

High-Throughput Virtual Screening of Small Molecule Inhibitors for SARS-CoV-2 Protein Targets with Deep Fusion Models

Garrett A. Stevenson¹, Derek Jones¹, Hyojin Kim¹, W. F. Drew Bennett¹, Brian J. Bennion¹, Monica Borucki¹, Feliza Bourguet¹, Aidan Epstein¹, Magdalena Franco¹, Brooke Harmon², Stewart He¹, Max P. Katz³, Daniel Kirshner¹, Victoria Lao¹, Edmond Y. Lau¹, Jacky Lo¹, Kevin McLoughlin¹, Richard Mosesso², Deepa K. Muruges¹, Oscar A. Negrete², Edwin A. Saada¹, Brent Segelke¹, Maxwell Stefan², Marisa W. Torres¹, Dina Weilhammer¹, Sergio Wong¹, Yue Yang¹, Adam Zemla¹, Xiaohua Zhang¹, Fangqiang Zhu¹, Felice C. Lightstone¹, Jonathan E. Allen^{1*}
allen99@llnl.gov

¹Lawrence Livermore National Laboratory
Livermore, California, USA
²Sandia National Laboratories
Livermore, California, USA
³NVIDIA Corporation
Santa Clara, California, USA

ABSTRACT

Structure-based Deep Fusion models were recently shown to outperform several physics- and machine learning-based protein-ligand binding affinity prediction methods. As part of a multi-institutional COVID-19 pandemic response, over 500 million small molecules were computationally screened against four protein structures from the novel coronavirus (SARS-CoV-2), which causes COVID-19. Three enhancements to Deep Fusion were made in order to evaluate more than 5 billion docked poses on SARS-CoV-2 protein targets. First, the Deep Fusion concept was refined by formulating the architecture as one, coherently backpropagated model (Coherent Fusion) to improve binding-affinity prediction accuracy. Secondly, the model was trained using a distributed, genetic hyper-parameter optimization. Finally, a scalable, high-throughput screening capability was developed to maximize the number of ligands evaluated and expedite the path to experimental evaluation. In this work, we present both the methods developed for machine learning-based high-throughput screening and results from using our computational pipeline to find SARS-CoV-2 inhibitors.

KEYWORDS

deep learning, hyper-parameter optimization, SARS-CoV-2, COVID-19, HPC, GPU, AI

ACM Reference Format:

Garrett A. Stevenson¹, Derek Jones¹, Hyojin Kim¹, W. F. Drew Bennett¹, Brian J. Bennion¹, Monica Borucki¹, Feliza Bourguet¹, Aidan Epstein¹, Magdalena Franco¹, Brooke Harmon², Stewart He¹, Max P. Katz³, Daniel

Kirshner¹, Victoria Lao¹, Edmond Y. Lau¹, Jacky Lo¹, Kevin McLoughlin¹, Richard Mosesso², Deepa K. Muruges¹, Oscar A. Negrete², Edwin A. Saada¹, Brent Segelke¹, Maxwell Stefan², Marisa W. Torres¹, Dina Weilhammer¹, Sergio Wong¹, Yue Yang¹, Adam Zemla¹, Xiaohua Zhang¹, Fangqiang Zhu¹, Felice C. Lightstone¹, Jonathan E. Allen¹. 2021. High-Throughput Virtual Screening of Small Molecule Inhibitors for SARS-CoV-2 Protein Targets with Deep Fusion Models. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

The COVID-19 disease caused by the severe acute respiratory syndrome coronavirus (SARS-CoV-2) is responsible for the most recent, severe pandemic in modern human history [33]. At the onset of the COVID-19 pandemic, a worldwide effort began to identify and provide target proteins for vaccine and drug development to neutralize the virus. Two distinct proteins were rapidly solved; the trimeric spike protein (spike), which binds to human ACE2 to enter human cells [3] and the main protease (M^{Pro}), which plays a pivotal role in viral gene expression and replication [59]. In response to the pandemic, we participated in a large-scale multi-institutional effort to virtually screen, experimentally test, and optimize therapeutic leads targeting the spike and M^{Pro} SARS-CoV-2 protein targets. Two different binding sites from the spike protein (denoted spike1, spike2) and two different conformations of the M^{Pro} active site (denoted protease1, protease2) were used in the high-throughput screening calculations.

Experimental tests of drug candidates are expensive and serve as a fundamental bottleneck in drug discovery. As a result, computational chemistry is used extensively to accelerate the discovery process by screening drugs and nominating the strongest candidates for experimental validation [64]. High Performance Computing (HPC) plays a critical role in virtual screening [6, 28, 30, 69, 70] by accelerating computationally expensive calculations and providing the scalability necessary to screen large numbers of candidate molecules. This is crucial, as the chemical space of potential

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference'17, July 2017, Washington, DC, USA

© 2021 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

molecules has been estimated to be on the order of 10^{60} [8, 49]. Accurately estimating protein-ligand binding affinities is an important step in drug discovery. However, even computationally expensive, biophysics-based scoring methods find predicting binding free energy a difficult task [27, 64, 72]. Deep learning methods represent an alternative, rapid approach to binding affinity prediction which alleviates the dependence on hand-curated features, which may not capture the mechanism of binding [1, 2].

The two leading deep learning approaches to structure-based binding affinity prediction fall into two categories: 3-dimensional Convolutional Neural Networks (3D-CNNs) [26, 48, 67] and Spatial Graph Convolutional Neural Networks (SG-CNNs) [13, 37, 71]. Fundamentally, 3D-CNN models exploit a voxelized representation of atoms in a 3D grid, which portrays protein-ligand compounds in a Euclidean space for inference. On the other hand, SG-CNN approaches leverage a graph representation of the protein-ligand complex, allowing for multiple "edge-types" to be encoded in the representation (e.g., distinct distance thresholds corresponding to covalent or non-covalent interactions) to sub-select groups of atoms and evaluate their pairwise interactions.

These methods are significantly different in the way they represent compounds and their mechanisms of inference. However, they seek to achieve the same goal: accurate prediction of binding free energy on novel compounds. This observation led to the hypothesis that the 3D-CNN and SG-CNN likely have complementary strengths, which could be exploited by fusing the latent spaces of each model's learned features. This approach to "Fusion" modeling was explored and shown to achieve superior generalization performance in predicting protein-ligand binding on X-ray crystallographic structures and virtually-docked poses of protein-ligand compounds in hold-out test sets [27]. In this work, we detail improvements to Fusion, describe its utility in high-throughput screening, and evaluate its application to the SARS-CoV-2 drug discovery problem. A key innovation reported here is the potential to train a fusion model "coherently". While Coherent Fusion comes with the increased computational burden of training a more complex model, the increase in complexity is mitigated by using HPC to perform parallel, distributed training.

2 DEEP FUSION

2.1 Fusion Modeling in Computational Chemistry

Machine learning, specifically deep learning, approaches to protein-ligand binding affinity prediction, represent a promising new development in drug discovery [1, 2, 13, 16, 26, 56, 67, 71, 72]. At a high level, the deep-learned models being proposed for binding affinity prediction are single-pass, feed-forward systems. This fundamental model formulation results in a computational advantage, in that, the models quickly predict binding affinity in one pass over their input. The simplicity and speed of deep learning prediction relative to biophysics-based computations make them especially attractive in the context of massive virtual drug screens.

The concept of fusion is a recent development in deep learning, initially applied to computer vision problems [35, 50, 60, 66]. In fusion, models are combined by integrating different modes of data or approaches for more predictive power. This concept was recently

applied to computational chemistry in the form of Fusion models for Atomic and molecular Structures (FAST) [27], where fusion of the two leading deep learning models (3D-CNN, SG-CNN) was shown to improve binding affinity prediction. Specifically, the Late Fusion and Mid-level Fusion models are shown as approximately equivalent or superior to individual SG-CNNs, individual 3D-CNNs, other state-of-the-art deep learning models [26, 56], and physics-based approaches including both Autodock Vina [58] and Molecular Mechanics - Generalized Born / Surface Area (MM/GBSA) which has been shown to improve ligand pose ranking for certain target proteins [19, 64]. The Late Fusion and Mid-level Fusion FAST models were specifically described and made publicly available. In the context of this work, we retrain and optimize the models on data from PDBbind-2019 [62].

$$pK = -\log(K); \text{ where } K = K_i, K_d, \text{ or } IC_{50} \quad (1)$$

The Late Fusion approach is simple, but its superior performance compared to individual SG-CNNs and 3D-CNNs shows the potential of fusion modeling. In Late Fusion, SG-CNN and 3D-CNN models were separately trained to predict absolute binding affinity (Equation 1), where absolute binding affinity is defined as the negative logarithm of a binding constant K . Binding affinity data in this study is measured as an inhibitory constant K_i , disassociation constant K_d , or inhibitory activity (IC_{50}), where these measurements are treated as equivalent labels in our calculations. Late Fusion takes the unweighted arithmetic mean across the individually-predicted binding affinity values from the SG-CNN and 3D-CNN models.

The Mid-level Fusion model is also described in [27]. The model is defined by extracting the latent space feature vector from $Layer^{N-3}$ of an N-layer SG-CNN and $Layer^{M-1}$ of an M-layer 3D-CNN. Each vector of gradients is then processed through model-specific dense layers, concatenated with the originally extracted vectors, and passed through two fusion dense layers for a final prediction. In contrast to Late Fusion, Mid-level Fusion is a non-linear combination of the respective models and its performance has been shown to outperform Late Fusion in some cases.

2.2 Advancing Fusion

Faced with the imperative task of virtually screening molecules for inhibition of SARS-CoV-2 activity, we sought to apply and improve the fusion deep learning approach. First, we updated the previously trained models with data from the latest version of PDBbind (2019) [62], which provides approximately 4000 more compounds than the 2016 version used to train the original models and, thereby, offers the potential for improved generalization

Our second step looked to improve the 3D-CNN and SG-CNN models through recent developments in the optimization of hyper-parameters. Deep learning models are highly sensitive to a human's definition of their hyper-parameters "model architecture, loss function, and optimization algorithm, among others" [24]. Automated hyper-parameter optimization was first addressed with parallel searches (grid or random), followed by sequential optimization methods such as Bayesian optimization [4, 21, 53, 55]. Hyper-parameter optimization algorithms have improved over time in their ability to scale [54]. Evolutionary Algorithms (EAs) have been shown to further improve the optimization process [5, 20, 24, 34].

