

Ensembler: Enabling high-throughput molecular simulations at the superfamily scale

Daniel L. Parton,¹ Patrick B. Grinaway,¹ Sonya M. Hanson,¹ Kyle A. Beauchamp,¹ and John D. Chodera^{1,*}

¹Computational Biology Program, Sloan Kettering Institute,
Memorial Sloan Kettering Cancer Center, New York, NY 10065

(Dated: October 15, 2015)

The rapidly expanding body of available genomic and protein structural data provides a rich resource for understanding protein dynamics with biomolecular simulation. While computational infrastructure has grown rapidly, simulations on an *omics* scale are not yet widespread, primarily because software infrastructure to enable simulations at this scale has not kept pace. It should now be possible to study protein dynamics across entire (super)families, exploiting both available structural biology data and conformational similarities across homologous proteins. Here, we present a new tool for enabling high-throughput simulation in the genomics era. **Ensembler** takes any set of sequences—from a single sequence to an entire superfamily—and shepherds them through various stages of modeling and refinement to produce simulation-ready structures. This includes comparative modeling to all relevant PDB structures (which may span multiple conformational states of interest), reconstruction of missing loops, addition of missing atoms, culling of nearly identical structures, assignment of appropriate protonation states, solvation in explicit solvent, and refinement and filtering with molecular simulation to ensure stable simulation. The output of this pipeline is an ensemble of structures ready for subsequent molecular simulations using computer clusters, supercomputers, or distributed computing projects like Folding@home. **Ensembler** thus automates much of the time-consuming process of preparing protein models suitable for simulation, while allowing scalability up to entire superfamilies. A particular advantage of this approach can be found in the construction of kinetic models of conformational dynamics—such as Markov state models (MSMs)—which benefit from a diverse array of initial configurations that span the accessible conformational states to aid sampling. We demonstrate the power of this approach by constructing models for all catalytic domains in the human tyrosine kinase family, using all available kinase catalytic domain structures from any organism as structural templates.

Ensembler is free and open source software licensed under the GNU General Public License (GPL) v2. It is compatible with Linux and OS X. The latest release can be installed via the `conda` package manager, and the latest source can be downloaded from <https://github.com/choderalab/enssembler>.

Keywords: molecular dynamics simulation; comparative modeling; distributed simulation

I. INTRODUCTION

Recent advances in genomics and structural biology have helped generate an enormous wealth of protein data at the level of amino-acid sequence and three-dimensional structure. However, proteins typically exist as an ensemble of thermally accessible conformational states, and static structures provide only a snapshot of their rich dynamical behavior. Many functional properties—such as the ability to bind small molecules or interact with signaling partners—require transitions between states, encompassing anything from reorganization of sidechains at binding interfaces to domain motions to large scale folding-unfolding events. Drug discovery could also benefit from a more extensive consideration of protein dynamics, whereby small molecules might be selected based on their predicted ability to bind and trap a protein target in an inactive state [1].

Molecular dynamics (MD) simulations have the capability, in principle, to describe the time evolution of a protein in atomistic detail, and have proven themselves to be a useful tool in the study of protein dynamics. A number of mature software packages and forcefields are now available, and much recent progress has been driven by advances in computing architecture. For example, many MD

packages are now able to exploit GPUs [2, 3], which provide greatly improved simulation efficiency per unit cost relative to CPUs, while distributed computing platforms such as Folding@home [4], Copernicus [5, 6], and GPGPUGrid [7], allow scalability on an unprecedented level. In parallel, methods for building human-understandable models of protein dynamics from noisy simulation data, such as Markov state modeling (MSM) approaches, are now reaching maturity [8–10]. MSM methods in particular have the advantage of being able to aggregate data from multiple independent MD trajectories, facilitating parallelization of production simulations and thus greatly alleviating overall computational cost. **There also exist a number of mature software packages for comparative modeling of protein structures, in which a target protein sequence is modeled using one or more structures as templates [11–14].**

However, it remains difficult for researchers to exploit the full variety of available protein sequence data (in simulating groups of related proteins) and structural data (exploiting multiple structures for each protein and its homologs/orthologs) in simulation studies in molecular simulations, largely due to limitations in software architecture. For example, the preparation of a biomolecular simulation is typically performed manually, encompassing a series of fairly standard (yet time-consuming) steps such as the choice of protein sequence construct and starting structure(s), addition of missing residues and atoms, solvation with explicit water and counterions (and potentially buffer

* Corresponding author; john.chodera@choderalab.org

components and cosolvents), choice of simulation parameters (or parameterization schemes for components where parameters do not yet exist), system relaxation with energy minimization, and one or more short preparatory MD simulations to equilibrate the system and relax the simulation cell. Due to the laborious and manual nature of this process, simulation studies typically consider only one or a few proteins and starting configurations, though notable exceptions exist, such as the Dymeomics effort of Daggett and coworkers in which over 100 proteins have been simulated so far using a single initial configuration for each [15]. Worse still, studies (or collections of studies) that *do* consider multiple proteins often suffer from the lack of consistent best practices in this preparation process, making comparisons between related proteins unnecessarily difficult.

The ability to fully exploit the large quantity of available protein sequence and structural data in biomolecular simulation studies could open up many interesting avenues for research, enabling the study of entire protein families or superfamilies within a single organism or across multiple organisms. The similarity between members of a given protein family could be exploited to generate arrays of conformational models for related sequences, which could be used as starting configurations to aid sampling in MD simulations. The conformations captured in structures of related members has been shown to provide useful information about the conformations accessible to all members of the family [16, 17], though energetic differences between individuals will modify the populations and dynamics of individual conformational states. This approach would be highly beneficial for many MD methods, such as MSM construction, which require global coverage of the conformational landscape to realize their full potential, and would also be particularly useful in cases where structural data is present for only a subset of the members of a protein family. It would also aid in studying protein families known to have multiple metastable conformations—such as kinases—for which the combined body of structural data for the family may cover a large range of these conformations, while the available structures for any individual member might encompass only one or two distinct conformations.

Here, we present the first steps toward bridging the gap between biomolecular simulation software and *omics*-scale sequence and structural data: a fully automated open source framework for building simulation-ready protein models in multiple conformational substates scalable from single sequences to entire superfamilies. **Ensembler** provides functions for selecting target sequences and homologous template structures, and (by interfacing with a number of external packages) performs pairwise alignments, comparative modeling of target-template pairs, and several stages of model refinement. As an example application, we have constructed models for the entire set of human tyrosine kinase (TK) catalytic domains, using all available structures of protein kinase domains (from any species) as templates. This results in a total of almost 400,000 models, and we demonstrate that these provide wide-ranging coverage of known functionally relevant conformations. By us-

ing these models as starting configurations for highly parallel MD simulations, we expect their structural diversity to greatly aid in sampling of conformational space. We further suggest that models with high target-template sequence identity are the most likely to represent native metastable states, while lower sequence identity models would aid in sampling of more distant regions of accessible phase space. It is also important to note that some models (especially low sequence identity models) may not represent natively accessible conformations. However, MSM methods benefit from the ability to remove outlier MD trajectories which start from non-natively accessible conformations, and which would thus be unconnected with the phase space sampled in other trajectories. These methods essentially identify the largest subset of Markov nodes which constitute an ergodic network [18, 19].

We anticipate that **Ensembler** will prove to be useful in a number of other ways. For example, the generated models could represent valuable data sets even without subsequent production simulation, allowing exploration of the conformational diversity present within the available structural data for a given protein family. Furthermore, automation of simulation preparation provides an excellent opportunity to make concrete certain "best practices", such as the choice of simulation parameters, approach to the treatment of protonation states, treatment of cofactors and structural ions, and pre-simulation refinement and equilibration procedures. While the current version of **Ensembler** only codifies some of these choices as default parameters, its modular nature allows additional stages to be easily added in the future.

II. DESIGN AND IMPLEMENTATION

Ensembler is written in Python, and can be used via a command-line tool (`ensembl`) or via a flexible Python API to allow integration of its components into other applications. All command-line and API information in this article refers to the [version 1.0.6 release of Ensembler](#). Up-to-date documentation can be found at [ensembl.readthedocs.org](#).

The **Ensembler** modeling pipeline comprises a series of stages which are performed in a defined order. A visual overview of the pipeline is shown in Fig. 1. The various stages of this pipeline are described in detail below.

A. Target selection and retrieval

The first stage entails the selection of a set of *target* protein sequences—the sequences for which the user is interested in generating simulation-ready structural models. This may be a single sequence—such as a full-length protein or a construct representing a single domain—or a collection of sequences, such as a particular domain from an entire family of proteins. The output of this stage is a FASTA-formatted text file containing the desired target sequences

