

# Higher order molecular organisation as a source of biological function

Thomas Gaudelet, Noel Malod-Dognin and Natasa Przulj\*

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

## Abstract

**Motivation:** Molecular interactions have widely been modelled as networks. The local wiring patterns around molecules in molecular networks are linked with their biological functions. However, networks model only pairwise interactions between molecules and cannot explicitly and directly capture the higher order molecular organisation, such as protein complexes and pathways. Hence, we ask if *hypergraphs* (*hypernetworks*), that directly capture entire complexes and pathways along with protein-protein interactions (PPIs), carry additional functional information beyond what can be uncovered from networks of pairwise molecular interactions. The mathematical formalism of a hypergraph has long been known, but not often used in studying molecular networks due to the lack of sophisticated algorithms for mining the underlying biological information hidden in the wiring patterns of molecular systems modelled as hypernetworks.

**Results:** We propose a new, *multi-scale*, protein interaction *hypernetwork model* that utilizes hypergraphs to capture different scales of protein organization, including PPIs, protein complexes and pathways. In analogy to graphlets, we introduce *hypergraphlets*, small, connected, non-isomorphic, induced sub-hypergraphs of a hypergraph, to quantify the local wiring patterns of these multi-scale molecular hypergraphs and to mine them for new biological information. We apply them to model the multi-scale protein networks of baker's yeast and human and show that the higher order molecular organisation captured by these hypergraphs is strongly related to the underlying biology. Importantly, we demonstrate that our new models and data mining tools reveal different, but complementary biological information compared to classical PPI networks. We apply our hypergraphlets to successfully predict biological functions of uncharacterised proteins.

**Availability:** Code and data are available online at <http://www0.cs.ucl.ac.uk/staff/natasa/hypergraphlets>

**Contact:** natasa@cs.ucl.ac.uk

## 1 Introduction

Deciphering the complex patterns of interactions between macromolecules in a cell is of crucial importance. Graph theory offers mathematical abstractions to represent and study molecular interactions. Simple *graphs* (also called *networks*) have been widely used to model the interactions between pairs of molecules. For instance, in Protein-Protein Interaction (PPI) networks, each node represents a protein and each edge connects a pair of proteins that can bind to each other [45, 18, 42, 39]. Exact comparison of networks is a hard problem due to the NP-completeness of the underlying subgraph isomorphism problem [9]. Thus, simple heuristics have been used to study PPI and other molecular networks, such as degree distribution and centralities [30]. *Graphlets* quantify the local topology of a network. They are small, non-isomorphic, induced subgraphs of a larger network, which precisely characterise the local wiring patterns around each node [36, 35]. Graphlets and their statistics have since been used to compare biological networks [49], to uncover their functional organisation [36, 35, 31, 49], to guide network alignment algorithms [22, 29], or to relate the wiring patterns of genes in these networks with their biological functions [31, 49, 10].

However, in biological systems, molecules do not interact solely in a pairwise fashion. Hence, simple graphs do not capture the multi-scale organisation of these systems [23, 21]. In the example

---

\*natasa@cs.ucl.ac.uk

in Figure 1, we observe that the simple graph representation, on the right, of the system on the left blurs the higher-order organisation of the system. Given only the network representation on the right, one might, for instance, falsely assume that the nodes b, c, and d form a complex of three elements, while it is true that b and d form a complex, b and c form a complex, and c, d and e form a complex.

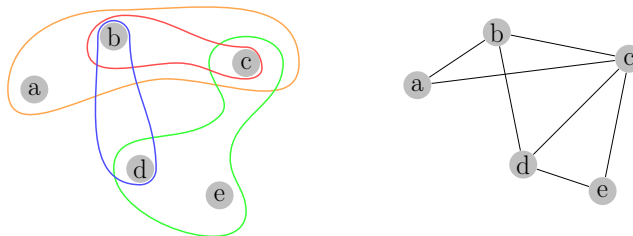


Figure 1: Illustration of a system with higher order interactions (left) and its simple graph representation (right).

A solution to overcome this limitation is to model a molecular system using hypergraphs. A *hypergraph* is defined by a set of nodes,  $V$ , and a set of edges,  $E$ , called *hyperedges*, where each hyperedge corresponds to a set of interacting nodes of any size [3]. This means that a simple graph is a special case of a hypergraph in which all hyperedges are sets of two nodes. The representation of the system in Figure 1 (left) is a hypergraph. To analyse data modelled as hypergraphs, it is necessary to develop methods to mine the structure of hypergraphs. A number of simple measures from graph theory have already been extended to hypergraphs, e.g., the clustering coefficient [12], degree distribution [24], and centralities [12, 33]. Approaches such as percolation and random walks [1, 32] have also been extended to study hypergraphs. Hypergraphs have also been used for learning tasks, such as clustering and nodes classification [44, 34]. However, hypergraphs lack more advanced descriptors of local topology. Hence, we introduce hypergraphlets, an extension of graphlets to hypernetworks.

We investigate biological hypernetworks in which nodes are proteins and hyperedges capture PPIs, protein complexes, or signaling pathways. A protein complex connects two or more proteins that bind together. A pathway connects together any number of proteins whose interactions, including (but not limited to) PPIs, leads to a certain product or change in a cell. The main aim is to check if the topology of these hypernetwork representations of the data carries biological information that goes beyond the information that can be obtained from PPI networks. We use hypergraphlets in this investigation.

## 2 Contributions

We motivate studying the higher order molecular interactions as models that capture additional and different biological information than the widely studied PPI networks. We introduce hypergraphlets as a new tool that unveils the pioneering observation of the close link between the multi-scale molecular organisation and biological function and that can serve as an underlying methodology for many new tools that will be developed to further study the multi-scale organisation of molecular systems.

We analyse the hypergraph representation of protein interactions of yeast *saccharomyces cerevisiae* and human and show that proteins that are similarly wired in a hypernetwork, independently of their location in the hypernetwork, tend to have similar biological functions. Also, we use the Canonical Correlation Analysis (CCA) [16] to correlate hypergraphlets around proteins in these networks with their biological functions. The results confirm the link between the local wiring patterns of the multi-scale molecular organisation of the cell and biological functions. We use these findings to predict biological functions of uncharacterised proteins from the wiring patterns of the multi-scale molecular organisation. We validate our predictions in the literature.

## 3 Materials & Methods

### 3.1 Data

We consider six different networks across two species, human and baker’s yeast. For each species, we consider the protein-protein interaction (PPI) network and two hypernetworks corresponding to protein complexes and biological pathways. In all networks, nodes correspond to proteins. In a PPI network, an edge between two proteins represents a physical interaction. Depending on the hypernetwork considered, a hyperedge represents either a protein complex or a biological pathway. These data are used jointly to build hypernetworks capturing multi-scale organisation of proteins in a cell, as detailed in Section 3.5 below.

The PPI data is obtained from the BioGRID database [8] (version 3.4.145). Both pathways hypernetworks come from the Reactome database [13] (accessed in April 2017). The human protein complexes are downloaded from the CORUM database [40, 41] (in May 2017), while the yeast protein complexes are collected from the CYC2008 database [37] (last updated in 2009). Table 1 gives an overview of the sizes of the data sets.

	Database	# proteins	# (hyper-) interactions
Human	CORUM	3,145	2,138
	Reactome	9,466	1,461
	PPI	16,008	216,865
Yeast	Reactome	1,465	400
	Cyc2008	1,607	406
	PPI	5,931	87,225

Table 1: Sizes of the data.

To investigate the links between networks and biological functions, we collect gene annotations from the Gene Ontology Consortium (GO) database [5] (downloaded at the end of January 2017). For each protein, we keep only the most specific annotations that are experimentally derived. We separate the annotations based on the three categories: Biological Process (BP), Molecular Function (MF), and Cellular Component (CC).

### 3.2 Hypergraphlets: the local topology of hypergraphs

We define *hypergraphlets* as small, connected, non-isomorphic, induced sub-hypergraphs of larger hypergraphs. [3] defines an induced sub-hypergraph of a hypergraph  $H = (V, E)$  on a set of nodes  $A \subset V$  as the hypergraph  $H_A$  with set of nodes  $A$  and set of unique hyperedges

$$E_{H_A} = \{e \cap A \mid e \in E, e \cap A \neq \emptyset\}. \quad (1)$$

Note that with this definition, hyperedges containing only one node exist for each node. With this definition, an induced hypergraph is simple, i.e. it has no duplicated edges.