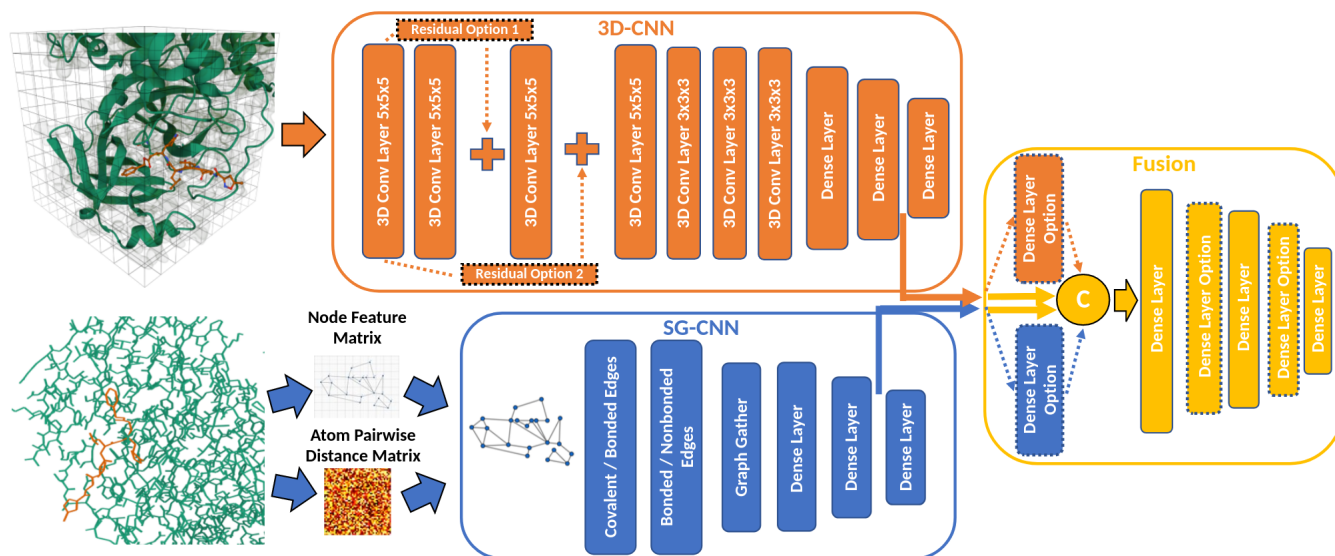


Figure 1: Fusion model architecture with voxelized and spatial graph inputs of COVID-19 MP^{ro} (PDB: 6LU7, denoted protease1, green) in complex with an N3 inhibitor (red) [51]. Components of the 3D-CNN (orange), SG-CNN (blue), and Fusion layers (yellow) which were given as options to the hyper-parameter optimization are shown with dashed lines/borders.

Recently, a leading population-based EA, Population-Based Bandits (PB2) [24], was improved by formulating hyper-parameter optimization as a Gaussian Process (GP) bandit optimization of a time-varying function [45].

Fusion modeling in particular is marked by a significant exposure to hyper-parameters. Both the 3D-CNN and SG-CNN models have their own hyper-parameters which enable them to learn the binding affinity prediction problem optimally in isolation. Additionally, the fusion layers require another set of hyper-parameters necessary to find an optimal non-linear combination of the two models. With this in mind, we saw an opportunity for improving fusion significantly by using a PB2 automated optimization algorithm [45] on Lassen, one of the most powerful high-performance computers in the world [39]. The libraries used to define the model and optimization architectures are PyTorch [46], Pytorch Geometric [14], and Ray/Ray[Tune] [43].

Finally, a new formulation of fusion, the Coherent Fusion model, was developed as a potential improvement to the previous Late and Mid-level Fusion models. In both the existing fusion approaches, a 3D-CNN and SG-CNN are individually optimized to minimize the mean squared error (MSE) between their predictions of binding free energy and ground-truth experimental values. The existing Mid-level approach then combines the independently optimized models to form a stronger predictor, by learning the latent space strengths of each model. However, in both Late and Mid-level Fusion, the 3D-CNN and SG-CNN weights are unaltered and remain in the state that was optimal for isolated prediction. Given the Late and Mid-level Fusion models’s superior performance compared to their isolated components, we hypothesized that fusion might be further improved by coherently backpropagating gradient through both the fusion layers and the separate models. In doing so, the Coherent Fusion model fine-tunes both the 3D-CNN and SG-CNN

heads to cooperatively exploit their strengths in a joint optimization. The drawbacks of Coherent Fusion are both an increased hyper-parameter search space and number of trainable parameters. To address this, we developed a parallel, distributed hyper-parameter optimization training architecture. Compared with the models in [27], the combination of these modifications to the concept of Fusion led to significant differences in the hyper-parameters for the 3D-CNN, SG-CNN, and Fusion layers of the model.

3 OPTIMIZATION AND EVALUATION

3.1 Data

PDBbind-2019 is a curated subset of the larger Protein Data Bank (PDB) [61], which is widely used to tune biophysics- and machine learning-based methods [13, 26, 27, 37, 64]. The PDBbind data set is comprised of crystal structures arranged into two groups (*general* and *refined*) based on size (where protein-ligand compounds containing a ligand with molecular weight >1000 Daltons (Da) are excluded from the refined set), data quality (where compounds with a measured IC_{50} but no K_i or K_d measurements are excluded from the refined set), and resolution of the crystal structure (<2.5 Angstroms). From the *refined* set, a third, *core* set is extracted using a clustering protocol based on protein-sequence similarity. The *core* set is compiled to represent a valid test for scoring methods by creating a high-quality subset of compounds sufficiently different from the *general* and *refined* sets. As such, we use the *core* set as a primary means for evaluating the Fusion methods considered and comparing against published literature.

In this study, we employ the quintile sub-sampling method from [27] to formulate training and validation sets from the PDBbind-2019 *general* and *refined* groupings. The sub-sampling is done independently on the *general* and *refined* sets and 10% of the examples from each are withdrawn to form the validation set. Quintile

sub-sampling guarantees both the training and validation sets to represent the full range of binding affinity values across PDBbind, where simple random sampling holds the risk of training and validating models on different sub-spaces of affinity values [10]. The outcome is a training set of 15,631 complexes, a validation set of 1,731 complexes, and the 290 PDBbind *core* set complexes are held-out for evaluation. Details of pre-processing and feature extraction for the PDBbind data can be found in [27], where here the same tools [25, 41, 44, 47] and sequence of operations are used.

As a means for additional evaluation, we supplement the *core* set of 290 crystal structures from PDBbind with virtually-docked representations of the complexes. In practice, docking pose data is used for large-scale virtual screening, but is noisy and error prone since the correct ligand pose is not known until a co-complex is crystallized experimentally. Therefore, a scoring function’s performance in the docking space is critical in gauging its robustness to noise and its pragmatic utility.

We leverage the ConveyorLC toolchain [68–70] to produce all docking complex data, as it is used in our high-throughput virtual screening pipeline. ConveyorLC generates docking poses using the Vina scoring function [58], then re-scores up to 10 best docking poses using MM/GBSA on a subset of the larger virtual screen; only a subset is re-scored because MM/GBSA is orders of magnitude more computationally expensive than docking. This sequence of down selecting to limit the search space, accompanied by increasingly complex analyses, is frequently used in drug discovery pipelines, and even molecular dynamics (MD) simulations can be used before finalizing candidates for physical experimentation. The opportunity for machine learning models, like Deep Fusion, is to replace or supplement a more costly stage of a drug discovery pipeline with either improved accuracy or speed.

3.2 Training Architecture

Our approach to train the various individual and fusion models was executed iteratively. Lassen uses the IBM Spectrum LSF Job Scheduler [22], which necessitates pausing, rescheduling, and resuming training jobs after a maximum run-time. As the hyper-parameter optimization began to converge, the range of hyper-parameter values was adjusted, when possible, to ensure the lower and upper-bounds of the search space were not limiting factors in model performance. The full scope of hyper-parameters and ranges evaluated for each model are provided in Table 1, where the ranges are binary (T/F), a list of options, uniformly sampled continuous variables, or not applicable (N/A).

We used the Ray/Ray[Tune] [43] Python library extensively as the foundation for running trials of individual hyper-parameter combinations within the context of a PB2 optimization. Together, Ray and PyTorch provide the ability to accelerate the training process by distributing individual trials across multiple nodes/GPUs. Each of Lassen’s 792 GPU nodes is made up of 44 3.45 GHz Power9 CPU cores, 4 NVIDIA Volta V100 GPUs each with 16GB of memory, and 256 GB of main memory. Depending on the complexity of each model, we distributed individual hyper-parameter configurations between 1 and 12 ranks (1 rank = 1 GPU, 10 CPU cores, 64 GB memory) to expedite the training process. Each rank also utilized 24 data workers running in parallel to pre-load future

Table 1: Hyper-parameters for each model and their range of values considered by the PB2 optimization

Hyper-parameter	3D-CNN	SG-CNN	Fusion
Optimizer	Adam [29]	Adam [29]	Adam [29], AdamW [40], RMSprop [17], Adadelta [9]
Activation function	ReLU	ReLU	ReLU, LReLU [65], SELU [31]
Batch size	8,12,24	4,8,12,16	1,2,4,5,8,12,16,24,28,34,38,48,56
Learning rate	$1e^{-6}$ - $1e^{-4}$	$2e^{-4}$ - $2e^{-2}$	$1e^{-8}$ - $1e^{-3}$
Model-Specific Fusion Layers	N/A	N/A	T/F
Epochs	0-150	0-350	0-500
Pre-trained	F	F	T/F
Batch norm.	T/F	F	T/F
Dropout 1 (early)	0.25	0	0 - 0.50
Dropout 2 (mid)	0.125	0	0 - 0.25
Dropout 3 (late)	0	0	0 - 0.125
# of Fusion Layers	N/A	N/A	3,4,5
# of Dense Nodes	40,64,88,104,128	N/A	8,24,40,64,88,104,128
Residual Option 1	T/F	N/A	T/F
Residual Option 2	T/F	N/A	T/F
# of Conv. Filters 1	32,64,96	N/A	32,64,96
# of Conv. Filters 2	64,96,128	N/A	64,96,128
Non-covalent / Covalent K	N/A	2,3,4,5,6,7,8	2,3,4,5,6,7,8
Non-covalent / Covalent Neighbor Threshold	N/A	1.2Å-5.9Å	1.2Å-5.9Å
Non-covalent / Covalent Gather Width	N/A	8,24,40,64,88,104,128	8,24,40,64,88,104,128

batches. The combination of distributed model training and parallel data loading was central to the feasibility of an experiment with this size/scope.

