

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

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

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

5 (Dated: July 14, 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

6 I. INTRODUCTION

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

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

29 packages are now able to exploit GPUs [2, 3], which provide
30 greatly improved simulation efficiency per unit cost relative
31 to CPUs, while distributed computing platforms such as
32 Folding@home [4], Copernicus [5, 6], and GPGPU [7], allow
33 scalability on an unprecedented level. In parallel, methods
34 for building human-understandable models of protein
35 dynamics from noisy simulation data, such as Markov state
36 modeling (MSM) approaches, are now reaching maturity [8–
37 10]. MSM methods in particular have the advantage of being
38 able to aggregate data from multiple independent MD
39 trajectories, facilitating parallelization of production simulations
40 and thus greatly alleviating overall computational cost. There
41 also exist a number of mature software packages for comparative
42 modeling of protein structures, in which a target protein
43 sequence is modeled using one or more structures as templates
44 [11, 12]. One such piece of software, MODELLER, has also
45 been used recently to study protein allostery by generating
46 and refining configurational models, sampled by interpolating
47 between two user-defined metastable structures [13].

48 However, it remains difficult for researchers to exploit the
49 full variety of available protein sequence and structural data
50 in simulation studies, largely due to limitations in software
51 architecture. For example, the set up of a biomolecular
52 simulation is typically performed manually, encompassing a
53 series of fairly standard (yet time-consuming) steps such as
54 the choice of protein sequence construct and starting structure(s),
55 addition of missing residues and atoms, solvation

* Corresponding author; john.chodera@choderalab.org

57 with explicit water and counterions (and potentially buffer
58 components and cosolvents), choice of simulation param-
59 eters (or parameterization schemes for components where
60 parameters do not yet exist), system relaxation with energy
61 minimization, and one or more short preparatory MD sim-
62 ulations to equilibrate the system and relax the simulation
63 cell. Due to the laborious and manual nature of this pro-
64 cess, simulation studies typically consider only one or a few
65 proteins and starting configurations. Worse still, studies (or
66 collections of studies) that *do* consider multiple proteins of-
67 ten suffer from the lack of consistent best practices in this
68 preparation process, making comparisons between related
69 proteins unnecessarily difficult.

70 The ability to fully exploit the large quantity of available
71 protein sequence and structural data in biomolecular sim-
72 ulation studies could open up many interesting avenues for
73 research, enabling the study of entire protein families or su-
74 perfamilies within a single organism or across multiple or-
75 ganisms. The similarity between members of a given pro-
76 tein family could be exploited to generate arrays of confor-
77 mational models, which could be used as starting configu-
78 rations to aid sampling in MD simulations. This approach
79 would be highly beneficial for many MD methods, such as
80 MSM construction, which require global coverage of the con-
81 formational landscape to realize their full potential, and
82 would also be particularly useful in cases where structural
83 data is present for only a subset of the members of a pro-
84 tein family. It would also aid in studying protein families
85 known to have multiple metastable conformations—such as
86 kinases—for which the combined body of structural data for
87 the family may cover a large range of these conformations,
88 while the available structures for any individual member
89 might encompass only one or two distinct conformations.

90 Here, we present the first steps toward bridging the
91 gap between biomolecular simulation software and *omics*-
92 scale sequence and structural data: a fully automated open
93 source framework for building simulation-ready protein
94 models in multiple conformational substates scalable from
95 single sequences to entire superfamilies. **Ensembler** pro-
96 vides functions for selecting target sequences and homolo-
97 gous template structures, and (by interfacing with a num-
98 ber of external packages) performs pairwise alignments,
99 comparative modeling of target-template pairs, and several
100 stages of model refinement. As an example application, we
101 have constructed models for the entire set of human tyro-
102 sine kinase (TK) catalytic domains, using all available struc-
103 tures of protein kinase domains (from any species) as tem-
104 plates. This results in a total of almost 400,000 models,
105 and we demonstrate that these provide wide-ranging cov-
106 erage of known functionally relevant conformations. By us-
107 ing these models as starting configurations for highly par-
108 allel MD simulations, we expect their structural diversity to
109 greatly aid in sampling of conformational space. We further
110 suggest that models with high target-template sequence
111 identity are the most likely to represent native metastable
112 states, while lower sequence identity models would aid
113 in sampling of more distant regions of accessible phase
114 space. It is also important to note that some models (es-

115 pecially low sequence identity models) may not represent
116 natively accessible conformations. However, MSM meth-
117 ods benefit from the ability to remove outlier MD trajec-
118 tories which start from non-natively accessible conforma-
119 tions, and which would thus be unconnected with the phase
120 space sampled in other trajectories. These methods essen-
121 tially identify the largest subset of Markov nodes which con-
122 stitute an ergodic network [14, 15].

123 We anticipate that **Ensembler** will prove to be useful in
124 a number of other ways. For example, the generated mod-
125 els could represent valuable data sets even without subse-
126 quent production simulation, allowing exploration of the
127 conformational diversity present within the available struc-
128 tural data for a given protein family. Furthermore, the au-
129 tomation of simulation set up provides an excellent oppor-
130 tunity to make concrete certain "best practices", such as the
131 choice of simulation parameters.

132 II. DESIGN AND IMPLEMENTATION

133 **Ensembler** is written in Python, and can be used via a
134 command-line tool (`ensembl`) or via a flexible Python
135 API to allow integration of its components into other
136 applications. All command-line and API information in
137 this article refers to the [version 1.0.2 release of Ensembler](#).
138 Up-to-date documentation can be found at [ensem-
139 bler.readthedocs.org](#).

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

144 A. Target selection and retrieval

145 The first stage entails the selection of a set of *target* pro-
146 tein sequences—the sequences for which the user is in-
147 terested in generating simulation-ready structural models.
148 This may be a single sequence—such as a full-length pro-
149 tein or a construct representing a single domain—or a col-
150 lection of sequences, such as a particular domain from an
151 entire family of proteins. The output of this stage is a FASTA-
152 formatted text file containing the desired target sequences
153 with corresponding arbitrary identifiers.

154 The `ensembl` command-line tool allows targets to
155 be selected from UniProt—a freely accessible resource for
156 protein sequence and functional data ([uniprot.org](#)) [16]—
157 via a UniProt search query. To retrieve target sequences
158 from UniProt, the subcommand `gather_targets` is used
159 with the `--query` flag followed by a UniProt query string
160 conforming to the same syntax as the search function
161 available on the UniProt website. For example, `--query`
162 `'mnemonic:SRC_HUMAN'` would select the full-length hu-
163 man Src sequence, while the query shown in Box 1 would
164 select all human tyrosine protein kinases which have been
165 reviewed by a human curator. In this way, the user may se-
166 lect a single protein, many proteins, or an entire superfam-

