

1 **Bimodal evolution of Src and Abl kinase substrate specificity revealed using**  
2 **mammalian cell extract as substrate pool**

3  
4  
5 **Authors**

6 Patrick Finneran<sup>1</sup>, Margaret Soucheray<sup>2</sup>, Christopher Wilson<sup>1,\*</sup>, Renee Otten<sup>1</sup>, Vanessa  
7 Buosi<sup>1\*</sup>, Nevan J. Krogan<sup>2</sup>, Danielle L. Swaney<sup>2</sup>, Douglas L. Theobald<sup>3</sup>, Dorothee Kern<sup>1#</sup>

8  
9 *<sup>1</sup> Department of Biochemistry and Howard Hughes Medical Institute, Brandeis University,*  
10 *Waltham, Massachusetts 02453, United States.<sup>2</sup> Department of Cellular and Molecular*  
11 *Pharmacology, University of California, San Francisco, San Francisco, CA 94158, USA;*  
12 *Quantitative Biosciences Institute (QBI), University of California, San Francisco, San*  
13 *Francisco, CA 94158, USA; Gladstone Institute of Data Science and Biotechnology, San*  
14 *Francisco, CA, USA, San Francisco, CA 94158, USA, <sup>3</sup>Department of Biochemistry,*  
15 *Brandeis University, Waltham, Massachusetts 02453, United States.*

16  
17  
18 *\* Current addresses: C.W. Merkin Institute of Transformative Technologies in Healthcare,*  
19 *Broad Institute of Harvard and MIT, Department of Chemistry and Chemical Biology,*  
20 *Harvard University Cambridge, MA, USA*  
21 *V.B. Sanofi Pasteur, 1541 Avenue Marcel Mérieux, 69280 MARCY L'ETOILE, France*

22  
23  
24 # corresponding author

25  
26 **Abstract**

27 The specificity of phosphorylation by protein kinases is essential to the integrity of  
28 biological signal transduction. While peptide sequence specificity for individual kinases  
29 has been examined previously, here we explore the evolutionary progression that has led  
30 to the modern substrate specificity of two non-receptor tyrosine kinases, Abl and Src. To  
31 efficiently determine the substrate specificity of modern and reconstructed ancestral

32 kinases, we developed a method using mammalian cell lysate as the substrate pool,  
33 thereby representing the naturally occurring substrate proteins. We find that the oldest  
34 tyrosine kinase ancestor was a promiscuous enzyme that evolved through a more specific  
35 last common ancestor into a specific human Abl. In contrast, the parallel pathway to  
36 human Src involved a loss of substrate specificity, leading to general promiscuity. These  
37 results add a new facet to our understanding of the evolution of signaling pathways, with  
38 both subfunctionalization and neofunctionalization along the evolutionary trajectories.

39

## 40 **Introduction**

41 The human genome contains 32 non-receptor tyrosine kinases (NRTKs) that are  
42 tightly involved in a multitude of cellular processes including differentiation, apoptosis,  
43 and proliferation<sup>1-3</sup>. The interaction between each NTRK and its substrate comprises a  
44 fundamental cellular signal that consequently required the evolution of specificity between  
45 signaling pathways. To prevent unwanted signaling 'crosstalk', NRTKs have evolved two  
46 main strategies to ensure substrate insulation<sup>4,5</sup>: kinase localization<sup>6-12</sup> and active-site  
47 peptide specificity<sup>13-18</sup>. The localization process is achieved through binding interactions  
48 on the NTRKs' SH2/SH3 domains, which complex with either phosphotyrosines or poly-  
49 prolines, respectively<sup>19-22</sup>. The differences in the active site of NRTK kinase domains  
50 results in specificity where only a subset of substrates can bind and thus get  
51 phosphorylated.

52 Unlike the serine/threonine family of kinases, NRTKs possess relatively  
53 promiscuous active-site peptide specificities with a broad range of potential substrates<sup>5</sup>.  
54 In high-throughput substrate screens, catalytic domains of NRTK members  
55 phosphorylated hundreds of distinct peptide sequences, highlighting the promiscuity of  
56 these kinases. Nevertheless, comparisons within the family show unique sequence  
57 preferences and a consequent range of substrate selectivity. For two of the most well-  
58 studied members, Abl and Src, narrow and broad selectivities are reported,  
59 respectively<sup>13,17</sup>. Because substrates bind at the active site in an elongated fashion, the  
60 primary peptide sequence largely dictates the description of selectivity<sup>5</sup>. Abl has a clear  
61 preference for hydrophobic residues flanking the tyrosine of interest (I/L/V<sub>-1</sub>, A<sub>+1</sub>, and  
62 P<sub>+3</sub>)<sup>14,23</sup>. In contrast, little sequence selectivity is observed for Src other than the relatively

63 weaker preferences for a bulky aliphatic residue (I/V/L-1) at the residue preceding the  
64 phosphoacceptor, a phenylalanine three residues away from the phosphoacceptor ( $F_{+3}$ ),  
65 and a negatively charged residue on the N-terminal side of the phosphoacceptor (D/E-4, -  
66 3, -2)<sup>17,24,25</sup>.

67 As Src and Abl are sister clades within the NRTK phylogeny, their distinct  
68 sequence preferences beget the question: how did peptide selectivity arise throughout  
69 evolution? Herein we answer this question using ancestral sequence reconstruction  
70 (ASR) for the catalytic domains. ASR uses modern sequences and an evolutionary model  
71 to infer the sequences of internal nodes in a phylogenetic tree (Figure 1A, Figure 1 — figure  
72 supplement 1)<sup>26-31</sup>. Resurrection of ancestral kinases bridging Src and Abl allows the  
73 evolution of peptide selectivity to be directly traced. The differences in protein sequence  
74 between modern kinases and ancestors are spread through the kinase domain and the  
75 oldest ancestor only has ~65% similarity with Src (Figure 1B).

76 Src and Abl display lower sequence selectivity than other kinases such as Aurora  
77 A and B-RAF serine/threonine kinases<sup>5,24,32</sup>, and, consequently, obtaining an accurate  
78 description of the primary sequence determinants for each tyrosine kinase is a greater  
79 statistical challenge<sup>33</sup>. Comparison of ancestral and modern kinases requires a  
80 comprehensive library of substrates since ancestral kinases likely refined their sequence  
81 preferences over time. To ensure biological relevance, the peptide library should ideally  
82 be composed of naturally occurring proteins. To construct such a library, we took  
83 advantage of the diversity of sequences present in mammalian whole cell lysate (HEK293  
84 cell line)<sup>34,35</sup>. After endogenous kinases are covalently inhibited, the proteome of  
85 mammalian cells presents a convenient substrate library containing thousands of  
86 potential protein substrates. Here we use this comprehensive library to examine the  
87 evolution of substrate selectivity in Abl/Src tyrosine kinases. We find that kinase substrate  
88 preferences evolved in a complex manner involving two different modes: a promiscuous  
89 progenitor specialized into the modern specific Abl, whereas evolution of Src involved  
90 relaxing selectivity via a specific ancestral intermediate. We find that kinase substrate  
91 preferences evolved in a complex manner involving two different modes: a promiscuous  
92 progenitor specialized into the modern specific Abl (subfunctionalization), whereas  
93 evolution of Src involved relaxing selectivity via a specific ancestral intermediate

94 (neofunctionalization). Therefore, our results shed light into a critical open question in  
95 signaling, how new protein kinases with novel substrate specificities have evolved.

96

## 97 **Results**

### 98 **Whole-cell lysate phosphoproteomics-based approach**

99 A large and diverse protein library is necessary to readily screen the primary  
100 sequence determinants of ancestral and modern NRTKs. An easily accessible, cheap,  
101 and biologically relevant pool of substrates was created by inactivating endogenous  
102 kinases in HEK293 lysate using the covalent, nonspecific inhibitor 5'-[p-  
103 (fluorosulfonyl)benzoyl]adenosine (FSBA) (Figure 1C). After dialyzing out unreacted  
104 inhibitor, purified kinases were added to the treated lysate and phosphorylation was  
105 initiated with the addition of Mg<sup>2+</sup> and ATP. To assess the required reaction time, total  
106 phosphorylation was monitored by western dot blot experiments. Constant  
107 phosphorylation levels were found to occur between two to four hours (Figure 1D,E). After  
108 protein digestion with trypsin, peptide fragments were enriched for phosphorylation using  
109 immobilized metal affinity chromatography (IMAC) and then analyzed by liquid  
110 chromatography-mass spectrometry (LC/MS/MS, Figure 1C). Peptides were associated  
111 with their full protein sequence based on the known HEK293 proteome, and results were  
112 focused on a 15-amino acid sequence window centered on the phosphorylated tyrosine.

113 To determine a kinase's sequence specificity, each amino acid frequency must be  
114 compared between the phosphorylated dataset and the background HEK293 proteome<sup>33</sup>.  
115 A background dataset was generated by proteomic analysis of lysate that was treated  
116 with kinase inhibitor, but where no kinase was added and no phosphorylation enrichment  
117 was performed<sup>36</sup>. From this analysis, the amino acid frequencies of all surrounding  
118 tyrosine residues (not only phosphorylated tyrosines) were calculated (Figure 1 — figure  
119 supplements 2 and 3).

120 To determine the extent of endogenously phosphorylated peptides present, we  
121 analyzed a lysate sample that was enriched for phosphorylated peptides but lacked  
122 exogenous kinase. A total of 26 phosphorylated tyrosines were found in this control and  
123 these peptides were excluded from the ancestral/modern NRTK list of substrates in which  
124 kinase was added (Figure 1 — figure supplement 4).

125           In the data set where Src was added to the cell lysate, 8208 unique phosphorylated  
126 sequences were identified (Figure 2A). These peptides include the characteristic  
127 preferences for large aliphatic residues directly preceding the phosphotyrosine (V<sub>-1</sub>),  
128 negatively charged residues in multiple positions N-terminal to the phosphotyrosine (D/E-  
129 <sub>3, -2</sub>), and a glycine following the phosphorylation site (G<sub>+1</sub>). An inclination for other  
130 aliphatic residues preceding the phosphotyrosine (I/P/T<sub>-1</sub>) is also seen (Figure 2B).  
131 Notably, a preference for proline at the -1 position identified in our data was not observed  
132 previously. For Abl, specificity for large aliphatic residues preceding the phosphotyrosine  
133 is found (I/V<sub>-1</sub>), as well as the canonical proline at the +3 position (Figure 2B). Additionally,  
134 Abl exhibits a high preference for proline at the -2 position, which had not been identified  
135 previously.

136           We can compare the results obtained here with the substrate specificity that is  
137 observed when only natural substrates are considered. PhosphoSitePlus is a database  
138 which annotates all known phosphorylation sites *in vivo* and *in vitro* that a given kinase  
139 phosphorylates<sup>24</sup>. Overall, our results confirm the previously found descriptions for both  
140 Src and Abl's substrate specificities based on substrates in the PhosphoSitePlus  
141 database, but additionally identify a few new preferences (Figure 2B, C). Our HEK293  
142 lysate has a much larger number of phosphorylated substrates than the PhosphoSitePlus  
143 database, which allows us to ascertain the residues dictating phosphorylation specificity  
144 with greater accuracy and statistical significance than is possible with the  
145 PhosphoSitePlus data (Figure 2B,C). Some of the differences may be due to the larger  
146 number of substrates in our whole cell lysate. However, many of the observed  
147 discrepancies likely result from differences in experimental design. PhosphoSitePlus is  
148 based on *in vivo* substrates for the full-length kinase, whereas we are interested in the  
149 intrinsic specificity of the kinase domain. In our experiments, we do not have the full-  
150 length kinases and, therefore, we only find substrates that are selected by the kinase  
151 domain itself. In contrast, phosphorylation within the cellular framework, as reported by  
152 PhosphoSitePlus, is strongly determined by regulation and co-localization events, and  
153 intrinsic kinase domain specificity plays a relatively smaller role. Indeed, Shah *et al.*  
154 studied the specificity of NRTK kinase domains with a high-throughput, cell-surface based

155 experiment and found similar discrepancies between PhosphoSitePlus based logos and  
156 their experimentally determined sequence determinants<sup>17</sup>.

