# Molecular BioSystems

Volume 1 | Number 1 | Jan 2013 | Pages 1–100

**Molecular Biosystems**

www.rsc.org/molecularbiosystems

THE BIOLOGY OF PLAGUE

ROYAL SOCIETY OF CHEMISTRY

Human genes with greater transcript variants are more likely to play functionally important roles such as cellular maintenance and survival.
160x74mm (150 x 150 DPI)

1

# Human genes with a greater number of transcript variants tend to show biological features of housekeeping and essential genes

Jae Yong Ryu,[a] Hyun Uk Kim[abd] and Sang Yup Lee[*abcd]

[a] Metabolic and Biomolecular Engineering National Research Laboratory, Department of Chemical and Biomolecular Engineering (BK21 Plus Program), Center for Systems and Synthetic Biotechnology, Institute for the BioCentury, Korea Advanced Institute of Science and Technology (KAIST), Daejeon 305-701, Republic of Korea.

[b] BioInformatics Research Center, KAIST, Daejeon 305-701, Republic of Korea.

[c] BioProcess Engineering Research Center, KAIST, Daejeon 305-701, Republic of Korea.

[d] The Novo Nordisk Foundation Center for Biosustainability, Technical University of Denmark, Hørsholm, Denmark.

† Electronic supplementary information (ESI) available: Figures S1-S4 and Tables S1-S8.

Corresponding author: Sang Yup Lee. Mailing address: Department of Chemical and Biomolecular Engineering, KAIST, 291 Daehak-ro, Yuseong-gu, Daejeon 305–701, Republic of Korea. Phone: 82-42-350-3930. Fax: 82-42-350-3910.

[*] e-mail: leesy@kaist.ac.kr

1  **Abstract**

2  Alternative splicing is a process observed in gene expression that results in a multi-exon gene

3  to produce multiple mRNA variants which might have different functions and activities.

4  Although physiologically important, many aspects of genes with different number of

5  transcript variants (or splice variants) still remain to be characterized. In this study, we

6  provide bioinformatic evidences that genes with a greater number of transcript variants are

7  more likely to play functionally important roles in cells, compared with those having fewer

8  transcript variants. Among 21,983 human genes, 3,728 genes were found to have a single

9  transcript, and the remaining genes had 2 to 77 transcript variants. The genes with more

10  transcript variants exhibited greater frequencies of acting as housekeeping and essential genes

11  rather than tissue-selective and non-essential genes. They were found to be more conserved

12  among 64 vertebrate species as orthologs, subjected to regulations by transcription factors

13  and microRNAs, and showed hub node-like properties in the human protein-protein

14  interaction network. These findings were also confirmed by metabolic simulations of 60

15  cancer metabolic models. All these results indicate that genes with a greater number of

16  transcript variants play biologically more fundamental roles.

17

1

## Introduction

3 During the expression of a multi-exon gene, alternative splicing results in generation of

4 multiple transcripts.[1, 2] In human, 92-97% of the multi-exon genes undergo alternative

5 splicing.[3] These alternatively spliced variants from a gene can have important implications in

6 mammalian physiology and have been a source of functional diversity of many human genes

7 by providing multiple protein products with alternative functional domains.[4] In particular,

8 correlations between the number of splice variants of a gene and its functional role have been

9 important topics of human genomic studies. In recent years, advent of next-generation

10 sequencing technology such as RNA-Seq with high resolution has facilitated elucidation of

11 functional features of splice variants of genes.[4, 5] RNA-Seq data revealed that alternative

12 splicing events are differentially regulated in human tissues, leading to tissue-specifically

13 coordinated splicing events.[6] Such tissue-specific alternative splicing events allow the same

14 gene to have different combinations of exons (i.e., splice variants) across the tissues, and

15 therefore tissue-specifically generated splice variants can have differentiated protein

16 structures and functions.[7] Importantly, the protein isoforms from the same gene can have

17 different degrees of disorder (i.e., lack of a well-defined three-dimensional structure)

18 depending on the inclusion of tissue-specific exons. Such protein isoforms can rewire the

19 overall protein-protein interaction (PPI) network by interacting with different proteins.

20 Recent studies on generic PPI[8] and tissue-specific PPI networks of human[9] revealed that

21 proteins encoded by genes with a greater number of splice variants tend to have more

22 neighbor nodes and higher centralities in contrast to those encoded by genes with fewer splice

23 variants. The number of neighbor nodes and node centralities are indicators of biologically

24 important functions, and their greater values tend to get greater for the functionally important

3

1    nodes (e.g., proteins).[10] Interestingly, tissue-specific exons, which are often observed in

2    proteins with large values of neighbor nodes and node centralities, also appeared to be more

3    associated with post-translational modifications and evolutionary conservations than

4    constitutive exons.[7] More functional features of splice variants remain to be elucidated

5    through combined experimental (i.e., RNA-Seq analysis) and theoretical studies.

6         In characterizing functional roles of genes, their expression patterns and essentiality

7    are important criteria to consider. Genes can typically be categorized into housekeeping (HK)

8    and tissue-selective (TS) genes depending on their expression patterns[11, 12]; the former being

9    defined as genes expressed across all tissues to maintain cellular functions and the latter

10    being expressed in only certain tissues. Essential (ES) genes are those that are critical to cell

11    growth and survival, whereas non-essential (NE) genes are not. There was a recent study on

12    identifying human essential genes based on the essential orthologs of mouse.[13] Evaluation of

13    the expression patterns in different tissues and essentiality of genes based on the different

14    number of splice variants can be useful in determining their biological importance.

15         In this study, we provide systemic evidences through bioinformatic analyses that

16    genes with a greater number of transcript variants (or splice variants) have a greater chance of

17    playing biologically important roles than those with fewer transcript variants. First, genes

18    were grouped based on the number of their transcript variants in order to identify correlations

19    between the number of their transcript variants and their expression patterns (as HK and TS

20    genes) / essentiality (as ES and NE genes). For the comparative analyses of genes with the

21    different number of transcript variants, a series of analyses were carried out to elucidate the

22    degree of their functional conservations *via* ortholog analysis across genomes of vertebrate

23    species, regulations by transcription factors and microRNAs, and central hub-like network

24    properties in human PPI network. Finally, we used 60 cancer metabolic models for

4

1   essentiality simulation of human metabolic genes upon their knockouts in order to further

2   validate our findings on correlations between the number of transcript variants and gene

3   functions. The present system-wide study provides additional evidences on the biological

4   importance of transcript variants of human genes.

