

References

- AHMAD, S. & SARAI, A. (2005). PSSM-based prediction of DNA binding sites in proteins. *BMC Bioinformatics*, **6**, 33. [26](#)
- ALBER, T., GILBERT, W.A., PONZI, D.R. & PETSKO, G.A. (1983). The role of mobility in the substrate binding and catalytic machinery of enzymes. *Ciba Found Symp*, **93**, 4–24. [111](#)
- AMICO, M., FINELLI, M., ROSSI, I., ZAULI, A., ELOFSSON, A., VIKLUND, H., von HEIJNE, G., JONES, D., KROGH, A., FARISELLI, P., MARTELLI, P.L. & CASADIO, R. (2006). PONGO: a web server for multiple predictions of all-alpha transmembrane proteins. *Nucleic Acids Res*, **34**, W169–W172. [137](#)
- ANDERSEN, J.S., LYON, C.E., FOX, A.H., LEUNG, A.K.L., LAM, Y.W., STEEN, H., MANN, M. & LAMOND, A.I. (2002). Directed proteomic analysis of the human nucleolus. *Curr Biol*, **12**, 1–11. [113](#)
- ANDERSEN, J.S., LAM, Y.W., LEUNG, A.K.L., ONG, S.E., LYON, C.E., LAMOND, A.I. & MANN, M. (2005). Nucleolar proteome dynamics. *Nature*, **433**, 77–83. [113](#)

REFERENCES

- BAILEY, T.L. & ELKAN, C. (1994). Fitting a mixture model by expectation maximization to discover motifs in biopolymers. *Proc Int Conf Intell Syst Mol Biol*, **2**, 28–36. [3](#)
- BAILEY, T.L. & ELKAN, C. (1995). The value of prior knowledge in discovering motifs with meme. *Proc Int Conf Intell Syst Mol Biol*, **3**, 21–29. [18](#), [27](#), [64](#)
- BAIROCH, A. & APWEILER, R. (1996). The SWISS-PROT protein sequence data bank and its new supplement TREMBL. *Nucleic Acids Res*, **24**, 21–25. [12](#), [40](#)
- BAIROCH, A. & APWEILER, R. (2000). The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res*, **28**, 45–48. [12](#)
- BANNAI, H., TAMADA, Y., MARUYAMA, O., NAKAI, K. & MIYANO, S. (2002). Extensive feature detection of N-terminal protein sorting signals. *Bioinformatics*, **18**, 298–305. [7](#), [9](#)
- BATEMAN, A., COIN, L., DURBIN, R., FINN, R.D., HOLLICH, V., GRIFFITHS-JONES, S., KHANNA, A., MARSHALL, M., MOXON, S., SONNHAMMER, E.L.L., STUDHOLME, D.J., YEATS, C. & EDDY, S.R. (2004). The Pfam protein families database. *Nucleic Acids Res*, **32**, D138–D141. [11](#), [25](#), [58](#)
- BENDTSEN, J.D., JENSEN, L.J., BLOM, N., HEIJNE, G.V. & BRUNAK, S. (2004a). Feature-based prediction of non-classical and leaderless protein secretion. *Protein Eng Des Sel*, **17**, 349–356. [65](#)

REFERENCES

- BENDTSEN, J.D., NIELSEN, H., VON HEIJNE, G. & BRUNAK, S. (2004b). Improved prediction of signal peptides: SignalP 3.0. *J Mol Biol*, **340**, 783–795.
- 7, 8
- BHASIN, M. & RAGHAVA, G.P.S. (2004). ESLpred: SVM-based method for subcellular localization of eukaryotic proteins using dipeptide composition and PSI-BLAST. *Nucleic Acids Res*, **32**, W414–W419. 7, 10
- BHASIN, M., GARG, A. & RAGHAVA, G.P.S. (2005). PSLpred: prediction of subcellular localization of bacterial proteins. *Bioinformatics*, **21**, 2522–2524.
- 10
- BIOJAVA (2007). <http://www.biojava.org>. 27, 75, 148, 150
- BRAMEIER, M., KRINGS, A. & MACCALLUM, R.M. (2007). NucPred—predicting nuclear localization of proteins. *Bioinformatics*, **23**, 1159–1160. 110
- BURGARD, A.P., MOORE, G.L. & MARANAS, C.D. (2001). Review of the TEIRESIAS-based tools of the IBM Bioinformatics and Pattern Discovery Group. *Metab Eng*, **3**, 285–288. 26
- BURGES, C.J.C. (1998). A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, **2**, 121–167. 23, 24
- BURLEY, S.K. & PETSKO, G.A. (1985). Aromatic-aromatic interaction: a mechanism of protein structure stabilization. *Science*, **229**, 23–28. 132
- CARNINCI, P., KASUKAWA, T., KATAYAMA, S., GOUGH, J., FRITH, M.C., MAEDA, N., OYAMA, R., RAVASI, T., LENHARD, B., WELLS, C., KODZ-