### 157 **Evolution of specificity between Src and Abl**

158 Having established now the accuracy and statistics of our methodology on the  
159 modern kinases, we next chose to determine the sequence specificity of three resurrected  
160 ancestral kinases (Figure 1A). Anc-AS and Anc-S1 were previously resurrected for  
161 investigating the mechanism of Abl selectivity for Gleevec<sup>29</sup>, while the newly resurrected  
162 Anc-AST (Figure 1 — figure supplement 1) is the oldest common ancestor of the Abl/Src  
163 branch and the Tec family. Using our whole cell lysate phosphoproteomics-based  
164 method, we identified a total of 12,056 unique sequences phosphorylated by the ancestral  
165 and modern proteins (Figure 3A, Figure 3 – figure supplement 2). The common ancestor  
166 of Src and Abl, Anc-AS, phosphorylated the least number of substrates (2495), which was  
167 comparable to Abl (3073). The relative dearth of substrates for Anc-AS hinted that this  
168 ancestor might be more specific than the promiscuous Src, which phosphorylated a total  
169 of 8208 substrates. In contrast, the ancestors preceding (Anc-AST) and following (Anc-  
170 S1) the common ancestor of Src and Abl each phosphorylated a significantly greater  
171 number of substrates (8189 and 7242, respectively), indicating these ancestors were  
172 likely more promiscuous.

173 General kinase specificity was then quantified by calculating the substrate  
174 sequence entropy for each position in the 15-residue window (Figure 3B)<sup>37</sup>. Lower  
175 entropy, with fewer potential amino acid possibilities, indicates higher specificity. As  
176 expected from the sequence logos (Figure 2B), Abl possesses the lowest entropy at  
177 residues close to the phosphorylated tyrosine (-3, -2, -1, +1, and +3), with the most  
178 specificity occurring at the signature proline position (+3). In contrast, the promiscuity of  
179 Src manifests itself as increased entropy across almost all positions and a more limited  
180 region of high specificity (-2, -1, and +1). The ancestral proteins give entropy plots that  
181 agree with what was suggested by the observed substrate counts: The common ancestor  
182 (Anc-AS) possessed entropy akin to Abl, albeit with a higher entropy at +3 and lower  
183 entropy at positions -2 and +5. The two additional ancestors, Anc-AST and Anc-S1, both  
184 exhibited 'hybrid' specificity with higher entropy than Abl, but less than Src. Notably, only  
185 Src lacks specificity at the +3 position.

186 To analyze each enzyme's positional specificity in more detail, specificity heat  
187 maps were created to illustrate the relative specificity for each amino acid at every position  
188 in the 15-residue window (shown as a 20x15 matrix of positional normalized amino acid  
189 log probabilities, Figure 3C)<sup>17,38</sup>. Qualitatively, Abl's specificity is apparent from the high-  
190 intensity signals for both preferred residues (red, P<sub>+3</sub> and I<sub>-1</sub>) and unfavorable residues  
191 (blue, S<sub>-1</sub>, P<sub>+1</sub>, and D/E<sub>+3</sub>), while Src has more white space and overall less intense signal.  
192 We note that under substrate saturating conditions the enzyme would phosphorylate even  
193 the less favorable substrates, which could result in an apparent low specificity. To ensure  
194 that the observed promiscuity is not due to substrate saturation, experiments with Src  
195 were repeated with a much shorter incubation period of the cell extract and the kinase  
196 (10 minutes versus 4 hours). In this control, less phosphorylation was observed (Figure  
197 1 C,D); however, the same primary sequence determinants were found (Figure 3 — figure  
198 supplements 1 and 2), validating our findings.

199 Tracing individual amino acid preferences at specific positions provides a clear  
200 picture of how specificity evolved for Abl (Figure 4A). Focusing on residues which are  
201 preferred in Abl, but disfavored in Src (P<sub>+3</sub>, A/V<sub>+2</sub>, A<sub>+1</sub>, and P<sub>-2</sub>), we see that moderate  
202 preference is already observed in ancestral kinases. The evolutionary path from the  
203 oldest ancestor (Anc-AST) to Abl involves further increasing specificity either at the Anc-  
204 AS node (P<sub>-2</sub> and A/V<sub>-1</sub>) or in the final transition to Abl (A<sub>+1</sub> and P<sub>+3</sub>). In contrast, the  
205 pathway from the more specific Anc-AS to Src involves a corresponding loss of specificity  
206 for each of these residues, with Anc-S1 possessing intermediary preferences (Figure 4A).

207 The evolution of the few positions favored by Src followed a different trend. The  
208 preference for P<sub>-1</sub> appears late, only in Src. Such recent evolution is in agreement with  
209 the lack of preference for proline at position -1 of its close homolog Lck<sup>17</sup>. Other Src  
210 sequence preferences were already present in the oldest ancestor, then lost in Anc-AS  
211 and regained in Anc-S1 and Src (Figure 4B).

212 The increased promiscuity of Src is revealed in overall lower log probability values  
213 than that of Abl's specific residues. Despite these differences in substrate specificity, there  
214 are multiple positions where all modern and ancestral kinases prefer the same residues,  
215 most of which are well-known features of NRTKs (e.g., I<sub>-1</sub>, D/E<sub>-3</sub>, and S<sub>+2</sub>) (Figure 4C). As  
216 these characteristics are observed in all ancestral and modern proteins in our study and

217 are common among most NRTKs, they likely represent the oldest features of substrate  
218 specificity for the NRTK family.

## 219 **Validation of evolutionary trends of primary sequence determinants via enzyme** 220 **kinetics**

221 Having determined the primary sequence determinants for the ancestral and  
222 modern kinases, *in vitro* peptide enzyme turnover experiments were performed to relate  
223 these bulk specificity experiments to quantitative enzymatic parameters. Compressing  
224 thousands of substrates into residue-by-residue descriptions is compelling (i.e.,  
225 preference for P<sub>+3</sub>), but how these preferences relate to enzymatic properties remains  
226 unclear. We therefore measured the Michaelis-Menten kinetics of four distinct peptide  
227 substrates with each of the five ancestral and modern kinases.

228 Previous microarray specificity experiments had determined optimized substrates  
229 for Src and Abl, known as Srctide and Abltide, respectively<sup>13,18,25</sup>. While both substrates  
230 are ideal for their respective modern kinases, each also has residues favored by the  
231 opposite kinase, which allow it to be phosphorylated to a certain degree by each kinase.  
232 Therefore, we designed modified versions of Srctide and Abltide, called Srctide2 and  
233 Abltide2, to test our evolutionary trends (Figure 5A,B). Srctide2 was intended to be  
234 favored by both Src and Anc-S1, with changes made to Srctide to include residues that  
235 occurred more frequently in substrates for these two kinases (D<sub>-2</sub>, A<sub>+2</sub>, and I<sub>+3</sub>). Abltide2  
236 was designed to be preferred by both Abl and Anc-AS by mutating the alanine at position  
237 -2 into a proline. P<sub>-2</sub> is favored by all the kinases, except for Src (Figure 5A, B).

238 As substrates in signaling cascades are generally present at low concentrations *in*  
239 *vivo*, the  $k_{cat}/K_M$  likely represents a more fundamentally important parameter than  $k_{cat}$  for  
240 substrate specificity. As can be seen from the measured Michaelis-Menten curves (Figure  
241 5C), the measured differences in the kinetics corroborate the evolutionary trends found  
242 before but suggest additional features for substrate specificities. Starting with a  
243 promiscuous Anc-AST, Anc-AS becomes more selective, particularly for substrates with  
244 P<sub>-2</sub>, Abltide2 (which was identified as a highly preferred residue for Anc-AS and Abl from  
245 our phosphoproteomics data, Figure 5B). Moving to Abl, the specificity for Abltide and  
246 Abltide2 further increases as seen with high increases in  $k_{cat}/K_M$ . This is primarily due to  
247 the strong preference for P<sub>+3</sub> (Figure 5B), which are both present in Abltide and Abltide2.



248 At higher concentrations of these two well-optimized peptides we observe partial  
249 inhibition, due to the negative cooperativity of ATP and peptide substrate found for Abl<sup>39</sup>.

250 Following the evolutionary branch towards Src, Anc-S1 becomes more  
251 promiscuous, mainly due to its ability to catalyze Src-preferred substrates in addition to  
252 the Abltide substrates. Furthermore, the strong preference of I<sub>+3</sub> observed in the  
253 proteomics data (Figure 5B) can be directly recapitulated by the preference for Srctide2  
254 (Figure 5C). The strong preference of Srctide and Srctide2 over the Abltide substrates  
255 only appears in Src, primarily due to a complete loss of preference for P<sub>+3</sub> leading to poor  
256 activity for the Abltide substrates, combined with a subtle preference for F<sub>+3</sub>.

## 257 **Discussion**

258 There are multiple methods described in the literature for assessing substrate  
259 specificity of protein kinases, including non-receptor tyrosine kinases and  
260 serine/threonine kinases. Several studies look only at known natural substrates, such as  
261 in PhosphoSitePlus, to determine which residues occur more frequently in known  
262 phosphorylation sites<sup>24,40</sup>. Other studies, including the one presented here, use the whole  
263 cell lysate as a pool of substrates to test the specificity of kinases<sup>34,35</sup>. The HAKA-MS  
264 method reported by Muller *et al.* used a similar whole cell lysate, yet only found a P<sub>+3</sub>  
265 preference from 104 Abl substrates (Figure 2 – figure supplement 2)<sup>35</sup>. In contrast, our  
266 study found ~30-fold more substrates and several residues that are preferred or  
267 disfavored at multiple positions. The large difference in detected phosphorylated  
268 substrates could be due to different methods to enrich for phosphorylated peptides: Muller  
269 *et al.* use multiple phosphotyrosine binding antibodies, whereas we use IMAC. We also  
270 tried using phosphotyrosine binding antibodies for enrichment, but we found peptide bias  
271 and artifacts using this method. Lastly, one of the most popular methods uses peptide  
272 libraries to determine the preference of kinases<sup>13,14,16-18,38</sup>.

273 Our experimental method, which exploits the substrate-rich, kinase-inactive whole  
274 cell lysate and has improved statistical significance, discovered new determinants of Abl  
275 and Src specificity in addition to those previously reported<sup>17</sup>: Src shows significant  
276 favorability for proline or threonine at the -1 position and a serine at the -2 location, while  
277 Abl shows a preference for proline at the -2 position and a serine at the -2 position. These  
278 new features were subsequently verified by *in vitro* enzyme kinetic experiments. The most

279 important advantage of the increased sensitivity of our assay has been the ability to  
280 explore how Src and Abl kinases evolved their sequence preferences.

281 How these kinases have differing specificities can be partially rationalized based on  
282 the details of kinase sequence and structure. Y569 in Abl has previously been determined  
283 to be required for its preference for proline in the +3 position of substrates<sup>41</sup>. Leucine at  
284 the homologous position (L475) in Src was previously shown to disfavor proline<sup>41</sup>. A437  
285 has been proposed to be responsible for Src's preference for phenylalanine at position  
286 +3, along with L475<sup>17</sup>. All ancestors contain an isoleucine at position 475 and show  
287 intermediate specificity towards P<sub>+3</sub>. Moreover, Anc-S1 differs from Src in both positions,  
288 I475 and L437 (Figure 5 – figure supplement 1), identical to another Src family kinase  
289 member, Lck (I412 and I450 in Lck). Interestingly, when comparing the Anc-S1 substrate  
290 preference to that of Lck, as investigated by Shah *et al.*<sup>17</sup>, we see a high similarity in the  
291 +3 position preferences. Both Anc-S1 and Lck show a strong preference for L<sub>+3</sub> and P<sub>+3</sub>,  
292 suggesting that Anc-S1 is more Lck-like, and that the substitution to the less bulky A437  
293 in Src causes its preference for F<sub>+3</sub> (Figure 5 – figure supplement 1). These structural  
294 differences between Src and the ancestors explain why the Src tide is less effective at  
295 being phosphorylated by other kinases. Anc-S1 was unable to effectively phosphorylate  
296 Src tide potentially due to F<sub>+3</sub>, wherein its substitution to the less bulky isoleucine in  
297 Src tide2 is favored by Anc-S1. We elect to steer clear from additional structural  
298 explanations for other detected specificities, as these are just coarse models of  
299 kinase/substrate complexes, and more collective, long-range effects often underlie such  
300 specificity changes. For example, in an appealing study of the evolution of CMGC  
301 kinases, Howard *et al.* identified a key residue for imparting specificity at the +1 position<sup>28</sup>.  
302 Tests of their hypothesis via mutations in the corresponding modern kinases resulted in  
303 partial changes in specificities. Since the authors were unable to achieve a full swap in  
304 specificity, they concluded that there must be additional residues in play that are not  
305 readily apparent by looking at the differences in active-site residues.

