

# Measuring the Impact of the Protein Data Bank

---

Helen M Berman

February 28, 2017

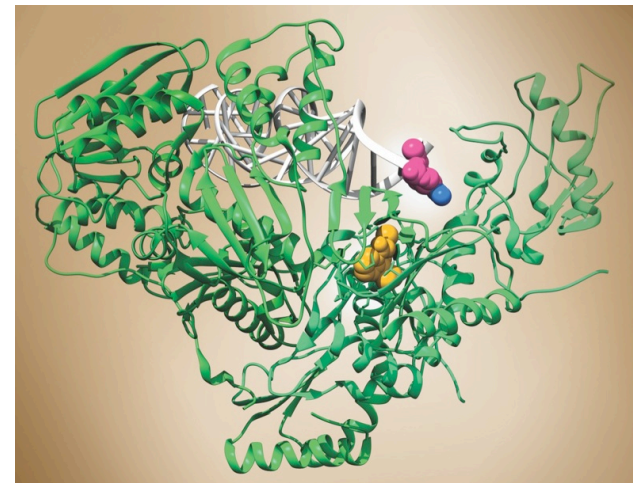
# Protein Data Bank

- First open access digital resource for biology data (est. 1971 with 7 entries)
- Single global archive of experimental 3D structures of biological macromolecules (>126,000 entries)
  - Primary data for structural biology, computational biology, drug discovery, ...
  - Complements GenBank and UniProt sequence database
- All data made freely available (primary users scientists and educators around the globe)
- Global archive of experimental macromolecular structure data central to biomedical research



**ABL tyrosine-kinase inhibited by Imatinib for treatment of chronic myeloid leukemia (CML).**

PDB ID 2hyy Cowan-Jacob et al. (2007) *Acta Crystallographica D*63: 80-93.



**HIV-1 reverse transcriptase complex with DNA and nevirapine**

PDB ID 3v81 Das et al. (2012) *Nature Structural and Molecular Biology*19: 253-259.

# Organizational Structure/Funding



- Partners share “Data In” responsibilities
  - Biocurate new depositions
  - Define deposition and annotation policies
  - Resolve data representation issues
  - Implement community validation standards
- Partners independently funded by each region
- Overseen by a wwPDB Advisory Committee
- Partners compete on “Data Out” resources

# PDB Archive Facts and Figures

## ■ Archival Contents

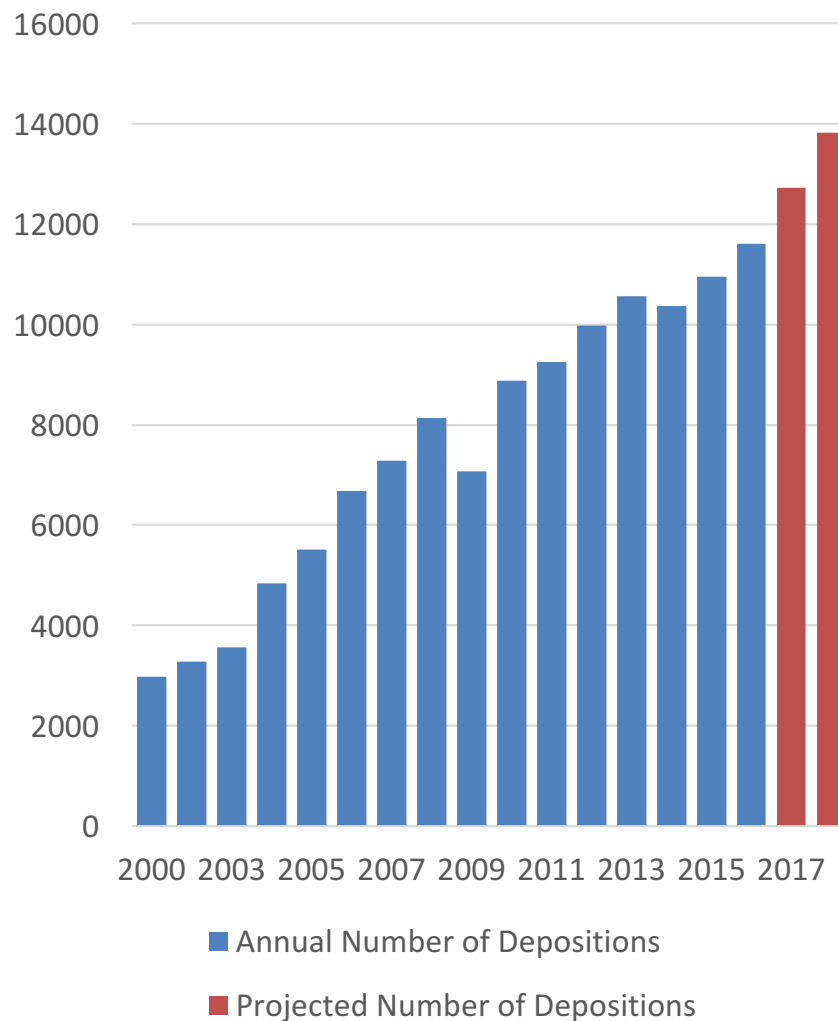
- ~126,000 Structures Released since 1971
- ~11,000 New Structures Deposited/Year

