This is a repository copy of *Multi-level functional genomics data integration as a tool for understanding physiology: A network biology perspective*.

White Rose Research Online URL for this paper:
http://eprints.whiterose.ac.uk/91929/

Version: Accepted Version

**Article:**

# Multi-level functional genomics data integration as a tool for understanding physiology: A network biology perspective

Peter K. Davidsen[1], Nil Turan[2], Stuart Egginton[3] and Francesco Falciani[1*]

[1]Institute of Integrative Biology, University of Liverpool, Crown Street, Liverpool, UK.

[2]School of Biosciences, University of Birmingham, Birmingham, UK.

[3]School of Biomedical Sciences, Faculty of Biological Sciences, University of Leeds, U.K.

[*]Correspondent author

**Contact Information:**

Prof. Francesco Falciani

Centre for Computational Biology and Modelling, Institute for Integrative Biology

University of Liverpool

Crown Street

L69 7ZB

Liverpool, UK

Tel: +44-151-795-4558

Email: f.falciani@liverpool.ac.uk

**Running Head:** Multi-level network integration to understand physiology.

26  ABSTRACT

27  The overall aim of physiological research is to understand how living systems

28  function in an integrative manner. Consequently, the discipline of physiology has

29  since its infancy attempted to link multiple levels of biological organization.

30  Increasingly this has involved mathematical and computational approaches, typically

31  to model a small number of components spanning several levels of biological

32  organization. With the advent of omics technologies, which can characterise the

33  molecular state of a cell or tissue (intended as the level of expression and/or activity

34  of its molecular components), the number of molecular components we can quantify

35  has increased exponentially. Paradoxically, the unprecedented amount of

36  experimental data has made it more difficult to derive conceptual models underlying

37  essential mechanisms regulating mammalian physiology.

38      We present an overview of state-of-the-art methods currently used to

39  identifying biological networks underlying genome-wide responses. These are based

40  on a data-driven approach that relies on advanced computational methods designed to

41  'learn' biology from observational data. In this review, we illustrate an application of

42  these computational methodologies using a case study integrating an in vivo model

43  representing the transcriptional state of hypoxic skeletal muscle with a clinical study

44  representing muscle wasting in COPD patients. The broader application of these

45  approaches to modelling multiple levels of biological data in the context of modern

46  physiology is discussed.

47

48
49  INTRODUCTION

50  **Modelling in physiological sciences**

51   Physiology has evolved as a series of sub-disciplines attempting to understand

52   organismal function as a combination of interacting components and systems. The last

53   decade or so has witnessed the development of Systems Biology as an investigative

54   approach, and its application in different areas of biology, ranging from

55   engineering/synthetic biology (e.g. design of bacterial strains with improved

56   properties) to health sciences (e.g. disease biomarker identification). Despite the lack

57   of a concise definition acceptable to the majority of the community (30, 32), Systems

58   Biology is frequently understood to be the study of complex regulatory interactions in

59   biological systems using a holistic approach. This is often achieved by integrating

60   different experimental approaches within the conceptual framework of a

61   computational model (i.e. a mathematical representation of a system that allows

62   simulation of its behaviour). Physiology is probably one of the few research areas in

63   biological sciences that have traditionally adopted such an approach. It has long

64   sought to understand the behaviour of complex biological processes and cellular

65   systems using an integrative approach, and has extensively adopted mathematical

66   modelling in its tool set. Classical examples include August Krogh's tissue cylinder

67   model of oxygen transport to skeletal muscle (33), and Huxley's two-state cross-

68   bridge model of muscle contraction (26), which are still used by investigators today.

69   Indeed, this shows that using modelling to study a system as a whole has been a key

70   component of physiology from its early days.

71   As often happens when a distinct discipline branches out of another, there developed

72   over time a separation of ideas based in part on confusion arising from use of esoteric

73   terminology – similar concepts masked by unfamiliar language. There is therefore a

74   need for an overview of this relatively new discipline, to both emphasise the essential

75   links with basic physiological principles and de-mystify the approach such that the

76   available tools may become more widely adopted in physiological research. The

77   overall aim of this opinion-based review is to describe, using concepts that will be

78   intuitive to physiology researchers, different key methodologies available from the

79   Systems Biology community. In addition, we provide a practical step-by-step guide

80   for integrating multi-level data within an analysis pipeline based around inferred

81   interactions of variables, modelled as a network based on statistical correlations, using

82   a worked example in the field of physiological sciences.

83

84   **The advent of Functional Genomics: a challenge for physiological modelling**

85   It is now clear that much of the complex mammalian physiology or pathophysiology

86   cannot be understood in sufficient detail through a reductionist approach alone.

87   Although this approach has proved valuable in explaining broad phenomena and

88   individual mechanisms, linking multiple mechanisms and effects has proved

89   challenging. For example, a disease phenotype is rarely caused by a single

90   dysfunctional gene or protein. Instead, genetic variability, epigenetic modifications,

91   post-transcriptional regulation mechanisms etc. all act in concert to determine a

92   specific high-level phenotypic response (43). The potential for such complex

93   interaction makes data interpretation much more complicated than originally

94   envisioned, highlighting the need to move away from the widespread 'candidate

95   gene' approach (39).

96          Triggered by the advent of genome sequencing, inspired by the Human

97   Genome Project, dramatic technological advances within the last decade or so have

98   led to increased throughput in genome-wide molecular analyses (i.e. genomics,

99   epigenomics, transcriptomics, proteomics, metabolomics). The comprehensive data

100  acquisition tools developed to cope with large datasets have allowed investigators to

101     determine the molecular state of cells, tissues or even entire organs in a single

102     experiment. Such cost-effective omics approaches are now becoming prevalent in

103     biological and medical research, and consequently have been responsible for the

104     generation of an incredibly large amount of multivariate molecular data. A large

105     proportion of this data is available in the public domain via different online databases

106     (e.g. NCBI Gene Expression Omnibus (5), EBI ArrayExpress (7), and PRIDE (29)).

107         For example, mRNA microarray technology and more recently mRNA

108     sequencing, has provided insight into the transcriptional response of skeletal muscle

109     to prolonged endurance exercise training, highlighting a pronounced inter-individual

110     variation at the molecular level that is consistent with the heterogeneous response

111     observed in a population of individuals at the physiology level (31, 59). Statistical

112     models built to explain such variation as a function of gene expression data can be

113     exploited to identify underlying mechanisms controlling tissue homeostasis. The

114     transcriptional signatures identified in such studies likely explain, at least in part, why

115     some people show great improvements in aerobic capacity ($VO_{2max}$) whereas others

116     only experience smaller benefits, despite completing the same supervised exercise

117     training program. Another example of applying omics technology to better understand

118     human physiology concerns the quantification of individual levels of different

119     proteins in health and disease; by use of proteomics methodology, Holloway et al.

120     (24) were the first to investigate adaptations in human muscle protein content to long-

121     term exercise training on a large scale.

122         While such omics-based studies hint at the potential of a data-driven approach,

123     they also illustrate the difficulty in deriving conceptual models underlying the

124     essential mechanisms regulating physiology, as most are restricted to only one aspect

125     of regulation. Perhaps surprisingly, the exponential growth in publicly available omics

126    data (34, 37) has not resulted in a paradigm shift in our understanding of biology. The

127    main reason is the continuing challenge of integrating multivariate datasets spanning

128    multiple organization levels in a way that allows the identification of discrete, small

129    biomolecular networks that are truly important in the context of a specific biological

130    response (47). Such a task cannot be achieved simply using unaided human

131    interpretation. Rather, complex computational techniques are needed that are able to

132    integrate and automatically 'learn' the structure of a biological system. Such a

133    modelling framework is very different from what physiological sciences have

134    traditionally employed.

135

136    **Towards data-driven predictive biology**

137    Although the modelling approach traditionally used by physiologists has been

138    extremely successful, it suffers from severe limitations when challenged with

139    extensive omics data. For example, physiological modelling relies to various degrees

140    on a mechanistic understanding of the biological system of interest (16), which

141    automatically limits the number of components that can be included due to gaps in our

142    current knowledge (19, 47). Moreover, estimation of model parameters, which is

143    usually a challenging task because of experimental limitations (e.g. due to limited

144    amount and quality of data), makes the approach difficult to scale up to a larger

145    number of components and their interactions. Perhaps the most comprehensive

146    example to date is modelling the cardiac cycle based on ion channel kinetics (44).

