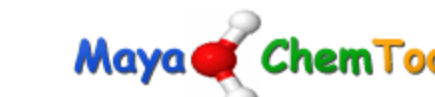


MayaChemTools: An open source package for computational discovery

MayaChemTools.org

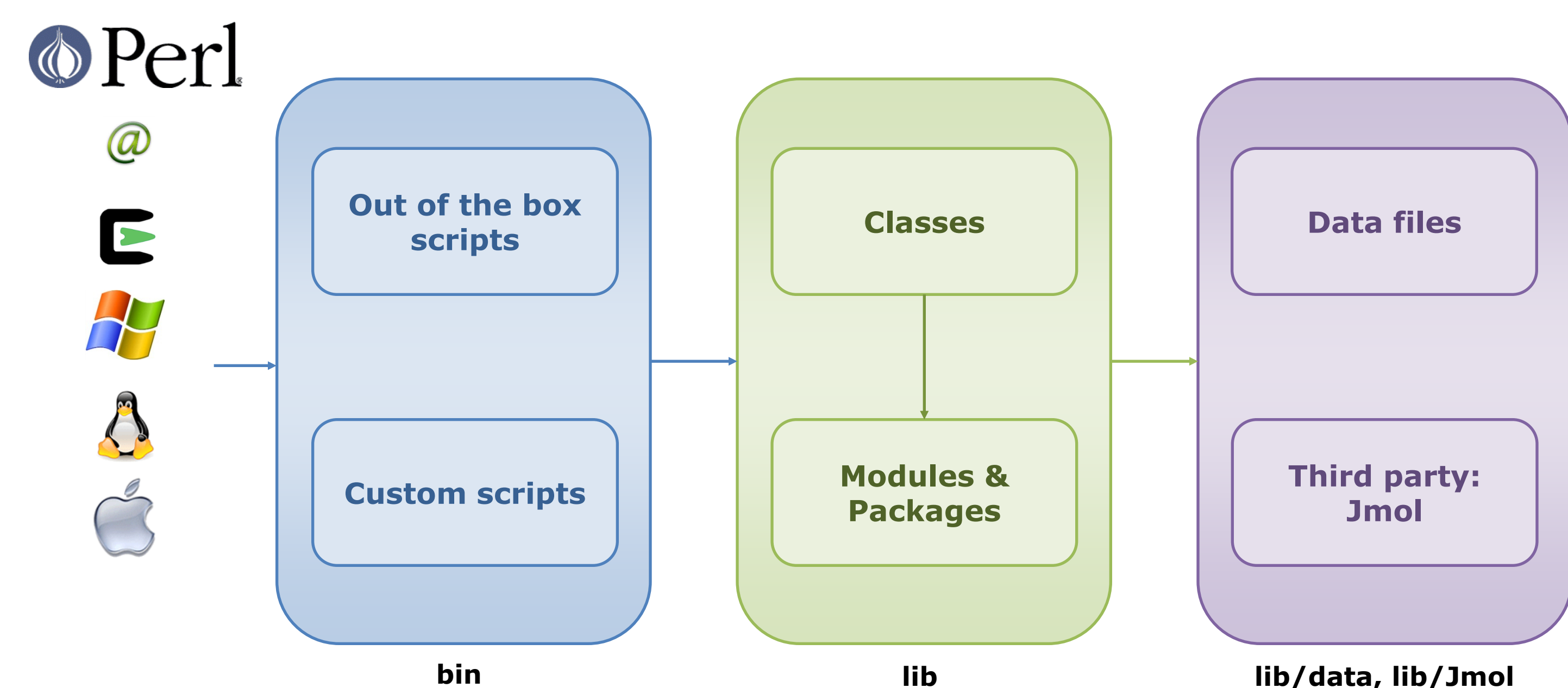
Manish Sud
COMP #306, 243rd ACS National Meeting & Exposition, March 25-29 2012, San Diego, CA



