# The Protein Data Bank: A Computer-based Archival File for Macromolecular Structures

The Protein Data Bank is a computer-based archival file for macromolecular structures. The Bank stores in a uniform format atomic co-ordinates and partial bond connectivities, as derived from crystallographic studies. Text included in each data entry gives pertinent information for the structure at hand (e.g. species from which the molecule has been obtained, resolution of diffraction data, literature citations and specifications of secondary structure). In addition to atomic co-ordinates and connectivities, the Protein Data Bank stores structure factors and phases, although these latter data are not placed in any uniform format. Input of data to the Bank and general maintenance functions are carried out at Brookhaven National Laboratory. All data stored in the Bank are available on magnetic tape for public distribution, from Brookhaven (to laboratories in the Americas), Tokyo (Japan), and Cambridge (Europe and worldwide). A master file is maintained at Brookhaven and duplicate copies are stored in Cambridge and Tokyo. In the future, it is hoped to expand the scope of the Protein Data Bank to make available co-ordinates for standard structural types (e.g. α-helix, RNA double-stranded helix) and representative computer programs of utility in the study and interpretation of macromolecular structures.

The Protein Data Bank† (1971,1973) was established in 1971 as a computer-based archival file for macromolecular structures. The purpose of the Bank is to collect, standardize, and distribute atomic co-ordinates and other data from crystallographic studies. As the number of solved protein and nucleic acid structures has grown to the point where some $10^7$ characters are necessary to represent the co-ordinate information currently held, the need for such a computer-readable file has become very clear, and demands for the Bank's services have increased accordingly. The Protein Data Bank is one of several data base activities in the field of crystallography, e.g. the Bibliographic (Kennard *et al.*, 1972) and Structural (Allen *et al.*, 1973) Data Files for organic and organometallic compounds, the Atlas of Macromolecular Structure on Microfiche (AMSOM) (Feldmann, 1977), the Bond Index to the Determination of Inorganic Crystal Structures (BIDICS)‡ and the Powder Diffraction File.§

## (a) Scope

The Protein Data Bank covers atomic co-ordinates, structure factors and phases from diffraction studies of macromolecules. Since most of this information is not generally published in the primary literature, the Bank depends for comprehensiveness on data supplied directly by the investigators. It is essentially a depository of data, held in computer-readable form, in contrast to other data banks that are based

† Protein Data Bank is a misnomer of historical origin, since the file now contains entries for a nucleic acid.

‡ I. D. Brown, Bond Index to the Determination of Inorganic Crystal Structures, McMaster University, Hamilton, Ontario, Canada, L8S 4M1.

§ American Society for Testing Materials, 1916 Race St., Philadelphia, PA, 19103, U.S.A.

