

# An open library of human kinase domain constructs for automated bacterial expression

Daniel L. Parton,<sup>1</sup> Sonya M. Hanson,<sup>1</sup> Lucelenie Rodríguez-Laureano,<sup>1</sup> Steven K. Albanese,<sup>2</sup> Scott Gradia,<sup>3</sup> Chris Jeans,<sup>3</sup> Markus Seeliger,<sup>4</sup> and John D. Chodera<sup>1,\*</sup>

<sup>1</sup>Computational Biology Program, Sloan Kettering Institute,  
Memorial Sloan Kettering Cancer Center, New York, NY 10065<sup>†</sup>

<sup>2</sup>Gerstner Sloan Kettering Graduate School, Memorial Sloan Kettering Cancer Center, New York, NY 10065<sup>‡</sup>

<sup>3</sup>QB3 MacroLab, University of California, Berkeley, CA 94720<sup>§</sup>

<sup>4</sup>Department of Pharmacological Sciences, Stony Brook University Medical School, Stony Brook, NY 11794<sup>¶</sup>

(Dated: February 2, 2016)

Kinases play a critical role in cellular signaling pathways. Human kinase dysregulation has been linked to a number of diseases, such as cancer, diabetes, and inflammation, and as a result, much of the effort in developing treatments (and perhaps 30% of *all* current drug development effort) has focused on shutting down aberrant kinases with targeted inhibitors. While insect and mammalian expression systems are frequently utilized for the expression of human kinases, they cannot compete with the simplicity and cost-effectiveness of bacterial expression systems, which historically had found human kinases difficult to express. Following the demonstration that phosphatase coexpression could give high yields of Src and Abl kinase domains in inexpensive bacterial expression systems [1], we have performed a large-scale expression screen to generate a library of His-tagged human kinase domain constructs that express well in a simple automated bacterial expression system where phosphatase coexpression (YopH for Tyr kinases, lambda for Ser/Thr kinases) is used. Starting from 96 kinases with crystal structures and any reported bacterial expression, we engineered a library of human kinase domain constructs and screened their coexpression with phosphatase, finding 52 kinases with yields greater than 2 mg/L culture. All sequences and expression data are provided online at <https://github.com/choderalab/kinase-ecoli-expression-panel>, and the plasmids are in the process of being made available through AddGene.

## I. INTRODUCTION

Kinases play a critical role in cellular signaling pathways. Perturbations to these pathways due to mutation, translocation, or upregulation events can cause one or more kinases to become highly active and cease responding normally to regulatory signals, often with disastrous consequences. Kinase dysregulation has been linked to a number of diseases, such as cancer, diabetes, and inflammation. Cancer alone is the second leading cause of death in the United States, accounting for nearly 25% of all deaths; in 2015, over 1.7 million new cases were diagnosed, with over 580,000 deaths [2]. Much of the effort in developing treatments (and perhaps 30% of *all* current drug development effort) has focused on shutting down aberrant kinases with targeted inhibitors.

The discovery of imatinib, which specifically targets the Abl kinase dysregulated in chronic myelogenous leukemia (CML) patients to abate disease progression, was transformative in revealing the enormous therapeutic potential of selective kinase inhibitors, kindling hope that this remarkable success could be recapitulated for other cancers and diseases [3]. While there are now 31 FDA-approved selective kinase inhibitors, these molecules were approved for target-

ing only 13 out of ~500 human kinases, with the vast majority targeting just a handful of kinases; the discovery of therapeutically effective inhibitors for other kinases has proven remarkably challenging.

The ability to probe human kinase biochemistry, biophysics, and structural biology in the laboratory is essential to making rapid progress in the understanding of kinase regulation and the design of selective inhibitors. While human kinase expression in baculovirus-infected insect cells can achieve high success rates [4, 5], it cannot compete in cost or convenience with bacterial expression. While a survey of 62 full-length non-receptor human kinases found that over 50% express well in *E. coli* [4], there is often a desire to express and manipulate only the soluble kinase domains, since these are the molecular targets of therapy for targeted kinase inhibitors and could be studied even for receptor-type kinases. While removal of regulatory domains can negatively impact expression, coexpression with phosphatase was shown to greatly enhance bacterial kinase expression in Src and Abl tyrosine kinases, presumably by ensuring that kinases remain in an unphosphorylated inactive form [1].

