# Accurate multiplexing and filtering for high-throughput amplicon-sequencing

Esling Philippe[1,2,*], Lejzerowicz Franck[1] and Pawlowski Jan[1]

[1] Department of Genetics and Evolution, University of Geneva, Switzerland
[2] IRCAM, UMR 9912, Université Pierre et Marie Curie, Paris, France

* To whom correspondence should be addressed. Tel: +41223793077; Fax: +41223793340; Email: philippe.esling@unige.ch

Present Address: [Philippe Esling], Department of Genetics and Evolution, University of Geneva, Sciences 3, 30, Quai Ernest Ansermet, CH-1211 Geneva 4, Switzerland
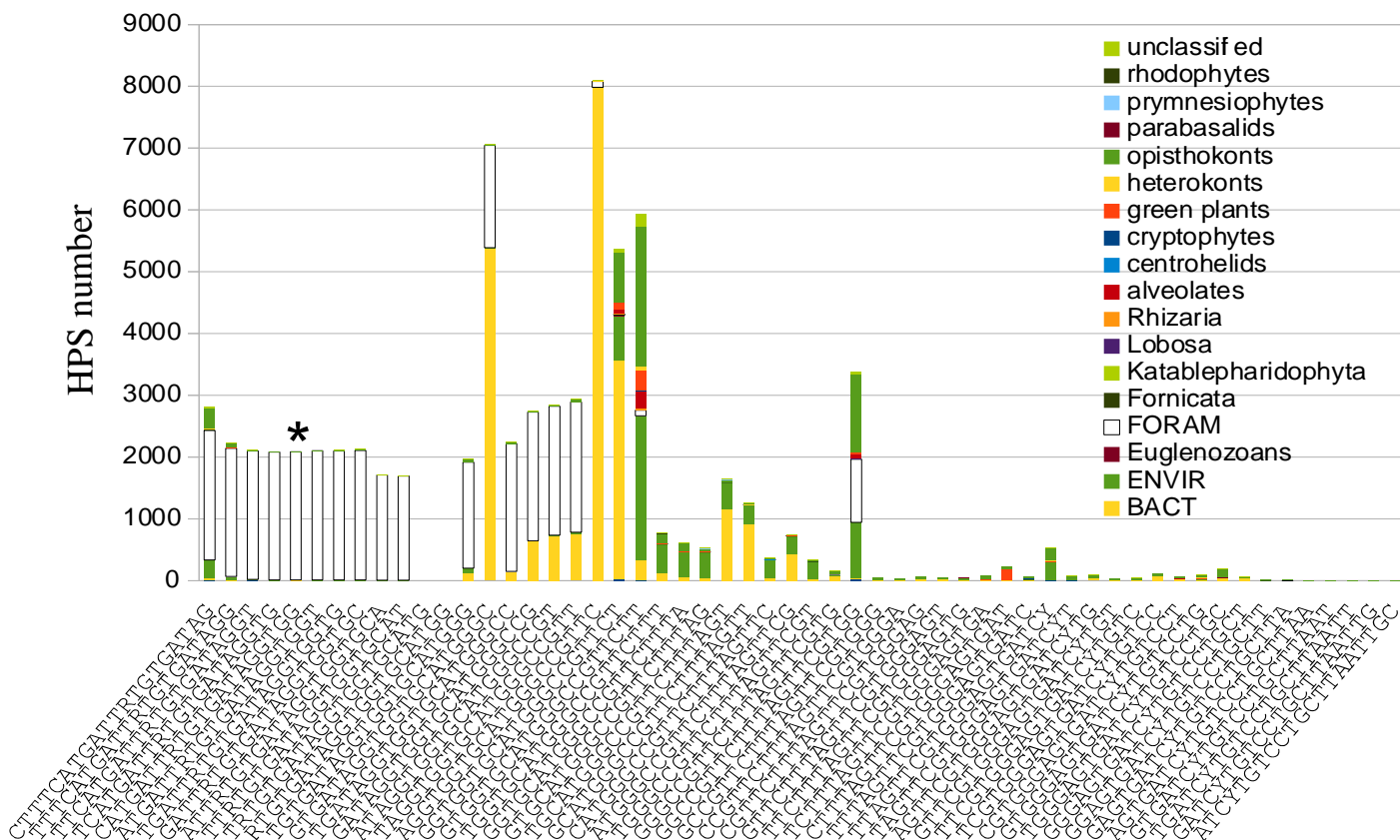
## *Supplementary Methods*

### Selection of the clone sample sets

We calculated all Needleman-Wunsch pairwise distances among the sequences of the single-sequence clone samples obtained as previously explained. Based on the resulting distance matrix, we performed clusters using average-linkage hierarchical clustering at decreasing sequence dissimilarity threshold ranging from 20 to 4 % dissimilarity. For each of the two sequencing runs, we manually assigned cluster reference sequences to the run libraries. We started by distributing the cluster reference sequences obtained at the 20 % divergence thresholds. If too few clusters exist at 20 % to bin enough sequences for our experiments, we continued the sequence distribution at 19 % dissimilarity, and continued until 4 %. This way, we ensure that we only put together samples (i.e. sequences) divergent enough to allow unambiguous assignment during analysis. We display inter- and intra-library samples divergences (Supplementary Figure 3), showing how we optimized the selection of the samples to be multiplexed per library experiment.
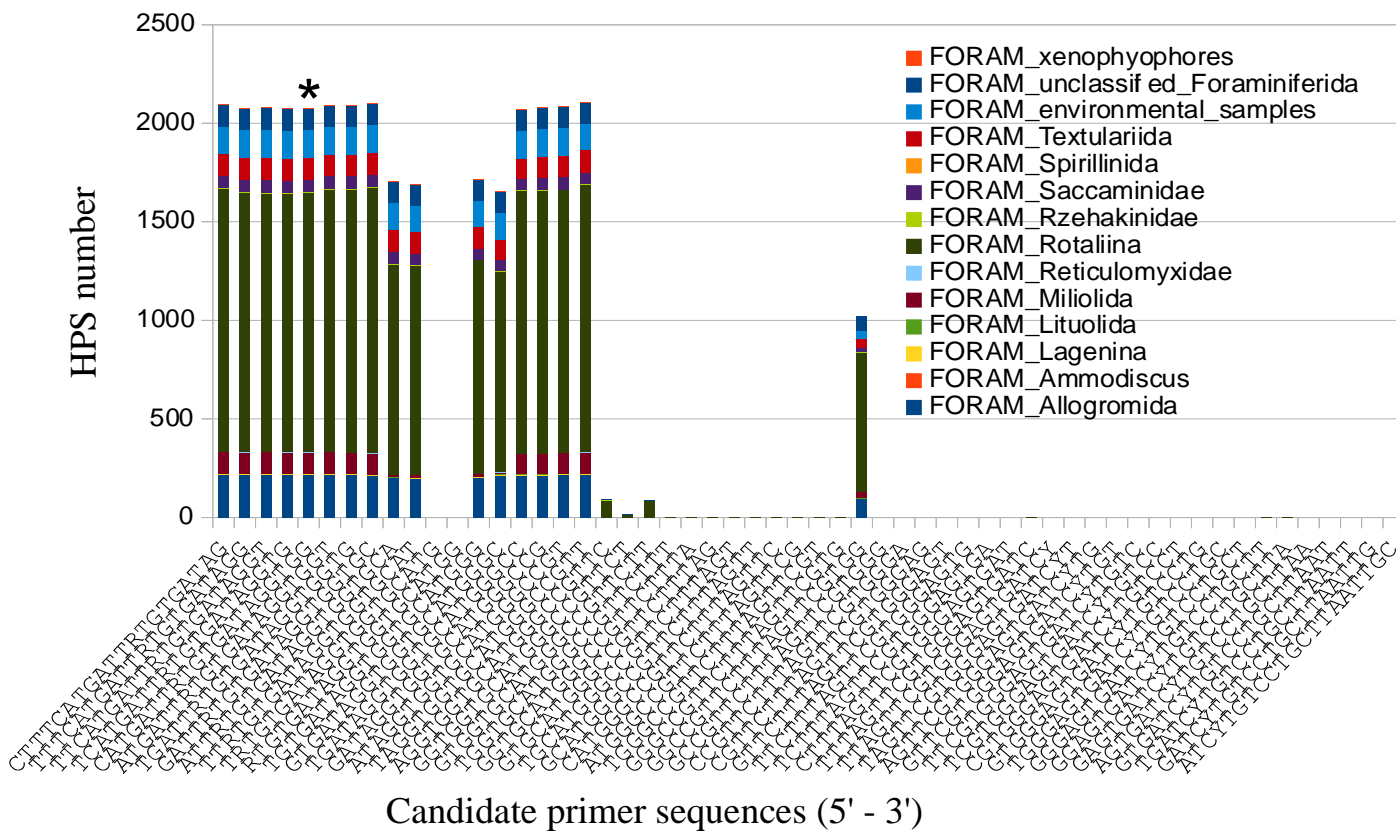
# *Supplementary Figures*

**Supplementary Figure 1**. Taxonomic specificity of the reverse foraminiferal primer s15. For each 20-nucleotide long candidate sequence, we show the results of extensive BLASTn searches against the NBCI's nt database (see online methods). The taxonomy retrieved for each HSP is displayed both at the phylum level (A) and at the foraminifera level (B). The s15 primer covering most of the foraminiferal diversity while avoiding most of the other phyla is indicated by a star.
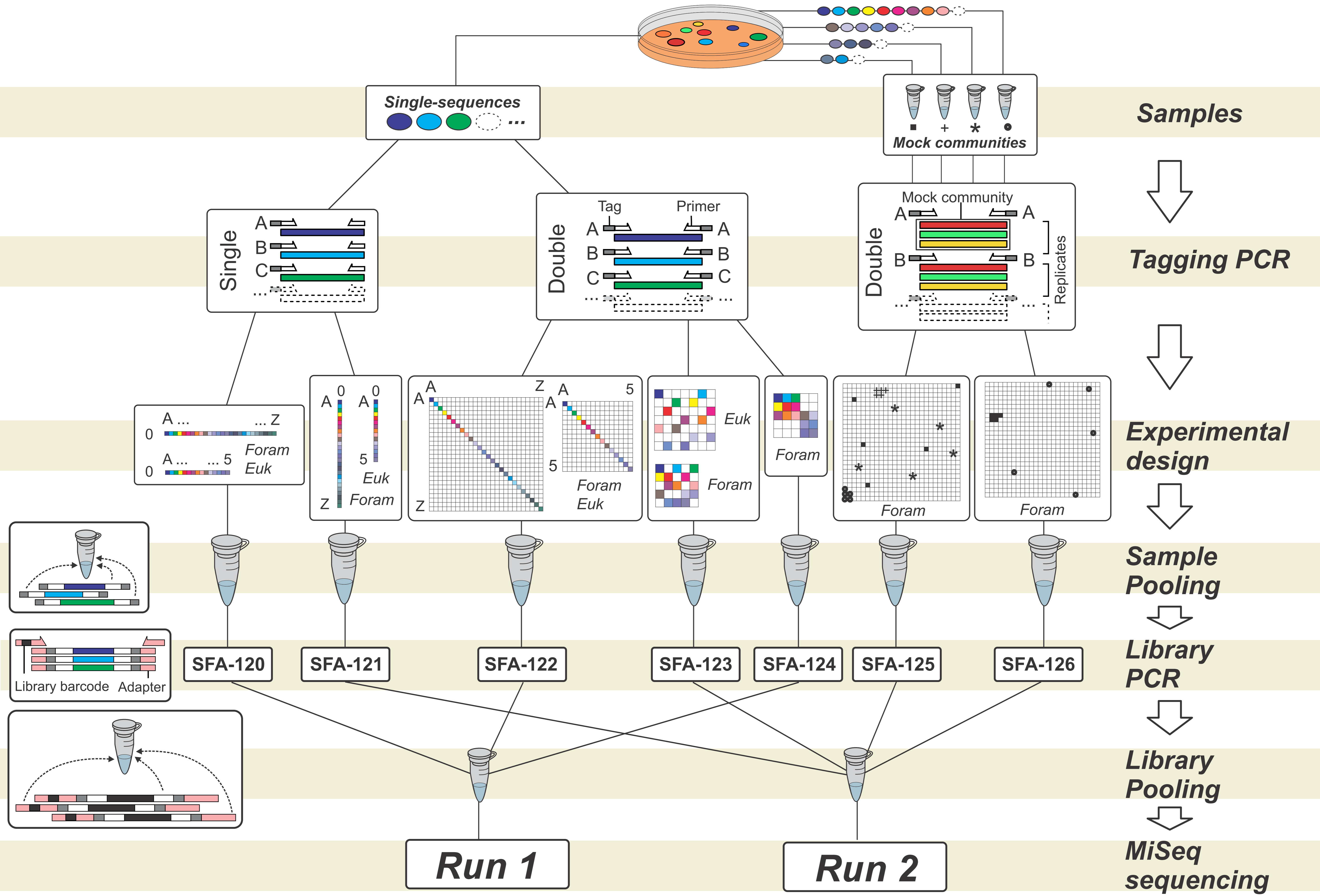
A

B

Candidate primer sequences (5' - 3')

**Supplementary Figure 2**. Experimental designs and molecular workflow. We use a library of Sanger-sequenced clones to provide either *Single-sequence* samples or *Mock communities* samples corresponding to single-sequence clone amplicons pooled in controlled ratios. These samples are labelled by PCR amplification either using one out of the two primers tagged (*Single*) or the two primers tagged (*Double*). We deployed these tagged primers according to each *Experimental design*, represented by the vectors and matrices (rows: forward primers, columns: reverse primers). The samples labelled by the deployed combinations of tagged primers are indicated both for single-sequence samples (colored blocks) and mock communities (black symbols). After the tagging PCR, the labelled samples are pooled in equimolar ratios (*Sample pooling*) and a TruSeq Nano sequencing library (from SFA-120 to SFA-126) is prepared for each pool (*Library PCR*). The resulting libraries are then distributed in two mixes as indicated (*Library pooling*) and sequenced (*MiSeq sequencing*).

*Clones library*

**Samples**

Single-sequences

Mock communities

**Tagging PCR**

Single

Double

Tag        Primer

Double

Mock community

Replicates

**Experimental design**

A ... ... Z

Foram
Euk

Euk
Foram

Foram
Euk

Euk

Foram

Foram

Foram

**Sample Pooling**

**Library PCR**

Library barcode   Adapter

SFA-120   SFA-121   SFA-122   SFA-123   SFA-124   SFA-125   SFA-126

**Library Pooling**

**MiSeq sequencing**

*Run 1*          *Run 2*

**Supplementary Figure 3**. Pairwise distance networks among Sanger-sequenced clone sequences per taxon and per run libraries. The distances among Foraminiferal clones are displayed according to their deployment in both the run 1 (a) and the run 2 (c). The distances among Eukaryotic clones are also displayed according to their deployment in both the run 1 (b) and the run 2 (d). The clones are represented by labeled vertices colored according to the library where they are used. The pairwise distances are represented by edges. Intra- and inter-library distances are materialized by plain and dotted edges, respectively. All distances were measured from exact, pairwise Needleman-Wunsch global alignments and by counting end gaps as well as each internal gaps as differences. Note that the minimum distance between clone samples pooled in a same library was always above 4 %.

**Foraminiferas**

**Eukaryotes**

**Run 1**

a

b

SFA-125
SFA-124
SFA-122
SFA-120
0.1
0.1
0.05

SFA-122
SFA-120
SFA-120_SFA-122

**Run 2**

c

d

SFA-123
SFA-121
SFA-126
0.05
0.1
0.1
0.05

SFA-123
SFA-121
SFA-121_SFA-123

**Supplementary Figure 4**. Clone-to-sample heat maps for per taxon and run. The numbers of reads associated with all the sequences assigned to each foraminiferal clone used in the first (a) run and second run (b) as well as to each eukaryotic clone in the first run (c) are displayed. Only the true samples are presented labeled with the primers combinations. The samples are grouped by library (color code in upper bars and legends) and the libraries are sorted according to their incremental order of preparation. The clones are sorted according to the samples.

**Supplementary Figure 5**. Single tagging mayhem (SFA-121). Mistagging events are displayed in the chord diagrams separately for foraminiferal (a) and eukaryotic (b) data. The central parts represent critical mistags as red links indicating the amount of reads when a sample targeted by a specific tag (one extremity of the string) is found labelled with another tag (other extremity). These central parts would be completely empty in the absence of mistags. For each expected tagged primer, joint barplots indicate the amounts of ISUs (light colors) and reads (dark colors) binned into several categories, including good (expected sample), critical (unexpected sample), non-critical (spurious combination), chimera, dimers and unknown, sequences. The legend to the colors is the same used for Figure 2.

a

b

**Supplementary Figure 6**. Primer-to-primer mistagging events for each taxon in each single-tagging library. For three sequence abundance thresholds, three networks displaying the numbers of mistagged reads above each threshold are displayed for Foraminifera in SFA-120 (a, b, c) and in SFA-121 (d, e, f) as well as for Eukaryota in SFA-120 (g, h, i) and in SFA-121 (j, k, l). The threshold values associated with each network are indicated at the bottom.

Foraminifera (SFA-120)

a  > 132
b  > 330
c  > 463

Foraminifera (SFA-121)

d  > 0
e  > 82
f  > 494

Eukaryota (SFA-120)

g  > 124
h  > 311
i  > 498

Eukaryota (SFA-121)

j  > 76
k  > 229
l  > 382

**Supplementary Figure 7**. Comparison of the 10 PCR product samples sequenced on two separate runs. Each of the 10 PCR products re-sequenced in either a LSD (SFA-123) or a Saturated Design (SFA-124) correspond to one sample and a clone (**a**: F1–B +15–R, **b**: F1–B + 15–T, **c**: F1–C + 15–R, **d**: F1–C + 15–S, **e**: F1–D + 15–S, **f**: F1–D + 15–T, **g**: F1–D + 15–V, **h**: F1–E + 15–U, **i**: F1–E + 15–V, **j**: F1–F + 15–U). The top row displays venn-euler diagrams of the assignments recovered in each sample (purple circle: SFA-123, green circle: SFA-124). The compositions of the re-sequenced sample in terms of relative read abundance are detailed in the vertical bars. For each sample, the correct clone used as template is not included in the bars. The read abundance of these clones are displayed in the pie charts (upper: SFA-123, below: SFA-124) relatively to all the other reads (black). The correct clones are boxed in the legend.

Legend:

- foram14
- foram63
- foram6
- foram35
- foram20
- foram24
- foram36
- foram73
- foram72
- foram1
- foram86
- foram50
- foram83
- foram76
- foram2
- others

**Supplementary Figure 8**. Box plots of the number of reads per ISU assigned to a clone with 1, 2, 3 or more than 4 differences. The results are shown separately for each library and mock community. At the position of each clone name corresponds the group of ISU with the lowest number of difference(s) to this clone and the number of reads in the ISU perfectly matching this clone (blue dot). All the clones found in each mock community are displayed, including the expected clones (black) and the clones resulting from a critical mistagging event (red). The numbers of reads are displayed on a $\log_{10}$ scale.

**Supplementary Figure 9**. Mistagging cohorts for each individual sequence unit (ISU) assigned with less than 2 differences to each expected clone of each mock community sample. For each ISU, the distributions are shown in two separate heat maps in the framework of their correct samples. One heat map shows the ISU perfectly corresponding to the original clone sequence (large, top heat map) and all ISUs matching this sequence with 1 difference (lower left). The numbers of reads are indicated according to a green-to-red scale. Each clone name is indicated above this scale. The tagged primer pairs used for the replicate PCRs of the library indicated in the upper-right box are colored per combination. The mock community into which the clone is expected is indicated in the box in red. The relative abundance of each clone belonging to the mock communities of SFA-125 are indicated by a red letter in parentheses ("l": low; "m": medium; "h", high and "H": very high relative abundances). The numbers of reads per correct and non-critical mistag ISU are indicated on the lower right panel. Further details on mock community compositions are provided in Supplementary Table 2.

