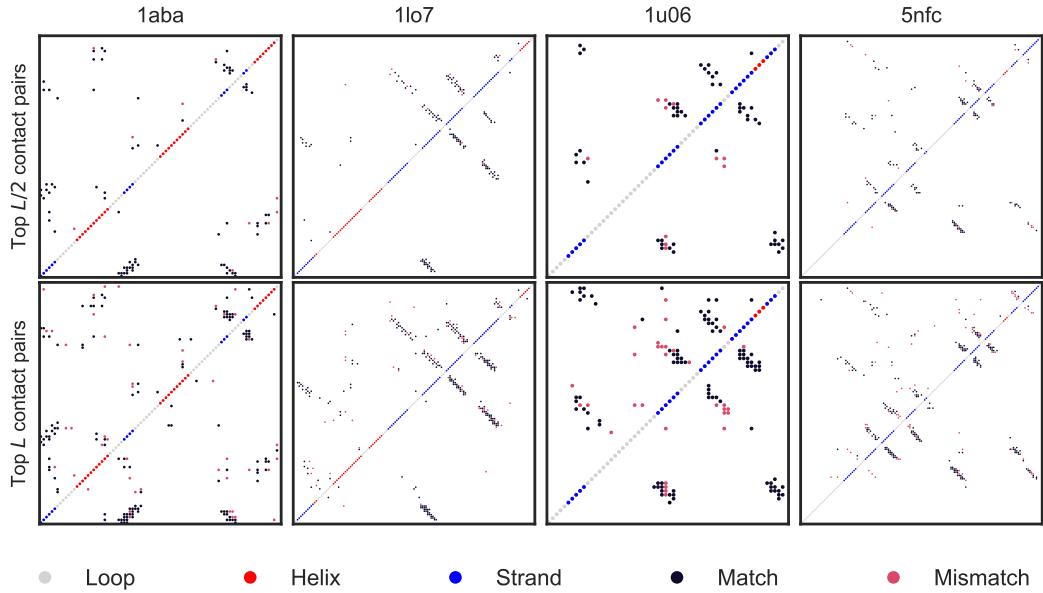


Figure 7.2: Comparison of $L/2$ and L correctly and incorrectly predicted contact pairs for four FLIB targets. Contacts were predicted using METAPSICOV [3] STAGE 1 (top left) and STAGE 2 (bottom right). True and false positive contact pairs were identified using a 8Å cutoff between Ca (C β in case of GLY) atoms of a reference crystal structure. PSIPRED [39] secondary structure prediction provided along the diagonal.

An analysis of the MAE of torsion angles between the SPIDER2 [40] prediction and a corresponding reference crystal structure highlights accurate predictions for three of four targets (Fig. 7.3). The largest MAE_ϕ across the four target sequences is 24.347° , and the largest MAE_ψ is 45.459° (MAE values for PDB entry 1u06). The smallest MAE_ϕ is 13.822° (PDB ID: 1aba) and smallest MAE_ψ is 17.273° (PDB ID: 1lo7). Segments in sequence space with regular secondary structure, as predicted by PSIPRED [39], result primarily in low MAE values of torsion angles. In contrast, unstructured regions highlight much larger MAE values indicating the difficulty of predicting these regions. Noticeably, the MAE_ψ appears to be much larger in those regions than the MAE_ϕ for the same residue.

In summary, all target sequences have FLIB input data of good quality, which should allow FLIB to select fragments of suitable accuracy for MR.

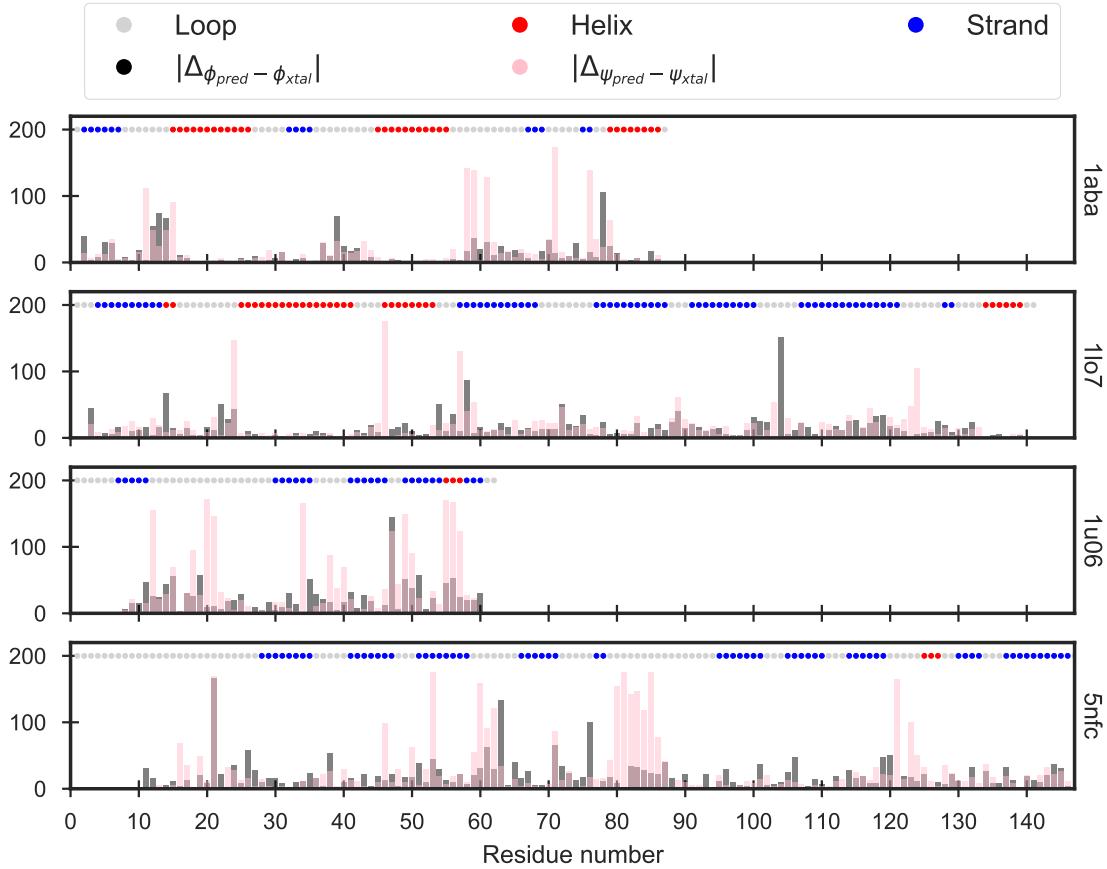


Figure 7.3: Comparison of MAE of torsion angles predicted by SPIDER2 and extracted from a corresponding PDB structure. PSIPRED [39] secondary structure prediction provided alongside the MAE values.

7.3.2 FLIB fragment picking

Sixteen FLIB fragment libraries were picked for each protein target in this study. Each fragment library consisted of one permutation of one of two contact prediction files and altering input parameters.

Across all four targets, the FLIB algorithm selected a total of 8,535,458 fragments (Table 7.2). The fragment libraries show similar statistics across the four protein targets despite the diversity in fold and chain lengths. The mean FLIB score is 3,200 score units with a mean RMSD of 9.00Å. Fragments for the alpha-spectrin SH3 domain (PDB ID: 1u06) scored the lowest mean FLIB score with 3,034 units; however, the same target scored the worst by mean RMSD with an average of 9.47Å. In contrast, fragments picked for the sequence of the bacteriophage T4 glutaredoxin (PDB ID: 1aba) achieved the best mean

RMSD of 7.85Å given the second highest mean FLIB score of 3,217 units (Table 7.2).

Table 7.2: Summary of fragment statistics for FLIB libraries selected for four protein targets. Count_H corresponds to the count of fragments extracted from homologs.

Target	Count	Count _H	FLIB score			RMSD		
			Median	Mean	Std Dev	Median	Mean	Std Dev
1aba	2,091,321	45,133	3,061	3,217	1,405	7.70	7.85	3.81
1lo7	2,497,813	23,396	3,187	3,371	1,497	9.00	9.43	4.61
1u06	1,133,517	60,159	2,901	3,034	1,306	9.51	9.47	3.94
5nfc	2,812,807	48,828	2,982	3,127	1,316	8.89	9.16	4.18
Total	8,535,458	177,516	3,049	3,208	1,397	8.68	8.96	4.25

A split of the per-target fragment libraries by input options highlights the better fragment library quality under certain conditions with regards to the mean FLIB score and RMSD (Fig. 7.4). In particular, top- L (6 residues sequence separation) METAPSICOV STAGE 1 contact predictions yielded the lowest for both metrics across all targets. A comparison of the sequence separation, i.e. using all contact pairs or medium- and long-range ones only, strongly suggests much lower and thus more favourable scores for using short-, medium- and long-range contact pairs. A very similar difference is noticeable for METAPSICOV STAGE 2 contact predictions (Fig. 7.4).

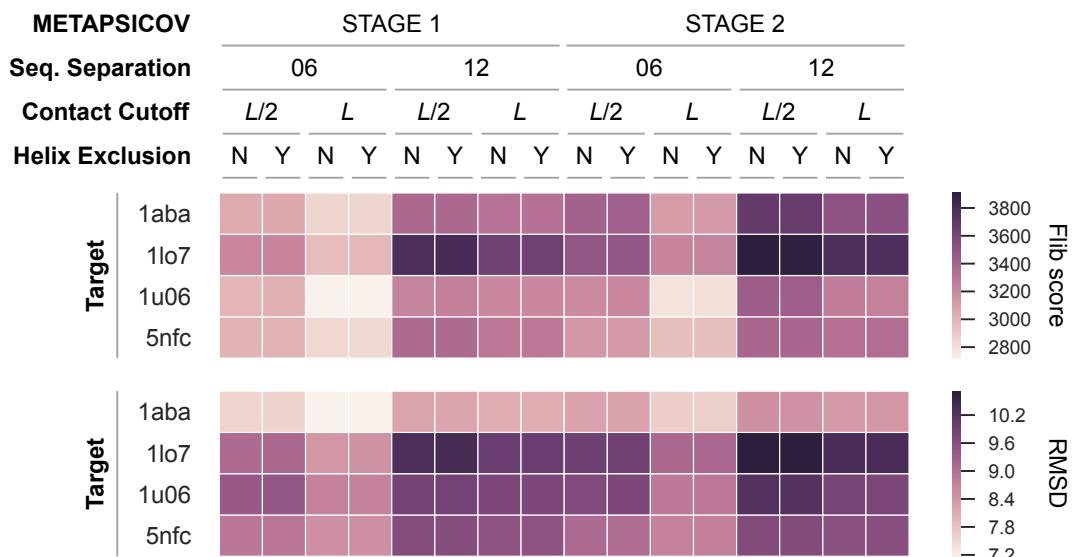


Figure 7.4: FLIB fragment library comparison for four targets highlighting the differences in mean FLIB score and RMSD by starting with different subsets of contact predictions. L refers to the number of residues per target sequence. Y refers to idealised α -helical fragment exclusion during fragment picking; N refers to treating those fragments like all others.