The PB2 hyper-parameter optimization was initialized with a quantile fraction ($\lambda\%$) of 50, a time scale (T) in Epochs, a perturbation interval (t_{ready}) of 100 Epochs, and an objective function (Q) of minimum validation set MSE loss [24, 45]. The procedure begins with a population of initial, randomly sampled hyper-parameter hypotheses. As every trial reaches the perturbation interval t_{ready} , PB2 looks at the model’s performance and determines if it is above or below the quantile fraction ($\lambda\%$). The best performing trials (above $\lambda\%$) continue, while the under-performing trials clone a top-performing configuration (exploiting) and modify it using a

Table 2: Final hyper-parameters for the SG-CNN

Hyper-parameter	Value
Epochs	213
Batch size	16
Learning rate	$2.66e^{-3}$
Non-covalent K	3
Covalent K	6
Non-covalent Neighbor Threshold	5.22Å
Covalent Neighbor Threshold	2.24Å
Non-covalent Gather Width	128
Covalent Gather Width	24

Table 3: Final hyper-parameters for the 3D-CNN

Hyper-parameter	Value
Epochs	75
Batch size	12
Learning rate	$4.90e^{-5}$
Batch normalization	<i>F</i>
# of Dense Nodes	128
# of Conv. Filters 1	32
# of Conv. Filters 2	64
Residual Option 1	<i>F</i>
Residual Option 2	<i>T</i>

parallel GP-bandit optimization (exploring). The training process produces both an optimal model and important information about how the hyper-parameters considered effect performance.

3.3 Model Architectures

We draw heavily from the original FAST network architectures in [27], which holds detailed descriptions of pre-processing, feature extraction, voxel grid sizing and atom propagation, which were unaltered. The following focuses on updates to the models and the final optimized hyper-parameter configurations for each component. For brevity, we list only the final optimized hyper-parameter values, where the advantage of PB2 is in its ability to learn a schedule of hyper-parameters to converge in an end-state [45].

3.3.1 Individual Models. The SG-CNN in this work is structurally unaltered from [27], which uses the PotentialNet [13] architecture based on Gated Graph Sequence Neural Networks [36]. The only notable difference is the size of the dense layers were set according to the Non-covalent Gather Width, such that it was sequentially reduced in size by a factor of 1.5 and then 2. A population of 90 SG-CNN trials produced the final model and hyper-parameter configuration given in Table 2.

The 3D-CNN model is slightly modified from the architecture in [27]. The model has dropout above the first two dense layers, 2 additional convolutional layers, the filter sizes begin at 5x5x5 and reduce to 3x3x3, the residual options shown in Figure 1 were fed to the hyper-parameter optimization, and similar to the SG-CNN, the second dense layer size was determined by the optimization and then sequentially reduced by a factor of 2. Again, a population

of 90 trials was used, the final hyper-parameter values are given in Table 3, where the optimization converged to using the second residual connection shown in Figure 1 and 32 to 64 filters for the 5x5x5 and 3x3x3 convolutional layers respectively. With this larger 3D-CNN architecture (deeper than in [27]) we found it beneficial to augment the input matrices for the training set by randomly rotating the input data in X, Y, and Z, each with a 10% probability of occurring. This random rotational augmentation was applied only to the voxelized representation of a compound. While the compound is fundamentally the same, altering its presentation to the model helps to prevent overtraining (e.g., learning rotation-dependent features) and to increase the effective size of the training data set.

Table 4: Final hyper-parameters for Mid-level Fusion

Hyper-parameter	Value
Epochs	64
Batch size	1
Learning rate	$4.03e^{-4}$
Batch normalization	<i>F</i>
Optimizer	Adam [29]
Activation function	SELU [31]
Residual Fusion Layers	<i>T</i>
Dropout Rate 1 (early)	0.251
Dropout Rate 2 (mid)	0.125
Dropout Rate 3 (late)	≈ 0
# of Fusion Layers	5

3.3.2 Late / Mid-level Fusion. The Late Fusion method was implemented the same as in [27]. On the other hand, the optimization led the Mid-level Fusion model to a modified structure. For Mid-level Fusion, every optional layer (dashed lines) in the yellow Fusion block of Figure 1 was turned on. Table 4 gives the final hyper-parameters for Mid-level Fusion, which are the output of a 180 individual trial population. The other minor differences are a SELU [31] activation was selected over the previous Leaky-ReLU activation [65], a final batch size of 1, and the usage of light dropout instead of none.

Table 5: Final hyper-parameters for Coherent Fusion

Hyper-parameter	Value
Pre-trained	<i>T</i>
Epochs	18
Batch size	48
Learning rate	$1.08e^{-4}$
Batch normalization	<i>F</i>
Optimizer	Adam [29]
Activation function	SELU [31]
Residual Fusion Layers	<i>F</i>
Dropout Rate 1 (early)	0.386
Dropout Rate 2 (mid)	0.247
Dropout Rate 3 (late)	0.055
# of Fusion Layers	4

Table 6: Performance of Fusion models on the PDBbind core set crystal structures

Model	RMSE	MAE	R ²	Pearson R	Spearman R
Pafnucy [56]	1.42	1.13	-	0.78	-
Mid-level Fusion	1.38	1.10	0.596	0.778	0.757
Late Fusion	1.33	1.07	0.623	0.813	0.805
Coherent Fusion	1.30	1.05	0.640	0.807	0.802
KDeep [26]	1.27	-	-	0.82	0.82

3.3.3 Coherent Fusion. In developing the Coherent Fusion model, it was unclear whether the same 3D-CNN and SG-CNN hyperparameter configurations found to be optimal in isolation would also be ideal for their collaborative prediction. As such, we gave the optimization the option to load the models individually trained for prediction or re-define their structure and train each head from scratch. Using the pre-trained models led to a significant improvement in validation loss. Therefore, Table V gives the final hyperparameters for the best performing Coherent Fusion model, which loads the weights from the SG-CNN in Table 2 and 3D-CNN model described in Table 3.

The Coherent Fusion model experiment optimized a population of 270 individual trials to produce a best performer. Interestingly, the Coherent Fusion model converged to exclude the model-specific dense layer options the Mid-level Fusion model uses (Figure 1) and used a simpler (4 fusion layers) architecture overall. Additionally, the Coherent Fusion used a larger batch size of 48 and significantly stronger dropout. Across the board, the Coherent Fusion model preferred a simpler Fusion architecture with significantly stronger regularization. Our intuition for this phenomenon is that the Coherent Fusion model adjusting a larger set of learned parameters allows for a simpler architecture, faster convergence, and heavier regularization compared to the Mid-level Fusion model, which serves as preliminary evidence of a stronger predictor.

3.4 Evaluation Results

Over 60,000 Lassen GPU hours were used to optimize the various models. All model iterations and intervals were not run across the same number of nodes, but our training architecture was run at its peak across 66 Lassen nodes capable of over 7,300 TFLOPS using 2904 CPU cores and 264 GPUs to train in parallel.

In Table 6, the Coherent Fusion model is shown to outperform the Late and Mid-level Fusion methods on the PDBbind core set of 290 compounds. While the difference between the Late and Coherent Fusion methods is only 0.03 RMSE, the genetic optimization of Coherent Fusion produced several nearly identical models, which consistently performed better than Late Fusion in all evaluated metrics. Importantly, the Coherent Fusion model converged to a model structure using an automated process that exceeds the performance of the hand-crafted fusion architecture used in the original Mid-level Fusion [27]. Additionally, we provide a comparison to two other deep learning approaches (KDeep [26] and Pafnucy [56]) to

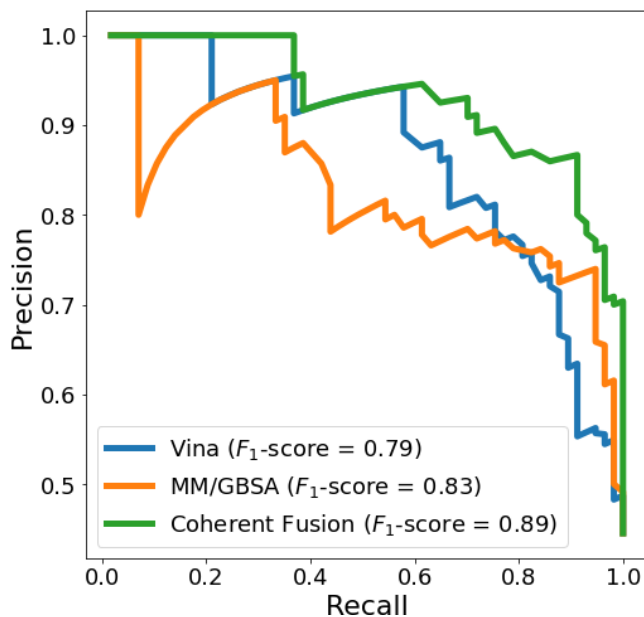


Figure 2: Binary classification of 128 docked complexes from the PDBbind core set, where the positive, “stronger” binder class represents 57 compounds with experimental pK_i or $pK_d > 8$ and the negative, “weaker” class consists of 71 compounds with pK_i or $pK_d < 6$.

view Coherent Fusion’s performance in a wider scope. While the PDBbind core set is a standard benchmark for machine learning methods [26, 27, 56], high-throughput virtual screening overwhelmingly relies to docked poses of compounds for drug discovery. To follow suit, we leveraged ConveyorLC [70] to compare Coherent Fusion against physics-based scoring functions in the docking space. The noisier docking data also provides insight into whether the machine learning model was over-trained and how robust it is to noise when scoring more realistic data.