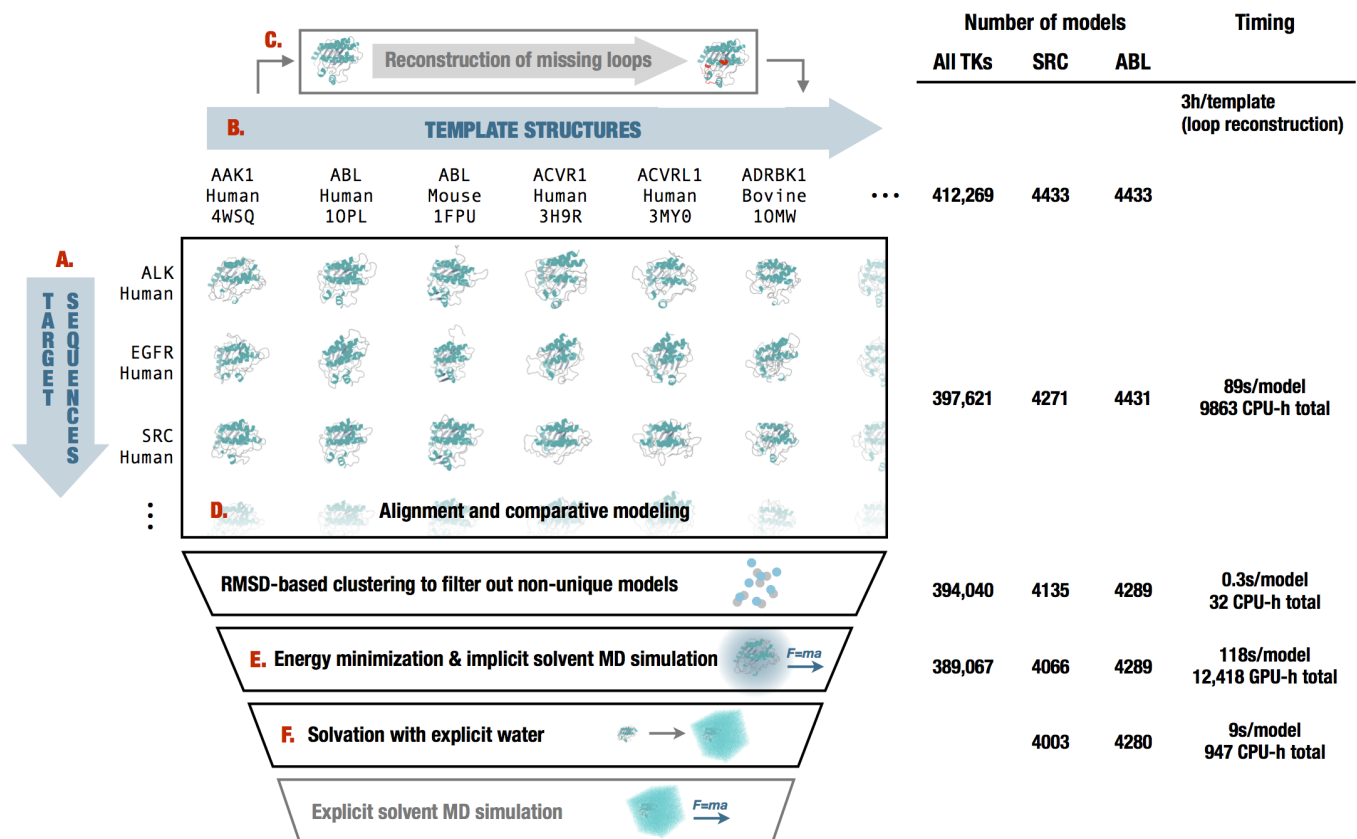


FIG. 1. Diagrammatic representation of the stages of the Ensembler pipeline and illustrative statistics for modeling all human tyrosine kinase catalytic domains. On the left, the various stages of the **Ensembler** pipeline are shown. The red labels indicate the corresponding text description provided for each stage in the Design and Implementation section. On the right, the number of viable models surviving each stage of the pipeline is shown for the 93 target TK domains and for two representative individual TK domains (*SRC* and *ABL*). Typical timings on a computer cluster (containing Intel Xeon E5-2665 2.4GHz hyperthreaded processors and NVIDIA GTX-680 or GTX-Titan GPUs) is reported to illustrate resource requirements per model for modeling the entire set of tyrosine kinases. Note that *CPU-h* denotes the number of hours consumed by the equivalent of a single CPU hyperthread and *GPU-h* on a single GPU—parallel execution via MPI reduces wall clock time nearly linearly.

with corresponding arbitrary identifiers.

The `enssembler` command-line tool allows targets to be selected from UniProt—a freely accessible resource for protein sequence and functional data (uniprot.org) [20]—via a UniProt search query. To retrieve target sequences from UniProt, the subcommand `gather_targets` is used with the `--query` flag followed by a UniProt query string conforming to the same syntax as the search function available on the UniProt website. For example, `--query 'mnemonic:SRC_HUMAN'` would select the full-length human Src sequence, while the query shown in Box 1 would select all human tyrosine protein kinases which have been reviewed by a human curator. In this way, the user may select a single protein, many proteins, or an entire superfamily from UniProt. The program outputs a FASTA file, setting the UniProt mnemonic (e.g. `SRC_HUMAN`) as the identifier for each target protein.

In many cases, it will be desirable to build models of an isolated protein domain, rather than the full-length protein. The `gather_targets` subcommand allows protein domains to be selected from UniProt data by passing a regu-

lar expression string to the `--uniprot_domain_regex` flag. For example, the above `--query` flag for selecting all human protein kinases returns UniProt entries with domain annotations including "Protein kinase", "Protein kinase 1", "Protein kinase 2", "Protein kinase; truncated", "Protein kinase; inactive", "SH2", "SH3", etc. The regular expression shown in Box 1 selects only domains of the first three types. If the `--uniprot_domain_regex` flag is used, target identifiers are set with the form `[UniProt mnemonic]_D[domain index]`, where the latter part represents a 0-based index for the domain—necessary because a single target protein may contain multiple domains of interest (e.g. `JAK1_HUMAN_D0`, `JAK1_HUMAN_D1`).

Target sequences can also be defined manually (or from another program) by providing a FASTA-formatted text file containing the desired target sequences with corresponding arbitrary identifiers.

B. Template selection and retrieval

Ensembler uses comparative modeling to build models, and as such requires a set of structures to be used as templates. The second stage thus entails the selection of templates and storage of associated sequences, structures, and identifiers. These templates can be specified manually, or using the `enssembler gather_templates` subcommand to automatically select templates based on a search of the Protein Data Bank (PDB) or UniProt. **A recommended approach is to select templates from UniProt which belong to the same protein family as the targets, guaranteeing homology and some degree of sequence identity between targets and templates.**

The `enssembler gather_templates` subcommand provides methods for selecting template structures from either UniProt or the PDB (<http://www.rcsb.org/pdb>), specified by the `--gather_from` flag. **Both methods select templates at the level of PDB chains—a PDB structure containing multiple chains with identical sequence spans (e.g. for crystals with non-crystallographic symmetry giving rise to independent conformations of the protein within the asymmetric unit) would thus give rise to multiple template structures.**

Selection of templates from the PDB simply requires passing a list of PDB IDs as a comma-separated string, e.g. `--query 2H8H,1Y57`. Specific PDB chain IDs can optionally also be selected via the `--chainids` flag. The program retrieves structures from the PDB server, as well as associated data from the SIFTS service (www.ebi.ac.uk/pdbe/docs/sifts) [21], which provides residue-level mappings between PDB and UniProt entries. The SIFTS data is used to extract template sequences, retaining only residues which are resolved and match the equivalent residue in the UniProt sequence—non-wildtype residues are thus removed from the template structures. Furthermore, PDB chains with less than a given percentage of resolved residues (default: 70%) are filtered out. Sequences are stored in a FASTA file, with identifiers of the form `[UniProt mnemonic]_D[UniProt domain index]_[PDB ID]_[PDB chain ID]`, e.g. `SRC_HUMAN_DO_2H8H_A`. Matching residues then extracted from the original coordinate files and stored as PDB-format coordinate files.

Selection of templates from UniProt proceeds in a similar fashion as for target selection; the `--query` flag is used to select full-length proteins from UniProt, while the optional `--uniprot_domain_regex` flag allows selection of individual domains with a regular expression string (Box 1). The returned UniProt data for each protein includes a list of associated PDB chains and their residue spans, and this information is used to select template structures, using the same method as for template selection from the PDB. Only structures solved by X-ray crystallography or NMR are selected, thus excluding computer-generated models available from the PDB. If the `--uniprot_domain_regex` flag is used, then templates are truncated at the start and end of the domain sequence.

Templates can also be defined manually. Manual specification of templates simply requires storing the sequences and arbitrary identifiers in a FASTA file, and the structures as PDB-format coordinate files with filenames matching the identifiers in the sequence file. The structure residues must also match those in the sequence file.

C. Template refinement

Unresolved template residues can optionally be modeled into template structures with the `loopmodel` subcommand, which employs a kinematic closure algorithm provided via the `loopmodel` tool of the Rosetta software suite [22, 23]. We expect that in certain cases, pre-building template loops with Rosetta `loopmodel` prior to the main modeling stage (with MODELLER) may result in improved model quality. Loop remodeling may fail for a small proportion of templates due to spatial constraints imposed by the original structure; the subsequent modeling step thus automatically uses the remodeled version of a template if available, but otherwise falls back to using the non-remodeled version. Furthermore, the Rosetta `loopmodel` program will not model missing residues at the termini of a structure—such residue spans are modeled in the subsequent stage.

D. Alignment and comparative modeling

In the modeling stage, structural models of the target sequence are generated from the template structures, with the goal of modeling the target in a variety of conformations that could be significantly populated under equilibrium conditions.

Modeling is performed using the `automodel` function of the MODELLER software package [24, 25] to rapidly generate a single model of the target sequence from each template structure. MODELLER uses simulated annealing cycles along with a minimal forcefield and spatial restraints—generally Gaussian interatomic probability densities extracted from the template structure with database-derived statistics determining the distribution width—to rapidly generate candidate structures of the target sequence from the provided template sequence [24, 25].

While MODELLER's `automodel` function can generate its own alignments automatically, a standalone function was preferable for reasons of programming convenience. **As such, we implemented pairwise alignment functionality using the BioPython `pairwise2` module [26]—which uses a dynamic programming algorithm—with the PAM 250 scoring matrix of Gonnet *et al.* [27], though other choices of scoring matrices available within the module can be selected.** The alignments are carried out with the `align` subcommand, prior to the modeling step which is carried out with the `build_models` subcommand. The `align` subcommand also writes a list of the sequence identities for each template to a text file, and this can be used to select models from a desired range of sequence identities. The `build_models`

subcommand and all subsequent pipeline functions have a `--template_seqid_cutoff` flag which can be used to select only models with sequence identities greater than the given value. We also note that alternative approaches could be used for the alignment stage. For example, multiple sequence alignment algorithms [28], allow alignments to be guided using sequence data from across the entire protein family of interest, while (multiple) structural alignment algorithms such as MODELLER's `salign` routine [24, 25], PRO-MALS3D [29], and Espresso and 3DCoffee [30, 31], can additionally exploit structural data. **Ensembler's** modular architecture facilitates the implementation of alternative alignment approaches, and we plan to implement some of these in future versions, to allow exploration of the influence of different alignment methods on model quality.