In this study, predicted contact information was used to further guide fragment selection. The FLIB algorithm only selected fragment for positions of the target sequence with at least one contact pair. Given this scenario, an analysis of the coverage of the target sequence with respect to each picking strategy further demonstrates the benefits of starting with METAPSICOV STAGE 1, i.e. noisier contact predictions (Fig. 7.5). Coverage is more evenly spread across the target sequences compared to missing regions especially for target 4-hydroxybenzoyl CoA thioesterase (PDB ID: 1lo7) when starting with METAPSICOV STAGE 2 predictions. Noticeably, none of the picking strategies yielded any fragments for the C-termini of α -spectrin SH3 domain (PDB ID: 1u06) and galectin-3 CRD (PDB ID: 5nfc) (Fig. 7.5). Furthermore, an analysis of the precision of fragments in each library strongly supports the benefits of starting with top- L (6 residues sequence separation) METAPSICOV STAGE 1 contact pairs. Across all four targets, the coverage of correct fragments (classed by $\text{RMSD} < 1.5\text{\AA}$ to the reference structure) is highest for this condition. This is of particular importance for α -spectrin SH3 domain (PDB ID: 1u06) and galectin-3 CRD (PDB ID: 5nfc), for which most strategies picked very few to no correct fragments. Excluding idealised α -helical fragments does not affect the quality of the FLIB libraries greatly. A consideration of differences in mean FLIB and RMSD scores shows Δ differences of 25.68 and 0.06 between the comparable libraries, i.e. with and without idealised α -helical fragments.



Figure 7.5: Summary of the coverage and precision of FLIB fragment libraries according to their target sequence. The coverage of all fragments with respect to their target-aligned sequence register are shown in red bars, and fragments with RMSD < 1.5 Å to the reference structure in blue. The predicted secondary structure of each target sequence is given at the top: α -helices (red), β -strands (blue), and loops (gray). Contact prediction information is illustrated using black bars. The fragment frequency is shown using a log-scale.

Given that FLIB uses co-evolution data to help select fragments, it is little surprise that higher degrees of TP fragments co-localise with high-density contact pair regions along the target sequence (Fig. 7.5). This characteristic explains less TP fragments in top- $L/2$ fragment libraries because less contacts (compared to top- L) are available during fragment selection. The resulting selection is purely based on the FLIB score which might not yield high-accuracy fragments ($\text{RMSD} < 1\text{\AA}$) as frequently. Therefore, the co-localisation of TP FLIB fragments and regions of high-density contact predictions highlights the importance of adding this additional source of information to pick fragments.

7.3.3 FLIB fragment selection for Molecular Replacement

One of the most important aspects of bypassing *ab initio* structure prediction and using the relevant fragments directly as MR search models is the selection of the fragments with the highest similarity between fragment and target structure.

A fragment’s FLIB score — its cumulative absolute error of predicted torsion angles — has the highest correlation with the RMSD of a fragment compared to all other scores used in the FLIB protocol [38]. To validate this finding, all non-homologous fragments in this study were tested for a correlation between their FLIB scores and RMSD values. The Spearman’s rank-order correlation coefficient analysis confirms the correlation between a fragment’s FLIB and RMSD scores (Fig. 7.6). However, the strength of the correlation varies greatly between different fragment libraries and targets. The optimal fragment picking strategy — top- L (6 residues sequence separation) METAPSICOV STAGE 1 — results in the strongest correlations across all targets. The same contact pair selection with METAPSICOV STAGE 2 predictions results in the second greatest correlations. Noticeably, the bacteriophage T4 glutaredoxin (PDB ID: 1aba) fragment libraries show much more positive correlations than the remaining targets. The fragments selected for α -spectrin SH3 domain (PDB ID: 1u06) show the overall weakest correlations. It is worth noting that the both targets, PDB IDs 1aba and 1lo7, are classed as mixed α - β targets, and thus the strength of this correlation might be fold dependent.

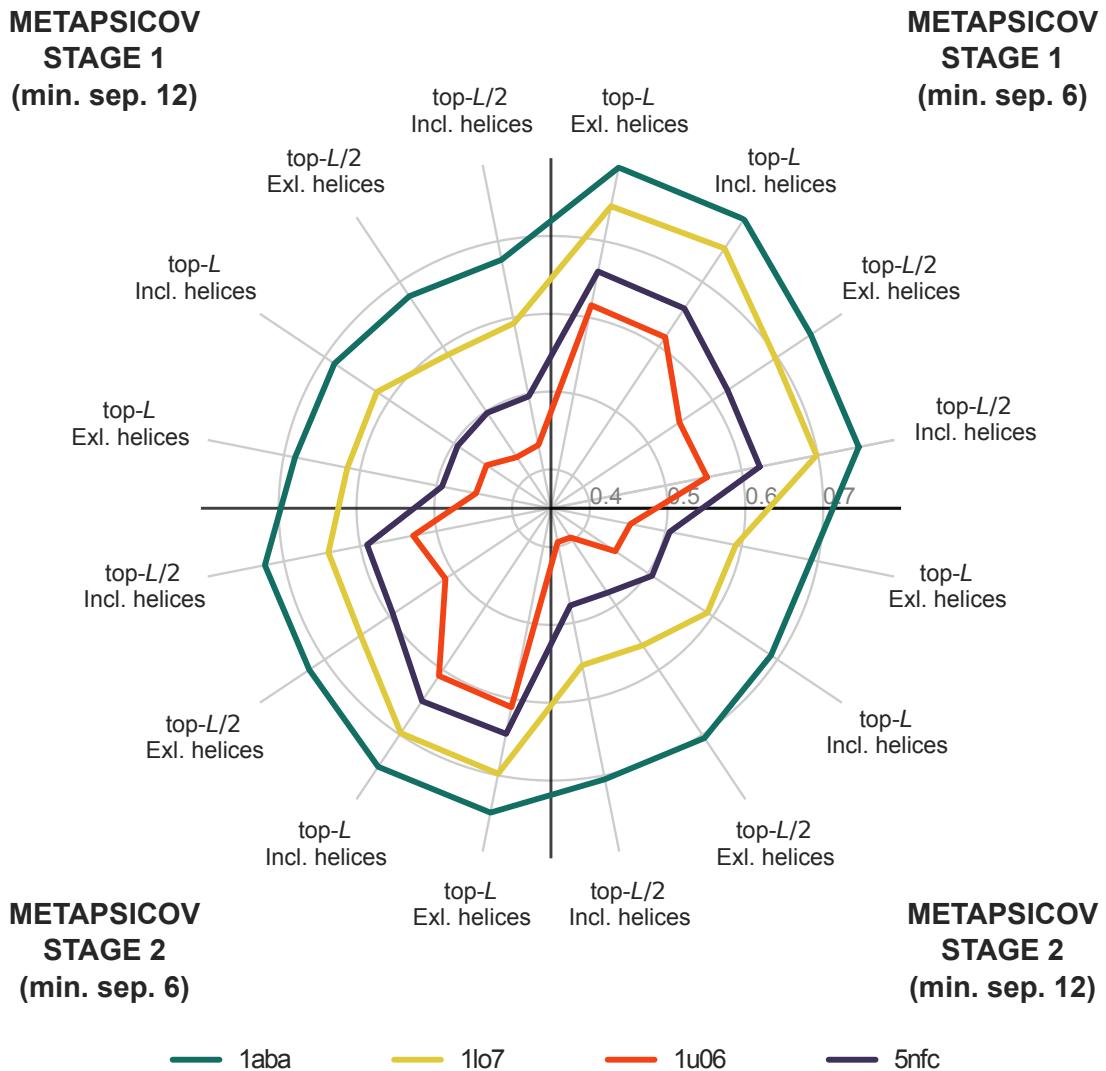


Figure 7.6: Spearman rank-order correlation coefficient analysis of FLIB fragments' FLIB score and RMSD value given the 16 unique fragment picking strategies across four targets. P-values of all Spearman correlations are < 0.001 and not shown for simplicity of the plot.

Further inspection of the fragments and the relationship between each fragment's FLIB score and RMSD value reveals a small subset of outliers in each fragment library. These fragments (hereafter referred to as outlier fragments) are sparse in each library with an overall mean count of $< 0.2\%$. An analysis for unique characteristics of these outliers, which would allow for their exclusion, reveals no unique feature. These fragments contain all secondary structure types, span the entire target sequence and range over all peptide lengths. Furthermore, they occur in all fragment libraries, irrelevant of their original picking strategy. The only characteristic setting these outlier fragments apart from the remaining set is a RMSD value of $> 30\text{\AA}$. Nevertheless, it appears that these outlier

fragments with unusually high RMSD values are never included in the final fragment search model set, given that their overall FLIB_{min} score is 796 units (one order of magnitude more than the overall minimum for the remaining fragments).

An analysis of the fragment metrics in the final MR set (6,547 fragments) further supports the positively linear relationship between a fragment's FLIB score and RMSD (Fig. 7.7a). However, the best FLIB fragments by RMSD show much less spread compared to the best fragments by FLIB score (Fig. 7.7b). Furthermore, the size of the fragments also positively correlates with the the FLIB ($\rho_{Spearman} = 0.860, p < 0.001$) and RMSD ($\rho_{Spearman} = 0.697, p < 0.001$) values. Longer fragments with higher dissimilarity with respect to the target show higher FLIB scores and RMSD values (Fig. 7.7a).