REFERENCES

IUS, R., SHIMOKAWA, K., BAJIC, V.B., BRENNER, S.E., BATALOV, S., FORREST, A.R.R., ZAVOLAN, M., DAVIS, M.J., WILMING, L.G., AIDINIS, V., ALLEN, J.E., AMBESI-IMPIOMBATO, A., APWEILER, R., ATURALIYA, R.N., BAILEY, T.L., BANSAL, M., BAXTER, L., BEISEL, K.W., BERSANO, T., BONO, H., CHALK, A.M., CHIU, K.P., CHOUDHARY, V., CHRISTOFFELS, A., CLUTTERBUCK, D.R., CROWE, M.L., DALLA, E., DALRYMPLE, B.P., DE BONO, B., GATTA, G.D., DI BERNARDO, D., DOWN, T., ENGSTROM, P., FAGIOLINI, M., FAULKNER, G., FLETCHER, C.F., FUKUSHIMA, T., FURUNO, M., FUTAKI, S., GARIBOLDI, M., GEORGII-HEMMING, P., GINGERAS, T.R., GOJOBORI, T., GREEN, R.E., GUSTINCICH, S., HARBERS, M., HAYASHI, Y., HENSCH, T.K., HIROKAWA, N., HILL, D., HUMINIECKI, L., IACONO, M., IKEO, K., IWAMA, A., ISHIKAWA, T., JAKT, M., KANAPIN, A., KATOH, M., KAWASAWA, Y., KELSO, J., KITAMURA, H., KITANO, H., KOLLIAS, G., KRISHNAN, S.P.T., KRUGER, A., KUMMERFELD, S.K., KUROCHKIN, I.V., LAREAU, L.F., LAZAREVIC, D., LIPOVICH, L., LIU, J., LIUNI, S., MCWILLIAM, S., BABU, M.M., MADERA, M., MARCHIONNI, L., MATSUDA, H., MATSUZAWA, S., MIKI, H., MIGNONE, F., MIYAKE, S., MORRIS, K., MOTTAGUI-TABAR, S., MULDER, N., NAKANO, N., NAKAUCHI, H., NG, P., NILSSON, R., NISHIGUCHI, S., NISHIKAWA, S., NORI, F., OHARA, O., OKAZAKI, Y., ORLANDO, V., PANG, K.C., PAVAN, W.J., PAVESI, G., PESOLE, G., PETROVSKY, N., PIAZZA, S., REED, J., REID, J.F., RING, B.Z., RINGWALD, M., ROST, B., RUAN, Y., SALZBERG, S.L., SANDELIN, A., SCHNEIDER, C., SCHNBACH, C., SEKIGUCHI, K., SEMPLE, C.A.M., SENO, S., SESSA, L., SHENG, Y.,

REFERENCES

- SHIBATA, Y., SHIMADA, H., SHIMADA, K., SILVA, D., SINCLAIR, B., SPERLING, S., STUPKA, E., SUGIURA, K., SULTANA, R., TAKENAKA, Y., TAKI, K., TAMMOJA, K., TAN, S.L., TANG, S., TAYLOR, M.S., TEGNER, J., TEICHMANN, S.A., UEDA, H.R., VAN NIMWEGEN, E., VERARDO, R., WEI, C.L., YAGI, K., YAMANISHI, H., ZABAROVSKY, E., ZHU, S., ZIMMER, A., HIDE, W., BULT, C., GRIMMOND, S.M., TEASDALE, R.D., LIU, E.T., BRUSIC, V., QUACKENBUSH, J., WAHLESTEDT, C., MATTICK, J.S., HUME, D.A., KAI, C., SASAKI, D., TOMARU, Y., FUKUDA, S., KANAMORI-KATAYAMA, M., SUZUKI, M., AOKI, J., ARAKAWA, T., IIDA, J., IMAMURA, K., ITOH, M., KATO, T., KAWAJI, H., KAWAGASHIRA, N., KAWASHIMA, T., KOJIMA, M., KONDO, S., KONNO, H., NAKANO, K., NINOMIYA, N., NISHIO, T., OKADA, M., PLESSY, C., SHIBATA, K., SHIRAKI, T., SUZUKI, S., TAGAMI, M., WAKI, K., WATAHIKI, A., OKAMURA-OHO, Y., SUZUKI, H., KAWAI, J., HAYASHIZAKI, Y., CONSORTIUM, F.A.N.T.O.M., GROUP, R.I.K.E.N.G.E.R. & GROUP), G.S.G.G.N.P.C. (2005). The transcriptional landscape of the mammalian genome. *Science*, **309**, 1559–1563. [114](#)
- CHANG, C.C. & LIN, C.J. (2001). LIBSVM: a library for support vector machines. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>. [23](#), [24](#), [79](#), [108](#), [118](#)
- CHOTHIA, C. & LESK, A.M. (1986). The relation between the divergence of sequence and structure in proteins. *EMBO J*, **5**, 823–826. [13](#)
- CLAROS, M.G. & VON HEIJNE, G. (1994). TopPred II: an improved software for membrane protein structure predictions. *Comput Appl Biosci*, **10**, 685–686.

REFERENCES

- 68, 137
- COEYTAUX, K. & POUPEON, A. (2005). Prediction of unfolded segments in a protein sequence based on amino acid composition. *Bioinformatics*, **21**, 1891–1900. 112
- COKOL, M., NAIR, R. & ROST, B. (2000). Finding nuclear localization signals. *EMBO Rep.*, **1**, 411–415. 7, 8, 110
- COMON, P. (1994). Independent component analysis, a new concept? *Signal Process.*, **36**, 287–314. 19
- COZZETTO, D., MATTEO, A.D. & TRAMONTANO, A. (2005). Ten years of predictions ... and counting. *FEBS J*, **272**, 881–882. 112, 159
- COZZETTO, D., GIORGETTI, A., RAIMONDO, D. & TRAMONTANO, A. (2007). The evaluation of protein structure prediction results. *Mol Biotechnol.* 112, 159
- DANG, C.V. & LEE, W.M. (1989). Nuclear and nucleolar targeting sequences of c-erb-a, c-myb, N-myc, p53, HSP70, and HIV tat proteins. *J Biol Chem*, **264**, 18019–18023. 108, 120
- DEMPSTER, A., LAIRD, N. & RUBIN, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *J Royal Statistical Society, Series B*, **39**. 18
- DINGWALL, C. & LASKEY, R.A. (1991). Nuclear targeting sequences— a consensus? *Trends Biochem Sci*, **16**, 478–481. 73

