

Word and sentence embedding tools to measure semantic similarity of Gene Ontology terms by their definitions

Dat Duong^{1,*}, Wasi Uddin Ahmad¹, Eleazar Eskin^{1,3}, Kai-Wei Chang¹, Jingyi Jessica Li^{2,3,*}

¹ Department of Computer Science, University of California, Los Angeles, United States of America. ² Department of Statistics, University of California, Los Angeles, United States of America. ³ Department of Human Genetics, University of California, Los Angeles, United States of America.

Availability: github.com/datduong/NLPMethods2CompareGOterms

***Contact:** datdb@cs.ucla.edu, jli@stat.ucla.edu

Abstract

The Gene Ontology (GO) database contains GO terms that describe biological functions of genes. Previous methods for comparing GO terms have relied on the fact that GO terms are organized into a tree structure. In this paradigm, the locations of two GO terms in the tree dictate their similarity score. In this paper, we introduce two new solutions for this problem, by focusing instead on the definitions of the GO terms. We apply neural network based techniques from the natural language processing (NLP) domain. The first method does not rely on the GO tree, whereas the second indirectly depends on the GO tree. In our first approach, we compare two GO definitions by treating them as two unordered sets of words. The word similarity is estimated by a word embedding model that maps words into an N-dimensional space. In our second approach, we account for the word-ordering within a sentence. We use a sentence encoder to embed GO definitions into vectors and estimate how likely one definition entails another. We validate our methods in two ways. In the first experiment, we test the model’s ability to differentiate a true protein-protein network from a randomly generated network. In the second experiment, we test the model in identifying orthologs from randomly-matched genes in human, mouse, and fly. In both experiments, a hybrid of NLP and GO-tree based method achieves the best classification accuracy.

1 Introduction

The Gene Ontology (GO) project founded in 1998 is a collaborative effort that has been providing consistent descriptions of genes and proteins across different data sources and species^[4]. The GO database is similar to a dictionary; it contains terms referred to as GO terms. Each GO term has a definition describing some biological event.

The GO database is divided into three categories: cellular components (CC), molecular functions (MF) and biological processes (BP). The CC ontology contains terms describing the components of the cell and can be used to locate a protein. The MF category contains

terms describing chemical reactions such as *catalytic activity* or *receptor binding*. These terms do not specify the genes or proteins involved in the reactions or the locations of the events. The BP category contains terms describing a series of events which may contain cellular locations and molecular reactions. For example, the BP term GO:0006874 has the definition "Any process involved in the maintenance of an internal steady state of calcium ions at the level of a cell."

In each category, the GO terms are organized into a tree where there is only one *root* node^[4]. In this tree of GO terms (or GO tree), a more generic term (i.e. lyase activity) is closer to the root, whereas a more specific term (i.e. carboxy-lyase activity) is closer to a leaf node.

Because there are three GO categories in the database, there are three GO trees. Interestingly, there are edges connecting terms in different GO trees (Figure 1). The GO database can be represented as three connected GO trees.

One application of the GO database is the comparison of two genes by first comparing the similarity of the GO terms that annotate them^[4]. To this end, we need a good metric for comparing GO terms. To solve this problem, we need to focus on the GO trees and the definitions of GO terms.

Because of the GO trees, GO terms with a direct ancestor (i.e. sibling nodes) are deemed to be more related than GO terms with a distal ancestor. Moreover, because of this design, existing methods to measure the similarity of two GO terms entirely rely on the GO trees^[12]. To our knowledge, no study has yet to directly compare the definitions of GO terms.

In this paper, we introduce two new solutions to this problem, by focusing instead on the definitions of the GO terms. We apply neural network based techniques from the natural language processing (NLP) domain.

First, we compare words by converting them into word embeddings. We train the Word2vec model using open access articles on PubMed, so that we can represent a word as an N-dimensional vector^[13]. Cosine similarity is used to compare two words. To compare two GO terms, we treat their definitions as two unordered sets of words, and use the weighted Modified Hausdorff Distance to measure the distance between two sets^[2]. We name this metric w2vGO. w2vGO is entirely independent of the GO trees.

Second, we consider the word-ordering within the GO definitions. We note that *entailment* relationships exist in the GO tree (i.e. two GO terms are linked by a directed edge) (Figure 1). We train the sentence encoder InferSent¹ using the definitions of child-parent and randomly-matched GO terms^[1]. InferSent embeds sentences into an N-dimensional vector space and computes the probability that one sentence entails another. We name this approach InferSentGO. InferSentGO needs the GO trees for the training phase. Once the model is trained, only the definitions of GO terms are required for calculating their semantic similarities.

We compare W2vGO and InferSentGO against two tree-based methods Resnik and Aggregate Information Content (AIC)^[18,20]. We also consider an ensemble approach AicInferSentGO by averaging the AIC and InferSentGO scores.

¹The original authors did not name their method; hence, we use their GitHub name. The NLP community refers to it as bi-directional long short-term memory with max-pooling.

To assess the performance for these five metrics, we conduct two experiments. In the first experiment, we test these metrics in differentiating the human protein-protein interaction network from a network where edges are randomly assigned. In the second experiment, we test the metrics in identifying orthologs against randomly-matched gene pairs for human, mouse, and fly. Our results show that the hybrid AicInferSentGO attains the best classification in terms of the area under the receiver operating characteristic (ROC) curve. Our software, data, and results are available at our GitHub².

