

1 **Proteogenomic analysis of granulocyte macrophage colony- stimulating**  
2 **factor autoantibodies in the blood of a patient with autoimmune**  
3 **pulmonary alveolar proteinosis**

4

5 Atsushi Hashimoto<sup>1†</sup>, Shiho Takeuchi<sup>2†</sup>, Ryo Kajita<sup>3</sup>, Akira Yamagata<sup>4</sup>, Ryota  
6 Kakui<sup>4</sup>, Takahiro Tanaka<sup>1</sup>, and Koh Nakata<sup>1\*</sup>

7 <sup>1</sup> Niigata University Medical & Dental Hospital, Niigata, Japan

8 <sup>2</sup> Division of Cancer Genome Informatics, Niigata University Graduate School of  
9 Medical and Dental Sciences, Niigata, Japan

10 <sup>3</sup> Bruker Japan K.K.

11 <sup>4</sup> Propoenix Division, Food and Life-Science Laboratory, Idea Consultants,  
12 Inc., Osaka-city, Osaka, Japan

13 † These authors contributed equally.

14

15 Corresponding author: Koh Nakata, MD, PhD

16 Director,

17 Clinical and Translational Research Center,

18 Niigata University Medical and Dental Hospital, 1-754 Asahimachi-dori, Chuo-  
19 ku, Niigata 951-8520, Japan

20 Phone: +81-25-227-2029

21 FAX: +81-25-227-0377

22 E-mail address: [radical@med.niigata-u.ac.jp](mailto:radical@med.niigata-u.ac.jp)

23 ORCID: 0000-0002-4980-3467

24

25     **Abstract**

26     Recently, attempts to reveal the structures of autoantibodies  
27     comprehensively using improved proteogenomics technology, have become  
28     popular. This technology identifies peptides in highly purified antibodies by  
29     using an Orbitrap device to compare spectra from liquid chromatography–  
30     tandem mass spectrometry against a cDNA database obtained through next-  
31     generation sequencing. In this study, we first analyzed granulocyte-  
32     macrophage colony-stimulating factor (GM-CSF) autoantibodies in a patient  
33     with autoimmune pulmonary alveolar proteinosis, using the trapped ion  
34     mobility spectrometry coupled with quadrupole time-of-flight (timsTOF) Pro  
35     mass spectrometry. The timsTOF Pro identified peptides that partially  
36     matched sequences in up to 156 out of 162 cDNA clones. Complementarity-  
37     determining region 3 (CDR3) was fully and partially detected in nine and 132  
38     clones, respectively. Moreover, we confirmed one unique framework region 4  
39     (FR4) and at least three unique across CDR3 to FR4 peptides via *de novo*  
40     peptide sequencing. This new technology may thus permit the  
41     comprehensive identification of autoantibody structure. (150 words)

42

43

## 44 **Introduction**

45 Pulmonary alveolar proteinosis (PAP) is a rare lung disorder  
46 characterized by the accumulation of surfactant material in the alveoli and  
47 terminal bronchioli<sup>1,2,3</sup>.

48 PAP is classified into three different forms: autoimmune, congenital, and  
49 secondary<sup>1,2,3</sup>. Autoimmune PAP (aPAP) accounts for 91% of all PAP cases<sup>1,2,3</sup>.  
50 It is thought to be caused by granulocyte-macrophage colony-stimulating factor  
51 (GM-CSF) autoantibodies (GMAbs), which interfere with alveolar macrophage  
52 function and maturation, resulting in impaired surfactant catabolism. However,  
53 the mechanism involved in the excessive production of GMAbs remains  
54 unclear.

55 According to a previous study, GMAb is a polyclonal antibody with  
56 immunoglobulin M (IgM-), IgG-, and IgA-isotypes<sup>4</sup>. IgG-GMAb is pathogenic, but  
57 IgM-GMAb is thought to be a bystander etiologically, because its binding avidity  
58 to GM-CSF is 100-fold lower than that of IgG-GMAb, and its neutralizing  
59 capacity is extremely weak, i.e., it has a 20,000-fold higher IC<sub>50</sub> than IgG-  
60 GMAb<sup>4</sup>. The currently available evidence regarding the effects of IgA-GMAb in  
61 the peripheral blood is limited, and its pathogenic role remains unclear.

62 Serum GMAbs are normally present in healthy individuals, although  
63 their concentration is much lower than that measured in patients with aPAP<sup>5</sup>.

64 Our previous study have reported that circulating GM-CSF autoreactive B cells  
65 (GMARBs) producible of GMAb by EBV transformation are consistently  
66 detected in healthy individuals and patients with aPAP<sup>6</sup>. These are potential  
67 precursor cells for GMAb-producing plasma cells and B cells.

68 Determining the relationship between the structures and functions of  
69 GMAbs is important for deepening our understanding of the pathogenesis and  
70 etiology of aPAP. To date, the structure of GMAbs has been investigated by  
71 analyzing monoclonal antibodies of clones isolated from patients' B cells, which  
72 have been transformed using the Epstein–Barr virus<sup>7</sup>, or generated via single-  
73 cell polymerase chain reaction (PCR). Circulating GMARBs harboring surface  
74 GMAb for B-cell receptors were detected in peripheral blood samples from both  
75 healthy individuals and patients with aPAP<sup>4</sup>. Piccoli *et al.* (2015) established  
76 three non-cross-competing monoclonal GMAbs bound simultaneously to a  
77 single molecule of GM-CSF to form a high-molecular-weight immune complex,  
78 resulting in efficient *in vitro* neutralization of GM-CSF biological activity and  
79 promoted the rapid degradation of GM-CSF-containing immune complexes *in*  
80 *vivo* in an Fc-dependent manner<sup>8</sup>. Wang *et al.* generated 19 monoclonal  
81 GMAbs from six patients, these clones used multiple V genes, excluding  
82 preferred-V-gene use as an etiology, and targeted at least four non-overlapping  
83 epitopes on GM-CSF<sup>9</sup>. This suggested that it is GM-CSF that drives the

84 production of autoantibodies, and not a B-cell epitope in a pathogen cross-  
85 reaction with GM-CSF<sup>9</sup>.

86           However, none of these studies could overcome the issue of selection  
87 bias. Creating a bias-free cDNA database of antibody sequences from sorted  
88 GMARBs using next-generation sequencers developed after 2007 would solve  
89 this problem. However, not all sorted cells are necessarily differentiated into  
90 plasma cells; thus, identifying the full range of antibodies in the sera of patients  
91 with aPAP using these technologies is challenging.

92           Proteogenomics technology identifies peptides produced from  
93 enzymatically digested proteins by comparing tandem mass spectrometry  
94 (MS/MS) spectra against a cDNA database that has been derived from next-  
95 generation sequencing<sup>10</sup>. This technology has been applied to the identification  
96 of sequences of monoclonal or polyclonal antibodies induced by antigen  
97 challenge in animals<sup>11</sup>. It identified 77% of the whole variable region sequence  
98 that was encoded in the cDNA database at the peptide/protein level for antigen-  
99 induced polyclonal antibodies<sup>11</sup>. Immune induction may render the identification  
100 of low-abundance regions, especially CDRs.

