

Construction of a Phylogenetic Tree Based on the Average Amino Acid Composition of Exomes

Esumi, Genshiro

Department of Pediatric Surgery, Hospital of the University of Occupational and Environmental Health, Kitakyushu, Japan

Abstract

All extant cellular organisms are believed to have descended from a common ancestor, and numerous phylogenetic trees have been constructed using various genetic and molecular data. In this study, I hypothesized that the average amino acid composition of an exome—computed across all exons—could serve as an index reflecting an organism’s characteristic amino acid usage and could be used to construct a phylogenetic tree. To test this hypothesis, I analyzed publicly available exome data from 81 species. For each species, I counted the fractional composition of each amino acid in each exon, and then calculated the average amino acid composition. I measured the pairwise distances between species using the angular distance based on these average compositions. Hierarchical clustering with Ward’s method was applied to construct a phylogenetic tree. The resulting tree showed a reasonable degree of similarity to previously established phylogenies, suggesting that this exome-based amino acid composition approach may offer some utility for inferring evolutionary relationships. To my knowledge, this is the first demonstration of constructing a phylogenetic tree solely from average exome-wide amino acid composition.

Keywords: Phylogenetic tree, Amino acid composition, Exome, Distance function, Evolution

Email: esumi@clnc.uoeh-u.ac.jp

Background

Throughout the history of evolutionary biology, numerous phylogenetic trees have been constructed under the assumption that all organisms are derived from a common ancestor, and that evolutionary change proceeds at a pace slow enough to be traced over time. Early efforts to infer evolutionary relationships focused largely on morphological phenotypes. However, with advancements in molecular biology, current approaches increasingly rely on evaluating changes in protein amino acid sequences, genes, and entire genome sequences [1].

Thanks to the completion of whole-genome sequencing projects for many species, the exomes of these organisms have become well-characterized. In this study, I hypothesize that the average amino acid composition of each exon within an exome can serve as a set of indices reflecting the organism's amino acid usage. I further assume that these indices have gradually diverged from a common ancestor. Under this assumption, I investigate whether a phylogenetic tree with at least some degree of validity can be reconstructed solely by calculating pairwise distance matrices of these average amino acid compositions across species.

Materials and Methods: Phylogenetic Analysis Based on Average Amino Acid Compositions of Exomes

Exome sequence data from 81 species listed in a publicly available dataset named “reference proteomes” on the EMBL-EBI website were used to calculate the average amino acid composition for each exome [2]. Briefly, all exon-derived amino acid sequences within the exome of each species were retrieved from this dataset. For each exon, the number of each amino acid was counted, and these counts were then summed exon by exon. Subsequently, to obtain the average amino acid composition of each exome, the total counts of each amino acid were normalized so that the sum of all fractions equals 1, by dividing by the total number of residues in each exon, and then averaged across all exons.

Next, an 81×81 distance matrix was constructed by computing the angular distance between each pair of species, based on their respective average amino acid compositions. The angular distance is defined as:

$$d(x, y) = \arccos\left(\frac{x \cdot y}{\|x\| \|y\|}\right).$$

This distance matrix was subsequently used to perform hierarchical clustering with Ward's method using JMP Pro 18 software (SAS Institute Inc., Cary, NC, USA). The resulting dendrogram was interpreted as a phylogenetic tree, illustrating the relationships among the 81 species based on their exome-wide amino acid composition patterns.

Results

Table 1 provides the scientific names of the 81 species examined in this study, the short labels used for the phylogenetic trees, and the number of annotated proteins registered for each species. Table 2 presents the average exome-wide composition of the 20 standard amino acids in each species' proteins. I generated a horizontal dendrogram based on these amino acid compositions and a radial dendrogram from the same dataset, shown in Figure 1 and Figure 2, respectively. In the figures, archaea are depicted in blue, bacteria in red, and eukaryotes in green. In addition, within the radial dendrogram (Figure 2), unicellular organisms are represented by small dots, whereas multicellular organisms are highlighted with larger dots and bold labels.

As a result of these procedures, a phylogenetic tree that appears reasonably valid at first glance was obtained. The validity of this tree will be evaluated in the subsequent Discussion section.

Discussion

In this study, I investigated whether it is possible to infer—or more precisely, restore—a phylogenetic tree based solely on the average amino acid composition of organismal exomes. Although my analysis was limited to 81 species, the resulting tree showed that groups of closely related organisms, as previously understood from existing phylogenies [3], largely clustered together. For example, mammals formed a distinct cluster, with primates (human, monkey, chimpanzee, and gorilla) and rodents (mouse and rat) each grouping on their own respective branches. These observations suggest that exome amino acid composition—at least in terms of its average values—shifts slowly enough to preserve certain phylogenetic relationships, allowing for the construction of such a tree (Figure 1, 2).

On the other hand, I also observed instances where organisms from entirely different lineages appeared on adjacent branches. One likely explanation is that while traditional phylogenetic reconstructions rely on highly multidimensional data (e.g., morphological or genetic diversity), the amino acid composition of exomes is comparatively low-dimensional. This may increase the chance

of coincidental similarities or convergence under evolutionary pressure. As a result, the overall validity or information content of the exome-based tree is not sufficient to replace or override existing phylogenies.

