

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

Daniel L. Parton,¹ Patrick B. Grinaway,¹ and John D. Chodera^{1, *}

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

(Dated: March 21, 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 should run on all major operating systems, and has been tested on 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], GPUGrid [5], and Copernicus [6] 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 [7–9]. 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 [10, 11].

However, it remains difficult for researchers to exploit the full variety of available protein sequence and structural data in simulation studies, largely due to limitations in software architecture. For example, the set up 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 components and cosolvents), choice of simulation parameters (or parameterization schemes for components where

* Corresponding author; john.chodera@choderalab.org

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. 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, which could be used as starting configurations to aid sampling in MD simulations. 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 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 using these models as starting configurations for highly parallel MD simulations, we expect their structural diversity to greatly aid in sampling of conformational space.

We anticipate that the tool 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, the automation of simulation set up provides an excellent opportunity to make concrete certain "best practices", such as the

choice of simulation parameters. [JDC: Can we also add the URL of where to get the code and TK models here?]

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.

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.

Target selection and retrieval

The first stage entails the selection of a set of *target* protein sequences—the sequences the user is interested in generating simulation-ready structural models for. 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 with corresponding arbitrary identifiers.

The `ensembl` command-line tool allows targets to be selected from UniProt—a freely accessible resource for protein sequence and functional data (uniprot.org) [JDC: Cite one of these UniProt citations?]*—*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 `--query 'domain:"Protein kinase" AND taxonomy:9606 AND reviewed:yes'` would select all human 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 regular expression string to the `--domains` 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. To select only domains of the first three types, the following regular expression could be used: `^Protein kinase(?!;`



FIG. 1. Diagrammatic representation of the various stages of the Ensembler pipeline. The number of viable models surviving each stage of the pipeline are shown, either for all tyrosine kinases (*All TKs*) or representative individual kinases (*SRC* and *ABL*). In addition, the typical timing on a cluster (containing Intel Xeon E5-2665 2.4GHz hyperthreaded processors and NVIDIA GTX-680 or GTX-Titan GPUs) is reported to exemplify the resources required per model and for modeling the entire set of tyrosine kinases. Note that *CPU-h* denotes the number of hours consumed by the equivalent of a single hyperthread and *GPU-h* on a single GPU—parallel execution can reduce wall clock time nearly linearly.

truncated) (?!; inactive)'. In this case, 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.

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 some degree of homology 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 crystal unit cells with multiple asymmetric units) 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) (CITE: Velankar Nucleic Acids Res 2013), 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 `--domains` flag allows selection of individual domains with a regular expression string. 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 `--domains` 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 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.

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 [12, 13]. Because fewer loops need to be built during the subsequent target model-building stage, we find that prebuilding template loops tends to provide higher-quality models after completion of the **Ensembler** pipeline. [JDC: Should we cite our evidence for this with the TKs, or maybe tone back the claim a bit to say that it is possible this could make things easier?]

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 residues spans are modeled in the subsequent stage.

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 [14, 15] 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 [14, 15].

While Modeller’s automodel function can generate its own alignments automatically, we utilize the BioPython `pairwise2` module [CITE: BioPython]—which uses a dynamic programming algorithm—with the PAM 250 scoring matrix of Gonnet *et al.* [16], which we have empirically found to produce better quality alignments for purposes of high-throughput model building. [JDC: Can we give some evidence to support this claim, or soften our language on this?]

Models are output as PDB-format coordinate files. A list of all model identifiers sorted by sequence identity is also written to a text file, so that users may select models from the desired range of sequence identities. [JDC: I wonder if we should add an option “exclude models with sequence identities poorer than X”, since this could be something people might want to do.] 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.

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. A fast clustering scheme is used to identify models differing by a user-specified minimum RMSD. The `mdtraj` [17] Python library is used to calculate RMSD (for C_α atoms only) with a fast quaternion characteristic polynomial (QCP) [18–20] implementation, and the leader algorithm [JDC: Citation for leader algorithm?] is then used to populate clusters. A minimum distance cutoff (which defaults to 0.6 Å) is used to retain only a single model per cluster.

Refinement of models

This stage entails the use of molecular dynamics simulations to refine the models built in the previous step. This helps to improve model quality and also prepares models for subsequent production simulation, including solvation with explicit water molecules, if desired.

