



OPEN

Feasibility and application of machine learning enabled fast screening of poly-beta-amino-esters for cartilage therapies

Stefano Perni & Polina Prokopovich✉

Despite the large prevalence of diseases affecting cartilage (e.g. knee osteoarthritis affecting 16% of population globally), no curative treatments are available because of the limited capacity of drugs to localise in such tissue caused by low vascularisation and electrostatic repulsion. While an effective delivery system is sought, the only option is using high drug doses that can lead to systemic side effects. We introduced poly-beta-amino-esters (PBAEs) to effectively deliver drugs into cartilage tissues. PBAEs are copolymer of amines and di-acrylates further end-capped with other amine; therefore encompassing a very large research space for the identification of optimal candidates. In order to accelerate the screening of all possible PBAEs, the results of a small pool of polymers ($n=90$) were used to train a variety of machine learning (ML) methods using only polymers properties available in public libraries or estimated from the chemical structure. Bagged multivariate adaptive regression splines (MARS) returned the best predictive performance and was used on the remaining ($n=3915$) possible PBAEs resulting in the recognition of pivotal features; a further round of screening was carried out on PBAEs ($n=150$) with small variations of structure of the main candidates from the first round. The refinements of such characteristics enabled the identification of a leading candidate predicted to improve drug uptake > 20 folds over conventional clinical treatment; this uptake improvement was also experimentally confirmed. This work highlights the potential of ML to accelerate biomaterials development by efficiently extracting information from a limited experimental dataset thus allowing patients to benefit earlier from a new technology and at a lower price. Such roadmap could also be applied for other drug/materials development where optimisation would normally be approached through combinatorial chemistry.

Biomaterials and drug design are regarded as a very resource (physical, economical and time) intensive operations¹; the process can be constructed into sequential stages (discovery, preclinical, clinical and pharmacovigilance) named Phase0 to Phase4. During Phase0, traditional bench experiments are carried out to identify optimal candidates that are screened through further developmental stages; while further clinical trials progressively assess toxicity, efficacy and long term safety (Phase1 to Phase4)². The overall development process can take from a minimum of 5 up to 15 years with an estimated total development cost per approved drug of \$2168 million in 2018³. However, the actual costs are generally a commercial confidential information and, therefore, such estimates may not fully capture the complete investments required⁴. The try-and-error approach to molecule development, particularly during the initial design and make phases of the design-make-test-analyse (DMTA) discovery cycle, is often directed by human intuition, which is inherently biased and limited in knowledge, thus slowing drug development⁵. In such contest, the ability of data-driven in-silico prediction tools to model outcomes without the need to physically prepare candidates and run experiments would enable a fast throughput screening of candidate molecules and thus reducing both the time and monetary investments required to identify lead candidates⁶⁻⁹. This can be achieved by establishing correlations between certain properties of the molecules (inputs, also known as descriptors) and outcomes of interest using experimentally generated data on

School of Pharmacy and Pharmaceutical Sciences, Cardiff University, Redwood Building, King Edward VII Avenue, Cardiff CF10 3NB, UK. ✉email: prokopovichp@cf.ac.uk

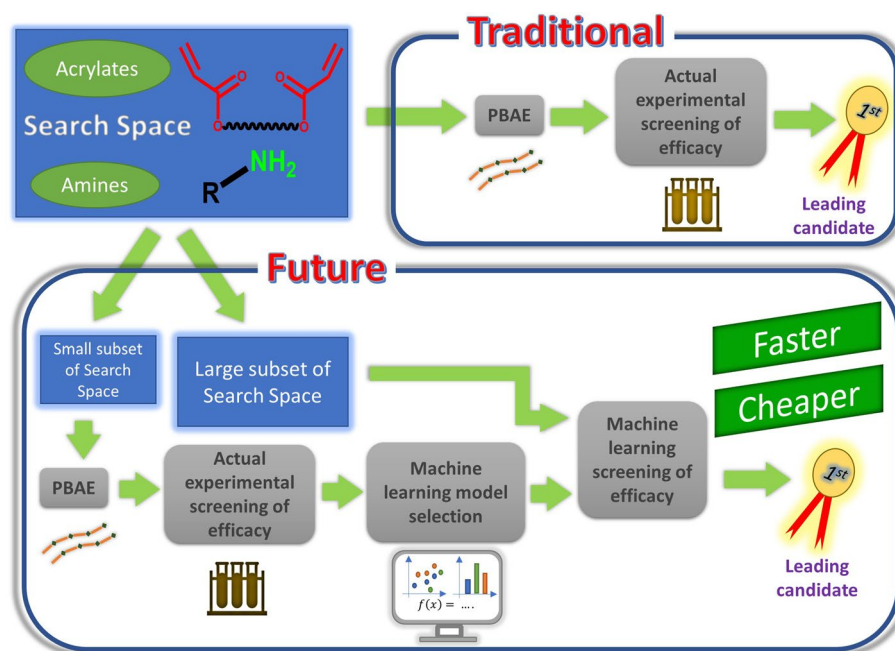


Figure 1. Schematic representation of machine learning driven drug development process.

a subset of relevant compounds; the established model would then be used to predict outcomes on the wider molecule search space¹⁰.

Machine learning (ML) based regression techniques are becoming wide spread in many areas of data analysis in the chemical^{11,12} and pharmaceutical sector^{13–16}; they have recently been employed in drug development^{17–19}, diagnostic²⁰, treatment algorithm optimisation²¹, drug repurposing²² and material discovery^{23,24}; however such applications are still quite limited despite being very promising^{25,26}. Another application of ML technologies in drug discovery is during compound screening or hit/lead generation and optimization enabling a virtual screening platform that offers a quicker and cheaper alternative to classic testing of large compounds libraries^{27,28}; virtual screening can be generally classified in ligand-based or structure-based²⁸. Compound optimisation using ML enabled virtual screening has been successfully applied to drug development for Alzheimer's disease²⁹, Class B G protein-coupled receptors (GPCRs)³⁰ and antiviral³¹. Figure 1 depicts how ML could be deployed to accelerate the biomaterial development process through virtual screening. Despite the flexibility of ML techniques, material design and optimisation involving numerous parameters are situations more likely to benefit from the development of machine learning predictive models.

Osteoarthritis (OA) is a thinning or loss of the cartilage layer covering the surfaces of joints reducing articular mobility, causing pain and inflammation. Although OA is not a life threatening disease, it has a great impact on the quality of life of patients and their ability to perform regular activities resulting in a great burden to society and health care providers. Worldwide, 303.1 million of people live with hip or knee osteoarthritis³²; furthermore, OA prevalence is expected to grow as consequence of the ageing population and overnutrition (two critical risk factors for OA). An effective treatment is still missing, current therapies (anti-inflammatory and analgesics) are only managing symptoms. This lack of therapeutic options is compounded by the inability of delivering the active molecule where is needed because of the obstacles posed by the low vascularisation and high electrostatic repulsion of cartilage tissues; these factors limit the amount of drug effectively available to the targeted cells³³. In order to achieve drug localisation, without a delivery system, high concentrations of drugs are used in the synovial fluid as mass transfer is governed by concentration differences (Fick's law)^{34–36}. Such approach has some problematic drawbacks; firstly, it is a wasteful use of the drug as only a minimal amount is actually therapeutic, with consequences on treatment acquisition costs. Secondly, drug washout lead to systemic exposure with possible side effects, as in case of steroids³⁷.