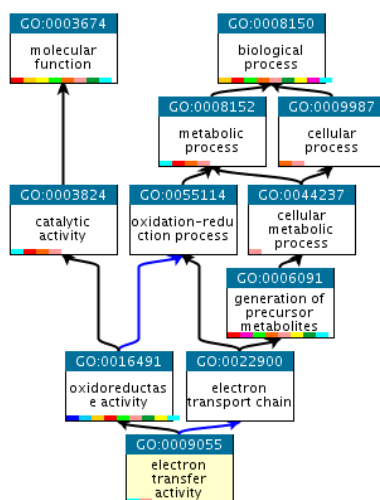


Figure 1: Terms shown are ancestors of GO:0009055 (yellow). GO:0003674 and GO:0008150 are root nodes for the MF and BP trees, respectively. Colors denote *part of* (blue) and *is a* (black) relationship. Snapshot is downloaded from ebi.ac.uk/QuickGO.

2 Method

2.1 Node-based methods to measure similarity between two GO terms

Broadly speaking, existing tree-based methods are divided into two types: node-based or edge-based^[11]. The key focus in node-based methods is the evaluation of the information content of the common ancestors for two GO terms. In brief, the information content of a GO term measures the usefulness of the term by evaluating how often the term is used to annotate a gene. Terms that are used sparingly have high information content because they are specific at distinguishing genes. Node-based methods have been shown to work well^[10,11,20]. However, we show in section 2.4 that we can improve the accuracy by including the semantic similarities of the GO definitions.

Unlike node-based methods, edge-based approaches measures the distance (or the average distance when more than one path exists) between two GO terms. Edge-based methods have one serious problem. Despite being organized into a tree, the GO terms at the same level do not always have the same specificity because different gene properties

²github.com/datduong/NLPMethods2CompareGOterms

require different levels of detailed explanation. Thus, edge-based methods suffer from the problem of *shallow annotation*: terms separating by the same distance are assigned the same similarity score regardless of their positions in the GO tree^[20]. Different schemes to weigh the edges according to their positions in the GO tree failed to fully resolve the problem^[10]. Node-based methods have been shown to be more successful than edge-based methods^[10,11,20]. Therefore, we do not compare our models against edge-based methods in our paper.

In this paper, we choose the node-based methods Resnik and AIC as the baseline for the following reasons. Resnik is a classical approach for quantifying the similarity between two GO terms^[18]. Despite being simple, Resnik has been shown to outdo some of its extensions on several test datasets^[11,16]. AIC is recent. In their paper, Song et al.^[20] showed that AIC outperforms well known node-based methods by Jiang and Conrath^[5], Lin^[9], Resnik^[18], and Wang et al.^[22].

2.1.1 Resnik method

The most basic node-based method introduced by Resnik in 1999 relies on the information content (IC) of a GO term^[18]. IC of a GO term t is computed as $IC(t) = -\log(p(t))$ where $p(t)$ is the probability of observing a term t in the ontology. $p(t)$ is computed as $p(t) = \frac{\text{freq}(t)}{\text{freq}(\text{root})}$. $\text{freq}(t)$ is defined as the cumulative count of term t and its descendants, where $\text{freq}(t) = \text{count}(t) + \sum_{c \in \text{child}(t)} \text{freq}(c)$. $\text{count}(t)$ is the number of genes annotated with the term t , and $\text{child}(t)$ are the children of t . Based on this definition, $IC(\text{root}) = 0$, and a node near the leaves has higher IC than nodes at upper levels. To compute a similarity score of the GO terms a and b , one finds the most informative common ancestor of these two terms.

$$\text{Resnik}(a, b) = \max_{p \in \{\text{par}(a) \cap \text{par}(b)\}} IC(p), \quad (1)$$

where $\text{par}(t)$ denotes all the ancestors of term t . $\text{Resnik}(a, b)$ ranges from 0 to ∞ because the probability $p(t)$ ranges from 0 to 1.

In this model, the similarity score between a GO term t and itself is not 1. Second, when a, b have only the *root* as a common ancestor, then $\text{Resnik}(a, b) = 0$. This is problematic because leaf nodes are more informative than other types of nodes. Consider the example in Song et al.^[20]. Here, the *root* is the only common ancestor of the pair a, b and the pair c, d . Next, suppose that a, b are leaf nodes, c is the parent of a , d is the parent of b , and *root* is the parent of both c, d . One would then expect that $\text{Resnik}(a, b) < \text{Resnik}(c, d)$; however, one would obtain $\text{Resnik}(a, b) = \text{Resnik}(c, d)$.

2.1.2 Aggregate Information Content (AIC) Method

The AIC method by Song et al.^[20] amends the two problems in the Resnik method. To encode the fact that leaf nodes are more informative, AIC defines a knowledge function of term t as $k(t) = 1/IC(t)$ which is used to measure its semantic weight $sw(t) = 1/(1 + \exp(-k(t)))$. Here $sw(\text{root}) = 1$. Semantic value $sv(t)$ of t is then defined as $sv(t) = \sum_{p \in \text{path}(t)} sw(p)$. Function $\text{path}(t)$ contains every ancestor of t and the term t itself. Usually,

$sv(a) < sv(b)$ when term a is nearer to the root than b . Song *et al.*^[20] define their similarity score of two GO terms a, b as

$$AIC(a, b) = \frac{2 \sum_{p \in \{\text{path}(a) \cap \text{path}(b)\}} sw(p)}{sv(a) + sv(b)}. \quad (2)$$

$AIC(a, b)$ ranges from 0 to 1. In this model, $AIC(a, a) = 1$. When a, b have only the *root* as the common ancestor, then $AIC(a, b) = 2/(sv(a) + sv(b))$ which depends on where a, b are on the GO tree.

