# PaccMann^RL: Designing anticancer drugs from transcriptomic data via reinforcement learning

Jannis Born[*1,2], Matteo Manica[*1], Ali Oskooei[*1], Joris Cadow[1], Karsten Borgwardt[2], and María Rodríguez Martínez[1]

[1]*IBM Research Zurich, Switzerland*
[2]*Machine Learning & Computational Biology Lab, ETH Zurich, Switzerland*

April 20, 2020

[*] Equal contributions

**125-character summary:** A framework to bridge systems biology and drug discovery for designing anticancer drugs from transcriptomic data

### Abstract

With the advent of deep generative models in computational chemistry, in silico drug design has undergone an unprecedented transformation. While state-of-the-art deep learning approaches have shown potential in generating compounds with desired chemical properties, they disregard the cellular environment and biomolecular properties of the target disease. Here, we introduce a novel framework for de-novo molecular design that systematically leverages systems biology information into the drug discovery process. Embodied through two separate Variational Autoencoders (VAE), the drug generation is driven through a disease context (transcriptomic profiles of cancer cells) deemed to represent the target environment in which the drug has to act. Showcased at the challenging task of de-novo anticancer drug discovery, our conditional generative model is demonstrated to be capable of tailoring anticancer compounds to target a specific biomolecular profile, according to the critic. Without incorporating explicit information about anticancer drugs, we demonstrate how the molecule generation, starting from a random point in a chemical space, can be biased towards compounds with high predicted inhibitory effect against individual cell lines or cell lines from specific cancer sites. We verify our approach by investigating candidate drugs generated against specific cancer types and find the highest structural similarity to existing compounds with known efficacy against these cancer types. Despite no direct optimization of other pharmacological properties, we report good agreement with known cancer drugs in metrics like drug-likeness, synthesizability and solubility. We envision our approach to be a step towards increasing success rates in lead compound discovery and finding more targeted medicines by leveraging the cellular environment of the disease.

## 1 Introduction

Eroom's Law describes the observation that the productivity of the drug discovery pipeline, as measured by the number of FDA approved drugs per invested billion US dollar, halves every

nine years since the 1950s (Scannell et al., 2012). Indeed, only a minimal portion of synthesized drug candidates obtain market approval (less than 0.01%), with an estimated 10-15 years until market release and costs that range between one (Scannell et al., 2012) to three billion dollars per drug (Schneider, 2019). This low efficiency has been attributed to the high dropout rate of candidate molecules in the early stages of the pipeline, highlighting the need for more accurate in silico and in vitro models that produce more potent candidate drugs. In addition to the initial wet-lab validations, the discovery pipeline involves a sequential process that builds upon high-throughput screenings, ADMET-assessments and a lengthy phase of clinical trials. The costs of the experimental and clinical phase can be prohibitive and any solution that helps to reduce the number of required experimental assays can provide a competitive advantage and reduce time to market. The problem's linchpin is on how to improve the exploration and navigation through the chemical space that has been estimated to contain $\sim 10^{30}$-$10^{60}$ drug-like molecules with bioactive properties (Polishchuk et al., 2013). For this task, deep learning methods have recently gained popularity (Chen et al., 2018) and many have demonstrated the feasibility of in silico design of novel candidate compounds with desired chemical properties (Popova et al., 2018; Gomez-Bombarelli et al., 2018; You et al., 2018). In all of these models, the generative process is controlled via a structurally driven evaluator (or critic) that biases the generation of a chemical to satisfy the required chemical structural properties. While very effective in generating compounds with desired chemical properties, these methods disregard system-level information, e.g. about the cellular environment in which the drug is intended to act. However, the two main causes of the increasing attrition rate in drug design are lacking efficacy against the specific disease of interest and off-target cytotoxicity (Wehling, 2009), calling to bridge systems biology closer with drug discovery. To this end, we present a novel framework for anticancer molecule generation based on deep generative models and reinforcement learning that, for the first time, enables generation of candidate compounds while taking into account the disease context encoded in the form of gene expression profile (GEP) of the tumor cell (for a graphical illustration see Figure 1A).

Related methodology has been used for protein-targeting de-novo generation (Zhavoronkov et al., 2019; Aumentado-Armstrong, 2018; Grechishnikova, 2019). These contributions attempt to utilize deep learning methods for de-novo design of compounds to specifically target a protein that has been implicated in tumor proliferation or treatment response (e.g. gene-knockout study). For example, the study by Zhavoronkov et al. (2019) curated and utilized, amongst others, patent data and several datasets about molecules (unspecific bioactive compounds, kinase inhibitors, DDR1 kinase inhibitors, molecules targeting non-kinase targets) specifically to develop DDR1 inhibitors. They synthesised and tested six drug candidates in cell assays, where two were found to be active and one was even successfully validated in animal models. Envisioning a precision or even personalized medicine perspective, identifying protein targets is challenging, whereas sequencing and omics data are straightforward to gather.

Very recently, Méndez-Lucio et al. (2020) proposed a method for de-novo design of molecules against desired targets, represented by the gene expression signatures of knocked-out (suspected) targets. However, 97% of anticancer candidate drugs fail in clinical trials and never receive FDA approval, questioning the current approaches of target identification for the discovery of pharmaceuticals (Wong et al., 2019). Taking as many as 10 drug-indication-pairs from ongoing clinical trials, Lin et al. (2019) found that the proposed mechanism of action (MOA) of *all* of them were incorrect; knocking out the target genes did not ever hamper cancer fitness. While the wrong target genes were identified through RNA interference with siRNA seemingly silencing essential off-target genes, all drugs retained their anticancer effect through target-independent mechanisms. Off-target cytotoxicity being a common MOA of anticancer drugs in clinical trials corroborates to scrutinize current lead compound discovery strategies and calls to develop novel methodology with unconventional approaches. It is for this

reason that we herein propose a novel framework to generate lead compound candidates solely based on a tumour's metabolic signature, as opposed to attempting to target a specific protein or incorporating information about potential targets directly into the design process. Acting as metabolic signature, we instead guide the learning process solely by transcriptomic profiles of cancer cells, since transcriptomic data has been successfully used for de-novo drug identification (Verbist et al., 2015; De Wolf et al., 2018) and has been advocated for a pivotal role for the future of drug discovery (Dopazo, 2014).
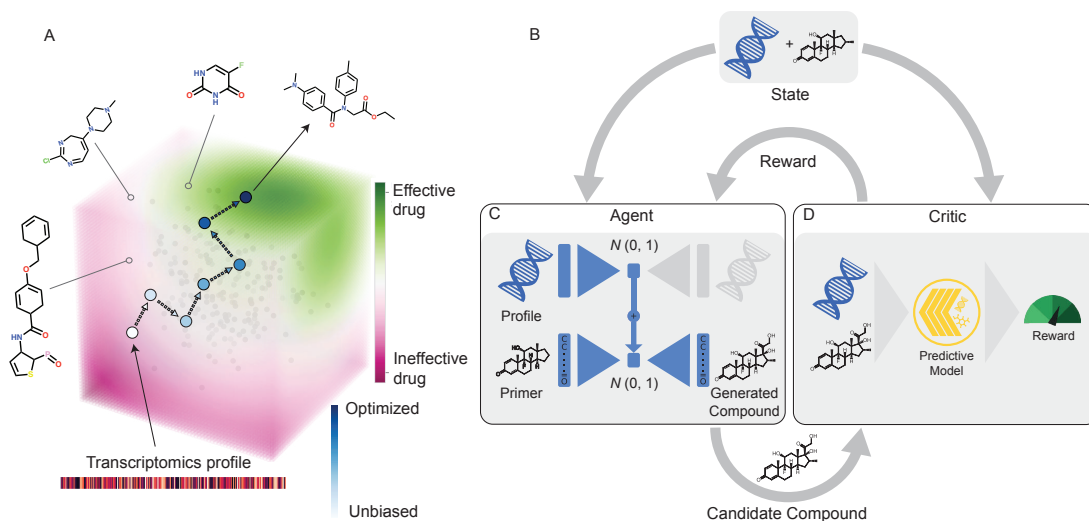


Figure 1: **The proposed framework for anticancer compound design against specific cancer profiles. A)** Conceptually, the model performs a guided walk through the chemical space in order to find effective compounds. Starting from an unbiased molecule generator (trained only on a dataset of bioactive compounds without any information about cancer), compounds are sampled and screened in-silico against the transcriptomic profiles of interest. The outcome of the screening guides the generator towards sampling from manifolds with more effective compounds. **B)** The training process depicted in more detail. The conditional compound generator (called "agent") is embodied through two, initially separate VAEs. The compound generation starts with a biomolecular profile of interest e.g. a transcriptomic profile. Through a pretrained omics VAE, the profile is encoded into the latent space of gene expression profiles. The latent representation of the profile is decoded through the molecular decoder of a separately pretrained molecule VAE to produce a candidate compound (see **C**). This generative process can optionally be "primed" through encoding a known, effective compound or a functional group with the molecular encoder. The proposed compound is then evaluated by a critic, where our critic is represented by a multimodal drug sensitivity prediction model that ingests the compound and the target profile of interest (see **D**). The IC50 efficacy as predicted by the critic, is interpreted as reward and is subject to maximization during the RL based optimization. Over the course of training, the generator will thus learn to produce candidate compounds with higher and higher efficacy. `<START>` is the start and `<END>` is the end token.

Our framework is depicted in Figure 1B and consists of a conditional molecule generator (embodied by two separate VAEs) and a critic module that evaluates the efficacy of proposed compounds on the target profile (see Figure 1D). The training procedure splits into two stages. In the first stage, the models are trained independently; one VAE is trained on gene expression data from TCGA (Weinstein et al., 2013), another VAE is trained on bioactive small molecules from ChEMBL (Bento et al., 2013) (see Figure 1C). As critic, a multimodal drug sensitivity prediction model is fetched from previous work (Manica et al., 2019). In the second stage, the encoder of the profile VAE is combined with the decoder of the molecule VAE and exposed to a joint retraining that is optimized in a policy gradient regime with a reward coming from the critic module. The goal of the optimization is to

tune the generative model such that it generates (novel) compounds that have maximal efficacy against a given biomolecular profile that is characteristic for a cancer site, a patient subgroup or even an individual. By *efficacy*, we refer to predicted cellular IC50 (i.e. the micromolar concentration necessary to inhibit 50% of the cells) as opposed to e.g. enzymatic IC50. Importantly, this efficacy is a joint property of a drug-cell-pair; as treatment response to a compound heavily varies depending on the tumor's genomic and transcriptomic makeup (Geeleher et al., 2016). In this work, we emphasize profile-specific compound generation and optimize the generator with IC50 as sole critic.

Here, we present a first step towards our vision of a generic framework where the molecule generation can be conditioned on possibly multimodal context information such as a (multi)omics profiles, a primed drug or a drug scaffold, a target protein or any combination thereof. The resulting latent spaces, consisting of semantically distinct ontological entities, could then be jointly explored by machine learning techniques designed to operate on sets instead of fixed-length vectors, such as permutation-invariant operations (Zaheer et al., 2017).

## 2    Results

**Pretraining Profile VAE and SMILES VAE.**    In the first phase of training, the two components of Figure 1C were trained independently. The profile VAE (PVAE) consisted of a set of stacked dense layers and was trained as a denoising VAE to enhance generalization abilities. The purpose of the PVAE was to find a lower dimensional representation of the cell profiles that maintains structural similarity and later allows a fusion with the latent representation of molecules. The encoder of the PVAE learned to embed gene expression profiles (bulk RNA-Seq from TCGA (Weinstein et al., 2013)) meaningfully into a latent space, such that the decoder could reconstruct the profiles, but also generate novel, seemingly realistic gene expression profiles (GEP).

