

RESEARCH ARTICLE

Bacterial Protein Interaction Networks: Connectivity is Ruled by Gene Conservation, Essentiality and Function

Maddalena Dilucca^{1,*}, Giulio Cimini² and Andrea Giansanti³

¹Dipartimento di Fisica, Sapienza University of Rome, 00185, Rome, Italy; ²Dipartimento di Fisica, Tor Vergata University of Rome, 00133, Rome, Italy Istituto dei Sistemi Complessi CNR UoS, Rome, Italy; ³Dipartimento di Fisica, Sapienza University of Rome, 00185, Rome, Italy INFN Roma1 Unit, Rome, Italy

Abstract: Background: Protein-protein interaction (PPI) networks are the backbone of all processes in living cells. In this work, we relate conservation, essentiality and functional repertoire of a gene to the connectivity k (i.e. the number of interactions, links) of the corresponding protein in the PPI network.

Methods: On a set of 42 bacterial genomes of different sizes, and with reasonably separated evolutionary trajectories, we investigate three issues: i) whether the distribution of connectivities changes between PPI subnetworks of essential and nonessential genes; ii) how gene conservation, measured both by the evolutionary retention index (ERI) and by evolutionary pressures, is related to the connectivity of the corresponding protein; iii) how PPI connectivities are modulated by evolutionary and functional relationships, as represented by the Clusters of Orthologous Genes (COGs).

Results: We show that conservation, essentiality and functional specialisation of genes constrain the connectivity of the corresponding proteins in bacterial PPI networks. In particular, we isolated a core of highly connected proteins (connectivities $k \geq 40$), which is ubiquitous among the species considered here, though mostly visible in the degree distributions of bacteria with small genomes (less than 1000 genes).

Conclusion: The genes that support this highly connected core are conserved, essential and, in most cases, belong to the COG cluster J, related to ribosomal functions and the processing of genetic information.

ARTICLE HISTORY

Received: April 14, 2020
Revised: August 13, 2020
Accepted: August 27, 2020

DOI:
10.2174/1389202922666210219110831

Keywords: Protein-protein interactions, gene essentiality, evolutionary retention index, clusters of orthologous genes, bacterial genomes, cellular processes.

1. INTRODUCTION

To operate biological activities in living cells, proteins work in association with other proteins, often assembled in large complexes. Hence, knowing the interactions of a protein is important to understand its cellular functions. Moreover, a comprehensive description of the stable and transient protein-protein interactions (PPIs) within a cell would facilitate the functional annotation of all gene products, and provide insight into the higher-order organization of the proteome [1, 2]. Several methodologies have been developed to detect PPIs, and have been adapted to chart interactions at the proteome-wide scale. These methods, combining different technologies, experiments and computational analyses, generate PPI networks of sufficient reliability, enabling the assignment of several proteins to functional categories [3, 4]. Moreover, the statistical study of bacterial PPIs over sev-

eral species (meta-interactomes) has brought important knowledge about protein functions and cellular processes [5, 6].

Our aim here is to shed some light on the relationships among conservation, essentiality and functional annotation at the genetic level and connectivities of PPI networks, at the protein level. We extend here our previous observations made on the PPI of *E. coli* which suggested a strong correlation between the connectivity of PPI networks on the one hand, and codon bias, gene conservation and essentiality on the other hand [7, 8]. In the next two paragraphs, it is worth making more precise what is usually meant by gene essentiality and gene conservation. Individual genes in the genome differently contribute to the survival of an organism. According to their known functional profiles and based on experimental evidence, genes can be divided into two categories: essential and nonessential ones [9, 10]. Essential genes are not dispensable for the survival of an organism in the environment it lives in [10, 11]. Nonessential genes are instead those which are dispensable [12], being related to

*Address correspondence to this author at the Dipartimento di Fisica, Sapienza University of Rome Piazzale Aldo Moro, 5, 00185, Rome, Italy; Tel: 3475407737; E-mail: maddalena.dilucca@gmail.com

functions that can be silenced without compromising the survival of the organism. Naturally, each species has adapted to one or more evolving environments and, plausibly, genes that are essential for one species may not be essential for another one.

It has been argued many times that essential genes are more conserved than nonessential ones [13-17]. The term “conservation” has, however, at least two meanings. On the one hand, a gene is conserved if orthologous copies of it are found in the genomes of many species, as measured by the Evolutionary Retention Index (ERI) [9, 18]. On the other hand, a gene is (evolutionarily) conserved when it is subject to a purifying, selective evolutionary pressure, which disfavors mutations. A common measure of evolutionary pressure is K_a/K_s , the ratio of the number of non-synonymous substitutions per non-synonymous site to the number of synonymous substitutions per synonymous site. In this second meaning, a conserved gene is, in a nutshell, a slowly evolving gene, a gene that hardly incorporates mutations [13, 19]. To measure the evolutionary pressures exerted on the genes of low, intermediate and high connectivity bacterial proteins, we use here K_a/K_s , and to measure evolutionary patterns of codon bias, we use the Effective Number of Codons (ENC) plots. The main finding of this work is the presence of a functional transition in bacterial PPI networks, ruled by degree connectivity k . The genes of proteins with high connectivities are under selective pressure, conserved, and essential. Below the transition ($k < 50$), the functional repertoire of low connectivity proteins is heterogeneous, whereas the genes of proteins with $k > 50$ mainly belong to the Cluster of Orthologous Genes (COG) J (related to translation, ribosomal structure and biogenesis), with just a few interesting hubs belonging to COGs I (Lipid transport and metabolism), K (Transcription) and L (Replication, recombination and repair). Moreover, we show here that in the degree distribution of each bacterial PPI network, there is a ubiquitous trace of an almost-invariant structure of conserved hubs, essentially due to the ribosomal protein complexes, mostly visible in the networks of bacteria with small genomes.

2. MATERIALS AND METHODS

2.1. Bacterial Dataset and Protein-protein Interaction Networks

We consider a set of 42 bacterial genomes (that we have previously investigated in [8]), shown in Table 1. Nucleotide sequences were downloaded from the FTP server of the National Center for Biotechnology Information [20]. These genomes were chosen in order to have a reasonably broad coverage of data concerning conservation, essentiality and selective pressure.

PPIs are obtained from the STRING database (Known and Predicted Protein-Protein Interactions, <https://string-db.org/>) [21]. We have chosen STRING because of its quite broad coverage of different bacterial species, useful to extend to multiple species we studied [7]. In STRING, each interaction is assigned with a confidence level or probability

w , evaluated by comparing predictions obtained by different techniques [22-24] with a set of reference associations, namely the functional groups of KEGG (Kyoto Encyclopedia of Genes and Genomes) [25]. In this way, interactions with high w are likely to be true positives, whereas a low w possibly corresponds to a false positive. As usually done in the literature, we consider only interactions with $w \geq 0.9$, a threshold that provides a fair balance between coverage and interaction reliability (see, for instance, the case study on *E. coli* reported in reference [7]). We denote the *degree* (number of connections) associated to each protein in each PPI network after the thresholding procedure by k . It is to be noted also that after applying the cut-off, we are left, for each network, with a number of isolated proteins (singletons, with no connections) that grow as n (where n is the number of proteins in the genome). These isolated proteins are not considered in the network analysis and are regarded as stemming from statistical noise or just appear isolated because the PPI data is incomplete.

It is known that PPIs of some species in our dataset might be much better known than others (*e.g.* *E. coli*). To take into account a potential bias in the dataset, we checked in Fig. (S1) of the Supplementary Information (bottom panel) that the densities of PPIs are high for small genomes and tend to be constant and not so different from that of *E. coli* in bacteria with bigger genomes, among which we collect here highly investigated pathogens.