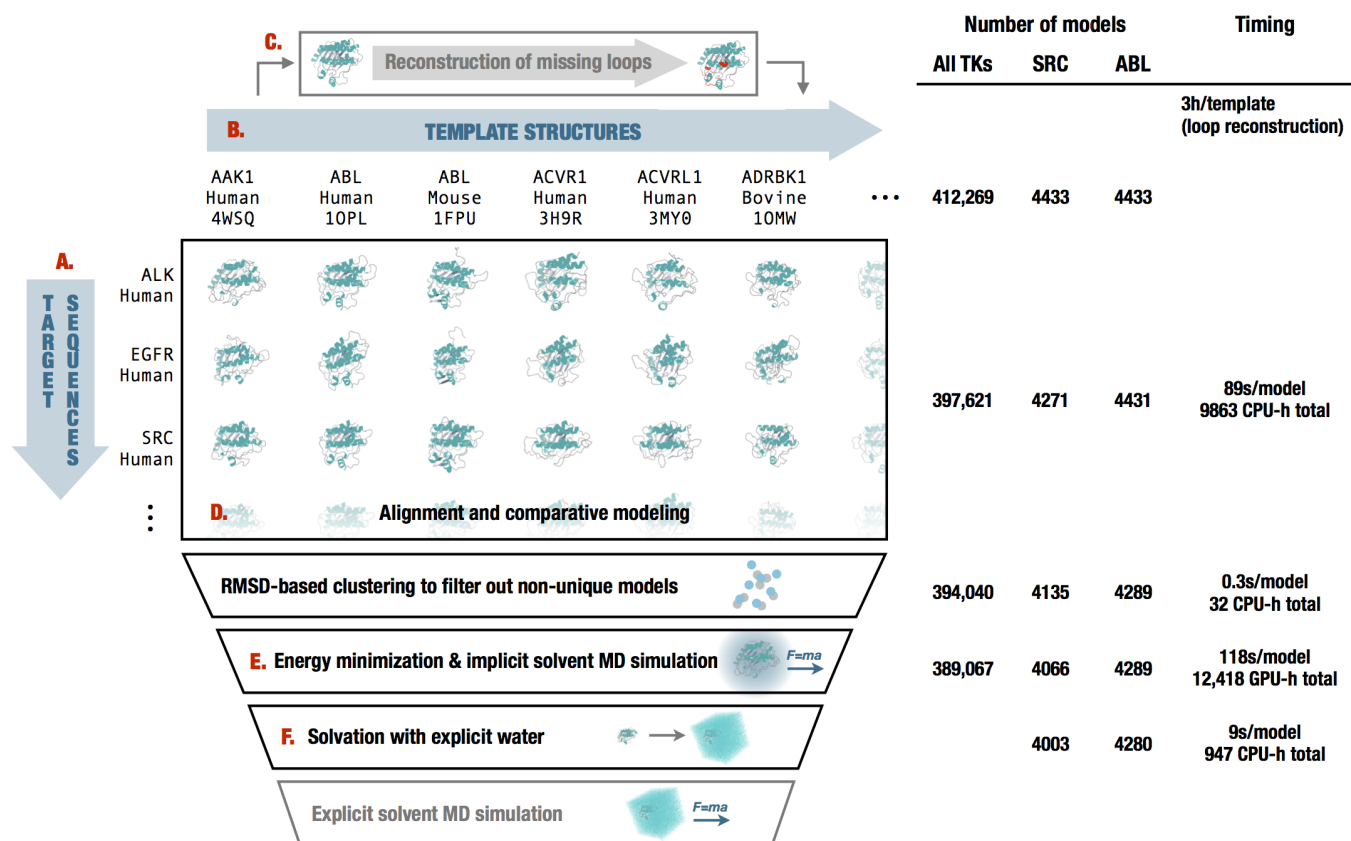


FIG. 1. Diagrammatic representation of the stages of the Ensembler pipeline and illustrative statistics for modeling all human tyrosine kinase catalytic domains. On the left, the various stages of the **Ensembler** pipeline are shown. The red labels indicate the corresponding text description provided 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.

167 ily from UniProt. The program outputs a FASTA file, setting
 168 the UniProt mnemonic (e.g. SRC_HUMAN) as the identifier for
 169 each target protein.

170 In many cases, it will be desirable to build models of an
 171 isolated protein domain, rather than the full-length pro-
 172 tein. The `gather_targets` subcommand allows protein
 173 domains to be selected from UniProt data by passing a regu-
 174 lar expression string to the `--uniprot_domain_regex` flag.
 175 For example, the above `--query` flag for selecting all hu-
 176 man protein kinases returns UniProt entries with domain
 177 annotations including "Protein kinase", "Protein kinase 1",
 178 "Protein kinase 2", "Protein kinase; truncated", "Protein ki-
 179 nase; inactive", "SH2", "SH3", etc. The regular expression
 180 shown in Box 1 selects only domains of the first three types.
 181 If the `--uniprot_domain_regex` flag is used, target identi-
 182 fiers are set with the form `[UniProt mnemonic]_D[domain`
 183 `index]`, where the latter part represents a 0-based index for
 184 the domain—necessary because a single target protein may
 185 contain multiple domains of interest (e.g. JAK1_HUMAN_D0,
 186 JAK1_HUMAN_D1).

187 Target sequences can also be defined manually (or from

188 another program) by providing a FASTA-formatted text file
 189 containing the desired target sequences with corresponding
 190 arbitrary identifiers.

191 B. Template selection and retrieval

192 **Ensembler** uses comparative modeling to build models,
 193 and as such requires a set of structures to be used as tem-
 194 plates. The second stage thus entails the selection of tem-
 195 plates and storage of associated sequences, structures, and
 196 identifiers. These templates can be specified manually, or
 197 using the `enssembler gather_templates` subcommand to
 198 automatically select templates based on a search of the
 199 Protein Data Bank (PDB) or UniProt. A recommended ap-
 200 proach is to select templates from UniProt which belong to
 201 the same protein family as the targets, guaranteeing some
 202 degree of homology between targets and templates.

203 The `enssembler gather_templates` subcommand pro-
 204 vides methods for selecting template structures from either
 205 UniProt or the PDB (<http://www.rcsb.org/pdb>), speci-

fied 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) [17], 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 [18, 19]. 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. 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 [20, 21] 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 [20, 21].

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 [22]—which uses a dynamic programming algorithm—with the PAM 250 scoring matrix of Gonnet *et al.* [23]. 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 [24], 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 `saAlign` routine [20, 21], PROMALS3D [25], and Espresso and 3DCoffee [26, 27], 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 ad-

ditional 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. 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 [14], which uses `mdtraj` [28] to calculate RMSD (for C_α atoms only) with a fast quaternion characteristic polynomial (QCP) [29–31] implementation. A minimum distance cutoff (which defaults to 0.6 \AA) is used to retain only a single model per cluster.

E. Refinement of models

A number of refinement methods have been developed to help guide comparative modeling techniques toward more "native-like" and physically consistent conformations [32, 33], of which MD simulations are an important example. While long-timescale unrestrained MD simulations (on the order of $100 \mu\text{s}$) have been found to be ineffective for recapitulating native-like conformations, possibly due to forcefield issues [34], even relatively short simulations can be useful for relaxing structural elements such as sidechain orientation [33].

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 [35], 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 [36], along with a modified generalized Born solvent model [37] 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 [38].

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 \AA), 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 [39] is used by default, but any of the other explicit water models available in OpenMM, such as TIP4P-Ew [40], 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 [36] and TIP3P water [39] 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 implemented in OpenMM, with a default pressure of 1 atm and

423 a period of 50 timesteps. The remaining simulation param-
424 eters have default values set to the same as for the implicit
425 solvent MD refinement.