2.2 Word2vec model

Here, we describe our first metric W2vGO which focuses strictly on the GO definitions and ignores the GO-trees. To compare two GO terms, this metric compares their definitions by treating the definitions as two unordered sets of words. To solve this problem, we want to first be able to compare two words. To this end, we use the word embedding model Word2vec^[13].

The Word2vec model converts a word into an N -dimensional vector³. These vectors are known as word embeddings. Word2vec transforms similar words into similar vectors, thus enabling us to quantify the similarity between two words by computing the Euclidean distance or cosine similarity. At the heart of the Word2vec is a neural network model with one input layer, one hidden layer, and one output layer^[13].

Loosely speaking, one can view the Word2vec model as a prediction problem^[19]. First, a word w from the input layer is mapped into an N -dimensional vector at the hidden layer. Word2vec learns the values for the hidden layer based on the co-concurrences between w and its neighboring words. The key idea is to predict the vectors for the surrounding context words based on the vector for w .

The purpose of this paper is not to dissect the Word2vec model; we are interested in adopting this model to measure the similarity of GO terms and compare it with other methods. Interested readers are encouraged to read the original paper by Mikolov *et al.*^[13], and the introduction by Rong^[19].

2.2.1 Measure similarity of two words using Word2vec

The training data influences the application of word embedding model. Existing pre-trained Word2vec models are often made by corpora collected from news, books, or the Internet. To obtain suitable word vectors, in this work, we train the Word2vec to recognize biological words. We set the dimension $N = 300$ and use 20 GB of data from open access articles on PubMed. The raw count of unique and repeated words is 14,526,527,855. We remove words which appear less than 25 times in the whole training data, thus reducing the final number of unique words to 986,615. We keep stop-words and symbols like + and – in the data because they may have important biological meanings. We use the Python library gensim to train the Word2vec model^[17]. A simple Python user interface is available at our GitHub.

³The user specifies N ; in practice, people often set $N = 300$.

There are two important details here. First, the training data do not contain definitions of GO terms found in the GO database. This helps us avoid data reusing. Second, theoretically speaking, Word2vec model can be trained on the GO terms in the PubMed data, so that one can convert a GO term into a vector. Unfortunately, the IDs of the GO terms are not used too often in published papers, and detecting definitions of GO terms in papers is a different type of research problem^[21]. For these reasons, we use the Word2vec model as a metric to compare two biological words.

To compare two vector representations of two words, we use the cosine similarity, because it is bounded; whereas, Euclidean distance is not. We define the function $w2v(z, v)$ as the similarity score of two words z, v .

2.2.2 Measuring similarity of two GO terms using Word2vec

A GO term comes with a definition which is usually one or two sentences describing some biological feature. For example, GO:0003700 has the definition: "Interacting selectively and non-covalently with a specific DNA sequence in order to modulate transcription. The transcription factor may or may not also interact selectively with a protein or macromolecular complex."

When a GO term definition has more than one sentence, we concatenate these sentences into the same sentence by ignoring the period symbol. For example, the two sentences for GO:0003700 is considered as one long sentence.

Thus, the task to compare two GO terms reduces to the problem of comparing their definitions which are two sentences. Suppose that GO terms a, b have sentences Z, V as their definitions respectively. We treat two sentences $Z = "z_1 z_2 z_3 \dots z_N"$ and $V = "v_1 v_2 v_3 \dots v_M"$ as two unordered sets of words $Z = \{z_1, z_2 \dots z_N\}$ and $V = \{v_1, v_2 \dots v_M\}$.

To measure the similarity of sentences (or sets) Z and V , we use the metric

$$w2vGO(Z, V) = \text{mean} \left\{ \sum_{i=1 \dots N} \text{content}(z_i) \max_{j=1 \dots M} w2v(z_i, v_j), \right. \\ \left. \sum_{j=1 \dots M} \text{content}(v_j) \max_{i=1 \dots N} w2v(z_i, v_j) \right\} \quad (3)$$

where $\text{content}(w)$ is the weight of the word w and is often used to distinguish common words from rare ones. The weights of words can help avoid the influence of hub-words (i.e. words such as *cell*, *DNA*, *activity*) that are ubiquitously associated with many other words^[7]. $\text{content}(w)$ is very similar to the IC function^[8]

$$\text{content}(w) = -\log \left(\frac{\text{frequency word } w \text{ in training data}}{\text{training data size}} \right). \quad (4)$$

Because GO definitions are often short, we hope that the accuracy of $w2vGO$ does not suffer too much from the removal of word-ordering. Surprisingly, this assumption holds true in several instances (Table 2). There are more sophisticated models that consider the word-ordering in the sentences; we will consider one of these methods in the next section.

In any case, we have defined a metric to measure the two GO terms a, b with definitions Z, V under the Word2vec paradigm. Because a GO term and its definition are two equivalent entities, we define $w2vGO(a, b) = w2vGO(Z, V)$ to be the similarity score of a, b . $w2vGO(a, b)$ ranges from -1 to 1 because $w2v(z, v)$ ranges from -1 to 1 .

2.3 InferSent model

In the previous section, we have ignored the word-ordering in the sentences, treating them as sets of words. In this section, we briefly explain InferSent, a model that focuses not only on the word embeddings but also on the word-ordering in the sentences^[1]. Loosely speaking, InferSent is similar in spirit to Word2vec; instead of words, InferSent identifies the relationship between two sentences. One training sample for InferSent consists of two sentences having some relationship R . Consider a simple classification where $R = \textit{entailment}$ or *neutral*. In this case, InferSent expect two classes of input. The first class *entailment* will have two sentences where the first *entails* the second; for example, "cat is napping on the mat" *entails* "cat is not running". *Entailment* is a one-directional relationship, because "cat is not running" does not *entail* "cat is napping on the mat". The second class *neutral* will have two unrelated sentences like "cat is napping on the mat" and "boy is watching the cat". Unlike *entailment*, *neutral* is a bi-directional relationship.

