

ANOLEA: ATOMIC NON-LOCAL ENVIRONMENT ASSESSMENT

Francisco Melo and Ernest Feytmans.

[Laboratory of Structural Molecular Biology, University of Namur, Belgium.](#)

● Assessing Protein Structures

The assessment of protein structures is a delicate matter. Currently, there is not a single method able to consistently and accurately predict the three-dimensional structure of a protein. Moreover, our lack of knowledge about protein folding and stability makes difficult to develop such a method.

Similarly, there is no single method able to consistently and accurately predict the errors in a protein structure. For that reason, we strongly recommend to use as many as possible different methods when assessing a protein structure (below we provide links to some well-known and widely used methods). Different methods use different approaches, thus they can complement to each other. In that way, you can have more confidence about the predicted error of specific regions in the protein.

● A List of 'Well-known' and 'Widely-used' Methods to Assess Protein Structures

- [PROCHECK](#)
- [PROSA II](#)
- [SURVOL](#)
- [VERIFY3D](#)
- [WHAT CHECK](#)

● What is ANOLEA ?

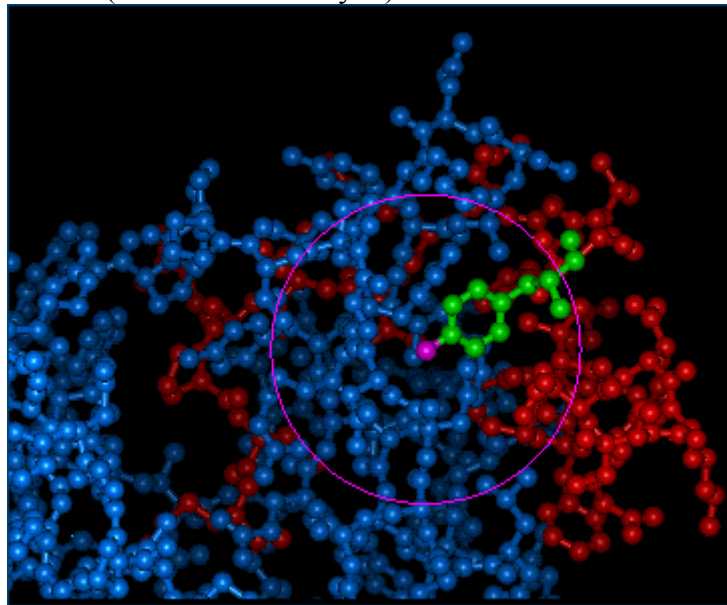
ANOLEA (Atomic NON-Local Environment Assessment) is a server to assess protein structures. It performs energy calculations on a protein chain, evaluating the "[non-local environment](#)" (NLE) of each heavy atom in the molecule. The NLE of one heavy atom is defined as all the heavy atoms, within an Euclidean distance of 7 Å, that belong to an amino acid that is farther than 11 residues in the chain or that belong to another chain. The energy of each atom pairwise interaction is taken from a [distance-dependent knowledge-based mean force potential](#) (DD-KB-MFP) that has been derived from a

[database](#) of 147 non-redundant proteins chains with a sequence identity below 25% and solved by x-ray crystallography with a resolution lower than 3 Å . If you want to see the organization of this server click [here](#).

ANOEA (Atomic Non-Local Environment Assessment) is a server that performs energy calculations on a protein chain, evaluating the "[Non- Local Environment](#)" (NLE) of each heavy atom in the molecule. The energy of each pairwise interaction in this non-local environment is taken from a [distance-dependent knowledge-based mean force potential](#) that has been derived from a [database](#) of 147 non-redundant protein chains with a sequence identity below 25% and solved by X-Ray crystallography with a resolution lower than 3 Å .