426 Packaging

427 **Ensembler** provides a packaging module which
428 can be used to prepare models for other uses. The
429 `package_models` subcommand currently provides func-
430 tions (specified via the `--package_for` flag) for com-
431 pressing models in preparation for data transfer, or for
432 organizing them with the appropriate directory and file
433 structure for production simulation on the distributed
434 computing platform Folding@home [4]. The module could
435 easily be extended to add methods for preparing models
436 for other purposes. For example, production simulations
437 could alternatively be run using Copernicus [5, 6]—a frame-
438 work for performing parallel adaptive MD simulations—
439 or GPUGrid [7]—a distributing computing platform which
440 relies on computational power voluntarily donated by the
441 owners of nondedicated GPU-equipped computers.

442 Other features

443 *Tracking provenance information*

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

450 *Rapidly modeling a single template*

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

462 III. RESULTS

463 Modeling of all human tyrosine kinase catalytic domains

464 As a first application of **Ensembler**, we have built mod-
465 els for the human TK family. TKs (and protein kinases in

466 general) play important roles in many cellular processes and
467 are involved in a number of types of cancer [41]. For exam-
468 ple, a translocation between the TK Abl1 and the pseudok-
469 inase Bcr is closely associated with chronic myelogenous
470 leukemia [42], while mutations of Src are associated with
471 colon, breast, prostate, lung, and pancreatic cancers [43].
472 Protein kinase domains are thought to have multiple acces-
473 sible metastable conformation states, and much effort is di-
474 rected at developing kinase inhibitor drugs which bind to
475 and stabilize inactive conformations [44]. Kinases are thus
476 a particularly interesting subject for study with MSM meth-
477 ods [45], and this approach stands to benefit greatly from
478 the ability to exploit the full body of available genomic and
479 structural data within the kinase family, e.g. by generating
480 large numbers of starting configurations to be used in highly
481 parallel MD simulation.

482 We selected all human TK domains annotated in UniProt
483 as targets, and all available structures of protein kinase do-
484 mains (of any species) as templates, using the commands
485 shown in Box 1. This returned 93 target sequences and
486 4433 template structures, giving a total of 412,269 target-
487 template pairs. The templates were derived from 3028 indi-
488 vidual PDB entries and encompassed 23 different species,
489 with 3634 template structures from human kinase con-
490 structs.

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

494 Ensembler modeling statistics

495 Crystallographic structures of kinase catalytic domains
496 generally contain a significant number of missing residues
497 (median 11, mean 14, standard deviation 13, max 102) due to
498 the high mobility of several loops (Fig. 2, top), with a number
499 of these missing spans being significant in length (median 5,
500 mean 7, standard deviation 6, max 82; Fig. 2, bottom). To re-
501 duce the reliance on the MODELLER rapid model construc-
502 tion stage to reconstruct very long unresolved loops, un-
503 resolved template residues were first remodeled using the
504 `loopmodel` subcommand. Out of 3666 templates with one
505 or more missing residues, 3134 were successfully remod-
506 eled by the Rosetta loop modeling stage (with success de-
507 fined simply as program termination without error); most
508 remodeling failures were attributable to unsatisfiable spa-
509 tial constraints imposed by the original template structure.
510 There was some correlation between remodeling failures
511 and the number of missing residues (Fig. 2, top); templates
512 for which remodeling failed had a median of 20 missing
513 residues, compared to a median of 14 missing residues for
514 templates for which remodeling was successful.

515 Following loop remodeling, the **Ensembler** pipeline was
516 performed up to and including the implicit solvent MD re-
517 finement stage, which completed with 389,067 (94%) sur-
518 viving models across all TKs. To obtain statistics for the sol-
519 vation stage without generating a sizeable amount of coordi-
520 nate data (with solvated PDB coordinate files taking up

```

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)?'
    
```

Box 1. Ensembler command-line functions used to select targets and templates. The commands retrieve target and template data by querying UniProt. The query string provided to the `gather_targets` command selects all human tyrosine protein kinases which have been reviewed by a curator, while the query string provided to the `gather_templates` command selects all reviewed protein kinases of any species. The `--uniprot_domain_regex` flag is used to select a subset of the domains belonging to the returned UniProt protein entries, by matching the domain annotations against a given regular expression. In this example, domains of type "Protein kinase", "Protein kinase 1", and "Protein kinase 2" were selected, while excluding many other domain types such as "Protein kinase; truncated", "Protein kinase; inactive", "SH2", "SH3", etc. Target selection simply entails the selection of sequences corresponding to each matching UniProt domain. Template selection entails the selection of the sequences and structures of any PDB entries corresponding to the matching UniProt domains.

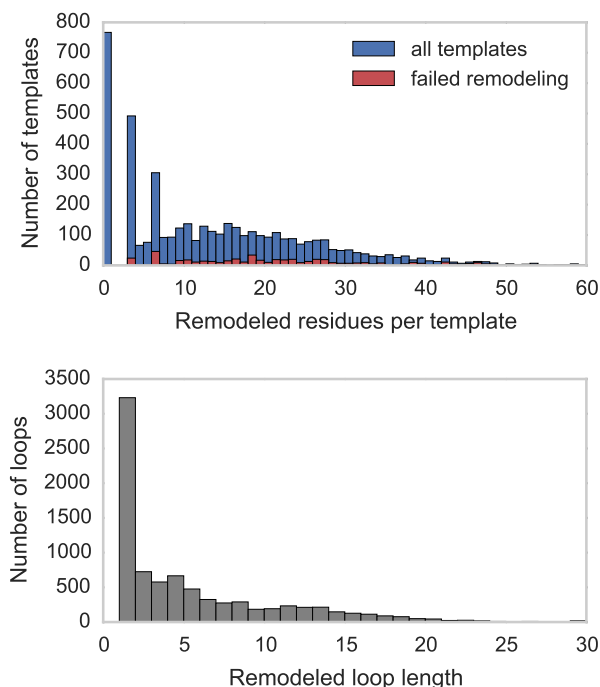


FIG. 2. Distributions for the number of missing residues in the TK templates. The upper histograms show the number of missing residues per template, for all templates (blue) and for only those templates for which template remodeling with the `loopmodel` subcommand failed (red). The lower histogram shows the number of residues in each missing loop, for all templates.

521 about 0.9 MB each), the `solvate` subcommand was per-
 522 formed for two representative individual kinases (*Src* and
 523 *Abl1*).

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

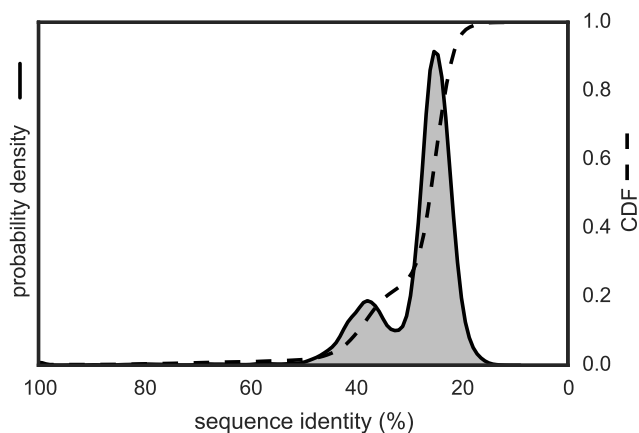


FIG. 3. Template-target sequence identity distribution for human tyrosine kinase catalytic domains. Sequence identities are calculated from all pairwise target-template alignments, where targets are human kinase catalytic domain sequences and templates are all kinase catalytic domains from any organism 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.

532 intensive.

