

Is structure based drug design ready for selectivity optimization?

Steven K. Albanese^{1,2,†}, John D. Chodera², Andrea Volkamer³, Simon Keng⁴, Robert Abel⁴, Lingle Wang^{4*}

¹Louis V. Gerstner, Jr. Graduate School of Biomedical Sciences, Memorial Sloan Kettering Cancer Center, New York, NY 10065; ²Computational and Systems Biology Program, Sloan Kettering Institute, Memorial Sloan Kettering Cancer Center, New York, NY 10065; ³Charité – Universitätsmedizin Berlin, Charitéplatz 1, 10117 Berlin; ⁴Schrödinger, New York, NY 10036

*For correspondence: lingle.wang@schrodinger.com (LW)

Present address: [†]Schrödinger, New York, NY 10036

Abstract

Alchemical free energy calculations are now widely used to drive or maintain potency in small molecule lead optimization with a roughly 1 kcal/mol accuracy. Despite this, the potential to use free energy calculations to drive optimization of compound *selectivity* among two similar targets has been relatively unexplored in published studies. In the most optimistic scenario, the similarity of binding sites might lead to a fortuitous cancellation of errors and allow selectivity to be predicted more accurately than affinity. Here, we assess the accuracy with which selectivity can be predicted in the context of small molecule kinase inhibitors, considering the very similar binding sites of human kinases CDK2 and CDK9, as well as another series of ligands attempting to achieve selectivity between the more distantly related kinases CDK2 and ERK2. Using a Bayesian analysis approach, we separate systematic from statistical error and quantify the correlation in systematic errors between selectivity targets. We find that, in the CDK2/CDK9 case, a high correlation in systematic errors suggests free energy calculations can have significant impact in aiding chemists in achieving selectivity, while in more distantly related kinases (CDK2/ERK2), the correlation in systematic error suggests fortuitous cancellation may even occur between systems that are not as closely related. In both cases, the correlation in systematic error suggests that longer simulations are beneficial to properly balance statistical error with systematic error to take full advantage of the increase in apparent free energy calculation accuracy in selectivity prediction.

Free energy methods have proven useful in aiding structure-based drug design by driving the optimization or maintenance of potency in lead optimization. Alchemical free energy calculations allow for the prediction of ligand binding free energies, including all enthalpic and entropic contributions [1]. Advances in atomistic molecular mechanics simulations and free energy methodologies [2–5] have allowed free energy methods to reach a level of accuracy sufficient for predicting ligand potencies [6]. These methods have been applied prospectively to develop inhibitors for Tyk2 [7], Syk [8], BACE1 [9], GPCRs [10], and HIV protease [11]. A recent large-scale review of the use of FEP+ [12] to predict potency for 92 different projects and 3 021 compounds determined that predicted binding free energies had a median root mean squared error (RMSE) of 1.0 kcal/mol [13].

Selectivity is an important consideration in drug design

In addition to potency, selectivity is an important property to consider in drug development, either in the pursuit of an inhibitor that is maximally selective [14, 15] or possesses a desired polypharmacology [16–

20]. Controlling selectivity can be useful not only in avoiding off-target toxicity (arising from inhibition of unintended targets) [21, 22], but also in avoiding on-target toxicity (arising from inhibition of the intended target) by selectively targeting disease mutations [23]. In either paradigm, considering the selectivity of a compound is complicated by the biology of the target. For example, kinases exist as nodes in complex signaling networks [24, 25] with feedback inhibition and cross-talk between pathways. Careful consideration of which off-targets are being inhibited can avoid off-target toxicity due to alleviating feedback inhibition and inadvertently reactivating the targeted pathway [24, 25] or the upregulation of a secondary pathway by alleviation of cross-talk inhibition [26, 27]. Off-target toxicity can also be caused by inhibiting unrelated targets, such as gefitinib, an EGFR inhibitor, inhibiting CYP2D6 [21] and causing hepatotoxicity in lung cancer patients. In a cancer setting, on-target toxicity can be avoided by considering the selectivity for the oncogenic mutant form of the kinase over the wild type form of the kinase [28–30], exemplified by a number of first generation EGFR inhibitors. Selective binding to multiple kinases can also lead to beneficial effects: Imatinib, initially developed to target BCR-Abl fusion proteins, is also approved for treating gastrointestinal stromal tumors (GIST) [31] due to its activity against receptor tyrosine kinase KIT.

The use of physical modeling to predict selectivity is relatively unexplored. While engineering compound selectivity is important for drug discovery, the utility of free energy methods for predicting this selectivity with the aim of reducing the number of compounds that must be synthesized to achieve a desired selectivity profile has been relatively unexplored in published studies. If there is fortuitous cancellation of systematic errors for closely related systems, free energy methods may be much more accurate than expected given the errors made in predicting the potency for each individual target. Such systematic errors might arise from force field parameters uncertainty, force field parameters assignment, protein or ligand protonation state assignment, or even from systematic errors arising in the target experimental data, where for example poor solubility of a particular compound might lead to a spuriously poor reported binding affinity for that compound, among other effects.

Molecular dynamics and free energy calculations have been used extensively to investigate the biophysical origins of the selectivity of imatinib for Abl kinase over Src [32, 33] and within a family of non-receptor tyrosine kinases [34]. This work focused on understanding the role reorganization energy plays in the exquisite selectivity of imatinib for Abl over the highly related Src despite high similarity between the cocrystallized binding mode and kinase conformations, and touches on neither the evaluation of the accuracy of these methods nor their application to drug discovery on congeneric series of ligands. Previous work predicting the selectivity of three bromodomain inhibitors across the bromodomain family achieved promising accuracy for single target potency of roughly 1 kcal/mol, but does not explicitly evaluate any selectivity metrics [35] or quantify the correlation in the errors made in predicting affinities for each bromodomain. Previous work using FEP+ to predict selectivity between pairs of phosphodiesterases (PDEs) showed promising performance but did not evaluate correlation in the error made in predicting affinities for each PDE [36]

Kinases are an important and particularly challenging model system for selectivity predictions. Kinases are a useful model system to work with for assessing the utility of free energy calculations to predict inhibitor selectivity in a drug discovery context. With the approval of imatinib for the treatment of chronic myelogenous leukemia in 2001, targeted small molecule kinase inhibitors (SMKIs) have become a major class of therapeutics in treating cancer and other diseases. Currently, there are 52 FDA-approved SMKIs [37], and it is estimated that kinase targeted therapies account for as much as 50% of current drug development [38], with many more compounds currently in clinical trials. While there have been a number of successful drug approvals, the current stable of FDA-approved kinase inhibitors targets only a small fraction of kinases implicated in disease, and the design of new selective kinase inhibitors for novel targets remains a significant challenge.

Achieving selective inhibition of kinases is quite challenging, as there are more than 518 protein kinases [39, 40] sharing a highly conserved ATP binding site that is targeted by the majority of SMKIs [41]. While kinase inhibitors have been designed to target kinase-specific sub-pockets and binding modes to achieve selectivity [42–47], previous work has shown that both Type I (binding to the active, DFG-in confor-

91 mation) and Type II (binding to the inactive, DFG-out conformation) inhibitors are capable of achieving a
92 range of selectivities [48, 49], often exhibiting significant binding to a number of other targets in addition
93 to their primary target. Even FDA-approved inhibitors—often the result of extensive drug development
94 programs—bind to a large number of off-target kinases [50]. Kinases are also targets of interest for devel-
95 oping polypharmacological compounds, or inhibitors that are specifically designed to inhibit multiple kinase
96 targets. Resistance to MEK inhibitors in KRAS-mutant lung and colon cancer has been shown to be driven
97 by ErbB3 upregulation [51], providing a rationale for dual MEK/ERBB family inhibitors. Similarly, combined
98 MEK and VEGFR1 inhibition has been proposed as a combinatorial approach to treat KRAS-mutant lung
99 cancer [52]. Developing inhibitors with a desired polypharmacology means navigating more complex se-
100 lectivity profiles, presenting a problem where physical modeling has the potential to dramatically speedup
101 drug discovery.

102 The correlation coefficient measures how useful predictions are in achieving selectivity
103 Since the prediction of selectivity depends on predicting the change of affinities to two or more targets
104 (or the change of affinities between pairs of related molecules for multiple targets), a spectrum of possi-
105 bilities exists for how accurately selectivity can be predicted even when the error in predicting individual
106 target affinities is fixed. In well-behaved kinase systems, for example, free energy calculation potency pre-
107 dictions have achieved root-mean-square of less than 1.0 kcal/mol [7, 12]. This residual error likely arises
108 from a variety of contributions. Systematic contributions to the residual error may include forcefield pa-
109 rameterization deficiencies, protein and ligand protonation assignment errors, and discrepancies between
110 the crystallographic construct protein and the assay construct protein. Likewise, unbiased contributions
111 to the observed residual error likely arises from incompletely converged sampling. Lastly, it should not be
112 forgotten that the target experimental value will have its own systematic and random errors.

113 In the best-case scenario, correlation in the systematic errors for predicting the interactions of a given
114 ligand with two related protein targets might exactly cancel out, allowing selectivity to be predicted much
115 more accurately than potency. On the other hand, if the uncorrelated random error dominates the residual
116 error between two protein targets, predictions of selectivity will be *less accurate* than potency predictions.
117 Real-world systems are likely to fall somewhere between these two extremes, and quantifying the *degree* to
118 which error in multiple protein targets is correlated, its implications for the use of free energy calculations
119 for prioritizing synthesis in the pursuit of selectivity, the ramifications for optimal calculation protocols, and
120 rough guidelines governing which systems we might expect good selectivity prediction is the primary focus
121 of this work.

122 In particular, in this work, we investigate the magnitude of the correlation (ρ) in error for the predicted
123 binding free energy differences between related compounds ($\Delta\Delta G_{ij}$) for two different targets, assessing
124 the utility of alchemical free energy calculations for the prediction of selectivity. We employ state of the
125 art relative free energy calculations [12, 13] to predict the selectivities of two different congeneric ligand
126 series [53, 54], and construct simple numerical models that allow us to quantify the potential utility in se-
127 lectivity optimization expected for different combinations of per target systematic errors and correlation
128 coefficients. To make a realistic assessment of our confidence in this correlation coefficient derived from
129 a limited number of experimental measurements, we develop a new Bayesian approach to quantify the
130 uncertainty in the correlation coefficient in the predicted change in selectivity on ligand modification, incor-
131 porating all sources of uncertainty and correlation in the computation to separate statistical from systematic
132 error. We find that in the closely related systems of CDK2 and CDK9, a high correlation of systematic errors
133 suggests that free energy methods can have a significant impact on speeding up selectivity optimization.
134 Even in the more distantly related case (CDK2/ERK2), correlation in the systematic errors allows free energy
135 calculations to speedup selectivity optimization, suggesting that these methodologies can impact drug dis-
136 covery even when comparing systems that are less closely related. We also present a model of the impact of
137 per target statistical error at different levels of systematic error correlation, suggesting that it is worthwhile
138 to expend more effort sampling in systems with high correlation.

139 Results

140 Alchemical free energy methods can be used to predict compound selectivity

141 While the potency of a ligand i for a single target is often quantified as a free energy of binding ($\Delta G_{i,\text{target}}$),
 142 there are a number of different metrics for quantifying compound selectivity [55, 56]. Here, we consider
 143 the selectivity S_i between one target and another (an *antitarget*) as the difference in free energy of binding
 144 for a given ligand i between the two,

$$S_i \equiv \Delta G_{i,\text{target 2}} - \Delta G_{i,\text{target 1}} \quad (1)$$

145 While in the optimization of potency we are concerned with $\Delta \Delta G_{ij} \equiv \Delta G_j - \Delta G_i$, the relative free energy
 146 of binding of ligands i and j to a single target, in the optimization of selectivity, we are concerned with
 147 $\Delta S_{ij} \equiv S_j - S_i$, which reflects the change in selectivity between ligand i and a related ligand j ,

$$\begin{aligned} \Delta S_{ij} &\equiv S_j - S_i & (2) \\ &= (\Delta G_{j,\text{target 2}} - \Delta G_{j,\text{target 1}}) - (\Delta G_{i,\text{target 2}} - \Delta G_{i,\text{target 1}}) \\ &= \Delta \Delta G_{ij,\text{target 2}} - \Delta \Delta G_{ij,\text{target 1}} \end{aligned}$$

148 To predict the change in selectivity, ΔS_{ij} , between two related compounds, we developed a protocol that
 149 uses a relative free energy calculation (FEP+) [12] to construct a map of alchemical perturbations between
 150 ligands in a congeneric series, as described in detail in the **Methods**. The calculation is repeated for each
 151 target of interest, with identical perturbations (edges) between each ligand (nodes). Each edge represents a
 152 relative alchemical free energy calculation that quantifies the $\Delta \Delta G$ between the ligands (nodes) for a single
 153 target. From these calculations, we can then compute the change in selectivity between the two targets of
 154 interest, ΔS_{ij} , achieved by transforming ligand i into ligand j .

155 Previous work has demonstrated that FEP+ can achieve an accuracy (σ_{target}) of roughly 1 kcal/mol in
 156 single-target potency prediction, which reflects a combination of systematic error and random statistical
 157 error [12]. However, it is possible that the systematic error for a given perturbation between ligands i
 158 and j ($\sigma_{\text{sys},ij,\text{target}}$) in two different systems may fortuitously cancel when computing ΔS_{ij} , resulting in the
 159 systematic contribution to the selectivity error ($\sigma_{\text{selectivity}}$) being significantly lower than its contribution to
 160 single-target potency error (σ_{target}). This systematic error may cancel between the two systems for a variety
 161 of reasons. For example, a ligand force field parameter assignment error might lead to an spuriously large
 162 solvation free energy for a particular compound, which will cancel in the selectivity analysis. Likewise, a spar-
 163 ingly soluble compound might have a similar experimental measurement error for the on-target protein as
 164 the off-target protein. Similar cancellation of systematic errors might be observed in ligand and/or protein
 165 protonation state assignment error, or systematic differences existing between the protein constructs used
 166 for crystallographic studies and biochemical or biophysical assays.