Different drug delivery systems have been developed for the localisation of drugs in cartilage in the attempt to overcome such barriers; poly-beta-amino-ester (PBAEs)^{38,39} and avidin³⁴ are two examples of these delivery systems. While no particular optimisation of the delivery system based on avidin performance is feasible as this a well-defined protein; there are, instead, essentially ∞^2 possible PBAEs as these are copolymers of an amine and a di-acrylate⁴⁰. Moreover, when PBAEs end-capping is also considered, the possible combinations rise to ∞^3 . In light of the performance of PBAE as cartilage drug delivery system being extremely dependent on the polymer backbone; ML algorithms predicting the efficacy of the drug delivery in cartilage from the polymer's constituents' properties would provide a high throughput screening for the optimisation of the PBAE driven cartilage drug localisation technology, reducing the cost and time to select the most promising candidate. We have previously demonstrated how the uptake of dexamethasone (DEX) (a drug routinely administered in clinics through intra-articular injections to reduce OA symptoms) in cartilage tissue, through a poly-beta-amino-ester drug delivery system, could be modelled using partial least square regression^{38,39}. The inputs of this model are the physical

properties of the polymers and co-polymeric units (di-acrylate and amine) along with some experimentally obtained parameters such as the diffusion coefficient of the polymer through cartilage, the drug loading in the delivery system and the molecular weights (M_w and M_n) of the polymer chain³⁹. Through this previous work, we identified a polymer (current lead candidate obtained from screening the combination of 3 acrylates and 15 amines) that increased DEX uptake in cartilage about 8 times compared to the clinical formulation³⁹. Despite the ability of predicting uptake, this model, in order to make predictions on new candidates, still requires inputs generated by experiments (such as M_w , M_n and diffusion coefficient) thus not fully able to completely substitute lab-based work. With the purpose of accelerating the optimisation of the PBAE structure for the cartilage delivery system through a systematic screening of a large library of both acrylates and amines, we hypothesised that machine learning algorithms, utilising only predictors available in public libraries or calculated from the compound structure, namely the physico-chemical properties of the PBAE components, could be employed to fully predict the performance of the delivery system without the need for any experimentally originated data. Drug uptake data experimentally obtained from a subset of a large polymer library were utilised to train and optimise 25 machine learning models (e.g. Random Forests, Kth nearest neighbour (kNN), support vector machine (SVM), neural network and multivariate adaptive regression splines (MARS)) and investigated their predictive performance to identify the most accurate algorithm. This model was then employed to screen the PBAEs research space (round1) representing acrylates and amines with a wide range of structural features and moieties; key features in the amine and acrylate structure were recognised in the PBAEs predicted to return the greatest drug uptake, further elucidating correlations between PBAE structural properties and drug uptake. A further round of ML predictions (round2) was conducted to refine and improve efficacy, screening a new set of PBAE exhibiting structures with small variations of the core features of the main candidates identified in the first round. The most promising candidate identified at the end of round2 had a predicted 3 folds efficacy improvement over the previous best performing candidate (round1). Finally, the actual efficacy and safety of the predicted best candidate were also experimentally determined.

Results

Machine learning model selection. Amine 1 to 20, acrylates A to F and end-capping e-1 and e-2 were used to generate the library of PBAE-DEX used for the experimental determination of DEX uptake in cartilage; in total $15 \times 6 \times 2 = 180$ unique PBAE were synthesised, doubling the size of the experimentally tested PBAE. After random splitting, the train set included 70 PBAEs, while the remaining 20 PBAEs constituted the test set. As the ultimate purpose of modelling is being able to estimate outcomes (in our work the uptake of DEX in cartilage) on previously unseen predictors, a split of the initial dataset into train and test set was implemented to be able to identify the model with the greatest predicting ability that is not necessarily the one that return the most accurate fit of the data used to calculate the model parameters (i.e. regression coefficients)^{41,42}. For the same reason, data split in training and test set was stratified based on PBAEs thus experimental data of DEX uptake for different exposure duration and related to PBAE with different end-capping all belonged to one set only. The 25–75% split also is in the typical range to provide sufficient data points for both model parameters estimation (training set) and testing^{43–46}. Therefore, it was expected that all models performed better on the training set than on the test set (Fig. 2).

Bagged multivariate adaptive regression splines (bagged MARS) returned the lowest Root Mean Squared Error (RMSE) on the test set (0.072). Random Forest had the lowest RMSE on the training set (0.036) but the second lowest on the test set (0.073). Furthermore, regressions based on decision trees/random forests do not allow for extrapolation of the measured outcome beyond the training set and such would limit the possibility of identifying PBAE performing better the experimentally observed optimal candidate. Linear regression (forward, backward or stepwise) had the highest RMSE on the training datasets, 0.128, 0.128 and 0.124, respectively. The difference in model performance between train and test set depended on the algorithm used; for example elastic regulation had RMSE of 0.080 and 0.081 for train and test set respectively, while Bayesian additive regression trees returned RMSE of 0.043 and 0.149 on train and test set, respectively. The small difference between the RMSE on train and test set observed for the elastic regulation model is a consequence of the penalties assigned to predictors in the algorithm that reduce the risk of overfitting^{42,47}. Moreover, boosting and bagging improved model performance (Fig. 2), for example RMSE of bagged MARS was lower than MARS and random forests had lower RMSE than decision tree. This was expected as such approaches have been developed to improve on model performance^{42,47}. Bagging is the process of resampling from the same data set to generate numerous new datasets then used to fit the model, this bootstrapping reduces overfitting and model variability^{42,47}; on the other hand, boosting employs weak predictors to improve on the predictions of other predictors⁴⁸.

The optimisation of the bagged MARS model hyper-parameters showed that with increasing number of bagged samples, mean RMSE during cross-validation decreased; averaging 75 resamples gave the lowest RMSE (Fig. 3a) while the number of pruned parameters increased model performance monotonically, but RMSE marginally decreased with the combinations of more than 10 (Fig. 3b). Moreover, performance of bagged MARS improved when the degree of interaction between parameters increased from 1 to 2 (Fig. 3b). The optimal bagged MARS model was made of a combination of 75 MARS models with a median number of predictors of 9 and a median number of terms of 15.

DEX uptake predicted by the optimised bagged MARS model against the actual data for the test data set (Fig. 4a) revealed a general good agreement between prediction and actual data regardless of the PBAE end-capping agent while the residual distribution exhibited a gaussian distribution (Fig. 4b). Similar patterns were observed when the model was applied on the train set (Fig. 4c and d); however, the residuals were smaller resulting in a narrower distribution. Modelled uptake curves of DEX in cartilages with PBAE in the test set demonstrated a general good fit of the experimental data (Fig. 4e).

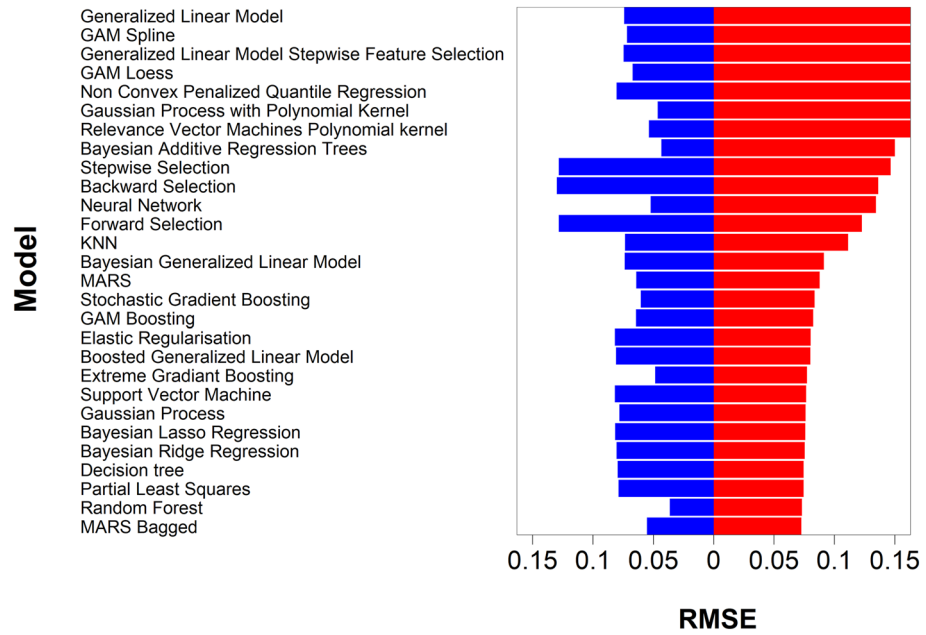


Figure 2. Comparison of the different performance of the tested algorithms on the train (blue) and test set (red).