197 compounds from the PDBbind core set were successfully evaluated by ConveyorLC with the physics-based Autodock Vina algorithm [58] and MM/GBSA methods for comparison with Coherent Fusion. Each compound was then filtered by *RMSD*, where each of the 197 compounds were checked for a pose with *RMSD* < 1Å such that a correct pose was found and sufficiently similar to the crystal structure from PDBbind. Using the binding affinity values from PDBbind as ground-truth for the docking poses, Vina achieved a Pearson correlation coefficient of .579, MM/GBSA scored .591, and Coherent Fusion reached .745. To further examine the performance difference between the three methods, binding affinity prediction can be cast as a binary classification problem [27]. Figure 2 shows the results on a subset of the 197 core set compounds. Positive and negative classes were created from 57 stronger binders and 71 weaker binders, respectively. Because the set of strong vs. weak docked poses is small (128 total), we elected to compare the different methods using a Precision-Recall Curves and F_1 -scores which give a much more direct picture of how each model is performing than a ROC curve provides.

The nominally small scoring improvements such as, Coherent vs. Late Fusion (0.02 MAE) or Coherent Fusion vs. MM/GBSA (0.06 F_1 -score), have an amplified value in large-scale screening. For example, consider a hypothetical virtual screen evaluating 1 million compounds to subsequently purchase 100 compounds (0.01%) for experimentation. If the results from classifying the PDBBind core set docking complexes translated to the top 100 of 1 million candidates with a factor of 10 decrease in precision, a virtual screen using MM/GBSA would produce 7 true positives and Coherent Fusion-based screen would produce 9 true positives in $< 1/100^{\text{th}}$ of the computational time (Table 7). While significant caveats apply, any additional true positive binders are valuable, as in practice, inhibitory compounds are hard to come by and must meet additional pharmacokinetic and safety requirements. With this analysis, we considered the Coherent Fusion model valuable and validated for use in screening for SARS-CoV-2 inhibitors.

4 HIGH-THROUGHPUT SCREENING

Our SARS-CoV-2 effort screened over 500 million compounds against each of the 4 M^{PTO} and spike targets, drawing compounds from four public virtual compound libraries [33]. The ZINC database [57] was used to create a set of “world-approved 2018” drugs from a list of FDA-approved and “world-not-FDA” approved drugs. An additional 1.5 million compounds were selected from ChEMBL [15]; 18 million compounds were drawn from eMolecules [11], and the remaining compounds came from Enamine’s list of drug-like compounds estimated to be synthetically feasible [12].

SMILES strings [63] from the eMolecules and Enamine databases and the 2D SDF structures from the ZINC and ChEMBL libraries were downloaded, respectively. Both forms of input were imported to the MOE program [18] to remove salts and metal-containing ligands, then the protonation states of compounds were set to the dominant form at pH 7. 3D structures of compounds were then generated and energetically minimized. The selected MOE descriptors were calculated and the final structures were exported from MOE as SDF files. These structures were then processed by the ligand preparation in the AMBER tool-chain utilizing antechamber and the GAFF force field [51]. The 3D SDF files were also converted to PDBQT format for the docking calculations by the Open Babel toolbox [44]. In sum, over 5 billion docking poses were generated and evaluated.

4.1 Physics-based Screening Pipeline

As with our comparison between Coherent Fusion and physics-based binding affinity scoring functions, we leveraged the existing ConveyorLC [70] tool chain to search for candidate inhibitors of the two binding sites from spike and active sites for two SARS-CoV-2 protease crystal structures. ConveyorLC is made up of four parallelized programs each designed to handle a specific task in the molecular docking and re-scoring processes. ConveyorLC uses CDT1Receptor to perform protein preparation, CDT2Ligand for ligand preparation, CDT3Docking performs the molecular docking, and finally CDT4mmgsa handles MM/GBSA re-scoring. Further details on the exact execution, pre-processing, and parameters used in the docking simulations can be found in previous studies [33, 70].

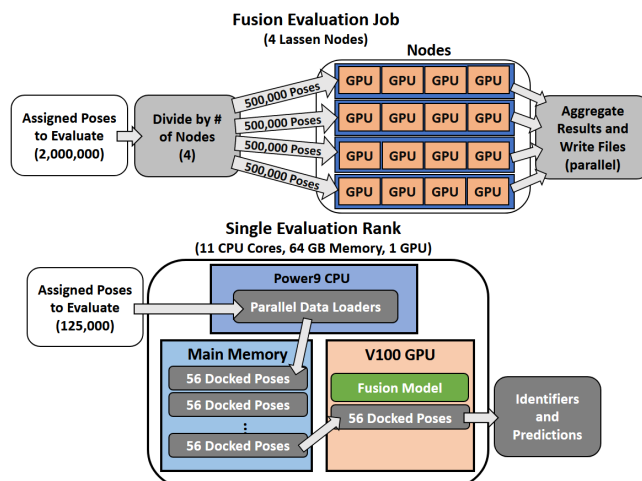


Figure 3: Structure of a single Fusion scoring job (top). A job begins with 2 million poses to score, divides them per node, then each node assigns poses to its ranks and scores. Individual ranks (bottom) take their assigned poses, begin loading batches into memory and feeding them to the GPU for inference. Finally, identifiers and predictions are collected and written in parallel.

The Vina scoring used in CDT3Docking, operates at approximately one minute per compound per CPU core. On a single Lassen node with 40 CPU cores (each core has 4 hardware threads) using 8 Monte Carlo simulations per compound, Vina is able to dock ≈ 10 docking poses per second. In contrast, a single-point MM/GBSA score takes 10 minutes per docking pose per CPU core. Because of its computational cost, MM/GBSA is often used as a re-scoring function to refine an already filtered set of compounds. [64]. Even on a Lassen node, MM/GBSA is only capable of re-scoring ≈ 0.067 poses per second.

4.2 Distributed Fusion Predictions

In order to screen millions of compounds against SARS-CoV-2, we developed a scalable architecture around the Coherent Fusion model for rapid evaluation (Figure 3). The Coherent Fusion model occupies 1.5 GB GPU memory, which fits on each 16GB NVIDIA Volta V100 GPU. The remainder of the GPU memory is used to simultaneously load 56 individual docked poses into a batch alongside each model. The 4 model instances on each node were given 12 parallel data loaders to accelerate inference. Each model in every job is assigned a subset of compounds to evaluate and its data loaders complete all file reading and pre-processing operations to prepare batches of data in an individual node’s 256 GB of memory, which are subsequently loaded onto the GPU. After evaluating a batch, the screening code unloads the compound, target, and pose identifiers along with the model’s predicted binding affinity. Once a job completes evaluation, the identifiers and predictions are gathered across MPI ranks and distributed across the individual ranks to be written in parallel to HDF5 files.

In the context of Lassen’s Job Scheduler LSF [22], we formulated Fusion evaluation jobs as many, individual 4 node processes,

each assigned to evaluate an independent set of 2 million poses, which is approximately 200,000 compounds. This format was also a response to our encountering a wide range of errors (bad metadata, node failure, broken pipe errors, etc...), which led to our pipeline being tailored for fault tolerance. With this architecture, when a job fails it has minimal impact on overall throughput (another job takes its place), the reason for failure is easier to pinpoint (log files are smaller and easier to parse), and only a small set of compounds are affected or need to be rescheduled.

Table 7: Throughput for Fusion prediction single job (2 million poses) and peak performance (125 parallel jobs)

Metric	Single Job	Peak
Avg. Startup	20 min.	"
Avg. Evaluation	280 min.	"
Avg. File Output	6.5 min.	"
Poses per sec.	108	13,594
Poses per hour	338,800	48,600,000
Compounds per hour	33,880	4,860,000

To create each 4 node job, we relied heavily on Horovod [52], which is based on MPI concepts and uses MPI for cross-node communication. With 4 GPUs on each node, each job is a 16-rank distributed process. Each rank runs a Python script and is given a specific GPU, CPU cores, and memory allocation to execute the evaluation. At the beginning of a job, we simply divide the set of compounds assigned to the job by the number of ranks and assign each rank the subset with its index. When evaluation completes, the ranks use *allgather* to compile the results and subsequently write out HDF5 files. File output was identified as a bottleneck to the evaluation process early on, which was mitigated by assigning each rank compounds to be written in the same files and directories. The output file format was designed to mirror the output format of ConveyorLC's CDT3Docking process [70] for interpretation with existing tools and further evaluation of pharmacokinetic and safety properties [33].

Table 7 includes a breakdown of how time is spent in an individual job, where the initial 20 minutes are consumed by loading HPC modules, an Anaconda [23] environment, initializing the Horovod ranks, loading an instance of the Fusion model onto each GPU, and pre-loading the initial batches of data for evaluation. The bulk of a job is, as expected, spent in an evaluation period, where batches are loaded, evaluated, and predictions are stored in parallel across all ranks. Finally, once the ranks gather together. The file writing process begins and completes about 6.5 minutes later yielding an average total run-time of approximately 5.1 hours.

