

## Articles

### Flexible Alignment of Small Molecules

Paul Labute\* and Chris Williams

Chemical Computing Group Inc., 1010 Sherbrooke Street West, 9th floor, Montreal, Quebec H3A 2R7, Canada

Miklos Feher,\* Elizabeth Sourial, and Jonathan M. Schmidt

Nanodesign Inc., Suite 300, Research Park Centre, 150 Research Lane, Guelph, Ontario N1G 4T2, Canada

Received June 20, 2000

A method is presented for flexibly aligning small molecules. The method accepts a collection of small molecules with 3D coordinates as input and computes a collection of alignments. Each alignment is given a score, which quantifies the quality of the alignment both in terms of internal strain and overlap of molecular features. The results of several computational experiments on pairs of compounds with known binding conformations are used to systematically and objectively tune the parameters for the method. The results indicate the method's utility for the elucidation of pharmacophores and comparative field analysis.

#### Introduction

Often, atomic-level details of the structures of pharmaceutically relevant receptors are not available. In such cases, 3D alignment (or superposition) of putative ligands can be used to deduce structural requirements for biological activity. For example, CoMFA<sup>1</sup> also uses a 3D molecular alignment as input. Another strategy is pharmacophore elucidation in which several ligands are aligned and a small collection of essential molecular features required for biological activity is derived from the alignment (e.g., Martin et al.<sup>2</sup>). Still another strategy is to use 3D molecular alignments to search a database for ligands that have one or more conformations that superpose well with a query molecule (e.g., Miller et al.<sup>3</sup>). Methodologies based upon 3D alignment for finding biologically active ligands generally make use of the qualitative assumption that *if two ligands have similar biological activity and bind in similar modes, then the bound conformations of the two ligands align well and inferences can be made about the nature of the receptor*. Each of the cited examples make use of this assumption (or its converse); moreover, the quality of their output depends to a large extent on the quality and significance of the calculated 3D alignments.

Many methods for calculating and/or evaluating 3D molecular alignments have been proposed<sup>3–16</sup> and have been reviewed by Lemmen and Lengauer.<sup>17</sup> The methods differ in their treatment of conformational flexibility (e.g., rigid molecules, systematic or stochastic conformational search, etc.) and the definitions of molecular similarity (e.g., ligand feature coincidence/similarity, projected receptor feature coincidence/similarity, ligand electrostatic or steric field similarity, etc.). Methods for

producing molecular alignments are generally judged according to a number of criteria:

**1. Scoring.** Given an experimental alignment of two ligands of similar activity that bind in the same mode, will the method score this alignment best (or near best) over all other alignments? This determines whether top-scoring alignments can be used to make inferences about the nature of the receptor.

**2. Completeness.** If a high-scoring molecular alignment exists, will the method find it? This determines if a molecule can be rejected (e.g., judged inactive) if no alignment is produced.

**3. Running Time.** What is the average running time per pair of molecules? This determines whether the method is suitable for database searching.

Although important for practical considerations, criterion 3 is, perhaps, of secondary importance: a fast method that is incomplete or has a poor scoring function may not be reliable enough for practical application. Criterion 2, the completeness of the method, can often be achieved with longer run times for both stochastic and systematic search methods. Most important is criterion 1, the scoring function: *if irrelevant alignments consistently score better than the relevant or "real" alignments, then the inferences drawn from the high-scoring alignments will likely be incorrect* (no matter how much computing power is allocated to the search).

Assessing the quality of an alignment scoring function is not straightforward. If a collection of molecules is ranked for similarity with respect to a query molecule (using search and scoring procedures), then a comparison with random selection methods can produce statistical measures of quality of scoring and completeness (but not scoring alone). Alternatively, X-ray crystallographic coordinates can be used to infer whether the alignment scoring function assigns the highest scores to corresponding experimentally determined align-

\* To whom correspondence should be addressed. P.L.: tel (514) 3931055; e-mail paul@chemcomp.com. M.F.: tel (519) 8239088; e-mail mfeher@nanodesign.com.

ments. This latter alternative was used in the present work, although it must be remembered that experimental and modeling errors must be taken into account when making any assessment of scoring function quality (and the degree to which bound conformations can be calculated by 3D molecular alignment).

Here we describe a method used to calculate and score 3D molecular alignments, as well as the results of several computational experiments to tune the parameters and assess the quality of the alignment scoring function.

## Methods

**Property Densities.** An isotropic (spherically symmetric) Gaussian, or normal, probability density has the functional form

$$f(x) = s^3(2\pi)^{-3/2} \exp\left\{-\frac{1}{2}|x - x_0|^2 s^2\right\}$$

where  $s^2$  is the inverse variance along each axis and  $x_0$  is the center (and mean) of the density. An atom with nucleus located at  $x_0$  with van der Waals radius  $r$  is represented by the following probability density

$$(a/r)^2(2\pi)^{-3/2} \exp\left\{-\frac{1}{2}|x - x_0|^2 (a/r)^2\right\}$$

where  $a$  and  $r$  determine the breadth of the density. It is worth noting that, in principle, two parameters are not required per density, only the ratio  $a/r$ . The reason for retaining both in the equation was that we wanted to avoid optimizing per-atom parameters, as it is difficult to assemble a sufficiently large collection of "experimentally" determined alignments. In the current work,  $r$  was fixed using force field parameters, and  $a$  was used to scale these radii simultaneously. For example, if  $a = 2$  then approximately 90% of the density will be contained within the van der Waals radius. In general, the value of  $a$  can be used to approximate molecular volumes and to fill in the gaps between atoms.<sup>18</sup> Let  $x_1, \dots, x_n$  denote the 3D positions of  $n$  atoms of a molecule in a given conformation. Let the non-negative  $w_i$  denote the degree to which atom  $i$  has some property  $P$  (for example,  $P$  could be "is aromatic" or "is donor"). We assign, to each point in space,  $x$ , a density of property  $P$  with a sum-of-Gaussians density

$$f_P(x; x_1, \dots, x_n) = \sum_{i=1}^n \frac{w_i}{n} \left( \frac{a^2}{2\pi r_i^2} \right)^{3/2} \exp\left\{-\frac{a^2}{2r_i^2} |x - x_i|^2\right\}$$

This density is called the  $P$ -density of the conformation. For example, if  $P$  is the property "is aromatic" so that  $w_i$  is 1 when atom  $i$  is aromatic and 0 otherwise, then the function  $f_{\text{aro}}$  is the *aromatic density* of the conformation. In the present work we shall consider the following atomic properties:

**1. Volume.** A volume feature is just the presence of an atom. In the  $P$ -density, this corresponds to using a value of 1 for each  $w_i$ .

