

Protein complex similarity based on Weisfeiler-Lehman labeling

Bianca K. Stöcker^{1,2}, Till Schäfer⁴, Petra Mutzel⁴, Johannes Köster^{1,2,3}, Nils Kriege⁴, and Sven Rahmann^{1,4}

¹Genome Informatics, Institute of Human Genetics, University Hospital Essen, University of Duisburg-Essen, 45147 Essen, Germany

²Algorithms for Reproducible Bioinformatics, Genome Informatics, Institute of Human Genetics, University Hospital Essen, University of Duisburg-Essen, 45147 Essen, Germany

³Medical Oncology, Dana-Farber Cancer Institute, Harvard Medical School, Boston MA 02215, USA

⁴Department of Computer Science, TU Dortmund University, 44221 Dortmund, Germany

Corresponding author:

Sven Rahmann¹

Email address: Sven.Rahmann@uni-due.de

ABSTRACT

Being able to quantify the similarity between two protein complexes is essential for numerous applications. Prominent examples are database searches for known complexes with a given query complex, comparison of the output of different protein complex prediction algorithms, or summarizing and clustering protein complexes, e.g., for visualization. While the corresponding problems have received much attention on single proteins and protein families, the question about how to model and compute similarity between protein complexes has not yet been systematically studied. Because protein complexes can be naturally modeled as graphs, in principle general graph similarity measures may be used, but these are often computationally hard to obtain and do not take typical properties of protein complexes into account. Here we propose a parametric family of similarity measures based on Weisfeiler-Lehman labeling. We evaluate it on simulated complexes of the extended human integrin adhesome network. Because the connectivity (graph topology) of real complexes is often unknown and hard to obtain experimentally, we use both known protein-protein interaction networks and known interdependencies (constraints) between interactions to simulate more realistic complexes than from interaction networks alone. We empirically show that the defined family of similarity measures is in good agreement with edit similarity, a similarity measure derived from graph edit distance, but can be much more efficiently computed. It can therefore be used in large-scale studies and simulations and serve as a basis for further refinements of modeling protein complex similarity.

INTRODUCTION

Proteins fulfill manifold tasks in living cells, but they rarely act alone. Indeed, most cellular functions are enabled only when proteins physically interact with other proteins, forming protein complexes. DNA transcription is a typical example, where RNA polymerase II, general transcription factors, cell type specific transcription regulators and mediator proteins interact.

Understanding protein complex formation and function is one of the big challenges of cell biology, approached by both experimental techniques and computational modeling. While the constituent protein sequences can be obtained from the genome (even that can be challenging), the computational prediction of real protein complexes from protein interaction networks appears to be much more difficult as evidenced by the recent literature on the topic; see Bhowmick and Seah (2016) for a survey, or Srihari et al. (2017) for a textbook introduction. Fortunately, new experimental technologies are about to enhance our understanding of complexes significantly in the near future, e.g. high-resolution protein-protein docking (Park et al., 2015; Vakser, 2014; Kozakov et al., 2017; Wass et al., 2011). Large scale generation of libraries of cell

47 lines having two or more endogenously tagged fluorescent proteins (Boutros et al., 2015) and recent
48 high-throughput and multiplexed implementations of fluorescence correlation spectroscopy allow us to
49 systematically measure endogenous concentrations, binding constants and high-order complexes in such
50 libraries of cell lines (Hwang et al., 2006; Wobma et al., 2012; Grecco et al., 2016; Wachsmuth et al.,
51 2015).

52 When studying biological entities such as protein sequences or protein complexes, a fundamental
53 task is to define a measure of similarity between two such entities. For protein sequences, there is a
54 well-established theory based on scoring matrices and alignment scores (Pearson, 2013). For protein
55 complexes, it appears that no systematic effort to quantify similarity has been made yet. The purpose of
56 the present article is therefore to discuss the different options to define a similarity measure on protein
57 complexes and to propose a reasonable and computationally tractable definition of protein complex
58 similarity.

59 Establishing a similarity measure is not only important fundamentally, but there are many immediate
60 applications, of which we mention the following three examples.

61 **Database search:** In the *database search problem* we are given a query complex and a large collection
62 (database) of complexes, and the task is to find the complexes in the database that are most similar
63 to the query. Obviously, a meaningfully defined similarity measure is essential for this task.

64 **Comparing predictions:** Several complex prediction methods predict putative complexes by locating
65 dense regions in a protein interaction network (Drew et al., 2017; Hernandez et al., 2017; Ma and
66 Gao, 2012; Pellegrini et al., 2016), and for comparing complexes predicted by different algorithms,
67 it is of interest to compute a maximum-weight matching between the output of two algorithms,
68 where the weighting is given by a similarity function.

69 **Summarizing and clustering:** When stochastically simulating complex formation based on available
70 knowledge such as possible interactions and interaction constraints, it is helpful to aggregate the
71 simulation output to focus on frequently seen or typical complexes, ignoring small differences.
72 Aggregation or clustering by similarity thereby reduces data size and complexity. Such a task first
73 and foremost requires a way to quantify the similarity between two protein complexes.

74 When there are many (say, tens of thousands of) different complexes subject to pairwise comparison, a
75 similarity measure must be efficiently computable to be of practical interest.

76 **Models for protein complexes.** We first need to discuss models for protein complexes on different
77 levels of detail, namely the *set*, *multiset*, and *graph* models.

78 While intuition suggests that protein complexes can be naturally described as graphs with proteins
79 as vertices and physical interactions as edges, there are in fact different ways to formally describe a
80 protein complex. In the following, we briefly mention the most prominent ones with their advantages
81 and disadvantages. We start with a given set P of all proteins of an organism, the building blocks of the
82 complexes.

83 **Set:** In its most simple form, a protein complex can be defined as a set (in the mathematical sense,
84 i.e., without multiplicities) of proteins, i.e., as a subset of P . We use the standard notation
85 $\{p_1, p_2, \dots, p_n\}$ for sets. Sets neither capture the multiplicities nor the nature of the physical
86 interactions between the constituent proteins of a complex. However, some experimental techniques
87 (e.g. immunoprecipitation with mass spectrometry) only give such set-type information, and several
88 existing databases only provide this type of information, e.g. the CORUM database provided by the
89 Munich Information center for Protein Sequences MIPS (Ruepp et al., 2010).

90 **Multiset:** Formally, a multiset is a function $C : P \rightarrow \mathbb{N}_0$ that assigns a multiplicity to each protein $p \in P$
91 with $C(p) = 0$ for proteins p that are not part of the complex. We also use the multiset notation
92 $C = \{\{p_1, p_1, p_2\}\}$ to express that $C(p_1) = 2$, $C(p_2) = 1$ and $C(p) = 0$ for all other $p \in P$. Defining
93 a protein complex as a multiset of proteins gives a more accurate representation of the complex, but
94 still does not consider the interaction topology.

95 **Graph:** To add more information, we can define a protein complex as an undirected graph $C = (V, E, \ell)$
96 with labeled vertices V , such that each vertex $v \in V$ represents a protein and hence has a label

97 $\ell(v) \in P$, each edge $e \in E \subseteq V \times V$ represents a physical interaction between the corresponding
 98 proteins, such that E is symmetric and C is connected. The graph description provides the interaction
 99 topology. We call this representation a *protein complex graph* and define its *size* as $|C| := |V| + |E|$.
 100 (This representation could be further refined by considering the different domains of each protein
 101 and specifying precisely which domains interact.)

