

High-Throughput Modeling of Human G-Protein Coupled Receptors: Amino Acid Sequence Alignment, Three-Dimensional Model Building, and Receptor Library Screening

Caterina Bissantz,[†] Antoine Logean, and Didier Rognan*

Bioinformatics Group, Laboratoire de Pharmacochimie de la Communication Cellulaire (CNRS UMR 7081),
74 Route du Rhin, B.P.24, F-67400 Illkirch, France

Received August 19, 2003

The current study describes the development of a computer package (GPCRmod) aimed at the high-throughput modeling of the therapeutically important family of human G-protein coupled receptors (GPCRs). GPCRmod first proposes a reliable alignment of the seven transmembrane domains (7 TMs) of most druggable human GPCRs based on pattern/motif recognition for each of the 7 TMs that are considered independently. It then converts the alignment into knowledge-based three-dimensional (3-D) models starting from a set of 3-D backbone templates and two separate rotamer libraries for side chain positioning. The 7 TMs of 277 human GPCRs have been accurately aligned, unambiguously clustered in three different classes (rhodopsin-like, secretin-like, metabotropic glutamate-like), and converted into high-quality 3-D models at a remarkable throughput (ca. 3s/model). A 3-D GPCR target library of 277 receptors has consequently been setup. Its utility for “in silico” inverse screening purpose has been demonstrated by recovering among top scorers the receptor of a selective GPCR antagonist as well as the receptors of a promiscuous antagonist. The current GPCR target library thus constitutes a 3-D database of choice to address as soon as possible the “virtual selectivity” profile of any GPCR antagonist or inverse agonist in an early hit optimization process.

INTRODUCTION

G protein-coupled receptors (GPCRs) constitute a superfamily of membrane receptors of outmost importance in pharmaceutical research.¹ Hence, GPCRs are the macromolecular targets of ca. 30% of marketed drugs, with 26 out of the top 100 selling drugs targeting this protein family.² The first draft of the human genome suggests that over 800 genes encode for a GPCR,³ out of which only 30 are currently addressed by marketed drugs. If one excludes the family of sensory receptors which significantly differ from other GPCRs, about 400 receptors are potentially druggable with ca. 120 proteins being still considered as orphan targets.² Most interesting GPCRs can be classified into three families or classes, depending on their amino acid sequence.^{4,5} Class I GPCRs belong to the family of rhodopsin-like receptors⁵ recognized by biogenic amines (dopamine, serotonin, and histamine, etc.) and small peptides (chemokines and neuropeptides, etc.). They are characterized by a small extracellular N-terminal domain, a canonical seven transmembrane (7 TM) domain, and a long intracellular C-terminal domain. In most of the cases, the ligand binding cavity is delimited by the 7 TM domain, though peptide-specific receptors may use two of the three extracellular loops to encompass the peptide binding site. This class is believed to contain the vast majority of GPCRs (about 240 nonolfactory receptors).⁵ Class II GPCRs belong to the family of secretin-like receptors and recognize protein hormones (vasointestinal peptide and

glucagon, etc.). Although the heptahelical 3-D fold of class II GPCRs is believed to be similar to that of class I receptors, they significantly differ from the latter class by a much larger N-terminal domain that delimits the hormone-binding site and a large C-terminal domain. About 60–65 GPCRs have been postulated to contribute to class II.² A recent phylogenetic analysis of human GPCRs suggests that this class may contain a family of adhesion receptors.⁵ Last, class III GPCRs belong to the family of metabotropic-like receptors. They recognize low molecular weight charged ligands (glutamate, calcium, and γ -aminobutyric acid, etc.) through a very large N-terminal domain composed of two symmetric lobes. Apart from the conserved 7 TM domain, they are characterized by rather short intracellular loops and a large C-terminal domain. Current estimates suggest that only a few GPCRs (about 15) may contribute to class III.⁵ Interestingly, although the ligand binding site is located in the N-terminal domain, class III GPCRs can be modulated by allosteric ligands (agonists, antagonists) binding to the 7 TM cavity.⁶

Traditionally, the first stage in the design of GPCR ligands has focused on the potency of the ligands for the selected receptor target. Selectivity toward the host receptor is usually considered once potency has already been reached. It would however be highly desirable to consider selectivity as soon as possible in the design process. Due to the high identity in amino acid sequence between different subtypes of the same receptor (e.g. muscarinic M₁ and M₂ amino acid sequences present an identity of 83% over 189 residues of the 7 TM domain), designing a selective ligand can be a quite cumbersome task.⁷ Ideally, one would like to consider the universe of GPCRs for designing a ligand with the desired

* Corresponding author phone: +33-3-90 24 42 35; fax: +33-3-90 24 43 10; e-mail: didier.rognan@pharma.u-strasbg.fr.

[†] Current address: F. Hoffmann-La Roche, Molecular Structure and Design, CH-4070 Basel, Switzerland.

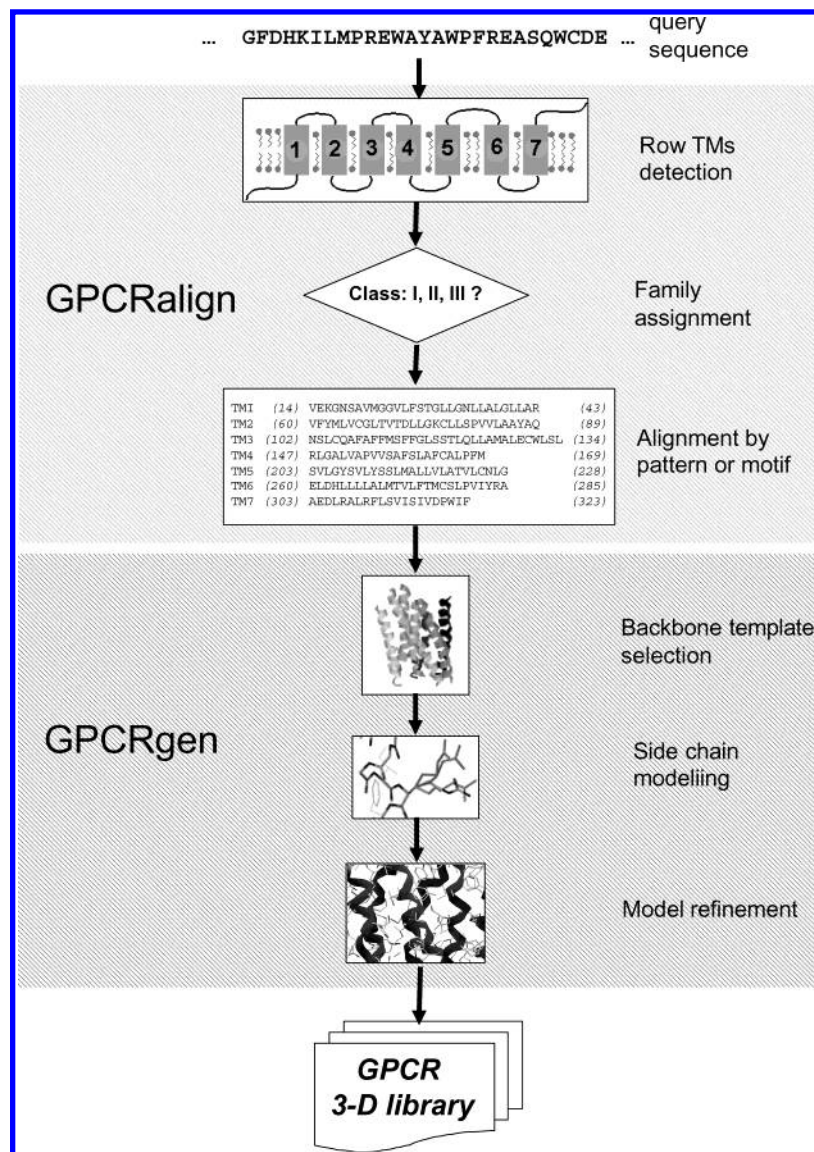


Figure 1. Overall flow chart of the GPCRmod program.

selectivity profile. As addressing this issue by high-throughput screening is currently impossible, “in silico” screening⁸ could provide a reasonable start. Indeed the recently described 2.8-Å resolution X-ray structure of bovine rhodopsin⁹ provides a possible 3-D template for modeling other GPCRs. It has been recently demonstrated that GPCR homology models are indeed accurate enough to identify known antagonists seeded in a randomly chosen “drug-like” library^{10,11} and to exactly map antagonist-binding sites.^{12,13} In the present paper, we present a software package (GPCRmod) for the high-throughput modeling of GPCRs, which provides high-quality ground-state models of the 7 TM domain. The resulting GPCR target library can be screened by an inverse docking tool to predict the most likely receptor(s) of any putative GPCR antagonist or inverse agonist.

METHODS

GPCRmod Architecture (Figure 1). GPCRmod is composed of two modules. The first one (GPCRalign) is a Perl module that performs the sequence alignment in a three-steps procedure: row location of the TMs within the target sequence, assignment of the target to one of the three GPCR

families and final refinement of the alignment of each TM. Perl was chosen here for both its high portability and ability to handle sequences in form of ASCII files. The second module (GPCRgen) is a Java library that ensures the structural part of the modeling procedure. The application programming interface (API) proposed by the library defines an object oriented description of the protein in either Cartesian or internal coordinates. Different tools of the library permit the user to manipulate the macromolecular models in an intuitive way. They are used by GPCRgen to automatically read from a SQL database the amino acid sequences of the 7 TM domains (given by GPCRalign) as well as the Cartesian coordinates of 9 GPCR templates. The 3-D template structures are then converted in internal coordinates, fragmented into rigid bodies and used by GPCRgen to set the two different knowledge-based libraries: one for the backbones, the other for the side chains. Other tools of the library will then be used to reconstruct each GPCR 3-D model by piecing together the different rigid bodies.

Both modules can be run independently. The highly portability of Perl and Java ensures GPCRmod to be run on many currently available operating systems: GPCRmod was

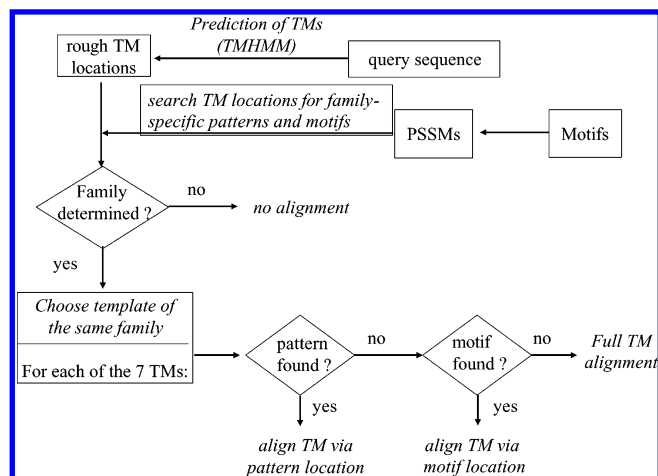


Figure 2. Overall flow chart of the alignment tool (GPCRalign) of the GPCRmod program.

successfully tested on Irix 6.5, Linux 2.4 (Suse 8.0), and Windows XP. Excepted for the last minimization step performed by AMBER6, both modules of GPCRmod require few CPU resources and can be run on any currently available computer.

Aligning the Amino Acid Sequences of the 7-TMs. (a) Flow Chart of the Alignment Procedure Used by GPCRmod (Figure 2). Given the amino acid sequence of the protein to align, the first step is the prediction of the rough location of the 7 TM helices using the TMHMM (transmembrane hidden Markov model) algorithm,¹⁴ a membrane protein topology prediction method based on a hidden Markov model which is currently considered as the best performing transmembrane prediction program.¹⁵ Predicting the rough location of the TMs in advance presents the advantage of focusing the alignment on short and ungapped amino acid sequences. If 7 TM domains are found, GPCRmod tries to assign the given sequence to one of the three GPCR families by searching the predicted rough TM locations for family characteristic patterns and motifs, as defined in the PRINTS database.¹⁶ First, the program uses only a pattern search to determine the family. If this does not result in a clear determination of the family, it tries to determine the family in a second step by searching for motifs using position specific scoring matrices (PSSM)¹⁷ defined for 19 out of the possible 21 TMs encompassing the 3 GPCR classes taken into account in the current study. The family is considered determined if either (i) two patterns (motifs) of one family are found and none of another family, or (ii) if 3 patterns (motifs) of one family are found and not more than 1 pattern of the other families, or (iii) if 4 or more patterns (motifs) of one family are found and not more than 2 patterns of the other families. The definition of the patterns are chosen strictly so that in this step no "false positives" are found, accepting that rather some receptors might be missed. Those receptors that did not match the patterns can then be classified with the more time-consuming, but more sensitive, motif search.

