

PROTEIN DATA BANK

ATOMIC COORDINATE AND BIBLIOGRAPHIC ENTRY FORMAT DESCRIPTION

February 1992

TABLE OF CONTENTS

	<u>Page</u>
Summary of Record Types and Their Sequence	3
Record Formats	5
Appendix A - Coordinate Systems and Transformations	23
Appendix B - Atom Names	26
Appendix C - Standard Residue Names and Abbreviations	36
Appendix D - Protein Data Bank Conventions	39
Appendix E - Formats for Literature Citations	42
Appendix F - Formulas and Molecular Weights for Standard Residues	44

SUMMARY OF RECORD TYPES AND THEIR SEQUENCE

For each atomic coordinate and bibliographic entry, the file consists of records each of 80 characters. The record sequence is as follows:

HEADER	:	Date entered into Data Bank; identification code
OBSLTE	:	Identifies entries which have been replaced
COMPND	:	Name of molecule and identifying information
SOURCE	:	Species, organ, tissue, and mutant from which the molecule has been obtained, where applicable
EXPDTA	:	Experimental technique of structure determination
AUTHOR	:	Names of contributors
REVDAT	:	Revision date; identifies current modification level
SPRSDE	:	Identifies entries which have replaced others
JRNL	:	Literature citation that defines coordinate set
REMARK	:	General remarks
SEQRES	:	Residue sequence
FTNOTE	:	Footnotes relating to specific atoms or residues
HET	:	Identification of non-standard groups or residues (heterogens)
FORMUL	:	Chemical formulas of non-standard groups
HELIX	:	Identification of helical substructures
SHEET	:	Identification of sheet substructures
TURN	:	Identification of hairpin turns
SSBOND	:	Identification of disulfide bonds
SITE	:	Identification of groups comprising the various sites
CRYST1	:	Unit cell parameters, space group designation
ORIGX	:	Transformation from orthogonal Å coordinates to submitted coordinates
SCALE	:	Transformation from orthogonal Å coordinates to fractional crystallographic coordinates
MATRIX	:	Transformations expressing non-crystallographic symmetry

TVECT	:	Translation vector for infinite covalently connected structures
MODEL	:	Specification of model number for multiple structure models in a single data entry
ATOM	:	Atomic coordinate records for "standard" groups
HETATM	:	Atomic coordinate records for "non-standard" groups
SIGATM	:	Standard deviations of atomic parameters
ANISOU	:	Anisotropic temperature factors
SIGUIJ	:	Standard deviations of anisotropic temperature factors
TER	:	Chain terminator
ENDMDL	:	End-of-model flag for multiple structure models in a single data entry
CONNECT	:	Connectivity records
MASTER	:	Master control record with checksums of total number of records in the file, for selected record types
END	:	End-of-entry record

In describing record formats it will be convenient to use the punched-card analogy and refer to column numbers. Although up to six characters are used in the record tag words above, only the first four are needed to define the record type uniquely.

Records are present in each entry in the order specified above with the following exceptions:

- (i) ATOM and HETATM records appear in the order appropriate to the structure.
- (ii) TER records may appear among ATOM or HETATM records as appropriate.
- (iii) SIGATM, ANISOU and SIGUIJ records, when present, directly follow the corresponding ATOM (or HETATM) record in the order SIGATM, ANISOU, SIGUIJ.
- (iv) A MODEL record precedes, and an ENDMDL record follows, the set of ATOM, HETATM, and TER records for each model among a series of multiple structure models in a single data entry. MODEL and ENDMDL records generally are employed only for NMR entries.

Note: In the future we anticipate making some changes in the format of our atomic coordinate entries. As far as possible, we will do this by adding new types of records; user programs should thus be prepared to handle or ignore record types that are not currently defined. Any revisions to existing record types will be announced well in advance via the Protein Data Bank Newsletter. Because of user requests, we are reserving a class of records for users' definition and use. All records beginning with the four letters USER are reserved for user definition and will be ignored by our programs.

RECORD FORMATS

1. **HEADER** Cols. 1 - 6 HEADER
- 11 - 50 Functional classification of macromolecule
- 51 - 59 Date of deposition into Data Bank⁽ⁱ⁾
- 63 - 66 Identification code⁽ⁱⁱ⁾

FORMAT (6A1,4X,40A1,9A1,3X,A4)

Note: (i) The date is given in the form dd-mmm-yy (e.g., December 1, 1983 is given as 01-DEC-83).

(ii) Each macromolecule is assigned an identification code. The code consists of 1 numeric and 3 alphanumeric characters in that order.

2. **OBSLTE** Cols. 1 - 6 OBSLTE
- 9 - 10 Continuation field (this field will be blank for the first OBSLTE record in each entry and numbered 2, 3, etc. for continuation records)
- 12 - 20 Date this entry was replaced
- 22 - 25 Identification code of this entry which is now obsolete
- 32 - 35 Identification code of a new entry which has replaced this old entry
- 37 - 40 .
- . .
- . .
- 67 - 70 Identification code of a new entry which has replaced this old entry

FORMAT (6A1,2X,2A1,1X,9A1,1X,4A1,5X,8(1X,4A1))

Note: This record will be inserted only in archived entries that are no longer distributed.

3. **COMPND** Cols. 1 - 6 COMPND
- 9 - 10 Continuation field (this field will be blank for the first COMPND record in each entry and numbered 2, 3, etc. for continuation records)
- 11 - 70 Name of macromolecule

FORMAT (6A1,4X,60A1)

Note: For enzymes the E.C. number is given in the form (E.C.n.n.n.n) with no internal blanks and without splitting over two lines. If an enzyme has not had an E.C. number assigned, the string (E.C. NUMBER NOT ASSIGNED) will be used. The Enzyme Commission numbers are obtained from *Enzyme Nomenclature 1984*, published for the International Union of Biochemistry by Academic Press, Inc.

4. **SOURCE** Cols. 1 - 6 SOURCE
- 9 - 10 Continuation field (this field will be blank for the first SOURCE record in each entry and numbered 2, 3, etc. for continuation records)
- 11 - 70 Species, organ, tissue, and mutant from which the macromolecule has been obtained. The systematic name of the species is given in parentheses.

FORMAT (6A1,4X,60A1)

5. **AUTHOR** Cols. 1 - 6 AUTHOR
- 9 - 10 Continuation field (this field will be blank for the first AUTHOR record in each entry and numbered 2, 3, etc. for continuation records)
- 11 - 70 Name(s) of contributor(s)

FORMAT (6A1,4X,60A1)

6. **EXPDTA** Cols. 1 - 6 EXPDTA
- 11 - 70 Experimental technique

FORMAT (6A1,4X,60A1)

Note: EXPDTA records may be used to specify the experimental technique employed to determine the structure in question, e.g., FIBER DIFFRACTION, NEUTRON DIFFRACTION, ELECTRON DIFFRACTION, NMR, THEORETICAL MODEL. If no EXPDTA record is present, the default X-RAY DIFFRACTION is assumed.

7. **REVDAT** Cols. 1 - 6 REVDAT
- 8 - 10 Modification number⁽ⁱ⁾
- 11 - 12 Continuation field⁽ⁱⁱ⁾
- 14 - 22 Date⁽ⁱⁱⁱ⁾
- 24 - 28 Identification name used for the correction

32 Modification type^(iv)
 40 - 70 Record types that were corrected.

FORMAT (6A1,1X,I3,2A1,1X,9A1,1X,5A1,3X,I1,7X,31A1)

- Notes: (i) Each revision will be given a modification number assigned in increasing numerical order but inserted in the entry in decreasing numerical order. New entries will be assigned the modification number 1.
- (ii) For each modification, more than one REDVAT record may be supplied. This field will be blank on the first record of each modification, and numbered 2, 3, etc. for continuation records.
- (iii) For new entries this date will be the date when the entry was released for distribution rather than the date of deposition which appears in the HEADER record.
- (iv) The following integer values will be used to identify the modification type: (In case of revisions with more than one possible type, the highest value applicable will be assigned).
- 3 - Used for modifications affecting the coordinates or their transforms. To be used for entries with revisions to any of the following records.
- a) CRYST1
 - b) ORIGX
 - c) SCALE
 - d) MTRIX
 - e) TVECT
 - f) ATOM
 - g) HETATM
 - h) SIGATM
- 2 - Used for modifications to the CONECT records.
- 1 - Used for all other types of modifications, mainly typographical in nature.
- 0 - Initial entry. New entries will contain a REVDAT record flagged as modification type 0.

8.	<u>SPRSDE</u>	Cols.	1 - 6	SPRSDE
			9 - 10	Continuation field (this field will be blank for the first SPRSDE record in each entry and numbered 2, 3, etc. for continuation records)
			12 - 20	Date that this entry superseded an older one
			22 - 25	Identification code of this entry which is replacing an older one

32 - 35 Identification codes of the entries which are being replaced
 •
 •
 •
 67 - 70 Identification codes of the entries which are being replaced

FORMAT (6A1,2X,2A1,1X,9A1,1X,4A1,5X,8(1X,4A1))

9. **JRNL** Cols. 1 - 4 JRNL
 11 - 70 Literature citation that defines the coordinate set

FORMAT (6A1,4X,60A1)

Note: If the coordinate set held is identified in the literature, the paper containing the definition is cited here. If an article defines more than one coordinate set, the particular designation assigned is given in the REMARKs section. The format of literature citations for both the JRNL and REMARK 1 records is given in Appendix E.

10. **REMARK** Cols. 1 - 6 REMARK
 8 - 10 Remark number
 12 - 70 Text of remark

FORMAT (6A1,1X,I3,1X,59A1)

Note: The first REMARK has serial number 1, the second 2, etc. REMARK 1 lists the important papers relating to a structure which originate from the depositor's laboratory. These papers are usually listed in inverse chronological order, except if a particular article (or series of articles) is considered to be a definitive description, in which case it may appear first. If any citation is given in the JRNL records, it is not repeated here. The format of literature citations for both JRNL and REMARK 1 records is given in Appendix E.

REMARKs 2 and 3 are reserved for statements relating to the resolution and refinement of the structure analysis. Other general commentary is given in higher numbered REMARKs.

11. **SEQRES** Cols. 1 - 6 SEQRES
 9 - 10 Serial number of SEQRES record for current chain
 12 Chain identifier
 14 - 17 Number of residues in this chain

20 - 22 Residue name
 24 - 26 Residue name
 .
 .
 68 - 70 Residue name

FORMAT (6A1,I4,1X,A1,1X,I4,1X,13(1X,A3))

12. **FTNOTE** Cols. 1 - 6 FTNOTE
 8 - 10 Footnote number
 12 - 70 Footnote text

FORMAT (6A1,1X,I3,1X,59A1)

Note: FTNOTE records are used to describe details which are specific to certain atoms or residues. These footnotes are keyed to particular atoms by the footnote number here and in cols. 68-70 of the ATOM/HETATM record. Any individual footnote may run over several FTNOTE records (each with the same footnote number). A maximum of 999 footnotes are allowed.