147        With such large multivariate datasets, and little knowledge about the way

148    biomolecules are connected with each other and to key phenotypic switches, the

149    fundamental question is whether or not we can 'learn' the structure of biological

150    interaction networks from high-throughput data. Clearly, there is a need for

151     sophisticated computational tools that are able to i) integrate genome-wide

152     measurements spanning multiple levels of biological organization (ranging from

153     subcellular to organ level), ii) identify key biomolecular components of the system,

154     and finally iii) statistically infer the way that these biomolecules interact in a pairwise

155     manner to generate an observed biological response.

156         Central to these approaches is the concept of interaction networks, a

157     mathematical representation of a system of biomolecules. Networks are commonly

158     used to describe biological systems at different levels of complexity (e.g. metabolic

159     and signal transduction networks). They can be descriptive models built using a wide

160     spectrum of qualitative data (e.g. biological knowledge of protein-protein interactions,

161     transcription factor binding, etc.) or they can be inferred from quantitative

162     measurements using complex computational models. In this case they can be used to

163     predict the behaviour of the system when perturbed.

164         In the following section, we summarise specific methodologies that can be

165     applied to achieve such tasks.

166

167     COMPUTATIONAL APPROACHES FOR THE ANALYSIS OF COMPLEX

168     DATASETS

169     The process of modelling a biological system from complex multi-level datasets can,

170     for the sake of convenience, be divided into four conceptually distinct yet

171     interconnected approaches (**Figure 1**).

172

173     [Figure 1 to be inserted here]

174

175    The first approach is biomarker discovery (**Figure 1A**), which perhaps is most

176    widely used in the analysis of functional genomics datasets. Here the objective is to

177    identify measurable variables that are predictive of a given outcome (e.g. the response

178    to physical training in a population of individuals). Such measurements can be

179    molecular (e.g. gene expression, protein levels, metabolite concentrations, genetic

180    mutations) and/or more traditional physiological endpoints (e.g. endurance, $VO_{2max}$).

181    The identification of predictive biomarkers can be achieved by use of univariate and

182    multivariate variable selection strategies that aim to identify the most relevant

183    explanatory measurement(s), while developing a computational model that can

184    accurately predict an outcome (60). Univariate methods will test every variable (e.g.

185    expression of a given gene) on its own, whereas multivariate methods test

186    combinations of variables for their ability to explain a given outcome. Clearly,

187    multivariate approaches better resemble the complex nature of biological networks,

188    and therefore are more likely to provide insights into the mechanisms underlying a

189    complex phenotypic trait. Consistent with this notion, multi-gene biomarkers are often

190    required for robust predictions in independent datasets.

191    The second approach (**Figure 1B**) consists of 'reverse engineering'

192    biomolecular networks from observational data (i.e. infer regulatory interactions

193    between quantified biomolecules based on mathematical principles). Here the overall

194    aim is to reconstruct the underlying structure of interactions between biological

195    molecules profiled using omics tools (ideally from multiple data sources) and rigorous

196    statistics. Such a network inference framework can be achieved by applying a

197    multitude of approaches with varying underlying data assumptions and modelling

198    principles, including ordinary differential-equation (ODE)-based methods (3),

199    probabilistic modelling techniques (e.g. Bayesian theory models) (42, 64), state-space

200    representation models (23), and correlation-based methods. Note, while the first three

201    approaches are able to infer directed networks, their capability is currently limited to

202    inferring smaller networks with few variables due to increased computational

203    complexity than possible with correlation approaches.

204         Importantly, this network inference part may potentially benefit from a

205    biomarker discovery phase, since it has been shown that identified predictive

206    variables are more likely to be directly controlling important physiological processes,

207    and therefore are good candidates to include in a network (47). Similarly, whole

208    networks can be used as an input for biomarker discovery procedures. It has been

209    shown that often the overall 'activity' of a biological network (e.g. a specific

210    signalling pathway) is a better predictor than a few key individual genes, proteins

211    and/or metabolites. This implies that in the coming years predictive biomarkers are

212    more likely to consist of a relatively large panel of measurements, possibly spanning

213    multiple levels of complexity within a pathway. Current omics platforms are

214    experiencing a rapid development as well as drop in costs, making routine collection

215    of large datasets a feasible option. Once a robust biological network has been inferred

216    this may serve as a good basis for developing a more conventional modelling

217    approach to provide explanations for observed phenomena that requires a mechanistic

218    understanding of the system (**Figure 1C**).

219    Finally, multiple computational models that initially were developed independently

220    can be integrated into a larger and more complex models, which allow responses to

221    physiological/pathological challenges to be simulated, thus integrating effects across

222    multiple organs and/or pathways. These complex models are often referred to as

223    decision support systems because of their potential to provide information about the

224    expected outcome of a therapeutic intervention (**Figure 1D**).

225        Several large international projects aiming at the development of such

226        technology into Systems Medicine integrated frameworks have been established so

227        far, e.g. the Virtual Physiological Human (VPH) project funded by the European

228        Commission $7^{th}$ Framework Programme, which aims to aid clinically relevant

229        research by establishing a framework for handling and integrating various mechanistic

230        models spanning different levels of organizational complexity (ranging from

231        molecular components to organ function). By unifying the modelling languages

232        employed across the different mathematical models included, parameters of a

233        particular model in the hierarchy can be processed by other appropriate models at a

234        lower hierarchical level. These global initiatives should be considered long-term

235        goals, aiming at understanding human physiology quantitatively as a dynamic system.

236        Developing a comprehensive model of a biological system requires integrating

237        mechanistic and probabilistic inferences. The mathematics for performing such a task

238        is in its infancy, and more development is needed. However, a successful example is

239        illustrated by the anatomically based model of human heart ventricles (44). In the

240        following sections we aim to provide an overview of some of the methodologies that

241        can be used to infer biomolecular networks, as well as introduce one particular

242        approach we have found useful in our research.

243

244        **Inference of biological networks from observational data**

245        Reverse engineering is an evolving field within network-based Systems Biology. The

246        rapid accumulation of omics data in the post-genomic era has made it possible to infer

247        (aka 'reverse engineer') models of cellular systems with the overall aim of deducing

248        the regulatory structure at a sub-cellular level. Most of the network-based approaches

249        that have been developed are in fact general and can be applied to any type of

250 experimental data. However, because the mRNA expression profiling technology is

251 the most mature omics discipline, most applications have been developed to

252 reconstruct transcriptional networks (i.e. decode the mechanisms of transcriptional

253 control). However, recently it has become apparent that, irrespective of the

254 methodology used to generate data, in order to be able to recapitulate the complex

255 behaviour of a biological system it is essential to integrate multiple types and scales of

256 experimental data (e.g. transcriptomic, proteomic, metabolomic).

257

258 **Static vs. dynamic networks**

259 Biological networks can be reconstructed from two different types of experimental

260 studies: either cross-sectional, e.g. representing a population of individuals at a given

261 time (i.e. steady-state measurements following an experimental perturbation), or

262 prospective, where the experimental data is available across a defined time-course. In

263 reverse engineering, statistical inference of biological causality is an important goal

264 (56). A simple example of causality could, for example, be a transcription factor

265 regulating the expression of several target genes. Since determining cause and effects

266 implies a direction (i.e. the cause precedes the effect), inference of causality from

267 cross-sectional studies presents a challenge due to their static nature, one that is less

268 difficult when a time-course is available. However, it must be stressed that both

269 approaches are often used in combination to, for example, integrate clinical cross-

270 sectional studies (thereby providing the researcher with a static network

271 representation) and experimental intervention studies that can provide dynamic

272 (prospective) models of the process being studied. At present, most of the developed

273 techniques infer regulatory networks without any causality information (likely due to

274 the scarcity of time-course datasets due to their higher costs). However, a small

275　number of causality detection techniques have been proposed in the literature such as

276　dynamic Bayesian networks (48) and Granger causality (46). It is also important to

277　point out that true time-course datasets can only be developed when the sequence of

278　events is measured within the same cells/tissues. This is for example achieved with

279　imaging techniques that require complex molecular probes, and can typically be only

280　applied to measure a relatively small number of system components (14). Omics

281　technologies unfortunately are disruptive, so time course data derived using these

282　approaches are in fact a sequence of independent snapshots, which clearly limits the

283　potential use of dynamical modelling tools.