The SMILES VAE (SVAE) was pretrained for 10 epochs on ∼1.4 million structures from ChEMBL (Bento et al., 2013). Both encoder and decoder consisted of stack-augmented gated-recurrent units as used in (Popova et al., 2018). The purpose of the SVAE was to learn the syntax of SMILES and general semantics about bioactive compounds. The novelty and diversity of the generated molecules was validated by sampling 10,000 molecules through decoding random points from the latent space. 96.2% of the 10,000 generated molecules were valid molecular structures (surpassing the results of Popova et al. (2018) who used the same stack-augmented GRUs and reported 95% SMILES validity) and 99.72% of the valid molecules were unique across the 10,000 generations. Comparing the Tanimoto similarity (i.e. the jaccard index, a well-established chemical structure similarity measure Tanimoto (1958)) of the molecular fingerprints (ECFP Rogers and Hahn (2010)) of 1000 generated molecules with the training and test data from ChEMBL, we find that the vast majority had a Tanimoto similarity ($\tau$) between 0.2 and 0.6 (on average $0.41 \pm 0.1$ for training and $0.38 \pm 0.08$ for testing molecules) suggesting that our model learned to propose novel molecular structures from the chemical space of about $10^{30}$ to $10^{60}$ molecules (Polishchuk et al., 2013). In addition, a visualization of the chemical space of ChEMBL as well as generated compounds through the TMAP algorithm (a library that visualizes high dimensional data through minimum spanning trees (Probst and Reymond, 2019)) showed that the generated molecules mix well with the training molecules into the chemical space.

For detailed results of both the PVAE and the SVAE, please see the appendix (S5).

## 2.1 Disease-specific compound generation

Herein, we present the results of our drug generator conditioned on gene expression profiles of cancer subtypes. As a proof of concept, we show results for cancer in four different sites: breast
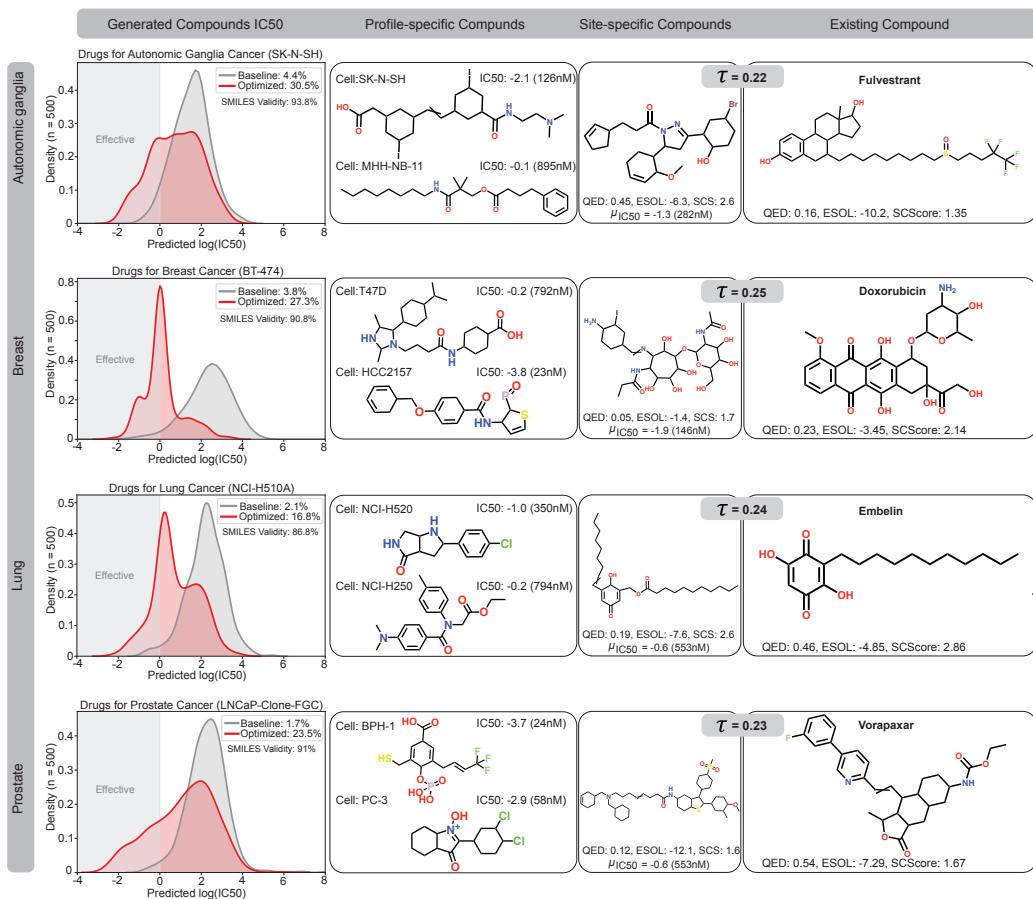


Figure 2: **Sample results for profile-driven model optimization and anticancer compound generation.** Each row illustrates the results of training the RL pipeline on cell lines from a specific cancer type: autonomic ganglia, breast, lung and prostate cancer. The first column compares the distributions of IC50 predictions given by the critic model for a set of $n$=500 drug candidates generated with RL optimization and without RL optimization. As demonstrated by the density plots, the RL optimization process leads to candidate compounds with a lower mean IC50 for the target cancer – highlighting the successful optimization of the generative model towards the design of more effective compounds. The second column presents candidate compounds with a high predicted efficacy (low IC50) against a particular cell line that was not seen during training (this corresponds to a "personalized medicine" regime) The third column showcases generated compounds that were optimized to be effective against all cell-line profiles of the given cancer type in each row (corresponds to a "precision medicine" regime). In the fourth column, we present an *existing* anticancer compound (approved against at least one type of cancer), that was in the top-3 neighborhood of the generated compound in the third column. The existing and generated compounds are compared in terms of Tanimoto structural similarity as well as three chemical scores crucial in drug design namely, druglikeness (QED, 0 worst, 1 best), synthetic complexity (SCS or SCScore, 1 best, 5 worst) and solubility (ESOL, given in $M/L$).

(carcinoma), lung (carcinoma), prostate (carcinoma) and autonomic ganglia (neuroblastoma). The conditional generator was initialized as the SVAE, i.e. sampling from the unbiased generator yielded

random molecules from the chemical space as learned from the ChEMBL data. For the evaluation, all generated compounds with a predicted IC50 value below $1\mu M$ were considered to be *effective*. Moreover, within each cancer type (or site), 80% of the cell lines (breast: 50, lung: 169, prostate: 7, autonomic ganglia: 56) were considered as training cell lines and used to optimize the parameters $\Theta$ of the conditional generator. We observed that over time the generator learned to produce more drugs with higher predicted efficacy according to the critic. To test both the generalization abilities and whether the generator actually utilized the omics-profile for the generation, we used the remaining 20% of cell lines to verify whether conditioning the generator on unseen cell lines of the same site also leads to compounds with low IC50. As presented in Figure 2 (left column), our model learned to produce compounds with lower IC50 values, for unseen cell lines from the given cancer site. In other words, the IC50 distribution of candidate compounds proposed by the generative model were successfully shifted towards higher efficacy (lower IC50). The baseline model corresponds to the pretrained SVAE from which $n = 500$ molecules were randomly sampled. In all four cases, a significant portion (between 17% and 30%) of molecules generated from the optimized model were assigned a IC50 value below $1\mu M$, whereas only 1-4% of the candidates generated by the baseline model (i.e. the SVAE) were classified as effective. Moreover, in all cases the generator maintained almost an equal SMILES validity (87%-94%) compared to the baseline, much higher than what Méndez-Lucio et al. (2020) reported based on gene expression (8-9%). The second column of Figure 2 shows generated molecules that are predicted as being effective against unseen cell lines from the respective cancer site. As opposed to the personalized regime in the second column, the third column of Figure 2 showcases a precision medicine regime. Here, novel molecules were designed specifically for each cancer site i.e. a single, characteristic GEP. In all cases, the model generated compounds that exhibited high efficacy against the average cellular profile of the target site while maintaining efficacy against the majority of individual cell lines for that site.

**Investigation of nearest neighbors**  For a more quantitative assessment, the last column of Figure 2 compares the four cancer type-specific candidate compounds with one of their top-3 neighbors using the Tanimoto similarity score, $\tau$, from several hundreds of existing anticancer compounds. It is well known that Tanimoto similarity across compounds is highly correlated with their induced sensitivity patterns on cancer cell lines (Shivakumar and Krauthammer, 2009). The candidate compound proposed against breast cancer (Figure 2 second row, third column) has Doxorubicin, a commonly used chemotherapeutical against breast cancer (Lao et al., 2013), as one of the top-3 nearest neighbors. The generated compound against lung cancer (Figure 2, third row, third column) presents similarities to Embelin, an existing anticancer compound from the GDSC database. Comparing the two structures, it is evident that the generated compound and Embelin share a long carbon chain and a single six-membered fully carbonic ring. Embelin was tested against 965 cell lines from GDSC/CCLE from which the highest reported efficacy is against a lung cell line (NT2-D1). Embelin is also known to be the only known non-peptide inhibitor of XIAP (Poojari, 2014), a protein that plays an important role in lung cancer development (Cheng et al., 2010). The closest neighbor of the prostate-specific generated compound (Figure 2 fourth row, third column) is Vorapaxar. Its efficacy is highest against a prostate cancer cell line (DU_145) according to GDSC/CCLE. Vorapaxar is an antagonist of a protease-activated receptor (PAR-1) that is known to be overexpressed in various types of cancer, including prostate (Zhang et al., 2009). Lastly, the third closest neighbor of the generated compound against neuroblastoma (Figure 2 first row, third column) is Fulvestrant, an antagonist/modulator of $ER\alpha$ which has recently been proposed as a novel anticancer agent for neuroblastoma (Gorska et al., 2016). To summarise, for all four investigated cancer types, the proposed compounds some high structural similarity to anticancer drugs that are, for each specific cancer type, either 1) already FDA approved (breast), 2) known inhibitors of relevant targets (lung,

prostate) or 3) have been advocated for (neuroblastoma). This result is remarkable, specially as the generator was never exposed to any anticancer compounds. Indeed only the critic had seen two out of the four compounds during training, highlighting the fact that the generator has *de novo* learn the structural characteristics that make a compound efficacious against a particular cancer type.

In the above literature review, the search space was restricted to compounds with known anticancer properties. To investigate whether the proposed compounds had generally a higher similarity to drugs associated with cancer, we carried out a comparison with compounds from a broader pool of chemicals, namely ChEMBL (Bento et al., 2013), a database of $> 1.5$ million bioactive molecules with drug-like properties. The nearest neighbour ($\tau = 0.54$) of our breast compound in the ChEMBL database is CHEMBL1093122, a conjugate of plumbagin and phenyl-2-amino-1-thioglucoside that inhibits the synthesis of mycothiol (Gammon et al., 2010). Plumbagin itself and many of its derivatives are widely studied anti breast cancer compounds (Zhang et al., 2016; Kawiak et al., 2017; Dandawate et al., 2014). For our lung cancer compound, the nearest neighbor ($\tau = 0.48$) is polyoxyethylene dioleate, a surfactant that has been patented for the treatment of eight types of cancer including three types of lung cancer (small lung cell cancer, lung adenocarcinoma and metastatic lung cancer, Girsh (2007)). It is also utilised in targeted drug delivery systems for drug-resistant lung cancer (Kaur et al., 2016). The nearest neighbour ($\tau = 0.31$) of the prostate cancer compound is Clinolamide (or Linolexamide) which is included in a patent of diagnostic and/or therapeutically active compounds for several types of cancer, including prostate cancer (Klaveness et al., 2004). The nearest neighbor ($\tau = 0.35$) of the proposed neuroblastoma compound is NSC-715466. NSC-715466 has been evaluated for anticancer effects in the NCI-60 database (Shoemaker, 2006) and inhibits cancer cell growth by $65\% \pm 15\%$ across all tested cell lines, with a below-averge inhibition for cancer in the central nervous system ($57\% \pm 9\%$). Regarding its efficacy, it only falls in the 51st percentile of all 53,217 compounds tested in NCI-60, which presumably prevented further investigations. The four discussed ChEMBL compounds as well as the analysis of NSC-715466 can be found in the appendix (S6). It is promising to observe that the molecules with the highest Tanimoto similarity to our compounds are associated with cancer (some even to the specific types of cancer our compounds were optimized for) even using a larger database of bioactive compounds. However, it is worth keeping in mind that a high Tanimoto score to a known cancer drug is not necessary for anticancer drug efficacy, as some cancer drugs used for the same cancer type or even sharing the same mechanism of action exhibit low Tanimoto similarity, e.g. TKIs used for NSCLC such as Crizotinib and Erlotinib, $\tau = 0.11$. Across all anticancer compounds in GDSC and CCLE databases, we note that the average Tanimoto similarity ($\tau = 0.149 \pm 0.05$) is not much lower than the average similarity of two compounds of a given site ($\tau = 0.154 \pm 0.06$). For the following results and discussion, the anticancer compounds from GDSC/CCLE were associated with the site where they had the highest average IC50 efficacy.

