

## Notes on the Residue Table

**SARvision | Biologics** uses three files to create a data set. These include sequences and data (project file), residues (monomers used to build the biopolymer) and modifiers (transformations that are possible for resides in the residue table, such as methylation, pegylation, I to d). The two last ones, are control files that are modified occasionally as new residue types become part of the projects. Arguably, you can simply use the Residue Table and minimize the use of the Modifiers Table. Therefore, the Residue Table is at the heart of all the capabilities in SARvision | Biologics and deserves some further comments.

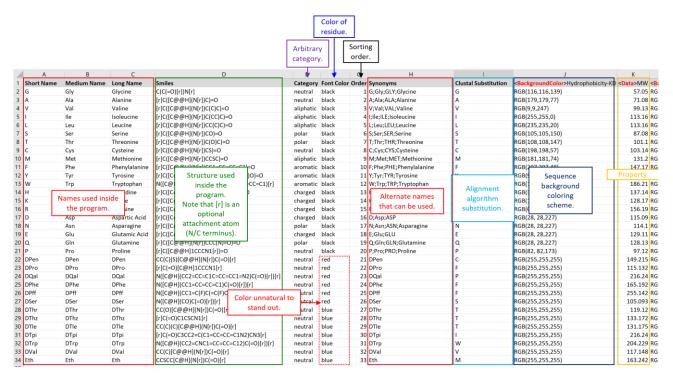
To find these control files: SARvision: main menu: tools: open resource folder... These notes relate to the Monomer File that is

1. The Residue (Monomer) table example is ResidueTable.csv which is a comma delimited ascii file. It contains names, colors, structures (smiles) for each building block. This example is the 20 natural building blocks plus a number of unnatural amino acids.

There are a number of important columns:

- a. Three name columns that contain a short, medium and long name. These are used interchangeably inside **SARvision | Biologics** to optimize the look of tables and views. Note that these can all be the same or of any string length. Ideally, 'Short Name' < 'Medium Name' < 'Long Name'.
- b. The <u>Smiles</u> string is a molecule encoding for the chemical structure of the monomer. The column is necessary, but the individual entries can be empty. Note that Smiles strings can be generated for molecules using most chemically aware programs.
- c. <u>Category</u> is an arbitrary field used to group monomers inside the program for convenience. The labels are completely arbitrary; however, 'polar', 'aliphatic', 'charged', 'neutral', 'aromatic', 'crosslinking'.... are convenient.
- d. The <u>Font Color</u> is the color of the font for this residue in the program. These are usually black with red, blue, green.... To designate unnatural or otherwise interesting residues. Note that RGB(##,##,##) can be used instead of the common color names.
- e. <u>Order</u> column tells the program how to sort these residues. These numbers are completely arbitrary and up to the user's description. If the user sorts an alignment column then this number is used for the sorting order for the column.
- f. The <u>Synonyms</u> column contains other names for this residue that may appear in the data. For example Glycine may appear as 'GLY', 'Gly', or 'G'. I most well curated databases this would be a single entry.
- g. The <u>Clustal Substitution</u> is the residue that should be used for alignments. Note that most alignment matrices (Blossum, PAM....) only work with natural amino acids. There is capabilities to work with alignment matrices with more than 20 entries. Contact us for details. Note that 'X' designates unknown and is handled by the algorithms.
- h. Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be any number of <a href="Note">Note that there can be used to color residues in the alignment table by hydrophobicity; or to color all the aromatic residues one color and the hydrophilic residues another color (e.g. clustal coloring).





2. These can be generated in a chemistry aware program such as SARvision-SM. It can be easily edited in excel and/or exported out of a database. Contact us if you want to directly integrate an oracle data table directly into SARvision-Biologics, such that the residue table is easily kept up to date with new monomer entries.

Outs Table 1								1	1								
	Structure	Short Name	Medium Name	Long Name	Category	Font Color	Order	Smles	Synonyms	Clustal Substitution	Comment	<backgroundcolor> Hydrophobicity-KD</backgroundcolor>	<backgroundcolor> Hydrophobicity-HW</backgroundcolor>	<backgroundcolor> ChemicalStructure</backgroundcolor>	<backgroundcolor> ChouFasman</backgroundcolor>	KData>Kyte Doolittle	<da V</da 
1	∠NH ~	G	Gly	Glycine	neutral	black	1	C(C(=0)[r]]N[r]	G;Gly;GLY;Glycine	G		RGB(116,116,139)	RGB(120,120,136)	RGB(0,255,255)	green	-0.4	
2	∠ <sub>N</sub> ~	А	Ala	Alanine	neutral	black	2	C(C(=O)}-[)(N}-[]C	A:Ala:ALA:Alanine	A		RGB(179,179.77)	RGB(139,139,116)	RGB(0,255,255)	red	1.8	
3	Z H	٧	Val	Valine	aliphatic	black	3	C(C(=O)[r])(N[r])C(C)C	V;Val;VAL;Valine	V		RGB(9.9,247)	RGB(179,179,76)	RGB(0,255,255)	blue	4.2	
4	NH C	ī	le	Isoleucine	aliphatic	black	4	)(C(=0)(-)(N(-))(C(C))	I;lle;lLE;lsoleucine	î		RGB(255,255,0)	RGB(191,191,64)	RGB(0,255,255)	blue	4.5	
5	HN CO	L	Leu	Leucine	aliphatic	black	5	O(C(=O)[F]](N[F]]CC(C)	L;Leu;LEU;Leucine	Ĺ		RGB(235,235,20)	RGB(191,191,64)	RGB(0,255,255)	red	3.8	

3. The Modifier table is more simple. The entries can be pseudo-molecules ('Peg-30') or actual molecules (acetyl or

). Note that there can be any number of <Data> columns. These entries are largely annotative in nature for visual purposes. These are not used for model building or mutation cliffs.

