Electronic Supplementary Material (ESI) for Molecular BioSystems. This journal is © The Royal Society of Chemistry 2017

# Text S1

## 1. Algorithm of our proposed pipeline

The instance-feature file (see main article and Table S1), containing 384 essential and 3120 nonessential reaction-gene pairs ( $R_a\_G_b$ ) was given to the following algorithm. It is noteworthy to mention that, our methodology is not specific to only this problem of classification and hence, can be applied for classifying any other kinds of datasets. Further, 1000 randomized balanced datasets (equal number of positive and negative classes) were generated and given to the integrated pipeline.

## Algorithm for choosing best feature combination (Part 1)

The following algorithm was used to choose best features from a total of 64 features:

// BD: Set of 1000 Balanced training Datasets

// **TRF:** Set of Top Ranking Features

// PM: Set of Performance Metrics (auROC)

// perf: Set of best performing metrics

// BD<sub>temp</sub>: Set of Balanced Datasets giving high performance

for i=1: length(BD)

TRF[] = Features ranked (descending) using SVM-RFE

#### for j=1: length(TRF[ ])

PM[i][j] = auROC measured using SMO

end loop j

BFC[i] = best features combination set from PM[i][j] which gives best auROC

for k=1: length(BD)

BFC\_BD[i][k] = performance metric with BFC[i] measured using SMO end loop  $_{k}$ 

end Loop<sub>i</sub>

for m=1: length(BD)

auROC[m] = sum(BFC\_BD[m]) / length(BFC\_BD[m]);

end loop m

Sort (descending) auROC[]

Select Best feature combination set  $(BFC_{best})$  which gives highest average performance (auROC[1])

Algorithm for parameter optimization (Part 2)

Training with BFC<sub>best</sub> (obtained from Part 1) and tuning complexity parameter C

C[] = {0.01, 0.1, 1, 10, 100}; for i = 1 : length(C[]) for j = 1 : length(BD) PM[i][j] = performance metrics measured using SMO; end loop j Sort (descending) PM[i] and choose corresponding BD[j] perf[i] = PM[1]; // best auROC for C[i] BD<sub>temp</sub>[i] = BD[1]; // dataset giving best auROC for C[i] avePERF\_C[i]=sum(PM[i])/length(PM[i]); end loop i Sort (descending) perf[] and choose corresponding BD<sub>temp</sub>[] Sort (descending) perf[] and choose corresponding BD<sub>temp</sub>[] Sort (descending) avePERF\_C[] and choose corresponding C[]

**Output** = Best feature combination, best penalty parameter, best dataset ( $BFC_{best}$ ,  $C_{best}$ ,  $BD_{best}$ ) which gives highest performance

### **Algorithm for Model testing**

The unbalanced instance-feature file again was given as the total master test set. The testing algorithm returns two results -

1) Predictions from the best model (C<sub>best</sub>, BFC<sub>best</sub>, BD<sub>best</sub>)

2) Predictions of the model (Cbest, BFCbest) with respect to the 1000 random datasets -

Given - C<sub>best</sub>, BFC<sub>best</sub>

```
for i=1: length(Ra_Gb)
count = 0;
    for j=1: length(BD)
        if Ra_Gb[i] == "E" // "E" means essential
            count=count+1;
            end if
            end loop j
PercentagePrediction = [count / length(BD)]*100;
End
```

#### 2. Curated features for *E. coli* K-12 MG1655

Sequence-based, gene expression-based, metabolic network and flux coupled subnetwork based features was assembled for each reaction-gene combination within *E. coli* K-12 MG1655 metabolism. Thus, a total of 64 features were obtained for each combination (Table A, Table S1).

*Nucleotide content and coding sequence length* - Previous studies showed that GC content of bacterial genomes were either correlated with environmental niche in which the bacterium survives or as a proxy for horizontally transferred genes.<sup>53</sup> Thus, the underlying GC content within a considered genome can be an appropriate representative of gene essentiality. Frequency of A, T, G, and C nucleotides at the 3<sup>rd</sup> synonymous position of codons and percentage GC content at all the 3 codon positions in a gene and length of each coding sequence (A3, T3, G3, C3, GC1, GC2, GC3, CDSlen) were calculated using an in-house code.

*Codon usage* - Codon usage was previously used as a predictor of protein abundance.<sup>54–56</sup> Abundant highly expressing proteins might have a functional importance in metabolism and hence, can be essential. Codon usage is also strongly associated to GC content within *E. coli* and is a distinctive signature between genomes.<sup>56–58</sup> Codon usage features like Codon Adaptation Index (CAI),<sup>55</sup> and Effective Number of Codons (ENC)<sup>57</sup> were extracted using EMBOSS package version 6.6.0-1.<sup>59</sup> Total Number of codons (Num\_codons) in a coding sequence was calculated using an in-house code.

