

Evidence of different metabolic phenotypes in humans

Michael Assfalg^{*†}, Ivano Bertini^{*§¶}, Donato Colangiuli^{||**}, Claudio Luchinat^{††}, Hartmut Schäfer^{††}, Birk Schütz^{††}, and Manfred Spraul^{††}

^{*}ProtEra, Viale delle Idee 22, 50019 Sesto Fiorentino, Italy; [†]Magnetic Resonance Center and ^{**}Center of Metabolomics, University of Florence, Via Luigi Sacconi 6, 50019 Sesto Fiorentino, Italy; [§]Department of Chemistry, University of Florence, Via della Lastruccia 3, 50019 Sesto Fiorentino, Italy; [¶]FiorGen Foundation, Via Luigi Sacconi 6, 50019 Sesto Fiorentino, Italy; ^{||}Department of Agricultural Biotechnology, University of Florence, Via Maragliano 75-77, 50144 Florence, Italy; and ^{††}Bruker BioSpin, Silberstreifen, D-76287 Rheinstetten, Germany

Edited by Joan Selverstone Valentine, University of California, Los Angeles, CA, and approved December 12, 2007 (received for review June 18, 2007)

The study of metabolic responses to drugs, environmental changes, and diseases is a new promising area of metabolomic research. Metabolic fingerprints can be obtained by analytical techniques such as nuclear magnetic resonance (NMR). In principle, alterations of these fingerprints due to appearance/disappearance or concentration changes of metabolites can provide early evidences of, for example, onset of diseases. A major drawback in this approach is the strong day-to-day variability of the individual metabolic fingerprint, which should be rather called a metabolic “snapshot.” We show here that a thorough statistical analysis performed on NMR spectra of human urine samples reveals an invariant part characteristic of each person, which can be extracted from the analysis of multiple samples of each single subject. This finding (i) provides evidence that individual metabolic phenotypes may exist and (ii) opens new perspectives to metabolomic studies, based on the possibility of eliminating the daily “noise” by multiple sample collection.

biofluids | metabolomics | metabonomics | NMR | urine

The global analysis of the type and quantity of metabolites in biological fluids, tissues, or related biological samples—i.e., the study of the metabolome, or metabolomics (1)—is a promising area of research, because of the potential relevance for human health of the study of metabolic responses to pathophysiological stimuli or genetic modifications—termed metabonomics (2). The relevance of metabonomics could be greatly enhanced if it were possible to identify an invariant part of the individual metabolome of a “healthy” subject with respect to, for example, pathological states, in such a way as to perform prediction, early diagnosis and prognosis of pathologies. Indeed, traditional biomedical/clinical approaches are limited by the number of parameters as well as in their efficiency, and they provide only a fragmented perspective on the health status of an individual.

Differences in experimental metabolic profiles due to genetic strain differences in animal models have been observed, leading to the suggestion that each individual or group of individuals may be characterized by a different metabotype, defined as a “multiparametric description of an organism in a given physiological state,” based on metabolomic data (3). The availability of metabotypes characteristic of an individual and stable over time could be fundamental in nutrigenomics (4, 5), in evaluation of drug efficacy, in pharmacometabonomics (6), and in studies of personalized nutrition aimed at maintaining metabolic health and avoiding loss of homeostasis or correcting homeostasis dysregulations.

A major problem is that the experimental metabolic profile is influenced not only by the genotype but also by age, lifestyle, environmental factors, nutritional status, assumption of drugs (7, 8), and other metabolites from symbiotic organisms (i.e., the gut microflora) (9–11). Consequently, changes in the metabolic profile of biologically complex organisms (like humans) in

response to pathological stimuli may be difficult to distinguish from normal physiological variations.

To overcome the problem of variability in studies with humans, attempts were made to minimize the variations by means, for example, of standardized diet, avoiding any vigorous activity and excluding smokers (12–14). Other studies examined the influence of perturbing factors on the metabolic profile of animals and humans such as diet, specific food, aging, and multiple intrinsic and extrinsic physiological parameters (15–19).