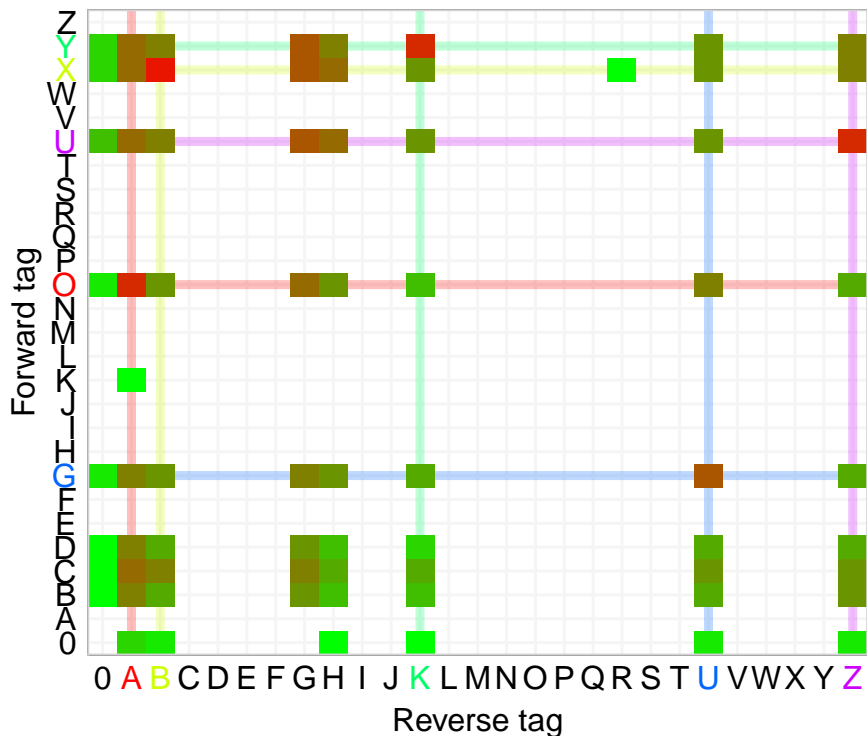**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram25

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 19]
]19 – 36]
]36 – 68]
]68 – 128]
]128 – 242]
]242 – 457]
]457 – 864]
]864 – 1633]
]1633 – 3088]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

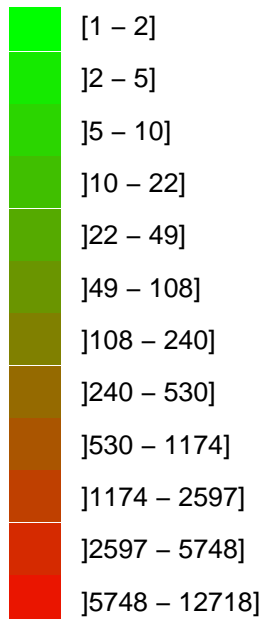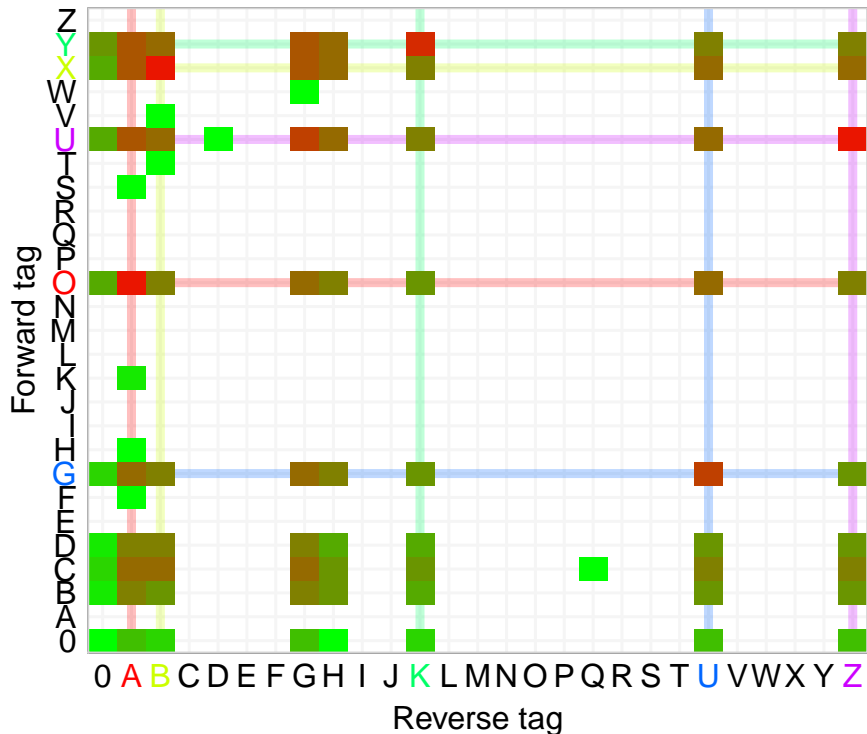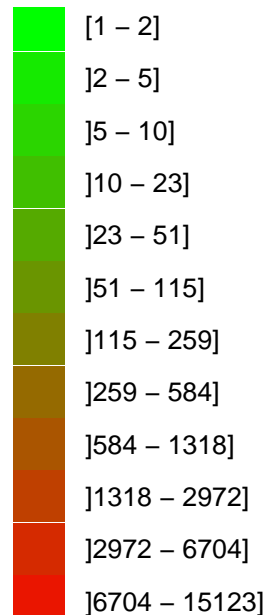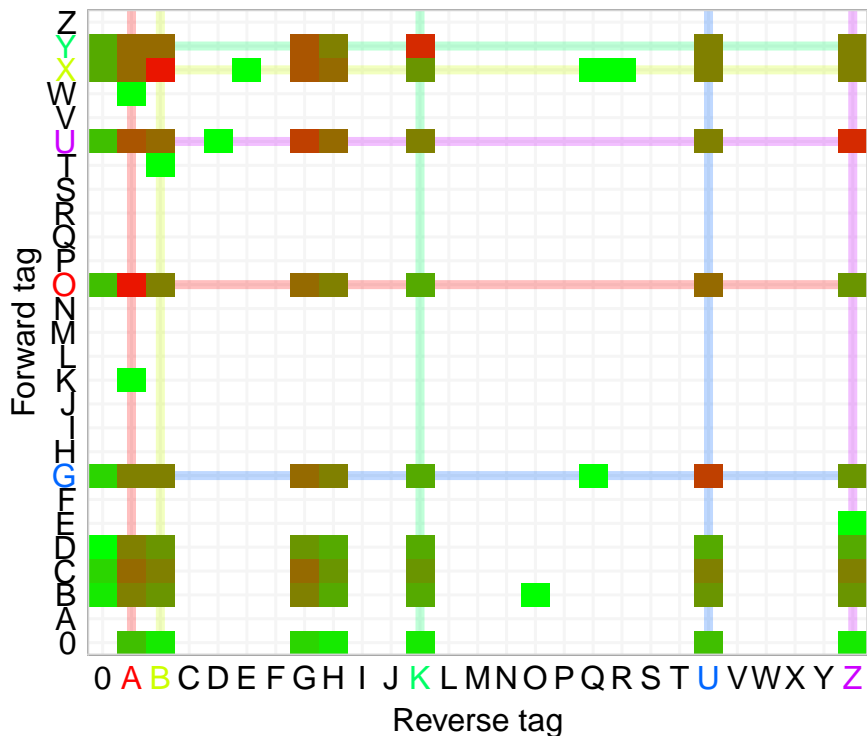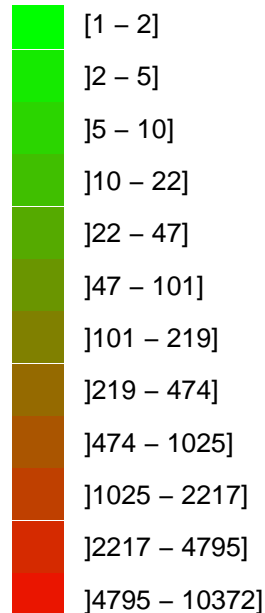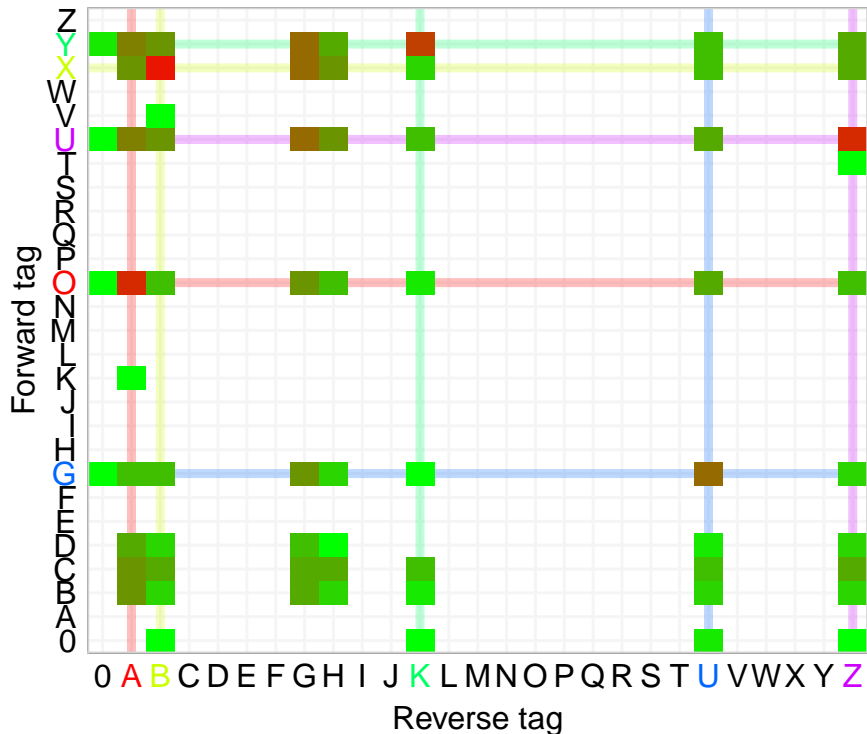# Perfect match (0 difference ISU)

**SFA−126**
mock: random

foram59

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 50] |
| | ]50 – 111] |
| | ]111 – 248] |
| | ]248 – 554] |
| | ]554 – 1236] |
| | ]1236 – 2758] |
| | ]2758 – 6155] |
| | ]6155 – 13736] |

Forward tag

Reverse tag

**1 difference ISUs**

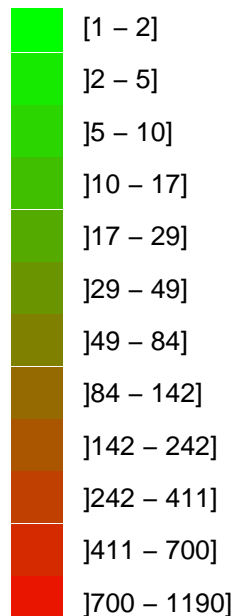○ non−critical mistag  ● correctly labelled

Number of reads per sample
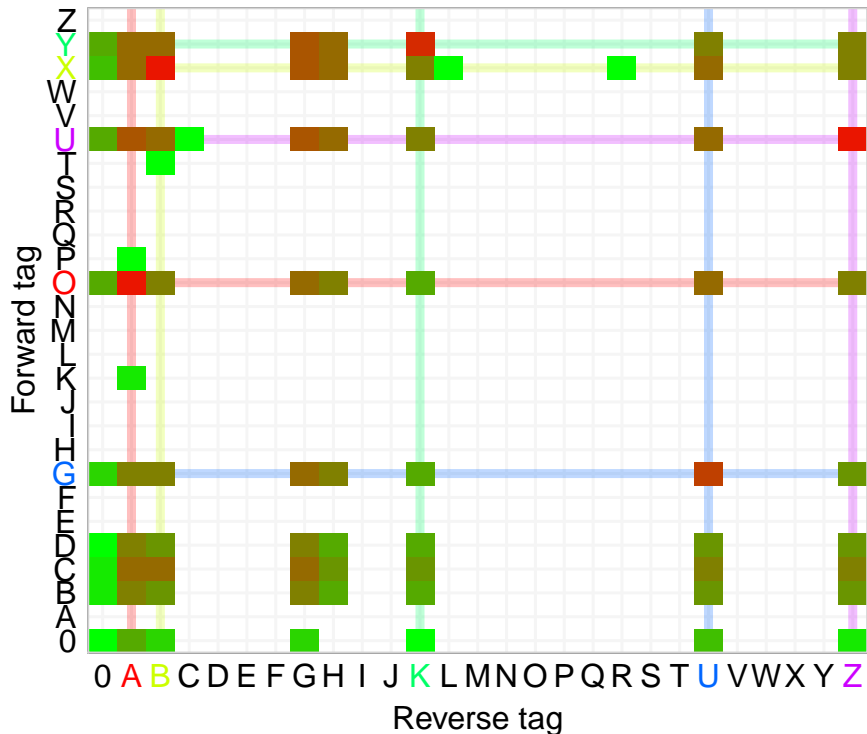
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram58

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 49] |
| | ]49 – 110] |
| | ]110 – 245] |
| | ]245 – 544] |
| | ]544 – 1210] |
| | ]1210 – 2690] |
| | ]2690 – 5982] |
| | ]5982 – 13302] |

Forward tag

Reverse tag

**1 difference ISUs**
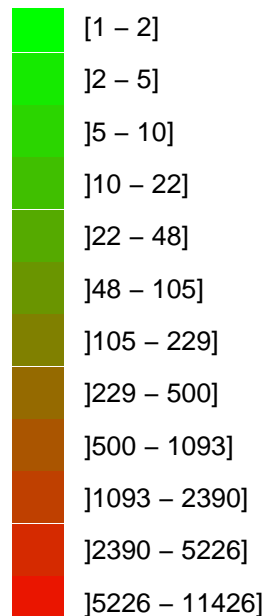
○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram55

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 20]
]20 – 41]
]41 – 84]
]84 – 170]
]170 – 345]
]345 – 700]
]700 – 1422]
]1422 – 2887]
]2887 – 5861]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag ● correctly labelled

Number of reads per sample

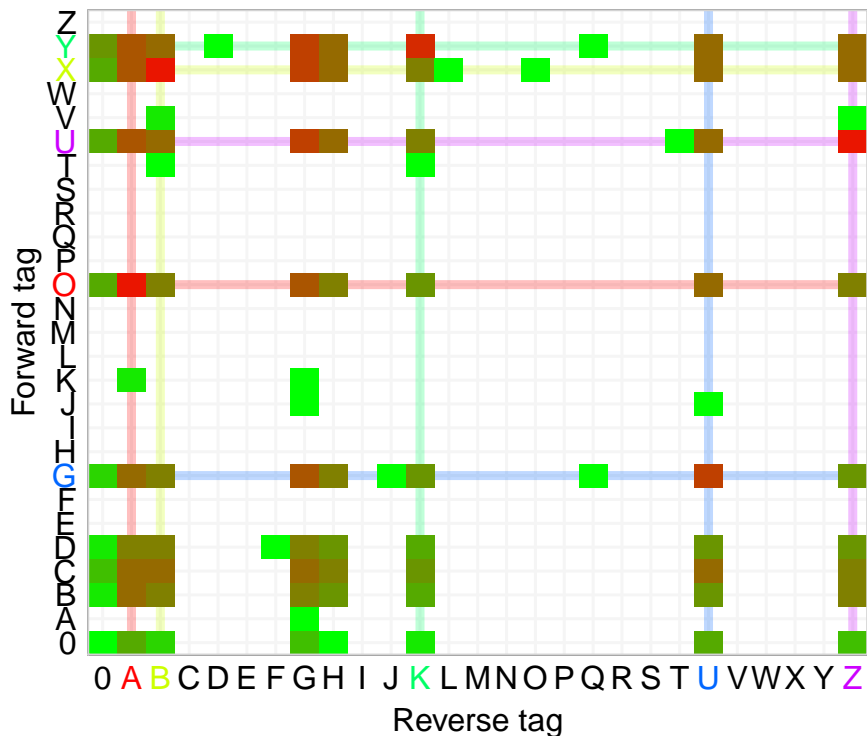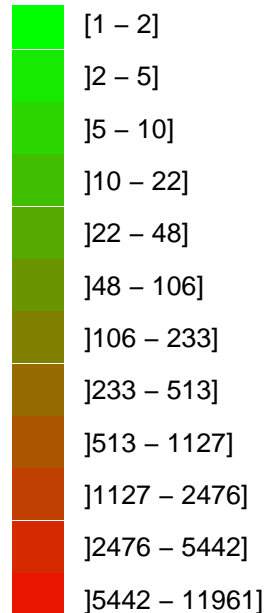**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram54

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 21]
]21 − 45]
]45 − 95]
]95 − 202]
]202 − 428]
]428 − 907]
]907 − 1922]
]1922 − 4073]
]4073 − 8634]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
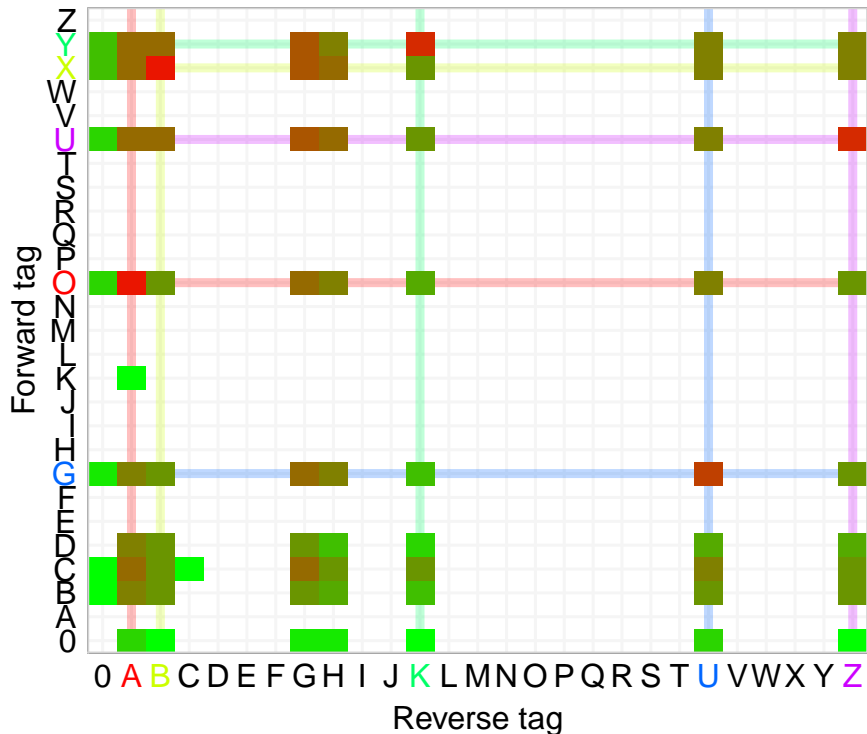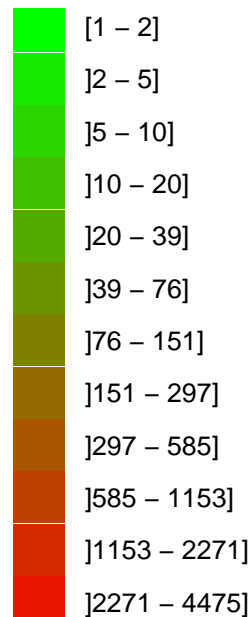
# Perfect match (0 difference ISU)

**SFA−126**
mock: random

foram57

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 48] |
| | ]48 – 107] |
| | ]107 – 235] |
| | ]235 – 517] |
| | ]517 – 1138] |
| | ]1138 – 2504] |
| | ]2504 – 5512] |
| | ]5512 – 12134] |

Forward tag

Reverse tag

## 1 difference ISUs

○ non−critical mistag   ● correctly labelled

Number of reads per sample
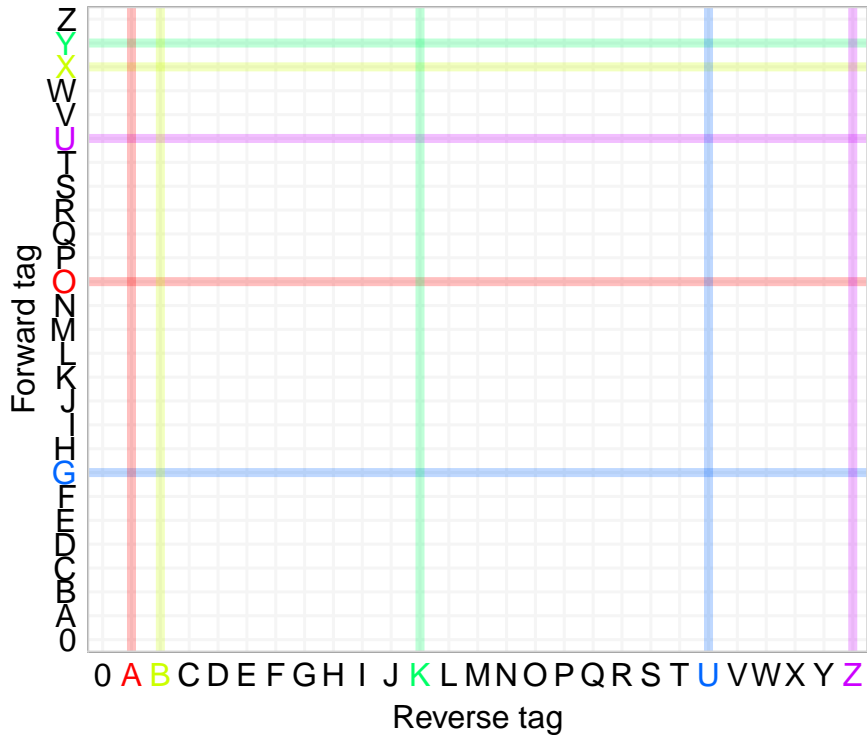
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram56

- [1 − 2]
- ]2 − 5]
- ]5 − 10]
- ]10 − 20]
- ]20 − 42]
- ]42 − 86]
- ]86 − 176]
- ]176 − 361]
- ]361 − 740]
- ]740 − 1516]
- ]1516 − 3107]
- ]3107 − 6366]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag   ● correctly labelled

