



wwPDB NMR Structure Validation Summary Report ⓘ

May 29, 2020 – 07:57 am BST

PDB ID : 5O9B
Title : Solution NMR structure of human GATA2 C-terminal zinc finger domain
Authors : Nurmohamed, S.S.; Broadhurst, R.W.; May, G.; Enver, T.
Deposited on : 2017-06-16

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

We welcome your comments at 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
buster-report : 1.1.7 (2018)
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

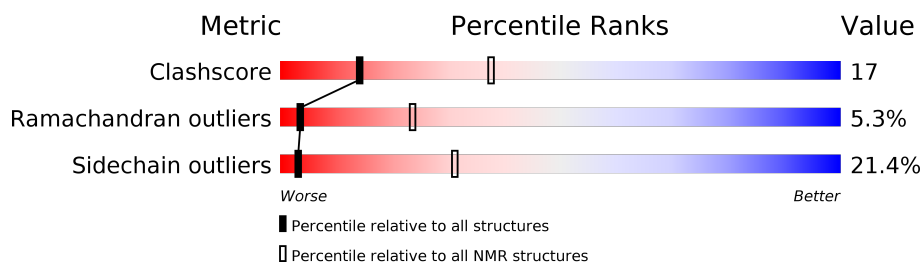
1 Overall quality at a glance

The following experimental techniques were used to determine the structure:

SOLUTION NMR

The overall completeness of chemical shifts assignment is 76%.

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 ≥ 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	68	<div> <div>21%</div> <div>32%</div> <div>.</div> <div>44%</div> </div>

2 Ensemble composition and analysis

This entry contains 20 models. Model 2 is the overall representative, medoid model (most similar to other models). The authors have identified model 1 as representative, based on the following criterion: *lowest energy*.

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:5-A:42 (38)	0.15	2

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 3 clusters and 2 single-model clusters were found.

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

3 Entry composition

There are 2 unique types of molecules in this entry. The entry contains 1056 atoms, of which 519 are hydrogens and 0 are deuteriums.

- Molecule 1 is a protein called Endothelial transcription factor GATA-2.

Mol	Chain	Residues	Atoms						Trace
1	A	68	Total	C	H	N	O	S	0
			1055	326	519	109	94	7	

There are 19 discrepancies between the modelled and reference sequences:

Chain	Residue	Modelled	Actual	Comment	Reference
A	-18	MET	-	initiating methionine	UNP P23769
A	-17	ALA	-	expression tag	UNP P23769
A	-16	HIS	-	expression tag	UNP P23769
A	-15	HIS	-	expression tag	UNP P23769
A	-14	HIS	-	expression tag	UNP P23769
A	-13	HIS	-	expression tag	UNP P23769
A	-12	HIS	-	expression tag	UNP P23769
A	-11	HIS	-	expression tag	UNP P23769
A	-10	SER	-	expression tag	UNP P23769
A	-9	SER	-	expression tag	UNP P23769
A	-8	GLY	-	expression tag	UNP P23769
A	-7	LEU	-	expression tag	UNP P23769
A	-6	GLU	-	expression tag	UNP P23769
A	-5	VAL	-	expression tag	UNP P23769
A	-4	LEU	-	expression tag	UNP P23769
A	-3	PHE	-	expression tag	UNP P23769
A	-2	GLN	-	expression tag	UNP P23769
A	-1	GLY	-	expression tag	UNP P23769
A	0	PRO	-	expression tag	UNP P23769

- Molecule 2 is ZINC ION (three-letter code: ZN) (formula: Zn) (labeled as "Ligand of Interest" by author).

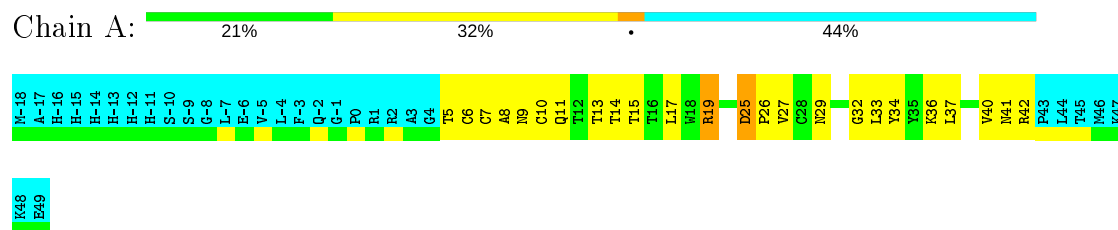
Mol	Chain	Residues	Atoms	
2	A	1	Total	Zn
			1	1

4 Residue-property plots

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.

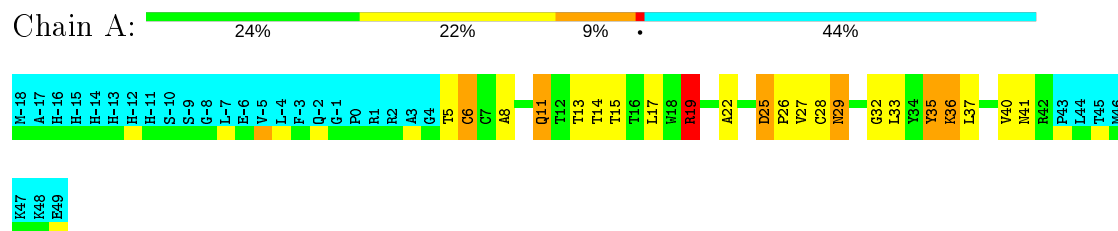
- Molecule 1: Endothelial transcription factor GATA-2



4.2 Residue scores for the representative (medoid) model from the NMR ensemble

The representative model is number 2. Colouring as in section 4.1 above.

