

# SCIENTIFIC DATA

CONFIDENTIAL

COPY OF SUBMISSION FOR PEER REVIEW ONLY

Tracking no: SDATA-21-00719B

## *InvitroSPI and a large database of proteasome-generated spliced and non-spliced peptides*

**Authors:** Hanna Roetschke (Max-Planck-Institute for Biophysical Chemistry), Guillermo Rodriguez-Hernandez (King's College London), John Cormican (Max-Planck-Institute for Multidisciplinary Sciences), Xioping Yang (King's College London), Steven Lynham (King's College London), Michele Mishto (King's College London), and Juliane Liepe (Max-Planck-Institute for Biophysical Chemistry)

### Abstract:

Noncanonical epitopes presented by Human Leucocyte Antigen class I (HLA-I) complexes to CD8+ T cells attracted the spotlight in the research of novel immunotherapies against cancer, infection and autoimmunity. Proteasomes, which are the main producers of HLA-I bound antigenic peptides, can catalyze both peptide hydrolysis and peptide splicing. The prediction of proteasome-generated spliced peptides is an objective that still requires a reliable (and large) database of non-spliced and spliced peptides produced by these proteases. Here, we present an extended database of proteasome-generated spliced and non-spliced peptides, which was obtained by analyzing in vitro digestions of 80 unique synthetic polypeptide substrates, measured by different mass spectrometers. Peptides were identified through invitroSPI method, which was validated through in silico and in vitro strategies. The peptide product database contains 16,631 unique peptide products (5,493 non-spliced, 6,453 cis-spliced and 4,685 trans-spliced peptide products), and a substrate sequence variety that is a valuable source for predictors of proteasome-catalyzed peptide hydrolysis and splicing. Potential artefacts and skewed results due to different identification and analysis strategies are discussed.

### Datasets:

Repository Name	Dataset Title	Dataset Accession Number	URL	Reviewer Passcode
Figshare	Database, supplementary figures and scripts from 'InvitroSPI and a large database of proteasome-generated spliced and non-spliced peptides'	267e3183af100c0b64ef	<a href="https://figshare.com/s/267e3183af100c0b64ef">https://figshare.com/s/267e3183af100c0b64ef</a>	
PRIDE Archive	Digestion of TSN2 and TSN89 synthetic peptides by proteasomes	PXD025995	<a href="https://www.ebi.ac.uk/pride/archive/projects/PXD025995/private">https://www.ebi.ac.uk/pride/archive/projects/PXD025995/private</a>	reviewer_pxd025995@ebi.ac.uk; Password: FgFAsRhP

figshare Fileset <https://figshare.com/s/ee096cd7ed165c9c6318>

# 1 invitroSPI and a large database of proteasome-generated spliced and non-spliced 2 peptides

3  
4 Hanna P. Roetschke<sup>1,2</sup>, Guillermo Rodriguez-Hernandez<sup>2,3</sup>, John A. Cormican<sup>1</sup>, Xiaoping  
5 Yang<sup>4</sup>, Steven Lynham<sup>4</sup>, Michele Mishto<sup>2,3,§,\*</sup>, Juliane Liepe<sup>1,§,\*</sup>

6  
7 <sup>1</sup> Max-Planck-Institute for Multidisciplinary Sciences (MPI-NAT), 37077 Göttingen, Germany

8 <sup>2</sup> Centre for Inflammation Biology and Cancer Immunology (CIBCI) & Peter Gorer  
9 Department of Immunobiology, King's College London (KCL), SE1 1UL London, United  
10 Kingdom

11 <sup>3</sup> Francis Crick Institute, NW1 1AT London, United Kingdom

12 <sup>4</sup> Proteomics Core Facility, James Black Centre, King's College London (KCL), SE5 9NU  
13 London, UK

14  
15 \* These authors contributed equally to this work.

16 § Correspondence to: Michele Mishto ([michele.mishto@kcl.ac.uk](mailto:michele.mishto@kcl.ac.uk)) & Juliane Liepe  
17 ([jliepe@mpinat.mpg.de](mailto:jliepe@mpinat.mpg.de)).  
18

## 19 20 Abstract

21 Noncanonical epitopes presented by Human Leucocyte Antigen class I (HLA-I) complexes to  
22 CD8<sup>+</sup> T cells attracted the spotlight in the research of novel immunotherapies against cancer,  
23 infection and autoimmunity. Proteasomes, which are the main producers of HLA-I-bound  
24 antigenic peptides, can catalyze both peptide hydrolysis and peptide splicing. The prediction of  
25 proteasome-generated spliced peptides is an objective that still requires a reliable (and large)  
26 database of non-spliced and spliced peptides produced by these proteases. Here, we present  
27 an extended database of proteasome-generated spliced and non-spliced peptides, which was  
28 obtained by analyzing *in vitro* digestions of 80 unique synthetic polypeptide substrates,  
29 measured by different mass spectrometers. Peptides were identified through invitroSPI method,  
30 which was validated through *in silico* and *in vitro* strategies. The peptide product database  
31 contains 16,631 unique peptide products (5,493 non-spliced, 6,453 *cis*-spliced and 4,685 *trans*-  
32 spliced peptide products), and a substrate sequence variety that is a valuable source for  
33 predictors of proteasome-catalyzed peptide hydrolysis and splicing. Potential artefacts and  
34 skewed results due to different identification and analysis strategies are discussed.  
35

## 36 37 Background & Summary

38 Despite being well known as proteolytic enzymes for four decades, the ability of proteasomes  
39 to catalyze the reverse reaction – namely, proteasome-catalyzed peptide splicing (PCPS) – was  
40 only identified in 2004, when two independent groups identified the first examples of tumor-  
41 associated spliced epitopes <sup>1,2</sup>. The proteolytic activity of these proteases, which is mediated  
42 by peptide hydrolysis (**Fig. 1a**), has been investigated from many angles and in many  
43 experimental and translational settings. Indeed, proteasomes degrade most of the cytoplasmic  
44 proteins - including transcription factors, obsolete, damaged or wrongly transcribed proteins -  
45 and changes in their proteolytic activity have been associated with many pathological  
46 conditions. Much less is known about PCPS, which comprises the ligation of two non-  
47 contiguous peptide fragments (*i.e.*, splice-reactants) of the same molecule (*cis*-spliced peptides;  
48 **Fig. 1b,c**) or from two distinct molecules (*trans*-spliced peptides; **Fig. 1d**) <sup>3</sup>. Although *trans*-  
49 spliced peptides have been identified in both *in vitro* experiments with purified proteasomes <sup>4-8</sup>,  
50 *in cellula* <sup>9</sup>, and in HLA-I immunopeptidomes - *i.e.*, in the pool of peptides bound to HLA-I  
51 complexes <sup>10</sup> - their immunological relevance is still an enigma. In contrast, the immunological  
52 relevance of *cis*-spliced peptides has been evident since their first identification and has likely

53 been a major driver for the development of methods for their identification. From few pioneering  
54 studies we know that many *cis*-spliced peptides are produced by proteasomes and presented  
55 by HLA-I molecules of various cells<sup>10-14</sup>. They can target CD8<sup>+</sup> T cell responses against  
56 otherwise neglected bacterial antigens *in vivo*, in a mouse model of *Listeria monocytogenes*  
57 infection<sup>15</sup>. They can activate CD8<sup>+</sup> T cells specific for *Listeria monocytogenes* or Human  
58 Immunodeficiency virus (HIV) through cross-recognition *ex vivo*<sup>14,16</sup>. Preliminary *in silico*  
59 studies suggest that *cis*-spliced peptides may not play an immunologically significant role in  
60 CD8<sup>+</sup> T cell tolerance, although potential cases of viral-human epitope mimicry associated with  
61 autoimmune diseases cannot be excluded<sup>17,18</sup>. *Cis*-spliced peptides can carry cancer-specific  
62 mutations<sup>6,19</sup>, and are recognized by CD8<sup>+</sup> T cells in peripheral blood of melanoma patients  
63<sup>11,20</sup> and healthy donors<sup>20,21</sup>. A melanoma patient with metastasis was cured through adoptive  
64 T cell therapy using an autologous tumor-infiltrating lymphocyte clone, which was proved, in a  
65 later study, to be specific for a *cis*-spliced epitope derived from a melanoma-associated antigen  
66<sup>22,23</sup>.

67 The location of the catalytic sites within the inner chamber of the proteasome barrel can be one  
68 of the reasons for efficient PCPS activity<sup>24</sup>, although proteases with different structures can  
69 catalyze peptide splicing as well<sup>25-28</sup>.

70 Both peptide hydrolysis and peptide splicing can be catalyzed by different proteasome isoforms,  
71 such as 20S standard-, immuno-, and thymo-proteasomes, as well as by 20S proteasomes  
72 coupled to regulatory subunits, such as 26S proteasomes<sup>3,5,7,8,29,30</sup>. Both catalytic reactions  
73 seem to be highly tuned mechanisms, wherein the residues surrounding the substrate  
74 cleavage- and splice-sites, as well as catalytic dynamics, may play a pivotal role<sup>4,5,25,31-34</sup>. This  
75 implies that, by dissecting these driving factors, we may predict which spliced and non-spliced  
76 peptides are produced by proteasomes. PCPS predictors may be integrated in some of the  
77 pipelines that have been proposed for a targeted epitope discovery and immunotherapies<sup>6,15,35-  
78 41</sup>.

79 Such PCPS predictors should be trained on robust, validated databases of non-spliced and  
80 spliced peptides produced by proteasomes. These databases should be large and diverse  
81 enough to ensure the generalizability of the obtained predictions.

82 The identification of spliced peptides in HLA-I immunopeptidomes has a number of technical  
83 hurdles. This has ignited an intense controversy and the proliferation of identification methods  
84 with discordant performance, and thereby divergent estimation of spliced peptide frequency in  
85 HLA-I immunopeptidomes (for more details see<sup>24,42-44</sup>). Theoretically, these technical hurdles  
86 are less pronounced in a controlled experimental set up, such as *in vitro* digestion of synthetic  
87 polypeptides by purified proteasomes, measured by mass spectrometry (MS). Indeed, this kind  
88 of assay requires a much smaller spliced peptide reference database and hence results in a  
89 significantly smaller theoretical search space in the MS data analysis compared to HLA-I  
90 immunopeptidome analysis. Correspondence between *in vitro* experiments carried out with  
91 purified 20S proteasomes and *in cellula* and *in vivo* experiments has been demonstrated in  
92 various studies investigating both viral and tumor epitopes<sup>2,15,20,23,29,45-53</sup>. 20S proteasomes can  
93 degrade intrinsically disordered proteins *in vitro* and *in cellula*<sup>54-56</sup>. Recently, Specht *et al.*<sup>8</sup> and  
94 Paes *et al.*<sup>57</sup> published the first two datasets of *in vitro* digested synthetic polypeptides, and  
95 systematically identified non-spliced and spliced peptides produced by proteasomes through  
96 the analysis of MS measurements by methods specifically developed for this purpose. Our  
97 study<sup>8</sup> investigated the degradation of 55 synthetic polypeptide substrates ('Specht dataset'),  
98 whereas the dataset published by Paes *et al.*<sup>57</sup> contained 25 substrates ('PB dataset'). Despite  
99 the attempts, both datasets were too small for a statistically robust analysis of the sequence  
100 motifs (see Technical Validation Section), which we suggest being the cornerstone of any PCPS  
101 predictor development. The two studies applied different methods for the identification of non-  
102 spliced and spliced peptides. The outcomes, in term of spliced peptide frequency and features,  
103 diverged, thereby rendering unwise the merging of the two databases of non-spliced and spliced  
104 peptides produced by proteasomes. Indeed, since the objective of our study was the generation

105 of a database of non-spliced and spliced peptides produced by proteasomes through the  
106 degradation of 80 synthetic polypeptides, all digestions should be analyzed with a single peptide  
107 identification method to avoid biases arising from differences in the respective identification  
108 algorithms. Therefore, we developed an improved version of our method – namely, *in vitro*  
109 Spliced Peptide Identifier (invitroSPI; **Fig. S1**) – and implemented Paes' method (referred to as  
110 invitroPB method; **Fig. S2**); then, we applied both of them to a new small dataset (namely,  
111 'gp100 Fusion dataset') and then to the larger PB dataset, and compared their outcome by  
112 using state-of-the-art methods for the evaluation of MS2 spectra and other MS features (**Fig.**  
113 **1e**). Based on the latter outcomes, we then applied invitroSPI to the whole dataset containing  
114 the Specht dataset, the PB dataset, and the new gp100 Fusion dataset. Thereby, we generated  
115 a database of non-spliced (n = 5,493), *cis*-spliced (n = 6,453) and *trans*-spliced (n = 4,685)  
116 peptides (ProteasomeDB) - produced by proteasomes, derived from 80 synthetic polypeptide  
117 substrates and analyzed through the same method - which may be informative enough for  
118 PCPS predictor development.

119

120

## 121 **Methods**

122

### 123 **Statistical analysis.**

124 All statistical tests were done in R. Differences in distributions have been tested using either a  
125 two-sided Student's t-test, a two-samples Wilcoxon test or a Kolmogorov-Smirnov test,  
126 depending on the data distribution. Bootstrapping was applied by sampling 80 % of the data  
127 repeatedly (n = 200 iterations) and calculating the 90% confidence interval over all bootstrap  
128 results.

129

### 130 **Peptide synthesis and proteasome purification.**

131 All synthetic peptides used for MS2 spectrum comparison were synthesized using Fmoc solid  
132 phase chemistry. The 20S standard proteasomes used in this study were purified from K562  
133 cell line, as described elsewhere<sup>8</sup>. Proteasome concentration was measured by Bradford  
134 staining and verified by Coomassie staining of an SDS-Page gel, as shown elsewhere<sup>58</sup>. The  
135 purity of the proteasome preparation using this protocol has previously been shown<sup>30</sup>. The  
136 Specht dataset was generated using human 20S and 26S standard- and immuno-proteasomes<sup>8</sup>,  
137 the PB dataset was produced using human 20S standard proteasomes<sup>57</sup>, and the gp100  
138 Fusion dataset was produced using human 20S standard proteasomes (**Fig. 1e**).

139

### 140 ***In vitro* digestions and MS measurements.**

141 As part of the gp100 Fusion dataset, the synthetic polypeptides TSN2  
142 [VSRQLRTKAWNRQLYPEWTEAQR] and TSN89 [RTKAWNRQLYPEW] (final concentration of  
143 40 μM) were digested for different time points (0, 2, 4, 20 h) at 37°C by either 0.75 μg (TSN2)  
144 or 1.5 μg (TSN89) 20S proteasomes in 40 μl TKMD buffer (50 mM Tris/HCl-pH 7.8, 20 mM KCl,  
145 5 mM MgAc, 1 mM DTT). Reactions were stopped by acidification. *In vitro* digestions were  
146 measured through Orbitrap Fusion Lumos spectrometer at Centre of Excellence of MS (CEMS)  
147 at King's College London (KCL) as follows: either 5 μl of *in vitro* digestion samples or 2 μl gp100-  
148 PMM\_210325 synthetic peptide library were injected using an Ultimate 3,000 RSLC nano pump  
149 (both from ThermoFisherScientific). Briefly, peptides were loaded and separated by a nanoflow  
150 HPLC (RSLC Ultimate 3000) on an Easy-spray C18 nano column (50 cm length, 75 μm internal  
151 diameter; ThermoFisherScientific). Peptides were eluted with a linear gradient of 5%–55%  
152 buffer B (80% ACN, 0.1% formic acid) at a flow rate of 300 nl/min over 100 min at 45°C. The  
153 instrument was programmed within Xcalibur 4.4 to acquire MS data using a "Universal" method  
154 by defining a 3 s cycle time between a full MS scan and MS2 fragmentation. We acquired one  
155 full-scan MS spectrum at a resolution of 120,000 at 200 m/z with a normalized automatic gain  
156 control (AGC) target (%) of 250 and a scan range of 300~1,600 m/z. The MS2 fragmentation

157 was conducted using HCD collision energy (35%) with an orbitrap resolution of 30,000 at 200  
158 m/z. The AGC target (%) was set up as 200 with a max injection time of 128 ms. A dynamic  
159 exclusion of 30 s and 1-7 included charged states were defined within this method.

160 Gp100-PMM\_210325 synthetic peptide library contained peptides and splice-reactants  
161 previously identified (or just investigated) in TSN2 and TSN89 substrate degradations<sup>5,20</sup>. Each  
162 peptide was present in a concentration of 0.4  $\mu$ M (**Table S1**).

163 The Specht<sup>8</sup> and PB<sup>57</sup> datasets were originally measured through either LTQ XL, Q Exactive  
164 Plus and Q Exactive Orbitrap or Fusion Lumos Orbitrap mass spectrometers, respectively.

