

A vision for a biomedical cloud

■ R. L. Grossman^{1,2} & K. P. White^{1,2,3}

From the ¹Institute for Genomics and Systems Biology; ²Departments of Medicine; and ³Human Genetics, The University of Chicago, Chicago, IL, USA

Abstract. Grossman RL, White KP (Institute for Genomics and Systems Biology; The University of Chicago, Chicago) A vision for a biomedical cloud (Key Symposium). *J Intern Med* 2012; **271**: 122–130.

We present a vision for a Biomedical Cloud that draws on progress in the fields of Genomics, Systems Biology and biomedical data mining. The successful fusion of these areas will combine the use of biomarkers, genetic variants, and environmental variables to build predictive models that will drastically increase the specificity and timeliness of diagnosis for a wide range of common diseases, whilst delivering accurate predictions about the efficacy of treatment options. However, the amount of data being gener-

ated by each of these fields is staggering, as is the task of managing and analysing it. Adequate computing infrastructure needs to be developed to assemble, manage and mine the enormous and rapidly growing corpus of 'omics' data along with clinical information. We have now arrived at an intersection point between genome technology, cloud computing and biological data mining. This intersection point provides a launch pad for developing a globally applicable cloud computing platform capable of supporting a new paradigm of data intensive, cloud-enabled predictive medicine.

Keywords: cloud computing, data mining, electronic medical records, genomics, Systems Biology.

Technology drivers

Data intensive science

The power of computer processors doubles in less than 18 months, as does the capacity of computer storage discs. This exponential growth impacts not only the power of computers, but also the power of scientific instruments and has resulted in an exponentially growing amount of scientific data [1]. In this section, we look at three important changes that have occurred in the last decade regarding data: (i) the explosion in the amount of data *produced*, (ii) a fundamental change in how data are *managed and processed* and (iii) new algorithms for how data are *analysed* using data mining. Together, these changes are beginning to change biology into a *data intensive science*.

Ubiquitous sequencing: an explosion of data

Genomics reveals the 'parts lists' and genetic variations that compose each unique human genome. Genome technology has been producing DNA sequence data ever faster, cheaper and in tremendous volumes. The amount of DNA sequence in the public domain has been growing exponentially over the last two decades. Low-cost, high-throughput DNA sequencing is a disruptive technology that promises to have a major impact on our understand-

ing of biology and our treatment of human diseases. Within a few years, we will be in an era of ubiquitous sequencing, where genome sequencing will become routine for both research and clinical applications. Recently, Stein pointed out that from 1990 until 2004 sequencing output doubled approximately every 19 months, but from 2005 until the present the doubling rate decreased to 5 months because of the development of 'NextGen' sequencing technologies [2]. However, one can also extrapolate more generally based on Genbank data from 1990 to 2005. In Fig. 1, we fit a logarithmic function to these data and extrapolate the curve into the present day and beyond (this curve was initially calculated in 2007 by KPW for a presentation to the National Science Foundation). The 2011 estimates extrapolated from this curve are approximately 30 Terabases of finished genome that corresponds to 10 000 human genomes. In this sense, the impact of 'NextGen' sequencing technology was foreseeable. The total number of complete human genomes sequenced by the end of 2011 worldwide is, in fact, likely to be >10 000, in addition to genome sequencing of other species and 'partial genomes' from procedures such as exome capture sequencing, RNA sequencing and chromatin immunoprecipitation sequencing (ChIP-seq). Taking this curve as a tentative estimate of future world capacity, we can speculate that by 2015 more than a million human genomes will be sequenced. Storing

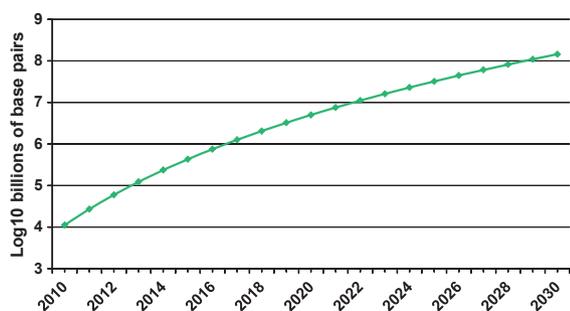


Fig. 1 Projected growth of DNA sequence data in the 21st century.

just the finished sequence will require several Petabytes of disc space, but it is reasonable to assume that dedicated data centers will be needed due to the substantial associated biomarker and phenotypic data with each genome.