*Homology based features* - A gene might be more important if it is conserved across evolutionarily related organisms residing in different environments. In bacteria, essential genes were observed to be more evolutionarily conserved as compared to non-essential genes irrespective of the environment.<sup>60,61</sup> Phyletic Retention (PR) is defined as the number of organisms in which ortholog of a given gene is present.<sup>8</sup> For computing PR, protein orthologs amongst 710 bacterial genomes (leaving *E. coli* K-12 substr MG1655) available from the 2014 update of the COG database<sup>36</sup> was searched. An ortholog was defined such that, it is the only bidirectional best hit of the query gene in an organism and possessed atleast 40% identity with the query gene along with an E-value cut-off  $10^{-7}$ . Bi-directional best hits for each gene were identified using BLAST version 2.2.26 along with the above parameters. Further, the number of

homologs in the 710 genomes with respect to the hits obtained with different E-value cut-offs ranging from 10<sup>-3</sup> to 10<sup>-30</sup> (H3, H5, H7, H10, H20, H30) was also calculated <sup>15</sup>.

*Peptide sequence features* - Biased amino acid usage is a property of essential genes in bacterial genomes. The amino acid usage within a protein sequence is largely dependent on the physicochemical properties that an amino acid provides. In *E. coli*, proteins from essential genes might have specific structural properties important for its function.<sup>27</sup> A total of 20 peptide sequence-based features with respect to the frequencies of the 20 amino acids for each protein related to a particular gene (Gly<sub>f</sub>, Met<sub>f</sub>, Ala<sub>f</sub>, Val<sub>f</sub>, Leu<sub>f</sub>, Ile<sub>f</sub>, Pro<sub>f</sub>, Phe<sub>f</sub>, Trp<sub>f</sub>, Tyr<sub>f</sub>, Asn<sub>f</sub>, Arg<sub>f</sub>, His<sub>f</sub>, Glu<sub>f</sub>, Gln<sub>f</sub>, Ser<sub>f</sub>, Thr<sub>f</sub>, Asp<sub>f</sub>, Cys<sub>f</sub>, Lys<sub>f</sub>) was calculated. Similarly, protein length is also an important factor for determining function of a gene in *E. coli*. Protein length and amino acid usage was calculated using EMBOSS package version 6.6.0-1.<sup>59</sup>

Gene expression features - Essential genes tend to express at higher rates as compared to nonessential genes across bacteria.<sup>45</sup> To calculate gene expression based features, 101 microarray experimental samples (from 16 different studies given in Table B) for E. coli K-12 substr. MG1655 that were performed under different environmental stress conditions was collected. The microarray studies were curated from E. coli gene expression database.<sup>62</sup> Microarray studies carried out on mutant strains was not considered, as our aim was to predict essential genes in a wild type strain, subject to an array of environmental conditions. Also, the microarray experiments related to gene expression based features were chosen such that it covers the expression profiles of genes under various environmental stress conditions to get a universal definition of gene irrespective of the environment. Average mRNA Expression of a gene (aveEXP) and mRNA Expression Fluctuation (mEF) which is the standard deviation of log2 normalized gene expression values of the Cy3/Cy5 intensity ratio of each gene from the above samples were calculated. From previous studies, it was reported that a gene might be important if it co-regulated with many other genes.<sup>63</sup> Hence, the Number of Genes with Similar Expression (NGSE), (number of gene pairs having a Pearson correlation coefficient: r < -0.8 and r > 0.8)<sup>15</sup> was also calculated.

Reaction and Flux-coupled sub-network features - The E. coli iJO1366 metabolic network was transformed as a static undirected reaction network representation, in which each node represents

an enzyme (reaction) and the edge signifying the connection between two reaction such that product of previous reaction is the substrate of the next reaction. Different commonly used network features that emphasize biological importance of an enzyme with respect to its position and connectivity to other enzymes in a network were calculated.<sup>64</sup> A highly connected and central enzyme in biological networks is necessarily essential as it represents either a hub or bottleneck within the network.<sup>65</sup> Each network feature would represent the importance of a reaction to interconvert substrates/products given by/to a subsequent reaction. Features like Degree Centrality (DC),<sup>14,15</sup> Closeness Centrality (CC),<sup>14,15</sup> Betweenness Centrality (BC),<sup>14,15</sup> Eccentricity Centrality (EC),<sup>65</sup> Page Rank (PageRank),<sup>66</sup> Eigenvector Centrality (EVC),<sup>15</sup> Modularity (M),<sup>67</sup> Clustering Coefficient (ClustCoef),<sup>14,15</sup> Number of triangles (Num\_triangles),<sup>68</sup> Hub (HS)<sup>69</sup> and Authority (AS)<sup>69</sup> Scores were calculated for the reaction network using Gephi version 0.8.6<sup>70</sup>. The aforementioned network topological features were calculated for the flux-coupled sub-graph as well.

