

Figure S1

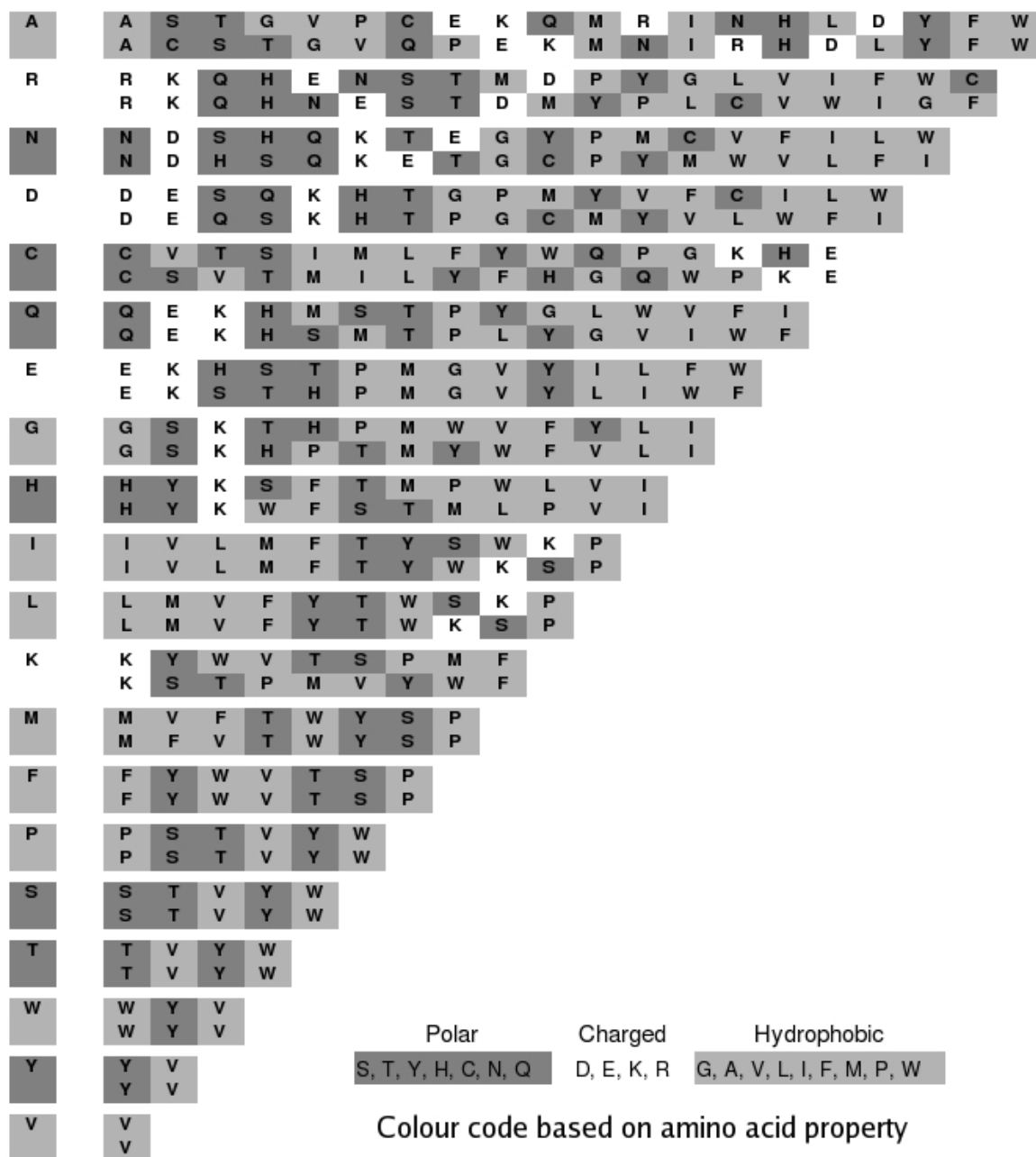


Figure S1. A half-diagonal representation of BLOSUM90 and Smat80 matrices showing differences in their substitution preferences The boxes in the left represent each of the 20 amino acids pairing with all other amino acids (adjacent rows) in the BLOSUM90 (odd rows of alphabets) and Smat80 (even rows of amino acid alphabets) matrices compared here. The amino acid pairs are arranged in a decreasing order of their lodscore values i.e. most preferred followed by the least preferred. Colour code is used to represent different classes of amino acids i.e. dark grey fill for polar amino acids, no fill for the charged amino acids and light grey fill for the hydrophobic amino acids.