In this perspective, it is crucial to assess whether it is possible to eliminate the noise due to random daily variations and obtain a “natural,” stable, and invariant metabolic profile that is typical of a given subject, even if not necessarily unique. This includes the need to separate within-group from between-group variations (each group of samples representing a single person), a fundamental objective of many metabolomics studies (see refs. 20–22 and references therein). The possibility that such an “individual metabolic phenotype” could exist is suggested by studies reporting differentiation between individuals on the basis of a subject-specific response to particular stimuli [such as kinds of diet (4, 23) or food (18, 19), drug treatment (24), or vitamin intake (25)].

High-resolution NMR spectroscopy is a technique of choice for the investigation of the metabolome (26–28) and has been shown to provide a wealth of metabolic information that can be related to physiological states or pathological conditions of an organism (29). The sample under investigation typically contains many metabolites, and simple one-dimensional ¹H NMR spectra provide several envelopes of not fully resolved signals (30, 31). One approach is to look only for selected metabolites (1, 31, 32) for diagnostic or prognostic purposes. This approach is similar to the classic clinical analyses. Another approach is to try and identify as many metabolites as possible (1, 31, 32), but this is a difficult task and presently not enough rewarding. A general approach is that of dividing the one-dimensional ¹H NMR spectrum in slices for statistical analysis (buckets), and to look at the overall NMR fingerprint (a reflection of the whole detectable metabolome in that particular biofluid) to characterize a single sample.

To assess the existence of an individual metabotype, in this study a panel of healthy subjects (with no evident chronic

Author contributions: M.A., I.B., D.C., C.L., and M.S. designed research; M.A., D.C., H.S., and B.S. performed research; H.S., B.S., and M.S. contributed new reagents/analytic tools; M.A., D.C., H.S., and B.S. analyzed data; and I.B., D.C., C.L., H.S., and M.S. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

[†]Present address: Scientific and Technological Department, University of Verona, 37100 Verona, Italy.

[¶]To whom correspondence should be sent at the † address. E-mail: bertini@cerm.unifi.it.

This article contains supporting information online at www.pnas.org/cgi/content/full/0705685105/DC1.

© 2008 by The National Academy of Sciences of the USA

pathologies) was chosen, and a multiple sample collection of urine samples from each donor was performed. Metabolic fingerprints were generated by one-dimensional ^1H NMR spectroscopy. Innovative multivariate models (based on sequences of well established methodologies) were developed with the aim of maximizing intersubject differences and minimizing intrasubject variability to establish whether an individual metabolic phenotype fingerprint could be found. We demonstrate the existence of such fingerprint and show that it constitutes such a strong characteristic of each donor as to allow its identification with 100% probability. By a projection/back-projection approach it was also possible to obtain this “core” profile free from random daily noise factors, opening new perspectives in metabonomic studies. It is important to highlight that this “individual metabolic phenotype” is not a consequence of a particular stimulus but is an intrinsic general characteristic of a subject.

Results and Discussion

In this study, we collected ≈ 40 urine samples of 22 healthy individuals (11 females and 11 males for a total of 873 samples) with ages in the 25–50 range over a period of 3 months. The corresponding one-dimensional ^1H NMR spectra were measured on a spectrometer operating at 600 MHz proton Larmor frequency following standardized procedures [see [supporting information \(SI\) Methods](#)]. The spectra were first normalized to the NMR signal intensity of the CH_3 -group of creatinine to compensate for the large variations in urine metabolite concentrations and then segmented into N consecutively integrated spectral regions (buckets) of fixed width. If not stated differently, the bucketing ranges were set to a 0.5- to 9.5-ppm chemical-shift range, and the bucket width was set to 0.02 ppm. The 4.5- to 6.0-ppm chemical-shift region was left out of the analysis to remove the effects of variations in the suppression of the water resonance and variations in the urea signal caused by partial cross-solvent saturation through solvent-exchanging protons.