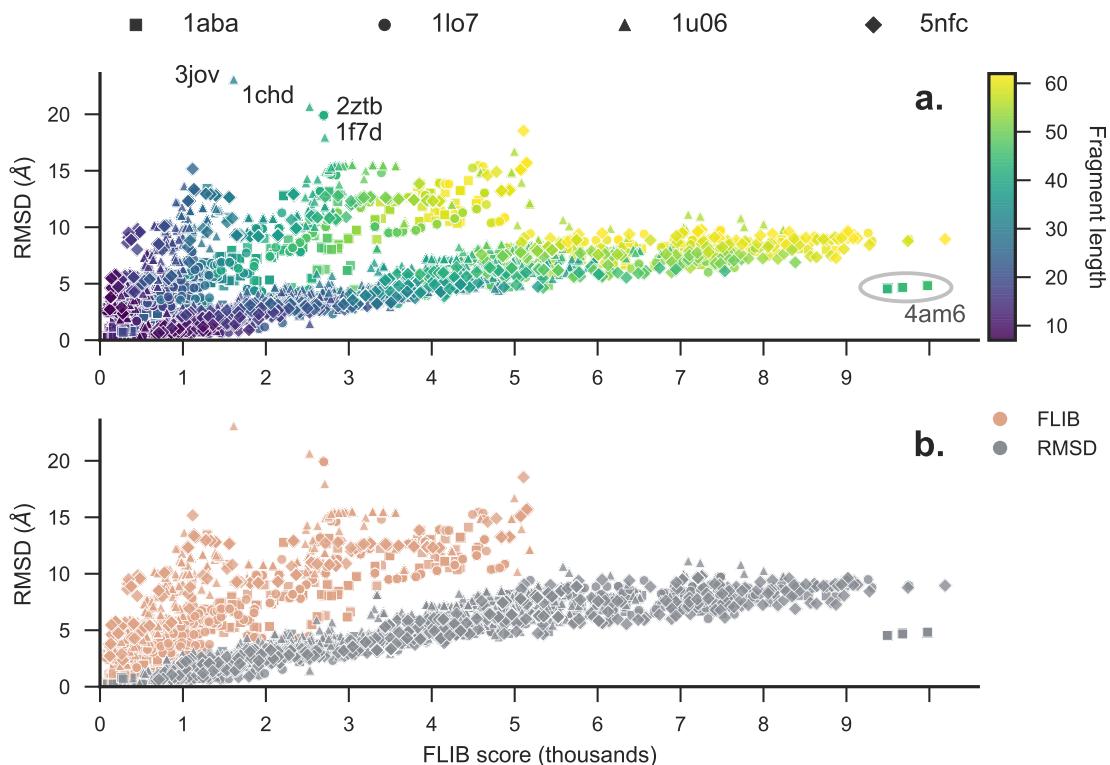


Figure 7.7: Scatterplot highlighting the positive correlation between fragment FLIB scores and RMSD values. The plot contains all fragments independent of target or picking strategy. **a.** The colour of each scatter point illustrates the fragment length. All extreme outlier fragments are highlighted with their PDB identifiers as labels. **b.** The colour codes indicate the sorting strategy to select the top FLIB fragments for each fragment peptide length bin.

Notably, a cluster of large fragments with some of the highest FLIB scores in the set show a reasonable similarity to their target structure (Fig. 7.7a). All fragments in this

cluster were picked for the bacteriophage T4 glutaredoxin sequence (PDB ID: 1aba) and extracted from the same region of the crystal structure of the actin-related protein ARP8 (PDB ID: 4am6). In comparison, some smaller fragments with peptide lengths < 50 residues and lower FLIB scores of < 3000 show the highest RMSD values in the final set.

One further unique aspect of this study compared to other fragment-MR approaches is the use of residue-residue contact information to select fragments during picking, only selecting fragments for target-sequence residues with at least one contact pair in the predicted set (Saulo de Oliveira, personal communication). In the final set 39% of all fragments satisfy at least one, 26% at least two and 20% at least three contact pairs. Across the four targets, 50% of all fragments selected for 4-hydroxybenzoyl CoA thioesterase (PDB ID: 1lo7) satisfy at least one predicted contact pair (Fig. 7.8). In comparison, 28% of fragments selected for the α -spectrin SH3 domain (PDB ID: 1u06) satisfy at least one contact pair.

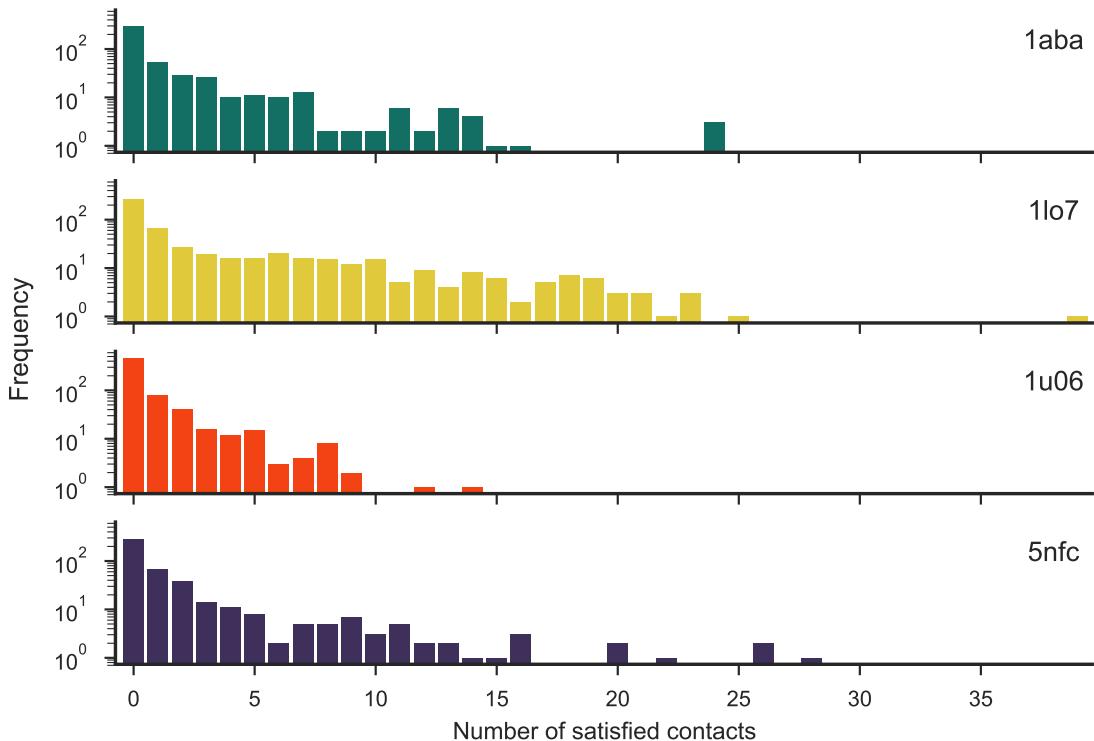


Figure 7.8: Distribution of contact precision for FLIB fragments selected as MR search models separated on a per-target basis.

Thus, the final set of FLIB fragment MR search models spans a wide range of peptide lengths, RMSD values, contact precision scores, and generally secondary structure make-

up. To illustrate the latter, a random selection of sample fragments is illustrated in Fig. 7.9. Importantly, not a single super-secondary structure motif dominates the set, increasing the sampling diversity to be undertaken during MR.

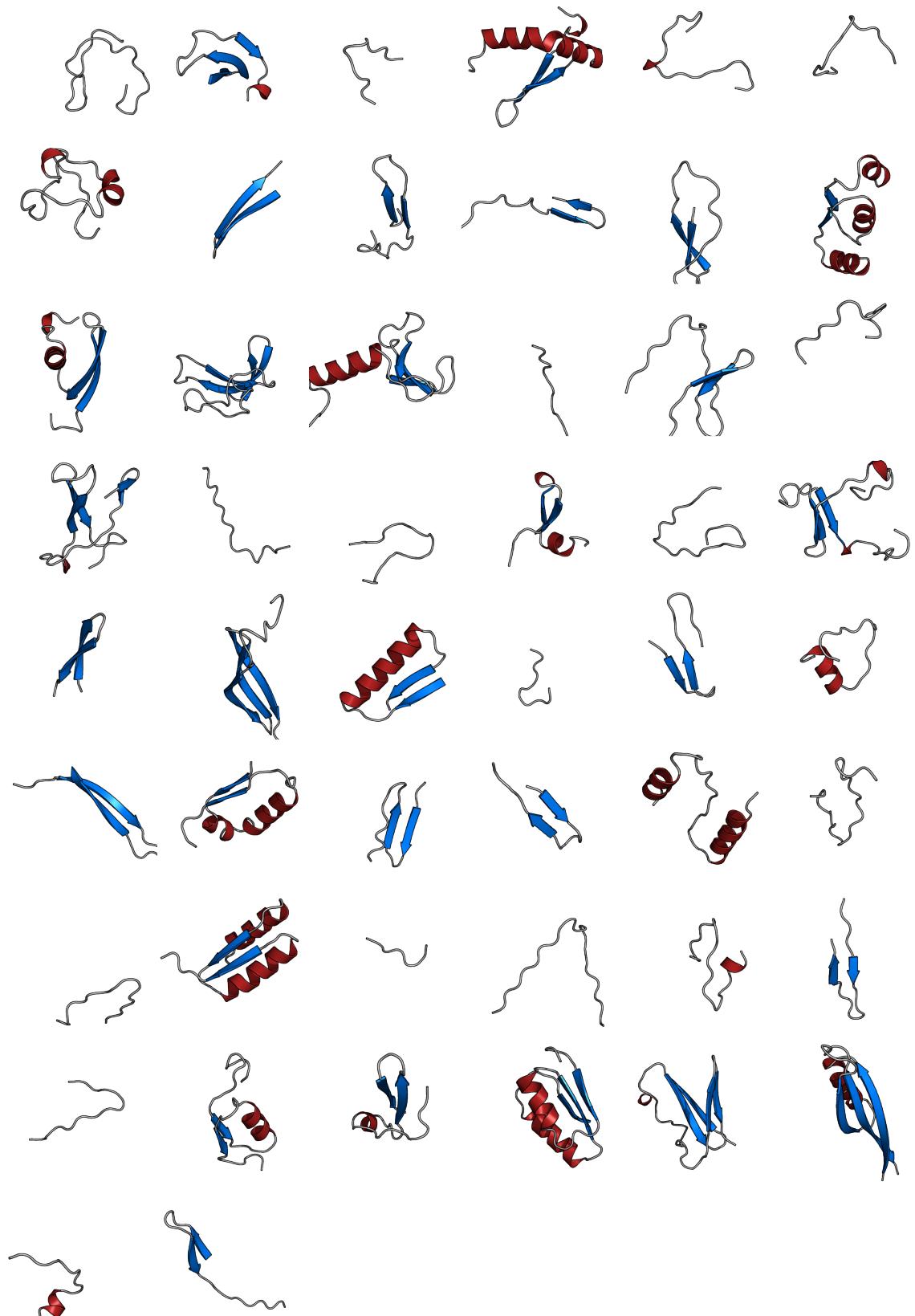


Figure 7.9: Non-redundant sample of FLIB fragment search models selected for four different protein targets. Secondary structure defined by and visualisation done in PyMOL [50]. Unpaired β -strands rendered using the loop style.

7.3.4 Molecular Replacement using FLIB fragments

FLIB fragments picked for four target sequences using a variety of FLIB input options generated $> 6,500$ fragments, which were subjected to the MR pipeline MRBUMP with their corresponding target experimental data. Given that each fragment was trialled with three different PHASER RMSD values, a total of 19,716 MR attempts were made across four target structures. Out of nearly 20,000 MR attempts, 299 led to the structure solutions of two targets, namely the T4 glutaredoxin (PDB ID: 1aba) and α -spectrin SH3 domain (PDB ID: 1u06) (Fig. 7.10).

