
Thesis proposal

An alpha version of the draft

Nil Adell Mill
Winter 2020
Institute of Neuroinformatics

1 Introduction/Background

Drug discovery—the process by which a potential new medicine is identified—is a complex process that encompasses the intersection of several fields (such as biology, statistics, chemistry or pharmacology). The entire process is a long and costly endeavour, with a typical time-frame of 10 to 20 years till market release and an estimated cost between 2 and 3 billion USD [Schneider, 2019, Scannell et al., 2012]. With just a small quantity of the initially identified compounds actually becoming an approved medicine. Many of these dropouts happening at the early stages of the entire pipeline.

It exists, then, a need for better mechanisms for detecting better candidates. One of the most promising directions is to improve the *in-silico* methods—computational simulations are relatively cheap and quick run that makes them an interesting solution. *In-silico* simulations then cover two main aspects: **predictive modelling**, meaning modelling the dynamics of the human body—such that any effect relevant to the drug or the disease will be captured by it—and **generative modelling**, i.e. methods to generate good candidates which, in part, means methods that are effective at exploring the vast space of possible compounds, estimated to be on the order of 10^{60} [Reymond et al., 2012].

Among the different computational approaches that have been used in the process of drug discovery deep learning (DL) has shown signs to be a potential game changer [Dargan et al., 2019]. DL has been able to capitalize on the exponential growth of data and the higher availability of computational resources. For example, DL has had a remarkable success on computer vision (CV) [?] and natural language processing (NLP) [?], and has become the go-to solution for any problem in these two fields. It is, at the same time, penetrating into other fields, drug-discovery being one of them [Chen et al., 2018].

When we deal with this biological and molecular data, it exists a challenge on how to deal with the intrinsic structure of the data. If we look at the case of deep learning for CV, where we deal with images, a key element of any architecture for its success was the use of convolutional layers—one will mostly observe convolutional neural networks (CNNs) when analyzing the state of the art in CV—which introduce a structural prior based on the structure of the data[?]. A similar case can be made for NLP[?]. For that reason, there exists a strong signal to look for models that can leverage the structural equivalent when in molecule or protein data, i.e. leverage graph structures. On that direction there has been a rising field on the use of Graph Convolutional Neural Networks (GCNNs) [Wu et al., 2019], and in fact, there have been several models as such being proposed in the drug design literature [Sun et al., 2019]. [NOTE: Rewrite this last section]

Another of the big challenges is to unify all the aspects of drug-discovery and be able to incorporate all the relevant biological information when designing possible candidate molecules. An initial success story on that line is a recently paper [Zhavoronkov et al., 2019] where the authors describe a deep learning method by which they are able to discover, synthesise, and test in an animal model, inhibitors of discoidin domain receptor 1 (DDR1)—a kinase implicated in fibrosis—in just 21 days.

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

[I am not completely sure about this paragraph but I leave it here so I don't forget for now] Even more, in the case of diseases like cancer, an heterogeneous population may appear within a single patient. So the same variant effects arise inside a dynamic ecosystem, where a drug that just targets a subpopulation may lead to an evolutionary pressure complicating further the treatment outlook [reference paper of evolutionary perspective to cancer].

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

Last of all, in a more holistic view, it is of interest to develop multi-scale models that capture system complexity at the different levels. For instance, a model that is able to learn protein-compound interactions—commonly known as the docking problem—while at the same time use this information to predict effects of the introduction of the compound on the larger protein-protein interaction (PPI) network. [ref? [Sun et al., 2019]]

Aim & Methods

[Should I separate em in two different sections?]

The aim of this thesis will be two fold. One the one side, analyze how the explicit use of GCNNs may open new oportunities when dealing with biological and checmical data. On the other side, explore how modelling the biology at different levels (e.g. molecular structure v.s. molecular interacton network) may help with our understanding [of the biology? of compounds interaction?] and create better models. Furthermore, evaluate how these different scales may be integrated toguether.