The raw dataset was then analyzed by applying standard statistical methods or combinations thereof, as detailed in *Statistical Methods* and *SI Methods*. On the descriptive level, principal component analysis (PCA), hierarchical cluster analysis (HCA), and canonical analysis (CA) were used, where multivariate analysis of variance (MANOVA) was used to determine the dimensionality of the relevant CA subspaces. On the predictive level, classification rules were derived on the basis of soft independent modeling of class analogy (SIMCA) and K-nearest neighbor (K-NN) methods. Several classification approaches for multiclass problems were used combining statistical techniques, such as a multi-model SIMCA classification (MM-SIMCA), a PCA subspace K-NN classification (PCA/K-NN), and a PCA/CA subspace K-NN classification (PCA/CA/K-NN). The identification of an individual person was performed either from a single new sample/spectrum (single vote) or from a set of new samples/spectra according to the majority vote—i.e., majority rule classification (MRC). The predictivity of the various classification methods was assessed by using test-set validation (TSV) combined with a Monte Carlo (MC) approach.

The first objective of the analysis was to ascertain whether the spectral line-patterns of the individual samples carry unique features that are donor-specific. For this purpose, the individual objects of the data matrix were labeled according to donor identity. For dimension reduction, the data were projected into a PCA subspace explaining 99.9% of the variance in the data. The respective score submatrix was used as input for MANOVA to obtain the dimensionality of the multivariate group means, which resulted to be 21 (with all p values $< 10^{-8}$). The score matrix was projected into the respective 21-dimensional subspace with maximum group discrimination provided by CA. The results are illustrated in Fig. 1, where a convex hull for each

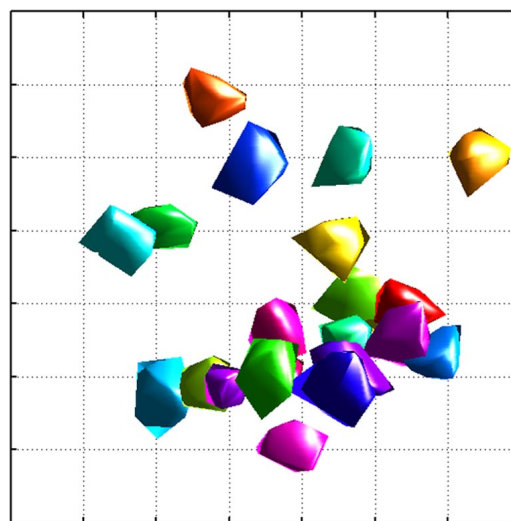


Fig. 1. Projection of the one-dimensional ^1H NMR spectral buckets into PCA/CA subspace in the three most significant dimensions. Convex hulls of the donor-specific point clusters (37–41 points each) are shown for better visualization. Each individual donor has his/her own color code.

donor, rather than the 37–41 individual data points enclosed, is reported in the space of the first three canonical variables.

It is immediately clear (Fig. 1) that a very strong individual characterization is possible, allowing even discrimination with respect to individual donors [even if not looked for, some “natural” gender-clustering is also obtained (see *SI Fig. 6*)]. It should be noted that no separation according to donors can be found when looking at simple PCA scores plots (*SI Fig. 7*).

Although a limited overlap within a few pairs of clusters is present in Fig. 1, it should be kept in mind that the discriminating subspace is 21-dimensional according to the MANOVA output, so that a 3-dimensional plot is not sufficient to provide the complete picture. This becomes apparent from the results of HCA applied to the canonical variables, where the complete 21-dimensional CA subspace is analyzed. The dendrogram in Fig. 2 illustrates how, with the only exception of two spectra, all of the remaining spectra are clearly clustered according to donor identity. It is striking to notice how intergroup distances are by far larger than intragroup distances.

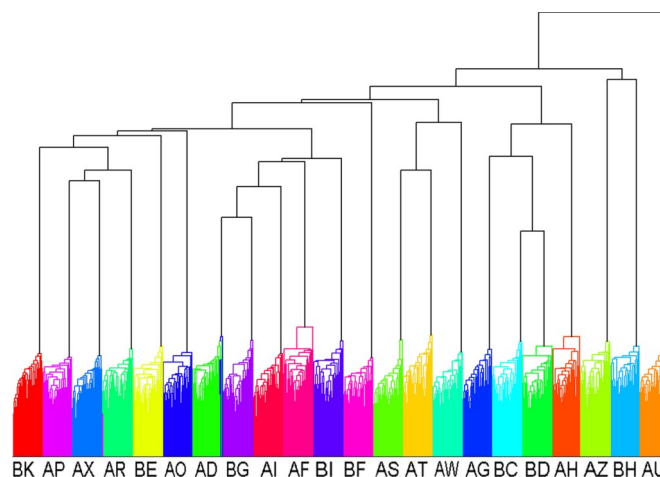


Fig. 2. Dendrogram relative to cluster analysis on the 21-dimensional PCA/CA subspace. Clustering according to donor is obvious.