With jobs designed for scalability, we regularly ran more than 10 at a time during the SARS-CoV-2 screening effort. However, at set times the majority of Lassen nodes were made available to accelerate evaluation. The peak of which was an allotment of 500 nodes for Fusion screening. The impact of a large number of nodes is clearly seen in Table 7, where throughput was increased more than 100 times. Ultimately, during several hours of evaluation at scale, the Coherent Fusion model used more than 14,010 TFLOPS of Lassen's compute power to screen nearly 5 million compounds

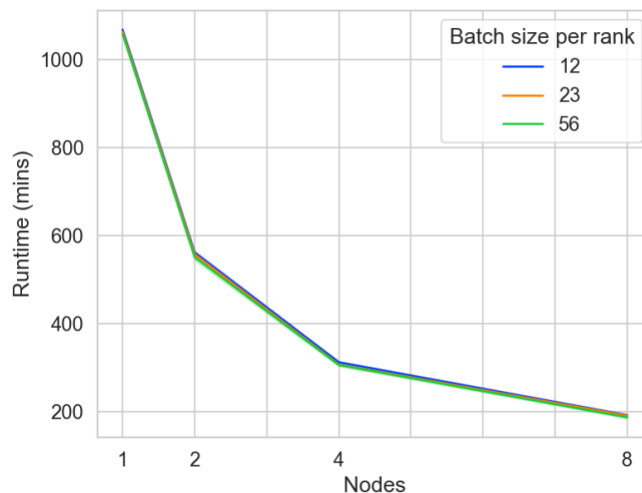


Figure 4: Strong scaling of a single Coherent Fusion job scoring 2 million poses at different batch sizes per rank (12, 23, 56) and number of nodes (1, 2, 4, 8).

per hour. The throughput advantage of Fusion is clear and compared with Vina and MM/GBSA, the Fusion model scoring code provides a 2.7x and 403x speed increase, respectively.

4.3 Single Job Scalability & Bottlenecks

While fault tolerance encouraged the use of many individual jobs each with a small number of nodes, the optimal scale of an individual job was not immediately obvious. Several factors contributed to selecting 4 nodes as optimal including: the 12 hour job run-time limit on Lassen, the startup overhead of a job, the benefit of additional nodes, and the stability of the Python libraries used.

The first parameter explored in development of the Fusion screening capability was the number of poses to evaluate per job. We found it possible to complete up to 5 million scored poses on 4 nodes under the 12 hour Lassen time limit. However, the prevalence of unpredictable errors in the docking data, featurization steps, and inter-node/rank communication, led us to instead assign 2 million poses to each job. While 2 million poses is less efficient (*i.e.*, startup overhead is a larger percentage of each job's run-time), in practice this decision led to less wasted computational time, as the Fusion scoring code does not write results until it finishes scoring all poses. In future work, efficiency will be improved by creating a separate, parallel process per rank to write results as they are computed, but due to the urgent need for SARS-CoV-2 predictions we elected to mitigate unpredictable errors by narrowing the size of each Fusion job to balance fault-tolerance and computational efficiency.

The next two parameters explored were the batch size per rank and number of nodes per job. The batch size per rank effects the number of times data is transferred from CPU->GPU and predictions from GPU->CPU. For the M^{Pro} and spike target sites, we found up to 56 poses (each with voxelized and connectivity representations) could consistently fit on the NVIDIA V100 GPU. Figure 4 displays the effect three different batch sizes had on a single job.

In practice, the performance difference between each batch size was small with a batch size of 56 yielding a ≈ 10 minute run-time advantage over batch size 12.

The Fusion scoring code under-utilized the Lassen GPUs, which led to consistent and relatively small offsets in run-time by batch size. This is due to the computational cost of pre-processing (*e.g.*, file reading and data featurization), which is the most significant bottleneck in the evaluation process. Despite using 12 parallel data loaders per rank, the GPU is intermittently waiting to evaluate more poses. In future work, further optimization of the parallel data loaders will increase GPU utilization and improve throughput, which increases the value of Fusion for large-scale virtual screening.

To determine the optimal number of nodes per job and performance benefit of additional nodes, we evaluated the run-times of Fusion jobs using 1, 2, 4, and 8 nodes per job. Figure 4 also displays the performance of each number of nodes where the same set of 10 jobs was run for every point evaluated. The variance between sets is also plotted surrounding each line, but was found to be small (< 5 minutes) and as a result is not clearly visible.

A significant factor in choosing the number of nodes for each job was stability of the Python libraries used by the Fusion model. Inter-node and inter-rank communication errors were increasingly prevalent as the number of nodes/ranks in a job increased and the percentage of failed jobs for each number of nodes was $\approx 2\%$ for 1 and 2 nodes, $\approx 3\%$ for 4 nodes, and $\approx 20\%$ of jobs failed when using 8 nodes. The instability was caused by the specific combination of Horovod [52] and PyTorch [46] used on the POWER9 architecture, which has since been updated. However, a 20% job failure rate eliminated 8 nodes per job as a candidate configuration.

The results of our scalability experiments led to a final selection of 2 million poses per job using 4 nodes. This was the result of several different factors and adjustments including the observation that when using 500 Lassen nodes at the same time, the LSF scheduler encountered problems simultaneously running 250 2-node jobs, which was solved by using 125 4-node Fusion jobs instead.

5 SARS-COV-2 RESULTS

The high-throughput virtual drug screening pipeline described in Section 4 produced several computational results for each compound screened. The Fusion model's binding-affinity prediction was one of the three energy calculations (Vina, MM/GBSA, Fusion), which were used as a component of a hand-tailored cost function designed to filter which compounds to purchase for experimental evaluation and which were less likely to be successful. Full details of the ranking and reasoning may be found in [33] and computational predictions are made available at <https://url-excluded> [38]. Virtual screening output on the computational side fed directly into an experimental process to physically interrogate candidate molecules.

5.1 Experimental Validation

Experimental testing of the candidate binders which were screened and purchased to target M^{pro} used a fluorescence resonance energy transfer (FRET) based activity assay or a SDS-PAGE gel protein cleavage assay. After assay optimization, additional screens were run in order to down select compounds. For example, compounds from the ZINC database [57] were down selected to an additional

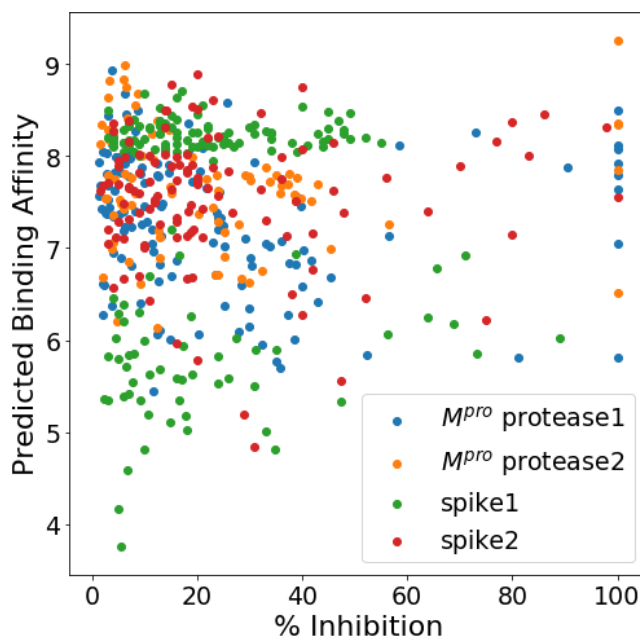


Figure 5: Coherent Fusion predicted binding affinity vs. experimental percentage of inhibition at 100 μ M for 130 compounds against M^{pro} protease1 (blue) and 81 compounds against M^{pro} protease2 (orange). The spike assays were evaluated at 10 μ M and include 151 compounds against spike1 (green) and 113 compounds against spike2 (red). Compounds which exhibited $\leq 1\%$ inhibition (no experimental binding activity) are excluded.

testing of 19 compounds, which yielded 4 candidates inhibiting the activity of M^{pro} at 100 micro-Molar (μ M) concentrations. The four identified compounds include: candestartan cilexetil, FAD disodium, tigecycline, and tetracycline [33].

On the other hand, compounds predicted to inhibit the SARS-CoV-2 spike protein were screened by both a pseudo-typed virus assay and a biolayer interferometry competitive assay (BLI). Here the candidate compounds are being evaluated for their ability to inhibit ACE2-spike binding and in parallel, the spike binding candidates were screened using a cell-based infection assay at 10 μ M. Further details of the experimental design, assays, results, and discussion can be found in [33].

5.2 Connecting Predictions and Experimental Results

Given the ground-truth experimental values for physically tested compounds, we're enabled to retrospectively evaluate the accuracy of each computational approach which generated a prediction. Some of the obvious questions to ask are: "Which method was most correlated with the experimental results?", "Are the most accurate scoring functions the same for all four M^{pro} and spike targets?", and "Which of the scoring methods is most accurate for the strongest experimental inhibitors?". Each experimentally prosecuted compound can be traced back to its virtually docked poses for either

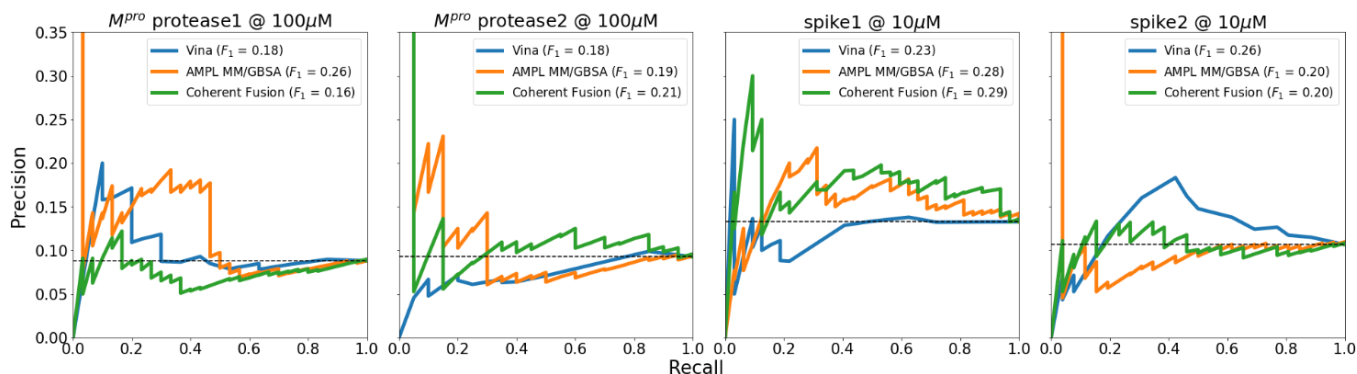


Figure 6: Precision/Recall Curves and F1-scores by SARS-CoV-2 protein target at 33% experimental inhibition. M^{Pro} protease1 (far left) shows results for 30 positive and 311 negative binders. M^{Pro} protease2 (middle left) includes 20 positive and 196 negative binders. The spike1 site (middle right) includes 32 positive and 209 negative binders. Finally, spike2 (far right) includes 26 positive and 218 negative binders. The black horizontal dashed line indicates the performance of a random classifier.

the two M^{Pro} or two spike binding targets. This means each scoring method may have predicted binding affinity values for up to 40 poses per compound (10 poses maximum per binding site). While in the computational domain we have predicted binding affinity values, the output from the M^{Pro} and spike assays is a percentage of inhibition normalized between 0 and 100%. These values are produced at a given concentration of the candidate drug (100 μ M for M^{Pro} targets and 10 μ M for spike targets), which gives context to the overall strength of a binder.

