

Mining and Visualizing Large Anticancer Drug Discovery Databases[†]

Leming M. Shi,^{*,‡,⊥} Yi Fan,^{‡,||} Jae K. Lee,[‡] Mark Waltham,[‡] Darren T. Andrews,[‡] Uwe Scherf,[‡]
Kenneth D. Paull,^{§,⊗} and John N. Weinstein^{*,‡,▽}

Laboratory of Molecular Pharmacology, Division of Basic Sciences, and Information Technology Branch,
Division of Cancer Treatment, Diagnosis, and Centers, National Cancer Institute, National Institutes of Health,
Bethesda, Maryland 20892-4255

Received July 30, 1999

In order to find more effective anticancer drugs, the U.S. National Cancer Institute (NCI) screens a large number of compounds in vitro against 60 human cancer cell lines from different organs of origin. About 70 000 compounds have been tested in the program since 1990, and each tested compound can be characterized by a vector (i.e., “fingerprint”) of 60 anticancer activity, or $-\log(GI_{50})$, values. GI_{50} is the concentration required to inhibit cell growth by 50% compared with untreated controls. Although cell growth inhibitory activity for a *single* cell line is not very informative, activity *patterns* across the 60 cell lines can provide incisive information on the mechanisms of action of screened compounds and also on molecular targets and modulators of activity within the cancer cells. Various statistical and artificial intelligence methods, including principal component analysis, hierarchical cluster analysis, stepwise linear regression, multidimensional scaling, neural network modeling, and genetic function approximation, among others, can be used to analyze this large activity database. Mining the database can provide useful information: (a) for the development of anticancer drugs; (b) for a better understanding of the molecular pharmacology of cancer; and (c) for improvement of the drug discovery process.

INTRODUCTION

Arrays of biological assays are becoming increasingly important to pharmaceutical research and development.¹ Information encoded by the output of such screening systems can be used not only to accelerate the processes of drug discovery and development but also to gain essential insights into the differences in pharmacological activity among drug candidates directed against different biological targets. In-depth analysis of the experimental data using various statistical and artificial intelligence (AI) techniques is crucially important to the success of such drug discovery programs (a list of abbreviations used in this paper is given in Table 1). A few examples of such analyses follow:

Ebert et al. used principal component analysis (PCA) to investigate the relative information contents of 15 routine antitumor tests on 13 para-substituted aryl dimethyl triazenes in order to show whether these compounds exhibit selective antimetastatic activity.² A two-component PCA model was able to explain 71% of the variance of the 13×15 activity matrix. Nendza and Seydel used multiple linear regression (MLR) and PCA to analyze ecotoxicity data obtained in 11 biological test systems (3 in bacteria, 4 in yeast, 2 in algae, 1 in protoplasts, and 1 in daphnia) for more than 50 phenols,

Table 1. Abbreviations Used in This Paper

2D:	two-dimensional
3D:	three-dimensional
AI:	artificial intelligence
ANN:	artificial neural network
CNS:	central nervous system
CV:	cross validation
DIS:	drug information system
DMSO:	dimethyl sulfoxide
GA:	genetic algorithm
GFA:	genetic function approximation
GI_{50} :	concentration required to inhibit cell growth by 50%
HCA:	hierarchical cluster analysis
HIV:	human immunodeficiency virus
HTS:	high throughput screening
MARS:	multivariate adaptive regression spline
MDS:	multidimensional scaling
MLR:	multiple linear regression
NCI:	National Cancer Institute
NTP:	National Toxicity Program
PC(s):	principal component(s)
PCA:	principal component analysis
PCR:	principal component regression
QSAR:	quantitative structure–activity relationship
SRB:	sulforhodamine B

anilines, and hydrocarbons.³ MLR showed that the relative toxicities of these compounds were reasonably well-determined in all test systems and could be described as a function of lipophilicity ($\log P$): the toxicity increased with increasing lipophilicity. This finding was corroborated and extended by PCA: The first principal component (PC), which explained 76% of the variance of the activity matrix, showed a significant correlation with toxicity. The authors concluded that other PCs had accounted mainly for the experimental errors and that assessment of the toxicity of these chemicals by only one of these assays was sufficient. By using PCA,

[†] In memory of Dr. Kenneth D. Paull.

[‡] Laboratory of Molecular Pharmacology, Division of Basic Sciences.

[§] Information Technology Branch, Division of Cancer Treatment, Diagnosis, and Centers.

[⊥] Current address: American Cyanamid Company, Quakerbridge and Clarksville Rd, P.O. Box 400, Princeton, NJ 08543. E-mail: shil@pt.cyanamid.com.

^{||} Current address: Wyeth-Ayerst Research, CN8000, 865 Ridge Rd., Monmouth Junction, NJ 08852. E-mail: fank@war.wyeth.com.

[⊗] Deceased Jan 29, 1998.

[▽] Telephone: (301) 496-9571. E-mail: weinstein@dtapx2.ncifcrf.gov.

Moret and Janssen analyzed the cytostatic activities of 38 benzoquinones tested in various *in vitro* and *in vivo* systems.⁴ The correlation among *in vitro* and *in vivo* tests was poor because of the large number of biological mechanisms. The first and second PCs, which explained about 50% and 33% of the variance, respectively, were significant. Because the authors put structural descriptors (σ and molecular reflectance) and activity data together when performing PCA, the first PC may have explained principally the variance in activity, whereas the second PC may have explained principally the variance in structural descriptors. Benigni analyzed the relationships between *in vitro* mutagenicity assays reported by the U.S. National Toxicology Program (NTP).⁵ Franke et al. used factor analysis and correlation techniques to analyze antibacterial and pharmacokinetic data from parallel biological measurements.⁶ Ojama et al. recently analyzed the binding of 187 steroids to 5 steroid hormone receptors (estrogen receptor, progesterone receptor, androgen receptor, mineralocorticoid receptor, and glucocorticoid receptor) using correspondence factor analysis.^{7,8}

More recently, Telik, Inc. (South San Francisco, CA), developed a method for predicting ligand binding to proteins by affinity fingerprinting from a small panel of reference proteins.^{9–11} After preliminary testing of over 300 proteins from a variety of sources, the authors selected panels of 8–18 proteins that displayed the broadest binding affinities for a set of over 5000 compounds. MLR was used to build “computational surrogates” to predict the binding potencies of additional molecules. The affinity fingerprint database, which provides a rich source of data defining “operational similarities” among the proteins, can be used for efficient prescreening of a large number of compounds against target proteins in order to select promising candidates for further study.

The past few years have seen an explosion of interest in combinatorial chemistry and high throughput screening (HTS).¹² The capacity for both synthesis and screening of compounds has increased to levels unthinkable just a few years ago. The consequence is that pharmaceutical companies are able to screen a much larger number of drug candidates against a much larger number of biological end points or targets. A major pharmaceutical company may have an inventory of a million or more compounds for screening. With combinatorial chemistry, the number of compounds available for screening is becoming larger and larger. Some companies are screening compounds against hundreds of biological targets, including various enzymes and cell types. The pharmaceutical industry is faced with the challenge of managing and analyzing the huge amount of high-dimensional data being generated through combinatorial chemistry and HTS. The analytical methods discussed here can be useful for that purpose.

NCI ANTICANCER DRUG DISCOVERY PROGRAM

In order to find new anticancer drugs, the NCI has conducted extensive testing for possible activity against different forms of cancer. Tests have been done on natural product extracts, synthetic compounds, biotechnology products, and, more recently, combinatorial chemistry libraries. Before 1985, the NCI used mice bearing murine leukemia P388 cells to screen compounds. That strategy identified

agents active against leukemias, but relatively few of these compounds were effective against solid tumors, including the most common human carcinomas. In 1990, the NCI established a primary screen in which compounds are tested *in vitro* for their ability to inhibit growth of 60 human cancer cell lines from different organs.^{13–17} Currently included in the screen, as shown in Table 2, are melanomas, leukemias, and cancer cells of breast, prostate, lung, colon, ovary, kidney, and central nervous system (CNS) origin.

The purpose of this screen is to provide an initial evaluation of compounds for cytotoxic or growth inhibitory activity against a diverse panel of tumor types. This strategy was based on the hypothesis that selective *in vitro* activity against cancer cell lines from a particular organ would predict selective activity against corresponding cancers in humans. Compounds that show interesting activity patterns in the *in vitro* screen are selected for subsequent *in vitro* and *in vivo* testing. This concept is still being tested as agents progress through clinical trials. However, as shown in Figure 1, activity patterns observed in the screen, together with information on chemical structure descriptors and molecular targets in the 60 cell lines, have provided important sources of information for testing hypotheses about the molecular pharmacology of cancer and structure–activity relationships in drug discovery.¹⁸

In the past few years, we have been developing an “information-intensive” approach to the discovery of anticancer drugs and to the molecular pharmacology of cancer.^{18–24} Our approach is based on the three databases, or matrices, shown in Figure 1: **A**, anticancer activity patterns across the 60 human cancer cell lines for screened compounds; **S**, chemical structure information for an inventory of about 500 000 compounds, including the compounds screened; and **T**, information on possible targets or modulators of activity in the 60 cell lines.

To date, more than 70 000 chemical compounds, plus a larger number of natural product extracts, have been tested. As our focus is on the chemical compounds, we will be considering an **A** database of dimension 70 000 \times 60. Each entry in database **A** represents the anticancer activity of a particular compound against a particular cancer cell line. Each row of **A** represents the activity pattern, or fingerprint, of a compound. Similarity of compounds can be assessed in terms of the similarity of their activity patterns. The growth inhibitory activity for a *single* cell line is not very informative, but we have found that activity *patterns* across the 60 lines provide incisive information on the mechanisms of action of screened compounds and also on molecular targets and modulators of activity.^{16–24} Similarity in activity patterns very often indicates a similarity in the mechanism of action, mode of drug resistance, and molecular structure of tested compounds.^{16,20} Similarly, each column of **A** can be regarded as the sensitivity profile of a particular cancer cell line to the large number of screened chemicals. Therefore, we can assess the similarity of cancer cell lines in terms of their responses to thousands of chemicals, as will be discussed later.

