1
2

# Machine learning based detection of genetic and drug class variant impact on functionally conserved protein binding dynamics

3

4    Gregory A. Babbitt[1]*, Ernest P. Fokoue[3], Joshua R. Evans[1], Kyle I. Diller[1,2], Lily E. Adams[1]

5

6    [1]Thomas H. Gosnell School of Life Sciences, Rochester Institute of Technology, Rochester NY, USA
7    14623
8    [2]Golisano College for Computing and Information Science, Rochester NY, USA 14623
9    [3]Applied Statistics, Rochester Institute of Technology, Rochester NY, USA 14623
10
11    *email address for primary correspondence: gabsbi@rit.edu

12

13    **Contact Information** –Gregory A. Babbitt (gabsbi@rit.edu).

14

15

16

17

18

19

20

21

22

23

24

25

26

## Abstract

The application of statistical methods to comparatively framed questions about protein dynamics can potentially enable investigations of biomolecular function beyond the current sequence and structure based methods in comparative genomics. However, addressing this problem requires proper statistical inferences obtained from large ensembles of individual molecular dynamic (MD) simulations that represent the comparative functional states of a given protein under investigation. Meaningful interpretation of such large and temporally rich forms of data poses serious challenges to users of MD. Here, we announce the release of DROIDS v3.0, an Amber18/Chimera-based software package for comparative protein dynamics, incorporating many new features including maxDemon v1.0, a multi-method machine learning application that trains on large ensembles of MD comparisons generated by DROIDS and deploys learned classifications of dynamic states to newly generated protein simulations. Up to seven different machine learners can be deployed on the dynamics of each amino acid, and local canonical correlations in learning patterns generated from self-similar MD runs are used to identify regions of functionally conserved protein dynamics. The subsequent impacts of genetic and drug class variants on conserved dynamics can also be analyzed by deploying the trained classifiers on new MD runs. Still and moving images of comparative dynamics allow users to see both when and where a protein dynamic simulation displays a specific functional state defined by the functional protein comparison. Our software allows quantitative visualization of complex changes in functional dynamics caused by temperature changes, chemical mutations, or binding interactions with nucleic acids or a variety of small molecules. Here, we demonstrate the utility of DROIDS 3.0 + maxDemon in four case studies of the impact of genetic and drug class variation on functionally conserved protein dynamics in ubiquitin, TATA binding protein and Hsp90.

## Introduction

62

63 The physicist Richard Feynman is said to have once famously quipped, 'all biology is ultimately due to
64 the wiggling and jiggling of atoms'. Stated with more precision, Feynman's conjecture would imply that
65 all biological function at the molecular level can actually be understood by analyzing rapid atom motions
66 in biomolecular structures as they alter their functional state(s). Many decades later, functional shifts in
67 molecular dynamics are now being illuminated by structural and computational biology. Examples of
68 such functional dynamics include the destabilization of inter-residue contacts, in both disease
69 malfunction and signal activation during phosphorylation, as well as stabilization of inter-residue
70 contacts during protein folding, the formation of larger complexes, and various other binding
71 interactions to small molecules. And while the functional role of rapid vibrations revealed by short term
72 dynamic simulations has been debated in the past, more recent empirical and computational studies
73 have clearly demonstrated that differences in rapid and directed vibrations can drive longer term
74 functional conformational change (1, 2). From a broader perspective, if Feynman's conjecture is true,
75 then the specific details of a given protein systems biomolecular dynamics will represent a potentially
76 large source of latent variability in our functional understanding of the genome; a problem largely
77 ignored by the considerable amount of current analyses of static forms of 'omic' data (i.e. DNA
78 sequence, transcript level, and protein structure)(3). However, in the last decade, simultaneous
79 advances in the development of gaming graphics hardware and biomolecular force fields has elevated
80 our ability to computationally simulate molecular dynamics (MD) long enough to capture the most
81 biologically relevant timescales for protein function (4, 5) (i.e. ns to μs). And now, the application of
82 proper statistical methods to comparatively framed questions about protein dynamics can potentially
83 enable meaningful comparisons of large ensembles of shorter timed framed MD simulation (6). But
84 despite the rich visual complexity of data and movie images generated by MD software, the functional
85 interpretation of these large complex data sets created by modern MD simulations poses a serious
86 challenge to current users, especially where large ensembles of runs are collected and compared. A
87 potential solution to this problem exists with the application of machine learning to the feature
88 extraction and classification of the dynamic differences between ensembles of MD runs representing
89 the functional states of biomolecular systems (e.g. before/after chemical mutation or binding
90 interaction). In other words, high performance computation to generate simulated protein motions for
91 comparison can be effectively partnered with high performance methods for optimally extracting and
92 learning what the underlying dynamic feature differences are that define the functional protein states.
93 Although machine learning has recently been applied to individual MD studies for a variety of specific
94 tasks (7–9), there is no current software platform for the general application of machine learning to
95 comparative protein dynamics.

96        In 2018, we released DROIDS v1.2 and v2.0 (Detecting Relative Outlier Impacts from molecular
97 Dynamic Simulation), a GPU accelerated software pipeline designed for calculating and visualizing
98 statistical comparisons of protein dynamics drawn from large repeated ensembles of short dynamic
99 simulations representing two protein states (6). This application allowed simple visual and statistical
100 comparison of protein MD ensembles set up in any way the user wanted to define them. Here, we
101 announce the release of DROIDS v3.0, which comprises multiple pipelines specifically tailored for
102 specific functional comparisons. These include different temperatures, different binding states (i.e. to
103 DNA, drugs, toxins or natural ligands), or divergent genetic mutant states. Many other new
104 improvements to DROIDS are listed below, however we also include a major new machine learning tool,
105 maxDemon v1.0, a multi-machine learning post-processing application capable of training on the data
106 representing the normal functionally divergent dynamic states compared in DROIDS, and subsequently
107 identifying the effects on these learned dynamics upon the deployment of new MD simulations

108    representing genetic or drug class variants of interest. Thus, much like James Clerk Maxwell's mythical
109    creature (10), maxDemon derives important information from all atom resolution observation of
110    dynamic motion. The three primary features/aims of our newly expanded software is to (A) improve
111    user experience in comparative protein dynamics, (B) enable the local detection of functionally
112    conserved protein dynamics, and to (C) enable the assessment of the local dynamic impacts of both
113    genetic and drug class variants within the functional context of protein system of interest. Functionally
114    conserved dynamics is defined as 'sequence-dependent dynamics' discovered after training machine
115    learners on the functional state ensembles derived with DROIDS. Conserved dynamics are detected
116    through the deployment of learners on new MD simulation runs that were setup identically to the
117    dynamic state recognized or queried by the MD ensemble training sets. We expect that functionally
118    conserved dynamics will be sequence encoded and therefore should display a repeated position
119    dependent signature in our learned pattern profiles whenever MD runs are set up identically to MD
120    upon which learners were trained. Therefore, a significant local canonical correlation (i.e. Wilk's lambda)
121    between learning performance profiles of self-similar MD runs can be used to detect local regions of
122    conserved protein dynamics. By extension, mutational impacts of genetic or drug class variants on the
123    functionally conserved dynamics can be quantified by their effects that range significantly beyond that
124    observed in the self-similar runs. Thus when canonical correlations of variants differ significantly from
125    the self-correlation observed, according to a bootstrap or z-test, we plot the relative entropy or
126    magnitude of impact defining how the variant's dynamics differs from the self-similar dynamics of the
127    normal functioning protein. Because the machine learning algorithms are ultimately trained on MD data
128    representing this normal functioning state(s), this metric of impact is potentially very specific to how a
129    given mutation or drug impacts a specific protein. Thus, it potentially gives considerably more functional
130    relevance to our results when compared to more general database-derived metrics of mutational
131    tolerance (e.g. SIFT, PolyPhen2, MegaMD etc.).
132          In Table 1, we list five primary methods pipelines in available in DROIDS 3.0+maxDemon to
133    address functional questions in comparative protein dynamics. In our results and discussion here, we
134    present data on four case studies of functional protein dynamics that include feature extraction and
135    classification of (A) a simple temperature shift in ubiquitin dynamics, (B) mutational tolerance in
136    ubiquitin dynamics, (C) mutation specific impacts of DNA binding of TATA binding protein, and (D)
137    comparison of binding dynamics of drug class variants that mimic ATP binding in Hsp90.

