

Discrete Fourier Transform based Multivariate Image Analysis: Application to Modeling Aromatase Inhibitory Activity.

Stephen J. Barigye, Matheus P. Freitas, Priscila Ausina, Patricia Zancan, Mauro Sola-Penna, and Juan Alberto Castillo-Garit

ACS Comb. Sci., **Just Accepted Manuscript** • DOI: 10.1021/acscmbosci.7b00155 • Publication Date (Web): 03 Jan 2018

Downloaded from <http://pubs.acs.org> on January 5, 2018

Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a free service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are accessible to all readers and citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Discrete Fourier Transform based Multivariate Image Analysis: Application to Modeling Aromatase Inhibitory Activity.

Stephen J. Barigye*,[†]Matheus P. Freitas,[‡]Priscila Ausina,[£]Patricia Zancan,[§]Mauro Solapenna,[£]Juan A. Castillo-Garit.^{||}

[†]Department of Chemistry, McGill University, 801 Sherbrooke St. W., Montréal, QC, Canada.

[‡]Department of Chemistry, Federal University of Lavras, P.O. Box 3037, 37200-000, Lavras, MG, Brazil

[§]Laboratório de Oncobiologia Molecular (LabOMol), Departamento de Biotecnologia Farmacêutica, Faculdade de Farmácia, Universidade Federal do Rio de Janeiro, Rio de Janeiro, RJ, Brazil

[£]Laboratório de Enzimologia e Controle do Metabolismo (LabECoM), Departamento de Biotecnologia Farmacêutica, Faculdade de Farmácia, Universidade Federal do Rio de Janeiro, Rio de Janeiro, RJ, Brazil

^{||}Unidad de Toxicología Experimental, Universidad de Ciencias Médicas “Serafín Ruiz de Zárte Ruiz” Santa Clara, 50200, Villa Clara, Cuba.

1
2
3 **ABSTRACT:** We have recently generalized the formerly alignment-dependent **MIA-**
4 **QSAR** (acronym for **M**ultivariate **I**mage **A**nalysis applied to **Q**uantitative **S**tructure
5 **A**ctivity **R**elationships) method through the application of the Discrete Fourier Transform
6 (DFT), allowing for its application to non-congruent and structurally diverse chemical
7 compound datasets. Here, we report the first practical application of this method in the
8 screening of molecular entities of therapeutic interest, with the human aromatase inhibitory
9 activity as the case study. We develop an ensemble classification model based on 2D-DFT
10 MIA-QSAR descriptors, with which we screen 34 chemical compounds with possible
11 aromatase inhibitory activity from the NCI Diversity Set V (1593 compounds). These
12 compounds are docked into the aromatase active site and 10 most promissory compounds
13 selected for *in vitro* experimental validation. Of these compounds, **7419** (non-steroidal) and
14 **89201** (steroidal) demonstrate satisfactory antiproliferative and aromatase inhibitory
15 activities. The obtained results suggest that the 2D-DFT MIA-QSAR method may be useful
16 in the ligand-based virtual screening of NMEs of therapeutic utility.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

INTRODUCTION

Molecular modeling methods have over the years evolved, permeating all stages of the drug discovery process.¹⁻³ The extrapolation of concepts and/or methods particular to a wide range of disciplines such as biophysics, computer science, statistics, mathematics, structural biology and chemistry, to drug discovery paradigms, coupled with the availability of volumes of experimental data ranging from bioactivity parameters to crystallographic and NMR data, have fostered the advancement of existing computational tools to achieve more accurate models of chemical structures and their interactions with biological targets.⁴⁻

We have recently developed a generalization scheme for the **MIA-QSAR** (acronym for **Multivariate Image Analysis-Quantitative Structure Activity Relationship**) method based on the extrapolation of the 2D-Discrete Fourier Transform (2D-DFT) to chemical structural images.⁸⁻¹⁰ The MIA-QSAR method is based on the use of pixels comprising 2D chemical structure images as descriptors of the topo-chemical properties of molecules. Notwithstanding the successful application of the MIA-QSAR approach in the modeling of numerous bioactivities for over a decade, this method has been until recently exclusively applicable to congeneric datasets, due to the requirement that chemical structural images are aligned with respect to the basic molecular scaffold.¹⁰ This prerequisite as consequence rendered MIA-QSAR method inapplicable to virtual screening tasks for new molecular entities of therapeutic interest. We have demonstrated that the application of the 2D-DFT to chemical structural images eliminates the necessity of alignment of chemical images, allowing for the modeling of structurally diverse chemical datasets. Scheme 1 illustrates the workflow followed in the derivation of the 2D-DFT MIA-QSAR descriptors.

1
2
3 In the present report, we aim to demonstrate the possible utility of the 2D-DFT MIA-
4 QSAR approach in drug discovery process, via integrated ligand- and structure-based
5 virtual screening workflow for novel cytochrome P450 aromatase inhibitors, coupled with
6 posterior *in vitro* experimental validation. Human aromatase cytochrome P450 catalyzes
7 the final step of the conversion of androgens (*i.e.* androstenedione, 16 α -
8 hydroxytestosterone and testosterone) to estrogens (*i.e.* estrone, 17 β , 16 α -estriol and 17 β -
9 estradiol, respectively).¹¹ Clinical evidence has demonstrated that breast adenocarcinomas
10 of postmenopausal women possess much higher levels of 17 β -estradiol relative to plasma
11 concentrations and thus the inhibition of the aromatase catalyzed biosynthesis of estrogens
12 constitutes an attractive target for the prophylaxis and therapy of estrogen-dependent breast
13 cancer.^{12, 13}