Introduction

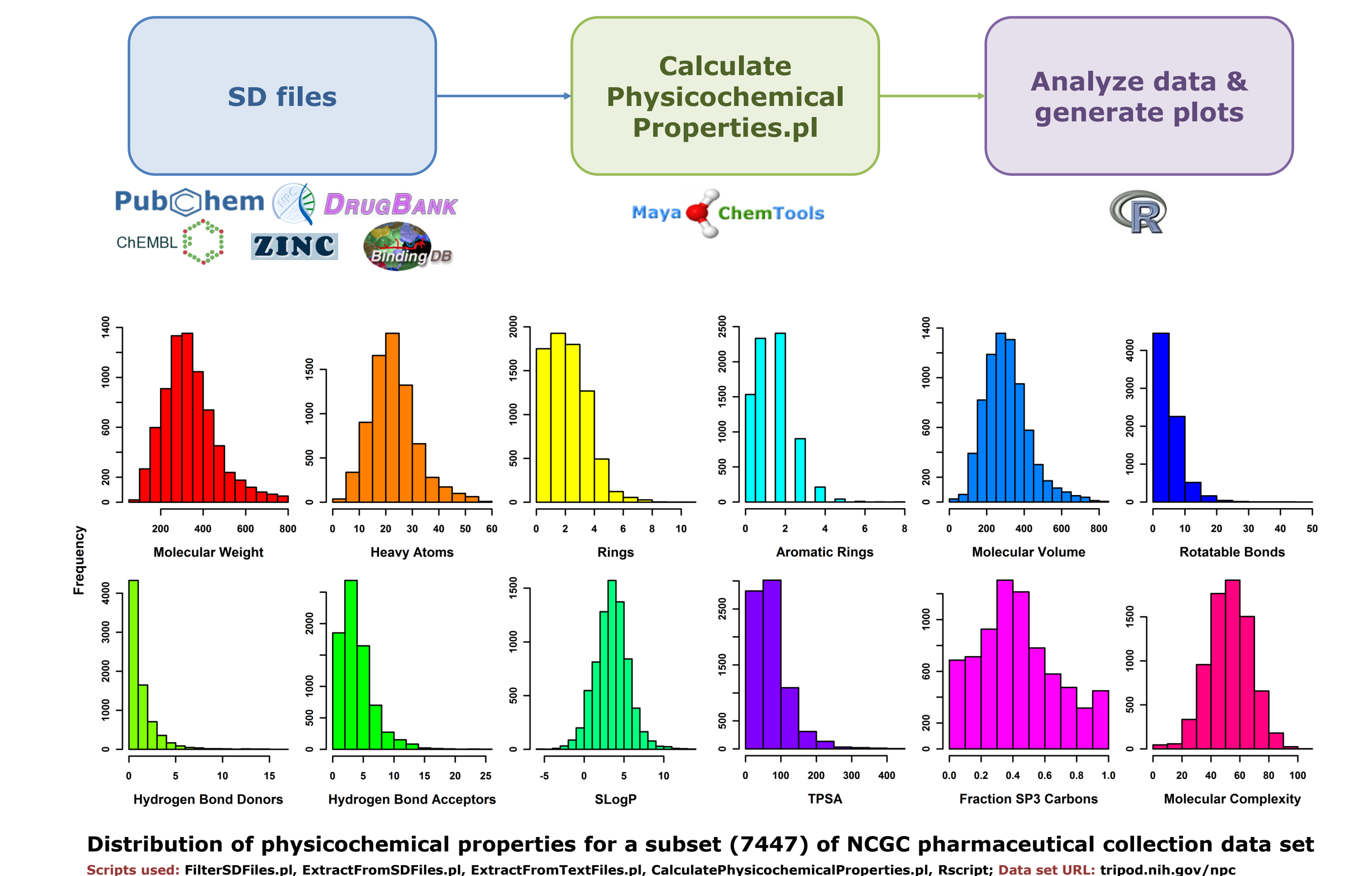
MayaChemTools is a growing collection of Perl scripts, modules and classes to support day-to-day computational drug discovery needs. It provides the following out of the box functionality: Manipulation and analysis of data in SD, CSV/TSV, sequence/alignments, PDB and fingerprints files; Properties of periodic table elements, amino acids and nucleic acids; Calculation of physicochemical properties such as hydrogen bond donors and acceptors, SLogP and topological polar surface area; Generation of fingerprints corresponding to atom neighborhoods, atom types, E-state indices, extended connectivity, MACCS keys, path lengths, topological atom pairs/triplets/torsions and topological pharmacophore atom pairs/triplets; Similarity searching and calculation of similarity matrices. An extensive set of modules and classes are also available for custom development. MayaChemTools is freely available under the terms of the LGPL license at www.MayaChemTools.org.