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

Evaluation of model quality and utility

All tyrosine kinases

540
 541
 542 To evaluate the variety of template sequence similarities
 543 relative to each target sequence, we calculated sequence

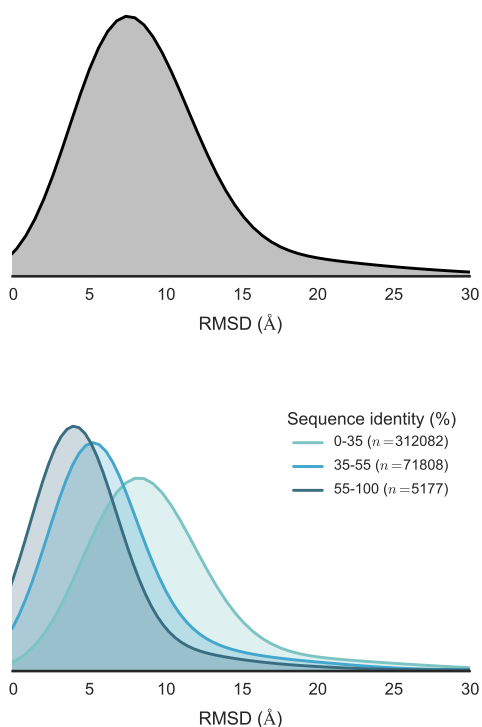


FIG. 4. Distribution of RMSDs to all TK catalytic domain models relative to the model derived from the highest sequence identity template. Distributions are built from data from all 93 TK domain targets. To better illustrate how conformational similarity depends on sequence identity, the lower plot illustrates 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.

544 identity distributions, as shown in Fig. 3. This suggests an
545 intuitive division into three categories, with 355,712 mod-
546 els in the 0–35% sequence identity range, 51,330 models in
547 the 35–55% range, and 5227 models in the 55–100% range.
548 We then computed the RMSD distributions for the models
549 created for each target (relative to the model derived from
550 the template with highest sequence identity) Fig. 4, to as-
551 sess the diversity of conformations captured by the mod-
552 eling pipeline. Furthermore, to understand the influence
553 of sequence identity on the conformational similarities of
554 the resulting models, the RMSD distributions were stratified
555 based on the three sequence identity categories de-
556 scribed above. This analysis indicates that higher sequence
557 identity templates result in models with lower RMSDs, while
558 templates with remote sequence identities result in larger
559 RMSDs on average.

560 We also analyzed the potential energies of the models
561 at the end of the implicit solvent MD refinement stage.
562 These ranged from -14180 kT to -3160 kT, with a median
563 of -9501 kT, mean of -9418 kT, and a standard deviation
564 of 1198 kT (with a simulation temperature of 300 K). The
565 distributions—stratified using the same sequence identity

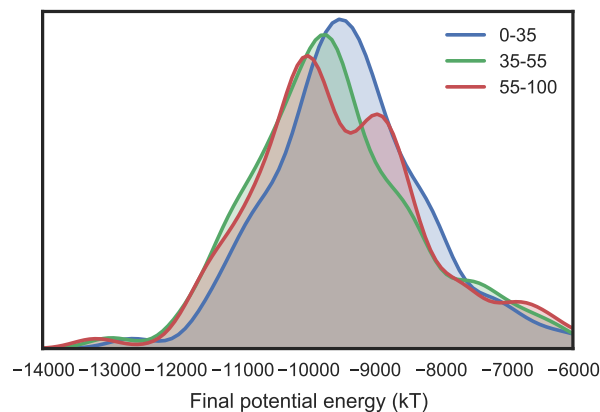


FIG. 5. Distribution of final energies from implicit solvent MD refinement of TK catalytic domain models. To illustrate how the energies are affected by sequence identity, the models are separated into three sequence identity classes: high identity (55–100%), moderate identity (35–55%), and remote identity (0–35%). The plotted distributions have been smoothed using kernel density estimation. Refinement simulations were carried out at the default temperature of 300 K.

566 ranges as above—are plotted in Fig. 5, indicating that higher
567 sequence identity templates tend to result in slightly lower
568 energy models. Of the 4973 models which failed to complete
569 the implicit refinement MD stage, all except 9 failed within
570 the first 1 ps of simulation.

571 *Src and Abl1*

572 To provide a more complete evaluation of the models
573 generated, we have analyzed two example TKs (*Src* and *Abl1*)
574 in detail. Due to their importance in cancer, these kinases
575 have been the subject of numerous studies, encompassing
576 many different methodologies. In terms of structural data,
577 a large number of crystal structures have been solved (with
578 or without ligands such as nucleotide substrate or inhibitor
579 drugs), showing the kinases in a number of different confor-
580 mations. These two kinases are thus also interesting targets
581 for MSM studies, with one recent study focusing on mod-
582 eling the states which constitute the activation pathway of
583 *Src* [45].

584 Fig. 6 shows a superposition of a set of representative
585 models of *Src* and *Abl1*. Models were first stratified into three
586 ranges, based on the structure of the sequence identity dis-
587 tribution (Fig. 3), then subjected to RMSD-based k -medoids
588 clustering (using the *msmbuilder* clustering package [14]) to
589 pick three representative models from each sequence iden-
590 tity range. Each model is colored and given a transparency
591 based on the sequence identity between the target and tem-
592 plate sequence. The figure gives an idea of the variance
593 present in the generated models. High sequence identity
594 models (in opaque blue) tend to be quite structurally sim-
595 ilar, with some variation in loops or changes in domain ori-

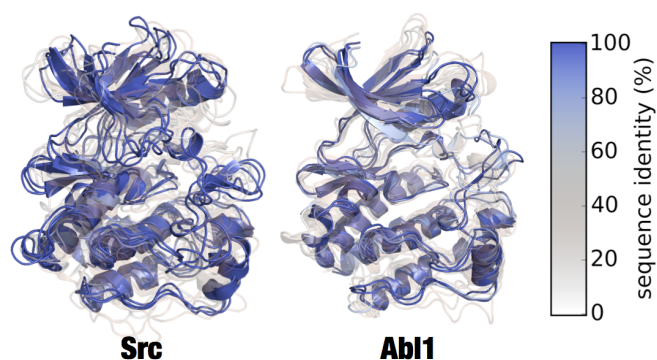


FIG. 6. Superposition of clustered models of Src and Abl1. Superposed renderings of nine models each for Src and Abl1, giving some indication the diversity of conformations generated by Ensembler. The models for each target were divided into three sequence identity ranges (as in Fig. 4), and RMSD-based k -medoids clustering was performed (using the msmbuilder clustering package [14]) 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.

596 entation.

597 The Abl1 renderings in Fig. 6 indicate one high sequence
598 identity model with a long unstructured region at one of
599 the termini, which was unresolved in the original template
600 structure. While such models are not necessarily incorrect
601 or undesirable, it is important to be aware of the effects they
602 may have on production simulations performed under periodic
603 boundary conditions, as long unstructured termini can
604 be prone to interact with a protein's periodic image. Lower
605 sequence identity models (in transparent white or red) indicate
606 much greater variation in all parts of the structure.
607 We believe the mix of high and low sequence identity models
608 to be particularly useful for methods such as MSM building,
609 which require thorough sampling of the conformational
610 landscape. The high sequence identity models could be
611 considered to be the most likely to accurately represent true
612 metastable states. Conversely, the lower sequence identity
613 models could be expected to help push a simulation into regions
614 of conformation space which might take intractably
615 long to reach if starting a single metastable conformation.

