

1 A Comprehensive Targeted Panel of 295 Genes: Unveiling Key
2 Disease Initiating and Transformative Biomarkers in Multiple
3 Myeloma

4 Vivek Ruhela^{a,b}, Rupin Oberoi^a, Ritu Gupta^{c, **}, Anubha Gupta^{a, **}

5 ^aSBILab, Deptt. of ECE & Centre of Excellence in Healthcare, Indraprastha Institute of Information
6 Technology-Delhi (IIIT-D), India

7 ^bDept. of Computational Biology, Indraprastha Institute of Information Technology-Delhi (IIIT-D), India

8 ^cLaboratory Oncology Unit, Dr.B.R.A.IRCH, All India Institute of Medical Sciences, Delhi, India

9 **Abstract**

Multiple myeloma (MM) is a haematological cancer that evolves from the benign precursor stage termed monoclonal gammopathy of undetermined significance (MGUS). Understanding the pivotal biomarkers, genomic events, and gene interactions distinguishing MM from MGUS can significantly contribute to early detection and an improved understanding of MM's pathogenesis. This study presents a curated, comprehensive, targeted sequencing panel focusing on 295 MM-relevant genes and employing clinically oriented NGS-targeted sequencing approaches. To identify these genes, an innovative AI-powered attention model, the *Bio-Inspired Graph Network Learning-based Gene-Gene Interaction* (BIO-DGI) model, was devised for identifying *Disease-Initiating* and *Disease-Transformative* genes using the genomic profiles of MM and MGUS samples. The BIO-DGI model leverages gene interactions from nine protein-protein interaction (PPI) networks and analyzes the genomic features from 1154 MM and 61 MGUS samples. The proposed model outperformed baseline machine learning (ML) and deep learning (DL) models on quantitative performance metrics. Additionally, the BIO-DGI model identified the highest number of MM-relevant genes in the post-hoc analysis, demonstrating its superior qualitative performance. Pathway analysis highlighted the significance of top-ranked genes, emphasizing their role in MM-related pathways. Encompassing 9417 coding regions with a length of 2.630 Mb, the 295-gene panel exhibited superior performance, surpassing previously published panels in detecting genomic disease-initiating and disease-transformative events. The panel also revealed highly influential genes and their interactions within MM gene communities. Clinical relevance was confirmed through a two-fold univariate

*Corresponding author

**Corresponding author

Email addresses: drritugupta@gmail.com (Ritu Gupta), anubha@iiitd.ac.in (Anubha Gupta)

survival analysis, affirming the significance of the proposed gene panel in understanding disease progression. The study's findings offer crucial insights into essential gene biomarkers and interactions, shaping our understanding of MM pathophysiology.

¹⁰ *Keywords:* AI in Cancer, Haematological malignancy, Multiple Myeloma, MGUS, Genomic
¹¹ Aberrations, ShAP, FAMD, PPI

¹² **1. Introduction**

¹³ Multiple Myeloma (MM) is a haematological malignancy characterized by the clonal proliferation
¹⁴ of plasma cells within the bone marrow. Monoclonal Gammopathy of Undetermined Significance
¹⁵ (MGUS) and MM are both plasma cell disorders, representing distinct stages along the disease
¹⁶ progression continuum. MM entails malignant plasma cell proliferation, organ damage, and clin-
¹⁷ ical symptoms, whereas MGUS is a precursor condition with no apparent clinical manifestations.
¹⁸ Progression from MGUS to MM occurs at a rate of 1% per year; thus, all MGUS patients do not
¹⁹ transition to overt MM during their lifetime [1]. In this context, identifying MGUS individuals
²⁰ likely to progress to MM is crucial for timely intervention and improved outcomes. The clinical
²¹ distinction between the two conditions primarily relies on tumour load, reflected in monoclonal
²² proteinemia, percentage bone marrow plasma cell infiltration, and end-organ damage [2]. This em-
²³ phasizes the need to delve into genomic markers and gene-gene interactions to enhance diagnostic
²⁴ accuracy.

²⁵ The advanced genomic profiling techniques like whole-exome sequencing (WES) and whole-
²⁶ genome sequencing (WGS) enable a thorough examination of genomic aberrations in MM and
²⁷ MGUS. They have proven crucial in identifying key events, including copy number variations
²⁸ (CNVs) and structural variations (SVs) [3]. Numerous MM-related genomic studies have re-
²⁹ ported significant CNVs, such as del(1p), gain(1q), del(13q), and del(17p), alongside key SVs
³⁰ like translocation involving IgH (e.g. t(4;14), t(11;14), t(14;16), t(8;14)), MYC rearrangement
³¹ like MYC-IGL, MYC-IGK rearrangements, shedding light on their association with MM prognosis
³² [4, 5, 6, 7, 8, 9, 10]. Moreover, recent studies have highlighted the impact of minor genomic alter-
³³ ations on MM patients' clinical outcomes [11, 12, 13, 14]. Recently, biallelic alterations in TP53

34 and DIS3 gene have been reported as high-risk markers in MM [15].

35 Integrating protein-protein interactions (PPI) with WES and WGS variant profiles can provide
36 crucial insights into genomic biomarkers essential for MGUS to MM progression. Distinguishing MM
37 from MGUS via gene-gene interactions can revolutionize clinical approaches, aiding early detection.
38 This can enable proactive monitoring and intervention for high-risk MGUS patients while sparing
39 low-risk MM cases from aggressive treatments [16].

40 Several targeted sequencing panels have been devised to comprehensively profile the genomics
41 complexity of MM [17, 18, 19]. These panels encompass critical genomic aberrations related to MM.
42 For instance, a 26-gene panel focused on prevalent mutations in previously published MM-relevant
43 genes [17], but lacked validation for SVs. Similarly, another panel of 182 genes validated for single
44 nucleotide variants (SNVs), CNVs, and specific translocations (related to IGH only) in previously
45 published MM-relevant genes [20]. A more extensive 228-gene panel covered various alterations,
46 including SNVs, CNVs, and translocations involving IgH and MYC genes [18]. In a similar quest for
47 comprehensive genomic profiling of MM, a 47-gene panel was crafted, encompassing dysregulated
48 and frequently mutated genes in MM and those targeted by common therapies, validated for SNVs
49 only [19]. Lastly, the largest gene panel of 465 genes was designed and validated for MM-related
50 SNVs, CNVs, and translocations related to the IGH gene only [21]. However, these panels were
51 designed using only MM samples and hence, lacked markers and interactions distinguishing MM
52 from MGUS that can give insights into MM pathogenesis.

53 Machine learning (ML) and deep learning (DL) advancements have revolutionized bioinformatics,
54 enabling precise biomarker discovery for early disease detection. Researchers utilize these tools
55 to predict Protein-Protein Interactions (PPIs) and unravel crucial gene interactions in cancer. Notably,
56 models like DeepPPI (Deep neural networks for Protein-Protein Interactions prediction)
57 predict gene interactions based on shared protein descriptors [22]. ML and DL-driven approaches
58 also facilitate inferring semantic similarity of gene ontology terms using PPIs [23, 24, 25]. Despite
59 these strides, no computational model was tailored to identify pivotal biomarkers and gene
60 interactions distinguishing MM from MGUS.

61 The fusion of genomic mutation profiles with the Protein-Protein Interaction (PPI) data re-

62 mains an underexplored domain. Graph Deep Learning (GDL) emerges as a potent tool in ge-
63 nomics, promising profound insights from intricate biological graph data structures. Genomic data
64 inherently manifests a graph-like structure, where nodes embody biological entities (genes, pro-
65 teins) and edges depict relationships. Graph Convolutional Networks (GCNs) play a pivotal role
66 in this realm. GCNs are instrumental in processing and analyzing graph-structured genomic data,
67 where biological entities and their relationships are effectively modelled as a graph, particularly for
68 disease classification [26]. Our previous work seamlessly integrated exonic mutation profiles with
69 PPI to unveil key distinguishing biomarkers of MM and MGUS through the bio-inspired BDL-SP
70 model [27]. We pursued heightened precision in identifying biomarkers and gene interactions by
71 integrating gene interactions from nine diverse PPI databases.

72 Motivated to bridge this gap, we envisioned a targeted sequencing panel for a thorough genomic
73 profiling of MM, aiming to capture the unique characteristics of MM and MGUS. To address this
74 challenge, we introduced a novel AI-powered attention model: *Bio-inspired Graph network learning*
75 *based on directed gene-gene interactions (BIO-DGI)*, employing graph network learning to discern
76 differentiating biomarkers and gene-gene interactions in MM and MGUS. In this proposed model,
77 we have integrated bio-inspired learning, utilizing the topological information gathered from nine
78 PPI networks and exomic mutational profiles. This empowered the BIO-DGI model to rank genes
79 and genomic features based on their role in disease progression more efficiently, with fewer graph
80 convolution network layers and multi-head attention modules compared to traditional machine
81 learning (ML) or deep learning (DL) models that relied solely on exomic mutational profiles. The
82 BIO-DGI model also helped us in identifying *Disease-Initiating* and *Disease-Transformative* genes
83 that can aid into understanding MM pathogenesis. This model outperformed exhaustive bench-
84 marking against several baseline ML and DL models, including both quantitative and qualitative
85 evaluations.

86 We further delved deeper and identified five distinct gene communities using the Leiden algo-
87 rithm [28] by utilizing the adjacency matrices derived from the five trained BIO-DGI classifiers.
88 This analysis shed light on the influential genes within these communities, quantified through Katz
89 centrality scores [29]. Importantly, we have highlighted genes that were observed to be located in

90 the central position in these gene communities and hence, might be playing a significant role in
91 MM pathogenesis. We analyzed various variant profiles, including SNVs, CNVs, SVs, and Loss of
92 Function (LOF) mutations. This detailed investigation, along with the post-hoc analysis via ShAP
93 (SHapley Additive exPlanations) algorithm [30] that identified top-ranked genes and genomic fea-
94 tures, led to the design of a clinically tailored 295-gene panel, aiming for comprehensive genomic
95 profiling of Multiple Myeloma (MM).

96 Pathway enrichment analysis of these 295 genes revealed enriched MM-related pathways, strongly
97 underscoring the pivotal role of these genes in MM pathogenesis. This discovery, along with the
98 survival analysis corresponding to these genes, underscores the clinical relevance and potential of
99 the targeted sequencing panel designed for comprehensive genomic profiling in MM.

100 **2. Materials & Methods**

101 *2.1. Whole-exome sequencing datasets of MM and MGUS patients*

102 In this study, we included tumour-normal pairs of bone marrow (BM) samples from an MM
103 cohort of 1154 samples and an MGUS cohort of 61 samples sourced from three global repositories
104 of whole-exome sequencing (WES) data. For the MM cohort, 1072 samples were acquired through
105 authorized access to the MMRF dbGaP study (phs000748; phs000348), predominantly comprising
106 American population samples [31]. We also downloaded processed MMRF datasets (version IA12)
107 containing CNVs, SVs, and clinical data from the MMRF Research Gateway. Additionally, we
108 included 82 MM samples from an AIIMS dataset representing the Indian population. In the MGUS
109 cohort, we incorporated 28 MGUS samples from the AIIMS dataset and 33 samples from the
110 European Genome-phenome Archive (EGA) data.

111 *2.2. Computational tools and software used for data analysis*

112 We utilized Python computational tools (version 3.9.13) for WES data analysis and visual-
113 ization. For training all deep learning (DL) models in this study, we employed PyTorch (version
114 1.12.0+cu113) [32]. Additionally, survival analysis was conducted using the statistical programming
115 language R (version 4.3.1) with the “survival” package [33] (version 3.5.5).

116 2.3. Whole exome sequencing Data: Identification of Significantly altered genes

117 The WES data obtained from AIIMS and EGA contained the raw fastq files, and the MMRF
118 dataset contained the processed VCF (Variant Call Format) files. The computational workflow
119 for the SNV identification, genomic annotation of SNVs, SNV filtration and grouping, and the
120 identification of significantly altered genes were taken from our previous study [27]. Briefly, raw
121 fastq files from AIIMS and EGA datasets were processed using the standard exome sequencing
122 pipeline [34]. Similar to the MMRF data, the SNVs in AIIMS and EGA WES data were extracted
123 using MuSE [35], Mutect2 [36], VarScan2 [37], and Somatic-Sniper [38] variant callers. The SNVs
124 in AIIMS, EGA and MMRF datasets were annotated using the ANNOVAR database [39]. The
125 annotated SNVs were categorized into three categories based on their functional significance, i.e.
126 synonymous SNVs, non-synonymous SNVs and other SNVs. The benign SNVs were filtered out
127 using FATHMM-XF [40]. Lastly, the annotated SNVs were pooled for MM and MGUS cohorts
128 separately and analyzed for identifying significantly altered genes using the ‘dndscv’ tool [41].
129 The union of significantly altered genes from all four variant callers for the MM cohort of 1154
130 samples and the MGUS cohort of 61 samples led to 617 and 362 genes, respectively. Union of these
131 significantly altered genes of MM and MGUS cohorts yielded a set of 824 genes.

132 For each of these 824 genes, the corresponding protein-protein interactions (PPIs) were extracted
133 from the nine PPI databases (BioGrid [42], Bio-Plex [43], FunCoup [44], HIPPIE [45], HumanNet
134 [46], IntegratedAssociationCorrNet [47], ProteomHD [48], Reactome [49], and STRING [50]) and
135 consolidated to form a merged adjacency matrix. These interactions were then combined to gen-
136 erate a consolidated adjacency matrix, where the criterion for consolidation was the presence of
137 interactions in at least one PPI dataset, denoted as 1 for present and 0 for absent. A total of 26
138 genes lacked interactions with other significantly altered genes and hence, were excluded resulting
139 in a final set of 798 genes (union of 351 genes from the MGUS cohort and 598 genes from the
140 MM cohort) that were used to construct the merged adjacency matrix (Table-S1, Supplemen-
141 tary File-1). Besides the PPI interactions, we extracted 26 genomic features (Table-S3, Supplemen-
142 tary File-3) that included the total number of synonymous SNVs (including UTR3 and UTR5
143 SNVs), non-synonymous SNVs (including start loss, stop loss, stop gain, exonic, ncRNA-exonic,

144 splicing, frameshift insertion, and frameshift deletion SNVs), and other SNVs (non-frameshift inser-
145 tion/deletion/substitution, intronic, intergenic, ncRNA-intronic, upstream, downstream, unknown,
146 and ncRNA-splicing SNVs), distributive statistics (median and standard deviation) of variant al-
147 lele frequency (VAF), allele depth (AD), and four variant conservation scores (GERP [51], PhyloP
148 [52], PhastCons [53], and Mutation assessor [54]). The details of the pre-processing workflow are
149 available in Supplementary File-2.

150 *2.4. Proposed Directed Gene-Gene Interaction Learning in Biological Network (BIO-DGI)*

151 This study is aimed towards identifying potential driver genes and uncover essential gene-gene
152 interactions responsible for the progression from MGUS to MM. We introduce an innovative Graph
153 Convolutional Network (GCN)-based attention model named “Bio-inspired network learning based
154 on directed gene-gene interactions” (BIO-DGI). The BIO-DGI model, depicted in Figure-1, har-
155 nesses the power of GCN to grasp pivotal gene-gene interactions and forecast potential driver genes.

156 We supplied two essential inputs to empower the BIO-DGI model: 1) an undirected PPI network
157 adjacency matrix sourced from PPI interaction databases and 2) the feature matrix derived from
158 the data. Two versions of the PPI network adjacency matrix were considered in this study. The
159 first version involved extracting PPI interactions solely from the STRING database, serving as the
160 basis for training the vanilla BIO-DGI model, denoted as BIO-DGI (PPI-STRING). In the second
161 version, we merged the PPI network adjacency matrix from nine distinct PPI databases. This
162 merged adjacency matrix was then utilized for training purposes of the model, referred to as BIO-
163 DGI (PPI9). In both versions of the adjacency matrix, each node corresponded to a significantly
164 altered gene, while the links represented interactions between these genes. Additionally, each node
165 was equipped with a feature vector of length 26, as illustrated in Figure-1. Consequently, the PPI
166 network comprising of 798 significantly altered genes, each associated with a feature vector of length
167 26, was integrated into the input layer of the BIO-DGI model.

168 The BIO-DGI model architecture contains 1) a multi-head attention module and 2) a GCN
169 Module. The multi-head attention modules contain three attention units to learn gene-gene inter-
170 actions, followed by an attention consensus module for taking the consensus of all three attention

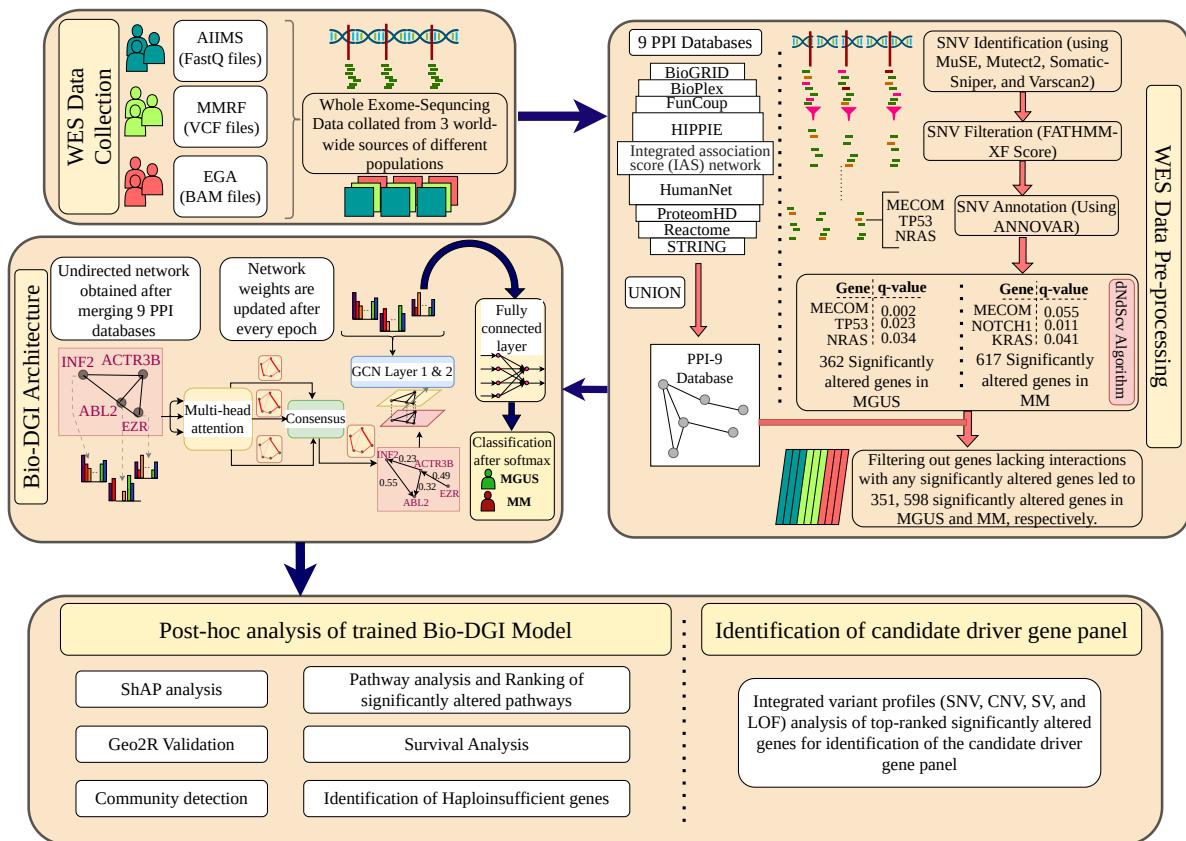


Figure 1: Infographic representation of proposed AI-based bio-inspired BIO-DGI model and post-hoc analysis for identifying pivotal genomic biomarkers that can distinguish MM from MGUS. In the proposed AI-based workflow, the BAM files sourced from EGA and AIIMS datasets, along with VCF files from the MMRF dataset, undergo processing to identify 798 notably altered genes utilizing the dndscv tool (as illustrated in the WES Data pre-processing block). Subsequently, interactions among these 798 genes are elucidated employing PPI networks from nine PPI databases (BioGRID, BioPlex, FunCoup, HIPPIE, IAS network, HumanNet, ProteomHD, Reactome, and STRING). A network is constructed with nodes representing the significantly altered genes and edges denoting interactions obtained after merging interactions from the nine above-mentioned PPI databases. Each node has 26 genomic features (Table-S3, Supplementary File-3) specific to its corresponding gene. The architecture of the BIO-DGI model contains a multi-head attention unit and a GCN layer followed by a fully connected layer. The feature matrix and adjacency matrix are provided as input to the BIO-DGI model. The multi-head attention unit in the BIO-DGI model updates the weights of gene interactions in the adjacency matrix, which are then integrated with the sample feature matrix to gain insights on distinguishing biomarkers that can differentiate MM from MGUS. The output of the fully connected layer is converted into the classification probabilities using the softmax activation function. Consequently, the WES data of each subject is analyzed, and feature vectors for all 798 genes are derived. These feature vectors, in conjunction with the subjects' MM/MGUS target class label, constitute the input for supervised training of the GCN. Following learning the BIO-DGI model for distinguishing MGUS from MM, the top genomic features and significantly altered signalling pathways are extracted utilizing the ShAP algorithm and cross-referencing with the Enrichr pathway database.