284

285　**A primer for network inference methods**

286　The simplest method for inferring statistical relationships between experimental

287　variables is computing the pairwise correlation coefficient across a large collection of

288　heterogeneous samples (8). Usually such an approach is not able to identify complex

289　non-linear dependencies, and does not discriminate between direct and indirect

290　connections. More complex methods, such as the mutual information (MI) based

291　Algorithm for the Reconstruction of Accurate Cellular Networks (ARACNE) (38),

292　also aim at establishing a statistical relationship between pairs of variables but have a

293　stronger theoretical foundation. Because of the added mathematical complexity they

294　can capture a broader range of biologically relevant dependencies between variables

295　including non-linear, non-monotonic relationships; importantly, they can distinguish

296　between direct and indirect relationships. ARACNE is a free tool for which a Java-

297　based graphical user interface exists; hence investigators do not need any

298　programming skills in order to use the software.

299        ARACNE relies on estimating the probability that a variable (e.g. the

300        expression of a gene or a protein) assumes a certain 'state' (i.e. abundance) given the

301        state of another biomolecule (conditional probability). A number of alternative MI-

302        based implementations have been proposed during the last decade (e.g. Context

303        Likelihood Relatedness (CLR) (13), Minimum Redundancy/Maximum Relevance

304        Networks (MRNET) (41)), which mainly differ by the way inferred indirect

305        relationships (so-called 'edges') are removed once the dependencies between all pairs

306        of variables have been mathematically formulated. In such analyses, unwanted

307        indirect interactions occur by default if there is strong correlation between

308        biomolecule 1 and biomolecule 2, and between biomolecule 1 and biomolecule 3 in a

309        three-node clique (i.e. a triplet of connected variables).

310        An MI value of zero means that there is no dependency (i.e. no information

311        flow) between two variables, whereas an MI value of 1 indicates a perfect association

312        between them, and therefore, a likely strong regulatory interaction between them. For

313        each inferred dependency, a P-value is calculated based on the distribution of MI

314        values between random permutations of the original dataset, thereby allowing the

315        elimination of all non-statistically relevant dependencies by thresholding using an

316        appropriate (user-defined) cut-off level. Importantly, the quality of the inferred

317        interaction network depends on the arbitrarily selected probability cut-off. A small

318        threshold (e.g. P=0.05) gives a high recall (i.e. fraction of true dependencies that

319        could be inferred) but low precision, whereas a higher threshold (e.g. $P=10^{-6}$) yields

320        better precision (i.e. fraction of inferred dependencies that really are in the network)

321        while suffering from a low recall. A further advantage of MI as an information-

322        theoretical measure of dependency between variables concerns its relatively low

323        computational requirements for building an interaction network. Hence, MI is able to

324　handle very large data matrices with thousands of experimental variables, whereas

325　most of the other more advanced techniques mentioned (e.g. Bayesian methods) can

326　only deal with much smaller numbers of variables (<100) because of the high

327　computational complexity.  However, in order to infer robust statistical associations

328　based on MI a fairly large sample size is required (> 50-100 biological replicates), due

329　to the required estimation of the (joint) frequency distribution of the connectivity.

330　Interaction networks derived from such reverse engineering methodologies can be

331　visualized and further analysed using various freeware software tools such as

332　Cytoscape (55), Pajak (6), and BioLayout (18). A comprehensive list of visualization

333　tools focused on interaction networks and their web-links has recently been reviewed

334　(17).

335　　　Up to now, these information-theoretic approaches have usually been

336　employed on gene expression data only, due to the wealth of such data available.

337　However, as physiologists have known for many decades, biological systems are

338　usually more complex and multi-layered. Indeed, despite some popularist science

339　writing to the contrary, genes on their own are merely permissive elements within

340　biological systems (43). Further, it has been shown that when multiple types of data

341　(e.g. copy number variants, protein or microRNA expression levels) are incorporated

342　in the network inference pipeline, the accuracy of the learned network topology

343　increases (49). Hence, at present there is a call for methodologies that can embed

344　multiple data sources in a single computational framework. Our recent work has

345　focused on methods that are able to handle large-scale, multi-dimensional genomic

346　datasets (9, 21).

347

**Topological analysis of inferred biological networks provides useful biological insight**

Up to now we have described some of the most widely used methodologies for inferring regulatory networks. However, an immediate challenge arises in interpreting these often large, complex networks that visually present as a 'hairball' (i.e. too dense a collection of connections to comprehend as a whole) (40). A simple solution, although not very objective, is to focus the analysis around a favourite gene(s). In this scenario, the investigator typically examines the manually selected sub-network in order to identify unknown or unexpected biological relationships, which in turn may be used to formulate new hypotheses. Such 'discovery-led' science may be useful when there is insufficient information to generate hypothesis.

Alternatively, the topological properties of the network can be used to identify interesting genes and sub-networks that can be interpreted. We and others have demonstrated the existence of a higher-level, modular organization in biological networks (47, 52, 54), i.e. components of biological systems that act in collaboration to carry out specific biological processes. Consequently, several modularization approaches have now been developed to help group subsets of cellular components based on a given property, such as topological structure or functional role. Such decomposition of a large complex network into relatively independent sub-networks (or 'modules') has been shown to be an effective way to deduce the underlying structure of the fully connected network containing many hundred variables (so-called 'nodes'), as each module can then be analysed independently. In addition, studies have demonstrated that such identified network modules can serve as better predictors of a physiological response than the classic biomarker discovery approach (see Figure 1).

373        In biomolecular interaction networks, as well as sub-networks, nodes have

374    different levels of connectivity (i.e. number of interactions with other nodes). It has

375    been shown that such interaction networks have so-called 'scale-free' structure

376    properties, as their node connectivity distribution fits a power law (4). Such a power

377    law degree distribution implies that most of the connections between biomolecules is

378    linked to a small number of highly connected nodes, such that a large proportion of

379    the molecular state of a cell can be explained by a small subset of biomolecules (so-

380    called 'hub' nodes; e.g. a transcription factor that regulates many more genes than

381    average). Hence, in biological networks a hub is often assumed to be a key component

382    of a regulatory networks, hence important for the function of a cell/tissue under

383    investigation. This assumption is supported by the fact that random node disruption

384    does not significantly affect the network architecture, whereas deletion of hub nodes

385    leads to a complete breakdown of the network structure (1). Hence, adjusting the

386    spatial position of each node according to its interconnectivity has been shown to be a

387    simple, yet effective way of visualizing large complex interaction networks (57).

388        More advanced methods to extract information from complex networks exist

389    that aim to identify functional modules (i.e. sub-networks of biomolecules that are

390    linked to the same biological function), e.g. by integrating both physical interactions

391    (i.e. experimentally validated protein-protein interactions) and mRNA expression data

392    (27). In this context, an identified functional module represents a putative multi-

393    protein complex that is transcriptionally regulated in a specific experimental condition

394    (e.g. treatment vs. control). Hence, by considering additional data on a different level

395    of organization, one can potentially infer a clearer composite picture of the underlying

396    biological function.

397        Finally, in order to generate objective hypotheses about biological processes

398    controlled by a specific hub node or sub-network, functional enrichment analysis can

399    be performed on all its direct neighbours (i.e. all the adjacent nodes that are directly

400    connected to the hub) (25). Such enrichment analysis aims at reducing complexity by

401    defining groups of molecules (represented by gene sets) that share similar biological

402    functions (e.g. a class of adhesion molecules). To accommodate latest advances in

403    knowledge, the different annotation databases used for this purpose (e.g. Gene

404    Ontology (2) and KEGG (45)) are frequently updated by curators. Using software

405    tools like the web-based application DAVID (11) or applications such as BiNGO (36)

406    developed specifically for use with software visualization tools like Cytoscape, one

407    can quickly determine whether any gene sets are statistically over-represented, thus

408    generating hypotheses on the biological processes controlled by those factors outlined

409    above.

