

Supplementary Information for: Between viral targets and differentially expressed genes in COVID-19: the sweet spot for therapeutic intervention

Carme Zambrana¹, Alexandros Xenos¹, René Böttcher¹, Noël Malod-Dognin^{1,2}, and Nataša Pržulj^{1,2,3,*}

¹Barcelona Supercomputing Center, Barcelona, Spain.

²Department of Computer Science, University College London, London, WC1E 6BT, United Kingdom.

³ICREA, Pg. Lluís Companys 23, Barcelona, Spain.

*natasha@bsc.es

Supplementary Results

The structure of PPI is preserved when merging it with GI and MI.

To validate that after merging the three networks the structure of the PPI network is preserved, we analyze the MIN and the constituent networks, PPI, GI and MI, by the following commonly used network properties: four centrality measures (degree, eigenvector, betweenness and closeness centrality) and clustering coefficient (for more details see section “Analysis of the molecular interaction network and its wiring patterns” in Methods). As shown in Supplementary Table S1, MI has highest values in all the network properties, especially for the clustering coefficient, which is expected since the MI connects all the genes that participate in the same metabolic pathway. The MIN has similar values to those for PPI network, which is expected since the PPI network is the biggest one compared to the other constituent networks, GI and MI (see Supplementary Figure S1A-B). Thus, in terms of centrality and clustering, the PPI network structure is preserved when merging it with the GI and MI networks.

To assess whether the wiring patterns of the PPI are preserved in the MIN, we use Graphlet Degree Vectors (GDVs) to compare the MIN and the three constituent networks, PPI, GI and MI (for more details see section “Analysis of the molecular interaction network and its wiring patterns” in Methods). As shown in Supplementary Figure S1C, GDV of the MI network is very different from the GDV of the rest of the networks, especially in the clique-orbits. Namely, orbits 3, 9, 10, 12 and 14. This is expected because of how the MI network is constructed, (i.e., by connecting all the genes that participate in the same metabolic pathway). The GDV of the MIN is very similar to that of the PPI, showing that the wiring patterns on the PPI are preserved after the merging.

To specify the relation between the genes that the data fusion process must conserve, we construct the MIN by merging three different interaction networks: protein-protein interaction (PPI), genetic interaction (GI) and metabolic interaction (MI) networks. By doing so, we want to add to the PPI different types of relations between genes that are key for the SARS-CoV-2 infection, (e.g., we add the MI network since it has been demonstrated that metabolic processes, such as glycolysis, promote SARS-CoV-2 replication) without losing its structure. When comparing the network properties and GDV of the constituent networks and the MIN, we showed that the PPI structure is preserved in the MIN although the MI has a very different structure due to how it is constructed. Therefore, we obtain a holistic view of the relationship between genes without losing the PPI structure, which is the most curated of the constituent networks.