101           More recently, liquid chromatography coupled with MS/MS (LC/MS/MS)  
102 using an Orbitrap device (offering high peptide sensitivity and resolution)  
103 enabled studies to determine the sequence diversity of polyclonal

104 autoantibodies in certain diseases. One study searched a cDNA database  
105 generated by next-generation sequencing and *de novo* peptide sequencing for  
106 mass spectral interpretation independent of a reference database<sup>12</sup>. This search  
107 identified, in the sera of patients with celiac disease, 20–75 clonotypes per  
108 patient containing the complex peptide sequence diversity necessary for IgA  
109 that reacts with a gluten-derived peptide or the autoantigen transglutaminase  
110 <sup>2</sup><sup>12</sup>. Using phage display combined with LC/MS/MS on an Orbitrap device, Chen  
111 *et al.* reported that the anti-desmoglein-1 or -3 autoantibodies in pemphigus  
112 were oligoclonal antibodies with less than 10 clonotypes<sup>13</sup>.

113 In this study, a new instrument offering a type of LC/MS/MS<sup>14</sup> termed  
114 trapped ion mobility spectrometry (TIMS) coupled with quadrupole time-of-flight,  
115 the timsTOF Pro, was used to identify trypsin-digested GMAb-derived peptides.  
116 In addition to accurate determination of the *m/z* in the TOF analyzer, TIMS  
117 offers separation based on size to charge ( $\Omega/z$ ) providing additional separation  
118 power and increased peak capacity. This is achieved by means of the parallel  
119 accumulation serial fragmentation (PASEF) method, which realizes a very high  
120 sequencing speed without a decrease in sensitivity and allows the selection of  
121 100–350 precursors per second<sup>14,15,16</sup>.

122           The purpose of this study was to assess the analytical accuracy of  
123   timsTOF Pro regarding the polyclonality of GMAb, by referring or not referring to  
124   a cDNA database created through genetic analysis.  
125



126 **Results**

127

128 ***Creation of a GMARB cDNA database***

129 Peripheral blood mononuclear cells (PBMCs) were purified from the patient's  
130 blood and GMARBs were isolated and subjected to total RNA extraction (Fig.  
131 1). From the total RNA, full-length cDNA with the SMARTer-oligo sequence at  
132 5' end was subjected to PCR to generate amplicons for variable regions of  
133 the IgG heavy chain (IgG-HV). Through 2 × 300 bp paired-end sequencing  
134 using an Illumina Miseq, 491,966 reads were obtained for the first raw reads,  
135 which were preprocessed to 56,928 high-quality reads using PRINSEQ version  
136 0.20.4, followed by MiXCR Immune Repertoire Analyzer version 2.0.  
137 Subsequently, reads were assigned to the V (D) J germline segments of the Ig  
138 sequences for annotation using the international ImMunoGeneTics/HighV-  
139 QUERy and STandardization (IMGT/HighV-QUEST) tool<sup>17</sup>, resulting in 54,619  
140 functional reads. These were further clustered to 162 clones, using the package  
141 Change-O to group clonotypes according to the V and J alleles and the  
142 nucleotide Hamming distance. These sequences from GMARBs were used to  
143 construct the reference databases and were used to analyze the LC/MS/MS  
144 spectra in the subsequent analysis.

145

146 ***Analysis of highly purified circulating GMAb via high-resolution***

147 ***LC/MS/MS and quadrupole TOF***

148 The peptide mixture derived from the highly purified GMAb was subjected to  
149 nanoElute ultra-high-performance LC (Bruker Daltonics, Bremen, Germany)  
150 coupled online with a hybrid timsTOF Pro with CaptiveSpray (Bruker Daltonics)  
151 ion source. The PEAKS-filtered result matched 4,000 spectra with 827 peptide  
152 sequences (Fig. 2), which shared some sequence similarity with IgG-HV in up  
153 to 156 out of the 162 clones listed in the cDNA database (Fig. 3a). Overall,  
154 median of no less than 63.4% [50.4%, 79.3%] of the total amino acid sequences  
155 were covered. Coverage appeared relatively lower for the CDR1 and CDR2  
156 regions (20% [0%,87.5%] and 0% [0%, 100%], respectively), whereas it was  
157 high for framework region 1 (FR1) (76% [0%, 82%]) and FR4 (100% [100%,  
158 100%]). In particular, a peptide (amino acid position of around 39–43 amino  
159 acids) – a region of five amino acids flanked by two tryptic cleavage sites – was  
160 hardly covered in 140 of the 156 clones. Six clones that were not covered by  
161 the PEAKS-filtered result at all exhibited a similar structure in terms of the  
162 putative peptide length of FRs and CDRs by tryptic digestion and candidate  
163 position of tryptic cleavage sites (Fig. 3b).

164 The PEAKS-filtered result covered the full length of CDR3 in nine of the  
165 156 clones, partially covered CDR3 in 132 of the clones, and did not cover it at

166 all in 15, giving a median coverage of 37.6% ([20.2%, 54.5%]; Table 1, Fig. 3c).  
167 There was a weak correlation between the length of the CDR3 region and the  
168 rate of coverage by the PEAKS-filtered result ( $\rho = -0.36$ ,  $p < 0.001$ ); however,  
169 the number and position of putative tryptic cleavage site (K and R) did not  
170 affect coverage. In the nine clones in which the whole sequence was covered,  
171 the CDR3 region comprised  $< 10$  amino acids. The second half of the FR3 had  
172 a relatively high coverage rate and the peptide consistently ended at the first  
173 tryptic cleavage site, near the beginning of the CDR3 region. In contrast, the  
174 beginning of the peptides leading to the FR4 region appeared to not be affected  
175 by the occurrence of tryptic cleavage sites located inside the CDR3 region. The  
176 diversity of IgG-HV was also confirmed by 2 dimensional development  
177 combined with timsTOF followed by PEAKS search (See Supplementary  
178 Result).

179

### 180 ***De novo CDR3-FR4 analysis***

181 We conducted *de novo* peptide sequence estimation from raw LC/MS/MS  
182 spectra not listed in the cDNA database to determine whether the patient's  
183 serum contained any GMAb clones other than those coded in the cDNA  
184 database. A total of 24,405 peptides remained unidentified by the PEAKS  
185 database search algorithms and are thus putative *de novo* sequences. Of

186 these, 127 contained the “ASTK” sequence at the C terminal end, which is  
187 characteristic of the beginning of the constant region of the IgG isotype. The  
188 PEAKS software assigns a local confidence score for each amino acid in a *de*  
189 *novo* sequence. The local confidence score ranges from 0% to 99%, indicating  
190 the confidence of the algorithm that a particular amino acid is correctly  
191 sequenced. Low-scoring amino acid (i.e., amino acid with a local confidence  
192 score of < 80%) were trimmed, and short-length peptides (< 5 amino acids)  
193 were excluded. Seventy-five peptides were further selected as quality control,  
194 as described in the Methods section (see Methods, Fig. 4). Of these, 10  
195 matched the J gene sequence when we searched for homology with J gene  
196 reference sequence data (TBLASTIN 2.2.16). One and at least three clones  
197 were confirmed to have *de novo* sequence of FR4 and across CDR3 to FR4  
198 regions, respectively, from the cDNA database (Table 2).

199

200 **Discussion**

201           We used proteome analysis by the newly developed LC/MS/MS  
202 timsTOF Pro followed by a PEAKS search to identify 156 of the total 162 clones  
203 listed in the IgG-HV cDNA database that matched enzymatically digested  
204 GMAb-derived peptides. This was supported by the fact that matched peptides  
205 were even identified for the hypervariable CDR3 region in 133 clones. The  
206 results suggested that most cDNA in GMARBs should actually be expressed as  
207 GMAb in serum. The results of the *de novo* peptide sequencing also suggested  
208 that some clones other than those in the cDNA database may be present in the  
209 patient's serum.

