

Isofunctional Protein Subfamily Detection using Data Integration and Spectral Clustering

Elisa Boari de Lima^{1,2,*}, Wagner Meira Júnior², Raquel Cardoso de Melo-Minardi²

1 Department of Biochemistry and Immunology, Federal University of Minas Gerais, Belo Horizonte, MG, Brazil

2 Department of Computer Science, Federal University of Minas Gerais, Belo Horizonte, MG, Brazil

* eblima@dcc.ufmg.br

S7 Text: Dividing the protein kinases into seven clusters

The second level of ASMC's [1] hierarchical clustering has seven clusters, whose logos and compositions according to the subfamily labels are presented in Fig. S7.1. One may easily note that, even after increasing the number of clusters, the subfamilies remained mixed. Only the EGFRs were mostly isolated from the other subfamilies, which had already been done in the previous hierarchy level. Furthermore, ASMC generated clusters with irrelevant differences, such as Clusters I.A and I.B. This hierarchy level's utility seems to have been limited to obtaining some relatively uniform Tyr kinase subgroups.

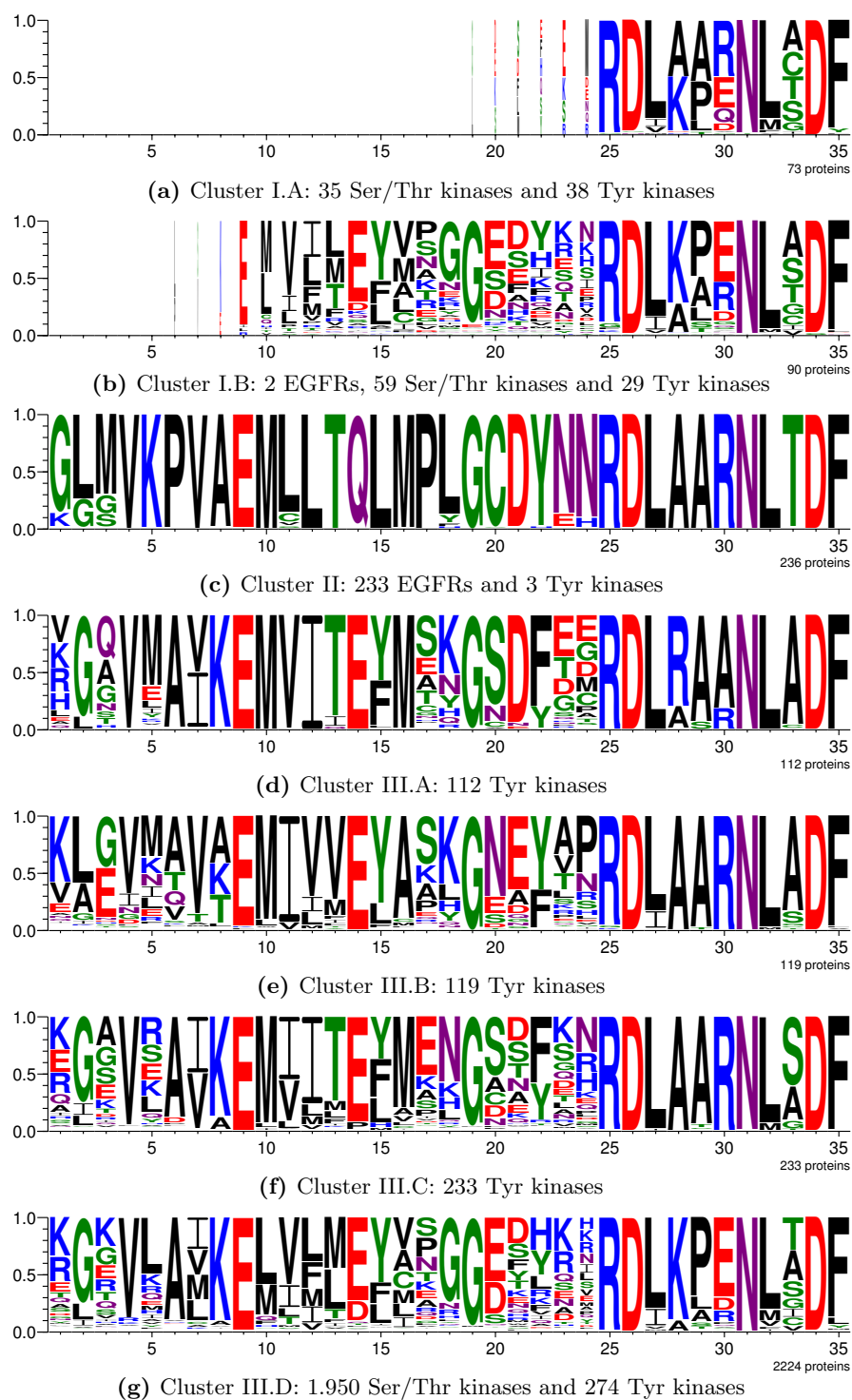


Figure S7.1. Protein kinase division into seven clusters in the second level of ASMC's hierarchical clustering.

Given ASMC's criteria for considering as specificity determining positions (SDPs) those positions with p-values smaller than 0.0001 [1], the SDPs per cluster for this clustering are presented in Table S7.1. Again, one may notice the presence of the known SDPs for the family, presented in the main text, although the subfamilies are mixed. This is due to the SDPs being valid for the family as a whole, and to ASMC's SDP criteria concerning positions, rather than specific residues in such positions.