5

## Results and Discussion

6

### Human genes with a greater number of transcript variants play biologically more

8   **important roles**

9   In order to examine the distribution of human genes showing different number of transcript

10  variants, the number of transcript variants for 21,983 human genes was examined (Table S1,

11  ESI†). These genes were downloaded from the Ensembl BioMart (release 78), and only the

12  protein-coding genes (covering both multi- and single-exon) and their transcripts including

13  transcript variants were considered. Meanwhile, we considered all types of transcript variants

14  for a gene, including both that have protein IDs and do not lead to protein products. The

15  reason is that all types of transcript variants have the chance to influence cellular physiology,

16  for instance in the form of microRNA sponge (see Experimental for details). On average,

17  there were 6.95 transcript variants per human gene. Among 21,983 human genes, 3,728 genes

18  were found to have a single transcript, and the remaining genes had 2 to 77 transcript variants

19  (Fig. 1). Overall, 83% of the human genes had 2-28 transcript variants, and the rest 0.01%

20  (219 genes) had 29 or more transcript variants. In order to investigate correlations between

21  the number of transcript variants and functions of human genes, human genes were

22  categorized into a total of 77 groups according to the number of transcript variants. Among

23  them, 60 groups had at least one or more genes, and none of the human genes had the

24  following numbers of transcript variants: 46, 51, 54, 59, 62, 63, 65, 66, 67, 69, 70, 71, 72, 73,

1  74, 75, and 76. Thus, there were 60 groups that were analyzed as below, excluding those 17

2  groups having no genes belonging to.

3       First, these grouped genes having different number of transcript variants were

4  analyzed with respect to the HK, TS, ES and NE gene categories by using gene expression

5  pattern data covering 3,804 HK and 2,293 TS genes[12, 14] and gene essentiality data for 2,472

6  ES and 3,811 NE genes[13] (Table S2, ESI†). Here, HK and ES genes can be considered to be

7  biologically important and fundamental genes, compared with their counterparts (i.e., TS and

8  NE genes). The numbers of transcript variants for HK, TS, ES and NE genes were

9  determined and compared (Fig. 2A). The results show that HK and ES genes tend to have a

10  greater number of transcript variants compared with TS and NE genes, respectively. Also, the

11  average numbers (9.12 and 9.29) of transcript variants for HK and ES genes are greater than

12  the average number (6.95) of transcript variants for all human genes. This observation

13  suggests that genes having important roles (HK and ES) tend to have a greater number of

14  transcript variants compared with their counterparts (6.95 for TS and 7.64 for NE). It should

15  be noted that the lines inside the boxplots in Figure 2A are median values, not averages. In

16  addition, our analysis on the correlation between the number of exons in all the genes

17  considered in this study and the number of their transcript variants revealed that they were

18  not significantly correlated (Fig. S1, ESI†; Pearson correlation coefficient = 0.39 in Fig. S1).

19  This observation manifests that the greater number of transcript variants for the HK and ES

20  genes was caused by various forms of alternative splicing events, not simply by a greater

21  number of exons in their genes. Splice variants can arise from several different mechanisms,

22  including exon skipping, mutual exclusion of exons, alternative 5' donor site, alternative 3'

23  acceptor site, and intron retention.[4]

24       The finding that the HK and ES genes overall generated a greater number of

6

1  transcript variants was further supported by increasing percentages of HK, TS, ES and NE

2  genes in each group as the number of transcript variants increased (Fig. 2B and C). The

3  percentages of TS and NE genes showed somewhat different patterns; they did not increase

4  as a function of the number of transcript variants. Statistical significances for the presence of

5  HK, TS, ES and NE genes in each group with different numbers of transcript variants were

6  calculated with Fisher's exact test, and are available in Table S3, ESI†.

7       In order to confirm that genes having more transcript variants are playing more

8  important roles, we analyzed expression levels of the genes belonging to 60 groups using a

9  recent proteomic study on 32 different human tissues by Uhlen et al.[15]; the percentages of

10  genes in 60 groups that are expressed and appeared in proteome data were calculated (Fig. 3).

11  Also, the percentages of HK, TS, ES, and NE genes in each tissue were calculated (Fig. 3). It

12  was found that genes with a greater number of transcript variants were more ubiquitously

13  expressed in all 32 different human tissues (red region in the heat map in Fig. 3), compared

14  with those having fewer transcript variants (blue and green regions in the heat map in Fig. 3).

15  As expected, the HK genes were ubiquitously expressed in all the 32 tissues, while 33-76%

16  of the TS and NE genes were expressed in the 32 tissues (Fig. 3). Also, greater than 75% of

17  the ES genes were expressed in all the tissues except for bone marrow and skeletal muscle

18  (Fig. 3). In order to clearly show that the tissue-specific expression patterns of the examined

19  genes were not affected by the presence of the HK, TS, ES and NE genes in each group, the

20  gene expression patterns were re-examined by excluding all the HK, TS, ES and NE genes

21  from each group, and the new results appeared to be consistent (Fig. S2, ESI†). These results

22  confirm that expression of those genes having more transcript variants is more demanded in

23  human cell compared with those having fewer transcript variants, which suggests that these

24  genes with more transcript variants are likely play more important functional roles in the cell.

7

1

2    **Analysis of orthologs in gene groups having different number of transcript variants**

3    Next, we examined the number of conserved orthologs across 64 vertebrate species in each

4    gene group (i.e., 60 groups having different number of transcript variants) to examine

5    whether the functional conservation is correlated with the number of transcript variants (or

6    splice variants). The number of orthologs would indicate the level of functional conservation

7    across the examined species.[16] First, the orthologs in all 60 groups having different number of

8    transcript variants were searched against genomes of 64 vertebrate species and counted (Fig.

9    4). Also, the orthologs of the human HK, TS, ES, and NE genes were searched for these

10   genomes. In this analysis, all the protein-coding genes known to be present in the 64

11   vertebrates were obtained from the OrthoDB.[17] Among genes in the vertebrates, human

12   orthologs were selected in order to examine the presence of conserved orthologs.

13           As a result, genes in groups having a large number of transcript variants appeared to

14   be more conserved in 64 vertebrates compared with those having fewer transcript variants

15   (Fig. 4). In particular, human orthologs were highly conserved in the orders such as

16   *Carnivora*, *Cetartiodactyla*, *Glires*, and *Primates*, whereas human orthologs including HK

17   and ES genes were not well conserved in the other orders such as *Ctenosquamata* and

18   *Saurischia* (Fig. 4). High-level conservation of human orthologs in *Carnivora*,

19   *Cetartiodactyla*, *Glires*, and *Primates* could be attributed to their common ancestor (the

20   magnorder *Boreoeutheria*) according to the NCBI Taxonomy database[18]. Furthermore, a

21   vertebrate species sharing orthologs with human to the greatest extent was olive baboon

22   (*Papio anubis*), which appeared to have all genes from 32 groups and 92.3-99.6% genes from

23   the remaining groups conserved in human. In contrast, sea lamprey (*Petromyzon marinus*)

24   had the lowest number of human orthologs, having all genes from 11 groups and 0-83.3%

8

1  from the remaining groups conserved in human. Sea lamprey (*Petromyzon marinus*) was

2  found to be phylogenetically located in the farthest distance from the rest of the vertebrate

3  species.[19] Groups having 42, 52, 55, 56 and 64 transcript variants, which are conserved

4  among all the 64 vertebrates, had 7 genes in total (i.e., *AKT2*, *EEF1D*, *MOK*, *MYB*, *NDRG4*,

