



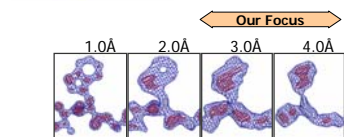
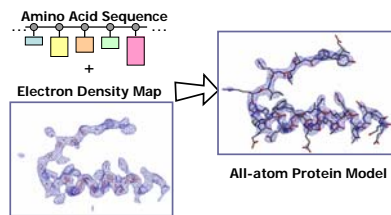
NEW APPROACH TO AUTOMATIC FITTING OF ELECTRON-DENSITY MAPS

Craig A. Bingman¹, Frank DiMaio^{2,3,4}, Ameet Soni^{1,2}, Jude W. Shavlik^{1,2}, Dmitry Kondrashov², Eduard Bitto^{1,5}, George N. Phillips, Jr.^{1,2}

University of Wisconsin-Madison, Department of Biochemistry, 433 Babcock Drive, Madison, WI USA 53706-1549 <http://www.uwstructuralgenomics.org>

Introduction

One bottleneck in high-throughput crystallography is the interpretation of a macromolecular model from the electron density map. When the map is well phased and at sufficiently high resolution, this interpretation can usually be completed more or less automatically by existing methods, such as Arp/warp (1), Textal (2), or Resolve (3). However, when the resolution is lower than about 2.5 Å or the phasing is not particularly accurate, it can sometimes take weeks or months of a crystallographer's time to complete a chain trace, registered with the sequence. This process requires substantial trial and error.

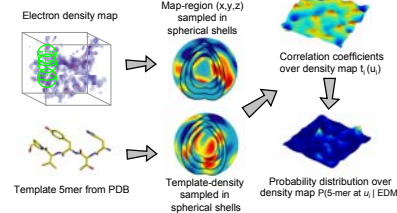


We have developed the automatic interpretation tool, ACMI (for *Automatic Crystallographic Map Interpreter*) (4-6). ACMI employs probabilistic inference to compute the probability distribution of each amino acid's 3D location given this density map.

Phase 1: Pentapeptide Matching

- Estimate probability of amino acid locations using pattern matching.
- Spherical Harmonic decomposition(6) allows for rapid search over all rotations for each amino acid template for each location in the map.

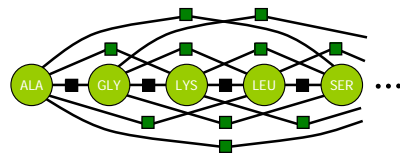
ACMI's observational potential is computed by searching for a pentapeptide centered at each amino acid location in the protein. That is, ACMI walks along the protein, one amino acid at a time, and considers the 5-amino-acid sequence centered at each position. It searches the PDB for all observed conformations of that particular sequence. Spherical harmonics are used to quickly compute the correlation coefficient between some region in the (unsolved) map and a pentapeptide fragment over all rotations for each location on the map.



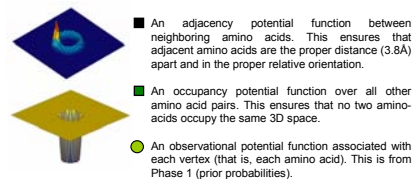
Phase 2: Probabilistic Protein Model

- Refine local search using protein structural constraints.
- Marginal probabilities approximated using Loopy Belief Propagation.

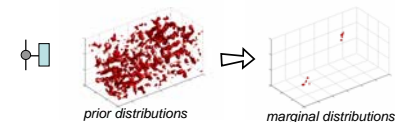
ACMI models a protein using a pairwise Markov random field, a probabilistic model where a probability distribution over some set of random variables is defined on an undirected graph. Specifically, the probability is given as the product of potential functions associated with vertices and edges in the graph.



To model a protein, ACMI constructs a graph where each vertex corresponds to an amino acid. A random variable associated with each node describes the position and orientation of that amino acid's C α . The joint probability of some configuration of C α 's is the product of ...

