

cg05883128 (DDX60) and NR3C2 (CD4 T-cells and B-cells) to Be the Genetic Roots of Systemic Lupus Erythematosus

Zhengjun Zhang

University of Chinese Academy of Sciences and University of Wisconsin

Abstract

Motivation: Systemic lupus erythematosus (SLE) is an autoimmune disease and a long-term condition affecting many body parts. Autoimmune diseases are affecting more people for reasons unknown, and the causes of these diseases remain a mystery.

Method: A newly introduced robust competing (risk) max-logistic regression classifier that can simultaneously perform subtype clustering and classification for disease diagnoses and predictions becomes a new hope to solve the mystery. We use this method in the study to discover critical DNA methylation CpG sites and genome genes, which lead to the highest accuracy and interpretability.

Results: The DNA methylation CpG site cg05883128 (DDX60) and gene NR3C2 are essentially responsible for SLE development. They can lead to 100% prediction accuracy together with a miniature set of other CpG sites and genes, respectively. cg05883128 (DDX60) reveals the LSE mechanism affecting many body parts. NR3C2 in CD4 T cells and B cells behaves reversely, leading to the cause of LSE and explaining the mechanism of the autoimmune disease.

Conclusions: This work represents a pioneering effort and intellectual discovery in applying the max-logistic competing risk factor model to identify critical genes for LSE, and the interpretability and reproducibility of the results across diverse populations suggest that the CpGs and DEGs identified can provide a comprehensive description of the transcriptomic features of SLE. The practical implications of this research include the potential for personalized risk assessment, precision diagnosis, and tailored treatment plans for patients.

Keywords: Site-site interaction effects; gene-gene interaction effects; Simpson's paradox; autoimmune disease.

33 **1 Introduction**

34 Systemic lupus erythematosus (SLE) is a long-term condition affecting many body parts —joints,
35 skin, kidneys, blood cells, brain, heart, and lungs. It occurs when a person's body's immune
36 system attacks her tissues and organs (autoimmune disease). The global SLE newly diagnosed
37 population is estimated to be around 400,000 annually, and most of these people are women [1].
38 In the one year in review 2023 [2], the authors reviewed 93 published papers in 2022 and
39 concluded new results and data have improved the understanding of SLE, although further
40 studies and research are needed to improve the knowledge we have on this complex disease.
41 Upon reviewing 120 published articles, the authors in [3] concluded that more robust
42 immunological biomarkers are needed to better understand disease progression in individuals
43 with SLE, including non-organ-specific SLE biomarkers and organ-specific SLE biomarkers. Since
44 no single biomarker can be sensitive and specific enough for SLE, multiple biomarkers combined
45 through mathematical models may be a good idea for assessing SLE. Moreover, advanced
46 computational methods are required to analyze large datasets and discover novel biomarkers.
47 This paper is not going to further conduct a literature review, and the readers are referred to two
48 excellent review papers [2,3]. This work aims to fill in the gap discussed in [3].

49 The significant contributions of this paper are four-fold. 1) a miniature set of genes is identified
50 at the genomic level to lead to 100% accuracy. This set of genes is mathematically and biologically
51 interpretable via different cohort studies and ethnicities. 2) a miniature set of CpGs is identified
52 at the DNA methylation level to lead to 100% accuracy. This set of CpGs shows universality and
53 why SLE is so complex at methylation levels. 3) The DNA methylation CpG site cg05883128
54 (DDX60) and gene NR3C2 are essential and responsible for SLE development. 4) This paper finds
55 NR3C2 in CD4 T-cells and B-cells behaves reversely compared to NR3C2 in whole blood and CD4+,
56 which causes a Simpson's paradox and the genetic discovery of the long puzzled autoimmune
57 disease (LSE). 5) The Simpson's paradox further leads to a new autoimmune disease theory and
58 sheds light on many other new biological studies and discoveries.

59

60 **2 Method**

61 We apply the newly proven max-linear competing logistic regression classifier method to the
62 confirmed LSE and healthy controls classifications. The new method is very different from other
63 classical statistical and modern machine learning methods, e.g., random forest, deep learning,
64 and support vector machine [4]. In addition, the new method has enhanced the interpretability
65 of results, consistency, and robustness, as shown in the literature in studies of COVID-19 and
66 biomarkers of several types of cancers [4-7]. This section briefly introduces the necessary
67 notations and formulas for self-containing due to the different data structures used in this work.
68 This new innovative approach can be classified as either an AI or machine learning algorithm.
69 However, this new approach has an explicit formula and is interpretable.

70 Suppose Y_i is the i th individual patient's LSE status ($Y_i = 0$ for LSE-free, $Y_i = 1$ for confirmed)
 71 and $X_i^{(k)} = (X_{i1}^{(k)}, X_{i2}^{(k)}, \dots, X_{ip}^{(k)})$, $k = 1, \dots, K$ are the CpG beta values or gene expression values,
 72 with $p = \#$ of CpG sites or genes in this study. Here, k stands for the k th type of beta
 73 (expression) values drawn based on K different biological sampling methodologies. Note that
 74 most published works set $K = 1$, and hence the superscript (k) can be dropped from the
 75 predictors. In this research paper, $K = 1$, for DNA methylation data, $K = 3$, for gene expression
 76 data, and as we have three datasets analyzed in Section 3. Using a logit link (or any monotone
 77 link functions), we can model the risk probability $p_i^{(k)}$ of the i th person's infection status as:

$$\log\left(\frac{p_i^{(k)}}{1 - p_i^{(k)}}\right) = \beta_0^{(k)} + X_i^{(k)}\beta^{(k)} \quad (1)$$

78 or alternatively, we write

$$p_i^{(k)} = \frac{\exp(\beta_0^{(k)} + X_i^{(k)}\beta^{(k)})}{1 + \exp(\beta_0^{(k)} + X_i^{(k)}\beta^{(k)})}$$

79 where $\beta_0^{(k)}$ is an intercept, $X_i^{(k)}$ is a $1 \times p$ observed vector, and $\beta^{(k)}$ is a $p \times 1$ coefficient vector
 80 which characterizes the contribution of each predictor (a CpG site or a gene, in this study) to the
 81 risk.