5  *RUNX1T1* and *SORBS2*). These genes were associated with fundamentally important

6  functions such as protein kinases (*AKT2* and *MOK*), eukaryotic translation elongation factor

7  (*EEF1D*), transcription factors (*MYB* and *RUNX1T1*), cell cycle progression (*NDRG4*), and

8  adaptor protein for signaling complex (*SORBS2*). The same consistent results were obtained

9  from the ortholog conservation analysis conducted with the same gene groups, but by

10 excluding all the HK, TS, ES and NE genes, in order to confirm that the conservation patterns

11 were affected by the number of transcript variants (or splice variants), not by the HK, TS, ES

12 and NE genes present in each group (Fig. S3, ESI†). Thus, analysis of the conserved

13 orthologs of human genes across the vertebrate species suggested another clue that genes

14 having many transcript variants are playing functionally more important roles (e.g.,

15 conserved functions across vertebrates) due to their greater level of conservation across the

16 examined species. On the contrary, genes with fewer transcript variants might play rather

17 species-specific roles as shown by the lower level of conservation among the examined

18 species for the group with a single transcript.

19

20 **Regulation of genes by transcription factors and microRNAs for the genes having**

21 **different number of transcript variants**

22 We next investigated to what extent genes with different number of transcript variants (or

23 splice variants) are subject to regulations by transcription factors and microRNAs, two

24 important intracellular regulators. Transcription factors activate or repress their target genes

1    by binding to their promoter regions, whereas microRNAs repress target genes by binding to

2    their seed sites (or microRNA-binding sites) in 3' UTR. Both transcription factors and

3    microRNAs can regulate multiple genes.[20]

4         The average numbers of transcription factors and microRNAs that regulate genes in

5    the 60 groups having different number of transcript variants, and the HK, TS, ES, and NE

6    genes were calculated. Data on target genes regulated by transcription factors and

7    microRNAs were obtained from HTRIdb[21] and miRTarBase,[22] respectively; both databases

8    provide information on experimentally validated target genes regulated by transcription

9    factors and microRNAs. Information on all the microRNAs and transcription factors

10   available in the abovementioned databases was used for this analysis in order to grasp overall

11   relationship between the average numbers of regulators and their target genes. A full list of

12   microRNAs, transcription factors and their target genes are available Table S4 (ESI†).

13        For transcription factors, genes with a single transcript appeared to be regulated by

14   1.49 transcription factors on average. The number of transcript variants and transcription

15   factors regulating the corresponding genes showed a positive correlation up to the group with

16   31 transcript variants; the genes in the group with 31 transcript variants were found to be

17   regulated by 3.21 transcription factors per gene on average. Correlations could not be inferred

18   for the groups having greater than 31 transcript variants because of the lack of sufficient

19   number of genes in these groups; less than 0.01% of human genes belong to these groups.

20   Nonetheless, the overall pattern observed was that genes with many transcript variants tend to

21   be subject to regulations by more transcription factors (Fig. 5A).

22        In a similar manner, gene regulations by microRNAs were examined. Genes with a

23   single transcript appeared to be subject to regulations by 0.92 microRNAs on average.

24   Positive correlations between the number of transcript variants and the number of

1    microRNAs regulating the corresponding genes were observed for the gene groups having up

2    to 30 transcript variants; the group with 30 transcript variants showed presence of 2.76

3    regulatory microRNAs per gene (Fig. 5B). Similarly to the transcription factor case, groups

4    having greater than 30 transcript variants could not be considered for inferring correlations

5    due to too few genes in these groups. Interestingly, three genes (*AKT2*, *MOK* and *MYB*) in a

6    group having 42 transcript variants appeared to be regulated by 7.33 transcription factors and

7    7.67 microRNAs on average; this group showed the greatest number of regulators among all

8    the gene groups having different number of transcript variants. In particular, *MYB,* known to

9    be an essential gene crucial in hematopoiesis[23]*,* was found to be regulated by 13 transcription

10   factors and 20 microRNAs.

11        Because transcription factors and microRNAs generate 9.12 and 9.29 transcript

12   variants on average, respectively, the number of transcription factors and microRNAs

13   regulating these gene sets were compared with genes in the group having 9 transcript variants.

14   Interestingly, the HK (regulated by 3.02 transcription factors and 3.23 microRNAs) and ES

15   (regulated by 3.54 transcription factors and 3.56 microRNAs) genes appeared to be more

16   regulated than the group having 9 transcript variants regulated by 2.91 transcription factors

17   and 2.17 microRNAs. This observation suggests that biologically more important genes such

18   as HK and ES genes tend to be subject to more complex regulations.

19        Taken together, these results confirm that functionally important genes such as those

20   with a greater number of transcript variants, and the HK and ES genes are subject to more

21   complex regulations by more transcription factors and microRNAs. Genes with multiple

22   transcript variants are likely to be involved in complex regulations through different promoter

23   binding and polyadenylations by creating alternative 5' and/or 3' exons of the variant

24   structures, and consequently help cells better adapt to environmental and/or genetic

11

1    perturbations.[24]

2

3    **Analysis of genes with different number of transcript variants from a network**

4    **perspective**

5    The above grouped genes (i.e., genes in 60 groups, and the HK, TS, ES, and NE genes) were

6    then analyzed at large-scale protein level by utilizing a human PPI network from the PINA

7    2.0 database.[25] This network consists of a total of 17,109 nodes and 166,776 edges, each

8    representing proteins and their interactions, respectively. In the PPI network, the degree is

9    defined as the number of interacting proteins. We examined average degrees of proteins

10   encoded by genes in the 60 groups having different number of transcript variants (or splice

11   variants), and the HK, TS, ES, and NE genes. It was hypothesized that genes having more

12   splice variants are likely to be central hubs that have a large number of connections with

13   other proteins, and are also related to cellular essentiality.[26] Consistent with the comparative

14   analyses presented above, central hubs were more frequently mapped to proteins encoded by

15   genes with a greater number of transcript variants, and the HK and ES genes (Fig. 6A). For

16   proteins encoded by the HK and ES genes, the average degrees of interactions were 36.10

17   and 49.07, respectively; these values are almost twice the average degree (20.29) of

18   interactions for proteins encoded by genes having 9 transcript variants (i.e., similar to the

19   average number of transcript variants for the HK and ES genes). Interestingly, genes with 14

20   transcript variants showed proteins with average degree of 53.68, which is a value

21   substantially greater than nearby gene groups. This outlier (group with 14 transcript variants)

22   is due to the presence of *UBC* gene encoding ubiquitin which interacts with 9,136 proteins in

23   the PPI network for protein degradation. Network hub nodes are in general known to be

24   essential because of a large number of their connections with other nodes and hence greater

1    damages to the network stability upon their removal.[26] The observation that proteins encoded

2    by the genes having a greater number of transcript variants are more likely to have central

3    hub-like properties is not strange because multiple proteins are generated from such genes,

4    and therefore allow more interactions with other proteins.[10] Consistently to these results, a

5    previous study revealed that the number of degrees of protein nodes in human generic and

6    tissue-specific PPI networks was positively correlated with the number of transcript variants

7    for their respective genes.[8, 9]

8            Finally, the correlation between the average degree of the PPI network and the

9    number of transcript variants was examined for the proteins with similar levels of disorder.