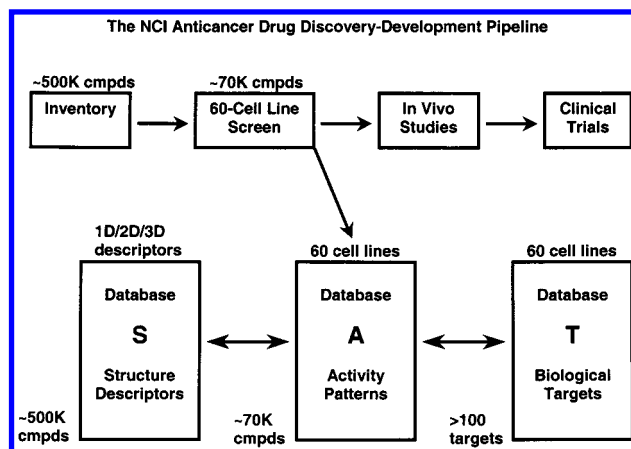
The second database is the chemical structure database, **S**. Each compound can be characterized by various chemical descriptors. The descriptors can be one-dimensional (e.g., log *P*, molecular weight, water solubility). They can also be two-dimensional, as in the case of the widely used MDL

Table 2. Sixty Human Cancer Cell Lines Used in the NCI Anticancer Drug Discovery Program

no.	cell type	cell name	doubling time, ^a h
1	leukemia	CCRF-CEM	26.7
2	leukemia	HL-60(TB)	28.8
3	leukemia	K-562	19.6
4	leukemia	MOLT-4	27.9
5	leukemia	RPMI-8226	33.5
6	leukemia	SR	28.7
7	nsclc ^b	A549/ATCC	22.9
8	nsclc	EKVX	43.6
9	nsclc	HOP-62	39.0
10	nsclc	HOP-92	79.5
11	nsclc	NCI-H226	61.0
12	nsclc	NCI-H23	33.4
13	nsclc	NCI-H322M	35.3
14	nsclc	NCI-H460	17.8
15	nsclc	NCI-H522	38.2
16	colon	COLO-205	23.8
17	colon	HCC-2998	31.5
18	colon	HCT-116	17.4
19	colon	HCT-15	20.6
20	colon	HT29	19.5
21	colon	KM12	23.7
22	colon	SW-620	20.4
23	cns ^c	SF-268	33.1
24	cns	SF-295	29.5
25	cns	SF-539	35.4
26	cns	SNB-19	34.6
27	cns	SNB-75	62.8
28	cns	U251	23.8
29	melanoma	LOX-IMVI	20.5
30	melanoma	MALME-3M	46.2
31	melanoma	M14	26.3
32	melanoma	SK-MEL-2	45.5
33	melanoma	SK-MEL-28	35.1
34	melanoma	SK-MEL-5	25.2
35	melanoma	UACC-257	38.5
36	melanoma	UACC-62	31.3
37	ovary	IGROV1	31.0
38	ovary	OVCAR-3	34.7
39	ovary	OVCAR-4	41.4
40	ovary	OVCAR-5	48.8
41	ovary	OVCAR-8	26.1
42	ovary	SK-OV-3	48.7
43	kidney	786-0	22.4
44	kidney	A498	66.8
45	kidney	ACHN	27.5
46	kidney	CAKI- I	39.0
47	kidney	RXF-393	62.9
48	kidney	SN12C	29.5
49	kidney	TK-10	51.3
50	kidney	UO-31	41.7
51	prostate	PC-3	27.1
52	prostate	DU-145	32.3
53	breast	MCF7	25.4
54	uncertain ^d	NCUIADR-RES	34.0
55	breast	MDA-MB-231/ATCC	41.9
56	breast	HS-578T	53.8
57	uncertain ^e	MDA-MB-435	25.8
58	uncertain ^e	MDA-N	22.5
59	breast	BT-549	53.9
60	breast	T-47D	45.5

^a Data from the NCI Developmental Therapeutics Program (<http://dtp.nci.nih.gov/>). ^b nsclc: non-small cell lung cancer. ^c cns: central nervous system. ^d Formerly designated as breast, but its origin is now considered uncertain. ^e Formerly designated as breast, but our analysis shows that they are most like melanomas or at least have the activity and protein expression signatures of melanoma.

MACCS 2D structural keys (MDL Information Systems, Inc., San Leandro, CA; <http://www.mdli.com>). They can also be three-dimensional, as in the case of the pharmacophore descriptors used in Chem-X (Chemical Design, Ltd., Ox-

**Figure 1.** Three types of databases generated by the NCI anticancer drug discovery program. Modified from Weinstein et al., ref 18.

fordshire, U.K.; <http://www.chemdesign.com/>). The NCI Drug Information System (DIS),^{25–29} a major resource for drug discovery, contains structural information for nearly 500 000 compounds, including the 70 000 tested to date. The NCI DIS 2D system offers a useful basis for 2D substructure (or fragment) searching. It was successfully used, for example, to search for agents that are able to reverse multidrug resistance in cancer.³⁰ In that study, the 2D queries (biophores and biophobes) were automatically generated from the CASE³¹ and MULTICASE³² programs. The NCI DIS 3D database²⁹ has been used successfully to identify inhibitors of protein kinase C,³³ HIV-I protease,³⁴ and HIV-1 integrase,^{35–38} through a process of pharmacophore identification and 3D database searching.^{39,40} Molecular structural fingerprints, such as the Daylight 2D hashed descriptor set (Daylight Chemical Information Systems, Inc., Santa Fe, NM; <http://www.daylight.com>), can also be used to represent chemical structural information. Since it is now feasible to create hundreds to thousands of descriptors for a given compound, the S matrix may contain hundreds of millions of numbers.

The third database is the molecular target matrix, T. It contains biological information on the 60 cell lines, including information at the DNA, RNA, protein, and functional levels. More than 100 “targets” (a shorthand term used here for any assayed cell characteristic) have been characterized for the 60 cell lines by different groups at the NCI and elsewhere. The number is increasing rapidly.^{18,41–51} Each value in the T matrix represents the expression level or the status of a particular biological target in a particular cancer cell line. With the application of new “omic”⁵² (i.e., genomic and proteomic) technologies, particularly the gene chips or DNA microarrays,⁵³ the size of the T matrix is becoming much larger. mRNA expression profiles for >8000 genes have been obtained using cDNA microarrays^{54,55} and for 6800 genes using oligonucleotide chips (unpublished data).

Several different algorithms have been introduced in order to use these data for the discovery of anticancer drugs and for study of the molecular pharmacology of cancer.¹⁸ The COMPARE program^{14,16,17,46,56} has proved useful in finding agents with activity patterns similar to that of a “seed” compound or in finding compounds with activity patterns that correlate significantly with the expression levels of a particular biological target across the 60 cell lines. Back-propagation neural networks,²⁰ Kohonen self-organizing

maps,²² principal component analysis (PCA),^{21,23} and cluster analysis^{18,23} have been used to predict mechanisms of action or to organize compounds into families based on their activity patterns. This "information-intensive" approach to the molecular pharmacology of cancer and anticancer drug discovery^{16,18–20,24,50} has proved useful in identifying subgroups of compounds with an apparent relationship to particular biological targets, e.g., the identification of "p53-inverse" agents.^{18,23,57,58}

We previously investigated the relationship between the A and T matrices.¹⁸ In this report, we will focus our analyses on the A matrix. Information in the S matrix has mainly been used for QSAR studies and for 2D substructure or 3D pharmacophore searching. Information in the T matrix that will be presented here relates to p53 and topoisomerase 1.

METHODS

Cell Screening and Activity Data Set. Details of the NCI cell screening protocols and reporting procedures have been described elsewhere.^{14–16,56} Briefly, DMSO (dimethyl sulfoxide) solutions of the compounds are routinely diluted by a factor of 500 with aqueous medium. Aliquots of the resulting aqueous solutions or suspensions are tested against the cancer cells in 96-well microtiter plates for 48 h of exposure, and cell growth is then assayed spectrophotometrically by staining for total cellular protein with sulforhodamine B (SRB). It should be noted that 2-day SRB staining is by no means the only method for obtaining an index of cell growth inhibition or cytotoxicity, but it is a standard one. Like any other, it has limitations and must be interpreted with care.

The growth inhibitory activity of a tested compound is expressed in terms of the quantity $-\log(\text{GI}_{50})$, where GI_{50} is the concentration (generally in moles/liter) required to inhibit cell growth by 50% in comparison with the untreated control. For each compound, 60 activity values (1 for each cell line) make up the activity fingerprint. Shown in Figure 2 are two formats for the anticancer activity pattern of paclitaxel (Taxol of Bristol-Myers Squibb). Figure 2A is drawn with the original activity data, i.e., the $-\log(\text{GI}_{50})$ values. The mean activity for Taxol across the 60 cell lines is 7.57 log units, and the standard deviation of the activity pattern is 0.59 log unit.

Shown in Figure 2B is the pattern of relative activity values presented in the "mean graph" representation introduced by Paull.^{14–17} The relative potency pattern is calculated by subtracting the mean activity from the original activity value for each cell line. An upward bar indicates that the particular cell line is more sensitive than the mean over all of the cell lines; a downward bar indicates that the cell line is less sensitive than the mean.

For most purposes, the mean graph provides a better way to visualize the selective activity of a tested compound. Both the original and relative activity patterns were analyzed in this study. As has been noted previously, the GI_{50} values for a single cell line are not very informative, but activity patterns across the 60 cell lines do encode rich information, both on the cells and on the tested compounds.

The analyses reported in this paper were based on a version of the activity database warehoused at the beginning of 1997.

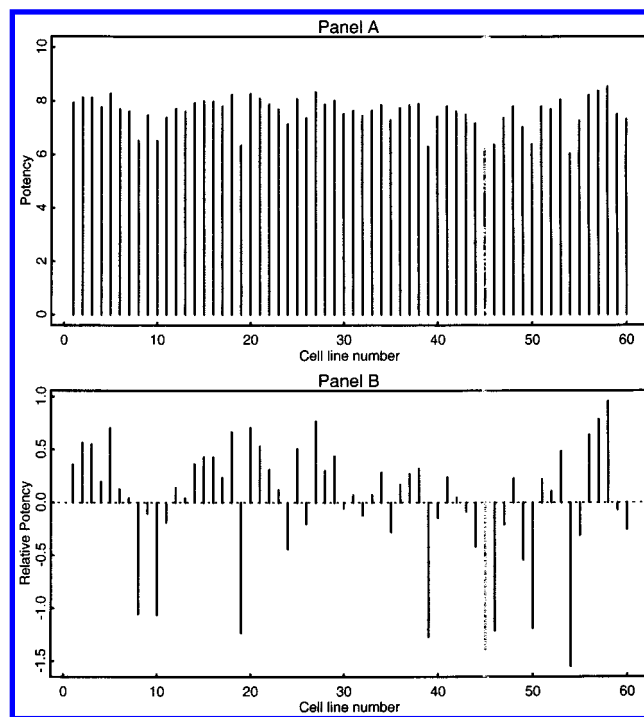


Figure 2. Anticancer activity pattern for paclitaxel (Taxol). (A) Original activity pattern; (B) relative (delta) activity pattern, or "mean graph" representation. After Paull et al., ref 16.