For each scoring method (Vina, MM/GBSA, Coherent Fusion) the results per compound were aggregated and represented by the strongest prediction across all poses for each binding site (maximum for Coherent Fusion, minimum for Vina and MM/GBSA). Because of the computational cost of MM/GBSA, we instead use the ATOM Modeling PipeLine (AMPL) MM/GBSA predicted MM/GBSA values, which have been shown to be highly correlated with actual MM/GBSA calculations and were trained to predict MM/GBSA scores on each specific target [42].

5.3 Analyzing Computational Predictions

Following this aggregation, the output is a single prediction per compound tied to a single percent inhibition. Each method can then be viewed as a scatter plot comparing predicted vs. actual experimental results as in Figure 5. We computed Pearson and Spearman correlation coefficients for each method across all experimentally tested M^{Pro} and spike compounds. However, most experimentally tested compounds are negatives ($\leq 1\%$ inhibition), which gives correlation coefficients near 0 for each of the three evaluated methods (table excluded for brevity). In an attempt to focus our analysis on the relative strengths and weaknesses of the different scoring methods and not the difficulty of the overall binding affinity prediction problem, we computed correlation coefficients for each method on the subset of compounds for which any experimental binding ($>1\%$) was observed.

Table 8 shows the correlations for all methods where the absolute value of the Vina and MM/GBSA scores are used. AMPL

Table 8: Correlation of predicted binding and percent inhibition on compounds with $> 1\%$ inhibition

Method	Target/Site	Pearson R	Spearman R
Vina	M^{Pro} /protease1	0.03	-0.08
AMPL MM/GBSA	M^{Pro} /protease1	0.08	0.01
Coherent Fusion	M^{Pro} /protease1	-0.06	-0.04
Vina	M^{Pro} /protease2	-0.08	-0.14
AMPL MM/GBSA	M^{Pro} /protease2	-0.05	-0.07
Coherent Fusion	M^{Pro} /protease2	0.04	0.04
Vina	spike/spike1	-0.02	0.06
AMPL MM/GBSA	spike/spike1	0.15	0.22
Coherent Fusion	spike/spike1	0.22	0.30
Vina	spike/spike2	0.13	0.27
AMPL MM/GBSA	spike/spike2	-0.02	-0.05
Coherent Fusion	spike/spike2	-0.02	-0.01

MM/GBSA gives the best correlation for the protease1 target, Coherent Fusion for the protease2 and spike1 targets, and Vina scores the spike2 binding site best. However, across the board, it is clear that even when limiting the analysis to $>1\%$ inhibition, the correlations for each method remain low and the interpretation of near-zero correlation coefficients is unavailing.

While removing all the non-inhibitors gives a glimpse into which methods are most correlated with the SARS-CoV-2 binders, it also removes the context of those predictions. That is, the overall prediction strength for each method is somewhat obscured as the range of each method's prediction values is limited to its minimum and maximum prediction in the smaller set of SARS-CoV-2 binders. With this in mind, we sought to answer the question of which scoring methods were most accurate for the strongest experimental inhibitors by including non-binding compounds and casting the prediction problem as a binary classification of compounds with $>33\%$ inhibition (positive class) and compounds with $\leq 33\%$ inhibition (negative class). A threshold of 33% was chosen to avoid severe class

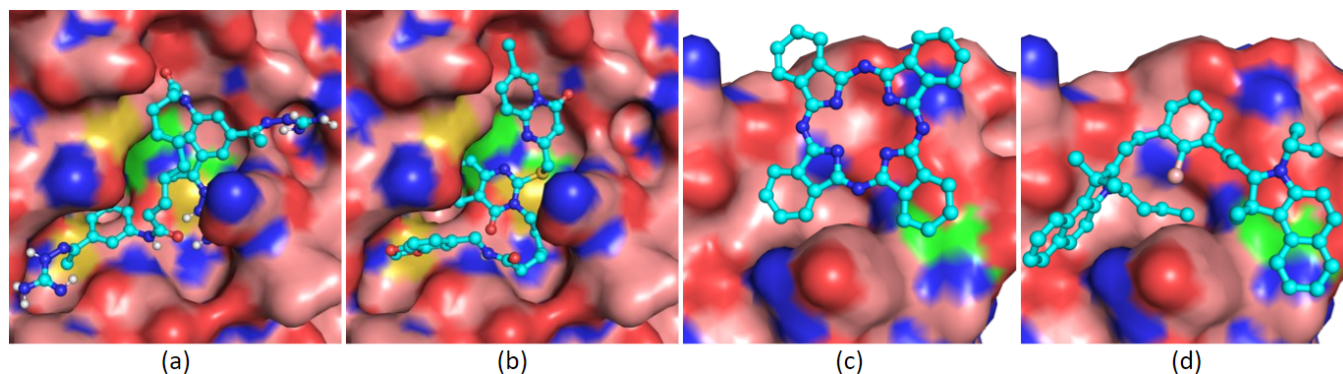


Figure 7: Four compounds from eMolecules [11] in complex with the M^{Pro}/protease1 (a, b) and spike/spike1 targets (c, d), where residues His41 and Cys145 in M^{Pro} and residue 501 in the spike RBD protein are green. Panels (a) and (b) show Compound IDs 76051337 and 24424612, respectively, which both had 100% inhibition at 100 μ M in the M^{Pro} assay. Panels (c) and (d) show Compound IDs 18594404 and 313102183, respectively, which reached 100% and 98% inhibition at 10 μ M in the spike assay.

imbalances caused by higher thresholds. This problem formulation is similar to that in Figure 2 where we set a threshold to separate stronger binders from weaker binders. This approach is applicable in practice, as virtual screens eventually down select to a small set of candidates for final analysis, purchase, and experimental testing.

Figure 6 displays Precision/Recall Curves and F_1 -scores using a threshold of 33% to separate the positive and negative binding examples for each target. The y-axis of each P/R curve is limited to 0.35 to observe different model behaviors. Although the curves for each target appear different from the P/R plots in Figure 2, they provide information about how each model performs in practice on a noisier experimental screen, where the cutoff separating active from inactive molecules is less clear.

While the P/R Curves in Figure 6 give low F_1 -scores and precisions, two important observations arise. First, for each of the four binding sites, we computed Cohen’s Kappa statistic (κ) to compare each model with a random classifier [7, 32]. Equation 2 gives the equation for computing κ ;

$$\kappa = \frac{\rho_o - \rho_e}{1 - \rho_e} \quad (2)$$

where ρ_o is the observed accuracy and ρ_e is the expected accuracy. A random classifier achieves a $\kappa = 0$ by guessing according to the frequency of the positive and negative classes. Every model on each target was measured to achieve a $\kappa > 0$ (with the exception of Vina on spike1), indicating the models are consistently better than random across the targets. In the context of P/R Curves, a random classifier’s performance is expressed as a horizontal line (Figure 6) with constant expected precision equal to the proportion of positive examples in the data set.

In practice, the three binding affinity prediction models served as input to a hand-crafted cost function [33], to select compounds for experimentation. The outcome of which produced 9 distinct compounds inhibiting 100% of M^{Pro} activity and 108 total compounds with $\geq 33\%$ inhibition. Using the 33% inhibition threshold, 108 of 1042 (10.4%) experimentally tested compounds inhibit activity, indicating the models have significant predictive power. In that, while 500+ million candidates were evaluated, only a fraction ($2.1e^{-6}\%$)

of the best candidates were tested experimentally, yet the models successfully yielded a 10.4% hit rate.

Figure 7 presents four top compounds bound to two of the four sites/configurations used in the virtual screen (M^{Pro}/protease1 and spike/spike1). The M^{Pro}/protease1 compounds reached 100% inhibition and had Coherent Fusion predicted binding affinities of 8.5 (Fig. 7a) and 8.1 (Fig. 7b). The spike RBD compounds evaluated at 10 μ M had predicted affinities and experimental inhibitions at 7.6/100% (Fig. 7c) and 8.3/98% (Fig. 7d).

For the M^{Pro} protease1 binding site, AMPL MM/GBSA is the best predictor of experimental binding. Additionally, between the two M^{Pro} binding sites, protease1 and protease2, the AMPL MM/GBSA predictions on protease1 conjointly achieve a higher F_1 -score and Pearson correlation coefficient than all other methods. In fact, for each of the four targets, not only do the best $>33\%$ F_1 -scores match with the best $>1\%$ Pearson/Spearman correlations, but the relative differences across target binding sites are also in agreement.