138

# Materials and Methods

140    Overview

141    The machine learning-based detection of variant impacts on functional protein dynamics presented here
142    is outlined schematically in Figure 1. Upon launch, the DROIDS graphical user interface (GUI) will help
143    the user write a control file for required working path directories on their system (first use only) and
144    then proceeds to a main GUI outlining the various types of comparisons that can be generated (as
145    detailed in Table 1) and the number of GPU available on the system. The next step is provides a user-
146    friendly GUI to control and schedule Amber16/18 GPU-accelerated MD simulation to generate
147    ensembles of short MD runs representing two functional protein states wanting to be compared. These
148    functional comparisons are not limited, but would typically entail the impact of mutation (comparing
149    dynamics before and after one or more amino acid replacements), the impact of an environmental
150    change (comparing two states of temperature of solvent set up), or the impact of a molecular
151    interaction (comparing bound to an unbound state). The DROIDS GUI will lead users through the

152    building of a structural alignment file using UCSF Chimera's MatchMaker and Match-Align tools. This will
153    be needed later by the graphics components of DROIDS to make sure that only homologous regions of
154    structures are being compared and analyzed.    In this application, where the user is primarily interested
155    in genetic or drug class variant impacts on an interactive signaling function, the typical training
156    ensembles generated by DROIDS for further analysis with maxDemon should represent the normal
157    binding function of the wild-type protein and therefore the bound vs unbound comparison would
158    typically be used. A PDB file of the bound state can be the starting point and an unbound PDB model can
159    be saved after deleting chains in the original file.    If a small molecule ligand interaction is under study
160    and requires application of an additional force fields such as GAFF, than an additional file representing
161    only the ligand should also be generated and saved for preparation with antechamber software prior to
162    building the solvent models using teLeAP. The GUI will pop open the .bat files that control more details
163    of the simulation setup allowing advanced users to write more lines into the teLeAP modeling prep (e.g.
164    to alter the water box dimensions, the water model itself, or to add additional ions beyond simple
165    charge neutralization). The user should read all warnings provided to the terminal at this stage by the
166    Amber software. Our GUI script will also double check the sizes of the files generated at this stage and
167    will supply a warning if teLeap failed altogether to set up the complete model system for simulation.
168    Upon successful setup the user can launch all the MD runs from the GUI.    The requested jobs are
169    automatically scheduled to each GPU one at a time by our software. When finished, the user can easily
170    generate rmsf data by using the GUI to setup and launch cpptraj software provided in Ambertools. Thus
171    the total process from file preparation, MD production and post-processing for DROIDS analysis by
172    simply working down the buttons on each GUI from top to bottom and subsequently following the
173    directions on the main terminal. After MD simulation and post-processing, DROIDS will take users to a
174    second GUI for generating R plots and analyses for statistically comparing the dynamics, and then to a
175    third GIU for visualization and movie generation. We refer users to our user manual and previous
176    publication for more details. This third GUI has buttons to optionally launch our new machine learning
177    application maxDemon if users wish to go beyond simple comparative protein dynamics and investigate
178    novel simulations utilizing the DROIDS MD ensembles as a training set for subsequent machine learning.
179    Users can select any of seven different machine learning classification algorithms including K-nearest
180    neighbors, naïve Bayes, linear discriminant analysis, quadratic discriminant analysis, random forest,
181    adaptive boosting and support vector machine (kernels include linear, tuned polynomial, laplace and
182    tuned radial basis function). We generally recommend users select all of them. The random forest and
183    adaboost algorithms are programmed to use all available CPU cores found on the system.

184        To detect functionally conserved dynamics after training on two MD ensembles representing a
185    query and reference functional state of a protein (e.g. bound and unbound), the machine learners are
186    deployed on two new MD simulations across space and time (i.e. fluctuations backbone atoms of
187    individual amino acids over subdivided time intervals) that are identical to the original query state. Any
188    sequence dependent or 'functionally conserved' dynamics can be recognized through a significant
189    canonical correlation in the profile of the overall learning performance along the amino acid positions
190    for the two similar state runs. In effect, this metric defines dynamics that are functionally conserved by
191    capturing a signal of significant self-similarity in dynamics that localizes to a specific part of the protein
192    backbone.

193    *Eqn. 1*                    $conserved_{dynamics} = significant\ (CC_{self})$

194  The impacts of dissimilar states caused by altered amino acid sequence or different binding
195  partners can be subsequently assessed through their local effect on this canonical correlation (i.e.
196  conserved dynamics). Here, we introduce a metric of relative entropy relating the canonical correlations
197  in both the self-similar and altered state. In essence, this is a metric of the 'impact' of a given genetic or
198  drug class variant within the context of normal functioning dynamics. For example, when trained on a
199  natural binding interaction (e.g. DROIDS analysis comparing a DNA binding protein in its bound and
200  unbound states), novel MD simulations with a variety of amino acid replacements can be deployed to
201  see whether the learners can recognize the functional dynamics in the mutant forms. In this case,
202  functionally tolerated mutations will result in dynamics that do not vary outside of the normal bounds of
203  the self-similar runs, whereas functionally intolerant mutations will result in significant deviations from
204  self-similarity of motion.

205  *Eqn. 2*  $$variant_{impact} = CC_{self} * log \frac{CC_{variant}}{CC_{self}}$$

206  Variant impacts are only plotted when they differ significantly the self-similar MD runs according to a
207  simple Z test.

208  More detailed instructions to users are included with our DROIDS 3.0 user manual available in the
209  GitHub repository.

210  **The main repository for DROIDS 3.0 and maxDemon 1.0 can be found here. Please follow the link to**
211  **"Releases" and download the latest release as .tar.gz or .zip file**

212  https://github.com/gbabbitt/DROIDS-3.0-comparative-protein-dynamics

213  and DOI: 10.5281/zenodo.3358976 concurrent with this publication

214  https://zenodo.org/record/3358976#.XURVkOhKiiM

215  We also post various videos of examples using DROIDS, video tutorials, and ongoing projects here

216  https://www.youtube.com/channel/UCJTBqGq01pBCMDQikn566Kw

217  To enhance the user experience and scientific utility, DROIDS v3.0 offers many new features beyond
218  earlier major release versions 1.2 and 2.0.    These are summarized below.