210           Proteomic analysis of patient-derived autoantibodies often encounters  
211 difficulties. Regions other than CDRs have a high sequence homology among  
212 clones, so it is difficult to determine clone-specific sequences, whereas CDRs  
213 are highly variable and the number of peptides in each clone is insufficient for  
214 detection by LC/MS/MS. A high-performance mass analyzer with high sensitivity  
215 and resolution is thus required to overcome this limitation. To date, a number of  
216 studies involving proteomic analysis of antibody structure have used an Orbitrap  
217 device that traps ions in electrostatic fields and converts them to a mass  
218 spectrum using a Fourier transform of the frequency signal<sup>18</sup>. In contrast, this  
219 study used a novel instrument that combines LC/MS/MS with trapped ion

220 mobility spectrometry (the timsTOF Pro), which offers additional separation  
221 power and increased peak capacity over instruments that do not perform  
222 trapped ion mobility separation<sup>14,16</sup>. There is little information available  
223 regarding comparison in performance between Orbitrap devices and the  
224 timsTOF Pro. One study that has compared their performance is that of Singhal  
225 *et al.*, who performed a proteomic analysis of proteins extracted from paraffin-  
226 embedded sections. Using timsTOF Pro with PASEF, they observed increased  
227 peptide and protein depth, due at least in part to the sensitivity and speed of the  
228 instrument: five times peptide than those identified by Orbitrap device<sup>19</sup>.

229         To investigate the structure of GMAb clones comprehensively using  
230 proteogenomics technology, it is necessary to prepare a bias-free IgG-HV  
231 cDNA database from circulating GMARBs. For this purpose, we applied the 5'  
232 rapid amplification of cDNA ends method<sup>20,21</sup> to total RNA extracted from the  
233 patient's GMARBs to generate a library of the full-length cDNA in the cells. We  
234 also applied high-throughput paired-end sequencing to maximize the number of  
235 reads obtained. Using pre- and postprocessing procedures, we succeeded in  
236 narrowing a sequence of approximately 500,000 reads down to 162 clones of  
237 antibody genes, using Change-O software to distinguish variations due to real  
238 mutations from those induced by PCR errors<sup>22</sup>.

239           However, the cDNA database itself does not necessarily reflect the  
240 presence of whole antibodies, because the 162 clones may not always  
241 differentiate into antibody-producing cells. Therefore, we investigated the  
242 existence of peptides deduced from the cDNA database using timsTOF Pro  
243 followed by a PEAKS search. Most of the clones in the cDNA database  
244 reflected the variety of GMABs in the serum, except for six that did not match  
245 any peptide assessed by the proteomic analysis. These six clones exhibited  
246 similar structural characteristics (e.g., regional length and tryptic cleavage sites)  
247 to the other 156 clones. They were therefore derived from memory B cells that  
248 may have been on the way to differentiating into GMAB-producing plasma cells.  
249 In contrast, *de novo* peptide sequencing confirmed that four clones may have  
250 differentiated into plasma cells in the past and did not currently exist in  
251 GMARBs.

252           In the present study the coverage of the peptides identified by timsTOF  
253 Pro analysis differed markedly among the IgG-HV regions. The coverage of the  
254 FR4 region of 136 of the 162 clones was 100%, whereas the small peptide area  
255 in the FR2 between the two tryptic cleavage sites was not covered at all for 128  
256 of the clones. This difference in coverage can be attributed to the characteristics  
257 of peptide identification by LC/MS/MS. Firstly, the peptide lengths that are likely  
258 to be detected by LC/MS/MS are 6–16 amino acids<sup>23,24</sup>. Secondly, the quantity

259 of certain peptides is a critical factor for detection (i.e., whether low quantities  
260 can be detected or not depends on the sensitivity of LC/MS/MS). Thirdly, the  
261 presence of acidic or hydrophilic amino acids in the peptide interferes with its  
262 detection. Considering these characteristics, it makes sense that the median  
263 coverage was as high as 100% for the FR4 region, because it is highly  
264 homologous and contains few acidic and/or hydrophilic amino acids. Even  
265 though amino acid sequences are highly conserved, but when the cleaved  
266 peptide is excessively short, detection of peptides was difficult as demonstrated  
267 by the narrow area in the FR2 region.

268         Of the regions in IgG-HV, the CDR3 region is especially important  
269 because it is the most hypervariable region and is responsible for antigen  
270 specificity, as well as being the primary determinant of clonality. However,  
271 previous studies using Orbitrap nano-electrospray ionization tandem mass  
272 spectrometry demonstrated consistently low coverage of this region<sup>25</sup>. In  
273 contrast, the CDR3 region was at least partially covered in 132 of the clones in  
274 this study, with a median coverage of 37.6%. This is conclusive because high  
275 sequencing speed with PASEF and additional TIMS separation can improve  
276 CDR3 peptide identifications if  $m/z$  is similar. This suggests that proteomic  
277 analysis of autoantibodies using the timsTOF Pro has even more potential if



278 pretreatment methods, such as enzymatic digestion of autoantibodies, or  
279 combined analyses using several different enzymes are improved<sup>11</sup>.  
280 By *de novo* peptide sequencing using the PEAKS-filtered result, we confirmed  
281 one unique FR4 and at least three unique across CDR3 to FR4 sequences that  
282 were not listed in the cDNA database. It is important to note that these numbers  
283 are minimum estimates. Seventy-five sequences were candidates for *de novo*  
284 sequences. However, most of these had < 80% local confidence for any amino  
285 acid in the sequence, so the true number of *de novo* sequences was unclear.  
286 Regardless of the true number, it is rather important that some clones with *de*  
287 *novo* peptide sequences, not found encoded in the cDNA database, were  
288 found. This suggests that the GMARB population changes over time and that  
289 plasma cells may remain and continue to produce GMAb, even after the  
290 differentiation of the original memory B cells into plasma cells.

291           In conclusion, using the timsTOF Pro, trypsin-digested peptides of  
292 GMABs partially matched the sequences of 156 of 162 cDNA clones of  
293 GMARBs generated by genetic analysis, suggesting that most GMARBs  
294 differentiate into plasma cells to generate GMABs in patients with aPAP.  
295 Proteomic analysis of divergent polyclonal antibodies, especially the CDR3  
296 region, is difficult. However, the fact that timsTOF Pro identified 81 sequences

297 in the CDR3 region may be attributed to the high sensitivity of this system. We

298 expect that this study will promote proteomic analysis of autoantibodies.

299 (2,700 words)

300

301 **Methods**

302

303 ***Subject***

304 In 2004, an inherently healthy male aged 44 years became aware of shortness  
305 of breath on exertion. High-resolution computed tomography revealed the  
306 presence of an abnormal chest shadow. The patient was diagnosed with aPAP  
307 following bronchoalveolar lavage and serological testing at a local hospital. The  
308 level of GMAB in his serum was 19.48 µg/mL. The Institutional Ethics Review  
309 Committee of Niigata University approved the study (No. 2015-2388). Written  
310 informed consent was provided by the patient prior to enrollment in the study.