82 Considering the complexity of LSE, it is natural to assume that the epigenetic structures have
 83 subtypes and subtypes can be different. Suppose that all subtypes of LSE may be related to G
 84 groups of CpG sites (genes):

$$\Phi_{ij}^{(k)} = (X_{i,j_1}^{(k)}, X_{i,j_2}^{(k)}, \dots, X_{i,j_{g_j}}^{(k)}), j = 1, \dots, G, g_j \geq 0, k = 1, \dots, K \quad (2)$$

85 where i is the i th individual in the sample, and g_j is the number of CpG sites (genes) in j th group.
 86 The competing (risk) factor classifier is defined as:

$$\log\left(\frac{p_i^{(k)}}{1 - p_i^{(k)}}\right) = \max(\beta_{01}^{(k)} + \Phi_{i1}^{(k)}\beta_1^{(k)}, \beta_{02}^{(k)} + \Phi_{i2}^{(k)}\beta_2^{(k)}, \dots, \beta_{0G}^{(k)} + \Phi_{iG}^{(k)}\beta_G^{(k)}) \quad (3)$$

87 where $\beta_{0j}^{(k)}$ s are intercepts, $\Phi_{ij}^{(k)}$ is a $1 \times g_j$ observed vector, and $\beta_j^{(k)}$ is a $g_j \times 1$ coefficient
 88 vector which characterizes the contribution of each predictor in the j group to the risk.

89 **Remark 1.** In (3), $p_i^{(k)}$ is mainly related to the largest component $CF_j = \beta_{0j}^{(k)} + \Phi_{ij}^{(k)}\beta_j^{(k)}$, $j =$
 90 $1, \dots, G$, i.e., all components compete to take the most significant effect.

91 **Remark 2.** Taking $\beta_{0j}^{(k)} = -\infty$, $j = 2, \dots, G$, (3) is reduced to the classical logistic regression, i.e.,
 92 the classical logistic regression is a special case of the new classifier. Compared with black-box
 93 machine learning methods (e.g., random forest, deep learning (convolutional) neural networks
 94 (DNN, CNN)) and regression tree methods, each competing risk factor in (3) forms a clear, explicit,
 95 and interpretable signature with the selected CpG sites (genes). The number of factors
 96 corresponds to the number of signatures, i.e., G . This model can be a bridge between linear
 97 models and more advanced machine learning methods (black box) models. However, (3) retains

98 *interpretability, computability, predictability, and stability properties. Note that this remark is*
 99 *similar to Remark 1 in Zhang (2021) [6].*

100 We have to choose a threshold probability value to decide a patient's class label in practice.
 101 Following the general trend in the literature, we set the threshold to be 0.5. As such, if $p_i^{(k)} \leq$
 102 0.5, the i th individual is classified as being disease-free; otherwise, the individual is classified as
 103 having the disease.
 104 Zhang (2021) [6] introduced a new machine learning classifier, the smallest subset and smallest
 105 number of signatures (S_4), for $K = 1$. We extended the S_4 classifier from $K = 1$ to $K = 3$ as
 106 follows:

$$\begin{aligned} (\hat{\beta}, \hat{S}, \hat{G}) = \operatorname{argmin}_{\beta, S_j \subset S, j=1,2,\dots,G} \{ & (1 + \lambda_1 + |S_u|)^{\sum_{k=1}^K} \sum_{i=1}^n (I(p_i^{(k)} \leq 0.5)I(Y_i=1) + I(p_i^{(k)} > 0.5)I(Y_i=0)) \\ & + \lambda_2 (|S_u| - \frac{|S_u| + G - 1}{(|S_u| + 1) \times G - 1}) \} \end{aligned} \quad (4)$$

107 where $I(\cdot)$ is an indicative function, $p_i^{(k)}$ is defined in Equation (3), $S = \{1, 2, \dots, n\}$ is the index
 108 set of all CpG sites (genes), $S_j = \{j_{j1}, \dots, j_{j,g_j}\}$, $j = 1, \dots, G$ are index sets corresponding to (2), S_u
 109 is the union of $\{S_j, j = 1, \dots, G\}$, $|S_u|$ is the number of elements in S_u , $\lambda_1 \geq 0$ and $\lambda_2 \geq 0$ are
 110 penalty parameters, and $\hat{S} = \{j_{j1}, \dots, j_{j,g_j}, j = 1, \dots, G\}$ and \hat{G} are the final CpG set (gene)
 111 selected in the final classifiers and the number of final signatures.

112
 113 **Remark 3.** When the S_4 classifier leads to 100% accuracy, the bioequivalence and DNA
 114 methylation (genome) geometry space can be established, which is a unique property established
 115 in (4) that does not appear in other classifiers in the literature [5].

