### **SNP analysis and binding site prediction**



Bioinformatics Department Prince Felipe Resarch Center (CIPF), Valencia, Spain Marc A. Marti-Renom



Thursday, September 23, 2010

# Program



# Objective

### TO UNDERSTAND THAT SNPs HAVE EFFECTS THAT CAN BE PREDICTED AND TO LEARN HOW-TO USE AutoDock FOR DOCKING SMALL MOLECULES IN THE SURFACE OF A PROTEIN

# Nomenclature

**SNP:** Single Nucleotide Polymorphism. A single change in the DNA sequence, which may or may not result in a change in the protein sequence.

Ligand: Structure (usually a small molecule) that binds to the binding site.

**Receptor**: Structure (usually a protein) that contains the active binding site.

**Binding site**: Set of aminoacids (residues) that physically interact with the lingad (usually within 6 Ångstroms).



#### Gene Sequence << +Protein Sequence << +Protein Structure

# Single Nucleotide Polymorphism

#### Single Nucleotide Polymorphism or SNP

is a DNA sequence variation occurring when a single nucleotide - A, T, C, or G - in the genome differs between members of the species. Usually one will want to refer to SNPs when the population frequency is  $\geq 1\%$ 

SNPs occur at any position and can be classified on the base of their locations.

Coding SNPs can be subdivided into two groups:

Synonymous: when single base substitutions do not cause a change in the resultant amino acid

Non-synonymous: when single base substitutions cause a change in the resultant amino acid.



http://www.ncbi.nlm.nih.gov

# **SNPs and disease**

Single nucleotide polymorphism are the most common type of genetic variations in human accounting for about 90% of sequence differences (Collins et al., 1998).

Studying SNPs distribution in different human populations can lead to important considerations about the history of our species (Barbujani and Goldstein, 2004; Edmonds et al., 2004).

SNPs can also be responsible of genetic diseases (Ng and Henikoff, 2002; Bell, 2004).



# **SNP databases**





## **Evolution and disease.**

Capriotti et al. Use of estimated evolutionary strength at the codon level improves the prediction of disease-related protein mutations in humans. Hum Mutat (2008) vol. 29 (1) pp. 198-204







PRINCIPE FELIPE CENTRO DE INVESTIGACION

# Sequence and evolutive - based predictors



Profile: MR and AS sequence profile information

Codon: omega, dS,dN: selective pressure at codon level, synonymous and non-synonymous rate at branch level.

# **Omidios method**

Omidios has higher accuracy than the previous two methods increasing the accuracy up to 82% and the correlation coefficient to 0.59.

	Q2	P[D]	Q[D]	P[N]	Q[N]	С
Omidios	82	88	84	68	76	0.59



Q2: Overall Accuracy C: Correlation Coefficient DB: Fraction of database that are predicted with a reliability ≥ the given threshold

## Comparison

Omidios results in higher accuracy and correlation than the other available methods covering the 100% of the dataset (see column %PM).

Omidios results in higher accuracy with respect to SIFT and although the quality of Omidios is comparable to PANTHER, when our prediction are selected by RI index the accuracy of our method is higher than PANTHER.

	Q2	P[D]	Q[D]	P[N]	Q[N]	С	PM
Omidios	82	89	84	68	76	59	100
SIFT	71	84	72	51	69	38	97
PANTHER	74	87	75	53	72	43	83

HM-Dic05: 8987 mutations

	Q2	P[D]	Q[D]	P[N]	Q[N]	С	PM
Omidios	74	65	79	83	72	48	100
SIFT	71	63	70	78	72	42	96
PANTHER	77	73	71	79	81	52	77

HM-Dic06: 2008 mutations

### **Omidios server**

#### http://sgu.bioinfo.cipf.es/services/Omidios

The Omidios server		
▲ ►	📀 ^ Q JMB Dopazo	0
Omidios (a.k.a SeqProfCod) ③   SWISS-PROT id:   Submit   Example: AQP2_HUMAN   HELP:  PLEASE NOTE. Our servers have been optimized for Firefox and Safari. If you are using Internet Explorer, the CSS may not be properly rendered.  To use Omidios you need to: - Enter a SWISS-PROT id of the sequence of interest.  NECC.	(SGU - HOME) DBAJI Eva-CM Omidios SARA TDIKernel	
The Omidios server is designed to take a query SWISS-PROT id and search for all annotated and predicted protein variants (nsSNP) in our database. The whole set of predictions is available for downloading: - MySQL dump format. - Tab separated format. Individual training and testing datasets for SeqProfCod are also available for downloading: - SP-Dec05 dataset. - SP-Dec06 dataset.		
13		

### Structural analysis of missense mutations in human BRCA1 BRCT domains

Mirkovic et al. Structure-based assessment of missense mutations in human BRCA1: implications for breast and ovarian cancer predisposition. Cancer Res (2004) vol. 64 (11) pp. 3790-7

#### ICANCER RESEARCH 64, 3790-3797, June 1, 2004

#### Structure-Based Assessment of Missense Mutations in Human BRCA1: Implications for Breast and Ovarian Cancer Predisposition

#### Nebojsa Mirkovic,<sup>1</sup> Marc A. Marti-Renom,<sup>2</sup> Barbara L. Weber,<sup>3</sup> Andrej Sali,<sup>2</sup> and Alvaro N. A. Monteiro<sup>4,5</sup>

<sup>1</sup>Laboratory of Molecular Biophysics, Pels Family Center for Biochemistry and Structural Biology, Rockefeller University, New York, New York; <sup>2</sup>Departments Ladvanny of docume unphysics, very landy Centre pol methodismity und automat doorsy, tookseptiter outering, very 16st, teer 16st, etcer 16st, Centre and California Istantian for Quantitative Biomedical Research California Istantian, tork Biopharmaeening Sciences and Pharmaceutical Consisty, and California Institute Grounditative Biomedical Research Francisco, California, 'Abramson Family Cancer Research Institute, University of Pensystemia, Philadelphia, Pensystemia, 'Srang Cancer Prevention Center, New York, New York, and 'Department of Cell and Developmental Biology, Weill Medical College Orael Diriversity, New York, New York

#### ABSTRACT

can be screened for the presence of mutations. However, the cancer association of most alleles carrying missense mutations is unknown, thus creating significant problems for genetic counseling. To increase our ability to identify cancer-associated mutations in BRCA1, we set out to use the principles of protein three-dimensional structure as well as the correlation between the cancer-associated mutations and those that abolish These observations suggest that abolishing the transcriptional activatranscriptional activation. Thirty-one of 37 missense mutations of known tion function of BRCA1 leads to tumor development and provides a impact on the transcriptional activation function of BRCA1 are readily rationalized in structural terms. Loss-of-function mutations involve non conservative changes in the core of the BRCA1 C-terminus (BRCT) fold or are localized in a groove that presumably forms a binding site involved in the transcriptional activation by BRCA1; mutations that do not abolish transcriptional activation are either conservative changes in the core or are on the surface outside of the putative binding site. Next, structurebased rules for predicting functional consequences of a given missense mutation were applied to 57 germ-line BRCA1 variants of unknown cancer association. Such a structure-based approach may be helpful in an integrated effort to identify mutations that predispose individuals to cancer.