A further question arises as to the significance of an exome-based phylogenetic tree. Since the average amino acid composition does not account for protein abundance, it has often been presumed not to directly reflect the amino acid composition that an organism actually utilizes. However, my earlier work found that exome-wide amino acid composition tends to follow an approximately binomial (bell-shaped) distribution [4]. If the exome indeed exhibits such a distribution, then the proteome—composed of proteins synthesized from these exons—may be constrained by the exome’s amino acid composition. In the same report, I also proposed that this bell-shaped pattern could itself be influenced by constraints imposed by the proteome [4]. Hence, the exome and proteome likely exist in a mutually constraining relationship that gives rise to the exome’s bell-shaped distribution. If so, the phylogenetic tree constructed here may reflect how each organism’s cytoplasmic or cellular amino acid composition has evolved into its present form—that is, the evolutionary trajectory as seen from the perspective of amino acid composition.

An intriguing observation is that maize and a pathogenic fungus infecting maize lie on adjacent branches. This could indicate that the fungus converged toward maize’s amino acid composition, its nutritional source. In other words, it could be said that the exome of the pathogenic fungus may have adapted to the proteome of maize—a form of convergent evolution—although further study is required to clarify this possibility.

It has already been pointed out that an organism’s genome and GC content can influence its exome’s amino acid composition [5]. Even so, my results demonstrate that exome-level amino acid composition not only exhibits such variation, but also—albeit partially—contains the information needed to restore certain aspects of a phylogenetic tree. Moreover, since I found that each organism’s exome tends toward a bell-shaped distribution—implying convergence—exomes will continue to shift while maintaining a reciprocal constraint with their proteomes under evolutionary pressures [4]. Even if this concept appears straightforward once the underlying premises are considered, to my knowledge it has not been explicitly emphasized before, and thus may represent a novel insight. In this light, I hope this study provides a small yet meaningful contribution to understanding life’s complex evolutionary tapestry.

Reference

1. Kapli, P., Yang, Z., & Telford, M. J. (2020). Phylogenetic tree building in the genomic age. *Nature Reviews Genetics*, 21(7), 428–444. <https://doi.org/10.1038/s41576-020-0233-0>
2. EMBL-EBI. (2024). *Reference Proteomes* (Release 2024_02). Retrieved January 7, 2025, from https://www.ebi.ac.uk/reference_proteomes/
3. Siepel, A. (2009). Phylogenomics of primates and their ancestral populations. *Genome Research*, 19(11), 1929–1941. <https://doi.org/10.1101/gr.084228.108>
4. Esumi, G. (2023). The Distributions of Amino Acid Compositions of Proteins in an Organism's Proteome Uniformly Approximate Binomial Distributions. *Jxiv*. <https://doi.org/10.51094/jxiv.408>
5. Du, M.-Z., Liu, S., Zeng, Z., Alemayehu, L. A., Wei, W., & Guo, F.-B. (2018). Amino acid compositions contribute to the proteins' evolution under the influence of their abundances and genomic GC content. *Scientific Reports*, 8(1). <https://doi.org/10.1038/s41598-018-25364-1>

