



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA



INFORMATION ENGINEERING DEPARTMENT  
MASTER'S DEGREE IN COMPUTER ENGINEERING

**APEROL from Networks:  
Analyzing Pipeline and Embedding Representations  
for Optimized Learning (from Networks)**

**Professor**  
Fabio Vandin

**Students**  
Marco Annunziata, 2160851  
Silvia Mondin, 2141201  
Sveva Turola, 2160852

ACADEMIC YEAR 2025-2026

# Contents

1	Motivation . . . . .	1
2	Datasets . . . . .	1
3	Methods . . . . .	1
4	Experiments . . . . .	1
5	References . . . . .	2
6	Contribution of Authors and AI Usage . . . . .	2

# 1 Motivation

Graphs have emerged as a natural model for the representation and analysis of data and complex systems in various domains. For instance, link prediction in social networks can help to understand the spreading process of rumors or epidemics [4], and in biological networks, it is commonly used for predicting novel interactions between proteins [2].

However, most traditional machine learning algorithms are not designed to work directly with graph-structured data, as they require numeric vectors or matrices as inputs. To address this challenge, graph embedding techniques have been developed to learn low-dimensional vector representations of nodes, edges, or entire graphs while preserving their topological properties.

Consequently, in order to design an effective graph embedding technique, it is necessary to consider several criteria [1]. Two of these are adaptability and scalability. An embedding is considered adaptable if it can be utilized for various data and tasks without the need for retraining. The scalability of an embedding is determined by its capacity to process large-scale networks within a reasonable amount of time.

The objective of this study is to evaluate known embeddings from the perspective of these two criteria.

From the perspective of adaptability, the objective is to determine the potential impact of the graph domain on the precision of an embedding. In other words, the objective is to ascertain whether there exist approaches that are better suited for a particular field, whether there is one embedding that consistently outperforms the others, or if the embeddings are comparable. Eventually, we would like to assess the adaptability of these embeddings when they are applied to different tasks.

In the context of scalability, the objective is to identify the embeddings that are better suited to process large graphs. The objective of this study is to evaluate the tradeoff between the accuracy of the predictions and the time required for training and inference.

## 2 Datasets

The datasets that will be used in this project are nine and they are divided into three categories: Road Networks, Social Networks, and Biological Networks. For each of these categories three datasets of different size were chosen: one small (~40k nodes), one medium (~100k nodes), and one large (~1M nodes). The datasets, with their main properties (number of nodes, number of edges, and type of graph) are reported in the table 1.

Network	$ V $	$ E $	Type
Pennsylvania [5]	1.088.092	1.541.898	Undirected
Padua (province)	122.680	164.737	Directed
Hong Kong (city)	43.620	91.542	Directed
Italian Covid-19 Retweet Network	0	800.000	Directed
Twitch	168.114	6.797.557	Undirected
GitHub	37.700	289.003	Undirected
Mus Musculus Protein Interactions	0	800.000	Undirected
Saccharomyces cerevisiae Protein Interactions	0	100.000	Undirected
Bio-grid-fission-yeast	2000	25.300	Undirected

Table 1: Datasets used in the project

## 3 Methods

## 4 Experiments

To evaluate the methods used in this project, two different evaluation measures will be used: the Area Under the Receiver Operating Characteristic curve (AUROC), and the Area Under the Precision-Recall curve (AUPR). AUROC is historically considered the primary performance metric used to evaluate the performances of Link Prediction methods [3]. However, AUROC favors accurate classification of positive examples, at the cost of misclassifying the negative ones. In a scenario like Link Prediction problem, which is inherently skewed towards the negative class, it may not be

suitable and can overestimate the performance of the methods. For this reason, the AUPR curve was also selected, as it can provide better evaluation of Link Prediction in the presence of class imbalance.

## 5 References

- [1] Anthony Baptista et al. “Zoo guide to network embedding”. In: *Journal of Physics: Complexity* 4.4 (2023), p. 042001.
- [2] Björn H Junker and Falk Schreiber. *Analysis of biological networks*. Vol. 2. Wiley Online Library, 2008.
- [3] I. Bhargavi Kalyani, A. Rama Prasad Mathi, and Niladri Sett. “Evaluating link prediction: new perspectives and recommendations”. In: *International Journal of Data Science and Analytics* 20.7 (2025), pp. 6855–6886. ISSN: 2364-4168. DOI: 10.1007/s41060-025-00858-0. URL: <https://doi.org/10.1007/s41060-025-00858-0>.
- [4] Romualdo Pastor-Satorras et al. “Epidemic processes in complex networks”. In: *Reviews of modern physics* 87.3 (2015), pp. 925–979.
- [5] Ryan A. Rossi and Nesreen K. Ahmed. “The Network Data Repository with Interactive Graph Analytics and Visualization”. In: *AAAI*. 2015. URL: <https://networkrepository.com>.

## 6 Contribution of Authors and AI Usage