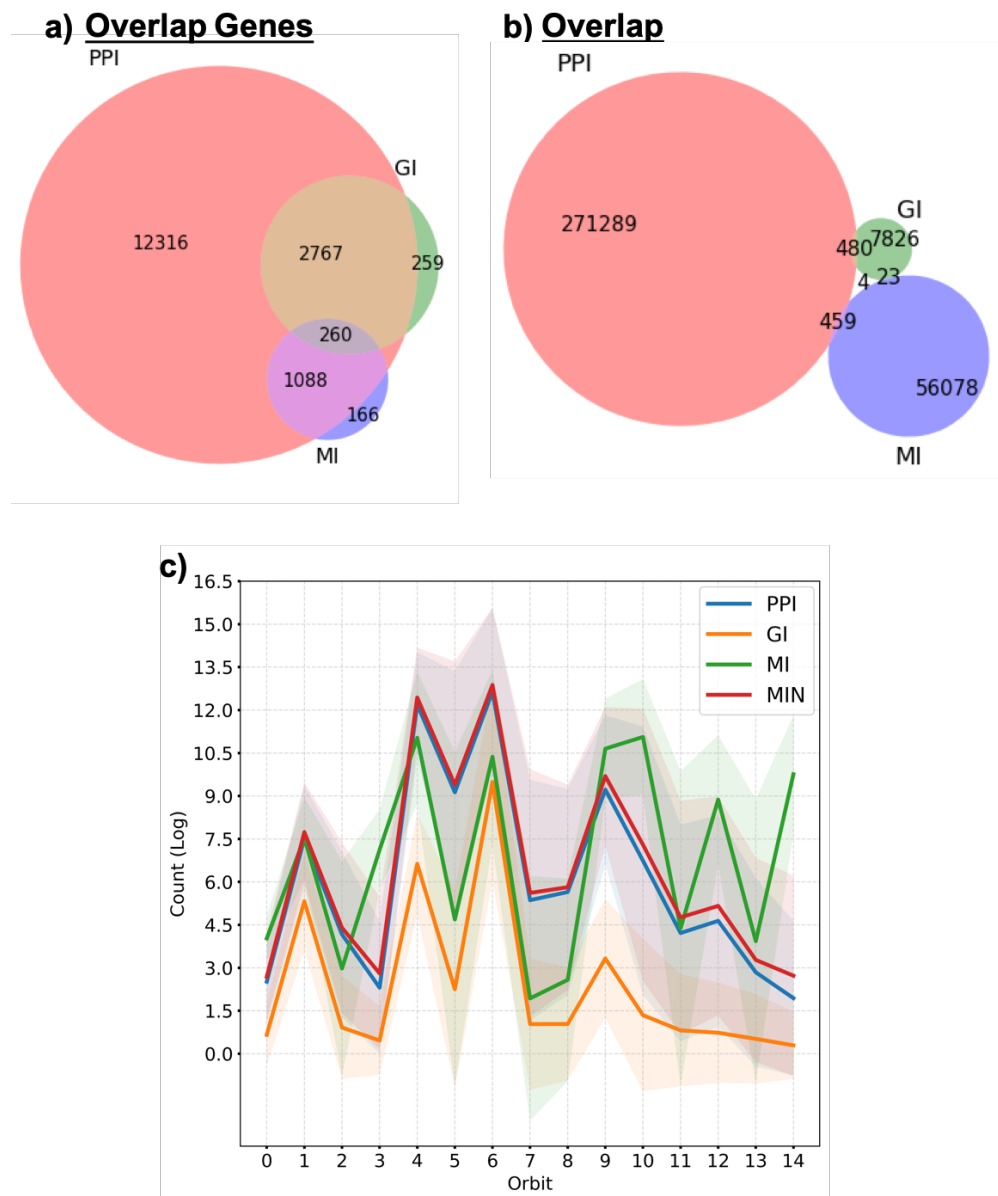
The holistic view of the human interactome increases the number of predicted DTIs

To assess which is the improvement when using the holistic view of the relationship between genes (i.e., MIN) instead of the PPI network, we applied the same framework only using the PPI network as the relation that must be conserved during the data fusion process.

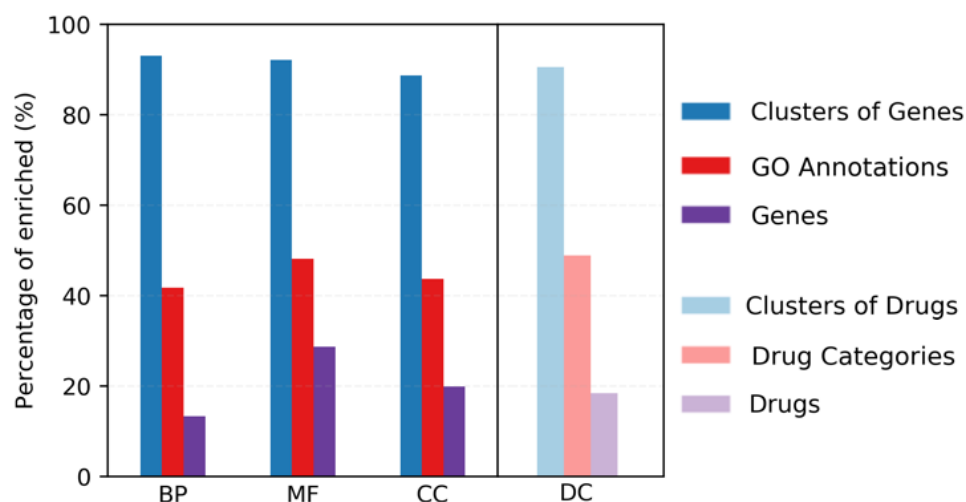
After obtaining the factor matrices, we cluster the genes and drugs by applying hard clustering to the corresponding matrix factors, G_2 and G_3 , respectively (for more details see “Extracting clusters of genes and drug” in Methods). As a validation step, we assess whether the genes and the drugs of the formed clusters are biologically and mechanistically related, respectively, by performing enrichment analysis with the genes in the G_2 matrix annotated with Gene Ontology (GO) terms and the drugs in the G_3 matrix annotated with their associated DrugBank “Drug Category” (DC) (for more details see “Enrichment analysis of gene and drug clusters” in Methods). As shown in Figure S2, more than 80% of the gene clusters show enrichments for each of the three GO domains (i.e. Biological Process, Cellular Component, Molecular Function), while the percentage of enriched genes in the clusters is at least 15%. Similarly, 90% of the drug clusters are enriched in DC (Figure S2). These results are very similar to those obtained when using the MIN (Figure 2 in the main document). Thus, the fusion framework successfully captures meaningful information encoded in the network either using the MIN or PPI network as the input of the genes relation that must be conserved.