To understand better whether the generated drugs mimic the space of cancer-specific anticancer drugs, Figure 3 shows visualizations of real and generated cancer drugs for one specific cancer type, using kernel PCA based on Tanimoto similarity (Schölkopf et al., 1998). In addition to the class belongings, the plots also depict the QED score (Bickerton et al., 2012), a quantitative estimate of drug-likeness (0 worst, 1 best), and SCScore (Coley et al., 2018), an estimated score of synthetic complexity (1 best, 5 worst). The fact that the generated molecules are well intermingled with the real drugs suggests that the generator proposes diversified structures that mimic some properties of anticancer compounds. It is also curios to see that several real drugs have low QED and/or high synthetic complexity scores (the same holds for the generated molecules). Moreover, we provide interactive TMAP visualizations of the site-specific, generated compounds (links in availability section).
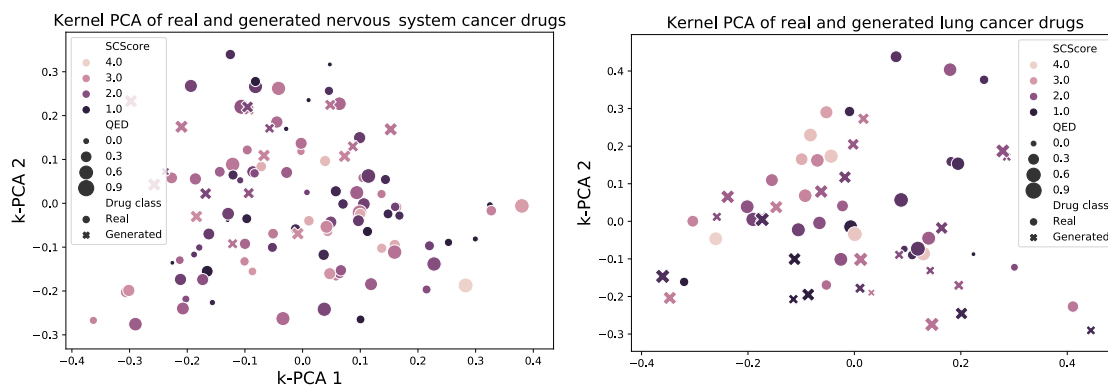
Figure 3: **Visualization of generated and real anticancer drugs.** A kernel PCA of real and generated molecules based on Tanimoto similarity. The size of the points is denoted the QED score while the coloring represents the synthetic complexity score (SCScore). Overall, both generated and existing molecules are heterogeneously distributed in the 2D projection and do not form clear clusters.

**Chemical properties of generated molecules**    In this work, the conditional generator is trained using PaccMann as sole critic. However, besides inhibitory efficacy, there is a myriad of properties of a candidate drug that crucially influence its potential for becoming an anticancer compound.

Some of these can be approximated in-silico, e.g. water solubility (ESOL, Delaney (2004)), drug-likeness (QED, Bickerton et al. (2012)) and synthesizability (SCScore, Coley et al. (2018)). Figure 4 gives an overview about the distribution of QED, ESOL and SCScore for sets of 1) known anticancer compounds (blue), 2) molecules from ChEMBL (orange), 3) compounds generated by the SVAE (red) and 4) compounds proposed by the conditional generator (green). Despite none of these properties was explicitly optimized, comparing the distributions reveals overall a good agreement. Interestingly, anticancer drugs exhibit, compared to the ChEMBL compounds, on average much less drug-like properties (lower QED) and seem easier to synthesize (lower SCScore). This tendency of anticancer drugs for synthetically less complex structures is likely to result from the high attrition rate in clinical trials and the corresponding cost reduction policies. It is also encouraging to see that the unbiased generator (SVAE) generates on average molecules with more desired properties compared to the data used for training (ChEMBL compounds have on average lower QED and higher SCScore).

Moreover, the cancer drugs have, on average a significantly lower QED than the other three sets ($0.45\pm0.2$ with $> 10\%$ of GDSC/CCLE drugs even having a QED $< 0.2$ whereas it is 0.55 for the other three sets). Indeed the QED scores of the other three sets were so similar that we failed to reject the null hypothesis that the QED scores of these three sets are from different distributions (Kruskal-Wallis test, $\alpha = 0.05$). Regarding synthetic complexity, both the unbiased and the biased generator fail to produce molecules with SCScores as low as the anticancer drugs (MWU, $p < 0.01$), but they produce structures that are estimated to be less complex than the ChEMBL molecules (MWU, $p < 0.01$). Overall, it can be seen that the biased generator produces molecules with less desired properties than the unbiased generator (SVAE). This is expected because the unbiased generator was pretrained to mimic the data from ChEMBL whereas no explicit optimization of chemical scores was performed during the RL optimization. For one generated compound, we exemplary show a possible synthesis route that was assigned a high confidence score by the retrosynthesis model (RXN, 2020) and consists of four reactions and 10 commercially available reactants (see appendix (S8)).

Savjani et al. (2012) reported that 40% of novel chemicals cause practical problems due to
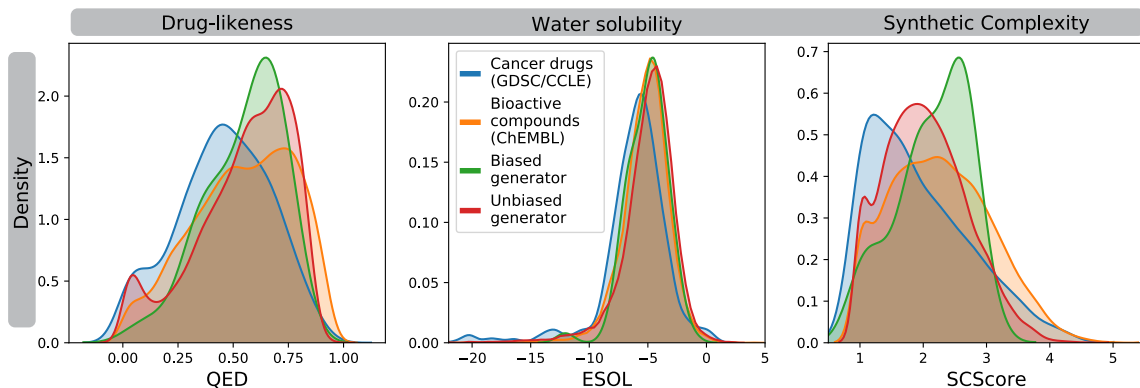
Figure 4: Comparison of chemical scores for real drugs in GDSC and CCLE database versus our generated compounds. We compared three chemical scores for druglikeness as assessed by QED score (0 worst, 1 best), for solubility as assessed via ESOL, given in $\log(M/L)$ (most drugs have a solubility between -8 and -2) and for synthetic accessibility as assessed by SAS (1 best, 10 worst). These three scores are computed for the panel of known anticancer drugs, bioactive molecules from ChEMBL, molecules generated before (red) and after (green) RL optimization.

insolubility. Water solubility remains challenging to approximate in-silico (Sorkun et al., 2019) and thus we treat with caution the good agreement in the ESOL scores Figure 4 (middle panel).

Finally, we would like to point out that utilizing a IC50 drug sensitivity prediction model as sole critic limits the performance of the entire pipeline, as the expressive power of the conditional generator is inherently upper bounded by the predictive power of the critic. Crucially, due to a lack of available data, the critic was only trained on anticancer drugs but not on compounds without inhibitory efficacy against cell lines. We therefore verified the generalization capabilities of the critic by comparing the predicted efficacy of cancer drugs (most of them were seen during training) and a "negative" set of molecules from ChEMBL across all 965 cell lines from GDSC. The results can be found in appendix S7 and show that while 15.2% of the virtual drug screenings with anticancer compounds were positive (IC50< 1$\mu$mol), only 2.2% of ChEMBL molecules showed potential anticancer effects. Moreover, the generated anticancer compounds were found to have a significantly higher Tanimoto similarity to anticancer drugs than to both ChEMBL molecules ($p < 0.01$, one-sided MWU) and molecules generated without the RL optimization, i.e. from the SVAE ($p < 0.01$, one-sided MWU). These two results are encouraging. They suggest that PaccMann can seemingly drive the molecule generation away from ordinary, bioactive compounds as in ChEMBL more towards mimicking the properties of actual anticancer drugs.

## 3 Discussion

We herein presented the first framework for anti-cancer compound generator that enables us to condition the molecular generation on the biomolecular profile (specifically we explored transcriptomic profiles) of the target cell or cancer site. We demonstrated, using a RL optimization framework, that our proposed generative model could be optimized to produce candidate compounds with high predicted efficacy (IC50) against a given target profile, even if this profile was never seen during training. Notably, this was achieved despite the fact that the generator was never exposed to anticancer drugs explicitly, but only pretrained on bioactive compounds from ChEMBL. The only component that has been trained on drugs with known anticancer effects is the critic, which only

communicates with the generator by providing a reward function.

An analysis of the generated compounds for four different cancer types demonstrated that the predicted compounds share many structural similarities with known anticancer compounds for the same cancer types that the generated compound was optimized for.

While our results are a promising stepping stone for profile-specific anticancer compound generation, further optimization must be done before it can be used a reliable tool for drug discovery. For instance, other properties of a candidate drug other than inhibitory efficacy that determine its potential for becoming a successful anticancer compound, for example water solubility, drug-likeness, synthesizability, environmental toxicity or off-target cytotoxicity are not directly optimized. However, despite not explicitly incorporating them into the reward function, we find that the produced molecules exhibit desired properties in terms of drug-likeness, water solubility and ease of synthesis.

Furthermore, the high attrition rate in drug discovery has been attributed to either a lack of efficacy or off-target cytotoxicity (Wehling, 2009), the latter implying that very often the mechanism of action of an active compound has been incorrectly characterized. In that respect, the design of new AI-enhanced drug design approaches that can bypass the need of a detailed characterization of drug targets and cytotoxicity mechanisms can greatly improve current drug discovery pipelines. Future work should focus on incorporating information in the reward function not only about drug efficacy but also about other drug-relevant chemical properties and predicted off-target cytotoxicity effects. The resulting multimodal objectives may be difficult to optimize due to possibly counteracting/interfering gradients. A possible approach to circumvent this challenge is by using explicit compensation techniques (Yu et al., 2020) or defining gradient-free global objectives (Häse et al., 2018). Another challenge that needs to be overcome to improve the reliability and accuracy of the critic is the expected distributional differences between the data used for training (cancer cell lines) and the targeted data (human data from clinical trials). A possible approach that can be explored is the exploitation of transfer learning techniques, as suggested in Sharifi Noghabi et al. (2020).