102 For the set and multiset models, a similarity measure is readily given by the *Jaccard similarity* (see
 103 Methods). For graphs, the *graph edit distance* has been proposed for pattern recognition tasks more than
 104 30 years ago (Sanfeliu and Fu, 1983). A graph edit distance between graphs C and C' measures the total
 105 costs of the edit operations required to transform C into C' . Defining similarity via graph edit operations
 106 appears intuitive, but has computational disadvantages, as edit distance computation on graphs is hard in
 107 general. More specifically, the graph edit distance generalizes the classical maximum common subgraph
 108 problem (Bunke, 1997), which is NP-complete (Garey and Johnson, 1979) and hard to approximate with
 109 given guarantees (Kann, 1992). Recently, a binary linear programming formulation for computing the
 110 graph edit distance has been proposed (Lerouge et al., 2017), which allows to compare graphs of moderate
 111 size using state-of-the-art general purpose solvers. However, when we want to compare many complexes,
 112 evaluating the edit distance between all pairs becomes infeasible in practice.

113 In this article, we therefore propose an efficient alternative: We define a family of similarity measures
 114 on graphs by resorting to the (efficiently computable) Jaccard similarity, while still taking the graph
 115 structure into account. This is achieved by so-called Weisfeiler-Lehman labeling of the vertices (Weisfeiler
 116 and Lehman, 1968), propagating vertex labels between neighbors. This approach is different from recent
 117 work that approximates and bounds the graph edit distance (Riesen et al., 2015) and has the advantage of
 118 scaling better to large-scale studies.

119 The remainder of the article is structured as follows. In the Methods section, we define a parametric
 120 family of similarity measures based on Weisfeiler-Lehman labeling and the precise definition of graph
 121 edit similarity we compare against. In the Results section, we describe how we obtain pairs of protein
 122 complexes, for which we compare Weisfeiler-Lehman similarity and edit similarity. The simulated protein
 123 complexes take known protein interaction networks and additionally constraints between interactions into
 124 account, and therefore should represent more realistic complexes than arbitrary connected subgraphs of
 125 protein interaction networks. Finally, we discuss limitations and possible extensions of this work.

126 METHODS

127 Our goal is to define a similarity measure between protein complexes that captures not only the (multisets
 128 of the) constituent proteins, but also the interaction topology (graph structure). Similarities derived from
 129 graph edit distance offer this property, but as mentioned above, they are hard to compute. Therefore,
 130 we introduce a parameterized family of similarity measures on protein complexes, which are based on
 131 multiset comparisons of vertex labels in the complex graph and take the local neighborhood of each
 132 protein into account by using Weisfeiler-Lehman labels.

133 Jaccard similarity of sets and multisets

134 To compare sets or multisets, Jaccard similarity coefficients are an established measure.

Let $M \subseteq U$ and $M' \subseteq U$ be two subsets of a common universe U . Then the *Jaccard similarity* between M and M' is defined as

$$J_{\text{set}}(M, M') := \frac{|M \cap M'|}{|M \cup M'|} \in [0, 1]. \quad (1)$$

This definition is extended to multisets as follows. Recall that multisets M and M' are functions $U \rightarrow \mathbb{N}_0$, assigning multiplicities $M(o)$ and $M'(o)$ to each object $o \in U$. (The set definition can be seen as the special case where the value set is only $\{0, 1\}$ instead of \mathbb{N}_0 .) Then the *Jaccard similarity* between M and M' is defined as

$$J_{\text{multiset}}(M, M') := \frac{\sum_{o \in U} \min\{M(o), M'(o)\}}{\sum_{o \in U} \max\{M(o), M'(o)\}} \in [0, 1]. \quad (2)$$

135 **A parametric family of protein complex similarity measures**

136 Instead of comparing the protein complexes directly by their graph topology and labeling, we extract and
 137 compare multisets of features of the protein complexes. Weisfeiler and Lehman (1968) developed an
 138 iterative label refinement procedure to derive a canonical graph representation for graph isomorphism
 139 testing. The same procedure is often used to define graph similarities or graph kernels (Shervashidze et al.,
 140 2011). For the latter purpose, the vertex labels of each Weisfeiler-Lehman iteration are used as features of
 141 the graphs. Initially, the feature multiset of a graph consists of the union of all vertex labels. After the
 142 initialization, the vertex labels are iteratively augmented by the labels of the neighboring vertices from the
 143 previous iteration, thereby encoding the (local) graph structure in the vertex labels. Let us now formally
 144 define the process.

Definition 1 (Weisfeiler-Lehman labeling of iteration i for a protein complex graph). Let $C = (V, E, \ell_0)$ be a protein complex graph with label function $\ell_0 : V \rightarrow L_0 := P$. Furthermore, let $N(v) := \{u \mid \{v, u\} \in E\}$ denote the neighbors of vertex $v \in V$. Then, the Weisfeiler-Lehman labeling of iteration i is defined as a re-labeling of the protein complex graph: It replaces the labeling function $\ell_0 : V \rightarrow L_0$ with a labeling function $\ell_i : V \rightarrow L_i$. The value of ℓ_i for a vertex $v \in V$ is recursively defined as

$$\ell_i(v) := (\ell_{i-1}(v), \{\{\ell_{i-1}(u) \mid u \in N(v)\}\}). \quad (3)$$

145 Note that the second component of the new label is a multiset.

146 To avoid that the length of labels increases in each iteration, label compression is performed after each
 147 step in practice. This is achieved by a one-to-one mapping of the labels $\{\ell_i(v) \mid v \in V\}$ to integer labels.

148 Given the Weisfeiler-Lehman labeling function of a protein complex graph for some iteration i , we
 149 can now define the multiset of Weisfeiler-Lehman features for iteration i .

150 **Definition 2** (Weisfeiler-Lehman feature set of iteration i for a protein complex graph). Let $C = (V, E, \ell_0)$
 151 be a protein complex graph with label function $\ell_0 : V \rightarrow L_0 = P$. Then, the Weisfeiler-Lehman features of
 152 iteration i are defined as multiset $WL_i(C) = \{\{\ell_i(v) \mid v \in V\}\}$. Note that $WL_0(C)$ corresponds to the initial
 153 multiset of protein names.

To compare two complexes C and C' , we compare the iteration sequences of Weisfeiler-Lehman features $(WL_i(C))_{i \geq 0}$ and $(WL_i(C'))_{i \geq 0}$, by computing a convex combination of the Jaccard similarities for each iteration. Let $w = (w_i)_{i \geq 0}$ be a weight sequence with $w_i \geq 0$ for all $i \geq 0$ and $\sum_{i \geq 0} w_i = 1$. For w as just defined, let

$$S_w(C, C') := \sum_{i \geq 0} w_i \cdot J_{\text{multiset}}(WL_i(C), WL_i(C')), \quad (4)$$

154 where J_{multiset} is given by Eq. (2). This defines a family of similarity measures between complexes with
 155 values in $[0, 1]$, parameterized by a convex combination $w = (w_0, w_1, \dots)$.

156 It is easy to see that, as long as $w_0 > 0$, we have $S_w(C, C') = 0$ if and only if the protein sets of C and C'
 157 are disjoint. If $S_w(C, C') < 1$, the protein complex graphs are not isomorphic. However, $S_w(C, C') = 1$
 158 does not necessarily imply that C and C' are isomorphic even if $w_i > 0$ for all i : There exist examples
 159 of non-isomorphic graphs G, G' with $WL_i(G) = WL_i(G')$ for all $i \geq 0$. (As a simple example, take G
 160 to be a cycle of six vertices, and G' to be two cycles of three vertices, all with the same label.) On the
 161 other hand, there exist classes of graphs, such as the so-called CR-graphs, for which the implication
 162 “ $S_w(C, C') = 1 \Rightarrow C, C'$ are isomorphic” is true if $w_i > 0$ for all i (Arvind et al., 2015).