**2. Aromatic.** The aromatic  $P$ -density is constructed by assigning  $w_i = 1$  if atom  $i$  is aromatic and 0 otherwise. The Hückel  $4n + 2$  rule was used to assign aromaticity to  $sp^2$  rings with no exocyclic double bonds.

**3. Donor.** The donor  $P$ -density was assigned by setting  $w_i = 1$  if atom  $i$  was of type "Donor", "Polar", or "Basic" under a pharmacophore atom typing scheme based upon the PATTY rules.<sup>19</sup>

**4. Acceptor.** The acceptor  $P$ -density was assigned by setting  $w_i = 1$  if atom  $i$  was of type "Acceptor", "Polar", or "Acidic" under a pharmacophore atom typing scheme based upon the PATTY rules.

**5. Hydrophobe.** The hydrophobic  $P$ -density was assigned by setting  $w_i = 1$  if atom  $i$  was of type "Hydrophobe" under a pharmacophore atom typing scheme based upon the PATTY rules.

**6. logP (Octanol/Water).** Each  $w_i$  is the (normalized) atomic contribution to logP as calculated using the Wildman and Crippen SlogP method<sup>20</sup> which was parametrized with atomic contributions in mind.

**7. Molar Refractivity.** Each  $w_i$  is the (normalized) atomic contribution to molar refractivity as calculated by the Wildman and Crippen SMR method<sup>20</sup> which was parametrized with atomic contributions in mind.

**8. Surface Exposure.** Each  $w_i$  is the percentage of the van der Waals surface area of atom  $i$  that is exposed (not contained in another atom).

**Similarity Measure.** Given two  $P$ -densities for two given molecules or conformations, the *overlap* of the two densities is, itself, a sum-of-Gaussians density in the interatomic distances:

$$F_P = \sum_{i=1}^n \sum_{j=1}^{n'} \frac{w_i w_j}{n n'} \left( \frac{a^2}{2\pi(r_i^2 + r_j^2)} \right)^{3/2} \exp\left[-\frac{a^2}{2} \frac{|x_i - x_j|^2}{r_i^2 + r_j^2}\right]$$

This overlap formulation generalizes to more than one property, or feature, through a weighted summation of the individual overlap equations. For example, given three properties defined by atomic property weights,  $u$ ,  $v$ ,  $w$ , we define the similarity of two molecular conformations to be

$$F = \sum_{i=1}^n \sum_{j=1}^{n'} \frac{C_u u_i u_j + C_v v_i v_j + C_w w_i w_j}{n n'} \left( \frac{a^2}{2\pi(r_i^2 + r_j^2)} \right)^{3/2} \times \exp\left[-\frac{a^2}{2} \frac{|x_i - x_j|^2}{r_i^2 + r_j^2}\right]$$

where  $C_u$ ,  $C_v$ , and  $C_w$  are positive weights intended to emphasize or de-emphasize particular  $P$ -densities. In this way, any number of features can be included into the similarity calculation without additional computational complexity: the pairwise weights need be calculated only once. This similarity measure also generalizes to more than two molecules through summation of the pairwise similarity function over all pairs of molecules in a collection. The SEAL method of Kearsley<sup>16</sup> is based on a similar approach with  $u$  representing atomic partial charge and  $v$  representing van der Waals volume. Along with the Gaussian exponent,  $a$ , the individual weights are tunable parameters that can be set according to the procedure detailed later in this paper.

**Search Procedure.** To simultaneously search the conformation space of each molecule and the alignment space of the collection for optimal alignments, we used a modified RIPS<sup>21</sup> procedure, summarized as follows:

0. [Initialize]. Set the values of the adjustable parameters: (a) the perturbation limit  $p$  used in step 1; (b) the "temperature"  $T$  used for refinement in step 2; (c) the Gaussian variance scale  $a$  used in the scoring function; (d)  $C_P$ , the weight used for property  $P$  in the similarity function (e.g., if four properties are used, such as donor, acceptor, volume, and aromatic, then four weights are used); (e) the RMSD threshold for duplicate testing in step 3; (f) the failure limit  $L$  used in step 3; (g) the iteration limit used in step 3; and (h) the average strain energy limit used in step 4.

1. [Perturb]. Set all rotatable bonds (nonterminal and nonring bonds) to random dihedral angles. Add a random number in the range  $[-p/2, p/2]$  to all atomic coordinates. Randomly orient all molecules by choosing three atoms randomly from each molecule, and superpose. (In the results that follow, a value of  $p = 1.0$  Å was used in order to ensure adequate sampling of ring conformations.)

2. [Optimize]. Minimize the objective function  $-kT \log F + U$  with respect to the coordinates of all of the atoms. Here,  $F$

is the similarity function, and  $U$  is the average potential energy of the molecules. (In the results that follow, a value of  $T = 30\,000$  was used.)

3. [Compare]. If the new configuration has not been seen before (RMSD greater than some threshold), then set  $k = 0$ , otherwise set  $k = k + 1$ . If  $k$  is greater than some predefined amount,  $L$ , then terminate the search and go to step 4; otherwise return to step 1. (In the results that follow, a heavy atom RMSD of 0.2 Å was used for duplicate configuration detection, and a failure threshold value of  $L = 1000$  along with an iteration limit of 1000 was used. This combination of parameters had the effect that, in all cases, 1000 attempts were made to generate new alignments.)

4. [Filter]. Prune the list of configurations by removing all configurations in which the average potential energy (of the alignment) is greater than the minimum observed average potential energy plus some predefined threshold. (In general, raw X-ray structures have large amounts of strain in force fields and can be over a 100 kcal/mol higher in energy than the nearest local minimum. To account for this effect, an energy cutoff of 200 kcal/mol was used.)

