**Figure 1.** Algorithm for diagnosis of interstitial lung disease (ILD) in non-small cell lung cancer (NSCLC) patients.

As an example, in our work with gefitinib, samples were taken after obtaining informed consent from a nested case-control study, i.e., a case-control study performed within a prospective pharmacoepidemiological cohort of several thousand patients with advanced or recurring NSCLC who had received at least one prior chemotherapy regimen, and who were to be treated with gefitinib or chemotherapy. The main objective of this study was to measure the relative risk of ILD in Japanese patients with NSCLC using gefitinib compared with conventional therapy, with the associated aims of determining the incidence rate of ILD in late stage NSCLC patients and the principal risk factors for this complication.

Central to both the case-control study and the proteomics analysis was the use of internationally agreed criteria for the diagnosis of ILD and an algorithm of diagnostic tests to exclude alternative diseases.[28] Principal investigators in the study were asked to assess all patients for possible ILD using the diagnostic algorithm (Figure 1). Two case review boards of experts from oncology, radiology, and pulmonary medicine were set up to independently establish a consistent final diagnosis of ILD. In addition, extensive standard clinical and demographic risk factor data were collected on all registered cases and controls.

This degree of rigor in establishing accurate phenotypic diagnosis is critical to develop a robust and reliable personal-

ized medicine test, as inaccuracies at this stage will affect all subsequent data analyses. The availability of clinical and risk factor data, and a rigorous epidemiological study design setting for the collection of proteomics samples is also of great value to fine-tune the statistical analysis.

## Is Proteomics Ready for Personalized Medicine Applications?

**The Human Proteome Map in Plasma.** The impetus to develop personalized medicine based on blood samples has encouraged proteomic profiling that identifies individual proteins and multiple "fingerprint" protein patterns. A remaining limitation has been the lack of integration of the technology of protein separation with bioinformatics and statistical methods. Extensive national and international[29,30] collaborations are being implemented to address some of these needs. An important component in this development is the Human Proteome Organization (HUPO; www.HUPO.org), a scientific consortium that supports various programmes to map the proteins expressed in various human tissues, disease states, etc.[31-33] One of these is the Plasma Proteome initiative started in 2002, aiming to annotate and catalog the many thousands of proteins and peptides[34-37] of the human plasma proteome. Recently results from the pilot phase with 35 collaborating laboratories from 13 countries[38-42] and multiple analytical

235

groups were made publicly available on the Internet (www.bioinformatics.med.umich.edu/hupo/ppp; www.ebi.ac.uk/pride). The combined efforts have generated 15 710 different MS/MS datasets that were linked to the International Protein Index (IPI) protein IDs, and an integration algorithm applied to multiple matches of peptide sequences yielded 9504 IPI proteins identified with one or more peptides[40] and characterized by Gene Ontology, InterPro, Novartis Atlas, and OMIM. Such advances provide an important platform for transforming proteomics from a technology to a useful biomarker tool applicable to personalized medicine.

**Protein Analysis in Blood—The Methods.** With respect to automated studies, multidimensional chromatography is the main technology used for protein analysis in blood. It is coupled to mass spectrometry either by electrospray ionization (ESI) for analysis in solution or matrix assisted laser desorption/ionization (MALDI) in solid phase applications.[39,41,43-47] Alternatively, ion-trap mass spectrometers are gaining recognition for high-throughput sequencing.[46,48-53] Linking a Fourier transform ion cyclotrone resonance (FTICR) unit to the linear trap can increase the resolution profoundly,[36,54-56] one of several novel principles for strengthening the assignment of protein annotations with the most commonly used protein search engines.[36,47,54-61] For protein annotation, the recent development of a human protein reference database complements these technologies.[61] Studies of protein expression in a variety of biological compartments ranging from sub-cellular to whole organisms have been undertaken with these analytic approaches.[62-70] Some key findings from the HUPO initiatives that impact on methodology include:

• For studies using blood samples, plasma rather than serum is preferred, with ethylenediaminetetraacetic acid (EDTA) as an anticoagulant.[40]

• The abundant proteins in plasma should be depleted prior to analysis.[40]

• Acceptance of protein annotation, i.e., accepted protein identities[39,40] should use standard criteria. These include having two identified peptide sequences from each protein, both with a statistical significance score high enough to ensure a correct sequence confirmation when compared with the corresponding gene sequence entity.[39]

Despite the advances in methodology, important hurdles to using proteomics in a personalized medicine context remain.

**Protein Expression Analysis in Blood—Some Important Hurdles.** Although protein profiling technology is highly automated and interfaced with database search engines to relate peptide sequences to protein identities and function,[39,40] there are many practical reasons why determining the relative abundance of proteins relevant for prediction purposes is difficult:

• About 90% of proteins are believed to be present only in low copy numbers, i.e., at medium and low abundance levels.[49]

• There can be variation both in the quantity and form of protein expression within normal physiological function.

• Between 300 000 and 3 million human protein species exist as direct gene products or post-translational modifications.[44]

• The relative abundance of the post-translational modifications occurring within the cell is called a Cell-Protein-Index Number (CPIN).[29,30] As an example, if one considers that there are 30 types of phosphorylation variants of a single phosphoprotein, and a hundred possible fold forms of glycosylation of a single glycoprotein, the theoretical CPIN varies considerably depending on the sample complexity.

• The dynamic range of protein expression within cells, between levels of most and least abundant proteins, is in the order of $10^8$–$10^{10}$.[34-36]

• In a typical clinical proteomics study the total cellular protein material in a sample seldom exceeds 10–20 milligrams. Therefore, the least abundant proteins would be present at starting levels not exceeding picograms.

• Recent studies use technology that can identify several thousand proteins in plasma samples,[29] but this still probably only represents a small fraction of the intermediate and processed protein forms. This is due to the current limitation of mass spectrometry not being able to ionize all amino acid sequences and protein modifications with equal efficiency. In most situations, a limited region of the full length protein is sequence annotated.

• The detection of differences in protein expression between groups of interest (e.g., cases and controls) takes place against a background of high variation between individuals within a group, within individuals over time and possible analytic run-to-run variation. Any method used to address this hurdle (which will involve "alignment" for spectral methods) directly impacts the ability to find good protein biomarkers.

Beyond the hurdles above, the fundamental challenge of protein biomarkers is to link the relative abundance of single markers or a fingerprint to clinically important biological processes based on some direct or indirect cause-effect link[29] related to normal or aberrant biological pathways.[47,49] In the following sections, we present the approach used for the identification of protein biomarkers potentially associated with development of ILD in NSCLC patients within the case-control study used as our motivating example. We build on the foundations described above and introduce further analytic developments and ideas relating to proteomic data generation, assaying and alignment to build a proteomics toolkit that can be applied today for personalized medicine approaches.

## A State of the Art Clinical Biomarker Analysis System

In the previous section, we described several challenges in proteomic analysis. Here we describe a system and analysis approaches that we have successfully implemented to address some of these issues.

**The Components of the Analysis System.** The analysis system (Figure 2) uses liquid chromatography-based high-resolution separation of peptides with an interface to tandem MS/MS, a technology which has been attracting great attention as the "shotgun" method of proteome analysis.[44,60-70] With this technology, after depletion of albumin and immunoglobulin G (IgG), all extracted plasma proteins are digested into their specific peptide components by proteolytic enzyme treatment.

The generated peptides are subjected to capillary reverse-phase submicro- to micro-flow liquid chromatography (capillary RP $\mu$LC), separated by retention times due to their physicochemical properties, and then detected and sequenced by a linear ion-trap tandem mass spectrometer[71] (LTQ, Thermo Fisher Scientific, San Jose, CA) interfaced with a spray needle tip for ESI of peptides.[70] A two-dimensional quadrupole ion trap mass spectrometer[71] is used, operated in a data-dependent acquisition mode with operational $m/z$ range limits set at 450–2000 (Figure 3, graphs A and B). Automatic switching to MS/MS acquisition mode is made in 1-second scanning cycles, controlled by the XCalibur software. The actual differences between annotated peptide fragment peaks shown in Figure 3, graph C, correspond to the amino acid residue mass, i.e.,
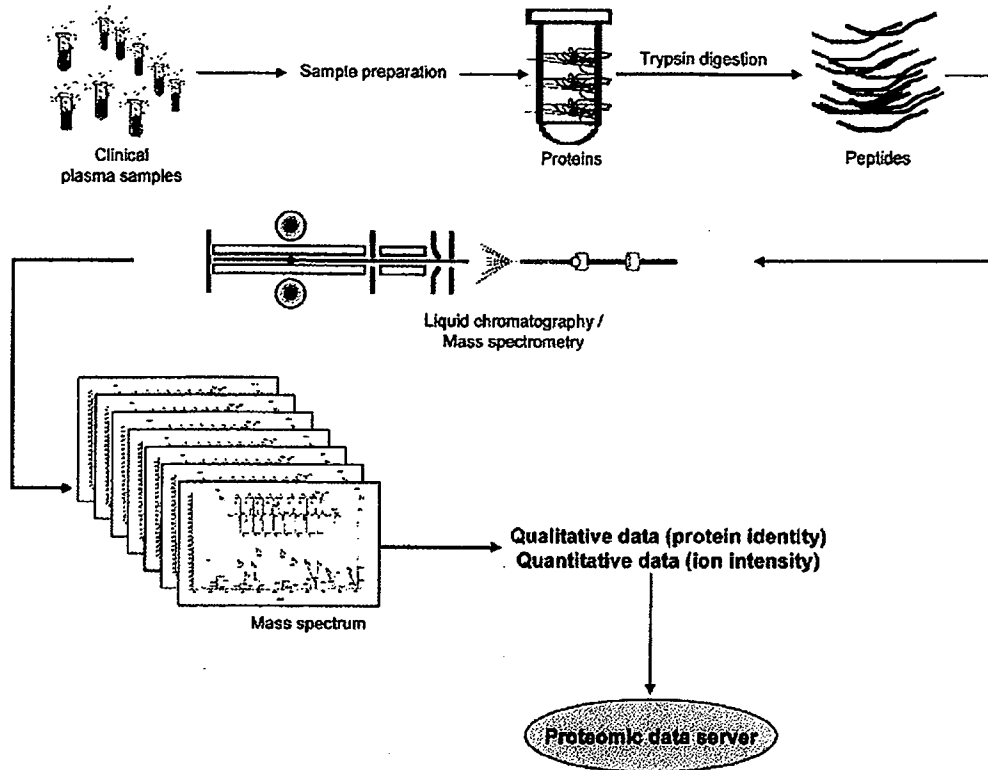
**Figure 2.** Schematic illustration of the clinical proteomics screening process.
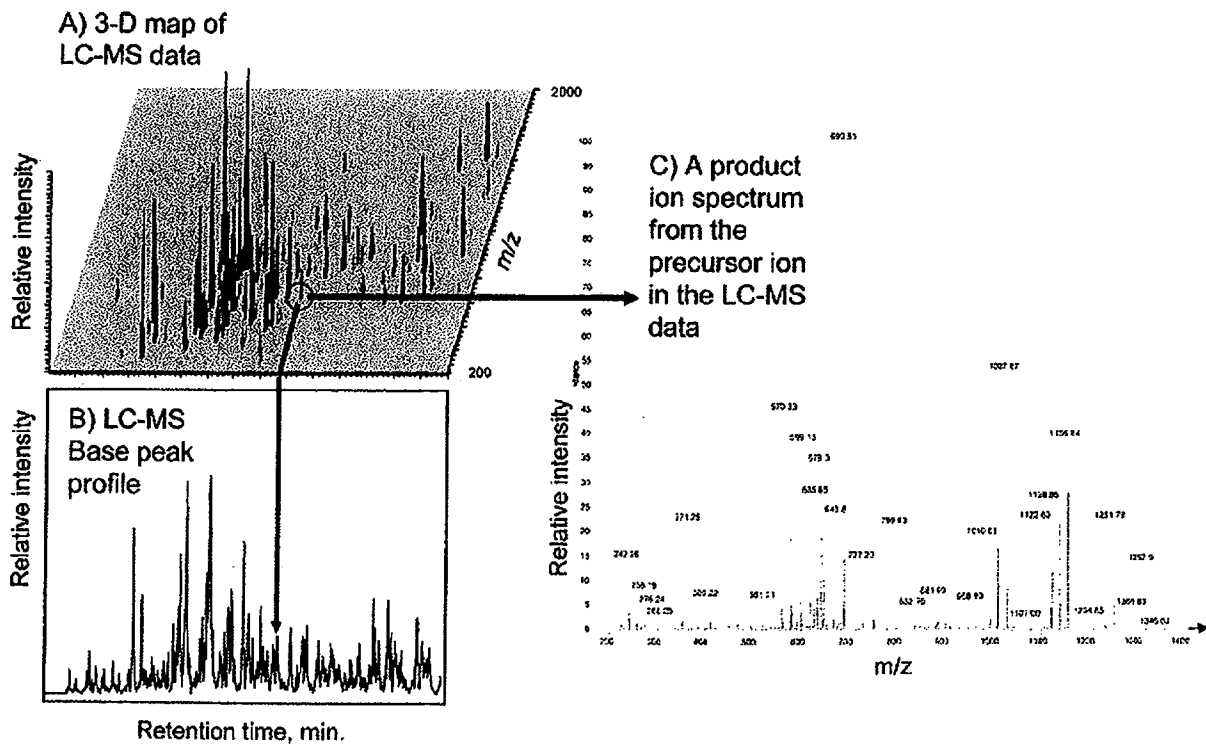


**Figure 3.** Profile of LC–MS data: (a) the three-dimensional view of LC-MS data, (b) the base-peak mass chromatogram, and (c) a product ion spectrum measured for a precursor ion in data-dependent acquisition mode (with MS acquisition operational *m/z* range set at 450–2000).
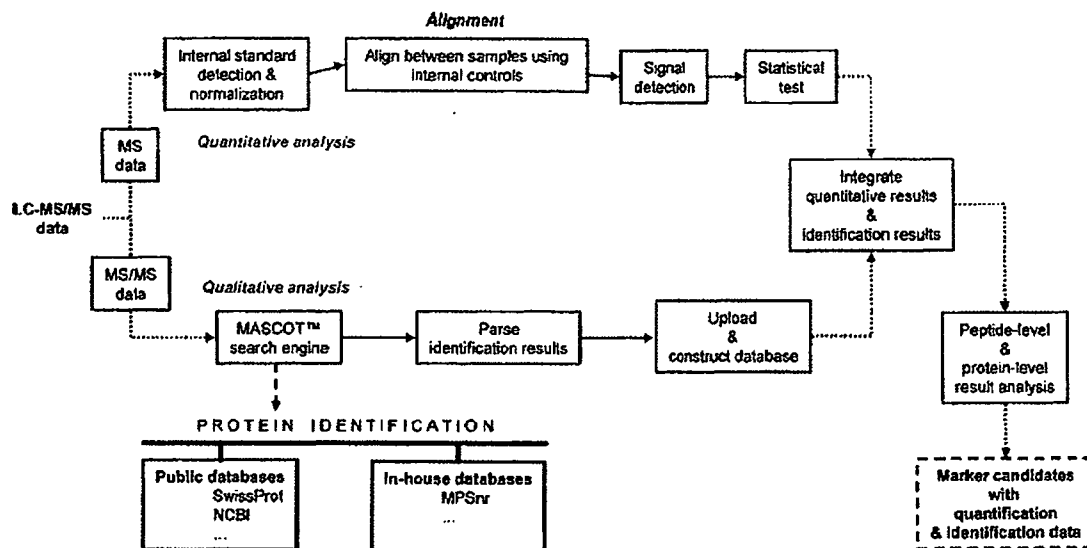
**Figure 4.** Overview of the data acquisition and database mining process developed within the gefitinib biomarker study.

identify the correct amino acid sequence. Internal standards are used for alignment of retention-times.