13. **HET** Cols. 1 - 3 HET
 8 - 10 Non-standard group (heterogen) identifier
 13 Chain identifier
 14 - 17 Sequence number
 18 Insertion code
 21 - 25 Number of atoms in non-standard group
 31 - 70 Text

FORMAT (6A1,1X,A3,2X,A1,I4,A1,2X,I5,5X,40A1)

Note: HET records are used to describe non-standard residues, prosthetic groups, inhibitors, solvent molecules (except water), etc. The Protein Data Bank attempts to use uniform atom nomenclature for HET groups, as illustrated for some commonly occurring groups in Appendix B. All non-standard groups (i.e., those not assigned a standard code in Appendix C) are defined in HET records. If there is insufficient space in the text portion of the HET record to properly define a non-standard group then the definition will be given in a REMARK and referenced here. An entry of -999 in the sequence number field is used to indicate that the HET group occurs too frequently, in the present entry, to use a separate HET record for each occurrence.

14.	<u>FORMUL</u>	Cols.	1 - 6	FORMUL
			9 - 12	Component number ⁽ⁱ⁾
			13 - 15	Non-standard group (HET) identifier
			17 - 18	Continuation number ⁽ⁱⁱ⁾ (blank on first record)
			19	*if this component is to be excluded from the molecular weight calculation ⁽ⁱⁱⁱ⁾
			20 - 70	Formula of non-standard group ^(iv)

FORMAT (6A1,2X,I2,2X,A3,1X,I2,A1,51A1)

- Notes: (i) Component numbers are assigned serially. Each component represented by a set of SEQRES records is counted first and then each HET group is assigned a component number in sequence. If a HET group is contained within a chain represented by a set of SEQRES records, the component number assigned is that of the chain involved.
- (ii) If a HET group is composed of more than one distinct part, then the formulas for these parts will occur on separate FORMUL cards, each with the same component number and HET identifier. All except the last of these records will be terminated with a period.
- (iii) Solvent molecules and certain other components are normally excluded. The molecular weight may be used as a key for automatic searching of the file.
- (iv) Each component defined in a HET record for which a standard chemical formula can be written is defined accordingly here. Atoms which are known to be present but not located in the crystallographic analysis (e.g., hydrogen atoms) are represented in the formula. Formulas are written as C, H, N, O with other elements following in alphabetical order. The repeat count of each atom type present immediately follows the chemical symbol. A repeat count of the entire group is indicated by enclosing the formula in parentheses and prefacing the string with the count. The oxidation state of metals is given when it is known. For the two heme groups of ferrihemoglobin the FORMUL record would have HEM 2(C34 H32 N4 O4 FE1 +++).

15.	<u>HELIX</u>	Cols.	1 - 6	HELIX
			8 - 10	Serial number (helix number)
			12 - 14	Helix identifier (right justified) ⁽ⁱ⁾

16 - 18	Residue name	}	Initial residue of helix ⁽ⁱⁱ⁾
20	Chain identifier		
22 - 25	Residue seq. no.		
26	Insertion code		
28 - 30	Residue name	}	Terminal residue of helix ⁽ⁱⁱ⁾
32	Chain identifier		
34 - 37	Residue seq. no.		
38	Insertion code		
39 - 40	Class of helix ⁽ⁱⁱⁱ⁾		
41 - 70	Comment		

FORMAT (6A1,1X,I3,1X,A3,2(1X,A3,1X,A1,1X,I4,A1),I2,30A1)

Notes: (i) Additional records with different serial numbers and identifiers occur if more than one helix is present.

(ii) The initial residue has a lower sequence number than the terminal residue.

(iii) Helices are classified as:

1 Right-handed α (default)	2 Right-handed ω
3 Right-handed π	4 Right-handed γ
5 Right-handed 3_{10}	6 Left-handed α
7 Left-handed ω	8 Left-handed γ
9 2_7 ribbon/helix	10 Polyproline

16. SHEET	Cols.	1 - 5	SHEET
		8 - 10	Strand number ⁽ⁱ⁾ (v)
		12 - 14	Sheet identifier ⁽ⁱ⁾ (right justified)
		15 - 16	Number of strands

18 - 20	Residue name	}	Initial residue ⁽ⁱⁱ⁾
22	Chain identifier		
23 - 26	Residue seq. no.		
27	Insertion code		
29 - 31	Residue name	}	Terminal residue ⁽ⁱⁱ⁾
33	Chain identifier		
34 - 37	Residue seq. no.		
38	Insertion code		
39 - 40	Sense of this strand with respect to previous strand ⁽ⁱⁱⁱ⁾		
42 - 45	Atom name	}	In current strand. Registration ^(iv)
46 - 48	Residue name		
50	Chain identifier		
51-54	Residue seq. no.		
55	Insertion code		
57 - 60	Atom name	}	In previous strand. Registration ^(iv)
61 - 63	Residue name		
65	Chain identifier		
66 - 69	Residue seq. no.		
70	Insertion code		

FORMAT (6A1,1X,I3,1X,A3,I2,2(1X,A3,1X,A1,I4,A1),I2,2(1X,A4,A3,1X,A1,I4,A1))

- Notes: (i) Different strands are described in subsequent records which bear the same sheet identifier but different strand numbers.
- (ii) The initial residue of a strand has a lower sequence number than the terminal residue.
- (iii) Parallelism or anti-parallelism of strand n with respect to strand $n - 1$ is denoted by 1 or -1 , respectively. Strand 1 has sense indicator 0.
- (iv) Registration of the strand n with respect to strand $n - 1$ may be specified by a particular hydrogen bond between the indicated atoms. One donor and one acceptor should be specified. These fields will be blank for strand 1.

- (v) Strand numbers are reset to 1 for the first strand of each new sheet. A closed sheet (β barrel) is indicated by having the first and last strands identical.

17.	<u>TURN</u>	Cols.	1 - 5	TURN	
			8 - 10	Sequence number (turn number)	
			12 - 14	Turn identifier (3 characters)	
			16 - 18	Residue name	} Residue i
			20	Chain identifier	
			21 - 24	Residue seq. no.	
			25	Insertion code	
			27 - 29	Residue name	} Residue i+3 (or i+2 for γ bend)
			31	Chain identifier	
			32 - 35	Residue seq. no.	
			36	Insertion code	
			41 - 70	Comment	

FORMAT (6A1,1X,I3,1X,A3,1X,A3,1X,A1,I4,A1,1X,A3,1X,A1,I4,A1,4X,30A1)

Note: These records identify the hairpin turns (β and γ bends) in the structure which do not occur in helices.

18.	<u>SSBOND</u>	Cols.	1 - 6	SSBOND
			8 - 10	Sequence number
			12 - 14	Residue name (CYS)
			16	Chain identifier
			18 - 21	Residue seq. no.
			22	Insertion code
			26 - 28	Residue name (CYS)
			30	Chain identifier

32 - 35 Residue seq. no.
 36 Insertion code
 41 - 70 Comment

FORMAT (6A1,1X,I3,1X,A3,1X,A1,1X,I4,A1,3X,A3,1X,A1,1X,I4,A1,4X,30A1)

19.	<u>SITE</u>	Cols.	1 - 4	SITE	
			8 - 10	Sequence number (i)	
			12 - 14	Site identifier ⁽ⁱⁱ⁾ (right justified)	
			16 - 17	Number of residues comprising site ⁽ⁱⁱⁱ⁾	
			19 - 21	Residue name	} First residue comprising site
			23	Chain identifier	
			24 - 27	Residue seq. no.	
			28	Insertion code	
			30 - 32	Residue name	} Second residue comprising site
			34	Chain identifier	
			35 - 38	Residue seq. no.	
			39	Insertion code	
			41 - 43	Residue name	} Third residue comprising site
			45	Chain identifier	
			46 - 49	Residue seq. no.	
			50	Insertion code	
			52 - 54	Residue name	} Fourth residue comprising site
			56	Chain identifier	
			57 - 60	Residue seq. no.	
			61	Insertion code	

FORMAT (6A1,1X,I3,1X,A3,1X,I2,4(1X,A3,1X,A1,I4,A1))

- Notes: (i) Sequence numbers are reset to 1 for each new site.
(ii) Site identifiers should be fully explained in the REMARKs.
(iii) If a site is comprised of more than four residues, these may be specified on additional records bearing the same site identifier.

20.	<u>CRYST1</u>	Cols.	1 - 6	CRYST1
			7 - 15	a(Å)
			16 - 24	b(Å)
			25 - 33	c(Å)
			34 - 40	α (deg.)
			41 - 47	β (deg.)
			48 - 54	γ (deg.)
			56 - 66	Space group symbol (left justified)
			67 - 70	Z

FORMAT (6A1,3F9.3,3F7.2,1X,11A1,I4)

Note: If the data entry describes a structure determined by a technique other than crystallography, CRYST1 will contain $a=b=c=1.0$, $\alpha=\beta=\gamma=90^\circ$; the space group symbol will contain P 1 and Z=1.

21.	<u>ORIGX</u>	Cols.	1 - 6	11-20	21-30	31-40	46-55
			ORIGX1	O ₁₁	O ₁₂	O ₁₃	T ₁
			ORIGX2	O ₂₁	O ₂₂	O ₂₃	T ₂
			ORIGX3	O ₃₁	O ₃₂	O ₃₃	T ₃

FORMAT (6A1,4X,3F10.6,5X,F10.5)

Note: Let the original submitted coordinates be X_{sub} , Y_{sub} , Z_{sub} and the orthogonal Å coordinates contained in the data entry be X,Y,Z. Then

$$\begin{aligned} X_{\text{sub}} &= O_{11}X + O_{12}Y + O_{13}Z + T_1 \\ Y_{\text{sub}} &= O_{21}X + O_{22}Y + O_{23}Z + T_2 \\ Z_{\text{sub}} &= O_{31}X + O_{32}Y + O_{33}Z + T_3 \end{aligned}$$

Even if this is an identity transformation (unit matrix, null vector) it is supplied. See below under SCALE for a definition of the default orthogonal Å system.

Appendix A details the derivation of the ORIGX coordinate transformation.