167 If we presume that the systematic errors for both targets are distributed according to a bivariate nor-
 168 mal distribution with correlation coefficient ρ quantifying the *degree* of correlation (with $\rho = 0$ denoting no
 169 correlation, $\rho = 1$ denoting perfect correlation, and $\rho = -1$ denoting perfect anti-correlation), and that the
 170 statistical errors for both targets ($\sigma_{\text{stat},ij,\text{target}}$) are completely independent because the simulations for each
 171 target are separate, we can model the error in predicting the ΔS_{ij} as $\sigma_{\text{selectivity}}$,

$$\sigma_{\text{selectivity}} \equiv \sqrt{\sigma_{\text{sys},ij,1}^2 + \sigma_{\text{sys},ij,2}^2 - 2\rho \sigma_{\text{sys},ij,1} \sigma_{\text{sys},ij,2} + \sigma_{\text{stat},ij,1}^2 + \sigma_{\text{stat},ij,2}^2} \quad (3)$$

172 $\sigma_{\text{selectivity}}$ can be split into two components: systematic error and statistical error. As more effort is spent on
 173 sampling, the per-target statistical error for a given transformation from ligand i to ligand j ($\sigma_{\text{stat},ij,\text{target}}$) will
 174 decrease, eventually becoming zero in the regime of infinite sampling. The correlation coefficient ρ can be
 175 both negative and positive. When the correlation coefficient ρ is positive, the systematic error ($\sigma_{\text{sys},ij,\text{target}}$)
 176 should cancel out, making $\sigma_{\text{selectivity}}$ smaller than expected. When the correlation coefficient ρ is negative,
 177 the systematic error ($\sigma_{\text{sys},ij,\text{target}}$) will be anti-correlated, making the $\sigma_{\text{selectivity}}$ larger than expected. As we
 178 shall see below, the quantitative value of the correlation coefficient ρ for the systematic error component
 179 has important ramifications for the accuracy with which selectivity can be predicted.

180 Correlation in systematic errors can significantly enhance accuracy of selectivity predictions

181 To demonstrate the potential impact the correlation coefficient ρ has on predicting selectivity using alchemi-
182 cal free energy techniques, we created a simple numerical model following Equation 3 which takes into
183 account each of the per-target systematic errors ($\sigma_{\text{sys},ij,1}$, $\sigma_{\text{sys},ij,2}$) expected from the methodology as well as
184 the correlation in those errors, while assuming infinite effort is spent on sampling to reduce the statistical
185 error component (σ_{stat}) to zero. As seen in Figure 1A, if the per target systematic errors are the same magni-
186 tude ($\sigma_{\text{sys},ij,1} = \sigma_{\text{sys},ij,2}$), $\sigma_{\text{selectivity}}$ approaches 0 as the correlation coefficient ρ approaches 1, even though the
187 single-target potency systematic error is nonzero. If the error for the free energy method is not the same
188 magnitude ($\sigma_{\text{sys},ij,1} \neq \sigma_{\text{sys},ij,2}$), $\sigma_{\text{selectivity}}$ gets smaller but approaches a non-zero value as ρ approaches 1.

189 To quantify the expected reduction in number of compounds that must be synthesized to achieve a de-
190 sired selectivity threshold (hereafter referred to as the *speedup* in selectivity optimization), we modeled the
191 change in selectivity with respect to a reference compound for a number of compounds a medicinal chemist
192 might suggest as a normal distribution centered around 0 with a standard deviation of 1 kcal/mol (Figure 1B,
193 black curve), reflecting the notion that most proposed modifications would not drive large changes in se-
194 lectivity. This assumption—that a synthetic chemist's proposal distribution can be modeled as a normal
195 distribution—is based on data-driven estimates from an Abbott Laboratories analysis of potency changes [57]

196 Further suppose that each compound is evaluated computationally with a free energy methodology that
197 has a per-target systematic error ($\sigma_{\text{sys},ij,\text{target}}$) of 1 kcal/mol, where we presume sufficient computational ef-
198 fort has been expended to make statistical error negligible. All compounds predicted to have a 1.4 kcal/mol
199 or greater improvement in selectivity (10x in ratio of affinities, or 1 \log_{10} unit) are synthesized and exper-
200 imentally tested (Figure 1B, colored curves), using an experimental technique with perfect measurement
201 accuracy. The fold-change in the proportion of compounds that are made that have a true 1.4 kcal/mol
202 improvement in selectivity compared to the original distribution can be calculated as a surrogate for the
203 expected speedup. For this 1.4 kcal/mol selectivity improvement threshold, a correlation coefficient $\rho = 0.5$
204 gives an expected speedup of 4.1x, which can be interpreted as needing to make 4.1x fewer compounds
205 to achieve a tenfold improvement in selectivity. This process can be extended for the even more difficult
206 proposition of achieving a hundredfold improvement in selectivity (Figure 1C), where 200–300x speedups
207 can be expected, depending on the single-target systematic error ($\sigma_{\text{sys},ij,\text{target}}$) for the free energy methodol-
208 ogy.

209 These estimates represent an ideal scenario, where the number of compounds scored and synthesized
210 is unlimited. In a more realistic discovery project, the number of compounds scored is limited by compu-
211 tational resources, and the number of compounds synthesized is limited by chemistry resources. In this
212 case, the observed speedup will depend not only on the correlation coefficient ρ and per-target systematic
213 error ($\sigma_{\text{sys},ij,\text{target}}$), but also the number of compounds scored and the synthesis rule, defined as the selectiv-
214 ity threshold a compound must be predicted to reach before being selected for synthesis. To model this
215 process, suppose a given number of compounds (Figure 1D, x-axis of each panel) are profiled with a free
216 energy method with a per-target systematic error ($\sigma_{\text{sys},ij,\text{target}}$) of 1 kcal/mol and some correlation coefficient
217 (ρ). The top compounds that are predicted to have an improvement in selectivity greater than a set "syn-
218 thesis rule" threshold (100x, 500x, or 1000x, Figure 1D, each curve) are synthesized, up to a maximum of
219 10 compounds. The expected speedup can then be calculated as the ratio of the number of synthesized
220 compounds that have a true selectivity improvement of 2.8 kcal/mol (100x or 2 log units) to the number
221 of compounds expected to have a true selectivity improvement of 2.8 kcal/mol had the same number of
222 compounds as were synthesized been drawn randomly from the underlying unit normal distribution.

223 As shown in Figure 1D, the more stringent synthesis rules combined with high correlation coefficients (ρ)
224 allow free energy calculations to have the highest impact in designing selectivity inhibitors, provided enough
225 compounds have been scored. Interestingly, at correlation coefficient $\rho=0.75$ and low numbers of scored
226 compounds, the 500x synthesis provides a greater speedup than 1000x synthesis rule. This is because
227 there is high probability no compounds meet the more 1000x stringent synthesis rule until many more
228 compounds are scored. This has implications for drug discovery efforts, where time and computational
229 effort may limit the number of compounds able to be profiled with free energy methods.

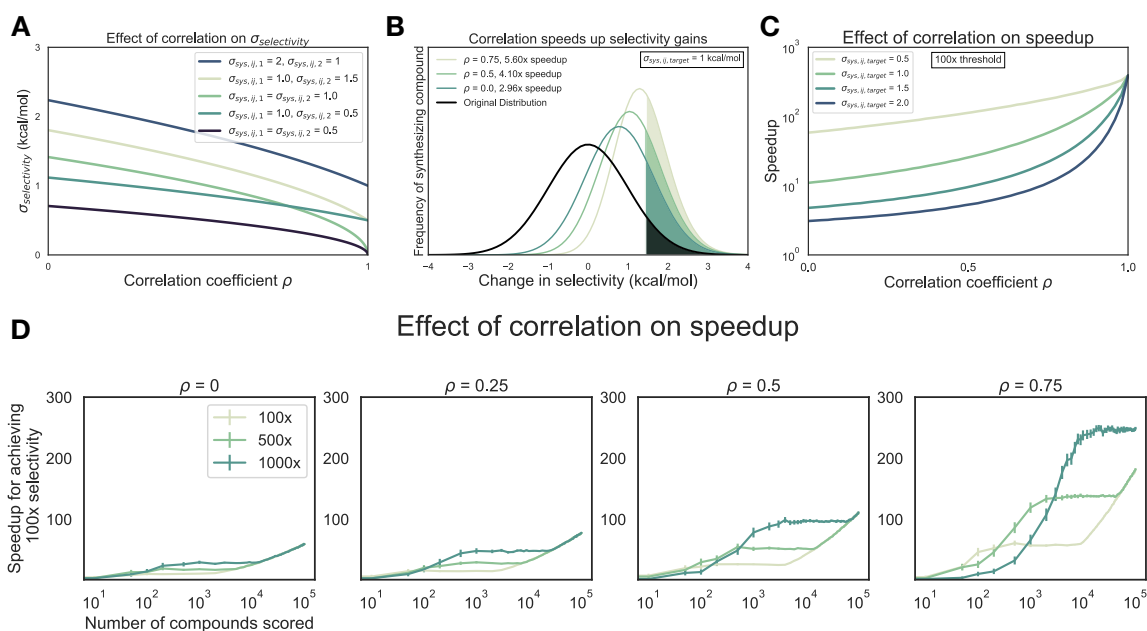


Figure 1. Free energy calculations can accelerate selectivity optimization. (A) The effect of correlation on expected errors for predicting selectivity ($\sigma_{\text{selectivity}}$) in kcal/mol when statistical error is negligible due to infinite sampling. Each curve represents a different combination of per target systematic errors ($\sigma_{\text{sys},ij,1}$ and $\sigma_{\text{sys},ij,2}$). (B) The change in selectivity for molecules proposed by medicinal chemists optimizing a lead candidate can be modeled by a normal distribution centered on zero with a standard deviation of 1 kcal/mol (black curve). Each green curve corresponds to the distribution of compounds made after screening for a 1 \log_{10} unit (1.4 kcal/mol) improvement in selectivity with a free energy methodology with a 1 kcal/mol per target systematic error and a particular correlation, in the regime of infinite sampling where statistical error is zero. The shaded region of each curve corresponds to the compounds with a real 1 \log_{10} unit improvement in selectivity. The speedup reflects the expected reduction in compounds that must be synthesized to reach a selectivity goal, and is calculated as the ratio of the percentage of compounds made with a real 1 \log_{10} unit improvement to the percentage of compounds that would be expected in the original distribution. (C) The speedup (y-axis, log scale) expected for 100x (2 \log_{10} units, or 2.8 kcal/mol) selectivity optimization as a function of correlation coefficient ρ . Each curve corresponds to a different value of $\sigma_{\text{sys},ij,\text{target}}$. (D) The speedup (y-axis) expected for 100x (2 \log_{10} units, or 2.8 kcal/mol) selectivity optimization as a function of number of compounds scored computationally (x-axis) and correlation coefficient ρ (each panel) for a method with per-target systematic error ($\sigma_{\text{sys},ij,\text{target}}$) of 1 kcal/mol in the regime of infinite sampling. After profiling, the top compounds that meet or surpass the synthesis rule (the predicted selectivity threshold a compound must reach to be triggered for synthesis, each curve) are synthesized, up to a maximum of 10 synthesized compounds. Error bars (y-axis) represent the 95% CI for 1000 replicates at each point. The expected speedup is calculated as the ratio of the number of synthesized compounds that have a true selectivity improvement of 2.8 kcal/mol (100x or 2 \log_{10} units) divided by the expectation of a compound showing a true selectivity improvement of 2.8 kcal/mol had the same number of compounds that were synthesized been drawn randomly from the underlying unit normal distribution. If no compounds were predicted to meet or surpass the synthesis rule, the speedup was assigned a default value of 1.

230 An experimental data set of CDK2/CDK9 inhibitors demonstrates the difficulty in achieving high
231 selectivity

232 To assess the correlation of errors in free energy predictions for selectivity, we set out to gather data sets
233 that met a number of criteria. We searched for data sets that contained binding affinity data for a number
234 of kinase targets and ligands in addition to crystal structures for each target with the same ligand.

235 This data set contains a congeneric series of ligands with experimental data for CDK2 and CDK9, with the
236 goal of potently inhibiting CDK9 and sparing CDK2. Based on a multiple sequence alignment of the 85 bind-
237 ing site residues identified in the kinase–ligand interaction fingerprints and structure (KLIFS) database [58,
238 59], CDK2 and CDK9 share 57% sequence identity (Table S1, Figure S1). For this CDK2/CDK9 data set [53],
239 ligand 12c was cocrystallized with CDK2/cyclin A (Figure 2A, left) and CDK9/cyclin T (Figure 2B, left), work that
240 was published in a companion paper [60]. In both CDK2 and CDK9, ligand 12c forms relatively few hydrogen
241 bond interactions with the kinase. Each kinase forms a pair of hydrogen bonds between the ligand scaffold
242 and a hinge residue (C106 in CDK9 and L83 in CDK2) that is conserved across all of the ligands in this se-
243 ries. CDK9, which has slightly lower affinity for ligand 12c (Figure 2C, right), forms an interaction between
244 the sulfonamide of ligand 12c and residue E107. On the other hand, CDK2 forms interactions between the
245 sulfonamide of ligand 12c and residues K89 and H84. The congeneric series of ligands contains a number
246 of difficult perturbations, particularly at substituent point R3 (Figure 2C, left). Ligand 12i also presented a
247 challenging perturbation, moving the 1-(piperazine-1-yl)ethanone from the *meta* to *para* location.