163 In practice, we may assume that most protein complexes are non-adversarial graphs with sufficiently
 164 simple structure such that their Weisfeiler-Lehman features are appropriate to characterize their similarity.
 165 In fact, we put forward the hypothesis that using a single iteration is frequently sufficient for practical
 166 purposes, and we set $w_i := 0$ for $i \geq 2$ in our computational experiments (see Results) and only have
 167 a single free parameter $w_0 \in [0, 1]$ that defines $w_1 := 1 - w_0$. In the following, we write ω for w_0 . In
 168 this case, S_ω is efficiently computable: A proof of the following lemma can be found in the work of
 169 Shervashidze et al. (2011).

Lemma 3. For $\omega \in [0, 1]$, each of the one-parameter similarity measures

$$S_\omega(C, C') := \omega \cdot J_{\text{multiset}}(WL_0(C), WL_0(C')) + (1 - \omega) \cdot J_{\text{multiset}}(WL_1(C), WL_1(C'))$$

170 can be computed in $O(|C| + |C'|)$ time, where $|C| = |V| + |E|$.

171 **A similarity measure based on the graph edit distance**

172 To compare the family of Weisfeiler-Lehman multiset-based similarity measures defined above with
 173 graph edit distance, we start with a formal definition of the edit-based similarity. We allow the following
 174 elementary operations to edit a graph: vertex deletion, vertex insertion, vertex relabeling, edge deletion,
 175 and edge insertion. A sequence (o_1, \dots, o_k) of such edit operations that transforms a graph G into another
 176 graph H is called an *edit path* from G to H . Each operation o is assigned a cost $c(o)$, which is zero for
 177 substituting vertices and edges with the same label. We use a cost of 1 for all operations except vertex
 178 relabeling which has a cost of 2, corresponding to one deletion and one insertion (leaving the edges in
 179 place). Note that deleting or inserting a vertex of degree k otherwise has cost $k + 1$ for deleting k edges
 180 and the vertex itself. We denote the set of all possible edit paths from G to H by $\Upsilon(G, H)$.

Definition 4. Let G and H be labeled graphs. The *graph edit distance* from G to H is defined by

$$d(G, H) = \min \left\{ \sum_{i=1}^k c(o_i) \mid (o_1, \dots, o_k) \in \Upsilon(G, H) \right\}. \quad (5)$$

Intuitively, the graph edit distance preserves a subgraph G' of G that is also contained in H using zero-cost substitutions, deletes the vertices and edges in G that are not in G' and then inserts vertices and edges to obtain an isomorphic copy of H . Therefore all non-zero costs can be attributed to the elements which are in one of the graphs, but not in their common subgraph. In this sense the graph edit distance is similar to the symmetric difference of two sets. This observation motivates the following normalized similarity measure derived from the graph edit distance. We define the *graph edit similarity* as

$$J_{\text{graph}}(G, H) := \frac{|G| + |H| - d(G, H)}{|G| + |H| + d(G, H)} \in [0, 1], \quad (6)$$

181 where $|G| := |V(G)| + |E(G)|$. Note that the graph edit distance between G and H is at most $|G| + |H|$,
 182 which is achieved by deleting all vertices and edges of G and inserting all vertices and edges of H . In
 183 this case the graph edit similarity is zero. Similarly, $J_{\text{graph}}(G, H) = 1$ if and only if $d(G, H) = 0$. In this
 184 respect the similarity measure resembles the Jaccard similarity. In fact, we can show a deeper relation to
 185 the multiset Jaccard similarity.

186 **Lemma 5.** For two protein complexes, let C, D denote their protein multisets and G, H their pro-
 187 tein complex graphs. For the edge-free graphs $G' = (V(G), \emptyset)$ and $H' = (V(H), \emptyset)$ it holds that
 188 $J_{\text{graph}}(G', H') = J_{\text{multiset}}(C, D)$.

Proof. An optimal graph edit path is obtained as follows: We substitute the vertices with common labels free of cost, which are $Z = \sum_{p \in P} \min\{C(p), D(p)\}$ in total. We delete the remaining $|G'| - Z$ vertices in G' and insert $|H'| - Z$ vertices to obtain an isomorphic copy of H' at a total cost of $|G'| + |H'| - 2Z = d(G, H)$. Instead we may also substitute up to $||G'| - |H' ||$ vertices, each at cost two, which results in the same total cost. Using the fact that $|G'| = \sum_{p \in P} C(p)$ and $|H'| = \sum_{p \in P} D(p)$, we obtain the result by calculating

$$\begin{aligned} J_{\text{graph}}(G', H') &= \frac{|G'| + |H'| - d(G', H')}{|G'| + |H'| + d(G', H')} = \frac{Z}{|G'| + |H'| - Z} = \frac{Z}{\sum_{p \in P} C(p) + \sum_{p \in P} D(p) - Z} \\ &= \frac{\sum_{p \in P} \min\{C(p), D(p)\}}{\sum_{p \in P} C(p) + D(p) - \min\{C(p), D(p)\}} = \frac{\sum_{p \in P} \min\{C(p), D(p)\}}{\sum_{p \in P} \max\{C(p), D(p)\}} = J_{\text{multiset}}(C, D). \end{aligned}$$

189

□

190 Lemma 5 shows that the graph edit similarity can indeed be seen as a natural extension of the multiset
 191 Jaccard similarity to graph structured data.

192 For our computations, we used a recent binary linear programming formulation to compute the graph
 193 edit distance exactly (Lerouge et al., 2017). The approach was implemented in Java, and all instances
 194 were solved using an academic license of Gurobi 7.5.2 on Linux x86-64.

195 RESULTS

196 Hypothesis

197 We hypothesize that the Weisfeiler-Lehman based family of similarity measures S_ω defined in Eq. (4)
 198 approximates well the edit distance based similarity defined in Eq. (6) for typical protein complexes. The
 199 similarity measures S_ω have the advantage that they can be efficiently computed (see Lemma 3).

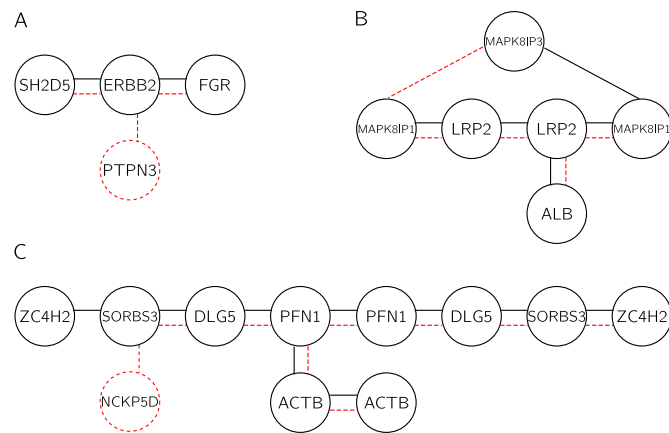


Figure 1. Three exemplary pairs of protein complexes: Each labeled node is a protein instance, each edge represents a protein interaction, and solid black vs. dashed red edges distinguish between the two complexes. **A:** Edit similarity 0.714; WL similarity in $[0.4, 0.75]$ depending on weight ω . **B:** Edit similarity 0.838; WL similarity 1.0 (independent of ω). **C:** Edit similarity 0.9; WL similarity in $[0.667, 0.818]$ depending on ω .

200 Data generation

201 As mentioned in the Introduction, obtaining real protein complex graphs is difficult at the moment,
 202 because experimental techniques that resolve the (graph) topology of the complexes are only being
 203 developed. Therefore we resort to the simulation of complexes, based on two types of knowledge:
 204 possible protein-protein interactions, formalized by a *protein interaction network*, and *constraints between*
 205 *protein interactions*.