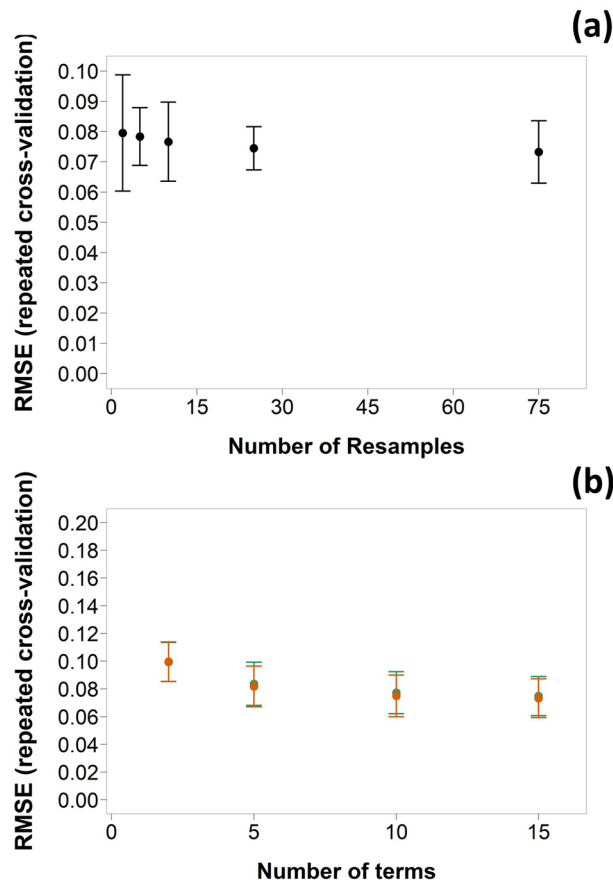


Figure 3. Relation for bag MARS models between RMSE (mean ± SD) for tenfold cross validation repeated 3 times and (a) number of resamples and (b) number of terms and degree of correlation ($n = 1$ ■, $n = 2$ ■) (number of resamples = 75).

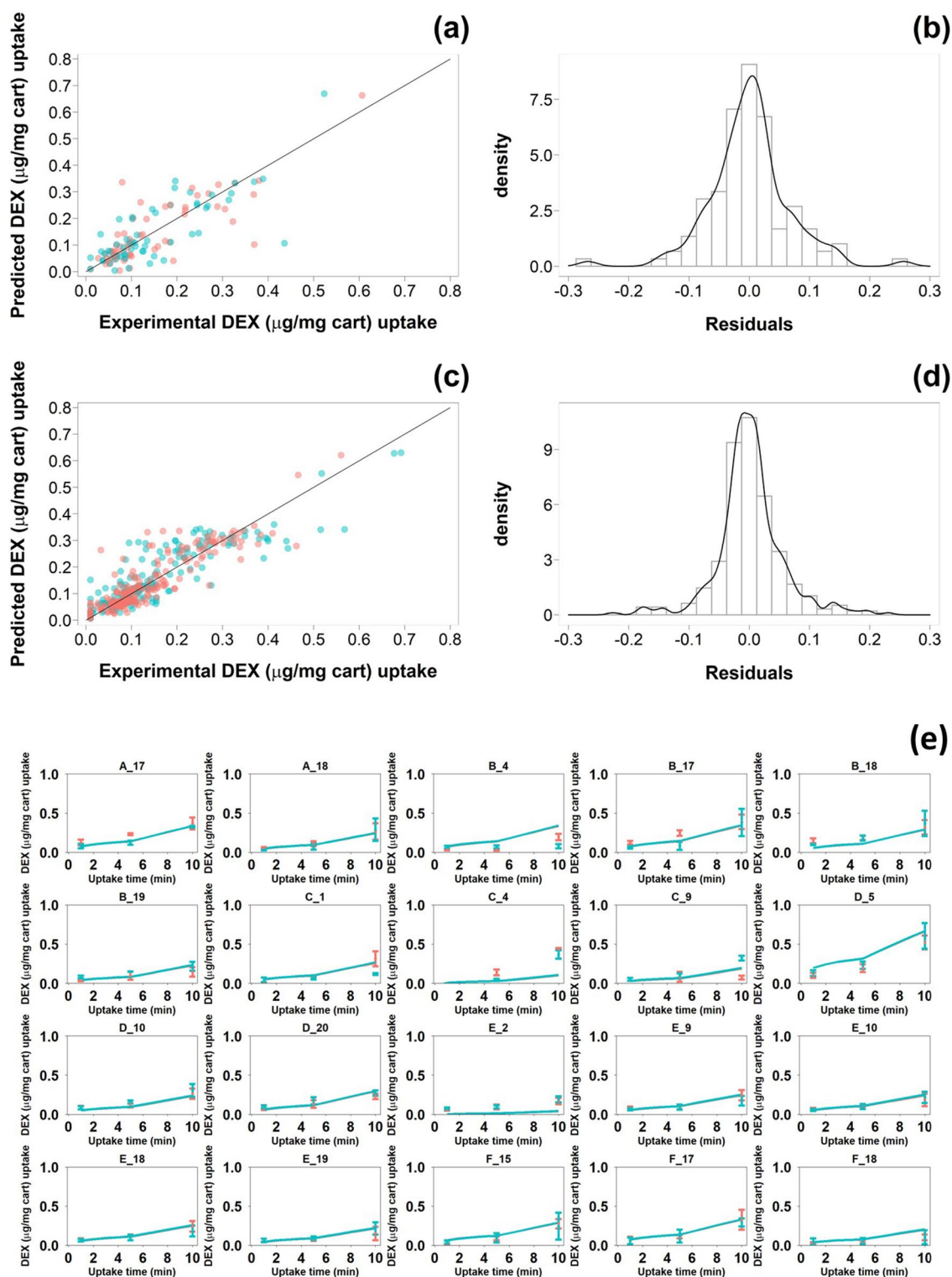


Figure 4. Comparison of predicted and experimental DEX uptake for PBAE-DEX (endcapped with e-1 ■ and e-2 ■) in the test (a) and train set (c); distribution of residuals of DEX uptake predictions for PBAE-DEX in the test (b) and train set (d). Comparison of time dependent DEX uptake (mean \pm SD) in cartilages predicted by optimised bag MARS model for PBAE in the test dataset (e).

The variable with the greatest importance in the bagged MARS model was ZStericQuad3D of the amine component, followed by the complexity of the amine component and the Henry's law coefficient of the PBAE repeated unit; the variable with the lowest importance was the molecular weight (MW) of the acrylate component

(Fig. 5a). In order to gain insights on the relations between the chemical and topological properties of the PBAE and the efficacy in localising DEX in cartilage, the specific dependence of the DEX uptake on the individual predictors was analysed on through the partial dependency plot (PDP).