The distinction between small and big genomes is a key emergent point in this work. We divided the set of 42 bacterial genomes into three groups, according to the number n of their genes: a) $n < 1000$, b) $1000 < n < 3000$ and c) $n > 3000$. In several figures in the Supplementary Information, we have addressed the dependence of various network properties on the size of the genome.

2.2. Gene Conservation

The Evolutionary Retention Index (ERI) [9] is a way of measuring the degree of conservation of a gene. In the present study, the ERI of a gene is the fraction of genomes, among those reported in Table 1, with at least an orthologous (same COG label) of the given gene. Then, as reminded in the Introduction, a low ERI value is related to a gene which is rather specific, common to a small number of genomes; whereas high ERI is characteristic of highly shared, putatively universal and essential genes.

We also make reference to another notion of gene conservation. Conserved genes are those which are subject to a purifying, conservative evolutionary pressure. To discriminate between genes subject to purifying selection and genes subject to positive selective Darwinian evolution, we use a classic but still widely used indicator, the ratio K_a/K_s between the number of nonsynonymous substitutions per nonsynonymous site (K_a) and the number of synonymous substitutions per synonymous site (K_s) [19]. Conserved genes are characterized by $K_a/K_s < 1$. We used K_a/K_s estimates by Luo [15] that are based on the method by Nej and Gojobori [26].

Table 1. Summary of the selected bacterial dataset. Organism name, abbreviation, class, RefSeq, STRING code, size of genome (number of genes *n*). Genomes annotated in the Database of Essential Genes (DEG) are highlighted with bold fonts. Classes are: Alphaproteobacteria(1), Betaproteobacteria(2), Gammaproteobacteria(3), Epsilonproteobacteria(4), Actinobacteria(5), Bacilli(6), Bacteroidetes(7), Clostridia(8), Deinococci(9), Mollicutes(10), Spirochaetales(11), Aquificae(12), Cyanobacteria(13), Chlamydiae(14), Fusobacteria(15), Thermotoga(16).

Organisms	Abbr.	Class	Ref Seq	STRING	n
<i>Mycoplasma genitalium</i> G37	myge	10	NC 000908	243273	475
<i>Buchnera aphidicola</i> Sg uid57913	busg	2	NC 004061	198804	546
<i>Mycoplasma pneumoniae</i> M129	mypn	10	NC 000912.1	272634	648
<i>Mycoplasma pulmonis</i> UAB CTIP	mypu	10	NC 002771	272635	782
<i>Chlamydia trachomatis</i> D/UW-3/CX	chtr	14	NC 000117.1	272561	894
<i>Treponema pallidum</i> Nichols	trpa	11	NC 000919.1	243276	1036
<i>Helicobacter pylori</i> 26695	hepy	4	NC 000915	85962	1469
<i>Aquifex aeolicus</i> VF5	aqae	12	NC 000918	224324	1497
<i>Campylobacter jejuni</i>	caje	4	NC 002163	192222	1572
<i>Haemophilus influenzae</i> Rd KW20	hain	3	NC 000907.1	71421	1610
<i>Streptococcus pyogenes</i> NZ131	stpy	6	NC 011375	471876	1700
<i>Francisella novicida</i> U112	frno	3	NC 008601	401614	1719
<i>Thermotoga maritima</i> MSB8	thma	16	NC 000853.1	243274	1858
<i>Neisseria gonorrhoeae</i> FA 1090 uid57611	nego	2	NC 002946	242231	1894
<i>Fusobacterium nucleatum</i> ATCC 25586	funu	15	NC 003454.1	190304	1983
<i>Brucella melitensis</i> bv. 1 str. 16M	brme	1	NC 003317.1	224914	2059
<i>Porphyromonas gingivalis</i> ATCC 33277	pogi	7	NC 010729	431947	2089
<i>Streptococcus sanguinis</i>	stsa	6	NC 009009	388919	2270
<i>Vibrio cholerae</i> N16961	vich	3	NC 002505	243277	2534
<i>Staphylococcus aureus</i> N315	stau	6	NC 002745.2	158879	2582
<i>Deinococcus radiodurans</i> R1	dera	9	NC 001263.1	243230	2629
<i>Agrobacterium tumefaciens</i> (fabrum)	agtu	1	NC 003062	176299	2765
<i>Xylella fastidiosa</i> 9a5c	xyfa	3	NC 002488	160492	2766
<i>Staphylococcus aureus</i> NCTC 8325	stau	6	NC 007795	93061	2767
<i>Listeria monocytogenes</i> EGD-e	limo	6	NC 003210.1	169963	2867
<i>Synechocystis</i> sp. PCC 6803	syp	13	NC 000911.1	1148	3179
<i>Burkholderia thailandensis</i> E264	buth	2	NC 007651	271848	3276
<i>Sinorhizobium meliloti</i> 1021	sime	1	NC 003047.1	266834	3359
<i>Burkholderia pseudomallei</i> K96243	bups	3	NC 006350	272560	3398
<i>Ralstonia solanacearum</i> GM11000	raso	2	NC 003295.1	267608	3436
<i>Clostridium acetobutylicum</i> ATCC 824	clac	8	NC 003030.1	272562	3602
<i>Caulobacter crescentus</i>	cacr	1	NC 011916	565050	3885
<i>Mycobacterium tuberculosis</i> H37Rv	mytu	5	NC 000962.3	83332	3936
<i>Escherichia Coli</i> K-12 MG1655	esco	3	NC 000913.3	511145	4004
<i>Shewanella oneidensis</i> MR-1	shon	3	NC 004347	211586	4065
<i>Bacillus subtilis</i> 168	basu	6	NC 000964	224308	4175
<i>Salmonella enterica</i> serovar Typhi	saen	3	NC 004631	209261	4352
<i>Bacteroides thetaiotaomicron</i> VPI-5482	bath	7	NC 004663	226186	4778
<i>Sphingomonas wittichii</i> RW1	spwi	1	NC 009511	392499	4850
<i>Pseudomonas aeruginosa</i> UCBPP-PA14	psae	3	NC 008463	208963	5892
<i>Mesorhizobium loti</i> MAFF303099	melo	1	NC 002678.2	266835	6743
<i>Rickettsia prowazekii</i> str. Madrid E	ripr	1	NC 000963.1	272947	8433

2.3. Gene Essentiality

We used the Database of Essential Genes (DEG, www.essentialgene.org) [15], which classifies a gene as either essential or nonessential, on the basis of a combination

of experimental evidence (null mutations or transposons) and general functional considerations. DEG collects genomes from Bacteria, Archaea and Eukarya, with different degrees of coverage [27, 28]. Of the 42 bacterial

genomes we considered, only 23 are covered-in total or partially-by DEG, as indicated in Table 1.

2.4. K_d/K_s

K_d/K_s is the ratio of nonsynonymous substitutions per nonsynonymous site (K_a) to the number of synonymous substitutions per synonymous site (K_s) [19]. This parameter is widely accepted as a straightforward and effective way of separating genes subject to purifying evolutionary selection ($K_d/K_s < 1$) from genes subject to positive selective Darwinian evolution ($K_d/K_s > 1$). There are different methods to evaluate this ratio, though the alternative approaches are quite consistent among themselves. For the sake of comparison, we have used here the K_d/K_s estimates by Luo *et al.* [15], which are based on the Nej and Gojobori method [26]. It must be noted that each genome has a specific average level of K_d/K_s [7]. Average values of K_d/K_s are shown for low, intermediate, and high connectivity bins of genes.