This precise work will be focused around exploring all these concepts in the context of designing anti-cancer drugs. The work will be done in colaboration with the Computational Systems Biology group at IBM Research (Zurich). The group is currently focused on individualised paediatric cure (iPC), so an end goal of this project is for the end results of it to help in that effort, for instance in contibuting to the ongoing research in neuroblastoma.

As mentioned previously, the idea of using GCNNs for drug discovery is not a new one in the literature [Sun et al., 2019]. My project will build upon those ideas presented in the literature, expand them and test their feasibility by implementing them into a wider framework for drug design [Born et al., 2019]. In that context two main areas of application appear. One of them would be to re-desing the generative model, for instance by reframing the vairational autoencoders, used for molecule generation, to architectures that operate over graphs [Simonovsky and Komodakis, 2018, Li et al., 2018a, Li et al., 2018b]. The second area would be to find better ways to asses the activity of these molecules, and in a wider context, assess their relevance as drug candidates, i.e. improve the predictive model. In the concrete case of the mentioned framework it is done by using a critic network [Manica et al., 2019]. This could be expanded on a set of different fronts: usign structural data instead of SMILES [Li et al., 2017, Do et al., 2019], by using GCNNs to cover a much wider network of genes [Oskooei et al., 2019, Wang et al., 2019], or by introducing particular scores (rewards) based in the interaction of the compound to certain targets [Yingkai Gao et al., 2018, Zhavoronkov et al., 2019] or the combination of the compound with other drugs [Zitnik et al., 2018]—a common practice in patients with cancer.

All these possible changes on the critic model would apply at different abstraction levels. That opens the door to seek for ways to integrate the representations learnt at those different stages [Ying et al., 2018, Ma and Zhang, 2019, Huang et al., 2019]. On top of that information extracted from here could be then leveraged on the drug generation part of the framework.

References

- [Born et al., 2019] Born, J., Manica, M., Oskooei, A., Cadow, J., and Martínez, M. R. (2019). PaccMann^{RL}: Designing anticancer drugs from transcriptomic data via reinforcement learning.
- [Chen et al., 2018] Chen, H., Engkvist, O., Wang, Y., Olivecrona, M., and Blaschke, T. (2018). The rise of deep learning in drug discovery.
- [Dargan et al., 2019] Dargan, S., Kumar, M., Ayyagari, M. R., and Kumar, G. (2019). A Survey of Deep Learning and Its Applications: A New Paradigm to Machine Learning. *Archives of Computational Methods in Engineering*.

90 [Do et al., 2019] Do, K., Tran, T., and Venkatesh, S. (2019). Graph transformation policy network for
 91 chemical reaction prediction. In *Proceedings of the ACM SIGKDD International Conference on Knowledge*
 92 *Discovery and Data Mining*, pages 750–760.

93 [Huang et al., 2019] Huang, J., Li, Z., Li, N., Liu, S., and Li, G. (2019). AttPool: Towards Hierarchical
 94 Feature Representation in Graph Convolutional Networks via Attention Mechanism. *The IEEE International*
 95 *Conference on Computer Vision (ICCV)*, pages 6480–6489.

96 [Li et al., 2017] Li, J., Cai, D., and He, X. (2017). Learning Graph-Level Representation for Drug Discovery.
 97 Technical report.

98 [Li et al., 2018a] Li, Y., Vinyals, O., Dyer, C., Pascanu, R., and Battaglia, P. (2018a). Learning Deep
 99 Generative Models of Graphs.

100 [Li et al., 2018b] Li, Y., Zhang, L., and Liu, Z. (2018b). Multi-objective de novo drug design with conditional
 101 graph generative model. *Journal of Cheminformatics*, 10(1).

102 [Ma and Zhang, 2019] Ma, T. and Zhang, A. (2019). Incorporating Biological Knowledge with Factor Graph
 103 Neural Network for Interpretable Deep Learning.

