

FACULDADE DE ENGENHARIA DA UNIVERSIDADE DO PORTO

Graph Reinforcement Learning for Improving Smart Grid Services

António Bernardo Linhares Oliveira



Mestrado em Engenharia Informática e Computação

Supervisor: Prof. António Costa & Prof. Rosaldo Rossetti

May 23, 2024

Graph Reinforcement Learning for Improving Smart Grid Services

António Bernardo Linhares Oliveira

Mestrado em Engenharia Informática e Computação

May 23, 2024

Resumo

Os grafos são representações que descrevem problemas e os seus objectos em domínios orientados para as redes, como as redes eléctricas, os transportes ou as redes sociais. Estas representações podem captar não só os conceitos e as suas respectivas propriedades, mas também as relações entre esses conceitos, resultando em estruturas de dados frequentemente complexas e esparsas que são especialmente úteis para representar problemas em que a topologia de uma rede desempenha um papel importante. No contexto da utilização destas estruturas em algoritmos de aprendizagem computacional, o seu desempenho torna-se dependente não só da sua conceção e seleção dos atributos relevantes e parâmetros, mas também das representações subjacentes utilizadas para captar a essência das estruturas em grafo.

A Aprendizagem por Reforço em Grafos ou *Graph Reinforcement Learning* (GRL) é um tópico que tem merecido grande atenção por parte dos académicos nos últimos anos. Ao permitir que as técnicas de Aprendizagem por Reforço aprendam e otimizem processos de decisão sequenciais em ambientes baseados em grafos, os sistemas podem ser melhorados de forma a tirar partido das características da topologia dos grafos em domínios de aplicação associados a redes. Com os avanços no final da década de 2010 em Redes Neurais de Grafos, ou *Graph Neural Networks*, na aprendizagem e extração de representações eficientes de grafos, foram propostos métodos mais sofisticados de GRL e o tópico começou a atrair mais curiosidade dos académicos. Embora, nos últimos anos, muito trabalho tenha sido feito nesta área, a pesquisa à volta do GRL ainda é considerada estar em fase inicial.

Além disso, considerando os actuais desafios globais associados à sustentabilidade e aos sistemas energéticos, há uma necessidade crescente de avanços em sistemas inteligentes focados em modernizar as redes de distribuição e transmissão de energia. Atualmente, as fontes de energia renováveis desempenham um papel importante na redução da dependência dos combustíveis fósseis, o que altera a topologia dos sistemas de distribuição de energia à medida que os consumidores adquirem a capacidade de gerar energia renovável. Com as melhorias na Inteligência Artificial e Aprendizagem Computacional, os sistemas podem ser adaptados à descentralização da produção e gerir eficientemente a monitorização, distribuição e transmissão da energia. Neste trabalho, a ênfase principal reside na melhoria dos algoritmos de GRL que serão aplicados no contexto do problema da distribuição dinâmica e económica de energia como seu principal domínio de aplicação, considerando fontes de energia renováveis e sistemas de armazenamento de energia.

Desta forma, esta dissertação tem como objetivo fazer avançar a investigação existente sobre técnicas de Aprendizagem por Reforço em Grafos através de: (1) realizar uma revisão exaustiva da literatura recente relativa às várias abordagens de GRL propostas e de sistemas de distribuição dinâmica e económica de energia, de modo a obter uma perspetiva global das técnicas recentes mais avançadas e das suas limitações; (2) realizar um estudo empírico comparativo e sistemático das diferentes técnicas de GRL no problema de distribuição dinâmica e económica, considerando cenários de estudo de caso de diferentes dimensões modelados por uma simulação de uma rede de distribuição de energia eléctrica; (3) propor um modelo que melhor integre as capacidades das

técnicas de *Deep Reinforcement Learning* e *Graph Neural Networks*, com base nos resultados do estudo empírico e com melhorias no desempenho e escalabilidade face aos modelos propostos pela literatura.

Abstract

Graphs are structures that depict problems and their objects in network-oriented domains such as power grids, transport or social networks. These representations can capture not only the concepts and their respective properties but the intricate relationships between those concepts, resulting in often complex and sparse data structures that are especially useful for representing problems where network topology plays a major role. In the context of using these structures in machine learning algorithms, their performance becomes not only dependent on its design and the selection of relevant features and parameters, but also on the underlying representations used to capture the essence of graph structures.

Graph Reinforcement Learning (GRL) is a topic that has earned significant attention from academics in the last few years. By enabling Reinforcement Learning techniques to learn and optimize sequential decision-making processes in graph-based environments, systems can be improved and gain the ability to leverage graph topology features in network-oriented application domains. With the advancements in the late 2010s on Graph Neural Networks on learning how to extract efficient graph representations from a given scenario, more sophisticated methods of GRL were proposed and the topic started to attract the curiosity of scholars. Although, in recent years, a lot of work has been done in this area, research on techniques is still considered to be in an early stage.

Furthermore, considering the current global challenges associated with sustainability and energy systems, there is an increasing need for advancements in energy-focused intelligent systems to modernize the current power grids. In the present, renewable energy sources play a major role in reducing the reliance on fossil fuels, which changes the topology of energy distribution systems as consumers gain the capability to generate renewable power. With the improvements in Artificial Intelligence and Machine Learning, systems can be adapted to the decentralization of energy production and efficiently manage the monitoring, distribution and transmission of energy systems. In this work, the primary emphasis lies on improving GRL algorithms which will be applied to solve the dynamic economic power dispatch problem as its main application domain, considering renewable energy sources and energy storage systems.

In this manner, this dissertation aims to advance the existing research on Graph Reinforcement Learning techniques by: (1) conducting a thorough review of the recent literature regarding various proposed GRL approaches and Dynamic Economic Dispatch Systems to gain an overall perspective of the recent state-of-the-art techniques and their limitations; (2) performing a comparative and systematic empirical study on the different GRL techniques on the Dynamic Economic Power Dispatch problem, considering case study scenarios of different sizes modelled by a power distribution grid simulation (3) propose a model that better integrates the capabilities of Deep Reinforcement Learning Agents with Graph Neural Networks, based on the results of the empirical study, with improvements in performance and scalability.

Keywords: Graph Reinforcement Learning, Graph Neural Networks, Deep Reinforcement Learning, Smart Grid, Dynamic Economic Dispatch

ACM Classification: Computing Methodologies → Machine Learning → Learning Paradigms → Reinforcement Learning

Acknowledgements

To my supervisors for all the guidance and teachings during the development of this work
To my parents who always supported me and expected me to become the best version of myself

António Oliveira

*“Man is not worried by real problems
so much as by his imagined anxieties about real problems”*

Epictetus

Contents

1	Introduction	1
1.1	Context	1
1.2	Motivation	1
1.3	Objectives	2
1.4	Report Structure	2
2	Background Knowledge	3
2.1	Artificial Neural Networks	3
2.1.1	Feedforward Neural Networks	3
2.2	Reinforcement Learning	5
2.2.1	Markov Decision Process	5
2.2.2	Rewards and Returns	7
2.2.3	Policies and Value Functions	8
2.2.4	Types of Reinforcement Learning (RL)	8
2.3	Graph Representation Learning	10
2.4	Graph Neural Networks	10
2.4.1	Graph Convolutional Network	11
2.4.2	Graph Attention Network	13
2.5	Smart Grid Services	15
3	Literature Review	16
3.1	Research Methodology	16
3.2	Graph Reinforcement Learning Approaches	17
3.2.1	Plain GCN-Based GRL	17
3.2.2	Attention-based GRL	19
3.2.3	Other Approaches	20
3.3	Dynamic Economic Dispatch Systems	21
3.4	Conclusions	22
4	Methodological Approach	24
4.1	Problem Statement	24
4.1.1	Scope	25
4.2	Problem Formalization	25
4.2.1	Dynamic Economic Dispatch	25
4.2.2	Markov Decision Process (MDP)	28
4.3	Requirements	28
4.3.1	Functional	28
4.3.2	Non-functional	29

4.3.3	Structural	29
4.4	Architecture	29
4.5	Methodology	30
4.5.1	Evaluation Metrics	32
4.6	Work Plan	33
5	Experimental Setup	36
5.1	Simulation Environment	36
5.2	Scenarios	36
5.3	Experiments	36
5.3.1	GRL	36
5.3.2	DRL vs. GRL	36
5.3.3	Scalability	36
5.3.4	Adaptability	36
6	Conclusions	37
6.1	Expected Contributions	37
6.2	SWOT Analysis	38
6.3	SMART Analysis	38

List of Figures

2.1	The Perceptron [?]	4
2.2	Architecture of a Feedforward Neural Network [?]	4
2.3	Agent-environment interaction in a MDP [?]	6
2.4	Taxonomy of algorithms in modern RL [?]	9
2.5	Smart Grid Capabilities Pyramid [?]	15
4.1	Solution Architecture	30
4.2	Grid2Op <i>l2rpn_idf_2023</i> 118-bus test case [?]	31
4.3	Project Work Plan	35
6.1	SWOT Analysis	38

List of Tables

3.1	GCN-Based GRL Literature	19
3.2	Attention-Based GRL Literature	20
3.3	Other GRL Approaches	21
3.4	Dynamic Economic Dispatch RL Systems	22
4.1	Functional Requirements	28
4.2	Non-functional Requirements	29
4.3	Structural Requirements	29
4.4	Test Case Sizes	32

Acronyms and Abbreviations

IT	Information Technology
ANN	Artificial Neural Network
MLP	Multilayer Perceptron
CNN	Convolutional Neural Network
RL	Reinforcement Learning
MDP	Markov Decision Process
DRL	Deep Reinforcement Learning
GNN	Graph Neural Network
GRL	Graph Reinforcement Learning
SAC	Soft Actor-Critic
DDPG	Deep Deterministic Policy Gradient
DQN	Deep Q-Network
PPO	Proximal Policy Optimization
GCN	Graph Convolutional Network
GAT	Graph Attention Network
GIN	Graph Isomorphism Network
ADN	Active Distribution Network
DED	Dynamic Economic Dispatch
ESS	Energy Storage System
RES	Renewable Energy Source

Chapter 1

Introduction

In this introductory chapter, the context and motivation regarding this dissertation, as well as its key objectives are presented in sections 1.1, 1.2 and 1.3, respectively. Additionally, the structure of this report is exposed and its logical divisions are described in section 1.4.

