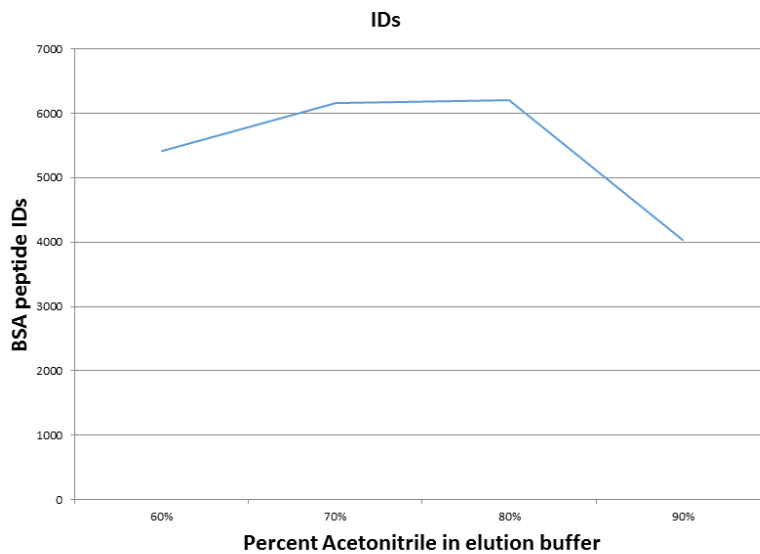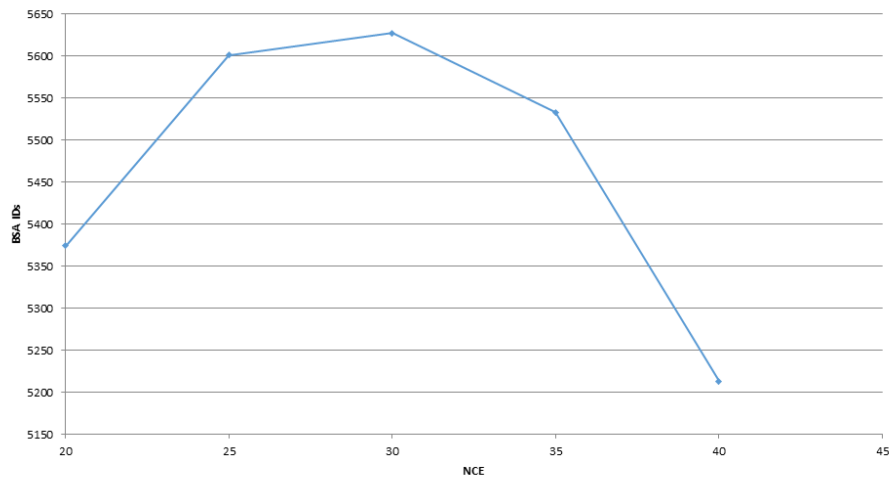**Supplementary Figure S1**- Optimization of microwave assisted acid hydrolysis time. A stock solution of 0.25 µg\µl reduced and alkylated BSA in 3 M HCl were prepared and divided into 15 tubes, each containing 10 µg BSA (40 µl). Tubes were microwaved together and at 2, 4, 8, 12 and 20 min, microwaving was briefly paused. Three tubes were transferred at that time point and placed on ice. Once all tubes were collected, they were subjected to identical solid phase extraction, LC-MS/MS analysis and database searching against the BSA sequence and common lab contaminants. The average number of BSA peptides identified in the database search of the 3 technical replicates was plotted against the microwaving time.

IDs

**Supplementary Figure S2**- Optimization of solid phase extraction. A stock solution of BSA peptides produced by 4 min MAAH was divided into 5 tubes. Each sample was subject to solid phase extraction by loading onto Oasis HLB sorbent in 96-well extraction plate, according to the manufacturer's instructions. Elution was performed with 60%, 70%, 80%, or 90% acetonitrile with 0.1% formic acid. Resulting elution was subjected to LC-MS/MS analysis and database searching against the BSA sequence and common lab contaminants. The number of BSA peptides identified in the database search was plotted against the acetonitrile percentage in the elution buffer.

**Supplementary Figure S3**- Optimization of normalized collision energy (NCE). A sample of BSA peptides produced by 4 min MAAH and solid phase extraction (elution with 60% acetonitrile) was analyzed by LC-MS/MS using 20, 25, 30, 35 or 40 NCE, with separate injections for each analysis. Resulting spectra were subjected to database searching against the BSA sequence and common lab contaminants. The number of BSA peptides identified in the database search is plotted against the NCE.

