



# Full wwPDB NMR Structure Validation Report ⓘ

May 29, 2020 – 07:45 am BST

PDB ID : 5MHD  
Title : Biosynthetic engineered A22S-B3K-B31R human insulin monomer structure in water/acetonitrile solutions.  
Authors : Bocian, W.; Kozerski, L.; Bednarek, E.; Sitkowski, J.  
Deposited on : 2016-11-24

This is a Full wwPDB NMR Structure Validation Report for a publicly released PDB entry.

We welcome your comments at [validation@mail.wwpdb.org](mailto:validation@mail.wwpdb.org)

A user guide is available at

<https://www.wwpdb.org/validation/2017/NMRValidationReportHelp>

with specific help available everywhere you see the ⓘ symbol.

---

The following versions of software and data (see [references ⓘ](#)) were used in the production of this report:

Cyrange : Kirchner and Güntert (2011)  
NmrClust : Kelley et al. (1996)  
MolProbity : 4.02b-467  
Percentile statistics : 20191225.v01 (using entries in the PDB archive December 25th 2019)  
RCI : v\_1n\_11\_5\_13\_A (Berjanski et al., 2005)  
PANAV : Wang et al. (2010)  
ShiftChecker : 2.11  
Ideal geometry (proteins) : Engh & Huber (2001)  
Ideal geometry (DNA, RNA) : Parkinson et al. (1996)  
Validation Pipeline (wwPDB-VP) : 2.11

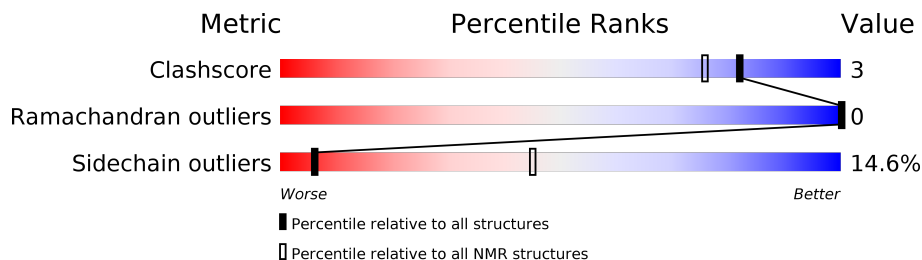
# 1 Overall quality at a glance

The following experimental techniques were used to determine the structure:

*SOLUTION NMR*

The overall completeness of chemical shifts assignment is 48%.

Percentile scores (ranging between 0-100) for global validation metrics of the entry are shown in the following graphic. The table shows the number of entries on which the scores are based.



Metric	Whole archive (#Entries)	NMR archive (#Entries)
Clashscore	158937	12864
Ramachandran outliers	154571	11451
Sidechain outliers	154315	11428

The table below summarises the geometric issues observed across the polymeric chains and their fit to the experimental data. The red, orange, yellow and green segments indicate the fraction of residues that contain outliers for  $\geq 3$ , 2, 1 and 0 types of geometric quality criteria. A cyan segment indicates the fraction of residues that are not part of the well-defined cores, and a grey segment represents the fraction of residues that are not modelled. The numeric value for each fraction is indicated below the corresponding segment, with a dot representing fractions  $\leq 5\%$

Mol	Chain	Length	Quality of chain
1	A	22	 68% 14% 5% 14%
2	B	31	 52% 6% 42%

## 2 Ensemble composition and analysis i

This entry contains 20 models. Model 16 is the overall representative, medoid model (most similar to other models). The authors have identified model 1 as representative, based on the following criterion: *fewest violations*.

The following residues are included in the computation of the global validation metrics.

Well-defined (core) protein residues			
Well-defined core	Residue range (total)	Backbone RMSD (Å)	Medoid model
1	A:2-A:20, B:4-B:19, B:23-B:24 (37)	0.12	16

Ill-defined regions of proteins are excluded from the global statistics.

Ligands and non-protein polymers are included in the analysis.

The models can be grouped into 2 clusters and 1 single-model cluster was found.