Early in the third decade of the 21st century, we can speculate that the capacity to sequence up to 1 billion people will have been realized. Approximately 3000 PB of storage will be needed to house 1 billion human genomes. Are these speculations realistic? Perhaps a billion genomes is a bit far-fetched when most of the world's population does not have even basic access to medical professionals. However, with the cost of sequencing likely to drop below \$100 per genome in the next decade, and the expected increase in utility of genome information for diagnosing and treating diseases, we can expect genome sequencing to become an important aspect of health care in both developed and developing countries.

Importantly, genome sequencing is simply a baseline. Understanding both normal and disease states requires whole genome expression data at different time-points during normal conditions and during disease progression, as well as for different treatment conditions. Similarly, proteomic and metabolomics data monitoring the state of cells are growing increasingly important in identifying disease states and indicating treatments [3–6]. In the future, clinical trials will sample and generate a variety of 'omics' profiles for multiple time-points and different conditions (e.g. plus or minus treatment) during normal and diseased states. As the number of temporal points and treatment conditions grows, so too will the amount of data to be sifted. A person's state of health will be assessed in the light of his or her genome sequence, but also the cellular state as measured by hundreds of

thousands of analytes. Before these data can be compressed into diagnostic/prognostic panels that are representative of the 'whole', it will be necessary to analyse and mine them in their totality. Because of the enormous scale of data in this new era of ubiquitous sequencing, and concomitant developments of other data intensive 'omics' technologies, managing and analysing these data corpus is a formidable challenge.

Cloud computing

The last 10 years or so has seen an important set of advances in how data are managed. A decade ago, an important split occurred between the types of computing systems used by companies providing Internet services, such as Google, and those used by scientists and academics working in high-performance computing. By and large, researchers in high-performance computing were mainly concerned with writing specialized programs for high-end computers that were used for simulation. To maximize the performance of simulations, the systems are designed to minimize latency; to achieve this, specialized connections were used between the various processors, and libraries were used to pass messages, so that the different processors could exchange relevant information. Writing the code was difficult, and managing the (sometimes) large data that simulations produced was an after thought. Little, if any, of the work in high-performance computing has been relevant to the clinical enterprise, although it has had tremendous impact on the basic quantitative sciences.

In contrast, companies such as Google began collecting very large amounts of data (measured in Petabytes instead of Terabytes) and analysing it to optimize search results as well as the placement of text ads, which provide the bulk of Google's revenue. To manage and analyse these very large amounts of data (that literally fill a data centre), Google developed storage and compute services that scaled to all the computers that filled their data centres and were easy for their software engineers to program [7]. To keep costs low, which is by and large not a factor in high-performance computing, systems were designed to use commodity computers and to replicate data and redistribute workloads when the components failed (as they do frequently). Instead of designing systems to minimize latency, they designed systems to maximize the rate that data could be processed. This was done by making sure that there were very large numbers of discs and that each disc was matched with a processor, and, later with a core. Such a system

would be ideal for mining large-scale electronic medical records with genome data.

In 2006, Amazon introduced another important innovation in computing infrastructure that they called the Amazon Elastic Compute Cloud (EC2). The fundamental idea with EC2 was that a user could use a web portal and obtain one or more virtual machines that they could pay for by the hour. A virtual machine in this context is a process that appears to the user as an actual computer but is in fact one of several such virtual machines that are managed by a single physical computer. This provides two advantages: First, for many data centre usage patterns, it was most cost effective to use multiple virtual machines than a single actual physical machine. Secondly, it is easier to set up and tear down a virtual machine than a physical machine.

With EC2, the cost of one virtual machine for 100 h is the same as the cost of 100 virtual machines for 1 h. With this model, it suddenly became practical for scientists to perform computations on 100 machines for just the time required and let someone else manage the data centre infrastructure required to support it.

These types of systems developed by Google and Amazon are often referred to as cloud computing systems and promise to have just as large an impact on data intensive biology as they have had on business. In particular, as biology enters the era of large data, cloud computing systems will be one of the tools that will enable biologists to manage and analyse the data produced.