165 All collected MS RAW files were converted to the Mascot Generic Format (MGF) using  
166 ProteoWizard msconvert, employing the vendor peak picking option. RAW files that contained  
167 XL Ion Trap and XL Orbitrap scans were split into separate files for each mass analyzer type.  
168 Afterwards, headers containing search parameters were added to the MGF files and matched  
169 using Mascot v2.7.01 and PEAKS v10.5 (and PEAKS v8.5) with a mass tolerance of either 10  
170 ppm (for XL mass spectrometer), 6 ppm (for Q Exactive Orbitrap mass spectrometer) or 5 ppm  
171 (for Orbitrap Fusion Lumos mass spectrometer) on precursor masses. Mass tolerance of  
172 fragment ions was set at either 0.5 Da (for Iontrap XL mass spectrometer in CID mode), 20 ppm  
173 (for XL and Q Exactive Orbitrap mass spectrometers in HCD mode), 0.02 Da (for the Orbitrap  
174 Fusion Lumos mass spectrometer in HCD mode at Proteomics Core Facility, KCL), and 0.03  
175 Da (for the Orbitrap Fusion Lumos mass spectrometer in HCD mode at Proteomics Core  
176 Facility, University of Oxford). All MS measurements derived from a given synthetic polypeptide  
177 substrate were analyzed together in all investigated methods.

178

#### 179 ***In vitro* digestion datasets and peptide product database.**

180 In the Specht dataset (55 synthetic polypeptide substrates), *in vitro* digestions of 48 synthetic  
181 substrates have been measured by XL MS at Charité Shared Facility for MS, 4 and 10 synthetic  
182 substrates have been measured by Q Exactive Orbitrap at Charité Shared Facility for MS and  
183 by Q Exactive Orbitrap at MPI-NAT Core Facility for Proteomics, respectively. *In vitro* digestions  
184 of 47 synthetic substrates have been carried out with human 20S standard proteasomes for  
185 4 h. For four synthetic substrates, *in vitro* digestions have also been carried out with human 20S  
186 immunoproteasomes. For one synthetic substrate, *in vitro* digestions have also been carried  
187 out with human 20S and 26S standard- and immuno-proteasomes<sup>8</sup>.

188 In the original PB dataset, *in vitro* digestions of 25 synthetic substrates have been measured by  
189 Orbitrap Fusion Lumos at the MS Centre of Jenner Institute (University of Oxford)<sup>57</sup>. To note,  
190 no product sequences were detected in the control PP9 (TSN108) substrate of the original PB  
191 dataset using Mascot search engine. Thus, potential synthesis errors and contaminants related  
192 to the TSN108 substrate could not be identified and removed in the final peptide product  
193 database (see invitroSPI and invitroPB method description below).

194 In the peptide product database published by Paes *et al.*<sup>57</sup>, *cis*-spliced peptides were detected  
195 in only 16 out of 25 synthetic substrates after 2 h digestion. After applying downstream filtering  
196 steps that were described by Paes *et al.*<sup>57</sup>, *i.e.*, removing all peptides carrying the substrate's  
197 N- or C-termini, the original peptide product database that contained *cis*-spliced peptides was  
198 restricted to 12 synthetic polypeptide substrates. This final peptide product database has been  
199 used for the latter part of the Technical Validation section (see below).

200 For the present study, we generated the gp100 Fusion dataset, which contained the gp100-  
201 derived TSN2 and TSN89 substrate digestions that have been measured through Orbitrap  
202 Fusion Lumos at Proteomics Core Facility (KCL). TSN2 and TSN89 substrates were already  
203 present in the Specht dataset, although the experiments were performed in different conditions,  
204 and were measured through a different mass spectrometer (**Fig. 1e**).

205

#### 206 **Proteasome-generated peptide product database.**

207 Our whole peptide product database (ProteasomeDB) contains non-spliced and spliced  
208 peptides produced in proteasome-mediated *in vitro* digestions of 80 unique synthetic

209 polypeptide substrates. The latter is the whole dataset containing the three datasets described  
210 above (**Online-Table 1**). The peptide products were identified by applying invitroSPI method.  
211 In the entire study, we reported the number of 'unique peptides per substrate', which we  
212 speculate will be more useful for the development of proteasome activity predictors than the  
213 'unique peptides' unrelated to the substrate origin. Therefore, if a peptide sequence was  
214 generated, for example, from 2 distinct substrates, it was reported as two distinct unique  
215 peptides per substrate in this study. However, ProteasomeDB structure allows the user to adopt  
216 different strategies for the computation of unique peptides, depending on the user's goal.  
217 Peptides have been produced by various proteasome isoforms and conditions, in 0, 2, 4, 20/24  
218 h *in vitro* experiments at 37° C. Samples containing either only synthetic substrates – *i.e.*,  
219 without proteasomes – left for 20 h at 37° C, or synthetic substrates and proteasomes left for 0  
220 h at 37° C, have been used as negative control. For each substrate, 1-4 biological replicates  
221 have been carried out, and measured 1-5 times.  
222 The length of synthetic polypeptide substrates varies from 13 to 47 amino acids (**Online-Table**  
223 **1**). They have an amino acid frequency that is similar to the frequency present in the human  
224 proteome<sup>8,57</sup>. The polypeptides are derived from bacterial, viral and human proteins (largely  
225 antigens). In the Specht dataset (comprising 55 synthetic polypeptide substrates), there is a  
226 preponderant presence of tumor-associated or autoimmune disease-associated antigens. In  
227 the PB dataset (comprising 25 synthetic polypeptide substrates), there is a preponderant  
228 presence of HIV antigens. The species of origin of the substrate and unique identifier of the  
229 substrate sequences are attributes of our ProteasomeDB database (see **Table 1**).  
230 Experiments have been carried out with synthetic polypeptides rather than the entire protein  
231 because purified proteasomes have been shown to hardly process entire proteins *in vitro*, likely  
232 because ligases and cofactors are lost during 20S/26S proteasome purification<sup>59</sup>. However, a  
233 correspondence between *in vitro* experiments - with synthetic polypeptides and purified  
234 proteasomes - and *in cellula* and *in vivo* experiments has been widely demonstrated (see text  
235 above).  
236 Each digestion has been performed with a single polypeptide as substrate. Therefore, non-  
237 spliced and *cis*-spliced peptides could be produced by processing of a single molecule of the  
238 substrate (**Fig. 1a-c**), whereas *trans*-spliced peptides resulted from the ligation of two partially  
239 overlapping fragments derived from two molecules of the same substrate (**Fig. 1d**). *Trans*-  
240 spliced peptides with splice-reactants from two different substrate sequences were not possible  
241 because each *in vitro* digestion contained only one substrate rather than various substrates.  
242 *In vitro* digestions have been performed at 0, 2, 4 and 20/24 h, and peptide products have been  
243 identified by applying invitroSPI method, which removed synthesis artefacts from the final list of  
244 identified peptide products (see Technical Validation section). To note, peptide synthesis  
245 artefacts can arise due to synthesis errors during the Fmoc solid phase chemistry, peptides  
246 that contaminated the samples during their preparation, or other forms of contaminations. Both  
247 types of contaminations are termed synthesis errors, in this study. The synthesis errors  
248 generated during the peptide synthesis by Fmoc solid phase chemistry could be: (i) truncated  
249 peptides that are shorter than the cognate synthetic polypeptide substrate at the N- and/or C-  
250 terminus, (ii) peptides lacking one or more residues within their sequence (*i.e.* not at the  
251 substrate termini), (iii) peptides containing the duplication of one (or more) amino acid. The  
252 example (i) could result in the wrong assignment of both non-spliced and spliced peptides, the  
253 example (ii) in the wrong assignment of *cis*-spliced peptides, and the example (iii) in the wrong  
254 assignment of *trans*-spliced peptides (**Fig. 1f**).  
255 Since substrate degradation rates varied from substrate to substrate, from proteasome  
256 preparation to preparation, *in vitro* reaction conditions were set up to have the 2-4 h time points,  
257 wherein substrate molecules were still present in the reaction, and 20/24 h time point, wherein  
258 most of the substrate molecules have been processed by proteasomes. The presence of intact  
259 substrate molecules in the reaction can be determined by analyzing the MS RAW files linked to  
260 our database (see Data Record section).

261 Compared to the previous version of the peptide product database <sup>8</sup>, in ProteasomeDB, we  
262 expanded the number of substrates, their sequence variety and origin, as well as we strongly  
263 increased the number of digestion samples measured with high accuracy Orbitrap MS. In fact,  
264 ProteasomeDB contains proteasome-generated peptide products of 80 synthetic polypeptide  
265 substrates that have been measured with Orbitrap mass spectrometers with a mass tolerance  
266 of 5 - 6 ppm on precursor masses, and 20 ppm or 0.02 – 0.03 Da for fragment ions (**Online-**  
267 **Table 1**). Furthermore, the improved performance of invitroSPI increased the precision of  
268 peptide identification (see below).

269 ProteasomeDB is a CSV table, which contains 26 columns describing features of the identified  
270 peptides, the original substrate sequence, sample processing and instrument parameters (see  
271 **Table 1** for a detailed description of the database columns/attributes). Additional to the  
272 information provided in the Specht database of peptide products <sup>8</sup>, this new database contains  
273 all possible multi-mapper peptides with their correct splice-type annotation (see Technical  
274 Validation section).

275

### 276 **Prediction of MS2 spectra.**

277 Prosit version 2020 <sup>60,61</sup> allows prediction of the MS2 spectra given a peptide sequence,  
278 precursor charge and calibrated collision energy. A predicted MS2 spectrum can be compared  
279 to the detected MS2 spectrum by computing a similarity score. In this study, we used the  
280 spectral angle between the L2 normalized spectra, also known as normalized spectral contrast  
281 angle <sup>62</sup>, which ranges from 0 (very bad match between MS2 spectra) to 1 (perfect match  
282 between MS2 spectra). The spectral angle consists of a transformation on the normalized dot  
283 product and corresponds to the loss metric on which Prosit was trained.

284

### 285 **Generation and analysis of simulated background databases.**

286 In order to identify proteasome specificities, a simulated background database containing a  
287 subset of all theoretically possible spliced and non-spliced peptides was generated, similar to  
288 what was previously described in Specht *et al.* <sup>8</sup>. The simulated background database reflected  
289 the peptide products that one would expect to be detected in absence of any proteasome  
290 specificities, *i.e.*, under the assumption that each theoretically possible spliced peptide is  
291 generated with the same probability. The simulated background database was obtained by  
292 sampling uniformly a subset of all theoretically possible spliced and non-spliced peptides (*i.e.*,  
293 a subset of the custom reference database that was also used for the MS search). In that sense,  
294 the peptide products were randomized. This simulated background database could then be  
295 compared to the database of experimentally identified peptide products. Thereby, we could  
296 verify whether the identification of spliced and non-spliced peptide characteristics (*e.g.*, splice-  
297 reactant, intervening sequence and peptide lengths, as well as amino acid frequencies) arose  
298 from theoretical database structure – and thus were potential analysis artefacts - or from  
299 biochemical drivers of the catalytic reaction. In this study we made use of the simulated  
300 background database to investigate amino acid preferences of forward and reverse *cis*-spliced  
301 peptides.

302

### 303 **Mapping of peptide sequences. Identification of peptides containing N- or C-termini of** 304 **substrates. Identification of spliced peptides with one amino acid long splice-reactant.**

305 Peptide sequences were mapped to a substrate sequence by exact string matching of the  
306 complete peptide product sequence. If this was not possible, the peptide product sequence was  
307 split into two splice-reactants at each possible position. Each pair of splice-reactants was then  
308 matched against the substrate sequence. If both splice-reactants could be matched to the  
309 substrate sequence, the respective locations within the substrate were recorded.

310 If a peptide sequence could be explained by multiple locations, all locations have been reported  
311 in the final database. However, when we computed frequency and features of product types,  
312 we applied the following rules: (i) if a sequence could be both a non-spliced and a spliced

313 peptide, we defined it as non-spliced peptide; (ii) if a sequence could be both a *cis*-spliced and  
 314 a *trans*-spliced peptide, we defined it as *cis*-spliced peptide; (iii) if a sequence could be both a  
 315 forward *cis*-spliced and a reverse *cis*-spliced peptide, we defined it as forward/reverse *cis*-  
 316 spliced peptide (*i.e.*, multi-mapper *cis*-spliced peptides). Implications of such multi-mapper  
 317 peptides, *i.e.*, peptides that map to multiple locations in the substrate, are discussed below.  
 318 A peptide with several potential substrate origins was assigned to the category “peptides  
 319 containing N- or C-termini of their cognate synthetic polypeptide substrate” only in case all  
 320 possible peptide locations contained the substrate’s N- or C-terminus. Analogously, a peptide  
 321 with several potential substrate origins was assigned to the category “spliced peptides with one  
 322 amino acid long splice-reactant” only if none of the possible origins resulted in longer splice-  
 323 reactants.

324

325 **Calculation of all possible *cis*-spliced and non-spliced peptide products to investigate**  
 326 **length and presence of substrate’s N- or C-termini.**

327 The number of possible unmodified spliced and non-spliced peptides that could be derived from  
 328 a protein sequence in sequence-agnostic fashion formed the theoretical sequence search  
 329 space. The number  $X$  of non-spliced peptides of length  $N$  that could theoretically arise from a  
 330 substrate of length  $L$  was:

331 
$$X_{non-spliced} = L - N + 1$$

332 To derive the theoretical number of all spliced peptides, we defined four indices  $i, j, k$  and  $n$  that  
 333 denoted the positions of the first ( $i, j$ ) and second ( $k, n$ ) splice-reactant, respectively. The  
 334 corresponding number of peptides was calculated via summing over interval ranges that form  
 335 valid spliced peptides. *Cis*-spliced peptides could be formed via forward or reverse ligation. The  
 336 number of all forward *cis*-spliced peptides of length  $N$  that could theoretically arise from a  
 337 substrate of length  $L$  was:

338 
$$X_{fwd. cis-spliced} = \sum_{i=1}^{L-N} \sum_{j=i+L_{ext}-1}^{N-L_{ext}+i-1} \sum_{k=j+2}^{L-N+j-i+2} 1 = \frac{1}{2}(N - 2L_{ext} + 1)(L - N)(L - N + 1)$$

339  $L_{ext}$  denoted the minimal splice-reactant length and was set to 1 per default. In case a peptide  
 340 was located at either of the substrate’s termini ( $i = 1$  or  $n = L$ ), the number of forward *cis*-  
 341 spliced peptides was calculated according to:

342 
$$X_{fwd. cis-spliced at termini} = (N - 2L_{ext} + 1)(L - N)$$

343 Analogously, the number of theoretically possible reverse *cis*-spliced peptides was calculated  
 344 as:

345 
$$X_{rev. cis-spliced} = \sum_{k=1}^{L-N+1} \sum_{n=k+L_{ext}-1}^{N-L_{ext}+k-1} \sum_{i=j+1}^{L-N+n-k+2} 1 = \frac{1}{2}(N - 2L_{ext} + 1)(L - N + 1)(L - N + 2)$$

346 
$$X_{rev. cis-spliced at termini} = (N - 2L_{ext} + 1)(L - N + 1)$$

347 To calculate the number of theoretical *trans*-spliced peptides in an *in vitro* scenario where a  
 348 single synthetic polypeptide substrate was digested with purified proteasome, the following  
 349 formula was derived:

350 
$$X_{trans-spliced} = -1 + \frac{2}{3}L_{ext}^3 + L_{ext}^2(-1 - N) + \frac{5}{6}N + N^2 - \frac{5}{6}N^3 + L(-1 + L_{ext}(2 - 2N) + N^2)$$
  
 351 
$$+ L_{ext} \left( \frac{7}{3} - 3N + 2N^2 \right)$$

352 
$$X_{trans-spliced at termini} = N(N - 2L_{ext} + 1)$$

353 To note, the number of non-spliced peptides of length  $N$  that could be derived from either of the  
 354 substrate’s termini was 2.