Within a given hypergraph, automorphic nodes are nodes whose labels can be exchanged without changing adjacency relationships. Formally these nodes can be mapped to each other by an *automorphism*, which is an isomorphism of a hypergraph with itself. An *isomorphism* is a mapping of nodes of the hypergraph that preserves the adjacency of the nodes [6]. A set of automorphic nodes form what is called an *orbit*. Here, we consider all 1- to 4-node hypergraphlets, which contain a total of 6,369 different orbits. For 5-node hypergraphlets, we estimate that there are more than a hundred thousands orbits, hence we restrict ourselves to 4-node hypergraphlets. In Figure 2, we illustrate all 65 orbits that occur in the 1- to 3-node hypergraphlets.

Analogous to graphlets, we use hypergraphlet orbits to quantify the wiring patterns around each node in a hypergraph. For each orbit  $i$  in hypergraphlet  $h$ , we define the  $i^{\text{th}}$  *hypergraphlet degree* of a node in the hypergraph  $H$  as the number of hypergraphlet orbits  $i$  that the node touches.

For each node in a hypergraph, we compute all 6,369 hypergraphlet degrees resulting in a 6,369-dimensional vector where entry  $i$  corresponds to the  $i^{\text{th}}$  hypergraphlet degree of the node. We term this vector capturing the local wiring around a node the *Hypergraphlet Degree Vector (HDV)*.

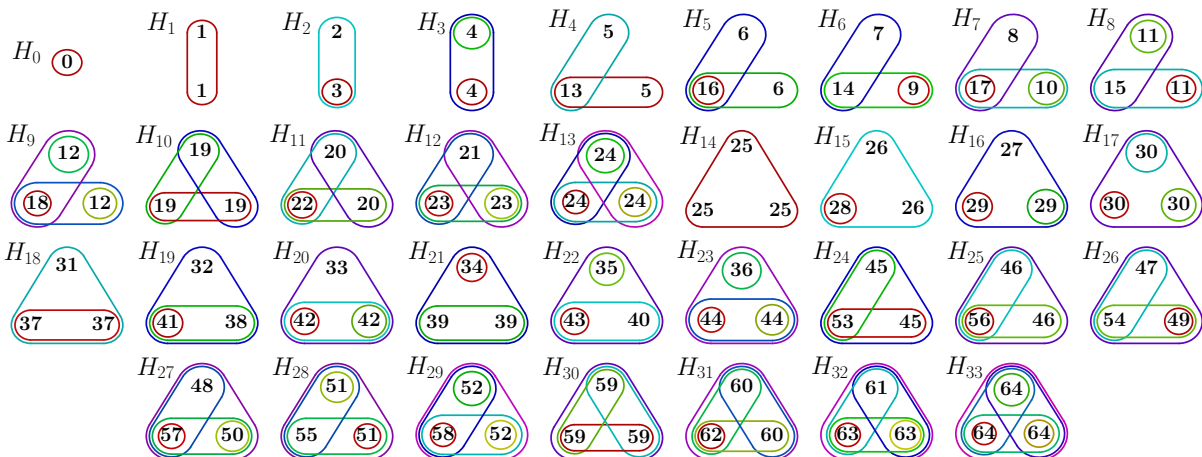


Figure 2: Illustration of all 1- to 3-node hypergraphlets ( $H_0$  to  $H_{33}$ ) and the 65 orbits. Each closed set corresponds to a hyperedge and each node is represented by an integer between 0 and 64 corresponding to the orbit it belongs to.

Considering a hypergraph with  $n$  nodes, with maximal hyperedge of size  $l$  and with maximal degree of a node  $d$ , where the *degree* of a node corresponds to the number of hyperedges that contain it, an upper bound on the complexity of counting all 1- to  $k$ -node hypergraphlets is  $O(n(ld)^{k-1})$ .

[27] introduced an alternative definition of hypergraphlets in the context of binary classification problems. They define kernels based on their definition of hypergraphlets and use support vector machines to classify the proteins. The key difference with our definition of hypergraphlets is that they do not consider the hypergraphlets of a hypergraph as *induced* sub-hypergraphs, thus ignoring some overlaps between hyperedges [27]. In particular, in the first step, they ignore all hyperedges containing more than four nodes. Instead, hyperedges with more than four nodes are taken into consideration independently in the second step, which decomposes a hyperedge of size  $n > 4$  into the  $\binom{n}{4}$  subsets of four nodes. Hence, with their definition and counting process, an important part of the topology of the hypernetwork is overlooked and therefore topological information is lost, which motivates our redefinition that is also a direct extension of the definition of graphlets for simple graphs. However, we could not compare the two approaches, as their implementation is not publicly available and they recently agreed with us that their definition needed to be changed to alleviate these issues<sup>1</sup>.

### 3.3 Topological distance

We define a distance measure to compare the wiring patterns of two nodes in a hypernetwork (or network, depending on the model considered) as follows. Consider a set of proteins  $P = \{p_1, p_2, \dots, p_m\}$  and let  $M$  be the matrix representing our data where row  $i$  corresponds to the HDV (or GDV) of protein  $p_i$ . Then, we define the distance,  $\delta$ , between two proteins  $p_i$  and  $p_j$  as

$$\delta(p_i, p_j) = \left[ \sum_{k \in K} \left( \frac{\log(M_{ik} + 1) - \log(M_{jk} + 1)}{\sigma_k} \right)^2 \right]^{\frac{1}{2}}, \quad (2)$$

where  $K$  corresponds to the set of orbits considered,  $M_{ik}$  denotes the entry of  $M$  on the  $i^{\text{th}}$  row and  $k^{\text{th}}$  column, and  $\sigma_k$  denotes the standard deviation of the distribution of the  $k^{\text{th}}$  hypergraphlet (or graphlet) orbit degree across our set of data value. Note that to reduce the impact of very large orbit counts we apply to  $M$  an element-wise log transformation.

### 3.4 Linking local structure to function

We explore two ways to evaluate the link between the local structure of a molecular network and the biological functions of its molecules. First, we cluster the nodes based on the similarity of their

<sup>1</sup>Personal communication.

wiring patterns defined in Section 3.3, and we do the enrichment analysis of the resulting clusters (Section 3.4.1). Second, we use CCA to test if biological functions tend to be characterised by specific wiring patterns (Section 3.4.2).

### 3.4.1 Cluster enrichment

We cluster proteins that are similarly wired in a graph or a hypergraph as measured by distance  $\delta$  (see Equation 2) and test if the proteins within the same cluster share GO functions.

Clusters are obtained by using k-means method [17] based on the distance defined in Equation 2. For each of various numbers of clusters,  $k$ , we run the clustering algorithm 20 times to account for the randomness in the k-means algorithm. For each clustering, we compute the enrichment of clusters in biological annotations for each GO category with correction for multiple hypothesis testing [2]. We consider a cluster enriched if at least one GO annotation is significantly enriched in the cluster (p-value < 5%). For each value of  $k$ , we also compute the average of Sum of Squared Error (SSE) and the Normalised Mutual Information (NMI) [47] considering all 20 repeats. SSE gives a measure of how close proteins within a cluster are on average according to our similarity measure, while NMI evaluates the stability of the clustering across the 20 runs, i.e. if proteins are consistently clustered together or apart. Then, we use “the elbow” analysis of the SSE and NMI with respect to  $k$  to choose the optimal number of clusters. For the resulting number of clusters, we select the clustering giving the highest percentage enrichment across the 20 runs of k-means for each GO category. We test the significance of the enrichment with random permutation tests: we keep the same number and size of clusters and randomly assign proteins to each cluster and measure the enrichments of the resulting clusters. We repeat this process 1,000 times and compute the significance.

To see whether the two models, networks and hypernetworks, harbour the same or different but complementary biological information, at least to the extent that it can be uncovered by the proposed methodologies, we measure Adjusted Mutual Information (AMI) [47] of the clusters and Jaccard Index [19] of the enriched annotations in the clusters. AMI is a variation of Mutual Information (MI) used to compare two clusterings. It measures if any pair of proteins is consistently clustered together or apart in both clusterings adjusting for chance. The Jaccard Index gives a measure of the overlap between the two sets of GO annotations.