# pTA_BSA_Rep1_DiPS

**Parameters**

acid input files:

H:\Proteomics\Projects\RnD\Hydrolysis\manuscript\Assembly
\BSA\BSA_Rep1_DiPS_de-Novo_Peaks-output.csv
trypsin input files:

kmer size: 7
kmer min overlap: 5
unite min overlap: 5
unite min extension: 7
merge min quality: 0.7

**Final Contig1**

**Sequence with potential ambiguities/variants**

■ - Replacements of N->D and Q->E [acid]

■ - Potential variants

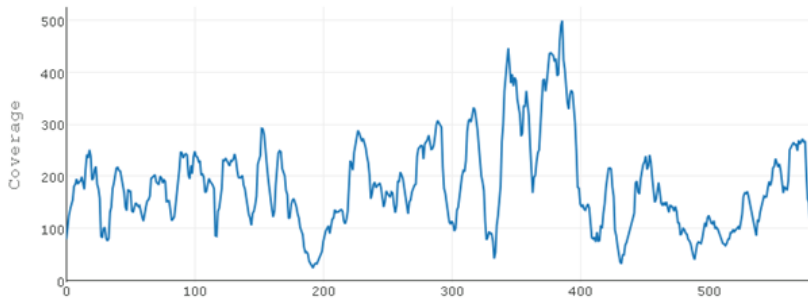E* - Pyroglutamate (Post-translational modification of N-terminal Q or E )

DTHKSELAHRFK[D:198,N:1]LGEEHFKGLVLLAFS[E:81,Q:1]YLE[E:75,Q:1
]CPFDEHVKLV[D:168,N:28]ELTEFAKTCVADESHAGCEKSLHTLFGD[E:183,Q:
1]LCKVASLR[E:133,Q:2]TYGDMADCCEK[E:240,Q:3][E:240,Q:1]P[E:18
8,Q:7]R[D:171,N:34][E:230,Q:8]CFLSHK[D:201,N:4][D:195,N:3]SP
DLPKLK[D:86,P:76][P:85,D:63]P[D:113,N:26]TLCDEFKAD[E:226,Q:3
]KKFWGKYLYELARRHPYFYAP[E:239,Q:1]LLYYA[D:160,N:38]KY[D:121,N
:16]GVF[E:212,Q:2][E:244,Q:1]CC[E:209,Q:6]AEDKGACLLPKLETMR[E
:62,Q:2]KVLTSSAR[E:32,Q:1]RLRCASL[E:95,Q:1]KFG[E:101,Q:2]RAL
KAWSVARLS[E:130,Q:3]KFPKA[E:261,Q:2]FV[E:281,Q:2]VTKLVTDLTKV
HKECCHGDLLECA[D:168,N:2][D:165,N:3]RA[D:139,N:4]LAKYLCD[D:11
9,N:18][E:127,Q:1]DTLSSKLKECC[D:232,N:1]KPLLEKSHCLAEVEKDALPE
[D:102,N:6]LPPLTADFA[E:209,Q:2][D:208,N:1]K[D:285,N:1]VCK[D:
252,N:71]Y[E:325,Q:2]EAKDAFLGSFLY[E:87,Q:3]YSRRHP[E:180,Q:1]
YAVSVLLRLAKEYEATL[E:334,Q:2]ECCAK[D:216,N:1][D:167,N:1]PHACY
STVF[D:362,N:1]KLKHLV[D:418,N:2][E:425,Q:1]P[E:389,Q:7][D:37
6,N:84]LLK[E:398,Q:7][D:316,N:58]C[D:325,N:4]EFEKLG[E:177,Q:
2]YGF[E:144,Q:2][D:117,N:20]ALLVRYTRKVP[E:76,Q:1]VSTPTLVEVSR
SLGKVGTRCCTKPESERMPCTE[D:222,N:1]YLSLLL[D:165,N:33]RLCVLHEKT
PVSEKVTKCCT[E:109,Q:1]SLV[D:81,N:17]RRPCFSALTPDETYVPKAFDEKLF
TFHADLCTLPDTEKELKKETALVELLKHKPKATEE[E:140,Q:1]LKTVM[E:188,Q:
2][D:152,N:30]FVAFVDKCCAA[D:176,N:1][D:172,N:1]KEACFAVEGPKLV
VST[E:145,Q:4]TALA