Models are output as PDB-format coordinate files. To minimize file storage requirements, **Ensembler** uses the Python `gzip` library to apply compression to all sizeable text files from the modeling stage onwards. The restraints used by MODELLER could potentially be used in alternative additional refinement schemes, and **Ensembler** thus provides a flag (`--write_modeller_restraints_file`) for optionally saving these restraints to file. This option is turned off by default, as the restraint files are relatively large (e.g. ~400 kB per model for protein kinase domain targets), and are not expected to be used by the majority of users.

At this time, the alignment and modeling functions cannot be used to model non-standard amino acids, though we plan to be able to provide this functionality in future versions. Note that the **Ensembler** functions for target and template selection only include standard amino acids which match the UniProt canonical isoform sequence, and thus any set of targets and templates selected this way should be compatible with the **Ensembler** alignment and modeling functions.

Filtering of nearly identical models

Because **Ensembler** treats individual chains from source PDB structures as individual templates, a number of models may be generated with very similar structures if these individual chains are nearly identical in conformation. For this reason, and also to allow users to select for high diversity if they so choose, **Ensembler** provides a way to filter out models that are very similar in RMSD. The `cluster` subcommand can thus be used to identify models which differ from other models in terms of RMSD distance by a user-specified cutoff. Clustering is performed using the regular spatial clustering algorithm [9], as implemented in the MSM-Builder Python library [18], which uses `mdtraj` [32] to calculate RMSD (for C_α atoms only) with a fast quaternion characteristic polynomial (QCP) [33–35] implementation. A minimum distance cutoff (which defaults to 0.6 Å) is used to retain only a single model per cluster.

E. Refinement of models and filtering of poor models by simulation

A number of refinement methods have been developed to help guide comparative modeling techniques toward more "native-like" and physically consistent conformations [36, 37]. Both short [37] and long [38] molecular dynamics simulations have been employed for this purpose. Here, we utilize short molecular dynamic simulations for two purposes: both to slightly relax the initial comparative models and to eliminate those comparative models that result in highly implausible conformations. This is especially critical here due to the inclusion of even very low sequence identity template structures. We stress that the limited refinement by molecular simulation here is primarily intended as initial relaxation and filtering stages, where implausible models might cause simulations to immediately fail, crash, or generate implausibly high energies or unstable dynamics. Exploration of conformational dynamics to derive MSMs, for example, will inevitably require orders of magnitude more simulation effort—very likely tens of microseconds to milliseconds of aggregate dynamics [8, 10].

Ensembler thus includes a refinement module, which uses short molecular dynamics simulations to refine the models built in the previous step. As well as improving model quality, this also prepares models for subsequent production MD simulation, including solvation with explicit water molecules, if desired.

Models are first subjected to energy minimization (using the L-BFGS algorithm [39], followed by a short molecular dynamics (MD) simulation with an implicit solvent representation. This is implemented using the OpenMM molecular simulation toolkit [2], chosen for its flexible Python API, and high performance GPU-accelerated simulation code. The simulation is run for a default of 100 ps, which in our example applications has been sufficient to filter out poor models (i.e. those with atomic overlaps unresolved by energy minimization, which result in an unstable simulation), as well as helping to relax model conformations. As discussed in the Results section, our example application of the **Ensembler** pipeline to the human tyrosine kinase family indicated that of the models which failed implicit solvent MD refinement, the vast majority failed within the first 1 ps of simulation.

The simulation protocol and default parameter values have been chosen to represent current "best practices" for the refinement simulations carried out here. As such, the simulation is performed using Langevin dynamics, with a default force field choice of Amber99SB-ILDN [40], along with a modified generalized Born solvent model [41] as implemented in the OpenMM package [2]. Any of the other force fields or implicit water models implemented in OpenMM can be specified using the `--ff` and `--water_model` flags respectively. The simulation length can also be controlled via the `--simlength` flag, and many other important simulation parameters can be controlled from either the API or CLI (via the `--api_params` flag). The default values are set as follows—timestep: 2 fs; temperature: 300 K; Langevin collision rate: 20 ps⁻¹; pH (used

by OpenMM for protonation state assignment): 7. We also draw attention to a recent paper which indicates that lower Langevin collision rates may result in faster phase space exploration [42].

For some studies, it may be useful to specify the protonation states of individual amino acids, rather than rely only on automatic protonation state assignment by OpenMM. The user can do this by listing the residue numbers and their protonation states in a configuration file (`manual_overrides.yaml`). The necessary formatting for the configuration file is specified in the software documentation, and a template file is written when initializing an **Ensembler** project. Protonation states are specified by naming the appropriate residue variant type in the force field, e.g. ‘ASH’ for an aspartic acid residue, as opposed to the aspartate base ‘ASP’. Any residues which do not have specific protonation states listed in the configuration file will have protonation states assigned automatically by OpenMM. Note that **Ensembler** currently only supports residue definitions provided by the forcefield definition files—it does not yet have the ability to derive new forcefield parameters for uncommon amino acids, cofactors, or ions not provided by the forcefield.

F. Solvation and NPT equilibration

While protein-only models may be sufficient for structural analysis or implicit solvent simulations, **Ensembler** also provides a stage for solvating models with explicit water and performing a round of explicit-solvent MD refinement/equilibration under isothermal-isobaric (NPT) conditions. The solvation step solvates each model for a given target with the same number of waters to facilitate the integration of data from multiple simulations, which is important for methods such as the construction of MSMs. The target number of waters is selected by first solvating each model with a specified padding distance (default: 10 Å), then taking a percentile value from the distribution (default: 68th percentile). This helps to prevent models with particularly long, extended loops—such as those arising from template structures with unresolved termini—from imposing very large box sizes on the entire set of models. The TIP3P water model [43] is used by default, but any of the other explicit water models available in OpenMM, such as TIP4P-Ew [44], can be specified using the `--water_model` flag. Models are resolvated with the target number of waters by first solvating with zero padding, then incrementally increasing the box size and resolvating until the target is exceeded, then finally deleting sufficient waters to match the target value. The explicit solvent MD simulation is also implemented using OpenMM, using the Amber99SB-ILDN force field [40] and TIP3P water [43] by default. The force field, water model, and simulation length can again be specified using the `--ff`, `--water_model`, and `--simlength` flags respectively. Further simulation parameters can be controlled via the API or via the CLI `--api_params` flag. Pressure control is performed with a Monte Carlo barostat as im-

plemented in OpenMM, with a default pressure of 1 atm and a period of 50 timesteps. The remaining simulation parameters have default values set to the same as for the implicit solvent MD refinement.

Model validation with MolProbity score

Ensembler provides a function for validating model quality and filtering models using MolProbity [45, 46]—a widely used tool for validation of protein models, which provides a numerical score derived from features such as steric clashes between atoms, bond geometry, Ramachandran angles, sidechain rotamer outliers, backbone deviations, and the presence of cis-peptides. This function is accessed via the `validate` subcommand, which for a given target will output a text file containing a list of model IDs sorted by validation score. The optional `--modeling_stage` flag specifies which of the three main Ensembler modeling stages to validate—the initial comparative modeling stage, the implicit MD refinement stage, or the explicit MD refinement stage. If this flag is not used, Ensembler defaults to selecting the latest stage for which models have been generated. The output text file can be used to filter models based on validation score, for example by using the `package_models` subcommand. Protein model validation is a challenging problem and an active area of research for many groups, including the developers of **Ensembler**. We plan to implement further validation methods in future versions of Ensembler.

Packaging

Ensembler provides a packaging module which can be used to prepare models for subsequent downstream use, such as the use of distributed or cluster computing resources for the generation of MSMs [8–10]. The `package_models` subcommand currently provides functions (specified via the `--package_for` flag) for compressing models in preparation for data transfer, or for organizing them with the appropriate directory and file structure for production simulation on the distributed computing platform Folding@home [4]. The module could easily be extended to add methods for preparing models for other purposes. For example, production simulations could alternatively be run using Copernicus [5, 6]—a framework for performing parallel adaptive MD simulations—or GPUGrid [7]—a distributing computing platform which relies on computational power voluntarily donated by the owners of nondedicated GPU-equipped computers.

An important use of the packaging stage is to filter models based on model quality. At the current time, the available filtering options are based on either target-template sequence identity or MolProbity validation score. The `package_models` subcommand includes optional flags for specifying a sequence identity cutoff (so that only models with a target-template sequence identity above the specified percentage are chosen), a MolProbity validation score

cutoff (to choose only models with lower validation scores, which indicate better model quality), or a MolProbity validation score percentile (to choose only models with validation scores lower than the value at the given percentile).

Models can also be exported into trajectory files for the purpose of performing structural analyses across model ensembles using tools like MDTraj [32]. This is done using the `mktraj` subcommand, which writes model coordinates for a given target to a Gromacs [47, 48] XTC format trajectory (chosen for its wide usage and data compression). Each frame in the trajectory represents a single model, and models are sorted in descending order of target-template sequence identity. Also output for each target are a PDB coordinate file (for use as a topology input file) and a CSV file containing model IDs (in the same order as the frames in the trajectory file) and other data such as target-template sequence identity. Using the `--modeling_stage` flag, models can be selected from any of three Ensembler modeling stages - after the initial comparative modeling stage, after implicit MD refinement, or after explicit MD refinement. If this flag is not used, Ensembler defaults to selecting the latest stage for which models have been generated.