30 RESULTS AND DISCUSSION

31
32
33 **Classification Model Building.** The built base classifier models for predicting the
34 aromatase inhibitory activity, based on the 2D-DFT MIA-QSAR descriptors, on the whole
35 demonstrated high accuracy, sensitivity and specificity, for the training and test sets,
36 respectively, with each classifier exhibiting a particular strength. Table 1 shows the
37 statistical parameters for the built classifier models. Data matrix employed in the model
38 building is provided as supporting information, S1. For the training set, the Boosted Trees
39 (BT) model yielded the highest accuracy (96.47%), sensitivity (93.75%) and specificity
40 (98.62%) values. As for the test set, the Artificial Neural Networks (ANN) model produced
41 the best accuracy (88.74%) and sensitivity (86.27%), followed by the k-Nearest Neighbor
42 (k-NN) model. On the other hand, the highest specificity was obtained with Linear
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Discriminant Analysis (LDA) model (95.00%), notwithstanding its much lower sensitivity
4 (71.57%), followed by the Support Vector Machine (SVM) model (92.50%). In this sense,
5
6 given the overall good performance of these base classifiers, a majority vote ensemble was
7
8 built to balance out the weaknesses of individual base classifiers, and thus enhancing the
9
10 robustness and predictive power of the classification models. The accuracy, sensitivity and
11
12 specificity values for the built ensemble were 86.94%, 81.37% and 91.67%, respectively
13
14 (Table 1). Moreover, the 2D-DFT MIA-QSAR models were validated using a set of decoys,
15
16 yielding high specificity (93.79%) for the ensemble, as well as the base classifiers BT
17
18 (98.81%), SVM(90.09%) and LDA (87.84%), while moderate specificity was obtained for
19
20 ANN (59.41%) and k-NN (60.84%). The generally good performance of these models
21
22 demonstrates the utility of the 2D-DFT MIA-QSAR descriptors in codifying vital chemical
23
24 structural information and may thus be considered as suitable in the virtual screening of
25
26 chemical compounds of therapeutic interest. Detailed data for the predictions of each
27
28 classifier is provided as Supporting Information, S2 (modeling set) and S3 (decoy set).
29
30
31
32
33
34

35 **Table 1 comes about here**

36
37 **Virtual Screening for Novel Aromatase Inhibitors.** Using the built multi-classifier
38
39 model, 34 compounds were screened from NCI Diversity Set V as potential aromatase
40
41 inhibitors. The screened compounds were docked into the aromatase active site to assess
42
43 their possible interaction modes and their binding affinities estimated according to the
44
45 PLANTS_{PLP} empirical scoring function.¹⁴ Based on the obtained scores, 10 chemical
46
47 compounds were selected for *in vitro* MCF-7 antiproliferative and aromatase inhibitory
48
49 activity assays. The structures of the screened compounds, as well as the corresponding
50
51 docking scores, are shown in Table 2 and demonstrate the structural diversity of the
52
53 selected compounds.
54
55
56
57
58
59
60

Table 2 comes about here

Inhibition of MCF-7 Cells Proliferation and Aromatase Activity. First, we evaluated the viability of TST-stimulated MCF-7 cells treated for 24 hours with 50 μM for each of the 10 chemical compounds. Tyrphostin A9 was used as a positive control. As shown in Figure 1a, 5 compounds (*i.e.* **7419**, **54709**, **89201**, **310354** and **661122**) were able to reduce cell viability by more than 50%. Subsequently, the inhibitor concentrations were lowered to 10 μM and the corresponding cell viability determined. At this concentration, only **7419** and **89201** negatively affected the viability of MCF-7 breast cancer cells (Figure 1b). In order to attest the cytochrome P450 aromatase inhibitory activity of the screened compounds, an evaluation of the enzymatic activity was performed using a cell-free extract of MCF-7 cells in presence of 10 μM for each of the compounds. As can be observed in Figure 2a, compounds **7419**, **89201** and Tyrphostin A9 exhibited aromatase inhibitory activity. Based on this result, **7419**, **89201** and Tyrphostin A9 were evaluated for their anti-proliferative activity against TST stimulated MCF-7 cells at 10 μM concentrations and all the three compounds strongly inhibited the proliferation of MCF-7 cells (Figure 2b). Finally, dose response curves were obtained for the antiproliferative activity of compounds **7419** [2-chloro-N-(2-(diethylamino)ethyl)-2-(perchlorocyclopenta-2,4-dien-1-ylidene)acetamide], **89201** [(8R,9S,13S,14S,17S)-17-hydroxy-13-methyl-7,8,9,11,12,13,14,15,16,17-decahydro-6H-cyclopenta[a]phenanthren-3-yl bis(2-chloroethyl)carbamate] and **Tyrphostin A9**, and the IC_{50} values were determined to be $4.9 \pm 1.1 \mu\text{M}$, $4.1 \pm 0.9 \mu\text{M}$ and $0.2 \pm 0.1 \mu\text{M}$, respectively (Figure 3). Note that compound **89201** is in fact a well-known alkylating agent (*i.e.* estramustine) whose phosphate ester is a clinically approved drug that has been employed for decades in prostate cancer therapy.¹⁵

¹⁶ On the other hand, although off-target activity evaluations are beyond the objectives of

1
2
3 the present study, it should be noted that compound **7419** is in fact a cyclic vinyl chloride
4
5 derivative. This class of compounds is notorious for presenting several nonspecific
6
7 interactions and therefore caution should be exercised in pursuing **7419** as a clinically
8
9 viable candidate.
10

11
12 **Figure 1 comes about here**

13
14 **Figure 2 comes about here**

15
16
17 **Figure 3 comes about here**
18
19
20
21
22

23 **CONCLUSION**

24
25 The results obtained with the 2D-DFT based generalization scheme of the MIA-
26
27 QSAR method in the screening of novel aromatase inhibitors, albeit the modest chemical
28
29 space explored, demonstrate the utility of this approach in the modeling of the bioactivities
30
31 of chemical datasets and in the virtual screening of chemical compounds with desired
32
33 therapeutic profiles. Moreover, in the *in silico* screening experiment, we employed a
34
35 combination of ligand-based and structure-based virtual screening approaches based on the
36
37 notion that one method helps to minimize the weaknesses of the other. In a broader sense,
38
39 the ligand-based VS techniques are based on the similarity principle in that similar
40
41 compounds demonstrate similar bioactivity profiles. However, this concept may be flawed
42
43 given that similar compounds may not yield similar interaction profiles with a given active
44
45 site and thus structure-based VS serves as a validation of the ligand-based VS results. On
46
47 the other hand, compounds with high binding affinity to an active site may not yield the
48
49 desired activity (*i.e.* agonists instead of antagonists and *vice versa*). In this sense, a
50
51 combined ligand and structure-based approach helps to limit the weaknesses of each
52
53
54
55
56
57
58
59
60