2.5. ENC Plot

The ENC-plot is a well-known tool to investigate the patterns of synonymous codon usage in which the ENC (Effective Number of Codons) values are plotted against GC_3 Guanine and Cytosine Content at the third codon position. The formula of ENC values expected under the hypothesis of pure mutational bias (no selection) is given by:

$$ENC = 2 + s + \frac{29}{s^2 + (1-s)^2} \quad (1)$$

where s represents the value of GC_3 [29]. When the corresponding points fall near the expected neutral curve, mutations that enforce the typical mutational bias of the species are the main factor affecting the observed codon diversity. Whereas when the corresponding points fall considerably below the expected curve, the observed codon usage bias of the species is mainly affected by natural selection. To quantitatively represent the balance between mutational bias and selective natural pressure, we parametrise the ENC formula to be used in non-linear fits to the experimental data:

$$ENC = a + b * s + \frac{29}{s^2 + d*(1-s)^2} \quad (2)$$

ENC plots of genes corresponding to low, intermediate and high connectivity proteins are shown in Fig. (S5) of the Supplementary Information. The best-fit parameters for the three groups of genes are collected in Table S1.

2.6. Clusters of Orthologous Proteins

We use the functional annotation given in the database of orthologous groups of proteins (COGs) from Koonin's group, available at <http://ncbi.nlm.nih.gov/COG/> [30, 31]. We consider 15 functional COG categories Table 2, excluding the generic categories R and S for which functional annotation is either too general or missing.

3. RESULTS AND DISCUSSION

Degree Distribution of PPI Networks. We start by studying the degree distributions $P(k)$ observed in bacterial PPIs. We first recall that such distribution was found to be scale-free in *E. coli* [7, 32-34], meaning that the corresponding PPI network features a large number of poorly connected proteins and a relatively small number of highly connected hubs. In order to assess the generality of this observation, we compute $P(k)$ for each genome in Table 1 (plots are reported in Figs. (S3 and S4) of the Supplementary Information). Note that, despite the fact that PPI networks of different bacteria have different sizes and densities, their average connectivity and the support of their $P(k)$ are very similar. Thus, we can superpose all the considered bacterial degree distributions without the need to normalise the support of each $P(k)$. When doing so, we observe two distinct regimes (Fig. 1). For low values of $k < 40$, the distribution is approximately scale-free: $P(k) \sim k^{-\gamma}$ ($\gamma = 2.48$). This scaling behaviour is consistent with previous studies on the genomes of yeast, worms and flies [35] and on co-conserved PPIs in some bacteria [36]. The scale-free nature of bacterial PPIs is still a matter of debate, and a rough discussion of the origin of this feature is out of the scope of this paper. In this work, we generally confirm that, as said above, there is, as expected, a large number of poorly connected proteins and a small number of hubs. Remarkably, for higher values of k , the distribution deviates from a power law, and a bump with a Gaussian-like shape emerges.

Table 2. Functional classification of COG clusters.

COG ID	Functional Classification
Information Storage and Processing	
J	Translation, ribosomal structure and biogenesis K Transcription
L	Replication, recombination and repair
Cellular Processes and Signaling	
D	Cell cycle control, cell division, chromosome partitioning T Signal transduction mechanisms
M	Cell wall/membrane/envelope biogenesis N Cell motility
O	Post-translational modification, protein turnover, chaperones
Metabolism	
C	Energy production and conversion
G	Carbohydrate transport and metabolism
E	Amino acid transport and metabolism
F	Nucleotide transport and metabolism
H	Coenzyme transport and metabolism
I	Lipid transport and metabolism
P	Inorganic ion transport and metabolism

This feature, visible for $k > 40$ may be due to the contribution of proteins belonging to large complexes [37]. From the whole set of observations presented in this paper, the bump in the $P(k)$ is due to the complexity of ribosomal interactions. Indeed, if one recalculates the degree distribution of a dataset in which the ribosomal proteins are removed, the bump is not present (Fig. (1), empty dots). Moreover, if we consider the separate contribution of essential and nonessen-

tial genes to the $P(k)$ (for DEG-annotated genomes), we see that the bump is present only in the degree distribution of essential genes. It is to be noted also that the degree distributions for essential and nonessential genes are well separated and the average degree is systematically higher for essential genes than for nonessential ones, consistently with previous findings [35]. Remarkably, we have shown in a previous paper [8] that the number of essential genes in bacteria is close to 500 and does not depend on the size of the genome. To correctly interpret the emergence of the bump in the average $P(k)$ in Fig. (1), it is worth pointing out the distinction between small and not so small genomes. In the small genomes, almost all the genes are essential, and among the essential genes, those belonging to COG J (functions related to translation and ribosomal structure and biogenesis) play a major and ubiquitous role. In (Fig. S2), we have checked that the bump that emerges in Fig. (1) as a feature of essential and conserved genes is quite visible in the $P(k)$ of small genomes, whereas, there seems a confusion in the case of bigger genomes. This might be interpreted as a dilution effect; in the networks of bigger genomes, there are a lot of specific interactions besides the essential ones. Then, averaging $P(k)$ over small, intermediate and big genomes (Fig. S2 in Supplementary Information), we can safely interpret the bump as an emerging feature due to a core of highly connected proteins (connectivities $k \geq 40$), which is mostly contributed, in the average, by degree distributions from PPIs of bacteria with small genomes Figs. S2-S4. From all the considerations above, we exclude that this bump, observed here for the first time, could emerge just because that part of the PPI is much more investigated than other subnetworks. It is there because the ribosome is there, in all bacteria Table 3.

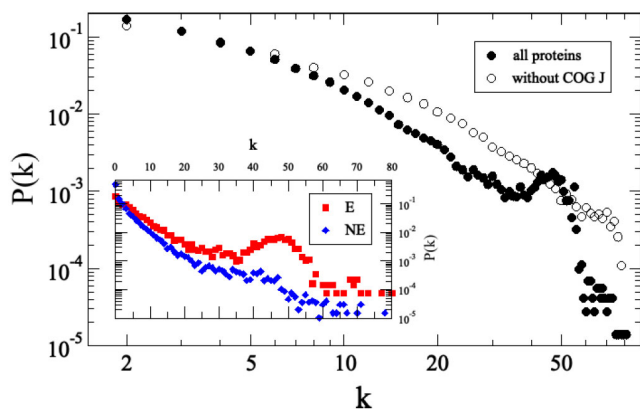


Fig. (1). Probability distribution $P(k)$ for the number of connections k of each protein averaged over the bacterial species considered in Table 1 (full dots), compared with the degree distribution after removal of the proteins corresponding to genes in COG J, related to translational processes (empty dots). Inset: $P(k)$ for essential (E) and nonessential (NE) genes, averaged over DEG-annotated genomes. Note that the average degree is higher for essential genes than for nonessential ones, and the two probability distributions are quite distinct. The region of the curve for low k can be well approximated by a power law [38]. (A higher resolution / colour version of this figure is available in the electronic copy of the article).

3.1. PPI Connectivity and Gene Conservation

We now investigate whether the connectivity k of a protein in a PPI network drives a transition in the degree of conservation (as measured by ERI) of the corresponding genes. Fig. (2) displays the average value and the spread of ERI in genes relative to bins of proteins that are iso-connected in the PPIs of different species. As a general feature, we observe that, on average, the genes of highly connected proteins are highly conserved among the bacterial species we consider that constitute a reasonably wide sample of different evolutionary adaptations. The same Fig. (2) shows that if $k < 50$, then the ERI highly fluctuates between different samples of proteins with the same k , in different species. For high connectivities (above $k = 50$), the ERI is close to 1, with a drastic drop in the fluctuation (as shown in the inset). This observation points to the existence of an almost-invariant structure of conserved hubs, in each bacterial PPI, sustained by highly conserved genes. We can conclude, as a rule of thumb, that a protein with connectivity degree of 40 or more is likely to be coded by a gene shared by at least 80% of the species in a generic pool of bacteria. At the moment, we do not have a general explanation for this apparent threshold. Let us just propose, as a heuristic observation, the existence of an almost-critical value of connectivity to be set between 40 and 50, that corresponds to the connectivity of the core of proteins specifically involved, as we have alluded to in the previous paragraph, to the ubiquitous ribosomal functions Tables 4 and 5.