311

312 ***Isolation of GMARBs from patient blood***

313 PBMCs were isolated through density gradient centrifugation with Ficoll-Paque  
314 Plus (GE Healthcare, Chicago, IL, USA), and B cells were isolated using a Pan  
315 B cell isolation kit (Miltenyi Biotec, Bergisch Gladbach, Germany). The purity of  
316 isolated B cells was confirmed by flow cytometry. GMARBs were further purified  
317 by incubation with 500 ng/mL of biotinylated rhGM-CSF for 30 min on ice,  
318 followed by magnetic isolation of GMARBs using the MACS MS system.

319

320 ***Preparation of full-length cDNA extracted from GMARBs***

321 Total RNA was extracted from the isolated GMARBs using the RNeasy Micro kit  
322 (QIAGEN, Venlo, Netherlands) according to the instructions provided by the  
323 manufacturer. The RNA was reverse-transcribed into full-length cDNA using the  
324 SMARTer cDNA Synthesis Kit (Takara Bio, Shiga, Japan) to avoid amplification  
325 bias related to primer variability. Briefly, the modified oligo (dT) primer 3'  
326 SMART CDS Primer II A (Takara Bio) was used for the first-strand single-strand  
327 cDNA synthesis reaction. When SMART Scribe Reverse Transcriptase (Takara  
328 Bio) reaches the 5' end of the mRNA, the enzyme's terminal transferase activity  
329 adds a few nucleotides to the 3' end of the single-strand cDNA. The SMARTer  
330 Oligonucleotide base pairs with a non-template nucleotide stretch are added,  
331 creating an extended template. SMART Scribe Reverse Transcriptase switches  
332 the templates and continues replicating to the end of the oligonucleotide<sup>21,26,27</sup>.  
333 Double-strand cDNA was then synthesized using the Advantage 2 PCR  
334 Enzyme kit (Takara Bio).

335

336 ***Preparation of a library of IgGH variable regions for use on next-***  
337 ***generation sequencing platforms***

338 The first PCR amplification for variable regions of the H chain was performed  
339 using KOD-Plus-Neo DNA polymerase (Toyobo, Osaka, Japan) with a  
340 SMARTer primer for the common 5' primer and 3' primer. The amplicons were

341 purified and subjected to high-throughput 2 × 300 bp paired-end sequencing  
342 using the MiSeq v3 reagent kit (Illumina) with a 50% PhiX spike on the Illumina  
343 MiSeq platform, according to the recommendations provided by the  
344 manufacturer.

345

346 ***Pre-processing of raw reads obtained through high-throughput***

347 ***sequencing***

348 Quality control of raw reads obtained by high-throughput sequencing was  
349 performed as follows. For the H chain, the 3' low-quality bases were trimmed;  
350 total sequences with lengths of < 100 nucleotides were filtered out using paired-  
351 end reads, as were raw reads with a mean Phred quality score of < 30. These  
352 steps were conducted using the quality control and data preprocessing  
353 functions of PRINSEQ version 0.20.4<sup>28</sup>. The two paired-end reads were  
354 subsequently assembled into a complete Ig sequence, the forward and reverse  
355 primers were removed, and the data were de-multiplexed for each isotype using  
356 the MiXCR Immune Repertoire Analyzer version 2.0<sup>29</sup>. The software assigned  
357 the reads to regions while discriminating between IgM and IgG, using the  
358 reference database (GenBank). Variable region sequence reads with  
359 incomplete lengths, unassembled read sequences, or read sequences with  
360 missing primers were also removed.

361

362 ***V (D) J assignment and annotation***

363 After quality control, reads were assigned to the V (D) J germline segments of  
364 the Ig sequences for annotation using IMGT/HighV-QUEST software  
365 (<http://imgt.org>, v1.5.0)<sup>17</sup> and the IMGT gene database (version: 7 January  
366 2016).

367

368 ***Additional quality control and clonal clustering***

369 For post-processing of the IMGT/HighV-QUEST output, additional quality  
370 control and clonal clustering were performed using several software programs,  
371 Change-O and SHazaM, and custom scripts within the R statistical computing  
372 environment (version 3.3.3).

373

374 ***Autoantibody purification***

375 Serum (6 mL) from the same patient collected on the same day as described  
376 above was filtered with a 0.2 µm filter, diluted 8-fold with tris-buffered  
377 saline (TBS) (pH 7.4), and processed using a HiTrap rProtein A FF column (1  
378 mL of resin). After washing with TBS, the IgG proteins were eluted with 4 mL of  
379 100 mM glycine-hydrochloride (glycine-HCl) (pH 2.5), followed by neutralization

380 with 1 M Tris-HCl (pH 8.6), dialysis against TBS for 2 d at 4 °C, and application  
381 to a GM-CSF-coupled HiTrap NHS-activated HP column.

382 After washing with 10 mL of TBS, nonspecifically bound IgG was further  
383 washed out via the addition of 10 mL of 50 mM ammonium acetate (pH 5.0).

384 The IgG proteins binding to GM-CSF were eluted with 100 mM glycine-HCl (pH  
385 2.5). The eluate was neutralized by adding 1 M Tris-HCl (pH 8.6), concentrated  
386 with Ultrafree-0.5 (5 kDa molecular weight cutoff) at 15,000 × g for 30 min, and  
387 loaded into a 10% sodium dodecyl sulfate-polyacrylamide gel for  
388 electrophoresis (Fig. 2).

389

### 390 ***In-gel digestion and mass spectrometric identification of proteins***

391 Protein bands were excised from the electrophoretic gels and dehydrated in  
392 acetonitrile. The gel pieces were rehydrated in a digestion solution consisting of  
393 50 mM ammonium bicarbonate and 0.01 µg/µl modified sequence-grade trypsin  
394 (Promega, Madison, WI, USA)<sup>30,31</sup>. After overnight incubation at 37 °C, the  
395 digested peptides were extracted using 1% trifluoroacetic acid in 60%  
396 acetonitrile, and the extracted peptides were dried in a vacuum centrifuge<sup>32,33,34</sup>.  
397 The peptides were purified with ZipTip (Millipore, Burlington, MA, USA)  
398 according to the protocol provided by the manufacturer.

399

400 ***LC/MS/MS analysis***

401 The peptide mixture was loaded onto a C18 column (25 cm × 75 μm, 1.6 μm;  
402 IonOpticks, Melbourne, Australia). The mobile phases consisted of (A) 0.1%  
403 formic acid and (B) acetonitrile with 0.1% formic acid (volume per volume). The  
404 nano-flow LC gradient was delivered at 400 nL/min and consisted of a linear  
405 gradient of mobile phase B increasing from 2% to 17% B in 60 min, followed by  
406 increases to 25% B in 30 min, 37% in 10 min, and 95% B in 10 min. Ions were  
407 collected in the TIMS device over 100 ms and MS and MS/MS data were  
408 acquired over an *m/z* range of 100–1,700. During the collection of MS/MS data,  
409 the TIMS cycle was adjusted to 1.1 s and included 1 MS plus 10 PASEF-  
410 MS/MS scans, each containing on average 12 MS/MS spectra (>100 Hz)<sup>15,16</sup>.