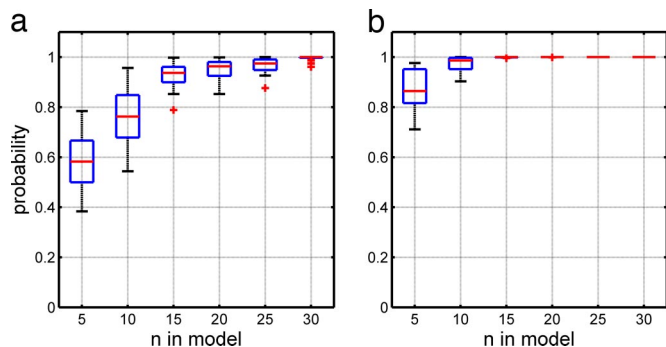


Fig. 4. PCA/CA/K-NN learning curves: Box plots reporting the probability of correct classification as a function of number of spectra in model set. (a) Single-spectrum classification. (b) Classification using seven spectra and MRC. Probabilities were determined by MC/TSV using 1,000 individual runs for averaging.

of the donor that is consistent over time—at least over the collection time (2–3 months) used. Such donor-specific traits can in principle be reconstructed from an existing sample set coming from a fixed set of donors. To illustrate this concept, the original dataset (bucket table) and a reconstructed bucket table are shown in Fig. 5 *a* and *b*, respectively. For reconstruction of each bucket spectrum in Fig. 5*b*, its representation in the 21-

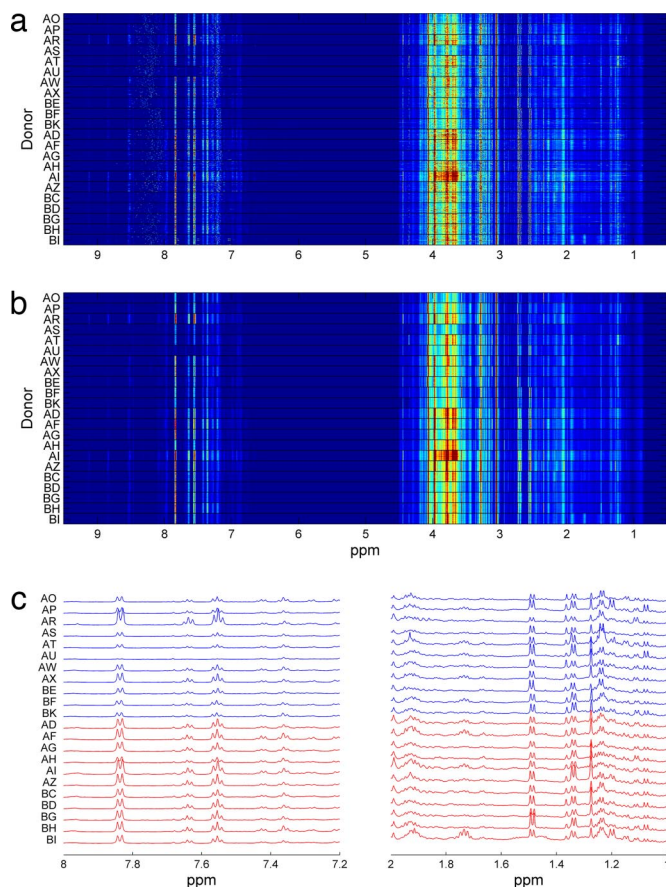


Fig. 5. Reconstruction of the one-dimensional ^1H NMR spectra by back-projection into the original space of their representations in the 21-dimensional discriminating subspace. (a) Image of the original bucket table. (b) Image of the table reconstructed from the bucket spectra representation in the 21-dimensional donor-discriminating PCA/CA subspace. (c) Median reconstructed spectra for individual donors (red, female; blue, male).

dimensional discriminating subspace described above (the dimension that was obtained, as already indicated, by MANOVA output) was back-projected into the original space. Comparison between original and reconstructed datasets confirms that the procedure had just a “cleaning” effect without creating any new artificial features.