219  -  New GUI organization directs users to specific comparative tasks/applications in Table 1
220  -  A new control file builder for managing path dependencies in Linux is included
221  -  Amber16/18 support has been beta tested and is defined via paths.ctl file
222  -  Single or dual GPU user options are available for faster analyses
223  -  Automated structure prep (dry and reduce) via pdb4amber is now included in the GUI. The
224      'reduce' variable is optional allowing users to either setup their own protonation states ahead of
225      DROIDS, or simply allow DROIDS to hydrogenate the input structures entirely.
226  -  Program/package dependency installer script named 'DROIDSinstaller.pl' is included. It will lead
227      users through all dependencies required after a fresh Linux build, including CUDA libraries and
228      tools required for Nvidia GPU accelerated Amber in the Linux environment
229  -  KL divergence (= relative entropy) definition of dFLUX is now included as an option providing a
230      richer color mapping of dFLUX in images and movies than the simple averaging algorithm
231      offered in earlier DROIDS versions

**dFLUX collected as an angstrom average**

$$dFLUX_{aa} = \left( \sum_{i=1}^{4} FLUX_{atom} \right)_{query} - \left( \sum_{i=1}^{4} FLUX_{atom} \right)_{reference}$$

$$dFLUX_{chain} = \sum_{i=1}^{L} |dFLUX_{aa}|$$ where L = number of structurally homologous amino acids in reference chain

$i = 4$ or avg atom flux on $N, O, CA$ and $C$ backbone atoms

**dFLUX collected as symmetric KL divergence**

$$dFLUX_{aa} = \left[ \begin{array}{c} D_{KL}(FLUX_{query}|FLUX_{reference}) \\ + D_{KL}(FLUX_{reference}|FLUX_{query}) \end{array} \right] \Big/ 2$$

where

$$D_{KL}(FLUX_{query}|FLUX_{reference}) = \sum_i FLUX_{query}(i) \log \frac{FLUX_{query}(i)}{FLUX_{reference}(i)}$$

$$D_{KL}(FLUX_{reference}|FLUX_{query}) = \sum_i FLUX_{reference}(i) \log \frac{FLUX_{reference}(i)}{FLUX_{query}(i)}$$

and

$i = 4$ or avg atom flux on $N, O, CA$ and $C$ backbone atoms

- Binding interaction analysis for both protein-DNA and protein-ligand systems is now offered with dedicated GUI for these comparisons. Protein-ligand system setup includes QMMM preprocessing in Antechamber and SQM.
- LeAP control files for explicit solvent runs are now presented for advanced user modifications (e.g. changing ion concentration, water model, water box dimension of volume).
- Dedicated GUI allowing genetic mutation placement (on DNA or AA) are included for setting up variants to analyze
- Self-stability and temperature shift analysis has its own dedicated GUI, allowing users to copy the input pdb file to compare MD ensembles generated on identical structures at the same of at different temperatures
- MaxDemon 1.0 - machine learning based detection of functionally conserved dynamic regions
- MaxDemon 1.0 - machine learning based impact assessment of variants (genetic, structural or binding)
- Dynamic visualization and movie rendering of machine learning classification performance
- Virtual reality and ChimeraX compatibility is also supported (additional information and download code can be found here https://cxtoolshed.rbvi.ucsf.edu/apps/moleculardynamicsviewer https://github.com/kdiller713/ChimeraX_MolecularDynamicViewer

To demonstrate the performance and utility of DROIDS 3.0 + maxDemon 1.0, we ran the following four comparative case studies using the PDB IDs mentioned below. Bound and unbound files were created by deleting binding partners in UCSF Chimera and resaving PDBs (e.g. 3t0z_bound.pdb, 3t0z_unbound and 3t0z_ligand). Each MD run ensemble consisted of 200 production runs at 0.5ns each explicitly solvated in a size 12 octahedral water box using TIP3P solvent model with constant temperature under an Anderson thermostat. The models were charge neutralized with both Na+ and Cl- ions. The heating and equilibration runs prior to production were 0.3ns and 10ns respectively. Prior to heating 2000 steps of energy minimization were also performed.    All seven available machine learning classifiers were trained on the functional MD ensembles and deployed upon new 5 ns production runs for each variant analyzed.

Case study 1 (figure 2) – PDB ID = 1ubq – to analyze self-stability and effect of temperature shift in ubiquitin

Case study 2 (figure 3)– PDB ID = 2oob – to analyze functional binding of ubiquitin to ubiquitin ligase and impacts of several tolerance pre-classified genetic variants

265  Case study 3 (figure 4) – PDB ID = 1cdw – to analyze functional binding of TATA binding protein to DNA
266  and impacts of several genetic variants

267  Case study 4 (figure 5) – PDB ID = 3t0z – to analyze functional ATP binding in Hsp90 and subsequent
268  impacts of six inhibitor drug variants

269

## Results and Discussion

271  To demonstrate the variety of comparative analyses that can be addressed with the new release of
272  DROIDS 3.0 and maxDemon 1.0, we chose four different case studies of comparative protein dynamics.
273  These included (A) an analysis of self-stability and temperature effects in free ubiquitin, (B) a functional
274  genetic variant analysis of ubiquitin and ubiquitin ligase binding interaction, (C) a functional genetic
275  variant analysis of DNA binding in TATA binding protein, and (D) a drug class variant analysis of
276  compounds targeting ATP binding region of Hsp90 heat shock protein.

### *Machine learning analysis of impacts due to simple environmental temperature shift*

278           We first ran a null comparison as a 'sanity check' by running a query and reference ubiquitin (11)
279  MD at the same temperatures (both 300K) and same solvent conditions. The DROIDS analysis (Figure 2A-
280  C) showed identical atom fluctuation profiles along the backbone and a random dFLUX profile indicative
281  of nonsignificant differences due to small random local thermal differences in the training sets. The
282  machine learning classification plots on new MD runs vary randomly around 0.5 reflecting the fact that
283  the learning algorithms had no features to train on (Figure 2D). As expected, no significantly conserved
284  dynamics were identified either (Figure 2E). By contrast, a protein dynamic comparison run with a 50K
285  temperature difference (Figure 2 F-H) shows a much higher machine learner performance upon
286  deployment (i.e. 70-80% successful classification – Figure 2I). Because environmental temperature shifts
287  are not expected to reflect evolutionary conserved dynamics (i.e. are not position dependent), they
288  subsequently do not result in canonical correlations in the learning profiles (Figure J). Representative
289  time slices of the positional classifications in each of these experiments are shown in K and L resp and
290  indicate that our machine learning is capable of extracting and identifying simple differences in
291  dynamics due to temperature. Another interesting observation here was the slightly higher learning
292  performance of the simpler machine learning methods QDA and LDA over others at all sites in the
293  temperature shifted example. We interpret this to be related to the fact that underlying rmsf
294  distributions are probably Gaussian, a critical assumption of these two models, with unequal variances
295  caused by steric hindrances on the backbone. This would predict that QDA might outperform other
296  learners in this situation and it appears that it does. We note that where more complex functional
297  dynamics are concerned, the more sophisticated learning methods such as support vector machine and
298  adaboost often perform slightly better than others.    However, we also note that these performance
299  differences are usually quite small and that all learning methods generally come to similar local
300  conclusions about functional dynamics. We now move on to demonstrate and examine machine
301  learning performance regarding more functional binding dynamics in ubiquitin and two other protein
302  systems.

303

304 ***Machine learning analysis of impacts of genetic variants on functional binding interactions***

305    To examine functional dynamics in ubiquitin, we conducted a DROIDS analysis comparing its two
306    functional states, bound and unbound to the ubiquitin associated binding domain of ubiquitin ligase
307    (12)(Figure 3A-D). This binding domain is highly conserved among many other proteins that interact
308    directly with ubiquitin. The binding interaction greatly reduces the atom fluctuation in ubiquitin at 3
309    characteristic positions, two loop structures centered at LEU 8 and ALA 46 and a portion of beta sheet at
310    the C terminus (Figure 3C). These three regions also drive significant differences in dynamics across the
311    whole protein. In novel self-similar MD runs on the bound state, we successfully detect significant
312    canonical correlations indicating conserved dynamics in these three regions with a broad expanse in
313    conserved dynamics (Figure 3E and F) across the UBA region (Figure 3G). We tested a set of 24
314    mutations that included sites with the most and least tolerated effects on growth rate in vivo in yeast
315    according to a study by Roscoe et al. (13).    In this study in vivo, nearly all mutations at E18 and G53 are
316    tolerated while nearly all mutations at K48 and R72 are not. Ultimately, the causes of tolerance in these
317    variants are not known, and do not necessarily invoke functional problems in dynamics. However, the
318    impacts that we did observed in simulation were on average twice as strong in the intolerant
319    backgrounds when compared to the mutation tolerant backgrounds. And, while we did not see large
320    differences in mutational impacts on dynamics between tolerated and non-tolerated mutant groups, the
321    24 mutations analyzed all show a general trend of dynamic impact falling outside of most of the
322    functional binding region (Figure 3H-K), suggesting that ubiquitin may have evolved a tertiary structure
323    that allosterically translates dynamic impacts to less functional regions of the protein. One exception to
324    this rule was demonstrated by the very large impact of R72D, centered squarely in the functionally
325    conserved binding region of ubiquitin, and would obviously heavily disrupt electrostatic charge
326    interactions as well.