To predict new, previously unobserved drug-gene interactions, we use the matrix completion property of the reconstructed drug-target relation matrix, $\widehat{R}_{23} \approx G_2 H_{23} G_3^\top$ (for more details see “Prediction of new drug-target interactions for drug re-purposing” in Methods). Each entry of the reconstructed matrix, corresponding to a drug-gene pair, contains an association score s_A which can be interpreted as a relative measure of confidence for each drug-gene pair. For assessing the accuracy of our predictions, we create precision-recall (PR) and receiver operating characteristic (ROC) curves, by using as ground truth the known DTIs (for more detail see “Prediction of new drug-target interactions for drug re-purposing” in Methods). As shown in Supplementary Figure S3, these PR and ROC curves are very similar when using MIN or PPI, having the same ROC-AUC ($ROC-AUC = 0.997$) and almost identical PR-AUC ($PR-AUC_{PPI} = 0.704$; $PR-AUC_{MIN} = 0.696$). To select a threshold for predicting new DTIs, we utilize the F1-score, which is the harmonic mean of precision and recall. The best F1-score ($F_1 = 0.733$) is associated with a threshold of $s_A = 0.340$, yielding 533 newly predicted DTIs with 399 drugs targeting 131 genes (Supplementary Table S3). The list of predicted DTI using MIN and PPI, have an overlap of 500 DTIs (see Supplementary Figure S3), meaning that by using PPI only 61.43% of the DTIs predicted using the MIN were also predicted when using the PPI. Moreover, only 33 out of the 533 DTIs predicted by using the PPI were not predicted by using the MIN (23 targeted by FDA-approved drugs and 10 by experimental ones). In particular, these 33 DTIs have small association scores (i.e., they are at the bottom of the list). Therefore, we obtain more putative DTIs by enforcing that the framework preserved not only PPI between genes but also GI and MI.