206 Formally, a protein interaction network is an undirected graph $N = (P, I)$, where P is the set of protein
 207 types of a cell (or an organism), and $I \subset P \times P$ indicates the pairs of protein types that may potentially
 208 physically interact. Since N describes the entirety of possible interactions, any protein complex can be
 209 seen as a connected subgraph of N .

210 It is important to realize that protein interactions are not independent of each other, but interdependent.
 211 Those interaction dependencies are generated by two major mechanisms. On the one hand there is
 212 allosteric regulation, in which the capability of a protein to bind other proteins is affected by a confor-
 213 mational change upon one interaction (Laskowski et al., 2009). The other key mechanism is steric hindrance
 214 that prevents proteins from binding simultaneously to too close or identical protein domains leading to
 215 mutual exclusiveness of interactions (Sánchez Claros and Tramontano, 2012). The dependencies between
 216 interactions constrain the set of possible protein complexes and their assembly paths. Therefore, for
 217 understanding the design and function of intracellular protein networks it is important to consider the
 218 dependencies between protein interactions. One possible model for this are *constrained protein interaction*
 219 *networks*, where the protein interaction network is enhanced by the interaction dependencies (*constraints*)
 220 modeled as propositional logic formulas (Stöcker et al., 2017).

221 With constrained protein interaction networks, we can stochastically simulate complex formation
 222 based on the available knowledge and obtain a detailed interaction topology (which proteins physically
 223 interact) for each complex.

224 To evaluate the Weisfeiler-Lehman based similarity (“WL similarity”) against the edit distance based
 225 similarity (“edit similarity”), we computed both similarity measures on selected pairs of 100 000 simulated
 226 protein complexes from the extended human adhesome network as presented by Stöcker et al. (2017).
 227 Since edit similarity computations are computationally costly, we only computed the edit similarity on
 228 500 000 candidate pairs from these simulated complexes. These candidate pairs were generated for all
 229 pairs of complexes that have at most 20 proteins (larger complexes are so rare that high similarities
 230 are unlikely), that have a size difference of protein multisets of at most 10, and that share at least one
 231 protein. The candidate pairs were sorted descendingly after the number of shared proteins and then
 232 the edit distance based similarity was computed on the first 500 000 candidate pairs. The resulting edit
 233 similarity values were classified into bins of width 0.1. Because most pairs of complexes share a small

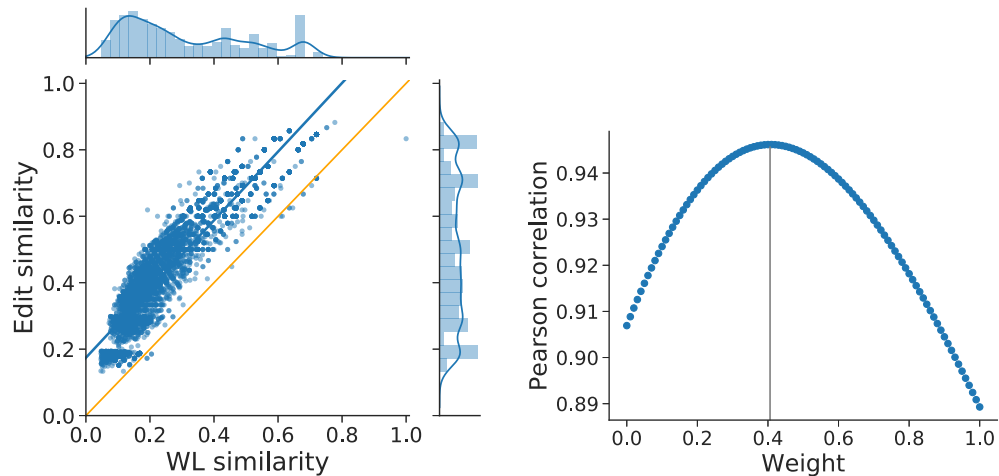


Figure 2. **Left:** Scatterplot comparison of edit similarity and WL similarity for weight $\omega = 0.41$, including marginal distributions and least-squares regression line. **Right:** Pearson correlation between edit similarity and WL similarity as a function of weight ω . The maximum correlation (0.946) occurs for $\omega = 0.41$, shown on the left side.

234 number of proteins, we find many pairs with small edit similarity (but none in the range $[0.0, 0.1[$ because
 235 we required one common protein) and comparatively few pairs with edit similarity above 0.5. To achieve
 236 a uniform distribution among bins for the comparison, we randomly selected 1000 pairs from each bin,
 237 excluding the bin $[0.9, 1.0[$ which contained a single pair. This yielded 8000 pairs of complexes from
 238 8 bins.

239 Similarity comparison

240 We first consider three exemplary pairs (Figure 1 A–C) with edit similarities of approximately 0.7, 0.8
 241 and 0.9, respectively, the latter being the most similar observed pair.

242 In example A, an additional protein (PTPN3) is added to an existing complex, a linear chain of 3
 243 proteins. The edit similarity is $10/14 = 0.714$, the WL similarity is between 0.75 for $\omega = 1$ and 0.4
 244 for $\omega = 0$. Because the edit similarity is between the extreme WL similarities, there exists a unique
 245 weight $\omega^* \approx 0.898$, for which WL and edit similarities agree for this particular complex pair. Example B
 246 is an noteworthy case, because the WL similarity is 1.0, independent of ω , because the vertex labels are
 247 identical even after the first Weisfeiler-Lehman iteration. (Further iterations would show a difference.)
 248 The edit similarity is $20/24 = 0.83$, which is obtained by attaching ALB to the other LRP2 protein. In
 249 example C, one protein is replaced by another one in a fairly large complex. The edit similarity (0.905)
 250 is relatively high and outside the WL similarity range between 0.667 for $\omega = 0$ and 0.818 for $\omega = 1$.

251 Because most protein complexes are small and do not exhibit properties of examples B or C, the
 252 overall agreement between WL similarity and edit similarity is high: For each of the selected com-
 253 plex pairs, we computed the exact edit similarity and the WL similarity for each weight $\omega \in W :=$
 254 $\{0.0, 0.01, 0.02, \dots, 1.0\}$. Let e be the vector of edit similarity values and $s(\omega)$ the corresponding vector
 255 of WL similarity values using weight ω . To compare the similarity measures, we calculated both the
 256 Pearson correlation coefficient and the cosine similarity of e and $s(\omega)$ for all $\omega \in W$. As can be seen from
 257 Figure 2, the highest values occur for ω between 0.38 and 0.44 and the maximum Pearson correlation
 258 coefficient is obtained for $\omega = 0.41$. For the cosine similarity, the maximum value is reached for weight
 259 $\omega = 0.69$, but the function is less peaked, and values above 0.4 lead to high agreement (Figure 3).

260 Overall, we find good agreement between edit similarity and WL similarity for sufficiently large
 261 values of ω , i.e., if the Jaccard similarity of the constituent protein multiset has sufficiently high weight.