Cluster number	Models
1	1, 3, 4, 5, 6, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 20
2	2, 7, 8
Single-model clusters	19

### 3 Entry composition

There are 2 unique types of molecules in this entry. The entry contains 829 atoms, of which 406 are hydrogens and 0 are deuteriums.

- Molecule 1 is a protein called Insulin.

Mol	Chain	Residues	Atoms					Trace	
			Total	C	H	N	O		S
1	A	22	323	102	154	26	37	4	0

There is a discrepancy between the modelled and reference sequences:

Chain	Residue	Modelled	Actual	Comment	Reference
A	22	SER	-	expression tag	UNP P01308

- Molecule 2 is a protein called Insulin.

Mol	Chain	Residues	Atoms					Trace	
			Total	C	H	N	O		S
2	B	31	506	166	252	44	42	2	0

There is a discrepancy between the modelled and reference sequences:

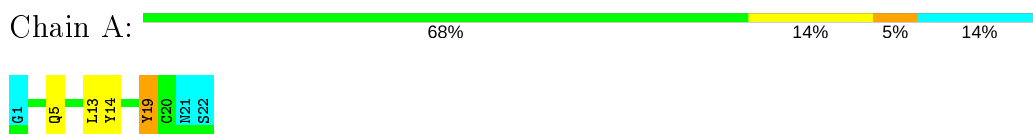
Chain	Residue	Modelled	Actual	Comment	Reference
B	3	LYS	ASN	conflict	UNP P01308

## 4 Residue-property plots [i](#)

### 4.1 Average score per residue in the NMR ensemble

These plots are provided for all protein, RNA and DNA chains in the entry. The first graphic is the same as shown in the summary in section 1 of this report. The second graphic shows the sequence where residues are colour-coded according to the number of geometric quality criteria for which they contain at least one outlier: green = 0, yellow = 1, orange = 2 and red = 3 or more. Stretches of 2 or more consecutive residues without any outliers are shown as green connectors. Residues which are classified as ill-defined in the NMR ensemble, are shown in cyan with an underline colour-coded according to the previous scheme. Residues which were present in the experimental sample, but not modelled in the final structure are shown in grey.

- Molecule 1: Insulin



- Molecule 2: Insulin

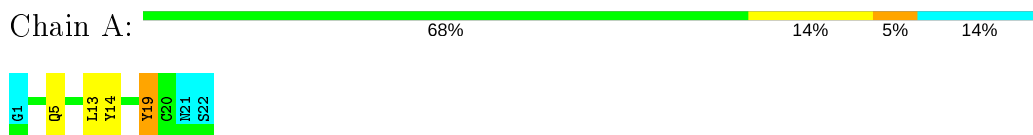


### 4.2 Scores per residue for each member of the ensemble

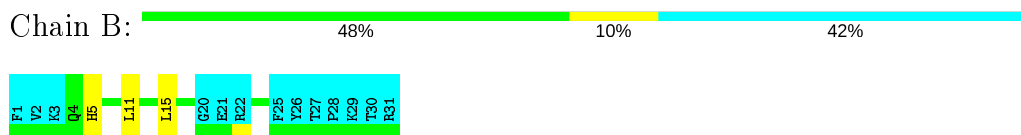
Colouring as in section 4.1 above.