116 **Remark 4.** *The case of $K = 1$ corresponds to the classifier introduced in Zhang (2021) [6]. The*
 117 *case of $K = 1$ and $\lambda_2 = 0$ corresponds to the classifier introduced in Zhang (2021) [4].*

118 **Remark 5.** *The computational details are referred to [7].*

119 3 Data Descriptions, Results, and Interpretations

120 3.1 The data

121 We apply the optimization equation (4) to eight gene expression datasets (DEGs): GSE81622 [8]
 122 had Platform GPL10558: Illumina HumanHT-12 V4.0 expression beadchip. In this dataset, PBMC
 123 of 30 SLE patients, including 15 with LN (SLE LN+) and 15 without LN (SLE LN-), and 25 normal
 124 controls (NC) using HumanHT-12 Beadchips and Illumina Human Methy450 chips, were
 125 performed whole genome transcription. GSE97263 [9] had Platform GPL16791: Illumina HiSeq
 126 2500 (Homo sapiens), mRNA extraction from isolated blood CD4+ T cells from 14 SLE with
 127 active disease, 16 SLE with less active disease and 14 controls. GSE177029 (two datasets FPKM
 128 and count) [10,11], with Platforms (1) GPL24676: Illumina NovaSeq 6000 (Homo sapiens),
 129 examined 8 placentas from systemic lupus erythematosus and 8 from normal full-term
 130 pregnancies. GSE50722 [12] had Platform: GPL570 [HG-U133_Plus_2] Affymetrix Human
 131 Genome U133 Plus 2.0 Array with 61 SLE patients and 20 healthy controls. GSE144390 [13] had

132 GPL6244 [HuGene-1_0-st] Affymetrix Human Gene 1.0 ST Array [transcript (gene) version]
 133 with three SLEs and three healthy controls. We included one DNA methylation dataset (CpGs):
 134 GSE82218 [8] with Platform PL13534 Illumina HumanMethylation450 BeadChip
 135 (HumanMethylation450_15017482). GSE4588 [14] (two datasets CD4 T cells and B cells sorted
 136 from total PBMC) has Platform: GPL570 [HG-U133_Plus_2] Affymetrix Human Genome
 137 U133 Plus 2.0 Array. CD4 T and B cells were sorted by flow cytometry from PBMC of patients
 138 with SLE, RA and healthy controls. GeneChip® Human genome U133 Plus 2.0 arrays were
 139 hybridized in monoplicates and the genes differentially expressed among the three groups of
 140 patients were identified using ANOVA tests with corrections for multiple comparisons. There
 141 are seven SLEs, nine healthy controls with B cells sorted, and nine SLEs and ten healthy controls
 142 with T cells sorted.

143 3.2 DEGs analysis

144 Following the computational procedure described in the earlier work [4-7], from 47322, 58037,
 145 127550, 19943, 33298, and 54675 genes or gene IDs, six genes were identified to lead to 100%
 146 accuracy in six datasets GSE81622, GSE97263, GSE 177029, GSE50722, GSE144390, GSE4588.
 147 Table 1 lists the classifiers and the coefficients of the sites.

148 Table 1. Performance of individual classifiers and combined max-competing classifiers using blood
 149 sampled data GSE81622, GSE97263, GSE177029, GSE50772, GSE144390, GSE4588 to classify SLE
 150 patients and non-SLE patients into their respective groups. CF1, 2 are two different classifiers.
 151 The numbers are fitted coefficient values.

Classifier	Intercept	NR3C2	SNX20	RAB11FIP5	ARL4C	EPHB2	BCCIP	Accuracy	Sensitivity	Specificity
GSE81262 PMBCs										
CF1	-11.6719				-0.8875	2.6136		78.18%	60%	100%
CF2	6.8769			-7.4137			3.8740	90.91%	83.33%	100%
CFmax								100%	100%	100%
GSE97263 CD4+										
CF1	-3.9346	-0.0061	0.0035					96.43%	92.86%	100%
CF2	4.4926	-0.0097			0.0004			92.86%	85.71%	100%
CFmax								100%	100%	100%
GSE177029 FPKM whole blood										
CF1	3.5101			-2.5039		7.8064		68.75%	37.50%	100%
CF2	23.0284		-14.262		-0.8878			87.50%	75.00%	100%
CFmax								100%	100%	100%
GSE177029 Count whole blood										
CF1	2.0421	-0.0604				0.0373		81.25%	62.50%	100%
CF2	6.2307	-0.0302			-0.0108			87.50%	75.00%	100%
CFmax								93.8%	87.5%	100%
GSE50772 whole blood										
CF1	-16.745		-0.0311	0.0045	0.0064		0.0087	81.48%	77.05%	95%
CF2	11.8033	-0.0101	-0.0409	-0.0014			-0.0072	83.95%	81.97%	90%
CFmax								95.06%	96.72%	90%
GSE144390 whole blood										
CF1	-62.7958						12.3162	100%	100%	100%
GSE4588 CD4 B cells sorted from total PBMC										

CF1	5.3068			-8.5092	-8.0866		6.7556*	81.25%	57.14%	100%	
CF2	20.5721	6.8079			2.1477		-10.062	81.25%	57.14%	100%	
CFmax								100%	100%	100%	
GSE4588 CD4 T cells sorted from total PBMC											
CF1	18.3908		0.0162				-0.0033*	-0.0429	88.89%	75%	100%
CF2	3.369	0.019		0.0729			-0.0038		83.33%	62.50%	100%
CFmax								100%	100%	100%	

152 *BCCIP in GSE4588 B cells has two different probe set IDs: 227896_at (CF1) and 218264_at (CF2). ARL4C in GSE4588 T cells also
 153 has two different probe set IDs: 202206_at (CF1) and 202207_at (CF2).