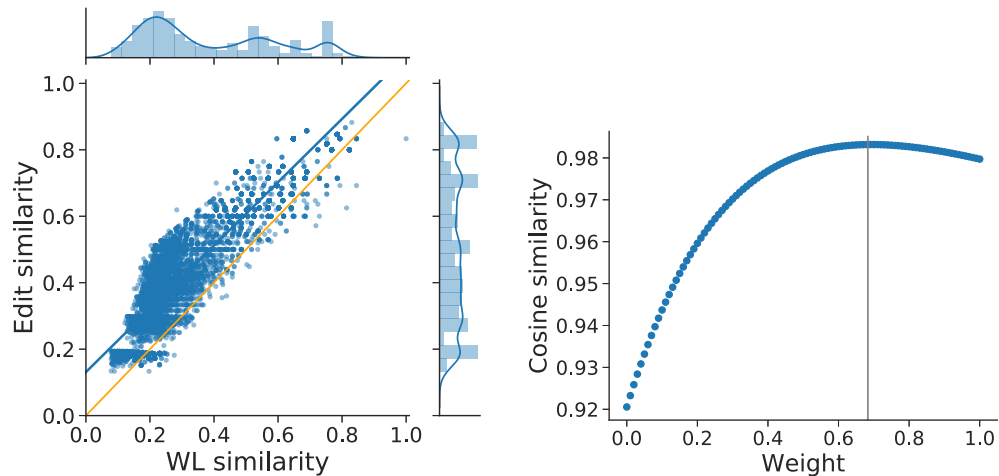


Figure 3. Left: Scatterplot comparison of edit similarity and WL similarity for weight $\omega = 0.69$, including marginal distributions and least-squares regression line. **Right:** Cosine similarity between edit similarity and WL similarity as a function of weight ω . The maximum cosine similarity (0.983) occurs for $\omega = 0.69$ shown on the left side.

Reproducibility

262

The performed data analysis is available as a reproducible Snakemake (Köster and Rahmann, 2012) workflow¹.

263

264

DISCUSSION

265

Our motivation to consider protein complex similarity was to reduce the complexity of the simulation output of our constrained protein interaction network simulator (Stöcker et al., 2017), and we were surprised to see that apparently, no similarity measures have been proposed in the literature. Depending on the underlying representation (set, multiset or graph), different alternatives suggest themselves. However, most graph-based measures are both theoretically and practically hard to compute for larger complexes or for large amounts of complexes. While different tractable graph similarity measures have been proposed, e.g. by Conte et al. (2004) or by Vishwanathan et al. (2010), or approximate graph edit distance (Riesen et al., 2015), none of these appear to be specifically tailored to the properties of protein complexes (often less than ten vertices; sparse). Our proposal to define the similarity as a convex combination of two Jaccard coefficients (protein label multiset and Weisfeiler-Lehman label multiset after one iteration) has two additional advantages. First, using Jaccard coefficients allows to efficiently pre-filter for high similarity using locality sensitive hashing. Second, for weight $\omega = 1$ of the 0-th WL iteration, our measure reduces to the natural similarity measure of the multiset representation. Our framework hence allows for a smooth transition between multiset and graph representation. The comparison to an edit-based similarity seems to indicate that the protein label multiset plays an important role if one wants to approximate the edit similarity.

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

From a biological point of view, a high similarity between two complexes should indicate a high probability that they share the same function and can substitute each other in a cellular process. If such information were available, we could evaluate each similarity measure with regard to how it relates to common function. At present, when not even the interaction topology of most complexes has been determined, the corresponding data is out of reach, and such an evaluation is not feasible. In this situation, we suggest that edit similarity is a measure that corresponds to intuition about similarity and that any reasonable similarity measure should be close to edit similarity. The measure we propose has this property (for any weight $\omega \in [0, 1]$) but offers the advantage that it can be quickly computed and scales to millions of complex pairs.

Both WL similarity and edit similarity, as previously defined, have limitations from a biological point

¹<https://doi.org/10.5281/zenodo.1178084>

292 of view in the sense that they do not consider similarities between proteins: Two proteins are either equal
293 or distinct. However, if two proteins are closely related, should they be treated as equal or distinct? In
294 the former case, we lose resolution. In the latter case, we would benefit from a fine-grained similarity
295 function between proteins (e.g. a modification of p is very similar to p , a protein with some common
296 domains is somewhat similar to p , but a completely disjoint protein in terms of domains has similarity
297 zero). In this sense, the question of how to best measure protein complex similarity is far from settled.

298 FUNDING

299 This work was supported by Deutsche Forschungsgemeinschaft (DFG) Collaborative Research Center
300 (SFB) 876, projects A6 and C1, and by Mercator Research Center Ruhr (MERCUR), project Pe-2013-0012
301 (UA Ruhr professorship).

302 REFERENCES

- 303 Arvind, V., Köbler, J., Rattan, G., and Verbitsky, O. (2015). On the power of color refinement. In
304 Kosowski, A. and Walukiewicz, I., editors, *Fundamentals of Computation Theory - 20th International*
305 *Symposium, FCT 2015, Gdańsk, Poland, August 17–19, 2015, Proceedings*, volume 9210 of *Lecture*
306 *Notes in Computer Science*, pages 339–350. Springer.
- 307 Bhowmick, S. S. and Seah, B. (2016). Clustering and summarizing protein-protein interaction networks:
308 A survey. *IEEE Trans. Knowl. Data Eng.*, 28(3):638–658.
- 309 Boutros, M., Heigwer, F., and Laufer, C. (2015). Microscopy-based high-content screening. *Cell*,
310 163(6):1314–25.
- 311 Bunke, H. (1997). On a relation between graph edit distance and maximum common subgraph. *Pattern*
312 *Recognition Letters*.
- 313 Conte, D., Foggia, P., Sansone, C., and Vento, M. (2004). Thirty years of graph matching in pattern
314 recognition. *International Journal of Pattern Recognition and Artificial Intelligence*, 18(3):265–298.
- 315 Drew, K., Lee, C., Huizar, R. L., Tu, F., Borgeson, B., McWhite, C. D., Ma, Y., Wallingford, J. B., and
316 Marcotte, E. M. (2017). Integration of over 9,000 mass spectrometry experiments builds a global map
317 of human protein complexes. *Mol Syst Biol*, 13(6):932.
- 318 Garey, M. and Johnson, D. (1979). *Computers and Intractability: A Guide to the Theory of NP-*
319 *Completeness*. W. H. Freeman.
- 320 Grecco, H. E., Imtiaz, S., and Zamir, E. (2016). Multiplexed imaging of intracellular protein networks.
321 *Cytometry A*, 89(8):761–75.
- 322 Hernandez, C., Mella, C., Navarro, G., Olivera-Nappa, A., and Araya, J. (2017). Protein complex
323 prediction via dense subgraphs and false positive analysis. *PLoS One*, 12(9):e0183460.
- 324 Hwang, L. C., Gösch, M., Lasser, T., and Wohland, T. (2006). Simultaneous multicolor fluorescence
325 cross-correlation spectroscopy to detect higher order molecular interactions using single wavelength
326 laser excitation. *Biophys J*, 91(2):715–27.
- 327 Kann, V. (1992). On the approximability of the maximum common subgraph problem. In *Proceedings of*
328 *the 9th Annual Symposium on Theoretical Aspects of Computer Science (STACS'92)*, volume 577 of
329 *LNCS*, pages 377–388, London, UK, UK. Springer-Verlag.
- 330 Kozakov, D., Hall, D. R., Xia, B., Porter, K. A., Padjhony, D., Yueh, C., Beglov, D., and Vajda, S. (2017).
331 The cluspro web server for protein-protein docking. *Nat Protoc*, 12(2):255–278.
- 332 Köster, J. and Rahmann, S. (2012). Snakemake—a scalable bioinformatics workflow engine. *Bioinfor-*
333 *matics*, 28(19):2520–2522.
- 334 Laskowski, R. A., Gerick, F., and Thornton, J. M. (2009). The structural basis of allosteric regulation in
335 proteins. *FEBS Letters*, 583(11):1692–1698.
- 336 Lerouge, J., Abu-Aisheh, Z., Raveaux, R., Héroux, P., and Adam, S. (2017). New binary linear program-
337 ming formulation to compute the graph edit distance. *Pattern Recognition*.
- 338 Ma, X. and Gao, L. (2012). Discovering protein complexes in protein interaction networks via exploring
339 the weak ties effect. *BMC Syst Biol*, 6 Suppl 1:S6.
- 340 Park, H., Lee, H., and Seok, C. (2015). High-resolution protein-protein docking by global optimization:
341 recent advances and future challenges. *Curr Opin Struct Biol*, 35:24–31.
- 342 Pearson, W. R. (2013). Selecting the Right Similarity-Scoring Matrix. *Curr Protoc Bioinformatics*,
343 43:1–9.