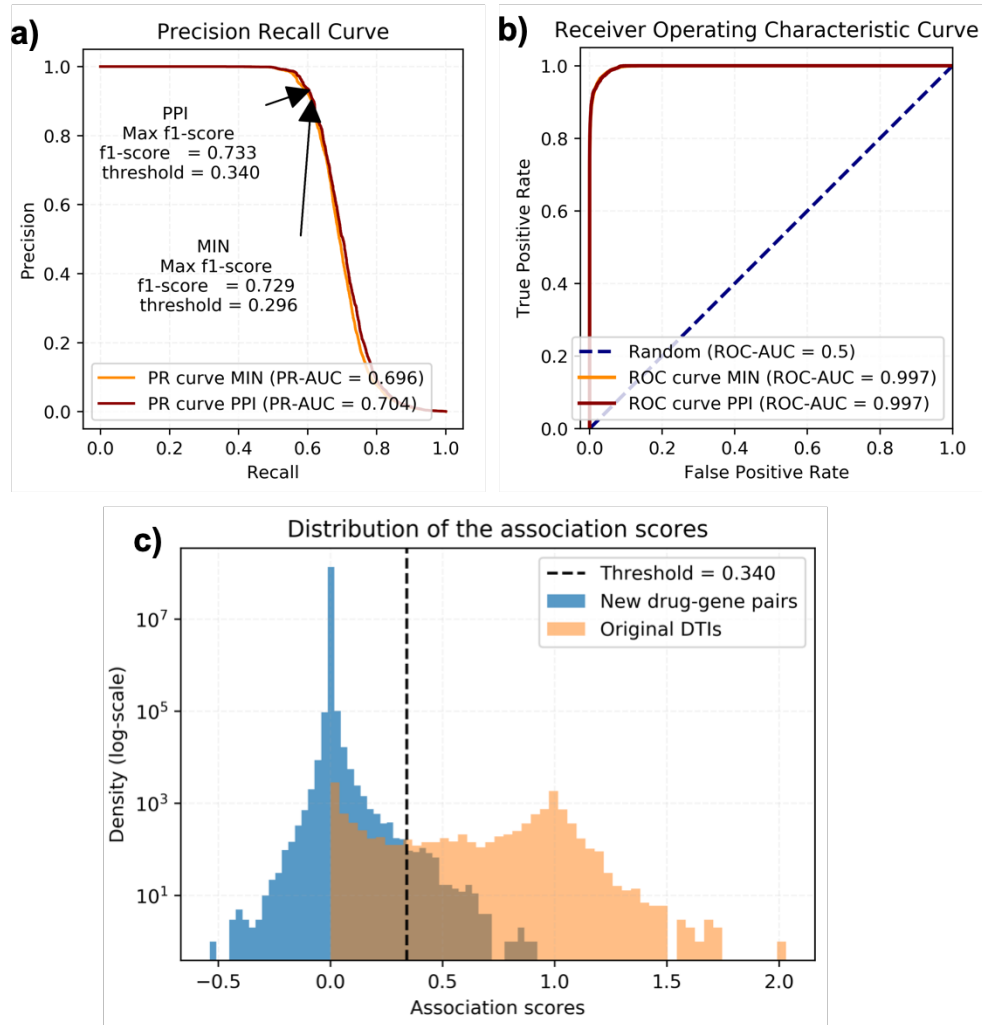
Supplementary Figures



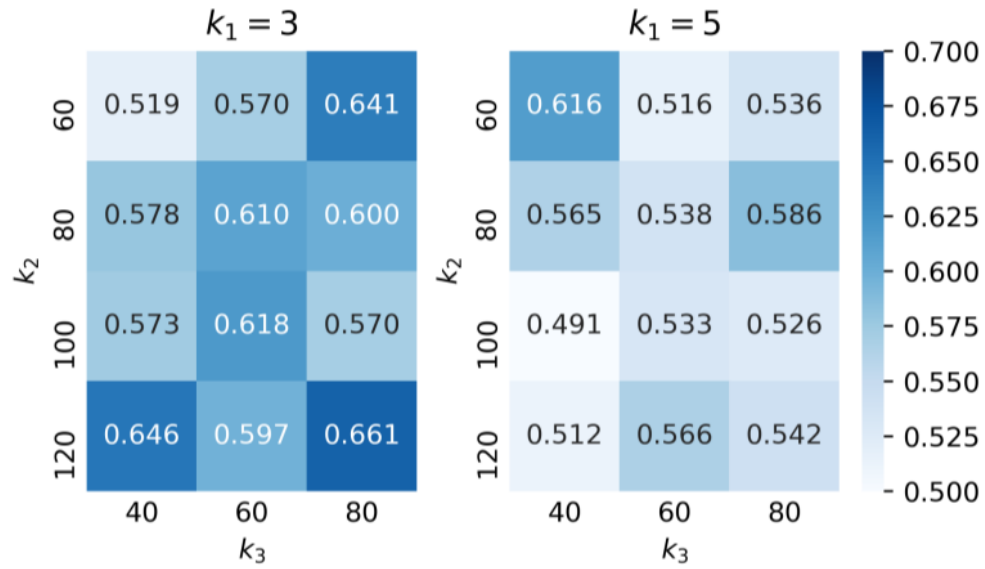
Supplementary Figure S1. Comparison between the molecular interaction network (MIN) and its constituent networks: protein-protein interactions (PPI), genetic interactions (GI) and metabolic interactions (MI) networks. A-B) Overlap of the genes and interactions of the constituent networks, respectively. C) GDV signature for the constituent networks and the MIN; counts (on the vertical axis) of the orbits (denoted by 0 to 14 on the horizontal axis).