22.	<u>SCALE</u>	Cols. 1 - 6	11-20	21-30	31-40	46-55
	SCALE1		S ₁₁	S ₁₂	S ₁₃	U ₁
	SCALE2		S ₂₁	S ₂₂	S ₂₃	U ₂
	SCALE3		S ₃₁	S ₃₂	S ₃₃	U ₃

FORMAT (6A1,4X,3F10.6,5X,F10.5)

Note: Let the orthogonal Å coordinates be X, Y, Z. Let the fractional cell coordinates be x_{frac} , y_{frac} , z_{frac} . Then

$$x_{\text{frac}} = S_{11}X + S_{12}Y + S_{13}Z + U_1$$

$$y_{\text{frac}} = S_{21}X + S_{22}Y + S_{23}Z + U_2$$

$$z_{\text{frac}} = S_{31}X + S_{32}Y + S_{33}Z + U_3$$

The SCALE transformation provides a means of generating fractional coordinates from the orthogonal Å coordinates contained in the data entry.

The standard orthogonal Å coordinate system is related to the axial system of the unit cell supplied (CRYST1 record) by the definition below. (Non-standard coordinate systems are generally explained in the REMARKs.)

If \vec{a} , \vec{b} , \vec{c} describe the crystallographic cell edges and \vec{A} , \vec{B} , \vec{C} are unit vectors in the default orthogonal Å system, then

(i) \vec{A} , \vec{B} , \vec{C} and \vec{a} , \vec{b} , \vec{c} have the same origin.

(ii) \vec{A} is parallel to \vec{a} .

(iii) \vec{B} is parallel to $\vec{C} \times \vec{A}$.

(iv) \vec{C} is parallel to $\vec{a} \times \vec{b}$ (i.e., \vec{c}^*).

Appendix A details the derivation of the SCALE coordinate transformation.

23.	<u>MTRIX</u> ^(i,ii)	Cols. 1-6	8-10	11-20	21-30	31-40	46-55	59-60
	MTRIX1	Ser.no. ⁽ⁱ⁾		M ₁₁	M ₁₂	M ₁₃	V ₁	IGIVEN ⁽ⁱⁱⁱ⁾
	MTRIX2	Ser.no. ⁽ⁱ⁾		M ₂₁	M ₂₂	M ₂₃	V ₂	IGIVEN ⁽ⁱⁱⁱ⁾
	MTRIX3	Ser.no. ⁽ⁱ⁾		M ₃₁	M ₃₂	M ₃₃	V ₃	IGIVEN ⁽ⁱⁱⁱ⁾

FORMAT (6A1,1X,I3,3F10.6,5X,F10.5,3X,I2)

- Notes: (i) One trio of MTRIX records with a constant serial number is given for each non-crystallographic symmetry operation that is defined.
- (ii) The MTRIX transformations operate on the stored coordinates to yield equivalent representations of the molecule in the same space.
- (iii) If coordinates for the representations which are approximately related by the transformation in question are contained in the file, the quantity IGIVEN is set to 1. Otherwise this field will be blank.

24.	<u>TVECT</u>	Cols.	1 - 5	TVECT
			8 - 10	Serial number
			11 - 20	} Components of translation vector
			21 - 30	
			31 - 40	
			41 - 70	Comment

FORMAT (6A1,1X,I3,3F10.5,30A1)

Note: For structures not comprised of discrete molecules (e.g., infinite polysaccharide chains) the Protein Data Bank entry will contain a fragment which can be built into the full structure by the simple translation vectors of TVECT records.

25.	<u>MODEL</u>	Cols.	1 - 5	MODEL
			11 - 14	Model serial number, for multiple structure models in a single data entry. Generally employed only for NMR structures.

FORMAT (5A1,5X,I4)

Note: ATOM, HETATM, SIGATM, ANISOU, SIGUIJ, and TER records for each structure model, as appropriate, will occur between MODEL and ENDMDL records.

26.	<u>ATOM</u> <u>HETATM</u>	Atomic coordinate records for "standard" groups		
		Atomic coordinate records for "non-standard" groups		
		Cols.	1 - 4	ATOM
		or	1 - 6	HETATM
			7 - 11	Atom serial number ⁽ⁱ⁾
	13 - 16	Atom name ⁽ⁱⁱ⁾		

17	Alternate location indicator ⁽ⁱⁱⁱ⁾	
18 - 20	Residue name ^(iv,v)	
22	Chain identifier, e.g., A for hemoglobin α chain	
23 - 26	Residue seq. no.	
27	Code for insertions of residues, e.g., 66A, 66B, etc.	
31 - 38	X	} Orthogonal Å coordinates
39 - 46	Y	
47 - 54	Z	
55 - 60	Occupancy	
61 - 66	Temperature factor ^(vi)	
68 - 70	Footnote number	

FORMAT (6A1,I5,1X,A4,A1,A3,1X,A1,I4,A1,3X,3F8.3,2F6.2,1X,I3)

- Notes: (i) Residues occur in order of their sequence numbers which always increase starting from the N-terminal residue for proteins and the 5'-terminal residue for nucleic acids. Within each residue the atoms are ordered as indicated in Appendix B. If the residue sequence is known, certain atom serial numbers may be omitted to allow for future insertion of any missing atoms. If the sequence is not reliably known, these serial numbers are simply ordinals.
- (ii) See Appendix B
- (iii) Alternate locations for atoms may be denoted here by A, B, C, etc.
- (iv) Standard residue names are given in Appendix C; other components are defined in HET records.
- (v) HETATM records are used for water molecules and atoms contained in HET groups.
- (vi) Normally, the isotropic B value appears here. However, if anisotropic temperature factors have been provided, the temperature factor field of the corresponding ATOM or HETATM record will contain the equivalent U-isotropic [U(eq)] which is calculated by

$$U(\text{eq}) = 1/3[U(1,1) + U(2,2) + U(3,3)] \times 10^{-4}$$

27.	<u>SIGATM</u>	Cols.	1 - 6	SIGATM	
			7 - 27	Identical to corresponding ATOM/HETATM record	
			31 - 38	} Standard deviations of the stored coordinates (Å)	
			39 - 46		
			47 - 54		
			55 - 60		Standard deviation of occupancy
			61 - 66	Standard deviation of temperature factor	
			68 - 70	Footnote number	

FORMAT (6A1,I5,1X,A4,A1,A3,1X,A1,I4,A1,3X,3F8.3,2F6.2,1X,I3)

28.	<u>ANISOU</u>	Cols.	1 - 6	ANISOU	
			7 - 27	Identical to corresponding ATOM/HETATM record	
			29 - 35	U(1,1)	} Anisotropic temperature factors × 10 ⁴ (Å ²) ^(i,ii)
			36 - 42	U(2,2)	
			43 - 49	U(3,3)	
			50 - 56	U(1,2)	
			57 - 63	U(1,3)	
			64 - 70	U(2,3)	

FORMAT (6A1,I5,1X,A4,A1,A3,1X,A1,I4,A1,1X,6I7)

Notes: (i) If anisotropic temperature factors have been provided, the temperature factor field of the corresponding ATOM or HETATM record will contain the equivalent U – isotropic [U(eq)] which is calculated by

$$U(\text{eq}) = 1/3[U(1,1) + U(2,2) + U(3,3)] \times 10^{-4}$$

(ii) The anisotropic temperature factors will be stored in the same coordinate frame as the atomic coordinate records.

29.	<u>SIGUIJ</u>	Cols.	1 - 6	SIGUIJ	
			7 - 27	Identical to corresponding ATOM/HETATM record	
			29 - 35	Sigma U(1,1)	} Standard deviations of anisotropic temperature factors $\times 10^4(\text{\AA}^2)$
			36 - 42	Sigma U(2,2)	
			43 - 49	Sigma U(3,3)	
			50 - 56	Sigma U(1,2)	
			57 - 63	Sigma U(1,3)	
			64 - 70	Sigma U(2,3)	

FORMAT (6A1,I5,1X,A4,A1,A3,1X,A1,I4,A1,1X,6I7)

30.	<u>TER</u>	Cols.	1 - 3	TER
			7 - 11	Serial number
			18 - 20	Residue name
			22	Chain identifier
			23 - 26	Residue seq. no.
			27	Insertion code

FORMAT (6A1,I5,6X,A3,1X,A1,I4,A1)

Note: TER records occur among the ATOM records, and are placed after the terminal atom of each chain. For a protein the residue defined on these TER records is the carboxy-terminal residue of the chain in question. For a nucleic acid it is the 3'-terminal residue.

31.	<u>ENDMDL</u>	Cols.	1 - 6	ENDMDL
-----	----------------------	-------	-------	--------

FORMAT (6A1)

Note: ENDMDL records follow ATOM, HETATM, SIGATM, ANISOU, SIGUIJ and TER records for each structure model, for data entries with multiple structure models. Generally employed only for NMR structures.

32.	<u>CONNECT</u>	Connectivity records	
	Cols.	1 - 6	CONNECT
		7 - 11	Serial number
		12 - 16	} Covalent bond connectivity (serial numbers of bonded atoms)
		17 - 21	
		22 - 26	
		27 - 31	
		32 - 36	Hydrogen bond
		37 - 41	Hydrogen bond
		42 - 46	Salt bridge
		47 - 51	Hydrogen bond
		52 - 56	Hydrogen bond
		57 - 61	Salt bridge

} in which the atom specified in cols. 7-11 acts as donor

} the atom specified in cols. 7-11 has an excess of negative charge

} in which the atom specified in cols. 7-11 acts as acceptor

} the atom specified in cols. 7-11 has an excess of positive charge

FORMAT (6A1,11I5)

Note: Serial numbers are identical to those in cols. 7-11 of the appropriate ATOM/HETATM records, and connectivity entries correspond to these serial numbers. A second CONNECT record, with the same serial number in cols. 7-11, may be used if necessary. Either all or none of the covalent connectivity of an atom must be specified, and if hydrogen bonding is specified the covalent connectivity is included also.

The occurrence of a negative atom serial number on a CONNECT record denotes that a translationally equivalent copy (see TVECT records) of the target atom specified is linked to the origin atom of the record.

33.	<u>MASTER</u>	Cols.	1 - 6	MASTER
			11 - 15	Number of REMARK records
			16 - 20	Number of FTNOTE records
			21 - 25	Number of HET records
			26 - 30	Number of HELIX records
			31 - 35	Number of SHEET records

36 - 40	Number of TURN records
41 - 45	Number of SITE records
46 - 50	Number of coordinate transformation records (ORIGX + SCALE + MTRIX)
51 - 55	Number of atomic coordinate records (ATOM + HETATM)
56 - 60	Number of TER records
61 - 65	Number of CONECT records
66 - 70	Number of SEQRES records

FORMAT (6A1,4X,12I5)

Note: The MASTER record gives checksums of the number of records in the data entry, for selected record types.

34. **END** End-of-entry record
 Cols. 1 - 3 END

FORMAT (6A1)

APPENDIX A - COORDINATE SYSTEMS AND TRANSFORMATIONS

The coordinates stored in the Protein Data Bank give the atomic positions measured in Ångstroms along three orthogonal directions. Unless otherwise specified, the default axial system (detailed below) will be assumed.