InferSent is a classification model based on the neural network architecture; its full description is at Conneau et al.^[1]. Here, we briefly mention the first layer of this network. InferSent's first layer is the word vectors for words in the entire training dataset. Conneau et al.^[1] uses the GloVe word vectors^[14]. However, to obtain the best result, these word vectors should be specific to biology. In this paper, we use the Word2vec vectors in section 2.2.1.

2.3.1 Measuring similarity of two GO terms using InferSent

InferSent takes two sentences as one training sample. In this paper, the two sentences will be the definitions of two GO terms a, b . We will define two categories *entailment* and *neutral*, and estimate the probability $\mathbb{P}(a \textit{ entails } b)$. This metric allows us to gauge the semantic similarity for a and b . We choose this option because *entailment* relationship exists in the GO tree. Child-parent GO terms are linked by a one-directional relationship like "is a", "part of", "regulates", "negatively regulates", and "positively regulates". For example, the term GO:1900237 "positive regulation of induction of conjugation with cellular fusion" *entails* the term GO:0010514 "induction of conjugation with cellular fusion."

To prepare the *entailment* dataset, for all three ontology categories, we randomly pair each GO term with one of its parents. To ensure that these child-parent GO terms are indeed similar in meaning, we compute the median AIC score for each category and retain pairs having scores above the median. Our final dataset contains 17,226 pairs. We treat the three ontology categories as one single dataset when training the InferSent model.

To create the *neutral* dataset, we make two types of unrelated pairs. For the first type, we randomly pick about half the number of GO terms in the *entailment* dataset. For each term c in this set, we pair it with a randomly chosen GO term d in the same ontology

category. For the second type, we pair the same term d with another randomly chosen term e . This sampling scheme improves the training by allowing some GO terms to be seen more than once under different circumstances.

Two unrelated GO terms should have $\mathbb{P}(\text{term}_1 \text{ entails term}_2)$ and $\mathbb{P}(\text{term}_2 \text{ entails term}_1)$ near zero. For each *neutral* pair, we create two different samples for InferSent. The first sample will be the pair (definition of term_1 , definition of term_2), and the second sample will be the pair (definition of term_2 , definition of term_1). Our final *neutral* dataset contains 35,044 pairs of sentences. Since the *neutral* dataset is nearly twice as large as the *entailment* dataset, when training InferSent, we weigh ratio 1:2 for the class *neutral* and *entailment*.

The code to train InferSent is available at our GitHub. We attained 96.93% accuracy in the validation set. We emphasize that InferSent relies on the GO tree only for its training phase. Once the model is trained, it requires only the GO definitions as the input for prediction.

When terms a and b are unrelated, we expect $\mathbb{P}(a \text{ entails } b) \approx \mathbb{P}(b \text{ entails } a) \approx 0$. When GO term a is the child of term b , we expect $\mathbb{P}(a \text{ entails } b)$ to be high, whereas $\mathbb{P}(b \text{ entails } a)$ may be low. However, in this case, we still want the similarity between a, b to be high. For this reason, to measure the semantic similarity for two GO terms, we use the metric

$$\text{InferSentGO}(a, b) = \max \{ \mathbb{P}(a \text{ entails } b), \mathbb{P}(b \text{ entails } a) \}. \quad (5)$$

$\text{InferSentGO}(a, b)$ ranges from 0 to 1.

2.4 Combining node-based and NLP methods

In this section, we discuss a few examples and motivate the need for combining node-based and NLP methods. In essence, both Resnik and AIC focus on the fraction of shared ancestors and weigh this ratio by the IC values. Sometimes, even when two terms are very similar in meaning, they can have a low score. For example, consider the terms GO:0005887 and GO:0016021 (Table 1). This is a child-parent pair, having almost identical definition. However, the Resnik and AIC score are not high enough as compared to W2vGO and InferSentGO.

In the second example, GO:0006814 and GO:0006874 are unrelated, sharing only the root node. Interestingly, one can argue that both terms are not entirely distinct because they both mention the regulation of ions. Similarly, in the third example, GO:0005829 and GO:0005615 share only the root node, but both mention fluid containing protein complexes. In both examples, W2vGO on its own or an average of AIC and InferSentGO may give a more satisfying score. InferSentGO and AIC on their own may underestimate and overestimate the similarity, respectively. Resnik on its own is not the best because when terms share only the root node, the similarity score is the IC of the root which is 0.

In the final example, GO:0004620 and GO:0019905 share only the root node and truly are different in meaning. Here, Resnik and InferSentGO give more reasonable scores than AIC and W2vGO do.

These examples suggest that no one method is always the best, and that we need to combine node-based and NLP approaches. To this end, taking an average of AIC and InferSentGO is reasonable. Empirically, from the examples, we have seen that this average

Table 1: A few examples to compare GO similarity scores. *Fraction of shared ancestors, with 0 indicates terms share only the root node.