1
2
3 method. Future initiatives include the application of the inverse DFT to identify functional
4 groups or substructures vital for modeled bioactivities. Moreover, other multivariate image
5 processing techniques such as Discrete Cosine, Wavelets and Walsh-Hadamard transforms
6 will be explored for their utility in the MIA-QSAR context.
7
8
9

10 11 12 13 14 **MATERIALS AND METHODS**

15
16
17 **Chemical Dataset and 2D-DFT MIA-QSAR.** For this study, we extracted from the
18 literature a dataset of 973 chemical compounds with known aromatase inhibitory activity
19 profiles.¹⁷ For each of these chemical structures, their corresponding structural images were
20 retrieved and transformed into magnitude spectra whose coordinates correspond to the
21 spatial frequencies and the pixel values the intensities (*i.e.* Fourier coefficient magnitudes)
22 of these frequencies. Posteriorly, consistent with the MIA-QSAR approach, the pixels of
23 the magnitude spectra were considered as the variables codifying topo-chemical
24 information for each of these structures (scheme 1). For a detailed treatise on the theoretical
25 aspects of the 2D-DFT MIA-QSAR method, see references.^{8,9} Given the characteristic
26 symmetry presented by Fourier spectra, only half of the data points were employed. As
27 anticipated, considering the pixels of magnitude spectra as descriptors generates high
28 dimensional data matrices and thus dimensionality reduction procedures are necessary in
29 order to work with manageable matrix dimensions. To this end, a combination of
30 unsupervised and supervised feature selection procedures were performed using our freely
31 available software, IMMAN (acronym for **I**nformation theory-based **cheMoMetrics**
32 **AN**alysis),¹⁸ and a low-dimensional data was obtained with which the posterior
33 experiments were performed.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **Data Set Splitting and Classification Model Building.** The chemical compounds
4
5 were split into actives ($pIC_{50} > 6.0$) and inactives ($pIC_{50} < 6.0$), comprising of 390 and
6
7 483 chemical structures, respectively. Consistent with good practices of QSAR modeling,
8
9 each group was divided into training and test sets, using hierarchical and k -means cluster
10
11 analysis methods, with Ward's algorithm as the amalgamation rule and the squared
12
13 Euclidean distance as the proximity function, respectively. Consequently, training and test
14
15 set sizes of 651 and 222 chemical compounds, respectively, were obtained. We built a
16
17 majority (hard) voting classifier comprising of 5 conceptually distinct base classifiers, *i.e.*
18
19 Linear Discriminant Analysis (LDA), Artificial Neural Networks (ANN), Support Vector
20
21 Machines (SVM), Boosted Trees (BT) and k-Nearest Neighbor (k-NN). For the ANN
22
23 approach, the MLP (multi-layer perceptron) model was employed with network
24
25 architecture of 14-10-2, while the logistic and SOS functions were employed as the hidden
26
27 activation and error functions, respectively. For the SVM model, the RBF kernel was
28
29 employed with the tunable parameters γ and C equal to 0.55 and 13, respectively. The
30
31 meta-parameters shrinkage (ν) and training set fraction (η) for the BT model were set to
32
33 0.13 and 0.30, respectively. For the LDA model, the tolerance parameter was fixed at 0.01
34
35 and the forward-step wise selection procedure was employed as the variable selection
36
37 strategy. The dataset splitting and model building was performed using STATISTICA
38
39 software.¹⁹ Several reports in the literature have demonstrated that consensus models yield
40
41 the most accurate predictions, in addition to eliminating the dependence on training and set
42
43 statistics in the selection of models for virtual screening tasks.²⁰⁻²²

44
45 **Virtual Screening and Molecular Docking Validation.** The NCI Diversity Set V
46
47 (1593 compounds) was screened for compounds with potential inhibitory activity against
48
49 CYP19A1 using the built ensemble model. The choice of this dataset was attributed to its
50
51
52
53
54
55
56
57
58
59
60

1
2
3 inherent diversity and the availability of vialled samples from the DTP repository for
4
5 posterior experimental validation. In the generation of this diversity dataset, priority was
6
7 given to chemical compounds with pharmacologically desirable characters (*i.e.* devoid of
8
9 weakly bonded heteroatoms, obvious leaving groups, polycyclic aromatic hydrocarbons,
10
11 etc.) and thus constitutes an interesting source of possible novel lead compounds,
12
13
14 notwithstanding the small size. The majority vote criterion was considered in the prediction
15
16 of the aromatase inhibitory activity of the screened compounds [*i.e.* compounds were
17
18 considered as active or inactive if the probability (active or inactive) ≥ 0.6].
19
20

21
22 On the other hand, while ligand-based virtual screening is in a broader sense, based
23
24 on the principle of chemical structural similarity, in that compounds with similar chemical
25
26 structural characteristics should exhibit similar bioactivity, favorable protein-ligand
27
28 interactions play a central role in determining the possible activity of screened compounds.
29
30 In this sense, the chemical compounds predicted as promissory aromatase inhibitors were
31
32 posteriorly screened using a protein-ligand docking procedure and the ligands with the
33
34 lowest free-energy binding were selected. For the docking procedure, the CLC Drug
35
36 Discovery Workbench was used and a 2.75 Å crystal structure of human placental
37
38 aromatase cytochrome P450 in complex with androstenedione (PDB code 3S79) was
39
40 obtained from the Protein Data Bank. The carboxylate moiety of the Asp309 side chain of
41
42 the aromatase structure was protonated consistent with previous reports that validated its
43
44 role as a proton donor in the aromatization reaction.^{23, 24} The most promissory chemical
45
46 compounds based on the integrated ligand-based and structure-based virtual screening were
47
48 posteriorly submitted to experimental validation to attest their antiproliferative and
49
50 inhibitory activity. Moreover, retrospective validation of the built 2D-DFT MIA-QSAR
51
52 models was performed using a decoy set comprising of 3348 chemical compounds
53
54
55
56
57
58
59
60

1
2
3 (supporting information, S3) generated based on the 50 most active aromatase inhibitors,
4
5 using the Directory of Useful Decoys, Enhanced (DUD-E) utility.²⁵
6