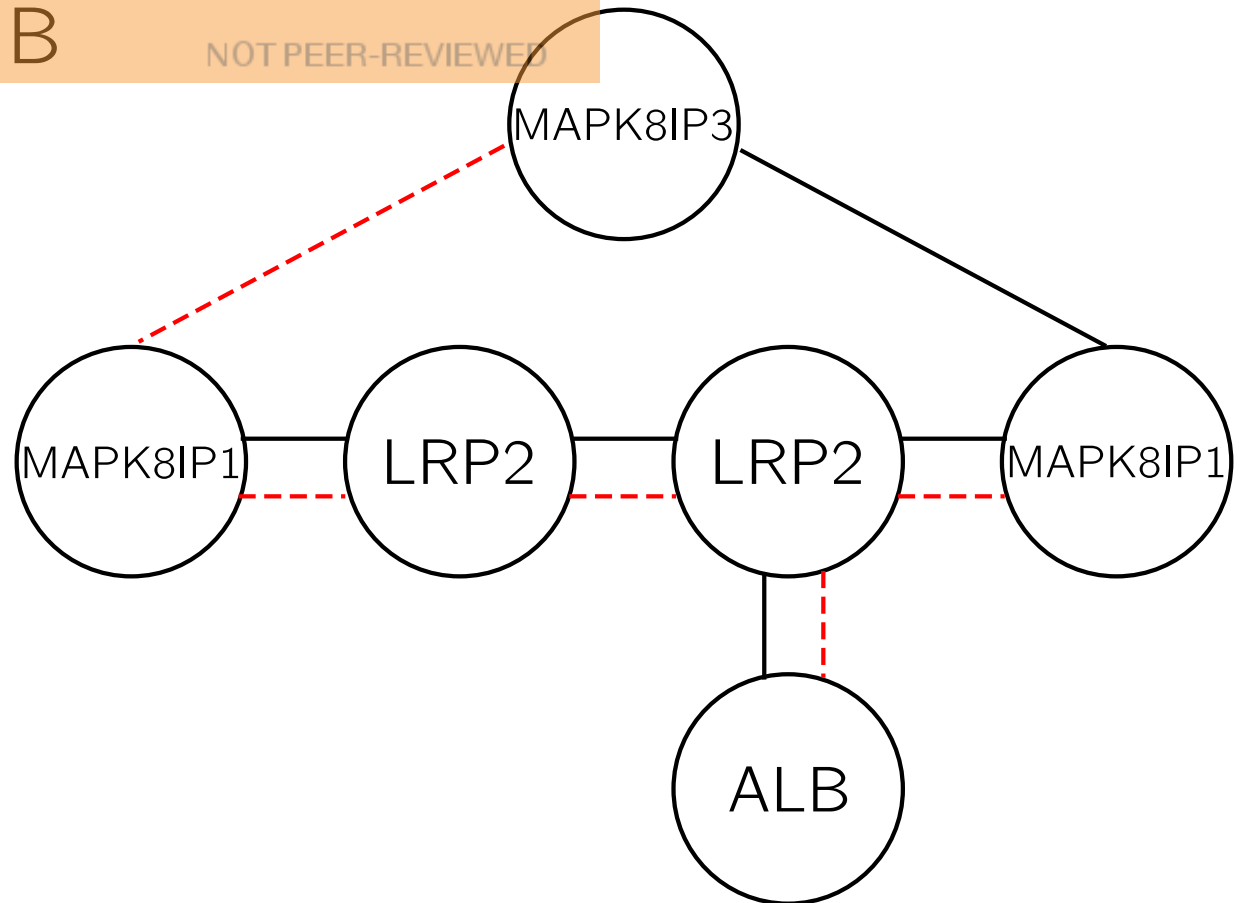
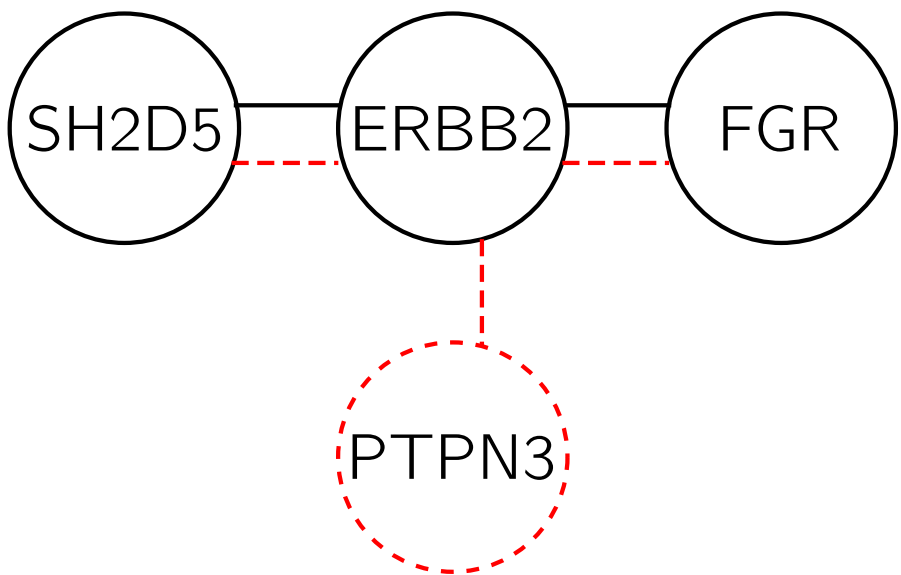
- 344 Pellegrini, M., Baglioni, M., and Geraci, F. (2016). Protein complex prediction for large protein protein
345 interaction networks with the core&peel method. *BMC Bioinformatics*, 17(Suppl 12):372.
- 346 Riesen, K., Ferrer, M., and Bunke, H. (2015). Approximate Graph Edit Distance in Quadratic Time.
347 *IEEE/ACM Trans Comput Biol Bioinform*. Epub ahead of print.
- 348 Ruepp, A., Waagele, B., Lechner, M., Brauner, B., Dunger-Kaltenbach, I., Fobo, G., Frishman, G.,
349 Montrone, C., and Mewes, H.-W. (2010). CORUM: the comprehensive resource of mammalian protein
350 complexes – 2009. *Nucleic acids research*, 38(suppl 1):D497–D501.
- 351 Sánchez Claros, C. and Tramontano, A. (2012). Detecting mutually exclusive interactions in protein-
352 protein interaction maps. *PLoS One*, 7(6):e38765.
- 353 Sanfeliu, A. and Fu, K. S. (1983). A distance measure between attributed relational graphs for pattern
354 recognition. *IEEE Transactions on Systems, Man, and Cybernetics*.
- 355 Shervashidze, N., Schweitzer, P., van Leeuwen, E. J., Mehlhorn, K., and Borgwardt, K. M. (2011).
356 Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research*, 12:2539–2561.
- 357 Srihari, S., Yong, C. H., and Wong, L. (2017). *Computational Prediction of Protein Complexes from
358 Protein Interaction Networks*. Association for Computing Machinery and Morgan & Claypool,
359 New York, NY, USA.
- 360 Stöcker, B. K., Köster, J., Zamir, E., and Rahmann, S. (2017). Modeling and simulating networks of
361 interdependent protein interactions. *bioRxiv*.
- 362 Vakser, I. A. (2014). Protein-protein docking: from interaction to interactome. *Biophys J*, 107(8):1785–93.
- 363 Vishwanathan, S. V. N., Schraudolph, N. N., Kondor, R., and Borgwardt, K. M. (2010). Graph kernels.
364 *Journal of Machine Learning Research*, 11:1201–1242.
- 365 Wachsmuth, M., Conrad, C., Bulkescher, J., Koch, B., Mahen, R., Isokane, M., Pepperkok, R., and
366 Ellenberg, J. (2015). High-throughput fluorescence correlation spectroscopy enables analysis of
367 proteome dynamics in living cells. *Nat Biotechnol*, 33(4):384–9.
- 368 Wass, M. N., Fuentes, G., Pons, C., Pazos, F., and Valencia, A. (2011). Towards the prediction of protein
369 interaction partners using physical docking. *Molecular systems biology*, 7:469.
- 370 Weisfeiler, B. and Lehman, A. A. (1968). A reduction of a graph to a canonical form and an algebra
371 arising during this reduction. *Nauchno-Tekhnicheskaya Informatsiya*, 2(9):12–16. In Russian.
- 372 Wobma, H. M., Blades, M. L., Grekova, E., McGuire, D. L., Chen, K., Chan, W. C. W., and Cramb,
373 D. T. (2012). The development of direct multicolour fluorescence cross-correlation spectroscopy:
374 towards a new tool for tracking complex biomolecular events in real-time. *Phys Chem Chem Phys*,
375 14(10):3290–4.