We stress that, despite evidence suggesting that there is a correspondence between solution-state dynamics and structural diversity of related template proteins [16], all models—especially those derived from low sequence identity templates—are not necessarily representative of conformations thermally accessible to the template proteins of interest. Care must be exercised in the use and analysis of these models.

Other features

Tracking provenance information

To aid the user in tracking the provenance of each model, each pipeline function also outputs a metadata file, which helps to link data to the software version used to generate it (both **Ensembler** and its dependencies), and also provides timing and performance information, and other data such as hostname.

Rapidly modeling a single template

For users interested in simply using **Ensembler** to rapidly generate a set of models for a single template sequence, **Ensembler** provides a command-line tool `quickmodel`, which performs the entire pipeline for a single target with a small number of templates. For larger numbers of models (such as entire protein families), modeling time is greatly reduced by using the main modeling pipeline, which is parallelized via MPI, distributing computation across each model (or across each template, in the case of the loop reconstruction code), and scaling (in a “pleasantly parallel” manner) up to the number of models generated.

III. RESULTS

Modeling of all human tyrosine kinase catalytic domains

As a first application of **Ensembler**, we have built models for the human TK family. TKs (and protein kinases in general) play important roles in many cellular processes and are involved in a number of types of cancer [49]. For example, a translocation between the TK Abl1 and the pseudokinase Bcr is closely associated with chronic myelogenous leukemia [50], while mutations of Src are associated with colon, breast, prostate, lung, and pancreatic cancers [51]. Protein kinase domains are thought to have multiple accessible metastable conformation states, and much effort is directed at developing kinase inhibitor drugs which bind to and stabilize inactive conformations [52]. Kinases are thus a particularly interesting subject for study with MSM methods [53], and this approach stands to benefit greatly from the ability to exploit the full body of available genomic and structural data within the kinase family, e.g. by generating large numbers of starting configurations to be used in highly parallel MD simulation.

We selected all human TK domains annotated in UniProt as targets, and all available structures of protein kinase domains (of any species) as templates, using the commands shown in Box 1. This returned 93 target sequences and 4433 template structures, giving a total of 412,269 target-template pairs. The templates were derived from 3028 individual PDB entries and encompassed 23 different species, with 3634 template structures from human kinase constructs.

The resultant models are available as part of a supplementary dataset which can be downloaded from the Dryad Digital Repository (DOI: [10.5061/dryad.7fg32](https://doi.org/10.5061/dryad.7fg32)).

Figure 2 shows the number of PDB structures available for each of the 93 target TK domains. While a number of experimental structures are available for some TK domains, many TKs have few or no structures. **Ensembler** thus helps to overcome this unequal distribution of structural information when building protein models for simulation by exploiting homologous structural data from a wider range of protein kinase domains and species.

Ensembler modeling statistics

Crystallographic structures of kinase catalytic domains generally contain a significant number of missing residues (median 11, mean 14, standard deviation 13, max 102) due to the high mobility of several loops (Fig. 3, top), with a number of these missing spans being significant in length (median 5, mean 7, standard deviation 6, max 82; Fig. 3, bottom). To reduce the reliance on the MODELLER rapid model construction stage to reconstruct very long unresolved loops, unresolved template residues were first remodeled using the `loopmodel` subcommand. Out of 3666 templates with one or more missing residues, 3134 were successfully remodeled by the Rosetta loop modeling stage (with success de-

```
ensembl gather_targets --query 'family:"tyr protein kinase family" AND organism:"homo sapiens" AND reviewed:yes'
                        --uniprot_domain_regex 'Protein kinase(?:; truncated)?(?:; inactive)'
ensembl gather_templates --gather_from uniprot --query 'domain:"Protein kinase" AND reviewed:yes'
                        --uniprot_domain_regex 'Protein kinase(?:; truncated)?(?:; inactive)'
```

633 fined simply as program termination without error); most
634 remodeling failures were attributable to unsatisfiable spa-
635 tial constraints imposed by the original template structure.
636 There was some correlation between remodeling failures
637 and the number of missing residues (Fig. 3, top); templates
638 for which remodeling failed had a median of 20 missing
639 residues, compared to a median of 14 missing residues for
640 templates for which remodeling was successful (when ex-

Following loop remodeling, the **Ensembler** pipeline was performed up to and including the implicit solvent MD refinement stage, which completed with 389,067 (94%) surviving models across all TKs. To obtain statistics for the solvation stage without generating a sizeable amount of coordinate data (with solvated PDB coordinate files taking up about 0.9 MB each), the `solvate` subcommand was per-

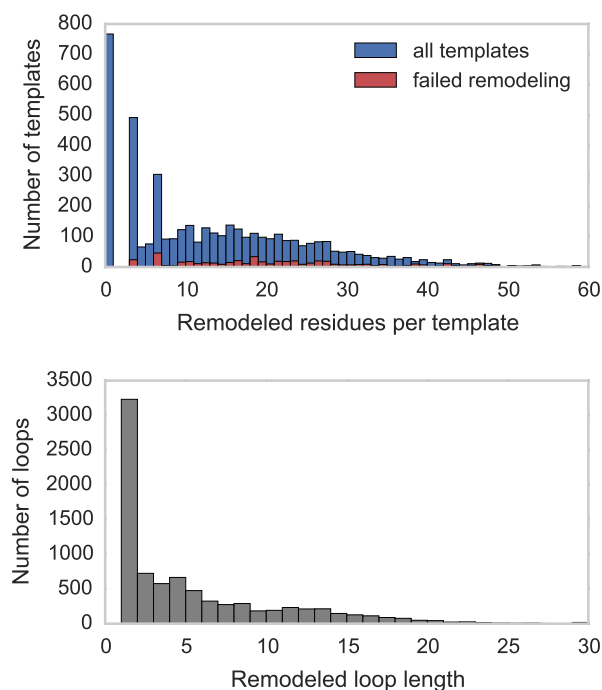


FIG. 3. Distributions for the number of missing residues in the TK templates. *Upper:* The number of missing residues per template, for all templates (blue) and for only those templates for which template remodeling with the `loopmodel` subcommand failed (red). Templates for which remodeling failed had a median of 20 missing residues, compared to a median of 14 missing residues for templates for which remodeling was successful. *Lower:* The number of residues in each missing loop, for all templates.

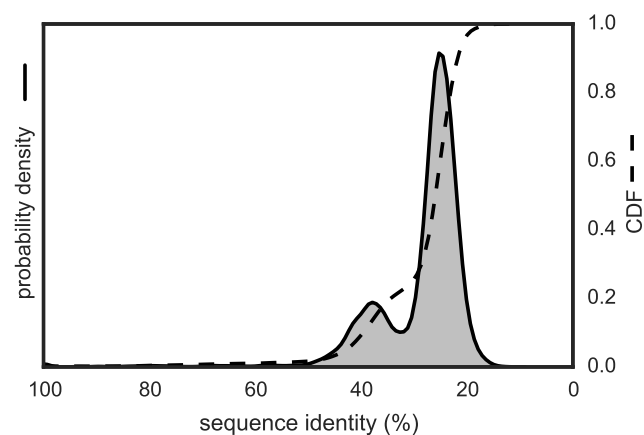


FIG. 4. Template-target sequence identity distribution for human tyrosine kinase catalytic domains. Sequence identities are calculated from all pairwise target-template alignments, where targets are human kinase catalytic domain sequences and templates are all kinase catalytic domains from any organism with structures in the PDB, as described in the text. A kernel density estimate of the target-template sequence identity probability density function is shown as a solid line with shaded region, while the corresponding cumulative distribution function is shown as a dashed line.

Evaluation of model quality and utility

All tyrosine kinases

formed for two representative individual kinases (*Src* and *Abl1*).

The number of models which survived each stage are shown in Fig. 1, indicating that the greatest attrition occurred during the modeling stage. The number of refined models for each target ranged from 4046 to 4289, with a median of 4185, mean of 4184, and standard deviation of 57. Fig. 1 also indicates the typical timing achieved on a cluster for each stage, showing that the `build_models` and `refine_implicit_md` stages are by far the most compute-intensive.

The files generated for each model (up to and including the implicit solvent MD refinement stage) totaled ~116 kB in size, totalling 0.5 GB per TK target or 42 GB for all 93 targets. The data generated per model breaks down as 39 kB for the output from the modeling stage (without saving MODELLER restraints files, which are about 397 kB per model) and 77 kB for the implicit solvent MD refinement stage.

To evaluate the variety of template sequence similarities relative to each target sequence, we calculated sequence identity distributions, as shown in Fig. 4. This suggests an intuitive division into three categories, with 355,712 models in the 0–35% sequence identity range, 51,330 models in the 35–55% range, and 5227 models in the 55–100% range. We then computed the RMSD distributions for the models created for each target (relative to the model derived from the template with highest sequence identity) Fig. 5, to assess the diversity of conformations captured by the modeling pipeline. Furthermore, to understand the influence of sequence identity on the conformational similarities of the resulting models, the RMSD distributions were stratified based on the three sequence identity categories described above. This analysis indicates that higher sequence identity templates result in models with lower RMSDs, while templates with remote sequence identities result in larger RMSDs on average, recapitulating the observation made years ago by Chothia and Lesk [54].