See ModifierTable.csv



Α	В	C	D	E	F	G
<b>Font Color</b>	Name	Sort Order	Synonyms	Comment	<data>MW</data>	Smiles
red	acetyl	2	Ac;acetyl;Acetyl;		43.04	[r]C(C)=O
red	amide	3	amide;Amide		44.03	
red	lauroyl	4	Lauroyl;lauroyl		199.31	
red	myristoyl	5	Myristoyl;myristoyl		227.36	
red	palmitoyl	6	Pal;PAL;Palmitoyl;palmitoyl		255.42	
red	N-methyl	7	N-Me;N(Me);N(CH3);NNe;NCH3;N-methyl;N-Methyl		29.04	
black	N15	8	N15		1	
black	C14	9	C14		2	
black	C13	10	C13		1	
red	Peg-38atom	22	Peg-38 Atom;Peg-38a;Peg-38A;Peg-17atom;Peg-17Atom	MW=(12*2+1	5008	
red	Peg-17atom	23	Peg-17 Atom;Peg-17a;Peg-17A;Peg-17Atom;Peg-17atom	MW=(3*2+1)	1408	
	Names		Possible names or aliases		Data	Optiona Structur

4. Sequences can be read in many sequence formats and are usually in csv format with data. This example had an ID column, two possible sequence formats to read in, and 3 data columns. See GLP-1\_sequences.csv.

А	В		U	E	F	G
ID	Sequence	Sequences(HELMS)	IC50	EC50	REF	
SLP1-7137	HAEGTFTSDVSSYLEGQAAKEFIAWLVKGRG	PEPTIDE1{H.A.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.R.G}	0.31	0.9	EJMC39:4	73(2004)
8	H[d-A]EGTFTSDVSSYLEGQAAKEFIAWLVKGG	PEPTIDE1{H.dA.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.G}	0.33	1.4	EJMC39:4	73(2004)
1	HAEGTFTSDVSSYLEGQAAKEFIAWLVKGG	PEPTIDE1{H.A.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.G}	0.18	0.5	EJMC39:4	73(2004)
2	FAEGTFTSDVSSYLEGQAAKEFIAWLVKGRG	PEPTIDE1{F.A.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.R.G}	0.32	0.9	EJMC39:4	73(2004)
3	FAEGTFTSDVSSYLEGQAAKEFIAWLVKGG	PEPTIDE1{F.A.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.G}	1.6	7	EJMC39:4	73(2004)
4	WAEGTFTSDVSSYLEGQAAKEFIAWLVKGRG	PEPTIDE1{W.A.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.R.G}	3.3	127	EJMC39:4	73(2004)
5	WAEGTFTSDVSSYLEGQAAKEFIAWLVKGG	PEPTIDE1{W.A.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.G}	4.6	152	EJMC39:4	73(2004)
6	YAEGTFTSDVSSYLEGQAAKEFIAWLVKGRG	PEPTIDE1{Y.A.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.R.G}	2.7	5.4	EJMC39:4	73(2004)
7	H[d-A]EGTFTSDVSSYLEGQAAKEFIAWLVKGRG	PEPTIDE1{H.dA.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.R.G}	0.15	0.8	EJMC39:4	73(2004)
9	HSEGTFTSDVSSYLEGQAAKEFIAWLVKGRG	PEPTIDE1{H.S.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.R.G}	2.9	15	EJMC39:4	73(2004)
10	HSEGTFTSDVSSYLEGQAAKEFIAWLVKGG	PEPTIDE1{H.S.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.G}	3.8	17	EJMC39:4	73(2004)
11	HVEGTFTSDVSSYLEGQAAKEFIAWLVKGRG	PEPTIDE1{H.V.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.R.G}	0.47	2.5	EJMC39:4	73(2004)
12	HVEGTFTSDVSSYLEGQAAKEFIAWLVKGG	PEPTIDE1{H.V.E.G.T.F.T.S.D.V.S.S.Y.L.E.G.Q.A.A.K.E.F.I.A.W.L.V.K.G.G}	1.4	3.5	EJMC39:4	73(2004)
12	THE CALLED ACCAL ECON WALLETY WILL AND THE COLOR	DEDTIDE A THE COLUMN TO THE CO	F 7	100	FINAC20.4	72/2004\

5. Complex residues having more than a single letter designation should be added in brackets. Note that the single letters or the names in brackets should match the names in the Residue Table: Synonyms: column. Please contact us if you have a sequence format routinely used other than Fasta with brackets or HELMS. We can help with that.

ID	Name	Sequences	Sequences(HEI
Peptide-1	(Pyr1)-ape	[Pyr1]QRPRLSHKGPMPF	PEPTIDE1{[Pyr:
1064		[pyrE]RPR[d-L]SHKGPMTY	PEPTIDE1{[pyrl
1008		[pyrE]RPRLSKKG	PEPTIDE1{[pyrl
1001		RPRLDHKGPM	PEPTIDE1{R.P.I
1002		RPRLSKKGPM	PEPTIDE1{R.P.I
1003		RPKLSHKGPM	PEPTIDE1{R.P.I
1004		RPR <mark>[d-L]</mark> SHKGPM	PEPTIDE1{R.P.I