Table S7.1. **Cluster SDPs for the seven protein kinase clusters produced by ASMC.**

| Cluster | Positions |
|---------|--|
| I.A | - |
| I.B | - |
| II | 1 ₁₃₀ , 2 ₁₃₂ , 3 ₁₃₃ , 5 ₁₄₂ , 6 ₁₅₆ , 7 ₁₅₇ , 8 ₁₅₈ , 10 ₁₈₁ , 11 ₁₉₀ , 12 ₂₀₃ , 13 ₂₀₅ , 14 ₂₀₆ , 15 ₂₀₇ , 16 ₂₀₈ , 17 ₂₀₉ , 18 ₂₁₀ , 20 ₂₁₂ , 21 ₂₁₅ , 22 ₂₁₆ , 23 ₂₁₉ , 24 ₂₂₀ , 28₂₅₄ , 29₂₅₅ , 30₂₅₆ , 33 ₂₆₉ |
| III.A | 5 ₁₄₂ , 12 ₂₀₃ , 18 ₂₁₀ , 20 ₂₁₂ , 28₂₅₄ , 29₂₅₅ , 30₂₅₆ |
| III.B | 2 ₁₃₂ , 8 ₁₅₈ , 11 ₁₉₀ , 12 ₂₀₃ , 13 ₂₀₅ , 18 ₂₁₀ , 20 ₂₁₂ , 21 ₂₁₅ , 24 ₂₂₀ , 28₂₅₄ , 29₂₅₅ , 30₂₅₆ |
| III.C | 10 ₁₈₁ , 12 ₂₀₃ , 13 ₁₃₃ , 16 ₂₀₈ , 17 ₂₀₉ , 18 ₂₁₀ , 20 ₂₁₂ , 22 ₂₁₆ , 28₂₅₄ , 29₂₅₅ , 30₂₅₆ , 33 ₂₆₉ |
| III.D | 2 ₁₃₂ , 3 ₁₃₃ , 5 ₁₄₂ , 6 ₁₅₆ , 8 ₁₅₈ , 10 ₁₈₁ , 11 ₁₉₀ , 13 ₂₀₅ , 16 ₂₀₈ , 18 ₂₁₀ , 20 ₂₁₂ , 21 ₂₁₅ , 22 ₂₁₆ , 23 ₂₁₉ , 28₂₅₄ , 29₂₅₅ , 30₂₅₆ |

Listed in order of active site position. Positions in bold correspond to known SDPs. Subscripted positions correspond to those in PDB structure 1U46:A.

Although the genetic programming (GP) system successfully separated the subfamilies into two and three clusters, it was run with seven clusters in order to compare results with ASMC's hierarchy's second level, which presented MI = 45.99. Obviously, this caused the subfamilies to be broken into subgroups. The best result presented MI = 50.70 and involved four data types, as shown in the main text. The logos and compositions for the seven clusters are presented in Fig. S7.2, while residues which most distinguish each cluster are presented in Table S7.2. One may notice that known SDPs are present for the cases in which the residues which occur at such positions distinguish the corresponding cluster from the others.