We also used the `ensampler validate` subcommand to subject the refined models to analysis with MolProbity [45, 46]. The MolProbity scores varied from 0.92 to 4.80, with a median of 3.84, a mean of 3.22, and a standard deviation of 1.07. Lower numbers represent better quality models. When stratified by the same sequence identity ranges as above, the mean scores were as follows: 2.96 (55–100%

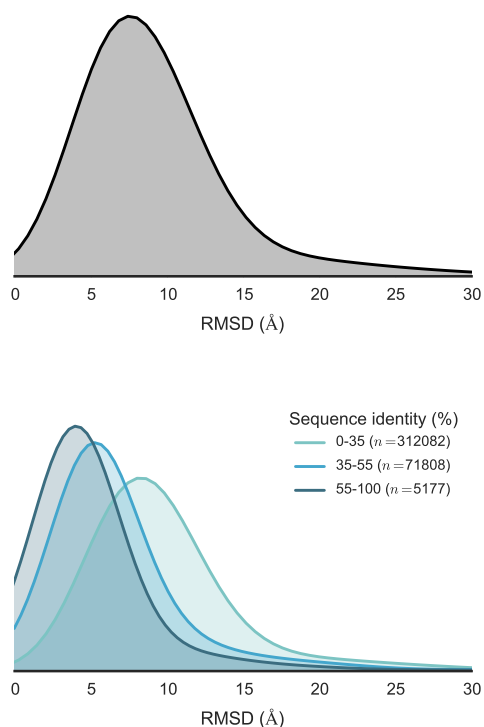


FIG. 5. Distribution of RMSDs to all TK catalytic domain models relative to the model derived from the highest sequence identity template. Distributions are built from data from all 93 TK domain targets. To better illustrate how conformational similarity depends on sequence identity, the lower plot illustrates the distributions as stratified into three sequence identity classes: high identity (55–100%), moderate identity (35–55%), and remote identity (0–35%). The plotted distributions have been smoothed using kernel density estimation.

sequence identity), 3.13 (35–55% sequence identity), 3.24 (0–35% sequence identity). This indicates that models with lower target-template sequence identities tend to be lower quality according to MolProbity analysis, as would be expected.

We also analyzed the potential energies of the models at the end of the implicit solvent MD refinement stage. These ranged from -14180 kT to -3160 kT, with a median of -9501 kT, mean of -9418 kT, and a standard deviation of 1198 kT (with a simulation temperature of 300 K). The distributions—stratified using the same sequence identity ranges again—are plotted in Fig. 6, indicating that higher sequence identity templates tend to result in slightly lower energy models. Of the 4973 models which failed to complete the implicit refinement MD stage, all except 9 failed within the first 1 ps of simulation.

Src and Abl

To provide a more detailed evaluation of the variety and utility of generated models, we have analyzed two specific

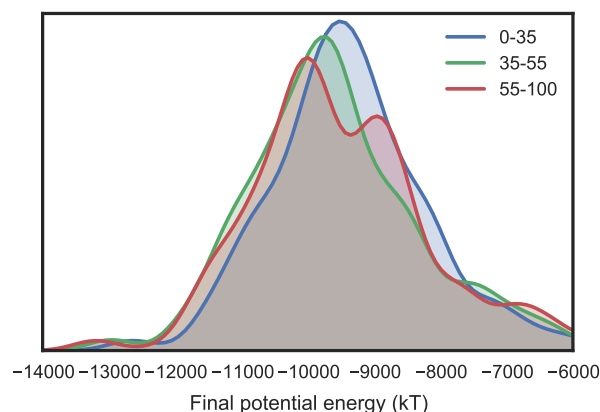


FIG. 6. Distribution of final energies from implicit solvent MD refinement of TK catalytic domain models. To illustrate how the energies are affected by sequence identity, the models are separated into three sequence identity classes: high identity (55–100%), moderate identity (35–55%), and remote identity (0–35%). The plotted distributions have been smoothed using kernel density estimation. Refinement simulations were carried out at the default temperature of 300 K.

TKs (*Src* and *Abl*) in depth. Due to their importance in cancer, these kinases have been the subject of numerous detailed structural and simulation studies. In terms of structural data, a large number of crystal structures have been solved (with or without ligands such as nucleotide substrate mimetics or small-molecule inhibitors), revealing a variety of conformations accessible to these kinases. A recent large-scale MSM study has also studied the activation pathway of *Src* [53], while a separate study employed biased sampling techniques to dissect the role of conformational changes in selectivity and affinity of imatinib recognition of *Abl* [55].

Visualizing model structural diversity. Fig. 7 shows a superposition of a set of representative models of *Src* and *Abl*. Models were first stratified into three ranges, based on the structure of the sequence identity distribution (Fig. 4), then subjected to RMSD-based *k*-medoids clustering (using the msmbuilder clustering package [18]) to pick three representative models from each sequence identity range. Each model is colored and given a transparency based on the sequence identity between the target and template sequence. The figure gives an idea of the variance present in the generated models. High sequence identity models (in opaque blue) tend to be quite structurally similar, with some variation in loops or changes in domain orientation.

The *Abl* renderings in Fig. 7 indicate one high sequence identity model with a long unstructured region at one of the termini, which was unresolved in the original template structure. While such models are not necessarily incorrect or undesirable, it is important to be aware of the effects they may have on production simulations performed under periodic boundary conditions, as long unstructured termini can be prone to interact with a protein's periodic image. Lower sequence identity models (in transparent white or red) in-

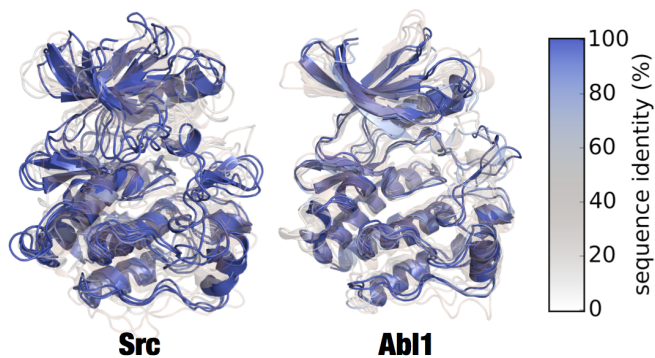


FIG. 7. Superposition of clustered models of Src and Abl1. Superposed renderings of nine models each for Src and Abl1, giving some indication the diversity of conformations generated by Ensembler. The models for each target were divided into three sequence identity ranges (as in Fig. 5), and RMSD-based k -medoids clustering was performed (using the msmbuilder clustering package [18]) to select three clusters from each. The models shown are the centroids of each cluster. Models are colored and given transparency based on their sequence identity, so that high sequence identity models are blue and opaque, while lower sequence identity models are transparent and red.

dicating much greater variation in all parts of the structure. We believe the mix of high and low sequence identity models to be particularly useful for methods such as MSM building, which require thorough sampling of the conformational landscape. The high sequence identity models could be considered to be the most likely to accurately represent true metastable states. Conversely, the lower sequence identity models could be expected to help push a simulation into regions of conformation space which might take intractably long to reach if starting a single metastable conformation.

Comparison with known biochemically relevant conformations. To evaluate the models of Src and Abl1 in the context of the published structural biology literature on functionally relevant conformations, we have focused on two residue pair distances thought to be important order parameters for the regulation of protein kinase domain activity. We use the residue numbering schemes for chicken Src (commonly employed in the literature even in reference to human Src) [56, 57] and human Abl1 isoform A [58–60] respectively; the exact numbering schemes are provided in Appendix 1.

Fig. 8 shows two structures of Src believed to represent inactive (PDB code: 2SRC) [56] and active (PDB code: 1Y57) [57] states. One notable feature which distinguishes the two structures is the transfer of an electrostatic interaction of E310 from R409 (in the inactive state) to K295 (in the active state), brought about by a rotation of the α C-helix. These three residues are also well conserved [61], and a number of experimental and simulation studies have suggested that this electrostatic switching process plays a role in a regulatory mechanism shared across the protein kinase family [53, 62, 63]. As such, we have plotted the distance between these two residue pairs for the Ensembler models for Src and Abl1, as well as Flt4 (Fig. 9). The models all

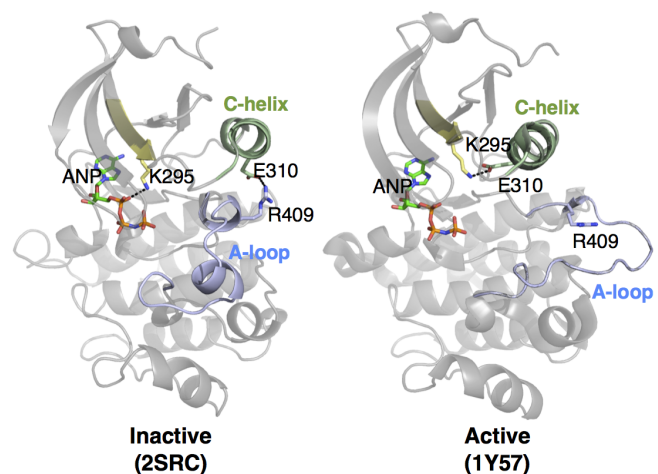


FIG. 8. Two structures of Src, indicating certain residues involved in activation. In the inactive state, E310 forms a salt bridge with R409. During activation, the α C-helix (green) moves and rotates, orienting E310 towards the ATP-binding site and allowing it to instead form a salt bridge with K295. This positions K295 in the appropriate position for catalysis. Note that ANP (phosphoaminophosphonic acid-adenylate ester; an analog of ATP) is only physically present in the 2SRC structure. To aid visualization of the active site in 1Y57, it has been included in the rendering by structurally aligning the surrounding homologous protein residues.

show strong coverage of regions in which either of the electrostatic interactions is fully formed (for models across all levels of target-template sequence identity), as well as a wide range of regions in-between (mainly models with low sequence identity). We thus expect that such a set of models, if used as starting configurations for highly parallel MD simulation, could greatly aid in sampling of functionally relevant conformational states. The Flt4 models (Fig. 8c) are of particular note, as there are no available crystal structures of the kinase domain of this TK protein (which is involved in tumor angiogenesis and lymphangiogenesis [64]), yet the models generated here include structural motifs which are conserved and of known importance to other proteins of the same family.