306 The different trajectories we find in the evolution of Src and Abl substrate specificity  
307 add a new facet to our understanding of the evolution of signaling pathways. It has been  
308 postulated that much of biological diversity, including metazoan complexity, has been  
309 driven by the evolution of new regulatory networks and signaling pathways, such as those

310 controlled by the post-translational modification of kinase phosphorylation<sup>42,43</sup>. A key  
311 evolutionary challenge in creating a new signaling pathway is ensuring kinase specificity  
312 to minimize crosstalk with other pathways, many of which are vital to cellular fitness. One  
313 critical open question is how new protein kinases with novel substrate specificities have  
314 evolved.

315 Gene duplication has been the major force driving the evolutionary diversity of  
316 signaling pathways and kinase specificity<sup>42-44</sup>. There are two main ways that gene  
317 duplication can evolve enzymes with new and different functions: (1)  
318 “subfunctionalization”, the specialization of previously existing functions, and (2)  
319 “neofunctionalization”, the creation of a novel function through the accumulation of  
320 beneficial, gain-of-function mutations<sup>45,46</sup>. The evolution of sequence specificity does not  
321 cleanly fall in either of these categories, because specificity is not an “all-or-nothing” gain  
322 or loss of a function. Nevertheless, for kinases, subfunctionalization most closely aligns  
323 with evolution from a non-specific ancestor (which can bind and use many different  
324 substrates) to descendants with differing, specialized specificities (which can bind and  
325 use only a subset of the ancestral substrates), whereas neofunctionalization involves  
326 evolution from a specific ancestor to a descendant with a broad, promiscuous specificity  
327 or with a different specificity. Though neofunctionalization is the older of the two  
328 hypotheses<sup>47,48</sup>, it is now widely viewed as relatively improbable and hence less frequent  
329 in evolution than subfunctionalization mechanisms<sup>49-53</sup>. For instance, in the evolution of  
330 specificity in CMGC protein kinases, the ancestral kinase was a promiscuous bispecific  
331 enzyme in respect to the +1 position, unlike the modern kinases which are specific for a  
332 single amino acid<sup>28</sup>. Currently, the evolutionary mechanisms by which gene duplications  
333 evolve new functions are controversial, and there are relatively few examples of classic  
334 neofunctionalization<sup>54-56</sup>.

335 Intriguingly, we see both mechanisms of gene duplication in the evolution of Abl  
336 and Src. With Abl, subfunctionalization converted a promiscuous, non-specific ancestor  
337 (ANC-AST) into specific descendants (ANC-AS and modern Abl); with Src,  
338 neofunctionalization transformed a surprisingly specific ancestor (ANC-AS, the last  
339 common ancestor of the Abl and Src kinase families) into progressively more

340 promiscuous descendants (ANC-S1 followed by modern Src). In fact, the particular  
341 lineage leading from ANC-AST to modern Src appears to involve both mechanisms.

342 In the toxin-antitoxin signaling systems of bacteria, Aakre *et al.*<sup>57</sup> found that the  
343 evolution of some enzymes passes through promiscuous intermediates before  
344 developing strong substrate specificity. Our results suggest that eukaryotic protein  
345 kinases similarly evolve through waves of promiscuous and specific effectors. Periods of  
346 increased promiscuity may allow kinases to access new substrates while maintaining their  
347 current function, which could explain the redundancy of kinase networks<sup>24</sup>. To insulate  
348 individual pathways from others, sometimes this promiscuity may be selected out. In other  
349 cases, the overlap may provide a fitness advantage resulting in modern protein kinases  
350 with overlapping substrates.

351 The present work has addressed one component of kinase substrate specificity –  
352 the intrinsic specificity of the kinase domain – but it is important to emphasize the rich  
353 literature about the crucial role of the additional regulatory domains found within most  
354 NRTKs for kinase specificity. Proteins containing SH2 or SH3 domain binding sites are  
355 more likely to be phosphorylation targets for NRTKs, due to the selective activation of the  
356 kinases<sup>7,20,58-60</sup>. SH2 and SH3 domains have their own unique specificity<sup>61-65</sup>, and bring  
357 the kinases to their target substrates. Studies using kinase domains in isolation have not  
358 identified all the known substrates that are found in databases like PhosphoSitePlus. For  
359 example, Shah *et al.* limited their library to only known phosphorylation sites, yet they  
360 could not see a preference for all known kinase substrates for a given kinase<sup>17</sup>. We fully  
361 agree with their discussion<sup>17</sup> that intrinsic kinase domain specificity<sup>13,14,25</sup> acts in concert  
362 with the selective activation and localization provided by the SH2 and SH3 domains in a  
363 cellular context<sup>4,7,8,12,20,21,64,66-70</sup> to provide the full specificity of the NRTKs.

## 364 **Funding**

365 This work was supported by the Howard Hughes Medical Institute (to D.K.), and  
366 NIH grant R01GM107671 to N.J.K.. C.W. is the Marion Abbe Fellow of the Damon  
367 Runyon Cancer Research Foundation (DRG-2343-18). R.O. was supported as an HHMI  
368 Fellow of the Damon Runyon Cancer Research Foundation (DRG-2114-12).

## 369 **Materials and Methods**

## 370 **Kinase Specificity in Whole Cell Lysate Experiment**

371 HEK293 cells were grown in CytoOne 150X20mm TC dishes with DMEM (High  
372 Glucose, No Glutamine; Fisher Sci) containing HyClone bovine growth serum (Fisher  
373 Sci), glutamine (Fisher Sci), fungizone (Fisher Sci), and penicillin-streptomycin (Fisher  
374 Sci). At ~90% confluency, cells were washed with 5mL of PBS before being harvested.  
375 Cells were centrifuged at 4500g for 6 minutes to pellet. PBS was then decanted from the  
376 pellet. Pelleted cells were washed by repeating the previous step. The pellet was then  
377 resuspended in ~3 mL of assay buffer (20 mM Tris, Fisher; 500 mM NaCl, Fisher; 1 mM  
378 MgCl<sub>2</sub>, Fisher; 1mM TCEP, Fisher; pH 8) per 10 plates harvested. Cells were lysed by  
379 sonication followed by centrifugation at 30,000g. The supernatant was pipetted from the  
380 pellet and 20 mM 5'-(4-Fluorosulfonylbenzoyl)adenosine hydrochloride (FSBA; Sigma  
381 Aldrich) was added to the supernatant. The lysate was incubated with FSBA at 25 °C for  
382 2 hours. Lysate was dialyzed in 2 L of assay buffer for 5 hours at room temperature  
383 followed by a second dialysis at 4 °C overnight in 2 L of assay buffer. The protein  
384 concentration in the lysate was then calculated with BCA assay Kit (~9 mg/mL; Pierce).

385 Wilson et al.<sup>29</sup> previously published the construction of an alignment and  
386 phylogenetic model using BALi-Phy<sup>71</sup>, which was used to resurrect ancestral protein  
387 sequences with PAML<sup>72</sup>. The robustness of these ancestral kinases have been previously  
388 tested by investigating activities of alternate sequences of the same nodes<sup>29</sup>. Ancestral  
389 and modern kinases were expressed and purified as previously reported<sup>29</sup>. The reaction  
390 was setup with 10 μM kinase, 20 mM MgCl<sub>2</sub>, 10 mM ATP, Phosphatase Inhibitor Cocktail  
391 #2 (Cal Biotech), and ~1.3 mL of lysate (~12 mg of protein). Reaction went for up to 4  
392 hours at 25 °C. This was repeated three times for each kinase along with a background  
393 sample where no kinase was added. For the western dot blot time course of Src, 5 μL of  
394 sample was quenched with 15 μL of 8 M urea (Fisher Sci) to make the 4-fold dilutions.  
395 Other dilutions were made accordingly for 2-fold and 8-fold dilution samples. 1μL samples  
396 were loaded directly onto a nitrocellulose membrane (GenScript). The protocol was  
397 followed for the iBind Automated Western Systems (ThermoFisher) with Phosphotyrosine  
398 antibody (P-Tyr-1000 MultiMab™ Rabbit mAb mix; Cell Signaling Technology) as the  
399 primary antibody (1:2000 dilution) and ScanLater anti-rabbit antibody (Molecular Devices)  
400 as the secondary antibody (1:5000 dilution). The western dot blot was then imaged on a

401 SpectraMax i3x Multi-Mode Microplate Reader with a ScanLater Western Blot cartridge  
402 (Molecular Devices).

### 403 **Mass Spectrometry of Phosphorylated Kinase-Inactive Lysate**

404 For analysis of phosphorylated peptides, ~1 mL of lysate (~10 mg of protein) which  
405 has been reacted with the kinase previously was quenched with a 1:1 ratio of 8 M urea  
406 (Fisher). Followed by digestion with MS grade trypsin protease (Pierce) as instructed by  
407 the manufacturer. The reaction was then quenched with 10% trifluoroacetic acid (TFA,  
408 Sigma) bringing the final concentration to 0.5% TFA. Samples were desalted under  
409 vacuum using Sep Pak tC18 cartridges (Waters). Each cartridge was activated with 1 mL  
410 80% acetonitrile (ACN)/0.1% TFA, then equilibrated with 3 × 1 mL of 0.1% TFA. Following  
411 sample loading, cartridges were washed with 4 × 1 mL of 0.1% TFA, and samples were  
412 eluted with 4 × 0.5 mL 50% ACN/0.25% formic acid (FA). 20 µg of each sample was kept  
413 for protein abundance measurements, and the remainder was used for phosphopeptide  
414 enrichment. Samples were dried by vacuum centrifugation. The digested peptides were  
415 enriched for phosphopeptides using ion metal affinity column (IMAC). FeCl<sub>3</sub>-NTA beads  
416 were prepared from Ni-NTA super flow slurry beads (QIAGEN) by first stripping the beads  
417 by incubating with 100 mM EDTA in a vacuum manifold three times. The beads were then  
418 washed with water before incubating with 15 mM FeCl<sub>3</sub> (Sigma) for one minute three  
419 times. Excess FeCl<sub>3</sub> was washed with water before rinsing the beads with 0.5% formic  
420 acid (FA). A slurry was prepared by adding water to the beads. 60 µL of slurry was added  
421 into a C18 NEST column that had been equilibrated with 150 µL of 80% ACN, 0.1% TFA.  
422 100 µL of 50% ACN was added to the lyophilized lysate pellet to dissolve it. 100 µL of  
423 100% ACN and 3 µL of 10% TFA was added after lysate was dissolved. 1 mg of peptide  
424 is added to Fe-NTA beads in the desalting tip and then incubated for 1-2 minutes followed  
425 by mixing and incubation for another 1-2 minutes. After incubating liquid is drained. Then  
426 200 µL of 80% ACN and 0.1% TFA was used to wash the beads 3 times. 200 µL of 0.5%  
427 FA is used twice to wash the beads. the beads are incubated for 2-3 minutes with 200 µL  
428 of 500 mM phosphate buffer pH7 before eluting peptides to C18 column. This is repeated  
429 one more time to fully elute the peptides from the beads. The beads are incubated for 15  
430 seconds with 200 µL 0.5% FA before the C18 column is used in a centrifuge to elute the

431 phosphorylated peptides with 75  $\mu$ L of 50% ACN and 0.1% TFA twice before the mass  
432 spectrometry run.

433 The LC/MS/MS was performed on all the samples prepared with phosphopeptide  
434 enrichment and a background sample with no enrichment. We used a 90-minute  
435 separation by nano reversed-phase HPLC gradient over a 75- $\mu$ m ID X 25-cm precolumn  
436 packed with Reprosil C18 1.9- $\mu$ m (Waters). The sample was run on a Q-Exactive Plus  
437 mass spectrometer (ThermoFisher) and the top 20 ions were selected for MS2  
438 sequencing. Resulting data was then searched with MaxQuant<sup>73</sup> to against the human  
439 proteome to identify phosphorylated peptides. The mass spectrometry proteomics data  
440 have been deposited to the ProteomeXchange Consortium via the PRIDE partner  
441 repository<sup>74</sup> with the dataset identifier PXD020299.