Figure 1(on next page)

Three exemplary pairs of protein complexes.

Each labeled node is a protein instance, each edge represents a protein interaction, and solid black vs. dashed red edges distinguish between the two complexes.

A



C

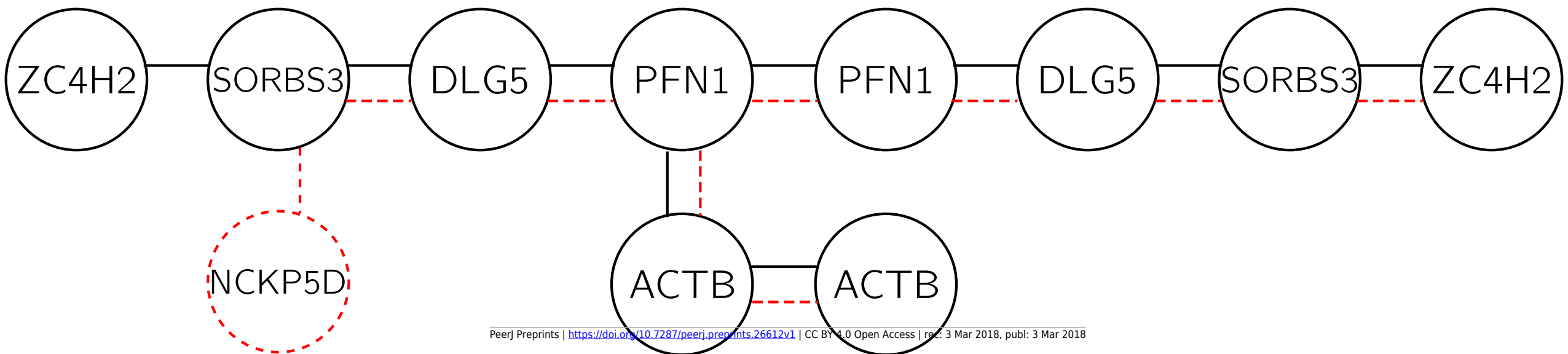


Figure 2 (on next page)

Scatterplot comparison of edit similarity and WL similarity for weight $\omega=0.41$, including marginal distributions and least-squares regression line.

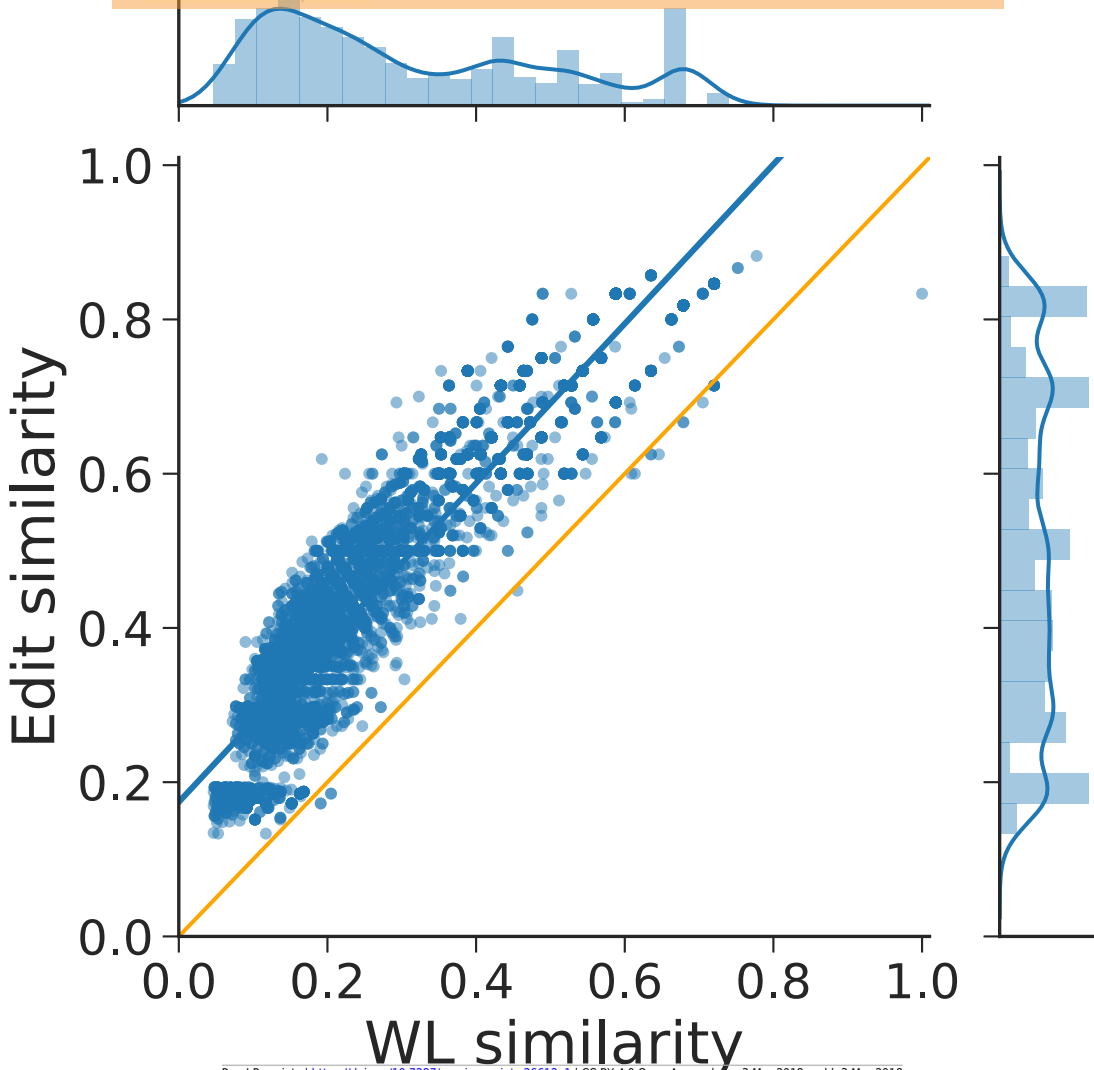


Figure 3(on next page)

Pearson correlation between edit similarity and WL similarity as a function of weight $\sim \omega$. The maximum correlation (0.946) occurs for $\omega=0.41$, shown on the left side.

