

1 **Community Research article**

2

3 **Title**

4 Finding the LMA needle in the wheat proteome haystack.

5

6 **Authors**

7 Delphine Vincent 1, \*, AnhDuyen Bui 1, Vilnis Ezernieks 1, Saleh Shahinfar 1, Timothy Luke  
8 1, Doris Ram 1, Nicholas Rigas 2, Joe Panozzo 2,3, Simone Rochfort 1,4, Hans Daetwyler 1,4  
9 and Matthew Hayden 1,4

10

11 **Affiliations**

12 1 Agriculture Victoria Research, AgriBio, Center Centre for AgriBioscience, 5 Ring  
13 Road, Bundoora, VIC 3083, Australia; AnhDuyen.Bui@agriculture.vic.gov.au (A.B.);  
14 Doris.Ram@agriculture.vic.gov.au (D.R.); Vilnis.Ezernieks@agriculture.vic.gov.au (V.E.);  
15 saleh.shahinfar@agriculture.vic.gov.au (S.S.); tim.luke@agriculture.vic.gov.au (T.L);  
16 Simone.Rochfort@agriculture.vic.gov.au (S.R.); nicholas.rigas@agriculture.vic.gov.au  
17 (N.R.), joe.panozzo@agriculture.vic.gov.au (J.P.); hansdd@gmail.com (H.D.);  
18 matthew.hayden@agriculture.vic.gov.au (M.H.)

19 2 Agriculture Victoria Research, 110 Natimuk Road, Horsham, VIC 3400, Australia;  
20 Joe.Panozzo@agriculture.vic.gov.au (J.P.); pankaj.maharjan@agriculture.vic.gov.au (P.M.)

21 3 Centre for Agricultural Innovation, University of Melbourne, Parkville, VIC 3010,  
22 Australia

23 4 School of Applied Systems Biology, La Trobe University, Bundoora, VIC 3083,  
24 Australia

25 \* Correspondence: delphine.vincent@agriculture.vic.gov.au

26

27 **Abstract**

28 Late maturity alpha-amylase (LMA) is a wheat genetic defect causing the synthesis of high  
29 isoelectric point (pI) alpha-amylase in the aleurone as a result of a temperature shock during  
30 mid-grain development or prolonged cold throughout grain development leading to an  
31 unacceptable low falling numbers (FN) at harvest or during storage. High pI alpha-amylase is  
32 normally not synthesized until after maturity in seeds when they may sprout in response to rain  
33 or germinate following sowing the next season's crop. Whilst the physiology is well  
34 understood, the biochemical mechanisms involved in grain LMA response remain unclear. We

35 have employed high-throughput proteomics to analyse thousands of wheat flours displaying a  
36 range of LMA values. We have applied an array of statistical analyses to select LMA-  
37 responsive biomarkers and we have mined them using a suite of tools applicable to wheat  
38 proteins. To our knowledge, this is not only the first proteomics study tackling the wheat LMA  
39 issue, but also the largest plant-based proteomics study published to date. Logistics,  
40 technicalities, requirements, and bottlenecks of such an ambitious large-scale high-throughput  
41 proteomics experiment along with the challenges associated with big data analyses are  
42 discussed. We observed that stored LMA-affected grains activated their primary metabolisms  
43 such as glycolysis and gluconeogenesis, TCA cycle, along with DNA- and RNA binding  
44 mechanisms, as well as protein translation. This logically transitioned to protein folding  
45 activities driven by chaperones and protein disulfide isomerase, as well as protein assembly via  
46 dimerisation and complexing. The secondary metabolism was also mobilised with the up-  
47 regulation of phytohormones, chemical and defense responses. LMA further invoked cellular  
48 structures among which ribosomes, microtubules, and chromatin. Finally, and unsurprisingly,  
49 LMA expression greatly impacted grain starch and other carbohydrates with the up-regulation  
50 of alpha-gliadins and starch metabolism, whereas LMW glutenin, stachyose, sucrose, UDP-  
51 galactose and UDP-glucose were down-regulated. This work demonstrates that proteomics  
52 deserves to be part of the wheat LMA molecular toolkit and should be adopted by LMA  
53 scientists and breeders in the future.

54

## 55 **Keywords**

56 *Triticum aestivum*; large-scale high-throughput workflow; bottom-up shotgun proteomics;  
57 LC-MS/MS; late maturity alpha-amylase, LMA; big data; statistics, data mining

58

## 59 **Introduction**

60 Common bread wheat (*Triticum aestivum* L.) is the dominant crop in temperate regions,  
61 currently covering more than 220 million hectares worldwide, exceeding 749 million tons in  
62 production annually (1) and predicted to reach 835 million tons by 2030 (2). From the most  
63 primitive form of wheat 10,000 years ago in the Fertile Crescent to the species currently grown  
64 all over the world, desirable characteristics have been selected and improved upon by human  
65 societies (3). Hexaploid *T. aestivum* (AABBDD;  $2n = 6x = 42$ ) originated from two  
66 polyploidization events. The first event associated diploid *Triticum urartu* (AA;  $2n = 2x = 14$ )  
67 which provided the A genome with the other yet unknown species from the *Sitopsis* section of  
68 *Triticum* genus which provided the B genome to produce the allotetraploid wild emmer wheat

69 (*Triticum turgidum*; AABB;  $2n = 4x = 28$ ). The second event associated *T. turgidum* with  
70 *Aegilops tauschii* (DD) (4, 5). The chromosomes from each closely related progenitor are  
71 grouped into homeologous groups. Because of the shared ancestry, genes may be common to  
72 all members of a homeologous group, albeit exhibiting high allelic variation and differences in  
73 gene count due gene duplication or silencing (2). Millennia of domestication have accrued an  
74 enormous genetic diversity in this species, with potentially more than 50,000 *T. aestivum*  
75 cultivars (6). Wheat owes its success to adaptability to temperate, Mediterranean, and  
76 subtropical climates, high yields, storability, but above all to the unique properties of doughs,  
77 which can be processed into a vast range of foods (3, 5). Wheat grains are not only a major  
78 source of carbohydrate in the form of starch, but also a great source of protein. The endosperm  
79 prolamins proteins comprise gliadins and glutenins; they are the main components of gluten  
80 and together confer unique viscoelastic and rheological properties to flour mixed with water.  
81 Indeed, hydrated gliadins largely determine the viscosity and extensibility of the dough, while  
82 the cohesive properties of hydrated glutenins essentially govern the strength and elasticity of  
83 the dough (5). Wheat seeds also contribute essential amino acids, minerals, vitamins, beneficial  
84 phytochemicals and dietary fibre components to the human diet. Beside nutritional benefits,  
85 different parts of the wheat plant confer advantageous medicinal uses such as anticancer  
86 properties of wheat bran and antimicrobial activities of wheat sprouts (7).

87 Current breeding programs mainly aim at sustaining wheat production and quality with reduced  
88 agrochemical inputs, as well as developing new disease-resistant and stress-tolerant varieties  
89 with enhanced quality for specific end-uses (8). Wheat research and breeding must accelerate  
90 genetic gain to keep augmenting crop yield while maintaining or improving grain quality traits  
91 if the demands of the growing human population are to be met (9). A critical element in the  
92 equation was the sequencing and functional annotation of the genome. Sequencing the  
93 hexaploid bread wheat genome was a gigantic achievement proportionate to its large size,  
94 abundance of repetitive DNA and the immense difficulty of discerning homoeologs from  
95 subgenomes A, B and D. Whilst this required the commitment of 20 countries collaborating as  
96 a consortium (International Wheat Genome Sequencing Consortium IWGSC) and a lot of  
97 strategizing from 2005 onward, including sequencing diploid and tetraploid ancestors, it was  
98 the advent of next generation sequencing technologies producing long but error-prone or  
99 accurate yet short reads that made this massive endeavour successful (10). A 13-year effort  
100 ensued, drafting *T. aestivum* genome in 2014 based on key breakthrough short read  
101 technologies by NRGene (4), and culminating in 2018 with the release of the long-awaited  
102 fully assembled and annotated 14.5 Gb reference genome, cataloguing 107,891 high-

103 confidence genes along 21 chromosome-like sequence assemblies (IWGSC RefSeq v1.0) (9).  
104 This helped bridge the gap on wheat research relative to other cereal model species such as  
105 rice, sorghum, corn and barley whose genomes had been sequenced years ago, and propelled  
106 wheat post-genomics studies forward with a continuous increase in publications since 2011.  
107 Both the numbers of high confidence protein-coding genes from subgenomes A, B, and D and  
108 their composition were largely similar (9). Transcriptomics analyses of genes present in all  
109 three subgenomes not only showed comparable expression levels for 72.5% of them, especially  
110 those located in syntenic regions, but also unveiled the lack of significant subgenome  
111 expression dominance (11, 12). As valuable such an asset was, it did not capture the extent of  
112 the wheat genomic diversity as only one cultivar, Chinese Spring, was chosen as a template. In  
113 fact, no single genome assembly can be sufficient to model the wheat proteome due to the high  
114 allelic and gene copy number variability (2). This shortfall was addressed in 2020 when 15  
115 hexaploid wheat lines from different regions, growth habits and breeding programs were  
116 sequenced and annotated against IWGSC reference genome (13). Such pan-genomic  
117 comparative analysis outlined extensive structural rearrangements, introgressions from wild  
118 relatives and differences in gene content arising from complex breeding events to boost  
119 resistance to biotic and abiotic stresses, as well as grain yield and quality. Unfortunately, fasta  
120 sequences of annotated proteins are not publicly available for these 15 assemblies. A refined  
121 version of the reference genome using optical mapping and long sequence reads was recently  
122 released (IWGSC RefSeq v2.1) (14). With such worthwhile genomic resources in store, wheat  
123 can now be instated as a model for plant genetic research and employed to tackle complex  
124 biological questions on evolution, domestication, polyploidization, as well as genetic and  
125 epigenetic interaction between homoeologous genes and genomes (10). Moreover, genome  
126 annotations paves the way to investigate pathways and biochemical attributes behind bread  
127 wheat quality using transcriptomics (15) or proteomics (2) approaches.

128 The industry will equally benefit from these latest scientific developments since processing  
129 companies, markets and food industries demand not only high yielding and resistant varieties,  
130 but also those with specific end-use qualities (1, 3). Market requirements have influenced wheat  
131 breeding as not to neglect essential protein content and quality. Because wheat is generally  
132 traded according to grain protein content and hardness, standards must be abided to by  
133 producers and distributors. Intact starch polymers provide the gelatinization and retrogradation  
134 needed for an acceptable product. Failure to meet receival standards for milling grades due to  
135 starch degradation measured in the wheat industry using the Hagberg–Perten falling number  
136 (FN) method (16) leads to grain discount and downgrading to animal feed, which incurs a loss

137 of profit (17). The low FN values manifests as a loss of viscosity upon mixing starch-degraded  
138 flour with water can alter appearance and texture of end-products (18), however, it might not  
139 deteriorate baking functionality (19) and could be used instead in alternate preparations (20).  
140 There are multiple causes of low FN symptomatic of starch degradation including preharvest  
141 sprouting, late maturity alpha-amylase (LMA), and variation in kernel starch and protein (21).  
142 LMA is a wheat genetic defect causing the synthesis of high isoelectric point (pI) alpha-  
143 amylase in the aleurone as a result of a temperature shock during mid-grain development or  
144 prolonged cold throughout grain development leading to an unacceptable low FN at harvest or  
145 during storage (22-24). High pI alpha-amylase is normally not synthesized until after maturity  
146 in seeds when they may sprout in response to rain or germinate following sowing the next  
147 season's crop (25).

148 Four alpha-amylase isoforms have been identified to date in wheat. Several  $\alpha$ -amylase 1  
149 (TaAMY1) loci have been localized on the long arm of group 6 chromosomes (26). In LMA-  
150 prone wheat genotypes and under given temperatures, Amy-1 genes are transcribed in isolated  
151 cells or cell islands distributed throughout the aleurone system of grains with a 50-60%  
152 moisture content before they have reached physiological maturity (25). Appearance of high pI  
153 a-amylase protein is preceded by a short-lived transient period of mRNA synthesis leading to  
154 a stable enzyme and retained through to seed maturity (22, 27). Multiple alpha-amylase 2  
155 (TaAMY2) loci are positioned on the long arm of the group 7 chromosomes and produce a low  
156 pI alpha-amylase in the pericarp of the developing grain (28). A single locus encodes alpha-  
157 amylase 3 (TaAMY3) on group 5 chromosomes and is transcribed throughout the grain  
158 development suggesting a role in grain development and maturation (29). Similar to TaAMY2,  
159 TaAMY3 enzyme mainly appears during grain development in the pericarp and would be the  
160 predominant alpha-amylase enzyme throughout grain development (30). Despite its shorter  
161 length and elevated pI, TaAMY3 displays equal numbers of calcium-binding and active sites  
162 relative to the other three isoforms; however, the distance between key AA residues and the  
163 last two active site residues is shortened (31). Overexpressing TaAMY3 in the endosperm of  
164 developing grain to levels of up to 100-fold higher than the wild-type results in low FN similar  
165 to those seen in LMA-affected grains, yet has no detrimental effect neither on starch structure,  
166 flour composition and baking quality of bread (32), nor on noodle colour or firmness (33). A  
167 fourth isoform alpha-amylase 4 (TaAMY4) is also encoded by a single locus on group 5  
168 chromosomes and is co-expressed with TaAMY1 in LMA-affected grains (31). Comparison of  
169 the four isoforms revealed that they contain 385-439 AAs, with a molecular mass between  
170 45.4-48.3 kD, and a pI ranging from 5.5 to 8.6. All isoforms differ slightly in their 3-D protein

171 structure including the presence of additional sugar binding domains hinting to various  
172 enzymatic properties (31, 34).

173 Although LMA expression correlates with measurable changes in both hormone content and  
174 transcript profiles during grain maturation, there are no obvious visual effects on grain  
175 appearance, development, or morphology (24), hence the need to perform assays to test for its  
176 activity (16). ELISA (35) and RT-qPCR (36) assays were developed to specifically target  
177 TaAMY1, the main enzyme involved in LMA. One limitation to the RT-qPCR method relates  
178 to the apparent short life of the high pI  $\alpha$ -amylase mRNA (22). Commonly employed is the  
179 colorimetric Ceralpha assay (37) whereby the alpha-amylase activity is expressed in terms of  
180 Ceralpha units per gram of flour (u/g). A single unit corresponds to the amount of enzyme  
181 required to release 1  $\mu$ M p-nitrophenyl in the presence of excess quantities of alpha-glucosidase  
182 in 1 min at 40°C (38). Such measurements have revealed that LMA is more prevalent than  
183 originally thought, with reports arising from North America, Australia, Japan, Canada, South  
184 Africa, China, Mexico, Germany, and the United Kingdom (39). The presence of LMA in  
185 breeding populations could be attributed to unexplained positive effects on grain  
186 production/quality or alternately simply manifest the lack of significant selection pressure  
187 against this trait (24). Both a cool temperature shock near physiological maturity or continuous  
188 cool maximum temperatures during grain development can induce LMA synthesis in wheat  
189 (23). The prediction of LMA occurrence during LMA dedicated field trial is impeded by the  
190 stochastic nature of LMA expression resulting from specific genetics, climatic conditions, and  
191 developmental stages.

192 LMA has a genetic (G) component (alpha-amylase gene required), yet it is only expressed and  
193 enzymatically active under particular environmental (E) conditions (temperature shock) at a  
194 given developmental stage making it the product of a GxE interaction, which lends itself to  
195 post-genomic quantitative studies to shed some lights into the biological mechanisms  
196 underpinning LMA expression. Yet, to date, only one LMA-related transcriptomics study has  
197 been published and no proteomics work has been attempted despite the potential this  
198 technology offers to help improve bread wheat quality (2). Using microarray technology,  
199 Barrero and colleagues reported that LMA resulted from very narrow and transitory peak of  
200 expression of genes encoding high-pI alpha-amylase during grain development (22).  
201 Furthermore, the LMA phenotype triggered elevated levels of gibberellins such as GA19 and  
202 much lower levels of auxin in the de-embryonated fraction of grains sampled shortly after the  
203 initiation of LMA synthesis. A recent report questions this hormonal response since, unlike  
204 alpha-amylase synthesis by aleurone during germination or following treatment with



205 exogenous GA, alpha-amylase synthesis by wheat aleurone during grain development appears  
206 to be independent of gibberellin (40). Even though on one hand genomics can catalogue genes  
207 present in a sample and possibly the biological context of their expression and on the other  
208 hand transcriptomics can validate expression levels, only proteomics can measure the actual  
209 protein abundance, record post-translational modification (PTM), as well as identify interacting  
210 proteins (2). We have developed a high-throughput proteomics method to rapidly profile *T.*  
211 *aestivum* grains and data mine their proteome (41). In the present study, we have applied our  
212 optimised procedure to a collection of in excess of 4,000 wheat cultivars and germplasm whose  
213 LMA content ranged from 0 to 8 u/g of flour. We have applied an array of statistical analyses  
214 to our big data to select LMA-responsive biomarkers and we have mined them using a suite of  
215 tools applicable to wheat proteins, yet not necessarily embraced by grain scientists. To our  
216 knowledge, this is not only the first proteomics study tackling the wheat LMA issue but also  
217 the largest plant-based proteomics study published to date. Logistics, technicalities,  
218 requirements, and bottlenecks of such an ambitious large-scale high-throughput proteomics  
219 experiment along with the challenges associated with big data analyses are discussed.

220

## 221 **2. Materials and Methods**

### 222 **2.1. Wheat Cultivation, Sampling, and Storage**

223 The wheat collection used in this study represents a diverse range of cultivars and germplasm  
224 sourced through the Australian Grains Genebank and representing global genetic diversity. The  
225 wheat was grown in field trials at Horsham Victoria and harvested using a mechanical small-  
226 plot harvester. The threshed grain was stored in seal containers at 20°C. The environmental  
227 conditions (rain and temperature) at the trial site were monitored throughout the growing  
228 season. No preharvest rainfall was recorded and therefore any  $\alpha$ -amylase activity was non-  
229 germinative but associated with LMA.

230 The list of wheat samples is supplied in Supplementary Table S1.

### 231 **2.2. LMA assay**

232 The alpha-amylase assay was performed using the Megazyme assay according to the procedure  
233 reported by McCleary and Sheehan (42) on 3,773 grain samples (Supplementary Table S1).

234 The distribution of LMA values was plotted as a histogram in Microsoft Excel. Various  
235 transformations were performed to achieve a normal distribution such as standardisation, log  
236 natural, log 2, inverse and standardisation of inversed values (data not shown). The transformed  
237 values were also plotted as histograms to check for gaussian distribution.

### 238 **2.3. Wheat Grain Processing for proteomics analyses**

239 Sample preparation was optimised and thoroughly described (41); it is schematised in Figure  
240 1. All sample packages were mixed together in a box for randomisation and assigned a unique  
241 number as they were processed. QR codes on sample bags and tubes were scanned and  
242 consigned to the Excel spreadsheet using a handheld barcode scanner (model 1902 GHD-2,  
243 Honeywell Australia, Matraville, NSW). All microtubes were pre-labelled with unique  
244 numbers and sample IDs, both also consigned to a QR code, using a handheld label maker (PT-  
245 E550WVP, Brother, Australia) controlled by the P-touch editor software (Brother, Australia)  
246 fitted with 12mm white laminated tape.

247 The grains were ground in 50 mL jars containing two 8 mm and two 3 mm metal grinding balls  
248 using an automated tissue homogeniser and cell lyser (Geno/Grinder® 2010, SPEX  
249 SamplePrep, Metuchen, NJ, USA) and pulverised into fine flour twice for 2 min at 1,500 rpm  
250 with a 15 s break in between. A total of 600 jars were employed in a rotation. Dirty jars and  
251 balls were rinsed to remove excess flour and soaked in 1% decon 90 surfactant (Decon, Hove,  
252 UK) for 2 hours followed by a thorough wash in a dishwasher with RO water and left to air dry  
253 prior to being used again. A wheat quality control (QC) sample was prepared by sampling 50  
254 mg ( $\pm 0.05$  mg) from each of the 96 flour samples described in (41) and mixing them all  
255 thoroughly.

256 A 20 mg ( $\pm 0.2$  mg) aliquot of flour was weighed in a 1.5 mL microtube and resuspended in 0.5  
257 mL Gnd-HCl buffer (6 M Guanidine hydrochloride, 0.1 M Bis-Tris, 10 mM DTT, 5.37 mM  
258 sodium citrate tribasic dihydrate) using a MS 1.5 sonicator probe (Ultrasonic Homogeniser  
259 SONOPULS mini 20, Bandelin, Berlin, Germany) for 30 s with 90% amplitude. The tubes  
260 were briefly vortexed (5 sec each, RAVM1 Ratek Vortex Mixers, Ratek, Boronia, VIC,  
261 Australia) and incubated for 60 min in a thermoblock (Digital Dry Bath/Block Heater, Thermo  
262 Scientific, Scoresby, VIC, Australia) at 60°C. The tubes were left to cool to room temperature  
263 for 5 min and 10  $\mu$ L of 1 M iodoacetamide was added to each tube. The tubes were thoroughly  
264 mixed for 30 s using a rack vortex mixer (MTV1 Multi Tube Vortex Mixer, Ratek, Boronia,  
265 VIC, Australia) at high speed and left to incubate at room temperature in the dark for 30 min.  
266 The tubes were centrifuged using a benchtop centrifuge (5415D Digital Microfuge, Eppendorf,  
267 Macquarie Park, NSW, Australia) at 13,000 rpm for 15 min at room temperature and the  
268 supernatant was transferred into a fresh 1.5 mL microtube pre-labelled with the QR code.

269 Two vials of trypsin/Lys-C mix (100 $\mu$ g, V5078, Promega, Alexandria, NSW, Australia) were  
270 dissolved into 1 mL of the resuspension buffer (50mM acetic acid) supplied by the  
271 manufacturer and kept on ice until use to digest 192 wheat samples at a time. Aliquots of 10  
272  $\mu$ L aliquot of protein extracts were transferred into two 96-well plates (Strata 96-well collection



273 plate, 350  $\mu$ L conical polypropylene, Phenomenex, Lane Cove, NSW, Australia), diluted 6  
274 times with 50 mM ammonium bicarbonate and digested with 5  $\mu$ L aliquots of the trypsin/Lys-  
275 C solution prepared earlier. Plates were sealed with silicone covers (pierceable sealing mats,  
276 96-square well, Phenomenex, Lane Cove, NSW, Australia) and vortexed for 30 s using a rack  
277 vortex mixer (MTV1 Multi Tube Vortex Mixer, Ratek, Boronia, VIC, Australia) at high speed.  
278 Plates were incubated at 37°C for 17 hours. Aliquots of 7  $\mu$ L 10% formic acid (FA)/water were  
279 added to stop the digestion. An internal standard (IS, [Glu1]-fibrinopeptide B human, F3261,  
280 Sigma, Port Melbourne, VIC, Australia) was added at a final concentration of 1  $\mu$ g.  
281 Protein digests were cleaned using 96-wells solid phase extraction (SPE) plates (Strata C18-E  
282 100 mg P/N 8E-S001-EGB, Phenomenex, Lane Cove, NSW, Australia) and fully evaporated  
283 as described in (Vincent, Bui et al. 2022). Peptide digests were reconstituted by adding 70  $\mu$ L  
284 of 0.1% FA/water to each well. The digests were dissolved by shaking the plates for 50 min at  
285 medium speed using a rack vortex mixer (MTV1 Multi Tube Vortex Mixer, Ratek, Boronia,  
286 VIC, Australia) at room temperature. The collection plates were sealed with a silicone lid and  
287 stored at -80 °C until LC-MS analysis.