Software architecture



Physicochemical properties profiling

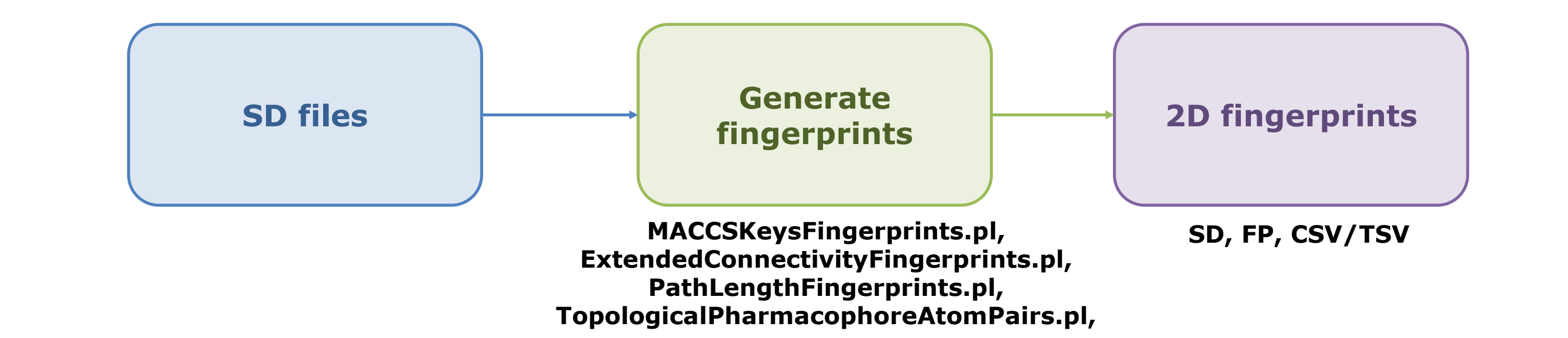
Name	Description
Molecular Weight	Sum of atomic weights
Heavy Atoms	Number of non-hydrogen atoms
Rings, Aromatic Rings	Number of rings and aromatic rings (aromaticity detection using Hückel's rule)
Rotatable bonds	Number of non-ring single bonds involving only non-hydrogen atoms with the option to exclude: terminal bonds; attached to triple bonds; amide, thioamide and sulfonamide bonds
van der Waals Molecular Volume	Sum of atomic volumes corresponding to van der Waals atomic radii with adjustments for number of bonds, aromatic and non-aromatic rings
Hydrogen Bond Donors & Acceptors	Type1 - Donor: Any N and O with implicit/explicit H; Acceptor: Any N without implicit/explicit H and any O Type2 - Donor: Any N and O with implicit/explicit H; Acceptor: Any N and O
LogP & Molar Refractivity (SLogP & SMR)	Sum of atomic contributions from pre-defined atom types corresponding to specific structure fragments
Topological Polar Surface Area (TPSA)	Sum of atomic contributions from pre-defined N and O atom types corresponding to specific structure fragments with option to include P and N atom types
Fraction of SP3 Carbons (FSP3Carbons)	Number of SP3Carbons divided by the total number of carbons
Molecular Complexity	Number of bits-set or unique keys in 2D fingerprints. Supported fingerprints: atom types, extended connectivity, MACCS keys, path lengths, topological atom pairs/triplets/torsions and topological pharmacophore atom pairs/triplets



