

1 A New Differential Gene Expression Based Simulated Annealing for 2 Solving Gene Selection Problem: A Case Study on Eosinophilic 3 Esophagitis and Few Other Gastro-Intestinal Diseases

4 Koushiki Sinha¹, Sanchari Chakraborty¹, Arohit Bardhan¹, Riju Saha¹, Srijan Chakraborty¹, Surama Biswas^{1,*}

5 **Abstract:**

6 *Background:*

7 Identifying the set of disease-causing genes is crucial for understanding pathogenesis and
8 developing therapies. This is particularly important to understand the pathophysiology of
9 Eosinophilic Esophagitis (EoE) and other gastrointestinal diseases. Comparing and contrasting gene
10 selection methods across these diseases can enhance our knowledge to identify potential therapeutic
11 targets.

12 *Methods:*

13 This study introduces two approaches for gene selection in gastrointestinal diseases: the Ranked
14 Variance (RV) method and Differential Gene Expression Based Simulated Annealing (DGESA).
15 RV acts as an initial screener by prioritizing genes based on variance. DGESA refines gene selection
16 further by employing simulated annealing with differential expression data. We compared the
17 outcomes of both methods through a case study on EoE and other gastrointestinal diseases.

18 *Results:*

19 Result finds greater number of genes with negative fold changes compared to those with positive
20 fold change in differential EoE dataset. RV Ranks top 40 genes with high variance of EoE which
21 overlaps with the disease-causing gene set of EoE from DGESA. 40 gene pathways for each of
22 EoE, Crohn's Disease (CD), and Ulcerative Colitis (UC) were identified as execution outcome of
23 our method DGESA. Among these, 10 genes for EoE, 8 for CD, and 7 for UC were confirmed in
24 the literature for their connection with respective diseases. For EoE, 10 such confirmed genes
25 include KRT79, CRISP2, IL36G, SPRR2B, SPRR2D and SPRR2E. For CD, the literature
26 confirmed set encompasses NPDC1, SLC2A4RG, LGALS8, CDKN1A, XAF1, and CYBA. The
27 validated genes in UC final gene set includes TRAF3, BAG6, CCDC80, CDC42SE2, and HSPA9.

28 *Conclusion:*

29 The RV method, serving as an initial screener, and the more refined DGESA both effectively
30 elucidate molecular signatures in gastrointestinal diseases. Identifying and validating genes like
31 SPRR2B, SPRR2D, SPRR2E and STAT6 for EoE showcase efficacy of DGESA. Other genes in
32 the same pathway are interesting targets for future laboratory validation.

33 **Introduction:**

34 Computational Genomics is one of the most emerging areas of biological research. It is empowered by the fusion of
35 powerful computational algorithms and vast genomic datasets. Its pivotal role in answering complex biological
36 questions is evident in its application across diverse domains, from tracing the origins of diseases like SARS to provide
37 insights on cancer immunotherapy [1]. By utilizing statistical methods and computational tools, bioinformaticians
38 decode the language of genomes, offer understanding of the evolutionary mechanisms, genomic variations and gene
39 interactions. Recent strides in high-throughput sequencing technologies have amplified the volume and complexity of
40 genomic data, necessitating sophisticated computational frameworks for data processing, analysis, and visualization.
41 These frameworks enable researchers to mine large datasets effectively, discovering novel mechanistic insights into
42 fundamental biological processes [2]. Moreover, as evidenced by the publicly available genomic databases and freely

1. Department of CSE, Meghnad Saha Institute of Technology, Behind Urbana Complex Near Ruby General Hospital, Anandapur Rd, Uchhepota, Kolkata, West Bengal 700150.

* Corresponding author: Dr. Surama Biswas, email: (surama.biswas@gmail.com)

43 downloadable program suits like Genomic Multi-tool (GEM), the democratization of computational tools fosters
44 collaboration and innovation in the scientific community. Computational Genomics not only shades lights on the
45 genetic factors of diseases but also escalates advancements in personalized medicine, drug discovery, and agricultural
46 resilience to climate change [3]. As universities introduce interdisciplinary courses to train the next generation of
47 bioinformaticians, the field continues to evolve, embracing emerging technologies like Machine Learning, Artificial
48 Intelligence and other modern computational methodologies. Through continuous refinement of computational
49 methodologies and the integration of diverse datasets, Computational Genomics promises to revolutionize our
50 understanding of life's fundamental processes, driving forward the frontiers of both biological discovery and
51 technological innovation [4].

52 Understanding biological processes requires an understanding of gene expression data, which provide insight
53 into how genes function within cells. There are many public repositories available. A notable instance of such a
54 database is the Gene Expression Omnibus (GEO) from National Center for Biotechnological Information (NCBI).
55 GEO provides high-throughput gene expression and functional genomics datasets in a consolidated manner, making
56 important information accessible to researchers across the globe. Many genomic data categories, such as, chromatin
57 structure, genome methylation, genome-protein interactions and gene expression studies are gathered in this database
58 [5]. GEO guarantees the availability of raw data and metadata by following community-driven reporting standards,
59 which promotes reliable research outputs. Scholars employ GEO's extensive collection to investigate a range of
60 biological inquiries, capitalizing on its intuitive interface and web-based instruments to facilitate streamlined data
61 retrieval, illustration, and examination. Scientific advancement and creativity are fueled by gene expression data from
62 repositories like GEO [6], which is used to drive discoveries in domains including developmental biology, cancer
63 research, and personalized medicine.

64 Gastrointestinal diseases encompass a diverse array of conditions affecting the digestive tract, ranging from
65 inflammatory bowel diseases like Crohn's disease and ulcerative colitis to eosinophilic esophagitis (EoE) [7]. EoE,
66 characterized by chronic immune-mediated inflammation of the esophagus, has emerged as a prominent source of
67 upper gastrointestinal morbidity in recent years. With an estimated prevalence of 34.4/100,000 in Europe and North
68 America, EoE presents symptoms such as esophageal strictures, dysphagia, and food impaction, affecting both
69 children and adults [8]. Unlike other conditions associated with esophageal eosinophilia, EoE diagnosis requires
70 symptoms of esophageal dysfunction alongside esophageal biopsies demonstrating at least 15 eosinophils per high-
71 power field. Genetic and environmental factors, including early exposure to antibiotics, are implicated in its etiology.
72 Current treatment modalities for EoE include proton pump inhibitors, dietary therapy, topical steroid formulations,
73 and endoscopic dilatation, tailored to individual patient needs and disease severity. As our understanding of EoE
74 evolves, further research into its pathogenesis, natural history, and optimal management strategies remains essential
75 for improving patient outcomes and quality of life [9 -11].

76 Genetic factors play a significant role in the pathogenesis of gastrointestinal diseases, with Eosinophilic
77 Esophagitis (EoE) standing out as a prime example. EoE, a chronic allergic condition characterized by eosinophilic
78 infiltration of the esophageal mucosa, is influenced by both hereditary and environmental factors [12]. Studies have
79 highlighted the substantial familial component of EoE, with evidence suggesting a higher likelihood of the condition
80 in family members of affected individuals [15]. Environmental risk factors also contribute to modulating genetic risk
81 in EoE, particularly through early-life events [13]. While rare genetic variations may account for a small subset of
82 EoE cases, the majority of genetic risk is mediated by common genetic variations [15]. Genome-wide association
83 studies (GWAS) have identified specific risk loci associated with EoE susceptibility, such as variants in genes like
84 TSLP and CAPN14, shedding light on the molecular mechanisms underlying the disease [14]. Interestingly, many of
85 these risk loci are located in non-coding regions of the genome, suggesting a role for gene regulation in EoE
86 pathogenesis [15]. Understanding the genetic architecture of EoE not only enhances our comprehension of its
87 molecular basis but also holds promise for the development of targeted therapeutic interventions and personalized
88 treatment strategies based on individual genetic profiles [14]. Further research into the intricate interplay between

89 genetic predisposition, environmental triggers, and immune dysregulation is essential for advancing our understanding
90 of EoE and improving patient care [12].

91 The prevalence of eosinophilic esophagitis (EoE) in the Indian population appears to be gradually gaining
92 recognition, although with limited data available. A study conducted in a hospital in the northern region of India
93 reported a prevalence of 3.2% among patients with symptoms suggestive of gastroesophageal reflux disease (GERD)
94 [16]. Similarly, another study noted an increase in diagnosed cases of EoE in their center over recent years, with 17
95 out of 73 patients being diagnosed with EoE based on clinical, endoscopic, and histopathologic features [17]. These
96 findings suggest that EoE exists in India and is gaining importance of clinical suspicion and diagnostic evaluation for
97 its identification. However, the prevalence of EoE in the Indian population remains to be fully explained, highlighting
98 the necessity for large-scale, multi-centric population-based studies to provide a more comprehensive understanding
99 of the disease burden in the country [17].

100 Gene selection, a critical task in gene expression studies, aims to identify subsets of genes that are relevant
101 for distinguishing between disease and normal conditions. Traditional methods for gene selection often face challenges
102 such as high within-class variation and the selection of an optimal subset of genes that can efficiently differentiate
103 between classes. Recent advancements in artificial intelligence (AI) and machine learning (ML) have introduced novel
104 approaches to address these challenges. For instance, in a study [18], a novel criterion was proposed for assessing the
105 significance of individual genes based on their mean and standard deviation of distances from each sample to the class
106 centroid. This method not only effectively tackles within-class variation but also offers a smaller number of genes
107 without compromising discriminating power, thus supporting further biological and clinical research. Similarly, the
108 utilization of ML techniques, such as random forest, has been demonstrated to be effective in gene selection for
109 microarray data classification [19]. Random forest excels in handling large numbers of variables and noisy data,
110 providing accurate predictions while simultaneously offering small sets of genes for classification. Additionally,
111 machine learning-based approaches have been employed to enhance the stability of gene selection techniques under
112 sample fluctuations. For example, a study by [20] introduced a framework of sample weighting to increase the stability
113 of representative feature selection algorithms, leading to the identification of more stable gene signatures.
114 Furthermore, the application of ML techniques extends beyond gene selection to various aspects of cancer research,
115 including cancer subtype classification, prognosis prediction, and identification of biomarkers. Studies such as [21,
116 22, 23, 24] highlight the utility of ML algorithms like XGboost and support vector machine (SVM) in classifying
117 cancer subtypes and identifying potential biomarkers for hepatocellular carcinoma (HCC) and Crohn's disease (CD),
118 respectively. These approaches leverage large-scale genomic data to facilitate personalized treatment strategies and
119 improve patient outcomes. In summary, the application of AI and ML techniques holds great promise in addressing
120 the gene selection problem by providing efficient, accurate, and stable methods for identifying disease-relevant genes
121 and advancing our understanding of complex diseases [25, 26, 27].