7
8 **Chemical Compounds and MCF-7 Cell Line.** The screened chemical compounds
9
10 were kindly provided by Development Therapeutics Program (DTP), National Cancer
11
12 Institute, USA. The human breast cancer cell line MCF-7 was obtained from the Cell Bank
13
14 of Hospital Universitário Clementino Fraga Filho, UFRJ, Brazil, and was maintained in
15
16 DMEM (Dulbecco's modified Eagle's medium; Invitrogen, São Paulo, SP, Brazil)
17
18 supplemented with 10% (v/v) FBS (fetal bovine serum; Invitrogen) and L-glutamine.²⁶ For
19
20 all the experiments, cells growth was stimulated by 1 nM testosterone (TST). The TST -
21
22 stimulated cell viability and cell proliferation experiments, the difference between the
23
24 results in the presence and in the absence of TST were considered.
25
26

27
28 **Cell Viability.** The MCF-7 cells viability was assessed through the MTT assay as
29
30 described previously.²⁷ Cells were seeded in 96-well plates (2×10^3 cells/well) and
31
32 incubated at 37 °C and 5% CO₂ for 48 hours. After this initial incubation, the medium was
33
34 replaced by a fresh medium containing the 1 nMTST and the desired concentration of the
35
36 drugs. This procedure was repeated 5-times at each 2 days of treatment. Then, the medium
37
38 was removed, fresh medium was added, and the cells were returned to the incubator in the
39
40 presence of different drugs used. In the end of the treatment, the medium was removed and
41
42 the cells were incubated with 5 mg/mL MTT reagent (3,4,5-dimethylazol-2,5-
43
44 diphenyltetrazolium bromide, Sigma-Aldrich Co., St. Louis, MO, USA) for 3 h. Thereafter,
45
46 the formazan crystals formed were dissolved in DMSO, and the absorbance at 560 nm was
47
48 evaluated using a VICTOR3 multilabel microplate reader (PerkinElmer, Waltham, MA,
49
50 USA) with subtraction of the background absorbance at 670 nm.²⁷
51
52
53
54
55
56
57
58
59
60

1
2
3 **Cell Proliferation Assay.** The MCF-7 cells proliferation was assessed through the
4
5 CyQuant kit (Thermo Fisher Scientific, Waltham, MA, USA) that evaluates the DNA
6
7 content of the cells as a measurement of cell number. For this, cells were seeded in 96-well
8
9 plates (2×10^3 cells/well) and incubated at 37 °C and 5% CO₂ for 48 hours. After this initial
10
11 incubation, the medium was replaced by a fresh medium containing the 1 nMTST and the
12
13 desired concentration of the drugs. This procedure was repeated 5-times at each 2 days of
14
15 treatment. Then, the medium was removed, a fresh medium added and the cells were
16
17 returned to the incubator in the presence of different drugs used. At the end of the
18
19 treatment, the medium was removed and 200 μL of the freshly prepared CyQuant reagent
20
21 was added to each well and after 5 min the fluorescence emission at 520 nm was read in a
22
23 VICTOR3 multilabel microplate reader (PerkinElmer, Waltham, MA, USA).
24
25

26
27
28 **Aromatase Inhibitory Activity.** Aromatase activity was assessed in a MCF-7 free
29
30 cell homogenate. For this, MCF-7 cells (10^5 cells) were seeded in 75 cm² culture flasks and
31
32 grown in DMEN supplemented with 10% FBS at 37 °C and 5% CO₂. When cells reached
33
34 approximately 80% of confluence, cells were hashed and homogenized in Cell Lysis Buffer
35
36 (Cell Signaling Technology, #9803, Danvers, MA, USA) in a proportion of 5×10^6 cells in
37
38 10 mL of buffer. Aromatase activity was assayed in a reaction mixture containing 50
39
40 mM Tris-HCl (pH 7.4), 5 mM MgCl₂, 120 mM KCl, 0.1 mM TST and 0.2 mM NADPH.
41
42 The reaction was initiated by the addition of 100 μL of the MCF-7 cell free homogenate.
43
44 Appropriate controls were performed to consider the effects of TST vehicle (ethanol) and
45
46 the drugs vehicle (DMSO). Oxidation of NADPH to NADP⁺ was followed measuring the
47
48 fluorescence emission at 340 nm in a VICTOR3 multilabel microplate reader (PerkinElmer,
49
50 Waltham, MA, USA). The aromatase activity was calculated discounting the slope of the
51
52
53
54
55
56
57
58
59
60

1
2
3 fluorescence decay in the presence of TST from its counterpart in the absence of TST. All
4
5 experiments were run in triplicate in a series of three independent experiments (n = 3).
6
7
8
9

10 **ABBREVIATIONS**

11
12 TST, Testosterone; DFT, Discrete Fourier Transform; MIA-QSAR, Multivariate Image
13
14 Analysis applied to Quantitative Structure–Activity Relationships
15
16
17
18

19 **ACKNOWLEDGMENT**

20
21 The present work was supported by grants from Fundação Carlos Chagas Filho de Apoio à
22
23 Pesquisa do Estado do Rio de Janeiro (FAPERJ), Fundação de Amparo à Pesquisa do
24
25 Estado de Minas Gerais (FAPEMIG), and Conselho Nacional de Desenvolvimento
26
27 Científico e Tecnológico (CNPq). The authors would also like to thank the Drug Synthesis
28
29 and Chemistry Branch (DSCB), Development Therapeutics Program (DTP), National
30
31 Cancer Institute, USA for providing the chemical compounds for the *in vitro* assays.
32
33
34
35
36
37

38 **Supporting Information.** Data matrix employed in the model building, activity predictions
39
40 for each compound in employed dataset. This material is available free of charge via the
41
42 Internet at <http://pubs.acs.org>
43
44
45
46