### 3.4.2 Canonical correlation analysis

CCA is used to infer correlations between two sets of features,  $X$  and  $Y$ . Consider features  $X = (X_1, \dots, X_n)$  and  $Y = (Y_1, \dots, Y_m)$  over the same elements. Then CCA will identify  $K$  pairs  $(\mathcal{L}_X^k, \mathcal{L}_Y^k)$ , called *canonical variates*, of linear combinations of features of  $X$  and of features of  $Y$ , with  $K = \min(m, n)$ , such that the correlations of  $\mathcal{L}_x^k$  and  $\mathcal{L}_y^k$  are maximal over all  $k$ . Each canonical variate is associated a score corresponding to the correlation between its two linear combinations.

In our case, the elements are proteins, the first set of features corresponds to the wiring patterns of proteins in networks or hypernetworks, and the second to the biological functions of proteins from GO. As mentioned above, each protein (node) has a GDV from the PPI network and an HDV from the hypernetwork. Hence, we have two matrices of topological features where entries  $(i, j)$  correspond to the  $j^{\text{th}}$  orbit degree of protein  $i$ . Also, we associate to each protein three vectors of GO annotations, one for each of the categories: BP, MF, and CC. In each of these vectors, an entry is equal to 1 if the gene is annotated with the corresponding GO term, and 0 otherwise. Hence, we form three matrices of biological features, where entries  $(i, j)$  correspond to the presence or absence of GO annotation  $j$  for protein  $i$ .

We compute CCA for each combination of topological features and biological annotations to uncover topology-function relationships in the data.

## 3.5 Summary of the analysis

As stated above, our main aim is to examine if modelling the higher order of molecular organisation harbours additional biological information and to demonstrate that the wiring patterns of biological hypernetworks are strongly linked to the underlying biology.

We compute vectors containing topological information around proteins in the molecular networks: we use graphlets on PPI networks and hypergraphlets on hypergraphs, as described above. To validate our approach, we focus on parts of PPI networks that we know are rich in biological information: protein complexes and pathways. Clearly, not all proteins in a PPI network belong

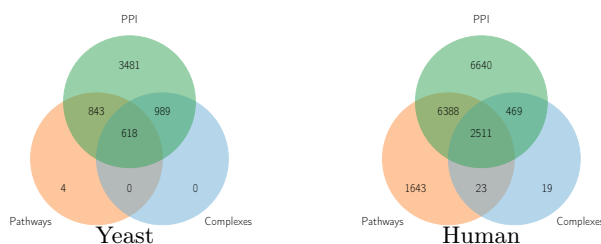


Figure 3: The overlaps of the protein sets of baker’s yeast (left) and human (right). Left: 3,481 proteins participate in PPIs only, 843 in PPIs and pathways, 618 in PPIs, pathways and complexes, 989 in PPIs and complexes, while 4 are in pathways only. Right: 6,640 proteins participate in PPIs only, 6,388 in PPIs and pathways, 2,511 in PPIs, pathways and complexes, 469 in PPIs and complexes, 23 in complexes and pathways, while 1,643 are in pathways only and 19 in complexes only.

to complexes, or pathways (illustrated in Figure 3). Hence to validate our method, we consider four sets of proteins: those belonging to pathways in human (human-pathways), those belonging to pathways in yeast (yeast-pathways), those belonging to complexes in human (human-complexes), and those belonging to complexes in yeast (yeast-complexes). For each protein in each of these sets, we have two topological signatures: one from the standard graphlets counted on the entire PPI network and one from the hypergraphlet counts in the hypergraph (HG) that we constructed by using only protein complexes (and equivalently pathways). That is, in an HG, nodes are proteins and each hyperedge represent a protein complex (or pathway) and contains the proteins that belongs to the complex (pathway). For each protein, we also have three biological signatures corresponding to the three levels of GO annotations: BP, MF, and CC. We use these as input into the methods described in Sections 3.4.1 and 3.4.2. The results of these validations are presented in Sections 4.1.1 and 4.1.2.

The reason for doing these validations on the sets of data for which we know that they are very enriched in biological information (i.e., pathways and complexes) is to demonstrate that our new model and method can correctly identify the biological information. After these validations of the methodology, we use it to perform the analysis of multi-scale protein interaction network data of yeast and human and uncover new biological information. In particular, for each species, we construct a hypergraph that contains all of its PPIs, all of its protein complexes, and all of its pathways; i.e., nodes are proteins and hyperedges correspond to PPIs, protein complexes, and pathways. The results of analysing these hypergraphs with our methods are presented in Section 4.2.

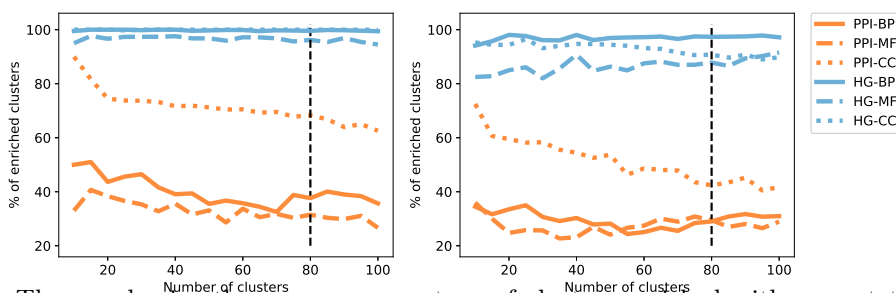
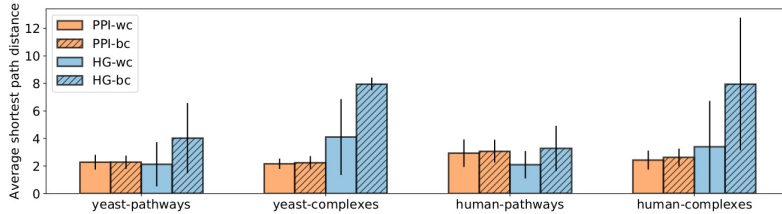


Figure 4: The panels give the average percentage of clusters enriched with respect to the total number of clusters for yeast-complexes (left) and yeast-pathways (right), the standard deviation is not represented to avoid overcrowding the panels. The colors represent the models from which the clustering is obtained: HG in blue and PPI in orange. The type of line represents the category of GO annotations: BP are full lines, MF are dashed lines, and CC are dotted line. The black vertical lines signal the number of clusters selected from the set of NMI and SSE curves according to the procedure described in Section 3.4.1.

	Biological Process		Molecular Function		Cellular Component	
	HG	PPI	HG	PPI	HG	PPI
Yeast-complexes	100% (51)	53.75% (80)	100% (49)	51.25% (80)	100% (51)	74.7% (79)
Yeast-pathways	100% (71)	45% (80)	95.2% (63)	37.5% (80)	95.4% (65)	56.25% (80)
Human-complexes	94.3% (105)	40.3% (119)	82.7% (98)	47.5% (120)	95.2% (105)	60.8% (120)
Human-pathways	98.2% (111)	59.2% (120)	98.3% (115)	70.8% (120)	96.6% (118)	63.3% (120)



	Biological Process	Molecular Function	Cellular Component
Yeast-complexes	0.11 (0.0)	0.1 (0.0)	0.1 (0.0)
Yeast-pathways	0.07 (0.0)	0.07 (0.0)	0.07 (0.0)
Human-complexes	0.07 (0.01)	0.07 (0.02)	0.08 (0.06)
Human-pathways	0.05 (0.07)	0.06 (0.1)	0.06 (0.12)

Figure 5: The top table presents the maximum enrichment measured across clusterings obtained with the “optimal” number of clusters (80 for yeast and 120 for human). The number in parenthesis is the number of non-empty clusters. The color indicates the statistical significance of the maximum enrichment with respect to random permutation tests: black indicates a significant value, grey a non-significant one. The middle panel gives, for each type of model (HG in blue and PPI in orange), the average of the shortest path lengths within the clusters (wc) and between clusters (bc) of the best clustering obtained for GO-BP annotations. The results are similar for other GO categories and are not presented here due to space limitations. The bottom table presents the results of comparing the obtained clusterings. We use the HG clustering as baseline and compute the Adjusted Mutual Information (AMI) between the clusterings and the Jaccard Index (in parenthesis) between the sets of enriched GO terms.

