

Table 2. Median protein lengths in eukaryotic, bacterial and archaeal organisms

Species <sup>a</sup>	All species		Classified in COG		Classified in Pfam-A	
	Number <sup>b</sup>	Median <sup>c</sup>	Number <sup>b</sup>	Median <sup>c</sup>	Number <sup>b</sup>	Median <sup>c</sup>
<b>Eukarya</b>	<b>104 394</b>	<b>361</b>	<b>5177</b>	<b>471</b>	<b>71 584</b>	<b>419</b>
HUMAN	33 869	375	–	–	21 686	416
DROME	14 226	373	3092	492	13 091	475
CAEEL	21 124	344	–	–	13 316	391
YEAST	6315	379	2085	438	3953	448
ARATH	28 860	356	–	–	19 538	407
<b>Bacteria</b>	<b>191 541</b>	<b>267</b>	<b>83 513</b>	<b>304</b>	<b>131 915</b>	<b>306</b>
ECOLI	4289	<u>278</u>	3289	<u>309</u>	3483	<u>303</u>
SALTY	4553	<u>271</u>	3408	<u>303</u>	3527	<u>300</u>
SALTI	4767	<u>253</u>	3258	<u>300</u>	3118	<u>300</u>
YERPE	4083	<u>268</u>	2991	<u>299</u>	3003	<u>304</u>
SHIFL	4180	<u>261</u>	–	–	2613	<u>304</u>
WIGBR	654	<u>268</u>	–	–	571	<u>291</u>
BUCAI	574	<u>282</u>	558	<u>284</u>	544	<u>285</u>
BUCAP	545	<u>279</u>	–	–	536	<u>285</u>
HAEin	1709	<u>262</u>	1470	<u>286</u>	750	<u>314</u>
PASMU	2014	<u>286</u>	1740	<u>302</u>	780	<u>289</u>
XANCP	4181	<u>291</u>	–	–	2976	<u>325</u>
XANAC	4312	<u>286</u>	–	–	3056	<u>326</u>
XYLFA	2832	<u>201</u>	1549	<u>305</u>	1544	<u>305</u>
PSEAE	5565	<u>291</u>	4355	<u>310</u>	4309	<u>309</u>
VIBCH	3828	<u>259</u>	2794	<u>315</u>	2731	<u>312</u>
Chromosome 1	2736	<u>273</u>	2133	<u>314</u>	–	–
Chromosome 2	1092	225	661	316	–	–
SHEON	4778	245	–	–	2913	<u>308</u>
NEIMB	2025	239	1448	<u>291</u>	1310	<u>305</u>
RALSO	5116	<u>276</u>	–	–	3518	<u>310</u>
Chromosome 1	3440	<u>271</u>	–	–	–	–
Chromosome 2	1676	296	–	–	–	–
AGRT5	5402	<u>280</u>	3984	<u>316</u>	4062	<u>307</u>
Circular chr.	2785	<u>258</u>	2098	<u>305</u>	–	–
Linear chr.	1876	302	1424	<u>329</u>	–	–
Plasmids	741	273	462	316	–	–
RHILO	7275	<u>269</u>	5184	<u>305</u>	5107	<u>303</u>
Chromosome	6746	<u>270</u>	4888	<u>304</u>	–	–
Plasmids	529	243	296	<u>327</u>	–	–
RHIME	6205	<u>281</u>	4614	<u>312</u>	4669	<u>308</u>
Chromosome	3341	<u>276</u>	2602	<u>302</u>	–	–
Plasmid A	1294	265	890	310	–	–
Plasmid B	1570	303	1122	330	–	–
BRUME	3198	<u>263</u>	–	–	2322	<u>300</u>
Chromosome 1	2059	<u>252</u>	–	–	–	–
Chromosome 2	1139	279	–	–	–	–
BRUSU	3264	<u>254</u>	–	–	2351	<u>301</u>
Chromosome 1	2116	<u>239</u>	–	–	–	–
Chromosome 2	1148	278	–	–	–	–
CAUCR	3737	<u>275</u>	2551	<u>317</u>	2686	<u>312</u>
RICPR	834	<u>283</u>	687	<u>295</u>	672	<u>299</u>
RICCN	1374	173	861	<u>247</u>	769	<u>264</u>
HELPHY	1566	<u>266</u>	1083	<u>303</u>	1052	<u>315</u>
CAMJE	1634	<u>268</u>	1309	<u>294</u>	1197	<u>298</u>
MYCTU	3918	<u>287</u>	2554	<u>322</u>	2213	<u>326</u>
MYCLE	1605	<u>282</u>	1145	<u>326</u>	1138	<u>324</u>
STRCO	7897	<u>278</u>	–	–	5330	<u>317</u>
CORGL	2993	<u>275</u>	1954	<u>307</u>	1985	<u>314</u>
COREF	2950	<u>287</u>	–	–	1961	<u>323</u>
BIFLO	1729	<u>321</u>	–	–	1286	<u>341</u>
OCEIH	3496	<u>261</u>	–	–	2583	<u>295</u>
BACSU	4100	<u>256</u>	2818	<u>298</u>	2974	<u>297</u>
BACHD	4066	<u>261</u>	2838	<u>303</u>	2916	<u>300</u>
STAAN	2625	<u>257</u>	1801	<u>300</u>	1542	<u>293</u>
LISIN	3043	<u>255</u>	2176	<u>289</u>	2234	<u>291</u>
LISMO	2846	<u>267</u>	2206	<u>289</u>	2211	<u>292</u>
LACLA	2266	<u>251</u>	1602	<u>288</u>	1690	<u>281</u>
STRA5	2124	<u>254</u>	–	–	1480	<u>290</u>
STRMU	1960	<u>250</u>	–	–	1445	<u>282</u>
STRPN	2043	243	1465	<u>287</u>	1409	<u>291</u>
STRPY	1696	<u>263</u>	1178	<u>294</u>	1159	<u>299</u>
CLOAB	3848	<u>262</u>	2487	<u>298</u>	2634	<u>299</u>