47 **AUTHOR INFORMATION**

48 **Corresponding Author**

49
50
51 *Phone: 1-514-660-9351. E-mail: stephen.barigye@mcgill.ca
52
53
54
55

56 **REFERENCES**

1. Jorgensen, W. L., The Many Roles of Computation in Drug Discovery. *Science* **2004**, *303*, 1813-1818.
2. Kar, S.; Roy, K., How far can virtual screening take us in drug discovery? *Expert Opin Drug Discov.* **2013**, 245-261.
3. Passeri, G. I.; Trisciuzzi, D.; Alberga, D.; Siragusa, L.; Leonetti, F.; Mangiatordi, G. F.; Nicolotti, O., Strategies of Virtual Screening in Medicinal Chemistry. *Int. J. Quant. Struct. Prop. Relat.* **2018**, *3*, 134-160.
4. Todeschini, R.; Consonni, V., *Molecular Descriptors for Chemoinformatics*. Wiley-VCH: Weinheim, 2009; Vol. 1, p 1265.
5. Barigye, S. J.; Marrero-Ponce, Y.; Pérez-Giménez, F.; Bonchev, D., Trends in information theory-based chemical structure codification. *Mol. Divers.* **2014**, *18*, 673-686.
6. Sliwoski, G.; Kothiwale, S.; Meiler, J.; Lowe, E. W., Computational methods in drug discovery. *Pharmacol. Rev.* **2014**, *66*, 334-395.
7. Cherkasov, A.; Muratov, E. N.; Fourches, D.; Varnek, A.; Baskin, I. I.; Cronin, M.; Dearden, J.; Gramatica, P.; Martin, Y. C.; Todeschini, R., QSAR modeling: where have you been? Where are you going to? *J. Med. Chem.* **2014**, *57*, 4977-5010.
8. Barigye, S. J.; Freitas, M. P., 2D-Discrete Fourier Transform: Generalization of the MIA-QSAR strategy in molecular modeling. *Chemom. Intell. Lab. Syst.* **2015**, *143*, 79-84.
9. Barigye, S. J.; Freitas, M. P., Is molecular alignment an indispensable requirement in the MIA-QSAR method? *J Comp Chem* **2015**, *36*, 1748-1755.
10. Barigye, S. J.; de Freitas, M. P., Ten years of the MIA-QSAR strategy: historical development and applications. *Int. J. Quant. Struct. Prop. Relat.* **2016**, *1*, 64-77.
11. Ghosh, D.; Griswold, J.; Erman, M.; Pangborn, W., Structural basis for androgen specificity and oestrogen synthesis in human aromatase. *Nature* **2009**, *457*, 219-223.

- 1
2
3 12. Yadav, M. R.; Barmade, M. A.; Tamboli, R. S.; Murumkar, P. R., Developing
4 steroidal aromatase inhibitors-an effective armament to win the battle against breast cancer.
5
6 *Eur. J. Med. Chem.* **2015**, *105*, 1-38.
7
8
9
10 13. Ghosh, D.; Lo, J.; Egbuta, C., Recent progress in the discovery of next generation
11 inhibitors of aromatase from the structure–function perspective. *J. Med. Chem.* **2016**, *59*,
12 5131-5148.
13
14
15
16 14. Korb, O.; Stutzle, T.; Exner, T. E., Empirical scoring functions for advanced
17 protein– ligand docking with PLANTS. *J. Chem. Inf. Model.* **2009**, *49*, 84-96.
18
19
20 15. Di Lorenzo, G., Estramustine in prostate cancer: new look at an old drug. *Lancet*
21 *Oncol.* **2007**, *8*, 959-961.
22
23
24
25 16. Rautio, J.; Kumpulainen, H.; Heimbach, T.; Oliyai, R.; Oh, D.; Järvinen, T.;
26 Savolainen, J., Prodrugs: design and clinical applications. *Nat. Rev. Drug Discov.* **2008**, *7*,
27 255-270.
28
29
30
31 17. Worachartcheewan, A.; Mandi, P.; Prachayasittikul, V.; Toropova, A. P.; Toropov,
32 A. A.; Nantasenamat, C., Large-scale QSAR study of aromatase inhibitors using SMILES-
33 based descriptors. *Chemom. Intell. Lab. Syst.* **2014**, *138*, 120-126.
34
35
36
37 18. Urias, R. W. P.; Barigye, S. J.; Marrero-Ponce, Y.; García-Jacas, C. R.; Valdes-
38 Martini, J. R.; Perez-Gimenez, F., IMMAN: free software for information theory-based
39 chemometric analysis. *Mol. Divers.* **2015**, *19*, 305-319.
40
41
42
43 19. Statsoft-Team *STATISTICA. Data Analysis Software System*, version 6.0; StatSoft:
44 Tulsa, 2001.
45
46
47 20. Rokach, L., Ensemble-based classifiers. *Artif Intell Rev* **2010**, *33*, 1-39.
48
49
50
51 21. Marrero-Ponce, Y.; Siverio-Mota, D.; Gálvez-Llompert, M.; Recio, M. C.; Giner, R.
52 M.; García-Domènech, R.; Torrens, F.; Arán, V. J.; Cordero-Maldonado, M. L.; Esguera,
53
54
55
56
57
58
59
60

1
2
3 C. V., Discovery of novel anti-inflammatory drug-like compounds by aligning in silico and
4 in vivo screening: the nitroindazolinone chemotype. *Eur. J. Med. Chem.* **2011**, *46*, 5736-
5 5753.
6
7

8
9
10 22. Ventura, C.; Latino, D. A.; Martins, F., Comparison of multiple linear regressions
11 and neural networks based QSAR models for the design of new antitubercular compounds.
12 *Eur. J. Med. Chem.* **2013**, *70*, 831-845.
13
14

15
16
17 23. Di Nardo, G.; Breitner, M.; Bandino, A.; Ghosh, D.; Jennings, G. K.; Hackett, J. C.;
18 Gilardi, G., Evidence for an elevated aspartate pKa in the active site of human aromatase. *J.*
19 *Biol. Chem* **2015**, *290*, 1186-1196.
20
21

22
23
24 24. Adhikari, N.; Amin, S. A.; Saha, A.; Jha, T., Combating breast cancer with non-
25 steroidal aromatase inhibitors (NSAIs): Understanding the chemico-biological interactions
26 through comparative SAR/QSAR study. *Eur. J. Med. Chem.* **2017**, 365-438.
27
28

29
30
31 25. Mysinger, M. M.; Carchia, M.; Irwin, J. J.; Shoichet, B. K., Directory of useful
32 decoys, enhanced (DUD-E): better ligands and decoys for better benchmarking. *J. Med.*
33 *Chem.* **2012**, *55*, 6582-6594.
34
35