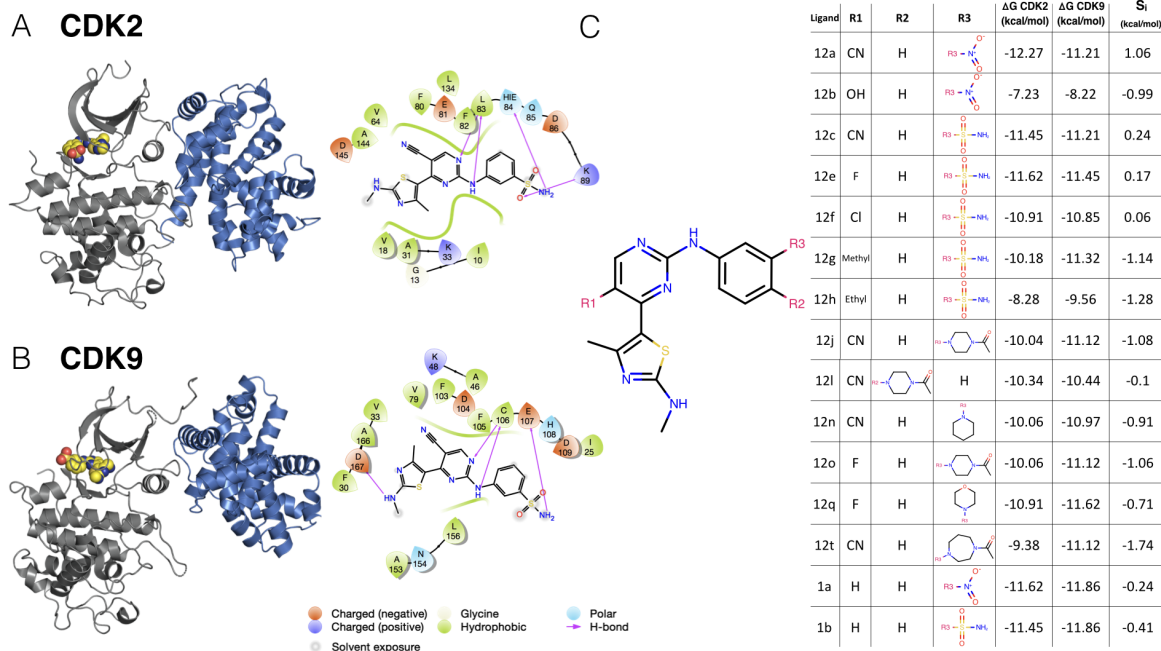
248 This congeneric series of ligands also highlights two of the challenges of working from publicly avail-
249 able data: First, the dynamic range of selectivity is incredibly narrow, with a mean S (CDK9 - CDK2) of -
250 0.65 kcal/mol, and a standard deviation of only 0.88 kcal/mol; the total dynamic range of this data set is 2.8
251 kcal/mol. Second, experimental uncertainties are not reported for the experimental measurements. This
252 data set reported K_i values calculated from measured IC_{50} , using the K_m (ATP) for CDK2 and CDK9 and [ATP]
253 from the assay using the Cheng-Prusoff equations [61]. Thus, for this and subsequent sets of ligands, the
254 random experimental uncertainty is assumed to be 0.3 kcal/mol based on previous work done to summa-
255 rize uncertainty in experimental data, assuming there is no systematic experimental error. While K_i values
256 are reported, these values are derived from IC_{50} measurements. A number of studies report on the re-
257 producibility of intra-lab IC_{50} measurements. These values range from as low as 0.22 kcal/mol [62], from
258 public data, to as high as 0.4 kcal/mol [6], which was estimated from internal data at Abbott Laboratories.
259 The assumed value of 0.3 kcal/mol falls within this range, and agrees well with the uncertainty reported
260 from Novartis for two different ligand series [63].

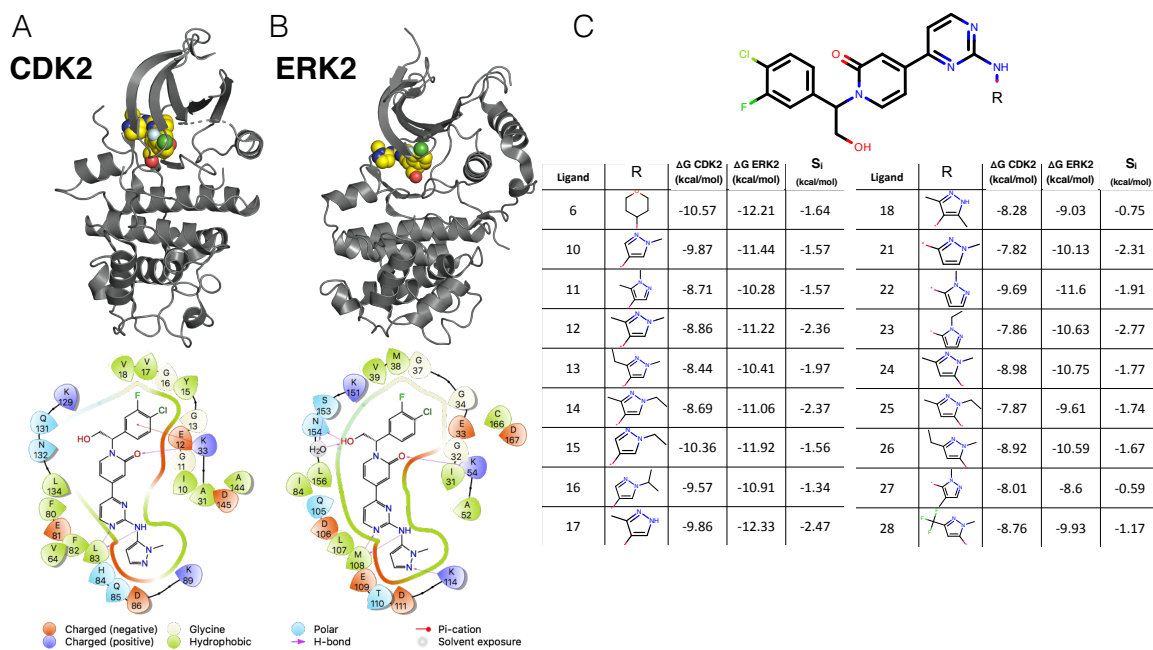
261 An experimental data set of CDK2/ERK2 inhibitors where greater selectivity was achieved

262 The CDK2/ERK2 data set from Blake et al. [54] also met the criteria described above, with the goal of
263 developing a potent ERK2 inhibitor. Based on a multiple sequence alignment of the KLIFs binding site
264 residues [58, 59], CDK2 and ERK2 share 52% sequence identity (Table S1, Figure S1), making them slightly
265 less closely related than CDK2 and CDK9 (57%). Note that while all three kinases belong to the CMGC fam-
266 ily and are closely related in the phylogenetic Manning tree, CDK2 and CDK9 belong to the CDK (Cyclin-
267 dependent kinase) subfamily, while ERK2 is part of the nearby MAPK (Mitogen-activated protein kinases)
268 subfamily. From a structural point of view, the two kinase pdb pairs used in this study are also very similar.
269 Binding site superposition revealed that both pdb pairs align well, only a marginally lower RMSD of 0.81 Å
270 was obtained for the CDK2/CDK9 pair compared to 0.92 Å for CDK2/ERK2 pair.

271 Crystal structures for both CDK2 (Figure 3A, top) and ERK2 (Figure 3B, top) were available with ligand 22
272 (according to the manuscript numbering scheme) co-crystallized. Of note, CDK2 was not crystallized with
273 cyclin A, despite cyclin A being included in the affinity assay reported in the paper [54].

274 CDK2 in this crystal structure (4BCK) adopts a DFG-in conformation with the α C helix rotated out, away
275 from the ATP binding site and breaking the conserved salt bridge between K33 and E51 (Figure S2A), indica-
276 tive of an inactive kinase [44, 64]. By comparison, the CDK2 structure from the CDK2/CDK9 data set adopts
277 a DFG-in conformation with the α C helix rotated in, forming the ionic bond between K33 and E51 indicative
278 of an active kinase, due to allosteric activation by cyclin A. While missing cyclins have caused problems for
279 free energy calculations in prior work, it is possible that the fully active, cyclin-bound conformation con-





280 tributes equally to binding affinity for all of the ligands in this series, and the high accuracy of the potency
 281 predictions (Figure 4, top left) is the result of fortuitous cancellation of errors.

282 The binding mode for this series is similar between both kinases. There is a set of conserved hydrogen
 283 bonds between the scaffold of the ligand and the backbone of one of the hinge residues (L83 for CDK2 and
 284 M108 for ERK2). The conserved lysine (K33 for CDK2 and K54 for ERK2), normally involved in the formation
 285 of a ionic bond with the αC helix, forms a hydrogen bond with the scaffold (Figure 3A and 3B, bottom) in
 286 both CDK2 and ERK2. However, in the ERK2 structure, the hydroxyl engages a crystallographic water as well
 287 as N154 in a hydrogen bond network that is not present in the CDK2 structure. The congeneric ligand series
 288 features a single solvent-exposed substituent. This helps to explain the narrow distribution of selectivities,
 289 with a mean selectivity of -1.74 kcal/mol (ERK2 - CDK2) and standard deviation of 0.56 kcal/mol; the total
 290 dynamic range of this data set is 2.2 kcal/mol. While the small standard deviation suggests that selectivity is
 291 difficult to drive with R-group substitution, the total dynamic range demonstrates that R-group substitutions
 292 can provide significant selectivity enhancements.

293 FEP+ calculations show smaller than expected errors for CDK2/CDK9 ΔS_{ij} predictions

294 Three replicates of FEP+ calculations were run on each target for both experimental data sets described
 295 above. The FEP+ predictions of the relative free energy of binding between ligands i and a reference com-
 296 pound (ref) for each target ($\Delta \Delta G_{i,ref,target}$) showed good accuracy and consistent results for all three replicates.
 297 The results for replicate 1 are reported in Figure 4 for both the CDK2 and ERK2 data set (bottom) and the
 298 CDK2/CDK9 data set (top), $\Delta \Delta G_{i,ref,target}$ is defined for each ligand i using a consistent reference compound
 299 within data sets.

$$\Delta \Delta G_{i,ref,target} = \Delta G_{i,target} - \Delta G_{reference,target} \quad (4)$$

300 The reference compounds (Compound 6 for CDK2/ERK2 and Compound 1a for CDK2/CDK9) were se-
301 lected because they were the initial compounds from which the reported synthetic studies were started.
302 Replicate 1 of the CDK2/ERK2 calculations is shown on the bottom of Figure 4, with an RMSE of $0.95^{1.25}_{0.63}$ and
303 $0.97^{1.22}_{0.70}$ kcal/mol to CDK2 and ERK2, respectively (where the lower and upper values indicate a 95% confi-
304 dence interval). The RMSE reported here is calculated for all of the $\Delta\Delta G_{i,ref,target}$ that were predicted. All of
305 the CDK2 and ERK2 $\Delta\Delta G_{i,ref,target}$ s were predicted within 1 log unit of the experimental value. The change
306 in selectivity (ΔS_{ij}) predictions show an RMSE of $1.41^{1.75}_{1.07}$ kcal/mol, with all the confidence intervals of the
307 predictions falling within 1 log unit of the experimental values (Figure 4, top right panel). This RMSE is com-
308 parable to the expected RMSE of 1.36, assuming the error from the CDK2/ERK2 calculations behaves in an
309 uncorrelated manner (Equation 3 where the correlation coefficient ρ is zero). This was consistent across all
310 three replicates of the calculations (Figure S6). The narrow dynamic range for selectivity combined with high
311 experimental and computational uncertainty highlight the challenges for predicting selectivity. When the
312 error of the calculated selectivity is comparable to the dynamic range of selectivity, then the calculations
313 cannot predict with statistical confidence whether any compound is more selective than the other.

314 Replicate 1 of the CDK2/CDK9 calculations are shown in the top panel of Figure 4. The CDK2 and CDK9
315 data sets show higher errors in $\Delta\Delta G_{i,ref,target}$ predictions, with an RMSE of $1.15^{1.59}_{0.67}$ and $2.10^{2.65}_{1.47}$ kcal/mol re-
316 spectively. This higher RMSE is driven by the reference compound, (Compound 1a) being poorly predicted,
317 particularly in CDK9. There are a number of outliers that fall outside of 1 \log_{10} unit from the experimental
318 value for CDK9. While the higher per target errors make predicting potency more difficult, the selectivity
319 predictions show an RMSE of $1.37^{1.66}_{1.04}$ kcal/mol. This observed RMSE is lower than what would be expected if
320 the error were completely uncorrelated between CDK2 and CDK9, propagated as in Equation 3 where the
321 correlation coefficient ρ is zero to get an expected value of 2.38 kcal/mol. This suggests that some correla-
322 tion in the error is leading to fortuitous cancellation of the systematic error, leading to more accurate than
323 expected predictions of ΔS_{ij} . These results were consistent across all three replicates of the calculation
324 (Figure S4).

325 Correlation of systematic errors accelerates selectivity optimization

326 To quantify the correlation coefficient (ρ) of the systematic error between targets, we built a Bayesian graphi-
327 cal model to separate the systematic error from the statistical error and quantify our confidence in estimates
328 of ρ (described in depth in Methods). Briefly, we modeled the absolute free energy (G) of each ligand in
329 each thermodynamic phase (ligand-in-complex and ligand-in-solvent, with G determined up to an arbitrary
330 additive constant for each phase) as in Equation 15. The model was chained to the FEP+ calculations by pro-
331 viding the $\Delta G_{phase,ij,target}^{calc}$ for each edge from the FEP+ maps (where j is now not necessarily the reference
332 compound) as observed data, as in Equation 17. As in Equation 19, the experimental data was modeled as
333 a normal distribution centered around the true free energy of binding ($\Delta G_{i,target}^{true}$) corrupted by experimental
334 error, which is assumed to be 0.3 kcal/mol from previous work done to quantify the uncertainty in publicly
335 available data [6]. ΔG values derived from reported IC_{50} s or K_i s, as described in the methods section, were
336 treated as data observations (Equation 19) and the $\Delta G_{i,target}^{true}$ was assigned a weak normal prior (Equation 20).

337 The correlation coefficient ρ was calculated for each Bayesian sample from the model posterior accord-
338 ing to equation 22. The CDK2/CDK9 calculations show strong evidence of correlation, with a correlation
339 coefficient of $0.72^{0.83}_{0.58}$ (Figure 5A, right) for replicate 1. The rest of the replicates showed strong agreement
340 (Figure S4). The joint marginal distribution of the error (ϵ) for each target (Figure 5A, left) is more diagonal
341 than symmetric, which is expected for cases in which ρ is high (Figure S3).

342 To quantify the expected speedup of selectivity with this level of correlation in the systematic errors for
343 CDK2/CDK9, we first calculated the per target systematic error $\sigma_{sys,ij,target}$ by taking the mean of the absolute
344 value of $\epsilon_{ij,target}$ where j is the reference compound 1a. Combining these estimates for the correlation coeffi-
345 cient (ρ) and the per target systematic errors ($\sigma_{sys,ij,target}$), we can compute $\sigma_{selectivity}$ and the expected speedup
346 in the regime of infinite sampling effort where there is no statistical error when the number of compounds
347 scored and synthesized is unlimited. The high correlation in errors for the CDK2/CDK9 calculations leads
348 to a speedup of 3x for 1 \log_{10} unit selectivity optimization and 10x for 2 \log_{10} unit selectivity optimization
349 (Figure 5A, right), despite the much high per target systematic errors ($\sigma_{sys,ij,target}$).