## 442 **Specificity Calculations**

443 The results from MaxQuant<sup>73</sup> were analyzed with an in-house script written in  
444 python. For each kinase, a set of substrate sequences was generated from the  
445 phosphorylated peptides found in at least one of the three trials. To generate the set of  
446 substrates, first the substrate peptides would be extended or shortened to 7 residues on  
447 each side of the phosphorylation site. If the sequence was too close to the beginning or  
448 end of a protein it would be rejected immediately. Next, if the sequence is already in the  
449 set of substrate sequences or was found in the control experiment, where no kinase was  
450 added, it would be rejected. Lastly, the localization probability must be greater than or  
451 equal to 70% and the MS intensity must be greater than 0. A background dataset was  
452 generated by applying the same rules to tyrosine containing peptides, from samples that  
453 were not enriched for phosphorylation. The background sequence logo was generated  
454 using WebLogo<sup>75</sup>.

455 Heatmaps were calculated by taking the log frequency of an amino acid occurring  
456 at a specific position minus the log of the frequency in the background. Each position and  
457 amino acid pair were tested using the logs odd ratio estimate used in pLogos and any  
458 significant value was marked with a black square in the heatmap.<sup>10</sup> This allows for  
459 intensities of significant residues to be compared between data sets accurately. Residues  
460 below a 1.3-fold effect size were masked even if significant to focus only on residues

461 which have the largest effect on specificity. Using the amino acid frequencies for each  
462 position the sequence entropy was calculated using the SciPy stats module<sup>76</sup>.

### 463 **Activity Assay**

464 Initial rates were measured by a continuous colorimetric assay<sup>77</sup>. The reactions  
465 contained 20-100 nM of purified kinase along with 20 mM MgCl<sub>2</sub> (Fisher), 525 μM β-  
466 Nicotinamide adenine dinucleotide (Sigma), 4 mM phosphoenolpyruvate (Sigma), 2.5 μL  
467 of PK/LDH (PK 600 U/mL – 1000 U/mL, LDH 900 U/mL-1400 U/mL; Sigma Aldrich), 0.3  
468 mg/mL Bovine Serum Albumin (Fisher), and substrate peptide (GenScript). Reactions  
469 were initiated with 5 mM ATP, by pipetting the solution up and down, and then the  
470 absorbance was read at 340 nm for the course of the reaction (20 minutes) at 25 °C.

### 471 **Homology Model of Bound Peptide**

472 An initial homology model was created from a crystal structure of Abl bound to an  
473 ATP-peptide conjugate (PDB: 2G2I). Using PyRosetta<sup>78</sup>, the initial structure was mutated  
474 to the sequence for either Src or Anc-S1. The bound peptide was then mutated to the  
475 sequence of Src tide. The backbone of the protein and peptide was set to be constrained  
476 before running the Fast Relax protocol using the ref2015 score function.

477

### 478 **References**

- 479 1 Robinson, D. R., Wu, Y. M. & Lin, S. F. The protein tyrosine kinase family of the human  
480 genome. *Oncogene* **19**, 5548-5557, doi:10.1038/sj.onc.1203957 (2000).
- 481 2 Parsons, S. J. & Parsons, J. T. Src family kinases, key regulators of signal transduction.  
482 *Oncogene* **23**, 7906-7909, doi:10.1038/sj.onc.1208160 (2004).
- 483 3 Wang, J. Y. J. The Capable ABL: What Is Its Biological Function? *Molecular and Cellular*  
484 *Biology* **34**, 1188-1197, doi:10.1128/mcb.01454-13 PMID - 24421390 (2014).
- 485 4 Ubersax, J. A. & Ferrell, J. E., Jr. Mechanisms of specificity in protein phosphorylation.  
486 *Nat Rev Mol Cell Biol* **8**, 530-541, doi:10.1038/nrm2203 (2007).
- 487 5 Pinna, L. A. & Ruzzene, M. How do protein kinases recognize their substrates?  
488 *Biochimica et Biophysica Acta (BBA) - Molecular Cell Research* **1314**, 191-225,  
489 doi:10.1016/s0167-4889(96)00083-3 (1996).
- 490 6 Mayer, B. J. & Baltimore, D. Mutagenic analysis of the roles of SH2 and SH3 domains in  
491 regulation of the Abl tyrosine kinase. *Mol Cell Biol* **14**, 2883-2894,  
492 doi:10.1128/mcb.14.5.2883 (1994).
- 493 7 Mayer, B. J., Hirai, H. & Sakai, R. Evidence that SH2 domains promote processive  
494 phosphorylation by protein-tyrosine kinases. *Current Biology* **5**, 296-305,  
495 doi:10.1016/s0960-9822(95)00060-1 (1995).



- 496 8 Yadav, S. S. & Miller, W. T. The evolutionarily conserved arrangement of domains in SRC  
497 family kinases is important for substrate recognition. *Biochemistry* **47**, 10871-10880,  
498 doi:10.1021/bi800930e (2008).
- 499 9 Faux, M. C. & Scott, J. D. More on target with proteinphosphorylation: conferring  
500 specificity by location. *Trends Biochem Sci* **21**, 312-315, doi:10.1016/s0968-  
501 0004(96)10040-2 PMID - 8772386 (1996).
- 502 10 Bhattacharyya, R. P., Remenyi, A., Yeh, B. J. & Lim, W. A. Domains, motifs, and scaffolds:  
503 the role of modular interactions in the evolution and wiring of cell signaling circuits.  
504 *Annu Rev Biochem* **75**, 655-680, doi:10.1146/annurev.biochem.75.103004.142710  
505 (2006).
- 506 11 Remenyi, A., Good, M. C. & Lim, W. A. Docking interactions in protein kinase and  
507 phosphatase networks. *Curr Opin Struct Biol* **16**, 676-685, doi:10.1016/j.sbi.2006.10.008  
508 (2006).
- 509 12 Hantschel, O. *et al.* Structural basis for the cytoskeletal association of Bcr-Abl/c-Abl. *Mol*  
510 *Cell* **19**, 461-473, doi:10.1016/j.molcel.2005.06.030 (2005).
- 511 13 Songyang, Z. *et al.* Catalytic specificity of protein-tyrosine kinases is critical for selective  
512 signalling. *Nature* **373**, 536-539, doi:10.1038/373536a0 (1995).
- 513 14 Till, J. H., Annan, R. S., Carr, S. A. & Miller, W. T. Use of synthetic peptide libraries and  
514 phosphopeptide-selective mass spectrometry to probe protein kinase substrate  
515 specificity. *The Journal of Biological Chemistry* **269**, 7423-7428 (1994).
- 516 15 Scott, M. P. & Miller, W. T. A peptide model system for processive phosphorylation by  
517 Src family kinases. *Biochemistry* **39**, 14531-14537, doi:10.1021/bi001850u (2000).
- 518 16 Schmitz, R., Baumann, G. & Gram, H. Catalytic specificity of phosphotyrosine kinases Blk,  
519 Lyn, c-Src and Syk as assessed by phage display. *J Mol Biol* **260**, 664-677,  
520 doi:10.1006/jmbi.1996.0429 (1996).
- 521 17 Shah, N. H., Lobel, M., Weiss, A. & Kuriyan, J. Fine-tuning of substrate preferences of the  
522 Src-family kinase Lck revealed through a high-throughput specificity screen. *Elife* **7**,  
523 e35190, doi:10.7554/eLife.35190 (2018).
- 524 18 Rychlewski, L., Kschischo, M., Dong, L., Schutkowski, M. & Reimer, U. Target specificity  
525 analysis of the Abl kinase using peptide microarray data. *J Mol Biol* **336**, 307-311,  
526 doi:10.1016/j.jmb.2003.12.052 (2004).
- 527 19 Pawson, T. & Gish, G. D. SH2 and SH3 domains: from structure to function. *Cell* **71**, 359-  
528 362, doi:10.1016/0092-8674(92)90504-6 (1992).
- 529 20 Pellicena, P., Stowell, K. R. & Miller, W. T. Enhanced phosphorylation of Src family kinase  
530 substrates containing SH2 domain binding sites. *J Biological Chem* **273**, 15325-15328,  
531 doi:10.1074/jbc.273.25.15325 (1998).
- 532 21 Machiyama, H., Yamaguchi, T., Sawada, Y., Watanabe, T. M. & Fujita, H. SH3 domain of  
533 c-Src governs its dynamics at focal adhesions and the cell membrane. *Febs J* **282**, 4034-  
534 4055, doi:10.1111/febs.13404 (2015).
- 535 22 Brehme, M. *et al.* Charting the molecular network of the drug target Bcr-Abl. *Proc*  
536 *National Acad Sci* **106**, 7414-7419, doi:10.1073/pnas.0900653106 (2009).
- 537 23 Wu, J. J., Afar, D. E., Phan, H., Witte, O. N. & Lam, K. S. Recognition of multiple substrate  
538 motifs by the c-ABL protein tyrosine kinase. *Comb Chem High Throughput Screen* **5**, 83-  
539 91, doi:10.2174/1386207023330516 (2002).

- 540 24 Hornbeck, P. V. *et al.* PhosphoSitePlus, 2014: mutations, PTMs and recalibrations.  
541 *Nucleic Acids Res* **43**, D512-520, doi:10.1093/nar/gku1267 (2015).
- 542 25 Deng, Y. *et al.* Global analysis of human nonreceptor tyrosine kinase specificity using  
543 high-density peptide microarrays. *J Proteome Res* **13**, 4339-4346,  
544 doi:10.1021/pr500503q (2014).
- 545 26 Harms, M. J. & Thornton, J. W. Evolutionary biochemistry: revealing the historical and  
546 physical causes of protein properties. *Nat Rev Genet* **14**, 559-571, doi:10.1038/nrg3540  
547 (2013).
- 548 27 Krishnan, N. M., Seligmann, H., Stewart, C. B., De Koning, A. P. & Pollock, D. D. Ancestral  
549 sequence reconstruction in primate mitochondrial DNA: compositional bias and effect  
550 on functional inference. *Mol Biol Evol* **21**, 1871-1883, doi:10.1093/molbev/msh198  
551 (2004).
- 552 28 Howard, C. J. *et al.* Ancestral resurrection reveals evolutionary mechanisms of kinase  
553 plasticity. *Elife* **3**, doi:10.7554/eLife.04126 (2014).
- 554 29 Wilson, C. *et al.* Kinase dynamics. Using ancient protein kinases to unravel a modern  
555 cancer drug's mechanism. *Science* **347**, 882-886, doi:10.1126/science.aaa1823 (2015).
- 556 30 Benjamin, D. R. & Marc, A. S. Joint Bayesian Estimation of Alignment and Phylogeny.  
557 *Systematic Biol* **54**, 401-418, doi:10.1080/10635150590947041 PMID - 16012107 (2005).
- 558 31 Williams, P. D., Pollock, D. D., Blackburne, B. P. & Goldstein, R. A. Assessing the accuracy  
559 of ancestral protein reconstruction methods. *PLoS Comput Biol* **2**, e69,  
560 doi:10.1371/journal.pcbi.0020069 (2006).
- 561 32 Ferrari, S. *et al.* Aurora-A site specificity: a study with synthetic peptide substrates.  
562 *Biochem J* **390**, 293-302, doi:10.1042/BJ20050343 (2005).
- 563 33 O'Shea, J. P. *et al.* pLogo: a probabilistic approach to visualizing sequence motifs. *Nat*  
564 *Methods* **10**, 1211-1212, doi:10.1038/nmeth.2646 (2013).
- 565 34 Knight, J. D. *et al.* A novel whole-cell lysate kinase assay identifies substrates of the p38  
566 MAPK in differentiating myoblasts. *Skelet Muscle* **2**, 5, doi:10.1186/2044-5040-2-5  
567 (2012).
- 568 35 Muller, A. C. *et al.* Identifying Kinase Substrates via a Heavy ATP Kinase Assay and  
569 Quantitative Mass Spectrometry. *Sci Rep* **6**, 28107, doi:10.1038/srep28107 (2016).
- 570 36 Schwartz, D. & Gygi, S. P. An iterative statistical approach to the identification of protein  
571 phosphorylation motifs from large-scale data sets. *Nat Biotechnol* **23**, 1391-1398,  
572 doi:10.1038/nbt1146 (2005).
- 573 37 Fuchs, J. E. *et al.* Cleavage entropy as quantitative measure of protease specificity. *PLoS*  
574 *Comput Biol* **9**, e1003007, doi:10.1371/journal.pcbi.1003007 (2013).
- 575 38 Shah, N. H. *et al.* An electrostatic selection mechanism controls sequential kinase  
576 signaling downstream of the T cell receptor. *Elife* **5**, e20105, doi:10.7554/eLife.20105  
577 (2016).
- 578 39 Foda, Z. H., Shan, Y., Kim, E. T., Shaw, D. E. & Seeliger, M. A. A dynamically coupled  
579 allosteric network underlies binding cooperativity in Src kinase. *Nat Commun* **6**, 5939,  
580 doi:10.1038/ncomms6939 (2015).
- 581 40 Colicelli, J. ABL tyrosine kinases: evolution of function, regulation, and specificity. *Sci*  
582 *Signal* **3**, re6, doi:10.1126/scisignal.3139re6 (2010).