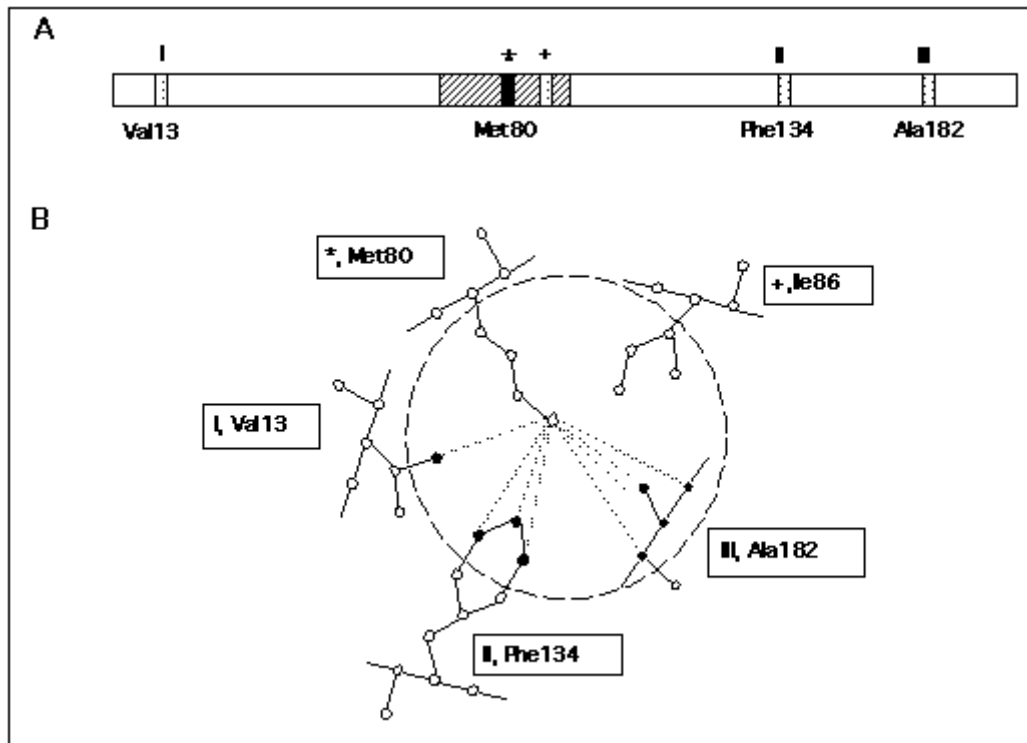
• Non local environment

The "non local environment" (NLE) of **one heavy atom** (atom shown in magenta) is defined as all the heavy atoms, within an Euclidean distance of 7 Å , that belong to an amino acid that is farther than 11 residues in the chain (atoms shown in blue) or that belong to another chain (atoms shown in cyan). This is illustrated in the Figure below.



The atoms shown in red and green are not considered in the calculation. These atoms belong to an amino acid that is closer than 12 residues in the chain.