Biomedical data mining

Biomedical data mining seeks to connect phenotypic data to biomarker profiles and therapeutic treatments, with the goal of creating predictive models of disease detection, progression and therapeutic response.

During the last decade, data mining of biological data has become an increasingly important technique. Biological data mining includes mining a wide variety of biological data, including (i) mining genomic data (and data from other high-throughput technologies such as DNA sequencing, RNA expression, proteomic data, metabolomic data and small molecule screening), (ii) text mining of the biological literature, medical records, etc., and (iii) image mining across a number of modalities, including X-rays, functional MRI, new types of scanning microscopes.

Biomedical-infomics synthesis

Our thesis is that the explosion of genomic, proteomic and other 'omic' data, the ability of cloud computing to process and analyse data at the scale of a data centre, and new algorithms from data mining and Systems Biology create what might be called a biomedical-omics-infomatics synthesis. We will refer to this more simply as the biomedical-infomics synthesis.

As context for this new scientific paradigm, we discuss three major shifts in biological thinking and approach over the last century that also constituted up-endings of the existing paradigms of their time: the NeoDarwinian synthesis, the Molecular Biology revolution and Systems Biology. We argue that a biomedical-infomatics synthesis is an important emerging component of Systems Biology.

NeoDarwinian synthesis

In the middle part of the last century, the field of Genetics underwent a transformation known as the NeoDarwinian synthesis [8]. During this period, there was a fusion of the theoretical genetics that explained the inheritance behaviour of (mostly) individual genes with simple allelic variants and the evolutionary ideas of Darwin and his intellectual progeny. Remarkably, the NeoDarwinian synthesis happened largely in the absence of the understanding of DNA as the heritable material of all biology and without any but the vaguest notion of the molecular nature of genes. The ramifications of this mid-century synthesis have been immense, leading to improvements in crops and livestock whilst setting a foundation for interpretation and understanding of the molecular basis of life in the second half of the 20th century as biologists turned to unravelling the basis of DNA and the products it encodes.

Molecular biology revolution

This subsequent Molecular biology revolution was catalysed by the discovery of the double helix in 1953 and subsequently dominated biomedicine and much of biological thinking in general [9, 10]. Great benefits to society were stimulated by the Molecular biology revolution as well, including the development of drugs such as synthetic insulin, humanized antibodies directed against tumours, modification of crops for pest resistance or increased yield and many other extremely meaningful contributions that relied on the understanding or engineering of one or only a few genes at a time.

However, as geneticists have known for 100 years, most of biology is more complex and involves many genes acting in the context of heterogeneous environments. The molecular biologists for decades chose to ignore this unpleasant reality, going to great lengths to eliminate genetic variability in their model systems, and the geneticists who did not become molecular biologists were largely sidelined academically and in biomedicine or focused on practical pursuits such as improving agriculture. Only in the last 10–15 years has the genetics of complex traits moved back into the mainstream of biology, enabled by the advent of genomic technologies that act as the first instrument that can allow scientists to see the totality of variation that contributes to complex traits. The major challenge in all of biology at the beginning of this century is to figure out how complex traits work at a molecular level.

Systems Biology

With a renewed appreciation for the complexity of biological systems, a modern version of a field known as Systems Biology has affirmed that understanding of complex traits requires their integrated study at the molecular, cellular and organismal levels. Whilst the field has its historical roots largely in metabolic flux analysis, neuronal modelling and bioengineering, during the last 10 years, Systems Biology has come to encompass much of modern biology and professes to usher a new era where biological theory and experiment become unified [11]. It presently connotes two major areas of investigation.

First, systems biologists use genomic scale data to analyse molecular networks, typically by integrating multiple heterogeneous data types that represent different aspects of cellular biology and genetics, and to make predictions about network structures and which network substructures (and individual genes or gene products) are crucial for a given phenotype being analysed. For example, gene expression networks that integrate RNA expression profiling, transcriptional factor binding to the genome and other data types have now been generated for a vast breadth of traits that range from yeast metabolism [12, 13], to embryonic development of fruit flies [14–16], to human cancers [17, 18] and to dozens of other complex traits.

