

# All Proteins Have a Basic Molecular Formula

Homa Torabizadeh

**Abstract**—This study proposes a basic molecular formula for all proteins. A total of 10,739 proteins belonging to 9 different protein groups classified on the basis of their functions were selected randomly. They included enzymes, storage proteins, hormones, signalling proteins, structural proteins, transport proteins, immunoglobulins or antibodies, motor proteins and receptor proteins. After obtaining the protein molecular formula using the ProtParam tool, the H/C, N/C, O/C, and S/C ratios were determined for each randomly selected sample. In this case, H, N, O, and S coefficients were specified per carbon atom. Surprisingly, the results demonstrated that H, N, O, and S coefficients for all 10,739 proteins are similar and highly correlated. This study demonstrates that despite differences in the structure and function, all known proteins have a similar basic molecular formula  $C_nH_{1.58 \pm 0.015n}N_{0.28 \pm 0.005n}O_{0.30 \pm 0.007n}S_{0.01 \pm 0.002n}$ . The total correlation between all coefficients was found to be 0.9999.

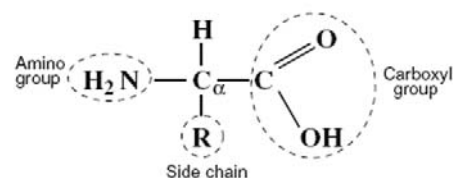
**Keywords**—Protein molecular formula, Basic unit formula, ProtParam tool.

## I. INTRODUCTION

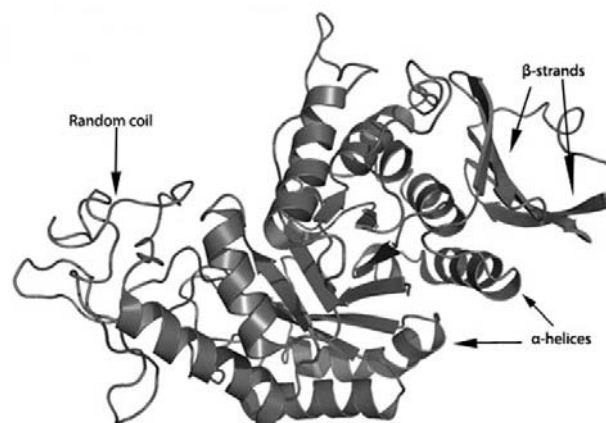
PROTEINS are macromolecules involved in many important biological structures and functions in living organisms. Examples of structural proteins are ligaments, fingernails, and hair; examples of functional proteins are digestive enzymes and muscle proteins that bring about metabolism and locomotion, respectively. Proteins are biopolymers constructed of amino acid monomers. The sequence of amino acids of a specific protein is determined by the sequence of the nucleic acid bases in the gene that encodes the protein.

The chemical properties of the constituent amino acids determine the biological activity of the protein (Fig. 1) [1]. Proteins play extremely important roles in biological systems [2]. The proteins present in single cell microorganisms, plants, animals, and humans differ greatly. The different number and sequence of amino acids result in the unique shapes, structures, and functions of proteins. Proteins are task and location-specific, and are broadly classified on the basis of structure and function. Structural classifications are based on protein folding, motifs, and protein family information. Functional classifications are based on biochemical and cellular roles, metabolic pathways, subcellular localization, and molecular interactions. Functional classifications are

H. Torabizadeh is with the Iranian Research Organization for Science and Technology, Institute of Chemical Technology, Department of Food Science and Technology, Enghelab Ave., No. 27, Forsat St., 15819 Tehran, Iran (corresponding author to provide phone: +98 88828052-7; fax:+98 2282276265; e-mail: htoraby@ alumni.ut.ac.ir).



The chemical structure of an amino acid



Secondary structure configurations

Fig. 1. Chemical structure of an amino acid and the secondary structure configurations of a protein.

