

This is a preprint of an article accepted for publication in *Physiological Genomics*. Yang X., Ye Y., Wang G., Huang H., Liang S. (2011). VeryGene: linking tissue-specific genes to diseases, drugs and beyond for knowledge discovery. *Physiological Genomics*, 43(8):457-460(doi:10.1152/physiolgenomics.00178.2010).

VeryGene: Linking Tissue-Specific Genes to Diseases, Drugs and Beyond for Knowledge Discovery

Xiaoqin Yang¹, Yun Ye^{1,2}, Guiping Wang^{1,3}, Hong Huang⁴, Dekuang Yu⁵, Shuang Liang¹

¹Institute of Genetic Engineering, Southern Medical University, Guangzhou, Guangdong Province, P.R. China

²Department of Biological and Chemical Engineering, Guangxi University of Technology, Liuzhou, Guangxi Province, P.R. China

³Guangzhou Medical College, Guangzhou, Guangdong Province, P.R. China

⁴School of Information, University of South Florida, Tampa, FL, USA

⁵Southern Medical University, Guangzhou, Guangdong Province, P.R. China

Corresponding author: Shuang Liang

Email: itshuang@gmail.com

45 In addition to many other genes, tissue-specific genes (TSGs) represent a set of genes of great
46 importance for human physiology. However, the links among TSGs, diseases and potential
47 therapeutic agents are often missing, hidden or too scattered to find. There is a need to establish a
48 knowledgebase for researchers to share such and additional information in order to speed up
49 discovery and clinical practice. As an initiative toward systems biology, VeryGene web server was
50 developed to fill this gap. A significant effort has been made to integrate TSGs from two
51 large-scale data analyses with respective information on subcellular localization, Gene Ontology,
52 Reactome and KEGG pathway, MGI Mammalian Phenotype, disease association and targeting
53 drugs. The current release carefully selected 3960 annotated TSGs derived from 127 normal
54 human tissues and cell types, including 5672 gene-disease and 2171 drug-target relationships.
55 Other than being a specialized source for TSGs, VeryGene can be used as a discovery tool by
56 generating novel inferences. Some inherently useful but hidden relations among genes, diseases,
57 drugs and other important aspects can be inferred to form testable hypotheses. VeryGene is
58 available online at <http://www.verygene.com>.

59

60 Keywords: tissue specificity; disease; targeting drug

61

62 Human tissues exhibit distinct characteristics in spite of differentiating from a common origin to
63 fulfil the different needs of our body. This kind of diversity is contributed largely by the
64 coordinated expression of different tissue-specific genes, in addition to other genes. The
65 tissue-specific expression pattern of a gene implies not only its physiological function(s), but also
66 where it plays roles in transcriptional regulation, development, stress-response and even disease
67 etiology. Evidences gathered through mining tissue specificity, gene connectivity and disease
68 association suggest that many disease-associated genes are likely to show specific expression in
69 the tissues from which the diseases originate (9, 18). Furthermore, several studies had utilized
70 tissue specificity as an important factor when characterizing therapeutic/drug targets (5, 27). Other
71 areas for the use of tissue-specificity include, but not limited to, pathogenic mechanism, diagnosis,
72 or therapeutic applications (17, 21, 23).

73

74 A number of databases have been created to facilitate studies of TSGs. For example, BioGPS (24),
75 TiGER (11), COXPRESdb (15) and TiSGeD (26) databases can be used to query human gene
76 expression in various tissues. However, most of the above databases focus on the specific
77 expression patterns of TSGs whereas other important biological aspects are not much emphasized.
78 For those who would like to study protein-function, protein-localization, protein-disease or
79 drug-target association altogether, the above databases could not serve to the users' best interest
80 alone. This hinders the practical use of TSGs in medical research and the development of human
81 systems biology. Therefore, a discovery tool dedicated to linking and providing all the above
82 information is highly desirable.

83

84 Herein, we present a web-accessible tissue-specific gene knowledge discovery tool, VeryGene. It
85 is the result of a systematic effort to integrate TSGs surveyed across a large panel of normal
86 human tissues with other important information including subcellular localization, functional
87 annotation, disease/drug relation and so forth. VeryGene serves as a TSG-specific knowledgebase
88 and a discovery tool to generate testable hypotheses for basic and clinical research.