By that time, the activity patterns of 47 702 compounds met the minimal criteria of an activity pattern for the DISCOVERY programs.^{18,50} One of the criteria used was that the standard deviation of the 60 GI_{50} values had to be greater than a preset threshold value. Patterns without significant variation may encode mainly experimental error, not real differences among cell lines. Among the 47 702 compounds, 25 357 were open (i.e., nonconfidential) and 22 345 were discreet (i.e., confidential). On the basis of agreements between the NCI drug discovery program and suppliers of compounds, the NCI is obliged to keep information on the discreet compounds confidential. The molar concentrations of 331 open compounds were unavailable because the molecular weights were not given. These compounds, and the discreet ones, were excluded from further analysis. In the present study, we have therefore focused our analyses on the 25 026 open compounds for which concentrations are given in molar units. This will be referred to here as the "open data set". Overall, this set included 14.6% missing values, each of which was replaced by the mean activity value for the drug in question. A missing value in the database means that a particular compound was not tested against a particular cell line or that tests done failed to meet quality control criteria.

Computational Methods. Most of the statistical analyses were performed by writing appropriate scripts in the S-Plus statistical package (StatSci Division, MathSoft, Inc., Seattle, WA; <http://www.mathsoft.com>). Additional statistical analyses were performed using SAS or JMP (SAS Institute Inc., SAS Campus Drive, Cary, NC; <http://www.sas.com>). Molecular modeling packages Cerius² 3.0 (Molecular Simulation, Inc., San Diego, CA; <http://www.msi.com>) and Sybyl 6.2 (Tripos Associates, Inc., St. Louis, MO; <http://www.tripos.com>) were used both for the calculation of molecular descriptors and for QSAR studies.

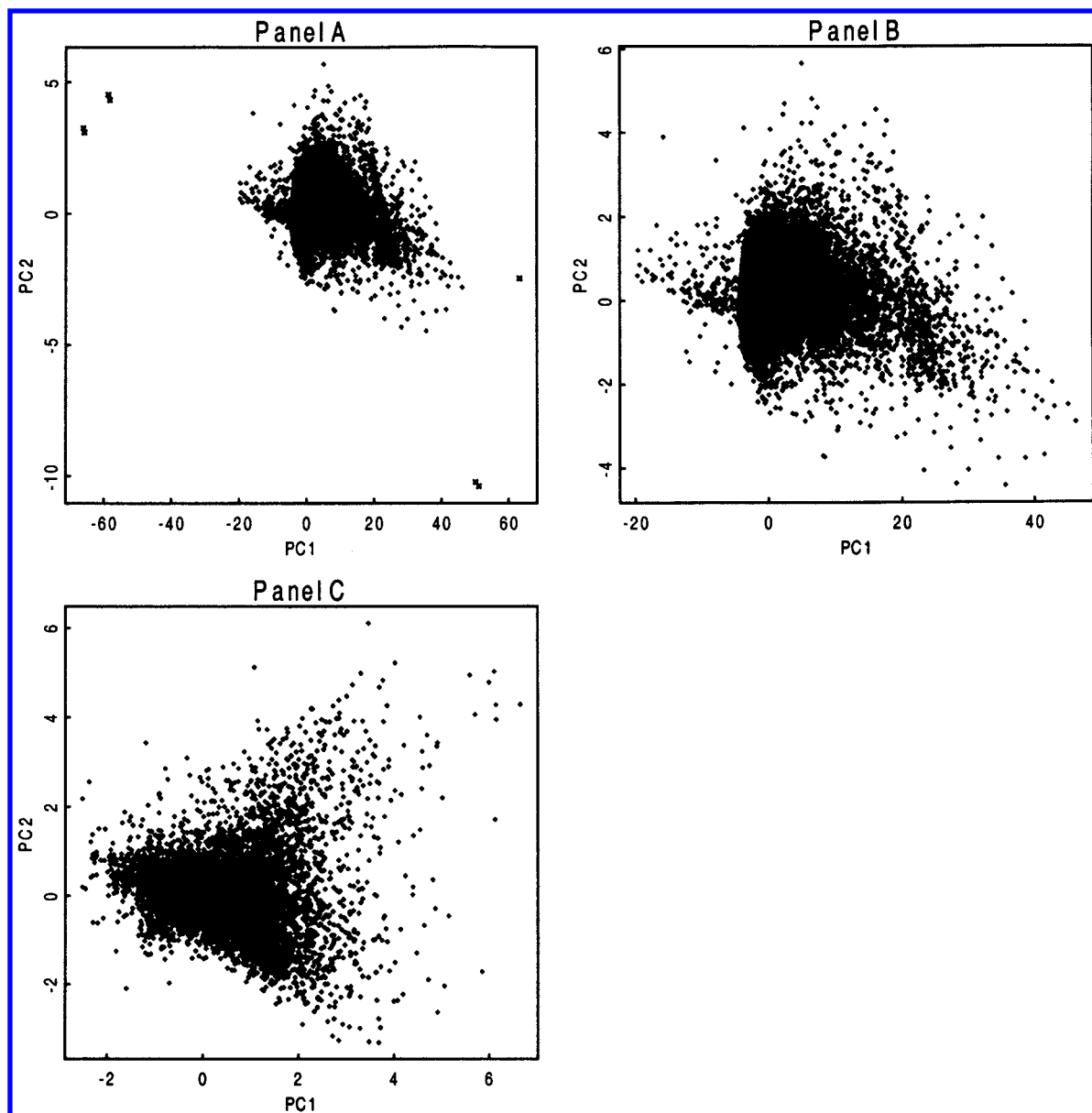


Figure 3. Principal component score plots. (A) PCA score plot based on the database of original activity patterns of 25 026 open compounds, indicating some outlier samples (×). (B) PCA score plot based on a “cleaned” database of original activity patterns of 25 023 open compounds. The first PC explains most of the variance (89.4%). (C) PCA score plot based on a “cleaned” database of relative activity patterns of 25 023 open compounds. The first PC explains only 10.6% of the variance.

RESULTS AND DISCUSSION

Principal Component Analysis (PCA). PCA is a technique that takes linear combinations of the original variables (columns) of a data matrix such that the first PC explains as much of the overall variation as possible, the second PC explains the next most variation subject to being orthogonal to the first, and so on.⁵⁹ PCA is often used as a data reduction technique, sometimes in conjunction with regression, i.e., in principal component regression (PCR). Meglen has given an excellent description of the use of PCA to examine large databases.⁶⁰

Koutsoukos et al.²¹ have used PCA to analyze the NCI cell screen data for a set of 141 “standard” anticancer agents for which the mechanisms of action are well-defined.²⁰ PCA score plots showed distinct clusters of compounds for some of the mechanisms of action examined. Here we extend the analysis to the entire open data set. Figure 3A, a PCA plot based on original activity patterns, i.e., $-\log(\text{GI}_{50})$ values,

shows the PCA results. PCA score plotting is a very good tool for identifying outliers, influential points, and data entry errors.⁶⁰ Figure 3A shows that 7 of the 25 026 compounds appeared at the extremes. By checking the original database, we found that the GI_{50} values of the four compounds (NSC numbers 619989, 620130, 624589, and 626674) located at the upper left corner of Figure 3A were simply miscoded in the database: the exponents of concentrations that should have been -4 or -6 were miscoded as $+4$ or $+6$. These apparent data entry errors for the four compounds were corrected before further data analysis. Two compounds (NSC 676495 and 676497) that are located at the lower right corner and another one (NSC 332598) that is located at the far right side of Figure 3A appeared to have extremely high activities, with $-\log(\text{GI}_{50})$ values up to 13 (i.e., 50% growth inhibition concentration of 10^{-13} mol/L). We were not able to confirm whether there had been miscoding for these three compounds; therefore, they were excluded from further analysis. A PCA