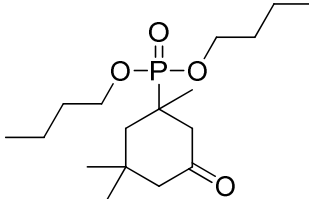
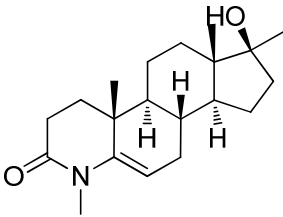
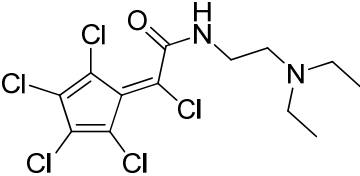
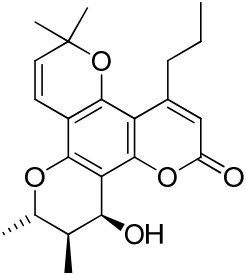
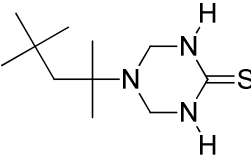
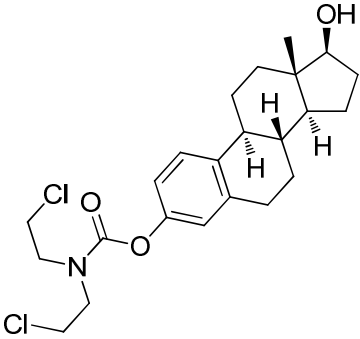
36
37
38 26. Zancan, P.; Sola-Penna, M.; Furtado, C. M.; Da Silva, D., Differential expression of
39 phosphofructokinase-1 isoforms correlates with the glycolytic efficiency of breast cancer
40 cells. *Mol. Genet. Metab.* **2010**, *100*, 372-378.
41
42

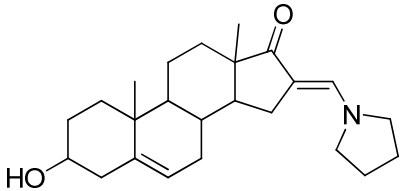
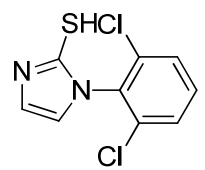
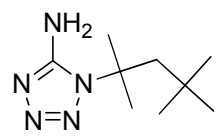
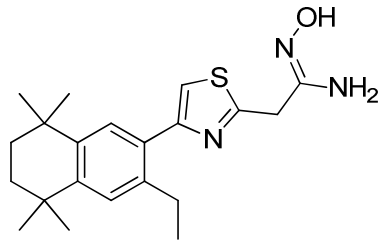
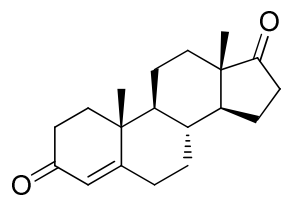
43
44
45 27. Spitz, G. A.; Furtado, C. M.; Sola-Penna, M.; Zancan, P., Acetylsalicylic acid and
46 salicylic acid decrease tumor cell viability and glucose metabolism modulating 6-
47 phosphofructo-1-kinase structure and activity. *Biochem. Pharmacol.* **2009**, *77*, 46-53.
48
49

Table 1. Statistical parameters of the base classifier models and majority-vote ensemble model.

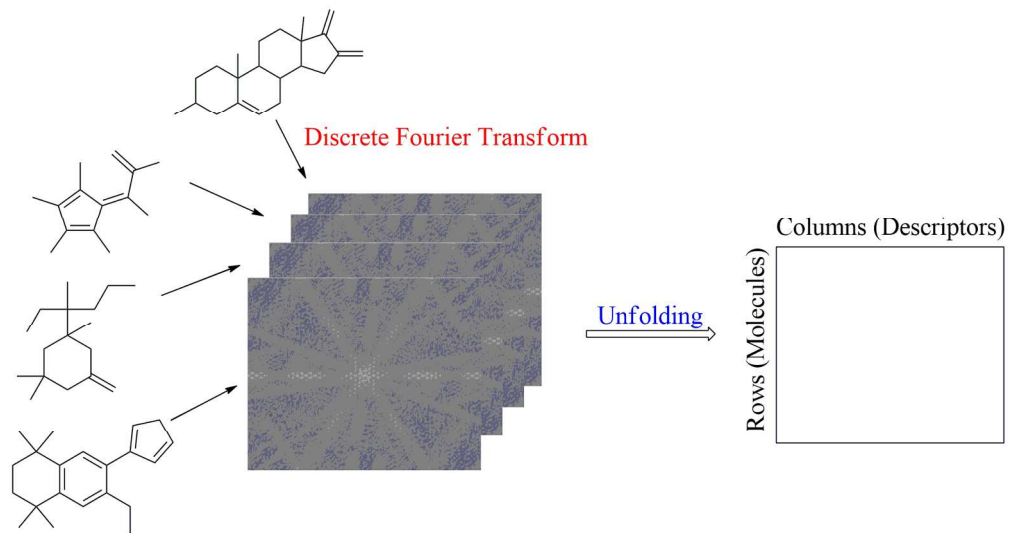
Classification Models		ANN	LDA	SVM	BT	k-NN	Ensemble
Accuracy	Training	87.71%	80.34%	87.71%	96.47%	-	-
	Test	88.74%	84.23%	85.14%	84.68%	85.59%	86.94%
Sensitivity	Training	87.15%	71.88%	80.56%	93.75%	-	-
	Test	86.27%	71.57%	76.47%	81.37%	82.35%	81.37%
Specificity	Training	88.15%	87.05%	93.39%	98.62%	-	-
	Test	90.83%	95.00%	92.50%	87.50%	88.33%	91.67%
	Decoy	59.41%	87.84%	90.09%	98.81%	60.84%	93.79%

Table 2. Chemical structures, docking scores and experimental IC₅₀ values for screened compounds with possible aromatase inhibitory activity, based on the 2D-DFT MIA-QSAR ensemble classification model and the molecular docking technique.