further subdivided on the basis of similarities in enzyme reaction mechanisms, participation in biochemical pathways, functional roles, and cellular localization [3]. Although structural classifications are probably well defined on the basis of the criterion of molecular similarity, their overlap is surprisingly limited. On the other hand, functional classifications encompass many processes and elements, ranging from pathways to cellular compartments. These functional classifications have been shown to overlap considerably with each other [4]. Functionally, proteins are classified into the following: enzymes (proteins that catalyze chemical and biochemical reactions inside and outside living cells), storage proteins (proteins involved in storing energy that can be released during metabolic processes), hormones (proteins responsible for the regulation of many biochemical processes in organisms), signalling proteins (proteins involved in signal translation processes), structural proteins (these proteins are maintain structures of other biological components, like cells and tissues), transport proteins (proteins involved in transporting or storing chemical compounds and ions), immunoglobulins or antibodies (proteins involved in the immune response of an organism against large foreign molecules, which are introduced, for

example, by an infection), motor proteins (proteins involved in conversion of chemical energy to mechanical energy), and receptor proteins (proteins responsible for signal detection and translation into other type of signals). Each of these classified proteins acts on its own particular activities [5]. This study compares the coefficients of hydrogen, nitrogen, oxygen, and sulfur with respect to the quantity of carbon atoms in the protein. Data were obtained from bioinformatics databases using data extraction tools for more than 10,000 random protein samples. The results of this investigation indicate that a single basic molecular formula can be used to describe the relative amounts of elements present in all proteins.

## II. METHODS

### A. Databases and Definition of the Protein Formula

In this study, the protein formula was determined for each type of protein by using bioinformatics tools and databases. It was performed by first obtaining the amino acid sequence and entering it in an available calculation tool. This allows for the formula of the individual protein to be determined. All analyses were performed using the ExPasy select as a set of experimentally determined protein profiles in the UniProtKB

UniProtKB > UniProtKB Downloads Contact Documentation

A2QFN2 (A2QFN2\_ASPNC) ★ Unreviewed, UniProtKB/TrEMBL  
Last modified November 30, 2010. Version 26 History...

Clusters with 100%, 90%, 50% identity | Third-party data text xml rdf/xml gff |

Names Attributes Ontologies Sequences References Cross-refs Entry info Customize order

Names and origin

Protein names	Submitted name: Similarity to several putative oxidoreductases from various species. EC=1.-.-.- EMBL CAH37992.1
Gene names	ORF Names: An03g00140 EMBL CAH37992.1
Organism	Aspergillus niger (strain CBS 513.88 / FGSC A1513) [Complete proteome]

Sequences

Sequence	Length	Mass (Da)	Tools
A2QFN2-1 [UniParc]	FASTA 344	37,230	Blast go

Last modified March 6, 2007. Version 1.  
Checksum: 0FC08EE0E965957B

```

10 20 30 40 50 60
MSETARALI QPTAESNADD LTLQTDVVVN SNPARGHELI HVRACSPCAG ELLWPRNFPF
70 80 90 100 110 120
PFRKLLIPCP DVSGVVSASP PGSPFPQGA E IYARTSYARP GNARDYTIAT TDELARKPES
130 140 150 160 170 180
LTWVEAAVVP VSAETAHQEL FINAGIVPPE AVMDIPRAKL ANEKGKRLIVT AASGGVGLVW
190 200 210 220 230 240
VQLAARVLGV EVIGTCGPDN VDLVRSMGAK EAVNYRATDL KAWVEEEEE EGRRVVDVVD
250 260 270 280 290 300
CVGGRALEDA WWTVDGGVY VSIFQPKPES CPWKDPSYRG VKDIFVMEF SQNQLGAVTE
310 320 330 340
LIELGKCGQR VDSVNPLEHF RSAFERLATG HARGKIVFDL SLNR
    
```

Fig. 2. Input form of the searching tool UniProtKB, for oxidoreductase from *Aspergillus niger* (Accession code: A2QFN2) available at <http://www.uniprot.org/uniprot/A2QFN2>.

(Universal Protein Resource Knowledgebase) database <http://www.expasy.org/sprot/and> <http://www.expasy.org/tools> located in the ExPasy server. For example, in the case of oxidoreductase enzymes, one of the 800 randomly selected enzymes was the oxidoreductase from *Aspergillus niger* (UniProt ID: A2QFN2) (Fig. 2). This enzyme has 344 amino acids and a molecular mass of 37,230

Dalton. At first, the amino acid sequence of the enzyme was selected and copied. The sequence was then pasted onto the calculation window of the ProtParam tool for determination of the protein formula (Fig. 3).

ExPasy Proteomics Server

You are here: ExPasy CH > Tools > Primary structure analysis > ProtParam

ProtParam tool

ProtParam (References | Documentation) is a tool which allows the computation of various physical and chemical parameters for a given protein stored in Swiss-Prot or TrEMBL, molecular weight, theoretical pI, amino acid composition, atomic composition, extinction coefficient, estimated half-life, instability index, aliphatic index and grand average of hydropathy.