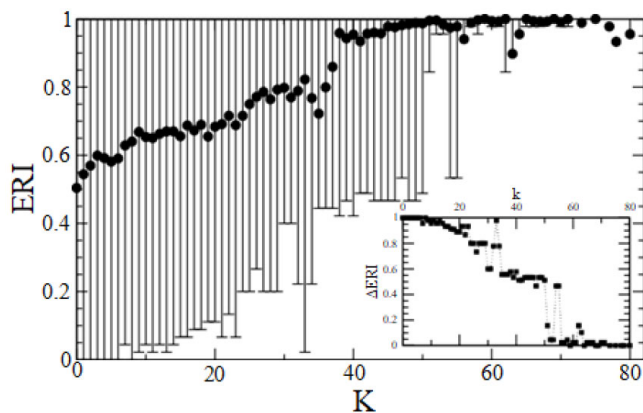


Fig. (2). Average ERI values of bacterial genes as a function of the degrees k of the corresponding proteins, for all the considered genomes. Error bars are standard deviations of ERI values associated to a given k value. Inset: amplitude of the error bar (ΔERI) as a function of k .

3.2. Evolutionary Pressure and PPI Connectivity

We then look at the evolutionary pressure exerted on genes whose proteins have different connectivities. The graph in Fig. (3) shows the ratio K_d/K_s for groups of genes binned by the connectivity k of the corresponding proteins, for all the 42 bacterial species in Table 1. As is well known, this ratio K_d/K_s provides a straightforward indication of the balance between a positive driving *Darwinian selection*

(when the numerator prevails) and a *purifying*, stabilising selection (acting against change in genes for which the denominator prevails).

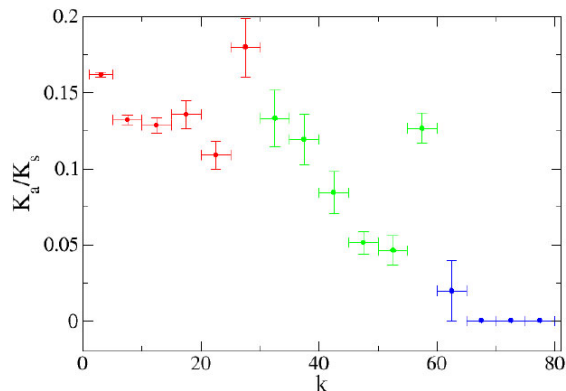


Fig. (3). (K_d/K_s) of groups of genes corresponding to proteins with different connectivity degrees k . As in the following Fig. (4) and in Fig. S5, low connectivities are shown in red, intermediate in green and high connectivities in blue. (A higher resolution / colour version of this figure is available in the electronic copy of the article).

We see that the more connected proteins correspond to genes that are subject to an increasing purifying evolutionary pressure. Indeed, the ratio (K_d/K_s) is less than 1 in all bins of connectivity and systematically decreases as a function of k . A decreasing ratio generally indicates an increasing role of purifying, conservative, Darwinian, evolutionary pressure on the corresponding set of genes. This is a reasonable result, pointing out that the groups of genes that support conserved structures of connectivity in the PPIs are more constrained, in evolution than the genes of less interacting proteins.

To add evidence to this observation, we have also considered ENC plots for sets of genes binned by the connectivities of the corresponding proteins. Interestingly, the ENC data in Fig. (S5) of Supplementary Information are fully consistent with those in Fig. (3). In the ENC plots, the points associated with low connectivity proteins (red) are closer to the so-called Wright's profile (represented there as solid black lines) than those associated to proteins with intermediate and high connectivities (green and blue lines). Fig. (4) stresses this observation in a more quantitative way by showing that in the ENC plots, the average distance from Wright's profile monotonously increases with k . Overall, the above results clearly indicate that codon bias and GC content of high connectivity genes are more under selective Darwinian pressure than genes coding for low-connectivity proteins, in which the rate of accepted mutations is mainly ruled by neutral mutational bias. These observations point out that the almost-invariant structure of protein hubs we alluded to in the previous paragraph, is supported by an underlying set of genes that are under strong mutational control; an expected result, perhaps, but clearly seen, here, as a general feature associated with ribosomal ubiquitous and conserved functions.

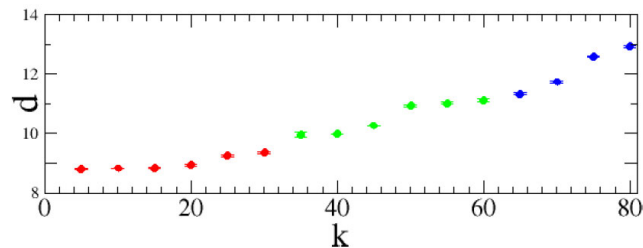


Fig. (4). ENC plot and connectivity. Each point in this graph represents a group of genes, characterised by the average connectivity k of the corresponding proteins in the PPI network and by the average euclidean distance d , in the ENC plot, from Wright's theoretical curve. Different groups of genes are represented with different colors as a function of k . As shown in the previous Fig. (3), red corresponds to low connectivities, green to intermediate and blue to high connectivities (Fig. S5). The distance from the curve clearly increases with k . Wright's curve corresponds, in the ENC plot, to pure mutational bias (Eq. 1), then higher connectivities of the proteins imply bigger evolutionary selective pressure on the corresponding group of genes. (A higher resolution / colour version of this figure is available in the electronic copy of the article).

PPI and Essentiality. To further investigate the relationship between gene essentiality and protein connectivities, we consider DEG-annotated genomes and classify interactions between proteins (links) making reference to the essentiality of the corresponding genes. We distinguish three sets of links: $|ee|$ (linking proteins from two essential genes), $|\bar{e}\bar{e}|$ (from two nonessential genes) and $|e\bar{e}|$ (from an essential gene and a nonessential one). We then compute the *density* of these sets of links respectively as:

$$\begin{aligned} \rho_{ee} &= \frac{|ee|}{\frac{1}{2}E(E-1)} \\ \rho_{\bar{e}\bar{e}} &= \frac{|\bar{e}\bar{e}|}{\frac{1}{2}NE(NE-1)} \\ \rho_{e\bar{e}} &= \frac{|e\bar{e}|}{\frac{1}{2}E+NE} \end{aligned} \quad (3)$$

where E and NE denote the number of essential and non-essential genes, respectively (self-connection are excluded in our analysis). The denominator is the maximum possible value of the numerator, corresponding to the fully-connected graph. Such densities are then compared with the overall density of the network-restricted to genes classified as either essential or nonessential:

$$\langle \rho \rangle = \frac{|ee|+|\bar{e}\bar{e}|+|e\bar{e}|}{\frac{1}{2}(E+NE)(E+NE-1)} \quad (4)$$

We use the ratios $r_{ee} = \frac{\rho_{ee}}{\rho}$, $r_{\bar{e}\bar{e}} = \frac{\rho_{\bar{e}\bar{e}}}{\rho}$, $r_{e\bar{e}} = \frac{\rho_{e\bar{e}}}{\rho}$ to assess the level of connectivity of the subnetworks with respect to the overall connectivity. Table 3 shows that subnetworks of essential genes are far denser than the overall networks, and that, in general, essential and nonessential genes tend to form network components that are weakly interconnected. This happens because many essential genes encode for ribo-

somal proteins, which in turn are localised in the ribosomal complex where they have a high probability of interacting [39] Table 3 of [8], which shows approximately 25% of essential genes fall into COG J. Figs. (S6 and S7) of the Supplementary Information collect the superposed adjacency matrices of the $|ee|$ (red dots), $|e\bar{e}|$ (violet dots) and $|\bar{e}\bar{e}|$ (blue dots) subnetworks that display such network features for each individual species. These graphs confirm the dominance of the interactions between the proteins of essential genes (red dots) in the small genomes. The adjacency matrices of bacteria with intermediate and big genomes are dominated by interactions involving proteins supported by non-essential genes (blue dots).