The second major area of modern Systems Biology echoes its modelling roots, focusing on discrete subsystems where enough data have been gathered to build predictive models that specify nontrivial outcomes of perturbing a given network. This area too

has been applied to a wide range of applications that span from modelling the stochastic behaviour of microbial chemosensing [19], to the aforementioned transcriptional networks controlling embryonic development in flies [20, 21] and to phosphorylation-based signalling networks such as those activated by MAP kinases and receptor tyrosine kinases in a wide range of cancers [22–24].

But Systems Biology, as it is often formulated, also seeks to understand emergent properties beyond the structures of biological networks and information flow within them. In fact many of the core themes of Systems Biology are identical to the themes identified many years ago by the geneticists that led the Neo-Darwinian synthesis, their contemporaries and predecessors. These properties are at the root of biology and include concepts such as emergence of three dimensional structure of cells, tissues and organisms from the simple materials of inheritance, cellular and developmental robustness, modularity of biological systems, group behaviours (such as schooling in fish or swarming in bees) and the process of organic evolution. These concepts, all representing emergent properties of complex systems, are driving much contemporary research in Systems Biology.

Sources of data for the biomedical-infomics synthesis

Systems Biology provides a natural conceptual framework for launching a biomedical-infomics synthesis. Systems Biology is the modern intellectual home for an integrated view of Biology, undivided into its dozens of subfields and specialties; genomic technologies play a major role; and Systems Biology draws from almost every scientific discipline to address fundamental problems, with particular avidity for computing and engineering. Most importantly for our present thesis, Systems Biology is the site of information integration about biological systems, and the field is extremely active. As mentioned earlier, data production will grow not only because of genome sequencing but also because of sequencing and otherwise measuring gene products under many different conditions. Here we will discuss several major types of data that may be used to generate predictive networks.

Experimental perturbations

The perturbation approach of hypothesis testing is the basis of modern experimental biology. A system's output is measured under different experimental or naturally occurring conditions (normal vs. disease,

mutant vs. wild type, hormone vs. control treatment, etc.). In the last decade, a transition has occurred where investigations have gone from measuring one variable (e.g. a gene product) at a time to measuring the output of the entire genome – in other words, we have transitioned to assaying the state of the entire system. Simultaneously, through miniaturization technologies and development of high-throughput screening approaches, perturbation analyses themselves have seen the same scale of expansion whereby the entire contents of genomes are routinely perturbed, and then, traits (phenotypes) are measured. For example, data matrices with 20 000 perturbations \times 10 GB of gene expression and corresponding measurements are now foreseeable (200 PB). Already it is practical to generate a matrix with hundreds to several thousand perturbations and whole genome measurements. With this transition to larger and larger matrices of ‘omics’ scale data, it becomes essential to employ statistical and computational methods to determine which variables are important for a given phenotype being analysed. More advanced studies have relied on building networks and distilling testable hypotheses from those networks, for example using the types of probabilistic algorithms aforementioned [13, 14, 25–28].

Genetic associations

The second successful experimental approach in Systems Biology has relied on genetic associations. This approach, descended from the same genetic thinking that drove the complex trait geneticists in the 20th century, takes full advantage of genomic technologies. A powerful implementation of this approach is to generate gene expression data alongside genotypic data to determine which genetic variants are affecting which genes’ expression. This is known as the expression quantitative trait locus (eQTL) approach [29–33]. By comparing populations that are sick versus healthy, this eQTL approach can identify genes that are hubs in networks that are associated with the trait.

More generally, understanding associations between whole genomic variants and phenotypes across populations will be an important source of data for the biomedical-infomics synthesis.

Evolutionary conservation

A third approach that is potentially extremely powerful but is just beginning to show its effectiveness in the context of Systems Biology is the evolutionary, or comparative, approach. As an example, perturbation

analysis was used in the model organism *Drosophila* to map a network that controls early developmental pattern formation in the embryo. Using the knowledge that many components of this network are evolutionarily conserved in humans and involved in diseases such as cancer, the human counterparts were screened for their disease association [14]. A key conserved factor (predicted by a network centrality metric) was associated with human kidney cancer, and subsequent studies have verified that this gene product can cause cancer in mice and that inhibiting it can kill cancer cells, thus making it a promising drug target for a disease that currently is resistant to both chemotherapy and radiation (Li *et al.*, submitted). More generally, using data gleaned from public databases, networks in model organisms can be built that map to orthologous networks associated with human diseases. Such networks can help to identify novel candidate genes involved in human disease processes or suggest genes for study in model organisms that could yield human disease insights [28].