355

356 **InvitroSPI and invitroPB pipelines.**

357 The computational pipelines of invitroSPI and invitroPB differ as follow (**Fig. 2**):



- 358 a) Both invitroSPI and invitroPB adopted conservative approaches by favoring the assignment  
359 of non-spliced over spliced peptides to counteract the imbalance of the theoretical sequence  
360 search space (**Fig. S1, S2**). Indeed, the theoretical search space computed from the 80  
361 substrate sequences of the whole dataset is 400-fold larger for spliced compared to non-  
362 spliced peptides, and significantly larger for *trans*- vs. *cis*-spliced peptides (**Fig. S3**).  
363 InvitroPB can identify only non-spliced and *cis*-spliced peptides, whereas invitroSPI can also  
364 identify *trans*-spliced peptides. Their inclusion in the final peptide product database could  
365 enrich the information that may be used to understand proteasome catalytic activities.  
366 InvitroSPI can identify peptide-spectrum matches (PSMs) that may be *trans*-spliced  
367 peptides, assign them if there is no better non-spliced candidate and the scan fulfills all  
368 quality criteria described above. Although this strategy may lead to a higher FDR for *trans*-  
369 spliced peptides compared to non-spliced peptides (see below), it may avoid the  
370 misassignment of MS2 spectra to non-spliced peptides, which, in reality, are *trans*-spliced  
371 peptides;
- 372 b) InvitroSPI applies a general threshold of at least 5 amino acid length for all peptides, and  
373 therefore does not apply different restrictions of peptide length between product types.  
374 InvitroPB, on the contrary, sets a different minimal length threshold for *cis*-spliced (8 amino  
375 acids) and non-spliced (5 amino acids) peptide candidates;
- 376 c) InvitroSPI allows the identification of spliced peptides with a splice-reactant length of one  
377 amino acid. These peptides could not be identified through invitroPB. To note, proteasomes  
378 can perform a second cleavage on a spliced peptide, thereby reducing the length of a splice-  
379 reactant to one amino acid after the PCPS reaction. This event was described *in vitro* and  
380 *in cellula* for a gp100-derived *cis*-spliced epitope by Michaux and colleagues<sup>51</sup>. That *cis*-  
381 spliced epitope was also demonstrated to be recognized by CD8<sup>+</sup> T cells of melanoma  
382 patients<sup>20</sup>, thereby confirming that *cis*-spliced peptides with a one amino acid long splice-  
383 reactant can be produced by proteasomes, and in an amount sufficient to be presented and  
384 to trigger a CD8<sup>+</sup> T cell response;
- 385 d) InvitroSPI allows the identification of non-spliced and spliced peptides with two post-  
386 translational modifications (PTMs) – *i.e.*, N/Q deamidation and M oxidation. On the contrary,  
387 the implemented invitroPB method does not allow the identification of peptides with PTMs,  
388 although it removes any query matched through PEAKS-PTM for the downstream *cis*-  
389 spliced peptide identification. PEAKS-PTM performs an open search of 313 PTMs, which  
390 could have similar statistical challenges as the identification of spliced peptides, since they  
391 both strongly increase the peptide sequence search space for MS2 spectrum assignment.  
392 In addition, many PTMs could not occur during the synthesis and *in vitro* digestions of the  
393 polypeptide substrates, such as ubiquitination or phosphorylation; therefore, we think that  
394 their prioritization over spliced peptides is not supported by biological evidence and may  
395 reduce the method's recall of *cis*-spliced peptides (see below).
- 396 e) Both invitroSPI and invitroPB adopt approaches to tackle the issue of the synthesis errors  
397 inherent in the synthetic polypeptides. The synthesis errors may appear as the product of  
398 PCPS if amino acids are skipped or added more than once during synthesis or may arise  
399 through hydrolysis of a contamination (here referred to as synthesis errors; **Fig. 1f**). In  
400 contrast to invitroPB method, invitroSPI adopts a more conservative approach since it  
401 removes spliced peptides not only if they are identified as such in control samples, but also  
402 if any longer spliced peptide containing the same splice-site is identified in control samples.

403  
404

#### Technical aspects of invitroSPI (Fig. S1).

405 The invitroSPI method is an improvement on the method previously described in Specht *et al.*  
406<sup>8</sup>, which was developed to tackle the issue of synthesis errors and the large number of  
407 theoretically possible spliced peptides that could be derived from one substrate in the database.  
408 Briefly, MS RAW files were converted to MGF with ProteoWizard msconvert, using the vendor  
409 peak picking option. Data have been searched against a reference custom database containing

410 all theoretically possible *cis*-spliced, *trans*-spliced and non-spliced peptides derived from the  
411 substrate of interest and with a minimal length of at least 5 amino acids. The custom reference  
412 databases were generated in FASTA format as previously described<sup>63</sup>. Briefly, we generated  
413 all possible spliced and non-spliced peptides as follows: (i) in the case of non-spliced peptides,  
414 by applying a single sliding window over the substrate sequence. The sliding window could vary  
415 in its size, reflecting a variable length of the peptide product; (ii) in the case of spliced peptides,  
416 we applied two sliding windows, which were *in silico* ligated if they could form a valid spliced  
417 peptide, as determined by their substrate origins.

418 The following variable modifications have been set whilst applying invitroSPI method to both  
419 non-spliced and spliced peptides: asparagine (N) and glutamine (Q) deamidation and  
420 methionine (M) oxidation. All ranked PSMs suggested by the Mascot Server for a single MS2  
421 scan (query) were mapped to all potential origins in the substrate sequence, thereby  
422 considering the redundancy of leucine (L) and isoleucine (I) (see 'Mapping of peptide  
423 sequences' section below). Subsequently, PSMs have been evaluated based on product type  
424 (spliced vs non-spliced) and differences in ion scores to determine the most probable peptide  
425 sequence and origin. Scans that did not allow for the high-confidence identification of a single  
426 peptide were not assigned and removed from further analysis. For all PSMs, the mandatory  
427 condition for the peptide identification was: (i) the Mascot ion score was higher than 20, (ii) the  
428 Mascot q-value was lower than 0.05. In the case that the top-ranked peptide was a spliced  
429 peptide, it was considered a correct PSM if the difference in Mascot ion score between the first-  
430 ranked and the second-ranked peptide (either non-spliced or spliced) was larger than 30%, *i.e.*,  
431 the delta score was larger than 0.3. This optimal delta score was determined by FDR estimation  
432 (see below, **Fig. S4**). In case there were several non-redundant sequences with identical scores  
433 identified, the scan was assigned only if there was a single, non-ambiguous non-spliced peptide  
434 among them that passed all other criteria mentioned above. This approach favors the  
435 assignment of non-spliced peptides over spliced peptides to counteract the imbalance of the  
436 large theoretical number of spliced and non-spliced peptides in the MS search space.

437 To select the best delta score for invitroSPI in the datasets investigated in this study, we applied  
438 invitroSPI to the PB dataset repeatedly, while varying the delta-score in a range from 0 to 0.5.  
439 The identified PSMs were subsequently compared to the predicted MS2 spectrum by  
440 application of ProSIT<sup>61</sup> and computation of spectral angles between normalized MS2 spectra for  
441 each product type. The spectral angle distribution for the identified non-spliced peptides  
442 represented the 'gold-standard' to which all other spectral angle distributions resulting from  
443 spliced peptides were compared to. If we assumed 1% False Discovery Rate (FDR) among the  
444 non-spliced peptides, we could determine a spectral angle cut-off, for which 1% of the non-  
445 spliced peptide PSMs fall below this cut-off and 99% of the non-spliced peptide PSMs fall above  
446 this cut-off. The same cut-off applied to spliced peptides allowed us to estimate the FDR for *cis*-  
447 spliced peptides, *trans*-spliced peptides or all spliced peptides compared to non-spliced  
448 peptides. For each investigated delta score, the FDRs for each product type were estimated  
449 and the delta score resulting in lowest FDR for spliced peptides was selected (*i.e.*, delta score  
450 = 0.3; **Fig. S4**).

451 Spliced peptides generated by ligation of three or more fragments were not allowed and  
452 therefore are not included in our database.

453 As in Specht *et al.*<sup>8</sup>, for each substrate digestion, peptide synthesis artefacts identified in control  
454 samples (either 0 h digestion time or samples with substrates and no proteasomes) were  
455 removed as follows: any non-spliced peptide identified in control samples was removed from  
456 the final list of identified non-spliced peptides. Any spliced peptide in the control samples,  
457 containing the same splice-site as an identified peptide (thus, either identified as such or  
458 identified as a longer precursor in control samples) was removed from the final list of identified  
459 spliced peptides (**Fig. 1f**, **Fig. S1**).

460 InvitroSPI is available as a user-friendly and readily executable tool on GitHub (see Code  
461 Availability section).

462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513

### Technical aspects of invitroPB (Fig. S2).

We implemented Paes' method based on the information provided in the original publications and source code <sup>14,57</sup>. Briefly, in invitroPB, MS data were first searched against a reference custom database containing only a given substrate sequence using PEAKS DB (PEAKS v10.5). Additionally, an open search for PTMs using PEAKS PTM (313 variable PTMs included) was performed. Although PSMs of non-spliced peptides with PTMs were not further considered, their corresponding MS2 spectra were dismissed and not further investigated. To note, while the original method described by Paes *et al.* <sup>57</sup> discarded all PTM-labelled non-spliced peptides during assignment, our invitroPB implementation recorded PTM-labelled non-spliced peptides. Those peptides were, however, not considered for downstream analyses; recording them served solely the purpose of dissecting the outcomes of the steps of method's strategy.

MS2 spectra not assigned as non-spliced peptides (with or without PTMs) with 5% PEAKS-computed FDR were re-searched using PEAKS De novo (without PTMs), which also converted all possible I to L amino acids. For the following analysis the top 100 *de novo* candidate sequences per MS2 spectrum with an ALC score equal or larger than 50 were exported, but only those *de novo* sequences within the top 5 ALC scores were further considered. All *de novo* sequences within the top 5 ALC scores were screened to determine if they could be generated through PCPS from the given substrate sequence upon exchange of all Is with Ls. All sequences that could be explained as non-spliced peptide sequences were removed. Among the remaining sequences, the implemented method computed those that could be *cis*-spliced peptides with splice-reactant length larger than 1 amino acid and a peptide length larger than 7 amino acids, which were then kept. Therefore, invitroPB could not identify *trans*-spliced peptides, *cis*-spliced peptides with a 1 amino acid long splice-reactant, and *cis*-spliced peptides with a length smaller than 8 amino acids. If more than one *cis*-spliced peptide candidate per MS2 spectrum was listed, only the peptide sequence with highest ALC score was kept and considered as the assigned sequence to that MS2 spectrum. Non-spliced and *cis*-spliced peptides, identified in the samples containing substrates but not proteasomes, were removed to exclude peptide products that may arise from peptide synthesis errors (Fig. 1f). As downstream filtering steps, Paes *et al.* <sup>57</sup> and invitroPB did not further consider non-spliced and *cis*-spliced peptides carrying the N- or C-termini of synthetic polypeptide substrates within their sequence.

As technical validation of our implementation, we applied invitroPB to the PB dataset, and obtained a partially different non-spliced and *cis*-spliced peptide list than published by Paes *et al.* <sup>57</sup> (Fig. S5a). This difference could in part be explained by a different PEAKS version applied by Paes *et al.* <sup>57</sup> – *i.e.*, PEAKS v8.0 – and by invitroPB (PEAKS v10.5) <sup>64</sup>. Indeed, when we applied invitroPB method - using either PEAKS v8.5 or v10.5 – on *in vitro* digestions of six substrates of the PB dataset, we observed some differences in the non-spliced and *cis*-spliced peptides list (Fig. S5b). Similarly, we noted a difference between the spliced and non-spliced peptides published by Paes *et al.* <sup>57</sup> and the spliced and non-spliced peptides derived using invitroPB method using PEAKS v8.5 (Fig. S5c), which could have been explained in a *corrigendum* by the same authors published during the revision of the current manuscript <sup>65</sup>. Nonetheless, in invitroPB, which used the better performing PEAKS v10.5 <sup>66</sup>, the pipeline and filtering steps of the original study were conserved, which allowed a proof-of-principle comparison of invitroPB and invitroSPI.

### Data Records

The MS files (.RAW, .mgf and search result files) of the Specht dataset <sup>8</sup> are available at the PRIDE repository <sup>67</sup> with the dataset identifier PXD016782 <sup>68</sup>.

The MS .RAW and .mgf files of the PB dataset are available at the PRIDE repository with the dataset identifier PXD021339 and PDX025893 <sup>57</sup>.

514 The MS files (.RAW, .mgf and search result files) of the gp100 Fusion dataset are available at  
515 the PRIDE repository with the dataset identifier PXD025995<sup>69</sup>. The reference custom  
516 databases that contain all theoretically possible spliced and non-spliced peptides and that were  
517 used to perform the MS search are available in a Figshare repository<sup>70</sup>.

518 The final database – *i.e.* ProteasomeDB - with all identified spliced and non-spliced peptide  
519 products, as well as their substrate sequences, is provided as CSV file, and is available in a  
520 Figshare repository<sup>70</sup>. In ProteasomeDB, all I/s of identified peptide products were replaced by  
521 Ls, whereas the substrate sequence contains the original I/L amino acids.

522 All 'online-figures' and 'online-tables' reported are available in a Figshare repository<sup>70</sup>.

523

524

## 525 Technical Validation

526

### 527 Comparison and validation of invitroSPI and invitroPB methods in gp100 Fusion dataset.

528 Our aim was to create ProteasomeDB – a database of non-spliced and spliced peptides  
529 produced *in vitro* by proteasomes and reliably identified by a single method with the highest  
530 recall of peptide products. Hence, we initially compared invitroSPI and invitroPB, to then select  
531 a single method and apply it to the whole dataset, thereby generating ProteasomeDB. Due to  
532 its dependence on *de novo* peptide sequencing, which relies on high-precision MS data,  
533 invitroPB could not be applied to the vast majority of digestions in the Specht dataset. Therefore,  
534 we initially validated and compared invitroSPI and invitroPB through the analysis of the PB  
535 dataset and the gp100 Fusion dataset by investigating methods' features and performances.  
536 We put particular attention in dissecting the several filtering steps of the two methods (**Fig. 2**,  
537 **Fig. S1-S2**) and their impact on PSM identifications.

538 The gp100 Fusion dataset contained two substrates, TSN2 and TSN89 (**Fig. 1e**). TSN89 is a  
539 subsequence of TSN2, which is the gp100<sub>35-53</sub> sequence. Two spliced epitopes immunogenic  
540 in melanoma patients have been identified within this sequence<sup>2,20</sup>, including the first *cis*-spliced  
541 epitope initially described by Vigneron and colleagues<sup>2</sup>. We measured the *in vitro* digestions  
542 (0, 2, 4, 20 h) of the synthetic polypeptide substrates with human 20S standard proteasomes  
543 through highly-sensitive Orbitrap Fusion Lumos mass spectrometers. We then applied  
544 invitroSPI and invitroPB method to the MS files. For each scan, both methods aim to assign the  
545 most likely PSM. A single unique peptide sequence can be assigned to multiple MS2 scans.  
546 InvitroSPI identified a larger number of unique non-spliced peptides, *cis*-spliced and *trans*-  
547 spliced peptides as compared to invitroPB (**Table 2**). This generally reflected what we observed  
548 at PSM level, albeit invitroPB method assigned more PSMs to non-spliced peptides compared  
549 to invitroSPI (**Fig. 3a**). InvitroSPI discarded, as synthesis errors, hundreds of PSMs of potential  
550 *cis*-spliced peptides, whereas invitroPB method eliminated only 10 of them (**Fig. 3a**, **Online-**  
551 **Table 2**). Upon removal of the synthesis errors, over 700 PSMs were discarded by invitroPB  
552 method because they were suggested to be PTM-modified non-spliced peptides by PEAKS-  
553 PTM (**Fig. 3b**). One of them was assigned to a spliced peptide sequence by invitroSPI. Both  
554 methods assigned fewer PSMs to forward than reverse *cis*-spliced peptides (**Fig. 3c**). InvitroSPI  
555 assigned over a hundred PSMs to spliced peptides with a one amino acid long splice-reactant,  
556 and over 500 PSMs to spliced peptides containing N- or C-termini of the substrates. These  
557 peptides were not identified by invitroPB analysis because of the different strategy of this  
558 method (**Fig. 3d,e** and **Online-Table 2**).

559 We also compared the MS2 spectra assigned by the two identification methods to the MS2  
560 spectra of a pool of synthetic non-spliced, *cis*-spliced, and *trans*-spliced peptides, which have  
561 been previously investigated<sup>2,5,9,20,29</sup> (**Table S1**). Among these peptides, both invitroSPI and  
562 invitroPB method identified many non-spliced and *cis*-spliced peptides, in addition to *trans*-  
563 spliced peptides, which could be identified only by invitroSPI (**Online-Fig. 1**, **Online-Table 2**).  
564 Both methods identified the *cis*-spliced epitopes TSN89<sub>1-3/6-13</sub> (gp100<sub>40-42/47-52</sub>) [RTK][QLYPEW]

565 and TSN2<sub>13-18/6-8</sub> (gp100<sub>47-52/40-42</sub>) [QLYPEW][RTK] (**Fig. 3f,g**), which have been proven to be  
566 produced by proteasomes and presented by HLA-I complexes of cancer cell lines <sup>2,5,9,20,29</sup>.  
567 When considering the single time points, *i.e.*, 2, 4 and 20 h, of the digestion kinetics, invitroSPI  
568 identified 458 unique peptides upon removal of the synthesis errors whereas invitroPB identified  
569 231 unique peptides (**Table 2**). Although overall more peptides were identified at later time  
570 points compared to earlier time points by invitroSPI (**Fig. 3h**), the frequency of spliced and non-  
571 spliced peptides remained constant over time (**Fig. 3i**), in agreement with our previous  
572 observation in Specht *et al.* <sup>8</sup>.