Models are first subjected to energy minimization (using the L-BFGS algorithm [21], 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. By default, the Amber99SB-ILDN force field [22] is used with a modified generalized Born solvent model [23] as implemented in the OpenMM package [2]. The **Ensembler** API allows the use of any of the other force fields implemented in OpenMM. The simulation is run for a default of 100 ps to filter out poor quality models (where atomic overlaps that cannot be resolved by energy minimization would cause the simulation to explode) and help relax models for subsequent production simulation. [JDC: What criteria were applied to filter out poor models? Do we only look for thrown exceptions or NaNs? Or do we use an energy filtering criteria too?] [DLP: We currently just filter out models which throw exceptions or NaNs.]

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, 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). [JDC: Would be useful to explain why we are doing this.] [DLP: Addressed.] 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. 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 [22] and TIP3P water [24] by default. Other force fields or water models such as TIP4P-Ew [?] can be specified via the **Ensembler** API. [JDC: We should allow other water models in OpenMM too, such as TIP4P-Ew?] [DLP: I forgot to mention this in the text previously - any of the OpenMM force fields can be chosen via the API.]

I've updated the text accordingly. Is this functionality sufficient? I guess it's ok to leave ff choice as an "advanced" feature which requires use of the API? Otherwise I could add a `--water_model` flag to the CLI, for example. [JDC: Would be good to let at least water model and forcefield name be selectable via the CLI!]

Packaging

Ensembler provides a packaging module which can be used to compress models in preparation for data transfer, or to prepare models with the appropriate directory and file structure for subsequent production simulations on the distributed computing platform Folding@home (CITE: F@H). The module could be easily extended to add methods for preparing models for use with other software, such as the Copernicus platform for running automated, distributed MD simulations. [JDC: Is there a way we can make this more generally useful to others? For example, is there a different system they might want to use, such as Copernicus?] [DLP: addressed]

Provenance

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 all 90 human tyrosine kinase (TK) domains listed in UniProt. [JDC: What query was used to yield this set?] [JDC: Is there a complete list of these somewhere? Maybe reference supplementary data?] TKs (and protein kinases

in general) play important roles in many cellular processes and are involved in a number of types of cancer. [JDC: CITE] For example, mutations of Src are associated with colon, breast, and prostate cancer [CITE: Src cancer involvement], while a translocation between the TK Abl1 and the pseudokinase Bcr is closely associated with chronic myelogenous leukemia [CITE: Abl1 cancer involvement]. Protein kinase domains are thought to have multiple accessible metastable conformation states, with a single active conformation, and much effort is directed at developing kinase inhibitor drugs which bind to and stabilize inactive conformations [CITE: Lee and Craik Science 2009]. [JDC: Lee and Craik do not discuss kinases, I don't believe; you'll have to find an accurate reference on kinase conformations.] Kinases are thus a particularly interesting subject for study with MSM methods [CITE: recent kinase MSM papers], 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 available structures of protein kinase domains (of any species) as templates, for a total of 4433 (398,970 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.

Ensembler modeling statistics

Crystallographic structures of kinase catalytic domains generally contain a significant number of missing residues (median 11, standard deviation 13, max 102) due to the high mobility of several loops (Fig. 2, top), with a number of these missing spans being significant in length (Fig. 2, bottom). [JDC: Can you add statistics (median, stddev, max) for loop lengths too?] 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 defined simply as program termination with out error); most remodeling failures were attributable to unsatisfiable spatial constraints imposed by the original template structure. There was some correlation between remodeling failures and the number of missing residues (Fig. 2, top); 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.

Following loop remodeling, the **Ensembler** pipeline was performed up to and including the implicit solvent MD refinement stage, which completed with 373,513 (94%) surviving models. To obtain statistics for the solvation stage without generating a sizeable amount of coordinate data, the `solvate` subcommand was performed for two repre-

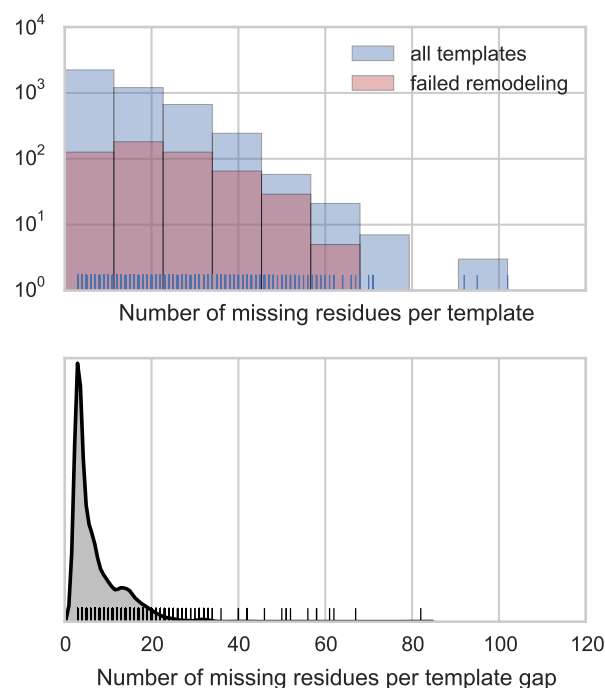


FIG. 2. Distributions for the number of missing residues in the TK templates.