1.1 Context

Several real-world problems and their objects can be instinctively represented by graph structures. These representations not only capture the main properties of a given domain but also the intricate topology of relationships in a network-oriented problem. Graph representations are often sparse and complex and to appropriately leverage their various topology features machine learning algorithms require underlying methods to efficiently generalize and produce adequate representations from these structures, considering the trade-off data completeness and computational efficiency.

In the case of sequential decision-making problems, the same is verified. Learning how to map good sequences of decisions in network-oriented domains can depend, in some cases, on accounting for the environment topological features [?, ?, ?, ?]. As the main paradigm of machine learning that addresses sequential decision-making problems, RL algorithms need to be adapted to reflect these considerations, which establishes the foundations of Graph Reinforcement Learning (GRL). This comprises the main focus of this work, which will be applied in the application domain of *Smart Grid* Services.

In other regards, this report is written in the context of the course of Dissertation Preparation (PDIS) inserted in the Master's Degree in Informatics and Computing Engineering (MEIC) of the Faculty of Engineering of the University of Porto (FEUP). In addition, this project is accommodated in the Artificial Intelligence and Computer Science Laboratory (LIACC).

1.2 Motivation

Reflecting on the current global issues associated with the energetic crisis and climate change, there is an increasing need for sustainable and economic energetic systems to modernize the cur-

rent power grids. This modernization is translated into the transition to the *smart grid*, a power grid equipped with intelligent control and monitoring systems to efficiently manage power distribution [?, ?], voltage regulation, system restoration [?], grid reliability [?] or other associated processes. Currently, renewable energy sources play a major role in reducing the reliance on fossil fuels [?], which changes the topology of energy distribution systems as consumers gain the capability to generate renewable power. Furthermore, investment in energy storage is becoming a priority in the energy sector [?], resulting in improvements on storage capacity and approximating the current solutions to the average consumer [?].

The exposed issues serve as the prime motivation for performing this work on the application domain of *smart grid* services, with the expectation that by proposing concrete improvements in GRL techniques a well-performing solution to the dynamic economic dispatch problem can be presented. Beyond this, we hope the proposed solution and its architecture is also adaptable and applicable to other smart grid problems, resulting in a significant contribution to these services.

Furthermore, the main reasons for addressing GRL algorithms lies on their novelty and complexity, the lack of well-documented literature regarding these approaches, and the need for systematic and comparative studies confronting the different proposed techniques and architectures.

1.3 Objectives

Considering this work's context and motivations, we define its main objectives:

1. Perform a review of literature regarding GRL approaches and Dynamic Economic Dispatch (DED) systems
2. Conduct a comparative and systematic empirical study of different GRL solutions of the DED problem
3. Propose a GRL model and concrete improvements facing the literature proposed models

The first goal addresses the analysis of existent techniques, as well as its limitations, by reviewing the relevant research on GRL approaches and DED systems. Secondly, we will focus on implementing the different observed approaches to solve the DED problem and perform a comparative study to analyse and confront the gathered results. Lastly, we hope the accomplishment of the first to goals to enable the proposition of a state-of-the-art GRL model and specific improvements to these techniques.

1.4 Report Structure

This report is organized as follows: (1) an introductory chapter; (2) a chapter explaining the relevant background concepts to perform this study; (3) the review of the literature regarding GRL techniques and DED systems; (4) the statement of the main problem and the presentation of the proposed solution and finally, (5) the main conclusions, reflections and expected contributions of this work.

Chapter 2

Background Knowledge

In this chapter, we will present the underlying concepts related to this dissertation. This knowledge consists in important context for the rest of this work by explaining the its background concepts.

This chapter is divided into section 2.1 addressing Artificial Neural Networks, section 2.2 regarding Reinforcement Learning problem and algorithms, 2.3 explaining Graph Representation Learning, 2.4 exposing the current Graph Neural Network (GNN) approaches and 2.5 addressing Smart Grid Services.

2.1 Artificial Neural Networks

Artificial Neural Networks (ANNs) are a class of machine learning algorithms based on the neural process of biological learning. The simplest form of an ANN is a Multilayer Perceptron (MLP), also called a feedforward network, whose main objective is to approximate to function f that models relationships between input data x and output data y of considerate complexity [?, ?]. It defines a mapping $y = f(x; \theta)$ and learns the best composition of parameters θ to approximate it to the unknown model. The MLP serves as a fundamental part of developing the other more complex types of neural networks [?].

2.1.1 Feedforward Neural Networks

The main building block of a MLP is the *Perceptron*, pictured in figure 2.1, a simple computational model initially designed as a binary classifier that mimics biological neurons' behaviour [?]. A neuron might have many inputs x and has a single output y . It contains a vector of *weights* $w = (w_1 \dots w_m)$, each associated with a single input, and a special weight b called the *bias*. In this context, a perceptron defines a computational operation formulated as equation 2.1 portrays [?].

$$f(x) = \begin{cases} 1 & \text{if } b + \mathbf{w} \cdot \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.1)$$

Functions that compute $b + \mathbf{w} \cdot \mathbf{x} > 0$ are called *linear units* and are identified with Σ [?, ?]. An activation function g was introduced to enable the output of non-linear data. The default

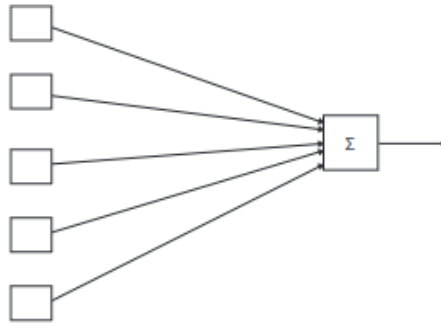


Figure 2.1: The Perceptron [?]

recommendation is the *Rectified Linear Unit (ReLU)*, with the Logistic curve (sigmoid) also being very common [?].

Feedforward networks are composed of an input layer formed by the vector of input values, an arbitrary number of hidden layers and an output layer, which is the last layer of neurons [?]. The greater the amount of layers the higher the *depth* of the network [?, ?].

On its own, the model amounts only to a complex function. Still, with real-world correspondence between input values and associated outputs, a feedforward network can be trained to approximate the unknown function of the environment. In more concrete terms, this involves updating all of the different weight and bias values of each neuron to achieve an output as close as possible to the real or desired value or minimize the total loss, which indicates how distant the network model is to the real function to approximate [?, ?]. *Loss functions* are used to calculate this value.

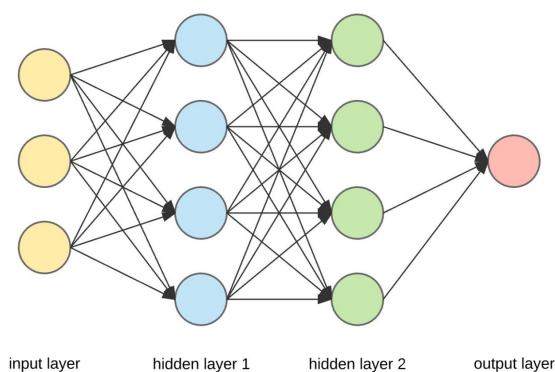


Figure 2.2: Architecture of a Feedforward Neural Network [?]

2.2 Reinforcement Learning

RL consists of a field and a class of machine learning algorithms that study how to learn to take good sequences of actions to achieve a goal associated with a maximizing received numerical reward [?]. The main objective is to maximize the received cumulative reward by trying between the available actions and discovering which ones yield the most reward [?]. This sequential decision-making process becomes more complex when a delayed reward is considered, given that an action with immediate reward may not always reflect the delayed consequences of that decision [?]. It's also the learner's job to consider this during the learning process. These concepts of *delayed reward* and *trial-and-error search* make up the most important characteristics of Reinforcement Learning [?]. The classic formalisation of this problem is the MDP through defining the agent-environment interaction process, explained in the following subsection 2.2.1.

A major challenge in this machine learning paradigm is the trade-off between *exploring* new unknown actions and *exploiting* the already known "good" actions [?]. To choose the sequence of actions that return the highest reward, the agent must choose actions it found effective in similar past situations or **exploit** what it learned from experience. Furthermore, given that the agent may not know the action-reward mappings initially, it has to *explore* possible actions that were not selected previously or may initially seem to yield a low reward to compute accurate reward estimates. The main problem is that neither exploitation nor exploration can be favoured exclusively without failing at the task [?]. Additionally, an agent's environment is uncertain, and changes in the environment's dynamics may also involve re-estimating action rewards.

In conclusion, RL techniques enable the implementation of sequential decision-making agents that seek to maximize a reward signal analogous to an explicit (complex) goal. The agents need to balance between actions that yield a reward on posterior time steps and actions that produce immediate rewards. In addition, these agents are also faced with the task of balancing the exploitation of information from past experiences and the exploration of new decision paths that could potentially return a higher reward down the road [?].

2.2.1 Markov Decision Process

Markov Decision Processes (MDPs) are a classical formalization of a sequential decision-making process, constituting the mathematical definition of the RL problem [?, ?]. Beyond estimating potential rewards for the available actions, the problem defined by MDPs involves learning which actions are optimal in specific situations, i.e. learning a mapping between states of the environment and actions [?].

The central component of MDPs is the agent, which acts as a decision-maker and learns from interactions with the environment it's inserted. In a continuous process, the agent takes actions that affect the environment's state, which in turn presents new situations [?]. The environment also responds with the reward signals which the agent aims to maximize over time through its decision process.

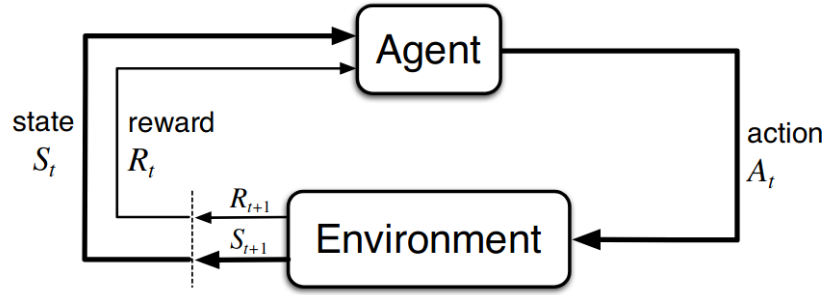


Figure 2.3: Agent-environment interaction in a MDP [?]