These plots represent the predicted outcomes against a single varying input variable while maintaining the remaining constant at their mean values. PDP revealed ZStericQuad3D returned a maximum DEX uptake at ~ 0.83 ; while complexity of the amine decreased DEX uptake for values up to 50, for greater amine complexity predicted DEX uptake increased monotonically but was lower than the maximum (complexity = 0) for the maximum amine complexity in the library tested (Fig. 5b). As the models were trained on transformed values the relations between variables and drug uptake does not appear linear on the back-transformed predictions.

PBAE structure optimisation. The optimised bagged MARS model was applied on the remaining PBAE search space constituted by 3915 un-synthesised polymers to predict DEX uptake in cartilage after 10 min of exposure to PBAE-DEX when end-capped with e-1 or e-2. The results of this round1 screening identified 3192 PBAEs, regardless of the end-capping (end-capping agent e-1 returning predominantly higher drug uptake than e-2 on the same PBAE backbone), with an expected DEX uptake greater than the commercial formulation. Furthermore, 11 polymers with a predicted uptake greater than the previous leading candidate, which returned a drug uptake about 8 times that of DEX commercial formulation, were identified through the model. These PBAEs clustered very closely according to the dendrogram determined using the chemico-physical properties of the polymers and were made mainly by acrylate AAA (Phenylmethanediol diacrylate) or XX (1,4-Phenylene diacrylate) and amine 69 (2-Amino-5-(cyclopropyl)pyrazine) or 70 (2-Amino-6-propylpyrazine). PBAE candidate XX-69 was predicted to exhibit the greatest uptake among the full PBAE library tested, about 13 folds greater than the commercial formulation (Figs. 6 and S5).

1,4-phenylene diacrylate and (acrylate XX) and phenyl-methanediol diacrylate (acrylate (AAA) are the only acrylates tested exhibiting a benzene group where the electron of the oxygen atoms forming the di-acrylate groups can resonate reducing the impact of the electrostatic repulsion between some areas of the PBAE backbone and glycosaminoglycans (GAG) constituents of cartilage. Similarly, the presence of pyrazine in the amine constituent can increase the availability of the electron pair in the nitrogen resulting in higher positive charge. These two features were assumed to be key properties for effective drug delivery in cartilage and a refinement of the PBAE structure was carried out screening further acrylates ($n = 3$) exhibiting at least a benzene group in proximity of the acrylate moiety along with amines ($n = 50$) with a pyrazine in their structure (Fig. S6) in Round2. 17 of the 150 PBAEs tested in round2 had an estimated DEX uptake greater than XX-69 (best performer in round1); the presence of a further tertiary amine bound to the pyrazine ring resulted in greater DEX uptake in cartilage; moreover, two benzene groups (Bisacrylic acid oxybis(4,1-phenylene) ester) improved on the drug delivery (Figs. 7 and S7). The most effective PBAE (DDD-114) identified in round2 had a predicted DEX uptake about 21 time greater than the commercial formulation.

Ex-vivo performance of best candidate. DEX uptake in cartilage using DDD_114 increased with increasing exposure time; after 10 min the amount of drug retrieved from the samples using the PBAE based drug delivery system was over 20-folds the commercial DEX-P formulation confirming the model predictions (Fig. 8).

Mitochondrial activity of chondrocytes was not affected by the presence of the polymer (DDD_114_e1) ($p > 0.05$) (Fig. 9).

Discussion

The key to accurate predictions through mathematical models is the size of the data set used for the estimation of the model parameters⁴⁹. As our previous work hinted to the possibility of modelling cartilage drug uptake achieved by PBAEs conjugated to DEX³⁹, the machine learning models in this work were trained using a dataset³⁹ doubled in size with further polymers to reach a sufficient level of confidence in the model estimates. The work presented here considers only two end-capping agents treated as a categorical variable; the actual properties of the compounds were not considered as the number of molecules did not allow to capture such parameters.

Majority of research dedicated to implement ML in drug discovery/chemistry employs a very narrow range of potential models, even just one⁴⁹⁻⁵¹, without a clear rationale for the selection of the algorithms included in the pool assessed^{5,18,52-55}. Here instead, we purposely screened a large number of potential algorithms based on different approaches (e.g. decision tree, linear regression, SVM and neural network) in order to maximise the strength and transferability of the results while, simultaneously, increase the likelihood of identify a satisfactory predictive model.

MARS are an extension of linear models that can account for nonlinearities between input and output values through the use of hinge functions and interactions between variables combining flexibility and interpretability of results^{42,47}. The overall regression model “goodness of fit” depends on hyperparameters such as the number of pruned parameters and the degree of interaction between predictors. Bagged MARS is an ensemble of MARS constructed on a randomly generated bootstrapped set of data. Although it was expected that aggregating further resamples would improve model predictive performance, no more than 75 resamples were implemented in this work as the reduction in RMSE from 50 to 75 resamples was already minimal and a further increase of the resamples would also impact computational time.

The efficacy end-point experimentally assessed in this study was the amount of drug retrieved from a cartilage sample after contact with the PBAE based delivery system; this could not itself differentiate between actual drug penetration or accumulation on the external cartilage surface. However, previous evidence demonstrated that PBAE drug delivery systems diffuse inside the cartilage tissue underlying the validity of the approach^{38,56}. The