#### 4.2.1 Score per residue for model 1

- Molecule 1: Insulin

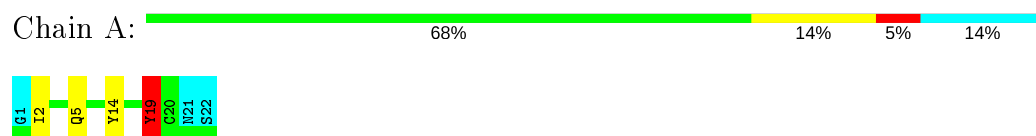


- Molecule 2: Insulin



### 4.2.2 Score per residue for model 2

- Molecule 1: Insulin



- Molecule 2: Insulin



### 4.2.3 Score per residue for model 3

- Molecule 1: Insulin

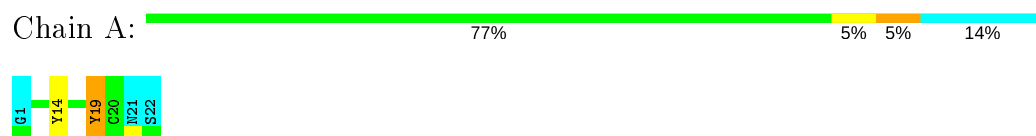


- Molecule 2: Insulin

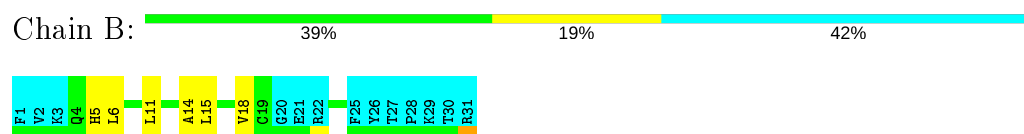


### 4.2.4 Score per residue for model 4

- Molecule 1: Insulin

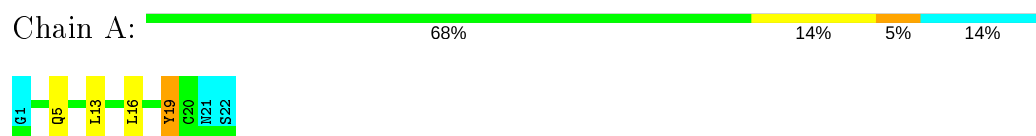


- Molecule 2: Insulin



### 4.2.5 Score per residue for model 5

- Molecule 1: Insulin



- Molecule 2: Insulin



### 4.2.6 Score per residue for model 6

- Molecule 1: Insulin



- Molecule 2: Insulin

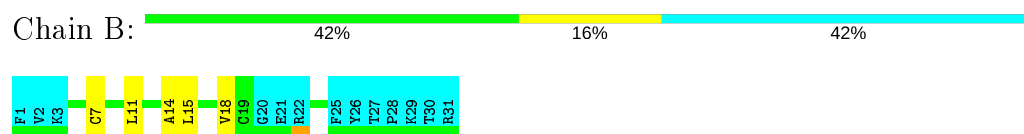


### 4.2.7 Score per residue for model 7

- Molecule 1: Insulin



- Molecule 2: Insulin

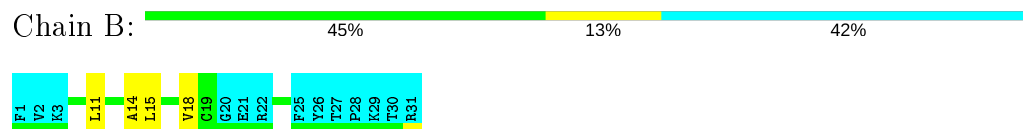


### 4.2.8 Score per residue for model 8

- Molecule 1: Insulin



- Molecule 2: Insulin

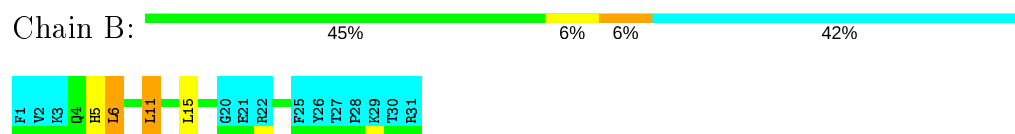


### 4.2.9 Score per residue for model 9

- Molecule 1: Insulin

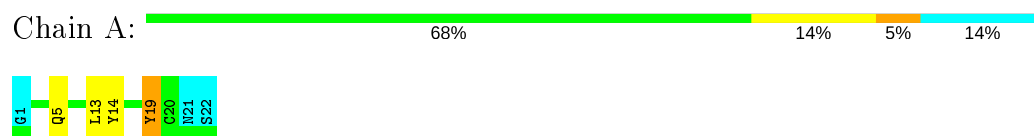


- Molecule 2: Insulin

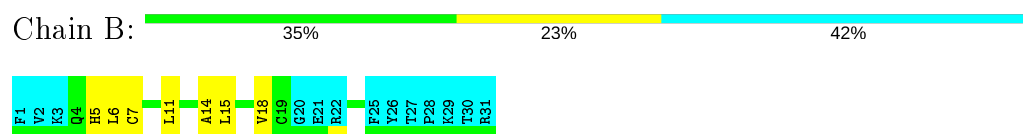


### 4.2.10 Score per residue for model 10

- Molecule 1: Insulin



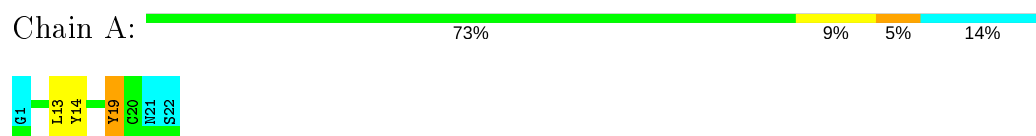
- Molecule 2: Insulin





### 4.2.11 Score per residue for model 11

- Molecule 1: Insulin

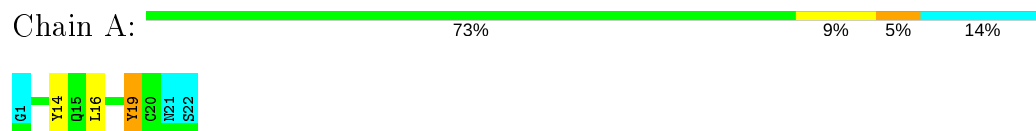