No.	Scientific Name	Short Label	Domain	Cell Organization	Protein Count
1	<i>Halobacterium salinarum</i> (strain ATCC 700922 / JCM 11081 / NRC-1) (Halobacterium halobium)	Halobacterium	archaea	unicellular	2427
2	<i>Thermococcus kodakarensis</i> (strain ATCC BAA-918 / JCM 12380 / KOD1) (Pyrococcus kodakaraensis (strain KOD1))	Thermococcus	archaea	unicellular	2301
3	<i>Methanoscarcina acetivorans</i> (strain ATCC 35395 / DSM 2834 / JCM 12185 / C2A)	Methanoscarcina	archaea	unicellular	4468
4	<i>Methanocaldococcus jannaschii</i> (strain ATCC 43067 / DSM 2661 / JAL-1 / JCM 10045 / NBRC 100440) (Methanococcus jannaschii)	Methanocaldococcus	archaea	unicellular	1787
5	<i>Saccharolobus solfataricus</i> (strain ATCC 35092 / DSM 1617 / JCM 11322 / P2) (Sulfolobus solfataricus)	Saccharolobus	archaea	unicellular	2937
6	<i>Korarchaeum cryptofolium</i> (strain OPF8)	Korarchaeum	archaea	unicellular	1602
7	<i>Nitrosopumilus maritimus</i> (strain SCM1)	Nitrosopumilus	archaea	unicellular	1795
8	<i>Mycobacterium tuberculosis</i> (strain ATCC 25618 / H37Rv)	M. tuberculosis	bacteria	unicellular	3999
9	<i>Escherichia coli</i> (strain K12)	E. coli	bacteria	unicellular	4416
10	<i>Helicobacter pylori</i> (strain ATCC 700392 / 26695) (Campylobacter pylori)	Helicobacter	bacteria	unicellular	1554
11	<i>Streptomyces coelicolor</i> (strain ATCC BAA-471 / A3(2) / M145)	Streptomyces	bacteria	unicellular	8039
12	<i>Neisseria meningitidis</i> serogroup B (strain MC58)	Meningococcus	bacteria	unicellular	2001
13	<i>Leptospira interrogans</i> serogroup Icterohaemorrhagiae serovar Lai (strain 56601)	Leptospira	bacteria	unicellular	3676
14	<i>Fusobacterium nucleatum</i> subsp. <i>nucleatum</i> (strain ATCC 25586 / DSM 15643 / BCRC 10681 / CIP 101130 / JCM 8532 / KCTC 2640 / LMG 13131 / VPI 4355)	Fusobacterium	bacteria	unicellular	2046
15	<i>Pseudomonas aeruginosa</i> (strain ATCC 15692 / DSM 22644 / CIP 104116 / JCM 14847 / LMG 12228 / IC / PRS 101 / PA01)	Pseudomonas	bacteria	unicellular	5564
16	<i>Bacillus subtilis</i> (strain 168)	B. subtilis	bacteria	unicellular	4267
17	<i>Aquifex aeolicus</i> (strain VF5)	Aquifex	bacteria	unicellular	1553
18	<i>Bradyrhizobium diazoefficiens</i> (strain JCM 10833 / BCRC 13528 / IAM 13628 / NBRC 14792 / USDA 110)	Bradyrhizobium	bacteria	unicellular	8253
19	<i>Bacteroides thetaiotaomicron</i> (strain ATCC 29148 / DSM 2079 / JCM 5827 / CCUG 10774 / NCTC 10582 / VPI-5482 / E50)	Bacteroides	bacteria	unicellular	4782
20	<i>Rhodopirellula baltica</i> (strain DSM 10527 / NCIMB 13988 / SH1)	Rhodopirellula	bacteria	unicellular	7271
21	<i>Deinococcus radiodurans</i> (strain ATCC 13939 / DSM 20539 / JCM 16871 / CCUG 27074 / LMG 4051 / NBRC 15346 / NCIMB 9279 / VKM B-1422 / R1)	Deinococcus	bacteria	unicellular	3084
22	<i>Geobacter sulfurreducens</i> (strain ATCC 51573 / DSM 12127 / PCA)	Geobacter	bacteria	unicellular	3402
23	<i>Mycoplasma genitalium</i> (strain ATCC 33530 / DSM 19775 / NCTC 10195 / G37) (Mycoplasmodium genitalium)	Mycoplasma	bacteria	unicellular	483
24	<i>Thermotoga maritima</i> (strain ATCC 43589 / DSM 3109 / JCM 10099 / NBRC 100826 / MSB8)	Thermotoga	bacteria	unicellular	1852
25	<i>Gloeobacter violaceus</i> (strain ATCC 29082 / PCC 7421)	Gloeobacter	bacteria	unicellular	4406
26	<i>Chlamydia trachomatis</i> (strain D/UW-3/Cx)	Chlamydia	bacteria	unicellular	895
27	<i>Thermodesulfovibrio yellowstonii</i> (strain ATCC 51303 / DSM 11347 / YP87)	Thermodesulfovibrio	bacteria	unicellular	1982
28	<i>Chloroflexus aurantiacus</i> (strain ATCC 29366 / DSM 635 / J-10-ff)	Chloroflexus	bacteria	unicellular	3850
29	<i>Dictyoglomus turgidum</i> (strain DSM 6724 / Z-1310)	Dictyoglomus	bacteria	unicellular	1743
30	<i>Synechocystis</i> sp. (strain PCC 6803 / Kazusa)	Synechocystis	bacteria	unicellular	3508
31	<i>Chlamydomonas reinhardtii</i> (Chlamydomonas smithii)	Chlamydomonas	eukaryota	unicellular	18832
32	<i>Physcomitrium patens</i> (Spreading-leaved earth moss) (Physcomitrella patens)	Physcomitrella	eukaryota	multicellular	47782
33	<i>Arabidopsis thaliana</i> (Mouse-ear cress)	Arabidopsis	eukaryota	multicellular	41596
34	<i>Zea mays</i> (Maize)	Maize	eukaryota	multicellular	63281
35	<i>Leishmania major</i>	Leishmania	eukaryota	unicellular	8038
36	<i>Paramecium tetraurelia</i>	Paramecium	eukaryota	unicellular	39461
37	<i>Caenorhabditis elegans</i>	C. elegans	eukaryota	multicellular	28553
38	<i>Helobdella robusta</i> (California leech)	Leech	eukaryota	multicellular	23328