573

#### 574 **Comparison and validation of invitroSPI and invitroPB methods in PB dataset.**

575 To compare and evaluate the performance of the two methods on a larger dataset, we next  
576 applied invitroSPI and invitroPB to the PB dataset of 25 synthetic polypeptides, digested for 2 h  
577 and 20 h with 20S standard proteasomes (**Fig. 1e**). Control samples were left for 20 h without  
578 proteasomes, but otherwise in the same conditions of the digestion kinetics. Overall, the  
579 analysis of the PB dataset confirmed what was observed on the smaller gp100 Fusion dataset.  
580 Indeed, invitroSPI identified more unique non-spliced, *cis*-spliced and *trans*-spliced peptides  
581 than invitroPB (3,413 peptides identified by invitroSPI and 2,245 peptides identified by  
582 invitroPB; **Table 2**), which was also observed at PSM level (**Online-Table 3**). As observed in  
583 the analysis of the gp100 Fusion dataset, both methods discarded PSMs of potential *cis*-spliced  
584 peptides as synthesis errors, although this filtering step was more stringent in invitroSPI (**Fig.**  
585 **4a**). After synthesis error removal in both methods, invitroPB, using PEAKS-PTM, identified and  
586 discarded over 3,000 putative PTM-labelled non-spliced peptides (**Fig. 4b**). InvitroSPI assigned  
587 around 250 PSMs of those discarded PSMs to spliced peptides. A distribution of PTMs identified  
588 at the PEAKS-PTM step of invitroPB is shown in **Fig. 4c**. Both methods assigned more PSMs  
589 to forward than reverse *cis*-spliced peptides in the PB dataset (**Fig. 4d**). In contrast to invitroPB,  
590 over 700 PSMs were assigned by invitroSPI to spliced peptides with one amino acid long splice-  
591 reactant, and over 2,000 PSMs to spliced peptides containing N- or C-termini of the substrates  
592 (**Fig. 4e,f**).

593 The two methods showed a high similarity between measured and predicted MS2 spectra  
594 (reflected by high spectral angles) for all peptide groups (**Fig. 4g**, **Fig. 5**, **Fig. S6**, **Online-Fig.**  
595 **2-3**), thereby confirming their reliable and comparable identification of PSMs. MS2 spectra were  
596 predicted by applying Prosit <sup>61</sup>. In this analysis, we considered non-spliced and spliced peptides  
597 which did not contain any cysteine residues (C), did not exceed a charge of +6 and were  
598 between 7 and 12 amino acids long, because Prosit showed a progressive decrease of its  
599 prediction performance on non-spliced peptides for longer peptides and/or peptides with higher  
600 charges (**Fig. S6**), in agreement with previously described analyses <sup>61</sup>.

601 As last step of method validation, we estimated the FDRs of invitroSPI and invitroPB in PB  
602 dataset by using the spectral angle analysis. We chose a spectral angle cut-off of 0.7 as  
603 approximative threshold to estimate the FDRs, with high-quality PSMs having a spectral angle  
604 above this threshold (**Fig. 5a**). The percentage of PSMs below this cut-off and identified as non-  
605 spliced peptides by invitroSPI was 1.4%, which could be interpreted as an estimated 1.4% FDR  
606 (**Fig. 5b**). By applying the same strategy for the computation of the FDR of spliced peptides, we  
607 estimated that invitroSPI had a 4.2% FDR for *cis*-spliced peptides and a significantly larger  
608 6.8% FDR for *trans*-spliced peptides (**Fig. 5b**). For invitroPB, the estimated FDRs were higher  
609 than those of invitroSPI for non-spliced peptides (statistical significance was reached only for  
610 non-spliced peptides). Indeed, 3.8 % of the PSMs assigned to non-spliced had a spectral angle  
611 below 0.7, which increased to a 5.3 % for *cis*-spliced peptides (**Fig. 5c**). The estimated FDRs  
612 for both non-spliced and spliced peptides identified by both methods should be considered  
613 critical in the use and evaluation of ProteasomeDB.

614

615 **ProteasomeDB – a non-spliced and spliced peptide product database computed through**  
616 **the application of invitroSPI on the whole dataset.**

617 Our comparison of invitroSPI and invitroPB on these two datasets showed that both methods  
618 successfully identified non-spliced and *cis*-spliced peptides produced by 20S proteasomes in *in*  
619 *vitro* digestions of synthetic polypeptides. However, invitroSPI systematically identified more  
620 unique non-spliced and *cis*-spliced peptides per substrate than invitroPB (Table 2), in addition  
621 to the identification of *trans*-spliced peptides. The FDR estimation hinted toward a lower FDR  
622 for invitroSPI compared to invitroPB for both non-spliced and *cis*-spliced peptides (Fig. 5).  
623 Furthermore, invitroSPI was - contrary to invitroPB – applicable to various kinds of MS and does  
624 not rely on high-precision instruments. Therefore, invitroSPI represented a suitable method for  
625 the analysis of the whole dataset of proteasome-catalyzed *in vitro* digestions of synthetic  
626 polypeptides.

627 Through the application of invitroSPI on the whole dataset of 80 substrates - derived from the  
628 combination of the PB dataset (25 substrates), the Specht dataset (55 substrates), and the  
629 gp100 Fusion dataset (TSN2 and TSN89 substrates) (Fig. 1e) - we identified non-spliced (n =  
630 5,493), *cis*-spliced (n = 6,453) and *trans*-spliced (n = 4,685) unique peptides (Table 2). They  
631 represented 33% (non-spliced peptides), 39% (*cis*-spliced peptides), and 28% (*trans*-spliced  
632 peptides) of the 16,631 unique peptides of the whole peptide product database (Table 2).

633 While the overall frequency of spliced peptides may appear high at first glance, it is worthwhile  
634 considering the number of theoretical peptide sequences here. The generation efficiency on  
635 qualitative level takes the theoretical search space - *i.e.* the number of peptides that could be  
636 theoretically produced by proteasomes – into account (Fig. S3). If we defined the generation  
637 efficiency as number of detected peptides over the theoretical number of peptides in each  
638 peptide product type, PCPS had, on average, a 280-fold lower generation efficiency than  
639 peptide hydrolysis in the whole dataset. Indeed, on average per substrate, 27.2 % of all non-  
640 spliced peptides were produced by 20S proteasomes and detected by MS. In contrast, 0.16 %  
641 of all theoretically *cis*-spliced and 0.06 % of all theoretically *trans*-spliced peptides were  
642 produced by 20S proteasomes and detected by invitroSPI (Fig. 6).

643 Among the unique peptides per substrate reported in the new ProteasomeDB, 1,031 non-  
644 spliced, 2,549 *cis*-spliced and 1,517 *trans*-spliced peptides were not reported in Specht and  
645 Paes databases of peptide products. In addition, 4,462 non-spliced, 3,904 *cis*-spliced and 3,168  
646 *trans*-spliced peptides originally reported in Specht and Paes databases of peptide products  
647 were confirmed by the application of invitroSPI to the cognate datasets, bearing in mind that  
648 invitroPB could not identify *trans*-spliced peptides, which therefore were not detected in the  
649 original Paes databases of peptide products (Fig. S7a-c). To note, in this study we reported the  
650 number of unique peptides per substrate. Therefore, since Specht and PB datasets had no  
651 common substrates, they also had no common unique peptides per substrate. We adopted this  
652 strategy because we speculated that further analysis and eventual prediction of proteasome-  
653 catalyzed peptide hydrolysis and peptide splicing would, in most cases, consider the peptide  
654 sequence as well as its substrate origin. ProteasomeDB structure, however, also allows to  
655 obtain a list of unique peptide sequences regardless of their substrate origin, depending on the  
656 user's choices and analysis goals.

657  
658 **Illustrative analysis of the ProteasomeDB and the whole dataset: focus on 20S standard**  
659 **proteasome and early time point digestions**

660 So far, we selected and compared subsets of the whole dataset, as well as the outcome of  
661 different identification methods. ProteasomeDB (generated through the application of invitroSPI  
662 on the whole dataset) could, however, be large enough to carry out analyses on the catalytic  
663 nature of proteasome-catalysed peptide splicing and hydrolysis. As proof of principle, we here  
664 analyzed *in vitro* digestions carried out for 2/4 h with 20S standard proteasomes and their  
665 corresponding controls. The analysis of these time points, for instance, could minimize the  
666 peptide product re-entry events in proteasomes; and the focus on 20S standard proteasome  
667 digestion could be a strategy to limit the variance due to the different dynamics of proteasome  
668 isoforms<sup>33</sup>. In addition, in this illustrative analysis, we compared the features of the unique

669 peptides either identified by applying invitroSPI and invitroPB methods to the 2h PB dataset (24  
670 substrates), or by applying invitroSPI to the 4h Specht dataset (47 substrates) and the whole  
671 2/4h dataset of 71 substrates (white inlets in **Fig. 1e**). A comparison of invitroSPI with invitroPB  
672 on the whole 2/4h dataset of 71 substrates could not be carried out, because invitroPB required  
673 high-precision MS data due to its dependence on de novo peptide sequencing, and many  
674 substrate digestions present in Specht dataset were measured by MS instruments with lower  
675 precision (**Fig. 1e**).

676 In this illustrative analysis, by applying invitroSPI to the PB dataset, we identified more unique  
677 *cis*-spliced (and of course *trans*-spliced) peptides than invitroPB, both considering the total  
678 number of unique peptides (**Table 3**) and the relative frequency of peptides per substrate (**Fig.**  
679 **7a**). By applying invitroSPI to all three investigated datasets, we identified *cis*-spliced and *trans*-  
680 spliced peptides with a similar frequency (**Table 3, Fig. 7a**). InvitroSPI identified a sizeable  
681 portion of non-spliced and spliced peptides that contained the N- or C-termini of the substrates  
682 (**Fig. 7a-c**). These peptides were excluded in the analysis carried out by Paes *et al.*<sup>57</sup>, with  
683 consequences discussed below (**Fig. 7a-b**). InvitroSPI also identified a sizeable portion of  
684 spliced peptides with a one amino acid long splice-reactant, which could not be identified by  
685 invitroPB (**Fig. 7a**). Through the application of invitroSPI to the PB dataset, we did not observe  
686 a narrower length distribution of *cis*-spliced compared to non-spliced peptides (**Fig. 7b**), which  
687 was described by Paes *et al.*<sup>57</sup>. In all datasets analyzed by invitroSPI, non-spliced peptides  
688 were, on average, shorter than *cis*-spliced peptides, in contrast to what was described by Paes  
689 *et al.*<sup>57</sup>. Furthermore, in the whole dataset analyzed by invitroSPI, *cis*-spliced peptides were  
690 shorter than *trans*-spliced peptides (**Fig. 7b**), in agreement with what was previously described  
691<sup>8</sup>. Because of multi-mapper spliced peptides and the features of the simulated background  
692 databases (see below), we avoided a more in-depth analysis of spliced peptide features.  
693 However, since Paes *et al.*<sup>57</sup> suggested that the length of the N- and C-terminal splice-  
694 reactants of *cis*-spliced peptides differed, we preliminary investigated this aspect, focusing only  
695 on *cis*-spliced peptides that could be unequivocally assigned to a unique splice-reactant length.  
696 Although both methods identified *cis*-spliced peptides with, on average, shorter N-terminal  
697 spliced-reactants than the C-terminal ones in the PB dataset, this phenomenon was not  
698 confirmed in the larger Specht dataset and in the whole dataset. Indeed, in these two largest  
699 datasets analyzed through invitroSPI, N- and C-terminal splice-reactants of *cis*-spliced peptides  
700 had a similar length distribution (**Fig. 7c**). As discussed below, however, for an unbiased  
701 analysis, all biochemical characteristics of peptide product types should be compared to a  
702 simulated background database, to identify features that are specific for peptide hydrolysis and  
703 peptide splicing reactions.

704

#### 705 **Potential pitfalls in data analysis: overview.**

706 For an appropriate investigation of sequence motifs and features of non-spliced and spliced  
707 peptides produced by proteasomes during the degradation of synthetic polypeptides, some  
708 factors play, in our opinion, a pivotal role: (i) the amino acid frequency should be normalized  
709 against an appropriate simulated background database to account for biases in the substrate  
710 amino acid composition; (ii) the database of identified peptides and digested substrates should  
711 be large enough to account for the large number of possible amino acid combinations; (iii-vi)  
712 non-spliced and spliced peptide identification algorithms could bias the features of the identified  
713 peptide pools, and, hence, methodological limitations should be considered during the analysis;  
714 (vii) for many spliced peptides, multiple splice-reactant locations are possible (multi-mapper  
715 peptides), thereby impinging upon the confidence in the computation of the features of splice-  
716 reactants, intervening sequences and PCPS splice-sites.

717

#### 718 **Potential pitfalls in data analysis: (i) normalization strategy.**

719 One potential use case of ProteasomeDB is the analysis of amino acid preferences at the splice  
720 sites, *i.e.* sP1 and sP1' (the two amino acid residues that are ligated together during PCPS; see

721 **Fig. 1b-d).** In such analysis one should carefully consider the expected amino acid frequencies  
722 in sP1 and sP1' observed by chance due to the limited sequence variety and amino acid  
723 composition of the substrates studied.

724 To this end, we computed the joint amino acid frequencies at sP1 and sP1' based on the  
725 theoretical possible spliced peptides that could be derived from all studied substrates (simulated  
726 background database). The resulting frequency matrix represented the splice site background  
727 distribution (**Fig. 8a**), which in part reflected the natural amino acid frequency in the studied  
728 substrates. This non-uniform background distribution must be considered when analyzing *in*  
729 *vitro* digested spliced and non-spliced peptide products generated from polypeptides, especially  
730 when dealing with a small peptide product database with limited sequence diversity. Therefore,  
731 we suggest that all observed amino acid frequencies have to be normalized by their respective  
732 frequency in a simulated background database, and not only by amino acid frequencies  
733 occurring in the substrate sequences as done by others<sup>57</sup>. An example of the use of the  
734 simulated background databases for normalization is illustrated in the following section.

735

### 736 **Potential pitfalls in data analysis: (ii) peptide product database size.**

737 The second factor in our list of potential pitfalls refers to the peptide product database size. To  
738 investigate the impact of the peptide product database size on the statistical analysis of PCPS  
739 features, we compared the amino acid frequency at sP<sub>1</sub> and sP<sub>1</sub>' sites between the peptide  
740 sequences originally published by Paes *et al.*<sup>57</sup> and ProteasomeDB. In both peptide product  
741 databases, the obtained amino acid frequencies were normalized by the frequencies in the  
742 respective simulated background database (discussed in the section above). This was done to  
743 account for potential biases introduced both through natural variation of amino acid frequency  
744 and substrate composition (see above). Paes *et al.* compared the splice-site signature of  
745 forward *cis*- and reverse *cis*-spliced peptides based on 130 *cis*-spliced peptides included in their  
746 analysis of 2 h *in vitro* degradation of 23 synthetic substrates. They concluded that forward and  
747 reverse *cis*-PCPS had a different preference for amino acids in sP1 and sP1'<sup>57</sup>. Their analysis  
748 was based on 63 forward and 67 reverse *cis*-spliced peptides from 15 substrates in the original  
749 Paes' peptide product database, since they did not identify *cis*-spliced peptides from 8 synthetic  
750 substrates<sup>57</sup>. The corresponding subset of the ProteasomeDB - restricted to 2/4 h *in vitro*  
751 degradation of 71 synthetic substrates with 20S standard proteasomes - included 1,674 forward  
752 and 1,080 reverse *cis*-spliced peptide products.

753 We repeatedly sampled a subset of peptide sequences (*i.e.*, applied 200 bootstrapping  
754 iterations on 80% of the data) and calculated the normalized amino acid frequency in each  
755 sampling iteration. In general, the 90% confidence interval of all bootstrap iterations results in  
756 an estimation of both the amino acid frequency at the sP1 and sP1', and the robustness of this  
757 estimation. Accordingly, large confidence intervals indicate low reliability of the obtained amino  
758 acid frequencies.

759 The confidence intervals of the original Paes' peptide product database<sup>57</sup> were always larger  
760 than the ProteasomeDB subset (*e.g.*, see A, C, H, Q amino acids in sP<sub>1</sub>' of reverse *cis*-spliced  
761 peptides; **Fig. 8b**). For many amino acids, the original Paes' peptide product database showed  
762 almost zero frequency at sP<sub>1</sub> and sP<sub>1</sub>', which may suggest that these amino acids were not  
763 used by proteasomes as splice-sites. This was not confirmed on the ProteasomeDB subset  
764 (*e.g.*, see N, S, T, V amino acids in sP<sub>1</sub>' of reverse *cis*-spliced peptides; **Fig. 8b**). At last, for  
765 most of the amino acids, the normalized frequency computed from the original Paes' peptide  
766 product database<sup>57</sup> and the ProteasomeDB subset did not match (**Fig. 8b**). All these analyses  
767 point toward the risk of overinterpretation of results obtained from small spliced peptide product  
768 databases, which may also explain the different results obtained in the three datasets shown in  
769 **Fig. 7** and **Fig. 8b**.

770 From this preliminary analysis, we observed pools of amino acids that were either favored or  
771 disfavored as sP1 and sP1', thereby suggesting that PCPS has peptide sequence preferences.  
772 Future studies, perhaps using grouping strategies based on chemophysical features of amino



773 acids, could use ProteasomeDB to decipher the peptide sequence preferences of both peptide  
774 hydrolysis and peptide splicing catalyzed by 20S proteasomes.