IV. AVAILABILITY AND FUTURE DIRECTIONS

Availability

The code for Ensembler is hosted on the collaborative open source software development platform GitHub (github.com/choderalab/enssembler). The latest release can be installed via the conda package manager for Python (conda.pydata.org), using the commands shown in Box 2. This will install all required dependencies, though a license must first be obtained for MODELLER. The optional dependencies Rosetta and MolProbity are not available through the conda package manager, and thus must be installed separately by the user according to the instructions for those packages.. The latest source can be downloaded from

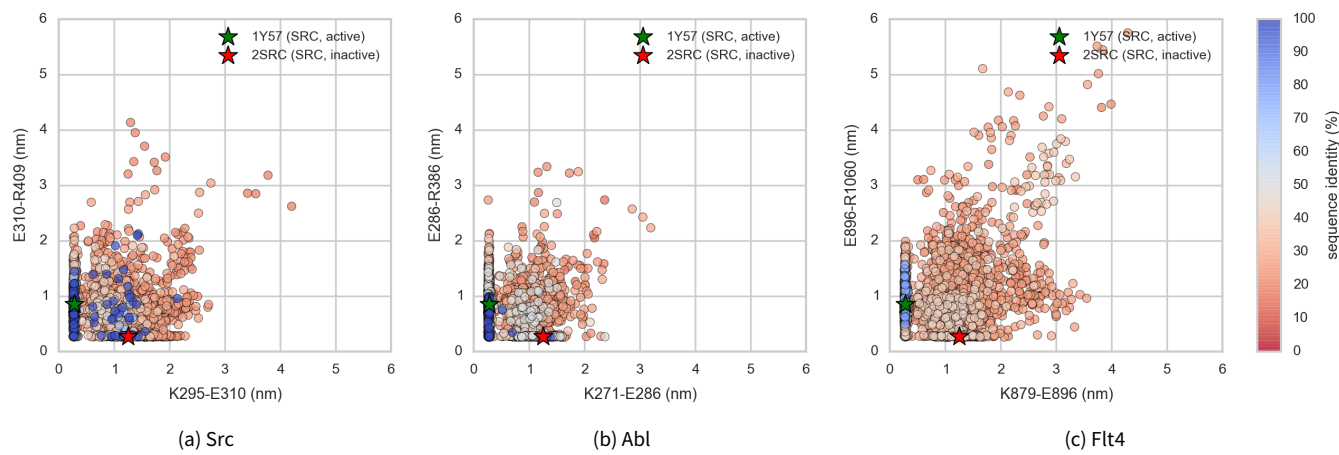


FIG. 9. Src, Abl, and Flt4 models projected onto the distances between two conserved residue pairs, colored by sequence identity. Two Src structures (PDB entries 1Y57 [57] and 2SRC [56]) are projected onto the plots for reference, representing active and inactive states respectively. These structures and the residue pairs analyzed here are depicted in Fig. 8. Distances are measured between the center of masses of the three terminal sidechain heavy atoms of each residue. The atom names for these atoms, according to the PDB coordinate files for both reference structures, are—Lys: NZ, CD, CE (ethylamine); Glu: OE1, CD, OE2 (carboxylate); Arg: NH1, CZ, NH2 (part of guanidine).

```
export KEY_MODELLE="Modeller license key"
conda config -add channels https://conda.anaconda.org/omnia
conda config -add channels https://conda.anaconda.org/salilab
conda install ensembler
```

Box 2. Ensembler installation using conda.

the GitHub repository, which also contains up-to-date instructions for building and installing the code. Documentation can be found at ensembl.readthedocs.org. A supplementary dataset can also be downloaded from the Dryad Digital Repository (DOI: [10.5061/dryad.7fg32](https://doi.org/10.5061/dryad.7fg32)). This contains the TK models described in the III section, general information on the targets and templates, plus a script and instructions for regenerating the same dataset.

Future Directions

We recognize that the current version of **Ensembler** has a number of limitations that bound its domain of applicability: support for nonnatural amino acids is currently rudimentary and confined to those already appearing in the forcefield; cofactors cannot currently be automatically modeled in; ligands, cofactors, and nonnatural amino acids cannot yet be automatically parameterized; protonation state assignment is limited to selection of the most populated state based on the intrinsic pK_a or user-specified overrides; the modeling of missing loops is rudimentary, relying on the subsequent dynamics for relaxation; there is not yet support for modeling of distinct domains from different templates, or the use of multiple templates to model a single domain. Nevertheless, there are a great number of use cases for this first version of an automated tool for simula-

tion preparation at the superfamily scale. To expand this domain of applicability, there are a number of obvious additions and improvements which we plan to implement in future versions of **Ensembler**.

Template remodeling. The lack of crystallographically-resolved regions of template structures presents a challenge to deriving structures from these templates by comparative modeling, especially in kinases, where loops are frequently unresolved. Improvements over the Rosetta-based strategy described here are likely possible, especially given the number of modeling failures observed in the template refinement stage (Fig. 3). An alternative approach could be to re-refine complete-chain template structures to the experimentally-derived electron density or scattering data deposited in the RCSB using methods capable of exploiting the scattering data and crystallographic symmetry [?]. Even if definitive placement of these unresolved regions is impossible, plausible locations constrained by weak scattering data and strong steric exclusion of crystallographic neighbors may provide a great deal of useful information, especially when combined with forcefield priors [?].

Comparative modeling. Comparative protein modeling can be approached in a number of different ways, with varying degrees of complexity. The comparative modeling stage of **Ensembler** currently uses MODELLER, but a number of excellent alternatives—such as RosettaCM [13] and the I-TASSER Suite [14]—can be added as user-selectable alternative choices. Additional options could be added to allow more expensive loop-modeling approaches to be employed to handle long insertions.

Protonation states. Some amino acids can exist in different protonation states, depending on pH and on their local environment. These protonation states can have important effects on biological processes. For example, long

timescale MD simulations have suggested that the conformation of the DFG motif of the TK Abl1—believed to be an important regulatory mechanism [65]—is controlled by protonation of the aspartate [66]. Currently, protonation states are assigned simply based on pH (a user-controllable parameter). At neutral pH, histidines have two protonation states which are approximately equally likely, and in this situation the selection is therefore made based on which state results in a better hydrogen bond. It would be highly desirable to instead use a method which assigns amino acid protonation states based on a rigorous assessment of the local environment. We thus plan to implement an interface and command-line function for assigning protonation states with MCCE2 [67–69], which uses electrostatics calculations combined with Monte Carlo sampling of side chain conformers to calculate pKa values.

Cofactors, structural ions, and ligands. Many proteins require the presence of various types of non-protein atoms and molecules for proper function, such as metal ions (e.g. Mg^{2+}), cofactors (e.g. ATP) or post-translational modifications (e.g. phosphorylation, methylation, glycosylation, etc.), and we thus plan for **Ensembler** to eventually have the capability to include such entities in the generated models. Binding sites for metal ions are frequently found in proteins, often playing a role in catalysis. For example, protein kinase domains contain two binding sites for divalent metal cations, and display significantly increased activity in the presence of Mg^{2+} [70], the divalent cation with highest concentration in mammalian cells. Metal ions are often not resolved in experimental structures of proteins, but by taking into account the full range of available structural data, it should be possible in many cases to include metal ions based on the structures of homologous proteins. We are careful to point out, however, that metal ion parameters in classical MD force fields have significant limitations, particularly in their interactions with proteins [71]. Cofactors and post-translational modifications are also often not fully resolved in experimental structures, and endogenous cofactors are frequently substituted with other molecules to facilitate experimental structural analysis. Future extensions to **Ensembler** could transfer cofactor and ion coordinates from homologous proteins in which these components are resolved.

Post-translationally modified amino acids and other molecules without forcefield parameters. A major challenge in the preparation of simulations of proteins of interest is the wide variety of post-translational modifications possible that are often functionally or structurally relevant. Often, forcefields lack parameters for these residues, or for other cofactors or ligands that might be vital to probing the relevant structural dynamics of these systems. While tools such as Antechamber [72, 73] can rapidly generate small molecule parameters in an automated manner, the parameterization of polymeric residues or covalently attached cofactors is much more challenging. In addition, small molecule forcefields are generally tied to specific corresponding protein and nucleic acid forcefields, meaning that different procedures may be needed to generate con-

sistent parameters.

Long insertions and deletions. Another limitation with the present version of **Ensembler** involves the treatment of members of a protein family with especially long residue insertions or deletions. For example, the set of all human protein kinase domains listed in UniProt have a median length of 265 residues (mean 277) and a standard deviation of 45, yet the minimum and maximum lengths are 102 and 801 respectively. The latter value corresponds to the protein kinase domain of serine/threonine-kinase *greatwall*, which includes a long insertion between the two main lobes of the catalytic domain. In principle, such insertions could be excluded from the generated models, though a number of questions would arise as to how best to approach this.

Markov state model (MSM) construction and model utility. We are actively utilizing **Ensembler**-generated models to seed the construction of Markov state models (MSMs) [8, 10]. While the observation that high sequence identity templates are likely to reflect accessible solution-phase conformations suggests that a number of these models occupy thermally accessible regions of configuration space [16], many models—especially those derived from very low sequence identity templates—are likely to be highly unrepresentative of conformations populated at equilibrium by the target protein. It is likely that even with hundreds of microseconds to milliseconds of aggregated dynamics, many of these poor quality models will remain trapped in inaccessible and irrelevant regions of configuration space. Standard approaches to MSM construction now employ an *ergodic trimming* step [18, 19] to prune away disconnected minor regions of configuration space, and this step is expected to be essential in the successful construction of MSMs using **Ensembler**-derived models.

Conclusion

We believe **Ensembler** to be an important first step toward enabling computational modeling and simulation of proteins on the scale of entire protein families, and suggest that it could likely prove useful for tasks beyond its original aim of providing diverse starting configurations for MD simulations. The code is open source and has been developed with extensibility in mind, in order to facilitate its customization for a wide range of potential uses by the wider scientific community.