**How the Methodology Overcomes Some of the Hurdles.** The system described above addresses some of the hurdles noted previously. The digestion of all extracted plasma proteins into peptides will reduce the complexity by combining high-resolution nanoflow chromatographic fractionation with the separation power of modern mass spectrometry, performing automated and unattended shotgun sequencing in plasma.[35] Peptides are also more soluble and easier to handle than intact proteins. In addition, the two-dimensional quadrupole ion trap mass spectrometer[71] operates with a high-volume quadrupole electric field that makes it highly efficient to trap ions. The result is high sensitivity, high scanning speed, and better quantification over a wide dynamic range in comparison with the conventional three-dimensional ion-trap instruments.

Finding signals against a background of high variation is a further challenge, and the next section describes some approaches for addressing these.

## Initial Data Handling, Processing, and Analysis

Proteomic data analysis process can be considered as consisting of two components (Figure 4). *Quantitative analysis* is used to discover significant differences in peptide signal intensities by comparing two (or more) sample groups. This process uses data collected from an entire MS run to quantify the amount of peptide ions by their respective ion signal intensity. *Qualitative analysis* is used to identify the amino acid sequence of each peptide ion, from the respective product ion spectra. To maximize their value, the results from the two component analyses should be considered in combination.

A typical quantitative analysis may consist of several steps:

1. Normalization: To account for differences in the original sample concentrations. Typically, the total signal intensity is scaled to a constant value for each analyzed sample.

2. Alignment: Correcting for nonlinear fluctuation in retention time between different samples. A variety of methodologies are available for aligning LC—MS data sets. We have found the i-OPAL algorithm (Patent # WO 2004/090526 AI), which is based on the single linkage clustering algorithm[72] and which makes

use of internal standard signals, to perform well. Other alignment algorithms include xcms.[73]

3. Peak picking or signal detection: Identifying individual peptide ions within the data.

4. Identify discriminating peptides: A number of methods can be used, often in combination. A common approach is to apply a Student's $t$-test and select peptides which are significant, i.e., with a $p$-value less than the chosen cutoff, and which also show a fold-change or intensity ratio greater than another criterion. Further developments of this aspect are discussed in the Principled Statistical Analysis section.

A popular choice for qualitative analysis is the MASCOT MS/MS ion search program.[74] This may be run against a number of different peptide sequence databases,' for example the NCBI Nr, Refseq, Gene Ontology, HUGO, and Swiss-Prot sequence databases. The results of the quantitative analysis can then be combined with the qualitative analysis so that, for example, a peptide must be both discriminating and have annotation—i.e., have achieved a high MASCOT score showing confidence in identification—to be considered a candidate biomarker.

The approaches we have discussed above are focused on finding potentially discriminating proteins of clinical utility. In the following section, we describe the next stage in our thinking, namely how we could rapidly deploy in the clinic a viable method for exploiting a predictive proteomic fingerprint.

## A Proposal for Proteomics in the Clinical Setting: Mass Spectrometric Biomarker Assays - MSBA

Although today's technology allows for high-throughput analyses of many proteins rather than a single protein,[30] the details of how such multiplexing assays will be adapted for clinical use have not been well clarified. The Mass Spectrometric Biomarker Assay (MSBA) platform described here was conceived as one example of a rapid and seamless method to progress from identification of a diagnostic more directly to a clinically useful test. MSBA requires only a minute sample amount (5—20 $\mu$L) to obtain a read-out from a handful of quantified protein biomarkers (typically 3—35) and automatically analyzes proteins using liquid-phase separation and tandem mass spectrometry with simultaneous quantitation and identification.
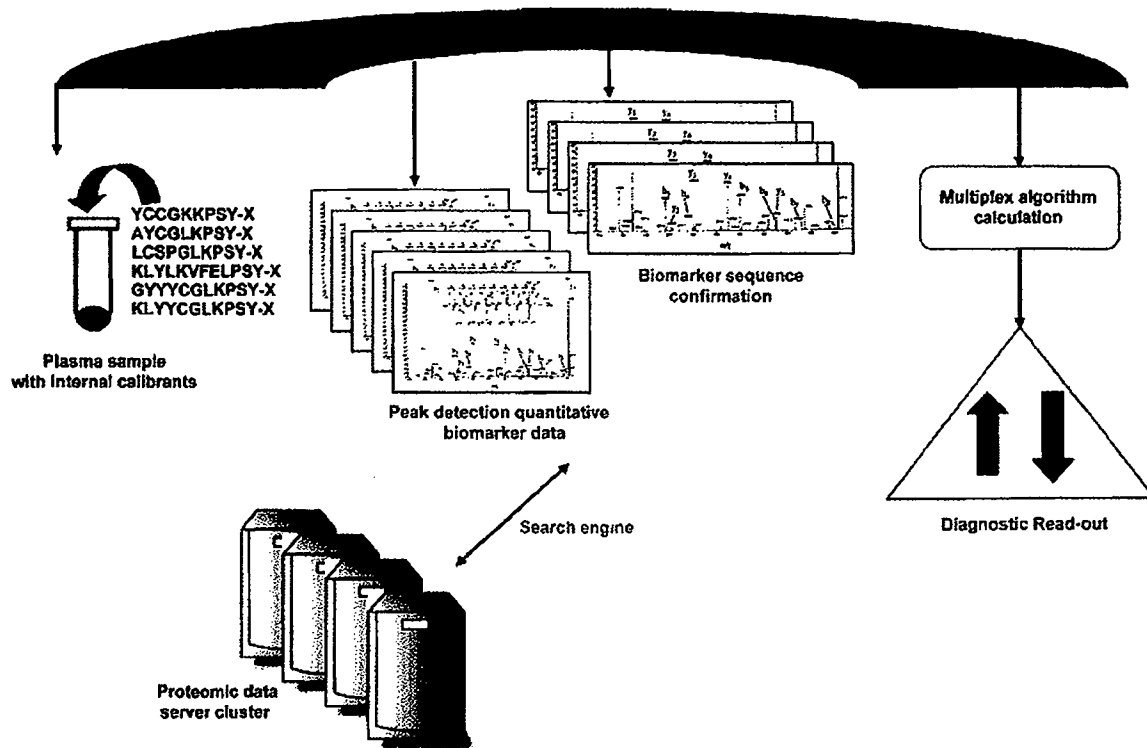
**Figure 5.** Entire flow of the operational components of Mass Spectrometric Biomarker Assays (MSBA).

The MSBA builds on a pre-defined Multiplex Biomarker list, which is stored within the MSBA database. Each marker entity has the values of masses and the relative retention time index with tolerance parameters. In running a patient sample, the predefined biomarker list is scanned to pick up patient sample signals that match with one of the predefined biomarker signals by satisfying the tolerance criteria (in general $\pm 1$ for $m/z$ value and $\pm 2\%$ for relative retention time index). The selected candidate signals are further confirmed using the product ion spectrum. That is, the product ion spectrum is represented as a vector by binning (grouping) the $m/z$ ratio values. Using the cosine correlation between the sample vectors and the reference vectors, we can confirm whether the selected candidate signals are truly assigned as target biomarkers. (A standard threshold value of the cosine correlation is 0.8.)

The process steps within the MSBA cycle are outlined in Figure 5. The calculation of the final multiplex biomarker assay read-out from all of the individual markers can be performed by a variety of applications, as discussed in more detail in the Principled Statistical Modeling Approach section. Figures 6A and B illustrate one approach, calculating a distance score which indicates to what extent a measured sample is distant from the case or control template in terms of predefined multiplex biomarkers.

$$S_{\text{case or control}} = \sqrt{\left[\frac{1}{n(n-2)}\right]\left[n\sum_i y_i^2 - \left(\sum_i y_i\right)^2 - \frac{\left[n\sum_i x_iy_i - \left(\sum_i x_i\right)\left(\sum_i y_i\right)\right]^2}{n\sum_i x_i^2 - \left(\sum_i x_i\right)^2}\right]}$$

If the ratio of $S_{\text{case}}$ and $S_{\text{control}}$ exceeds an MSBA threshold parameter, then the test sample is predicted to be a patient susceptible to develop ILD (ILD case); if not, the test sample is predicted to be a non-susceptible patient (control). We are currently evaluating the MSBA approach in practice.

## A Principled Statistical Modeling Approach

We have described an analytical approach based on proteomic data, with various novel developments. However, additional insight is needed to further improve model discrimination and to broaden the focus from the proteomic data to the ultimate goal of prediction using combinations of data. Statistical analysis can be used to provide further refinement by combining information from the full clinical and laboratory datasets.

An advantage of a multiple biomarker approach (e.g., proteomics) compared with standard single biomarker development is the capability to combine information from many different entities. An example is illustrated in Figure 7A. Considering each biomarker alone fails to separate the two groups of subjects, as there is considerable overlap for both biomarkers. Use of two biomarkers in combination completely separates the two groups.

We can also use clinical variables to advantage in the analysis of the peptide patterns. For example, the efficacy of gefitinib appears to be greater in non-smokers, women, patients of Asian origin, and patients with adenocarcinomas.[8] Figure 7B illustrates how, instead of two protein biomarkers, the combination of clinical data (e.g., age) and a proteomic biomarker is able to separate two groups.

On this basis, we propose using a principled statistical analysis approach to first explore and understand the data and
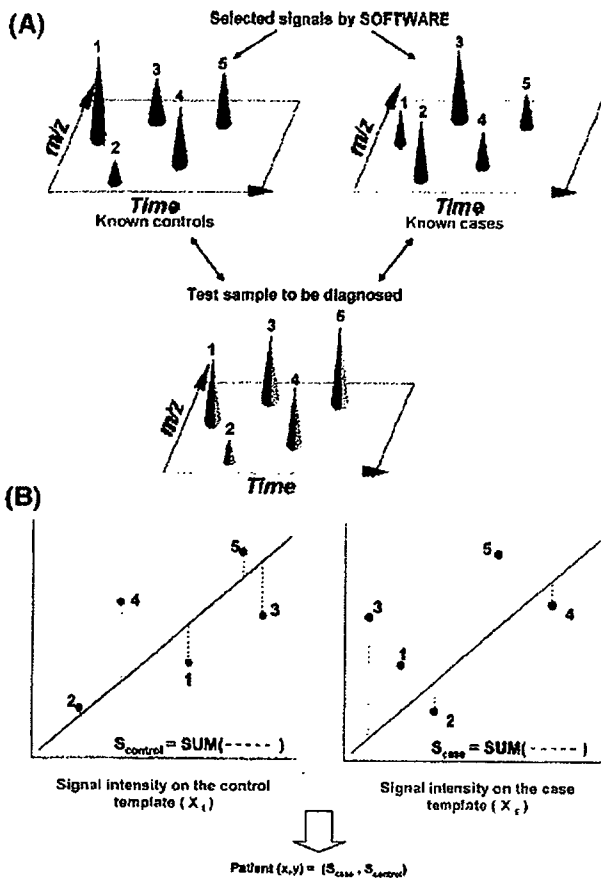
**(A)** Selected signals by SOFTWARE

Time
Known controls

Time
Known cases

Test sample to be diagnosed

Time

**(B)**

$S_{control} = SUM(----)$

$S_{case} = SUM(----)$

Signal intensity on the control template ($X_t$)

Signal intensity on the case template ($X_c$)

Patient $(x,y) = (S_{case}, S_{control})$

**Figure 6.** (A) Peptide signal comparison that MSBA (Mass Spectrometric Biomarker Assays) performs of the generated ions from the sample. The comparison is made both with the pattern of the controls and with the pattern of the case group for the corresponding signals. (B) Illustration of the regression model application of the MSBA where control templates and case templates are compared to that of the sample template generated in the analysis process.

then to model it and understand the quality of any models produced. A first step is to perform exploratory data analysis (EDA), for example using principal components analysis (PCA), to understand the major sources of data variation and the covariation between clinical parameters and protein intensity measures. The next step is univariate modeling for each protein marker individually, for example using analysis of covariance (ANCOVA), and an assessment of the effect of clinical parameters across the whole set of protein biomarkers using, for example, the False Discovery Rate as a tool.[75] This provides an understanding of key clinical variables and sources of variation within the data.

The next step is to perform multivariate predictive modeling using the proteins and clinical variables identified as being potentially important. There are a number of mathematical methods described in the literature for performing supervised classification, for example Support Vector Machines,[76] Random Forests,[77] PAM,[78] all of which have been successfully applied to high dimensional genomics data.[79] It remains an important unanswered question which modeling approach, or combination of modeling approaches, will generate the most predictive and robust models for data generated using this technology within a prospective study of this design.

Finally, to confirm that we have a practical prediction, the predictive power of a model must be assessed on a different set of patients from that used to generate the model. There are a number of approaches for external validation given a limited size dataset, for example the sequential approach of building a model based upon currently available data and testing on data from new patients when they become available, or withholding an arbitrary selection of subjects from the modeling as a test set and testing the model on these subjects. Internal validation approaches such as cross-validation or related bootstrapping methods may also be useful to assess the model selection *procedure*, but tend to overestimate the performance of a specific predictive model in subsequent external validation.[80,81] The key properties to consider when selecting an assessment method are to ensure that it will provide both precise and unbiased information regarding the prediction error rate of the potential model to be tested for clinical use. As well as assessing an overall predictive rate, it is also useful to separately assess the predictive rate for both the cases and controls and to consider the relative costs of making these false predictions within a clinical setting. Finally, the prevalence of the condition in question (here ILD) is also a critical factor in estimating what proportion of people predicted to be at risk are truly at risk, and this should also be borne in mind when evaluating a model for potential clinical use. The recently published FDA concept paper on drug-diagnostic co-
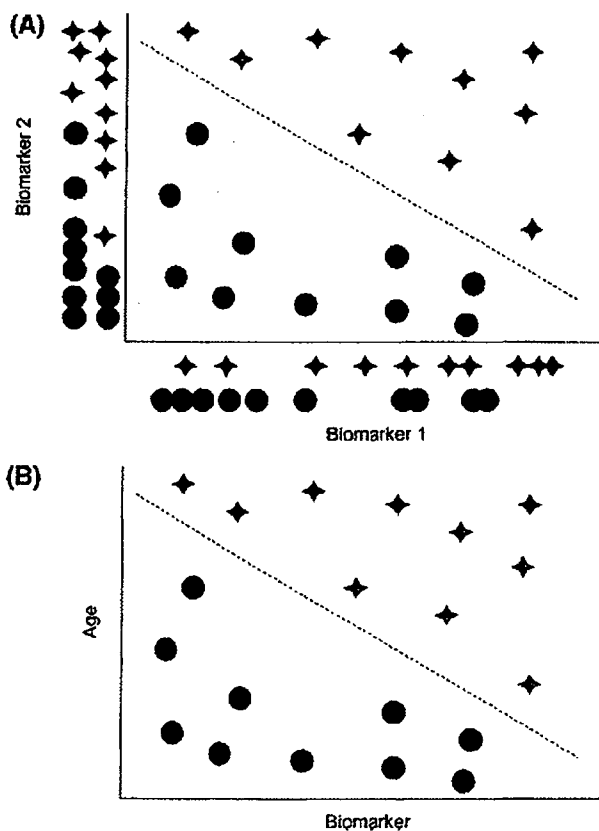


**(A)** Biomarker 2 / Biomarker 1

**(B)** Age / Biomarker

**Figure 7.** (A) Hypothetical example of the combined disease-linkage effect of two protein biomarkers. (Stars signify affected case individuals, circles non-affected control individuals). (B) Hypothetical example of the combined disease-linkage effect of a biomarker and a clinical variable. (Stars signify affected case individuals, circles non-affected control individuals).
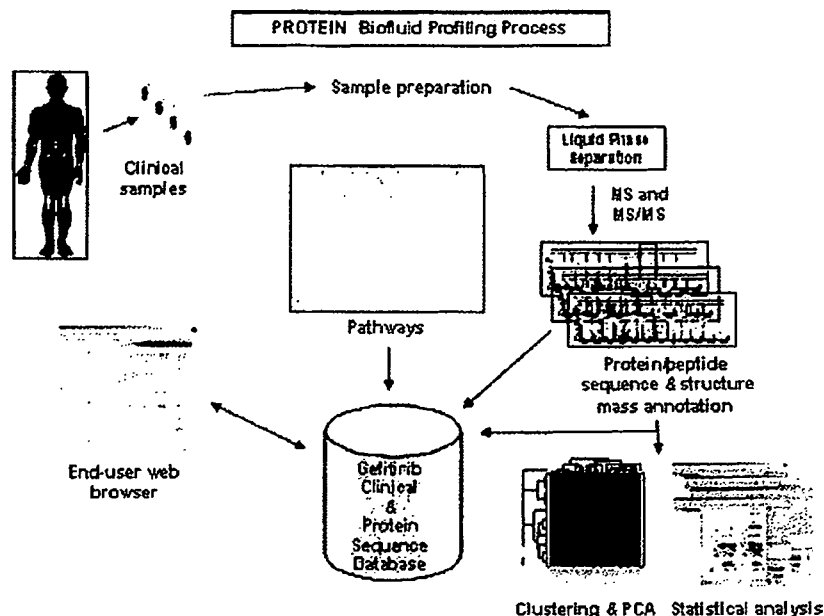
240

**Figure 8.** Illustration of the bioinformatics and data processing structure within which MSBA (Mass Spectrometric Biomarker Assays) data are captured, modified and analyzed.

development discusses many of the issues around validating predictive biomarkers.[82]

Finally, it is preferable to be able to assign a biological rationale to the biomarkers. Confidence in the reliability of a biomarker is greatly enhanced if we can correctly understand how it relates to the mechanism and progression of the disease of interest. Figure 8 illustrates a bioinformatics and data processing structure that we have developed to allow us to both conduct interactive exploratory and statistical analyses, and also investigate the disease and pathway linkage of discovered biomarker proteins through direct access to reference databases.