If the GPCR family determination has been successful, the query sequence can consequently be aligned to the TMs of one of the three template sequences used for matching (template for class I GPCRs, bovine rhodopsin; template for class II GPCRs, human calcitonin receptor; template for class III GPCRs, human extracellular calcium-sensing receptor).

Table 1. GPCR Specific Patterns Classified by Family and TM Domain^a

family	TM	pattern
class I	2	LA..D
	3	[D/E]R[Y/H]
	5	[F/Y]..P.....Y
	6	[F/Y]...W.P
	7	[N/D]P..Y
class II	1	G...S...L
	2	H.[H/N/Q]....[F/Y]..[N/R/K]
	3	W...E...L
	4	GW..P
	6	[K/R]....L. P..G
	7	QG.....C
	7	QG.....C
class III	2	[K/R]....[E/D].[C/S][F/Y]
	3	[S/A]...KT
	4	Q.....[W/L]
	5	Y...L...C
	6	E.[K/R]...F. M.....W....P
	6	E.[K/R]...F. M.....W....P

^a A single dot stands for any amino acid; amino acids in brackets are different options for a single position.

Each TM of the query for which a pattern or motif has been detected can be aligned by simply matching the found pattern (motif) on the respective pattern (motif) in the template sequence. If a pattern was found, pattern matching is used for the alignment; otherwise (if no pattern, but a motif was found) a motif matching is applied. TMs for which neither a pattern nor a motif has been detected are aligned by applying a full TM alignment. The templates of classes II and III have been previously aligned to the sequence of bovine rhodopsin, using ClustalW.¹⁸ Thus we can present here an alignment of the class II and III receptors to receptors of class I. Last, if the whole sequence cannot be assigned to any of the GPCR families by either pattern or motif recognition, no alignment is carried out and the corresponding protein discarded from the GPCR library.

(b) Definition of GPCR Patterns and Motifs. A pattern is here defined as a short sequence of continuous or discontinuous highly conserved amino acids, typical for one TM and one GPCR class. As example, the [E/D]R[Y/H] pattern is characteristic of TM3 for class I rhodopsin-like GPCRs.¹⁹ We defined patterns for 5 TMs of both the rhodopsin-like and metabotropic glutamate-like family and for 6 TMs of the secretin-like family (Table 1). Searching a sequence for the presence of a pattern can be simply carried out using regular expressions giving as the nonambiguous result "present" or "not present".

A motif is an ungapped multiple-sequence alignment of a short sequence region (21–27 amino acids) that includes conserved amino acids. GPCRmod currently uses a collection of 19 transmembrane GPCR motifs (Figure 3A) taken from the PRINTS database.²⁰ The motifs of the rhodopsin-like receptors are built from the PRINTS multiple alignment of 739 sequences (thereby 128 human GPCRs), the motifs of the secretin-like receptors from 59 sequences (14 human GPCRs), and the motifs of the metabotropic glutamate-like receptors from 32 receptors (9 human GPCRs). A motif is currently defined for each of the 7 TMs of both the rhodopsin-like and secretin-like family, but only for 5 TMs of the metabotropic glutamate-like family (TMs 2–6). Each motif has been converted into a position-specific scoring matrix (Figure 3B) calculated as described by Henikoff et al.,¹⁷ from the multiple alignment of all available amino acid sequences of the corresponding TM. For a motif of n

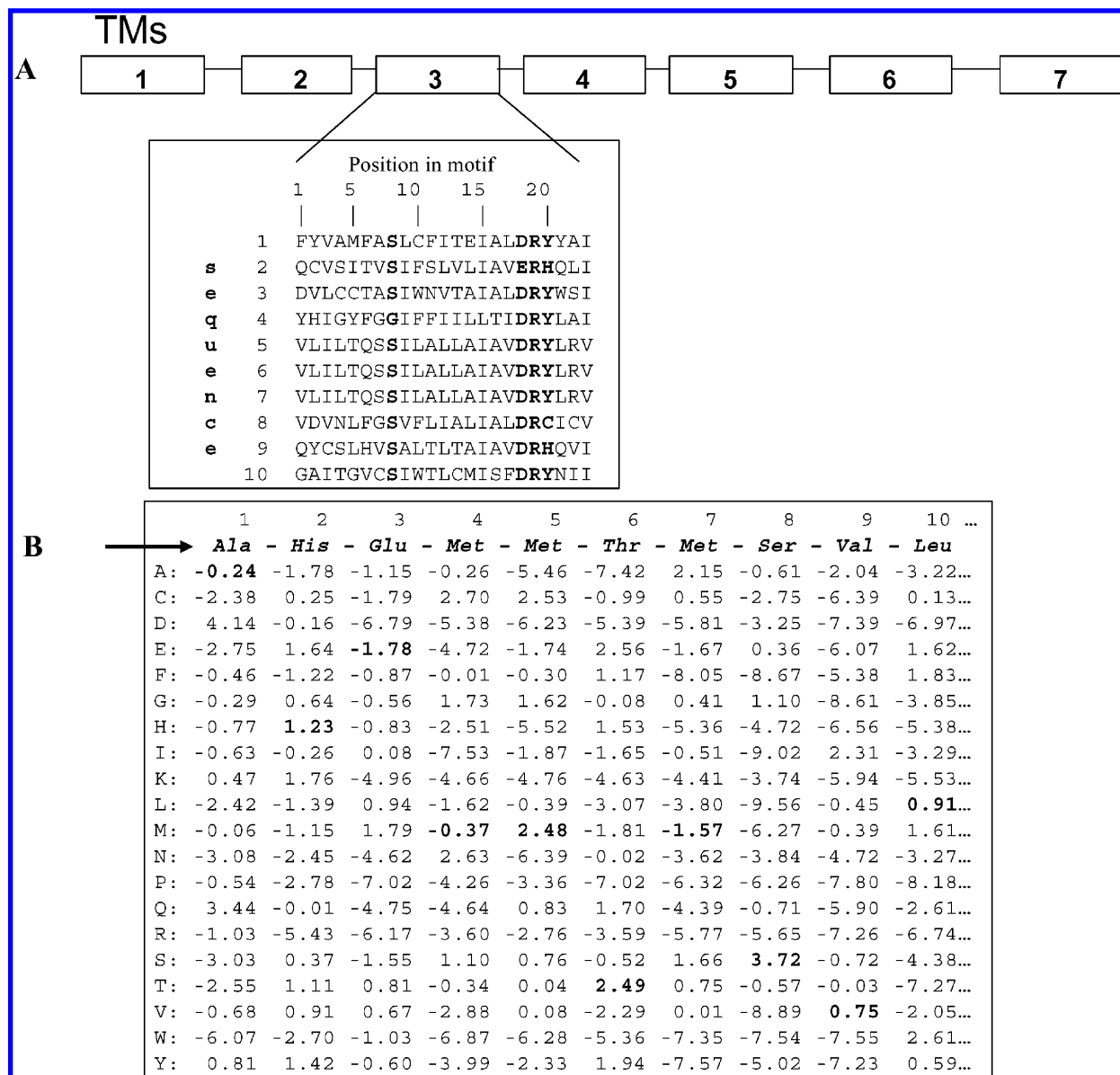


Figure 3. Definition of motifs and position-specific scoring matrices. (A) A motif is a multiple alignment of a short sequence region including highly conserved amino acids, here on the example of TM3 of the rhodopsin-like family. The most conserved positions in the motif (Ser at position 8 and E/D-R-Y/H at positions 18–20) are displayed in bold. (B) Position-specific scoring matrix (PSSM). Each column represents one of the positions in the motif. For every column, a value is given for each of the 20 standard amino acids (rows). Here are shown the positions 1–10 of the TM3 motif for the rhodopsin-like family. The more favored an amino acid at a certain position, the more positive its score for this position. For example, the alignment score of a putative “AHMMTMSVL” sequence to the first 10 amino acids of a motif is derived by summing the individual position-specific score (displayed in bold) of all amino acids of the target sequence ($\text{Score}_{\text{align}} = -0.24 + 1.23 - 1.78 - 0.37 + 2.48 + 2.49 - 1.57 + 3.72 + 0.75 + 0.91 = 7.62$). Of course, the real score is computed over all amino acids encompassing the motif.

sequences with a length of m residues, the PSSM will be a matrix of m columns (corresponding to the length of the motif) and 20 lines (one line for each of the 20 amino acids). Each element w_{ca} of the matrix is given by

$$w_{ca} = \log_2 \left(\frac{f_{ca}}{f_a} \right), \quad c = (1, 2, 3, \dots, m),$$

$$a = (1, 2, 3, \dots, 20) \quad (1)$$

where f_{ca} is the frequency of amino acid a at position c of the motif and f_a the overall frequency of amino acid a in a reference data set of protein sequences (background frequency). f_{ca} could simply be calculated by dividing the total number of counts n_{ca} for amino acid a at position c by the

total number of counts N_c (number of aligned amino acids at position c). However, to account for unobserved frequencies and to solve the problem of having to calculate a PSSM score for an amino acid never appearing at a defined position c in the used alignment ($n_{ca} = 0$), pseudocounts¹⁷ are added to the background frequency f_{ca} term as follows:

$$f_{ca} = \frac{n_{ca} + b_{ca}}{N_c + B_c}, \quad \text{with } 0 \leq n_{ca} \leq n \quad (2)$$

where n_{ca} is the total number of counts over the n motif sequences for amino acid a at position c , b_{ca} the pseudocount for amino acid a at position c , N_c the total number of counts, and B_c the total number of pseudocounts at position c . The

Table 2. Template Amino Acid Sequences for TM Domains

TM	receptor	sequence
1	OPSD_BOVIN ^a	WQFSMLAAYMFLIMLGFPINFLTLYVTQ
	CALR_HUMAN ^b	AYVLYLAIVGHSLSIPTLVISLGI FVFFR
	CASR_HUMAN ^c	SWTEPFGIALTLFAVLGIFLTAFVLGVFIK
2	OPSD_BOVIN	PLNYILLNLAVADLFMVFGGFTTTLTSLH
	CALR_HUMAN	QRVTLHKNMFLTYILNSMIIIIHLVEVVPN
	CASR_HUMAN	IVKATNRELSYLLLFSLCCFSSSLFFIGE
3	OPSD_BOVIN	PTGCNLEGGFATLGGEIALWSLVVLAIERVVV
	CALR_HUMAN	PVSKILHFFHQYMMACNYFWMLCEGIYLTLI
	CASR_HUMAN	DWTCRLRQPAFGISFVLCISCILVKTNRVLLVF
4	OPSD_BOVIN	NHAIMGVAFTWVMALACAAPPLV
	CALR_HUMAN	QRLRWYLLGWGFPLVPTTIHAI
	CASR_HUMAN	VFLCTFMQIVICVIWLYTAPPSS
5	OPSD_BOVIN	NESFVIYMFVVFHFIPLIVIFFCYGQ
	CALR_HUMAN	VETHLLYIIHGPVMAALVNVFFLLN
	CASR_HUMAN	SLMALGFLIGYTCLLAICFFFAFKS
6	OPSD_BOVIN	EKEVTRMVIIMVIAFLICWLPHYAGVA
	CALR_HUMAN	YLKAVKATMILVPLLGIQFVVPWRP
	CASR_HUMAN	NFNEAKFITFSMLIFFIIVWISFIPAY
7	OPSD_BOVIN	IFMTIPAFFAKTSAVYNPVIY
	CALR_HUMAN	IYDYVMHSLIHFGGFFVATIIY
	CASR_HUMAN	KFVSAVEVIAILAASFGLLAC

^a OPSD_BOVIN: bovine rhodopsin (rhodopsin-like class I family). ^b CALR_HUMAN: human calcitonin receptor (secretin-like class II family).