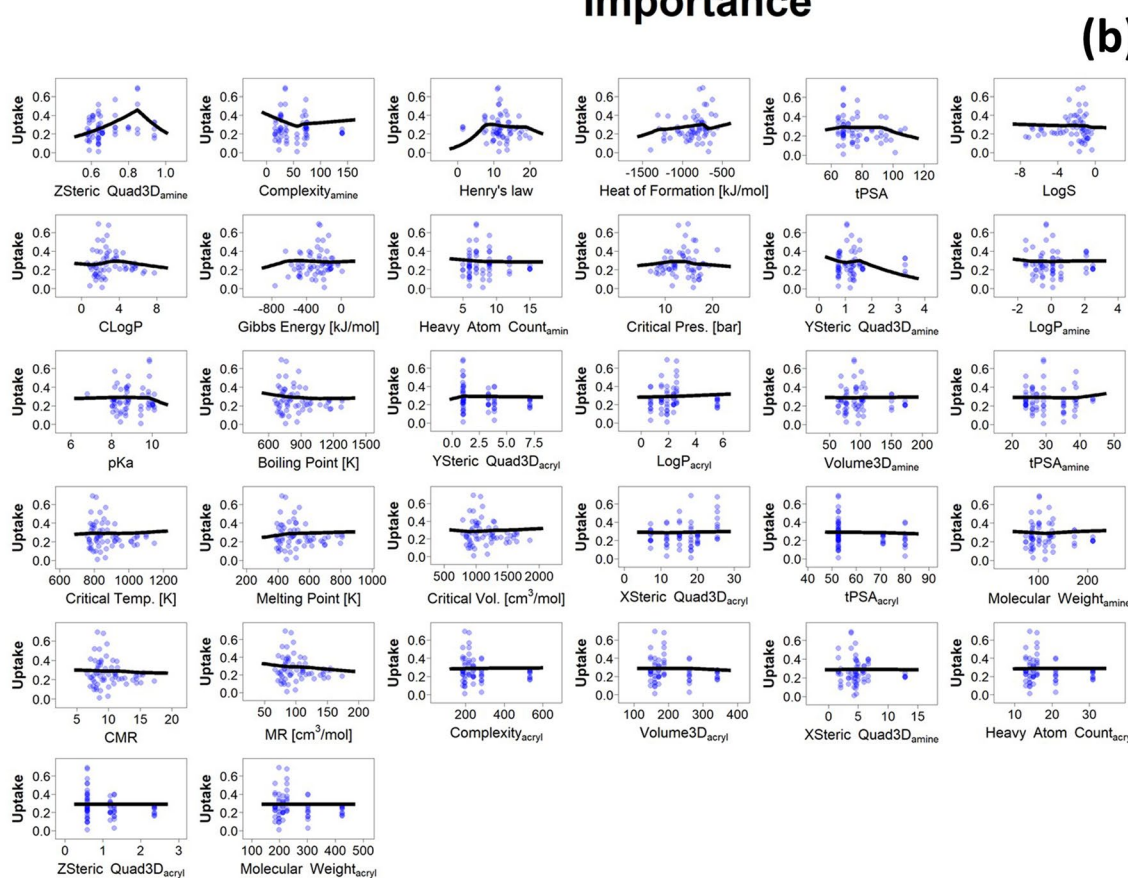
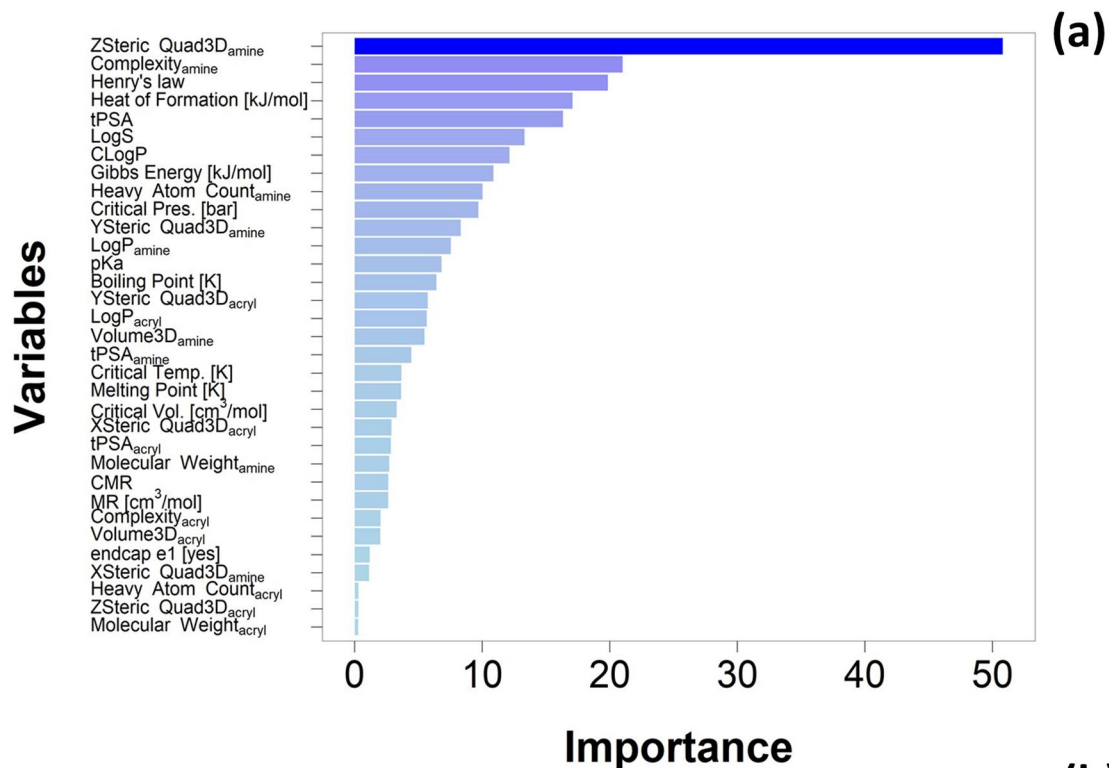


Figure 5. Variable importance in optimised bagged MARS model (a) and partial dependency plot of optimised bagged MARS model compared to experimentally obtained data for 10 min uptake of DEX into cartilage (b).

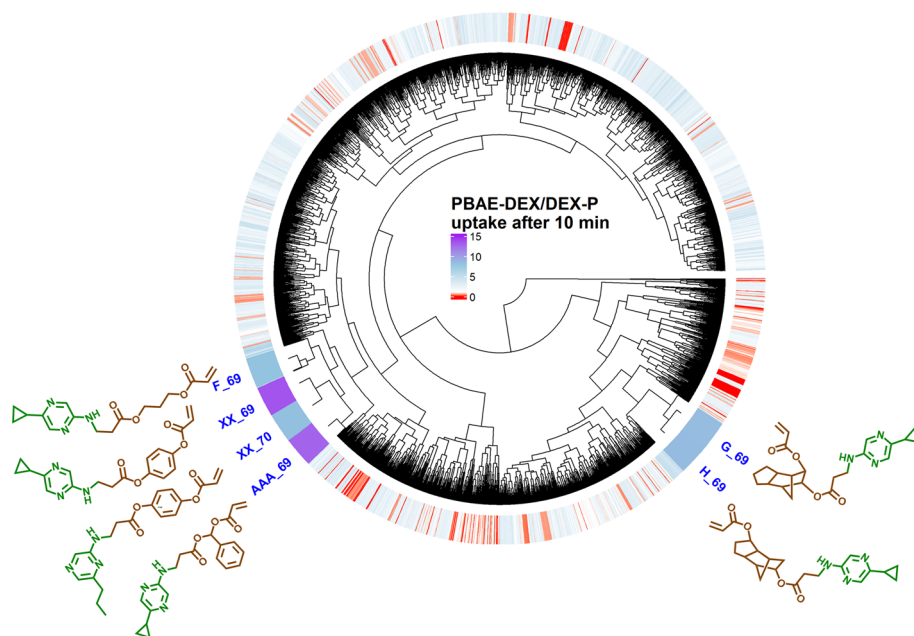


Figure 6. Heatmap of predicted ratio of DEX uptake for PBAE endcapped with e1 conjugated with DEX over commercial formulation of DEX after 10 min of exposure and structure of PBAE repeated unit with predicted drug uptake superior to experimental found candidate.

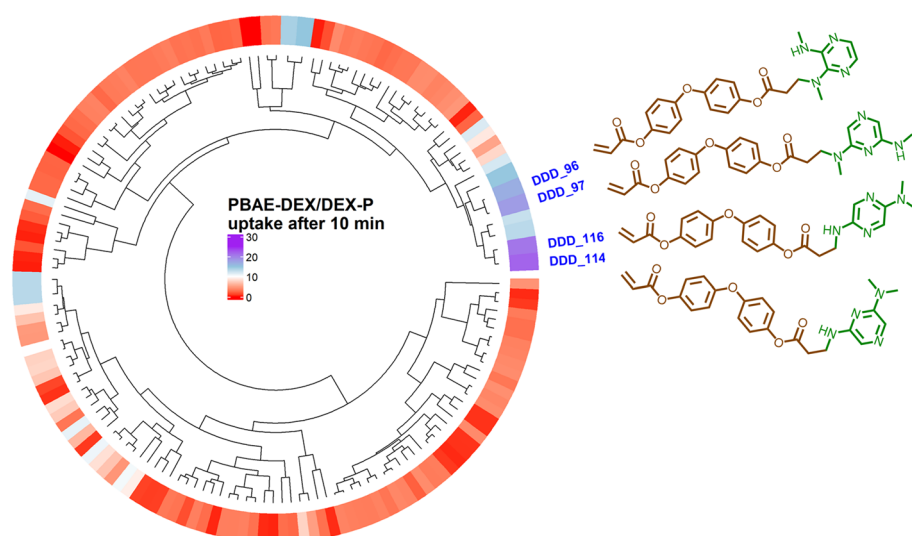


Figure 7. Heatmap of predicted ratio of DEX uptake for PBAE during round2 endcapped with e1 conjugated with DEX over commercial formulation of DEX after 10 min of exposure and structure of PBAE repeated unit with predicted drug uptake superior to best candidate in round1.

electrostatic interactions between positively charged PBAE and cartilage tissue components (predominantly the highly negatively charged GAGs) are the key mechanism of action of the delivery system under the presented investigation. The ranking of the PBAE properties variables showing quadrupole on the Z axis of the amine component as the key parameter demonstrated by the analysis of variable importance is in agreement with the mechanisms of action and it was also found to be one the key parameters when PLS regression was carried out using not only chemico-physical properties but also experimentally determined characteristics (diffusion coefficient, zeta potential and molecular weight of the polymer)³⁹. These PBAEs properties were not explicitly considered in the work as it was assumed that they depend on the properties of the amine and acrylate constituents and that the ML models would capture the correlation between drug uptake and polymer properties such as MW implicitly.