410

411    CASE    STUDY:    INFERENCE    OF    OXYGEN-DEPENDENT

412    PATHWAYS IN SKELETAL MUSCLE

413    The main purpose of this case study is to illustrate in a step-by-step manner the

414    application of reverse engineering to integrate supra-cellular physiological measures

415    and genome-wide expression profiling. From a more biological perspective we aim to

416    identify a clinically relevant signature of hypoxia in skeletal muscles.

417    This analysis uses two different datasets. The first is a publicly available dataset

418    (GSE27536) representing a cohort of COPD patients and healthy controls matched for

419    age and smoking history (10) (see **Table S1** for subject characteristics), which

420    includes gene expression profiling in vastus lateralis muscle and whole-body

421    physiological variables (e.g. $VO_{2max}$, minute ventilation, $PaO_2$) (50)(61). The second

422  dataset represents an unpublished, genome-wide transcriptional response of mouse

423  soleus muscle to a gradual decline in atmospheric oxygen concentration (GSE64076).

424  Using the first dataset, representing the transcriptional state of skeletal muscles in a

425  COPD cohort (**Figure 2A**), we first show how to infer connections between oxygen

426  availability (e.g. $VO_2max$), oxidative stress (protein carbonylation) and gene

427  expression signatures (**Figure 2A-C**).

428  Having defined an oxygen-related signature in the disease setting we then transpose

429  these findings in a mouse model of gradual hypoxia (second dataset, **Figure 2D-E**).

430  Here we use a different computational approach to develop a hierarchical dynamical

431  model explaining the transcriptional response of oxidative leg muscles to a prolonged

432  gradual reduction in blood oxygenation (hypoxaemia) (**Figure 2F-G**). The model we

433  describe below validates the notion that the signature identified using the clinical

434  study may be truly triggered by changes in oxygen availability. Moreover, the model

435  contributes to the understanding of the transient events following oxygen depletion

436  that cannot be observed using a cross-sectional clinical study.

437

438

439  **Step 1. Linking physiological measurements and gene expression data in the**

440  **COPD cohort**

441  In order to reconstruct an interaction network spanning multiple levels of

442  organization, we have utilised the following strategy that was developed earlier (61).

443  1. <u>Combining measurements from different data sources</u>

444  In order to combine gene expression data with whole-body physiological readouts, all

445  variables need to have the same units of measurement (as the range of e.g. VEGF

446  mRNA expression values is very different from that of $VO_{2max}$). All such raw scale

447   units can be unified by simply 'transforming' each experimental variable to have the

448   same dynamic range, e.g. this can be achieved by standardising measurements across

449   samples to have a mean of 0 with a standard deviation of 1. Such an established

450   approach, called z scoring, enables us to treat the physiological indicators as

451   individual 'nodes' in the inferred interaction network with states (just as each gene on

452   the array is treated).

453   Definition of a biological framework for data-driven network inference

454   The outcome of data-driven reverse engineering of biological networks, in the

455   absence of any biological assumption(s), often provides results that are difficult to

456   interpret due to the large number of inferred significant interactions. Thus, to reduce

457   complexity of the problem, we decided to focus the analysis on the set of

458   physiological parameters and genes encoding for enzymes in the central bioenergetic

459   pathways (i.e. TCA cycle, oxidative phosphorylation, glycolysis) (see **Table S2** for

460   the full list of variables). The latter choice is reasonable considering the paramount

461   importance of these molecular pathways in skeletal muscle adaptation. The overall

462   strategy is therefore to identify biomolecules that are highly correlated (based on MI)

463   with biologically important experimental variables. Such a focused analysis will

464   generate multiple network modules of interacting biomolecules, each with a

465   bioenergetic hub gene or physiological measurement at its centre. Two modules will

466   be linked together if a specific gene is statistically linked to both hubs.

467

468   2. <u>Reverse engineering.</u>

469   In order to infer robust regulatory relationships between variables in the integrated

470   multi-level dataset, we used the ARACNE algorithm. This choice was based on the

471   large number of measured variables to be considered by the mathematical framework.

By combining all genes expressed in human skeletal muscle (>10,000 mRNAs) with the list of physiological variables we far exceed the number of variables that can be handled by more advanced network inference methods (e.g. Bayesian methods). Hence, we infer a static network without any obvious hierarchical organization. The result of an ARACNE run is an 'adjacency matrix' containing MI values for all pairwise interactions above the specified MI threshold, which can be visualized automatically in Cytoscape.

After calculating MI-based dependencies between all the different variables in our multi-level data matrix, all those inferred regulatory interactions with an MI value below 0.22 (corresponding to a P-value cut-off of $10^{-6}$) were removed. Such filtering of weaker statistical dependencies is an important step in the generation of a more sparse interaction network, which can more easily be interpreted by the investigator. The stringent P-value cut-off means that the remaining associations have been inferred with high precision at the cost of a lower recall rate.

3. <u>Network visualization</u>

Data visualized as a network are often easier to interpret than long lists of biomolecules and their associated statistical dependencies. Hence, the numeric output of ARACNE, which contains MI values for all pairwise associations, was imported into Cytoscape for visualization, a conventional way of analysing interaction networks. Briefly, we reconstruct the network neighbourhood of each of the bioenergetic 'seed' genes listed in **Table S2** (i.e. all variables directly connected to them). The neighbouring variables can either be genes expressed in muscle and/or physiological variables. **Figure 3** summaries key regulatory associations (based on MI) between this seed set of genes and their immediate neighbours.

497

498 4. <u>Functional analysis of the network hubs</u>

499 We further explored whether the direct interacting neighbours of each central

500 metabolism pathway mapped to functional categories (i.e. GO terms) as well as

501 KEGG pathways. Notably, a marked enrichment of the different bioenergetic

502 compartments was observed (**Figure 3**, boxes A-C) that clearly highlights the

503 interconnected nature of the bioenergetic machinery, i.e. functionally related genes

504 appear to be co-expressed.

505

506 [Figure 3 to be inserted here]

507

508 5. <u>Biological interpretation</u>

509 The most important finding of the current analysis is that among the direct neighbours

510 to each bioenergetic pathway, particularly the two oxidative ones, we noted a

511 statistical over-representation of genes encoding histone deacetylase enzymes (i.e.

512 HDAC and SIRT mRNAs). This observation is consistent with previous studies that

513 have highlighted the importance of sirtuins in regulating metabolism (15, 22, 28).

514 Further, the protein deacetylase SIRT3 that primarily is localized in the mitochondrial

515 matrix was also significantly positively correlated to both arterial oxygen tension

516 (PaO$_2$) and oxygen uptake (VO$_{2max}$). In support of deacetylation being an important

517 control point, it was recently shown that Sirt3 knockout mice exhibit decreased

518 oxygen consumption, thus affecting cellular respiration (28). Hence, besides the

519 obvious oxygen-driven effect on aerobic pathways (as indirect measures of oxygen

520 availability such as VO$_{2max}$ are linked to key genes in oxidative phosphorylation), the

521 present network-based Systems Biology approach points to tissue hypoxia as being a

522    potential important player in modifying expression of deacetylase modifying enzymes

523    in severe COPD patients with a muscle wasting phenotype. Our Systems Biology

524    approach also negatively links protein carbonylation (an established proxy measure

525    for oxidative stress; (58)) to Complex 1 and 3 in the electron transport chain (**Figure**

526    **3**, bottom left part). The validity of such an association is further strengthened via

527    functional enrichment analysis using DAVID, as a significant fraction of direct

528    neighbouring genes to protein carbonylation is statistically associated to gene

529    ontology (GO) terms representing cellular respiration.

530    If we then focus on the genes in the glycolytic pathway (**Figure 3**, top right

531    part), a high proportion of pro-inflammatory mediators/receptors (e.g. IL1B, IL1R1

532    and TNFRSF21) are among the direct neighbours, as indicated by the enrichment of

533    the 'inflammatory response' GO term (**Figure 3**, box A). Hence, hypoxia is pro-

534    inflammatory, as seen by more traditional observation methods (20).