#### INTRODUCTION

Many germ-line mutations in the human BRCA1 gene are associ-The multiple structure-based alignment of the native structures of the ated with inherited breast and ovarian cancers (1, 2). This information has allowed clinicians and genetic counselors to identify individuals at high risk for developing cancer. However, the disease association of voer 350 missense mutations remains unclear, primarily because their relatively low frequency and ethnic specificity limit the usefulness of protein (1KZY; Ref. 7), human p53-binding the population-based statistical approaches to identifying cancer-causing mutations. To address this problem, we use here the threeing mutations. Io address this problem, we use nere the inree-dimensional structure of the human BRCA1 BRCT domains to assess the transcriptional activation functions of BRCA1 mutants. Our study is made nossible by the recently determined sequences (3–6) and is made possible by the recently determined sequences (3-6) and three-dimensional structures of the BRCA1 homologs (7, 8). In addition, we benefited from prior studies that attempted to rationalize and predict functional effects of mutations in various proteins (9-12), three-dimensional model for each of the 94 mutants. The crystallographic including those of BRCA1 (13, 14).

tates DNA damage repair (15, 16). The tandem BRCT domains at the structure (1694 and 1817–1819) were modeled *de novo* (27). All of the models

Received 9/24/03; revised 1/30/04; accepted 3/15/04. Grant support: This work was supported by Lee Kaplan Foundation, the Fashion Footward Association of New York/QVC; United States Army award DAMD[7:99-1. 989 and NH A OS2200 (A.N. A.M.); the Mathers Foundation, Sandler Family States Army award DAMD[7:99-1. 1880 and NH A OS2200 (A.N. A.M.); the Mathers Foundation, Sandler Family States Army award DAMD[7:99-1. 1980 and NH A OS2200 (A.N. A.M.); the Mathers Foundation, Sandler Family States Army award DAMD[7:99-1. 1980 and NH A OS2200 (A.N. A.M.); the Mathers Foundation, Sandler Family States Army award DAMD[7:99-1. 1990 And Mathers Foundation, Sandler Family States Army award DAMD[7:99-1. 1990 And Mathers Army award DAMD[7:99-1. 1991 And Mathers Army award DAMD[7:99-1. 1992 And Mathers Army award DAMD[7:99-1. 1991 And Mathers Army award DAMD[7:99-1. 1991 And Mathers Army award DAMD[7:99-1. 1992 And Mathers Army award DAMD[7:99-1. 1994 And Mathers Army award Army award DAMD[7:99-1. 1994 And Mathers Army award Ar

COOH-terminus of BRCA1 are involved in several of its functions including modulation of the activity of several transcription factors The *BRCAI* gene from individuals at risk of breast and ovarian cancers (15), binding to the RNA polymerase II holoenzyme (17), and activating transcription of a reporter gene when fused to a heterologous DNA-binding domain (18, 19). Importantly, cancer-associated mutations in the BRCT domains, but not benign polymorphisms, inactivate transcriptional activation and binding to RNA polymerase II (18-21) genetic framework for characterization of BRCA1 BRCT variants.

#### MATERIALS AND METHODS

The multiple sequence alignment (MSA) of orthologous BRCA1 BRCT domains from seven species, including Homo sapiens (GenBank accession number U14680), Pan troglodytes (AF207822), Mus musculus (U68174), Rattus norvegicus (AF036760), Gallus gallus (AF355273), Canis familiaris (U50709), and Xenopus laevis (AF416868), was obtained by using program ClustalW (22) and contains only one gapped position (Supplementary Fig. 1). According to PSI-BLAST (23), the latter six sequences are the only sequence in the nonredundant protein sequence database at National Center for Biotech nology Information that have between 30% and 90% sequence identity to the

human BRCA1 BRCT domains (residues 1649–1859). The multiple structure-based alignment of the native structures of the XRCC1 protein (1CDZ; Ref. 13). Structure variability was defined by the -square deviation among the superposed  $C\alpha$  positions, as calculated may point to putative functional site(s) on the surface of BRCT repeats.

Comparative protein structure modeling by satisfaction of spatial restraints, implemented in the program MODELLER-6 (26), was used to produce a structure of the human wild-type BRCA1 BRCT domains was used as the BRCA1 is a nuclear protein that activates transcription and facili-template for modeling (8). The four residues missing in the crystallographic are available in the BRCA1 model set deposited in our ModBase database of comparative protein structure models (28).6

For the native structure of the human BRCT tandem repeat and each of the 94 mutant models, a number of sequence and structure features were calculated. These features were used in the classification tree in Fig. 3 (values for

porting Foundation, Sun, IBM and Intel (A. S.); and NIH (MI 54762 (IM61390 (A. S.); and the Resart Cancer Research Status of the Status of the

3790



Received 9/24/03; revised 1/30/04; accepted 3/15/04

### Human BRCA1 and its two BRCT domains



Williams, Green, Glover. Nat.Struct.Biol. 8, 838, 2001

CONFIDENTIAL



#### BRACAnalysis <sup>™</sup> Comprehensive BRCA1-BRCA2 Gene Sequence Analysis Result



#### Interpretation

#### GENETIC VARIANT OF UNCERTAIN SIGNIFICANCE

The BRCA2 variant H2116R results in the substitution of arginine for histidine at amino acid position 2116 of the BRCA2 protein. Variants of this type may or may not affect BRCA2 protein function. Therefore, the contribution of this variant to the relative risk of breast or ovarian cancer cannot be established solely from this analysis. The observation by Myriad Genetic Laboratories of this particular variant in an individual with a deleterious truncating mutation in BRCA2, however, reduces the likelihood that H2116R is itself deleterious.

Authorized Signature:

Thomas S. Frank, M.D. Medical Director

Brian E. Ward, Ph.D. Laboratory Director

These testrepuls should only be used in conjunction with the pacent's clinical history and any previous analysis of appropriate family members. It is strongly recommended that these results be communicated to the patient in a setting that includes appropriate family. The accomptonying Technical Specifications summary describes the analysis, method, performance characteristics, nomenclanue, and interpretive optima of this test. This test may be considered investigational by some states. This test was developed and its performance characteristics determined by Myriad Genetic Laboratories. It has not been reviewed by the U.S. Food and Orug Administration. The FDA has determined that such clearance or approval is not necessary.