616 To evaluate the models of *Src* and *Abl1* in the context of the
617 published structural biology literature on functionally relevant
618 conformations, we have focused on two residue pair
619 distances thought to be important for the regulation of protein
620 kinase domain activity. We use the residue numbering
621 schemes for chicken *Src* (which is commonly used in the literature
622 even in reference to human *Src*) [46, 47] and human
623 *Abl1* isoform A [48–50] respectively; the exact numbering
624 schemes are provided in Appendix 1.

625 Fig. 7 shows two structures of *Src* believed to represent
626 inactive (PDB code: 2SRC) [46] and active (PDB code: 1Y57)
627 [47] states. One notable feature which distinguishes the two
628 structures is the transfer of an electrostatic interaction of E310
629 from R409 (in the inactive state) to K295 (in

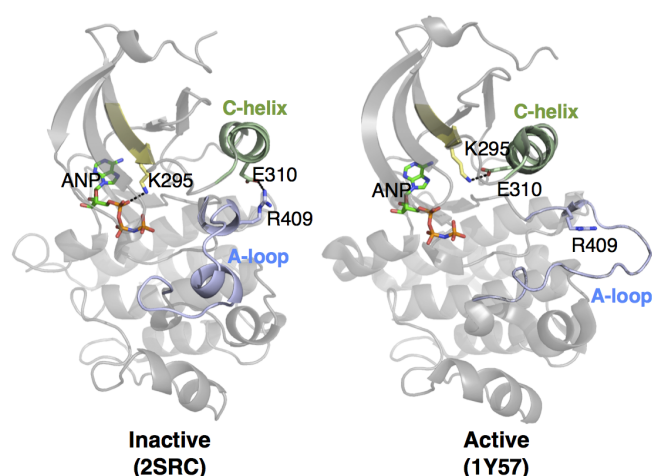


FIG. 7. Two structures of Src, indicating certain residues involved in activation. In the inactive state, E310 forms a salt bridge with R409. During activation, the α C-helix (green) moves and rotates, orienting E310 towards the ATP-binding site and allowing it to instead form a salt bridge with K295. 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.

630 the active state), brought about by a rotation of the α C-
631 helix. These three residues are also well conserved [51], and
632 a number of experimental and simulation studies have suggested
633 that this electrostatic switching process plays a role in a
634 regulatory mechanism shared across the protein kinase family
635 [45, 52, 53]. As such, we have projected the **Ensembler**
636 models for *Src* and *Abl1* onto a space consisting of the
637 distances between these two residue pairs (Fig. 8). The models
638 show strong coverage of regions in which either of the
639 electrostatic interactions is fully formed (for models across
640 all levels of target-template sequence identity), as well as a
641 wide range of regions in-between (mainly models with low
642 sequence identity). We thus expect that such a set of models,
643 if used as starting configurations for highly parallel MD
644 simulation, could greatly aid in sampling of functionally relevant
645 conformational states.

646 IV. AVAILABILITY AND FUTURE DIRECTIONS

647 Availability

648 The code for **Ensembler** is hosted on the collaborative
649 open source software development platform GitHub
650 (github.com/choderalab/enssembler). The latest release can
651 be installed via the conda package manager for Python
652 (conda.pydata.org), using the two commands shown in
653 Box 2. This will install all dependencies except for
654 MODELLER and Rosetta, which are not available through the
655 conda package manager, and thus must be installed separately
656 by the user. The latest source can be downloaded

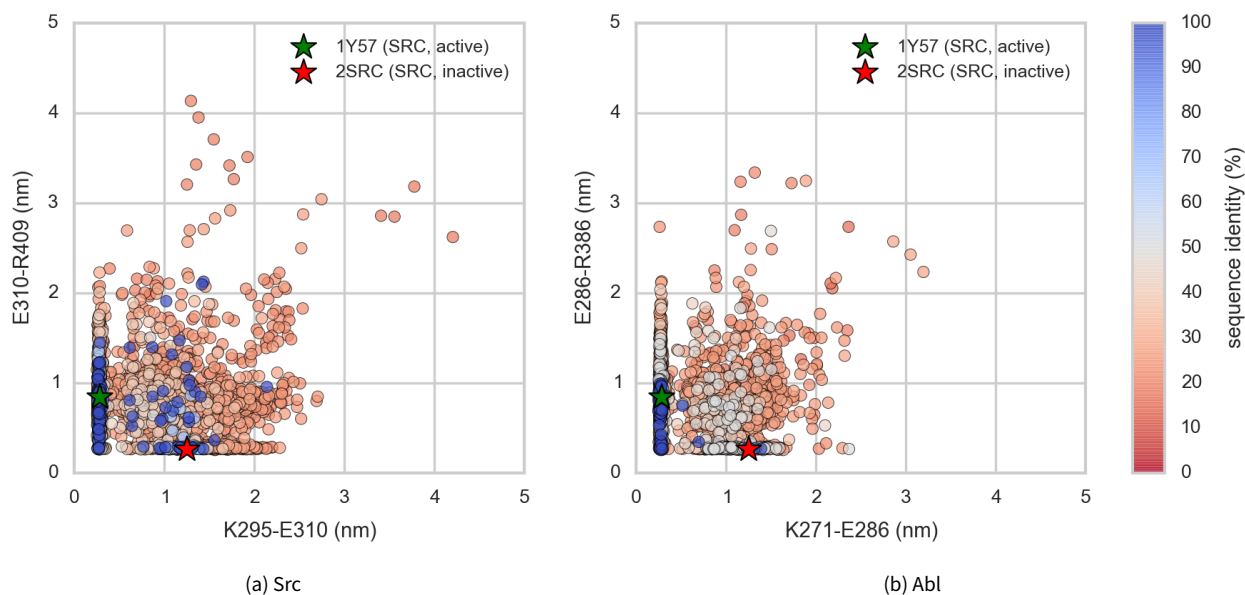


FIG. 8. Src and Abl1 models projected onto the distances between two conserved residue pairs, colored by sequence identity. Two Src structures (PDB entries 1Y57 [47] and 2SRC [46]) are projected onto the plots for reference, representing active and inactive states respectively. These structures and the residue pairs analyzed here are depicted in Fig. 7. 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).

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

Box 2. Ensembler installation using conda.

657 from the GitHub repository, which also contains up-to-date
658 instructions for building and installing the code. Documenta-
659 tion can be found at ensembler.readthedocs.org.

660 A supplementary dataset can also be downloaded from
661 the Dryad Digital Repository (DOI: [10.5061/dryad.7fg32](https://doi.org/10.5061/dryad.7fg32)).
662 This contains the TK models described in the III section, gen-
663 eral information on the targets and templates, plus a script
664 and instructions for regenerating the same dataset.

Future Directions

666 Comparative protein modeling and MD simulation set-up
667 can be approached in a number of different ways, with vary-
668 ing degrees of complexity, and there are a number of obvi-
669 ous additions and improvements which we plan to imple-
670 ment in future versions of **Ensembler**.

671 Some amino acids can exist in different protonation
672 states, depending on pH and on their local environment.
673 These protonation states can have important effects on bi-
674 ological processes. For example, long timescale MD simu-
675 lations have suggested that the conformation of the DFG mo-
676 tif of the TK Abl1—believed to be an important regulatory
677 mechanism [54]—is controlled by protonation of the aspar-

678 tate [55]. Currently, protonation states are assigned simply
679 based on pH (a user-controllable parameter). At neutral pH,
680 histidines have two protonation states which are approxi-
681 mately equally likely, and in this situation the selection is
682 therefore made based on which state results in a better hy-
683 drogen bond. It would be highly desirable to instead use a
684 method which assigns amino acid protonation states based
685 on a rigorous assessment of the local environment. We thus
686 plan to implement an interface and command-line function
687 for assigning protonation states with MCCE2 [56–58], which
688 uses electrostatics calculations combined with Monte Carlo
689 sampling of side chain conformers to calculate pKa values.