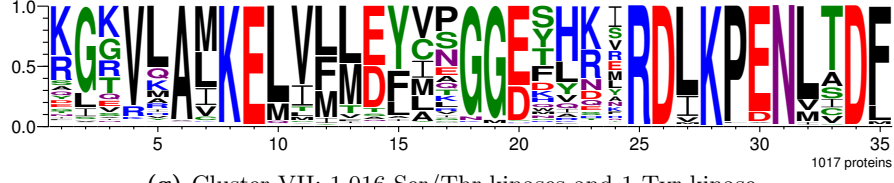
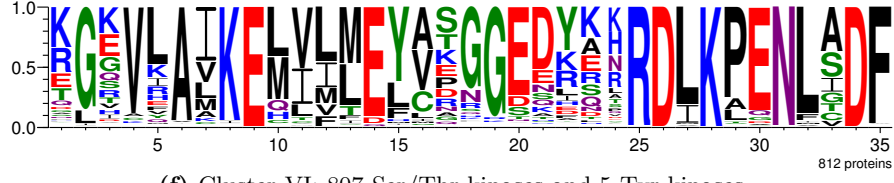
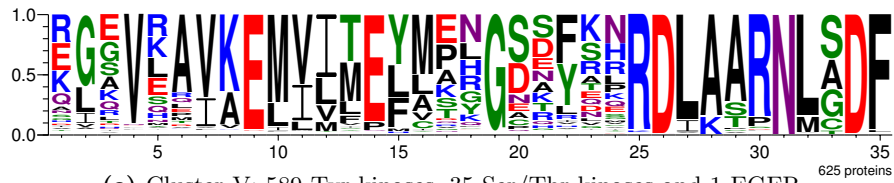
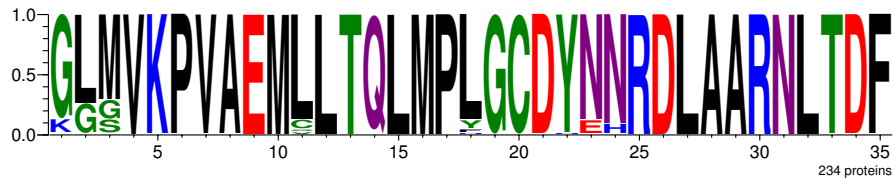
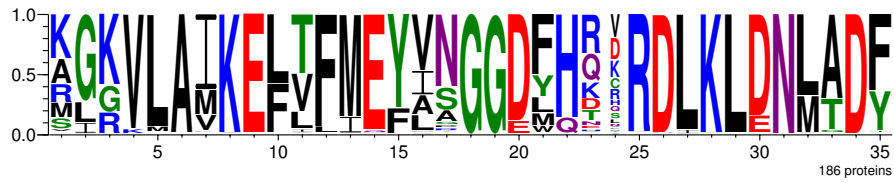
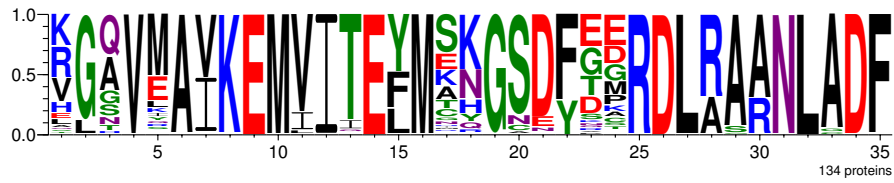
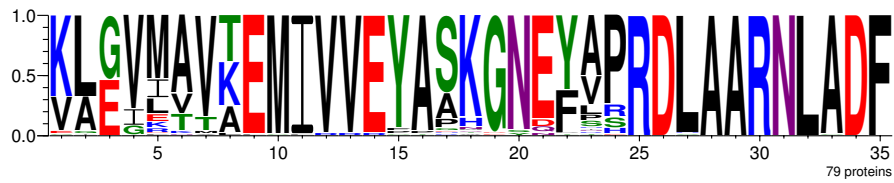


Figure S7.2. Protein kinase division into seven clusters by the GP system.