REFERENCES

- DOĞRUCL, M., DOWN, T. & HUBBARD, T. (2008). NestedMICA as an ab initio protein motif discovery tool. *BMC Bioinformatics*, **9**, 19. [3](#), [25](#), [64](#), [69](#), [115](#)
- DOWN, T.A. & HUBBARD, T.J.P. (2002). Computational detection and location of transcription start sites in mammalian genomic DNA. *Genome Res*, **12**, 458–461. [3](#), [27](#), [64](#), [73](#), [74](#), [76](#)
- DOWN, T.A. & HUBBARD, T.J.P. (2004). What can we learn from noncoding regions of similarity between genomes? *BMC Bioinformatics*, **5**, 131. [74](#)
- DOWN, T.A. & HUBBARD, T.J.P. (2005). NestedMICA: sensitive inference of over-represented motifs in nucleic acid sequence. *Nucleic Acids Res*, **33**, 1445–1453. [2](#), [27](#), [30](#), [32](#), [33](#), [34](#), [36](#), [64](#)
- DUNKER, A.K., OBRADOVIC, Z., ROMERO, P., GARNER, E.C. & BROWN, C.J. (2000). Intrinsic protein disorder in complete genomes. *Genome Inform Ser Workshop Genome Inform*, **11**, 161–171. [111](#), [168](#)
- DURBIN, R., EDDY, S.R., KROGH, A. & MITCHISON, G. (1999). *Biological Sequence Analysis : Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press. [16](#), [18](#), [30](#)
- ECK, R.V. & DAYHOFF, M.O. (1966). *Atlas of Protein Sequence and Structure*, vol. 3. National Biomedical Research Foundation, Silver Spring, Maryland. [68](#)
- EDWARDS, R.J., DAVEY, N.E. & SHIELDS, D.C. (2007). SLiMFinder: a probabilistic method for identifying over-represented, convergently evolved, short linear motifs in proteins. *PLoS ONE*, **2**, e967. [26](#)

REFERENCES

- EMANUELSSON, O., NIELSEN, H. & VON HEIJNE, G. (1999). ChloroP, a neural network-based method for predicting chloroplast transit peptides and their cleavage sites. *Protein Sci*, **8**, 978–984. [7](#), [8](#), [66](#), [85](#), [136](#)
- EMANUELSSON, O., NIELSEN, H., BRUNAK, S. & VON HEIJNE, G. (2000). Predicting subcellular localization of proteins based on their N-terminal amino acid sequence. *J Mol Biol*, **300**, 1005–1016. [7](#), [40](#), [56](#), [59](#)
- EMANUELSSON, O., BRUNAK, S., VON HEIJNE, G. & NIELSEN, H. (2007). Locating proteins in the cell using TargetP, SignalP and related tools. *Nat Protoc*, **2**, 953–971. [8](#)
- ENDRES, M., NEUPERT, W. & BRUNNER, M. (1999). Transport of the ADP/ATP carrier of mitochondria from the TOM complex to the TIM22.54 complex. *EMBO J*, **18**, 3214–3221. [65](#)
- FAVOROV, A.V., GELFAND, M.S., GERASIMOVA, A.V., RAVCHEEV, D.A., MIRONOV, A.A. & MAKEEV, V.J. (2005). A Gibbs sampler for identification of symmetrically structured, spaced DNA motifs with improved estimation of the signal length. *Bioinformatics*, **21**, 2240–2245. [27](#)
- FINK, J.L., ATURALIYA, R.N., DAVIS, M.J., ZHANG, F., HANSON, K., TEASDALE, M.S., KAI, C., KAWAI, J., CARNINCI, P., HAYASHIZAKI, Y. & TEASDALE, R.D. (2006). LOCATE: a mouse protein subcellular localization database. *Nucleic Acids Res*, **34**, D213–D217. [12](#), [114](#)

REFERENCES

- GAO, Y. & MEHTA, K. (2007). N-linked glycosylation of CD38 is required for its structure stabilization but not for membrane localization. *Mol Cell Biochem*, **295**, 1–7. [67](#)
- GARDY, J.L., SPENCER, C., WANG, K., ESTER, M., TUSNÁDY, G.E., SIMON, I., HUA, S., DEFAYS, K., LAMBERT, C., NAKAI, K. & BRINKMAN, F.S.L. (2003). PSORT-B: Improving protein subcellular localization prediction for Gram-negative bacteria. *Nucleic Acids Res*, **31**, 3613–3617. [9](#)
- GARDY, J.L., LAIRD, M.R., CHEN, F., REY, S., WALSH, C.J., ESTER, M. & BRINKMAN, F.S.L. (2005). PSORTb v.2.0: expanded prediction of bacterial protein subcellular localization and insights gained from comparative proteome analysis. *Bioinformatics*, **21**, 617–623. [9](#), [13](#)
- GARNER, ROMERO, DUNKER, BROWN & OBRADOVIC (1999). Predicting binding regions within disordered proteins. *Genome Informatics*, **10**, 41–50. [112](#)
- GOLDFARB, D.S., GARIPY, J., SCHOOLNIK, G. & KORNBERG, R.D. (1986). Synthetic peptides as nuclear localization signals. *Nature*, **322**, 641–644. [108](#), [116](#)
- GOULD, S.G., KELLER, G.A. & SUBRAMANI, S. (1987). Identification of a peroxisomal targeting signal at the carboxy terminus of firefly luciferase. *J Cell Biol*, **105**, 2923–2931. [65](#)
- GOULD, S.J., KELLER, G.A., HOSKEN, N., WILKINSON, J. & SUBRAMANI, S. (1989). A conserved tripeptide sorts proteins to peroxisomes. *J Cell Biol*, **108**, 1657–1664. [65](#)

REFERENCES

- GRIBSKOV, M., MCLACHLAN, A.D. & EISENBERG, D. (1987). Profile analysis: detection of distantly related proteins. [17](#)
- GUAN, J.L., MACHAMER, C.E. & ROSE, J.K. (1985). Glycosylation allows cell-surface transport of an anchored secretory protein. *Cell*, **42**, 489–496. [67](#)
- GUDA, C. & SUBRAMANIAM, S. (2005). pTARGET [CORRECTED] a new method for predicting protein subcellular localization in eukaryotes. *Bioinformatics*, **21**, 3963–3969. [7](#), [10](#), [13](#), [36](#), [70](#), [71](#)
- HANNINK, M. & DONOGHUE, D.J. (1986). Cell surface expression of membrane-anchored v-sis gene products: glycosylation is not required for cell surface transport. *J Cell Biol*, **103**, 2311–2322. [67](#)
- HENIKOFF, S. & HENIKOFF, J.G. (1992). Amino acid substitution matrices from protein blocks. *Proc Natl Acad Sci U S A*, **89**, 10915–10919. [68](#)
- HERTZ, G.Z. & STORMO, G.D. (1999). Identifying DNA and protein patterns with statistically significant alignments of multiple sequences. *Bioinformatics*, **15**, 563–577. [26](#)
- HINSBY, A.M., KIEMER, L., KARLBERG, E.O., LAGE, K., FAUSBLL, A., JUNCKER, A.S., ANDERSEN, J.S., MANN, M. & BRUNAK, S. (2006). A wiring of the human nucleolus. *Mol Cell*, **22**, 285–295. [114](#)
- HIROKAWA, T., BOON-CHIENG, S. & MITAKU, S. (1998). SOSUI: classification and secondary structure prediction system for membrane proteins. *Bioinformatics*, **14**, 378–379. [68](#), [137](#)