● Non local environment scheme

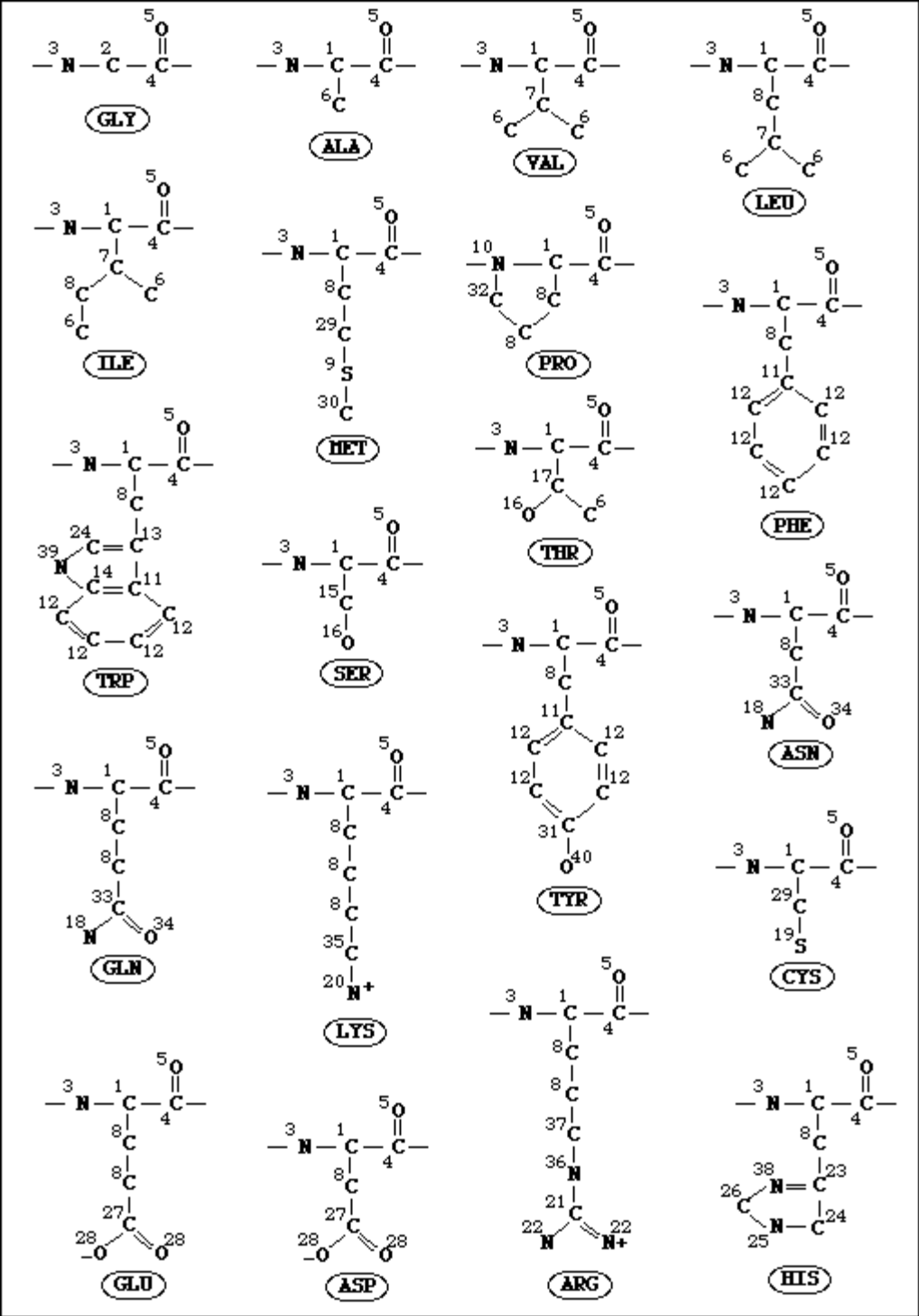


Non-local environment of an atom. (A) Linear representation of a protein sequence. The asterisk represents the amino acid (Met80) containing the atom for which non-local interaction energies are to be evaluated (Ce). The hatched area in the sequence represents the amino acids that are considered as local for Met80 (i.e. not farther than 11 residues in the sequence) and that are not considered in the calculation. Regions I, II and III represent non-local interacting amino acids for Met80, because some of their atoms are closer than 7 Å to Ce of Met80. (B) Two-dimensional structural representation of part of the protein chain illustrated in A. The dashed circle represents a sphere of 7 Å radius centered on Ce of Met80. All the atoms inside this sphere that belong to Met80 (Cb, Cg, Sd) and the local amino acid Ile86 (Cb, Cg1, Cg2, Cd) are shown as open circles and are not considered in the calculation. All the atoms of non-local amino acids that are located within the sphere (Val13, Cg2; Phe134, Ce1, Ce2, Cx; Ala182, N, Ca, C, Cb) are shown as filled circles. All these atoms are included in the non-local interaction energy calculation of the Met80 Ce atom. The Euclidean distances between each one of these atoms and the Met80 Ce atom are calculated, and the corresponding energy values are taken from the energy functions of the atomic mean force potential. The total energy of the Met80 Ce atom is the sum of all the pairwise non-local interaction energy values.