Term 1	Term 2	Ancestor*	Resnik	AIC	W2vGO	InferSentGO
GO:0016021 The component of a membrane consisting of the gene products and protein complexes having at least some part of their peptide sequence embedded in the hydrophobic region of the membrane.	GO:0005887 The component of the plasma membrane consisting of the gene products and protein complexes having at least some part of their peptide sequence embedded in the hydrophobic region of the membrane.	5/9	2.493	0.588	0.982	0.999
GO:0006814 The directed movement of sodium ions (Na+) into, out of or within a cell, or between cells, by means of some agent such as a transporter or pore.	GO:0006874 Any process involved in the maintenance of an internal steady state of calcium ions at the level of a cell.	0	0	0.107	0.590	0.920
GO:0005829 The part of the cytoplasm that does not contain organelles but which does contain other particulate matter, such as protein complexes.	GO:0005615 That part of a multicellular organism outside the cells proper, usually taken to be outside the plasma membranes, and occupied by fluid.	0	0	0.212	0.449	0.845
GO:0004620 Catalysis of the hydrolysis of a glycerophospholipid.	GO:0019905 Interacting selectively and non-covalently with a syntaxin, a SNAP receptor involved in the docking of synaptic vesicles at the presynaptic zone of a synapse.	0	0	0.264	0.333	0.001

produces a reasonable value. Theoretically, only AIC and InferSentGO scores are in the same $[0, 1]$ range; whereas Resnik and W2vGO range are $[0, \infty]$ and $[-1, 1]$.

For two GO terms a, b , we introduce the metric

$$\text{AicInferSentGO}(a, b) = \frac{\text{AIC}(a, b) + \text{InferSentGO}(a, b)}{2}. \quad (6)$$

From AIC’s perspective, AicInferSentGO improves the distinction between child-parent and randomly-matched GO terms, especially in the CC and MF ontology (Figure 2). From InferSentGO’s perspective, AicInferSentGO gives a more continuous score, allowing for a better resolution when comparing terms. InferSentGO on its own tends to give a stiff 0/1 score.

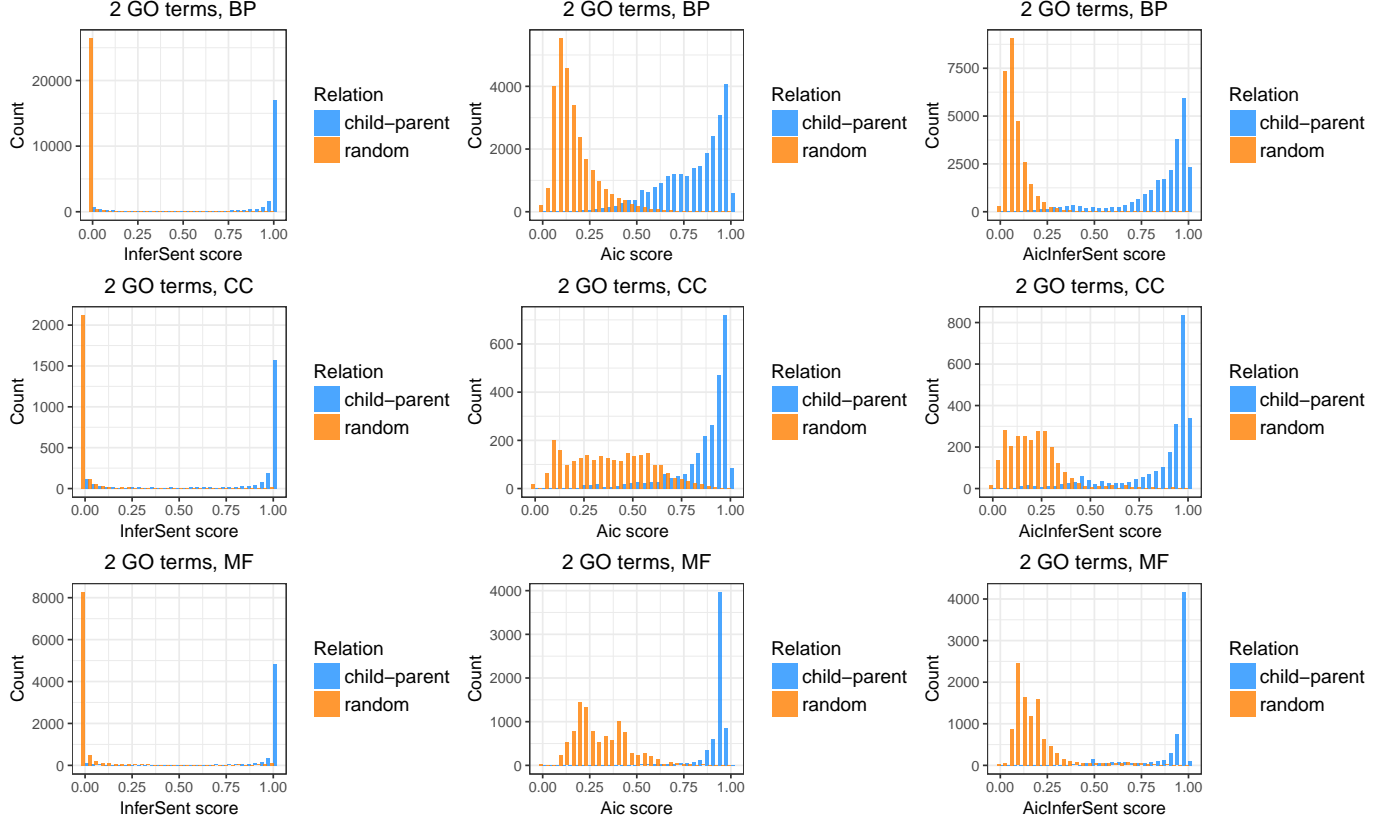


Figure 2: Similarity scores for 33020 child-parent GO terms, and 40419 randomly-matched GO terms. The number of pairs are 51931, 4986, and 16522 for BP, CC, and MF ontology respectively. From AIC’s perspective, the hybrid method AicInferSentGO improves the distinction between child-parent and randomly-matched GO terms. From InferSentGO’s perspective, AicInferSentGO gives a more continuous score.