89

90 METHODS

91 Although there are several TSG sets identified from other independent studies, the respective
92 coverage of sample and tissue number is often limited. This makes it harder to conclude whether
93 or not such TSGs are truly expressed in a tissue-selective/specific pattern. In 2004 and 2006, Su *et*
94 *al* (20) and our previous study (10) independently generated a tissue-specific/selective mRNA
95 expression matrix of thousands of genes across a large panel of biological samples (~4000
96 samples combined) and tissue types (~130 tissue types combined) from normal human subjects
97 through microarray expression profiling analysis. Therefore, only these two datasets were selected
98 for integration because of their extended coverage.

99

100 Analytically, searching for tissue-specific genes amounts to comparing gene expression over many
101 tissue types. To determine the tissue distribution for a given gene i across K tissue types, there
102 exists $P = K(K - 1)/2$ pair-wise comparisons for K tissue types. In our previous analyses (10), a
103 modified Tukey-Kramer's honest significant difference (HSD) test with an Enrichment Score
104 (equation 1) was proposed to overcome the type I error from multiple tests. A HSD test generates
105 one Q value (difference of means between tissue pairs over standard deviation) per pair-wise
106 comparison for each gene. An ES ($ES \in (0,1)$) takes into account of Q values from P pairwise
107 comparisons in one HSD test to represent tissue selectivity of a gene. The higher the value of ES
108 for a gene, the more selective it would be. To minimize type I error from multiple HSD tests, each
109 and every TSG $_i$ was identified by observing an ES greater than that by chance alone (estimated by
110 permutation). z_{ij} -scores (equation 2) were calculated to represent the relative level of a given TSG $_i$
111 expressed in one particular tissue j ($j=1$ to K) with regard to the mean expression of TSG $_i$ across
112 all K tissues. The product of z_{ij} and ES_i , denoted as τ_{ij} (equation 3), was computed to account for
113 both tissue specificity/selectivity and relative expression level of a TSG $_i$ in a given tissue j . A large
114 τ_{ij} specifies that a TSG $_i$ is highly specific and significant to a tissue j . In accordance with this
115 quantitative index, not only genes specific to a tissue, but also tissues in which a gene selectively
116 expressed could be ranked (available only in tissue view).

117

$$118 \quad ES_j = \frac{1}{P-1} \sum_{p=1}^P \left(1 - \frac{Q_p - \text{Min}(Q_p)}{\text{Max}(Q_p) - \text{Min}(Q_p)}\right) \quad (1)$$

$$119 \quad z_{ij} = \frac{x_{ij} - \mu_i}{\sigma_i} \quad (2)$$

$$120 \quad \tau_{ij} = z_{ij} * ES_i \quad (3)$$

121

122 Probe IDs were mapped to Entrez Gene IDs. Tissue names were carefully unified according to
123 standard anatomical terms, and redundant tissue affiliations were merged according to the mean
124 value of τ . Finally, 3960 tissue-specific genes were identified through expression profiling of a
125 panel of 127 human tissue and cell types. These TSGs express selectively in ~ 2 tissues on
126 average.

127
 128 To elucidate the functional aspects of these TSGs, detailed annotations were collected. Features of
 129 each specific gene are available at six levels: subcellular localization, Gene Ontology annotation,
 130 biological pathways, mammalian phenotype linkage, disease association and targeting drugs.
 131 Subcellular localization information for these tissue-specific genes was retrieved from LOCATE
 132 (19), supplemented with cellular component annotation of the Gene Ontology (GO) database (6).
 133 Molecular function and Biological process were also obtained from GO. The pathway and reaction
 134 information came from KEGG (8) and Reactome (12), respectively. Mammalian Phenotype
 135 information derived from MGI (Mouse Genome Informatics) (3). Gene-disease relationships were
 136 gathered from Gene2MeSH (1), OMIM information (7) and Swiss-Prot (14). Non-standard disease
 137 names were associated with MeSH IDs and mapped to the MeSH tree categories. Gene-targeting
 138 drug relationships were obtained from DrugBank (22). By integrating these data, 5672
 139 gene-disease relationship and 2171 gene-drug relationship have been collected (Table 1).

140
 141 Table 1. VeryGene data status.