10   Intrinsically disordered proteins are known to interact with more diverse proteins than

11   ordered proteins because of their structural flexibility, and they also have regions enriched for

12   alternative splicing.[27] Therefore, it was important to confirm that the observed average

13   degrees were purely caused by the number of transcript variants (or splice variants), and not

14   the level of protein disorder. For this analysis, disorder levels of all the proteins in the PPI

15   network were calculated using MobiDB 2.0.[28] As a result, the proteins with 0% and 50-70%

16   disorders all consistently showed that their average degrees and the number of their transcript

17   variants were correlated in a positive manner (Fig. 6B and C). The results were also similar

18   for the proteins with > 50% and >70% disorders (Fig. S4, ESI†). Taken together, we found

19   that genes having a greater number of transcript variants indeed followed patterns of the HK

20   and ES genes. This should be useful additional information for better characterization of the

21   human PPI network.

22

23   **Characterizing essentiality of metabolic genes having different number of transcript**

24   **variants using *in silico* genome-scale metabolic models**

13

Comprehensive human genome-scale metabolic models have proven useful in human metabolic studies including the understanding of physiological phenomena[29, 30], prediction of disease-specific biomarkers,[31] and drug targeting.[32, 33] To this end, we used recently reported 60 different NCI-60 cancer cell line-specific metabolic models[34] to further validate that metabolic prediction outcomes are consistent with the observed functional characteristics of genes having different number of transcript variants (or splice variants). Here, cancer cell metabolic models, instead of generic metabolic or normal cell type-specific models, were used in simulations. This is because the objective of cancer cell can be assumed to be biomass maximization, while that of normal cell cannot be.[35]

In order to get the number of metabolic genes in each group having different number of transcript variants, metabolic genes in the human generic model Recon 2 were searched against all the genes in 60 groups. Recon 2 is the latest version of the large-scale human metabolic model that has information on 1,789 metabolic genes, which appear to be present in human genome and correspond to 7,440 reactions and 2,626 metabolites.[36] Metabolic genes were found to have 8.78 transcript variants on average, which is a value greater than the average of all human genes (6.95 transcript variants). The percentage of metabolic genes to the total genes in each group increased as the number of transcript variants increases (Fig. 7A); this observation is reasonable because metabolism plays an important role in cellular growth through energy and biomass generation.

In order to confirm the aforementioned finding that genes with a greater number of transcript variants more frequently followed behaviors shown by the HK and ES genes, we next performed essentiality simulation for the metabolic genes using 60 cancer metabolic models; resulting growth rates of the cancer metabolic models were predicted using constraint-based flux analysis with each gene knocked out individually (see Experimental). In

1    each model, the deleted genes were considered to be essential if the resulting predicted

2    growth rate is lower than 5% of the maximum growth rate. As a result, none of the genes with

3    a single transcript were predicted to be essential, while genes having 2 to 7 transcript variants

4    were increasingly predicted to essential; however, the percentages of essential genes in the

5    groups having 2 to 7 transcript variants were low (Fig. 7B and Table S5, ESI†). Taken

6    together, the results from the simulation of 60 cancer cell metabolic models support our

7    hypothesis that genes having a greater number of transcript variants are more likely

8    associated with cellular essentiality.

9

10    **Conclusions**

11    In this study, we examined functional characteristics of genes according to the number of

12    transcript variants (or splice variants) at genome-scale. It was found that genes having a

13    greater number of transcript variants showed characteristics more similar to those of the HK

14    and ES genes, suggesting that these genes play biologically more important roles. Biological

15    importance of these genes with a greater number of transcript variants was further supported

16    by greater conservation of orthologs across vertebrates, more complex regulations by greater

17    number of transcription factors and microRNAs, and more hub-like properties in the human

18    PPI network compared with genes having fewer transcript variants. Finally, we employed 60

19    cancer genome-scale metabolic models to further examine correlation between the

20    essentiality of genes and the number of transcript variants. Genes having a greater number of

21    transcript variants caused more deleterious effects on cell essentiality upon their knockouts.

22    In summary, several different genome-wide analyses on the genes having different number of

23    transcript variants consistently suggested that those genes having greater number of transcript

24    variants indeed play biologically more important roles, and thus these genes and various

1  transcript variants produced from these genes should receive much more attention in

2  biological studies.

3

4  **Experimental**

5  **Sources of data on human genes used for various comparative analyses**

6  Data on a total of 21,983 protein-coding genes (covering both multi- and single-exon) and

7  their transcripts including splice variants were downloaded from the Ensembl BioMart

8  (release 78).[37] Only the protein-coding genes, not pseudogenes, were considered, but in case

9  of their transcript variants (or splice variants), those given any category of the Transcript

10  Support Level (TSL) were considered because they all have the chance to influence cellular

11  physiology. When only transcripts having the TSL category of tsl1 were counted for the HK,

12  TS, ES and NE genes, it was not possible to observe the differences in the number of their

13  transcript variants; this contrasts with the data presented in Figure 2A. In fact, the percentage

14  of the transcript variants with the tsl1 category was only 27.1% among all the transcript

15  variants theoretically and/or experimentally identified in human genes. Therefore, it was

16  considered reasonable to treat all the transcript variants to more precisely grasp hidden

17  features of transcript variants of the human genes.

18      Information on 3,804 HK and 2,293 TS genes were obtained from Eisenberg et al.

19  (2013)[14] and Chang et al. (2011),[12] respectively. Information on 2,472 ES and 3,811 NE genes

20  was collected from Georgi et al. (2013).[13] As to the analysis of metabolic genes, those defined

21  in the human generic metabolic model Recon 2 were considered.[36] Finally, the orthologs data

22  in 64 vertebrates were obtained from OrthoDB, and among them, only human orthologs were

23  selected.[17] Tissue-specific proteome expression data were obtained from Uhlén et al (2015),[15]

24  which were used to analyze tissue-specific expressions of protein-associated genes in 60

1    groups, and the HK, TS, ES and NE genes. The 32 human tissues were considered in this

2    study.

3

**Analysis of microRNA and transcription factor regulating human genes**

5    Information on target genes regulated by transcription factors and microRNAs was obtained

6    from two experimentally validated databases, HTRIdb[21] and miRTarBase,[22] respectively.

7    Ensembl gene IDs for genes obtained from the Ensembl BioMart were next converted to

8    Entrez gene IDs used in the HTRIdb and miRTarBase using gene2ensembl available at the

9    NCBI FTP Site (Feb. 2015) in order to map the genes onto those regulated by microRNAs

10   and transcription factors. As a result of the gene ID conversion, 17,362 genes were

11   considered for this analysis, and the numbers of microRNAs and transcription factors

12   regulating them were counted. For the statistical significances, one-sided Wilcoxon rank sum

13   tests were performed for each pair of the groups with different number of transcript variants

14   presented in Figure 5 (Table S6 and S7, ESI†).