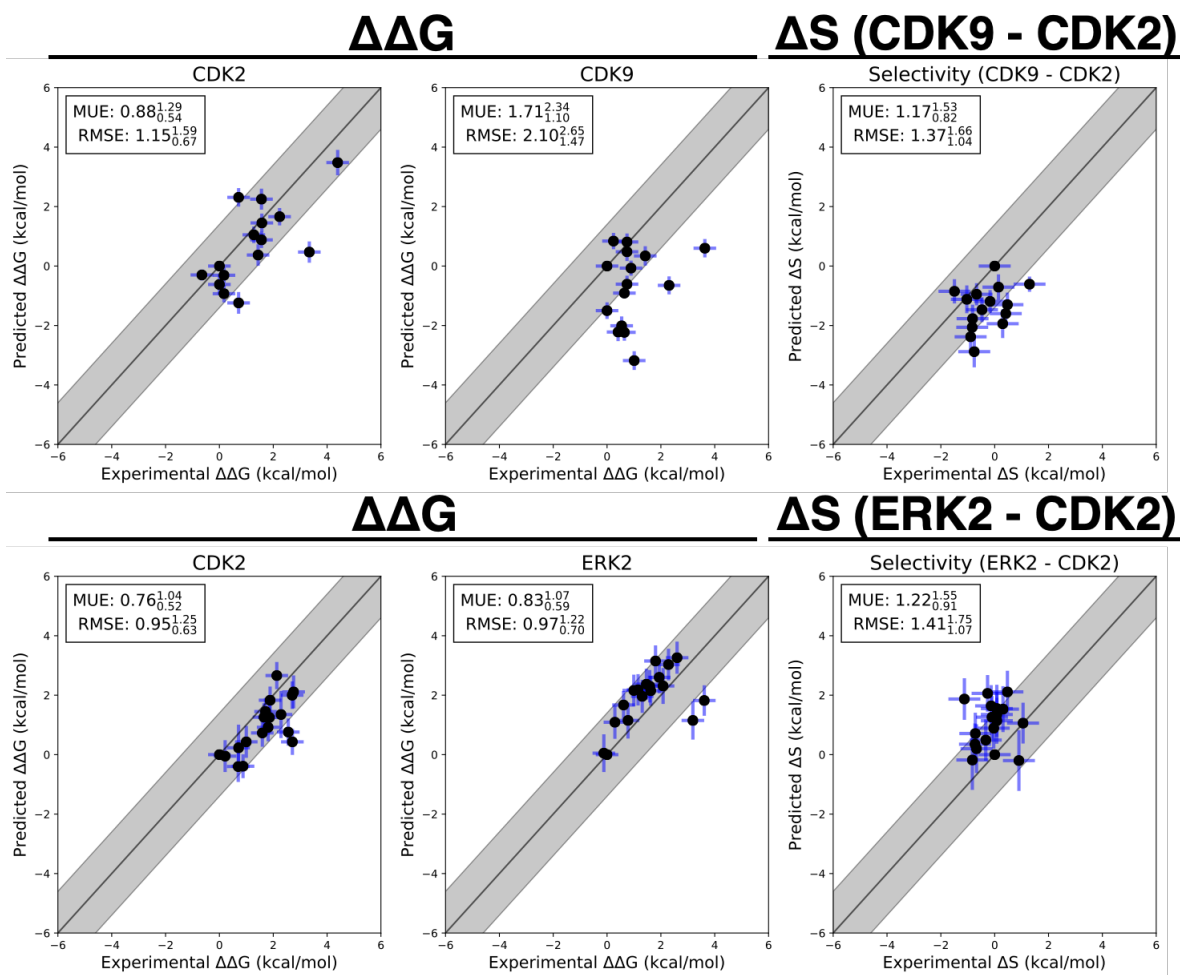


Figure 4. Selectivity predictions suggest correlation in systematic error

$\Delta\Delta G_{i,\text{ref},\text{target}}$ and $\Delta S_{i,\text{ref}}$ predictions for CDK2/CDK9 (top) from the Shao data sets and CDK2/ERK2 from the Blake data sets (bottom). The experimental values are shown on the X-axis and calculated values on the Y-axis. Each data point corresponds to a transformation between a ligand i to a set reference ligand (ref) for a given target. All values are shown in units of kcal/mol. The horizontal error bars show to the $\delta\Delta\Delta G_{ij}^{\text{exp}}$ based on the assumed uncertainty of 0.3 kcal/mol [6, 63] for each ΔG_i^{exp} . We show the estimated statistical error ($\sigma_{\text{stat},ij,\text{target}}$) as vertical blue error bars, which are one standard error. For selectivity, the errors were propagated under the assumption that they were completely uncorrelated. $\sigma_{\text{stat},ij,\text{target}}$ was estimated by calculating the standard deviation of $\Delta\Delta G_{ij,\text{target}}^{\text{FEP}}$ from the Bayesian model described in depth in **Methods**, where j is the reference compound. The black line indicates agreement between calculation and experiment, while the gray shaded region represent 1.36 kcal/mol (or 1 log₁₀ unit) error. The mean unsigned error (MUE) and root-mean squared error (RMSE) are shown on each plot with bootstrapped 95% confidence intervals.

350 The correlation coefficient ρ for replicate 1 of the CDK2/ERK2 calculations was quantified to be $0.44^{0.70}_{0.12}$,
351 indicating that the errors are moderately correlated between ERK2 and CDK2 (Figure 5B, right); this was
352 consistent with the distribution for ρ in replicate 3 (Figure S7), while the confidence interval of ρ for replicate
353 2 is much wider, indicating the correlation is weak.

354 Considering the speedup model where the number of compounds scored and synthesized is unlimited,
355 the modest correlation and low per target systematic errors for the CDK2/ERK2 calculations allow for a pre-
356 dicted 4–5x speedup for 1 \log_{10} unit selectivity optimization, and a 30–40x speedup for 2 \log_{10} unit selectivity
357 optimization (Figure 5B, right).

358 Using the correlation coefficient (ρ), $\sigma_{\text{stat},ij,\text{target}}$, and $\sigma_{\text{sys},ij,\text{target}}$ quantified from the Bayesian model for each
359 set of calculations, we can now calculate the y-axis error bars for the ΔS panels of Figure 4 according to the
360 proposed $\sigma_{\text{selectivity}}$ equation (Eq 3). Shown in Figure S9, we can see that $\sigma_{\text{selectivity}}$ accounts for most of the
361 disagreement between the predicted ΔS_{ij} and the experimental ΔS_{ij} .

362 Expending more effort to reduce statistical error can be beneficial in selectivity optimization

363 Up to this point, we have considered only systematic error in quantifying the speedup free energy calcula-
364 tions can enable for selectivity optimization, by assuming enough sampling is done to reduce the statistical
365 error for each target to zero. To begin understanding how statistical error impacts this speedup, we mod-
366 ified the model of speedup by additionally considering the per target statistical error ($\sigma_{\text{stat},\text{target}}$), which we
367 define in Equation 7 such that at the baseline effort, N , $\sigma_{\text{stat},ij,\text{target}}$ is 0.2 kcal/mol. In this definition, it takes
368 4x the sampling, or effort, to reduce statistical error by a factor of 2x. We assume that statistical error is
369 uncorrelated when propagating to two targets, and that $\sigma_{\text{sys},ij,\text{target}}$ is ≈ 1.0 kcal/mol for both targets [4, 62].
370 As shown in Figure 6, expending effort to reduce $\sigma_{\text{stat},ij,\text{target}}$ when ρ is less than 0.5 does not change the
371 expected speedup for the 100x selectivity threshold in meaningful way, suggesting that it is not worth run-
372 ning calculations longer than the default protocol in this case. However, when $\rho > 0.5$, the curves do start
373 to separate, particularly the 1/4x, 1x, and 4x effort curves. This suggests that when the correlation is high,
374 running longer calculations can produce net improvements in selectivity optimization speed. Interestingly,
375 the 16x, 48x, and ∞ effort curves do not differ greatly from the 4x effort curve, indicating that there are
376 diminishing returns to running longer calculations.

377 The estimated correlation coefficient is robust to Bayesian model assumptions

378 In order to better understand the statistical error in our calculations, we performed three replicates of our
379 calculations, and calculated the standard deviation of the cycle closure corrected $\Delta\Delta G$ for each edge of
380 the map, and compared that value to the cycle closure errors and Bennett errors reported for each edge
381 (Figure S8). For each set of calculations, the standard deviation suggests that the statistical error is between
382 0.1 and 0.3 kcal/mol, which is in good agreement with the reported Bennett error (Figure S8). However,
383 hysteresis in the closed cycles in the FEP map as reflected by the cycle closure error estimates indicate much
384 larger sampling errors than those estimated by the Bennett method or standard deviations of multiple runs,
385 suggesting that both the Bennett errors and standard deviation of multiple replicates are underestimating
386 the statistical error for these calculations. Based on this observation, we include a scaling parameter α in
387 the Bayesian error model (Eq. 16) to account for the BAR errors underestimating the cycle closure statistical
388 uncertainty. We also considered using a distribution with heavier tails, such as a Student's t-distribution, but
389 found the quantification of the correlation coefficient ρ insensitive to the use of either a scaling parameter
390 or heavier-tailed distributions (Figure S10).

391 Discussion and Conclusions

392 S is a useful metric for selectivity in lead optimization

393 There are a number of different metrics for quantifying the selectivity of a compound [55], which look at
394 selectivity from different views depending on the information trying to be conveyed. One of the earliest
395 metrics was the standard selectivity score, which conveyed the number of inhibited kinase targets in a broad
396 scale assay divided by the total number of kinases in the assay [65]. The Gini coefficient is a method that
397 does not rely on any threshold, but is highly sensitive to experimental conditions because it is dependent on

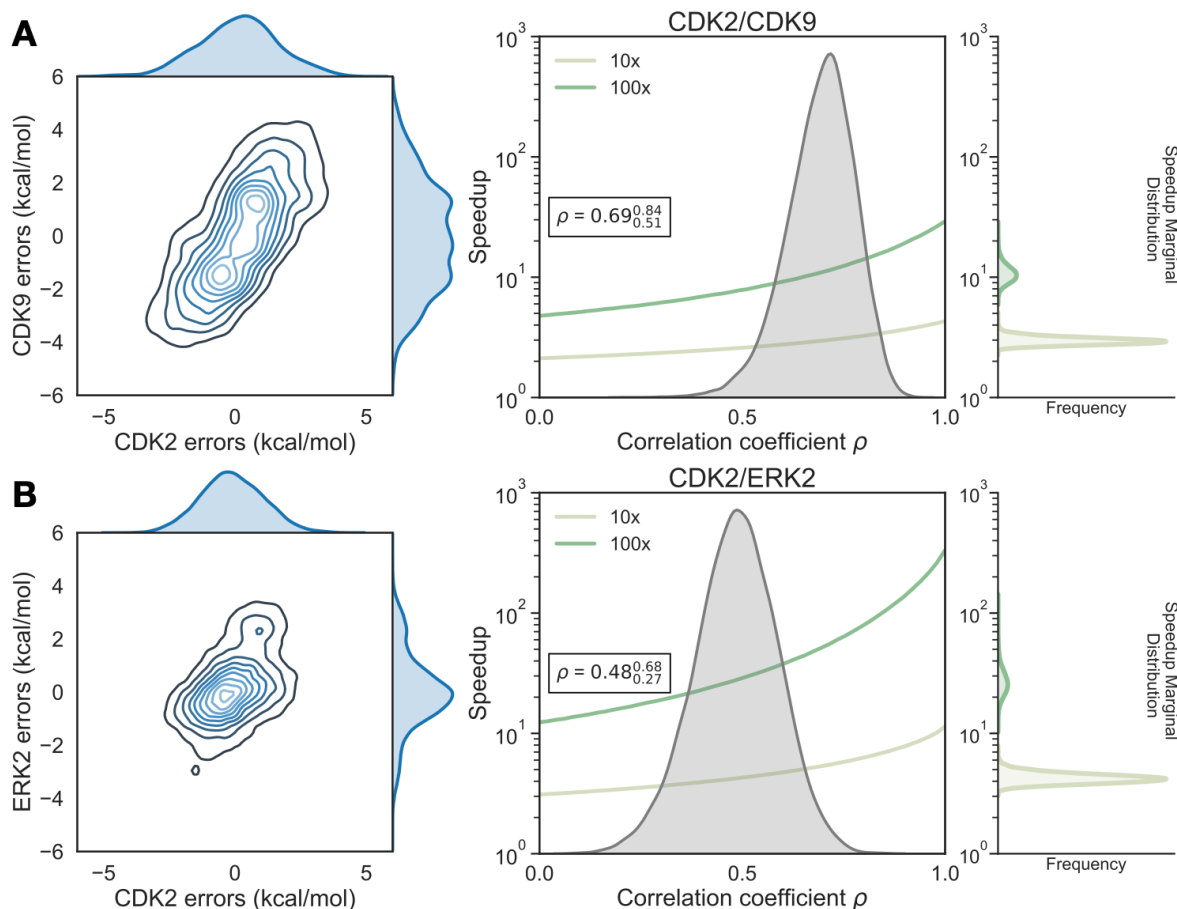


Figure 5. Correlation in systematic errors between targets can significantly accelerate selectivity optimization (A, left) The joint posterior distribution of the prediction errors for the more distantly related CDK2 (x-axis) and CDK9 (y-axis) from the Bayesian graphical model. **(A, right)** Speedup in selectivity optimization (y-axis), which estimates the reduction in compounds that must be synthesized to achieve a target selectivity when aided by free energy calculations, using the model where the number of compounds scored and synthesized is unlimited, as a function of correlation coefficient (x-axis). To calculate $\sigma_{\text{selectivity}}$, we calculate the per target systematic error ($\sigma_{\text{sys},ij,\text{target}}$) by taking the mean of $\epsilon_{ij,\text{target}}$ where j is the reference compound 1a. The posterior marginal distribution of the correlation coefficient (ρ) is shown in gray, while the expected speedup is shown for 100x (green curve) and 10x (yellow curve) selectivity optimization. The inserted box shows the mean and 95% confidence interval for the correlation coefficient. The marginal distribution of speedup is shown on the right side of the plot for both 100x (green) and 10x (yellow) selectivity optimization speedups. **(B)** As above, but for the more closely related CDK2/ERK2 selectivity data set using compound 6 as the reference.