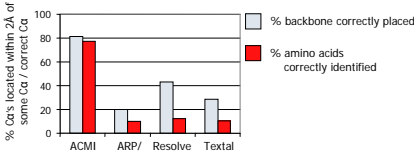


The marginal probability distributions are intractable to compute; Belief Propagation (Pearl 1988) is a message passing algorithm that provides an approximation (in loopy graphs) that is not guaranteed to converge, but works well in practice. Final result is an approximate marginal distribution of each amino acid's C α position/orientation.



Phase 2 Results

We compare ACMI's backbone trace to three popular tools on a set of 10 experimentally-phased density maps over a wide range of phase error and resolutions. We compare the resultant models (C α locations only) both in terms of average RMS error and in terms of correctly-identified amino acids. ACMI is producing a more accurate and more complete C α trace.



Phase 3: Producing an All-Atom Model

- Sampling using particle filtering produces a set of physically feasible all-atom models.
- Phase 2's belief and domain knowledge guide the search.

Given the approximate marginal distributions provided by ACMI, we want to compute an ensemble of physically feasible all-atom protein models. Our approach uses a probabilistic method known as particle filtering (PF).

A particle here refers to one specific 3D layout of all the non-hydrogen atoms in a contiguous subsequence of the protein. PF represents the distribution of some subsequence's layout using a set of particles and their weights

$$p(x_{i:k} | y_{i:k}) \approx \sum_{i=1}^N w_i^{(k)} \delta(x_{i:k} - x_{i:k}^{(i)})$$

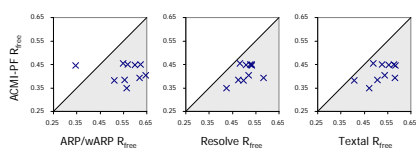
Algorithm
place "seeds" b_i^j for each particle $i=1..N$
while amino-acids remain
(A) place b_{i+1}^j / b_{i+1}^k given b_i^j for each $i=1..N$
(B) place s_i^j given b_{i+1}^j for each $i=1..N$
optionally resample N particles
end while

(A) Backbone sampling (one particle)
(1) Sample L b_{i+1}^j 's from $b_{i+1}^j - b_{i+1}^k$
(2) Weight each sample by its ACMI-computed approximate marginal
(3) Select b_{i+1}^j with probability proportional to sample weight
(4) Particle weight is sum of sample weights $w_{i+1}^j \propto \sum_{k=1}^L p(b_{i+1}^j | b_i^k) \rightarrow w_{i+1}^j$

(B) Sidechain sampling
(1) Sample s_i^j from a database of sidechain conformations
(2) For each sidechain conformation, compute probability of map given the sidechain
(3) Select sidechain conformation from this weighted distribution
(4) Particle weight is sum of sample weights $w_{i+1}^j \propto \sum_{s_i^j} p(\text{map} | s_i^j) \rightarrow w_{i+1}^j$

Phase 3 Results

ACMI's all-atom model using particle filtering produces physically feasible along with fewer chains and better RMSD than Phase 2 alone. Compared to current methods, ACMI-PF produces lower R_{free} values on all maps with only one exception.

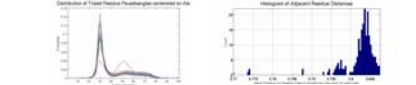


Introducing Domain Knowledge

Importance sampling's performance is dependent on the accuracy of our importance distribution (the sampling function - angle/bond constraints for backbones and PDB instances for sidechains). The likelihood of keeping a sample is dependent on the accuracy of importance evaluation function (quality of match to Phase 2 beliefs for the backbone and density correlation for sidechains).

Sampling Function

Currently, we assume C α -C α distances are the same for all pairs of residues (≈ 3.8 Å with a tight distribution). In actuality, specific pairs exhibit different distances and variation (proline in particular). In addition, sampling pseudorings for C α -C α -C α triplets from a bimodal distribution hides the fact that not all residues exhibit the same distribution of pseudorings (vary by secondary structure, surface accessibility, etc.). By accumulating data from the PDB, we can capture higher order information such as secondary structure into the sampling function.