<center>TABLE 1</center>

<center>*Protein data bank holdings*</center>

| IDENT CODE | MOLECULE | DEPOSITOR | STATUS CODE |
|---|---|---|---|
| 1ADK | ADENYLATE KINASE | G. SCHULZ | A |
| 1ADH | ALCOHOL DEHYDROGENASE (ADP-RIB) | C.-I. BRANDEN | |
| 2ADH | ALCOHOL DEHYDROGENASE (ORTHOPHEN) | C.-I. BRANDEN | |
| 2CHA | ALPHA-CHYMOTRYPSIN (TOSYL) | D. BLOW | R |
| 3CHA | ALPHA-CHYMOTRYPSIN | A. TULINSKY | |
| 1FAB | ANTIGEN BINDING FRAGMENT (NEW) | R. POLJAK | |
| 1REI | BENCE-JONES IMMUNOGLOBULIN REI | O. EPP, R. HUBER | |
| 1CPV | CALCIUM-BINDING PARVALBUMIN SET 6A | R. KRETSINGER | |
| 2CPV | CALCIUM-BINDING PARVALBUMIN SET 6H | R. KRETSINGER | |
| 3CPV | CALCIUM-BINDING PARVALBUMIN SET 6I | R. KRETSINGER | |
| 1CAB | CARBONIC ANHYDRASE B | K. KANNAN | |
| 1CAC | CARBONIC ANHYDRASE C | K. KANNAN | |
| 1CPA | CARBOXYPEPTIDASE A | W. LIPSCOMB | |
| 1CHG | CHYMOTRYPSINOGEN | J. KRAUT | |
| 2CNA | CONCANAVALIN A | G. REEKE, G. EDELMAN | N |
| 3CNA | CONCANAVALIN A | K. HARDMAN | R |
| 1B5C | CYTOCHROME B5 | F. S. MATHEWS | |
| 1CYT | CYTOCHROME C (ALBACORE, OXIDIZED) | R. DICKERSON | |
| 2CYT | CYTOCHROME C (ALBACORE, REDUCED) | R. DICKERSON | |
| 1CYC | CYTOCHROME C (BONITO, HEART) | M. KAKUDO | |
| 1C2C | CYTOCHROME C2 | J. KRAUT | |
| 155C | CYTOCHROME C550 | R. TIMKOVICH | |
| 1EST | ELASTASE | H. WATSON | |
| 1FDX | FERREDOXIN | L. JENSEN | |
| 1FXN | FLAVODOXIN (CLOSTRIDIUM MP) | M. LUDWIG | |
| 1GCH | GAMMA-CHYMOTRYPSIN | COHEN,DAVIES,SILVERTON | P |
| 1GPD | GLYCERALDEHYDE-3-P-DEHYDROGENASE(LOBSTR) | M. ROSSMANN | N |
| 2MHB | HEMOGLOBIN (HORSE, AQUO MET) | LADNER, HEIDNER, PERUTZ | RP |
| 1DHB | HEMOGLOBIN (HORSE, DEOXY) | M. PERUTZ, G. FERMI | |
| 1HHB | HEMOGLOBIN (HUMAN, DEOXY) | M. PERUTZ, G. FERMI | |
| 1FDH | HEMOGLOBIN (HUMAN, FETAL, DEOXY) | J. FRIER | |
| 1LHB | HEMOGLOBIN (LAMPREY) | W. HENDRICKSON | |
| 1YHX | HEXOKINASE (YEAST) BIII | T. STEITZ | B |
| 1HIP | HIGH POTENTIAL IRON PROTEIN | J. KRAUT | |
| 2LDH | LACTATE DEHYDROGENASE | M. ROSSMANN | PD |
| 3LDH | LACTATE DEHYDROGENASE/NAD/PYRUVATE | M. ROSSMANN | PD |
| 1LYZ | LYSOZYME (HEN EGG-WHITE, SET W2) | R. DIAMOND | P |