- Molecule 2: Insulin



### 4.2.12 Score per residue for model 12

- Molecule 1: Insulin



- Molecule 2: Insulin

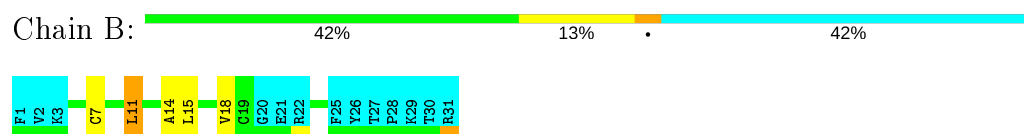


### 4.2.13 Score per residue for model 13

- Molecule 1: Insulin

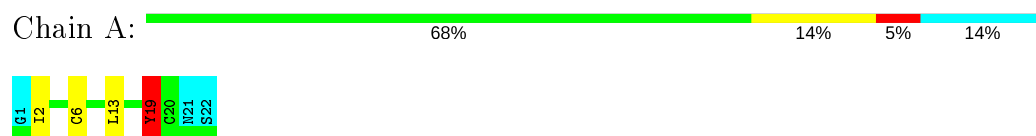


- Molecule 2: Insulin

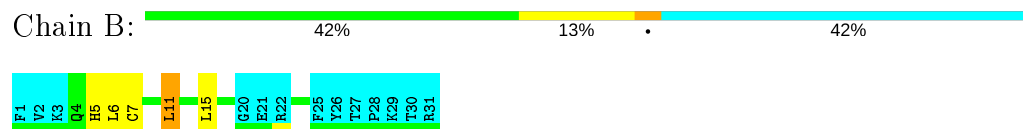


#### 4.2.14 Score per residue for model 14

- Molecule 1: Insulin



- Molecule 2: Insulin



#### 4.2.15 Score per residue for model 15

- Molecule 1: Insulin



- Molecule 2: Insulin



#### 4.2.16 Score per residue for model 16 (medoid)

- Molecule 1: Insulin

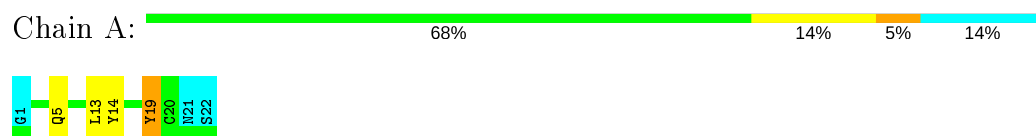


- Molecule 2: Insulin

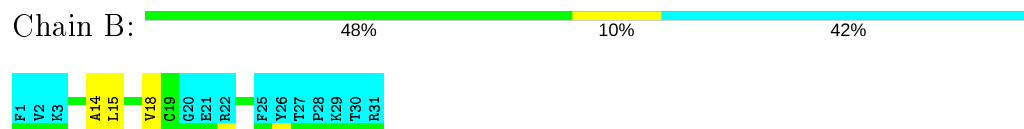


#### 4.2.17 Score per residue for model 17

- Molecule 1: Insulin

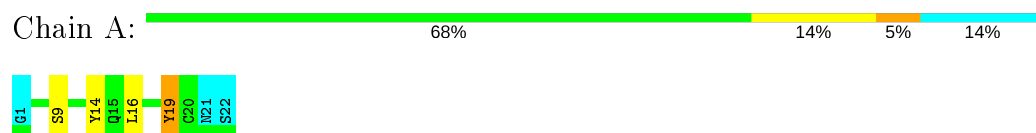


- Molecule 2: Insulin

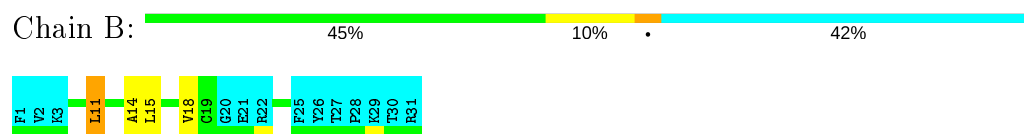


#### 4.2.18 Score per residue for model 18

- Molecule 1: Insulin

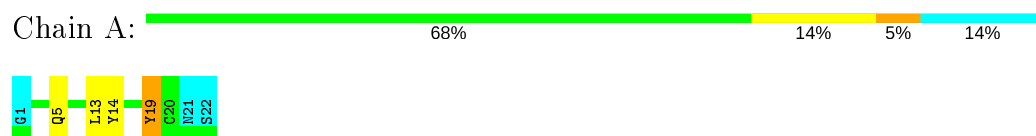


- Molecule 2: Insulin



#### 4.2.19 Score per residue for model 19

- Molecule 1: Insulin

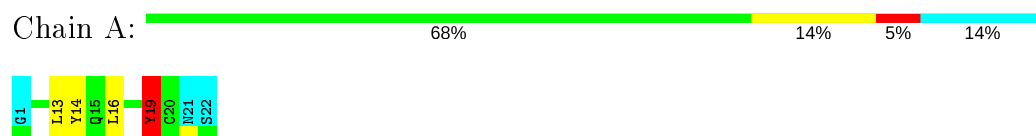


- Molecule 2: Insulin



#### 4.2.20 Score per residue for model 20

- Molecule 1: Insulin



- Molecule 2: Insulin



## 5 Refinement protocol and experimental data overview (i)

The models were refined using the following method: *simulated annealing*.

Of the 100 calculated structures, 20 were deposited, based on the following criterion: *20 structures for lowest energy*.