411

412 ***cDNA database search and data processing***

413 All raw LC/MS/MS data were submitted to PEAKS Studio 8.5 (Bioinformatics  
414 Solutions, Waterloo, Canada) for data processing<sup>35</sup>. The cDNA database  
415 described above was used as a reference, and a decoy database was  
416 constructed for use in this analysis to control false discovery rates<sup>36</sup>. Search  
417 parameters included trypsin as the enzyme, with up to one missed cleavage  
418 allowed. Carbamidomethylation of cysteine residues was set as a fixed  
419 modification, and oxidation of methionine was set as a variable modification.



420 Parent mass error tolerance was set to 10 ppm, while fragment mass error  
421 tolerance was set to 0.05 Da. False discovery rates were adjusted to 1% at the  
422 peptide spectrum matches.

423

#### 424 ***De novo peptide sequencing***

425 *De novo* sequencing<sup>37</sup> was used to identify peptides harboring the CDR3–FR4  
426 regions, characterized as amino acid sequences ending in ASTK (IgG)<sup>12</sup>. The  
427 list of exclusively *de novo* peptides includes high-quality peptide sequences  
428 detected by *de novo* sequencing that remain unidentified by the PEAKS  
429 database search. The peptides containing ASTK sequences were extracted  
430 from this list. Low-quality amino acids (local confidence < 80%) were trimmed.  
431 Peptide sequences with lengths of less than five peptides were excluded. The  
432 novel proteins that required *de novo* sequencing and lacked similar proteins (for  
433 the CDR3–FR4 regions) in our cDNA databases were defined via the following  
434 two steps. Firstly, the datasets were searched against the J gene reference  
435 sequences obtained from the IMGT/LIGM-DB reference sequences database to  
436 identify the J gene. Next, a second search was performed against the in-house  
437 cDNA next-generation sequencing database assembled for the identified  
438 CDR3–FR4 region sequence.

439

440 ***Statistical analysis***

441 Non-normally distributed data are reported as medians with interquartile [25%,  
442 75%]. Spearman's correlation coefficient was used to estimate the relationship  
443 between two parameters. Statistical analyses were performed on a  
444 microcomputer using JMP (12.0.1) software (SAS Institute, Cary, NC, USA).  
445

446 **Data availability**

447 NGS sequenced data that support the findings of this study have been deposited  
448 in the DDBJ Japanese Genotype-phenotype Archive (JGA,  
449 <http://trace.ddbj.nig.ac.jp/jga>)<sup>38</sup> with the accession codes JGAS00000000163.  
450 The mass spectrometry data have been deposited to the ProteomeXchange  
451 Consortium (<http://proteomecentral.proteomexchange.org>) via the jPOSTrepo  
452 (<https://repository.jpostdb.org/>)<sup>39</sup> with the accession numbers PXD015420  
453 /JPST000674.

454

455 **Acknowledgements**

456 We thank Dr. Syujiro Okuda for providing instructions regarding the post-  
457 processing methods for data obtained from next-generation sequencing. This  
458 work was supported by the Japan Society for the Promotion of Science (JSPS)  
459 KAKENHI Grant Number JP20390230 (KN), JP24390208 (KN), and  
460 JP15H04829 (KN).

461

462 **Authorship Contributions**

463 KN, AH, and ST designed this study and wrote the manuscript. AH and ST  
464 mainly performed the experiment and analyzed the data. Ryo K, AY, Ryota K  
465 and TT assisted in the experiments.

466

467 **Conflicts of Interest Disclosure**

468 The authors declare no competing financial interests.

469

470 **References**

- 471 1. Seymour, J. F. & Presneill, J. J. Pulmonary alveolar proteinosis: Progress  
472 in the first 44 years. *American Journal of Respiratory and Critical Care*  
473 *Medicine* **166**, 215–235 (2002).
- 474 2. Trapnell, B. C., Whitsett, J. A. & Nakata, K. Pulmonary Alveolar  
475 Proteinosis. *N. Engl. J. Med.* **349**, 2527–2539 (2003).
- 476 3. Borie, R. *et al.* Pulmonary alveolar proteinosis. *Eur. Respir. Rev.* **20**, 98–  
477 107 (2011).
- 478 4. Nei, T. *et al.* IgM-type GM-CSF autoantibody is etiologically a bystander  
479 but associated with IgG-type autoantibody production in autoimmune  
480 pulmonary alveolar proteinosis. *Am. J. Physiol. Cell. Mol. Physiol.* **302**,  
481 L959–L964 (2012).
- 482 5. Uchida, K. *et al.* Granulocyte/macrophage-colony-stimulating factor  
483 autoantibodies and myeloid cell immune functions in healthy subjects.  
484 *Blood* **113**, 2547–2556 (2009).
- 485 6. Nei, T. *et al.* Memory B cell pool of autoimmune pulmonary alveolar  
486 proteinosis patients contains higher frequency of GM-CSF autoreactive B  
487 cells than healthy subjects. *Immunol. Lett.* **212**, 22–29 (2019).
- 488 7. Revoltella, R. P. *et al.* Antibodies binding granulocyte-macrophage colony  
489 stimulating factor produced by cord blood-derived B cell lines

- 490 immortalized by Epstein-Barr virus in vitro. *Cell. Immunol.* **204**, 114–127  
491 (2000).
- 492 8. Piccoli, L. *et al.* Neutralization and clearance of GM-CSF by  
493 autoantibodies in pulmonary alveolar proteinosis. *Nat. Commun.* **6**,  
494 (2015).
- 495 9. Wang, Y. *et al.* Characterization of pathogenic human monoclonal  
496 autoantibodies against GM-CSF. *Proc. Natl. Acad. Sci.* **110**, 7832–7837  
497 (2013).
- 498 10. Sheynkman, G. M., Shortreed, M. R., Cesnik, A. J. & Lloyd, M.  
499 Proteogenomics: Integrating Next-Generation Sequencing and Mass  
500 Spectrometry to Characterize Human Proteomic Variation. *Annu Rev Anal*  
501 *Chem (Palo Alto Calif)* **12**, 521–545 (2016).
- 502 11. Cheung, W. C. *et al.* A proteomics approach for the identification and  
503 cloning of monoclonal antibodies from serum. *Nat. Biotechnol.* **30**, 447–  
504 452 (2012).
- 505 12. Iversen, R. *et al.* Strong Clonal Relatedness between Serum and Gut IgA  
506 despite Different Plasma Cell Origins Article Strong Clonal Relatedness  
507 between Serum and Gut IgA despite Different Plasma Cell Origins. *Cell*  
508 *Rep.* 2357–2367 (2017). doi:10.1016/j.celrep.2017.08.036

- 509 13. Chen, J. *et al.* Proteomic Analysis of Pemphigus Autoantibodies Indicates  
510 a Larger , More Diverse , and More Dynamic Repertoire than Determined  
511 by B Cell Genetics Resource Proteomic Analysis of Pemphigus  
512 Autoantibodies Indicates a Larger , More Diverse , and More Dynamic .  
513 *CellReports* **18**, 237–247 (2017).
- 514 14. Vasilopoulou, C. G. *et al.* Trapped ion mobility spectrometry (TIMS) and  
515 parallel accumulation - serial fragmentation (PASEF) enable in-depth  
516 lipidomics from minimal sample amounts. *BioRxiv* (2019).
- 517 15. Meier, F. *et al.* Online Parallel Accumulation–Serial Fragmentation  
518 (PASEF) with a Novel Trapped Ion Mobility Mass Spectrometer. *Mol. Cell.*  
519 *Proteomics* **17**, 2534–2545 (2018).
- 520 16. Meier, F. *et al.* Parallel Accumulation – Serial Fragmentation (PASEF):  
521 Multiplying Sequencing Speed and Sensitivity by Synchronized Scans in  
522 a Trapped Ion Mobility Device. *J. Proteome Res* **14**, 5378–5387 (2015).
- 523 17. Alamyar, E., Duroux, P., Lefranc, M. P. & Giudicelli, V. IMGT® tools for  
524 the nucleotide analysis of immunoglobulin (IG) and t cell receptor (TR) V-  
525 (D)-J repertoires, polymorphisms, and IG mutations: IMGT/V-QUEST and  
526 IMGT/HighV-QUEST for NGS. *Methods Mol. Biol.* **882**, 569–604 (2012).