2D Fingerprints

Type	Values Type	Key Default Parameters/Description
Atom Neighborhoods	Vector	Values: Alphanumerical vector; MinNeighborhoodRadius: 0; MaxNeighborhoodRadius: 2; AtomIdentifierType: AtomicInvariants (AS, X, BO, H, FC)
Atom Types	Bit-vector or vector	Values: Numerical vector; AtomIdentifierType: AtomicInvariants (AS, X, BO, H, FC)
E-state Indices	Vector	Values: Numerical vector; EStateAtomTypesSetSize: Arbitrary
Extended Connectivity	Bit-vector or vector	Values: Alphanumerical vector; NeighborhoodRadius: 2; AtomIdentifierType: AtomicInvariants (AS, X, BO, H, FC, MN)
MACCS Keys	Bit-vector or vector	Values: Bit-vector; Size: 166; Available sizes: 166 and 322; Keys count available
Path Lengths	Bit-vector or vector	Values: Bit-vector; Size: 1024; AtomIdentifierType: AtomicInvariants (AS); MinPathLength: 1; MaxPathLength: 8; Paths count available
Topological Atom Pairs	Vector	Values: Numerical vector; AtomIdentifierType: AtomicInvariants (AS, X, BO, H, FC); MinDistance: 1; MaxDistance: 10
Topological Atom Triplets	Vector	Values: Numerical vector; AtomIdentifierType: AtomicInvariants (AS,X,BO,H,FC); MinDistance: 1; MaxDistance: 10; TriangleInequality: No
Topological Atom Torsions	Vector	Values: Numerical vector; AtomIdentifierType: AtomicInvariants (AS, X, BO, H, FC)
Topological Pharmacophore Atom Pairs	Vector	Values: Numerical vector; AtomTypes: HBD, HBA, PI, NI, H; MinDistance: 1; MaxDistance: 10; AtomTypesWeight: None; FuzzifyAtomPairsCount: No
Topological Pharmacophore Atom Triplets	Vector	Values: Numerical vector; AtomTypes: HBD, HBA, PI, NI, H, Ar; MinDistance: 1; MaxDistance: 10; DistanceInequality: No

Atom Identifier atom types: Atomic invariants, Functional class, DREIDING, EState, HMPPF94, SLogP, SYBYL, TPSA and UFF
Atomic invariants: AS(Atom symbol), X(Num of heavy atom neighbors), BO(Sum of bond orders to heavy atoms), BOO(Largest bond order to heavy atoms), SB(Num of single bonds to heavy atoms), DB(Num of double bonds to heavy atoms), TB(Num of triple bonds to heavy atoms), H(Num of implicit and explicit hydrogens), Ar (Aromatic), RA(Ring atom), FC(Formal charge), MN(Mass number), SM(Multiplicity)
Functional class: HBD(Hydrogen bond donor), HBA(Hydrogen bond acceptor), PI(Positively ionizable), NI(Negatively ionizable), Ar(Aromatic), Hal(Halogen), H(Hydrophobic), RA(RingAtom), CA(ChainAtom)



Fingerprints comparisons

Name	Formula	Name	Formula
Baroni Urbani & Buser	$(SQR((Nc*Nd) + Nc) / (SQR((Nc*Nd) + Nc) + (Na - Nc) + (Nb - Nc)))$	Matching	$(Nc + Nd) / Nc$
Cosine & Ochiai	$Nc / SQR((Na*Nb))$	McConnaughey	$(Nc**2 - (Na - Nc)*(Nb - Nc)) / ((Na*Nb))$
Dice	$2*Nc / (Na + Nb)$	Pearson	$((Nc*Nd) - ((Na - Nc)*(Nb - Nc))) / (SQR((Na*Nb*(Na - Nc + Nd)*(Nb - Nc + Nd)))$
Dennis	$(Nc*Nd - ((Na - Nc)*(Nb - Nc))) / SQR((Nc*Na*Nb))$	Rogers Tanimoto	$(Nc + Nd) / (Na + Nb - 2*Nc + Nt)$
Forbes	$Nc*Nc / Na*Nb$	Russell Rao	Nc / Nt
Fossum	$(Nt*((Nc - 0.5)**2)) / (Na*Nb)$	Simpson	$Nc / MIN(Na, Nb)$
Hamann	$((Nc + Nd) - (Na - Nc) - (Nb - Nc)) / Nt$	Skoal Sneath	1: $Nc / (2*Na + 2*Nb - 3*Nc)$ 2: $(2*Nc + 2*Nd) / (Nc + Nd + Nt)$ 3: $(Nc + Nd) / (Na + Nb - 2*Nc)$
Jaccard & Tanimoto	$Nc / ((Na - Nc) + (Nb - Nc) + Nc) = Nc / (Na + Nb - Nc)$	Tversky	$Nc / (\alpha*(Na - Nb) + Nb)$
Kulczynski	1: $Nc / ((Na + Nb - 2*Nc) * 0.5**2)$ 2: $0.5**2 * (Nc / Na + Nc / Nb)$	Yule	$((Nc*Nd) - ((Na - Nc)*(Nb - Nc))) / ((Nc*Nd) + ((Na - Nc)*(Nb - Nc)))$