## Future Perspectives

Within this paper we have discussed many of the issues that need to considered in developing a personalized medicine approach. A key starting point is that rigorous steps are taken to ensure accurate diagnosis and the careful gathering of both clinical and proteomic data to facilitate the search for peptide patterns.

There are many challenges in performing protein analysis in blood, but mass spectrometry equipment and methods can now be used to generate peptide data with high sensitivity, high scanning speed, and improved quantification. Data handling and processing techniques for steps such as peak alignment and the subsequent methodologies for statistical modeling and analysis are now far enough developed to generate high quality data and robustly analyze these data with confidence.

We have provided details of the MSBA method that can be used to easily translate protein intensities into a practical multiplex assay which can be exploited in the clinic without the need to develop anti-bodies for ELISA. We have also described how an expanded statistical analysis can be used to allow for the individual variance of protein expression to enable us to focus on the proteomic patterns that are actually related to ILD. Finally, we have emphasized the importance of validat-

ing the predictive power of a biomarker tool in a way that reflects the real-life setting of intended clinical use.

Hopefully, this combination of developments over a range of different areas brings us one step closer to a practical personalized medicine.

IRESSA is a trademark of the AstraZeneca group of companies.

## References

(1) Thatcher, N.; Chang, A.; Parikh, P.; Pereira, J. R.; Ciuleanu, T.; von Pawel, J.; et al. Gefitinib plus best supportive care in previously treated patients with refractory advanced non-small-cell lung cancer: results from a randomised, placebo-controlled, multicentre study (Iressa Survival Evaluation in Lung Cancer). *Lancet* **2005**, *366*, 1527–1537.

(2) Hirsch, F. R.; Varella-Garcia, M.; McCoy, J.; West, H.; Xavier, A. C.; Gumerlock, P.; et al. Increased epidermal growth factor receptor gene copy number detected by fluorescence in situ hybridization associates with increased sensitivity to gefitinib in patients with bronchioloalveolar carcinoma subtypes: a Southwest Oncology Group study. *J. Clin. Oncol.* **2005**, *23*, 6838–6845.

241

(3) Cappuzzo, F.; Varella-Garcia, M.; Shigematsu, H.; Domenichini, I.; Bartolini, S.; Ceresoli, G. L.; et al. Increased HER2 gene copy number is associated with response to gefitinib therapy in epidermal growth factor receptor-positive non-small-cell lung cancer patients. *J. Clin. Oncol.* 2005, *23*, 5007–5018.

(4) Araki, J.; Okamoto, I.; Suto, R.; Ichikawa, Y.; Sasaki, J. Efficacy of the tyrosine kinase inhibitor gefitinib in a patient with metastatic small cell lung cancer. *Lung Cancer* 2005, *48*, 141–144.

(5) Kim, K. S.; Jeong, J. Y.; Kim, Y. C.; Na, K. J.; Kim, Y. H.; Ahn, S. J.; et al. Predictors of the response to gefitinib in refractory non-small cell lung cancer. *Clin. Cancer Res.* 2005, *11*, 2244–2251.

(6) Lynch, T. J.; Bell, D. W.; Sordella, R.; Gurubhagavatula, S.; Okimoto, R. A.; Brannigan, B. W.; et al. Activating mutations in the epidermal growth factor receptor underlying responsiveness of non-small-cell lung cancer to gefitinib. *N. Engl. J. Med.* 2004, *350*, 2129–2139.

(7) Paez, J. G.; Jänne, P. A.; Lee, J. C.; Tracy, S.; Greulich, H.; Gabriel, S.; et al. EGFR mutations in lung cancer: correlation with clinical response to gefitinib therapy. *Science* 2004, *304*, 1497–1500.

(8) Shigematsu, H.; Lin, L.; Takahashi, T.; Nomura, M.; Suzuki, M.; Wistuba II; et al. Clinical and biological features associated with epidermal growth factor receptor gene mutations in lung cancers. *J. Natl. Cancer Inst.* 2005, *97*, 339–346.

(9) American Thoracic Society: American Thoracic Society/European Respiratory Society International Multidisciplinary Consensus Classification of the Idiopathic Interstitial Pneumonias. This joint statement of the American Thoracic Society (ATS), and the European Respiratory Society (ERS) was adopted by the ATS Board of Directors, June 2001 and by The ERS Executive Committee, June 2001. *Am. J. Respir. Crit. Care Med.* 2002, *165*, 277–304.

(10) Raghu, G.; Nyberg, F.; Morgan, G. The epidemiology of interstitial lung disease and its association with lung cancer. *Br. J. Cancer* 2004, *91* (Suppl. 2), S3–S10.

(11) Asada, K.; Mukai, J.; Ougushi, F. Characteristics and management of lung cancer in patients with idiopathic pneumonia. *Jap. J. Thor. Dis.* 1992, *51*, 214–219.

(12) Hubbard, R.; Venn, A.; Lewis, S.; Britton, J. Lung cancer and cryptogenic fibrosing alveolitis. A population-based cohort study. *Am. J. Respir. Crit. Care Med.* 2000, *161*, 5–8.

(13) Matsushita, H.; Tanaka, S.; Saiki, Y.; Hara, M.; Nakata, K.; Tanimura, S.; et al. Lung cancer associated with usual interstitial pneumonia. *Pathol. Int.* 1995, *45*, 925–932.

(14) Ogura, T.; Kondo, A.; Sato, A.; Ando, M.; Tamura, M. Incidence and clinical features of lung cancer in patients with idiopathic interstitial pneumonia. *Nihon Kyobu Shikkan Gakkai Zasshi* 1997, *35*, 294–299.

(15) Takeuchi, E.; Yamaguchi, T.; Mori, M.; Tanaka, S.; Nakagawa, M.; Yokota, S.; et al. Characteristics and management of patients with lung cancer and idiopathic interstitial pneumonia. *Nihon Kyobu Shikkan Gakkai Zasshi* 1996, *34*, 653–658.

(16) Turner-Warwick, M.; Lebowitz, M.; Burrows, B.; Johnson, A. Cryptogenic fibrosing alveolitis and lung cancer. *Thorax* 1980, *35*, 496–499.

(17) Baumgartner, K. B.; Samet, J. M.; Stidley, C. A.; Colby, T. V.; Waldron, J. A. Cigarette smoking: a risk factor for idiopathic pulmonary fibrosis. *Am. J. Respir. Crit. Care Med.* 1997, *155*, 242–248.

(18) Britton, J.; Hubbard, R. Recent advances in the aetiology of cryptogenic fibrosing alveolitis. *Histopathology* 2000, *37*, 387–392.

(19) Iwai, K.; Mori, T.; Yamada, N.; Yamaguchi, M.; Hosoda, Y. Idiopathic pulmonary fibrosis. Epidemiologic approaches to occupational exposure. *Am. J. Respir. Crit. Care Med.* 1994, *150*, 670–675.

(20) Nagai, S.; Hoshino, Y.; Hayashi, M.; Ito, I. Smoking-related interstitial lung diseases. *Curr. Opin. Pulm. Med.* 2000, *6*, 415–419.

(21) Lilly. Gemcitabine prescribing information. http://pi.lilly.com/gemzar.pdf, 2003.

(22) Kunitoh, H.; Watanabe, K.; Onoshi, T.; Furuse, K.; Niitani, H.; Taguchi, T. Phase II trial of docetaxel in previously untreated advanced non-small-cell lung cancer: a Japanese cooperative study. *J. Clin. Oncol.* 1996, *14*, 1649–1655.

(23) Merad, M.; Le Cesne, A.; Baldeyrou, P.; Mesurolle, B.; Le Chevalier, T. Docetaxel and interstitial pulmonary injury. *Ann. Oncol.* 1997, *8*, 191–194.

(24) Wang, G.-S.; Yan, K.-Y.; Perng, R.-P. Life-threatening hypersensitivity pneumonitis induced by docetaxel (taxotere). *Br. J. Cancer* 2001, *85*, 1247–1250.

(25) Erasmus, J. J.; McAdams, H. P.; Rossi, S. E. Drug-induced lung injury. *Semin. Roentgenol.* 2002, *37*, 72–81.

(26) Aviram, G.; Yu, E.; Tai, P.; Lefcoe, M. S. Computed tomography to assess pulmonary injury associated with concurrent chemoradiotherapy for inoperable non-small cell lung cancer. *Can. Assoc. Radiol. J.* 2001, *52*, 385–391.

(27) Yoshida, S. The results of gefitinib prospective investigation. *Med. Drug J.* 2005, *41*, 772–789.

(28) Mueller, N. L.; White, D. A.; Jiang, H.; Gemma, A. Diagnosis and management of drug-associated interstitial lung disease. *Br. J. Cancer* 2004, *91*, S24–S30.

(29) Marko-Varga, G.; Fehniger, T. E. Proteomics and disease—the challenges for technology and discovery. *J. Proteome Res.* 2004, *3*, 167–178.

(30) Marko-Varga, G.; Lindberg, H.; Lofdahl, C. G.; Jonsson, P. H. L.; Dahlback, M.; Lindquist, E.; et al. Discovery of biomarker candidates within disease by protein profiling: principles and concepts. *J. Proteome Res.* 2005, *4*, 1200–1212.

(31) Omenn, G. S. The Human Proteome Organization Plasma Proteome Project pilot phase: reference specimens, technology platform comparisons, and standardized data submissions and analyses. *Proteomics* 2004, *4*, 1235–1240.

(32) Omenn, G. S. Advancement of biomarker discovery and validation through the HUPO plasma proteome project. *Dis. Markers* 2004, *20*, 131–134.

(33) Orchard, S.; Hermjakob, H.; Binz, P. A.; Hoogland, C.; Taylor, C. F.; Zhu, W.; et al. Further steps towards data standardisation: the Proteomic Standards Initiative HUPO 3(rd) annual congress, Beijing 25–27(th) October, 2004. *Proteomics* 2005, *5*, 337–339.

(34) Anderson, N. G.; Matheson, A.; Anderson, N. L. Back to the future: the human protein index (HPI) and the agenda for post-proteomic biology. *Proteomics* 2001, *1*, 3–12.

(35) Anderson, N. L.; Anderson, N. G. The human plasma proteome: history, character, and diagnostic prospects. *Mol. Cell. Proteomics* 2002, *1*, 845–867.

(36) Jacobs, J. M.; Adkins, J. N.; Qian, W. J.; Liu, T.; Shen, Y.; Camp, D. G.; et al. Utilizing human blood plasma for proteomic biomarker discovery. *J. Proteome Res.* 2005, *4*, 1073–1085.

(37) Anderson, N. G.; Anderson, L. The Human Protein Index. *Clin. Chem.* 1982, *28*, 739–748.

(38) Haab, B. B.; Geierstanger, B. H.; Michailidis, G.; Vitzthum, F.; Forrester, S.; Okon, R.; et al. Immunoassay and antibody microarray analysis of the HUPO Plasma Proteome Project reference specimens: systematic variation between sample types and calibration of mass spectrometry data. *Proteomics* 2005, *5*, 3278–3291.

(39) Martens, L.; Hermjakob, H.; Jones, P.; Adamski, M.; Taylor, C.; States, D.; et al. PRIDE: the proteomics identifications database. *Proteomics* 2005, *5*, 3537–3545.

(40) Omenn, G. S.; States, D. J.; Adamski, M.; Blackwell, T. W.; Menon, R.; Hermjakob, H.; et al. Overview of the HUPO Plasma Proteome Project: results from the pilot phase with 35 collaborating laboratories and multiple analytical groups, generating a core dataset of 3020 proteins and a publicly-available database. *Proteomics* 2005, *5*, 3226–3245.

(41) Patterson, S. D. Data analysis–the Achilles heel of proteomics. *Nat. Biotechnol.* 2003, *21*, 221–222.

(42) Rahbar, A. M.; Fenselau, C. Integration of Jacobson's pellicle method into proteomic strategies for plasma membrane proteins. *J. Proteome Res.* 2004, *3*, 1267–1277.

(43) Ho, Y.; Gruhler, A., Heilbut, A.; Bader, G. D.; Moore, L.; Adams, S. L.; et al. Systematic identification of protein complexes in Saccharomyces cerevisiae by mass spectrometry. *Nature* 2002, *415*, 180–183.

(44) Aebersold, R.; Mann, M. Mass spectrometry-based proteomics. *Nature* 2003, *422*, 198–207.

(45) Anderson, N. L.; Polanski, M.; Pieper, R.; Gatlin, T.; Tirumalai, R. S.; Conrads, T. P.; et al. The human plasma proteome: a nonredundant list developed by combination of four separate sources. *Mol. Cell. Proteomics* 2004, *3*, 311–326.

(46) Olsen, J. V.; Mann, M. Improved peptide identification in proteomics by two consecutive stages of mass spectrometric fragmentation. *Proc. Natl. Acad. Sci. U.S.A.* 2004, *101*, 13417–13422.

(47) Sadygov, R. G.; Liu, H.; Yates, J. R. Statistical models for protein validation using tandem mass spectral data and protein amino acid sequence databases. *Anal. Chem.* 2004, *76*, 1664–1671.

(48) Fujii, K.; Nakano, T.; Kanazawa, M.; Akimoto, S.; Hirano, T.; Kato, H.; et al. Clinical-scale high-throughput human plasma proteome analysis: lung adenocarcinoma. *Proteomics* 2005, *5*, 1150–1159.

242

(49) Campbell, J. M.; Collings, B. A.; Douglas, D. J. A new linear ion trap time-of-flight system with tandem mass spectrometry capabilities. *Rapid Commun. Mass Spectrom.* 1998, *12*, 1463–1474.

(50) Cha, B. C.; Blades, M.; Douglas, D. J. An interface with a linear quadropole ion guide for an electrospray-ion trap mass spectrometer system. *Anal. Chem.* 2000, *72*, 5647–5654.

