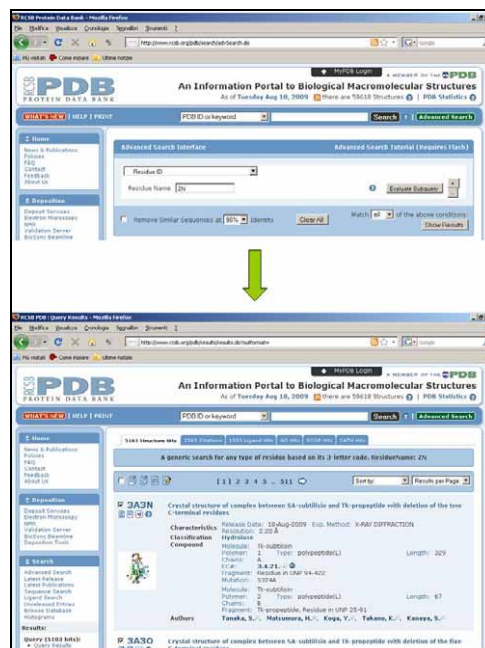
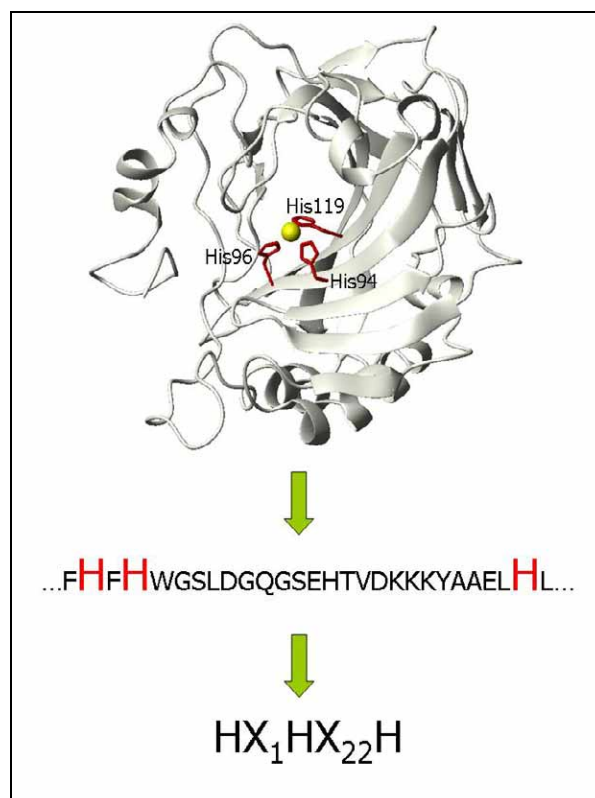


ESI. Example protocol for the identification of zinc proteins based on our method.

Step 1. All protein structures that bind zinc are retrieved from the PDB. A comprehensive retrieval strategy involves identifying all HET groups defined in PDB that contain zinc (see http://deposit.pdb.org/het_dictionary.txt) and then querying the PDB for structures that contain any of those HET groups. The Figure shows the result of a query for structures containing HET groups named “ZN”.



Step 2. After removing proteins where zinc is not present under physiological conditions, zinc-binding PDB structures are analysed to identify the amino acid residues that coordinate the metal. The zinc ligands are then mapped onto the protein sequence to define a metal-binding pattern. The Figure shows the zinc ligands and the corresponding pattern for human carbonic anhydrase II (PDB code 2ILI). Note that the removal of non-physiological zinc proteins relies on the analysis of literature, which is much facilitated by grouping PDB structures according to CATH and/or SCOP databases.



Step 3. Zinc-binding PDB structures are analysed to identify the protein domains that contain the zinc-binding sites. This is done by comparing their amino acid sequences against the Pfam database. The Figure shows that human carbonic anhydrase II (PDB code 2ILI) contains the Pfam domain named “Carb_anhydrase”, which contains the zinc-binding site. The “Carb_anhydrase” domain is thus taken as zinc-binding and associated with the HX₁HX₂₂H pattern (see Step 2). In this way, a library of zinc-binding Pfam domains associated with zinc-binding patterns is built.

Step 4. The Pfam database is queried for all domains whose annotation contains the word “zinc” (the Figure shows the result of such a query) or the word “Zn”. By these queries, one retrieves most (but not all) domains identified in Step 3, other domains which have been characterized as zinc-binding though a structure is not available, and other domains which are not in fact zinc-binding though the word “zinc” or “Zn” is mentioned in their annotation.

The screenshot shows the Pfam website interface. At the top, there are navigation links: HOME, SEARCH, BROWSE, FTP, HELP. The main heading is "Sequence search results". Below this, there is a table titled "Significant Pfam-A Matches".

Pfam-A	Description	Entry type	Sequence Start	End	HMM From	To	Bits score	E-value	Alignment mode	Predicted active sites	Show/hide alignment
Carb_anhydrase	Eukaryotic-type carbonic anhydrase	Domain	4	258	1	280	646.1	3.3e-191	ls	H63, N66, V155	Show

Below the table, there is a section for "Insignificant Pfam-A Matches" with a table showing one entry: VHS (VHS domain, Family) with a bits score of 7.0 and an E-value of 0.13.