The upper plot shows the distribution of the total 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). The raw data points for all templates are shown as a rug plot. The lower plot shows the distribution of the number of residues per template gap, normalized and smoothed using kernel density estimation. The raw data points are shown as a rug plot. [JDC: Some ideas for cleaning this up: Either the tick marks are being misrendered in that they are not taller if there are multiple data points with the same number or the data is really funky, since I would expect there to be a few examples in some bins. Also, is there a big drawback to making the top histogram bin size unity, since the values are integral? I don't think transparency is needed for the histogram bars either. We can also ditch the semilogy axis. I would also make the x-axes for the top and bottom plots different, since the data ranges are different. I'd see if a histogram with unit bin size might be more appropriate for the bottom plot as well—the KDE just doesn't feel right for this kind of data, since we are trying to report exact statistics from a specific example rather than estimate a general density for problems of this sort. Finally, I like “remodeled loop length” or “missing loop length” much better than “Number of missing residues per template gap”, which seems unnecessarily verbose.]

sentative 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 4005 to 4248, with a median of 4160 and standard deviation of 60. 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.

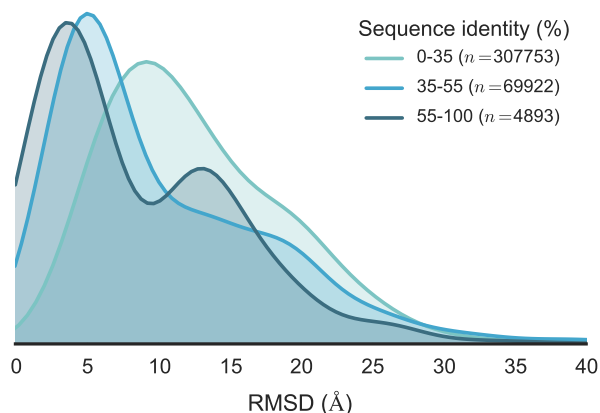


FIG. 3. Distribution of RMSDs to all tyrosine kinase catalytic domain models relative to model derived from highest sequence identity template. Distributions are averaged over all 90 tyrosine kinase catalytic domains. [JDC: Is this accurate?] To better illustrate how conformational similarity depends on sequence identity, these were 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. [JDC: Are these computed for the templates or the resulting models?] [JDC: Can we also show the overall distribution without stratification (e.g. in grey in a separate panel)?]

Each model generated about 116 KB of file data (up to and including the implicit solvent MD refinement stage), totalling 0.5 GB per TK target or 41 GB for all 90 TKs. The data generated per model breaks down as 39 kB for the output from the modeling stage (without saving Modeller restraints to file) and 77 kB for the implicit solvent MD refinement stage. [JDC: Can you rework this paragraph to specify total space requirements without saving MODELLER restraint files (since this is now the default), plus how much each restraint file adds if they are retained as well? I took a stab at it, but couldn't figure out how the numbers worked out.]

Evaluation of model quality and utility

To evaluate the diversity of conformations captured for each target sequence, we first computed the RMSD distributions for all models for each target (relative to the model derived from the highest-identity template) are shown in Fig. 3. To better understand the influence of sequence identity on the conformational similarities of resulting models, the sequence identities were stratified based on the sequence identity distribution plotted in Fig. 4, which suggests an intuitive division into three categories, with 307,753 models in the 0–35% sequence identity range, 69,922 models in the 35–55% range, and 4893 models in the 55–100% range. It is clear that higher sequence identity templates result in models with lower RMSDs, while templates with remote se-

quence identities result in larger RMSDs on average.

To provide a more complete evaluation of the models generated, we have analyzed two example TKs (*Src* and *Abl1*) in detail. Due to their importance in cancer, as outlined above, these kinases have been the subject of numerous studies, encompassing many different methodologies. In terms of structural data, a large number of crystal structures have been solved (with or without ligands such as nucleotide substrate or inhibitor drugs), showing the kinases in a number of different conformations. These two kinases are thus also interesting targets for MSM studies, with one recent study focusing on modeling the states which constitute the activation pathway of *Src* [25]. Fig. 5 shows a superposition of a set of representative models of *Src* and *Abl1*. Models were first stratified into three ranges, based on the structure of the sequence identity distribution (Fig. 4), then subjected to *k*-medoids clustering to pick three representative models from each sequence identity range. [JDC: Explain how *k*-medoids clustering was done either here or in figure caption.] 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 *Abl1* renderings 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) indicate 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.