The following table shows the software used for structure solution, optimisation and refinement.

Software name	Classification	Version
DYANA	refinement	
Amber	structure calculation	14

The following table shows chemical shift validation statistics as aggregates over all chemical shift files. Detailed validation can be found in section 6 of this report.

Chemical shift file(s)	input_cs.cif
Number of chemical shift lists	1
Total number of shifts	351
Number of shifts mapped to atoms	351
Number of unparsed shifts	0
Number of shifts with mapping errors	0
Number of shifts with mapping warnings	0
Assignment completeness (well-defined parts)	48%

No validations of the models with respect to experimental NMR restraints is performed at this time.

COVALENT-GEOMETRY INFOmissingINFO

### 5.1 Too-close contacts (i)

In the following table, the Non-H and H(model) columns list the number of non-hydrogen atoms and hydrogen atoms in each chain respectively. The H(added) column lists the number of hydrogen atoms added and optimized by MolProbity. The Clashes column lists the number of clashes averaged over the ensemble.

Mol	Chain	Non-H	H(model)	H(added)	Clashes
1	A	150	138	138	1±1
2	B	138	132	132	1±1
All	All	5760	5400	5400	28

The all-atom clashscore is defined as the number of clashes found per 1000 atoms (including hydrogen atoms). The all-atom clashscore for this structure is 3.