39	<i>Ixodes scapularis</i> (Black-legged tick) (Deer tick)	Tick	eukaryota	multicellular	20496
40	<i>Tribolium castaneum</i> (Red flour beetle)	Tribolium	eukaryota	multicellular	18505
41	<i>Anopheles gambiae</i> (African malaria mosquito)	Mosquito	eukaryota	multicellular	14411
42	<i>Drosophila melanogaster</i> (Fruit fly)	Fruit fly	eukaryota	multicellular	23539
43	<i>Ciona intestinalis</i> (Transparent sea squirt) (Ascidia intestinalis)	Sea squirt	eukaryota	multicellular	17311
44	<i>Branchiostoma floridae</i> (Florida lancelet) (Amphioxus)	Amphioxus	eukaryota	multicellular	38648
45	<i>Lepisosteus osseus</i> (Spotted gar)	Spotted gar	eukaryota	multicellular	22463
46	<i>Danio rerio</i> (Zebrafish) (Brachydanio rerio)	Zebrafish	eukaryota	multicellular	46840
47	<i>Oryzias latipes</i> (Japanese rice fish) (Japanese killifish)	Medaka	eukaryota	multicellular	36138
48	<i>Xenopus laevis</i> (African clawed frog)	X. laevis	eukaryota	multicellular	61769
49	<i>Xenopus tropicalis</i> (Western clawed frog) (Silurana tropicalis)	X. tropicalis	eukaryota	multicellular	37693
50	<i>Gallus gallus</i> (Chicken)	Chicken	eukaryota	multicellular	43968
51	<i>Macaca mulatta</i> (Rhesus macaque)	Rhesus	eukaryota	multicellular	44416
52	<i>Gorilla gorilla gorilla</i> (Western lowland gorilla)	Gorilla	eukaryota	multicellular	44726
53	<i>Pan troglodytes</i> (Chimpanzee)	Chimpanzee	eukaryota	multicellular	48794
54	<i>Homo sapiens</i> (Human)	Human	eukaryota	multicellular	104573
55	<i>Canis lupus familiaris</i> (Dog) (Canis familiaris)	Dog	eukaryota	multicellular	43672
56	<i>Bos taurus</i> (Bovine)	Cow	eukaryota	multicellular	37871
57	<i>Mus musculus</i> (Mouse)	Mouse	eukaryota	multicellular	63289
58	<i>Rattus norvegicus</i> (Rat)	Rat	eukaryota	multicellular	49582
59	<i>Monodelphis domestica</i> (Gray short-tailed opossum)	Opossum	eukaryota	multicellular	36221
60	<i>Thalassiosira pseudonana</i> (Marine diatom) (Cyclotella nana)	Thalassiosira	eukaryota	unicellular	11612
61	<i>Daphnia magna</i>	Daphnia	eukaryota	multicellular	26600
62	<i>Plasmidium falciiparum</i> (isolate 3D7)	Plasmodium	eukaryota	unicellular	5369
63	<i>Oryza sativa</i> subsp. <i>japonica</i> (Rice)	Rice	eukaryota	multicellular	49224
64	<i>Dictyostelium discoideum</i> (Social amoeba)	Dictyostelium	eukaryota	unicellular	12746
65	<i>Nematostella vectensis</i> (Starlet sea anemone)	Sea anemone	eukaryota	multicellular	24445
66	<i>Monosiga brevicollis</i> (Choanoflagellate)	Monosiga	eukaryota	unicellular	9156
67	<i>Phytophthora ramorum</i> (Sudden oak death agent)	Phytophthora	eukaryota	unicellular	15349
68	<i>Giardia intestinalis</i> (strain ATCC 50803 / WB clone C6) (Giardia lamblia)	Giardia	eukaryota	unicellular	4900
69	<i>Cryptococcus neoformans</i> var. <i>neoformans</i> serotype D (strain JEC21 / ATCC MYA-565) (Filobasidiella neoformans)	Cryptococcus	eukaryota	unicellular	6746
70	<i>Candida albicans</i> (strain SC5314 / ATCC MYA-2876) (Yeast)	Candida	eukaryota	unicellular	6037
71	<i>Ustilago maydis</i> (strain S21 / FGSC 9021) (Corn smut fungus)	Corn smut	eukaryota	unicellular	6805
72	<i>Yarrowia lipolytica</i> (strain CLIB 122 / E 150) (Yeast) (Candida lipolytica)	Yarrowia	eukaryota	unicellular	6454
73	<i>Schizosaccharomyces pombe</i> (strain 972 / ATCC 24843) (Fission yeast)	S. pombe	eukaryota	unicellular	5132
74	<i>Phaeosphaeria nodorum</i> (strain SN15 / ATCC MYA-4574 / FGSC 10173) (Glume blotch fungus) (Parastagonospora nodorum)	Phaeosphaeria	eukaryota	unicellular	15998
75	<i>Aspergillus fumigatus</i> (strain ATCC MYA-4609 / CBS 101355 / FGSC A1100 / Af293) (Neosartorya fumigata)	Aspergillus	eukaryota	unicellular	9648
76	<i>Neurospora crassa</i> (strain ATCC 24698 / 74-OR23-1A / CBS 708.71 / DSM 1257 / FGSC 987)	Neurospora	eukaryota	unicellular	10266
77	<i>Trichomonas vaginalis</i> (strain ATCC PRA-98 / G3)	Trichomonas	eukaryota	unicellular	50190
78	<i>Puccinia graminis</i> f. sp. <i>tritici</i> (strain CRL 75-36-700-3 / race SCCL) (Black stem rust fungus)	Stem rust	eukaryota	unicellular	15808
79	<i>Saccharomyces cerevisiae</i> (strain ATCC 204508 / S288c) (Baker's yeast)	S. cerevisiae	eukaryota	unicellular	6091
80	<i>Sclerotinia sclerotiorum</i> (strain ATCC 18683 / 1980 / Ss-1) (White mold) (Whetzelinia sclerotiorum)	Sclerotinia	eukaryota	unicellular	14445
81	<i>Batrachochytrium dendrobatidis</i> (strain JAM81 / FGSC 10211) (Frog chytrid fungus)	Batrachochytrium	eukaryota	unicellular	8610
		Total			154730