To evaluate the models of *Src* and *Abl1* in the context of the published structural biology literature, we have focused on two residue pair distances thought to be important for the regulation of protein kinase domains. We use the residue numbering schemes for chicken *Src* (which is commonly used in the literature even in reference to human *Src*) [CITE: 2SRC, 1Y57] and human *Abl1* isoform A [CITE: 2F4J, 2HYI, 2GIT] respectively; the exact numbering schemes are provided in Supporting Information S1. Fig. 6 shows two structures of *Src* believed to represent inactive (PDB code: 2SRC) [CITE: 2SRC] and active (PDB code: 1Y57) [CITE: 1Y57] states. One notable feature which distinguishes the two structures is the transfer of an electrostatic interaction of E310 from

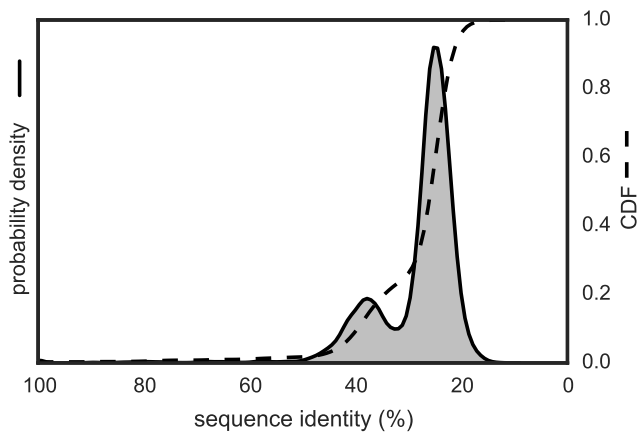


FIG. 4. Template-target sequence identity distribution for human tyrosine kinase catalytic domains. Sequence identities are calculated from all pairwise target-template alignments, where targets are human kinase catalytic domain sequences and templates are all kinase catalytic domains from any organism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

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 [CITE Kannan Neuwald JMB 2005], 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 [25] [CITE Foda Shan Seeliger Src Nat Commun 2015; Ozkirimli Post Prot Sci 2008]. As such, we have projected the **Ensembler** models for *Src* and *Abl1* onto a space consisting of the distances between these two residue pairs (Fig. 7). The models show strong coverage of regions in which either of the electrostatic interactions is formed, as well as a wide range of regions inbetween. 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 the activation process.

IV. AVAILABILITY AND FUTURE DIRECTIONS

Availability

The code for **Ensembler** is hosted on the collaborative open source software development platform GitHub, <http://github.com/choderalab/ensembler>. The latest release of **Ensembler** can be installed via the conda package manager for Python [<http://conda.pydata.org/>]:

```
# conda config -add channels https://conda.binstar.org/omnia
```

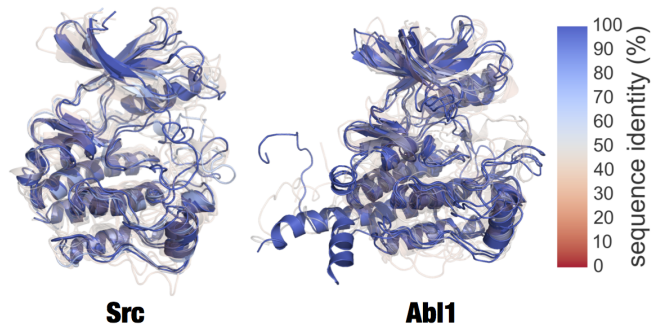


FIG. 5. Superposition of clustered models of Src and Abl1. Superposed renderings of nine models each for Src and Abl1, [JDC: Src and Abl, or Src and Abl1? The description should match the captions above.] [DLP: Addressed. Using Abl1, as this is the HGNC recommended symbol.] giving some indication the diversity of conformations generated by Ensembler. The models for each target were divided into three sequence identity ranges (as in Fig. 3), and RMSD-based k -medoids clustering was performed 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.