690 Many proteins require the presence of various types of
691 non-protein atoms and molecules for proper function, such
692 as metal ions (e.g. Mg^{+2}), cofactors (e.g. ATP) or post-
693 translational modifications (e.g. phosphorylation, methyl-
694 ation, glycosylation, etc.), and we thus plan for **Ensembler**
695 to eventually have the capability to include such entities
696 in the generated models. Binding sites for metal ions are
697 frequently found in proteins, often playing a role in cataly-
698 sis. For example, protein kinase domains contain two bind-
699 ing sites for divalent metal cations, and display significantly
700 increased activity in the presence of Mg^{2+} [59], the diva-
701 lent cation with highest concentration in mammalian cells.
702 Metal ions are often not resolved in experimental structures
703 of proteins, but by taking into account the full range of avail-
704 able structural data, it should be possible in many cases
705 to include metal ions based on the structures of homolo-
706 gous proteins. We are careful to point out, however, that
707 metal ion parameters in classical MD force fields have signif-

708 icant limitations, particularly in their interactions with pro-
709 teins [60]. Cofactors and post-translational modifications
710 are also often not fully resolved in experimental structures,
711 and endogenous cofactors are frequently substituted with
712 other molecules to facilitate experimental structural analy-
713 sis. Again, **Ensembler** could exploit structural data from a
714 set of homologous proteins to model in these molecules, al-
715 though there will likely be a number of challenges to over-
716 come in the design and implementation of such function-
717 ality.

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

731 Conclusion

732 We believe **Ensembler** to be an important first step to-
733 ward enabling computational modeling and simulation of

734 proteins on the scale of entire protein families, and suggest
735 that it could likely prove useful for tasks beyond its original
736 aim of providing diverse starting configurations for MD sim-
737 ulations. The code is open source and has been developed
738 with extensibility in mind, in order to facilitate its customiza-
739 tion for a wide range of potential uses by the wider scientific
740 community.

741 V. ACKNOWLEDGMENTS

742 The authors are grateful to Robert McGibbon (Stanford)
743 and Arien S. Rustenburg (MSKCC) for many excellent soft-
744 ware engineering suggestions. The authors thank Nicholas
745 M. Levinson (University of Minnesota), Markus A. Seeliger
746 (Stony Brook), Diwakar Shukla (Stanford), and Avner Sch-
747 lessinger (Mount Sinai) for helpful scientific feedback on
748 modeling kinases. The authors are grateful to Benjamin
749 Webb and Andrej Šali (UCSF) for help with the MODELLER
750 package, Peter Eastman and Vijay Pande (Stanford) for as-
751 sistance with OpenMM, and Marilyn Gunner (CCNY) for assis-
752 tance with MCCE2. All authors acknowledge support from
753 the Sloan Kettering Institute. JDC, KAB, and DLP acknowl-
754 edge partial support from NIH grant P30 CA008748. JDC
755 and DLP also acknowledge the generous support of a Louis
756 V. Gerstner Young Investigator Award. KAB was also sup-
757 ported in part by Starr Foundation grant I8-A8-058. PBG ac-
758 knowledges partial funding support from the Weill Cornell
759 Graduate School of Medical Sciences.

-
- 760 [1] G. M. Lee and C. S. Craik, *Science* **324**, 213 (2009).
761 [2] P. Eastman, M. S. Friedrichs, J. D. Chodera, R. J. Radmer, C. M.
762 Bruns, J. P. Ku, K. A. Beauchamp, T. J. Lane, L.-P. Wang, D.
763 Shukla, and V. S. Pande, *J. Chem. Theory Comput.* **9**, 461
764 (2012).
765 [3] R. Salomon-Ferrer, A. W. Götz, D. Poole, S. L. Grand, and R. C.
766 Walker, *J. Chem. Theor. Comput.* **9**, 3878 (2013).
767 [4] M. Shirts and V. S. Pande, *Science* **290**, 1903 (2000).
768 [5] S. Pronk, P. Larsson, I. Pouya, G. R. Bowman, I. S. Haque, K.
769 Beauchamp, B. Hess, V. S. Pande, P. M. Kasson, and E. Lindahl,
770 in *Proceedings of 2011 International Conference for High Per-*
771 *formance Computing, Networking, Storage and Analysis, SC '11*
772 *(ACM, New York, NY, USA, 2011)*, pp. 60:1–60:10.
773 [6] S. Pronk, I. Pouya, M. Lundborg, G. Rotskoff, B. Wesén, P. M.
774 Kasson, and E. Lindahl, *Journal of Chemical Theory and Com-*
775 *putation* (2015).
776 [7] I. Buch, M. J. Harvey, T. Giorgino, D. P. Anderson, and G.
777 De Fabritiis, *Journal of Chemical Information and Modeling*
778 **50**, 397 (2010).
779 [8] V. S. Pande, K. Beauchamp, and G. R. Bowman, *Methods* **52**,
780 99 (2010).
781 [9] J.-H. Prinz, H. Wu, M. Sarich, B. Keller, M. Fischbach, M. Held,
782 J. D. Chodera, C. Schütte, and F. Noé, *J. Chem. Phys.* **134**,
783 174105 (2011).
784 [10] J. D. Chodera and F. Noé, *Curr. Opin. Struct. Biol.* **25**, 135
785 (2014).
786 [11] J. Moulton, K. Fidelis, A. Kryshchuk, T. Schwede, and A. Tra-
787 montano, *Proteins: Structure, Function, and Bioinformatics*
788 **82**, 1 (2014).
789 [12] D. Baker and A. Šali, *Science* **294**, 93 (2001).
790 [13] P. Weinkam, J. Pons, and A. Šali, *Proceedings of the National*
791 *Academy of Sciences of the United States of America* **109**,
792 4875 (2012).
793 [14] K. A. Beauchamp, G. R. Bowman, T. J. Lane, L. Maibaum, I. S.
794 Haque, and V. S. Pande, *Journal of Chemical Theory and Com-*
795 *putation* **7**, 3412 (2011).
796 [15] R. Scalco and A. Caflisch, *The Journal of Physical Chemistry.*
797 *B* **115**, 6358 (2011).
798 [16] T. U. Consortium, *Nucleic Acids Research* **43**, D204 (2015).
799 [17] S. Velankar, J. M. Dana, J. Jacobsen, G. van Ginkel, P. J. Gane,
800 J. Luo, T. J. Oldfield, C. O'Donovan, M.-J. Martin, and G. J. Kley-
801 wegt, *Nucleic Acids Research* **41**, D483 (2013).
802 [18] B. Qian, S. Raman, R. Das, P. Bradley, A. J. McCoy, R. J. Read,
803 and D. Baker, *Nature* **450**, 259 (2007).
804 [19] C. Wang, P. Bradley, and D. Baker, *Journal of Molecular Biol-*
805 *ogy* **373**, 503 (2007).
806 [20] A. Fiser, R. K. G. Do, and A. Šali, *Protein Science* **9**, 1753 (2000).
807 [21] A. Šali and T. L. Blundell, *Journal of Molecular Biology* **234**,
808 779 (1993).
809 [22] P. J. A. Cock, T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, A.
810 Dalke, I. Friedberg, T. Hamelryck, F. Kauff, B. Wilczynski, and
811 M. J. L. de Hoon, *Bioinformatics (Oxford, England)* **25**, 1422
812 (2009).