Table 1. Species Information, Short Labels, and Protein Counts

This table lists the scientific names of the 81 species analyzed in this study, along with their corresponding short labels (used for phylogenetic tree visualization), domain classification (archaea, bacteria, eukaryotes), cell organization (unicellular or multicellular), and the total number of annotated proteins (“Protein Count”) obtained from publicly available datasets. The color codes in the “Domain” column correspond to blue for archaea, red for bacteria, and green for eukaryotes.

No.	Short Label	Ala	Cys	Asp	Glu	Phe	Gly	His	Ile	Lys	Leu	Met	Asn	Pro	Gln	Arg	Ser	Thr	Val	Trp	Tyr
1	Halobacterium	0.12407	0.00889	0.0888	0.06979	0.03058	0.08201	0.02265	0.03687	0.01821	0.08469	0.0191	0.02153	0.0464	0.02691	0.06561	0.05464	0.06794	0.0949	0.01126	0.02514
2	Thermococcus	0.07304	0.00605	0.04605	0.09242	0.04258	0.07421	0.01543	0.07038	0.0718	0.10734	0.02497	0.02987	0.04216	0.01729	0.06113	0.04855	0.04367	0.08395	0.01264	0.03645
3	Methanosaerina	0.06629	0.01453	0.04976	0.07941	0.04731	0.06994	0.01726	0.07562	0.07015	0.05989	0.02711	0.04246	0.03945	0.02558	0.04721	0.06736	0.05082	0.068	0.01016	0.03571
4	Methanococcoides	0.05447	0.01407	0.05329	0.08678	0.04244	0.06184	0.01392	0.071	0.0671	0.09512	0.02504	0.05102	0.0327	0.01459	0.03919	0.04432	0.03922	0.06867	0.00701	0.04246
5	Saccharolobus	0.05343	0.00756	0.04685	0.07053	0.0438	0.06169	0.01298	0.09485	0.08221	0.11399	0.02327	0.04882	0.03588	0.02106	0.04918	0.06635	0.04584	0.07477	0.00995	0.04701
6	Korarchaeum	0.06846	0.00938	0.04895	0.08472	0.03745	0.07669	0.01348	0.08305	0.05775	0.10988	0.02782	0.02675	0.04342	0.01481	0.07305	0.07177	0.03527	0.07344	0.01108	0.03318
7	Nitrosopumilus	0.05963	0.01131	0.05959	0.07541	0.04374	0.06223	0.01819	0.08494	0.09008	0.08403	0.02842	0.04776	0.03691	0.03234	0.03552	0.06945	0.05444	0.06666	0.00989	0.03041
8	M. tuberculosis	0.13255	0.01007	0.05865	0.04819	0.02838	0.09087	0.02317	0.04236	0.0218	0.09666	0.02099	0.02233	0.05782	0.03129	0.07809	0.05515	0.05916	0.08598	0.01504	0.02057
9	E. coli	0.09305	0.01318	0.05002	0.0576	0.03944	0.07025	0.02335	0.0617	0.04742	0.10621	0.03054	0.03904	0.04268	0.04376	0.05607	0.05831	0.0534	0.07084	0.01514	0.028
10	Helicobacter	0.06867	0.01225	0.04646	0.06903	0.05481	0.05555	0.02189	0.07188	0.09247	0.11448	0.02528	0.05466	0.03216	0.03653	0.0363	0.06709	0.0417	0.05719	0.00745	0.03593
11	Streptomyces	0.13599	0.00888	0.06061	0.05768	0.02622	0.04979	0.0239	0.02903	0.02075	0.10077	0.0185	0.01645	0.06188	0.02648	0.08583	0.05005	0.06105	0.08621	0.01524	0.01972
12	Meningococcus	0.06555	0.01234	0.05037	0.06095	0.04411	0.07361	0.02219	0.05997	0.05986	0.0982	0.02674	0.04048	0.04249	0.0315	0.05665	0.05563	0.05087	0.06652	0.01217	0.03026
13	Leptospira	0.05031	0.01006	0.04632	0.07204	0.0582	0.06183	0.01719	0.08184	0.08268	0.10395	0.02057	0.05206	0.03772	0.03269	0.04483	0.07716	0.0494	0.05664	0.01145	0.03517
14	Fusobacterium	0.05234	0.00858	0.0527	0.0799	0.05173	0.05921	0.01188	0.01035	0.01049	0.09525	0.02538	0.06217	0.02491	0.02151	0.03207	0.05801	0.04602	0.06088	0.00634	0.04486
15	Pseudomonas	0.1635	0.0114	0.05248	0.06156	0.03557	0.08272	0.02222	0.04153	0.02988	0.12462	0.02175	0.0253	0.05084	0.04232	0.07704	0.05463	0.04113	0.06885	0.01504	0.02486
16	B. subtilis	0.07373	0.00897	0.05087	0.07343	0.0458	0.06647	0.02288	0.07454	0.07458	0.09641	0.03031	0.04034	0.03454	0.03849	0.04144	0.06179	0.05295	0.0672	0.01024	0.03541
17	Aquifex	0.05813	0.00893	0.04224	0.09674	0.05124	0.06659	0.01555	0.07351	0.09662	0.10561	0.0205	0.0356	0.04013	0.02066	0.0495	0.04767	0.04139	0.07947	0.00937	0.04055
18	Bradyrhizobium	0.12569	0.01047	0.05365	0.05283	0.03696	0.08185	0.02124	0.05213	0.0374	0.0763	0.02633	0.04524	0.05292	0.0316	0.07496	0.05847	0.05299	0.07343	0.01345	0.02154
19	Bacteroides	0.06771	0.0138	0.05323	0.06664	0.04671	0.06457	0.0188	0.06206	0.06924	0.0925	0.0263	0.04971	0.03669	0.03432	0.04655	0.06205	0.05549	0.06357	0.01271	0.04485