775

776 **Potential pitfalls in data analysis: (iii) synthesis errors.**

777 The third factor in our list of potential pitfalls refers to a confounding element in this type of  
778 sample, *i.e.* the presence of synthesis errors and their elimination. Both invitroSPI and invitroPB  
779 developed (different) strategies for synthesis error removal (**Fig. S1, S2**). Between the two  
780 methods, invitroSPI has likely the most stringent strategy to eliminate synthesis errors, which  
781 might have been assigned as spliced peptides. Indeed, it discards not only peptides identified  
782 as such in the control sample (either 0 h digestions or samples with synthetic substrates and  
783 no proteasomes) but also any putative spliced peptide with the same splice-site as a synthesis  
784 error identified in the control samples (see **Fig. 1f, Fig. 3a, Fig. 4a**). The latter step, which is not  
785 present in the invitroPB pipeline, may result in the elimination of spliced peptides that are  
786 produced by proteasomes although a peptide with the same splice-site is also present as  
787 synthesis error. In addition, both methods do not eliminate spliced peptides that have, as C-  
788 terminal splice-reactants, non-spliced peptides, which are present in the control samples and,  
789 thus, are assigned as synthesis errors. For example, we computed that a fraction of spliced  
790 peptides (median 6.4 % per substrate) contained a potential non-spliced synthesis error as their  
791 C-terminal splice-reactant in ProteasomeDB. In the same peptide product database, the fraction  
792 of spliced peptides that contain a non-spliced peptide – that is not a synthesis error – as their  
793 C-terminal splice-reactant is 36.5 %. The former circumstance might be considered for further  
794 analysis, depending on the user's choice and analysis goal. Indeed, for these cases, no  
795 unequivocal statement about the origin of the C-terminal splice-reactant can be made (they  
796 could be generated as non-spliced peptides by proteasomes even if they are also present as  
797 synthesis errors), and, regardless of the splice-reactant's origin, all these splice-reactants  
798 underwent PCPS to generate a splice-peptide. Therefore, the splice-reactants matched the  
799 catalytic requirements of PCPS in terms of sequence length and composition, and, hence, could  
800 be included in the downstream analysis depending on the analysis objective.

801 To this end, we have a feature in ProteasomeDB, which denotes whether a spliced peptide  
802 potentially contains a C-terminal splice-reactant that could be a synthesis error (**Tab. 1**). This  
803 information is limited to C-terminal splice-reactants because, according to the transpeptidation  
804 model of PCPS the N-terminal splice-reactant needs to be first cleaved by proteasomes to form  
805 the acyl-enzyme intermediate, and, thus, it cannot be a synthesis error. There are only few  
806 examples of PCPS via other reaction mechanisms such as condensation<sup>20,25</sup>.

807

808 **Potential pitfalls and missed opportunities in data analysis: (iv) restriction in peptide and**  
809 **splice-reactant length.**

810 The fourth factor in our list of potential pitfalls refers to the restriction of features of spliced and  
811 non-spliced peptides that can be identified. Both invitroSPI and invitroPB developed different  
812 strategies for non-spliced and spliced peptide identification, which can impinge upon the pool  
813 of identified peptide products. For example, Paes *et al.*<sup>57</sup> restricted the identification to *cis*-  
814 spliced peptides longer than 7 amino acids, whereas non-spliced peptides had no restriction of  
815 this kind. This could explain why Paes *et al.*<sup>57</sup> described a narrower length distribution for *cis*-  
816 spliced peptides than non-spliced peptides. In contrast, invitroSPI, which applies the same  
817 identification strategy for non-spliced and spliced peptides regarding their length restrictions for  
818 spliced and non-spliced peptides (5 residues or longer; see **Fig. 2**), did not confirm the result of  
819 Paes *et al.*<sup>57</sup>. In fact, the analysis of the whole dataset using invitroSPI showed that  
820 proteasome-generated *cis*-spliced peptides, and *trans*-spliced peptides, are on average longer  
821 than non-spliced peptides in the ProteasomeDB (**Fig. 7b**), as previously shown in the Specht  
822 database of peptide products<sup>8</sup>.

823 InvitroPB also forbids the identification of spliced peptides with a splice-reactant length of one  
824 amino acid (**Fig. 2**). This strategy was in part based on a single example of a *cis*-spliced epitope

825 previously described by Michaux *et al.*<sup>51</sup> (*i.e.*, gp100<sub>195-202/192</sub>). It has been demonstrated *in vitro*  
826 and *in cellula* that this specific *cis*-spliced epitope is not spliced as such, but as a C-terminal  
827 extended precursor, with a splice-reactant that is three amino acids long. However, upon PCPS,  
828 the spliced epitope precursor can be further processed by proteasomes, thereby generating the  
829 *cis*-spliced epitope that is recognized by CD8<sup>+</sup> T cells of melanoma patients<sup>20,51</sup>. In contrast to  
830 invitroPB, invitroSPI identifies *cis*-spliced peptides with a one amino acid-long splice-reactant  
831 as final products (**Fig. 2**), since they could be the result of an initial PCPS event followed by  
832 peptide hydrolysis. In our analysis, these *cis*-spliced peptides had a comparable MS2 spectrum  
833 quality as the other identified peptides (**Fig. 4g**), thereby supporting their reliable identification.  
834 It is also worth noting that in *in vitro* digestions of synthetic polypeptides with proteasomes, MS  
835 measurements and the downstream analysis identify the final products of the PCPS reaction  
836 rather than intermediate products, and, hence, we can never exclude that an identified peptide  
837 is the outcome of multiple peptide catalytic reaction. As a consequence, it is non-trivial to draw  
838 conclusions about the minimal length of splice-reactants in such datasets of proteasome-  
839 generated spliced peptides. Spliced peptides can re-bind to the proteasome's active site and  
840 be cleaved, thus resulting in shorter splice-reactants than in the original transpeptidation  
841 reaction. By including spliced peptides with a splice-reactant length of one amino acid in  
842 ProteasomeDB, we could, however, carry out a simple analysis to understand if the pioneering  
843 observation of Michaux *et al.*<sup>51</sup> could be generalized. To this end, we investigated the length of  
844 (i) N-terminal splice-reactants of forward *cis*-spliced peptides that originate from the substrate's  
845 N-terminus and (ii) C-terminal splice-reactants of forward *cis*-spliced peptides that originate  
846 from the substrate's C-terminus. The frequency of short fragments as splice-reactants that were  
847 located at the substrate's termini could allow conclusions to be drawn about the minimal splice-  
848 reactant lengths required for PCPS since these splice-reactants could not be derived from  
849 trimming of longer fragments (**Fig. S8a**). In ProteasomeDB, we identified many forward *cis*-  
850 spliced peptides with N-terminal splice-reactants located at the substrate's N-terminus, among  
851 which around 4.5 % were one amino acid long and 9.2 % were two amino acids long (**Fig. S8b**).  
852 This analysis suggested that N-terminal splice-reactants of one amino acid length could be  
853 efficiently used as such for PCPS. On the contrary, forward *cis*-spliced peptides with a one  
854 amino acid long C-terminal splice-reactant located at the substrate's C-terminus were identified  
855 far less frequent (**Fig. S8b**), thereby suggesting that the C-terminal splice-reactants of at least  
856 2 amino acids length were required for an efficient PCPS. Overall, this result confirmed the  
857 initial observation of Michaux *et al.*<sup>51</sup>, which was limited to C-terminal splice-reactants, although  
858 exceptions have been reported in ProteasomeDB. Furthermore, the relative frequency of  
859 spliced peptide products with a one amino acid long splice-reactant seemed to be smaller in 2/4  
860 h vs 20/24 h digestion experiments (**Fig. S8c**). Similarly, the splice-reactant length distribution  
861 appeared narrower at later digestion time points compared to earlier time points, although not  
862 statistically significant (**Fig. S8d**). These data could be due a re-entry of spliced peptides  
863 followed by peptide hydrolysis in the late time point of the reactions, as well a change in  
864 proteasome dynamics over time as shown in other experimental set up<sup>33</sup>.

865  
866 **Potential pitfalls and missed opportunities in data analysis: (v) restriction in peptide**  
867 **identification based on their location within the substrate.**

868 InvitroSPI and invitroPB differed in another aspect of the peptide product identification strategy,  
869 which impinged upon the features of the identified peptide pool. InvitroSPI allowed the  
870 identification of non-spliced and spliced peptides carrying N- or C-termini of synthetic  
871 polypeptide substrates. These peptides were also identified by InvitroPB, which, however,  
872 excluded them in the downstream analysis.

873 To understand if their exclusion could bias the analysis of the identified peptide products, we  
874 computed *in silico* the theoretical fraction of non-spliced and *cis*-spliced peptides carrying the  
875 N- or C-termini of the substrates of the whole dataset, by applying the peptide length restrictions  
876 applied by invitroSPI and invitroPB. The fraction of theoretical peptides carrying the substrate's

877 N- or C-termini strongly depended on the substrate length (**Fig. S9**). This theoretically expected  
878 frequency was not confirmed among the fractions of experimentally identified peptides which  
879 carried the substrate's N- or C-termini, analyzing the PB dataset and the whole dataset by  
880 applying the two identification methods (invitroSPI and invitroPB), respectively (**Fig. S9** and **Fig.**  
881 **7a**). We observed that while the fraction of identified non-spliced peptides that carry either of  
882 the substrate's termini laid within the theoretically expected range, the fraction of spliced  
883 peptides with this property was much higher than expected by chance (**Fig. S9**). The similarity  
884 between measured and predicted MS2 spectra of spliced peptides with or without substrate's  
885 N- or C-termini did not differ among the peptides identified in the PB dataset (**Fig. 4g**), hence  
886 suggesting that their identification was equally reliable. Therefore, by removing spliced and non-  
887 spliced peptides carrying the substrate's N- or C-termini, one would not only remove a large  
888 portion of peptides produced by proteasomes *in vitro*, but also introduce a bias in the analysis  
889 by artificially constraining the spliced peptide pool. Furthermore, and in line with our  
890 observations, there is preliminary evidence of preferential processing of protein termini by  
891 proteasomes in living cells. Indeed, a larger frequency of non-spliced peptides produced by  
892 proteasomes by peptide hydrolysis of the termini of proteins compared to their central area has  
893 been shown *in cellula* by the pioneering work of Wolf-Levy and colleagues <sup>71</sup>. It would be  
894 worthwhile to verify in the same kind of samples, *i.e.*, peptides eluted from proteasomes *in*  
895 *cellula*, if this holds also true for spliced peptides.

#### 896 897 **Potential pitfalls in data analysis: (vi) non-spliced peptides with PTMs.**

898 Another example of a different strategy between invitroSPI and invitroPB, which could have an  
899 impact on the features of the identified peptide pool, is related to chemical PTMs. InvitroSPI  
900 allowed the identification of both non-spliced and spliced peptides carrying three chemical  
901 modifications (see 'Technical aspects of invitroSPI' chapter), and thus treated non-spliced and  
902 spliced peptides equally for this aspect. In contrast, invitroPB filtered out PTM-labelled non-  
903 spliced peptides, and introduced a specific filter only for *cis*-spliced peptides. We believe that  
904 the exclusion of PTM-labelled non-spliced peptides in the final list of identified peptides was a  
905 specific strategy adopted by Paes *et al.* <sup>57</sup> for the comparison of non-spliced and *cis*-spliced  
906 peptides in that study, and, thus, could potentially be omitted in future applications of invitroPB.  
907 Conversely, PTMs had a key role in the identification of *cis*-spliced peptides by invitroPB: the  
908 method excluded MS2 spectra potentially assigned to *cis*-spliced peptides if they might have  
909 been non-spliced peptides tagged with any of the 313 PTMs considered by PEAKS-PTM. The  
910 original objective of Paes *et al.* <sup>57</sup> was to reduce the risk of miss-assignment of MS2 spectra to  
911 *cis*-spliced peptides. This step of invitroPB, which was embedded in the method pipeline, may  
912 have achieved the original objective, although it may have also resulted in a reduced recall of  
913 *cis*-spliced peptides. As indicated in **Fig. 3b** and **Fig. 4b-c**, invitroSPI assigned several PSMs  
914 to spliced peptide sequences, which were dismissed by invitroPB because they were identified  
915 as non-spliced peptides with PTMs by PEAKS-PTM. The competition of different peptide  
916 sequences for the assignment of a MS2 spectrum in the presence of a large search space is  
917 an issue that has been addressed with various strategies <sup>72</sup> and benchmarking approaches <sup>44,73</sup>.  
918 It is worth noting that the sequence search space of PEAKS-PTM, which considers 313 PTMs  
919 (maximum two PTMs allowed per peptide), may be even larger than the spliced peptide  
920 sequence search space, and thus be tangled to similar statistical issues. Therefore, the *a priori*  
921 exclusion of MS2 spectra for spliced peptide identification because they might be non-spliced  
922 peptides with unlikely PTMs (**Fig. 3c**) may not be directly supported by statistical  
923 considerations. In our opinion, any PTM-modified peptide assigned by PEAKS-PTM should be  
924 revisited to understand if it could occur in the specific experimental context. In the present study,  
925 technical modifications such as formylations could be explained through the use of formic acid  
926 in the MS buffer. On the contrary, biological modifications such as phosphorylation,  
927 ubiquitination and others, although suggested by PEAKS-PTM, are most likely false positive  
928 assignments (**Fig. 3c**). Nevertheless, an interesting avenue of further research could be to

929 investigate to what extent PTMs occur before or after the splicing/hydrolysis reaction, and to  
930 what extent they influence the reaction towards either splicing or hydrolysis.

931

### 932 **Potential pitfalls in data analysis: (vii) multi-mapper peptides.**

933 The last issue that we would like to mention is the presence of peptide sequences that may  
934 have different locations within the substrate sequence, *i.e.*, multi-mapper peptides. In invitroSPI,  
935 we imposed a hierarchical strategy, which gives preference to non-spliced over spliced  
936 peptides, and *cis*-spliced over *trans*-spliced peptides (see Methods section). Nonetheless,  
937 many *cis*-spliced peptide sequences may be both forward and reverse *cis*-spliced peptides;  
938 many spliced peptide sequences may be spliced peptides with different splice-reactant lengths,  
939 and hence different splice-sites. This issue has not been considered by previous studies on  
940 both *in vitro* digestions of synthetic polypeptides by proteasomes<sup>8,57</sup>, and HLA-I  
941 immunopeptidomes<sup>10-14</sup>. These studies adopted simple random assignment strategies, which  
942 may lead to artefacts. We think that more elaborated biochemical approaches should be used  
943 to better define the origin of these multi-mapper spliced peptides to avoid bias in the  
944 development of PCPS predictors. ProteasomeDB could be a cornerstone of such studies.

945

946

### 947 **Usage Notes**

948 The whole peptide product database - ProteasomeDB - is provided as CSV file, which can be  
949 opened in Excel or any text editor.

950

### 951 **Code Availability**

952 The algorithm generating all possible *cis* and *trans* spliced peptides was originally described by  
953 Liepe *et al.*<sup>63</sup>.

954 InvitroSPI method has been implemented with Snakemake in the Conda environment and is  
955 available at GitHub (<https://github.com/QuantSysBio/invitroSPI>).

956 The analysis scripts (written in R) and implementation of invitroPB are available on Figshare  
957 online repository<sup>70</sup>.

958 Analyses were carried out in R v4.1.1.

959 Figures have been generated in R and postprocessing was done with Adobe Illustrator v25.2.3.

960 The new *in vitro* TSN2 and TSN89 digestion samples were measured on Fusion Lumos  
961 Orbitrap, and acquired using Xcalibur v4.4.

962

963

### 964 **Acknowledgements**

965 We thank the CRUK-KHP Cancer Centre and the CEMS of KCL for the MS measurements, A.  
966 Mansurkhodzhaev and W.T. Soh (MPI-NAT) for PRIDE's file uploading, W. Paes and P. Borrow  
967 (Oxford) to provide technical information related to their original paper.

968 This work was financed in part by: (i) Cancer Research UK [C67500/A29686] and National  
969 Institute for Health Research (NIHR) Biomedical Research Centre based at Guy's as well as St  
970 Thomas' NHS Foundation Trust and King's College London and/or the NIHR Clinical Research  
971 Facility to MM. (ii) ERC-StG 945528 IMAP to JL. HPR was funded by the Manfred Eigen-  
972 Förderstiftung (Principles of Cancer Research - Stipend for an exceptional, independently  
973 working young scientist), and by King's College London as part of the "Neuro-Immune  
974 Interactions in Health & Disease Wellcome Trust PhD Programme". JAC is supported by the  
975 International Max-Planck Research School (IMPRS) for Genome Science

976

### 977 **Author contributions**

978 HPR, MM and JL developed the project, performed and/or supervised the data analysis and  
979 data generation and wrote the manuscript. GRH performed the *in vitro* digestions. JAC

980 performed the Prosit analysis and proofread the manuscript. XY and SL optimized the MS  
981 method, measured the samples via MS and edited the manuscript.