(51) Hager, J. W. Product ion spectral simplification using time-delayed fragment ion capture with tandem linear ion traps. *Rapid Commun. Mass Spectrom.* 2003, *17*, 1389–1398.

(52) Syka, J. E.; Marto, J. A.; Bai, D. L.; Horning, S.; Senko, M. W.; Schwartz, J. C.; et al. Novel linear quadrupole ion trap/FT mass spectrometer: performance characterization and use in the comparative analysis of histone H3 post-translational modifications. *J. Proteome Res.* 2004, *3*, 621–626.

(53) Shen, Y.; Zhao, R.; Belov, M. E.; Conrads, T. P.; Anderson, G. A.; Tang, K.; et al. Packed capillary reversed-phase liquid chromatography with high-performance electrospray ionization Fourier transform ion cyclotron resonance mass spectrometry for proteomics. *Anal. Chem.* 2001, *73*, 1766–1775.

(54) Wu, S. L.; Kim, J.; Hancock, W. S.; Karger, B. Extended Range Proteomic Analysis (ERPA): a new and sensitive LC-MS platform for high sequence coverage of complex proteins with extensive post-translational modifications-comprehensive analysis of beta-casein and epidermal growth factor receptor (EGFR). *J. Proteome Res.* 2005, *4*, 1155–1170.

(55) Olsen, J. V.; de Godoy, L. M.; Li, G.; Macek, B.; Mortensen, P.; Pesch, R.; et al. Parts per million mass accuracy on an Orbitrap mass spectrometer via lock mass injection into a C-trap. *Mol. Cell. Proteomics* 2005, *4*, 2010–2021.

(56) Yates, J. R.; Cociorva, D.; Liao, L.; Zabrouskov, V. Performance of a linear ion trap-Orbitrap hybrid for peptide analysis. *Anal. Chem.* 2006, *78*, 493–500.

(57) Anderson, D. C.; Li, W.; Payan, D. G.; Noble, W. S. A new algorithm for the evaluation of shotgun peptide sequencing in proteomics: support vector machine classification of peptide MS/MS spectra and SEQUEST scores. *J. Proteome Res.* 2003, *2*, 137–146.

(58) Carr, S.; Aebersold, R.; Baldwin, M.; Burlingame, A.; Clauser, K.; Nesvizhskil, A. The need for guidelines in publication of peptide and protein identification data: working group on publication guidelines for peptide and protein identification data. *Mol. Cell. Proteomics* 2004, *3*, 531–533.

(59) Fenyo, D.; Beavis, R. C. A method for assessing the statistical significance of mass spectrometry-based protein identifications using general scoring schemes. *Anal. Chem.* 2003, *75*, 768–774.

(60) Nesvizhskii, A. I.; Keller, A.; Kolker, E.; Aebersold, R. A statistical model for identifying proteins by tandem mass spectrometry. *Anal. Chem.* 2003, *75*, 4646–4658.

(61) Peri, S.; Navarro, J. D.; Kristiansen, T. Z.; Amanchy, R.; Surendranath, V.; Muthusamy, B.; et al. Human protein reference database as a discovery resource for proteomics. *Nucleic Acids Res.* 2004, *32*, D497–D501.

(62) Kratchmarova, I.; Blagoev, B.; Haack-Sorensen, M.; Kassem, M.; Mann, M. Mechanism of divergent growth factor effects in mesenchymal stem cell differentiation. *Science* 2005, *308*, 1472–1477.

(63) Dreger, M.; Bengtsson, L.; Schöneberg, T.; Otto, H.; Hucho, F. Nuclear envelope proteomics: novel integral membrane proteins of the inner nuclear membrane. *Proc. Natl. Acad. Sci. U.S.A.* 2001, *98*, 11943–11948.

(64) Giot, L.; Bader, J. S.; Brouwer, C.; Chaudhuri, A.; Kuang, B.; Li, Y.; et al. A protein interaction map of Drosophilia melanogaster. *Science* 2003, *302*, 1727–1736.

(65) Johnson, J. R.; Florens, L.; Carucci, D. J.; Yates, J. R., III. Proteomics in malaria. *J. Proteome Res.* 2004, *3*, 296–306.

(66) Hirsch, J.; Hansen, K. C.; Burlingame, A. L.; Matthay, M. A. Proteomics: current techniques and potential applications to lung disease. *Am. J. Physiol. Lung Cell Mol. Physiol.* 2004, *287*, L1–L23.

(67) Malmström, J.; Larsen, K.; Hansson, L.; Löfdahl, C.-G.; Norregard-Jensen, O.; Marko-Varga, G.; et al. Proteoglycan and proteome profiling of central human pulmonary fibrotic tissue utilizing minaturized sample preparation: A feasibility study. *Proteomics* 2002, *2*, 394–404.

(68) Malmström, J.; Larsen, K.; Malmström, L.; Tufvesson, E.; Parker, K.; Marchese, J.; et al. Proteome annotations and identifications of the human pulmonary fibroblast. *J. Proteome Res.* 2004, *3*, 525–537.

(69) Oh, P.; Li, Y.; Yu, J.; Durr, E.; Krasinska, K. M.; Carver, L. A.; et al. Subtractive proteomic mapping of the endothelial surface in lung and solid tumours for tissue-specific therapy. *Nature* 2004, *429*, 629–635.

(70) Fujii, K.; Nakano, T.; Kawamura, T.; Usui, F.; Bando, Y.; Wang, R.; et al. Multidimensional protein profiling technology and its application to human plasma proteome. *J. Proteome Res.* 2004, *3*, 712–718.

(71) Schwartz, J. C.; Senko, M. W.; Syka, J. E. A two-dimensional quadrupole ion trap mass spectrometer. *J. Am. Soc. Mass Spectrom.* 2002, *13*, 659–669.

(72) Sneath, P. H. A.; Sokai, R. R. *Numerical Taxonomy, The principles and practice of numerical classification*; W. H. Freeman and Co.: San Francisco, 1973.

(73) Smith, C. A.; Want, E. J.; O'Maille, G.; Abagyan, R.; Siuzdak, G. XCMS: processing mass spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and identification. *Anal. Chem.* 2006, *78*, 779–787.

(74) Perkins, D. N.; Pappin, D. J.; Creasy, D. M.; Cottrell, J. S. Probability-based protein identification by searching sequence-databases using mass spectrometry data. *Electrophoresis* 1999, *20*, 3551–3567.

(75) Storey, J. A direct approach to false discovery rates. *J. R. Stat. Soc. Ser. B* 2002, *64*, 479.

(76) Vapnik, V. *Statistical Learning Theory*; Wiley: Chichester, UK, 1998.

(77) Breiman, L. Random forests. *Mach. Learn.* 2001, *45*, 5–32.

(78) Tibshirani, R.; Hastie, T.; Narasimhan, B.; Chu, G. Diagnosis of multiple cancer types by shrunken centroids of gene expression. *Proc. Natl. Acad. Sci. U.S.A.* 2002, *99*, 6567–6572.

(79) Lee, J. W.; Lee, J. B.; Park, M.; Song, S. H. An extensive comparison of recent classification tools applied to microarray data. *Comp. Stat. Data Anal.* 2005, *48*, 869–885.

(80) Steyerberg, E. W.; Harrell, F. E., Jr.; Borsboom, G. J.; Eijkemans, M. J.; Vergouwe, Y.; Habbema, J. D. Internal validation of predictive models: efficiency of some procedures for logistic regression analysis. *J. Clin. Epidemiol.* 2001, *54*, 774–781.

(81) Bleeker, S. E.; Moll, H. A.; Steyerberg, E. W.; Donders, A. R.; Derksen-Lubsen, G.; Grobbee, D. E.; et al. External validation is necessary in prediction research: a clinical example. *J. Clin. Epidemiol.* 2003, *56*, 826–832.

(82) Food and Drug Administration (FDA): Drug-diagnostic co-development concept paper. Draft—not for implementation. http://www.fda.gov/cder/genomics/pharmacoconceptfn.pdf, 2005.

# Personalized Medicine and Proteomics: Lessons from Non-Small Cell Lung Cancer

György Marko-Varga,[†,‡] Atsushi Ogiwara,[†,§,‖] Toshihide Nishimura,[§,‖] Takeshi Kawamura,[§,‖]
Kiyonaga Fujii,[§,‖] Takao Kawakami,[§,‖] Yutaka Kyono,[‖] Hsiao-kun Tu,[‖] Hisae Anyoji,[‖]
Mitsuhiro Kanazawa,[‖] Shingo Akimoto,[‖] Takashi Hirano,[⊥] Masahiro Tsuboi,[⊥] Kazuto Nishio,[§]
Shuji Hada,[#] Haiyi Jiang,[×] Masahiro Fukuoka,[∧] Kouichiro Nakata,[♦] Yutaka Nishiwaki,[÷]
Hideo Kunito,[§] Ian S. Peers,[○] Chris G. Harbron,[○] Marie C. South,[⁓] Tim Higenbottam,[∇,£]
Fredrik Nyberg,[*,∇,☆] Shoji Kudoh,[ɩ] and Harubumi Kato[⊥]

Respiratory Biological Sciences, AstraZeneca R&D Lund, SE-221 87 Lund, Sweden, Clinical Proteome Center,
Tokyo Medical University, Shinjuku Sumitomo Building 17F, 2-6-1 Nishishinjuku, Shinjuku, Tokyo 163-0217,
Japan, Medical ProteoScope Company, Limited, Shinjuku Sumitomo Building 17F, 2-6-1 Nishishinjuku,
Shinjuku, Tokyo 163-0217, Japan, Department of Surgery, Tokyo Medical University, 6-7-1 Nishishinjuku,
Shinjuku, Tokyo 160-0023, Japan, Clinical Division, Research & Development, AstraZeneca K.K., Umeda Sky
Building Tower East, 1-88, 1-chome, Ohyodo-naka, Kita-ku, Osaka 531-0076, Japan, Clinical Science
Department, Research & Development, AstraZeneca K.K., Umeda Sky Building Tower East, 1-88, 1-chome,
Ohyodo-naka, Kita-ku, Osaka 531-0076, Japan, Department of Medical Oncology, Kinki University School of
Medicine, 377-2,Ohno-higashi, Osakasayama-city 589-8511, Osaka, Japan, Department of Respiratory Diseases,
Toho University School of Medicine, 6-11-1, Omori-nishi, Ota-ku, Tokyo 143-8541, Japan, Dept. of Thoracic
Oncology, National Cancer Centre Hospital East, 6-5-1, Kashiwanoha, Kashiwa-city, Chiba 277-8577, Japan,
Statistical Sciences, AstraZeneca R&D Alderley Park, Cheshire, UK, Cancer & Infection Statistics, AstraZeneca
R&D Alderley Park, Cheshire, UK, Medicine & Science, AstraZeneca R&D Charnwood, Loughborough LE11 5RH,
Leicestershire, UK, Sheffield University, Sheffield, UK, Epidemiology, AstraZeneca R&D Mölndal, SE-431 83
Mölndal, Sweden, Institute of Environmental Medicine, Karolinska Institute, Box 210,
SE-171 77 Stockholm, Sweden, and 4th Department of Internal Medicine, Nippon Medical School,
1-1-5, Sendagi, Bunkyo-ku, Tokyo 113-8603, Japan

Personalized medicine allows the selection of treatments best suited to an individual patient and disease
phenotype. To implement personalized medicine, effective tests predictive of response to treatment or
susceptibility to adverse events are needed, and to develop a personalized medicine test, both high
quality samples and reliable data are required. We review key features of state-of-the-art proteomic
profiling and introduce further analytic developments to build a proteomic toolkit for use in personalized
medicine approaches. The combination of novel analytical approaches in proteomic data generation,
alignment and comparison permit translation of identified biomarkers into practical assays. We further
propose an expanded statistical analysis to understand the sources of variability between individuals
in terms of both protein expression and clinical variables and utilize this understanding in a predictive
test.

**Keywords:** personalized medicine • gefitinib • therapy • interstitial lung disease • non-small cell lung cancer •
biomarkers • predictive test • mass spectrometry • statistical analysis • proteomics

* To whom correspondence should be addressed. Epidemiology, Astra-
Zeneca R&D Mölndal, SE-431 83 Mölndal, Sweden; Tel, +46 31 706 5203;
Fax, +46 31 776 3828; E-mail, Fredrik.Nyberg@astrazeneca.com.
† György Marko-Varga and Atsushi Ogiwara made equal contributions to
this manuscript.
‡ Respiratory Biological Sciences, AstraZeneca R&D Lund.
§ Clinical Proteome Center, Tokyo Medical University.
‖ Medical ProteoScope Co., Ltd.
⊥ Department of Surgery, Tokyo Medical University.
∧ Clinical Division, Research & Development, AstraZeneca K.K.
⁓ Clinical Science Department, Research & Development, AstraZeneca K.K.

∇ Department of Medical Oncology, Kinki University School of Medicine.
♦ Department of Respiratory Diseases, Toho University School of
Medicine.
⁓ Dept. of Thoracic Oncology, National Cancer Centre Hospital East.
○ Statistical Sciences, AstraZeneca R&D Alderley Park.
⁓ Cancer & Infection Statistics, AstraZeneca R&D Alderley Park.
∇ Medicine & Science, AstraZeneca R&D Charnwood.
£ Medical School, Sheffield University.
☆ Epidemiology, AstraZeneca R&D Mölndal.
× Institute of Environmental Medicine, Karolinska Institute.
ɩ Fourth Department of Internal Medicine, Nippon Medical School.

# perspectives

## Introduction

A personalized medicine approach uses appropriate biomarkers to select treatments best suited for an individual patient and disease phenotype. A multiple biomarker approach (e.g., proteomics) has the advantage over conventional single biomarkers of combining many different pieces of information. Here, we review the key features of state-of-the-art proteomic profiling and introduce recent analytic developments to build a proteomic toolkit for use in personalized medicine, and we describe how these may be applied in a viable method for exploiting predictive proteomic fingerprints in the clinic. The potential of our proteomics toolkit hopefully brings us one step closer to a practical personalized medicine.

Cancer therapy is moving toward individually selected treatments, chosen not only according to tumor cell type but also based on the patient's predicted responsiveness to different classes of therapy or susceptibility to therapeutic adverse events. This emerging personalized medicine approach draws on both genotype and phenotype information, including protein expression. To implement personalized medicine, we need to develop effective biomarker tests predictive of response to treatment or susceptibility to adverse events. The benefits of personalized medicine are exemplified by considering interstitial lung disease (ILD) among non-small cell lung cancer (NSCLC) patients, which is associated with various kinds of chemotherapy treatment. A personalized medicine approach, using a simple blood test to predict those NSCLC patients at risk of developing ILD, would clearly be of great value.

We review current thinking and present some novel developments in a number of areas that have to be integrated to develop and then practically apply such tests in a clinical setting:

- The large scale collection of reliable and high quality phenotypic and clinical data and blood samples.
- Protein analysis in blood.
- Data acquisition, handling, combining and analysis.
- Interpretation and utilization of results in a clinical setting.

## Clinical Background

### A Motivating Example: Gefitinib (IRESSA) Treatment of NSCLC.
The concepts of proteomics-based personalized medicine discussed in this article are very generally applicable. A motivating example that we will refer to in order to illustrate the potential benefits of personalized medicine is ongoing work in attempting to develop a simple blood test to address the potential occurrence of ILD in seriously ill NSCLC patients, the target group for the NSCLC treatment gefitinib.