PPI Connectivity and Functional Specialisation. For each PPI network, we define the conditional probability (Bayes' theorem) that a protein with degree k belongs to a given COG as:

$$P(\text{COG}|k) = P(k|\text{COG})P(\text{COG})/P(k), \tag{5}$$

where $P(k)$ is the degree distribution in the PPI network, $P(\text{COG})$ is the frequency of that COG in the proteome, and $P(k|\text{COG})$ is the degree distribution restricted to that COGs. Fig. (5) shows the COG spectrum as a function of k over all the bacteria here considered. Interestingly, we again note a marked transition. Below k 40, the COG spectrum is quite heterogeneous: genes corresponding to proteins with low connectivity are spread over several COGs, which correspond to different functions (Table 2).

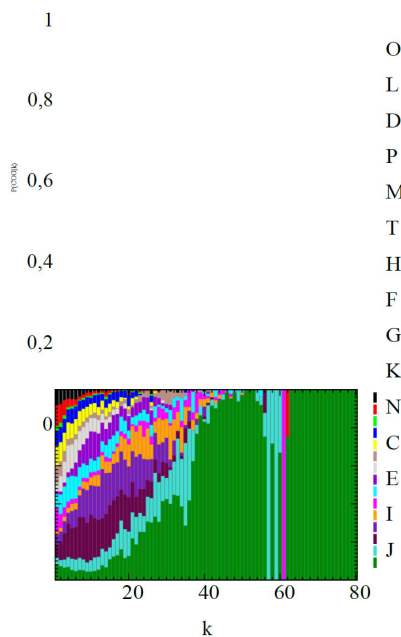


Fig. (5). Probability distribution $P(\text{COG } k)$ of belonging to a given COG for proteins with degree k , overall considered genomes. Proteins with low connectivity have a very heterogeneous COG composition, whereas those with high k basically belong only to COG J. (A higher resolution / colour version of this figure is available in the electronic copy of the article).

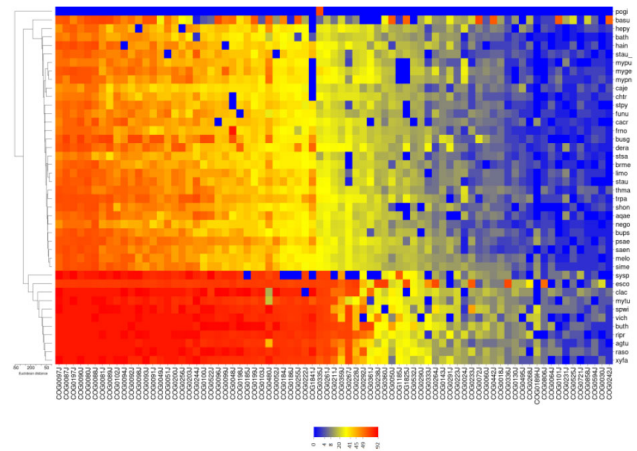


Fig. (6). Heat map of the connectivity degree of the protein as distributed over the COG J genes with ERI=1, in each species. Genes are sorted by decreasing average degree. We note that those genes which correspond to degrees bigger than 40 are conserved for all species. Details of these genes are in Table 5. (A higher resolution / colour version of this figure is available in the electronic copy of the article).

The transition shows that proteins with more than 40 interactions are likely to be coded by genes belonging to COG J. There are yet a handful of outliers, hubs with connectivities between 57 and 62, that belong to COG I (related to lipid transport and metabolism) and K and L (which, together with J, define the functional class of information storage and processing). The list of these outliers is reported in Table 4. Interestingly, they correspond to RNA polymerases and to enzymes involved in acetate metabolism. But, which are the genes of COG J that drive the transition? In the next Fig. (6), we show which genes are the main characters in the transition. We then investigate the connectivities of the highly conserved (ERI=1, shared by all the species in Table 1) genes belonging to COG J and whose proteins have connectivities bigger than 40. These highly shared genes corresponding to cores of highly connected ribosomal proteins are listed in Table 5. In the heat map of Fig. (6), we sort each gene in the COG J in order of descending degree, species by species, and we see there is a core of genes (in red, lower left sector) that correspond to highly connected proteins, which are also highly shared (ERI =1, see Table 5) among all the species we considered. It is quite clear that in the heat map of Fig. (6) the 42 species in this study can be split into at least two groups (see the cladogram on the left). In one group the group of species at the Bottom in Fig. (5) there is a shared set of genes (the red band at the bottom-left side of the heat map) corresponding to a common core of highly connected ribosomal proteins. This remarkable observation suggests that the species in this group (namely, *Synechocystis* sp. PCC 6803, *Escherichia coli* K-12 MG1655, *Clostridium acetobutylicum* ATCC 824, *Mycobacterium tuberculosis* H37Rv, *Sphingomonas wittichii* RW1, *Vibrio cholerae* N16961, *Burkholderia thailandensis* E264, *Rickettsia prowazekii* str. Madrid E, *Agrobacterium tumefaciens* (fab-

Table 3. Relative density values r for PPI subnetworks between essential genes (r_{ee}), between nonessential genes ($r_{\bar{e}\bar{e}}$) and between essential and nonessential genes $r_{e\bar{e}}$, for each DEG-annotated bacterial genome.

Organisms	r_{ee}	$r_{\bar{e}\bar{e}}$	$r_{e\bar{e}}$
basu	44.46	0.80	0.11
bath	20.07	0.76	0.25
bups	6.21	0.83	0.27
buth	18.69	0.70	0.22
cacr	18.40	0.70	0.15
caje	3.65	0.82	0.32
esco	2.91	0.88	0.31
frno	9.84	0.52	0.18
hain	1.65	1.15	0.27
hepy	2.91	0.78	0.38
myge	1.42	0.29	0.08
mypu	3.42	0.22	0.12
mytu	8.09	0.78	0.23
pogi	11.03	0.41	0.21
psae	9.85	0.92	0.16
saen	28.80	0.81	0.12
shon	6.50	0.64	0.16
spwi	15.47	0.74	0.22
stau	23.05	0.58	0.23
stau	21.89	0.64	0.16
stpy	9.30	0.73	0.23
tsa	30.65	0.61	0.22
vich	8.37	0.81	0.19

Table 4. Specific hubs. In this table we detail which proteins populate the few bins of connectivity around $k = 60$ in Fig. (5).

k	COG	Gene	Protein
57	1250I	paaH	3-hydroxyadipyl-CoA dehydrogenase, NADdependent
	0365I	acs	acetyl-CoA synthetase
58	0222J	rpL	50S ribosomal subunit protein L7/L12
	0335J	rpL	50S ribosomal subunit protein L19
	0267J	rpmG	50S ribosomal subunit protein L33
	0365I	acs	acetyl-CoA synthetase
59	0183I	paaJ	3-oxoadipyl-CoA3-oxo-5,6-dehydrosuberyl-CoA thio-lase
	1960I	ydiO	putative acyl-CoA dehydrogenase
	0183I	atoB	acetyl-CoA acetyltransferase
60	0197J	rpIP	50S ribosomal subunit protein L16
	0088J	rpID	50S ribosomal subunit protein L4
	0197J	rpIP	50S ribosomal subunit protein L16
	0087J	rpIC	50S ribosomal subunit protein L3
	1960I	aidB	putative acyl-CoA dehydrogenase
61	0085K	rpoB	RNA polymerase, beta subunit
	0202K	rpoA	RNA polymerase, alpha subunit
62	0087J	rpIC	50S ribosomal subunit protein L3
	0052J	rpsB	30S ribosomal subunit protein S2
	2965L	PriB	ribosomal replication protein

Table 5. Genes belonging to COG J with average degree bigger than 40 Fig. (6). All these genes are conserved, common to all species (ERI=1), and drive the transition shown in Fig. (5).