327    We also conducted a DROIDS analysis comparing human TATA binding protein (TBP) (14) in its
328    functionally bound and unbound states (Figure 4A-C). TBP exhibits a characteristic signature of
329    dampening of atom fluctuation throughout its entire structure with most pronounced effects in two
330    loop regions that interact with the minor groove of DNA (arrows in Figure 4A and 4C). Canonical
331    correlations in new self-similar MD runs marking increased performance in classification were observed
332    in these regions (Figure 4D) along with corresponding regions of conserved dynamics identified by
333    significant Wilk's lamda (Figure 4E). Conserved dynamics from these loop areas are connected through
334    the chains in the beta sheet region of TBP spanning the DNA major groove contact. Mutational impacts
335    of four variants affecting the binding loop most proximal to the C terminal exhibited followed our
336    expectation of increasing impact ordering from R192Q, R192K, R192polyD, and R192polyW (Figure 4G
337    and 4H). The polyD and polyW mutations incorporated 5 sequential ASP or TRP residues centered at
338    R192, both causing the loop region to become more rigid (causing increased negative dFLUX). We
339    expected the strong functional binding affect observed across nearly all residues in this system would
340    make it relatively highly tolerant to single amino acid substitutions, even when located in the most
341    functional binding loop. In accordance with this idea, we found the most impactful multiple mutation
342    (i.e. R192polyW) significantly affected the dynamics of nearly 6 times more local residues than the least
343    impactful single substitution (i.e. R192Q).

344

345

346     ***Machine learning analysis of impacts of drug class variants targeting the ATP binding region of Hsp90***

347     Lastly, we conducted a DROIDS analysis comparing the dynamics of Hsp90 chaperone, a
348     common drug target for inhibitors in many cancer therapies, in both its ATP bound and unbound states.
349     The binding of ATP was discovered to significantly destabilize three co-localized alpha helical regions of
350     the protein adjacent to and extending from the ATP binding site (Figure 5A-D). MaxDemon analysis
351     confirmed the dynamics of this region to be highly conserved in new MD runs (Figure 5D-G). We also
352     analyzed the impacts of the six drug class variants targeting the ATP site (15–19), but interacting
353     differently with residues in this region (Figure 5H). The contacts in the ATP binding site are shown in
354     Figure 5I. While the localized patterns of impacts of the drug variants were all quite similar to ATP
355     (Figure 5J), the drug variants that most closely mimicked the contacts of ATP (i.e. geldanamycin) had far
356     less impact on dynamics than variants that interacted very differently with the binding pocket (i.e.
357     benzamide SNX1321 and inhibitor FJ1(Figure H-I). We feel that this finding demonstrates that while it is
358     important to be able to target a druggable protein binding site (20), researchers should also consider
359     how these various chemicals might alter the natural dynamics of the system. In situations where a drug
360     might too closely mimic the dynamic effects of a natural activator like ATP, a hyperactivation response
361     might occur in non-tumor cells leading to secondary cancer (21–23). Alternatively, other situations may
362     require drug targeting that does not alter the natural dynamic behavior too much, potentially activating
363     proteolytic systems in the cell. Our software allows more detailed investigations of these potential
364     dynamic impacts of drug class variants.

365     ***Conclusion***

366     We provide a well demonstrated method and user-friendly software pipeline for conducting statistically
367     sound comparative studies of large ensembles of comparative protein dynamics. The method/software
368     also now provides machine learning based extrapolations of effects on novel MD simulations
369     representing various functional variants of interest to the user. While there currently is at least one
370     other software allowing users to connect sequence-based evolutionary metrics to protein dynamics (24),
371     our method/software is unique in that regions of functional conservation are identified by analyzing
372     self-similar features of dynamics themselves rather than traditional sequence-based approaches, which
373     do not necessarily assume that a protein function has a strong dynamic component. Another advantage
374     to our method/software is that our functional impacts (i.e. mutational tolerance) are defined solely
375     within the context of protein dynamic system being simulated. This provides a much deeper look into
376     protein specific function than current genomic and proteomic database methods of predicting
377     mutational tolerance (25, 26) currently allow. As GPU technology continues to advance at a rapid pace
378     over the next few years, our method/software may have profound potential application to the
379     development of precision and personalized medicine, where understanding the detailed interaction
380     between genetic and drug class variants within the context of specific protein dynamic systems will be
381     greatly needed.

382

# Supporting Material

383

384     The main repository for DROIDS 3.0 and maxDemon 1.0 can be found at the GitHub repository link
385     below. Please follow the link to "Releases" and download the latest release as .tar.gz or .zip file

386   https://github.com/gbabbitt/DROIDS-3.0-comparative-protein-dynamics

387   We also post various videos of examples using DROIDS, video tutorials, and ongoing projects here

388   https://www.youtube.com/channel/UCJTBqGq01pBCMDQikn566Kw

389   A 'live version' of the figures in this manuscript is also available on our YouTube channel.

390

## Author Contribution

392   GAB and EPF conceived the project and method. All authors contributed to the code base. GAB and LEA
393   worked on beta testing and debugging.

394

## Acknowledgements

399

## References

401   1.   Henzler-Wildman, K.A., M. Lei, V. Thai, S.J. Kerns, M. Karplus, and D. Kern. 2007. A hierarchy of
402        timescales in protein dynamics is linked to enzyme catalysis. Nature. 450: 913–916.

403   2.   Niessen, K.A., M. Xu, A. Paciaroni, A. Orecchini, E.H. Snell, and A.G. Markelz. 2017. Moving in the
404        Right Direction: Protein Vibrations Steering Function. Biophys. J. 112: 933–942.

405   3.   Babbitt, G.A., E.E. Coppola, M.A. Alawad, and A.O. Hudson. 2016. Can all heritable biology really be
406        reduced to a single dimension? Gene. 578: 162–168.

407   4.   Götz, A.W., M.J. Williamson, D. Xu, D. Poole, S. Le Grand, and R.C. Walker. 2012. Routine
408        Microsecond Molecular Dynamics Simulations with AMBER on GPUs. 1. Generalized Born. J. Chem.
409        Theory Comput. 8: 1542–1555.

410   5.   Salomon-Ferrer, R., A.W. Götz, D. Poole, S. Le Grand, and R.C. Walker. 2013. Routine Microsecond
411        Molecular Dynamics Simulations with AMBER on GPUs. 2. Explicit Solvent Particle Mesh Ewald. J.
412        Chem. Theory Comput. 9: 3878–3888.

413   6.   Babbitt, G.A., J.S. Mortensen, E.E. Coppola, L.E. Adams, and J.K. Liao. 2018. DROIDS 1.20: A GUI-
414        Based Pipeline for GPU-Accelerated Comparative Protein Dynamics. Biophys. J. 114: 1009–1017.

415   7.   Díaz, Ó., J.A.R. Dalton, and J. Giraldo. 2019. Artificial Intelligence: A Novel Approach for Drug
416        Discovery. Trends Pharmacol. Sci. .