## 4 Results & Discussion

### 4.1 Validation of our methodology

#### 4.1.1 Enrichment Analysis

Having computed the topological vectors from both network models (PPI and HG) for each protein of each of the four sets of proteins described in Section 3.5 (human-pathways, human-complexes, yeast-pathways and yeast-complexes), we apply the methodology detailed in Section 3.4.1 to investigate if similarly wired proteins have similar functions. Interestingly, the percentage of enriched clusters is relatively stable as we increase the number of clusters. Hence, any partitioning of the proteins based on the local wiring patterns in a network, quantified by using graphlets or hypergraphlets, captures the underlying biological information (see Figure 4). This underlines the crucial role played by the way proteins interact in determining protein function without any information about their sequence, or interacting partners. Furthermore, when examining the clusterings obtained at a specific number of clusters,  $k$  (see Section 3.4.1 for details on how  $k$  is chosen), we observe that the enrichments (top table in Figure 5) are all statistically significant, except for the one in gray. Importantly, clusters obtained from HG models are more enriched than those obtained from PPI networks. This result validates the relevance of our HG modelling in capturing the underlying biological information and underlines the potential of hypergraphlets for mining molecular hypernetworks.

To further investigate the clusterings, we compute for each the average shortest path distances, in the corresponding (hyper-)network, between pairs of proteins belonging to the same clusters (“within-clusters”) and between pairs of proteins which are in different clusters (“between-clusters”; see middle panel in Figure 5). We observe a larger gap between within-cluster and between-clusters average shortest path lengths for clustering obtained from higher order molecular organisation

than from clusterings obtained from PPI networks. Hence, proteins that are topologically similar in the HG model in addition to sharing biological functions tend to be at shorter distance from each other. This result is consistent with the literature on “guilt by associations”, which predicts protein functions from their neighbourhoods in molecular networks [46].

Finally, we observe that the clusterings obtained from the PPI model are different from those obtained from the HG model both in terms of GO annotations that are enriched and in terms of clustered proteins (see bottom table in Figure 5). This is because a Jaccard Index close to 0 means that the sets of the enriched GO terms in the PPI and HG clusterings tend not to overlap. Also, AMI scores below 0.1 mean that pairs of proteins belonging to the same clusters in one clustering are typically in different clusters in the other clustering. This demonstrates that modelling the interactomes by hypergraphs will uncover new biological information that cannot be uncovered from the analysis of PPI networks. Also, it demonstrates the complementarity of the two representations and that the two are capturing different underlying biological information.

#### 4.1.2 Canonical Correlation Analysis

We investigate the existence of specific topology-function links, i.e. the connection between specific hypergraphlets (or graphlets) and GO annotations by using CCA described in Section 3.4.2. We apply it on the same PPI and HG of yeast and human used in the clustering and enrichment analysis (Section 4.1.1): for each set of proteins, we compute the CCA between the topology-containing vectors of each of the associated models (PPI and HG) and the vector of GO annotations for each category (BP, MF, and CC). Due to space limitations, we present only the results obtained for yeast and GO-BP annotations. We obtain similar results in all other cases and the discussion below holds for them as well.

We observe that each model has a number of canonical variates with correlation close to 1 (Figure 6), which indicates a strong topology-function relationship in these data that was previously highlighted in the context of economic network data [49]. In particular, this means that some functions are strongly linked to specific wiring patterns and thus, local topology can potentially be used for predicting protein functions. For that purpose, hypergraphlets of HGs have a strong advantage over graphlets of PPI networks in the number of canonical variates with a score close to 1, which is 3 to 13 times more variates with HGs. This is also expected, since we chose our hypernetworks to model already function rich parts of molecular networks, protein complexes and pathways, and it validates our methodology.

In Figure 7, we take a closer look at the most significant CCA variate. The variate score of 1.0 links a linear combination of GO annotations to a linear combination of hypergraphlets orbits. For instance, this means that a gene annotated with positive regulation of barrier spectrum assembly (GO:0010973) will likely have a relatively large 2,644<sup>th</sup> orbit degree in the hypernetwork. Why these specific orbits are linked to these functions is a question that is outside of the scope of this study and that needs to be further investigated. We find that the GO terms identified here are also biologically coherent: each of the GO-BP terms denoted in blue text in Figure 7 is annotating at least one protein conjointly with at least one other annotation, that is also denoted in blue text in Figure 7, according to QuickGO search engine [4]. Furthermore, the only remaining annotation, cell cycle arrest (GO:0007050), has been linked to the MAPK pathway in the literature [38], as have been most of the other terms [28, 15]. Hence, the entire set of GO annotations presented in Figure 7 is biologically coherent, which validates the relevance of the canonical variate and of our hypergraph-based methodology in capturing functional information.

## 4.2 Analysing multi-scale molecular organisation

To explicitly capture the multi-scale organisation of protein interactions, we model them by a hypernetwork containing all PPIs, all protein complexes and all biological pathways as hyperedges (detailed in Section 3.5). To assess if the wiring patterns in our new HG model capture the biological functions of proteins, we do the clustering and enrichment analysis (Section 3.4.1), as well as the canonical correlation analysis (Section 3.4.2) on these hypernetworks of baker’s yeast and human. We compare the results with those that we obtain by applying the same methodologies to PPI networks. In these unifying HG models of multi-scale molecular organisation, we observe that clusterings of the proteins based on their topological vectors in a network, obtained by using graphlets or hypergraphlets, capture the underlying biological information (see the top panels of Figure 8). Furthermore, the clusters obtained from the hypernetwork topology lead to higher



enrichments in GO-BP, GO-MF, and GO-CC annotations. This shows that our newly proposed model, regardless of the choice of the total number of clusters,  $k$ , captures more protein biological function in its topology than the standard PPI networks.

When choosing the number of clusters,  $k$ , according to the criteria detailed in Section 3.4.1, we observe that all enrichments are statistically significant and that the HG models allow for an increase of over 15% in the number of enriched clusters when compared to the PPI networks. This finding underlines the link between multi-scale interaction patterns and biological functions. Interestingly, when investigating the clusters, we observe that a majority of the proteins in the non-enriched clusters only have reported PPIs, but not any pathways or complexes that they belong to. This is true for 59% of the proteins in the HG model of yeast and 38% of the proteins in the HG model of human. This might be due to incompleteness of the pathways and protein complexes data. Our results indicate that when more complete data on complexes and pathways becomes available, our methodology will be able to extract additional biological information.

We observe that proteins clustered using topological features derived from representations of multi-scale molecular organisation tend to also be closer in terms of shortest path distances compared to those obtained by clusterings based on the topology of PPI networks (see bottom left panel in Figure 8). Interestingly, most proteins clustered together in the HG models are direct neighbours or second neighbours. Hence, the fact that we obtain enriched biological functions in those clusters is consistent with empirical evidences showing that 70-80% of interacting proteins share at least one function. Those evidences were the motivation for the *majority rule* used in the literature for functional prediction [46].

Finally, we observe that the clusterings obtained from the PPI models are different from those obtained from the HG models both in terms of GO annotations that are enriched, with a Jaccard Index below 0.25, and in terms of similarity of clusters, with an AMI below 0.35 (see bottom right panel in Figure 8). This confirms that our multi-scale model is not equivalent to the standard PPI network and uncover additional biological information complementary to that of the PPI network.

Using CCA (Section 3.4.2), we observe that each model has high scoring canonical variates, which indicates that some functions are strongly linked to specific wiring patterns (see Figure 9). For that purpose, hypergraphlets of our new HG models have an advantage over graphlets of PPI networks in the number of canonical variates with high correlation score: it has over 300 canonical variates with score greater than 0.9 compared to only 10 for PPI networks. This indicates that the HG model's local wiring patterns are more correlated with the underlying biology than those of the PPI networks.

Finally, we use the clusterings to investigate the potential of our newly proposed models in conjunction with our hypergraphlets to predict protein functions. As demonstrated above, we identified clusters of proteins with significantly enriched GO annotations. We use these clusters to predict the functions of proteins. For each GO category, we identify two disjoint sets of proteins in each of our hypernetworks: the set of proteins that are experimentally annotated with at least one of the enriched GO terms in their cluster (on which the enrichment computations are based) and the set of proteins that have some predicted annotation in the GO database.