If $\vec{a}, \vec{b}, \vec{c}$ describe the crystallographic cell edges and $\vec{A}, \vec{B}, \vec{C}$ are unit vectors in the default orthogonal Å system, then

- (i) $\vec{A}, \vec{B}, \vec{C}$ and $\vec{a}, \vec{b}, \vec{c}$ have the same origin.
- (ii) \vec{A} is parallel to \vec{a} .
- (iii) \vec{B} is parallel to $\vec{C} \times \vec{A}$.
- (iv) \vec{C} is parallel to $\vec{a} \times \vec{b}$ (i.e., \vec{c}^*).

The matrix which premultiplies the column vector of fractional crystallographic coordinates ($x_{\text{frac}}, y_{\text{frac}}, z_{\text{frac}}$) to yield coordinates in the $\vec{A}, \vec{B}, \vec{C}$ system, i.e., (X,Y,Z) is

$$\begin{bmatrix} a & b(\cos\gamma) & c(\cos\beta) \\ 0 & b(\sin\gamma) & c(\cos\alpha - \cos\beta \cos\gamma) / \sin\gamma \\ 0 & 0 & V/(ab \sin\gamma) \end{bmatrix}$$

where $V = abc(1 - \cos^2\alpha - \cos^2\beta - \cos^2\gamma + 2\cos\alpha \cos\beta \cos\gamma)^{1/2}$

If the submitted coordinates are either fractions of the unit cell edges or are with respect to the default orthogonal system, the ORIGX and SCALE transformations will be given default values.

In general the depositor will have supplied:

- (i) The original submitted coordinates,

$$\text{i.e.,} \quad \vec{X}_{\text{sub}}$$

- (ii) A transformation from \vec{X}_{sub} to the orthogonal Å coordinates stored in the Data Bank (\vec{X}),

$$\text{i.e.,} \quad O_{\text{sub}}\vec{X}_{\text{sub}} + \vec{T}_{\text{sub}} = \vec{X}$$

- (iii) A transformation from \vec{X}_{sub} to fractional crystallographic coordinates \vec{x}_{frac}

$$\text{i.e.,} \quad S_{\text{sub}}\vec{X}_{\text{sub}} + \vec{U}_{\text{sub}} = \vec{x}_{\text{frac}}$$

(iv) A set of transformations expressing any approximate or exact non-crystallographic symmetry elements in the structure

$$\text{i.e., } \mathbf{M}_{\text{sub}} \vec{\mathbf{X}}_{\text{sub}} + \vec{\mathbf{V}}_{\text{sub}} = \vec{\mathbf{X}}'_{\text{sub}}$$

Note: The notation $\vec{\mathbf{X}}_{\text{sub}}$ is used for the column vector $X_{\text{sub}}, Y_{\text{sub}}, Z_{\text{sub}}, \text{etc.}$

Since it is desirable for the stored ORIGX, SCALE and MTRIX transformations to operate on the stored rather than the submitted coordinates, some manipulation of the supplied quantities is performed in order to obtain the stored quantities.

The stored quantities are:

(i) The coordinates in orthogonal Ångstroms ($\vec{\mathbf{X}}$)

$$\vec{\mathbf{X}} = \mathbf{O}_{\text{sub}} \vec{\mathbf{X}}_{\text{sub}} + \vec{\mathbf{T}}_{\text{sub}}$$

(ii) The ORIGX transformation from stored to original coordinates ($\mathbf{O}, \vec{\mathbf{T}}$).

$$\text{From above } \vec{\mathbf{X}} = \mathbf{O}_{\text{sub}} \vec{\mathbf{X}}_{\text{sub}} + \vec{\mathbf{T}}_{\text{sub}}$$

$$\text{whence } \mathbf{O}_{\text{sub}} \vec{\mathbf{X}}_{\text{sub}} = \vec{\mathbf{X}} - \vec{\mathbf{T}}_{\text{sub}}$$

$$\therefore \vec{\mathbf{X}}_{\text{sub}} = \mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{X}} + (-\mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}})$$

$$\text{Thus } \mathbf{O} = \mathbf{O}_{\text{sub}}^{-1}$$

$$\text{and } \vec{\mathbf{T}} = -\mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}}$$

(iii) The SCALE transformation from stored to fractional coordinates ($\mathbf{S}, \vec{\mathbf{U}}$).

$$\text{From above } \vec{\mathbf{x}}_{\text{frac}} = \mathbf{S}_{\text{sub}} \vec{\mathbf{X}}_{\text{sub}} + \vec{\mathbf{U}}_{\text{sub}}$$

$$\text{but } \vec{\mathbf{X}}_{\text{sub}} = \mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{X}} + (-\mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}})$$

$$\therefore \vec{\mathbf{x}}_{\text{frac}} = \mathbf{S}_{\text{sub}} [\mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{X}} + (-\mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}})] + \vec{\mathbf{U}}_{\text{sub}}$$

$$\text{i.e., } \vec{\mathbf{x}}_{\text{frac}} = \mathbf{S}_{\text{sub}} \mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{X}} + (-\mathbf{S}_{\text{sub}} \mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}}) + \vec{\mathbf{U}}_{\text{sub}}$$

$$\therefore \mathbf{S} = \mathbf{S}_{\text{sub}} \mathbf{O}_{\text{sub}}^{-1}$$

$$\text{and } \vec{\mathbf{U}} = -(\mathbf{S}_{\text{sub}} \mathbf{O}_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}}) + \vec{\mathbf{U}}_{\text{sub}}$$

(iv) The MTRIX transformation(s) expressing non-crystallographic symmetry in the space of the stored coordinates ($\vec{\mathbf{M}}, \vec{\mathbf{V}}$).

$$\begin{aligned}
 \vec{\mathbf{X}}'_{\text{sub}} &= \mathbf{M}_{\text{sub}} \vec{\mathbf{X}}_{\text{sub}} + \vec{\mathbf{V}}_{\text{sub}} \\
 \vec{\mathbf{X}}' &= O_{\text{sub}} \vec{\mathbf{X}}'_{\text{sub}} + \vec{\mathbf{T}}_{\text{sub}} \\
 &= O_{\text{sub}} \{ \mathbf{M}_{\text{sub}} \vec{\mathbf{X}}_{\text{sub}} + \vec{\mathbf{V}}_{\text{sub}} \} + \vec{\mathbf{T}}_{\text{sub}} \\
 \text{but } \vec{\mathbf{X}}_{\text{sub}} &= (O_{\text{sub}}^{-1} \vec{\mathbf{X}} + (-O_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}})) \\
 \therefore \vec{\mathbf{X}}' &= O_{\text{sub}} \{ \mathbf{M}_{\text{sub}} [O_{\text{sub}}^{-1} \vec{\mathbf{X}} + (-O_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}})] + \vec{\mathbf{V}}_{\text{sub}} \} + \vec{\mathbf{T}}_{\text{sub}} \\
 \text{and } \mathbf{M} &= O_{\text{sub}} \mathbf{M}_{\text{sub}} O_{\text{sub}}^{-1} \\
 \vec{\mathbf{V}} &= -O_{\text{sub}} \mathbf{M}_{\text{sub}} O_{\text{sub}}^{-1} \vec{\mathbf{T}}_{\text{sub}} + O_{\text{sub}} \vec{\mathbf{V}}_{\text{sub}} + \vec{\mathbf{T}}_{\text{sub}}
 \end{aligned}$$

In summary the stored coordinates and transformations are:

$\vec{\mathbf{X}}$	(ATOM, HETATM records)
$\mathbf{O}, \vec{\mathbf{T}}$	(ORIGX records)
$\mathbf{s}, \vec{\mathbf{U}}$	(SCALE records)
$\mathbf{M}, \vec{\mathbf{V}}$	(MTRIX records)

APPENDIX B - ATOM NAMES

A. Amino Acids

These atom names follow the IUPAC-IUB rules¹ except:

- (i) Greek letter remoteness codes are transliterated as follows: α -A, β -B, γ -G, δ -D, ϵ -E, ζ -Z, η -H
- (ii) Atoms for which some ambiguity exists in the crystallographic results are designated A. This will usually apply only to the terminal atoms of asparagine and glutamine and to the ring atoms of histidine.

Within each residue the atoms occur in the order specified by the superscripts (following figure).

The extra oxygen atom of the carboxy terminal amino acid is designated OXT.

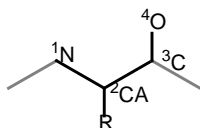
Four characters are reserved for these atom names. They are assigned as follows:

- 1-2 Chemical symbol - right justified
- 3 Remoteness indicator (alphabetic)
- 4 Branch designator (numeric)

- (iii) For protein coordinate sets containing hydrogen atoms, the IUPAC-IUB rules¹ have been followed. Recommendation rule number 4.4 has been modified as follows: When more than one hydrogen atom is bonded to a single non-hydrogen atom, the hydrogen atom number designation is given as the first character of the atom name rather than as the last character (e.g. H ^{β 1} is denoted as 1HB). Exceptions to these rules may occur in certain data sets at the depositors' request. Any such exceptions will be delineated clearly in FTNOTE and REMARK records.

¹IUPAC-IUB Commission on Biochemical Nomenclature. "Abbreviations and Symbols for the Description of the Conformation of Polypeptide Chains. Tentative Rules (1969)", *J. Biol. Chem.* 245, 6489 (1970).

The 1974 recommendations on the "Nomenclature of α -Amino Acids (*Biochemistry*, 14, 449 (1975)) provides a scheme based on normal rules for organic compounds, but this scheme will not be used here.

backbone

<u>Name</u>	<u>Side Chain</u>	<u>Name</u>	<u>Side Chain</u>
Alanine	— 5CB	Leucine	
Arginine		Lysine	— 5CB — 6CG — 7CD — 8CE — 9NZ
Asparagine		Methionine	— 5CB — 6CG — 7SD — 8CE
Aspartic Acid		Phenylalanine	
Cysteine/ Cystine	— 5CB — 6SG	Proline	
Glutamic Acid		Serine	— 5CB — 6OG
Glutamine		Threonine	
Glycine	— null	Tryptophan	
Histidine		Tyrosine	
Hydroxyproline		Valine	
Isoleucine			

**ATOM NAMES, REMOTENESS CODES, AND ORDER INDICATORS
FOR THE COMMON AMINO ACIDS.**

B. Nucleic Acids

Atom names employed for polynucleotides generally follow the precedents set for mononucleotides. The following points are worthy of note.