Feature Name	Description	
Sequence Features	A3, T3, G3, C3, GC1, GC2, GC3, CDSlen, CAI, ENC, Num_Codon, Gly <sub>f</sub> , Met <sub>f</sub> , Ala <sub>f</sub> , Val <sub>f</sub> , Leu <sub>f</sub> , Ile <sub>f</sub> , Pro <sub>f</sub> , Phe <sub>f</sub> , Trp <sub>f</sub> , Tyr <sub>f</sub> , Asn <sub>f</sub> , Arg <sub>f</sub> , His <sub>f</sub> , Glu <sub>f</sub> , Gln <sub>f</sub> , Ser <sub>f</sub> , Thr <sub>f</sub> , Asp <sub>f</sub> , Cys <sub>f</sub> , Lys <sub>f</sub> , PL, PR, H3, H5, H7, H10, H20, H30	39
Expression features	NGSE, aveEXP, mEF	3
Network features		
Reaction Network (RN)	RN_DC, RN_EC, RN_CC, RN_ BC, RN_EvC, RN_ HS, RN_AS, RN_PageRank, RN_ClustCoeff, RN_Num_triangles, RN_M	11
Flux Coupled Sub Network (FCA)	FCA_DC, FCA_EC, FCA_CC, FCA_BC, FCA_EvC, FCA_HS, FCA_AS, FCA_PageRank, FCA_ClustCoeff, FCA_Num_triangles, FCA_M	11

Table A. Curated features for *E. coli* K-12 MG1655

*Experiments curated for calculating gene expression-based features* - Gene expression based features like average mRNA Expression (aveEXP), mRNA Expression Fluctuation (mEF) and the gene co-expression feature NGSE were calculated from 16 microarray experiments.

Accession	Conditions	# Samples	References
GSE1730	mRNA abundance from wild type grown in LB or M9	12	71
GSE1981	Gene expression profiles of <i>E.coli</i> grown in LB and minimal medium (DM) at OD 600 =1	3	72
GSE4344	Transcript abundance in LB and L9 at 30 <sup>o</sup> C and OD 600=0.8	4	73
GSE4359	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	6	74
GSE4363	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	11	74
GSE4364	Responses of genes over time as they recover from one stationary phase In rich media at OD 0.4	7	74
GSE4370	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	7	74
GSE4371	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	7	74
GSE4373	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	6	74
GSE4374	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	6	74
GSE4375	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	6	74
GSE4376	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	6	74
GSE4380	Responses of genes over time as they recover from one stationary phase in rich media at OD 0.4	6	74
GSE4706	Response and adaptation to growth with low glucose	6	75

 Table B. Microarray experiments curated for E. coli K-12 substr. MG1655 from E. coli gene

 expression database

	concentration		
GSE4735	Transcriptome dynamic analysis during transition	4	76
	Transcriptome dynamic analysis during transition		
GSE6644	from aerobic and anaerobic conditions	4	76

# 3. Flux coupling analysis

Flux coupling analysis (FCA) is a flux-based optimization procedure that calculates reaction subsets that are either coupled with each other via flux or represent a set of block reactions, given specific environmental exchange constraints.<sup>31,32</sup> Let  $v_1$  and  $v_2$  be fluxes through reactions R<sub>1</sub> and R<sub>2</sub>. Keeping either  $v_1$  or  $v_2$  as objective functions to be optimized, if a non-zero flux in  $v_1$ imposes a non-zero flux in  $v_2$  or vice versa, the two reaction fluxes are termed to be coupled with each other. If zeroing the flux of one reaction does not produce any effect on any other reaction within the metabolic network, then the reaction is termed to be uncoupled. If maximum or minimum of a particular reaction flux objective functions, the coupled reactions can be classified into:

- 1) **Fully coupled:** If  $v_1 = 0$  implies  $v_2 = 0$  and if  $v_2 = 0$  implies  $v_1 = 0$ , and  $v_1 = v_2$ , then the reaction pair is fully coupled.
- 2) <u>Directionally coupled</u>: If  $v_1 = 0$  implies  $v_2 = 0$  but if  $v_2 = 0$  does not imply  $v_1 = 0$ , then the reaction pair is directionally coupled.
- 3) **Partially coupled:** If  $v_1 = 0$  implies  $v_2 = 0$  and if  $v_2 = 0$  implies  $v_1 = 0$ , and  $v_1 \neq v_2$ , then the reaction pair is partially coupled.

Performing FCA on the iJO1366 metabolic network, 1527 fully, 41049 directionally, and 7438 partially coupled reaction pairs and 865 blocked reactions were obtained. As our aim was to find a flux-coupled subnetwork, the nature/property of each reaction pair can be represented within an adjacency matrix (1718 x 1718) where each reaction pair can be given a value of 1 or 0 corresponding to whether they are either coupled or not. Here, we give a value of 0 to both uncoupled and blocked reaction pairs. The adjacency matrix represents a flux-coupled subgraph, which can be used to extract biologically relevant topological features dependent on predicted physiological flux relationships.