The termination criteria can be interpreted as follows. Upon termination, there were  $L$  consecutive attempts to generate a new configuration, and each has failed. By way of analogy to coin tossing, we can estimate the probability that there exists a configuration not yet seen. Using the coin tossing analogy, a biased coin has been tossed  $L$  times, and each time "heads" was observed. The Bayes estimate for the probability of observing "heads" is  $(\text{number of heads} + 1)/(\text{number of tosses} + 2)$ . If  $L$  "heads" are observed in  $L$  tosses, then the probability of observing "heads" is  $(L + 1)/(L + 2)$ ; thus the probability of observing "tails" (or, analogously, a new configuration) is  $1/(L + 2)$ . At  $L = 18$  this probability is 5%.

**Materials and Software.** The foregoing methods were implemented in the SVL programming language of Chemical Computing Group Inc.'s Molecular Operating Environment (MOE) version 1999.05.<sup>22</sup> The calculations were performed on a 195 MHz Silicon Graphics Octane running IRIX 6.2 as well as 200 and 350 MHz Intel Pentium II processors running Windows NT. Nonlinear optimization was carried out with the MOE Truncated Newton optimizer preceded by two steps of Steepest Descent<sup>23</sup> and terminated when the RMS gradient fell below 0.001. The MOE implementation of MMFF94<sup>24</sup> force field was used to measure the internal strain of each molecule. Chirality was preserved using signed volume restraints on all chiral centers. Pharmacophore atom type assignment was performed using the MOE implementation of the Daylight SMARTS pattern matching language. The determination of the RMSD between structures used the MOE Superpose functionality, which calculates an optimal global superposition by minimizing a weighted least squares error function. The CPU time required for the calculation of a single alignment varies according to the total number of atoms in all of the molecules. In the case of the DHF/DLS, a system with 108 atoms in total, 32 s were required (on average) to produce an alignment on a Hewlett-Packard C3600 Visualize workstation with an HP9000 processor and 42 s on a PC with an 800 MHz Intel Pentium III processor running under Windows NT. Another way of looking at the speed is how long it takes to generate the alignment closest to the crystal structure. For the examples described below, we obtained the following timings: DHF/DLS 2.6 min, CPX/CPA 2.3 min, ERT/ERR 8.5 min, and ERE/ERR 3.1 min. These recorded times were for a sample run and will obviously vary due to the random nature of the search algorithm. However, it has been our experience with the overlay of a number of drug-size molecules that the best alignments are usually found within a few minutes.

**Parameter Optimization.** In this study, Gaussian distance dependence,  $a$ , as well as the weights of the main similarity properties (volume, aromatic, donor, and acceptor) were optimized. The "temperature" parameter  $T$ , used to balance the competition between overlap and the underlying force field, was not optimized with the other parameters, and the value of  $T = 30\,000$  K was used in the calculations. This

was determined independently by examining the output of a collection of alignment runs on different molecules (not all of which aligned well). From these, a value of  $T$  was selected that produced approximately 1% frequency of occurrence of strain energies greater than 20 kcal/mol. The frequency of occurrence of highly strained alignments did not appear sensitive to the precise value of  $T$ , only its order of magnitude. Therefore, we did not feel that it was necessary to vary  $T$  with the other parameters.

The initial parametrization tests were carried out for the heterocyclic ring portion of dihydrofolic acid and methotrexate. It was found that the information contained in volume and aromatic weights is largely confounded, as no well-defined optimum could be identified in the separate optimization of these two factors. This is not surprising as they both describe steric effects. In a similar manner, acceptor and donor weights could not be optimized separately as they are both responsible for the "electronic" effects. Hence in this work these were treated as compound steric ( $w_s$ ) and electronic ( $w_e$ ) interactions, respectively, with the weight of the constituent factors kept equal. These two molecule fragments were also used to establish the importance of additional physicochemical feature functions in the fit, including the role of logP, hydrophobicity, and molar refractivity, as well as the effect of using the exposed van der Waals surface area instead of volume.

The final parameters were optimized using four ligand pairs: (1) the entire methotrexate and dihydrofolic acid molecules, (2) L-benzylsuccinate and glycyl-L-tyrosine, (3) raloxifene and 4-OH-tamoxifen, and (4) estradiol and raloxifene. The final orientations from the flexible alignment program were compared to the reference superpositions retrieved from the Protein Data Bank,<sup>25</sup> which were generated using the following approach. First, the common residues were identified in the protein pair. In the case of the molecules estradiol, raloxifene, and 4-OH-tamoxifen, the region called 'helix 12' was also ignored in the alignment, as its conformation is known to be substantially different for different bound ligands.<sup>26</sup> The  $\alpha$ -carbons of these residues were then aligned using rigid body superposition (RMSD < 0.35 Å in all cases). This process led to the alignment of the ligands, the absolute coordinates of which were extracted and used as a reference. This will be referred to as the crystal alignment in this work.

There are different possibilities for evaluating the obtained solutions. In a number of previous studies, several high-scoring solutions were considered simultaneously, and the one closest to the X-ray structure was selected.<sup>3,6,8</sup> In contrast, in this study only the solution with the highest score was considered. In cases when experimental structural information is unavailable, this is, in fact, the only objective selection. The validity of the predicted flexible alignment solutions was verified by calculating the heavy atom root-mean-square distance between the predicted alignment and the reference crystal alignment. Although the best parameter settings in all four cases were found to be similar (see Table 1), the final parameter values were determined by calculating the lowest RMSD mean over the four ligand pairs as defined by Klebe et al.<sup>4</sup>

$$\text{RMS}_{\text{mean}} = \sqrt{\sum_i \text{rms}_i^2 / n}$$

where  $n$  is the number of considered ligand pairs ( $n = 4$ ).