#### 288 **2.4. LC-MS analyses**

289 All 4,061 wheat and QC samples were processed using the LC-MS method listed below.  
290 Liquid chromatography (LC) was optimised (41). Our chosen LC method applied 0.2 mL/min  
291 flow rate, 38 min LC run duration, 6% B for 2.5 min, 6–36% B gradient for 30.5 min, increased  
292 up to 98% B gradient for 0.1 min, 98% B for 5 min, drop down to 3% B in 0.1 min, 6% B for  
293 5 min. The LC system used was a Vanquish Flex Binary UHPLC System (Vanquish UHPLC+  
294 focused, ThermoFisher Scientific, Scoresby, VIC, Australia). Mobile phase A was 0.1%  
295 FA/water and mobile phase B was 0.1% FA/acetonitrile (ACN). The needle wash solution was  
296 80% isopropanol (IPA)/water, and the rear seal wash solution was 10% IPA/water. The needle  
297 wash solution was 10% IPA/water. The needle was washed after each injection. The rack types  
298 were specified as DeepWell96 in the LC-MS method and the SamplerModule tab of Xcalibur  
299 Direct Control software (version 3.0.63, ThermoFisher Scientific, Scoresby, VIC, Australia)  
300 with a 29,000  $\mu$ m injection depth. Blanks (0.1% FA/water) and QC were injected from two 10  
301 mL vials. Peptides were separated using a RP-LC column (bioZen 1.7  $\mu$ m Peptide XB-C18,  
302 100 Å, LC column 150  $\times$  2.1 mm, Phenomenex, Lane Cove, NSW, Australia) using a 60°C  
303 oven temperature. The blank, IS and QC samples were injected every 48 samples for  
304 normalisation purposes. The IS was used to check for mass accuracy (<50ppm). The LC  
305 separation column was changed with a new one when peak resolution degraded (every 1000  
306 samples or so).

307 The UHPLC was online with an Orbitrap Velos hybrid ion trap–Orbitrap mass spectrometer  
308 (ThermoFisher Scientific, Scoresby, VIC, Australia) fitted with a heated electrospray  
309 ionisation (HESI) source. Every three weeks, the instrument was mass calibrated, and the  
310 source sweeping cone and the heated capillary were cleaned. HESI parameters were: needle at  
311 3.9 kV, 100  $\mu$ A, sheath gas flow 20, auxiliary gas flow 7, sweep gas flow 2, source heated to  
312 200°C, capillary heated to 275°C, and S-Lens RF level 55%. Spectra were acquired using the  
313 full MS scan mode of the Fourier transform (FT) orbitrap mass analyser (FTMS) in positive  
314 ion mode at a resolution of 15,000 along a 300–2000  $m/z$  mass window in profile mode with 3  
315 microscans.

316 The sequence lists were prepared in advance in Excel as .csv files and imported into Xcalibur  
317 data acquisition software (version 3.0.63); five sequences were needed as Xcalibur only  
318 accommodated a maximum of 1000 lines. Because samples had been randomised, 96-well  
319 plates were analysed consecutively. Throughout the LC-MS run, the RAW files were  
320 individually visualised using Xcalibur Qual Browser (version 3.0.63,). Files that failed to pass  
321 our check (loss of peak resolution, incomplete run, no signal, mass accuracy > 50 ppm, etc...) were rerun  
322 concomitantly to when LC-MS was interrupted for maintenance.

## 323 **2.5. LC-MS/MS analyses**

324 For protein identification, 400 random samples (10% samples, 4 plates) were used following  
325 the LC-MS1 analysis. LC, HESI and full scan FTMS parameters were as indicated above. MS2  
326 data was acquired using ITMS in positive mode as centroid values and applied various methods  
327 summarised below. In an attempt to maximise the number of peptides sequenced, several  
328 passes were performed with inclusion and exclusion lists, and various parameters.

329 Pass 1: FTMS parameters were as specified above. Using the Nth order double play method,  
330 MS/MS spectra were acquired in data-dependent mode. Singly charged peptides were ignored.  
331 In the linear ion trap, the 10 most abundant peaks with charge state >2 and a minimum signal  
332 threshold of 3,000 were fragmented using collision-induced dissociation (CID) with a  
333 normalised collision energy of 35%, 0.25 activation Q, and activation time of 10 ms. The  
334 precursor isolation width was 2  $m/z$ . Dynamic exclusion was activated, and peptides selected  
335 for fragmentation more than once within 30 s were excluded from selection for 180 s. No  
336 inclusion or exclusion list was used; however, a list of MS2 event was produced by exporting  
337 the “Scan Filters” of the RAW file in Xcalibur Qual Browser (ThermoFisher Scientific,  
338 Scoresby, VIC, Australia) and to be used in Pass 2 as an exclusion list containing 2,000 unique  
339  $m/z$  values (maximum number allowed in Xcalibur). This method was run in duplicate.

340 Pass 2: Same method as Pass 1, except that the list of MS2 events generated in Pass 1 was  
341 uploaded in the Data Dependent Settings as a Reject Mass List. Like in Pass 1, a list of MS2  
342 event was produced by exporting the “Scan Filters” of the RAW file and to be used in Pass 3  
343 as an exclusion list containing 1,997 unique  $m/z$  values. This method was run in triplicate.

344 Pass 3: Same method as Pass 2, except that the list of MS2 events generated in Pass 2 was  
345 uploaded in the Data Dependent Settings as a Reject Mass List. Like in Pass 2, a list of MS2  
346 event was produced by exporting the “Scan Filters” of the RAW file and to be used in Pass 4  
347 as an exclusion list containing 1,998 unique  $m/z$  values. This method was run in duplicate.

348 Pass 4: Same method as Pass 3, except that the list of MS2 events generated in Pass 3 was  
349 uploaded in the Data Dependent Settings as a Reject Mass List. This was the last exclusion list  
350 used in this study. This method was run in duplicate.

351 Pass 5: Same method as Pass 1, except that the threshold was dropped from 3,000 down to 500  
352 to perform MS2 on peptides of low abundance. This method was run in duplicate.

353 Pass 6: Same method as Pass 1, except with a Parent Mass List (i.e. an inclusion list) made out  
354 of the 2,000 most abundant peptides. This method was run in duplicate.

355 For the following methods, LC-MS1 reproducible peptides for which intensity exceeded  
356 0.0001 (19,956 peptides in total) were randomised along retention time (RT) and divided into  
357 10 lists (inclusion lists 1 to 10 containing <2,000  $m/z$  values each).

358 Pass 7: FTMS parameters were as specified above. Using the global MS/MSn method, MS/MS  
359 spectra were acquired in non-data dependent mode. ITMS parameters were: CID with a  
360 normalised collision energy of 35%, 0.25 activation Q, isolation width of 1. and activation time  
361 of 10 ms. Inclusion list 1 was uploaded in the inclusion global MS/MS mass list tab of the  
362 Global Non-Data Dependent Settings. All remaining nine parent lists were loaded to individual  
363 pass 7 methods.

364 Pass 8: FTMS parameters were as specified above. ITMS parameters were: CID with a  
365 normalised collision energy of 35%, 0.25 activation Q, and activation time of 10 ms. The  
366 precursor isolation width was 2  $m/z$ . The signal threshold was 500. Inclusion list 1 was  
367 uploaded in the parent mass list of the data-dependent settings. All remaining nine parent lists  
368 were loaded to individual pass 8 methods.

369 Pass 9: Same method as Pass 8, except that the precursor isolation width was 1  $m/z$  to increase  
370 the mass accuracy the  $m/z$  values targeted in the parent mass list. All remaining nine parent  
371 lists were loaded to individual pass 9 methods.

372 Pass 10: Same method as Pass 8, except that the precursor isolation width was 0.5 m/z to further  
373 increase the mass accuracy the m/z values targeted in the parent mass list. All remaining nine  
374 parent lists were loaded to individual pass 10 methods.

375 Pass 11: Same method as Pass 8, except that the precursor isolation width was 0.2 m/z to target  
376 the parent masses as accurately as possible. All remaining nine parent lists were loaded to  
377 individual pass 11 methods.

378 All the Xcalibur parameters of the various MS/MS methods can be found in Supplementary  
379 File SF1. Exclusion and inclusion lists can be found in Supplementary File SF2. A total of 63  
380 LC-MS2 files were thus acquired; they are available from the MassIVE repository  
381 (<https://massive.ucsd.edu/ProteoSAFe/static/massive.jsp>, MSV000090572).

## 382 **2.6. Proteomics data processing**

383 The LC-MS RAW files of the 4,061 wheat samples along with the 86 QC and IS replicates  
384 (injected once every 48 wheat samples) were processed in the Refiner MS module of Genedata  
385 Expressionist® 13.0 (Genedata AG, Basel, Switzerland. To process all files in one batch, a  
386 stepwise workflow was devised (Supplementary Figure S1A-B).

387 In the first step, a repetition activity was used (processing one file at a time) in which the  
388 consecutive sub-activities were performed: 1/ Load from File, 2/ RT Structure Removal with a  
389 minimum of 4 scans and m/z Structure Removal with a minimum of 8 points, 3/ Chromatogram  
390 Smoothing using a 3 scan RT window and a Moving Average estimator, 4/ RT Structure  
391 Removal with a minimum of 5 scans, and 5/ Save Snapshot to export all the processed files  
392 individually. The files were individually checked for inconsistencies that would invalidate the  
393 subsequent quantitative analyses. Inadequate files were removed from the dataset leaving 3,990  
394 reproducible wheat files. In the second step (Supplementary Figure S1C), the activities applied  
395 were: 1/ Load from File on the left for all the samples and on the right for the QCs, 2/ Adaptive  
396 Grid with 10 m/z scan counts, 3/ Average across Experiments (files) using the arithmetic mean,  
397 4/ Reference Grid joining both sides, 5/ Chromatogram RT Alignment applying a maximum  
398 RT shift of 50 scans (30 s), 6/ Chromatogram Peak Detection using a 12 scan Summation  
399 Window, Minimum Peak Size of 8 scans, Maximum Merge Distance of 5 points and  
400 Boundaries Merge Strategy, 10% Gap/Peak ratio for Peak RT Splitting, 3 points for m/z  
401 Smoothing, Ascent-based Peak Detection with 3 points Isolation Threshold, Local Maximum  
402 Centre Computation and Maximum Curvature Boundary Determination, 7/ Chromatogram  
403 Isotope Clustering with 0.1 min RT Tolerance and 20 ppm m/z Tolerance, the Peptide Isotope  
404 Shaping method with Protonation Ionisation, Minimum Charge of 2 and Maximum Charge of  
405 10, Maximum Log-Ratio Distance of 0.8, and Variable Charge Dependency for Cluster Size

406 Restriction, 8/ Singleton Filter, 9/ Metadata Import, 10/ Save Snapshot, and 11/ Export Analyst  
407 of the Clusters using the Integrated Maximum Intensity.

408 LC-MS processed quantitative data and metadata (sample description, LMA measurements,  
409 sample preparation technical steps, LC-MS sequence, instrument maintenance, etc...) were  
410 exported into Genedata Analyst (version 13, Genedata AG, Basel, Switzerland) for  
411 normalisation purposes (Supplementary Figure S1D). Data file normalisation with three  
412 consecutive steps was reported (41). In brief, first the quantities were normalised using the  
413 flour weights (1% accuracy) to account for sample preparation variation, second the IS cluster  
414 was used to normalise peptide abundances in order to take into consideration post-digestion  
415 technical variation, and third QCs and injection order were taken into account to correct  
416 instrument variation over time. The normalised quantitative data was exported as a CSV file  
417 for further processing. The CSV file contained 44,444 rows (peptide clusters) and 3,990  
418 columns (wheat samples).

419 The effects of technical biases on the LC-MS spectra were quantified using ANOVA  
420 simultaneous component analysis (ASCA), a generalisation of ANOVA which quantifies the  
421 variation induced by fixed experimental factors on complex multivariate datasets (43). Firstly,  
422 the normalised data were imported into R where clusters containing 100% missing values were  
423 removed ( $n = 12,108$ ), leaving 32,336 peptide clusters. The resulting dataset was a 3,990 x  
424 32,336 matrix with each row being an individual sample, and each column an LC-MS cluster.  
425 All remaining missing values were then imputed to a value zero. A separate metadata matrix  
426 (3990 x 4) which contained information on the technical conditions in the LC-MS run for each  
427 sample was compiled. These metadata were 1/ LC separation column – Categorical variable  
428 with 4 levels, 2/ Mass Calibration – Categorical variable with 6 levels, and 3/ Source heated  
429 capillary – categorical variable with 2 levels. A total of 3,090 samples had complete data (LC-  
430 MS spectra and corresponding metadata). This complete dataset was then analysed using  
431 ASCA in MatLab v.R2017b (Mathworks, Natick, WA, USA) utilising the PLS Toolbox v.  
432 8.5.2 (Eigenvector Research Inc., Manson, WA, USA) to see which, if any, of the fixed  
433 experimental effects had a significant impact on the LC-MS cluster data. The statistical  
434 significance of the impact of each fixed experimental effect was estimated by calculating a p-  
435 value from permutation testing with 100 iterations.

436 The impact of experimental factors with a significant effect on LC-MS cluster data was then  
437 accounted for by correcting the data using multiple linear regression in R (44) as described in  
438 (45). The linear model was fitted as follows:

439  $Y_{ijkl} = u + \text{Column}_i + \text{MassCal}_j + \text{Cap}_k + e_{ijkl}$



440 Where  $y$  is the signal intensity of a given cluster,  $u$  is the overall mean, Column is the  $i^{\text{th}}$  LC  
441 column (4 levels), MassCal is the  $j^{\text{th}}$  Mass calibration (6 levels), Cap is  $k^{\text{th}}$  Source heated  
442 capillary (2 levels), and  $e_{ijkl}$  is the random error term. The “corrected data” was a matrix of  
443 the residuals of the above model, which was run iteratively for each of the 32,336 peptide  
444 clusters. PCA plots were produced using R (44) and the gg2plot package.

## 445 **2.7. Protein identification**

446 The 63 RAW LC-MS2 files were processed in the Refiner MS module of Genedata  
447 Expressionist@ 13.0 using a stepwise workflow similar the one described for LC-MS1 data,  
448 with the exception of additional activities pertaining to protein database search (Supplementary  
449 Figure S2A-C).

450 RAW files were searched using Mascot program (version: 2.6.1, Matrix Science Ltd, London,  
451 UK) within Genedata Refiner. The wheat database searched was retrieved from three  
452 independent sources. The first source was UniProtKB  
453 (<https://www.uniprot.org/uniprot/?query=triticum%20aestivum&fil=organism%3A%22Triticum+aestivum+%28Wheat%29+%5B4565%5D%22&sort=score>) with 142,969 *T. aestivum*  
454 protein sequences (accessed on 26 February 2020, (41)). The second source was the  
455 EnsemblPlants repository hosting the *T. aestivum* genome initially sequenced by the  
456 International Wheat Genome Sequencing Consortium (IWGSC (9)) and containing 143,241  
457 Traes AA sequences ([http://ftp.ensemblgenomes.org/pub/plants/release-52/fasta/triticum\\_aestivum/pep/](http://ftp.ensemblgenomes.org/pub/plants/release-52/fasta/triticum_aestivum/pep/)). A contaminant database was also retrieved (common  
459 Repository of Adventitious Proteins (cRAP); <ftp://ftp.thegpm.org/fasta/cRAP>). All the FASTA  
460 files were combined and redundant sequences removed by following the GalaxyP tutorial  
461 “Protein FASTA Database Handling” ([https://training.galaxyproject.org/training-](https://training.galaxyproject.org/training-material/topics/proteomics/tutorials/database-handling/tutorial.html)  
462 [material/topics/proteomics/tutorials/database-handling/tutorial.html](https://training.galaxyproject.org/training-material/topics/proteomics/tutorials/database-handling/tutorial.html)) (46, 47). The decoy  
463 database was created by reversing all the sequences and appending them using the GalaxyP  
464 tool “DecoyDatabase” (<https://github.com/galaxyproteomics>). Our Galaxy workflow is  
465 available in Supplementary File SF1. The final FASTA file was imported and indexed in  
466 Mascot. It contained 286,482 protein sequences and 1,647,476,761 AA residues; its longest  
467 sequence bore 5,359 residues.

469 All MS2 files were searched in one batch using Mascot Daemon (version 2.6.1, Matrix Science  
470 Ltd, London, UK) and the following parameters: MS/MS ions search, Mascot generic data  
471 format, ESI-TRAP instrument, trypsin enzyme, 9 maximum missed cleavages,  
472 carbamidomethyl (C) as fixed modification, guanidyl (K) and oxidation (M) as variable  
473 modifications, quantitation none, monoisotopic mass, 2+, 3+ and 4+ peptide charge, 10 ppm



474 peptide tolerance, 0.5 Da MS/MS tolerance, and error tolerant search (Supplementary Figure  
475 S2D). Results were exported as .csv files into Excel.

476 The 32,336 peptide clusters from the corrected dataset produced by the LC-MS analyses were  
477 matched in R (44) (version 4.1.0-foss-2021a) to the 29,908 peptide clusters generated by the  
478 LC-MS/MS analyses using their respective RT,  $m/z$  and mass values with  $\pm 0.1$  accuracy, and  
479 then linked to the Mascot identification results. The identification results of the peptide clusters  
480 whose RT shifted by more than 1 min were not included.

481

## 482 **2.8. Statistical analyses of proteomics data**

483 Out of the 4,061 grains samples processed in this work, 3,990 yielded reproducible LC-MS  
484 data for 32,336 peptide clusters. The full quantitative data is available from the MassIVE  
485 repository (<https://massive.ucsd.edu/ProteoSAFe/static/massive.jsp>, MSV000090572). The  
486 corrected dataset with Mascot identification results were imported into Genedata Analyst  
487 (version 13, Genedata AG, Basel, Switzerland). LMA measurements were obtained on 3,773  
488 (out of 3,990) wheat samples. Whilst LMA trait characterised the wheat samples, we also  
489 wanted to analyse it along with the peptides to facilitate biomarker discovery. To this end, we  
490 used the inverse function to normally distribute the LMA values (Inv(LMA)) and transposed  
491 them as a row to incorporate them into the LC-MS dataset under the label “Cluster\_AAA”  
492 along with all the other 32,336 peptides, thus bringing the total number of clusters to 32,337.  
493 This “Cluster\_AAA” row was used in the subsequent statistical analyses to isolate peptides  
494 displaying profiles similar to that of LMA.

### 495 *2.8.1. Principal Component Analysis (PCA)*

496 A PCA was performed on the full dataset (3,990 samples x 32,336 peptides) in R using the  
497 `prcomp()` function of the `stats` package. The eigenvalues were plotted using the `screeplot()`  
498 function.

### 499 *2.8.2. Checking the distribution of LC-MSI data*

500 To redistribute data normally, the corrected dataset rows (peptides and Cluster\_AAA) were z-  
501 transformed and plotted as a histogram in R. The `hist()` function was used to plot the corrected  
502 and z-transformed dataset as histograms in R (version 4.1.0-foss-2021a). One-sample  
503 Kolmogorov-Smirnov tests were applied to check the normality of the distribution of both  
504 corrected and z-transformed datasets using the `ks.test()` function and “`pnorm`” argument in R.  
505 All the subsequent statistical analyses were performed on the z-transformed dataset.

### 506 *2.8.3. Subsampling wheat samples to eliminate the bias towards low LMA values*

507 LMA values spanned 0 to 8 u/g with the vast majority (95%) below 0.2 u/g (which corresponds  
508 to FN 300 s (18)); therefore, the LMA distribution was greatly skewed towards low LMA  
509 values. To eliminate this bias, a subset of wheat samples was selected as follows: all the  
510 samples bearing a  $LMA \geq 0.17$  were selected (467 samples in total) and an equivalent number  
511 of samples (467) with  $LMA < 0.17$  were randomly selected among the 3,306 remaining wheat  
512 samples. This subset of 934 wheat samples was no longer skewed towards low LMA values  
513 and is referred to as “unbiased samples” hereafter.

#### 514 *2.8.4. Partial Least Squares (PLS) to subset LMA-responding peptides*

515 In Genedata Analyst, a PLS 2-D plot was created using the 934 unbiased samples and all the  
516 32,346 peptides resolved in this study. The parameters were: LMA as a response, 3 latent  
517 factors, 10% valid values, and row mean imputation. Both score and loading plots were  
518 exported along with the variable importance in projection (VIP) scores. The higher the score,  
519 the greater the contribution of the peptide to the PLS and the closer to LMA response. These  
520 VIP scores were used to select meaningful subsets of peptides for the subsequent statistical  
521 analyses.

#### 522 *2.8.5. Univariate Partial Least Square (PLS) Regression to impute LMA missing values*

523 The missing LMA values were predicted using a univariate PLS regression model in Genedata  
524 Analyst. First a model was developed using the 934 unbiased samples and 2,996 peptides with  
525 PLS high VIP scores ( $> 1.5$ ). Second, among the 934 wheat samples, 179 were randomly  
526 chosen so that LMA evenly spanned 0 to 5 and those LMA values were erased. Several PLSR  
527 models were tested to accurately predict erased LMA values. The most accurate model applied  
528 the following parameters: LMA as a response, 20% valid values, and 20 latent factors. The  
529 model was then applied to the 217 missing LMA values against the 934 unbiased wheat  
530 samples.

#### 531 *2.8.6. Self-Organising Maps (SOM) Clustering*

532 In Genedata Analyst, a SOM was created using the 934 unbiased samples and 7,254 peptides  
533 with VIP scores above 1 (including Cluster\_AAA) and the following parameters: 6 rows, 8  
534 columns, positive correlation distance, 50 maximum iterations, and 10% valid values.

#### 535 *2.8.7. K-Means*

536 In Genedata Analyst, a k-means was performed using the 934 unbiased samples and 7,254  
537 peptides with VIP scores above 1 (including Cluster\_AAA) and the following parameters:  
538  $k=20$ , positive correlation distance, mean centroid calculation, 10% valid values, and 50  
539 maximum iterations.

#### 540 *2.8.8. Divisive Hierarchical Clustering Analysis (HCA) and agglomerative HCA*

541 An HCA was produced in Genedata Analyst a divisive HCA using the 934 unbiased samples  
542 and 7,254 peptides with VIP scores above 1 (including Cluster\_AAA) and the following  
543 parameters: clustering peptides, tree with tile plot, positive correlation distance, Ward linkage,  
544 10% valid values, k-means cluster profile, and split by size. The outcome of this analysis  
545 enabled us to sort the peptides based on their accumulation patterns in wheat samples.

546 Still in Genedata Analyst, we also performed an agglomerative HCA using the all the 934  
547 unbiased samples and 532 LMA-related biomarkers (including Cluster\_AAA) and the  
548 following parameters: clustering samples, tree, positive correlation distance, Ward linkage,  
549 50% valid values. The outcome of this analysis allowed us to sort the grain samples according  
550 to their LC-MS molecular similarity which was then exploited in a heat map.

#### 551 *2.8.9. Correlation*

552 An annotation correlation was performed in Genedata Analyst using the full dataset including  
553 Cluster\_AAA (3,990 samples x 32,337 peptides) against standardised LMA values. This  
554 produced R squared (R<sup>2</sup>) values.

#### 555 *2.8.10. Simple linear mixed regression*

556 The full dataset including Cluster\_AAA (3,990 samples x 32,337 peptides) was used to run a  
557 linear regression in Genedata Analyst with one explanatory variable using the following model:  
558  $y = \text{Inv}(\text{LMA}) + \varepsilon$ , in which  $\text{Inv}(\text{LMA})$  is the normal inverse function of LMA measurements.  
559 The false discovery rates were computed according to the Benjamini-Hochberg estimates as q-  
560 values.

#### 561 *2.8.11. Peptide expression profiles along 2 or 8 LMA bins*

562 Our data matrix of 3,990 columns by 32,337 rows contained 129,024,630 quantities which  
563 posed representation challenges. We adopted a data reduction strategy involving binning the  
564 samples into 8 or 2 arbitrary bins based on their LMA values in order to produce simpler more  
565 legible graphs for individual peptide profiling.

566 In the first instance, we sorted all 3,990 wheat samples based on an increasing order of LMA  
567 values, and then split them into 8 arbitrary bins of 499 samples each. The last bin  
568 ( $0.17132 < \text{LMA} < 7.95442$ ) contained all the 266 unsound grains ( $\text{LMA} > 0.2$ ).

569 In the second instance and using the 934 unbiased wheat samples, we created 2 bins based on  
570 LMA value threshold of 0.17. The bin containing 467 samples with  $\text{LMA} < 0.17$  only  
571 comprised sound grains. All the 266 unsound grains ( $\text{LMA} > 0.2$ ) were comprised in the bin  
572 containing 467 samples with  $\text{LMA} \geq 0.17$ .