Please note that you may only fill out one of the following fields at a time.

Enter a Swiss-Prot/TrEMBL accession number (AC) (for example P05130) or a sequence identifier (ID) (for example KPC1\_DROME):  
A2QFN2

Or you can paste your own sequence in the box below:

```

MSETARALI QPTAESNADD LTLQTDVVVN SNPARGHELI HVRACSPCAG ELLWPRNFPF
PFRKLLIPCP DVSGVVSASP PGSPFPQGA E IYARTSYARP GNARDYTIAT TDELARKPES
GNARDYTIAT TDELARKPES LTVVEAAVVP VSAETAHQEL FINAGIVPPE
ANMDIPRAKL ANEKGKRLIVT AASGGVGLVW VQLAARVLGV EVIGTCGPDN
VDLVRSMGAK EAVNYRATDL KAWVEEEEE EGRRVVDVVD CVGGRALEDA
WWTVDGGVY VSIFQPKPES CPWKDPSYRG VKDIFVMEF SQNQLGAVTE
LIELGKCGQR VDSVNPLEHF RSAFERLATG HARGKIVFDL SLNR
    
```

RESET Compute parameters

Fig. 3. Input form of the characterization tool ProtParam, available at <http://www.expasy.org/tools/protparam.html>.

The computed protein formula was then extracted (Fig. 4).

ExPasy Proteomics Server

You are here: ExPasy CH > Tools > Primary structure analysis > ProtParam

ProtParam

User-provided sequence:

```

10 20 30 40 50 60
MSETARALI QPTAESNADD LTLQTDVVVN SNPARGHELI HVRACSPCAG ELLWPRNFPF
70 80 90 100 110 120
PFRKLLIPCP DVSGVVSASP PGSPFPQGA E IYARTSYARP GNARDYTIAT TDELARKPES
130 140 150 160 170 180
LTWVEAAVVP VSAETAHQEL FINAGIVPPE AVMDIPRAKL ANEKGKRLIVT AASGGVGLVW
190 200 210 220 230 240
VQLAARVLGV EVIGTCGPDN VDLVRSMGAK EAVNYRATDL KAWVEEEEE EGRRVVDVVD
250 260 270 280 290 300
CVGGRALEDA WWTVDGGVY VSIFQPKPES CPWKDPSYRG VKDIFVMEF SQNQLGAVTE
310 320 330 340
LIELGKCGQR VDSVNPLEHF RSAFERLATG HARGKIVFDL SLNR
    
```

References and documentation are available.

Please note the modified algorithm for extinction coefficient.

Number of amino acids: 344

Molecular weight: 37230.4

Theoretical pI: 5.34

Amino acid composition: CSV format

Total number of negatively charged residues (Asp + Glu): 43  
Total number of positively charged residues (Arg + Lys): 36

Atomic composition:

Carbon	C	1656
Hydrogen	H	2620
Nitrogen	N	460
Oxygen	O	494
Sulfur	S	11

Formula: C<sub>1656</sub>H<sub>2620</sub>N<sub>460</sub>O<sub>494</sub>S<sub>11</sub>  
Total number of atoms: 5241

Fig. 4. Input form of the characterization tool ProtParam, available at <http://www.expasy.org/cgi-bin/protparam>.

In this study, 10,739 different proteins were selected randomly from among 9 groups of proteins. Starting from the N-terminus of each protein sequence, a running window of n amino acids was selected. The specified amino acid sequence of each protein was then copied and pasted onto the calculation window of the ProtParam program. Eventually, in addition to the calculated profiles of the desired proteins, a

related molecular formula was determined. The total numbers of hydrogen, nitrogen, oxygen, and sulfur atoms were then divided by the total number of carbon atoms. In this manner, the coefficients of carbon, hydrogen, nitrogen, oxygen, and sulfur were obtained. These coefficients indicate the number of hydrogen, nitrogen, oxygen, and sulfur atoms in the protein structure with respect to the number of carbon atoms. This provides the ability to assemble a basic formula for each protein sample. Table I is an example of resulted formula for oxidoreductase enzyme.