122 Gene expression variance plays a crucial role in computational genomics, influencing our understanding of
123 complex genetic diseases and population genetics. Studies such as [28] have highlighted the significance of genetic
124 variants in Alzheimer's disease (AD) risk, with known single nucleotide polymorphisms (SNPs) explaining only a
125 portion of the phenotypic variance. This represents the importance of exploring additional sources of genetic variation,
126 such as rare or unknown SNPs, to take into account the full genetic landscape of AD. Similarly, research on the
127 ALDH2 gene variant, as discussed in [29], sheds light on how specific genetic variations can influence susceptibility
128 to various diseases and physiological traits. Understanding the mechanisms underlying these associations is essential
129 for precision medicine and disease prevention strategies. Furthermore, investigations into gene expression variance,
130 as described in [30], provide insights into the regulatory mechanisms governing gene expression across different
131 tissues and conditions. By observing patterns of transcriptional variance and its relationship with gene function,
132 computational genomics can identify key regulatory elements and pathways underlying complex traits and diseases.
133 Overall, the study of gene expression variance in computational genomics enhances our understanding of genetic

134 diversity, disease susceptibility, and molecular mechanisms, paving the way for more effective diagnostic and
135 therapeutic interventions.

136 Simulated Annealing (SA) draws inspiration from metallurgy's annealing process, where material is
137 gradually cooled to a stable state. Introduced by Kirkpatrick, Gelatt, and Vecchi [31], SA optimizes by navigating a
138 solution space, occasionally accepting moves that increase the cost function value (in minimization). As it progresses,
139 the algorithm gradually reduces the likelihood of accepting worse solutions, akin to cooling. SA's extensions include
140 adaptations to combinatorial optimization and nonconvex problems [32-33]. In genomic feature selection, SA aids in
141 identifying relevant genes from high-dimensional datasets, such as microarray gene expression data. Recent
142 advancements like curious simulated annealing [34] and Simulated Annealing aided Genetic Algorithm (SAGA) [35]
143 address SA's convergence limitations. They strike an optimal balance between exploring the search space and
144 exploiting potential solutions.

145 The gene selection problem, which focuses on identifying a set of genes that collectively cause a disease, is
146 crucial for understanding complex medical conditions. Though many complex formulations are already available, a
147 very simple, efficient and biologically plausible fact that optimization process needs guidance from the
148 differentiability of diseased to normal genomic profiles was overlooked. This study introduces a new Simulated
149 Annealing based algorithm called Differential Gene Expression Based Simulated Annealing (DGESA) where a
150 specially designed objective function has been introduced which aims to maximize the collective differentiability of
151 the obtained genes in their diseased and normal genomic profiles. For an initial guess of differentiability in the gene
152 expression data, an approach, termed here as Ranked Variance (RV) has been introduced that prioritize genes based
153 on their variance. Through a case study on Eosinophilic Esophagitis (EoE) and other gastrointestinal diseases, we
154 compare the outcomes of both methods. Notably, we find 10 common genes between RV and DGESA in EoE,
155 indicating their complementary nature. Validation analyses show that 10 of the 40 final genes identified by DGESA
156 for EoE are supported by existing literature, confirming their biological relevance. Similarly, for Ulcerative Colitis
157 (UC) and Crohn's Disease (CD), 8 and 7 of the 40 genes, respectively, are validated by literature. Ten confirmed genes
158 for EoE are as follows: KRT79, CRISP2, IL36G, SPRR2B, SPRR2D, and SPRR2E. The collection of CD that has
159 been confirmed by literature includes NPDC1, SLC2A4RG, LGALS8, CDKN1A, XAF1, and CYBA. The final gene
160 set from UC contains the validated genes TRAF3, BAG6, CCDC80, CDC42SE2, and HSPA9. These results
161 underscore the efficacy of our framework and specifically DGESA in identifying significant molecular signatures
162 associated with gastrointestinal diseases.

163 **Methodology:**

164 The method section of this study details two distinct approaches employed for gene selection and analysis: the RV
165 method and DGESA. The RV method prioritizes genes based on their variance, providing an initial perspective on
166 potential biomarkers. In contrast, DGESA utilizes simulated annealing to identify sets of genes exhibiting significant
167 differences in expression between diseased and normal states, facilitating the discovery of disease-associated genetic
168 signatures. Each method offers unique insights into gene selection and contributes to our understanding of molecular
169 mechanisms underlying disease pathogenesis (see Figure 1).

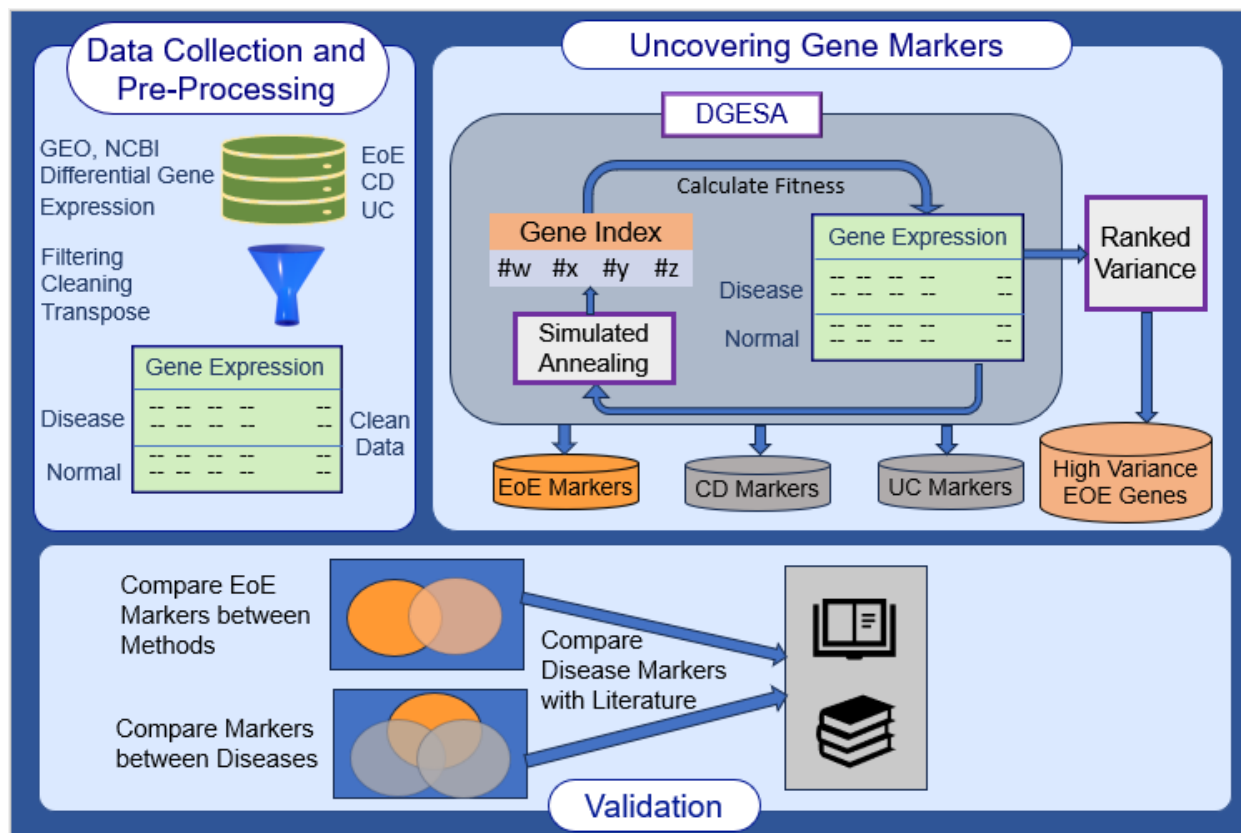
170 Prior to analysis, several preprocessing steps were implemented to ensure data quality and compatibility.
171 Firstly, rows lacking valid gene names were removed to maintain consistency across datasets. Subsequently, a
172 normalization procedure was applied to each dataset, wherein gene expression values ($e_{i,j}$) were mapped to the range
173 [0, 1]. This normalization step helped maintain potential biases arising from variations in gene expression magnitude
174 across samples. Finally, transposition of the datasets was performed to prepare the data matrix (denoted as X) for
175 subsequent processing, facilitating the application of required analytical techniques. These preprocessing steps
176 collectively ensured that the gene expression data were standardized and conducive to meaningful analysis and
177 interpretation.

178 The Ranked Variance (RV) method employed in this study focused on using gene expression variability as a
179 means of discerning potential biomarkers associated with disease states. By computationally analyzing the variance
180 of gene expression across samples, the RV method identified genes exhibiting significant variations in expression
181 levels. This approach facilitated the separation of disease-associated genes from those with relatively stable expression
182 patterns, thereby providing valuable insights into the molecular mechanisms underlying disease pathogenesis.
183 Moreover, the identification of genes with pronounced expression variations enabled subsequent association studies,
184 wherein these genes could be further investigated for their roles in disease development, progression, and potential
185 therapeutic targeting. Overall, the RV method served as a powerful tool for examining the genetic signatures
186 associated with various diseases, contributing to our understanding of their underlying biological processes and aiding
187 in the discovery of novel biomarkers.

188 DGESA is a methodology devised to address gene selection challenges in the context of biological diseases.
189 At its core, DGESA operates on a transposed gene expression matrix (denoted as X), where each row represents an
190 observation (diseased or normal person) and each column corresponds to a gene. The method begins by defining a
191 candidate solution, represented as a set of gene indices, which is iteratively refined through a simulated annealing
192 process. During each iteration, a perturbation is applied to the current solution by randomly altering a gene index from
193 the gene expression matrix X that is not already present in the solution set s . The fitness of each candidate solution is
194 evaluated using a devised fitness function, represented by Equation (1):

195
$$\left| \sum_{i=1}^g (\overline{D}(e_i)) - (\overline{Norm}(e_i)) \right| \quad (1)$$

196 Here, g represents the number of genes in the candidate solution s . For each gene index i in s , the expression
197 profile (e_i) is considered. The mean expression profile of the i -th gene in the candidate solution s for diseased samples
198 is denoted by $(\overline{D}(e_i))$, while the mean expression profile for normal samples is denoted by $(\overline{Norm}(e_i))$. The fitness
199 function computes the absolute difference between the mean expression profiles of diseased and normal samples
200 across all genes in the candidate solution. This difference serves as a measure of the discriminative power of the
201 selected genes in distinguishing between diseased and normal states.