Number of reads per sample

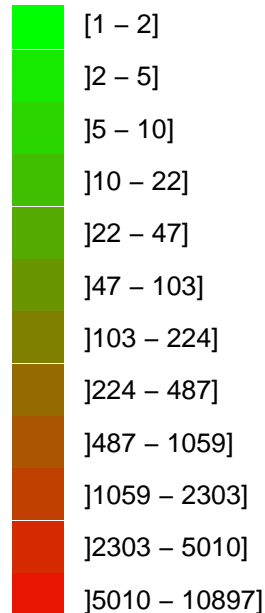**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram51

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 20]
]20 – 40]
]40 – 79]
]79 – 157]
]157 – 312]
]312 – 620]
]620 – 1233]
]1233 – 2452]
]2452 – 4878]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag ● correctly labelled

Number of reads per sample
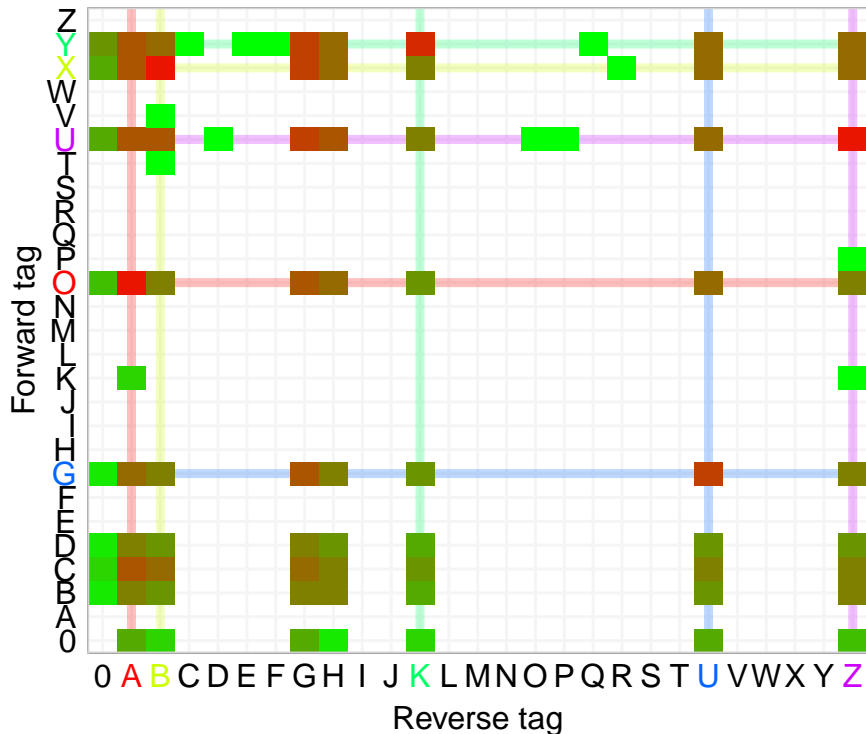
**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram53

]1 – 2]
]2 – 5]
]5 – 10]
]10 – 22]
]22 – 49]
]49 – 108]
]108 – 240]
]240 – 530]
]530 – 1174]
]1174 – 2597]
]2597 – 5748]
]5748 – 12718]

Forward tag

Reverse tag

**1 difference ISUs**

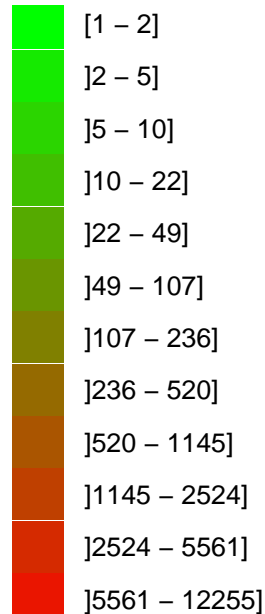non–critical mistag ● correctly labelled

Number of reads per sample
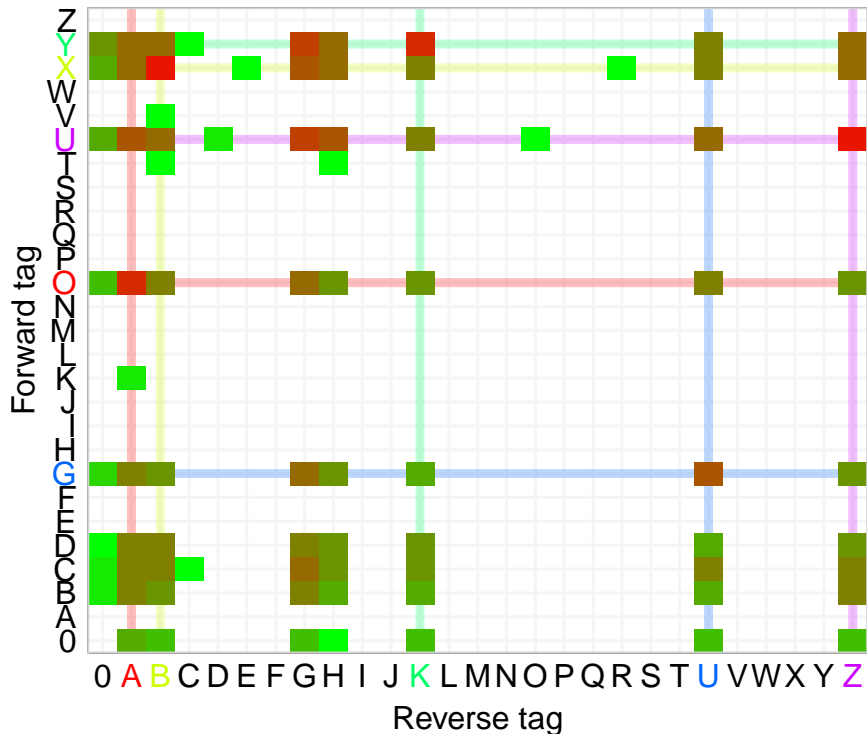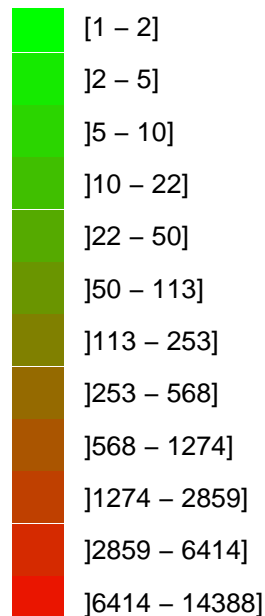
# Perfect match (0 difference ISU)

**SFA−126**
mock: random

foram52

| Color | Range |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 23] |
| | ]23 – 51] |
| | ]51 – 115] |
| | ]115 – 259] |
| | ]259 – 584] |
| | ]584 – 1318] |
| | ]1318 – 2972] |
| | ]2972 – 6704] |
| | ]6704 – 15123] |

Forward tag

Reverse tag

## 1 difference ISUs

○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

**SFA−126**
mock: random

foram88

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 47] |
| | ]47 – 101] |
| | ]101 – 219] |
| | ]219 – 474] |
| | ]474 – 1025] |
| | ]1025 – 2217] |
| | ]2217 – 4795] |
| | ]4795 – 10372] |

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag   ● correctly labelled

Number of reads per sample
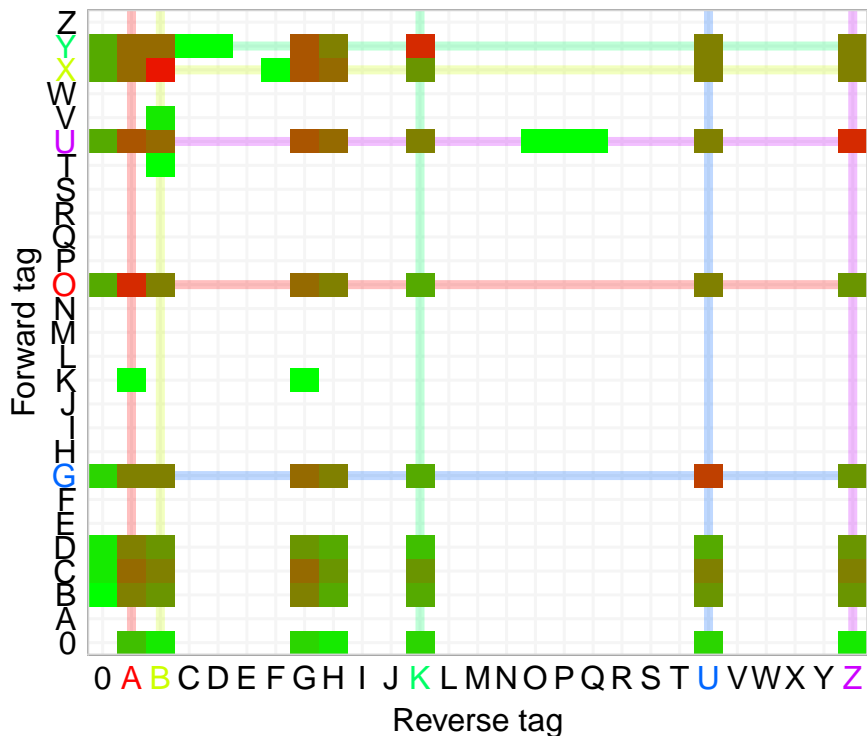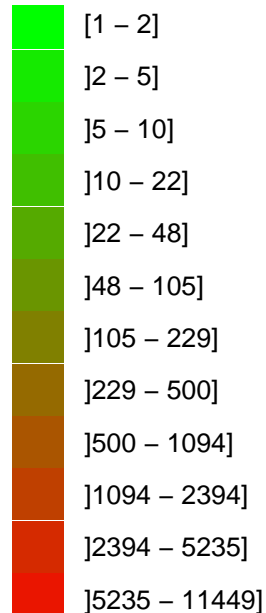
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram89

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 17]
]17 − 29]
]29 − 49]
]49 − 84]
]84 − 142]
]142 − 242]
]242 − 411]
]411 − 700]
]700 − 1190]

Forward tag

Reverse tag

**1 difference ISUs**

o non−critical mistag  • correctly labelled

Number of reads per sample

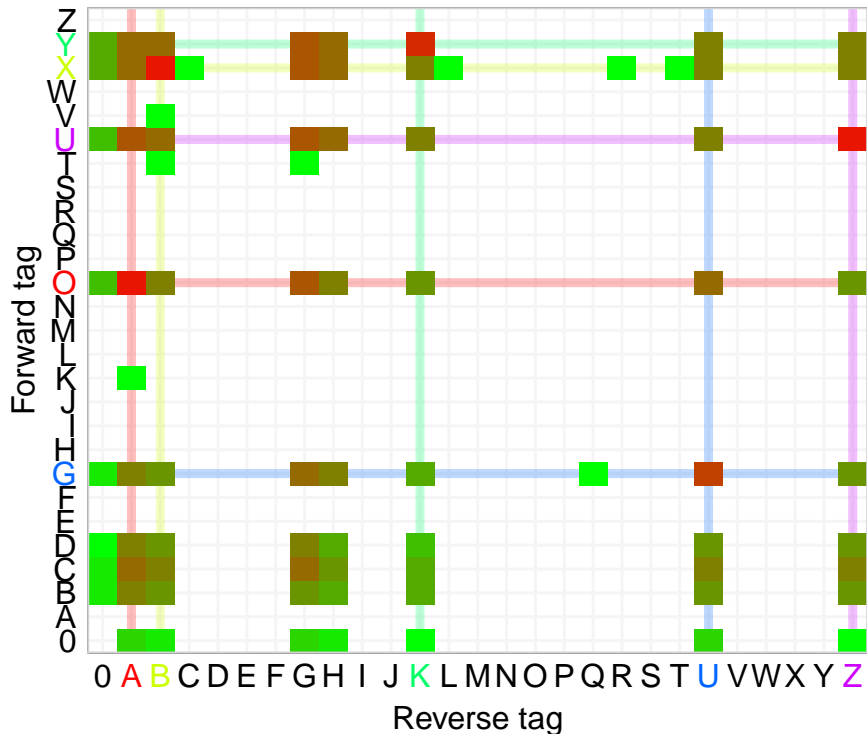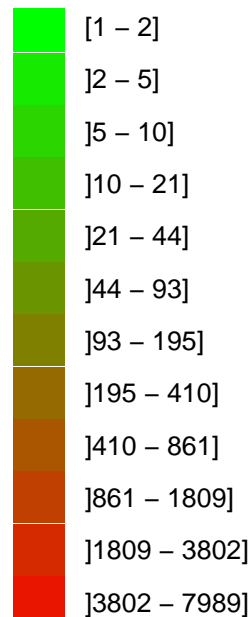**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram82

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 48] |
| | ]48 – 105] |
| | ]105 – 229] |
| | ]229 – 500] |
| | ]500 – 1093] |
| | ]1093 – 2390] |
| | ]2390 – 5226] |
| | ]5226 – 11426] |

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
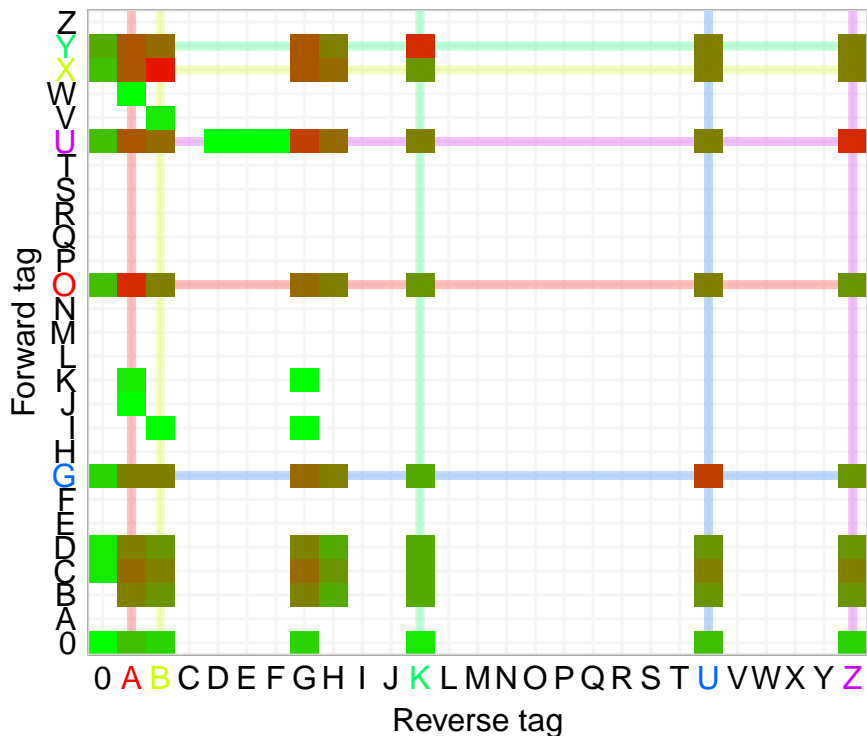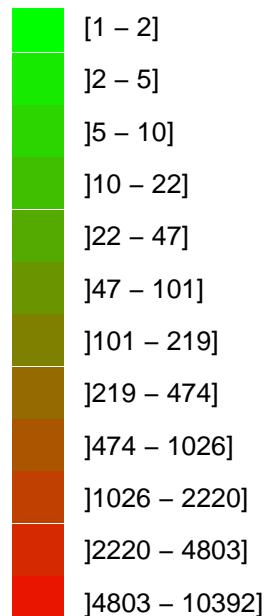
**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram80

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 48] |
| | ]48 – 106] |
| | ]106 – 233] |
| | ]233 – 513] |
| | ]513 – 1127] |
| | ]1127 – 2476] |
| | ]2476 – 5442] |
| | ]5442 – 11961] |

Forward tag

Reverse tag

**1 difference ISUs**

○ non–critical mistag   ● correctly labelled

Number of reads per sample
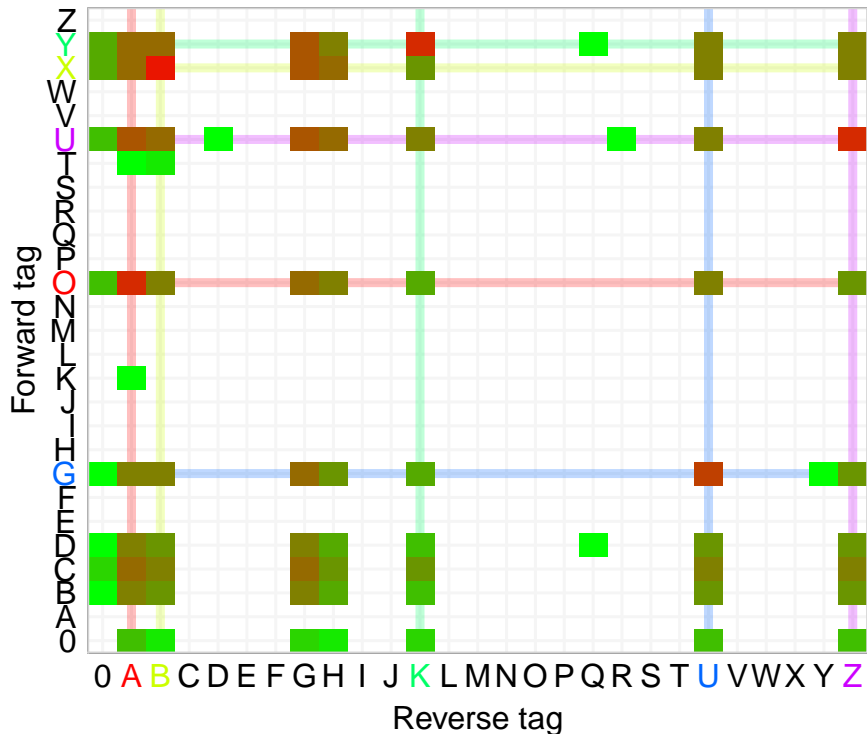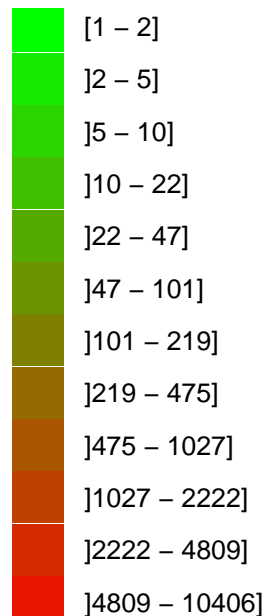
**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram81

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 20]
]20 – 39]
]39 – 76]
]76 – 151]
]151 – 297]
]297 – 585]
]585 – 1153]
]1153 – 2271]
]2271 – 4475]

Forward tag

Reverse tag

**1 difference ISUs**

non–critical mistag    correctly labelled

Number of reads per sample

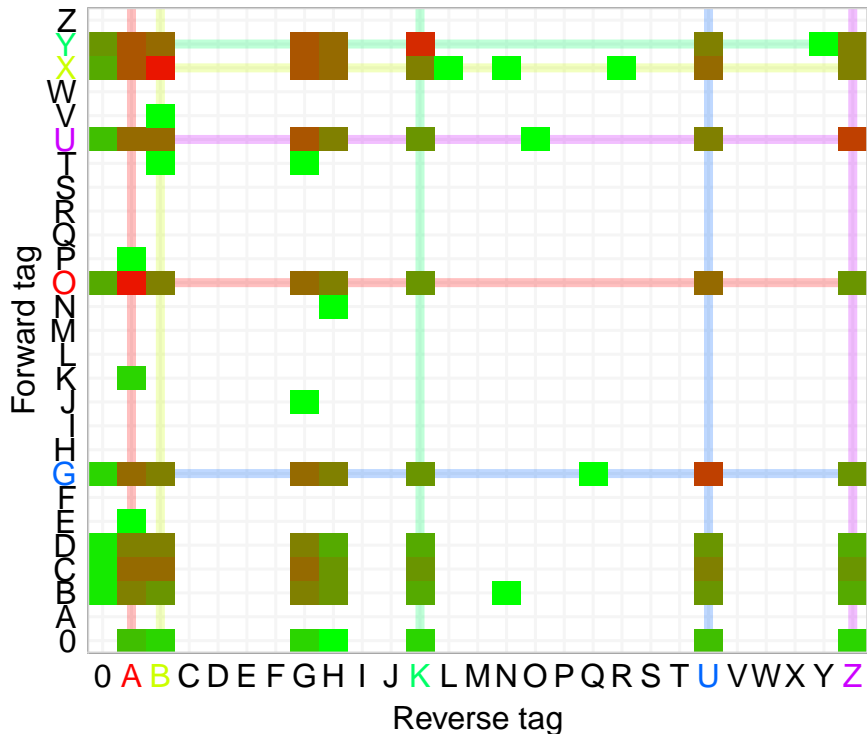**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram85

[1 − 2]
]2 − 5]
]5 − 10]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram38

