

# Supporting information

## Representability of algebraic topology for biomolecules in machine learning based scoring and virtual screening

Zixuan Cang<sup>1</sup>, Lin Mu<sup>2</sup>, and Guo-Wei Wei<sup>1,3,4</sup> \*

<sup>1</sup> Department of Mathematics

Michigan State University, MI 48824, USA

<sup>2</sup> Computer Science and Mathematics Division

Oak Ridge National Laboratory, Oak Ridge, USA

<sup>3</sup> Department of Biochemistry and Molecular Biology

Michigan State University, MI 48824, USA

<sup>4</sup> Department of Electrical and Computer Engineering

Michigan State University, MI 48824, USA

November 8, 2017

Feature Protein cluster	$F^S$ Alpha complex							$F^C$ Alpha complex						
	1	2	3	4	5	6	7	1	2	3	4	5	6	7
NOH	0.64(1.73)	0.79(1.30)	0.63(1.64)	<b>0.63(1.40)</b>	<b>0.82(1.36)</b>	0.63(1.68)	0.67(1.37)	0.53(1.96)	<b>0.76(1.39)</b>	0.45(2.05)	0.50(1.63)	<b>0.76(1.58)</b>	0.39(2.12)	0.57(1.58)
OSH	0.64(1.74)	0.79(1.30)	0.61(1.66)	<b>0.63(1.39)</b>	<b>0.79(1.45)</b>	0.56(1.79)	<b>0.78(1.17)</b>	0.45(2.10)	<b>0.75(1.41)</b>	0.37(2.16)	0.49(1.64)	0.74(1.63)	0.52(1.88)	0.64(1.47)
NOSH	0.65(1.72)	0.78(1.34)	0.64(1.61)	<b>0.64(1.38)</b>	<b>0.81(1.40)</b>	0.61(1.71)	0.67(1.37)	0.54(1.96)	0.74(1.44)	0.51(1.91)	<b>0.58(1.48)</b>	0.74(1.62)	0.40(2.09)	0.58(1.57)
CNS	0.62(1.77)	<b>0.82(1.23)</b>	0.65(1.61)	0.52(1.54)	<b>0.79(1.46)</b>	<b>0.65(1.65)</b>	0.70(1.32)	0.41(2.24)	0.72(1.53)	0.55(1.82)	0.50(1.70)	<b>0.75(1.64)</b>	<b>0.54(1.99)</b>	<b>0.69(1.38)</b>
CNO	0.62(1.77)	<b>0.82(1.23)</b>	<b>0.66(1.58)</b>	0.60(1.44)	0.78(1.48)	<b>0.64(1.65)</b>	0.60(1.49)	<b>0.58(1.90)</b>	<b>0.75(1.43)</b>	<b>0.58(1.79)</b>	<b>0.58(1.53)</b>	0.73(1.67)	0.51(1.99)	0.50(1.73)
OH	0.63(1.77)	0.79(1.30)	0.58(1.73)	0.60(1.44)	<b>0.80(1.43)</b>	0.56(1.79)	<b>0.78(1.16)</b>	0.51(1.98)	0.73(1.46)	0.43(2.03)	0.41(1.75)	<b>0.79(1.45)</b>	0.52(1.88)	<b>0.68(1.38)</b>
CN	0.63(1.77)	0.80(1.27)	0.64(1.63)	0.52(1.55)	0.78(1.48)	0.61(1.72)	<b>0.73(1.27)</b>	0.48(2.10)	0.74(1.47)	<b>0.58(1.78)</b>	0.41(1.83)	<b>0.76(1.59)</b>	0.51(2.07)	<b>0.71(1.35)</b>