| 2LYZ | LYSOZYME (HEN EGG-WHITE, SET RS5D) | R. DIAMOND | P |
| 3LYZ | LYSOZYME (HEN EGG-WHITE, SET RS6A) | R. DIAMOND | P |
| 4LYZ | LYSOZYME (HEN EGG-WHITE, SET RS9A) | R. DIAMOND | P |
| 5LYZ | LYSOZYME (HEN EGG-WHITE, SET RS12A) | R. DIAMOND | P |
| 6LYZ | LYSOZYME (HEN EGG-WHITE, SET RS16) | R. DIAMOND | P |
| 1MDH | MALATE DEHYDROGENASE | L. BANASZAK | A |
| 1MBN | MYOGLOBIN (SPERM WHALE) | H. WATSON | |
| 2MBN | MYOGLOBIN (SPERM WHALE, MET) | T. TAKANO | |
| 3MBN | MYOGLOBIN (SPERM WHALE, DEOXY) | T. TAKANO | |
| 3PTI | PANCREATIC TRYPSIN INHIBITOR | R. HUBER | R |
| 8PAP | PAPAIN, NATIVE | J. DRENTH | R |
| 2PAP | PAPAIN (ACE-ALA-ALA-PHE-ALA, CYS-25) | J. DRENTH | |
| 3PAP | PAPAIN (CYS DERIV OF CYS-25) | J. DRENTH | |
| 4PAP | PAPAIN (OXIDIZED CYS-25) | J. DRENTH | |
| 5PAP | PAPAIN (TOS-LYS, CYS-25) | J. DRENTH | |
| 6PAP | PAPAIN (BZOXY-GLY-PHE-GLY, CYS-25) | J. DRENTH | |
| 7PAP | PAPAIN (BZOXY-PHE-ALA,CYS-25) | J. DRENTH | |
| 1PGK | PHOSPHOGLYCERATE KINASE (YEAST) | H. WATSON | A |
| 2PGK | PHOSPHOGLYCERATE KINASE (HORSE) | P. EVANS, D. PHILLIPS | B |
| 1PAB | PREALBUMIN (HUMAN, PLASMA) | S. OATLEY, D. PHILLIPS | |
| 1RNS | RIBONUCLEASE S | H. WYCKOFF | |
| 2RXN | RUBREDOXIN | L. JENSEN | ND |
| 1SNS | STAPHYLOCOCCAL NUCLEASE | F. A. COTTON, E. HAZEN | |
| 1SGB | STREPTOMYCES GRISEUS PROTEINASE B | M. JAMES | A |
| 1SBT | SUBTILISIN BPN' | J. KRAUT | |
| 2SBT | SUBTILISIN NOVO | J. DRENTH | |
| 1SOD | SUPEROXIDE DISMUTASE | J. AND D. RICHARDSON | A |
| 1TLN | THERMOLYSIN (UNREFINED) | B. MATTHEWS | |
| 2TLN | THERMOLYSIN (REFINED) | B. MATTHEWS | |
| 1SRX | THIOREDOXIN | B.-O. SODERBERG | A |
| 1TNA | TRANSFER RNA (YEAST, PHE) | J. SUSSMAN, S.-H. KIM | N |
| 2TNA | TRANSFER RNA (YEAST, PHE) | M. SUNDARALINGAM | P |
| 3TNA | TRANSFER RNA (YEAST, PHE) | JACK, LADNER, KLUG | P |
| 1TIM | TRIOSE PHOSPHATE ISOMERASE | I. WILSON, D. PHILLIPS | |
| 1PTN | TRYPSIN (NATIVE, PH8) | FEHLHAMMER,BODE,SCHWAGER | N |
| 2PTB | TRYPSIN(BENZAMIDINE INHIBITED, PH7) | FEHLHAMMER,BODE,SCHWAGER | RN |
| 1PTC | TRYPSIN/TRYPSIN INHIBITOR COMPLEX | BODE ET AL. | N |