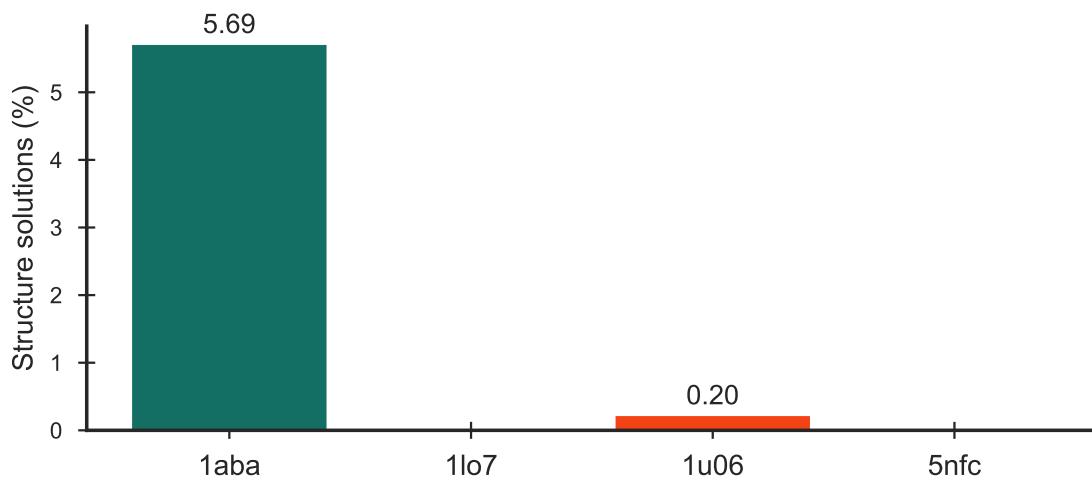


Figure 7.10: Distribution of structure solutions by FLIB target. All MR attempts total to 19,716, out of which 299 are structure solutions. Values above each bar indicates percentage search models successful out of the corresponding set.

The total of 299 MR structure solutions were achieved by 70 sequence-unique fragments. Sixty-nine of those fragments were picked from 60 unique structures for the T4 glutaredoxin (PDB ID: 1aba) leading to 97% of all structure solutions. In comparison, a single fragment, selected from three different fragment libraries, led to 9 structure solutions of the α -spectrin SH3 domain (PDB ID: 1u06). The largest FLIB fragment leading to a structure solution contained 37 residues and the smallest 10.

A division of FLIB-fragment search models by their respective origin libraries provides strong evidence that METAPSICOV STAGE 1 contact predictions allows for the selection of the most accurate fragments (Fig. 7.4), which directly translates into the structure solution count (Fig. 7.11). Furthermore, this division also highlights and supports the quality

of fragment libraries picked with top- L (6 residues sequence separation) METAPSICOV STAGE 1 predictions. Trialling the optimal fragment picking strategy with and without helical fragments ($> 90\%$ α -helical content assigned using DSSP) resulted in the library without outperforming the other (Fig. 7.11, 3rd and 4th bars).

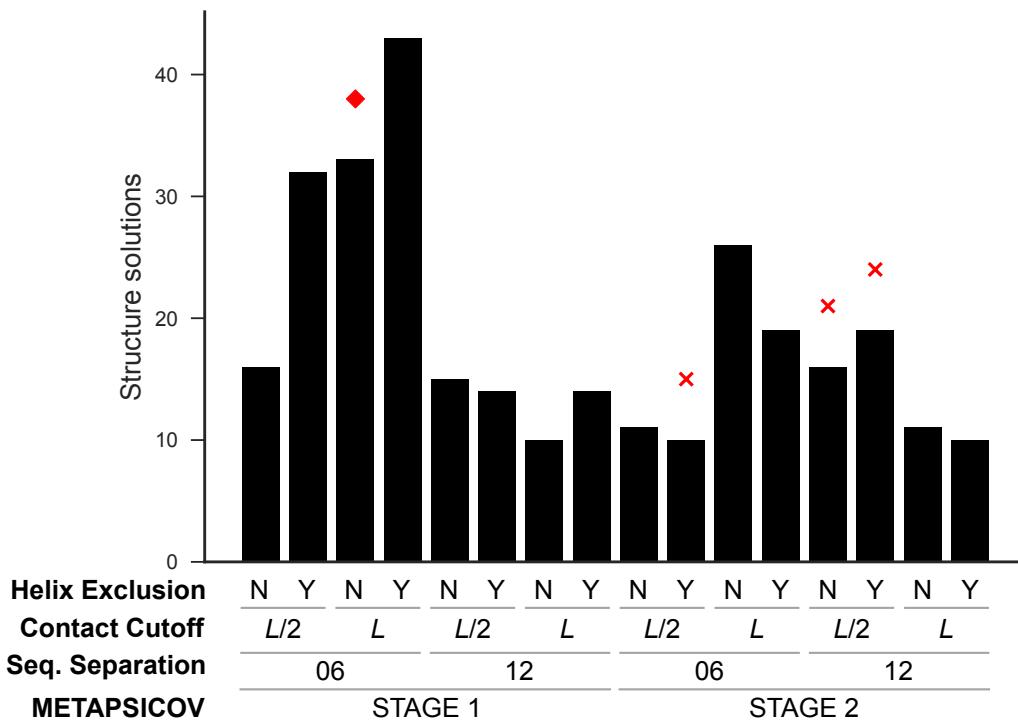


Figure 7.11: Distribution of structure solutions by FLIB library configuration. The optimal fragment picking strategy, as assessed by FLIB values, is highlighted with a red diamond to illustrate that the method that picks the best fragments is close to, but not the absolute best for ultimate structure solution. Fragment picking strategies leading to solutions of α -spectrin SH3 domain (PDB ID: 1u06) are highlighted with red crosses.

An analysis of the binned results by fragment-ranking or PHASER RMSD value confirms the expected outcome: the top fragments selected by fragment RMSD score result in more structure solutions than their FLIB score counterparts (Fig. 7.12). To reiterate, all FLIB fragments were grouped by their peptide length, and the top fragment in each group selected when sorted by either FLIB or RMSD values. When separating the total number of structure solutions by the score that made each fragment the best in its original library, it becomes clear that two-thirds of solutions were achieved with fragments scoring best by RMSD. However, the structure of α -spectrin SH3 domain (PDB ID: 1u06) was only solved with fragments that scored best in their FLIB fragment libraries by FLIB score. A further subdivision of successful fragments, sorted either by FLIB scores or RMSD values,

highlights that a larger proportion of successful RMSD-sorted fragments satisfied at least 1 contact (FLIB-sorted: 7%; RMSD-sorted: 13%). A separation of attempts by PHASER input RMSD value suggests a value of 0.1 to be the most favourable.

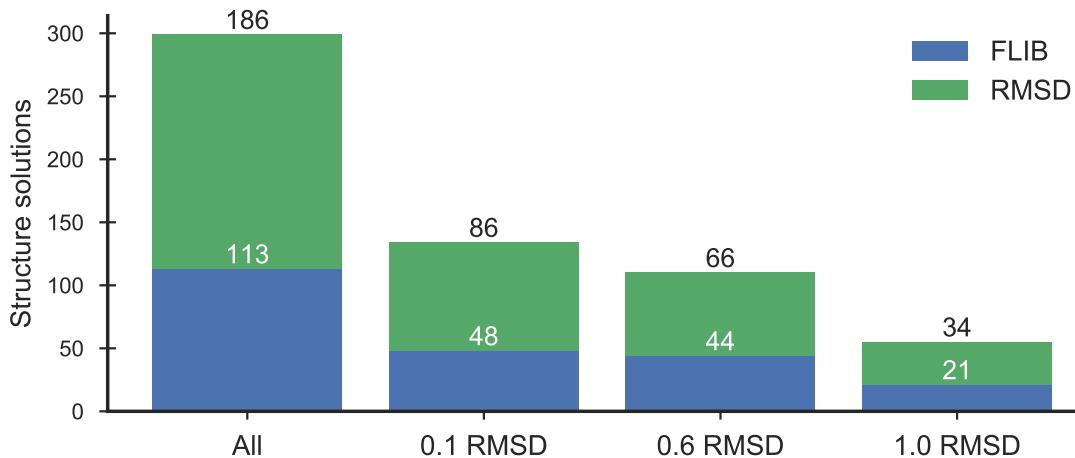


Figure 7.12: Distribution of structure solutions by fragment and MRBUMP configuration. The structure solution count is provided above each bar.

In MR, the correct placement of very small structural fragment may not always be detectable by the output metrics of underlying software. In benchmarking exercises, the Residue-Independent Overlap (RIO) metric has shown to be a very useful and powerful metric to detect such situations [14, 18]. Given that the peptide lengths of FLIB fragments in this study range from 6 to 63 residues, the RIO score is most suitable in validating the correct placement of FLIB-fragment search models. Indeed, all fragments with SHELXE CC ≥ 25 and ACL ≥ 10 contain at least 3 correctly placed C α atoms (i.e. a RIO score ≥ 3). Furthermore, the RIO metric indicates that more than 500 fragments have C α atoms placed within 1.5 \AA of any atom in the target structure. However, only 4 residues are on average placed correctly, which was not enough to achieve structure solution (Fig. 7.13). All successful FLIB fragments have a minimum model- and target-normalised RIO scores of 29.7% and 9.2% (Fig. 7.13, green markers).