N = Num of bits set to "1" in A
Nb = Num of bits set to "1" in B
Nc = Num of bits set to "1" in both A and B
Nd = Num of bits set to "0" in both A and B
Nt = Num of bits set to "1" or "0" in A and B
Na - Nc = Num of bits set to "1" in A not in B
Nb - Nc = Num of bits set to "1" in B not in A

Name	Algebraic Form	Binary Form	Set Theoretic Form
City Block, Hamming & Manhattan Distance	$SUM(ABS(Xai - Xbi))$	$Na + Nb - 2*Nc$	$SUM(Xai) + SUM(Xbi) - 2*(SUM(MIN(Xai, Xbi)))$
Cosine & Ochiai Similarity	$SUM(Xai*Xbi) / SQR(SUM(Xai**2) * SUM(Xbi**2))$	$Nc / SQR((Na*Nb))$	$SUM(MIN(Xai, Xbi)) / SQR(SUM(Xai) * SUM(Xbi))$
Czekanowski, Dice & Sorenson Similarity	$(2*(SUM(Xai*Xbi))) / (SUM(Xai**2) + SUM(Xbi**2))$	$2*Nc / (Na + Nb)$	$2*(SUM(MIN(Xai, Xbi))) / (SUM(Xai) + SUM(Xbi))$
Euclidean Distance	$SQR(SUM((Xai - Xbi)**2))$	$SQR((Na + Nb - 2*Nc))$	$SQR(SUM(Xai) + SUM(Xbi) - 2*(SUM(MIN(Xai, Xbi))))$
Jaccard & Tanimoto Similarity	$SUM(Xai*Xbi) / (SUM(Xai**2) + SUM(Xbi**2) - SUM(Xai*Xbi))$	$Nc / (Na + Nb - Nc)$	$SUM(MIN(Xai, Xbi)) / (SUM(Xai) + SUM(Xbi) - SUM(MIN(Xai, Xbi)))$
Soergel Distance	$SUM(ABS(Xai - Xbi)) / SUM(MAX(Xai, Xbi))$	$(Na + Nb - 2*Nc) / (Na + Nb - Nc)$	$(SUM(Xai) + SUM(Xbi) - 2*(SUM(MIN(Xai, Xbi)))) / (SUM(Xai) + SUM(Xbi) - SUM(MIN(Xai, Xbi)))$