STATUS CODES

| | |
|---|---|
| BLANK | STANDARD ENTRY AVAILABLE FOR DISTRIBUTION |
| A | ALPHA CARBON ATOMS ONLY |
| B | BACKBONE ONLY |
| D | NEW DATA HAS BEEN PROMISED |
| N | NEW ENTRY WITH DEPOSITOR FOR APPROVAL |
| P | IN PREPARATION |

on data abstracted from scientific publications. The Bank contains 77 atomic co-ordinate entries for 47 macromolecules (Table 1),† and 13 sets of structure factors and phases. The atomic co-ordinate entries, which include descriptive text and partial bond connectivities, conform to a uniform format (see below), but the structure factors and phases are stored in the format received from depositors. All co-ordinate entries are referred to depositors for verification, before being made available publicly through the Bank.

### (b) *Record structure of atomic co-ordinate entries*

Atomic co-ordinate entries consist of records each of 80 characters.‡ Using the punched card analogy, columns 1 to 6 contain a record type identifier, and columns 7 to 70 contain data.§ Columns 71 to 80 are normally blank, but may contain sequence information which is added by the library-file management program UPDATE¶ used to maintain the file on the Brookhaven CDC CYBER 70/76 computing system. In order to facilitate retrieval of data from the file, the first four characters of each record define the unique record type, and the syntax of each record is independent of the order of records within any entry for a particular macromolecule. (In the master file, this order is always fixed.) Atomic co-ordinate data contributed by depositors are processed into the standard format with program MACMOL,‖ which also subjects the data to certain nomenclature and connectivity checking procedures.

A sample partial entry for the protein ribonuclease S is shown in Table 2.†† The unique code 1RNS identifying this entry is given in the HEADER record, along with the date these data were entered into the Bank, and a provisional classification based on function, intended for future use in indexing and subdividing the file. Text giving the name of molecule, species from which it has been obtained, authors, literature citations, and other general description are presented in records COMPND through REMARK. SEQRES gives the amino acid sequence, and FTNOTE records are footnotes keyed to particular residues or atoms. Records HELIX through TURN describe the secondary structure as stated or approved by the depositor. Record CRYST1 defines the unit cell, while ORIGX and SCALE respectively give trans-formations relating the orthogonal Ångström co-ordinates stored in the file to those originally supplied by the depositor (these frequently are referred to an oblique or non-isometric system) and to standard crystallographic fractional co-ordinates. ATOM records give the IUPAC-IUB (1969) standard atom names (IUPAC-IUB, 1970), and residue abbreviations (IUPAC-IUB, 1971), along with sequence identifiers (cf. SEQRES, above), co-ordinates in Ångström units, and occupancies and thermal

---

† In addition to current co-ordinate entries shown in Table 1, the Bank contains obsolete entries (for adenylate kinase tosyl, α-chymotrypsin, concanavalin A, lactate dehydrogenase, horse methemoglobin, papain, rubredoxin, benzamidine-inhibited trypsin and pancreatic trypsin inhibitor), which have been superseded by later, more accurate data. These obsolete data are available on special request.