All unique clashes are listed below, sorted by their clash magnitude.

Atom-1	Atom-2	Clash(Å)	Distance(Å)	Models	
				Worst	Total
2:B:14:ALA:O	2:B:18:VAL:HG23	0.52	2.05	7	7
1:A:13:LEU:HG	2:B:18:VAL:HG22	0.49	1.85	7	1
1:A:16:LEU:HD11	2:B:11:LEU:CD2	0.48	2.38	13	7
1:A:13:LEU:HB3	2:B:18:VAL:HG22	0.47	1.86	8	1
1:A:16:LEU:HB2	2:B:18:VAL:HG21	0.46	1.86	8	1
1:A:2:ILE:CG2	1:A:19:TYR:CZ	0.43	3.02	3	6
1:A:11:CYS:SG	1:A:16:LEU:HD21	0.42	2.54	3	1
1:A:6:CYS:CB	2:B:11:LEU:HD21	0.42	2.44	14	1
1:A:19:TYR:CE1	2:B:15:LEU:HD22	0.41	2.50	20	1
1:A:2:ILE:HG21	1:A:19:TYR:CE2	0.41	2.51	3	1
1:A:16:LEU:HD11	2:B:11:LEU:HD23	0.40	1.93	18	1

## 5.2 Torsion angles [i](#)

### 5.2.1 Protein backbone [i](#)

In the following table, the Percentiles column shows the percent Ramachandran outliers of the chain as a percentile score with respect to all PDB entries followed by that with respect to all NMR entries. The Analysed column shows the number of residues for which the backbone conformation was analysed and the total number of residues.

Mol	Chain	Analysed	Favoured	Allowed	Outliers	Percentiles	
1	A	19/22 (86%)	17±1 (92±4%)	2±1 (8±4%)	0±0 (0±0%)	100	100
2	B	18/31 (58%)	18±1 (98±3%)	0±1 (2±3%)	0±0 (0±0%)	100	100
All	All	740/1060 (70%)	702 (95%)	38 (5%)	0 (0%)	100	100

There are no Ramachandran outliers.

### 5.2.2 Protein sidechains [i](#)

In the following table, the Percentiles column shows the percent sidechain outliers of the chain as a percentile score with respect to all PDB entries followed by that with respect to all NMR entries. The Analysed column shows the number of residues for which the sidechain conformation was analysed and the total number of residues.

Mol	Chain	Analysed	Rotameric	Outliers	Percentiles	
1	A	19/21 (90%)	17±1 (88±4%)	2±1 (12±4%)	8	51
2	B	15/27 (56%)	12±1 (82±8%)	3±1 (18±8%)	4	38
All	All	680/960 (71%)	581 (85%)	99 (15%)	6	45

All 10 unique residues with a non-rotameric sidechain are listed below. They are sorted by the frequency of occurrence in the ensemble.

Mol	Chain	Res	Type	Models (Total)
2	B	15	LEU	19
2	B	11	LEU	18
1	A	14	TYR	18
1	A	13	LEU	15
1	A	5	GLN	11
2	B	7	CYS	7
2	B	5	HIS	5
2	B	6	LEU	4
1	A	17	GLU	1
1	A	9	SER	1

### 5.2.3 RNA [i](#)

There are no RNA molecules in this entry.

### 5.3 Non-standard residues in protein, DNA, RNA chains [i](#)

There are no non-standard protein/DNA/RNA residues in this entry.

### 5.4 Carbohydrates [i](#)

There are no carbohydrates in this entry.

### 5.5 Ligand geometry [i](#)

There are no ligands in this entry.

### 5.6 Other polymers [i](#)

There are no such molecules in this entry.

### 5.7 Polymer linkage issues [i](#)

There are no chain breaks in this entry.

## 6 Chemical shift validation

The completeness of assignment taking into account all chemical shift lists is 48% for the well-defined parts and 46% for the entire structure.

### 6.1 Chemical shift list 1

File name: input\_cs.cif

Chemical shift list name: *SKRR\_ChemicalShift.txt*

#### 6.1.1 Bookkeeping

The following table shows the results of parsing the chemical shift list and reports the number of nuclei with statistically unusual chemical shifts.