573 The peptide quantities were then averaged per bin to produce mean expression profiles along  
574 2 or 8 bins.

575 *2.8.12. T test with effect size and volcano plot*

576 Using the unbiased biomarker dataset (934 samples x 532 peptides including Cluster\_AAA), a  
577 t test was performed with the LMA threshold of 0.17 as a factor (see sections 2.8.3 and 2.9.1)  
578 and the following parameters: bootstrap with 10 repeats and balanced permutations, effect size  
579 based on group means, and 90% valid values. This produced a volcano plot.

580

581 **2.9. Proteomics data mining**

582 The LC-MS2 experiments followed by Mascot search produced identification results for 5,414  
583 peptide clusters which matched 8,044 protein accessions. These identification results were  
584 mined using the databases and tools described below. Resulting outputs were consigned to  
585 Supplementary Tables S3.

586 *2.9.1. UniProt database and Gene Ontology (GO)*

587 The list of 8,044 UniProt accessions identified in this study was uploaded in the Retrieve/ID  
588 mapping tool of UniProt (<https://www.uniprot.org/uploadlists/> accessed on May 2022) (48) to  
589 retrieve protein descriptions, FASTA sequences, GO terms, and TRAES accession IDs. Out of  
590 the 8,044 UniProt accessions, 5,960 UniProt accessions corresponded to 6,622 TRAES  
591 accessions. TRAES accessions were needed to interrogate ShinyGO and BreadwheatCyc  
592 databases (described below).

593 *2.9.2. Kyoto Encyclopedia of Genes and Genomes (KEGG) database and pathway maps*

594 The 8,044 FASTA sequences were uploaded into the Assign KO tool  
595 ([https://www.kegg.jp/kegg/mapper/assign\\_ko.html](https://www.kegg.jp/kegg/mapper/assign_ko.html) accessed on May 2022) (49) by specifying  
596 the Poaceae family to retrieve KEGG ORTHOLOGY (KO) identifiers. KO identifiers were  
597 then mapped using the KEGG Mapper Reconstruct tool  
598 (<https://www.genome.jp/kegg/mapper/reconstruct.html> accessed on May 2022) to list  
599 pathways, Brites and modules involving identified proteins.

600 *2.9.3. ShinyGO, Functional Category enrichment and chromosomal positions*

601 The list of 6,622 TRAES accessions was uploaded into ShinyGO  
602 (<http://bioinformatics.sdstate.edu/go/>) (50) to generate Functional Category enrichments, dot  
603 plots, tree, networks, as well as retrieve chromosomal positions. Positions were obtained for  
604 4,571 TRAES accessions which were used in Circos plots (detailed below).

605 *2.9.4. Pathway Tools, BreadwheatCyc and perturbed pathways*

606 The list of 6,622 TRAES accessions along with quantitative data along 8 bins was uploaded  
607 into the Pathway Tools software (51) and run online via the BreadwheatCyc database  
608 (<https://pmn.plantcyc.org/organism-summary?object=BREADWHEAT> accessed on June

609 2022) via the Plant Metabolic Network server (52) using the Omics Dashboard  
610 (<https://pmn.plantcyc.org/dashboard/dashboard-intro.shtml> accessed on June 2022), the  
611 Cellular Overview tools  
612 (<https://pmn.plantcyc.org/overviewsWeb/celOv.shtml?orgid=BREADWHEAT> accessed on  
613 June 2022) to generate Pathway Perturbation Scores (PPS).

614 The Chrome extension Veed.io was used to create a film capturing the Cellular Overview  
615 animation.

#### 616 *2.9.5. Circos and chromosomal position*

617 The 4,571 TRAES accessions whose chromosomal positions were known from ShinyGO were  
618 charted along a Circos plot invented by Krzywinski and colleagues (53) and recently wrapped  
619 in the Galaxy platform by Rasche and colleagues ([https://usegalaxy.eu/?tool\\_id=circos](https://usegalaxy.eu/?tool_id=circos)) (46,  
620 54, 55). The details of the various layers are indicated in the figure's legend.

#### 621 *2.9.6. R and Power BI Desktop*

622 Most identified peptide matched several UniProt accessions which corresponded to several  
623 TRAES IDs, and GO terms. This produced wide tables. In R (version 4.1.0-foss-2021a) (44),  
624 wide tables were converted to long tables using the `pivot_longer()` function from `tidyr` package.  
625 Long tables were merged using the `merge()` function of the R base package using peptide  
626 Cluster IDs as unique references.

627 Wheat sample metadata, peptide metadata and quantitative dataset and identities for the  
628 biomarkers were imported into Microsoft Power BI Desktop (Version: 2.106.883.0 64-bit June  
629 2022) and linked via the Clusters names to produce dashboards using multiple visuals (word  
630 clouds, tree maps, histograms, scatterplots, waterfall plots, pie charts, violin plots and ribbon  
631 charts).

632

### 633 **3. Results and Discussion**

#### 634 **3.1. Resources for scientific studies on wheat**

##### 635 *3.1.1. Wheat resources*

636 A total 858 wheat genotypes, sourced from all over the world, grown over 8 years since 2012  
637 and stored in optimal conditions amounting to 4,061 grain samples were analysed in this work  
638 (Supplementary Table S1). Because LMA measurements occurred simultaneously to the  
639 proteomics analyses in 2019, we did not consider storage time for the statistics. We also did  
640 not statistically test for varietal differences which was outside the focus of this study.

641

642 *3.1.2. High-throughput proteomics workflow to efficiently process and analyse thousands of*  
643 *samples*

644 We have developed a high-throughput proteomics LC-MS method (41) that was applied to  
645 4,061 wheat grain samples following the workflow described in Figure 1. The technical aspects  
646 pertaining to sample preparation/tracking and data acquisition steps that ensured a high-  
647 throughput workflow are available in Supplementary File SF1. Overall, the LC-MS continuous  
648 run lasted for 143 days (20.4 weeks or 4.5 months) and included regular system maintenance  
649 (mass calibration, source cleaning, HPLC column swapping). A total of 4,370 RAW files were  
650 acquired. A Gantt chart illustrates the timeline of the workflow steps along with data  
651 accumulation (Figure 2).

652 The wet experiment bottlenecks were resolved where possible as explained in (41). Most time  
653 was spent grinding, transferring, weighing and extracting the samples as there was no option  
654 to greatly up-scale those steps (Figure 2). The workflow became much faster when 96-well  
655 plates were introduced (from digestion step onward) allowing for high throughput  
656 multipipetting and multidispensing activities, as well as minimising the footprint of sample  
657 freezer storage. Although steps were sequential, they could overlap with two experimenters  
658 operating in a staggered fashion from one lab workstation to the next.

659 LC-MS1 acquisition started when enough plates were ready to ensure continuous instrument  
660 run while samples processing was still happening. Data acquisition was completed 18 days  
661 after the last wheat sample was fully processed, demonstrating minimum time loss (Figure 2).  
662 The Genedata Refiner workflow used to process LC-MS1 data was previously optimised (41)  
663 (Supplementary Figure S1 described in section 3.1.2); its first step was applied to batches of  
664 ~200 LC-MS1 files. The time limiting factor was the server computing ability.

665 Overall, all 4,061 wheat samples were processed and analysed (from receiving the samples to  
666 processing the LC-MS1 data) in 334 days (~11 months). Purchasing all required consumable  
667 ahead, keeping track of the samples, good logistics by setting up working stations for each wet  
668 lab step, as well as overlapping activities across experimenters guaranteed efficient time  
669 management. Stowing samples in the freezer in-between steps allowed to safely interrupt the  
670 sample preparation procedure to accommodate equipment/experimenter downtime without  
671 compromising the quality of the samples processed so far.

672 The subsequent steps had to follow one another. LC-MS2 acquisition necessitated LC-MS1  
673 data processing to be finished in order to produce parent mass lists and consequently had to be  
674 performed post-hoc. Whilst LC-MS2 acquisition was rapid (2 weeks), its processing took  
675 longer (3 months) because it required another Genedata Refiner workflow (Supplementary



676 Figure S2 described in section 3.1.3), a more recent non redundant database with decoy  
677 sequences, testing several Mascot parameters (data not shown), and linking LC-MS2 clusters  
678 to LC-MS1 clusters (data not shown).

679 The final bottleneck in the workflow pertained to statistical analyses and data mining (8  
680 months) which necessitated trying different statistical methods with multiple trial and error  
681 stages working out optimal parameters, testing and using different data mining tools which  
682 required training and a lot of strategising on how best to present big data. Running such large  
683 datasets proved computationally taxing, necessitated extensive dwell times; it often ran out of  
684 memory and triggered server crashes.

685 One way to increase the throughput and therefore shrink the timeline would be to use an  
686 automated sample preparation station. A robot (Bravo Automated Liquid Handling Platform  
687 from Agilent) was used to automate peptide clean-up and phosphopeptide enrichment from  
688 wheat and maize vegetative samples (56). We could not find any other high throughput method  
689 in wheat or cereals.

690

691 *3.1.2. LC-MS1 quantitative data processing, normalisation, correction and standardisation to*  
692 *remove technical biases*

693 The Genedata Refiner workflow described in (41) was applied to 4,147 LC-MS1 files (4,061  
694 wheat + 86 QCs; Supplementary Figure S1). Step 1 covered noise subtraction nodes that could  
695 be run on individual data file. It was performed throughout LC-MS1 acquisition activity on  
696 weekly batches (~230 files) to optimise server dwell time. Step 1 helped assess data  
697 reproducibility and non-reproducible files (71 samples) were omitted from the remainder of  
698 the processing, leaving 3,990 wheat and 86 QC data files. Step 2 encapsulated all alignment,  
699 peak detection and quantitation, as well as isotope clustering and singleton filtering activities.  
700 This step had to be performed on all 4,076 reproducible data files simultaneously and therefore  
701 could only be attempted when the LC-MS1 run was finalised. The experiment metadata  
702 captured in Excel was associated to the quantitative data and exported to Genedata Analyst for  
703 data normalisation purpose.

704 The data was normalised as described in (41) following three steps: using flour weights, IS  
705 cluster and QC replicates along with LC-MS injection order (Figure 3).

706 Raw data displayed a clear sample grouping based on injection order during the LC-MS1 run  
707 (Figure 3A) and mirrored the instrument maintenance events (mass calibration, etc...). Two  
708 large groups appeared that could not be explained by any experimental steps. Normalising  
709 using flour weight accuracy of 1% helped creating tighter wheat sample groups with four

710 outliers, and isolated QCs (Figure 3B). The two larger groups of samples were less distinct.  
711 This first normalisation step did not significantly impact the peptide distribution as can be seen  
712 on the PCA loading plots (Figure S3G,H). Normalising against the IS shifted the sample groups  
713 around but did not combine or homogenise them (Figure 3C). The two larger sample groups  
714 observed in panels A-B became indistinguishable in panel C. This normalisation step also  
715 affected peptide distribution assuming a more oval shape on the loading plot (Figure S3I). The  
716 final normalisation step further scattered the samples more widely across the PCA plot and  
717 accentuated the technical variation gradually expanding overtime during the instrument run  
718 (Figure 3D). Yet at the peptide level, this last normalisation activity further shrunk the grouping  
719 assuming a more circular distribution with less outliers (Figure S3J). The benefits of  
720 normalisation were discussed before (41) with respect to precise sample weights mandated in  
721 metabolomics (57), spiking IS post-digestion to alleviate for sample to sample variations (58,  
722 59), and QCs to account for batch differences over time and minimise cross run effects (59-  
723 61). In their ground-breaking study to assess and ameliorate the reproducibility of large-scale  
724 proteomics experiments, Poulos and colleagues have highlighted the decrease over time in  
725 mass analyser sensitivity in-between cleaning events and how technical replicates, such as  
726 QCs, help remove unwanted variation (62). Despite all the normalisation steps applied to our  
727 data, not all technical biases could be removed, thus necessitating further data correction.  
728 The fully normalised dataset of 3,990 wheat samples and 32,336 reproducible peptides was  
729 exported as a CSV file and imported into R to run a linear model fitting the technical factors  
730 that bore the greatest variance and were associated with LC-MS maintenance. The  
731 experimental variation was successfully eradicated as illustrated by PCA (Figure 3E,K). The  
732 results showed that while instrument mass calibration had a much bigger effect, all three  
733 technical factors had a significant effect ( $P < 0.05$  based on permutation testing with 100  
734 iterations) on the spectral data (data not shown). This correction method was initially developed  
735 in a metabolomics study to account for uncontrollable environmental effects (45). Quantitative  
736 geneticists routinely exploit linear models to measure the influence of systematic  
737 environmental effects (fixed effects) which impact phenotypic variation and unscramble  
738 genetic from non-genetic factors (63). To our knowledge, this is the first time such correction  
739 method was applied to proteomics data.  
740 The final data transformation step involved a z-transformation (scaling and centring) to level  
741 out extreme quantities and facilitate the comparison and clustering of peptide profiles during  
742 statistical analyses. Finding linear combinations of predictors based on how much variation  
743 they explain is achieved by centring to a mean of 0 and scaling to a standard deviation of 1

744 (64). Such mathematical transformation is common practice in post-genomics expression  
745 studies, and MS-proteomics is no exception (65, 66). In our study, z-transformation radically  
746 modified the data from an homogenous plot to defined groups stretching in four main directions  
747 (Figure 3F,L), which could not be attributed to any of our metadata. Peptide quantities that  
748 originally ranged from 0 to  $1 \times 10^7$  ultimately spanned a mere -22 to 63 scale.

749

750 *3.1.3. A non-redundant wheat database to annotate LC-MS2 results and identify post-*  
751 *translational modifications (PTMs)*

752 A *T. aestivum* database was created by combining all the protein sequences publicly available  
753 from UniProt and IWGSC EnsemblPlants repositories. Because protein annotations from the  
754 IWGSC (hereafter called TRAES sequences) referred to UniProt, we used the latter as a  
755 template to eliminate AA sequence redundancy. This completely removed all IWGSC TRAES  
756 sequences (data not shown) from our merged data file indicating they were all included in the  
757 UniProt repository. The database was reversed to create a decoy database which was then  
758 concatenated to the latter. This way, not only a single file has to be interrogated in Mascot  
759 system, but also false positives are only recorded when a match from the decoy sequences  
760 exceeds any match from the target sequences (67, 68). All LC-MS2 files were processed in  
761 Genedata Refiner and searched using the Mascot algorithm with an error tolerant search to  
762 maximise PTM discovery. The search outputs were merged into a single file and exported to  
763 Excel (Supplementary Figure S2E).

764 Our strategy to quickly identify as many peptides as possible was to multiply the number of  
765 data-dependent LC-MS2 methods rather than multiplying the number of samples analysed. We  
766 thus pooled 10% of the wheat samples randomly chosen into one tube and subjected this pooled  
767 sample to 11 methods (passes) with replicates, varied ITMS parameters and 10 unique parent  
768 lists of 2,000 ions each. Each method had a drastic impact of the selection of precursor ion,  
769 with some areas being thoroughly samples whilst others were ignored (Supplementary Figure  
770 S3).

771 A total of 63 LC-MS2 files were thus obtained. The LC-MS2 methods varied in their  
772 efficiencies, identifying as few as 104 peptides (pass 7) up to 11,662 peptides (pass 8),  
773 irrespective of the number of MS2 events (Supplementary Figure S4).

774 Passes 8-10 yielded by far the largest identity counts across all 10 parent lists, even though  
775 they did not feature the highest MS2 event counts (Supplementary Figure S4). Key MS  
776 parameters to maximise peptide identifications were the inclusion of the parent lists into the  
777 data-dependent settings (passes 8-11) albeit not the at the global level (pass 7) as well as

778 allowing for wider mass tolerance window during precursor selection. The widest tolerance (2  
779 m/z) achieved the greatest counts (pass 8, Supplementary Figure S4). Overall, a total of  
780 315,934 peptides were identified, comprising only 6,550 unique peptides which matched  
781 10,437 unique wheat proteins, 277 decoy accessions, and 3 contaminant proteins. The huge  
782 peptide redundancy was explained by the fact that a single pooled sample (from 400 individual  
783 samples) was repeatedly analysed using various LC-MS2 methods. Pooling digests erased  
784 sample-to-sample variation. More protein identities could have been realised with a diverse  
785 sample set subject to all the methods developed here but that would have extended the data  
786 acquisition, analysis and mining by many more months. A greater proteome coverage was  
787 achieved in our method optimisation study yielding 13,165 identified peptides even though far  
788 less samples were analysed because two extraction protocols and three orthogonal digestions  
789 were applied which produced more diverse LC-MS profiles (41). An array of strategies can be  
790 employed to increase the proteome coverage of plant seeds, including depletion and pre-  
791 fractionation strategies as well as exploring different organs, developmental stages, and cell  
792 cultures (69, 70). However, these additional experimental steps are time-consuming, labour-  
793 intensive, as well as costly thus unsuitable for large-scale high-throughput experiments like  
794 ours. Our strategy was first to rapidly and reproducibly quantify digested peptides from  
795 thousands of wheat samples using a label-free LC-MS approach and apply robust statistical  
796 analyses to detect potential trait-related biomarkers, and second to quickly identify as many  
797 peptides as possible using LC-MS2. Large-scale proteomics studies have been applied to  
798 human (71); to our knowledge, this is the largest plant proteomics study carried out to date.  
799 In this study, we opted for an error-tolerant search which accrued a plethora of modifications  
800 (Supplementary Table S2). A total of 21,486 carbamidomethylations of Cys residues were  
801 identified as fixed modifications. This was expected to occur during our denaturing protein  
802 extraction procedure. The most prevalent dynamic modifications were non-specific cleavages  
803 (5,480), followed by N-terminal ammonia losses (907), and conversion from N-terminal Gln  
804 to pyroGlu (815). During the digestion process involving trypsin, proteomics studies have often  
805 reported the formation of semi-tryptic and non-specific peptides besides cleavages after Arg or  
806 Lys residues (72). Therefore, some of our non-specific peptides could have resulted from the  
807 digestion step, but we cannot rule out that non-tryptic peptides were naturally present on our  
808 stored grains, resulting from residual enzymatic activities. Ammonia losses are neutral losses  
809 commonly triggered by CID upon creating b and y ions, and can be detected by high resolution  
810 mass analysers such as FTMS instruments (73). C-terminal Arg or Lys of tryptic peptides often  
811 leads to abundant y ions with ammonia loss (74) and as well as b ions specific enough to detect

812 the presence of Gln, Asn, His, Lys, and Arg residues (73). PyroGlu formation is a common  
813 cyclization side reaction of Glu and/or Gln residues in peptides and proteins that occurs when  
814 those residue are located at the N-terminus and under slightly acidic conditions (75), such as  
815 our experimental conditions therefore this PTM could also be a process artifact. Other frequent  
816 PTMs in our study were N-terminal ethylation (265 occurrences), deamidation (147  
817 occurrences), guanidylation (141 occurrences), the latter of which could have been triggered  
818 during protein resuspension in Guanidine-HCl solution as discussed in (41), as well as  
819 oxidation of Met (100 occurrences) (Supplementary Table S2). Numerous PTMs have been  
820 identified in plants (69) and cereals in particular (76), including barley (77), and wheat (2, 78,  
821 79). Deamidations of glutamine residues in glutenins have been reported (5), along with C-  
822 terminal loss of tyrosine potentially facilitating protein sorting during seed maturation (2).  
823 Starch content and starch-related proteins are prominent in wheat grain; PTMs involved in  
824 starch quality have been reviewed (80). Our study lists numerous potential PTMs; this warrants  
825 more experiments to validate them and decipher their role in LMA response. Future proteomics  
826 experiments should endeavour to explore the relationship between structure and functionality  
827 of gluten proteoforms arising from key PTMs in response to LMA phenotype.

828

#### 829 *3.1.4. Linking LC-MS1 and LC-MS2 data to annotate quantities with identities*

830 LC-MS1 files resolved 32,336 reproducible clusters which had to be matched to 29,908 clusters  
831 from LC-MS2 data files. Using tolerances of 20 ppm for m/z and mass and 1 min for retention  
832 times, 16,874 (52%) peptide clusters were matched across both datasets, of which 5,414 bore  
833 peptide identification results. These identified peptides matched 8,044 *T. aestivum* protein  
834 accessions. Our experimental results are summarised in Table 1; number of identified peptide  
835 numbers aside, they compared well with our previous findings during method optimisation  
836 (41).

837

838 **Table 1: Experiment summary.**

Items quantified	Occurrences
Number of wheat genotypes	858
Number of wheat samples	4061
Sampling years	8 (2012-2019)
Trait (LMA)	1
Digestion types	1
Number of reproducible LC-MS1 files	3990
Number of LC-MS1 peaks	137669
Number of reproducible LC-MS1 clusters	32336
Cluster size range	2 - 10
Cluster charge range	2 - 7
Cluster m/z range	300.13 - 1921.55
Cluster mass range	598.26 - 6527.06
Base peak range	120 - 520083
Number of clusters with peptide identity	5414
Number of identified accessions	8044
Range of peptides/accession	1 - 64
Range of accessions/peptide	1 - 212

839

840 Our strategy was to consider all 8,044 protein hits identified from the 5,414 sequenced peptides  
841 irrespective of their homology. We thus turned the 5,414 x 212 wide table into a long table  
842 containing 32,347 rows of peptides and replicated the quantitative data accordingly for  
843 statistical analysis purposes. The list of all identities is captured in Supplementary Table S3.  
844 Up to 64 unique peptides matched a particular protein with an average of 4 peptides per hit  
845 (Supplementary Figure S5A-B).

846 A given peptide matched to up to 212 protein accessions with an average of 6 hits per peptide  
847 (Cluster\_29452, VLQQLNPCK, Supplementary Figure S5C-D). This mirrored the high  
848 frequency of homoeologous proteins in the hexaploid wheat samples expressed from three  
849 similar subgenomes, A, B and D (81). Another compounding factor was that wheat protein  
850 accessions were created from genomic sequences, resulting in multiple accessions bearing  
851 identical sequences but arising from different gene accessions (2). This created on one hand  
852 protein accessions labelled as “fragments” despite having a complete coding region and, on the  
853 other hand, other accessions lacking this tag despite having an incomplete coding region  
854 (Supplementary Table S3). Finally, the vast number of PTMs identified here also contributed



855 to boosting hits against a particular peptide AA sequence. The most dominant wheat grain  
856 proteins are storage proteins such as gliadins and glutenins, which featured prominently in our  
857 proteome (Supplementary Figure S5E, Supplementary Table S3), despite the fact that their low  
858 Lys/Arg content makes them less prone to trypsin digestion (2). Other major proteins  
859 comprised histones, beta-D-glucosidases, and ubiquitin. This list of identified proteins  
860 compared well with our previous methodological work (41). Other recent studies on mature  
861 wheat seed proteome using gel-based or gel-free technologies also published comparable list  
862 of identities (82-84).

863

### 864 **3.2. Application to a wheat industry problem: Late maturity alpha-amylase (LMA)**

865 Wheat marketing for milling grades dictates that below a certain FN value, grains are no longer  
866 suitable for human diet and must then be discounted causing significant financial losses to the  
867 suppliers (17). FN assesses starch degradation resulting from LMA activity which can be  
868 assayed in flour samples using the Ceralpha method (37) for instance. Even though LMA trait  
869 is a genetic defect, it persists in wheat germplasm implying that it is either not selected against  
870 or alternatively imparts unbeknown beneficial attributes to LMA-prone varieties (24). By  
871 unravelling the genetic, biochemical, and physiological mechanisms that lead to LMA  
872 expression, scientists strive to understand and eliminate LMA from wheat breeding programs  
873 (39). Surprisingly, post-genomics is not one of the strategies adopted by researchers to close  
874 the biological knowledge gap, with only one transcriptomics study registered so far (22). Our  
875 study constitutes the first proteomics experiment performed to decipher the mechanisms  
876 involved. Machine learning was performed on the complete dataset to distinguish LMA-  
877 susceptible from non-susceptible wheat genotypes without success (data not shown). Results  
878 from statistics and data mining are described and discussed below.