### **Missense mutations in BRCT domains by function**

	cancer associated	not cancer associated	?
no transcription activation	C1697R R1699W A1708E S1715R P1749R M1775R		M1652K       L1705PS       F1761S         L1657P       L1705PS       M1775E         E1660G       1715NS1       M1775K         H1686Q       722FF17       L1780P         R1699Q       34LG173       I1807S         K1702E       8EG1743       V1833E         Y1703HF       RA1752P       A1843T         1704S       F1761I       F1761I
transcription activation		M1652I A1669S	V1665M D1692N G1706A D1733G M1775V P1806A
?			M1652TW1718SR1751PC1787SA1823TV1653MT1720AR1751QG1788DV1833ML1664PW1730SR1758GG1788VW1837RT1685AF1734SL1764PG1803AW1837GT1685IE1735KL1766SV1804DS1841NM1689RV1736AI1766SV1808AA1843PD1692YG1738RP1771LV1809AT1852SF1695LD1739ET1773SV1809FP1856TV1696LD1739GP1776SV1810GP1859RR1699LD1739YD1778NQ1811RP1859RW1718CH1746ND1778GP1812SN1819S



### **Putative binding site on BRCA1**



Williams et al. 2004 Nature Structure Biology. June 2004 11:519

Mirkovic et al. 2004 Cancer Research. June 2004 64:3790

### Supervised learning approach

Karchin et al. Functional Impact of Missense Variants in BRCA1 Predicted by Supervised Learning. PLoS Comput Biol (2007) vol. 3 (2) pp. e26

OPEN a ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

#### Functional Impact of Missense Variants in BRCA1 Predicted by Supervised Learning

Rachel Karchin<sup>1,2\*</sup>, Alvaro N. A. Monteiro<sup>3</sup>, Sean V. Tavtigian<sup>4</sup>, Marcelo A. Carvalho<sup>3</sup>, Andrej Sali<sup>5,6\*</sup>

1 Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland, United States of America, 2 Institute of Computational Medicine, Johns Hopkins Lorepartment of balmetoria chypterenity, Johns robustis University, Balmider, marylands, University, Baltimore, Maryland, United States of America, **3** Risk Assessment, Detection, and Intervention Program, H. Lee Moffitt Cancer Center and Research Institute, Tampa, Forliad, United States of America, **4** International Agency for Research on Cancer, Lyon, France, **5** Department of Pharmaceutical Chemistry, University of California San Francisco, San Francisco, California, United States of America, **6** California Institute for Quantitative Biomedical Research, University of California San Francisco, San Francisco, San Francisco, California, United States of America, **6** California Institute for Quantitative Biomedical Research, University of California San Francisco, San Francisco, California, United States of America

Many individuals tested for inherited cancer susceptibility at the BRCA1 gene locus are discovered to have variants o unknown clinical significance (UCVs). Most UCVs cause a single amino acid residue (missense) change in the BRCA1 protein. They can be biochemically assayed, but such evaluations are time-consuming and labor-intensive. Computational methods that classify and suggest explanations for UCV impact on protein function can complement functional tests. Here we describe a supervised learning approach to classification of *BRCA1* UCVs. Using a novel combination of 16 predictive features, the algorithms were applied to retrospectively classify the impact of 36 BRCA1 C-terminal (BRCT) domain UCVs biochemically assayed to measure transactivation function and to blindly classify 54 documented UCVs. Majority vote of three supervised learning algorithms is in agreement with the assay for more than 94% of the UCVs. Two UCVs found deleterious by both the assay and the classifiers reveal a previously uncharacterized putative binding site. Clinicians may soon be able to use computational classifiers such as those described here to better inform patients. These classifiers can be adapted to other cancer susceptibility genes and systematically applied to prioritize the growing number of potential causative loci and variants found by large-scale disease association studies.

Citation: Karchin R, Monteiro ANA, Tavtigian SV, Carvalho MA, Sali A (2007) Functional impact of missense variants in BRCA1 predicted by supervised learning. PLoS Comput Biol 3(2): e26. doi:10.1371/journal.pcbi.0030026

0268

#### Introduction

The BRCA1 gene encodes a large multifunction protein involved in cell-cycle and centrosome control, transcriptional counseling their patients. regulation, and in the DNA damage response [1-3]. Inherited mutations in this gene have been associated with an increased lifetime risk of breast and ovarian cancer (6–8 times that of the general population) [4]. There are several thousand known deleterious BRCA1 mutations that result in frameshifts and/or premature stop codons, producing a truncated protein product [5]. In contrast, the functional impact of ost missense variants that result in a single amino acid ing both conservation and location of variant amino acid residue change in BRCA1 protein is not known. The Breast residues in an X-ray crystal structure [12]. Variants were Cancer Information Core database (http://research.nhgri.nih.- predicted deleterious if their properties were similar to gov/bic/), a central repository of BRCA1 and BRCA2 mutations identified in genetic tests, currently contains 487 unique missense BRCA1 variants (April 2006), of which only 17 have sufficient genetic/epidemiological evidence to be classified as deleterious (Clinically Important) and 33 as neutral or of little clinical importance (Not Clinically Important). As genetic 28, 2006 (doi:10.1371/journal.pcbi.0030026.eor). testing for inherited disease predispositions becomes more commonplace, predicting the clinical significance of missense variants and other UCVs will be increasingly important for risk assessment.

Because most UCVs in BRCA1 and BRCA2 occur at very low population frequencies (<0.0001) [6], direct epidemiological easures, such as familial cosegregation with disease, are often not sufficiently powerful to identify the variants

protein function and bioinformatics analysis [6-8]. In the future, physicians and genetic counselors may be able to rely on all these sources of information about UCVs when Previous bioinformatics analysis of BRCA1 UCVs has

related proteins from other organisms [9-11]. Two groups have attempted to include information about BRCA1 protein structure. Williams et al. predicted the impact of 25 missense variants in BRCA1's C-terminal BRCT domains by consider-

Editor: Greg Tucker-Kellogg, Lilly Systems Biology, Singapore Received September 5, 2006; Accepted December 27, 2006; Published February 16, 2007

Copyright: © 2007 Karchin et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Abbreviations: Align-GVGD, Align Grantham Variation Grantham Deviation; AUC, area under the ROC curve; BIC breast information core database; BBCT, BRCA1 C-terminal domain; BBCT-L, BBCT - Leterminal domain; BBCT-N BRCT Naternial domain; GD, Grantham Deviation; GV, Grantham Variation; ROC, receiver operating characteristic; Rule-based decision tree, empirically derived rules encoded in a decision tree; SIFT, Sorting Intolerant from Tolerant; UCV, variant of unknown Patient displayment.

associated with cancer predisposition. A promising approach clinical significance is to supplement epidemiological and clinical analysis of UCVs with indirect approaches such as biochemical studies of (RK; sul@sallab.org (AS)