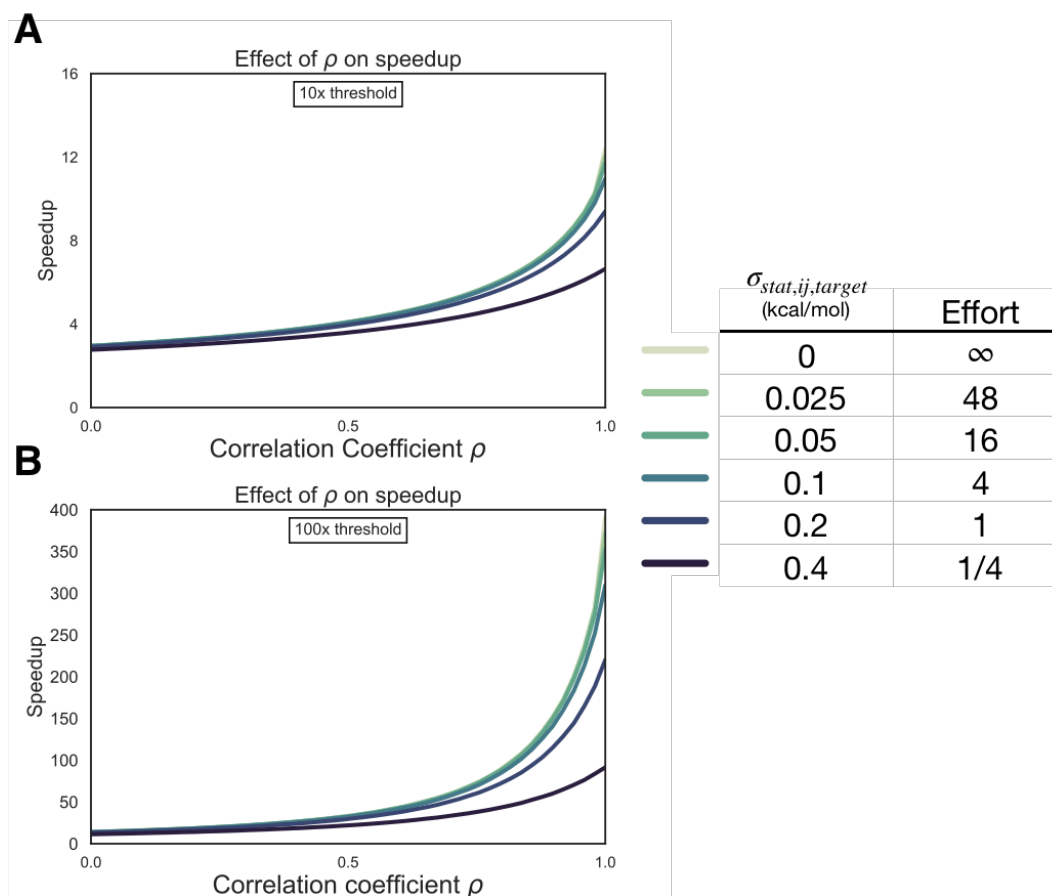


Figure 6. Reducing statistical uncertainty when systematic error correlation is high improves the speedup in selectivity optimization achievable with free energy calculations. (*left*) The speedup in selectivity (Y-axis) as a function of correlation coefficient (X-axis). Each curve represents a different per target statistical error ($\sigma_{stat,ij,target}$) for 10× (1 \log_{10} unit) (**A**) and 100× (2 \log_{10} unit) (**B**) thresholds (*right*) Table with the per target statistical error ($\sigma_{stat,ij,target}$), kcal/mol) corresponding to each curve on the left and a rough estimate of the generic amount of computational effort it would take to achieve that statistical uncertainty.

398 percent inhibition [66]. Other metrics take a thermodynamic approach to kinase selectivity and are suitable
399 for smaller panel screens [67, 68]. Here, we propose a more granular, thermodynamic view of selectivity
400 that is straightforward to calculate using free energy methods: the change in free energy of binding for a
401 given ligand between two different targets (S). S is a useful metric of selectivity in lead optimization once a
402 single, or small panel, of off-targets have been identified and the goal is to use physical modeling to either
403 improve or maintain selectivity within a lead series.

404 Systematic error correlation can accelerate selectivity optimization

405 We have demonstrated, using a simple numerical model that assumes unlimited synthetic and computa-
406 tional resources, the impact that free energy calculations with even weakly correlated systematic errors
407 can have on speeding up the optimization of selectivity in small molecule kinase inhibitors. While the ex-
408 pected speedup is dependent on the per target systematic error of the method ($\sigma_{\text{sys},ij,\text{target}}$), the speedup is
409 also highly dependent on the correlation of errors made for both targets. Unsurprisingly, free energy meth-
410 ods have greater impact as the threshold for selectivity optimization goes from 10x to 100x. While 100x
411 selectivity optimization is difficult to achieve, the expected benefit from free energy calculations is also quite
412 high, with speedups of one or two orders of magnitude possible. In a more realistic scenario, where the
413 number of compounds scored and synthesized is limited by resources, we have demonstrated using the
414 same numerical model that more stringent synthesis rules results in increased speedup from free energy
415 calculations. This holds true across different correlation coefficients (ρ), provided enough compounds are
416 scored. As our model shows, it is possible for stringent synthesis rules to provide benefits similar to oper-
417 ating with high systematic error correlation coefficients (ρ).

418 Two pairs of kinase test systems suggest systematic errors can be correlated

419 To quantify the correlation of errors in two example systems, we gathered experimental data for two con-
420 generic ligand series with experimental data for CDK2 and ERK2, as well as CDK2 and CDK9. These data
421 sets, which had crystal structures for both targets with the same ligand co-crystallized, exemplify the diffi-
422 culty in predicting selectivity. The dynamic range of selectivity for both systems is very narrow, with most
423 of the perturbations not having a major impact on the overall selectivity achieved. Further, the data was
424 reported without reliable experimental uncertainties, which makes quantifying the errors made by the free
425 energy calculations difficult. This issue is common when considering selectivity, as many kinase-oriented
426 high throughput screens are carried out at a single concentration and not highly quantitative.

427 The CDK9 calculations contained an outlier, compound 12h, that drove much of the prediction error for
428 that set. Compound 12e ($R1 = F$) is the most potent against CDK9 of the compounds in with a sulfonamide
429 at R3 (Figure 2). The addition of a single methyl group decreases the potency against CDK9 (compound 12g)
430 and while only slightly changing the affinity for CDK2. However, adding on another methyl group (compound
431 12h) results in an order of magnitude decrease in K_i for both CDK9 and CDK2. Crystal structures for both
432 kinases show that R1 points into a pocket formed by the backbone, and the sidechains of a Valine and
433 Phenylalanine. While ethyl at R1 in compound 12h is bulkier, the magnitude of the decrease in affinity for
434 both kinases is larger than might be expected, given that the pocket suggests an ethyl group would be well
435 accommodated in terms of fit and the hydrophobicity of the sidechains. For both kinases, the free energy
436 calculations predict that this addition should *improve* the potency, suggesting that it is possible that the
437 model is missing a chemical detail that might explain the trend seen in the experimental data. We expect
438 that these types of errors, which would be troubling when predicting potency alone, will drive the correlation
439 of systematic errors and fortuitously cancel when predicting selectivity.

440 Despite CDK2 and ERK2 belonging to different kinase subfamilies, the calculated correlation in the sys-
441 tematic error for two of the replicates suggests that fortuitous cancellation of errors may be applicable in a
442 wider range of scenarios than closely related kinases within the same subfamily. This may be driven by rela-
443 tively high binding site sequence identity between CDK2 and ERK2 (52% compared for 57% for CDK2/CDK9).
444 However, the confidence interval of the correlation is quite broad, including 0 for the lower bound for the
445 third replicate, suggesting that errors for more distantly related proteins will have only moderate, if any,
446 correlation.

447 Reducing statistical error is beneficial when systematic errors are correlated

448 In order to better understand if there are situations where it is beneficial to run longer calculations to mini-
449 mize statistical error to achieve a larger speedup in the synthesis of selective compounds, we built a numer-
450 ical model of the impact of statistical error in the context of different levels of systematic error correlation.
451 Our results suggest that unless the correlation coefficient ρ is highly positive for the two targets of interest,
452 there is not much benefit in running longer calculations. However, when the systematic error is reduced by
453 correlation, longer calculations can help realize large increases in speedup to achieve selectivity goals. Keep-
454 ing a running quantification of ρ for free energy calculations as compounds are made and the predictions
455 can be tested will allow for decisions to be made about whether running longer calculations is worthwhile.
456 It will also allow for an estimate of $\sigma_{\text{selectivity}}$, which is useful for estimating expected systematic error for
457 prospective predictions. Importantly, we expect that correlation will be modeling protocol dependent and
458 any changes to the way the system is modeled over the course of discovery program are expected to change
459 the observed correlation in the systematic error.

460 Larger data sets with a wide range of protein targets will enable future work

461 The data sets gathered here were limited by the total number of compounds, the small dynamic range for
462 selectivity (S), and the lack of reliable experimental uncertainties. The small size of the data set makes it
463 difficult to draw broad conclusions about the correlation in systematic errors. Understanding the degree
464 of correlation *a priori* based on structural or sequence similarity requires study on a larger range of targets
465 than the two pairs presented in this study. A larger data set that contained many protein targets, crystal
466 structures, and quantitative binding affinity data would be ideal to draw conclusions about the broader
467 prevalence of systematic error correlation.

468 This work demonstrates that correlation in the systematic errors can allow free energy calculations to
469 facilitate significant speedups in selectivity optimization for drug discovery projects. This is particularly im-
470 portant in kinase systems, where considering multiple targets is an important part of the development
471 process. The results suggest that free energy calculations can be particularly helpful in the design of kinase
472 polypharmacological agents, especially in cases where there is high correlation in the systematic errors
473 between multiple targets.

474 Methods

475 Numerical model of selectivity optimization speedup

476 To model the impact correlation of systematic error would have on the expected uncertainty for selectivity
477 predictions, $\sigma_{\text{selectivity}}$ was calculated using Equation 3 for 1000 evenly spaced values of the correlation coef-
478 ficient (ρ) from 0 to 1, for a number of combinations of per target systematic errors ($\sigma_{\text{sys},ij,1}$ and $\sigma_{\text{sys},ij,2}$). In
479 the regime of infinite sampling and zero statistical error, the second term reduces to zero.

$$\sigma_{\text{selectivity}} = \sqrt{\sigma_{\text{sys},ij,1}^2 + \sigma_{\text{sys},ij,2}^2 - 2\rho\sigma_{\text{sys},ij,1}\sigma_{\text{sys},ij,2} + \sigma_{\text{stat},ij,1}^2 + \sigma_{\text{stat},ij,2}^2} \quad (3)$$

480 The speedup in selectivity optimization that could be expected from using free energy calculations of a par-
481 ticular per target systematic error ($\sigma_{\text{sys},ij,\text{target}}$) was quantified as follows using NumPy (v 1.14.2). An original,
482 true distribution for the change in selectivity of 200 000 000 new compounds proposed with respect to a
483 reference compound was modeled as a normal distribution centered around 0 with a standard deviation of
484 1 kcal/mol. This assumption was made on the basis that the majority of selectivity is driven by the scaffold,
485 and R group modifications will do little to drive changes in selectivity. The 1 kcal/mol distribution is sup-
486 ported by the standard deviations of the selectivity in the experimental data sets referenced in this work,
487 which are all less than, but close to, 1 kcal/mol.

488 In this model, we suppose that each of proposed compound is triaged by a free energy calculation and
489 only proposed compounds predicted to increase selectivity by $\Delta S_{ij} \geq 1.4$ kcal/mol (1 \log_{10} unit) with respect
490 to a reference compound would be synthesized. Based on reported estimates in the literature, we pre-
491 sume that relative free energy calculations have a per-target systematic error $\sigma_{\text{sys},ij,\text{target}} \approx 1$ kcal/mol [4], and
492 explore the impact of the correlation coefficient ρ governing the correlation of these predictions between

493 targets. The standard error in predicted selectivity, $\sigma_{\text{selectivity}}$, is given by Equation 3. When sampling is infi-
 494 nite and $\sigma_{\text{stat},ij,\text{target}}$ is zero, $\sigma_{\text{selectivity}}$ is driven entirely by the systematic error component ($\sigma_{\text{sys},ij,\text{target}}$), resulting
 495 in the error in predicted change in selectivity ΔS_{ij} modeled as a normal distribution centered around 0 with
 496 a standard deviation of $\sigma_{\text{sys},ij,\text{target}}$ and added to the "true" ΔS_{ij} ,

$$\Delta S_{ij, \text{predicted}} = \Delta S_{ij, \text{true}} \left(\mathcal{N}_{\text{true}}(\mu = 0, \sigma^2 = 1) \right) + \Delta S_{\text{systematic error}} \left(\mathcal{N}_{\text{error}}(\mu = 0, \sigma_{\text{sys},ij,\text{target}}^2(\rho)) \right) \quad (5)$$

497 We ignore the potential complication of finite experimental error in this thought experiment, presuming
 498 the experimental uncertainty is sufficiently small as to be negligible.

499 The *speedup* in synthesizing molecules that reach this 10x selectivity gain threshold is calculated, as a
 500 function of ρ , as the ratio of the number of compounds that exceed the selectivity threshold in the case that
 501 molecules predicted to fall below this threshold by free energy calculations were triaged and not synthe-
 502 sized, divided by the number of compounds that exceeded the selectivity threshold without the benefit of
 503 free energy triage. This process was repeated for a 100x (2.8 kcal/mol, 2 log₁₀ unit) selectivity optimization
 504 and 50 linearly spaced values of the correlation coefficient (ρ) between 0 and 1, for four values of $\sigma_{\text{selectivity}}$,
 505 using a sample size of 4×10^7 compounds.

506 The above model assumes that the number of compounds scored and synthesized is essentially unlim-
 507 ited. To assess the impact these methods might have on real drug discovery projects, where the number
 508 of compounds scored and synthesized are limited by computational and chemistry resources, we altered
 509 the above model to consider the number of compounds scored, the number of compounds triggered for
 510 synthesis, and the threshold a compound needed to reach in order to be considered for synthesis.