Medical text

Today, there are trillions of pages of scholarly text available in the world libraries, and although science-focused text mining is a formidable intellectual challenge, it is beginning to be used to extract new discoveries from this stored text.

The typical stages of text mining include identification of named entities (gene and disease names, organizations and geographic locations), capturing relations amongst the named entities (such as a reported association between a genetic polymorphism and a disease or interaction between a protein and a small molecule) and then computational reasoning to construct complex semantic networks from these entities and their relations. Text mining at this scale is challenging because of the sheer volume of the data and because of the difficulty extracting what are often quite complex and subtle assertions from the data.

Today, as new experimental data are generated, the process of analysing the data and deriving text assertions summarizing it, such as identifying the specific genetic polymorphism associated with a disease, is largely manual. Manual processes like these do not scale to the amount of data that will be produced to capture millions of variants that must be analysed from thousands or millions of genomes. Biological data mining is beginning to automate this process, which is essential because deriving these types of text-based assertions from newly generated experimental

data provides access to historical context for new data analysis, to previously formulated and supported, untested, or rejected hypotheses, and to legacy observations from multiple research communities.

Electronic medical records

Another important source of data for the biomedical-infomics synthesis is electronic medical records (EMR). More and more hospitals and medical research centres are implementing EMR, in part because of financial advantages they offer. Once an EMR system is in place, it is natural to create what is usually called a clinical research data warehouse, so that multiple years of EMRs over entire populations of patients can be analysed. Patient data for which consent has been provided can be analysed along with genomic data. This allows phenotype data to be correlated with genomic data for patient populations that sometimes hundreds or thousands in size.

Large-scale genomic studies

Finally, a very important source of data for the biomedical-infomics synthesis is the increasing number of large-scale studies. For example, the Encyclopaedia of DNA Elements (ENCODE) project is mapping functionality in the human genome [34, 35]. The 1000 Genomes project has mapped already multiple thousands of genomes ([36], <http://www.1000genomes.org/>, <http://www.1000genomes.org/>). The Cancer Genome Atlas is identifying the variation associated with more than a dozen types of cancer (<http://cancergenome.nih.gov/>, <http://cancergenome.nih.gov/>). Various consortiums are sequencing the genomes of patients afflicted with most major complex diseases, and in addition, a large-scale effort is underway to map the microbial genomes associated with humans (The Human Microbiome Project) (<https://commonfund.nih.gov/hmp/>). Finally, sequencing of wild and domestic species abounds, for example with the initiation of the 10,000 vertebrates sequencing project (<http://genome10k.soe.ucsc.edu/>). The end result, in terms of data, is that tens of thousands of human genomes and the genomes of our commensals, our parasites, and even our pets are flooding databases around the world.

A biomedical cloud

Although Petabytes of genomic, biological, clinical, text and related data are available for downloading to-

day, there is no way currently to compute over *all* of these data to make discoveries, and, more importantly, there is no conceptual framework to integrate all this data.

It is important to note that just over a decade ago, the same could have been said about the all the data available from web sites. But during this period, companies such as Google filled data centres with this data, developed software to manage and process this data, and then computed over *all* of these data to improve algorithms for search, online advertising and related areas.

We argue that an interconnected network of data centre scale facilities (loosely speaking 'clouds') with the appropriate security architecture and a rich set of secure intercloud services is the proper foundation for the biomedical-infomics synthesis. Call this the *Biomedical Cloud*. It is an example of a community cloud [37]. It could be filled with *all* publicly available data relevant to biology, medicine and health care. Like the data in commercial search engines, such a repository would be accessible to individuals via personal devices, with the compute intensive operations being performed in data centres and associated high-performance computing facilities.

Moreover, secure private clouds and clinical research data warehouses located at medical centres and hospitals containing EMRs, and other data with protected health information (PHI) information could be enriched with data from the community Biomedical Cloud. Figure 2 illustrates at a high level how Biomedical Clouds might be deployed to interact with a variety of users, including medical professionals, researchers and the general public.