COG	Genes Name	<k >
COG0097J	50S ribosomal protein L6	60.24
COG0087J	50S ribosomal protein L3	60.19
COG0197J	50S ribosomal protein L16	60.19
COG0090J	50S ribosomal protein L2	60.14
COG0080J	50S ribosomal protein L11	60.12
COG0088J	50S ribosomal protein L4	60.12
COG0081J	50S ribosomal protein L1	58.19
COG0089J	50S ribosomal protein L23	57.88
COG0102J	50S ribosomal protein L13	57.45
COG0094J	50S ribosomal protein L5	57.21
COG0092J	30S ribosomal protein S3	57.12
COG0098J	30s ribosomal protein S5	57.10
COG0093J	50S ribosomal protein L14	57.00
COG0091J	50S ribosomal protein L22	56.24
COG0049J	30S ribosomal protein S7	55.31
COG0051J	30S ribosomal protein S10	55.24
COG0200J	50S ribosomal protein L15	55.12
COG0256J	50S ribosomal protein L18	54.86
COG0203J	50S ribosomal protein L17	54.43
COG0244J	50S ribosomal Protein L10	54.19
COG0100J	30S ribosomal protein S11	53.76
COG0522J	30S ribosomal protein S4	53.43
COG0096J	30S ribosomal protein S8	53.10
COG0099J	30S ribosomal protein S13	52.88
COG0048J	30S ribosomal protein S12	52.14
COG0198J	50S ribosomal protein L24	50.83
COG0185J	30S ribosomal protein S19	50.52
COG0199J	30S ribosomal protein S14	50.45
COG0103J	30S ribosomal protein S9	49.45
COG0480J	tetracycline resistance protein. tetM	47.90
COG0052J	30S ribosomal protein S2	47.69
COG0184J	30S ribosomal protein S15	45.95
COG0186J	30S ribosomal protein S17	44.60
COG0255J	50S ribosomal protein L29	43.95
COG0222J	50S ribosomal protein L7/L12	42.43
COG1841J	50S ribosomal protein L30	40.71

rum), *Ralstonia solanacearum* GMI1000, *Xylella fastidiosa* 9a5c) should have a common structural and functional organisation of their ribosomes, an interesting point to be further investigated. In the rest of the species, the connectivity of the proteins, corresponding to the highly shared COG J genes, with $k > 40$ is more heterogeneous. We can conclude that the abrupt transition shown in Fig. (5) is driven by a subset of COG J genes which are highly conserved among a sub-

set of species and are listed in Table 5. As one can see, these genes correspond to a specific subset of ribosomal proteins in the small and large subunits that should be further investigated in their functional and structural role.

CONCLUSION

Connectivity analysis of biological networks, such as protein-protein interaction or metabolic networks, has demonstrated that structural features of network subgraphs are correlated with biological functions [40, 41]. For instance, it was shown that highly connected patterns of proteins in a PPI are fundamental to cell viability [42]. In this work, we have shown the existence of a functional transition in bacterial species, ruled by the connectivity of proteins in the PPI networks (Fig. 5). The critical threshold in k of the transition is located between $k=40$ and $k=50$. Proteins that have connectivities above the threshold are mostly encoded by genes that are conserved, under selective pressure (as measured both by ERI and K_e/K_s) and essentiality. Moreover, the functional repertoire above the threshold mainly focuses on the COG J (translation, ribosomal structure and biogenesis), with just a few interesting hubs belonging to COGs I (Lipid transport and metabolism), K (Transcription) and L (Replication, recombination and repair).

Indeed, the PPI network of each bacterial species is characterised by a highly connected core of conserved ribosomal proteins, the components of multi-subunit complexes whose corresponding genes are mostly essential [32, 36] and code for supra-molecular complexes that pile up in the bump we have observed for the degree distribution (Fig. 1). Hence, what we see here is essentially the ribosome and related protein complexes such as RNA Polymerase. Indeed, the ribosome is the only molecular machine in bacteria in which a given protein could legitimately have 40 or more protein binding partners, with the help of rRNA mediating interactions [43].

It is reasonable to admit that, since there are bacterial species that are much more investigated than others, comparative statistical studies of bacterial PPIs might be particularly biased by the choice of the sample of genomes to be included in the study. Our dataset is no exception. In order to address this hard to settle problem in our study, we have checked Fig. (S1) that in our study, we have included small genomes (*i.e.* less than 1000 genes) whose PPIs have densities (a rough proxy for the coverage of the interactions in the network) that are higher than those of bigger genomes. The group of small genomes comprises Buchnera, Chlamydia, and Mycoplasmas, whereas bigger genomes refer mostly to illustrious pathogens that are surely among the most investigated bacterial species. The densities of the networks of these species are quite similar and comparable with that of *E. coli*. As a general rule, and quite obviously, the networks relative to small genomes are better covered in the STRING database (after the application of a conservative cutoff $w = 900$) than those relative to bigger genomes. Interestingly, we have shown Figs. (6 and 7) in Supplementary Information) that, indeed, the PPI adjacency matrices of bacteria with

small genomes are dominated by the interactions constituting the ribosomal complex. In the adjacency matrices of the PPIs of bacteria with bigger genomes, the cloud of interactions between the proteins of nonessential genes tends to superpose to the ever-present ribosomal core. In conclusion, we believe to have convincingly shown that bacterial PPIs are characterised by the presence of a highly connected structure, associated with the ribosomal functions, and particularly visible in bacteria with small genomes. We believe that the observations we have presented here could be of some utility for the prediction of gene essentiality, based on the knowledge of PPI networks, and for the prediction of interactions between proteins, based on genetic information [44, 45]. It is interesting to note that our results are consistent with a previous study based on inferred bacterial conserved networks based on phylogenetic profiles [36]. This work suggests to further and systematically investigate how the structure of the PPI networks is correlated with multiple networks at the genetic level, at least in unicellular organisms. In particular, we believe that a recent approach based on the introduction of the multiple-layer networks could be of great potential interest (*e.g.* to search for a general scheme behind antimicrobial resistance [46-50]).

ETHICS APPROVAL AND CONSENT TO PARTICIPATE

Not applicable.

HUMAN AND ANIMAL RIGHTS

No Animals/Humans were used for studies that are basis of this research.

CONSENT FOR PUBLICATION

Not applicable.

AVAILABILITY OF DATA AND MATERIALS

Not applicable.

FUNDING

None.

CONFLICT OF INTEREST

The authors declare no conflict of interest, financial or otherwise.

ACKNOWLEDGEMENTS

Declared none.

SUPPLEMENTARY MATERIAL

Supplementary material is available on the publisher's website along with the published material.

REFERENCES

- [1] Drewes, G.; Bouwmeester, T. Global approaches to protein-protein interactions. *Curr. Opin. Cell Biol.*, **2003**, *15*(2), 199-205.