**Final consensus sequence**

Position: [ 0 ]    Weak ▬▬▬▬▬▬▬▬▬▬ Strong

```
  0 DTHKSELAHRFKDLGEEHFKGLVLLAFSEYLEECPFDEHVKLVNELTEFAKTCVADESHA 59
 60 GCEKSLHTLFGDELCKVASLRETYGDMADCCEKEEPQRNQCFLSHKDDSPDLPKLKDPPN 119
120 TLCDEFKADEKKFWGKYLYELARRHPYFYAPELLYYANKYNGVFEECCQAEDKGACLLPK 179
180 LETMRQKVLTSSARQRLRCASLEKFGERALKAWSVARLSEKFPKAEFVEVTKLVTDLTKV 239
240 HKECCHGDLLECADDRADLAKYLCDNEDTLSSKLKECCDKPLLEKSHCLAEVEKDALPED 299
300 LPPLTADFAEDKDVCKNYEEAKDAFLGSFLYQYSRRHPEYAVSVLLRLAKEYEATLEECC 359
360 AKDDPHACYSTVFDKLKHLVDEPENLLKENCDEFEKLGEYGFENALLVRYTRKVPEVSTP 419
420 TLVEVSRSLGKVGTRCCTKPESERMPCTEDYLSLLLNRLCVLHEKTPVSEKVTKCCTESL 479
480 VNRRPCFSALTPDETYVPKAFDEKLFTFHADLCTLPDTEKELKKETALVELLKHKPKATE 539
540 EELKTVMENFVAFVDKCCAADDKEACFAVEGPKLVVSTQTALA                  582
```

Coverage Graph

**Supplementary Figure S4**- Example of pTA final report. pTA reports the final assembled contigs. For each reported contig, there are 3 panels displayed: Top panel- "Sequence with potential ambiguities/variants" reports the assembled sequence including both optional residues and their coverage at positions where there is evidence supporting different residues. Middle panel- "Final consensus sequence" reports the final sequence decided upon where the selected residue in case of non-conclusive assignments is color coded for confidence in assignment. If proteolytic digestion data is also used, in the top panel the ratios are reported separately for MAAH and proteolysis in order to separate MAAH-derived chemical modifications from other causes of sequence variation. Bottom panel- "Coverage Graph" displays the coverage (number of peptide tags covering the position) at each position. In this example (BSA replicate #1), the final sequence decided upon contained a "swap" of two residues ('DP' instead of 'PD') at positions 117-118 as evident from alignment to the known BSA sequence (see Figure 3 in the main text of the manuscript). These positions are colored red (low confidence) in the "Final consensus sequence" panel, and show about equal amount of evidence for both options ("PD" or "DP") when examined in the top panel.

b)  Bovine serum albumin (583 AA long, excluding signal peptide)

| Replicate | Coverage | Accuracy | # contigs | Incorrect Q/E assignment without trypsin data | Incorrect N/D assignment without trypsin data | Incorrect Q/E assignment with trypsin data | Incorrect N/D assignment with trypsin data |
|---|---|---|---|---|---|---|---|
| 1 | 100% | 100% | 1 | 21 | 1 | 3 | 0 |
| 2 | 100% | 100% | 1 | 18 | 2 | 5 | 0 |
| 3 | 100% | 100% | 1 | 17 | 0 | 4 | 0 |

**Supplementary Figure S5** - DiPS results of BSA using MAAH and trypsin digestion data. BSA was cleaved using MAAH in triplicates and analyzed by nanoLC-MS/MS. A single BSA tryptic digest was also analyzed by nanoLC-MS/MS. Tryptic peptide tags were used as input for pTA in addition to each of the MAAH replicates. pTA's output sequences were aligned to the known sequence of each protein using Clustal-Omega (http://www.ebi.ac.uk/Tools/msa/clustalo/). Alignment mismatches that are the result of isobaric I/L, deamidated-Q/E, or deamidated-N/D ambigouities were not counted as sequencing mistakes in the accuracy calculations. Yellow – Q/E mismatches N/D mismatches. Grey – signal peptide or N terminal methionine.

**Supplementary Figure S6**- DiPS results of AR37 antibody light chain. AR37 was subjected to DiPS using two LC-MS/MS experiments of a single MAAH preparation and one experiment of a tryptic digest. Parameters for *de novo* analysis of MS/MS spectra included glutamine and glutamate conversion to pyroglutamate as variable modifications. (i) The resulting pTA assembled contig 1 was aligned to G0YP42 (its homolog identified by BLAST search). Domains (Frames 1-4, CDRs 1-3, and constant region) and tryptic peptides identified by a database search of the tryptic digest against DiPS determined sequence are marked on the sequence. (ii) pTA output web report for contig 1. Pyroglutamate is denoted by E* in the top panel.

**Supplementary Figure S7**- DiPS results of AR37 antibody heavy chain. (i) The resulting pTA assembled contigs 2 and 3 were aligned to I6L985 (the homolog protein to contig 2 identified by BLAST search). Domains (Frames 1-4 and CDRs 1-3) and tryptic peptides identified by a database search of the tryptic digest against DiPS determined sequence are marked on the sequence. (ii) pTA output web report for contig 3. Pyroglutamate is denoted by E* in the top panel.