^c CASR_HUMAN: extracellular calcium-sensing receptor (metabotropic glutamate-like class III family).

pseudocount b_{ca} is obtained by multiplying the number of different amino acids at position c by a positive real number α defined by

$$b_{ca} = B_c * \alpha \quad (3)$$

with

$$\alpha = \sum_{i=1}^{20} \frac{f_{ci} q_{ia}}{N_c Q_i}$$

where

$$Q_i = \sum_{a=1}^{20} q_{ia}$$

and where f_{ci} is the frequency of amino acid i at position c and q_{ia} the probability for amino acid i to replace amino acid a according to the Blosum62 matrix.²¹

(c) Alignment by Pattern or Motif. For consistency in the further 3-D model building step, the TMs of all GPCRs, whatever the class, are always assigned the same length than the TMs in the X-ray structure of bovine rhodopsin.⁹ Thus, all TMs of the GPCR templates for class II and class III receptors have been aligned to that of bovine rhodopsin. The search for the presence of a pattern in a predicted TM region is a simple regular-expression search that gives as non-ambiguous result “present” or “not present”. If a pattern is detected, a straightforward alignment to a subfamily-specific template (Table 2) is performed by matching the common pattern. Pattern matching is therefore the preferred alignment method. If no pattern is found, the transmembrane region is searched for the presence of a motif in order to allow motif matching. To search the protein sequence for the presence of a motif, the TMHMM-predicted TM sequence is slid along the PSSM. For each possible alignment, a score ($\text{Score}_{\text{align}}$) is calculated by summing the position-specific scores of every amino acid in the sequence window (Figure 3B). The score

of each alignment is used to calculate the probability (odds) of the respective alignment ($\text{odds} = 2^{\text{Score}_{\text{align}}}$). The odds of all possible alignments are then summed ($\text{Odds}_{\text{total}}$) to determine the percentage score ($\text{Score}_{\text{per}}$) of any single alignment as

$$\text{Score}_{\text{per}} = (\text{odds}/\text{Odds}_{\text{total}}) \times 100 \quad (4)$$

A motif is declared as “found” for alignments where the percentage score $\text{Score}_{\text{per}}$ is higher than 30%. The best alignment is then saved and consequently defines the corresponding TM of the sequence query. Due the necessity to define a similarity score and a cutoff value (here 30%) which has to be reached from a sequence in order to match the motif, motif searching is more complicated than pattern searching. Therefore, we use pattern matching whenever possible whereas motif searching is only applied when no pattern could be detected.

(d) Full TM Alignment. In cases where no alignment can be proposed from pattern or motif detection, the TMHMM-predicted TM is aligned to the corresponding TM of the family-specific template (Table 2) with a simple alignment algorithm using the Blosum series as scoring matrix.²¹ Which matrix of the Blosum series is used depends on the maximal sequence identity between the template TM and the sequence of the predicted rough TM location to align. This algorithm consists therefore of two steps:

(1) The template TM sequence is slid along the protein sequence region in which the respective TM is predicted to be located. For each alignment, the sequence identity percentage is calculated as the number of identities divided by the number of residues compared. The maximal sequence identity found determines the scoring matrix used for scoring the possible alignments (Blosum30 if max identity $\leq 30\%$, Blosum45 if $30\% < \text{max identity} \leq 60\%$, Blosum62 if $60\% < \text{max identity} \leq 80\%$, and Blosum80 if max identity $> 80\%$).

(2) For each alignment, the score (log odds) for each amino acid pair is looked up in the Blosum matrix, and the log

TMI			TMII		
OPSD_BOVIN	35	WQFSMLAAYMFLIMLGFPINFLTLYVTVQ	71	PLNYILLNLAVADLFMVFGGFTTTLTSLH	
D2DR_HUMAN	32	PHYNYATLLTLLIAVIVFGNVLVCMASVR	68	TTNYLIVSLAVADLLVATLVMPWVYLEVY	
D3DR_HUMAN	27	RPHAYYALSYCALILAIVFGNGLVCMVLK	63	TNYLVVSLAVADLLVATLVMPWVYLEVY	
EDG2_HUMAN	47	TVSKLVMLGITVCIFIMLANLLVMVAIYV	83	PIYYLMANLAAADFFAGLAYFYLMFNTGPN	
V1AR_HUMAN	49	ELAKLEIAVLAVTFFAVAVLGNSSVLLALXR	85	RMHLFIRHLSLADLAVAFFQVLPQMCWDIT	
ACM1_HUMAN	23	WQVAFIGITTTGLLSLATVTGNLLVLISFKV	59	VNNYFLLSLACADLIIGTFSMNLYTTYLLM	
OPSR_HUMAN	49	GLKVTIVGLYLAVCVGGLGNCLVMYVILR	85	ATNIYIFNLALADTLVLLTLPFQGTNILLG	
SMO_HUMAN	9	DMHSYIAAFGAVTGLCTFLTATFVADWRN	46	ILFYVNACFFVGSIGWLAQFMDGARREIVC	
CASR_HUMAN	607	SWTEPFGIALTLFAVLGIFLTAFLVGVFIK	642	IVKATNR ¹ ELSY ² LLLSLLCCFSSSLFFIGE	
TMIII			TMIV		
OPSD_BOVIN	107	PTGCNLEGGFATLGGELWLSVLAIERYVVV	151	NHAIMGVAFTWVMALACAAPPLV	
D2DR_HUMAN	104	RIHCDIFVTLDVMMCTASILNLCAISIDRYTAV	150	RRVTVMISIVWVLSFTISCPLLF	
D3DR_HUMAN	100	RICCDVFVTLDMCTASILNLCAISIDRYTAV	148	RRVALMITAVWVLAFAVSCPLLF	
EDG2_HUMAN	118	VSTWLLRQGLIDTSLTASVANLLAIAIERHITV	162	RRVVVVIVVIWTMAIVMGAIPSV	
V1AR_HUMAN	121	DWLCRVVKHLQVFGMFASAYMLVMTADRYIAV	165	RRSRLMIAAAWVLSFVLSTPQYF	
ACM1_HUMAN	95	TLACDLWLALDYVASNASVMNLLISFDYFSV	140	RRAALMIGLAWLVSVFLWAPAIL	
OPSR_HUMAN	120	NALXKTVIAIDYNNMFTSTFTLTAMSVDRYVAI	165	SKAQAVNVAIWALASVVGVPVAI	
SMO_HUMAN	91	TLSCVLIIFVIVYALMAGVWVFLTYAWHTSF	135	GKTSYFHLTLWSLPFVLTVAILA	
CASR_HUMAN	674	DWTCRLRQPAFGISFVLCIS ³ CI ⁴ LV ⁵ KTNRVLLVF	728	VFLCTFM ⁶ IVICV ⁷ WLYTAPSS	
TMV			TMVI		
OPSD_BOVIN	200	NESFVIYMFVVFHFI ¹ PLIVIFFCYGQ	247	EKEVTRMVIIMVIAFLIC ² WL ³ PYAGVA	
D2DR_HUMAN	186	NPAFVVYSSIVSFYV ⁴ PFITLLVYIK	368	EKKATQMLAIVLG ⁵ VFI ⁶ CWL ⁷ PPFFITH	
D3DR_HUMAN	185	NPDFVIYSSVVSFYLP ⁸ FGVTVLVYAR	324	EKKATQMVAVLVGAF ⁹ IVCWL ¹⁰ PPFFLTH	
EDG2_HUMAN	202	YSDSYLVFWAIFNLVTFVVMVVL ¹¹ YAH	253	MMSLLKT ¹² VVIVLGAF ¹³ IICW ¹⁴ T ¹⁵ PGLVLL	
V1AR_HUMAN	213	SRAYVTWMTGGIFVAPV ¹⁶ VILGTCYGF	286	KIRTVKMTFVIVTAY ¹⁷ IVCW ¹⁸ APFFIIQ	
ACM1_HUMAN	185	QPIITFGTAMAAFYLPV ¹⁹ TVMCTLYWR	360	EKKAARTLSAILLAF ²⁰ ILTW ²¹ TPYNIMV	
OPSR_HUMAN	212	GPVFAICIFLTSEFIV ²² PLVISVCYSL	258	LRRITRLVLVVAVF ²³ VG ²⁴ CW ²⁵ TPVQVFS	
SMO_HUMAN	180	RAGFVLAPIGLVLIVGGYFLIRGVMT	216	SKINETMLRLGIFG ²⁶ FLAF ²⁷ GFVLITFS	
CASR_HUMAN	769	SLMALGFLIGYTCL ²⁸ AAI ²⁹ CF ³⁰ FFAFKS	800	NFNEAKFIT ³¹ ESML ³² LIFFIV ³³ ISFI ³⁴ AY	
TMVII					
OPSD_BOVIN	286	IFMTIPAFFAKTSAVYNPVIY			
D2DR_HUMAN	406	VLYSAFTWLGYN ¹ SAVNPIIY			
D3DR_HUMAN	363	ELYSATTWLGYN ² SALNPVIY			
EDG2_HUMAN	291	AYEKFFLLAEFNSAMNPIIY			
V1AR_HUMAN	328	PTITITALLGSLN ³ SCNPNPIY			
ACM1_HUMAN	398	TLWELGYWLCYVN ⁴ STINPMCY			
OPSR_HUMAN	299	AILRFCTALGYVN ⁵ SCNPIIY			
SMO_HUMAN	291	EKINLFAMFGTGIAM ⁶ STVWWT			
CASR_HUMAN	831	KFVSAVEVIAILAAS ⁷ FGLLAC			

Figure 4. Alignment of the 7 TMs of the 9 structural templates (one X-ray structure, 8 homology models) used to build two knowledge-based libraries (*B-lib* and *PSR-lib*). Numbers at the beginning of each block line indicate the starting position of each helix, based on the SwissProt numbering. Class I and III patterns are enclosed by light and dark gray boxes, respectively.

odds of all positions are summed up to obtain the alignment score. The alignment with the largest positive log odds is accepted as the final solution.

Matrices such as the Blosum series are however not ideal for aligning transmembrane regions since they have been developed using soluble proteins. Thus, one has to be critical if the proposed alignment is correct. Therefore, we use this procedure only as last option, when both pattern and motif matching failed.

Automated 3-D Model Generation of TM Domains. (a) Setting-up Knowledge-Based Libraries. Eight GPCR models (dopamine D₂ and D₃ receptors; muscarinic M₁ receptor; EDG-2 receptor; smoothened, nociceptin ORL-1 receptor; vasopressin V_{1a} receptor; calcium sensing receptor) modeled by homology to the X-ray structure of bovine rhodopsin⁹ according to a previously reported procedure¹⁰ were retained as templates for the proposed high-throughput comparative modeling. The accuracy of these 3-D templates has been validated by either side-directed mutagenesis,^{13,22} covalent labeling,²² or their capacity to discriminate true antagonists from randomly chosen druglike molecules.¹⁰ Extraction of

the 7 TM residues for the X-ray structure (bovine rhodopsin) as well as for the eight above-reported models afforded nine 3-D structural templates (Figure 4) which were fitted together with the "Magic Fit" tool of Swiss-PDBviewer.²³ This subset of nine template structures is fragmented by GPCRmod into two classes of rigid bodies used to define two knowledge-based libraries. The first class of rigid bodies is composed by the nine backbone structures. It constitutes the first knowledge-based library, *B-lib*, containing the Cartesian coordinates of the backbone atoms. The second class of rigid bodies stored in a second knowledge-based library, *PSR-lib*, is a position-specific rotamer library storing the amino acid and its rotameric state (expressed in internal coordinates) at each position of the 7 TMs.