‡ Originally, the Bank used a 140-character format, similar to that employed in the protein refinement programs of Diamond (1966,1971). The 140-character format has been superseded by the 80-character format.

§ A detailed description of the file formats is available from Brookhaven on request.

¶ Control Data Corporation, UPDATE Reference Manual, Publication No. 60342500, Control Data Corporation, Arden Hills, Minnesota, 1974.

‖ G. J. B. Williams, unpublished. For the 140-character data, program PROIN by E. F. Meyer was utilized.

†† The file is organized in a similar way for proteins and nucleic acids, although certain differences exist, e.g. with regard to details of atom and residue names.

F. C. BERNSTEIN *ET AL.*

## TABLE 2

*Abbreviated sample atomic co-ordinate entry (ribonuclease S)*

```
HEADER     HYDROLASE (PHOSPHORIC DIESTER, RNA)      01-APR-73   1RNS
COMPND     RIBONUCLEASE-S (E.C. 3.1.4.22)
SOURCE     BOVINE (BOS TAURUS) PANCREAS
AUTHOR     F. M. RICHARDS AND H. W. WYCKOFF
JRNL       R.J. FLETTERICK AND H. W. WYCKOFF, PRELIMINARY REFINEMENT
JRNL       OF PROTEIN COORDINATES IN REAL SPACE, ACTA CRYST., VOL. A31,
JRNL       P698 (1975).
REMARK   1
REMARK   1 REFERENCE 1. F. M. RICHARDS AND H. W. WYCKOFF, ATLAS OF
REMARK   1  STRUCTURES FOR MOLECULAR BIOLOGY, VOL. 1. RIBONUCLEASE-S,
REMARK   1  CLARENDON PRESS (1973).
REMARK   1 REFERENCE 2. F. M. RICHARDS AND H. W. WYCKOFF, BOVINE
REMARK   1  PANCREATIC RIBONUCLEASE, THE ENZYMES, EDITED BY P. D.
REMARK   1  BOYER, VOL. IV, THIRD EDITION, P647, ACADEMIC PRESS (1971)
REMARK   1 REFERENCE 3. F. M. RICHARDS, H. W. WYCKOFF, W. D. CARLSON,
REMARK   1  N. M. ALLEWELL, B. LEE AND Y. MITSUI, PROTEIN STRUCTURE,
REMARK   1  RIBONUCLEASE-S AND NUCLEOTIDE INTERACTIONS, COLD SPRING
REMARK   1  HARBOR SYMPOSIA ON QUANTITATIVE BIOLOGY, VOL. XXXVI, P35
REMARK   1  (1971).
REMARK   1 REFERENCE 4. N. M. ALLEWELL AND H. W. WYCKOFF,
REMARK   1  CRYSTALLOGRAPHIC ANALYSIS OF THE INTERACTION OF CUPRIC
REMARK   1  ION WITH RIBONUCLEASE S, J. BIOL. CHEM., VOL. 246, P4657
REMARK   1  (1971).
REMARK   1 REFERENCE 5. H. W. WYCKOFF, D. TSERNOGLOU, A. W. HANSON,
REMARK   1  J. R. KNOX, B. LEE AND F. M. RICHARDS, THE THREE-
REMARK   1  DIMENSIONAL STRUCTURE OF RIBONUCLEASE-S. INTERPRETATION
REMARK   1  OF AN ELECTRON DENSITY MAP AT A NOMINAL RESOLUTION OF 2
REMARK   1  ANGSTROMS. J. BIOL. CHEM., VOL. 245, P305 (1970).
REMARK   1 REFERENCE 6. H. W. WYCKOFF, K. D. HARDMAN, N. M. ALLEWELL,
REMARK   1  T. INAGAMI, D. TSERNOGLOU, L. N. JOHNSON AND F. M.
REMARK   1  RICHARDS, THE STRUCTURE OF RIBONUCLEASE-S AT 6 ANGSTROM
REMARK   1  RESOLUTION, J. BIOL. CHEM., VOL. 242, P3749 (1967).
REMARK   2
REMARK   2 RESOLUTION. 2.0 ANGSTROMS.
REMARK   3
REMARK   3 REFINEMENT. BY A STEEPEST-DESCENTS PROCEDURE. REFER TO THE
REMARK   3  JRNL CITATION ABOVE.
REMARK   4
REMARK   4 THIS COORDINATE SET IS DESIGNATED 6D BY THE DEPOSITOR.
REMARK   5
REMARK   5 THE *S-PEPTIDE* (RESIDUES 1-20) WHICH FORMS A SEPARATE
REMARK   5 CHAIN FROM THE REMAINDER OF THE MOLECULE IS GIVEN THE
REMARK   5 CHAIN IDENTIFIER S.
SEQRES   1 S   20   LYS GLU THR ALA ALA ALA LYS PHE GLU ARG GLN HIS MET
SEQRES   2 S   20   ASP SER SER THR SER ALA ALA
SEQRES   1     104   SER SER SER ASN TYR CYS ASN GLN MET MET LYS SER ARG
SEQRES   2     104   ASN LEU THR LYS ASP ARG CYS LYS PRO VAL ASN THR PHE
SEQRES   3     104   VAL HIS GLU SER LEU ALA ASP VAL GLN ALA VAL CYS SER
SEQRES   4     104   GLN LYS ASN VAL ALA CYS LYS ASN GLY GLN THR ASN CYS
SEQRES   5     104   TYR GLN SER TYR SER THR MET SER ILE THR ASP CYS ARG
SEQRES   6     104   GLU THR GLY SER SER LYS TYR PRO ASN CYS ALA TYR LYS
SEQRES   7     104   THR THR GLN ALA ASN LYS HIS ILE ILE VAL ALA CYS GLU
SEQRES   8     104   GLY ASN PRO TYR VAL PRO VAL HIS PHE ASP ALA SER VAL
FTNOTE   1
FTNOTE   1 THE MAIN CHAIN AND MOST OF THE ASSOCIATED SIDE CHAINS ARE
FTNOTE   1 NOT WELL-DEFINED IN THE REGIONS OF RESIDUES 2, 65-72 AND
FTNOTE   1 119-123.
FTNOTE   2
FTNOTE   2 THE MAIN CHAIN IS VERY POORLY DEFINED OR NOT VISIBLE AT ALL
FTNOTE   2 IN THE ELECTRON DENSITY MAP IN THE REGIONS OF RESIDUES 1,
FTNOTE   2 18-20,21-23 AND   124.
HELIX    1  H1 THR S     3 MET S   13  1
HELIX    2  H2 ASN      24 ASN      34  1
```