PLoS Computational Biology | www.ploscompbiol.org

February 2007 | Volume 3 | Issue 2 | e26

(50) PLo5 Computational Biology (www.plascomphiclarg. 0205

•	B	В	В	В	倒	E	B	В	B	B	B	В	В	В	B	В	B	В	
	1655 S-→F	1664 L <b>→</b> P	1665 V <b>→</b> M	1669 A→S	1685 T <b>-</b> >I	1692 D→N	1689 M→R	1697 C→R	1699 R→W	1699 R-∋Q	1700 T→A	1706 G→E	1706 G→A	1708 A→E	1713 V→A	1715 S-→R	1715 S→N	1718 W-→C	
N S	•	0	0		•							•							
R A		Ŏ	Ö	Ö		00					ě				Ó				
U	Ŏ		ŏ		ě	ŏ	ě	ě	ě	ě	ě	ě	ě	ě	Ŏ	ě	ě	ě	
F	ž	ò	ĕ		ě	ŏ				ě	ě		0	Ŏ	ě		0	ě	
D	ŏ	ŏ	ŏ	ĕ	ŏ	Ō	ŏ	ĕ	ĕ	ŏ	ŏ	ŏ	ŏ	ĕ	ŏ	ŏ	ĕ	ĕ	
	E	B	В	В	В	E	B	В	B	B	E	⟨₿	∖₿	E	Ē	Ē	B	В	
	1720 T <del>-}</del> A	1736 V→A	1738 G→R	1738 G→E	1752 A-→P	1753 R <del>→</del> T	1764 L→P	1766 I <del>.)</del> S	1775 M <del>→</del> R	1775 M→K	1785 Q <del>.)</del> H	1788 G→V	1788 G→D	1794 E <del>→</del> D	1804 V→D	1806 P→A	1809 V <b>→</b> F	1837 W→R	MCC 0 0.5
N S														0					N S
R A		•			•									0					R A
U T		•			00		00								0				U T
F B		•	•	0				0	•	•				0		0	0		F B
D			0	0	0			•			0			0		0			D
					)-shee bend H-bond oop x-helix	t ded tur	B bu E ex	ried posed	(	Correc	tly or ir tly or ir D Insu		tly clas tly clas confid	sified a sified a ence to	as dele as neu o class	eteriou ıtral sify	S		

20

### Predictors are combined in support vector machine supervised learning

3.1

120

30

8

1.5

0.12

21

7

45

-3.2

17

1.5

#### X<sub>1</sub>..X<sub>k</sub>= TRAINING

relative entropy
Grantham
distance
solvent
accessibility
methyl(ene)
groups
volume change

\_\_\_







### Features

Feature Category	Feature Description
Structural	Solvent Accessiblity of wild-type amino acid residue (A <sup>2</sup> )
	Solvent Accessibility of wild-type residue normalized by maximum exposed Sol-
	vent Accessibility of that residue type in a GLY-X-GLY tripeptide, using values gi-
	ven by Rose et al. [80]
	Solvent Accessibility of variant residue
	Normalized Solvent Accessibility of variant residue
	Number of methyl(ene) groups within 6 Å of the variant sidechain [81]
	Number of unsatisfied spatial restraints in the MODELLER objective function after
	in silico mutation and simulated annealing refinement of the variant <sup>a</sup>
	$\Phi$ and $\Psi$ backbone dihedral angles at the mutated position
	Whether the mutation results in buried charge
Physiochemical differences between wild-type and variant amino acid residues	Change in formal charge
	Change in volume (Å <sup>3</sup> ) [82]
	Change in polarity [83]
	Grantham difference [37]
Evolutionary conservation of amino acid residues in protein orthologs	Relative entropy estimated by amino acids in the variant's alignment column [84]
	Positional hidden Markov model conservation score based on the probabilities of
	the wild-type, variant, and most probable amino acid residue in the variant's alignment column <sup>b</sup> [24]

<sup>a</sup>Violated restraints suggest that the mutated sidechain introduced steric clashes or unusual geometries into the protein model. Examples of violated restraints include extreme values of the Lennard-Jones 6–12 potential [85], bond angle potential, bond length potential, sidechain dihedral angle restraints, and nonbonded restraints. Two thresholds are used to identify violated restraints yielding two features.

<sup>b</sup>The probabilities are estimated by a hidden Markov model built with SAM-T2K and the w0.5 script [23].

PHC = log(|p(Wild-type) - p(Variant)|) + log(p(Wild-type)) + log(P(Most Probable)) - log (p(Variant)))

The features were computed for 618 TP53 missense variants, 36 BRCA1 BRCT missense variants biochemically characterized in our companion paper [14], and 54 BRCA1 BRCT UCVs found in BIC.

doi:10.1371/journal.pcbi.0030026.t002

### Results



The β-sheetB buriedBendE exposed	Correctly or incorrectly classified as deleterious
H-bonded turn	Correctly or incorrectly classified as neutral
🕽 loop	O Insufficient confidence to classify
<b>U</b> α-helix	-





### **LS-SNP Large Scale SNP analysis**

### http://salilab.org/LS-SNP/



### Protein function from structure ab-initio localization of binding sites

Rossi. Localization of binding sites in protein structures by optimization of a composite scoring function. Protein Science (2006) vol. 15 (10) pp. 2366-2380

Downloaded from www.proteinscience.org on September 18, 2006

Localization of binding sites in protein structures by optimization of a composite scoring function

ANDREA ROSSI, MARC A. MARTI-RENOM, AND ANDREJ SALI Departments of Biopharmaceutical Sciences and Pharmaceutical Chemistry, California Institute for Quantitative Biomedical Research, University of California, San Francisco, California 94143-2552, USA (RECEIVED March 28, 2006; FINAL REVISION July 10, 2006; ACCEPTED July 11, 2006)

#### Abstract

The rise in the number of functionally uncharacterized protein structures is increasing the demand for structure-based methods for functional annotation. Here, we describe a method for predicting the location of a binding site of a given type on a target protein structure. The method begins by constructing a scoring function, followed by a Monte Carlo optimization, to find a good scoring patch on the protein surface. The scoring function is a weighted linear combination of the z-scores of various properties of protein structure and sequence, including amino acid residue conservation, compactness, protrusion, convexity, rigidity, hydrophobicity, and charge density; the weights are calculated from a set of previously identified instances of the binding-site type on known protein structures. The scoring function can easily incorporate different types of information useful in localization, thus increasing the applicability and accuracy of the approach. To test the method, 1008 known protein structures were split into 20 different groups according to the type of the bound ligand. For nonsugar ligands, such as various nucleotides, binding sites were correctly identified in 55%–73% of the cases. The method is completely automated (http://salilab.org/patcher) and can be applied on a large scale in a structural genomics setting.

Keywords: protein function annotation; small ligand binding-site localization

chosen because of their function, but rather by their HEADER record of their PDB files. In contrast, only 174 location in the protein sequence-structure space (Burley (0.5%) of the 35,199 protein structures solved outside of et al. 1999; Brenner 2000, 2001; Sali 2001; Vitkup et al. 2001; Chance et al. 2002; Goldsmith-Fischman and Honig 2003). Therefore, the number of functionally To class uncharacterized protein structures is growing. Of the 36,606 entries in the Protein Data Bank (PDB) (Kouranov

Many protein targets of structural biologists are no longer of which had an unknown function according to the

et al. 2006) as of February 23, 2006, 1407 structures were on the known structures, automated structure-based funcdeposited by structural genomics consortia, 985 (70%) tional annotation is required (Wallace et al. 1996, 1997; Kleywegt 1999; Thornton et al. 2000; Babbitt 2003; Reprint requests to: Andrea Rossi or Andrej Sali, Departments of iopharmaceutical Sciences and Pharmaceutical Chemistry, California Biopharmaceutucai sciences and rharmaceutucai chemisny, camorina Institute for Quantiative Biomedical Research, University of California, San Francisco Byers Hall, Office 503B, 1700 4th Street, San Francisco, CA