TABLE I  
EXAMPLE OF RESULTED BASIC FORMULA FOR OXIDOREDUCTASE ENZYMES

Enzymes, Oxidoreductases: 800 samples		
Accession Code	Protein Formula	Basic Molecular Formula
A5D4R8	C <sub>3142</sub> H <sub>5015</sub> N <sub>903</sub> O <sub>907</sub> S <sub>30</sub>	C <sub>n</sub> H <sub>1.60n</sub> N <sub>0.29n</sub> O <sub>0.29n</sub> S <sub>0.010n</sub>
A5D560	C <sub>1926</sub> H <sub>3107</sub> N <sub>583</sub> O <sub>548</sub> S <sub>17</sub>	C <sub>n</sub> H <sub>1.61n</sub> N <sub>0.30n</sub> O <sub>0.28n</sub> S <sub>0.009n</sub>
A5D616	C <sub>1712</sub> H <sub>2741</sub> N <sub>517</sub> O <sub>481</sub> S <sub>18</sub>	C <sub>n</sub> H <sub>1.60n</sub> N <sub>0.30n</sub> O <sub>0.28n</sub> S <sub>0.011n</sub>
A5D684	C <sub>1791</sub> H <sub>2811</sub> N <sub>483</sub> O <sub>521</sub> S <sub>13</sub>	C <sub>n</sub> H <sub>1.57n</sub> N <sub>0.27n</sub> O <sub>0.29n</sub> S <sub>0.007n</sub>
A5D685	C <sub>1226</sub> H <sub>1925</sub> N <sub>347</sub> O <sub>364</sub> S <sub>16</sub>	C <sub>n</sub> H <sub>1.57n</sub> N <sub>0.28n</sub> O <sub>0.30n</sub> S <sub>0.013n</sub>
A5D686	C <sub>879</sub> H <sub>1426</sub> N <sub>250</sub> O <sub>260</sub> S <sub>3</sub>	C <sub>n</sub> H <sub>1.62n</sub> N <sub>0.28n</sub> O <sub>0.30n</sub> S <sub>0.003n</sub>
A5E8U7	C <sub>1944</sub> H <sub>3137</sub> N <sub>557</sub> O <sub>560</sub> S <sub>11</sub>	C <sub>n</sub> H <sub>1.61n</sub> N <sub>0.29n</sub> O <sub>0.29n</sub> S <sub>0.006n</sub>
A6SYB5	C <sub>1630</sub> H <sub>2616</sub> N <sub>448</sub> O <sub>477</sub> S <sub>4</sub>	C <sub>n</sub> H <sub>1.60n</sub> N <sub>0.27n</sub> O <sub>0.29n</sub> S <sub>0.002n</sub>
A6TA35	C <sub>1588</sub> H <sub>2514</sub> N <sub>418</sub> O <sub>474</sub> S <sub>15</sub>	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.26n</sub> O <sub>0.30n</sub> S <sub>0.009n</sub>
A6TAB8	C <sub>1765</sub> H <sub>2748</sub> N <sub>518</sub> O <sub>517</sub> S <sub>2</sub>	C <sub>n</sub> H <sub>1.56n</sub> N <sub>0.29n</sub> O <sub>0.29n</sub> S <sub>0.001n</sub>
A7K3A5	C <sub>980</sub> H <sub>1546</sub> N <sub>258</sub> O <sub>297</sub> S <sub>9</sub>	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.26n</sub> O <sub>0.30n</sub> S <sub>0.009n</sub>
A7K3A6	C <sub>870</sub> H <sub>1384</sub> N <sub>240</sub> O <sub>259</sub> S <sub>6</sub>	C <sub>n</sub> H <sub>1.59n</sub> N <sub>0.28n</sub> O <sub>0.30n</sub> S <sub>0.007n</sub>
A7K531	C <sub>1741</sub> H <sub>2722</sub> N <sub>482</sub> O <sub>521</sub> S <sub>12</sub>	C <sub>n</sub> H <sub>1.56n</sub> N <sub>0.28n</sub> O <sub>0.30n</sub> S <sub>0.007n</sub>
A7WHY8	C <sub>1323</sub> H <sub>2091</sub> N <sub>333</sub> O <sub>379</sub> S <sub>10</sub>	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.25n</sub> O <sub>0.29n</sub> S <sub>0.008n</sub>
A7ZQY9	C <sub>1264</sub> H <sub>2029</sub> N <sub>357</sub> O <sub>389</sub> S <sub>8</sub>	C <sub>n</sub> H <sub>1.61n</sub> N <sub>0.28n</sub> O <sub>0.31n</sub> S <sub>0.006n</sub>
A8A422	C <sub>1262</sub> H <sub>2027</sub> N <sub>357</sub> O <sub>387</sub> S <sub>8</sub>	C <sub>n</sub> H <sub>1.61n</sub> N <sub>0.28n</sub> O <sub>0.31n</sub> S <sub>0.006n</sub>
A9HMF9	C <sub>2681</sub> H <sub>4274</sub> N <sub>802</sub> O <sub>791</sub> S <sub>18</sub>	C <sub>n</sub> H <sub>1.59n</sub> N <sub>0.30n</sub> O <sub>0.30n</sub> S <sub>0.007n</sub>
A9HRN4	C <sub>2565</sub> H <sub>3995</sub> N <sub>747</sub> O <sub>760</sub> S <sub>23</sub>	C <sub>n</sub> H <sub>1.56n</sub> N <sub>0.29n</sub> O <sub>0.30n</sub> S <sub>0.009n</sub>
A9HX35	C <sub>1804</sub> H <sub>2804</sub> N <sub>516</sub> O <sub>532</sub> S <sub>11</sub>	C <sub>n</sub> H <sub>1.55n</sub> N <sub>0.29n</sub> O <sub>0.29n</sub> S <sub>0.006n</sub>
A9HXI4	C <sub>1629</sub> H <sub>2597</sub> N <sub>467</sub> O <sub>475</sub> S <sub>8</sub>	C <sub>n</sub> H <sub>1.59n</sub> N <sub>0.29n</sub> O <sub>0.29n</sub> S <sub>0.005n</sub>