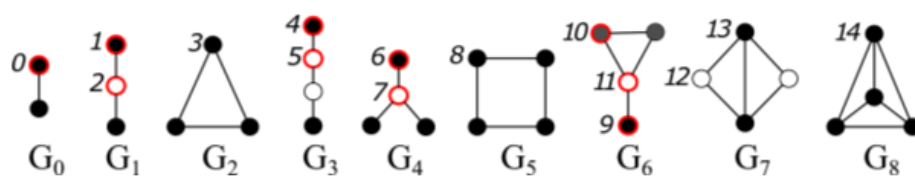
Supplementary Figure S2. Enrichment analysis for assessing the functional relevance of the gene and drug clusters obtained by the framework by only using the PPI. The gene clusters are analyzed by using GO term annotations for the three domains: Biological Process (BP), Molecular Function (MF) and Cellular Component (CC); and the drug clusters are analyzed by using “Drug Categories” (DC) from DrugBank (horizontal axis). The probability that an annotation is enriched in a cluster was computed using a hypergeometric test. Then, we computed three types of enrichments: the percentage of clusters of genes (drugs) having at least one of their genes (drugs) enriched (in blue); the percentage of GO annotations (Drug Categories) enriched (in red); and the percentage of genes (drugs) having at least one of their annotations enriched in their clusters over all annotated genes (drugs) (in purple).



Supplementary Figure S3. Prediction of new DTIs using only PPI as relation between genes. A-B) Comparison of the Precision-Recall (PR) and Receiver Operating Characteristic (ROC) curves for the framework using the MIN or only the PPI as the relation between genes. AUC - area under the curve. **C)** Distribution of the association scores of the reconstructed matrix when only using the PPI as the relation between the genes; for the original DTIs (orange) and new drug-gene pairs obtained due to the matrix completion property of GNMTF (blue). New drug-gene pairs on the right side of the threshold (dashed line) were considered to be newly predicted DTIs.



Supplementary Figure S4. Mean of the three dispersion coefficients (ρ_1, ρ_2, ρ_3) for all the values explored for choosing the parameters k_1, k_2 and k_3 . Coefficients for $k_1 = 3$ are on the left and for $k_1 = 5$ on the right. The most stable clustering was achieved by $k_1 = 3, k_2 = 120$ and $k_3 = 80$ ($mean_{\rho_1, \rho_2, \rho_3} = 0.661$).



Supplementary Figure S5. Illustration of graphlets up to 4-nodes and their 15 automorphism orbits. The ten non-redundant orbits, whose counts cannot be derived from the counts of the other orbits, are highlighted in red.

Supplementary Tables

Note: Supplementary Tables S2, S3, S4, S5, S7, S8, S9, S10 and S11 are provided as comma-separated values (csv) files due to its extension.

Supplementary Table S1. Network properties of the molecular interaction network (MIN) and its constituent networks: protein-protein interaction (PPI), genetic interaction (GI) network and metabolic interaction (MI) network.

The four networks are compared by the following commonly used network properties: four centrality measures (degree, eigenvector, betweenness and closeness centrality) and clustering coefficient.

	Average Degree	Eigenvector Centrality	Clustering Coefficient	Betweenness Centrality	Closeness Centrality
PPI	33.14	0.0032	0.1105	0.0001	0.3269
GI	5.05	0.0036	0.0623	0.0008	0.2552
MI	73.94	0.0137	0.8445	0.0011	0.344
MIN	39.85	0.0034	0.153	0.0001	0.3317

Supplementary Table S6. Functional enrichment analysis of the VI-unique neighbors genes using the gprofiler2 python-package. The table includes the adjusted p-values (P-value), the database and its domain (Source), the id and description of the term (Term id and Term name), number of genes associated with the term (Term size), number of the genes in the query gene set that are found in the database (Query size), as well as their intersection (Intersection size), precision and recall are defined as intersection_size/query_size and intersection_size/term_size, respectively.

P-value	Source	Term id	Term name	Term size	Query size	Intersection size	Precision	Recall
$4.98 \cdot 10^{-2}$	CORUM	CORUM:2018	IL12A-IL12B complex	2	51	2	0.0392	1
$4.98 \cdot 10^{-2}$	CORUM	CORUM:5548	IL12A-IL12B complex	2	51	2	0.0392	1
$4.98 \cdot 10^{-2}$	CORUM	CORUM:6102	ZP3-IZUMO1 complex	2	51	2	0.0392	1
$2.91 \cdot 10^{-2}$	GO:CC	GO:0098591	external side of apical plasma membrane	5	403	3	0.0074	0.6
$4.99 \cdot 10^{-2}$	KEGG	KEGG:00563	Glycosylphosphatidylinositol (GPI)-anchor biosynthesis	25	126	4	0.0317	0.16