The optimal components identified here are structurally very different from those found as optimal copolymers for PBAE application in DNA vector^{40,57} and a direct consequence of the different mechanisms involved

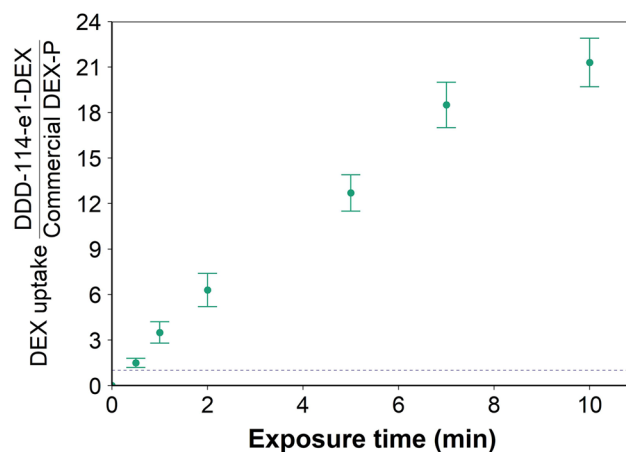


Figure 8. Uptake profile of DEX in cartilage using predicted best performing PBAE (DDD_114_e1-DEX).

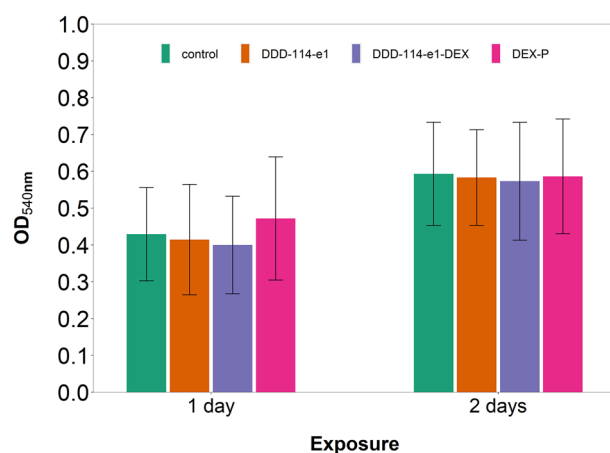


Figure 9. Mitochondrial activity of chondrocytes in cartilage explants cultured with basal media or medium containing predicted best performing PBAE (DDD_114_e1-DEX).

in the technology (DNA binding and cell membrane penetration vs. electrostatic attraction toward negatively charged GAG chains in cartilage).

The application of ML to PBAE structure optimisation for drug delivery in cartilage presented in this work can also potentially act as blueprint for the optimisation of other applications of PBAE such as drug releasing degradable coatings⁵⁸, non-viral DNA vectors for gene therapy⁴⁰ and mRNA vaccines⁵⁷ fast-tracking products to patients where, to date, only a lab based combinatorial chemistry approach to optimisation has been undertaken⁵⁹. The expected reduction in the time required to screen numerous polymers will also be coupled with monetary saving in the drug development costs with clear benefits not only to patients but also to health care providers.

We demonstrated an ML guided drug design optimisation approach that accurately predicts the relation between structure/property and outcome requiring only 2% of the compositional space (90 out of at least 3915 copolymers) to be experimentally explored. Our work led to the discovery of several PBAEs expected to result in a higher drug uptake than those of previously reported candidates. The actual efficacy was also determined and found to be very close to that predicted by the model (Fig. 8). Additionally, no negative impact on chondrocytes viability (Fig. 9) was detected for such PBAE as in line with the well-known safety of such polymers^{40,59}. Moreover, the trends uncovered between properties and efficacy of the polymers, along with the non-intuitive optimal design elements of PBAE for cartilage delivery identified in this study, such as the presence pyrazine in the amine constituent (likely related to the increased hydrogenation of the nitrogen atom), are also critical in the search for next-generation polymer driven cartilage delivery systems.

Methods

PBAE are denoted throughout the text with a code containing letters referring to the diacrylate (Fig. S1) and numbers (Fig. S2) referring to the amine; for example, A5 is the polymer made from 1,4-Butanediol diacrylate and 3-(dimethylamino)propylamine. The polymer backbone code is followed by e1 for PBAE end-capped with ethylene-diamine and e2 for PBAE end-capped with diethylene-triamine.

Data analysis. All models were fitted through R⁶⁰ and all other necessary packages necessary to perform regression with the “caret” package⁶¹.

Manhattan distance between PBAEs and complete distance between clusters were used for generating dendrograms.

Datasets and descriptors. Two PBAE uptake datasets were used to develop predictions, the publicly available set³⁹ was expanded with a purposely obtained new set collected after the inclusion of further acrylate monomers in the library.

Drug uptake predictions were performed utilizing physical and chemical parameters of amine and acrylate components of each PBAE obtained from PubChem library (Mw, logP, tPSA, Complexity, Heavy Atom Count, Volume 3D, X_Steric Quadrupole 3D, Y_Steric Quadrupole 3D, Z_Steric Quadrupole 3D); along with parameters related to the repeated polymeric unit (amine + acrylate) calculated through ChemDraw. The later included boiling point, melting point, critical volume and pressure, Gibb's free energy, logP (partition-coefficient between two immiscible phases at equilibrium which is proportional to hydrophobicity), solubility (logS), pKa, molar refractivity (CMR), heat of formation and the topological polar surface area (tPSA), which represent the total area of all polar atoms (mainly oxygen and nitrogen) including their affixed hydrogen atoms.

Kth nearest neighbour imputation was employed to handle missing data⁴⁷.

Machine learning algorithms training and predictions. Outcome data were transformed ($1/y^4$) to achieve a distribution of the drug uptake closer to a gaussian profile; moreover, input values for possible predictive variables were centred and scaled using mean and standard deviation.

A random split of the PBAEs into training (75%) and test (25%) datasets was applied. Weights to each point were assigned proportionally based on the distance from the median. Classic Machine Learning methods, such as Bernoulli Naive Bayes, Elastic regularisation, kNN, generalised additive models (GAM), Decision Tree, Random Forests, Neural Networks and SVM were employed to establish correlations between predictors and drug uptake. Tuning and hyper parameters search for each model were conducted through tenfold cross validation repeated 3 times on the training dataset; final model selection was based on minimisation of RMSE. The same training and test data set were employed for all models tested.

The best performing predictive model was used to estimate the drug uptake of the PBAE not previously experimentally tested in a two-steps approach. During the first round, amine and acrylates exhibiting a variety of structural features and moieties was employed to recognise critical patterns. In round2, variations of the pivotal characteristics observed in round1 were explored to further refine the optimal candidate.