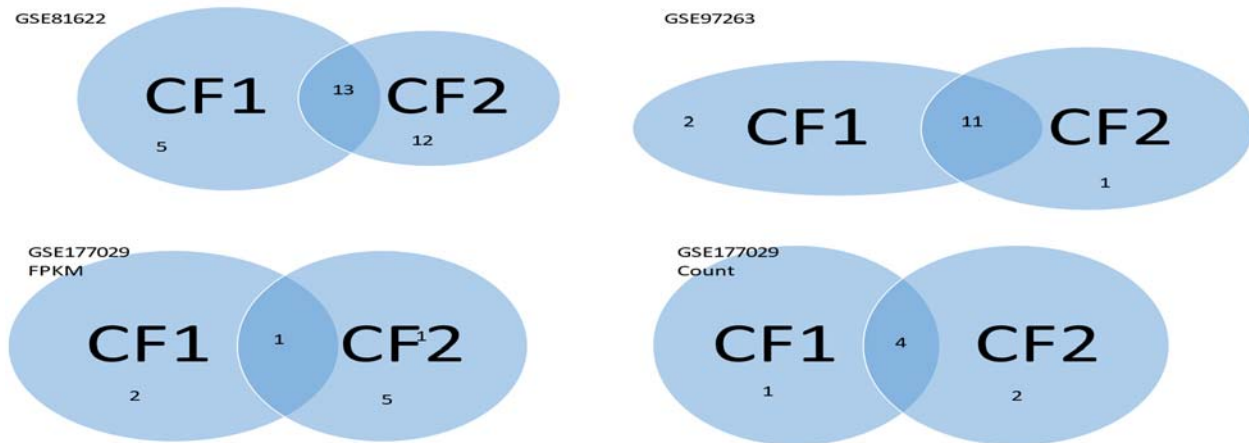
154 In the table, each CFi is expressed as (using GSE81262 CF2 as an example):

155
$$6.8769 - 7.4137 * RAB11F1P5 + 3.874 * BCCIP$$

156 The six genes in genecords.org have the following descriptions. NR3C2 (Nuclear Receptor
 157 Subfamily 3 Group C Member 2) is a Protein Coding gene. Diseases associated with NR3C2 include
 158 Pseudohypoaldosteronism, Type I, Autosomal Dominant and Hypertension, Early-Onset,
 159 Autosomal Dominant, With Severe Exacerbation In Pregnancy. Among its related pathways are
 160 Gene expression (Transcription) and ACE Inhibitor Pathway, Pharmacodynamics. SNX20 (Sorting
 161 Nexin 20) is a Protein Coding gene. Diseases associated with SNX20 include Immunodeficiency
 162 73A With Defective Neutrophil Chemotaxis And Leukocytosis and Inflammatory Bowel
 163 Disease. RAB11FIP5 (RAB11 Family Interacting Protein 5) is a Protein Coding gene. Diseases
 164 associated with RAB11FIP5 include Neonatal Lupus Erythematosus and Mucocutaneous
 165 Leishmaniasis. Among its related pathways are wtCFTR and delta508-CFTR traffic / Generic
 166 schema (norm and CF). Gene Ontology (GO) annotations related to this gene include small
 167 GTPase binding and gamma-tubulin binding. An important paralog of this gene is RAB11FIP1.
 168 ARL4C (ADP Ribosylation Factor Like GTPase 4C) is a Protein Coding gene. Diseases associated
 169 with ARL4C include Vulvar Leiomyosarcoma. Among its related pathways are NR1H2 and NR1H3-
 170 mediated signaling and ESR-mediated signaling. Gene Ontology (GO) annotations related to this
 171 gene include GTP binding and alpha-tubulin binding. An important paralog of this gene is ARL4A.
 172 EPHB2 (EPH Receptor B2) is a Protein Coding gene. Diseases associated with EPHB2 include
 173 Bleeding Disorder, Platelet-Type, 22 and Prostate Cancer/Brain Cancer Susceptibility. Among its
 174 related pathways are GPCR Pathway and EPH-Ephrin signaling. Gene Ontology (GO) annotations
 175 related to this gene include transferase activity, transferring phosphorus-containing groups and
 176 protein tyrosine kinase activity. An important paralog of this gene is EPHB1. BCCIP (BRCA2 And
 177 CDKN1A Interacting Protein) is a Protein Coding gene. Gene Ontology (GO) annotations related
 178 to this gene include RNA binding and kinase regulator activity.

179 Notably, except GSE50772 and GSE4588, all coefficients associated with NR3C2 and RAB11FIP5
 180 are negative, which means increasing their expression values will decrease the risk of SLE
 181 symptoms. On the other hand, all coefficients associated with EPHB2 and BCCIP are positive,
 182 which means decreasing their expression values will decrease the risk of SLE symptoms. Two
 183 other genes do not have uniform coefficient signs, which complicates the SLE symptoms,
 184 leading to SLE subtypes, i.e., CF1 and CF2 can subdivide the patients into three subtypes in each
 185 dataset. Note that the three datasets were collected from UK (GSE97263), and China (two

186 others). The subtypes can be different from different geographical regions. We can have a Venn
 187 diagram, as shown in Figure 1.



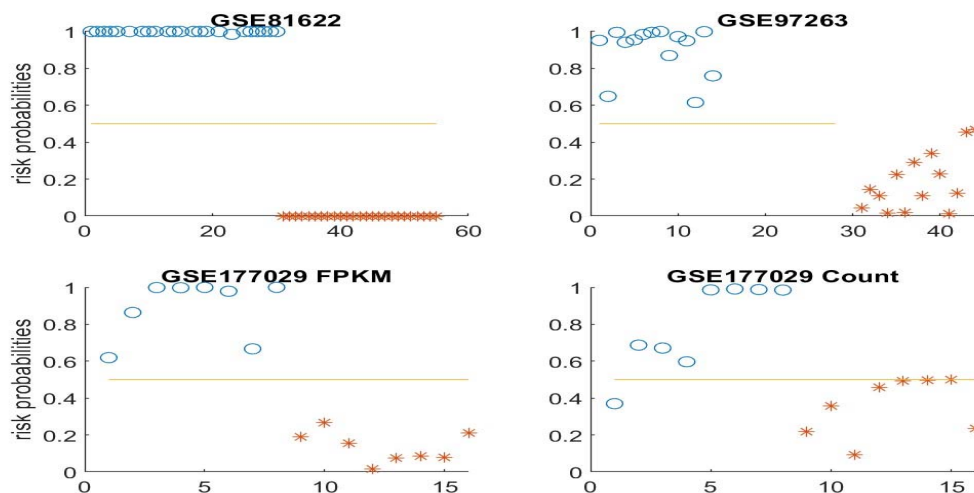
188

189 **Figure 1.** Venn diagrams are established from individual classifiers for each dataset.