(b) 3-D Model Building. To build the backbone of each GPCR target, GPCRmod selects first in the *B-lib* library one of the 9 templates presenting the highest sequence identity in the 7 TMs with the target sequence. Following the previously determined target-template alignment, each side chain of the target is constructed (Figure 5). Information about the initial placement of the side chains comes from

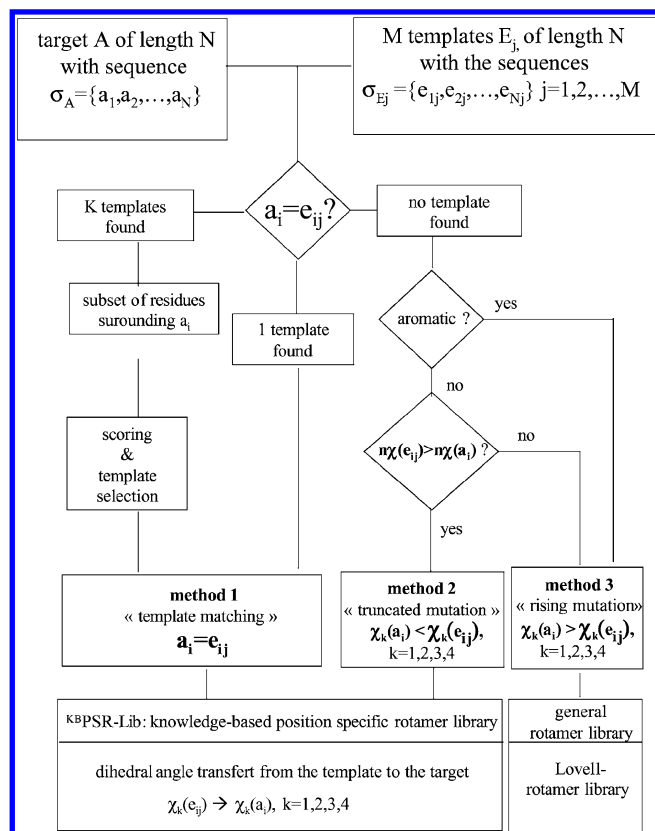


Figure 5. Overall flow chart of the 3-D building tool (GPCRgen) of the GPCRmod program.

the combination of two rotamer libraries: the previously described knowledge-based position-specific rotamer library (*PSR-lib*) as well as a backbone-independent rotamer library.²⁴ Internal coordinates are used to build the side chains. While bond lengths and angles are derived from AMBER 6.0,²⁵ dihedral angles are taken from one of the above-mentioned rotamer libraries. To construct each target residue, GPCRmod uses the following set of “rule-based” modeling procedures (Figure 5). Starting from a target protein *A* of length *N* with a sequence $\sigma_A = \{a_1, a_2, \dots, a_N\}$ and a library of *M* templates *E_j* with their *M* corresponding sequences $\sigma_{Ej} = \{e_{1j}, e_{2j}, \dots, e_{Nj}\}$, *j* = 1, 2, ..., *M*, GPCRmod selects among the *M* templates *E_j* those carrying the same residue at position *i* as in the target *A* ($e_{ij} = a_i$). If more than one template carrying the right residue at the desired position *i* is found, GPCRmod first determines for the target *A* a subset of *Q* residues $\sigma_A^i = \{sa_1^i, sa_2^i, \dots, sa_Q^i\}$ surrounding *a_i*. The subset σ_A^i is formed by the *Q* residues having at least one heavy atom in a 5-Å sphere radius centered on *a_i*. The corresponding subsets $\sigma_{Ej}^i = \{se_{1j}^i, se_{2j}^i, \dots, se_{Qj}^i\}$ are considered for each of the *k* selected templates *E_j* for which $e_{ij} = a_i$. For each of the *k* pair ($\sigma_A^i, \sigma_{Ej}^i$), a similarity score $s(\sigma_A^i, \sigma_{Ej}^i)$ between the target and the template subsets surrounding *a_i* is determined as follows:

$$s(\sigma_A^i, \sigma_{Ej}^i) = \sum_{n=1}^Q D(sa_n^i, se_{nj}^i) \quad \text{with} \quad D(sa_n^i, se_{nj}^i) = D(i, j) = D_{ij} \quad (5)$$

where *D_{ij}* is a 20 × 20 residue substitution matrix²⁶ that gives a measure of the chemical similarity between the residue *i*

and *j*, *i, j* being one of the 20 natural amino acids. The template *i* having the subset σ_{Ej}^i presenting the highest similarity with the target pocket σ_A^i is selected and its dihedral angles are taken from *PSR-lib* and used to define the side chain of the target residue *a_i*. If just one template is found, the dihedral angles of the template can be directly assigned to the target side chain. These cases where one or more template side chains are found in the knowledge-based rotamer library *PSR-lib* is referred in GPCRmod as “method 1” of side-chain positioning called “template matching”, as shown in Figure 5.

If no templates are found, GPCRmod checks first if *a_i* presents an aromatic side chain (Phe, Tyr, Trp). Indeed, aromatic residues adopt most of the time a χ_2 angle of $\pm 90^\circ$ differing from the most frequently observed -60° (gauche[−])/ 180° (trans) χ_2 values of other amino acids. Thus, if *a_i* is aromatic, a general backbone-independent rotamer library²⁴ is used to set the dihedral angles of *a_i*. These cases where the side chains are aromatic are referred to in GPCRmod as “method 3” (Figure 5). If *a_i* is not aromatic, GPCRmod tries to find in the *PSR-lib* library residues at positions *i* (*e_{ij}*) having a dihedral degree of freedom higher or equal to that of *a_i* ($n\chi(e_{ij}) \geq n\chi(a_i)$, with $n\chi(e_{ij})$ and $n\chi(a_i)$ being the number of rotatable bonds of *e_{ij}* and *a_i*, respectively). If more than one template is found, the previously defined substitution matrix²⁶ is used to select the *e_{ij}* amino acid that presents the higher physicochemical similarity with *a_i*. The dihedral angles of the selected *e_{ij}* are used to set those of *a_i*. These cases constitute in GPCRmod the “method 2” of side-chain positioning and are called “truncated mutation” (Figure 5). If none of the *M* templates *E_j* can be used to model *a_i*, GPCRmod selects the top-ranked rotamer (found with the higher probability in the PDB) of the general backbone-independent rotamer library.²⁴ These cases constitute “method 3” of side-chain positioning, referred to by GPCRmod as “rising mutations” (Figure 5).

(c) Model Refinement. Once the side chains have been built, the model is refined in a two-step procedure involving a first refinement aimed at removing main steric clashes, followed by a force-field energy minimization. For each residue of the 3-D model, a list of neighboring heavy atoms closer than 2.5 Å to any atom of the inspected residue is defined. All residues whose list contains more than 10 atoms are considered separately. This value of 10 atoms, determined by a trial-and-error procedure, is the best compromise to detect almost all clashes, mainly those involving tyrosine, phenylalanine, and tryptophane. Going through the 7 helices (N-terminus to C-terminus of helix I, then N-terminus to C-terminus of helix II, etc.), for each problematic residue (having at least 10 heavy atoms closer than 2.5 Å), all rotameric states of the general backbone-independent library²⁴ are checked, beginning with the top-ranked rotamer (found with the higher probability in the PDB). For each possible rotameric state, the list of neighboring atoms is recalculated and the one presenting the less neighboring close atoms is finally selected. For all GPCRs modeled in the current study, most of the clashes were associated with aromatic residues and could be successfully resolved by applying this row and fast refinement protocol.

In a second step, hydrogen atoms are added using AMBER6 geometries and the model is relaxed using the AMBER6 force field.²⁵ The model is first refined by 1000

steps of the descent method followed by a maximum of 1000-steps conjugate gradient minimization, unless the root-mean-square of the potential energy gradient converged to a threshold of $0.25 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{\AA}^{-1}$. In practice, the upper limit of 2000 energy minimization steps was never reached as the refinement always converge during the conjugate gradient minimization. Energy refinement was performed under vacuum using a distance-dependent dielectric function ($\epsilon = 4r$) and a twin cutoff (10.0 and 15.0 \AA) to calculate nonbonded interactions.

Inverse Screening of the GPCR Target Database. (a) Customizing the GOLD2.1 Docking Program for Inverse Screening. The necessary GOLD²⁷ input files used for screening a single ligand against a library of protein targets is generated by an in-house Perl module (InvGOLD). All GPCR entries are first stored in mol2²⁸ format in a single directory. An additional file stores the center of mass of the TM cavity of all GPCR entries. A configuration file (gold.conf) is then defined for each GPCR entry in an entry-specific directory with the corresponding protein file name and center of mass.

(b) Setting-up Ligand Coordinates. Starting from Isis/Draw²⁹ 2-D sketches, a 3-D structure of the ligand is generated with Concord.³⁰ A quick energy-minimization protocol is then used to refine the Concord structure, using the TRIPOS force field³¹ and 1000 steps of Newton–Raphson energy refinement. Final ligand coordinates are stored in TRIPOS mol2 format.

(c) Ligand Docking. Seven speed-up settings of the GOLD software²⁷ were used in the current study. For each of the 10 independent genetic algorithm (GA) runs, a maximum number of 10 000 GA operations were performed on a single population of 100 individuals. Operator weights for crossover, mutation, and migration in the entry box were set to 100, 100, and 0, respectively. To allow poor nonbonded contacts at the start of each GA run, the maximum distance between hydrogen donors and fitting points was set to 4 \AA , and nonbonded van der Waals energies were cutoff at a value equal to $9 k_{ij}$ (well depth of the van der Waals energy for the atom pair i,j). To further speed up the calculation, the GA docking was stopped when the top 3 solutions were within 1.5- \AA root-mean-square deviation (rmsd) of each other. The GOLD output files are then treated by InvGOLD to generate a ranking list (maximum fitness value, average fitness) for all GPCR targets and target-specific docked coordinates of the investigated ligand.

RESULTS

Alignment of Human GPCR Amino Acid Sequences. The GPCRmod program has been applied to align 277 human GPCRs from the SWISSPROT/TREMBL database.³² An example of the GPCRmod alignment output is displayed in Figure 6.

(a) Aligning Receptors of the Rhodopsin-like Family (Class I). 208 of the 235 rhodopsin-like receptors could be unambiguously classified via pattern search (see Table 3). About 80% of the rhodopsin-like GPCRs possess at least 3 class-I-specific patterns. Patterns from other classes have been detected in only four cases (Swiss Prot id: B2AR, PD2R, EDG4, O14804), all of them originating from class II fingerprints. The remaining 27 rhodopsin-like receptors

that were not classified via pattern search were all identified as rhodopsin-like by the motif search. 25 out of these 27 GPCRS depict 5 or more characteristic class I motifs. Class II and class III motifs can be accidentally detected for the small subset of receptors that escaped pattern detection but never enough to perturb family assignment.

For 14 class I receptors, at least one TM was aligned using the full TM alignment protocol (see Methods) since neither a pattern nor a motif was found. 7 of these 14 cases (Figure 6) belong to the family of the prostaglandine/thromboxane receptors. In 5 prostaglandine receptors, TM5 had to be aligned with this protocol; in 2 cases TM7 was concerned. It has already been reported that the prostaglandine/thromboxane receptors only partially match the motifs of the rhodopsin-like family, lacking notably the motif in TM5.²⁰

It should be stated that any alignment of a transmembrane region using general matrices such as the Blosum series always has to be regarded critically since these matrices have been developed using soluble proteins and not membrane proteins. The higher the similarity between the query sequence and the template, the higher the probability of obtaining the correct alignment. If possible, we therefore repeated the alignment of those GPCRs for which the full TM alignment protocol had to be used, using sequences of TMs that belong to the same (sub-)family as the query and that were already unambiguously aligned via a pattern or motif as additional templates. All alignments obtained with the full TM alignment protocol were then compared with solutions suggested from other alignment protocols such as ClustalW. We then kept the most realistic alignment in terms of (i) properties of the amino acids aligned to residues in the rhodopsin sequence that are in other conserved GPCRs and (ii) the length of the resulting loops connecting the respective TM with the preceding and following TMs.

For 10 of the 14 receptors for which the full TM alignment protocol was applied, the original alignment proposed by our algorithm was kept since suggested solutions from other programs were identical. For 4 receptors of the prostaglandine family (SwissProt id: PD2R, PE22, PE24, TA2R), repeating the alignment of TM5 using the four already aligned prostaglandine receptors as additional templates gave different results that are now equivalent to the one obtained by applying ClustalW¹⁸ for aligning the TM5 of the prostaglandine receptors. Therefore, these alignments were kept thereafter.

(b) Aligning Receptors of the Secretin-like Family (Class II). 27 sequences were classified via the pattern search (Table 3). Not a single pattern from the other two GPCR classes was found (see Table 3). Two additional receptors (SwissProt id: CLR3, Q9UL61) could be classified into the class II GPCR family by detection of at least 3 class II-specific motifs (Table 3). 14 of the secretin-like sequences have at least one TM for which neither a pattern nor a motif could be identified. These TM sequences had thus to be aligned using the full TM alignment. As described for class I, we compared these alignments with the alignment given by ClustalW and corrected them in 4 cases (SwissProt ids: BAI1, BAI3, CD97, EMR1) where a significant improvement was obtained by using the other already aligned TMs of the same family as additional templates.