The protein databank (PDB) now contains over 100 human kinases that—according to the PDB data records—were expressed in bacteria. Since bacterial expression is often complicated by the need to tailor expression and purification protocols individually for each protein expressed, we wondered whether a simple, uniform, automatable expression and purification protocol could be used to express a large number of human kinases to produce a convenient bacterial expression library to facilitate kinase research and selective inhibitor development. As a first step toward this goal, we developed a structural informatics pipeline to use

\* Corresponding author; [john.chodera@choderalab.org](mailto:john.chodera@choderalab.org)

<sup>†</sup> [daniel.parton@choderalab.org](mailto:daniel.parton@choderalab.org)

<sup>‡</sup> [steven.albanese@choderalab.org](mailto:steven.albanese@choderalab.org)

<sup>§</sup> Current address: Caribou Biosciences, Berkeley, CA 94720; [sgradia@cariboubio.com](mailto:sgradia@cariboubio.com)

<sup>¶</sup> [markus.seeliger@stonybrook.edu](mailto:markus.seeliger@stonybrook.edu)

65 available kinase structural data and associated metadata  
66 to select constructs from available human kinase libraries  
67 to clone into a standard set of vectors intended for phos-  
68 phatase coexpression. Automated expression screening in  
69 Rosetta2 cells found that 52 human kinase domains express  
70 with yields greater than 2 mg/L culture, which should be use-  
71 able for biochemical, biophysical, screening, and structural  
72 biology studies.

73 All code and source files used in this project can  
74 be found at [https://github.com/choderalab/](https://github.com/choderalab/kinase-ecoli-expression-panel)  
75 [kinase-ecoli-expression-panel](https://github.com/choderalab/kinase-ecoli-expression-panel), and a con-  
76 venient sortable table of results can be viewed at  
77 [http://choderalab.github.io/kinome-data/](http://choderalab.github.io/kinome-data/kinase_constructs-addgene_hip_sgc.html)  
78 [kinase\\_constructs-addgene\\_hip\\_sgc.html](http://choderalab.github.io/kinome-data/kinase_constructs-addgene_hip_sgc.html).

## 79 II. METHODS

### 80 A. Semi-automated selection of kinase construct sequences 81 for E. coli expression

#### 82 1. Selection of human protein kinase domain targets

83 Human protein kinases were selected by querying the  
84 UniProt API for any human protein with a domain contain-  
85 ing the string "protein kinase", and which was manually  
86 annotated and reviewed (i.e. a Swiss-Prot entry). The query  
87 string used was:

```
88 taxonomy:"Homo sapiens (Human) [9606]" AND  
89 domain:"protein kinase" AND reviewed:yes  
90 Data was returned by the UniProt API in XML format and  
91 contained protein sequences and relevant PDB structures,  
92 along with many other types of genomic and functional  
93 information. To select active protein kinase domains, the  
94 UniProt domain annotations were searched using the reg-  
95 ular expression ^Protein kinase(?!; truncated)(?!;  
96 inactive), which excludes certain domains annotated  
97 "Protein kinase; truncated" and "Protein kinase; inactive".  
98 Sequences for the selected domains were then stored. The  
99 sequences were derived from the canonical isoform as  
100 determined by UniProt.
```

#### 101 2. Matching target sequences with relevant PDB constructs

102 Each target kinase gene was matched with the same gene  
103 in any other species where present, and UniProt data was  
104 downloaded for those genes also. The UniProt data in-  
105 cluded a list of PDB structures which contain the protein,  
106 as well as their sequence spans in the coordinates of the  
107 UniProt canonical isoform. This information was used to  
108 filter out PDB structures which did not include the pro-  
109 tein kinase domain; structures were kept if they included  
110 the protein kinase domain sequence less 30 residues at  
111 each end. PDB coordinate files were then downloaded for  
112 each PDB entry. The coordinate files contain various meta-  
113 data, including an EXPRESSION\_SYSTEM annotation, which

114 was used to filter PDB entries to keep only those which in-  
115 clude the phrase "ESCHERICHIA COLI". The majority of PDB  
116 entries returned had an EXPRESSION\_SYSTEM tag of "ES-  
117 CHERICHIA COLI", while a small number had "ESCHERICHIA  
118 COLI BL21" or "ESCHERICHIA COLI BL21(DE3).

