
Thesis proposal

My Thesis

Anonymous Author(s)

Affiliation

Address

email

1 Introduction

2 Drug discovery—the process by which a potential new medicine is identified—is a complex process that
3 encompasses the intersection of several fields (such as biology, statistics, chemistry or pharmacology).
4 The entire process is a long and costly endeavour, with a typical time-frame of 10 to 20 years till
5 market release and an estimated cost between 1 and 2 billion USD. With just a small quantity of the
6 initially identified compounds actually becoming an approved medicine. Many of these dropouts
7 happening at the early stages of the entire pipeline.

8 It exists, then, a need for better mechanisms for detecting better candidates. One of the most promising
9 directions is to improve the currently used and develop new *in-silico* methods—computational
10 simulations are relatively cheap and quick run that makes them an interesting solution. In-silico
11 simulations then cover two main aspects: modelling the dynamics of the human body, such that
12 any effect relevant to the drug or the disease will be captured by it; and a method to generate good
13 candidates that is effective at exploring the vast space of possible compounds.

14 Among the different computational approaches that have been used in the process of drug discovery
15 deep learning (DL) has shown signs to be a potential game changer. DL has been able to capitalize on
16 the exponential growth of data and the higher availability of computational resources. For example, DL
17 has had a remarkable success on computer vision (CV) and natural language processing (NLP), and
18 has become the go-to solution for any problem in these two fields. It is, at the same time, penetrating
19 into other fields, drug-discovery being one of them [1].

20 One of the big challenges is to unify all the aspects of drug-discovery and be able to incorporate all
21 the relevant biological information when designing possible candidate molecules. A success story
22 on that line is the recently paper published by Zhavoronkov et al. [3] where the authors describe a
23 deep learning method by which they are able to discover inhibitors of discoidin domain receptor 1
24 (DDR1)—a kinase implicated in fibrosis—in just 21 days.

25 Those promising results, albeit encouraging, are just the tip of the iceberg. There is still a long way
26 till a model can satisfactorily capture the biological complexity of any arbitrary target and produce
27 promising candidates. On top of that, there is an added dimension as such model should account for
28 the variability from patient to patient and be able to generate a molecule that accommodates for all the
29 genotypic and phenotypic variants, or generate different candidates for each of the genetic populations
30 of interest. [need a ref here]

31 Even more, in the case of diseases like cancer, an heterogeneous population may appear within a
32 single patient. So the same variant effects arise inside a dynamic ecosystem, where a drug that just
33 targets a subpopulation may lead to an evolutionary pressure complicating further the treatment
34 outlook [reference paper of evolutionary perspective to cancer].

35 There is then a great need to develop models that can be conditioned based on a large set of biological
36 [conditions?] and meaningfully account for this variatos when generating a compound or/and
37 evaluating a compounds effect when administered.

38 —

39 When it comes to the methodology it exists an orthogonal problem [regarding structured data] when
40 we deal with biological data. If we look at the case of deep learning for CV, where we deal with
41 images, a key element of any architecture for its success was the use of convolutional layers—one
42 will mostly observe convolutional neural networks (CNNs) when analyzing the state of the art in
43 CV—which introduce a structural a priori based on the structure of the input. A similar case can be
44 made for NLP. For that reason, there exists a strong signal to look for models that can leverage the
45 structural equivalent when in molecule or protein data, i.e. leverage graph structures.

46 [There is already a literature on this and I'm gonna talk about it [2]]

47 —

48 The aim of this thesis will be to explore several ways to tackle this challenge.

49 —

50 Deep learning has been applied very successfully in many fields, notably computer vision (CV) and
51 natural language processing (NLP), among those fields

52 graph convolutional neural networks. We plan to

53 —

54 Analyzing graph like structures allows us to

55 **Aim**

56 The aim of my thesis will be to explore how deep learning applied to the domain of graphs can help
57 capture better biological aspects

58 **References**

- 59 [1] Hongming Chen, Ola Engkvist, Yinhai Wang, Marcus Olivecrona, and Thomas Blaschke. The
60 rise of deep learning in drug discovery, 2018.
- 61 [2] Mengying Sun, Sendong Zhao, Coryandar Gilvary, Olivier Elemento, Jiayu Zhou, and Fei Wang.
62 Graph convolutional networks for computational drug development and discovery. *Briefings in*
63 *Bioinformatics*, 2019(0):1–17, 2019.
- 64 [3] Alex Zhavoronkov, Yan A Ivanenkov, Alex Aliper, Mark S Veselov, Vladimir A Aladinskiy,
65 Anastasiya V Aladinskaya, Victor A Terentiev, Daniil A Polykovskiy, Maksim D Kuznetsov,
66 Arip Asadulaev, Yury Volkov, Artem Zholus, Rim R Shayakhmetov, Alexander Zhebrak, Lidiya I
67 Minaeva, Bogdan A Zagribelnyy, Lennart H Lee, Richard Soll, David Madge, Li Xing, Tao
68 Guo, and Alán Aspuru-Guzik. Deep learning enables rapid identification of potent DDR1 kinase
69 inhibitors. *Nature Biotechnology*, 37(9):1038–1040, 2019.