171 unit weights to get the updated learned adjacency matrix. The multi-head attention module aimed
172 to learn and update the adjacency matrix to get a weighted PPI adjacency matrix. Similarly, in

173 the GCN module, the input layer is followed by one hidden layer of GCN that is further followed by
174 one fully connected layer of 798 neurons to 2 neurons, giving output through log-softmax activation
175 function for sample class classification (MM vs MGUS).

176 Our study had 95% MM samples and 5% MGUS samples, which made the data highly im-
177 balanced. Hence, the BIO-DGI model was trained using a cost-sensitive negative log-loss (NLL)
178 function to account for the data imbalance. The BIO-DGI model was trained using a five-fold cross-
179 validation technique that led to the training of five best-performing classifiers. All five classifiers
180 with learned adjacency matrices were saved for further post-hoc analysis. We used the ShAP algo-
181 rithm for post-hoc analysis of BIO-DGI model classifiers to get top-performing genes and genomic
182 features that were further used for pathway enrichment analysis, gene-community identification and
183 candidate driver gene panel. The setting of layers, hyperparameters used to train the BIO-DGI
184 model, and mathematical description of the BIO-DGI model are available in Supplementary File-2.

185 *2.5. Quantitative benchmarking of BIO-DGI model with traditional Machine learning classifiers*

186 In our quantitative benchmarking analysis, we conducted a comprehensive comparison of the
187 BIO-DGI (PPI9) model involving three key performance metrics: balanced accuracy, area under the
188 curve (AUC), and area under the precision and recall curve (AUPRC). This evaluation encompassed
189 the five-fold cross-validation of six established baseline cost-sensitive machine learning models:
190 random forest, decision tree, logistic regression, XGBoost, CatBoost, and SVM from the scikit-learn
191 library [55]. Further, we also included two cost-sensitive DL models: BDL-SP and BIO-DGI (PPI-
192 STRING) models for quantitative benchmarking. We incorporated a tailored cost-sensitive loss
193 function to enhance the models' sensitivity to class imbalance. This function implements weighted
194 penalization for sample misclassifications, with the weighting being directly proportional to the
195 class imbalance ratio. This strategic implementation of weighted penalization ensures unbiased
196 learning outcomes for major and minor classes, fostering a more equitable predictive capability.

197 *2.6. Identification of gene communities using learned adjacency matrices of BIO-DGI*

198 We employed a five-fold cross-validation training strategy to our proposed BIO-DGI (PPI9)
199 model. Subsequent to training the model, we retained the learned adjacency matrix from each clas-

200 sifier, yielding five distinct learned adjacency matrices. We individually identified gene communities
201 from these adjacency matrices using the Leiden algorithm. This process yielded 5, 5, 6, 5, and 6
202 gene communities across the learned adjacency matrices. From the communities extracted from
203 each adjacency matrix, the top 3 communities were selected based on the number of OGs, TSGs,
204 ODGs, and AGs. Consequently, we retained the weights of only those genes in the adjacency matrix
205 that were a part of these top 3 communities, while the links of rest of the genes were dropped by
206 assigning zero weight. We called this modified adjacency matrix as the *Gene Community Adjacency*
207 (GCA) matrix. This process was carried out for all the five learned adjacency matrices (one matrix
208 learned from the training of each fold classifier). Next, we computed the mean GCA matrix by
209 calculating the mean weight of a gene across all GCAs. Lastly, we applied the Leiden algorithm for
210 community detection on the mean GCA matrix to obtain the consolidated gene communities.

211 *2.7. Qualitative application-aware post-hoc benchmarking of BIO-DGI model*

212 The ShAP (SHapley Additive exPlanations) algorithm is a powerful tool for gauging the sig-
213 nificance of attributes in a model's predictions. It achieves this by assigning scores to attributes
214 based on their individual contributions. In this context, ShAP played a pivotal role in enhancing
215 the post-hoc explainability of the BIO-DGI (PPI9) model. This process unearthed the most influ-
216 ential genomic features and the genes that experienced significant alterations, both at the cohort
217 level (MGUS or MM) and at the level of individual samples. The ShAP algorithm was applied to
218 each trained classifier obtained after a rigorous five-fold cross validation was carried out during the
219 model's training phase. This enabled the identification of significant genomic attributes (genes and
220 features) for every sample. It is important to note that a ShAP score can have both positive and
221 negative values, wherein a positive ShAP score for a specific attribute highlights its contribution
222 to the model's prediction for the MM class (positive class). Conversely, a negative score indicates
223 its role in the model's prediction for the MGUS class (negative class). Consequently, the magni-
224 tude of the ShAP score directly correlates with the attribute's impact on the model's positive class
225 outcome. Furthermore, the extraction of ShAP interpretability was limited to samples correctly
226 predicted by at least one of the five classifiers. This approach ensured a robust basis for deriving

227 insights through ShAP analysis.

228 We estimated the best ShAP scores on a per-sample basis: 1) for all 798 significantly altered
229 genes and 2) for all 26 genomic features. To this end, for each sample in the MM and MGUS
230 cohort, class predictions were taken from all the five trained classifiers of the BIO-DGI (PPI9)
231 model. Next, the inference of the ShAP algorithm was taken for only the classifiers that made
232 correct predictions for that sample. ShAP scores were collected both at the classifier and sample
233 levels for all genomic attributes. Next, for all the significantly altered genes, the ShAP scores of the
234 26 genomic features were grouped by their positive and negative signs. The best ShAP score for
235 each gene was determined by comparing the absolute values of these grouped scores, considering the
236 highest absolute value among all classifiers as the best possible score. Similarly, for each genomic
237 feature, the ShAP scores of all 798 genes were grouped and assessed in a similar manner, resulting in
238 the best ShAP score. Following this process, the most highly ranked genes and genomic attributes
239 were identified at both the cohort and sample levels.

240 We extended our analysis by comparing the BIO-DGI (PPI9) model's top-ranked significantly
241 altered genes with those reported in previous studies, aiming to identify genes previously reported
242 to be associated with disease progression or suppression. We validated and analyzed our model
243 using information from multiple databases such as OncoKB [56], IntoGen [57], COSMIC [58], and
244 TargetDB [59] at the gene level. For model validation, we extracted 1064 cancer genes from the
245 OncoKB database for oncogenes and tumour-suppressor genes. From the COSMIC database, we
246 utilized 318 oncogenes and 320 tumor-suppressor genes.

247 We utilized the IntoGen database (<https://www.intogen.org/>) and MM-related studies [12, 60]
248 to compile a catalogue of MM driver genes. Additionally, 180 actionable genes from COSMIC
249 and 135 from TargetDB helped infer actionable genes. Using the above information, we regrouped
250 our top-ranked significantly altered genes into Oncogenes (OGs), Tumor-Suppressor genes (TSGs),
251 Onco-driver genes (ODGs), and Actionable genes (AGs). This comprehensive approach facilitated
252 a thorough exploration of genomic features in the post-hoc interpretability analysis of the BIO-
253 DGI (PPI9) model, providing valuable insights into their roles in disease contexts. Furthermore,
254 we introduced another way of understanding the disease pathogenesis by assessing whether a gene

255 was observed to be significantly altered in MM or MGUS or both. Genes found to be significantly
256 mutated in only MM (and not in MGUS) were designated as *Disease-Transformative* genes, while
257 those significantly altered in both MM and MGUS were labelled as *Disease-Initiating* genes. This
258 new way of categorization deepened our understanding of these genes, shedding light on their bio-
259 logical functions and specific relevance to MM and MGUS. This comprehensive approach facilitated
260 a thorough exploration of genomic features in the post-hoc interpretability analysis of the BIO-DGI
261 (PPI9) model, providing valuable insights into their roles in the context of disease.

262 *2.8. Identification of CNVs, SVs and LOFs for top 500 significantly altered genes*

263 Our exploration into significantly altered genes underwent expansion to encompass a broader
264 array of genomic profiles, including copy number variants, structural variants, and loss-of-function
265 events. This extended analysis allowed us to delve into the impact of these variants at the gene
266 level, shedding light on their influence on disease progression. For the MMRF dataset, we leveraged
267 the segment data obtained from MMRF CoMMpass to identify copy number variants (CNVs) using
268 the CNVkit [61] tool. To ensure consistency in our CNV identification workflow, we applied CNVkit
269 to detect CNVs in the WES samples of both AIIMS and EGA datasets. Within this framework, we
270 filtered out genes with a copy number (CN) value of 2 across all samples, focusing on genes with CN
271 values that deviated from 2. Next, we utilized all the structural variants identified in WGS samples
272 from the MMRF dataset (i) by the Translational Genomics Research Institute (TGen) through
273 their in-house SV identification workflow and (ii) by the Delly tool [62]. Our analysis centred on
274 significantly altered genes ranked in the top 500, whose genomic regions were affected by structural
275 variants, spanning insertions, inversions, deletions, duplications, and translocations.

276 Furthermore, our investigation extended to encompass genes marked by loss-of-function aber-
277 rations. Loss-of-function refers to a disruption in a gene's normal functioning, hindering the pro-
278 duction of the typical gene product or rendering it ineffective. We assessed every transcript of a
279 gene in a sample to ascertain if it satisfied any of the following conditions: deletion of over half
280 of the coding sequence, deletion of the start codon, deletion of the first exon, deletion of a splice
281 signal, or deletion causing a frameshift, it was considered to exhibit loss-of-function [63]. If at least

282 one condition was satisfied by all the transcripts of a gene in a sample, that sample was labelled to
283 exhibit LoF in that gene [63]. This evaluation was conducted for all the MM and MGUS samples
284 to identify genes featuring loss-of-function.

285 *2.9. Geo2R Validation of top 500 significantly altered genes*

286 We conducted a thorough validation using the Geo2R tool [64] to validate 295 genes against
287 previously published studies focused on multiple myeloma (MM). We utilized a total of eleven
288 micro-array and miRNA data-based MM studies for the validation [65, 66, 67, 68, 69, 70, 71, 72,
289 73, 74, 75, 76, 77]. To ensure rigorous assessment, we exclusively considered genes that displayed
290 significant dysregulation and maintained an adjusted *p*-value ≤ 0.05 in each validation instance.

291 *2.10. Workflow for the design of targeted sequencing gene panel*

292 We devised an innovative workflow to identify potential driver genes for designing the targeted
293 sequencing gene panel for MM. This extensive process integrated various genomic profiles, including
294 single nucleotide variations (SNVs), copy number variations (CNVs), structural variations (SVs),
295 and loss-of-function (LoF) mutations, alongside validated datasets from Geo2R. Initially, we fo-
296 cused on the top 500 genes, extracting relevant data from their SNV, CNV, SV, and LoF profiles.
297 Following this, we implemented a rigorous filtering criteria to identify pivotal genes within each
298 variant profile. This involved generating box plots for the profiling feature across all samples and
299 subsequently applying filtering criteria based on statistical analysis. For instance, for LOF, we
300 made the box plot of number of samples having LoF for every gene of the MM cohort and retained
301 only those genes that were found to have LOF in the number of samples greater than the 3rd
302 quartile of this box plot. Similar criteria were applied on the other variant profiles, detailed in
303 Figure-2. The resulting gene list was consolidated, retaining those meeting criteria in at least one
304 variant profile and had at least one dataset validation in Geo2R analysis. with Geo2R validation.
305 We specifically retained disease-transformative and disease-initiating genes, dropping those genes
306 that were significant only in monoclonal gammopathy of undetermined significance (MGUS) but
307 not in the MM cohort, yielding a final list of 282 genes.

308 Since the analysis originated from SNVs extracted from the WES data and certain key MM
 309 genes like IGH and MYC, known for translocations in MM, were not initially visible, we incorpo-
 310 rated twelve additional well-established MM biomarkers (ATM, CCND1, IGH, IGL, IGK, CKS1B,
 311 HIST1H1E, JAK1, MAF (or c-MAF), MAFB, MYC, and PRDM1), yielding our panel of 295 genes.
 312 These genes, recognized for CNV or SV profiles, were reported in at least two previously published
 313 MM-related panels studied in this work. The detailed workflow for potential driver gene identifi-
 314 cation is presented in Figure-2. Lastly, we assessed the major molecular aberrations for each gene
 315 in this panel of 295 genes and examined coding regions and genomic locations for altered regions
 316 using the UCSC Genome database [78] to understand the genomic spectrum of MM.

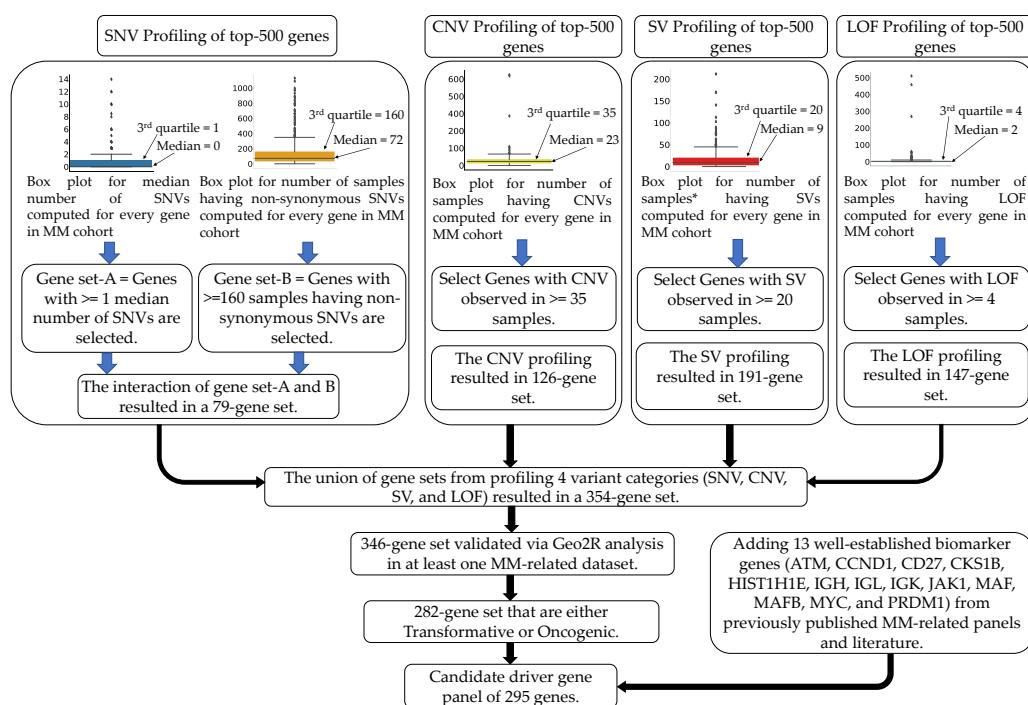


Figure 2: Workflow for the Identification of Biomarker Genes for the Proposed 295-Gene Panel. The workflow integrates variant profiles, including SNVs, CNVs, SVs, and LOF, to discern genes relevant to MM. The combination of these profiles yielded a gene set of 354 candidates, which was further refined to 346 genes after Geo2R validation. From this set, disease-initiating genes (significantly altered in both MM and MGUS) and disease-transformative genes (significantly altered only in MM) were selected for inclusion in the targeted sequencing panel, resulting in a final list of 282 genes. Additionally, thirteen well-established biomarker genes from previously published MM-related panels and literature were incorporated, culminating in the completion of the final 295-gene panel.

317 *2.11. Identification of CNVs, SVs and LOFs of the proposed gene panel and Geo2R Validation*

318 Building upon the genomic profile analysis of CNVs, SVs, and LOFs in the top 500 genes
319 presented in Section 2.8, we further investigated the proposed 295-gene panel by examining its
320 CNV, SV, and LOF landscape. This provided deeper insights into the panel's suitability for MM
321 diagnosis and risk stratification. Additionally, we leveraged Geo2R to validate the panel's relevance
322 against existing MM-related studies, bolstering its potential clinical utility.

323 *2.12. Workflow for the comprehensive survival analysis of the proposed gene panel*

324 Building upon the targeted sequencing gene panel designed in Section-2.10, we conducted a novel
325 survival analysis to investigate the impact of gene variant profiles on patient survival in Multiple
326 Myeloma (MM). Employing two distinct approaches outlined below, we sought to discern the genes
327 that exert a statistically significant impact on the survival outcomes of these patients.

328 In the first approach, univariate survival analysis was conducted for all 295 genes individually,
329 considering each variant profile (SNV, CNV, SV, and LOF) as a singular prognostic factor. For
330 the SNV profile, we utilized the total count of (non-synonymous + other) single nucleotide variants
331 (SNVs) as the prognostic factor in the univariate survival analysis. Analogously, for CNV, SV,
332 and LOF profiles, we constructed categorical vectors (yes/no) indicating the presence or absence
333 of copy number variations, structural variations, and loss-of-function mutations in the multiple
334 myeloma (MM) sample for each gene. Subsequently, we performed univariate survival analysis for
335 each variant profile separately. Genes with a p -value ≤ 0.05 in the univariate survival analysis for
336 individual variant profiles were retained.

337 In a parallel vein, the second approach amalgamated all four variant profiles for each gene, with
338 an aim to elucidate the cumulative impact of gene variant profiles on clinical outcomes. Factor
339 Analysis of Mixed Data (FAMD) [79] was employed for dimensionality reduction in this process.
340 Subsequently, univariate survival analysis were executed on each of the 295 genes in the panel,
341 utilizing the first FAMD component as the prognostic factor. Genes with a p -value ≤ 0.05 in the
342 univariate survival analysis of the first FAMD component were retained. Finally, we considered
343 the union of genes identified as clinically relevant (p -value ≤ 0.05) through the aforementioned two

344 approaches.

345 *2.13. Identification of significantly altered pathways and pathway ranking using the gene panel*

346 Out of 295 genes, the noteworthy 282 genes highlighted by the BIO-DGI (PPI9) model as
347 instrumental in distinguishing MM from MGUS and included in the 295-genes panel via the workflow
348 of Figure-2 were cross-referenced with the significant gene lists derived separately for MM and
349 MGUS cohorts using the dndscv tool. The MM cohort genes were specifically employed in the
350 pathway analysis for MM, while the MGUS cohort genes were utilized for MGUS pathway analysis.
351 Notably, thirteen genes included in the 295-gene panel due to their association with translocations
352 specific to MM disease were integrated with the gene list of the MM cohort for the aforementioned
353 pathway analyses.