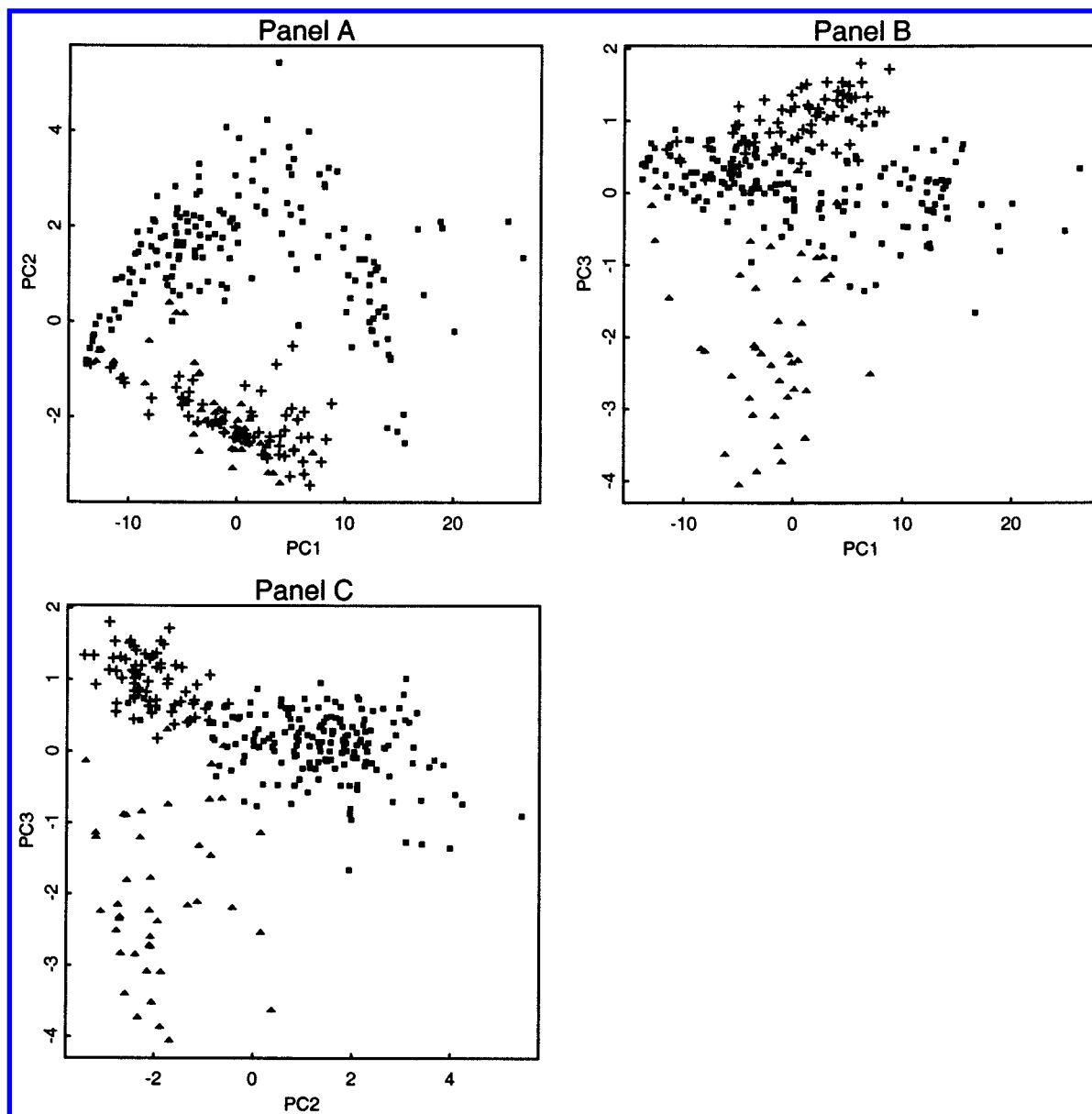


Figure 4. PCA score plots of the ellipticine–ellipticinium–camptothecin data set using original activity patterns. The first PC explains about 90% of the variance but is not informative with respect to mechanism of action, as shown in plots A and B. Plot C with PC2 and PC3 gives the best separation. (+) Ellipticines; (▲) ellipticiniums; (■) camptothecins.

was performed on the remaining “cleaned” open data set of 25 023 compounds, and the score plot for the first 2 PCs is shown in Figure 3B. This cleaned data set can be found at <http://discover.nci.nih.gov>.

The first five PCs accounted for 89.4%, 0.97%, 0.73%, 0.56%, and 0.47% of the variance in activity, respectively. Sixteen PCs were required to account for 95.0% of the activity variance. An immediate conclusion one might draw from these data is that the 60 cell lines (variables) are quite similar to each other such that the various cell lines do not provide much independent information. This suggestion is further supported by the very high Pearson correlation coefficients (mean = 0.894, median = 0.896, maximum = 0.981, minimum = 0.777, SD = 0.034) among the 60 cell lines in terms of their responses to the 25 023 open compounds. The high similarity of the cell lines reflects the fact that the cytotoxic response to a given compound is generally not specific to particular cell types. In fact, the first PC was determined almost entirely by the mean activity,

which is not specific or informative, as can be seen in the following discussion.

Another PCA was performed, this time on the matrix of relative, rather than absolute, activities. In this case, the first five PCs explained 10.6%, 7.0%, 5.3%, 3.8%, and 3.4% of the variance, respectively. Preprocessing of the original activity patterns by subtracting the mean decreased the correlation among cell lines numerically (mean = -0.016 , median = 0.026 , maximum = 0.734 , minimum = -0.307 , SD = 0.124)—i.e., centering the distribution of correlation coefficients at approximately zero.

As stated above, the relationship among compounds can be explored graphically by making PCA score plots using the first two or three PCA scores.^{21,60} Compounds similar in activity pattern should group together on the score plots, whereas dissimilar compounds should scatter over the plots. A particular example is the PCA analysis on a data set of 279 analogues of ellipticine, ellipticinium, and camptothecin. These three types of compounds have shown interesting

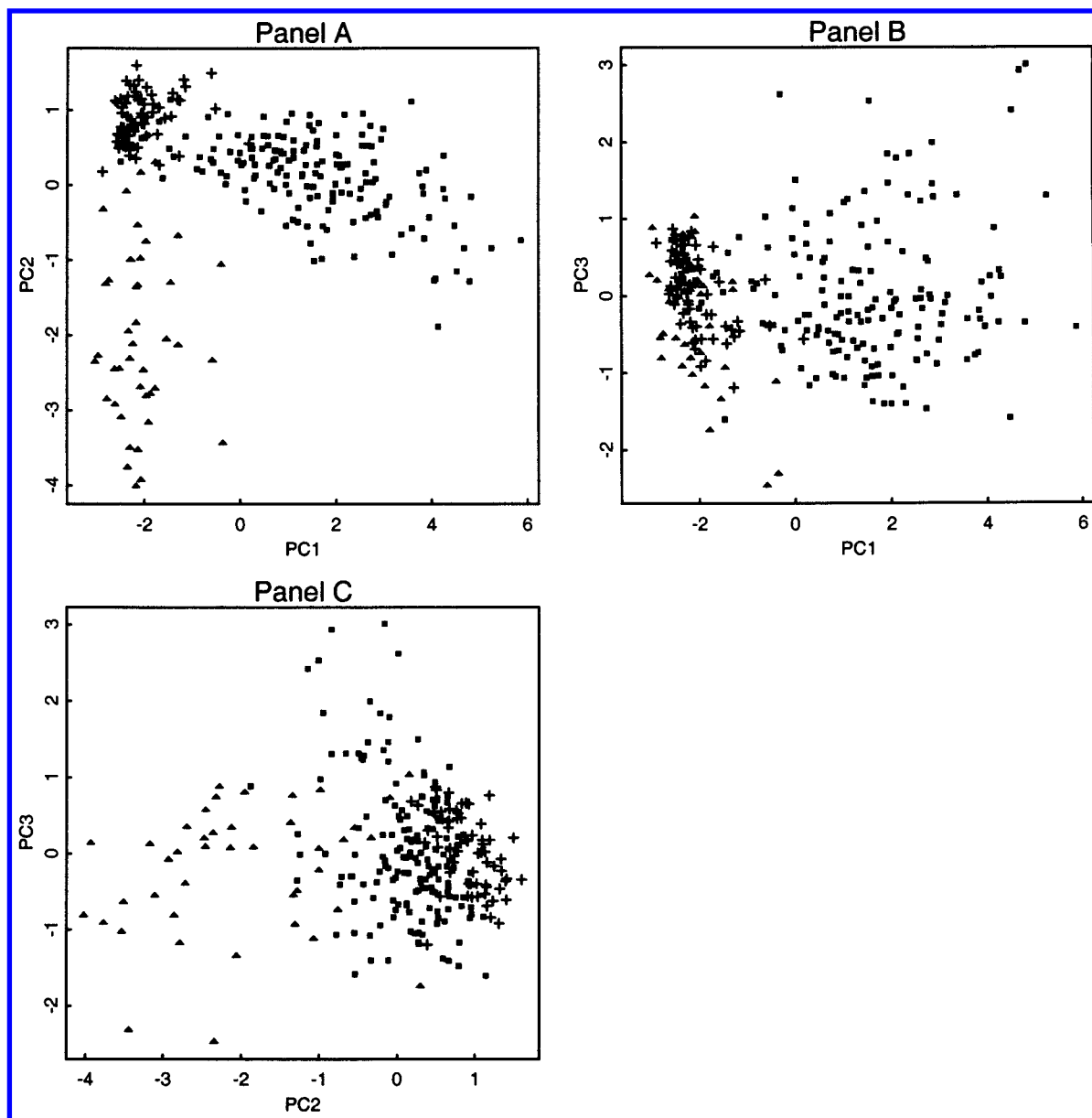


Figure 5. PCA score plots of the ellipticine–ellipticinium–camptothecin data set using relative activity patterns. Plot A with PC1 and PC2 gives the best separation, very similar to that seen in Figure 4C. (+) Ellipticines; (▲) ellipticiniums; (■) camptothecins.

activity patterns in the NCI screen and are believed to exhibit anticancer activity through different mechanisms of action. These compounds have been studied in our laboratory using molecular modeling and computational chemistry methods.^{23,24,58,61,62} Shown in Figure 4 are score plots for the first three PCs obtained using the original activity values, $-\log(\text{GI}_{50})$. When the first PC is used for plotting (Figure 4A and 4B), the grouping of compounds is not as good as that obtained by plotting the second and third PCs (Figure 4C). Our explanation is that the first PC, although it explains almost 90% of the variance in activity, is not very informative; it describes principally differences in mean activity of the compounds over all of the cells, and the mean activities themselves provide little information on the mechanism of action.

PCA was also performed on the relative activity patterns of the same data set. The score plots for the first three PCs are shown in Figure 5. Although the first two PCs explain only 36.4% and 12.0% of the variance, respectively, they clearly separate the data set into three groups. This grouping

is a clear reflection of the different mechanisms of action of compounds in the data set: camptothecins are thought to act through a unique mechanism by forming a DNA–topoisomerase 1–camptothecin ternary cleavable complex (as reviewed in ref 61); ellipticines and ellipticiniums act through other distinct mechanisms of action that are not as well-defined.²³

PCA loading plots can be used analogously to show the grouping of cell lines. Similar cell lines cluster together on the loading plot, whereas dissimilar cell lines scatter over the plot (data not shown).