- [http://dx.doi.org/10.1016/S0955-0674\(03\)00005-X](http://dx.doi.org/10.1016/S0955-0674(03)00005-X) PMID: 12648676
- [2] Golemis, E.; Adams, P.D. *Protein-protein Interactions: A Molecular Cloning Manual*; Cold Spring Harbor Laboratory Press, **2005**.
- [3] von Mering, C.; Krause, R.; Snel, B.; Cornell, M.; Oliver, S.G.; Fields, S.; Bork, P. Comparative assessment of large-scale data sets of protein-protein interactions. *Nature*, **2002**, *417*(6887), 399-403.
<http://dx.doi.org/10.1038/nature750> PMID: 12000970
- [4] Tong, A.H.; Drees, B.; Nardelli, G.; Bader, G.D.; Brannetti, B.; Castagnoli, L.; Evangelista, M.; Ferracuti, S.; Nelson, B.; Paoluzi, S.; Quondam, M.; Zucconi, A.; Hogue, C.W.V.; Fields, S.; Boone, C.; Cesareni, G. A combined experimental and computational strategy to define protein interaction networks for peptide recognition modules. *Science*, **2002**, *295*(5553), 321-324.
<http://dx.doi.org/10.1126/science.1064987> PMID: 11743162
- [5] Shatsky, M.; Allen, S.; Gold, B.L.; Liu, N.L.; Juba, T.R.; Reveco, S.A.; Elias, D.A.; Prathapam, R.; He, J.; Yang, W.; Szakal, E.D.; Liu, H.; Singer, M.E.; Geller, J.T.; Lam, B.R.; Saini, A.; Trotter, V.V.; Hall, S.C.; Fisher, S.J.; Brenner, S.E.; Chhabra, S.R.; Hazen, T.C.; Wall, J.D.; Witkowska, H.E.; Biggin, M.D.; Chandonia, J.M.; Butland, G. Bacterial interactomes: Interacting protein partners share similar function and are validated in independent assays more frequently than previously reported. *Mol. Cell. Proteomics*, **2016**, *15*(5), 1539-1555.
<http://dx.doi.org/10.1074/mcp.M115.054692> PMID: 26873250
- [6] Harry, C.; Wimble, J.C.; Shary, S.; Wuchty, S.; Uetz, P. Bacterial protein meta-interactomes predict cross-species interactions and protein function. *BMC Bioinformatics*, **2017**, *18*(1), 171.
<http://dx.doi.org/10.1186/s12859-017-1585-0> PMID: 28298180
- [7] Dilucca, M.; Cimini, G.; Semmoloni, A.; Deiana, A.; Giansanti, A. Codon bias patterns of *E. coli*'s interacting proteins. *PLoS One*, **2015**, *10*(11), e0142127.
<http://dx.doi.org/10.1371/journal.pone.0142127> PMID: 26566157
- [8] Dilucca, M.; Cimini, G.; Giansanti, A. Essentiality, conservation, evolutionary pressure and codon bias in bacterial genomes. *Gene*, **2018**, *663*, 178-188.
<http://dx.doi.org/10.1016/j.gene.2018.04.017> PMID: 29678658
- [9] Gerdes, S.Y.; Scholle, M.D.; Campbell, J.W.; Balázsi, G.; Ravasz, E.; Daugherty, M.D.; Somera, A.L.; Kyrpides, N.C.; Anderson, I.; Gelfand, M.S.; Bhattacharya, A.; Kapatral, V.; D'Souza, M.; Baev, M.V.; Grechkin, Y.; Mseeh, F.; Fonstein, M.Y.; Overbeek, R.; Barabási, A.L.; Oltvai, Z.N.; Osterman, A.L. Experimental determination and system level analysis of essential genes in *Escherichia coli* MG1655. *J. Bacteriol.*, **2003**, *185*(19), 5673-5684.
<http://dx.doi.org/10.1128/JB.185.19.5673-5684.2003> PMID: 13129938
- [10] Fang, G.; Rocha, E.; Danchin, A. How essential are nonessential genes? *Mol. Biol. Evol.*, **2005**, *22*(11), 2147-2156.
<http://dx.doi.org/10.1093/molbev/msi211> PMID: 16014871
- [11] Peng, C.; Gao, F. Protein localization analysis of essential genes in prokaryotes. *Sci. Rep.*, **2014**, *4*, 6001.
<http://dx.doi.org/10.1038/srep06001> PMID: 25105358
- [12] Lin, Y.; Gao, F.; Zhang, C.-T. Functionality of essential genes drives gene strand-bias in bacterial genomes. *Biochem. Biophys. Res. Commun.*, **2010**, *396*(2), 472-476.
<http://dx.doi.org/10.1016/j.bbrc.2010.04.119> PMID: 20417622
- [13] Hurst, L.D.; Smith, N.G. Do essential genes evolve slowly? *Curr. Biol.*, **1999**, *9*(14), 747-750.
[http://dx.doi.org/10.1016/S0960-9822\(99\)80334-0](http://dx.doi.org/10.1016/S0960-9822(99)80334-0) PMID: 10421576
- [14] Jordan, I.K.; Rogozin, I.B.; Wolf, Y.I.; Koonin, E.V. Essential genes are more evolutionarily conserved than are nonessential genes in bacteria. *Genome Res.*, **2002**, *12*(6), 962-968.
<http://dx.doi.org/10.1101/gr.87702> PMID: 12045149
- [15] Luo, H.; Gao, F.; Lin, Y. Evolutionary conservation analysis between the essential and nonessential genes in bacterial genomes. *Sci. Rep.*, **2015**, *5*, 13210.
<http://dx.doi.org/10.1038/srep13210> PMID: 26272053
- [16] Ish-Am, O.; Kristensen, D.M.; Rupp, E.; David, M.K.; Rupp, E. Evolutionary conservation of bacterial essential metabolic genes across all bacterial culture media. *PLoS One*, **2015**, *10*(4), e0123785.
<http://dx.doi.org/10.1371/journal.pone.0123785> PMID: 25894004
- [17] Alvarez-Ponce, D.; Sabater-Muñoz, B.; Toft, C.; Ruiz-González, M.X.; Fares, M.A.; Sabater-Munoz, B.; Toft, C.; Ruiz-González, M.X.; Fares, M.A. Essentiality is a strong determinant of protein rates of evolution during mutation accumulation experiments in *Escherichia coli*. *Genome Biol. Evol.*, **2016**, *8*(9), 2914-2927.
<http://dx.doi.org/10.1093/gbe/evw205> PMID: 27566759
- [18] Bergmiller, T.; Ackermann, M.; Silander, O.K. Patterns of evolutionary conservation of essential genes correlate with their compensability. *PLoS Genet.*, **2012**, *8*(6), e1002803.
<http://dx.doi.org/10.1371/journal.pgen.1002803> PMID: 22761596
- [19] Hurst, L.D. The Ka/Ks ratio: diagnosing the form of sequence evolution. *Trends Genet.*, **2002**, *18*(9), 486-487.
[http://dx.doi.org/10.1016/S0168-9525\(02\)02722-1](http://dx.doi.org/10.1016/S0168-9525(02)02722-1) PMID: 12175810
- [20] Benson, D.A.; Cavanaugh, M.; Clark, K.; Karsch-Mizrachi, I.; Lipman, D.J.; Ostell, J.; Sayers, E.W. Genbank *Nucleic Acids Res.*, **2013**, *41*(D1), D36-D42.
- [21] Szklarczyk, D.; Morris, J.H.; Cook, H.; Kuhn, M.; Wyder, S.; Simonovic, M.; Santos, A.; Doncheva, N.T.; Roth, A.; Bork, P.; Jensen, L.J.; von Mering, C. The STRING database in 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids Res.*, **2017**, *45*(D1), D362-D368.
<http://dx.doi.org/10.1093/nar/gkw937> PMID: 27924014
- [22] Chien, C.T.; Bartel, P.L.; Sternglanz, R.; Fields, S. The two-hybrid system: a method to identify and clone genes for proteins that interact with a protein of interest. *Proc. Natl. Acad. Sci. USA*, **1991**, *88*(21), 9578-9582.
<http://dx.doi.org/10.1073/pnas.88.21.9578> PMID: 1946372
- [23] Phizicky, E.M.; Fields, S. Protein-protein interactions: methods for detection and analysis. *Microbiol. Rev.*, **1995**, *59*(1), 94-123.
<http://dx.doi.org/10.1128/MR.59.1.94-123.1995> PMID: 7708014
- [24] Puig, O.; Caspar, F.; Rigaut, G.; Rutz, B.; Bouveret, E.; Bragado-Nilsson, E.; Wilm, M.; Séraphin, B. The tandem affinity purification (TAP) method: a general procedure of protein complex purification. *Methods*, **2001**, *24*(3), 218-229.
<http://dx.doi.org/10.1006/meth.2001.1183> PMID: 11403571
- [25] Kanehisa, M.; Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.*, **2000**, *28*(1), 27-30.
<http://dx.doi.org/10.1093/nar/28.1.27> PMID: 10592173
- [26] Nei, M.; Gojobori, T. Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Mol. Biol. Evol.*, **1986**, *3*(5), 418-426.
PMID: 3444411
- [27] Zhang, R.; Lin, Y. DEG 5.0, a database of essential genes in both prokaryotes and eukaryotes. *Nucleic Acids Res.*, **2009**, *37*(Database issue)(Suppl. 1), D455-D458.
<http://dx.doi.org/10.1093/nar/gkn858> PMID: 18974178
- [28] Luo, H.; Lin, Y.; Gao, F.; Zhang, C.-T.; Zhang, R. DEG 10, an update of the database of essential genes that includes both protein-coding genes and noncoding genomic elements. *Nucleic Acids Res.*, **2014**, *42*(Database issue), D574-D580.
<http://dx.doi.org/10.1093/nar/gkt1131> PMID: 24243843
- [29] Wright, F. The 'effective number of codons' used in a gene. *Gene*, **1990**, *87*(1), 23-29.
[http://dx.doi.org/10.1016/0378-1119\(90\)90491-9](http://dx.doi.org/10.1016/0378-1119(90)90491-9) PMID: 2110097
- [30] Roman, L.; Tatusov, D.A.; Natale, I.V. Garkavtsev, T.A. Tatusova, U.T. Shankavaram, Rao, B.S.; Kiryutin, B.; Galperin, M.Y.; Fedorova, N.D.; Koonin, E.V. The cog database: New developments in phylogenetic classification of proteins from complete genomes. *Nucleic Acids Res.*, **2001**, *29*(1), 22.
PMID: 11125040
- [31] Michael, Y.; Makarova, K.S.; Wolf, Y.I.; Koonin, E.V. Expanded microbial genome coverage and improved protein family annotation in the cog database. *Nucleic Acids Res.*, **2015**, *43*(D1), D261.
<http://dx.doi.org/10.1093/nar/gku1223>
- [32] Butland, G.; Peregrin-Alvarez, J.M.; Li, J.; Yang, W.; Yang, X.; Canadien, V.; Starostine, A.; Richards, D.; Beattie, B.; Krogan, N.; Davey, M.; Parkinson, J.; Greenblatt, J.; Emili, A. Interaction network containing conserved and essential protein complexes in *Escherichia coli*. *Nature*, **2005**, *433*(7025), 531-537.
<http://dx.doi.org/10.1038/nature03239> PMID: 15690043