Several issues need to be considered when the quality of the fits is judged using the RMS deviations from the crystal alignment. First, these RMS distances were consistently calculated for all heavy atoms of the studied ligand pair. The consideration of only selected fit centers (such as in Klebe<sup>4</sup>) or the consideration of only one of the molecules ignoring the other and their relative positions (such as in Nissink et al.<sup>5</sup>) is likely to lead to lower RMS values without a real improvement in the alignment. Second, equilibrium geometries using a force field are substantially different from cocrystallized geometries in the crystal. The RMS distance between the two gives a lower bound for the achievable accuracy in the alignment (see Table 1). Third, as we can see in all the



**Table 1.** Results of the Optimization of Adjustable Parameters in the Flexible Alignment of Ligands, Pairwise Binding to the Same Protein<sup>a</sup>

	lowest RMSD (Å)	$a/w_s/w_e$ parameters at lowest RMSD	RMSD (Å) at final parameters ( $a/w_s/w_e$ of 2.5/3/1)	RMSD force field error <sup>b</sup>
1DHF/1DLS (heterocycles only)	0.34	2.5/7/1	0.36	0.09
1DHF/1DLS	1.4	2.5/3/1	1.4	0.68
1CBX/3CPA	1.1	2.5/1/1	1.3	0.65
3ERT/1ERR	1.6	2.5/3/1	1.6	0.73
1ERE/1ERR <sup>c</sup>	1.3	2.5/0.25/1	1.4 (0.65 biased)	0.78

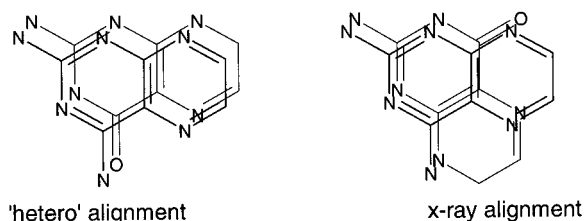
<sup>a</sup> The parameters were  $a$ ,  $w_s$  (steric weight, comprising equal aromatic and volume weights),  $w_e$  (electrostatic weight, comprising equal acceptor and donor weights). The root-mean-square deviations are given between the alignment in the crystal and the first flexible alignment solution. The final parameters,  $a/w_s/w_e = 2.5/3/1$ , were obtained by minimizing the mean RMS error (as calculated from the RMS errors for the four ligand pairs). The best mean RMS error was 1.39 Å. <sup>b</sup> The RMS distance between the cocrystallized ligand conformations before and after geometry optimization with the MMFF94 force field. <sup>c</sup> This error was calculated by excluding the flexible side chain of raloxifene from the fit. When the side chain was considered, the RMS deviation at the optimized parameter setting was 2.2 Å, whereas its lowest value was 2.1 Å. The biased value was obtained by excluding the pair of oxygen atoms from the fit, which are known from the crystal structure to be displaced considerably.

examples below, the cocrystallized ligands are often visibly misaligned, the distance between the corresponding pharmacophoric groups in the pair being as much as 0.5 Å. This arises for spacious binding pockets where the only constraint in the position of the ligand is that the corresponding groups interact with the same centers in the receptor. This 'misalignment' cannot be reproduced in computed overlays, as these attempt to get all centers as close to each other as possible. Fourth, calculated alignments consider the overlap of all corresponding centers, unless prior binding information is taken into account. In contrast, only some of the centers are overlaid in real receptors, as others may not be involved in interactions. This may introduce significant differences between the calculated and experimental alignments as the corresponding noninteracting centers may be several angstroms apart. Fifth, one of the molecules from the pair may contain large regions that have no equivalent volume in the other molecule. In these cases, the alignment program does not have sufficient information on these volumes, and hence it is no surprise that these differ in their orientations from that in the crystal. An example for this is shown in the estradiol/raloxifene alignment. In such cases, it is only meaningful to consider the alignment of the common parts of the volume. Finally, the alignment of the receptor structures also has an error (of the order of about 0.4 Å). These effects all lead to distortions with respect to the experimental alignment and increase the RMSD of the experimental and the aligned ligands.

## Results and Discussion

**Alignments.** The results for the best flexible overlays of the studied ligand pairs are given in Table 1. The Gaussian distance factor,  $a$ , was changed systematically from 1 to 6, with the ratio  $w_s/w_e$  being varied simultaneously between 0.25 and 10. Some values outside these ranges were also tested. In principle, a direct numerical optimization of these parameters is feasible; however, it was felt that such a procedure was not warranted in light of the small size of the training set. For the same reason, it is not possible to attach much meaning to the resulting values of the optimized parameters other than the obvious "these are the values that produce the best results". For all these parameter settings, the RMS<sub>mean</sub> for the four ligand pairs were calculated. The optimum was found when  $a$  was set to 2.5 and the  $w_s/w_e$  ratio had a value of 3 (this set of parameters will be denoted as  $a/w_s/w_e$  2.5/3/1). The RMS deviations at this final parameter setting are given in Table 1. As described below, additional physicochemical properties had little effect on the quality of alignments. Next, the results are described for the studied molecule pairs.

The initial parameter optimizations were carried out on the heterocyclic rings of methotrexate and dihydro-



**Figure 1.** Schematic representation of the possible overlays of the heterocyclic rings in methotrexate and dihydrofolic acid. The first alignment is intuitive but incorrect, while the second overlay corresponds to the experimental crystal structure.

folic acid. This pair was chosen because it is probably the most common test example for alignment programs. Also, the ease of calculation on this system allowed a very detailed analysis of the effect of different parameters.

One can imagine two possible overlays of the two heterocycles: an intuitive one with the heterocycles on top of each other and the alignment based on the X-ray structure, originally suggested by Bolin et al.<sup>27</sup> These are shown schematically in Figure 1 and will be referred to as the 'hetero' and the 'X-ray' alignments, respectively. In judging the quality of the alignments, the actual X-ray coordinates were also used (PDB codes of the ligands cocrystallized with the dihydrofolate reductase enzyme are 1DLS and 1DHF). As these two molecular fragments are quite rigid, once the proper alignment was found the RMS deviation from the crystal structure hardly changed (to within 0.03 Å). Hence the separation of the first and the second alignment was used to judge the quality of the fits. This function had a clear optimum at  $a = 2.5$ . It was more difficult to establish the steric and electrostatic weights, as a high weight on the acceptor-donor properties invariably produced the correct alignment, whereas a high steric weight led to the 'hetero' alignment. Nonetheless, the weights selected for the four ligand pairs,  $w_s/w_e = 3/1$  (vide infra), sufficiently separated the first two solutions, the second being the 'hetero' alignment. Under these conditions, the predicted and the X-ray alignment agreed well (RMSD of 0.35 Å).