Coherent Fusion reaches the maximal F_1 and correlation coefficients for the M^{Pro} protease2 and spike1 binding sites. Across all model / binding site combinations, Coherent Fusion’s performance on the spike1 protein is the top performer, though followed closely by AMPL MM/GBSA on the same target site. This is unexpected, as the protease targets are much larger and nominally thought to provide better opportunity for drug-like compounds to bind.

However, target specific strengths are not unexpected [64]. Especially where the M^{Pro} sites are large protein pockets and the spike targets are much smaller. Interestingly, the maximal $>1\%$ correlations and $>33\%$ F_1 -scores across all targets favor the spike proteins. The differing concentrations between the M^{Pro} and spike experiments is important to note, as M^{Pro} binders were evaluated at 100 μ M, which is a higher drug concentration and, therefore, allows for weaker binders to exhibit higher observed percentages of inhibition.

6 CONCLUSION AND FUTURE WORK

Deep Fusion was improved by coherent backpropagation and distributed, genetic hyper-parameter optimization. The optimization automatically produced a version of the Coherent Fusion model, which was shown to exceed the performance of hand-crafted Fusion model variants and alternative machine learning methods on crystal structures from the PDBbind *core* set benchmark. In evaluating noisier docked poses of the same test set, Coherent Fusion also achieved an improvement in correlation and F_1 -score relative to the physics-based Vina and MM/GBSA methods.

We utilized Coherent Fusion to screen over 500 million compounds across four SARS-CoV-2 binding sites. This was achieved by designing a high-throughput distributed architecture for fault-tolerant scoring at scale. Using parallel data loaders and Lassen's 4 GPUs per node, the Coherent Fusion model is currently capable of screening ≈ 30 poses per second per node. Coherent Fusion was used as input to a weighted cost function across binding affinity and other scoring models in order to sub-select compounds for experimentation.

Among the nearly 1000 compounds tested experimentally, several inhibited activity of M^{pro} and spike. In analyzing which binding affinity methods were most correlated with the experimental results, we found the optimal performer to vary by target. However, two of the four targeted binding sites were best predicted by Coherent Fusion both in terms of $>1\%$ inhibiting correlation with the experimental results and $>33\%$ inhibiting binary classification. Although overall predictive power of all of the scoring functions is limited, any opportunity for computational enrichment of strong binders is needed to alleviate the otherwise prohibitive time and cost of the physical experimental screens.

In future work, we aim to use our baseline Coherent Fusion model from this work to fine tune and predict for specific protein target types and binding sites. We believe introducing target specificity to the models and thereby reducing the scope of the binding affinity prediction problem will increase the value of relative differences in the model's binding affinity predictions. Our high-throughput architecture can also be accelerated, as the GPUs on each Lassen node were observed to be under-utilized. Increased numbers of parallel data loaders were observed to decrease the overall stability of each individual evaluation job indicating some refactoring is necessary.

ACKNOWLEDGMENTS

The authors gratefully acknowledge Professor Xinquan Wang (Tsinghua University) for providing early access to his crystal structure of the ACE2-RBD complex. The authors gratefully acknowledge extensive computer time provided by Livermore Computing. Part of this research was supported by the DOE Office of Science through the National Virtual Biotechnology Laboratory, a consortium of DOE national laboratories focused on response to COVID-19, with funding provided by the Coronavirus CARES Act. The authors thank Lawrence Livermore National Laboratory for funding Laboratory Directed Research and Development projects 20-ERD-065 and 20-ERD-062. Part of this research was also supported by the American Heart Association under CRADA TC02274-4 and the National Nuclear Security Administration through the Accelerating

Therapeutics for Opportunities in Medicine (ATOM) Consortium under CRADA TC02349. This work was funded in part by DTRA under award HDTRA1036045. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525. All work performed at Lawrence Livermore National Laboratory is performed under the auspices of the U.S. Department of Energy under Contract DE-AC52-07NA27344. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

REFERENCES

- [1] Qurrat Ul Ain, Antoniya Aleksandrova, Florian D Roessler, and Pedro J Ballester. 2015. Machine-learning scoring functions to improve structure-based binding affinity prediction and virtual screening. *Wiley Interdisciplinary Reviews: Computational Molecular Science* 5, 6 (2015), 405–424.
- [2] Pedro Ballester and John Mitchell. 2010. A machine learning approach to predicting protein-ligand binding affinity with applications to molecular docking. *Bioinformatics (Oxford, England)* 26 (03 2010), 1169–75. <https://doi.org/10.1093/bioinformatics/btq112>
- [3] Sandrine Belouzard, Jean K Millet, Beth N Licitra, and Gary R Whittaker. 2012. Mechanisms of coronavirus cell entry mediated by the viral spike protein. *Viruses* 4, 6 (2012), 1011–1033.
- [4] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. 2011. Algorithms for hyper-parameter optimization. In *25th annual conference on neural information processing systems (NIPS 2011)*, Vol. 24. Neural Information Processing Systems Foundation.
- [5] PA Castillo, J Carpio, JJ Merelo, Alberto Prieto, V Rivas, and Gustavo Romero. 2000. Evolving multilayer perceptrons. *Neural Processing Letters* 12, 2 (2000), 115–128.
- [6] Tiejun Cheng, Qingliang Li, Zhigang Zhou, Yanli Wang, and Stephen H Bryant. 2012. Structure-based virtual screening for drug discovery: a problem-centric review. *The AAPS journal* 14, 1 (2012), 133–141.
- [7] Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement* 20, 1 (1960), 37–46.
- [8] Kurt LM Drew, Hakim Baiman, Prashanna Khwaounjoo, Bo Yu, and Jóhannes Reynisson. 2012. Size estimation of chemical space: how big is it? *Journal of Pharmacy and Pharmacology* 64, 4 (2012), 490–495.
- [9] John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research* 12, 7 (2011).
- [10] Sally R Ellingson, Brian Davis, and Jonathan Allen. 2020. Machine learning and ligand binding predictions: a review of data, methods, and obstacles. *Biochimica et Biophysica Acta (BBA)-General Subjects* 1864, 6 (2020), 129545.
- [11] eMolecules. 2021. eMolecules. <https://www.emolecules.com/>
- [12] Enamine. 2021. Enamine. <https://enamine.net/>
- [13] Evan N Feinberg, Debnil Sur, Zhenqin Wu, Brooke E Husic, Huanghao Mai, Yang Li, Saisai Sun, Jianyi Yang, Bharath Ramsundar, and Vijay S Pande. 2018. PotentialNet for molecular property prediction. *ACS central science* 4, 11 (2018), 1520–1530.
- [14] Matthias Fey and Jan Eric Lenssen. 2019. Fast graph representation learning with PyTorch Geometric. *arXiv preprint arXiv:1903.02428* (2019).
- [15] Anna Gaulton, Louisa J Bellis, A Patricia Bento, Jon Chambers, Mark Davies, Anne Hersey, Yvonne Light, Shaun McGlinchey, David Michalovich, Bissan Al-Lazikani, et al. 2012. ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic acids research* 40, D1 (2012), D1100–D1107.
- [16] Joseph Gomes, Bharath Ramsundar, Evan N Feinberg, and Vijay S Pande. 2017. Atomic convolutional networks for predicting protein-ligand binding affinity. *arXiv preprint arXiv:1703.10603* (2017).
- [17] Alex Graves. 2013. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850* (2013).
- [18] Chemical Computing Group ULC. 2021. Molecular Operating Environment. https://www.chemcomp.com/MOE-Molecular_Operating_Environment.htm
- [19] Tingjun Hou, Junmei Wang, Youyong Li, and Wei Wang. 2011. Assessing the performance of the molecular mechanics/Poisson Boltzmann surface area and molecular mechanics/generalized Born surface area methods. II. The accuracy of ranking poses generated from docking. *Journal of computational chemistry* 32, 5 (2011), 866–877.