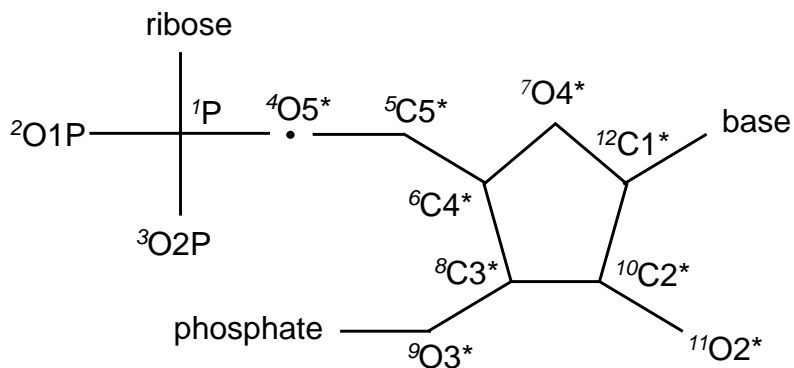
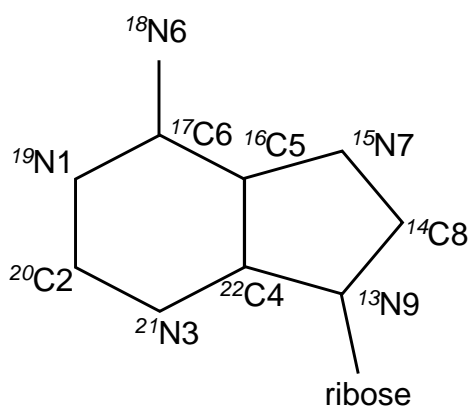
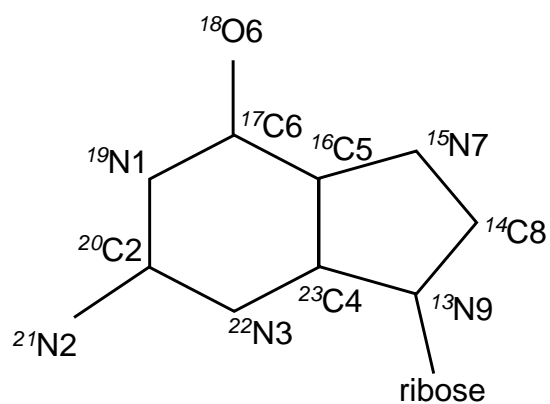
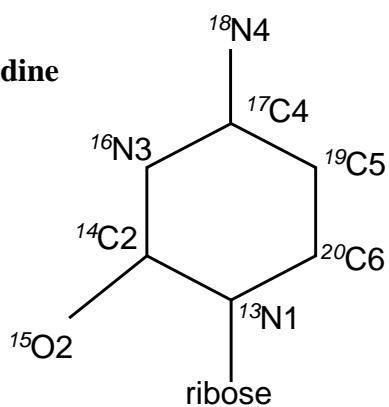
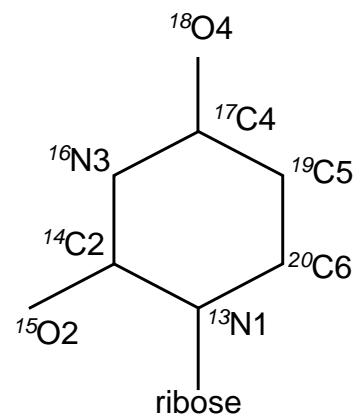
- (i) The prime character (') commonly used to denote atoms of the ribose originally was avoided because of non-uniformity of its external representation. An asterisk (*) therefore was used in its place, in entries released through January 1992.
- (ii) Of the four characters reserved for atom names the leftmost two are reserved for the chemical symbol (right justified) and the remaining two denote the atom's position.
- (iii) Atoms exocyclic to the ring systems have the same position identifier as the atom to which they are bonded except if this would result in identical atom names. In this case an alphabetic character is used to avoid ambiguity.
- (iv) The ring-oxygen atom of the ribose is denoted O4 rather than O1.
- (v) The extra oxygen atom at the free 5' phosphate terminus is designated OXT. This atom will be placed first in the coordinate set.

For nucleotides which are simple derivatives (e.g., methyl or acetyl) of the parent nucleotide the modifying atoms or groups occur immediately after the atom to which they are bonded. In the case of an acetyl modifier, the three atoms are ordered carbonyl carbon, carbonyl oxygen, methyl carbon.

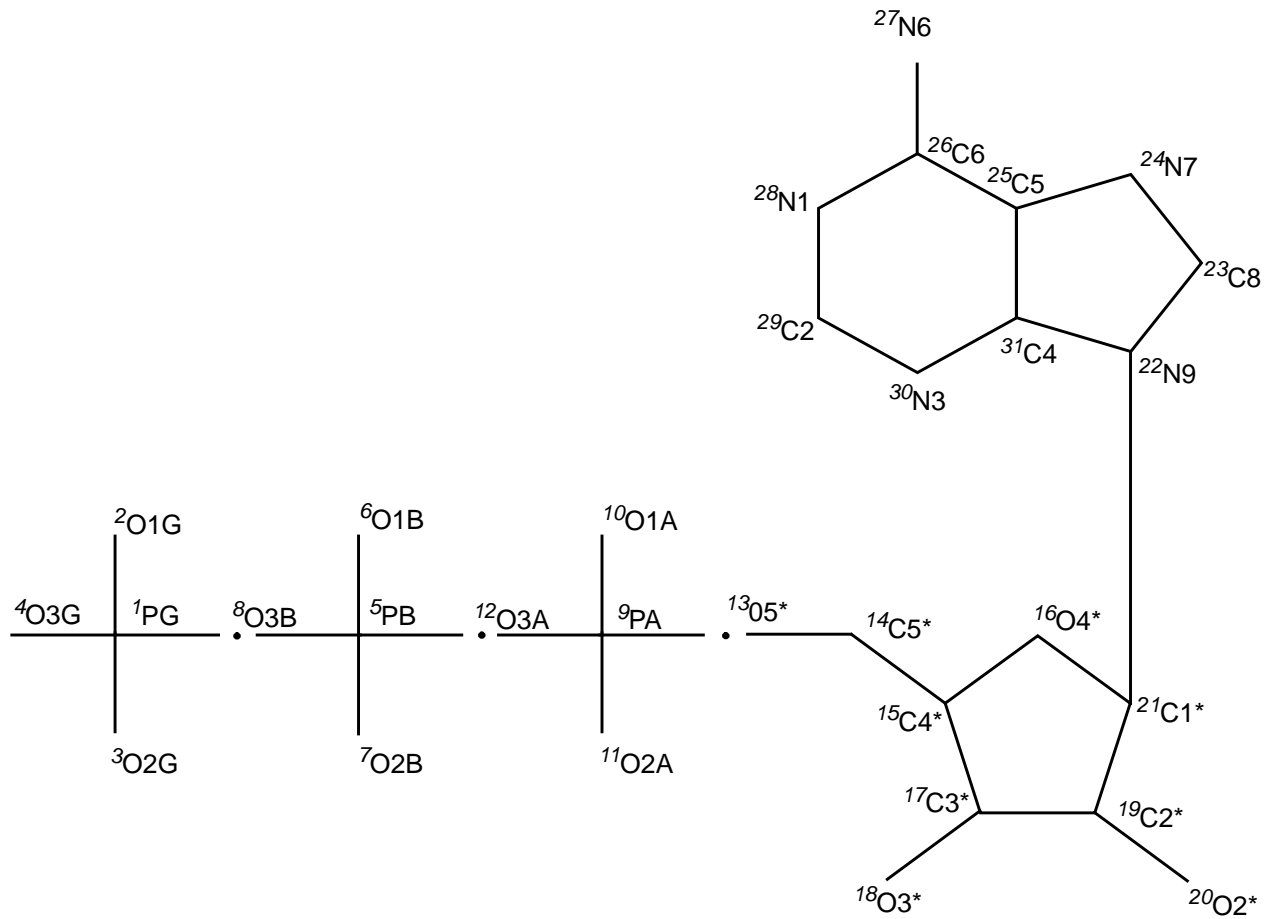
C. Non-Standard (HET) Groups

Because of the repeated occurrence of certain cofactors, prosthetic groups, etc., the almost complete lack of uniformity in the nomenclature assigned by depositors, and in the absence of any authoritative precedent, the Data Bank has assigned a standard nomenclature and ordering of the atoms in some of these groups. These assignments appear on the subsequent pages, for the following groups:

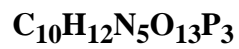
	<u>page</u>
ATP	30
Coenzyme A	31
Flavin mononucleotide (FMN)	32
Heme	33
Methotrexate	34
NAD	35

backbone**bases (names according to nucleoside)****Adenosine****Guanosine****Cytidine****Uridine****ATOM NAMES AND ORDER INDICATORS FOR THE COMMON RIBONUCLEOTIDES**

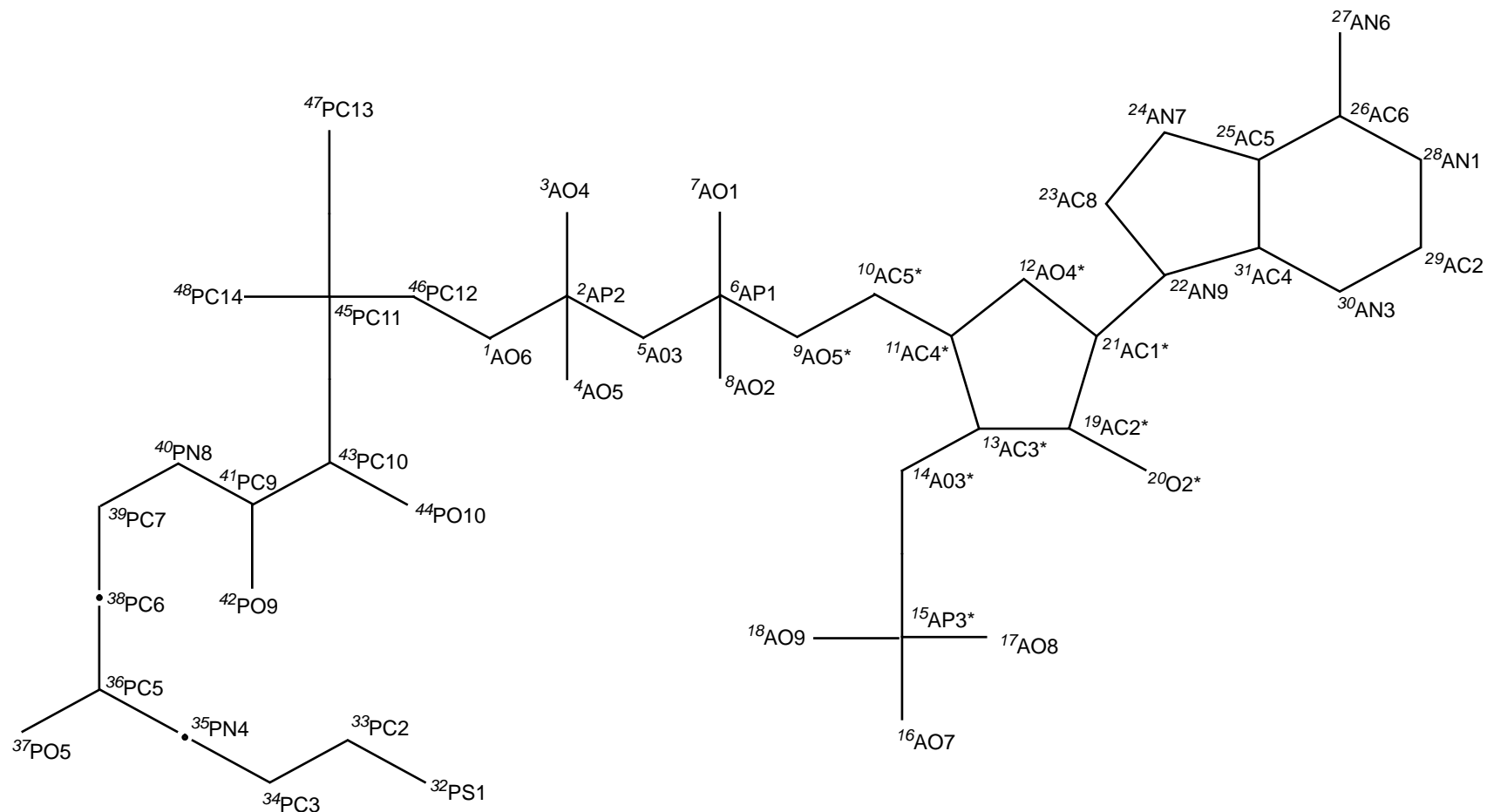
(Order indicators are given as preceding superscripts.)