#### 4. Definitions of Performance metrics

In our data set, two classes essential (positive class) and non-essential (negative class) genes are considered. The classifier has four outcomes:

True positive (TP): Number of essential instances correctly predicted by the classifier.

False positive (FP): Number of non-essential instances wrongly predicted as essential.

*True negative (TN)*: Number of non-essential instances correctly predicted by the classifier.

False negative (FN): Number of essential instances wrongly predicted as non-essential.

With respect to the above model outcomes, a set of model performance metrics can be calculated (See Table 1 of Main Article for formulae). The metrics are defined as,

*True Positive Rate (TPR) or Sensitivity*: It is defined as the proportion of positive (essential) instances predicted correctly by the model.

*False Positive Rate (FPR)*: It is defined as the proportion of negative (non-essential) instances predicted as positive by the model.

*Precision*: It determines the measure of correctness (i.e., how many instances predicted as essential class really belongs to positive class.

*Recall*: It measures the proportion of essential instances correctly predicted by the model.

*F-measure*: This performance metric is defined as the harmonic mean between precision and recall. A high value of F-measure suggests that the predictive performance better on essential class.

*Matthews Correlation Coefficient (MCC)*: This performance metric was proposed by biochemist Brain W. Matthews in 1975.<sup>77</sup> This measure is less influenced by imbalanced data sets.

*Area under Receiver Operating Characteristic curve (auROC)*: Area which is calculated from ROC curve. Wilcoxon-Mann-Whitney test statistic is used to calculate auROC.

#### 5. Comparison with other available methods – Proof of training set independence

Apart from Hwang et al., 2009,<sup>14</sup> our strategy was also compared with other recent supervised classification studies on essential gene identification.<sup>16,33</sup> To compare the performance of our strategy with these classification methods, training (*Escherichia coli* genes) and test dataset (*Bacillus subtilis* genes) considered in these studies were provided to our methodology for generating a best SVM model and for further testing, respectively. The best model generated from our methodology using the previously available training dataset (consisting of sequence features of *E. coli* genes) was further tested with test dataset (sequence-based features of *B. subtilis* genes) of the available methods.

Testing results in the form of auROC and precision, indicate that the best model generated through our strategy outperforms both the available supervised classification methods (Table C). The achieved sensitivity and specificity from our methodology is the highest suggesting an enhanced model performance, irrespective of the given input training dataset.

Table C. Comparison of our proposed strategy with methods proposed by Song et al.201433 and Deng et al. 201116

Performance metric	Our Method	<b>Song et al. 2014</b> <sup>33</sup>	<b>Deng et al. 2011</b> <sup>16</sup>
auROC	0.966	0.930	0.800
Precision	0.970	0.730	0.540

#### 6. List of selected features

After applying SVM-RFE algorithm for feature selection, combinations of 26 features that give highest model performance were shortlisted. The ranks (contributions to classification) of these features based on SVM-RFE are given in Table D. Also, after performing Wilcoxon rank-sum test, 21 features out of the total 26 have significantly different median values between essential and non-essential genes. *P*-values calculated by comparing the distributions of features between essential and non-essential classes (Fig. 3 of Main Article) are also indicated in Table D.

Feature Name	Rank (SVM-RFE)	<i>P</i> -value
PR	1	2.2E-16
НЗ	2	2.2E-16
RN_CC	3	2.2E-16
Н5	4	2.2E-16
Ala <sub>f</sub>	5	0.000052
Arg <sub>f</sub>	6	0.00996
Н7	7	2.2E-16
CAI	8	0.2447*
FCA_AS	9	2.2E-16
FCA_Num_triangles	10	2.2E-16
FCA_Hub	11	2.2E-16
Leu <sub>f</sub>	12	0.2669*
RN_BC	13	0.000309
GC1	14	2.2E-16
Т3	15	1.35E-10
NGSE	16	0.2399*
Gly <sub>f</sub>	17	2.2E-16
Val <sub>f</sub>	18	0.5293*
Ser <sub>f</sub>	19	0.009535
FCA_ClusteringCoef	20	3.08E-12
Phe <sub>f</sub>	21	2.21E-15
C3	22	1.65E-05
$\operatorname{Glu}_{\mathrm{f}}$	23	0.001545
GC2	24	0.2302*
Asn <sub>f</sub>	25	2.2E-16
FCA_BC	26	2.93E-15

Table D: The 26 selected best features

\* *P*-values insignificant at a threshold of P < 0.05

References cited in this ESI Text S1 are provided in the Reference list of the Main Article.