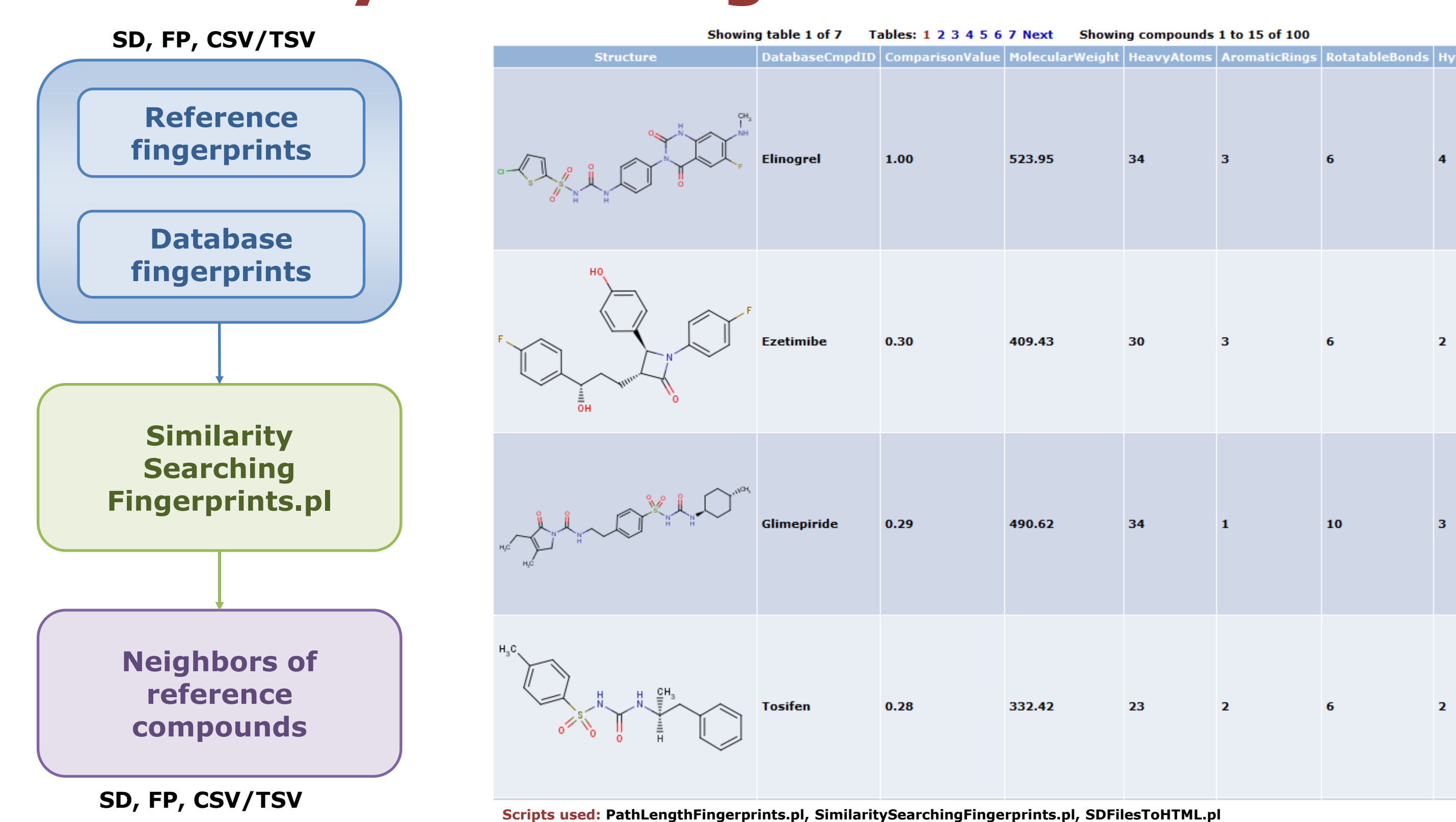
N = Num of values
SUM = Sum over values
Xa = Values of vector A
Xai = Value of ith element in A
Xb = Values of vector B
Xbi = Value of ith element in B
SetIntersectionXaXb = SUM(MIN(Xai, Xbi))
SetDifferenceXaXb = SUM(Xai) + SUM(Xbi) - SUM(MIN(Xai, Xbi))
Nc = Num of bits set to "1" in both A and B = SUM(1 - Xai - Xbi + Xai*Xbi)
Nd = Num of bits set to "0" in both A and B = SUM(1 - Xai - Xbi + Xai*Xbi)