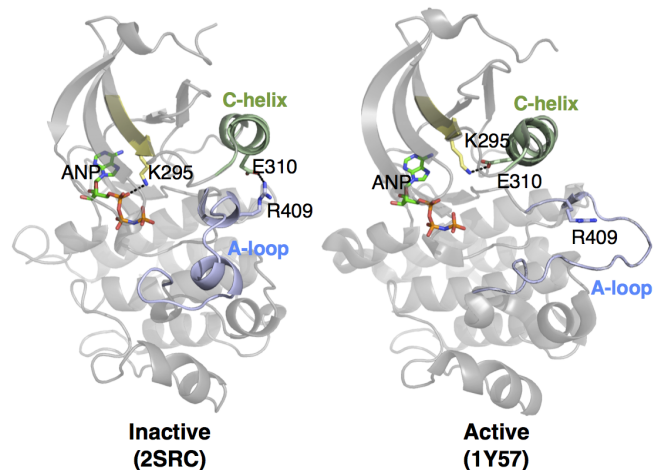


FIG. 6. Two structures of Src, indicating certain residues involved in activation. In the inactive state, E310 forms a salt bridge with R409. During activation, the α C-helix (green) moves and rotates, orienting E310 towards the ATP-binding site and allowing it to instead form a salt bridge with K295. This positions K295 in the appropriate position for catalysis.

```
# conda ensembler
```

This will install all dependencies except for Modeller and Rosetta, which are not available through the conda package manager, and thus must be installed separately by the user. The latest source can be downloaded from the GitHub repository, which also contains up-to-date instructions for building and installing the code. [JDC: Where does documentation reside? And example inputs for generating the results in this paper?]

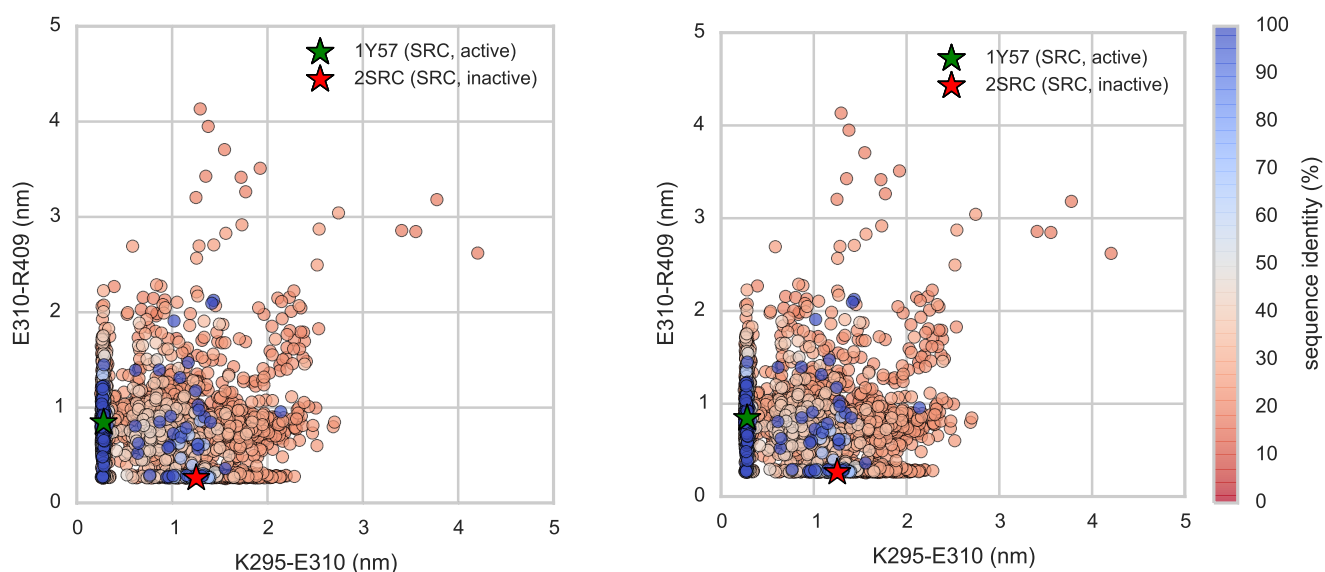


FIG. 7. Src and Abl1 models projected onto the distances between two conserved residue pairs, colored by sequence identity. Two Src structures (PDB entries 1Y57 [CITE] and 2SRC [CITE]) are projected onto the plots for reference, representing active and inactive states respectively. These structures and the residue pairs analyzed here are depicted in Fig. 6. 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).