15

**Analysis of protein-protein interaction network for the genes with a single transcript**

**and multiple transcript variants**

18   Protein interactome data were downloaded from the Protein Interaction Network Analysis

19   (PINA) 2.0 database.[25] This network contains a total of 17,109 nodes and 166,776 edges,

20   each representing proteins and their interactions, respectively. NetworkX version 1.8

21   (http://networkx.lanl.gov/) python package was used to calculate degree distributions of protein

22   nodes. The same statistical procedure used for the target genes regulated by microRNAs and

23   transcription factors was used for this analysis to obtain statistical significances (Fig. 6). For

24   the analysis of correlations between the degree of protein disorder and the number of

17

1    transcript variants, the degree of protein disorder was calculated using a python script

2    available at MobiDB 2.0.[28]

3

4    *In silico* **genome-scale metabolic simulation**

5    Metabolic simulations are typically conducted by using an optimization technique for a

6    metabolic model that has stoichiometric coefficients of all the metabolites in metabolic

7    reactions that appear to be present in an organism.[38] The genome-scale metabolic models are

8    usually underdetermined systems for which the optimization is needed, and the objective

9    function is typically set to maximization of biomass formation for human cancer cells and

10   microorganisms.[35] In contrast to kinetic modeling, this genome-scale metabolic modeling

11   does not require kinetic parameters, but optionally can take omics data which can be set as

12   optimization constraints for a human system.[39] In this study, recently reported 60 NCI-60

13   cancer cell line-specific metabolic models were used for the metabolic simulations.[34] Gene

14   essentiality simulation was conducted using minimization of metabolic adjustment

15   (MOMA).[40] Essential genes were defined as genes causing the cellular growth rate lower than

16   5 % of its maximum value upon their knockouts. All the metabolic simulations were

17   conducted under the COBRApy environment[41] with Gurobi Optimizer (Gurobi Optimization,

18   Inc., Houston, TX).

19

20   **Acknowledgements**

24

1    **Competing financial interests**

2    The authors declare no competing financial interests.

3

4    **References**

5    1.    T. Maniatis, *Science*, 1991, **251**, 33-34.
6    2.    Y. Barash, J. A. Calarco, W. Gao, Q. Pan, X. Wang, O. Shai, B. J. Blencowe and B. J.
7          Frey, *Nature*, 2010, **465**, 53-59.
8    3.    Q. Pan, O. Shai, L. J. Lee, B. J. Frey and B. J. Blencowe, *Nat. Genet.*, 2008, **40**, 1413-
9          1415.
10   4.    H. D. Li, R. Menon, G. S. Omenn and Y. Guan, *Trends Genet.*, 2014, **30**, 340-347.
11   5.    Z. Wang, M. Gerstein and M. Snyder, *Nat. Rev. Genet.*, 2009, **10**, 57-63.
12   6.    E. T. Wang, R. Sandberg, S. Luo, I. Khrebtukova, L. Zhang, C. Mayr, S. F. Kingsmore,
13         G. P. Schroth and C. B. Burge, *Nature*, 2008, **456**, 470-476.
14   7.    M. Buljan, G. Chalancon, S. Eustermann, G. P. Wagner, M. Fuxreiter, A. Bateman and
15         M. M. Babu, *Mol Cell*, 2012, **46**, 871-883.
16   8.    A. Sinha and H. A. Nagarajaram, *J. Proteome. Res.*, 2013, **12**, 1980-1988.
17   9.    A. Sinha and H. A. Nagarajaram, *Proteomics*, 2014, **14**, 2242-2248.
18   10.   C. J. Tsai, B. Ma and R. Nussinov, *Trends Biochem. Sci.*, 2009, **34**, 594-600.
19   11.   A. J. Butte, V. J. Dzau and S. B. Glueck, *Physiol. Genomics*, 2001, **7**, 95-96.
20   12.   C. W. Chang, W. C. Cheng, C. R. Chen, W. Y. Shu, M. L. Tsai, C. L. Huang and I. C.
21         Hsu, *PLoS One*, 2011, **6**, e22859.
22   13.   B. Georgi, B. F. Voight and M. Bucan, *PLoS Genet.*, 2013, **9**, e1003484.
23   14.   E. Eisenberg and E. Y. Levanon, *Trends Genet.*, 2013, **29**, 569-574.
24   15.   M. Uhlen, L. Fagerberg, B. M. Hallstrom, C. Lindskog, P. Oksvold, A. Mardinoglu, A.
25         Sivertsson, C. Kampf, E. Sjostedt, A. Asplund, I. Olsson, K. Edlund, E. Lundberg, S.
26         Navani, C. A. Szigyarto, J. Odeberg, D. Djureinovic, J. O. Takanen, S. Hober, T. Alm,
27         P. H. Edqvist, H. Berling, H. Tegel, J. Mulder, J. Rockberg, P. Nilsson, J. M. Schwenk,
28         M. Hamsten, K. von Feilitzen, M. Forsberg, L. Persson, F. Johansson, M. Zwahlen, G.
29         von Heijne, J. Nielsen and F. Ponten, *Science*, 2015, **347**, 1260419.
30   16.   R. L. Tatusov, E. V. Koonin and D. J. Lipman, *Science*, 1997, **278**, 631-637.
31   17.   E. V. Kriventseva, N. Rahman, O. Espinosa and E. M. Zdobnov, *Nucleic Acids Res.*,
32         2008, **36**, D271-275.
33   18.   S. Federhen, *Nucleic Acids Res.*, 2012, **40**, D136-143.
34   19.   J. J. Smith, S. Kuraku, C. Holt, T. Sauka-Spengler, N. Jiang, M. S. Campbell, M. D.
35         Yandell, T. Manousaki, A. Meyer, O. E. Bloom, J. R. Morgan, J. D. Buxbaum, R.
36         Sachidanandam, C. Sims, A. S. Garruss, M. Cook, R. Krumlauf, L. M. Wiedemann, S.
37         A. Sower, W. A. Decatur, J. A. Hall, C. T. Amemiya, N. R. Saha, K. M. Buckley, J. P.
38         Rast, S. Das, M. Hirano, N. McCurley, P. Guo, N. Rohner, C. J. Tabin, P. Piccinelli, G.
39         Elgar, M. Ruffier, B. L. Aken, S. M. Searle, M. Muffato, M. Pignatelli, J. Herrero, M.
40         Jones, C. T. Brown, Y. W. Chung-Davidson, K. G. Nanlohy, S. V. Libants, C. Y. Yeh,
41         D. W. McCauley, J. A. Langeland, Z. Pancer, B. Fritzsch, P. J. de Jong, B. Zhu, L. L.
42         Fulton, B. Theising, P. Flicek, M. E. Bronner, W. C. Warren, S. W. Clifton, R. K.