Formally, the agent-environment interactions, as figure 2.3 entails, occur in a sequence of discrete time steps t , where at each step, the agent receives a representation of the state of the environment $S_t \in \mathcal{S}$ which is used to select an appropriate action $A_t \in \mathcal{A}(s)$, where \mathcal{S} is the set of possible states called the *state space* and $\mathcal{A}(s)$ is the set of available actions for state s [?, ?]. In the next step, the agent receives, as a consequence of its decision, a numerical reward signal $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$ and is faced with a new state S_{t+1} [?]. Ultimately, the MDP agent follows a logical sequence that occurs as equation 2.2 states. The collection of a state S_t , action taken A_{t+1} , reward R_{t+1} received and next state S_{t+1} constitutes an *experience tuple* [?].

$$S_0, A_0, R_1, S_1, A_2, R_2, S_2, A_3, R_3, \dots \quad (2.2)$$

In addition, when the set of possible actions, states and rewards (\mathcal{A} , \mathcal{S} and \mathcal{R}) are finite, the MDP is said to be *finite* [?]. This results in S_t and R_t having well-defined discrete probability distributions in function of the preceding state and chosen action [?]. Therefore, the probability of receiving a particular reward and state given the previous state and selected action, which characterizes a finite MPD's dynamics, may be characterized by function p defined in equation 2.3

$$p(s', r | s, a) \doteq \Pr\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\} \quad (2.3)$$

For all $s, s' \in \mathcal{S}$, $r \in \mathcal{R}$ and $a \in \mathcal{A}(s)$, where \doteq denotes a mathematical formal definition. This encompasses the assumption that the probability of each possible state, S_t , and reward, R_t , pair is only dependent on the preceding state, S_{t-1} , and action taken, A_{t-1} [?]. Instead of observing this as a restriction on the decision process, it's more convenient to view it as a constraint on the state variable, considering that it must contain all the necessary information from experience to make a valuable decision in the immediate step. If this condition is satisfied, the state is declared to have the *markov property* [?].

From function p in equation 2.3, the state-transition probabilities, also called the *transition function*, can be computed as described by equation 2.4 [?, ?].

$$p(s' | s, a) \doteq \Pr\{S_t = s' | S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in \mathcal{R}} p(s', r | s, a) \quad (2.4)$$

In addition, the expected rewards can be calculated for state-action pairs (equation 2.5) or state-action-next-action triples (equation 2.6) [?, ?].

$$r(s, a) \doteq \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a) \quad (2.5)$$

$$r(s, a, s') \doteq \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in \mathcal{R}} r \frac{p(s', r | s, a)}{p(s' | s, a)} \quad (2.6)$$

2.2.2 Rewards and Returns

As stated in the previous subsections, the main goal of a RL agent defined by the numeric reward signal, $R_t \in \mathbb{R}$, it receives from the environment [?]. In this context, the agent's objective is to maximize the total reward it receives, considering not only immediate but also the cumulative reward over time. In the ubiquitous work of [?], the *reward hypothesis* is stated as follows:

That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward). [?]

This also entails that the process of reward maximization from the agent has to be closely tied to it achieving its defined goals in a practical sense. Otherwise, the agent will fail at fulfilling the desired objectives [?].

Formally, the goal of an RL agent can be defined by the maximization of the cumulative reward received of time called the *expected return*, Return_t [?].

$$\text{Return}_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \cdots + R_T \quad (2.7)$$

T describes the final time step. This definition can be applied in domains with a natural notion of a terminal state or final time step. In these cases, the agent-environment interaction process can be broken into logically independent subsequences called *episodes* [?]. Each episode ends in a special state, called the terminal state, restarting a new sequence of states and actions completely independent from the previous episode [?]. In this context, episodes can be considered to end in the same terminal state, with different accumulated rewards for the different outcomes [?].

In contrast, there are situations where the decision-making process doesn't divide itself into logically identifiable episodes but goes on indefinitely. In this case, $T = \infty$ and according to equation 2.7, the expected return the agent aims to maximize would be infinite [?]. In this manner, another concept is added in the expected return definition called the *discount rate*, γ where $0 \leq \gamma \leq 1$, representing how strongly the agent should account for future rewards in the expected return calculations, as equation 2.8 [?].

$$\text{Return}_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (2.8)$$

From this equation, we can compute the expected discounted return on a given time step t in the function of the immediate reward signal received and the expected return for the next time step $t + 1$, which eases the job of calculating expected returns for reward sequences [?]. This is entailed by equation 2.9.

$$\text{Return}_t = R_{t+1} + \gamma G_{t+1} \quad (2.9)$$

In this manner, a MDP can be defined by a tuple with a state space \mathcal{S} , an action space \mathcal{A} , a transition function p , a reward function r and a discount factor γ , as equation 2.10 portrays [?].

$$M = (\mathcal{S}, \mathcal{A}, p, r, \gamma) \quad (2.10)$$

2.2.3 Policies and Value Functions

RL techniques typically involve the estimation of what is understood as *value functions*, functions that estimate the expected return based on the current state value or state-action pair. This characterizes how good is for an agent to be in a specific state or to take an action in a specific state, respectively, using the expected return to characterize the overall *goodness* of these scenarios [?, ?]. These functions are tied to a specific way of determining the action in a given state. Formally, this is defined as a *policy* π , that defines the probability $\pi(a|s)$ of taking action a in state s [?]. In this context, the *state-value function*, $v_\pi(s)$ and *action-value functions* $q_\pi(s, a)$ for policy π can be defined by equations 2.11 and 2.14, respectively [?].

$$v_\pi(s) \doteq \mathbb{E}_\pi[\text{Return}_t | S_t = s], \forall s \in \mathcal{S} \quad (2.11)$$

$$q_\pi(s, a) \doteq \mathbb{E}_\pi[\text{Return}_t | S_t = s, A_t = a] \quad (2.12)$$

The utility of such functions rely on the possibility of estimating them with regard to past experience of the agent [?]. A fundamental property of value functions is that it can, as was the case with the expected return (equation 2.9), satisfy recursive relationships with the next immediate value as equation 2.13 entails [?]. This equation is called the *Bellman equation for v_π* , which characterizes the relationship between the value of current and subsequent states.

$$v_\pi \doteq \sum_{a \in \mathcal{A}(s)} \pi(a|s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r | s, a) [r + \gamma v_\pi(s')] \quad (2.13)$$

$$q_\pi(s, a) \doteq \sum \mathbb{E}_\pi[\text{Return}_t | S_t = s, A_t = a] \quad (2.14)$$

2.2.4 Types of RL

Regarding RL algorithms, they can be divided into model-free and model-based techniques [?]. These categories are distinguished by whether an agent uses a provided or learned *model* of the set of transition and reward functions, another optional element of RL techniques [?, ?]. In the positive case, the method is said to be model-based, otherwise, it's model-free. Having an accurate model of the environment allows the RL agent to focus on planning ahead by calculating future scenarios

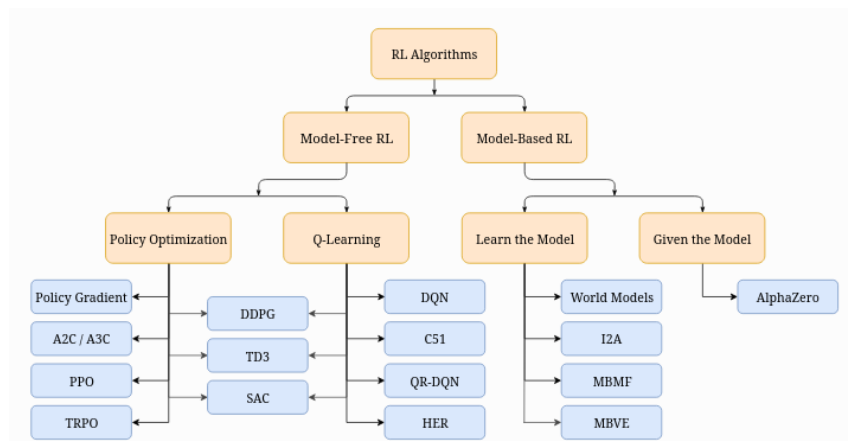


Figure 2.4: Taxonomy of algorithms in modern RL [?]

and creating policies based on the results of the planning process. An example of a famous system of this kind is AlphaZero [?]. However, in most cases, agents can't access a ground-truth model of the environment, leaving only the scenario where an agent learns a model purely from experience. This creates several challenges, the most prominent of which relies on the fact that the model, in most times, doesn't fully capture the environment's transition dynamics, equipping it with bias in relation to the actual dynamics. With this, learning how to generalise the model to real-world environments so that the bias is not over-exploited becomes a very complex task [?]. Model-free algorithms can be also further divided into Q-learning and Policy Approximation techniques.

Furthermore, algorithms can also be subdivided into on-policy and off-policy methods. [?] On-policy algorithms evaluate and improve a single policy used to determine the agent's behaviour [?]. The methods under the policy optimization category such as A2C and A3C [?] or the Proximal Policy Optimization (PPO) [?] almost always fall into this label. In contrast, off-policy algorithms learn how to improve a different target policy based on the results that arise from the policy used to determine the system behaviour initially [?]. Such approaches include Q-Learning algorithms such as Deep Q-Networks (DQNs) [?, ?]. Deep Deterministic Policy Gradient (DDPG) [?] combines policy optimization with q-learning, consisting of an off-policy method that learns both a q-function and a policy. DDPG constitutes the adaption of q-learning methods to continuous action spaces [?]. Another example of an off-policy algorithm is the Soft Actor-Critic (SAC) [?] method, which bridges stochastic policy optimization with the DDPG approach and has entropy regularization as one of its central features, which translates into training a policy that maximizes the expected return and entropy, a measure of randomness in the policy [?].

Lastly, with the advent of deep learning becoming one of the most ubiquitous techniques in machine learning, RL algorithms have evolved beyond the traditional tabular methods [?]. Traditional RL has evolved to Deep Reinforcement Learning (DRL), which studies how to use deep neural networks in RL problems to leverage their generalization abilities for solving more complex problems.

2.3 Graph Representation Learning

Several objects and problems can be naturally expressed in the real world using graphs, such as social networks, power grids, transportation networks, recommendation systems or drug discovery. The usefulness of such representations is tied to how they instinctively represent the complex relationships between objects. However, graph data is often very sparse and complex, and their sophisticated structure is difficult to deal with [?, ?].

Furthermore, the performance of machine learning models strongly relies not only on their design but also on good representations of the underlying information [?]. Ineffective representations, on the one hand, can lack important graph features and, on the other, can carry vast amounts of redundant information, affecting the algorithms' performance in leveraging the data for different analytical tasks [?, ?].