Future Directions

[JDC: In the Discussion, let's be sure to talk about the limitations and what could be improved or added in the future. For example, we don't yet handle counterions (e.g. structural Zn^{2+}), prosthetic groups (e.g. heme), or cofactors (e.g. ATP) yet. We don't handle post-translational modifications either (such as phosphorylation, methylation, glycosylation, etc.). It's a good idea to suggest that this is an important first step toward enabling superfamily- and genomics-scale modeling, but there's a lot of work yet to be done.]

Comparative protein modeling and MD simulation set-up can be approached in a number of different ways, with varying degrees of complexity, and there are a number of obvious additions and improvements which we plan to implement in future versions of **Ensembler**.

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[CITE: Abl1 DFG flip evidence]—is controlled by protonation of the aspartate [CITE: Shan Shaw Proton-dependent switch Abl1 PNAS 2009]. 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 [26–28], which uses electrostatics calculations combined with Monte Carlo sampling of side chain conformers to calculate pKa values. [JDC: I think we may want to consider doing that at this stage. Let's discuss.]

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+} [CITE: Adams Taylor Protein Sci 1993], 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 [CITE: Sousa Ramos chapter 11 of Kinetics and Dynamics: From Nano- to Bio-Scale, Springer, 2010]. 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. Again, **Ensembler** could exploit structural data from a set of homologous proteins to model in these molecules, although there will be likely be a number of challenges to overcome in the design and implementation of such functionality.

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 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 *great-wall*, 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.

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 Kyle A. Beauchamp (MSKCC), Robert McGibbon (Stanford), Arien S. Rustenburg (MSKCC) for many excellent software engineering suggestions. The authors thank Sonya M. Hanson (MSKCC), 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 package, Peter Eastman and Vijay Pande (Stanford) for assistance with OpenMM, and Marilyn Gunner (CCNY) for assistance with MCCE2. DLP and this work was supported in part by the generous support of a Louis V. Gerstner Young Investigator Award. [\[Add PBG support statement.\]](#)

- [1] G. M. Lee and C. S. Craik, *Science* **324**, 213 (2009).
- [2] P. Eastman, M. S. Friedrichs, J. D. Chodera, R. J. Radmer, C. M. Bruns, J. P. Ku, K. A. Beauchamp, T. J. Lane, L.-P. Wang, D. Shukla, and V. S. Pande, *J. Chem. Theory Comput.* **9**, 461 (2012).
- [3] R. Salomon-Ferrer, A. W. Götz, D. Poole, S. L. Grand, and R. C. Walker, *J. Chem. Theor. Comput.* **9**, 3878 (2013).
- [4] M. Shirts and V. S. Pande, *Science* **1903** (2000).
- [5] I. Buch, M. J. Harvey, T. Giorgino, D. P. Anderson, and G. De Fabritiis, *Journal of Chemical Information and Modeling* **50**, 397 (2010).
- [6] S. Pronk, P. Larsson, I. Pouya, G. R. Bowman, I. S. Haque, K. Beauchamp, B. Hess, V. S. Pande, P. M. Kasson, and E. Lindahl, in *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis, SC '11* (ACM, New York, NY, USA, 2011), pp. 60:1–60:10.
- [7] V. S. Pande, K. Beauchamp, and G. R. Bowman, *Methods* **52**, 99 (2010).
- [8] J.-H. Prinz, H. Wu, M. Sarich, B. Keller, M. Fischbach, M. Held, J. D. Chodera, C. Schütte, and F. Noé, *J. Chem. Phys.* **134**, 174105 (2011).
- [9] J. D. Chodera and F. Noé, *Curr. Opin. Struct. Biol.* **25**, 135 (2014).
- [10] J. Moulton, K. Fidelis, A. Kryshchuk, T. Schwede, and A. Tramontano, *Proteins: Structure, Function, and Bioinformatics* **82**, 1 (2014).
- [11] D. Baker and A. Šali, *Science* **294**, 93 (2001).
- [12] B. Qian, S. Raman, R. Das, P. Bradley, A. J. McCoy, R. J. Read, and D. Baker, *Nature* **450**, 259 (2007).
- [13] C. Wang, P. Bradley, and D. Baker, *Journal of Molecular Biology* **373**, 503 (2007).
- [14] A. a. Fiser, R. K. G. Do, and A. Šali, *Protein Science* **9**, 1753 (2000).
- [15] A. Šali and T. L. Blundell, *Journal of Molecular Biology* **234**, 779 (1993).
- [16] G. H. Gonnet, M. A. Cohen, and S. A. Brenner, *Science* **256**, 1443 (1992).
- [17] R. T. McGibbon, K. A. Beauchamp, C. R. Schwantes, L.-P. Wang, C. X. Hernández, M. P. Harrigan, T. J. Lane, J. M. Swails, and V. S. Pande, *bioRxiv* (2014).
- [18] D. L. Theobald, *Acta Cryst. A* **61**, 478 (2005).
- [19] P. Liu, D. K. Agrafiotis, and D. L. Theobald, *J. Comput. Chem.* **31**, 1561 (2010).
- [20] P. Liu, D. K. Agrafiotis, and D. L. Theobald, *J. Comput. Chem.* **32**, 185 (2011).
- [21] D. C. Liu and J. Nocedal, *Mathematical Programming* **45**, 503 (1989).
- [22] K. Lindorff-Larsen, S. P. anad Kim Palmo, P. Maragakis, J. L. Klepeis, R. O. Dror, and D. E. Shaw, *Proteins* **78**, 1950 (2010).
- [23] A. Onufriev, D. Bashford, and D. A. Case, *Proteins* **55**, 383 (2004).
- [24] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, *Journal of Chemical Physics* **79**, 926 (1983).
- [25] D. Shukla, Y. Meng, B. Roux, and V. S. Pande, *Nature Commun.* **5**, 3397 (2014).
- [26] E. G. Alexov and M. R. Gunner, *Biophys. J.* **72**, 2075 (1997).
- [27] R. E. Georgescu, E. G. Alexov, and M. R. Gunner, *Biophys. J.* **83**, 1731 (2002).
- [28] Y. Song, J. Mao, and M. R. Gunner, *J. Comput. Chem.* **30**, 2231 (2009).