### III. RESULTS

The resulting coefficients allow a basic formula to be extracted for each selected protein. Finally, a basic unit formula was obtained for each classification of protein (e.g., enzymes) by determining the average of the coefficients. The results are shown in Table II.

TABLE II  
FINAL AVERAGED FORMULA FOR EACH OF THE 9 PROTEIN CLASSIFICATIONS

Proteins	Samples	Basic Molecular Formula	
Oxidoreductase	800	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.27n</sub> O <sub>0.29n</sub> S <sub>0.008n</sub>	
Transferase	1030	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.28n</sub> O <sub>0.29n</sub> S <sub>0.008n</sub>	
Enzymes	Hydrolase	1030	C <sub>n</sub> H <sub>1.56n</sub> N <sub>0.27n</sub> O <sub>0.29n</sub> S <sub>0.007n</sub>
	Lyase	1050	C <sub>n</sub> H <sub>1.59n</sub> N <sub>0.28n</sub> O <sub>0.30n</sub> S <sub>0.008n</sub>
	Isomerase	280	C <sub>n</sub> H <sub>1.57n</sub> N <sub>0.275n</sub> O <sub>0.30n</sub> S <sub>0.018n</sub>
Ligase	500	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.28n</sub> O <sub>0.30n</sub> S <sub>0.009n</sub>	
Storage proteins	380	C <sub>n</sub> H <sub>1.57n</sub> N <sub>0.27n</sub> O <sub>0.30n</sub> S <sub>0.009n</sub>	
Hormones	1020	C <sub>n</sub> H <sub>1.57n</sub> N <sub>0.27n</sub> O <sub>0.296n</sub> S <sub>0.016n</sub>	
Signalling proteins	450	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.28n</sub> O <sub>0.30n</sub> S <sub>0.011n</sub>	
Structural proteins	1030	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.27n</sub> O <sub>0.30n</sub> S <sub>0.01n</sub>	
Transport proteins	1100	C <sub>n</sub> H <sub>1.59n</sub> N <sub>0.27n</sub> O <sub>0.29n</sub> S <sub>0.008n</sub>	
Immunoglobulines	612	C <sub>n</sub> H <sub>1.57n</sub> N <sub>0.28n</sub> O <sub>0.31n</sub> S <sub>0.007n</sub>	
Motor Proteins	693	C <sub>n</sub> H <sub>1.62n</sub> N <sub>0.28n</sub> O <sub>0.31n</sub> S <sub>0.008n</sub>	
Receptor Proteins	764	C <sub>n</sub> H <sub>1.56n</sub> N <sub>0.27n</sub> O <sub>0.29n</sub> S <sub>0.01n</sub>	
Average Formula	10739	C <sub>n</sub> H <sub>1.58n</sub> N <sub>0.28n</sub> O <sub>0.30n</sub> S <sub>0.01n</sub>	

The average of the coefficients that indicate the carbon, hydrogen, nitrogen, oxygen, and sulfur contents were determined, present in Table III. Then the total average of mean coefficients for each atom was calculated. As it is shown in table III. In this table, the total correlation between the coefficients was 0.9999.