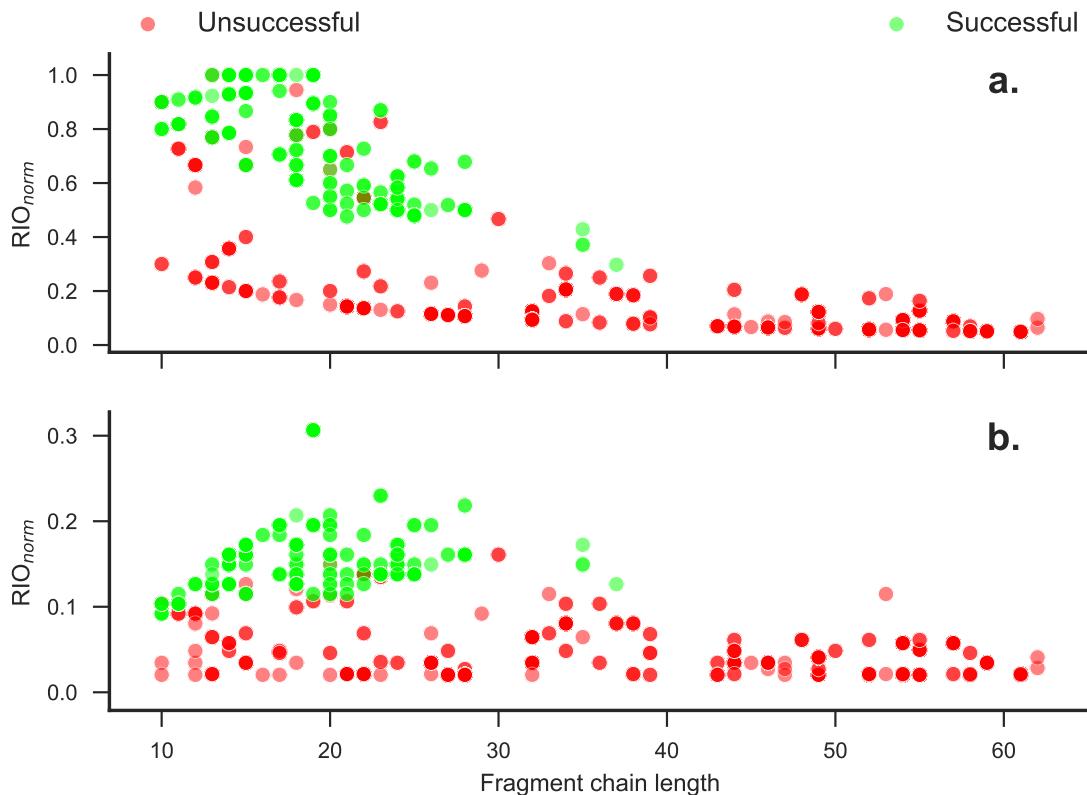


Figure 7.13: Dependence of normalised Residue-Independent Overlap (RIO_{norm}) score on the fragment chain length. The two plots show RIO scores normalised by the chain lengths of (a) the fragment and (b) the target. Colour coding indicates if the FLIB-fragment search model resulted in a structure solution. Each plot contains 890 fragment points; however, not all points are visible due to the superposition of individual scatter points because the same fragment was scored under different MR conditions.

In 33 MR attempts more than 60% of a fragment's residues were placed correctly, yet structure solution was not achieved. These trials affect exclusively fragments picked for the target sequences of T4 glutaredoxin (PDB ID: 1aba) and 4-hydroxybenzoyl CoA thioesterase (PDB ID: 1lo7). Overall, the 33 MR attempts made were done with 17 fragments extracted from 15 templates containing between 10 and 23 amino acids. The fragments' RMSD values range from 0.19 to 2.72 Å with a mean RMSD of 1.10 Å. Surprisingly, almost all of these fragments contain primarily α -helices. Given the presence of helices in the fold of both targets (Fig. 7.1) and the success of idealised fragments to solve such targets with data resolution < 2.0 Å, it is a surprise to not see more structure solutions from these fragments.

Finally, the co-evolution data used in this study select fragments is a novelty in the field. Thus, it is of great interest to identify if fragments leading to structure solution satisfy

many predicted residue-residue contacts. Eighty-seven percent ($n = 61$) of all unique fragments leading to structure solutions for either target satisfy no predicted residue-residue contact. The remaining nine fragments, all of which lead to structure solutions of T4 glutaredoxin (PDB ID: 1aba), satisfy either one ($n = 4$), two ($n = 4$) or 24 ($n = 1$) predicted contacts.

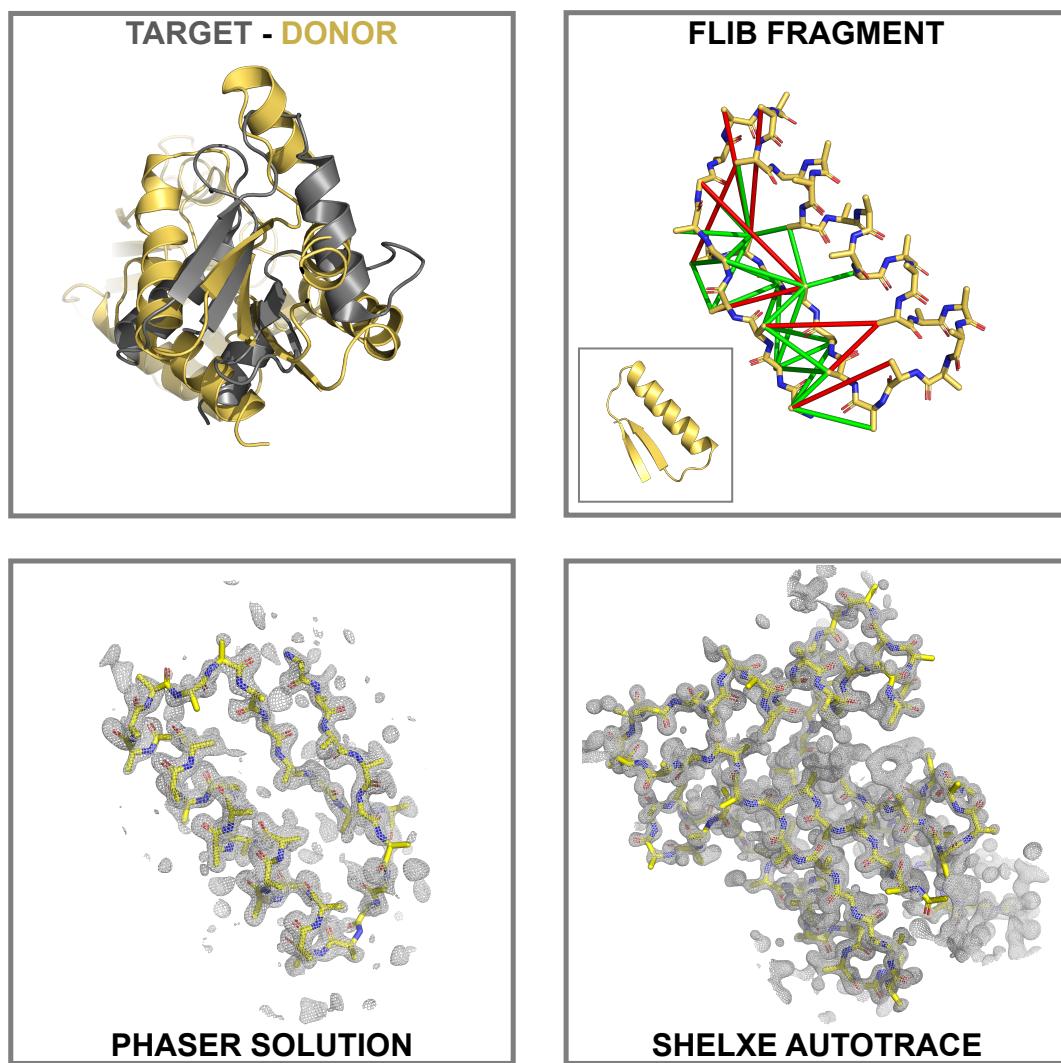


Figure 7.14: Intermediary steps from donor structure to SHELXE main-chain autotrace for a fragment derived from cobalt chelatase found in *Salmonella typhimurium* (PDB ID: 1qgo). The structure solution was obtained against the target crystallographic data of T4 glutaredoxin (PDB ID: 1aba). METAPSICOV STAGE 2 predicted contacts, against which the fragment was selected, are illustrated with True Positive (green) and False Positive (red) contacts (distance cutoff of 8 Å). 2mFo-DFc electron density maps shown at 2.0 sigma and radius around the peptide atoms of 5 Å. The RMSD between the sequence-independently superposed structures of target and donor is 10.384 Å (computed with the `super` command in PyMOL [50]).

The fragment with 24 satisfied contacts is a particularly striking example of the value of the approach explored in this study (Fig. 7.14). The fragment was derived from the template structure of cobalt chelatase found in *Salmonella typhimurium* (PDB ID: 1qgo). The picked fragment contains 35 residues and its supersecondary structure consists of a two-strand β -sheet packing against a single α -helix. The majority of satisfied contact pairs are between C β atoms of the β -strands; however, a small number of individual contact pairs also identifies the packing of one β -strand against the α -helix (Fig. 7.14, top-right). Although not considered at this stage in the FLIB algorithm, this particular fragment satisfies 75% of all relevant contact pairs. Most importantly though, this fragment was derived from an entirely unrelated protein structure, and thus illustrated the value in *ab initio* structure prediction fragments as MR search models.

7.4 Discussion

The main objective of this study was to investigate the application of FLIB structural fragments to MR. Four experimental datasets were chosen and 16 FLIB fragment libraries built per target sequence varying primarily in the predicted residue-residue contact information. A selection of highest scoring fragments were then forwarded to MRBUMP to trial each fragment as MR search model. The findings in this study validate the concept of this approach. Firstly, a positive correlation between a fragment's FLIB score and RMSD value was identified. These correlations were target-independent and found, with various strengths, in all FLIB fragment libraries. Furthermore, this work has identified top- L (6 residue sequence separation) METAPSICOV STAGE 1 contact pairs to be the optimal selection of contact pairs for the FLIB algorithm when starting with METAPSICOV predictions. The additional noise, typically filtered in the second STAGE of the METAPSICOV algorithm [3], allowed for the selection of more accurate fragments across the entire target sequence. Lastly, trialling a selection of high-scoring FLIB fragments in routine MR showed the usefulness of such fragments in attempting to solve protein structures. Two out of four targets were successfully solved albeit only trialling a small proportion of FLIB fragments per library (mean MRBUMP runtime of 10.5 CPU hours per fragment).

Intuitively, most crystallographers would declare the limitations of this approach to be the size and quality of the selected FLIB fragments as well as the resolution of the crystallographic data. Although the former was long-thought to be a major limitation, more recent work highlighted the success of likelihood-based MR methods (i.e., PHASER [47]) with very small search models. McCoy et al. [51] demonstrated the successful *ab initio* MR structure solution of aldose reductase starting from as little as two correctly placed atoms. Furthermore, automated MR pipelines, such as AMPLE [52], ARCIMBOLDO [16], BORGES [17], FRAGON [53] or FRAP [54], also successfully demonstrated MR successes with search models comprising a fraction of the target structure. Thus, MR structure solutions with FLIB fragments as short as 6 residues should be considered possible, especially when high resolution data is available and the fragment size is proportionally large compared to target size.

MR search models need to be sufficiently accurate to derive phase information for successful structure solution. The findings in this study highlight the success of identifying accurate fragments solely by the fragment's FLIB score. Given that the FLIB implementation used in this study only selected fragments for positions with at least one available contact pair, future research is required to identify the potential benefits of specifically selecting fragments that satisfy at least one contact pair. Furthermore, it is important to understand the potentially beneficial implications of using the contact satisfaction score in the FLIB score metric of a given fragment. In theory, higher precision scores should imply a closer match of the overall tertiary structure of the trialled region. Alternatively, selecting secondary structure motifs or substructures of templates by means of searching with a predicted contact map could be an attractive alternative. Recent studies indicated success in identifying sub-folds by means of Contact Map Overlap (CMO) [22, 55]. Further work also needs to explore the benefits of considering the expected Log-Likelihood Gain (eLLG) as a conceptual framework to identify the linked effects of the fragment search model size, its accuracy and the resolution on the solvability of a target structure McCoy et al. [51].