Figure S2

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
	0	-2	-2	0	-2	-1	-2	-1	-1	0	-2	-1	0	-1	-1	0	-2	0	-2	0	C
		-1	-1	0	0	0	0	0	0	0	0	1	0	1	1	1	1	2	1	1	S
C	9		0	1	1	1	1	0	1	0	1	0	0	1	1	1	0	1	2	1	T
S	-2	5		0	1	1	0	0	0	1	1	1	0	2	1	1	1	1	1	1	P
T	-2	1	6		0	1	0	0	1	0	1	1	1	0	1	1	0	1	1	1	A
P	-4	-2	-2	8		-1	0	1	1	0	1	2	1	2	2	2	1	1	1	2	G
A	-1	1	0	-1	5		0	0	0	0	0	0	0	0	1	0	0	1	0	-1	N
G	-4	-1	-3	-3	0	6		0	-1	0	0	0	0	1	2	1	1	2	1	0	D
N	-4	0	0	-3	-2	-1	7		0	1	1	1	0	1	1	1	1	1	0	1	E
D	-5	-1	-2	-3	-3	-2	1	7		0	1	1	0	1	0	0	0	0	0	1	Q
E	-6	-1	-1	-2	-1	-3	-1	1	6		-1	1	0	0	1	0	1	0	0	-2	H
Q	-4	-1	-1	-2	-1	-3	0	-1	2	7		-1	0	1	1	1	1	1	0	0	R
H	-5	-2	-2	-3	-2	-3	0	-2	-1	1	8		0	1	0	1	1	1	1	0	K
R	-5	-1	-2	-3	-2	-3	-1	-3	-1	1	0	6		-1	0	0	0	-1	0	0	M
K	-4	-1	-1	-2	-1	-2	0	-1	0	1	-1	2	6		0	0	0	0	1	0	I
M	-2	-2	-1	-3	-2	-4	-3	-4	-3	0	-3	-2	-2	7		0	0	0	0	0	L
I	-2	-3	-1	-4	-2	-5	-4	-5	-4	-4	-4	-4	-4	1	5		0	0	0	1	V
L	-2	-3	-2	-4	-2	-5	-4	-5	-4	-3	-4	-3	-3	2	1	5		0	0	-1	F
V	-2	-2	-1	-3	-1	-5	-4	-5	-3	-3	-4	-3	-3	0	3	0	5		0	0	Y
F	-3	-3	-3	-4	-3	-5	-4	-5	-5	-4	-2	-4	-4	-1	-1	0	-2	7		-1	W
Y	-4	-3	-2	-4	-3	-5	-3	-4	-4	-3	1	-3	-3	-2	-2	-2	-3	3	8		
W	-4	-4	-4	-5	-4	-4	-5	-6	-5	-3	-3	-4	-5	-2	-4	-3	-3	0	2	11	
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	