Table 2. Continued

Species <sup>a</sup>	All species		Classified in COG		Classified in Pfam-A	
	Number <sup>b</sup>	Median <sup>c</sup>	Number <sup>b</sup>	Median <sup>c</sup>	Number <sup>b</sup>	Median <sup>c</sup>
CLOPE	2723	<u>268</u>	–	–	1997	<u>303</u>
THETN	2588	<u>269</u>	–	–	1903	<u>306</u>
UREPA	614	<u>286</u>	409	<u>298</u>	395	<u>303</u>
MYCGE	484	<u>292</u>	384	<u>292</u>	375	<u>304</u>
MYCPN	677	<u>286</u>	407	<u>299</u>	507	<u>299</u>
MYCPU	782	<u>297</u>	489	<u>302</u>	498	<u>320</u>
MYCPE	1037	<u>304</u>	–	–	664	<u>315</u>
FUSNN	2067	<u>261</u>	–	–	1432	<u>303</u>
SYNY3	3169	<u>274</u>	2141	<u>318</u>	2344	<u>306</u>
ANASP	6129	<u>256</u>	–	–	3600	<u>320</u>
SYNEL	2475	<u>272</u>	–	–	1759	<u>315</u>
CHLTR	894	<u>289</u>	615	<u>316</u>	639	<u>327</u>
CHLMU	916	<u>290</u>	644	<u>327</u>	641	<u>320</u>
CHLPN	1052	<u>289</u>	646	<u>324</u>	716	<u>333</u>
BORBU	1637	<u>220</u>	635	<u>318</u>	981	<u>265</u>
Chromosome	850	286	–	–	–	–
Plasmids	787	179	–	–	–	–
TREPA	1031	<u>293</u>	708	<u>331</u>	691	<u>337</u>
LEPIN	4727	<u>207</u>	–	–	2243	<u>309</u>
Chromosome 1	4360	206	–	–	–	–
Chromosome 2	367	223	–	–	–	–
DEIRA	3182	<u>264</u>	2249	<u>303</u>	2050	<u>307</u>
Chromosome 1	2629	<u>257</u>	1873	<u>294</u>	–	–
Chromosome 2	368	304	265	347	–	–
Plasmids	185	303	111	336	–	–
CHLTE	2252	239	–	–	1431	<u>311</u>
THEMA	1846	<u>284</u>	1509	<u>303</u>	1459	<u>304</u>
AQUAE	1560	<u>272</u>	1321	<u>297</u>	1231	<u>297</u>
Archaea	<b>37 141</b>	<b><u>247</u></b>	<b>18 219</b>	<b><u>283</u></b>	<b>24 067</b>	<b><u>288</u></b>
PYRAB	1765	265	1450	281	1407	282
PYRHO	1801	257	1398	283	1312	285
PYRFU	2065	253	1627	273	1477	281
ARCFU	2420	243	1887	270	1720	276
THEAC	1482	269	1233	293	1083	301
THEVO	1499	259	1247	287	1074	304
METTH	1869	242	1382	273	1325	277
METJA	1770	241	1298	266	1260	272
METAC	4540	256	–	–	2677	306
METMA	3371	255	–	–	2141	294
METKA	1691	257	–	–	1067	272
HALN1	2622	242	1746	297	1471	303
AERPE	1840	239	1191	293	1067	301
PYRAE	2603	208	–	–	1411	267
SULSO	2977	251	1983	294	1917	293
SULTO	2826	226	1777	284	1658	279

<sup>a</sup>See Table 1 for abbreviations.