119 The PDB coordinate files also contain SEQRES  
120 records, which should contain the protein se-  
121 quence used in the crystallography or NMR ex-  
122 periment. According to the PDB documentation  
123 (<http://deposit.rcsb.org/format-faq-v1.html>),  
124 "All residues in the crystal or in solution, including residues  
125 not present in the model (i.e., disordered, lacking electron  
126 density, cloning artifacts, HIS tags) are included in the  
127 SEQRES records." However, we found that these records  
128 are very often misannotated, instead representing only the  
129 crystallographically resolved residues. Since expression  
130 levels can be greatly affected by insertions or deletions  
131 of only one or a few residues at either terminus [6], it is  
132 important to know the full experimental sequence, and  
133 we thus needed a way to measure the authenticity of a  
134 given SEQRES record. We developed a crude measure by  
135 hypothesizing that a) most crystal structures would be  
136 likely to have at least one or a few unresolved residues at  
137 one or both termini and b) the presence of an expression  
138 tag (which is typically not crystallographically resolved)  
139 would indicate an authentic SEQRES record. To achieve  
140 this, unresolved residues were first defined by comparing  
141 the SEQRES sequence to the resolved sequence, using  
142 the SIFTS service to determine which residues were not  
143 present in the canonical isoform sequence. Then regular  
144 expression pattern matching was used to detect common  
145 expression tags at the N- or C-termini. Sequences with a  
146 detected expression tag were given a score of 2, while those  
147 with any unresolved sequence at the termini were given  
148 a score of 1, and the remainder were given a score of 0.  
149 This data was not used to filter out PDB structures at this  
150 stage, but was stored to allow for subsequent selection of  
151 PDB constructs based on likely authenticity. Also stored for  
152 each PDB sequence was the number of residues extraneous  
153 to the target kinase domain, and the number of residue  
154 conflicts with the UniProt canonical isoform within that  
155 domain span.

#### 156 3. Plasmid libraries

157 As a source of kinase DNA sequences, we purchased three  
158 kinase plasmid libraries: the [addgene Human Kinase ORF](#)  
159 [kit](#), a kinase library from the Structural Genomics Consor-  
160 tium (SGC), Oxford (<http://www.thesgc.org>), and a ki-  
161 nase library from the [PlasmID Repository](#) maintained by  
162 the Dana-Farber/Harvard Cancer Center. The aim was to  
163 subclone the chosen sequence constructs from these plas-  
164 mids, though we did not use the same vectors. Annotated  
165 data for the kinases in each library was used to match them  
166 against the human protein kinases selected for this project.  
167 A Python script was written which translated the plasmid  
168 ORFs into protein sequences, and aligned them against the

169 target kinase domain sequences from UniProt. Also calcu-  
170 lated were the number of extraneous protein residues in the  
171 ORF, relative to the target kinase domain sequence, and the  
172 number of residue conflicts.

#### 173 4. Selection of sequence constructs for expression

174 Of the kinase domain targets selected from UniProt, we  
175 filtered out those with no matching plasmids from our avail-  
176 able plasmid libraries and/or no suitable PDB construct se-  
177 quences. For this purpose, a suitable PDB construct se-  
178 quence was defined as any with an authenticity score > 0, i.e.  
179 those derived from SEQRES records with no residues out-  
180 side the span of the resolved structure. Plasmid sequences  
181 and PDB constructs were aligned against each target do-  
182 main sequence, and various approaches were then consid-  
183 ered for selecting a) the sequence construct to use for each  
184 target, and b) the plasmid to subclone it from. Candidate se-  
185 quence constructs were drawn from two sources - PDB con-  
186 structs and the SGC plasmid library. The latter sequences  
187 were included because the SGC plasmid library was the only  
188 one of the three libraries which had been successfully tested  
189 for E. coli expression.