Figure S2. The difference in the lodscore values of BLOSUM90 and Smat80 matrices. The lower half-diagonal represents the BLOSUM90 lodscore values and the upper half-diagonal shows the difference in the lodscore values of BLOSUM90 & Smat80 (i.e. BLOSUM90-Smat80). The half-bit values have been considered here. Comparison is made with BLOSUM90 as the entropy and scaling is similar to Smat80.


```

      170      180      190      200      210      220
P14223 PTDLS-IHETAWGLARYASICQQNRLVPIVEPEILADGPHSIEVCAVVTQKVLSCVFKAL
      . . . . . : : . . . . . : : : : : . . . . . : : : : :
gi|197 PAGIARVVEQQFEVA--AQVVAAG-LIPIIEPEVDINNVDKVQCEEILRDEIRKHL-NAL
      150      160      170      180      190      200

      230      240      250      260      270      280
P14223 QENG-VLLEGALLKPNMVTAGYECTAKTTTQDVGFLTVRTLRRTVPPALPGVVFLSGGQS
      . . . . . : : . . . . . : : . . . . . : : : : : : : : : :
gi|197 PETSINVMLK--LTLPT-VENLYEEFTKHP-----RVVR-----VVALSGGYS
      210      220      230      240

      290      300      310      320      330      340
P14223 EEEASVNLNSINALGPHPWALTFSYGRALQASVLNTWQGKKE-NVAKAREVLLQRAEANS
      . . . . . : : . . . . . : : . . . . . : : : : : : : : : .
gi|197 REKAN-DILSKNK-----GVIASFSRALTEG-LSAQQTDEEFNKTLAASIDGIYEASVK
      250      260      270      280      290

```

Figure S3: Alignment of fructose-bisphosphate aldolase from AT-rich genomes of *P. falciparum* (P14223) and *Fusobacterium nucleatum* at a gap opening and extension penalties of 12 & 2 respectively. a) Alignment obtained with the best performing standard matrix, Blosum50 b) Alignment extension (highlighted in grey) obtained with PfFSmat60 that extends left and right of the alignment obtained with BLOSUM50, with a good conservation of residues. The alignment length improved 2 fold with an increase in similarity and identity. The query sequence overlap with Pam2 and BLOSUM100 for this example was 81-88 and 79-118 residues respectively.

Figure S4

a)

```

>>P38013|AHP1_YEAST Peroxiredoxin type-2 - Saccharomyces (176 aa)
initn: 213 init1: 183 opt: 222 Z-score: 593.9 bits: 116.8 E(): 2.8e-31
Smith-Waterman score: 222; 30.081% identity (67.480% similar) in 123 aa overlap
(101-218:46-164)

      80      90      100      110      120      130
MAL7P1 MIDVRNMNNSD TDGSPNDFTSIDTHELFNNK KILLISLPGAFPTCSTKMIPGYEEEEYD
      . . . . . : : . . . . . : : : : : . . . . . : : : : :
P38013 FQYIAISQSDADSESKMPQTVEWSKLI SENKKVIITGAPAAFSP TCTVSHIPGYINYLD
      20      30      40      50      60      70

      140      150      160      170      180
MAL7P1 YFIKENNFDDIYCITNNDIYVLKSWFKSMDIK---KIKYISDGNSSFTESMN--MLVDKS
      . . . . . : : . . . . . : : . . . . . : : : : : . . . . . : : : : :
P38013 ELVKEKEVDQVIVVTVDNPFANQAWAKSLGVKDTTHIKFASDPGCAFTK SIGFELAVGDG
      80      90      100      110      120      130

      190      200      210      220      230      240
MAL7P1 NFFMGMRPWR FVAIVENNILVKMFQEKDKQHNIQTD PYDISTVNNVKEFLKNNQL
      . . : : . . : : : : . . : : . .
P38013 VYWSG----RWAMVVENGIVTYAAKETNPGTDVTVSSVESVLAHL
      140      150      160      170

```


Figure S6

a)

```
110      120      130      140      150      160
gi|16805184| VIPNLFKIIYKEKIPSMLDGIFAGVISDKKNNTFFAFRDPIGICPLYIGYAADGSIWFSSE
          : : : : : :
gi|15607948| HGAADSTLEQAALDLLPTVRGAFCLTFMDENTLYACRDPYGVRPLSLGRLDRGWVVASET
          180      190      200      210      220      230
```