The following eight requirements seem to be necessary for a Biomedical Cloud:

- 1 *Appropriate security.* A Biomedical Cloud would contain a mixture of data, some of it public, some of it restricted to research collaborations and some of it restricted because it is human genome data or contains PHI information. At one extreme, for the most restricted data, specialized secure private clouds will be required where the clouds are designed so that the data remains within the required organization and all appropriate regulations are followed. At the other extreme, public clouds will be used to analyse and distribute publically available data.

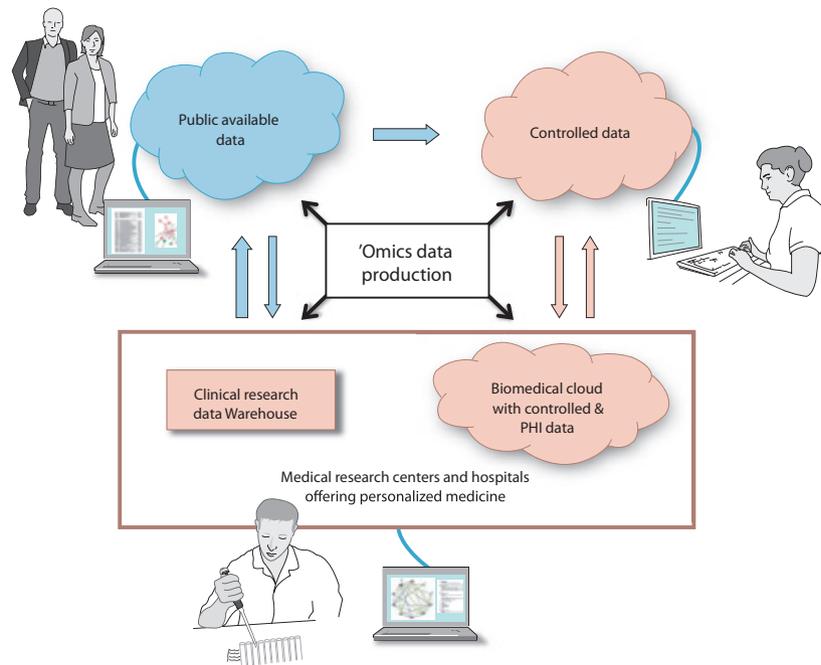


Fig. 2 Biomedical Clouds will interact with medical centres and hospitals providing personalized medicine, as well as researchers and the general public. The blue cloud represents publicly accessible data, whilst pink is used for data that are controlled and have limited access, because of the presence of protected health information, whole human genomes and similar data. Data production may come from many sources, including private and public. There will be multiple Biomedical Clouds associated with medical centres and hospitals that contain controlled data, and these private Biomedical Clouds will be able to ingest data from public Biomedical Clouds. Private clouds may also provide data for research studies via de-identified data sets to public clouds and controlled data to clouds designed to hold such data.

- 2 *Secure communications with private clouds.* Secure private clouds and data warehouses containing human data and PHI data will be located at medical research centres and contain data critical to the biomedical-infomics synthesis. With the collections of the necessary consents, approval of suitable protocols and appropriate secure communications, data in these systems will be analysed along with data in the Biomedical Cloud.
- 3 *On-demand and scalable storage.* A Biomedical Cloud should scale, so that it can manage and archive *all* the data relevant to the biomedical-infomics synthesis.
- 4 *On-demand and scalable analysis.* A Biomedical Cloud should scale, so that it could analyse all the data relevant to the biomedical-infomics synthesis without moving the data out of the cloud.
- 5 *Scalable ingestion of data.* A Biomedical Cloud should support the scalable ingestion of biological, medical and healthcare data, including the ingestion of data from next-generation sequencers, genomics databases and other clouds.
- 6 *Support data liberation.* A Biomedical Cloud should provide both long-term storage for data as well as a mechanism for exporting data, so that it can be moved to another cloud or facility.
- 7 *Peer with other private and community clouds.* A Biomedical Cloud should interoperate, preferably via peering, with other clouds, so that complex analyses can be carried out using data that span multiple clouds. By peering, we mean that the two clouds can exchange data without paying a charge per GB of data transported.
- 8 *Peer with public clouds.* A Biomedical Cloud should interoperate, preferably via peering, with public clouds, so that an investigator can analyse data within the genomic cloud or using public clouds.