## ■ Global User Base

- ~30,000 Depositors Worldwide
- >1 Million Unique Visitors/Year from 192/195 UN-recognized sovereign nations

## ■ Impacts all of Biology and Medicine

- >590 Million Data Files Downloaded/Year
- ~1.6 Million Data Files Downloaded/Day
- >200 derived data resources repackage PDB data



# Download Statistics

Year	Total	Total FTP Archive	Total Website	RCSB PDB FTP Archive	RCSB PDB Website	PDBe FTP Archive	PDBe Website	PDBj FTP Archive	PDBj Website
2010	294,326,976	213,180,966	81,146,010	159,248,214	64,569,658	34,383,219	14,017,349	19,549,533	2,559,003
2011	383,131,048	276,952,286	106,178,762	204,939,406	81,560,098	40,960,368	18,515,245	31,052,512	6,103,419
2012	376,944,070	255,837,735	121,106,335	213,510,347	90,438,501	21,601,103	23,982,801	20,726,285	6,685,033
2013	441,262,210	296,176,290	145,085,920	215,331,908	97,549,580	43,684,850	37,762,496	37,159,532	9,773,844
2014	512,227,251	339,193,721	173,033,530	237,168,615	110,115,316	52,362,370	48,031,414	49,662,736	14,886,800
2015	534,339,871	368,244,766	166,095,105	255,346,630	111,802,897	48,544,330	41,127,219	64,353,806	13,164,989
2016	591,876,087	366,677,897	225,198,190	293,648,366	161,208,456	30,274,284	44,432,830	42,755,247	19,556,904

More than 1.6 million / day



Geographic origins of FTP downloads, 2012-2015

# Impact: Primary RCSB PDB Publication

© 2000 Oxford University Press

Nucleic Acids Research, 2000, Vol. 28, No. 1 235–242

Cited by 21459

Cited ~1500 times/year



## The Protein Data Bank

Helen M. Berman<sup>1,2\*</sup>, John Westbrook<sup>1,2</sup>, Zukang Feng<sup>1,2</sup>, Gary Gilliland<sup>1,3</sup>, T. N. Bhat<sup>1,3</sup>, Helge Weissig<sup>1,4</sup>, Ilya N. Shindyalov<sup>4</sup> and Philip E. Bourne<sup>1,4,5,6</sup>

<sup>1</sup>Research Collaboratory for Structural Bioinformatics (RCSB), <sup>2</sup>Department of Chemistry, Rutgers University, 610 Taylor Road, Piscataway, NJ 08854-8087, USA, <sup>3</sup>National Institute of Standards and Technology, Route 270, Quince Orchard Road, Gaithersburg, MD 20899, USA, <sup>4</sup>San Diego Supercomputer Center, University of California, San Diego, 9500 Gilman Drive, La Jolla, CA 92093-0505, USA, <sup>5</sup>Department of Pharmacology, University of California, San Diego, 9500 Gilman Drive, La Jolla, CA 92093-0500, USA and <sup>6</sup>The Burnham Institute, 10901 North Torrey Pines Road, La Jolla, CA 92037, USA

Received September 20, 1999; Revised and Accepted October 17, 1999

NEWS FEATURE

## THE TOP 100 PAPERS

Nature explores the most-cited research of all time.

BY RICHARD VAN WOODEN,  
RICHARD MARSH AND REGINA REEDS

**T**he discovery of high-temperature superconductors, the determination of DNA's double-helix structure, the first observation that the expansion of the Universe is accelerating – all of these breakthroughs were valued prizes and international acclaim. Yet none of the papers that announced them comes anywhere close to ranking among the 100 most highly cited papers of all time.