REFERENCES

- HOBOM, U., SCHARF, M., SCHNEIDER, R. & SANDER, C. (1992). Selection of representative protein data sets. *Protein Sci*, **1**, 409–417. [13](#), [15](#), [41](#), [56](#)
- HÖGLUND, A., DÖNNES, P., BLUM, T., ADOLPH, H.W. & KOHLBACHER, O. (2006). MultiLoc: prediction of protein subcellular localization using N-terminal targeting sequences, sequence motifs and amino acid composition. *Bioinformatics*, **22**, 1158–1165. [7](#), [10](#), [11](#), [13](#), [41](#), [70](#), [71](#), [97](#), [98](#), [99](#), [180](#)
- HORTON, P. & NAKAI, K. (1997). Better prediction of protein cellular localization sites with the k nearest neighbors classifier. *Proc Int Conf Intell Syst Mol Biol*, **5**, 147–152. [7](#), [9](#)
- HORTON, P., PARK, K.J., OBAYASHI, T., FUJITA, N., HARADA, H., ADAMS-COLLIER, C.J. & NAKAI, K. (2007). WoLF PSORT: protein localization predictor. *Nucleic Acids Res*, **35**, W585–W587. [7](#), [9](#)
- HSU, C.W. & LIN, C.J. (2002). A simple decomposition method for support vector machines. *Machine Learning*, **46**, 291–314. [24](#), [79](#)
- HUA, S. & SUN, Z. (2001). Support vector machine approach for protein subcellular localization prediction. *Bioinformatics*, **17**, 721–728. [7](#), [10](#)
- HUBER, R. (1979). Conformational flexibility in protein molecules. *Nature*, **280**, 538–539. [111](#)
- HUGHES, J.D., ESTEP, P.W., TAVAZOIE, S. & CHURCH, G.M. (2000). Computational identification of cis-regulatory elements associated with groups of

REFERENCES

- functionally related genes in *Saccharomyces cerevisiae*. *J Mol Biol*, **296**, 1205–1214. [27](#)
- HULO, N., BAIROCH, A., BULLIARD, V., CERUTTI, L., CASTRO, E.D., LANGENDIJK-GENEVAUX, P.S., PAGNI, M. & SIGRIST, C.J.A. (2006). The PROSITE database. *Nucleic Acids Res*, **34**, D227–D230. [11](#), [25](#), [40](#)
- HWANG, S., GOU, Z. & KUZNETSOV, I.B. (2007). DP-Bind: a web server for sequence-based prediction of DNA-binding residues in DNA-binding proteins. *Bioinformatics*, **23**, 634–636. [26](#)
- JIN, Y. & DUNBRACK, R.L. (2005). Assessment of disorder predictions in CASP6. *Proteins*, **61 Suppl 7**, 167–175. [113](#)
- JOACHIMS, T. (1999). Making large-scale support vector machine learning practical. 169–184, MIT Press. [23](#), [79](#)
- JONES, D.T. (2007). Improving the accuracy of transmembrane protein topology prediction using evolutionary information. *Bioinformatics*, **23**, 538–544. [137](#), [159](#)
- JONES, D.T., TAYLOR, W.R. & THORNTON, J.M. (1994). A model recognition approach to the prediction of all-helical membrane protein structure and topology. *Biochemistry*, **33**, 3038–3049. [142](#)
- JUNCKER, A.S., WILLENBROCK, H., HEIJNE, G.V., BRUNAK, S., NIELSEN, H. & KROGH, A. (2003). Prediction of lipoprotein signal peptides in Gram-negative bacteria. *Protein Sci*, **12**, 1652–1662. [7](#), [8](#)

REFERENCES

- KÄLL, L., KROGH, A. & SONNHAMMER, E.L.L. (2004). A combined transmembrane topology and signal peptide prediction method. *J Mol Biol*, **338**, 1027–1036. [137](#), [139](#), [159](#), [160](#)
- KÄLL, L., KROGH, A. & SONNHAMMER, E.L.L. (2007). Advantages of combined transmembrane topology and signal peptide prediction—the Phobius web server. *Nucleic Acids Res*, **35**, W429–W432. [138](#), [164](#)
- KAPLAN, H.A., WELPLY, J.K. & LENNARZ, W.J. (1987). Oligosaccharyl transferase: the central enzyme in the pathway of glycoprotein assembly. *Biochim Biophys Acta*, **906**, 161–173. [67](#)
- KELLEY, L.P. & KINSELLA, B.T. (2003). The role of N-linked glycosylation in determining the surface expression, G protein interaction and effector coupling of the alpha isoform of the human thromboxane A(2) receptor. *Biochim Biophys Acta*, **1621**, 192–203. [67](#)
- KIEMER, L., BENDTSEN, J.D. & BLOM, N. (2005). NetAcet: prediction of N-terminal acetylation sites. *Bioinformatics*, **21**, 1269–1270. [38](#)
- KISSINGER, C.R., PARGE, H.E., KNIGHTON, D.R., LEWIS, C.T., PELLETIER, L.A., TEMPCZYK, A., KALISH, V.J., TUCKER, K.D., SHOWALTER, R.E. & MOOMAW, E.W. (1995). Crystal structures of human calcineurin and the human FKBP12-FK506-calcineurin complex. *Nature*, **378**, 641–644. [132](#)
- KLEIN, P., KANEHISA, M. & DELISI, C. (1984). Prediction of protein function from sequence properties. Discriminant analysis of a data base. *Biochim Biophys Acta*, **787**, 221–226. [68](#)