190 For most of the kinase domain targets, multiple candi-  
191 date sequence constructs were available. To select the most  
192 appropriate sequence construct, we sorted them first by au-  
193 thenticity score, then by the number of conflicts relative  
194 to the UniProt domain sequence, then by the number of  
195 residues extraneous to the UniProt domain sequence span.  
196 The top-ranked construct was then chosen. In cases where  
197 multiple plasmids were available, these were sorted first by  
198 the number of conflicts relative to the UniProt domain se-  
199 quence, then by the number of residues extraneous to the  
200 UniProt domain sequence span, and the top-ranked plas-  
201 mid was chosen.

202 This process resulted in a set of 96 kinase domain con-  
203 structs, which (by serendipity) matched the 96-well plate  
204 format we planned to use for parallel expression testing. We  
205 therefore selected these construct sequences for expression  
206 testing.

207 A sortable table of results can be viewed at  
208 [http://choderalab.github.io/kinome-data/  
209 kinase\\_constructs-addgene\\_hip\\_sgc.html](http://choderalab.github.io/kinome-data/kinase_constructs-addgene_hip_sgc.html).

#### 210 5. Other notes

211 While much of this process was performed programmat-  
212 ically using Python, many steps required manual supervi-  
213 sion and intervention. We hope eventually to develop a fully  
214 automated software package for the selection of expression  
215 construct sequences for a given protein family, but this was  
216 not possible within the scope of this article.

## 217 B. Expression testing

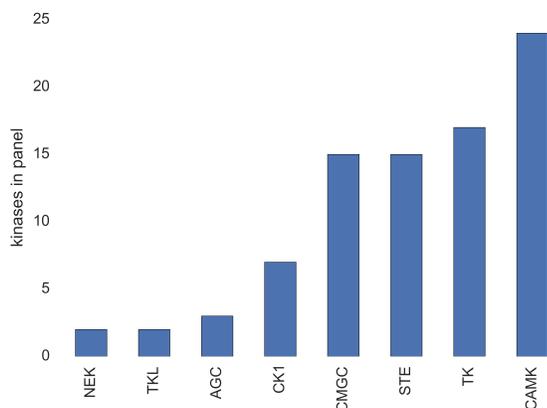
218 For each target, the selected construct sequence was  
219 subcloned from the selected DNA plasmid. Expression  
220 testing was performed by the QB3 MacroLab (QB3 Macro-  
221 Lab, University of California, Berkeley, CA 94720) [[http:  
222 //qb3.berkeley.edu/qb3/macrolab/](http://qb3.berkeley.edu/qb3/macrolab/)], a core facility of-  
223 fering automated gene cloning and recombinant protein ex-  
224 pression and purification services.

225 Each kinase domain was tagged with a N-terminal His10-  
226 TEV and coexpressed with either the truncated YopH164 for  
227 Tyr kinases or lambda phosphatase for Ser/Thr kinases. All  
228 construct sequences were cloned into the 2BT10 plasmid,  
229 an AMP resistant ColE1 plasmid with a T7 promoter, using  
230 LIC (ligation-independent cloning). The inserts were gen-  
231 erated by PCR using the LICv1 forward and reverse tags  
232 on the primers (LICv1 FW= TACTCCAATCCAATGCA; LICv1  
233 RV= TTATCCAATCCAATGTTATTA). Gel purified PCR prod-  
234 ucts were LIC treated with dCTP. Plasmid was linearized, gel  
235 purified and LIC treated with dGTP. LIC-treated plasmid and  
236 insert were mixed together and transformed into XL1-Blues  
237 for plasmid preps.

238 Expression was performed in Rosetta2 cells grown with  
239 Magic Media (Invitrogen autoinducing medium), 100  $\mu$ g/mL  
240 of carbenicillin and 100  $\mu$ g/mL of spectinomycin. Single  
241 colonies of transformants were cultivated with 900  $\mu$ L of  
242 MagicMedia into a gas permeable sealed 96-well block. The  
243 cultures were incubated at 37°C for 4 hours and then at 16°C  
244 for 40 hours while shaking. Next, cells were centrifuged and  
245 the pellets were frozen at -80°C overnight. Cells were lysed  
246 on a rotating platform at room temperature for an hour us-  
247 ing 700  $\mu$ L of SoluLyse (Genlantis) supplemented with 400  
248 mM NaCl, 20 mM imidazole, 1  $\mu$ g/mL pepstatin, 1  $\mu$ g/mL leu-  
249 peptin and 0.5 mM PMSF.