Data Type	Number
Tissue/Cell types	127
TSGs	3960
TSG - Disease relationships	5672
TSG - Drug relationships	2171
TSG - Subcellular Localization	3687
TSG - GO annotation	47418
TSG - Pathway	6359
TSG - Mammalian Phenotype	32397

142
 143 **APPLICATION**
 144 The VeryGene server was implemented in PHP/SQL and is web-accessible through an intuitive
 145 interface. The data contents were configured into two basic views: Gene View and Tissue View to
 146 allow users to conveniently retrieve information relevant to a single gene and tissue/subcellular
 147 localization of interest respectively. Of particular note, Batch View, which evaluates the
 148 enrichment of tissue specificity, subcellular localization, pathway, Gene Ontology, phenotype,
 149 disease and drug for many genes in a single query, is also provided for users to analyze genes of
 150 interest. Batch view is useful to find hidden links and to generate hypotheses. Multiple View is
 151 also developed to allow users to conduct richer combinatorial queries meeting several biological
 152 characteristics simultaneously. This can be used to look up complex relationships and facilitates
 153 discoveries such as "Which proteins of pathway X OR subcellular location Y are tissue-specific"
 154 and so forth. The resultant genes can subsequently be used to perform enrichment analysis with
 155 Batch View. Wildcard search is supported under suitable circumstances. Results from all views, as
 156 well as the entire dataset used to build VeryGene are downloadable for offline use.

157
 158 The TSGs closely related to a specific disease could have hidden links to other biomarkers or
 159 therapeutic targets/agents. VeryGene allows us to identify these unexpected links in order to
 160 generate new hypotheses. In the following example, 8 TSGs for periodontitis (MeSH:D010518)
 161 could be found from Multiple View. Batch View analysis shows that 5 genes among them are also

162 related to rheumatoid arthritis (MeSH:D001172). These 5 TSGs are enriched in such biological
163 processes as immune response and inflammatory response. In addition, they share some common
164 biological pathways, such as cytokine-cytokine receptor interaction and toll-like receptor
165 signalling pathway for the two diseases. Indeed, these findings are consistent with emerging
166 evidence of periodontitis and rheumatoid arthritis sharing many pathological features and
167 biological links (4, 13, 25). Batch View result also suggests that certain TNF inhibitors (e.g.
168 Etanercept and Adalimumab) suitable for one medical condition might be useful for another. For
169 instance, recent studies showed that periodontal therapy using these inhibitors reduced the severity
170 of active rheumatoid arthritis in patients (Figure 1) (16). Besides, it is well known that drug
171 development is time-consuming and very expensive. Finding new indications of existing drugs
172 may help to capitalize the use of such drugs to remedy other medical conditions. Another example
173 presented here is regarding Simvastatin (DrugBank:DB00641), which is a hypolipidemic drug
174 used to control hypercholesterolemia and to prevent cardiovascular disease. A sequential Multiple
175 View/Batch View analysis indicates that 10 TSGs enriched with Simvastatin also enriched with
176 eight most significant MeSH/MIM terms ($p = 0$). Most of these terms can be broadly classified
177 into vascular or inflammatory diseases, among which Endometriosis (MeSH:D004715) (Figure 2)
178 distinguishes itself from others as being suggested to be an autoimmune disease. The potential
179 protective effect of Simvastatin on Endometriosis was preliminarily verified by study in nude
180 mouse model (2). The above examples clearly demonstrate the power of VeryGene to reveal the
181 hidden links some of the earlier databases failed to capture. Many questions such as “How many
182 pathways are enriched in tissue A and what are they? Are they disease-specific? What are the
183 mitochondrial proteins involved in apoptosis in tissue X? Is leukemia linked to any neural disorder?
184 What are the drugs targeting pathway Y?” and so forth can thus be addressed similarly.

185

186 **CONCLUSIONS**

187 We have integrated rich information associated with human TSGs from multiple sources in a
188 web-accessible form to reveal many hidden links beyond tissue-specificity. This makes it a
189 potentially useful source for many applications: for instance, screening for therapeutic targets or
190 biomarkers by tissue, subcellular localization or gene-drug relationship, or looking up for
191 functional enrichment of similarly localized genes or genes participates in a common
192 pathway/disease or vice versa. And most importantly, some hypotheses for pathogenic mechanism,
193 diagnosis and therapeutic researches, could be inferred based on the biological links of TSGs.
194 Much of our effort will be geared toward the understanding of how TSGs play their roles in
195 development, differentiation, stress response and pathology. Study on tissue-specific
196 transcriptional regulation is under way. We also expect to generate many testable hypotheses to
197 maximize VeryGene’s practical value as a knowledge discovery tool.

198

199 Future development of VeryGene will aim to expand and update the extent of current dataset,
200 ensure data quality control, and to enhance user experience as well as data query capability to
201 enable visualization of the complex data relationship. Besides, expression profiles of diseased
202 tissues will also be considered.