V. ACKNOWLEDGMENTS

The authors are grateful to Robert McGibbon (Stanford) and Arien S. Rustenburg (MSKCC) for many excellent software engineering suggestions. The authors thank Nicholas M. Levinson (University of Minnesota), Markus A. Seeliger (Stony Brook), Diwakar Shukla (Stanford), and Avner Schlessinger (Mount Sinai) for helpful scientific feedback on modeling kinases. The authors are grateful to Benjamin Webb and Andrej Šali (UCSF) for help with the MODELLER

977 package, Peter Eastman and Vijay Pande (Stanford) for as-
 978 sistance with OpenMM, and Marilyn Gunner (CCNY) for as-
 979 sistance with MCCE2, **as well as the anonymous referees for**
 980 **constructive feedback on this manuscript.** All authors ac-
 981 knowledge support from the Sloan Kettering Institute. JDC,
 982 KAB, and DLP acknowledge partial support from NIH grant
 983 P30 CA008748. JDC and DLP also acknowledge the gener-
 984 ous support of a Louis V. Gerstner Young Investigator Award.
 985 KAB was also supported in part by Starr Foundation grant
 986 I8-A8-058. PBG acknowledges partial funding support from
 987 the Weill Cornell Graduate School of Medical Sciences.

-
- 988 [1] G. M. Lee and C. S. Craik, *Science* **324**, 213 (2009).
 989 [2] P. Eastman, M. S. Friedrichs, J. D. Chodera, R. J. Radmer, C. M.
 990 Bruns, J. P. Ku, K. A. Beauchamp, T. J. Lane, L.-P. Wang, D.
 991 Shukla, and V. S. Pande, *J. Chem. Theory Comput.* **9**, 461
 992 (2012).
 993 [3] R. Salomon-Ferrer, A. W. Götz, D. Poole, S. L. Grand, and R. C.
 994 Walker, *J. Chem. Theor. Comput.* **9**, 3878 (2013).
 995 [4] M. Shirts and V. S. Pande, *Science* **290**, 1903 (2000).
 996 [5] S. Pronk, P. Larsson, I. Pouya, G. R. Bowman, I. S. Haque, K.
 997 Beauchamp, B. Hess, V. S. Pande, P. M. Kasson, and E. Lin-
 998 dahl, in *Proceedings of 2011 International Conference for High*
 999 *Performance Computing, Networking, Storage and Analysis, SC*
 1000 *'11* (ACM, New York, NY, USA, 2011), pp. 60:1–60:10.
 1001 [6] S. Pronk, I. Pouya, M. Lundborg, G. Rotskoff, B. Wesén, P. M.
 1002 Kasson, and E. Lindahl, *Journal of Chemical Theory and Com-*
 1003 *putation* (2015).
 1004 [7] I. Buch, M. J. Harvey, T. Giorgino, D. P. Anderson, and G. De Fab-
 1005 ritiis, *Journal of Chemical Information and Modeling* **50**, 397
 1006 (2010).
 1007 [8] V. S. Pande, K. Beauchamp, and G. R. Bowman, *Methods* **52**,
 1008 99 (2010).
 1009 [9] J.-H. Prinz, H. Wu, M. Sarich, B. Keller, M. Fischbach, M. Held,
 1010 J. D. Chodera, C. Schütte, and F. Noé, *J. Chem. Phys.* **134**,
 1011 174105 (2011).
 1012 [10] J. D. Chodera and F. Noé, *Curr. Opin. Struct. Biol.* **25**, 135
 1013 (2014).
 1014 [11] J. Moul, K. Fidelis, A. Kryshchuk, T. Schwede, and A. Tra-
 1015 montano, *Proteins: Structure, Function, and Bioinformatics*
 1016 **82**, 1 (2014).
 1017 [12] D. Baker and A. Šali, *Science* **294**, 93 (2001).
 1018 [13] Y. Song, F. DiMaio, R. Y.-R. Wang, D. Kim, C. Miles, T. Brunette,
 1019 J. Thompson, and D. Baker, *Structure* **21**, 1735 (2013).
 1020 [14] J. Yang, R. Yan, A. Roy, D. Xu, J. Poisson, and Y. Zhang, *Nature*
 1021 *Methods* **12**, 7 (2015).
 1022 [15] M. W. van der Kamp, R. D. Schaeffer, A. L. Jonsson, A. D.
 1023 Scouras, A. M. Simms, R. D. Toofanny, N. C. Benson, P. C. An-
 1024 derson, E. D. Merkley, S. Rysavy, D. Bromley, D. A. C. Beck, and
 1025 V. Daggett, *Structure* **18**, 423 (2010).
 1026 [16] G. D. Friedland, N.-A. Lakomek, C. Griesinger, J. Meiler, and T.
 1027 Kortemme, *PLoS Comput. Biol.* **5**, e1000393 (2009).
 1028 [17] P. Weinkam, J. Pons, and A. Šali, *Proceedings of the National*
 1029 *Academy of Sciences of the United States of America* **109**,
 1030 4875 (2012).
 1031 [18] K. A. Beauchamp, G. R. Bowman, T. J. Lane, L. Maibaum, I. S.
 1032 Haque, and V. S. Pande, *Journal of Chemical Theory and Com-*
 1033 *putation* **7**, 3412 (2011).
 1034 [19] R. Scalco and A. Caflich, *The Journal of Physical Chemistry.*
 1035 *B* **115**, 6358 (2011).
 1036 [20] T. U. Consortium, *Nucleic Acids Research* **43**, D204 (2015).
 1037 [21] S. Velankar, J. M. Dana, J. Jacobsen, G. van Ginkel, P. J. Gane,
 1038 J. Luo, T. J. Oldfield, C. O'Donovan, M.-J. Martin, and G. J. Kley-
 1039 wegt, *Nucleic Acids Research* **41**, D483 (2013).
 1040 [22] B. Qian, S. Raman, R. Das, P. Bradley, A. J. McCoy, R. J. Read,
 1041 and D. Baker, *Nature* **450**, 259 (2007).
 1042 [23] C. Wang, P. Bradley, and D. Baker, *Journal of Molecular Biol-*
 1043 *ogy* **373**, 503 (2007).
 1044 [24] A. Fiser, R. K. G. Do, and A. Šali, *Protein Science* **9**, 1753 (2000).
 1045 [25] A. Šali and T. L. Blundell, *Journal of Molecular Biology* **234**,
 1046 779 (1993).
 1047 [26] P. J. A. Cock, T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, A.
 1048 Dalke, I. Friedberg, T. Hamelryck, F. Kauff, B. Wilczynski, and
 1049 M. J. L. de Hoon, *Bioinformatics (Oxford, England)* **25**, 1422
 1050 (2009).
 1051 [27] G. H. Gonnet, M. A. Cohen, and S. A. Brenner, *Science* **256**, 1443
 1052 (1992).
 1053 [28] J. D. Thompson, B. Linard, O. Lecompte, and O. Poch, *PLoS*
 1054 *ONE* **6**, e18093 (2011).
 1055 [29] J. Pei, B.-H. Kim, and N. V. Grishin, *Nucleic Acids Research* **36**,
 1056 2295 (2008).
 1057 [30] F. Armougom, S. Moretti, O. Poirot, S. Audic, P. Dumas, B.
 1058 Schaeli, V. Keduas, and C. Notredame, *Nucleic Acids Research*
 1059 **34**, W604 (2006).
 1060 [31] O. Poirot, K. Suhre, C. Abergel, E. O'Toole, and C. Notredame,
 1061 *Nucleic Acids Research* **32**, W37 (2004).
 1062 [32] R. T. McGibbon, K. A. Beauchamp, C. R. Schwantes, L.-P. Wang,
 1063 C. X. Hernández, M. P. Harrigan, T. J. Lane, J. M. Swails, and V. S.
 1064 Pande, *bioRxiv* (2014).
 1065 [33] D. L. Theobald, *Acta Cryst. A* **61**, 478 (2005).
 1066 [34] P. Liu, D. K. Agrafiotis, and D. L. Theobald, *J. Comput. Chem.*
 1067 **31**, 1561 (2010).
 1068 [35] P. Liu, D. K. Agrafiotis, and D. L. Theobald, *J. Comput. Chem.*
 1069 **32**, 185 (2011).
 1070 [36] J. L. MacCallum, A. Pérez, M. J. Schnieders, L. Hua, M. P. Jacob-
 1071 son, and K. A. Dill, *Proteins: Structure, Function, and Bioinform-*
 1072 *atics* **79**, 74 (2011).
 1073 [37] Y. Zhang, *Current Opinion in Structural Biology* **19**, 145 (2009).
 1074 [38] A. Raval, S. Piana, M. P. Eastwood, R. O. Dror, and D. E. Shaw,
 1075 *Proteins: Structure, Function, and Bioinformatics* **80**, 2071
 1076 (2012).
 1077 [39] D. C. Liu and J. Nocedal, *Mathematical Programming* **45**, 503
 1078 (1989).
 1079 [40] K. Lindorff-Larsen, S. P. anad Kim Palmo, P. Maragakis, J. L.
 1080 Klepeis, R. O. Dror, and D. E. Shaw, *Proteins* **78**, 1950 (2010).
 1081 [41] A. Onufriev, D. Bashford, and D. A. Case, *Proteins* **55**, 383
 1082 (2004).
 1083 [42] J. E. Basconi and M. R. Shirts, *Journal of Chemical Theory and*
 1084 *Computation* **9**, 2887 (2013).
 1085 [43] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey,
 1086 and M. L. Klein, *Journal of Chemical Physics* **79**, 926 (1983).
 1087 [44] H. W. Horn, W. C. Swope, J. W. Pitera, J. D. Madura, T. J.
 1088 Dick, G. L. Hura, and T. Head-Gordon, *The Journal of Chem-*
 1089 *ical Physics* **120**, 9665 (2004).
 1090 [45] I. W. Davis, A. Leaver-Fay, V. B. Chen, J. N. Block, G. J. Kapral, X.
 1091 Wang, L. W. Murray, W. B. Arendall, J. Snoeyink, J. S. Richard-
 1092 son, and D. C. Richardson, *Nucleic Acids Research* **35**, W375