535    Multi-scale network inference approaches, similar to that illustrated in Figure

536    3, have proven very effective in generating robust hypotheses (e.g. 45). However,

537    statistical associations may not represent causality, particularly when the inferred

538    associations stem from steady-state measures. Thus, in order to validate our

539    hypothesis that varying oxygen levels (represented by $VO_{2max}$ and $PaO_2$) control the

540    expression of epigenetic modifiers, we used a more sophisticated network inference

541    algorithm that can learn the structure of networks from time-course data. We applied

542    this dynamic inference approach to a murine model of hypoxia (Step 2).

543

544    **Step 2. Gene expression dynamics in response to tissue hypoxia**

545    Animal models are commonly used for studying the in vivo effects of hypoxia, for

546    ethical reasons, where severe or prolonged hypoxaemia is induced and invasive

547 samples are required to explore mechanisms. Importantly, hindlimb skeletal muscles

548 have been reported to alter metabolic phenotype and reduce fibre size in response to

549 prolonged hypoxic stress in mice (53, 63), highlighting their potential relevance as a

550 pre-clinical model of muscle wasting in COPD patients. In order to experimentally

551 test the hypothesis derived from the clinical COPD network presented in **Figure 3**, we

552 therefore exposed adult male C57/Bl6 mice to chronic systemic hypoxia for up to 2

553 weeks, in order to simulate levels of hypoxaemia reported in COPD patients with

554 advanced respiratory insufficiency. To capture the temporal effect of reduced oxygen

555 tension on gene regulation, we sampled and gene profiled the soleus muscle (n=4) at 3

556 different time-points (day 3, 7 and 14) following initiation of the gradual hypoxic

557 insult (i.e. the $O_2$ level was gradually lowered to 10% over the first week and kept

558 stable during the second week) (**Figure 2**, bottom part).

559 First, a high-level representation of the temporal transcriptional changes was

560 performed using a variable reduction technique called principal component analysis

561 (PCA) (**Figure 4B**). When plotting replicates of two variables against each other, it is

562 relatively easy to see which is a better discriminating factor; visual inspection

563 becomes increasingly difficult as the number of variables increase, hence the need for

564 PCA. In essence, this method aims at 'tilting' the axes through the multidimensional

565 data space, such that the first principal component accounts for as much of the

566 variation in the original dataset as possible (the assumption is that the most important

567 dynamics in the dataset are the ones with the largest variation). Our PCA revealed that

568 the early dynamics of hypoxia is captured by the first principal component whereas

569 the 2[nd] most important principal component (in terms of variance captured) separated

570 the later time-points. Further, functional enrichment analysis of the differentially

571 expressed genes (ANOVA, P<0.05) using DAVID (**Figure 4A**), highlighted several

572 important pathways/ontologies. Most striking was the enrichment of protein catabolic

573 process and ubiquitin-mediated proteolysis among genes up-regulated at day 7 and 14,

574 clearly suggestive of a transcriptionally regulated muscle wasting phenotype driven

575 by the experimentally induced hypoxaemic state.

576      State space models (SSMs) are a class of probabilistic graphical models

577 (Koller and Friedman, 2009). SSM provides a general framework for analyzing

578 deterministic and stochastic dynamical systems that can be measured/observed

579 through a stochastic process. The SSM framework has been successfully used for the

580 analysis of gene expression data (23, 51). In its simpler application the model

581 formalises the effect of hidden, unmeasurable factors in specifying observed gene

582 expression changes over time. The inclusion of these hidden factors is important since

583 we cannot hope to measure all possible factors contributing to genetic regulatory

584 interactions (e.g. levels of regulatory proteins as well as effects of mRNA and protein

585 degradation).

586 The next step was to apply state-space modelling to reverse engineer transcriptional

587 network modules (i.e. representing discrete temporal dimensions) from our replicated

588 murine time-course dataset. Such module-based reduction in complexity allows

589 analysis of hundreds or even thousands of genes, as those with a similar temporal

590 expression profile are aggregated into a transcriptional module. To allow construction

591 of a near genome-level model, we took advantage of a newer approach that

592 incorporates this concept of modularization (23).

593      A state space model can reconstruct the topology of a network representing the

594 systems dynamics, despite a relatively small number of time-points, by using

595 biological replicates for each time-point (23). In order to reduce complexity, variables

596 that do not change significantly are excluded from the modelling process. In this case

597 study, genes deemed to be significant by ANOVA at a 1% significance level, as well

598 as all hub genes listed in Table S3, were included (931 variables in total). The hub

599 genes were chosen to represent the different components in our interpretative model

600 derived from the clinical COPD dataset (**Figure 3**). Finally, the experimentally set

601 oxygen level was used as an independent variable.

602

603 [Figure 5 to be inserted here]

604

605 Based on unsupervised clustering using HOPACH within the software

606 programming environment R (35), we identified 8 distinct gene clusters with similar

607 expression profiles. Hence, to model the effect of hypoxaemia on the skeletal muscle

608 transcriptome the hidden state dimension was set to 4, as each inferred module

609 contains both a positive (+) and a negative (-) component.

610 The hierarchical dynamic model in 4 temporal dimensions shows that modules

611 1 and 2, which sit on the highest level of hierarchy (i.e. precede others in time), were

612 enriched in GO terms related to muscle contraction, bioenergetic pathways, and

613 inflammation among others (**Figure 5**). Interestingly, the experimental oxygen

614 concentration was represented in module 1(-) whereas two deacetylases SIRT3 and

615 SIRT5 were found in module 2(-). A negative influence is observed of module 1 on

616 module 3, which is located further down the temporal hierarchy. Module 3(+) is

617 highly enriched in inflammatory processes whereas its negative counterpart mainly

618 represents two key signalling pathways (mTOR and insulin). At the lowest temporal

619 level we find module 4, which is enriched in GO terms related to muscle

620 differentiation, tissue remodelling and blood vessel development. Interestingly, three

621 HDACs are represented in module 4(+) (**Figure 5**). **Figure 6** represents a more

622    focused version of **Figure 5**, highlighting the most significant interactions between

623    components in the four inferred modules from Figure 5.

624

625    [Figure 6 to be inserted here]

626

627        We therefore conclude that the inferred dynamic model using a state space

628    modelling approach appropriately recapitulates the interpretative model advanced in

629    **Figure 3**. In addition, it identifies oxygen at the highest level of hierarchy, whereas

630    key effector functions controlled by oxygen such as inflammation and muscle

631    differentiation are downstream in the temporal hierarchy.

632

633    CONCLUSIONS

634    The aim of this brief review is to provide an intuitive overview on data-driven

635    'learning' of biological pathways, linking molecular and physiological readouts. We

636    used a case study to make it easier for experimental biologists to see the potential of

637    computational biology to provide interpretative models of complex patterns, and

638    stress that the identification of general properties of a system from a genome wide

639    analysis of a molecular state of a system is a very powerful approach.

640    The ability to generate omics data with relatively accessible technologies offer an

641    unprecedented opportunity to study how genetic information is used to control

642    complex biological processes and their interaction. Until now we have only been able

643    to understand a fraction of that complexity. The computational methods described in

644    this review are designed to support this effort in the measure that they help isolate

645    from these large datasets molecular signatures that correlate to phenotypic outcome.

646    With the help of computational biology, we are therefore able to develop hypothesis,

647 which can be experimentally validated. In this context data-driven biology is not in

648 contraposition with hypothesis-driven research. Instead it is a tool that support

649 hypothesis generation in the event that the data is too complex to be interpreted solely

650 using common sense. This approach is well developed in other areas of science, such

651 as cancer biology where there is a vast literature showing that important hypothesis

652 can be generated from modelling of these large datasets (12, 62).

653 In this manuscript, we demonstrate the development of an integrative workflow that

654 incorporates measurements from different levels of cellular and molecular

655 organization using a case study representing muscle wasting in COPD. The outline

656 provides an exemplar where individual steps can be modified according to the type of

657 data at hand and addition data types added. For example, in contrast to established

658 gene expression microarrays, techniques for proteomics and especially metabolomics

659 are still under development. Once it is possible to measure the whole proteome and

660 metabolome of a sample, systems identification pipelines will clearly benefit from

661 these omics techniques.