1        Wilson and W. Li, *Nat. Genet.*, 2013, **45**, 415-421, 421e411-412.
2    20.    M. S. Ebert and P. A. Sharp, *Cell*, 2012, **149**, 515-524.
3    21.    L. A. Bovolenta, M. L. Acencio and N. Lemke, *BMC Genomics*, 2012, **13**, 405.
4    22.    S. D. Hsu, Y. T. Tseng, S. Shrestha, Y. L. Lin, A. Khaleel, C. H. Chou, C. F. Chu, H. Y.
5           Huang, C. M. Lin, S. Y. Ho, T. Y. Jian, F. M. Lin, T. H. Chang, S. L. Weng, K. W. Liao,
6           I. E. Liao, C. C. Liu and H. D. Huang, *Nucleic Acids Res.*, 2014, **42**, D78-85.
7    23.    M. L. Mucenski, K. McLain, A. B. Kier, S. H. Swerdlow, C. M. Schreiner, T. A.
8           Miller, D. W. Pietryga, W. J. Scott, Jr. and S. S. Potter, *Cell*, 1991, **65**, 677-689.
9    24.    D. D. Licatalosi and R. B. Darnell, *Nat. Rev. Genet.*, 2010, **11**, 75-87.
10   25.    M. J. Cowley, M. Pinese, K. S. Kassahn, N. Waddell, J. V. Pearson, S. M. Grimmond,
11          A. V. Biankin, S. Hautaniemi and J. Wu, *Nucleic Acids Res.*, 2012, **40**, D862-865.
12   26.    H. Jeong, S. P. Mason, A. L. Barabasi and Z. N. Oltvai, *Nature*, 2001, **411**, 41-42.
13   27.    M. Buljan, G. Chalancon, A. K. Dunker, A. Bateman, S. Balaji, M. Fuxreiter and M.
14          M. Babu, *Curr Opin Struct Biol*, 2013, **23**, 443-450.
15   28.    E. Potenza, T. Di Domenico, I. Walsh and S. C. Tosatto, *Nucleic Acids Res.*, 2015, **43**,
16          D315-320.
17   29.    T. Shlomi, T. Benyamini, E. Gottlieb, R. Sharan and E. Ruppin, *PLoS Comput. Biol.*,
18          2011, **7**, e1002018.
19   30.    A. Mardinoglu, R. Agren, C. Kampf, A. Asplund, M. Uhlen and J. Nielsen, *Nat.*
20          *Commun.*, 2014, **5**, 3083.
21   31.    T. Shlomi, M. N. Cabili and E. Ruppin, *Mol. Syst. Biol.*, 2009, **5**, 263.
22   32.    R. Agren, A. Mardinoglu, A. Asplund, C. Kampf, M. Uhlen and J. Nielsen, *Mol. Syst.*
23          *Biol.*, 2014, **10**, 721.
24   33.    K. Yizhak, O. Gabay, H. Cohen and E. Ruppin, *Nat. Commun.*, 2013, **4**, 2632.
25   34.    K. Yizhak, S. E. Le Devedec, V. M. Rogkoti, F. Baenke, V. C. de Boer, C. Frezza, A.
26          Schulze, B. van de Water and E. Ruppin, *Mol. Syst. Biol.*, 2014, **10**, 744.
27   35.    O. Folger, L. Jerby, C. Frezza, E. Gottlieb, E. Ruppin and T. Shlomi, *Mol. Syst. Biol.*,
28          2011, **7**, 501.
29   36.    I. Thiele, N. Swainston, R. M. Fleming, A. Hoppe, S. Sahoo, M. K. Aurich, H.
30          Haraldsdottir, M. L. Mo, O. Rolfsson, M. D. Stobbe, S. G. Thorleifsson, R. Agren, C.
31          Bolling, S. Bordel, A. K. Chavali, P. Dobson, W. B. Dunn, L. Endler, D. Hala, M.
32          Hucka, D. Hull, D. Jameson, N. Jamshidi, J. J. Jonsson, N. Juty, S. Keating, I.
33          Nookaew, N. Le Novere, N. Malys, A. Mazein, J. A. Papin, N. D. Price, E. Selkov, Sr.,
34          M. I. Sigurdsson, E. Simeonidis, N. Sonnenschein, K. Smallbone, A. Sorokin, J. H.
35          van Beek, D. Weichart, I. Goryanin, J. Nielsen, H. V. Westerhoff, D. B. Kell, P.
36          Mendes and B. O. Palsson, *Nat. Biotechnol.*, 2013, **31**, 419-425.
37   37.    R. J. Kinsella, A. Kahari, S. Haider, J. Zamora, G. Proctor, G. Spudich, J. Almeida-
38          King, D. Staines, P. Derwent, A. Kerhornou, P. Kersey and P. Flicek, *Database*
39          *(Oxford)*, 2011, **2011**, bar030.
40   38.    J. D. Orth, I. Thiele and B. O. Palsson, *Nat. Biotechnol.*, 2010, **28**, 245-248.
41   39.    J. Y. Ryu, H. U. Kim and S. Y. Lee, *Integr. Biol.*, 2015, DOI: 10.1039/c5ib00002e.
42   40.    D. Segre, D. Vitkup and G. M. Church, *Proc. Natl. Acad. Sci. USA*, 2002, **99**, 15112-
43          15117.
44   41.    A. Ebrahim, J. A. Lerman, B. O. Palsson and D. R. Hyduke, *BMC Syst. Biol.*, 2013, **7**,
45          74.

46

1

## Figures

2 **Fig. 1** Distribution of genes with respect to the number of their transcript variants. Among

3 21,983 genes obtained from Ensembl BioMart[37], 3,728 genes were found to have a single

4 transcript, while the remaining genes had 2 to 77 transcript variants. There were 6.95

5 transcript variants per human gene on average. The inset shows distribution of 219 genes,

6 each having more than 29 transcript variants, which represents 0.01% of human genes.

7

8

9 **Fig. 2** Correlations between the number of transcript variants of human genes and the

10 functional characteristics (i.e., HK, housekeeping; TS, tissue-selective; ES, essential; and NE,

11 non-essential). (A) Distribution of the number of transcript variants for the HK ($n = 3,804$),

12 TS ($n = 2,293$), ES ($n = 2,472$), and NE genes ($n = 3,811$). Boxes represent the $25^{th}$-$75^{th}$

13 percentiles, while whiskers represent the $5^{th}$-95th percentiles. The line inside the box

14 indicates the median value of the distribution. (B) The percentages of the HK and TS genes,

15 and (C) ES and NE genes among all the genes present in each group according to the number

16 of transcript variants. The $x$-axis is the group name corresponding to the number of transcript

17 variants and the $y$-axis is the percentages of HK, TS, ES, and NE genes in each group.

18 Statistical significances for the presence of HK, TS, ES and NE genes in each group with

19 different number of transcript variants were calculated using Fisher's exact test (Table S3,

20 ESI†).

21

22 **Fig. 3** A heat map showing percentage of the number of expressed genes in each group

23 against each tissue. Tissue-specific expression data were obtained from proteomics studies on

24 32 different human tissues.[15] The percentage represents the number of expressed genes

21

1    among all the genes in each tissue. Tissue names are shown in the *x*-axis, and group names

2    corresponding to the number of transcript variants, and the HK, TS, ES, and NE genes are

3    indicated on the *y*-axis. Abbreviations are: HK, housekeeping; TS, tissue-selective; ES,

4    essential; NE, non-essential.