417    8.    Plante, A., D.M. Shore, G. Morra, G. Khelashvili, and H. Weinstein. 2019. A Machine Learning
418         Approach for the Discovery of Ligand-Specific Functional Mechanisms of GPCRs. Mol. Basel Switz.
419         24.

420    9.    Terayama, K., H. Iwata, M. Araki, Y. Okuno, and K. Tsuda. 2018. Machine learning accelerates MD-
421         based binding pose prediction between ligands and proteins. Bioinforma. Oxf. Engl. 34: 770–778.

422    10.    Maxwell, J.C. 2001. Theory of Heat. 9 Reprint edition. Mineola, NY: Dover Publications.

423    11.    Vijay-Kumar, S., C.E. Bugg, and W.J. Cook. 1987. Structure of ubiquitin refined at 1.8 A
424         resolution. J. Mol. Biol. 194: 531–544.

425    12.    Peschard, P., G. Kozlov, T. Lin, I.A. Mirza, A.M. Berghuis, S. Lipkowitz, M. Park, and K. Gehring.
426         2007. Structural basis for ubiquitin-mediated dimerization and activation of the ubiquitin protein
427         ligase Cbl-b. Mol. Cell. 27: 474–485.

428    13.    Roscoe, B.P., K.M. Thayer, K.B. Zeldovich, D. Fushman, and D.N.A. Bolon. 2013. Analyses of the
429         effects of all ubiquitin point mutants on yeast growth rate. J. Mol. Biol. 425: 1363–1377.

430    14.    Nikolov, D.B., H. Chen, E.D. Halay, A. Hoffman, R.G. Roeder, and S.K. Burley. 1996. Crystal
431         structure of a human TATA box-binding protein/TATA element complex. Proc. Natl. Acad. Sci. U. S.
432         A. 93: 4862–4867.

433    15.    Stebbins, C.E., A.A. Russo, C. Schneider, N. Rosen, F.U. Hartl, and N.P. Pavletich. 1997. Crystal
434         structure of an Hsp90-geldanamycin complex: targeting of a protein chaperone by an antitumor
435         agent. Cell. 89: 239–250.

436    16.    Li, J., L. Sun, C. Xu, F. Yu, H. Zhou, Y. Zhao, J. Zhang, J. Cai, C. Mao, L. Tang, Y. Xu, and J. He. 2012.
437         Structure insights into mechanisms of ATP hydrolysis and the activation of human heat-shock
438         protein 90. Acta Biochim. Biophys. Sin. 44: 300–306.

439    17.    Immormino, R.M., Y. Kang, G. Chiosis, and D.T. Gewirth. 2006. Structural and quantum chemical
440         studies of 8-aryl-sulfanyl adenine class Hsp90 inhibitors. J. Med. Chem. 49: 4953–4960.

441    18.    Fadden, P., K.H. Huang, J.M. Veal, P.M. Steed, A.F. Barabasz, B. Foley, M. Hu, J.M. Partridge, J.
442         Rice, A. Scott, L.G. Dubois, T.A. Freed, M.A.R. Silinski, T.E. Barta, P.F. Hughes, A. Ommen, W. Ma,
443         E.D. Smith, A.W. Spangenberg, J. Eaves, G.J. Hanson, L. Hinkley, M. Jenks, M. Lewis, J. Otto, G.J.
444         Pronk, K. Verleysen, T.A. Haystead, and S.E. Hall. 2010. Application of chemoproteomics to drug
445         discovery: identification of a clinical candidate targeting hsp90. Chem. Biol. 17: 686–694.

446    19.    Austin, C., S.N. Pettit, S.K. Magnolo, J. Sanvoisin, W. Chen, S.P. Wood, L.D. Freeman, R.J.
447         Pengelly, and D.E. Hughes. 2012. Fragment screening using capillary electrophoresis (CEfrag) for hit
448         identification of heat shock protein 90 ATPase inhibitors. J. Biomol. Screen. 17: 868–876.

449    20.    Vajda, S., D. Beglov, A.E. Wakefield, M. Egbert, and A. Whitty. 2018. Cryptic binding sites on
450         proteins: definition, detection, and druggability. Curr. Opin. Chem. Biol. 44: 1–8.

451    21.    Poulikakos, P.I., C. Zhang, G. Bollag, K.M. Shokat, and N. Rosen. 2010. RAF inhibitors
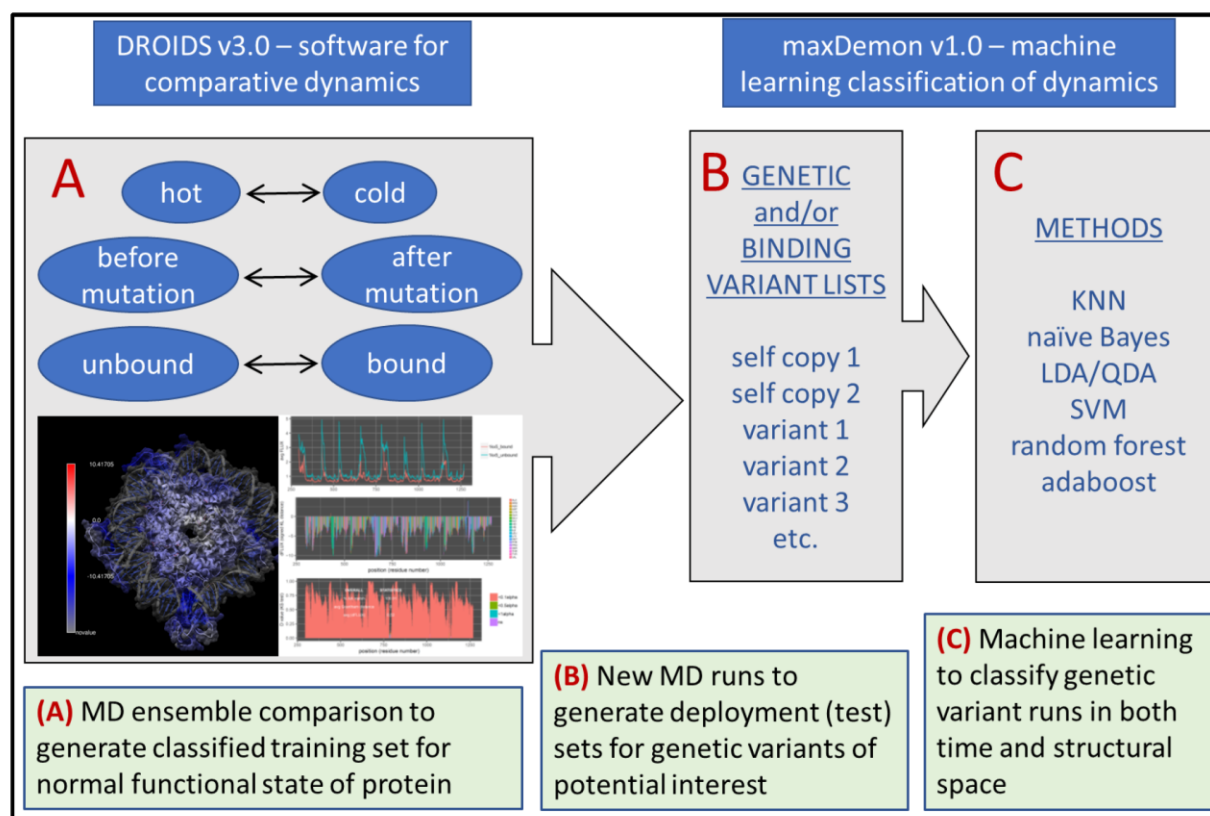452         transactivate RAF dimers and ERK signalling in cells with wild-type BRAF. Nature. 464: 427–430.