94143-2552, USA; e-mail: andrea@salilab.org or sali@salilab.org; fax: (415) 514-4231. Article published online ahead of print. Article and publication date are at http://www.proteinscience.org/cgi/doi/10.1110/ps.062247506. To classify the functions of thousands of uncharacter-

Laskowski et al. 2003). In particular, we need to be able to identify the locations and types of binding sites on a given structure, because the binding sites define the

The most principled computational approach to pre dicting the molecular function is to dock a large library of potential ligands against the surface of the protein. In

Protein Science (2006), 15:1-15. Published by Cold Spring Harbor Laboratory Press. Copyright © 2006 The Protein Society



25

# For 20% protein structures function is *unknown*

	Structural Genomics*	Traditional methods
Annotated**	654	28,342
Not Annotated	506 (43.6%)	6,815 (19,4%)
Total deposited	1,160	35,157

\* annotated as STRUCTURAL GENOMICS in the header of the PDB file \*\*annotated with either CATH, SCOP, Pfam or GO terms in the MSD database 36,317 protein structures, as of August 8th, 2006

# Representation







# Ligand fingerprints

	Compactness	Conservation	Charge density	<b>B-factor</b>	Protrusion coefficient	Convexity score	Hydrophobicity
ADP	-1.266	-2.009	0.447	-0.414	-1.521	-1.388	-0.118
AMP	-1.62	-1.962	0.341	-0.381	-1.909	-1.944	-0.518
ANP	-1.007	-2.227	0.176	-0.392	-1.706	-1.595	-0.14
АТР	-1.122	-2.156	0.228	-0.274	-1.845	-1.768	0.038
BOG	-2.067	-0.012	0.552	-0.465	-0.356	-0.49	-0.781
CIT	-2.948	-1.58	0.563	-0.527	-0.922	-0.838	-0.113
FAD	0.505	-2.108	0.366	-0.702	-1.735	-1.725	-0.75
FMN	-1.132	-1.98	0.382	-0.387	-1.803	-1.886	-0.695
FUC	-3.43	0.016	-0.295	-0.123	0.002	0.132	0.459
GAL	-3.186	-0.538	-0.234	-0.068	-0.906	-0.987	0.298
GDP	-1.061	-1.471	0.409	-0.81	-1.472	-1.423	0.182
GLC	-2.813	-1.247	-0.207	-0.399	-1.247	-1.337	-0.089
HEC	-0.172	-0.912	0.286	-0.325	-1.153	-1.27	-1.282
HEM	-0.65 I	-1.571	0.683	-0.51	-1.797	-1.937	-1.47
MAN	-3.72	0.131	0.105	-0.52	-0.605	-0.509	0.405
MES	-3.049	-0.24	-0.338	-0.479	-0.714	-0.926	0.296
NAD	-0.005	-1.852	0.156	-0.232	-1.775	-1.804	-0.858
NAG	-3.419	-0.46	-0.126	-0.154	-0.341	-0.523	-0.078
NAP	-0.009	-1.898	0.612	-0.321	-1.587	-1.656	-0.336
NDP	0.217	-1.741	0.535	-0.312	-1.463	-1.562	-0.498

# Ligand fingerprints



### **Prediction accuracy**



### **Protein function from structure**

#### Comparative annotation. AnnoLite and AnnoLyze.

Marti-Renom et al. The AnnoLite and AnnoLyze programs for comparative annotation of protein structures. BMC Bioinformatics (2007) vol. 8 (Suppl 4) pp. S4





### **DBAliv2.0 database**

http://www.dbali.org



# AnnoLyze

MO2				48 49 52 62 63 66 67 113 116	A
	20.			23 29 31 37 44 48 49 83 85 94 96 103 121	RUSPI .
80G	20.	.00	0.111	19 20 21 48 49 51 96 98 136	(STA
ACY	15.	.87	0.163	23 29 31 37 44 45 81 83 85 94 96 98 103 121 135	
Partner	Av. binding site seq. id.	Av. residue conservatio	e n	Residues in predicted binding site (size proportional to the local conservation)	-A
<u>d.113.1.1</u>	23.68	<u>0.948</u>	19205 81828	0 51 52 53 54 55 56 57 58 77 78 79 80 3 84 85 93 95 97 99 134 135 138 142 145	Car

# Benchmark

	Number of chains
Initial set*	78,167
LigBase**	30,126
Non-redundant set***	4,948 (8,846 ligands)

\*all PDB chains larger than 30 aminoacids in length (8th of August, 2006) \*\*annotated with at least one ligand in the LigBase database

\*\*\*not two chains can be structurally aligned within 3A, superimposing more than 75% of their Ca atoms, result in a sequence alignment with more than 30% identity, and have a length difference inferior to 50aa

	Number of chains
Initial set*	78,167
<b>πBase</b> **	30,425
Non-redundant set***	4,613 (11,641 partnerships)

\*all PDB chains larger than 30 aminoacids in length (8th of August, 2006)

\*\*annotated with at least one partner in the  $\pi$ Base database

\*\*\*not two chains can be structurally aligned within 3A, superimposing more than 75% of their Ca atoms, result in a sequence alignment with more than 30% identity, and have a length difference inferior to 50aa

#### AnnoLyze

# Method



Inherited I	nherited ligands: 4				
Ligand	Av. binding site seq. id.	Av. residue conservation	Residues in predicted binding site (size proportional to the local conservation)		
MO2	59.03	0.185	48 49 52 62 63 66 67 113 116		
CRY	20.00	0.111	23 29 31 37 44 48 49 83 85 94 96 103 121		
<u>80G</u>	20.00	0.111	19 20 21 48 49 51 96 98 136		
ACY	15.87	0.163	23 29 31 37 44 45 81 83 85 94 96 98 103 121 135		



Inherited partners:1			
Partner	Av. binding site seq. id.	Av. residue conservation	Residues in predicted binding site (size proportional to the local conservation)
<u>d.113.1.1</u>	23.68	<u>0.948</u>	19 20 50 51 52 53 54 55 56 57 58 77 78 79 80 81 82 83 84 85 93 95 97 99 134 135 138 142 145


AnnoLyze

# **Scoring function**

#### Ligands

#### Partners



Aloy et al. (2003) J.Mol.Biol. 332(5):989-98.