CNOSPFClBrH	0.64(1.74)	0.81(1.25)	<b>0.65(1.60)</b>	0.55(1.51)	0.73(1.63)	0.60(1.73)	<b>0.74(1.26)</b>	0.49(2.06)	0.71(1.55)	0.33(2.31)	0.35(1.94)	0.67(1.82)	<b>0.54(1.91)</b>	0.57(1.57)
COSH	<b>0.66(1.70)</b>	0.80(1.29)	<b>0.65(1.59)</b>	0.54(1.52)	0.74(1.61)	0.55(1.80)	0.72(1.28)	<b>0.60(1.83)</b>	0.69(1.58)	0.46(1.99)	<b>0.55(1.59)</b>	0.63(1.90)	0.48(2.02)	0.42(1.80)
CS	0.62(1.78)	<b>0.82(1.23)</b>	0.57(1.81)	0.59(1.45)	0.74(1.60)	<b>0.66(1.63)</b>	0.66(1.39)	0.50(2.09)	0.67(1.65)	0.52(1.90)	0.44(1.78)	0.72(1.71)	0.51(2.06)	0.60(1.50)
NH	0.64(1.74)	0.79(1.30)	0.61(1.66)	<b>0.63(1.39)</b>	0.78(1.48)	0.62(1.69)	0.57(1.53)	<b>0.62(1.79)</b>	0.72(1.50)	0.54(1.87)	<b>0.59(1.47)</b>	0.74(1.62)	0.48(1.94)	0.49(1.69)
CNOSH	0.64(1.75)	0.80(1.28)	<b>0.66(1.58)</b>	0.57(1.48)	0.73(1.62)	0.60(1.72)	0.64(1.42)	0.48(2.06)	0.73(1.50)	0.30(2.33)	0.33(1.94)	0.67(1.81)	<b>0.55(1.91)</b>	0.36(1.95)
NSH	0.63(1.75)	0.79(1.31)	0.63(1.63)	<b>0.65(1.37)</b>	0.76(1.55)	0.59(1.74)	0.57(1.53)	<b>0.63(1.78)</b>	0.70(1.54)	0.54(1.87)	<b>0.61(1.44)</b>	<b>0.75(1.59)</b>	<b>0.54(1.85)</b>	0.49(1.69)
CNOH	0.64(1.74)	0.80(1.28)	<b>0.65(1.59)</b>	0.56(1.49)	0.73(1.62)	0.59(1.74)	0.64(1.43)	0.53(1.98)	0.69(1.58)	0.53(1.87)	0.42(1.80)	0.66(1.83)	<b>0.55(1.89)</b>	0.42(1.83)
CNOS	0.58(1.85)	<b>0.82(1.23)</b>	0.62(1.66)	0.58(1.47)	0.78(1.50)	<b>0.65(1.65)</b>	0.59(1.51)	0.47(2.09)	0.73(1.50)	0.47(1.98)	0.47(1.72)	0.72(1.72)	0.52(1.97)	0.58(1.63)
CO	<b>0.65(1.72)</b>	<b>0.82(1.22)</b>	0.60(1.70)	0.55(1.50)	0.75(1.57)	0.59(1.76)	0.62(1.47)	0.58(1.89)	<b>0.74(1.44)</b>	<b>0.60(1.74)</b>	0.50(1.70)	0.65(1.89)	0.51(2.01)	0.58(1.58)
C	0.61(1.79)	0.80(1.27)	0.59(1.75)	0.52(1.55)	0.75(1.56)	0.64(1.68)	0.68(1.36)	0.54(2.00)	0.70(1.59)	0.51(1.89)	0.41(1.86)	0.70(1.76)	0.44(2.24)	0.55(1.59)
COH	<b>0.67(1.69)</b>	0.78(1.32)	0.61(1.65)	0.55(1.51)	0.71(1.67)	0.52(1.85)	0.73(1.28)	<b>0.58(1.87)</b>	0.70(1.55)	0.47(2.07)	0.39(1.81)	0.67(1.80)	0.44(2.08)	0.46(1.82)
CSH	0.65(1.72)	0.79(1.32)	0.56(1.74)	0.56(1.49)	0.71(1.68)	0.56(1.80)	<b>0.75(1.24)</b>	0.53(1.98)	0.70(1.54)	0.40(2.12)	0.53(1.62)	0.61(1.97)	0.43(2.12)	0.62(1.51)