- 813 [23] G. H. Gonnet, M. A. Cohen, and S. A. Brenner, *Science* **256**, 1443
814 (1992).
- 815 [24] J. D. Thompson, B. Linard, O. Lecompte, and O. Poch, *PLoS*
816 *ONE* **6**, e18093 (2011).
- 817 [25] J. Pei, B.-H. Kim, and N. V. Grishin, *Nucleic Acids Research* **36**,
818 2295 (2008).
- 819 [26] F. Armougom, S. Moretti, O. Poirot, S. Audic, P. Dumas, B.
820 Schaeli, V. Keduas, and C. Notredame, *Nucleic Acids Research*
821 **34**, W604 (2006).
- 822 [27] O. Poirot, K. Suhre, C. Abergel, E. O'Toole, and C. Notredame,
823 *Nucleic Acids Research* **32**, W37 (2004).
- 824 [28] R. T. McGibbon, K. A. Beauchamp, C. R. Schwantes, L.-P. Wang,
825 C. X. Hernández, M. P. Harrigan, T. J. Lane, J. M. Swails, and
826 V. S. Pande, *bioRxiv* (2014).
- 827 [29] D. L. Theobald, *Acta Cryst. A* **61**, 478 (2005).
- 828 [30] P. Liu, D. K. Agrafiotis, and D. L. Theobald, *J. Comput. Chem.*
829 **31**, 1561 (2010).
- 830 [31] P. Liu, D. K. Agrafiotis, and D. L. Theobald, *J. Comput. Chem.*
831 **32**, 185 (2011).
- 832 [32] J. L. MacCallum, A. Pérez, M. J. Schnieders, L. Hua, M. P. Jacob-
833 son, and K. A. Dill, *Proteins: Structure, Function, and Bioinforma-*
834 *tics* **79**, 74 (2011).
- 835 [33] Y. Zhang, *Current Opinion in Structural Biology* **19**, 145 (2009).
- 836 [34] A. Raval, S. Piana, M. P. Eastwood, R. O. Dror, and D. E. Shaw,
837 *Proteins: Structure, Function, and Bioinformatics* **80**, 2071
838 (2012).
- 839 [35] D. C. Liu and J. Nocedal, *Mathematical Programming* **45**, 503
840 (1989).
- 841 [36] K. Lindorff-Larsen, S. P. anad Kim Palmo, P. Maragakis, J. L.
842 Klepeis, R. O. Dror, and D. E. Shaw, *Proteins* **78**, 1950 (2010).
- 843 [37] A. Onufriev, D. Bashford, and D. A. Case, *Proteins* **55**, 383
844 (2004).
- 845 [38] J. E. Basconi and M. R. Shirts, *Journal of Chemical Theory and*
846 *Computation* **9**, 2887 (2013).
- 847 [39] W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey,
848 and M. L. Klein, *Journal of Chemical Physics* **79**, 926 (1983).
- 849 [40] H. W. Horn, W. C. Swope, J. W. Pitera, J. D. Madura, T. J.
850 Dick, G. L. Hura, and T. Head-Gordon, *The Journal of Chem-*
851 *ical Physics* **120**, 9665 (2004).
- 852 [41] D. S. Krause and R. A. Van Etten, *New England Journal of*
853 *Medicine* **353**, 172 (2005).
- 854 [42] E. K. Greuber, P. Smith-Pearson, J. Wang, and A. M. Pender-
855 gast, *Nature Reviews Cancer* **13**, 559 (2013).
- 856 [43] L. C. Kim, L. Song, and E. B. Haura, *Nature Reviews Clinical*
857 *Oncology* **6**, 587 (2009).
- 858 [44] Y. Liu and N. S. Gray, *Nature Chemical Biology* **2**, 358 (2006).
- 859 [45] D. Shukla, Y. Meng, B. Roux, and V. S. Pande, *Nature Commun.*
860 **5**, 3397 (2014).
- 861 [46] W. Xu, A. Doshi, M. Lei, M. J. Eck, and S. C. Harrison, *Molecular*
862 *Cell* **3**, 629 (1999).
- 863 [47] S. W. Cowan-Jacob, G. Fendrich, P. W. Manley, W. Jahnke, D.
864 Fabbro, J. Liebetanz, and T. Meyer, *Structure* **13**, 861 (2005).
- 865 [48] M. A. Young, N. P. Shah, L. H. Chao, M. Seeliger, Z. V. Milanov,
866 W. H. Biggs, D. K. Treiber, H. K. Patel, P. P. Zarrinkar, D. J. Lock-
867 hart, C. L. Sawyers, and J. Kuriyan, *Cancer Research* **66**, 1007
868 (2006).
- 869 [49] S. W. Cowan-Jacob, G. Fendrich, A. Floersheimer, P. Furet, J.
870 Liebetanz, G. Rummel, P. Rheinberger, M. Centeleghe, D. Fab-
871 bro, and P. W. Manley, *Acta Crystallographica Section D: Bio-*
872 *logical Crystallography* **63**, 80 (2006).
- 873 [50] N. M. Levinson, O. Kuchment, K. Shen, M. A. Young, M. Koldob-
874 skiy, M. Karplus, P. A. Cole, and J. Kuriyan, *PLoS Biol* **4**, e144
875 (2006).
- 876 [51] N. Kannan and A. F. Neuwald, *Journal of Molecular Biology*
877 **351**, 956 (2005).
- 878 [52] Z. H. Foda, Y. Shan, E. T. Kim, D. E. Shaw, and M. A. Seeliger,
879 *Nature Communications* **6**, 5939 (2015).
- 880 [53] E. Ozkirimli, S. S. Yadav, W. T. Miller, and C. B. Post, *Protein*
881 *Science : A Publication of the Protein Society* **17**, 1871 (2008).
- 882 [54] B. Nagar, O. Hantschel, M. A. Young, K. Scheffzek, D. Veach, W.
883 Bornmann, B. Clarkson, G. Superti-Furga, and J. Kuriyan, *Cell*
884 **112**, 859 (2003).
- 885 [55] Y. Shan, M. A. Seeliger, M. P. Eastwood, F. Frank, H. Xu, M. Å.
886 Jensen, R. O. Dror, J. Kuriyan, and D. E. Shaw, *Proceedings of*
887 *the National Academy of Sciences* **106**, 139 (2009).
- 888 [56] E. G. Alexov and M. R. Gunner, *Biophys. J.* **72**, 2075 (1997).
- 889 [57] R. E. Georgescu, E. G. Alexov, and M. R. Gunner, *Biophys. J.* **83**,
890 1731 (2002).
- 891 [58] Y. Song, J. Mao, and M. R. Gunner, *J. Comput. Chem.* **30**, 2231
892 (2009).
- 893 [59] J. A. Adams and S. S. Taylor, *Protein Science* **2**, 2177 (1993).
- 894 [60] S. F. Sousa, R. A. Fernandes, and M. J. Ramos, in *Kinetics*
895 *and Dynamics: From Nano- to Bio-Scale*, Vol. 12 of *Challenges*
896 *and Advances in Computational Chemistry and Physics*, edited
897 by P. a. D.-D. A. Paneth (Springer Science & Business Media,
898 Berlin, 2010), p. 530.