Oftentimes, medical chemists do not start the drug design from scratch, but from the scaffold of an approved drug. The goal of the scaffold hopping is to find a drug with similar effects (e.g. increased efficacy or reduced side effects). While our framework enables users to incorporate prior knowledge into the design process by priming the latent code, we have not yet explored the full potential of this idea herein. As the decoded molecule is not guaranteed to maintain similarity to the primer, the recently proposed "deep generative scaffold decorator" could be integrated into our framework to facilitate a more systematic exploration and the possibility of adding fragments to established drug scaffolds (Arús-Pous et al., 2020). An alternative future endeavor is to explore graph-based instead of sequential representations of molecules to directly generate a molecular graph from the context set using a conditional structure generation framework (Yang et al., 2019).

## 4 Methods

**Conditional generator ($\mathcal{G}$).** Our conditional generator is a molecule generator that produces a candidate drug structure represented as its SMILES sequences (Weininger, 1988). SMILES sequences are preferable over (functional) fingerprint-based representations of molecules (e.g. ECFP (Rogers and Hahn, 2010)) since they have shown to be superior in both predictive (Jastrzebski et al., 2016; Manica et al., 2019) and generative models for molecules (Bjerrum and Sattarov, 2018). In our use case, the generative process is conditioned on a target biomolecular profile, e.g. from a patient or a disease. Inspired by Gomez-Bombarelli et al. (2018), we concluded that variational autoencoders (VAE) are the ideal generative model for our task since by design they bring about a structurally ordered latent space that simplifies the combination of different information sources. Our conditional generator combines two VAEs that are trained independently prior to being fused together: 1) a denoising VAE for cancer profile encoding/generation (called PVAE, Figure 5A) and 2) a sequential VAE (SVAE) for SMILES sequence generation (Figure 5B). The mathematical formulation of VAEs can be found in Kingma and Welling (2013); Sohn et al. (2015). PVAE is pretrained on gene expression profiles (GEP) to learn a consistent latent representation for biomolecular signatures. SVAE is pretrained on
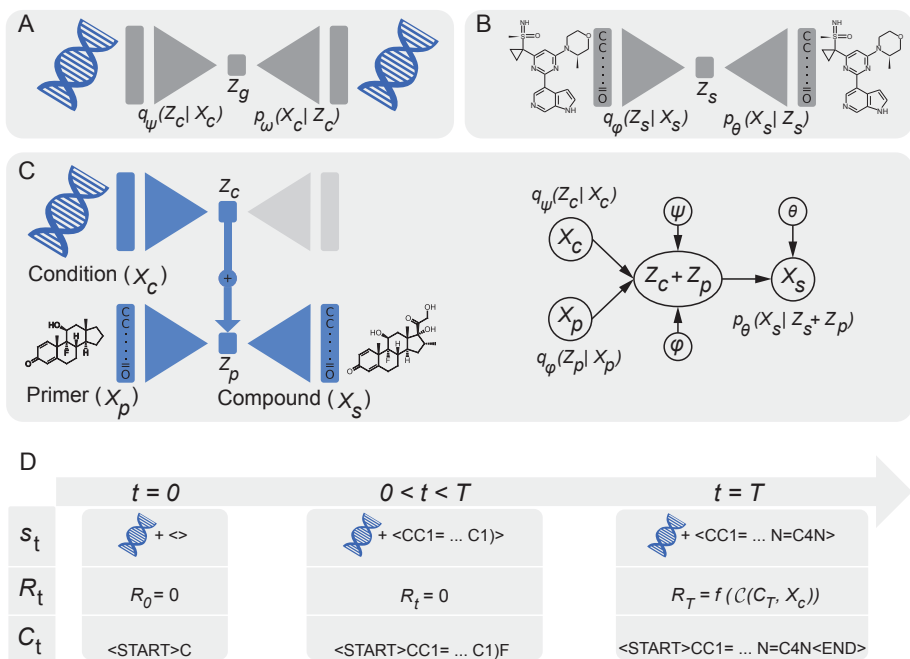
Figure 5: **Architectural details of the conditional drug generator.**
**A)** A biomolecular profile VAE (PVAE) was pretrained on RNA-Seq data from TCGA to encode a transcriptomic profile $X_c$ into a latent code $Z_c$, before attempting to decode $X_c$ from it. **B)** Similarly, a sequential compound generator VAE (SVAE) was trained to encode and decode SMILES representations $X_s$ of molecules. **C)** PVAE and SVAE are combined to obtain a conditional molecule generator. As shown in the graphical model, the combination is achieved by using a permutation-invariant operation (e.g. addition) to fuse the latent spaces of omics profiles and molecules to a joint, multimodal representation. **D)** Molecules are generated directly as SMILES sequences and are assembled in a sequential process, one token at a time. A full cycle of this process includes a state ($s_t$, where $s_0 = X_c$, i.e. a TCGA RNA-Seq transcriptomic profile), a reward ($R_t$) and a generated candidate compound ($C_t$).

bioactive drug-like molecules to learn the syntax of valid SMILES and general molecular semantics. Generative models of SMILES sequences necessitate the ability to *count* Models that process SMILES sequences greatly benefit from the ability to *count* the ring opening and closing symbols in a molecule, as a single mistake in the sequential generation of a SMILES renders the entire string invalid. To circumvent that standard recurrent and convolutional networks lack the proficiency to count, we utilize a stack memory (Hopcroft and Ullman, 1969), in our case implemented through stack-augmented GRUs as proposed by Joulin and Mikolov (2015) (for equations and other details of the SVAE and the stack see the appendix (S1)). Thereafter, the encoder of the PVAE is fused with the decoder of the SVAE via their latent space (Figure 5C). The combination of the two models enables to learn a latent space that links biomolecular profiles and chemical structures providing an effective way to sample novel compounds given a specific GEP. In the RL optimization phase, the weights of the fused model (which were pretrained independently) are fine-tuned using a reward from the critic.

**Critic ($\mathcal{C}$).** The critic is a multimodal drug sensitivity prediction model that evaluates the efficacy of any given candidate compound against a biomolecular profile of interest, e.g. gene expression of a cancer cell line. The critic outputs a non-negative reward that depends on the candidate compound predicted IC50 for the target profile, such that low IC50 values associated with higher compound efficacy receive higher rewards than high IC50 values. The reward is then used in a RL framework to update the conditional generator. Following the most recent advances for multimodal drug sensitivity prediction we herein utilize *PaccMann* as a critic Manica et al. (2019).

**The RL framework.** The conditional generator is retrained in combination with the critic in a RL-based optimization process to tailor molecules towards the given GEP. First, the GEP is encoded into a latent space, $Z_c$ (see

11

Figure 5C). This embedding is then added to the latent encoding of a primer compound or substructure ($Z_p$). The advantage of using a primer is that it enables injection of prior knowledge into the model by starting the generative process from an existing and proven effective compound or functional group – instead of designing a compound from scratch. However, this priming is optional and we do not only sample closely around existing compounds but we instead sample a larger fraction of the chemical space. As can be seen in the graphical model in Figure 5C the molecule generation is conditioned on a context $\mathcal{Z}$, where in this work $\mathcal{Z} = \{Z_c, Z_p\}$. $Z_c$ and $Z_p$ reflect embeddings learned from semantically different data modalities (gene expression and molecules). To combine these (latent) representations, we use summation because it is a permutation invariant operation and has been proposed to combine a variable set of unstructured latent encodings (Zaheer et al., 2017). Alternatives include mixup functions such as weighted sums or dimension-wise sampling from a categorical (Bernoulli) distribution (Beckham et al., 2019). Our additive latent representation is similar in concept to the conditional VAE with additive Gaussian encoding space (Wang et al., 2017). Intuitively, this fusion presumably warps the latent space from encoding structural similarity (of molecules or GEP) into functional similarity so as to aggregate molecules with similar predicted efficacy for a given cell line (Gomez-Bombarelli et al., 2018). Note that using a primer compound or substructure is optional and if no priming compound is used, simply the latent space representation of the `<START>` token is added to the latent encoding of the target GEP.

Next, the conditional generator decodes the latent encoding, $Z_c + Z_p$, and generates a molecular structure that, in combination with the GEP, is fed to the critic to produce a certain reward for the generated compound, as illustrated in Figure 5C. Following the notation of Popova et al. (2018), the conditional generator, $\mathcal{G}$, acts as the *agent* and PaccMann (the multimodal IC50 prediction model, $\mathcal{C}$) represents the *critic*. The weights of $\mathcal{C}$ are fixed. We aim to optimize $\Theta$, the parameters of $\mathcal{G}$, to produce candidate compounds, $C_T$, that target a specific GEP, $X_c$. In contrast to Popova et al. (2018), we define the set of states $\mathcal{S}$ as all possible SMILES strings (with length $\leq T$) paired with the target GEP. The set of possible actions $a$ that $\mathcal{G}$ can take is a set $\mathcal{A}$, which is a vocabulary of all characters and symbols of the canonical SMILES language. As depicted in Figure 5D, molecules are generated by $\mathcal{G}$ by sampling an action $a_t$ at each step($0 < t < T$) from $p(a_t|s_{t-1})$, where $s_{t-1} = (C_{t-1}, X_c)$ and $C_0$ is simply the `<START>` token. Terminal states $S^* \subset S$ are reached when either $t = T$ or when the terminal action $a_T = $ `<END>` has been sampled. $\mathcal{G}$ is trained to learn a policy, $\Pi(\Theta)$, by maximizing:

$$\Pi(\Theta) = \sum_{s_T \in S^*} P_\Theta(s_T) R(s_T) \tag{1}$$

where $P_\Theta(s_T) := \prod_{t=0:T} p(a_t|s_{t-1})$ and the state $s_T = (C_T, X_c)$ is a tuple of the candidate compound $C_T$ and the cell profile $X_c$ and the reward $R(s_T) = f(\mathcal{C}(C_T, X_c))$ is the output of the critic $C$ scaled by a reward function $f$. In our experiments, all intermediate rewards $R(s_t)$ are set to 0, since $C_t$ (the intermediate SMILES string) will in almost all cases not resemble a valid molecule. The sum is approximated using policy gradients, specifically the REINFORCE algorithm (Williams, 1992) and the reward function $f$ for determining the reward from the IC50 prediction, $\mathcal{C}(C_T, X_c)$, is computed by $f(IC50) = \exp\left(\frac{-IC50}{\alpha}\right)$ where $\alpha = 5$ in this work (see details in the appendix (S2)).

**Data.** For the PVAE, we employed a training dataset of 11,592 (standardized) RNA-Seq GEPs from healthy and cancerous human tissue from the TCGA database and validated it on 1,289 samples from the same database (Weinstein et al., 2013). Since the dataset was too small to train on the full cohort of 20,000 genes and most genes are correlated to a subset of landmark genes, the number of genes was reduced to the same 2,128 genes as used in Manica et al. (2019), following the network propagation procedure described in Oskooei et al. (2019). The SVAE was pretrained on the SMILES representation of 1,576,904 compounds (10% were held out for performance validation) from the ChEMBL database (Gaulton et al., 2016). For RL optimization of $\mathcal{G}$, we used GEPs publicly available from GDSC (Yang et al., 2012) and CCLE (Barretina et al., 2012) databases. Since the RNA-Seq of these cancer cell line databases were passed through the PVAE (pretrained on human samples from TCGA (Weinstein et al., 2013)), we compared the standardized gene expression distributions for the selected genes across these databases and found good agreement (see appendix (S3)), in alignment with the reported consensus between transcriptomic data in CCLE and TCGA (Ghandi et al., 2019). To train the critic ($\mathcal{C}$), IC50 drug sensitivity data from GDSC and CCLE was utilized. The hyperparameter and details on the utilized hardware and software can be found in the appendix (S4).