5

6    **Fig. 4** A heat map showing the percentages of orthologs in gene groups having different

7    transcript variants in 64 vertebrate species. Data on orthologs were obtained from OrthoDB.[17]

8    The percentage represents the number of orthologs among all the genes present in the

9    corresponding gene group and species. The *x*-axis shows the vertebrate species, which were

10   clustered according to their order; the order names are shown only if it has more than 3

11   relevant species. The *y*-axis is the group name corresponding to the number of transcript

12   variants, and the HK, TS, ES, and NE genes. Abbreviations are: HK, housekeeping; TS,

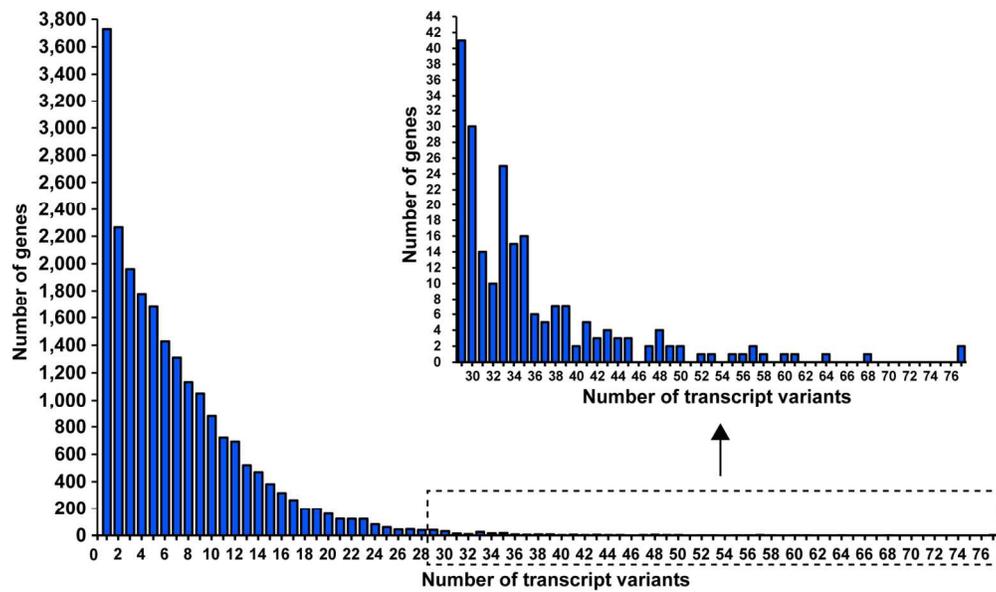13   tissue-selective; ES, essential; NE, non-essential.

14

15   **Fig. 5** Bubble plots representing average numbers of (A) transcription factors and (B)

16   microRNAs regulating genes in each group. Information on target genes regulated by

17   transcription factors and microRNAs was obtained from HTRIdb[21] and miRTarBase,[22]

18   respectively. Average numbers of transcription factors and microRNAs increased for the

19   genes having a greater number of transcript variants. Red and green bubbles represent groups

20   having the HK and ES genes, respectively. Blue bubbles represent 60 groups classified by the

21   number of transcript variants. The TS and NE genes (6.95 and 7.64 transcript variants on

22   average, respectively) appeared to be regulated by 2.36 and 2.87 transcription factors, and

23   1.19 and 1.94 microRNAs, respectively. The bubbles for these genes are not shown because

24   they block those of genes in 60 groups. Bubble size indicates the number of genes in each

1  group. Statistical significances calculated for each pair of the groups using Wilcoxon rank

2  sum test are available in Table S6 and S7, ESI†. Abbreviations are: HK, housekeeping; TS,

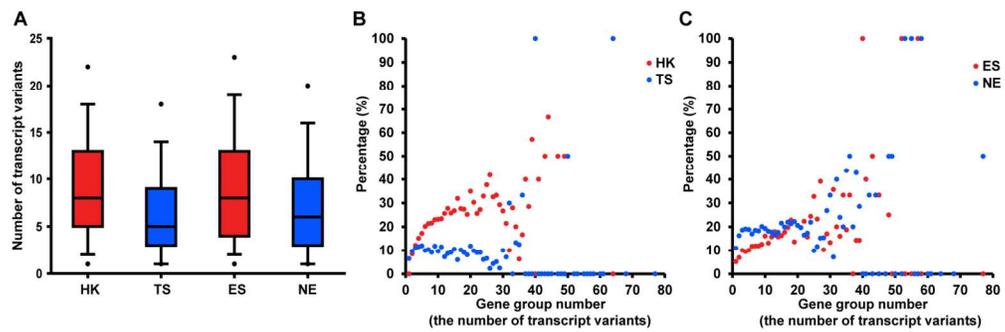3  tissue-selective; ES, essential; NE, non-essential.

4

5  **Fig. 6** Bubble plot showing the average degrees of proteins encoded by human genes in 60

6  groups, and the HK and ES genes in human PPI network. Protein interactome data were

7  downloaded from PINA 2.0 database.[25] This network contains 17,109 nodes and 166,776

8  edges. The three bubble plots were presented for (A) all the proteins, and proteins with (B)

9  0% disorder only and with (C) 50-70% disorder only. Red and green bubbles represent

10  groups having the HK and ES genes, respectively. Blue bubbles represent 60 groups

11  classified by the number of transcript variants. Proteins encoded by the TS and NE genes (on

12  average 6.95 and 7.64 transcript variants, respectively) had average degrees of 13.87 and

13  22.08, respectively. The bubbles for these genes are not shown because they block those of

14  genes in 60 groups. Bubble size indicates the number of genes in each group. Statistical

15  significances calculated for each pair of the groups using Wilcoxon rank sum test are

16  available in Table S8, ESI†. Abbreviations are: HK, housekeeping; TS, tissue-selective; ES,

17  essential; NE, non-essential.