Data availability

The datasets used and/or analysed during the current study are available from the corresponding author on reasonable request.

Received: 5 March 2022; Accepted: 9 August 2022

Published online: 20 August 2022

References

- Gupta, R. *et al.* Artificial intelligence to deep learning: Machine intelligence approach for drug discovery. *Mol. Divers.* **25**, 1315–1360. <https://doi.org/10.1007/s11030-021-10217-3> (2021).
- Réda, C., Kaufmann, E. & Delahaye-Duriez, A. Machine learning applications in drug development. *Comput. Struct. Biotechnol. J.* **18**, 241–252. <https://doi.org/10.1016/j.csbj.2019.12.006> (2020).
- Deloitte Centre for Health Solutions - Embracing the future of work to unlock RD productivity. <https://www2.deloitte.com/content/dam/Deloitte/uk/Documents/life-sciences-health-care/deloitte-uk-measuring-roipharma.pdf>.
- Morgan, S., Grootendorst, P., Lexchin, J., Cunningham, C. & Greyson, D. The cost of drug development: A systematic review. *Health Policy* **100**, 4–17. <https://doi.org/10.1016/j.healthpol.2010.12.002> (2011).
- Reis, M. *et al.* Machine-learning-guided discovery of 19F MRI agents enabled by automated copolymer synthesis. *J. Am. Chem. Soc.* **143**, 17677–17689. <https://doi.org/10.1021/jacs.1c08181> (2021).
- Ekins, S. *et al.* Exploiting machine learning for end-to-end drug discovery and development. *Nat. Mater.* **18**, 435–441. <https://doi.org/10.1038/s41563-019-0338-z> (2019).
- Struble, T. J. *et al.* Current and future roles of artificial intelligence in medicinal chemistry synthesis. *J. Med. Chem.* **63**, 8667–8682. <https://doi.org/10.1021/acs.jmedchem.9b02120> (2020).
- Paul, D. *et al.* Artificial intelligence in drug discovery and development. *Drug Discov. Today* **26**, 80–93. <https://doi.org/10.1016/j.drudis.2020.10.010> (2021).
- Kimber, T. B., Chen, Y. & Volkamer, A. Deep learning in virtual screening: Recent applications and developments. *Int. J. Mol. Sci.* **22**, 4435. <https://doi.org/10.3390/ijms22094435> (2021).
- Moosavi, S. M., Jablonka, K. M. & Smit, B. The role of machine learning in the understanding and design of materials. *J. Am. Chem. Soc.* **142**, 20273–20287. <https://doi.org/10.1021/jacs.0c09105> (2020).
- Baum, Z. J. *et al.* Artificial intelligence in chemistry: Current trends and future directions. *J. Chem. Inf. Model.* **61**, 3197–3212. <https://doi.org/10.1021/acs.jcim.1c00619> (2021).
- Butler, K. T., Davies, D. W., Cartwright, H., Isayev, O. & Walsh, A. Machine learning for molecular and materials science. *Nature* **559**, 547–555. <https://doi.org/10.1038/s41586-018-0337-2> (2018).
- Stephenson, N. *et al.* Survey of machine learning techniques in drug discovery. *Curr. Drug Metab.* **20**, 185–193. <https://doi.org/10.2174/1389200219666180820112457> (2019).
- Khan, S. R., Al Rijjal, D., Piro, A. & Wheeler, M. B. Integration of AI and traditional medicine in drug discovery. *Drug Discov. Today* **26**, 982–992. <https://doi.org/10.1016/j.drudis.2021.01.008> (2021).
- Rohall, S. L. *et al.* An Artificial intelligence approach to proactively inspire drug discovery with recommendations. *J. Med. Chem.* **63**, 8824–8834. <https://doi.org/10.1021/acs.jmedchem.9b02130> (2020).
- Yi, Z. *et al.* Mapping drug-induced neuropathy through in-situ motor protein tracking and machine learning. *J. Am. Chem. Soc.* **143**, 14907–14915. <https://doi.org/10.1021/jacs.1c07312> (2021).