(c) Aligning Receptors of the Metabotropic Glutamate-like Family (Class III). This family comprises the metabo-

SwissProt entry	TM	Sequence	Family	TM Start	TM End
PD2R_HUMAN	TM1	VEKGN SAVMGGVLFSTGLLGNLLALGLLAR	RHODOPSIN	14	43
PD2R_HUMAN	TM2	VFYMLVCGLTVDLLGKCLLSPVVLAAYAQ	RHODOPSIN	60	89
PD2R_HUMAN	TM3	NSLCQAF AFMSFFGLSSTLQLLAMALECWLSL	RHODOPSIN	102	134
PD2R_HUMAN	TM4	RLGALVAPVVSASFSLAFCALPFM	RHODOPSIN	147	169
PD2R_HUMAN	TM5	SVLGYSVLYSSLMALLVLATVLCNLG	RHODOPSIN	203	228
PD2R_HUMAN	TM6	ELDHL LLLALMTVLFTMCSLPVIYRA	RHODOPSIN	260	285
PD2R_HUMAN	TM7	AEDLRALRFLSVISIVDPWIF	RHODOPSIN	303	323
PE21_HUMAN	TM1	PPSGASPALPIFSMTLGAVSNLLALALLAQ	RHODOPSIN	30	59
PE21_HUMAN	TM2	TFLLFVASLLATDLAGHVI PGALVLRLYTA	RHODOPSIN	72	101
PE21_HUMAN	TM3	GGACHFLGGCMVFVGLCP LLLGCGMAVERCVGV	RHODOPSIN	107	139
PE21_HUMAN	TM4	ARARLALA AAVAVALAVALLPLA	RHODOPSIN	152	174
PE21_HUMAN	TM5	RQALLAGLFASLGLVALLAALVCNTL	RHODOPSIN	199	224
PE21_HUMAN	TM6	DVEMVGQLVGIMVWSCICWSPMLVLV	RHODOPSIN	292	317
PE21_HUMAN	TM7	RPLFLAVRLASWNQILDWPVY	RHODOPSIN	331	351
PE22_HUMAN	TM1	LPPGESPAISSVMFSAGVLGNLIALALLAR	RHODOPSIN	19	48
PE22_HUMAN	TM2	LFHVLVTELVTDL LGTCLISPVVLASYAR	RHODOPSIN	66	95
PE22_HUMAN	TM3	SRACTYFAFAMTFFSLATMLMLFAMALERYLSI	RHODOPSIN	106	138
PE22_HUMAN	TM4	SGGLAVLPVIYAVSLLFCSLPLL	RHODOPSIN	151	173
PE22_HUMAN	TM5	GRTAYLQLYATLLLLLLIVSVLACNFS	RHODOPSIN	193	218
PE22_HUMAN	TM6	ETDHLILLAIMTITFAVCSLPFTIFA	RHODOPSIN	259	284
PE22_HUMAN	TM7	KWDLQALRFLSINSIIDPWFV	RHODOPSIN	295	315
PE23_HUMAN	TM1	DCGSVSVAFPI TMLLTGFVGNALAMLLVSR	RHODOPSIN	46	75
PE23_HUMAN	TM2	SFLLCIGWLALTDLVGQLLTTPVVI VVYLS	RHODOPSIN	87	116
PE23_HUMAN	TM3	GRLCTFFGLTMTVFGLSSLFIASAM AVERALAI	RHODOPSIN	127	159
PE23_HUMAN	TM4	RATRAVLLGVWLAVLAFALLPVL	RHODOPSIN	172	194
PE23_HUMAN	TM5	GNLFFASAF AFLGLLALT VTFSCNLA	RHODOPSIN	226	251
PE23_HUMAN	TM6	TTETAIQLMGIMCVLSVCWSP LLI MM	RHODOPSIN	277	302
PE23_HUMAN	TM7	NFFLIAVRLASLNQILDWPVY	RHODOPSIN	326	346
PE24_HUMAN	TM1	DRLNSPVTIPAVMFIFGVVGNLVAIVVLCK	RHODOPSIN	15	44
PE24_HUMAN	TM2	TFYTLVCGLA VTDLLGTLLVSPVTIATYMK	RHODOPSIN	53	82
PE24_HUMAN	TM3	QPLCEYSTFILLFFSLSGLSI ICAMSV ERYLAI	RHODOPSIN	89	121
PE24_HUMAN	TM4	RLAGLTLFAVYASNVLFCA LPMN	RHODOPSIN	134	156
PE24_HUMAN	TM5	AHAAYS MYAGFSSFLILATVLCNVL	RHODOPSIN	182	207
PE24_HUMAN	TM6	BIQM VILLIATSLVVLICSIPLVVRV	RHODOPSIN	267	292
PE24_HUMAN	TM7	NPD LQAIRIASVNPILD PVIY	RHODOPSIN	309	329
PF2R_HUMAN	TM1	TENRLSVFFSVIFMTV GILSNSLAIAILMK	RHODOPSIN	24	53
PF2R_HUMAN	TM2	SFLLLASGLVITDFFGH LINGAIAVFVYAS	RHODOPSIN	65	94
PF2R_HUMAN	TM3	NVLCSIFGICMVFSGLCP LLLGSVMAIERCIGV	RHODOPSIN	105	137
PF2R_HUMAN	TM4	KHVKMMLSGVCLFAVFIALLPIL	RHODOPSIN	150	172
PF2R_HUMAN	TM5	EDRFYLLLF SFLGLLALGVSLLCNAI	RHODOPSIN	197	222
PF2R_HUMAN	TM6	HLEMVIQLLAIMCVSCICWSPFLVTM	RHODOPSIN	244	269
PF2R_HUMAN	TM7	ETTLFALRMATWNQILD PWFVY	RHODOPSIN	284	304
PI2R_HUMAN	TM1	VRGSVGPATSTLMFVAGVVG NGLALGILSA	RHODOPSIN	11	40
PI2R_HUMAN	TM2	AFAVLVTGLAATD LLGTSF LSPAVFVAYAR	RHODOPSIN	48	77
PI2R_HUMAN	TM3	PALCDAF AFAMTFFGLASMLILFAM AVERCLAL	RHODOPSIN	89	121
PI2R_HUMAN	TM4	RCARLALPAIYAF CVLFCALPLL	RHODOPSIN	134	156
PI2R_HUMAN	TM5	GGAAFSLAYAGLVALLVAAIFLCNGS	RHODOPSIN	180	205
PI2R_HUMAN	TM6	EVDHLILLALMTVVM AVCSLP LTI RC	RHODOPSIN	234	259
PI2R_HUMAN	TM7	MGDLLAFRFYAFNPILD PWFV	RHODOPSIN	272	292

Figure 6. GPCRmod output on the example of prostaglandine receptors. For each receptor, indexed by its SwissProt entry name, the GPCR family and the amino acid sequence of the 7 TMs are displayed along with the residue numbers delimiting the TM domains.

tropic glutamate receptors (8 members), the extracellular calcium-sensing receptor and the γ -aminobutyric acid ((GABA) type B subunit 1 and 2 receptors. Additionally, we could classify two orphan receptors (SwissProt id: O75205, O95357) into this family.

Class III GPCRs lack both pattern and motif for TM1 and TM7, which thus always have to be aligned by the full TM alignment algorithm. The TM2–TM6 of the metabotropic

glutamate receptors and the extracellular calcium-sensing receptor—the receptors used for PRINTS motif definition—were all aligned via pattern location. The alignments of TM1 and TM7 of the metabotropic glutamate receptors and the extracellular calcium-sensing receptor that were generated using the full TM alignment were again treated as described for class I. For 5 of these receptors, the alignments generated with the full TM alignment protocol were accepted. In the

Table 3. Statistics for the Pattern and Motif-Based Alignment of 277 GPCRs^a

family	pattern search (all receptors)			classified via pattern search	motif search (receptors not classified via pattern search)			classified via motif search
	class I	class II	class III		class I	class II	class III	
class I (<i>n</i> = 235)	66(5) ^b 62(4) 42(3) 38(2) 23(1) 4(0)	4(1)	0	208	17(7) 4(6) 4(5) 1(4) 1(3)	9(2) 6(1)	2(2) 4(1)	27
class II (<i>n</i> = 29)	0	13(6) 2(5) 9(3) 3(2) 1(1) 1(0)	0	27	1(1)	1(4) 1(3)	0	2
class III (<i>n</i> = 13)	1(1)	0	9(5)	9	1(2) 2(1)	1(1)	2(4) 2(3)	4

^a The motif search applies only for receptors that have not been classified in the previous pattern search. ^b *m*(*n*) indicates that *m* receptors have been assigned by pattern/motif detection over *n* transmembrane domains.

remaining 4 cases (SwissProt id: MGR4, MGR6, MGR7, MGR8), the alignment of TM7 was adjusted by accepting the alignment obtained when using the already-aligned sequences as additional templates. The GABA receptors play a special role in this family. Though grouped together with the metabotropic glutamate and extracellular calcium-sensing receptors into the family of metabotropic glutamate-like receptors, they do not share, according to the PRINTS database, the same motifs. GPCRmod was nevertheless able to correctly classify them as class III specific by detecting motifs for TM3–TM6.

In summary, for 234 (85.7%) of the 277 human GPCRs, all 7 TMs were aligned using either a pattern or motif. The alignment of these receptors was therefore completely automated and did not require any manual intervention. In 41 cases, at least one TM did not contain any pattern or motif and required therefore the full TM alignment protocol, which we decided to control manually since such general matrices are not optimized for transmembrane regions. This manual check was however only performed for 2.9% (57 out of 1939) of the aligned TMs.

Three-Dimensional Model Building. The virtual GPCR target library is composed of 277 3-D models (235 class I, 29 class II, and 13 class III GPCR 3-D models). For class I targets, the nociceptin receptor model was used as the backbone template in 45% of the cases (Figure 7). As no class II template is currently present, class II GPCRs have been built from various class I backbone structures (Figure 7). Due to their extreme amino acid sequence peculiarity, TM domains of the class III receptor were all built from the only class III template (calcium-sensing receptor).

(a) Knowledge-Based Side-Chain Positioning. Side-chain positioning has been achieved using two rotamer libraries. A knowledge-based position specific rotamer library (*PSR-lib*), covering a total of 1701 rotameric states for 20 residues derived from 9 GPCR templates, is used to find the target side chain by either direct template matching (same side chain at the same position; method 1, Figure 5) or by performing a truncated mutation (target side chain with a lower dihedral degree of freedom; method 2, Figure 5). A general backbone-independent rotamer library²⁴ is finally used to build side chains whose rotational degree of freedom is higher than that of available templates stored in the *PSR-*

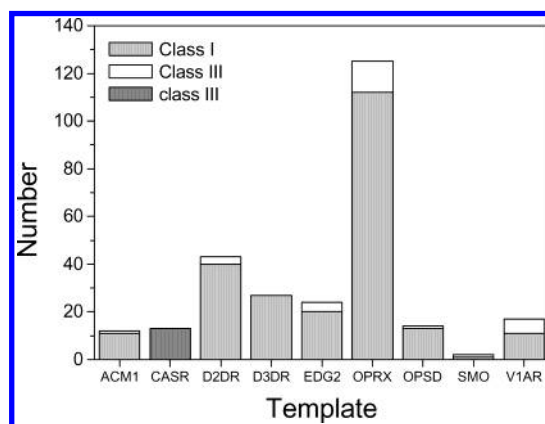


Figure 7. Statistical analysis of backbone templates used to generate 277 GPCR models. Class I, II, and III GPCRs are indicated by light gray, white, and dark gray bars, respectively.

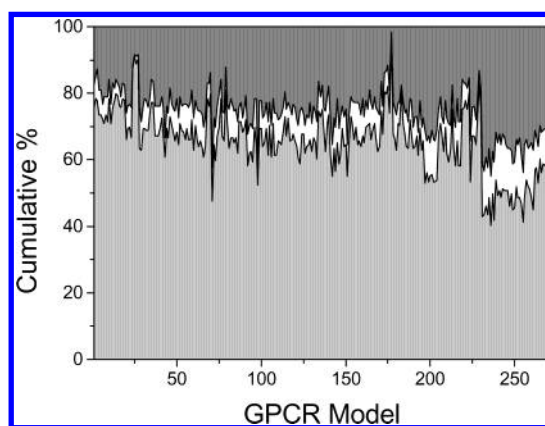


Figure 8. Cumulative percentage of the 3 methods (template matching, light gray surface; truncated mutation, white surface; rising mutation, dark gray surface) used by GPCRmod for the side-chain positioning. The analysis has been performed on 228 class I (models 1–228), 29 class II (models 229–258), and 12 class III (models 259–271) GPCRs. The 8 human GPCR templates used to assist 3-D model generation have been excluded here.

lib library (method 3, Figure 5). Despite the relatively low sequence identity between any of the 269 new GPCR targets and the 9 templates (about 20–25% for the 7 TM domain), direct template matching was possible for ca. 65% of target side chains (Figure 8). In 10% of all cases, the “truncated mutation” method could be used still using the knowledge-

Table 4. Comparison of "Leave-One-Out" GPCRmod Models with Templates

model	GPCR ^a	templates ^b	Procheck score ^c	% of χ^1 < 40° ^d	% of χ^2 < 40° ^e	% of (χ^1, χ^2) < 40° ^f
1	OPSD_B	2-8	0.16	87	71	69
2	D2DR_H	1, 4-8	0.09	86	78	72
3	D3DR_H	1, 4-8	0.06	87	83	76
4	EDG2_H	1-3, 5-8	0.04	80	66	60
5	CASR_H	1-4, 6-8	-0.09	75	63	53
6	V1AR_H	1-8	0.12	79	67	58
7	ACM1_H	1-6, 8	0.00	83	74	66
8	OPRX_H	1-7	-0.01	78	78	66
9	SMO_H	1-8	0.02	73	66	54

^a SwissProt entry (B, bovin; H, human). ^b Handmade GPCR templates used by GPCRmod. ^c Overall Procheck³³ score. ^d Percentage of χ^1 dihedral angles predicted within 40° of the corresponding template. ^e Percentage of χ^2 dihedral angles predicted within 40° of the corresponding template. ^f Percentage of χ^1, χ^2 dihedral angles predicted within 40° of the corresponding template.

based GPCR-derived position-specific rotamer library. For only ca. 25% of all target side chains, the general rotamer library was used (Figure 8).