With a Biomedical Cloud satisfying these requirements, all the Petabytes of available data relevant to the biomedical-informatics synthesis could be collocated in one place, and more importantly, algorithms could be used to integrate and process these data on a continuous basis. With the proper architecture, appropriate security and scalable algorithms, Terabytes of new data would be added to the Biomedical Cloud each day, processed each night and available for search each morning. In this way, we would continuously update the connections between phenotypic data, biomarker profiles and therapeutic treatments, with the goal of creating predictive models of disease detection, progression and therapeutic response. The successful fusion of these areas will create a predictive matrix of biomarkers, genetic variants and environmental variables that will drastically increase the specificity and timeliness of diagnosis for a wide range of common diseases, whilst delivering accurate predictions about the efficacy of treatment options. Of course, there are many logistical and technical challenges to be surmounted if such a Biomedical Cloud will come into existence. How will personal 'omics' and medical data be protected from being decoded in such an environment? How will physicians make best use of such a powerful resource? However, some version of an organically evolving biomedical informatics machine is likely to arise in the not too distant future. The ingredients are all in place.

Conflict of interest statement

The authors have no conflicts of interest to disclose.

References

- 1 Szalay A, Gray J. Science in an exponential world. *Nature* 2006; **440**: 413–4.
- 2 Stein LD. The case for cloud computing in genome informatics. *Genome Biol* 2010; **11**: 207.
- 3 Ponten F, Jirstrom K, Uhlen M. The Human Protein Atlas—a tool for pathology. *J Pathol* 2008; **216**: 387–93.
- 4 Uhlen M, Oksvold P, Fagerberg L *et al*. Towards a knowledge-based Human Protein Atlas. *Nat Biotechnol* 2010; **28**: 1248–50.
- 5 Nicholson JK, Lindon JC. Systems biology: metabonomics. *Nature* 2008; **455**: 1054–6.
- 6 Rubakhin SS, Romanova EV, Nemes P, Sweedler JV. Profiling metabolites and peptides in single cells. *Nat Methods* 2011; **8**(4 Suppl.): S20–9.
- 7 Dean J, Ghemawat S. MapReduce: simplified data processing on large clusters. *Commun ACM* 2008; **51**: 107–13.
- 8 Dobzhansky T. *Genetics and the Origin of Species*, 3rd edn. New York, NY: Columbia University Press, 1951.
- 9 Watson JD, Crick FH. Genetical implications of the structure of deoxyribonucleic acid. *Nature* 1953; **171**: 964–7.
- 10 Watson JD, Crick FH. Molecular structure of nucleic acids; a structure for deoxyribose nucleic acid. *Nature* 1953; **171**: 737–8.
- 11 Ideker T, Galitski T, Hood L. A new approach to decoding life: systems biology. *Annu Rev Genomics Hum Genet* 2001; **2**: 343–72.
- 12 Herrgard MJ, Swainston N, Dobson P *et al*. A consensus yeast metabolic network reconstruction obtained from a community approach to systems biology. *Nat Biotechnol* 2008; **26**: 1155–60.
- 13 Ideker T, Thorsson V, Ranish JA *et al*. Integrated genomic and proteomic analyses of a systematically perturbed metabolic network. *Science* 2001; **292**: 929–34.
- 14 Liu J, Ghanim M, Xue L *et al*. Analysis of *Drosophila* segmentation network identifies a JNK pathway factor overexpressed in kidney cancer. *Science* 2009; **323**: 1218–22.
- 15 Roy S, Ernst J, Kharchenko PV *et al*. Identification of functional elements and regulatory circuits by *Drosophila* modENCODE. *Science* 2010; **330**: 1787–97.
- 16 Zinzen RP, Girardot C, Gagneur J, Braun M, Furlong EE. Combinatorial binding predicts spatio-temporal cis-regulatory activity. *Nature* 2009; **462**: 65–70.
- 17 Lamb J, Crawford ED, Peck D *et al*. The Connectivity Map: using gene-expression signatures to connect small molecules, genes, and disease. *Science* 2006; **313**: 1929–35.
- 18 Hua S, Kallen CB, Dhar R *et al*. Genomic analysis of estrogen cascade reveals histone variant H2A.Z associated with breast cancer progression. *Mol Syst Biol* 2008; **4**: 188.
- 19 Korobkova E, Emonet T, Vilar JM, Shimizu TS, Cluzel P. From molecular noise to behavioural variability in a single bacterium. *Nature* 2004; **428**: 574–8.
- 20 Jaeger J, Surkova S, Blagov M *et al*. Dynamic control of positional information in the early *Drosophila* embryo. *Nature* 2004; **430**: 368–71.
- 21 Janssens H, Hou S, Jaeger J *et al*. Quantitative and predictive model of transcriptional control of the *Drosophila melanogaster* even-skipped gene. *Nat Genet* 2006; **38**: 1159–65.
- 22 Ciaccio MF, Wagner JP, Chuu CP, Lauffenburger DA, Jones RB. Systems analysis of EGF receptor signaling dynamics with microwestern arrays. *Nat Methods* 2010; **7**: 148–55.
- 23 Jones RB, Gordus A, Krall JA, MacBeath G. A quantitative protein interaction network for the ErbB receptors using protein microarrays. *Nature* 2006; **439**: 168–74.
- 24 Morris MK, Saez-Rodriguez J, Clarke DC, Sorger PK, Lauffenburger DA. Training signaling pathway maps to biochemical data with constrained fuzzy logic: quantitative analysis of liver cell responses to inflammatory stimuli. *PLoS Comput Biol* 2011; **7**: e1001099.
- 25 Amit I, Garber M, Chevrier N *et al*. Unbiased reconstruction of a mammalian transcriptional network mediating pathogen responses. *Science* 2009; **326**: 257–63.
- 26 Costanzo M, Baryshnikova A, Bellay J *et al*. The genetic landscape of a cell. *Science* 2010; **327**: 425–31.
- 27 Krogan NJ, Cagney G, Yu H *et al*. Global landscape of protein complexes in the yeast *Saccharomyces cerevisiae*. *Nature* 2006; **440**: 637–43.
- 28 McGary KL, Park TJ, Woods JO, Cha HJ, Wallingford JB, Marcotte EM. Systematic discovery of nonobvious human disease models through orthologous phenotypes. *Proc Natl Acad Sci USA* 2010; **107**: 6544–9.