453  22.      Hatzivassiliou, G., K. Song, I. Yen, B.J. Brandhuber, D.J. Anderson, R. Alvarado, M.J.C. Ludlam, D.
454          Stokoe, S.L. Gloor, G. Vigers, T. Morales, I. Aliagas, B. Liu, S. Sideris, K.P. Hoeflich, B.S. Jaiswal, S.
455          Seshagiri, H. Koeppen, M. Belvin, L.S. Friedman, and S. Malek. 2010. RAF inhibitors prime wild-type
456          RAF to activate the MAPK pathway and enhance growth. Nature. 464: 431–435.

457  23.      Cichowski, K., and P.A. Jänne. 2010. Drug discovery: inhibitors that activate. Nature. 464: 358–
458          359.

459  24.      Bakan, A., A. Dutta, W. Mao, Y. Liu, C. Chennubhotla, T.R. Lezon, and I. Bahar. 2014. Evol and
460          ProDy for bridging protein sequence evolution and structural dynamics. Bioinformatics. 30: 2681–
461          2683.

462  25.      Kumar, P., S. Henikoff, and P.C. Ng. 2009. Predicting the effects of coding non-synonymous
463          variants on protein function using the SIFT algorithm. Nat. Protoc. 4: 1073–1081.

464  26.      Adzhubei, I., D.M. Jordan, and S.R. Sunyaev. 2013. Predicting functional effect of human
465          missense mutations using PolyPhen-2. Curr. Protoc. Hum. Genet. Chapter 7: Unit7.20.

466

467

468     Table 1. Common learner assisted comparative protein dynamic investigations enabled by

469     DROIDS 3.0 + maxDemon 1.0.

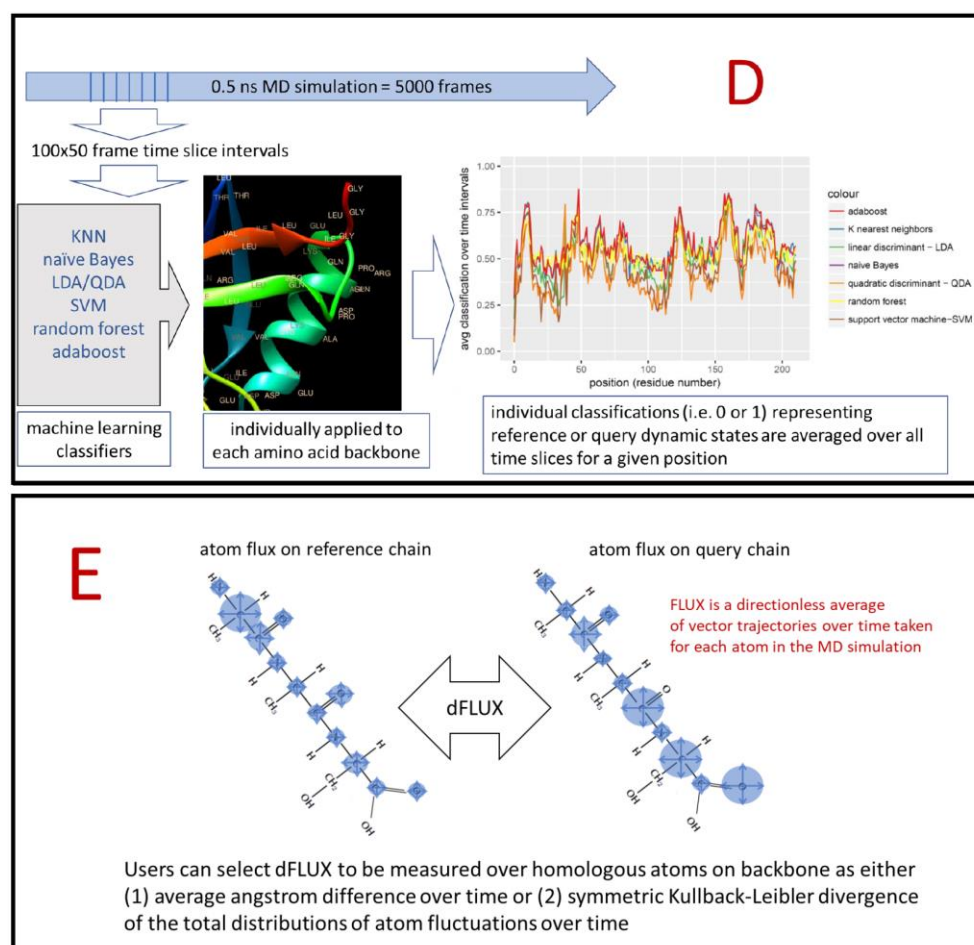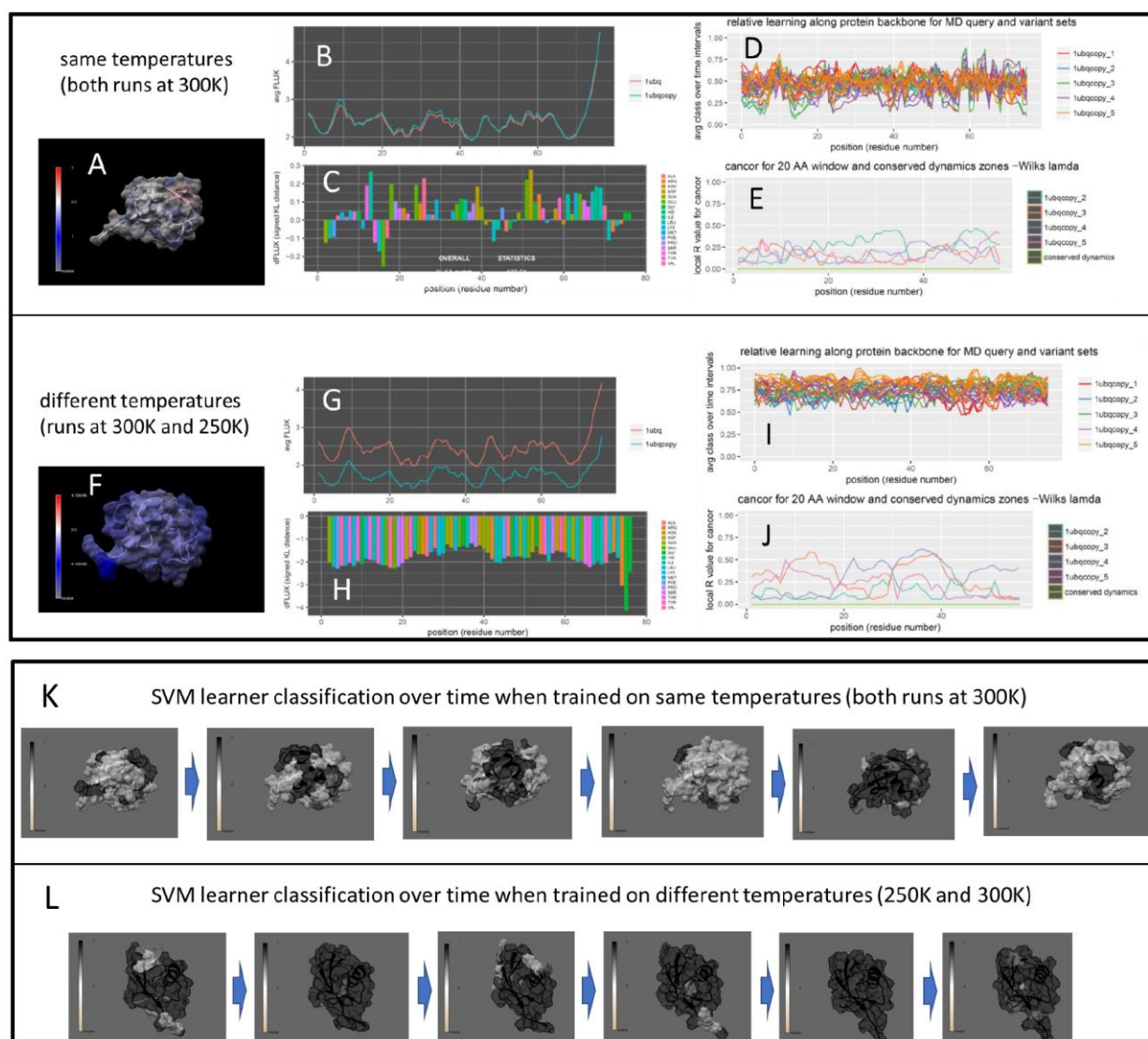| QUESTION | DROIDS 3.0 training comparison | Deployment of learners in maxDemon | Important notes |
|---|---|---|---|
| Measure dynamic tolerances of single protein to various genetic mutations | Two sets (ensembles) of MD on the same protein at the same temperature | MD run on one or more genetic mutant structures | Isolates MD impacts of mutation(s) from natural variability in self-similar dynamics |
| Measure dynamic tolerances of DNA binding interaction to genetic mutation(s) | MD ensembles comparing both the unbound and   DNA bound protein | MD run on one or more unbound genetic mutant structures | Isolates MD impacts of mutation from natural binding function of the system |
| Measure dynamic tolerances of individual genetic differences to a given drug | MD ensembles comparing both the unbound and drug bound protein | MD run on one or more drug-bound genetic mutant structures | Isolates MD impacts of mutation from novel drug binding function of the system |
| Measure dynamic similarities of different drug candidates to | MD ensembles comparing both the unbound and ligand bound protein | MD run on one or more drug variant bound structures | Isolates MD impacts of drug candidates from the natural binding function of the ligand |