- 527 18. Downard, K. M., Morrissey, B. & Schwahn, A. B. Orbitrap mass  
528 spectrometry: instrumentation, ion motion and applications. *Mass*  
529 *Spectrosc. Rev.* **28**, 35–49 (2009).
- 530 19. Chris, M., Ryan, D. & Chien, A. S. Deep proteome mining of FFPE tissue  
531 with PASEF technology. *6th ASMS, San Diego, California, ThP785, June*  
532 *7* (2018).
- 533 20. Dunn-Walters, D., Townsend, C., Sinclair, E. & Stewart, A.  
534 Immunoglobulin gene analysis as a tool for investigating human immune  
535 responses. *Immunological Reviews* **284**, 132–147 (2018).
- 536 21. Matz, M. *et al.* Amplification of cDNA ends based on template-switching  
537 effect and step-out PCR. *Nucleic Acids Res.* **27**, 1558–1560 (1999).
- 538 22. Gupta, N. T. *et al.* Change-O: A toolkit for analyzing large-scale B cell  
539 immunoglobulin repertoire sequencing data. *Bioinformatics* **31**, 3356–  
540 3358 (2015).
- 541 23. Uchida, Y. *et al.* A study protocol for quantitative targeted absolute  
542 proteomics (QTAP) by LC-MS/MS: Application for inter-strain differences  
543 in protein expression levels of transporters, receptors, claudin-5, and  
544 marker proteins at the blood-brain barrier in ddY, FVB, and. *Fluids*  
545 *Barriers CNS* **10**, 1 (2013).

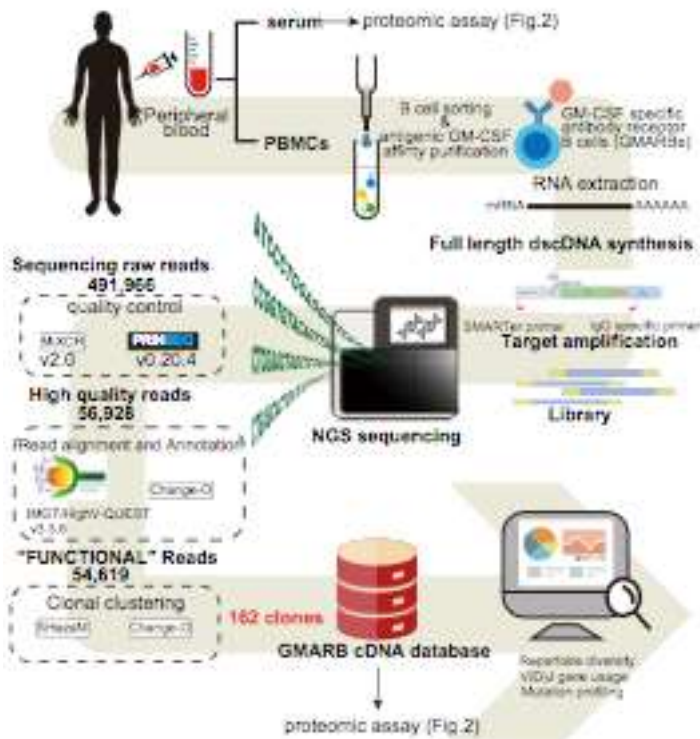


- 546 24. Kamiie, J. *et al.* Quantitative atlas of membrane transporter proteins:  
547 Development and application of a highly sensitive simultaneous  
548 LC/MS/MS method combined with novel in-silico peptide selection criteria.  
549 *Pharm. Res.* **25**, 1469–1483 (2008).
- 550 25. Boutz DR, Horton AP, Wine Y, Lavinder JJ, Georgiou G, M. E. Proteomic  
551 identification of monoclonal antibodies from serum. *Anal Chem.* **20**, 4758–  
552 66 (2014).
- 553 26. Zhu, Y. Y., Machleder, E. M., Chenchik, A., Li, R. & Siebert, P. D.  
554 Reverse transcriptase template switching: A SMART™ approach for full-  
555 length cDNA library construction. *Biotechniques* **30**, 892–897 (2001).
- 556 27. Cocquet, J., Chong, A., Zhang, G. & Veitia, R. A. Reverse transcriptase  
557 template switching and false alternative transcripts. *Genomics* **88**, 127–  
558 131 (2006).
- 559 28. Schmieder, R. & Edwards, R. Quality control and preprocessing of  
560 metagenomic datasets. *Bioinformatics* **27**, 863–864 (2011).
- 561 29. Bolotin, D. A. *et al.* MiXCR: Software for comprehensive adaptive  
562 immunity profiling. *Nat. Methods* **12**, 380–381 (2015).
- 563 30. Rabilloud, T., Valette, C. & Lawrence, J. J. Sample application by in - gel  
564 rehydration improves the resolution of two - dimensional electrophoresis

- 565 with immobilized pH gradients in the first dimension. *Electrophoresis* **15**,  
566 1552–1558 (1994).
- 567 31. Sanchez, J.-C. *et al.* Improved and simplified in-gel sample application  
568 using. *Electrophoresis* **18**, 324–327 (1997).
- 569 32. Gharahdaghi F , Weinberg CR , Meagher DA , Imai BS, M. S. Mass  
570 spectrometric identification of proteins from silver- stained polyacrylamide  
571 gel : a method for the removal of silve ... Mass spectrometric identification  
572 of proteins from silver-stain. *Electrophoresis* **20**, 1–6 (1999).
- 573 33. Kristensen, D. B. *et al.* Mass spectrometric approaches for the  
574 characterization of proteins on a hybrid quadrupole time-of-flight ( Q-  
575 TOF ) mass spectrometer Proteomics and 2-DE. *Electrophoresis* **21**,  
576 430–439 (2000).
- 577 34. Lopez, M. F. *et al.* A comparison of silver stain and SYPRO Ruby Protein  
578 Gel Stain with respect to protein detection in two-dimensional gels and  
579 identification by peptide mass profiling. *Electrophoresis* **21**, 3673–3683  
580 (2000).
- 581 35. Ma, B. *et al.* PEAKS : powerful software for peptide de novo sequencing  
582 by tandem mass spectrometry. *Rapid Commun Mass Spectrom* **17**,  
583 2337–2342 (2003).