• **Distance-dependent knowledge-based mean force potential**

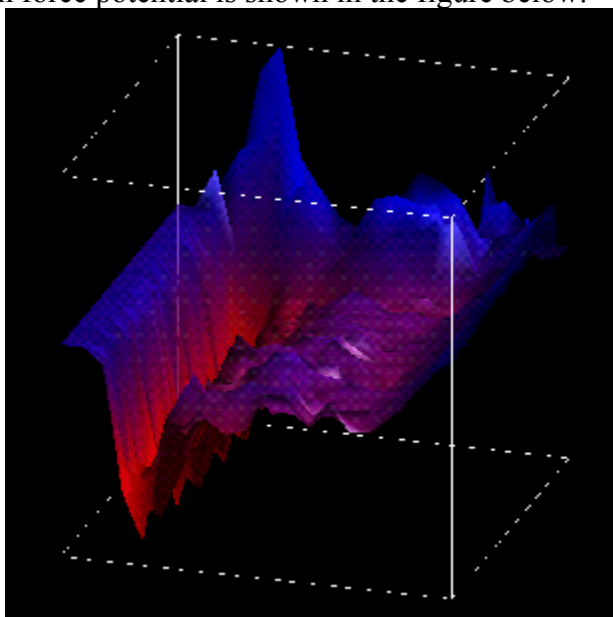
• **Atom type definitions**

ANOLEA uses a mean force potential at the atomic level, which is based on a particular definition of atom types. Forty different atom types are defined for all the heavy atoms of the 20 existing amino acids. In a strict physico-chemical point of view, all the atoms with different environments, connectivities and chemical nature, would be different. In the 20 amino acids the total number of heavy atoms is 167 and the number of non-equivalent heavy atoms is 98. We defined a total number of 40 different atom types for all the heavy atoms of the 20 amino acids (see figure below). The atom type definition is based on its connectivity, chemical nature and location level (side chain or backbone). For example, one type of heavy atom is the carbon of a methyl group bonded to a carbon with sp^3 hybridization. This type represents the beta carbon of alanine side chain, the gamma1 and gamma2 carbons of valine, the delta1 and delta2 carbons of leucine, etc.



● Mean force potential derivation

This mean force potential was calculated to evaluate the non-local environment of each atom in a protein molecule. The [non-local environment](#) of any given atom is defined as all the heavy atoms that are included within a 7 Å radius sphere centered on the given atom, and that belong to amino acids that are more distant than 11 residues in the same chain or that belong to another chain. This definition is based on the observation of the energy curves of a previous mean force potential that included local and non-local interactions, where sequence separations higher than 11 do not affect the shape of the energy functions for the different pairwise interactions. As an example, one energy curve of this previous mean force potential is shown in the figure below.



The current potential has been derived using 35 different classes of distance (l), ranging from 0 to 7 Å by steps of 0.2 Å . All pairwise occurrences out of this range were excluded. We considered all the pairs of atoms occurring within an Euclidean distance of 7 Å , that belong to amino acids that are farther than 11 residues in the chain or that belong to different chains (if you want to see a scheme of this click [here](#)). A sequence separation of 2 means that two amino acids are adjacent in the sequence.

The potential was obtained from a [database](#) of 147 non-redundant protein chains. The calculation was performed over all these chains. In the case of multimers, all the atom contacts with the other chains were considered.

We have compiled each pairwise contact symmetrically, obtaining 820 ($40 \times 41 / 2$) different [atom pairwise distributions](#). The calculation of pairwise pseudo-energy terms has been carried out as described by Sippl (1990). The following expression was used:

$$\Delta E^{ij}(l) = RT \ln [M_{ij} \times \sigma] - RT \ln \left[M_{ij} \times \sigma \times \frac{f^{ij}(l)}{f^{xx}(l)} \right]$$

where \mathbf{M}_{ij} is the total number of observations for the atomic pair \mathbf{ij} and corresponds to:

$$M_{ij} = \sum_{l=1}^n f(i, j, l)$$

σ represents the weight given to each observation. We have used $\sigma = 1/50$ as proposed by Sippl (1990), such that for 50 observations, $\mathbf{F}_{ij}(\mathbf{l})$ and $\mathbf{F}_{xx}(\mathbf{l})$ have the same weight in the calculation of $\mathbf{E}_{ij}(\mathbf{l})$.

$\mathbf{F}_{ij}(\mathbf{l})$ is the relative frequency of occurrence of the atomic pair \mathbf{ij} in the class of distance \mathbf{l} and corresponds to:

$$f^{ij}(\mathbf{l}) = \frac{f(i, j, l)}{M_{ij}}$$

$\mathbf{F}_{xx}(\mathbf{l})$ is the relative frequency of occurrence of all the atomic pairs in the class of distance \mathbf{l} and can be expressed as:

$$f^{xx}(\mathbf{l}) = \frac{\sum_{i=1}^n \sum_{j=1}^n f(i, j, l)}{\sum_{i=1}^n \sum_{j=1}^n \sum_{l=1}^n f(i, j, l)}$$

The temperature was set to 293K, so that \mathbf{RT} is equal to 0.582 kcal/mol.