479

**Figure 1. Schematic overview of DROIDS 3.0 + maxDemon 1.0 software for machine learning-based detection of variant impacts on functionally conserved protein dynamics**. The pipeline starts with (A) generation of two large ensembles of molecular dynamic (MD) simulations that represent a functional comparison of protein states (e.g. mutation, binding or environmental change). The root mean square fluctuations (rmsf) of protein backbone atoms in these ensembles are comparatively analyzed/visualized (i.e. using DROIDS) and are also later used as pre-classified training data sets for machine learning (i.e. using maxDemon). Note: in the pictured DROIDS analysis of nucleosome shows overall dampening of rmsf in the histone core with maximal dampening where the histone tails cross the DNA helix (B) New MD simulations are generated on two structures self-similar to the query state of training as well as a list of functional variants, and (C) up to seven machine learning methods are employed to classify the MD in the self-similar and variant runs according to the functional comparison defined by the initial training step. (D) The relative efficiency of learning is defined by average value of classification (i.e. 0 or 1) over 50 frame time slices for each amino acid position and regions of functionally conserved dynamics are later identified by significant canonical correlations in this learning efficiency (i.e. Wilk's lamda) in self-similar MD runs. The impacts of variants are defined by relative entropy of variant MD compared to the MD in the self-similar runs and plotted when this entropy is significantly different from the variation in self-similarity (i.e. bootstrapped z-test). (E) A visual representation of the difference in local rmsf (dFLUX) is typically calculated using symmetric Kullback-Leibler (KL) divergence between the two distributions of rmsf in the training MD ensembles.

499



500

501

**Figure 2. Analysis of environmental temperature change on non-functional ubiquitin dynamics.**
DROIDS image and analysis of random ubiquitin dynamics compared at the same (A-E) and different (F-J) temperatures. Note: blue color quantifies damped rmsf at temperature lowered by 50K. Note that relative learning is much higher when a temperature difference is modeled (D and I resp), however, as expected, neither comparison offers the machine learners a sequence-dependent profile by which to establish a signal of conserved dynamics (E or J). The learner classifications for the best performing learner in this case (quadratic discriminant function: QDA) is shown imaged on the ubiquitin structure over time in both the (K) random dynamics and (L) temperature dampened dynamics. (Movies of this can be observed in supplemental file A)
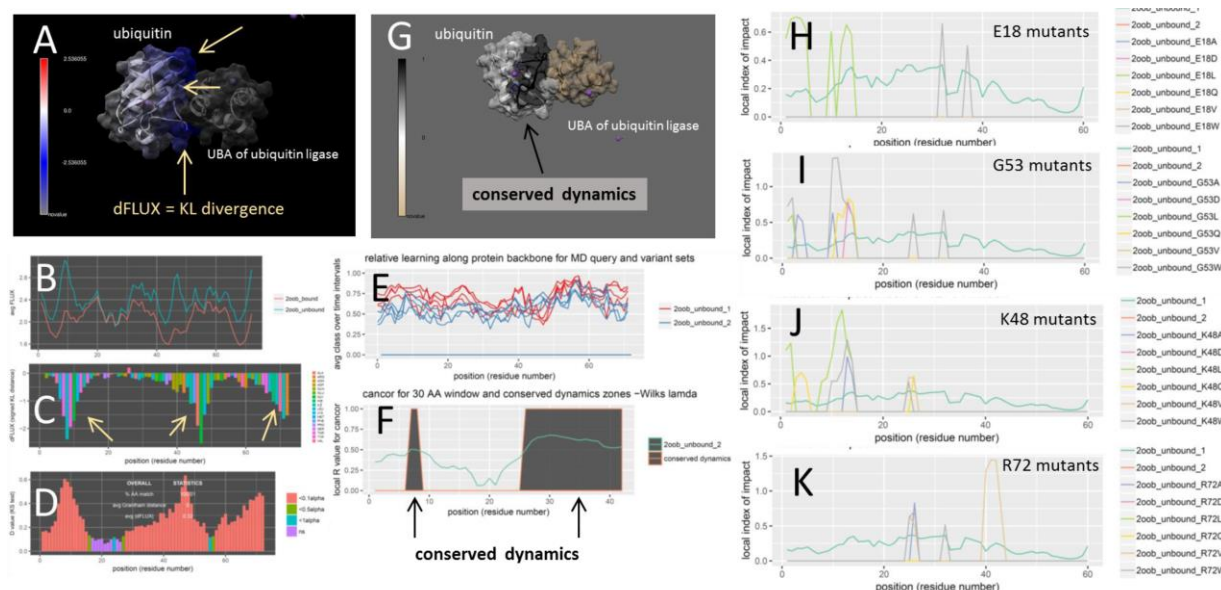
511

512

513

514

515



516

**Figure 3. Analysis of mutational impact and tolerance on functional ubiquitin dynamics.** (A) DROIDS image and analysis of ubiquitin bound to the ubiquitin associated binding domain (UBA) of ubiquitin ligase. Note: blue color quantifies damped rmsf at binding interface. (i.e. negative dFLUX) also by the (B) respective rmsf profiles of bound and unbound training states and (C) the KL divergence or dFLUX profile colored by residue. Arrows indicate most prominent dampening of rmsf near loops at THR 9, ALA 46 and C terminus. (D) Significant differences in these rmsf profiles is determined by multiple-test corrected two sample KS test. (E) Local relative learning efficiency of each machine learning method in self-similar testing runs are shown color-coded by run and regions of functionally conserved dynamics, determined via significant local canonical correlation are shown in dark gray in both (F) traditional N to C terminal plot as well as (G) structural image. The mutational impacts of 24 genetic variants (H-K: six variants at each or four sites) are shown all demonstrating lack of impact in functionally conserved regions of the binding interaction.