Table S7.2. Most important residues for the seven protein kinase clusters produced by the GP system.

| Cluster | Residues |
|---------|--|
| I | V13 ₂₀₅ , N20 ₂₁₂ , V12 ₂₀₃ , K18 ₂₁₀ , E21 ₂₁₅ , P24 ₂₂₀ , A16 ₂₀₈ |
| II | R28 ₂₅₄ , A30 ₂₅₆ , I12 ₂₀₃ , S20 ₂₁₂ , M16 ₂₀₈ , F22 ₂₁₆ , T13 ₂₀₅ |
| III | L29 ₂₅₅ , D30 ₂₅₆ , F12 ₂₀₃ , D20 ₂₁₂ , H22 ₂₁₆ , N17 ₂₀₉ , F10 ₁₈₁ |
| IV | C20 ₂₁₂ , P6 ₁₅₆ , Q14 ₂₀₆ , G1 ₁₃₀ , L18 ₂₁₀ , A8 ₁₅₈ , M3 ₁₃₃ , N23 ₂₁₉ , K5 ₁₄₂ , N24 ₂₂₀ |
| V | R30 ₂₅₆ , A28 ₂₅₄ , A29 ₂₅₅ , V7 ₁₅₇ , F22 ₂₁₆ , I12 ₂₀₃ , S20 ₂₁₂ |
| VI | E30 ₂₅₆ , K28 ₂₅₄ , E20 ₂₁₂ , P29 ₂₅₅ , D21 ₂₁₅ |
| VII | P29 ₂₅₅ , L10 ₁₈₁ , G18 ₂₁₀ , E30 ₂₅₆ , K28 ₂₅₄ , D14 ₂₀₆ , H22 ₂₁₆ , T33 ₂₆₉ , M7 ₁₅₇ |

Listed in decreasing order of partial MI value. Residues in bold correspond to known SDPs. Subscripted positions correspond to those in PDB structure 1U46:A.

In Cluster V, whose majority is labeled as Tyr kinases, there are 35 Ser/Thr kinase-labeled proteins and one labeled as EGFR (Q2HZD7). The latter is unreviewed and automatically annotated with InterPro domain IPR016245 (*Tyr protein kinase, EGF/ERP/XmrK receptor*), i.e., despite being a Tyr kinase as all EGFRs are, it should have been clustered along with the other EGFRs. Among the 35 labeled as Ser/Thr kinases, two (Q54TM7 and A7J1T0) have been manually reviewed and annotated with InterPro domain IPR008271 (*serine/threonine-protein kinase, active site*), and also with GO term *protein serine/threonine kinase activity*. Only one of the other 33 unreviewed proteins, namely Q5AUJ7, lacks an annotation with the Ser/Thr kinase active site InterPro domain, yet is annotated with domain IPR002290 (*Ser/Thr/dual specificity protein kinase, catalytic domain*), which is present either in Ser/Thr kinases or dual specificity proteins. Such annotations indicate the GP system wrongly placed these 36 proteins in a Tyr kinase cluster.

In Cluster VI, of majority Ser/Thr kinases, there are five Tyr-kinase labeled proteins, all unreviewed. Two of them (Q69U56 and A2DGV6) are annotated with InterPro domain IPR008271 (*serine/threonine-protein kinase, active site*), while the other three (Q5RAR7, Q4T0K5, and A0JN96) are annotated with domain IPR002290 (*Ser/Thr/dual specificity protein kinase, catalytic domain*). Furthermore, among the five, only Q5RAR7 has not been annotated with GO term *protein serine/threonine kinase activity*. Hence, these annotations suggest these five subfamily labels are incorrect and that they were correctly placed into a Ser/Thr kinase cluster by the GP system. Lastly, Cluster VII shows a Tyr kinase-labeled protein among over a thousand Ser/Thr kinases: unreviewed protein Q6K3D4, which lacks subfamily-related annotations, so we cannot affirm if its label is incorrect or if the GP system clustered it improperly.

Despite wrongly positioning 36 proteins, which corresponds to less than 1.2% of the family, the clusters generated by the GP system were shown to be in much more agreement with the subfamily labels, as well as more contrasting to each other, than those produced by ASMC.

References

1. Melo-Minardi RC, Bastard K, Artiguenave F. Identification of subfamily-specific sites based on active sites modeling and clustering. *Bioinformatics*. 2010 Dec;26(24):3075–3082.