Gefitinib is a "small molecule" inhibitor of the enzyme tyrosine kinase of the epidermal growth factor receptor (EGFR) family, such as erbB1. It is an approved therapy for advanced NSCLC in many countries and offers important clinical benefits (tumor shrinkage and improvement in disease-related symptoms) for "end-stage" patients. The large phase III ISEL (IRESSA Survival Evaluation in Lung Cancer) trial demonstrated some improvement in survival with gefitinib which failed to reach statistical significance compared with placebo in the overall population and in patients with adenocarcinoma.[1] However, in preplanned subgroup analyses, a significant increase in survival was shown with gefitinib in patients of Asian ethnicity and in patients who had never smoked.[1]

Analysis of the biomarker data from a subset of patients in the ISEL study suggested that patients with pretreated advanced

NSCLC who have tumors with a high EGFR gene copy number (detected by fluorescent in situ hybridization (FISH)) have a higher likelihood of increased survival when treated with gefitinib compared with placebo.[2] Increased HER2 gene copy number has also been seen in tumors from patients who are responsive to gefitinib.[3] Somatic-activating mutations of EGFR in tumor tissue have also been associated with increased gefitinib responsiveness in patients with NSCLC.[4-7] Such mutations are more commonly found in tumor samples from patients of Asian origin and non-smokers.[8]

Following the ISEL subgroup analyses, and the biomarker evidence, it has become important to clarify which patients are more suitable for treatment with gefitinib. Analyses for both somatic-activating mutations and gene copy number require tumor tissue, which is not always available from the time of diagnosis; therefore, a blood test may represent a more versatile option and be of great value to clinicians.

With respect to tolerability, the search for a blood test that might include both genetic and proteomic biomarkers to define patients at risk of adverse effects from a drug, for example interstitial lung disease with gefitinib, is a focus of research.

### Interstitial Lung Disease as a Complication in NSCLC Patients.
ILD is a disease that afflicts the parenchyma or alveolar region of the lungs.[9] The alveolar septa (the walls of the alveoli) become thickened with fibrotic tissue. Associated with drug use, it can present precipitously with acute diffuse alveolar damage (DAD). The lungs show so-called "ground glass" shadowing on chest radiology, and patients complain of severe breathlessness. There are no effective treatments but patients can be supported by oxygen supplementation, corticosteroid therapy, or assisted ventilation. The process of alveolar damage is however fatal in some patients. ILD is a co-morbidity in patients with NSCLC.[10-16] Both diseases are associated with cigarette smoking,[17-20] and ILD is also considered to be associated with various kinds of lung cancer chemotherapy.[21-26]

In the ISEL study of gefitinib in NSCLC mentioned above, ILD-type events occurred in 1% of both placebo and gefitinib-treated patients.[1] Most ILD-type events occurred in patients of Asian origin, where placebo and treated patients had similar prevalences of respectively 4% and 3%. The rate observed in the gefitinib-treated arm was in line with earlier safety data from Japan and a large gefitinib post-marketing surveillance study in Japan (3322 patients), where the reported rate of ILD-type events was 5.8%.[27]

A simple blood test to predict the potential occurrence of ILD in seriously ill NSCLC patients before initiating treatments would clearly be of great value. This article describes the personalized medicine approach, which could be used to provide such a test. Consequently, the proteomics objectives of the preliminary phase of the study we describe were to verify the protein expression alterations in blood plasma from case patients (who developed ILD) and control patients (without ILD) treated by gefitinib, using a liquid chromatography–mass spectrometry/mass spectrometry (LC–MS/MS) proteomics platform.

## Data and Sample Collection

To develop a personalized medicine test, it is essential to have access to an adequately sized collection of high quality tissue samples on which to perform proteomics analysis, with corresponding reliable diagnostic and clinical data.
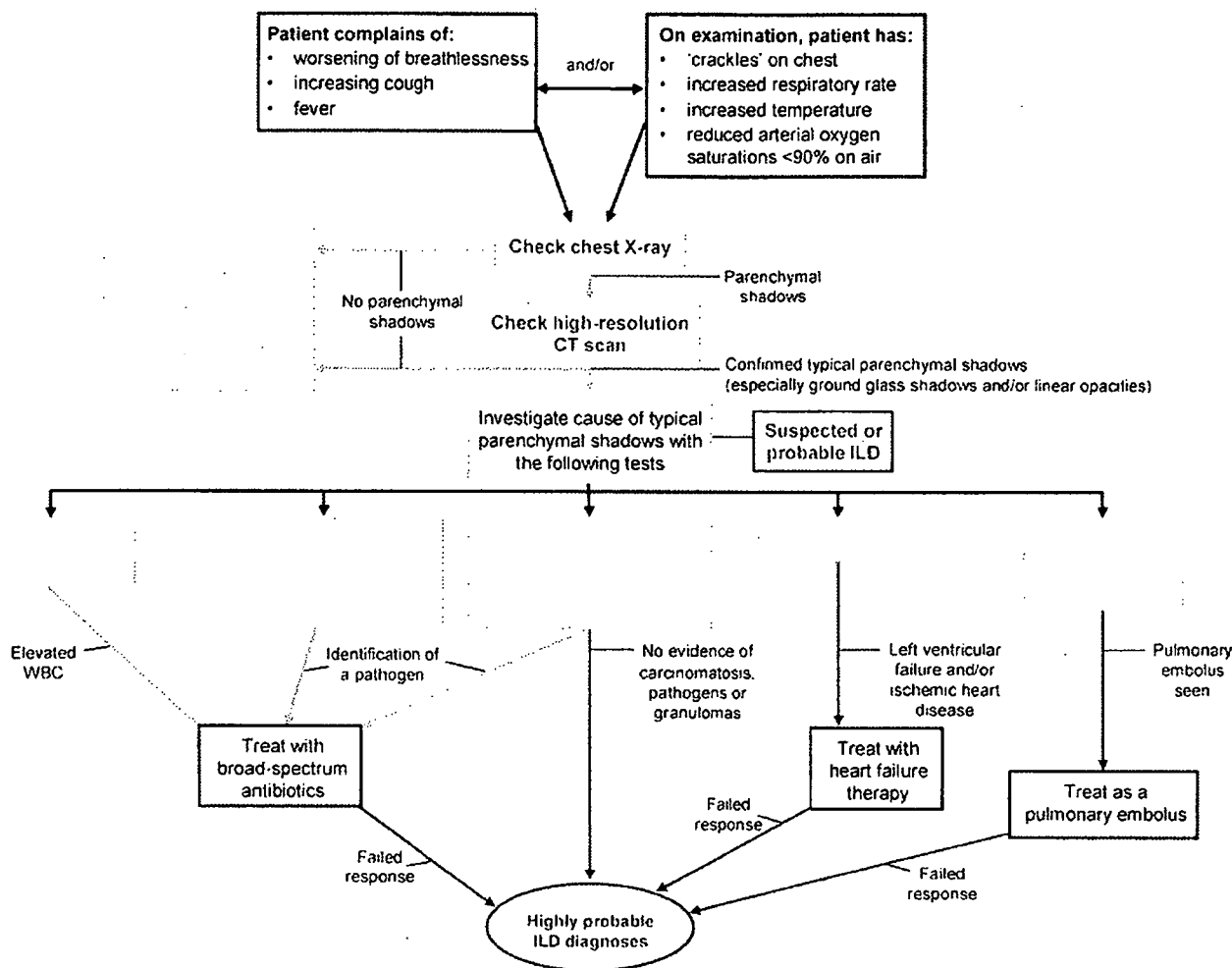
Patient complains of:
- worsening of breathlessness
- increasing cough
- fever

and/or

On examination, patient has:
- 'crackles' on chest
- increased respiratory rate
- increased temperature
- reduced arterial oxygen saturations <90% on air

Check chest X-ray

No parenchymal shadows

Parenchymal shadows

Check high-resolution CT scan

Confirmed typical parenchymal shadows (especially ground glass shadows and/or linear opacities)

Investigate cause of typical parenchymal shadows with the following tests

Suspected or probable ILD

Elevated WBC

Identification of a pathogen

No evidence of carcinomatosis, pathogens or granulomas

Left ventricular failure and/or ischemic heart disease

Pulmonary embolus seen

Treat with broad-spectrum antibiotics

Treat with heart failure therapy

Treat as a pulmonary embolus

Failed response

Failed response

Failed response

Highly probable ILD diagnoses

**Figure 1.** Algorithm for diagnosis of interstitial lung disease (ILD) in non-small cell lung cancer (NSCLC) patients.

As an example, in our work with gefitinib, samples were taken after obtaining informed consent from a nested case-control study, i.e., a case-control study performed within a prospective pharmacoepidemiological cohort of several thousand patients with advanced or recurring NSCLC who had received at least one prior chemotherapy regimen, and who were to be treated with gefitinib or chemotherapy. The main objective of this study was to measure the relative risk of ILD in Japanese patients with NSCLC using gefitinib compared with conventional therapy, with the associated aims of determining the incidence rate of ILD in late stage NSCLC patients and the principal risk factors for this complication.

Central to both the case-control study and the proteomics analysis was the use of internationally agreed criteria for the diagnosis of ILD and an algorithm of diagnostic tests to exclude alternative diseases.[28] Principal investigators in the study were asked to assess all patients for possible ILD using the diagnostic algorithm (Figure 1). Two case review boards of experts from oncology, radiology, and pulmonary medicine were set up to independently establish a consistent final diagnosis of ILD. In addition, extensive standard clinical and demographic risk factor data were collected on all registered cases and controls.

This degree of rigor in establishing accurate phenotypic diagnosis is critical to develop a robust and reliable personal-ized medicine test, as inaccuracies at this stage will affect all subsequent data analyses. The availability of clinical and risk factor data, and a rigorous epidemiological study design setting for the collection of proteomics samples is also of great value to fine-tune the statistical analysis.

## Is Proteomics Ready for Personalized Medicine Applications?

**The Human Proteome Map in Plasma.** The impetus to develop personalized medicine based on blood samples has encouraged proteomic profiling that identifies individual proteins and multiple "fingerprint" protein patterns. A remaining limitation has been the lack of integration of the technology of protein separation with bioinformatics and statistical methods. Extensive national and international[29,30] collaborations are being implemented to address some of these needs. An important component in this development is the Human Proteome Organization (HUPO; www.HUPO.org), a scientific consortium that supports various programmes to map the proteins expressed in various human tissues, disease states, etc.[31–33] One of these is the Plasma Proteome initiative started in 2002, aiming to annotate and catalog the many thousands of proteins and peptides[34–37] of the human plasma proteome. Recently results from the pilot phase with 35 collaborating laboratories from 13 countries[38–42] and multiple analytical

groups were made publicly available on the Internet (www.bioinformatics.med.umich.edu/hupo/ppp; www.ebi.ac.uk/pride). The combined efforts have generated 15 710 different MS/MS datasets that were linked to the International Protein Index (IPI) protein IDs, and an integration algorithm applied to multiple matches of peptide sequences yielded 9504 IPI proteins identified with one or more peptides[40] and characterized by Gene Ontology, InterPro, Novartis Atlas, and OMIM. Such advances provide an important platform for transforming proteomics from a technology to a useful biomarker tool applicable to personalized medicine.

**Protein Analysis in Blood—The Methods.** With respect to automated studies, multidimensional chromatography is the main technology used for protein analysis in blood. It is coupled to mass spectrometry either by electrospray ionization (ESI) for analysis in solution or matrix assisted laser desorption/ionization (MALDI) in solid phase applications.[39,41,43–47] Alternatively, ion-trap mass spectrometers are gaining recognition for high-throughput sequencing.[46,48–53] Linking a Fourier transform ion cyclotrone resonance (FTICR) unit to the linear trap can increase the resolution profoundly,[36,54–56] one of several novel principles for strengthening the assignment of protein annotations with the most commonly used protein search engines.[36,47,51–61] For protein annotation, the recent development of a human protein reference database complements these technologies.[61] Studies of protein expression in a variety of biological compartments ranging from sub-cellular to whole organisms have been undertaken with these analytic approaches.[62–70] Some key findings from the HUPO initiatives that impact on methodology include:

- For studies using blood samples, plasma rather than serum is preferred, with ethylenediaminetetraacetic acid (EDTA) as an anticoagulant.[49]

- The abundant proteins in plasma should be depleted prior to analysis.[40]

- Acceptance of protein annotation, i.e., accepted protein identities[39,40] should use standard criteria. These include having two identified peptide sequences from each protein, both with a statistical significance score high enough to ensure a correct sequence confirmation when compared with the corresponding gene sequence entity.[39]

Despite the advances in methodology, important hurdles to using proteomics in a personalized medicine context remain.

**Protein Expression Analysis in Blood—Some Important Hurdles.** Although protein profiling technology is highly automated and interfaced with database search engines to relate peptide sequences to protein identities and function,[39,40] there are many practical reasons why determining the relative abundance of proteins relevant for prediction purposes is difficult:

- About 90% of proteins are believed to be present only in low copy numbers, i.e., at medium and low abundance levels.[19]

- There can be variation both in the quantity and form of protein expression within normal physiological function.

- Between 300 000 and 3 million human protein species exist as direct gene products or post-translational modifications.[44]

- The relative abundance of the post-translational modifications occurring within the cell is called a Cell-Protein-Index Number (CPIN).[29,30] As an example, if one considers that there are 30 types of phosphorylation variants of a single phosphoprotein, and a hundred possible fold forms of glycosylation of a single glycoprotein, the theoretical CPIN varies considerably depending on the sample complexity.

- The dynamic range of protein expression within cells, between levels of most and least abundant proteins, is in the order of $10^8$–$10^{10}$.[34–36]

- In a typical clinical proteomics study the total cellular protein material in a sample seldom exceeds 10–20 milligrams. Therefore, the least abundant proteins would be present at starting levels not exceeding picograms.

- Recent studies use technology that can identify several thousand proteins in plasma samples,[39] but this still probably only represents a small fraction of the intermediate and processed protein forms. This is due to the current limitation of mass spectrometry not being able to ionize all amino acid sequences and protein modifications with equal efficiency. In most situations, a limited region of the full length protein is sequence annotated.

- The detection of differences in protein expression between groups of interest (e.g., cases and controls) takes place against a background of high variation between individuals within a group, within individuals over time and possible analytic run-to-run variation. Any method used to address this hurdle (which will involve "alignment" for spectral methods) directly impacts the ability to find good protein biomarkers.

Beyond the hurdles above, the fundamental challenge of protein biomarkers is to link the relative abundance of single markers or a fingerprint to clinically important biological processes based on some direct or indirect cause-effect link[29] related to normal or aberrant biological pathways.[47,49] In the following sections, we present the approach used for the identification of protein biomarkers potentially associated with development of ILD in NSCLC patients within the case-control study used as our motivating example. We build on the foundations described above and introduce further analytic developments and ideas relating to proteomic data generation, assaying and alignment to build a proteomics toolkit that can be applied today for personalized medicine approaches.

## A State of the Art Clinical Biomarker Analysis System

In the previous section, we described several challenges in proteomic analysis. Here we describe a system and analysis approaches that we have successfully implemented to address some of these issues.

**The Components of the Analysis System.** The analysis system (Figure 2) uses liquid chromatography-based high-resolution separation of peptides with an interface to tandem MS/MS, a technology which has been attracting great attention as the "shotgun" method of proteome analysis.[44,68–70] With this technology, after depletion of albumin and immunoglobulin G (IgG), all extracted plasma proteins are digested into their specific peptide components by proteolytic enzyme treatment.