Bibliography

- Sippl, M.J. (1990) "Calculation of conformational ensembles from potentials of mean force". *J. Mol. Biol.* **213**, 859-883.
- Melo, F. and Feytmans, E. (1997) "Novel knowledge-based mean force potential at atomic level". *J. Mol. Biol.* **267**, 207-222.
- Melo, F. and Feytmans, E. (1998) "Assessment of protein structures based on the non-local energy". *J. Mol. Biol.* **277**, 1141-1152.

● Non-redundant dataset

A set of 147 protein chains with complete atomic coordinates was used to perform the calculations. This set was selected from the October 1995 release of PDB representative list (Hobohm et al., 1992; Hobohm & Sander, 1994), excluding all the proteins with duplicated or missing atoms, structural gaps, or with a total number of residues lower than 100. All these proteins share a sequence identity below 25% and have been solved by X-ray crystallography with a resolution better than 3.0 Angstroms. The list of this dataset using the PDB chain identifiers is as follows:

1AAJ_01_
1ADD_01_
1AEP_01_
1ALKA_02_A_B
1AMG_01_
1AMP_01_
1AORA_02_A_B
1AOZA_02_A_B
1ARB_01_
1ASZA_02_A_B
1ATNA_02_A_D
1AZCB_02_A_B
1BABB_04_A_B_C_D
1BBPA_04_A_B_C_D
1BET_01_
1BMTA_02_A_B
1CAUA_02_A_B
1CAUB_02_A_B
1CCR_01_
1CELA_01_A
1CEWI_01_I
1CFB_01_
1CID_01_
1CMCA_02_A_B
1COLA_02_A_B
1CPCA_04_A_B_K_L
1CPCB_04_A_B_K_L
1CRL_01_
1CSN_01_
1CTM_01_
1CTN_01_
1CHMA_02_A_B
1DAAA_02_A_B
1DHR_01_
1DPRA_02_A_B
1DSBA_01_A
1DYNA_02_A_B
1EDE_01_
1EDT_01_
1EFT_01_
1FBAA_04_A_B_C_D
1FCDC_04_A_B_C_D
1FKF_01_
1FNC_01_

1FRUB 06 A B C D E F
1GBS_01_
1GHSA 02 A B
1GKY_01_
1GMFA 02 A B
1GOF_01_
1GPR_01_
1HDCA 04 A B C D
1HEX_01_
1HFC_01_
1HMY_01_
1HNF_01_
1HSLA 02 A B
1HTMD 06 A B C D E F
1HUCB 02 A B
1HVD_01_
1IAE_01_
1IAG_01_
1IGP_01_
1KAB_01_
1L92_01_
1LBA_01_
1LKI_01_
1LPE_01_
1LTSA 07 A C D E F G H
1LTSD 07 A C D E F G H
1MSAA 04 A B C D
1NAL1 04 1 2 3 4
1NAR_01_
1NBAA 04 A B C D
1NDH_01_
1NNT_01_
1NPK_01_
1PBE_01_
1PBP_01_
1PCRH 03 H L M
1PHP_01_
1PHR_01_
1PMY_01_
1PNE_01_
1POA_01_
1POC_01_
1PPI_01_
1PRCC 04 C L M H
1PRCL 04 C L M H
1PRTC 12 A B C D E F G H I J K L
1PRTD 12 A B C D E F G H I J K L
1PSPA 02 A B
1RCB_01_
1RIBA 02 A B
1RTP1 03 1 2 3
1RVAA 02 A B
1SACA 05 A B C D E
1SCUA 04 A B D E
1SCUB 04 A B D E
1SRYA 02 A B