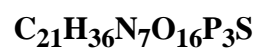
PROTEIN DATA BANK STANDARD NOMENCLATURE FOR ADENOSINE TRIPHOSPHATE (ATP)



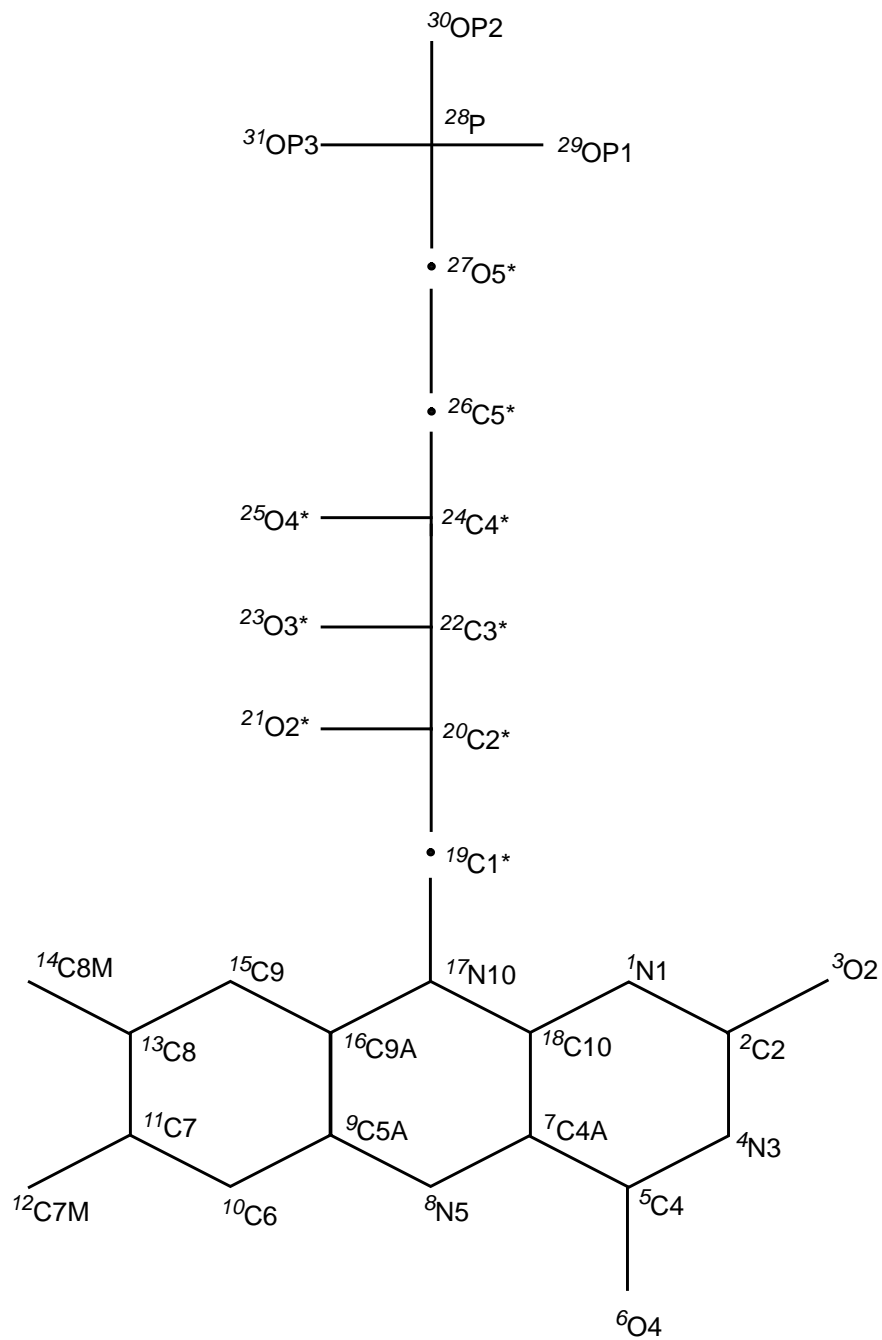
(Order indicators are given as preceding superscripts.)



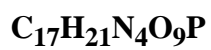
PROTEIN DATA BANK STANDARD NOMENCLATURE FOR COENZYME A



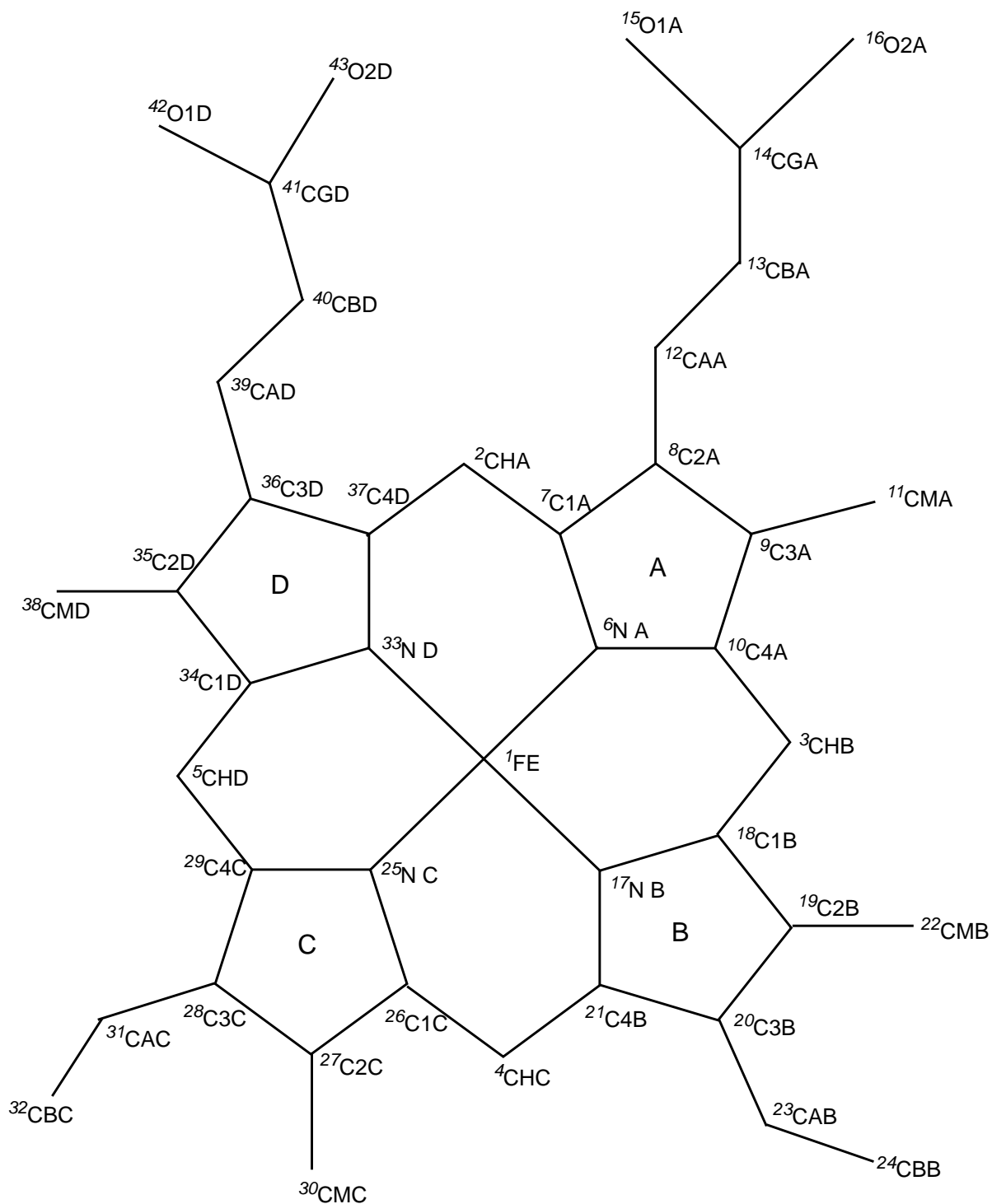
(Order indicators are given as preceding superscripts.)