CNOSPFClBrI	0.59(1.83)	0.81(1.24)	0.62(1.66)	0.56(1.50)	0.73(1.62)	<b>0.65(1.64)</b>	0.61(1.48)	0.44(2.15)	0.72(1.52)	0.50(1.94)	0.36(1.86)	0.64(1.91)	0.42(2.21)	0.53(1.71)
CNSH	0.64(1.74)	0.77(1.35)	0.61(1.66)	0.53(1.53)	0.72(1.65)	0.62(1.70)	0.66(1.40)	0.56(1.91)	0.73(1.47)	0.49(1.96)	0.49(1.68)	0.65(1.88)	0.43(2.15)	0.62(1.53)
H	0.63(1.76)	0.79(1.30)	0.58(1.71)	0.56(1.49)	0.78(1.49)	0.53(1.83)	0.66(1.39)	0.54(1.95)	<b>0.78(1.32)</b>	0.44(2.11)	0.53(1.56)	0.70(1.72)	0.40(2.05)	0.54(1.63)
COS	0.61(1.79)	0.82(1.23)	0.60(1.72)	0.57(1.48)	0.71(1.68)	0.63(1.68)	0.60(1.50)	0.52(2.00)	0.71(1.54)	0.54(1.85)	0.50(1.65)	0.63(1.98)	0.49(2.03)	0.54(1.64)
SH	0.62(1.78)	0.79(1.31)	0.58(1.74)	0.61(1.43)	0.74(1.60)	0.54(1.82)	0.65(1.41)	0.55(1.91)	0.73(1.46)	0.47(1.98)	0.50(1.62)	0.69(1.75)	0.46(1.95)	0.59(1.52)
CNH	<b>0.65(1.72)</b>	0.77(1.36)	0.60(1.69)	0.53(1.53)	0.71(1.67)	0.58(1.76)	0.66(1.39)	0.50(2.04)	0.70(1.56)	0.43(2.10)	0.51(1.62)	0.56(2.08)	0.35(2.33)	0.62(1.49)
CH	0.63(1.76)	0.77(1.36)	0.56(1.73)	0.54(1.52)	0.70(1.70)	0.52(1.86)	0.73(1.27)	0.46(2.11)	0.64(1.68)	0.41(2.07)	0.54(1.59)	0.62(1.91)	0.33(2.29)	<b>0.66(1.42)</b>
NO	0.65(1.72)	0.72(1.47)	0.62(1.65)	0.51(1.55)	0.59(1.95)	0.57(1.80)	0.64(1.44)	0.52(2.07)	0.62(1.75)	0.57(1.80)	0.50(1.61)	0.54(2.12)	0.38(2.28)	<b>0.64(1.49)</b>
NOS	0.73(1.46)	0.50(1.93)	0.49(1.59)	0.56(2.01)	0.57(1.80)	0.65(1.44)	0.53(2.03)	0.60(1.78)	<b>0.57(1.79)</b>	0.48(1.66)	0.48(2.27)	0.34(2.36)	0.64(1.49)	
NS	0.50(2.00)	0.40(1.97)	0.32(2.03)	0.41(1.67)	0.41(2.18)	0.44(1.98)	0.36(1.75)	0.50(1.99)	0.42(1.97)	0.37(1.95)	0.34(1.84)	0.23(2.45)	0.30(2.19)	0.41(1.70)
OS	0.54(1.94)	0.39(1.96)	0.52(1.84)	0.21(1.78)	0.04(2.41)	0.17(2.15)	0.60(1.52)	0.51(2.07)	0.37(2.00)	<b>0.57(1.75)</b>	0.18(1.83)	0.13(2.38)	0.15(2.21)	0.57(1.58)
N	0.45(2.04)	0.39(1.98)	-	0.38(1.69)	0.38(2.22)	0.42(2.01)	0.29(1.78)	0.41(2.09)	0.43(1.95)	-	0.39(1.73)	0.28(2.35)	0.33(2.12)	0.23(1.81)
O	0.55(1.92)	0.40(1.95)	0.37(1.96)	0.18(1.79)	0.00(2.43)	0.09(2.17)	0.59(1.55)	0.57(1.94)	0.35(2.02)	0.16(2.23)	0.12(1.82)	0.01(2.45)	0.12(2.16)	0.55(1.81)