511 We repeated the mode detailed above, this time scoring only the following numbers of compounds: 10,
 512 50, 100, 200, 500, the range from 1000 to 10000 in steps of 1000, and the range from 10000 to 100 000 in
 513 steps of 2000. Compounds were drawn from a true distribution of $\Delta S_{ij, \text{true}} \left(\mathcal{N}_{\text{true}}(\mu = 0, \sigma^2 = 1) \right)$ and triaged
 514 using a free energy method as detailed above with a per-target systematic error ($\sigma_{\text{sys},ij,\text{target}}$) of 1 kcal/mol.
 515 The top predicted compounds that meet or surpass a synthesis rule, up to a maximum of 10 compounds,
 516 are selected for synthesis. Here, we consider synthesis rules of 100x, 500x and 1000x when trying to design
 517 100x (2.8 kcal/mol, 2 log₁₀ unit) improvements in selectivity. The *speedup* was calculated as the number of
 518 synthesized compounds whose $\Delta S_{ij, \text{true}}$ reaches the desired 100x threshold divided by the expected value
 519 ($E_{\text{selective}}$) for a selective compound given the number of synthesized compounds. This expectation can be
 520 calculated as,

$$E_{\text{selective}} = P(\Delta S_{ij} > \text{threshold} | \mathcal{N}_{\text{true}}) * n_{\text{synthesized}} \quad (6)$$

521 Where $P(\Delta S_{ij} > \text{threshold} | \mathcal{N}_{\text{true}})$ is the probability $\Delta S_{ij, \text{true}}$ for some compound is better than a particular
 522 selectivity threshold given the distribution of $\Delta S_{ij, \text{true}} \left(\mathcal{N}_{\text{true}}(\mu = 0, \sigma^2 = 1) \right)$ for 100 000 000 compounds, and
 523 $n_{\text{synthesized}}$ is the number of compounds synthesized. If no compounds were predicted to meet or surpass
 524 the synthesis rule, the speedup was assigned a default value of 1. We performed 1000 replicates of each
 525 condition and report the mean and 95 % CI in Figure 1D.

526 Numerical model of impact of statistical error on selectivity optimization

527 To model the impact of finite statistical error in the alchemical free energy calculations, a similar scheme was
 528 used with the following modifications: Each proposed compound was triaged by a free energy calculation
 529 with a per target systematic error ($\sigma_{\text{sys},ij,\text{target}}$) of 1.0 kcal/mol [4] and a specified correlation coefficient ρ . A
 530 $\sigma_{\text{selectivity}}$ was calculated according to Equation 3, this time considering the statistical terms as non-negligible.
 531 The per target statistical error ($\sigma_{\text{stat},ij,\text{target}}$) was defined as,

$$\sigma_{\text{stat},ij,\text{target}} = \frac{\sigma_{\text{stat},\text{base}}}{\sqrt{N}} \quad (7)$$

532 where N is the relative effort put into running sampling the calculation and $\sigma_{\text{stat},\text{base}}$ is such that when $N = 1$,
 533 $\sigma_{\text{stat},ij,\text{target}} = 0.2$ kcal/mol. The statistical error is propagated assuming it is uncorrelated, as independent sets
 534 of calculations are used for each target, giving us the second set of terms in 3. This gives an updated model

535 for the error in predicted change in selectivity ΔS_{ij} . The systematic and statistical errors were modeled as
536 Gaussian noise added to the true distribution,

$$\Delta S_{ij,\text{predicted}} = \Delta S_{ij,\text{true}} \left(\mathcal{N}_{\text{true}}(\mu = 0, \sigma^2 = 1) \right) + \Delta S_{\text{systematic error}} \left(\mathcal{N}_{\text{systematic}}(\mu = 0, \sigma_{\text{sys},ij,\text{target}}^2(\rho)) \right) \quad (8) \\ + \Delta S_{\text{statistical error}} \left(\mathcal{N}_{\text{statistical}}(\mu = 0, \sigma_{\text{stat},ij,\text{target}}^2) \right)$$

537 Any compound predicted to have an improvement in selectivity of above the threshold (either 1.4 kcal/mol
538 ($1 \log_{10}$ units) or 2.8 kcal/mol ($2 \log_{10}$ units)) would then be made and have its selectivity experimentally
539 measured, using an experimental method with perfect accuracy. The speedup value for each value of ρ is
540 calculated as previously described.

541 Binding Site Similarity analysis

542 To quantify the similarity between the different kinase pairs, a structure-informed binding site sequence
543 comparison was performed. In the KLIFS database, the binding site of typical human kinases is defined
544 by 85 residues, comprising known kinase motives (DFG, hinge, G-loop, aC-helix, ...), which cover potential
545 interactions with type I-IV inhibitors [58, 59]. KLIFS provides a multiple sequence alignment in which each
546 kinase sequence is mapped to these 85 binding site residues. This mapping was used to calculate the
547 sequence identify between the three kinases CDK2, CDK9, and ERK2 used in this study (Figure S1 and Table
548 S1). The score shows the percentage of identical residues between two kinases with respect to the 85
549 positions.

550 For structural comparison, the respective pdb's of the two kinases were downloaded from the pdb (CDK2-
551 4bci/CDK9-4bck) and CDK2-5k4j/ERK2-5k4i). PyMol v.2.3.0 was used for preprocessing and alignment of the
552 structures. For all structures only chain A was kept. Additionally, for structure 4bck alternate location C
553 was chosen only. Next, binding sites were selected as all residues within 10 Å of the co-crystallized ligand,
554 yielding. Finally, the respective binding site pairs were aligned using PyMol's default align function and the
555 RMSD was returned. The following is an example command: `create [pdb]_bs, byres [pdb]_A within 10`
556 `of ([pdb]_A and resn [lig_name])`

557 Extracting the binding free energy ΔG from reported experimental data

558 K_i values were derived from IC_{50} measurements reported for the ERK2/CDK2 data set (Figure 3), assuming
559 Michaelis-Menten binding kinetics for an ATP-competitive inhibitor,

$$K_i = \frac{IC_{50}}{1 + \frac{[S_0]}{K_m}} \quad (9)$$

560 Where the Michaelis-Menten constant for ATP (K_m (ATP)) is much larger than the initial concentration of ATP,
561 S_0 , so that $IC_{50} \approx K_i$.

562 These K_i values were then used to calculate a ΔG (Equation 10),

$$\Delta G = -k_B T \ln K_i \quad (10)$$

563 Here, k_B is the Boltzmann constant and T is absolute temperature (taken to be room temperature, $T \sim$
564 300K).

565 For the CDK2/CDK9 data set, the authors note that the assumption K_m (ATP) $\gg S_0$ does *not* hold, and
566 report K_i s derived from their IC_{50} measurements using the K_m (ATP) for each kinase, as well as the S_0 from
567 their assay. These values were then converted to ΔG using Equation 10. For both data sets, these derived
568 ΔG were used to calculate $\Delta\Delta G$ between ligands for each kinase target.

569 As mentioned above, the assumption that K_m (ATP) $\gg S_0$ may not always hold, and changes in IC_{50} may
570 be driven by factors other than changes in ligand binding affinity. However, these assumptions have been
571 used successfully to estimate relative free energies previously [62, 69]. Further, data was taken from the
572 same lab and assay for each target. By using assays with the same kinase construct and ATP concentration,
573 the relative free energies ($\Delta\Delta G_{ij}$) should be well determined for compounds within the assay. Even if the
574 absolute free energies (ΔG_i) are off due to uncertainties in K_m (ATP) or S_0 , they will be off by the same
575 constant, which will cancel when calculating $\Delta\Delta G_{ij}$.

576 Structure Preparation

577 Structures from the Shao [53] (CDK2/CDK9), Hole [60] (CDK2/CDK9), and Blake [54] (CDK2/ERK2) papers
578 were downloaded from the PDB [70], selecting structures with the same co-ligand crystallized.

579 For the Shao (CDK2/CDK9) data set, PDB IDs 4BCK (CDK2) and 4BCI (CDK9) were selected, which have
580 ligand 12c cocrystallized. For the Blake data set (ERK2/CDK2), 5K4J (CDK2) and 5K4I (ERK2) were selected,
581 cocrystallized with ligand 21. The structures were prepared using Schrodinger's Protein Preparation Wiz-
582 ard [71] (Maestro, Release 2017-3). This pipeline modeled in internal loops and missing atoms, added hy-
583 drogens at the reported experimental pH (7.0 for the Shao data set, 7.3 for the Blake data set) for both the
584 protein and the ligand. All crystal waters were retained. The ligand was assigned protonation and tautomer
585 states using Epik at the experimental pH \pm 2, and hydrogen bonding was optimized using PROPKA at the
586 experimental pH \pm 2. Finally, the entire structure was minimized using OPLS3 with an RMSD cutoff of 0.3Å.

587 Ligand Pose Generation

588 Ligands were extracted from the publication entries in the BindingDB as 2D SMILES strings. 3D conforma-
589 tions were generated using LigPrep with OPLS3 [4]. Ionization state was assigned using Epik at experimen-
590 tal pH \pm 2. Stereoisomers were computed by retaining any specified chiralities and varying the rest. The
591 tautomer and ionization state with the lowest Epik state penalty was selected for use in the calculation. Any
592 ligands predicted to have a positive or negative charge in its lowest Epik state penalty was excluded, with
593 the exception of Compound 9 from the Blake data set. This ligand was predicted to have a +1 charge for its
594 lowest state penalty state. The neutral form the ligand was include for the sake of cycle closure in the FEP+
595 map, but was ignored for the sake any analysis afterwards. Ligand poses were generated by first aligning
596 to the co-crystal ligand using the Largest Common Bemis-Murcko scaffold with fuzzy matching (Maestro,
597 Release 2017-3). Ligands that were poorly aligned or failed to align were then aligned using Maximum
598 Common Substructure (MCSS). Finally, large R-groups conformations were sampled with MM-GBSA using a
599 common core restraint, VSGB solvation model, and OPLS3 force field. No flexible residues were defined for
600 the protein.

601 Free Energy Calculations

602 The FEP+ panel (Maestro, Release 2017-3) was used to generate perturbation maps. FEP+ calculations were
603 run using the FEP+ panel from Maestro release 2018-3 in order to take advantage of the newest force
604 field (OPLS3e) parameters available at the time. Any missing ligand torsions were fit using the automated
605 FFbuilder protocol [7]. Custom charges were assigned using the OPLS3e force field using input geometries,
606 according to the automated FEP+ workflow in Maestro Release 2018-3. Neutral perturbations were run for
607 15 ns per replica, using an NPT ensemble and water buffer size of 5Å. The SPC water model was used. A
608 GCMC solvation protocol was used to sample buried water molecules in the binding pocket prior to the
609 calculation, which discards any retained crystal waters.

610 Statistical Analysis of FEP+ calculations

611 To quantify the overall errors in the FEP+ calculations, we computed the mean unsigned error (MUE),

$$\text{MUE} = \frac{\sum_0^n |\Delta\Delta G_{i,\text{ref},\text{target}}^{\text{calc}} - \Delta\Delta G_{i,\text{ref},\text{target}}^{\text{exp}}|}{n} \quad (11)$$

612 and the root mean squared error (RMSE)

$$\text{RMSE} = \sqrt{\frac{\sum_0^n (\Delta\Delta G_{i,\text{ref},\text{target}}^{\text{calc}} - \Delta\Delta G_{i,\text{ref},\text{target}}^{\text{exp}})^2}{n}} \quad (12)$$

613 MUE and RMSE were computed for $\Delta\Delta G_{ij,\text{target}}$. For each ligand i , $\Delta\Delta G_{i,\text{ref},\text{target}}$ is defined where ref is a
614 reference compound.

$$\Delta\Delta G_{i,\text{ref},\text{target}} = \Delta G_{i,\text{target}} - \Delta G_{\text{reference},\text{target}} \quad (13)$$

615 For the CDK2/CDK9 data set, compound 1a was used as the reference compound, as it was the first com-
616 pound from which the others in the series were derived. For the CDK2/ERK2 data set, compound 6 was

617 used as the reference compound, since it was the compound from which the investigation was launch. A
 618 metabolite of compound 6 (not included in the data set here) was used as the starting compound from
 619 which the rest were derived. To account for the finite ligand sample size, we used 10 000 replicates of boot-
 620 strapping with replacement to estimate 95% confidence intervals. The code used to bootstrap these values
 621 is available on GitHub [<https://github.com/choderalab/selectivity>].

622 To compute the per-target statistical error ($\sigma_{\text{stat},ij,\text{target}}$) for each i,ref pair of ligands, we used the standard
 623 deviation of $\Delta\Delta G_{ij,\text{target}}^{\text{FEP}}$, where j is the reference compound, from the Bayesian model described in depth
 624 below in the **Methods** section. To compute the per target systematic error ($\sigma_{\text{sys},ij,\text{target}}$), we calculated the
 625 mean of $\epsilon_{ij,\text{target}}$, where j is the reference compound, described in equation 21 in the Bayesian Model section
 626 of the **Methods**.

627 **Quantification of the correlation coefficient ρ**

628 To quantify ρ , we built a Bayesian graphical model using `pymc3` 3.5 [72] and `theano` 1.0.3 [73]. All code for
 629 this model is available on GitHub [<https://github.com/choderalab/selectivity>].

630 For each phase (complex and solvent), the prior for the absolute free energy (G) of ligand i (up to an arbi-
 631 trary additive constant for each thermodynamic phase, ligand-in-complex or ligand-in-solvent), was treated
 632 as a normal distribution (Equation 15).

$$G_{i,\text{target}}^{\text{phase}} \sim \mathcal{N}(\mu = 0, \sigma = 25.0 \text{ kcal/mol}) \quad (14)$$

633 To improve sampling efficiency, for each phase, one ligand was chosen as the reference, and pinned to an
 634 absolute free energy of $G = 0$, with a standard deviation of 1 kcal/mol.

$$G_{1,\text{target}}^{\text{phase}} \sim \mathcal{N}(\mu = 0, \sigma = 1.0 \text{ kcal/mol}) \quad (15)$$

635 For each edge of the FEP map (ligand $i \rightarrow$ ligand j), there is a contribution from dummy atoms, that was
 636 modeled as in Equation 16. Note that here, unlike what was done in Figure 4, ligand j is not necessarily a
 637 reference compound.

$$c_{i,j} \sim \mathcal{N}(\mu = 0, \sigma = 25.0 \text{ kcal/mol}) \quad (16)$$

638 The model was conditioned by including data from the FEP+ calculation.

$$\Delta G_{\text{phase},ij,\text{target}}^{\text{calc}} \sim \mathcal{N}(G_{j,\text{target}}^{\text{phase}} - G_{i,\text{target}}^{\text{phase}} - \alpha \delta^2 \Delta G_{\text{phase},ij,\text{target}}^{\text{BAR}}) \quad (17)$$

639 where $\delta^2 \Delta G_{\text{phase},ij,\text{target}}^{\text{BAR}}$ is the reported BAR uncertainty from the calculation, and $\Delta G_{\text{phase},ij,\text{target}}^{\text{calc}}$ is the BAR
 640 estimate of the free energy for the perturbation between ligands i and j in a given phase. α is a scaling
 641 parameter shared by all $\Delta G_{\text{phase},ij,\text{target}}^{\text{calc}}$ for each target. Such scaling is necessary to account for the BAR
 642 statistical uncertainty underestimating cycle closure statistical uncertainty of our calculations, shown by
 643 Figure S8.