The effect of including other physicochemical properties into the  $P$ -density calculation was investigated for this ligand fragment pair at the final parameters (i.e.,  $a/w_s/w_e = 2.5/3/1$ ). It was found that the inclusion of either hydrophobicity or logP did not have any effect

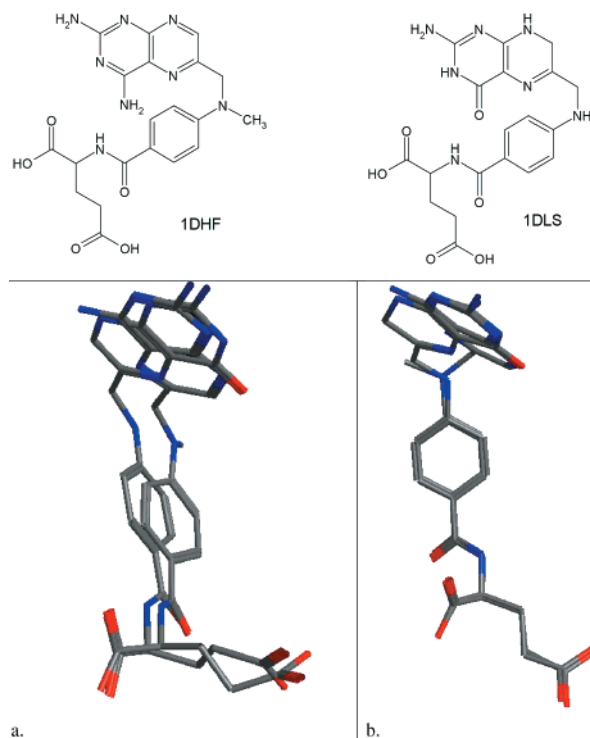
on the alignments if their respective weights were set to 1. On increasing the weight of these two effects further, the fits deteriorated. On including molar refractivity, it had no effect on the fits up to a weight of  $\sim 0.2$  when the quality of fits started deteriorating. Hence we can conclude that inclusion of either of these properties does not lead to any improvement in the alignments. These findings are similar to those of Klebe et al.,<sup>4</sup> in which the inclusion of hydrophobicity and refractivity into the SEAL alignment function did not significantly improve the overlay of ligands. The role of hydrophobicity was also tested at a later stage for two ligand pairs, 1CBX-3CPA and 1ERE-1ERR. At hydrophobic weights ( $w_h$ ) set to 1 and 5 (at  $a/w_s/w_e = 2.5/3/1$ ), no improvement in the RMSD of the best scoring solutions was observed.

The effect of using exposed surface area instead of volume in the  $P$ -density was similarly tested using the heterocyclic portions of the methotrexate and dihydrofolic acid fragment pair. The ratio of the molecular surface and aromatic weights had to be kept constant because of difficulties in separately optimizing these factors. The lowest RMS deviation from the crystal was found again at  $a = 2.5$ , with the optimum value of  $w_s/w_e$  being 0.25/1 (RMSD of 0.34 Å). This RMSD, however, is not significantly different from the optimum found when volume was used (RMSD of 0.36 Å). Furthermore, although at most parameter settings the crystal alignment was obtained when volume was applied, the use of solvent accessible surface area favored the hetero alignment. Hence only the volume terms were considered for the remainder of this study.

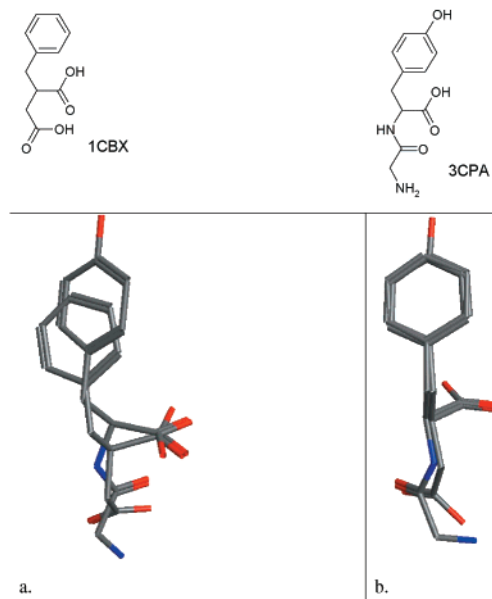
The best flexible overlay of the entire methotrexate and dihydrofolic acid molecules is displayed in Figure 2. This alignment is consistent with the experimental overlay (RMSD of 1.4 Å): the corresponding pharmacophoric groups are all aligned, and the relative position of the fused heterocycles (X-ray alignment) as well as their angle to the rest of the molecule are correct in the overlays.

It is difficult to compare our observed RMS distance of 1.4 Å to other works that investigated the alignment of the same ligand pair. Although the RMSD values from the alignments by Klebe et al.<sup>4</sup> (RMSD of 0.963 Å) and Nissink et al.<sup>5</sup> (RMSD of 0.87 Å) appear better, these were obtained by using only selected fit centers and not all heavy atom coordinates. Similarly, it should be borne in mind that, with the FlexS method of Lemmen et al.,<sup>6</sup> the quoted RMSD's of 1.39 and 1.68 Å (depending on the order of the ligands) were both obtained as the fourth best scoring solutions, whereas we only considered our best scoring alignment solution.

Next, the overlay of L-benzylsuccinate and glycyl-L-tyrosine was examined. In this example, a number of potential pharmacophoric centers are present. The ligands, cocrystallized with the carboxypeptidase-A enzyme (PDB reference codes 1CBX and 3CPA, respectively), as well as the flexible alignment solution, are displayed in Figure 3. It can be seen in Figure 3a that the pharmacophoric centers in the crystal are not well aligned, even though the experimental resolution in both cases is 2 Å. This arises as the phenyl groups of the two molecules are located in a large pocket of the enzyme with no nearby residues in close proximity.

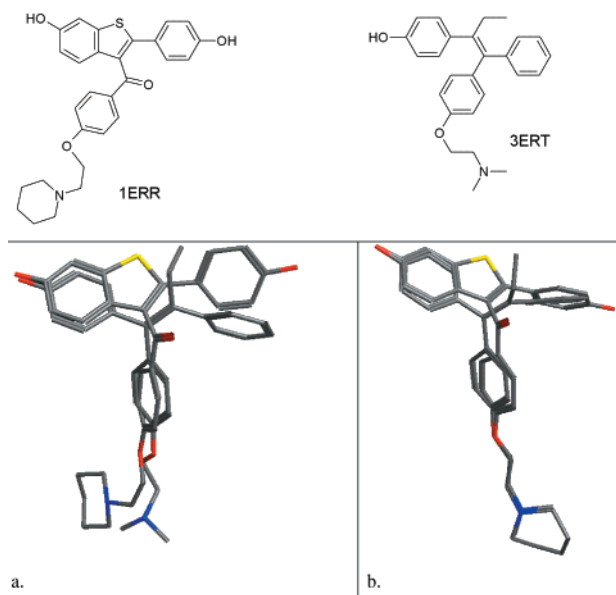


**Figure 2.** Overlay of methotrexate and dihydrofolic acid: (a) alignment in the crystal; (b) best calculated flexible alignment, obtained at the final set of parameters optimized for all ligand pairs (RMSD of 1.4 Å).