18

19  **Fig. 7** Bubble plots showing (A) the percentage of metabolic genes to the total genes found in

20  each group having different number of transcript variants and (B) the percentage of predicted

21  essential genes for each group in 60 NCI-60 cancer cell line-specific metabolic models. Red

22  and green bubbles represent groups having the HK and ES genes, respectively. Blue bubbles

23  represent 60 groups classified by the number of transcript variants. The HK and ES genes

24  (9.12 and 9.29 transcript variants on average, respectively) had 13.3% and 10.6% metabolic
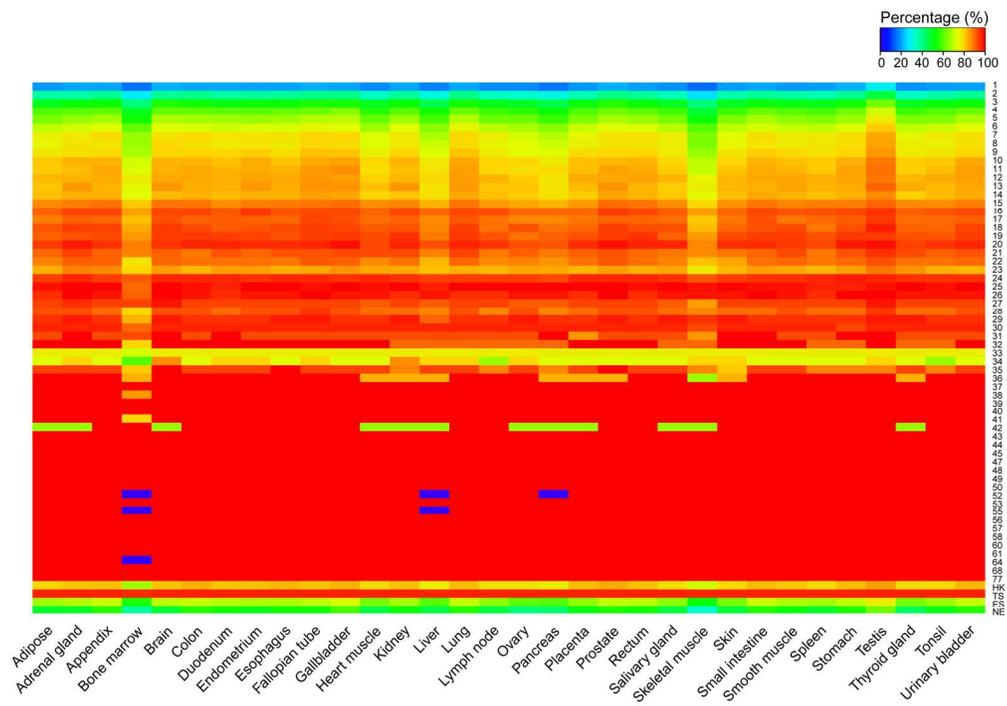
23

1    genes, and 4.5% and 4.4% essential genes, respectively. The TS and NE genes (6.95 and 7.64

2    transcript variants on average, respectively) had 14.8% and 10.9% metabolic genes, and 1.4%

3    and 1.2% essential genes, respectively. The bubbles for the TS and NE genes are not shown

4    because they block those of genes in 60 groups. Bubble size indicates the number of genes in

5    each group. Abbreviations are: HK, housekeeping; TS, tissue-selective; ES, essential; NE,

6    non-essential.
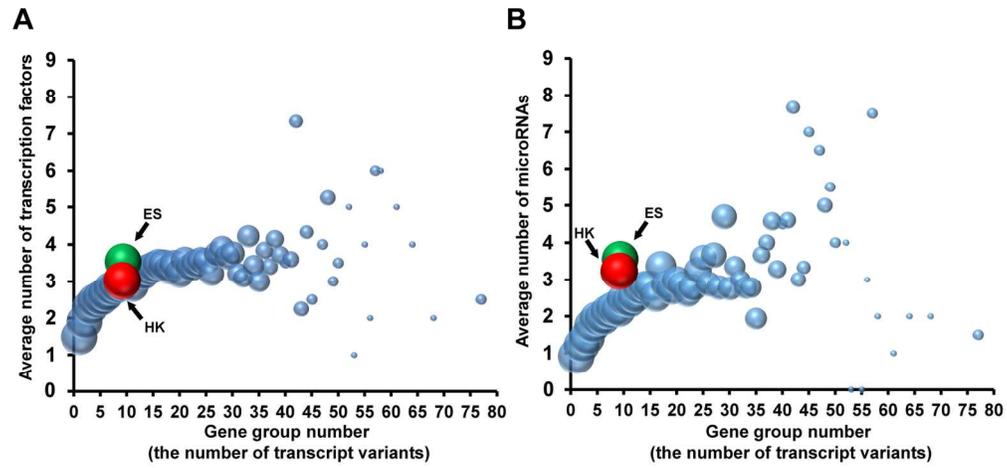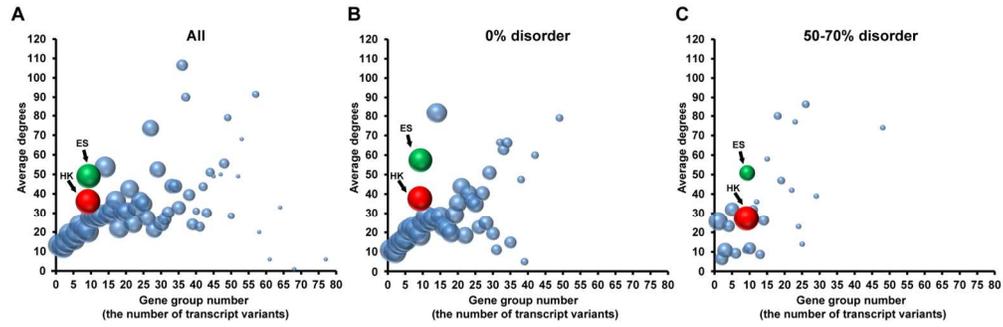
342x209mm (150 x 150 DPI)

343x110mm (150 x 150 DPI)

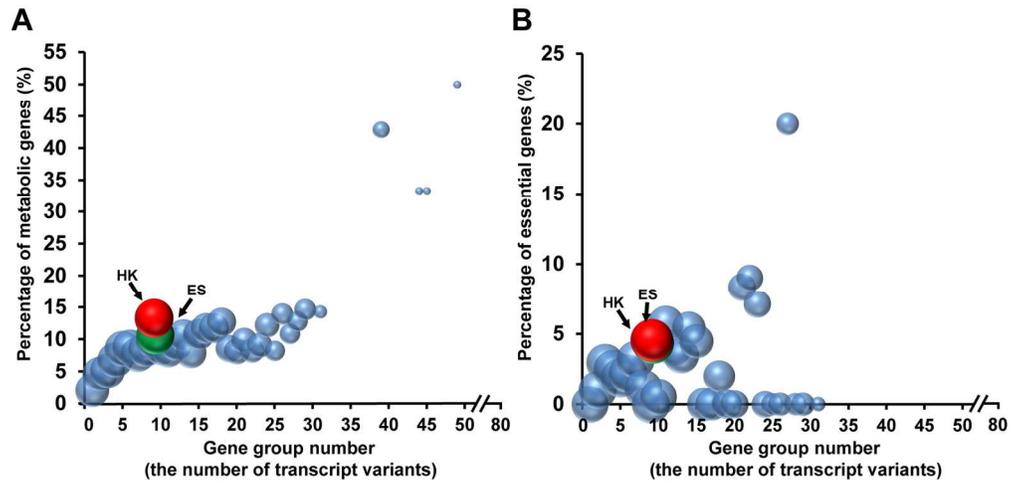360x252mm (150 x 150 DPI)

360x229mm (150 x 150 DPI)

349x161mm (150 x 150 DPI)

355x113mm (150 x 150 DPI)

360x172mm (150 x 150 DPI)