2.5 Measuring similarity of two genes

To assess the performance of Resnik, AIC, W2vGO, InferSentGO and AicInferSentGO, we will use them to compare genes. A gene is annotated with several GO terms from the three GO categories. For example, the gene HOXD4, which is important for morphogenesis, is annotated by these GO terms GO:0003677, GO:0003700, and GO:0006355. Thus, we can view any gene A as a set of GO terms. A GO term a is in the set A (i.e. $a \in A$) if a is used to annotate A .

To assess the similarity between two genes A and B , we must compare two sets of GO terms. There are many metrics for this task^[12]. Here, we use the Modified Hausdorff Distance (MHD) and the Best Max Average distance (BMA). MHD is a traditional metric for comparing two sets, and was often used in image processing (i.e. to compare two sets of pixels)^[2]. BMA has been shown to be better than taking the maximum or minimum of all pairwise distances for the elements in the two sets^[15].

$$\text{MHD}(A, B) = \min \left\{ \frac{1}{|A|} \sum_{a \in A} \max_{b \in B} s(a, b), \frac{1}{|B|} \sum_{b \in B} \max_{a \in A} s(a, b) \right\} \quad (7)$$

$$\text{BMA}(A, B) = \text{mean} \left\{ \frac{1}{|A|} \sum_{a \in A} \max_{b \in B} s(a, b), \frac{1}{|B|} \sum_{b \in B} \max_{a \in A} s(a, b) \right\} \quad (8)$$

In the above, the function $s(a, b)$ is a generic placeholder for measuring the similarity of GO terms a and b . For example, if one uses Resnik, AIC, or W2vGO metric then $s(a, b) = \text{Resnik}(a, b)$, $\text{AIC}(a, b)$ or $\text{w2vGO}(a, b)$ respectively. $\text{MHD}(A, B)$ and $\text{BMA}(A, B)$ have ranges $[0, 1]$ for AIC, InferSentGO and AicInferSentGO, $[0, \infty]$ for Resnik, and $[-1, 1]$ for w2vGO.

3 Results

We compare the five GO metrics. Because genes are annotated by GO terms, good GO metrics should differentiate similar genes from unrelated genes well. Hence, we conduct two experiments. First, we test the GO metrics in identifying true protein-protein interactions. Second, we test the metrics in identifying orthologs in human, mouse, and fly. We download the GO term definitions and GO annotations from the Gene Ontology (geneontology.org). We download the orthologs from Ensembl (ensembl.org/biomart). The source code, data, and results in this section are available at our GitHub.

3.1 Protein-protein interaction network

We use the 6031 protein-protein interaction (PPI) data prepared by Mazandu and Mulder^[11]. We trim this data further, keeping only human proteins that can be mapped to some genes via UniProt (uniprot.org). Next, to avoid data reusing, we remove electronically inferred GO terms (removing terms with tag IEA, NAS, NA, NR)^[15]. To keep only genes that are well studied, we retain only genes with at least one GO term in each ontology (BP, CC, and MF). The final data has 2593 pairs.

Like in Mazandu and Mulder^[11], we want to compare how well each metric differentiates a true PPI network (positive set) from a randomly made PPI network (negative set). We follow the procedure by Mazandu and Mulder^[11]. We make the positive and negative sets to have the same number of edges. For the negative set, we randomly assign edges between proteins that do not interact in the real PPI network. The real and random PPI network have the same proteins; we only require that they have different interacting partners. For each PPI network, to compute the similarity scores of the edges (i.e. pairs of proteins), we use Eq. 7 and 8 with $s(a, b)$ being one of the five metrics Resnik, AIC, w2vGO, InferSentGO, and AicInferSentGO

To compare the performance of the five metrics, we find the area under the curve (AUC) of the Receiver Operative Characteristic (ROC) curve. The real and random PPI

networks serve as a basis to calculate the true positive and false negative rate, respectively. The AUC is computed by plotting the true positive versus false negative rates at different thresholds and estimating the area under this curve. AUC value goes from 0 to 1, with 1 being the best prediction power.

We judge the GO metrics based on their AUC values. From Fig. 2, because node-based methods adequately distinguish related BP terms from unrelated ones, the NLP methods' improvement is best seen in the CC and MF ontologies (Table 2).

On average, each gene in this experiment is annotated by 20.66 GO terms, with the composition of 46.25% BP, 24.07% CC, and 29.68% MF terms. Because of these fractions, when using only the BP ontology to compare GO metrics, on average, we are using only 46.25% of the full description for a gene. The same argument can be made for using only the CC or MF ontology to compare GO metrics. For this reason, we conduct a joint analysis. Here, when comparing two genes, we use all the GO terms in their annotations, allowing for comparison of GO terms across different ontologies. This approach aligns with the observation that GO terms in different categories are connected (Fig. 1). Table 2 and 3 indeed show that the joint analyses yield the highest classification accuracy for all GO metrics. We have two explanations for this outcome.

First, intuitively, when using all BP, CC, and MF terms in the gene annotation, one can better understand the genes' functionality. For example, when looking at the CC ontology terms alone, arguably proteins in the same part of the cell do not necessarily interact. However, when we consider not only the locations but also the biological and molecular events taking place, then we can accurately compare the two genes.

Second, empirically, in 46,967 randomly chosen child-parent pairs, we count 2060 pairs (4.38%) having terms in different ontologies. A few terms having parents in different ontologies are GO:0009055, GO:0035514, GO:0102496, GO:1903198, and GO:1903934. The fraction of these terms 4.38%, despite being small, has a nontrivial repercussion.