#### 879 *3.2.1. Getting the quantitative data ready for statistical analyses*

##### 880 *3.2.1.1. Assessing the normality of LC-MS1 datasets*

881 To assess whether our LC-MS1 datasets following the correction and z-transformation steps  
882 was normally distributed, we plotted the data as histogram and boxplot. We further performed  
883 the nonparametric one-sample Kolmogorov-Smirnov (K-S) test (85) well suited to analysing  
884 big data (86). Both histogram and boxplot of the corrected data were asymmetrical with most  
885 values being on the low range (Supplementary Figure 6A-B), which revealed that this dataset  
886 was not normally distributed. This was confirmed by the high K-S statistics (D) of 0.41 and a  
887 very low p-value ( $< 2.2 \times 10^{-16}$ ).

888 Using the z-transformed data, the histogram and boxplot were more symmetrical  
889 (Supplementary Figure 6C-D). Whilst the K-S statistics (D) was reduced to 0.27, it was still  
890 too high to conclude to normality. Even though we did not achieve a gaussian distribution by  
891 standardising the data, we managed to make it more even which improved statistical analyses  
892 for biomarker discovery.

893

#### 894 3.2.1.2. PLS of unbiased samples to select a meaningful set of LMA-responsive peptides

895 Analysing such a large dataset (3,990 columns x 32,337 rows) was computationally taxing,  
896 necessitating extensive dwell times to finalise statistical analyses, and often triggering  
897 Genedata sever crashes due to out-of-memory failures despite recent upgrades. Consequently,  
898 we devised a strategy to select a subset of relevant peptides via the supervised cluster method  
899 PLS. Using the 934 unbiased samples and all 32,337 peptides (including Cluster\_AAA), we  
900 executed a PLS analysis with LMA trait as a response. The score plot of the first two  
901 components showed that the PLS successfully pulled out the grain samples exhibiting high  
902 LMA activities (Supplementary Figure S7A).

903 The corresponding loading plot allowed us to categorise peptides according to their  
904 contribution to the PLS model via their Variable Importance in Projection (VIP) scores. The  
905 most-contributing peptides (i.e. exhibiting the highest VIP score) were located in the plot area  
906 equivalent to that of high LMA samples (Supplementary Figure S7B).

907 VIP scores indicated the importance of each variable (peptide) in the projection used in the  
908 PLS model. Peptide VIP scores were calculated as weighted sums of the squared correlations  
909 between the PLS components and the original peptides; weights were inferred from the  
910 percentage variation explained by the PLS component in the model (87). VIP scores greater  
911 than 0.5, 1.0, and 1.5 segregated 14,440 (45%), 7,252 (22%), and 2,996 (9%) peptides,  
912 respectively. By setting up three VIP score thresholds of increasing stringency, we thus created  
913 three subsets of peptides of decreasing sizes that could be used in more computationally  
914 demanding processes.

915

#### 916 3.2.1.3. Wheat subsampling to create an unbiased dataset and transforming LMA trait profile 917 to achieve normal distribution

918 In the 3,990 reproducible wheat samples, 3,773 featured LMA measurements that ranged from  
919 0.04 to 7.95 u/g (Supplementary Table S1), albeit mostly on the low scale with 88% of the  
920 values recording less than 0.2 u/g (Figure 4A), which corresponds to the receival threshold of  
921 FN 300 s (18, 23).

922 Our range far exceeded those reported earlier, spanning either 0.08 to 0.67 u/g across 33 spring  
923 wheat cultivars grown across 18 field sites (88), 0.023 to 1.417 u/g over 39 varieties grown  
924 under controlled and triggering LMA-conditions (23), or 0.002 to 1.977 u/g among 196  
925 genotypes from three experimental locations (19). We chose a threshold of 0.17 as a tipping  
926 point to delineate between grain samples displaying either low (3,306 samples) or high (467  
927 samples) alpha-amylase activity. The LMA profiles below and above this arbitrary value  
928 showed a slow gradual increase of enzyme activity up to 3.2 units where datapoints became  
929 more scattered (Figure 4B-C). Because the LMA distribution was significantly skewed towards  
930 low values and to restore balance to the trait profile, we retained all the wheat samples with an  
931 LMA above 0.17 (467 samples) and randomly selected 467 samples (out of 3,306) for which  
932 LMA fell below this threshold. The LMA profile of this unbiased subset of 934 samples (Figure  
933 4D) was very similar to the complete distribution (Figure 4A).

934 When LMA measurements were plotted as a histogram, it confirmed the skewness towards low  
935 activities and highlighted that most values fell between 0.068 and 0.203 u/g (Figure 4E). A  
936 natural logarithm transformation did not make the data gaussian (Figure 4F); nor did other  
937 logarithmic bases (data not shown). A binary logarithm function was used to transform LMA  
938 data to ascertain the significant negative correlation with Falling Numbers (FN) (19, 23). FNs  
939 inferior to 300 sec, which is the commercial trade cut-off manifesting significant alpha-amylase  
940 activity, corresponded to log<sub>2</sub> LMA value of -3 (23). In our work, an inverse function normally  
941 distributed LMA values, albeit as a slightly asymmetrical bell curve (Figure 4G). This  
942 INV(LMA) data was further standardised (centred around zero and scaled down to comparable  
943 variance) when it was incorporated at the peptide level which did not compromise its gaussian  
944 distribution (Figure 4H).

945

#### 946 3.2.1.4. Predicting LMA missing values

947 Out of the 3,990 reproducibly processed grain samples, 217 were not measured for LMA. We  
948 employed a univariate PLS regression strategy to impute them. Using our 2,996 peptide set  
949 with the highest VIP scores (see section 3.2.1.2), we tested various PLS regression models  
950 (data not shown) using a random selection of 179 samples out of the 934 unbiased sample set  
951 which ranged from 0.5 to 4.9. This testing set was analysed against the remainder of the  
952 unbiased set (755 samples). The best regression model utilised 20% of the valid values and 20  
953 latent factors; it predicted the 179 tested values with 93% accuracy (Supplementary Figure  
954 S8A).

955 This model was not accurate for small LMA values with a  $R^2$  of 6%, even imputing negative  
956 values (Supplementary Figure S8B). Yet, it was 98% accurate for LMA measurements greater  
957 than 0.17 u/g (Supplementary Figure S8C). It was more critical to faithfully estimate high LMA  
958 values given that it was the criteria for grain soundness; our PLS regression (PLSR) model  
959 fulfilled this. We applied the model's parameters to predict the 217 LMA missing values  
960 against the unbiased set of 934 samples; the imputations ranged from -0.29 to 0.63 u/g  
961 (Supplementary Figure S8D). The negative values were converted to zeros. LMA predictions  
962 are reported in Supplementary Table S1.

963 The simplest method for imputing missing data relied on single value imputation, such as the  
964 mean (89), whilst more complex methods were based on regression (90) or K-Nearest  
965 Neighbours (KNN) which estimates a missing data point using distances calculated from its  
966 most similar neighbours (91). Invented in 1966 (92), PLS regression has become very popular  
967 notably in the fields of bioinformatics (93) and spectroscopy (94). Nengsih and colleagues  
968 demonstrated that while computation times increased with the proportion of missing data, up  
969 to 30% missing values could be imputed using PLSR (95). In our study, LMA was the single  
970 trait provided to analyse LC-MS1 data. Not imputing missing LMA measurements meant that  
971 5.4% (217/3,990 samples) of our dataset would have been useless, therefore it was a  
972 worthwhile effort. Along with PLSR, we have also tested multivariate linear regression (MLR),  
973 univariate polynomial regression and KNN imputation by varying several parameters including  
974 valid value percentage, number of latent factors, number of parameters (for MLR), as well as  
975 distance computation and number of K (for KNN), albeit without success (data not shown).

976

#### 977 3.2.1.5. Incorporating LMA trait at the peptide level for biomarker discovery

978 Because we only had a single trait to make biological sense of our big data, we introduced all  
979 3,990 LMA values (including the predicted values) which characterised wheat samples at the  
980 peptide level by transposing it and renaming "Cluster\_AAA". This added one extra row to our  
981 dataset of 32,336 peptides to make a final matrix of 3,990 columns (wheat samples) and 32,337  
982 rows. This way, we could apply statistical analyses that would group peptides that behaved  
983 similarly or conversely to our LMA trait thereby facilitating biomarker discovery. To permit  
984 the comparison between LMA and grain peptides, we first needed to normalise and standardise  
985 LMA values, as detailed above in section 3.2.1.3, prior to their transposition.

986 Having LMA incorporated with wheat grain peptides (as Cluster\_AAA) further helped us  
987 assess the relevance of the statistical tests carried out by validating anticipated results. For  
988 instance, when performing a correlation analysis with LMA, as expected Cluster\_AAA

989 achieved a positive correlation of 1. In another instance, when executing a one factor linear  
990 model with LMA as a covariate, Cluster\_AAA was confirmed to yield a q-value of 0. Finally,  
991 when performing multivariate clustering analyses (HCA, SOM, k-means), this strategy assisted  
992 us in finding peptides with profiles similar to that of Cluster\_AAA.

993

### 994 *3.2.2. Statistical analyses to discover LMA-responsive biomarkers*

995 Big data produced by gene expression studies are too large to analyse by mere sorting in  
996 spreadsheets or plotting on few charts. Multivariate data analyses such as clustering and  
997 correlating methods are required to make sense of the data (96, 97). Yet, as helpful these  
998 multivariate analyses are, they are not as statistically robust as uni- or bivariate analyses (96)  
999 to test the relationship between peptides and LMA. We thus performed a few uni-, bi- and  
1000 multivariate analyses to explore our large dataset against our single LMA trait.

#### 1001 *3.2.2.1. Unsupervised multivariate clustering analyses (SOM, k-means, HCA) for pattern* 1002 *recognition and peptide profiling of LMA phenotype*

1003 As multivariate analyses handle integral datasets and iteratively impute many statistics, they  
1004 incur heavy computational costs. Suffering multiple Genedata server crashes, we could only  
1005 apply such methods to a subset of our data. Using the unbiased set of 934 wheat samples and  
1006 the list of 7,254 peptides with LMA-responsive VIP scores above 1 (see section 3.2.1.2), we  
1007 have performed three unsupervised clustering analyses, SOM, k-means and divisive HCA.  
1008 Because we had incorporated the LMA trait at the peptide level as Cluster\_AAA, we could  
1009 look for groups resulting from these analyses which assembled peptides behaving similarly to  
1010 Cluster\_AAA. Clustering or cluster analysis corresponds to a set of learning methods grouping  
1011 observations that share similar characteristics. Within a set of related values of the variables  
1012 analysed, these methods find feature patterns which generate clusters that group similar  
1013 observations (98). Unsupervised clustering analyses are commonly employed in gene  
1014 expression studies (97).

1015 In our experiment, the SOM model yielded 48 groups comprising 8 to 555 peptides with mean  
1016 distances from 0.09 to 0.80. The group including Cluster\_AAA (4,3) contained 26 biomarker  
1017 peptides; its distance from the group centre ranged from 0.00-0.83 with a mean of 0.38 and a  
1018 SD of 0.31 (Supplementary Table S4). Cluster\_AAA stood 0.70 from the group centre. While  
1019 SOM has been widely used in exploratory data analyses in diverse fields (99), it has only been  
1020 applied to proteomics in the context of animal cell culture (100), GPI anchor prediction (101),  
1021 transmembrane helix predictor (102) protein conformation (103) or protein-protein interaction  
1022 (104), never in plant grains.



1023 We tested different number of neighbours (k) and observed that the larger k the greater the  
1024 variance explained by the k-means model (data not shown). Applying the biggest k possible  
1025 (20) produced a model that overall explained 71.1% of the variance. Group 14 with a variance  
1026 of 35% contained 93 biomarker peptides spanning a distance of 0.12 to 0.94, including  
1027 Cluster\_AAA whose distance was 0.79 (Supplementary Table S4). K-means clustering was  
1028 well adopted by the proteomics community to group gene products of similar profiles, notably  
1029 in plants such as bamboo (105), nightshade (106), or grape (107), but to our knowledge not in  
1030 wheat. In developing corn grains, coordinated protein expression associated with different  
1031 functional categories was revealed by a k-means clustering analysis (108).

1032 We successfully applied an agglomerative 2-D HCA to cluster both samples and peptides (data  
1033 not shown) but failed to select individual groups to retain the one hosting Cluster\_AAA.  
1034 Instead, we performed a divisive HCA which ordered the peptides into clusters that could then  
1035 be chosen individually. Cluster\_AAA belonged to a group of 33 biomarker peptides (order  
1036 1915-1947, Supplementary Table S4). We could not find in the literature any proteomics study  
1037 which resorted to divisive HCA; conversely, classic (agglomerative) HCA created in 1998  
1038 (109) and its extension 2-D HCA (110) are widely used by the community, including wheat  
1039 scientists (111-115). Using agglomerative HCA on 2-DE-resolved proteins, Tasleem-Tahir  
1040 distinguished nine expression profiles throughout wheat grain growth, from anthesis to  
1041 maturity (115). In their gel-free iTRAQ analysis of early developing wheat endosperms (from  
1042 7-28 days post-anthesis (DPA)), Ma and colleagues employed HCA to delineate starch  
1043 processes (113). Similarly, five major protein expression patterns across developmental stages  
1044 4-12 DPA were outlined using HCA (116). HCA was also employed to explore the change in  
1045 expression of embryo and endosperm proteomes during wheat seed germination (117). In their  
1046 comprehensive proteomics and proteogenomics study of key developmental stages of 24 wheat  
1047 organs and tissues, Duncan and colleagues showed that HCA faithfully assigned samples to  
1048 three main clusters corresponding to first photosynthetic tissues (leaves, bracts and other green  
1049 organs), second non-photosynthetic, developmental and reproductive organs (pollen, stem,  
1050 anther, coleoptiles, roots, immature spike), and third grain (developmental series, embryo,  
1051 pericarp, endosperm) (111). More recently, Cao and colleagues discriminated differentially  
1052 expressed proteins in two wheat lines using HCA (82). All these reports demonstrate that  
1053 genotype-, sample- and tissue-specificity of protein profiles can be highlighted using  
1054 unsupervised clustering tools.

1055 3.2.2.2. Bivariate analyses (correlation and linear regression) to consider each individual  
1056 peptide against LMA



1057 As bivariate analyses handle only two variables at a time, they are not computationally taxing.  
1058 We were thus able to apply such methods on our complete dataset comprising 3,990 samples  
1059 and 32,337 peptides (including Cluster\_AAA). Due to the quantitative nature of LMA trait, we  
1060 could not perform an analysis of variance (ANOVA). We have thus carried out two bivariate  
1061 analyses, a correlation and a linear model. Because we had incorporated the LMA trait at the  
1062 peptide level as Cluster\_AAA, we could assess the validity of our analyses based on the outputs  
1063 produced by the latter.

1064 In our experiment, correlation coefficients ranged from -0.07 to 0.3, except for Cluster\_AAA  
1065 which as expected attained absolute positive correlation with a  $R^2$  of 1 (Supplementary Table  
1066 S4). Our coefficients do not show a strong relationship between peptide profiles and LMA. We  
1067 arbitrarily chose an absolute value of 0.15 to retain any LMA-associated peptide which  
1068 excluded all negatively-correlated features but included 28 positively-correlated biomarkers.  
1069 Correlation analyses are frequently employed in proteomics to unravel proteins underpinning  
1070 particular sample types, conditions or traits (118), and wheat is no exception (119-127).  
1071 Concordance of transcript and protein profiles in wheat grain were assessed via correlation  
1072 coefficients, which increased with seed maturity (120, 126). Grain yield and grain protein  
1073 content were observed to be negatively correlated, yet both also positively correlated to  
1074 nitrogen availability in a wheat genotype-specific manner (128).

1075 The q-value for the linear regression slope indicates whether changes in the explanatory  
1076 variable are significantly linked with changes in the outcome. In our work, we looked for  
1077 significant relationships between the 32,337 peptides (including Cluster\_AAA) and the inverse  
1078 function of LMA which assumed normality as a covariate factor. Q-values ranged from  $6 \times 10^{-8}$   
1079 to 1, with the exception of Cluster\_AAA which exhibited a q-value of 0 as expected  
1080 (Supplementary Table S4). We arbitrarily applied a 5% q-value threshold to consider 494  
1081 biomarker peptides whose change in expression profiles were significantly linked to variation  
1082 in LMA measurements. Linear mixed models are regularly employed by the proteomics  
1083 community for biomarker discovery approaches (129-132), but as far as we know not on wheat  
1084 grains.

1085  
1086 3.2.2.3. Compiling all statistical analyses to generate a list of candidate peptides and binning  
1087 LMA values for biomarker profiling and t test

1088 In this study, LMA-responsive biomarkers were selected based on the statistical analyses  
1089 presented above and had to fulfill at least one of the following criterium: belong to SOM group  
1090 (4,3), be included in k-means group 14, bear a divisive HCA order from 1915 to 194, exhibit a

1091 correlation  $R^2$  greater than 15%, or display a q-value inferior to 5%. This created a list of 531  
1092 biomarkers, most of which fulfilled several statistical criteria and all of them exhibiting a VIP  
1093 score for the LMA-responsive PLS greater than 1 (Supplementary Table S4).

1094 When attempting to chart the biomarker profiles, we were faced with the challenge of plotting  
1095 3,990 datapoints per gene product which ruled out typical line graphs, scatter plots, histograms  
1096 or utilising oversized illegible heat maps to represent all data points simultaneously (data not  
1097 shown). We consequently adopted a data reduction strategy involving binning the samples into  
1098 8 or 2 arbitrary bins based on their LMA values.

1099 The 8-bin profiling comprised all 3,990 samples sorted by increasing LMA measurements and  
1100 partitioned into 8 groups of equal sample size (~499 samples/bin, Supplementary Table S1).  
1101 Plotting the average of each bin as a line chart faithfully maintained the pattern of LMA  
1102 measurement observed in Figure 4A with a flat profile for the first 7 bins followed by a steep  
1103 increase in the last bin (Supplementary Figure S9A).

1104 This profiling strategy was not used for statistical purpose but proved useful during data mining  
1105 of all identified 5,514 peptides upon using tools that offered quantitative charting such as  
1106 Pathway Tools and Circos (see below).

1107 The 2-bin profiling only featured the 934 unbiased samples separated according to an arbitrary  
1108 0.17 u/g threshold (Supplementary Table S1). Plotting the average of each bin as a histogram  
1109 clearly displayed a marked quantitative increase from bin 1 to bin 2 (Supplementary Figure  
1110 S9B). This simple representation tool allowed us to categorise the 531 biomarkers as being  
1111 either up-regulated when bin 2 was taller than bin 1 denoting an accumulation in samples with  
1112  $LMA > 0.17$  u/g or down-regulated when bin 1 was taller than bin 2 denoting an accumulation  
1113 in samples with  $LMA < 0.17$  u/g.

1114 This oversimplified binning scheme allowed us to perform one last statistical analysis on the  
1115 532 (including Cluster\_AAA) biomarkers using the unbiased set of 934 samples, namely a  
1116 Student's t test with an effect size. We generated a volcano plot based on the p-values and the  
1117 directed effect size (i.e. fold change) which clearly delineated the biomarkers according to their  
1118 accumulation in bin 1 or 2 (Figure 5A).

1119 More LMA-related biomarkers were up-regulated (325) than down-regulated (206) according  
1120 to our 2-bin profiling. This was explained by the fact that all our statistical analyses, bar the  
1121 PLS and linear model, favoured peptides behaving similarly to Cluster\_AAA a proxy to LMA  
1122 actual measurements. Some exemplary patterns are displayed as histograms with error bars and  
1123 compared to that of Cluster\_AAA to expose the assortment of up- and down regulation profiles  
1124 (Figure 5B). Because the 2-bin representation was very reductive, we also present a heat map

1125 of all the intensities of the 532 biomarkers (including Cluster\_AAA) sorted by directed effect  
1126 size (i.e. fold change) in each of 934 unbiased wheat samples organised by HCA cluster order  
1127 (Figure 5C). No strong differential expression trend appeared apart from a horizontal gradient  
1128 of colours from left to right denoting the change from up- to down-regulation of the biomarkers  
1129 and a swap in colour vertically suggesting that samples were efficiently classified by the HCA.  
1130 Despite merely featuring a small subset (934x532) of our global dataset (3,990x32,337), the  
1131 heat map looked noisy and remained very hard to interpret due to an excessive number of data  
1132 points (469,888 quantities) and the lack of visually striking pattern. This further reinforced the  
1133 need to devise simple representations tools such as a Volcano plot when reporting results on  
1134 big data.

1135 To our knowledge, volcano plots have not been widely adopted by the proteomics community,  
1136 let alone wheat grain scientists with only one report so far (84), unlike heat maps which are  
1137 frequently reported in proteomics publications (133). In our work, we sorted the 531 biomarker  
1138 peptides according to their 2-bin fold changes and wheat sample based on their LC-MS  
1139 molecular similarity (Figure 5C). Zang and colleagues have adopted heat maps to profile the  
1140 proteins underpinning seed tissue organogenesis (134).

1141

### 1142 *3.2.3. Mining biomarkers to make biological sense of the data*

1143 Among the 531 biomarkers that exhibited significance levels in response to LMA  
1144 measurements, 390 were identified by LC-MS2 and matched 3,798 protein accessions  
1145 (Supplementary Table S5). This list included the most abundant and homoeologous proteins  
1146 such as the prominent storage and starch-related proteins, gliadins, glutenins, avenins, and  
1147 starch synthases as well as constitutive proteins such as histones, protein disulfide isomerases,  
1148 and tubulin, or else stress-related proteins such as heat shock and 14-3-3 proteins. We did not  
1149 identify any peptides belonging to LMA in this study, likely because we did not target high  
1150 LMA samples. To visualise our peptides of interest in a biological context, we have undertaken  
1151 a series of data mining steps. We have also made use of our 8- or 2-bins profiling strategy when  
1152 using quantitative mapping tools. The 2-bin profiling is hereafter referred to it as up- or down-  
1153 regulated gene products. The data mining tools presented below suited wheat proteins. Many  
1154 other in silico tools are freely available online which we encourage the community to employ;  
1155 however, we would not recommend using String or PlantReactome which in our hands yielded  
1156 very little results.

#### 1157 *3.2.3.1 Protein descriptions and GO terms from UniProtKB*

1158 Out of the 8,044 identities, 7,939 could be mapped in UniProtKB which flagged 6,457 GOMF  
1159 terms, 3,769 GOCC terms, 3,991 GOBP terms, as well as 1,385 unique protein names  
1160 (Supplementary Table S3). Power BI proved very useful to mine identified peptides and  
1161 simultaneously plot some of their features as histogram, scatterplot, pie chart, violin plot, tree  
1162 map and word cloud into a single dashboard (Supplementary Figure S10A) and then drill down  
1163 on some aspects, for instance inhibitor (Supplementary Figure S10B) or deamidation  
1164 (Supplementary Figure S10C).

1165 The protein names were turned into word clouds and the most frequent GO terms for each  
1166 category were presented as tree maps. Standing out from the cloud were the words “protein”,  
1167 “containing”, “domain”, “subunit”, “glutenin”, “LMW”, “molecular”, and “weight”,  
1168 confirming the preponderance of LMW glutenin subunits and domain-containing proteins such  
1169 as AAI domain-containing protein homoeologous to alpha-amylase inhibitors (Supplementary  
1170 Figure S11B-D). Also predominant among identified proteins were the words “alpha” and  
1171 “gliadin”. Word cloud is a text processing method that offers an efficient and compact  
1172 visualization of the most frequent terms in a text (135), yet it seldom appears in the scientific  
1173 literature. It has been cleverly used to categorise moonlighting proteins (136) or depict the  
1174 history of GOMF terms (137), but not in the wheat proteome. Representing our 390 identified  
1175 LMA-responsive biomarkers as word clouds revealed that up-regulated peptides belonged  
1176 predominantly to alpha-gliadins whereas down-regulated peptides mostly matched LMW  
1177 glutenins (Figure 6A,F).