End table

Supplementary Table S2. Predicted DTIs obtained by the data fusion framework. The list contains 814 newly predicted DTIs with 565 drugs targeting 172 genes. For each DTI the table contains the following information: the targeted gene (Gene); the drug involved in the DTI (DrugBank ID and Drug Name); association score (Score); DrugBank status of the drug, either FDA-approved or experimental (Drug Status); the gene set to which the targeted gene belongs, either VI, DEG, VI-unique neighbor, DEG-unique neighbor, common neighbor or background (Gene Description); and the databases in which the DTI was validated, either DrugCentral, CTD, TTD or PharmGKB (External Database).

Supplementary Table S3. Predicted DTIs obtained by the data fusion framework using only the PPI. The list contains 533 newly predicted DTIs with 399 drugs targeting 131 genes. For each DTI the table contains the following information: the targeted gene (Gene); the drug involved in the DTI (DrugBank ID and Drug Name); association score (Score); DrugBank status of the drug, either FDA-approved or experimental (Drug Status).

Supplementary Table S4. GDV counts (signature) comparison between the different gene sets. Each orbit (first column) is compared pair-wisely through all the gene sets (rest of the columns). We used the Mann-Whitney U test (with a significance level of 0.05) for each pair of orbits. The gene sets are: Viral interactors (VI), differentially expressed genes after infection (DEG), overlap of the direct network neighbors these two sets (common neighbors), neighbors of the VI and DEG gene set that were not in the common neighbors gene set (VI-unique neighbors and DEG-unique neighbors), and the rest of the genes in the MIN (background genes).

Supplementary Table S5. Functional enrichment analysis of the common neighbor genes using the gprofiler2 python-package. The table includes the adjusted p-values (P-value), the database and its domain (Source), the id and description of the term (Term id and Term name), number of genes associated with the term (Term size), number of the genes in the query gene set that are found in the database (Query size), as well as their intersection (Intersection size), precision and recall are defined as $\text{intersection_size}/\text{query_size}$ and $\text{intersection_size}/\text{term_size}$, respectively.

Supplementary Table S7. Functional enrichment analysis of the DEG-unique neighbor genes using the gprofiler2 python-package. The table includes the adjusted p-values (P-value), the database and its domain (Source), the id and description of the term (Term id and Term name), number of genes associated with the term (Term size), number of the genes in the query gene set that are found in the database (Query size), as well as their intersection (Intersection size), precision and recall are defined as $\text{intersection_size}/\text{query_size}$ and $\text{intersection_size}/\text{term_size}$, respectively.

Supplementary Table S8. Functional enrichment analysis of the background genes using the gprofiler2 python-package. The table includes the adjusted p-values (P-value), the database and its domain (Source), the id and description of the term (Term id and Term name), number of genes associated with the term (Term size), number of the genes in the query gene set that are found in the database (Query size), as well as their intersection (Intersection size), precision and recall are defined as $\text{intersection_size}/\text{query_size}$ and $\text{intersection_size}/\text{term_size}$, respectively.

Supplementary Table S9. List of the 5870 common neighbor genes in the MIN. The common neighbor genes are at the same time neighbors of the human proteins targeted by SARS-CoV-2 proteins (viral interactors (VI)) and neighbors of the genes that are differentially expressed after the infection (DEGs).

Supplementary Table S10. Predicted DTIs containing FDA-approved drugs targeting common neighbors (common neighbors DTIs). The list contains 185 newly predicted DTIs targeting 49 common neighbors with 149 drugs. For each DTI the table contains the following information: the targeted gene (Gene), the drug involved in the DTI (DrugBank ID and Drug Name), association score (Score), the databases in which the DTI was validated (External Database), whether the drugs were previously found in COVID-19 context (CORDITE) and number of interventional clinical trials for COVID-19 in which the drugs are currently studied (Clinical Trials).

Supplementary Table S11. Functional enrichment analysis of 49 genes involved in the 185 common neighbor DTIs using the gprofiler2 python-package. The table includes the adjusted p-values (P-value), the database and its domain (Source), the id and description of the term (Term id and Term name), number of genes associated with the term (Term size), number of the genes in the query gene set that are found in the database (Query size), as well as their intersection (Intersection size), precision and recall are defined as $\text{intersection_size}/\text{query_size}$ and $\text{intersection_size}/\text{term_size}$, respectively.