In this context, **Graph Representation Learning** studies how to learn the underlying features of graphs to extract a minimal but sufficient representation of the graph attributes and structure [?, ?, ?]. Currently, the improvements in deep learning allow representation learning techniques consisting of the composition of multiple non-linear transformations that yield more abstract and, ultimately, more useful representations of graph data [?].

2.4 Graph Neural Networks

In the present, deep learning and ANN have become one of the most prominent approaches in Artificial Intelligence research [?]. Approaches such as recurrent neural networks and convolutional networks have achieved remarkable results on Euclidean data, such as images or sequence data, such as text and signals [?]. Furthermore, techniques regarding deep learning applied to graphs have also experienced rising popularity among the research community, more specifically **GNNs** that became the most successful learning models for graph-related tasks across many application domains [?, ?].

The main objective of GNNs is to update node representations with representations from their neighbourhood iteratively [?]. Starting at the first representation $H^0 = X$, each layer encompasses two important functions:

- **Aggregate**, in each node, the information from their neighbours
- **Combine** the aggregated information with the current node representations

The general framework of GNNs, outlined in [?], can be defined mathematically as follows:

Initialization: $H^0 = X$

For $k = 1, 2, \dots, K$

$$\begin{aligned} a_v^k &= \text{AGGREGATE}^k \{H_u^{k-1} : u \in N(v)\} \\ H_v^k &= \text{COMBINE}^k \{H_u^{k-1}, a_v^k\} \end{aligned}$$

Where $N(v)$ is the set of neighbours for the v -th node. The node representations H^K in the last layer can be treated as the final representations, which sequentially can be used for other downstream tasks [?].

2.4.1 Graph Convolutional Network

A **Graph Convolutional Network (GCN)** [?] is a popular architecture of GNNs praised by its simplicity and effectiveness in a variety of tasks [?, ?]. In this model, the node representations in each layer are updated according to the following convolutional operation:

$$H^{k+1} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^k W^k) \quad (2.15)$$

$= A + I$ - Adjacency Matrix with self-connections

$I \in \mathbb{R}^{N \times N}$ - Identity Matrix

\tilde{D} - Diagonal Matrix, with $\tilde{D}_{ii} = \sum_j i_j$

σ - Activation Function

$W^k \in \mathbb{R}^{F \times F'}$ - Laywise linear transformation matrix (F and F' are the dimensions of node representations in the k -th and $(k+1)$ layer, respectively)

$W^k \in \mathbb{R}^{F \times F'}$ is a layerwise linear transformation matrix that is trained during optimization [?]. The previous equation 2.15 can be dissected further to understand the *AGGREGATE* and *COMBINE* function definitions in a GCN [?]. For a node i , the representation updating equation can be reformulated as:

$$H_i^k = \sigma\left(\sum_{j \in \{N(i) \cup i\}} \frac{\tilde{A}_{ij}}{\sqrt{\tilde{D}_{ii} \tilde{D}_{jj}}} H_j^{k-1} W^k\right) \quad (2.16)$$

$$H_i^k = \sigma\left(\sum_{j \in N(i)} \frac{A_{ij}}{\sqrt{\tilde{D}_{ii} \tilde{D}_{jj}}} H_j^{k-1} W^k\right) + \frac{1}{\tilde{D}_i} H_i^{k-1} W^k \quad (2.17)$$

In the second equation, the *AGGREGATE* function can be observed as the weighted average of the neighbour node representations [?]. The weight of neighbour j is defined by the weight of the edge (i, j) , more concretely, A_{ij} normalized by the degrees of the two nodes [?]. The *COMBINE* function consists of the summation of the aggregated information and the node representation itself, where the representation is normalized by its own degree [?].

Spectral Graph Convolutions

Regarding the connection between GCNs and spectral filters defined on graphs, spectral convolutions can be defined as the multiplication of a node-wise signal $x \in \mathbb{R}^N$ with a convolutional filter $g_\theta = \text{diag}(\theta)$ in the *Fourier domain* [?, ?], formally:

$$g_\theta \star x = U_{g_\theta} U^T x \quad (2.18)$$

$\theta \in \mathbb{R}^N$ - Filter parameter

U - Matrix of eigenvectors of the normalized graph Laplacian Matrix $L = I_N - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$

The eigendecomposition of the Laplacian matrix can also be defined by $L = U \Lambda U^T$ with Λ serving as the diagonal matrix of eigenvalues and $U^T x$ is the graph Fourier transform of the input signal x [?]. In a practical context, g_θ is the function of eigenvalues of the normalized graph Laplacian matrix L , that is $g^\theta(\Lambda)$ [?, ?]. Computing this is a problem of quadratic complexity to the number of nodes N , something that can be circumvented by approximating $g_\theta(\Lambda)$ with a truncated expansion of Chebyshev polynomials $T_k(x)$ up to K -th order [?, ?]:

$$g_{\theta'}(\Lambda) = \sum_{k=0}^K \theta'_k T_k(\tilde{\Lambda}) \quad (2.19)$$

$$\tilde{\Lambda} = \frac{2}{\lambda_{\max}} \Lambda - I$$

λ_{\max} - Largest eigenvalue of L

$\theta' \in \mathbb{R}^N$ - Vector of Chebyshev coefficients

$T_k(x)$ - Chebyshev polynomials

$T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x)$ with $T_0(x) = 1$ and $T_1(x) = x$

By combining this with the previous equation, the first can be reformulated as:

$$g_\theta \star x = \sum_{k=0}^K \theta'_k T_k(\tilde{L}) x \quad (2.20)$$

$$\tilde{L} = \frac{2}{\lambda_{\max}} L - I$$

From this equation, it can be observed that each node depends only on the information inside the K -th order neighbourhood and with this reformulation, the computation of the equation is reduced to $O(|\xi|)$, linear to the number of edges ξ in the original graph G .

To build a neural network with graph convolutions, it's sufficient to stack multiple layers defined according to the previous equation, each followed by a nonlinear transformation. However, the authors of GCN [?] proposed limiting the convolution number to $K = 1$ at each layer instead of limiting it to the explicit parametrization by the Chebyshev polynomials. This way, each level only defines a linear function over the Laplacian Matrix L , maintaining the possibility of handling complex convolution filter functions on graphs by stacking multiple layers [?, ?]. This means the model can alleviate the overfitting of local neighbourhood structures for graphs whose node degree distribution has a high variance [?, ?].

At each layer, it can further considered $\lambda_{\max} \approx 2$, which the neural network parameters could accommodate during training [?]. With these simplifications, the equation is transformed into:

$$g_{\theta'} \star x \approx \theta'_0 x + \theta'_1 x (L - I_N) = \theta'_0 x - \theta'_1 D^{-\frac{1}{2}} A D^{-\frac{1}{2}} \quad (2.21)$$

θ'_0 and θ'_1 - Free parameters that can be shared over the entire graph

The number of parameters can, in practice, be further reduced, minimising overfitting and minimising the number of operations per layer as well [?] as equation 2.22 entails.

$$g_\theta \star x \approx \theta(I + D^{-\frac{1}{2}}AD^{-\frac{1}{2}})x \quad (2.22)$$

$$\theta = \theta'_0 = -\theta'_1$$

One potential problem is the $I_N + D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ matrix whose eigenvalues fall in the $[0, 2]$ interval. In a deep GCN, the repeated utilization of the above function often leads to an exploding or vanishing gradient, translating into numerical instabilities [?, ?]. In this context, the matrix can be further renormalized by converting $I + D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ into $\tilde{D}^{-\frac{1}{2}}\tilde{A}\tilde{D}^{-\frac{1}{2}}$ [?, ?]. In this case, only the scenario where there is one feature channel and one filter is considered which can be then generalized to an input signal with C channels $X \in \mathbb{R}^{N \times C}$ and F filters (or hidden units) [?, ?]:

$$H = \tilde{D}^{-\frac{1}{2}}\tilde{A}\tilde{D}^{-\frac{1}{2}}XW \quad (2.23)$$

$W \in \mathbb{R}^{C \times F}$ - Matrix of filter parameters

H - Convolved Signal Matrix

2.4.2 Graph Attention Network

Graph Attention Network (GAT) [?] is another type of GNNs that focuses on leveraging an attention mechanism to learn the importance of a node's neighbours. In contrast, the GCN uses edge weight as importance, which may not always represent the true strength between two nodes [?, ?].

The Graph Attention Layer defines the process of transferring the hidden node representations at layer $k - 1$ to the next node presentations at k . To ensure that sufficient expressive power is attained to allow the transformation of the lower-level node representations to higher-level ones, a linear transformation $W \in \mathbb{R}^{F \times F'}$ is applied to every node, followed by the self-attention mechanism, which measures the attention coefficients for any pair of nodes through a shared attentional mechanism $a : \mathbb{R}^{F'} \times \mathbb{R}^{F'} \rightarrow \mathbb{R}$ [?, ?]. In this context, relationship strength e_{ij} between two nodes i and j can be calculated by:

$$e_{ij} = a(WH_i^{k-1}, WH_j^{k-1}) \quad (2.24)$$

$H_i^{k-1} \in \mathbb{R}^{N \times F'}$ - Column-wise vector representation of node i at layer $k - 1$ (N is the number of nodes and F the number of features per node)

$W \in \mathbb{R}^{F \times F'}$ - Shared linear transformation

$a : \mathbb{R}^{F'} \times \mathbb{R}^{F'} \rightarrow \mathbb{R}$ - Attentional Mechanism

e_{ij} - Relationship Strength between nodes i and j

Theoretically, each node can attend to every other node on the graph, although it would ignore the graph's topological information in the process. A more reasonable solution is to only attend nodes in the neighbourhood [?, ?]. In practice, only first-order node neighbours are used, including the node itself, and to make the attention coefficients comparable across the various nodes, they are normalized with a *softmax* function:

$$\alpha_{ij} = \text{softmax}_j(\{e_{ij}\}) = \frac{\exp(e_{ij})}{\sum_{l \in N(i)} \exp(e_{il})}$$

Fundamentally, α_{ij} defines a multinomial distribution over the neighbours of node i , which can also be interpreted as a transition probability from node i to each node in its neighbourhood [?]. In the original work [?], the attention mechanism is defined as a single-layer Feedforward Neural Network that includes a linear transformation with weigh vector $W_2 \in \mathbb{R}^{1 \times 2F'}$ and a LeakyReLU nonlinear activation function with a negative input slope $\alpha = 0.2$ [?, ?]. More formally, the attention coefficients are calculated as follows:

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(W_2[WH_i^{k-1} || WH_j^{k-1}]))}{\sum_{l \in N(i)} \exp(\text{LeakyReLU}(W_2[WH_i^{k-1} || WH_l^{k-1}]))} \quad (2.25)$$

$||$ - Vector concatenation operation

The novel node representation is a linear composition of the neighbouring representations with weights determined by the attention coefficients [?, ?], formally:

$$H_i^k = \sigma\left(\sum_{j \in N(i)} \alpha_{ij} WH_j^{k-1}\right) \quad (2.26)$$

Multi-head Attention

Multi-head attention can be used instead of self-attention, determining a different similarity function over the nodes. An independent node representation can be obtained for each attention head according to the equation bellow [?, ?]. The final representation is a concatenation of the node representations learned by different heads, formally:

$$H_i^k = \left\|_{t=1}^T \sigma\left(\sum_{j \in N(i)} \alpha_{ij}^t W^t H_j^{k-1}\right)\right\|$$

T - Number of attention heads

α_{ij}^t - attention coefficient computed from the t -th attention head

W^t - Linear transformation matrix of the t -th attention head

Lastly, the author also mentions that other pooling techniques can be used in the final layer for combining the node representations from different heads, for example, the average node represen-

tations from different attention heads [?, ?].

$$H_i^k = \sigma\left(\frac{1}{T} \sum_{t=1}^T \sum_{j \in N(i)} \alpha_{ij}^t W^t H_j^{k-1}\right) \quad (2.27)$$

2.5 Smart Grid Services

Given the global ecological emergency and the increasing energetic crisis, there is a necessity for advancements in energy distribution and transmission systems now more than ever. To fulfil the need for energy sustainability, traditional centralized distribution grids must be adapted to, on the one hand, accommodate the rise of distributed renewable energy sources in corporate and domestic consumers and, on the other, to make more efficient and reliable distribution of energy resources [?, ?].

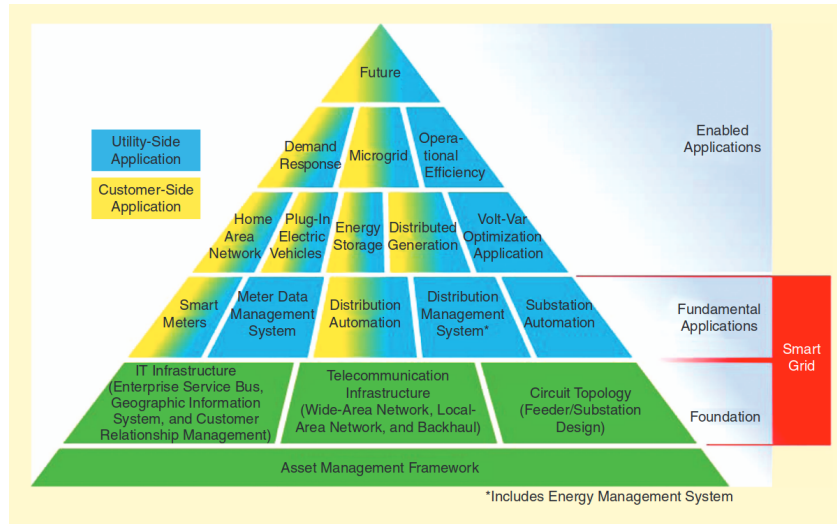


Figure 2.5: Smart Grid Capabilities Pyramid [?]

The *Smart Grid* or the *Smart Power Grid* conceptualizes this modernization of the electricity network by leveraging the technological advancements in information technology and communication science to create intelligent systems that manage and monitor the distributed generation of energy [?, ?]. Figure 2.5 describes the smart grid pyramid, which has asset management at its base. On this foundation, the foundation of the smart grid is laid out by the circuit topology, IT systems and telecommunications infrastructure, the basic ingredients for the emergence of fundamental applications such as smart meters and distribution automation [?]. In turn, these serve as building blocks for creating more intelligent systems that leverage upper-layer applications, enabling the true smart grid capabilities [?].

Chapter 3

Literature Review

In this chapter, the literature review regarding GRL approaches and DED systems is documented. This unit is divided into three sections with the first exposing the research methodology and the others reviewing GRL approaches and the DED Systems, the second is further subdivided into plain GCN, attention-based and other approaches.

3.1 Research Methodology

This literature review focuses on analysing the existent literature around the main objective of this dissertation, which is improving GRL techniques in the context of this work's application domain, smart grid services. In this manner, the main research questions are presented as:

- RQ1 - *How can RL algorithms be adapted to effectively solve sequential decision-making problems on graph-based environments?*
 - RQ1.1 - *What are the existent GRL approaches?*
 - RQ1.2 - *What are their limitations?*
- RQ2 - *How can GRL algorithms improve smart grid services?*
 - RQ2.1 - *How can GRL algorithms improve Dynamic Economic Dispatch (DED) in power distribution grids?*

These interrogations aggregate the relevant topics we aim to address in this review. The main requirement on the analyzed literature was to only consider research from the past five years, with the exception of seminal works. We also exclusively considered literature published in scientific conferences and top-tier journals.

In this manner, the initial exploratory research was conducted, primarily using *Scopus* and *Web of Science* for searching the published literature and alternatively using *Google Scholar* for finding cross-references when necessary. The research questions were translated into fundamental search queries where the research process was based.

- "Graph Reinforcement Learning" OR "Reinforcement Learning on Graphs"
- "Reinforcement Learning" AND "Power" AND "Dispatch"
- "Graph Reinforcement Learning" OR "Reinforcement Learning on Graphs" AND "Power" AND "Dispatch"

The gathered literature was screened and the relevant works were thoroughly reviewed, the following sections present the main findings.

3.2 Graph Reinforcement Learning Approaches

GRL or Reinforcement Learning on Graphs is a relatively new area in the broader field of machine learning. GRL techniques have shown significant progress in solving problems with underlying graph-based representations such as power grid management [?, ?], smart transport [?, ?] or task offloading [?, ?]. In this work, the main focus lies on studying the development of GRL techniques and subsequent application to smart grid services such as dynamic economic energy dispatch systems [?, ?], residential electricity behaviour identification and energy management [?], or Volt-VAR regulation [?].

Research on this topic has significantly increased in the last few years with the improvements of DRL techniques and the developments in GNNs in the mid-2010s [?, ?, ?, ?]. GNNs became the state-of-the-art for solving numerous data mining tasks involving graph-structured data, excelling at classification, link prediction and representation learning [?, ?]. This advancement brought more sophisticated RL applications on graphs and the surge of a new field studying how to combine the improvements of graph mining and reinforcement learning techniques [?, ?].

3.2.1 Plain GCN-Based GRL

A common approach in Graph Reinforcement Learning model implementation is the use of graph convolutions with the GCNs architecture for leveraging graph-based structures to extract and aggregate the essential features of data in hand and improve the performance of RL agents in those environments. The techniques listed in this subsection constitute approaches that integrate a GCN with RL algorithms. The gathered literature can be observed in table 3.1.

[?] implements a GRL system to improve the decision quality of economic dispatch under high penetration of distributed energy generations. To accomplish this, a SAC system is employed with the main objective of finding the optimal action policy for minimizing generation cost with the appropriate reliability concerns. This problem is represented by an undirected graph with nodes describing the power grid elements with their respective attributes and edges describing the underlying energy connections between those units. To extract the structural features of the graph, this work implements a full connected layer to perform feature transformation with a two-layer GCN followed by three full connected layers for the non-linear mapping of state-to-action policy in both actor and critic modules. [?] develops a similar approach, with both concluding that

it significantly reduces learning time for achieving better training results in comparison to plain SAC and showing significant improvement on economy and flexibility of the system on more complex and sparse state graphs. The use of GCNs enables the system to adapt to changes in the state space by leveraging the network's generalization ability.

In [?] a three-layer GCN is used to extract node feature and graph topology information and is integrated into a Rainbow-based [?] DRL algorithm for electric vehicle charging guidance. In this article, the testing results show promising performance in reducing operation costs for electric vehicle users, portraying a model with good generalization ability in untrained scenarios. This work and [?] further tested their proposed model's performance when inducing topology changes, yield promising results in the adaptability of the two GRL algorithms in scenarios where the environment suffers dynamic changes.

Another interesting implementation of this approach is [?], which studies and compares different solutions for optimizing autonomous exploration under uncertain environments. It analyses combinations of a single agent Deep Q-Network (DQN) and Advantageous Actor-Critic (A2C) with Graph Convolutional Networks, Gated Graph Recurrent Networks and Graph U-Nets. The algorithms are executed in test cases of various sizes and the paper reports that the GCN-DQN was the model that achieved the highest reward during policy training, followed by the GGNN-A2C model, although in the end, it concludes that the second showed improved scalability in relation to the first model. Other similar approaches include [?] for residential electricity behaviour identification and energy management with a behaviour correlation graph and [?] for line flow control in a power grid simulation.

The reviewed literature fails to consider methods for analysing the scalability of proposed models in relatively larger scenarios, except [?], with some proposing this as relevant future work [?, ?]. Beyond that, only [?], [?] [?] defined methods for evaluating the performance of the algorithms under topology variations. This shows a gap for more complete studies in plain GCN-based approaches that also focus on studying the scalability of GRL to large scenarios and the adaptability of the models to dynamic topology variations. In most approaches, namely in [?, ?, ?, ?], the presented results portray GRL techniques as significantly more effective than plain DRL models without a GCN already suggesting that these techniques are a potential solution for solving sequential decision-making problems with graph-based environments.

Table 3.1: GCN-Based GRL Literature

Reference	DRL Algorithm	GNN Algorithm	Application Domain
[?]	MDP	GCN	Interpret GNNs at model-level
[?]	DDPG	GCN	Automatic transistor sizing
[?]	A3C	GCN	Automatic Virtual Network Embeddings
[?]	Deep Q-Learning	GCN	Task Offloading in Edge Computing
[?]	Rainbow	GCN	Electrical Vehicle Charging Guidance
[?], [?]	SAC	GCN	Dynamic economic energy dispatch
[?]	SAC	GCN	Multi-access Edge Computing
[?]	DQN	GCN	Line flow control
[?]	DQN	GCN	Autonomous Exploration under uncertainty
[?]	DQN	GCN	Residential electricity behavior identification and energy management

3.2.2 Attention-based GRL

Another effective approach in extracting relevant topology and graph features relies on using attention mechanisms to weigh different nodes' contributions dynamically. While this encompasses techniques that use the GAT architecture, which is a GNN design with the attention mechanism at its core, various scholars propose GCN approaches integrated with attention mechanisms such as [?] and [?]. In the top part of table 3.2 the single-agent reviewed attention-based approaches can be observed, while at the bottom some relevant multi-agent approaches are listed.