TABLE III  
EXTRACTED DATASETS BASED ON THE AVERAGED COEFFICIENTS OF C, H, N, O, AND S FOR EACH OF THE 9 PROTEIN CLASSIFICATIONS

Proteins	Coefficients					
	C	H	N	O	S	
Enzymes	Oxidoreductase	1	1.58	0.27	0.30	0.008
	Transferase	1	1.58	0.28	0.29	0.008
	Hydrolase	1	1.56	0.27	0.29	0.007
	Lyase	1	1.59	0.28	0.30	0.008
	Isomerase	1	1.57	0.28	0.30	0.018
	Ligase	1	1.58	0.28	0.30	0.009
Storage proteins	1	1.57	0.27	0.30	0.009	
Hormones	1	1.57	0.27	0.30	0.016	
Signalling proteins	1	1.58	0.28	0.30	0.011	
Structural proteins	1	1.58	0.27	0.30	0.010	
Transport proteins	1	1.59	0.27	0.29	0.008	
Immunoglobulines	1	1.57	0.28	0.31	0.007	
Motor Proteins	1	1.62	0.28	0.31	0.008	
Receptor Proteins	1	1.56	0.27	0.29	0.010	
Average	1	1.58	0.28	0.30	0.010	
Total correlation	0.9999					

Now, the question may come to mind that, is it possible that the high similarity between the coefficients of these atoms be a result of the similarity between the coefficients of the 21 known amino acids which are present as building blocks in all proteins? To finding the answer of this question and evaluate the success of the prediction, it was necessary to compare the determined coefficients of protein elements (H, N, O, and S) with the coefficients that were obtained from the formula of individual amino acids. Calculations performed on the molecular formula of individual amino acids indicate that they are collectively not described by a single basic formula. Table IV, provides the molecular and basic formula based on the calculated coefficients for each of the amino acids.