A green arrow points from the "Carb_anhydrase" entry in the table to a green box containing the following information:

Pfam domain	Associated pattern
Carb_anhydrase	HX ₁ HX ₂₂ H

The screenshot shows the Pfam website interface. At the top, there are navigation links: HOME, SEARCH, BROWSE, FTP, HELP. The main heading is "Pfam 23.0 (July 2008, 10340 families)". Below this, there is a section for "Keyword search results".

The search results show 778 unique results for the query "Zinc". A table lists the database sections and the number of hits found in each one:

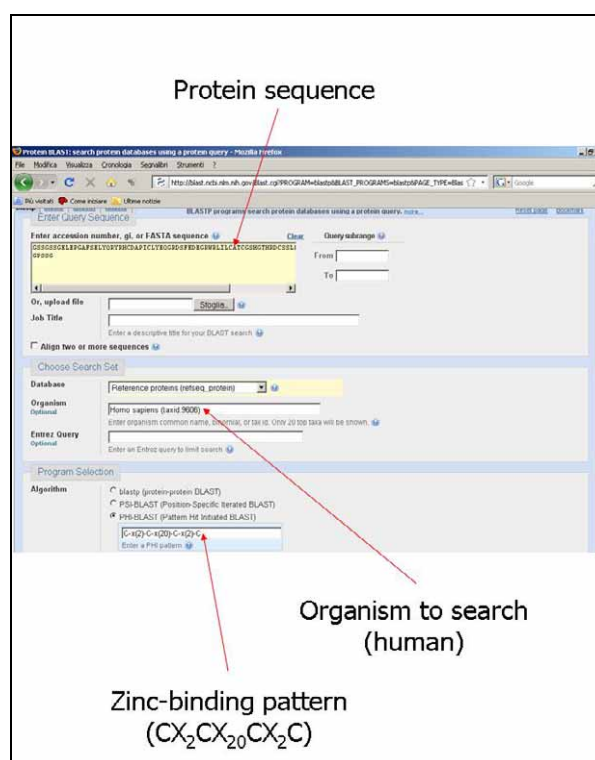
Section	Description	Number of hits
Pfam	Text fields for Pfam entries	205
Seq_info	Sequence description and species fields	2125
PROB	HEADER and TITLE records from PROB entries	540
GO	Gene ontology IDs and terms	105
Interpro	InterPro entry abstracts	301

Below this table, there is a table listing the Pfam domains found in the search results:

Accession	ID	Description	Pfam	Seq_Info	Prob	GO	Interpro
PF00908	Zn-CSD2	Zinc finger, CSD-type	✓	✓	✓	✓	✓
PF00927	Zn-CSD4	Zinc finger, CSD4-type (RING finger)	✓	✓	✓	✓	✓
PF00908	Zn-CSD3	Zinc finger, CSD3-type	✓	✓	✓	✓	✓
PF00246	Esptidase_M14	Zinc carboxypeptidase	✓	✓	✓	✓	✓
PF00320	GalT	GATA zinc finger	✓	✓	✓	✓	✓
PF00383	QOMP_cyt_deam_1	Cytidine and deoxycytidylate deaminase zinc-binding region	✓	✓	✓	✓	✓

Step 7. For proteins containing a zinc-binding domain associated with a zinc-binding pattern, sequences are checked for the occurrence of the pattern, and those that lack the pattern are filtered out. The Figure shows that, out of the 27 human proteins containing a Carb_anhydrase domain (see Step 6), 13 do not have the zinc-binding pattern, and are thus filtered out. These proteins correspond to the so-called carbonic anhydrase-related proteins which, despite their sequence homology to the catalytic isozymes, cannot bind zinc and are devoid of CO₂ hydration activity.

Step 8. Zinc proteins which do not contain any known Pfam domain, but have a zinc-binding site similar to that of a protein with known structure are searched using PHI-BLAST. PHI-BLAST input consists of the amino acid sequence of the zinc protein with known structure plus the zinc-binding pattern extracted from that structure. The Figure shows PHI-BLAST settings to search zinc proteins with sites similar to that of mouse PHD finger protein 7 (PDB code 1weq).



Step 9. PHI-BLAST hits retrieved in Step 8 are evaluated based on the I_d^{Global} parameter, defined as the ratio between the number of amino acids aligned by PHI-BLAST and the length of the sequence of the query protein. Proteins with $I_d^{\text{Global}} > 0.2$ are taken as zinc proteins. The Figure shows the PHI-BLAST alignment between mouse PHD finger protein 7 (PDB code 1weq) and human protein NP_060239.2, which is taken as a zinc protein ($I_d^{\text{Global}} = 0.39$) despite no zinc-binding Pfam domains could be detected.

Step 10. The predicted zinc proteome is obtained by collecting all the zinc proteins identified by (i) a zinc-binding domain with an associated zinc-binding pattern (e.g., the proteins with the Carb_anhydrase domain and the HX₁HX₂₂H pattern retrieved in Step 7), (ii) a zinc-binding domain only (e.g., the proteins with the zf-MYND domain retrieved in Step 6), and (iii) a zinc-binding pattern only (e.g., the NP_060239.2 protein retrieved in Step 9). In human, predicted zinc proteins represent about 10% of the entire proteome.