Appendix 1: Sequences and residue numbering schemes for Src and Abl1

Kinase catalytic domains are highlighted in red, and the conserved residues analyzed in the main text (Figs. 6 and 7) are highlighted with yellow background. [JDC: Very cool! Did you write a script to generate this, or did you have to do this by hand?]

Human Abl1 sequence

1	MLEICLKLVG	CKSKKGLSSS	SSCYLEEALQ	RPVASDFEPQ	GLSEAAWNS	KENLLAGPSE	60
61	NDPNLFVALY	DFVASGDNTL	SITKGEKLRV	LGYNHNGEWC	EAQTKNGQGW	VPSNYITPVN	120
121	SLEKHSWYHG	PVSRNAAEYL	LSSGINGSFL	VRESESSPGQ	RSISLRYEGR	VYHYRINTAS	180
181	DGKLYVSSSES	RFNTLAELVH	HHSTVADGLI	TTLHYPAPKR	NKPTVYGVSP	NYDKWEMERT	240
241	DITMKHKLGG	GQYGEVYEGV	WKYSLTVAV	KTLEDTMEV	EEFLKEAAVM	KEIKHPNLVQ	300
301	LLGVCTREPP	FYIITEFMTY	GNLLDYLREC	NRQEVNAVVL	LYMATQISSA	MEYLEKKNFI	360
361	HRDLAARNCL	VGENHLVKVA	DFGLSRMLTG	DTYTAHAGAK	FPIKWTAPES	LAYNKFSIKS	420
421	DVWAFGVLLW	EIATYGMSPY	PGIDLSQVYE	LLEKDYRMER	PEGCPEKVYE	LMRACWQWNP	480
481	SDRPSFAEIH	QAFETMFQES	SISDEVEKEL	GKQGVRGAVS	TLLQAPELPT	KTRTSRRAAE	540
541	HRDITDVPDM	PHSKGQGQES	PLDHEPAVSP	LLPRKERGPP	EGGLNEDERL	LPKDKKTNLF	600
601	SALIKKKKKT	APTPPKRSSS	FREMDGQPER	RGAGEEEGRD	ISNGALAFPT	LDTADPAKSP	660
661	KPSNGAGVPN	GALRESGGSG	FRSPHLWKKS	STLTSSRLAT	GEEEGGGSSS	KRFLRSCSAS	720
721	CVPHGAKDTE	WRSVTLPRLD	QSTGRQFDSS	TFGGHKSEKP	ALPRKRAGEN	RSDQVTRGTV	780
781	TPPPRLVKKN	EEAADEVFKD	IMESSPGSSP	PNLTPKPLRR	QVTVPASGL	PHKEEAGKGS	840
841	ALGTPAAAEF	VTPTSAGSG	APGGTSKGPA	EESRVRHKKH	SSESPGRDKG	KLSRLKPAPP	900
901	PPPAASAGKA	GGKPSQSPSQ	EAAGEAVLGA	KTATSLVDA	VNSDAAKPSQ	PGEGLKPPVL	960
961	PATPKPQSAK	PSGTPISPAP	VPSTLPSASS	ALAGDQPSST	AFIPLISTRV	SLRKRTPPE	1020
1021	RIASGAITKG	VVLDESTALC	LAISRNSEQM	ASHSAVLEAG	KNLYTFCVSY	VDSIQQMRNK	1080
1081	FAFREAINKL	ENNLRELQIC	PATAGSGPAA	TQDFSLLSS	VKEISDIVQR		1130