662 The specific findings in the case study relate to the definition of an oxygen dependent

663 signature in COPD. Such signature (exemplified in **Figure 3**) is static and entirely

664 based on statistical inference. The model is therefore based only on correlation

665 between a series of patient biopsie snapshots, and therefore does not allow any

666 inference of causality. The use of a mouse model of gradual hypoxia allowed us to

667 demonstrate that a signature inferred from the clinical cohort is indeed modulated by

668 experimental reduction in oxygen levels. Moreover, the development of a

669 mathematical model identifies oxygen as the most upstream event as an emergent

670 property. This may appear an obvious finding but, from a methodological perspective,

671 validates the analytical approach.

672    The data we have used in this case study is gene expression profiling, and as such is

673    representative of available datasets. This has several limitations. The first is that

674    models including multiple levels in the expression of genetics information (e.g.

675    epigenetics, microRNA, proteomics, metabolomics, etc.) may better represent

676    biological complexity. However, current computational methods are inadequate to

677    represent properly the interaction between these levels. Moreover, time course data

678    that rely on disruptive sampling strategies are not true time course experiments. As

679    the new functional genomics technologies develop further, as well as novel

680    approaches to model the interaction between different layer of biological organisation,

681    we expect that the efficacy of data-driven approaches will increase further.

682

**Figure legends**

684 **Figure 1:** Schematic representation of the process involved in modelling a biological

685 system by integrating knowledge from various sources, and complex multi-level

686 datasets. The process can be conceptually subdivided into four distinct yet

687 interconnected approaches (A-D). The experimental data used can either be novel

688 multivariate data generated in your own (wet) laboratory or taken from a public

689 repository. These may then be used to identify predictive biomarkers, i.e. variables

690 that are predictive of a defined outcome (e.g. response to exercise training), and also

691 to inform development of important networks that infer such outcomes; experimental

692 data and other source of biological knowledge may also be useful in refining these

693 representation of complex interactions. Such networks may in turn aid biomarker

694 discovery, but are an essential precursor to computational models that are able to

695 explore underlying molecular mechanisms; again, knowledge of specific biological

696 issues may help in their refinement. Finally, incorporation of these models into larger

697 scale analyses offer the potential for in silico experimentation, whereby e.g. the effect

698 of different therapeutic interventions on disease outcome may be tested.

699

700 **Figure 2:** Schematic representation of the analysis strategy used in the case study,

701 highlighting how the inferred static multi-scale network from the clinical COPD

702 cohort (Fig 2A-C) can be bridged to the inference of a dynamical network

703 representing the temporal progression of events following an experimental challenge

704 (hypoxic exposure) in a murine animal model (Fig 2D-G). Having identified a clinical

705 condition with known outcome (exercise intolerance in patients with respiratory

706 disease), we could target unknown mechanisms by focussing on one likely source of

707 functional limitation  (skeletal muscle dysfunction $\pm$ central limitation on $O_2$ supply),

708    and generate data characterising the phenotype. Both genomic and physiological

709    readouts were used to construct a network of inferred interactions, which was then

710    interrogated to identify statistically robust linkages among broad biological functions.

711    While very useful in providing a list of useful biomarkers, there remains a potential

712    limitation with single point associations. The dynamic nature of relationships is

713    captured by repeated measures across a suitable time scale (which will vary for

714    different molecular, physiological and structural responses) using an animal model of

715    respiratory distress, where the transcriptome-based model demonstrated the central

716    importance of oxygen in the response.

717

718    **Figure 3:** Graphical representation highlighting putative regulatory associations

719    (significant correlation between two factors is shown as a dotted line) that likely

720    represent robust interactions, based on high mutual information values. The focus is

721    on central metabolism pathways (i.e. glycolysis, TCA and OXPHOS, respectively)

722    and their immediate neighbours. The grey boxes define functional enrichment of the

723    different bioenergetic compartments based on direct neighbours. Individual genes of

724    relevance are grouped into modules with others of related function, as are

725    physiological readouts that may be treated in a similar manner for statistical analysis.

726    C1-5: the different complexes in the electron transport chain. The value of such an

727    approach is in providing a detailed overview of a complex interaction network,

728    reducing the huge number of potential factors into groups of defined function, and

729    offering a limited number of candidates whose utility as biomarkers or therapeutic

730    targets may be experimentally verified.

731

732     **Figure 4:** High-level representation of temporal transcriptional changes in the murine

733     model of hypoxia. A) Graphical representation of the pre-clinical experimental design.

734     The oxygen level was gradually decreased from 21% to 10% during the first week and

735     mice were housed for another week at this oxygen concentration. B) Principal

736     component plot highlighting the transcriptional dynamics caused by the hypoxic

737     challenge. C) Hierarchical clustering using mRNA expression levels of genes

738     modulated by hypoxia (P<0.05). Each row represents a transcript and each column

739     represents a sample. Red and green colours indicate expression levels above and

740     below the median value of the distribution of signal, respectively. Using solid yellow

741     lines we have subdivided genes into overall trends in order to help the reader.

742     Enriched functional terms within these are listed next to the heatmap.

743

744     **Figure 5:** The hierarchical dynamic state-space model identified 4 modules (x-axes

745     define length of hypoxic exposure), each characterised by two separate transcriptional

746     profiles: plus and minus, representing up- and down-regulation, respectively. The

747     hierarchical position of the modules represents the estimated temporal structure of the

748     network. Functionally enriched GO terms  (regular text) as well as key genes (italics)

749     are identified next to the relevant module. Blue arrows represent temporal repression

750     whereas red arrows represent temporal induction. The numeric value next to each

751     arrow represents the estimated coefficient.

752

753     **Figure 6:** A higher resolution representation of Figure 5, highlighting the most

754     significant gene interactions between components in the four inferred modules. Lines

755     represent factor interactions based on mutual information (blue represents temporal

756     repression, red represents temporal induction). Genes are colour coded for broad

757    functional categories (red=cytokines; blue=epigenetic modifiers; green=aerobic

758    metabolism; purple=muscle differentiation; yellow=cell-interaction).

759

760

**References**:


761

762

763    1.    **Alderson D**, **Doyle JC**, **Li L**, **Willinger W**. Towards a Theory of Scale-Free
764          Graphs: Definition, Properties, and Implications. Internet Math 2: 431–523,
765          2005.
766

767    2.    **Ashburner M**, **Ball CA**, **Blake JA**, **Botstein D**, **Butler H**, **Cherry JM**, **Davis**
768          **AP**, **Dolinski K**, **Dwight SS**, **Eppig JT**, **Harris MA**, **Hill DP**, **Issel-Tarver L**,
769          **Kasarskis A**, **Lewis S**, **Matese JC**, **Richardson JE**, **Ringwald M**, **Rubin**
770          **GM**, **Sherlock G**. Gene ontology: tool for the unification of biology. The Gene
771          Ontology Consortium. Nat Genet 25: 25–9, 2000.
772

773    3.    **Bansal M**, **Della Gatta G**, **di Bernardo D**. Inference of gene regulatory
774          networks and compound mode of action from time course gene expression
775          profiles. Bioinformatics 22: 815–22, 2006.
776

777    4.    **Barabasi A**, **Albert R**. Emergence of scaling in random networks [Online].
778          Science 286: 509–12, 1999. http://www.ncbi.nlm.nih.gov/pubmed/10521342
779          [26 Aug. 2015].
780

781    5.    **Barrett T**, **Suzek TO**, **Troup DB**, **Wilhite SE**, **Ngau W-C**, **Ledoux P**,
782          **Rudnev D**, **Lash AE**, **Fujibuchi W**, **Edgar R**. NCBI GEO: mining millions of
783          expression profiles--database and tools. Nucleic Acids Res 33: D562–6, 2005.
784

785    6.    **Batagelj V**, **Mrvar A**. Pajek - Program for large network analysis.
786          Connections 21: 47–57, 1998.
787

788    7.    **Brazma A**, **Parkinson H**, **Sarkans U**, **Shojatalab M**, **Vilo J**,
789          **Abeygunawardena N**, **Holloway E**, **Kapushesky M**, **Kemmeren P**, **Lara**
790          **GG**, **Oezcimen A**, **Rocca-Serra P**, **Sansone S-A**. ArrayExpress--a public
791          repository for microarray gene expression data at the EBI. [Online]. Nucleic
792          Acids Res 31: 68–71, 2003.
793          http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=165538&tool=pmc
794          entrez&rendertype=abstract [18 Jun. 2013].
795

796    8.    **Butte AJ**, **Kohane IS**. Mutual information relevance networks: functional
797          genomic clustering using pairwise entropy measurements. [Online]. Pac. Symp.
798          Biocomput.  http://www.ncbi.nlm.nih.gov/pubmed/10902190 [23 Jul. 2013].
799

800    9.    **Cassese A**, **Guindani M**, **Tadesse MG**, **Falciani F**, **Vannucci M**. A
801          HIERARCHICAL BAYESIAN MODEL FOR INFERENCE OF COPY
802          NUMBER VARIANTS AND THEIR ASSOCIATION TO GENE
803          EXPRESSION. Ann Appl Stat 8: 148–175, 2014.