It is even possible to obtain, by applying a similar procedure (see *SI Methods*), median reconstructed spectra for individual donors (Fig. 5*c*). Such spectra may be viewed as representing an image of the individual metabolic fingerprint (somewhat related to “metabotype”) and by which it is possible to assess easily, even from visual inspection, donor-specific spectral differences. For example, it is possible to note that subject AU is characterized by persistently low intensity of the doublet at ≈ 7.36 ppm and of the two triplets at ≈ 7.45 and 7.74 ppm (i.e., very low levels of hippuric acid despite the absence of antibiotic treatment), whereas subject BI presents higher intensity of the multiplet at ≈ 1.73 ppm (higher lysine levels) than the other individuals. A graphical representation of the variability of each spectral point and of its discrimination ability is given in *SI Fig. 9*. However, identification of just one particular metabolite for discrimination of a single person seems to be the absolute exception (see *SI Fig. 10*). When analyzing the concentration levels of the 12 selected metabolites, apart from donor AU, no single metabolite can be used for unique identification of any other donor in the panel. Instead, the unique fingerprints are coded in the levels of multiple markers. For example, donor AZ is unique because of a combination of high levels in seven metabolites—i.e., 3-hydroxyisovalerate, citrate, glycine, isoleucine, lactate/threonine, leucine, and valine. In another example, median concentrations of donor AO with respect to creatinine, 3-hydroxyisovalerate, alanine, citrate, dimethylamine, glycine, isoleucine, lactate/threonine, leucine, and valine are consistently at the lower end within the panel. In both examples, the combined concentration values are quite unique, leading to a high single-sample correct classification rate even on the basis of the 12 selected metabolites (AZ, 90.4%; AO, 87.8%), and in the MRC approach AZ and AO can be correctly recognized in practically 100% of the cases. Unfortunately, available metadata are not sufficient for any conclusion on the origin of this particular metabolic finding.

Conclusions

So far, metabonomic studies have concentrated on metabolic “snapshots” of individual samples, which were ostensibly influenced by uncontrolled variables. The possibility, demonstrated here, of eliminating the daily noise by multiple sample collection hints at a new perspective in metabonomic studies, based on the definition and identification of an individual metabolic fingerprint constituted by the invariant part of multiple samples of a single subject.

Although the genotype surely is a unique characteristic of each and every living organism, the existence of a truly individual metabolic phenotype still needs to be assessed. We report experimental evidence of the possibility of recognizing with 100% probability a subject within a group of individuals (none of which being subjected to particular imposed treatment or diet) by the metabolic fingerprint of one of its biofluids. This fingerprint is linked to the genotype but probably also to the general lifestyle and persisting environmental factors. It seems, however, to be independent from random daily variation. NMR confirms itself as a valuable tool for obtaining such fingerprints.

It is interesting to consider that a subject could be already characterized by only using the data from 12 metabolites, and identified with a confidence not much lower than in the case of the analysis performed on the complete fingerprint, provided that MRC is used. If one considers that in biofluids there are hundreds (cerebrospinal fluid) to thousands (urine) of metab-

olites, there is a strong possibility of being able to distinguish individual fingerprints even in a much larger panel study. From the investigation of the 12 metabolites, we conclude that a significant part of the individual metabolic fingerprint seems to be “coded” in the concentration levels and ratios of the metabolites rather than “coded” just by their pure presence or absence. It is the strength of NMR to provide access to multiparametric metabolic fingerprints in a fast, untargeted, and highly reproducible manner providing precise quantitative information even on the smallest concentration changes of multiple metabolites at the same time, despite a large dynamic range covering several orders of magnitude.