# 5 Availability of software and materials

The omics data used to pretrain the PVAE, the molecular data for the SVAE and the cell profiles used in the RL regime as well as the pretrained models can be found on https://ibm.box.com/v/paccmann-pytoda-data. To assess the critic, please see Manica et al. (2019). All code to reproduce the experiments is publicly available on https://github.com/PaccMann/. For a detailed example see https://github.com/PaccMann/paccmann_rl.

The interactive TMAP visualizations of the molecules generated by the (unbiased) SVAE can be found on https://paccmann.github.io/rl/unbiased.html. The TMAPs of the cancer-site-specific candidate compounds are accessible on https://paccmann.github.io/.

# Acknowledgements

# References

Josep Arús-Pous, Atanas Patronov, Esben Jannik Bjerrum, Christian Tyrchan, Jean-Louis Reymond, Hongming Chen, and Ola Engkvist. Smiles-based deep generative scaffold decorator for de-novo drug design. 2020.

Tristan Aumentado-Armstrong. Latent molecular optimization for targeted therapeutic design. *arXiv preprint arXiv:1809.02032*, 2018.

Jordi Barretina, Giordano Caponigro, Nicolas Stransky, Kavitha Venkatesan, et al. The cancer cell line encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature*, 483(7391):603, 2012.

Christopher Beckham, Sina Honari, Vikas Verma, Alex M Lamb, Farnoosh Ghadiri, R Devon Hjelm, Yoshua Bengio, and Chris Pal. On adversarial mixup resynthesis. In *Advances in Neural Information Processing Systems*, pages 4348–4359, 2019.

A. Patrícia Bento, Anna Gaulton, Anne Hersey, Louisa J. Bellis, et al. The ChEMBL bioactivity database: an update. *Nucleic Acids Research*, 42(D1):D1083–D1090, 11 2013. ISSN 0305-1048. doi: 10.1093/nar/gkt1031. URL https://doi.org/10.1093/nar/gkt1031.

G Richard Bickerton, Gaia V Paolini, Jérémy Besnard, Sorel Muresan, and Andrew L Hopkins. Quantifying the chemical beauty of drugs. *Nature chemistry*, 4(2):90, 2012.

Esben Bjerrum and Boris Sattarov. Improving chemical autoencoder latent space and molecular de novo generation diversity with heteroencoders. *Biomolecules*, 8(4):131, 2018.

Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.

Hongming Chen, Ola Engkvist, Yinhai Wang, Marcus Olivecrona, and Thomas Blaschke. The rise of deep learning in drug discovery. *Drug discovery today*, 2018.

Yow-Jyun Cheng, Hang-Shiang Jiang, Shih-Lan Hsu, Li-Chiung Lin, Chieh-Liang Wu, Vithal K Ghanta, and Chi-Mei Hsueh. Xiap-mediated protection of h460 lung cancer cells against cisplatin. *European journal of pharmacology*, 627 (1-3):75–84, 2010.

Connor W Coley, Luke Rogers, William H Green, and Klavs F Jensen. Scscore: synthetic complexity learned from a reaction corpus. *Journal of chemical information and modeling*, 58(2):252–261, 2018.

Prasad Dandawate, Aamir Ahmad, Jyoti Deshpande, K Venkateswara Swamy, Ejazuddin M Khan, Madhukar Khetmalas, Subhash Padhye, and Fazlul Sarkar. Anticancer phytochemical analogs 37: synthesis, characterization, molecular docking and cytotoxicity of novel plumbagin hydrazones against breast cancer cells. *Bioorganic & medicinal chemistry letters*, 24(13):2900–2904, 2014.

Hans De Wolf, Laure Cougnaud, Kirsten Van Hoorde, An De Bondt, Joerg K Wegner, Hugo Ceulemans, and Hinrich Göhlmann. High-throughput gene expression profiles to define drug similarity and predict compound activity. *Assay and drug development technologies*, 16(3):162–176, 2018.

John S. Delaney. Esol: Estimating aqueous solubility directly from molecular structure. *Journal of Chemical Information and Computer Sciences*, 44(3):1000–1005, 2004. doi: 10.1021/ci034243x.

Joaquin Dopazo. Genomics and transcriptomics in drug discovery. *Drug discovery today*, 19(2):126–132, 2014.

David W Gammon, Daniel J Steenkamp, Vuyo Mavumengwana, Mohlopheni J Marakalala, Theophilus T Mudzunga, Roger Hunter, and Muganza Munyololo. Conjugates of plumbagin and phenyl-2-amino-1-thioglucoside inhibit mshb, a deacetylase involved in the biosynthesis of mycothiol. *Bioorganic & medicinal chemistry*, 18(7):2501–2514, 2010.

Anna Gaulton, Anne Hersey, Michał Nowotka, A Patrícia Bento, Jon Chambers, David Mendez, Prudence Mutowo, Francis Atkinson, Louisa J Bellis, Elena Cibrián-Uhalte, et al. The chembl database in 2017. *Nucleic acids research*, 45(D1):D945–D954, 2016.

Paul Geeleher, Nancy J Cox, and R Stephanie Huang. Cancer biomarker discovery is improved by accounting for variability in general levels of drug sensitivity in pre-clinical models. *Genome biology*, 17(1):190, 2016.

Mahmoud Ghandi, Franklin W Huang, Judit Jané-Valbuena, Gregory V Kryukov, Christopher C Lo, E Robert McDonald, Jordi Barretina, Ellen T Gelfand, Craig M Bielski, Haoxin Li, et al. Next-generation characterization of the cancer cell line encyclopedia. *Nature*, 569(7757):503, 2019.

Leonard Girsh. Lipid-containing compositions and methods of using them, February 15 2007. US Patent App. 11/501,380.

Rafael Gomez-Bombarelli, Jennifer N Wei, David Duvenaud, Jose Miguel Hernandez-Lobato, et al. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.

Magdalena Gorska, Alicja Kuban-Jankowska, Ryszard Milczarek, and Michal Wozniak. Nitro-oxidative stress is involved in anticancer activity of $17\beta$-estradiol derivative in neuroblastoma cells. *Anticancer research*, 36:1693–8, 04 2016.

Daria A Grechishnikova. Transformer neural network for protein specific de novo drug generation as machine translation problem. *BioRxiv*, page 863415, 2019.

Florian Häse, Loïc M Roch, and Alán Aspuru-Guzik. Chimera: enabling hierarchy based multi-objective optimization for self-driving laboratories. *Chemical science*, 9(39):7642–7655, 2018.

John E. Hopcroft and Jeffrey D. Ullman. *Formal Languages and Their Relation to Automata*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1969.

Stanisław Jastrzebski, Damian Leśniak, and Wojciech Marian Czarnecki. Learning to smile (s). *arXiv preprint arXiv:1602.06289*, 2016.

Armand Joulin and Tomas Mikolov. Inferring algorithmic patterns with stack-augmented recurrent nets. In *Advances in neural information processing systems*, pages 190–198, 2015.

Prabhjot Kaur, Tarun Garg, Goutam Rath, RSR Murthy, and Amit K Goyal. Surfactant-based drug delivery systems for treating drug-resistant lung cancer. *Drug delivery*, 23(3):717–728, 2016.

Anna Kawiak, Anna Domachowska, Anna Jaworska, and Ewa Lojkowska. Plumbagin sensitizes breast cancer cells to tamoxifen-induced cell death through grp78 inhibition and bik upregulation. *Scientific reports*, 7:43781, 2017.

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

Jo Klaveness, Pål Rongved, Anders Høgset, Helge Tolleshaug, Alan Cuthbertson, Aslak Godal, Lars Hoff, Geir Gogstad, Klaus Bryn, Anne Naevestad, et al. Diagnostic/therapeutic agents, January 20 2004. US Patent 6,680,047.

Juan Lao, Julia Madani, Teresa Puértolas, María Álvarez, Alba Hernández, Roberto Pazo-Cid, Ángel Artal, and Antonio Antón Torres. Liposomal doxorubicin in the treatment of breast cancer patients: a review. *Journal of drug delivery*, 2013, 2013.

Ann Lin, Christopher J Giuliano, Ann Palladino, Kristen M John, Connor Abramowicz, Monet Lou Yuan, Erin L Sausville, Devon A Lukow, Luwei Liu, Alexander R Chait, et al. Off-target toxicity is a common mechanism of action of cancer drugs undergoing clinical trials. *Science translational medicine*, 11(509):eaaw8412, 2019.

Matteo Manica, Ali Oskooei, Jannis Born, Vigneshwari Subramanian, Julio Saez-Rodriguez, and Maria Rodriguez Martinez. Toward explainable anticancer compound sensitivity prediction via multimodal attention-based convolutional encoders. *Molecular Pharmaceutics*, 2019. doi: 10.1021/acs.molpharmaceut.9b00520. PMID: 31618586.

Oscar Méndez-Lucio, Benoit Baillif, Djork-Arné Clevert, David Rouquié, and Joerg Wichard. De novo generation of hit-like molecules from gene expression signatures using artificial intelligence. *Nature Communications*, 11(1):1–10, 2020.

14

Ali Oskooei, Matteo Manica, Roland Mathis, and María Rodríguez Martínez. Network-based biased tree ensembles (netbite) for drug sensitivity prediction and drug sensitivity biomarker identification in cancer. *Scientific reports*, 9 (1):1–13, 2019.

Pavel G Polishchuk, Timur I Madzhidov, and Alexandre Varnek. Estimation of the size of drug-like chemical space based on gdb-17 data. *Journal of computer-aided molecular design*, 27(8):675–679, 2013.

Radhika Poojari. Embelin–a drug of antiquity: shifting the paradigm towards modern medicine. *Expert opinion on investigational drugs*, 23(3):427–444, 2014.

Mariya Popova, Olexandr Isayev, and Alexander Tropsha. Deep reinforcement learning for de novo drug design. *Science advances*, 4(7):eaap7885, 2018.

Daniel Probst and Jean-Louis Reymond. Visualization of very large high-dimensional data sets as minimum spanning trees. *arXiv preprint arXiv:1908.10410*, 2019.

David Rogers and Mathew Hahn. Extended-Connectivity Fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754, May 2010. ISSN 1549-9596. doi: 10.1021/ci100050t. URL https://doi.org/10.1021/ci100050t.

RXN. Ibm rxn for chemistry. https://rxn.res.ibm.com/, 2020. Accessed: 2020-02-02.

Ketan T Savjani, Anuradha K Gajjar, and Jignasa K Savjani. Drug solubility: importance and enhancement techniques. *ISRN pharmaceutics*, 2012, 2012.

Jack W Scannell, Alex Blanckley, Helen Boldon, and Brian Warrington. Diagnosing the decline in pharmaceutical r&d efficiency. *Nature reviews Drug discovery*, 11(3):191, 2012.

Gisbert Schneider. Mind and machine in drug design. *Nature Machine Intelligence*, page 1, 2019.

Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural computation*, 10(5):1299–1319, 1998.

Hossein Sharifi Noghabi, Shuman Peng, Olga Zolotareva, Colin C Collins, and Martin Ester. Aitl: Adversarial inductive transfer learning with input and output space adaptation for pharmacogenomics. *bioRxiv*, 2020. doi: 10.1101/2020.01.24.918953. URL https://www.biorxiv.org/content/early/2020/01/25/2020.01.24.918953.

Pavithra Shivakumar and Michael Krauthammer. Structural similarity assessment for drug sensitivity prediction in cancer. In *BMC bioinformatics*, volume 10, page S17. Springer, 2009.

Robert H Shoemaker. The nci60 human tumour cell line anticancer drug screen. *Nature Reviews Cancer*, 6(10):813, 2006.

Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. In *Advances in neural information processing systems*, pages 3483–3491, 2015.

Murat Cihan Sorkun, Abhishek Khetan, and Süleyman Er. Aqsoldb, a curated reference set of aqueous solubility and 2d descriptors for a diverse set of compounds. *Scientific data*, 6(1):1–8, 2019.

Taffee T Tanimoto. Elementary mathematical theory of classification and prediction. *IBM Internal Report*, 1958.

Bie Verbist, Günter Klambauer, Liesbet Vervoort, Willem Talloen, Ziv Shkedy, Olivier Thas, Andreas Bender, Hinrich WH Göhlmann, Sepp Hochreiter, QSTAR Consortium, et al. Using transcriptomics to guide lead optimization in drug discovery projects: Lessons learned from the qstar project. *Drug discovery today*, 20(5):505–513, 2015.

Liwei Wang, Alexander Schwing, and Svetlana Lazebnik. Diverse and accurate image description using a variational auto-encoder with an additive gaussian encoding space. In *Advances in Neural Information Processing Systems*, pages 5756–5766, 2017.

Martin Wehling. Assessing the translatability of drug projects: What needs to be scored to predict success? *Nature reviews. Drug discovery*, 8:541–6, 07 2009. doi: 10.1038/nrd2898.

David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.

John N Weinstein, Eric A Collisson, Gordon B Mills, Kenna R Mills Shaw, Brad A Ozenberger, Kyle Ellrott, Ilya Shmulevich, Chris Sander, Joshua M Stuart, Cancer Genome Atlas Research Network, et al. The cancer genome atlas pan-cancer analysis project. *Nature genetics*, 45(10):1113, 2013.

Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.

Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989.

Chi Heem Wong, Kien Wei Siah, and Andrew W Lo. Estimation of clinical trial success rates and related parameters. *Biostatistics*, 20(2):273–286, 2019.

Carl Yang, Peiye Zhuang, Wenhan Shi, Alan Luu, and Pan Li. Conditional structure generation through graph variational generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 1338–1349, 2019.

Wanjuan Yang, Jorge Soares, Patricia Greninger, Edelman, et al. Genomics of drug sensitivity in cancer (gdsc): a resource for therapeutic biomarker discovery in cancer cells. *Nucleic acids research*, 41(D1):D955–D961, 2012.

Jiaxuan You, Bowen Liu, Zhitao Ying, Vijay Pande, and Jure Leskovec. Graph convolutional policy network for goal-directed molecular graph generation. pages 6412–6422, 2018.

Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. *arXiv preprint arXiv:2001.06782*, 2020.

Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan R Salakhutdinov, and Alexander J Smola. Deep sets. In *Advances in neural information processing systems*, pages 3391–3401, 2017.

Xiaotun Zhang, Wenbin Wang, Lawrence D. True, Robert L. Vessella, and Thomas K. Takayama. Protease-activated receptor-1 is upregulated in reactive stroma of primary prostate cancer and bone metastasis. *The Prostate*, 69(7): 727–736, 2009. doi: 10.1002/pros.20920.

XQ Zhang, CY Yang, XF Rao, and JP Xiong. Plumbagin shows anti-cancer activity in human breast cancer cells by the upregulation of p53 and p21 and suppression of g1 cell cycle regulators. *European journal of gynaecological oncology*, 37(1):30–35, 2016.

Alex Zhavoronkov, Yan A Ivanenkov, Alex Aliper, Mark S Veselov, Vladimir A Aladinskiy, Anastasiya V Aladinskaya, Victor A Terentiev, Daniil A Polykovskiy, Maksim D Kuznetsov, Arip Asadulaev, et al. Deep learning enables rapid identification of potent ddr1 kinase inhibitors. *Nature biotechnology*, 37(9):1038–1040, 2019.

# PaccMann^RL: Designing anticancer drugs from transcriptomic data via reinforcement learning - Appendix

## S1 SMILES VAE architecture with StackGRU

To enable neural networks to count, Joulin and Mikolov (2015) introduced stack-augmented RNN. Stack-RNNs complement RNNs with a differentiable push-down stack operated through learnable controllers, $op_t$ at step $t$, that involve three operations: PUSH, POP and NO-OP (see Figure S1).

$$op_t = \mathbf{s}(W_{op}h_t), \tag{2}$$

where $h_t$ is the hidden state, $W_{op}$ is a $3 \times H$ matrix (H being the dimension of hidden state) and $\mathbf{s}$ is the softmax function. At each time step the controller probabilities are determined from Equation 2 and the stack memory is updated using the learned controller via a multiplicative gating mechanism:

$$\begin{cases} S_t[0] & = op_t[\text{PUSH}]\mathbf{s}(W_{so}h_t) + op_t[\text{POP}]S_{t-1}[1] + \\ & \quad op_t[\text{NO-OP}]S_{t-1}[0] \\ h_t & = \mathbf{s}(W_i X_t + W_R h_{t-1} + W_{si} S_{t-1}) \end{cases} \tag{3}$$

where $S_t$ is the stack, $W_{so}$ is a $1 \times H$ matrix and $W_{si}$ is a $H \times N$ matrix ($N$ being the stack height). $W_i$ is the input matrix applied to the sequence and $W_R$ is the recurrent matrix. It should be noted that for the sake of brevity, we only show the update equation for the topmost element of the stack in Equation 3.
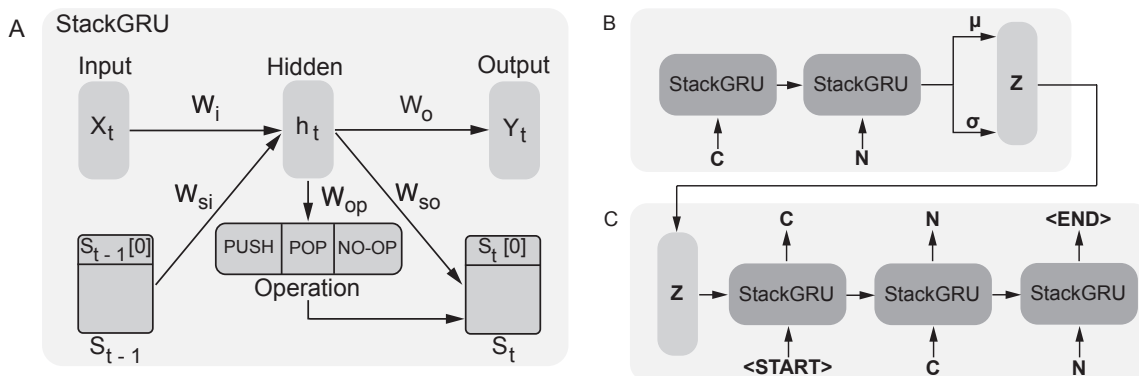


Figure S1: (A) The StackGRU architecture adopted in the SVAE. The stack-augmented GRU (StackGRU) architecture complements a regular GRU with a stack that allows one out of three possible operations at each time-step: PUSH, POP and NO-OP. The operation vector is determined through a softmax from the hidden state of each time step.
(B) and (C) are encoder and decoder of the SVAE architecture. (B) encodes the SMILES sequences into multivariate Gaussians with parameters $\mu$ and $\sigma$. (C) The decoder StackGRU units reconstruct the SMILES sequence from a latent representation ($Z_p$) sampled from the multivariate Gaussian.

## S2 Reward function

A reward function $f$ was used to map the logarithmic micromolar IC50 values predicted by the critic to a reward that was subject to maximization in our adopted RL framework (see Figure S2).. It is computed by $f(IC50) = \exp\left(\frac{-IC50}{\alpha}\right)$, where $IC50 = \mathcal{C}(C_T, X_c)$ and $\alpha \in \mathbb{R}^+$ is a tunable hyperparameter that determines how much the generator is rewarded for designing effective versus ineffective compounds, i.e., a smaller $\alpha$ leads to a greedier generator.
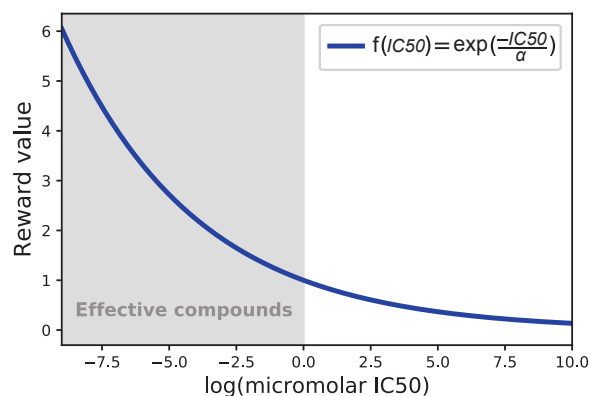
Figure S2: Reward function to map the predicted IC50 of the critic (PaccMann) to a reward being fed to the conditional generator. To produce the plot, $\alpha$ was set to 5.

## S3  Gene expression in human samples and cancer cell lines

Comparing the standardized gene expression values of GDSC (Yang et al., 2012) and CCLE (Barretina et al., 2012) with the one from human samples from TCGA (Weinstein et al., 2013) reveals a similarity (Figure S3). This justifies our choice of utilizing the encoder of the PVAE for cell line data during the RL regime, although it was initially pretrained on human samples from TCGA.



Figure S3: Distribution of standardized gene expression values across the cancer cell line databases CCLE and GDSC as well as the human sample database TCGA.

## S4  Implementation and training details

All models were implemented in `PyTorch` 1.0 and trained on a cluster equipped with `POWER8` processors and a `NVIDIA Tesla P100`.

**PVAE.**  The model consisted of four dense layers of [1024, 512, 256 and 200] units with ReLU activation function and dropout of $p = 0.2$ in both, the encoder and the decoder. The dimensionality of the latent space ($n$) was 128.

We minimized the variational loss, consisting of the reconstruction loss and KL divergence, using Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 1e-8$) and a decreasing learning rate starting at 0.001 Kingma and Ba (2014). To further regularize the PVAE, denoising methods were employed by 1) applying a dropout of 0.1 on the input genes and 2) adding noise to gene expression values ($\varepsilon \sim \mathcal{N}(0, 0.1)$). The model was trained with a batch size of 64 for a maximum of 2000 epochs.

**SVAE.** The model was trained on molecules provided in SMILES notation, the longest molecules had 1423 tokens. Both encoder and decoder consisted of two layers of bidirectional GRU (hidden size of 128, dropout of 0.1 at the first layer), each complemented with 50 parallel memory stacks with the depth of 50. The latent space of SVAE had the same dimensionality as the PVAE (128) to enable the addition of encodings. Similar optimization parameters as PVAE were used. This model further utilized teacher forcing Williams and Zipser (1989), i.e., the model's output is conditioned on the previous ground truth sample as opposed to its generated output. Whilst this significantly simplifies learning, it may drive the generator to predominantly rely on the decoder (thus neglecting the latent encoding). This so called posterior collapse was resolved by applying a token dropout rate of 0.1 during teacher forcing as suggested by Bowman et al. (2015). In addition to token dropout, KL cost-annealing Bowman et al. (2015) was employed during training, The model was trained with a batch size of 128 for a maximum (early stopping) of $\sim$ 110,000 steps (i.e., exactly 10 epochs) During training, KL cost-annealing as described in Bowman et al. (2015) was explored in order to trade-off reconstruction and KL loss.

**Critic.** The critic was trained using the parameters reported in Manica et al. (2019) and replicating the best performing architecture based on multiscale convolutional encoders.

**RL training.** In order to maximize Equation 1, we employed Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 1e-4$, weight decay 1e$-$4) and a decreasing learning rate starting at 1e$-$5. The gradients were clipped to 2 to prevent $\mathcal{G}$ from destroying its chemical knowledge about SMILES syntax obtained through pretraining on ChEMBL. The reward function hyperparameter $\alpha$ was set to 5.