ID	Chemical Structure	PLANTS _{PLP} Score	IC ₅₀ , μM
14506		-73.43	> 50.0
36923		-68.31	> 50.0
7419		-64.24	4.9±1.1
661122		-62.36	> 10.0
319034		-61.69	> 50.0
89201		-61.43	4.1±0.9

1				
2				
3				
4				
5	54709		-57.03	> 10.0
6				
7				
8				
9				
10				
11	321506		-54.69	> 50.0
12				
13				
14				
15				
16	11128		-53.39	> 50.0
17				
18				
19				
20				
21				
22	310354		-51.46	> 10.0
23				
24				
25				
26				
27				
28				
29				
30	Androstenedione		-84.54	-
31				
32				
33				
34				
35				
36				
37				
38				
39				
40				
41				
42				
43				
44				
45				
46				
47				
48				
49				
50				
51				
52				
53				
54				
55				
56				
57				
58				
59				
60				

1
2
3 **Scheme1.** Workflow followed in the derivation of the 2D-DFT MIA-QSAR descriptors.
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Scheme 1

168x90mm (300 x 300 DPI)

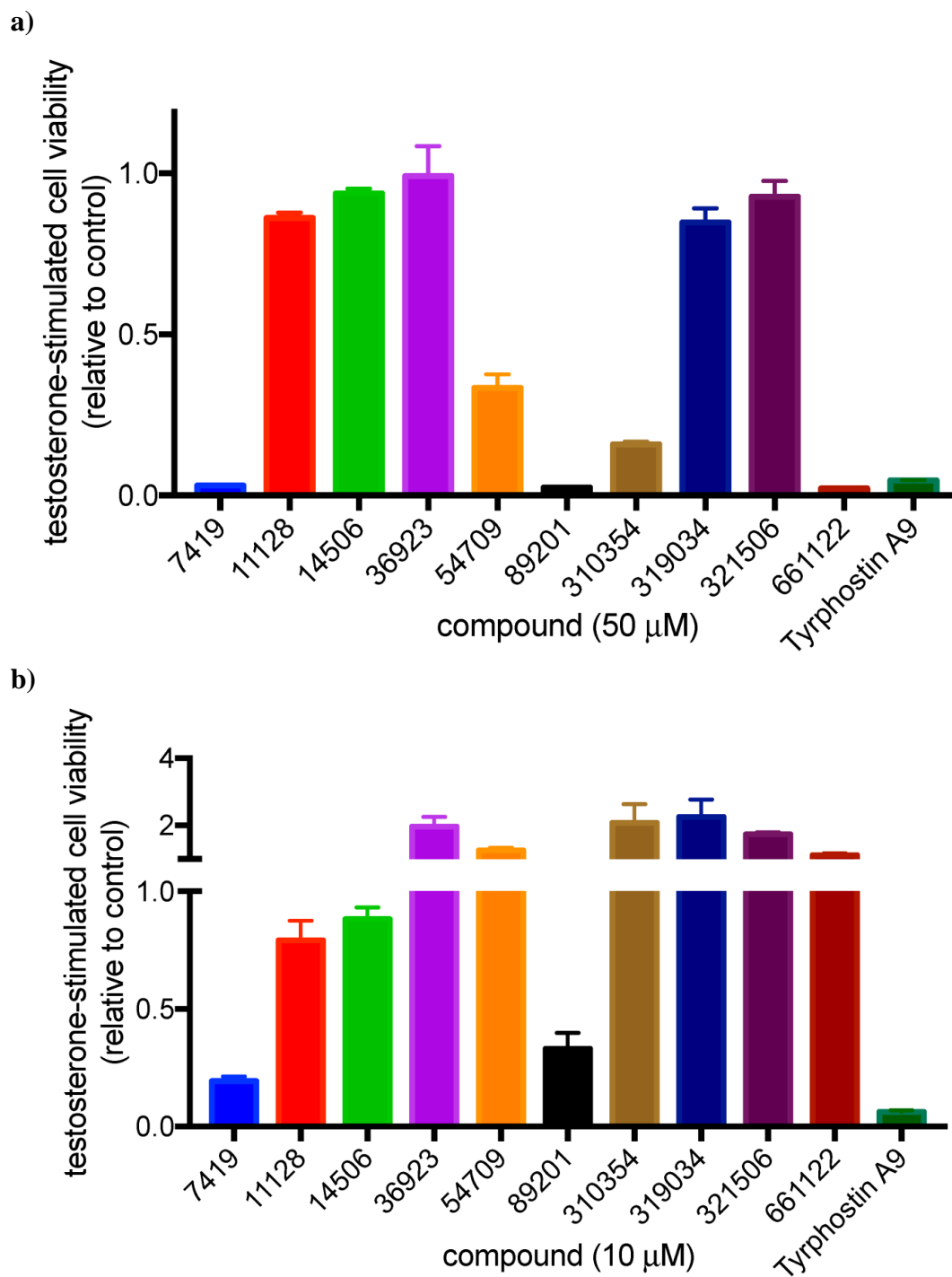


Figure 1. Measurement of TSTstimulated MCF-7 cell anti-proliferative activity of 10 compounds screened from NCI Diversity Set V at a) 50 μ M concentrations b) 10 μ M concentrations. Tyrphostin A9 is employed as a positive control.