202

203 *Figure 1: Overview of the Study Workflow.* This figure illustrates the comprehensive workflow of the study, beginning
204 with data acquisition and preprocessing steps that include cleaning, normalization, and transposition to produce the gene
205 expression matrix (samples as rows, genes as columns). The gene marker identification methodologies encompassing
206 RV and DGESA, are applied to datasets of EOE, CD, and UC to identify disease-specific gene markers. The final step
207 involves the validation process, which compares the effectiveness of the two methods across the three diseases and
208 validates the identified markers against existing literature.

209 In DGESA, the Temperature (T) is first initialized into a very high value. Then few simulated annealing steps
210 have been performed until the stopping criteria are not met. In each such step, L neighborhood search iterations,
211 followed by a reduction of T by a fraction of Cooling Factor C have been performed. In each neighborhood search
212 state, we select a neighbor s' from the neighborhood of s and replace s by s' if s' is better than s in terms of fitness
213 value otherwise s may be replaced by s' depending on a small probability. Through this iterative optimization process,
214 DGESA aims to identify a set of genes that collectively exhibit significant differences in expression patterns between
215 diseased and normal samples. The output of DGESA, denoted as s^* , represents the final selection of genes maximizing
216 the discrimination between diseased and normal conditions in gene expressions, thereby facilitates the identification
217 of potential biomarkers and provides insights into disease mechanisms.

218 The methods applied in this study, including RV method and DGESA, have provided valuable insights into
219 gene selection and analysis within the context of gastrointestinal diseases. The RV method effectively identified genes
220 with significant expression variations, aiding in disease gene separation and association studies. On the other hand,
221 DGESA upgrades simulated annealing to pinpoint genes exhibiting differential expression patterns between diseased
222 and normal samples, thereby contributing to the discovery of disease-associated genetic signatures. By employing
223 these complementary methodologies, the understanding of molecular mechanisms underlying disease pathogenesis
224 have been advanced and a robust framework for biomarker discovery and disease classification is presented.

DGESA

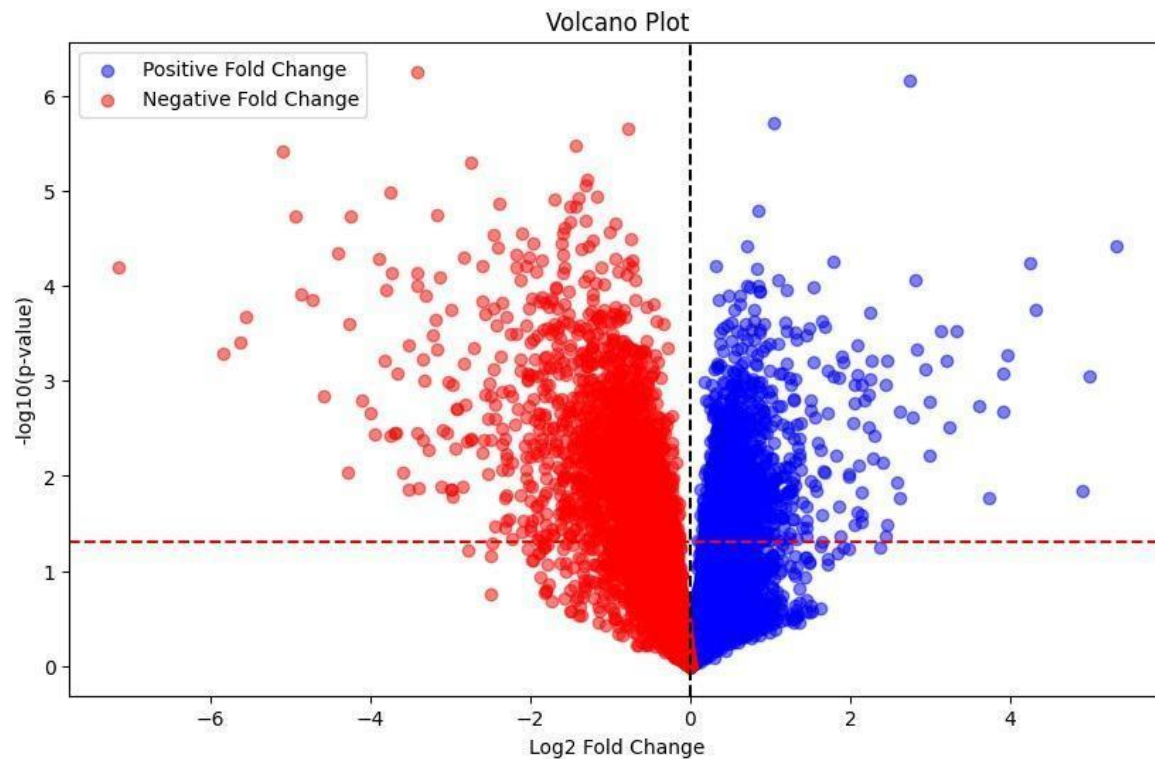
```
1. Initialize temperature T,
Iterations P and Cooling rate C;
2. Choose initial candidate solution
s with random gene index;
3. Calculate fitness f(s) using
Eq. (1);
4. Repeat
  A. for i = 1 to P do
    ◦ Randomly select  $s' \in N(s)$ ;
    // N(s): Neighbor of s
    ◦ if  $f(s') \geq f(s)$  then
       $s \leftarrow s'$ ;
    ◦ else
       $s \leftarrow s'$  with
      probability  $e^{-(f(s')-f(s))/T}$ ;
    ◦ end if
  B. end for
  C.  $T = T \times C$ ;
5. Until stopping criteria not met
6. end
```

225

226 *Figure 2: DGESA Algorithm.* This figure depicts the details of the DGESA algorithm. The process involves simulated annealing
227 with a specially designed objective function to effectively select gene markers from differential gene expression data.

228 Results:

229 The data collection for this study involved retrieving two gene expression datasets from GEO, NCBI. The first dataset,
230 GSE228083, comprised samples from patients with EoE compared to normal samples, facilitating the investigation of
231 gene expression patterns specific to this condition. The second dataset, GSE24287, encompassed gene expression
232 profiles from patients with UC, CD, and normal samples. From GSE24287, two distinct datasets were prepared by
233 segregating samples into UC vs. Normal and CD vs. Normal categories.



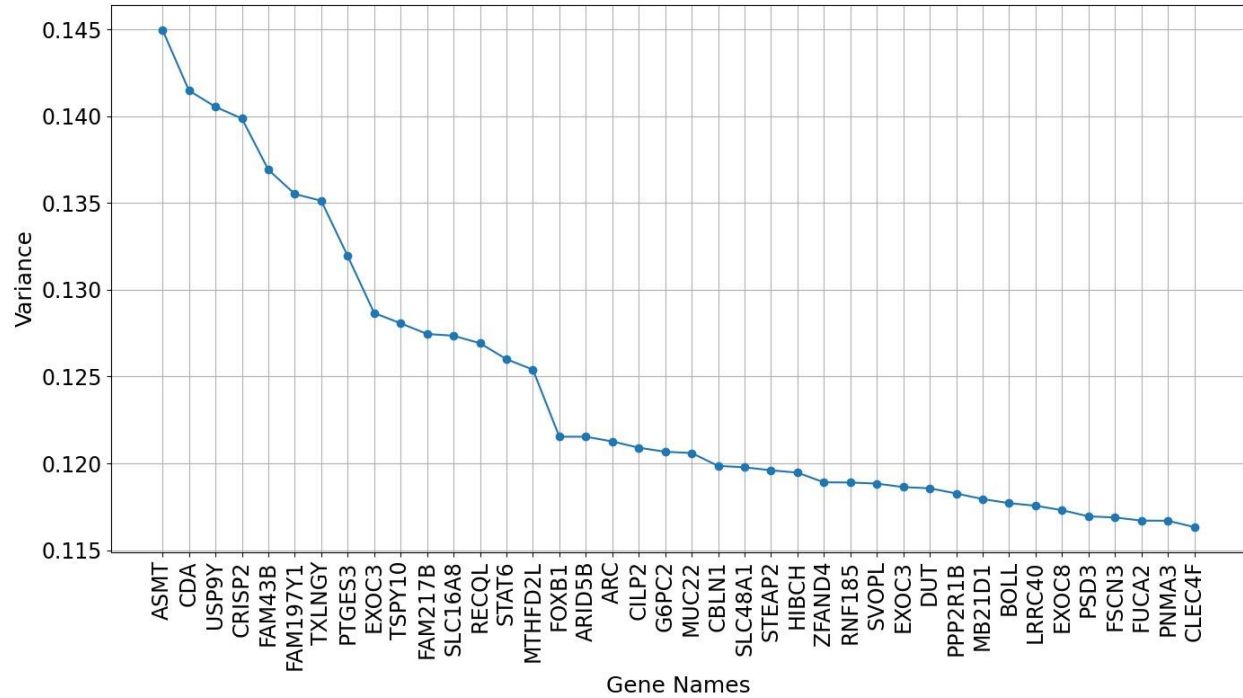
234

235 *Figure 3. Volcano Plot of Gene Expression Data for EOE:* This figure presents a volcano plot illustrating the gene expression data
236 for EOE. The x-axis represents the \log_2 fold change, while the y-axis denotes the $-\log_{10}$ P-value. The plot highlights a greater
237 number of genes with negative fold changes compared to those with positive fold changes, indicating differential gene expression
238 patterns in EOE.

239 Hyper-parameter tuning of DGESEA is a critical aspect of optimizing its performance in gene selection tasks.
240 This iterative process involves systematically adjusting parameters that control the learning process, known as hyper-
241 parameters, through experimentation with different configurations. In the case of DGESEA, key hyper-parameters
242 include the number of genes (g) in the candidate solution, the maximum number of iterations, the initial temperature
243 (T), and the cooling rate (C). By systematically adjusting these hyper-parameters, the DGESEA model can be fine-
244 tuned to enhance its efficiency in identifying disease-associated genes. In this study, after thorough experimentation
245 and analysis of resulting performance of the algorithm to convergence, the final hyper-parameter configurations were
246 determined as follows: $g = 40$ genes in the candidate solution, 200,000 iterations, $T = 10^6$, and a cooling rate of 0.9.
247 These optimized hyper-parameters ensure the effectiveness of DGESEA in identifying relevant genetic signatures
248 associated with gastrointestinal diseases, thereby advancing our understanding of disease mechanisms and aiding in
249 biomarker discovery.