899

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

900 Kinase catalytic domains are highlighted in red, and the conserved residues analyzed in the main text (Figs. 7 and 8) are
 901 highlighted with yellow background.

902

Human Abl1 sequence

903 1 MLEICLKLVG CKSKKGLSS SSCYLEEALQ RVPASDFEPQ GLSEARWNS KENLLAGPSE 60
 904 61 NDPNLFVALY DFVASGDNTL SITKGEKLRV LGYNHNGEWC EAQTKNGQGW VPSNYITPVN 120
 905 121 SLEKHSWYHG PVSRNAEAYL LSSGINGSFL VRESESSPGQ RSISLRYEGR VYHYRINTAS 180
 906 181 DGKLYVSSSE RFNTLAELVH HHSTVADGLI TTLHYPAPKR NKPTVYGVSP NYDKWEMERT 240
 907 241 **DITMKHKLGG GQYGEVYEGV WKKYSLTVAV KTLKEDTMEV EEFLKEAAVM KEIKHPNLVQ** 300
 908 301 **LLGVCTREPP FYIITEFMTY GNLLDYLREC NRQEVNAVVL LYMATQISSA MEYLEKKNFI** 360
 909 361 **HRDLAARNCL VGENHLVKVA DFGLSRMLMTG DTYTAHAGAK FPIKWTAPES LAYNKFSIKS** 420
 910 421 **DVWAFGVLLW EIATYGMSPY PGIDLSQVYE LLEKDYRMER PEGCPEKVYE LMRACQWNP** 480
 911 481 **SDRPSFAEIH QAFETMFQES SISDEVEKEL GKQGVRGAVS TLLQAPELPT KTRTSRRAAE** 540
 912 541 HRDITDVPPEM PHSKGGQGESD PLDHEPAVSP LLPRKERGPP EGGLNEDERL LPKDKKTNLF 600
 913 601 SALIKKKKKT APTPPKRSSS FREMDGQPER RGAGEEGRD ISNGALAFPT LDTADPAKSP 660
 914 661 KPSNGAGVPN GALRESGGSG FRSPHLWKKS STLTSSRLAT GEEEGGSSS KRFLRSCSAS 720
 915 721 CVPHGAKDTE WRSVTLPRDL QSTGRQFDSS TFGGHKSEKP ALPRKRAGEN RSDQVTRGTV 780
 916 781 TPPPRLVKKK EEAADDEVFKD IMESSPGSSP PNLTPKPLRR QVTVAPASGL PHKEEAGKGS 840
 917 841 ALGTPAAAEPT VPTSKAGSG APGGTSKQPA EESRVRHKKH SSESPPGRDKG KLSRLKPAPP 900
 918 901 PPPAASAGKA GKPSQSPSQ EAAGEAVLGA KTKATSLVDA VNSDAAKPSQ PGEGLKQPVL 960
 919 961 PATPKPQSAK PSGTPIAPVP VPSTLPSASS ALAGDQPSST AFIPLISTRV SLRKRTRQPE 1020
 920 1021 RIASGAIKTKG VVLDSTEALC LAISRNSEQM ASHAVLEAG KNLYTFCVSY VDSIQQMRNK 1080
 921 1081 FAFREAINKL ENNRELQIC PATAGSGPAA TQDFSKLLSS VKEISDIVQR 1130

922

Sequences for human and chicken Src, aligned using Clustal Omega

923 SRC_HUMAN 1 MGSNKSHPKD ASQRRRSLEP AENVHGAGGG AFPASQTPSK PASADGHRGP SAAFAPAAAE 60
 924 SRC_CHICK 1 MGSSKSKPKD PSQRRRSLEP PDSTH--HG GFPASQTPNK TAAPDTHRTP SRSFGTVATE 57
 925 ***.***** ***** :.* * .*****.* *: * * * * * :.* .:*
 926 SRC_HUMAN 61 PKLFGGFNSS DTVTSPQRAG PLAGGVTTFFV ALYDYESRTE TDLSFKKGER LQIVNNTTEGD 120
 927 SRC_CHICK 58 PKLFGGFNTS DTVTSPQRAG ALAGGVTTFFV ALYDYESRTE TDLSFKKGER LQIVNNTTEGD 117
 928 *****.* ***** ***** ***** ***** ***** ***** *****
 929 SRC_HUMAN 121 WWLAHSLSTG QTGYIPSNYV APSDSIQAE WYFGKITRRE SERLLLNAEN PRGTFLVRES 180
 930 SRC_CHICK 118 WWLAHSLTTG QTGYIPSNYV APSDSIQAE WYFGKITRRE SERLLNPNEN PRGTFLVRES 177
 931 *****.* ***** ***** ***** ***** ***** ***** *****
 932 SRC_HUMAN 181 ETTKGAYCLS VSDFDNAKGL NVKHYKIRKL DSGGFYITSR TQFNSLQQLV AYYSKHADGL 240
 933 SRC_CHICK 178 ETTKGAYCLS VSDFDNAKGL NVKHYKIRKL DSGGFYITSR TQFSSLQQLV AYYSKHADGL 237
 934 ***** ***** ***** ***** ***** ***** ***** *****
 935 SRC_HUMAN 241 CHRLTTVCPT SKPQTQGLAK DAWAIPRESL RLEVKLGGQC FGEVWMTWN GTTRVAIKTL 300
 936 SRC_CHICK 238 CHRLTNVCPT SKPQTQGLAK DAWAIPRESL RLEVKLGGQC FGEVWMTWN GTTRVAIKTL 297
 937 *****.* ***** ***** ***** ***** ***** ***** *****
 938 SRC_HUMAN 301 KPGTMSPEAF LQEAQVMKKL RHEKLVQLYA VVSEEPYIV TEYMSKGSLL DFLKGETGKY 360
 939 SRC_CHICK 298 KPGTMSPEAF LQEAQVMKKL RHEKLVQLYA VVSEEPYIV TEYMSKGSLL DFLKGEMGKY 357
 940 ***** ***** ***** ***** ***** ***** ***** *****
 941 SRC_HUMAN 361 LRLPQLVDMA AQIASGMAYV ERMNYVHRDL RAANILVGEN LVCKVADFGL ARLIEDNEYT 420
 942 SRC_CHICK 358 LRLPQLVDMA AQIASGMAYV ERMNYVHRDL RAANILVGEN LVCKVADFGL ARLIEDNEYT 417
 943 ***** ***** ***** ***** ***** ***** ***** *****
 944 SRC_HUMAN 421 ARQGAKFPIK WTAPEAALYG RFTIKSDVWS FGILLTELTT KGRVPYPGMV NREVLQDQVER 480
 945 SRC_CHICK 418 ARQGAKFPIK WTAPEAALYG RFTIKSDVWS FGILLTELTT KGRVPYPGMV NREVLQDQVER 477
 946 ***** ***** ***** ***** ***** ***** ***** *****
 947 SRC_HUMAN 481 GYRMPCPPEC PESLHDLMCQ CWRKEPEERP TFEYLQAFLE DYFTSTEPQY QPGENL 536
 948 SRC_CHICK 478 GYRMPCPPEC PESLHDLMCQ CWRKDPEERP TFEYLQAFLE DYFTSTEPQY QPGENL 533
 949 ***** ***** *****.* ***** ***** ***** ***** *****