Citations, in which one paper refers to earlier works, are the standard measure by which authors acknowledge the source of their methods, ideas and findings, and are often used as a rough measure of a paper's importance. Fifty years ago, Eugene Garfield published the Science Citation Index (SCI), the first systematic effort to track citations in the scientific literature. To mark the anniversary, Nature asked Thomson Reuters, which now owns the SCI, to list the 100 most highly cited papers of all time. Over the fall that a web site (entitled top100) The search covered all of Thomson Reuters' Web of Science, an online version of the SCI that also includes databases covering the social sciences, arts and humanities, conference proceedings and more books. It lists papers published from 1900 to the present day. The exercise revealed some surprises, not least that it takes a staggering 12.19 citations to rank in the top 100 – and that many of the world's most famous papers do not make the cut. A few findings, such as the first observation

to other scientists what kind of work one is doing. Another common practice in science courses that truly foundational discoveries – Einstein's special theory of relativity, for instance – get fewer citations than they might deserve: they are so important that they quickly enter the textbooks and are incorporated into the main text of papers or terms devoted so familiar that they do not need a citation.

Citation counts are riddled with other confounding factors. The volume of citations has increased, for example – yet older papers have had more time to accrue citations. Biologists tend to cite one another's work more frequently than, say, physicists. And not all fields produce the same number of publications. Modern bibliometrics therefore resort from methods as crude as simply counting citations when they want to measure a paper's value. Instead, they prefer to compare counts for papers of similar age, and in comparable fields.

Not is Thomson Reuters' list the only ranking system available. Google Scholar compiled its own top-100 list for Nature. It is based on many more citations because the search engine calls references from a much greater (although poorly characterized) literature base, including from a large range of books. In that list, available at [www.citation.com/hot100](http://www.citation.com/hot100), economics papers have more prominence. Google Scholar's list also features books, which Thomson Reuters did not analyse. But among the science papers, many of the same titles show up.

Yet even with all the caveats, the old-fashioned hall of fame still has value. If nothing else, it serves as a reminder of the nature of scientific knowledge. To make existing advances, researchers rely on relatively young papers to learn the experimental methods, data and software.

Here Nature takes some of the key methods that tens of thousands of citations have bestowed on the top 100 research publications – essential, but rarely cited into the limelight.

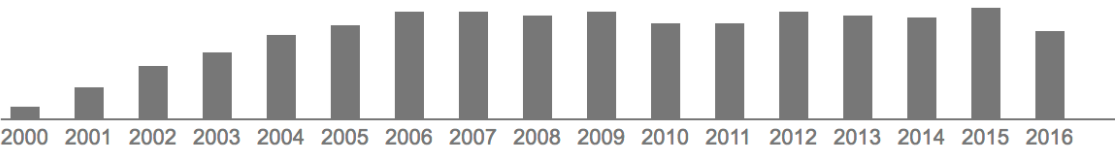
**BIOMOLECULAR TECHNIQUES**


For decades the top 100 list has been dominated by protein biochemistry. The 1951 paper 'describing the early method for purifying protein remains practically unreachably at number 1, even though many subsequent papers would represent just 1 continue at the peak. Only 1,499 papers – roughly a mere and a half's worth – have more than 1,000 citations (see 'The paper revolution'). Meanwhile, the foothills comprise works that have been cited only once, or at all – a group that comprises roughly half of the list.

Delicately fully understands what distinguishes the elite at the top from papers that are merely very well known – but researchers' custom explain some of it. Paul Wentworth, director of the Centre for Science and Technology Studies in London, the Netherlands, says that many methods papers 'become a standard reference that one cites in order to make clear



50 • NATURE • VOL 314 • 16 OCTOBER 2014  
© 2014 Macmillan Publishers Limited. All rights reserved



 Field: Research Areas	Record Count	% of 15137	Bar Chart
BIOCHEMISTRY MOLECULAR BIOLOGY	7907	52.236 %	<div></div>
CHEMISTRY	3075	20.314 %	<div></div>
BIOPHYSICS	2823	18.650 %	<div></div>
COMPUTER SCIENCE	2310	15.261 %	<div></div>
PHARMACOLOGY PHARMACY	1962	12.962 %	<div></div>
MATHEMATICAL COMPUTATIONAL BIOLOGY	1596	10.544 %	<div></div>
BIOTECHNOLOGY APPLIED MICROBIOLOGY	1258	8.311 %	<div></div>
SCIENCE TECHNOLOGY OTHER TOPICS	810	5.351 %	<div></div>
CRYSTALLOGRAPHY	762	5.034 %	<div></div>
PHYSICS	695	4.591 %	<div></div>
MATHEMATICS	648	4.281 %	<div></div>
CELL BIOLOGY	609	4.023 %	<div></div>
GENETICS HEREDITY	390	2.576 %	<div></div>
ENGINEERING	371	2.451 %	<div></div>
LIFE SCIENCES BIOMEDICINE OTHER TOPICS	240	1.586 %	<div></div>
MICROBIOLOGY	177	1.169 %	<div></div>
IMMUNOLOGY	166	1.097 %	<div></div>
MATERIALS SCIENCE	159	1.050 %	<div></div>
RESEARCH EXPERIMENTAL MEDICINE	120	0.793 %	<div></div>
PLANT SCIENCES	118	0.780 %	<div></div>
SPECTROSCOPY	112	0.740 %	<div></div>
POLYMER SCIENCE	96	0.634 %	<div></div>