982

983 **Competing interests**

984 The authors have no conflicts of interest.

985

986

987  
988

**Figures and Tables**

<b>Column name</b>	<b>Description</b>
<b>sampleID</b>	Unique identifier for every sample
<b>sampleName</b>	Sample Name used during experiment
<b>filename</b>	Mascot search result file name (available on PRIDE)
<b>runID</b>	Technical replicate number
<b>protIsoType</b>	Proteasome isoform used for digestion
<b>digestTime</b>	Elapsed digestion time (hours) at time of measurement
<b>proteasomeSpecies</b>	Species origin of used proteasomes
<b>sampleDate</b>	Sample date
<b>instrument</b>	Instrument used for measurement
<b>fragmentation</b>	Fragmentation method used for measurement
<b>location</b>	Measurement location
<b>substrateSeq</b>	Amino acid sequence of substrate
<b>substrateOrigin</b>	Protein origin of substrate
<b>substrateSpecies</b>	Species origin of substrate
<b>substrateID</b>	Unique identifier for a substrate sequence
<b>pepSeq</b>	Amino acid sequence of peptide products
<b>scanNum</b>	Scan number listed in the RAW file
<b>rank</b>	Peptide rank assigned by Mascot Server
<b>ionScore</b>	Ion score assigned by Mascot Server
<b>qValue</b>	q-value assigned by Mascot Server
<b>productType</b>	PCP: non-spliced peptide; PSP: spliced peptide
<b>spliceType</b>	cis: forward <i>cis</i> -spliced peptide; revCis: reverse <i>cis</i> -spliced peptide; trans: <i>trans</i> -spliced peptide; N/A: non-spliced peptide
<b>positions</b>	Location(s) of the peptide sequence in the synthetic polypeptide substrate
<b>synErrSR2</b>	Indication whether the C-terminal splice-reactant of a spliced peptide matches a non-spliced synthesis error; N/A: non-spliced peptide
<b>charge</b>	Ion charge
<b>PTM</b>	Post-translational modifications

989  
990  
991  
992  
993  
994

**Table 1. ProteasomeDB database description.** Listed are the column names (attributes) in the ProteasomeDB database and their corresponding explanations.

	Datasets analyzed by invitroSPI				Datasets analyzed by invitroPB	
	gp100 Fusion	PB	Specht	whole	gp100 Fusion	PB
<b>Peptide types:</b>						
<b>Non-spliced</b>	68	1,196	4,288	5,493	46	1,185
<b>Cis-spliced</b>	275	1,403	4,915	6,453	185	1,060
<b>Trans-spliced</b>	115	814	3,781	4,685	0	0
<b>Forward cis-spliced</b>	96	838	2,828	3,716	71	701
<b>Reverse cis-spliced</b>	174	476	1,876	2,435	101	316
<b>Forward/reverse cis-spliced (multi-mapper)</b>	5	89	211	302	13	43
<b>Spliced with 1 amino acid splice-reactant</b>	67	409	943	1,390	0	0
<b>Non-spliced with N- or C-terminal residues</b>	11	298	634	932	1	104
<b>spliced with N- or C-terminal residues</b>	204	1,530	5,247	6,876	89	653

995  
996 **Table 2. Number of unique peptides identified in the various datasets by applying**  
997 **different identification methods.** Number of unique peptides identified through the application  
998 of invitroSPI and invitroPB to the PB, Specht and whole datasets. In this table, all substrates,  
999 all proteasome types, and time points, have been included.  
1000  
1001

1002

	Datasets analyzed by invitroSPI				Datasets analyzed by invitroPB	
	gp100 Fusion	PB	Specht	whole	gp100 Fusion	PB
<b>Peptide types:</b>						
<b>Non-spliced</b>	54	823	2,996	3,837	40	864
<b><i>Cis</i>-spliced</b>	171	759	2,363	3,240	101	617
<b><i>Trans</i>-spliced</b>	77	410	2,058	2,536	0	0
<b>Forward <i>cis</i>-spliced</b>	70	451	1,400	1903	43	404
<b>Reverse <i>cis</i>-spliced</b>	100	258	868	1,191	53	186
<b>Forward/reverse <i>cis</i>-spliced (multi-mapper)</b>	1	50	95	146	5	27
<b>Spliced with 1 amino acid splice-reactant</b>	37	232	456	711	0	0
<b>Non-spliced with N- or C-terminal residues</b>	10	224	514	740	1	88
<b>spliced with N- or C-terminal residues</b>	148	839	2,859	3,800	54	383

1003

1004

1005

1006

1007

1008

1009

**Table 3. Number of unique peptides identified by applying different identification methods and focusing on 2/4 h digestions with 20S standard proteasomes.** Number of unique peptides identified through the application of invitroSPI and invitroPB to the PB, Specht and the whole datasets. In this table only substrates digested with 20S standard proteasomes for 2/4 h have been included.



## 1010 Figure Legends

1011

1012 **Figure 1. Proteasome-generated non-spliced and spliced peptides, and overview of**  
1013 **method and dataset application.** Proteasomes form: (a) non-spliced peptides via peptide  
1014 hydrolysis, (b-d) spliced peptides through ligation of two non-contiguous splice-reactants either  
1015 derived from the same polypeptide molecule (*cis*-spliced peptides, b,c) or from two distinct  
1016 molecules of the same protein or two distinct proteins (*trans*-spliced peptides, d). In b-c, peptide  
1017 fragment ligation can occur in forward order, *i.e.*, following the orientation from N- to C-terminus  
1018 of the parental protein (forward *cis*-peptide splicing; b), or in reverse order (reverse *cis*-peptide  
1019 splicing; c). The two ligated fragments are named splice-reactants, and their junction is named  
1020 splice-site. The C-terminus of the first (N-terminal) splice-reactant is named sP<sub>1</sub>, whilst the N-  
1021 terminus of the second (C-terminal) splice-reactant is named sP<sub>1</sub>'. The sequence segment  
1022 between two splice-reactants is called the intervening sequence. Arrows represent the  
1023 substrate cleavage sites used by proteasome catalytic Thr1. (e) Overview of methods and  
1024 datasets described in this study. (f) Substrate synthesis errors. Various forms of synthesis errors  
1025 could result in alleged non-spliced and/or spliced peptides. Those synthesis errors are captured  
1026 using control measurements. Furthermore, alleged spliced synthesis errors can be trimmed by  
1027 the proteasome. All such spliced peptides of which a precursor is identified in control  
1028 measurements are removed by invitroSPI but not by invitroPB.

1029

1030 **Figure 2. Difference in the peptide identification strategy and downstream analysis**  
1031 **adopted by invitroSPI and invitroPB.**

1032

1033 **Figure 3. Comparison of invitroSPI and invitroPB methods applied to the gp100 Fusion**  
1034 **dataset. a-e)** Number of PSMs assigned to: (a) non-spliced, *cis*-spliced, *trans*-spliced peptides,  
1035 and related synthesis error peptides, (b) PTM-labelled peptides, (c) forward and reverse *cis*-  
1036 spliced peptides, (d) spliced peptides with one amino acid long splice-reactant, (e) spliced  
1037 peptides containing substrate's N- or C-termini. Assignment was carried out by applying  
1038 invitroSPI and invitroPB methods to *in vitro* digestions of TSN2 and TSN89 substrates with  
1039 proteasomes. PTM-modified non-spliced peptides identified by PEAKS-PTM are reported,  
1040 although they are not kept in the final list of identified peptides by invitroPB. In invitroSPI-  
1041 identifications, PTM-modified peptides are included. In (b-e), PSMs assigned to synthesis errors  
1042 have been removed. In (c), forward/reverse *cis*-spliced peptides, *i.e.* multi-mapping *cis*-spliced  
1043 peptides, are not shown. f,g) MS2 spectra of the *cis*-spliced epitopes (f) [RTK][QLYPEW] and  
1044 (g) [QLYPEW][RTK] identified in *in vitro* digestions of (f) TSN89 and (g) TSN2 substrates, and  
1045 of their cognate synthetic peptides. Detected *m/z* and charges in the MS2 spectra shared  
1046 between *in vitro* digestion samples and synthetic peptides are indicated in red. Other assigned  
1047 *m/z* are indicated in blue. In MS2 spectra, charged b-, a- and y-ions are reported. Double  
1048 charged ions are marked as \*\*. Ions' neutral loss of ammonia is symbolized by \*. Extracted ion  
1049 chromatograms of target peptides in *in vitro* digestion and synthetic peptides are plotted in the  
1050 right panels and indicate matching retention times and absence of a biologically meaningful  
1051 peak in the 0 h digestion. MS ion chromatograms correspond to the *m/z* = 610.80-610.84 (+2;  
1052 f) and 407.53-407.57 (+3; g). h) number of unique peptide sequences identified by invitroSPI  
1053 in the gp100 Fusion dataset shown for 2h, 4h and 20h. i) frequency of spliced and non-spliced  
1054 peptides over time identified by invitroSPI in the gp100 Fusion dataset comprising two  
1055 substrates.

1056 In (a-e,h-i) *in vitro* digestion samples (0, 2, 4, 20 h) and cognate synthetic peptides were  
1057 measured by Orbitrap Fusion Lumos (KCL-CEMS) by using the same MS method. For MS2  
1058 spectrum references, (f): file 20210422\_WB2\_2h\_TSN89\_FusionCEMS, charge +2, scan 5897  
1059 (upper panel); file 20210422\_GP100\_mix\_FusionCEMS, charge +2, scan 5208 (lower panel).  
1060 (g): file 20210422\_WA4\_20h\_TSN2\_FusionCEMS, charge +3, scan 6115 (upper panel); file  
1061 20210422\_GP100\_mix\_FusionCEMS, charge +3, scan 4936 (lower panel).

1062

1063

1064 **Figure 4. Comparison of invitroSPI and invitroPB methods applied to the PB dataset. a,b)**  
1065 Number of PSMs assigned to: (a) non-spliced, *cis*-spliced, *trans*-spliced peptides, and either  
1066 related synthesis error peptides, or (b) PTM-labelled peptides. c) Frequency of PTMs among  
1067 PTM-labelled non-spliced peptides suggested by PEAKS-PTM as part of invitroPB. d-f) Number  
1068 of PSMs assigned to: (d) forward and reverse *cis*-spliced peptides (multi-mapper  
1069 forward/reverse *cis*-spliced peptides are not shown), (e) spliced peptides with one amino acid  
1070 long splice-reactant, and (f) spliced peptides containing substrate's N- or C-termini. Assignment  
1071 was carried out by applying invitroSPI and invitroPB methods to the PB dataset. In invitroSPI-  
1072 identified peptides, also PTM-modified peptides are included. In (b and d-f), PSMs assigned to  
1073 synthesis errors have been removed. g) Spectral angle distribution computed between  
1074 measured and predicted MS2 spectra identified by invitroSPI (red) and invitroPB methods  
1075 (grey). Only PSMs of unmodified non-spliced and spliced peptide that do not contain any  
1076 cysteine (C) residues, do not exceed a charge of 6 and are 7-12 amino acid long are here  
1077 included, since Prosit cannot predict PTM-modified peptide's MS2 spectra and Prosit  
1078 performance is influenced by peptide length (Fig. S3). In the violin plots, horizontal black lines  
1079 represent the median. The number of PSMs for each group is reported. In (a-g), *in vitro*  
1080 digestion samples (2 h and 20 h digestions with proteasomes and 20 h without proteasomes)  
1081 were measured by Orbitrap Fusion Lumos (Oxford proteomics centre).

1082

1083

1084 **Figure 5. FDR estimation for invitroSPI and invitroPB in PB dataset. a,b)** Spectral angle  
1085 distribution of non-spliced, *cis*-spliced and *trans*-spliced peptide identified by either (a)  
1086 invitroSPI or (b) invitroPB in the PB dataset. c) Estimated FDRs based on spectral angle  
1087 distributions, choosing a spectral angle cut-off of 0.7 (dash line) reported in (a,b). The bars  
1088 represent the relative frequency of PSMs below the cut-off in each peptide strata. Statistically  
1089 significant p values < 0.05 (two-samples Wilcoxon test) are reported in (c), and they refer to the  
1090 comparison of the spectral angle distribution shown in (a,b).

1091

1092

1093 **Figure 6. Generation efficiency of spliced and non-spliced peptides.** Violin plots show the  
1094 distribution of generation efficiencies for peptide hydrolysis and splicing. Generation efficiencies  
1095 were calculated as the number of detected over the number of theoretically possible peptides  
1096 for each substrate. Calculations were carried out on the peptide products and substrate  
1097 sequences in the whole dataset digested with 20S standard proteasome (80 substrates). The  
1098 generation efficiency differs significantly between spliced and non-spliced peptides and, among  
1099 spliced peptides, between *cis*- and *trans*-spliced peptides. Significant p values of a two-samples  
1100 Wilcoxon test are reported.

1101

1102

1103 **Figure 7. Features of unique peptides identified in all datasets. a,b)** Frequency (a) and  
1104 length (b) of unique peptides per substrate. c) Length of N- and C-terminal splice-reactant of  
1105 *cis*-spliced peptides that could unequivocally be assigned to a single position within a substrate.  
1106 In (a-c), analysis has been carried out in the 2/4 h *in vitro* digestions with 20S standard  
1107 proteasomes, derived from the PB dataset (24 substrates) analyzed by invitroSPI and invitroPB,  
1108 as well as from the Specht dataset (47 substrates) and the whole dataset (71 substrates)  
1109 analyzed by invitroSPI. Here, PTM-tagged peptides identified by invitroSPI are added to the  
1110 unmodified peptides. In (a-c), all peptides that could not be unambiguously annotated as either  
1111 forward or reverse *cis*-spliced peptides (*i.e.* the multi-mapper forward/reverse *cis*-spliced  
1112 peptides) were removed. Spliced peptides containing a single amino acid residue splice-  
1113 reactant or the substrate's N- or C-termini were labelled as such only if that was the only

1114 explanation out of all possible peptide origins within the polypeptide substrate. In (c), multi-  
1115 mapper peptides that could be assigned unambiguously to a spliced peptide type were  
1116 subsequently checked for the length of their splice-reactants. Among multi-mapper spliced  
1117 peptides, only those that had a single and unambiguous splice-reactant length are included.

1118

1119

1120 **Figure 8. Potential pitfalls in data analysis related to peptide product database size. a)**

1121 Normalization strategies. Heatmaps display the joint frequency of amino acid combinations at  
1122 the splice-site (formed by sP1 and sP1') in the simulated background databases normalized by  
1123 the amino acid frequency of the investigated substrates. Simulated background databases were  
1124 computed from the PB dataset (n = 25 substrates) and from the whole dataset (n = 80  
1125 substrates). Frequencies were then normalized by the frequency of the amino acids within the  
1126 substrate sequences. White spots indicate combinations that are impossible to derive from the  
1127 given set of substrate sequences. Low frequencies are depicted in red, whereas high  
1128 frequencies are shown in blue. b) Amino acid frequencies at sP1 and sP1' sites of forward and  
1129 reverse *cis*-spliced peptides in the whole database of unique peptide products identified through  
1130 invitroSPI, as well as those sequences originally published by Paes *et al.* The frequency in the  
1131 true dataset was normalized by the frequency of the respective simulated background database  
1132 as well as by the sum of all values. To verify the robustness of the frequency estimation, 200  
1133 bootstrap iterations were performed, each time sampling 80 % of the splice-sites. The 90 %  
1134 confidence intervals of the resulting frequency estimations are displayed. Large confidence  
1135 intervals indicate low robustness of the frequency estimation.