529

530

531

532

533

534

535

536

537



**Figure 4. Analysis of mutational impact and tolerance on DNA binding in Tata Binding Protein (TBP).**
DROIDS image and analysis of TBP in DNA-bound and unbound states showing (A) colored TBP structure, (B) respective rmsf profiles and (C) KL divergence (dFLUX) plot. Note: arrows indicate functional binding loops in the DNA minor groove red color indicates dampened rmsf. maxDemon analysis (D-E) identifying conserved dynamics supporting both minor groove binding loops and (F) connecting them through the central region of the beta sheet in the main body of TBP closest to the DNA. Mutational impacts of 4 genetic variants with increasing impact one of the functional loops are also shown (G) plotted and (H) on the TBP structure. They are R192K, R192D, R192Q and polyW centered at R192 in 1cdw.pdb and and position 161 (red arrow) in plots (Note: 31 position offset is due to DNA in the original file).

548

I. Ligand contacts of drug binding variants on protein dynamics of Hsp90 in the ATP binding pocket

| ligand | PDB | % sim | Val 186 | Thr 184 | Trp 162 | Tyr 139 | Phe 138 | Lys 112 | Alal 111 | Leu 107 | Asn 106 | Asp 102 | Met 98 | Gly 97 | Ile 96 | Asp 93 | Lys 58 | Ala 55 | Asp 54 | Asn 51 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ATP | 3t0z | 100 | x | | | | x | | | | x | | x | x | | x | | x | | x |
| Geldanamycin | 1yet | 87.5 | | x | | | x | x | | | x | | | x | | x | | x | | x |
| Inhibitor PU-H64 | 2fwy | 75 | | x | x | | x | | | x | | | x | | | x | | x | | x |
| Inhibitor PU-DZ8 | 2h55 | 75 | | x | x | | x | | | x | | | x | | x | x | x | x | | x |
| Radicicol | 4egk | 62.5 | x | x | | | x | | | x | | | | | x | x | | x | x | x |
| Benzamide SNX-1321 | 3mnr | 37.5 | | x | x | x | | | | x | | x | | | x | x | x | x | | |
| Inhibitor FJ1 | 4hy6 | 0 | | | x | x | | | x | | | | | | | | | | | |

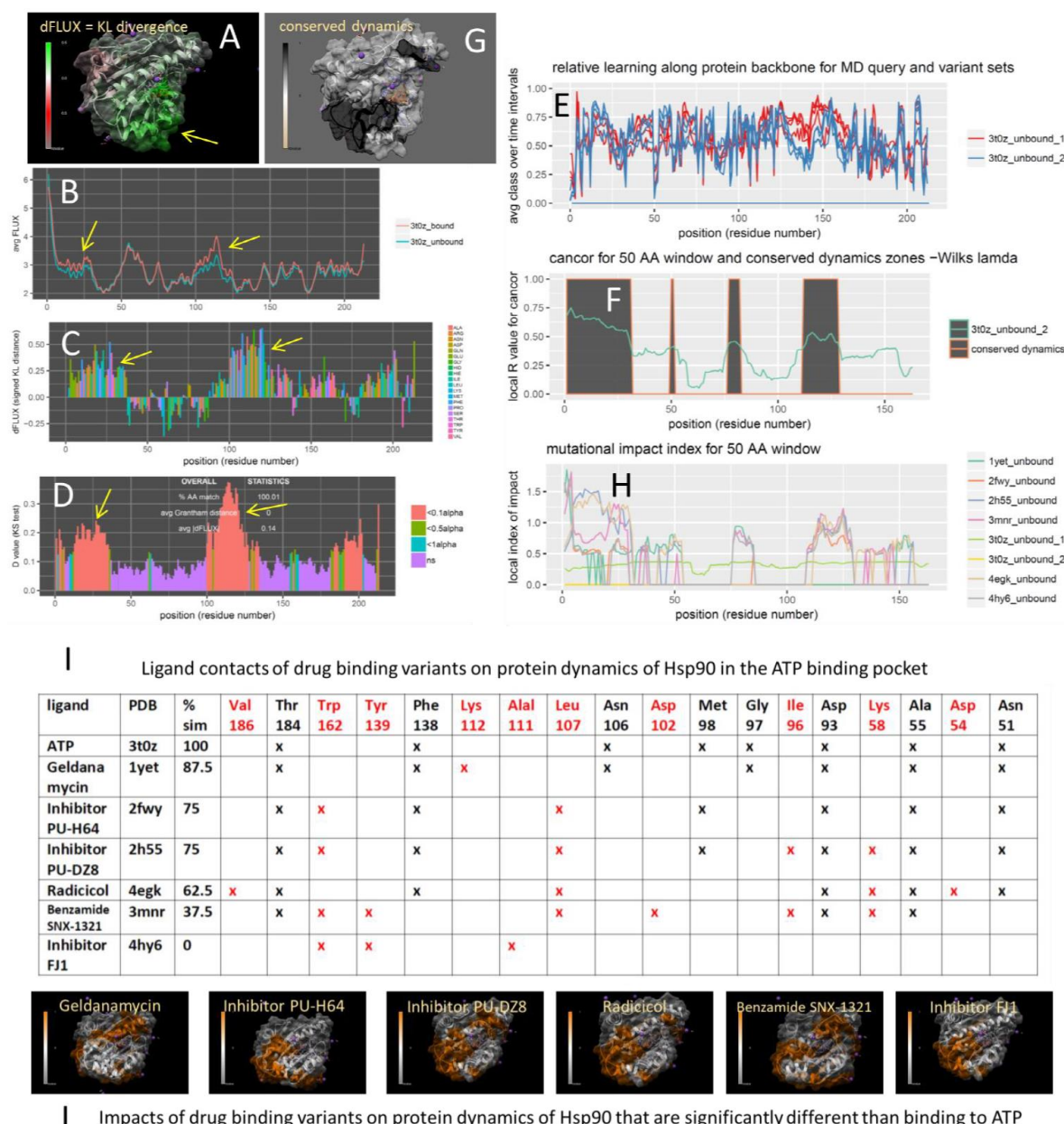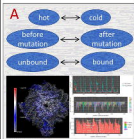J. Impacts of drug binding variants on protein dynamics of Hsp90 that are significantly different than binding to ATP

549

**Figure 5. Analysis of drug class variant binding in the ATP-binding domain of Hsp90.** DROIDS image and analysis of Hsp90 in ATP-bound and unbound states showing (A) colored Hsp90 structure, (B) respective rmsf profiles and (C) KL divergence (dFLUX) plot and (D) significant differences in dynamics determined via the KS test. Note: arrows and green color indicate regions where rmsf is amplified in response to ATP binding. maxDemon analysis (E-G) identifying conserved dynamics connecting the ATP binding pocket and region of amplified rmsf. (H) Mutational impacts of 6 drug class variants targeting the ATP binding pocket of Hsp90 are plotted and (I) ordered by number of differences in structural contacts within the binding pocket. (J) Mutational impacts of these variants are demonstrated to predominantly impact the functionally conserved region of amplified rmsf thus mimicking the dynamic effect of functional ATP binding.

| DROIDS v3.0 – software for comparative dynamics | maxDemon v1.0 – machine learning classification of dynamics |

**A**

- hot — cold
- before mutation — after mutation
- unbound — bound

**B** GENETIC and/or BINDING VARIANT LISTS

- self copy 1
- self copy 2
- variant 1
- variant 2
- variant 3
- etc.

**C** METHODS

- KNN
- naive Bayes
- LDA/QDA
- SVM
- random forest
- adaboost

(A) MD ensemble comparison to generate classified training set for normal functional state of protein

(B) New MD runs to generate deployment (test) sets for genetic variants of potential interest

(C) Machine learning to classify genetic variant runs in both time and structural space

conserved dynamics

A. dl LUX – RL divergence

B.

C.

D. relative learning along protein backbone for MB query and variant sets

E. sensor for 20 AA window and conserved dynamics stress – Holic series

F. conserved dynamics – block

H. R192X

R150X

R192polyW

R192polyD

G. mutational impact metric for 20 AA window