1178 Rather than adopting a pie chart or histogram to plot the GO terms of all identified proteins as  
1179 commonly reported, we opted for tree maps which were initially implemented for microarray  
1180 data (138, 139) and later integrated into the web server REVIGO (140) used during our wheat  
1181 method optimisation (41). For all 8,044 identified proteins in the present study, we generated  
1182 the tree maps for all three GO classes using Power BI as it afforded more display options than  
1183 REVIGO. The most frequent biological processes (GOBP) were “polysaccharide catabolic  
1184 process” (5,643), “starch biosynthetic process” (3,688), “nucleosome assembly” (3,626),  
1185 “protein folding” (2,950) and “protein refolding” (2,499) (Supplementary Figure S11E).  
1186 “Cytoplasm” (11,888), “extracellular region” (9,964), and nucleus” (7,478) were the most  
1187 common cellular components (GOCC); recording 3,687 entries, the amyloplast was listed in  
1188 6<sup>th</sup> position (Supplementary Figure S11F). With 37,308 occurrences, the “nutrient reservoir  
1189 activity” was by far the most recurrent molecular function (GOMF), followed by “ATP  
1190 binding” (7,012) and “serine-type endopeptidase inhibitor activity” (5,811) (Supplementary  
1191 Figure S11G). The list of dominant proteins and associated GO terms in this work pointed to a

1192 storage organ such as the wheat seed and confirmed what has previously been reported in wheat  
1193 grain (41, 126, 134, 141-143). All GO terms against the 390 identified LMA-related biomarkers  
1194 are listed in Supplementary Table S5. The 207 up-regulated biomarkers came mostly from  
1195 cytoplasmic and chloroplastic proteins involved in protein translation and folding, with ATP  
1196 binding activities (Figure 6B). The 183 down-regulated peptides predominantly belonged to  
1197 cytoplasmic and cytosolic proteins acting in protein folding and TCA cycle and bearing ATP  
1198 binding activity (Figure 6G).

1199

### 1200 3.2.3.2. KEGG to retrieve Pathway, Brite and Module names

1201 From the 8,044 fasta sequences, 677 unique KEGG Orthologs (KOs) could be retrieved which  
1202 mapped to 327 KEGG pathways, 41 brites and 117 modules and annotated 11,888 peptides  
1203 (Supplementary Table S3). Identified proteins belonged to 179 (26%) KEGG metabolic  
1204 pathways with 109 (16%) KOs involved in the biosynthesis of secondary metabolites  
1205 (Supplementary Figure S12A), including sugar-related enzymes such as amylases, sucrose  
1206 synthases, hexokinases, fructokinase and beta-glucosidases.

1207 Half of KOs pointed to enzymes (336), then exosomes (71, 10%), ribosomes (62, 9%), and  
1208 chromosome-associated proteins (60, 9%) (Supplementary Figure S12B). Primary  
1209 metabolisms such as glycolysis, TCA cycle and gluconeogenesis were prominent KEGG  
1210 modules (Supplementary Figure S12C). Unexpectedly, 62 KOs (exclusively ribosomal  
1211 proteins) were associated with “Coronavirus disease – COVID 19” pathway. Similarly, many  
1212 proteins were linked with other human-related afflictions (e.g. sclerosis, neurodegeneration,  
1213 Parkinson, Huntington, Alzheimer and prion diseases; Supplementary Figure S12A). This  
1214 demonstrated the limitations of using generalist databases like KEGG that are mostly relevant  
1215 to human research to map plant proteins. While KEGG plant interface exists  
1216 (<https://www.genome.jp/kegg/genome/plant.html>) (144), plant-related datasets are dispersed  
1217 throughout the whole KEGG server so that one cannot exclusively mine plant-specific entries.  
1218 There is a need for future KEGG iterations to restrict searches to relevant taxa. Notwithstanding  
1219 non-plant hits, pathways symptomatic of grains were accurately captured in this experiment  
1220 such as the carbon metabolism (42, 6%), glycolysis/gluconeogenesis (25, 4%), as well as the  
1221 starch and sucrose metabolism (18, 3%) (Supplementary Figure S12D-F). Despite the  
1222 constraint raised above, KEGG remains a database widely employed to explore plant  
1223 proteomes, including wheat grain proteins (41, 145-147). Mapping our 390 LMA-associated  
1224 biomarkers (Supplementary Table S5) highlighted that many up-regulated peptides came from



1225 ribosomal proteins (Figure 6D) while several down-regulated peptides belonged to enzymes  
1226 acting in the biosynthesis of AAs (Figure 6I).

1227

1228 3.2.3.3. ShinyGO to retrieve enriched functional categories and chromosomal positions

1229 Multiple online tools exist to efficiently mine GO terms, however only a few cater for non-  
1230 model species, let alone plants (148-150). When looking for relevant mining tools during our  
1231 method development stage, we resorted to AgriGO online program which specifically focused  
1232 on agricultural species and offered valuable illustrations to display enrichment sets (41).  
1233 Unfortunately, AgriGO server is no longer available. We have found instead ShinyGO (50),  
1234 recently developed, which surpassed AgriGO not only in terms of enrichment visualisations  
1235 but also provided wheat protein chromosomal positions, desirable for Circos plots. A downside  
1236 of ShinyGO was that it did not perform well with UniProt accession IDs, hence the prerequisite  
1237 to retrieve TRAES IDs from UniProtKB. A total of 6,622 TRAES accessions corresponding to  
1238 the 8,044 UniProt proteins were thus retrieved, of which 4,571 could be mapped by ShinyGO  
1239 (Supplementary Table S6). An enrichment analysis ensued and could be visualised as a chart,  
1240 tree, network and chromosomal map; density plots and histograms were also produced  
1241 (Supplementary Figure S13).

1242 The most enriched category was the TCA cycle with a fold enrichment in excess of 12.5 and  
1243 the most significant GO classes were translation and peptide biosynthesis with an FDR inferior  
1244 to  $e^{-160}$  (Supplementary Figure S13A,E). Protein folding and ribonucleoprotein complex  
1245 biogenesis stood out as well among the proteins identified in this study (Supplementary Figure  
1246 S13B). Identities covered the whole genome with lower density around centromeres  
1247 (Supplementary Figure S13F). ShinyGO and other online data mining algorithms were  
1248 employed to predict genetic components systems implicated in the plant model species  
1249 *Arabidopsis* in response to high light from transcriptomics datasets publicly available (151).  
1250 Our results exemplify the relevance of ShinyGO for non-model plant species; we could not  
1251 find other cereal reports making use of it, probably due to its recent emergence (50). A fold  
1252 enrichment exceeding 200 was found among the 207 up-regulated peptides from gene products  
1253 involved in protein folding in endoplasmic reticulum (Figure 6C), followed by glycogen  
1254 metabolism, energy reserve and starch biosynthesis. ShinyGO enrichment analysis produced  
1255 very different results for our 183 down-regulated peptides, mostly invoking chromatin  
1256 assembly and remodeling, nucleosome assembly and organisation, DNA packaging and  
1257 conformation change, as well as protein-DNA complex assembly and organisation (Figure 6H).

1258



1259 3.2.3.4. Pathway Tools to retrieve differentially perturbed pathways based on 8-bin profiling  
1260 As useful as the program described above are, they yet do not accommodate quantitative data,  
1261 unlike Pathway Tools (51) made available online by the Plant Metabolic Network server and  
1262 curating the PlantCyc databases encapsulating 126 plant and algae species  
1263 (<https://plantcyc.org/>), including BreadwheatCyc (52). We could thus display protein  
1264 expression data on pathway diagrams in a dynamic and interactive way. Using the 6,622  
1265 TRAES accessions corresponding to the proteins identified in this study and the quantitative  
1266 data averaged along 8 bins, we mapped 1,432 proteins in the *T. aestivum* Pathway Tools  
1267 website (Supplementary Figure S14A).

1268 The change in expression profiles along the 8 bins was recorded and showed that all peptide  
1269 quantities varied across sample groups with multiple trends throughout the whole cellular  
1270 overview (Supplementary Video SV1). As previously reported (41), the primary and secondary  
1271 metabolisms were well covered. Overall quantities of homoeologous wheat proteins involved  
1272 in TCA and glyoxylate cycles declined along 8 bin expression profiles (Supplementary Figure  
1273 S14B).

1274 Also featured was plant hormone biosynthesis (Supplementary Figure S14C) which was  
1275 lacking in the other exploratory tools, thus demonstrating the superiority of *T. aestivum*  
1276 Pathway Tools over other databases (41). The 8 bin-profiling hinted an accumulation of  
1277 proteins related to auxin, cytokinin and gibberellin biosynthesis and a reduction of enzymes  
1278 participating in 5-deoxystrigol, brassinosteroid, and jasmonate synthesis in LMA-rich samples.  
1279 Hormonal response was flagged as one of the biochemical mechanisms of LMA expression, in  
1280 particular gibberellin and ABA signalling (22, 25, 152). Focussing on the ent-kaurene  
1281 biosynthesis, expression patterns accumulated in low LMA samples at the initial step of the  
1282 pathway and diminished in high LMA samples at the last step (Supplementary Figure S14D-  
1283 E). The first biosynthetic step is controlled by ent-copalyl disphosphate synthase (TaCSP)  
1284 which was reported to be associated with LMA via a major locus on wheat chromosome 7B  
1285 accordingly renamed as LMA-1 (153). TaCSP (Cluster\_22809 in Supplementary Figure S13F)  
1286 was one of our biomarkers. Even though databases such as Pathway Tools mapped TaCSP to  
1287 the gibberellin metabolism, its function with this phytohormone was recently contested and it  
1288 was suggested that high pI alpha-amylase synthesis in the aleurone of developing wheat grains  
1289 would be independent of gibberellins during LMA response (40). Other biomarkers matching  
1290 phytohormone-associated proteins included a cytokinin dehydrogenase whose decreasing  
1291 pattern picked up in the bin containing all the wheat sample registering high LMA  
1292 (Cluster\_24683 in Supplementary Figure S14F), and a Responsive to ABA (Rab) protein

1293 whose expression profile closely resembled that of Cluster\_AAA (Cluster\_36748 in  
1294 Supplementary Figure S14F). Interestingly, Cluster\_24621 with an increasing expression  
1295 profile belonged to an uncharacterised protein annotated with GO terms “Response to Auxin”  
1296 and “Response to ethylene” (Supplementary Figure S14F).

1297 Because Pathway Tools handles quantitative data, it produced lists of differentially perturbed  
1298 pathways (DPPS) for each set of up- and down-regulated biomarkers. Pathways characterising  
1299 wheat grains with high LMA measurements were degradations of aminobutanoate, glutamate,  
1300 and stachyose, as well as biosynthesis of UDP-galactose, UDP-glucose and sucrose (Figure  
1301 6E). DPPS differentiating samples with low LMA activities were AA metabolisms (A, K, T,  
1302 and M) rubsico shunt, superoxide radical degradation, starch biosynthesis, gluconeogenesis, S-  
1303 adenosyl-M cycle and glycolysis (Figure 6J). Our method study aside (41), we could not find  
1304 any other wheat gene expression study utilising this impressive PlantCyc database. However,  
1305 work on other plant species have amply demonstrated its value (154-159).

1306

1307 3.2.3.5. Circos plot to visualise chromosomal positions, expression profile and statistics of  
1308 identified proteins and biomarkers

1309 Invented over a decade ago (53), Circos plots have proven so valuable to efficiently represent  
1310 qualitative and quantitative information that a multitude of emulations have since arisen,  
1311 including its packaging within the Galaxy server (55) which we took advantage of here. When  
1312 the IWGSC released *T. aestivum* genome and published their findings, the genomic features  
1313 were elegantly and succinctly captured in a circular plot which highlighted homeologous genes  
1314 and translocated chromosomal regions (9). Being infinitely flexible, Circos plots can chart any  
1315 data as multiple concentric circular layers provided the correct file format is applied. We opted  
1316 to chart proteins encoded by genes we could locate on the genome (chromosomal positions  
1317 retrieved from ShinyGO analysis) and overlay their expression profiles, along with some  
1318 statistics of candidate LMA-responsive biomarkers (Figure 7).

1319 Proteins identified in this experiment aligned with the full genome, densely covering each  
1320 chromosome albeit less so around centromeric regions (Figure 7B). Overall, expression profiles  
1321 along 8-bin accumulated in bins 1-6 corresponding to wheat samples with low LMA and  
1322 decreased in bins 7-8 characterised by high LMA samples (Figure 7C). LMA-related  
1323 biomarkers were evenly dispersed on all chromosomes (Figure 7D). Plotting their effect size  
1324 (fold changes, Figure 7E) outlined that most genome areas hosted both up- and down-regulated  
1325 biomarkers bar a few exceptions on chromosomes 4, 6 and 7 for all 3 genomes A, B, and D.  
1326 Only up-regulated biomarkers could be seen on chromosome 4A region 300-500 x 10<sup>6</sup> cM and

1327 chromosome 7A region 300-480 x 10<sup>6</sup> cM (replicated on genomes B and D). They matched  
1328 three uncharacterised proteins, a 60S ribosomal protein L18a, a glucose-1-phosphate  
1329 adenylyltransferase, a polyadenylate-binding protein, a 14-3-3 protein and a protein disulfide  
1330 isomerase (Supplementary Table S5). Conversely, chromosome 6A region 300-410 x 10<sup>6</sup> cM  
1331 (replicated on genomes B and D) exclusively located down-regulated biomarkers matching a  
1332 glyceraldehyde-3-phosphate dehydrogenase, a glutathione peroxidase, a tripeptidyl-peptidase  
1333 II and an uncharacterised protein. Charting biomarker correlation values with LMA as links  
1334 failed to isolate stretches of genomic areas specific to LMA-responding proteins (Figure 7I).  
1335 This could be explained by the fact that LMA expression in our experiment elicited a complex  
1336 metabolic response involving many gene products independent of their genomic position. LMA  
1337 is indeed a multigenic trait; associated quantitative trait loci (QTLs) have been located across  
1338 all three genomes and would contribute to the LMA phenotype in an independently effective  
1339 and additive fashion (39).

1340

#### 1341 **Concluding remarks**

1342 For the first time, LMA phenotype was explored via proteomics. All the differentially regulated  
1343 biological processes highlighted in this study by the various data mining means have been  
1344 condensed into one summarising diagram and organised into broad functional categories  
1345 (Figure 8).

1346 In this work, stored LMA-affected grains activated their primary metabolisms such as  
1347 glycolysis and gluconeogenesis, TCA cycle. It also including DNA- and RNA binding  
1348 mechanisms, as well as protein translation. This logically transitioned to protein folding  
1349 activities driven by chaperones and protein disulfide isomerase, as well as protein assembly via  
1350 dimerisation and complexing. The secondary metabolism was also flagged notably with the  
1351 up-regulation of phytohormones, chemical and defense responses. LMA further invoked  
1352 cellular structures among which ribosomes, microtubules, and chromatin. Finally, and  
1353 unsurprisingly, LMA expression greatly impacted grain starch and other carbohydrates with  
1354 the up-regulation of alpha-gliadins and starch metabolism, while LMW glutenin, stachyose,  
1355 sucrose, UDP-galactose and UDP-glucose were down-regulated. This work demonstrates that,  
1356 whilst we did not find the LMA needle in the proteome haystack, proteomics deserves to be  
1357 part of the wheat LMA molecular toolkit and should be adopted by LMA scientists and breeders  
1358 in the future.

1359

#### 1360 **Abbreviations**

<b>Abbreviation</b>	<b>Full name</b>
ABA	abscisic acid
ACN	acetonitrile
AA	amino acid
AMY	amylase
ANOVA	analysis of variance
ASCA	ANOVA simultaneous component analysis
BP	biological process
CC	cellular component
cM	centimorgan
CID	collision-induced dissociation
CSV	comma separated value
cRAP	common Repository of Adventitious Proteins
DPA	day post anthesis
DNA	deoxyribonucleic acid
DPPS	differentially perturbed pathways
TaCSP	ent-copalyl disphosphate synthase from <i>Triticum aestivum</i>
ELISA	enzyme-linked immunosorbent assay
FN	falling number
FA	formic acid
FTMS	Fourier transform orbitrap mass analyser
GO	gene ontology
GxE	genetic by environment interaction
GA	gibberellic acid
Gnd-HCl	guanidine hydrochloric acid
HESI	heated electrospray ionisation
HCA	hierarchical clustering analysis
HMW	high molecular weight
HPLC	high performance liquid chromatography
ID	identity
IS	internal standard
IWGSC	International Wheat Genome Sequencing Consortium
ITMS	ion trap orbitrap mass analyser
pI	isoelectric point
IPA	isopropanol
KO	KEGG orthology
kD	kiloDalton
KNN	K-Nearest Neighbours
K-S	Kolmogorov-Smirnov

KEGG	Kyoto Encyclopedia of Genes and Genomes
LMA	late maturity alpha-amylase
LC	liquid chromatography
LMW	low molecular weight
MS or MS1	mass spectrometry
m/z	mass to charge ratio
mRNA	messenger ribonucleic acid
MF	molecular function
MLR	multivariate linear regression
ppm	part per million
PLS	partial least squares
PLSR	partial least squares regression
PTM	post-translational modification
PC	principal component
PCA	principal component analysis
QC	quality control
QTL	quantitative trait locus
QR code	quick response code
RT	retention time
Rab	Responsive to abscisic acid
RO	reverse osmosis
RT-qPCR	reverse transcription quantitative real-time polymerase chain reaction
SOM	self-organising map
SPE	solid phase extraction
MS/MS or MS2	tandem mass spectrometry
3-D	three-dimensional
TCA	trichloroacetic acid
T. aestivum	Triticum aestivum (common bread wheat)
TRAES	Triticum aestivum accession
2-DE	two-dimensional electrophoresis
2-D	two-dimensional
UTR	untranslated region
UDP	uridine diphosphate
VIP	variable importance in projection

1361

1362

1363 **Declarations**

1364 **Ethics approval and consent to participate**

1365 Not applicable.

1366 **Consent for publication**

1367 Not applicable.

1368 **Availability of data and materials**

1369 The LC-MS1 dataset and raw LC-MS2 data generated and analysed during the current study  
1370 are available in the MassIVE repository, <ftp://massive.ucsd.edu/MSV000090572>. All data  
1371 generated or analysed during this study are included in this published article and its  
1372 supplementary information files.

1373 **Competing interests**

1374 The authors declare that they have no competing interests.

1375 **Funding**

1376 This research was funded by the Grains Research and Development Corporation (GRDC),  
1377 Project DJP2001-008RTX.

1378 **Authors' contributions**

1379 Conceptualisation, M.H., H.D., J.P., D.V.; plant materials: J.P., LMA assays: N.R.; grain  
1380 grinding: D.V., A.B., D.R.; sample processing, D.V., A.B. ; LC-MS maintenance: D.V. and  
1381 V.E.; LC-MS data acquisition: D.V., and A.B.; LC-MS and LC-MS/MS data acquisition and  
1382 analysis: D.V.; LC-MS matching with LC-MS/MS in R: S.S.; technical bias removal: T.L.;  
1383 statistical analyses: D.V. and S.R.; data mining and figures, D.V.; investigation, D.V.;  
1384 resources, S.R.; data curation, D.V.; writing—original draft preparation, D.V.; review and  
1385 editing, D.V., T.L., J.P., S.R., and H.D.; visualization, D.V.; logistics: D.V.; supervision, S.R.;  
1386 project administration, D.V., S.R., H.D., and M.H; funding acquisition, M.H. and H.D. All  
1387 authors have read and agreed to the published version of the manuscript.

1388 **Acknowledgements**

1389 We thank Mr Pankaj Maharjan for retrieving all wheat samples from storage. We are grateful  
1390 for advice on MS/MS targeted methods from Drs Aaron Elkins, Priyanka Reddy from AVR,  
1391 and Dr Enzo Huang from Thermo Scientific. We are grateful to Carl Thomas and Piotr Malicki  
1392 from AVR for upgrading Genedata and Mascot servers, as well as maintaining the  
1393 Bioinformatics Advanced Scientific Computing cluster. We thank Dr Gabriel Keeble-Gagnere  
1394 from AVR for his critical review of the manuscript.

1395



1396 **Citations**

- 1397 1. Hussain B, Akpinar BA, Alaux M, Algharib AM, Sehgal D, Ali Z, et al. Capturing Wheat  
1398 Phenotypes at the Genome Level. *Front Plant Sci.* 2022;13:851079.
- 1399 2. Bacala R, Hatcher DW, Perreault H, Fu BX. Challenges and opportunities for proteomics and  
1400 the improvement of bread wheat quality. *J Plant Physiol.* 2022;275:153743.
- 1401 3. de Sousa T, Ribeiro M, Sabenca C, Igrejas G. The 10,000-Year Success Story of Wheat!  
1402 *Foods.* 2021;10(9).
- 1403 4. International Wheat Genome Sequencing C. A chromosome-based draft sequence of the  
1404 hexaploid bread wheat (*Triticum aestivum*) genome. *Science.* 2014;345(6194):1251788.
- 1405 5. Shewry PR. Wheat. *J Exp Bot.* 2009;60(6):1537-53.
- 1406 6. Shewry PR. Do ancient types of wheat have health benefits compared with modern bread  
1407 wheat? *J Cereal Sci.* 2018;79:469-76.
- 1408 7. Moshawih S, Abdullah Juperi RNA, Paneerselvam GS, Ming LC, Liew KB, Goh BH, et al.  
1409 General Health Benefits and Pharmacological Activities of *Triticum aestivum* L. *Molecules.*  
1410 2022;27(6).
- 1411 8. Venske E, Dos Santos RS, Busanello C, Gustafson P, Costa de Oliveira A. Bread wheat: a  
1412 role model for plant domestication and breeding. *Hereditas.* 2019;156:16.
- 1413 9. International Wheat Genome Sequencing C, investigators IRp, Appels R, Eversole K, Feuillet  
1414 C, Keller B, et al. Shifting the limits in wheat research and breeding using a fully annotated reference  
1415 genome. *Science.* 2018;361(6403).
- 1416 10. Guan J, Garcia DF, Zhou Y, Appels R, Li A, Mao L. The Battle to Sequence the Bread Wheat  
1417 Genome: A Tale of the Three Kingdoms. *Genomics Proteomics Bioinformatics.* 2020;18(3):221-9.
- 1418 11. Pfeifer M, Kugler KG, Sandve SR, Zhan B, Rudi H, Hvidsten TR, et al. Genome interplay in  
1419 the grain transcriptome of hexaploid bread wheat. *Science.* 2014;345(6194):1250091.
- 1420 12. Ramirez-Gonzalez RH, Borrill P, Lang D, Harrington SA, Brinton J, Venturini L, et al. The  
1421 transcriptional landscape of polyploid wheat. *Science.* 2018;361(6403).
- 1422 13. Walkowiak S, Gao L, Monat C, Haberer G, Kassa MT, Brinton J, et al. Multiple wheat  
1423 genomes reveal global variation in modern breeding. *Nature.* 2020;588(7837):277-83.
- 1424 14. Zhu T, Wang L, Rimbart H, Rodriguez JC, Deal KR, De Oliveira R, et al. Optical maps refine  
1425 the bread wheat *Triticum aestivum* cv. Chinese Spring genome assembly. *Plant J.* 2021;107(1):303-  
1426 14.
- 1427 15. Henry RJ, Furtado A, Rangan P. Wheat seed transcriptome reveals genes controlling key  
1428 traits for human preference and crop adaptation. *Curr Opin Plant Biol.* 2018;45(Pt B):231-6.
- 1429 16. Hagberg S. A Rapid Method for Determining Alpha-Amylase Activity. *Cereal Chemistry.*  
1430 1960;37(218-222).