- 584 36. Pevzner, N. G. and P. A. False Discovery Rates of Protein Identifications  
585 A Strike against the Two-Peptide Rule. *J Proteome Res* **8**, 4173–4181  
586 (2012).
- 587 37. Zhang, J. *et al.* PEAKS DB: De Novo Sequencing Assisted Database  
588 Search for Sensitive and Accurate Peptide Identification. *Mol. Cell.*  
589 *Proteomics* **11**, M111.010587 (2012).
- 590 38. Kodama, Y. *et al.* The DDBJ Japanese genotype-phenotype archive for  
591 genetic and phenotypic human data. *Nucleic Acids Res.* **43**, D18–D22  
592 (2015).
- 593 39. Okuda, S. *et al.* JPOSTrepo: An international standard data repository for  
594 proteomes. *Nucleic Acids Res.* **45**, D1107–D1111 (2017).  
595  
596

597



598

599 **Figure 1.**

600 Basic workflow for the generation of full-length cDNA of GMARBs, followed by

601 next-generation sequencing of immunoglobulin G heavy chain.

602 PBMCs and serum were simultaneously extracted from whole blood (20 mL)

603 obtained from a patient with aPAP. GMARBs were isolated, and the full-length

604 cDNA library was generated. This was followed by the synthesis of bias-free

605 PCR amplicons coding an IgG heavy chain variable region, and next-generation

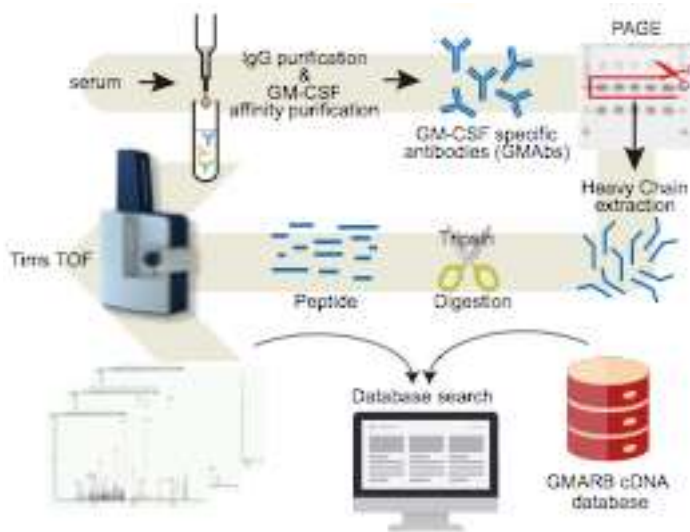
606 sequencing. The raw data were pre-processed by quality trimming using FASTX

607 and applied to read alignment by IMG2 High-V-QUEST, to identify the

608 functional reads. Numbers in the dotted squares indicate the number of

609 nucleotide sequence reads. Thereafter, classification of the functional reads into  
610 clones was performed using the clone clustering software Change-O. The  
611 methods used for clustering are described in the Methods.  
612

613



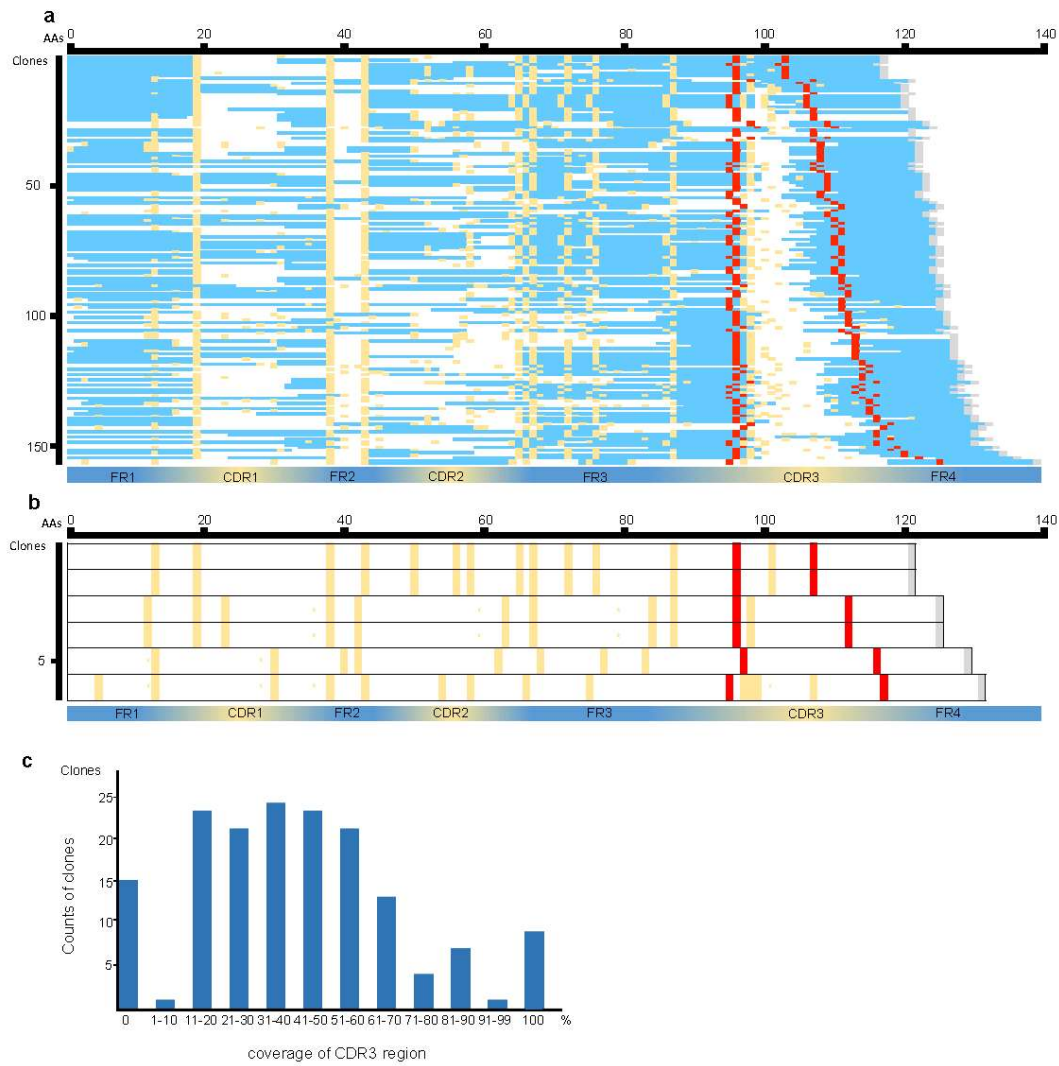
614

615 **Figure 2.**

616 Proteomic platform for the identification of peptides derived from purified GM-  
617 CSF autoantibodies (GMABs). Serum from a patient was processed for the  
618 purification of IgG, followed by isolation of GM-CSF autoantibodies using a GM-  
619 CSF-coupled column. Isolated GMABs were loaded on SDS-PAGE. The band  
620 for the IgG-heavy chain was cut and subjected to in-gel digestion using trypsin.  
621 The resultant peptides were purified with a C18 column and subjected to nano-  
622 elute UH-PLC-coupled timsTOF pro MS. Raw MS/MS data were subjected to  
623 analysis using PEAKS Studio 8.5 for data processing.

624

625



626

627 **Figure 3.**

628 Identification of peptide sequences coded in clones listed in the cDNA database

629 using PEAKS-filtered peptide sequences.

630 a: Coverage of 156 cDNA clones partially matching the PEAKS-filtered

631 peptides. The peptides in which the amino acids were covered are shown as

632 blue bars. The red bars indicate CDR3 anchor amino acids. The yellow bars

633 indicate candidates for trypsin cleavage sites (i.e., arginine or lysine). The  
634 vertical and horizontal axes indicate the number of clones and the amino acid  
635 positions from the N-terminal end, respectively.