Characterizing individual metabolic fingerprints may allow researchers to (i) better plan personalized therapy and nutrition; (ii) perform studies of pharmacometabonomics to better predict and assess drug efficacy and toxicity; (iii) follow phenotype changes as a function of disease progression, possibly leading to earlier diagnosis and prognosis; (iv) perform cost-effective screenings on large human populations; and (v) address how possible long-term changes may be related to aging, because the metabolic fingerprint found is consistent over at least the collection period ($\approx 2\text{--}3$ months).

Statistical Methods

Data Analysis. The software package MATLAB (MathWorks; Version 5.3.1 R11.1) was used for data preparation, data preprocessing, and statistical analysis, using standard MATLAB and routines developed in-house.

Multivariate Statistical Analysis. In the current study, several classification approaches for multiclass problems were used by combining, in an innovative way, classic methods (see *SI Methods*) used for data analysis.

PCA subspace K-NN classification (PCA/K-NN). PCA was applied on a set of model data for dimension reduction, first defining a relevant PCA subspace. New test-set data were classified by applying the K-NN classification based on measured distances between representations of test-set objects and model set objects in the PCA subspace defined by the respective previous PCA on the model set.

PCA/CA subspace K-NN classification (PCA/CA/K-NN). PCA on the model data were initially applied as in the PCA/K-NN approach. MANOVA and CA were then applied to the model set representations in the relevant PCA subspace to define a further reduced subspace with optimum group separation. Its dimensionality was chosen according to the MANOVA output of the dimen-

sionality of the respective group means. In practice, in the PCA/CA/K-NN, two subsequent coordinate transformations and dimension reductions were performed to identify a subspace with optimum multigroup discrimination. The final classification procedure of new test-set objects was similar to the PCA/K-NN approach. They were first projected in the discriminating subspace defined by the model set and then the K-NN classification was applied.

Majority rule classification (MRC). It was introduced for the current study such that a predefined group of test samples was classified according to the relative majority of the single-sample classification results. The reasoning is that an improved donor identification from spectral metabolic fingerprints of body fluids is expected if several samples per donor are used for testing rather than performing the identification on the basis of a single sample. In the context of this study, 7 spectra were used for MRC. This choice was made to account for the fact that 30 spectra per donor were needed to establish the respective models, so that the maximum number available for all donors for testing the MRC performance was 7 (considering that the samples available for the AH subject were 37).

Test-set validation (TSV). The dataset was subdivided into a training set and a test set. Class membership of training- and test-set samples were known in advance. The respective classification rule and model under evaluation was designed from the training-set data exclusively and applied to the test-set samples that were left out during the modeling process. The classification result for the test-set samples was compared with the known true class, and the fractions of correct and false classifications were recorded. If not stated differently, in each TSV approach, 30 spectra were used as training set per donor as well as the remaining spectra in the test set.

MC embedded TSV (MC/TSV). This validation was performed as a generalization of a cross-validation procedure, looping an MC test scenario over multiple TSVs (e.g., 1,000 times). In each MC iteration step, the dataset was divided into a training set and a test set, randomly assigning a predefined number of objects from the dataset to either of both subsets. Then, models and classification rules were derived on the basis of the respective training set and applied to the current test set. For evaluation of the classification performance, several parameters were recorded—i.e., the group into which each test-set sample was classified (possibly multiple groups in the case of MM-SIMCA), the fraction of test-set samples classified correctly, and the MRC result on all combined samples from the respective test set. Once the classification was finished, the next MC step was initiated by a new random split of the data into a training set and a test set, and so on. Classification results were recorded for each MC step.

ACKNOWLEDGMENTS. We thank Alessandro Quattrone for discussions that lead to this project's conception and advice during the early stages of sample collection. This work was supported by grants from Ente Cassa di Risparmio di Firenze to the FiorGen Foundation and by a fellowship from Boehringer Ingelheim Italia (to D.C. through the FiorGen Foundation).