Table A: Pearson correlation coefficients with RMSE (kcal/mol) in parentheses for predictions by different element type combinations for small-molecule characterization using alpha complex validated with 10-fold cross validation within each of the seven protein clusters in S1322. All experiments are repeated 20 times and the median Pearson correlation coefficient (RMSE in kcal/mol) are reported. The element combinations are ordered according to their average Pearson correlation coefficient in each row of the table. The top five best performing combinations for each protein cluster are highlighted in bold in each column.

\*Address correspondences to Guo-Wei Wei. E-mail:wei@math.msu.edu



	C	N	O	S	CN	CO	NO	CNO
C	0.49(2.63)	0.45(2.77)	0.41(2.80)	0.21(3.08)	0.52(2.56)	0.54(2.51)	0.41(2.74)	0.51(2.56)
N	0.31(2.97)	<b>0.28(3.23)</b>	<b>0.28(3.11)</b>	0.22(2.90)	0.29(3.00)	0.36(2.91)	<b>0.35(2.94)</b>	0.34(2.91)
O	0.22(3.06)	0.21(3.16)	0.24(3.11)	0.01(3.07)	0.27(3.00)	0.20(3.09)	0.22(3.12)	0.27(2.98)
S	0.19(2.95)	<b>0.20(3.01)</b>	0.16(3.07)	<b>0.23(2.89)</b>	0.15(3.00)	0.12(3.03)	0.10(3.10)	0.10(3.12)
P	0.09(2.96)	0.12(2.94)	0.07(2.99)	-0.02(2.98)	0.10(2.96)	0.07(2.98)	0.03(2.98)	0.06(2.98)
F	0.11(2.97)	0.09(3.00)	0.07(3.02)	0.09(2.97)	0.08(2.99)	0.15(2.94)	0.17(2.92)	0.18(2.92)
Cl	0.15(2.94)	0.02(3.05)	0.03(3.05)	0.09(2.97)	0.17(2.93)	0.10(2.98)	0.04(3.03)	0.09(2.99)
Br	0.06(2.96)	0.07(2.96)	0.01(2.99)	0.07(2.96)	0.06(2.97)	0.01(2.99)	-0.01(3.03)	0.07(2.96)
I	0.08(2.96)	0.01(2.98)	0.07(2.96)	0.11(2.95)	0.02(2.97)	0.09(2.95)	0.09(2.95)	0.07(2.96)
CN	0.48(2.62)	0.47(2.72)	0.45(2.71)	0.31(2.92)	0.53(2.53)	<b>0.59(2.40)</b>	0.51(2.58)	0.49(2.59)
CO	0.54(2.52)	0.49(2.68)	0.54(2.54)	0.27(2.96)	<b>0.59(2.40)</b>	<b>0.60(2.39)</b>	0.55(2.52)	0.51(2.59)
CS	0.51(2.56)	0.44(2.80)	0.43(2.76)	0.23(3.11)	0.54(2.53)	<b>0.56(2.47)</b>	0.38(2.83)	0.50(2.59)
NO	0.38(2.87)	0.34(3.02)	0.28(3.08)	0.15(3.07)	0.39(2.84)	0.30(2.98)	0.30(2.96)	0.30(2.97)
NS	0.35(2.90)	0.26(3.24)	0.28(3.08)	0.09(3.01)	0.34(2.95)	0.41(2.80)	0.40(2.83)	0.40(2.79)
OS	0.17(3.16)	0.20(3.19)	0.27(3.17)	-0.06(3.15)	0.23(3.09)	0.17(3.14)	0.15(3.21)	0.29(2.94)
CNO	<b>0.56(2.51)</b>	0.51(2.60)	0.48(2.67)	0.33(2.91)	<b>0.61(2.37)</b>	0.56(2.49)	0.55(2.51)	<b>0.61(2.36)</b>
CNS	0.51(2.57)	0.45(2.76)	0.44(2.73)	0.31(2.90)	0.49(2.60)	<b>0.59(2.40)</b>	0.48(2.62)	0.53(2.53)
COS	<b>0.56(2.47)</b>	0.52(2.61)	<b>0.57(2.47)</b>	0.27(2.97)	<b>0.57(2.45)</b>	<b>0.61(2.37)</b>	0.55(2.51)	0.54(2.53)
NOS	0.43(2.77)	0.37(2.95)	0.32(3.04)	0.25(2.90)	0.40(2.79)	0.31(2.95)	0.32(2.94)	0.35(2.88)
CNOS	<b>0.58(2.44)</b>	0.51(2.64)	0.49(2.65)	0.31(2.93)	<b>0.62(2.34)</b>	0.54(2.53)	<b>0.56(2.47)</b>	<b>0.59(2.42)</b>

Table D: Pearson correlation coefficients with RMSE (kcal/mol) in parentheses for predictions of the PDBBind v2016 core set by different combinations of element types for protein (columns) and ligand (rows) using Rips complex with interactive persistent homology and  $F^C$  features based on the bins,  $\{[0, 2.5], [2.5, 3], [3, 3.5], [3.5, 4.5], [4.5, 6], [6, 12]\}$ . The models are trained with the PDBBind v2016 refined set, excluding the PDBBind v2016 core set. Only Betti-0 information is collected. The top 16 combinations are marked in bold.