# 3166 Patents Mention “protein data bank”

1 [9,476,035](#)  [Recombinant polymerases with increased phototolerance](#)

2 [9,475,886](#)  [Recombinant antibody composition](#)

3 [9,475,881](#)  [Antibody variants with enhanced complement activity](#)

4 [9,475,862](#)  [Neutralizing GP41 antibodies and their use](#)

5 [9,475,851](#)  [High MAST2-affinity polypeptides and uses thereof](#)

6 [9,475,847](#)  [Insecticidal proteins and methods for their use](#)

7 [9,474,759](#)  [Broad-spectrum antivirals against 3C or 3C-like proteases of picornavirus-like supercluster: picornaviruses, caliciviruses and coronaviruses](#)

8 [9,469,684](#)  [Therapeutic and diagnostic cloned MHC-unrestricted receptor specific for the MUC1 tumor associated antigen](#)

9 [9,468,660](#)  [Antinematodal methods and compositions](#)

10 [9,464,311](#)  [Method for identifying modulators of ubiquitin ligases](#)

11 [9,464,280](#)  [Beta-lactamases with improved properties for therapy](#)

12 [9,458,470](#)  [Recombinant influenza virus-like particles \(VLPs\) produced in transgenic plants expressing hemagglutinin](#)

13 [9,458,434](#)  [Mutant enzyme and application thereof](#)

14 [9,458,229](#)  [Immunogenic proteins and compositions](#)

15 [9,453,236](#)  [Polynucleotides and polypeptides involved in post-transcriptional gene silencing](#)

16 [9,453,224](#)  [MiRNA modulators of thermogenesis](#)

17 [9,453,019](#)  [Linked purine pterin HPPK inhibitors useful as antibacterial agents](#)

18 [9,452,222](#)  [Nucleic acids encoding modified relaxin polypeptides](#)

19 [9,452,210](#)  [Influenza virus-like particles \(VLPS\) comprising hemagglutinin produced within a plant](#)

20 [9,451,783](#)  [Phytase variants](#)

21 [9,447,157](#)  [Nitration shielding peptides and methods of use thereof](#)

22 [9,447,156](#)  [Methods and compositions for inhibiting neddylation of proteins](#)

23 [9,447,127](#)  [Synthetic lung surfactant and use thereof](#)

24 [9,446,121](#)  [Cloning of honey bee allergen](#)

25 [9,446,116](#)  [Peptide sequences and compositions](#)

26 [9,443,017](#)  [System and method for displaying search results](#)

[USPTO PATENT FULL-TEXT AND IMAGE DATABASE](#)

[Home](#)

[Quick](#)

[Advanced](#)

[Pat Num](#)

[Help](#)

[Next List](#)

[Bottom](#)

[View Cart](#)

Searching US Patent Collection...

Results of Search in US Patent Collection db for:

"protein data bank": 3166 patents.