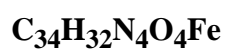
PROTEIN DATA BANK STANDARD NOMENCLATURE FOR FLAVIN MONONUCLEOTIDE (FMN)



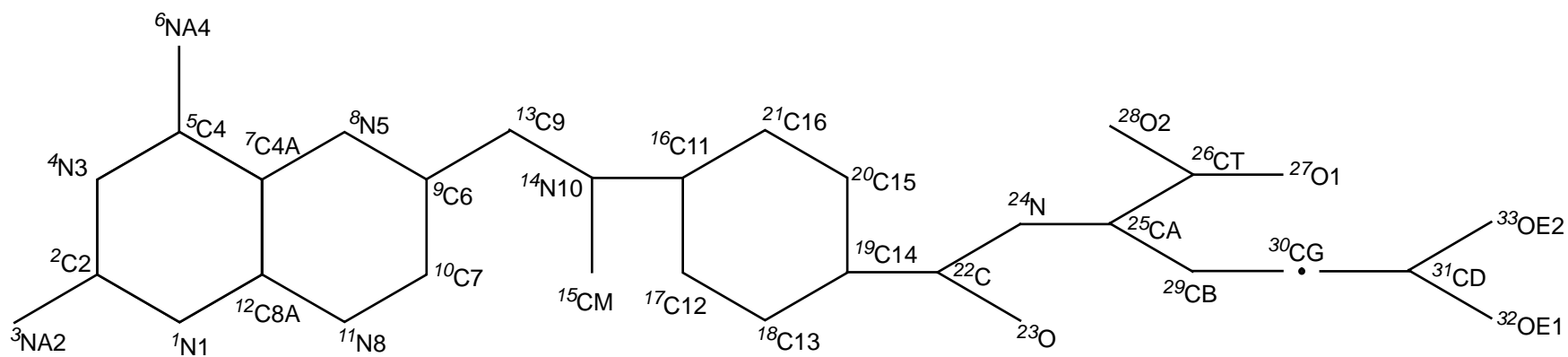
(Order indicators are given as preceding superscripts.)



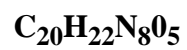
PROTEIN DATA BANK STANDARD NOMENCLATURE FOR A HEME GROUP



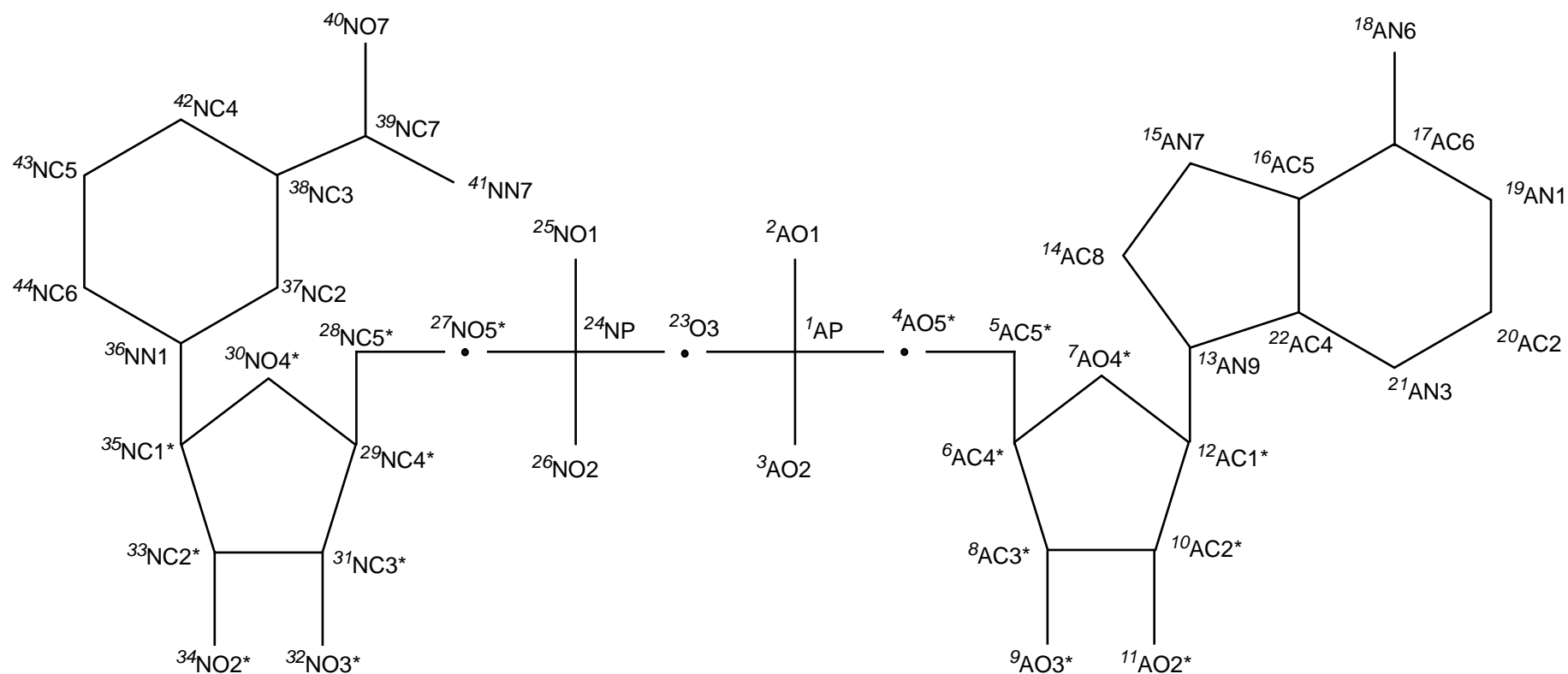
(Order indicators are given as preceding superscripts. Non-protein ligands of the iron atom are listed after those atoms given above.)



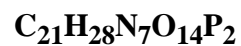
PROTEIN DATA BANK STANDARD NOMENCLATURE FOR METHOTREXATE



(Order indicators are given as preceding superscripts.)



PROTEIN DATA BANK STANDARD NOMENCLATURE FOR NICOTINAMIDE ADENINE DINUCLEOTIDE (NAD)



(Order indicators are given as preceding superscripts. For an NADP molecule the atoms of the extra phosphate group will be listed after those above.)

APPENDIX C - STANDARD RESIDUE NAMES AND ABBREVIATIONS

A. Amino Acids

<u>Residue</u>	<u>Abbr.</u>	<u>Synonym</u>
γ -Aminobutyric acid	ABU	
Acidic unknown	ACD	
Alanine.....	ALA	A
β -Alanine	ALB	
Aliphatic unknown.....	ALI	
Arginine	ARG	R
Aromatic unknown	ARØ	
Asparagine	ASN	N
Aspartic acid	ASP	D
ASP/ASN ambiguous	ASX	B
Basic unknown.....	BAS	
Cysteine	CYS.....	C,CYH,CSH
Cystine	CYS.....	C,CSS,CYX
Glutamine.....	GLN	Q
Glutamic acid.....	GLU	E
GLU/GLN ambiguous	GLX	Z
Glycine.....	GLY	G
Histidine.....	HIS	H
Hydroxyproline.....	HYP	
Isoleucine	ILE	I,ILU
Leucine.....	LEU.....	L
Lysine.....	LYS.....	K
Methionine	MET	M
Pyrrolidone carboxylic acid..... (pyroglutamate)	PCA.....	PGA
Phenylalanine.....	PHE.....	F
Proline.....	PRØ.....	P,PR0,PRZ
Sarcosine	SAR	
Serine	SER	S
Threonine	THR	T
Tryptophan.....	TRP	W,TRY
Tyrosine	TYR	Y
Valine.....	VAL	V

Notes:

- 1) Standard residue abbreviations conform to the IUPAC-IUB rules in *J. Biol. Chem.* 241, 527, 2491 (1966).
- 2) Recognizable synonyms, such as those above, will be changed to the standard abbreviation.
- 3) Non-standard residues (metals, prosthetic groups, etc.) are given a three-character designation which is defined in a special HET record. See page 9.
- 4) To avoid confusion here within residue abbreviations, the alphabetic character is written "Ø" and the numeric "0". This convention is *not* observed elsewhere throughout these specifications.

B. Nucleic Acids

Abbreviations conform to the IUPAC-IUB recommendations (*J. Biol. Chem.* 245, 5171 (1970)) for nucleosides with some extensions to cover the modified nucleosides and alterations because of character-set limitations.

Currently, the following abbreviations are in use for the indicated residues.

<u>Residue</u>	<u>Abbreviation</u>
Adenosine	A
1-Methyladenosine.....	1MA
Cytidine.....	C
5-Methylcytidine.....	5MC
2'-O-Methylcytidine.....	ØMC
Guanosine	G
1-Methylguanosine	1MG
N(2)-Methylguanosine.....	2MG
N(2)-Dimethylguanosine	M2G
7-Methylguanosine	7MG
2'-O-Methylguanosine	ØMG
Wybutosine	YG
Inosine.....	I
Thymidine.....	T
Uridine	U
Modified Uridine	+U
Dihydrouridine.....	H2U
Ribosylthymidine.....	5MU
Pseudouridine.....	PSU

Note: To avoid confusion here within residue abbreviations, the alphabetic character is written "Ø" and the numeric "0". This convention is *not* observed elsewhere throughout these specifications.

C. Miscellaneous

The following residue names are used to identify other commonly occurring groups.

<u>Residue</u>	<u>Abbr.</u>	<u>Synonym</u>
Acetyl.....	ACE	
Formyl.....	FØR	
Water.....	HØH.....	H2Ø,WAT,ØH2
Unknown.....	UNK	

Note: To avoid confusion here within residue abbreviations, the alphabetic character is written "Ø" and the numeric "0". This convention is *not* observed elsewhere throughout these specifications.

APPENDIX D - PROTEIN DATA BANK CONVENTIONS

Subsequent to the original adoption of the format described in this document, it has been decided to tighten the rules under which certain categories of information are presented. Specifications for the bibliographic citations given in the JRNL and REMARK 1 records are given in Appendix E. Concurrent with these changes it was deemed desirable to allow the availability of both upper- and lowercase characters to be exploited by inserting certain typesetting codes.

In addition to the detailed specifications given below the following general rules apply:

- (i) No word is to be hyphenated and split over two records.
- (ii) Only the surname of an author or editor is given in full; other names are indicated by initials only, e.g., A.B.Cooper.
- (iii) Blanks and hyphens are used in author lists only if they are properly part of a name (e.g., C.-I.Branden, C.J.Birkett-Clews, L.Riva di Sansaverino).
- (iv) The word Junior is written out in full.

Typesetting codes are kept to a minimum by a judicious choice of default conventions. In the text strings of COMPND, SOURCE, REF, TITL and PUBL records, all letters are lowercase unless preceded by one of the following characters: blank, comma, period, left parenthesis, or asterisk. The occurrence of a slash forces all succeeding letters to be uppercase until column 70 is reached or either a dollar sign or a hyphen (minus sign) is encountered.

Superscripts are initiated and terminated by double equal signs, e.g., S==2+==.

Subscripts are initiated and terminated by single equals signs, e.g., F=c=.

For author lists all characters are lowercase unless they are adjacent to a period, a comma, or preceded by an asterisk (*). A dollar sign (\$) is used to separate a lowercase character from a period or comma which otherwise would force uppercase.

Comments for specific record types follow:

- HEADER.** In cols. 11-50 of these records an attempt is made to assign the macromolecule to some functional class. No general classification scheme for biological macromolecules according to function yet exists (except for enzymes) and so the designation given here is intended to be informative rather than definitive. Its future use in indexing and subdividing the Data Bank is envisioned.