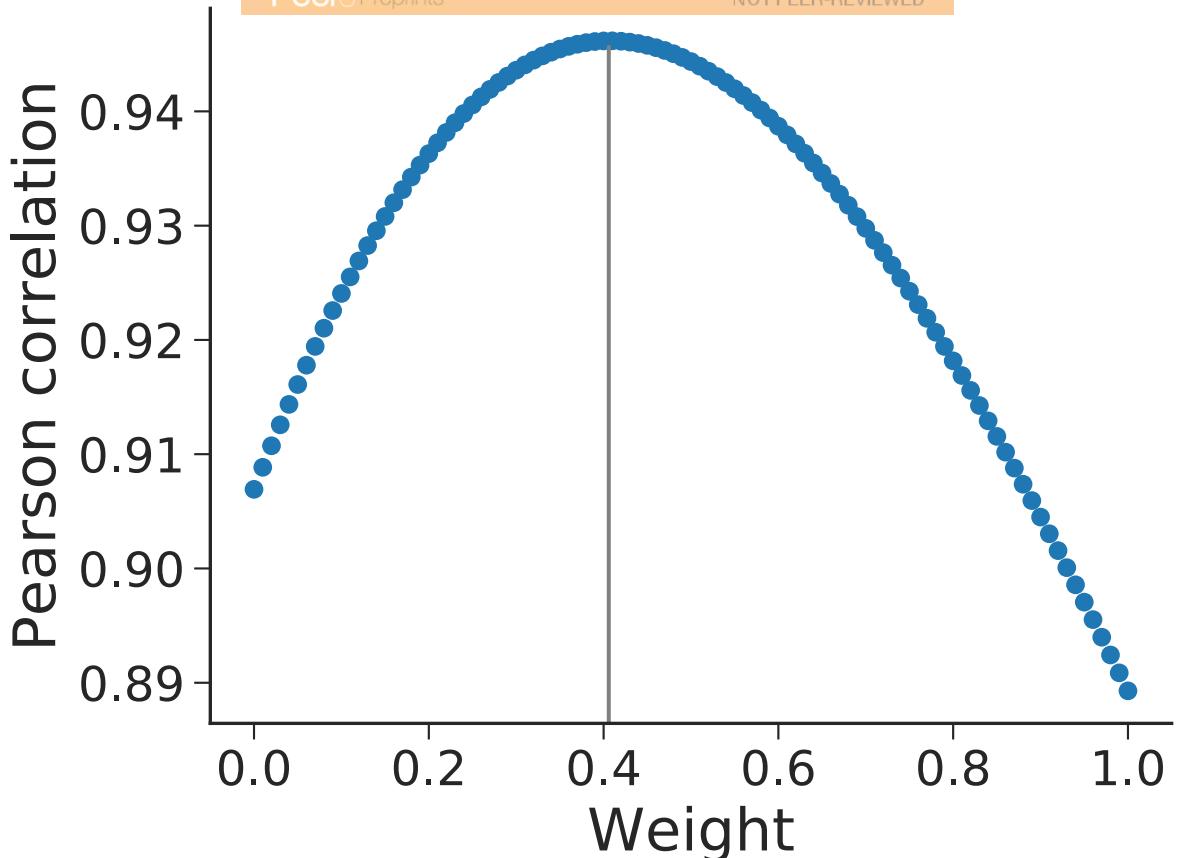


Figure 4(on next page)

Scatterplot comparison of edit similarity and WL similarity for weight $\omega=0.69$, including marginal distributions and least-squares regression line.

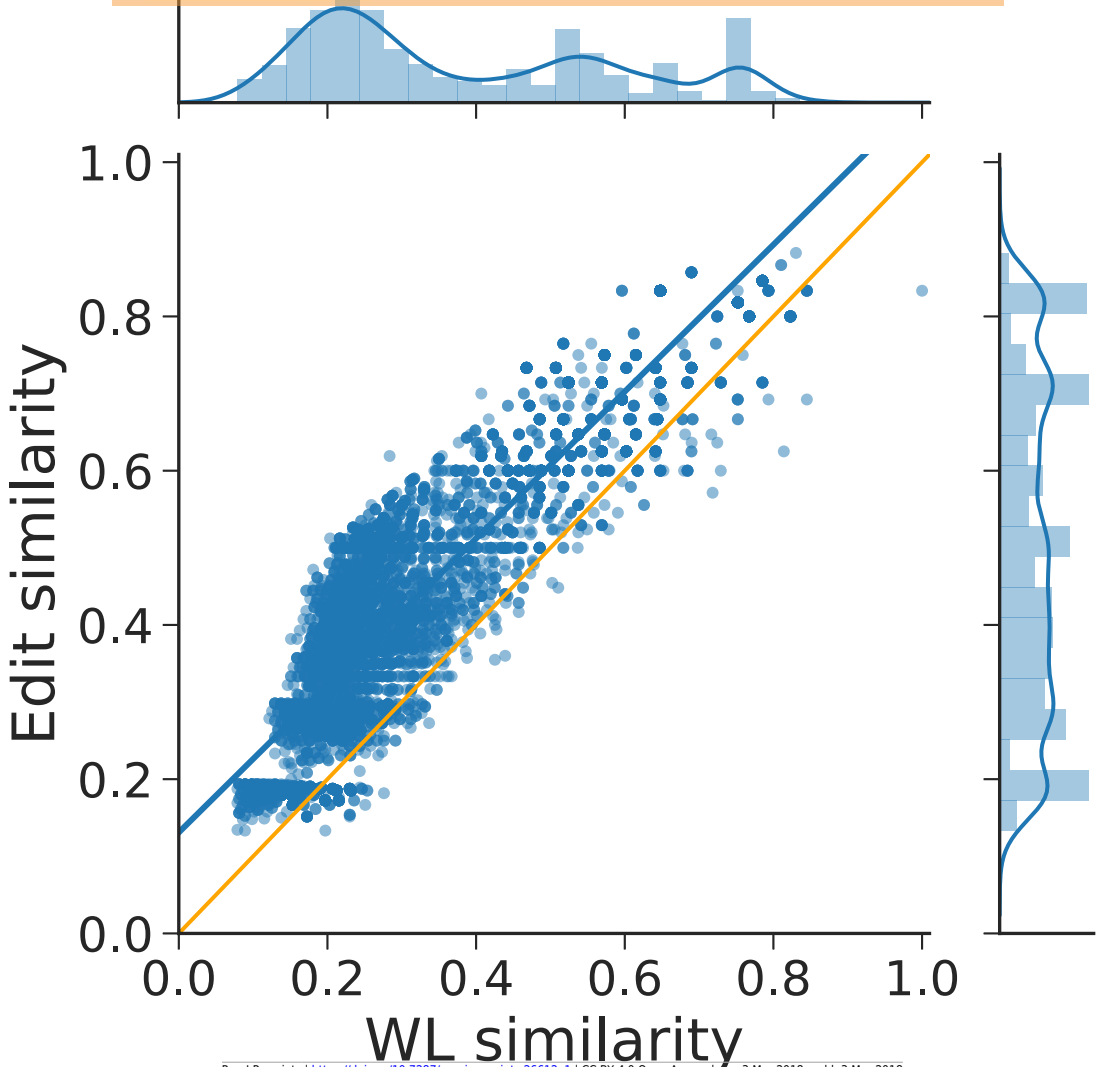


Figure 5(on next page)

Cosine similarity between edit similarity and WL similarity as a function of weight $\sim \omega$. The maximum cosine similarity (0.983) occurs for $\omega=0.69$ shown on the left side.