190 Using GSE81622 as an example. The SLE patients are divided into three subtypes. Type 1
 191 contains 5 patients who can only be detected by CF1; Type 2 contains 12 patients who can only
 192 be detected by CF2; Type 3 contains 13 patients whom both CF1 and CF2 can detect. Other
 193 Venn diagrams are interpreted similarly.

194 We note that GSE177029 contains two types of data: FPKM values and raw counts. It contains
 195 127550 ensembl transcript (ENST) IDs. For illustration, we have GSE81622 containing 47322
 196 ILMN IDs, and GSE97263 contains 58037 ensembl gene (ENSG) IDs. Some ENSG IDs can have
 197 multiple ENST IDs, which makes multiple solutions possible. Nevertheless, we still reached
 198 100% accuracy with FPKM data.

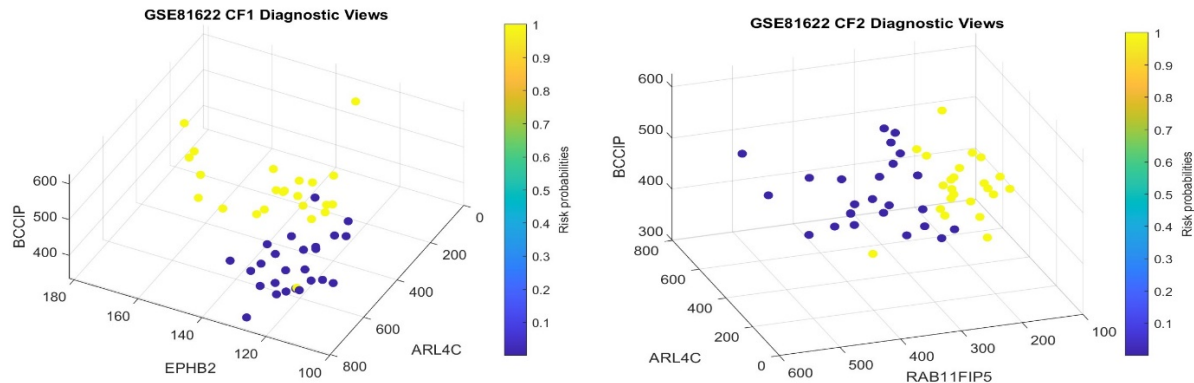
199 Figure 2 presents risk probabilities evaluated for each individual in four datasets.



200

201 **Figure 2.** Risk probabilities for all individuals (SLE and controls). Blue circles correspond to SLE
 202 patients, while asterisks are for controls.

203 Figure 3 presents gene-gene interaction clinical view plots for GSE81622.



204

205 **Figure 3.** Clinical view of GSE81622.

206 Figures 2 and 3 demonstrate that the identified genes show the greatest separability between
 207 SLEs and controls, and these genes can be treated as critical biomarkers.

208 In GSE50772, there are only 19943 genes recorded. The gene EPHB2 was not found in the data.
 209 Nevertheless, the remaining five genes still led to an accuracy of 95.06%, a sensitivity of
 210 96.72%, and a specificity of 90%. However, the coefficients associated with the genes
 211 RAB11FIP5 and BBCIP in GSE50772 are inconsistent with the signs in other datasets, making the
 212 interpretations of gene expression levels to the disease status not straightforward. Due to the
 213 missing EPHB2, its functions have to be allocated to other genes, and as a result, the signs of
 214 other genes may be affected, which does not disclose the truth. Such a phenomenon reveals
 215 that gene-gene interactions play pivotal roles in disease development, and finding driver genes
 216 is crucial to success.

217 In GSE4588, there are 54675 probe set IDs. Among them, (1,3,1,3,4,3) probe set IDs correspond
 218 to genes (NR3C2, SNX20, RAB11FIP5, ARL4C, EPHB2, BCCIP), respectively. In the footnote of
 219 Table 1, we listed the probe set IDs linked to ARL4C and BCCIP. The probe set IDs for NR3C2,
 220 SNX20, RAB11FIP5, EPHB2 used in Table 1 are 205259_at, 228869_at, 210879_s_at, 230088_at,
 221 respectively. At first glance, we see that the coefficient signs associated with genes in GSE4588
 222 are reversed from the signs in other datasets, and we may conclude Simpson's paradox exists in
 223 cell-level gene expression analysis. Note that in GSE4588, the expression values between CD4 T
 224 cells and B cells within several probe set IDs are significantly different. Figure 4 illustrates the
 225 differences. We can immediately see that DEGs at cell levels can be significantly different. As a
 226 result, it is not appropriate to put them in one model. Nevertheless, the genes identified for
 227 other cohorts still work perfectly with this cohort, indicating that these genes represent genetic
 228 insights into SLE. Likewise, the interpretations of gene functions in developing LSE have to
 229 consider its interactions with other genes and subtypes jointly.

cg05883128	DDX60	-25.71	-21.40	-18.07	-19.49	-18.35	-32.58	-34.21	-25.27			
cg09914304	PRF1		25.13									
cg12424383	NA									-26.06		
cg13210595				13.62								
cg14286514	DDX58				-15.70						-32.52	
cg14898177	CACNA1C					31.73						
cg15293582	PRF1						24.12					18.09
cg18686270	PLSCR1										-33.18	
cg18923906								18.65				
cg19371652	OAS2											-31.04
cg22365240	LINC01114								21.28	38.82		
	Sensitivity	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
	Specificity	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
	Accuracy	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%