#### AnnoLyze

# Sensitivity .vs. Precision

	Optimal cut-off	Sensitivity (%) Recall or TPR	Precision (%)
Ligands	30%	71.9	13.7
Partners	40%	72.9	55.7

Sensitivity = 
$$\frac{TP}{TP + FN}$$
 Precision =  $\frac{TP}{TP + FP}$ 

#### Example (2azwA) Structural Genomics Unknown Function

#### Molecule: MutT/nudix family protein



## AnnoLyze http://www.dbali.org

000	DBAli v2.0 tools page	
< > < <	🕨 🛃 🐐 http://salilab.org/DBAli/?page=tools&action=f_annotatechai 📀 ^ Q- JMB Dopa	zo 😡
CIPFISGU Labi UCSFISali La DBAliv20 Tools	I DBAIL I MAMMOTH	date 2007
Home Search	DBAII. Tools associated to the database.	
Structural Genomics Help DBAII ALERTI 06/17/08 - Due to changes of disks, some DBAII tools may not be properly functioning. We are working to solve such problems.	<ul> <li>DBAlit! Compare your own structure to the whole PDB (temporarily not available)</li> <li>AnnoLite: Fast annotation of a chain</li> <li>AnnoLyze: Annotate a chain</li> <li>ModClus: Cluster a list of chains</li> <li>ModClus: Cluster from a chain</li> <li>ModDom: Define domains from a chain</li> <li>SALIGN: Get a multiple structure alignment of a list of chains</li> </ul> Annotate a given chain using the DBAli, LigBase, PiBase and ModBase databases. Chain: 4cpal AnnoLyze Clear Mn Seq. Id.: 20 Y Max Seq. Id.: 100 Y Mn RMSD: 0 Y Max PMSD: 4 Y Mn % Eqpos: 75 Y Max % Eqpos: 100 Y Mn P-value: 4 Y Max P-value: 40 Y Homology based data A Homology based data A Y bornain data A Domain data A	
	Please note:  A permisive selection may result in significant server delay and incorrect annotation.  Running ModDom to obtain domain based data may result in significant server delay.  The annotation of a chain takes significant CPU time. Expect delays of about 2 minutes when selecting all available options.  Size May = Reference + Developed = Statistice + Servertices = Record a recking. Values 544	4 14 (6
	<ul> <li>Huming woodcom to cotain comain bas? Usits may result in significant sorror delay.</li> <li>The annotation of a chain takes significant CPU time. Expect delays of about 2 minutes when selecting all available options.</li> <li>Service Enforcement Environd - Servicite - Servicites - Expect delays - Value - Marcul Marc</li></ul>	



### **Docking of small molecules. Vina.**



Marc A. Marti-Renom <a href="http://bioinfo.cipf.es/squ/">http://bioinfo.cipf.es/squ/</a>

t t



Structural Genomics Unit Bioinformatics Department Prince Felipe Resarch Center (CIPF), Valencia, Spain

## **DISCLAIMER!**

#### Credit should go to Dr. Oleg Trott, Dr. Ruth Huey and Dr. Garret M. Morris



http://autodock.scripps.edu
http://vina.scripps.edu

<b>TTT1</b>	•		1	•	
What	<b>1S</b>	D	OC	kin	<u>g</u> :

Thursday,	September	23,	2010	)
-----------	-----------	-----	------	---

Concernence of the concernence		Software News AutoDock Vina: Improving Docking with a New Scou Optimization, and	s and Update the Speed and Accuracy of ring Function, Efficient 1 Multithreading	
Abstract: Analyses May, a new program for molecular docking under greening. Advectory for used in a diverse of the maining mode prediction, gives main diverse of the maining mode prediction, by using under series in a warran series in the main diverse of the maining mode prediction, by using under series in the main diverse of the main din diverse din diverse of the main din diverse of the ma		OLEG TROTT, AR Department of Molecular Biology, The Scrip, Received 3 March 2009; / DOI 10.1002 Published online in Wiley InterScien	THUR I OISON ps Research Institute, La Jolla, California Accepted 21 April 2009 Dfcc:21334 nee (www.interscience.wiley.com).	
P2009 Wikey Periodicals, Er: J Comput Chen 00: 000-000, 2009      Ary works: AutoDock; molecular docking: virtual screening: computer sided drug design; multithreading: scoring fractions     increases         Introduction          Molecular docking is a computational procedure that attempts to freed to monovalent biologi of macromolecules on, more frequently, of macromolecule (receptor) and a small molecule its proteints of an accomolecule (receptor) and a small molecule its proteint is due to configure attempts and molecules to proteint of the macromolecules on the binding affinity.     The prediction of binding of small molecule its proteint is the molecule is that on to change between, and the standing affinity.     The prediction of binding of small molecule is travious to change between, there is the comment to a particular proteint is the molecule is that on the change between, there is the comment on particular proteint is obtained bid free to ensergio fraction states. Additionally, docking generally assumes in the molecule is that on the change between, the change is the comment of molecule is due to not change batteen.     The molecule is molecule to obtain leads for community due to the science is the contract function, which can be the compary of the particular protections and compare the the due to particular protections and compare the particular protections and compare the top that because are used for running due to the science function, which can be the compary of the particular protection is the compary of the particular protection is the compary of the particular protection is the comment of particular protection is the comment of the science is the compary of the particular science is the comment of the science is the compary of the particular science is the compary of the partite		Abstract: AutoDock Vina, a new program for molecular de achieres an approximately two orders of magnitude speed-ap- developed in our lak (AutoDock 4), while also significantly judging by our tests on the training set used in AutoDock 4 de by using multithreading on multicore machines. AutoDock V results in a way transparent to the user.	ocking and virtual screening, is presented. AutoDock Vina usampared with the malecular docking software previously improving the accuracy of the binding mode predictions, evolopment. Further speech up is achieved from parallelism, fina automatically calculates the grid maps and clusters the	
Indexedual docking is a computational procedure that attempts to predict noncovalent hinding of macromolecules on more frequently.       I detail for computational speed. <sup>3</sup> Index of the source		© 2009 Wiley Periodicals, Inc. J Comput Chem 00: 000–00 Key words: AutoDock; molecular docking: virtual screening function	00, 2009 g: computer-aided drug design; multithreading; scoring sar be seen as making an increasing trade-off of the representational	
drug devolopment. Docking can also be used to ty to predict the bound conformation of known binders, when the experimental blob structures are unavailable. <sup>1</sup> ultinately interested in reproducing chemical potentials, which determine the bound conformation preference and the free energy of binding. It is a qualitatively different concept governed not only by the minimizing the computer time they take, because the computational resources spent on docking are considerable. For example, hundreds of thousands of computer vare used for running docking in F-ightAIDS#Home and similar projects. <sup>2</sup> ultinately interested in reproducing the computer vare used for running docking in F-ightAIDS#Home and similar projects. <sup>2</sup> ultinately interested in the standard chemical potentials of the system. When the superficially system-based terms like the 6-12 van der Wash interactions and Coulomb energies are used in the scoring function, they need to be significantly enginically enginically enginically enginically engines and molecular dynamics with explicit solvent, the molecular dynamics with explicit solvent, and       weighted, in part, to account for this difference between energies and free energies. <sup>4,3</sup> 0 comelecular docking       Correspondence ter X J. Olson; e-mail: olson@ scripps.edu       Correspondence ter X J. Olson; e-mail: olson@ scripps.edu		Molecular docking is a computational procedure that attempts to predict noncovalent binding of macronolecules or, more frequently, of a macromolecule (receptor) and a small molecule (lignal) effi- ciently, starting with their unbound structures, structures obtained from MD simulations, or homology modeling, etc. The goal is to predict the bound conformations and the binding affinity. The prediction of binding of small molecules to proteins is of particular practical importance because it is used to screen vir- tual libraries of dura-like molecules to obtain leads for further	detail for computational speed. <sup>3</sup> Among the assumptions made by these approaches is the com- mitment to a particular protonation state of and charge distribution in the molecules that do not change between, for example, their bound and umbound states. Additionally, docking generally assumes much or all of the receptor rigid, the covalent lengths, and angles constant, while considering a chosen set of covalent bonds freely rotatable (refered to as active rotatable bonds here). Importantly, although molecular dynamics directly deals with energies. Irrefered to as a scive rotatable.	
In the spectrum of computational approaches to modeling receptor- ligand binding. a. molecular dynamics with explicit solvent, b. molecular dynamics and molecular mechanics with implicit solvent, and c. molecular docking O 2000 Wiley Periodical Inc		drug development. Docking can also be used to try to predict the bound conformation of known binders, when the experimental holo structures are unavailable. <sup>1</sup> One is intersted in maximizing the accuracy of these predictions while minimizing the computer time they take, because the compu- tational resources speet on docking are considerable. For example, hundreds of thousands of computers are used for running docking in P 2;pH2DB#HOme and similar projects. <sup>2</sup> <b>Theory</b>	ultimately interested in reproducing chemical potentials, which determine the bound conformation preference and the free energy of binding. It is a qualitatively different concept governed not only by the minima in the energy profile but also by the shape of the profile and the temperature. <sup>4,5</sup> Docking programs generally use a scoring function, which can be seen as an attempt to approximate the standard chemical potentials of the system. When the superficially physica-based terms like the 6–12 van der Waals interactions and Coulomb energies are used in the scoring function, they need to be significantly empirically	
1 ALIN WINN PRIVATE IN		In the spectrum of computational approaches to modeling receptor- ligand binding. a. molecular dynamics with explicit solvent, b. molecular dynamics and molecular mechanics with implicit solvent, and c. molecular docking	weigniedi, in part, to account for fuis aniertence between energies and free energies. <sup>4,5</sup>	
		<ul> <li>account for all</li> <li>CORANY AL</li> </ul>	Conjuntation of the conjugate sales (processing) constrained prior for constrained sales (processing) endependent	(2222)
C accessita accessita c accessit	Trott, A.	I. Olson, Journal of (	Computational Chen	ııstry (2009)