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 22]
]22 – 47]
]47 – 103]
]103 – 224]
]224 – 487]
]487 – 1059]
]1059 – 2303]
]2303 – 5010]
]5010 – 10897]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
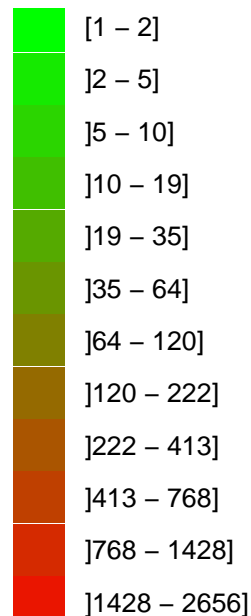
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram33

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 49] |
| | ]49 – 107] |
| | ]107 – 236] |
| | ]236 – 520] |
| | ]520 – 1145] |
| | ]1145 – 2524] |
| | ]2524 – 5561] |
| | ]5561 – 12255] |

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag    ● correctly labelled

Number of reads per sample
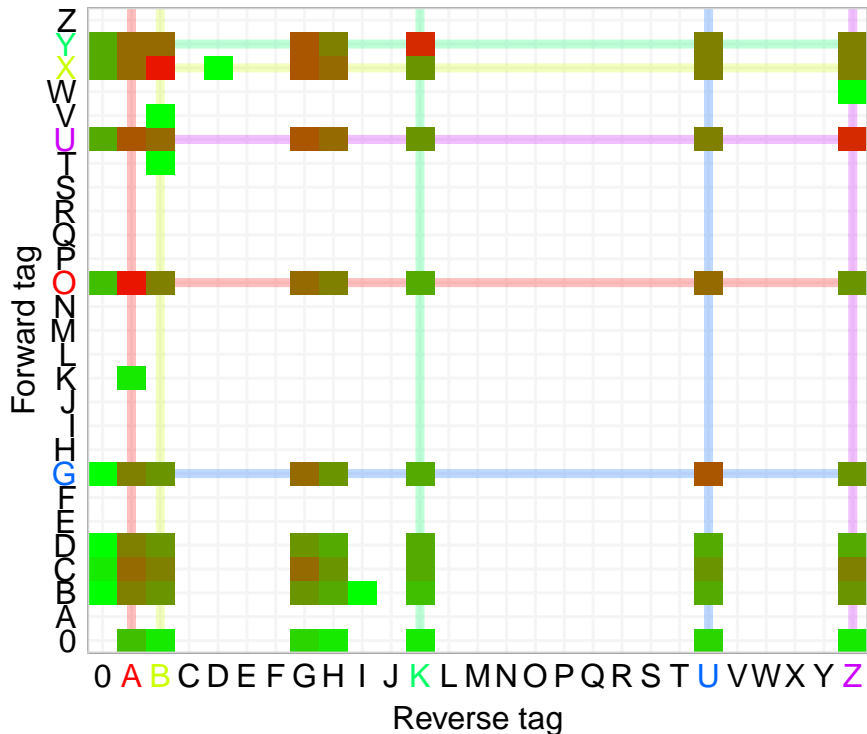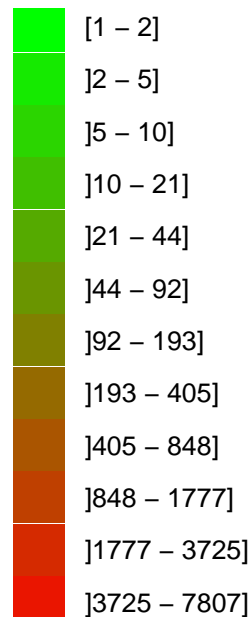
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram31

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 22]
]22 – 50]
]50 – 113]
]113 – 253]
]253 – 568]
]568 – 1274]
]1274 – 2859]
]2859 – 6414]
]6414 – 14388]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
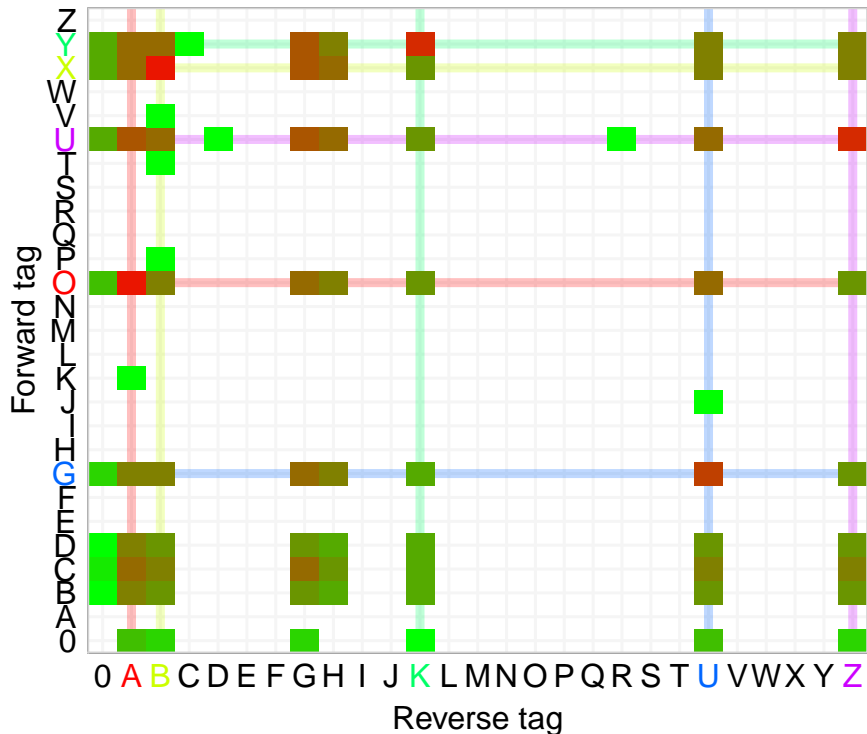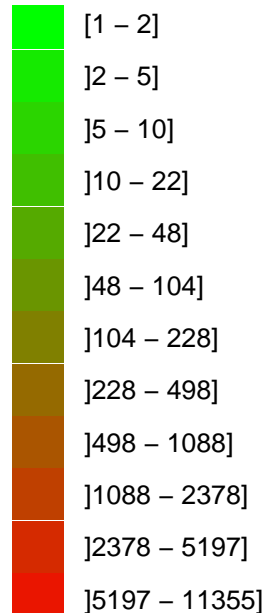
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram37

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 48] |
| | ]48 – 105] |
| | ]105 – 229] |
| | ]229 – 500] |
| | ]500 – 1094] |
| | ]1094 – 2394] |
| | ]2394 – 5235] |
| | ]5235 – 11449] |

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
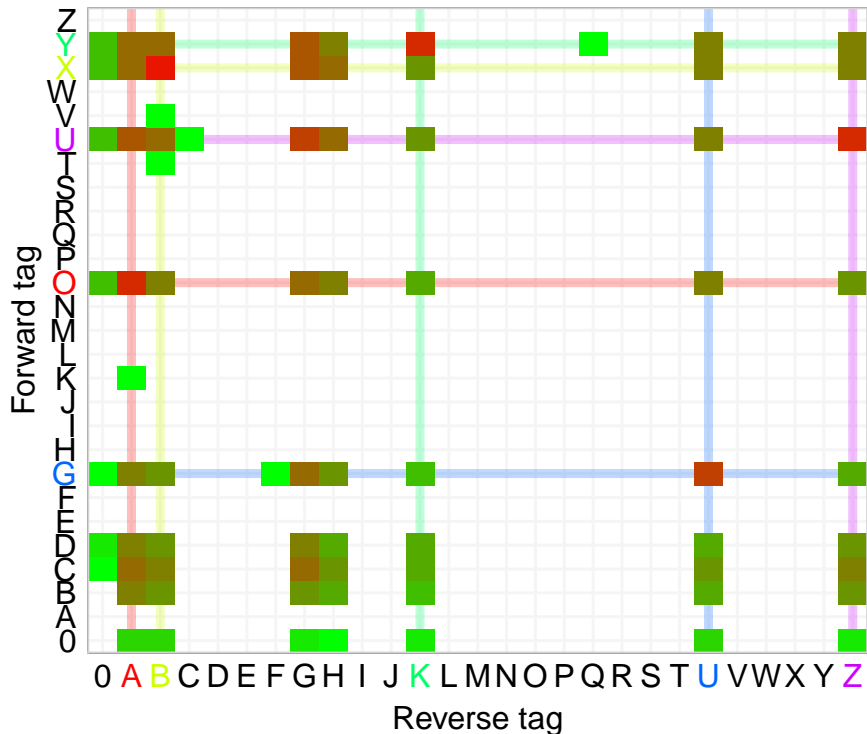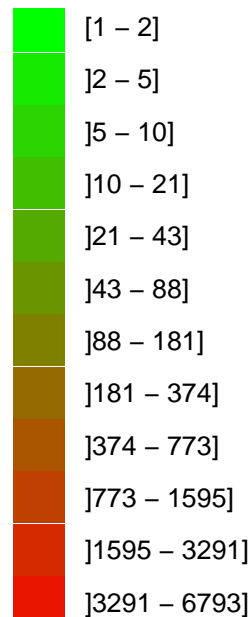
**Perfect match (0 difference ISU)**

SFA-126
mock: random

foram60

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 21]
]21 – 44]
]44 – 93]
]93 – 195]
]195 – 410]
]410 – 861]
]861 – 1809]
]1809 – 3802]
]3802 – 7989]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
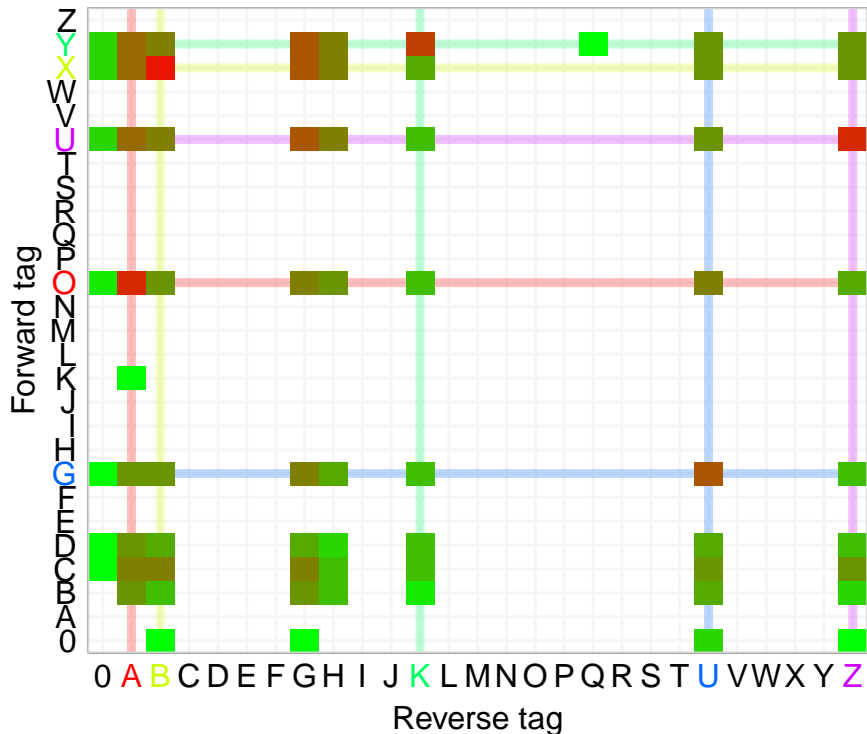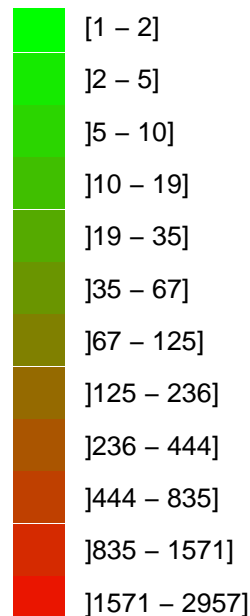
**Perfect match (0 difference ISU)**

SFA-126
mock: random

foram61

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 22]
]22 – 47]
]47 – 101]
]101 – 219]
]219 – 474]
]474 – 1026]
]1026 – 2220]
]2220 – 4803]
]4803 – 10392]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
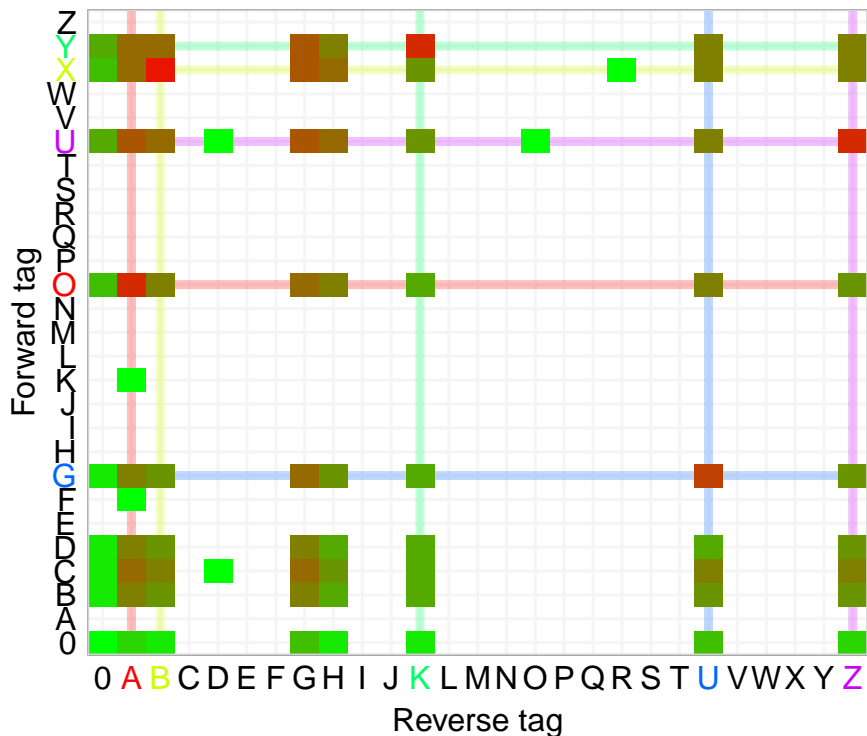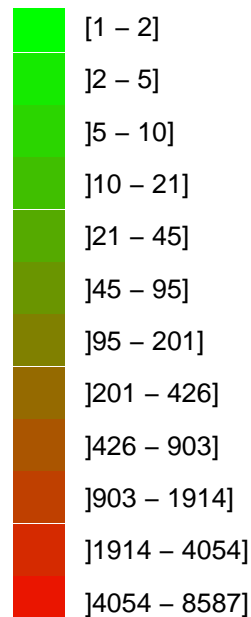
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram62

| | |
|---|---|
| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 22] |
| | ]22 – 47] |
| | ]47 – 101] |
| | ]101 – 219] |
| | ]219 – 475] |
| | ]475 – 1027] |
| | ]1027 – 2222] |
| | ]2222 – 4809] |
| | ]4809 – 10406] |

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

# Perfect match (0 difference ISU)

**SFA−126**
**mock: random**

foram64

| | |
|---|---|
| | [1 − 2] |
| | ]2 − 5] |
| | ]5 − 10] |
| | ]10 − 22] |
| | ]22 − 50] |
| | ]50 − 113] |
| | ]113 − 255] |
| | ]255 − 573] |
| | ]573 − 1287] |
| | ]1287 − 2892] |
| | ]2892 − 6498] |
| | ]6498 − 14601] |

Forward tag

Reverse tag

## 1 difference ISUs

○ non−critical mistag　● correctly labelled

Number of reads per sample
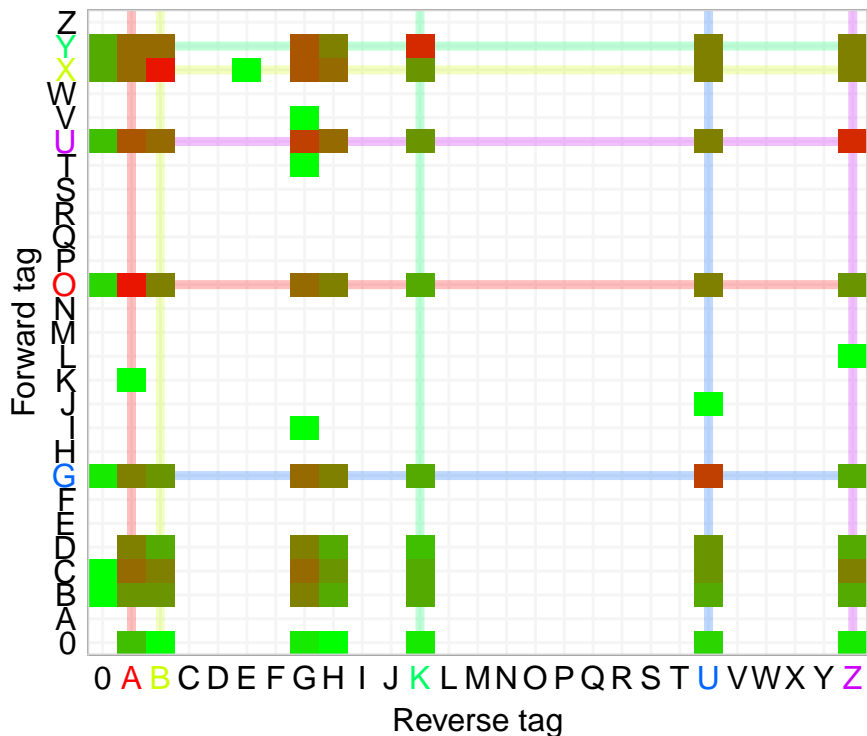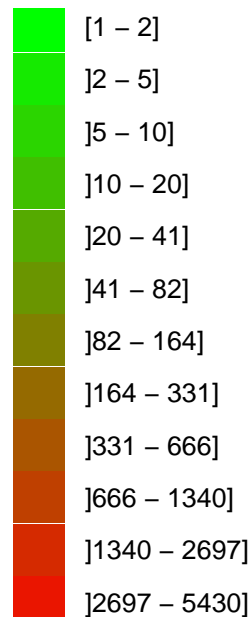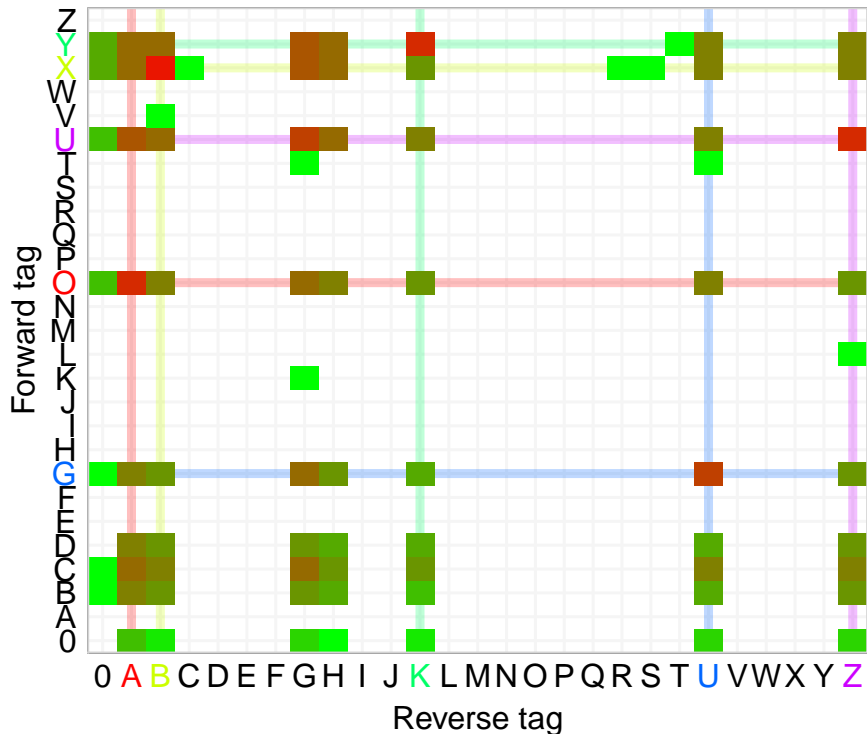
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram65

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 19]
]19 – 35]
]35 – 64]
]64 – 120]
]120 – 222]
]222 – 413]
]413 – 768]
]768 – 1428]
]1428 – 2656]

Forward tag

Reverse tag

**1 difference ISUs**
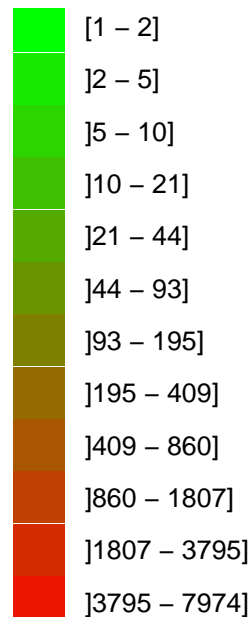
○ non−critical mistag ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram66