17. Espinoza, G. Z., Angelo, R. M., Oliveira, P. R. & Honorio, K. M. Evaluating deep learning models for predicting ALK-5 inhibition. *PLoS ONE* **16**, e0246126. <https://doi.org/10.1371/journal.pone.0246126> (2021).
18. Stokes, J. M. *et al.* A deep learning approach to antibiotic discovery. *Cell* **180**, 688–702.e613. <https://doi.org/10.1016/j.cell.2020.01.021> (2020).
19. Bess, A. *et al.* Artificial intelligence for the discovery of novel antimicrobial agents for emerging infectious diseases. *Drug Discov. Today* <https://doi.org/10.1016/j.drudis.2021.10.022> (2021).
20. Kundu, S. *et al.* Enabling early detection of osteoarthritis from presymptomatic cartilage texture maps via transport-based learning. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 24709–24719. <https://doi.org/10.1073/pnas.1917405117> (2020).
21. Tsigelny, I. F. Artificial intelligence in drug combination therapy. *Brief. Bioinform.* **20**, 1434–1448. <https://doi.org/10.1093/bib/bby004> (2019).
22. Patel, L., Shukla, T., Huang, X., Ussery, D. W. & Wang, S. Machine learning methods in drug discovery. *Molecules* **25**, 5277. <https://doi.org/10.3390/molecules25225277> (2020).
23. Gao, C. *et al.* Innovative materials science via machine learning. *Adv. Funct. Mater.* **32**, 2108044. <https://doi.org/10.1002/adfm.202108044> (2022).
24. Yin, Z.-W. *et al.* Advanced electron energy loss spectroscopy for battery studies. *Adv. Funct. Mater.* **32**, 2107190. <https://doi.org/10.1002/adfm.202107190> (2022).
25. Miljković, F., Rodríguez-Pérez, R. & Bajorath, J. Impact of artificial intelligence on compound discovery, design, and synthesis. *ACS Omega* <https://doi.org/10.1021/acscomega.1c05512> (2021).
26. Tkatchenko, A. Machine learning for chemical discovery. *Nat. Commun.* **11**, 4125. <https://doi.org/10.1038/s41467-020-17844-8> (2020).
27. Ripphausen, P., Nisius, B., Peltason, L. & Bajorath, J. Quo vadis, virtual screening? A comprehensive survey of prospective applications. *J. Med. Chem.* **53**, 8461–8467. <https://doi.org/10.1021/jm101020z> (2010).
28. Kimber, T. B., Chen, Y. & Volkamer, A. Deep learning in virtual screening: Recent applications and developments. *Int. J. Mol. Sci.* <https://doi.org/10.3390/ijms22094435> (2021).
29. Gautam, V., Gaurav, A., Masand, N., Lee, V. S. & Patil, V. M. Artificial intelligence and machine-learning approaches in structure and ligand-based discovery of drugs affecting central nervous system. *Mol. Divers.* <https://doi.org/10.1007/s11030-022-10489-3> (2022).
30. Mizera, M. & Latek, D. Ligand-receptor interactions and machine learning in GCGR and GLP-1R drug discovery. *Int. J. Mol. Sci.* <https://doi.org/10.3390/ijms22084060> (2021).
31. Gawriljuk, V. O. *et al.* Development of machine learning models and the discovery of a new antiviral compound against yellow fever virus. *J. Chem. Inf. Model.* **61**, 3804–3813. <https://doi.org/10.1021/acs.jcim.1c00460> (2021).
32. Safiri, S. *et al.* Global, regional and national burden of osteoarthritis 1990–2017: A systematic analysis of the Global Burden of Disease Study 2017. *Ann. Rheum. Dis.* **79**, 819–828. <https://doi.org/10.1136/annrheumdis-2019-216515> (2020).
33. Buckwalter, J. A., Mankin, H. J. & Grodzinsky, A. J. Articular cartilage and osteoarthritis. *Instr. Course Lect.* **54**, 465–480 (2005).
34. Bajpayee, A. G., Wong, C. R., Bawendi, M. G., Frank, E. H. & Grodzinsky, A. J. Avidin as a model for charge driven transport into cartilage and drug delivery for treating early stage post-traumatic osteoarthritis. *Biomaterials* **35**, 538–549. <https://doi.org/10.1016/j.biomaterials.2013.09.091> (2014).
35. Geiger, B., Grodzinsky, A. & Hammond, P. - Designing Drug Delivery Systems for Articular Joints - May 2018 Chemical Engineering Progress (CEP) - American Institute of Chemical Engineers (AIChE)
36. Geiger, B. C., Wang, S., Padera, R. F., Grodzinsky, A. J. & Hammond, P. T. Cartilage-penetrating nanocarriers improve delivery and efficacy of growth factor treatment of osteoarthritis. *Sci. Transl. Med.* **10**, eaat8800. <https://doi.org/10.1126/scitranslmed.aat8800> (2018).
37. Jacobs J.W.G. & Bijlsma J.W.J. Glucocorticoid therapy. in *Kelley's Textbook of Rheumatology* 7th edn. 870–874 (Elsevier Saunders, 2005).
38. Perni, S. & Prokopovich, P. Poly-beta-amino-esters nano-vehicles based drug delivery system for cartilage. *Nanomedicine* **13**, 539–548. <https://doi.org/10.1016/j.nano.2016.10.001> (2017).
39. Perni, S. & Prokopovich, P. Optimisation and feature selection of poly-beta-amino-ester as a drug delivery system for cartilage. *J. Mater. Chem. B* **8**, 5096–5108. <https://doi.org/10.1039/c9tb02778e> (2020).
40. Green, J. J., Langer, R. & Anderson, D. G. A combinatorial polymer library approach yields insight into nonviral gene delivery. *Acc. Chem. Res.* **41**, 749–759. <https://doi.org/10.1021/ar7002336> (2008).
41. Burger, S. V. Introduction to machine learning with R: Rigorous mathematical analysis. (2018).
42. Friedman, J., Hastie, J. & Tibshirani, R. The elements of statistical learning. (2009).
43. Russo, D. P., Zorn, K. M., Clark, A. M., Zhu, H. & Ekins, S. Comparing multiple machine learning algorithms and metrics for estrogen receptor binding prediction. *Mol. Pharm.* **15**, 4361–4370. <https://doi.org/10.1021/acs.molpharmaceut.8b00546> (2018).
44. Korotcov, A., Tkachenko, V., Russo, D. P. & Ekins, S. Comparison of deep learning with multiple machine learning methods and metrics using diverse drug discovery data sets. *Mol. Pharm.* **14**, 4462–4475. <https://doi.org/10.1021/acs.molpharmaceut.7b00578> (2017).
45. Fan, Y. *et al.* Investigation of machine intelligence in compound cell activity classification. *Mol. Pharm.* **16**, 4472–4484. <https://doi.org/10.1021/acs.molpharmaceut.9b00558> (2019).
46. Guan, X. *et al.* Clinical and inflammatory features based machine learning model for fatal risk prediction of hospitalized COVID-19 patients: results from a retrospective cohort study. *Ann. Med.* **53**, 257–266. <https://doi.org/10.1080/07853890.2020.1868564> (2021).
47. Kuhn, M. & Johnson, K. Applied Predictive Modeling. (2013).
48. Zhou, Z. H. *Ensemble Methods: Foundations and Algorithms* (Chapman and Hall/CRC, 2012).
49. Fanourgakis, G. S., Gkagkas, K., Tylianakis, E. & Froudakis, G. E. A universal machine learning algorithm for large-scale screening of materials. *J. Am. Chem. Soc.* **142**, 3814–3822. <https://doi.org/10.1021/jacs.9b11084> (2020).
50. Chen, J. *et al.* Machine learning aids classification and discrimination of noncanonical DNA folding motifs by an arrayed host: Guest sensing system. *J. Am. Chem. Soc.* **143**, 12791–12799. <https://doi.org/10.1021/jacs.1c06031> (2021).
51. Jang, J., Gu, G. H., Noh, J., Kim, J. & Jung, Y. Structure-based synthesizability prediction of crystals using partially supervised learning. *J. Am. Chem. Soc.* **142**, 18836–18843. <https://doi.org/10.1021/jacs.0c07384> (2020).
52. Guo, Y. *et al.* Machine-learning-guided discovery and optimization of additives in preparing Cu catalysts for CO₂ reduction. *J. Am. Chem. Soc.* **143**, 5755–5762. <https://doi.org/10.1021/jacs.1c00339> (2021).
53. Xie, Y. *et al.* Machine learning assisted synthesis of metal-organic nanocapsules. *J. Am. Chem. Soc.* **142**, 1475–1481. <https://doi.org/10.1021/jacs.9b11569> (2020).
54. Hatakeyama-Sato, K., Tezuka, T., Umeki, M. & Oyaizu, K. AI-assisted exploration of superionic glass-type Li⁺ conductors with aromatic structures. *J. Am. Chem. Soc.* **142**, 3301–3305. <https://doi.org/10.1021/jacs.9b11442> (2020).
55. Tiihonen, A. *et al.* Predicting antimicrobial activity of conjugated oligoelectrolyte molecules via machine learning. *J. Am. Chem. Soc.* **143**, 18917–18931. <https://doi.org/10.1021/jacs.1c05055> (2021).
56. Saeedi, T. & Prokopovich, P. Poly beta amino ester coated emulsions of NSAIDs for cartilage treatment. *J. Mater. Chem. B* **9**, 5837–5847. <https://doi.org/10.1039/d1tb01024g> (2021).

57. Capasso Palmiero, U., Kaczmarek, J. C., Fenton, O. S. & Anderson, D. G. Poly(β -amino ester)-co-poly(caprolactone) Terpolymers as nonviral vectors for mRNA delivery in vitro and in vivo. *Adv. Healthc. Mater.* **7**, e1800249. <https://doi.org/10.1002/adhm.201800249> (2018).
58. Moskowitz, J. S. *et al.* The effectiveness of the controlled release of gentamicin from polyelectrolyte multilayers in the treatment of *Staphylococcus aureus* infection in a rabbit bone model. *Biomaterials* **31**, 6019–6030. <https://doi.org/10.1016/j.biomaterials.2010.04.011> (2010).
59. Anderson, D. G., Lynn, D. M. & Langer, R. Semi-automated synthesis and screening of a large library of degradable cationic polymers for gene delivery. *Angew. Chem. Int. Ed.* **42**, 3153–3158. <https://doi.org/10.1002/anie.200351244> (2003).
60. R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. (2019).
61. Kuhn, M. The caret package. *J. Stat. Softw.* **28**, 1–26 (2012).

Acknowledgements

The work has been supported by Pathfinder Fund and by Wellcome Trust.

Author contributions

S.P.: Conceptualization, Methodology, Software, Investigation, Formal analysis, Data Curation, Visualization, Writing—Original Draft. P.P.: Validation, Methodology, Funding acquisition, Resources, Writing—Review & Editing, Project administration.

Competing interests

SP is named inventors on patents related to the use of PBAE as drug delivery systems. PP is named inventors on patents related to the use of PBAE as drug delivery systems.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-022-18332-3>.

Correspondence and requests for materials should be addressed to P.P.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022