REFERENCES

- KROGH, A. (2002). www.binf.ku.dk/~krogh/docs/labeled_fasta_format.html.
[148](#)
- KROGH, A., LARSSON, B., VON HEIJNE, G. & SONNHAMMER, E.L. (2001). Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *J Mol Biol*, **305**, 567–580. [6](#), [68](#), [79](#), [123](#), [124](#), [134](#), [137](#), [138](#), [140](#)
- KUZNETSOV, I.B., GOU, Z., LI, R. & HWANG, S. (2006). Using evolutionary and structural information to predict DNA-binding sites on DNA-binding proteins. *Proteins*, **64**, 19–27. [27](#)
- LA COUR, T., KIEMER, L., MØLGAARD, A., GUPTA, R., SKRIVER, K. & BRUNAK, S. (2004). Analysis and prediction of leucine-rich nuclear export signals. *Protein Eng Des Sel*, **17**, 527–536. [110](#)
- LAO, D.M., OKUNO, T. & SHIMIZU, T. (2002). Evaluating transmembrane topology prediction methods for the effect of signal peptide in topology prediction. *In Silico Biol*, **2**, 485–494. [69](#), [137](#)
- LEUNG, A.K.L., ANDERSEN, J.S., MANN, M. & LAMOND, A.I. (2003). Bioinformatic analysis of the nucleolus. *Biochem J*, **376**, 553–569. [111](#), [125](#)
- LI, ROMERO, RANI, DUNKER & OBRADOVIC (1999). Predicting protein disorder for N-, C-, and internal regions. *Genome Informatics*, **10**, 30–40. [112](#)

REFERENCES

- LI, H. & BINGHAM, P.M. (1991). Arginine/serine-rich domains of the su(wa) and tra RNA processing regulators target proteins to a subnuclear compartment implicated in splicing. *Cell*, **67**, 335–342. [108](#)
- LI, W. & GODZIK, A. (2006). CD-HIT: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics*, **22**, 1658–1659. [14](#), [15](#), [36](#), [70](#), [71](#), [115](#)
- LI, W., JAROSZEWSKI, L. & GODZIK, A. (2001). Clustering of highly homologous sequences to reduce the size of large protein databases. *Bioinformatics*, **17**, 282–283. [14](#)
- LI, W., JAROSZEWSKI, L. & GODZIK, A. (2002). Tolerating some redundancy significantly speeds up clustering of large protein databases. *Bioinformatics*, **18**, 77–82. [14](#), [15](#)
- LINDING, R., JENSEN, L.J., DIELLA, F., BORK, P., GIBSON, T.J. & RUSSELL, R.B. (2003a). Protein disorder prediction: implications for structural proteomics. *Structure*, **11**, 1453–1459. [112](#)
- LINDING, R., RUSSELL, R.B., NEDUVA, V. & GIBSON, T.J. (2003b). GlobPlot: Exploring protein sequences for globularity and disorder. *Nucleic Acids Res*, **31**, 3701–3708. [112](#)
- LOBLEY, A., SWINDELLS, M.B., ORENGO, C.A. & JONES, D.T. (2007). Inferring function using patterns of native disorder in proteins. *PLoS Comput Biol*, **3**, e162. [102](#), [112](#), [168](#)

REFERENCES

- LU, Z., SZAFRON, D., GREINER, R., LU, P., WISHART, D.S., POULIN, B., ANVIK, J., MACDONELL, C. & EISNER, R. (2004). Predicting subcellular localization of proteins using machine-learned classifiers. *Bioinformatics*, **20**, 547–556. [7](#), [10](#)
- LUMB, M.J., PURDUE, P.E. & DANPURE, C.J. (1994). Molecular evolution of alanine/glyoxylate aminotransferase 1 intracellular targeting. Analysis of the feline gene. *Eur J Biochem*, **221**, 53–62. [104](#)
- MACHAMER, C.E., FLORKIEWICZ, R.Z. & ROSE, J.K. (1985). A single N-linked oligosaccharide at either of the two normal sites is sufficient for transport of vesicular stomatitis virus G protein to the cell surface. *Mol Cell Biol*, **5**, 3074–3083. [67](#)
- MAEDA, N., KASUKAWA, T., OYAMA, R., GOUGH, J., FRITH, M., ENGSTRM, P.G., LENHARD, B., ATURALIYA, R.N., BATALOV, S., BEISEL, K.W., BULT, C.J., FLETCHER, C.F., FORREST, A.R.R., FURUNO, M., HILL, D., ITOH, M., KANAMORI-KATAYAMA, M., KATAYAMA, S., KATOH, M., KAWASHIMA, T., QUACKENBUSH, J., RAVASI, T., RING, B.Z., SHIBATA, K., SUGIURA, K., TAKENAKA, Y., TEASDALE, R.D., WELLS, C.A., ZHU, Y., KAI, C., KAWAI, J., HUME, D.A., CARNINCI, P. & HAYASHIZAKI, Y. (2006). Transcript annotation in FANTOM3: mouse gene catalog based on physical cDNAs. *PLoS Genet*, **2**, e62. [114](#)
- MAEDA, Y., HISATAKE, K., KONDO, T., HANADA, K., SONG, C.Z., NISHIMURA, T. & MURAMATSU, M. (1992). Mouse rRNA gene transcription