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 21]
]21 − 44]
]44 − 92]
]92 − 193]
]193 − 405]
]405 − 848]
]848 − 1777]
]1777 − 3725]
]3725 − 7807]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
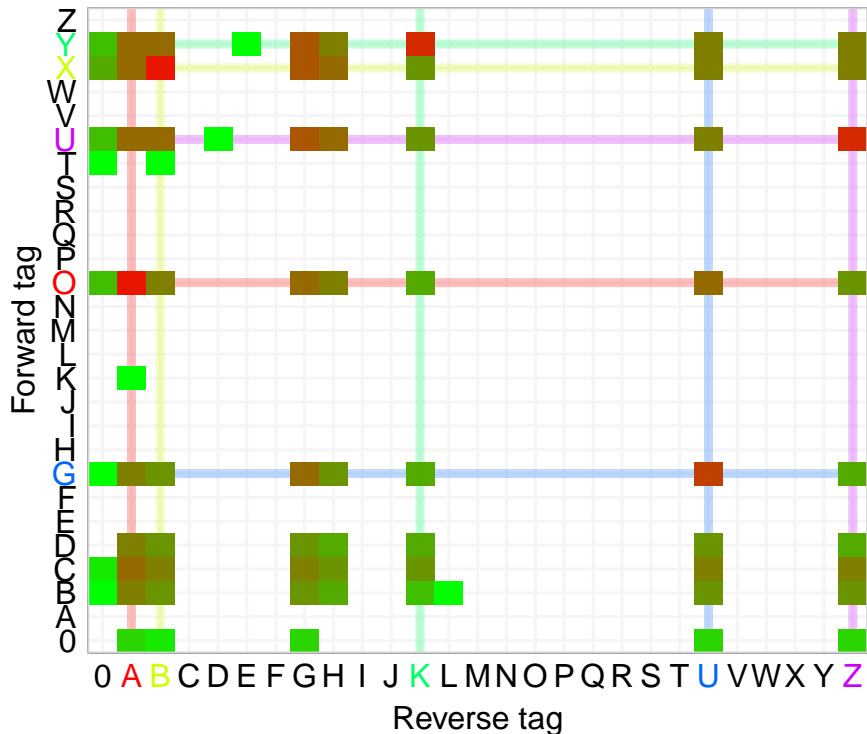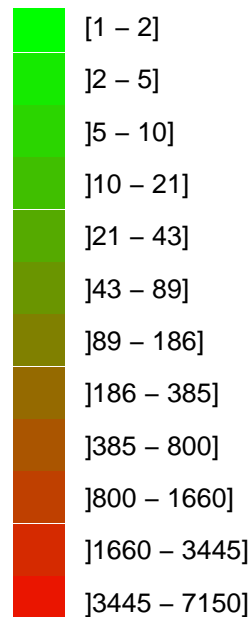
**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram67

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 22]
]22 – 48]
]48 – 104]
]104 – 228]
]228 – 498]
]498 – 1088]
]1088 – 2378]
]2378 – 5197]
]5197 – 11355]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag    ● correctly labelled

Number of reads per sample
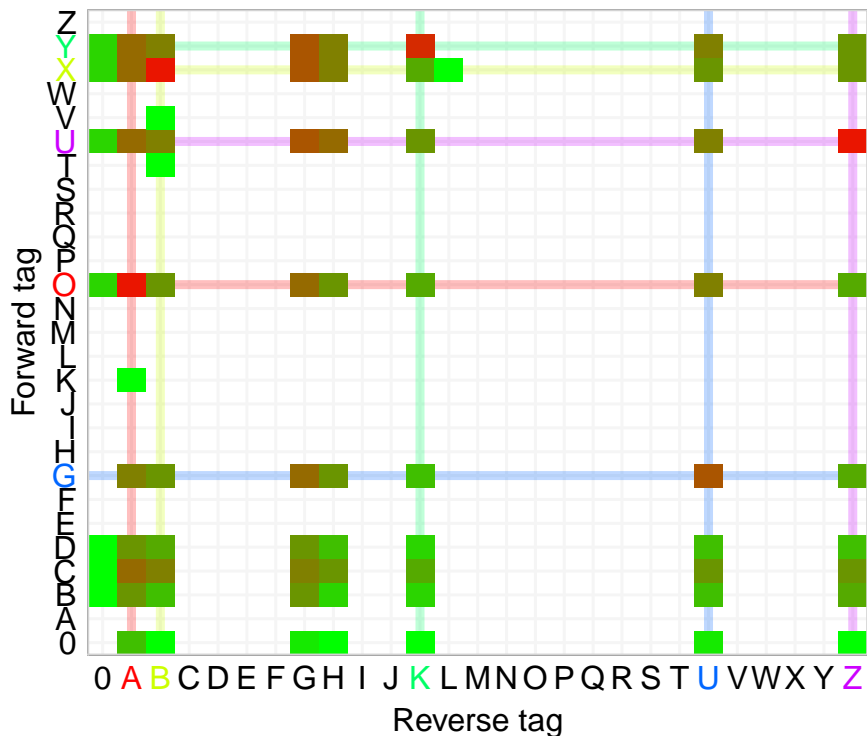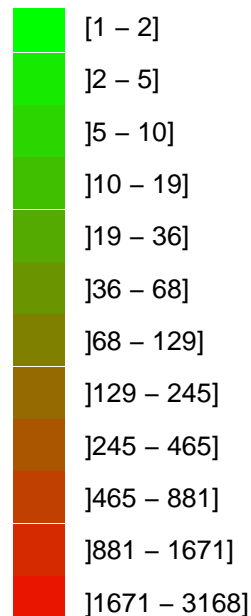
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram68

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 21]
]21 − 43]
]43 − 88]
]88 − 181]
]181 − 374]
]374 − 773]
]773 − 1595]
]1595 − 3291]
]3291 − 6793]

**1 difference ISUs**

○ non−critical mistag   ● correctly labelled

Number of reads per sample
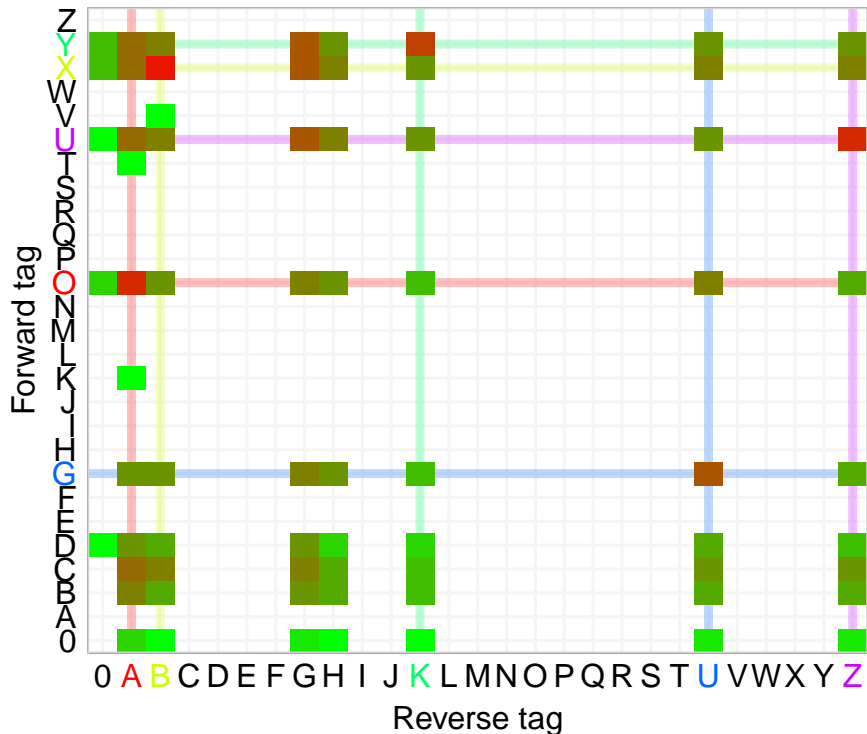
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram69

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 19]
]19 − 35]
]35 − 67]
]67 − 125]
]125 − 236]
]236 − 444]
]444 − 835]
]835 − 1571]
]1571 − 2957]

Forward tag

Reverse tag

**1 difference ISUs**

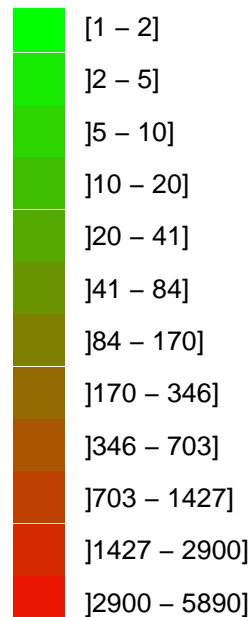○ non−critical mistag    ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram48

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 21]
]21 − 45]
]45 − 95]
]95 − 201]
]201 − 426]
]426 − 903]
]903 − 1914]
]1914 − 4054]
]4054 − 8587]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag   ● correctly labelled

Number of reads per sample
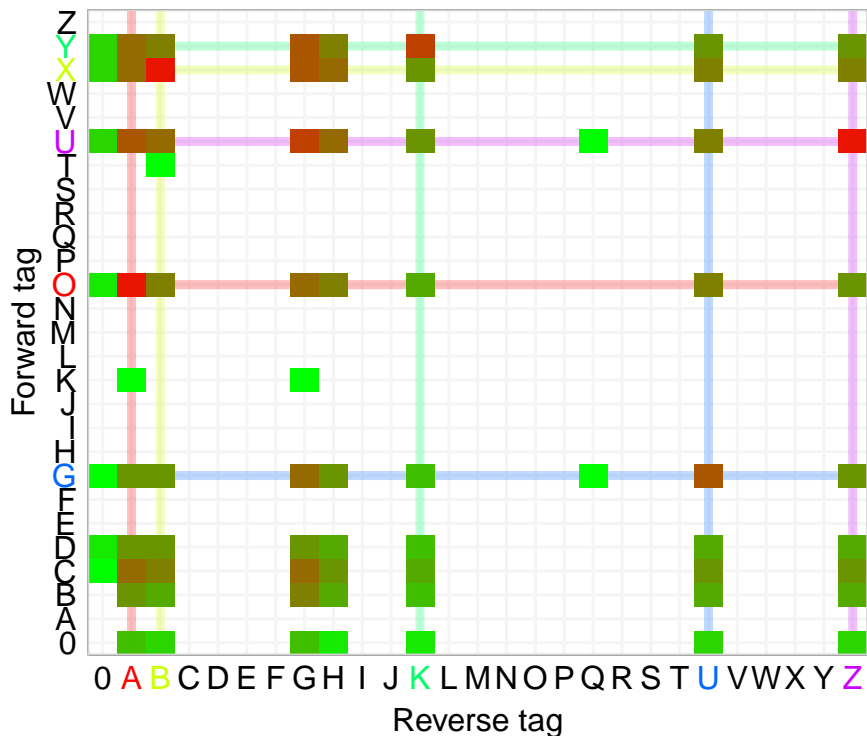
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram46

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 20]
]20 − 41]
]41 − 82]
]82 − 164]
]164 − 331]
]331 − 666]
]666 − 1340]
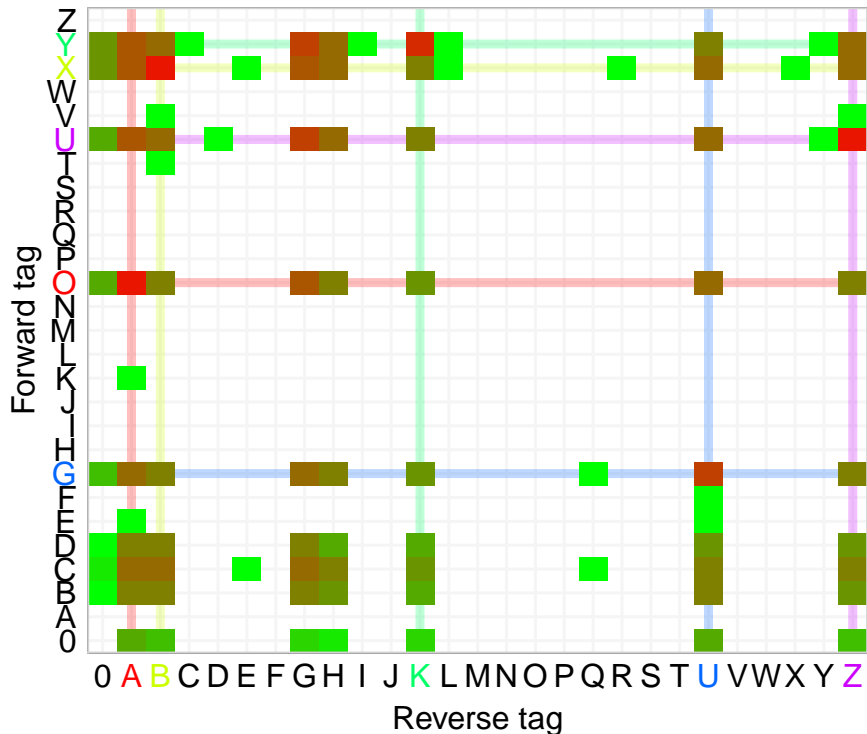]1340 − 2697]
]2697 − 5430]

Forward tag

Reverse tag

**1 difference ISUs**
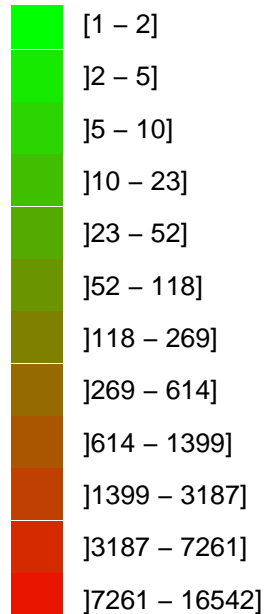
○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: random
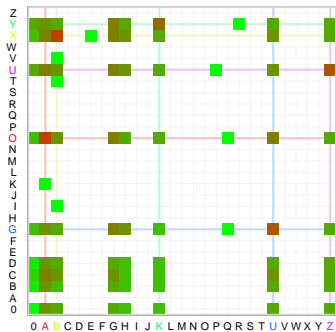
foram44

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 21]
]21 − 44]
]44 − 93]
]93 − 195]
]195 − 409]
]409 − 860]
]860 − 1807]
]1807 − 3795]
]3795 − 7974]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag   ● correctly labelled

Number of reads per sample
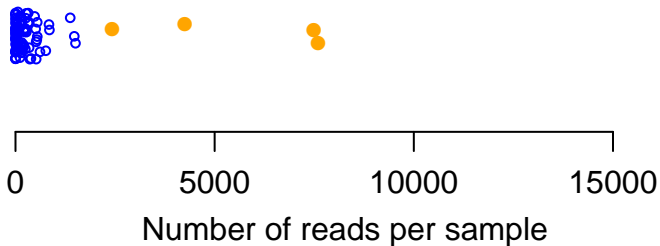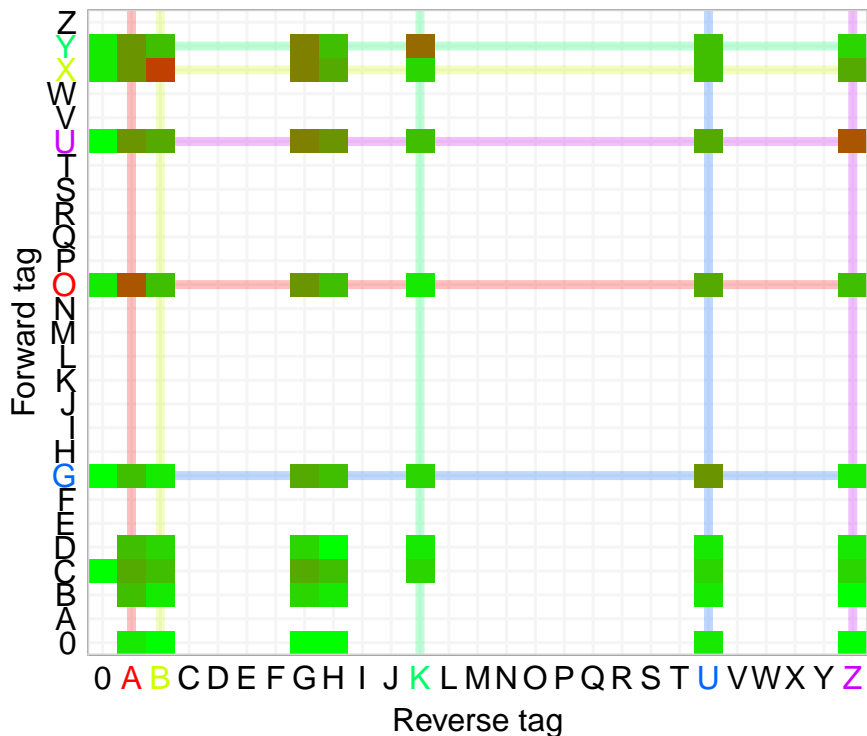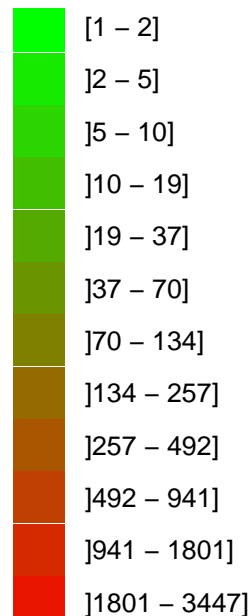
**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram45

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 21]
]21 – 43]
]43 – 89]
]89 – 186]
]186 – 385]
]385 – 800]
]800 – 1660]
]1660 – 3445]
]3445 – 7150]

Forward tag

Reverse tag

**1 difference ISUs**

non−critical mistag • correctly labelled

Number of reads per sample

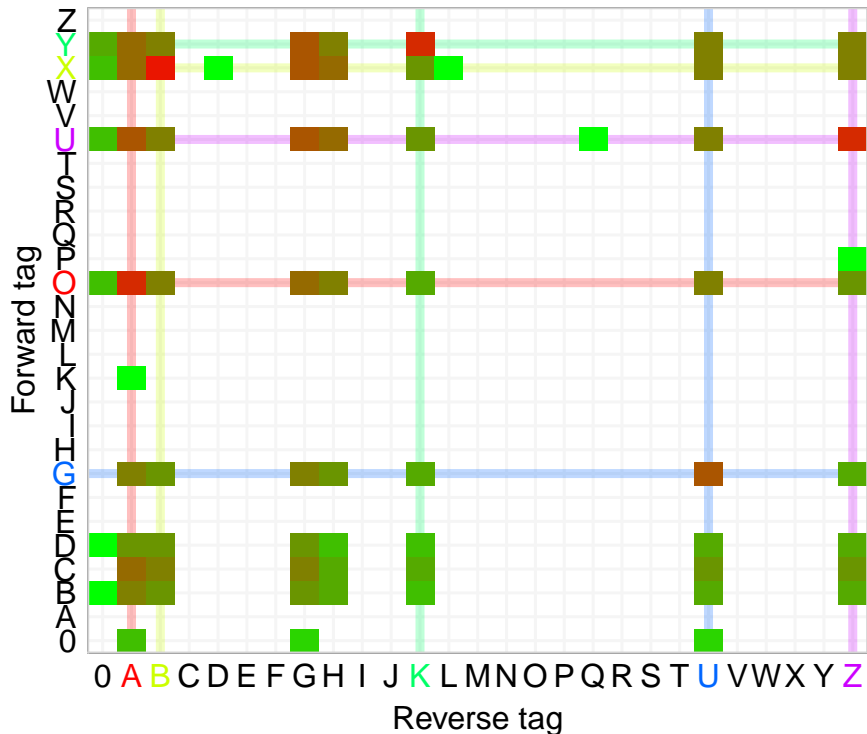**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram77

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 19]
]19 – 36]
]36 – 68]
]68 – 129]
]129 – 245]
]245 – 465]
]465 – 881]
]881 – 1671]
]1671 – 3168]

Forward tag

Reverse tag

**1 difference ISUs**

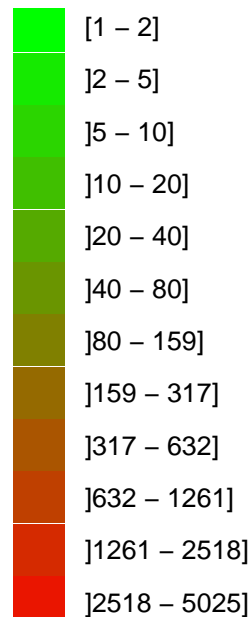○ non−critical mistag    ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA–126
mock: random

foram75

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 20]
]20 – 41]
]41 – 84]
]84 – 170]
]170 – 346]
]346 – 703]
]703 – 1427]
]1427 – 2900]
]2900 – 5890]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag   ● correctly labelled