- 583 41 Till, J. H., Chan, P. M. & Miller, W. T. Engineering the substrate specificity of the Abl  
584 tyrosine kinase. *J Biological Chem* **274**, 4995-5003, doi:10.1074/jbc.274.8.4995 (1999).
- 585 42 Moses, A. M. & Landry, C. R. Moving from transcriptional to phospho-evolution:  
586 generalizing regulatory evolution? *Trends Genet* **26**, 462-467,  
587 doi:10.1016/j.tig.2010.08.002 (2010).
- 588 43 Manning, G., Whyte, D. B., Martinez, R., Hunter, T. & Sudarsanam, S. The protein kinase  
589 complement of the human genome. *Science* **298**, 1912-1934,  
590 doi:10.1126/science.1075762 (2002).
- 591 44 Copley, S. D. Evolution of new enzymes by gene duplication and divergence. *Febs J* **287**,  
592 1262-1283, doi:10.1111/febs.15299 (2020).
- 593 45 Conant, G. C. & Wolfe, K. H. Turning a hobby into a job: How duplicated genes find new  
594 functions. *Nature Reviews Genetics* **9**, 938-950, doi:10.1038/nrg2482 PMID - 19015656  
595 (2008).
- 596 46 Innan, H. & Kondrashov, F. The evolution of gene duplications: classifying and  
597 distinguishing between models. *Nat Rev Genetics* **11**, 97-108, doi:10.1038/nrg2689  
598 PMID - 20051986 (2010).
- 599 47 Muller, H. J. Bar Duplication. *Science* **83**, 528-530, doi:10.1126/science.83.2161.528-a  
600 (1936).
- 601 48 Ohno, S. Evolution by Gene Duplication. doi:10.1007/978-3-642-86659-3 (1970).
- 602 49 Force, A. *et al.* Preservation of duplicate genes by complementary, degenerative  
603 mutations. *Genetics* **151**, 1531-1545 (1999).
- 604 50 Hughes, A. L. The evolution of functionally novel proteins after gene duplication. *Proc*  
605 *Biol Sci* **256**, 119-124, doi:10.1098/rspb.1994.0058 (1994).
- 606 51 Lynch, M. & Conery, J. S. The evolutionary fate and consequences of duplicate genes.  
607 *Science* **290**, 1151-1155, doi:10.1126/science.290.5494.1151 (2000).
- 608 52 Walsh, J. B. How often do duplicated genes evolve new functions? *Genetics* **139**, 421-  
609 428 (1995).
- 610 53 Wheeler, L. C., Anderson, J. A., Morrison, A. J., Wong, C. E. & Harms, M. J. Conservation  
611 of Specificity in Two Low-Specificity Proteins. *Biochemistry* **57**, 684-695,  
612 doi:10.1021/acs.biochem.7b01086 (2018).
- 613 54 Siddiq, M. A., Hochberg, G. K. & Thornton, J. W. Evolution of protein specificity: insights  
614 from ancestral protein reconstruction. *Curr Opin Struct Biol* **47**, 113-122,  
615 doi:10.1016/j.sbi.2017.07.003 (2017).
- 616 55 Boucher, J. I., Jacobowitz, J. R., Beckett, B. C., Classen, S. & Theobald, D. L. An atomic-  
617 resolution view of neofunctionalization in the evolution of apicomplexan lactate  
618 dehydrogenases. *Elife* **3**, e02304, doi:10.7554/eLife.02304 (2014).
- 619 56 McKeown, A. N. *et al.* Evolution of DNA specificity in a transcription factor family  
620 produced a new gene regulatory module. *Cell* **159**, 58-68, doi:10.1016/j.cell.2014.09.003  
621 (2014).
- 622 57 Aakre, C. D. *et al.* Evolving new protein-protein interaction specificity through  
623 promiscuous intermediates. *Cell* **163**, 594-606, doi:10.1016/j.cell.2015.09.055 PMID -  
624 26478181 (2015).

- 625 58 Pellicena, P. & Miller, W. T. Processive phosphorylation of p130Cas by Src depends on  
626 SH3-polyproline interactions. *J Biological Chem* **276**, 28190-28196,  
627 doi:10.1074/jbc.M100055200 (2001).
- 628 59 Filippakopoulos, P. *et al.* Structural coupling of SH2-kinase domains links Fes and Abl  
629 substrate recognition and kinase activation. *Cell* **134**, 793-803,  
630 doi:10.1016/j.cell.2008.07.047 (2008).
- 631 60 Lorenz, S., Deng, P., Hantschel, O., Superti-Furga, G. & Kuriyan, J. Crystal structure of an  
632 SH2-kinase construct of c-Abl and effect of the SH2 domain on kinase activity. *Biochem J*  
633 **468**, 283-291, doi:10.1042/BJ20141492 (2015).
- 634 61 Songyang, Z. *et al.* Specific motifs recognized by the SH2 domains of Csk, 3BP2, fps/fes,  
635 GRB-2, HCP, SHC, Syk, and Vav. *Mol Cell Biol* **14**, 2777-2785, doi:10.1128/mcb.14.4.2777  
636 (1994).
- 637 62 Lim, W. A., Richards, F. M. & Fox, R. O. Structural determinants of peptide-binding  
638 orientation and of sequence specificity in SH3 domains. *Nature* **372**, 375-379,  
639 doi:10.1038/372375a0 (1994).
- 640 63 Bradshaw, J. M., Mitaxov, V. & Waksman, G. Mutational investigation of the specificity  
641 determining region of the src SH2 domain 1 Edited by J. A. Wells. *Journal of Molecular*  
642 *Biology* **299**, 523-537, doi:10.1006/jmbi.2000.3765 PMID - 10860756 (2000).
- 643 64 Marengere, L. E. *et al.* SH2 domain specificity and activity modified by a single residue.  
644 *Nature* **369**, 502-505, doi:10.1038/369502a0 (1994).
- 645 65 Waksman, G. & Kuriyan, J. Structure and specificity of the SH2 domain. *Cell* **116**, S45-48,  
646 43 p following S48, doi:10.1016/s0092-8674(04)00043-1 (2004).
- 647 66 Ren, R., Ye, Z. S. & Baltimore, D. Abl protein-tyrosine kinase selects the Crk adapter as a  
648 substrate using SH3-binding sites. *Gene Dev* **8**, 783-795, doi:10.1101/gad.8.7.783 PMID -  
649 7926767 (1994).
- 650 67 Moran, M. F. *et al.* Src homology region 2 domains direct protein-protein interactions in  
651 signal transduction. *Proc National Acad Sci* **87**, 8622-8626, doi:10.1073/pnas.87.21.8622  
652 (1990).
- 653 68 Boggan, T. J. & Eck, M. J. Structure and regulation of Src family kinases. *Oncogene* **23**,  
654 7918-7927, doi:10.1038/sj.onc.1208081 (2004).
- 655 69 Shah, N. H., Amacher, J. F., Nocka, L. M. & Kuriyan, J. The Src module: an ancient scaffold  
656 in the evolution of cytoplasmic tyrosine kinases. *Crit Rev Biochem Mol Biol* **53**, 535-563,  
657 doi:10.1080/10409238.2018.1495173 (2018).
- 658 70 Miller, W. T. Determinants of substrate recognition in nonreceptor tyrosine kinases. *Acc*  
659 *Chem Res* **36**, 393-400, doi:10.1021/ar020116v (2003).
- 660 71 Suchard, M. A. & Redelings, B. D. BAli-Phy: simultaneous Bayesian inference of  
661 alignment and phylogeny. *Bioinformatics* **22**, 2047-2048,  
662 doi:10.1093/bioinformatics/btl175 (2006).
- 663 72 Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**, 1586-  
664 1591, doi:10.1093/molbev/msm088 (2007).
- 665 73 Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized  
666 p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nat Biotechnol*  
667 **26**, 1367-1372, doi:10.1038/nbt.1511 (2008).

668 74 Perez-Riverol, Y. *et al.* The PRIDE database and related tools and resources in 2019:  
669 improving support for quantification data. *Nucleic acids research* **47**, D442-D450,  
670 doi:10.1093/nar/gky1106 PMID - 30395289 (2018).

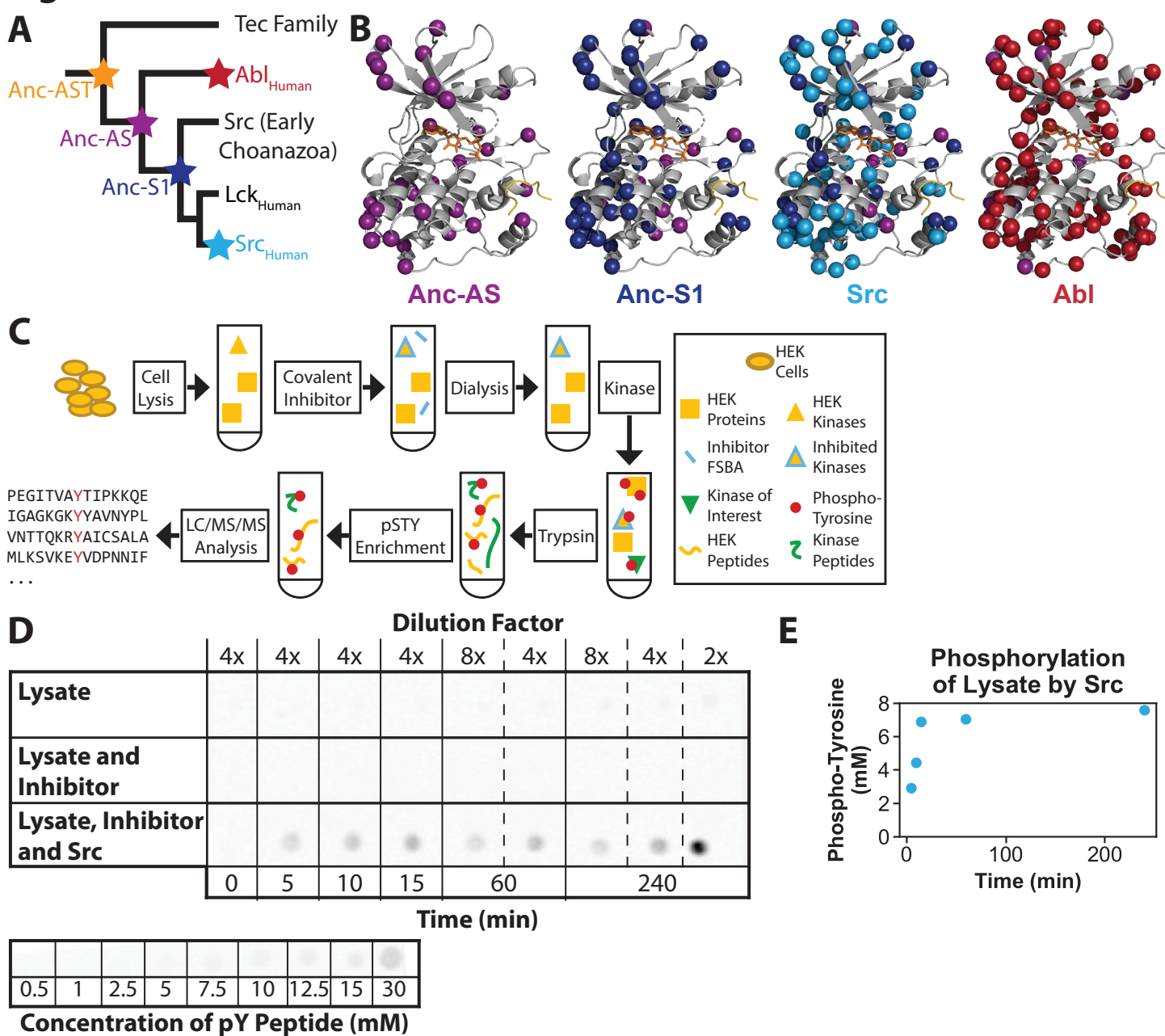
671 75 Crooks, G. E., Hon, G., Chandonia, J.-M. & Brenner, S. E. WebLogo: A Sequence Logo  
672 Generator. *Genome Research* **14**, 1188-1190, doi:10.1101/gr.849004 PMID - 15173120  
673 (2004).