TABLE IV

CLASSIFICATION AND BASIC FORMULA EXTRACTED FROM THE MAIN AMINO ACIDS IN PROTEINS

Amino acids	Molecular formula	Basic molecular formula of amino acids
Glycine	C <sub>2</sub> H <sub>5</sub> NO <sub>2</sub>	C <sub>n</sub> H <sub>2.5n</sub> N <sub>0.5n</sub> O <sub>n</sub>
Alanine	C <sub>3</sub> H <sub>7</sub> NO <sub>2</sub>	C <sub>n</sub> H <sub>2.33n</sub> N <sub>0.33n</sub> O <sub>0.67n</sub>
Valine	C <sub>5</sub> H <sub>11</sub> NO <sub>2</sub>	C <sub>n</sub> H <sub>2.2n</sub> N <sub>0.2n</sub> O <sub>0.4n</sub>
Leucine	C <sub>6</sub> H <sub>13</sub> NO <sub>2</sub>	C <sub>n</sub> H <sub>2.17n</sub> N <sub>0.17n</sub> O <sub>0.33n</sub>
Isoleucine	C <sub>6</sub> H <sub>13</sub> NO <sub>2</sub>	C <sub>n</sub> H <sub>2.17n</sub> N <sub>0.17n</sub> O <sub>0.33n</sub>
Serine	C <sub>3</sub> H <sub>7</sub> NO <sub>3</sub>	C <sub>n</sub> H <sub>2.33n</sub> N <sub>0.33n</sub> O <sub>n</sub>
Threonine	C <sub>4</sub> H <sub>9</sub> NO <sub>3</sub>	C <sub>n</sub> H <sub>2.25n</sub> N <sub>0.25n</sub> O <sub>0.75n</sub>
Cysteine	C <sub>3</sub> H <sub>7</sub> NO <sub>2</sub> S	C <sub>n</sub> H <sub>2.33n</sub> N <sub>0.33n</sub> O <sub>0.67n</sub> S <sub>0.33n</sub>
Cystine	C <sub>6</sub> H <sub>12</sub> N <sub>2</sub> O <sub>4</sub> S <sub>2</sub>	C <sub>n</sub> H <sub>2.0n</sub> N <sub>0.33n</sub> O <sub>0.67n</sub> S <sub>0.33n</sub>
Methionine	C <sub>5</sub> H <sub>11</sub> NO <sub>2</sub> S	C <sub>n</sub> H <sub>2.2n</sub> N <sub>0.2n</sub> O <sub>0.4n</sub> S <sub>0.2n</sub>
Phenylalanine	C <sub>9</sub> H <sub>11</sub> NO <sub>2</sub>	C <sub>n</sub> H <sub>1.22n</sub> N <sub>0.11n</sub> O <sub>0.22n</sub>
Tyrosine	C <sub>9</sub> H <sub>11</sub> NO <sub>3</sub>	C <sub>n</sub> H <sub>1.22n</sub> N <sub>0.11n</sub> O <sub>0.33n</sub>
Tryptophan	C <sub>11</sub> H <sub>12</sub> N <sub>2</sub> O <sub>2</sub>	C <sub>n</sub> H <sub>1.09n</sub> N <sub>0.18n</sub> O <sub>0.18n</sub>
Proline	C <sub>5</sub> H <sub>9</sub> NO <sub>2</sub>	C <sub>n</sub> H <sub>1.80n</sub> N <sub>0.2n</sub> O <sub>0.4n</sub>
Asparagine	C <sub>4</sub> H <sub>8</sub> N <sub>2</sub> O <sub>3</sub>	C <sub>n</sub> H <sub>2.0n</sub> N <sub>0.25n</sub> O <sub>0.75n</sub>
Glutamine	C <sub>5</sub> H <sub>10</sub> N <sub>2</sub> O <sub>3</sub>	C <sub>n</sub> H <sub>2.0n</sub> N <sub>0.2n</sub> O <sub>0.6n</sub>
Aspartic acid	C <sub>4</sub> H <sub>7</sub> NO <sub>4</sub>	C <sub>n</sub> H <sub>1.75n</sub> N <sub>0.25n</sub> O <sub>n</sub>
Glutamic acid	C <sub>5</sub> H <sub>9</sub> NO <sub>4</sub>	C <sub>n</sub> H <sub>1.80n</sub> N <sub>0.2n</sub> O <sub>0.8n</sub>
Lysine	C <sub>6</sub> H <sub>14</sub> N <sub>2</sub> O <sub>2</sub>	C <sub>n</sub> H <sub>2.33n</sub> N <sub>0.33n</sub> O <sub>0.33n</sub>
Histidine	C <sub>6</sub> H <sub>9</sub> N <sub>3</sub> O <sub>2</sub>	C <sub>n</sub> H <sub>1.5n</sub> N <sub>0.5n</sub> O <sub>0.33n</sub>
Arginine	C <sub>6</sub> H <sub>14</sub> N <sub>4</sub> O <sub>2</sub>	C <sub>n</sub> H <sub>2.33n</sub> N <sub>0.66n</sub> O <sub>0.33n</sub>

The numbers presented in Table V were obtained from coefficients that were determined for each atom in the molecular formula for each amino acid. Then, the average of each column was specified and total correlation between the resulting coefficients for each of the amino acids was calculated. This correlation was found to be 0.896, (Table V).

TABLE V

EXTRACTED DATASETS BASED ON THE COEFFICIENTS OF C, H, N, O, AND S FOR EACH OF THE 21 AMINO ACIDS AND THEIR TOTAL CORRELATIONS

Amino acids	Coefficients				
	C	H	N	O	S
Glycine	1	2.50	0.50	1.00	-
Alanine	1	2.33	0.33	0.67	-
Valine	1	2.20	0.20	0.40	-
Leucine	1	2.17	0.17	0.33	-
Isoleucine	1	2.17	0.17	0.33	-
Serine	1	2.33	0.33	1.00	-

TABLE V  
CONTINUE

Amino acids	Coefficients				
	C	H	N	O	S
Threonine	1	2.25	0.25	0.75	-
Cysteine	1	2.33	0.33	0.67	0.33
Cystine	1	2.00	0.33	0.67	0.33
Methionine	1	2.20	0.20	0.40	0.20
Phenylalanine	1	1.22	0.11	0.22	-
Tyrosine	1	1.22	0.11	0.33	-
Tryptophan	1	1.09	0.18	0.18	-
Proline	1	1.80	0.20	0.40	-
Asparagine	1	2.00	0.25	0.75	-
Glutamine	1	2.00	0.20	0.60	-
Aspartic acid	1	1.75	0.25	1.00	-
Glutamic acid	1	1.80	0.20	0.80	-
Lysine	1	2.33	0.33	0.33	-
Histidine	1	1.50	0.50	0.33	-
Arginine	1	2.33	0.66	0.33	-
Average	1	1.943	0.266	0.569	0.287
Total correlation	0.896				

The standard deviations of averaged coefficients of the amino acids and the averaged coefficients for all 10,739 proteins are shown in Table VI.