- 1431 17. Hu Y, Sjoberg SM, Chen CJ, Hauvermale AL, Morris CF, Delwiche SR, et al. As the number  
1432 falls, alternatives to the Hagberg-Perten falling number method: A review. *Compr Rev Food Sci Food*  
1433 *Saf.* 2022;21(3):2105-17.
- 1434 18. Steber CM. Avoiding problems in wheat with low falling numbers. *Crops & Soils.*  
1435 2017;50(2):22.
- 1436 19. Newberry M, Zwart AB, Whan A, Mieog JC, Sun M, Leyne E, et al. Does Late Maturity  
1437 Alpha-Amylase Impact Wheat Baking Quality? *Front Plant Sci.* 2018;9:1356.
- 1438 20. Neoh GKS, Dieters MJ, Tao K, Fox GP, Nguyen PTM, Gilbert RG. Late-Maturity Alpha-  
1439 Amylase in Wheat (*Triticum aestivum*) and Its Impact on Fresh White Sauce Qualities. *Foods.*  
1440 2021;10(2).
- 1441 21. Sjoberg SM, Carter AH, Steber CM, Garland-Campbell KA. Unraveling complex traits in  
1442 wheat: Approaches for analyzing genotype  $\times$  environment interactions in a multienvironment study of  
1443 falling numbers. *Crop science.* 2020;60(6):3013-26.
- 1444 22. Barrero JM, Mrva K, Talbot MJ, White RG, Taylor J, Gubler F, et al. Genetic, hormonal, and  
1445 physiological analysis of late maturity alpha-amylase in wheat. *Plant Physiol.* 2013;161(3):1265-77.
- 1446 23. Derkx AP, Mares DJ. Late-maturity alpha-amylase expression in wheat is influenced by  
1447 genotype, temperature and stage of grain development. *Planta.* 2020;251(2):51.
- 1448 24. Mares DJ, Mrva K. Wheat grain preharvest sprouting and late maturity alpha-amylase. *Planta.*  
1449 2014;240(6):1167-78.
- 1450 25. Mrva K, Wallwork M, Mares DJ. alpha-Amylase and programmed cell death in aleurone of  
1451 ripening wheat grains. *J Exp Bot.* 2006;57(4):877-85.
- 1452 26. Ainsworth CC, Doherty P, Edwards KG, Martienssen RA, Gale MD. Allelic variation at  
1453 alpha-Amylase loci in hexaploid wheat. *Theor Appl Genet.* 1985;70(4):400-6.
- 1454 27. Mrva K, Mares D. Late-maturity alpha-amylase: Low falling number in wheat in the absence  
1455 of preharvest sprouting. *Journal of Cereal Science.* 2008;47:6-17.
- 1456 28. Gale MD, Law CN, Chojecki AJ, Kempton RA. Genetic control of alpha-Amylase production  
1457 in wheat. *Theor Appl Genet.* 1983;64(4):309-16.
- 1458 29. Baulcombe DC, Huttly AK, Martienssen RA, Barker RF, Jarvis MG. A novel wheat alpha-  
1459 amylase gene (alpha-Amy3). *Mol Gen Genet.* 1987;209(1):33-40.
- 1460 30. Whan A, Dielen AS, Mieog J, Bowerman AF, Robinson HM, Byrne K, et al. Engineering  
1461 alpha-amylase levels in wheat grain suggests a highly sophisticated level of carbohydrate regulation  
1462 during development. *J Exp Bot.* 2014;65(18):5443-57.
- 1463 31. Mieog JC, Janeček S, Ral JF. New insight in cereal starch degradation: identification and  
1464 structural characterization of four  $\alpha$ -amylases in bread wheat. *Amylase.* 2017;1:35-49.
- 1465 32. Ral JP, Whan A, Larroque O, Leyne E, Pritchard J, Dielen AS, et al. Engineering high alpha-  
1466 amylase levels in wheat grain lowers Falling Number but improves baking properties. *Plant*  
1467 *Biotechnol J.* 2016;14(1):364-76.

- 1468 33. Ral JF, Sun M, Mathy A, Pritchard J, Konik-Rose C, Larroque O, et al. A biotechnological  
1469 approach to directly assess the impact of elevated endogenous  $\alpha$ -amylase on Asian white-salted  
1470 noodle quality. *Starch/Stärke*. 2018;70(1700089):1-10.
- 1471 34. Cockburn D, Nielsen MM, Christiansen C, Andersen JM, Rannes JB, Blennow A, et al.  
1472 Surface binding sites in amylase have distinct roles in recognition of starch structure motifs and  
1473 degradation. *Int J Biol Macromol*. 2015;75:338-45.
- 1474 35. Verity JC, K., Hac L, Skerritt JH. Development of a Field Enzyme-Linked Immunosorbent  
1475 Assay (ELISA) for Detection of  $\alpha$ -Amylase in Preharvest-Sprouted Wheat. *Cereal Chemistry*.  
1476 1999;76(5):673-81.
- 1477 36. Mieog JC, Howitt CA, Ral JP. Fast-tracking development of homozygous transgenic cereal  
1478 lines using a simple and highly flexible real-time PCR assay. *BMC Plant Biol*. 2013;13:71.
- 1479 37. McCleary BV. Measurement of polysaccharide degrading enzymes using chromogenic and  
1480 colorimetric substrates. *Chemistry in Australia*. 1991;58:398-401.
- 1481 38. McCleary BV, McNally M, Monaghan D, Mugford DC. Measurement of alpha-amylase  
1482 activity in white wheat flour, milled malt, and microbial enzyme preparations, using the Ceralpha  
1483 assay: collaborative study. *J AOAC Int*. 2002;85(5):1096-102.
- 1484 39. Cannon AE, Marston EJ, Kiszonas AM, Hauvermale AL, See DR. Late-maturity alpha-  
1485 amylase (LMA): exploring the underlying mechanisms and end-use quality effects in wheat. *Planta*.  
1486 2021;255(1):2.
- 1487 40. Mares D, Derkx A, Cheong J, Zaharia I, Asenstorfer R, Mrva K. Gibberellins in developing  
1488 wheat grains and their relationship to late maturity alpha-amylase (LMA). *Planta*. 2022;255(6):119.
- 1489 41. Vincent D, Bui A, Ram D, Ezernieks V, Bedon F, Panozzo J, et al. Mining the Wheat Grain  
1490 Proteome. *Int J Mol Sci*. 2022;23(2):713.
- 1491 42. McCleary BV, Sheehan H. Measurement of cereal  $\alpha$ -amylase: A new assay procedure.  
1492 *Journal of Cereal Science*. 1987;6(3):237-51.
- 1493 43. Smilde AK, Jansen JJ, Hoefsloot HC, Lamers RJ, van der Greef J, Timmerman ME.  
1494 ANOVA-simultaneous component analysis (ASCA): a new tool for analyzing designed metabolomics  
1495 data. *Bioinformatics*. 2005;21(13):3043-8.
- 1496 44. R Core Team. R: A language and environment for statistical computing. R Foundation for  
1497 Statistical Computing, Vienna, Austria. 2021.
- 1498 45. Luke TDW, Pryce JE, Elkins AC, Wales WJ, Rochfort SJ. Use of Large and Diverse Datasets  
1499 for (1)H NMR Serum Metabolic Profiling of Early Lactation Dairy Cows. *Metabolites*. 2020;10(5).
- 1500 46. Batut B, Hiltemann S, Bagnacani A, Baker D, Bhardwaj V, Blank C, et al. Community-  
1501 Driven Data Analysis Training for Biology. *Cell Syst*. 2018;6(6):752-8 e1.
- 1502 47. Protein FASTA Database Handling (Galaxy Training Materials). [Internet]. 2021. Available  
1503 from: [https://training.galaxyproject.org/training-material/topics/proteomics/tutorials/database-](https://training.galaxyproject.org/training-material/topics/proteomics/tutorials/database-handling/tutorial.html)  
1504 [handling/tutorial.html](https://training.galaxyproject.org/training-material/topics/proteomics/tutorials/database-handling/tutorial.html)

- 1505 48. UniProt C. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res.*  
1506 2021;49(D1):D480-D9.
- 1507 49. Kanehisa M. The KEGG database. *Novartis Found Symp.* 2002;247:91-101; discussion -3,  
1508 19-28, 244-52.
- 1509 50. Ge SX, Jung D, Yao R. ShinyGO: a graphical gene-set enrichment tool for animals and  
1510 plants. *Bioinformatics.* 2020;36(8):2628-9.
- 1511 51. Karp PD, Latendresse M, Paley SM, Krummenacker M, Ong QD, Billington R, et al. Pathway  
1512 Tools version 19.0 update: software for pathway/genome informatics and systems biology. *Brief*  
1513 *Bioinform.* 2016;17(5):877-90.
- 1514 52. Hawkins C, Ginzburg D, Zhao K, Dwyer W, Xue B, Xu A, et al. Plant Metabolic Network  
1515 15: A resource of genome-wide metabolism databases for 126 plants and algae. *J Integr Plant Biol.*  
1516 2021;63(11):1888-905.
- 1517 53. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, et al. Circos: an  
1518 information aesthetic for comparative genomics. *Genome Res.* 2009;19(9):1639-45.
- 1519 54. Visualisation with Circos (Galaxy Training Materials). [Internet]. 2021. Available from:  
1520 <https://training.galaxyproject.org/training-material/topics/visualisation/tutorials/circos/tutorial.html>.
- 1521 55. Rasche H, Hiltemann S. Galactic Circos: User-friendly Circos plots within the Galaxy  
1522 platform. *Gigascience.* 2020;9(6).
- 1523 56. He M, Wang J, Herold S, Xi L, Schulze WX. A Rapid and Universal Workflow for Label-  
1524 Free-Quantitation-Based Proteomic and Phosphoproteomic Studies in Cereals. *Curr Protoc.*  
1525 2022;2(6):e425.
- 1526 57. Wu Y, Li L. Sample normalization methods in quantitative metabolomics. *J Chromatogr A.*  
1527 2016;1430:80-95.
- 1528 58. Li H, Han J, Pan J, Liu T, Parker CE, Borchers CH. Current trends in quantitative proteomics  
1529 - an update. *J Mass Spectrom.* 2017;52(5):319-41.
- 1530 59. O'Rourke MB, Town SEL, Dalla PV, Bicknell F, Koh Belic N, Violi JP, et al. What is  
1531 Normalization? The Strategies Employed in Top-Down and Bottom-Up Proteome Analysis  
1532 Workflows. *Proteomes.* 2019;7(3).
- 1533 60. Mitra V, Smilde AK, Bischoff R, Horvatovich P. Tutorial: Correction of shifts in single-stage  
1534 LC-MS(/MS) data. *Anal Chim Acta.* 2018;999:37-53.
- 1535 61. Mizuno H, Ueda K, Kobayashi Y, Tsuyama N, Todoroki K, Min JZ, et al. The great  
1536 importance of normalization of LC-MS data for highly-accurate non-targeted metabolomics. *Biomed*  
1537 *Chromatogr.* 2017;31(1):e3864.
- 1538 62. Poulos RC, Hains PG, Shah R, Lucas N, Xavier D, Manda SS, et al. Strategies to enable  
1539 large-scale proteomics for reproducible research. *Nat Commun.* 2020;11(1):3793.
- 1540 63. Mrode RA. *Linear Models for the Prediction of Animal Breeding Values.* 3rd ed. CABI,  
1541 editor. Wallingford, UK2014. 362 p.

- 1542 64. Lin H, Li M. Introduction to Data Science: bookdown; 2021. Available from:  
1543 <https://scientistcafe.com/ids/index.html>.
- 1544 65. Calderon-Celis F, Encinar JR, Sanz-Medel A. Standardization approaches in absolute  
1545 quantitative proteomics with mass spectrometry. *Mass Spectrom Rev.* 2018;37(6):715-37.
- 1546 66. Geyer PE, Voytik E, Treit PV, Doll S, Kleinhempel A, Niu L, et al. Plasma Proteome  
1547 Profiling to detect and avoid sample-related biases in biomarker studies. *EMBO Mol Med.*  
1548 2019;11(11):e10427.
- 1549 67. Elias JE, Haas W, Faherty BK, Gygi SP. Comparative evaluation of mass spectrometry  
1550 platforms used in large-scale proteomics investigations. *Nat Methods.* 2005;2(9):667-75.
- 1551 68. Wang G, Wu WW, Zhang Z, Masilamani S, Shen RF. Decoy methods for assessing false  
1552 positives and false discovery rates in shotgun proteomics. *Anal Chem.* 2009;81(1):146-59.
- 1553 69. Chen Y, Wang Y, Yang J, Zhou W, Dai S. Exploring the diversity of plant proteome. *J Integr*  
1554 *Plant Biol.* 2021;63(7):1197-210.
- 1555 70. Min CW, Gupta R, Agrawal GK, Rakwal R, Kim ST. Concepts and strategies of soybean  
1556 seed proteomics using the shotgun proteomics approach. *Expert Rev Proteomics.* 2019;16(9):795-804.
- 1557 71. Adhikari S, Nice EC, Deutsch EW, Lane L, Omenn GS, Pennington SR, et al. A high-  
1558 stringency blueprint of the human proteome. *Nat Commun.* 2020;11(1):5301.
- 1559 72. Burkhardt JM, Schumbrutzki C, Wortelkamp S, Sickmann A, Zahedi RP. Systematic and  
1560 quantitative comparison of digest efficiency and specificity reveals the impact of trypsin quality on  
1561 MS-based proteomics. *J Proteomics.* 2012;75(4):1454-62.
- 1562 73. Savitski MM, Kjeldsen F, Nielsen ML, Zubarev RA. Relative specificities of water and  
1563 ammonia losses from backbone fragments in collision-activated dissociation. *J Proteome Res.*  
1564 2007;6(7):2669-73.
- 1565 74. Sun S, Yu C, Qiao Y, Lin Y, Dong G, Liu C, et al. Deriving the probabilities of water loss  
1566 and ammonia loss for amino acids from tandem mass spectra. *J Proteome Res.* 2008;7(1):202-8.
- 1567 75. Yang Y. Intramolecular Cyclization Side Reactions. In: Yang Y, editor. *Side Reactions in*  
1568 *Peptide Synthesis*: Academic Press; 2016. p. 119-61.
- 1569 76. Ghatak A, Chaturvedi P, Weckwerth W. Cereal Crop Proteomics: Systemic Analysis of Crop  
1570 Drought Stress Responses Towards Marker-Assisted Selection Breeding. *Front Plant Sci.* 2017;8:757.
- 1571 77. Kerr ED, Caboche CH, Pegg CL, Phung TK, Gonzalez Viejo C, Fuentes S, et al. The post-  
1572 translational modification landscape of commercial beers. *Sci Rep.* 2021;11(1):15890.
- 1573 78. Gao F, Ayele BT. Functional genomics of seed dormancy in wheat: advances and prospects.  
1574 *Front Plant Sci.* 2014;5:458.
- 1575 79. Komatsu S, Kamal AH, Hossain Z. Wheat proteomics: proteome modulation and abiotic  
1576 stress acclimation. *Front Plant Sci.* 2014;5:684.
- 1577 80. Adegoke TV, Wang Y, Chen L, Wang H, Liu W, Liu X, et al. Posttranslational Modification  
1578 of Waxy to Genetically Improve Starch Quality in Rice Grain. *Int J Mol Sci.* 2021;22(9).



- 1579 81. Zhou C, Dong Z, Zhang T, Wu J, Yu S, Zeng Q, et al. Genome-Scale Analysis of  
1580 Homologous Genes among Subgenomes of Bread Wheat (*Triticum aestivum* L.). *Int J Mol Sci*.  
1581 2020;21(8).
- 1582 82. Cao H, Duncan O, Islam S, Zhang J, Ma W, Millar AH. Increased Wheat Protein Content via  
1583 Introgression of an HMW Glutenin Selectively Reshapes the Grain Proteome. *Mol Cell Proteomics*.  
1584 2021;20:100097.
- 1585 83. Di Francesco A, Saletti R, Cunsolo V, Svensson B, Muccilli V, Vita P, et al. Qualitative  
1586 proteomic comparison of metabolic and CM-like protein fractions in old and modern wheat Italian  
1587 genotypes by a shotgun approach. *J Proteomics*. 2020;211:103530.
- 1588 84. Maignan V, Bernay B, Geliot P, Avicé JC. Biostimulant impacts of Glutacetine(R) and  
1589 derived formulations (VNT1 and VNT4) on the bread wheat grain proteome. *J Proteomics*.  
1590 2021;244:104265.
- 1591 85. Dimitrova DS, Kaishev VK, Tan S. Computing the Kolmogorov-Smirnov Distribution When  
1592 the Underlying CDF is Purely Discrete, Mixed, or Continuous. *Journal of Statistical Software*.  
1593 2020;95(10):1-42.
- 1594 86. Lazariv T, Lehmann C. Goodness-of-Fit Tests for Large Datasets. *arXiv*.  
1595 2018:arXiv:1810.09753v1.
- 1596 87. Banerjee P, Ghosh S, Dutta M, Subramani E, Khalpada J, Roychoudhury S, et al.  
1597 Identification of key contributory factors responsible for vascular dysfunction in idiopathic recurrent  
1598 spontaneous miscarriage. *PLoS One*. 2013;8(11):e80940.
- 1599 88. Rasul G, Glover KD, Krishnan PG, Wu J, Berzonsky WA, Fofana B. Genetic analyses using  
1600 GGE model and a mixed linear model approach, and stability analyses using AMMI bi-plot for late-  
1601 maturity alpha-amylase activity in bread wheat genotypes. *Genetica*. 2017;145(3):259-68.
- 1602 89. Troyanskaya O, Cantor M, Sherlock G, Brown P, Hastie T, Tibshirani R, et al. Missing value  
1603 estimation methods for DNA microarrays. *Bioinformatics*. 2001;17(6):520-5.
- 1604 90. Horton NJ, Lipsitz SR. Multiple imputation in practice: Comparison of software packages for  
1605 regression models with missing variables. *The American Statistician*. 2001;55(3):244-54.
- 1606 91. Dixon JK. Pattern recognition with partly missing data. *IEEE Transactions on Systems, Man,  
1607 and Cybernetics*. 1979;9(10):617-21.
- 1608 92. Wold H. Estimation of principal components and related models by iterative least squares. .  
1609 In: Krishnajah PR, editor. *Multivariate analysis*. New York: Academic Press; 1966. p. 391-420.
- 1610 93. Nguyen DV, Roche DM. On partial least squares dimension reduction for microarray-based  
1611 classification: a simulation study. *Computational Statistics & Data Analysis*. 2004;46(3):407-524.
- 1612 94. Oleszko A, Hartwich J, Wojtowicz A, Gasior-Glogowska M, Huras H, Komorowska M.  
1613 Comparison of FTIR-ATR and Raman spectroscopy in determination of VLDL triglycerides in blood  
1614 serum with PLS regression. *Spectrochim Acta A Mol Biomol Spectrosc*. 2017;183:239-46.



- 1615 95. Nengsih TA, Bertrand F, Maumy-Bertrand M, Meyer N. Determining the number of  
1616 components in PLS regression on incomplete data set. *Stat Appl Genet Mol Biol*. 2019;18(6).
- 1617 96. Sherlock G. Analysis of large-scale gene expression data. *Current Opinion in Immunology*.  
1618 2000;12(2):201-5.
- 1619 97. Wang K, Wang W, Li M. A brief procedure for big data analysis of gene expression. *Animal*  
1620 *Model Exp Med*. 2018;1(3):189-93.
- 1621 98. Cresta Morgado P, Carusso M, Alonso Alemany L, Acion L. Practical foundations of  
1622 machine learning for addiction research. Part I. Methods and techniques. *Am J Drug Alcohol Abuse*.  
1623 2022;48(3):260-71.
- 1624 99. Kohonen T. Essentials of the self-organizing map. *Neural Netw*. 2013;37:52-65.
- 1625 100. Liu Z, Dai S, Bones J, Ray S, Cha S, Karger BL, et al. A quantitative proteomic analysis of  
1626 cellular responses to high glucose media in Chinese hamster ovary cells. *Biotechnol Prog*.  
1627 2015;31(4):1026-38.
- 1628 101. Fankhauser N, Maser P. Identification of GPI anchor attachment signals by a Kohonen self-  
1629 organizing map. *Bioinformatics*. 2005;21(9):1846-52.
- 1630 102. Yu D, Shen H, Yang J. SOMRuler: a novel interpretable transmembrane helices predictor.  
1631 *IEEE Trans Nanobioscience*. 2011;10(2):121-9.
- 1632 103. Fraccalvieri D, Tiberti M, Pandini A, Bonati L, Papaleo E. Functional annotation of the  
1633 mesophilic-like character of mutants in a cold-adapted enzyme by self-organising map analysis of  
1634 their molecular dynamics. *Mol Biosyst*. 2012;8(10):2680-91.
- 1635 104. Madani S, Faez K, Aminghafari M. Identifying similar functional modules by a new hybrid  
1636 spectral clustering method. *IET Syst Biol*. 2012;6(5):175-86.
- 1637 105. Tu M, Wang W, Yao N, Cai C, Liu Y, Lin C, et al. The transcriptional dynamics during de  
1638 novo shoot organogenesis of Ma bamboo (*Dendrocalamus latiflorus* Munro): implication of the  
1639 contributions of the abiotic stress response in this process. *Plant J*. 2021;107(5):1513-32.
- 1640 106. Bednarz H, Roloff N, Niehaus K. Mass Spectrometry Imaging of the Spatial and Temporal  
1641 Localization of Alkaloids in Nightshades. *J Agric Food Chem*. 2019;67(49):13470-7.
- 1642 107. Wang L, Sun X, Weiszmann J, Weckwerth W. System-Level and Granger Network Analysis  
1643 of Integrated Proteomic and Metabolomic Dynamics Identifies Key Points of Grape Berry  
1644 Development at the Interface of Primary and Secondary Metabolism. *Front Plant Sci*. 2017;8:1066.
- 1645 108. Yu T, Li G, Dong S, Liu P, Zhang J, Zhao B. Proteomic analysis of maize grain development  
1646 using iTRAQ reveals temporal programs of diverse metabolic processes. *BMC Plant Biol*.  
1647 2016;16(1):241.
- 1648 109. Eisen MB, Spellman PT, Brown PO, Botstein D. Cluster analysis and display of genome-wide  
1649 expression patterns. *Proc Natl Acad Sci U S A*. 1998;95(25):14863-8.

- 1650 110. Alon U, Barkai N, Notterman DA, Gish K, Ybarra S, Mack D, et al. Broad patterns of gene  
1651 expression revealed by clustering analysis of tumor and normal colon tissues probed by  
1652 oligonucleotide arrays. *Proc Natl Acad Sci U S A*. 1999;96(12):6745-50.
- 1653 111. Duncan O, Trosch J, Fenske R, Taylor NL, Millar AH. Resource: Mapping the *Triticum*  
1654 *aestivum* proteome. *Plant J*. 2017;89(3):601-16.
- 1655 112. Fercha A, Capriotti AL, Caruso G, Cavaliere C, Samperi R, Stampachiacchiere S, et al.  
1656 Comparative analysis of metabolic proteome variation in ascorbate-primed and unprimed wheat seeds  
1657 during germination under salt stress. *J Proteomics*. 2014;108:238-57.
- 1658 113. Ma C, Zhou J, Chen G, Bian Y, Lv D, Li X, et al. iTRAQ-based quantitative proteome and  
1659 phosphoprotein characterization reveals the central metabolism changes involved in wheat grain  
1660 development. *BMC Genomics*. 2014;15:1029.
- 1661 114. Singh RP, Runthala A, Khan S, Jha PN. Quantitative proteomics analysis reveals the  
1662 tolerance of wheat to salt stress in response to *Enterobacter cloacae* SBP-8. *PLoS One*.  
1663 2017;12(9):e0183513.
- 1664 115. Tasleem-Tahir A, Nadaud I, Chambon C, Branlard G. Expression profiling of starchy  
1665 endosperm metabolic proteins at 21 stages of wheat grain development. *J Proteome Res*.  
1666 2012;11(5):2754-73.
- 1667 116. Yang M, Gao X, Dong J, Gandhi N, Cai H, von Wettstein DH, et al. Pattern of Protein  
1668 Expression in Developing Wheat Grains Identified through Proteomic Analysis. *Front Plant Sci*.  
1669 2017;8:962.
- 1670 117. He M, Zhu C, Dong K, Zhang T, Cheng Z, Li J, et al. Comparative proteome analysis of  
1671 embryo and endosperm reveals central differential expression proteins involved in wheat seed  
1672 germination. *BMC Plant Biol*. 2015;15:97.
- 1673 118. Molendijk J, Parker BL. Proteome-wide Systems Genetics to Identify Functional Regulators  
1674 of Complex Traits. *Cell Syst*. 2021;12(1):5-22.
- 1675 119. Chen S, Chen J, Hou F, Feng Y, Zhang R. iTRAQ-based quantitative proteomic analysis  
1676 reveals the lateral meristem developmental mechanism for branched spike development in tetraploid  
1677 wheat (*Triticum turgidum* L.). *BMC Genomics*. 2018;19(1):228.
- 1678 120. Guo H, Zhang H, Li Y, Ren J, Wang X, Niu H, et al. Identification of changes in wheat  
1679 (*Triticum aestivum* L.) seeds proteome in response to anti-trx s gene. *PLoS One*. 2011;6(7):e22255.
- 1680 121. He X, Fang J, Li J, Qu B, Ren Y, Ma W, et al. A genotypic difference in primary root length  
1681 is associated with the inhibitory role of transforming growth factor-beta receptor-interacting protein-1  
1682 on root meristem size in wheat. *Plant J*. 2014;77(6):931-43.
- 1683 122. Islam N, Woo SH, Tsujimoto H, Kawasaki H, Hirano H. Proteome approaches to characterize  
1684 seed storage proteins related to ditelocentric chromosomes in common wheat (*Triticum aestivum* L.).  
1685 *Proteomics*. 2002;2(9):1146-55.