250 For protein purification, 500  $\mu$ L of the soluble lysate was  
251 added to a 25  $\mu$ L Ni-NTA resin in a 96-well filter plate. Nickel  
252 Buffer A (25 mM HEPES pH 7.5, 5% glycerol, 400 mM NaCl,  
253 20 mM imidazole, 1 mM BME) was added and the plate was  
254 shaken for 30 minutes at room temperature. The resin was  
255 washed with 2 mL of Nickel Buffer A. Target proteins were  
256 eluted by a 2 hour incubation at room temperature with 10  
257  $\mu$ g of TEV protease in 80  $\mu$ L of Nickel Buffer A per well and  
258 a subsequent wash with 40  $\mu$ L of Nickel Buffer A to maxi-  
259 mize protein release. Nickel Buffer B (25 mM HEPES pH 7.5,  
260 5% glycerol, 400 mM NaCl, 400 mM imidazole, 1 mM BME)  
261 was used to elute TEV resistant material remaining on the  
262 resin. Untagged protein eluted with TEV protease was run  
263 on a LabChip GX II Microfluidic system to analyze the major  
264 protein species present. Samples of total cell lysate, soluble  
265 cell lysate and Nickel Buffer B elution were run on a SDS-  
266 PAGE for analysis.

267 We are currently making the library of kinase domain  
268 constructs, generated in this work, available for distribu-  
269 tion through the plasmid repository [Addgene](http://Addgene). In the mean-  
270 time, requests for plasmids can be directed to [requests@  
271 choderalab.org](mailto:requests@choderalab.org).



**FIG. 1. Distribution of kinases in expression test panel by family.** Histogram of the 96 kinases in the expression test panel, separated out by kinase family.

### III. RESULTS

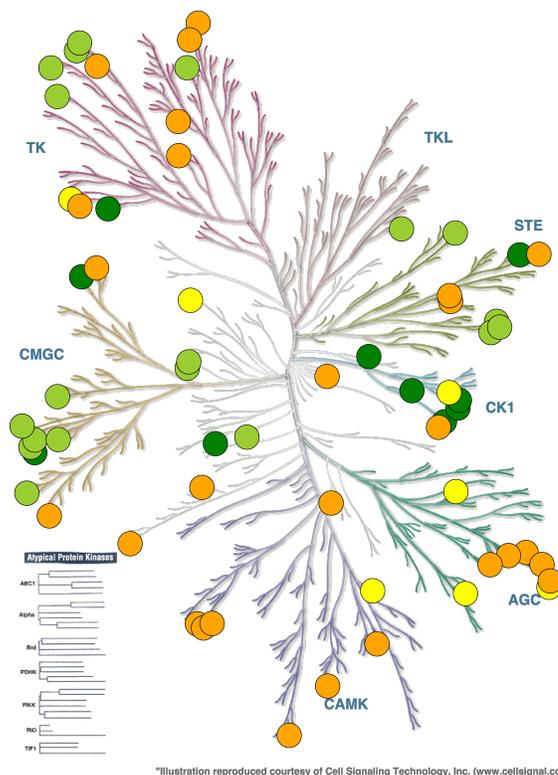
#### A. PDB mining results

Selecting the kinases and their constructs for this expression trial was primarily on the basis of expected success: these specific kinase constructs were bacterially expressed and purified to a degree that a crystal structure could be solved. While the expression protocols used to produce protein for crystallographic studies were likely tailored to maximize expression for individual proteins, we considered these kinases had a high chance of expressing in our semi-automated expression pipeline where the *same* protocol is utilized for all kinases. Statistics of the number of kinases obtained from the PDB mining procedure are shown in Figure 1. Surprisingly, the most highly sampled family was the CAMK family, suggesting that other researchers may have found this family particularly amenable to bacterial expression.