1136

1137 **References**

- 1138 1 Hanada, K., Yewdell, J. W. & Yang, J. C. Immune recognition of a human renal cancer  
 1139 antigen through post-translational protein splicing. *Nature* **427**, 252-256 (2004).  
 1140 2 Vigneron, N. *et al.* An antigenic peptide produced by peptide splicing in the  
 1141 proteasome. *Science* **304**, 587-590 (2004).  
 1142 3 Mishto, M. & Liepe, J. Post-Translational Peptide Splicing and T Cell Responses.  
 1143 *Trends Immunol* **38**, 904-915, doi:10.1016/j.it.2017.07.011 (2017).  
 1144 4 Berkers, C. R. *et al.* Definition of Proteasomal Peptide Splicing Rules for High-  
 1145 Efficiency Spliced Peptide Presentation by MHC Class I Molecules. *J Immunol* **195**,  
 1146 4085-4095 (2015).  
 1147 5 Mishto, M. *et al.* Driving Forces of Proteasome-catalyzed Peptide Splicing in Yeast  
 1148 and Humans. *Mol Cell Proteomics* **11**, 1008-1023 (2012).  
 1149 6 Mishto, M. *et al.* An in silico-in vitro Pipeline Identifying an HLA-A(\*)02:01(+) KRAS  
 1150 G12V(+) Spliced Epitope Candidate for a Broad Tumor-Immune Response in Cancer  
 1151 Patients. *Front Immunol* **10**, 2572, doi:10.3389/fimmu.2019.02572 (2019).  
 1152 7 Kuckelkorn, U. *et al.* Proteolytic dynamics of human 20S thymoproteasome. *J Biol*  
 1153 *Chem* **294**, 7740-7754, doi:10.1074/jbc.RA118.007347 (2019).  
 1154 8 Specht, G. *et al.* Large database for the analysis and prediction of spliced and non-  
 1155 spliced peptide generation by proteasomes. *Sci Data* **7**, 146, doi:10.1038/s41597-  
 1156 020-0487-6 (2020).  
 1157 9 Dalet, A., Vigneron, N., Stroobant, V., Hanada, K. & Van den Eynde, B. J. Splicing of  
 1158 distant Peptide fragments occurs in the proteasome by transpeptidation and  
 1159 produces the spliced antigenic peptide derived from fibroblast growth factor-5. *J*  
 1160 *Immunol* **184**, 3016-3024 (2010).  
 1161 10 Faridi, P. *et al.* A subset of HLA-I peptides are not genomically templated: Evidence  
 1162 for cis- and trans-spliced peptide ligands. *Sci Immunol* **3**, eaar3947,  
 1163 doi:10.1126/sciimmunol.aar3947 (2018).  
 1164 11 Faridi, P. *et al.* Spliced Peptides and Cytokine-Driven Changes in the  
 1165 Immunopeptidome of Melanoma. *Cancer Immunol Res* **8**, 1322-1334,  
 1166 doi:10.1158/2326-6066.CIR-19-0894 (2020).  
 1167 12 Liepe, J. *et al.* A large fraction of HLA class I ligands are proteasome-generated  
 1168 spliced peptides. *Science* **354**, 354-358 (2016).  
 1169 13 Liepe, J., Sidney, J., Lorenz, F. K. M., Sette, A. & Mishto, M. Mapping the MHC Class I-  
 1170 Spliced Immunopeptidome of Cancer Cells. *Cancer Immunol Res* **7**, 62-76,  
 1171 doi:10.1158/2326-6066.CIR-18-0424 (2019).  
 1172 14 Paes, W. *et al.* Contribution of proteasome-catalyzed peptide cis-splicing to viral  
 1173 targeting by CD8(+) T cells in HIV-1 infection. *Proc Natl Acad Sci U S A* **116**, 24748-  
 1174 24759, doi:10.1073/pnas.1911622116 (2019).  
 1175 15 Platteel, A. C. M. *et al.* Multi-level Strategy for Identifying Proteasome-Catalyzed  
 1176 Spliced Epitopes Targeted by CD8+ T Cells during Bacterial Infection. *Cell Rep* **20**,  
 1177 1242-1253, doi:10.1016/j.celrep.2017.07.026 (2017).  
 1178 16 Platteel, A. C. *et al.* CD8(+) T cells of *Listeria monocytogenes*-infected mice recognize  
 1179 both linear and spliced proteasome products. *Eur J Immunol* **46**, 1109-1118,  
 1180 doi:10.1002/eji.201545989 (2016).  
 1181 17 Mansurkhodzhaev, A., Barbosa, C. R. R., Mishto, M. & Liepe, J. Proteasome-  
 1182 Generated cis-Spliced Peptides and Their Potential Role in CD8(+) T Cell Tolerance.  
 1183 *Front Immunol* **12**, 614276, doi:10.3389/fimmu.2021.614276 (2021).  
 1184 18 Mishto, M., Mansurkhodzhaev, A., Rodriguez-Calvo, T. & liepe, J. Potential mimicry  
 1185 of viral and pancreatic beta cell antigens through non-spliced and cis-spliced zwitter  
 1186 epitope candidates in Type 1 Diabetes. *Front Immunol* **12**, 656451, doi:doi:  
 1187 10.3389/fimmu.2021.656461 (2021).

1188 19 Mishto, M., Rodriguez-Hernandez, G., Neefjes, J., Urlaub, H. & Liepe, J. Response:  
1189 Commentary: An In Silico-In Vitro Pipeline Identifying an HLA-A\*02:01+ KRAS G12V+  
1190 Spliced Epitope Candidate for a Broad Tumor-Immune Response in Cancer Patients.  
1191 *Front Immunol* **12**, 679836, doi:10.3389/fimmu.2021.679836 (2021).

1192 20 Ebstein, F. *et al.* Proteasomes generate spliced epitopes by two different  
1193 mechanisms and as efficiently as non-spliced epitopes. *Sci Rep* **6**, 24032 (2016).

1194 21 Kato, K. *et al.* Characterization of Proteasome-Generated Spliced Peptides Detected  
1195 by Mass Spectrometry. *J Immunol* **208**, 2856-2865, doi:10.4049/jimmunol.2100717  
1196 (2022).

1197 22 Robbins, P. F. *et al.* Recognition of tyrosinase by tumor-infiltrating lymphocytes from  
1198 a patient responding to immunotherapy. *Cancer Res* **54**, 3124-3126 (1994).

1199 23 Dalet, A. *et al.* An antigenic peptide produced by reverse splicing and double  
1200 asparagine deamidation. *Proc Natl Acad Sci U S A* **108**, E323-E331 (2011).

1201 24 Mishto, M. Commentary: Are there indeed spliced peptides in the  
1202 immunopeptidome? *Mol Cell Proteomics*, 100158, doi:10.1016/j.mcpro.2021.100158  
1203 (2021).

1204 25 Liepe, J., Ovaa, H. & Mishto, M. Why do proteases mess up with antigen  
1205 presentation by re-shuffling antigen sequences? *Curr Opin Immunol* **52**, 81-86,  
1206 doi:10.1016/j.coi.2018.04.016 (2018).

1207 26 Reed, B. *et al.* Lysosomal cathepsin creates chimeric epitopes for diabetogenic CD4 T  
1208 cells via transpeptidation. *J Exp Med* **218**, doi:10.1084/jem.20192135 (2021).

1209 27 Fuchs, A. C. D. *et al.* Archaeal Connectase is a specific and efficient protein ligase  
1210 related to proteasome beta subunits. *Proc Natl Acad Sci U S A* **118**,  
1211 doi:10.1073/pnas.2017871118 (2021).

1212 28 Berkers, C. R., de Jong, A., Ovaa, H. & Rodenko, B. Transpeptidation and reverse  
1213 proteolysis and their consequences for immunity. *Int J Biochem Cell Biol* **41**, 66-71  
1214 (2009).

1215 29 Dalet, A., Stroobant, V., Vigneron, N. & Van den Eynde, B. J. Differences in the  
1216 production of spliced antigenic peptides by the standard proteasome and the  
1217 immunoproteasome. *Eur J Immunol* **41**, 39-46 (2011).

1218 30 Mishto, M. *et al.* Proteasome isoforms exhibit only quantitative differences in  
1219 cleavage and epitope generation. *Eur J Immunol* **44**, 3508-3521 (2014).

1220 31 Groll, M. & Huber, R. Substrate access and processing by the 20S proteasome core  
1221 particle. *Int J Biochem Cell Biol* **35**, 606-616 (2003).

1222 32 Huber, E. M. *et al.* Immuno- and constitutive proteasome crystal structures reveal  
1223 differences in substrate and inhibitor specificity. *Cell* **148**, 727-738 (2012).

1224 33 Liepe, J. *et al.* Quantitative time-resolved analysis reveals intricate, differential  
1225 regulation of standard- and immuno-proteasomes. *Elife* **4**, e07545, doi:doi:  
1226 10.7554/eLife.07545 (2015).

1227 34 Ben-Nissan, G. & Sharon, M. Regulating the 20S proteasome ubiquitin-independent  
1228 degradation pathway. *Biomolecules* **4**, 862-884 (2014).

1229 35 Gubin, M. M. *et al.* Checkpoint blockade cancer immunotherapy targets tumour-  
1230 specific mutant antigens. *Nature* **515**, 577-581 (2014).

1231 36 Gonzalez-Duque, S. *et al.* Conventional and Neo-Antigenic Peptides Presented by  
1232 beta Cells Are Targeted by Circulating Naive CD8+ T Cells in Type 1 Diabetic and  
1233 Healthy Donors. *Cell Metab* **28**, 946-960, doi:10.1016/j.cmet.2018.07.007 (2018).

1234 37 Wu, J. *et al.* DeepHLApan: A Deep Learning Approach for Neoantigen Prediction  
1235 Considering Both HLA-Peptide Binding and Immunogenicity. *Front Immunol* **10**, 2559,  
1236 doi:10.3389/fimmu.2019.02559 (2019).

1237 38 Rizvi, N. A. *et al.* Cancer immunology. Mutational landscape determines sensitivity to  
1238 PD-1 blockade in non-small cell lung cancer. *Science* **348**, 124-128,  
1239 doi:10.1126/science.aaa1348 (2015).

1240 39 Riley, T. P. *et al.* Structure Based Prediction of Neoantigen Immunogenicity. *Front*  
1241 *Immunol* **10**, 2047, doi:10.3389/fimmu.2019.02047 (2019).

1242 40 Luksza, M. *et al.* A neoantigen fitness model predicts tumour response to checkpoint  
1243 blockade immunotherapy. *Nature* **551**, 517-520, doi:10.1038/nature24473 (2017).

1244 41 Balachandran, V. P. *et al.* Identification of unique neoantigen qualities in long-term  
1245 survivors of pancreatic cancer. *Nature* **551**, 512-516, doi:10.1038/nature24462  
1246 (2017).

1247 42 Faridi, P., Dorvash, M. & Purcell, A. W. Spliced HLA bound peptides; a Black-Swan  
1248 event in Immunology. *Clin Exp Immunol* **204**, 179-188, doi:10.1111/cei.13589 (2021).

1249 43 Admon, A. Are There Indeed Spliced Peptides in the Immunopeptidome? *Mol Cell*  
1250 *Proteomics* **20**, 100099, doi:10.1016/j.mcpro.2021.100099 (2021).

1251 44 Mishto, M. *et al.* Database search engines and target database features impinge  
1252 upon the identification of post-translationally cis-spliced peptides in HLA class I  
1253 immunopeptidomes. *Proteomics* **22**, e2100226, doi:10.1002/pmic.202100226  
1254 (2022).

1255 45 Chapiro, J. *et al.* Destructive cleavage of antigenic peptides either by the  
1256 immunoproteasome or by the standard proteasome results in differential antigen  
1257 presentation. *J Immunol* **176**, 1053-1061 (2006).

1258 46 Deol, P., Zaiss, D. M., Monaco, J. J. & Sijts, A. J. Rates of processing determine the  
1259 immunogenicity of immunoproteasome-generated epitopes. *J Immunol* **178**, 7557-  
1260 7562 (2007).

1261 47 Guillaume, B. *et al.* Two abundant proteasome subtypes that uniquely process some  
1262 antigens presented by HLA class I molecules. *Proc Natl Acad Sci U S A* **107**, 18599-  
1263 18604 (2010).

1264 48 Guillaume, B. *et al.* Analysis of the processing of seven human tumor antigens by  
1265 intermediate proteasomes. *J Immunol* **189**, 3538-3547 (2012).

1266 49 Tenzer, S. *et al.* Antigen processing influences HIV-specific cytotoxic T lymphocyte  
1267 immunodominance. *Nat Immunol* **10**, 636-646 (2009).

1268 50 Zanker, D., Waithman, J., Yewdell, J. W. & Chen, W. Mixed Proteasomes Function To  
1269 Increase Viral Peptide Diversity and Broaden Antiviral CD8+ T Cell Responses. *J*  
1270 *Immunol* **191**, 52-59 (2013).

1271 51 Michaux, A. *et al.* A spliced antigenic peptide comprising a single spliced amino acid  
1272 is produced in the proteasome by reverse splicing of a longer peptide fragment  
1273 followed by trimming. *J Immunol* **192**, 1962-1971 (2014).

1274 52 Platteel, A. C. *et al.* CD8 T cells of Listeria monocytogenes-infected mice recognize  
1275 both linear and spliced proteasome products. *Eur J Immunol* (2016).

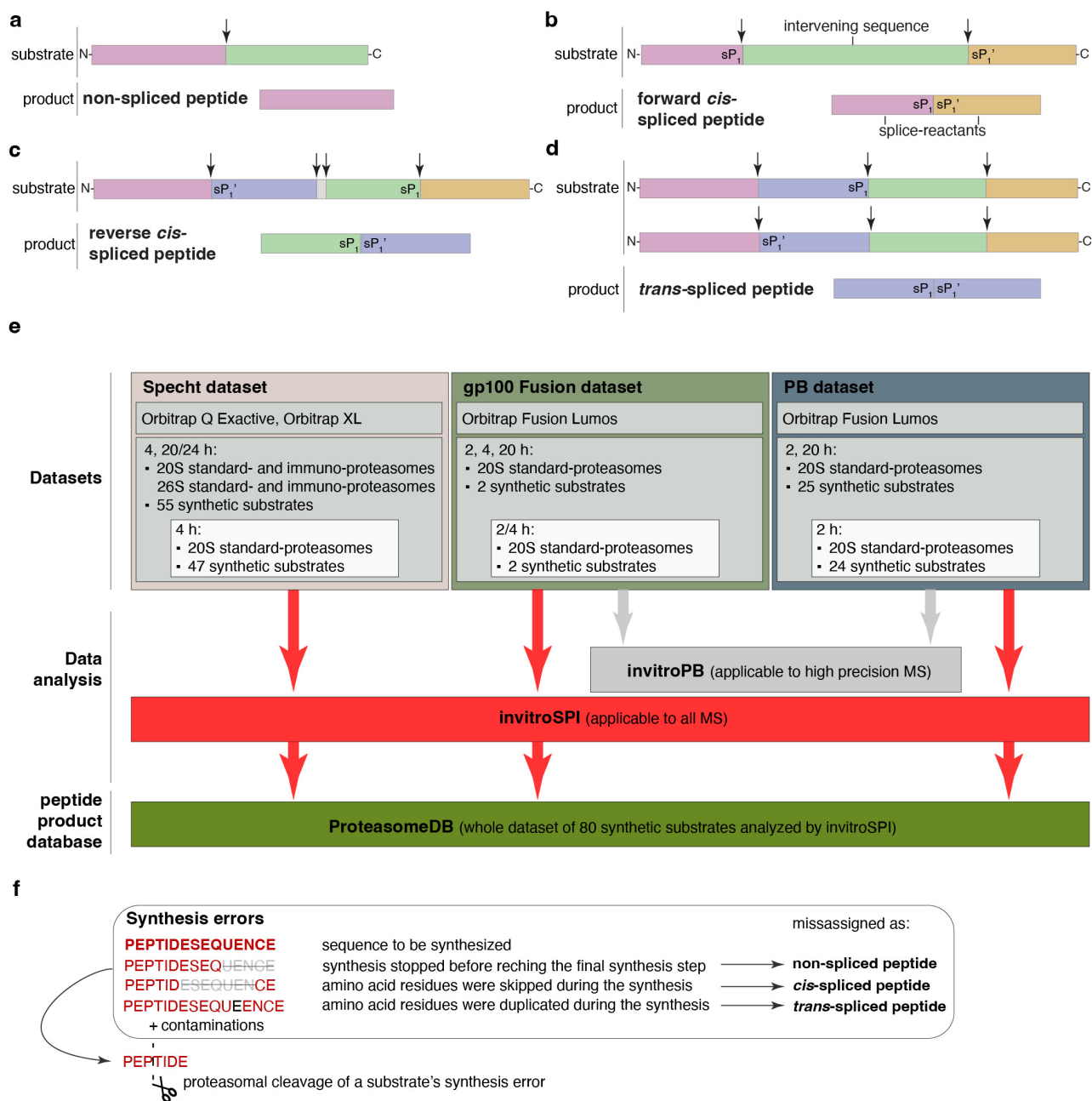
1276 53 Warren, E. H. *et al.* An antigen produced by splicing of noncontiguous peptides in the  
1277 reverse order. *Science* **313**, 1444-1447 (2006).

1278 54 Tsvetkov, P., Reuven, N., Prives, C. & Shaul, Y. Susceptibility of p53 unstructured N  
1279 terminus to 20 S proteasomal degradation programs the stress response. *J Biol Chem*  
1280 **284**, 26234-26242, doi:10.1074/jbc.M109.040493 (2009).

1281 55 Myers, N. *et al.* The Disordered Landscape of the 20S Proteasome Substrates Reveals  
1282 Tight Association with Phase Separated Granules. *Proteomics* **18**, e1800076,  
1283 doi:10.1002/pmic.201800076 (2018).

1284 56 Fabre, B. *et al.* Label-free quantitative proteomics reveals the dynamics of  
1285 proteasome complexes composition and stoichiometry in a wide range of human cell  
1286 lines. *J Proteome Res* **13**, 3027-3037, doi:10.1021/pr500193k (2014).