[?] proposes a DDPG-based algorithm improved with a GAT block with three graph Attention Layers for extracting and learning the topology information for achieve real-time optimal scheduling for Active Distribution Networks (ADNs). This paper compares the obtained test results against a GCN-DDPG model and shows increased performance over the GCN method in reducing cost and power loss. Beyond this, the work demonstrates that the GAT's attention mechanism enables the algorithm to focus on more important nodes and improve the signal-to-noise ratio compared to its GCN counterpart. [?] and propose a multi-agent approach to the same domain but more focused on voltage regulation with a multi-agent SAC instead of a single-agent DDPG algorithm.

In [?], another model for the electric vehicle charging guidance is proposed, consisting of a bi-level approach of a Rainbow-based algorithm with a GAT block. The upper level focuses on the decision-making process regarding charging, while the lower level handles routing. The proposed model proved to be more effective than a shortest distance path-based [?] and a DRL-based [?] approach. [?] develops a similar approach with a Double-prioritized DQN for the same application domain. In [?] and [?], the sequential distribution system restoration problem is addressed with a multi-agent RL algorithm equipped with a GCN with an attention mechanism. In the first case, multi-head attention is used as the convolution kernel for the GCN with a DQN algorithm. In the second, self-attention is used for improving the centralized training of the used multi-agent actor-critic algorithm, more concretely, by embedding it in the critic networks. At the same time, the GCN is integrated into the actor networks for extracting the graph features. Both solutions proved more efficient than traditional RL techniques, with the first highlighting its solution generalization ability and the second showing increased scalability facing the non-GRL techniques. In general, the literature regarding attention-based approaches is a lot sparser and scarcer than GCN-based approaches. Only three relevant single-agent works were found [?, ?, ?] which might either suggest that scholars have a slight preference for multi-agent systems when implementing GRL with attention mechanisms or that these approaches in single-agent systems still require further research. Nevertheless, some of the works showed models with good adaptability to topology variations [?, ?, ?] and scalability to large scenarios [?, ?, ?, ?]. Notably, some works suggest directly embedding the GNN in the RL framework to boost the methodology computation performance and robustness [?, ?].

Table 3.2: Attention-Based GRL Literature

Reference	DRL Algorithm	GNN Algorithm	Application Domain
[?]	DDPG	GAT	Optimal Scheduling for ADNs
[?]	Rainbow	GAT	Electric Vehicle Charging Guidance
[?]	DQN	GAT	Electric Vehicle Charging Guidance
[?]	DQN	GCN	Multi-agent Sequential Distribution System Restoration
[?]	DQN	GCN	Active Power Rolling Dispatch
[?]	AC	GCN	Multi-agent Service Restoration
[?]	SAC	GAT	Multi-agent Voltage Regulation

3.2.3 Other Approaches

This subsection includes other relevant and promising GRL approaches that combine of other GNNs architectures with RL algorithms . In [?], a GraphSAGE network [?] is used with a

Deep Dueling Double Q-Network (D3QN) for emergency control of Undervoltage load shedding for power systems with various topologies. The author presents promising results for the GraphSAGE-D3QN model compared to a GCN-D3QN, achieving higher cumulative reward and faster voltage recovery speed, although it required longer decision times. The proposed model performed excellently in the application domain and successfully generalized the learned knowledge to new topology variation scenarios. Another approach that showed good performance with the GraphSAGE architecture was [?] in the context of the dynamic economic dispatch problem.

[?] focused on solving the Job shop scheduling problem through a priority dispatching rule with a Graph Isomorphism Network (GIN) [?] and an actor-critic PPO algorithm where the GIN is shared between actor and critic networks. The method showed superior performance against other traditional manually designed priority dispatching rule baselines, outperforming them by a large margin.

Table 3.3: Other GRL Approaches

Reference	DRL Algorithm	GNN Algorithm	Application Domain
[?]	DQN	GraphSAGE	Undervoltage Load Shedding
[?]	PPO	GraphSAGE	Dynamic Economic Dispatch
[?]	PPO	GIN	Dispatch for Job Shop Scheduling

3.3 Dynamic Economic Dispatch Systems

RL algorithms are already regarded as a well established potential solution for solving economic dispatch problems [?]. In table 3.4 the relevant RL approaches to the problem, including GRL techniques, can be observed. In a general level, the reviewed literature regarding DED solutions take into account a wide range of considerations and constraints. The recent works show an almost ubiquitous study of Energy Storage Systems (ESSs) management and Renewable Energy Source (RES) curtailment, given the current global energetic issues and the relevance of these technologies for solving them.

Some systems [?, ?, ?] take further steps in ensuring grid stability by also considering voltage deviations while a notable paper considers battery degradation when calculating the cost of dispatch [?]. In general, GRL algorithms show superior performance in relation with plain DRL approaches [?, ?, ?] as also concluded in section 3.2.1, with studies successfully ensuring adaptability of the GRL models to variations in the scenario's topology [?, ?]. However, only [?] conducts tests in test cases of different sizes, which shows a gap for studies addressing the scalability limitations of the developed models.

Table 3.4: Dynamic Economic Dispatch RL Systems

Reference	Approach	Application
[?]	DDPG with Prioritized Experience Replay Mechanism and L2 Regularization	Integrated Energy Systems, Utility (Selling, Purchasing and Gas) and ESSs
[?]	DDPG	Thermal Power, Renewable Energy Sources and ESSs
[?]	SAC with Imitation learning	Thermal Power, Renewable Energy Sources and Voltage deviation
[?]	GCN-SAC with Replay Buffer	Thermal Power, Renewable Energy Sources, ESSs and Voltage Deviation
[?]	GCN-SAC	Utility (Time-of-Use), Thermal Power, Renewable Energy Sources and ESSs
[?]	GraphSAGE-PPO	Thermal Power, Renewable Energy Sources, ESSs and Voltage Deviation
[?]	Multi-Agent RL with Function Approximation and Diffusion Mechanism	Utility (Time-of-Use), Thermal Power (Diesel) and ESSs (Considering degradation)
[?]	Multi-Agent GAT-DQN	Thermal Power, Renewable Energy Sources and ESSs

3.4 Conclusions

In this chapter, we reviewed relevant literature for this dissertation's main research topic, GRL algorithms. GRL is very promising field, where several different applications and techniques were already studied. GNNs architectures such as GCN have been extensively applied with DRL algorithms for enabling feature extracting from graph-based state representations [?, ?]. Architectures such as the GraphSAGE and other attention-based have also been successfully applied with very promising results [?, ?] in comparison with GCNs. However, less research regarding their integration with DRL algorithms was discovered. This suggests that a possible improvement and research direction in the development of GRL techniques might be connected with exploring the use of different GNNs architectures and using the rising attention-based techniques.

- **RQ1 - How can RL algorithms be adapted to effectively solve sequential decision-making problems on graph-based environments?**
 - **RQ1.1 - What are the existent GRL approaches?** Existent approaches ubiquitously use DRLs with GNNs for efficiently extract graph features. They can be divided into plain GCN approaches, Attention-based approaches and others. GCNs compromise

a popular method while attention-based approaches showed to be less researched but more promising [?, ?]. Tables 3.1, 3.2 and 3.3 depict the reviewed GCN, attention-based and other implementations, respectively.

- **RQ1.2 - What are their limitations?** Research already depicts GRL techniques as better performing in relation to plain RL algorithms in graph-based environments, which portrays GRL is a potential solution to this problem [?, ?, ?, ?, ?, ?]. The reviewed literature proved to be quite scarce and sparse, probably due to the novelty of the field. After a thorough review, we failed to find works that listed specific limitations of GRL algorithms. This suggests the existence of a research gap for works with a thorough and well-documented scientific process as well as comparative and systematic studies between the different approaches, highlighting models performance in large scenarios and under topology variation. Furthermore, some papers suggested the directly embedding of the GNN into the RL framework [?, ?].
- **RQ2 - How can GRL algorithms improve smart grid services?**
 - **RQ2.1 - How can GRL algorithms improve Dynamic Economic Dispatch (DED) in power distribution grids?** RL algorithms are already regarded as a well established potential solution for solving economic dispatch problems [?]. We were able to find evidence of GRL being successfully implemented in DED systems such as [?] and [?] in a power grid simulation context, showing efficient performance and scalability in comparison with plain DRL models in the same issue. However, as it was the case with GRL algorithms, the sparseness of different considerations, constraints and formalizations of the DED problem, highlighted in section 3.3 bring added complexity when comparing different approaches, since models are build with different DED requirements.

Chapter 4

Methodological Approach

In this chapter, we will formally state the problem this work aims to address, as well as its the proposed solution. In section 4.1 the problem statement is exposed, section 4.3 lists the solutions main functional, non-functional and structural requirements, section 4.4 addresses the architecture of the solution, section 4.5 exposes the established methodology and section 4.6 contains the proposed work plan.

4.1 Problem Statement

The reviewed literature addressing solutions to sequential decision-making problems in graph-based environments is sparse and scarce, leading to a research gap for comparative and systematic GRL approaches analysing scalability under different sized scenarios and adaptability to topology variation.

As addressed in the previous chapter, graphs are ubiquitous representations that can serve to instinctively represent several problems and their objects. In some network-oriented domains, these representations reveal underlying features that can't be naturally represented by plain Euclidean data. This problem becomes even more difficult considering that graph data is complex and sparse, something that brings the need for methods that efficiently extract representations.

By conducting a thorough literature review of the relevant studies in this context, we observed that current RL algorithms are not as efficient as GRL techniques in handling such complex environments, because of not considering and generalizing environment topology features in the decision-making process. This deeply affects the performance of decision systems inserted in network-oriented domains where the intricate relationships between the objects may be relevant for mapping the observable environment states to optimal action policies.

More and more GRL attracts the curiosity of academics, only increasing the relevance of this problem. With the recent advancements of GNNs, the popularity around GRL has risen because of their excellent efficiency in creating optimal graph representations and other graph machine learning problems. However, with GRL being a field whose research is still in initial phase, the gathered literature is very sparse, with a lack of works addressing the benefits, disadvantages

and performance of the various proposed models. Moreover, the literature also highlights the importance of studying GRL models in scenarios under topology variations and of different sizes for analysing their scalability.