804

10.  **Davidsen PK**, **Herbert JM**, **Antczak P**, **Clarke K**, **Ferrer E**, **Peinado VI**,
     **Gonzalez C**, **Roca J**, **Egginton S**, **Falciani F**. A systems biology approach
     reveals a link between systemic cytokines and skeletal muscle energy
     metabolism in a rodent smoking model and human COPD. Genome Med. .

11.  **Dennis G**, **Sherman BT**, **Hosack DA**, **Yang J**, **Gao W**, **Lane HC**, **Lempicki
     RA**. DAVID: Database for Annotation, Visualization, and Integrated
     Discovery. [Online]. Genome Biol 4: P3, 2003.
     http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3720094&tool=pm
     centrez&rendertype=abstract [5 Jun. 2014].

12.  **Du W**, **Elemento O**. Cancer systems biology: embracing complexity to
     develop better anticancer therapeutic strategies. Oncogene 34: 3215–25, 2015.

13.  **Faith JJ**, **Hayete B**, **Thaden JT**, **Mogno I**, **Wierzbowski J**, **Cottarel G**,
     **Kasif S**, **Collins JJ**, **Gardner TS**. Large-scale mapping and validation of
     Escherichia coli transcriptional regulation from a compendium of expression
     profiles. PLoS Biol 5: e8, 2007.

14.  **Falati S**, **Gross P**, **Merrill-Skoloff G**, **Furie BC**, **Furie B**. Real-time in vivo
     imaging of platelets, tissue factor and fibrin during arterial thrombus formation
     in the mouse. Nat Med 8: 1175–1180, 2002.

15.  **Finkel T**, **Deng C-X**, **Mostoslavsky R**. Recent progress in the biology and
     physiology of sirtuins. Nature 460: 587–91, 2009.

16.  **Gavaghan D**, **Garny A**, **Maini PK**, **Kohl P**. Mathematical models in
     physiology. Philos Trans A Math Phys Eng Sci 364: 1099–106, 2006.

17.  **Gehlenborg N**, **O'Donoghue SI**, **Baliga NS**, **Goesmann A**, **Hibbs MA**,
     **Kitano H**, **Kohlbacher O**, **Neuweger H**, **Schneider R**, **Tenenbaum D**, **Gavin
     A-C**. Visualization of omics data for systems biology. Nat Methods 7: S56–68,
     2010.

18.  **Goldovsky L**, **Cases I**, **Enright AJ**, **Ouzounis CA**. BioLayout(Java): versatile
     network visualisation of structural and functional relationships. [Online]. Appl
     Bioinformatics 4: 71–4, 2005. http://www.ncbi.nlm.nih.gov/pubmed/16000016
     [15 Jul. 2014].

19.  **Gomez-Cabrero D**, **Compte A**, **Tegner J**. Workflow for generating
     competing hypothesis from models with parameter uncertainty. Interface Focus
     1: 438–49, 2011.

847

20. **Gonzalez NC**, **Wood JG**. Alveolar hypoxia-induced systemic inflammation: what low PO(2) does and does not do. Adv Exp Med Biol 662: 27–32, 2010.

21. **Gupta R**, **Stincone A**, **Antczak P**, **Durant S**, **Bicknell R**, **Bikfalvi A**, **Falciani F**. A computational framework for gene regulatory network inference that combines multiple methods and datasets. BMC Syst Biol 5: 52, 2011.

22. **He W**, **Newman JC**, **Wang MZ**, **Ho L**, **Verdin E**. Mitochondrial sirtuins: regulators of protein acylation and metabolism. Trends Endocrinol Metab 23: 467–76, 2012.

23. **Hirose O**, **Yoshida R**, **Imoto S**, **Yamaguchi R**, **Higuchi T**, **Charnock-Jones DS**, **Print C**, **Miyano S**. Statistical inference of transcriptional module-based gene networks from time course gene expression profiles by using state space models. Bioinformatics 24: 932–42, 2008.

24. **Holloway K V**, **O'Gorman M**, **Woods P**, **Morton JP**, **Evans L**, **Cable NT**, **Goldspink DF**, **Burniston JG**. Proteomic investigation of changes in human vastus lateralis muscle in response to interval-exercise training. Proteomics 9: 5155–74, 2009.

25. **Huang DW**, **Sherman BT**, **Lempicki RA**. Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists. Nucleic Acids Res 37: 1–13, 2009.

26. **HUXLEY AF**. Muscle structure and theories of contraction. [Online]. Prog Biophys Biophys Chem 7: 255–318, 1957. http://www.ncbi.nlm.nih.gov/pubmed/13485191 [15 Jul. 2014].

27. **Ideker T**, **Ozier O**, **Schwikowski B**, **Siegel AF**. Discovering regulatory and signalling circuits in molecular interaction networks. [Online]. Bioinformatics 18 Suppl 1: S233–40, 2002. http://www.ncbi.nlm.nih.gov/pubmed/12169552 [3 Jun. 2014].

28. **Jing E**, **Emanuelli B**, **Hirschey MD**, **Boucher J**, **Lee KY**, **Lombard D**, **Verdin EM**, **Kahn CR**. Sirtuin-3 (Sirt3) regulates skeletal muscle metabolism and insulin signaling via altered mitochondrial oxidation and reactive oxygen species production. Proc Natl Acad Sci U S A 108: 14608–13, 2011.

29. **Jones P**, **Côté RG**, **Martens L**, **Quinn AF**, **Taylor CF**, **Derache W**, **Hermjakob H**, **Apweiler R**. PRIDE: a public repository of protein and peptide identifications for the proteomics community. Nucleic Acids Res 34: D659–63,

890     2006.
891

892   30.   **Joyner MJ**, **Pedersen BK**. Ten questions about systems biology. J Physiol
893         589: 1017–30, 2011.
894

895   31.   **Keller P**, **Vollaard NBJ**, **Gustafsson T**, **Gallagher IJ**, **Sundberg CJ**,
896         **Rankinen T**, **Britton SL**, **Bouchard C**, **Koch LG**, **Timmons JA**. A
897         transcriptional map of the impact of endurance exercise training on skeletal
898         muscle phenotype. J Appl Physiol 110: 46–59, 2011.
899

900   32.   **Kohl P**, **Noble D**. Systems biology and the virtual physiological human. Mol
901         Syst Biol 5: 292, 2009.
902

903   33.   **Krogh A**. The number and distribution of capillaries in muscles with
904         calculations of the oxygen pressure head necessary for supplying the tissue.
905         [Online]. J Physiol 52: 409–15, 1919.
906         http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1402716&tool=pm
907         centrez&rendertype=abstract [4 Jul. 2014].
908

909   34.   **Kupershmidt I**, **Su QJ**, **Grewal A**, **Sundaresh S**, **Halperin I**, **Flynn J**,
910         **Shekar M**, **Wang H**, **Park J**, **Cui W**, **Wall GD**, **Wisotzkey R**, **Alag S**,
911         **Akhtari S**, **Ronaghi M**. Ontology-based meta-analysis of global collections of
912         high-throughput public data. PLoS One 5: e13066, 2010.
913

914   35.   **van der Laan MJ**, **Pollard KS**. Hybrid clustering of gene expression data with
915         visualization and the bootstrap. J Stat Plan Inference 117: 275–303, 2003.
916