243

244 From genecards.org, we can find that DDX60 (DEXD/H-Box Helicase 60) is a protein coding gene.
 245 Diseases associated with DDX60 include Plantar Wart and Lip Cancer. Gene Ontology (GO)
 246 annotations related to this gene include nucleic acid binding and hydrolase activity. An important
 247 paralog of this gene is DDX60L. DDX58 RIGI (RNA Sensor RIG-I) is a Protein Coding gene. Diseases
 248 associated with RIGI include Singleton-Merten Syndrome 2 and Singleton-Merten Dysplasia.
 249 Among its related pathways are DDX58/IFIH1-mediated induction of interferon-alpha/beta and
 250 SARS-CoV-2 Infection. An important paralog of this gene is IFIH1. PLSCR1 (Phospholipid
 251 Scramblase 1) is a Protein Coding gene. Diseases associated with PLSCR1 include Hepatitis B and
 252 Influenza. OAS2 (2'-5'-Oligoadenylate Synthetase 2) is a Protein Coding gene. Diseases associated
 253 with OAS2 include Microphthalmia With Limb Anomalies and Tick-Borne Encephalitis. Among its
 254 related pathways are Interferon gamma signaling and Antiviral mechanism by IFN-stimulated
 255 genes. PRF1(Perforin 1) is a Protein Coding gene. Diseases associated with PRF1 include
 256 Hemophagocytic Lymphohistiocytosis, Familial, 2 and Aplastic Anemia. Among its related
 257 pathways are Granzyme Pathway and TCR signaling in naïve CD4+ T cells. Gene Ontology (GO)
 258 annotations related to this gene include calcium ion binding and wide pore channel activity. An
 259 important paralog of this gene is C6. CACNA1 (Calcium Voltage-Gated Channel Subunit Alpha1 C)
 260 is a Protein Coding gene. Diseases associated with CACNA1C include Timothy Syndrome and Long
 261 Qt Syndrome 8. Among its related pathways are DREAM Repression and Dynorphin Expression
 262 and TCR Signaling (Qiagen). Gene Ontology (GO) annotations related to this gene include enzyme
 263 binding and monoatomic ion channel activity. An important paralog of this gene is CACNA1D.

264 LINC01114 (Long Intergenic Non-Protein Coding RNA 1114) is an RNA Gene affiliated with the
 265 lncRNA class. Diseases associated with LINC01114 include Childhood Ependymoma.

266 Table 2 shows universality and complexity. On the one hand, the CpG site cg05883128 (gene
 267 DDX60) shows its coefficient signs are all negative, which indicates the cause of SLE could be
 268 cg05883128 CpG site much under-methylated, and this site can be a target curable medical
 269 treatment DNA methylation site. On the other hand, cg05883128 can combine and interact with
 270 eight other CpG sites to reveal SLE disease, which greatly increases the SLE complexity. In the

271 combination classifiers CF1-3, CF5-8, increasing the methylation level of cg05883128 and
272 decreasing the methylation level of the other CpGs can benefit protecting SLE. It is interesting
273 that cg14286514 (DDX58) can benefit the SLE status via increasing its methylation level.

274 In addition, CF9, CF10, and CF11 are not related to cg05883128 (gene DDX60). Given their 100%
275 accuracy, it is clear SLE has complex subtypes at DNA methylation levels.

276 Notably, the gene PRF1 has two CpGs involving the development of SLE. LINC01114 can interact
277 with cg05883128 (gene DDX60) and cg12424383 for the development of SLE. Besides
278 cg05883128 (gene DDX60), cg14286514 (DDX58) can be a key site for SLE development if it is
279 under-methylated.

280 More importantly, Table 2 reveals why SLE can affect many parts of a person's body as there are
281 many gene (site) combinations, i.e., each can affect one part.

282 **5 Discussions and Conclusions**

283 **5.1 Discussions**

284 SLE is a complex disease, though not life-threatening. Given it can affect many parts of the body,
285 its indirect cause of death can be an issue. Understanding the SLE disease and its development
286 can help treat patients with multiple complex symptoms.

287 Among all genes identified in this paper, RAB11FIP5, DDX58, EPHB2 appeared in the list of
288 hundreds of genes in [15]. We note that our critical gene detection method is significantly
289 different from the literature approach. Our method focuses on site-site interaction, site-disease
290 subtype interaction, gene-gene interaction, and gene-disease subtype interaction. As a result,
291 the CpGs and genes identified in this paper can be interpreted using their fitted coefficients
292 (signs) and their combinations with other sites and genes. Equation (4) outperformed AI, machine
293 learning algorithm, and probabilistic algorithm in a COVID-19 study in [16] (Section 4.3).

294 In [2], the gene ITGA8 was discussed; see also [17]. However, this gene was not identified in our
295 computation. It is worth studying further. In [9,18-19], epithelial stromal interaction 1 (EPSTI1) is
296 associated with SLE. In [20-21], EPSTI1 and interferon (IFN)-alpha-inducible protein 27 (IFI27),
297 were found related to SLE. IFI27 and EPSTI1 were found to have pivotal links to SARS-CoV-2 in
298 [22]. Such a coincidence calls attention to COVID-19 infection-related SLE as so many people have
299 been infected with COVID-19.

300 Besides genes reported earlier, this set of genes can also be insightful: MGC40489, IL18RAP,
301 NUAK1, TOP1, MRPS18B, AL390026.1, ARPC4, LINC02762 as they also led to high and significant
302 accuracy in our computations.