## Summary

- INTRO
- DOCKING
- SEARCH METHODS
- EXAMPLE

## • Vina 1.1.1 with ADT

## What is docking?

#### Predicting the best ways two molecules interact.

- Obtain the 3D structures of the two molecules
- Locate the best binding site (Remember AnnoLyze? :-))
- Determine the best binding mode.



## What is docking?

Predicting the **best** ways two molecules interact.

- We need to quantify or rank solutions
- We need a good scoring function for such ranking



## What is docking?

Predicting the best ways two molecules interact.

- X-ray and NMR structures are just ONE of the possible solutions
- There is a need for a search solution



## BIOINFORMATICS

## REPRESENTATION SCORING SAMPLING

## REPRESENTATION



### SCORING AutoDock Vina

 $\Delta G_{binding} = \Delta G_{vdW} + \Delta G_{elec} + \Delta G_{hbond} + \Delta G_{desolv} + \Delta G_{tors}$ 

•  $\Delta G_{vdW}$ *G<sub>vdW</sub>* 12-6 Lennard-Jones potential



•  $\Delta G_{elec}$ 

Coulombic with Solmajer-dielectric

 $\varepsilon(r) = A + \frac{B}{1 + ke^{-\lambda Br}}$ 

•  $\Delta G_{hbond}$ 

12-10 Potential with Goodford Directionality

•  $\Delta G_{desolv}$ 

**Stouten Pairwise Atomic Solvation Parameters** 

•  $\Delta G_{tors}$ 

Number of rotatable bonds







### **PROBLEM!** Very CPU time consuming...



Dihidrofolate reductase with a metotrexate (4dfr.pdb)

### **N=T**<sup>360/i</sup>

N: number of conformations T: number of rotable bonds I: incremental degrees Metotrexato 10 rotable bonds 30° increments (discrete) 10<sup>12</sup> plausible conformations!

### **SOLUTION** Use of grid maps!



- Saves lots of time (compared to classical MM/MD)
- Need to map each atom to a grid point
- ♦ Limits the search space!

### AutoGrid Vina Use of grid maps!

#### Center of grid \*

- ♦ center of ligand
- center of receptor
- a selected atom or coordinate
- Box dimension \*
- Grid resolution (spacing)
  - default 0.375 Angstroms
- Number of grid points (dimension)
  - ouse ONLY even numbers

MAKE SURE ALL LIGAND IS INSIDE GRID AND CAN MOVE!

### With VINA much simplified (\*)



# Simulated Annealing



Φ

### Search algorithms Genetic Algorithm

Use of a Genetic Algorithm as a sampling method

- Each conformation is described as a set of rotational angles.
- 64 possible angles are allowed to each of the bond in the ligand.
- Each plausible dihedral angle is codified in a set of binary bits (2<sup>6</sup>=64)
- Each conformation is codified by a so called chromosome with  $4 \times 6$  bits (0 or 1)

111010.010110.001011.010010

 $\Phi_2$ 

 $\Phi_1 = 1 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 = 58^\circ$ 



Population (ie, set of chromosomes or configurations)



Genetic operators...



### Genetic operators...



# $\overset{H}{\longrightarrow} \overset{OH}{\longrightarrow} \overset{$



## 001010.010101.000101.010001 011010.010110.011010.010111

### Recombination

### 001010.010101.011010.010111 011010.010110. 000101.010001

Genetic operators...

011010.010110.011010.010111 111010.010110.001011.010010 001010.010101.000101.010001 101001.101110.101010.001000 001010.101000.011101.001011

**Migration** 

1111110.010010.0111110.010101 101010.110110.011011.01100 001010.010101.000101.010001 101101.101010.101011.001100 011010.100000.011001.101011

## AutoDock Example Discovery of a novel binding trench in HIV Integrase

Where patients come first 📀 MERCK		Patients & Caregivers   Healthcare Profession Quick Find V	
HOME   ABOUT MERCK   P	RODUCTS   NEWSROOM   INVESTOR RELATIONS   CA	REERS   RESEARCH   LICENSING	G   THE MERCK MANUALS
Newsroom	Product News		()
Product News			
Research & Development News			
Corporate News			
Financial News	FDA Approves ISENTRESS™ (raitegravit	r) Tablets, First-in-	ABOUT ISENTRESS
Corporate Responsibility	Class Oral HIV-1 Integrase Inhibitor	,	Eul Prescribing Information
News			2 Patient Product Information
Fact Sneet	WHITEHOUSE STATION, N.J., Oct. 12, 2007 - Men today that the U.S. Food and Drug Administration (	CK & Co., Inc., announced FDA) granted ISENTRESS™	
Wohcaste	(raltegravir) tablets accelerated approval for use in	combination with other	
VIOXX® (rofecoxib) Information Center	antiretroviral agents for the treatment of HIV-1 infect experienced adult patients who have evidence of strains resistant to multiple antiretroviral agents.	ction in treatment- viral replication and HIV-1	
<u> Contact Newsroom</u> <u>Podcast</u> <u>RSS</u>	This indication is based on analyses of plasma HM weeks in two controlled studies of ISENTRESS [pn studies were conducted in clinically advanced, thre [nucleoside reverse transcriptase inhibitors (NRTIs) transcriptase inhibitors (NNRTIs) and protease inh experienced adults. The use of other active agent associated with a greater likelihood of treatment re efficacy of ISENTRESS have not been established patients or pediatric patients. There are no study r effect of ISENTRESS on clinical progression of HM data will be required before the FDA can consider ISENTRESS.	V-1 RNA levels up through 24 onounced i-sen-tris]. These ee-class antiretroviral s), non-nucleoside reverse libitors (PIs)] treatment- s with ISENTRESS is esponse. The safety and i in treatment-naïve adult esults demonstrating the V-1 infection. Longer term traditional approval for	
	associated with a greater likelihood of treatment re efficacy of ISENTRESS have not been established patients or pediatric patients. There are no study r effect of ISENTRESS on clinical progression of HIV data will be required before the FDA can consider ISENTRESS.	sponse. The safety and in treatment-naive adult esults demonstrating the /-1 infection. Longer term traditional approval for	