Similarity matrices

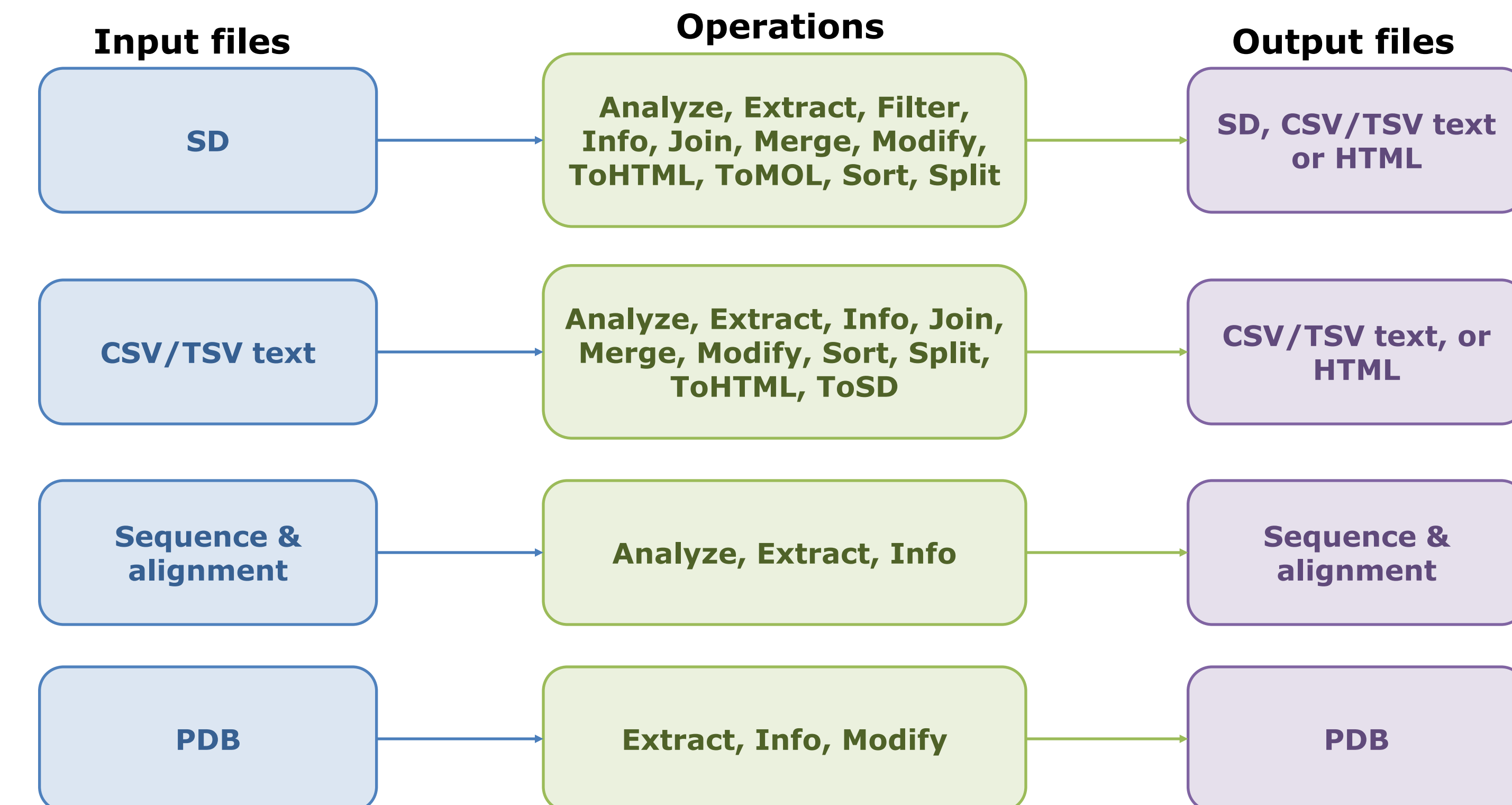


Scripts used: ExtendedConnectivityFingerprints.pl, SimilarityMatricesFingerprints.pl, TextFilesToHTML.pl

Similarity searching



File data info, manipulation & analysis



Data retrieval from databases



Elements, amino acids & nucleic acids