4.1.1 Scope

This dissertation will focus on studying this problem in the context of single-agent RL algorithms, given that multi-agent systems are significantly more complex to implement. Furthermore, in the context of the dissertation's application domain, which is smart grid services, the possible improvements in GRL techniques will be implemented to the Dynamic Economic Dispatch (DED) problem that studies solutions that optimize power generation cost while maintaining reliable grid stability. Additionally, the GRL proposed models may be also implemented to solve other smart grid systems such as Undervoltage Load Shedding and Volt-VAR Regulation.

4.2 Problem Formalization

In this section, the main problem of study of this dissertation is formally introduced. Firstly, the problem is approached from an application domain perspective, uncovering the details of the DED problem and its main features. In the sequent subsection, the DED problem is formalized as a dynamic sequential decision-making problem, and its characteristics are presented under the form of its corresponding MDP.

4.2.1 Dynamic Economic Dispatch

The DED problem addresses the issue of balancing the necessity for ensuring stability, security and reliability of a power grid while also minimizing its operating cost. This problem is hardened by the paradigm observed in current power systems where ESS and RES are available and also need to be taken in consideration when studying solutions for this problem.

Equation 4.1 highlights the main objective function for the DED problem and is composed of three main components: $F_G(t)$ represents the cost of energy produced by conventional generators, $F_{RES}(t)$ is the term that depicts the cost of wasted energy from RES and $F_{ESS}(t)$ the ESS operation cost.

$$\min \sum_{t=1}^T F_G(t) + F_{RES}(t) + F_{ESS}(t) \quad (4.1)$$

Conventional generation cost calculation is presented in equation 4.2. For the sake of simplicity, this component is reduced to a linear equation and represented by a static c_i term for generator i in €/MWh . Regarding the cost of wasted RES energy, a penalty term β_{RES} is introduced that consists in a

$$F_G(t) = \sum_{i=0}^N c_i P_i^G \Delta t \quad (4.2)$$

$$F_{\text{RES}}(t) = \sum_{i=0}^K \beta_{\text{RES}} (\overline{P_i^{\text{RES}}} - P_i^{\text{RES}}) \Delta t \quad (4.3)$$

$$F_{\text{ESS}}(t) = \sum_{i=0}^K c_i (\overline{P_i^{\text{RES}}} - P_i^{\text{RES}}) \Delta t \quad (4.4)$$

4.2.2 Markov Decision Process (MDP)

4.3 Requirements

4.3.1 Functional

Table 4.1: Functional Requirements

FR	Title	Description
F1	Dispatch Optimization	<p>The system should be able to optimize the dispatch power of generation resources in real-time by managing generator power levels to meet the time-varying load demands. They can be of the following types:</p> <ul style="list-style-type: none"> • Conventional Thermal Generation • Photovoltaic Cell • Wind Turbine
F2	ESS Management	The system is responsible for managing energy storage resources by controlling their discharge and charge operations.
F3	Voltage Regulation	The system needs to appropriately manage voltage levels of generators and ESSs to maintain grid stability
F4	Graph Features Extraction	The system must implement a graph machine learning mechanism for creating and generalising efficient representations of the environment
F5	Learning Capability	The system must learn from experience and optimize its dispatch policies over time
F6	Motorization	<p>The system must implement tracking mechanisms in order to gather the different metrics, namely:</p> <ul style="list-style-type: none"> • Convergence Rate • Training Time • Computation Time • CPU/GPU Utilization • Operating Cost • Average ESS charge • Loss of renewable energy active power <p>in order to enable the comparative and evaluative process between the different models.</p>

4.3.2 Non-functional

Table 4.2: Non-functional Requirements

NFR	Title	Description
NF1	Adaptability	The system should adapt to changes in the power grid topology or in load patterns
NF2	Reliability	The system should fulfil the various power grid constraints, namely power balance, generator constraints, ramp rate limits and ESS constraints
NF3	Scalability	The system must handle scenarios of different sizes, ranging from small to complex power distribution grids
NF4	Simulation Environment	The system will be implemented in the context of an IEEE bus system test case that simulated the operation of a real-world power distribution grid.

4.3.3 Structural

Table 4.3: Structural Requirements

SR	Title	Description
S1	Storage Space	Required for storing the test cases and the project's various artifacts
S2	CPU/GPU & Memory Resources	The system requires sufficient processing power and volatile memory for executing the models and power grid simulations
S3	Connectivity	This project requires a stable internet connection, specially for downloading the large scenarios
S4	Visualization	This system requires visualization capabilities in order to observe the grid state and metrics at real-time

4.4 Architecture

Considering the conclusions taken from the reviewed literature, in chapter 3, our proposed solution combines efficient Deep Reinforcement Learning (DRL) approaches with Graph Neural Networks (GNNs), the state-of-the-art approach for graph representation learning. Generally, the system receives graph-based representations of the environment and encodes them using the GNN algorithm. By also leveraging deep learning techniques, DRL maps the encoded embeddings to

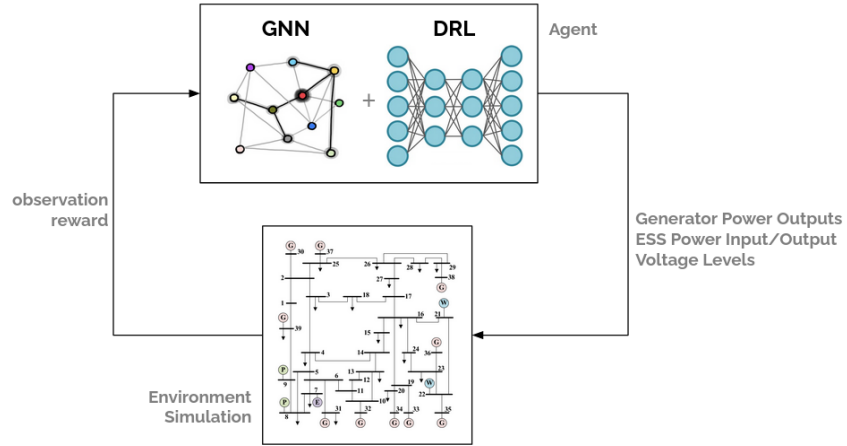


Figure 4.1: Solution Architecture

optimal action sequences with the goal of meeting the real-time load demand and reducing the operating cost of the power distribution grid. The general architecture of the solution can be observed in figure 4.1.

To fulfil this, the agent adjusts the generator power output and voltage levels and manages the ESS operation in real-time.

The solution is executed in the context of a power grid simulation that models the appropriate elements, properties, constraints and operating of the power distribution grid. This method is further described in the upcoming subsection 4.5.

4.5 Methodology

Regarding the established methodology for enabling the implementation of the proposed solution, the *grid2Op* framework [?] will be used for modelling the sequential decision-making process on simulated power distribution grid. *grid2Op* is designed by RTE (*Réseau de Transport d'Électricité*), the electricity transmission system operator of France, and is equipped with a variety of pre-defined scenarios used in coding competitions and based on real-world data [?]. This platform focuses on easing the job using the grid topology to control the power flows and also allowing for it to be reconfigured in real-time [?]. Additionally, enables to graphically plot the current observable state of the grid, as portrayed in figure 4.2. Apart from graph topology, *grid2op* also enables the manipulation of:

- the **voltages** by manipulating the set-point value of the generators;
- the **active generation** by mapping the received observations to optimal sequences of dispatch actions [?]

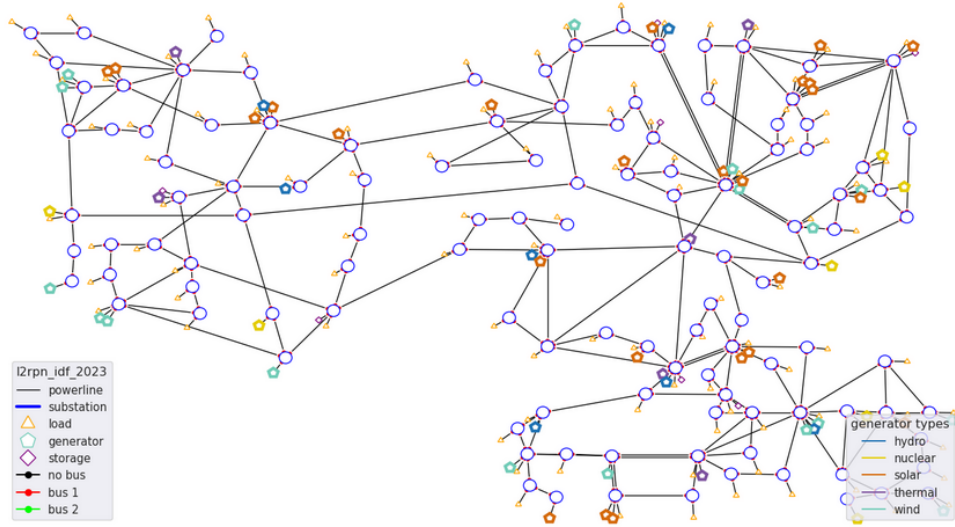


Figure 4.2: Grid2Op *l2rpn_idf_2023* 118-bus test case [?]

This framework is compatible with the *Gymnasium* framework [?], a widely used toolkit for developing RL algorithms, which will also be used in this work together with the Grid2Op simulation environment.

This platform focuses on describing the power distribution grids by modelling the following objects:

- **Buses** are the fundamental objects of the power grid, representing nodes where power sources, loads and other elements are connected[?].
- **Powerlines** represent edges in the power grid and connect the different buses together. They represent the physical transmission and distribution lines and allow power to flow from one part of the grid to another [?].
- **Generators** are critical grid elements connected to buses whose main role is to produce power and maintain grid stability by balancing the energy supply and demand. They can be Conventional Thermal Generators, Wind Turbines or Photovoltaic Cells [?].
- **Loads** consume power from the grid, simulating electricity use. They're also associated to an individual bus [?].
- **Storage Units** can act as both consumers and producers. They're able to retain energy from the power grid when production surpasses demand for later injecting back power when convenient. Storage units are bound by a maximum energy storage capacity [?].