250 To gain insight into the differential expression patterns of genes in the EoE dataset, a volcano plot was
251 generated, depicting the relationship between the log₂ fold change and the log₁₀ p-values of various genes. In this
252 plot, the x-axis represents the log₂ fold change, which quantifies the magnitude of gene expression differences
253 between EoE samples and normal samples. Meanwhile, the y-axis displays the -log₁₀ p-values, which serve as a
254 measure of the statistical significance of these expression differences. The volcano plot (see Figure 3) revealed that
255 the majority of genes exhibited negative fold changes, indicating under-expression in EoE compared to normal
256 samples. This observation suggests a potential downregulation of gene expression associated with EoE pathology.
257 However, it's essential to interpret these findings in conjunction with additional analyses to elucidate the specific genes
258 and biological pathways underlying the disease's pathogenesis and progression.

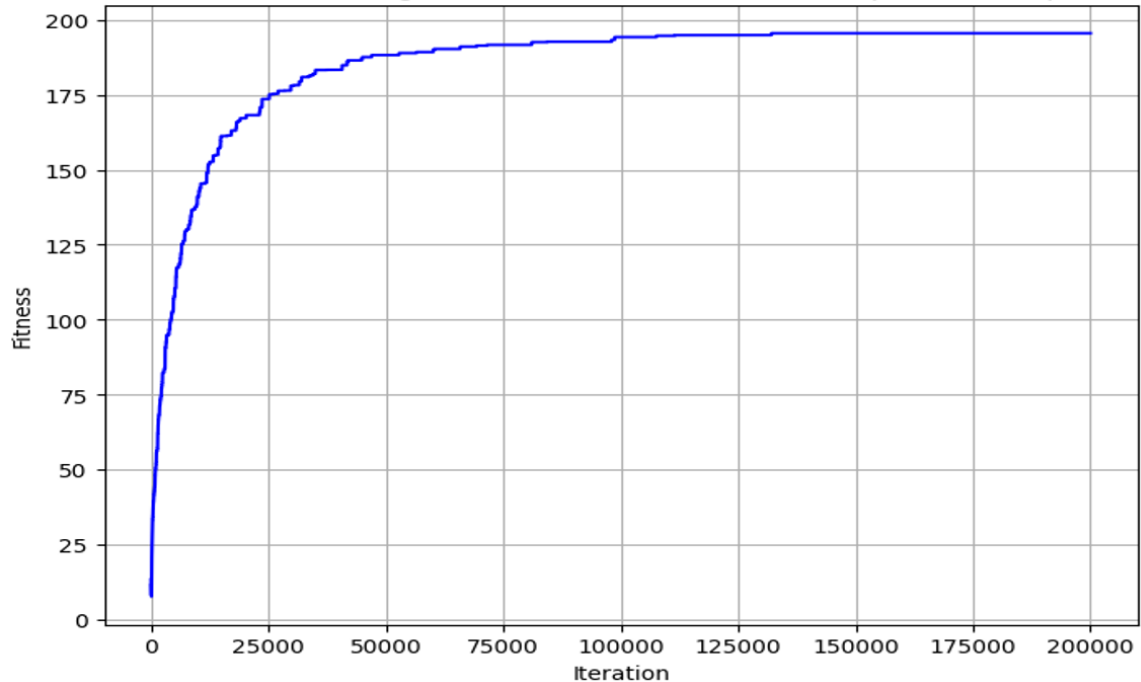
259 The application of the RV method to the EoE vs. normal dataset yielded insightful results regarding the
260 variability of gene expression across samples. By plotting the curve (see Figure 4) where the x-axis represents the
261 gene index and the y-axis denotes the corresponding variance, it was observed that approximately 40 genes exhibited
262 decreasing variance. This observation suggests a notable reduction in the variability of expression levels for these
263 genes in EoE samples compared to normal samples. Such a trend of decreasing variance may indicate a degree of
264 regulatory homogeneity or consistent downregulation of gene expression within this subset of genes in the context of
265 EoE pathology. These findings highlight the potential significance of these genes in contributing to the molecular
266 mechanisms underlying EoE development and progression. The variance graph of CD and UC are available on
267 Supplement 1 and 2 respectively.



268

269 *Figure 4. Variance of Top 40 Genes in EOE Gene Expression Data:* This figure shows a line plot of the variance of the top 40
270 genes in the EOE gene expression dataset. The x-axis represents the gene names with the highest variance, and the y-axis indicates
271 the variance values. The plot reveals that the initial few genes exhibit very high variance, which gradually flattens down for the
272 subsequent genes.

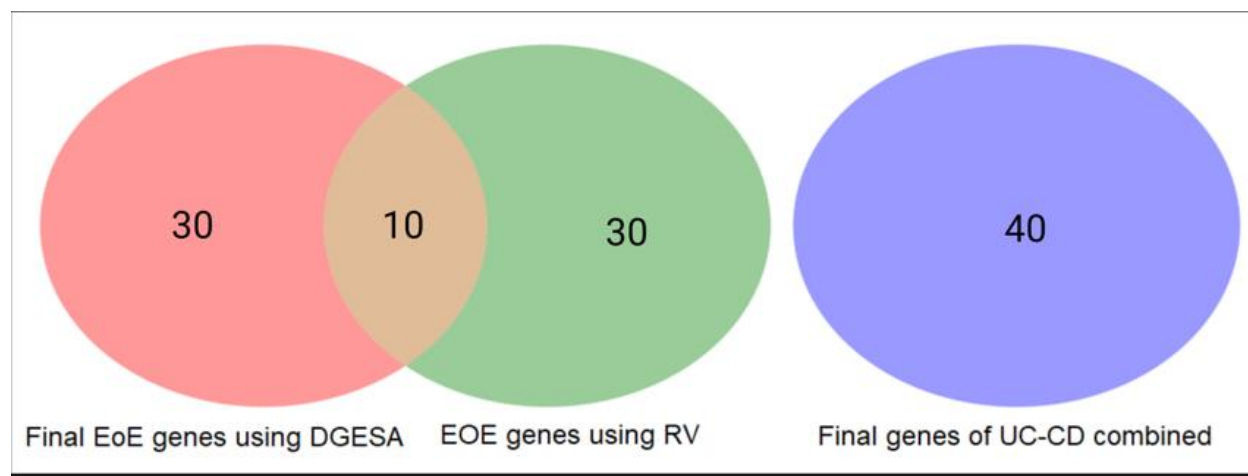
273 The application of the DGESA method to the EoE vs. normal dataset yielded a convergence curve that
274 provides valuable insights into the optimization process. In this curve, the x-axis represents the iterations, reflecting
275 the number of iterations or steps taken during the simulated annealing optimization procedure. Meanwhile, the y-axis
276 denotes the corresponding fitness values, which quantify the effectiveness of the candidate solutions at each iteration.
277 The observed convergence of the curve indicates that as the optimization progresses through iterations, the fitness
278 values gradually stabilize or improve, eventually reaching an optimal or near-optimal solution. This convergence
279 phenomenon signifies the effectiveness of DGESA in iteratively refining the selection of genes to maximize their
280 discriminative power between EoE and normal samples. The convergence curve represents the robustness and
281 efficiency of DGESA in (see Figure 5) identifying disease-associated genetic signatures and highlights its potential as
282 a valuable tool for biomarker discovery. Showcasing the consistency, the optimization curves obtained from CD vs.
283 normal and UC vs. normal datasets show convergence (see Supplement 3 and 4 respectively).



284
285 *Figure 5. Optimization Process for Identifying 40 Discriminative Genes in EOE:* This figure illustrates the optimization process
286 aimed at selecting a set of 40 genes that collectively maximize the discrimination between EOE and control samples in the gene
287 expression dataset. The curve shows the convergence of the optimization process, stabilizing near 130,000 epochs.

288 The Venn diagram analysis comparing the gene sets identified by the DGESA and RV methods in the context
289 of EoE versus normal samples revealed important findings (see Figure 6). Specifically, the diagram indicated that 10
290 genes were shared between DGESA and RV (see Table 1 for the common gene names), suggesting a degree of
291 consistency or agreement between the two methods in identifying potential biomarkers associated with EoE. However,
292 notably, no overlap was observed between the gene sets identified for EoE and the combined UC and CD versus
293 normal dataset. This absence of connection between the gene sets for EoE and UC-CD may reflect distinct molecular
294 mechanisms underlying these two gastrointestinal diseases. The lack of shared genes highlights the specificity of gene
295 expression profiles associated with each disease entity and emphasizes the importance of tailored approaches for
296 biomarker discovery and therapeutic targeting. Further investigation into the unique genetic signatures of EoE and
297 UC-CD could offer deeper insights into their pathogenesis and facilitate the development of more precise diagnostic
298 and therapeutic strategies.

299 The final gene set identified by the DGESA algorithm for EoE presents a compelling alignment with previous
300 literature, demonstrating its potential significance in the context of EoE pathology. Among the 40 genes listed (see
301 Table 2), 10 genes, including KRT79 [36], SPRR2E [37], CRISP2 [38], IL36G [39], SPRR2B [37, 40], SPRR2D [37,
302 40] and RORC [41] have been previously implicated in EoE and are highlighted in blue in Table 2 to denote their
303 strong confirmation of association with the disease. These genes represent key players in various molecular pathways
304 relevant to EoE, such as immune regulation, epithelial barrier function, and tissue remodeling. Interestingly, CCND1
305 [42] linked with allergy in cow milk, appeared in our result. It might have an indirect link with EoE because cow milk
306 allergy has been considered as one of the factors causing EoE. Additionally, the presence of other genes in the final
307 gene set, although not explicitly highlighted, suggests potential connections to EoE based on their co-appearance with
308 established EoE-associated genes. This comprehensive gene set derived from DGESA not only validates known
309 associations but also offers new insights into the molecular mechanisms underlying EoE pathogenesis, paving the way
310 for further research into diagnostic and therapeutic interventions for this complex disease. The unique genes of CD
311 and UC by DGESA are available on Table 3 and 4 respectively.



312
313 *Figure 6. Venn diagram of Common Genes Identified by RV and DGESA Methods in EoE and UC-CD Combined*
314 *Datasets.* This figure shows a Venn diagram comparing the common genes identified by Ranked Variance (RV) and Differential
315 Gene Expression Based Simulated Annealing (DGESA) methods in the EoE gene expression dataset (left side) and the common
316 genes identified using DGESA in the UC-CD combined dataset. The diagram reveals 10 common genes between the RV and
317 DGESA methods for EoE, but no overlapping genes between the EoE dataset and the UC-CD combined dataset.