Nevertheless, FLIB fragments with near-identical subfolds to the target might not be traceable by current means of assessing structure solutions. Commonly, MR success is judged by the combination of SHELXE CC and ACL scores [49]. However, it is known

that β -strands are notoriously difficult to trace, and thus SHELXE might not pick up on correctly placed search models. Although this study did not suffer from this problem for fragments containing primarily β -strands, it did have correctly placed α -helices without structure solutions. Thus, the approach taken in this study would benefit from improvements to the density modification and sequence tracing algorithms.

Finally, this work served primarily as proof-of-concept study, and thus attempted to explore a diversity of options. With a better understanding of input parameters future work could build on the work presented here and use a large-scale analysis to assess the suitability of this concept more thoroughly. Furthermore, improvements to the FLIB algorithm through the incorporation of co-evolution data should also improve the quality of *ab initio* structure predictions, which should result in a greater success rate of other MR pipelines, such as AMPLE [52].

Chapter 8

Conclusion

Appendix A

Appendix

Table A.1: Summary of the ORIGINAL dataset.

PDB ID	Molecule	ResolutionSpace (Å)	Chain Group	Chain ID	Chain Length	Molecules per ASU	Matthew's Solvent Content (%)	Fold	Citation
1a6m	Oxy-myoglobin	1.00	P2 ₁	A	151	1	1.90	36.00	all-α
1aba	T4 glutaredoxin	1.45	P2 ₁ 2 ₁ 2 ₁	A	87	1	2.22	44.62	mixed α/β
1bdo	Biotinyl domain of acetyl-coenzyme A carboxylase	1.80	P2 ₁ 2 ₁ 2	A	80	1	2.48	49.00	all-β
1bkr	Calponin Homology (CH) domain from β-spectrin	1.10	P2 ₁	A	109	1	2.04	39.80	all-α
1chd	CheB methyltransferase domain	1.75	P3 ₂ 2 ₁	A	203	1	2.35	47.65	mixed α/β
1e0s	G-protein Arf6-GDP	2.28	P6 ₁ 2 ₂	A	174	1	2.18	37.00	mixed α/β
1eaz	Phosphoinositol (3,4)-bisphosphate PH domain	1.40	C222 ₁	A	125	1	2.48	48.00	mixed α+β
1hh8	N-terminal region of P67Phox	1.80	P3 ₁	A	213	1	2.71	45.00	all-α
1kjl	Galectin-3 domain	1.40	P2 ₁ 2 ₁ 2 ₁	A	146	1	2.15	42.68	all-β
1kw4	Polyhomeotic SAM domain	1.75	P6 ₅	A	89	1	2.25	45.27	all-α
1l07	4-hydroxybenzoyl CCoA thioesterase	1.50	I222	A	141	1	2.06	40.22	mixed α+β
1mpu	Extracellular domain of murine PD-1	2.00	P2 ₁ 2 ₁ 2 ₁	A	117	1	1.67	25.80	all-β
1pnc	Poplar plastocyanin	1.60	P2 ₁ 2 ₁ 2 ₁	A	99	1	1.82	32.48	all-β
1tjx	Synaptotagmin I C2B domain	1.04	P3 ₂ 2 ₁	A	159	1	2.40	48.00	mixed α+β
1tlv	LicT PRD	1.95	P3 ₂ 2 ₁	A	221	1	2.80	50.00	all-α
2nuz	α-spectrin SH3 domain	1.85	P2 ₁ 2 ₁ 2 ₁	A	62	1	2.57	52.16	all-β
2cyj	Ankyrin	2.05	P6 ₁	A	166	1	2.28	45.99	all-α
3w56	C2 domain	1.60	I2	A	131	1	2.05	40.10	all-β
4cl9	N-terminal bromodomain of Brd4	1.40	P2 ₁ 2 ₁ 2 ₁	A	127	1	2.21	44.37	all-α
4n3h	FN3con	1.98	P4 ₁ 3 ₂	A	100	1	2.47	50.27	all-β
4w97	KstR2	1.60	C2	A	200	1	2.75	55.25	all-α

Table A.2: Summary of the PREDICTORS dataset.

PDB ID	Molecule	ResolutionSpace (Å)	Space Group	Chain ID	Chain Length	Molecules per ASU	Matthew's Solvent Content (%)	Fold	Citation
1fcy	Retinoic acid nuclear receptor HRAR	1.30	P4 ₁ 2 ₁ 2	A	236	1	2.25	45.50	all-α
1fgv	Peptide methionine sulfoxide reductase	1.60	C121	A	199	1	2.10	41.55	mixed α+β
1gm4	Cytochrome C3	2.05	P6 ₁ 22	A	107	1	2.48	50.43	all-α
1gy8	N-II domain of ovotransferrin	1.95	P3 ₁	A	159	1	2.24	45.00	mixed α/β
1k40	FAT domain of focal adhesion kinase	2.25	C121	A	126	1	2.21	44.40	all-α
1oee	Hypothetical protein YodA	2.10	C121	A	193	1	2.30	46.20	all-β
1oz9	Hypothetical protein AQ_1354	1.89	P4 ₃ 2 ₁ 2	A	150	1	2.76	55.07	mixed α+β
1q8c	Hypothetical protein MG027	2.00	P4 ₁	A	151	1	2.42	49.25	all-α
1rh	Conserved hypothetical protein	1.80	P6 ₃	A	173	1	2.12	41.98	mixed α+β
1s2x	Cag-Z	1.90	P2 ₁ 2 ₁ 2 ₁	A	206	1	2.74	54.70	all-α
1u61	Putative Ribonuclease III	2.15	I4 ₁ 32	A	138	1	6.50	80.80	all-α
1zxu	At5g01750 protein	1.70	P2 ₁ 2 ₁ 2 ₁	A	217	1	2.50	50.20	mixed α+β
2eum	Glycolipid transfer protein	2.30	C121	A	209	1	2.25	45.39	all-α
2o18	Outer surface protein A	1.90	P12 ₁ 1	O	249	1	2.19	43.87	all-β
2eqz	Sortase B	1.60	P12 ₁ 1	A	223	1	2.07	40.71	all-β
2x6u	T-Box transcription factor TBX5	1.90	P2 ₁ 2 ₁ 2 ₁	A	203	1	2.20	44.21	all-β
2y64	Xylanase	1.40	P2 ₁ 2 ₁ 2 ₁	A	167	1	2.15	43.00	all-β
2yjm	TtrD	1.84	C121	A	176	1	2.08	40.80	all-α
2yq9	2, 3-cyclic-nucleotide phosphodiesterase	3-	P2 ₁ 2 ₁ 2 ₁	A	221	1	2.10	41.70	mixed α+β
3dju	Protein BTG2	2.26	P2 ₁ 2 ₁ 2 ₁	B	122	1	1.98	37.73	mixed α+β
3g0m	Cysteine desulfurase protein suffE	1.76	P12 ₁ 1	A	141	1	1.88	34.58	mixed α+β
3cqz	Iron-regulated surface determinant protein A	1.30	P2 ₁ 2 ₁ 2	A	127	1	2.42	49.12	all-β
4aaj	N-(5-phosphoribosyl)anthranilate isomerase	1.75	P6 ₁	A	228	1	2.38	48.30	mixed α/β
4dbb	Amyloid-β A4 precursor protein-binding family A1	1.90	P4 ₁ 2 ₁ 2	A	162	1	3.25	62.10	all-β
4e9e	Methyl-CpG-binding domain protein 4	1.90	H3	A	161	1	2.42	49.23	all-α
4bj	Galectin-3	1.80	P2 ₁ 2 ₁ 2 ₁	A	138	1	2.09	41.01	all-β
4pg0	Hypothetical protein PF0907	2.30	P6 ₅ 22	A	116	1	3.25	62.10	all-β

Table A.3: Summary of the TRANSMEMBRANE dataset.

PDB ID	Molecule	Resolution (Å)	Space Group	Chain ID	Chain Length	Molecules per ASU	Matthew's Solvent Content (%)	Fold	Citation
1gn8 2bhw	Sensory rhodopsin II Chlorophyll A-B binding protein	2.27 2.50	C222 ₁ C121	A A	239 232	1 3	2.75 4.10	53.00 69.00	all- α all- α
AB80									[99] [100]
2evu	Aquaporin aquPM	2.30	I4	A	246	1	3.38	63.57	[101]
2o9g	Aquaporin Z	1.90	I4	A	234	1	3.34	63.19	[102]
2wie	ATP synthase C chain	2.13	P6 ₃ 22	A	82	5	3.41	68.00	[103]
2xov	Rhomboid protease GLPG	1.65	H32	A	181	1	3.50	64.92	[104]
3gd8	Aquaporin 4	1.80	P42 ₁ 2	A	223	1	2.73	54.97	[105]
3hap	Bacteriorhodopsin	1.60	C222 ₁	A	249	1	2.73	54.99	[106]
3ldc	Calcium-gated potassium channel	1.45	P42 ₁ 2	A	82	1	2.48	50.44	[107]
	nthK								
3ouf	Potassium channel protein	1.55	I2	A	97	2	2.40	48.76	[108]
	Leukotriene C4 synthase		F23	A	156	1	4.91	74.77	[109]
3pcv		1.90							
3rlb	ThiT	2.00	C121	A	192	2	3.89	68.39	[110]
3u2f	ATP synthase subunit C	2.00	P4 ₂ 22	K	76	5	2.32	46.92	[111]
4dye	Biotin transporter BioY	2.09	C121	A	198	3	3.27	62.40	[112]