Multidimensional Scaling (MDS). In multidimensional scaling, the dissimilarities or distances among members of a collection of objects are used to reduce dimensionality. Often the dimensionality is reduced to that of a two- or three-dimensional Euclidean space so that the relationships among objects can be visualized graphically by humans.^{63,64} The S-Plus “cmdscale” function is an implementation of classical metric multidimensional scaling. That is, locations of data points in the resulting lower-dimensional space are such that

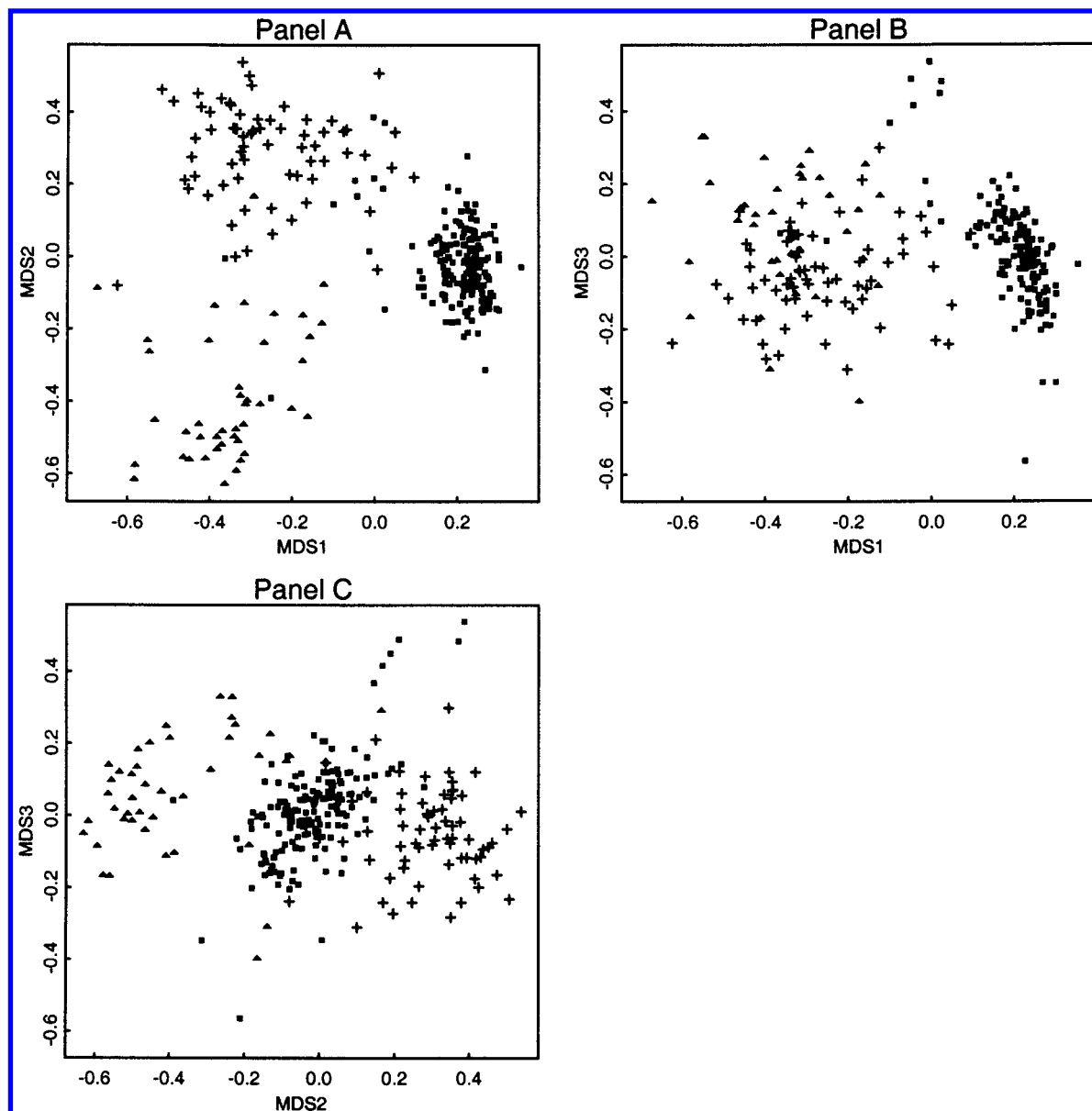


Figure 6. MDS plots of the ellipticine–ellipticinium–camptothecin data set showing clear groupings similar to those from PCA and cluster analysis. (+) Ellipticines; (▲) ellipticiniums; (■) camptothecins. The dimensionality was reduced from 60 to 3. Each panel represents a projection of the results onto the plane formed by two of the axes.

the overall distances among points in the reduced space are as close as possible globally to the original distances. There are various measures of goodness of fit for a solution. One of the commonly used criteria is the Stress function:

$$S = \sqrt{\frac{\sum_i \sum_j (d_{ij} - o_{ij})^2}{\sum_i \sum_j o_{ij}^2}}$$

where d_{ij} is the distance between data points i and j in the derived Euclidean space and o_{ij} is the original distance between i and j in the distance matrix. The aim in MDS is to minimize S . Many algorithms, e.g., the singular value decomposition method and Kruskal's negative gradient algorithm, have been developed to solve this problem.^{63,64} MDS is somewhat similar to PCA and, in fact, identical to it if one particular form of S is used.

We analyzed the ellipticine–ellipticinium–camptothecin data set using MDS. The original distance matrix for MDS was calculated by comparing the activity patterns. We found that better groupings of compounds were obtained when using the distance metric $(1 - r)$, where r is the Pearson correlation coefficient between the activity patterns of two compounds, to calculate the distance matrix. The results are shown in Figure 6. MDS appears to give better groupings than do the PCA score plots (Figures 4 and 5). The groupings are similar to those seen with hierarchical clustering (comparison data not shown). We also visualized the dissimilarities among the 60 cell lines in a 2D Euclidean space in terms of their responses to the 25 023 screened compounds (data not shown). The informational coherence of cell lines seen by cluster analysis or PCA could also be seen in MDS plots.

Hierarchical Cluster Analysis (HCA). Cluster analysis groups objects in terms of their similarities.⁶⁵ Hierarchical clustering produces a cluster tree in which similar objects are grouped together, whereas dissimilar objects are distant

from each other in the hierarchy. We used the agglomerative hierarchical clustering ("hclust") function in the S-Plus statistical package to cluster compounds in terms of their *in vitro* activity patterns across the 60 cell lines and to cluster the 60 cell lines in terms of their responses to 25 023 open compounds. At each step in the clustering process, the two clusters (or individuals) nearest to each other by some chosen criterion are combined to form one larger cluster. The procedure continues to aggregate clusters together until there is only one. The shape of the cluster tree formed by a particular data set is determined by a combination of the distance metric and clustering method. The distance metric can be Euclidean, maximum, Manhattan, binary, or some other user-selected measure. Clustering options implemented in S-Plus include average linkage, single linkage, and complete linkage, among others. For average linkage clustering, the distance between two clusters is the *average* of the distances between the data points in one cluster and the data points in the other. For single linkage clustering, the distance between two clusters is the *minimum* distance between any of the data points in the first cluster and any of the data points in the second. For complete linkage clustering, the distance is the *largest* distance between any of the data points in one cluster and any of the data points in the other.

Clustering of Compounds. We investigated different clustering algorithms but settled on average linkage clustering for this study. We used a distance metric of $(1 - r)$, where r is the Pearson correlation coefficient between the activity patterns of two compounds or clusters. In a previous study,²³ it appeared that the combination of average linkage clustering and $(1 - r)$ distance metric gave coherent results. All of the methods except single linkage yielded essentially the same regularities as those to be briefly shown here. When we used the average linkage hierarchical clustering algorithm with distance metric $(1 - r)$ to cluster the ellipticine–ellipticinium–camptothecin data set described earlier, compounds similar in activity pattern grouped together and compounds with dissimilar activity patterns were distant from each other in the cluster tree.²³ We also analyzed (data not shown) a data set of 131 "standard agents" whose mechanisms of action have been clearly classified.^{20,21} The following observations can be made from the cluster trees: (a) Compounds similar in chemical structure tend to group together. This is what one would expect because the general assumption in drug design is that similar structures lead to similar activities.^{16,17,20,21} (b) Compounds similar in mechanism of action tend to cluster together, even when they possess diverse chemical structures.^{16,17,20,21} (c) Compounds apparently similar in chemical structure but actually different in mechanism of action are distant from each other in the cluster tree.²³ These observations clearly corroborate the finding that, although activity data for a single cell line are not informative about the mechanism of action, activity *patterns* generated over the 60 cell lines encode incisive information about the mechanism.^{16–18,20,23}

Clustering of Cancer Cell Lines. We used the hclust (agglomerative hierarchical clustering) function in the S-Plus statistical package to cluster the 60 cell lines in terms of their responses to the 25 023 tested open compounds. Here, the cell lines were considered as "objects" and the compounds as "variables". In this study, we used complete linkage and the distance metric $(1 - r)$, where r is the

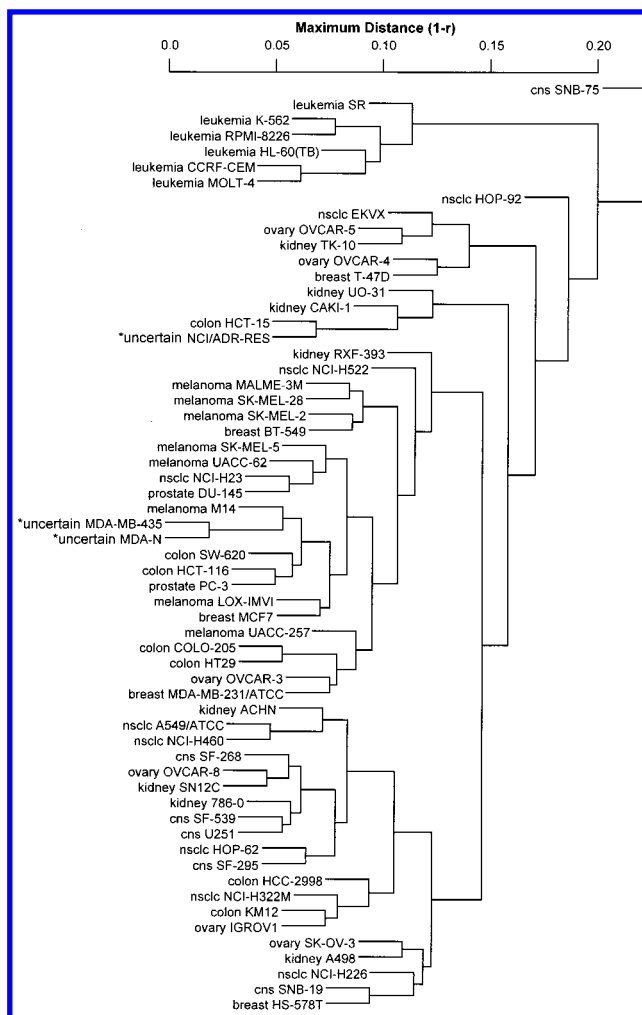


Figure 7. Clustering of the 60 human cancer cell lines based on original $-\log(GI_{50})$ activity patterns of 25 023 open compounds.