354 To further elucidate the functional implications, we employed the ‘Enrichr gene set enrichment
355 analysis web server’ [80, 81, 82], facilitating the identification of KEGG and Reactome pathways
356 associated with our proposed gene panel. Subsequent ranking of significantly altered pathways in
357 the MM and MGUS cohorts, based on their adjusted *p*-values, provided a comprehensive insight
358 into the primary pathways undergoing substantial alterations due to genomic aberrations in the
359 significantly altered genes.

360 *2.14. Identification of Haploinsufficient genes of the gene panel*

361 To assess the likelihood of genes exhibiting haploinsufficiency, we draw upon two previously
362 published haploinsufficiency prediction scores: the genome-wide haploinsufficiency score (GHIS)
363 [83] and the DECIPHER score [63]. The DECIPHER score amalgamates patient genomic data,
364 evolutionary profiles, and functional and network properties to predict the likelihood of haploin-
365 sufficiency. Meanwhile, the GHIS score draws from diverse large-scale datasets, encompassing gene
366 co-expression and genetic variation in over 6000 human exomes. These comprehensive methods
367 enhance identifying haplo-insufficient genes, revealing their crucial role in diseases. This deepens
368 our understanding of genes that lack proper function when only one copy is present.

369 **3. Results**

370 *3.1. Cohort Description*

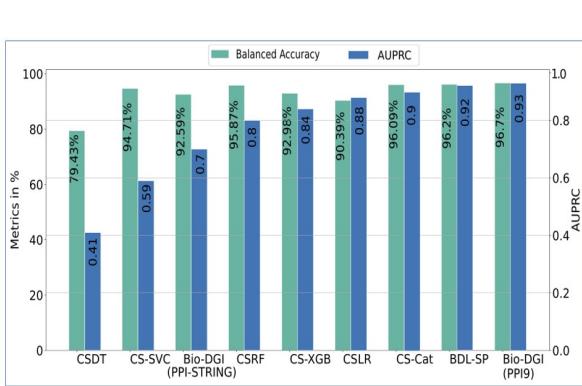
371 In this comprehensive study, we analyzed two distinct cohorts related to MM and MGUS,
372 encompassing a total of 1154 MM samples and 61 MGUS samples sourced from three globally
373 recognized datasets: AIIMS, EGA, and MMRF. Specifically, within the MM cohort, we examined
374 1072 samples from the MMRF dataset and 82 samples from the AIIMS dataset. Additionally,
375 in the MGUS cohort, we examined 28 samples from the AIIMS repository and 33 from the EGA
376 repository. Augmenting our analysis, we incorporated crucial clinical data, including overall survival
377 (OS) time and event data for MM samples retrieved from the MMRF and AIIMS datasets. This
378 enabled a thorough exploration of the clinical relevance of the proposed targeted sequencing panel,
379 underlining the significance of our findings.

380 *3.2. Identification of Significantly altered genes*

381 We employed the DNDSCV tool, as illustrated in the pre-processing block of Figure-1, to discern
382 genes exhibiting significant alterations within the MM and MGUS cohorts. A total of 598 and 351
383 significantly altered genes were identified in the MM and MGUS cohorts, respectively. Notably,
384 151 genes were found to be shared between the MM and MGUS cohorts. Subsequently, we pursued
385 the inference of pivotal genes and gene-gene interactions vital for discriminating MM from MGUS,
386 leveraging our innovative graph-based BIO-DGI (PPI9) model.

387 *3.3. Benchmarking of proposed BIO-DGI (PPI9) model*

388 Employing our AI-driven BIO-DGI workflow (depicted in Figure-1), we trained the BIO-DGI
389 (PPI9) model using 5-fold cross-validation and compared its performance with six standard cost-
390 sensitive machine learning and two deep learning models. Remarkably, the proposed BIO-DGI
391 (PPI9) model showcased superior performance in terms of balanced accuracy and AUPRC (area
392 under the precision-recall curve). Specifically, the BIO-DGI (PPI9) model achieved the highest
393 balanced accuracy of 96.7%. Following closely, the BDL-SP model attained a balanced accuracy
394 of 96.26%, and the cost-sensitive Catboost (CS-Cat) model achieved the third-best performance



Model	Balanced Accuracy	AUPRC	Confusion Matrix
Bio-DGI (PPI9)	96.7	0.93	{TP: 1099, FP: 1, FN: 55, TN: 60}
BDL-SP	96.2	0.92	{TP: 1087, FP: 1, TN: 60, FN: 67}
CS-Cat	96.09	0.9	{TP: 1120, FP: 3, TN: 58, FN: 34}
CSLR	90.39	0.88	{TP: 1141, FP: 11, TN: 50, FN: 13}
CS-XGB	92.98	0.84	{TP: 1122, FP: 7, TN: 54, FN: 32}
CS-RF	95.87	0.8	{TP: 1078, FP: 1, TN: 60, FN: 76}
Bio-DGI (PPI-STRING)	92.59	0.7	{TP: 1097, FP: 6, TN: 55, FN: 57}
CS-SVC	94.71	0.59	{TP: 1069, FP: 2, TN: 59, FN: 85}
CS-DT	79.43	0.41	{TP: 1132, FP: 24, TN: 37, FN: 22}

Figure 3: Quantitative benchmarking of proposed BIO-DGI(PPI9) model. (A) Comparison of balanced accuracy and AUPRC score of BIO-DGI(PPI9) model with other baseline ML and DL models, and (B) Confusion matrix of top-performing models including BIO-DGI (PPI-STRING) model. Notations: Cost-Sensitive Decision Tree (CS-DT); Cost-Sensitive Support Vector Classifier (CS-SVC); Cost-Sensitive Random Forest (CS-RF); Cost-Sensitive XGBoost (CS-XGB); Cost-Sensitive Logistic Regression (CS-LR); Cost-Sensitive CatBoost (CS-Cat).

395 with a balanced accuracy of 96.09%. The BIO-DGI (PPI9) model also outperformed other models
 396 in AUPRC, securing the highest AUPRC score of 0.93, while the AUPRC score for BDL-SP and
 397 CS-Cat models stood at 0.92 and 0.9, respectively. Notably, the BIO-DGI (PPI9) model correctly
 398 identified 1099 out of 1154 MM samples and 60 out of 61 MGUS samples, showcasing its superior
 399 performance.

400 These results affirm that, quantitatively, the BIO-DGI (PPI9) model performed superior with
 401 the BDL-SP model being the second best. For a comprehensive understanding of quantitative
 402 performance, refer to Figure-3 (A) and (B) for balanced accuracy and AUPRC scores, confusion
 403 matrices, and AUPRC curves, respectively.

404 Given the marginal difference in the balanced accuracy and AUPRC performance metrics among
 405 the top three models (BIO-DGI (PPI9), BDL-SP, and CS-Cat), we conducted post-hoc interpretabil-
 406 ity benchmarking by applying ShAP algorithm to identify the top-ranked genes for each of the top
 407 three performing models. Subsequently, we analyzed these genes to identify previously reported
 408 oncogenes (OG), tumour-suppressor genes (TSG), both oncogenes and driver genes (ODG), and

Table 1: Number of previously reported genes present in 798 significantly altered genes and qualitative benchmarking of top-performing models

(A) Types of four different gene categories (OG, TSG, ODG, and AG) and their counts in 798 significantly altered genes

Gene type based on functionality	Number of genes
Oncogenes (OGs)	31
Tumor-suppressor genes (TSGs)	43
Both oncogene and driver gene (ODGs)	10
Actionable genes (AGs)	19

(B) Counts of previously reported four categories of genes as found in the post-hoc analysis based on top-250 and top 500 genes of the top-3 models (BIO-DGI (PPI9), BDL-SP, CS-Cat, and BIO-DGI (PPI-STRING))

Top Genes	BIO-DGI (PPI9) (Top-performing model)				BDL-SP (Second best model)				CS-Cat (Third best model)				BIO-DGI (PPI-STRING) (Baseline version of BIO-DGI)			
	OG	TSG	ODG	AG	OG	TSG	ODG	AG	OG	TSG	ODG	AG	OG	TSG	ODG	AG
Top-250	23	26	8	14	20	21	7	11	0	0	0	0	18	24	7	13
top 500	28	41	9	19	27	37	8	17	0	0	0	0	28	41	9	19

The number of previously reported genes (OG/TSG/ODG/AG) obtained in each category (top-250/top 500) using the best performing model are highlighted in bold.

409 actionable genes (AG). Out of the total 798 genes, we identified 31 OGs (including *ABL2*, *BIRC6*,
 410 *FUBP1*, *IRS1*), 43 TSGs (including *APC*, *ARID1B*, *CYLD*, *PABPC1*, *ZFHX3*), 10 ODGs (in-
 411 cluding *BRAF*, *FGFR3*, *TP53*, *TRRAP*), and 19 AGs (including *ARID2*, *BRD4*, *MITF*, *NF1*,
 412 *TYRO3*) (Table-1A).

413 Our analysis revealed that the proposed BIO-DGI (PPI9) model exhibited the highest count of
 414 identified oncogenes (OG), tumor-suppressor genes (TSG), both oncogene and driver genes (ODG),
 415 and actionable genes (AG) in both the top-250 and top 500 gene lists (Table-1B). Specifically,
 416 the BIO-DGI (PPI9) model detected 23 and 28 oncogenes in the top 250 and top 500 gene list,
 417 respectively. Of the 43 known TSGs, the BIO-DGI (PPI9) model identified 26 genes in the top 250
 418 and 41 in the top 500 gene lists. Of the 10 known ODGs, the BIO-DGI (PPI9) model identified 8
 419 genes in top 250 and 9 genes in top 500 gene lists. Lastly, of the 19 known AGs, the BIO-DGI
 420 (PPI9) model identified 14 genes in top 250 and 19 genes in top 500 gene list.

421 We have considered only those genes in the top-250 or top 500 gene list that have a non-zero

422 ShAP score in the post-hoc explainability analysis. The total counts of previously reported genes as
423 found in the top-250 and top 500 genes of the top-four models (BIO-DGI(PPI9), BDL-SP, CS-Cat,
424 BIO-DGI (PPI-STRING)) is shown in Table-1B. Furthermore, the lists of top 500 genes obtained
425 using Bio-DGI (PPI9) and the previously reported genes ranked within the top 250 and top 500 by
426 the top-performing models are outlined in Table-S2, Supplementary File-1 and Table-2.

427 Given the BIO-DGI (PPI9) model's superior identification of previously reported OGs, TSGs,
428 ODGs, and AGs, it also stands out as the best-performing model in the post-hoc analysis and
429 was subsequently used to infer the top significantly altered genes, gene-gene interactions, genomic
430 features, and altered signalling pathways critical for distinguishing MM from MGUS. This analysis
431 underscores the importance of model interpretability within the application domain, particularly,
432 when similar quantitative results are obtained with different machine learning models.

433 *3.4. Interpretability of BIO-DGI (PPI9) model using ShAP algorithm*

434 We utilized the ShAP algorithm for post-hoc model explainability and rank genomic attributes
435 based on their influence on the model prediction. Each genomic attribute received a ShAP score,
436 representing its contribution to each class (MM/MGUS). Subsequently, the attributes were ranked
437 at the cohort level (MM versus MGUS) accordingly. This ShAP analysis provided post-hoc ex-
438 plainability of the trained model, following a methodology akin to that outlined in [27], enabling
439 the ranking of genes and genomic features at both cohort and sample levels.

440 By evaluating the ShAP scores assigned to each gene, we identified *MUC6*, *LILRA1*, and
441 *LILRB1* as the top three genes in MM and MGUS samples among the 798 significantly altered
442 genes. Furthermore, several previously reported oncogenes (e.g., *MUC16*, *USP6*, *BIRC6*, *VAV1*),
443 tumor-suppressor genes (e.g., *EP400*, *HLA-B/C*, *SDHA*, *MYH11*), both oncogenes and driver genes
444 (e.g., *PABPC1*, *KRAS*, *TRRAP*, *TP53*, *FGFR3*, *BRAF*), and actionable genes (e.g., *NOTCH1*,
445 *FANCD2*, *TYRO3*, *ARID1B*) were highlighted as top-ranked genes.

446 Similarly, we ranked genomic features based on their impact on the model's prediction using
447 their ShAP scores. In our model training for BIO-DGI (PPI9), a set of 26 genomic features was
448 employed. Notably, the PhyloP score of non-synonymous SNVs, allele depth of synonymous SNVs,

Table 2: List of 4 categories of previously reported genes as found in the post-hoc analysis based on top-250 and top 500 genes of the top-3 models (BIO-DGI (PPI9), BDI-SP, and BIO-DGI (PPI-STRING))

449 and the total number of other SNVs (that included non-frameinsertion/deletion/substitution,
450 intronic, intergenic, ncRNA-intronic, upstream, downstream, unknown, and ncRNA-splicing SNVs)
451 emerged as the top three genomic features. Figure-4 presents the beeswarm plot illustrating the
MGUS|MM



Figure 4: Genomic Feature Ranking using the ShAP Algorithm in MM and MGUS based on post-hoc explainability by the BIO-DGI model. Genomic features are ranked according to their ShAP scores. A positive ShAP score indicates the feature's contribution to MM, while a negative score represents its contribution to MGUS. Each dot in the scatter plot represents a sample colour-coded to reflect genomic feature values—dark blue for low and red for high values.

453 3.5. Analysis of CNVs, SVs and LOF of top-500 genes in MM

454 In addition to analyzing SNV profile, we comprehensively investigated CNVs, SVs, and LOF of
455 the top 500 genes of the MM cohort. CNV identification was performed using CNVkit on AIIMS
456 MM samples and on exome segment data from MMRF CoMMpass for MMRF samples. Processed
457 SV data from MMRF CoMMpass was utilized to identify key SVs in MM and 295-genes panel

458 designing. For identifying genes with LOF within a sample, we employed established criteria to
459 evaluate disruptions in gene transcripts due to deletion of essential coding segments, exons, splice
460 signals, or frameshift-inducing deletions [63]. We studied both CNVs and SNVs to identify genes
461 with LOF within each sample.

462 CNVs, SVs, and LOF analysis of the top 500 genes revealed crucial molecular aberrations in
463 MM. Chromosome-wise distribution analysis indicated that chr19 (19%), chr1 (17%), chr6 (8.6%),
464 and chr14 (7.1%) were notably affected by CNVs (Figure-5(A)). Similarly, chr1 (12.6%), chr6
465 (9.9%), chr12 (5.3%), and chr14 (5%) showed prominent SV involvement (Figure-5(B)), while chr19
466 (20%), chr1 (19.9%), chrX (13%), and chr14 (11.8%) were most affected by LOF (Figure-5(C)). The
467 majority of CNVs were gains (58.3%) and deletions (17%) (Figure-5(D)), while inversions (65%) and
468 translocations (13.1%) dominated the SV landscape (Figure-5(E)). Notable chromosomes impacted
469 by inversion SV included chr1, chr3, chr2, and chr7 (Figure-5(F)), and translocations mainly affected
470 chr7, chr21, chr1, and chr14 (Figure-5(G)). The distribution of CNV and SV types within each
471 chromosome highlighted their relative abundance (Figure-5(H) and Figure-5(I)).

472 *3.6. Design of 295-gene targeted sequencing panel*

473 To design an effective targeted sequencing panel, we refined the initially identified top-ranked
474 genes based on their significant alterations and the collective impact of their variant profiles in
475 MM. Firstly, we considered four critical variant profiles to identify the candidate driver gene panel:
476 1. SNV profile, 2. CNV profile, 3. SV profile, and 4. LOF profile. We also integrated the
477 Geo2R validation profile to include MM-relevant genes in the targeted sequencing panel. Finally,
478 we excluded genes that were neither disease-transformative nor disease-initiating. For the SNV
479 profiling of the top 500 significantly altered genes, we filtered based on the median SNV count and
480 the number of samples with non-synonymous SNVs, resulting in 79 genes. The features extracted for
481 SNV profile analysis are detailed in Table-S4, Supplementary File-3. The variant profiling for CNV,
482 SV, and LOF involved filtering genes based on the number of samples exhibiting that particular
483 variant type, yielding 126, 191, and 147 genes, respectively. The features extracted for CNV, SV,
484 and LOF profile analysis can be found in Table-S5, Table-S6, and Table-S7, Supplementary File-3.

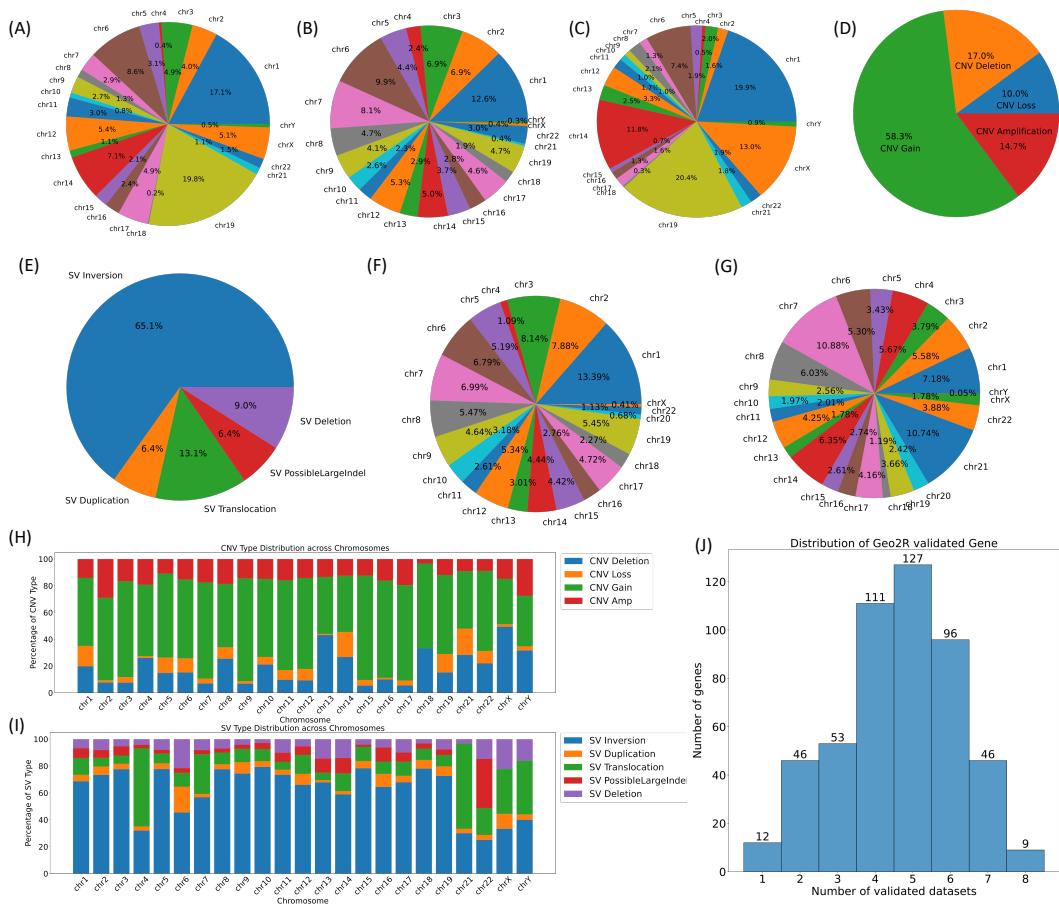


Figure 5: Genomic Aberrations Overview (CNVs, SVs, and LOF) in MM Samples from AIIMS and MMRF Repositories. The figure in panels (A)-(C) displays the chromosome-wise distribution of CNVs, SVs, and LOF. Panel (D) presents the distribution of CNV types identified in MM samples from both AIIMS and MMRF datasets. Similarly, panel (E) shows the distribution of SV types identified in MM samples from the MMRF dataset. Notably, SV analysis was conducted exclusively for MMRF samples due to the absence of WGS data in the AIIMS repository. Continuing SV analysis, panels (F) and (G) exhibit the chromosome-wise distribution of inversions and translocations found in MM samples. Panels (H) and (I) provide the distribution of CNV and SV types for each chromosome individually. (J) Distribution of the top 500 genes validated through MM-related studies using the Geo2R tool. The x-axis represents the number of MM-related studies validating the gene, while the y-axis indicates the count of genes.

485 By combining genes from SNV, CNV, SV, and LOF variant profiles, we arrived at a comprehensive
 486 set of 354 genes. To ensure relevance, we retained genes validated in at least one MM-related
 487 study using Geo2R validation. Of the 354 genes, 346 were validated through Geo2R validation
 488 analysis. In the final selection, we focused on 212 disease-transformative and 70 disease-initiating
 489 genes, resulting in a list of 282 genes. Further, we also included 13 well-established MM biomarker

490 genes that included nine disease-initiating genes (CCND1, CKS1B, IGH, IGK, IGL, MAF, MAFB,
491 MYC, and PRDM1) and four disease-transformative genes (ATM, CD27, HIST1H1E, and JAK1),
492 leading to the design of our proposed 295-gene panel (Table-3, Table-S16, Supplementary File-8).
493 The workflow for designing the 295-gene panel is illustrated in Figure-2.

494 In this panel, four genes, namely, HLA-A, HLA-B, HLA-DRB5, and RYR3, were heavily mutated
495 in all four variant profiles. Additionally, 122 and 32 genes were substantially mutated in at least
496 two and three variant profiles, respectively. Within the MM cohort, notable previously reported
497 significantly altered genes were present including *BRAF*, *IGLL5*, *IRF*, *IGH*, *MYC*, *JAK*, *MAF*,
498 *KRAS*, *TP53*, *TRAF2/3*, among others. Similarly, the MGUS cohort exhibited previously reported
499 genes like *HLA-B*, *LILRB1*, *PABPC1*, *PRSS3*, among others. Several previously reported genes
500 were found in both MM and MGUS cohorts, including *HLA-B*, *PRSS3*, *KMT2C*, among others,
501 illustrating their potential role as shared genomic features in the progression from MGUS to MM.

502 We determined each gene's most prevalent molecular aberration, such as CNV gain, CNV loss,
503 SV translocation, LOF, etc. We observed that CNV gain was the most frequent molecular aberration
504 found in 188 out of the 295 genes, while LOF was the least common, identified in 12 out of the
505 295 genes. To refine the targeted sequencing regions further, we assessed the most affected coding
506 regions using the UCSC Genome database. The targeted sequencing panel of 295 genes covered
507 9,417 coding regions in the human genome, spanning a genomic region with a total length of 2.630
508 Mb in the human genome (Table-S17, Supplementary File-8).

509 *3.7. Identification of Significantly altered pathways and ranking of pathway of 295-gene panel*

510 Of the 295-gene panel, only 70 genes were found to be significantly altered in the MGUS cohort,
511 while all 295 genes were part of the MM cohort. We utilized the Enrichr database to identify
512 significantly altered KEGG and Reactome signalling pathways associated with the 295 genes of the
513 MM cohort and those associated with the 70 genes of the MGUS cohort. A total of 39 KEGG
514 and 25 Reactome pathways exhibited significant alterations for the MGUS cohort (see Table-S8 in
515 Supplementary File-4), while 123 KEGG and 50 Reactome pathways were observed to be signifi-
516 cantly affected for the MM cohort (refer to Table-S9 in Supplementary File-4). We categorized the

Table 3: List of disease-transformative and disease-initiating genes in the proposed 295 gene panel. ‘Red’ color denotes oncogenes (OG), ‘black’ denotes genes which are both oncogenes and driver genes (ODG), ‘green’ denotes tumour suppressor genes (TSG), ‘blue’ denotes actionable genes (AG), and ‘magenta’ color is for the rest of the genes which are not previously reported as OG/TSG/ODG/AG in multiple myeloma.

Gene type	Genes
disease-initiating genes	ABCA3, ACTR3B, ADAM21, AHNAK, AHNAK2, AMER1, ANKRD36C, ASH1L, CACNA1B, CCND1, CKS1B, CSMD2, DHX35, DNAH6, DOCK8, FAT2, FCGBP, FLG2, FMN2, FRG2B, HELZ2, HLA-B, HLA-DQA2, HLA-DRB5, HUWE1, IGFN1, IGH, IGK, IGL, ITPR1, ITPR2, KCNJ12, KMT2C, KPRP, KRT38, KRT6B, KRTAP5-4, KRTAP9-9, LAMA3, LILRB1, MAF, MAFB, MAP3K10, METTL2B, MUC16, MUC4, MUC6, MYC, MYH7, NBPFI0, NBPFI9, NEB, OBSCN, OR8U1, PABPC1, PABPC3, PAK2, PLIN4, PLXNA3, PRAMEF11, PRDM1, PRSS3, RBMXL1, RLIM, RYR1, RYR3, SACS, SIRPA, SKA3, SLC25A5, SYNGAP1, TAS2R30, TBP, TTN, TUBA3C, TUBB8, U2AF2, WDFY4, ZNF676
Disease Transformative genes	ABCA1, ABCA7, ABL2, ACACB, ACVR1B, ADAMTS18, ALG13, ANKIB1, APBA2, ARHGAP4, ARID1B, ARID2, ATM, ATP12A, ATP2B2, ATP2B3, BRAF, BRD4, BSN, C3, C4B, CACNA1A, CACNA1F, CACNA2D2, CASKIN2, CCNT1, CD27, CDC42BPG, CELSR1, CFH, CHD5, CHN2, CHRHM3, CLIP1, COL14A1, COL5A1, CSPG4, CUL9, CYLD, CYP2A6, DENND1A, DIS3, DNAH1, DNAH17, DNAH5, DNAH7, DOCK2, DOCK3, DUSP2, DYSF, EGR1, EIF4EBP1, ELFN2, EML2, EPC1, EPPK1, EZR, F8, FAM104B, FAM178B, FAM186A, FAM46C, FASN, FBN3, FGFR3, FHOD3, FLNA, FRG2C, FRY, FUBP1, GEMIN2, GOLGA6L2, GON4L, GSN, HADHB, HERC1, HERC2, HIPK3, HIST1H1E, HLA-A, HLA-C, HLA-DQA1, HLA-DQB2, HRNR, IGLL5, INF2, INPP5D, IQSEC3, IRF1, IRF2BPL, IRS1, ITGA2, JAK1, KCNT1, KHDRBS1, KIF13A, KIF26B, KIR2DL1, KIR3DL2, KLC3, KMT2B, KMT2D, KRAS, KRT8, KRTAP5-10, L1CAM, LAMA2, LCORL, LILRA1, LILRA2, LILRA4, LILRB2, LTB, LYST, MAGEC1, MANEAL, MAP3K9, MAP4, MAX, MBNL1, MECOM, MED12L, MEI1, MITF, MLLT1, MMP16, MUC12, MUC20, MYH14, MYH6, MYO15A, MYO18B, MYO1C, MYO5B, NELFE, NF1, NFKBIA, NFX1, NIPBL, NNT, NOS1, NRAS, NUDT10, NUMBL, OTOG, PARP4, PDE1C, PGM5, PGR, PHF14, PKD1, PLCH2, PLEC, PLXNB3, PPARGC1B, PPFIA3, PRRC2A, PRSS1, PSORS1C1, PTK2B, PTK7, RAB12, RAD54B, RARG, RB1, RBM20, RBM23, RPL10, RPTOR, RTELI, RXRB, RYR2, SAMHD1, SCAF1, SDHA, SF3A3, SLAMF7, SLC12A3, SLC7A1, SP140, SPTA1, SPTB, STAB1, SVIL, TAF1, TAL1, TGM7, TNRC6A, TP53, TPSD1, TPTE2, TRAF2, TRAF3, TRPM2, TRPM7, TUBGCP6, UNC13A, UNC79, USP9X, VAV1, VCL, VPS13A, VPS13B, VWF, XCR1, YWHAZ, ZC3H4, ZFHX3, ZNF208, ZNF469, ZNF587, ZNF717, ZNF763, ZNF865, ZNFX1, ZZEF1

517 significantly altered pathways into four distinct groups according to their significance level changes
 518 during the MGUS to MM transition:

- 519 Category-1: Pathways increasing in significance in MGUS to MM progression.
 520 Category-2: Pathways decreasing in significance in MGUS to MM transition.
 521 Category-3: Pathways significantly altered in MM but not in MGUS.
 522 Category-4: Pathways significantly altered in MGUS but not in MM.

523 The complete list of significantly altered pathways for these categories is provided in Tables-S10
524 and Table-S11 of Supplementary File-4. A total of 32 KEGG and 13 Reactome pathways became
525 more significant as the disease progressed from MGUS to MM, while 5 KEGG pathways and 7
526 Reactome pathways displayed reduced significance with disease progression from MGUS to MM.
527 A total of 86 KEGG and 30 Reactome pathways were significantly altered only in MM and not
528 in MGUS. Notably, 37 out of 86 KEGG pathways and 5 out of 30 Reactome pathways showed no
529 overlapping genes, with 70 significantly altered genes in MGUS. Lastly, 2 KEGG pathways and 5
530 Reactome pathways were observed as significantly altered only in MGUS and not in MM.

531 To determine the top-ranked pathways, we ranked significantly altered pathways in MM based
532 on their adjusted *p*-values (refer to Table-S12, Supplementary File-4). This analysis revealed a
533 selection of MM-related signalling pathways, notably encompassing the antigen processing and
534 presentation, PI3K-AKT signalling pathways, and B-cell receptor prominently featured among the
535 top-ranking pathways.

536 *3.8. Analysis of identified gene communities with reference to the 295-gene panel*

537 We employed a five-fold cross-validation training strategy to obtain five distinct learned ad-
538 jacency matrices for five classifiers, each with a dimension of 798x798. We applied the Leiden
539 algorithm to the respective learned adjacency matrix for each classifier to identify gene commu-
540 nities. Consequently, we derived 5, 5, 6, 5, and 6 gene communities using the learned adjacency
541 matrices from the first, second, third, fourth, and fifth classifiers, respectively. We ranked the
542 communities within each classifier based on the number of previously reported genes present and
543 selected the top three gene communities for each. Subsequently, we merged these top three gene
544 communities for each classifier, resulting in five new distinct learned adjacency matrices with di-
545 mensions of 500x500, 500x500, 539x539, 500x500, and 422x422. In the following step, we merged
546 these five distinct newly learned adjacency matrices by computing the mean of gene-gene interac-
547 tions across the five classifiers. In cases where a specific gene-gene interaction was absent in any
548 fold, we assigned a zero weight for the corresponding interaction in that fold. This process yielded
549 a final adjacency matrix with dimensions of 690 x 690.

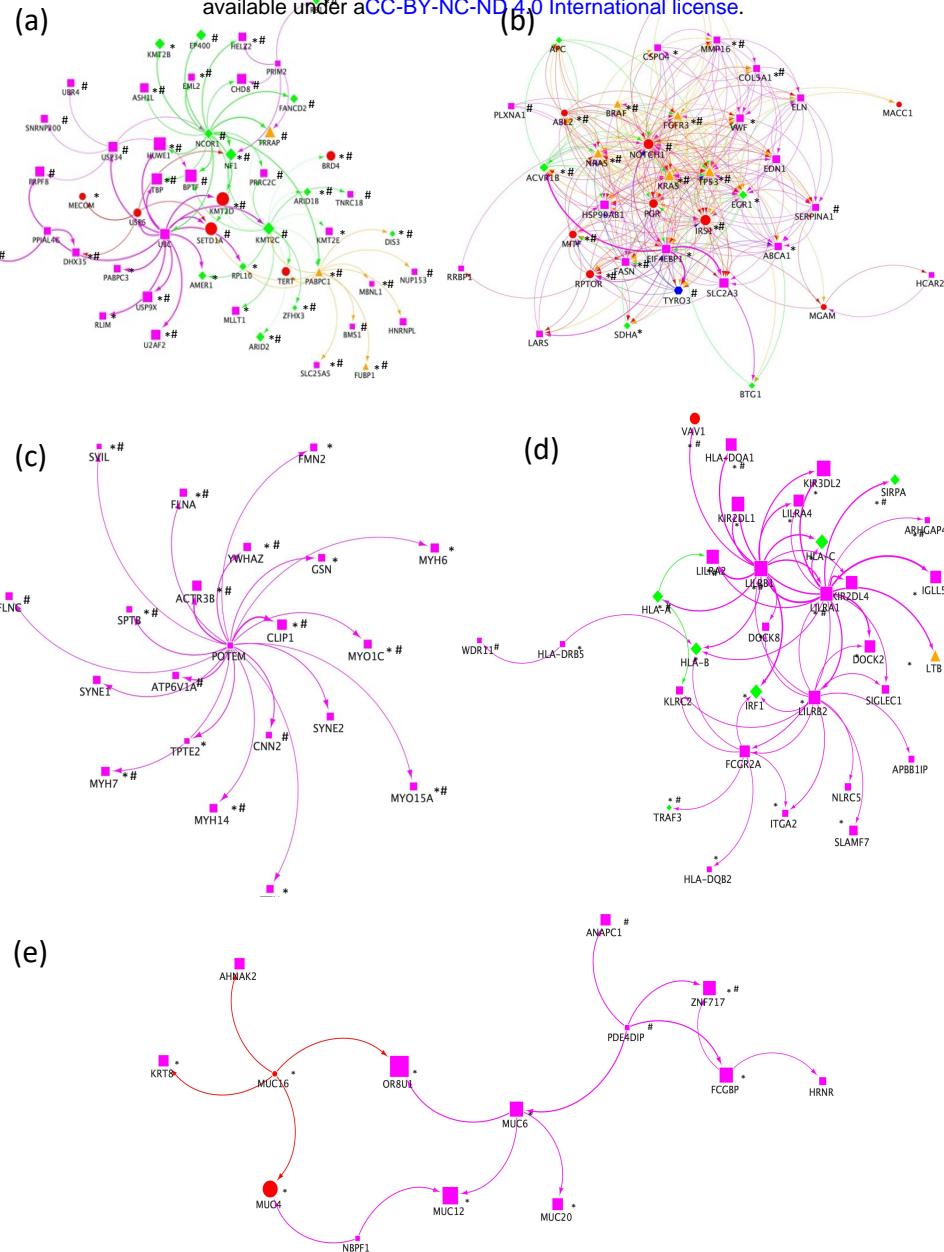


Figure 6: Directed gene community visualization using the learned adjacency matrix obtained from five trained BIO-DGI (PPI9) classifiers. In this figure, (a), (b), (c), (d), and (e) represent the top genes in the first, second, third, fourth and fifth gene communities, respectively. These figures showcase the previously reported genes (OG, TSG, ODG, AG) regardless of their rank, alongside other non-reported genes (in magenta colour) within the top 250 ranks, respectively. Genes marked with “**” are included in the 295-genes panel. Similarly, genes marked with “#” are also highly likely haploinsufficient, with a GHIS score > 0.52. **red circle** shape node represents oncogene, **green diamond** shape node represents TSG, **orange triangle** shape node represents genes that are both oncogene and driver gene, **blue hexagon** shape node actionable gene and **magenta square** shape node represents genes that are not previously reported OG/TSG/ODG/AG.

Finally, we identified five gene communities from the final learned adjacency matrix using the Leiden algorithm, yielding communities having 202, 125, 122, 104, and 21 genes. The pseudo

552 codes for community detection are provided in Supplementary File-5. The first gene community,
553 comprising of 202 genes, contained 11 OGs, 21 TSGs, 3 ODGs, and 8 AGs. Similarly, the second
554 gene community, with 125 genes, contained 14 OGs, 8 TSGs, 6 ODGs, and 10 AGs. The third gene
555 community, comprising 122 genes, did not include any OGs, TSGs, ODGs, or AGs. The fourth
556 gene community, with 104 genes, contained 4 OGs, 11 TSGs, 1 ODG, and 1 AG. Lastly, the fifth
557 gene community, comprising 21 genes, contained 2 OGs and no TSGs, ODGs, or AGs. The list of
558 genes present in all five gene communities and previously reported genes within each are provided
559 in Table-S13 and Table-S14 of Supplementary File-6. Visualization of all five gene communities,
560 including the top 250 genes and previously reported genes (regardless of their rank), is presented
561 in Figure-6(A)-(E).

562 *3.9. Analysis of CNVs, SVs and LOF associated with 295-gene panel*

563 In addition to analyzing SNV profiles, we comprehensively investigated CNVs, SVs, and LOF in
564 the MM cohort. CNV identification was performed using CNVkit on AIIMS MM samples and on
565 exome segment data from MMRF CoMMpass for MMRF samples. Processed SV data from MMRF
566 CoMMpass was utilized to identify key SVs in MM and 295-genes panel designing. For identifying
567 genes with LOF within a sample, we employed established criteria to evaluate disruptions in gene
568 transcripts due to deletion of essential coding segments, exons, splice signals, or frameshift-inducing
569 deletions [63]. We studied both CNVs and SNVs to identify genes with LOF within each sample.
570 CNVs, SVs, and LOF analysis in the 295 genes revealed crucial molecular aberrations in MM.
571 Chromosome-wise distribution analysis indicated that chr19 (18.03%), chrX (14.2%), and chr1
572 (9.02%) were notably affected by CNVs (7(A)). Similarly, chr1 (11.4%), chr6 (9.3%), and chr (6.8%)
573 showed prominent SV involvement (7(B)), while chrX (20.43%), chr16 (13.27%), and chr1 (12.52%)
574 were most affected by LOF (7(C)). The majority of CNVs were gains (50.8%) and deletions (16.9%)
575 (7(D)), while inversions (59.2%) and translocations (18.3%) dominated the SV landscape (7(E)).
576 Notable chromosomes impacted by inversion SV included chr1, chr3, chr7, and chr8 (7(F)), and
577 translocations mainly affected chr7, chr11, chr8, and chr1 (7(G)). The distribution of CNV and SV
578 types within each chromosome highlighted their relative abundance (7(H) and 7(I)).

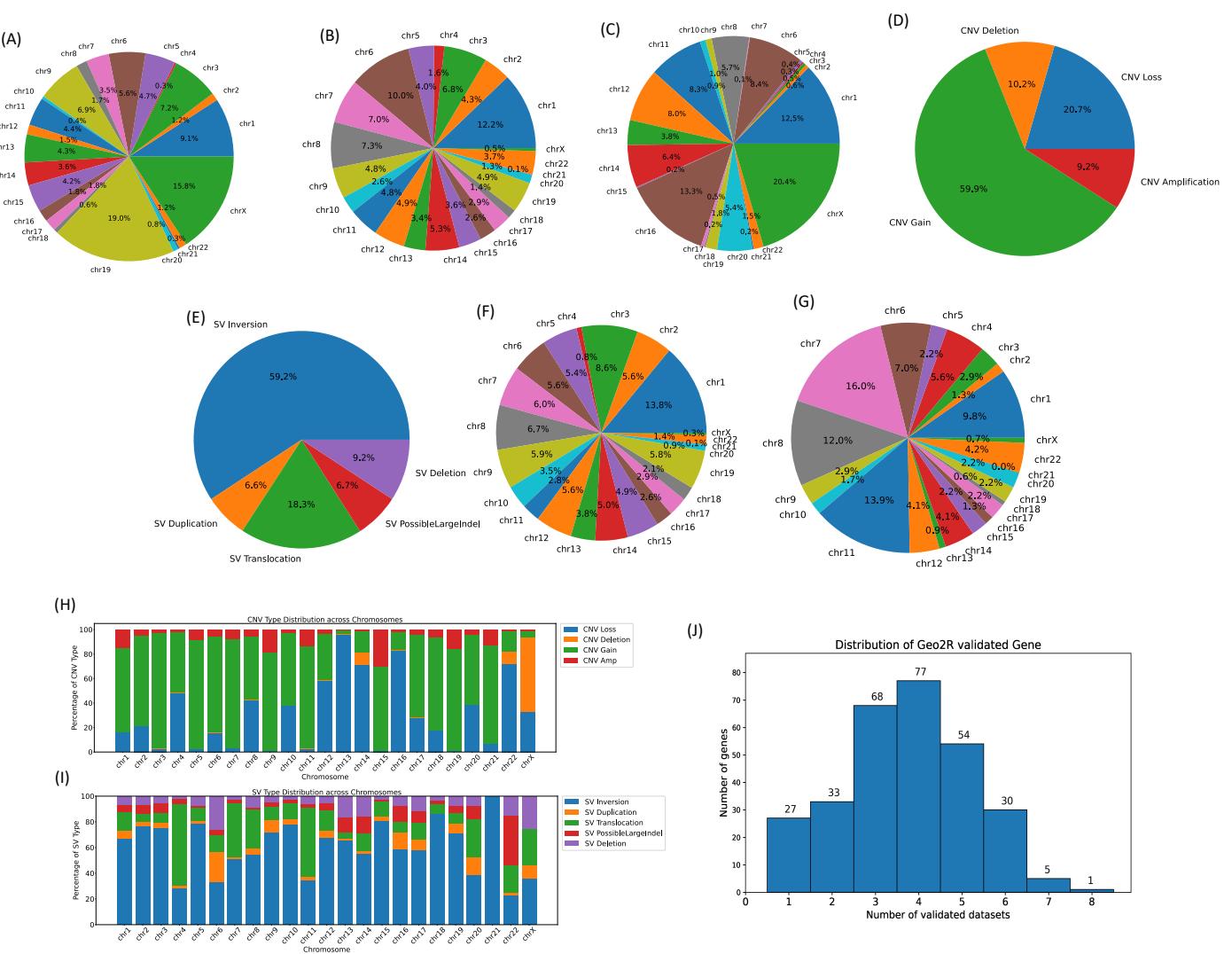


Figure 7: Overview of genomic Aberrations associated with 295 genes (CNVs, SVs, and LOF) in MM Samples from AIIMS and MMRF Repositories. The chromosome-wise distribution of genomic features is shown in panels (A) for CNVs, (B) for SVs, and (C) for LOF. Panel (D) presents the distribution of CNV types identified in MM samples from both AIIMS and MMRF datasets. Similarly, panel (E) shows the distribution of SV types identified in MM samples from the MMRF dataset. Notably, SV analysis was conducted exclusively for MMRF samples due to the absence of WGS data in the AIIMS repository. Continuing SV analysis, panels (F) and (G) exhibit the chromosome-wise distribution of inversions and translocations found in MM samples. Panels (H) and (I) provide the distribution of CNV and SV types for each chromosome individually. (J) Distribution of 295-gene panel validated through MM-related studies using the Geo2R tool. Here, the x-axis represents the number of MM-related studies validating the gene, while the y-axis indicates the count of genes.

579 *3.10. Geo2R validation of the proposed 295-gene panel*

580 We ascertained the relevance of the proposed 295-gene panel with reference to MM via results
581 of the existing MM-related studies through Geo2R validation. Geo2R tool is one of the most
582 widely used tools for identifying significantly dysregulated genes using gene expression or microarray
583 data from previously published studies. We considered 11 MM-related studies for this validation,
584 identifying significantly expressed genes with an adjusted *p*-value of ≤ 0.05 and compared them
585 with our top-ranked genes. Remarkably, out of 295 genes, 268 genes were validated in at least two
586 MM-related studies (Table-S17, Supplementary File-7). Moreover, 68 (23.05%), 77 (26.10%), and
587 54 (18.30%) genes were found to be significantly dysregulated in MM across datasets related to
588 three, four and five MM-related studies, respectively, as depicted in Figure-7(J).

589 *3.11. Clinical relevance of targeted sequencing 295-gene panel*

590 We performed a two-fold univariate survival analysis on a targeted sequencing panel comprising
591 295 genes to comprehend how gene variant profiles affect clinical outcomes in MM patients (Figure-
592 8). We utilised two distinct approaches to gauge the effect of gene variant profiles on MM sample
593 clinical outcomes. In the first approach, we individually assessed the impact of each variant profile
594 (SNV, CNV, SV, and LOF) on clinical outcomes using univariate survival analysis. Notably, 168
595 of the 295 genes significantly influenced clinical outcomes based on at least one variant profile. Of
596 these, 30, 88, 27, and 79 genes significantly impacted clinical outcomes based on SNV, CNV, SV,
597 and LOF variant profiles as prognostic factors, respectively (Table-S16, Supplementary File-8). In
598 the second approach, we amalgamated all four variant profiles into a single feature vector using
599 the FAMD method, leveraging the FAMD first component as a prognostic factor for univariate
600 survival analysis. Subsequently, we found that 188 of the 295 genes significantly influenced clinical
601 outcomes based on the FAMD first component. Combining the clinically relevant genes from the
602 two approaches mentioned above, we discovered that 226 of the 295 genes were clinically relevant
603 for MM. We analyzed the remaining 69 genes that did not show significance in any of the mentioned
604 approaches; we examined them and retained them in the proposed gene panel as these genes were
605 heavily mutated in at least one variant profile (Table-S16, Supplementary File-8). The forest plots

606 and survival curves of genes found significant in at least three variant profiles (e.g. WDFY4, EGR1,
607 IGF1, INF2, PRSS1, etc.) are shown in Supplementary File-10.

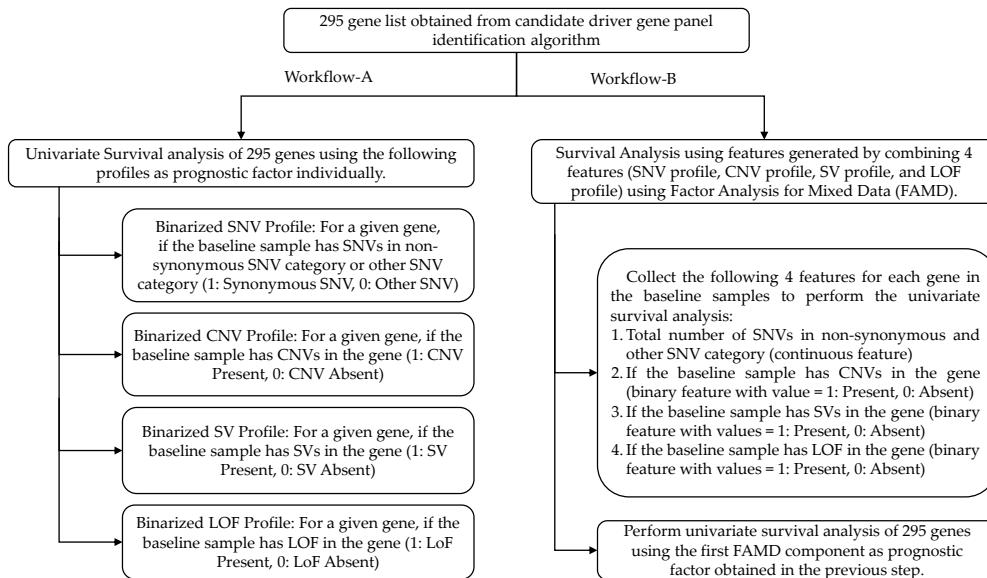


Figure 8: Workflow for two-fold survival analysis of proposed 295-gene panel. In this workflow, we estimated the clinical relevance of gene variant profiles on MM patient clinical outcomes using two distinct approaches. In the first approach (Workflow-A), We individually assessed the impact of each variant profile (SNV, CNV, SV, and LOF) on clinical outcomes. Univariate survival analysis was performed for each variant profile, providing insights into their impact. Using this approach, 167 out of the 295 genes significantly influenced clinical outcomes in univariate survival analysis based on at least one prognostic factor. In the second approach (Workflow-B), we amalgamated the four variant profiles (SNV, CNV, SV, and LOF) for each gene using Factor Analysis for Mixed Data (FAMD), enabling us to estimate a joint feature. Subsequently, we performed univariate survival analysis using the FAMD 1st component as a prognostic factor. In this approach, 187 out of 295 genes demonstrated significance in univariate survival analysis based on the FAMD 1st component (a combined feature generated by integrating the four variant profiles for each gene). Interestingly, 129 genes out of these 173 were also identified as significant in univariate survival analysis using Workflow-A. By combining both approaches, out of 295, 225 genes were found to influence the clinical outcomes of MM patients significantly.

608 4. Discussion

609 Multiple Myeloma (MM) is a malignancy that typically progresses from premalignant stages,
610 often starting with MGUS [84]. A targeted sequencing panel is important for the precise characteri-
611 zation of the genomic alterations to understand the risk of progression, enabling timely interventions
612 and ultimately improving patient outcomes. Recent studies have shed light on the genomic events
613 that drive the transformation from premalignant stages to MM [85, 86, 87, 88]. Moreover, a number

614 of studies have proposed targeted sequencing panels for molecular profiling of MM patients based
615 on previously identified genomic events in MM and MGUS [17, 18, 19, 20, 21]. However, none
616 of these studies have taken into account the design of the panel using biomarkers and gene-gene
617 interactions that have the potential to distinguish MM from MGUS.

618 In this study, we addressed this challenge by designing a targeted sequencing panel of 295 genes
619 hosting key genomic biomarkers. We designed an AI-powered attention-based bio-inspired BIO-
620 DGI (PPI9) model to identify the key genomic biomarkers and gene interactions for panel crafting.
621 The BIO-DGI (PPI9) model is biologically inspired, learning to identify distinguishing patterns of
622 MM and MGUS using gene-gene interactions and their corresponding genomic features. Genes with
623 a higher number of interactions are deemed more biologically relevant. We specifically considered
624 deleterious SNVs associated with MM and MGUS, resulting in highly MM-relevant, significantly
625 altered genes being ranked at the top. The inclusivity of three global repositories having MM and
626 MGUS cohorts with diverse ethnicities, the ability of the AI-based workflow to comprehend gene
627 inter-dependencies, extensive benchmarking, and rigorous post-hoc analysis collectively render the
628 BIO-DGI (PPI9) model innovative and highly efficient.

629 During classification, the functional significance of nonsynonymous SNVs, as quantified by Phy-
630 lop scores, emerged as the most prominent genomic feature. Following closely, the allele depth
631 of synonymous SNVs and the overall count of other SNVs (encompassing non-frame-shift inser-
632 tions/deletions/substitutions, intronic, intergenic, ncRNA_intronic, upstream, downstream, un-
633 known, and ncRNA_splicing SNVs) ranked as the second and third most influential genomic features,
634 respectively (Figure-4). These findings are in line with the literature because the impact of synony-
635 mous SNVs across various cancer types has been highlighted by various studies [89, 90, 91, 92, 93].

636 In the post-hoc analysis for model interpretability, we utilized the ShAP algorithm to identify
637 the top-ranked genes within the top-performing models. Table-1A, Table-1B, Table-2 provides an
638 overview of the total number of previously reported genes present in the 798 significantly altered
639 genes and those identified by top-performing models, presenting complete gene lists under top-250
640 and top 500 ranks. Notably, the BIO-DGI (PPI9) model outperformed by identifying the highest
641 number of previously reported genes, encompassing known oncogenes (OGs) such as BIRC6, MUC4,

642 NOTCH1, PGR, SETD1A, and VAV1, tumour suppressor genes (TSGs) like DIS3, EP400, MYH11,
643 SDHA, both oncogenes and driver genes (ODGs) such as KRAS, NRAS, TP53, TRRAP, and
644 actionable genes (AGs) including APC, ARID1B, MITF, NFKBIA, and TYRO3. Interestingly,
645 most of these genes (except ODGs) display high relevance to MM despite not being explicitly
646 reported as MM driver genes.

647 Additionally, our analysis identified MUC6, LILRA1, and LILRB1 as the top three genes con-
648 tributing significantly to the classification of MM and MGUS, although none of these have been
649 previously categorized as OGs, TSGs, ODGs, or AGs in the literature. Here, this is to note that
650 MUC6 gene is associated with the immune system pathway, playing a crucial role in MM devel-
651 opment and progression [94]. Similarly, the other two genes, LILRA1 and LILRB1, are associated
652 with the innate immune system pathway. Notably, LILRB1 has been reported to be associated with
653 MM pathogenesis as an inhibitory immune checkpoint for B-cell function in prior studies [95, 96].

654 We also employed the Geo2R tool to validate the top-ranked genes obtained from the post-hoc
655 analysis of the BIO-DGI (PPI9) model. We included 11 MM-related studies for validation and
656 observed that 488 out of 500 genes were found to be disrupted in MM. This finding ensures the
657 relevance of top-ranked genes in MM.

658 We curated a 295-gene panel by rigorously analysing variant profiles (SNVs, CNVs, SVs, and
659 LOF) of the top 500 genes. We specifically considered MM-relevant genes that were disrupted in at
660 least one previously published MM study. To identify pivotal genomic events responsible for MM
661 development and progression, we categorized them into different groups based on their occurrence
662 at specific disease stages (MM or both MM and MGUS). Genomic events observed in both MGUS
663 and MM such as translocations associated with the IGH and MYC genes [86, 13, 97, 98, 99, 100]
664 and amp(1q) [13] were labeled by us as “disease-initiating” events, while those that were observed
665 to be present in MM but not in MGUS including del(13q), del(16q), del(17p), etc. [13] were labeled
666 by us as “disease-transformative” events, and are shown in Table-4A and Table-4B.

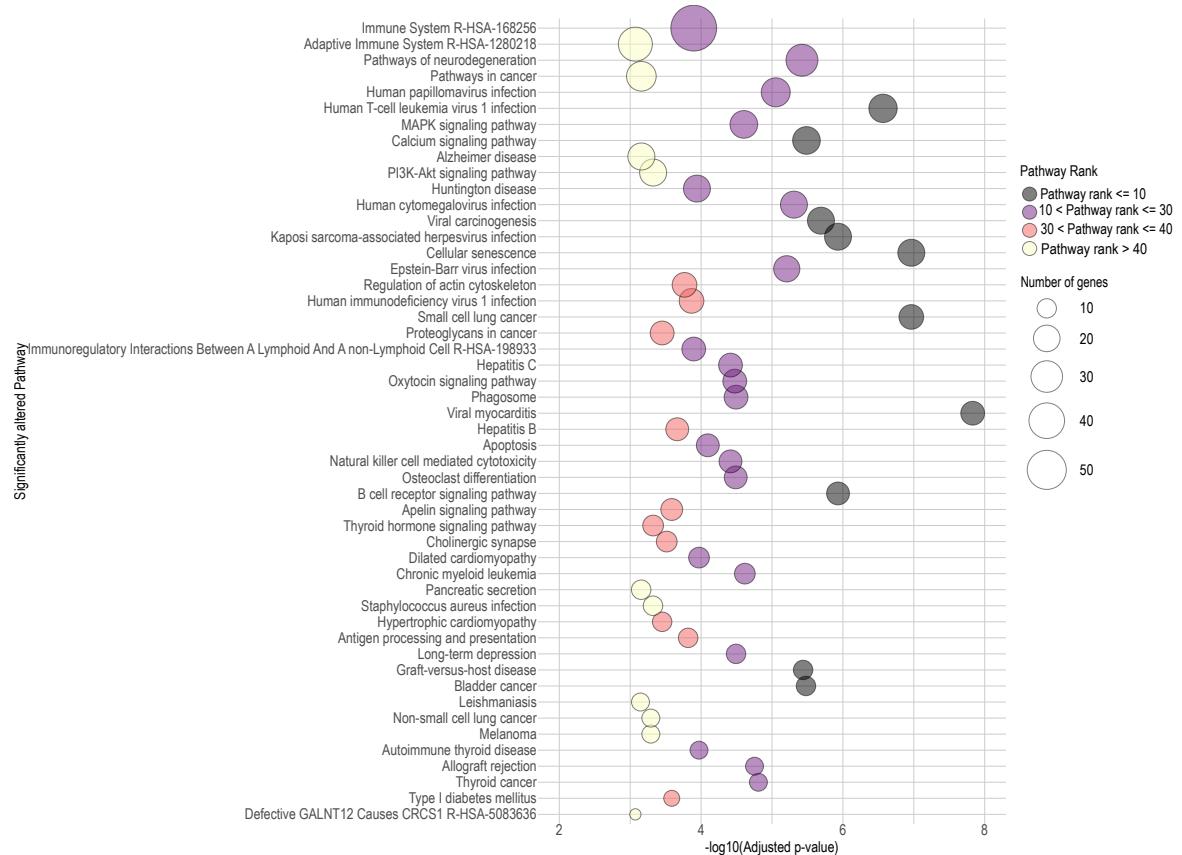


Figure 9: Bubble Plot illustrating the Top 50 significantly altered signalling pathways in MM associated with genes included in the proposed 295-gene panel. The bubble size indicates the number of significantly altered genes linked to the 295-gene panel, and colour signifies the pathway rank. The x-axis represents the $-\log_{10}$ (adjusted p -value) score of the pathway, and the y-axis displays the pathway names.

Table 4: List of key genomic events in MM and MGUS and overlapping of their associated genes with 295 genes panel.

(A) List of disease-initiating genomic events reported previously in both MM and MGUS and the overlapping of their associated genes with our proposed 295-genes panel

S.No.	A. Genomic Events in MM and MGUS	B. Genes associated with the event (Column-A)	C. References for the genes shown in column-B	D. Whether the genes of Column-B are present in 295-genes panel	E. If yes, list of genes common with Column-B	F. Associated gene-gene interactions found by BIO-DGI (PPI9) in the 295-gene panel
1	t(11;14)	CCND1, BCL-2	[86, 97, 98, 99]	Yes	CCND1	BRD4, IRS1, KRAS, NRAS, TP53
2	t(4;14)	FGFR3	[86, 13, 98, 99, 100]	Yes	FGFR3	CYLD, DIS3, KRAS, NRAS
3	t(14;16)	MAF	[98, 99, 101, 86]	Yes	MAF	FLNA
4	t(14;20)	MAFB	[86, 99, 101]	Yes	MAFB	HUWE1, USP9X
5	Amp(1q21) Del(17p13)	MCL1, CKS1B, ANP32E or BCL9 TP53 KRAS NRAS	[13] [102] [86] [86]	Yes	CKS1B	KPRP, BRD4, PLEC, USP9X
6	KRAS mutations	LTB	[86]	Yes	TP53	IRF4, KMT2B/C/D, KRAS
7	NRAS Mutations	DIS3	[86]	Yes	KRAS	MAX, BRAF, EGFR, NF1
8	LTB Mutations	EGR1	[86]	Yes	NRAS	ARID2, FGFR3, SP140
9	DIS3 mutations	IGH, IGL, IGK, NSMCE2, TXNDC5, FAM46C, FOXO3, IGJ, PRDM1	[10]	Yes	LTB	IRFL, NFKB1A, SP140
10	EGR1 mutations			Yes	DIS3	FGFR3, KRAS, NRAS
11	MYC Rearrangement			Yes	EGR1	CYLD, FGFR3, KRAS, NRAS
12				Yes	IGH, IGL, IGK,FAM46C	CYLD, DIS3, FGFR3, KRAS

(B) List of disease-transformative genomic events reported previously in MM but not in MGUS and the overlapping of their associated genes with our proposed 295 genes panel.

S.No.	A. Disease-Transformative Genomic Events	B. Genes associated with the event (Column-A)	C. References for the genes shown in column-B	D. Whether the genes of Column-B are present in 295-genes panel	E. If yes, genes present in the 295-genes panel	F. Associated gene-gene alterations found by BIO-DGI (PPI9) in the 295-genes panel
1	Del(13q14)	RBL	[103, 13]	Yes	RBL, DIS3	BRAF, FGFR3, KRAS
3	Del(16q23)	CYLD	[13]	Yes	CYLD	DIS3, KRAS, NRAS
4	Del(1p21)	FAM46C	[13, 104]	Yes	FAM46C	CYLD, DIS3, FGFR3, KRAS
5	Del(12p13)	CD27	[105]	Yes	CD27	TRAF2, TRAF3, ATP2B3
6	TP53 Mutations	TP53	[65, 106]	Yes	TP53	IRF4, KMT2B/C/D, KRAS
7	BRAF Mutations	BRAF	[107, 108]	Yes	BRAF	FGFR3, KRAS, NRAS
8	Gain(9q)	ABCA1, KCNT1, TRAF2, VPS13A	[109, 110]	Yes	ABCA1, KCNT1, TRAF2, VPS13A	BRAF, CYLD, IRF1
9	del(14q)	TRAF3	[13]	Yes	TRAF3	CYLD, DIS3, KMT2D
10	del(17p)	TP53	[111]	Yes	TP53	IRF4, KMT2B/C/D, KRAS
11	del(8p)	PTK2B, TP53	[112]	Yes	PTK2B	EGR1, FGFR3, KRAS

667 We comprehensively analysed CNVs, SVs, and LOFs identified in the 295 genes panel across
668 multiple myeloma (MM) samples obtained from both AIIMS and MMRF datasets. Our analysis
669 highlighted chr1, chr14, chr19, and chrX as the most affected chromosomes, displaying various CNV
670 genomic alterations. Notably, chr1 exhibited significant alterations, such as amp(1q), associated
671 with disease aggressiveness [13, 113], and del(1p), frequently observed in MGUS [13, 104]. Fur-
672 thermore, chr14 revealed prevalent translocations involving IGH, such as t(4;14), t(14;16), t(14;20),
673 established as biomarkers in MM [13]. Additionally, CNVs linked to chr19, such as gain(19p) and
674 gain(19q), were significantly more prevalent in MM than in MGUS [84]. Recently, it was shown
675 that abnormalities of chromosome X and MAGE-C1/CT7 expression are much more frequent events
676 in MM than previously reported [114]. The intricate interplay between alterations in these chro-
677 mosomes and other genetic events contributes to increased genomic instability, facilitating the
678 acquisition of additional mutations that promote MM aggressiveness [11].

679 Upon scrutinizing the clinical significance of the proposed panel consisting of 295 genes through
680 survival analysis across five variant profiles (SNV, CNV, SV, LOF, and FAMD 1st component of
681 amalgamation of four variant profiles), seven genes (*ACACB*, *ARHGAP4*, *ASKIN2*, *FAM186A*,
682 *IGFN1*, *NBPF9*, and *TPTE2*) demonstrated clinical significance in at least four variant pro-
683 files. Notably, *ACACB* and *TPTE2* were identified as vulnerable genes in Multiple Myeloma Cells
684 through RNA Interference Lethality Screening of the Druggable Genome [115]. *ACACB* might have
685 been playing a pivotal role in Multiple Myeloma (MM) progression because its top transcription
686 factor binding sites such as AP-1, C/EBPalpha, MAZR, RFX1, STAT1, STAT1 α , and STAT1 β play
687 a role in either cell proliferation, differentiation, apoptosis, oncogenesis or in regulating the immune
688 response. For example, the role of MAZR in cell proliferation, apoptosis, and tumorigenesis may
689 indicate its contribution to MM progression by affecting ACACB regulation. Additionally, RFX1,
690 influencing the cell cycle and immune response, could also play a role in ACACB modulation within
691 the context of MM. Lastly, the engagement of STAT1, along with its isoforms STAT1 α and STAT1 β
692 in immune response and tumour suppression might implicate ACACB in MM pathogenesis, poten-
693 tially linking aberrant fatty acid metabolism to the dysregulated immune responses characteristic
694 of the disease. The intricate interplay between ACACB and these transcription factors underscores

695 its multifaceted involvement in MM progression.

696 Similarly, *TPTE2*'s involvement in diverse cellular processes, including immune responses, in-
697flammation, and cell survival, critical aspects of MM progression, is suggested by its association
698with NF-kappaB and its subunits (*NF-kappaB1* and *NF-kappaB2*). Moreover, RelA binding sites
699signify *TPTE2*'s involvement in the activation of gene expression in response to stimuli. Altogether,
700*TPTE2*'s engagement with these transcription factors indicates its intricate participation in diverse
701cellular processes, potentially contributing to the complex pathogenesis of multiple myeloma. Fur-
702ther investigations are warranted to elucidate the precise molecular mechanisms and implications
703of *TPTE2* in MM progression.

704 We thoroughly evaluated our proposed 295-gene panel, comparing it with five previously pub-
705lished targeted sequencing panels used for MM genomic profiling. These panels were thoughtfully
706crafted based on MM-related literature and underwent validation using diverse methods such as
707FISH and analysis of whole-genome sequencing (WGS) data, etc. Upon analyzing the validated
708variant profiles, we noted that, alongside our proposed panel, Sudha et al. [18] also validated their
709panel on SNVs, CNVs, and SVs, encompassing translocations linked to IGH and MYC. However,
710Sudha et al.'s panel validation was carried out on WGS cohorts of MM samples and MM cell lines
711and did not account for potentially distinguishing genomic biomarkers of MGUS and MM. Further-
712more, none of the targeted panels developed so far analyzed the clinical significance of individual
713genes in their panels on survival outcomes. The present study is unique and adds valuable informa-
714tion on the potential impact of these genes on clinical outcomes in MM. Thus, the proposed panel
715is unique in that it not only helps identify transforming events in patients with MGUS, but also is
716powered to assess genomic features that impact treatment outcomes.

717 Moreover, our panel incorporated MM-relevant genes exhibiting loss-of-function (LOF), a critical
718consideration lacking in previous panels. Comparing the genes across the previously published
719panels, we found that 19 out of 47 (34%) genes were common with Kortum et al.'s, 26 out of 182
720(10.43%) with Bolli et al.'s, 47 out of 465 (8.38%) with White et al.'s, 18 out of 26 (57.69%) with
721Cutler et al.'s, and 42 out of 228 (14.5%) with Sudha et al.'s panels, respectively. The comprehensive
722gene list encompassing all genes from the five panels is provided in Table-S18, Supplementary File-

Table 5: Comparison of previously published targeted sequencing panels with our proposed 295-genes panel

S. No.	Panel Reference, Publication year	Total number of genes in the proposed gene panel	Number of samples used for panel validation	Data Type	Detected variant profiles	Overlapping with 295-genes panel
1	Kortum et al [19], 2015	47	22 NDMM, 3 pretreated MM samples	WES	SNVs, clonal evolution analysis	19
2	Bolli et al [20], 2016	182	5 MM samples	WGS	SNVs, CNVs, SVs*	26
3	White et al [21], 2018	465	110 MM samples 76 (20 MGUS, 3 SMM, 52 MM, and 1 PCL) samples	WGS	SNVs, CNVs, SVs*	47
4	Cutler et al [17], 2021	26		WGS	SNVs, CNVs, Clinical validation using survival analysis	18
5	Sudha et al [18], 2022	228	185 MM samples	WGS	SNVs, CNVs, SVs*	42
6	Vivek et al (Current study)	295	1215 (1154 MM and 61 MGUS) samples + 11 MM-datasets	WES, microarray, mRNA	SNVs, CNVs, SVs, clinical validation using two-fold survival analysis	-

723 9. Additionally, a detailed comparison of these panels is presented in Table-5 and Table-S19,
 724 Supplementary File-9.

725 In addition, we conducted pathway analysis using the Enrichr database to elucidate the signif-
 726 icantly altered pathways associated with the 295-gene panel. These pathways were subsequently
 727 ranked based on their statistical significance (adjusted *p*-value) to identify the top pathways ex-
 728 hibiting substantial alterations. Notably, a distinct pattern emerged when assessing the significance
 729 of altered pathways in relation to disease progression. Pathways associated with various cellular
 730 processes displayed significant alterations in MGUS, but their significance diminished as the dis-
 731 ease transitioned from MGUS to MM. In contrast, pathways specifically linked to multiple myeloma
 732 exhibited pronounced alterations as the disease advanced (Table-S10, S11, Supplementary File-4).
 733 Out of the 295 genes, 174 were implicated in significantly altered pathways. Key MM-related path-
 734 ways, including MAPK signaling, PI3K-AKT, B-cell receptor, Human papillomavirus, and immune
 735 system pathways, prominently featured among the significantly altered pathways. The top 50 path-
 736 ways associated with the proposed 295-gene panel, along with the number of significantly altered
 737 genes and their respective rankings for each pathway, are illustrated in the bubble plot presented in

738 Figure-9. These compelling findings warrant further investigation to determine whether the signifi-
739 cantly altered genes associated with these pathways could potentially serve as valuable biomarkers
740 of the development of MM during the early stages of the disease, particularly MGUS.

741 Using the interaction weights acquired from the BIO-DGI (PPI9) model, we identified five gene
742 communities and to enhance the information for each node within a gene community, we integrated
743 node influence determined by the Katz centrality score and likelihood of haploinsufficiency gauged
744 through the GHIS score. The genes surpassing the median GHIS score of 0.52 (Figure 6) notably
745 include, UBC, USP6, PRIM2, USP34, KMT2C, PABPC1, and NCOR1 in the first gene community
746 (Figure-6(A)), TP53, NRAS, IRS1, EIF4EBP1, HSP90AB1, and FGFR3 in the second gene com-
747 munity (Figure-6(B)), POTE in the third gene community (Figure-6(C)), and LILRA1, LILRB1,
748 LILRB2, FCGR2A in the fourth gene community (Figure 6(D)) and appear as central genes that
749 shows that these might have been playing a significant role in MM pathogenesis. This is to note
750 that these genes are already known to be associated with MM. Furthermore, we observed that, out
751 of the 295 genes, 67 displayed substantial node influence within the gene community, encompassing
752 various previously reported MM-relevant genes like *BRAF*, *HLA-A/B*, *FGFR*, *IRS1*, *NRAS*, and
753 *SDHA*. Additionally, 74 genes exhibited a high likelihood of haploinsufficiency, including several
754 previously reported MM-relevant genes such as *ARID1B*, *FGFR*, *NRAS*, *TRAF2*, and *ZNF717*.
755 Moreover, 32 genes displayed both high node influence and a high likelihood of haploinsufficiency,
756 including *FGFR*, *HUWE1*, *KRAS*, *KMT2C/D*, *TP53*, and *ZNF717*. We strongly recommend fur-
757 ther analysis on these central genes to unveil their role in disease progression.

758 While examining the gene communities and their involvement in key genomic events of MM,
759 we noted several genes with substantial node influence and likelihood of actively participating in
760 these events. For instance, in the first gene community (Figure-6(A)), seven genes (*BRD4*, *DIS3*,
761 *HUWE1*, *RB1*, *SLC25A5*, *RB1*, and *USP9X*) were associated with genomic events observed in both
762 MM and MGUS. Similarly, the second gene community (Figure-6(B)) included five genes (*EGR1*,
763 *IRS1*, *KRAS*, *NRAS*, and *TP53*) involved in genomic events observed in both MM and MGUS. In
764 the third gene community, *FLNA* was found to be associated with genomic events observed in both
765 MM and MGUS. The fourth community featured *LTB* associated with genomic events observed

766 in both MM and MGUS, while *TRAF3* was associated with genomic events observed in MM only.
767 Finally, the fifth gene community had no genes linked to the key genomic events shown in Table-
768 4A and Table-4B. The presence of genes actively participating in MM-related key genomic events,
769 displaying high node influence within the community, and a high likelihood of haploinsufficiency
770 underscores the relevance of our proposed targeted sequencing panel in MM and MGUS.

771 In this study, we investigated the interactions of the 295-gene panel with drugs by leveraging
772 the DGIdb database [116]. Our focus was on identifying drugs associated with these genes based
773 on the strength of evidence for interaction, considering factors such as the number of publica-
774 tions and sources supporting each claim. The top 15 drugs with the most robust gene interactions
775 were prioritized, involving well-known genes like *BRAF*, *KRAS*, *FGFR*, *TP53*, among others. The
776 resulting gene-drug interaction network, depicted in Figure 10, was both weighted and directed.
777 Among the top three drugs demonstrating interactions with key driver genes in multiple myeloma
778 were Dabrafenib, Trametinib, and Vemurafenib. Notably, Vemurafenib and Dabrafenib are rec-
779 ognized as BRAF inhibitors (BRAFi) [117] and have been considered in the context of multiple
780 myeloma treatment [118, 119]. Additionally, Cisplatin, a drug known for its inclusion in the potent
781 combination therapy VTD-PACE (bortezomib-thalidomide-dexamethasone-cisplatin-doxorubicin-
782 cyclophosphamide-etoposide) [120], demonstrated significant interactions. Furthermore, our anal-
783 ysis highlighted drugs commonly used in the treatment of other cancers, such as Carboplatin, Do-
784 cetaxel, and Fluorouracil, showing interactions with key driver genes in multiple myeloma. These
785 findings suggest the potential of these drugs as novel chemotherapeutic agents for multiple myeloma.

786 5. Conclusions

787 Distinguishing Multiple Myeloma from its precursor stage, MGUS, and identifying those at risk
788 of progression to overt MM poses a significant challenge due to overlapping genomic characteristics.
789 The present study addresses this challenge by leveraging the innovative AI-based BIO-DGI (PPI9)
790 model, incorporating gene interactions from nine PPI databases and exonic mutational profiles
791 from global MM and MGUS repositories (AIIMS, EGA, and MMRF). Our study demonstrates
792 superior quantitative and qualitative performance with application-aware interpretability. The

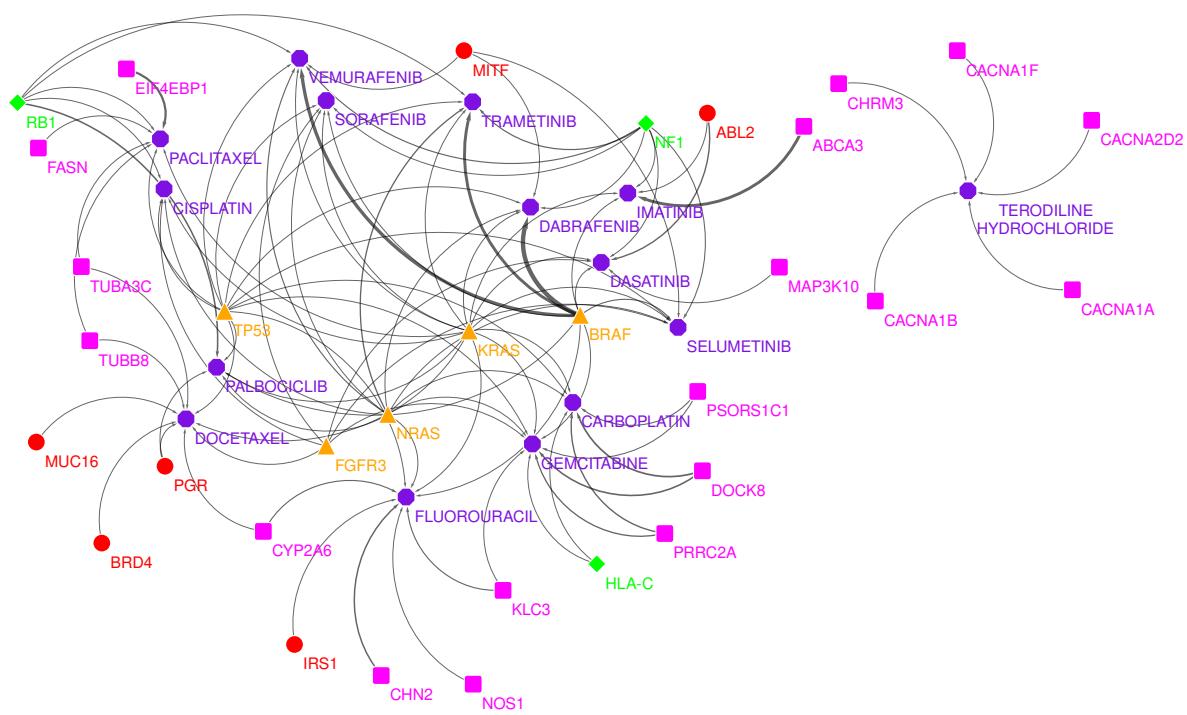


Figure 10: Gene-drug interactions for genes in 295 genes panel. In the above network, the colour and shape notations are as follows: **red circle** represent oncogenes, **orange triangle** represent genes that are both oncogenes and driver genes (ODG), **magenta square** represent genes that are not previously reported as OG/TSG/ODG/AG, **purple octagon** represent drugs. The edge width represents the mean of gene-drug interaction scores obtained from DGIdb 4.0 database [116].

model identified a substantial number of previously reported genes, including Oncogenes (OGs), Tumor Suppressor Genes (TSGs), Oncogene Databases (ODGs), and Associated Genes (AGs), known for their high relevance in MM. Geo2R validation of the 295-gene panel, coupled with an association with MM-relevant pathways, underscores the panel's pertinence to MM. Our analysis highlights the functional significance of non-synonymous mutations, allele depth of synonymous SNVs, and the total number of other SNVs as crucial genomic biomarkers in distinguishing MM from MGUS. Significant alterations in chromosomes 1, 14, 19, and X suggest the inclusion of genes associated with these chromosomes in genomic evaluation of myeloma patients. Within this panel, we identified genes with substantial node influence and prominent gene-gene interactions from five gene communities, shedding light on crucial gene biomarkers and their interactions pivotal to

803 MM pathogenesis. These findings hold immense potential for informed therapeutic interventions,
804 facilitating early detection and interception of disease progression in MM. The proposed panel,
805 driven by innovative AI modelling and comprehensive genomic analysis, emerges as a promising
806 tool for advancing our understanding of MM and improving patient outcomes through precision
807 medicine.

808 **6. Acknowledgements**

809 This work was supported by a grant from the Department of Biotechnology, Govt. of India
810 [Grant: BT/MED/30/SP11006/2015] and the Department of Science and Technology, Govt. of In-
811 dia [Grant: DST/ICPS/CPS-Individual/2018/279(G)]. Authors acknowledge dbGaP (Project#18964)
812 for providing authorized access to the MM datasets (phs000748 and phs000348). The Multiple
813 Myeloma Research Foundation provided funding support for the phs000348 study in collabora-
814 tion with the Multiple Myeloma Research Consortium. The Broad Institute Genome Sequencing,
815 Genetic Analysis, and Biological Samples Platforms provided assistance with data generation, pro-
816 cessing and analysis. The datasets used for the analyses described in this work were obtained from
817 dbGaP through dbGaP accession number phs000348.v1.p1. Data of study phs000748 were gen-
818 erated as part of the Multiple Myeloma Research Foundation CoMMpass [SM] (Relating Clinical
819 Outcomes in MM to Personal Assessment of Genetic Profile) study (www.themmrf.org). We also
820 acknowledge EGA (EGAD00001001901) for providing authorized access to the MGUS data. The
821 authors would also like to thank the Centre of Excellence in Healthcare, IIIT-Delhi, for their sup-
822 port in this research. Authors acknowledge the valuable insights provided by Dr Satish Sankaran,
823 Dr. Nandini Pal Basak, and Dr. Mohit Malhotra from Farcast Biosciences Pvt Ltd, India, that
824 greatly improved the gene panel designing workflow.

825 **Bibliography**

826 **References**

- 827 [1] R. A. Kyle, T. M. Therneau, S. V. Rajkumar, J. R. Offord, D. R. Larson, M. F. Plevak,
828 L. J. Melton III, A long-term study of prognosis in monoclonal gammopathy of undetermined
829 significance, *New England Journal of Medicine* 346 (8) (2002) 564–569.
- 830 [2] S. V. Rajkumar, Mgs and smoldering multiple myeloma: update on pathogenesis, natural
831 history, and management, *ASH Education Program Book* 2005 (1) (2005) 340–345.
- 832 [3] A. Laganà, I. Beno, D. Melnekoff, V. Leshchenko, D. Madduri, D. Ramdas, L. Sanchez,
833 S. Niglio, D. Perumal, B. A. Kidd, et al., Precision medicine for relapsed multiple myeloma
834 on the basis of an integrative multiomics approach, *JCO precision oncology* 2 (2018) 1–17.
- 835 [4] A. Palumbo, H. Avet-Loiseau, S. Oliva, H. M. Lokhorst, H. Goldschmidt, L. Rosinol,
836 P. Richardson, S. Caltagirone, J. J. Lahuerta, T. Facon, et al., Revised international staging
837 system for multiple myeloma: a report from international myeloma working group, *Journal
838 of clinical oncology* 33 (26) (2015) 2863.
- 839 [5] S. A. Holstein, P. L. McCarthy, Immunomodulatory drugs in multiple myeloma: mechanisms
840 of action and clinical experience, *Drugs* 77 (2017) 505–520.
- 841 [6] V. Shah, D. C. Johnson, A. L. Sherborne, S. Ellis, F. M. Aldridge, J. Howard-Reeves, F. Be-
842 gum, A. Price, J. Kendall, L. Chieccchio, et al., Subclonal tp53 copy number is associated with
843 prognosis in multiple myeloma, *Blood, The Journal of the American Society of Hematology*
844 132 (23) (2018) 2465–2469.
- 845 [7] A. Mikulasova, C. Ashby, R. G. Tytarenko, P. Qu, A. Rosenthal, J. A. Dent, K. R. Ryan, M. A.
846 Bauer, C. P. Wardell, A. Hoering, et al., Microhomology-mediated end joining drives complex
847 rearrangements and overexpression of myc and pvt1 in multiple myeloma, *Haematologica*
848 105 (4) (2020) 1055.

- 849 [8] N. Abdallah, L. B. Baughn, S. V. Rajkumar, P. Kapoor, M. A. Gertz, A. Dispenzieri, M. Q.
850 Lacy, S. R. Hayman, F. K. Buadi, D. Dingli, et al., Implications of myc rearrangements in
851 newly diagnosed multiple myeloma, *Clinical Cancer Research* 26 (24) (2020) 6581–6588.
- 852 [9] S. Manier, K. Salem, S. V. Glavey, A. M. Roccaro, I. M. Ghobrial, Genomic aberrations in
853 multiple myeloma, *Plasma Cell Dyscrasias* (2016) 23–34.
- 854 [10] M. Affer, M. Chesi, W. Chen, J. J. Keats, Y. N. Demchenko, K. Tamizhmani, V. Gar-
855 bitt, D. Riggs, L. Brents, A. Roschke, et al., Promiscuous myc locus rearrangements hijack
856 enhancers but mostly super-enhancers to dysregulate myc expression in multiple myeloma,
857 *Leukemia* 28 (8) (2014) 1725–1735.
- 858 [11] N. Bolli, H. Avet-Loiseau, D. C. Wedge, P. Van Loo, L. B. Alexandrov, I. Martincorena, K. J.
859 Dawson, F. Iorio, S. Nik-Zainal, G. R. Bignell, et al., Heterogeneity of genomic evolution and
860 mutational profiles in multiple myeloma, *Nature communications* 5 (1) (2014) 2997.
- 861 [12] B. A. Walker, K. Mavrommatis, C. P. Wardell, T. C. Ashby, M. Bauer, F. E. Davies, A. Rosen-
862 thal, H. Wang, P. Qu, A. Hoering, et al., Identification of novel mutational drivers reveals
863 oncogene dependencies in multiple myeloma, *Blood, The Journal of the American Society of*
864 *Hematology* 132 (6) (2018) 587–597.
- 865 [13] S. Manier, K. Z. Salem, J. Park, D. A. Landau, G. Getz, I. M. Ghobrial, Genomic complexity
866 of multiple myeloma and its clinical implications, *Nature reviews Clinical oncology* 14 (2)
867 (2017) 100–113.
- 868 [14] J. B. Egan, C.-X. Shi, W. Tembe, A. Christoforides, A. Kurdoglu, S. Sinari, S. Middha,
869 Y. Asmann, J. Schmidt, E. Braggio, et al., Whole-genome sequencing of multiple myeloma
870 from diagnosis to plasma cell leukemia reveals genomic initiating events, evolution, and clonal
871 tides, *Blood, The Journal of the American Society of Hematology* 120 (5) (2012) 1060–1066.
- 872 [15] B. A. Walker, K. Mavrommatis, C. P. Wardell, T. C. Ashby, M. Bauer, F. Davies, A. Rosen-
873 thal, H. Wang, P. Qu, A. Hoering, et al., A high-risk, double-hit, group of newly diagnosed
874 myeloma identified by genomic analysis, *Leukemia* 33 (1) (2019) 159–170.

- 875 [16] M.-V. Mateos, J. F. San Miguel, Management of multiple myeloma in the newly diagnosed
876 patient, Hematology 2014, the American Society of Hematology Education Program Book
877 2017 (1) (2017) 498–507.
- 878 [17] S. D. Cutler, P. Knopf, C. J. Campbell, A. Thoni, M. Abou El Hassan, N. Forward, D. White,
879 J. Wagner, M. Goudie, J. E. Boudreau, et al., Dmg26: A targeted sequencing panel for
880 mutation profiling to address gaps in the prognostication of multiple myeloma, The Journal
881 of Molecular Diagnostics 23 (12) (2021) 1699–1714.
- 882 [18] P. Sudha, A. Ahsan, C. Ashby, T. Kausar, A. Khera, M. H. Kazeroun, C.-C. Hsu, L. Wang,
883 E. Fitzsimons, O. Salminen, et al., Myeloma genome project panel is a comprehensive targeted
884 genomics panel for molecular profiling of patients with multiple myeloma, Clinical Cancer
885 Research 28 (13) (2022) 2854–2864.
- 886 [19] K. Kortüm, C. Langer, J. Monge, L. Bruins, Y. Zhu, C. Shi, P. Jedlowski, J. Egan, J. Ojha,
887 L. Bullinger, et al., Longitudinal analysis of 25 sequential sample-pairs using a custom multiple
888 myeloma mutation sequencing panel (m 3 p), Annals of hematology 94 (2015) 1205–1211.
- 889 [20] N. Bolli, Y. Li, V. Sathiaseelan, K. Raine, D. Jones, P. Ganly, F. Cocito, G. Bignell, M. A.
890 Chapman, A. Sperling, et al., A dna target-enrichment approach to detect mutations, copy
891 number changes and immunoglobulin translocations in multiple myeloma, Blood Cancer Jour-
892 nal 6 (9) (2016) e467–e467.
- 893 [21] B. S. White, I. Lanc, J. O’Neal, H. Gupta, R. S. Fulton, H. Schmidt, C. Fronick, E. A.
894 Belter Jr, M. Fiala, J. King, et al., A multiple myeloma-specific capture sequencing platform
895 discovers novel translocations and frequent, risk-associated point mutations in igll5, Blood
896 cancer journal 8 (3) (2018) 35.
- 897 [22] X. Du, S. Sun, C. Hu, Y. Yao, Y. Yan, Y. Zhang, Deepppi: boosting prediction of protein–
898 protein interactions with deep neural networks, Journal of chemical information and modeling
899 57 (6) (2017) 1499–1510.

- 900 [23] S.-B. Zhang, Q.-R. Tang, Protein–protein interaction inference based on semantic similarity
901 of gene ontology terms, *Journal of theoretical biology* 401 (2016) 30–37.
- 902 [24] S. R. Maetschke, M. Simonsen, M. J. Davis, M. A. Ragan, Gene ontology-driven inference of
903 protein–protein interactions using inducers, *Bioinformatics* 28 (1) (2012) 69–75.
- 904 [25] I. Ieremie, R. M. Ewing, M. Niranjan, Transformergo: predicting protein–protein interactions
905 by modelling the attention between sets of gene ontology terms, *Bioinformatics* 38 (8) (2022)
906 2269–2277.
- 907 [26] R. Schulte-Sasse, S. Budach, D. Hnisz, A. Marsico, Integration of multiomics data with graph
908 convolutional networks to identify new cancer genes and their associated molecular mecha-
909 nisms, *Nature Machine Intelligence* 3 (6) (2021) 513–526.
- 910 [27] V. Ruhela, L. Jena, G. Kaur, R. Gupta, A. Gupta, Bdl-sp: A bio-inspired dl model for the
911 identification of altered signaling pathways in multiple myeloma using wes data, *American*
912 *Journal of Cancer Research* 13 (4) (2023) 1155.
- 913 [28] V. A. Traag, L. Waltman, N. J. Van Eck, From louvain to leiden: guaranteeing well-connected
914 communities, *Scientific reports* 9 (1) (2019) 5233.
- 915 [29] L. Katz, A new status index derived from sociometric analysis, *Psychometrika* 18 (1) (1953)
916 39–43.
- 917 [30] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, *Advances in*
918 *neural information processing systems* 30 (2017).
- 919 [31] J. J. Keats, D. W. Craig, W. Liang, Y. Venkata, A. Kurdoglu, J. Aldrich, D. Auclair, K. Allen,
920 B. Harrison, S. Jewell, et al., Interim analysis of the mmrf comppass trial, a longitudinal study
921 in multiple myeloma relating clinical outcomes to genomic and immunophenotypic profiles
922 (2013).
- 923 [32] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,

- 924 N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learn-
925 ing library, *Advances in neural information processing systems* 32 (2019).
- 926 [33] T. Therneau, et al., A package for survival analysis in s, R package version 2 (7) (2015).
- 927 [34] G. A. Van der Auwera, M. O. Carneiro, C. Hartl, R. Poplin, G. Del Angel, A. Levy-Moonshine,
928 T. Jordan, K. Shakir, D. Roazen, J. Thibault, et al., From fastq data to high-confidence
929 variant calls: the genome analysis toolkit best practices pipeline, *Current protocols in bioin-*
930 *formatics* 43 (1) (2013) 11–10.
- 931 [35] Y. Fan, L. Xi, D. S. Hughes, J. Zhang, J. Zhang, P. A. Futreal, D. A. Wheeler, W. Wang, Muse:
932 accounting for tumor heterogeneity using a sample-specific error model improves sensitivity
933 and specificity in mutation calling from sequencing data, *Genome biology* 17 (1) (2016) 1–11.
- 934 [36] D. Benjamin, T. Sato, K. Cibulskis, G. Getz, C. Stewart, L. Lichtenstein, Calling somatic
935 snvs and indels with mutect2, *BioRxiv* (2019) 861054.
- 936 [37] D. C. Koboldt, Q. Zhang, D. E. Larson, D. Shen, M. D. McLellan, L. Lin, C. A. Miller, E. R.
937 Mardis, L. Ding, R. K. Wilson, Varscan 2: somatic mutation and copy number alteration
938 discovery in cancer by exome sequencing, *Genome research* 22 (3) (2012) 568–576.
- 939 [38] D. E. Larson, C. C. Harris, K. Chen, D. C. Koboldt, T. E. Abbott, D. J. Dooling, T. J. Ley,
940 E. R. Mardis, R. K. Wilson, L. Ding, Somaticsniper: identification of somatic point mutations
941 in whole genome sequencing data, *Bioinformatics* 28 (3) (2012) 311–317.
- 942 [39] K. Wang, M. Li, H. Hakonarson, Annovar: functional annotation of genetic variants from
943 high-throughput sequencing data, *Nucleic acids research* 38 (16) (2010) e164–e164.
- 944 [40] M. F. Rogers, H. A. Shihab, M. Mort, D. N. Cooper, T. R. Gaunt, C. Campbell, Fathmm-
945 xf: accurate prediction of pathogenic point mutations via extended features, *Bioinformatics*
946 34 (3) (2018) 511–513.
- 947 [41] I. Martincorena, K. M. Raine, M. Gerstung, K. J. Dawson, K. Haase, P. Van Loo, H. Davies,

- 948 M. R. Stratton, P. J. Campbell, Universal patterns of selection in cancer and somatic tissues,
949 Cell 171 (5) (2017) 1029–1041.
- 950 [42] R. Oughtred, C. Stark, B.-J. Breitkreutz, J. Rust, L. Boucher, C. Chang, N. Kolas,
951 L. O'Donnell, G. Leung, R. McAdam, et al., The biogrid interaction database: 2019 update,
952 Nucleic acids research 47 (D1) (2019) D529–D541.
- 953 [43] E. L. Huttlin, R. J. Bruckner, J. Navarrete-Perea, J. R. Cannon, K. Baltier, F. Gebreab,
954 M. P. Gygi, A. Thornock, G. Zarraga, S. Tam, et al., Dual proteome-scale networks reveal
955 cell-specific remodeling of the human interactome, Cell 184 (11) (2021) 3022–3040.
- 956 [44] E. Persson, M. Castresana-Aguirre, D. Buzzao, D. Guala, E. L. Sonnhammer, Funcoup 5:
957 functional association networks in all domains of life, supporting directed links and tissue-
958 specificity, Journal of Molecular Biology 433 (11) (2021) 166835.
- 959 [45] G. Alanis-Lobato, M. A. Andrade-Navarro, M. H. Schaefer, Hippie v2. 0: enhancing meaning-
960 fulness and reliability of protein–protein interaction networks, Nucleic acids research (2016)
961 gkw985.
- 962 [46] C. Y. Kim, S. Baek, J. Cha, S. Yang, E. Kim, E. M. Marcotte, T. Hart, I. Lee, Humannet v3:
963 an improved database of human gene networks for disease research, Nucleic acids research
964 50 (D1) (2022) D632–D639.
- 965 [47] F. Zheng, M. R. Kelly, D. J. Ramms, M. L. Heintschel, K. Tao, B. Tutuncuoglu, J. J. Lee,
966 K. Ono, H. Foussard, M. Chen, et al., Interpretation of cancer mutations using a multiscale
967 map of protein systems, Science 374 (6563) (2021) eabf3067.
- 968 [48] G. Kustatscher, P. Grabowski, T. A. Schrader, J. B. Passmore, M. Schrader, J. Rappsilber,
969 Co-regulation map of the human proteome enables identification of protein functions, Nature
970 biotechnology 37 (11) (2019) 1361–1371.
- 971 [49] M. Gillespie, B. Jassal, R. Stephan, M. Milacic, K. Rothfels, A. Senff-Ribeiro, J. Griss,

- 972 C. Sevilla, L. Matthews, C. Gong, et al., The reactome pathway knowledgebase 2022, Nucleic
973 acids research 50 (D1) (2022) D687–D692.
- 974 [50] D. Szklarczyk, R. Kirsch, M. Koutrouli, K. Nastou, F. Mehryary, R. Hachilif, A. L. Gable,
975 T. Fang, N. T. Doncheva, S. Pyysalo, et al., The string database in 2023: protein–protein
976 association networks and functional enrichment analyses for any sequenced genome of interest,
977 Nucleic acids research 51 (D1) (2023) D638–D646.
- 978 [51] C. D. Huber, B. Y. Kim, K. E. Lohmueller, Population genetic models of gerp scores suggest
979 pervasive turnover of constrained sites across mammalian evolution, PLoS genetics 16 (5)
980 (2020) e1008827.
- 981 [52] K. S. Pollard, M. J. Hubisz, K. R. Rosenbloom, A. Siepel, Detection of nonneutral substitution
982 rates on mammalian phylogenies, Genome research 20 (1) (2010) 110–121.
- 983 [53] A. Siepel, G. Bejerano, J. S. Pedersen, A. S. Hinrichs, M. Hou, K. Rosenbloom, H. Clawson,
984 J. Spieth, L. W. Hillier, S. Richards, et al., Evolutionarily conserved elements in vertebrate,
985 insect, worm, and yeast genomes, Genome research 15 (8) (2005) 1034–1050.
- 986 [54] B. Reva, Y. Antipin, C. Sander, Predicting the functional impact of protein mutations: ap-
987 plication to cancer genomics, Nucleic acids research 39 (17) (2011) e118–e118.
- 988 [55] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
989 P. Prettenhofer, R. Weiss, V. Dubourg, et al., Scikit-learn: Machine learning in python, the
990 Journal of machine Learning research 12 (2011) 2825–2830.
- 991 [56] D. Chakravarty, J. Gao, S. Phillips, R. Kundra, H. Zhang, J. Wang, J. E. Rudolph, R. Yaeger,
992 T. Soumerai, M. H. Nissan, et al., Oncokb: a precision oncology knowledge base, JCO preci-
993 sion oncology 1 (2017) 1–16.
- 994 [57] F. Martínez-Jiménez, F. Muiños, I. Sentís, J. Deu-Pons, I. Reyes-Salazar, C. Arnedo-Pac,
995 L. Mularoni, O. Pich, J. Bonet, H. Kranas, et al., A compendium of mutational cancer driver
996 genes, Nature Reviews Cancer 20 (10) (2020) 555–572.

- 997 [58] J. G. Tate, S. Bamford, H. C. Jubb, Z. Sondka, D. M. Beare, N. Bindal, H. Boutselakis, C. G.
998 Cole, C. Creatore, E. Dawson, et al., Cosmic: the catalogue of somatic mutations in cancer,
999 Nucleic acids research 47 (D1) (2019) D941–D947.
- 1000 [59] S. De Cesco, J. B. Davis, P. E. Brennan, Targetdb: A target information aggregation tool
1001 and tractability predictor, PloS one 15 (9) (2020) e0232644.
- 1002 [60] F. Maura, N. Bolli, N. Angelopoulos, K. J. Dawson, D. Leongamornlert, I. Martincorena,
1003 T. J. Mitchell, A. Fullam, S. Gonzalez, R. Szalat, et al., Genomic landscape and chronological
1004 reconstruction of driver events in multiple myeloma, Nature communications 10 (1) (2019)
1005 3835.
- 1006 [61] E. Talevich, A. H. Shain, T. Botton, B. C. Bastian, Cnvkit: genome-wide copy number
1007 detection and visualization from targeted dna sequencing, PLoS computational biology 12 (4)
1008 (2016) e1004873.
- 1009 [62] T. Rausch, T. Zichner, A. Schlattl, A. M. Stütz, V. Benes, J. O. Korbel, Delly: structural
1010 variant discovery by integrated paired-end and split-read analysis, Bioinformatics 28 (18)
1011 (2012) i333–i339.
- 1012 [63] N. Huang, I. Lee, E. M. Marcotte, M. E. Hurles, Characterising and predicting haploinsuffi-
1013 ciency in the human genome, PLoS genetics 6 (10) (2010) e1001154.
- 1014 [64] T. Barrett, S. E. Wilhite, P. Ledoux, C. Evangelista, I. F. Kim, M. Tomashevsky, K. A.
1015 Marshall, K. H. Phillippe, P. M. Sherman, M. Holko, et al., Ncbi geo: archive for functional
1016 genomics data sets—update, Nucleic acids research 41 (D1) (2012) D991–D995.
- 1017 [65] A. Farswan, A. Gupta, R. Gupta, G. Kaur, Imputation of gene expression data in blood
1018 cancer and its significance in inferring biological pathways, Frontiers in oncology 9 (2020)
1019 1442.
- 1020 [66] A. Garcia-Gomez, T. Li, C. de la Calle-Fabregat, J. Rodríguez-Ubreva, L. Ciudad, F. Català-
1021 Moll, G. Godoy-Tena, M. Martín-Sánchez, L. San-Segundo, S. Muntiόn, et al., Targeting

- 1022 aberrant dna methylation in mesenchymal stromal cells as a treatment for myeloma bone
1023 disease, *Nature communications* 12 (1) (2021) 421.
- 1024 [67] N. C. Gutiérrez, M. E. Sarasquete, I. Misiewicz-Krzeminska, M. Delgado, J. De Las Rivas,
1025 F. Ticona, E. Ferminan, P. Martin-Jimenez, C. Chillon, A. Risueno, et al., Deregulation of
1026 microrna expression in the different genetic subtypes of multiple myeloma and correlation
1027 with gene expression profiling, *Leukemia* 24 (3) (2010) 629–637.
- 1028 [68] Y. Zhou, L. Chen, B. Barlogie, O. Stephens, X. Wu, D. R. Williams, M.-A. Cartron, F. van
1029 Rhee, B. Nair, S. Waheed, et al., High-risk myeloma is associated with global elevation of
1030 mirnas and overexpression of eif2c2/ago2, *Proceedings of the National Academy of Sciences*
1031 107 (17) (2010) 7904–7909.
- 1032 [69] M. Lionetti, M. Biasiolo, L. Agnelli, K. Todoerti, L. Mosca, S. Fabris, G. Sales, G. L. Deliliers,
1033 S. Bicciato, L. Lombardi, et al., Identification of microrna expression patterns and definition
1034 of a microrna/mrna regulatory network in distinct molecular groups of multiple myeloma,
1035 *Blood, The Journal of the American Society of Hematology* 114 (25) (2009) e20–e26.
- 1036 [70] M. Biasiolo, G. Sales, M. Lionetti, L. Agnelli, K. Todoerti, A. Bisognin, A. Coppe, C. Ro-
1037 mualdi, A. Neri, S. Bortoluzzi, Impact of host genes and strand selection on mirna and mirna*
1038 expression, *PloS one* 6 (8) (2011) e23854.
- 1039 [71] N. Amodio, M. Di Martino, U. Foresta, E. Leone, M. Lionetti, M. Leotta, A. Gullà, M. Pitari,
1040 F. Conforti, M. Rossi, et al., mir-29b sensitizes multiple myeloma cells to bortezomib-induced
1041 apoptosis through the activation of a feedback loop with the transcription factor sp1, *Cell
1042 death & disease* 3 (11) (2012) e436–e436.
- 1043 [72] M. Bolzoni, P. Storti, S. Bonomini, K. Todoerti, D. Guasco, D. Toscani, L. Agnelli, A. Neri,
1044 V. Rizzoli, N. Giuliani, Immunomodulatory drugs lenalidomide and pomalidomide inhibit
1045 multiple myeloma-induced osteoclast formation and the rankl/opg ratio in the myeloma mi-
croenvironment targeting the expression of adhesion molecules, *Experimental hematology*
1046 41 (4) (2013) 387–397.

- 1048 [73] D. Ronchetti, K. Todoerti, G. Tuana, L. Agnelli, L. Mosca, M. Lionetti, S. Fabris, P. Co-
1049 lapietro, M. Miozzo, M. Ferrarini, et al., The expression pattern of small nucleolar and small
1050 cajal body-specific rnas characterizes distinct molecular subtypes of multiple myeloma, *Blood*
1051 cancer journal
- 2 (11) (2012) e96–e96.
- 1052 [74] I. P. N. Bong, C. C. Ng, N. Othman, E. Esa, Gene expression profiling and in vitro functional
1053 studies reveal rad54l as a potential therapeutic target in multiple myeloma, *Genes & Genomics*
1054 44 (8) (2022) 957–966.
- 1055 [75] J. R. Nair, J. Caserta, K. Belko, T. Howell, G. Fetterly, C. Baldino, K. P. Lee, Novel inhibition
1056 of pim2 kinase has significant anti-tumor efficacy in multiple myeloma, *Leukemia* 31 (8) (2017)
1057 1715–1726.
- 1058 [76] A. Sacco, C. Federico, K. Todoerti, B. Ziccheddu, V. Palermo, A. Giacomini, C. Ravelli,
1059 F. Maccarinelli, G. Bianchi, A. Belotti, et al., Specific targeting of the kras mutational land-
1060 scape in myeloma as a tool to unveil the elicited antitumor activity, *Blood*, The Journal of
1061 the American Society of Hematology 138 (18) (2021) 1705–1720.
- 1062 [77] D. Soncini, C. Martinuzzi, P. Becherini, E. Gelli, S. Ruberti, K. Todoerti, L. Mastracci,
1063 P. Contini, A. Cagnetta, A. Laudisi, et al., Apoptosis reprogramming triggered by splicing
1064 inhibitors sensitizes multiple myeloma cells to venetoclax treatment, *Haematologica* 107 (6)
1065 (2022) 1410.
- 1066 [78] J. Navarro Gonzalez, A. S. Zweig, M. L. Speir, D. Schmelter, K. R. Rosenbloom, B. J. Raney,
1067 C. C. Powell, L. R. Nassar, N. D. Maulding, C. M. Lee, et al., The ucsc genome browser
1068 database: 2021 update, *Nucleic acids research* 49 (D1) (2021) D1046–D1057.
- 1069 [79] J. Pagès, Multiple factor analysis by example using R, CRC Press, 2014.
- 1070 [80] M. V. Kuleshov, M. R. Jones, A. D. Rouillard, N. F. Fernandez, Q. Duan, Z. Wang, S. Koplev,
1071 S. L. Jenkins, K. M. Jagodnik, A. Lachmann, et al., Enrichr: a comprehensive gene set
1072 enrichment analysis web server 2016 update, *Nucleic acids research* 44 (W1) (2016) W90–
1073 W97.

- 1074 [81] Z. Xie, A. Bailey, M. V. Kuleshov, D. J. Clarke, J. E. Evangelista, S. L. Jenkins, A. Lachmann,
1075 M. L. Wojciechowicz, E. Kropiwnicki, K. M. Jagodnik, et al., Gene set knowledge discovery
1076 with enrichr, *Current protocols* 1 (3) (2021) e90.
- 1077 [82] E. Y. Chen, C. M. Tan, Y. Kou, Q. Duan, Z. Wang, G. V. Meirelles, N. R. Clark, A. Ma'ayan,
1078 Enrichr: interactive and collaborative html5 gene list enrichment analysis tool, *BMC bioin-*
1079 *formatics* 14 (1) (2013) 1–14.
- 1080 [83] J. Steinberg, F. Honti, S. Meader, C. Webber, Haploinsufficiency predictions without study
1081 bias, *Nucleic acids research* 43 (15) (2015) e101–e101.
- 1082 [84] L. Lopez-Corral, M. E. Sarasquete, S. Beà, R. García-Sanz, M. V. Mateos, L. Corchete,
1083 J. Sayagués, E. García, J. Bladé, A. Oriol, et al., Snp-based mapping arrays reveal high
1084 genomic complexity in monoclonal gammopathies, from mgus to myeloma status, *Leukemia*
1085 26 (12) (2012) 2521–2529.
- 1086 [85] B. A. Walker, C. P. Wardell, L. Melchor, A. Brioli, D. C. Johnson, M. F. Kaiser, F. Mirabella,
1087 L. Lopez-Corral, S. Humphray, L. Murray, et al., Intraclonal heterogeneity is a critical early
1088 event in the development of myeloma and precedes the development of clinical symptoms,
1089 *Leukemia* 28 (2) (2014) 384–390.
- 1090 [86] A. Mikulasova, C. P. Wardell, A. Murison, E. M. Boyle, G. H. Jackson, J. Smetana, Z. Kufova,
1091 L. Pour, V. Sandecka, M. Almasi, et al., The spectrum of somatic mutations in monoclonal
1092 gammopathy of undetermined significance indicates a less complex genomic landscape than
1093 that in multiple myeloma, *Haematologica* 102 (9) (2017) 1617.
- 1094 [87] A. Farswan, A. Gupta, L. Jena, V. Ruhela, G. Kaur, R. Gupta, Characterizing the mutational
1095 landscape of mm and its precursor mgus, *American Journal of Cancer Research* 12 (4) (2022)
1096 1919.
- 1097 [88] A. K. Dutta, J. L. Fink, J. P. Grady, G. J. Morgan, C. G. Mullighan, L. B. To, D. R. Hewett,
1098 A. C. Zannettino, Subclonal evolution in disease progression from mgus/smm to multiple
1099 myeloma is characterised by clonal stability, *Leukemia* 33 (2) (2019) 457–468.

- 1100 [89] D. Chu, L. Wei, Nonsynonymous, synonymous and nonsense mutations in human cancer-
1101 related genes undergo stronger purifying selections than expectation, BMC cancer 19 (1)
1102 (2019) 1–12.
- 1103 [90] Y. Sharma, M. Miladi, S. Dukare, K. Boulay, M. Caudron-Herger, M. Groß, R. Backofen,
1104 S. Diederichs, A pan-cancer analysis of synonymous mutations, Nature communications 10 (1)
1105 (2019) 2569.
- 1106 [91] T. Soussi, P. E. Taschner, Y. Samuels, Synonymous somatic variants in human cancer are not
1107 infamous: a plea for full disclosure in databases and publications, Human mutation 38 (4)
1108 (2017) 339–342.
- 1109 [92] H. Teng, W. Wei, Q. Li, M. Xue, X. Shi, X. Li, F. Mao, Z. Sun, Prevalence and architecture
1110 of posttranscriptionally impaired synonymous mutations in 8,320 genomes across 22 cancer
1111 types, Nucleic acids research 48 (3) (2020) 1192–1205.
- 1112 [93] F. Supek, B. Miñana, J. Valcárcel, T. Gabaldón, B. Lehner, Synonymous mutations frequently
1113 act as driver mutations in human cancers, Cell 156 (6) (2014) 1324–1335.
- 1114 [94] K. Nakamura, M. J. Smyth, L. Martinet, Cancer immunoediting and immune dysregulation
1115 in multiple myeloma, Blood, The Journal of the American Society of Hematology 136 (24)
1116 (2020) 2731–2740.
- 1117 [95] E. Lozano, T. Díaz, M.-P. Mena, G. Suñe, X. Calvo, M. Calderón, L. Pérez-Amill,
1118 V. Rodríguez, P. Pérez-Galán, G. Roué, et al., Loss of the immune checkpoint cd85j/lilrb1
1119 on malignant plasma cells contributes to immune escape in multiple myeloma, The Journal
1120 of Immunology 200 (8) (2018) 2581–2591.
- 1121 [96] X. Kang, J. Kim, M. Deng, S. John, H. Chen, G. Wu, H. Phan, C. C. Zhang, Inhibitory
1122 leukocyte immunoglobulin-like receptors: Immune checkpoint proteins and tumor sustaining
1123 factors, Cell cycle 15 (1) (2016) 25–40.

- 1124 [97] L. López-Corral, N. C. Gutiérrez, M. B. Vidriales, M. V. Mateos, A. Rasillo, R. García-Sanz,
1125 B. Paiva, J. F. San Miguel, The progression from mgus to smoldering myeloma and eventually
1126 to multiple myeloma involves a clonal expansion of genetically abnormal plasma cells, Clinical
1127 cancer research 17 (7) (2011) 1692–1700.
- 1128 [98] J. Bladé, Monoclonal gammopathy of undetermined significance, New England Journal of
1129 Medicine 355 (26) (2006) 2765–2770.
- 1130 [99] N. Korde, S. Y. Kristinsson, O. Landgren, Monoclonal gammopathy of undetermined signifi-
1131 cance (mgus) and smoldering multiple myeloma (smm): novel biological insights and develop-
1132 ment of early treatment strategies, Blood, The Journal of the American Society of Hematology
1133 117 (21) (2011) 5573–5581.
- 1134 [100] S. A. Van Wier, G. J. Ahmann, K. J. Henderson, P. R. Greipp, S. V. Rajkumar, D. M.
1135 Larson, A. Dispenzieri, M. A. Gertz, R. A. Kyle, R. Fonseca, The t (4; 14) is present in
1136 patients with early stage plasma cell proliferative disorders including mgus and smoldering
1137 multiple myeloma (smm)., Blood 106 (11) (2005) 1545.
- 1138 [101] F. M. Ross, L. Chiechino, G. Dagradá, R. K. Protheroe, D. M. Stockley, C. J. Harrison, N. C.
1139 Cross, A. J. Szubert, M. T. Drayson, G. J. Morgan, The t (14; 20) is a poor prognostic factor
1140 in myeloma but is associated with long-term stable disease in monoclonal gammopathies of
1141 undetermined significance, haematologica 95 (7) (2010) 1221.
- 1142 [102] K. Neben, A. Jauch, T. Hielscher, J. Hillengass, N. Lehnert, A. Seckinger, M. Granzow, M. S.
1143 Raab, A. D. Ho, H. Goldschmidt, et al., Progression in smoldering myeloma is independently
1144 determined by the chromosomal abnormalities del (17p), t (4; 14), gain 1q, hyperdiploidy,
1145 and tumor load, Journal of clinical oncology 31 (34) (2013) 4325–4332.
- 1146 [103] Z. He, J. O’Neal, W. C. Wilson, N. Mahajan, J. Luo, Y. Wang, M. Y. Su, L. Lu, J. B. Skeath,
1147 D. Bhattacharya, et al., Deletion of rb1 induces both hyperproliferation and cell death in
1148 murine germinal center b cells, Experimental hematology 44 (3) (2016) 161–165.

- 1149 [104] H. Chang, X. Qi, A. Jiang, W. Xu, T. Young, D. Reece, 1p21 deletions are strongly associated
1150 with 1q21 gains and are an independent adverse prognostic factor for the outcome of high-
1151 dose chemotherapy in patients with multiple myeloma, *Bone marrow transplantation* 45 (1)
1152 (2010) 117–121.
- 1153 [105] F. Li, Y. Xu, P. Deng, Y. Yang, W. Sui, F. Jin, M. Hao, Z. Li, M. Zang, D. Zhou, et al.,
1154 Heterogeneous chromosome 12p deletion is an independent adverse prognostic factor and
1155 resistant to bortezomib-based therapy in multiple myeloma, *Oncotarget* 6 (11) (2015) 9434.
- 1156 [106] K. K. Jovanović, G. Escure, J. Demonchy, A. Willaume, Z. Van de Wyngaert, M. Farhat,
1157 P. Chauvet, T. Facon, B. Quesnel, S. Manier, Deregulation and targeting of tp53 pathway in
1158 multiple myeloma, *Frontiers in oncology* 8 (2019) 665.
- 1159 [107] S. F. Mohamed, M. Khan, A. Quesada, J. Ma, P. Lin, C. C. Yin, K. Sasaki, G. Borthakur,
1160 N. Pemmaraju, Q. Bashir, et al., Disease characteristics of multiple myeloma involving braf
1161 mutations, *Blood* 138 (2021) 4755.
- 1162 [108] S. Pasca, C. Tomuleasa, P. Teodorescu, G. Ghiaur, D. Dima, V. Moisoiu, C. Berce, C. Ste-
1163 fan, A. Ciechanover, H. Einsele, Kras/nras/braf mutations as potential targets in multiple
1164 myeloma, *Frontiers in Oncology* 9 (2019) 1137.
- 1165 [109] P. Liebisch, C. Wendt, A. Wellmann, A. Kröber, G. Schilling, H. Goldschmidt, H. Einsele,
1166 C. Straka, M. Bentz, S. Stilgenbauer, et al., High incidence of trisomies 1q, 9q, and 11q in
1167 multiple myeloma: results from a comprehensive molecular cytogenetic analysis, *Leukemia*
1168 17 (12) (2003) 2535–2537.
- 1169 [110] P. Liebisch, D. Scheck, S. A. Erné, A. Wellmann, C. Wendt, S. Janczik, S. Kolmus, A. Kröber,
1170 H. Einsele, C. Straka, et al., Duplication of chromosome arms 9q and 11q: Evidence for a
1171 novel, 14q32 translocation-independent pathogenetic pathway in multiple myeloma, *Genes,*
1172 *Chromosomes and Cancer* 42 (1) (2005) 78–81.
- 1173 [111] A. Aktas Samur, S. Minvielle, M. Shammas, M. Fulciniti, F. Magrangeas, P. G. Richardson,

- 1174 P. Moreau, M. Attal, K. C. Anderson, G. Parmigiani, et al., Deciphering the chronology of
1175 copy number alterations in multiple myeloma, *Blood cancer journal* 9 (4) (2019) 39.
- 1176 [112] A. D. Duru, T. Sutlu, A. Wallblom, K. Uttervall, J. Lund, B. Stellan, G. Gahrton, H. Nahi,
1177 E. Alici, Deletion of chromosomal region 8p21 confers resistance to bortezomib and is asso-
1178 ciated with upregulated decoy trail receptor expression in patients with multiple myeloma,
1179 *PLoS One* 10 (9) (2015) e0138248.
- 1180 [113] T. M. Schmidt, R. Fonseca, S. Z. Usmani, Chromosome 1q21 abnormalities in multiple
1181 myeloma, *Blood cancer journal* 11 (4) (2021) 83.
- 1182 [114] J. Balcarkova, P. Flodr, N. Svobodova, T. Pika, P. Krhovska, T. Papajik, H. Urbankova, J. Mi-
1183 narik, Aberrations of chromosome x in patients with multiple myeloma, *Clinical Lymphoma,*
1184 *Myeloma and Leukemia* 19 (10) (2019) e56–e57.
- 1185 [115] R. E. Tiedemann, Y. X. Zhu, J. Schmidt, C. X. Shi, C. Sereduk, H. Yin, S. Mousses, A. K.
1186 Stewart, Identification of molecular vulnerabilities in human multiple myeloma cells by rna
1187 interference lethality screening of the druggable genome, *Cancer research* 72 (3) (2012) 757–
1188 768.
- 1189 [116] S. L. Freshour, S. Kiwala, K. C. Cotto, A. C. Coffman, J. F. McMichael, J. J. Song, M. Griffith,
1190 O. L. Griffith, A. H. Wagner, Integration of the drug–gene interaction database (dgidb 4.0)
1191 with open crowdsourcing efforts, *Nucleic acids research* 49 (D1) (2021) D1144–D1151.
- 1192 [117] N. Raje, I. Chau, D. M. Hyman, V. Ribrag, J.-Y. Blay, J. Tabernero, E. Elez, J. Wolf, A. J.
1193 Yee, M. Kaiser, et al., Vemurafenib in patients with relapsed refractory multiple myeloma
1194 harboring brafv600 mutations: a cohort of the histology-independent ve-basket study, *JCO*
1195 *Precision Oncology* 2 (2018).
- 1196 [118] V. Subbiah, R. J. Kreitman, Z. A. Wainberg, A. Gazzah, U. Lassen, A. Stein, P. Y. Wen,
1197 S. Dietrich, M. J. de Jonge, J.-Y. Blay, et al., Dabrafenib plus trametinib in brafv600e-mutated
1198 rare cancers: the phase 2 roar trial, *Nature medicine* (2023) 1–10.

- 1199 [119] National cancer institute (NCI). targeted therapy directed by genetic testing in treating pa-
1200 tients with advanced refractory solid tumors, lymphomas, or multiple myeloma (the match
1201 screening trial). nlm identifier: NCT02465060; 2020.
- 1202 [120] T. Togano, S. Andoh, M. Komuro, Y. Mitsui, S. Itoi, R. Hirai, M. Nakamura, A. Tan-
1203 imura, R. Sekine, M. Takeshita, et al., Bortezomib-thalidomide-dexamethasone-cisplatin-
1204 doxorubicin-cyclophosphamide-etoposide as a salvage and bridging regimen before hematopoietic
1205 stem cell transplantation for relapsed or refractory multiple myeloma, Internal Medicine
1206 61 (22) (2022) 3329–3334.