303 Many SLE research results at the genomic level have been published in the literature. These
304 published results explored the pathological causes of SLE from various aspects. Due to study
305 methodology limitations, some of the published results can hardly be cross-validated from cohort
306 to cohort. Our work at the genomic level was a comprehensive study with nearly perfect
307 performance. We didn't find any other method that led to 100% accuracy in the literature, not

308 even to mention interpretability. Many studies focused on only a single cohort whose
309 representativeness cannot be assessed.

310 The dataset GSE4588 led to a Simpson's paradox: using gene NR3C2 as an example. At the overall
311 genomic level, the higher the DEGs, the better protection of LSE. On the other hand, at both CD4
312 T cells and B cell levels, the lower the DEGs, the better the protection of LSE. This phenomenon
313 calls for further profound studies of how different cell types interact, e.g., cell clustering and
314 classifications and identifying significant clusters. Also, the data process and quality have to be
315 controlled.

316 We now discuss the most significant difference between our approach and the literature
317 approach in finding critical genes. Much attention has been paid to the individual effects of every
318 single gene in the literature due to the study design and available analysis methods. Our approach
319 is jointly studying site-site, gene-gene, and gene-subtype interactions, which were largely missed
320 in the literature. We can see from Tables 1-2 that the effects of each gene depend on other genes
321 in the combinations. As a result, our findings of interaction effects can be the key to discovering
322 the cause of SLE.

323 Many published results studied the functional effects of genes (including IFI27 and EPSTI1) based
324 on single gene expression value changes. They lack interaction effects study, mainly due to the
325 limitations of study methods. As a result, they lack accuracy and may not really be useful. Using
326 the gene RAB11FIP5 as an example, in CF2 in Table 1 GSE81262, it must be jointly studied with
327 another two genes EPHB2 and BCCIP, to fully understand its functional effects on SLE as its
328 functional effects in CF1 with GSE177029 are significantly different.

329 Our results are nearly perfect, with some cohort studies with 100% accuracy and others with 95%
330 or higher accuracy. In some scenarios, such nearly perfect results can be considered too good to
331 be true. In our earlier work [16], we argued that the traditional cross-validation method is not
332 applicable to our model (3). Instead, we apply cohort-to-cohort cross-validation in our earlier and
333 present work. We used a driver gene dataset to demonstrate the superiority of our model (3)
334 compared to those algorithms built for AI, machine learning, and deep learning. We found that
335 our results are with better precisions, and more importantly, our results are interpretable [7, 22].

336 **5.2 Conclusions**

337 The pathological knowledge of the cause of the complex SLE is still unknown until now. The new
338 findings of cg05883128 (DDX60) can shed light on new research and treatment plans. NR3C2 and
339 RAB11FIP5 can be applied in medical treatment and drug development.

340 The theory of the cause of SLE is due to the increase of NR3C2 expressions in CD4 T-cells and B-
341 cells will shed light on biological studies. The theory will completely change the ways genetic,
342 cell, immunological studies on diseases are done.

343 **Acknowledgments**

344 In memory of Kathy Chow Hoi-mei for playing the role of Zhou Zhiruo in "Heaven Sword and
345 Dragon Sabre". Chow had suffered from SLE disease and died on December 11, 2023, which
346 motivated this research.

347

348 **Data Availability and Supplementary materials**

349 The datasets are publicly available. The data links are stated in Section Data Description.
 350 Computing outputs are in a supplementary file available online
 351 <https://pages.stat.wisc.edu/~zjz/SLEmarkers.zip> during the review process, and the file will be
 352 submitted to the publisher after the paper has been accepted. The results presented in this paper
 353 are all verifiable by simply checking the Excel sheets and formulas in the file.

354 **Competing Interests**

355 The Authors declare no Competing Financial or Non-Financial Interests.

356 **Author Contributions**

357 Zhengjun Zhang is the sole author with 100% contributions to the article.

358 **Statement of ethics**

359 The authors conducted research based on published work. The new research does not need IRB
 360 approval and a statement of ethics.

361

362 **Limitation statements**

363 Compared with other studies the authors have conducted, the sample sizes in this study are
 364 relatively small. Given 100% accuracy in all datasets, which naturally and reasonably confirms
 365 the results are intrinsic. Of course, large-scale medical research can further verify the findings.
 366

367 **References**

368

- 369 1. Tian J, Zhang D, Yao X, Huang Y, Lu Q. Global epidemiology of systemic lupus
 370 erythematosus: a comprehensive systematic analysis and modelling study. *Ann Rheum*
 371 *Dis*. Published online October 14, 2022. doi:10.1136/ard-2022-223035
- 372 2. Zucchi D, Silvagni E, Elefante E, Signorini V, Cardelli C, Trentin F, Schilirò D,
 373 Cascarano G, Valevich A, Bortoluzzi A, Tani C. Systemic lupus erythematosus: one year
 374 in review 2023. *Clin Exp Rheumatol*. 2023 May;41(5):997-1008. doi:
 375 10.55563/clinexprheumatol/4uc7e8. Epub 2023 May 3. PMID: 37133502.
- 376 3. Yu H, Nagafuchi Y, Fujio K. Clinical and Immunological Biomarkers for Systemic
 377 Lupus Erythematosus. *Biomolecules*. 2021 Jun 22;11(7):928. doi:
 378 10.3390/biom11070928. PMID: 34206696; PMCID: PMC8301935.
- 379 4. Zhang Z. Five critical genes related to seven COVID-19 subtypes: A data science discovery.
 380 *Journal of Data Science*, 19(1):142-150, 2021. <https://doi.org/10.6339/21-JDS1005>.
- 381 5. Zhang Z. The existence of at least three genomic signature patterns and at least seven
 382 subtypes of COVID-19 and the end of the disease. *Vaccines*, 10, 761, 2022.
 383 <https://doi.org/10.3390/vaccines10050761>.
- 384 6. Zhang Z. Functional effects of four or fewer critical genes linked to lung cancers and new
 385 sub-types detected by a new machine learning classifier. *Journal of Clinical Trials*,
 386 11:S14:100001, 2021. <https://www.longdom.org/open-access/functional-effects-of-four-or-fewer-critical-genes-linked-to-lung-cancers-and-new-subtypes-detected-by-a-new-machine-learning-clas-88321.html>
 387
 388