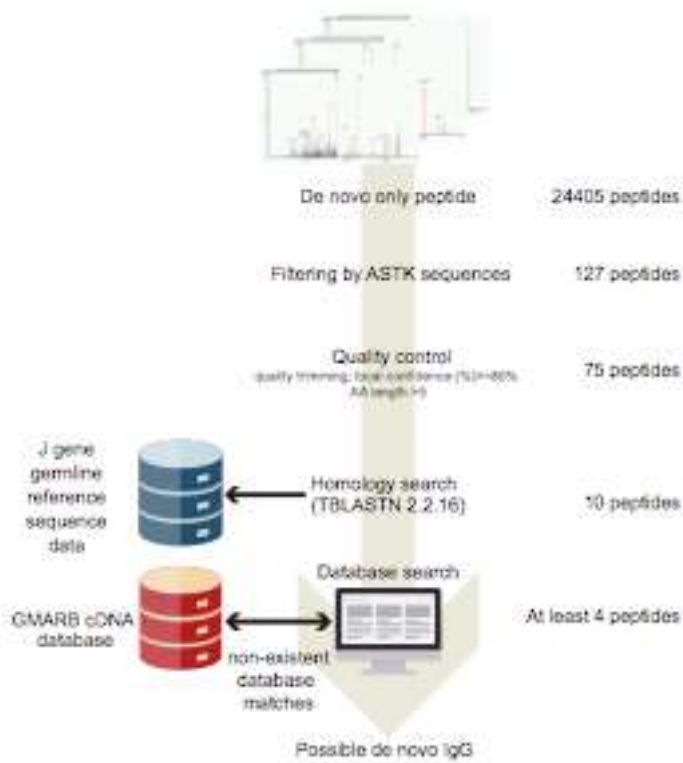
636 b: Amino acid sequences deduced from the six cDNA clones of which coverage  
637 by the PEAKS-filtered peptides was zero.

638 c: The number of clones with different amounts of coverage by the PEAKS-  
639 filtered peptides in the CDR3 region. The vertical and horizontal axes indicate  
640 the number of clones and the coverage, respectively.

641



642



643

644 **Figure 4.**

645 Procedures for the *de novo* peptide sequencing of the CDR3–FR4 region.

646 From the peptides unidentified by the PEAKS search, those containing “ASTK”

647 sequences were extracted. Seventy-five peptides with local confidence < 80%

648 were excluded. The peptides were further selected via a homology search for J

649 gene reference sequence data; four of them were not listed in the cDNA

650 database.

651

652

653 **Table 1.** Identification of IgG-HV clones by tims TOF Pro

Identification type	Clones
<b>Clones identified by timsTOF Pro [total]</b>	156
Clones with fully identified CDR3	9
Clones with partially identified CDR3	132
Clones with non-identified CDR3	15
<b>Clones not identified by timsTOF Pro</b>	6

654

655 IgG-HV were identified by timsTOF Pro followed by a PEAKS search.

656 Count of clones with CDR3 identified based on their peptides by timsTOF Pro in

657 the cDNA database. timsTOF Pro: trapped ion mobility spectrometry coupled

658 with quadrupole time-of-flight mass spectrometry.

659

**Table 2.** CDR3–FR4 sequences identified by *de novo* peptide sequencing

No.	Qualified peptides	TBLASTN_hit	FR4_matched	CDR3_matched	Putative <i>de novo</i>
1	<u>AHGSSTDH</u> <b>WGQ</b> GLTVSSAS	Yes	Yes	No	Yes
2	<u>AHGSS</u> <b>ESH</b> <b>W</b> WGQGLTVSSAS	Yes	Yes	No	Yes
3	<u>H</u> <b>W</b> WGQGLTVSSAST	Yes	Yes	No	(Yes)
4	<u>V</u> <b>W</b> WGQGLTVSSAS	Yes	Yes	N/A	
5	<u>DVP</u> <b>W</b> WGQGLTVSS	Yes	Yes	No	Yes
6	<u>L</u> <b>W</b> GLGLTVSSA	Yes	No	N/A	Yes
7	<u>Y</u> <b>W</b> WGQGLTVSS	Yes	Yes	N/A	
8	<b>W</b> WGQGLTVSS	Yes	Yes	N/A	
9	<u>Y</u> <b>W</b> WGQGPVTVS	Yes	N/A	N/A	
10	<b>W</b> WGQGLTVS	Yes	N/A	N/A	

662 N/A: not applicable

663 (Yes): “unclassified due to insufficient CDR3 length”.

664 The local confidence of each amino acid constituting the peptide was  $\geq 80\%$ .

665 VH sequences in the CDR3 region and CDR3 anchor amino acids are indicated

666 in black underlined and bold characters, respectively. “TBLASTN\_hit” indicates

667 whether the J gene germline reference sequence search via a TBLASTN

668 homology search found a match. “FR4\_matched” indicates whether the flag for

669 the FR4 region in each peptide matched the FR4 region in the cDNA database.

670 “CDR3\_matched” indicates whether the flag for the CDR3 region in each

671 peptide matched the CDR3 region in the cDNA database. “Putative *de novo*”

672 indicates whether the peptide is a putative *de novo* peptide.



674 **Abbreviations used in this article:**

675 aPAP, autoimmune PAP

676 CDR, complementarity determining region

677 FR, framework region

678 GM-CSF, granulocyte-macrophage colony-stimulating factor

679 GMAb, GM-CSF autoantibody

680 GMARBs, GM-CSF autoreactive B cells

681 IgG, immunoglobulin G

682 IgG-HV, variable region of immunoglobulin G heavy chain

683 IMGT, the international ImMunoGeneTics information system

684 LC/MS/MS, liquid chromatography coupled with tandem mass spectrometry

685 PAP, pulmonary alveolar proteinosis

686 PBMCs, peripheral blood mononuclear cells

687 PCR, polymerase chain reaction

688 TBS, tris-buffered saline

689 timsTOF, trapped ion mobility spectrometry coupled with quadrupole time-of-

690 flight mass spectrometry

691

G i d d ` Y a Y b h U f m ` ] b Z c f a U h ] c b `

.  
3URWHRJHQRPLF DQDO\VLV RI JUDQXORF\WH PDFURSKDJ  
IDFWRU DXWRDQWLERGLHV LQ WKH EORRG RI D SDWLHQ  
SXOPRQDU\ DOYHRODU SURWHLQRVLV

\$WVXVKL +DVK K K R R W R D N H X R K . I D M \$ M L D U D < D P D U R D W D D

.DNXL7DNDKLUR ~~7DQGD NDK~~ 1DNDWD

1LLJDWD 8QLYHUVLW\ 0HGLFDO 'HQWDO +RVSLWDO 1L

1LLJDWD 8QLYHUVLW\ \*UDGXDWH 6FKRRO RI 0HGLFDO D

-DSDQ

%UXNHU -DSDQ . . .

, ' ( \$ & R Q V X O W D Q W V , Q F 2 V D N D - D S D Q

, 7 K H V H D X W K R U V F R Q W U L E X W H G H T X D O O \

& R U U H V S R Q G L Q J D X W K R U . R K 1 D N D W D 0 ' 3 K '

( P D L O D G G U H V V U D G L F D O # P H G Q L L J D W D X D F M S