(b) 3-D Structure Check. Before assessing the correctness of each new 3-D model, the current modeling protocol has first been validated using a "leave-one-out" building approach. Each of the 9 available GPCR templates was reconstructed by GPCRmod using the remaining structures as templates. The GPCRmod models were found to be very close to state-of-the-art template models⁽¹⁰⁾ (Table 4) but are generated within a few seconds instead of a few hours. 82% of the χ^1 dihedral angles as well as 65% of the (χ^1, χ^2) dihedral angles are predicted within 40° of that of the reference structures. The overall stereochemical quality of each leave-one-out model, as considered by Procheck,³³ is correct. No source of errors (wrong stereochemistry, close contacts) could be detected. Hence, the average Procheck score for all 277 GPCR models is 0.08 ± 0.006 . Since the current high-throughput models are aimed at being screened against a single ligand, it is important that our modeling procedure proposes a reliable binding site cavity. Comparing the TM binding cavity of one of the starting templates (human dopamine D₃ receptor) carefully modeled in a previous study¹⁰ with the corresponding GPCRmod model (generated after removing the D₂ and D₃ receptors from the template list) shows that both cavities are rather similar (Figure 9), as exemplified by the observed low rmsd (0.89 Å for heavy atoms).

(c) Inverse Screening of the GPCR Target Database. The herein described InvGOLD script was used to recover, from the GPCR target database, either the known receptor of a selective purinergic P2Y₁ ligand (MRS-2179,³⁴ Figure 10A) or the known receptors of a promiscuous ligand (NAN-190,³⁵ Figure 10B) known to bind to several monoamine receptors with nanomolar affinities (Table 5).

When screening the protein library for putative receptors of MRS-2179, the P2Y₁ receptor is indeed ranked among the top scorers (7th, Figure 10A) with three related receptor subtypes (P2Y₆ ranked 5th, P2Y₅ ranked 9th, and P2Y₁₀ ranked 12th). The remaining five P2Y receptors (P2Y₂, P2Y₄, P2Y₇, P2Y₉, P2Y₁₁) present in the current GPCR database are all ranked beyond the 35th position.

5 out of the 9 known targets of NAN-190, the second ligand investigated herein, are ranked in the top 25 positions,

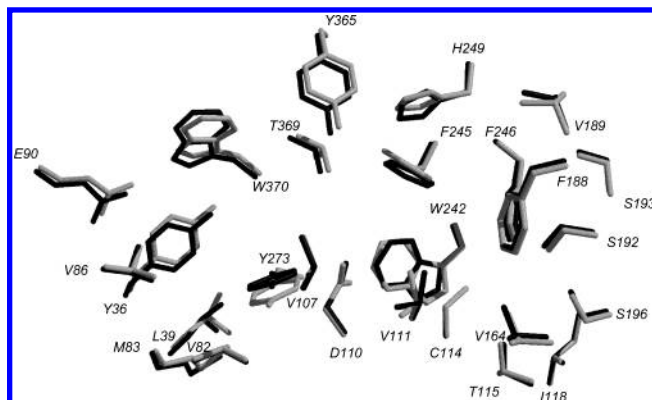


Figure 9. Close-up in the transmembrane cavity of the dopamine D₃ receptor. The template model¹⁰ and the GPCRmod model are displayed by dark and gray sticks, respectively. The root-mean-square deviation between the two models, calculated over heavy atoms, is 0.89 Å.

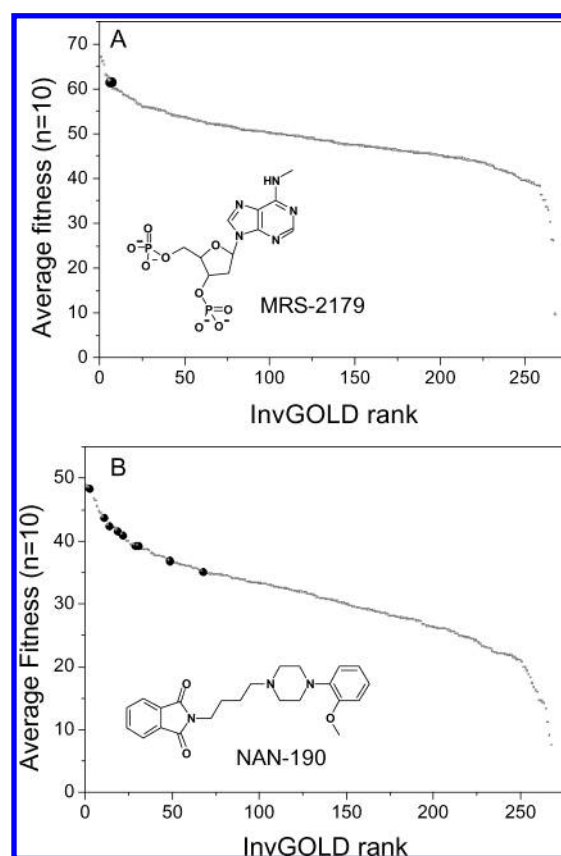


Figure 10. InvGOLD ranking of the true receptor(s) of selective ligands ((A) MRS-2179, high-affinity P2Y₁ receptor antagonist) and of a promiscuous ligand ((B) NAN-90, high-affinity antagonist of the dopamine D₂ receptor, serotonin 5-HT_{1A}, 5-HT_{1C}, 5-HT_{1D}, 5-HT_{2A} receptors, and adrenergic α_{1A} receptor). Known receptor(s) are indicated by a dark ball.

and 7 out of 9 in the top 31 positions (Figure 10B). The worst-ranked true receptor (5-HT_{1A}) is ranked 68th. As expected from many previous studies, there is a very weak correlation between the experimentally determined inhibition constant and the GOLD fitness score (Table 5). The fine selectivity profile for the whole 5-HT receptor family is unfortunately not fully addressed because the 12 5-HT receptor subtypes currently present in our database are all clustered among the top 68 positions.

Table 5. Inverse Screening of a Promiscuous 5-HT Receptor Ligand (NAN-190)

InvGOLD			InvGOLD		
receptor	rank	K_i , ^a nM	receptor	rank	K_i , ^a nM
5-HT _{1A}	68	3	5-HT ₇	6	79
5-HT _{1D}	11	275	D ₂	3	47
5-HT _{1F}	29	703	D ₃	19	3
5-HT _{2A}	31	708	α_{1A}	14	2
5-HT _{2C}	49	630			

^a Inhibition constant (K_i) values were taken from the PDSP K_i database.⁵⁸

DISCUSSION

Straightforward Aligning of Most Human GPCR Transmembrane Domains. Despite the common heptahehical architecture of their transmembrane domains, GPCRs are characterized by a relatively low sequence identity (less than 20%), especially when amino acid sequences of 2 GPCRs from different classes (class I, class II, class III) are compared. Moreover, the length of variable parts of the amino acid sequence (N- and C-terminal domains, extra- and intracellular loops) can vary dramatically. These observations explains our choice for (i) focusing the alignment on isolated TMs which are rather easy to detect by the TMHMM algorithm,¹⁰ (ii) separating GPCRs into homogeneous families, (iii) biasing the alignment procedure toward known GPCR fingerprints (patterns, motifs). 277 GPCRs from three different classes could then be unambiguously classified into one of the 3 main GPCR families and lead to reliable alignment of all GPCR amino acid sequences under investigation. The amino acid fingerprints (patterns, motifs) taken from the PRINTS database²⁰ are specific enough to avoid ambiguous alignments. Out of the 277 sequences investigated herein, only 5 present an ambiguous pattern (Table 3). However, this accidental match always appears only once in a single sequence and does not preclude for TM and class detection. The number of ambiguous motifs is slightly higher (26), but these motifs are never found more than twice in the same sequence. This observation confirms previous results,²⁰ suggesting that a maximum of 2 motifs can be randomly found in a wrong GPCR class. Importantly, these discrepancies did not induce any failure in our family assignment. The GPCRmod alignment is different from a whole-sequence-based alignment³⁶ trying to find a unique consensus sequence³⁷ which is obviously wrong as far as GPCRs from different families are compared (Figure 11). A clear drawback of our approach is that GPCRmod neglects intra- and extracellular loops for which a few specific fingerprints could also be found.³⁸

High-Throughput GPCR Models for Studying Receptor–Antagonist Interactions. A single high-resolution X-ray structure of a GPCR (bovine rhodopsin) is currently available.⁹ Because rhodopsin has been crystallized in its ground (antagonist-bound) state, the current models only pretend to represent receptor structures aimed at studying receptor–antagonist interactions. Hence, most recent reports agree to suggest that the latter X-ray structure is a good template for mapping GPCRs even for receptors that significantly differ from rhodopsin.^{12,13} However, it must be noticed that modeling GPCRs in their activated states^{39–41} is clearly out of the scope of the present report for the simple

CONSENSUS ^a	PTNYFLNLAVADLLVALTLPPFALYYALL
OPSD_BOVIN ^b	PLNYILLNLAVADLFMVFGGFTTTLTSLH
CASR_HUMAN ^c	488 NYSIINWHLSPEDGSIVFKEVGYYNVYAKK
CASR_HUMAN ^d	642 IVKATNRELSYLLLFSLCCFSSSLFFIGE

Figure 11. Comparison of a fingerprint-based alignment (GPCRmod) with a full sequence alignment⁵⁹ of two GPCRs from different classes. The second transmembrane domain (TM2) of bovine rhodopsin (OPSD·BOVIN) is aligned to that of the human calcium-sensing receptor (CASR·HUMAN). Whereas the full sequence alignment is able to properly match the bovine rhodopsin sequence with that of a GPCR consensus sequence (LA..D common match), the alignment proposed for the calcium-sensing receptor is forced to match the LA..D consensus and thus align a sequence from the N-terminal extracellular domain (488–517) with the TM2 of bovine rhodopsin. By opposition, GPCRmod uses a different template for class III GPCRs and outputs a reliable alignment based on the class-III-specific TM2-characteristic K....E.SY pattern (see Table 1).

reason that the X-ray structure the bovine rhodopsin is unlikely to be a good template in that case⁴²

Direct threading of the target model onto the rhodopsin structure generates 3-D models that are not well-suited for virtual screening purpose.¹⁰ Hence, the dimensions of the TM binding cavity in rhodopsin is obviously biased by the cocrystallized covalently bound ligand (retinal). We thus decided to generate other 3-D templates for comparative modeling. Eight additional targets (Figure 4) have been selected because the corresponding high-quality 3-D models were in our hands and proved useful to either (i) explain fine details of receptor–antagonist interactions (V_{1A} receptor²² and extracellular calcium-sensing receptor¹³) or (ii) discriminate true ligands from randomly chosen “druglike” molecules in protein-based virtual screening tests (M_1 , D_3 , V_{1A} , and ORL-1 receptors).¹⁰ The backbone coordinates used to generate the target model is thus chosen from 9 possible templates which present a similar 3-D fold but slightly different helix bundle assemblies. Furthermore, the position-specific side-chain library generated from these 9 templates is diverse enough to find, in 75% of the cases, a target side chain present at the same position of the same TM segment in any of the 9 templates (Figure 8).

The herein described 3-D modeling strategy is not novel by itself. However, we believe that the use of rotamer libraries customized from several experimentally validated GPCR models is a clear advantage with respect to generic homology modeling procedures⁴³ which are parametrized from soluble proteins that are very different from membrane receptors. Last, our modeling protocol does not require the knowledge and prior docking of a known antagonist in order to extend the binding cavity of the target GPCR.