1287 57 Paes, W. *et al.* Elucidation of the Signatures of Proteasome-Catalyzed Peptide  
1288 Splicing. *Front Immunol* **11**, 563800, doi:10.3389/fimmu.2020.563800 (2020).  
1289 58 Mishto, M. *et al.* The immunoproteasome beta5i subunit is a key contributor to  
1290 ictogenesis in a rat model of chronic epilepsy. *Brain Behav Immun* **49**, 188-196  
1291 (2015).  
1292 59 Collins, G. A. & Goldberg, A. L. The Logic of the 26S Proteasome. *Cell* **169**, 792-806,  
1293 doi:10.1016/j.cell.2017.04.023 (2017).  
1294 60 Gessulat, S. *et al.* Prosit: proteome-wide prediction of peptide tandem mass spectra  
1295 by deep learning. *Nat Methods* **16**, 509-518, doi:10.1038/s41592-019-0426-7 (2019).  
1296 61 Wilhelm, M. *et al.* Deep learning boosts sensitivity of mass spectrometry-based  
1297 immunopeptidomics. *Nat Commun* **12**, 3346, doi:10.1038/s41467-021-23713-9  
1298 (2021).  
1299 62 Toprak, U. H. *et al.* Conserved peptide fragmentation as a benchmarking tool for  
1300 mass spectrometers and a discriminating feature for targeted proteomics. *Mol Cell*  
1301 *Proteomics* **13**, 2056-2071, doi:10.1074/mcp.O113.036475 (2014).  
1302 63 Liepe, J. *et al.* The 20S Proteasome Splicing Activity Discovered by SpliceMet. *PLOS*  
1303 *Computational Biology* **6**, e1000830 (2010).  
1304 64 Tran, N. H., Zhang, X., Xin, L., Shan, B. & Li, M. De novo peptide sequencing by deep  
1305 learning. *Proc Natl Acad Sci U S A* **114**, 8247-8252, doi:10.1073/pnas.1705691114  
1306 (2017).  
1307 65 Paes, W. *et al.* Corrigendum: Elucidation of the Signatures of Proteasome-Catalysed  
1308 Peptide Splicing. *Front Immunol* **12**, 755002, doi:10.3389/fimmu.2021.755002  
1309 (2021).  
1310 66 Tran, N. H. *et al.* Deep learning enables de novo peptide sequencing from data-  
1311 independent-acquisition mass spectrometry. *Nat Methods* **16**, 63-66,  
1312 doi:10.1038/s41592-018-0260-3 (2019).  
1313 67 Perez-Riverol, Y. *et al.* The PRIDE database and related tools and resources in 2019:  
1314 improving support for quantification data. *Nucleic Acids Res* **47**, D442-D450,  
1315 doi:10.1093/nar/gky1106 (2019).  
1316 68 Specht, G. *et al.* Digestion of a variety of synthetic peptides by proteasomes. *PRIDE*  
1317 <https://www.ebi.ac.uk/pride/archive/projects/PXD016782>. (2020).  
1318 69 Roetschke, H. P., Mishto, M. & Liepe, J. Digestion of TSN2 and TSN89 synthetic  
1319 peptides by proteasomes. *PRIDE*  
1320 <https://www.ebi.ac.uk/pride/archive/projects/PXD025995> (2021).  
1321 70 Roetschke, H. P., Mishto, M. & Liepe, J. Database and scripts from 'InvitroSPI and a  
1322 large database of proteasome-generated spliced and non-spliced peptides'. *Figshare*  
1323 <https://doi.org/XXXXX> (2022).  
1324 71 Wolf-Levy, H. *et al.* Revealing the cellular degradome by mass spectrometry analysis  
1325 of proteasome-cleaved peptides. *Nat Biotechnol*, doi:10.1038/nbt.4279 (2018).  
1326 72 Verbruggen, S. *et al.* Spectral prediction features as a solution for the search space  
1327 size problem in proteogenomics. *Mol Cell Proteomics*, 100076,  
1328 doi:10.1016/j.mcpro.2021.100076 (2021).  
1329 73 Cormican, J. A., Soh, W. T., Mishto, M. & Liepe, J. iBench: A ground truth approach  
1330 for advanced validation of mass spectrometry identification method. *Proteomics*,  
1331 e2200271, doi:10.1002/pmic.202200271 (2022).  
1332



**Figure 1**



**invitroSPI****invitroPB****strategy of PSM assignment**

non-spliced peptide candidates are favored over spliced peptide candidates  
 non-spliced, *cis*-spliced and *trans*-spliced peptides can be assigned to MS2 spectra

non-spliced peptide candidates are favored over spliced peptide candidates  
 non-spliced and *cis*-spliced peptides can be assigned to MS2 spectra

**peptide length**

5 residues or longer for all peptides

5 residues or longer for non-spliced peptides

8 residues or longer for *cis*-spliced peptides

**splice-reactant length**

no limit

2 residues or longer

**PTMs**

M oxidation, N/Q deamidation for all peptides

no PTMs allowed.

List of non-spliced peptides with 313 PTMs (PEAKS-PTM) used to remove MS2 spectra that could be assigned to *cis*-spliced peptides

**removal of synthesis errors**

non-spliced peptides: removed if present in controls  
 spliced peptides: removed if present as such or as N-/C-terminal precursors in controls  
 spliced peptides: flagged in ProteasomeDB if their splice-reactants are present in controls

non-spliced peptides: removed if present in controls  
*cis*-spliced peptides: removed if present in controls

**peptide products**

non-spliced, *cis*-spliced and *trans*-spliced peptides

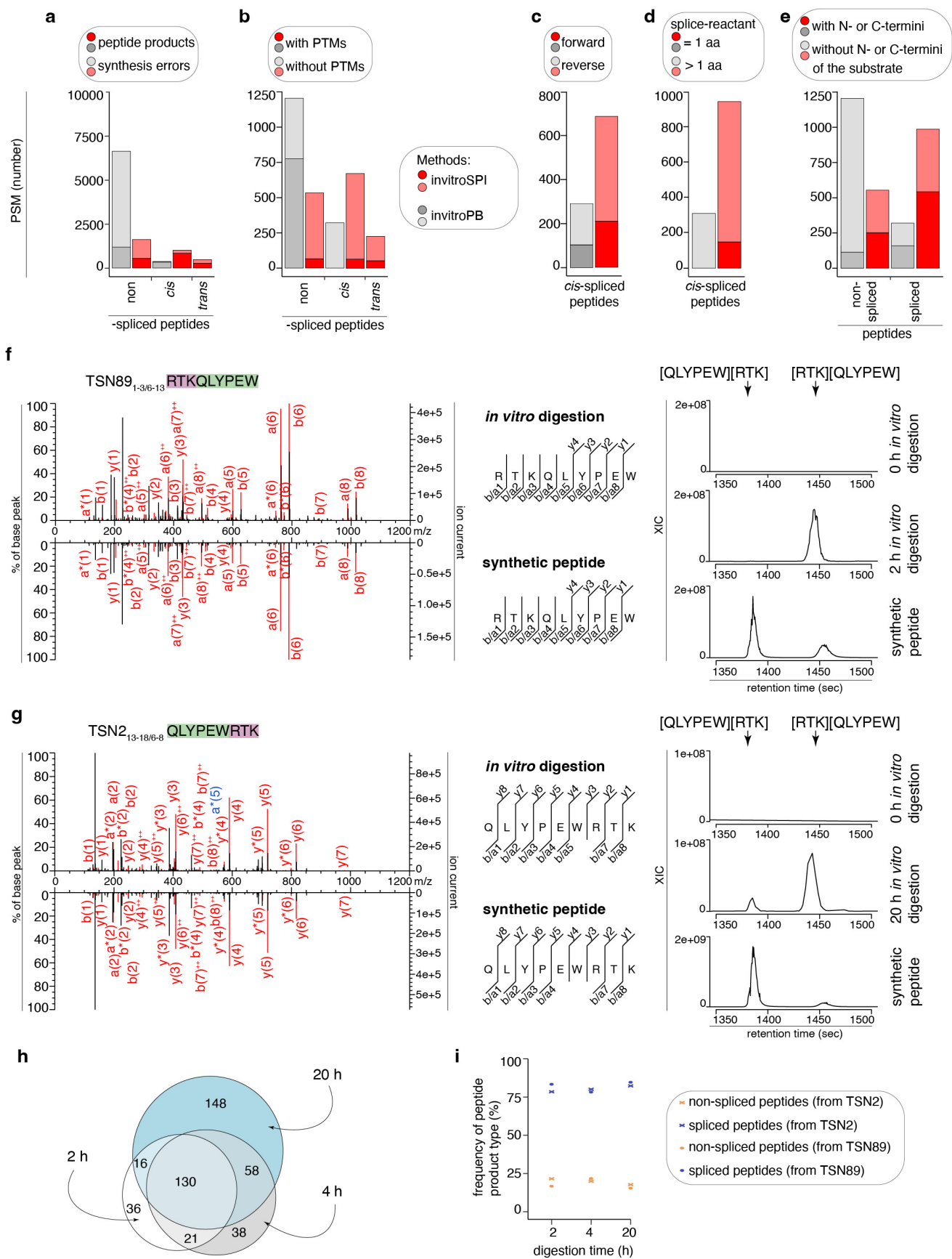
non-spliced and *cis*-spliced peptides

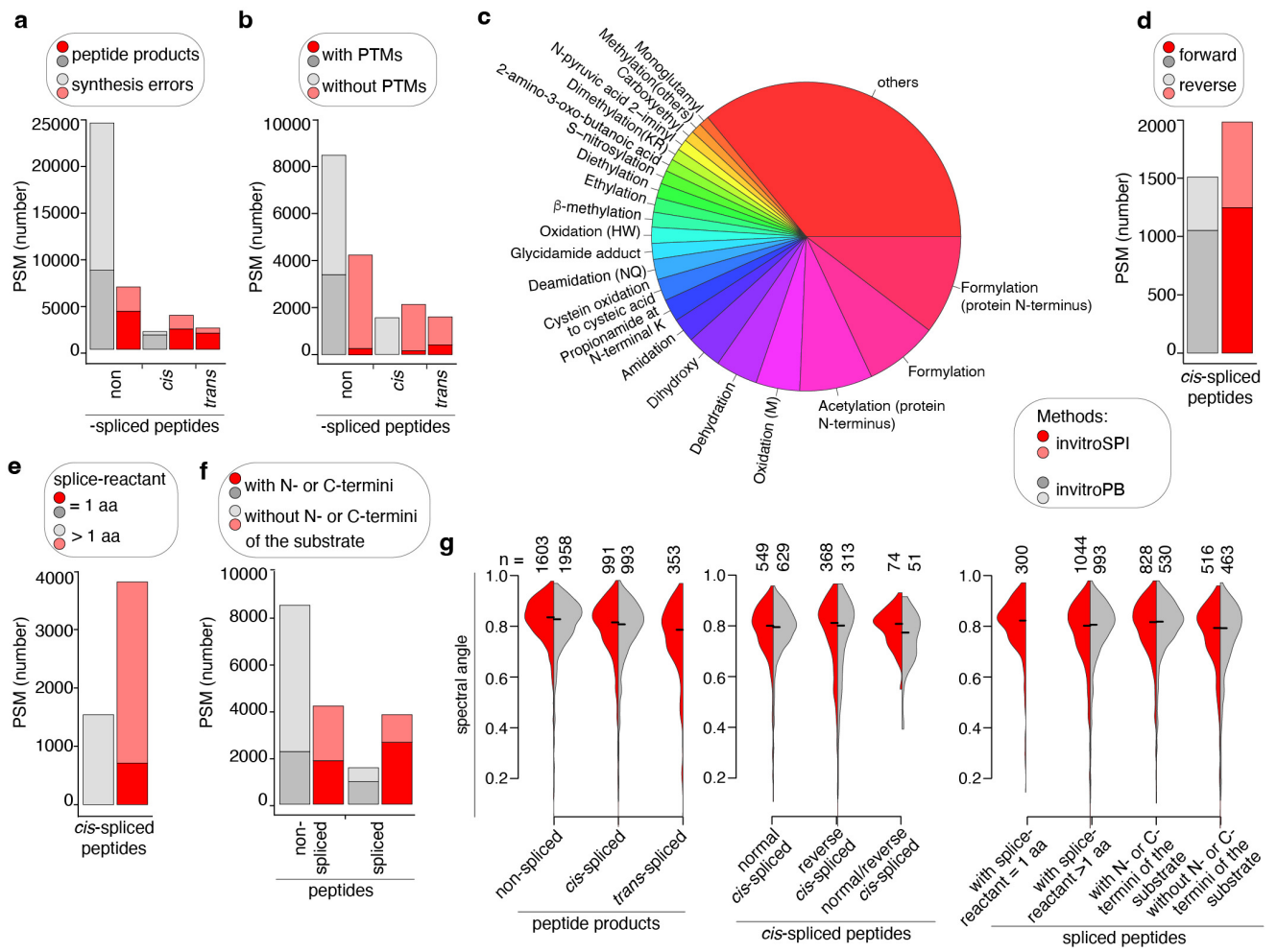
**downstream analysis**

peptides carrying the N-/C-terminus

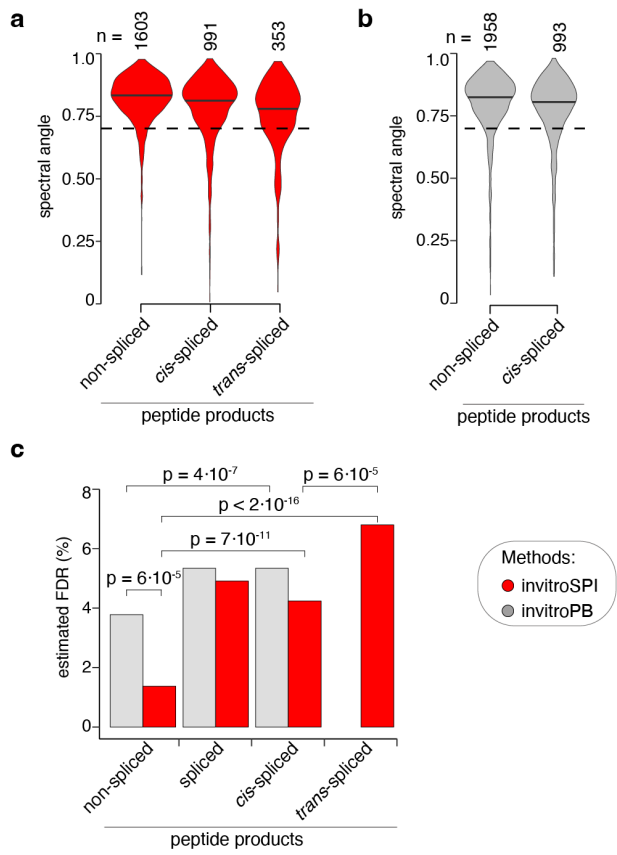
peptides carrying the N-/C-terminus: removed

**Figure 2**

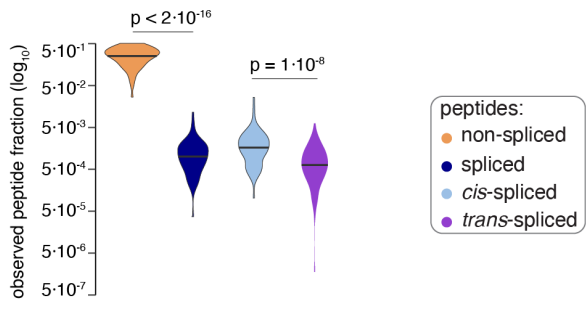




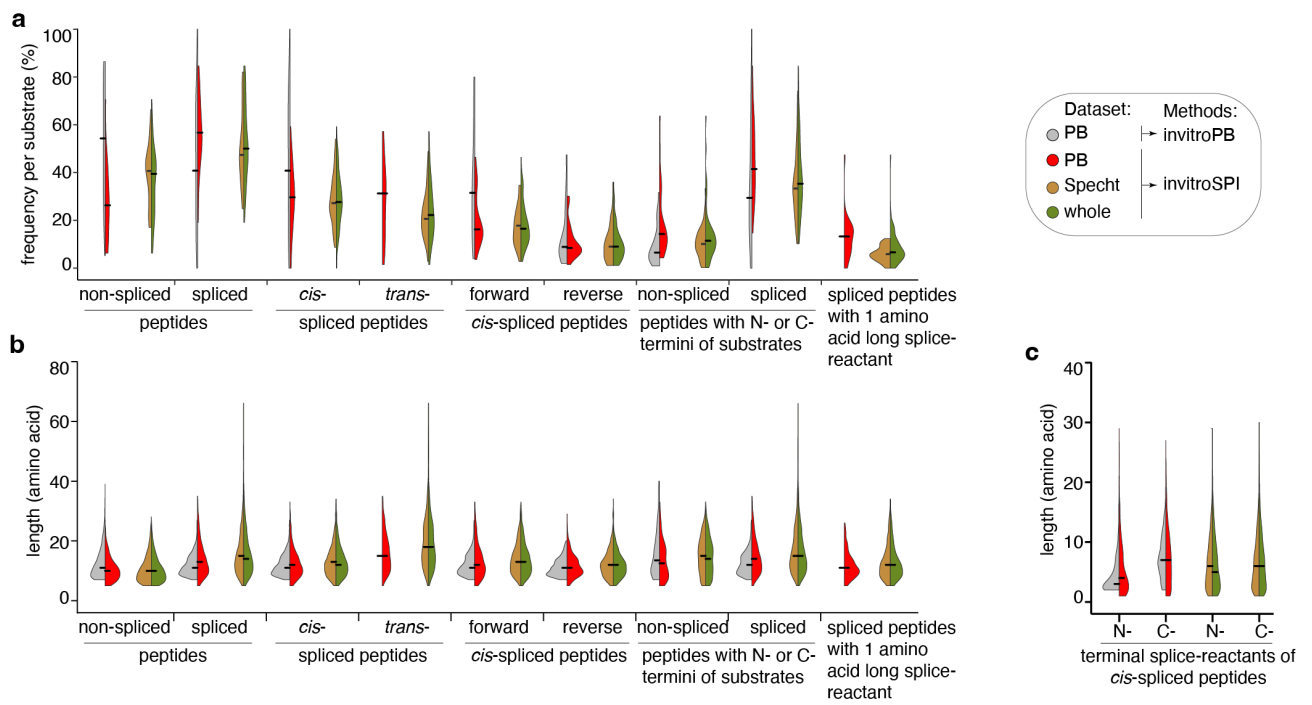
**Figure 4**



**Figure 5**



**Figure 6**



**Figure 7**