**Figure 3.** Overlay of L-benzylsuccinate and glycyl-L-tyrosine: (a) alignment in the crystal; (b) calculated flexible alignment at the final set of parameters optimized for all ligand pairs (RMSD of 1.3 Å).

Hence, the alignment of the aromatic rings in the pocket is not required for activity. According to the crystal structure, the groups establishing hydrogen bonds are as follows: in 3CPA, the two carboxylic oxygens and the primary amino group, but not the oxo and secondary amino groups; in 1CBX, all four carboxylic oxygens. These groups can still interact with the same receptor atoms, despite the obvious displacement of the corresponding centers. In contrast, in the flexible alignment these corresponding centers are forced together, with

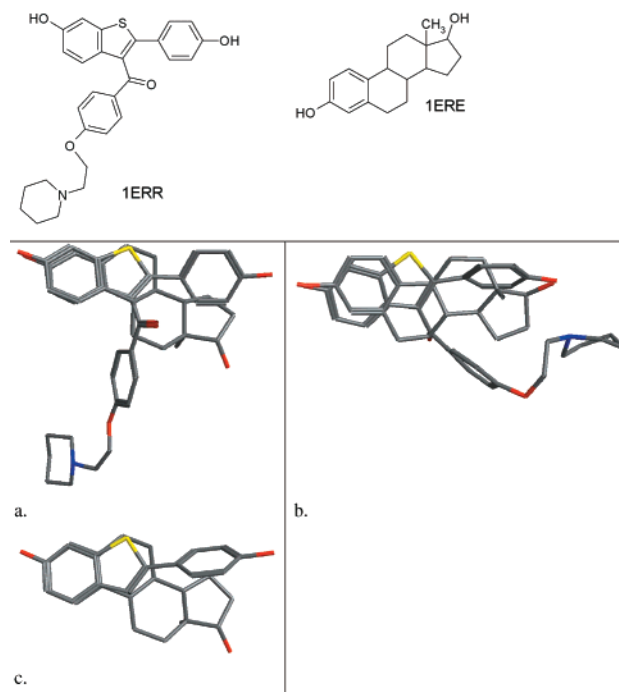


**Figure 4.** Overlay of 4-OH-tamoxifen and raloxifene: (a) alignment in the crystal; (b) best calculated flexible alignment, obtained at the final set of parameters optimized for all ligand pairs (RMSD of 1.6 Å).

the major features of the crystal overlay still preserved (see Figure 3b, RMSD of 1.3 Å). These results compare favorably with those obtained by Lemmen et al.<sup>6</sup> using the FlexS program (RMSD of 1.42 Å), where it was obtained as the third best solution.

In the next example, the overlay of two flexible molecules is presented, raloxifene and 4-OH tamoxifen (PDB codes of ligands cocrystallized with the estrogen receptor are 1ERR and 3ERT). As shown in Figure 4 the predicted alignment is in good qualitative agreement with the crystal alignment (RMSD of 1.6 Å). It can be observed that the long side chains are primarily responsible for the relatively high RMSD. For example, the corresponding  $sp^3$  oxygen and nitrogen atoms in the side chains are 0.87 and 1.97 Å apart, respectively, in the superposed crystal structures, whereas the corresponding distances in the flexible alignment are 0.12 and 0.18 Å, respectively. The misalignment of the side chains in the ligand–receptor complex arises as the end of the side chain protrudes from the active site and only interacts with the solvent, with its position not being well determined. Hence, the results of the flexible alignment cannot be improved in this case without the prior knowledge of the position of the side chain.

Finally, two spatially dissimilar compounds, estradiol and raloxifene, were overlaid (PDB reference codes of the ligands cocrystallized with the estrogen receptor are 1ERE and 1ERR, respectively). The alignment in the crystal and at the final set of parameters is shown in Figure 5. As can be seen in this figure, the predicted overlay solution appears to be somewhat different from the crystal alignment (RMSD of 2.1 Å). The obvious reason for this large discrepancy lies in the fact that raloxifene has a long side chain, which has no equivalent in estradiol. The major difference between the crystal and the aligned structures is the relative angle at which the side chain points. However, it is no surprise that the orientation of this side chain differs from that in the crystal, since the alignment program has no



**Figure 5.** Overlay of estradiol and raloxifene: (a) alignment in the crystal; (b) calculated flexible alignment, obtained at the final set of parameters optimized for all ligand pairs (RMSD without side chain of 1.4 Å); (c) best calculated flexible alignment, obtained by excluding those oxygens from the alignment that are displaced in the crystal (RMSD without side chain of 0.65 Å).

information on how to align this part of the molecule. The issue of nonoverlapping fragments has been described in detail by Lemmen et al.,<sup>6</sup> where these were identified using the common Connolly surface, and the geometry of the fragment outside this volume was set to generic values from crystal databases. Another procedure to detect this situation was described by Robinson et al.,<sup>7</sup> in which the outlying atoms were identified based on distances of each atom of the larger molecule to the closest atom in the smaller one. The solution in the present study was simply to exclude the side chain from the calculation of the RMSD value, which led to a considerable improvement (RMSD of 1.4 Å). On comparing the alignments in Figure 5a and 5b, one further difference can be observed. Whereas both hydroxyl groups of estradiol are aligned with both hydroxyls of raloxifene in the calculated overlay, one set of hydroxyl groups is totally misaligned in the crystal (the distance between the corresponding oxygens is 4.7 Å). When the hydroxyl oxygens were forced together by the alignment program, the best alignment of the rings was reached by also flipping estradiol along the axis connecting the hydroxyl groups. This solution was found at all parameter settings. Unfortunately, this situation cannot be easily remedied without prior knowledge about the binding of this particular set of molecules. Once such information is available, however, we can simply exclude from the fit the oxygens that are misaligned in the crystal. Figure 5c displays such a fit, achieved using the same set of parameters. This change indeed improved the alignment substantially (RMSD excluding the side chain is 0.65 Å). It must be noted that it is possible that our original alignment, shown



in Figure 5b, may occur experimentally. A second binding mode for estradiol has been postulated from docking experiments,<sup>28</sup> in which estradiol is rotated by about 180° in the plane to reverse the role of the two oxygens. Similarly, a second binding orientation of raloxifene was identified in flexible docking experiments,<sup>29</sup> in which raloxifene needs to be rotated by about 180° out of plane so that the tail section remains fixed. The combination of these two binding modes would result in the alignment in Figure 5b.