104 [Manica et al., 2019] Manica, M., Oskooei, A., Born, J., Subramanian, V., Saéz-Rodríguez, J., and Rodríguez
 105 Martínez, M. (2019). Toward Explainable Anticancer Compound Sensitivity Prediction via Multimodal
 106 Attention-Based Convolutional Encoders. Technical report.

107 [Oskooei et al., 2019] Oskooei, A., Manica, M., Mathis, R., and Martínez, M. R. (2019). Network-based
 108 Biased Tree Ensembles (NetBiTE) for Drug Sensitivity Prediction and Drug Sensitivity Biomarker Identifi-
 109 cation in Cancer. *Scientific Reports*, 9(1).

110 [Reymond et al., 2012] Reymond, J. L., Ruddigkeit, L., Blum, L., and van Deursen, R. (2012). The enumera-
 111 tion of chemical space. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 2(5):717–733.

112 [Scannell et al., 2012] Scannell, J. W., Blanckley, A., Boldon, H., and Warrington, B. (2012). Diagnosing the
 113 decline in pharmaceutical R&D efficiency. Technical Report 3.

114 [Schneider, 2019] Schneider, G. (2019). Mind and machine in drug design. *Nature Machine Intelligence*,
 115 1(3):128–130.

116 [Simonovsky and Komodakis, 2018] Simonovsky, M. and Komodakis, N. (2018). GraphVAE: Towards
 117 generation of small graphs using variational autoencoders. In *Lecture Notes in Computer Science (including*
 118 *subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 11139
 119 LNCS, pages 412–422.

120 [Sun et al., 2019] Sun, M., Zhao, S., Gilvary, C., Elemento, O., Zhou, J., and Wang, F. (2019). Graph
 121 convolutional networks for computational drug development and discovery. *Briefings in Bioinformatics*,
 122 2019(0):1–17.

123 [Wang et al., 2019] Wang, M., Yu, L., Zheng, D., Gan, Q., Gai, Y., Ye, Z., Li, M., Zhou, J., Huang, Q., Ma,
 124 C., Huang, Z., Guo, Q., Zhang, H., Lin, H., Zhao, J., Li, J., Smola, A., and Zhang, Z. (2019). Deep Graph
 125 Library: Towards Efficient and Scalable Deep Learning on Graphs. Technical report.

126 [Wu et al., 2019] Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., and Yu, P. S. (2019). A Comprehensive
 127 Survey on Graph Neural Networks. Technical report.

128 [Ying et al., 2018] Ying, R., Morris, C., Hamilton, W. L., You, J., Ren, X., and Leskovec, J. (2018). Hier-
 129 archical graph representation learning with differentiable pooling. In *Advances in Neural Information*
 130 *Processing Systems*, volume 2018-Decem, pages 4800–4810.

131 [Yingkai Gao et al., 2018] Yingkai Gao, K., Fokoue, A., Luo, H., Iyengar, A., Dey, S., and Zhang, P. (2018).
 132 Interpretable drug target prediction using deep neural representation. In *IJCAI International Joint Confer-*
 133 *ence on Artificial Intelligence*, volume 2018-July, pages 3371–3377.

134 [Zhavoronkov et al., 2019] Zhavoronkov, A., Ivanenkov, Y. A., Aliper, A., Veselov, M. S., Aladinskiy, V. A.,
 135 Aladinskaya, A. V., Terentiev, V. A., Polykovskiy, D. A., Kuznetsov, M. D., Asadulaev, A., Volkov, Y.,
 136 Zholus, A., Shayakhmetov, R. R., Zhebrak, A., Minaeva, L. I., Zagribelnyy, B. A., Lee, L. H., Soll, R.,
 137 Madge, D., Xing, L., Guo, T., and Aspuru-Guzik, A. (2019). Deep learning enables rapid identification of
 138 potent DDR1 kinase inhibitors. *Nature Biotechnology*, 37(9):1038–1040.

139 [Zitnik et al., 2018] Zitnik, M., Agrawal, M., and Leskovec, J. (2018). Modeling polypharmacy side effects
 140 with graph convolutional networks. In *Bioinformatics*, volume 34, pages i457–i466.