First, we consider the second set and investigate how many of those proteins have at least one of the enriched terms of their cluster as their predicted GO annotation [5]. For GO-BP, this set contains 11,686 proteins for human (4,161 for yeast). For GO-MF, it contains 7,243 proteins for human (3,586 for yeast). For GO-CC, it contains 6,589 proteins for human (3,510 for yeast). We show that out of these proteins, about 5% for yeast and 15-23% for human have been putatively annotated in GO with at least one of our enriched functions in their clusters (see Figure 10), which validates our approach.

Second, we focus on the proteins of the hypernetworks that are unannotated in GO database (this corresponds to 994 proteins for human and 97 proteins for yeast) and investigate the GO-BP annotations we predict for them. We predict function for each of these proteins by associating it with the enriched experimentally obtained GO term that annotates the most proteins in its cluster. We survey the literature to validate some of our predictions<sup>2</sup> for human (the top predictions correspond to the most statistically significantly enriched GO terms). We predict that HIST1H2AJ is involved in nucleosome assembly (GO:0006334), which is confirmed in the literature [11]. We further predict that XIST is linked to chromatin organization (GO:0006325), which has also been highlighted in past studies [7]. We also predict that NME1-NME2 (an unknown protein

---

<sup>2</sup>All predictions are available online at <http://www0.cs.ucl.ac.uk/staff/natasa/hypergraphlets/>

encoded between NME1 and NME2 in the DNA) is involved in cell proliferation (GO:0008283). The function of this protein is not yet established [25], however NME2 has been linked to reduction of cell proliferation [26] and proteins encoded in the neighbour locations of the DNA tend to have similar function [14]. For microRNA mir-3606, we predict a role in collagen fibril organization (GO:0030199). Collagen plays a key role in cell adhesion, which can involve integrin [20, 43] and mir-3606 has been linked to integrin in the literature as it has been suggested that mir-3606 can bind to ITGA4 (integrin subunit alpha 4) [48]. Finally, we propose that LOC101929876 (40S ribosomal protein S26) is involved in rRNA processing (GO:0006364), which is corroborated by the Reactome database in which the protein is associated with a major pathway of rRNA processing in the nucleolus and cytosol [13].

These results confirm the ability of our hypergraphlets to predict biological functions of proteins from the wiring patterns in our novel model capturing multi-scale organisation of proteins in a cell.

## 5 Conclusion

We highlight the importance of considering the higher order organisation of protein interactions in conjunction with the standard PPI networks. We propose a novel methodology, hypergraphlets, to quantify the local wiring patterns of hypergraphs. We apply it to biological hypernetworks representing protein complexes and pathways of yeast and human and demonstrate a strong link between hypernetwork structure and the function of the proteins. Our novel methodology is able to mine the biological information hidden in the multi-scale architecture of molecular organisation. Furthermore, our analysis highlights the superiority, in terms of uncovering the underlying biology, of our multi-scale model when compared to the standard PPI networks. Additionally, we demonstrate that our new hypernetwork model, combined with our hypergraphlets, can be used for functional predictions.

Despite a simple functional prediction approach, we obtain promising results when using hypergraphlets on our new multi-scale model for functional predictions. It would be interesting to train an advanced machine learning model, such as random forrest, using HDVs as features in an effort to improve predictions. Finally, we have demonstrated that the union of networks capturing the multi-scale molecular organisation is strongly linked to the underlying biology of the molecules. It would be interesting to further investigate if different data integration methods could lead to even more biologically relevant models.

## Funding

This work was supported by UCL Computer Science departmental funds, the European Research Council (ERC) Starting Independent Researcher Grant 278212, the European Research Council (ERC) Consolidator Grant 770827, the Serbian Ministry of Education and Science Project III44006, the Slovenian Research Agency project J1-8155 and the awards to establish the Farr Institute of Health Informatics Research, London, from the Medical Research Council, Arthritis Research UK, British Heart Foundation, Cancer Research UK, Chief Scientist Office, Economic and Social Research Council, Engineering and Physical Sciences Research Council, National Institute for Health Research, National Institute for Social Care and Health Research, and Wellcome Trust (grant MR/K006584/1) and UK Medical Research Council (MC\_U12266B).

## References

- [1] A. Bellaachia and M. Al-Dhelaan. Random Walks in Hypergraph. *Proceedings of the 2013 International Conference on Applied Mathematics and Computational Methods*, pages 187–194, 2013.
- [2] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the royal statistical society. Series B (Methodological)*, pages 289–300, 1995.
- [3] C. Berge. *Graphs and Hypergraphs*, volume 6. Amsterdam: North-Holland publishing company, 1973.

- [4] D. Binns, E. Dimmer, R. Huntley, D. Barrell, C. O'donovan, and R. Apweiler. Quickgo: a web-based tool for gene ontology searching. *Bioinformatics*, 25(22):3045–3046, 2009.
- [5] J. A. Blake, K. R. Christie, M. E. Dolan, H. J. Drabkin, D. P. Hill, L. Ni, D. Sitnikov, S. Burgess, T. Buza, C. Gresham, F. McCarthy, L. Pillai, H. Wang, S. Carbon, H. Dietze, S. E. Lewis, C. J. Mungall, M. C. Munoz-Torres, M. Feuermann, P. Gaudet, S. Basu, R. L. Chisholm, R. J. Dodson, P. Fey, H. Mi, P. D. Thomas, A. Muruganujan, S. Poudel, J. C. Hu, S. A. Aleksander, B. K. McIntosh, D. P. Renfro, D. A. Siegele, H. Attrill, N. H. Brown, S. Tweedie, J. Lomax, D. Osumi-Sutherland, H. Parkinson, P. Roncaglia, R. C. Lovering, P. J. Talmud, S. E. Humphries, P. Denny, N. H. Campbell, R. E. Foulger, M. C. Chibucos, M. G. Giglio, H. Y. Chang, R. Finn, M. Fraser, A. Mitchell, G. Nuka, S. Pesseat, A. Sangrador, M. Scheremetjew, S. Y. Young, R. Stephan, M. A. Harris, S. G. Oliver, K. Rutherford, V. Wood, J. Bahler, A. Lock, P. J. Kersey, M. D. McDowall, D. M. Staines, M. Dwinell, M. Shimoyama, S. Laulederkind, G. T. Hayman, S. J. Wang, V. Petri, P. D'Eustachio, L. Matthews, R. Balakrishnan, G. Binkley, J. M. Cherry, M. C. Costanzo, J. Demeter, S. S. Dwight, S. R. Engel, B. C. Hitz, D. O. Inglis, P. Lloyd, S. R. Miyasato, K. Paskov, G. Roe, M. Simison, R. S. Nash, M. S. Skrzypek, S. Weng, E. D. Wong, T. Z. Berardini, D. Li, E. Huala, J. Argasinska, C. Arighi, A. Auchincloss, K. Axelsen, G. Argoud-Puy, A. Bateman, B. Bely, M. C. Blatter, C. Bonilla, L. Bougueleret, E. Boutet, L. Breuza, A. Bridge, R. Britto, C. Casals, E. Cibrian-Uhalte, E. Coudert, I. Cusin, P. Duek-Roggli, A. Estreicher, L. Famiglietti, P. Gane, P. Garmiri, A. Gos, N. Gruaz-Gumowski, E. Hatton-Ellis, U. Hinz, C. Hulo, R. Huntley, F. Jungo, G. Keller, K. Laiho, P. Lemercier, D. Lieberherr, A. Macdougall, M. Magrane, M. Martin, P. Masson, P. Mutowo, C. O'Donovan, I. Pedruzzi, K. Pichler, D. Poggioli, S. Poux, C. Rivoire, B. Roechert, T. Sawford, M. Schneider, A. Shypitsyna, A. Stutz, S. Sundaram, M. Tognolli, C. Wu, I. Xenarios, J. Chan, R. Kishore, P. W. Sternberg, K. Van Auken, H. M. Muller, J. Done, Y. Li, D. Howe, and M. Westerfeld. Gene ontology consortium: Going forward. *Nucleic Acids Research*, 43(D1):D1049–D1056, jan 2015.
- [6] J. A. Bondy and U. S. R. Murty. *Graph theory with applications*, volume 290. London: Macmillan, 1976.
- [7] N. Brockdorff, A. Ashworth, G. F. Kay, V. M. McCabe, D. P. Norris, P. J. Cooper, S. Swift, and S. Rastan. The product of the mouse xist gene is a 15 kb inactive x-specific transcript containing no conserved orf and located in the nucleus. *Cell*, 71(3):515–526, 1992.
- [8] A. Chatr-Aryamontri, R. Oughtred, L. Boucher, J. Rust, C. Chang, N. K. Kolas, L. O'Donnell, S. Oster, C. Theesfeld, A. Sellam, C. Stark, B. J. Breitkreutz, K. Dolinski, and M. Tyers. The BioGRID interaction database: 2017 update. *Nucleic Acids Research*, 45(D1):D369–D379, jan 2017.
- [9] S. A. Cook. The complexity of theorem-proving procedures. In *Proceedings of the Third Annual ACM Symposium on Theory of Computing*, pages 151–158. ACM, 1971.
- [10] D. Davis, Ö. N. Yaveroglu, N. Malod-Dognin, A. Stojmirovic, and N. Przulj. Topology-function conservation in protein-protein interaction networks. *Bioinformatics*, 31(10):1632–1639, 2015.
- [11] C. Díaz-Jullien, A. Pérez-Estévez, G. Covelo, and M. Freire. Prothymosin  $\alpha$  binds histones in vitro and shows activity in nucleosome assembly assay. *Biochimica et Biophysica Acta (BBA)-Protein Structure and Molecular Enzymology*, 1296(2):219–227, 1996.
- [12] E. Estrada and J. A. Rodríguez-Velázquez. Subgraph centrality and clustering in complex hyper-networks. *Physica A: Statistical Mechanics and its Applications*, 364:581–594, 2006.
- [13] A. Fabregat, K. Sidiropoulos, P. Garapati, M. Gillespie, K. Hausmann, R. Haw, B. Jassal, S. Jupe, F. Korninger, S. McKay, L. Matthews, B. May, M. Milacic, K. Rothfels, V. Shamovsky, M. Webber, J. Weiser, M. Williams, G. Wu, L. Stein, H. Hermjakob, and P. D'Eustachio. The reactome pathway knowledgebase. *Nucleic Acids Research*, 44(D1):D481–D487, jan 2016.
- [14] A. Feuerborn and P. R. Cook. Why the activity of a gene depends on its neighbors. *Trends in Genetics*, 31(9):483–490, 2015.