- 389 7. Liu, Y., Xu, Y., Li, X., Chen, M., Wang, X., Zhang, N., Zhang, H., Zhang, Z. (2023), Towards
390 Biomarker Discovery and Precision Oncology: Four Less Known Genes and Their Unknown
391 Interactions as Highest-Performed Biomarkers for Colorectal Cancer. *npj Precision*
392 *Oncology*. Accepted.
- 393 8. Zhu H, Mi W, Luo H, Chen T et al. Whole-genome transcription and DNA methylation analysis of
394 peripheral blood mononuclear cells identified aberrant gene regulation pathways in systemic
395 lupus erythematosus. *Arthritis Res Ther* 2016 Jul 13;18:162. PMID: [27412348](#)
- 396 9. Buang N, Tapeng L, Gray V, Sardini A et al. Type I interferons affect the metabolic fitness of
397 CD8⁺ T cells from patients with systemic lupus erythematosus. *Nat Commun* 2021
398 Mar 31;12(1):1980. PMID: [33790300](#)
- 399 10. Li HH, Sai LT, Tian S, Liu Y et al. Sexual Dimorphisms of Protein-Coding Gene Profiles in Placentas
400 From Women With Systemic Lupus Erythematosus. *Front Med* (Lausanne) 2022;9:798907.
401 PMID: [35372436](#)
- 402 11. Li HH, Sai LT, Liu Y, Freel CI et al. Systemic lupus erythematosus dysregulates the expression of
403 long noncoding RNAs in placentas. *Arthritis Res Ther* 2022 Jun 14;24(1):142. PMID: [35701843](#)
- 404 12. Kennedy WP, Maciuca R, Wolslegel K, Tew W et al. Association of the interferon signature
405 metric with serological disease manifestations but not global activity scores in multiple
406 cohorts of patients with SLE. *Lupus Sci Med* 2015;2(1):e000080. PMID: 25861459
- 407 13. Chen, W. et al. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE144390> China
408 2020.
- 409 14. Lauwerys BR, Louahed J, Knoop L, Maudoux A, Wakeland EK, Van den Eynde BJ,
410 Houssiau FA. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE4588> 2019
- 411 15. Woodridge, L., Chocano, EC., Ashford, P., Robinson, G., Waddington, K., Rahman, A.,
412 Orengo, C., Jury, EC., Torra, IP. (2022), Unique monocyte transcriptomic profiles are
413 associated with preclinical atherosclerosis in women with systemic lupus erythematosus
414 (SLE), *medRxiv* 2020.08.05.20169136; doi: <https://doi.org/10.1101/2020.08.05.20169136>
- 415 16. Zhang, Z., (2023). Discovery of Initial SARS-CoV-2 as DNA Viruses and Reliable Interactive
416 COVID-19 DNA Methylation Markers and Druggable Targets and Potential Malignant
417 Diseases with Long Incubation Period, 03 October 2023, PREPRINT (Version 4) available at
418 *Research Square* [<https://doi.org/10.21203/rs.3.rs-2248912/v4>]
- 419 17. Okazaki Y, Taniguchi K, Miyamoto Y et al. (2022), Glucocorticoids increase the risk of
420 preterm premature rupture of membranes possibly by inducing ITGA8 gene expression in
421 the amnion. *Placenta* 2022; 128: 73-82. <https://doi.org/10.1016/j.placenta.2022.07.012>
- 422 18. Mo, JS., Chae, SC. *EPSTI1* polymorphisms are associated with systemic lupus
423 erythematosus. *Genes Genom* 39, 445–451 (2017). [https://doi.org/10.1007/s13258-](https://doi.org/10.1007/s13258-017-0515-x)
424 [017-0515-x](https://doi.org/10.1007/s13258-017-0515-x)
- 425 19. Demers-Mathieu V. Optimal Selection of IFN- α -Inducible Genes to Determine Type I
426 Interferon Signature Improves the Diagnosis of Systemic Lupus Erythematosus.
427 *Biomedicines*. 2023; 11(3):864. <https://doi.org/10.3390/biomedicines11030864>
- 428 20. Ishii T, Onda H, Tanigawa A, Ohshima S, Fujiwara H, Mima T, Katada Y, Deguchi H, Suemura
429 M, Miyake T, Miyatake K, Kawase I, Zhao H, Tomiyama Y, Saeki Y, Nojima H. Isolation and
430 expression profiling of genes upregulated in the peripheral blood cells of systemic lupus
431 erythematosus patients. *DNA Res*. 2005;12(6):429-39. doi: 10.1093/dnares/dsi020. Epub
432 2006 Feb 23. PMID: 16769699.

-
- 433 21. Zhao, X., Zhang, L., Wang, J. et al. Identification of key biomarkers and immune infiltration
434 in systemic lupus erythematosus by integrated bioinformatics analysis. *J Transl Med* 19, 35
435 (2021). <https://doi.org/10.1186/s12967-020-02698-x>
- 436 22. Zhang Z. Genomic Biomarker Heterogeneities Between SARS-CoV-2 and COVID-19.
437 *Vaccines* 2022, 10(10), 1657; <https://doi.org/10.3390/vaccines10101657>