- 1686 123. Kumar RR, Dubey K, Arora K, Dalal M, Rai GK, Mishra D, et al. Characterizing the putative  
1687 mitogen-activated protein kinase (MAPK) and their protective role in oxidative stress tolerance and  
1688 carbon assimilation in wheat under terminal heat stress. *Biotechnol Rep (Amst)*. 2021;29:e00597.
- 1689 124. Li HT, Sartika RS, Kerr ED, Schulz BL, Gidley MJ, Dhital S. Starch granular protein of high-  
1690 amylose wheat gives innate resistance to amylolysis. *Food Chem*. 2020;330:127328.
- 1691 125. Peng Z, Wang M, Li F, Lv H, Li C, Xia G. A proteomic study of the response to salinity and  
1692 drought stress in an introgression strain of bread wheat. *Mol Cell Proteomics*. 2009;8(12):2676-86.
- 1693 126. Tahir A, Kang J, Choulet F, Ravel C, Romeuf I, Rasouli F, et al. Deciphering carbohydrate  
1694 metabolism during wheat grain development via integrated transcriptome and proteome dynamics.  
1695 *Mol Biol Rep*. 2020;47(7):5439-49.
- 1696 127. Zhao Y, Zhang F, Mickan B, Wang D, Wang W. Physiological, proteomic, and metabolomic  
1697 analysis provide insights into *Bacillus* sp.-mediated salt tolerance in wheat. *Plant Cell Rep*.  
1698 2022;41(1):95-118.
- 1699 128. Yu Z, Islam S, She M, Diepeveen D, Zhang Y, Tang G, et al. Wheat grain protein  
1700 accumulation and polymerization mechanisms driven by nitrogen fertilization. *Plant J*.  
1701 2018;96(6):1160-77.
- 1702 129. Daly DS, Anderson KK, Panisko EA, Purvine SO, Fang R, Monroe ME, et al. Mixed-effects  
1703 statistical model for comparative LC-MS proteomics studies. *J Proteome Res*. 2008;7(3):1209-17.
- 1704 130. D'Angelo G, Chaerkady R, Yu W, Hizal DB, Hess S, Zhao W, et al. Statistical Models for the  
1705 Analysis of Isobaric Tags Multiplexed Quantitative Proteomics. *J Proteome Res*. 2017;16(9):3124-36.
- 1706 131. Goeminne LJ, Argentini A, Martens L, Clement L. Summarization vs Peptide-Based Models  
1707 in Label-Free Quantitative Proteomics: Performance, Pitfalls, and Data Analysis Guidelines. *J*  
1708 *Proteome Res*. 2015;14(6):2457-65.
- 1709 132. Klann K, Munch C. PBLMM: Peptide-based linear mixed models for differential expression  
1710 analysis of shotgun proteomics data. *J Cell Biochem*. 2022;123(3):691-6.
- 1711 133. Pleil JD, Stiegel MA, Madden MC, Sobus JR. Heat map visualization of complex  
1712 environmental and biomarker measurements. *Chemosphere*. 2011;84(5):716-23.
- 1713 134. Zhang S, Ghatak A, Bazargani MM, Bajaj P, Varshney RK, Chaturvedi P, et al. Spatial  
1714 distribution of proteins and metabolites in developing wheat grain and their differential regulatory  
1715 response during the grain filling process. *Plant J*. 2021;107(3):669-87.
- 1716 135. Ertl P, Rohde B. The Molecule Cloud - compact visualization of large collections of  
1717 molecules. *J Cheminform*. 2012;4(1):12.
- 1718 136. Khan IK, Bhuiyan M, Kihara D. DextMP: deep dive into text for predicting moonlighting  
1719 proteins. *Bioinformatics*. 2017;33(14):i83-i91.
- 1720 137. Caetano-Anolles G. The Compressed Vocabulary of Microbial Life. *Front Microbiol*.  
1721 2021;12:655990.

- 1722 138. McConnell P, Johnson K, Lin S. Applications of Tree-Maps to hierarchical biological data.  
1723 *Bioinformatics*. 2002;18(9):1278-9.
- 1724 139. Baehrecke EH, Dang N, Babaria K, Shneiderman B. Visualization and analysis of microarray  
1725 and gene ontology data with treemaps. *BMC Bioinformatics*. 2004;5:84.
- 1726 140. Supek F, Bosnjak M, Skunca N, Smuc T. REVIGO summarizes and visualizes long lists of  
1727 gene ontology terms. *PLoS One*. 2011;6(7):e21800.
- 1728 141. Daba SD, Liu X, Aryal U, Mohammadi M. A proteomic analysis of grain yield-related traits  
1729 in wheat. *AoB Plants*. 2020;12(5):plaa042.
- 1730 142. Sharma A, Garg S, Sheikh I, Vyas P, Dhaliwal HS. Effect of wheat grain protein composition  
1731 on end-use quality. *J Food Sci Technol*. 2020;57(8):2771-85.
- 1732 143. Yang M, Liu Y, Dong J, Zhao W, Kashyap S, Gao X, et al. Probing early wheat grain  
1733 development via transcriptomic and proteomic approaches. *Funct Integr Genomics*. 2020;20(1):63-74.
- 1734 144. Kanehisa M. KEGG Bioinformatics Resource for Plant Genomics and Metabolomics.  
1735 *Methods Mol Biol*. 2016;1374:55-70.
- 1736 145. Lv X, Zhang Y, Zhang Y, Fan S, Kong L. Source-sink modifications affect leaf senescence  
1737 and grain mass in wheat as revealed by proteomic analysis. *BMC Plant Biol*. 2020;20(1):257.
- 1738 146. Yadav R, Chakraborty S, Ramakrishna W. Wheat grain proteomic and protein-metabolite  
1739 interactions analyses provide insights into plant growth promoting bacteria-arbuscular mycorrhizal  
1740 fungi-wheat interactions. *Plant Cell Rep*. 2022;41(6):1417-37.
- 1741 147. Zhang Y, Pan J, Huang X, Guo D, Lou H, Hou Z, et al. Differential effects of a post-anthesis  
1742 heat stress on wheat (*Triticum aestivum* L.) grain proteome determined by iTRAQ. *Sci Rep*.  
1743 2017;7(1):3468.
- 1744 148. Soldatos TG, Perdigao N, Brown NP, Sabir KS, O'Donoghue SI. How to learn about gene  
1745 function: text-mining or ontologies? *Methods*. 2015;74:3-15.
- 1746 149. Canto-Pastor A, Mason GA, Brady SM, Provart NJ. Arabidopsis bioinformatics: tools and  
1747 strategies. *Plant J*. 2021;108(6):1585-96.
- 1748 150. Fridrich A, Hazan Y, Moran Y. Too Many False Targets for MicroRNAs: Challenges and  
1749 Pitfalls in Prediction of miRNA Targets and Their Gene Ontology in Model and Non-model  
1750 Organisms. *Bioessays*. 2019;41(4):e1800169.
- 1751 151. Bobrovskikh AV, Zubairova US, Bondar EI, Lavrekha VV, Doroshkov AV. Transcriptomic  
1752 Data Meta-Analysis Sheds Light on High Light Response in *Arabidopsis thaliana* L. *Int J Mol Sci*.  
1753 2022;23(8).
- 1754 152. Kondhare KR, Hedden P, Kettlewell PS, Farrell AD, Monaghan JM. Quantifying the impact  
1755 of exogenous abscisic acid and gibberellins on pre-maturity alpha-amylase formation in developing  
1756 wheat grains. *Sci Rep*. 2014;4:5355.

- 1757 153. Derkx A, Baumann U, Cheong J, Mrva K, Sharma N, Pallotta M, et al. A Major Locus on  
1758 Wheat Chromosome 7B Associated With Late-Maturity alpha-Amylase Encodes a Putative ent-  
1759 Copalyl Diphosphate Synthase. *Front Plant Sci.* 2021;12:637685.
- 1760 154. Machicao J, Filho HA, Lahr DJG, Buckeridge M, Bruno OM. Topological assessment of  
1761 metabolic networks reveals evolutionary information. *Sci Rep.* 2018;8(1):15918.
- 1762 155. Gupta V, Estrada AD, Blakley I, Reid R, Patel K, Meyer MD, et al. RNA-Seq analysis and  
1763 annotation of a draft blueberry genome assembly identifies candidate genes involved in fruit ripening,  
1764 biosynthesis of bioactive compounds, and stage-specific alternative splicing. *Gigascience.* 2015;4:5.
- 1765 156. Shi X, Sun H, Chen Y, Pan H, Wang S. Transcriptome Sequencing and Expression Analysis  
1766 of Cadmium (Cd) Transport and Detoxification Related Genes in Cd-Accumulating *Salix integra*.  
1767 *Front Plant Sci.* 2016;7:1577.
- 1768 157. Nadiya F, Anjali N, Thomas J, Gangaprasad A, Sabu KK. Transcriptome profiling of *Elettaria*  
1769 *cardamomum* (L.) Maton (small cardamom). *Genom Data.* 2017;11:102-3.
- 1770 158. Sobhani Najafabadi A, Naghavi MR. Mining *Ferula gummosa* transcriptome to identify  
1771 miRNAs involved in the regulation and biosynthesis of terpenes. *Gene.* 2018;645:41-7.
- 1772 159. Ganugi P, Miras-Moreno B, Garcia-Perez P, Lucini L, Trevisan M. Concealed metabolic  
1773 reprogramming induced by different herbicides in tomato. *Plant Sci.* 2021;303:110727.
- 1774
- 1775

1776 **Figure legends**

1777 **Figure 1. High-throughput workflow used on the 4061 wheat samples.** The snowflakes  
1778 indicate storage in -80°C freezers.

1779 **Figure 2. Gantt chart capturing the timeline for each step of the proteomics workflow**  
1780 **and file accumulation.**

1781 **Figure 3: Normalisation, correction and standardisation of the raw data visualised using**  
1782 **PCA projection plots of the samples (A-F) and loading plots of the peptides (F-K).** Samples  
1783 are coloured accordingly to LC-MS injection order from blue-green to yellow-orange-red.  
1784 (A,G) PC1 vs. PC2 plot based on unnormalised LC-MS1 quantitative data; (B,H) PC1 vs. PC2  
1785 plot based on data from panels A,G normalised using the sample weights; QCs are all  
1786 condensed in a tight group (C,I) PC1 vs. PC2 plot based on data from panels B,H normalised  
1787 using the IS cluster; (D,J) PC1 vs. PC2 plot based using data from panels C,I normalised using  
1788 the injection order and the ‘intensity drift’ algorithm; (E,K) PC1 vs. PC2 plot using normalised  
1789 data from panels D,J corrected using a linear model and keeping the residuals; (F,L) PC1 vs.  
1790 PC2 plot using corrected data from panels E,L and z-transformed per row (peptides).

1791 **Figure 4: Profiles of LMA measurements for each wheat sample sorted by increasing**  
1792 **values illustrated as scatterplots (A-D) and histograms (E-H).** (A) Scatterplot of LMA  
1793 values assayed in 3,773 wheat samples; (B) Scatterplot of LMA values less than 0.17 U/g in  
1794 3,306 wheat samples; (C) Scatterplot of LMA values equal to or greater than 0.17 U/g in 467  
1795 wheat samples; (D) Scatterplot of LMA values in unbiased set containing 934 samples (see  
1796 Section 2.8.2 for explanation); (E) Histogram of LMA values assayed in 3,773 wheat samples  
1797 along 30 bins; (F) Histogram of LMA values assayed in 3773 wheat samples and transformed  
1798 using a natural logarithm (LN) function along 30 bins; (G) Histogram of LMA values assayed  
1799 in 3,773 wheat samples and transformed using an inverse function ( $1/\text{LMA}=\text{INV}(\text{LMA})$ ) along  
1800 30 bins; (H) Histogram of LMA values assayed in 3,773 wheat samples and transformed  
1801 standardising the inversion function ( $\text{STD}(\text{INV}(\text{LMA}))$ ) from panel G along 30 bins.

1802 **Figure 5: Volcano plot from t test and heat map of up- and down-regulated 531**  
1803 **biomarkers using the unbiased set of 934 wheat samples.** (A) Volcano plot of the 325 up-  
1804 regulated and 206 down-regulated biomarkers. Numbers position exemplary peptides plotted  
1805 in panel B. Cluster\_AAA with coordinates (-1.2, -23.5) is an outlier in the upper left corner  
1806 and is not featured for display purpose; (B) Mean histograms along 2 bins of clusters illustrating  
1807 up- and down-regulation patterns and located with numbers on panel A. Standard errors are  
1808 depicted with the vertical bars. Bin 1 corresponds to 467 samples with LMA < 0.17 u/g and  
1809 bin 2 corresponds to 467 samples with LMA > 0.17 u/g; (C) Heat map corresponding to the



1810 Volcano plot in panel A with peptides sorted according to directed effect size and samples  
1811 sorted based on HCA cluster order.

1812 **Figure 6: Data mining of up- and down-regulated biomarkers.** (A, F) word cloud of protein  
1813 names; (B, G) tree maps of GO terms for BP, CC and MF categories; (C, H) dot plots from  
1814 ShinyGO; (D, I) most significant KEGG pathways, ribosomes for up-regulated biomarkers and  
1815 AA biosynthesis for down-regulated biomarkers; (E, J) differentially perturbed pathways  
1816 (DPPS) from Pathway Tools.

1817 **Figure 7: Circos plot of identified proteins and LMA-responsive biomarkers with**  
1818 **expression patterns and statistics.** (A) *T. aestivum* karyotype with chromosome length  
1819 marked each  $10^6$  cM and centromeres indicated by the change in shade. LMA is displayed as a  
1820 chromosome to portray the trait's 8-bin colour pattern in trace C; (B) chromosomal positions  
1821 of all identified proteins as highlights; (C) profiling of all identified proteins along 8 bins as  
1822 heatmaps. LMA pattern is provided as a reference; (D) chromosomal positions of all identified  
1823 LMA-responsive biomarkers as highlights; (E) Volcano plot effect size of biomarkers as  
1824 scatterplot. Red denotes down-regulation and green denotes up-regulation; (F) profiling of  
1825 biomarkers along 2 bins as stacked histogram; (G) profiling of biomarkers along 8 bins as  
1826 stacked histogram; (H) biomarker accession IDs as text labels; (I) positive (green) and negative  
1827 (red) correlation with LMA as links. Green and red tags under chromosomes 4ABD, 6ABD,  
1828 and 7ABD denote genomic regions exclusive to biomarkers up- and down-regulated,  
1829 respectively.

1830 **Figure 8: Synopsis of mechanisms involved in LMA response.**

1831

### 1832 **Supplementary Figure legends**

1833 **Supplementary Figure S1: Genedata Refiner workflow to process all wheat, IS and QC**  
1834 **LCMS1 RAW files and export them to Genedata Analyst.** A. Refiner Step 1; B. Refiner  
1835 Repetition node from Step 1; C. Refiner Step 2; D. Analyst setup. See Materials and Methods  
1836 for description.

1837 **Supplementary Figure S2: Genedata Refiner workflow to process all wheat LCMS2**  
1838 **RAW files and export them to Excel.** A. Step 1; B. Repetition node from Step 1; C. Step 2;  
1839 D. Mascot parameters; E. Excel output. See Materials and Methods for description.

1840 **Supplementary Figure S3: LC-MS2 RAW maps for each tandem pass.** X-axis delineates  
1841 300-2000 m/z. Y-axis delineates 1-35 min Retention Time. White dots represent MS2 events.  
1842 (A) LC-MS1 map of pooled sample; (B) LC-MS2 map of Pass 1 replicate 1 with 3000  
1843 threshold; (C) LC-MS2 map of Pass 2 replicate 1 with exclusion list of 2000 ions fragmented

1844 in Pass 1; (D) LC-MS2 map of Pass 3 replicate 1 with exclusion list of 2000 ions fragmented  
1845 in Pass 2; (E) LC-MS2 map of Pass 4 replicate 1 with exclusion list of 2000 ions fragmented  
1846 in Pass 3; (F) LC-MS2 map of Pass 5 replicate 1 (same as Pass 1 but with 500 threshold); (G)  
1847 LC-MS2 map of Pass 6 replicate 1 with inclusion list of 2000 most abundant ions from Pass 1;  
1848 (H) LC-MS2 map of Pass 7 with inclusion list 1 loaded Global mass tab and 2 m/z tolerance;  
1849 (I) LC-MS2 map of Pass 8 with inclusion list 1 loaded in data-dependent settings and 2 m/z  
1850 tolerance; (J) LC-MS2 map of Pass 9 with inclusion list 1 loaded in data-dependent settings  
1851 and 1 m/z tolerance; (K) LC-MS2 map of Pass 10 with inclusion list 1 loaded in data-dependent  
1852 settings and 0.5 m/z tolerance; (L) LC-MS2 map of Pass 11 with inclusion list 1 loaded in data-  
1853 dependent settings and 0.2 m/z tolerance. Maps from other replicates in Passes 1-6 or with  
1854 inclusion lists 2-10 for Passes 7-11 are not shown.

1855 **Supplementary Figure S4: Histogram of the number of peptides identified using Mascot**  
1856 **algorithm and number of MS2 events in each of the LC-MS2 file.** Black bars represent  
1857 peptide counts (y axis on the left) and orange dots depict MS/MS event counts (y axis on the  
1858 right).

1859 **Supplementary Figure S5: Histograms (A, C, E) and box plots (B, D) of the number of**  
1860 **peptides per accession (A-B, E) and number of accessions per peptides (C-D).** The orange  
1861 line in panels A and C represents cumulated counts in percent. Panel E displays the peptides  
1862 with the highest hit counts belonging either to low molecular weight glutenin subunit (LMW-  
1863 GS), alpha-gliadin (GLIA), or gamma-gliadin (GLIG).

1864 **Supplementary Figure S6: Distribution of LC-MS1 data across 3,990 wheat samples and**  
1865 **32,336 quantified peptides.** (A) Histogram of the corrected dataset using a linear model and  
1866 keeping the residuals; (B) Boxplot of corrected dataset log<sub>10</sub> transformed for display purpose;  
1867 (C) Histogram of the corrected dataset z-transformed per row of peptides; (D) Boxplot of z-  
1868 transformed dataset log<sub>10</sub> transformed for display purpose. Insets in panels A-B indicate one-  
1869 sample Kolmogorov-Smirnov (K-S) test results where D is the value of the K-S statistics.

1870 **Supplementary Figure S7: Partial Least Square (PLS) using LMA as a response on the**  
1871 **unbiased samples and the unbiased samples and all the quantified peptides.** (A) Score plot  
1872 of Component 1 vs Component 2 of the 934 unbiased samples coloured based on LMA  
1873 measurements; samples with high LMA are circled; (B) Loading plot of Component 1 vs  
1874 Component 2 of the 32,337 peptides coloured based on PLS VIP scores; peptides with high  
1875 LMA are circled; Cluster\_AAA resolves in the top right corner and contributes the most to the  
1876 PLS with a VIP score of 38.84.

1877 **Supplementary Figure S8: Partial least square regression (PLSR) to impute LMA**  
1878 **missing values.** (A) Full scatterplot of the measured vs. predicted LMA values of the testing  
1879 set containing 179 samples; (B) same as panel A but limiting LMA predicted values inferior to  
1880 0.17 u/g; (C) same as panel A but limiting LMA predicted values superior to 0.17 u/g; (D) Line  
1881 chart of the 217 LMA missing values and predicted by our PLSR model and sorted based on  
1882 increasing LMA.

1883 **Supplementary Figure S9: Binning strategies of wheat samples based on LMA**  
1884 **measurements.** (A) all 3990 wheat samples were sorted by increasing order of LMA values  
1885 and then split into 8 arbitrary bins of 499 samples each; the line chart displays bin averages;  
1886 (B) the 934 unbiased wheat samples were sorted by increasing order of LMA values and then  
1887 split into 2 arbitrary bins of 467 samples each based on a LMA value threshold of 0.17 u/g; the  
1888 histogram displays bin averages. Bins are listed in Supplementary Table S1.

1889 **Supplementary Figure S10: Mining identified proteins using Power BI.** (A) all identified  
1890 peptides plotted as peptide mass against Mascot peptide scores (dot histogram), peptide missed  
1891 cleavages (pie chart), peptide PTMs (tree map), peptide lengths (violin plot), peptide charges  
1892 (vertical bar plot), protein score against sequence coverage (scatterplot) and protein description  
1893 (word cloud); (B) same charts but drilled down on the term “inhibitor” in the word cloud of  
1894 protein descriptions; (C) same charts but drilled down on “deamidated” peptides in the tree  
1895 map of PTMs.

1896 **Supplementary Figure S11: Retrieval of protein descriptions and Gene Ontology (GO)**  
1897 **terms for Molecular Function (MF), Cellular Component (CC), and Biological Process**  
1898 **(BP) from UniProtKB using all 8,044 protein identities.** (A) UniprotKB output viewed by  
1899 GO; (B) word cloud of all protein names; (C) word cloud of protein names filtered as  
1900 “glutenin”; (D) word cloud of protein names filtered as “domain-containing”; (E) tree map of  
1901 the most abundant terms for GOBP category; (F) tree map of the most abundant terms for  
1902 GOCC category; (E) tree map of the most abundant terms for GOMF category.

1903 **Supplementary Figure S12: KEGG output using all 8,044 identified proteins matching**  
1904 **677 KOs.** (A) Histogram of the most frequent pathways; (B) Histogram of the most frequent  
1905 brite terms; (C) Histogram of the most frequent modules; (D) Carbon metabolism map; (E)  
1906 Glycolysis/gluconeogenesis map; (F) Starch and sucrose metabolism map. Proteins identified  
1907 in this study are highlighted in green in panels D-F.

1908 **Supplementary Figure S13: ShinyGO outputs using all 6,622 TRAES accessions**  
1909 **corresponding to the 8,044 UniProt proteins.** (A) dot plot of the GO categories sorted by  
1910 fold enrichment; (B) network of nodes representing enriched GO terms. Related GO terms are

1911 connected by a line, whose thickness reflects percent of overlapping genes. Node size  
1912 represents the number of genes; (C-D) statistical analysis on the genomic features. Chi-squared  
1913 and Student's t-tests are run to compare the user's genes to the *T. aestivum* genome. Results on  
1914 number of exons, transcript isoforms, GC content, untranslated region (UTR) length, and types  
1915 of genes (coding, non-coding, pseudogenes) are displayed as density scatterplots or histograms;  
1916 (E) hierarchical clustering tree of significant enriched pathways. Pathways that share many  
1917 genes are clustered together and dot size indicates q-values significance; (F) Plot of the  
1918 chromosomal positions of the genes encoding our identified proteins.

1919 **Supplementary Figure S14: Pathway Tools output using 6622 TRAES accessions and**  
1920 **quantitative data averaged along 8 bins.** (A) OMICS dashboard general view; (B) cellular  
1921 view zoomed in on TCA cycle II and glyoxylate cycle. Each expression profile points to a  
1922 unique TRAES accession, most of them being homologous. The whole cellular view is  
1923 available in Supplementary Video SV1; (C) OMICS dashboard zoomed in on hormone  
1924 biosynthesis; (D) OMICS dashboard zoomed in on gibberellin and gibberellin precursor  
1925 biosynthesis; (E) Pathway view of ent-kaurene biosynthesis from the Gibberellin biosynthesis  
1926 pathway further illustrating high homology of wheat proteins; (F) 8-bin profiles of peptide  
1927 biomarkers belonging to proteins involved in phytohormone biosynthesis; Cluster\_AAA is  
1928 displayed for comparison purpose.