203

204 **GRANT**

205 This project is supported by a special grant from the Higher Education Agency of Guangdong

206 Province, P.R. China.

207

208 ACKNOWLEDGMENTS

209 We thank Dr. Changqing Zuo for insightful suggestions in the design and implementation. Our
210 special thanks go to the authors and institutions of the named data sources used in building
211 VeryGene. We are also grateful to the anonymous reviewers' suggestions for enhancing
212 VeryGene's capabilities.

213

214 REFERENCES

- 215 1. **Ade A, Wright Z, and States D.** Gene2MeSH. *Ann Arbor (MI): National Center for*
216 *IntergrativeBiomedical Informatics.* 2007 [The database is available at <http://gene2mesh.ncibi.org/>].
- 217 2. **Bruner-Tran KL, Osteen KG, and Duleba AJ.** Simvastatin Protects against the Development of
218 Endometriosis in a Nude Mouse Model. *J Clin Endocrinol Metab* 94: 2489-2494, 2009.
- 219 3. **Bult CJ, Eppig JT, Kadin JA, Richardson JE, Blake JA, and the Mouse Genome Database**
220 **Group.** The Mouse Genome Database (MGD): mouse biology and model systems. *Nucl Acids Res* 36:
221 D724-728, 2008.
- 222 4. **Demmer RT, Behle JH, Wolf DL, Handfield M, Kebschull M, Celenti R, Pavlidis P, and**
223 **Papapanou PN.** Transcriptomes in Healthy and Diseased Gingival Tissues. *Journal of Periodontology*
224 79: 2112-2124, 2008.
- 225 5. **Dezso Z, Nikolsky Y, Sviridov E, Shi W, Serebriyskaya T, Dosymbekov D, Bugrim A,**
226 **Rakhmatulin E, Brennan R, Guryanov A, Li K, Blake J, Samaha R, and Nikolskaya T.** A
227 comprehensive functional analysis of tissue specificity of human gene expression. *BMC Biology* 6: 49,
228 2008.
- 229 6. **Gene Ontology Consortium.** The Gene Ontology (GO) database and informatics resource.
230 *Nucleic Acids Res* 32: D258-261, 2004.
- 231 7. **Hamosh A, Scott AF, Amberger JS, Bocchini CA, and McKusick VA.** Online Mendelian
232 Inheritance in Man (OMIM), a knowledgebase of human genes and genetic disorders. *Nucleic Acids*
233 *Res* 33: D514-517, 2005.
- 234 8. **Kanehisa M, Goto S, Furumichi M, Tanabe M, and Hirakawa M.** KEGG for representation
235 and analysis of molecular networks involving diseases and drugs. *Nucleic Acids Res* 38: D355-360,
236 2010.
- 237 9. **Lage K, Hansen NT, Karlberg EO, Eklund AC, Roque FS, Donahoe PK, Szallasi Z, Jensen**
238 **TS, and Brunak S.** A large-scale analysis of tissue-specific pathology and gene expression of human
239 disease genes and complexes. *Proceedings of the National Academy of Sciences of the United States of*
240 *America* 105: 20870-20875, 2008.
- 241 10. **Liang S, Li Y, Be X, Howes S, and Liu W.** Detecting and profiling tissue-selective genes.
242 *Physiol Genomics* 26: 158-162, 2006.
- 243 11. **Liu X, Yu X, Zack D, Zhu H, and Qian J.** TiGER: A database for tissue-specific gene
244 expression and regulation. *BMC Bioinformatics* 9: 271, 2008.
- 245 12. **Matthews L, Gopinath G, Gillespie M, Caudy M, Croft D, de Bono B, Garapati P, Hemish J,**
246 **Hermjakob H, Jassal B, Kanapin A, Lewis S, Mahajan S, May B, Schmidt E, Vastrik I, Wu G,**
247 **Birney E, Stein L, and D'Eustachio P.** Reactome knowledgebase of human biological pathways and
248 processes. *Nucleic Acids Res* 37: D619-622, 2009.
- 249 13. **Modi DK, Chopra VS, and Bhau U.** Rheumatoid arthritis and periodontitis: biological links and