Number of reads per sample
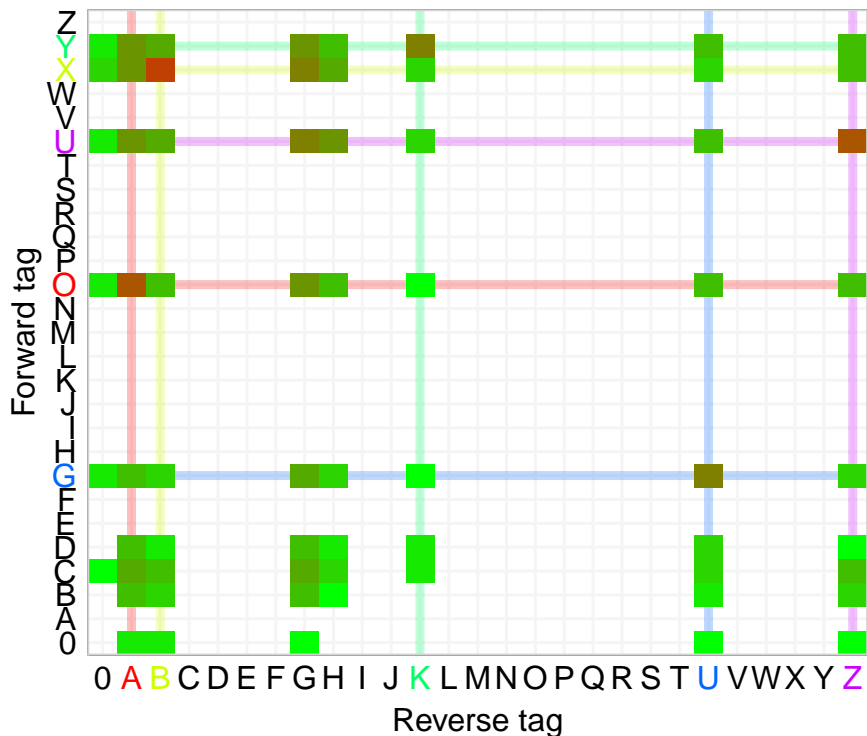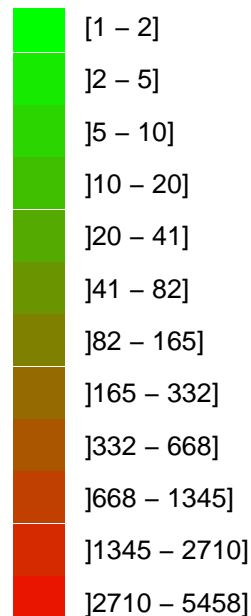
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram74

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 20]
]20 – 39]
]39 – 78]
]78 – 154]
]154 – 304]
]304 – 602]
]602 – 1192]
]1192 – 2361]
]2361 – 4674]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
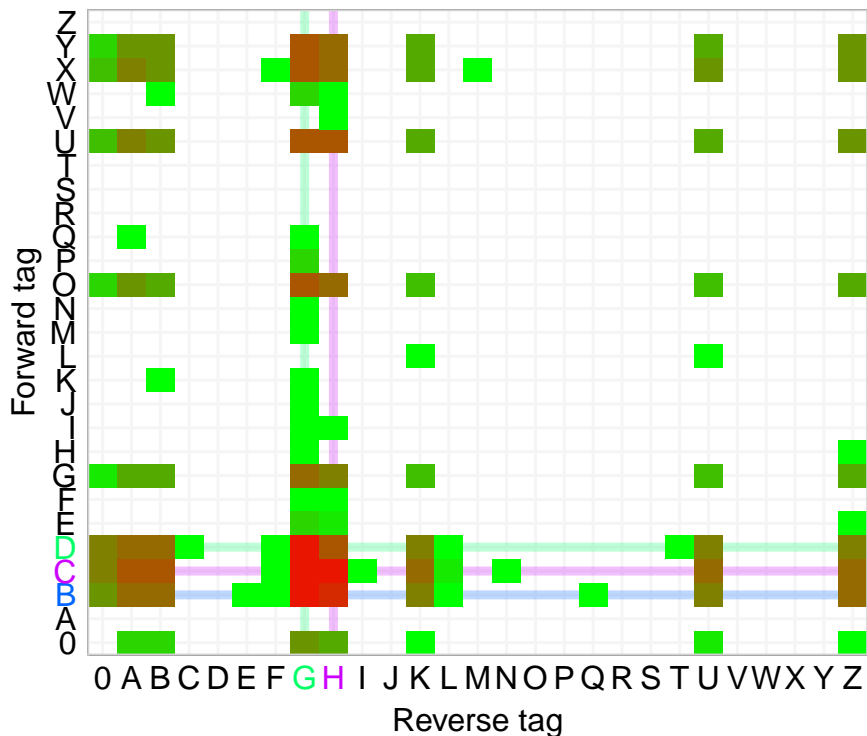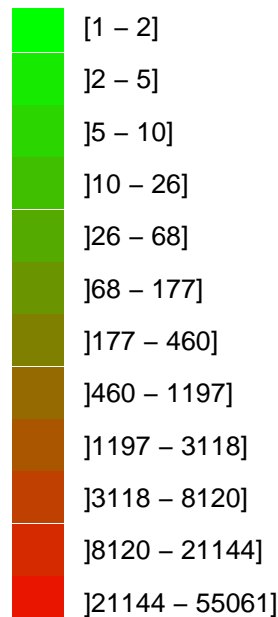
**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram71

| | |
|---|---|
| | [1 − 2] |
| | ]2 − 5] |
| | ]5 − 10] |
| | ]10 − 23] |
| | ]23 − 52] |
| | ]52 − 118] |
| | ]118 − 269] |
| | ]269 − 614] |
| | ]614 − 1399] |
| | ]1399 − 3187] |
| | ]3187 − 7261] |
| | ]7261 − 16542] |

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag    ● correctly labelled

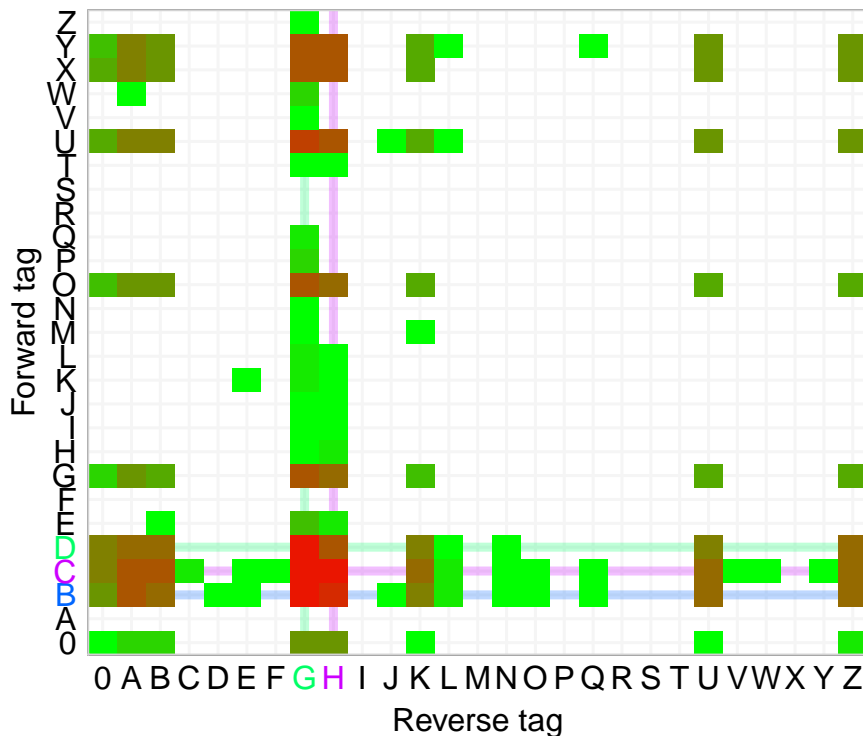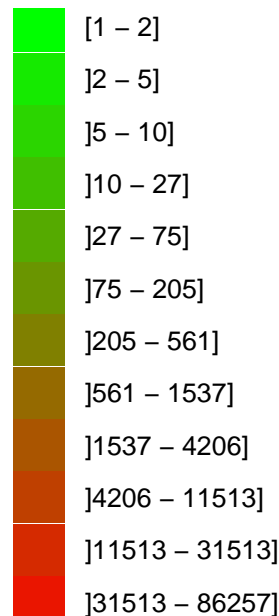Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram70

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 19]
]19 − 37]
]37 − 70]
]70 − 134]
]134 − 257]
]257 − 492]
]492 − 941]
]941 − 1801]
]1801 − 3447]

Forward tag

Reverse tag

**1 difference ISUs**

non−critical mistag    correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram79

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 20]
]20 − 40]
]40 − 80]
]80 − 159]
]159 − 317]
]317 − 632]
]632 − 1261]
]1261 − 2518]
]2518 − 5025]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: random

foram78

- [1 − 2]
- ]2 − 5]
- ]5 − 10]
- ]10 − 20]
- ]20 − 41]
- ]41 − 82]
- ]82 − 165]
- ]165 − 332]
- ]332 − 668]
- ]668 − 1345]
- ]1345 − 2710]
- ]2710 − 5458]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

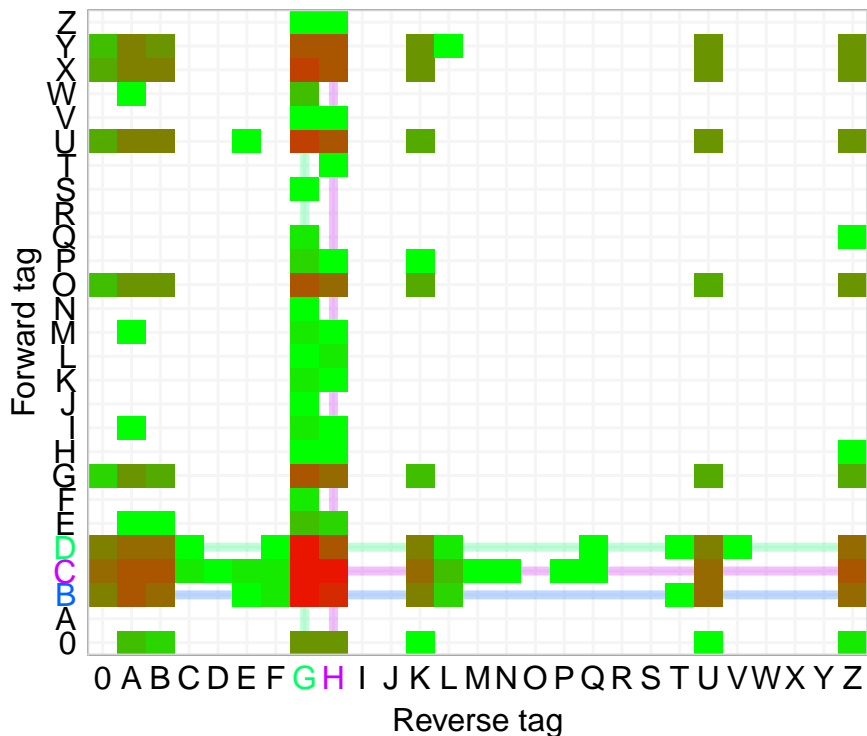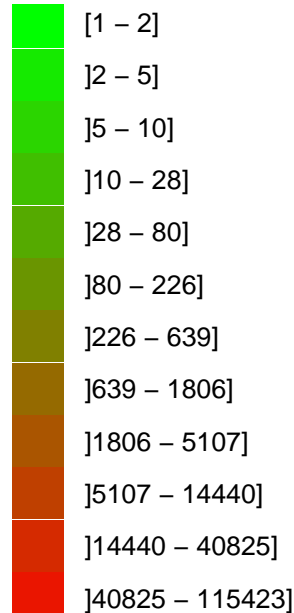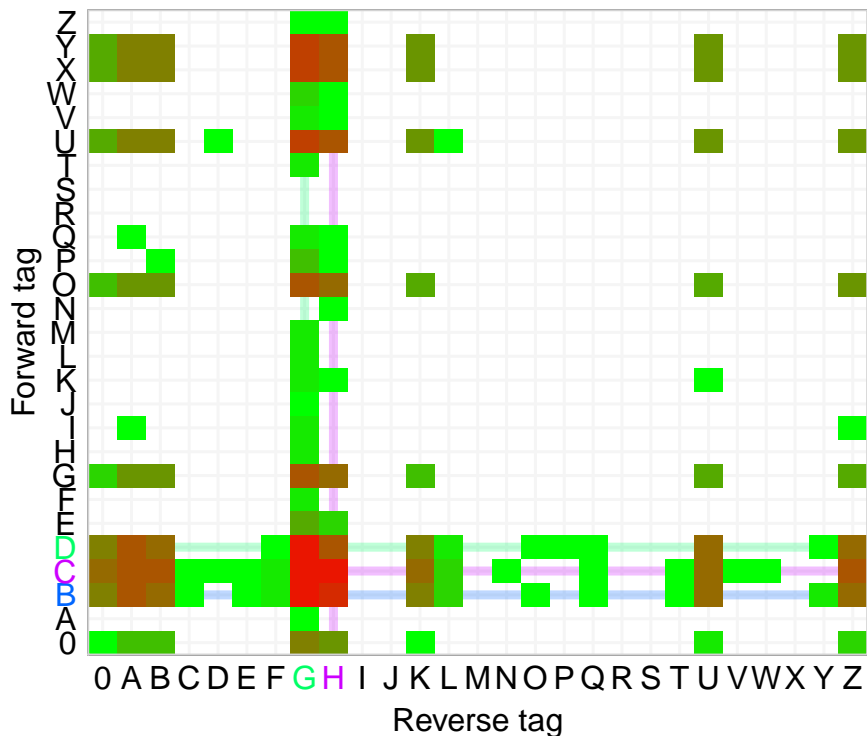**Perfect match (0 difference ISU)**

SFA−126
mock: even

foram26

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 26]
]26 − 68]
]68 − 177]
]177 − 460]
]460 − 1197]
]1197 − 3118]
]3118 − 8120]
]8120 − 21144]
]21144 − 55061]

Forward tag

Reverse tag
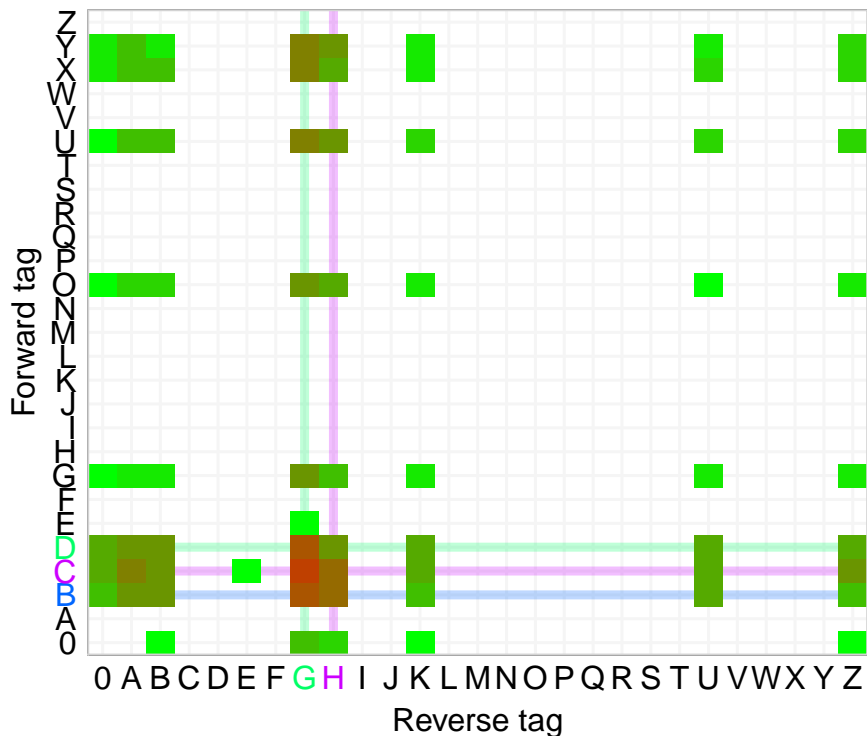
**1 difference ISUs**
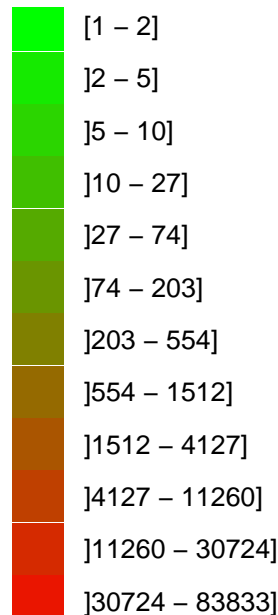
○ non−critical mistag   ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: even
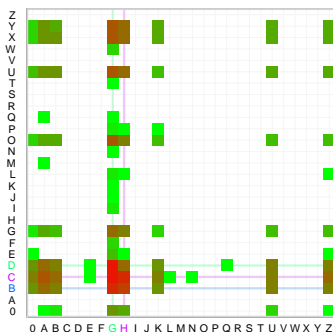
foram27

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 27]
]27 − 75]
]75 − 205]
]205 − 561]
]561 − 1537]
]1537 − 4206]
]4206 − 11513]
]11513 − 31513]
]31513 − 86257]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample
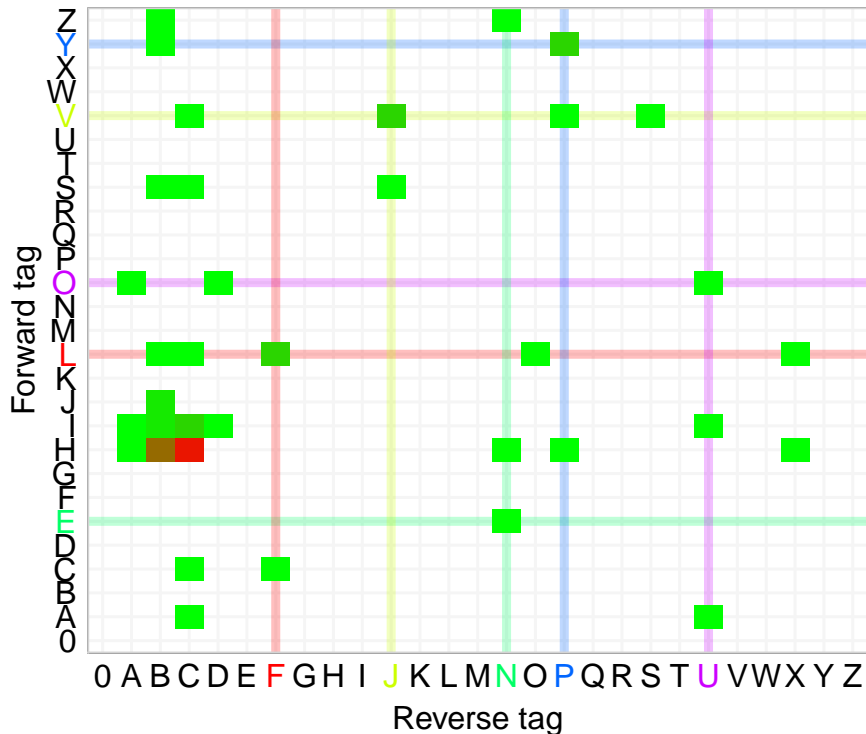
# Perfect match (0 difference ISU)

**SFA−126**
mock: even

foram23

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 28]
]28 – 80]
]80 – 226]
]226 – 639]
]639 – 1806]
]1806 – 5107]
]5107 – 14440]
]14440 – 40825]
]40825 – 115423]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag    ● correctly labelled

Number of reads per sample

# Perfect match (0 difference ISU)

Forward tag / Reverse tag

**SFA−126**
mock: even

foram7

| | |
|---|---|
| | [1 − 2] |
| | ]2 − 5] |
| | ]5 − 10] |
| | ]10 − 28] |
| | ]28 − 81] |
| | ]81 − 229] |
| | ]229 − 649] |
| | ]649 − 1844] |
| | ]1844 − 5234] |
| | ]5234 − 14859] |
| | ]14859 − 42182] |
| | ]42182 − 119747] |

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−126
mock: even

foram18

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 27]
]27 − 74]
]74 − 203]
]203 − 554]
]554 − 1512]
]1512 − 4127]
]4127 − 11260]
]11260 − 30724]
]30724 − 83833]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−125
mock: hhhhl
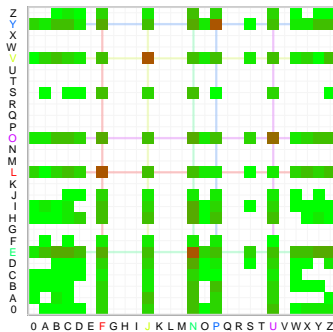
foram88 (I)