A clear drawback of the current study is that only TM domains have been taken into account to facilitate the amino acid sequence alignment. GPCRs also display extracellular residues (N-terminal domain, three loops) that may participate to ligand binding, notably for class I peptidergic receptors, hormone-binding class II receptors, and metabotropic-like GPCRs (class III). However, the scope of the current GPCR library is not to propose high-resolution all-atom 3-D models for all GPCRs but only models which are precise enough to identify receptor antagonists. As occupying the TM binding cavity is a common feature of most GPCR

antagonists,¹ we do think that models of the 7 TMs are sufficient to achieve this task. Furthermore, we believe that the high-throughput modeling of highly variable extracellular loops (especially the second one shown to fold back over the TM cavity in bovine rhodopsin)⁹ would provide more noise than real information because of the clear lack of structural data on the contribution of these loops to ligand binding. A recent site-directed mutagenesis of the second extracellular loop of the calcium-sensing receptor¹³ clearly shows that, despite a significant sequence identity to bovine rhodopsin, it cannot adopt the peculiar 3-D fold observed in the latter template.

Another limitation of the current modeling procedure is the omission of kinks, bends, and differential inclinations of helical axis specifically induced by either certain amino acids (proline or glycine)^{44,45} or non- α -helical 2-D structures (3₁₀-helices, π -helices).⁴⁶ For example, the chemokine CCR5 receptor presents a Thr-Xaa-Pro (Xaa being any amino acid) motif at the extracellular side of TM2, conserved all over the chemokine-receptor family and hypothesized to induce a bend critical for chemokine binding.⁴⁷ In the next release of our GPCR library in which we foresee adding 130 new GPCRs from the GPCRDB database,⁴⁸ we plan to use a rule-based method for bending TM helices after backbone selection of the template and accommodate, as much as possible, the above-described deformation.

We assume in the current alignment/building procedure a conservation of TM lengths for all GPCRs. It is possible that the beginning and the end of certain TM regions are slightly shifted for some peculiar GPCRs. However, as our alignment and building procedure is based on the existence of conserved amino acids at key points of the TM cavity, it is very unlikely that TM residues lining the antagonist binding site have been shifted. Discrepancies can only apply to capping residues of the TM helices that do not influence the antagonist binding site cavity. Because the primary goal of the present study is to provide a library of GPCR models presenting a consistent TM cavity, the above-noted discrepancy is unlikely to influence our usage of these models.

Last, incorporation of more class I templates and of at least one class II reference should allow a less biased backbone template selection. If we assume a conserved helix bundle assembly of the ground state for most GPCRs, as suggested by two recent site-directed-mutagenesis studies,^{12,13} the last bias should not influence that much the overall 3-D structure of the current high-throughput models.

Comparison of GPCRmod with Other GPCR Modeling Approaches. Various modeling strategies aimed at proposing reliable 3-D models of the most interesting GPCRs have already been described. They can be classified in 3 categories: *ab initio* folding techniques,^{49–51} distance–geometry based methods,⁵² and rhodopsin-based threading tools.^{36,53} *Ab initio* building tools present the advantage of being independent of any template and are currently able to reproduce the general architecture (helix bundle assembly) of GPCRs at a very low throughput. However, such models still are not accurate enough to guide a drug design approach. When applied to the prediction of the bovine rhodopsin structure, rms deviations of the *ab initio* model from the X-ray structure are typically larger than 3.0 Å.^{49–51} Mosberg et al. developed a distance–geometry-based method aimed at optimizing interhelical hydrogen bonds between buried

polar residues.⁵² The method has been applied to the construction of 26 GPCR models whose predicted binding pockets are in agreement with known experimental data. Whether the proposed models are able to discriminate known ligands from randomly chosen molecules has not been assessed. The herein described procedure is closer in its spirit to that used in WHAT IF³⁶ for generating rhodopsin-based homology models⁵³ (currently stored in the GPCRDB database⁵³). However, GPCRmod and WHAT IF models are significantly different for two main reasons: (i) Amino acid sequence alignment to rhodopsin are markedly different, especially in TM5 and TM6 which follow the variable third intracellular loop. This is a direct consequence of the different alignment methods used. (ii) GPCRmod models are derived from ligand-bound energy-minimized GPCR templates and not from the X-ray structure of rhodopsin itself. Thus, TM binding cavities of GPCRmod models are typically larger (less biased from the rhodopsin-bound retinal volume) than WHAT IF models.

Utility of the GPCR Target Library for Inverse Screening. The accuracy of the current GPCR models has been assessed by their ability to accommodate either a supposedly selective GPCR antagonist (MRS-2179, Figure 10A) or a known promiscuous ligand (NAN-190, Figure 10B) in cross-docking experiments. Remarkably, the true receptor(s) of both ligands is (are) ranked among the top 10% scorers in our inverse screening protocol. The proposed binding mode of MRS-2179 to the high-throughput model of the P2Y₁ receptor is in remarkable agreement with side-directed mutagenesis data.⁵⁴ The adenine moiety is embedded in an hydrophobic pocket delimited by TM6 and TM7, the ribose lies in another hydrophobic site between TM3 and TM6, and the diphosphate interacts through hydrogen-bond-assisted salt bridges to two arginine (Arg 68, Arg 310) and a lysine residue (Lys128). By comparison, docking the same antagonist to the WHAT IF model of the same receptor leads to a lower fitness score and a very different docking mode (between TM5 and TM6), which is not supported by known experimental data.⁵⁴ The promiscuous GPCR antagonist NAN-190 is similarly docked to all its known receptors (Table 5) with the basic nitrogen placed within a salt bridge distance to a conserved aspartic acid (TM3) as suggested by side-directed mutagenesis experiments.⁵⁵ The *n*-butyl spacer fills a gorge between TM3 and TM6, locating both aromatic moieties in two hydrophobic subsites (one between TM2, TM3, and TM7 and one between TM3 and TM5) in agreement with a previous model.⁵⁶

Several reasons may explain why the true receptor is not ranked first: (i) the full specificity profile of the two investigated antagonists is only partially known and the binding affinity of MRS-2179 and NAN-190 for their corresponding top-ranked GPCRs (Table 6) is still unknown, (ii) the contribution of amino acids from the extracellular loops (especially the second one) has been omitted in the current docking, and (iii) the fast scoring function utilized by GOLD cannot exactly reproduce binding free energies and consequently binding affinities.⁵⁷ However, it should be recalled that the main purpose of fast scoring functions is not to exactly predict absolute binding free energies (affinities) form protein–ligand coordinates but to be robust enough to clearly discriminate potential hits (targets in the present case) from unlikely solutions. As far as the investigated

Table 6. Top Scoring GPCRs for MRS-2171 and NAN-190

MRS-2171		NAN-190	
SwissProt Id ^a	rank	SwissProt Id ^a	rank
CCR5	1	CXC1	1
MAS	2	ACM4	2
IL8B	3	D2DR	3
B2AR	4	P2Y5	4
P2Y6	5	PAR1	5
5HT _{1F}	6	CCR5	6
P2YR	7	5H7	7
OPSB	8	HH1R	8
P2Y5	9	ACM2	9
GLP2	10	ACM1	10

^a CCR5, C-X-C chemokine receptor type 5; MAS, MAS protooncogene receptor; IL8B, interleukine 8 receptor type B; B2AR, β_2 adrenergic receptor; P2Y6, purinergic P2Y₆ receptor; 5HT_{1F}, serotonin 5-HT_{1F} receptor; P2YR, purinergic P2Y₁ receptor; OPSB, blue-sensitive opsin receptor; P2Y5, purinergic P2Y₅ receptor; GLP2, glucagon-like peptide 2 receptor; CXC1, chemokine XC receptor type 1; ACM4, acetylcholine muscarinic M₄ receptor; D2DR, dopamine D₂ receptor; PAR1, thrombin receptor; 5H7, serotonin 5-HT₇ receptor; HH1R, histamine H₁ receptor; ACM2, acetylcholine muscarinic M₂ receptor; ACM1, acetylcholine muscarinic M₁ receptor.

ligands have not been experimentally tested for binding to all GPCRs and especially the top-ranked receptors, it is impossible to unambiguously detect both false positives (overestimated receptors) and false negatives (underestimated receptors). However, we have to admit that the fine specificity for closely related GPCR subtypes is only partially addressed. For both ligands, ca. 80% of GPCRs closely related to the true target(s) (P2Y receptors for MRS-2171; 5-HT receptors for NAN-190) usually clustered in the top 20% scorers. Thus, the current inverse screening procedure is more aimed at identifying the likely receptor subfamily (dopamine, serotonin, adenosine, etc.) than precisely mapping the individual preference for highly related GPCR subtypes. It could thus be used as a computational filter to study the most likely targets when addressing the selectivity profile of a given compound or trying to identify the yet unknown receptor of a molecule showing promising in vivo biological effects.

CONCLUSIONS

The GPCRmod software package has been developed to enable the high-throughput modeling of most pharmaceutically interesting G protein coupled receptors. Starting from the amino acid sequence of 277 targets, an alignment of the seven transmembrane domains of all receptors is proposed on the basis of existing GPCR fingerprints. The multiple alignment has then been translated into reliable three-dimensional models using a knowledge-based threading algorithm based on eight different GPCR models and two side chains libraries. All models were concatenated into a target library of 277 receptors that can be screened electronically to identify the receptor(s) of known ligands. They were shown to be accurate enough to allow the recovery of the true targets among the top candidates. To the best of our knowledge, the current study is the first report of a direct use of high-throughput GPCR models for virtual screening purpose. The herein described comparative modeling and inverse screening procedures can easily be extended to other pharmaceutically interesting protein families (e.g. kinases,

phosphatases) as GPCRmod uses a target-independent Java library. More importantly, inverse screening of protein databases enables the "in silico" identification of the most plausible target(s) of any given ligand as well as the computation of ligand specificity profiles that might be used to assist lead development in a very early phase.

ACKNOWLEDGMENT

The authors thank the Strasbourg-Alsace-Lorraine Genopole, the Fondation pour la Recherche Médicale (FRM, Paris, France) for financial support and the Centre Informatique National de l'Enseignement Supérieur (CINES, Montpellier, France) for allocation of computing time. Prof. A. Krogh is acknowledged for providing the TMHMM code.

REFERENCES AND NOTES

- Schwalbe, H.; Wess, G. Dissecting G-Protein-Coupled Receptors: Structure, Function, and Ligand Interaction. *ChemBioChem* **2002**, *3*, 915–919.
- Wise, A.; Gearing, K.; Rees, S. Target Validation of G-Protein Coupled Receptors. *Drug Discovery Today* **2002**, *7*, 235–246.
- Venter, J. C.; et al. *Science* **2001**, *291*, 1304–1351.
- Gether, U. Uncovering Molecular Mechanisms Involved in Activation of G Protein-Coupled Receptors. *Endocr. Rev.* **2000**, *21*, 90–113.
- Fredriksson, R.; Lagerstrom, M. C.; Lundin, L. G.; Schioth, H. B. The G-Protein-Coupled Receptors in the Human Genome Form Five Main Families. Phylogenetic Analysis, Paralogon Groups, and Fingerprints. *Mol. Pharmacol.* **2003**, *63*, 1256–1272.
- Gasparini F.; Kuhn, R.; Pin, J. P. Allosteric Modulators of Group I Metabotropic Glutamate Receptors: Novel Subtype-Selective Ligands and Therapeutic Perspectives. *Curr. Opin. Pharmacol.* **2002**, *2*, 43–49.
- Klabunde, T.; Hessler, G. Drug Design Strategies for Targeting G-Protein-Coupled Receptors. *ChemBioChem* **2002**, *3*, 928–944.
- Walters, W. P.; Stahl, M. T.; Murcko, M. A. Virtual Screening—An Overview. *Drug. Discovery Today* **1998**, *3*, 160–178.
- Palczewski, K.; Kumasaka, T.; Hori, T.; Behnke, C.A.; Motoshima, H.; Fox, B.A.; Trong, I.L.; Teller, D.C.; Okada, T.; Stenkamp, R.E.; Yamamoto, M.; Miyano, M. Crystal Structure of Rhodopsin: A G Protein-Coupled Receptor. *Science* **2000**, *289*, 739–745.
- Bissantz, C.; Bernard, P.; Hibert, M.; Rognan, D. Virtual Screening of Chemical Databases. II. Are Homology Models of G-Protein Coupled Receptors Suitable Targets? *Proteins: Struct., Funct., Genet.* **2003**, *50*, 5–25.
- Becker, O. M.; Shacham, S.; Marantz, Y.; Noiman, S. Modeling the 3D Structure of GPCRs: Advances and Application to Drug Discovery. *Curr. Opin. Drug. Discovery Dev.* **2003**, *6*, 353–361.
- Malherbe, P.; Kratochwil, N.; Zenner, M. T.; Piusi, J.; Diener, C.; Kratzeisen, C.; Fischer, C.; Porter, R. H. Mutational Analysis and Molecular Modeling of the Binding Pocket of the Metabotropic Glutamate 5 Receptor Negative Modulator 2-Methyl-6-(phenylethynyl)pyridine. *Mol. Pharmacol.* **2003**, *64*, 823–832.
- Petrel, C.; Kessler, A.; Maslah, F.; Dauban, P.; Dodd, R. H.; Rognan, D.; Ruat, M. Modeling and Mutagenesis of the Binding Site of Calx 231, a Novel Negative Allosteric Modulator of the Extracellular Ca²⁺ Sensing Receptor. *J. Biol. Chem.* **2003**, *278*, 49487–49494.
- Krogh A.; Larsson B.; von Heijne, G.; Sonnhammer, E. L. L. Predicting Transmembrane Protein Topology with a Hidden Markov Model: Application to Complete Genomes. *J. Mol. Biol.* **2001**, *305*, 567–580.
- Moller, S.; Croning, M. D.; Apweiler, R. Evaluation of Methods for the Prediction of Membrane Spanning Regions. *Bioinformatics* **2001**, *17*, 646–653.
- Attwood, T. K. A Compendium of Specific Motifs for Diagnosing GPCR Subtypes. *Trends Pharmacol. Sci.* **2001**, *22*, 162–165.
- Henikoff, J. G.; Henikoff, S. Using Substitution Probabilities To Improve Position-Specific Scoring Matrices. *Comput. Appl. Biosci.* **1996**, *12*, 135–143.
- Thompson, J. D.; Higgins, D. G.; Gibson, T. J. CLUSTAL W: Improving the Sensitivity of Progressive Multiple Sequence Alignment through Sequence Weighting, Position-Specific Gap Penalties and Weight Matrix Choice. *Nucleic Acids Res.* **1994**, *22*, 4673–4680.
- Schwartz, T. W. Locating Ligand-Binding Sites in 7 TM Receptors by Protein Engineering. *Curr. Opin. Biotechnol.* **1994**, *5*, 434–444.
- Attwood, T. K.; Blythe, M.; Flower, D. R.; Gaulton, A.; Mabey, J. E.; Maudling, N.; McGregor, L.; Mitchell, A.; Moulton, G.; Paine,