Pearson correlation coefficient relating the sensitivities of two cell lines to the 25 023 compounds. Figure 7, a cluster tree for the 60 cancer cell lines based on the original $-\log(GI_{50})$ values, shows informational coherence for leukemia and melanoma lines, but not for those derived from lung, colon, central nervous system, ovary, kidney, prostate, or breast.

Preprocessing of the original data by subtracting the mean of each row cannot have any effect on the clustering of compounds if the distance metric $(1 - r)$ is used. As shown in Figure 8, clustering of cancer cell lines based on the relative, rather than absolute or original, activity patterns appears in several respects to be more informative.

First, in addition to the fact that the six leukemias continue to cluster side by side, the melanoma and colon lines cluster better. If we examine in more detail the subtree for the six leukemia cell lines, we find that it reveals the biological and/or pathological similarity of these cell types. Cell lines CCRF-CEM and MOLT-4 are derived from acute lymphocytic leukemia (ALL), and they show much more similarity to each other in response to anticancer agents than to other leukemias. For other types of cell lines, it is hard to see clear grouping by organ of origin with clustering based on this large drug data set.

Second, the two cell lines MDA-N and MDA-MB-435 show a dramatic similarity to each other. In fact, MDA-N is

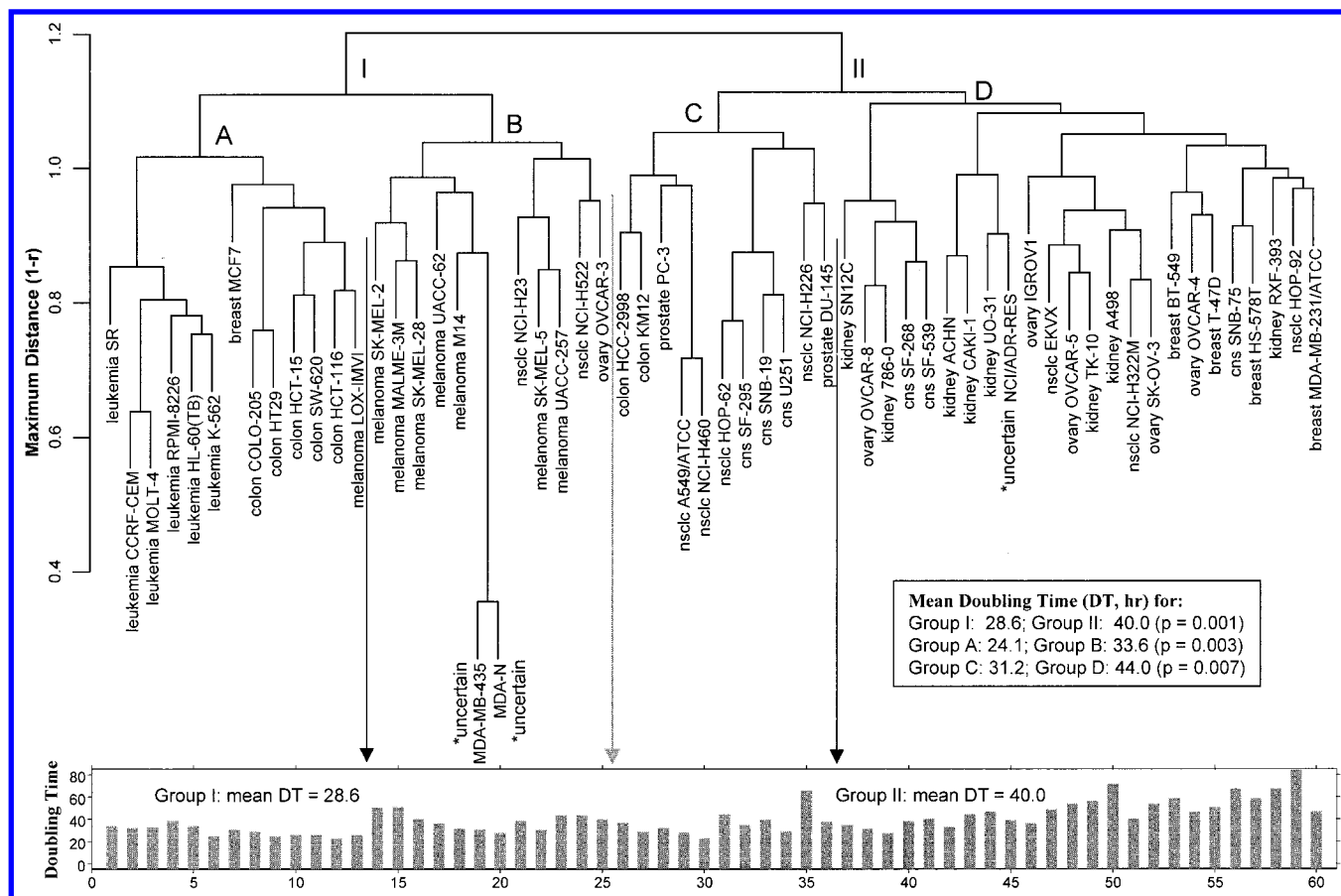


Figure 8. Clustering of 60 human cancer cell lines based on relative activity patterns of 25 023 open compounds. More apparent coherence based on organ of origin is seen in this tree than in that shown in Figure 7.

an ERB/B2-transfected derivative of MDA-MB-435. This is the only pair of cell lines in the NCI screen such that one is a genetically transfected form of the other. Therefore, they would be expected to be more similar to each other than other pairs of lines. It is reassuring to see from the screen data that they turn out to be so similar in terms of their responses to the 25 023 compounds.

Third, there is a significant difference in cell growth rates between the two groups of cell lines (see Figure 8). The doubling time for each line is listed in Table 2. Group I has an average doubling time of 28.6 h, whereas that for group II is 40.0 h. If one assumed independence, the two-sided Wilcoxon rank-sum p value for the null hypothesis of equal median growth rates between the two groups would be 0.0011. This difference is largely accounted for by the difference between groups A and D.

Genetic Function Approximation for Modeling Structure–Activity Relationships. Important goals in quantitative structure–activity relationship (QSAR) studies are (a) to develop a model for the relationships between molecular structure descriptors and biological responses for a set of training compounds and (b) to use the model to predict the biological responses of new compounds. Pharmaceutical research and development can be focused on those compounds with a predicted activity of the desired type. A number of statistical and artificial intelligence techniques, including multiple linear regression, principal component regression, partial least-squares regression, and artificial neural networks, have been applied to derive QSAR models. For one study,⁵⁸ we adopted the recently developed genetic

function approximation (GFA) algorithm^{66,67} to model quantitative structure–anticancer activity relationships.

GFA is a rational combination of MARS (multivariate adaptive regression splines)⁶⁸ and genetic algorithm⁶⁹ for modeling nonlinearity and for handling outliers. MARS is a regression method recently developed by J. H. Friedman.⁶⁸ It uses spline terms to partition the variable space in order to build a regression model locally. But MARS can handle only a moderate number of variables because it is computationally intensive if all possible combinations of variables are to be examined exhaustively.

GFA offers a means of exploring combinations of variables nonexhaustively but with a good chance of finding effective combinations even when the individual variables making up the combinations are not themselves highly predictive. Briefly, GFA first generates a population of QSAR models (e.g., 100) with descriptors (variables) that are randomly selected. The number of terms in these equations can be different. The terms can be linear, quadratic, or spline forms of a descriptor. Optimization of the choice of variables is done through an iterative crossover process. Two parent equations are cut at random points. Terms to the left of the cut in one parent are combined with terms to the right in the second parent to form a progeny equation. The progeny equation is evaluated using Friedman's "lack-of-fit" (LOF) measure:⁶⁸

$$\text{LOF} = \frac{\text{LSE}}{\left(1 - \frac{c + dp}{m}\right)^2}$$

where LSE is the least-squares error, c is the number of basis functions in the model, d is a smoothing parameter that controls the number of terms in the model equation (a larger value of d leads to fewer terms in the model), p is the number of features contained in all terms of the model, and m is the number of compounds in the training set. The LOF measure penalizes for the addition of terms to the equation (and consequent loss of degrees of freedom) in such a way as to resist overfitting.

If the progeny equation is better than the worst one in the population of 100 equations, it is added to the population and the worst equation is dropped out. At the end of evolution, the QSAR equations have stabilized, and the LOF value remains essentially unchanged over crossovers.

For a data set of 112 ellipticine analogues, we investigated their structure–anticancer activity relationships using GFA. The number of molecular descriptors was 69, including 49 Cerius² default descriptors, 17 partial atomic charges on the ellipticine ring-forming atoms, and 3 indicator variables that reflected our classification of the ellipticine analogues.²³ Seven activity indexes were studied separately. Overall, GFA gave much better results (fitting R^2 and cross-validated R^2) than did stepwise linear regression.⁵⁸

For a data set of camptothecin analogues, GFA also yielded better QSAR models than did stepwise linear regression or PLS.⁶² Furthermore, the GFA-derived QSAR models revealed information consistent with that from a comprehensive molecular modeling study on the DNA–topoisomerase 1 ternary cleavable complex with camptothecin.⁶¹ Information from both QSAR and molecular modeling studies has been successfully used to derive a pharmacophore with which to search the NCI DIS 3D database for topoisomerase 1 inhibitors (Fan et al., unpublished results).

CONCLUSIONS

We have used various kinds of multivariate statistical and artificial intelligence techniques to mine the large amounts of activity and structural data generated by the NCI anticancer drug discovery program. In agreement with analyses using various methods,^{16–24} the activity patterns clearly encode incisive information on the mechanisms of action of tested compounds, and coherence between chemical structures and activity patterns is often observed. Information gained from the mining of these databases can be used not only for drug discovery but also for basic studies in the molecular pharmacology of cancer. For example, analysis of the similarity of compounds in terms of their biological activity patterns can provide information on the mechanism of action.

When the cell lines were clustered in terms of their responses to the 25 023 open compounds, the relative activity patterns appeared to distinguish better among the cell lines than did the original $-\log(\text{GI}_{50})$ activity patterns. The mean-zero preprocessing procedure¹⁶ seemed to eliminate the noninformative “inherent” cytotoxicity, thus bringing out the informative differential cell responses. The same conclusion was also drawn from a series of PCA score plots. These observations reflect how important it is in any such data mining study to pay careful attention to the details of the algorithm, metric, and preprocessing steps used.

