

Chalmers University of Technology, Gothenburg, Sweden ² KTH Royal University of Technology, Stockholm, Sweden ³ Cambridge University, Cambridge, United Kingdom CHALMERS ⁴ Alan Turing Institute, London, United Kingdom **UNIVERSITY OF TECHNOLOGY**

Overview

We construct a Knowledge Graph (KG) about the yeast Saccharomyces cerevisiae combining information from various sources. Based on ontologies defining class hierarchies for the terms used box embeddings are found for the nodes in the graph. We demonstrate that Graph Neural Networks (GNNs), combining the KG and box embeddings can be used to predict the effect of gene deletions in S. cerevisiae. Insights from these models can drive scientific discovery, for example through the embeddings generated.

Predicting yeast properties from knowledge graphs using

Filip Kronström¹, Daniel Brunnsåker¹, Alexander H. Gower¹, levgeniia A. Tiukova¹² and Ross D. King¹³⁴



Knowledge Graph

A Knowledge Graph (KG) is created using data from Saccharomyces Genome Database, BioCyc, and The Cell Map. This KG contains information about genes and reactions in the yeast S. cerevisiae expressed in controlled vocabularies.

GNNs and box embeddings



Box embeddings

The terms in the KG are expressed as classes from various ontologies defining class hierarchies. We represent the nodes (classes) in our KG as axis-aligned hyperrectangles which are found by optimizing the conditional probabilities for class inclusions [2]. The classes in our KG can be divided into nine distinct class types for which separate box embeddings are found.

Interaction relations between genes are accompanied by a real number describing the growth when these genes are knocked out [1].



A class inclusion axiom, $A \sqsubseteq B$, corresponds to P(B|A) = 1 which is calculated as

$$P(B|A) = \frac{\text{Vol}(\text{Box}_A \cap \text{Box}_B)}{\text{Vol}(\text{Box}_A)}$$

 $A \sqsubseteq B$













Graph Neural Network



Scientific discovery

Heterogeneous GNNs are used to propagate relation information in the KG predict properties given the node embeddings and edges in the KG.

During the GNN training class box embeddings are altered, but penalised when not adhering to the class hierarchies from the ontologies.

These models can be used for, e.g., gene function prediction (link prediction) or, as shown below, growth prediction for double gene knock-outs.

Embeddings generated by the GNNs combine information from box representations and graph neighbours. Below clusters in the embeddings of genes and their interactions (the Hadamard product between two genes) are shown.

Embeddings can be useful tools for scientific discovery as they for instance allow for comparison of new knowledge to the facts asserted in the KG.



Costanzo, et al. 2016. A global genetic interaction network maps a [1] wiring diagram of cellular function. Science 353.

Dasgupta, et al. 2020. Improving local identifiability in [2] probabilistic box embeddings. NeurIPS '20.