**Comparison to Other Methods.** As many methodological details on other methods in the literature are unavailable, it is not possible to draw complete comparisons. In this section, we will make a qualitative comparison of our method with those reported previously.

**Completeness.** The cited rigid alignment approaches<sup>3,5,9</sup> rely on external methods to generate conformations; hence, their completeness regarding variation in conformation cannot be easily evaluated. In addition, the predetermined conformations in these rigid alignment methods are local energy minima, obtained by minimizing the potential energy of the isolated molecule. In contrast, the method in this work optimizes the objective function made up of the similarity score and the internal energy. Thus, in the resulting overlays, the conformations will be slightly more energetic but also more similar to each other than would be possible in the rigid methods.

Methods that rely on systematic searches exhibit good completeness only for smaller search spaces and become impractical for larger and more flexible molecules, as well as for more than two structures. This arises due to the combinatorial explosion in numbers of alignments and conformations. For this reason, the present method and most of the cited nonrigid methods use stochastic search techniques, which have a high probability of being complete (more samples lead to higher probabilities of completeness). Some methods<sup>4,6,8</sup> use biased conformation generation, usually based on dihedral angles predominant in crystal structures. By definition these only produce a small fraction of the available conformations, and therefore are a priori incomplete (especially for large flexible rings). Of the cited methods, only reference 11 incorporates an unbiased conformational search of dihedral angles simultaneously with the alignment search. Unfortunately, since the search is based on a genetic algorithm with no force field, unrealistic conformations may dominate the population and lead to infeasible alignments scoring best. The present method is unique in that it uses a completely unbiased all-atom force field and a conformational search that is complete with high probability even for alignment problems with many large flexible molecules.

**Scoring.** As was pointed out earlier, the alignment scoring function is, perhaps, the most important criterion; however, methodological differences in the literature make comparisons difficult. None of the methods in refs 3, 11, or 16 showed that "experimentally" determined alignments score best over alternatives. In refs 4,5,8 only selected centers were used when calculating RMSD, whereas the present work considered all heavy atoms in the RMSD calculation. In ref 6 success was judged based on the "correct" alignment's appear-

ance in the top few high-scoring solutions, whereas we always considered only the top-scoring solution. Clearly, ours is an unbiased and stringent criterion for judging the quality of the scoring function, and yet the solutions still appear to be sufficiently similar to the crystal structure.

**Speed.** In methods that rely on an external conformational search,<sup>3,5,9,16</sup> the run time will be a function of the number of applied conformations. Intuitively, the efficiency of these methods will deteriorate when presented with molecules with many conformational degrees of freedom. The situation is even worse in the overlay of more than two such molecules due to the combinatorial explosion in possible solutions. The methods that rely on biased conformation generation<sup>4,6,8</sup> have traded completeness for speed (especially for large flexible rings).

It is generally difficult to compare the speeds of methods numerically, unless identical molecules are run on identical machines. Nevertheless, we can compare average run times to assess the suitability of the method for different applications. Average run times of about 2 min/molecule were recorded for FlexS<sup>6</sup> using a SUN-Ultra-30 workstation with 296 MHz clock speed. Using the alignment based on a genetic algorithm,<sup>11</sup> run times below 10 min were reported using an SGI Indigo II machine with an R4000 processor. These speeds are similar to those in this work (see Materials and Software section). Future work will include further improvements to speed up the process. In addition, our method is well suited to alignment problems involving the overlay of two or more molecules (only ref 9 supports more than two molecules) and the overlay of flexible molecules (only ref 11 searches larger ring systems directly).

## Conclusions

We have presented a method for aligning a collection of small molecules in a flexible manner. The method produces a collection of alignments along with a score for each alignment based upon the internal energy of the molecules and a similarity score defined by an overlap of Gaussian feature densities. In principle, any numeric feature, or several, can be used in the similarity score (e.g., hydrophobicity, donor, acceptor, logP contribution, etc.) without additional computational complexity. Feature weights can be used to emphasize certain features over others. The determination of a few key feature weights and a tunable Gaussian exponent parameter was the focus of computational experiments. In particular, it was found that the volume, aromatic, donor, and acceptor feature densities were most important while other features such as logP, molar refractivity, hydrophobicity, and exposed surface area did not significantly improve alignments. It was determined that a relative weight of the steric features (volume and aromaticity) to the electronic features (donor and acceptor) of 3 to 1 and a Gaussian exponent parameter of 2.5 produce alignments which correspond well to experimental results.

The results presented confirm that by using the selected parameters the method can reproduce well the crystal alignment with some key limitations in mind. If molecules occupying greatly different volumes are

aligned, the conformation of the volume outside the common one will not be unambiguously determined in the process. Furthermore, without the crystal structure it is generally unknown which of the pharmacophoric groups of the two molecules overlay experimentally and which process will fit them all. Although these issues negatively impact the RMSD from the crystal alignment, this approach is objective, as it requires neither predefined fit-centers nor preorientation of molecules. The results indicate that the described flexible superposition method can lead to meaningful and unbiased alignments.