318 *Table 1: Common genes of EoE obtained using RV and DGESA*

CRISP2	FAM43B	PPP2R1B	MTHFD2L	EXOC3	RECQL	STAT6	ARC	FAM217B	ZFAND4
--------	--------	---------	---------	-------	-------	-------	-----	---------	--------

319 *Table 2: Final unique genes of EoE using DGESA:* Here the genes highlighted with blue have confirmed their connection with
320 state-of-the-art literature as related with CD. For example, KRT79 [36], SPRR2E [37], CRISP2 [38], IL36G [39], SPRR2B [37,
321 40], SPRR2D [37, 40], RORC [41] and CCND1 [42] are confirmed connection with EoE.

HSPA12A	DPCR1	FAM25G	FAM43B	IVD	PPP2R1B	MTHFD2L	SPRR2E	GGT6	KRT79
RORC	CRISP2	ZNF562	OAZ3	C18orf54	EXOC3	IL36G	TPPP2	ANXA8	CSN2
RECQL	RPAP3	SPINK13	TAF4B	LYPD6	COX6B1	CPB1	SPRR2B	SPRR2D	DMKN
FAM217B	HIP1	ARC	ZFAND4	CCND1	RMI1	LOC388780	CSNK1A1L	ADGRB1	STAT6

322 *Table 3: Final unique genes of CD using DGESA:* Here the genes highlighted with blue have confirmed their existence in state-
323 of-the-art literature as related with CD. For example, NPDC1 [43], SLC2A4RG [44], LGALS8 [45], CDKN1A [46], XAF1 [47],
324 *CYBA* [48] etc. are confirmed connection with Inflammatory Bowel Syndrome (IBD) like CD.

EPHX3	NRG1	PC	NPDC1	VTI1A	ZNF646	TCEA1	SLC2A4RG	LGALS8	VDAC3
G2E3	ATXN1	KTN1	PLEK2	CDKN1A	ACOT8	GLB1	XAF1	CYBA	WDR85
NF1	ANXA2P1	PSAT1	ABHD11	HPN	PEX19	MAPK1	MKLN1	PSMD5	RBMX2
PNPO	ARID3B	PRPS1	GJD3	TECPR1	CAPS	ZIM2	H3F3A	POPDC2	BHLHE22

325

326

327 *Table 4: Final unique genes of UC using DGESEA:* Here the genes highlighted with blue have confirmed their existence in state-of-
328 the-art literature as related with UC. For example, TRAF3 [49], BAG6 [50], CCDC80 [51], CDC42SE2 [52] and HSPA9 [53] are
329 confirmed connection with Inflammatory Bowel Syndrome (IBD) like UC.

C9orf25	TRAF3	CNOT1	BAG6	ZNF658	ARL3	KPNA3	NOP16	C1orf182	HOXD9
EME1	RPL22	XPO6	LNX1	CCDC80	DMXL2	QRSL1	WDR55	SLC4A7	HCP5
GPR137	PJA2	GOLGB1	BSG	CYP4F12	SLC35A2	RSPO3	PPP6C	C6orf115	TNPO1
SEC61G	CDC42SE2	FKBP1A	SRSF1	FOXO1	LOC401022	PRODH	PRPF18	HSPA9	C10orf125

330

331 Discussion

332 The identification of genes implicated in complex diseases is pivotal for advancing our understanding of their
333 underlying biological mechanisms. In this study, we introduce the Differential Gene Expression Based Simulated
334 Annealing (DGESEA) algorithm, which employs a novel objective function to enhance the differentiability of gene
335 expression profiles between diseased and normal states. This method addresses a significant gap in existing gene
336 selection strategies by emphasizing the importance of collective differentiability, a concept that has been overlooked
337 in previous research.

338 Our results from applying DGESEA to datasets of Eosinophilic Esophagitis (EoE), Ulcerative Colitis (UC),
339 and Crohn's Disease (CD) demonstrate its potential in uncovering biologically relevant genes. By comparing the
340 outcomes of DGESEA with the Ranked Variance (RV) approach, we found a noteworthy overlap, particularly in the
341 EoE case study, where 10 common genes were identified. This convergence underscores the complementary nature
342 of these methods and suggests that the integration of multiple approaches can enhance the robustness of gene selection
343 processes.

344 The validation of our findings against existing literature further supports the efficacy of DGESEA.
345 Specifically, 10 of the 40 genes identified for EoE, 8 for CD, and 7 for UC were corroborated by previous studies,
346 highlighting their biological relevance. For instance, genes such as KRT79, CRISP2, and IL36G in EoE, and CDKN1A
347 and CYBA in CD, have established roles in the pathophysiology of these diseases, which reinforces the credibility of
348 our method. These validated genes also provide valuable targets for future research and potential therapeutic
349 interventions.

350 The introduction of the Ranked Variance (RV) approach as an initial step in the DGESEA process is
351 particularly noteworthy. By prioritizing genes based on their variance, RV offers a biologically plausible preliminary
352 filter that simplifies the subsequent optimization process. This step not only enhances the efficiency of DGESEA but
353 also ensures that the selected genes are inherently variable and thus more likely to be differentially expressed between
354 diseased and normal states.

355 Our framework's ability to identify significant molecular signatures associated with gastrointestinal diseases
356 holds promise for broader applications. The DGESEA method can be adapted to various other diseases and datasets,
357 potentially uncovering critical genes that have been missed by traditional methods. Moreover, the integration of
358 DGESEA with other bioinformatics tools and databases could further enhance its utility and accuracy.

359 Despite the promising results, several limitations should be acknowledged. The reliance on existing literature
360 for validation, while necessary, may introduce bias, as genes not yet studied or published may be equally important
361 but remain unrecognized. Additionally, the performance of DGESEA should be evaluated on larger and more diverse
362 datasets to ensure its generalizability across different populations and disease contexts.

363 Future research should focus on refining the objective function and exploring alternative strategies for initial
364 gene prioritization. Integrating additional layers of biological data, such as protein-protein interactions and pathway
365 analyses, could provide a more comprehensive understanding of the selected genes' roles in disease. Furthermore,
366 experimental validation of the identified genes through laboratory studies will be crucial in confirming their functional
367 relevance and potential as therapeutic targets.

368 **Conclusion**

369 The identification of gene sets that collectively cause a disease, known as the gene selection problem, is a critical area
370 of study in understanding complex diseases. This research introduces two innovative approaches for gene selection in
371 the context of various gastrointestinal diseases: the RV method and DGESA. The RV method prioritizes genes based
372 on their variance, providing an initial perspective on potential biomarkers by identifying genes with significant
373 variability in expression between diseased and normal samples. DGESA, on the other hand, utilizes the principles of
374 simulated annealing to integrate differential gene expression data, refining the selection process by iteratively
375 optimizing gene sets to maximize their discriminative power.

376 Through a focused case study on EoE and other gastrointestinal diseases like CD and UC, we systematically
377 compare the outcomes of both methods. The RV method initially identifies genes with high variance, offering a broad
378 overview of potential candidates. In contrast, DGESA fine-tunes this selection by incorporating a fitness function that
379 assesses the difference in mean gene expression between diseased and normal states, thus honing in on genes with the
380 most significant impact.

381 Our results reveal a notable intersection between the two methods, with 10 common genes identified in EoE,
382 highlighting their complementary nature and the robustness of the selection process. Further validation analyses
383 demonstrate that 10 out of the 40 final genes identified by DGESA for EoE are confirmed by existing literature,
384 showcase their biological relevance and potential role in disease pathogenesis. Similarly, in the contexts of UC and
385 Crohn's Disease CD, 8 and 7 genes, respectively, from the final 40 identified by DGESA are supported by literature
386 evidence, indicating their significance in these diseases. KRT79, CRISP2, IL36G, SPRR2B, SPRR2D, and SPRR2E
387 are among the ten confirmed genes for EoE. NPDC1, SLC2A4RG, LGALS8, CDKN1A, XAF1, and CYBA are
388 included in the literature-confirmed CD set. TRAF3, BAG6, CCDC80, CDC42SE2, and HSPA9 are among the
389 validated genes in the UC final gene collection.

390 These findings underscore the efficacy of both RV and DGESA in elucidating molecular signatures associated
391 with gastrointestinal diseases. The complementary strengths of the RV method and DGESA provide a robust
392 framework for identifying key genetic contributors to disease, enhancing our understanding of disease mechanisms,
393 and identifying potential therapeutic targets. By integrating these approaches, we can more accurately pinpoint the
394 genes that play pivotal roles in disease development, paving the way for advancements in diagnostics and personalized
395 medicine.

396 **List of ORCID IDs for the authors:** <https://orcid.org/0009-0003-6825-3709>, <https://orcid.org/0009-0007-6196-2253>,
397 <https://orcid.org/0009-0003-1823-3574>, <https://orcid.org/0009-0000-1988-4430>, XXXX, <https://orcid.org/0000-0001-5979-3605>.

398 **Reference**

- 399 1. Cristianini, N., & Hahn, M. W. (2006). Introduction to computational genomics: a case studies approach.
400 Cambridge University Press.
- 401 2. Derrien, T., Estellé, J., Marco Sola, S., Knowles, D. G., Raineri, E., Guigó, R., & Ribeca, P. (2012). Fast
402 computation and applications of genome mappability. *PloS one*, 7(1), e30377.
- 403 3. Hackl, H., Charoentong, P., Finotello, F., & Trajanoski, Z. (2016). Computational genomics tools for
404 dissecting tumour-immune cell interactions. *Nature Reviews Genetics*, 17(8), 441-458.