20	Rhodopirellula	0.08859	0.01879	0.05437	0.05502	0.03876	0.07142	0.02544	0.04779	0.03579	0.09128	0.02706	0.03036	0.05542	0.03938	0.07983	0.0791	0.05724	0.06811	0.01593	0.01851
21	Deinococcus	0.12044	0.01765	0.0497	0.06138	0.09042	0.02157	0.03185	0.02746	0.11759	0.02005	0.02303	0.06085	0.0418	0.07576	0.05194	0.05744	0.07643	0.01467	0.02237	
22	Geobacter	0.10093	0.01443	0.05256	0.06639	0.03912	0.08291	0.02089	0.05658	0.04475	0.10136	0.02614	0.02788	0.0467	0.0273	0.07124	0.05314	0.05342	0.0785	0.0102	0.02556
23	Mycoplasma	0.05619	0.00988	0.04591	0.05393	0.06123	0.04645	0.08423	0.0179	0.07169	0.02089	0.05174	0.04429	0.0457	0.03399	0.05281	0.05272	0.06213	0.00941	0.03183	
24	Thermotoga	0.05742	0.00847	0.04814	0.08891	0.05199	0.06882	0.01575	0.07238	0.07894	0.11056	0.02553	0.03494	0.03891	0.01979	0.05662	0.05611	0.04436	0.0866	0.01079	0.03467
25	Gloeobacter	0.13994	0.01139	0.04863	0.06323	0.03664	0.08016	0.01917	0.04493	0.0308	0.11385	0.0197	0.02606	0.05579	0.04251	0.0758	0.0543	0.05052	0.07519	0.01542	0.02599
26	Chlamydia	0.07537	0.01766	0.04465	0.06744	0.05057	0.04805	0.06224	0.02665	0.06105	0.11303	0.02206	0.03484	0.04241	0.0167	0.05089	0.0795	0.04938	0.06615	0.00928	0.0294
27	Thermodesulfobivirid	0.05953	0.01214	0.04616	0.08053	0.04949	0.06168	0.01614	0.0201	0.09635	0.1066	0.02422	0.04223	0.0373	0.02707	0.04214	0.05656	0.04469	0.06233	0.00845	0.03577
28	Chloroflexus	0.11665	0.00913	0.04999	0.05489	0.03035	0.07508	0.02203	0.06084	0.01942	0.11375	0.0219	0.0257	0.05775	0.04189	0.07782	0.04814	0.05758	0.07602	0.01582	0.02663
29	Dictyoglycomus	0.05175	0.00781	0.04687	0.08251	0.05036	0.06613	0.01364	0.08922	0.08789	0.10536	0.02077	0.04564	0.03887	0.0212	0.04543	0.05539	0.04123	0.06691	0.01081	0.04221
30	Synechocystis	0.08285	0.01119	0.0493	0.06173	0.04029	0.07021	0.01908	0.0628	0.04545	0.1143	0.02225	0.03915	0.05043	0.04566	0.05283	0.0563	0.06169	0.0106	0.02931	
31	Chlamydomonas	0.15477	0.01565	0.04349	0.05059	0.02038	0.0551	0.0216	0.02349	0.02327	0.0880	0.02077	0.02148	0.07044	0.0449	0.065	0.0529	0.04813	0.0649	0.01302	0.01929
32	Phycisomittella	0.07484	0.02155	0.04843	0.06009	0.04002	0.06676	0.02582	0.04853	0.05496	0.09671	0.02654	0.03897	0.0497	0.03905	0.06068	0.08663	0.05225	0.06816	0.01383	0.02647
33	Arabidopsis	0.06273	0.01994	0.05249	0.0656	0.04369	0.0642	0.02281	0.05348	0.06493	0.09354	0.02598	0.0433	0.04843	0.0345	0.05462	0.09068	0.05136	0.06676	0.01238	0.02861
34	Maize	0.07537	0.01959	0.05225	0.05851	0.03506	0.07529	0.02417	0.04821	0.0923	0.02409	0.0375	0.0599	0.03554	0.04682	0.08228	0.04772	0.06834	0.01287	0.02485	
35	Leishmania	0.12123	0.01988	0.04876	0.0607	0.0325	0.06396	0.02691	0.03322	0.03921	0.0921	0.02512	0.02826	0.05251	0.04032	0.04032	0.07125	0.08224	0.05839	0.07266	0.01129
36	Paramaecium	0.03532	0.01676	0.04883	0.03576	0.01871	0.08293	0.09275	0.05977	0.02274	0.06816	0.03047	0.0417	0.0512	0.03745	0.0669	0.04875	0.04489	0.00769	0.04145	
37	C. elegans	0.064	0.0217	0.05195	0.06312	0.04829	0.05487	0.02308	0.06134	0.06279	0.02846	0.04835	0.04356	0.0496	0.0417	0.05239	0.08043	0.05844	0.06179	0.01113	0.03232
38	Leech	0.05147	0.02567	0.05646	0.05683	0.04471	0.04953	0.02566	0.06147	0.07335	0.08879	0.02498	0.07006	0.03854	0.04047	0.04827	0.08343	0.05604	0.06032	0.01101	0.03304
39	Tick	0.08020	0.02489	0.04887	0.05652	0.04033	0.06947	0.02688	0.03858	0.04984	0.05981	0.02322	0.03232	0.03288	0.03958	0.03782	0.07059	0.05572	0.07014	0.01288	0.02766
40	Tribolium	0.06016	0.0219	0.05065	0.0638	0.04046	0.05509	0.02475	0.05828	0.05699	0.09195	0.02324	0.0307	0.05071	0.03975	0.0555	0.07415	0.05709	0.06388	0.01131	0.0331
41	Mosquito	0.07664	0.02044	0.05163	0.06165	0.03768	0.06549	0.02553	0.04528	0.05058	0.08998	0.02429	0.04352	0.05282	0.04528	0.0567	0.07537	0.05846	0.06528	0.01046	0.03311
42	Fruit fly	0.07446	0.01998	0.05093	0.0634	0.02629	0.05018	0.03259	0.04417	0.0561	0.0954	0.02259	0.04768	0.0523	0.04249	0.05					