- [15] M. C. Gustin, J. Albertyn, M. Alexander, and K. Davenport. Map kinase pathways in the yeasts *Saccharomyces cerevisiae*. *Microbiology and Molecular Biology Reviews*, 62(4):1264–1300, 1998.
- [16] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor. Canonical Correlation Analysis: An Overview with Application to Learning Methods. *Neural Computation*, 16(12):2639–2664, dec 2004.
- [17] J. A. Hartigan and M. A. Wong. Algorithm as 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):100–108, 1979.
- [18] T. Ito, T. Chiba, R. Ozawa, M. Yoshida, M. Hattori, and Y. Sakaki. A comprehensive two-hybrid analysis to explore the yeast protein interactome. *Proceedings of the National Academy of Sciences*, 98(8):4569–4574, 2001.
- [19] P. Jaccard. The distribution of the flora in the alpine zone. *New phytologist*, 11(2):37–50, 1912.
- [20] J. Jokinen, E. Dadu, P. Nykvist, J. Käpylä, D. J. White, J. Ivaska, P. Vehviläinen, H. Reunanen, H. Larjava, L. Häkkinen, et al. Integrin-mediated cell adhesion to type I collagen fibrils. *Journal of Biological Chemistry*, 279(30):31956–31963, 2004.
- [21] S. Klamt, U. U. Haus, and F. Theis. Hypergraphs and cellular networks, may 2009.
- [22] O. Kuchaiev, T. Milenković, V. Memišević, W. Hayes, and N. Pržulj. Topological network alignment uncovers biological function and phylogeny. *Journal of the Royal Society Interface*, page rsif20100063, 2010.
- [23] V. Lacroix, L. Cottret, P. Thébaud, and M. F. Sagot. An introduction to metabolic networks and their structural analysis. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 5(4):594–617, oct 2008.
- [24] M. Latapy, C. Magnien, and N. D. Vecchio. Basic notions for the analysis of large two-mode networks. *Social Networks*, 30(1):31–48, 2008.
- [25] R. W. Li, C. Li, and T. T. Wang. Transcriptomic alterations in human prostate cancer cell line xenograft modulated by dietary phenethyl isothiocyanate. *Molecular Carcinogenesis*, 52(6):426–437, 2013.
- [26] Y.-f. Liu, A. Yang, W. Liu, C. Wang, M. Wang, L. Zhang, D. Wang, J.-f. Dong, and M. Li. Nme2 reduces proliferation, migration and invasion of gastric cancer cells to limit metastasis. *PloS one*, 10(2):e0115968, 2015.
- [27] J. Lugo-Martinez and P. Radivojac. Classification in biological networks with hypergraphlet kernels. *arXiv:1703.04823*, 2017.
- [28] H. D. Madhani and G. R. Fink. The control of filamentous differentiation and virulence in fungi. *Trends in Cell Biology*, 8(9):348–353, 1998.
- [29] N. Malod-Dognin and N. Pržulj. L-GRAAL: Lagrangian graphlet-based network aligner. *Bioinformatics*, 31(13):2182–2189, jul 2015.
- [30] O. Mason and M. Verwoerd. Graph theory and networks in biology. *IET Systems Biology*, 1(2):89–119, 2007.
- [31] T. Milenkovic and N. Pržulj. Uncovering Biological Network Function via Graphlet Degree Signatures. *Cancer Informatics*, pages 257–273, 2008.
- [32] N. Percy, N. Chuzhanova, and J. J. Crofts. Complexity and robustness in hypernetwork models of metabolism. *Journal of Theoretical Biology*, 406:99–104, 2016.
- [33] N. Percy, J. J. Crofts, and N. Chuzhanova. Hypergraph Models of Metabolism. *International Journal of Biological, Biomolecular, Agricultural, Food and Biotechnological Engineering*, 8(8):19–23, 2014.