644 From this, we can calculate the $\Delta G_{i,\text{target}}^{\text{FEP}}$ for each ligand and target,

$$\Delta G_{i,\text{target}}^{\text{FEP}} = G_{i,\text{target}}^{\text{complex}} - G_{i,\text{target}}^{\text{solvent}} \quad (18)$$

645 From $\Delta G_{i,\text{target}}^{\text{FEP}}$, we calculated $\Delta\Delta G_{ij,\text{target}}^{\text{FEP}}$ for each pair of ligands, filtering out pairs where i and j are the
 646 same ligand and where the reciprocal was already included.

647 The experimental binding affinity was treated as a true value ($\Delta G_{i,\text{target}}^{\text{true}}$) corrupted by experimental un-
 648 certainty, which is assumed to be 0.3 kcal/mol [6]. There are a number of studies that report on the re-
 649 producibility and uncertainty of intra-lab IC_{50} measurements, ranging from as small as 0.22 kcal/mol [62]
 650 to as high as 0.4 kcal/mol [6]. The assumed value falls within this range and is in good agreement with the
 651 uncertainty reported from multiple replicate measurements in internal data sets at Novartis [63].

652 The values reported in the papers ($\Delta G_{i,\text{target}}^{\text{obs}}$) were treated as observations from this distribution (Equa-
 653 tion 19),

$$\Delta G_{i,\text{target}}^{\text{obs}} \sim \mathcal{N}(\mu = \Delta G_{i,\text{target}}^{\text{true}}, \sigma = 0.3 \text{ kcal/mol}) \quad (19)$$

654 $\Delta G_{i,\text{target}}^{\text{true}}$ was assigned a weak normal prior, as in Equation 20,

$$\Delta G_{i,\text{target}}^{\text{true}} = \mathcal{N}(\mu = 0, \sigma = 50 \text{ kcal/mol}) \quad (20)$$

655 $\Delta \Delta G_{ij,\text{target}}^{\text{true}}$ for each pair of ligands was calculated from $\Delta G_{i,\text{target}}^{\text{true}}$, filtering out pairs where i and j are the
656 same ligand and where the reciprocal was already included as above.

657 The error for a given ligand was calculated as

$$\epsilon_{ij,\text{target}} = \Delta \Delta G_{ij,\text{target}}^{\text{FEP}} - \Delta \Delta G_{ij,\text{target}}^{\text{true}} \quad (21)$$

658 From these ϵ values, we calculated the correlation coefficient, ρ , from the sampled errors for the finite set
659 of molecules for which measurements were available,

$$\rho = \frac{\text{cov}(\epsilon_{\text{target}1}, \epsilon_{\text{target}2})}{\sigma_{\epsilon \text{ target}1} \sigma_{\epsilon \text{ target}2}} \quad (22)$$

660 where $\sigma_{\epsilon \text{ target}2}$ is the standard deviation of $\epsilon_{ij,\text{target}}$.

661 To quantify ρ from these calculations, the default NUTS sampler with `jitter+adapt_diag` initialization,
662 3 000 tuning steps, and the default target accept probability was used to draw 20 000 samples from the
663 model.

664 Calculating the marginal distribution of speedup

665 To quantify the expected speedup from the calculations we ran, we utilized 10^4 replicates of the scheme
666 detailed above to calculate the speedup given parameters ρ , $\sigma_{\text{sys},ij,1}$, and $\sigma_{\text{sys},ij,2}$, in the regime of infinite
667 effort and zero statistical error. Using Numpy 1.14.2, ρ was drawn from a normal distribution with the
668 mean and standard deviation from the posterior distribution of ρ from the Bayesian Graphical model. The
669 per-target systematic errors, $\sigma_{\text{sys},ij,1}$ and $\sigma_{\text{sys},ij,2}$, were estimated from the mean of the absolute value of $\epsilon_{ij,1}$
670 and $\epsilon_{ij,2}$, which are the magnitude of errors from the Bayesian graphical model. $\sigma_{\text{selectivity}}$ was calculated
671 using Equation 3. 10^6 molecules were proposed from true normal distribution, as above. The error of the
672 computational method was modeled as in Equation 5.

673 **Data Availability**

674 All curated starting structures, FEP+ results, and data analysis scripts and notebooks are available on GitHub:

675 <https://github.com/choderalab/selectivity>

676 **Acknowledgments**

677 The authors are grateful to Patrick Grinaway (ORCID: [0000-0002-9762-4201](https://orcid.org/0000-0002-9762-4201)) for useful discussions about
678 Bayesian statistics and Mehtap Işık (ORCID: [0000-0002-6789-952X](https://orcid.org/0000-0002-6789-952X)) for useful discussion about kinase in-
679 hibitor protonation states. SKA is grateful to Haoyu S. Yu, Wei Chen, and Dmitry Lupyan for advice on
680 running FEP+ calculations.

681 **Funding**

682 Research reported in this publication was supported by the National Institute for General Medical Sciences
683 of the National Institutes of Health under award numbers R01GM121505 and P30CA008748. SKA acknowl-
684 edges financial support from Schrödinger and the Sloan Kettering Institute. JDC acknowledges financial
685 support from Cycle for Survival and the Sloan Kettering Institute.

686 **Disclosures**

687 JDC was a member of the Scientific Advisory Board for Schrödinger, LLC during part of this study. JDC is a cur-
688 rent member of the Scientific Advisory Board of OpenEye Scientific Software and a consultant for Foresite
689 Labs. The Chodera laboratory receives or has received funding from multiple sources, including the National
690 Institutes of Health, the National Science Foundation, the Parker Institute for Cancer Immunotherapy, Re-
691 lay Therapeutics, Bayer, Entasis Therapeutics, Silicon Therapeutics, EMD Serono (Merck KGaA), AstraZeneca,
692 XtalPi, the Molecular Sciences Software Institute, the Starr Cancer Consortium, the Open systematic Consor-
693 tium, Cycle for Survival, a Louis V. Gerstner Young Investigator Award, and the Sloan Kettering Institute. A
694 complete funding history for the Chodera lab can be found at <http://choderalab.org/funding>

695 **Author Contributions**

696 Conceptualization: SKA, LW, RA, JDC; Methodology: SKA, LW, JDC; Formal Analysis: SKA, JDC, LW; Data Cura-
697 tion: SKA, SP; Investigation: SKA, SP, AV; Writing – Original Draft: SKA, JDC; Writing – Review & Editing: SKA,
698 JDC, LW, AV, RA; Visualization: SKA, JDC, LW; Supervision: LW, JDC, RA; Project Administration: SKA, LW, JDC,
699 RA; Funding Acquisition: RA, JDC; Resources: LW, JDC

References

- 700
701 [1] Chodera, J. D.; Mobley, D. L.; Shirts, M. R.; Dixon, R. W.; Branson, K.; Pande, V. S. Alchemical free energy
702 methods for drug discovery: progress and challenges. *Curr. Opin. Struct. Biol.* **2011**, *21*, 150–160.
- 703 [2] Huang, J.; MacKerell, A. D. CHARMM36 All-Atom Additive Protein Force Field: Validation Based on Com-
704 parison to NMR Data. *J. Comput. Chem.* **2013**, *34*, 2135–2145.
- 705 [3] Maier, J. A.; Martinez, C.; Kasavajhala, K.; Wickstrom, L.; Hauser, K. E.; Simmerling, C. ff14SB: Improving
706 the Accuracy of Protein Side Chain and Backbone Parameters from ff99SB. *J. Chem. Theory Comput.*
707 **2015**, *11*, 3696–3713.
- 708 [4] Harder, E. et al. OPLS3: A Force Field Providing Broad Coverage of Drug-like Small Molecules and Pro-
709 teins. *J. Chem. Theory Comput.* **2016**, *12*, 281–296.
- 710 [5] Cournia, Z.; Allen, B.; Sherman, W. Relative Binding Free Energy Calculations in Drug Discovery: Recent
711 Advances and Practical Considerations. *Journal of chemical information and modeling* **2017**, *57*, 2911–
712 2937.
- 713 [6] Brown, S. P.; Muchmore, S. W.; Hajduk, P. J. Healthy Skepticism: Assessing Realistic Model Performance.
714 *Drug Discov. Today* **2009**, *14*, 420 – 427.
- 715 [7] Abel, R.; Mondal, S.; Masse, C.; Greenwood, J.; Harriman, G.; Ashwell, M. A.; Bhat, S.; Wester, R.; Frye, L.;
716 Kapeller, R.; Friesner, R. A. Accelerating drug discovery through tight integration of expert molecular
717 design and predictive scoring. *Curr. Opin. Struct. Biol.* **2017**, *43*, 38–44.
- 718 [8] Lovering, F. et al. Imidazotriazines: Spleen Tyrosine Kinase (Syk) Inhibitors Identified by Free-Energy
719 Perturbation (FEP). *ChemMedChem* **2016**, *11*, 217–233.
- 720 [9] Ciordia, M.; Pérez-Benito, L.; Delgado, F.; Trabanco, A. A.; Tresadern, G. Application of Free Energy
721 Perturbation for the Design of BACE1 Inhibitors. *Journal of chemical information and modeling* **2016**, *56*,
722 1856–1871.
- 723 [10] Lenseink, E. B. et al. Predicting Binding Affinities for GPCR Ligands Using Free-Energy Perturbation. *ACS*
724 *omega* **2016**, *1*, 293–304.
- 725 [11] Jorgensen, W. L. Computer-aided discovery of anti-HIV agents. *Bioorganic & medicinal chemistry* **2016**,
726 *24*, 4768–4778.
- 727 [12] Wang, L. et al. Accurate and Reliable Prediction of Relative Ligand Binding Potency in Prospective Drug
728 Discovery by Way of a Modern Free-Energy Calculation Protocol and Force Field. *J. Am. Chem. Soc.* **2015**,
729 *137*, 2695–2703.
- 730 [13] Abel, R.; Wang, L.; Harder, E. D.; Berne, B. J.; Friesner, R. A. Advancing Drug Discovery through Enhanced
731 Free Energy Calculations. *Accounts of chemical research* **2017**, *50*, 1625–1632.
- 732 [14] Zhang, J.; Yang, P. L.; Gray, N. S. Targeting cancer with small molecule kinase inhibitors. *Nat. Rev. Cancer*
733 **2009**, *9*, 28–39.
- 734 [15] Huggins, D. J.; Sherman, W.; Tidor, B. Rational approaches to improving selectivity in drug design. *J.*
735 *Med. Chem.* **2012**, *55*, 1424–1444.
- 736 [16] Fan, Q.-W.; Cheng, C. K.; Nicolaidis, T. P.; Hackett, C. S.; Knight, Z. A.; Shokat, K. M.; Weiss, W. A. A
737 dual phosphoinositide-3-kinase alpha/mTOR inhibitor cooperates with blockade of epidermal growth
738 factor receptor in PTEN-mutant glioma. *Cancer Res.* **2007**, *67*, 7960–7965.
- 739 [17] Apsel, B.; Blair, J. A.; Gonzalez, B.; Nazif, T. M.; Feldman, M. E.; Aizenstein, B.; Hoffman, R.; Williams, R. L.;
740 Shokat, K. M.; Knight, Z. A. Targeted polypharmacology: discovery of dual inhibitors of tyrosine and
741 phosphoinositide kinases. *Nat. Chem. Biol.* **2008**, *4*, 691–699.