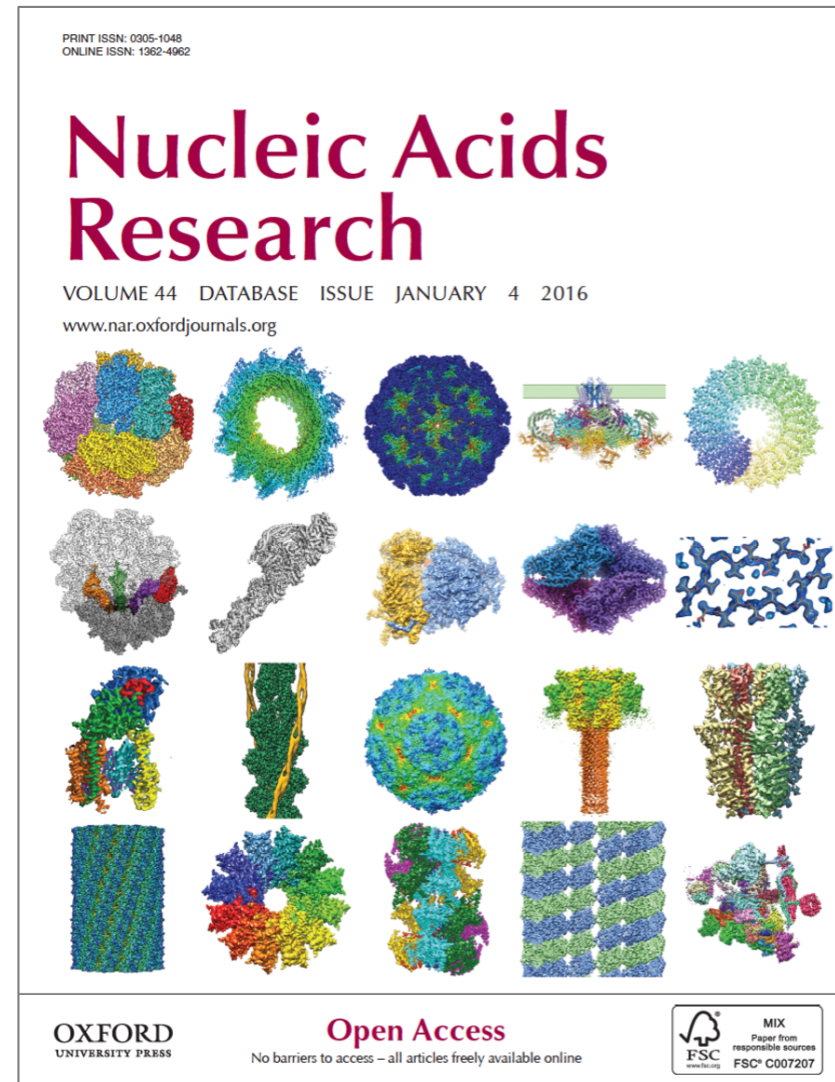
Hits 1 through 50 out of 3166

<http://patft.uspto.gov/>

Accessed October 26, 2016

# Impact: PDB Data Reuse

- PDB data used by >200 biological databases
  - Based on databases publishing in *NAR* 2011-2016
  - 11 Categories: Structure, Protein Sequence, Nucleotide Sequence, RNA Sequence, Genomics, Metabolic and Signaling, Human Genes and Diseases, Immunology, Proteomics, Plant, Other
- Since 2011, >25% of new databases utilize PDB data (119 out of 452 new databases)



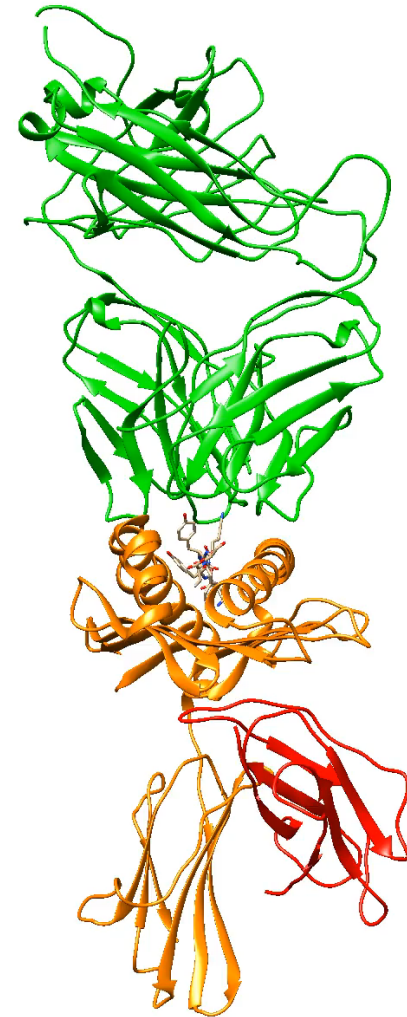


# Cost of Replicating the PDB Archive

- Data integrity and security are of paramount importance to the wwPDB partnership
- Estimated cost of replicating each PDB entry ranges from US\$50,000 to > US\$250,000
- Cost of replicating the PDB archive:  
**US\$12 billion**  
(assuming <unit cost>=US\$100,000)

# What Has the PDB Archive Enabled?

- Reproducibility and secure storage
- Accelerated structure determination technologies
- Understanding evolution in 3D
  - Structure classification and prediction
- Creation of structural bioinformatics as a discipline
- Structure-based drug discovery
- Functional understanding of biology at molecular and atomic levels



Antigen Presenting Cell meets the T-cell  
PDB 2CKB, Garcia *et al.* (1998)