<sup>b</sup>Number of proteins in each set.

<sup>c</sup>Median length.

Table 1. Proteomic collections

	Species	Abbreviation
Eukaryota	<i>Homo sapiens</i> <sup>a</sup>	HUMAN
	<i>Drosophila melanogaster</i>	DROME
	<i>Caenorhabditis elegans</i> <sup>a</sup>	CAEEL
	<i>Saccharomyces cerevisiae</i>	YEAST
Euryarchaeota	<i>Arabidopsis thaliana</i> <sup>a</sup>	ARATH
	<i>Pyrococcus abyssi</i>	PYRAB
	<i>Pyrococcus horikoshii</i>	PYRHO
	<i>Pyrococcus furiosus</i>	PYRFU
	<i>Archaeoglobus fulgidus</i>	ARCFU
	<i>Thermoplasma acidophilum</i>	THEAC
	<i>Thermoplasma volcanium</i>	THEVO
	<i>Methanothermobacter thermoautotrophicus</i>	METTH
	<i>Methanococcus jannaschii</i>	METJA
	<i>Methanosarcina acetivorans</i> <sup>a</sup>	METAC
	<i>Methanosarcina mazei</i> <sup>a</sup>	METMA
	<i>Methanopyrus kandleri</i> <sup>a</sup>	METKA
	<i>Halobacterium</i> sp. NRC-1	HALNI
Crenarchaeota	<i>Aeropyrum pernix</i>	AERPE
	<i>Pyrobaculum aerophilum</i> <sup>a</sup>	PYRAE
	<i>Sulfolobus solfataricus</i>	SULSO
	<i>Sulfolobus tokodaii</i>	SULTO
γ-Proteobacteria	<i>Escherichia coli</i> K12	ECOLI
	<i>Salmonella typhimurium</i>	SALTY
	<i>Salmonella enterica</i>	SALTI
	<i>Yersinia pestis</i> CO92	YERPE
	<i>Shigella flexneri</i> <sup>a</sup>	SHIFL
	<i>Wigglesworthia brevipalpis</i> <sup>a</sup>	WIGBR
	<i>Buchnera</i> sp. APS	BUCAI
	<i>Buchnera aphidicola</i> <sup>a</sup>	BUCAP
	<i>Haemophilus influenzae</i>	HAEIN
	<i>Pasteurella multocida</i>	PASMU
	<i>Xanthomonas campestris</i> <sup>a</sup>	XANCP
	<i>Xanthomonas axonopodis</i> <sup>a</sup>	XANAC
	<i>Xylella fastidiosa</i> <sup>a</sup>	XYLFA
	<i>Pseudomonas aeruginosa</i>	PSEAE
	<i>Vibrio cholerae</i>	VIBCH
	<i>Shewanella oneidensis</i> <sup>a</sup>	SHEON
β-Proteobacteria	<i>Neisseria meningitidis</i> MC58	NEIMB
	<i>Ralstonia solanacearum</i> <sup>a</sup>	RALSO
α-Proteobacteria	<i>Agrobacterium tumefaciens</i> C58 Cereon	AGRT5
	<i>Sinorhizobium loti</i>	RHILO
	<i>Sinorhizobium meliloti</i>	RHIME
	<i>Brucella melitensis</i> <sup>a</sup>	BRUME
	<i>Brucella suis</i> <sup>a</sup>	BRUSU
	<i>Caulobacter crescentus</i>	CAUCR
	<i>Rickettsia prowazekii</i>	RICPR
	<i>Rickettsia conorii</i>	RICCN
ε-Proteobacteria	<i>Helicobacter pylori</i> 26 695	HELPE
	<i>Campylobacter jejuni</i>	CAMJE
Actinobacteria	<i>Mycobacterium tuberculosis</i> H37 Rv	MYCTU
	<i>Mycobacterium leprae</i>	MYCLE
	<i>Streptomyces coelicolor</i> <sup>a</sup>	STRCO
	<i>Corynebacterium glutamicum</i>	CORGL
	<i>Corynebacterium efficiens</i> <sup>a</sup>	COREF
	<i>Bifidobacterium longum</i> <sup>a</sup>	BIFLO
Firmicutes	<i>Oceanobacillus ihyeensis</i> <sup>a</sup>	OCEIH
	<i>Bacillus subtilis</i>	BACSU
	<i>Bacillus halodurans</i>	BACHD
	<i>Staphylococcus aureus</i> N315	STAAN
	<i>Listeria innocua</i>	LISIN
	<i>Listeria monocytogenes</i>	LISMO
	<i>Lactococcus lactis</i>	LACLA
	<i>Streptococcus agalactiae</i> 2603 V/R <sup>a</sup>	STRA5
	<i>Streptococcus mutans</i> <sup>a</sup>	STRMU
	<i>Streptococcus pneumoniae</i> TIGR4	STRPN
	<i>Streptococcus pyogenes</i> M1	STRPY
	<i>Clostridium acetobutylicum</i>	CLOAB
	<i>Clostridium perfringens</i> <sup>a</sup>	CLOPE
	<i>Thermoanaerobacter tengcongensis</i> <sup>a</sup>	THETN
	<i>Ureaplasma urealyticum</i>	UREPA