#### B. Small-scale kinase expression test in *E. coli*

A panel containing the 96 kinase domain constructs selected through our semi-automated method, was tested for expression in *E. coli*. From this initial test, 52 kinase domains showed reasonable expression (yield of more than 2 ng/ $\mu$ L eluate, which corresponds to 2 mg/L culture) (Table I). While

the initial panel of 96 kinases was well-distributed across kinase families, the final most highly expressing (yield of more than 12 mg/L kinase) were not as evenly distributed (Figure 2). The 17 most highly expressing kinases showed relatively high purity after elution, though we note that eluting via TEV site cleavage results in a quantity of TEV protease in the eluate (Figure 3).



**FIG. 2. Representation of kinase domain expression results on phylogenetic tree.** Dark green circles represent kinases with expression above 50 mg/L yield. Light green circles represent kinases with expression between 50 and 12 mg/L yield. Yellow circles represent kinases with expression between 12 and 7 mg/L yield. Orange circles represent kinases with any expression (even below 2 mg/L) up to 7 mg/L yield. Image made with KinMap: <http://www.kinhub.org/kinmap>.

### IV. DISCUSSION

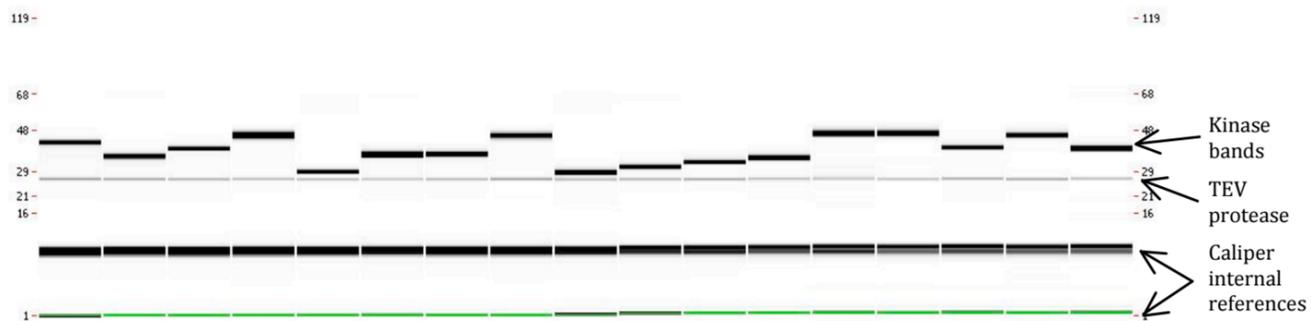
Bacterial coexpression of kinases appears to be a viable approach for studying a wide variety of human kinase domain constructs. We hope that other laboratories find these resources useful in their own work.

[1] M. A. Seeliger, M. Young, M. N. Henderson, P. Pellicena, D. S. King, A. M. Falick, and J. Kuriyan, *Protein Sci.* **14**, 3135 (2005).  
[2] American Cancer Society, *Cancer Facts & Figures 2015*, 2015.  
[3] F. Stegmeier, M. Warmuth, W. R. Sellers, and M. Dorsch, *Clin. Pharm. & Therap.* **87**, 543 (2010).

[4] S. P. Chambers, D. A. Austen, J. R. Fulghum, and W. M. Kim, *Protein Expression and Purification* **36**, 40 (2004).  
[5] L. Wang, M. Foster, Y. Zhang, W. R. Tschantz, L. Yang, J. Worrall, C. Loh, and X. Xu, *Protein Express. Pur.* **61**, 204 (2008).