b)

```
100      110      120      130      140      150
gi|16805184| NLNKLKSCSDCAVIPNLFKIIYKEKIPSMLDGIFAGVISDKKNNTFFAFRDPIGICPLYIG
          : : . . : : : : : : : :
gi|15607948| DSDILGALLAHGAADSTLEQAALDLLPTVRGAFCLTFMDE--NTLYACRDPYGVRPLSLG
          170      180      190      200      210      220

160      170      180      190      200      210
gi|16805184| YAADGSIWFSSEFKALDNCIRYVI-FPPGHYYKNNKNGEFVRYYNPNWWSLNNNSIPNN
          : . : : : : : : . . . : :
gi|15607948| RLDRGWV-VASETAALDIVGASFVRDIEPGELLAIDADGVRSTRFANPTPKGCVFEYVYL
          230      240      250      260      270      280
```

c)

```
10      20      30
gi|16805184| MCGILAIHFHSSIEKHRLRRKALNLSKILRHR
          : : : : : : : : : :
gi|15607948| SDYVTDRAAGSRQVTVTGQQPEQDLNSPREECGVFGVWAPGEDVAKLTYYGly---ALQHR
          10      20      30      40      50      60

40      50      60
gi|16805184| GPDWNGIIVVEEN-----DDG-TTNVLAHERLAIVD--VLSGH-----QP
          : . : : : : : : : : : : : : : : : : :
gi|15607948| GQEAAAGIAVADGSQVLVFKDLGLVSQVFDEQTLAAMQGHVAIGHCRYSTTGDTTWENAQP
          70      80      90      100      110      120

70      80      90      100      110
gi|16805184| LYDDEEE---VCLTINGEIYNHLELRKLIKEENL-----NKLKSCSDC-----AVI
          : : : : : : : : : : : : : : : :
gi|15607948| VFRNTAAGTGVALGHNGNLVNAALAAARADAGLIATRC PAPATT DSDILGALLAHGAAD
          130      140      150      160      170      180

120      130      140      150      160      170
gi|16805184| PNLFKIIYKEKIPSMLDGIFAGVISDKKNNTFFAFRDPIGICPLYIGYAADGSIWFSSEFK
          : : : : : : : : : : : : : : : :
gi|15607948| STLEQAALDLLPTV-RGAFCLTFMDE--NTLYACRDPYGVRPLSLGRLDRGWV-VASETA
          190      200      210      220      230

180      190      200      210      220
gi|16805184| ALKDNCIRYVI-FPPGHYYKNNKNGEFVRYYNPN----WWSLNNNSIPNNKVD FNEIRI
          : : . . : : : : : . . . : :
gi|15607948| ALDIVGASFVRDIEPGELLAIDADGVRSTRFANPTPKGCVFEYVYLARPDSTIAGRS--V
          240      250      260      270      280      290
```

```

      230      240      250      260      270      280
gi|16805184| HLEKAVI-KRIMMGDVPF-GILLSGGLDSSIIAAILAKHLNILDNKNKGKSIQNGNDKKSNN
      :... : ... : ... : ... : ... : ... : ...
gi|15607948| HAARVEIGRRLARECPVEADLVI GVPESGTPAAV-----GYAQESG---VPYG
      300      310      320      330      340

```

Figure S6. Alignment extension with PffSmat60 for *P. falciparum* putative asparagine synthetase and *M. tuberculosis* PurF protein. The sequences compared are *P. falciparum* putative asparagine synthetase (NCBI gi 16805184) (top line) and *M. tuberculosis* PurF protein (NCBI gi 15607948). a) Alignment with BLOSUM100 with a bit score of 23.2. b) The alignment with BLOSUM50 with a bit score of 42.5. c) The alignment region with PffSmat60 where the structural elements highlighted corresponds closely to the three-dimensional structural superposition of the common domain shared between PurF and asparagine synthase families. The bit score was equal to 369.6 at an E-value of $1.7e-106$. The known crystal structures of *E. coli* asparagine synthetase B (1CT9, chainA) and *B. subtilis* PurF protein (1GPH, chain 3) were used to assign the secondary structure elements. The beta sheets are shaded grey and the alpha elements are shown in white text and a black background.