Recent Activities

We used ACMI to construct atomic models on 20 crystallographic data sets with the initial phases available before any model building. The structures had been previously solved and deposited in the PDB, providing correct models to assess the performance of the automated model-building. The input files provided to the software are the protein sequence, the crystal symmetry and cell parameters, and the electron density map computed with the initial, experimental phases. The density maps provided different levels of challenge for model-building, depending on the accuracy of initial phasing and the available resolution. Two measures of quality and the number of residues in the asymmetric unit have been examined, with the phase error computed by comparing the initial phases with the calculated phases from the final model. In addition to the resolution, phase error is an important factor for interpretability of an electron density map: a high-resolution data set with poor initial phases map produce an indistinct map. Based on our experience and expert judgement of the quality of the initial density maps, we divided the data sets into two categories: easy to interpret and difficult. We now have the results of automated model building from ACMI and three other programs for the two categories of data sets. The model building performance is assessed for the backbone completeness and sidechain identification, in terms of prediction sensitivity (fraction of the correct model predicted) and specificity (fraction of the predicted model which is correct). Below is one sample table (manuscript in preparation). ACMI has also been used to solve a new 2.7 Å angstrom resolution structure that is currently in refinement.

Sidechain Identification Comparison for Difficult Data Sets

PDB ID	ACMI-PF	Arp/Warp	Resolve	Textal
2mf	73/5773.5	96/159.3	6/810.8	3/94.1
2q7a	80/784.8	2/07.0	35/147.1	28/9/31.9
3bu	89/5919.9	1/03.7	113/272.6	4/76.8
1wt	97/784.8	4/35.1	17/226.0	37/740.0
1tp	95/096.7	0/58.2	83/81.3	52/057.6
1vz	81/285.0	21/338.5	15/617.1	
2a3a	67/771.4	10/417.8	5/9/9.3	
2fu	63/079.9	0/20.4	5/912.9	3/87.6
2bdu	78/985.3	0/51.8	12/520.1	6/06.6
2ab1	56/070.1	0/94.9	3/45.9	3/97.5

*The two numbers reported are the percent of the correct structure predicted/percent of the predicted structure correct.

Other Improvements

Phase 1: Filtering, Fragment Selection

Phase 1 is the largest consumer of CPU time, although high parallelization mitigates the problem. The proportion of C α locations to points in the map is minuscule; thus most of the computations over the map are useless. A SVM trained to predict Pr (C α | EDM) allows a large proportion of the map to be eliminated *a priori* without losing any positives (80% precision at 95% recall). This leads to ~35% reduction in computation over using density values alone, and up to 5x improvement overall.

Fragment selection has been expanded to select larger n-mers if they exist in the DB. In addition, work is being done to search for larger homologous domains and more accurately select "good" 5mers (particularly in maps where only secondary structure can be determined).

Phase 2-3: Ensembles and Diversity

One key extension is to overcome error accumulation in BP by sampling more diverse distributions. Ensemble generation in both Phase 2 and 3 can allow greater sampling of the conformational space. These can be created via dynamic sampling in BP, early stopping of BP, selecting numerous starting points in PF (already done), and testing particle filtering variations.

References

- (1) R. Morris, A. Petrakis, and L. Lamzin. (2003). *Math Enz* 374: 229-244.
- (2) T. Joergers and J. Sacchettini. (2002). *Acta Cryst D5*, 2042-2054.
- (3) T.C. Terwilliger. (2003). *Acta Cryst D59*, 38-44.
- (4) F. DiMaio, J. Shavlik, and G. Phillips, Jr. (2006). *Bioinformatics* 22, 681-689.
- (5) F. DiMaio, D. Kondrashov, E. Bitto, A. Soni, C. Bingman, G. Phillips & J. Shavlik (2007). *Bioinformatics*. doi: 10.1093/bioinformatics/btm480
- (6) F. DiMaio, A. Soni, G. Phillips & J. Shavlik (2007). *BIBIM07*, Fremont, CA.