Find authenticated court documents without watermarks at docketalarm.com.

TABLE 2—*continued*

```
HELIX     3  H3 SER     50  ALA     56  1
SHEET     1  S1 3 LYS   41  HIS     48  0
SHEET     2  S1 3 MET   79  THR     87 -1   N  ASN    44   O  CYS    84
SHEET     3  S1 3 ALA   96  LYS    104 -1   N  ASP    83   O  THR   100
SHEET     1  S2 4 LYS   61  ALA     64  0
SHEET     2  S2 4 ASN   71  SER     75 -1   N  VAL    63   O  CYS    72
SHEET     3  S2 4 HIS  105  GLU    111 -1   N  TYR    73   O  VAL   108
SHEET     4  S2 4 VAL  116  VAL    124 -1   O  ALA   109   N  VAL   118
TURN      1  T1 VAL    54  VAL     57      PSEUDO 3/10 HELIX
TURN      2  T2 ALA    56  SER     59      PSEUDO 3/10 HELIX
TURN      3  T3 CYS    65  GLY     68      BETW S1RNDS 1,2 OF SHEET S2
TURN      4  T4 THR    87  SER     90      END OF STRAND 2 OF SHEET S1
CRYST1  44.650    44.650     97.150   90.00   90.00 120.00 P 31 2 1        6
ORIGX1       1.000000  0.000000  0.000000       0.000000
ORIGX2       0.000000  1.000000  0.000000       0.000000
ORIGX3       0.000000  0.000000  1.000000       0.000000
SCALE1        .022306   .012931  0.000000       0.000000
SCALE2       0.000000   .025861  0.000000       0.000000
SCALE3       0.000000  0.000000   .010293       0.000000
ATOM      1  N   LYS S   1      -15.394    7.914   20.202  1.00   0.00      2
ATOM      2  CA  LYS S   1      -15.145    7.636   18.730  1.00   0.00      2
ATOM      3  C   LYS S   1      -14.982    6.107   18.763  1.00   0.00      2
ATOM      4  O   LYS S   1      -15.145    5.351   19.732  1.00   0.00      2
ATOM      5  CB  LYS S   1      -13.872    8.244   18.185  1.00   0.00      2
ATOM      6  CG  LYS S   1      -12.693    7.654   18.794  1.00   0.00      2
```

```
ATOM    927  N   ASP   121      -6.795   -9.247    7.034  1.00   0.00      1
ATOM    928  CA  ASP   121      -5.813   -9.425    5.935  1.00   0.00      1
ATOM    929  C   ASP   121      -6.217  -10.156    4.789  1.00   0.00      1
ATOM    930  O   ASP   121      -5.828   -9.850    3.652  1.00   0.00      1
ATOM    931  CB  ASP   121      -4.529  -10.015    6.648  1.00   0.00      1
ATOM    932  CG  ASP   121      -3.471   -9.503    5.687  1.00   0.00      1
ATOM    933  OD1 ASP   121      -3.320   -8.082    5.636  1.00   0.00      1
ATOM    934  OD2 ASP   121      -2.718  -10.333    4.799  1.00   0.00      1
ATOM    935  N   ALA   122      -7.049  -11.201    5.013  1.00   0.00      1
ATOM    936  CA  ALA   122      -7.965  -12.086    4.084  1.00   0.00      1
ATOM    937  C   ALA   122      -8.554  -13.331    4.724  1.00   0.00      1
ATOM    938  O   ALA   122      -8.495  -13.636    5.925  1.00   0.00      1
ATOM    939  CB  ALA   122      -6.991  -12.510    2.881  1.00   0.00      1
ATOM    940  N   SER   123      -8.885  -13.915    3.717  1.00   0.00      1
ATOM    941  CA  SER   123      -9.758  -15.155    3.627  1.00   0.00      1
ATOM    942  C   SER   123      -8.915  -16.127    2.880  1.00   0.00      1
ATOM    943  O   SER   123      -8.372  -15.812    1.810  1.00   0.00      1
ATOM    944  CB  SER   123     -10.877  -14.659    2.597  1.00   0.00      1
ATOM    945  OG  SER   123     -10.157  -14.035    1.530  1.00   0.00      1
ATOM    946  N   VAL   124      -8.845  -17.415    3.438  1.00   0.00      2
ATOM    947  CA  VAL   124      -8.591  -18.490    2.596  1.00   0.00      2
ATOM    948  C   VAL   124      -9.235  -18.381    1.209  1.00   0.00      2
ATOM    949  O   VAL   124      -8.580  -17.735     .377  1.00   0.00      2
ATOM    950  CB  VAL   124      -8.937  -19.929    3.162  1.00   0.00      2
ATOM    951  CG1 VAL   124      -9.135  -20.905    2.012  1.00   0.00      2
ATOM    952  CG2 VAL   124      -7.784  -20.573    4.226  1.00   0.00      2
ATOM    953  OXT VAL   124     -10.419  -19.165    1.046  1.00   0.00      2
TER     954      VAL   124
CONECT  196  195  644
CONECT  312  311  729
CONECT  448  447  844
CONECT  498  497  549
CONECT  549  498  548
CONECT  644  196  643
CONECT  729  312  728
CONECT  844  448  843
MASTER       36   10    0    3    7    4    0    6  952    2    8   10
END
```

# DOCKET ALARM

# Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

## Real-Time Litigation Alerts

Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

## Advanced Docket Research

With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

## Analytics At Your Fingertips

Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

## API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

### LAW FIRMS
Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

### FINANCIAL INSTITUTIONS
Litigation and bankruptcy checks for companies and debtors.

### E-DISCOVERY AND LEGAL VENDORS
Sync your system to PACER to automate legal marketing.

fastcase®
Smarter legal research.