- 405 4. Emes, R. D., Pirooznia, M., Zou, Q., & Pellegrini, M. (2023). Insights in computational genomics: 2022.
406 *Frontiers in Genetics*, 14, 1256011.
- 407 5. Clough, E., & Barrett, T. (2016). The gene expression omnibus database. *Statistical Genomics: Methods and*
408 *Protocols*, 93-110.
- 409 6. <https://www.ncbi.nlm.nih.gov/>
- 410 7. Barmeyer, C., Schulzke, J. D., & Fromm, M. (2015, June). Claudin-related intestinal diseases. In *Seminars*
411 *in cell & developmental biology* (Vol. 42, pp. 30-38). Academic Press.
- 412 8. Lucas López, R., Grande Burgos, M. J., Gálvez, A., & Pérez Pulido, R. (2017). The human gastrointestinal
413 tract and oral microbiota in inflammatory bowel disease: a state of the science review. *Apmis*, 125(1), 3-10.
- 414 9. Cianferoni, A., & Spergel, J. (2016). Eosinophilic esophagitis: a comprehensive review. *Clinical reviews in*
415 *allergy & immunology*, 50, 159-174.
- 416 10. Dellon, E. S., & Hirano, I. (2018). Epidemiology and natural history of eosinophilic esophagitis.
417 *Gastroenterology*, 154(2), 319-332.
- 418 11. Muir, A., & Falk, G. W. (2021). Eosinophilic esophagitis: a review. *Jama*, 326(13), 1310-1318.
- 419 12. Saito, Y. A., Mitra, N., & Mayer, E. A. (2010). Genetic approaches to functional gastrointestinal disorders.
420 *Gastroenterology*, 138(4), 1276-1285.
- 421 13. Rothenberg, M. E. (2015). Molecular, genetic, and cellular bases for treating eosinophilic esophagitis.
422 *Gastroenterology*, 148(6), 1143-1157.
- 423 14. Kottyan, L. C., & Rothenberg, M. (2017). Genetics of eosinophilic esophagitis. *Mucosal immunology*, 10(3),
424 580-588.
- 425 15. Kottyan, L. C., Parameswaran, S., Weirauch, M. T., Rothenberg, M. E., & Martin, L. J. (2020). The genetic
426 etiology of eosinophilic esophagitis. *Journal of Allergy and Clinical Immunology*, 145(1), 9-15.
- 427 16. Baruah, B., Kumar, T., Das, P., Thakur, B., Sreenivas, V., Ahuja, V., ... & Makharia, G. K. (2017). Prevalence
428 of eosinophilic esophagitis in patients with gastroesophageal reflux symptoms: A cross-sectional study from
429 a tertiary care hospital in North India. *Indian journal of gastroenterology*, 36, 353-360.
- 430 17. Nagarajan, K. V., Krishnamurthy, A. N., Yelsangikar, A., Mallappa, R. B., Bhat, V., Narasimhamurthy, V.
431 M., & Bhat, N. (2023). Does eosinophilic esophagitis exist in India?. *Indian Journal of Gastroenterology*,
432 42(2), 286-291.
- 433 18. Cho JH, Lee D, Park JH, Lee IB (2003) New gene selection method for classification of cancer subtypes
434 considering within class variation. *FEBS Lett* 551(1–3):3–7.
- 435 19. Díaz-Uriarte R, Andrés SA (2006) Gene selection and classification of microarray data using random forest.
436 *BMC Bioinform*.
- 437 20. Yu L, Han Y, Berens ME (2012) Stable gene selection from microarray data via sample weighting.
438 *IEEE/ACM Trans Comput Biol Bioinform* 9(1):262–272
- 439 21. Libbrecht, M. W., & Noble, W. S. (2015). Machine learning applications in genetics and genomics. *Nature*
440 *Reviews Genetics*, 16(6), 321-332.
- 441 22. Kourou, K., Exarchos, T. P., Exarchos, K. P., Karamouzis, M. V., & Fotiadis, D. I. (2015). Machine learning
442 applications in cancer prognosis and prediction. *Computational and structural biotechnology journal*, 13, 8-
443 17.

- 444 23. Tabl, A. A., Alkhateeb, A., ElMaraghy, W., Rueda, L., & Ngom, A. (2019). A machine learning approach
445 for identifying gene biomarkers guiding the treatment of breast cancer. *Frontiers in genetics*, 10, 256.
- 446 24. Biswas, S., Dutta, S., & Acharyya, S. (2019). Identification of disease critical genes using collective meta-
447 heuristic approaches: an application to preeclampsia. *Interdisciplinary Sciences: Computational Life*
448 *Sciences*, 11, 444-459.
- 449 25. Esengönül, M., Marta, A., Beirão, J., Pires, I. M., & Cunha, A. (2022). A systematic review of artificial
450 intelligence applications used for inherited retinal disease management. *Medicina*, 58(4), 504.
- 451 26. Colak, C., Kucukakcali, Z., & Akbulut, S. (2023). Artificial intelligence-based prediction of molecular and
452 genetic markers for hepatitis C-related hepatocellular carcinoma. *Annals of Medicine and Surgery*, 85(10),
453 4674-4682.
- 454 27. Bao, W., Wang, L., Liu, X., & Li, M. (2023). Predicting diagnostic biomarkers associated with immune
455 infiltration in Crohn's disease based on machine learning and bioinformatics. *European Journal of Medical*
456 *Research*, 28(1), 255.
- 457 28. Ridge, P. G., Hoyt, K. B., Boehme, K., Mukherjee, S., Crane, P. K., Haines, J. L., ... & Reitz, C. (2016).
458 Assessment of the genetic variance of late-onset Alzheimer's disease. *Neurobiology of aging*, 41, 200-e13.
- 459 29. Chen, C. H., Kraemer, B. R., & Mochly-Rosen, D. (2022). ALDH2 variance in disease and populations.
460 *Disease Models & Mechanisms*, 15(6), dmm049601.
- 461 30. Wolf, S., Melo, D., Garske, K. M., Pallares, L. F., Lea, A. J., & Ayroles, J. F. (2023). Characterizing the
462 landscape of gene expression variance in humans. *PLoS genetics*, 19(7), e1010833.
- 463 31. Bertsimas, D., & Tsitsiklis, J. (1993). Simulated annealing. *Statistical science*, 8(1), 10-15.
- 464 32. Aarts, E., Korst, J., & Michiels, W. (2005). Simulated annealing. *Search methodologies: introductory*
465 *tutorials in optimization and decision support techniques*, 187-210.
- 466 33. Guilmeau, T., Chouzenoux, E., & Elvira, V. (2021, July). Simulated annealing: A review and a new scheme.
467 In *2021 IEEE Statistical Signal Processing Workshop (SSP)* (pp. 101-105). IEEE.
- 468 34. Koul, N., & Manvi, S. S. (2022). Feature selection from gene expression data using simulated annealing and
469 partial least squares regression coefficients. *Global Transitions Proceedings*, 3(1), 251-256.
- 470 35. Marjit, S., Bhattacharyya, T., Chatterjee, B., & Sarkar, R. (2023). Simulated annealing aided genetic
471 algorithm for gene selection from microarray data. *Computers in Biology and Medicine*, 158, 106854.
- 472 36. Kc, K., Rothenberg, M. E., & Sherrill, J. D. (2015). In vitro model for studying esophageal epithelial
473 differentiation and allergic inflammatory responses identifies keratin involvement in eosinophilic
474 esophagitis. *PloS one*, 10(6), e0127755.
- 475 37. Rochman, M., Azouz, N. P., & Rothenberg, M. E. (2018). Epithelial origin of eosinophilic esophagitis.
476 *Journal of Allergy and Clinical Immunology*, 142(1), 10-23.
- 477 38. Dellon, E. S., Selitsky, S. R., Genta, R. M., Lash, R. H., & Parker, J. S. (2018). Gene expression-
478 phenotype associations in adults with eosinophilic esophagitis. *Digestive and Liver Disease*, 50(8), 804-811.
- 479 39. Qin, X., Liu, M., Zhang, S., Wang, C., & Zhang, T. (2019). The role of IL-36 γ and its regulation in
480 eosinophilic inflammation in allergic rhinitis. *Cytokine*, 117, 84-90.
- 481 40. Morrison, H. A., Hoyt, K. J., Mounzer, C., Ivester, H. M., Barnes, B. H., Sauer, B., ... & Allen, I. C.
482 (2023). Expression profiling identifies key genes and biological functions associated with eosinophilic
483 esophagitis in human patients. *Frontiers in Allergy*, 4.

- 484 41. 39. Ding, J., Garber, J. J., Uchida, A., Lefkovith, A., Carter, G. T., Vimalathas, P., ... & Xavier, R. J. (2024).
485 An esophagus cell atlas reveals dynamic rewiring during active eosinophilic esophagitis and remission.
486 Nature Communications, 15(1), 3344.
- 487 42. Rangel, A. H. D. N., Sales, D. C., Urbano, S. A., GALVÃO, J. G. B., ANDRADE, J. C. D., & Macedo, C.
488 D. S. (2016). Lactose intolerance and cow's milk protein allergy. Food Science and Technology
489 (Campinas), 36(2), 179-187.
- 490 43. Frenkel, S., Bernstein, C. N., Sargent, M., Kuang, Q., Jiang, W., Wei, J., ... & Hu, P. (2019).
491 Genome-wide analysis identifies rare copy number variations associated with inflammatory bowel
492 disease. PLoS One, 14(6), e0217846.
- 493 44. <https://www.ncbi.nlm.nih.gov/gtr/genes/56731/>
- 494 45. Elding, H., Lau, W., Swallow, D. M., & Maniatis, N. (2013). Refinement in localization and
495 identification of gene regions associated with Crohn disease. The American Journal of Human
496 Genetics, 92(1), 107-113.
- 497 46. Gologan, S., Iacob, R., Iancu, D., Iacob, S., Cotruta, B., Vadan, R., ... & Diculescu, M. (2013).
498 Inflammatory gene expression profiles in Crohn's disease and ulcerative colitis: a comparative
499 analysis using a reverse transcriptase multiplex ligation-dependent probe amplification protocol.
500 Journal of Crohn's and Colitis, 7(8), 622-630.
- 501 47. Parackova, Z., Milota, T., Vrabcova, P., Smetanova, J., Svaton, M., Freiburger, T., ... & Sediva, A.
502 (2020). Novel XIAP mutation causing enhanced spontaneous apoptosis and disturbed NOD2
503 signalling in a patient with atypical adult-onset Crohn's disease. Cell death & disease, 11(6), 430.
- 504 48. Serra, E. G., Schwerd, T., Moutsianas, L., Cavounidis, A., Fachal, L., Pandey, S., ... & Anderson,
505 C. A. (2020). Somatic mosaicism and common genetic variation contribute to the risk of very-
506 early-onset inflammatory bowel disease. Nature communications, 11(1), 995.
- 507 49. Shen, J., Qiao, Y. Q., Ran, Z. H., & Wang, T. R. (2013). Up-regulation and pre-activation of TRAF3
508 and TRAF5 in inflammatory bowel disease. International journal of medical sciences, 10(2), 156.
- 509 50. Di'Narzo, A. F., Houten, S. M., Kosoy, R., Huang, R., Vaz, F. M., Hou, R., ... & Argmann, C.
510 (2022). Integrative analysis of the inflammatory bowel disease serum metabolome improves our
511 understanding of genetic etiology and points to novel putative therapeutic targets.
512 Gastroenterology, 162(3), 828-843.
- 513 51. Hao-Hua, W. A. N. G., Wan-Ying, L. U. O., Min, L. I. N., Xiao-Jing, L. I., Xiang, G. D., & D
514 TRIGANTI, S. (2021). Plasma asprosin, CCDC80 and ANGPTL4 levels are associated with
515 metabolic and cardiovascular risk in patients with inflammatory bowel disease. Physiological
516 Research, 70(2), 203.
- 517 52. Mo, A., Nagpal, S., Gettler, K., Haritunians, T., Giri, M., Haberman, Y., ... & Gibson, G. (2021).
518 Stratification of risk of progression to colectomy in ulcerative colitis via measured and predicted
519 gene expression. The American Journal of Human Genetics, 108(9), 1765-1779.
- 520 53. Jang, S., Jang, S., Ko, J., Bae, J. E., Hyung, H., Park, J. Y., ... & Ryoo, Z. Y. (2024). HSPA9
521 reduction exacerbates symptoms and cell death in DSS-Induced inflammatory colitis. Scientific
522 Reports, 14(1), 5908.

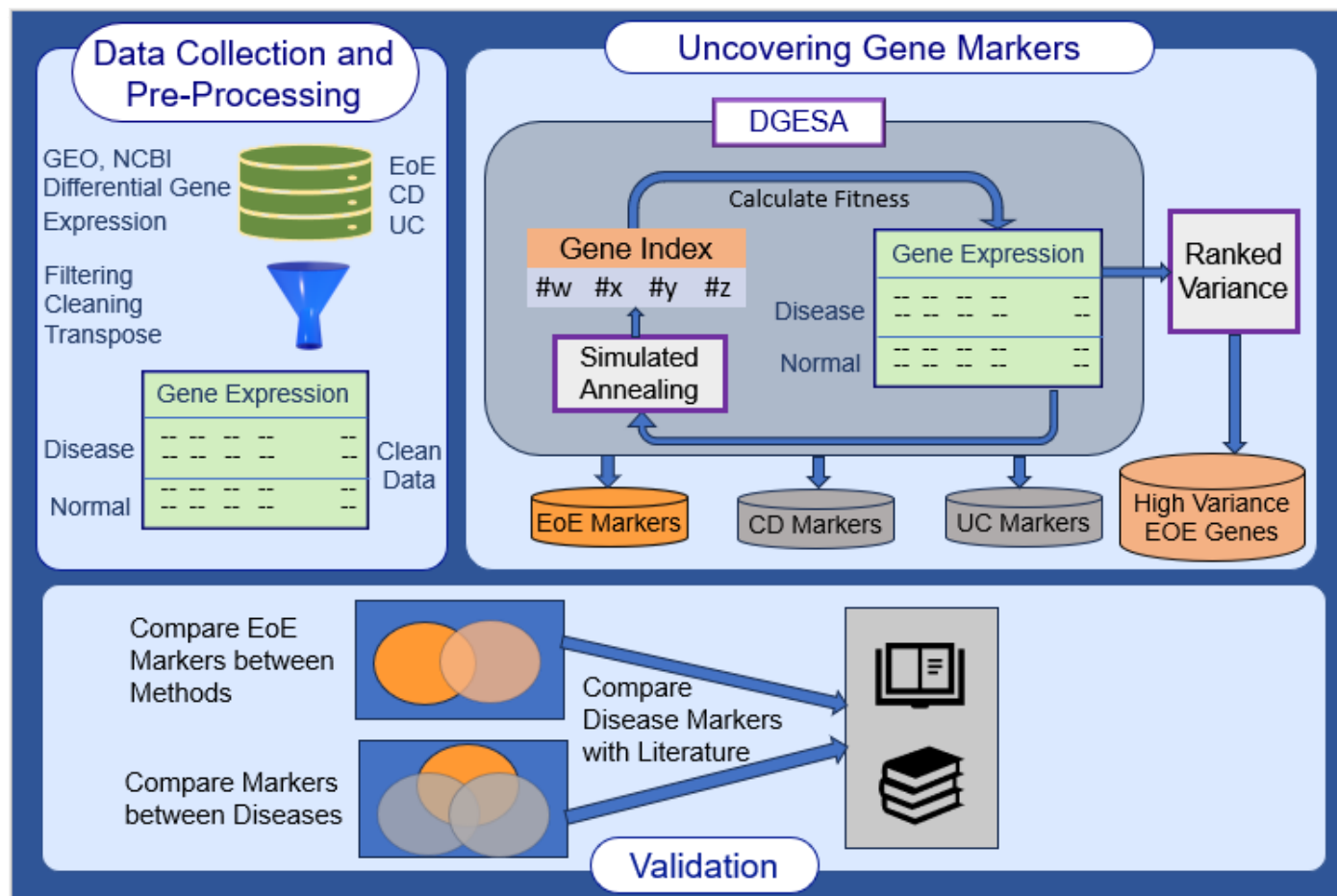


Figure 1: Overview of the Study Workflow. This figure illustrates the comprehensive workflow of the study, beginning with data acquisition and preprocessing steps that include cleaning, normalization, and transposition to produce the gene expression matrix (samples as rows, genes as columns). The gene marker identification methodologies encompassing RV and DGESA, are applied to datasets of EOE, CD, and UC to identify disease-specific gene markers. The final step involves the validation process, which compares the effectiveness of the two methods across the three diseases and validates the identified markers against existing literature.

DGESA

```
1. Initialize temperature T,  
   Iterations P and Cooling rate C;  
2. Choose initial candidate solution  
   s with random gene index;  
3. Calculate fitness f(s) using  
   Eq. (1);  
4. Repeat  
   A. for i = 1 to P do  
       ○ Randomly select  $s' \in N(s)$ ;  
         //  $N(s)$ : Neighbor of s  
       ○ if  $f(s') \geq f(s)$  then  
            $s \leftarrow s'$ ;  
       ○ else  
            $s \leftarrow s'$  with  
             probability  $e^{-(f(s')-f(s))/T}$ ;  
       ○ end if  
   B. end for  
   C.  $T = T \times C$ ;  
5. Until stopping criteria not met  
6. end
```

Figure 2. DGESA Algorithm: This figure depicts the details of the DGESA algorithm. The process involves simulated annealing with a specially designed objective function to effectively select gene markers from differential gene expression data.

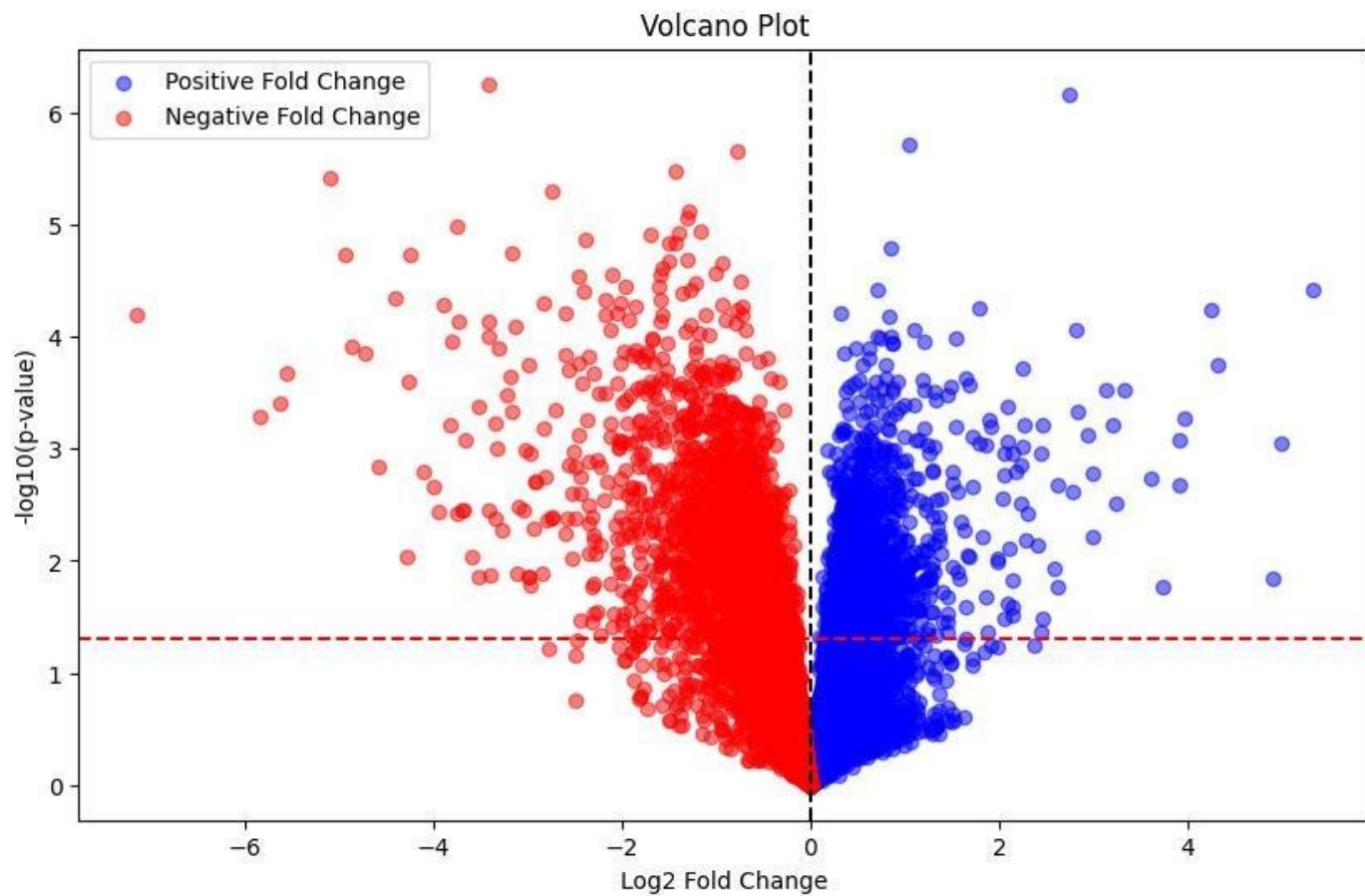


Figure 3. Volcano Plot of Gene Expression Data for EOE: This figure presents a volcano plot illustrating the gene expression data for EOE. The x-axis represents the log₂ fold change, while the y-axis denotes the -log₁₀ P-value. The plot highlights a greater number of genes with negative fold changes compared to those with positive fold changes, indicating differential gene expression patterns in EOE.

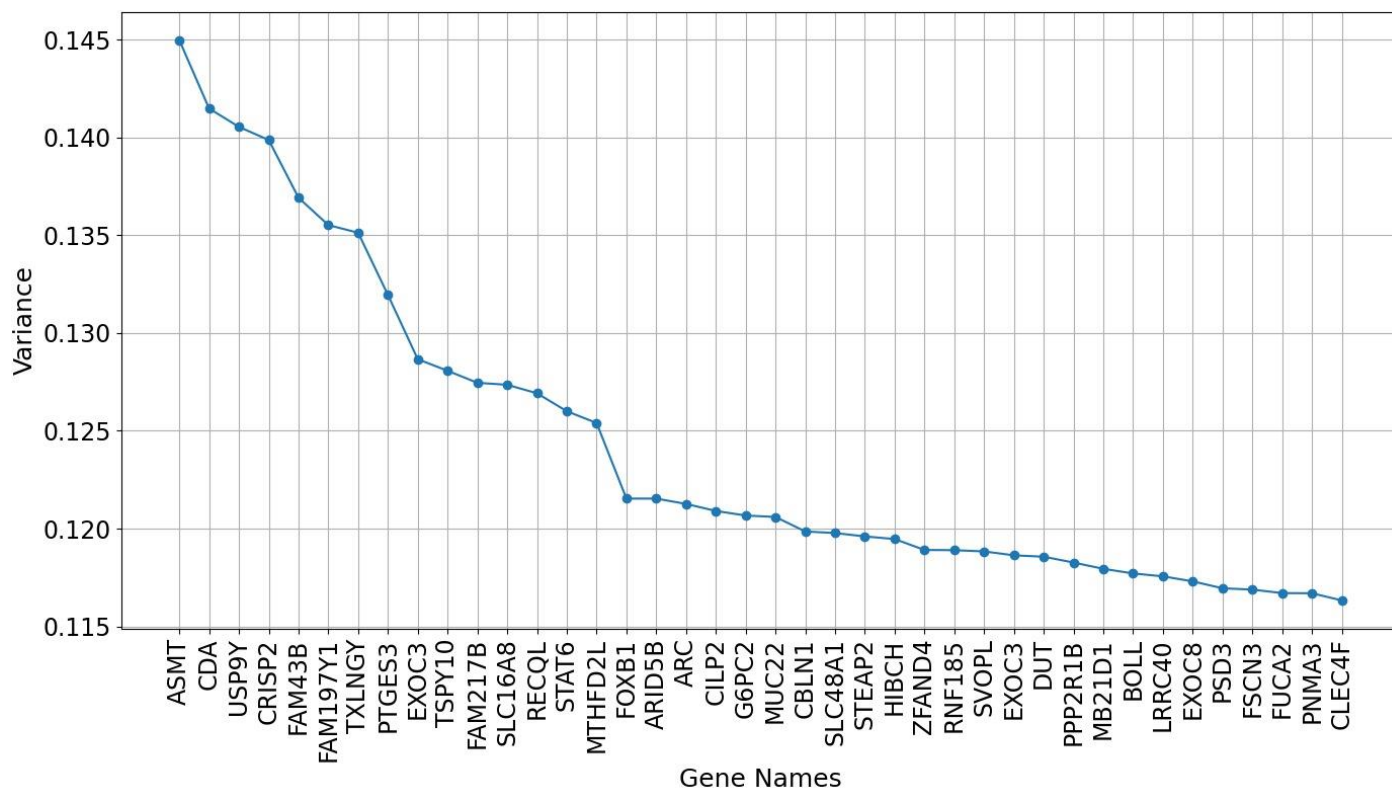


Figure 4. Variance of Top 40 Genes in EOE Gene Expression Data: This figure shows a line plot of the variance of the top 40 genes in the EOE gene expression dataset. The x-axis represents the gene names with the highest variance, and the y-axis indicates the variance values. The plot reveals that the initial few genes exhibit very high variance, which gradually flattens down for the subsequent genes.

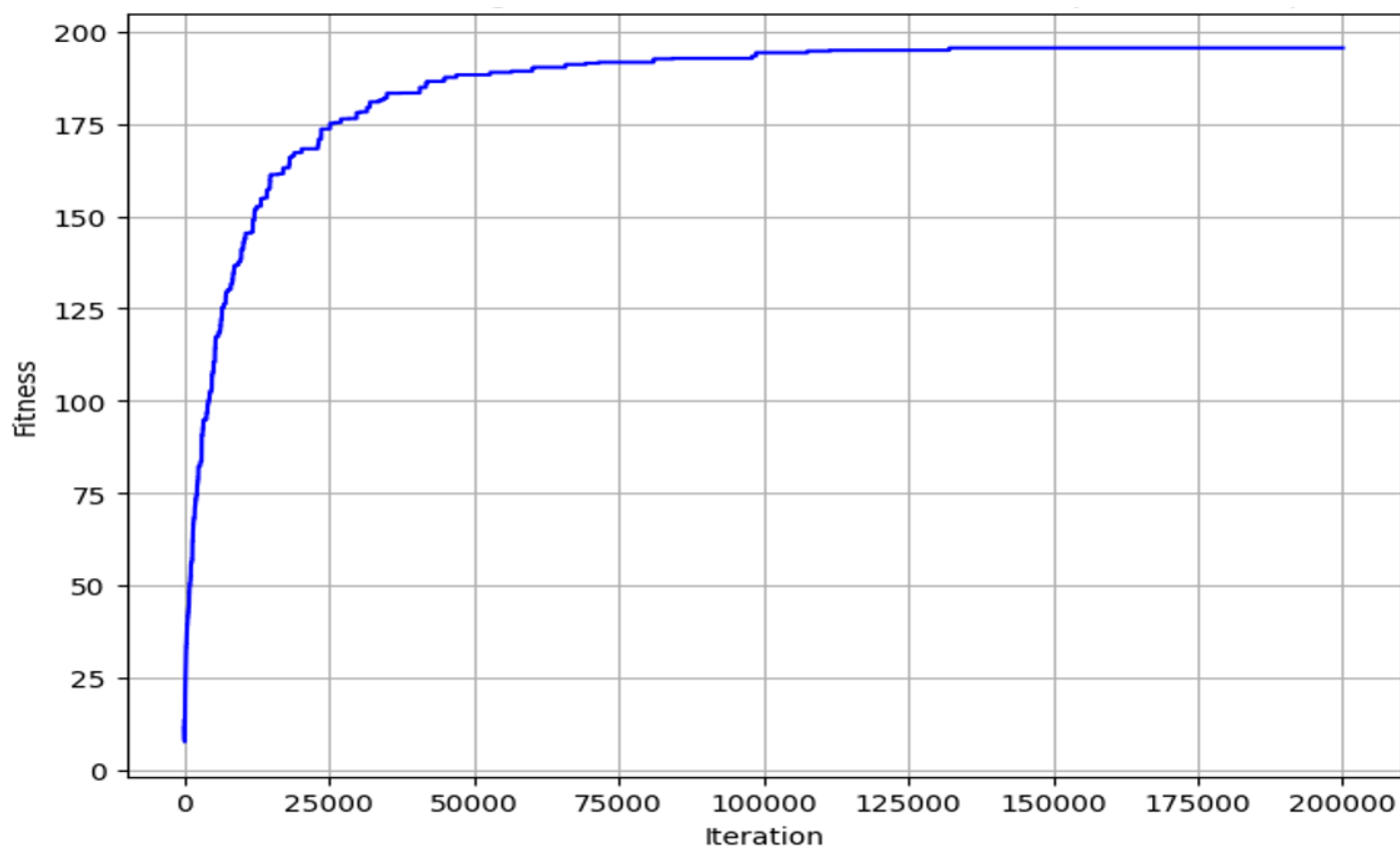


Figure 5. Optimization Process for Identifying 40 Discriminative Genes in EOE: This figure illustrates the optimization process aimed at selecting a set of 40 genes that collectively maximize the discrimination between EOE and control samples in the gene expression dataset. The curve shows the convergence of the optimization process, stabilizing near 130,000 epochs.

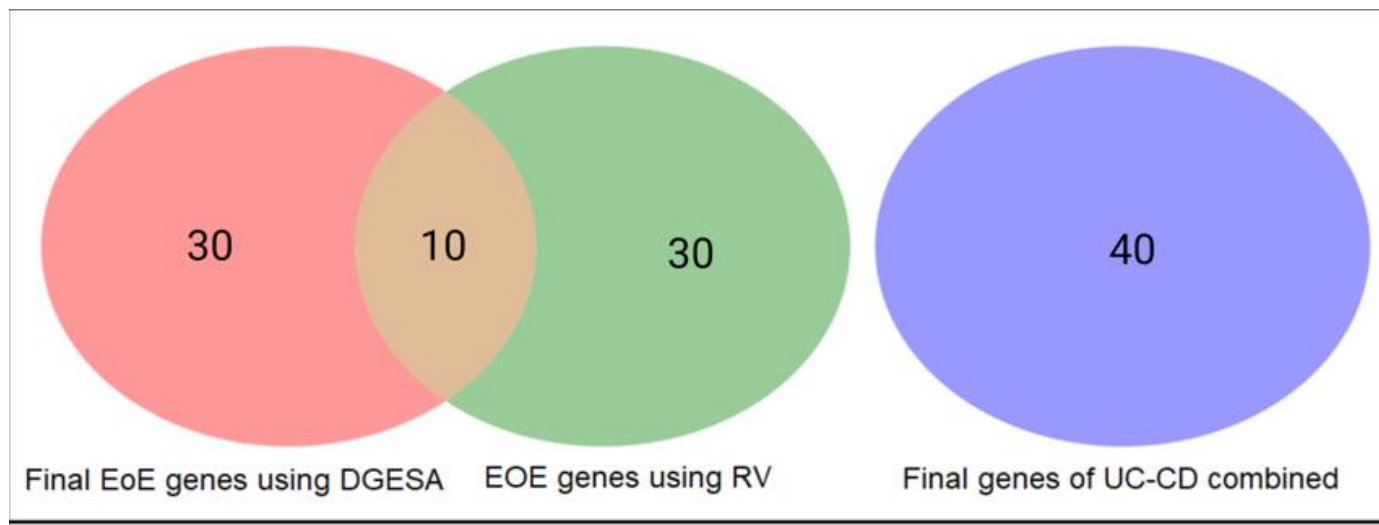


Figure 6. Venn diagram of Common Genes Identified by RV and DGESA Methods in EoE and UC-CD Combined Datasets. This figure shows a Venn diagram comparing the common genes identified by Ranked Variance (RV) and Differential Gene Expression Based Simulated Annealing (DGESA) methods in the EoE gene expression dataset (left side) and the common genes identified using DGESA in the UC-CD combined dataset. The diagram reveals 10 common genes between the RV and DGESA methods for EoE, but no overlapping genes between the EoE dataset and the UC-CD combined dataset.