## References

- (1) Cramer, R. D.; Patterson, D. E.; Bunce, J. D. Comparative Molecular Field Analysis (CoMFA). 1. Effect of Shape on Binding of Steroids to Carrier Proteins. *J. Am. Chem. Soc.* **1988**, *110*, 5959–5967.
- (2) Martin, Y. C.; Bures, M. G.; Danaher, E. A.; Delazzer, J.; Lico, I.; Pavlik, P. A. A Fast New Approach to Pharmacophore Mapping and its Application to Dopaminergic and Benzodiazepine Agonists. *J. Comput.-Aided Mol. Des.* **1993**, *7*, 83–102.
- (3) Miller, M. D.; Sheridan, R. P.; Kearsley, S. K. SQ: A Program for Rapidly Producing Pharmacophorically Relevant Molecular Superpositions. *J. Med. Chem.* **1999**, *42*, 1505–1514.
- (4) Klebe, G.; Mietzner, T.; Weber, F. Different Approaches Toward an Automatic Structural Alignment of Drug Molecules: Applications to Sterol Mimics, Thrombin and Thermolysin Inhibitors. *J. Comput.-Aid. Mol. Des.* **1994**, *8*, 751–778.
- (5) Nissink, J. W. M.; Verdonk, M. L.; Kroon, J.; Mietzner, T.; Klebe, G. Superposition of Molecules: Electron Density Fitting by Application of Fourier Transforms. *J. Comput. Chem.* **1997**, *18*, 638–644.
- (6) Lemmen, C.; Lengauer, T.; Klebe, G. FlexS: A Method for Fast Flexible Ligand Superposition. *J. Med. Chem.* **1998**, *41*, 4502–4520.
- (7) Robinson, D. D.; Lyne, P. D.; Richards, W. G. Alignment of 3D-Structures by the Method of 2D Projections. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 594–600.
- (8) Klebe, G.; Mietzner, T.; Weber, F. Methodological Developments and Strategies for a Fast Flexible Superposition of Drug-Sized Molecules. *J. Comput.-Aided Mol. Des.* **1999**, *13*, 35–49.
- (9) Feher, M.; Schmidt, J. Multiple Flexible Alignment with SEAL: A Study of Molecules Acting on the Colchicine Binding Site. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 495–502.
- (10) Lemmen, C.; Hiller, C.; Lengauer, T. RigFit: A New Approach to Superimposing Ligand Molecules. *J. Comput.-Aided Mol. Des.* **1998**, *12*, 491–502.
- (11) Jones, G.; Willett, P.; Glen, R. C. A Genetic Algorithm for Flexible Molecular Overlay and Pharmacophore Elucidation. *J. Comput.-Aided Mol. Des.* **1995**, *9*, 532–549.
- (12) Iwase, K.; Hirono, S. Estimation of Active Conformations of Drugs by a New Molecular Superposing Procedure. *J. Comput.-Aided Mol. Des.* **1999**, *13*, 499–512.
- (13) Handschuh, S.; Wagener, M.; Gasteiger, J. Superposition of Three-Dimensional Chemical Structures Allowing for Conformational Flexibility by a Hybrid Method. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 220–232.
- (14) Miller, M.; Fluder, E. M.; Castonguay, L. A.; Culberson, J. C.; Mosley, R. T.; Prendergast, K.; Kearsley, S. K.; Sheridan, R. P. MEGA-SQ: A Method Using the SQuEAL Function to Find the Optimal Superposition of Several Quasi-Flexible Molecules. *Med. Chem. Res.* **1999**, *9*, 513–534.
- (15) Blinn, J. R.; Rohrer, D. C.; Maggiora, G. M. Field-Based Similarity Forcing in Energy Minimization and Molecular Matching. *Pacific Symposium on Biocomputing 1999*; Altman et al., Eds.; World Scientific Publishing: Singapore.
- (16) Kearsley, S. An Alternative Method for the Alignment of Molecular Structures: Maximizing Electrostatic and Steric Overlap. *Tetrahedron Comput. Methodol.* **1990**, *3*, 615–633.
- (17) Lemmen, C.; Lengauer, T. Computational Methods for the Structural Alignment of Molecules. *J. Comput.-Aided Mol. Des.* **2000**, *14*, 215–232.
- (18) Schaefer, M.; Karplus, M. A Comprehensive Analytical Treatment of Continuum Electrostatics. *J. Phys. Chem.* **1996**, *100*, 1587–1599.
- (19) Bush, B. L.; Sheridan, R. P. PATTY: A Programmable Atom Typer and Language for Automatic Classification of Atoms in Molecular Databases. *J. Chem. Inf. Comput. Sci.* **1993**, *33*, 756–762.
- (20) Wildman, S. A.; Crippen, G. M. Prediction of Physicochemical Parameters by Atomic Contributions. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 868–873.
- (21) Ferguson, D. M.; Raber, D. J. A New Approach to Probing Conformational Space with Molecular Mechanics: Random Incremental Pulse Search. *J. Am. Chem. Soc.* **1989**, *111*, 4371–4378.
- (22) MOE software available from Chemical Computing Group Inc., Montreal, Canada. Consult <http://www.chemcomp.com> for further information.
- (23) Gill, P.; Murray, W.; Wright, M. H. *Practical Optimization*; Academic Press: New York, 1981.
- (24) Halgren, T. A. The Merck Molecular Force Field. I. Basis, Form, Scope, Parametrization, and Performance of MMFF94. *J. Comput. Chem.* **1996**, *17*, 490–519.
- (25) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Shindyalov, I. N.; Bourne, P. E. The Protein Data Bank. *Nucl. Acid Res.* **2000**, *28*, 235–242.
- (26) Anstead, G. M.; Carlson, K. E.; Katzenellenbogen, J. A. The Estradiol Pharmacophore: Ligand Structure Estrogen Receptor Binding Affinity Relationships and a Model for the Receptor Binding Site. *Steroids* **1997**, *62*, 268–303.
- (27) Bolin, J. T.; Filman, D. J.; Matthews, D. A.; Hamlin, R. C.; Kraut, J. Crystal Structures of Escherichia Coli and Lactobacillus Casei Dihydrofolate Reductase Refined at 1.7 Å Resolution. I. General Features and Binding of Methotrexate. *J. Biol. Chem.* **1982**, *257*, 13650–13662.
- (28) Wurtz, J.-M.; Egner, U.; Heinrich, N.; Moras, D.; Mueller-Fahrnow, A. Three-Dimensional Models of the Estrogen Receptor Ligand Binding Domain Complexes, Based on Related Crystal Structures and Mutational Structure–Activity Data. *J. Med. Chem.* **1998**, *41*, 1803–1814.
- (29) Schmidt, J.; Mercure, J.; Feher, M.; Dunn-Dufault, R.; Peter, M.; Redden, P. De Novo Design, Synthesis and Evaluation of Novel Non-Steroidal High Affinity Ligands for the Estrogen Receptor. Unpublished results.

JM0002634