REFERENCES

- factor mUBF requires both HMG-box1 and an acidic tail for nucleolar accumulation: molecular analysis of the nucleolar targeting mechanism. *EMBO J*, **11**, 3695–3704. [108](#)
- MARTELLI, P.L., FARISELLI, P. & CASADIO, R. (2003). An ENSEMBLE machine learning approach for the prediction of all-alpha membrane proteins. *Bioinformatics*, **19 Suppl 1**, i205–i211. [137](#)
- MATSUDA, K., ZHENG, J., DU, G.G., KLCKER, N., MADISON, L.D. & DALLOS, P. (2004). N-linked glycosylation sites of the motor protein prestin: effects on membrane targeting and electrophysiological function. *J Neurochem*, **89**, 928–938. [67](#)
- MATTHEWS, B.W. (1975). Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochim Biophys Acta*, **405**, 442–451. [38](#), [47](#), [80](#)
- MILSTEIN, C., BROWNLEE, G.G., HARRISON, T.M. & MATTHEWS, M.B. (1972). A possible precursor of immunoglobulin light chains. *Nat New Biol*, **239**, 117–120. [66](#)
- MOHRMANN, K., VAN EIJDHOVEN, M.A.J., SCHINKEL, A.H. & SCHELLENS, J.H.M. (2005). Absence of N-linked glycosylation does not affect plasma membrane localization of breast cancer resistance protein (BCRP/ABCG2). *Cancer Chemother Pharmacol*, **56**, 344–350. [67](#)
- MOTLEY, A., LUMB, M.J., OATEY, P.B., JENNINGS, P.R., ZOYSA, P.A.D., WANDERS, R.J., TABAK, H.F. & DANPURE, C.J. (1995). Mammalian ala-

REFERENCES

- nine/glyoxylate aminotransferase 1 is imported into peroxisomes via the PTS1 translocation pathway. Increased degeneracy and context specificity of the mammalian PTS1 motif and implications for the peroxisome-to-mitochondrion mistargeting of AGT in primary hyperoxaluria type 1. *J Cell Biol*, **131**, 95–109. [104](#)
- MÜLLER, S., CRONING, M.D. & APWEILER, R. (2001). Evaluation of methods for the prediction of membrane spanning regions. *Bioinformatics*, **17**, 646–653. [69](#), [137](#)
- MUNRO, S. (1995). An investigation of the role of transmembrane domains in Golgi protein retention. *EMBO J*, **14**, 4695–4704. [95](#)
- NAIR, R. & ROST, B. (2002a). Inferring sub-cellular localization through automated lexical analysis. *Bioinformatics*, **18 Suppl 1**, S78–S86. [7](#), [9](#)
- NAIR, R. & ROST, B. (2002b). Sequence conserved for subcellular localization. *Protein Sci*, **11**, 2836–2847. [7](#), [9](#)
- NAIR, R. & ROST, B. (2004). LOCnet and LOCtarget: sub-cellular localization for structural genomics targets. *Nucleic Acids Res*, **32**, W517–W521. [9](#), [110](#)
- NAIR, R., CARTER, P. & ROST, B. (2003). NLSdb: database of nuclear localization signals. *Nucleic Acids Res*, **31**, 397–399. [8](#), [11](#), [110](#)
- NAKAI, K. & HORTON, P. (1999). PSORT: a program for detecting sorting signals in proteins and predicting their subcellular localization. *Trends Biochem Sci*, **24**, 34–36. [7](#), [9](#), [110](#)

REFERENCES

- NAKAI, K. & KANEHISA, M. (1991). Expert system for predicting protein localization sites in gram-negative bacteria. *Proteins*, **11**, 95–110. [7](#), [9](#), [68](#)
- NAKAI, K. & KANEHISA, M. (1992). A knowledge base for predicting protein localization sites in eukaryotic cells. *Genomics*, **14**, 897–911. [9](#)
- NEDUVA, V. & RUSSELL, R.B. (2006). DILIMOT: discovery of linear motifs in proteins. *Nucleic Acids Res*, **34**, W350–W355. [26](#), [58](#)
- NG, P., NAGARAJAN, N., JONES, N. & KEICH, U. (2006). Apples to apples: improving the performance of motif finders and their significance analysis in the twilight zone. *Bioinformatics*, **22**, e393–e401. [34](#)
- NICKEL, W. (2003). The mystery of nonclassical protein secretion. A current view on cargo proteins and potential export routes. *Eur J Biochem*, **270**, 2109–2119. [66](#)
- NIELSEN, H. & KROGH, A. (1998). Prediction of signal peptides and signal anchors by a hidden Markov model. *Proc Int Conf Intell Syst Mol Biol*, **6**, 122–130. [160](#)
- NIELSEN, H., ENGELBRECHT, J., BRUNAK, S. & von HEIJNE, G. (1997a). Identification of prokaryotic and eukaryotic signal peptides and prediction of their cleavage sites. *Protein Eng*, **10**, 1–6. [160](#)
- NIELSEN, H., ENGELBRECHT, J., BRUNAK, S. & von HEIJNE, G. (1997b). A neural network method for identification of prokaryotic and eukaryotic signal

REFERENCES

- peptides and prediction of their cleavage sites. *Int J Neural Syst*, **8**, 581–599. [8](#), [138](#)
- NIELSEN, H., BRUNAK, S. & VON HEIJNE, G. (1999). Machine learning approaches for the prediction of signal peptides and other protein sorting signals. *Protein Eng*, **12**, 3–9. [7](#)
- OBENAUER, J.C., CANTLEY, L.C. & YAFFE, M.B. (2003). Scansite 2.0: Proteome-wide prediction of cell signaling interactions using short sequence motifs. *Nucleic Acids Res*, **31**, 3635–3641. [25](#)
- ODA, T., MIYAJIMA, H., SUZUKI, Y. & ICHIYAMA, A. (1987). Nucleotide sequence of the cDNA encoding the precursor for mitochondrial serine:pyruvate aminotransferase of rat liver. *Eur J Biochem*, **168**, 537–542. [104](#)
- Osumi, T., TSUKAMOTO, T., HATA, S., YOKOTA, S., MIURA, S., FUJIKI, Y., HIJIKATA, M., MIYAZAWA, S. & HASHIMOTO, T. (1991). Amino-terminal presequence of the precursor of peroxisomal 3-ketoacyl-coa thiolase is a cleavable signal peptide for peroxisomal targeting. *Biochem Biophys Res Commun*, **181**, 947–954. [65](#)
- PARK, K.J. & KANEHISA, M. (2003). Prediction of protein subcellular locations by support vector machines using compositions of amino acids and amino acid pairs. *Bioinformatics*, **19**, 1656–1663. [7](#), [10](#), [13](#)
- PAVESI, G., MEREGHETTI, P., MAURI, G. & PESOLE, G. (2004). Weeder Web: discovery of transcription factor binding sites in a set of sequences from co-regulated genes. *Nucleic Acids Res*, **32**, W199–W203. [27](#)