917   36.   **Maere S**, **Heymans K**, **Kuiper M**. BiNGO: a Cytoscape plugin to assess
918         overrepresentation of gene ontology categories in biological networks.
919         Bioinformatics 21: 3448–9, 2005.
920

921   37.   **Mah N**. A comparison of oligonucleotide and cDNA-based microarray
922         systems. Physiol Genomics 16: 361–370, 2004.
923

924   38.   **Margolin AA**, **Nemenman I**, **Basso K**, **Wiggins C**, **Stolovitzky G**, **Dalla**
925         **Favera R**, **Califano A**. ARACNE: an algorithm for the reconstruction of gene
926         regulatory networks in a mammalian cellular context. [Online]. BMC
927         Bioinformatics 7 Suppl 1: S7, 2006.
928         http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1810318&tool=pm
929         centrez&rendertype=abstract [19 Jul. 2012].
930

931   39.   **Mattson DL**. Functional Genomics. In: Integrative Physiology in the
932         Proteomics and Post-Genomics Age, edited by Walz W. Humana Press Inc.,

933    2005, p. 7–26.
934

935    40.   **Merico D**, **Gfeller D**, **Bader GD**. How to visually interpret biological data
936          using networks. Nat Biotechnol 27: 921–924, 2009.
937

938    41.   **Meyer PE**, **Kontos K**, **Lafitte F**, **Bontempi G**. Information-theoretic inference
939          of large transcriptional regulatory networks. EURASIP J. Bioinform. Syst. Biol.
940          ( January 2007). doi: 10.1155/2007/79879.
941

942    42.   **Neapolitan RE**. Learning Bayesian Networks. New Jersey: Pearson Prentice
943          Hall, 2004.
944

945    43.   **Noble D**. The Music of Life: Biology beyond genes. Oxford University Press,
946          2008.
947

948    44.   **Noble D**. Computational models of the heart and their use in assessing the
949          actions of drugs. [Online]. J Pharmacol Sci 107: 107–17, 2008.
950          http://www.ncbi.nlm.nih.gov/pubmed/18566519 [7 Jul. 2014].
951

952    45.   **Ogata H**, **Goto S**, **Sato K**, **Fujibuchi W**, **Bono H**, **Kanehisa M**. KEGG:
953          Kyoto Encyclopedia of Genes and Genomes. [Online]. Nucleic Acids Res 27:
954          29–34, 1999.
955          http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=148090&tool=pmc
956          entrez&rendertype=abstract [17 Jun. 2013].
957

958    46.   **Opgen-Rhein R**, **Strimmer K**. Learning causal networks from systems
959          biology time course data: an effective model selection procedure for the vector
960          autoregressive process. BMC Bioinformatics 8 Suppl 2: S3, 2007.
961

962    47.   **Ortega F**, **Sameith K**, **Turan N**, **Compton R**, **Trevino V**, **Vannucci M**,
963          **Falciani F**. Models and computational strategies linking physiological
964          response to molecular networks from large-scale data. Philos Trans A Math
965          Phys Eng Sci 366: 3067–89, 2008.
966

967    48.   **Perrin B-E**, **Ralaivola L**, **Mazurie A**, **Bottani S**, **Mallet J**, **d'Alché-Buc F**.
968          Gene networks inference using dynamic Bayesian networks. [Online].
969          Bioinformatics 19 Suppl 2: ii138–48, 2003.
970          http://www.ncbi.nlm.nih.gov/pubmed/14534183 [2 Nov. 2014].
971

972    49.   **Poultney CS**, **Greenfield A**, **Bonneau R**. Integrated inference and analysis of
973          regulatory networks from multi-level measurements. Methods Cell Biol 110:
974          19–56, 2012.
975

976 50. **Rabinovich RA**, **Bastos R**, **Ardite E**, **Llinàs L**, **Orozco-Levi M**, **Gea J**,
977     **Vilaró J**, **Barberà JA**, **Rodríguez-Roisin R**, **Fernández-Checa JC**, **Roca J**.
978     Mitochondrial dysfunction in COPD patients with low body mass index. Eur
979     Respir J 29: 643–50, 2007.
980

981 51. **Rangel C**, **Angus J**, **Ghahramani Z**, **Lioumi M**, **Sotheran E**, **Gaiba A**, **Wild**
982     **DL**, **Falciani F**. Modeling T-cell activation using gene expression profiling and
983     state-space models. Bioinformatics 20: 1361–1372, 2004.
984

985 52. **Ravasz E**, **Somera AL**, **Mongru DA**, **Oltvai ZN**, **Barabási AL**. Hierarchical
986     organization of modularity in metabolic networks. Science 297: 1551–5, 2002.
987

988 53. **Reinke C**, **Bevans-Fonti S**, **Drager LF**, **Shin M-K**, **Polotsky VY**. Effects of
989     different acute hypoxic regimens on tissue oxygen profiles and metabolic
990     outcomes. J Appl Physiol 111: 881–90, 2011.
991

992 54. **Segal E**, **Shapira M**, **Regev A**, **Pe'er D**, **Botstein D**, **Koller D**, **Friedman N**.
993     Module networks: identifying regulatory modules and their condition-specific
994     regulators from gene expression data. Nat Genet 34: 166–76, 2003.
995

996 55. **Shannon P**, **Markiel A**, **Ozier O**, **Baliga NS**, **Wang JT**, **Ramage D**, **Amin N**,
997     **Schwikowski B**, **Ideker T**. Cytoscape: a software environment for integrated
998     models of biomolecular interaction networks. Genome Res 13: 2498–504,
999     2003.
1000

1001 56. **De Smet R**, **Marchal K**. Advantages and limitations of current network
1002     inference methods. Nat Rev Microbiol 8: 717–29, 2010.
1003

1004 57. **Su G**, **Kuchinsky A**, **Morris JH**, **States DJ**, **Meng F**. GLay: community
1005     structure analysis of biological networks. Bioinformatics 26: 3135–7, 2010.
1006

1007 58. **Suzuki YJ**, **Carini M**, **Butterfield DA**. Protein carbonylation. Antioxid Redox
1008     Signal 12: 323–5, 2010.
1009

1010 59. **Timmons JA**, **Knudsen S**, **Rankinen T**, **Koch LG**, **Sarzynski M**, **Jensen T**,
1011     **Keller P**, **Scheele C**, **Vollaard NBJ**, **Nielsen S**, **Akerström T**, **MacDougald**
1012     **OA**, **Jansson E**, **Greenhaff PL**, **Tarnopolsky MA**, **van Loon LJC**, **Pedersen**
1013     **BK**, **Sundberg CJ**, **Wahlestedt C**, **Britton SL**, **Bouchard C**. Using molecular
1014     classification to predict gains in maximal aerobic capacity following endurance
1015     exercise training in humans. [Online]. J Appl Physiol 108: 1487–96, 2010.
1016     http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2886694&tool=pm
1017     centrez&rendertype=abstract [13 Jul. 2012].
1018

1019    60.    **Trevino V**, **Falciani F**. GALGO: an R package for multivariate variable
1020           selection using genetic algorithms. Bioinformatics 22: 1154–6, 2006.
1021

1022    61.    **Turan N**, **Kalko S**, **Stincone A**, **Clarke K**, **Sabah A**, **Howlett K**, **Curnow SJ**,
1023           **Rodriguez DA**, **Cascante M**, **O'Neill L**, **Egginton S**, **Roca J**, **Falciani F**. A
1024           systems biology approach identifies molecular networks defining skeletal
1025           muscle abnormalities in chronic obstructive pulmonary disease. PLoS Comput
1026           Biol 7: e1002129, 2011.
1027

1028    62.    **Werner HMJ**, **Mills GB**, **Ram PT**. Cancer Systems Biology: a peek into the
1029           future of patient care? Nat Rev Clin Oncol 11: 167–176, 2014.
1030

1031    63.    **Willmann G**. Transcriptional Regulation after Chronic Hypoxia Exposure in
1032           Skeletal Muscle. University of Cologne: 2013.
1033

1034    64.    **Yu J**, **Smith VA**, **Wang PP**, **Hartemink AJ**, **Jarvis ED**. Advances to
1035           Bayesian network inference for generating causal networks from observational
1036           biological data. Bioinformatics 20: 3594–603, 2004.
1037