## S5 Results for gene expression profile VAE and SMILES VAE

**Profile VAE (PVAE)** The pretraining results of the PVAE are presented in Figure S4A, B and C. As shown in Figure S4B, the reconstructed gene expression profiles (GEP), shown in blue, as well as the generated GEPs (green) accurately mimic the distribution of the original GEPs (red). Furthermore, the sampled GEPs follow the same lognormal distribution as the original data. Figure S4C shows that the generated GEPs exhibit a higher similarity to the testing than to the training sample. Overall, these results suggest that the PVAE learns to embed GEPs meaningfully into a latent space that allows both reconstruction and generation of new realistic GEPs of human cells.

**SMILES VAE (SVAE)** Figure S4D, E and F give a quantitative analysis of the SVAE results following pretraining for 10 epochs with $\sim$1.4 million structures from ChEMBL. To investigate the novelty and diversity of the generated molecules, we sampled 10,000 molecules by decoding random points from the latent space. Overall, 96.2% of the 10,000 generated molecules were valid molecular structures (assessed via `RDKit`) surpassing the results reported by Popova et al. (2018) who used the same stack-augmented GRUs trained on the ChEMBL database (95% SMILES validity). In addition, 99.72% of the valid generated molecules were unique across the 10,000 generations. The kernel density estimate (KDE) of the dimensions of the latent space (Figure S4E) validates that the SVAE fulfills the variational constraint as imposed by the Kullback-Leibler divergence in its loss function. We then utilized a well-established chemical structure similarity measure, the Tanimoto similarity (Tanimoto, 1958) to compare the ECFP (Rogers and Hahn, 2010) of a subset of 1000 generated molecules with the training and test data from ChEMBL. Figure S4F presents the distributions of the highest Tanimoto similarity between each generated compound and all compounds in training and test dataset respectively. Only a negligible fraction of the generated molecules existed in either of the datasets, whereas the vast majority had a Tanimoto similarity ($\tau$) between 0.2 and 0.6 suggesting that our model learned to propose novel molecular structures from the chemical space of about $10^{30}$ to $10^{60}$ molecules (Polishchuk et al., 2013). In addition, a visualization of the chemical space of ChEMBL as well as generated compounds through the TMAP algorithm (a library that visualizes high dimensional data through minimum spanning trees (Probst and Reymond, 2019)) showed that the generated molecules mix well with the training molecules into the chemical space. A snapshot of the interactive visualization is shown in Figure S6.

Figure S5A, showcases a panel of 12 generated molecules for qualitative assessment of the molecular structures. The generated molecules generally share drug-like structural features. To inspect the smoothness of the latent space of molecules, we encoded a reference molecule shown at the top of Figure S5B into the latent space and decoded four points in the vicinity of the reference molecule leading to the generation of structurally similar yet different compounds.
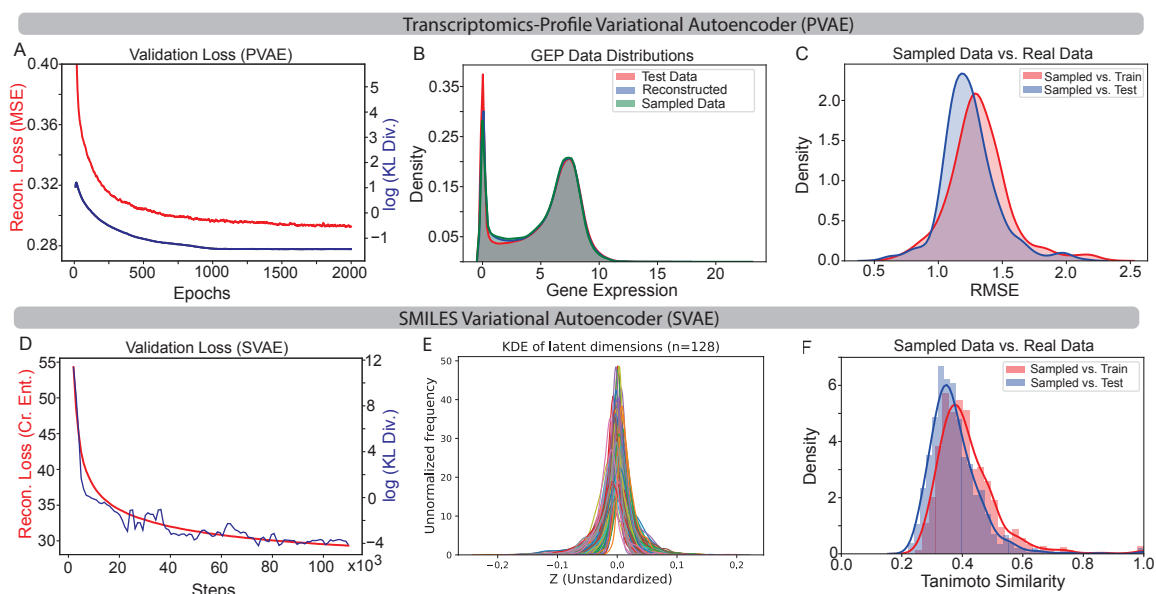
Figure S4: Results of pretrained PVAE and SVAE models. PVAE: (A) Development of validation error over the course of training. Reconstruction loss (MSE) and KL divergence are shown separately for comparison. (B) Distribution of gene expression values in real, reconstructed and generated samples. (C) Sampled (i.e. generated) data from the latent space of PVAE compared against training and test datasets from TCGA. SVAE: (D) Development of validation error over the course of training. Cross-entropy between target and generated SMILES is shown separately from the KL divergence (log scale for visual clarity). One epoch corresponds to ~11 000 training steps. (E) Kernel density estimates of all 128 latent dimensions before decoding the test samples. As enforced by the variational constraint, the latent variables follow Gaussian distributions. (F) The Tanimoto similarity between the Morgan fingerprints (ECFP) of the generated molecules and the structures from ChEMBL train and test datasets is used to verify that the generated compounds are sufficiently different from the training data.

Figure S5: Qualitative inspection of generated molecules. (A) A sample of 12 molecular structures produced with the SVAE. (B) The molecule depicted at the top was encoded into the latent space. The four molecules below show different decodings from the latent space in the vicinity of the starting molecule.
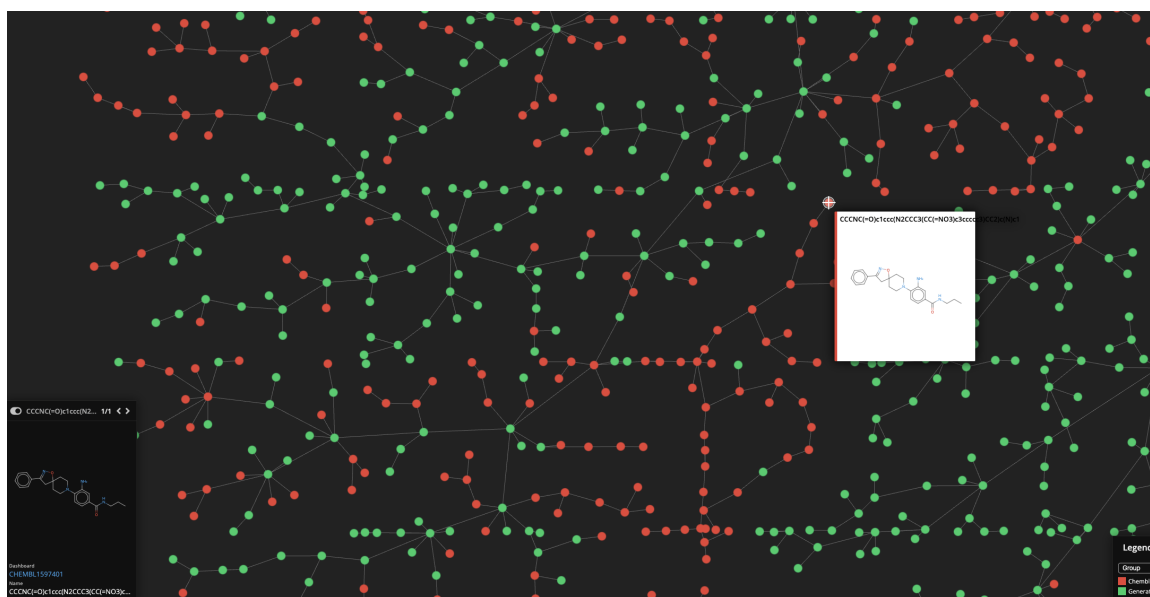


Figure S6: **Snapshot of the interactive TMAP visualization.** Generated molecules (green) and ChEMBL compounds (red) are shown through the TMAP algorithm which visualizes the chemical space by aggregating molecules with similar fingerprints (ECFP). The similarity in fingerprints is proportional to the distance of the repsective nodes on the spanning tree. To explore the visualization interactively, please visit https://paccmann.github.io/rl/unbiased.

## S6 Nearest neighbors in ChEMBL

To further validate the four site-specific compounds as proposed by our model, Figure S7 depicts the respective nearest neighbors (measured by Tanimoto similarity) in the ChEMBL database of bioactive compounds.
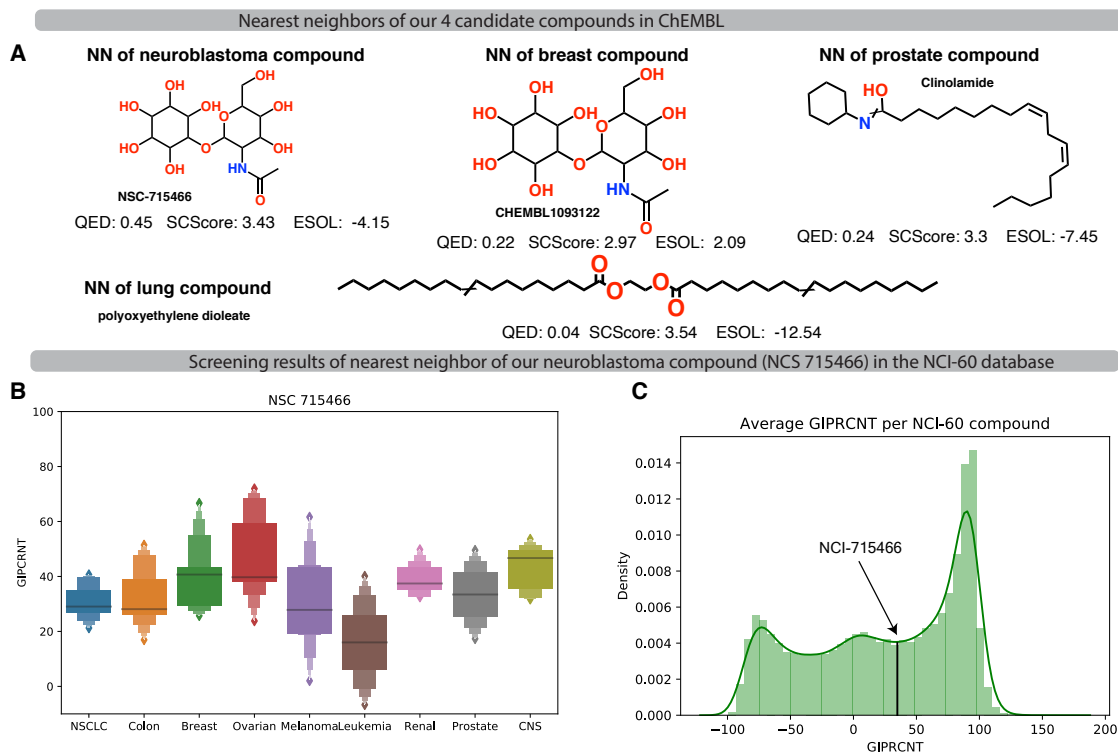


Figure S7: **A depiction of the nearest neighbors of our site-specific compounds in the ChEMBL database.**
**A)** For all four site-specific generated compounds, we performed a similarity search across all molecules in the ChEMBL database. The nearest neighbors are depicted together with relevant drug-like properties. **B)** NSC-715466 was tested as part of the NCI-60 database (Shoemaker, 2006), where it showed the strongest cell growth inhibition effect against leukemia cell lines. The GIPRCNT is a metric to measure cytotoxicity, where 100% refers to unchanged cell proliferation (identical to the control cells), 0% to complete stopping of cell proliferation and -100% to a full inhibition of all cells. **C)** NSC-715466 showed only moderate anticancer effects as reported in the NCI-60 database.