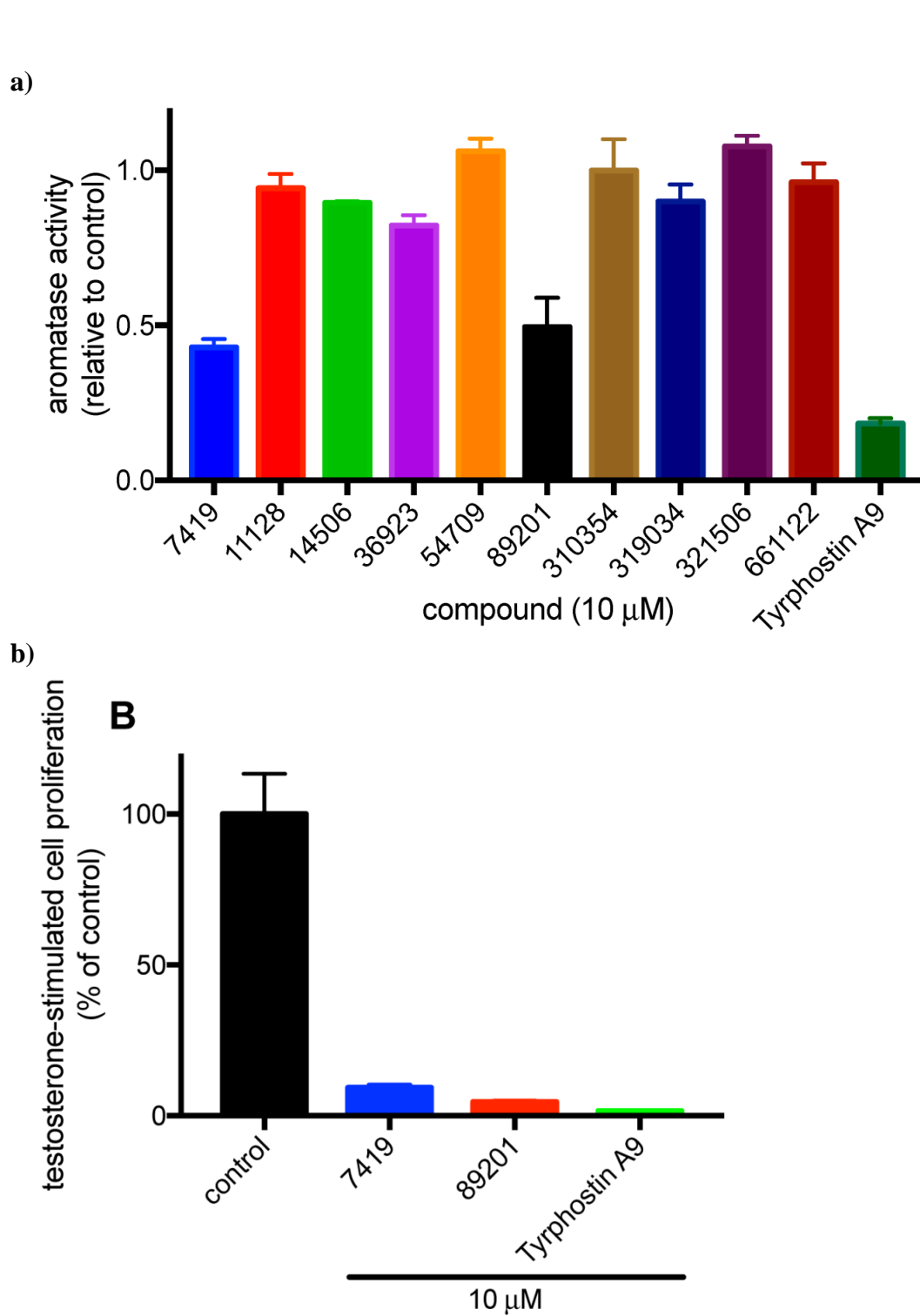


Figure 2.a) Aromatase inhibitory activity of 10 compounds screened from NCI Diversity Set V at 10 μ M concentrations. b) Anti-proliferative activity of compounds 7419, 89201 and Tyrphostin A9 at 10 μ M concentrations.

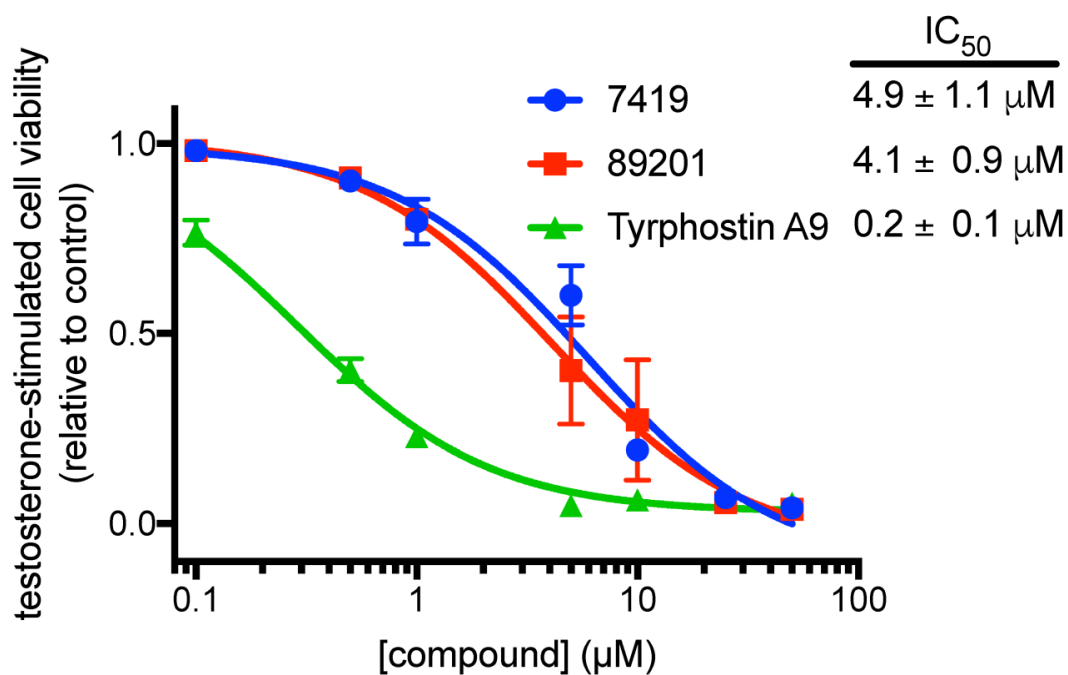
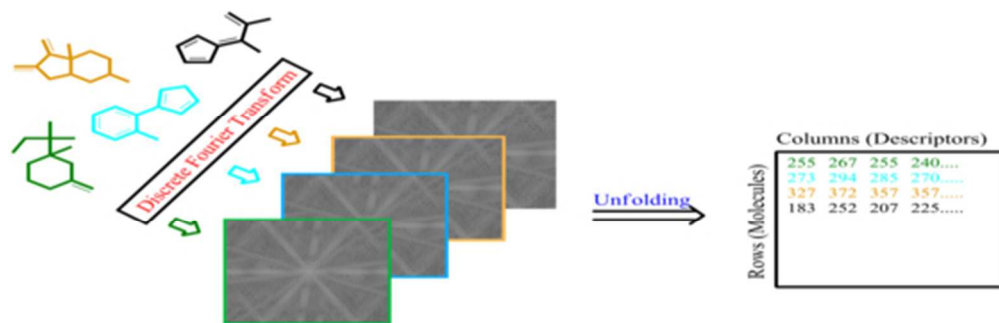


Figure 3. Dose response curve of TSTstimulated MCF-7 cell antiproliferative activity for 7419, 89201 and Tyrphostin A9. All experiments were run in triplicate in a series of three independent experiments.



TOC Graphic

148x47mm (96 x 96 DPI)