- [34] M. Pelillo. A Game-Theoretic Approach to Hypergraph Clustering. *Advances in Neural Information Processing Systems*, 35(6):1312–1327, 2013.
- [35] N. Pržulj. Biological network comparison using graphlet degree distribution. In *Bioinformatics*, volume 23, pages e177–e183, jan 2007.
- [36] N. Pržulj, D. G. Corneil, and I. Jurisica. Modeling interactome: Scale-free or geometric? *Bioinformatics*, 20(18):3508–3515, dec 2004.
- [37] S. Pu, J. Wong, B. Turner, E. Cho, and S. J. Wodak. Up-to-date catalogues of yeast protein complexes. *Nucleic Acids Research*, 37(3):825–831, feb 2009.
- [38] K. M. Pumiglia and S. J. Decker. Cell cycle arrest mediated by the mek/mitogen-activated protein kinase pathway. *Proceedings of the National Academy of Sciences*, 94(2):448–452, 1997.
- [39] T. Rolland, M. Taşan, B. Charlotteaux, S. J. Pevzner, Q. Zhong, N. Sahni, S. Yi, I. Lemmens, C. Fontanillo, R. Mosca, et al. A proteome-scale map of the human interactome network. *Cell*, 159(5):1212–1226, 2014.
- [40] A. Ruepp, B. Brauner, I. Dunger-Kaltenbach, G. Frishman, C. Montrone, M. Stransky, B. Waegle, T. Schmidt, O. N. Doudieu, and V. Stümpflen. CORUM: the comprehensive resource of mammalian protein complexes. *Nucleic Acids Research*, 36(Database):D646, jan 2007.
- [41] A. Ruepp, B. Waegle, M. Lechner, B. Brauner, I. Dunger-Kaltenbach, G. Fobo, G. Frishman, C. Montrone, and H. W. Mewes. CORUM: The comprehensive resource of mammalian protein complexes-2009. *Nucleic Acids Research*, 38(SUPPL.1):D497–D501, jan 2009.
- [42] U. Stelzl, U. Worm, M. Lalowski, C. Haenig, F. H. Brembeck, H. Goehler, M. Stroedicke, M. Zenkner, A. Schoenherr, S. Koeppen, et al. A human protein-protein interaction network: a resource for annotating the proteome. *Cell*, 122(6):957–968, 2005.
- [43] S. Testaz and J.-L. Duband. Central role of the  $\alpha4\beta1$  integrin in the coordination of avian truncal neural crest cell adhesion, migration, and survival. *Developmental Dynamics*, 222(2):127–140, 2001.
- [44] Z. Tian, T. Hwang, and R. Kuang. A hypergraph-based learning algorithm for classifying gene expression and arrayCGH data with prior knowledge. *Bioinformatics*, 25(21):2831–2838, 2009.
- [45] P. Uetz, L. Giot, G. Cagney, T. A. Mansfield, R. S. Judson, J. R. Knight, D. Lockshon, V. Narayan, M. Srinivasan, P. Pochart, et al. A comprehensive analysis of protein–protein interactions in *saccharomyces cerevisiae*. *Nature*, 403(6770):623, 2000.
- [46] A. Vazquez, A. Flammini, A. Maritan, and A. Vespignani. Global protein function prediction from protein-protein interaction networks. *Nature Biotechnology*, 21(6):697, 2003.
- [47] N. X. Vinh, J. Epps, and J. Bailey. Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance. *Journal of Machine Learning Research*, 11:2837–2854, 2010.
- [48] N. Wong and X. Wang. mirdb: an online resource for microRNA target prediction and functional annotations. *Nucleic Acids Research*, 43(D1):D146–D152, 2014.
- [49] Ö. N. Yaveroglu, N. Malod-Dognin, D. Davis, Z. Levnajic, V. Janjic, R. Karapandza, A. Stojmirovic, and N. Pržulj. Revealing the hidden language of complex networks. *Scientific Reports*, 4:4547, 2014.

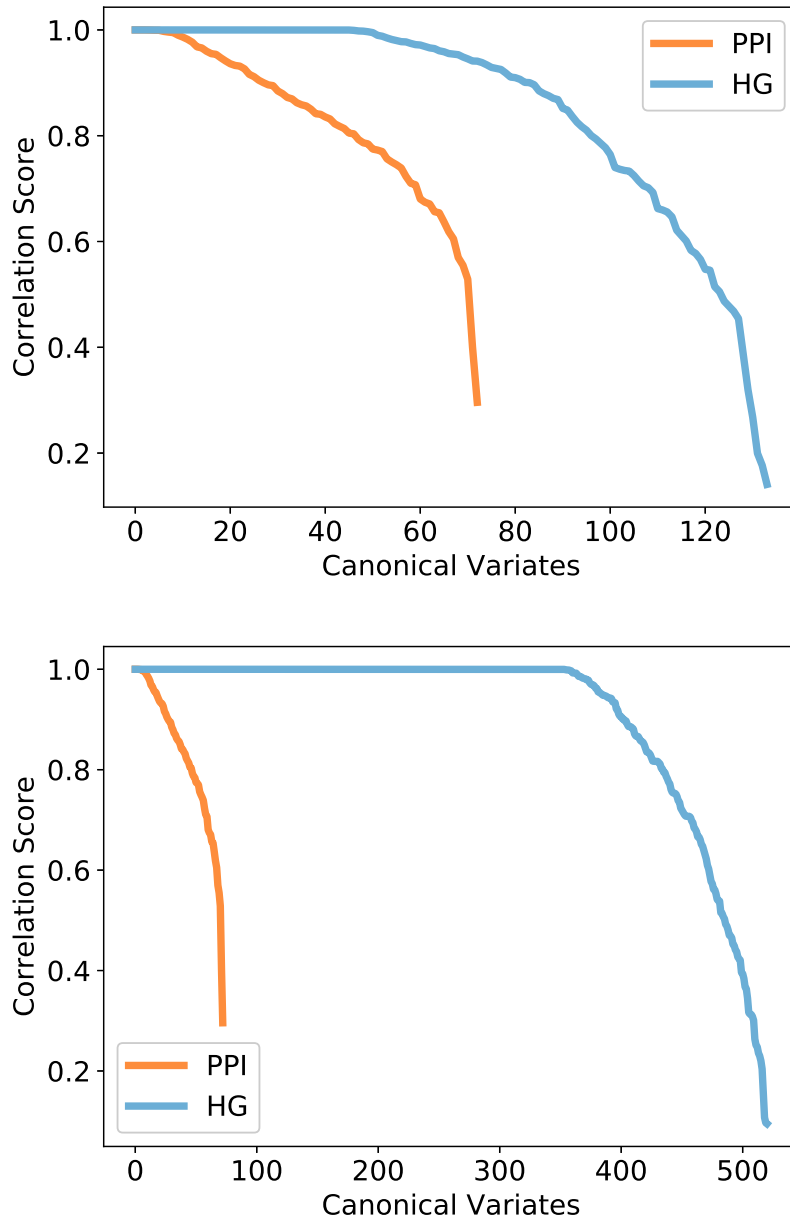
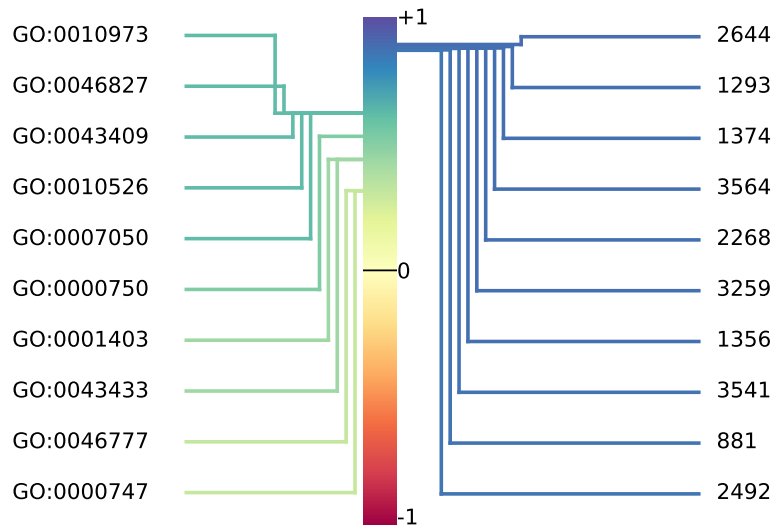


Figure 6: Canonical correlation score distribution for yeast-complexes (top) and yeast-pathways (bottom). The canonical variates represented are all statistically significant ( $p\text{-value} \leq 5\%$ ) and are sorted by correlation score. The colors represent the model and the topological signatures from which the canonical variates are obtained: HG in blue and PPI in orange.



- GO:0010973: positive regulation of barrier septum assembly
- GO:0046827: positive regulation of protein export from nucleus
- GO:0043409: negative regulation of MAPK cascade
- GO:0010526: negative regulation of transposition, RNA-mediated
- GO:0046777: protein autophosphorylation
- GO:0007050: cell cycle arrest
- GO:0000750: pheromone-dependent signal transduction involved in conjugation with cellular fusion
- GO:0001403: invasive growth in response to glucose limitation
- GO:0043433: negative regulation of sequence-specific DNA binding transcription factor activity
- GO:0000747: conjugation with cellular fusion

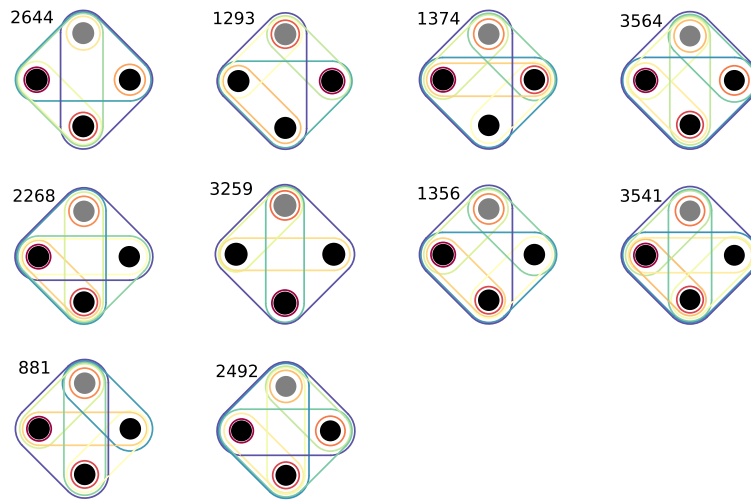


Figure 7: The most significant CCA variate between HDVs of the proteins of yeast-pathways and their GO-BP annotations. The correlation score between the linear combination of annotations and the linear combination of hypergraphlet orbits is 1. The annotations (orbits) illustrated above correspond to the 10 that have the highest Pearson's correlation scores with respect to the linear combinations of annotations (orbits). Each GO term in blue font is annotating at least one protein conjointly with at least one other annotation that is also denoted in blue font, according to QuickGO ontology search engine [4].