## S7 Validation of critic (PaccMann) on ChEMBL data

Our critic, PaccMann, is an anticancer drug sensitivity prediction model that has been trained *only* on anticancer compounds from GDSC. Since IC50 cell screening data for compounds with knowingly no anticancer effects are notoriously unavailable, PaccMann lacks a *negative training set* which would help extending its generalization across the space of known anticancer compounds. One could thus suspect that PaccMann is generally a flawed evaluator of compounds falling outside the space of anticancer drugs and that it would be biased towards predicting high efficacy for compounds without anticancer effects.

For that reason, Figure S8 shows the predicted efficacy of all anticancer compounds from GDSC (both training and testing data) as well as the predicted efficacy of a representative set of 1000 molecules from ChEMBL. The predicted logarithmic IC50 across all cancer drugs was $2.2 \pm 2.2$ whereas it was only $3.2 \pm 1.6$ for ChEMBL molecules.
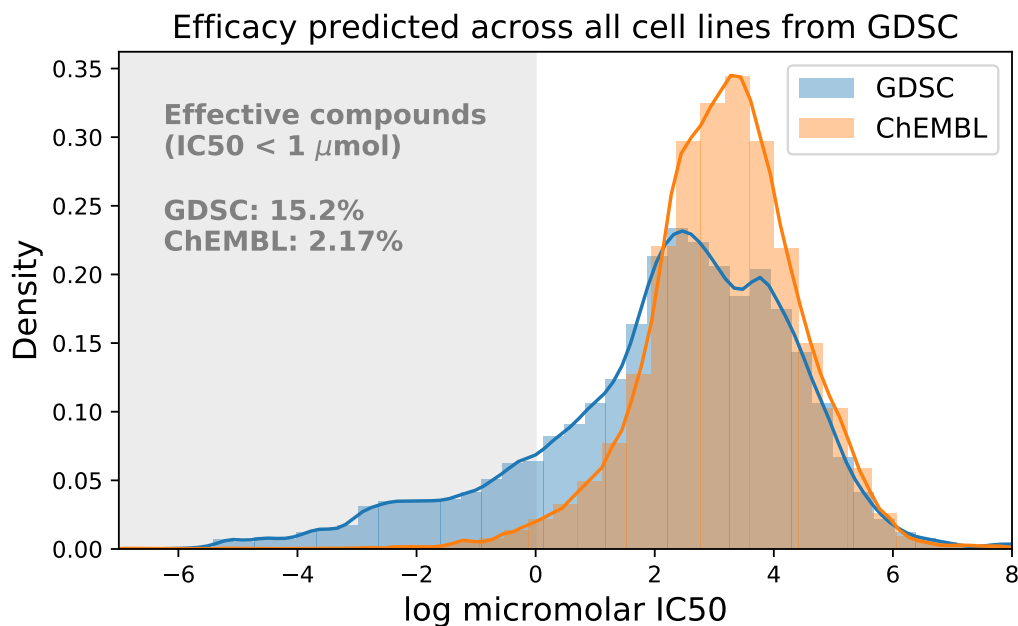
Figure S8: **A comparison of the predicted IC50 of GDSC drugs (known anticancer effects) and bioactive, drug-like compounds from ChEMBL.** The density plot of IC50 values as predicted by PaccMann from 210 cancer drugs from GDSC (blue) and 1000 molecules from ChEMBL (orange) across all 965 cell lines from the GDSC panel shows that only a minimal portion (2.17%) of ChEMBL compounds are predicted as effective against a given cell line, whereas this holds for a significantly larger fraction of GDSC compounds.

## S8    Possible synthesis routes for generated molecules

In order to assess the complexity of a potential synthesis of compound proposed by our model Figure S9 shows a predicted synthesis route for a compound gen rated against nervous system cancer. While the model proposes a panel of possible synthesis routes, the one depicted in Figure S9 decomposes the synthesis into four reactions and a total of ten commercially available reactants. The simplicity of the proposed route as well as the high confidence scores ($> 0.8$) associated to each reaction seem to be promising indicators towards a possible synthesis. Details of the synthesis route are depicted at the end of document.
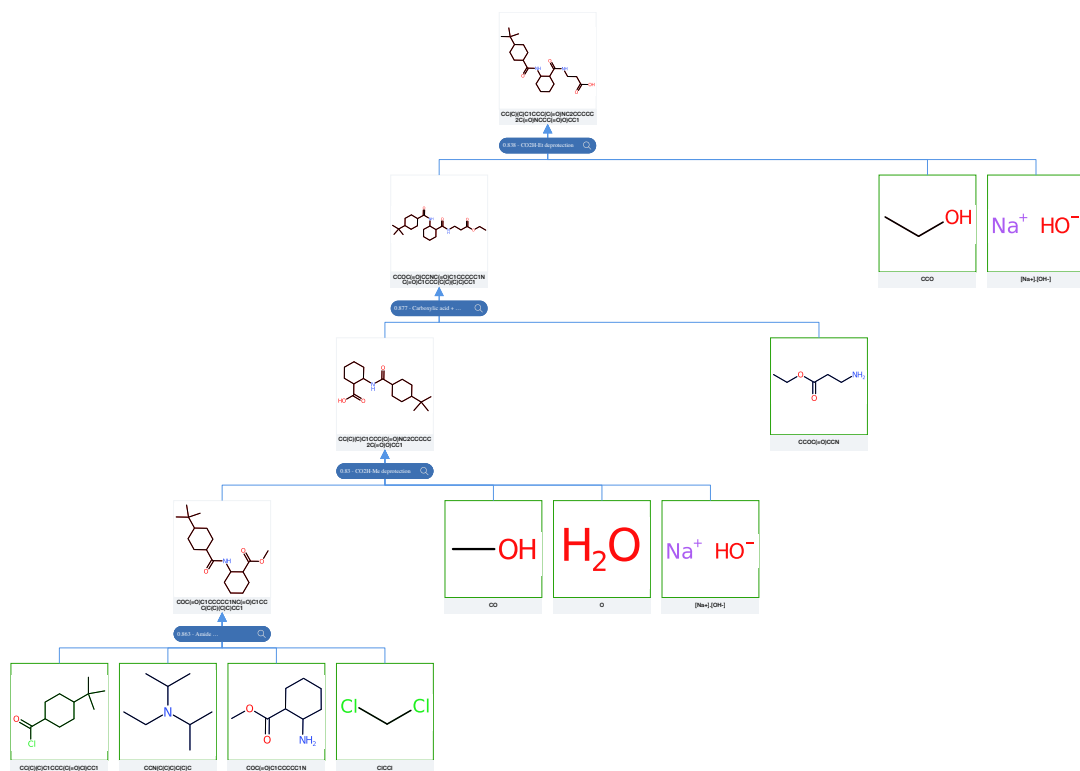
Figure S9: **A retrosynthesis route for a generated molecule.** A possible synthesis route, predicted by a molecular retrosynthesis model (RXN, 2020) is shown for a compound proposed against nervous system cancer (top middle). The predicted synthesis consists of four sequential reactions with a total 10 commercially available reactants (green).

# Information about the retrosynthesis

Created On: 2020-02-04T12:27:22.365000
Model: MolecularTransformer_v2.0_R-Inchi-MolecularTransformer_v2.0_F
Product: CC(C1CCC(C(NC2C(C(NCCC(O)=O)=O)CCCC2)=O)CC1)(C)C
MSSR: 15
FAP: 0.65
MRP: 50
SbP: 3
Available smiles:
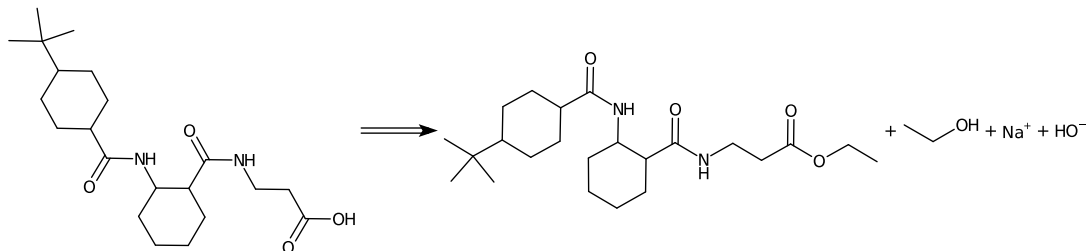Exclude smiles: CC(C1CCC(C(NC2C(C(NCCC(O)=O)=O)CCCC2)=O)CC1)(C)C
Exclude substructures:

# Sequence 0, Confidence: 0.527

## Step 1

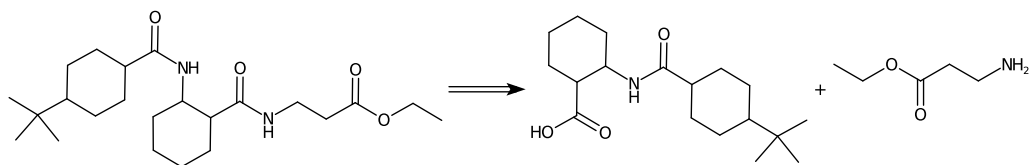*Type: CO2H-Et deprotection, Confidence: 0.838*

*CCOC(=O)CCNC(=O)C1CCCCC1NC(=O)C1CCC(C(C)(C)C)CC1.CCO.[Na+].[OH-]>>CC(C)(C)C1CCC(C(=O)NC2CCCCC2C(=O)NCCC(=O)O)CC1*



## Step 2

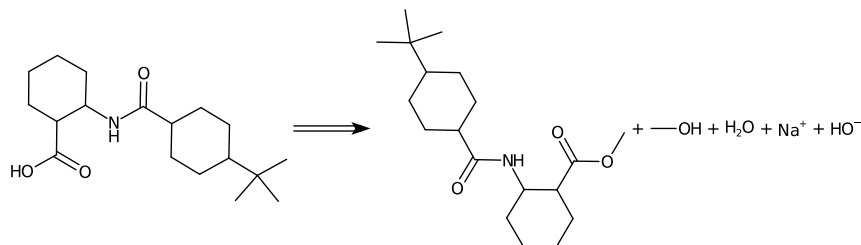*Type: Carboxylic acid + amine condensation, Confidence: 0.877*

*CC(C)(C)C1CCC(C(=O)NC2CCCCC2C(=O)O)CC1.CCOC(=O)CCN>>CCOC(=O)CCNC(=O)C1CCCCC1NC(=O)C1CCC(C(C)(C)C)CC1*



## Step 3

*Type: CO2H-Me deprotection, Confidence: 0.83*

*COC(=O)C1CCCCC1NC(=O)C1CCC(C(C)(C)C)CC1.CO.O.[Na+].[OH-]>>CC(C)(C)C1CCC(C(=O)NC2CCCCC2C(=O)O)CC1*



## Step 4

*Type: Amide Schotten-Baumann, Confidence: 0.863*

*CC(C)(C)C1CCC(C(=O)Cl)CC1.CCN(C(C)C)C(C)C.COC(=O)C1CCCCC1N.ClCCl>>COC(=O)C1CCCCC1NC(=O)C1CCC(C(C)(C)C)CC1*