One structure known with 5CITEP

- Not clear (low resolution)
- Sinding near to DNA interacting site
- Loop near the binding
- <sup>,</sup> Docking + Molecular Dynamics
  - AMBER snapshots
  - AutoDock flexible torsion thetetrazolering and indole ring.



F Α D

R=







Where patients come	first SMERCK	care Professionals   Worldwide		
HOME   ABOUT MERCK   PRODUCTS   NEWSROOM   INVESTOR RELATIONS   CAREERS   RESEARCH   LICENSING   THE MERCK MANUALS				
Newsroom	Product News			
Product News				
Research & Development News				
Corporate News				
Financial News	FDA Approves ISENTRESS™ (raitegravir) Tablets, First-in-	ABOUT ISENTRESS		
Corporate Responsibility	Class Oral HIV-1 Integrase Inhibitor	Eul Prescribing Information		
News Fact Sheet	WHITEHOUSE STATION N I Oct 12 2007 - Merck & Co. Inc. appounced	2 Patient Product Information		
Executive Speeches	today that the U.S. Food and Drug Administration (FDA) granted ISENTRESS™			
Webcasts	(raltegravir) tablets accelerated approval for use in combination with other			
VIOXX® (rofecoxib) Information Center	experienced adult patients who have evidence of viral replication and HIV-1 strains resistant to multiple antiretroviral agents.			
<ul> <li> <u>Contact Newsroom</u> <u>         Podcast</u> <u>         RSS         </u> </li> </ul>	This indication is based on analyses of plasma HIV-1 RNA levels up through 24 weeks in two controlled studies of ISENTRESS [pronounced i-sen-tris]. These studies were conducted in clinically advanced, three-class antiretroviral [nucleoside reverse transcriptase inhibitors (NRTIs), non-nucleoside reverse transcriptase inhibitors (NRTIs), non-nucleoside reverse transcriptase inhibitors (NRTIs) and protease inhibitors (PIs)] treatment-experienced adults. The use of other active agents with ISENTRESS is associated with a greater likelihood of treatment response. The safety and efficacy of ISENTRESS have not been established in treatment-naïve adult patients or pediatric patients. There are no study results demonstrating the effect of ISENTRESS on clinical progression of HIV-1 infection. Longer term data will be required before the FDA can consider traditional approval for ISENTRESS.			
	effect of ISENTRESS on clinical progression of HIV-1 infection. Longer term data will be required before the FDA can consider traditional approval for ISENTRESS.			

## Vina 1.1.1

Goodsell, D. S. and Olson, A. J. (1990), Automated Docking of Substrates to Proteins by Simulated Annealing Proteins:Structure, Function and Genetics., 8: 195-202. Morris, G. M., et al. (1996), Distributed automated docking of flexible ligands to proteins: Parallel applications of AutoDock 2.4 J. Computer-Aided Molecular Design, 10: 293-304. Morris, G. M., et al. (1998), Automated Docking Using a Lamarckian Genetic Algorithm and and Empirical Binding Free Energy Function J. Computational Chemistry, 19: 1639-1662. Huey, R., et al. (2007), A Semiempirical Free Energy Force Field with Charge-Based Desolvation J. Computational Chemistry, 28: 1145-1152.

# Vina 1.1.1



Goodsell, D. S. and Olson, A. J. (1990), Automated Docking of Substrates to Proteins by Simulated Annealing Proteins:Structure, Function and Genetics., 8: 195-202. Morris, G. M., et al. (1996), Distributed automated docking of flexible ligands to proteins: Parallel applications of AutoDock 2.4 J. Computer-Aided Molecular Design, 10: 293-304. Morris, G. M., et al. (1998), Automated Docking Using a Lamarckian Genetic Algorithm and and Empirical Binding Free Energy Function J. Computational Chemistry, 19: 1639-1662. Huey, R., et al. (2007), A Semiempirical Free Energy Force Field with Charge-Based Desolvation J. Computational Chemistry, 28: 1145-1152.

### Vina 1.1.1 Where to get help...



### Vina 1.1.1 Alternatives



### AutoDock 4.0 Why AutoDock over others



### AutoDock 4.0 Why AutoDock over others



### AutoDock 4.0 Why AutoDock over others


## Vina 1.1.1 Vina and ADT

#### Vina

#### AutoDock Tools

- ♦ 1990 (AutoDock)
- Number crunching (CPU expensive)
- ♦ Command-line!
- ♦ C & C++ compiled

♦ 2000

- Visualizing set-up
- ♦ Graphical user interphase
- Python interpreter



## AutoDock / Vina Practical considerations

- \* What problem does AutoDock solve?
  - *Flexible* ligands (4.0 *flexible* protein).
- \* What range of problems is feasible?
  - \* Depends on the search method:
    - \* LGA > GA >> SA >> LS
    - \* SA : can output trajectories, D < about 8 torsions.
    - \* LGA: D < about 8-32 torsions.
- \* When is AutoDock not suitable?
  - \* No 3D-structures are available;
  - \* Modelled structure of poor quality;
  - \* Too many (32 torsions, 2048 atoms, 22 atom types);
  - \* Target protein too flexible.

# Vina 1.1.1

### Things to know before using AutoDock

### Ligand:

- \* Add all hydrogens, compute Gasteiger charges, and merge non-polar H; also assign AutoDock 4 atom types
- \* Ensure total charge corresponds to tautomeric state
- \* Choose torsion tree root & rotatable bonds

### Macromolecule:

- \* Add all hydrogens, compute Gasteiger charges, and merge non-polar H; also assign AutoDock 4 atom types
- \* Assign Stouten atomic solvation parameters
- \* Optionally, create a flexible residues PDBQT in addition to the rigid PDBQT file
- \* Compute AutoGrid maps

# Vina 1.1.1

### Good that we have AutoDock Tools (ATD)



## Vina 1.1.1 Good we have a nice tutorial



# Acknowledgements

This presentation was based on:

#### "Using AutoDock 4 with ADT. A tutorial"

by Dr. Ruth Huey and Dr. Garret M. Morris

Vina Tutorial by Dr. Oleg Trott





What is Docking?