REFERENCES

- PELHAM, H.R. (1995). Sorting and retrieval between the endoplasmic reticulum and Golgi apparatus. *Curr Opin Cell Biol*, **7**, 530–535. [66](#), [88](#)
- PIERLEONI, A., MARTELLI, P.L., FARISELLI, P. & CASADIO, R. (2006). Ba-CelLo: a balanced subcellular localization predictor. *Bioinformatics*, **22**, e408–e416. [7](#), [10](#)
- PRILOUSKY, J., FELDER, C.E., ZEEV-BEN-MORDEHAI, T., RYDBERG, E.H., MAN, O., BECKMANN, J.S., SILMAN, I. & SUSSMAN, J.L. (2005). FoldIndex: a simple tool to predict whether a given protein sequence is intrinsically unfolded. *Bioinformatics*, **21**, 3435–3438. [112](#)
- PUNTERVOLL, P., LINDING, R., GEMND, C., CHABANIS-DAVIDSON, S., MATTINGSDAL, M., CAMERON, S., MARTIN, D.M.A., AUSIELLO, G., BRANNETTI, B., COSTANTINI, A., FERR, F., MASSELLI, V., VIA, A., CESARENI, G., DIELLA, F., SUPERTI-FURGA, G., WYRWICZ, L., RAMU, C., MCGUIGAN, C., GUDAVALLI, R., LETUNIC, I., BORK, P., RYCHLEWSKI, L., KSTER, B., HELMER-CITTERICH, M., HUNTER, W.N., AASLAND, R. & GIBSON, T.J. (2003). ELM server: A new resource for investigating short functional sites in modular eukaryotic proteins. *Nucleic Acids Res*, **31**, 3625–3630. [25](#)
- PURDUE, P.E., LUMB, M.J. & DANPURE, C.J. (1992). Molecular evolution of alanine/glyoxylate aminotransferase 1 intracellular targeting. Analysis of the marmoset and rabbit genes. *Eur J Biochem*, **207**, 757–766. [104](#)

REFERENCES

- QUEVILLON, E., SILVENTOINEN, V., PILLAI, S., HARTE, N., MULDER, N., APWEILER, R. & LOPEZ, R. (2005). InterProScan: protein domains identifier. *Nucleic Acids Res*, **33**, W116–W120. [58](#)
- RABINER, L.R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, **77**, 257–286. [16](#)
- RADIVOJAC, P., OBRADOVIC, Z., BROWN, C.J. & DUNKER, A.K. (2003). Prediction of boundaries between intrinsically ordered and disordered protein regions. *Pac Symp Biocomput*, 216–227. [112](#)
- RADIVOJAC, P., OBRADOVIC, Z., SMITH, D.K., ZHU, G., VUCETIC, S., BROWN, C.J., LAWSON, J.D. & DUNKER, A.K. (2004). Protein flexibility and intrinsic disorder. *Protein Sci*, **13**, 71–80. [112](#)
- RAMADASS, A.S. (2005). *Computational detection of regulatory signals in human genome sequence*. Ph.D. thesis, University of Cambridge. [74](#), [75](#)
- REINHARDT, A. & HUBBARD, T. (1998). Using neural networks for prediction of the subcellular location of proteins. *Nucleic Acids Res*, **26**, 2230–2236. [68](#)
- RICHARDSON, J.S. & RICHARDSON, D.C. (1988). Amino acid preferences for specific locations at the ends of alpha helices. *Science*, **240**, 1648–1652. [142](#)
- RIGOUTSOS, I. & FLORATOS, A. (1998). Combinatorial pattern discovery in biological sequences: The TEIRESIAS algorithm. *Bioinformatics*, **14**, 55–67. [26](#), [58](#)

REFERENCES

- ROBBINS, J., DILWORTH, S.M., LASKEY, R.A. & DINGWALL, C. (1991). Two interdependent basic domains in nucleoplasmin nuclear targeting sequence: identification of a class of bipartite nuclear targeting sequence. *Cell*, **64**, 615–623. [73](#)
- ROMERO, P., OBRADOVIC, Z. & DUNKER, A.K. (2004). Natively disordered proteins: functions and predictions. *Appl Bioinformatics*, **3**, 105–113. [112](#)
- SANDER, C. & SCHNEIDER, R. (1991). Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins*, **9**, 56–68. [14](#)
- SCHERL, A., COUT, Y., DON, C., CALL, A., KINDBEITER, K., SANCHEZ, J.C., GRECO, A., HOCHSTRASSER, D. & DIAZ, J.J. (2002). Functional proteomic analysis of human nucleolus. *Mol Biol Cell*, **13**, 4100–4109. [113](#)
- SCHREIBER, V., MOLINETE, M., BOEUF, H., DE MURCIA, G. & DE MURCIA, J.M. (1992). The human poly(ADP-ribose) polymerase nuclear localization signal is a bipartite element functionally separate from DNA binding and catalytic activity. *EMBO J*, **11**, 3263–3269. [73](#)
- SCHULTZ, J., MILPETZ, F., BORK, P. & PONTING, C.P. (1998). SMART, a simple modular architecture research tool: identification of signaling domains. *Proc Natl Acad Sci U S A*, **95**, 5857–5864. [58](#)
- SINHA, S. & TOMPA, M. (2003). YMF: A program for discovery of novel transcription factor binding sites by statistical overrepresentation. *Nucleic Acids Res*, **31**, 3586–3588. [27](#)