Sequences for human and chicken Src, aligned using Clustal Omega

SRC_HUMAN	1	MGSNKS	SKPKD	ASQRRRSLEP	AENVHGAGGG	AFPASQTPSK	PASADGHRGP	SAAFAPAAAE	60
SRC_CHICK	1	MGSSKSKPKD	PSQRRRSLEP	PDSTH--HG	GFPASQTPNK	TAAPDTHRTP	SRSFGTVATE		57
		.**	*****	:..*	* .*****.*	*: * * *	* :*. .*:*		
SRC_HUMAN	61	PKLFGGFNNS	DTVTSPQRAG	PLAGGVTTTFV	ALYDYESRTE	TDLSEFKKGER	LQIVNNTGEGD		120
SRC_CHICK	58	PKLFGGFNTS	DTVTSPQRAG	ALAGGVTTTFV	ALYDYESRTE	TDLSEFKKGER	LQIVNNTGEGD		117
		*****.*	*****	*****	*****	*****	*****		
SRC_HUMAN	121	WWLAHSLSTG	QTGYIPSNYV	APSDSIQAE	WYFGKITRRE	SERLLLNAEN	PRGTFLVRES		180
SRC_CHICK	118	WWLAHSLTTG	QTGYIPSNYV	APSDSIQAE	WYFGKITRRE	SERLLLNPN	PRGTFLVRES		177
		*****.*	*****	*****	*****	*****	*****		
SRC_HUMAN	181	ETTKGAYCLS	VSDFDNAKGL	NVKHYKIRKL	DSGGFYITSR	TQFNSLQQLV	AYYSKHADGL		240
SRC_CHICK	178	ETTKGAYCLS	VSDFDNAKGL	NVKHYKIRKL	DSGGFYITSR	TQFSSLQQLV	AYYSKHADGL		237
		*****	*****	*****	*****	***.*****	*****		
SRC_HUMAN	241	CHRLTTVCPT	SKPQTQGLAK	DAWEIPRESL	RLEVKLQGC	FGEVWMGTWN	GTTRVAIKTL		300
SRC_CHICK	238	CHRLTNVCPT	SKPQTQGLAK	DAWEIPRESL	RLEVKLQGC	FGEVWMGTWN	GTTRVAIKTL		297
		*****.***	*****	*****	*****	*****	*****		
SRC_HUMAN	301	KPGTMSPEAF	LQEAQVMKKL	RHEKLVQLYA	VVSEEEPIYIV	TEYMSKGSLL	DFLKGETGKY		360
SRC_CHICK	298	KPGTMSPEAF	LQEAQVMKKL	RHEKLVQLYA	VVSEEEPIYIV	TEYMSKGSLL	DFLKGEEMKY		357
		*****	*****	*****	*****	*****	*****		
SRC_HUMAN	361	LRLPQLVDMA	AQIASGMAYV	ERMNYVHRDL	RAANILVGEN	LVCKVADFGL	ALLIEDNEYT		420
SRC_CHICK	358	LRLPQLVDMA	AQIASGMAYV	ERMNYVHRDL	RAANILVGEN	LVCKVADFGL	ALLIEDNEYT		417
		*****	*****	*****	*****	*****	*****		
SRC_HUMAN	421	ARQGAKFPIK	WTAEAAALYG	RFTIKSDVWS	FGILLTELTT	KGRVPYPGMV	NREVLQDQVER		480
SRC_CHICK	418	ARQGAKFPIK	WTAEAAALYG	RFTIKSDVWS	FGILLTELTT	KGRVPYPGMV	NREVLQDQVER		477
		*****	*****	*****	*****	*****	*****		
SRC_HUMAN	481	GYRMPCPPEC	PESLHDLMCQ	CWRKEPEERP	TFEYLQAFLE	DYFTSTEPQY	QPGENL		536
SRC_CHICK	478	GYRMPCPPEC	PESLHDLMCQ	CWRKDPEERP	TFEYLQAFLE	DYFTSTEPQY	QPGENL		533

809

***** ***** *****:***** ***** ***** *****

810

Appendix 2: Figures