Bibliography

- [1] H. Kamisetty, S. Ovchinnikov, D. Baker, *Proceedings of the National Academy of Sciences Sept.* **2013**, *110*, 15674–15679.
- [2] M. J. Skwark, D. Raimondi, M. Michel, A. Elofsson, en, *PLoS Comput. Biol.* **Nov.** **2014**, *10*, e1003889.
- [3] D. T. Jones, T. Singh, T. Kosciolak, S. Tetchner, en, *Bioinformatics Apr.* **2015**, *31*, 999–1006.
- [4] J. Ma, S. Wang, Z. Wang, J. Xu, en, *Bioinformatics Nov.* **2015**, *31*, 3506–3513.
- [5] B. He, S. M. Mortuza, Y. Wang, H. B. Shen, Y. Zhang, en, *Bioinformatics Mar.* **2017**, *33*, 2296–2306.
- [6] M. Michel, M. J. Skwark, D. M. Hurtado, M. Ekeberg, A. Elofsson, en, *Bioinformatics Sept.* **2017**, *33*, 2859–2866.
- [7] S. Wang, S. Sun, Z. Li, R. Zhang, J. Xu, en, *PLoS Comput. Biol. Jan.* **2017**, *13*, e1005324.
- [8] D. S. Marks, L. J. Colwell, R. Sheridan, T. A. Hopf, A. Pagnani, R. Zecchina, C. Sander, en, *PLoS One Dec.* **2011**, *6*, e28766.
- [9] M. Michel, S. Hayat, M. J. Skwark, C. Sander, D. S. Marks, A. Elofsson, en, *Bioinformatics Sept.* **2014**, *30*, i482–8.
- [10] B. Adhikari, D. Bhattacharya, R. Cao, J. Cheng, en, *Proteins: Struct. Funct. Bioinf. Aug.* **2015**, *83*, 1436–1449.
- [11] S. Ovchinnikov, L. Kinch, H. Park, Y. Liao, J. Pei, D. E. Kim, H. Kamisetty, N. V. Grishin, D. Baker, en, *Elife Sept.* **2015**, *4*, e09248.
- [12] B. Adhikari, J. Cheng, en, *BMC Bioinformatics Jan.* **2018**, *19*, 22.
- [13] C. A. Rohl, C. E. M. Strauss, K. M. S. Misura, D. Baker, en, *Methods Enzymol.* **2004**, *383*, 66–93.
- [14] J. M. H. Thomas, F. Simkovic, R. M. Keegan, O. Mayans, Y. Zhang, D. J. Rigden, C. Zhang, Y. Zhang, D. J. Rigden, *Acta Crystallographica Section D Structural Biology Dec.* **2017**, *73*, 985–996.
- [15] J. Xu, Y. Zhang, en, *Bioinformatics Apr.* **2010**, *26*, 889–895.
- [16] D. Rodríguez, M. Sammito, K. Meindl, I. M. de Ilarduya, M. Potratz, G. M. Sheldrick, I. Usón, en, *Acta Crystallogr. D Biol. Crystallogr. Apr.* **2012**, *68*, 336–343.
- [17] M. Sammito, C. Millán, D. D. Rodríguez, I. M. De Ilarduya, K. Meindl, I. De Marino, G. Petrillo, R. M. Buey, J. M. De Pereda, K. Zeth, G. M. Sheldrick, I. Usón, en, *Nat. Methods Nov.* **2013**, *10*, 1099–1104.
- [18] J. M. H. Thomas, R. M. Keegan, J. Bibby, M. D. Winn, O. Mayans, D. J. Rigden, en, *IUCrJ Mar.* **2015**, *2*, 198–206.
- [19] D. Frishman, P Argos, en, *Proteins Dec.* **1995**, *23*, 566–579.

- [20] Q. Wuyun, W. Zheng, Z. Peng, J. Yang, *Brief. Bioinform.* **Oct. 2016**, bbw106.
- [21] S. H. P. De Oliveira, J. Shi, C. M. Deane, en, *Bioinformatics* **Feb. 2017**, *33*, 373–381.
- [22] S. Ovchinnikov, H. Park, N. Varghese, P.-S. Huang, G. A. Pavlopoulos, D. E. Kim, H. Kamisetty, N. C. Kyrpides, D. Baker, en, *Science* **Jan. 2017**, *355*, 294–298.
- [23] S. Seemayer, M. Gruber, J. Söding, en, *Bioinformatics* **Nov. 2014**, *30*, 3128–3130.
- [24] J. Yang, R. Jang, Y. Zhang, H. B. Shen, en, *Bioinformatics* **Oct. 2013**, *29*, 2579–2587.
- [25] T. Koscioletk, D. T. Jones, en, *PLoS One* **Mar. 2014**, *9*, e92197.
- [26] D. T. Jones, en, *Proteins: Structure Function and Genetics* **2001**, *Suppl 5*, 127–132.
- [27] D. Gront, D. W. Kulp, R. M. Vernon, C. E. M. Strauss, D. Baker, en, *PLoS One* **Aug. 2011**, *6*, e23294.
- [28] N. Fernandez-Fuentes, J. M. Dybas, A. Fiser, en, *PLoS Comput. Biol.* **Apr. 2010**, *6*, e1000750.
- [29] F. Delaglio, G. Kontaxis, A. Bax, *J. Am. Chem. Soc.* **Mar. 2000**, *122*, 2142–2143.
- [30] G. Kontaxis, F. Delaglio, A. Bax, en, *Methods Enzymol.* **2005**, *394*, 42–78.
- [31] T. A. Jones, S Thirup, en, *EMBO J.* **Apr. 1986**, *5*, 819–822.
- [32] J. Abbass, J.-C. Nebel, en, *BMC Bioinformatics* **Apr. 2015**, *16*, 136.
- [33] Y. Shen, G. Picord, F. Guyon, P. Tuffery, en, *PLoS One* **Nov. 2013**, *8*, e80493.
- [34] S. C. Li, D. Bu, X. Gao, J. Xu, M. Li, en, *Bioinformatics* **July 2008**, *24*, i182–9.
- [35] I. Kalev, M. Habeck, en, *Bioinformatics* **Nov. 2011**, *27*, 3110–3116.
- [36] D. Bhattacharya, B. Adhikari, J. Li, J. Cheng, en, *Bioinformatics* **July 2016**, *32*, 2059–2061.
- [37] T. Wang, Y. Yang, Y. Zhou, H. Gong, en, *Bioinformatics* **Mar. 2017**, *33*, 677–684.
- [38] S. H. P. de Oliveira, J. Shi, C. M. Deane, en, *PLoS One* **Apr. 2015**, *10*, e0123998.
- [39] D. T. Jones, en, *J. Mol. Biol.* **Sept. 1999**, *292*, 195–202.
- [40] R. Heffernan, K. Paliwal, J. Lyons, A. Dehzangi, A. Sharma, J. Wang, A. Sattar, Y. Yang, Y. Zhou, en, *Sci. Rep.* **June 2015**, *5*, 11476.
- [41] M. Remmert, A. Biegert, A. Hauser, J. Söding, en, *Nat. Methods* **Dec. 2011**, *9*, 173–175.
- [42] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, en, *J. Mol. Biol.* **Oct. 1990**, *215*, 403–410.
- [43] C. Camacho, G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, T. L. Madden, en, *BMC Bioinformatics* **Dec. 2009**, *10*, 421.
- [44] J. Söding, A. Biegert, A. N. Lupas, en, *Nucleic Acids Res.* **July 2005**, *33*, W244–8.
- [45] A. Biegert, C. Mayer, M. Remmert, J. Söding, A. N. Lupas, en, *Nucleic Acids Res.* **July 2006**, *34*, W335–9.
- [46] R. M. Keegan, S. J. McNicholas, J. M. H. Thomas, A. J. Simpkin, F. Simkovic, V. Uski, C. C. Ballard, M. D. Winn, K. S. Wilson, D. J. Rigden, *Acta Crystallographica Section D Structural Biology* **Mar. 2018**, *74*, 167–182.
- [47] A. J. McCoy, R. W. Grosse-Kunstleve, P. D. Adams, M. D. Winn, L. C. Storoni, R. J. Read, en, *J. Appl. Crystallogr.* **Aug. 2007**, *40*, 658–674.
- [48] G. N. Murshudov, P. Skubák, A. A. Lebedev, N. S. Pannu, R. A. Steiner, R. A. Nicholls, M. D. Winn, F. Long, A. A. Vagin, en, *Acta Crystallogr. D Biol. Crystallogr.* **Apr. 2011**, *67*, 355–367.
- [49] A. Thorn, G. M. Sheldrick, en, *Acta Crystallogr. D Biol. Crystallogr.* **Nov. 2013**, *69*, 2251–2256.
- [50] W. L. DeLano, The PyMOL Molecular Graphics System, <http://www.pymol.org>, **Nov. 2002**.

- [51] A. J. McCoy, R. D. Oeffner, A. G. Wrobel, J. R. M. Ojala, K. Tryggvason, B. Lohkamp, R. J. Read, en, *Proceedings of the National Academy of Sciences* **Apr.** **2017**, *114*, 3637–3641.
- [52] J. Bibby, R. M. Keegan, O. Mayans, M. D. Winn, D. J. Rigden, en, *Acta Crystallogr. D Biol. Crystallogr.* **Dec.** **2012**, *68*, 1622–1631.
- [53] H. T. Jenkins, *Acta Crystallographica Section D Structural Biology* **Mar.** **2018**, *74*, 205–214.
- [54] R. Shrestha, K. Y. J. Zhang, *Acta Crystallogr. D Biol. Crystallogr.* **2015**, *71*, 304–312.
- [55] D. W. A. Buchan, D. T. Jones, *Bioinformatics* **Sept.** **2017**, *33*, 2684–2690.
- [56] J. Vojtěchovský, K. Chu, J. Berendzen, R. M. Sweet, I. Schlichting, en, *Biophys. J. Oct.* **1999**, *77*, 2153–2174.
- [57] H. Eklund, M. Ingelman, B. O. Söderberg, T. Uhlin, P. Nordlund, M. Nikkola, U. Sonnertam, T. Joelson, K. Petratos, en, *J. Mol. Biol.* **Nov.** **1992**, *228*, 596–618.
- [58] F. K. Athappilly, W. A. Hendrickson, en, *Structure* **Dec.** **1995**, *3*, 1407–1419.
- [59] S. Bañuelos, M. Saraste, K. D. Carugo, en, *Structure* **Nov.** **1998**, *6*, 1419–1431.
- [60] A. H. West, E. Martinez-Hackert, A. M. Stock, en, *J. Mol. Biol.* **July 1995**, *250*, 276–290.
- [61] J. Ménétrey, E. Macia, S. Pasqualato, M. Franco, J. Cherfils, en, *Nat. Struct. Biol.* **June 2000**, *7*, 466–469.
- [62] C. C. Thomas, S. Dowler, M. Deak, D. R. Alessi, D. M. van Aalten, en, *Biochem. J. Sept.* **2001**, *358*, 287–294.
- [63] S. Grizot, F. Fieschi, M. C. Dagher, E. Pebay-Peyroula, en, *J. Biol. Chem.* **June 2001**, *276*, 21627–21631.
- [64] P. Sörme, P. Arnoux, B. Kahl-Knutsson, H. Leffler, J. M. Rini, U. J. Nilsson, en, *J. Am. Chem. Soc.* **Feb.** **2005**, *127*, 1737–1743.
- [65] C. A. Kim, M. Gingery, R. M. Pilpa, J. U. Bowie, en, *Nat. Struct. Biol.* **June 2002**, *9*, 453–457.
- [66] J. B. Thoden, H. M. Holden, Z. Zhuang, D. Dunaway-Mariano, en, *J. Biol. Chem.* **July 2002**, *277*, 27468–27476.
- [67] X. Zhang, J.-C. D. Schwartz, X. Guo, S. Bhatia, E. Cao, M. Lorenz, M. Cammer, L. Chen, Z.-Y. Zhang, M. A. Edidin, S. G. Nathenson, S. C. Almo, en, *Immunity* **Mar.** **2004**, *20*, 337–347.
- [68] B. A. Fields, H. H. Bartsch, H. D. Bartunik, F. Cordes, J. M. Guss, H. C. Freeman, en, *Acta Crystallogr. D Biol. Crystallogr.* **Sept.** **1994**, *50*, 709–730.
- [69] Y. Cheng, S. M. Sequeira, L. Malinina, V. Tereshko, T. H. Söllner, D. J. Patel, en, *Protein Sci.* **Oct.** **2004**, *13*, 2665–2672.
- [70] M. Graille, C. Z. Zhou, V. Receveur-Bréchot, B. Collinet, N. Declerck, H. Van Tilbeurgh, en, *J. Biol. Chem.* **Apr.** **2005**, *280*, 14780–14789.
- [71] T. Merz, S. K. Wetzel, S. Firbank, A. Plückthun, M. G. Grütter, P. R. E. Mittl, en, *J. Mol. Biol.* **Feb.** **2008**, *376*, 232–240.
- [72] D. A. K. Traore, A. J. Brennan, R. H. P. Law, C. Dogovski, M. A. Perugini, N. Lukyanova, E. W. W. Leung, R. S. Norton, J. A. Lopez, K. A. Browne, H. Yagita, G. J. Lloyd, A. Ciccone, S. Verschoor, J. A. Trapani, J. C. Whisstock, I. Voskoboinik, en, *Biochem. J. Dec.* **2013**, *456*, 323–335.
- [73] S. J. Atkinson, P. E. Soden, D. C. Angell, M. Bantscheff, C.-W. Chung, K. A. Giblin, N. Smithers, R. C. Furze, L. Gordon, G. Drewes, I. Rioja, J. Witherington, N. J. Parr, R. K. Prinjha, en, *Med. Chem. Commun.* **Feb.** **2014**, *5*, 342–351.
- [74] B. T. Porebski, A. A. Nickson, D. E. Hoke, M. R. Hunter, L. Zhu, S. McGowan, G. I. Webb, A. M. Buckle, en, *Protein Eng. Des. Sel.* **Mar.** **2015**, *28*, 67–78.
- [75] A. M. Crowe, P. J. Stogios, I. Casabon, E. Evdokimova, A. Savchenko, L. D. Eltis, en, *J. Biol. Chem.* **Jan.** **2015**, *290*, 872–882.