The generated peptides are subjected to capillary reverse-phase submicro- to micro-flow liquid chromatography (capillary RP μLC), separated by retention times due to their physicochemical properties, and then detected and sequenced by a linear ion-trap tandem mass spectrometer[71] (LTQ, Thermo Fisher Scientific, San Jose, CA) interfaced with a spray needle tip for ESI of peptides.[70] A two-dimensional quadrupole ion trap mass spectrometer[71] is used, operated in a data-dependent acquisition mode with operational $m/z$ range limits set at 450–2000 (Figure 3, graphs A and B). Automatic switching to MS/MS acquisition mode is made in 1-second scanning cycles, controlled by the XCalibur software. The actual differences between annotated peptide fragment peaks shown in Figure 3, graph C, correspond to the amino acid residue mass, i.e.,
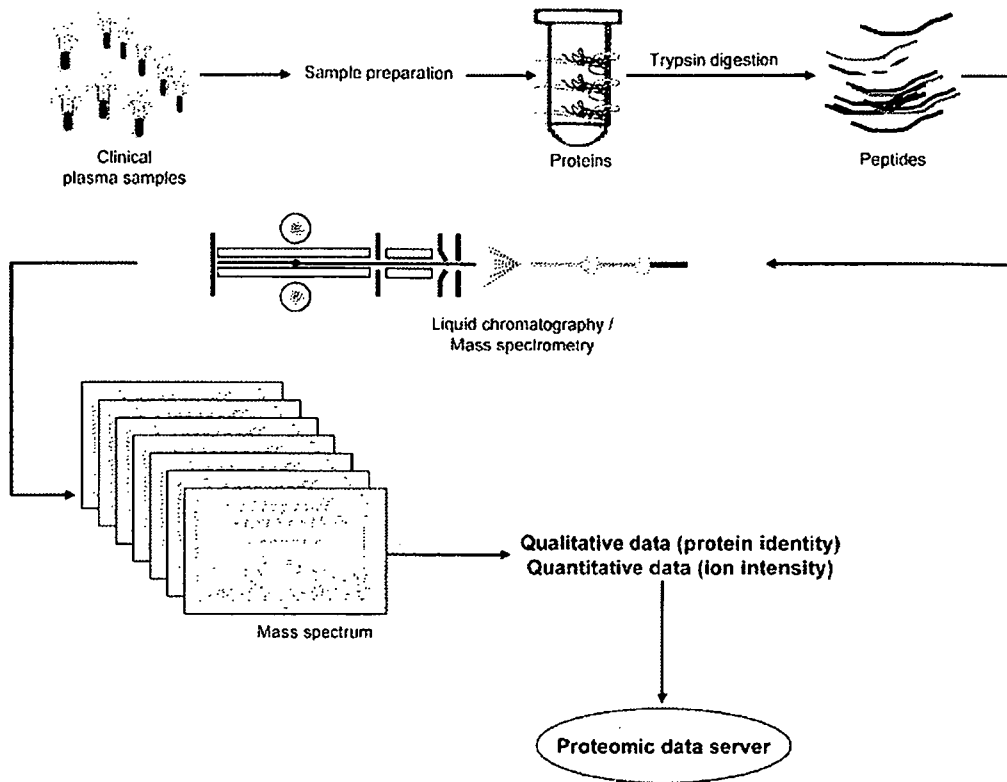
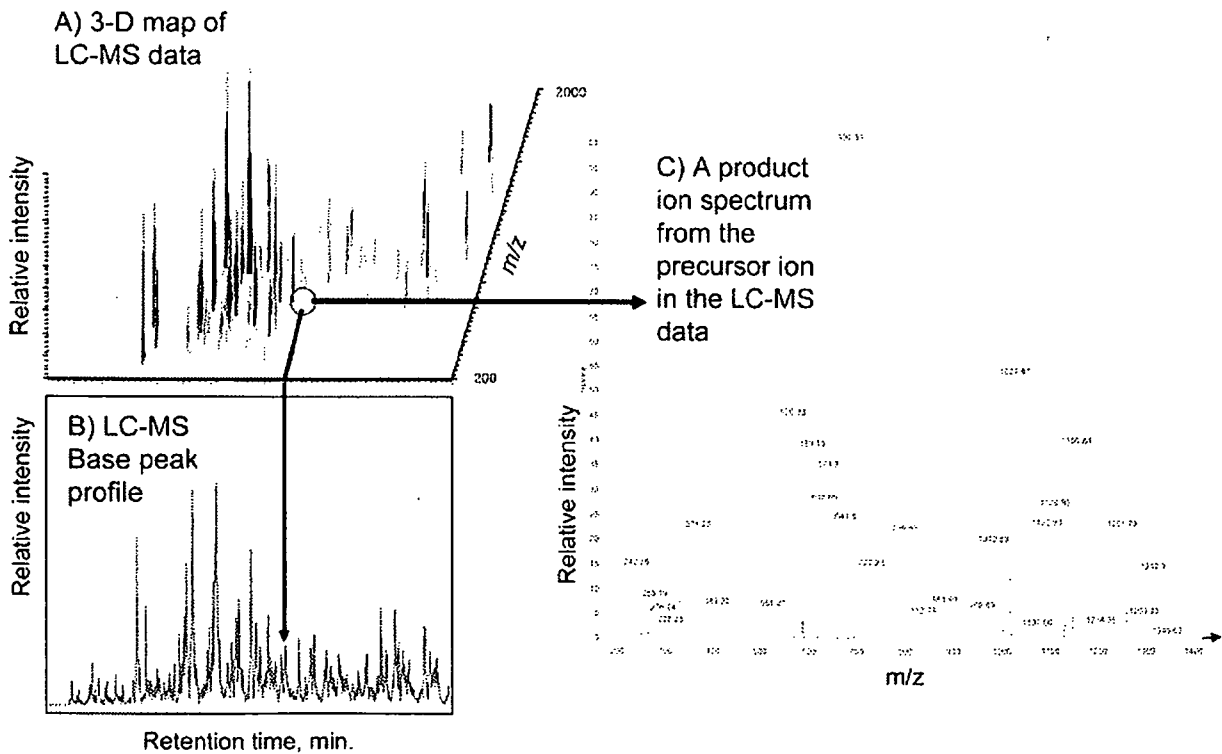**Figure 2.** Schematic illustration of the clinical proteomics screening process.



**Figure 3.** Profile of LC–MS data: (a) the three-dimensional view of LC-MS data, (b) the base-peak mass chromatogram, and (c) a product ion spectrum measured for a precursor ion in data-dependent acquisition mode (with MS acquisition operational *m/z* range set at 450–2000).
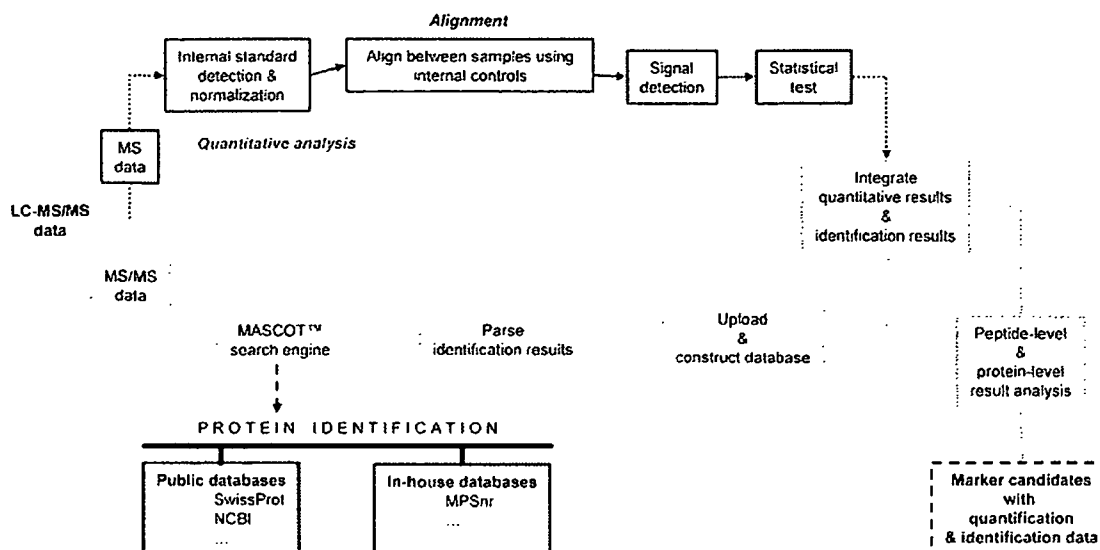
**Figure 4.** Overview of the data acquisition and database mining process developed within the gefitinib biomarker study.

identify the correct amino acid sequence. Internal standards are used for alignment of retention-times.

**How the Methodology Overcomes Some of the Hurdles.** The system described above addresses some of the hurdles noted previously. The digestion of all extracted plasma proteins into peptides will reduce the complexity by combining high-resolution nanoflow chromatographic fractionation with the separation power of modern mass spectrometry, performing automated and unattended shotgun sequencing in plasma.[35] Peptides are also more soluble and easier to handle than intact proteins. In addition, the two-dimensional quadrupole ion trap mass spectrometer[71] operates with a high-volume quadrupole electric field that makes it highly efficient to trap ions. The result is high sensitivity, high scanning speed, and better quantification over a wide dynamic range in comparison with the conventional three-dimensional ion-trap instruments.

Finding signals against a background of high variation is a further challenge, and the next section describes some approaches for addressing these.

## Initial Data Handling, Processing, and Analysis

Proteomic data analysis process can be considered as consisting of two components (Figure 4). *Quantitative analysis* is used to discover significant differences in peptide signal intensities by comparing two (or more) sample groups. This process uses data collected from an entire MS run to quantify the amount of peptide ions by their respective ion signal intensity. *Qualitative analysis* is used to identify the amino acid sequence of each peptide ion, from the respective product ion spectra. To maximize their value, the results from the two component analyses should be considered in combination.

A typical quantitative analysis may consist of several steps:

1. Normalization: To account for differences in the original sample concentrations. Typically, the total signal intensity is scaled to a constant value for each analyzed sample.

2. Alignment: Correcting for nonlinear fluctuation in retention time between different samples. A variety of methodologies are available for aligning LC—MS data sets. We have found the i-OPAL algorithm (Patent # WO 2004/090526 AI), which is based on the single linkage clustering algorithm[72] and which makes

use of internal standard signals, to perform well. Other alignment algorithms include xcms.[73]

3. Peak picking or signal detection: Identifying individual peptide ions within the data.

4. Identify discriminating peptides: A number of methods can be used, often in combination. A common approach is to apply a Student's $t$-test and select peptides which are significant, i.e., with a $p$-value less than the chosen cutoff, and which also show a fold-change or intensity ratio greater than another criterion. Further developments of this aspect are discussed in the Principled Statistical Analysis section.

A popular choice for qualitative analysis is the MASCOT MS/MS ion search program.[74] This may be run against a number of different peptide sequence databases, for example the NCBI Nr, Refseq, Gene Ontology, HUGO, and Swiss-Prot sequence databases. The results of the quantitative analysis can then be combined with the qualitative analysis so that, for example, a peptide must be both discriminating and have annotation—i.e., have achieved a high MASCOT score showing confidence in identification—to be considered a candidate biomarker.

The approaches we have discussed above are focused on finding potentially discriminating proteins of clinical utility. In the following section, we describe the next stage in our thinking, namely how we could rapidly deploy in the clinic a viable method for exploiting a predictive proteomic fingerprint.

## A Proposal for Proteomics in the Clinical Setting: Mass Spectrometric Biomarker Assays - MSBA

Although today's technology allows for high-throughput analyses of many proteins rather than a single protein,[30] the details of how such multiplexing assays will be adapted for clinical use have not been well clarified. The Mass Spectrometric Biomarker Assay (MSBA) platform described here was conceived as one example of a rapid and seamless method to progress from identification of a diagnostic more directly to a clinically useful test. MSBA requires only a minute sample amount (5—20 $\mu$L) to obtain a read-out from a handful of quantified protein biomarkers (typically 3—35) and automatically analyzes proteins using liquid-phase separation and tandem mass spectrometry with simultaneous quantitation and identification.
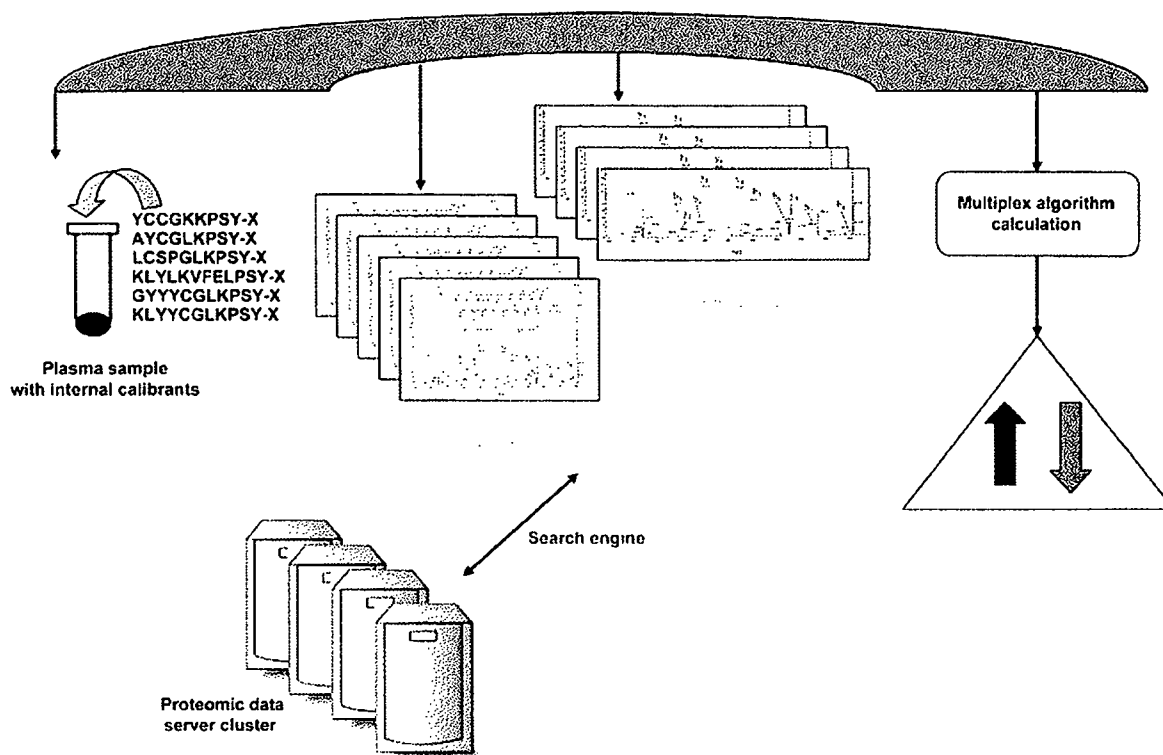
249

**Figure 5.** Entire flow of the operational components of Mass Spectrometric Biomarker Assays (MSBA).

The MSBA builds on a pre-defined Multiplex Biomarker list, which is stored within the MSBA database. Each marker entity has the values of masses and the relative retention time index with tolerance parameters. In running a patient sample, the predefined biomarker list is scanned to pick up patient sample signals that match with one of the predefined biomarker signals by satisfying the tolerance criteria (in general $\pm 1$ for $m/z$ value and $\pm 2\%$ for relative retention time index). The selected candidate signals are further confirmed using the product ion spectrum. That is, the product ion spectrum is represented as a vector by binning (grouping) the $m/z$ ratio values. Using the cosine correlation between the sample vectors and the reference vectors, we can confirm whether the selected candidate signals are truly assigned as target biomarkers. (A standard threshold value of the cosine correlation is 0.8.)

The process steps within the MSBA cycle are outlined in Figure 5. The calculation of the final multiplex biomarker assay read-out from all of the individual markers can be performed by a variety of applications, as discussed in more detail in the Principled Statistical Modeling Approach section. Figures 6A and B illustrate one approach, calculating a distance score which indicates to what extent a measured sample is distant from the case or control template in terms of predefined multiplex biomarkers.

$$S_{case\ or\ control} =$$

$$\sqrt{\left[\frac{1}{n(n-2)}\right]\left[n\sum_i y_i^2 - \left(\sum_i y_i\right)^2 - \frac{\left[n\sum_i x_iy_i - \left(\sum_i x_i\right)\left(\sum_i y_i\right)\right]^2}{n\sum_i x_i^2 - \left(\sum_i x_i\right)^2}\right]}$$

If the ratio of $S_{case}$ and $S_{control}$ exceeds an MSBA threshold parameter, then the test sample is predicted to be a patient susceptible to develop ILD (ILD case); if not, the test sample is predicted to be a non-susceptible patient (control). We are currently evaluating the MSBA approach in practice.