Table 1. Continued

	Species	Abbreviation
	<i>Mycoplasma genitalium</i>	MYCGE
	<i>Mycoplasma pneumoniae</i>	MYCPN
	<i>Mycoplasma pulmonis</i>	MYCPU
	<i>Mycoplasma penetrans</i> <sup>a</sup>	MYCPE
	<i>Fusobacterium nucleatum</i> <sup>a</sup>	FUSNN
Cyanobacteria	<i>Synechocystis</i> sp. PCC 6803	SYNY3
	<i>Nostoc</i> sp. PCC 7120 <sup>b</sup>	ANASP
	<i>Thermosynechococcus elongatus</i> <sup>a</sup>	SYNEL
Chlamydiae	<i>Chlamydia trachomatis</i>	CHLTR
	<i>Chlamydia muridarum</i>	CHLMU
	<i>Chlamydomonas pneumoniae</i> CWL029	CHLPN
Spirochaetes	<i>Borrelia burgdorferi</i>	BORBU
	<i>Treponema pallidum</i>	TREPA
	<i>Leptospira interrogans</i> <sup>a</sup>	LEPIN
Others	<i>Deinococcus radiodurans</i>	DEIRA
	<i>Chlorobium tepidum</i> <sup>a</sup>	CHLTE
	<i>Thermotoga maritima</i>	THEMA
	<i>Aquifex aeolicus</i>	AQUAE

<sup>a</sup>Proteins from these species are not classified in the COG database and are excluded from the functional group analyses.

<sup>b</sup>The COG classification of proteins from this species does not follow the standard coding and has been excluded from the COG analyses.

lengths were compared with respect to taxonomic, functional and ecological classes.

#### Datasets of selected proteins

We evaluated results using the set of proteins classified in the COG (Clusters of Orthologous Groups of proteins) database (13,14) and the set of genomic proteins included in the Pfam (15–17) database of functional/structural domain alignments verified by human intervention (Pfam-A). The COG database classifies orthologous proteins in functional groups. At the onset of this study, classification of proteins into COGs was available for 2 eukaryotic proteomes (yeast and *Drosophila melanogaster*), 12 archaeal proteomes and 44 bacterial proteomes (Tables 1 and 2). COG data were obtained for prokaryotic organisms and yeast from the corresponding tables (\*.ptt) available from the NCBI genomes ftp site (<ftp://ftp.ncbi.nih.gov/genomes>). Data for *D.melanogaster* were obtained from the classification table available at the COG website (<http://www.ncbi.nlm.nih.gov/COG>). All proteomes included in our analysis are also represented in the Pfam-A database.

#### Statistical significance evaluations

We compared median protein lengths (the midpoint of all lengths arranged in order of magnitude) rather than average lengths between two sets of proteins to reduce the effect of outliers. The statistical significance of the difference in median length between the proteins of two sets was evaluated estimating the distribution of median length differences between samples created by randomly redistributing all sequences from the two sets into two new sets of the original sizes. For each determination, from 200 to 1000 independent data shufflings were implemented. We highlight the differences in median length observed in <1% of all shuffled samples ( $P < 0.01$ , boldfaced in the tables) and observed in the range 1–10% of all shuffled samples ( $0.01 \leq P < 0.1$ , underlined in the tables).