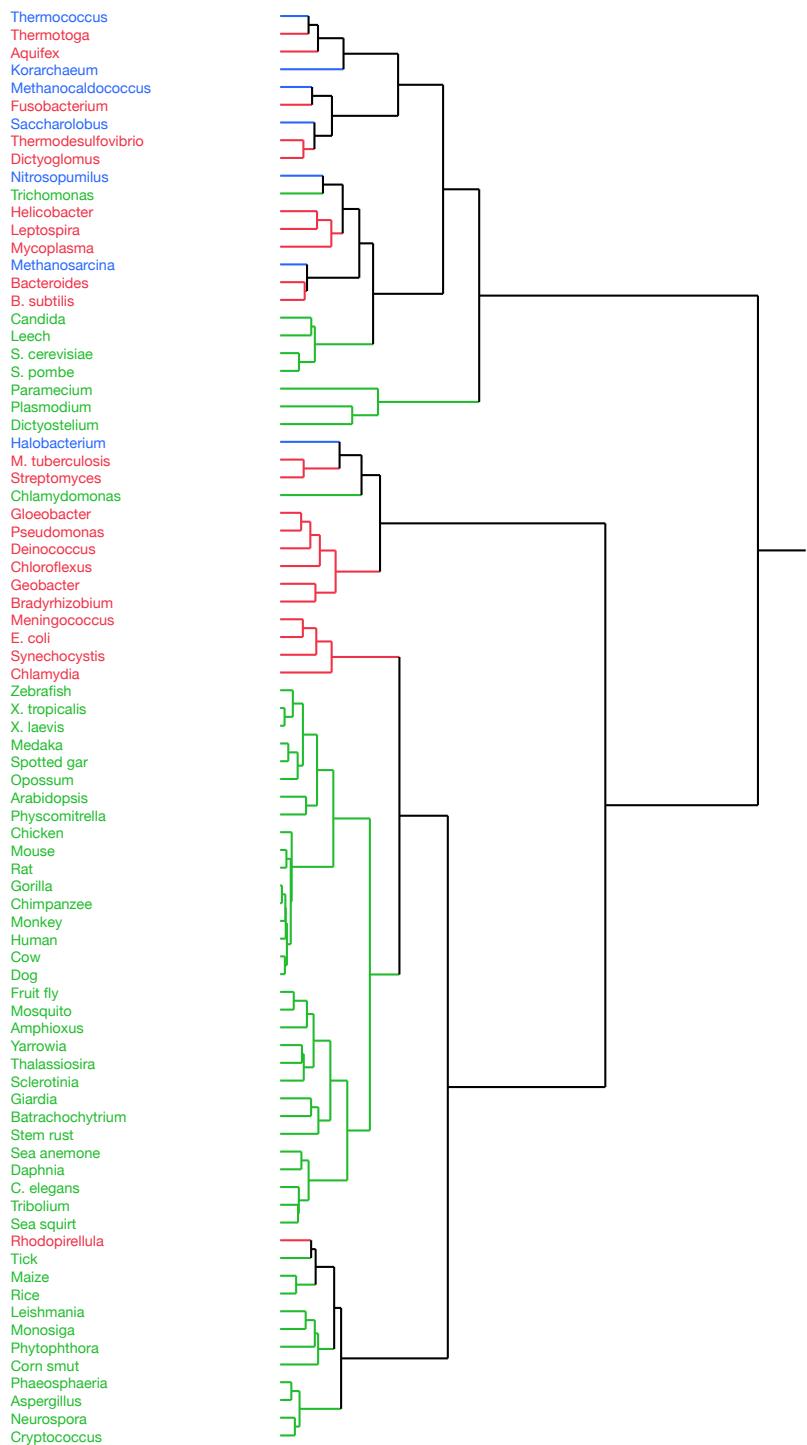


Figure 1. Horizontal Phylogenetic Tree Based on Average Exome-Wide Amino Acid Composition