1STD_01_
 1TADC 03 A B C
 1TCA_01_
 1THV_01_
 1TIB_01_
 1TLCA 02 A B
 1TLK_01_
 1TPH1 02 1 2
 1TRKA 02 A B
 1TSSA 03 A B C
 1ULA_01_
 1VAAA 03 A B P
 1XNB_01_
 1YPTB 02 A B
 1YTBA 02 A B
 2ACG_01_
 2BBKH 04 H L J M
 2CPL_01_
 2CTC_01_
 2CHR_01_
 2END_01_
 2FCR_01_
 2GSTA 02 A B
 2HNQ_01_
 2HPDA 02 A B
 2KAUB 09 A B C D E F G H I
 2KAUC 09 A B C D E F G H I
 2LIGA 02 A B
 2LIV_01_
 2MGE_01_
 2OHXA 02 A B
 2PGD_02_ A
 2PIA_01_
 2RN2_01_
 2SAS_01_
 2SCPA 02 A B
 2SIL_01_
 2SNV_01_
 3AAHA 04 A B C D
 3MDDA 02 A B
 4BLMA 02 A B
 4ENL_01_
 4FXN_01_
 5P21_01_
 6FABL 02 L H
 6TAA_01_
 8ACN_01_

The syntax of the list is as follows:

XXXY NN **A B C D E**

where "**XXX**" is the PDB identifier of the protein, "Y" is the non-redundant chain ("_" means that the chain has no identifier) , "NN" is the number of protein chains considered in the calculation of the mean force potential and "**A B C D E**" is the specification of the

chain identifiers considered in the calculation.

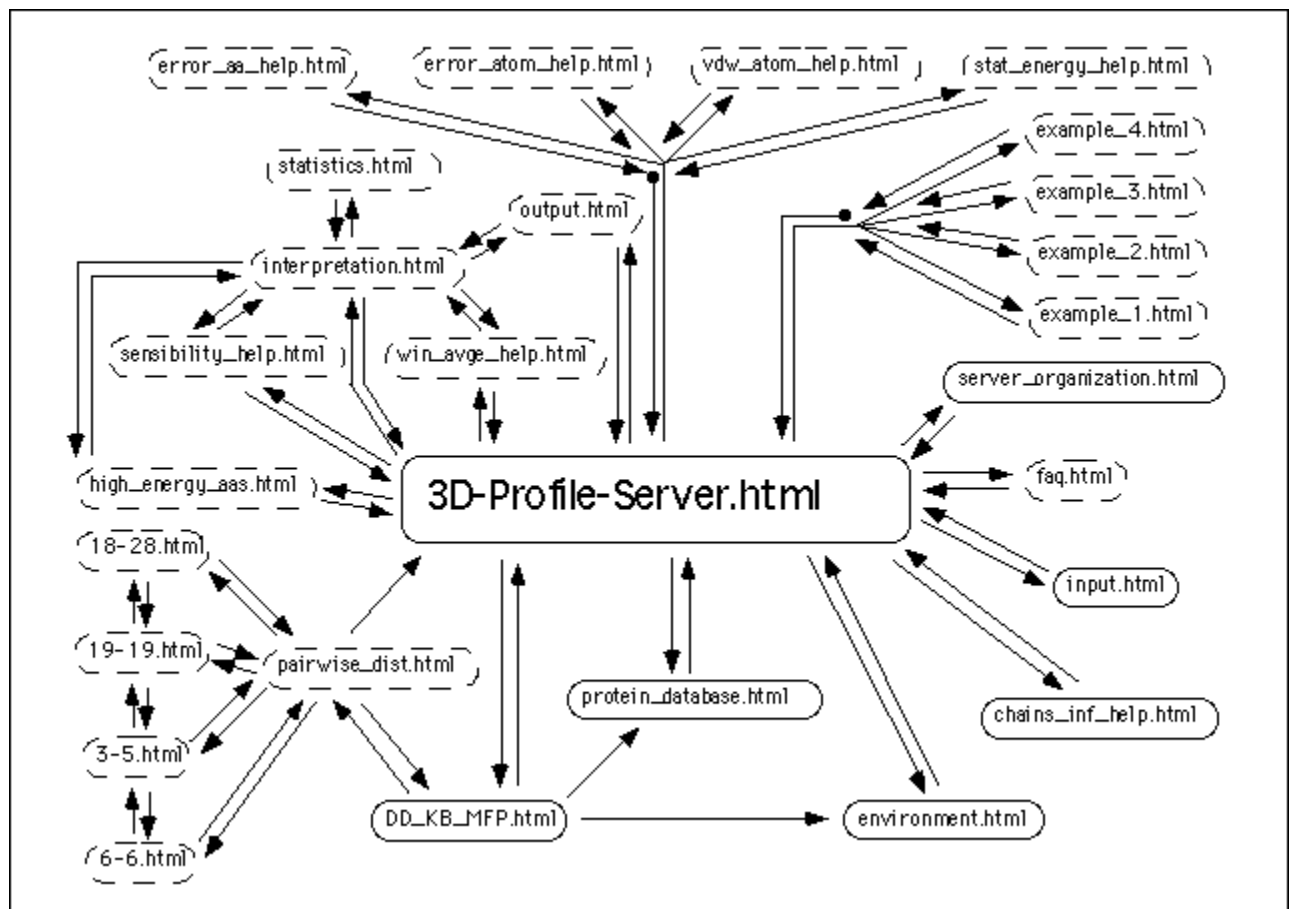
The calculation was performed on the non-redundant protein chain only (one per protein) but considering all the other heavy atoms of the other chains, if there are any.