- [33] Jin, Y.; Turaev, D.; Weinmaier, T.; Rattei, T.; Makse, H.A. The evolutionary dynamics of protein-protein interaction networks inferred from the reconstruction of ancient networks. *PLoS One*, **2013**, *8*(3), e58134. <http://dx.doi.org/10.1371/journal.pone.0058134> PMID: 23526967
- [34] Rajagopala, S.V.; Sikorski, P.; Kumar, A.; Mosca, R.; Vlasblom, J.; Arnold, R.; Franca-Koh, J.; Pakala, S.B.; Phanse, S.; Ceol, A.; H'ausser, R.; Siszler, G.; Wuchty, S.; Emili, A.; Babu, Mohan.; Aloy, P.; Pieper, R.; Uetz, P. The binary protein-protein interaction landscape of *Escherichia coli*. *Nat. Biotechnol.*, **2014**, *32*, 285-290.
- [35] Hahn, M.W.; Kern, A.D. Comparative genomics of centrality and essentiality in three eukaryotic protein-protein interaction networks. *Mol. Biol. Evol.*, **2005**, *22*(4), 803-806. <http://dx.doi.org/10.1093/molbev/msi072> PMID: 15616139
- [36] Karimpour-Fard, A.; Leach, S.M.; Hunter, L.E.; Gill, R.T. The topology of the bacterial co-conserved protein network and its implications for predicting protein function. *BMC Genomics*, **2008**, *9*(1), 313. <http://dx.doi.org/10.1186/1471-2164-9-313> PMID: 18590549
- [37] Wuchty, S.; Uetz, P. Protein-protein interaction networks of *E. coli* and *S. cerevisiae* are similar. *Sci. Rep.*, **2014**, *4*, 7187. <http://dx.doi.org/10.1038/srep07187> PMID: 25431098
- [38] Annibale, A.; Coolen, A.C.C.; Planell-Morell, N. **2015**.
- [39] Bader, G.D.; Hogue, C.W. An automated method for finding molecular complexes in large protein interaction networks. *BMC Bioinformatics*, **2003**, *4*(2), 2. <http://dx.doi.org/10.1186/1471-2105-4-2> PMID: 12525261
- [40] Balaji, S. *Novak, Antal F.; Flannick, Jason A.; Batzoglou, Serafim.; McAdams, Harley H. Integrated Protein Interaction Networks for 11 Microbes*; Springer Berlin Heidelberg, **2006**, pp. 1-14.
- [41] Rao, V.S.; Srinivas, K.; Sujini, G.N.; Kumar, G.N. Protein-protein interaction detection: methods and analysis. *Int. J. Proteomics*, **2014**, *2014*(147648), 147648. PMID: 24693427
- [42] Jeong, H.; Mason, S.P.; Barabási, A.-L.; Oltvai, Z.N. Lethality and centrality in protein networks. *Nature*, **2001**, *411*(6833), 41-42. <http://dx.doi.org/10.1038/35075138> PMID: 11333967
- [43] George, E. Fox. Origin and evolution of the ribosome. *Cold Spring Harb. Perspect. Biol.*, **2010**, *2*(9), a003483.
- [44] Hwang, Y.-C.; Lin, C.-C.; Chang, J.-Y.; Juan, H.-F.; Huang, H.-C. Predicting essential genes based on network and sequence analysis. *Mol. Biosyst.*, **2009**, *5*, 1672-1678.
- [45] Wei, W.; Ning, L.-W.; Ye, Y.-N.; Guo, F.-B. Geptop: a gene essentiality prediction tool for sequenced bacterial genomes based on orthology and phylogeny. *PLoS One*, **2013**, *8*(8), e72343. <http://dx.doi.org/10.1371/journal.pone.0072343> PMID: 23977285
- [46] Bardini, R.; Di Carlo, S.; Politano, G.; Benso, A. Modeling antibiotic resistance in the microbiota using multi-level petri nets. *BMC Syst. Biol.*, **2018**, *12*(6), 108. <http://dx.doi.org/10.1186/s12918-018-0627-1>
- [47] Terradot, L.; Noiro-Gros, M.-F. Bacterial protein interaction networks: puzzle stones from solved complex structures add to a clearer picture. *Integr. Biol.*, **2011**, *3*, 645-652. <http://dx.doi.org/10.1039/c0ib00023j>
- [48] Zoraghi, R.; Reiner, N.E. Protein interaction networks as starting points to identify novel antimicrobial drug targets. *Curr. Opin. Microbiol.*, **2013**, *16*(5), 566-572. <http://dx.doi.org/10.1016/j.mib.2013.07.010> PMID: 23938265
- [49] Sevimoglu, T.; Arga, K.Y. The role of protein interaction networks in systems biomedicine. *Comput. Struct. Biotechnol. J.*, **2014**, *11*(18), 22-27. <http://dx.doi.org/10.1016/j.csbj.2014.08.008> PMID: 25379140
- [50] Boccaletti, S.; Bianconi, G.; Criado, R.; Del Genio, C.I.; Gómez-Gardeñes, J.; Romance, M.; Sendiña-Nadal, I.; Wang, Z.; Zanin, M. The structure and dynamics of multilayer networks. *Phys. Rep.*, **2014**, *544*(1), 1-122. <http://dx.doi.org/10.1016/j.physrep.2014.07.001> PMID: 32834429