This effect is especially true for Resnik and AIC whose key ideas rely on the number of common ancestors. For example, consider the term GO:0009055 in the MF ontology with its parent GO:0022900 in the BP category (Fig. 1). When treating the BP and MF trees separately, GO:0009055 and GO:0022900 have AIC score zero because they will not have any shared ancestors. When treating the trees jointly, the AIC score is 0.7197. Thus, we can better estimate the similarity between genes containing not only these terms but their descendant terms.

For reasons explained above, in this paper, we select the metric with the highest AUC in the joint analysis to be the best method. Here, Table 2 shows that AicInferSentGO is the winner. We would like to note that W2vGO, despite being simple, works quite well on its own (2nd rank). We provide the ROC plots for the joint analyses at our GitHub.

3.2 Orthologs

Like in section 3.1, we remove electronically inferred GO terms from the gene annotation, and use genes with at least one GO term in each ontology. We test the following species: human/mouse and human/fly. For each pair, the positive set contains orthologs from the two species; whereas, the negative set contains randomly-matched genes. We set the

Table 2: AUCs for classifying human protein-protein interactions. Bold font indicates the best value in each column. *Joint analysis: When comparing genes, we keep their entire GO annotations, effectively treating the BP, CC, MF ontologies as connected GO trees.

Set metric	GO metric	BP	CC	MF	Joint*
MHD	Resnik	0.84194	0.77278	0.70436	0.86074
	AIC	0.84042	0.76218	0.69590	0.85815
	W2vGO	0.8246	0.79338	0.71915	0.86571
	InferSentGO	0.83877	0.76026	0.75343	0.85368
	AicInferSentGO	0.84830	0.78263	0.74489	0.86871
BMA	Resnik	0.85434	0.77871	0.70348	0.86766
	AIC	0.85423	0.77628	0.69902	0.87854
	W2vGO	0.84296	0.80035	0.71279	0.88170
	InferSentGO	0.84600	0.76548	0.76081	0.87503
	AicInferSentGO	0.85739	0.79011	0.74139	0.88987

Table 3: AUCs for classifying orthologs. Bold font indicates the best value in each column. *Joint analysis: When comparing genes, we keep their entire GO annotations, effectively treating the BP, CC, MF ontologies as connected GO trees.

Set metric	GO metric	BP	CC	MF	Joint*
MHD	Resnik	0.91815	0.89986	0.90938	0.95163
	AIC	0.92799	0.89843	0.91091	0.95443
	W2vGO	0.92630	0.92647	0.90735	0.95927
	InferSentGO	0.92616	0.88021	0.89532	0.95608
	AicInferSentGO	0.93165	0.90478	0.91196	0.96732
BMA	Resnik	0.92333	0.89358	0.90951	0.95306
	AIC	0.93335	0.90788	0.91700	0.95820
	W2vGO	0.93105	0.92107	0.91007	0.96134
	InferSentGO	0.92795	0.89061	0.89931	0.96010
	AicInferSentGO	0.93580	0.91464	0.91689	0.97071

human/mouse orthologs

Set metric	GO metric	BP	CC	MF	Joint*
MHD	Resnik	0.87291	0.84752	0.88819	0.93411
	AIC	0.85202	0.80557	0.85526	0.92685
	W2vGO	0.86360	0.83227	0.87130	0.92795
	InferSentGO	0.85878	0.75396	0.86402	0.90022
	AicInferSentGO	0.87325	0.80216	0.88138	0.93914
BMA	Resnik	0.89321	0.84995	0.89290	0.94240
	AIC	0.85973	0.83137	0.88504	0.94017
	W2vGO	0.87713	0.84203	0.88308	0.94013
	InferSentGO	0.84890	0.78171	0.87207	0.90712
	AicInferSentGO	0.87292	0.83415	0.89697	0.95296

human/fly orthologs

sizes of the positive set and negative set to be equal. For human/mouse dataset, we have 10,269 pairs for each set; for the human/fly dataset, we have 4932 pairs for each set.

Inherently, human genes are more similar together than genes from different species, making the classification in section 3.1 more difficult. Hence, the AUCs in Table 2 are

lower than those in Table 3.

Like before, we consider the best GO metric to be the one with the highest AUC for the joint analysis. AicInferSentGO is again the winner in this experiment (Table 3). We also conduct the classification for mouse/fly orthologs. In the joint analysis, BMA+AicInferSentGO metric gives the highest AUC score 91.52% (full table not shown).

4 Discussion

In our results, we do not aim to attain perfect classification; rather, we use the classification to rank the GO metrics. Other papers have used sequence similarity and co-expression data to evaluate GO metrics. However, sequence similarity correlates only with MF terms, and co-expression data works best with BP and CC ontology^[11,15,20]. Moreover, genes are expressed non-uniformly across different tissues^[3]. Depending on the data source, experiments using co-expression data can give highly varying outcomes.

The Word2vec has an extension Sentence2vec that converts a sentence into a vector^[6]. Theoretically, one can convert GO definitions into vectors. However, our Word2vec result contains 986,615 words; so, the number of sentences in the training dataset is exponentially larger than this number. We encounter computer memory problem in training Sentence2vec. Therefore, we opt for the InferSent model instead.

Arguably, the *entailment* relation in InferSent does not necessarily equate to a perfect similarity measurement. For example, one can argue that every term in the BP ontology *entails* the root node *biological processes*. Moreover, the NLP approaches in this paper are yet to fully recognize chemical equations. An expression like $2H_2 + O_2$ may not be seen as strictly equal to $2H_2O$ or the word *water*. For these reasons, we view NLP methods as ways to refine existing node-based GO metrics. In this paper, we have seen that InferSentGO improves the AIC scores. Moreover, InferSentGO does not need to be paired with AIC; it can work with any GO similarity metric that gives scores in the range $[0, 1]$.