- 250 the emergence of dual purpose therapies. *Indian J Dent Res* 20: 86-90, 2009.
- 251 14. **Mottaz A, Yip Y, Ruch P, and Veuthey A-L.** Mapping proteins to disease terminologies: from
252 UniProt to MeSH. *BMC Bioinformatics* 9: S3, 2008.
- 253 15. **Obayashi T, Hayashi S, Shibaoka M, Saeki M, Ohta H, and Kinoshita K.** COXPRESdb: a
254 database of coexpressed gene networks in mammals. *Nucleic Acids Res* 36: D77-82, 2008.
- 255 16. **Ortiz P, Bissada NF, Palomo L, Han YW, Al-Zahrani MS, Panneerselvam A, and Askari A.**
256 Periodontal Therapy Reduces the Severity of Active Rheumatoid Arthritis in Patients Treated With or
257 Without Tumor Necrosis Factor Inhibitors. *Journal of Periodontology* 80: 535-540, 2009.
- 258 17. **Pacak C, Sakai Y, Thattaliyath B, Mah C, and Byrne B.** Tissue specific promoters improve
259 specificity of AAV9 mediated transgene expression following intra-vascular gene delivery in neonatal
260 mice. *Genetic Vaccines and Therapy* 6: 13, 2008.
- 261 18. **Reverter A, Ingham A, and Dalrymple B.** Mining tissue specificity, gene connectivity and
262 disease association to reveal a set of genes that modify the action of disease causing genes. *BioData*
263 *Mining* 1: 8, 2008.
- 264 19. **Sprenger J, Lynn Fink J, Karunaratne S, Hanson K, Hamilton NA, and Teasdale RD.**
265 LOCATE: a mammalian protein subcellular localization database. *Nucleic Acids Res* 36: D230-233,
266 2008.
- 267 20. **Su AI, Wiltshire T, Batalov S, Lapp H, Ching KA, Block D, Zhang J, Soden R, Hayakawa M,**
268 **Kreiman G, Cooke MP, Walker JR, and Hogenesch JB.** A gene atlas of the mouse and human
269 protein-encoding transcriptomes. *Proceedings of the National Academy of Sciences of the United States*
270 *of America* 101: 6062-6067, 2004.
- 271 21. **Vasmataz G, Klee EW, Kube DM, Therneau TM, and Kosari F.** Quantitating tissue specificity
272 of human genes to facilitate biomarker discovery. *Bioinformatics* 23: 1348-1355, 2007.
- 273 22. **Wishart DS, Knox C, Guo AC, Cheng D, Shrivastava S, Tzur D, Gautam B, and Hassanali**
274 **M.** DrugBank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic Acids Res* 36:
275 D901-906, 2008.
- 276 23. **Wu C, Lin J, Hong M, Choudhury Y, Balani P, Leung D, Dang LH, Zhao Y, Zeng J, and**
277 **Wang S.** Combinatorial Control of Suicide Gene Expression by Tissue-specific Promoter and
278 microRNA Regulation for Cancer Therapy. *Molecular Therapy* 17: 2058-2066, 2009.
- 279 24. **Wu C, Orozco C, Boyer J, Leglise M, Goodale J, Batalov S, Hodge C, Haase J, Janes J, Huss**
280 **J, and Su A.** BioGPS: an extensible and customizable portal for querying and organizing gene
281 annotation resources. *Genome Biology* 10: R130, 2009.
- 282 25. **Wu G, Zhu L, Dent JE, and Nardini C.** A Comprehensive Molecular Interaction Map for
283 Rheumatoid Arthritis. *PLoS ONE* 5: e10137, 2010.
- 284 26. **Xiao S-J, Zhang C, Zou Q, and Ji Z-L.** TiSGeD: a database for tissue-specific genes.
285 *Bioinformatics* 26: 1273-1275, 2010.
- 286 27. **Zheng CJ, Han LY, Yap CW, Ji ZL, Cao ZW, and Chen YZ.** Therapeutic targets: Progress of
287 their exploration and investigation of their characteristics. *Pharmacological Reviews* 58: 259-279,
288 2006.

289
290
291

292 Fig. 1. An example of discovering TSG-targeting drugs by linking two different diseases,
293 susceptible TSGs and common targeting drugs.

294 Illustration of a disease-TSG-drug network with Cytoscape. Different nodes are specified
295 according to their categories as follows: disease, dark circle; drug, white square; common TSG,
296 gray circle; other TSG, white circle. Lines are drawn according to their categories as follows:
297 disease-TSG relationship, solid; drug-TSG relationship, dashed.

298

299 Fig. 2. An example of uncovering genes and diseases targeted by Simvastatin.

300 A filtered TSG-disease network visualized with Cytoscape where significant relationships ($p = 0$)
301 are displayed. Nodes are colored and shaped according to their categories as follows: disease, gray
302 circle; drug, dark square; TSG, white circle. Lines are drawn according to their categories as
303 follows: disease-TSG relationship, solid; drug-TSG relationship, dashed. Simvastatin-targeting
304 TSGs in the subset are related to Endometriosis and several other diseases.

305