## A Principled Statistical Modeling Approach

We have described an analytical approach based on proteomic data, with various novel developments. However, additional insight is needed to further improve model discrimination and to broaden the focus from the proteomic data to the ultimate goal of prediction using combinations of data. Statistical analysis can be used to provide further refinement by combining information from the full clinical and laboratory datasets.

An advantage of a multiple biomarker approach (e.g., proteomics) compared with standard single biomarker development is the capability to combine information from many different entities. An example is illustrated in Figure 7A. Considering each biomarker alone fails to separate the two groups of subjects, as there is considerable overlap for both biomarkers. Use of two biomarkers in combination completely separates the two groups.

We can also use clinical variables to advantage in the analysis of the peptide patterns. For example, the efficacy of gefitinib appears to be greater in non-smokers, women, patients of Asian origin, and patients with adenocarcinomas.[a] Figure 7B illustrates how, instead of two protein biomarkers, the combination of clinical data (e.g., age) and a proteomic biomarker is able to separate two groups.

On this basis, we propose using a principled statistical analysis approach to first explore and understand the data and
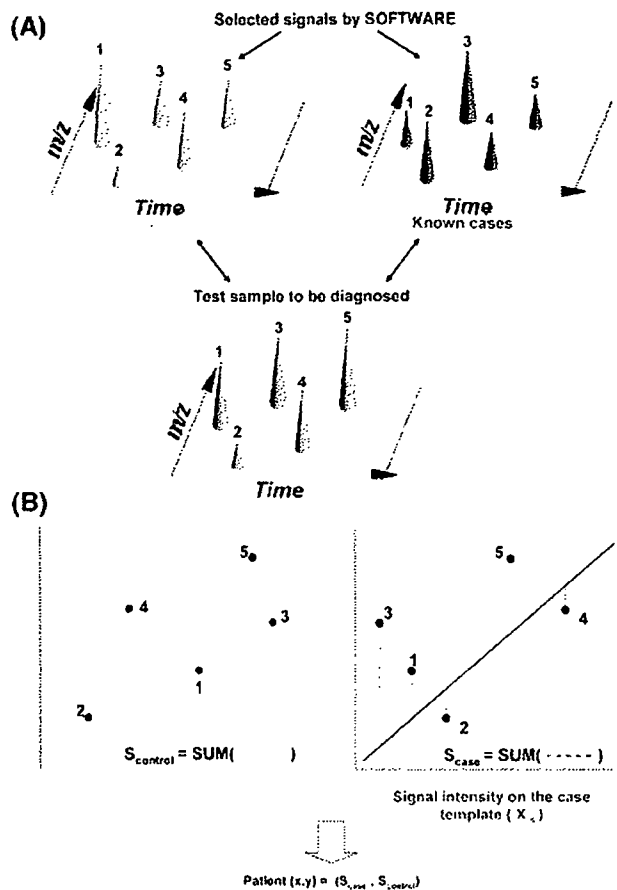
## (A)



## (B)



**Figure 6.** (A) Peptide signal comparison that MSBA (Mass Spectrometric Biomarker Assays) performs of the generated ions from the sample. The comparison is made both with the pattern of the controls and with the pattern of the case group for the corresponding signals. (B) Illustration of the regression model application of the MSBA where control templates and case templates are compared to that of the sample template generated in the analysis process.

then to model it and understand the quality of any models produced. A first step is to perform exploratory data analysis (EDA), for example using principal components analysis (PCA), to understand the major sources of data variation and the covariation between clinical parameters and protein intensity measures. The next step is univariate modeling for each protein marker individually, for example using analysis of covariance (ANCOVA), and an assessment of the effect of clinical parameters across the whole set of protein biomarkers using, for example, the False Discovery Rate as a tool.[75] This provides an understanding of key clinical variables and sources of variation within the data.

The next step is to perform multivariate predictive modeling using the proteins and clinical variables identified as being potentially important. There are a number of mathematical methods described in the literature for performing supervised classification, for example Support Vector Machines,[76] Random Forests,[77] PAM,[78] all of which have been successfully applied to high dimensional genomics data.[79] It remains an important unanswered question which modeling approach, or combination of modeling approaches, will generate the most predictive and robust models for data generated using this technology within a prospective study of this design.

Finally, to confirm that we have a practical prediction, the predictive power of a model must be assessed on a different set of patients from that used to generate the model. There are a number of approaches for external validation given a limited size dataset, for example the sequential approach of building a model based upon currently available data and testing on data from new patients when they become available, or withholding an arbitrary selection of subjects from the modeling as a test set and testing the model on these subjects. Internal validation approaches such as cross-validation or related bootstrapping methods may also be useful to assess the model selection procedure, but tend to overestimate the performance of a specific predictive model in subsequent external validation.[80,81] The key properties to consider when selecting an assessment method are to ensure that it will provide both precise and unbiased information regarding the prediction error rate of the potential model to be tested for clinical use. As well as assessing an overall predictive rate, it is also useful to separately assess the predictive rate for both the cases and controls and to consider the relative costs of making these false predictions within a clinical setting. Finally, the prevalence of the condition in question (here ILD) is also a critical factor in estimating what proportion of people predicted to be at risk are truly at risk, and this should also be borne in mind when evaluating a model for potential clinical use. The recently published FDA concept paper on drug-diagnostic co-
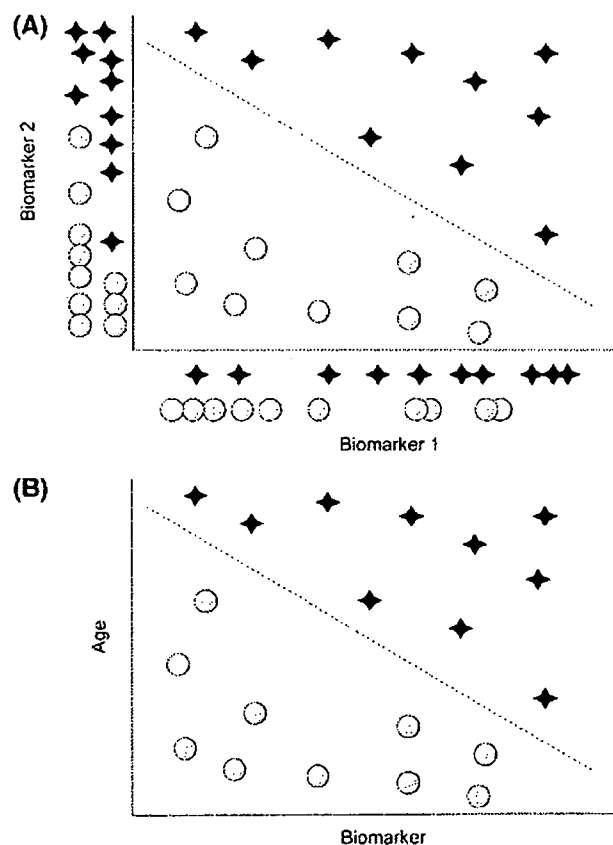
## (A)



## (B)



**Figure 7.** (A) Hypothetical example of the combined disease-linkage effect of two protein biomarkers. (Stars signify affected case individuals, circles non-affected control individuals). (B) Hypothetical example of the combined disease-linkage effect of a biomarker and a clinical variable. (Stars signify affected case individuals, circles non-affected control individuals).
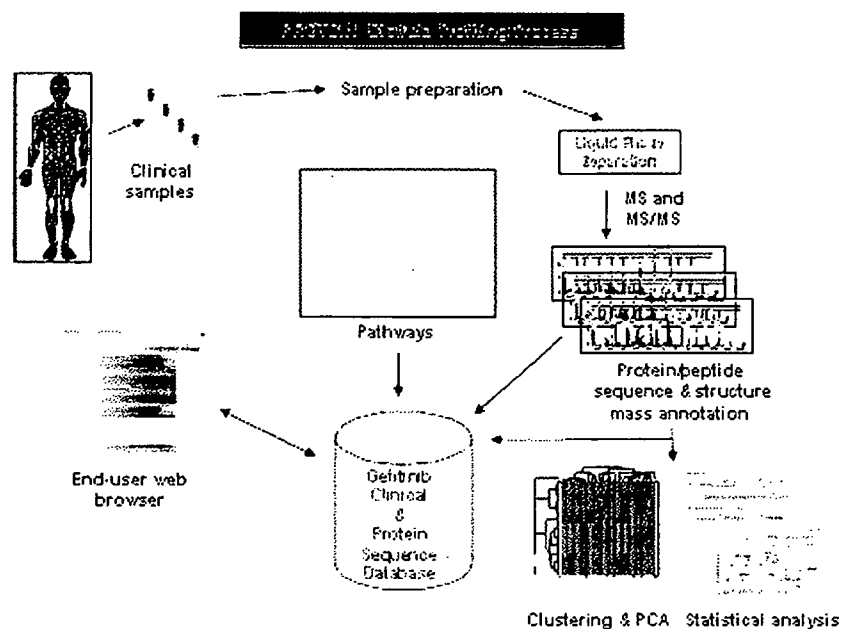
251

**Figure 8.** Illustration of the bioinformatics and data processing structure within which MSBA (Mass Spectrometric Biomarker Assays) data are captured, modified and analyzed.

development discusses many of the issues around validating predictive biomarkers.[82]

Finally, it is preferable to be able to assign a biological rationale to the biomarkers. Confidence in the reliability of a biomarker is greatly enhanced if we can correctly understand how it relates to the mechanism and progression of the disease of interest. Figure 8 illustrates a bioinformatics and data processing structure that we have developed to allow us to both conduct interactive exploratory and statistical analyses, and also investigate the disease and pathway linkage of discovered biomarker proteins through direct access to reference databases.

## Future Perspectives

Within this paper we have discussed many of the issues that need to be considered in developing a personalized medicine approach. A key starting point is that rigorous steps are taken to ensure accurate diagnosis and the careful gathering of both clinical and proteomic data to facilitate the search for peptide patterns.

There are many challenges in performing protein analysis in blood, but mass spectrometry equipment and methods can now be used to generate peptide data with high sensitivity, high scanning speed, and improved quantification. Data handling and processing techniques for steps such as peak alignment and the subsequent methodologies for statistical modeling and analysis are now far enough developed to generate high quality data and robustly analyze these data with confidence.

We have provided details of the MSBA method that can be used to easily translate protein intensities into a practical multiplex assay which can be exploited in the clinic without the need to develop anti-bodies for ELISA. We have also described how an expanded statistical analysis can be used to allow for the individual variance of protein expression to enable us to focus on the proteomic patterns that are actually related to ILD. Finally, we have emphasized the importance of validat-

ing the predictive power of a biomarker tool in a way that reflects the real-life setting of intended clinical use.

Hopefully, this combination of developments over a range of different areas brings us one step closer to a practical personalized medicine.

IRESSA is a trademark of the AstraZeneca group of companies.

## References

(1) Thatcher, N.; Chang, A.; Parikh, P.; Pereira, J. R.; Ciuleanu, T.; von Pawel, J.; et al. Gefitinib plus best supportive care in previously treated patients with refractory advanced non-small-cell lung cancer: results from a randomised, placebo-controlled, multicentre study (Iressa Survival Evaluation in Lung Cancer). *Lancet* **2005**, *366*, 1527—1537.

(2) Hirsch, F. R.; Varella-Garcia, M.; McCoy, J.; West, H.; Xavier, A. C.; Gumerlock, P.; et al. Increased epidermal growth factor receptor gene copy number detected by fluorescence in situ hybridization associates with increased sensitivity to gefitinib in patients with bronchioloalveolar carcinoma subtypes: a Southwest Oncology Group study. *J. Clin. Oncol.* **2005**, *23*, 6838—6845.

252

(3) Cappuzzo, F.; Varella-Garcia, M.; Shigematsu, H.; Domenichini, I.; Bartolini, S.; Ceresoli, G. L.; et al. Increased HER2 gene copy number is associated with response to gefitinib therapy in epidermal growth factor receptor-positive non-small-cell lung cancer patients. *J. Clin. Oncol.* 2005, *23*, 5007–5018.

(4) Araki, J.; Okamoto, I.; Suto, R.; Ichikawa, Y.; Sasaki, J. Efficacy of the tyrosine kinase inhibitor gefitinib in a patient with metastatic small cell lung cancer. *Lung Cancer* 2005, *48*, 141–144.

(5) Kim, K. S.; Jeong, J. Y.; Kim, Y. C.; Na, K. J.; Kim, Y. H.; Ahn, S. J.; et al. Predictors of the response to gefitinib in refractory non-small cell lung cancer. *Clin. Cancer Res.* 2005, *11*, 2244–2251.

(6) Lynch, T. J.; Bell, D. W.; Sordella, R.; Gurubhagavatula, S.; Okimoto, R. A.; Brannigan, B. W.; et al. Activating mutations in the epidermal growth factor receptor underlying responsiveness of non-small-cell lung cancer to gefitinib. *N. Engl. J. Med.* 2004, *350*, 2129–2139.

(7) Paez, J. G.; Jänne, P. A.; Lee, J. C.; Tracy, S.; Greulich, H.; Gabriel, S.; et al. EGFR mutations in lung cancer: correlation with clinical response to gefitinib therapy. *Science* 2004, *304*, 1497–1500.

(8) Shigematsu, H.; Lin, L.; Takahashi, T.; Nomura, M.; Suzuki, M.; Wistuba II; et al. Clinical and biological features associated with epidermal growth factor receptor gene mutations in lung cancers. *J. Natl. Cancer Inst.* 2005, *97*, 339–346.

(9) American Thoracic Society: American Thoracic Society/European Respiratory Society International Multidisciplinary Consensus Classification of the Idiopathic Interstitial Pneumonias. This joint statement of the American Thoracic Society (ATS), and the European Respiratory Society (ERS) was adopted by the ATS Board of Directors, June 2001 and by The ERS Executive Committee, June 2001. *Am. J. Respir. Crit. Care Med.* 2002, *165*, 277–304.

(10) Raghu, G.; Nyberg, F.; Morgan, G. The epidemiology of interstitial lung disease and its association with lung cancer. *Br. J. Cancer* 2004, *91* (Suppl. 2), S3–S10.

(11) Asada, K.; Mukai, J.; Ougushi, F. Characteristics and management of lung cancer in patients with idiopathic pneumonia. *Jap. J. Thor. Dis.* 1992, *51*, 214–219.

(12) Hubbard, R.; Venn, A.; Lewis, S.; Britton, J. Lung cancer and cryptogenic fibrosing alveolitis. A population-based cohort study. *Am. J. Respir. Crit. Care Med.* 2000, *161*, 5–8.

(13) Matsushita, H.; Tanaka, S.; Saiki, Y.; Hara, M.; Nakata, K.; Tanimura, S.; et al. Lung cancer associated with usual interstitial pneumonia. *Pathol. Int.* 1995, *45*, 925–932.

(14) Ogura, T.; Kondo, A.; Sato, A.; Ando, M.; Tamura, M. Incidence and clinical features of lung cancer in patients with idiopathic interstitial pneumonia. *Nihon Kyobu Shikkan Gakkai Zasshi* 1997, *35*, 294–299.

(15) Takeuchi, E.; Yamaguchi, T.; Mori, M.; Tanaka, S.; Nakagawa, M.; Yokota, S.; et al. Characteristics and management of patients with lung cancer and idiopathic interstitial pneumonia. *Nihon Kyobu Shikkan Gakkai Zasshi* 1996, *34*, 653–658.