•Bibliography

- Hobohm, U., Scharf, M., Schneider, R. & Sander, C. (1992) "Selection of representative protein data sets". *Protein Science* **1**, 409-417.
- Hobohm, U. & Sander, C. (1994) "Enlarged representative set of protein structures". *Protein Science* **3**, 522.

•Server organization

This server is organized as shown in the figure below. Continuous line frames represent already functional html pages linked with the server and dashed line frames represent not yet available pages, but they will be implemented soon. The link directions are shown by the arrows.



● How can I use ANOLEA ?

You must submit the Cartesian coordinates for each heavy atom of a protein (monomer or multimer). ANOLEA will calculate a non-local energy profile and you will obtain a list of the ["high energy"](#) amino acids for one chain of this protein, depending on the [sensibility](#) and the [window average](#) used to perform the calculations.

● What is the usefulness of ANOLEA ?

The high energy amino acid list obtained after calculate the energy profile could represent some of the following features:

- errors in a model or in an experimentally solved protein structure. [\(example\)](#)
- potential interacting zones between subunits in a multimeric protein. [\(example\)](#)
- potential interacting zones between different proteins [\(example\)](#)

● Window average help

As default, a moving window average of 5 residues is used. Using this value, a threshold of zero E/kT units is appropriated to define a high energy region. When changing the window average value, the user should be careful about the interpretations and the threshold defined. We provide the flexibility to fix the values for the window average and threshold. This is intended to be used carefully.

The user could want to see the energy profile using a window average of 1, after some regions have been identified as suspicious ones. The value of 1 in the window average parameter means that no average is used, thus each residue can be examined independently of its neighbors. With this analysis, the residues that are mainly responsible for the high energy regions obtained can be pointed out.

● Threshold help

The threshold value can be changed to any real number. As default, a value of zero is defined. In some cases, i.e. comparative models assessment, the user could be interested only in checking the existence of misalignments or regions that have very high energy scores. In that case, the threshold could be conveniently changed in such a way to display only those regions in the profile or in the graphic output.

● Chain information help

Suppose that you want to submit a protein that contains two chains, named A and B. Then you have four possibilities to obtain an energy profile:

1.- To fill the [chain information box](#) with the next sequence of characters:

AB

With this information, the profile will be performed on chain **A**, but considering all the atom contacts with the chain **B**.

2.- To fill the [chain information box](#) with the next sequence of characters:

BA

With this information, the profile will be performed on chain **B**, but considering all the atom contacts with the chain **A**.

3.- To fill the [chain information box](#) with the next sequence of characters:

A

With this information, the profile will be performed only on chain **A**, without considering the atoms of chain **B**, even if you submit the atom coordinates of this chain.

4.- To fill the [chain information box](#) with the next sequence of characters:

B

With this information, the profile will be performed only on chain **B**, without considering the atoms of chain **A**, even if you submit the atom coordinates of this chain.

If you are submitting a unique protein chain without identifier, you must leave the [chain information box](#) empty.

If you are submitting a unique protein chain with identifier, you must write it in the [chain information box](#).