Table 4.4: Test Case Sizes

	l2rpn_icaps_2021	l2rpn_idf_2023
Buses	36	118
Powerlines	59	186
Generators	22	62
Loads	37	99

This work will use the pre-defined *grid2op l2rpn_icaps_2021* and *l2rpn_idf_2023* test environments, a modified case studies based on the original IEEE 118-bus system test case [?] of 36 and 118 buses, respectively. The first test case is a subset of the original 118-bus system with 50 years worth of data divided into independent ¹ weekly scenarios, while in the second case the system was modified to accommodate the *possible energy mix* of France in 2035, containing 16 years of data [?]. The scenarios document the loads and productions at each time step, the grid layout (for display purposes), generator and ESS characteristics [?] for a weekly episode, the sizes of test cases can be further observed in table 4.4. In addition, the test cases data will be modified to reflect the defined requirements of this work.

Beyond *grid2op* to build and run the test cases and *gymnasium* as the RL framework, this project will use *PyTorch* [?] with *PyTorch Geometric* [?] library for developing the GNN because of its extensive list of available implemented models. The different combinations of algorithms will be applied to a set of modified scenarios that fulfil the settled requirements. For Deep Reinforcement Learning algorithms, solutions with plain SAC, DDPG and PPO approaches and combined with the GCN, GAT and GIN architectures.

Concerning result analysis, it's also important to point out the use of quantitative methods for evaluating the different implemented models, a topic that is further explored in the following section 4.5.1.

4.5.1 Evaluation Metrics

The solution described in the previous subsections 4.4 and 4.5 clearly involves intricate operational mechanisms, something that calls for a sophisticated evaluative process. Furthermore, given the comparative nature of this dissertation the evaluation and analysis methods will be key factors in studying and confronting the different combinations of GNNs and DRL, as well as possible improvements in the integration of these techniques. In this manner, we define the four dimensions for evaluating and analysing the different GRL models:

Learning efficiency This dimension assesses how effectively the models learn and improve their decision-making process over time. It involves evaluating how quickly they converge to optimal or near-optimal dispatch strategies through the convergence rate.

¹non-consecutive

Computational Efficiency It's crucial that the solution is able to perform well on real-time execution. In this context, not only it's important to assess the solution's decision computation performance but also to measure the time necessary for the model offline training, as well as the observed CPU/GPU resource utilization.

Dispatch Efficiency The performance of the GRL model in managing power distribution from the various generators will be mainly measured by the average operating cost (in *Euros*) derived from the agent's sequence. Furthermore, other system behaviours will be analysed such as Renewable Energy Source Penetration, through the average ratio of maximum power and real power output, and ESS utilization, through the average energy stored.

Scalability and Adaptation The solution will be applied to scenarios of several sizes for analysing its scalability. Furthermore, we will induce variations in the simulation scenarios to test the model's ability to handle topology changes.

4.6 Work Plan

In order to make this project possible, we propose a 40-week work plan divided into four stages, whose start was in the last week of September 2023 and expected conclusion is scheduled for the last week of June 2024, observable in figure 4.3. The first stage, ranging from the first week of the plan to the delivery of this report, encompasses the **dissertation preparation** phase. This stage mainly addressed the initial knowledge acquisition necessary to understand this work's context and gather a basic grasp on the concepts and topics it addresses, namely:

- Reinforcement Learning
- Graph Theory
- Graph Representation Learning
- Graph Reinforcement Learning
- Graph Neural Networks
- Smart Grids Services

This process is followed by the analysis of existing research and approaches in GRL and Smart Grid technologies which spans till the first week of January 2024, documented in chapter 3. Lastly, the preparation phase is also compromised of the PDIS report write-up until its delivery in the first week of February 2024, marking the end of this phase.

The preparation phase is followed by the **Development** phase. This stage encompasses all activities related to the system implementation, from the preparation of the simulation environment to the model training and calibration. The development phase spans from the first week of February 2024 to the second-to-last week of April 2024.

Next, the **Evaluation & Analysis** stage is set, comprising the time frame where the various implemented models will be evaluated in the context of the previously mentioned dimensions and the evaluation results will be analysed.

Lastly, the **Dissertation** write-up and subsequent revision work will be performed, comprising of the write-up, proofreading and correction tasks and ranging from the end of the evaluation and analysis phase in the last week of April 2024 until the last week of the plan.

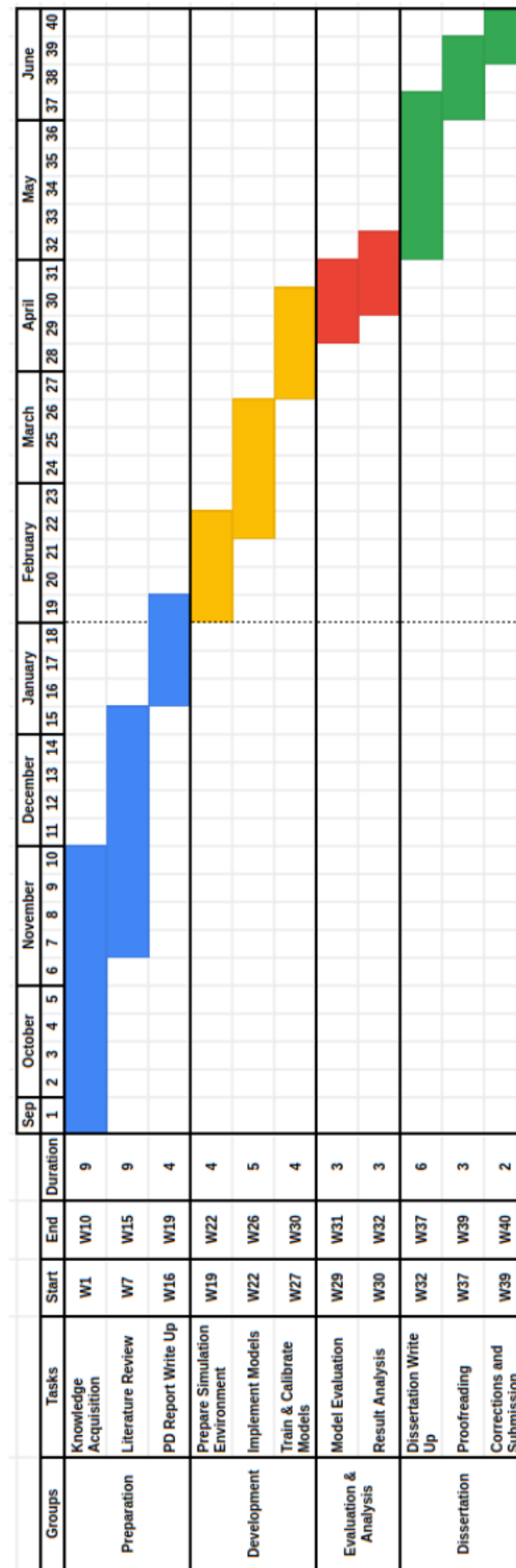


Figure 4.3: Project Work Plan

Chapter 5

Experimental Setup

5.1 Simulation Environment

5.2 Scenarios

5.3 Experiments

5.3.1 GRL

5.3.2 DRL vs. GRL

5.3.3 Scalability

5.3.4 Adaptability

Chapter 6

Conclusions

In conclusion, this report establishes a solid foundation for implementing and proposing concrete improvements to GRL techniques and efficiently solve the problem of DED in graph-based environments.

We concluded in the literature review, that GRL is already a promising solution for DED, yielding better performance than other approaches. In general, GNNs are an ubiquitous method, throughout the gathered literature, for extracting efficient environment topological representations. The identified limitations of existent implementations lie on their lack of: (1) scalability to larger environments, resulting in a significant decrease in computation performance; (2) a seamless integration between GNN and DRL algorithms and (3) adaptability to real-time topology changes.

Moreover, we aim to study how to tackle these limitations by implementing different GRL approaches for solving the DED problem and conducting a comparative study between the different techniques. The implementations will be compromised of various combinations of GNNs, for efficiently extracting the relevant environment features, and DRL algorithms for mapping the extracted representations into optimal action sequences. Case studies will be carried out within different size modified IEEE bus systems for testing the various models and confront their results. Lastly, the models will be analysed considering their training efficiency, computational performance, dispatch efficiency, scalability to large scenarios and adaptability to topological changes. Based on the conclusions drawn from the comparative study, we will propose a GRL that implements concrete enhancements.

6.1 Expected Contributions

In this section, we list the expected contributions derived from this work on a Scientific, Technological and Application levels:

Scientific A systematic and comparative study of different GRL approaches. This fills the research gap for systematic studies comparing different proposed techniques. Furthermore, this work will bring a clearer insight regarding the best practices on implementing GRL models

Technological A model resulting from this study with concrete improvements over the current GRL techniques proposed so far and tackling their limitations. Beyond this,

Application An efficient solution for the DED problem with GRL algorithms, compromising significant contribution to the research on DED systems in complex scenarios

6.2 SWOT Analysis

	Helpful	Harmful
Internal	Strenghts <ul style="list-style-type: none"> - Good understanding of the base knowledge inspiring this work - Well-defined evaluative dimensions - At least a third of the dissertation write-up is concluded - Well-maintained power grid simulation framework 	Weaknesses <ul style="list-style-type: none"> - Sparse and short supply of literature - Somewhat unclear picture of the concrete architectural improvements to be made in GRL algorithms
External	Opportunities <ul style="list-style-type: none"> - Deliver scientific contributions to the area of GRL - Propose a state-of-the-art GRL model for DED - Identify deficiencies in current GRL approaches - Compare different GRL techniques 	Threats <ul style="list-style-type: none"> - Significant training times for training models in large scenarios - High computational complexity of implemented models

Figure 6.1: SWOT Analysis

6.3 SMART Analysis

Lastly, we formulated this dissertation's objectives considering the criteria imposed by SMART goals. Recapitulating from chapter 1:

1. Perform a review of literature regarding GRL approaches and DED systems
2. Conduct a comparative and systematic empirical study of different GRL solutions of the DED problem
3. Propose a GRL model and concrete improvements facing the literature proposed models

In this manner, we analyse these goals from the SMART perspective to ensure that they represent a clear road map for conducting this study. Considering each dimension of the SMART criteria, we conclude that the defined objectives are:

Specific All of the goals refer to a concrete objective of this work and are further specified in the proposed work plan (figure 4.3)

Measurable The progress towards each objective is measurable because their time-bound

Achievable The objectives are realistic and attainable, considering the already existent research

Relevant Objectives 1 and 2 enable a better foundation for accomplishing the third objective, which compromises the main goal of this dissertation

Time-bound All of the objectives are time-bound by the proposed work plan in figure 4.3, divided into more specific tasks.