- Molecule 1: Endothelial transcription factor GATA-2



5 Refinement protocol and experimental data overview

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

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

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

Software name	Classification	Version
ARIA	refinement	2.3
ARIA	structure calculation	2.3

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	594
Number of shifts mapped to atoms	594
Number of unparsed shifts	0
Number of shifts with mapping errors	0
Number of shifts with mapping warnings	0
Assignment completeness (well-defined parts)	76%

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

COVALENT-GEOMETRY INFOmissingINFO

5.1 Too-close contacts

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	295	278	277	10±3
All	All	5920	5560	5540	199

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 17.

5 of 52 unique clashes are listed below, sorted by their clash magnitude.

Atom-1	Atom-2	Clash(Å)	Distance(Å)	Models	
				Worst	Total
1:A:17:LEU:HG	1:A:19:ARG:NH1	0.70	2.02	18	5
1:A:33:LEU:O	1:A:37:LEU:HG	0.70	1.86	6	17
1:A:32:GLY:O	1:A:36:LYS:HD2	0.69	1.88	16	3
1:A:36:LYS:HE2	1:A:36:LYS:N	0.67	2.05	14	1
1:A:32:GLY:O	1:A:36:LYS:HD3	0.67	1.89	3	15

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	38/68 (56%)	33±1 (87±3%)	3±1 (7±3%)	2±1 (5±1%)	3	23
All	All	760/1360 (56%)	664 (87%)	56 (7%)	40 (5%)	3	23

All 3 unique Ramachandran outliers are listed below. They are sorted by the frequency of occurrence in the ensemble.

Mol	Chain	Res	Type	Models (Total)
1	A	8	ALA	19
1	A	11	GLN	19
1	A	9	ASN	2

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	33/58 (57%)	26±1 (79±4%)	7±1 (21±4%)	3	31
All	All	660/1160 (57%)	519 (79%)	141 (21%)	3	31

5 of 15 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)
1	A	27	VAL	20
1	A	19	ARG	20
1	A	29	ASN	17
1	A	25	ASP	16
1	A	40	VAL	16

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](#)

Of 1 ligands modelled in this entry, 1 is monoatomic - leaving 0 for Mogul analysis.

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 76% for the well-defined parts and 61% for the entire structure.

6.1 Chemical shift list 1

File name: input_cs.cif

Chemical shift list name: CF_170605_csdep.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	594
Number of shifts mapped to atoms	594
Number of unparsed shifts	0
Number of shifts with mapping errors	0
Number of shifts with mapping warnings	0
Number of shift outliers (ShiftChecker)	7

6.1.2 Chemical shift referencing

The following table shows the suggested chemical shift referencing corrections.

Nucleus	# values	Correction \pm precision, ppm	Suggested action
$^{13}\text{C}_\alpha$	61	-0.20 ± 0.21	None needed (< 0.5 ppm)
$^{13}\text{C}_\beta$	55	0.04 ± 0.18	None needed (< 0.5 ppm)
$^{13}\text{C}'$	0	—	None (insufficient data)
^{15}N	55	-0.03 ± 0.46	None needed (< 0.5 ppm)

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 76%, i.e. 343 atoms were assigned a chemical shift out of a possible 452. 0 out of 5 assigned methyl groups (LEU and VAL) were assigned stereospecifically.

	Total	^1H	^{13}C	^{15}N
Backbone	148/188 (79%)	74/75 (99%)	38/76 (50%)	36/37 (97%)
Sidechain	163/228 (71%)	100/133 (75%)	56/78 (72%)	7/17 (41%)

Continued on next page...

Continued from previous page...

	Total	¹ H	¹³ C	¹⁵ N
Aromatic	32/36 (89%)	16/18 (89%)	15/15 (100%)	1/3 (33%)
Overall	343/452 (76%)	190/226 (84%)	109/169 (64%)	44/57 (77%)

6.1.4 Statistically unusual chemical shifts [i](#)

The following table lists the statistically unusual chemical shifts. These are statistical measures, and large deviations from the mean do not necessarily imply incorrect assignments. Molecules containing paramagnetic centres or hemes are expected to give rise to anomalous chemical shifts.

Mol	Chain	Res	Type	Atom	Shift, ppm	Expected range, ppm	Z-score
1	A	15	THR	HG22	-0.83	2.29 – -0.01	-8.6
1	A	15	THR	HG21	-0.83	2.29 – -0.01	-8.6
1	A	15	THR	HG23	-0.83	2.29 – -0.01	-8.6
1	A	15	THR	CG2	15.27	27.15 – 15.95	-5.6
1	A	5	THR	HG23	-0.10	2.29 – -0.01	-5.4
1	A	5	THR	HG22	-0.10	2.29 – -0.01	-5.4
1	A	5	THR	HG21	-0.10	2.29 – -0.01	-5.4

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

The image below reports *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:

