

Deepening the knowledge of rare diseases dependent on angiogenesis through semantic similarity clustering and network analysis

Raquel Pagano-Márquez, José Córdoba-Caballero, Beatriz Martínez-Poveda, Ana R. Quesada, Elena Rojano, Pedro Seoane, Juan AG. Ranea[†] and Miguel Ángel Medina[†]

Corresponding authors: Pedro Seoane, CIBERER, University of Malaga, 29010 Malaga, Spain. Tel.: +34-952132025; E-mail: seoanezonjic@uma.es; Elena Rojano, IBIMA Plataforma BIONAND, University of Malaga, 29010 Malaga, Spain. Tel.: +34-952132025; E-mail: elenarojano@uma.es.

[†]Juan A G Ranea and Miguel Ángel Medina contributed equally to this work.

Abstract

Background: Angiogenesis is regulated by multiple genes whose variants can lead to different disorders. Among them, rare diseases are a heterogeneous group of pathologies, most of them genetic, whose information may be of interest to determine the still unknown genetic and molecular causes of other diseases. In this work, we use the information on rare diseases dependent on angiogenesis to investigate the genes that are associated with this biological process and to determine if there are interactions between the genes involved in its deregulation.

Results: We propose a systemic approach supported by the use of pathological phenotypes to group diseases by semantic similarity. We grouped 158 angiogenesis-related rare diseases in 18 clusters based on their phenotypes. Of them, 16 clusters had traceable gene connections in a high-quality interaction network. These disease clusters are associated with 130 different genes. We searched for genes associated with angiogenesis through ClinVar pathogenic variants. Of the seven retrieved genes, our system confirms six of them. Furthermore, it allowed us to identify common affected functions among these disease clusters.

Availability: https://github.com/ElenaRojano/angio_cluster.

Contact: seoanezonjic@uma.es and elenarojano@uma.es

Keywords: rare diseases, systems biology, semantic similarity, disease clustering, angiogenesis

Introduction

Angiogenesis deregulation is associated with a large number of diseases, including different types of cancer, autoimmune and rare diseases [3, 39, 50]. This biological process is complex in molecular terms and its regulation is susceptible to changes in the genome [52]. Despite being an essential process for the maintenance of the organism there is much more to investigate about the genes and the regulation of this process.

It is known that genes involved in diseases described with similar phenotypes can be functionally related on

the molecular level [43]. Consequently, to deepen in the knowledge of the angiogenesis deregulation, it can be possible to analyze what phenotypic similarities there are between angiogenesis-dependent diseases. There are different approaches that use semantic similarity to calculate how similar two diseases are using Gene Ontology terms [25], the Disease Ontology [4] and the Human Phenotype Ontology (HPO) [24]. Here, we consider that phenotypically similar angiogenesis-dependent diseases can be grouped to determine which genes they have in common and to relate them to the angiogenesis deregulation.

Raquel Pagano-Márquez is a PhD student at the Department of Molecular Biology and Biochemistry, University of Malaga, Spain. Her research focuses on the analysis and interpretation of rare diseases dependent on angiogenesis.

José Córdoba-Caballero is a PhD student at the Department of Molecular Biology and Biochemistry, University of Malaga, Spain. His research focuses on the development of bioinformatics tools for the analysis of patients with rare diseases.

Beatriz Martínez-Poveda is an Associate Professor at the Department of Molecular Biology and Biochemistry, University of Malaga, Spain. Her research focuses on analyzing the molecular factors involved in angiogenesis-related diseases, including cancer and atherosclerosis.

Ana R. Quesada is a professor at the Department of Molecular Biology and Biochemistry, University of Malaga, Spain. She leads a research group focuses on the analysis of the molecular signaling pathways associated with angiogenesis and the search for modulating compounds of angiogenic activity.

Elena Rojano is a post-doctoral researcher in bioinformatics at the Department of Molecular Biology and Biochemistry, University of Malaga, Spain. Her research focuses on systems biology methods for associating pathological phenotypes to genomic variants from patients with rare diseases.

Pedro Seoane-Zonjic is a post-doctoral researcher in bioinformatics at the CIBER of Rare Diseases (CIBERER), University of Malaga, Spain. His current research focuses on developing software for the analysis of rare diseases from high-throughput genomic data.

Juan AG. Ranea is a professor at the Department of Molecular Biology and Biochemistry, University of Malaga, Spain. He leads a research group focused on bioinformatics and systems biology for the analysis of rare diseases.

Miguel Ángel Medina is a professor at the Department of Molecular Biology and Biochemistry, University of Malaga, Spain. He leads a research group focused on systems biology, angiogenesis and cancer research, as well as rare diseases.

Received: October 18, 2021. Revised: April 28, 2022. Accepted: May 11, 2022

© The Author(s) 2022. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

In this work, we rely on the analysis of angiogenesis-related rare diseases (A-RDs) to analyze the genetic factors involved in the angiogenesis deregulation. This heterogeneous group of diseases provide information of multiple pathological phenotypes and genes that is useful to help to understand the angiogenesis dysregulation.

In 2012, our research group performed a systemic review of A-RDs, with a manual search of A-RDs used to identify disease-associated genes and available drugs for their treatment using Orphanet resources [39]. Here, we update the list of A-RDs and use it to get a better understanding of the A-RDs and their associated genes. To do so, we use a semantic similarity measure and clustering analysis of these A-RDs. There are plenty of studies that describe several similarity measures [31] applied for disease analysis at the phenotypic level as the rare disease map (RDmap) [59] or for the differential diagnostics for common diseases [47]. Furthermore, the use of semantic similarity for data clustering is frequent for the stratification of patients in cohorts [48] or disease groups [2]. In addition, this methodology is also used for the identification of genes that could be involved in disease development. For example, there are essential resources such as the Monarch Initiative [28] and DisGeNET [33] that can be used to retrieve gene information, and tools like Priorit-T [36] that uses information from MEDLINE abstracts for gene prioritization.

In this way, we apply this knowledge to A-RDs designing a full analysis protocol to analyze a pool of diseases. Our motivation is to give a reliable and straightforward insight of the rare diseases related to angiogenesis, grouping them at a phenotypic level. With this stratification, we explored each A-RD group at a genetic level to analyze the biological functions of the associated genes and which other genes could be related. Finally, we identified close A-RD groups at both phenotypic and genetic levels, giving a reference point to A-RD researchers to elucidate the disease mechanisms.

Material and methods

A-RDs selection

Information concerning A-RDs was compiled with the criteria described in the work of Rodríguez-Caso and collaborators [39]. Following their procedure, we performed an advanced search for specific terms that emerged in publications of any year. This search was performed in the Web of Science (WOS) and PubMed databases. We personalized our search in the following way, according to the database consulted: in the case of WOS, we searched for terms related to rare diseases and angiogenesis '(TS = (rare disease AND angiogen*))', whereas in PubMed were used '((rare diseases [MeSH Terms]) OR (rare AND diseases) OR (rare diseases) OR (rare AND disease) OR (rare disease) AND angiogen*)'. Both searches were made in August 2021. We exported

the results of the articles corresponding to this search and eliminated the repeated records.

In the same line as Rodríguez-Caso and collaborators followed in their study, [39], to perform the A-RDs articles selection we made a search for terms in the title and abstract of all articles, and specifically in the abstract keywords and MESH terms of PubMed articles and in the author keyword and keyword plus of WOS articles. We searched for two groups of keywords: 'angiogen' and 'VEGF' to verify that the article had content about this biological process, or 'rare' and 'disease' to confirm that there were rare diseases mentioned in the article. All the articles that did not meet these search requirements were removed from the study.

We calculated a content score to prioritize articles according to where these terms were included in the publications. This score is calculated in the following way: if the searched terms are in the title, we add 3 points to the score, 2 if they are included in the keywords or MeSH terms and 1 if it is in the abstract. Using the content score, we focused on articles whose score was equal to or greater than 4 as we considered them as the most relevant.

We manually inspected these articles to verify that they were describing A-RDs. All diseases resulting from manual curation were searched in Orphanet to get an official ORPHA code. These ORPHA codes will be used to find additional information about A-RDs in different databases.

For each A-RD ORPHA code we retrieved both their phenotype and gene annotations. In the case of phenotypes, we selected all the HPOs related to each A-RD ORPHA codes from the HPO annotation website [19] (<http://purl.obolibrary.org/obo/hp/hpoa/phenotype.hpoa>). The genes associated with A-RDs were retrieved from the Monarch Initiative resources [28] (https://data.monarchinitiative.org/tsv/all_associations/gene_disease.all.tsv.gz). More information about the list of A-RDs, their HPOs and associated genes is available in Supplementary Table 4.

Disease workflow analysis overview

We developed a workflow to group diseases with similar phenotypes to determine the genetics and molecular processes common to these diseases. The workflow uses a list of disease codes, in this case A-RDs ORPHA codes, to retrieve their associated HPO terms and genes from the Monarch Initiative. Then, the following steps are performed: (1) grouping diseases in clusters by phenotypic similarity; (2) calculation of the average shortest path (ASP) between known genes associated with diseases in each cluster using STRING data; (3) for clusters that have all gene pairs with computable paths in the interaction network, we do an expansion of the cluster with available genes in all shortest paths; (4) enrichment analysis in the Gene Ontology (GO) for both raw and expanded gene clusters; and (5) gene cluster prioritization in the interaction network using the CRank algorithm.

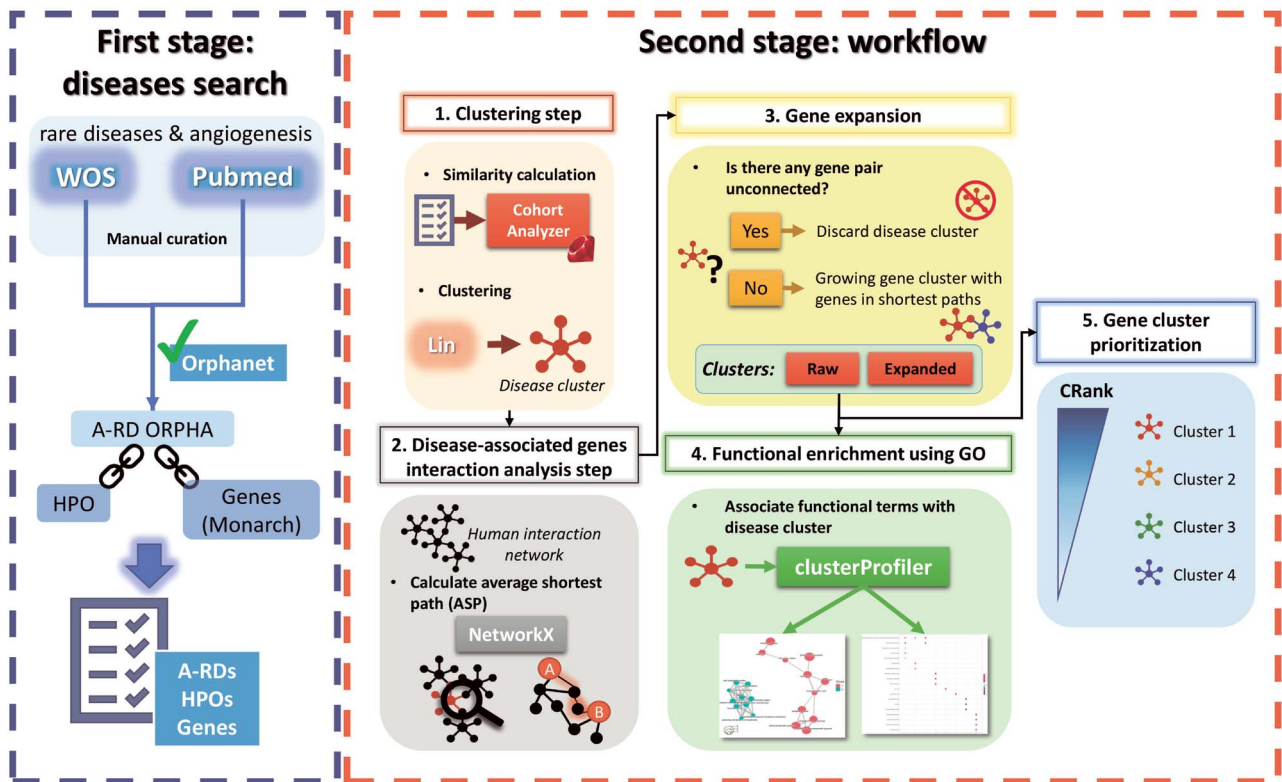


Figure 1. Representation of the main stages and steps followed in this work. The first stage consisted on the search for A-RDs through a systematic review of the literature. HPO terms associated with each A-RD were searched in the HPO annotation website and genes from the Monarch Initiative. The second stage describes the workflow developed in this work. The first step of this stage consists in retrieving the disease clusters. Then, in the second step, the ASP is calculated among genes for each disease cluster in the interaction network. For clusters with computable paths between its associated genes in the interaction network, they are expanded with the genes that are presented in the computed shortest paths. Finally, raw and expanded gene clusters are used in GO enrichment analysis and cluster prioritization in the interaction network.

The overview of the methodology is given in Figure 1. The workflow was developed in AutoFlow [45] and is available at https://github.com/ElenaRojano/angio_cluster.

Establishing A-RDs groups by phenotypes

We used the Cohort Analyzer tool [40] to group diseases using semantic similarity. This tool, included in the Patient Exploration Tools Suite [41], calculates different statistics in a cohort of patients or a pool of diseases and uses different semantic similarity methods to group them according to their HPO profiles [40].

The steps to obtain the disease clusters are performed as follows. First, the Lin similarity measure [31] is used to calculate a semantic similarity matrix among the A-RD HPO terms. This matrix is transformed to a dissimilarity matrix (the Lin measure ranges between 0 and 1, it is transformed with $1 - \text{similarity}$). This dissimilarity matrix is used to perform a hierarchical clustering with the R core function `hclust` using the Ward criterion [29]. The resulting dendrogram is split with the `cutreeDynamic` function included in the `dynamicTreeCut` R package [21] to get the final disease clusters. This algorithm is a hard clustering procedure that iterates the dendrogram analyzing the tips of the branches to identify possible tightly connected clusters. It is used with default settings

except for the `minClusterSize` and `deepSplit` parameters. The `minClusterSize` is the minimum items that can contain a cluster and we calculated it as the 1% of the diseases that have HPO terms. And the `deepSplit` parameter configures several internal parameters that controls how the branch partitioning and the clustering merging is performed. We set it to 2 following the `cutreeDynamic` authors recommendations.

Exploring the molecular mechanisms involved in the disease clusters

Once we obtained the disease clusters grouped by their phenotypic similarity, we analyzed their associated genes to explore the underlying molecular mechanisms. These clusters may have some variability in terms of the number of diseases they have. However, it is expected they share associated genes or at least having genes with similar functions. We eliminated clusters with a single disease as they were meaningless in this study.

Then, we selected the union of all disease-associated genes for each A-RD cluster and analyzed how they were connected in the protein interaction network. We downloaded all human interactions from the STRING database [51] (version 11.0b) and selected those with a combined score higher or equal to 900, which indicates a high confidence interaction between proteins [54].

In an additional step, we computed the degree for each node in the network. All the degree values were converted into Z-scores subtracting the mean degree and dividing by the standard deviation each node degree. We removed nodes with a Z-score greater or equal than 2.5 as they were considered as hubs excepting those that are listed in GO term angiogenesis or in the WikiPathway angiogenesis. In addition, we traced the paths between genes associated with the A-RD clusters by adjusting the methodology described for the Human Gene Connectome [15]. This approach is useful to detect paths between genes in clusters with a small number of genes. For this, we downloaded and modified the script ‘Gene-specific_connectome.py’ available at <https://lab.rockefeller.edu/casanova/HGC>. Then, once all the paths between genes have been determined for the genes presented in the clusters, we used the gene paths of the Human Gene Connectome to calculate the ASP to measure the closeness of the disease associated genes in each cluster. Disease clusters with a pair of associated genes without path in the Human Gene Connectome were removed from this study.

A gene expansion of the genes associated with the disease clusters was performed using the paths of the Human Gene Connectome. For this, we associated with each cluster all the genes that presented in the paths that connect each pair of disease cluster associated genes. This expansion allows us to find genes very close to disease-associated genes but that have not been initially described for the diseases of each cluster, and that potentially could be considered as possible genes involved in the disease development. Additionally, the CRank [61] algorithm was applied to the gene lists (raw or expanded) to rank the A-RD clusters in accordance with their network connectivity features at gene level. Furthermore, it allows us to evaluate the improvement of the expanded gene lists. This algorithm measures the magnitude of structural features and the robustness against noise for the clusters in the network using four different connectivity metrics: Likelihood, Density, Boundary and Allegiance. All these metrics are summarized in the CRank value, which ranges between 0 (the evaluated list is the most dispersed cluster from the clusters set in the network) and 1 (the evaluated list is the most connected and coherent cluster in the clusters set).

We finally performed the functional enrichment analysis for the genes associated with the disease clusters in their expanded form or not, using the clusterProfiler R package [60]. This enrichment analysis was performed in molecular function and biological process GO sub-ontologies. The P-value associated with each functional category was calculated using the Over Representation Analysis (ORA) algorithm and corrected by multiple testing with the Benjamini–Hochberg method. Functional categories with adjusted P-value equal or less than 0.01 are reported. As the functional categories belong to the GO, when a functional category and its parent are significant for the same clusters, the parental terms are

removed to simplify the interpretation of the results. The visualization of the enrichment results is generated with the dotplot function of the clusterProfiler R package.

When the enrichment analysis displays a large amount of functional terms, we use a summary representation. For it, the terms for each disease cluster are sorted by their adjusted P-value and the top N categories (custom threshold) with the lowest P-value are selected for each cluster. Then, using specific functions of clusterProfiler, we calculate a Wang semantic similarity matrix between the selected functional terms [56]. This matrix is hierarchically clustered using the hclust R function with the average method, and the resulting dendrogram is split with the treecut function setting h to $1 - S$, where S is a custom similarity threshold used to get the GO clusters. For each similarity cluster, the common ancestor in all GO terms is searched and used as a representative term of the cluster. All child terms are replaced by this representative term and it gets the A-RD cluster relations available in their children. A parental cleaning process is applied as previously described. The results of this analysis are plotted with the heatmaply R package [8], building a heat map that groups rows and columns by their similarity vector. We show only a dendrogram for columns, corresponding to the disease clusters.

Results and Discussion

Retrieving A-RDs

For this work, we performed a bibliographic search on A-RDs in PubMed and WOS databases. We performed an automated scoring of the found articles depending on where the search terms for these diseases appeared and to select those with the largest amount of information regarding A-RDs. From an original list of 1107 articles, 242 were related to A-RDs. We selected and inspected them manually, resulting in 158 A-RDs that were extracted from the Orphanet database. Of these diseases, 107 were characterized with HPO terms and 109 have associated genes in the Monarch Initiative database (Supplementary Table 4).

Characterization and clustering of A-RDs

We calculated with the Cohort Analyzer some statistics of our A-RD list. The full report is available in the GitHub repository at https://github.com/ElenaRojano/angio_cluster. The A-RD list includes a large number of different pathological phenotypes: 1476. Likewise, the average number of HPOs used to describe each disease is high: 28.36. This detailed description of the diseases will allow us to cluster the diseases in a more precise and informative way.

Cohort Analyzer computes the frequency for each phenotype in the disease list. In Table 1, we show the top 10 most frequent HPOs. We also check in the current bibliography its relationship with angiogenesis. For example, the terms HP: ‘Seizure’, HP: ‘Headache’

Table 1. Top 10 most frequent HPOs in the A-RD cohort

HPO	%
Seizure	28.03
Fatigue	25.23
Abdominal pain	20.56
Hepatomegaly	19.62
Splenomegaly	18.69
Weight loss	18.69
Thrombocytopenia	16.82
Fever	14.95
Headache	14.95
Hypertension	14.95

and HP: ‘Hypertension’ are quite related to endothelial dysfunctions in patients with preeclampsia and hypertensive encephalopathy, and have been associated with dysregulations of vascular endothelial growth factor (VEGF) in endothelial cells [20, 22, 26]. In fact, there are studies that relates VEGF-induced angiogenesis and the phenotypes HP: ‘Hepatomegaly’, HP: ‘Splenomegaly’ and HP: ‘Thrombocytopenia’ [57, 58]. Dysregulations affecting VEGF levels lead to blood vessel anomalies and consequently produce all these symptoms. The HP: ‘Fatigue’ term is mostly related to patients with cancer [13], and it also has been reported along with HP: ‘Weight loss’ and HP: ‘Fever’ in a patient with hemophagocytic lymphohistiocytosis, a rare immune disease [23]. Other phenotypes observed in patients with this rare syndrome include the top terms HP: ‘Splenomegaly’, HP: ‘Hepatomegaly’ and HP: ‘Abdominal pain’ [7]. Taken altogether, this information shows that top 10 most frequent HPOs are related to alterations of the angiogenic process.

Cohort Analyzer tool computes the semantic similarity of the HPO profiles associated with each A-RD by the Lin method. Then, these A-RDs are clustered using these similarity values. In this case, the tool generated 18 different clusters (Supplementary Table 4, ‘ClusterID’ column) with an average of diseases per cluster of 5.88. It is worth mentioning that from the initial list of 158 A-RDs, 47 do not have HPOs described. It also draws attention that most of them are different types of cancer, including retinoblastoma, various sarcomas such as liposarcoma and rhabdomyosarcoma and carcinomas including pancreatic and renal cell carcinomas, among others. However, the HPO has not yet included all the pathological phenotypes used to describe the different types of cancer available in Orphanet. It should be mentioned that the medical focus of the HPO in its early years was the phenotypic characterization of Mendelian diseases [12], and many types of cancer are produced by somatic mutations in individual cells that do not follow a pattern of inheritance [34]. This would explain why for many types of A-RDs there are no annotations found for this ontology and suggest that they are enriched in oncologic diseases.

The angiogenesis map of genes and diseases

In Supplementary Table 4, we show in which clusters the A-RDs have been grouped and the genes they have. It is worth mentioning that a gene can be associated with several or to few diseases of a cluster. In fact, in most cases, the genes are connected to only a disease of the cluster. As can be seen, from the 107 diseases with pathological phenotypes available for this study, 23 have no genes described. Diseases with HPO description were used to perform the clustering. This does not mean that this information is not valuable, but quite the opposite: it is possible to determine whether diseases within the same cluster participate in the same biological processes to extrapolate the information to diseases whose genetics are still unknown.

In Supplementary Figure 1 can be observed the robustness of the cluster procedure and the semantic similarity selection. Cohort Analyzer can use three measures: Resnik, Lin and Jiang-Conrath. We executed all of them and performed the correlation analysis shown in Supplementary Figure 1A. As can be seen, Lin and Resnik are very similar, in contrast with Jiang-Conrath. We performed the clustering using the Jiang-Conrath similarity and we obtained the results shown in Supplementary Figure 1B. The similarity matrix has very homogenous values, thus the clustering can only identify five clusters (colored segments in vertical bar), whereas the Lin similarity identifies 18 clusters. As the correlation analysis shows that Lin and Resnik are equivalent, we selected Lin similarity because it ranged between 0 and 1. Regarding the clustering robustness of the Lin matrix similarity, we have sampled the disease list 100 times selecting the 99%, 98%, 95% and 90% of the elements and performed the whole clustering procedure with each sample. Each disease partitioning was compared with the full disease set partitioning using the adjusted mutual information [53]. Values near 0 means that partitions are very different, whereas values near 1 means that both partitions are very similar. Supplementary Figure 1C shows the distribution for the 100 samples of each sampling. As it can be seen, most of the samples are accumulated between 0.85 and 0.87 of adjusted mutual information for selection of 99% but the other selections are decreased notably. These are good values and we must take into account that to calculate the adjusted mutual information, the members of the two comparisons must be the same. To overcome this problem, we have added a cluster with the removed elements to sample clusterings and this step decreases the adjusted mutual information.

We generated a network representation with Cytoscape [46] to create the angiogenesis map of genes and diseases (Figure 2), related by the computed A-RDs clusters. The frequency of occurrence for each gene can be observed in Supplementary Figure 2A. As can be seen, most genes are connected to a single disease. We observe in Figure 2 that most of the clusters (green circles) are connected between them by at least one gene (lilac circles). The connected clusters present at least one disease (salmon

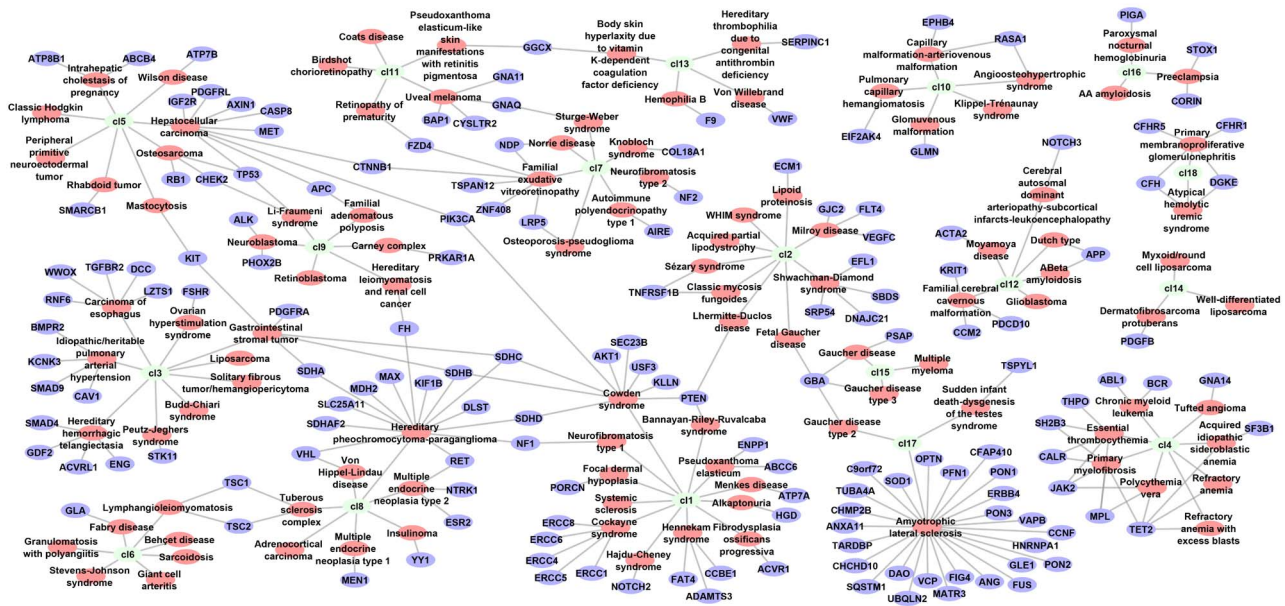


Figure 2. Angiogenesis map of genes and diseases representation. The salmon circles represent A-RDs and green circles in which cluster they belong. The lilac circles are the genes associated with each disease.

circles) that has been described with the same gene. Some isolated clusters are also observed, such as clusters 4, 10, 12, 14, 16 and 18. For example, cluster 4 is a tight gene-disease cluster with genes such as *TET2*, which plays a key role in erythropoiesis and its mutations are associated with anemia [10]. This gene has been described in five diseases, three of them are different types of anemia, and the other two are polycythemia vera (a blood cancer characterized by excessive production of red blood cells) and primary myelofibrosis (a bone marrow disease that affects the correct production of blood cells). As can be seen, these A-RDs share not only characteristics at the phenotypic level but also at the genetic level.

It is also remarkable that cluster 1 has 12 different diseases, each one described with at least one gene except the systemic sclerosis syndrome. This cluster overlaps with clusters 2, 3, 5 and 8 due mainly to Cowden syndrome whose genes connect with A-RDs characterized by the development of tissue tumors (malign or benign), a characteristic phenotypic type of Cowden syndrome [11].

In the case of overlapping clusters, cluster 5 is an interesting example of both similar phenotypic features and genetic basis. Intrahepatic cholestasis of pregnancy disease has two associated genes: the ATP binding cassette subfamily B member 4 (*ABCB4*) and the ATPase Phospholipid Transporting 8B1 (*ATP8B1*). The latter gene belongs to the same family as the gene associated with Wilson disease, the *ATP7B* gene [42]. This suggests a very similar genetic basis for both diseases, supported by a high phenotypic similarity. For this reason, this approach could be used to identify some of these genes as involved in diseases that have not genes associated yet.

In the case of clusters 1, 3 and 8 it is interesting that they are connected by the genes *SDHB* and *SDHC*

that are shared by three diseases: Cowden syndrome, gastrointestinal stromal tumor and hereditary pheochromocytoma–paraganglioma. Furthermore, from the same gene family, the genes *SDHD* and *SDHA* are shared for some pairs of these diseases, all of them characterized for generating benign overgrowths in different tissues and following an inheritance pattern [30].

In addition, we can find some diseases with a large number of genes, such as amyotrophic lateral sclerosis (ALS) in cluster 17. In fact, these genes are only connected to ALS. It is known that ALS is produced by mutations in a single or several genes at the same time [27] and this explains the large number of associated genes. Among them, we found angiogenin (*ANG*), a gene that stimulates angiogenesis in healthy and tumor tissue.

Altogether, the results shown and discussed in this section clearly show that the angiogenesis map of genes and diseases is very useful to extract new relations between genes and diseases. For this reason, we perform further analysis at gene interaction and functional levels.

Mapping the disease clusters onto the human interactions network

Once we have the A-RD clusters and the genes associated with the diseases, we can explore how these genes are related between them.

To measure the proximity of the genes for each cluster, we mapped them to a high-quality STRING human interaction network that includes interactions with a combined score higher or equal to 900 and removes hub nodes as described in Material and Methods. The distribution of degrees of this network is shown in the [Supplementary Figure 2B](#). This proximity measure is performed through the ASP calculation, which gives the number

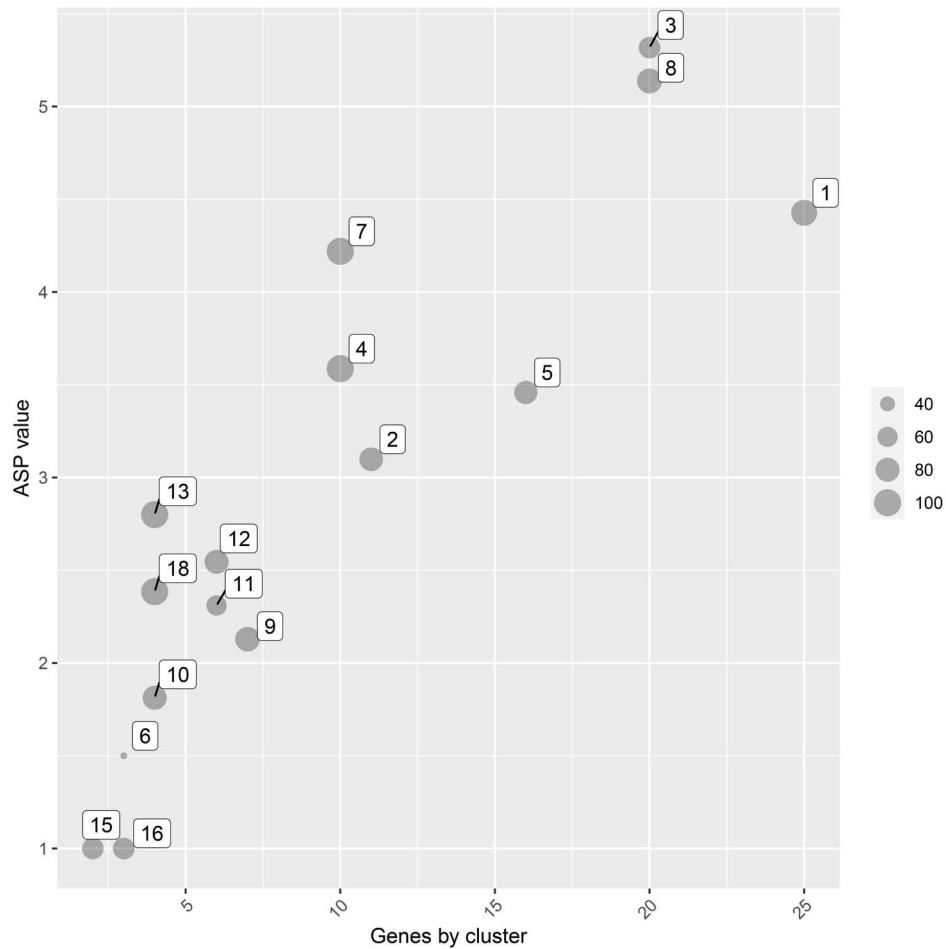


Figure 3. Scatter plot representing the ASP calculated among genes within each disease cluster. X-axis represents the number of genes by cluster and the Y-axis the ASP value. Dot size represents the number of diseases by cluster. Dots numbers are the identifier for each disease cluster.

of nodes between two genes. If only a single disease-associated gene appeared in the interaction network, as it happens with cluster 14, the cluster is discarded. Genes that were not found in the interactions network are available in the [Supplementary Table 3](#).

For each cluster gene list, we calculated the ASP values using the interaction network considering the interaction weight itself ([Supplementary Table 2](#), column ASP_value). We removed cluster 14 and 17 because no direct path could be established between all their associated genes. Thus, we considered these clusters as unconnected to the gene level. [Figure 3](#) shows for the remaining clusters the ASP, the number of disease-associated genes in the cluster and the proportion of the diseases in the cluster that have at least one associated gene by Monarch. First, in most of the clusters it can be observed that when the number of genes per cluster increases, the ASP values also increase. This makes sense because it is more difficult to find direct paths between multiple different genes than in small groups of genes. This trend can be observed in [Figure 3](#) for most of the clusters, except for cluster 1 which, despite having 25 genes, does not reach the ASP value of 5 that clusters 3 and 8 have with fewer genes (20). The same can be

observed for clusters 5 and 2, which have a higher number of genes per cluster (16 and 11 respectively) but a lower ASP than clusters 4 and 7 (10 genes). This indicates that the genes in these A-RD clusters are close in the interactions network and may point to similar pathways, being interesting to study in more detail. Besides, the distribution of clusters would indicate that they include different diseases at the phenotypic level but that the affected genes are close in the interaction network.

This could also suggest that the A-RDs within a cluster may have alterations in the same biological processes in which different genes are involved. It may be discussed that the gene closeness is due to the association of genes with a specific A-RD; however, in [Figure 3](#) it is shown that 60% to 80% of the A-RDs in a cluster have associated at least one gene. If we focus on specific A-RD clusters, three different groups can be observed. The first is composed of three clusters: to the right of the figure we observe clusters 3, 8 and 1 with the highest number of genes per cluster that ranges from 20 to 25, and their ASP values range from 4.42 to 5.31. In the middle of the graph, we can observe a second group of four clusters (7, 4, 2 and 5) with between 10 and 16 genes per cluster and ASP values between 3.09 and 4.21. These two groups of clusters

would point to coherent clusters at both phenotype and interaction levels. Finally, the rest of clusters have both the lowest number of associated genes (less than eight) and the lowest ASP values, from 1 to 2.80. Consequently, all these clusters have associated genes very close in the interaction network.

This evidence supports that the disease clusters are coherent at both phenotypic and interaction levels. The genes associated with each cluster are very close in the interactome and this suggests that they are involved in the same biological mechanism.

Functional analysis of the A-RD clusters

In the previous sections, we determined the disease similarity and the gene closeness in the human interaction network for each disease cluster. In this section, we focus on the functional perspective of these A-RD clusters. In this way, we used the gene lists with a computable average ASP value and performed a Gene Ontology (GO) enrichment analysis.

In [Figure 4](#), we found significant categories in GO molecular function for 12 clusters. It is shown a high specialization in the functions for each disease cluster. Clusters 1 and 3 have the largest gene lists (25 and 20, respectively) and several significant functional categories. Cluster 5, with 16 genes, is the one with the highest number of significant functional categories (15), including different kinase activities and transcription factor binding, among others.

The highest functional overlap is observed in clusters 1, 3 and 8, with the two categories, ‘oxidoreductase activity, acting on the CH-CH group of donors’ and ‘quinone binding’, both directly related to the electron transport chain in mitochondria. Looking closely at the genes for each cluster, we verified that all three clusters have genes that code for different subunits of the succinate dehydrogenase complex (*SDH* gene, [Figure 2](#)). Regarding the specific functions of these clusters, we found cluster 15 with functional terms that are known to be associated with angiogenesis, including ‘glycosphingolipid binding’ [17] and ‘glucosidase activity’ [32]. In the same way, cluster 6 has the single category ‘Hsp90 protein binding’ and this complex is known to be involved in angiogenesis as well [16].

Regarding GO biological process, in [Figure 5](#), we found significant terms for 14 clusters. We also observed the specialization in the functions for each disease cluster mentioned for annotations in GO molecular function. Clusters 1, 5 and 3 have the largest number of functions, likely due to their substantial and diverse gene lists. In fact, cluster 1 presents a high variability of functional annotations but also includes specific processes related to the VEGF signaling pathway, like ‘endothelial cell migration’ and ‘response to cadmium ions’ [18]. It is worth to mention that related to angiogenesis, when detailed results of biological process are observed, [Supplementary Figure 3](#), cluster 3 has the terms ‘negative

regulation of epithelial cell proliferation’ and ‘endothelial cell differentiation’.

This functional analysis supports the relationship between the selected diseases and the angiogenesis mechanism in which they relay. Furthermore, it allows to inspect the functional specialization for each disease cluster and which functions are shared by different groups of diseases reflecting the interconnection of the different angiogenesis related mechanisms.

A-RD clusters gene expansion to find unknown disease associated genes

We explored the phenotypic, interaction and functional levels of the clustered A-RDs and the evidence shown in this work avails the relationships between the A-RDs, as well as those between them and their associated genes. To deepen the results, we can identify new putative candidates and members of molecular mechanisms. To do this, we used the ASP calculation to take the genes in these short paths and expand the gene list for each disease cluster ([Supplementary Table 2](#), ASP_expanded_genes column).

In [Supplementary Figure 4](#), we can see how the number of genes associated with clusters 1, 3 and 8 range from 160 to > 250. This is clearly due to the number of associated genes in the Monarch Initiative to the cluster diseases, from 20 to 25 genes ([Figure 3](#)). With this number of genes, these clusters likely will be uninformative.

In addition, the expanded gene lists were explored to identify new functions and connections between the disease clusters, repeating the functional analysis.

When we explored the summary results for GO molecular function ([Figure 6](#)), we observed that a larger number of clusters (15) have significant functional categories than without gene expansion (12). Additionally, the gene expansion increased the number of significant functions and the functional overlap between disease clusters, specially highlighted when full results are inspected ([Supplementary Figure 5](#)). We also observed a functional specialization for all disease clusters in [Figure 6](#).

There are several functional categories that are presented by many clusters, such as ‘ubiquitin protein ligase binding’, shared by seven clusters, ‘p53 binding’, shared by six clusters and ‘RNA polymerase II transcription factor binding’, shared by five clusters. Being such general functions, it is not rare to find them in different clusters, although somehow they are related to angiogenesis. For example, p53 tumor suppressor is known to have a regulatory effect on VEGF expression and consequently on angiogenesis [5]. Another more specific function such as ‘phosphoprotein binding’ is shared in clusters 10, 4, 5, 3 and 8. There are proteins with this function that have been related to angiogenesis stimulation in pathogenic processes [9]. Clusters 1, 2 and 8 are those with a greater number of associated functional categories, being mostly related to the mitochondrial transport chain in clusters 1 and 8, while cluster 2 has functions mostly related to protein synthesis and cellular communication. It is



Figure 4. Dot plot for results obtained with the clusterProfiler R package in GO molecular function. X-axis includes the A-RD cluster identifiers and the number of genes by cluster between brackets. Y-axis represents each GO molecular function term associated with the genes for these clusters. Color scale represents the adjusted P-value (red: lower, blue: higher) and dot size indicates the proportion of genes in the functional category that are annotated in the cluster.

remarkable that cluster 3 has functional annotations especially related to angiogenesis. For example, in the case of the ‘S100 protein binding’ function, it is known that protein S100A4 is associated with metastasis and promotes angiogenesis [44]. Furthermore, the ‘platelet derived growth factor receptor binding’ function in the same cluster is related to angiogenesis and cell proliferation in injured tissues [35]. It is also worth to

mention that clusters 5 and 9 both share the ‘sequence-specific double-stranded DNA binding’ function. It is known that the double-stranded RNA-binding protein DRBP76/NF90 regulates the stability of the VEGF mRNA stability in breast cancer, and its repression is associated with a reduction of the angiogenic and tumorigenic process in breast cancer cells [55]. Another interesting angiogenesis-related function is ‘SH2 domain binding’,

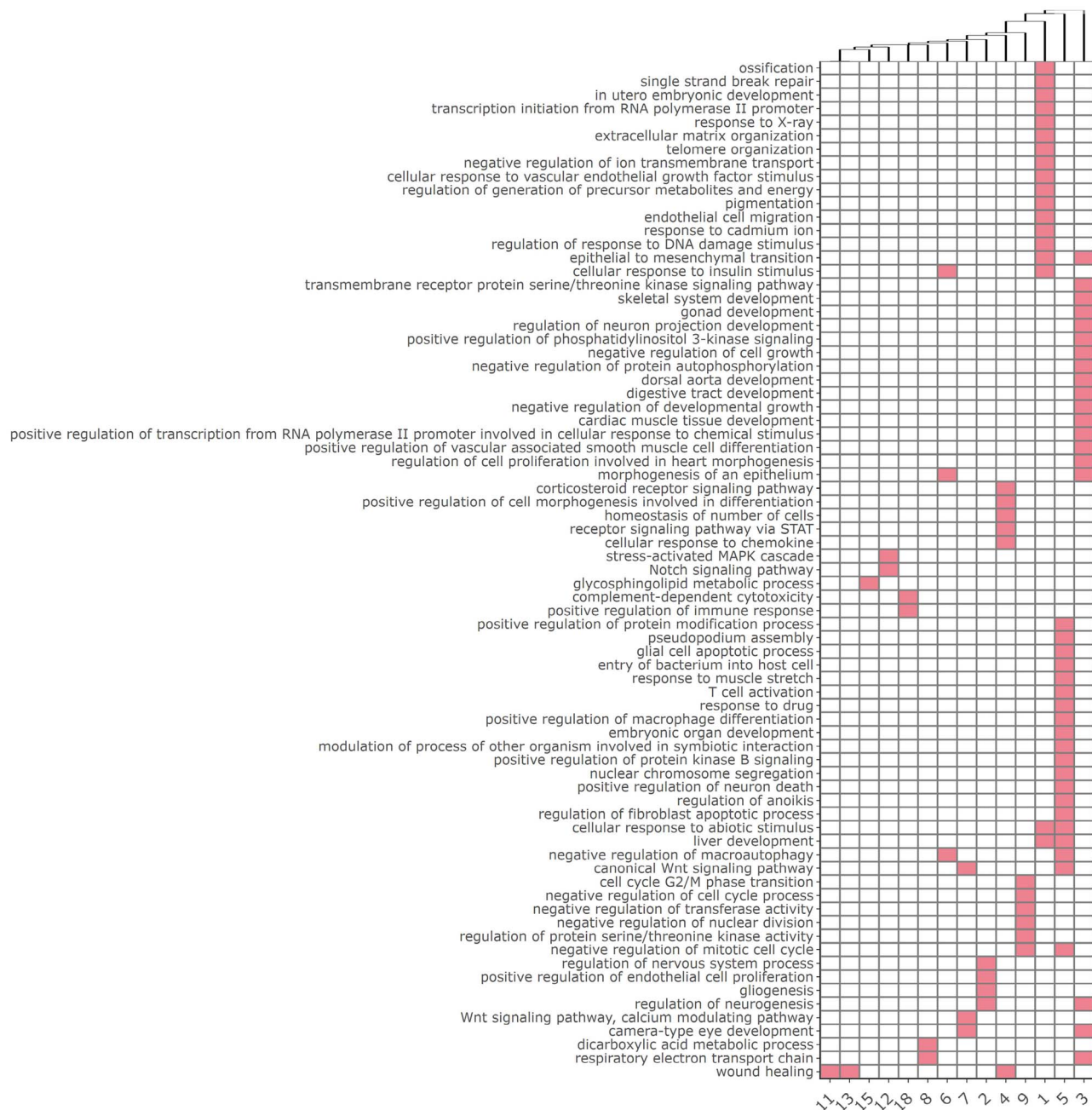


Figure 5. Heat map for results obtained with clusterProfiler R package in GO biological process. X-axis shows the A-RD cluster identifiers and Y-axis shows the summarized terms for the enrichment results as described in Material and methods, using 50 terms per cluster and a similarity threshold of 0.7.

shared by clusters 4 and 2. The growth factor receptor bound protein 2 (Grb2-SH2) domain binding is an antagonist of VEGF and blocks angiogenesis [49]. And finally, it is known that proteins with the ‘G protein-coupled receptor binding’ function, observed in clusters 11, 7 and 2, also have a regulatory effect in angiogenesis [37].

In the case of GO biological process (Supplementary Figure 6), due to the large number of functions associated with the genes for each cluster, we show only the summary results. As in the case of GO molecular function, here we can see again how the gene expansion finds functional terms for clusters that were not available with the original genes.

In any case, the gene expansion approach allows us to identify the participation of diseases or molecular mechanisms of previously not related genes, and to contribute to reveal the biological basis of these diseases.

A comparison of angiogenesis-related genes with ClinVar data and known angiogenesis gene sets

To illustrate the value of the relationships found between A-RDs, their associated genes and the inferred genes through interaction data, we compared our results with genomic known data. For this, we explored known pathogenic variants related to angiogenesis. We searched for the keyword ‘angiogen*’ in the ClinVar database and

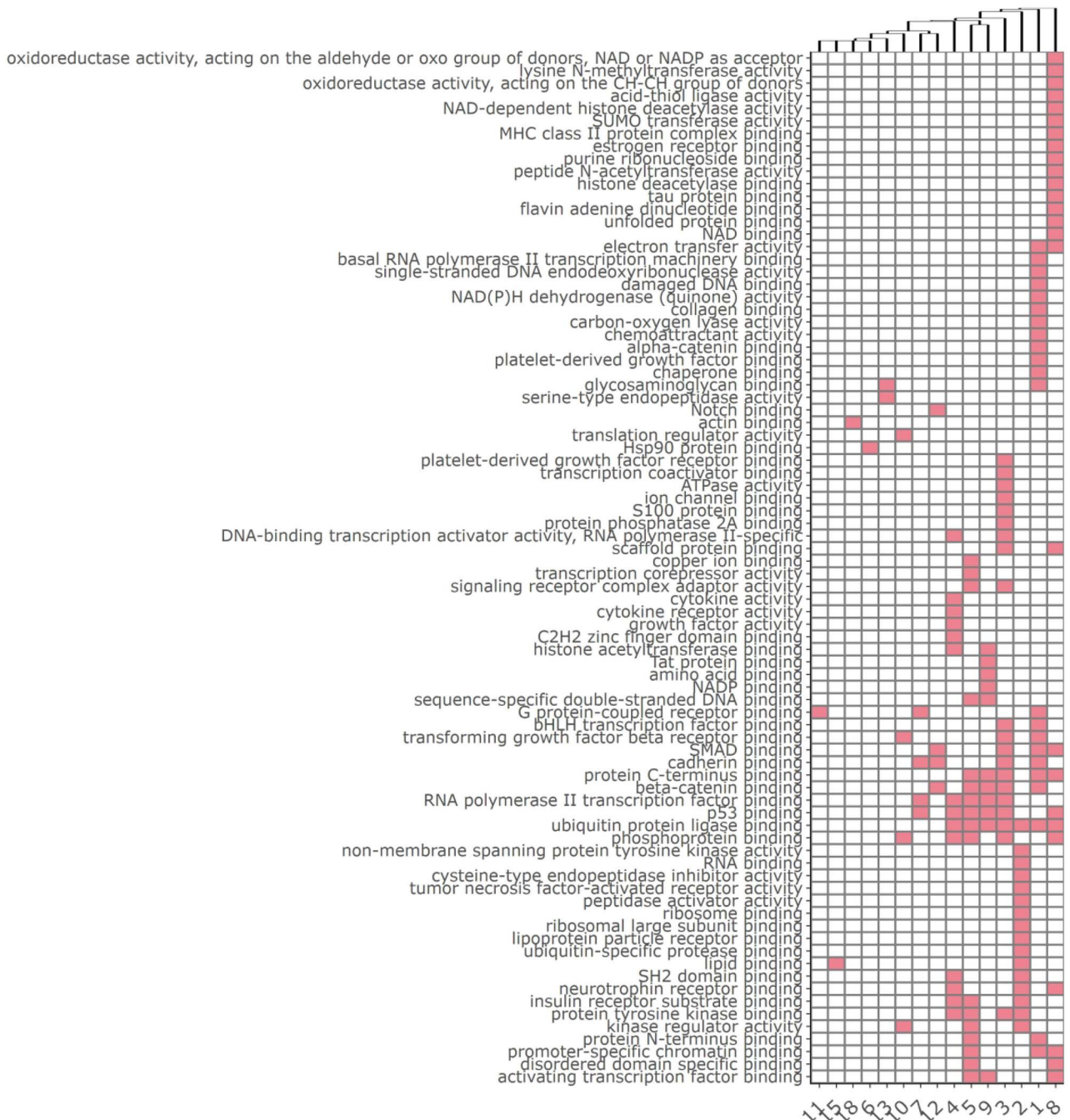


Figure 6. Heat map for results obtained with clusterProfiler R package in GO molecular function for expanded clusters. X-axis shows the A-RD cluster identifiers and Y-axis shows the summarized terms for the enrichment results as described in Material and methods, using 35 terms per cluster and a similarity threshold of 0.6.

selected variants with clinical significance defined as pathogenic and whose length was less or equal to 50 nt to ensure that they only affected a single gene. We got a list of 16 pathogenic variants included in the [Supplementary Table 4](#). From this pool of variants, we used the gene identifiers associated with ClinVar. As we are aware that ClinVar search engine could give inaccurate results, we checked if they were involved in angiogenesis by performing a bibliographic search ([Supplementary Table 5](#)). In [Table 2](#), we show the comparison between the angiogenesis-related genes from ClinVar and the gene lists obtained in this study. The genes *ANG*, *SDHA*, *SDHD*, *TP53* and *VEGFC* were found between the genes

associated with the A-RDs by the Monarch Initiative. The *F7* gene was found when the gene clusters were expanded with the interaction data. This *F7* gene is a coagulation factor related to angiogenesis, but the study [1] was overlooked by our bibliographic search. This work has all the features described in Material and methods to be included; however, it has no mention of the *VEGF* gene although it relates *F7* with angiogenesis. There are three genes that our study did not relate to angiogenesis: *MT-TE*, *RNASE4* and *AIMP*. In the first case, *MT-TE* is a mitochondrial gene that encodes a tRNA for glutamic acid. Mutations in this gene are known that produces myopathies [14] or diabetes mellitus [38], but

Table 2. List of genes affected by pathogenic variants associated with angiogenesis in ClinVar database. The Match column shows if the gene is identified in this study. *Associated genes* means that the gene is found in the list of genes retrieved from the Monarch Initiative, whereas *Expanded genes* means that the gene was found in the gene lists obtained with the STRING interaction data

Gene	Associated variants	Match
AIMP1	1	No
ANG	8	Associated genes
RNASE4	8	No
F7	1	Expanded genes
MT-TE	1	No
SDHA	1	Associated genes
SDHD	1	Associated genes
TP53	1	Associated genes
VEGFC	1	Associated genes

there are no studies relating this gene to angiogenesis at the disease level. This gene does not encode a protein, therefore the protein interaction data are useless. In the case of the RNASE4 gene, we found that its genomic coordinates were overlapping the ANG gene coordinates and its eight associated variants were also affecting the ANG gene. In fact, when the identifiers of these variants were inspected in [Supplementary Table 4](#) for these two genes, they referenced the ANG gene but not RNASE4. Consequently, we could consider that the pathogenic variants affect the ANG gene but not RNASE4, and ANG was identified as an angiogenesis-related gene. In the case of the AIMP gene, it was not identified at all due to the two following reasons: this gene causes the Pelizaeus–Merzbacher-like disease encoded as ORPHA:280293 but it does not have HPO terms described in the Monarch Initiative, and its related publication ([6]) does not include the ‘rare’ keyword. Consequently, our criteria ignores it although it lists five different OMIM entries. This highlights that our methodology can identify six of seven angiogenesis-related genes with known pathogenic variants.

In addition, we listed the genes associated with the GO term angiogenesis (GO:0001525) and the GSEA group wp_angiogenesis that is extracted from WikiPathways (pathway WP1539). We performed an enrichment analysis for each cluster only with these two categories. We selected results with adjusted P -value ≤ 0.05 (Table 3). Surprisingly, the angiogenesis GO category does not give any significant results but the wp_angiogenesis list is significant for 6 of 16 expanded clusters. This evidences that several clusters generated in this work are related to the angiogenesis pathway and the others that are not significant are related to angiogenesis-dependent processes, suggesting the genes presented in the expanded clusters are important for the angiogenesis process.

Prioritization of gene groups associated with A-RD clusters in the interaction network

Finally, we applied a prioritization approach to the A-RD clusters using their associated genes. In this way,

Table 3. List of clusters with adjusted P -value 0.05 or less for angiogenesis pathway in WikiPathway database

Cluster id	Adjusted P -value
1	5.86×10^{-4}
8	1.2×10^{-2}
5	9.13×10^{-4}
7	7.15×10^{-3}
9	8.86×10^{-3}
3	2.26×10^{-2}

we could rank the A-RD clusters in accordance with the network connectivity features of the associated genes in the interaction network, rewarding the clusters with high interconnectivity. This ranking allows us to select which clusters (raw or expanded) are promising candidates for downstream experiments. For this reason, we used the CRank algorithm [61] that uses Likelihood, Density, Boundary and Allegiance metrics to characterize network connectivity features for each cluster in an integrated way. These metrics measure the structural features and the robustness against noise of the clusters in the network structure. The integrated CRank value was computed for both raw and expanded clusters, as shown in Figure 7.

Regarding the A-RD clusters with the raw gene lists, clusters 18, 7 and 2 were in the top three with 1, 0.81 and 0.81 CRank values, respectively. When the gene lists were expanded with the ASP computation, six clusters highly increased their CRank measures, whereas seven clusters decreased their values to a lesser extent. Noteworthy, with the expanded gene list, clusters 18, 7 and 2 decreased their CRank, but three new clusters reach the top: clusters 4, 16 and 13 increased its Crank value to 0.84, 1 and 0.93, respectively. As can be seen, the top of prioritization changes if the gene clusters are expanded.

In this way, the gene association for the ASP calculation can be measured, giving the opportunity to choose which clusters need further investigation.

Concluding remarks

The approach presented in this work has allowed us to deepen our knowledge of A-RDs at the genetic, phenotypic and molecular levels. Starting from the work of Rodríguez-Caso and collaborators, we have characterized the phenotypes of A-RDs to be able to group them according to their semantic similarity. This allowed us to analyze disease clusters at the genetic level, exploring the molecular mechanisms involved in the development of different diseases. Likewise, we have demonstrated the coherence of the diseases within each cluster at the genetic level with the use of network characteristics such as the ASP. This approach allows the identification of clusters whose genes are very close in the network of interactions and that could be involved in related molecular mechanisms. Besides, we propose the CRank

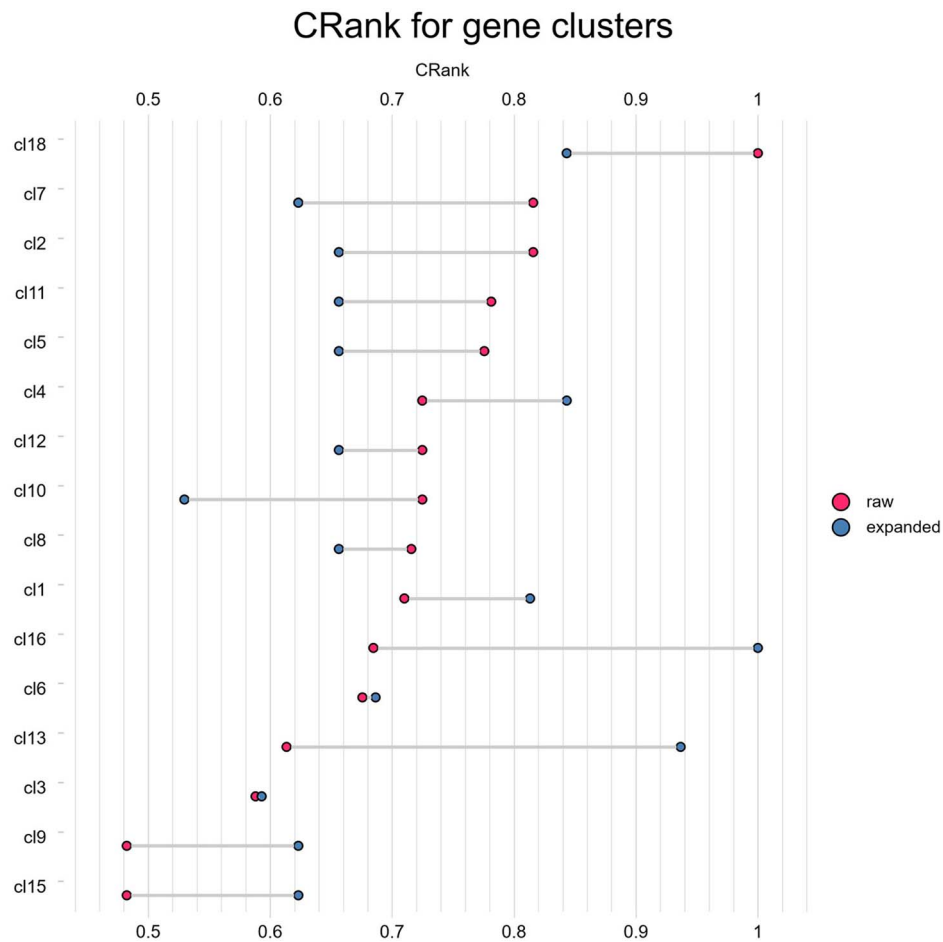


Figure 7. CRank values for the genes associated with each A-RD cluster, using the interaction network from STRING filtered by a combined score of 900. Series show the CRank for raw clusters (pink) and gene expanded clusters (blue).

measure to prioritize the A-RD clusters and to select candidates for downstream experiments in wet laboratories to validate the new gene associations.

In addition, this strategy has been essential to determine the common molecular mechanisms of these diseases. It also allowed us to explore the putative genes that could be associated with the A-RDs and whose function is not well characterized, considering them as possible genes involved in the disease development. To confirm the role of these genes in angiogenesis, an experimental validation is necessary. Furthermore, we have compared the results in this work with ClinVar angiogenesis-related data and we have achieved to list six of seven genes, one of them through the interaction data. Finally, our protocol can be extrapolated to the analysis of other diseases or biological processes.

Keys Points

- We propose a systemic methodology for the study of a set of rare diseases, grouping them according to their phenotypic similarity and analyzing them at a functional level using their disease-associated genes.

- This methodology is used to identify possible genes involved in angiogenesis-related rare diseases for those cases in which the genetic cause or functional impact is not known.
- We applied our methodology to the study of angiogenesis-related rare diseases, but it can be used to analyze other human genetic diseases.

Supplementary data

Supplementary data are available at *Briefings in Bioinformatics* online.

Author contributions statement

E.R. and P.S.Z. conceived the methodology. R.P.M., J.C.C., E.R. and P.S.Z., developed the software that implements the protocol. R.P.M., J.C.C., E.R., P.S.Z. and M.A.M. analyzed the results and provided interpretation. R.P.M., E.R. and P.S.Z. wrote the manuscript. B.M.P., J.A.G., A.R.Q. and M.A.M. were involved in planning of the study, contributed to the acquisition of funding for research

and headed the project. All authors read and approved the final version of the manuscript.

Acknowledgments

The authors thank the Supercomputing and Bioinnovation Center (SCBI) of the University of Malaga for their provision of computational resources and technical support (<http://www.scbi.uma.es/site>).

Funding

This work was supported by the Spanish Ministry of Science, Innovation and Universities (grant PID2019-105010RB-I00, grant PID2019-108096RB-C21), the Andalusian Government and FEDER (grants UMA18-FEDERJA-102, UMA18-FEDERJA-220, PY20_00257, PY20_00372, RH-0079-2021 and funds from the group PAIDI BIO 267); the Ramón Areces foundation, which funds project for the investigation of rare disease (National call for research on life and material sciences, XIX edition) and the University of Malaga (Ayudas del I Plan Propio). The 'CIBER de Enfermedades Raras' and 'CIBER de Enfermedades Cardiovasculares' are initiatives from the ISCIII (Spain). The funders had no role in the study design, data collection and analysis, decision to publish or preparation of the manuscript.

References

- Bernardi F, Mariani G. Biochemical, molecular and clinical aspects of coagulation factor VII and its role in hemostasis and thrombosis. *Haematologica* 2021;**106**(2):351–62.
- Buphamalai P, Kokotovic T, Nagy V, et al. Network analysis reveals rare disease signatures across multiple levels of biological organization. *Nat Commun* 2021;**12**(1):6306.
- Carmeliet P, Jain RK. Angiogenesis in cancer and other diseases. *Nature* 2000;**407**(6801):249–57.
- Cheng L, Li J, Ju P, et al. emFunSim: A New Method for Measuring Disease Similarity by Integrating Semantic and Gene Functional Association. *PLoS ONE* 2014;**9**(6).
- Farhang Ghahremani M, Goossens S, Nittner D, et al. p53 promotes VEGF expression and angiogenesis in the absence of an intact p21-Rb pathway. *Cell Death Differ* 2013;**20**(7):888.
- Feinstein M, Markus B, Noyman I, et al. Pelizaeus-merzbacher-like disease caused by AIMP1/p43 homozygous mutation. *Am J Hum Genet* 2010.
- Fisman DN. Hemophagocytic syndromes and infection. *Emerg Infect Dis* 2000;**6**(6):601.
- Galili T, O'Callaghan A, Sidi J, et al. Heatmaply: An R package for creating interactive cluster heatmaps for online publishing. *Bioinformatics* 2018;**34**(9).
- Gao Y, Yin Z, Qi Y, et al. Golgi phosphoprotein 3 promotes angiogenesis and sorafenib resistance in hepatocellular carcinoma via upregulating exosomal miR-494-3p. *Cancer Cell Int* 2022;**22**(1):1–17.
- Ge L, Zhang R-P, Wan F, et al. TET2 Plays an Essential Role in Erythropoiesis by Regulating Lineage-Specific Genes via DNA Oxidative Demethylation in a Zebrafish Model. *Mol Cell Biol* 2014;**34**(6):989.
- Gosein MA, Narinesingh D, Nixon CAAC, et al. Multi-organ benign and malignant tumors: Recognizing Cowden syndrome: A case report and review of the literature. 2016.
- Groza T, Köhler S, Moldenhauer D, et al. The Human Phenotype Ontology: Semantic Unification of Common and Rare Disease. *Am J Hum Genet* 2015;**97**(1):111.
- Himbert C, Ose J, Lin T, et al. Inflammation- and angiogenesis-related biomarkers are correlated with cancer-related fatigue in colorectal cancer patients: Results from the ColoCare Study. *Eur J Cancer Care* 2019;**28**(4).
- Horvath R, Kemp JP, Tuppen HA, et al. Molecular basis of infantile reversible cytochrome c oxidase deficiency myopathy. *Brain* 2009.
- Itan Y, Zhang SY, Vogt G, et al. The human gene connectome as a map of short cuts for morbid allele discovery. *Proc Natl Acad Sci U S A* 2013.
- Iwabayashi M, Taniyama Y, Sanada F, et al. Role of serotonin in angiogenesis: Induction of angiogenesis by sarpogrelate via endothelial 5-HT1B/Akt/eNOS pathway in diabetic mice. *Atherosclerosis* 2012;**220**(2):337–42.
- Jernigan PL, Makley AT, Hoehn RS, et al. The role of sphingolipids in endothelial barrier function. *Biol Chem* 2015;**396**(6–7):681.
- Kim J, Lim W, Ko Y, et al. The effects of cadmium on VEGF-mediated angiogenesis in HUVECs. *J Appl Toxicol* 2012;**32**(5):342–9.
- Köhler S, Vasilevsky NA, Engelstad M, et al. The Human Phenotype Ontology in 2017. *Nucleic Acids Res* 2017;**45**(D1):D865–76.
- Lamy C, Mas JL, Encephalopathy H. *Stroke* 2011;734–40.
- Langfelder P, Zhang B, Horvath S. Defining clusters from a hierarchical cluster tree: The Dynamic Tree Cut package for R. *Bioinformatics* 2008;**24**(5):719–20.
- Lenke L, Martínez de la Escalera G, Clapp C, et al. A Dysregulation of the Prolactin/Vasoinhibin Axis Appears to Contribute to Preeclampsia. *Front Endocrinol* 2020;**10**:893.
- Levy L, Nasereddin A, Rav-Acha M, et al. Prolonged Fever, Hepatosplenomegaly, and Pancytopenia in a 46-Year-Old Woman. *PLoS Med* 2009;**6**(4).
- Masino AJ, Dechene ET, Dulik MC, et al. Clinical phenotype-based gene prioritization: An initial study using semantic similarity and the human phenotype ontology. *BMC Bioinformatics* 2014;**15**(1):1–11.
- Mathur S, Dinakarandian D. Finding disease similarity based on implicit semantic similarity. *J Biomed Inform* 2012;**45**(2):363–71.
- Maynard SE, Karumanchi SA. Angiogenic Factors and Preeclampsia. *Semin Nephrol* 2011;**31**(1):33.
- Mejzini R, Flynn LL, Pitout IL, et al. ALS Genetics, Mechanisms, and Therapeutics: Where Are We Now? *Front Neurosci* 2019;**13**(1310).
- Mungall CJ, McMurry JA, Kohler S, et al. The Monarch Initiative: An integrative data and analytic platform connecting phenotypes to genotypes across species. *Nucleic Acids Res* 2017;**45**(D1):D712–22.
- Murtagh F, Legendre P. Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion? *Journal of Classification* 2014.
- Nazar E, Khatami F, Saffar H, et al. The Emerging Role of Succinate Dehydrogenase Genes (SDHx) in Tumorigenesis. *International Journal of Hematology-Oncology and Stem Cell Research* 2019;**13**(2):72.
- Pesquita C, Faria D, Bastos H, et al. Metrics for GO based protein semantic similarity: A systematic evaluation. *BMC Bioinformatics* 2008;**9**(SUPPL. 5):1–16.

32. Pili R, Chang J, Partis RA, et al. The α -Glucosidase I Inhibitor Castanospermine Alters Endothelial Cell Glycosylation, Prevents Angiogenesis, and Inhibits Tumor Growth. *Cancer Res* 1995;**55**(13):2920–6.
33. Piñero J, Ramírez-Anguita JM, Saüch-Pitarch J, et al. The DisGeNET knowledge platform for disease genomics: 2019 update. *Nucleic Acids Res* 2020;**48**(D1):D845–55.
34. Poduri A, Evrony GD, Cai X, et al. Somatic Mutation, Genomic Variation, and Neurological Disease. *Science (New York, NY)* 2013;**341**(6141):1237758.
35. Raica M, Cimpean AM. Platelet-derived growth factor (PDGF)/PDGF receptors (PDGFR) axis as target for antitumor and antiangiogenic therapy. *Pharmaceuticals* 2010;**3**(3):572–99.
36. Rao A, Joseph T, Saipradeep VG, et al. Piori-T: A tool for rare disease gene prioritization using MEDLINE. *PLoS ONE* 2020.
37. D. E. Richard, V. Vouret-Craviari, and J. Pouysségur. Angiogenesis and G-protein-coupled receptors: signals that bridge the gap. *Oncogene*, **20**(13):1556–62, 2001.
38. Rigoli L, Prisco F, Caruso RA, et al. Association of the T14709C mutation of mitochondrial DNA with maternally inherited diabetes mellitus and/or deafness in an Italian family, 2001.
39. Rodríguez-Caso L, Reyes-Palomares A, Sánchez-Jiménez F, et al. What is known on angiogenesis-related rare diseases? A systematic review of literature. *J Cell Mol Med* 2012;**16**(12):2872–93.
40. Rojano E, Córdoba-Caballero J, Jabato FM, et al. Evaluating, Filtering and Clustering Genetic Disease Cohorts Based on Human Phenotype Ontology Data with Cohort Analyzer. *J Pers Med* 2021;**11**(8):730.
41. Rojano E, Seoane-Zonjic P, Jabato FM, et al. Comprehensive Analysis of Patients with Undiagnosed Genetic Diseases Using the Patient Exploration Tools Suite (PETS). In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, Vol. **12108**, 2020, 775–86.
42. Roy S, McCann CJ, Ralle M, et al. Analysis of Wilson disease mutations revealed that interactions between different ATP7B mutants modify their properties. *Scientific Reports* 2020 10:1 2020;**10**(1):1–15.
43. Schlicker A, Lengauer T, Albrecht M. Improving disease gene prioritization using the semantic similarity of Gene Ontology terms. *Bioinformatics* 2010;**26**(18):i561–7.
44. Semov A, Moreno MJ, Onichtchenko A, et al. Metastasis-associated protein S100A4 induces angiogenesis through interaction with Annexin II and accelerated plasmin formation. *J Biol Chem* 2005;**280**(21):20833–41.
45. Seoane P, Ocaña S, Carmona R, et al. AutoFlow, a Versatile Workflow Engine Illustrated by Assembling an Optimised de novo Transcriptome for a Non-Model Species, such as Faba Bean (*Vicia faba*). *Current Bioinformatics* 2016;**11**(4):440–50.
46. Shannon P, Markiel A, Ozier O, et al. Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction Networks. *Genome Res* 2003;**13**(11):2498.
47. Slater LT, Karwath A, Williams JA, et al. Towards similarity-based differential diagnostics for common diseases. *Comput Biol Med* 2021;**133**(104360).
48. Slater LT, Williams JA, Karwath A, et al. Multi-faceted semantic clustering with text-derived phenotypes. *Comput Biol Med* 2021;**138**.
49. Soriano J, Liu N, Gao Y, et al. Inhibition of angiogenesis by growth factor receptor bound protein 2- Src homology 2 domain bound antagonists. *Mol Cancer Ther* 2004;**3**(10):1289–99.
50. Szekanecz Z, Koch AE. Mechanisms of Disease: angiogenesis in inflammatory diseases. *Nature Clinical Practice Rheumatology* 2007 3:11 2007;**3**(11):635–43.
51. Szklarczyk D, Gable AL, Lyon D, et al. STRING v11: protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res* 2019;**47**(D1):D607–13.
52. Trifonova EA, Swarovskaya MG, Ganzha OA, et al. The interaction effect of angiogenesis and endothelial dysfunction-related gene variants increases the susceptibility of recurrent pregnancy loss. *J Assist Reprod Genet* 2019;**36**(4):717–26.
53. Vinh NX, Epps J, Bailey J. Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance. *J Mach Learn Res* 2010.
54. von Mering C, Jensen LJ, Snel B, et al. STRING: known and predicted protein-protein associations, integrated and transferred across organisms. *Nucleic Acids Res* 2005;**33**(Database Issue):D433.
55. Vumbaca F, Phoenix KN, Rodriguez-Pinto D, et al. Double-Stranded RNA-Binding Protein Regulates Vascular Endothelial Growth Factor mRNA Stability, Translation, and Breast Cancer Angiogenesis. *Mol Cell Biol* 2008;**28**(2):772.
56. Wang JZ, Du Z, Payattakool R, et al. A new method to measure the semantic similarity of GO terms. *Bioinformatics* 2007;**23**(10):1274–81.
57. Xu Y, Xiao Y-Y, Simon M, et al. Plasma Vascular Endothelial Growth Factor (VEGF) Levels Correlate with Thrombocytopenia of Various Etiology. *Blood* 2014;**124**(21):4991–1.
58. Xue Y, Chen F, Zhang D, et al. Tumor-derived VEGF modulates hematopoiesis. *Journal of Angiogenesis Research* 2009;**1**(1):9.
59. Yang J, Dong C, Duan H, et al. RDmap: a map for exploring rare diseases. *Orphanet J Rare Dis* 2021;**16**(101).
60. Yu G, Wang LG, Han Y, et al. ClusterProfiler: An R package for comparing biological themes among gene clusters. *OMICS A Journal of Integrative Biology* 2012;**16**(5):284–7.
61. Zitnik M, Sosič R, Leskovec J. Prioritizing network communities. *Nature. Communications* 2018.