	C	N	O	CN	CO	NO	CNO	CNOS
C	0.76(1.98)	<b>0.79(1.96)</b>	0.77(2.01)	0.76(2.02)	0.75(2.02)	<b>0.78(1.95)</b>	0.76(2.01)	0.76(2.00)
N	0.75(2.02)	0.74(2.05)	0.75(2.07)	0.75(2.02)	0.75(2.02)	0.74(2.06)	0.75(2.05)	0.76(2.02)
O	0.75(2.02)	0.72(2.12)	0.74(2.07)	0.75(2.02)	0.73(2.09)	0.74(2.06)	0.73(2.08)	0.74(2.06)
S	0.70(2.15)	0.70(2.17)	0.71(2.16)	0.72(2.12)	0.70(2.16)	0.72(2.10)	0.71(2.13)	0.73(2.09)
CN	<b>0.78(1.94)</b>	<b>0.78(1.96)</b>	<b>0.78(1.98)</b>	0.76(1.99)	0.76(2.00)	<b>0.78(1.94)</b>	0.77(1.97)	<b>0.78(1.96)</b>
CO	0.76(2.00)	0.77(1.98)	0.75(2.05)	0.74(2.06)	0.75(2.05)	0.76(2.00)	0.75(2.04)	0.75(2.04)
CS	0.77(1.96)	<b>0.79(1.95)</b>	0.77(2.00)	0.76(2.01)	0.76(1.99)	<b>0.78(1.95)</b>	0.77(1.99)	0.77(1.98)
NO	0.73(2.09)	0.74(2.08)	0.71(2.16)	0.73(2.10)	0.72(2.13)	0.73(2.11)	0.72(2.11)	0.73(2.10)
NS	0.72(2.09)	0.72(2.12)	0.72(2.13)	0.71(2.12)	0.70(2.14)	0.72(2.11)	0.70(2.16)	0.71(2.13)
OS	0.74(2.03)	0.73(2.09)	0.75(2.04)	0.74(2.03)	0.72(2.10)	0.74(2.05)	0.73(2.07)	0.74(2.04)
CNO	0.75(2.03)	0.77(1.99)	0.76(2.05)	0.74(2.06)	0.73(2.08)	0.75(2.03)	0.75(2.04)	0.76(2.04)
CNS	<b>0.78(1.93)</b>	<b>0.78(1.95)</b>	<b>0.78(1.98)</b>	0.76(1.99)	0.76(1.99)	<b>0.78(1.94)</b>	<b>0.77(1.96)</b>	<b>0.78(1.95)</b>
COS	0.76(1.99)	<b>0.78(1.97)</b>	0.76(2.04)	0.75(2.04)	0.75(2.04)	0.76(1.99)	0.76(2.03)	0.76(2.02)
NOS	0.74(2.06)	0.74(2.07)	0.71(2.17)	0.73(2.09)	0.71(2.13)	0.73(2.10)	0.72(2.12)	0.73(2.10)
CNOS	0.75(2.03)	0.77(2.00)	0.75(2.06)	0.74(2.05)	0.73(2.08)	0.75(2.03)	0.75(2.03)	0.76(2.02)
CNOSFCIBrI	0.75(2.03)	0.77(1.99)	0.75(2.08)	0.74(2.05)	0.75(2.06)	0.75(2.04)	0.76(2.03)	0.76(2.02)

Table E: Pearson correlation coefficients with RMSE (kcal/mol) in parentheses for predictions of the PDBBind v2016 core set by different combinations of element types for protein (columns) and ligand (rows) using alpha complex with  $F^S$  features. The models are trained with the PDBBind v2016 refined set, excluding the PDBBind v2016 core set. Betti-0, Betti-1, and Betti-2 barcodes are considered. The top 16 combinations are marked in bold.