Portions of the NCI anticancer chemical structure and activity databases are searchable on the WWW at <http://cactus.cit.nih.gov/ncidb/> and <http://ntp.nci.nih.gov>. The cleaned drug data set, additional data sets related to the 60 cell lines, and tools for analysis of high-dimensional data can be found at <http://discover.nci.nih.gov>. The latter web site will be updated progressively with additional tools and genomic databases over the coming months.

ACKNOWLEDGMENT

We are grateful to Drs. K. W. Kohn and Y. Pommier for helpful discussions. We are also grateful to members of the Developmental Therapeutics Program, NCI, for generating the large amounts of chemosensitivity data analyzed in these studies. In particular, we are grateful to Dr. E. A. Sausville for his leadership of the program and to Drs. D. Scudiero and A. Monks and their staffs for conducting the drug screening program. These studies were made possible by the pioneering work of Dr. K. D. Paull, who died in 1998.

REFERENCES AND NOTES

- (1) Castell, J. V.; Gomez-Lechon, M. J., Eds. *In Vitro Methods in Pharmaceutical Research*; Academic Press: San Diego, 1997.
- (2) Ebert, C.; Lassiani, L.; Linda, P.; Nisi, C.; Alunni, S.; Clementi, S. Chemometric investigation of antitumor tests. *Quant. Struct.-Act. Relat.* **1984**, *3*, 143–147.
- (3) Nendza, M.; Seydel, J. K. Multivariate data analysis of various biological test systems used for quantification of ecotoxic compounds. *Quant. Struct.-Act. Relat.* **1988**, *7*, 165–174.
- (4) Moret, E. E.; Janssen, L. H. M. Principal component analysis of the cytostatic activity of 2,5-(bisaziridinyl)-1,4-benzoquinones. Rational approaches to the design of bioactive compounds. *Pharmacochemistry Library, QSAR*: 1991; Vol. 16, pp 381–384.
- (5) Benigni, R. Relationships between in vitro mutagenicity assays. *Mutagenesis* **1992**, *7*, 335–341.
- (6) Franke, R.; Gruska, A.; Presber, W. Combined factor and QSAR analysis for antibacterial and pharmacokinetic data from parallel biological measurements. *Pharmazie* **1994**, *49*, 600–605.
- (7) Ojasoo, T.; Raynaud, J. P. J. C. D. Affiliations among steroid receptors as revealed by multivariate analysis of steroid binding data. *J. Steroid Biochem. Mol. Biol.* **1994**, *48*, 31–46.
- (8) Ojasoo, T.; Raynaud, J. P., J. C. D. Correspondence factor analysis of steroid libraries. *Steroids* **1995**, *60*, 458–469.
- (9) Kauvar, L. M. Method to produce immunodiagnostic reagents. U.S. Patent 5300425. 1994, Terrapin Technologies, Inc.
- (10) Kauvar, L. M.; Higgins, D. L.; Villar, H. O.; Sportsman, J. R.; Engqvist-Goldstein, A.; Buckar, R.; Bauer, K. E.; Dilley, H.; Roche, D. M. Predicting ligand binding to proteins by affinity fingerprinting. *Chem. Biol.* **1995**, *2*, 107–118.
- (11) Villar, H. O. Affinity fingerprints: applications and implications. Network Science 1: <http://www.awod.com/netsci/Issues/Nov95/feature3.html>, 1995.
- (12) Warr, W. A. Combinatorial chemistry and molecular diversity: an overview. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 134–140.
- (13) Boyd, M. R. Status of the NCI preclinical antitumor drug discovery screen. In *Cancer: Principles and Practice of Oncology Update*; DeVita, V. T., Hellman, S., Rosenberg, S. A., Eds.; J. B. Lippincott: Philadelphia, 1989; Vol. 3, pp 1–12.
- (14) Boyd, M. R. The NCI in vitro anticancer drug discovery screen: concept, implementation, and operation, 1985–1995. In *Anticancer Drug Development Guide: Preclinical Screening, Clinical Trials, and Approval*; Teicher, B. A., Ed.; Humana Press: Totowa, NJ, 1997; pp 23–42.
- (15) Monks, A.; Scudiero, D. A.; Shoemaker, R. H.; Paull, K. D.; Vistica, D.; Hose, C.; Langley, J.; Cronise, P.; Vaigro-Wolff, A.; Gray-Goodrich, M.; Campell, H.; Mayo, J.; Boyd, M. R. Feasibility of a high-flux anticancer screen using a diverse panel of cultured human tumor lines. *J. Natl. Cancer Inst.* **1991**, *83*, 757–766.
- (16) Paull, K. D.; Shoemaker, R. H.; Hodes, L.; Monks, A.; Scudiero, D. A.; Rubinstein, L.; Plowman, J.; Boyd, M. R. Display and analysis of patterns of differential activity of drugs against human tumor cell lines: development of mean graph and COMPARE algorithm. *J. Natl. Cancer Inst.* **1989**, *81*, 1088–1092.
- (17) Paull, K. D.; Hamel, E.; Malspeis, L. Prediction of biochemical mechanism of action from the in vitro antitumor screen of the National