(16) Turner-Warwick, M.; Lebowitz, M.; Burrows, B.; Johnson, A. Cryptogenic fibrosing alveolitis and lung cancer. *Thorax* 1980, *35*, 496–499.

(17) Baumgartner, K. B.; Samet, J. M.; Stidley, C. A.; Colby, T. V.; Waldron, J. A. Cigarette smoking: a risk factor for idiopathic pulmonary fibrosis. *Am. J. Respir. Crit. Care Med.* 1997, *155*, 242–248.

(18) Britton, J.; Hubbard, R. Recent advances in the aetiology of cryptogenic fibrosing alveolitis. *Histopathology* 2000, *37*, 387–392.

(19) Iwai, K.; Mori, T.; Yamada, N.; Yamaguchi, M.; Hosoda, Y. Idiopathic pulmonary fibrosis. Epidemiologic approaches to occupational exposure. *Am. J. Respir. Crit. Care Med.* 1994, *150*, 670–675.

(20) Nagai, S.; Hoshino, Y.; Hayashi, M.; Ito, I. Smoking-related interstitial lung diseases. *Curr. Opin. Pulm. Med.* 2000, *6*, 415–419.

(21) Lilly. Gemcitabine prescribing information. http://pi.lilly.com/gemzar.pdf, 2003.

(22) Kunitoh, H.; Watanabe, K.; Onoshi, T.; Furuse, K.; Niitani, H.; Taguchi, T. Phase II trial of docetaxel in previously untreated advanced non-small-cell lung cancer: a Japanese cooperative study. *J. Clin. Oncol.* 1996, *14*, 1649–1655.

(23) Merad, M.; Le Cesne, A.; Baldeyrou, P.; Mesurolle, B.; Le Chevalier, T. Docetaxel and interstitial pulmonary injury. *Ann. Oncol.* 1997, *8*, 191–194.

(24) Wang, G.-S.; Yan, K.-Y.; Perng, R.-P. Life-threatening hypersensitivity pneumonitis induced by docetaxel (taxotere). *Br. J. Cancer* 2001, *85*, 1247–1250.

(25) Erasmus, J. J.; McAdams, H. P.; Rossi, S. E. Drug-induced lung injury. *Semin. Roentgenol.* 2002, *37*, 72–81.

(26) Aviram, G.; Yu, E.; Tai, P.; Lefcoe, M. S. Computed tomography to assess pulmonary injury associated with concurrent chemoradiotherapy for inoperable non-small cell lung cancer. *Can. Assoc. Radiol. J.* 2001, *52*, 385–391.

(27) Yoshida, S. The results of gefitinib prospective investigation. *Med. Drug J.* 2005, *41*, 772–789.

(28) Mueller, N. L.; White, D. A.; Jiang, H.; Gemma, A. Diagnosis and management of drug-associated interstitial lung disease. *Br. J. Cancer* 2004, *91*, S24–S30.

(29) Marko-Varga, G.; Fehniger, T. E. Proteomics and disease—the challenges for technology and discovery. *J. Proteome Res.* 2004, *3*, 167–178.

(30) Marko-Varga, G.; Lindberg, H.; Lofdahl, C. G.; Jonsson, P. H. L.; Dahlback, M.; Lindquist, E.; et al. Discovery of biomarker candidates within disease by protein profiling: principles and concepts. *J. Proteome Res.* 2005, *4*, 1200–1212.

(31) Omenn, G. S. The Human Proteome Organization Plasma Proteome Project pilot phase: reference specimens, technology platform comparisons, and standardized data submissions and analyses. *Proteomics* 2004, *4*, 1235–1240.

(32) Omenn, G. S. Advancement of biomarker discovery and validation through the HUPO plasma proteome project. *Dis. Markers* 2004, *20*, 131–134.

(33) Orchard, S.; Hermjakob, H.; Binz, P. A.; Hoogland, C.; Taylor, C. F.; Zhu, W.; et al. Further steps towards data standardisation: the Proteomic Standards Initiative HUPO 3(rd) annual congress, Beijing 25–27(th) October, 2004. *Proteomics* 2005, *5*, 337–339.

(34) Anderson, N. G.; Matheson, A.; Anderson, N. L. Back to the future: the human protein index (HPI) and the agenda for post-proteomic biology. *Proteomics* 2001, *1*, 3–12.

(35) Anderson, N. L.; Anderson, N. G. The human plasma proteome: history, character, and diagnostic prospects. *Mol. Cell. Proteomics* 2002, *1*, 845–867.

(36) Jacobs, J. M.; Adkins, J. N.; Qian, W. J.; Liu, T.; Shen, Y.; Camp, D. G.; et al. Utilizing human blood plasma for proteomic biomarker discovery. *J. Proteome Res.* 2005, *4*, 1073–1085.

(37) Anderson, N. G.; Anderson, L. The Human Protein Index. *Clin. Chem.* 1982, *28*, 739–748.

(38) Haab, B. B.; Geierstanger, B. H.; Michailidis, G.; Vitzthum, F.; Forrester, S.; Okon, R.; et al. Immunoassay and antibody microarray analysis of the HUPO Plasma Proteome Project reference specimens: systematic variation between sample types and calibration of mass spectrometry data. *Proteomics* 2005, *5*, 3278–3291.

(39) Martens, L.; Hermjakob, H.; Jones, P.; Adamski, M.; Taylor, C.; States, D.; et al. PRIDE: the proteomics identifications database. *Proteomics* 2005, *5*, 3537–3545.

(40) Omenn, G. S.; States, D. J.; Adamski, M.; Blackwell, T. W.; Menon, R.; Hermjakob, H.; et al. Overview of the HUPO Plasma Proteome Project: results from the pilot phase with 35 collaborating laboratories and multiple analytical groups, generating a core dataset of 3020 proteins and a publicly-available database. *Proteomics* 2005, *5*, 3226–3245.

(41) Patterson, S. D. Data analysis—the Achilles heel of proteomics. *Nat. Biotechnol.* 2003, *21*, 221–222.

(42) Rahbar, A. M.; Fenselau, C. Integration of Jacobson's pellicle method into proteomic strategies for plasma membrane proteins. *J. Proteome Res.* 2004, *3*, 1267–1277.

(43) Ho, Y.; Gruhler, A.; Heilbut, A.; Bader, G. D.; Moore, L.; Adams, S. L.; et al. Systematic identification of protein complexes in Saccharomyces cerevisiae by mass spectrometry. *Nature* 2002, *415*, 180–183.

(44) Aebersold, R.; Mann, M. Mass spectrometry-based proteomics. *Nature* 2003, *422*, 198–207.

(45) Anderson, N. L.; Polanski, M.; Pieper, R.; Gatlin, T.; Tirumalai, R. S.; Conrads, T. P.; et al. The human plasma proteome: a nonredundant list developed by combination of four separate sources. *Mol. Cell. Proteomics* 2004, *3*, 311–326.

(46) Olsen, J. V.; Mann, M. Improved peptide identification in proteomics by two consecutive stages of mass spectrometric fragmentation. *Proc. Natl. Acad. Sci. U.S.A.* 2004, *101*, 13417–13422.

(47) Sadygov, R. G.; Liu, H.; Yates, J. R. Statistical models for protein validation using tandem mass spectral data and protein amino acid sequence databases. *Anal. Chem.* 2004, *76*, 1664–1671.

(48) Fujii, K.; Nakano, T.; Kanazawa, M.; Akimoto, S.; Hirano, T.; Kato, H.; et al. Clinical-scale high-throughput human plasma proteome analysis: lung adenocarcinoma. *Proteomics* 2005, *5*, 1150–1159.

253

(49) Campbell, J. M.; Collings, B. A.; Douglas, D. J. A new linear ion trap time-of-flight system with tandem mass spectrometry capabilities. *Rapid Commun. Mass Spectrom.* 1998, *12*, 1463–1474.

(50) Cha, B. C.; Blades, M.; Douglas, D. J. An interface with a linear quadropole ion guide for an electrospray-ion trap mass spectrometer system. *Anal. Chem.* 2000, *72*, 5647–5654.

(51) Hager, J. W. Product ion spectral simplification using time-delayed fragment ion capture with tandem linear ion traps. *Rapid Commun. Mass Spectrom.* 2003, *17*, 1389–1398.

(52) Syka, J. E.; Marto, J. A.; Bai, D. L.; Horning, S.; Senko, M. W.; Schwartz, J. C.; et al. Novel linear quadrupole ion trap/FT mass spectrometer: performance characterization and use in the comparative analysis of histone H3 post-translational modifications. *J. Proteome Res.* 2004, *3*, 621–626.

(53) Shen, Y.; Zhao, R.; Belov, M. E.; Conrads, T. P.; Anderson, G. A.; Tang, K.; et al. Packed capillary reversed-phase liquid chromatography with high-performance electrospray ionization Fourier transform ion cyclotron resonance mass spectrometry for proteomics. *Anal. Chem.* 2001, *73*, 1766–1775.

(54) Wu, S. L.; Kim, J.; Hancock, W. S.; Karger, B. Extended Range Proteomic Analysis (ERPA): a new and sensitive LC-MS platform for high sequence coverage of complex proteins with extensive post-translational modifications-comprehensive analysis of beta-casein and epidermal growth factor receptor (EGFR). *J. Proteome Res.* 2005, *4*, 1155–1170.

(55) Olsen, J. V.; de Godoy, L. M.; Li, G.; Macek, B.; Mortensen, P.; Pesch, R.; et al. Parts per million mass accuracy on an Orbitrap mass spectrometer via lock mass injection into a C-trap. *Mol. Cell. Proteomics* 2005, *4*, 2010–2021.

(56) Yates, J. R.; Cociorva, D.; Liao, L.; Zabrouskov, V. Performance of a linear ion trap-Orbitrap hybrid for peptide analysis. *Anal. Chem.* 2006, *78*, 493–500.

(57) Anderson, D. C.; Li, W.; Payan, D. G.; Noble, W. S. A new algorithm for the evaluation of shotgun peptide sequencing in proteomics: support vector machine classification of peptide MS/MS spectra and SEQUEST scores. *J. Proteome Res.* 2003, *2*, 137–146.

(58) Carr, S.; Aebersold, R.; Baldwin, M.; Burlingame, A.; Clauser, K.; Nesvizhskil, A. The need for guidelines in publication of peptide and protein identification data: working group on publication guidelines for peptide and protein identification data. *Mol. Cell. Proteomics* 2004, *3*, 531–533.

(59) Fenyo, D.; Beavis, R. C. A method for assessing the statistical significance of mass spectrometry-based protein identifications using general scoring schemes. *Anal. Chem.* 2003, *75*, 768–774.

(60) Nesvizhskii, A. I.; Keller, A.; Kolker, E.; Aebersold, R. A statistical model for identifying proteins by tandem mass spectrometry. *Anal. Chem.* 2003, *75*, 4646–4658.

(61) Peri, S.; Navarro, J. D.; Kristiansen, T. Z.; Amanchy, R.; Surendranath, V.; Muthusamy, B.; et al. Human protein reference database as a discovery resource for proteomics. *Nucleic Acids Res.* 2004, *32*, D497–D501.

(62) Kratchmarova, I.; Blagoev, B.; Haack-Sorensen, M.; Kassem, M.; Mann, M. Mechanism of divergent growth factor effects in mesenchymal stem cell differentiation. *Science* 2005, *308*, 1472–1477.

(63) Dreger, M.; Bengtsson, L.; Schöneberg, T.; Otto, H.; Hucho, F. Nuclear envelope proteomics: novel integral membrane proteins of the inner nuclear membrane. *Proc. Natl. Acad. Sci. U.S.A.* 2001, *98*, 11943–11948.

(64) Giot, L.; Bader, J. S.; Brouwer, C.; Chaudhuri, A.; Kuang, B.; Li, Y.; et al. A protein interaction map of Drosophilia melanogaster. *Science* 2003, *302*, 1727–1736.

(65) Johnson, J. R.; Florens, L.; Carucci, D. J.; Yates, J. R., III. Proteomics in malaria. *J. Proteome Res.* 2004, *3*, 296–306.

(66) Hirsch, J.; Hansen, K. C.; Burlingame, A. L.; Matthay, M. A. Proteomics: current techniques and potential applications to lung disease. *Am. J. Physiol. Lung Cell Mol. Physiol.* 2004, *287*, L1–L23.

(67) Malmström, J.; Larsen, K.; Hansson, L.; Löfdahl, C.-G.; Norregard-Jensen, O.; Marko-Varga, G.; et al. Proteoglycan and proteome profiling of central human pulmonary fibrotic tissue utilizing minaturized sample preparation: A feasibility study. *Proteomics* 2002, *2*, 394–404.

(68) Malmström, J.; Larsen, K.; Malmström, L.; Tufvesson, E.; Parker, K.; Marchese, J.; et al. Proteome annotations and identifications of the human pulmonary fibroblast. *J. Proteome Res.* 2004, *3*, 525–537.

(69) Oh, P.; Li, Y.; Yu, J.; Durr, E.; Krasinska, K. M.; Carver, L. A.; et al. Subtractive proteomic mapping of the endothelial surface in lung and solid tumours for tissue-specific therapy. *Nature* 2004, *429*, 629–635.

(70) Fujii, K.; Nakano, T.; Kawamura, T.; Usui, F.; Bando, Y.; Wang, R.; et al. Multidimensional protein profiling technology and its application to human plasma proteome. *J. Proteome Res.* 2004, *3*, 712–718.

(71) Schwartz, J. C.; Senko, M. W.; Syka, J. E. A two-dimensional quadrupole ion trap mass spectrometer. *J. Am. Soc. Mass Spectrom.* 2002, *13*, 659–669.

(72) Sneath, P. H. A.; Sokai, R. R. *Numerical Taxonomy. The principles and practice of numerical classification*; W. H. Freeman and Co.: San Francisco, 1973.

(73) Smith, C. A.; Want, E. J.; O'Maille, G.; Abagyan, R.; Siuzdak, G. XCMS: processing mass spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and identification. *Anal. Chem.* 2006, *78*, 779–787.

(74) Perkins, D. N.; Pappin, D. J.; Creasy, D. M.; Cottrell, J. S. Probability-based protein identification by searching sequence-databases using mass spectrometry data. *Electrophoresis* 1999, *20*, 3551–3567.

(75) Storey, J. A direct approach to false discovery rates. *J. R. Stat. Soc. Ser. B* 2002, *64*, 479.

(76) Vapnik, V. *Statistical Learning Theory*; Wiley: Chichester, UK, 1998.

(77) Breiman, L. Random forests. *Mach. Learn.* 2001, *45*, 5–32.

(78) Tibshirani, R.; Hastie, T.; Narasimhan, B.; Chu, G. Diagnosis of multiple cancer types by shrunken centroids of gene expression. *Proc. Natl. Acad. Sci. U.S.A.* 2002, *99*, 6567–6572.

(79) Lee, J. W.; Lee, J. B.; Park, M.; Song, S. H. An extensive comparison of recent classification tools applied to microarray data. *Comp. Stat. Data Anal.* 2005, *48*, 869–885.

(80) Steyerberg, E. W.; Harrell, F. E., Jr.; Borsboom, G. J.; Eijkemans, M. J.; Vergouwe, Y.; Habbema, J. D. Internal validation of predictive models: efficiency of some procedures for logistic regression analysis. *J. Clin. Epidemiol.* 2001, *54*, 774–781.

(81) Blecker, S. E.; Moll, H. A.; Steyerberg, E. W.; Donders, A. R.; Derksen-Lubsen, G.; Grobbee, D. E.; et al. External validation is necessary in prediction research: a clinical example. *J. Clin. Epidemiol.* 2003, *56*, 826–832.

(82) Food and Drug Administration (FDA): Drug-diagnostic co-development concept paper. Draft–not for implementation. http://www.fda.gov/cder/genomics/pharmacoconceptfn.pdf, 2005.

PR070046S