- (2007).
- [46] V. B. Chen, W. B. Arendall, J. J. Headd, D. A. Keedy, R. M. Im-
mormino, G. J. Kapral, L. W. Murray, J. S. Richardson, and
D. C. Richardson, *Acta Crystallographica. Section D, Biologi-
cal Crystallography* **66**, 12 (2010).
- [47] H. J. C. Berendsen, D. van der Spoel, and R. van Drunen, *Comp.*
Phys. Comm. **91**, 43 (1995).
- [48] E. Lindahl, B. Hess, and D. van der Spoel, *J. Mol. Model.* **7**, 306
(2001).
- [49] D. S. Krause and R. A. Van Etten, *New England Journal of*
Medicine **353**, 172 (2005).
- [50] E. K. Greuber, P. Smith-Pearson, J. Wang, and A. M. Pender-
gast, *Nature Reviews Cancer* **13**, 559 (2013).
- [51] L. C. Kim, L. Song, and E. B. Haura, *Nature Reviews Clinical*
Oncology **6**, 587 (2009).
- [52] Y. Liu and N. S. Gray, *Nature Chemical Biology* **2**, 358 (2006).
- [53] D. Shukla, Y. Meng, B. Roux, and V. S. Pande, *Nature Commun.*
5, 3397 (2014).
- [54] C. Chothia and A. M. Lesk, *EMBO J.* **5**, 823 (1986).
- [55] Y.-L. Lin, Y. Meng, W. Jiang, and B. Roux, *Proc. Natl. Acad. Sci.*
USA **110**, 1664 (2013).
- [56] W. Xu, A. Doshi, M. Lei, M. J. Eck, and S. C. Harrison, *Molecular*
Cell **3**, 629 (1999).
- [57] S. W. Cowan-Jacob, G. Fendrich, P. W. Manley, W. Jahnke, D.
Fabbro, J. Liebetanz, and T. Meyer, *Structure* **13**, 861 (2005).
- [58] M. A. Young, N. P. Shah, L. H. Chao, M. Seeliger, Z. V. Milanov,
W. H. Biggs, D. K. Treiber, H. K. Patel, P. P. Zarrinkar, D. J. Lock-
hart, C. L. Sawyers, and J. Kuriyan, *Cancer Research* **66**, 1007
(2006).
- [59] S. W. Cowan-Jacob, G. Fendrich, A. Floersheimer, P. Furet, J.
Liebetanz, G. Rummel, P. Rheinberger, M. Centeleghe, D. Fab-
bro, and P. W. Manley, *Acta Crystallographica Section D: Bio-
logical Crystallography* **63**, 80 (2006).
- [60] N. M. Levinson, O. Kuchment, K. Shen, M. A. Young, M. Koldob-
skiy, M. Karplus, P. A. Cole, and J. Kuriyan, *PLoS Biol* **4**, e144
(2006).
- [61] N. Kannan and A. F. Neuwald, *Journal of Molecular Biology*
351, 956 (2005).
- [62] Z. H. Foda, Y. Shan, E. T. Kim, D. E. Shaw, and M. A. Seeliger,
Nature Communications **6**, 5939 (2015).
- [63] E. Ozkirimli, S. S. Yadav, W. T. Miller, and C. B. Post, *Protein*
Science : A Publication of the Protein Society **17**, 1871 (2008).
- [64] J.-L. Su, P.-C. Yang, J.-Y. Shih, C.-Y. Yang, L.-H. Wei, C.-Y. Hsieh,
C.-H. Chou, Y.-M. Jeng, M.-Y. Wang, K.-J. Chang, M.-C. Hung,
and M.-L. Kuo, *Cancer Cell* **9**, 209 (2006).
- [65] B. Nagar, O. Hantschel, M. A. Young, K. Scheffzek, D. Veach, W.
Bornmann, B. Clarkson, G. Superti-Furga, and J. Kuriyan, *Cell*
112, 859 (2003).
- [66] Y. Shan, M. A. Seeliger, M. P. Eastwood, F. Frank, H. Xu, M. Å.
Jensen, R. O. Dror, J. Kuriyan, and D. E. Shaw, *Proceedings of*
the National Academy of Sciences **106**, 139 (2009).
- [67] E. G. Alexov and M. R. Gunner, *Biophys. J.* **72**, 2075 (1997).
- [68] R. E. Georgescu, E. G. Alexov, and M. R. Gunner, *Biophys. J.* **83**,
1731 (2002).
- [69] Y. Song, J. Mao, and M. R. Gunner, *J. Comput. Chem.* **30**, 2231
(2009).
- [70] J. A. Adams and S. S. Taylor, *Protein Science* **2**, 2177 (1993).
- [71] S. F. Sousa, R. A. Fernandes, and M. J. Ramos, in *Kinetics*
and Dynamics: From Nano- to Bio-Scale, Vol. 12 of *Challenges*
and Advances in Computational Chemistry and Physics, edited
by P. a. D.-D. A. Paneth (Springer Science & Business Media,
Berlin, 2010), p. 530.
- [72] J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, and D. A.
Case, *J. Comput. Chem.* **25**, 1157 (2004).
- [73] J. Wang, W. Wang, P. A. Kollman, and D. A. Case, *J. Mol. Graph*
Model. **25**, 247260 (2006).

1160 Kinase catalytic domains are highlighted in red, and the conserved residues analyzed in the main text (Figs. 8 and 9) are
1161 highlighted with yellow background.

1162

1182

1183	SRC_HUMAN	1	MGSNKSXPKD	ASQRRRSLEP	AENVHGAGGG	AFPASQTPSK	PASADGHRGP	SAAFAPAAAAE	60
1184	SRC_CHICK	1	MGSSKSXPKD	PSQRRRSLEP	PDSTH--HG	GFPASQTPNK	TAAPDTHRTP	SRSFGTVATE	57
1185			***.*****	*****	:. . *	*****.*	*: * * *	*:*. .*:*	
1186	SRC_HUMAN	61	PKLFGGFNSS	DTVTSQRAG	PLAGGVITFV	ALYDYESRTE	TDLSFKKGER	LQIVNNTGD	120
1187	SRC_CHICK	58	PKLFGGFNTS	DTVTSQRAG	ALAGGVITFV	ALYDYESRTE	TDLSFKKGER	LQIVNNTGD	117
1188			*****:*	*****	*****	*****	*****	*****	
1189	SRC_HUMAN	121	WWLAHSLSTG	QTGYIPSNYV	APSDSIQAE	WYFGKITRE	SERLLNNAEN	PRGTFVRES	180
1190	SRC_CHICK	118	WWLAHSLTTG	QTGYIPSNYV	APSDSIQAE	WYFGKITRE	SERLLNPNEN	PRGTFVRES	177
1191			*****:**	*****	*****	*****	*****	*****	
1192	SRC_HUMAN	181	ETTKGAYCLS	VSDFDNAKGL	NVKHYKIRKL	DSGGFYITSR	TQFNSLQQLV	AYYSKHADGL	240
1193	SRC_CHICK	178	ETTKGAYCLS	VSDFDNAKGL	NVKHYKIRKL	DSGGFYITSR	TQFSSLQQLV	AYYSKHADGL	237
1194			*****	*****	*****	*****	*****	*****	
1195	SRC_HUMAN	241	CHRLTTVCPT	SKPQTQGLAK	DAWEIPRESL	RLEVKLGQGC	FGEVWMGTWN	GTTRVAIKTL	300
1196	SRC_CHICK	238	CHRLTNVCPT	SKPQTQGLAK	DAWEIPRESL	RLEVKLGQGC	FGEVWMGTWN	GTTRVAIKTL	297
1197			*****.****	*****	*****	*****	*****	*****	
1198	SRC_HUMAN	301	KPGTMSPEAF	LQEAQVMKKL	RHEKLVLQLYA	VVSEPIYIV	TEYMSKGSLL	DFLKGETGKY	360
1199	SRC_CHICK	298	KPGTMSPEAF	LQEAQVMKKL	RHEKLVLQLYA	VVSEPIYIV	TEYMSKGSLL	DFLKGEKGKY	357
1200			*****	*****	*****	*****	*****	*****	
1201	SRC_HUMAN	361	LRLPQLVDMA	AQIASGMAYV	ERMNYVHRDL	RAANILVGEN	LVCKVADFGL	ARLIEDNEYT	420
1202	SRC_CHICK	358	LRLPQLVDMA	AQIASGMAYV	ERMNYVHRDL	RAANILVGEN	LVCKVADFGL	ARLIEDNEYT	417
1203			*****	*****	*****	*****	*****	*****	
1204	SRC_HUMAN	421	ARQGAQFPIK	WTAPEAALYG	RFTIKSDVWS	FGILLTELTT	KGRVPYPGMV	NREVLQVER	480
1205	SRC_CHICK	418	ARQGAQFPIK	WTAPEAALYG	RFTIKSDVWS	FGILLTELTT	KGRVPYPGMV	NREVLQVER	477
1206			*****	*****	*****	*****	*****	*****	
1207	SRC_HUMAN	481	GYRMPCPPEC	PESLHDLMCQ	CWRKEPEERP	TFEYLQAFLE	DYFTSTEPQY	QPGENL	536
1208	SRC_CHICK	478	GYRMPCPPEC	PESLHDLMCQ	CWRKDPEERP	TFEYLQAFLE	DYFTSTEPQY	QPGENL	533
1209			*****	*****	*****	*****	*****	*****	