- Cancer Institute. In *Cancer Chemotherapeutic Agents*; Foye, W. O., Ed.; American Chemical Society: Washington, D.C., 1995; pp 9–45.
- (18) Weinstein, J. N.; Myers, T. G.; O'Connor, P. M.; Friend, S. H.; Fornace, A. J., Jr.; Kohn, K. W.; Fojo, T.; Bates, S. E.; Rubinstein, L. V.; Anderson, N. L.; Buolamwini, J. K.; van Osdol, W. W.; Monks, A. P.; Scudiero, D. A.; Sausville, E. A.; Zaharevitz, D. W.; Bunow, B.; Viswanadhan, V. N.; Johnson, G. S.; Wittes, R. E.; Paull, K. D. An information-intensive approach to the molecular pharmacology of cancer. *Science* **1997**, *275*, 343–349.
 - (19) Weinstein, J. N.; Myers, T. G.; Buolamwini, J.; Raghavan, K.; van Osdol, W.; Licht, J.; Viswanadhan, V. N.; Kohn, K. W.; Rubinstein, L. V.; Koutsoukos, A. D.; Monks, A.; Scudiero, D. A.; Anderson, N. L.; Zaharevitz, D.; Chabner, B. A.; Grever, M. R.; Paull, K. D. Predictive statistics and artificial intelligence in the U.S. National Cancer Institute's Drug Discovery Program for Cancer and AIDS. *Stem Cells* **1994**, *12*, 13–22.
 - (20) Weinstein, J. N.; Kohn, K. W.; Grever, M. R.; Viswanadhan, V. N.; Rubinstein, L. V.; Monks, A. P.; Scudiero, D. A.; Welch, L.; Koutsoukos, A. D.; Chiusa, A. J.; Paull, K. D. Neural computing in cancer drug development: Predicting mechanism of action. *Science* **1992**, *258*, 447–451.
 - (21) Koutsoukos, A. D.; Rubinstein, L. V.; Faraggi, D.; Simon, R. M.; Kalyandrug, S.; Weinstein, J. N.; Kohn, K. W.; Paull, K. D. Discrimination techniques applied to the NCI in vitro anti-tumour drug screen: predicting biochemical mechanism of action. *Stat. Med.* **1994**, *13*, 719–730.
 - (22) van Osdol, W. W.; Myers, T. G.; Paull, K. D.; Kohn, K. W.; Weinstein, J. N. Use of the Kohonen self-organizing map to study the mechanisms of action of chemotherapeutic agents. *J. Natl. Cancer Inst.* **1994**, *86*, 1853–1859.
 - (23) Shi, L. M.; Myers, T. G.; Fan, Y.; O'Connor, P. M.; Paull, K. D.; Friend, S. H.; Weinstein, J. N. Mining the National Cancer Institute Anticancer Drug Discovery Database: cluster analysis of ellipticine analogs with p53-inverse and central nervous system-selective patterns of activity. *Molec. Pharmacol.* **1998**, *53*, 241–251.
 - (24) Shi, L. M.; Fan, Y.; Myers, T. G.; Waltham, M.; Paull, K. D.; Weinstein, J. N. Mining the anticancer activity database generated by the NCI anticancer drug discovery program using statistical and artificial intelligence techniques. *Math. Modell. Sci. Comput.*, in press.
 - (25) Milne, G. W. A.; Miller, J. A. The NCI Drug Information System. 1. System overview. *J. Chem. Inf. Comput. Sci.* **1986**, *26*, 154–159.
 - (26) Milne, G. W. A.; Feldman, A.; Miller, J. A.; Daly, G. P.; Hammel, M. J. The NCI Drug Information System. 2. DIS pre-registry. *J. Chem. Inf. Comput. Sci.* **1986**, *26*, 159–168.
 - (27) Milne, G. W. A.; Feldman, A.; Miller, J. A.; Daly, G. P. The NCI Drug Information System. 3. The DIS chemistry module. *J. Chem. Inf. Comput. Sci.* **1986**, *26*, 168–179.
 - (28) Milne, G. W. A.; Miller, J. A.; Hoover, J. R. The NCI Drug Information System. 4. Inventory and shipping modules. *J. Chem. Inf. Comput. Sci.* **1986**, *26*, 179–185.
 - (29) Milne, G. W. A.; Nicklaus, M. C.; Driscoll, J. S.; Wang, S.; Zaharevitz, D. National Cancer Institute drug information system 3D database. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 1219–1224.
 - (30) Klopman, G.; Shi, L. M.; Ramu, A. Quantitative structure–activity relationship of multidrug resistance reversal agents. *Molec. Pharmacol.* **1997**, *52*, 323–334.
 - (31) Klopman, G. Artificial intelligence approach to structure–activity studies. Computer automated structure evaluation of biological activity of organic molecules. *J. Am. Chem. Soc.* **1984**, *106*, 7315–7320.
 - (32) Klopman, G. MULTICASE: a hierarchical computer automated structure evaluation program. *Quant. Struct.–Act. Relat.* **1992**, *11*, 176–184.
 - (33) Wang, S.; Zaharevitz, D. W.; Sharma, R.; Marquez, V. E.; Lewin, N. E.; Du, L.; Blumberg, P. M.; Milne, G. W. A. The discovery of novel, structurally diverse protein kinase C agonists through computer 3D-database search. Molecular modeling studies. *J. Med. Chem.* **1994**, *37*, 4479–4489.
 - (34) Wang, S.; Milne, G. W. A.; Yan, X.; Posey, I.; Nicklaus, M. C.; Graham, L.; Rice, W. G. Discovery of novel, non-peptide HIV-1 protease inhibitors by pharmacophore searching. *J. Med. Chem.* **1996**, *39*, 2047–2054.
 - (35) Neamati, N.; Hong, H.; Mazumder, A.; Wang, S.; Sunder, S.; Nicklaus, M. C.; Milne, G. W. A.; Proksa, B.; Pommier, Y. Depsides and depsidones as inhibitors of HIV-1 integrase: Discovery of novel inhibitors through 3D database searching. *J. Med. Chem.* **1997**, *40*, 942–951.
 - (36) Neamati, N.; Hong, H.; Sunder, S.; Milne, G. W. A.; Pommier, Y. Potent inhibitors of human immunodeficiency virus type 1 integrase: identification of a novel four-point pharmacophore and tetracyclines as novel inhibitors. *Molec. Pharmacol.* **1997**, *52*, 1041–1055.
 - (37) Hong, H.; Neamati, N.; Wang, S.; Nicklaus, M. C.; Mazumder, A.; Zhao, H.; Burke, T. R.; Pommier, Y.; Milne, G. W. A. Discovery of HIV-1 integrase inhibitors by pharmacophore searching. *J. Med. Chem.* **1997**, *40*, 930–936.
 - (38) Nicklaus, M. C.; Neamati, N.; Hong, H.; Mazumder, A.; Sunder, S.; Chen, J.; Milne, G. W. A.; Pommier, Y. HIV-1 integrase pharmacophore: Discovery of inhibitors through three-dimensional database searching. *J. Med. Chem.* **1997**, *40*, 920–929.
 - (39) Martin, Y. C.; Bures, M. G.; Willett, P. Searching databases of three-dimensional structures. In *Reviews in Computational Chemistry*; Lipkowitz, K. B., Boyd, D. B., Eds.; VCH: New York, 1990; Vol. 1.
 - (40) Good, A. C.; Mason, J. S. Three-dimensional structure database searches. In *Reviews in Computational Chemistry*; Lipkowitz, K. B., Boyd, D. B., Eds.; VCH: New York, 1996; Vol. 7.
 - (41) Wu, L.; Smythe, A. M.; Stinson, S. F.; Mullendore, L. A.; Monks, A.; Scudiero, D. A.; Paull, K. D.; Koutsoukos, A. D.; Rubinstein, L. V.; Boyd, M. R.; Shoemaker, R. H. Multidrug-resistant phenotype of disease-oriented panels of human tumor cell lines used for anticancer drug screening. *Cancer Res.* **1992**, *52*, 3029–3034.
 - (42) Bates, S. E.; Zhan, Z.; Dickstein, B.; Lee, J. S.; Scala, S.; Fojo, A. T.; Paull, K.; Wilson, W. Reversal of multidrug resistance. *Prog. Clin. Biol. Res.* **1994**, *389*, 33–37.
 - (43) Lee, J.-S.; Paull, K. D.; Alvarez, M.; Hose, C.; Monks, A.; Grever, M.; Fojo, A. T.; Bates, S. E. Rhodamine efflux patterns predict P-glycoprotein substrates in the National Cancer Institute drug screen. *Molec. Pharmacol.* **1994**, *46*, 627–638.
 - (44) Alvarez, M.; Paull, K. D.; Hose, C.; Lee, J.-S.; Weinstein, J. N.; Grever, M.; Bates, S.; Fojo, T. Generation of a drug resistance profile by quantitation of MDR-1/P-glycoprotein expression in the cell lines of the NCI anticancer drug screen. *J. Clin. Invest.* **1995**, *95*, 2205–2214.
 - (45) Bates, S. E.; Fojo, A. T.; Weinstein, J. N.; Myer, T. G.; Alvarez, M.; Paull, K. D.; Chabner, B. A. Molecular targets in the National Cancer Institute drug screen. *J. Cancer Res. Clin. Oncol.* **1995**, *121*, 495–500.
 - (46) Koo, H.-M.; Monks, A.; Mikheev, A.; Rubinstein, L. V.; Gray-Goodrich, M.; McWilliams, M. J.; Alvord, W. G.; Oie, H. K.; Gazdar, A. F.; Paull, K. D.; Zarbl, H.; Vande Woude, G. F. Enhanced sensitivity to 1-beta-D-arabinofuranosylcytosine and topoisomerase II inhibitors in tumor cell lines harboring activated ras oncogenes. *Cancer Res.* **1996**, *56*, 5211–5216.
 - (47) Izquierdo, M. A.; Shoemaker, R. H.; Flens, M. J.; Scheffer, G. L.; Wu, L.; Prather, T. R. Overlapping phenotypes of multidrug resistance among panels of human cancer-cell lines. *Int. J. Cancer* **1996**, *65*, 230–237.
 - (48) Freije, J. M.; Lawrence, J. A.; Hollingshead, M. G.; De la Rosa, A.; Narayanan, V.; Grever, M.; Sausville, E. A.; Paull, K.; Steeg, P. S. Identification of compounds with preferential inhibitory activity against low-Nm23-expressing human breast carcinoma and melanoma cell lines. *Nature Med.* **1997**, *3*, 395–401.
 - (49) Li, G.; Waltham, M.; Unsworth, E.; Treston, A.; Mushine, J.; Anderson, N. L.; Kohn, K. W.; Weinstein, J. N. Rapid protein identification from two-dimensional polyacrylamide gels by MALDI mass spectrometry. *Electrophoresis* **1997**, *18*, 391–402.
 - (50) Myers, T. G.; Waltham, M.; Li, G.; Buolamwini, J. K.; Scudiero, D. A.; Rubinstein, L. V.; Paull, K. D.; Sausville, E. A.; Anderson, N. L.; Weinstein, J. N. A protein expression database for the molecular pharmacology of cancer. *Electrophoresis* **1997**, *18*, 647–653.
 - (51) Wosikowski, K.; Schuurhuis, D.; Johnson, K.; Paull, K. D.; Myers, T. G.; Weinstein, J.; Bates, S. E. Identification of epidermal growth factor receptor and c-erbB2 pathway inhibitors by correlation with gene expression patterns. *J. Natl. Cancer Inst.* **1997**, *89*, 1505–1513.
 - (52) Weinstein, J. N. Fishing expeditions. *Science* **1998**, *282*, 627.
 - (53) Shi, L. M. Gene-chips (DNA microarrays). <http://www.gene-chips.com>, *Science* **1999**, *285*, 799.
 - (54) Ross, D. T.; Scherf, U.; Eisen, M. B.; Perou, C. M.; Spellman, P.; Iyer, V.; Jeffrey, S. S.; Van de Rijn, M.; Waltham, M.; Pergamenschikov, A.; Lee, J. C. F.; Lashkari, D.; Shalon, D.; Myers, T. G.; Weinstein, J. N.; Botstein, D.; Brown, P. O. Patterns of gene expression in sixty cell lines derived from tumors. *Nature Gen.*, in preparation.
 - (55) Scherf, U.; Ross, D. T.; Waltham, M.; Smith, L. H.; Lee, J. K.; Kohn, K. W.; Eisen, M. B.; Reinhold, W. C.; Myers, T. G.; Andrews, D. T.; Scudiero, D. A.; Sausville, E. A.; Pommier, Y.; Botstein, D.; Brown, P. O.; Weinstein, J. N. A cDNA microarray gene expression database for cancer drug discovery. *Nature Gen.*, in preparation.
 - (56) Boyd, M. R.; Paull, K. D. Some practical considerations and applications of the National Cancer Institute in vitro anticancer drug discovery screen. *Drug Dev. Res.* **1995**, *34*, 91–109.
 - (57) O'Connor, P. M.; Jackman, J.; Bae, I.; Myers, T. G.; Fan, S.; Mutoh, M.; Scudiero, D. A.; Monks, A.; Sausville, E. A.; Weinstein, J. N.; Friend, S.; Fornace, J.; A. J.; Kohn, K. W. Characterization of the p53-tumor suppressor pathway in cell lines of the National Cancer Institute anticancer drug screen and correlations with the growth-inhibitory potency of 123 anticancer agents. *Cancer Res.* **1997**, *57*, 4285–4300.

- (58) Shi, L. M.; Fan, Y.; Myers, T. G.; O'Connor, P. M.; Paull, K. D.; Friend, S. H.; Weinstein, J. N. Mining the NCI anticancer drug discovery databases: genetic function approximation for the QSAR study of anticancer ellipticine analogues. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 189–199.
- (59) *StatSci. S-PLUS Reference Manual*; MathSoft, Inc.: Seattle, 1993; Vol. 1.
- (60) Meglen, R. R. Examining large databases: a chemometric approach using principal component analysis. *J. Chemom.* **1991**, *5*, 163–179.
- (61) Fan, Y.; Weinstein, J. N.; Kohn, K. W.; Shi, L. M.; Pommier, Y. Molecular modeling studies of the DNA–topoisomerase I ternary cleavable complex with camptothecin. *J. Med. Chem.* **1998**, *41*, 2216–2226.
- (62) Fan, Y.; Shi, L. M.; Myers, T. G.; Weinstein, J. N. QSAR studies of anticancer camptothecins using genetic function approximation. *J. Comput.-Aided Mol. Des.*, in preparation.
- (63) Schiffman, S. C.; Reynolds, M. L.; Young, F. W. *Introduction to Multidimensional Scaling: Theory, Methods, and Applications*; Academic Press: New York, 1981.
- (64) Borg, I.; Groenen, P. *Modern Multidimensional Scaling: Theory and Applications*; Springer: New York, 1997.
- (65) Mirkin, B. G. *Mathematical Classification and Clustering*; Kluwer Academic: Boston, 1996.
- (66) Rogers, D.; Hopfinger, A. J. Application of genetic function approximation to quantitative structure–activity relationships and quantitative structure–property relationships. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 854–866.
- (67) Rogers, D. Some theory and examples of genetic function approximation with comparison to evolutionary techniques. In *Genetic Algorithms in Molecular Modeling*; Devillers, J., Ed.; Academic Press: London, 1996; pp 87–107.
- (68) Friedman, J. H. Multivariate adaptive regression splines (with discussion). *Ann. Stat.* **1991**, *19*, 1–141.
- (69) Holland, J. *Adaptation in Artificial and Natural Systems*; University of Michigan Press: Ann Arbor, 1975.

CI990087B