[1 − 2]
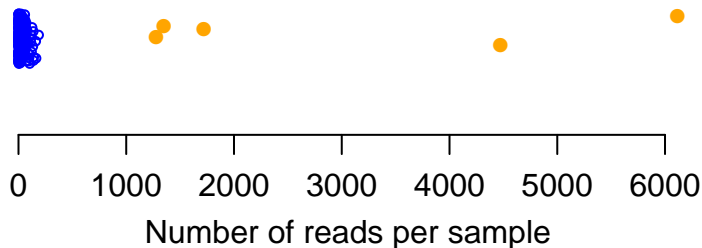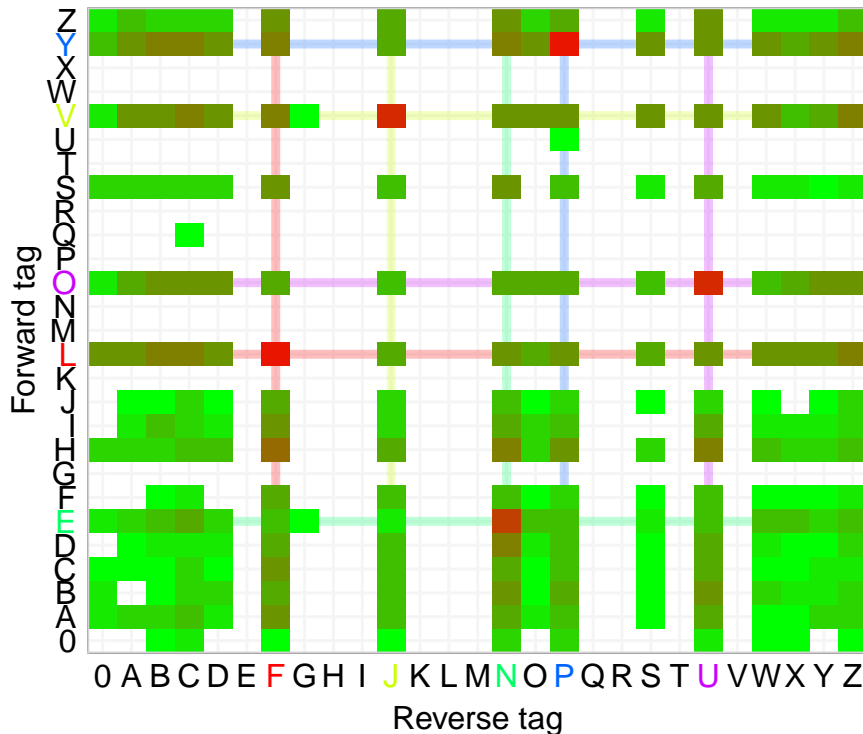]2 − 5]
]5 − 10]
]10 − 12]
]12 − 14]
]14 − 17]
]17 − 20]
]20 − 24]
]24 − 28]
]28 − 34]
]34 − 40]
]40 − 48]

Forward tag

Reverse tag

**1 difference ISUs**

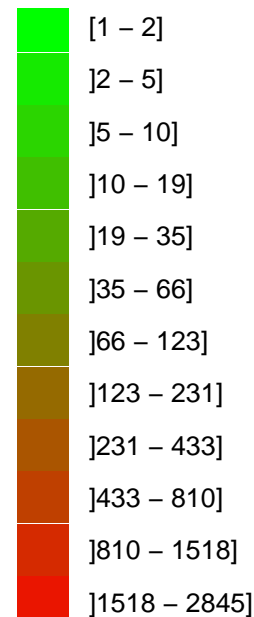non−critical mistag    correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−125

mock: hhhhl

foram33 (h)

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 21]
]21 − 42]
]42 − 87]
]87 − 178]
]178 − 366]
]366 − 751]
]751 − 1543]
]1543 − 3169]
]3169 − 6510]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA–125
mock: hhhhl

foram31 (h)
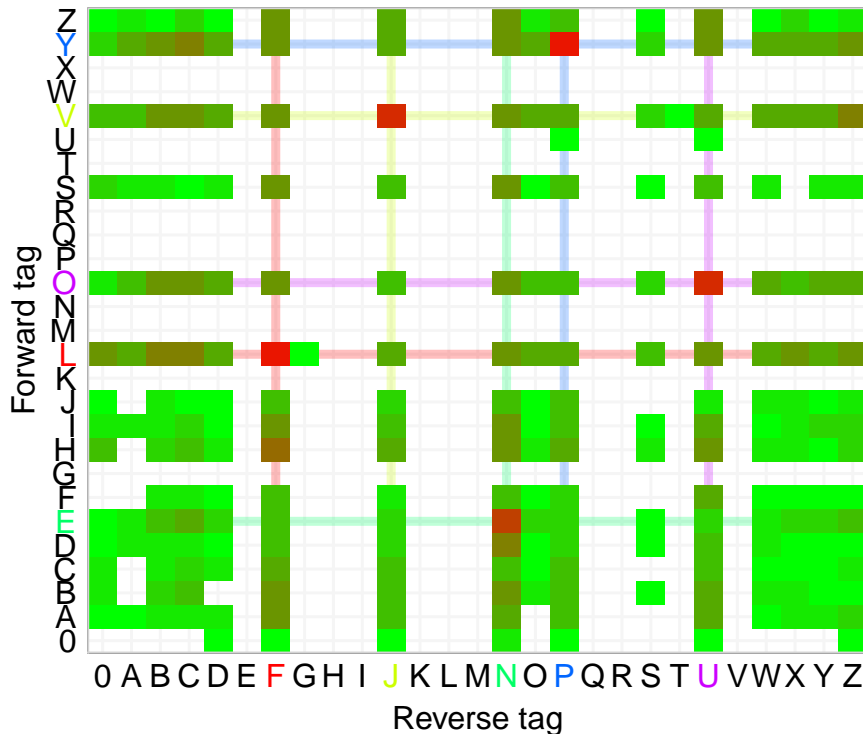
[1 – 2]
]2 – 5]
]5 – 10]
]10 – 20]
]20 – 42]
]42 – 85]
]85 – 173]
]173 – 353]
]353 – 720]
]720 – 1470]
]1470 – 2998]
]2998 – 6115]

Forward tag

Reverse tag

**1 difference ISUs**

○ non–critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA–125
mock: hhhhl

foram37 (h)
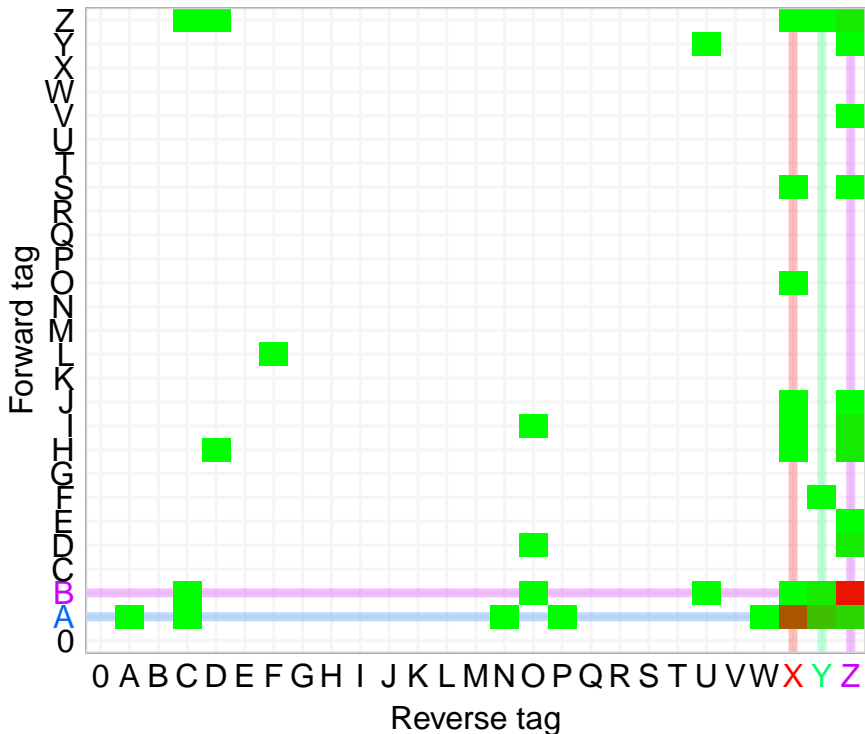
[1 – 2]
]2 – 5]
]5 – 10]
]10 – 19]
]19 – 35]
]35 – 66]
]66 – 123]
]123 – 231]
]231 – 433]
]433 – 810]
]810 – 1518]
]1518 – 2845]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−125
mock: hhhhl

foram45 (h)

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 18]
]18 − 34]
]34 − 63]
]63 − 116]
]116 − 213]
]213 − 394]
]394 − 726]
]726 − 1339]
]1339 − 2469]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−125
mock: hlIII

foram27  (I)
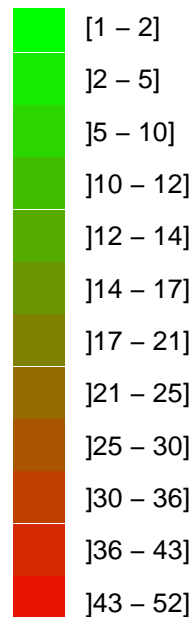
[1 − 2]
]2 − 5]
]5 − 10]
]10 − 12]
]12 − 14]
]14 − 17]
]17 − 21]
]21 − 25]
]25 − 30]
]30 − 36]
]36 − 43]
]43 − 52]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

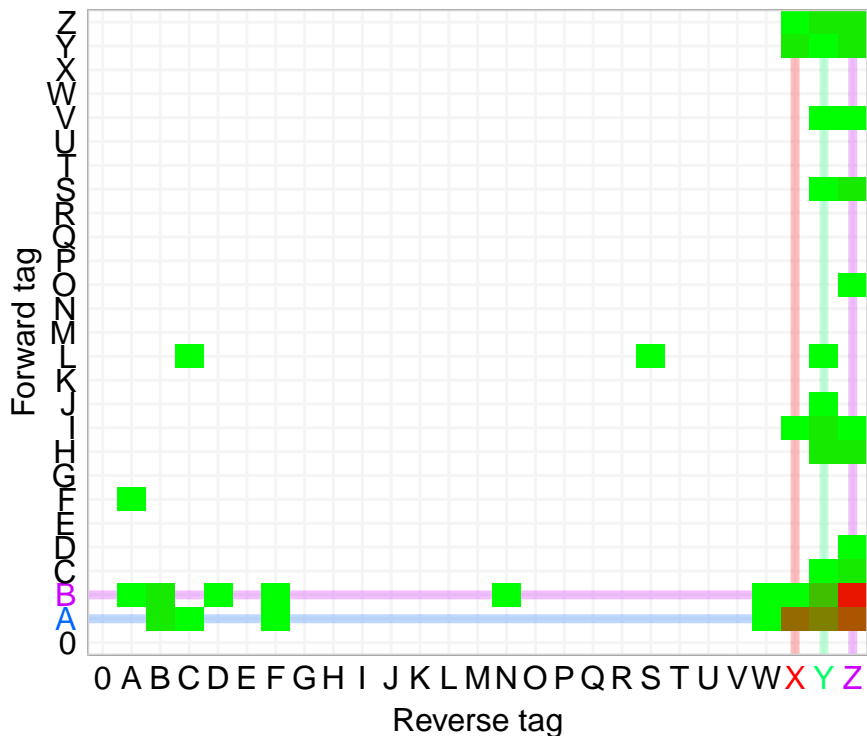Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−125
mock: hllll

foram23  (I)

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 12]
]12 − 14]
]14 − 17]
]17 − 20]
]20 − 24]
]24 − 29]
]29 − 35]
]35 − 42]
]42 − 50]

Forward tag

Reverse tag

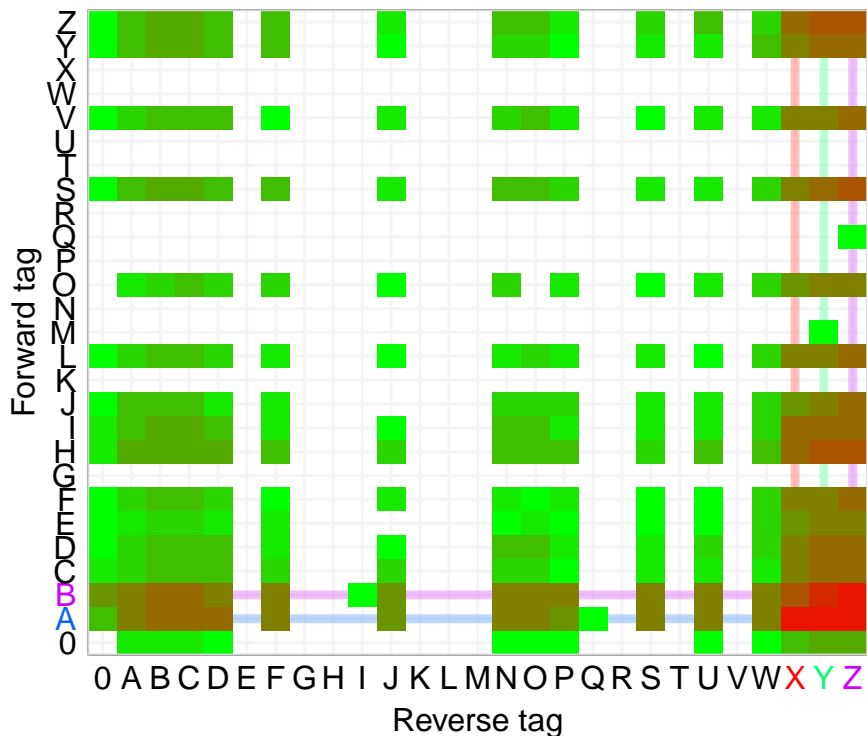**1 difference ISUs**

○ non−critical mistag  ● correctly labelled
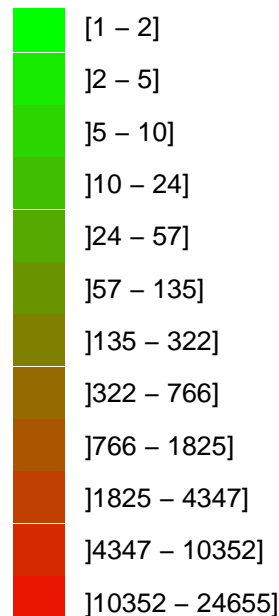
Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−125

mock: hlllI

foram7    (h)

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 24]
]24 − 57]
]57 − 135]
]135 − 322]
]322 − 766]
]766 − 1825]
]1825 − 4347]
]4347 − 10352]
]10352 − 24655]

Forward tag

Reverse tag

**1 difference ISUs**

o non−critical mistag ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−125
mock: hllll

foram18 (I)

[1 − 2]
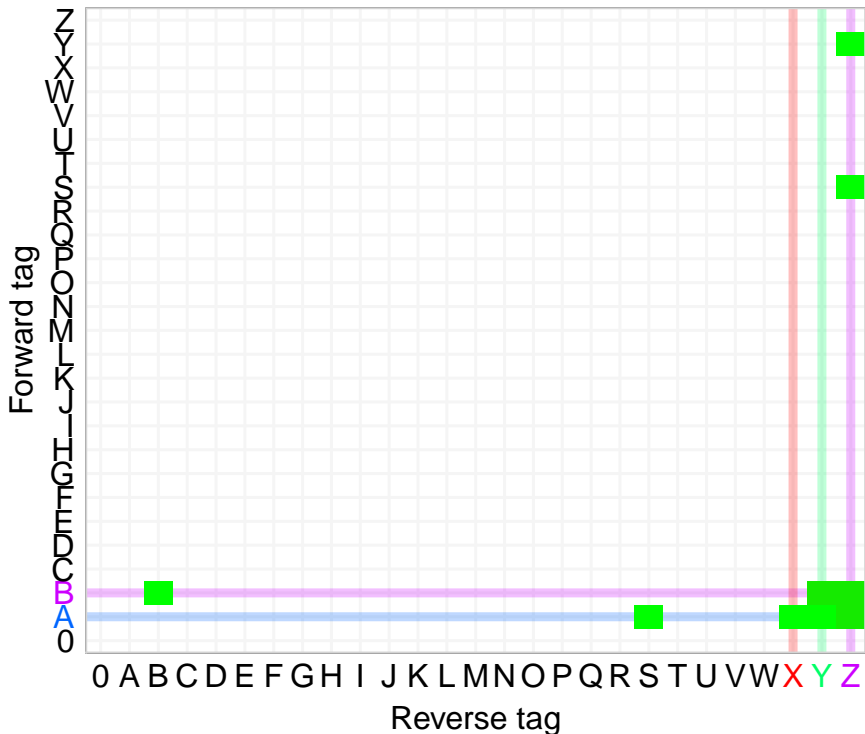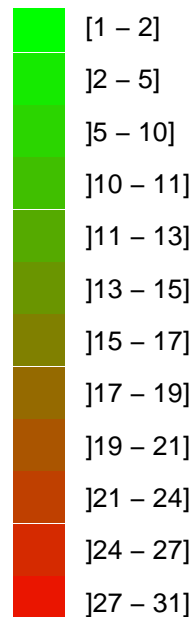]2 − 5]
]5 − 10]
]10 − 11]
]11 − 13]
]13 − 15]
]15 − 17]
]17 − 19]
]19 − 21]
]21 − 24]
]24 − 27]
]27 − 31]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**

SFA−125
mock: hhmll

foram59 (h)

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 23]
]23 – 55]
]55 – 127]
]127 – 297]
]297 – 693]
]693 – 1619]
]1619 – 3780]
]3780 – 8824]
]8824 – 20599]

**1 difference ISUs**

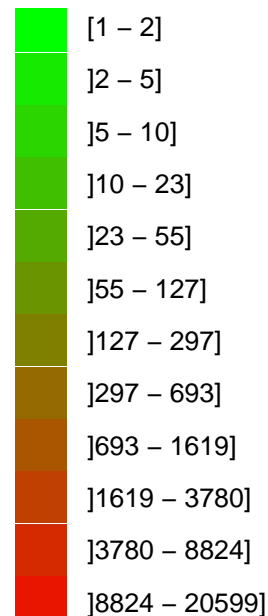○ non−critical mistag  ● correctly labelled

Number of reads per sample

## Perfect match (0 difference ISU)

SFA−125

mock: hhmll

foram54  (h)

[1 – 2]

]2 – 5]

]5 – 10]

]10 – 23]

]23 – 52]

]52 – 118]

]118 – 268]

]268 – 610]

]610 – 1389]

]1389 – 3161]

]3161 – 7193]

]7193 – 16369]

Forward tag

Reverse tag

### 1 difference ISUs

○ non−critical mistag   ● correctly labelled

Number of reads per sample

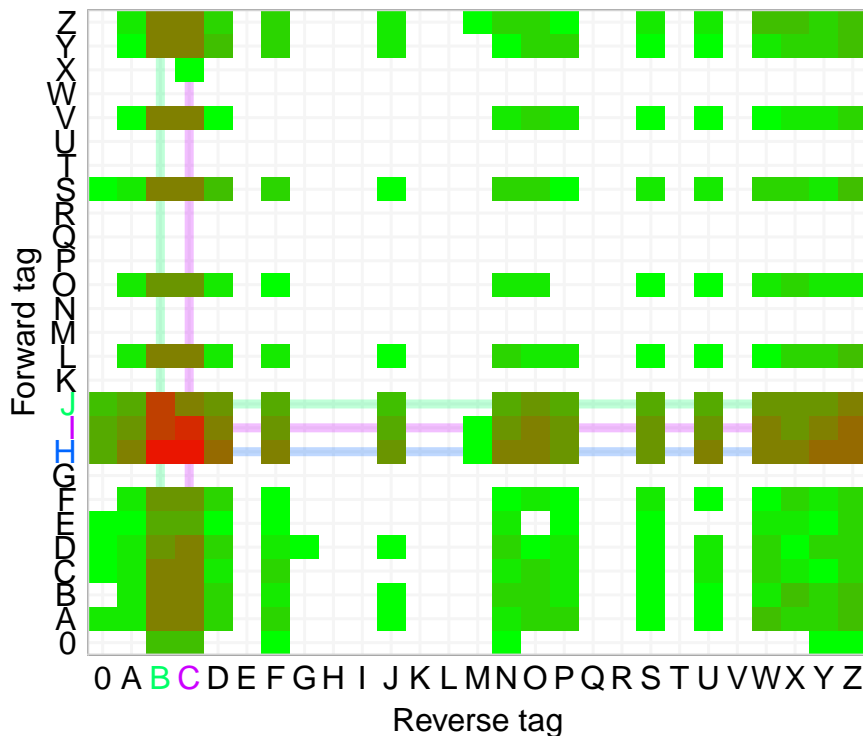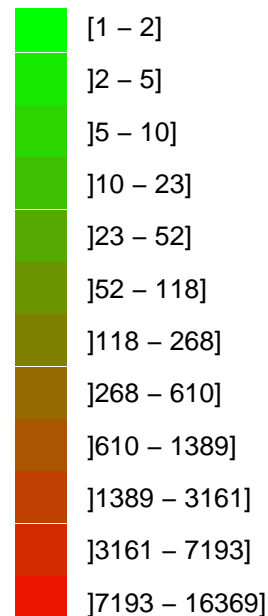**Perfect match (0 difference ISU)**

SFA−125
mock: hhmll
foram89  (I)

Forward tag / Reverse tag

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 11]