- [20] Michael Husken, Jens E Gayko, and Bernhard Sendhoff. 2000. Optimization for problem classes-neural networks that learn to learn. In *2000 IEEE Symposium on Combinations of Evolutionary Computation and Neural Networks. Proceedings of the First IEEE Symposium on Combinations of Evolutionary Computation and Neural Networks (Cat. No. 00)*. IEEE, 98–109.
- [21] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. 2011. Sequential model-based optimization for general algorithm configuration. In *International conference on learning and intelligent optimization*. Springer, 507–523.
- [22] IBM. 2021. IBM Spectrum LSF V10.1 documentation. https://www.ibm.com/support/knowledgecenter/en/SSWRJV_10.1.0/lfsf_welcome/lfsf_welcome.html
- [23] Anaconda Inc. 2021. Anaconda Software Distribution. <https://docs.anaconda.com/>
- [24] Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M Czarnecki, Jeff Donahue, Ali Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, et al. 2017. Population based training of neural networks. *arXiv preprint arXiv:1711.09846* (2017).
- [25] Araz Jakalian, Bruce L Bush, David B Jack, and Christopher I Bayly. 2000. Fast, efficient generation of high-quality atomic charges. AM1-BCC model: I. Method. *Journal of computational chemistry* 21, 2 (2000), 132–146.
- [26] José Jiménez, Miha Skalic, Gerard Martínez-Rosell, and Gianni De Fabritiis. 2018. K deep: protein–ligand absolute binding affinity prediction via 3d-convolutional neural networks. *Journal of chemical information and modeling* 58, 2 (2018), 287–296.
- [27] Derek Jones, Hyojin Kim, Xiaohua Zhang, Adam Zemla, Garrett Stevenson, W. F. Drew Bennett, Daniel Kirshner, Sergio E. Wong, Felice C. Lightstone, and Jonathan E. Allen. 0. Improved Protein–Ligand Binding Affinity Prediction with Structure-Based Deep Fusion Inference. *Journal of Chemical Information and Modeling* 0, 0 (0), null. <https://doi.org/10.1021/acs.jcim.0c01306> arXiv:<https://doi.org/10.1021/acs.jcim.0c01306> PMID: 33754707.
- [28] William L Jorgensen. 2004. The many roles of computation in drug discovery. *Science* 303, 5665 (2004), 1813–1818.
- [29] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [30] Douglas B Kitchen, Hélène Decornez, John R Furr, and Jürgen Bajorath. 2004. Docking and scoring in virtual screening for drug discovery: methods and applications. *Nature reviews Drug discovery* 3, 11 (2004), 935–949.
- [31] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. 2017. Self-normalizing neural networks. *arXiv preprint arXiv:1706.02515* (2017).
- [32] J Richard Landis and Gary G Koch. 1977. The measurement of observer agreement for categorical data. *biometrics* (1977), 159–174.
- [33] Edmond Y. Lau, Oscar A. Negrete, W. F. Drew Bennett, Brian J. Bennon, Monica Borucki, Feliza Bourguet, Aidan Epstein, Magdalena Franco, Brooke Harmon, Stewart He, Derek Jones, Hyojin Kim, Daniel Kirshner, Victoria Lao, Jacky Lo, Kevin McLoughlin, Richard Mosesso, Deepa K. Muruges, Edwin A. Saada, Brent Segelke, Maxwell Stefan, Garrett A. Stevenson, Marisa W. Torres, Dina Weillhammer, Sergio Wong, Yue Yang, Adam Zemla, Xiaohua Zhang, Fangqiang Zhu, Jonathan E. Allen, and Felice C. Lightstone. 2021. Discovery of Small-molecule Inhibitors of SARS-CoV-2 Proteins Using a Computational and Experimental Pipeline. *Frontiers in Molecular Biosciences, In Review* (2021).
- [34] Ang Li, Ola Spyra, Sagi Perel, Valentin Dalibard, Max Jaderberg, Chenjie Gu, David Budden, Tim Harley, and Pramod Gupta. 2019. A generalized framework for population based training. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1791–1799.
- [35] Huibin Li, Jian Sun, Zongben Xu, and Liming Chen. 2017. Multimodal 2D+ 3D facial expression recognition with deep fusion convolutional neural network. *IEEE Transactions on Multimedia* 19, 12 (2017), 2816–2831.
- [36] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. 2015. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493* (2015).
- [37] Jaechang Lim, Seongok Ryu, Kyubyong Park, Yo Joong Choe, Jiyeon Ham, and Woo Youn Kim. 2019. Predicting drug–target interaction using a novel graph neural network with 3D structure-embedded graph representation. *Journal of chemical information and modeling* 59, 9 (2019), 3981–3988.
- [38] LLNL. 2021. Lawrence Livermore National Laboratory Covid-19 Therapeutic Design Data Portal. <https://doi.org/10.11578/1608139>
- [39] LLNL. 2021. LLNL’s Lassen supercomputer leaps to No. 10 on TOP500 list, Sierra remains No. 2n. <https://www.llnl.gov/news/llnl-s-lassen-supercomputer-leaps-no-10-top500-list-sierra-remains-no-2>. Accessed: 2020-03-06.
- [40] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- [41] James A Maier, Carmenza Martinez, Koushik Kasavajhala, Lauren Wickstrom, Kevin E Hauser, and Carlos Simmerling. 2015. ff14SB: improving the accuracy of protein side chain and backbone parameters from ff99SB. *Journal of chemical theory and computation* 11, 8 (2015), 3696–3713.
- [42] Kevin McLoughlin. 2021. ATOM Modeling PipeLine (AMPL) MM/GBSA predicted MM/GBSA values. Personal Communication.
- [43] Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang, Melih Elibol, Zongheng Yang, William Paul, Michael I Jordan, et al. 2018. Ray: A distributed framework for emerging {AI} applications. 561–577.
- [44] Noel M O’Boyle, Michael Banck, Craig A James, Chris Morley, Tim Vandermeersch, and Geoffrey R Hutchison. 2011. Open Babel: An open chemical toolbox. *Journal of cheminformatics* 3, 1 (2011), 1–14.
- [45] Jack Parker-Holder, Vu Nguyen, and Stephen J Roberts. 2020. Provably efficient online hyperparameter optimization with population-based bandits. *Advances in Neural Information Processing Systems* 33 (2020).
- [46] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703* (2019).
- [47] Eric F Pettersen, Thomas D Goddard, Conrad C Huang, Gregory S Couch, Daniel M Greenblatt, Elaine C Meng, and Thomas E Ferrin. 2004. UCSF Chimera—A visualization system for exploratory research and analysis. *Journal of computational chemistry* 25, 13 (2004), 1605–1612.
- [48] Matthew Ragoza, Joshua Hochuli, Elisa Idrobo, Jocelyn Sunseri, and David Ryan Koes. 2017. Protein–ligand scoring with convolutional neural networks. *Journal of chemical information and modeling* 57, 4 (2017), 942–957.
- [49] Jean-Louis Reymond, Ruud Van Deursen, Lorenz C Blum, and Lars Ruddigkeit. 2010. Chemical space as a source for new drugs. *MedChemComm* 1, 1 (2010), 30–38.
- [50] Alina Roitberg, Tim Pollert, Monica Haurilet, Manuel Martin, and Rainer Stiefelhagen. 2019. Analysis of deep fusion strategies for multi-modal gesture recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 0–0.
- [51] Romelia Salomon-Ferrer, David A Case, and Ross C Walker. 2013. An overview of the Amber biomolecular simulation package. *Wiley Interdisciplinary Reviews: Computational Molecular Science* 3, 2 (2013), 198–210.
- [52] Alexander Sergeev and Mike Del Balso. 2018. Horovod: fast and easy distributed deep learning in TensorFlow. *arXiv preprint arXiv:1802.05799* (2018).
- [53] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. 2012. Practical bayesian optimization of machine learning algorithms. *arXiv preprint arXiv:1206.2944* (2012).
- [54] Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram, Mostofa Patwary, Mr Prabhat, and Ryan Adams. 2015. Scalable bayesian optimization using deep neural networks. In *International conference on machine learning*. PMLR, 2171–2180.
- [55] Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. 2009. Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995* (2009).
- [56] Marta M Stepniewska-Dziubinska, Piotr Zielenkiewicz, and Pawel Siedlecki. 2017. Pafnucy—A deep neural network for structure-based drug discovery. *stat* 1050 (2017), 19.
- [57] Teague Sterling and John J. Irwin. 2015. ZINC 15 — A Ligand Discovery for Everyone. *Journal of Chemical Information and Modeling* 55, 11 (2015), 2324–2337. <https://doi.org/10.1021/acs.jcim.5b00559> PMID: 26479676.
- [58] Oleg Trott and Arthur J Olson. 2010. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of computational chemistry* 31, 2 (2010), 455–461.
- [59] Sven Ullrich and Christoph Nitsche. 2020. The SARS-CoV-2 main protease as drug target. *Bioorganic & Medicinal Chemistry Letters* (2020), 127377.
- [60] Jörg Wagner, Volker Fischer, Michael Herman, and Sven Behnke. 2016. Multispectral Pedestrian Detection using Deep Fusion Convolutional Neural Networks.. In *ESANN*, Vol. 587. 509–514.
- [61] Renxiao Wang, Xueliang Fang, Yipin Lu, and Shaomeng Wang. 2004. The PDBbind Database: A Collection of Binding Affinities for Protein–Ligand Complexes with Known Three-Dimensional Structures. *Journal of Medicinal Chemistry* 47, 12 (2004), 2977–2980. <https://doi.org/10.1021/jm030580l> arXiv:<https://doi.org/10.1021/jm030580l> PMID: 15163179.
- [62] Renxiao Wang, Xueliang Fang, Yipin Lu, Chao-Yie Yang, and Shaomeng Wang. 2005. The PDBbind database: methodologies and updates. *Journal of medicinal chemistry* 48, 12 (2005), 4111–4119.
- [63] David Weininger. 1988. SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *Journal of chemical information and computer sciences* 28, 1 (1988), 31–36.
- [64] Sergio Wong. 2021. Large scale evaluation of MM-GB/SA rescoring on the PDBbind 2019 refined dataset. Personal Communication.
- [65] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. 2015. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853* (2015).
- [66] Fei Yang, Jian Yang, Zhong Jin, and Huan Wang. 2018. A fusion model for road detection based on deep learning and fully connected CRF. In *2018 13th Annual Conference on System of Systems Engineering (SoSE)*. IEEE, 29–36.
- [67] Haiping Zhang, Linbu Liao, Konda Mani Saravanan, Peng Yin, and Yanjie Wei. 2019. DeepBindRG: a deep learning based method for estimating effective protein–ligand affinity. *PeerJ* 7 (2019), e7362.

- [68] Xiaohua Zhang, Horacio Perez-Sanchez, and Felice C Lightstone. 2017. A comprehensive docking and MM/GBSA rescoring study of ligand recognition upon binding antithrombin. *Current topics in medicinal chemistry* 17, 14 (2017), 1631–1639.
- [69] Xiaohua Zhang, Sergio E Wong, and Felice C Lightstone. 2013. Message passing interface and multithreading hybrid for parallel molecular docking of large databases on petascale high performance computing machines. *Journal of computational chemistry* 34, 11 (2013), 915–927.
- [70] Xiaohua Zhang, Sergio E Wong, and Felice C Lightstone. 2014. Toward fully automated high performance computing drug discovery: a massively parallel virtual screening pipeline for docking and molecular mechanics/generalized Born surface area rescoring to improve enrichment.
- [71] Jingbo Zhou, Shuangli Li, Liang Huang, Haoyi Xiong, Fan Wang, Tong Xu, Hui Xiong, and Dejing Dou. 2020. Distance-aware Molecule Graph Attention Network for Drug-Target Binding Affinity Prediction. *arXiv preprint arXiv:2012.09624* (2020).
- [72] Fangqiang Zhu, Xiaohua Zhang, Jonathan E. Allen, Derek Jones, and Felice C. Lightstone. 2020. Binding Affinity Prediction by Pairwise Function Based on Neural Network. *Journal of Chemical Information and Modeling* 60, 6 (2020), 2766–2772. <https://doi.org/10.1021/acs.jcim.0c00026> PMID: 32338892.