- 29 Stranger BE, Forrest MS, Dunning M *et al*. Relative impact of nucleotide and copy number variation on gene expression phenotypes. *Science* 2007; **315**: 848–53.
- 30 Stranger BE, Nica AC, Forrest MS *et al*. Population genomics of human gene expression. *Nat Genet* 2007; **39**: 1217–24.
- 31 Brem RB, Yvert G, Clinton R, Kruglyak L. Genetic dissection of transcriptional regulation in budding yeast. *Science* 2002; **296**: 752–5.
- 32 Morley M, Molony CM, Weber TM *et al*. Genetic analysis of genome-wide variation in human gene expression. *Nature* 2004; **430**: 743–7.
- 33 Schadt EE, Lamb J, Yang X *et al*. An integrative genomics approach to infer causal associations between gene expression and disease. *Nat Genet* 2005; **37**: 710–7.
- 34 Birney E, Stamatoyannopoulos JA, Dutta A *et al*. Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* 2007; **447**: 799–816.
- 35 Myers RM, Stamatoyannopoulos J, Snyder M *et al*. A user's guide to the encyclopedia of DNA elements (ENCODE). *PLoS Biol* 2011; **9**: e1001046.
- 36 1000GenomesProject. A map of human genome variation from population-scale sequencing. *Nature* 2010; **467**: 1061–73.
- 37 Mell P, Grance T. *The NIST Definition of Cloud Computing (Draft): Recommendations of the National Institute of Standards and Technology*. Gaithersburg, MD: National Institute of Standards and Technology, 2011.

Correspondence: Robert L. Grossman and Kevin P. White, Institute for Genomics and Systems Biology, KCBBD 10100, The University of Chicago, Chicago, IL 60637, USA.
(fax: (773) 834 2877; e-mail: kpwhite@igsb.org and grossman@igsb.org). ■