Past methods to compare GO terms strictly rely on the GO tree, focusing on the fraction of shared common ancestors. Instead, we focus on the GO definitions themselves and apply neural network based NLP approaches. NLP methods can work together with node-based models to achieve higher accuracy. To our knowledge, this paper is the first to apply neural network based NLP techniques to compare the semantic meaning of GO terms. Our application suggests that there are great promises in developing NLP methods for this research area.

Funding

DD is supported by NIH-NLM National Cancer Institute T32LM012424. D.D. and E.E. are supported by National Science Foundation grants 0513612, 0731455, 0729049, 0916676, 1065276, 1302448, 1320589 and 1331176, and National Institutes of Health grants K25-HL080079, U01-DA024417, P01-HL30568, P01-HL28481, R01-GM083198, R01-ES021801, R01-MH101782 and R01-ES022282. WUA and KWC are supported by National Science Foundation Grant IIS-1657193. J.J.L. is supported by NIH/NIGMS R01GM120507 and NSF DMS-1613338.

Contributions

DD and JJJ came up with the problem. DD downloaded and cleaned the data. DD, JJJ, WUA and KWC came up with the solution. DD and WUA coded the solution. DD wrote the paper. JJJ gave advice on database to be used. DD and JJJ discussed the results. DD, JJJ, WUA, KWC and EE read the paper.

References

- [1] Conneau, A., Kiela, D., Schwenk, H., et al. (2017). Supervised learning of universal sentence representations from natural language inference data. arXiv preprint arXiv:1705.02364.
- [2] Dubuisson, M.-P. and Jain, A. K. (1994). A modified hausdorff distance for object matching. In Pattern Recognition, 1994. Vol. 1-Conference A: Computer Vision & Image Processing., Proceedings of the 12th IAPR International Conference on, volume 1, pages 566–568. IEEE.
- [3] Duong, D., Gai, L., Snir, S., et al. (2017). Applying meta-analysis to genotype-tissue expression data from multiple tissues to identify eqtls and increase the number of egenes. Bioinformatics, **33**(14), i67–i74.
- [4] Gene Ontology Consortium (2017). Expansion of the gene ontology knowledgebase and resources. Nucleic acids research, **45**(D1), D331–D338.
- [5] Jiang, J. J. and Conrath, D. W. (1997). Semantic similarity based on corpus statistics and lexical taxonomy. arXiv preprint cmp-lg/9709008.
- [6] Le, Q. and Mikolov, T. (2014). Distributed representations of sentences and documents. In Proceedings of the 31st International Conference on Machine Learning (ICML-14), pages 1188–1196.
- [7] Levy, O., Goldberg, Y., and Ramat-Gan, I. (2014). Linguistic regularities in sparse and explicit word representations. In CoNLL, pages 171–180.
- [8] Li, Y., McLean, D., Bandar, Z., et al. (2006). Sentence similarity based on semantic nets and corpus statistics. IEEE Transactions on Knowledge and Data Engineering, **18**(8), 1138–1150.
- [9] Lin, D. (1998). An information-theoretic definition of similarity.
- [10] Mazandu, G. K. and Mulder, N. J. (2012). A topology-based metric for measuring term similarity in the gene ontology. Advances in bioinformatics, **2012**.
- [11] Mazandu, G. K. and Mulder, N. J. (2014). Information content-based gene ontology functional similarity measures: Which one to use for a given biological data type? PLoS ONE, **9**(12), e113859.
- [12] Mazandu, G. K., Chimusa, E. R., and Mulder, N. J. (2016). Gene ontology semantic similarity tools: survey on features and challenges for biological knowledge discovery. Briefings in Bioinformatics, page bbw067.

- [13] Mikolov, T., Sutskever, I., Chen, K., et al. (2013). Distributed representations of words and phrases and their compositionality. In Advances in neural information processing systems, pages 3111–3119.
- [14] Pennington, J., Socher, R., and Manning, C. (2014). Glove: Global vectors for word representation. In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP), pages 1532–1543.
- [15] Pesquita, C., Faria, D., Bastos, H., et al. (2008). Metrics for go based protein semantic similarity: a systematic evaluation. BMC bioinformatics, 9(5), S4.
- [16] Pesquita, C., Faria, D., Falcão, A. O., et al. (2009). Semantic similarity in biomedical ontologies. PLoS Computational Biology, 5(7), e1000443.
- [17] Rehurek, R. and Sojka, P. (2011). Gensim–python framework for vector space modelling. NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic.
- [18] Resnik, P. (1999). Semantic similarity in a taxonomy: An information-based measure and its application to problems of ambiguity in natural language. J. Artif. Intell. Res.(JAIR), 11, 95–130.
- [19] Rong, X. (2014). word2vec parameter learning explained. arXiv preprint arXiv:1411.2738.
- [20] Song, X., Li, L., Srimani, P. K., et al. (2014). Measure the semantic similarity of GO terms using aggregate information content. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 11(3), 468–476.
- [21] Tuan, L. A., Kim, J.-j., and Ng, S.-K. (2013). Gene ontology concept recognition using cross-products and statistical methods. In BioCreative Challenge Evaluation Workshop vol., page 174.
- [22] Wang, J. Z., Du, Z., Payattakool, R., et al. (2007). A new method to measure the semantic similarity of go terms. Bioinformatics, 23(10), 1274–1281.