674 76 Virtanen, P. *et al.* SciPy 1.0--Fundamental Algorithms for Scientific Computing in Python.  
675 *Arxiv* **17**, 261-272, doi:10.1038/s41592-019-0686-2 PMID - 32015543 (2019).

676 77 Barker, S. C. *et al.* Characterization of pp60c-src tyrosine kinase activities using a  
677 continuous assay: autoactivation of the enzyme is an intermolecular  
678 autophosphorylation process. *Biochemistry* **34**, 14843-14851, doi:10.1021/bi00045a027  
679 (1995).

680 78 Chaudhury, S., Lyskov, S. & Gray, J. J. PyRosetta: a script-based interface for  
681 implementing molecular modeling algorithms using Rosetta. *Bioinformatics* **26**, 689-691,  
682 doi:10.1093/bioinformatics/btq007 PMID - 20061306 (2010).  
683

## Figure 1



**Analysis of the evolution of non-receptor tyrosine kinase specificity with HEK lysate based approach.** (A) Phylogenetic tree of non-receptor tyrosine kinase domains constructed with Bali-Phy (Suchard *et al.* 2006). The reconstructed nodes and modern kinases used herein are marked with stars. For complete tree, sequences, testing of alternate sequences, and statistics see (Wilson *et al.* 2015). (B) Crystal structure of Abl (PDB ID: 2G2I; Levinson *et al.* 2006) with peptide substrate (yellow) and ADP (orange) bound. The additive differences in primary sequence between Anc-AST and Anc-AS (purple; 87.3% identity), Anc-S1 (dark blue; 85.8% identity), Src (blue; 65.2% identity), and Abl (red; 67.0% identity) are shown. (C) Flow chart displaying how cell lysate was prepared and used for kinase specificity assays. See methods for a detailed description. (D) Western dot blots using anti-pY1000 antibody to detect phosphorylated tyrosine in proteins. (Top) Cell lysate incubated without Src does not show any phosphorylation, whereas lysate treated with the kinase does. Dilutions of the 60-minute and 240-minute time points illustrate that measurements are in the linear range of the dot-blot. (Bottom) Control of phosphorylated peptide blotted at a range of concentrations, diluted 4-fold. (E) Phosphorylation of the cell lysate over time by Src kinase.

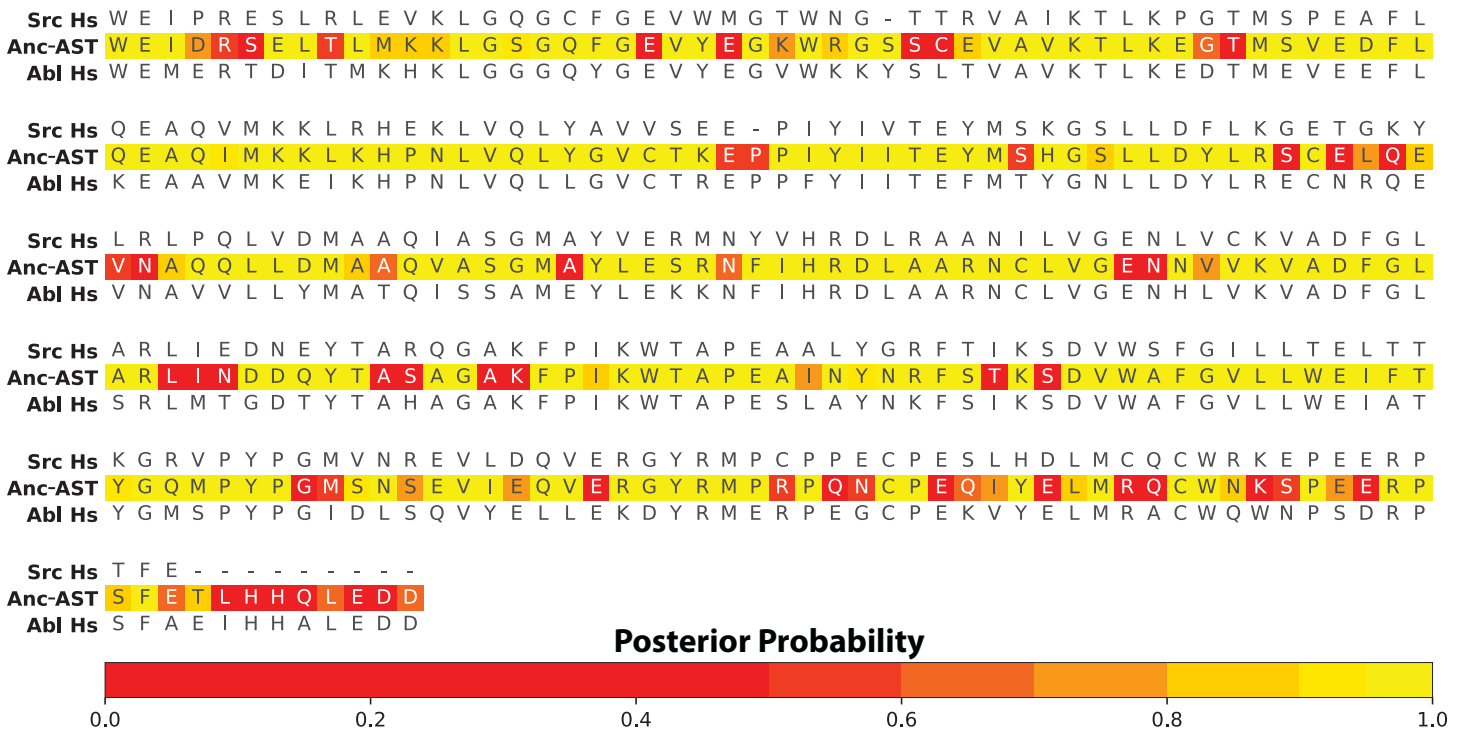


Figure 1 — figure supplement 1: **Posterior probabilities for Anc-AST reconstruction aligned with Src and Abl from *Homo sapiens*.** Wilson *et al.* (2015) performed ancestral sequence reconstruction using PAML (Yang *et al.* 2007) and reported all other ancestral statistics and alignments. Here we report the Anc-AST sequence which has an average posterior probability of 0.86 across all positions. Anc-AST is aligned to the Src and Abl kinase domains and each position displays the given residue's posterior probability from PAML.

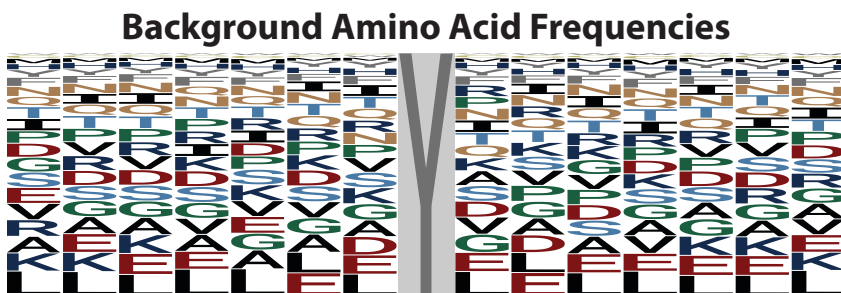


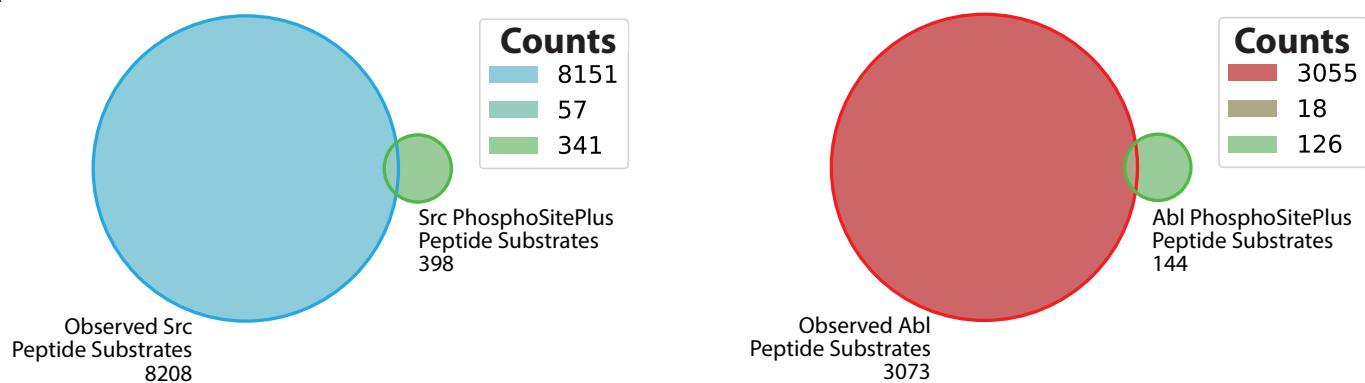
Figure 1 — figure supplement 2: **Frequencies of each amino acid from the unenriched background sample.** Sequence logo generated by WebLogo (Crooks *et al.* 2004) displaying frequency of amino acids in peptides containing a phosphotyrosine from an unenriched sample, that was also not treated with any target kinases.

Figure 1 — figure supplement 3: **BackgroundUnenrichedProteomics.csv** MaxQuant (Cox *et al.* 2008) results from a sample that was not enriched for phosphotyrosine. This data set was used to create the background amino acid frequencies.

Figure 1 — figure supplement 4: **ControlPhosphoproteomicsData.csv** MaxQuant (Cox *et al.* 2008) results from phosphoproteomics experiment where no kinase was added.