- 742 [18] Knight, Z. A.; Lin, H.; Shokat, K. M. Targeting the Cancer Kinome through Polypharmacology. *Nat. Rev.*
743 *Cancer* **2010**, *10*, 130.
- 744 [19] Hopkins, A. L.; Mason, J. S.; Overington, J. P. Can we rationally design promiscuous drugs? *Curr. Opin.*
745 *Struct. Biol.* **2006**, *16*, 127–136.
- 746 [20] Hopkins, A. L. Network pharmacology: the next paradigm in drug discovery. *Nat. Chem. Biol.* **2008**, *4*,
747 682–690.
- 748 [21] Kijima, T. et al. Safe and successful treatment with erlotinib after gefitinib-induced hepatotoxicity: dif-
749 ference in metabolism as a possible mechanism. *J. Clin. Oncol.* **2011**, *29*, e588–90.
- 750 [22] Liu, S.; Kurzrock, R. Toxicity of targeted therapy: Implications for response and impact of genetic poly-
751 morphisms. *Cancer Treat. Rev.* **2014**, *40*, 883–891.
- 752 [23] Rudmann, D. G. On-target and off-target-based toxicologic effects. *Toxicol. Pathol.* **2013**, *41*, 310–314.
- 753 [24] Mendoza, M. C.; Er, E. E.; Blenis, J. The Ras-ERK and PI3K-mTOR pathways: cross-talk and compensation.
754 *Trends Biochem. Sci.* **2011**, *36*, 320–328.
- 755 [25] Tricker, E. M.; Xu, C.; Uddin, S.; Capelletti, M.; Ercan, D.; Ogino, A.; Pratilas, C. A.; Rosen, N.; Gray, N. S.;
756 Wong, K.-K.; Jänne, P. A. Combined EGFR/MEK Inhibition Prevents the Emergence of Resistance in EGFR-
757 Mutant Lung Cancer. *Cancer Discov.* **2015**, *5*, 960–971.
- 758 [26] Bailey, S. T.; Zhou, B.; Damrauer, J. S.; Krishnan, B.; Wilson, H. L.; Smith, A. M.; Li, M.; Yeh, J. J.; Kim, W. Y.
759 mTOR Inhibition Induces Compensatory, Therapeutically Targetable MEK Activation in Renal Cell Car-
760 cinoma. *PLoS One* **2014**, *9*, e104413.
- 761 [27] Chandarlapaty, S.; Sawai, A.; Scaltriti, M.; Rodrik-Outmezguine, V.; Grbovic-Huezo, O.; Serra, V.; Ma-
762 jumder, P. K.; Baselga, J.; Rosen, N. AKT Inhibition Relieves Feedback Suppression of Receptor Tyrosine
763 Kinase Expression and Activity. *Cancer Cell* **2011**, *19*, 58–71.
- 764 [28] Pao, W.; Miller, V.; Zakowski, M.; Doherty, J.; Politi, K.; Sarkaria, I.; Singh, B.; Heelan, R.; Rusch, V.; Ful-
765 ton, L.; Mardis, E.; Kupfer, D.; Wilson, R.; Kris, M.; Varmus, H. EGF receptor gene mutations are common
766 in lung cancers from “never smokers” and are associated with sensitivity of tumors to gefitinib and er-
767 lotinib. *Proceedings of the National Academy of Sciences* **2004**, *101*, 13306–13311.
- 768 [29] Kim, Y.; Li, Z.; Apetri, M.; Luo, B.; Settleman, J. E.; Anderson, K. S. Temporal resolution of autophospho-
769 rylation for normal and oncogenic forms of EGFR and differential effects of gefitinib. *Biochemistry* **2012**,
770 *51*, 5212–5222.
- 771 [30] Juchum, M.; Günther, M.; Laufer, S. A. Fighting Cancer Drug Resistance: Opportunities and Challenges
772 for Mutation-Specific EGFR Inhibitors. *Drug Resist. Updat.* **2015**, *20*, 12–28.
- 773 [31] Din, O. S.; Woll, P. J. Treatment of gastrointestinal stromal tumor: focus on imatinib mesylate. *Ther. Clin.*
774 *Risk Manag.* **2008**, *4*, 149–162.
- 775 [32] Lin, Y.-L.; Meng, Y.; Jiang, W.; Roux, B. Explaining why Gleevec is a specific and potent inhibitor of Abl
776 kinase. *Proc. Natl. Acad. Sci. U. S. A.* **2013**, *110*, 1664–1669.
- 777 [33] Lin, Y.-L.; Meng, Y.; Huang, L.; Roux, B. Computational Study of Gleevec and G6G Reveals Molecular
778 Determinants of Kinase Inhibitor Selectivity. *J. Am. Chem. Soc.* **2014**, *136*, 14753–14762.
- 779 [34] Lin, Y.-L.; Roux, B. Computational Analysis of the Binding Specificity of Gleevec to Abl, c-Kit, Lck, and
780 c-Src Tyrosine Kinases. *J. Am. Chem. Soc.* **2013**, *135*, 14741–14753.
- 781 [35] Aldeghi, M.; Heifetz, A.; Bodkin, M. J.; Knapp, S.; Biggin, P. C. Predictions of Ligand Selectivity from
782 Absolute Binding Free Energy Calculations. *J. Am. Chem. Soc.* **2017**, *139*, 946–957.

- 783 [36] Moraca, F.; Negri, A.; de Oliveira, C.; Abel, R. Application of Free Energy Perturbation (FEP+) to Under-
784 standing Ligand Selectivity: A Case Study to Assess Selectivity Between Pairs of Phosphodiesterases
785 (PDE's). *Journal of Chemical Information and Modeling* **2019**, *59*, 2729–2740, PMID: 31144815.
- 786 [37] Robert Roskoski Jr, USFDA Approved Protein Kinase Inhibitors. **2017**, Updated 3 May 2017.
- 787 [38] Santos, R.; Ursu, O.; Gaulton, A.; Bento, A. P.; Donadi, R. S.; Bologa, C. G.; Karlsson, A.; Al-Lazikani, B.;
788 Hersey, A.; Oprea, T. I.; Overington, J. P. A Comprehensive Map of Molecular Drug Targets. *Nat. Rev.*
789 *Drug Discov.* **2016**, *16*, 19–34.
- 790 [39] Volkamer, A.; Eid, S.; Turk, S.; Jaeger, S.; Rippmann, F.; Fulle, S. Pocketome of human kinases: prioritizing
791 the ATP binding sites of (yet) untapped protein kinases for drug discovery. *J. Chem. Inf. Model.* **2015**, *55*,
792 538–549.
- 793 [40] Manning, G.; Whyte, D. B.; Martinez, R.; Hunter, T.; Sudarsanam, S. The Protein Kinase Complement of
794 the Human Genome. *Science* **2002**, *298*, 1912–1934.
- 795 [41] Wu, P.; Nielsen, T. E.; Clausen, M. H. FDA-approved small-molecule kinase inhibitors. *Trends Pharmacol.*
796 *Sci.* **2015**, *36*, 422–439.
- 797 [42] Cowan-Jacob, S. W.; Fendrich, G.; Floersheimer, A.; Furet, P.; Liebetanz, J.; Rummel, G.; Rheinberger, P.;
798 Centeleghe, M.; Fabbro, D.; Manley, P. W.; IUCr, Structural biology contributions to the discovery of
799 drugs to treat chronic myelogenous leukaemia. *Acta Crystallogr. D Biol. Crystallogr.* **2007**, *63*, 80–93.
- 800 [43] Seeliger, M. A.; Nagar, B.; Frank, F.; Cao, X.; Henderson, M. N.; Kuriyan, J. c-Src Binds to the Cancer Drug
801 Imatinib with an Inactive Abl/c-Kit Conformation and a Distributed Thermodynamic Penalty. *Structure*
802 **2007**, *15*, 299–311.
- 803 [44] Huse, M.; Kuriyan, J. The conformational plasticity of protein kinases. *Cell* **2002**, *109*, 275–282.
- 804 [45] Harrison, S. C. Variation on an Src-like theme. *Cell* **2003**, *112*, 737–740.
- 805 [46] Volkamer, A.; Eid, S.; Turk, S.; Rippmann, F.; Fulle, S. Identification and Visualization of Kinase-Specific
806 Subpockets. *J. Chem. Inf. Model.* **2016**, *56*, 335–346.
- 807 [47] Christmann-Franck, S.; van Westen, G. J. P.; Papadatos, G.; Beltran Escudie, F.; Roberts, A.; Overing-
808 ton, J. P.; Domine, D. Unprecedentedly Large-Scale Kinase Inhibitor Set Enabling the Accurate Prediction
809 of Compound–Kinase Activities: A Way toward Selective Promiscuity by Design? *Journal of chemical*
810 *information and modeling* **2016**, *56*, 1654–1675.
- 811 [48] Anastassiadis, T.; Deacon, S. W.; Devarajan, K.; Ma, H.; Peterson, J. R. Comprehensive assay of kinase
812 catalytic activity reveals features of kinase inhibitor selectivity. *Nat. Biotechnol.* **2011**, *29*, 1039–1045.
- 813 [49] Davis, M. I.; Hunt, J. P.; Herrgard, S.; Ciceri, P.; Wodicka, L. M.; Pallares, G.; Hocker, M.; Treiber, D. K.;
814 Zarrinkar, P. P. Comprehensive Analysis of Kinase Inhibitor Selectivity. *Nat. Biotechnol.* **2011**, *29*, 1046–
815 1051.
- 816 [50] Klaeger, S. et al. The target landscape of clinical kinase drugs. *Science* **2017**, 358.
- 817 [51] Sun, C. et al. Intrinsic resistance to MEK inhibition in KRAS mutant lung and colon cancer through
818 transcriptional induction of ERBB3. *CellReports* **2014**, *7*, 86–93.
- 819 [52] Manchado, E.; Weissmueller, S.; Morris, J. P.; Chen, C.-C.; Wullenkord, R.; Lujambio, A.; de Stanchina, E.;
820 Poirier, J. T.; Gainor, J. F.; Corcoran, R. B.; Engelman, J. A.; Rudin, C. M.; Rosen, N.; Lowe, S. W. A combi-
821 natorial strategy for treating KRAS-mutant lung cancer. *Nature* **2016**, *534*, 647–651.

- 822 [53] Shao, H.; Shi, S.; Huang, S.; Hole, A. J.; Abbas, A. Y.; Baumli, S.; Liu, X.; Lam, F.; Foley, D. W.; Fischer, P. M.;
823 Noble, M.; Endicott, J. A.; Pepper, C.; Wang, S. Substituted 4-(Thiazol-5-yl)-2-(phenylamino)pyrimidines
824 Are Highly Active CDK9 Inhibitors: Synthesis, X-ray Crystal Structures, Structure–Activity Relationship,
825 and Anticancer Activities. *J. Med. Chem.* **2013**, *56*, 640–659.
- 826 [54] Blake, J. F. et al. Discovery of (S)-1-(1-(4-Chloro-3-fluorophenyl)-2-hydroxyethyl)-4-(2-((1-methyl-1H-
827 pyrazol-5-yl)amino)pyrimidin-4-yl)pyridin-2(1H)-one (GDC-0994), an Extracellular Signal-Regulated Ki-
828 nase 1/2 (ERK1/2) Inhibitor in Early Clinical Development. *J. Med. Chem.* **2016**, *59*, 5650–5660.
- 829 [55] Bosc, N.; Meyer, C.; Bonnet, P. The use of novel selectivity metrics in kinase research. *BMC bioinformatics*
830 **2017**, *18*, 17.
- 831 [56] Cheng, A. C.; Eksterowicz, J.; Geuns-Meyer, S.; Sun, Y. Analysis of kinase inhibitor selectivity using a
832 thermodynamics-based partition index. *J. Med. Chem.* **2010**, *53*, 4502–4510.
- 833 [57] Shirts, M. R.; Mobley, D. L.; Brown, S. P. *Structure Based Drug Design*; Cambridge University Press, 2009.
- 834 [58] Kooistra, A. J.; Kanev, G. K.; van Linden, O. P. J.; Leurs, R.; de Esch, I. J. P.; de Graaf, C. KLIFS: a structural
835 kinase-ligand interaction database. *Nucleic acids research* **2016**, *44*, D365–D371.
- 836 [59] van Linden, O. P. J.; Kooistra, A. J.; Leurs, R.; de Esch, I. J. P.; de Graaf, C. KLIFS: a knowledge-based
837 structural database to navigate kinase-ligand interaction space. *Journal of medicinal chemistry* **2014**, *57*,
838 249–277.
- 839 [60] Hole, A. J.; Baumli, S.; Shao, H.; Shi, S.; Huang, S.; Pepper, C.; Fischer, P. M.; Wang, S.; Endicott, J. A.; No-
840 ble, M. E. Comparative Structural and Functional Studies of 4-(Thiazol-5-yl)-2-(phenylamino)pyrimidine-
841 5-carbonitrile CDK9 Inhibitors Suggest the Basis for Isotype Selectivity. *J. Med. Chem.* **2013**, *56*, 660–670.
- 842 [61] Yung-Chi, C.; Prusoff, W. H. Relationship between the inhibition constant (KI) and the concentration of
843 inhibitor which causes 50 per cent inhibition (I50) of an enzymatic reaction. *Biochemical Pharmacology*
844 **1973**, *22*, 3099 – 3108.
- 845 [62] Hauser, K.; Negron, C.; Albanese, S. K.; Ray, S.; Steinbrecher, T.; Abel, R.; Chodera, J. D.; Wang, L. Pre-
846 dicting resistance of clinical Abl mutations to targeted kinase inhibitors using alchemical free-energy
847 calculations. *Communications Biology* **2018**, *1*, 70.
- 848 [63] Kalliokoski, T.; Kramer, C.; Vulpetti, A.; Gedeck, P. Comparability of Mixed IC50 Data—a Statistical Analy-
849 sis. *PloS One* **2013**, *8*, e61007.
- 850 [64] Hari, S. B.; Merritt, E. A.; Maly, D. J. Sequence determinants of a specific inactive protein kinase confor-
851 mation. *Chemistry & biology* **2013**, *20*, 806–815.
- 852 [65] Davis, M. I.; Hunt, J. P.; Herrgard, S.; Ciceri, P.; Wodicka, L. M.; Pallares, G.; Hocker, M.; Treiber, D. K.;
853 Zarrinkar, P. P. Comprehensive analysis of kinase inhibitor selectivity. *Nat. Biotechnol.* **2011**, *29*, 1046–
854 1051.
- 855 [66] Graczyk, P. P. Gini coefficient: a new way to express selectivity of kinase inhibitors against a family of
856 kinases. *Journal of medicinal chemistry* **2007**, *50*, 5773–5779.
- 857 [67] Duong-Ly, K. C.; Devarajan, K.; Liang, S.; Horiuchi, K. Y.; Wang, Y.; Ma, H.; Peterson, J. R. Kinase Inhibitor
858 Profiling Reveals Unexpected Opportunities to Inhibit Disease-Associated Mutant Kinases. *CellReports*
859 **2016**, *14*, 772–781.
- 860 [68] Uitdehaag, J. C. M.; Zaman, G. J. R. A theoretical entropy score as a single value to express inhibitor
861 selectivity. *BMC bioinformatics* **2011**, *12*, 94.
- 862 [69] Michel, J.; Verdonk, M. L.; Essex, J. W. Protein-ligand binding affinity predictions by implicit solvent
863 simulations: a tool for lead optimization? *Journal of medicinal chemistry* **2006**, *49*, 7427–7439.

- 864 [70] Berman, H. M. et al. The Protein Data Bank. *Acta Crystallogr. D Biol. Crystallogr.* **2002**, 58, 899–907.
- 865 [71] Sastry, G. M.; Adzhigirey, M.; Day, T.; Annabhimoju, R.; Sherman, W. Protein and ligand preparation:
866 parameters, protocols, and influence on virtual screening enrichments. *J. Comput. Aided Mol. Des.* **2013**,
867 27, 221–234.
- 868 [72] Salvatier, J.; Wiecki, T. V.; Fonnesbeck, C. Probabilistic programming in Python using PyMC3. *PeerJ Com-*
869 *puter Science* **2016**, 2, e55.
- 870 [73] Al-Rfou, R. et al. Theano: A Python framework for fast computation of mathematical expressions. *arXiv*
871 *e-prints* **2016**, *abs/1605.02688*.