**1 difference ISUs**

○ non−critical mistag   ● correctly labelled

Number of reads per sample

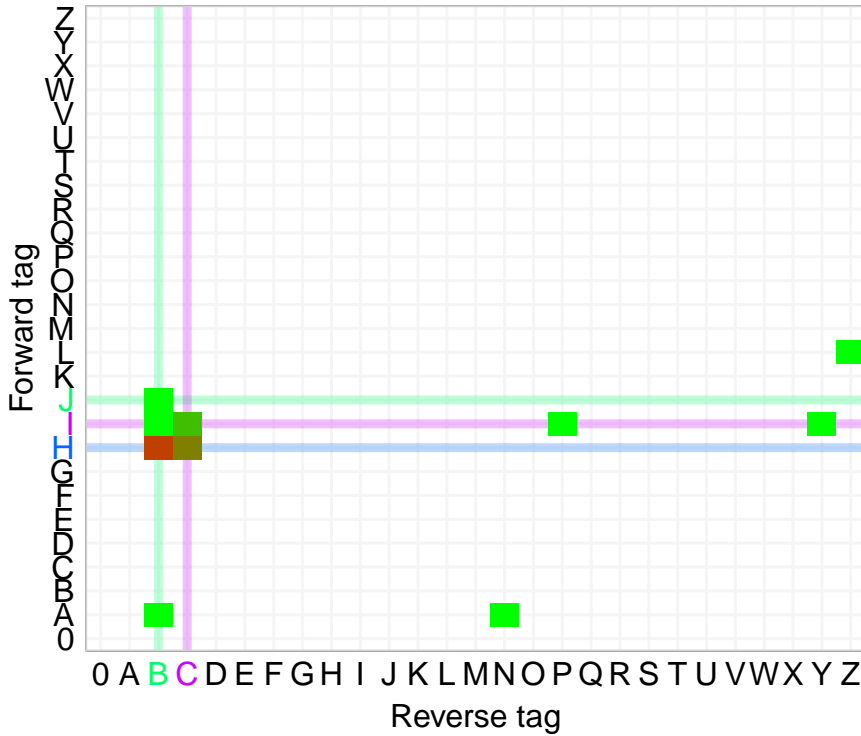**Perfect match (0 difference ISU)**

SFA−125
mock: hhmll

foram62 (m)

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 19]
]19 − 37]
]37 − 71]
]71 − 137]
]137 − 263]
]263 − 507]
]507 − 975]
]975 − 1876]
]1876 − 3608]

Forward tag
Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

**Perfect match (0 difference ISU)**
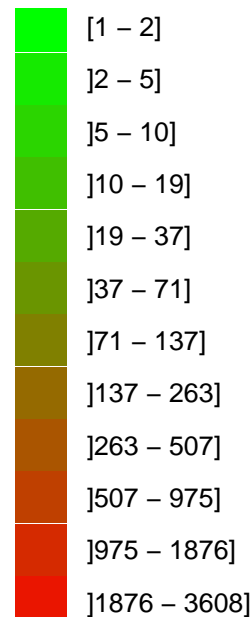
SFA−125

mock: hhmll

foram64  (I)

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 13]
]13 − 16]
]16 − 20]
]20 − 25]
]25 − 32]
]32 − 40]
]40 − 50]
]50 − 63]
]63 − 80]

Forward tag

Reverse tag

**1 difference ISUs**

○ non−critical mistag  ● correctly labelled

Number of reads per sample

# Perfect match (0 difference ISU)

**SFA–125**
**mock: Hhml**

foram52 (h)

[1 – 2]
]2 – 5]
]5 – 10]
]10 – 21]
]21 – 44]
]44 – 92]
]92 – 193]
]193 – 404]
]404 – 846]
]846 – 1772]
]1772 – 3714]
]3714 – 7781]

Forward tag

Reverse tag

## 1 difference ISUs

○ non–critical mistag  ● correctly labelled

Number of reads per sample

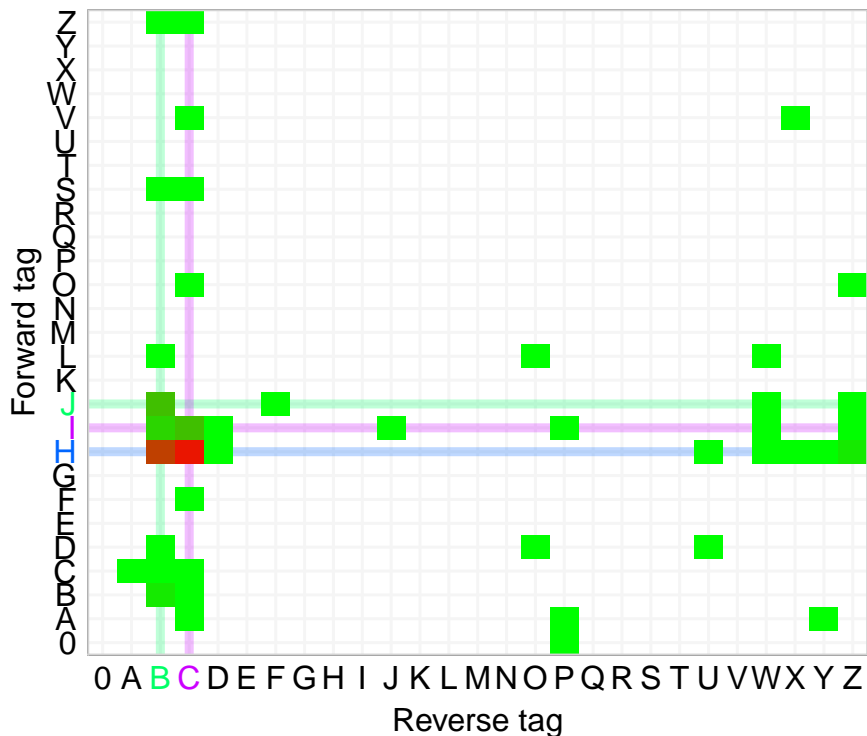**Perfect match (0 difference ISU)**

SFA−125
mock: Hhml

foram65 (m)

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 16]
]16 − 24]
]24 − 38]
]38 − 60]
]60 − 93]
]93 − 145]
]145 − 227]
]227 − 354]
]354 − 553]

Forward tag

Reverse tag

**1 difference ISUs**

non−critical mistag · correctly labelled

Number of reads per sample

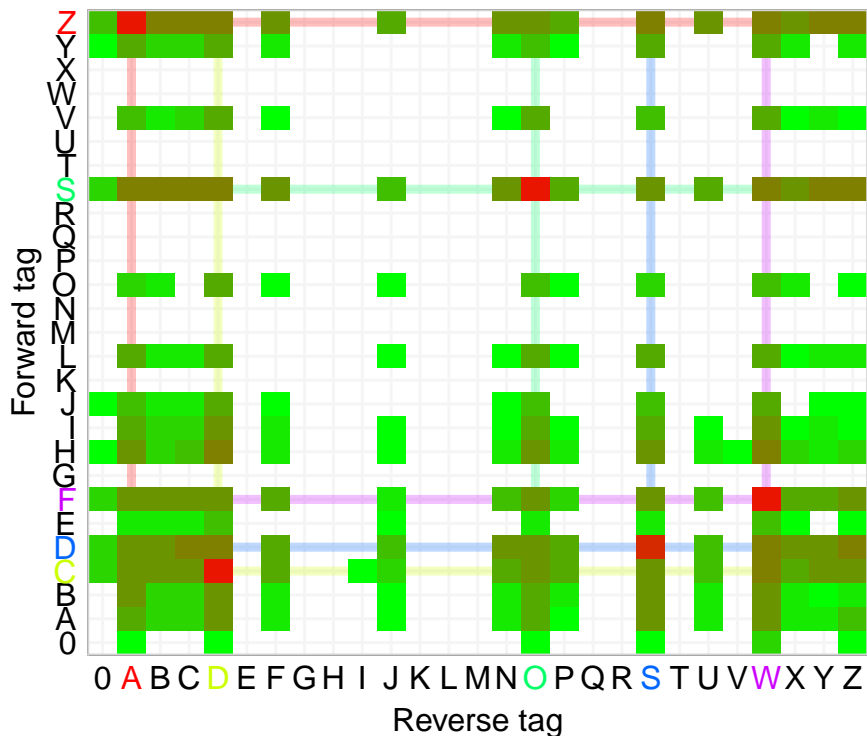# Perfect match (0 difference ISU)

SFA−125
mock: Hhml

foram46  (H)

Forward tag / Reverse tag

| | [1 – 2] |
| | ]2 – 5] |
| | ]5 – 10] |
| | ]10 – 24] |
| | ]24 – 57] |
| | ]57 – 135] |
| | ]135 – 320] |
| | ]320 – 762] |
| | ]762 – 1813] |
| | ]1813 – 4312] |
| | ]4312 – 10258] |
| | ]10258 – 24403] |

## 1 difference ISUs

○ non−critical mistag   ● correctly labelled

Number of reads per sample
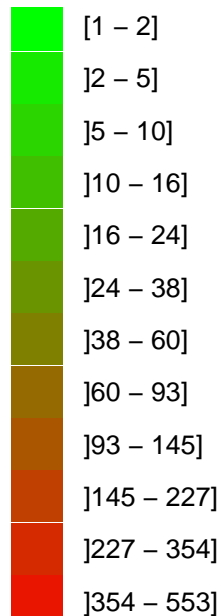
**Perfect match (0 difference ISU)**

SFA−125
mock: Hhml

foram79 (I)

[1 − 2]
]2 − 5]
]5 − 10]
]10 − 12]
]12 − 15]
]15 − 18]
]18 − 22]
]22 − 27]
]27 − 33]
]33 − 40]
]40 − 48]
]48 − 59]

Forward tag

Reverse tag
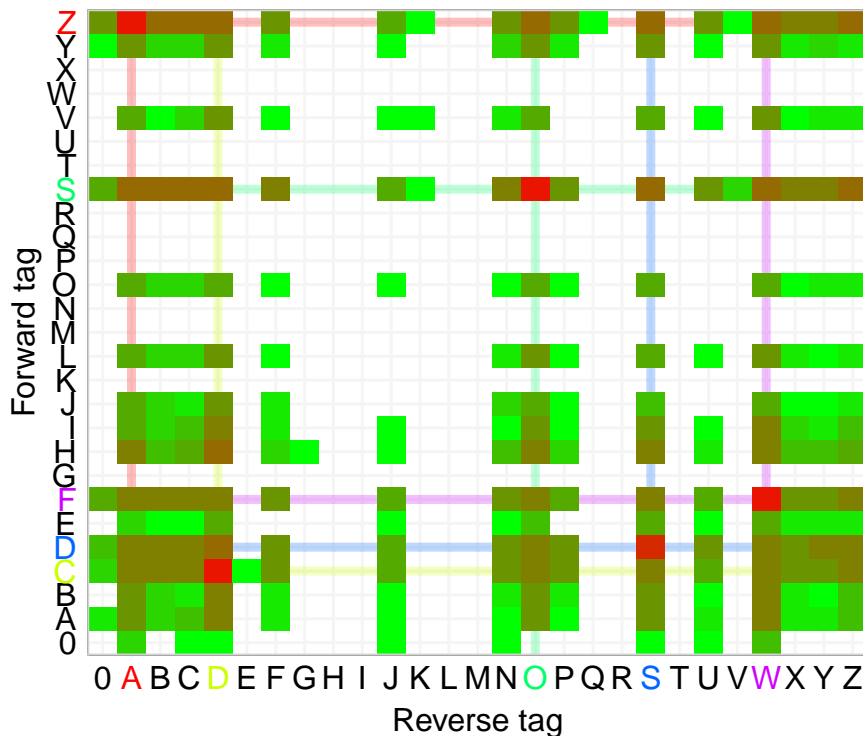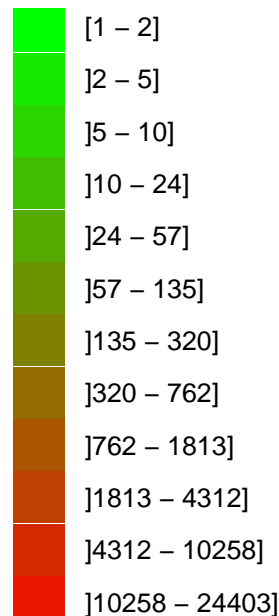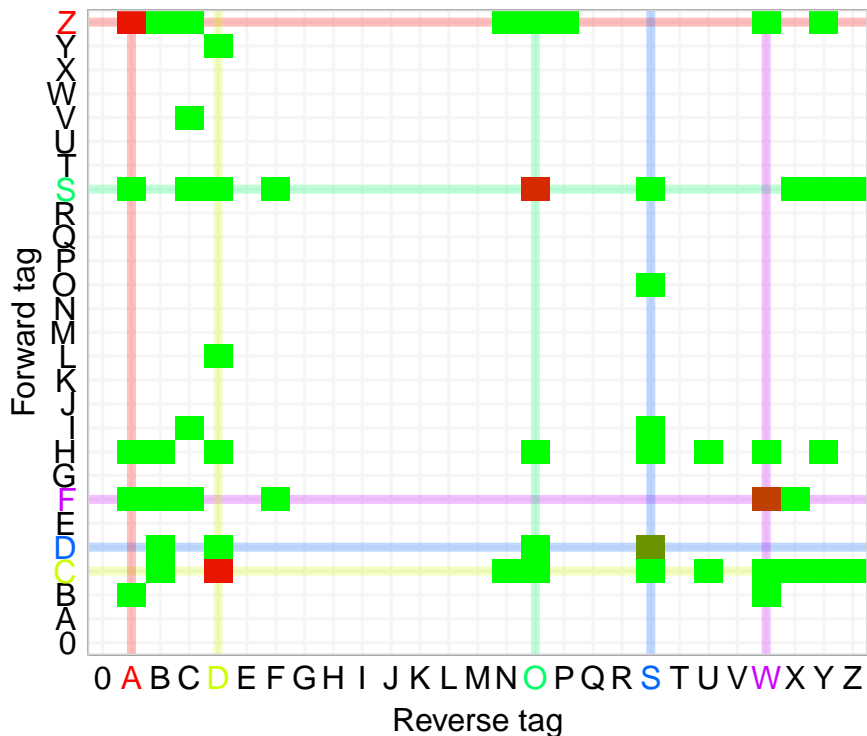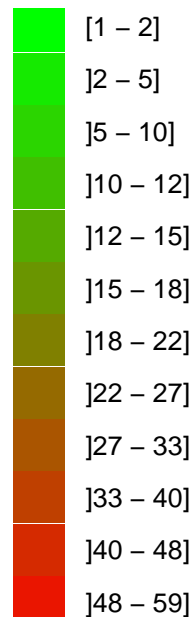
**1 difference ISUs**

○ non−critical mistag   ● correctly labelled

Number of reads per sample

**Supplementary Figure 10**. Distribution of ISUs in each possible intersection of PCR replicate samples. The ISU distributions are indicated by two 5-way Venn diagrams for each mock community: *hllll* (a, b), *hhhhl* (c, d), *hhmll* (e, f), *Hhml* (g, h), *even* (i, j) and *random* (k, l). In each intersection area of the "Inclusive" diagrams (left side), the numbers correspond to the total numbers of ISUs that would be found, including those that would also be found using more replicates. In the "Exclusive" diagrams (right side), an ISU already counted in a given intersection area is not re-counted in the intersection areas involving the same replicates.

Inclusive   Exclusive

hllll   hhhhl   hhmll   Hhml   even   Random

a b c d e f g h i j k l

**Supplementary Figure 11**. Distributions of ISUs and reads per category of number of replicates and per mock community. Each series of violin plots represents the data collected for each mock community: *hllll* (a), *hhhhl* (b), *hhmll* 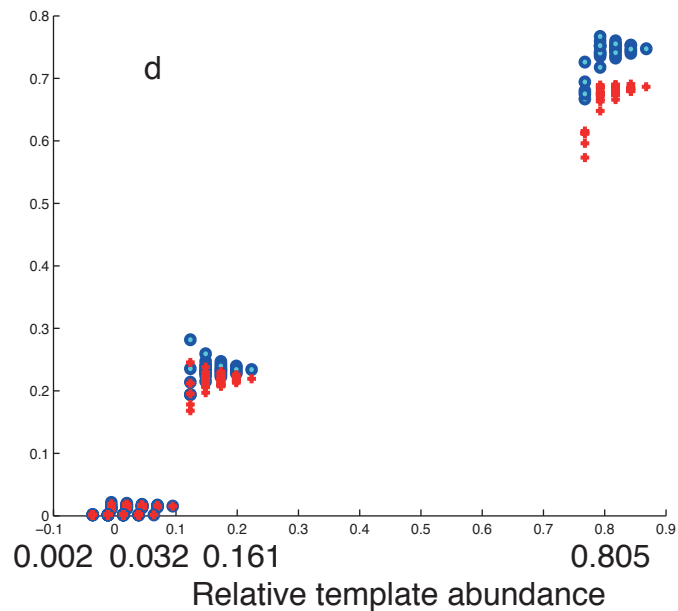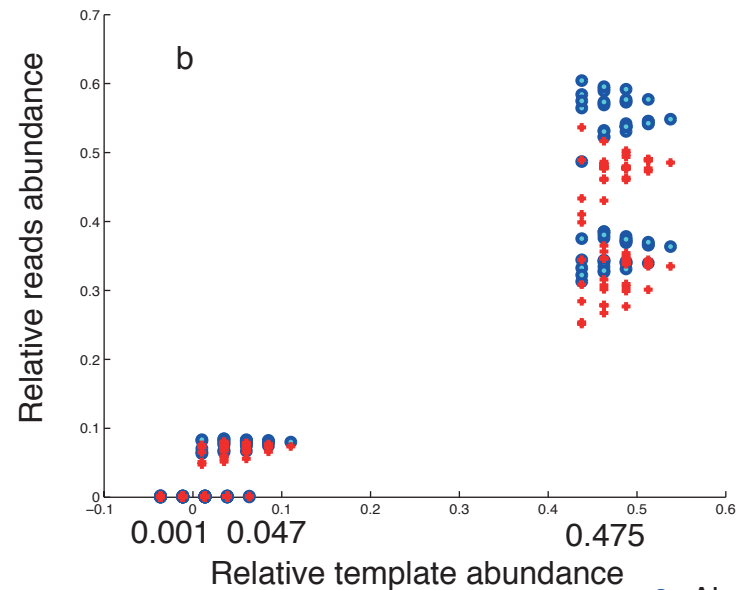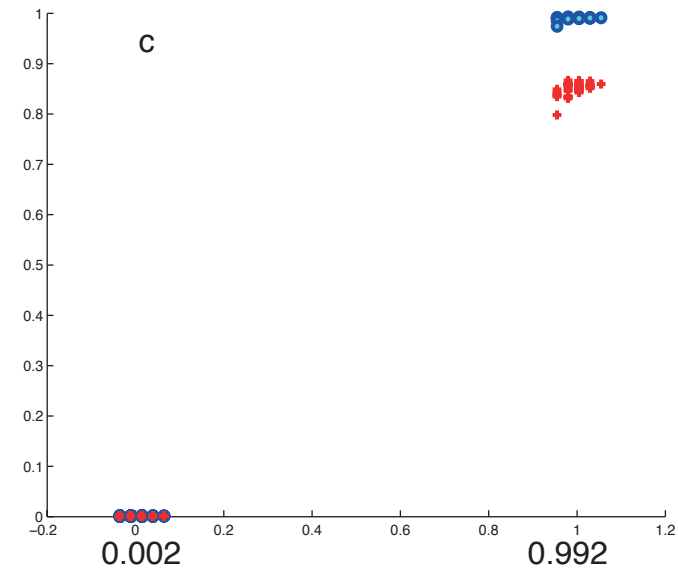(c), *Hhml* (d), *random* (e) and *even* (f). The series are split per color according to the number of replicates category (bottom colored box legends). For each category are represented both the density distribution of the number of ISUs per clone (left violin) and of the number of reads per ISU (right violin). The $\log_{10}$-transformed y-axis for the number of ISUs and for the number of reads are situated on the left (plain line) and right (dotted line) sides of each plot, respectively, on the Each violin separates the expected (left side) from the mistagging (right side) data. The median of each distribution is indicated (horizontal bars).

**Supplementary Figure 12**. Abundances comparisons between mock community sequence templates and the resulting ISUs. Abundances are displayed for each clone, but separately for each of the 4 mock communities of SFA-125, including (a) *hhhhl* containing 4 clones at high abundance and 1 clone at low abundance, (b) *hhmll* containing 2 clones at high abundance, 1 clone at medium abundance and 2 clones at low abundance, (c) *hllll* containing 1 clone at high abundance and 4 clones at low abundance and (d) *Hhml* containing 1 clone at very high abundance, 1 clone at high abundance, 1 clone at medium abundance and 1 clone at low abundance (see Supplementary Table 2). The clones (blue dots) and ISUs (red crosses) are organized in five columns according to the number of replicates intersection where it is found simultaneously. In each case, the exact value of the template abundance is located at the "3 replicates" position.

Number of replicates
1 2 3 4 5

Relative reads abundance (y-axis, panels a–d)

Relative template abundance (x-axis, panels a–d)

○ Abundances sum
+ Most abundant ISU