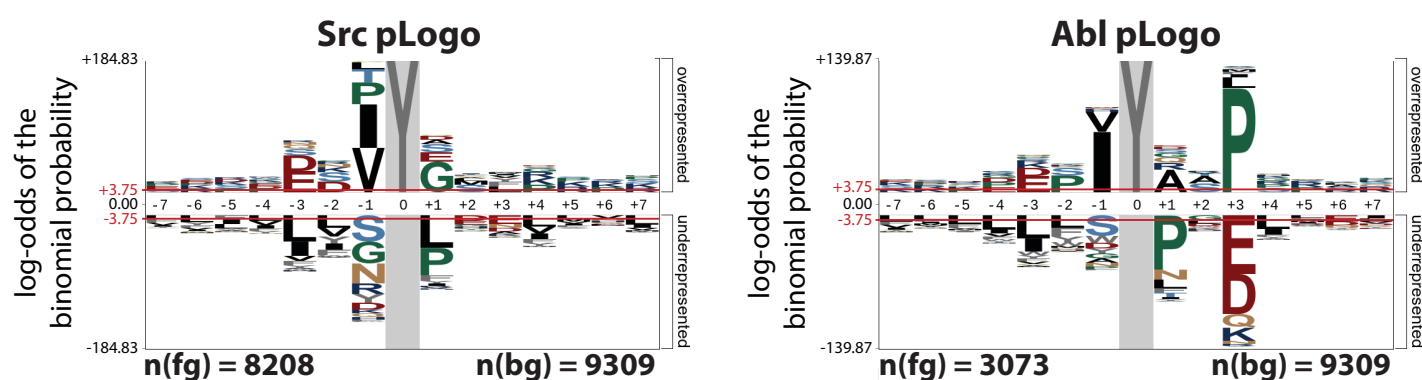
## Figure 2

**A**



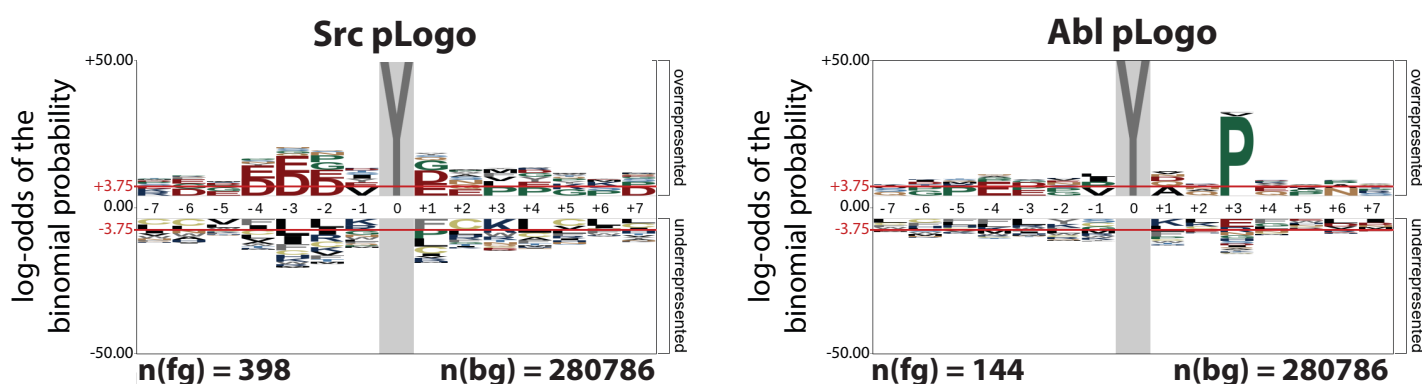
**B**

### Observed Substrates



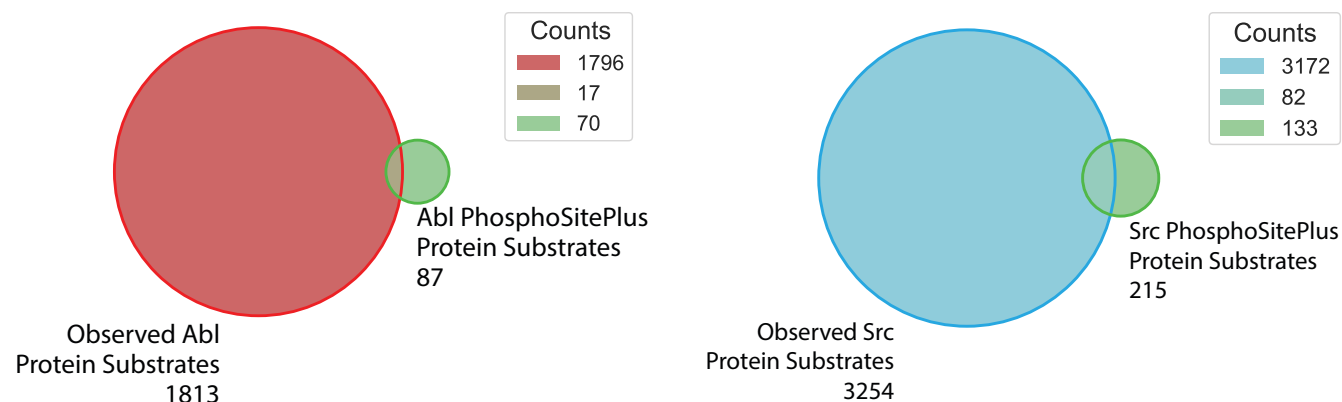
**C**

### PhosphoSitePlus Substrates

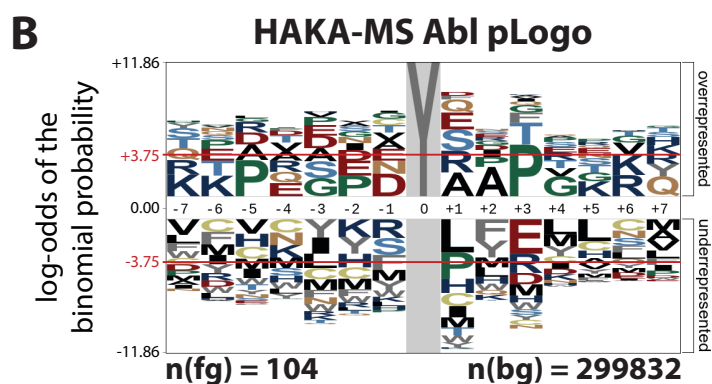
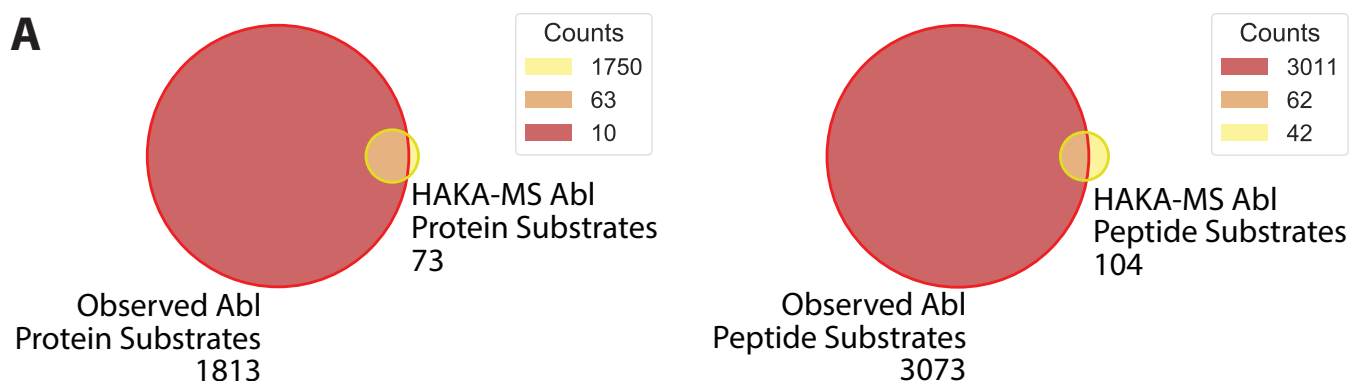


**HEK cell lysate library approach finds large number target substrates and recapitulates known specificities of Src and Abl with improved statistics. (A)** Venn diagrams of sequences that were observed in our Src and Abl kinase specificity experiments compared to the substrates listed in PhosphoSitePlus (Hornbeck *et al.* 2015) for the respective kinase. **(B)** pLogo's for the prevalence of each amino acid in Src and Abl's individual set of substrates (8208 and 3073 sequences, respectively) relative to the background experiment (9309 sequences). **(C)** In comparison, pLogo's generated from the list of human substrates on the PhosphoSitePlus database (Hornbeck *et al.* 2015). Statistical significance level is shown as red line in B and C.



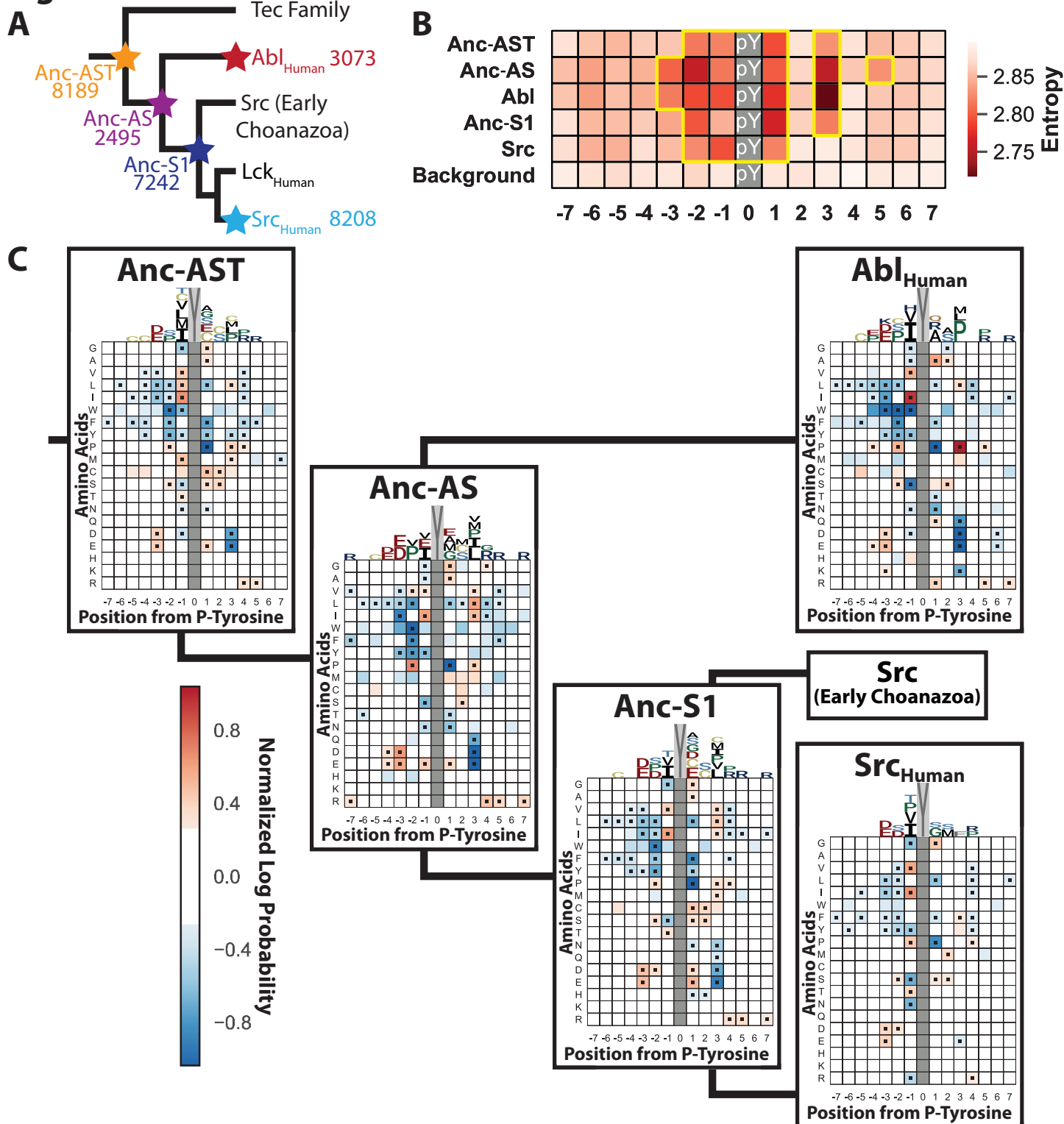


**Figure 2 — figure supplement 1: Natural protein substrates found in phosphoproteomics experiments.** Venn diagrams of proteins that were observed in our kinase specificity experiments compared to the proteins listed in PhosphoSitePlus (Hornbeck *et al.* 2015) for the respective kinase.

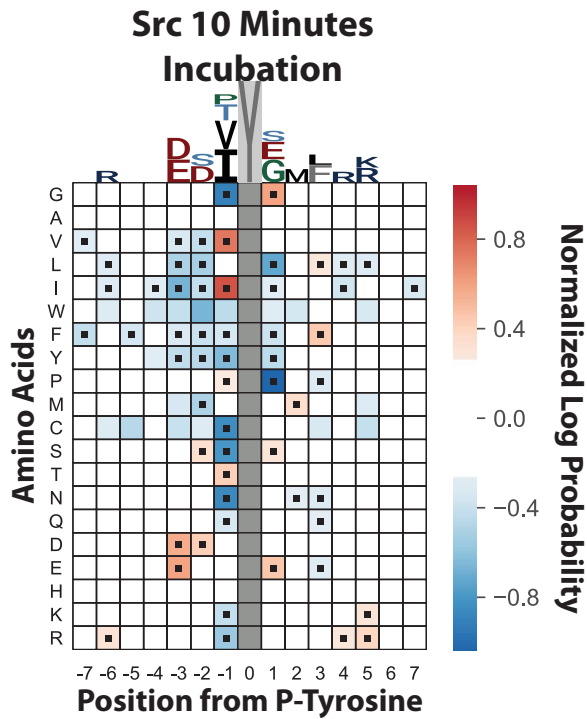


**Figure 2 — figure supplement 2: Comparison of our results with published data from the HAKA-MS experimental procedure (Muller *et al.* 2016).** (A) Comparison of how many substrates (proteins and peptide fragments) were found in either our whole cell lysate assay or the HAKA-MS method. (B) pLogo showing the preference of Abl based on the peptides from the HAKA-MS method. The background used in their study was all tyrosines in the human proteome. Only P<sub>+3</sub> was found to be statistically significant.