COMPND. }
SOURCE. } For these three records, the text portion of continuation lines begins in
AUTHOR. } col. 12, leaving col. 11 blank. Such continuation lines are numbered 2,3,
 etc. in col. 10. The first line in each of these records has col. 10 blank.

SEQRES. This set of records gives the number and sequence of residues in each chain of the particular macromolecule or complex under study. No cognizance of homologous molecules on which the residue sequence identifiers may be based is taken here. Residues which are present but not found in the structure analysis are listed, but residues removed from the chain termini (e.g., during an activation process) are not included. Residues excised from the chain (not at the termini, e.g., in α -chymotrypsin) are represented by EXC in the SEQRES records.

In general, if the macromolecule is composed of two or more chains which are commonly conceptualized as being logically separable, e.g., ribonuclease S, or papain with an oligopeptide inhibitor, then separate sets of SEQRES records are provided for each of these chains. If, however, these chains are usually thought of as comprising an integral unit (e.g., the three chains of α -chymotrypsin) a single set of SEQRES records is given.

If the residue sequence is unknown, the number of residues thought to comprise the molecule is entered in cols. 14-17, col. 10 contains '0' and cols. 20-22 contain 'UNK'.

SHEET. For the case of bifurcated sheets, or those containing split strands (i.e., one strand comprised of two distinct amino acid runs), sufficient redundant sheets are defined to accommodate the bifurcations. For the case illustrated below, two sheets would be given:

1 -----	
2 -----	-----7
3 -----	
4 -----	
5 -----	
6 -----	-----8
	-----9

In this sketch each line represents a continuous amino-acid run forming a strand of a sheet.

The strands labelled 1, 2, 3, 4, 5, 6 would comprise one sheet and those denoted 7, 3, 4, 5, 8, 9 another. This redundancy would be explicitly noted in a REMARK.

TURN. These records were originally set up to describe four-residue turns (β turns) but three-residue turns (γ turns) may be included with a notation.

SSBOND. Each pair of cysteine residues which participate in a disulfide bond is listed here. Intra-chain bonds are listed before inter-chain linkages. The amino acid with the lower sequence identifier is listed first in each intra-chain pair. For inter-chain pairs the cysteine which occurs first in the Data Bank entry is listed first.

CRYST1. The unit cell constants of the native crystals are given here unless explicitly stated otherwise. Native in this context means "underivatized" but if a derivative structure is solved as the native, e.g., tosyl elastase, then the cell constants of this pseudo-native macromolecule are given.

The Hermann-Mauguin space-group symbol is given without parentheses or slashes, e.g., P 43 21 2.

Confusion over the value to use for Z (number of molecules per cell) arises because of different conceptions of the meaning of "molecule". We have adopted the (crystallographic) convention that Z should equal the number of times the same polymeric chain is contained in the cell. In the case of different numbers of different chains per cell this will be explained in the REMARK section and Z will denote the number of the more populous species per cell.

- TVECT.** These records are used to denote the translation vectors which are used to build the infinite covalently connected structure of which the Protein Data entry is representative.
- ATOM.** The orthogonal Ångstrom coordinates stored are either those specified by the depositor or defined with respect to the default set of orthogonal axes (Appendix A). In the case that the stored coordinates are in orthogonal Ångstroms but not with respect to the default axial system, then this is explained in a REMARK. The occupancy and temperature factor fields will contain the default values 1.0 and 0.0 if these parameters were not deposited. Otherwise these fields will contain the supplied quantities in their original form, i.e., as fractional occupancy/isotropic thermal parameter (B) or electron count/atomic-radius form. If an atom is found in two or more locations (i.e., disordered) the records carrying the different coordinates for the atom in question occur together.
- HETATM.** Comments as above for the ATOM records apply. In order to avoid problems associated with the special characters ' and ", which are often employed for saccharide atomic nomenclature, the more standard characters * and \$ were employed in their place in entries released through January 1992. A uniform nomenclature and ordering (this may not be the same as that employed by the depositor) for the atoms of all non-standard groups is assigned. This nomenclature is illustrated for some commonly-occurring non-standard groups in Appendix B.
- TER.** These records are inserted after the carboxy-terminal (3'-terminal) residue of each polypeptide (nucleotide) chain *if the terminal residue is represented in the data set*. TER cards are also inserted to denote the ends of inhibitors or pseudo-substrates that are obtained by condensing like structural units present (e.g., peptides, oligonucleotides, oligosaccharides, etc.).
- CONNECT.** These records may be used to specify all linkages not implied by the primary structure. Bonds from the polymeric chain to any non-standard groups present are given here as are all covalent bonds within such groups. Cross-links between polymeric chains (e.g., disulfide bonds) are specified as are any other important linkages deemed worthy of inclusion by the depositor. The connectivity list given here is redundant in that each bond indicated is given twice, once with each of the two atoms involved specified in cols. 7-11. These CONNECT records occur in increasing order of the atom serial numbers they carry in cols. 7-11. The target-atom serial numbers carried on these records also occur in increasing order.

APPENDIX E - FORMATS FOR LITERATURE CITATIONS

References to published works from the depositor's laboratory and relating to the Data Bank entry may be carried in either the JRNL or REMARK 1 records. The subsidiary tag-words AUTH, TITL, EDIT, PUBL, and REFN are used as appropriate to indicate the information carried. The details of these specifications are identical for the JRNL and REMARK 1 records except that for each citation in the latter list a lead record is provided which carries the word REFERENCE in cols. 12-20 and a left-justified ordinal in cols. 22-23. The details are exemplified by a JRNL citation.

Cols.	1 - 4	JRNL
	13 - 16	AUTH (or EDIT)
	17 - 18	Continuation record number - blank for the first AUTH record of this citation and set to 2,3, etc. for succeeding records.
	20 - 70	Author list or editor list
Cols.	1 - 4	JRNL
	13 - 16	TITL
	17 - 18	Continuation record number
	20 - 70	Title of Article
Cols.	1 - 4	JRNL
	13 - 15	REF
	17 - 18	Continuation record number
	20 - 47	Name of publication (including section or series designation) ⁽ⁱ⁾
	50 - 51	V.
	53 - 55	Volume number
	57 - 61	First page number of article
	63 - 66	Year of publication
		If more than one REF record is necessary to carry the name of the publication, the volume number, page and date of publication is always carried on the first record.
Cols.	1 - 4	JRNL
	13 - 16	PUBL (this category is omitted for journal articles)

	17 - 18	Continuation record number
	20 - 70	Name of publisher and city of publication
Cols.	1 - 4	JRNL
	13 - 16	REFN
	20 - 23	ASTM
	25 - 30	Code from ASTM list ⁽ⁱ⁾
	33 - 34	Country of publication
	36 - 39	SSN or ISBN
	41 - 65	ISSN or ISBN number ⁽ⁱ⁾
	68 - 70	Code from Cambridge Crystallographic Data Centre CODEN list ⁽ⁱ⁾

⁽ⁱ⁾Note: A complete list of names, journals, other publications, and codes assigned to them is available upon request.

**APPENDIX F - FORMULAS AND MOLECULAR WEIGHTS
FOR STANDARD RESIDUES**

Note that these weights and formulas correspond to the unpolymerized state of the component. The elements of one water molecule are eliminated for each two components joined.

<u>Name</u>	<u>Code</u>	<u>Formula</u>	<u>Mol. Wt.</u>
Amino Acids			
Alanine.....	ALA	C3 H7 N1 O2	89.09
Arginine	ARG	C6 H14 N4 O2	174.20
Asparagine	ASN	C4 H8 N2 O3	132.12
Aspartic acid	ASP	C4 H7 N1 O4	133.10
ASP/ASN ambiguous	ASX	C4 H7 ^{1/2} N1 ^{1/2} O3 ^{1/2}	132.61
Cysteine	CYS	C3 H7 N1 O2 S1	121.15
Glutamine.....	GLN	C5 H10 N2 O3	146.15
Glutamic acid.....	GLU	C5 H9 N1 O4	147.13
GLU/GLN ambiguous	GLX	C5 H9 ^{1/2} N1 ^{1/2} O3 ^{1/2}	146.64
Glycine.....	GLY	C2 H5 N1 O2	75.07
Histidine.....	HIS	C6 H9 N3 O2	155.16
Isoleucine	ILE	C6 H13 N1 O2	131.17
Leucine.....	LEU	C6 H13 N1 O2	131.17
Lysine.....	LYS	C6 H14 N2 O2	146.19
Methionine	MET	C5 H11 N1 O2 S1	149.21
Phenylalanine.....	PHE	C9 H11 N1 O2	165.19
Proline.....	PRØ	C5 H9 N1 O2	115.13
Serine	SER	C3 H7 N1 O3	105.09
Threonine	THR	C4 H9 N1 O3	119.12
Tryptophan.....	TRP	C11 H12 N2 O2	204.23
Tyrosine	TYR	C9 H11 N1 O3	181.19
Valine.....	VAL	C5 H11 N1 O2	117.15
Undetermined.....	UNK	C5 H6 N1 O3	128.16

Nucleotides

Adenosine	A	C10 H14 N5 O7 P1	347.22
1-Methyladenosine.....	1MA	C11 H16 N5 O7 P1	361.25
Cytidine.....	C	C9 H14 N3 O8 P1	323.20
5-Methylcytidine.....	5MC	C10 H16 N3 O8 P1	337.23
2'-O-Methylcytidine	ØMC	C10 H17 N3 O8 P1	338.23
Guanosine	G	C10 H14 N5 O8 P1	363.22
1-Methylguanosine	1MG	C11 H16 N5 O8 P1	377.25
N(2)-Methylguanosine.....	2MG	C11 H16 N5 O8 P1	377.25
N(2)-Dimethylguanosine	M2G	C12 H18 N5 O8 P1	391.28
7-Methylguanosine	7MG	C11 H10 N5 O8 P1	377.25
2'-O-Methylguanosine.....	ØMG	C11 H16 N5 O8 P1	377.25
Wybutosine	YG	C21 H26 N6 O11 P1	587.48
Inosine.....	I	C10 H13 N4 O8 P1	348.21
Thymidine	T	C10 H15 N2 O8 P1	322.21
Uridine	U	C9 H13 N2 O9 P1	324.18
Dihydrouridine.....	H2U	C9 H15 N2 O9 P1	326.20
Ribosylthymidine.....	5MU	C10 H16 N2 O10 P1	355.22
Pseudouridine.....	PSU	C9 H13 N2 O9 P1	324.18

Miscellaneous

Acetic Acid	ACE	C2 H4 O2	60.05
Formic Acid	FØR	C1 H2 O2	46.03
Water.....	HØH	H2 O1	18.015

Note: To avoid confusion here within residue abbreviations, the alphabetic character is written "Ø" and the numeric "0". This convention is *not* observed elsewhere throughout these specifications.