- K.; Scordis, P. PRINTS and PRINTS-S Shed Light on Protein Ancestry. *Nucleic Acid Res.* **2002**, *30*, 239–241.
- (21) Henikoff, S.; Henikoff, J. G. Amino Acid Substitution Matrices from Protein Blocks *Proc. Natl. Acad. Sci. U.S.A.* **1992**, *89*, 10915–10919.
- (22) Tahtaoui, C.; Balestre, M.-N.; Klotz, P.; Rognan, D.; Barberis, C.; Mouillac, B.; Hibert, M. Identification of the Binding Sites of the SR49059 Nonpeptide Antagonist into the V1A Vasopressin Receptor Using Sulfidryl-Reactive Ligands and Cysteine Mutants as Chemical Sensors. *J. Biol. Chem.* **2003**, *278*, 40010–40019.
- (23) Guex, N.; Peitsch, M. C. SWISS-MODEL and the Swiss-PdbViewer: An Environment for Comparative Protein Modeling. *Electrophoresis* **1997**, *18*, 2714–2723.
- (24) Lovell, S. C.; Word, J. M.; Richardson, J. S.; Richardson, D. C. The Penultimate Rotamer Library. *Proteins: Struct., Funct., Genet.* **2000**, *40*, 389–408.
- (25) Case, D. A.; Pearlman, D. A.; Caldwell, J. W.; Cheatham, D. E., III; Ross, W. S.; Simmerling, C. L.; Darden, T. A.; Merz, K. M.; Stanton, R. V.; Cheng, A. L.; Vincent, J. J.; Crowley, M.; Tsui, V.; Radmer, R. J.; Duan, Y.; Pitera, J.; Massova, I.; Seibel, G. L.; Singh, U. C.; Weiner, P. K.; Kollman, P. A. *AMBER 6*; Univeristy of California: San Francisco, CA, 1999.
- (26) Grantham, R. Amino Acid Difference Formula To Help Explain Protein Evolution. *Science* **1974**, *185*, 862–864.
- (27) Jones, G.; Wilett, P.; Glen, R. C.; Leach, A. R.; Taylor, R. Development and Validation of a Genetic Algorithm for Flexible Docking. *J. Mol. Biol.* **1997**, *267*, 727–748.
- (28) SYBYL 6.81; Tripos Inc.: St. Louis, MO.
- (29) <http://www.mdli.com>, MDL Information Systems, Inc., San Leandro, CA.
- (30) Concord 4.0 is part of the SYBYL software distribution (<http://www.tripos.com>).
- (31) Clark, M.; Cramer, R. D., III; Van Opdenbosch, N. Validation of the General Purpose TRIPOS 5.2 Force Field. *J. Comput. Chem.* **1989**, *10*, 982–1012.
- (32) Boeckmann, B.; Bairoch, A.; Apweiler, R.; Blatter, M.-C.; Estreicher, A.; Gasteiger, E.; Martin, M. J.; Michoud, K.; O'Donovan, C.; Phan, I.; Pilboud, S.; Schneider, M. The SWISS-PROT Protein Knowledgebase and Its Supplement TrEMBL in 2003. *Nucleic Acids Res.* **2003**, *31*, 365–370.
- (33) Morris, A. L.; MacArthur, M. W.; Hutchinson, E. G.; Thornton, J. M. Stereochemical Quality of Protein Structure Coordinates. *Proteins: Struct., Funct., Genet.* **1992**, *12*, 345–364.
- (34) Boyer, J. L.; Mohanram, A.; Camaioni, E.; Jacobson, K. A.; Harden, T. K. Competitive and Selective Antagonism of P2Y1 Receptors by N6-Methyl 2'-Deoxyadenosine 3',5'-Bisphosphate *Br. J. Pharmacol.* **1998**, *124*, 1–3.
- (35) Glennon, R. A.; Naiman, N. A.; Pierson, M. E.; Titeler, M.; Lyon, R. A. NAN-190: An Arylpiperazine Analogue That Antagonizes the Stimulus Effects of the 5-HT_{1A} Agonist 8-Hydroxy-2-(di-*n*-propyl-amino)tetralin. *Eur. J. Pharmacol.* **1988**, *154*, 339–341.
- (36) Vriend, G. WHAT IF: A Molecular Modeling and Drug Design Program. *J. Mol. Graphics* **1990**, *8*, 52–56.
- (37) Oliveira, L.; Paiva, A. C.; Vriend, G. A Common Motif in G Protein-Coupled Seven Transmembrane Helix Receptors. *J. Comput.-Aided Mol. Des.* **1993**, *7*, 649–658.
- (38) Attwood, T. K.; Findlay, J. B. C. Fingerprinting G-Protein-Coupled Receptors. *Protein Eng.* **1994**, *7*, 195–203.
- (39) Gether, U.; Kobilka, B. K. G Protein-Coupled Receptors. II. Mechanism of Agonist Activation. *J. Biol. Chem.* **1998**, *273*, 17979–17982.
- (40) Hulme, E. C.; Lu, Z. L.; Ward, S. D. C.; Allman, K.; Curtis, C. A. The Conformational Switch in 7-Transmembrane Receptors: The Muscarinic Receptor Paradigm. *Eur. J. Pharmacol.* **1999**, *375*, 247–260.
- (41) Ghanouni, P.; Steenhuis, J. J.; Farrens, D. L.; Kobilka, B. K. Agonist-Induced Conformational Changes in the G-Protein-Coupling Domain of the β_2 Adrenergic Receptor. *Proc. Natl. Acad. Sci. U.S.A.* **2001**, *98*, 5997–6002.
- (42) Archer, E.; Maigret, B.; Escricut, C.; Pradayrol, L.; Fourmy, D. Rhodopsin Crystal: New Template Yielding Realistic Models of G-Protein-Coupled Receptors? *Trends Pharmacol. Sci.* **2003**, *24*, 36–40.
- (43) John, B.; Sali, A. Comparative Protein Structure Modeling by Iterative Alignment, Model Building and Model Assessment. *Nucleic Acids Res.* **2003**, *31*, 3982–3992.
- (44) Chakrabarti, P.; Chakrabarti, S. C-H...O Hydrogen Bond Involving Proline Residues in Alpha-Helices. *J. Mol. Biol.* **1998**, *284*, 867–873.
- (45) Punta, M.; Maritan, A. A Knowledge-Based Scale for Amino Acid Membrane Propensity. *Proteins: Struct., Funct., Genet.* **2003**, *50*, 114–121.
- (46) Riek, R. P.; Rigoutsos, I.; Novotny, J.; Graham, R. M. Non-alpha-helical Elements Modulate Polytopic Membrane Protein Architecture. *J. Mol. Biol.* **2001**, *306*, 349–362.
- (47) Govaerts, C.; Bondue, A.; Springael, J.-Y.; Olivella, M.; Deupi, X.; Le Poul, E.; Wodak, S. J.; Parmentier, M.; Pardo, L.; Blanpain, C. Activation of CCR5 by Chemokines Involves and Aromatic Cluster between Transmembrane Helices 2 and 3. *J. Biol. Chem.* **2003**, *278*, 1892–1903.
- (48) Hom, F.; Bettler, E.; Oliveira, L.; Campagne, F.; Cohen, F. E.; Vriend, G. GPCR Information System for G Protein-Coupled Receptors. *Nucleic Acid Res.* **2003**, *31*, 294–297.
- (49) Shacham, S.; Topf, M.; Avisar, N.; Glaser, F.; Marantz, Y.; Bar-Haim, S.; Noiman, S.; Naor, Z.; Becker, O. M. Modeling the 3D Structure of GPCRs from Sequence. *Med. Res. Rev.* **2001**, *21*, 472–483.
- (50) Nikiforovitch, G. V.; Galaktionov, S.; Balodis, J.; Marshall, G. R. Novel Approach to Computer Modeling of Seven-Helical Transmembrane Proteins: Current Progress in the Test Case of Bacteriorhodopsin. *Acta Biochim. Pol.* **2001**, *48*, 53–64.
- (51) Vaidehi, N.; Floriano, W. B.; Trabanino, R.; Hall, S. E.; Freddolino, P.; Choi, E. J.; Zamanakos, G.; Goddard, W. A., III. Prediction of Structure and Function of G Protein-Coupled Receptors. *Proc. Natl. Acad. Sci. U.S.A.* **2002**, *99*, 12622–12627.
- (52) Lomize, A. L.; Pogozheva, I. D.; Mosberg, H. I. Structural Organization of G-Protein-Coupled Receptors. *J. Comput.-Aided Mol. Des.* **1999**, *13*, 325–353.
- (53) <http://www.gpcr.org/7tm/models/vriend3/index.html>.
- (54) Moro, S.; Guo, D.; Camaioni, E.; Boyer, J. L.; Harden, T. K.; Jacobson, K. A. Human P2Y1 Receptor: Molecular Modeling and Site-Directed Mutagenesis as Tools To Identify Agonist and Antagonist Recognition Sites. *J. Med. Chem.* **1998**, *41*, 1456–1466.
- (55) Ho, B. Y.; Karschin, A.; Branchek, T.; Davidson, N.; Lester, H. A. The Role of Conserved Aspartate and Serine Residues in Ligand Binding and in Function of the 5-HT_{1A} Receptor: A Site-Directed Mutation Study. *FEBS Lett.* **1992**, *312*, 259.
- (56) Bronowska, A.; Chilmonezyck, Z.; Les, A.; Edvardson, O.; Ostensen, R.; Sylte, I. Molecular Dynamics of 5-HT_{1A} and 5-HT_{2A} Serotonin Receptors with Methylated Bupirone Analogues. *J. Comput.-Aided Mol. Des.* **2001**, *15*, 1005–1023.
- (57) Bissantz, C.; Folkers, G.; Rognan, D. Protein-Based Virtual Screening of Chemical Databases. 1. Evaluation of Different Docking/Scoring Combinations. *J. Med. Chem.* **2000**, *43*, 4759–4767.
- (58) Roth, B. L.; Kroeze, W. K.; Patel, S.; Lopez, E. The Multiplicity of Serotonin Receptors: Uselessly Diverse Molecules or an Embarrassment of Riches? *Neuroscientist* **2000**, *6*, 252–262.
- (59) <http://www.gpcr.org/seq/200/200.MSF.mview.html>.

CI034181A