### Figure 3



**Anc-AST and Anc-S1 are promiscuous, but are bridged by a relatively specific Anc-AS. (A)** Total number of phosphopeptides found for each of the five kinases are plotted onto the gene tree. **(B)** Positional entropy description for each kinase, where a lower entropy indicates higher specificity. Highlighted in yellow are values below the ~30th percentile for sequence entropy, illustrating the positions with the highest specificity. **(C)** 20x15 positional/amino acid heatmap displaying the normalized log probability of each amino acid at a given position for modern and ancestral kinases. Significant residues at each position are marked with black squares ( $p < 0.05$ ). The sequence logo above the heatmap displays positions with positive values only, and the height of the character is equal to the normalized log probability.

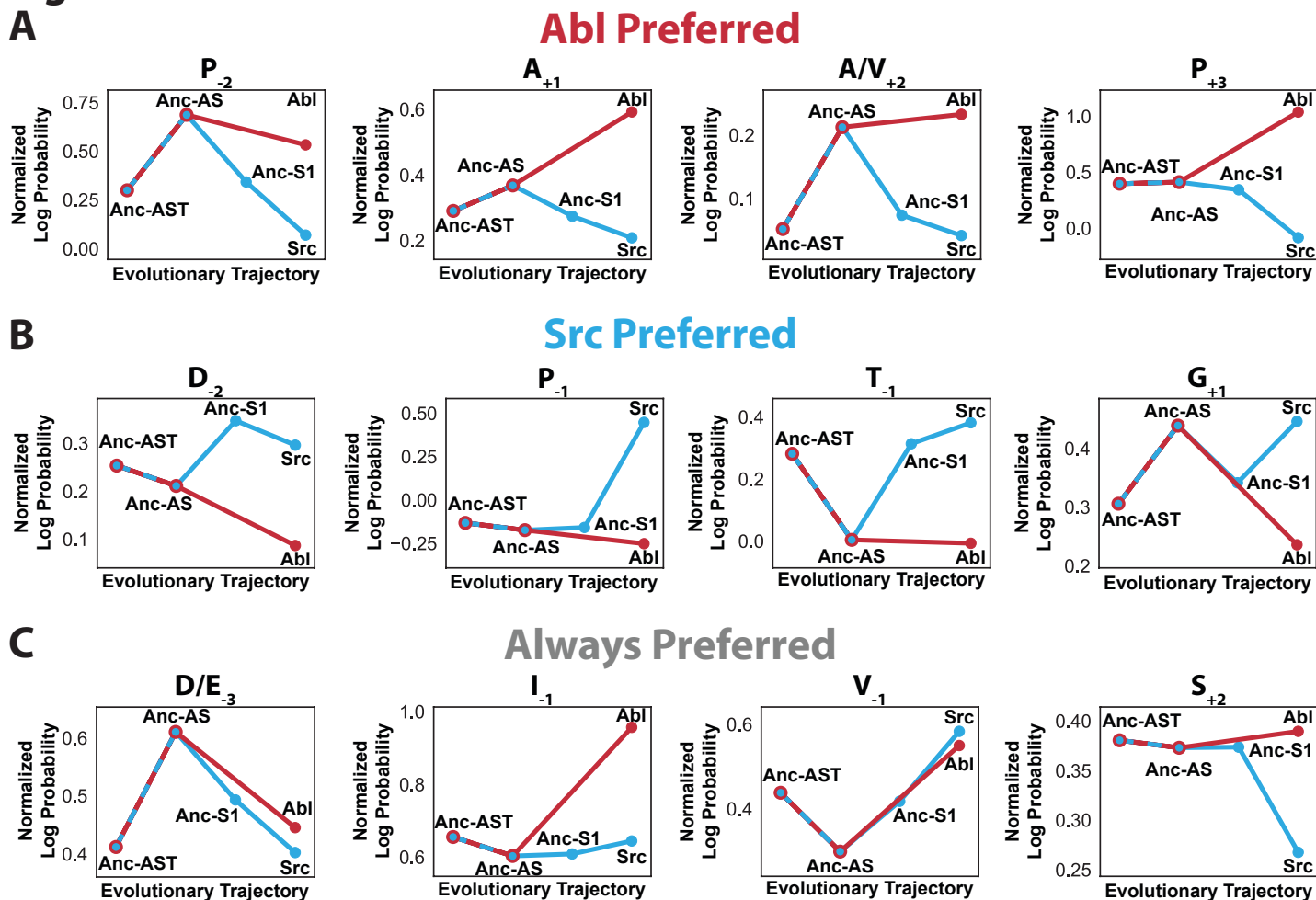


**Figure 3 — figure supplement 1: Substrate specificity is not affected by saturation phosphorylation.** Heatmap displaying enrichment of amino acids at each position in a set of Src substrates after only 10 minutes of incubation. This short time point retains all major features of the heatmap from Src's 4-hour time point in Figure 2C.

**Figure 3 — figure supplement 2: PhosphoproteomicsData.csv** MaxQuant results from phosphoproteomic experiments where kinase was added for the 4-hour incubation.

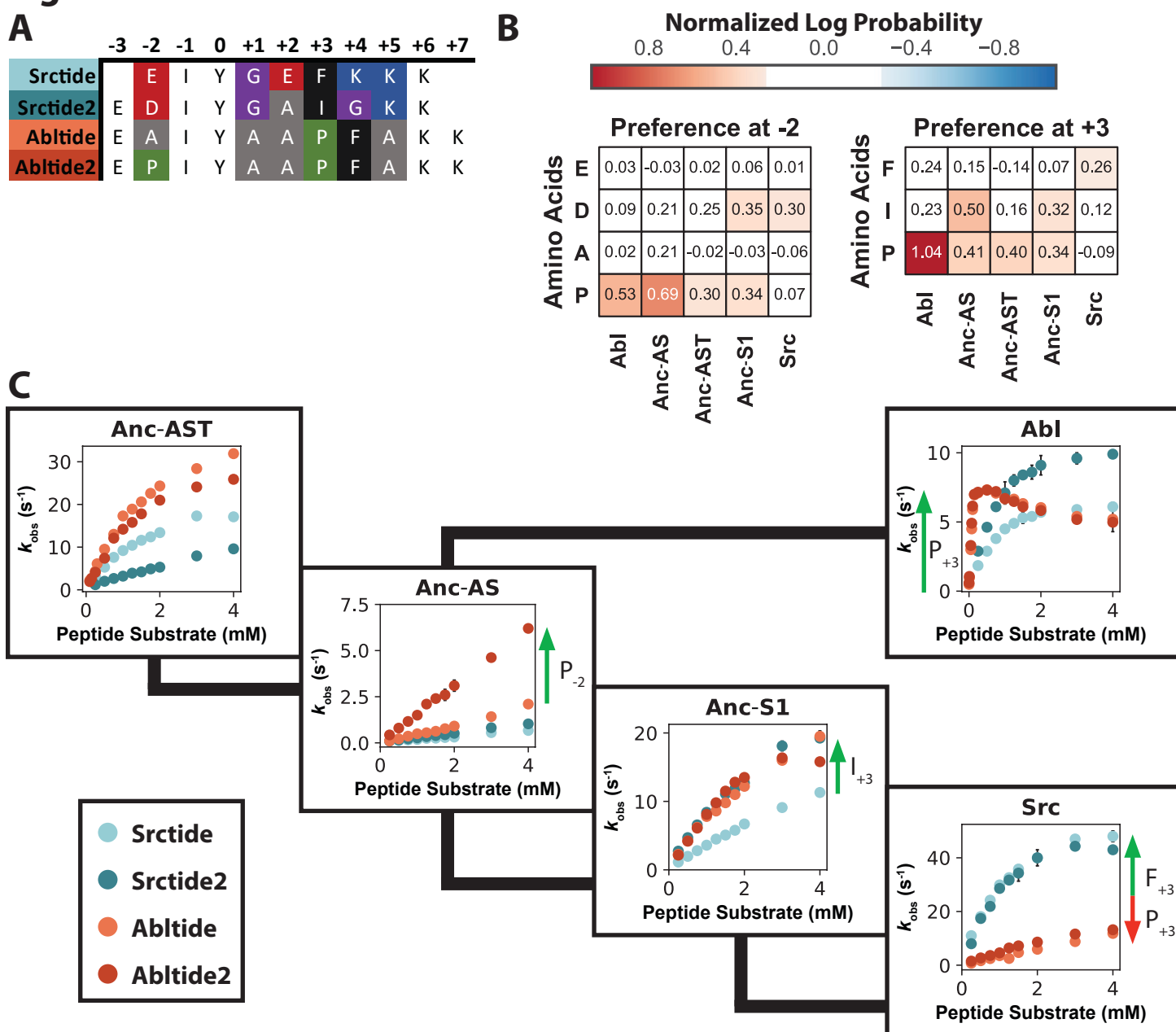
**Figure 3 — figure supplement 3: PhosphoproteomicsSrcShortTimePoint.csv** MaxQuant results from the phosphoproteomic experiment where Src was added for only 10 minutes.

## Figure 4

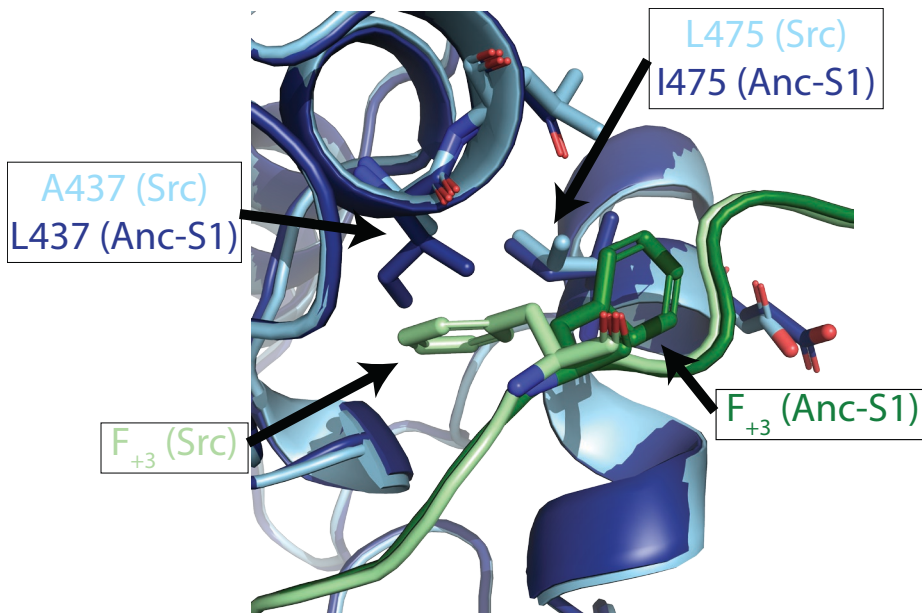


**Evolutionary trajectories of sequence specificity show both subfunctionalization and neofunctionalization.** The normalized log probability of an amino acid occurring in a pool of substrates demonstrates the evolutionary progression for **(A)** Abl specific residues, **(B)** Src specific residues, and **(C)** residue determinants common to all all kinases in this study.

## Figure 5



**Specificity of kinases correlate with individual peptide turnover parameters. (A)** Sequence alignment of the substrate peptides used with differences between peptides highlighted. Coloring of the peptides are used in the Figure 5C. **(B)** Normalized log probability of amino acid preferences for -2 and +3 positions for the five ancestral and modern kinases determined from the cell lysate experiments, see also Fig. 3C. **(C)** Michaelis-Menten curves of phosphorylation for the four peptide substrates and the five different kinases measured with a coupled assay for detecting ADP production (see methods). While it is not possible to saturate and measure accurate  $K_M$  values for many of the substrates, it is possible to observe how the rates are affected under  $k_{cat}/K_M$  conditions. Each kinase was assayed at 20-200 nM, and error bars represent the standard deviation from three measurements. Green arrows indicate the key primary sequence determinants in the different substrates responsible for the large changes in observed rates.



**Figure 5 — figure supplement 1: Homology modelling of protein/substrate complexes suggests amino acid differences between modern and ancestral proteins responsible for the differing substrate specificity in the +3 position.** A crystal structure of Abl containing a peptide bound to the active site (PDB: 2G2I; Levinson *et al.* 2006) was used to build homology models of Src and Anc-S1 bound to Src-tide. The bound peptide in PDB 2G2I already contained F<sub>+3</sub>, but the rest of the Src-tide sequence was modelled in. A Fast Relax protocol was ran in Rosetta with full constraints on the backbone. The zoom-in indicates how the L475I and A437L substitutions in Anc-S1 could prohibit to bind F<sub>+3</sub> in the same pocket as Src.