Target	Proteins excluded from training set of the target
ACE	ACE, ADA, COMT, PDE5, ALR2,
AChE	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, FXa,
ADA	ACE, ADA, COMT, PDE5,
ALR2	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, GART, ACE, RXRa, PPARg, AmpC, COX1, COX2,
AmpC	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, GART, ACE, RXRa, PPARg, ALR2, COX1, COX2,
AR	AR, TK, ADA, ALR2, PARR, PNP, SAHH, ER_agonist, ER_antagonist, GR, MR, PPARg, PR, RXRa,
CDK2	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2,
COMT	ACE, ADA, COMT, PDE5, RXRa, ALR2, AmpC, PNP,
COX1	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, ALR2, COX2,
COX2	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, HSP90, ALR2, PARR,
DHFR	GART, DHFR, PPARg,
EGFr	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2,
ER_agonist	ER_agonist, PNP, AR, ER_antagonist, GR, MR, PPARg, PR, RXRa,
ER_antagonist	AR, ER_agonist, ER_antagonist, GR, MR, PPARg, PR, RXRa,
FGF1	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2,
FXa	FXa, thrombin, trypsin, DHFR, GART,
GART	GART, DHFR, PPARg,
GPB	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, COMT,
GR	AR, ER_agonist, ER_antagonist, GR, MR, PPARg, PR, RXRa,
HIVPR	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH,
HIVRT	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, PNP,
HMGR	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, RXRa, ACE, GART, ALR2, AmpC, COX1,
HSP90	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2,
InhA	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH,
MR	AR, ER_agonist, ER_antagonist, GR, MR, PPARg, PR, RXRa, PARR,
NA	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, PPARg, thrombin, trypsin, ADA,
P38 MAP	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2,
PARP	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, COX1, PNP,
PDE5	ACE, ADA, COMT, PDE5, P38 MAP,
PDGFrb	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2,
PNP	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, TK, ADA, COMT, COX1, GPB, PARR, SAHH,
PPARg	AR, ER_agonist, ER_antagonist, GR, MR, PPARg, PR, RXRa,
PR	AR, ER_agonist, ER_antagonist, GR, MR, PPARg, PR, RXRa,
RXRa	AR, ER_agonist, ER_antagonist, GR, MR, PPARg, PR, RXRa, COX1,
SAHH	AChE, ALR2, AmpC, COX1, COX2, GPB, HIVPR, HIVRT, HMGR, InhA, NA, PARR, PNP, SAHH, TK, ADA, COMT, COX1, PARR, GPB, PNP,
SRC	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2, PDE5,
thrombin	FXa, thrombin, trypsin, DHFR, ER_antagonist,
TK	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2, ADA, COMT, ALR2, COX1, GPB, PARR, PNP, SAHH,
trypsin	FXa, thrombin, trypsin, PPARg, ADA, DHFR,
VEGF2	CDK2, EGFr, FGF1, HSP90, P38 MAP, PDGFrb, SRC, TK, VEGF2,

Table F: List of proteins that are excluded from the training set of each target in the DUD dataset.

Run #	EF <sub>2%</sub>	EF <sub>20%</sub>	AUC
1	8.63	3.42	0.830
2	8.57	3.47	0.832
3	8.61	3.43	0.830
4	8.52	3.42	0.831
5	8.36	3.47	0.834
6	8.39	3.45	0.833
7	8.59	3.43	0.831
8	8.94	3.47	0.832
9	8.77	3.42	0.832
10	8.62	3.46	0.832
Median	8.60	3.45	0.832
Std	0.16	0.02	1.2E-3

Table G: The DUD test results of TopVS-ML for each repeated run. Median and Std are the median value and the standard deviation across repeated runs.

Target	ADV	LIG	COM	ALL
AR	0.81	0.83	0.93	0.90
COX2	0.86	0.97	0.80	0.97
DHFR	0.82	0.95	0.94	0.96
ER <sub>agonist</sub>	0.84	0.69	0.91	0.81
MR	0.82	0.78	0.91	0.89
PPAR <sub>g</sub>	0.82	0.70	0.72	0.72
RXR <sub>α</sub>	0.95	0.74	0.91	0.79
SAHH	0.80	0.81	0.72	0.84
Average	0.84	0.81	0.86	0.86

Table H: The AUC for autodock vina, TopVS-ML with only compound features, TopVS-ML with only protein-compound complex features, and TopVS-ML with all features. The targets with high quality results by Autodock Vina are reported (AUC > 0.8)

Target	ADV	LIG	COM	ALL
ACE	0.42	0.85	0.78	0.81
ADA	0.49	0.89	0.89	0.89
AmpC	0.34	0.56	0.37	0.53
FGFr1	0.44	0.97	0.71	0.95
GPB	0.48	0.70	0.69	0.71
NA	0.37	0.79	0.82	0.84
PDGFrb	0.32	0.98	0.90	0.96
Average	0.41	0.82	0.74	0.81

Table I: The AUC for autodock vina, TopVS-ML with only compound features, TopVS-ML with only protein-compound complex features, and TopVS-ML with all features. The targets with low quality results by Autodock Vina are reported (AUC < 0.5)