This dendrogram was generated using hierarchical clustering with Ward's method, applied to the pairwise angular distance matrix calculated from the average amino acid compositions of 81 species' exomes. Taxonomic groups are color-coded: blue for archaea, red for bacteria, and green for eukaryotes. Short labels for each species are displayed on the left; see Table 1 for full scientific names. The branch lengths reflect the relative distances in amino acid composition between species.

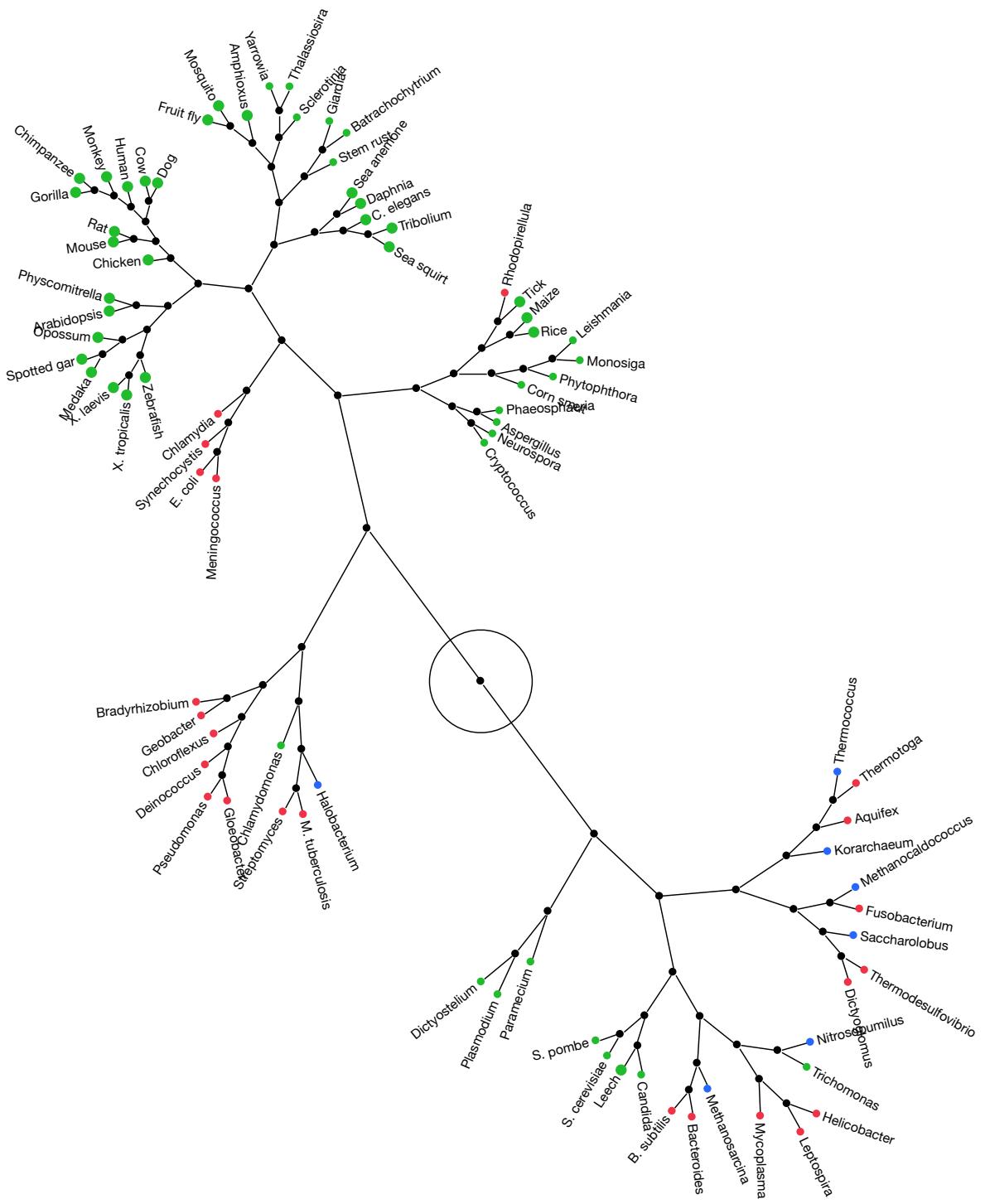


Figure 2. Radial Phylogenetic Tree Based on Average Exome-Wide Amino Acid Composition

This radial (divergent) dendrogram was generated using the same dataset and clustering approach described in Figure 1, but displayed in a circular layout to illustrate the evolutionary divergence among the 81 species. As in Figure 1, archaea are shown in blue, bacteria in red, and eukaryotes in green. Unicellular organisms are indicated with small nodes, whereas multicellular organisms are represented by larger nodes and bold labels. Short labels correspond to those in Table 1.