kinase expressed	phosphatase co-expressed	expected scale-up culture (mg/L)
MK14_HUMAN_D0	Lambda	70.7
VRK3_HUMAN_D0	Lambda	67.5
GAK_HUMAN_D0	Lambda	64.7
CSK_HUMAN_D0	Truncated YopH164	62.5
VRK1_HUMAN_D0	Lambda	62.3
KC1G3_HUMAN_D0	Lambda	56.3
FES_HUMAN_D0	Truncated YopH164	44.0
PMYT1_HUMAN_D0	Lambda	38.0
MK03_HUMAN_D0	Lambda	36.4
STK3_HUMAN_D0	Lambda	34.3
DYR1A_HUMAN_D0	Lambda	34.1
KC1G1_HUMAN_D0	Lambda	34.1
MK11_HUMAN_D0	Lambda	31.7
MK13_HUMAN_D0	Lambda	31.7
EPHB1_HUMAN_D0	Truncated YopH164	28.9
MK08_HUMAN_D0	Lambda	28.5
CDK16_HUMAN_D0	Lambda	26.9
EPHB2_HUMAN_D0	Truncated YopH164	25.1
PAK4_HUMAN_D0	Lambda	23.9
CDKL1_HUMAN_D0	Lambda	23.2
SRC_HUMAN_D0	Truncated YopH164	22.0
STK16_HUMAN_D0	Lambda	20.7
MAPK3_HUMAN_D0	Lambda	18.8
PAK6_HUMAN_D0	Lambda	18.0
CSK22_HUMAN_D0	Lambda	17.9
MERTK_HUMAN_D0	Truncated YopH164	16.8
PAK7_HUMAN_D0	Lambda	14.7
CSK21_HUMAN_D0	Lambda	14.5
EPHA3_HUMAN_D0	Truncated YopH164	14.1
BMPR2_HUMAN_D0	Lambda	14.1
M3K5_HUMAN_D0	Lambda	14.0
KCC2G_HUMAN_D0	Lambda	13.3
E2AK2_HUMAN_D0	Lambda	11.6
MK01_HUMAN_D0	Lambda	11.2
CSKP_HUMAN_D0	Lambda	10.1
CHK2_HUMAN_D0	Lambda	8.1
KC1G2_HUMAN_D0	Lambda	7.6
DMPK_HUMAN_D0	Lambda	7.6
KCC2B_HUMAN_D0	Lambda	7.1
FGFR1_HUMAN_D0	Truncated YopH164	6.1
KS6A1_HUMAN_D1	Lambda	5.7
DAPK3_HUMAN_D0	Lambda	4.0
STK10_HUMAN_D0	Lambda	3.7
KC1D_HUMAN_D0	Lambda	3.7
KC1E_HUMAN_D0	Lambda	3.5
NEK1_HUMAN_D0	Lambda	3.3
CDK2_HUMAN_D0	Lambda	3.1
ABL1_HUMAN_D0	Truncated YopH164	2.5
DAPK1_HUMAN_D0	Lambda	2.4
DYRK2_HUMAN_D0	Lambda	2.4
HASP_HUMAN_D0	Lambda	2.3
FGFR3_HUMAN_D0	Truncated YopH164	2.3
EPHB3_HUMAN_D0	Truncated YopH164	1.7
SLK_HUMAN_D0	Lambda	1.6
KCC2D_HUMAN_D0	Lambda	1.6
NEK7_HUMAN_D0	Lambda	1.3
PHKG2_HUMAN_D0	Lambda	1.3
VRK2_HUMAN_D0	Lambda	1.2
AAPK2_HUMAN_D0	Lambda	1.1
AURKA_HUMAN_D0	Lambda	1.1
MARK3_HUMAN_D0	Lambda	1.1
KAPCA_HUMAN_D0	Lambda	0.9
STK24_HUMAN_D0	Lambda	0.8
VGFR1_HUMAN_D0	Truncated YopH164	0.5
KCC4_HUMAN_D0	Lambda	0.4
KCC1G_HUMAN_D0	Lambda	0.3
KCC2A_HUMAN_D0	Lambda	0.3
FAK2_HUMAN_D0	Truncated YopH164	0.3

**TABLE I. Expression results by kinase.** Yield (determined by Caliper GX II quantitation of the expected size band) reported in mg/L culture, where total eluate volume was 120  $\mu$ l from 900  $\mu$ l bacterial culture.



**FIG. 3. Synthetic gel image rendering of highest expressing kinases.** Caliper GX II synthetic gel image rendering of kinases expressing > 25 mg/L culture from microfluidic capillary electrophoresis quantitation.

<sup>316</sup> [6] H. E. Klock, E. J. Koesema, M. W. Knuth, and S. A. Lesley, Pro-  
<sup>317</sup> teins: Structure, Function, and Bioinformatics **71**, 982 (2008).