TABLE VI

THE STANDARD DEVIATIONS OF AVERAGED COEFFICIENTS FOR ALL PROTEIN SAMPLES AND AMINO ACIDS

	Proteins	Amino acids
H coefficients average	1.580	1.940
STDEV of H coefficients	0.015	0.410
N coefficients average	0.280	0.270
STDEV of N coefficients	0.005	0.140
O coefficients average	0.300	0.570
STDEV of O coefficients	0.007	0.260
S coefficients average	0.010	0.290
STDEV of S coefficients	0.002	0.075

Comparing the standard deviations of proteins and amino acid coefficients reveal that the data points tend to be very close to the mean as it mentions in equations 1, 2, 3, and 4. Unlike to the proteins, these values for amino acid coefficients are spread out over a larger range of values around the mean.

$$\Sigma_H / \Sigma_C = 1.58 \pm 0.015 \quad (1)$$

$$\Sigma_N / \Sigma_C = 0.28 \pm 0.005 \quad (2)$$

$$\Sigma_O / \Sigma_C = 0.30 \pm 0.007 \quad (3)$$

$$\Sigma_S / \Sigma_C = 0.01 \pm 0.002 \quad (4)$$

Thus, for 99.9% of the proteins examined, the molecular formula is:

$$C_n H_{1.58 \pm 0.015n} N_{0.28 \pm 0.005n} O_{0.30 \pm 0.007n} S_{0.01 \pm 0.002n}$$

For example, if the carbon coefficient be 100, the protein formula would be:

$$C_{100} H_{156.50 - 159.50} N_{27.50 - 28.50} O_{29.30 - 30.70} S_{0.8 - 1.20}$$

#### IV. CONCLUSIONS

As a frontier research area, bioinformatics has developed substantially over the past few decades. This study takes advantage of the bioinformatics databases and tools to specify a basic molecular formula for proteins, which indicates the ratio of the constituent elements C, H, N, O, and S. This study demonstrates that despite differences in the structure and function, all known proteins are constructed on the basis of a single basic molecular formula. It is thought that, the formation of all proteins by translation from mRNA is based on a unique pattern in which the ratio of the coefficients of the contributing elements with respect to each other is kept constant. It is believed that, all proteins are formed based on a similar basic molecular formula. There is significant diversity in transcription and rRNA-specific translation at the ribosome for production of proteins based on the specific genetic code because certain sequences of amino acids have completely different structural and functional properties. The survey results were quite surprising because, despite this extremely high level of diversity among the molecular structures of all proteins, for each 1.0 carbon atom there are 1.57–1.60 hydrogen atoms, 0.28–0.29 nitrogen atoms, 0.29–0.31 oxygen atoms, and 0.01 sulfur atom in all proteins. As a result, the smallest and the largest protein molecules have the same basic molecular formula. The specified formula may be applied to protein modelling prediction tasks for scientific research. Furthermore, another interesting research topic appears to be the development of a tool for designing new proteins and peptides for therapeutic, pharmaceutical, and industrial purposes. This is not a saturated research area and further research is ongoing.

#### ACKNOWLEDGMENT

Many thanks to Pegah Abdollahi and Zohre Torabian for their precious assistance in protein molecular formula determination.

#### REFERENCES

- [1] L. Jin, W. Fang, "On A comment on 'Prediction of protein structural classes by a new measure of information discrepancy,'" *Comput. Biol. Chem.*, vol. 33, no. 6, pp. 469–470, December. 2009.
- [2] L. Palopoli, S. E. Rombo, G. Terracina, G. Tradigo, P. Veltri, "Improving protein secondary structure predictions by prediction fusion," *Info. Fusion.*, vol. 10, no. 3, pp. 217–232, July. 2009.
- [3] C.A. Ouzounis, R.M.R. Coulson, A.J. Enright, V. Kunin, J.B. Pereira-Leal, "Classification schemes for protein structure and function," *Nat. Rev. Genet.*, vol. 4, no. 7, pp. 508–519, July. 2003.
- [4] S.C. Rison, T.C. Hodgman, J.M. Thornton, "Comparison of functional annotation schemes for genomes," *Funct. Integr. Genomics.*, vol. 1, no. 1, pp. 56–69, May. 2000.
- [5] <http://proteincrystallography.org/protein/>