- Fiehn O (2001) *Comp Funct Genomics* 2:155–168.
- Nicholson JK, Lindon JC, Holmes E (1999) *Xenobiotica* 29:1181–1189.
- Gavaghan CL, Holmes E, Lenz E, Wilson ID, Nicholson JK (2000) *FEBS Lett* 484:169–174.
- Rezzi S, Ramadan Z, Fay LB, Kochhar S (2007) *J Proteome Res* 6:513–525.
- Kussmann M, Raymond F, Affolter M (2006) *J Biotechnol* 124:758–787.
- Clayton TA, Lindon JC, Cloarec O, Antti H, Charuel C, Hanton G, Provost JP, Le Net JL, Baker D, Walley RJ, et al. (2006) *Nature* 440:1073–1077.
- Rang HP, Dale MM, Ritter JM (1995) *Pharmacology* (Churchill Livingstone, Edinburgh).
- Tannock GW (1995) *Normal Microflora: An Introduction to Microbes Inhabiting the Human Body* (Chapman & Hall, London).
- Nicholson JK, Wilson ID (2003) *Nat Rev Drug Discovery* 2:668–676.
- Nicholson JK, Holmes E, Lindon JC, Wilson ID (2004) *Nat Biotechnol* 22:1268–1274.
- Tiret L (2002) *Proc Nutr Soc* 61:457–463.
- Lenz EM, Bright J, Wilson ID, Morgan SR, Nash AFP (2003) *J Pharmacol Biomed Anal* 33:1103–1115.
- Lenz EM, Bright J, Wilson ID, Hughes A, Morrisson J, Lindberg H, Lockton A (2004) *J Pharmacol Biomed Anal* 36:841–849.
- Zuppi C, Messana I, Forni F, Ferrari F, Cristina R, Giardina B (1998) *Clin Chim Acta* 278:75–79.
- Kochhar S, Jacobs DM, Ramadan Z, Berruex F, Fuerholz A, Fay LB (2006) *Anal Biochem* 352:274–281.
- Walsh MC, Brennan L, Malthouse JP, Roche HM, Gibney MJ (2006) *Am J Clin Nutr* 84:531–539.
- Bollard ME, Stanley EG, Lindon JC, Nicholson JK, Holmes E (2005) *NMR Biomed* 18:143–162.
- Solanky KS, Bailey NJ, Beckwith-Hall BM, Bingham S, Davis A, Holmes E, Nicholson JK, Cassidy A (2005) *J Nutr Biochem* 16:236–244.
- Solanky KS, Bailey NJ, Beckwith-Hall BM, Davis A, Bingham S, Holmes E, Nicholson JK, Cassidy A (2003) *Anal Biochem* 323:197–204.
- Jansen JJ, Hoefsloot HJ, van der Greef J, Timmerman ME, Smilde AK (2005) *Anal Chim Acta* 530:173–183.
- Jansen JJ, Hoefsloot HJ, Boelens HFM, van der Greef J, Smilde AK (2004) *Bioinformatics* 20:2438–2446.
- Smilde AK, Jansen JJ, Hoefsloot HJ, Lamers RJ, van der Greef J, Timmerman ME (2005) *Bioinformatics* 21:3043–3048.
- Stella C, Beckwith-Hall B, Cloarec O, Holmes E, Lindon JC, Powell J, Van der Ouderaa F, Bingham S, Cross AJ, Nicholson JK (2006) *J Proteome Res* 5:2780–2788.
- Idle JR, Mahgoub A, Lancaster R, Smith RL (1978) *Life Sci* 22:979–983.
- Kelly FJ, Lee R, Mudway IS (2004) *Ann NY Acad Sci* 1031:22–39.
- Nicholson JK, Connolly J, Lindon JC, Holmes E (2002) *Nat Rev Drug Discovery* 1:153–161.
- Lindon JC, Holmes E, Nicholson JK (2001) *Prog Nucl Magn Reson Spectrosc* 39:1–40.
- Griffin JL (2003) *Curr Opin Chem Biol* 7:648–654.
- Nicholson JK, Wilson ID (1989) *Prog Nucl Magn Reson Spectrosc* 21:449–501.
- Fiehn O (2002) *Plant Mol Biol* 48:155–171.
- Goodacre R, Vaidyanathan S, Dunn WB, Harrigan GG, Kell DB (2004) *Trends Biotechnol* 22:245–252.
- Claudino WM, Quattrone A, Biganzoli L, Pestrin M, Bertini I, Di Leo A (2007) *J Clin Oncol* 25:2840–2846.