REFERENCES

- SKILLING, J. (2004). Nested Sampling. In R. Fischer, R. Preuss & U.V. Toussaint, eds., *American Institute of Physics Conference Series*, 395–405. [20](#), [27](#), [28](#), [64](#), [69](#)
- SMITH, G.B. (1987). Preface to S. Geman and D. Geman, “Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images”. 562–563, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA. [18](#)
- SMOLA, A. & SCHOLKOPF, B. (1998). A tutorial on support vector regression. Tech. rep. [23](#)
- SORO, S. & TRAMONTANO, A. (2005). The prediction of protein function at CASP6. *Proteins*, **61 Suppl 7**, 201–213. [112](#), [159](#)
- SPRENGER, J., FINK, J.L. & TEASDALE, R.D. (2006). Evaluation and comparison of mammalian subcellular localization prediction methods. *BMC Bioinformatics*, **7 Suppl 5**, S3. [11](#)
- SWINKELS, B.W., GOULD, S.J., BODNAR, A.G., RACHUBINSKI, R.A. & SUBRAMANI, S. (1991). A novel, cleavable peroxisomal targeting signal at the amino-terminus of the rat 3-ketoacyl-coa thiolase. *EMBO J*, **10**, 3255–3262. [65](#)
- SZAFRON, D., LU, P., GREINER, R., WISHART, D.S., POULIN, B., EISNER, R., LU, Z., ANVIK, J., MACDONELL, C., FYSHE, A. & MEEUWIS, D. (2004). Proteome Analyst: custom predictions with explanations in a web-based tool for high-throughput proteome annotations. *Nucleic Acids Res*, **32**, W365–W371. [13](#)

REFERENCES

- TAKADA, Y., KANEKO, N., ESUMI, H., PURDUE, P.E. & DANPURE, C.J. (1990). Human peroxisomal L-alanine: glyoxylate aminotransferase. Evolutionary loss of a mitochondrial targeting signal by point mutation of the initiation codon. *Biochem J*, **268**, 517–520. [104](#)
- THOMAS, P.D. & DILL, K.A. (1996). An iterative method for extracting energy-like quantities from protein structures. *Proc Natl Acad Sci U S A*, **93**, 11628–11633. [77](#), [78](#), [95](#), [97](#)
- TIPPING, M.E. (2001). Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, **1**, 211–244. [75](#)
- TOMPA, M., LI, N., BAILEY, T.L., CHURCH, G.M., MOOR, B.D., ESKIN, E., FAVOROV, A.V., FRITH, M.C., FU, Y., KENT, W.J., MAKEEV, V.J., MIRONOV, A.A., NOBLE, W.S., PAVESI, G., PESOLE, G., RGNIER, M., SIMONIS, N., SINHA, S., THIJS, G., VAN HELDEN, J., VANDENBOGAERT, M., WENG, Z., WORKMAN, C., YE, C. & ZHU, Z. (2005). Assessing computational tools for the discovery of transcription factor binding sites. *Nat Biotechnol*, **23**, 137–144. [27](#)
- TUSNÁDY, G.E. & SIMON, I. (2001). The HMMTOP transmembrane topology prediction server. *Bioinformatics*, **17**, 849–850. [68](#), [137](#)
- VALENCIA, A. (2005). Protein refinement: a new challenge for CASP in its 10th anniversary. *Bioinformatics*, **21**, 277. [112](#), [159](#)
- VAPNIK, V. & LERNER, A. (1963). Pattern recognition using generalized portrait method. *Automation and Remote Control*, **24**, 774–780, 1963. [22](#)

REFERENCES

- VIKLUND, H. & ELOFSSON, A. (2004). Best alpha-helical transmembrane protein topology predictions are achieved using hidden Markov models and evolutionary information. *Protein Sci*, **13**, 1908–1917. [137](#)
- VON HEIJNE, G. (1986). A new method for predicting signal sequence cleavage sites. *Nucleic Acids Res*, **14**, 4683–4690. [82](#), [137](#)
- VON HEIJNE, G. (1990). The signal peptide. *J Membr Biol*, **115**, 195–201. [66](#), [136](#)
- WARD, J.J., SODHI, J.S., MCGUFFIN, L.J., BUXTON, B.F. & JONES, D.T. (2004). Prediction and functional analysis of native disorder in proteins from the three kingdoms of life. *J Mol Biol*, **337**, 635–645. [112](#)
- WEINREB, P.H., ZHEN, W., POON, A.W., CONWAY, K.A. & LANSBURY, P.T. (1996). NACP, a protein implicated in Alzheimer’s disease and learning, is natively unfolded. *Biochemistry*, **35**, 13709–13715. [111](#)
- WIEDEMANN, N., PFANNER, N. & RYAN, M.T. (2001). The three modules of ADP/ATP carrier cooperate in receptor recruitment and translocation into mitochondria. *EMBO J*, **20**, 951–960. [65](#)
- WORKMAN, C.T. & STORMO, G.D. (2000). ANN-Spec: a method for discovering transcription factor binding sites with improved specificity. *Pac Symp Biocomput*, 467–478. [27](#)

REFERENCES

- WRIGHT, P.E. & DYSON, H.J. (1999). Intrinsically unstructured proteins: re-assessing the protein structure-function paradigm. *J Mol Biol*, **293**, 321–331.
[112](#), [168](#)
- WU, L.F., CHANAL, A. & RODRIGUE, A. (2000). Membrane targeting and translocation of bacterial hydrogenases. *Arch Microbiol*, **173**, 319–324. [103](#)
- YAN, K., KHOSHNOODI, J., RUOTSALAINEN, V. & TRYGGVASON, K. (2002). N-linked glycosylation is critical for the plasma membrane localization of nephrin. *J Am Soc Nephrol*, **13**, 1385–1389. [67](#)
- YANG, Z.R., THOMSON, R., MCNEIL, P. & ESNOUF, R.M. (2005). RONN: the bio-basis function neural network technique applied to the detection of natively disordered regions in proteins. *Bioinformatics*, **21**, 3369–3376. [5](#), [102](#),
[112](#), [114](#), [118](#)
- YU, C.S., LIN, C.J. & HWANG, J.K. (2004). Predicting subcellular localization of proteins for gram-negative bacteria by support vector machines based on n-peptide compositions. *Protein Sci*, **13**, 1402–1406. [7](#), [10](#)
- YUAN, Z. & TEASDALE, R.D. (2002). Prediction of Golgi Type II membrane proteins based on their transmembrane domains. *Bioinformatics*, **18**, 1109–1115. [66](#)
- ZASLAVSKY, E. & SINGH, M. (2006). A combinatorial optimization approach for diverse motif finding applications. *Algorithms Mol Biol*, **1**, 13. [26](#)