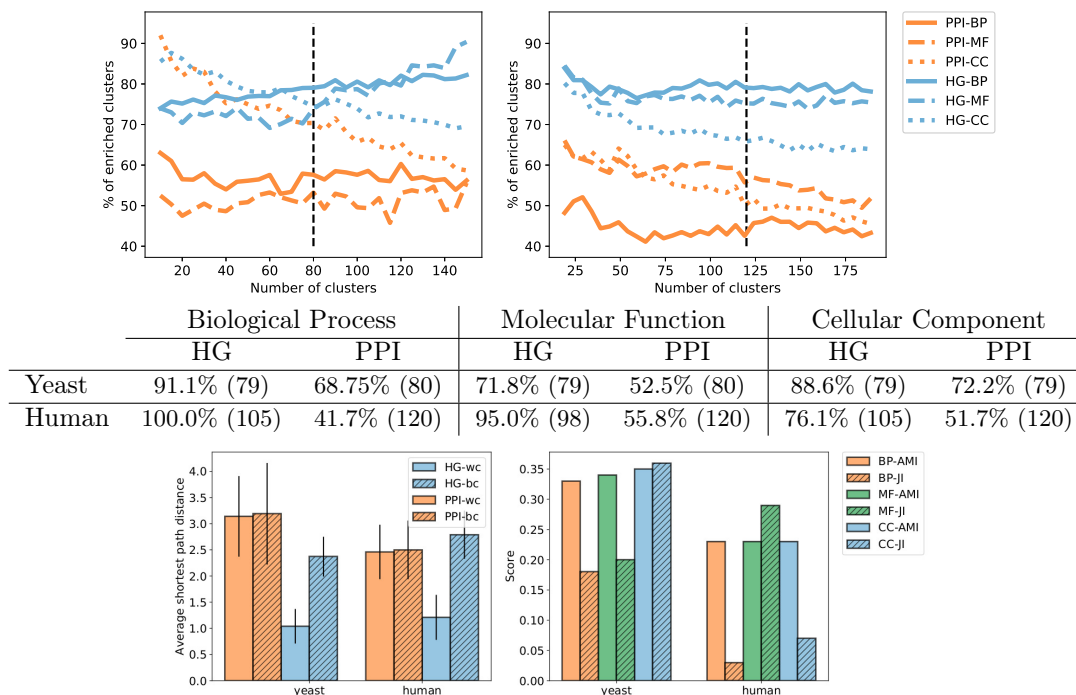


Figure 8: The top panels give the average percentages of clusters enriched with respect to the total number of clusters for yeast (left) and human (right), the standard deviation is not represented to avoid overcrowding the panels. The colors represent the models from which the clustering is obtained: HG in blue and PPI in orange. The type of line represents the category of GO annotations: BP are full lines, MF are dashed lines, and CC are dotted line. The black vertical lines denote the number of clusters selected from the set of NMI and SSE curves according to the procedure described in Section 3.4.1. The middle table presents the maximum enrichment measured across clusterings obtained with the “optimal” number of clusters (denoted by the black vertical lines in the top panels). The number in parenthesis is the number of non-empty clusters. All enrichments are significant. The bottom left panel gives, for each type of model (HG in blue and PPI in orange), the average of the shortest path lengths within the clusters (wc) and between clusters (bc) of the best clustering obtained for GO-BP annotations. The results are similar for other GO categories and are not presented here due to space limitations. The bottom right panel represents the results of the comparison of the obtained clusterings. We use the HG clustering as baseline and compute the Adjusted Mutual Information (AMI) between the clusterings and the Jaccard Index (JI) between the sets of enriched GO terms.



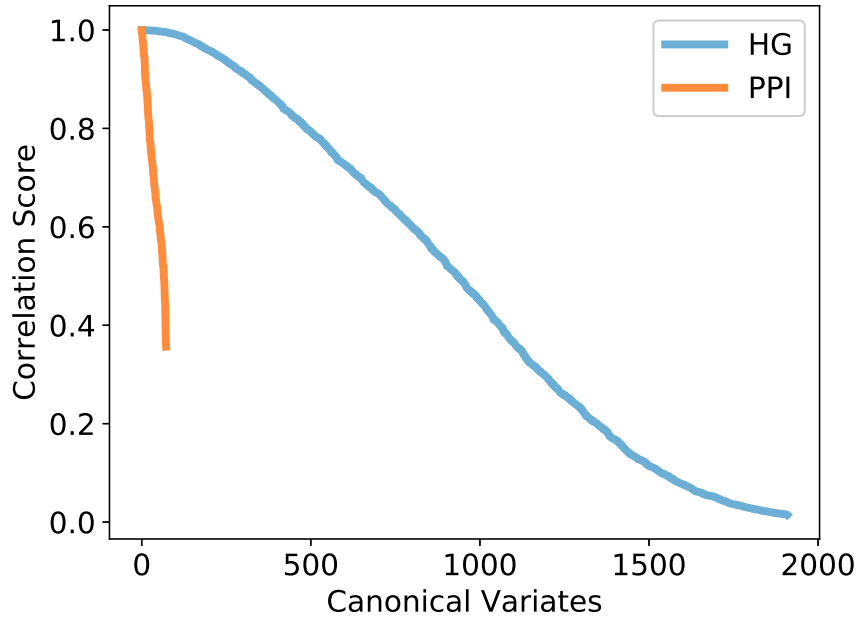


Figure 9: Canonical correlation score distribution for the human hypernetwork. The canonical variates represented are all statistically significant ( $p\text{-value} \leq 5\%$ ) and are sorted by correlation score. The colors represent the model and the topological signatures from which the canonical variates are obtained: HG in blue and PPI in orange.

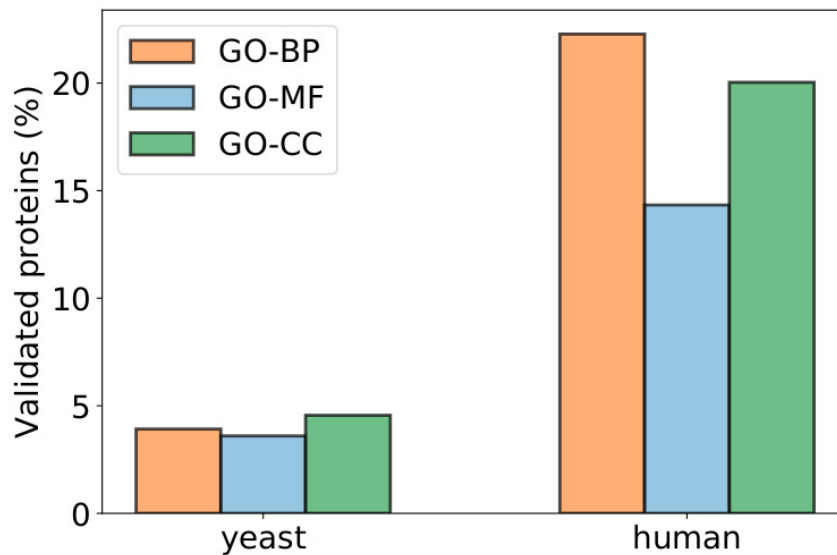


Figure 10: Percentages of proteins that have at least one of the enriched terms of their clusters in their set of predicted GO annotations (obtained from the GO database [5]). The values correspond to the number of such proteins out of the number of proteins that have at least one putative annotation in the GO database and are not experimentally annotated with any of the enriched terms of their clusters.