- [76] B. P. Klaholz, A. Mitschler, D. Moras, en, *J. Mol. Biol.* **Sept. 2000**, *302*, 155–170.
- [77] W. T. Lowther, N Brot, H Weissbach, B. W. Matthews, en, *Biochemistry* **Nov. 2000**, *39*, 13307–13312.
- [78] R. O. Louro, I. Bento, P. M. Matias, T. Catarino, A. M. Baptista, C. M. Soares, M. A. Carrondo, D. L. Turner, A. V. Xavier, en, *J. Biol. Chem.* **Nov. 2001**, *276*, 44044–44051.
- [79] P. Kuser, D. R. Hall, L. H. Mei, M. Neu, R. W. Evans, P. F. Lindley, en, *Acta Crystallogr. D Biol. Crystallogr.* **May 2002**, *58*, 777–783.
- [80] I. Hayashi, K. Vuori, R. C. Liddington, en, *Nat. Struct. Biol.* **Feb. 2002**, *9*, 101–106.
- [81] G. David, K. Blondeau, M. Schiltz, S. Penel, A. Lewit-Bentley, en, *J. Biol. Chem.* **Oct. 2003**, *278*, 43728–43735.
- [82] V. Oganesyan, D. Busso, J. BrandSEN, S. Chen, J. Jancarik, R. Kim, S. H. Kim, en, *Acta Crystallographica - Section D Biological Crystallography* **July 2003**, *59*, 1219–1223.
- [83] J. Liu, H. Yokota, R. Kim, S. H. Kim, en, *Proteins: Structure Function and Genetics* **June 2004**, *55*, 1082–1086.
- [84] L. Cendron, A. Seydel, A. Angelini, R. Battistutta, G. Zanotti, en, *J. Mol. Biol.* **July 2004**, *340*, 881–889.
- [85] L. Malinina, M. L. Malakhova, A. T. Kanack, M. Lu, R. Abagyan, R. E. Brown, D. J. Patel, en, *PLoS Biol.* **Nov. 2006**, *4*, 1996–2011.
- [86] K. Makabe, S. Yan, V. Tereshko, G. Gawlak, S. Koide, en, *J. Am. Chem. Soc.* **Nov. 2007**, *129*, 14661–14669.
- [87] A. W. Maresso, R. Wu, J. W. Kern, R. Zhang, D. Janik, D. M. Missiakas, M. E. Duban, A. Joachimiak, O. Schneewind, en, *J. Biol. Chem.* **Aug. 2007**, *282*, 23129–23139.
- [88] C. U. Stirnimann, D. Ptchelkine, C. Grimm, C. W. Müller, en, *J. Mol. Biol.* **July 2010**, *400*, 71–81.
- [89] L. Von Schantz, M. Håkansson, D. T. Logan, B. Walse, J. Österlin, E. Nordberg-Karlsson, M. Ohlin, M. Håkansson, D. T. Logan, B. Walse, J. Österlin, E. Nordberg-Karlsson, M. Ohlin, en, *Glycobiology* **July 2012**, *22*, 948–961.
- [90] S. J. Coulthurst, A. Dawson, W. N. Hunter, F. Sargent, en, *Biochemistry* **Feb. 2012**, *51*, 1678–1686.
- [91] M. Myllykoski, A. Raasakka, M. Lehtimäki, H. Han, I. Kursula, P. Kursula, en, *J. Mol. Biol.* **Nov. 2013**, *425*, 4307–4322.
- [92] X. Yang, M. Morita, H. Wang, T. Suzuki, W. Yang, Y. Luo, C. Zhao, Y. Yu, M. Bartlam, T. Yamamoto, Z. Rao, en, *Nucleic Acids Res.* **Dec. 2008**, *36*, 6872–6881.
- [93] J. C. Grigg, C. X. Mao, M. E. P. Murphy, en, *J. Mol. Biol.* **Oct. 2011**, *413*, 684–698.
- [94] H. Repo, J. S. Oeemig, J. Djupsjöbacka, H. Iwaï, P. Heikinheimo, en, *Acta Crystallogr. D Biol. Crystallogr.* **Nov. 2012**, *68*, 1479–1487.
- [95] M. F. Matos, Y. Xu, I. Dulubova, Z. Otwinowski, J. M. Richardson, D. R. Tomchick, J. Rizo, A. Ho, en, *Proc. Natl. Acad. Sci. U. S. A.* **Mar. 2012**, *109*, 3802–3807.
- [96] S. Moréra, I. Grin, A. Vigouroux, S. Couvé, V. Henriot, M. Saparbaev, A. A. Ishchenko, en, *Nucleic Acids Res.* **Oct. 2012**, *40*, 9917–9926.
- [97] P. M. Collins, K. Bum-Erdene, X. Yu, H. Blanchard, en, *J. Mol. Biol.* **Apr. 2014**, *426*, 1439–1451.
- [98] T. Weinert, V. Olieric, S. Waltersperger, E. Panepucci, L. Chen, H. Zhang, D. Zhou, J. Rose, A. Ebihara, S. Kuramitsu, D. Li, N. Howe, G. Schnapp, A. Pautsch, K. Bargsten, A. E. Prota, P. Surana, J. Kottur, D. T. Nair, F. Basilico, V. Cecatiello, S. Pasqualato, A. Boland, O. Weichenrieder, B. C. Wang, M. O. Steinmetz, M. Caffrey, M. Wang, en, *Nat. Methods* **Feb. 2015**, *12*, 131–133.
- [99] K. Edman, A. Royant, P. Nollert, C. A. Maxwell, E. Pebay-Peyroula, J. Navarro, R. Neutze, E. M. Landau, en, *Structure* **Apr. 2002**, *10*, 473–482.

- [100] J. Standfuss, A. C. T. Van Scheltinga, M. Lamborghini, W. Kühlbrandt, en, *EMBO J.* **Mar. 2005**, *24*, 919–928.
- [101] J. K. Lee, D. Kozono, J. Remis, Y. Kitagawa, P. Agre, R. M. Stroud, en, *Proc. Natl. Acad. Sci. U. S. A.* **Dec. 2005**, *102*, 18932–18937.
- [102] D. F. Savage, R. M. Stroud, en, *J. Mol. Biol.* **May 2007**, *368*, 607–617.
- [103] D. Pogoryelov, Ö. Yıldız, J. D. Faraldo-Gómez, T. Meier, en, *Nat. Struct. Mol. Biol.* **Oct. 2009**, *16*, 1068–1073.
- [104] K. R. Vinothkumar, K. Strisovsky, A. Andreeva, Y. Christova, S. Verhelst, M. Freeman, en, *EMBO J.* **Nov. 2010**, *29*, 3797–3809.
- [105] J. D. Ho, R. Yeh, A. Sandstrom, I. Chorny, W. E. C. Harries, R. A. Robbins, L. J. W. Miercke, R. M. Stroud, en, *Proceedings of the National Academy of Sciences* **May 2009**, *106*, 7437–7442.
- [106] N. H. Joh, A. Oberai, D. Yang, J. P. Whitelegge, J. U. Bowie, en, *J. Am. Chem. Soc.* **Aug. 2009**, *131*, 10846–10847.
- [107] S. Ye, Y. Li, Y. Jiang, en, *Nat. Struct. Mol. Biol.* **Aug. 2010**, *17*, 1019–1023.
- [108] M. G. Derebe, D. B. Sauer, W. Zeng, A. Alam, N. Shi, Y. Jiang, en, *Proceedings of the National Academy of Sciences* **Jan. 2011**, *108*, 598–602.
- [109] H. Saino, Y. Ukita, H. Ago, D. Irikura, A. Nisawa, G. Ueno, M. Yamamoto, Y. Kanaoka, B. K. Lam, K. F. Austen, M. Miyano, en, *J. Biol. Chem.* **May 2011**, *286*, 16392–16401.
- [110] G. B. Erkens, R. P. A. Berntsson, F. Fulyani, M. Majserowska, A. Vujičić-Žagar, J. Ter Beek, B. Poolman, D. J. Slotboom, en, *Nat. Struct. Mol. Biol.* **June 2011**, *18*, 755–760.
- [111] J. Symersky, V. Pagadala, D. Osowski, A. Krah, T. Meier, J. D. Faraldo-Gómez, D. M. Mueller, en, *Nat. Struct. Mol. Biol.* **Apr. 2012**, *19*, 485–91, S1.
- [112] R. P.-A. Berntsson, J. ter Beek, M. Majserowska, R. H. Duurkens, P. Puri, B. Poolman, D.-J. Slotboom, en, *Proceedings of the National Academy of Sciences* **Aug. 2012**, *109*, 13990–13995.