Total number of shifts	351
Number of shifts mapped to atoms	351
Number of unparsed shifts	0
Number of shifts with mapping errors	0
Number of shifts with mapping warnings	0
Number of shift outliers (ShiftChecker)	0

#### 6.1.2 Chemical shift referencing

No chemical shift referencing corrections were calculated (not enough data).

#### 6.1.3 Completeness of resonance assignments

The following table shows the completeness of the chemical shift assignments for the well-defined regions of the structure. The overall completeness is 48%, i.e. 207 atoms were assigned a chemical shift out of a possible 431. 0 out of 9 assigned methyl groups (LEU and VAL) were assigned stereospecifically.

	Total	<sup>1</sup> H	<sup>13</sup> C	<sup>15</sup> N
Backbone	73/185 (39%)	73/74 (99%)	0/74 (0%)	0/37 (0%)
Sidechain	114/199 (57%)	114/116 (98%)	0/79 (0%)	0/4 (0%)
Aromatic	20/47 (43%)	20/25 (80%)	0/20 (0%)	0/2 (0%)
Overall	207/431 (48%)	207/215 (96%)	0/173 (0%)	0/43 (0%)

The following table shows the completeness of the chemical shift assignments for the full structure. The overall completeness is 46%, i.e. 295 atoms were assigned a chemical shift out of a possible 646. 0 out of 10 assigned methyl groups (LEU and VAL) were assigned stereospecifically.



	Total	<sup>1</sup> H	<sup>13</sup> C	<sup>15</sup> N
Backbone	101/263 (38%)	101/105 (96%)	0/106 (0%)	0/52 (0%)
Sidechain	163/310 (53%)	163/183 (89%)	0/114 (0%)	0/13 (0%)
Aromatic	31/73 (42%)	31/39 (79%)	0/32 (0%)	0/2 (0%)
Overall	295/646 (46%)	295/327 (90%)	0/252 (0%)	0/67 (0%)

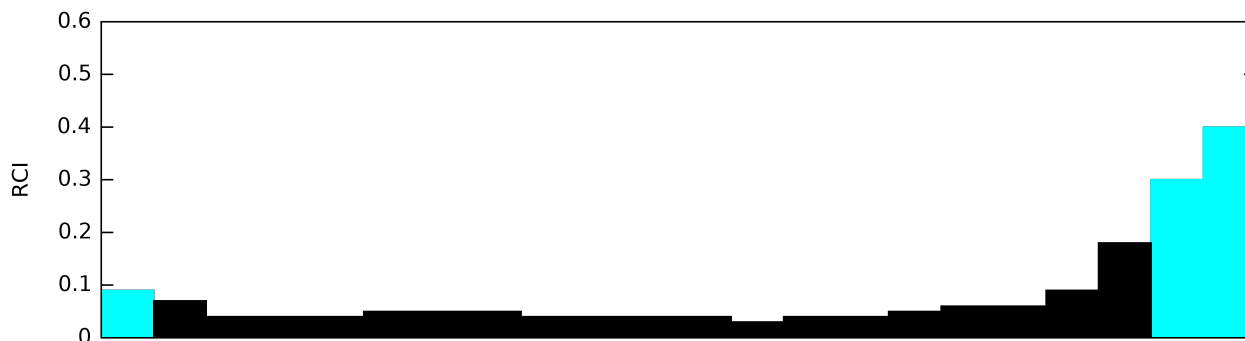
### 6.1.4 Statistically unusual chemical shifts [i](#)

There are no statistically unusual chemical shifts.

### 6.1.5 Random Coil Index (RCI) plots [i](#)

The images below report *random coil index* values for the protein chains in the structure. The height of each bar gives a probability of a given residue to be disordered, as predicted from the available chemical shifts and the amino acid sequence. A value above 0.2 is an indication of significant predicted disorder. The colour of the bar shows whether the residue is in the well-defined core (black) or in the ill-defined residue ranges (cyan), as described in section 2 on ensemble composition.

Random coil index (RCI) for chain A:



Random coil index (RCI) for chain B:

