



# wwPDB X-ray Structure Validation Summary Report ⓘ

May 14, 2020 – 11:06 am BST

PDB ID : 4CU9  
Title : Unravelling the multiple functions of the architecturally intricate *Streptococcus pneumoniae* beta-galactosidase, BgaA  
Authors : Singh, A.K.; Pluvinae, B.; Higgins, M.A.; Dalia, A.B.; Flynn, M.; Lloyd, A.R.; Weiser, J.N.; Stubbs, K.A.; Boraston, A.B.; King, S.J.  
Deposited on : 2014-03-17  
Resolution : 1.83 Å(reported)

This is a wwPDB X-ray Structure Validation Summary 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/XrayValidationReportHelp>

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:

MolProbity : 4.02b-467  
Mogul : 1.8.5 (274361), CSD as541be (2020)  
Xtriage (Phenix) : 1.13  
EDS : 2.11  
Percentile statistics : 20191225.v01 (using entries in the PDB archive December 25th 2019)  
Refmac : 5.8.0158  
CCP4 : 7.0.044 (Gargrove)  
Ideal geometry (proteins) : Engh & Huber (2001)  
Ideal geometry (DNA, RNA) : Parkinson et al. (1996)  
Validation Pipeline (wwPDB-VP) : 2.11

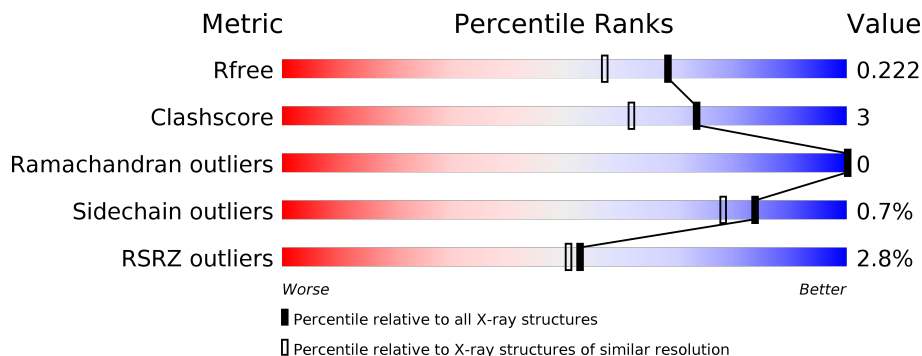
# 1 Overall quality at a glance

The following experimental techniques were used to determine the structure:

*X-RAY DIFFRACTION*

The reported resolution of this entry is 1.83 Å.

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)	Similar resolution (#Entries, resolution range(Å))
$R_{free}$	130704	4003 (1.86-1.82)
Clashscore	141614	4233 (1.86-1.82)
Ramachandran outliers	138981	4185 (1.86-1.82)
Sidechain outliers	138945	4186 (1.86-1.82)
RSRZ outliers	127900	3957 (1.86-1.82)

The table below summarises the geometric issues observed across the polymeric chains and their fit to the electron density. The red, orange, yellow and green segments on the lower bar indicate the fraction of residues that contain outliers for  $\geq 3$ , 2, 1 and 0 types of geometric quality criteria respectively. 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  $\leq 5\%$ . The upper red bar (where present) indicates the fraction of residues that have poor fit to the electron density. The numeric value is given above the bar.

Mol	Chain	Length	Quality of chain
1	A	186	
1	B	186	

## 2 Entry composition [i](#)

There are 4 unique types of molecules in this entry. The entry contains 3108 atoms, of which 0 are hydrogens and 0 are deuteriums.

In the tables below, the ZeroOcc column contains the number of atoms modelled with zero occupancy, the AltConf column contains the number of residues with at least one atom in alternate conformation and the Trace column contains the number of residues modelled with at most 2 atoms.

- Molecule 1 is a protein called BETA-GALACTOSIDASE.

Mol	Chain	Residues	Atoms					ZeroOcc	AltConf	Trace
			Total	C	N	O	S			
1	A	173	1378	859	244	271	4	4	1	0
1	B	185	1460	910	261	284	5	0	1	0

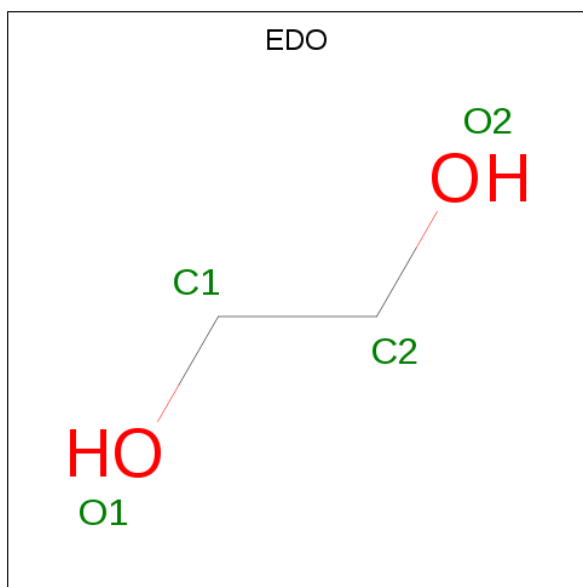
There are 20 discrepancies between the modelled and reference sequences:

Chain	Residue	Modelled	Actual	Comment	Reference
A	1813	LEU	-	expression tag	UNP Q8DQP4
A	1814	VAL	-	expression tag	UNP Q8DQP4
A	1815	PRO	-	expression tag	UNP Q8DQP4
A	1816	ARG	-	expression tag	UNP Q8DQP4
A	1817	GLY	-	expression tag	UNP Q8DQP4
A	1818	SER	-	expression tag	UNP Q8DQP4
A	1819	HIS	-	expression tag	UNP Q8DQP4
A	1820	MET	-	expression tag	UNP Q8DQP4
A	1821	ASN	-	expression tag	UNP Q8DQP4
A	1972	ALA	VAL	conflict	UNP Q8DQP4
B	1813	LEU	-	expression tag	UNP Q8DQP4
B	1814	VAL	-	expression tag	UNP Q8DQP4
B	1815	PRO	-	expression tag	UNP Q8DQP4
B	1816	ARG	-	expression tag	UNP Q8DQP4
B	1817	GLY	-	expression tag	UNP Q8DQP4
B	1818	SER	-	expression tag	UNP Q8DQP4
B	1819	HIS	-	expression tag	UNP Q8DQP4
B	1820	MET	-	expression tag	UNP Q8DQP4
B	1821	ASN	-	expression tag	UNP Q8DQP4
B	1972	ALA	VAL	conflict	UNP Q8DQP4

- Molecule 2 is CALCIUM ION (three-letter code: CA) (formula: Ca).

Mol	Chain	Residues	Atoms	ZeroOcc	AltConf
2	B	1	Total Ca 1 1	0	0
2	A	1	Total Ca 1 1	0	0

- Molecule 3 is 1,2-ETHANEDIOL (three-letter code: EDO) (formula: C<sub>2</sub>H<sub>6</sub>O<sub>2</sub>).



Mol	Chain	Residues	Atoms	ZeroOcc	AltConf
3	A	1	Total C O 4 2 2	0	0
3	B	1	Total C O 4 2 2	0	0
3	B	1	Total C O 4 2 2	0	0
3	B	1	Total C O 4 2 2	0	0
3	B	1	Total C O 4 2 2	0	0

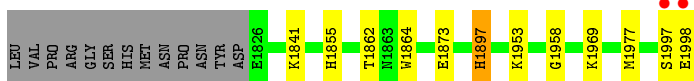
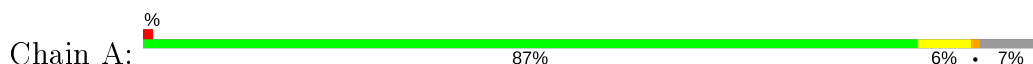
- Molecule 4 is water.

Mol	Chain	Residues	Atoms	ZeroOcc	AltConf
4	A	138	Total O 138 138	0	0
4	B	110	Total O 110 110	0	0

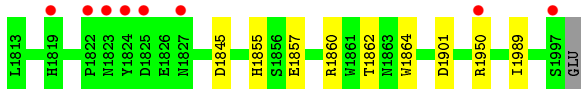
### 3 Residue-property plots [i](#)

These plots are drawn for all protein, RNA and DNA chains in the entry. The first graphic for a chain summarises the proportions of the various outlier classes displayed in the second graphic. The second graphic shows the sequence view annotated by issues in geometry and electron density. Residues are color-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. A red dot above a residue indicates a poor fit to the electron density ( $RSRZ > 2$ ). Stretches of 2 or more consecutive residues without any outlier are shown as a green connector. Residues present in the sample, but not in the model, are shown in grey.

- Molecule 1: BETA-GALACTOSIDASE



- Molecule 1: BETA-GALACTOSIDASE



## 4 Data and refinement statistics

Property	Value	Source
Space group	P 21 21 21	Depositor
Cell constants a, b, c, $\alpha$ , $\beta$ , $\gamma$	38.25Å 69.28Å 121.17Å 90.00° 90.00° 90.00°	Depositor
Resolution (Å)	19.86 – 1.83 19.86 – 1.83	Depositor EDS
% Data completeness (in resolution range)	99.6 (19.86-1.83) 99.7 (19.86-1.83)	Depositor EDS
$R_{merge}$	0.10	Depositor
$R_{sym}$	(Not available)	Depositor
$\langle I/\sigma(I) \rangle$ <sup>1</sup>	5.20 (at 1.82Å)	Xtrriage
Refinement program	REFMAC 5.7.0032	Depositor
R, $R_{free}$	0.169 , 0.217 0.177 , 0.222	Depositor DCC
$R_{free}$ test set	1481 reflections (5.07%)	wwPDB-VP
Wilson B-factor (Å <sup>2</sup> )	14.5	Xtrriage
Anisotropy	0.065	Xtrriage
Bulk solvent $k_{sol}$ (e/Å <sup>3</sup> ), $B_{sol}$ (Å <sup>2</sup> )	0.39 , 49.3	EDS
L-test for twinning <sup>2</sup>	$\langle  L  \rangle = 0.48$ , $\langle L^2 \rangle = 0.32$	Xtrriage
Estimated twinning fraction	No twinning to report.	Xtrriage
$F_o, F_c$ correlation	0.95	EDS
Total number of atoms	3108	wwPDB-VP
Average B, all atoms (Å <sup>2</sup> )	17.0	wwPDB-VP

Xtrriage's analysis on translational NCS is as follows: *The largest off-origin peak in the Patterson function is 10.91% of the height of the origin peak. No significant pseudotranslation is detected.*

<sup>1</sup>Intensities estimated from amplitudes.

<sup>2</sup>Theoretical values of  $\langle |L| \rangle$ ,  $\langle L^2 \rangle$  for acentric reflections are 0.5, 0.333 respectively for untwinned datasets, and 0.375, 0.2 for perfectly twinned datasets.

## 5 Model quality [i](#)

### 5.1 Standard geometry [i](#)

Bond lengths and bond angles in the following residue types are not validated in this section: CA, EDO

The Z score for a bond length (or angle) is the number of standard deviations the observed value is removed from the expected value. A bond length (or angle) with  $|Z| > 5$  is considered an outlier worth inspection. RMSZ is the root-mean-square of all Z scores of the bond lengths (or angles).

Mol	Chain	Bond lengths		Bond angles	
		RMSZ	# $ Z  > 5$	RMSZ	# $ Z  > 5$
1	A	1.13	2/1414 (0.1%)	0.96	4/1920 (0.2%)
1	B	0.64	0/1500	0.77	1/2038 (0.0%)
All	All	0.91	2/2914 (0.1%)	0.87	5/3958 (0.1%)

All (2) bond length outliers are listed below:

Mol	Chain	Res	Type	Atoms	Z	Observed(Å)	Ideal(Å)
1	A	1841	LYS	CD-CE	-27.44	0.82	1.51
1	A	1969	LYS	CD-CE	-21.72	0.96	1.51

All (5) bond angle outliers are listed below:

Mol	Chain	Res	Type	Atoms	Z	Observed(°)	Ideal(°)
1	A	1841	LYS	CG-CD-CE	15.78	159.23	111.90
1	A	1969	LYS	CG-CD-CE	12.86	150.46	111.90
1	A	1969	LYS	CD-CE-NZ	-12.67	82.56	111.70
1	A	1841	LYS	CD-CE-NZ	9.10	132.63	111.70
1	B	1860	ARG	NE-CZ-NH2	-5.71	117.44	120.30

There are no chirality outliers.

There are no planarity outliers.

### 5.2 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 the 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 within the asymmetric unit, whereas Symm-Clashes lists symmetry related clashes.

Mol	Chain	Non-H	H(model)	H(added)	Clashes	Symm-Clashes
1	A	1378	0	1286	10	0
1	B	1460	0	1353	8	0
2	A	1	0	0	0	0
2	B	1	0	0	0	0
3	A	4	0	6	0	0
3	B	16	0	24	2	0
4	A	138	0	0	2	0
4	B	110	0	0	1	0
All	All	3108	0	2669	18	0

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.

The worst 5 of 18 close contacts within the same asymmetric unit are listed below, sorted by their clash magnitude.

Atom-1	Atom-2	Interatomic distance (Å)	Clash overlap (Å)
1:A:1997:SER:HA	1:A:1998:GLU:C	1.99	0.83
1:B:1862:THR:CG2	1:B:1864:TRP:H	1.96	0.78
1:A:1862:THR:HG23	1:A:1864:TRP:H	1.51	0.74
1:A:1862:THR:CG2	1:A:1864:TRP:H	2.06	0.67
1:A:1997:SER:O	4:A:2137:HOH:O	2.15	0.64

There are no symmetry-related clashes.

## 5.3 Torsion angles [i](#)

### 5.3.1 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 X-ray entries followed by that with respect to entries of similar resolution.

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	172/186 (92%)	165 (96%)	7 (4%)	0	100	100
1	B	184/186 (99%)	174 (95%)	10 (5%)	0	100	100
All	All	356/372 (96%)	339 (95%)	17 (5%)	0	100	100

There are no Ramachandran outliers to report.



### 5.3.2 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 X-ray entries followed by that with respect to entries of similar resolution.

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	145/157 (92%)	144 (99%)	1 (1%)	84	78
1	B	152/157 (97%)	151 (99%)	1 (1%)	84	78
All	All	297/314 (95%)	295 (99%)	2 (1%)	84	78

All (2) residues with a non-rotameric sidechain are listed below:

Mol	Chain	Res	Type
1	A	1897	HIS
1	B	1950	ARG

Some sidechains can be flipped to improve hydrogen bonding and reduce clashes. All (2) such sidechains are listed below:

Mol	Chain	Res	Type
1	A	1897	HIS
1	B	1962	ASN

### 5.3.3 RNA [i](#)

There are no RNA molecules in this entry.

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

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

### 5.5 Carbohydrates [i](#)

There are no carbohydrates in this entry.

## 5.6 Ligand geometry

Of 7 ligands modelled in this entry, 2 are monoatomic - leaving 5 for Mogul analysis.

In the following table, the Counts columns list the number of bonds (or angles) for which Mogul statistics could be retrieved, the number of bonds (or angles) that are observed in the model and the number of bonds (or angles) that are defined in the Chemical Component Dictionary. The Link column lists molecule types, if any, to which the group is linked. The Z score for a bond length (or angle) is the number of standard deviations the observed value is removed from the expected value. A bond length (or angle) with  $|Z| > 2$  is considered an outlier worth inspection. RMSZ is the root-mean-square of all Z scores of the bond lengths (or angles).

Mol	Type	Chain	Res	Link	Bond lengths			Bond angles		
					Counts	RMSZ	# Z  > 2	Counts	RMSZ	# Z  > 2
3	EDO	B	3000	-	3,3,3	0.41	0	2,2,2	0.52	0
3	EDO	A	3000	-	3,3,3	0.57	0	2,2,2	0.39	0
3	EDO	B	3001	-	3,3,3	0.33	0	2,2,2	0.54	0
3	EDO	B	2999	-	3,3,3	0.51	0	2,2,2	0.26	0
3	EDO	B	3002	-	3,3,3	0.50	0	2,2,2	0.25	0

In the following table, the Chirals column lists the number of chiral outliers, the number of chiral centers analysed, the number of these observed in the model and the number defined in the Chemical Component Dictionary. Similar counts are reported in the Torsion and Rings columns. '-' means no outliers of that kind were identified.

Mol	Type	Chain	Res	Link	Chirals	Torsions	Rings
3	EDO	B	3000	-	-	1/1/1/1	-
3	EDO	A	3000	-	-	0/1/1/1	-
3	EDO	B	3001	-	-	1/1/1/1	-
3	EDO	B	2999	-	-	0/1/1/1	-
3	EDO	B	3002	-	-	1/1/1/1	-

There are no bond length outliers.

There are no bond angle outliers.

There are no chirality outliers.

All (3) torsion outliers are listed below:

Mol	Chain	Res	Type	Atoms
3	B	3002	EDO	O1-C1-C2-O2
3	B	3000	EDO	O1-C1-C2-O2
3	B	3001	EDO	O1-C1-C2-O2

There are no ring outliers.

1 monomer is involved in 2 short contacts:

Mol	Chain	Res	Type	Clashes	Symm-Clashes
3	B	3001	EDO	2	0

## 5.7 Other polymers [i](#)

There are no such residues in this entry.

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

There are no chain breaks in this entry.

## 6 Fit of model and data [i](#)

### 6.1 Protein, DNA and RNA chains [i](#)

In the following table, the column labelled ‘#RSRZ> 2’ contains the number (and percentage) of RSRZ outliers, followed by percent RSRZ outliers for the chain as percentile scores relative to all X-ray entries and entries of similar resolution. The OWAB column contains the minimum, median, 95<sup>th</sup> percentile and maximum values of the occupancy-weighted average B-factor per residue. The column labelled ‘Q< 0.9’ lists the number of (and percentage) of residues with an average occupancy less than 0.9.

Mol	Chain	Analysed	<RSRZ>	#RSRZ>2	OWAB(Å <sup>2</sup> )	Q<0.9
1	A	173/186 (93%)	-0.21	2 (1%) 79 79	8, 14, 24, 48	2 (1%)
1	B	185/186 (99%)	-0.09	8 (4%) 35 32	8, 15, 36, 62	1 (0%)
All	All	358/372 (96%)	-0.15	10 (2%) 53 51	8, 15, 32, 62	3 (0%)

The worst 5 of 10 RSRZ outliers are listed below:

Mol	Chain	Res	Type	RSRZ
1	B	1827	ASN	5.2
1	B	1824	TYR	4.7
1	A	1997	SER	3.1
1	B	1997	SER	2.9
1	B	1950	ARG	2.8

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

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

### 6.3 Carbohydrates [i](#)

There are no carbohydrates in this entry.

### 6.4 Ligands [i](#)

In the following table, the Atoms column lists the number of modelled atoms in the group and the number defined in the chemical component dictionary. The B-factors column lists the minimum, median, 95<sup>th</sup> percentile and maximum values of B factors of atoms in the group. The column labelled ‘Q< 0.9’ lists the number of atoms with occupancy less than 0.9.

Mol	Type	Chain	Res	Atoms	RSCC	RSR	B-factors( $\text{\AA}^2$ )	Q<0.9
3	EDO	B	3002	4/4	0.76	0.24	43,43,44,45	0
3	EDO	B	3000	4/4	0.85	0.17	44,45,47,48	0
3	EDO	B	2999	4/4	0.96	0.09	11,13,14,15	0
3	EDO	B	3001	4/4	0.96	0.08	22,22,23,24	0
3	EDO	A	3000	4/4	0.97	0.11	12,12,12,12	0
2	CA	A	2999	1/1	1.00	0.04	12,12,12,12	0
2	CA	B	2998	1/1	1.00	0.03	11,11,11,11	0

## 6.5 Other polymers [\(i\)](#)

There are no such residues in this entry.