1929

Graphical abstract

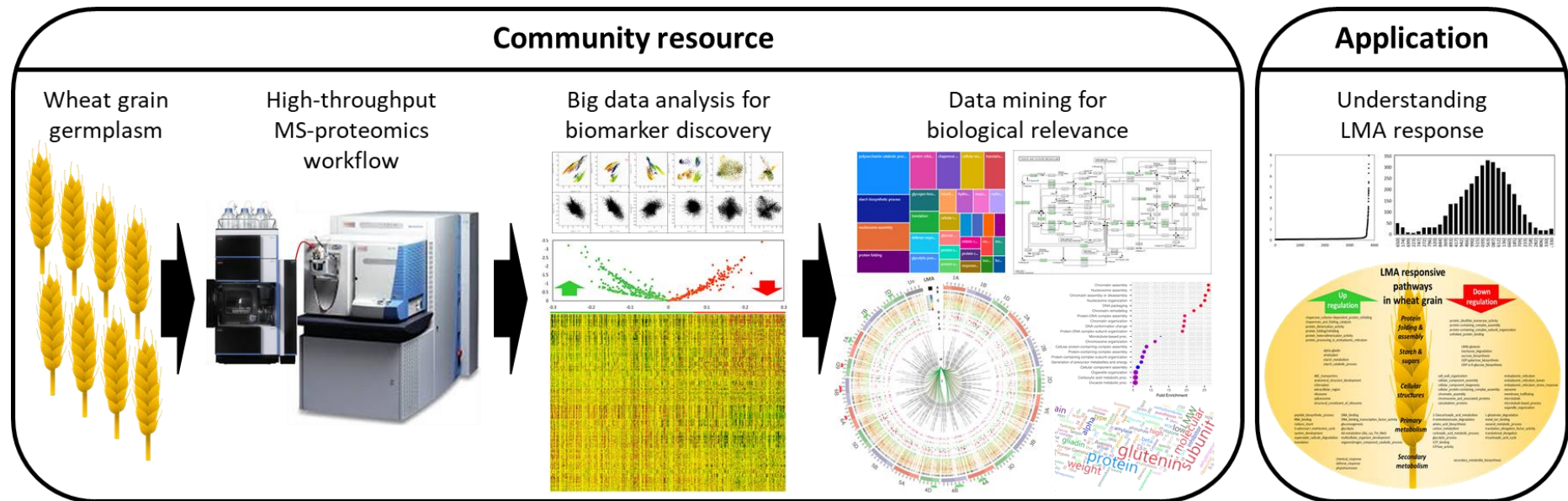




Figure 1

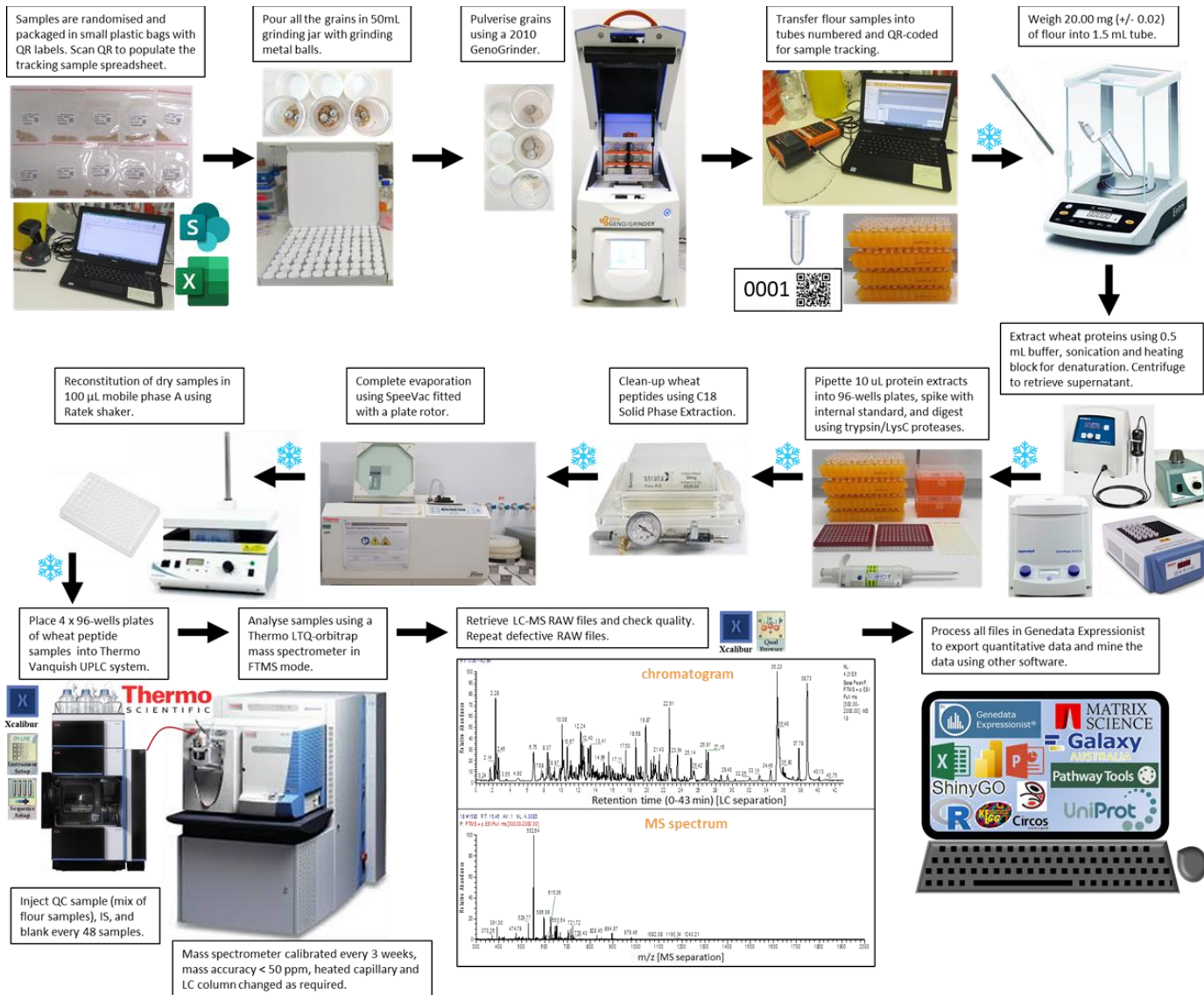




Figure 2

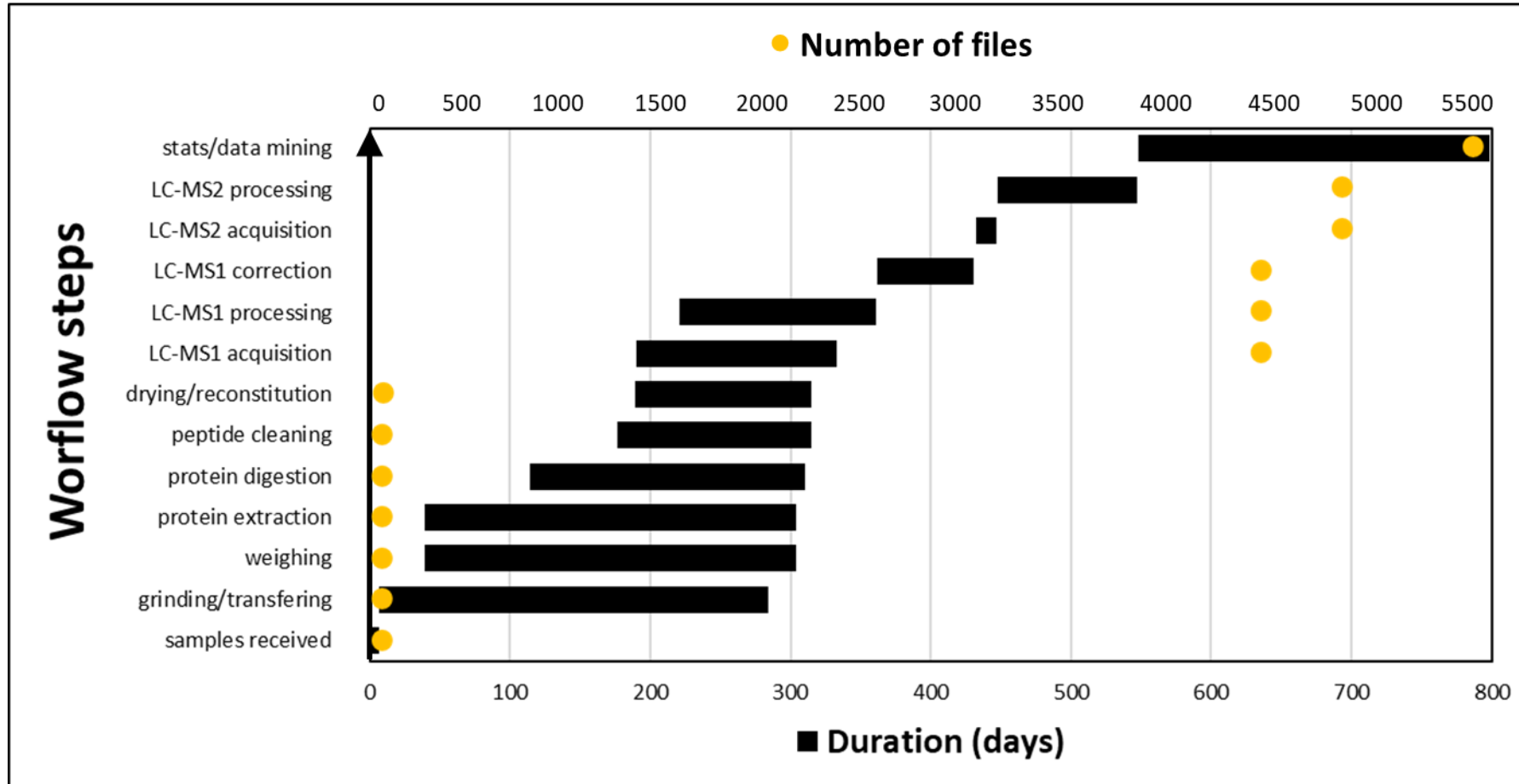


Figure 3

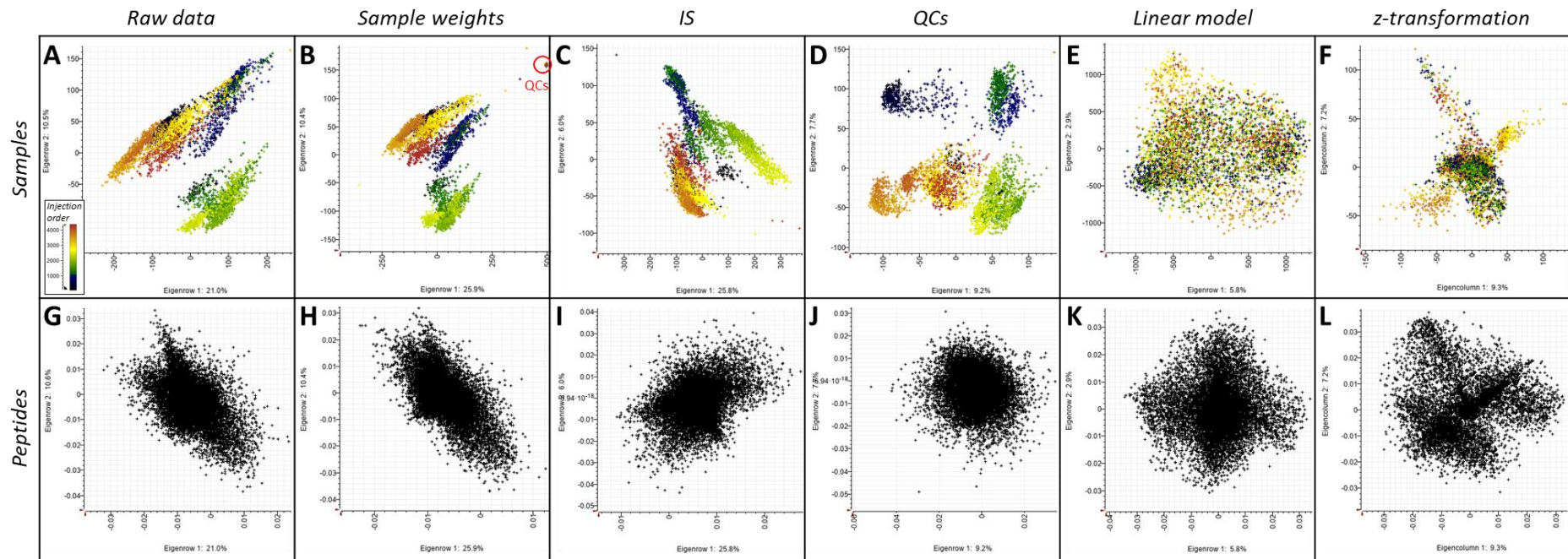


Figure 4

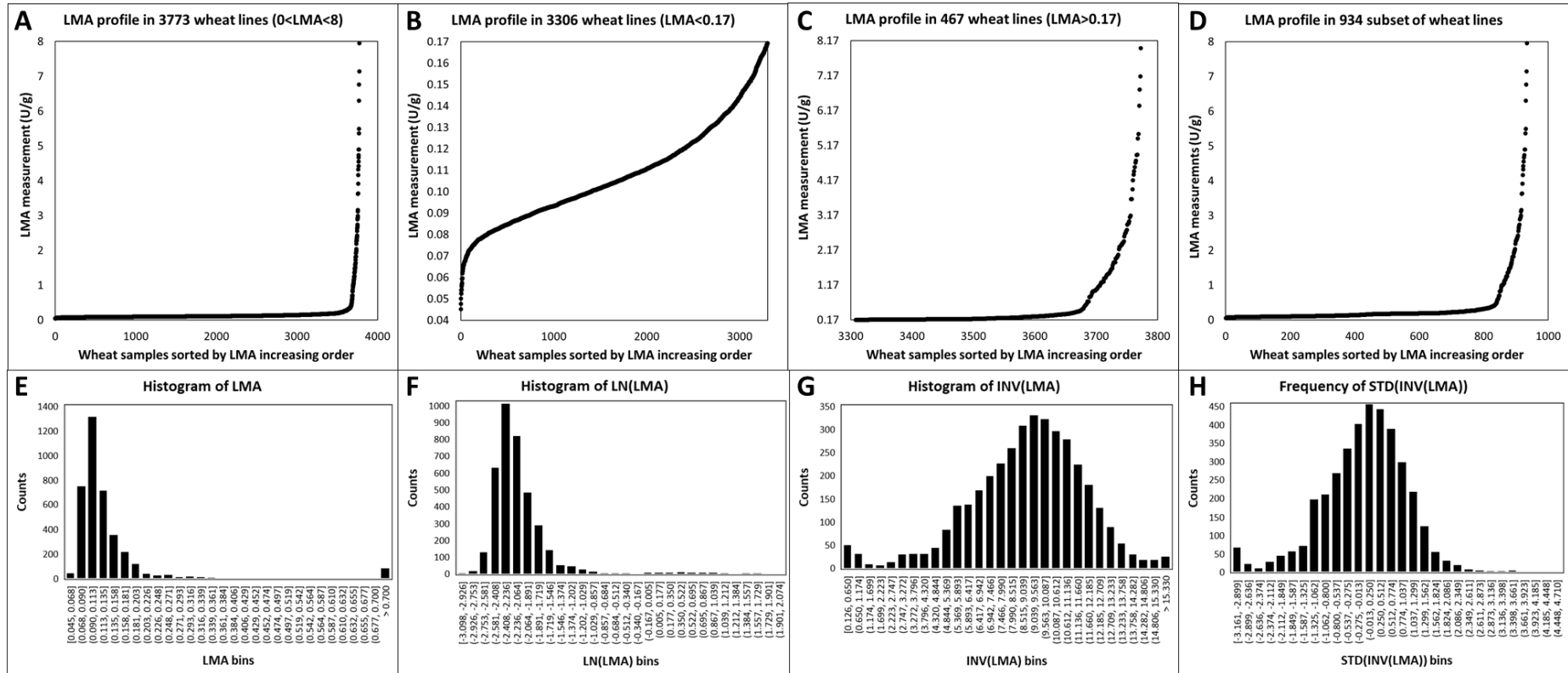


Figure 5

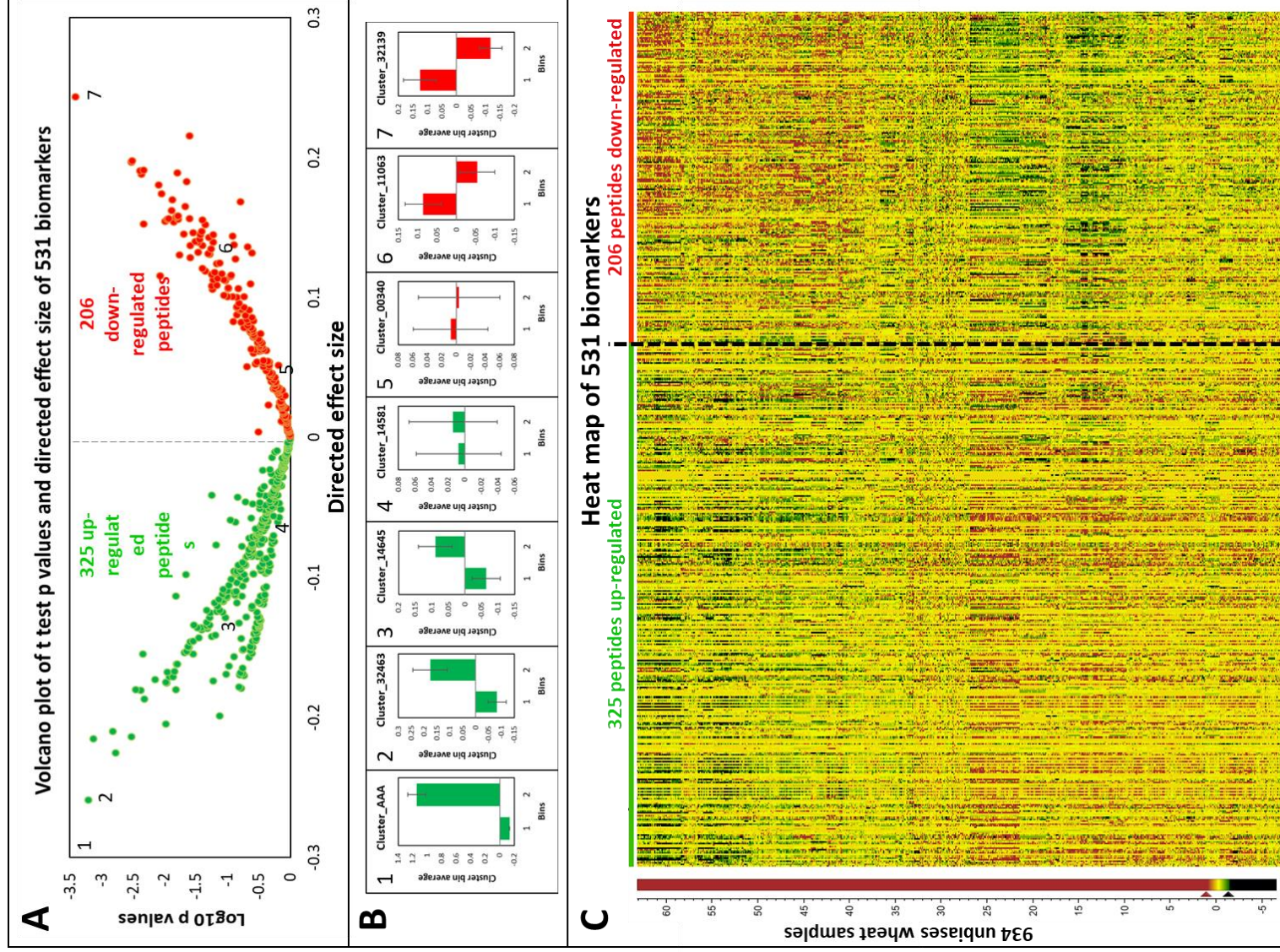




Figure 6

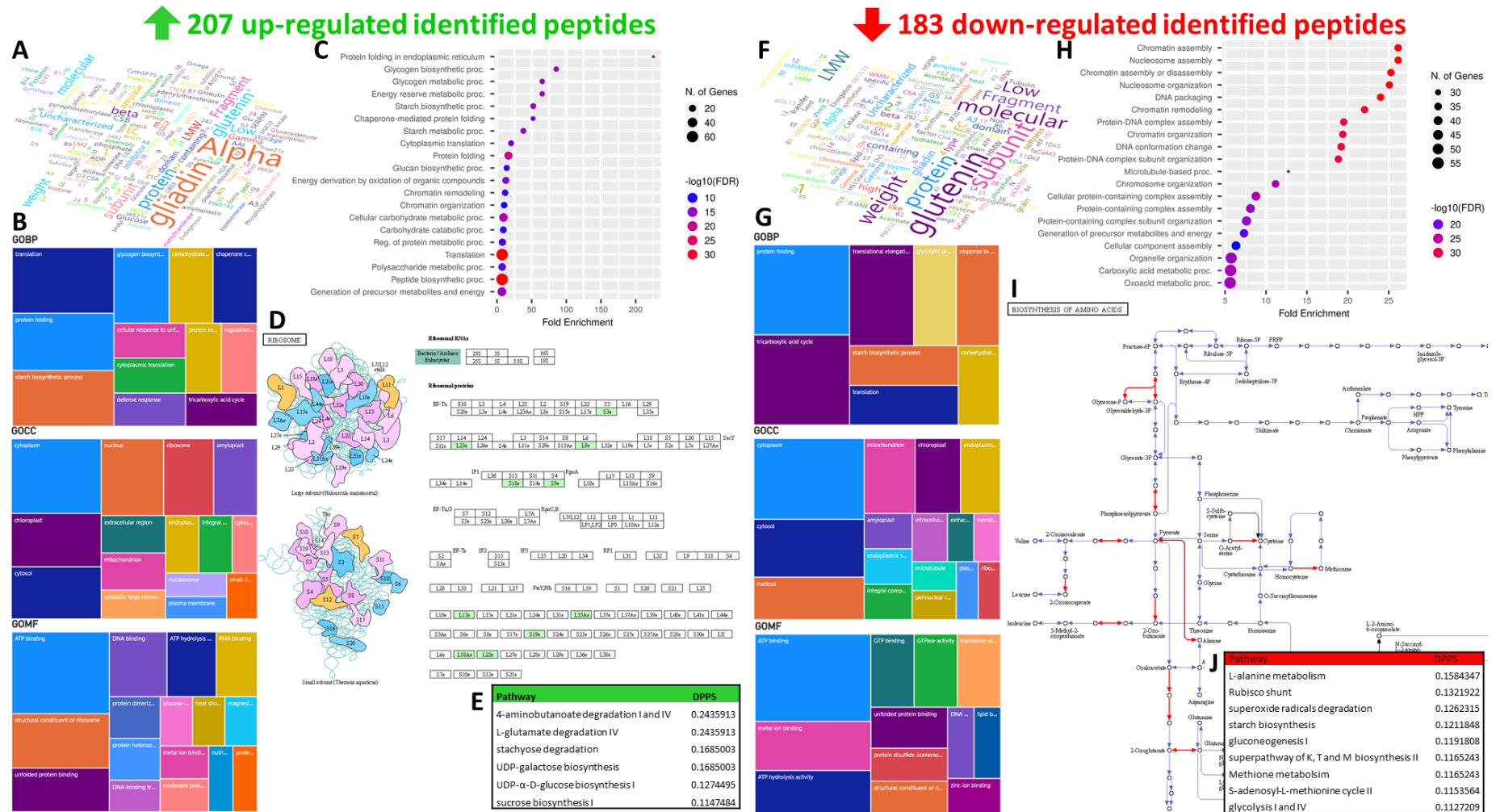
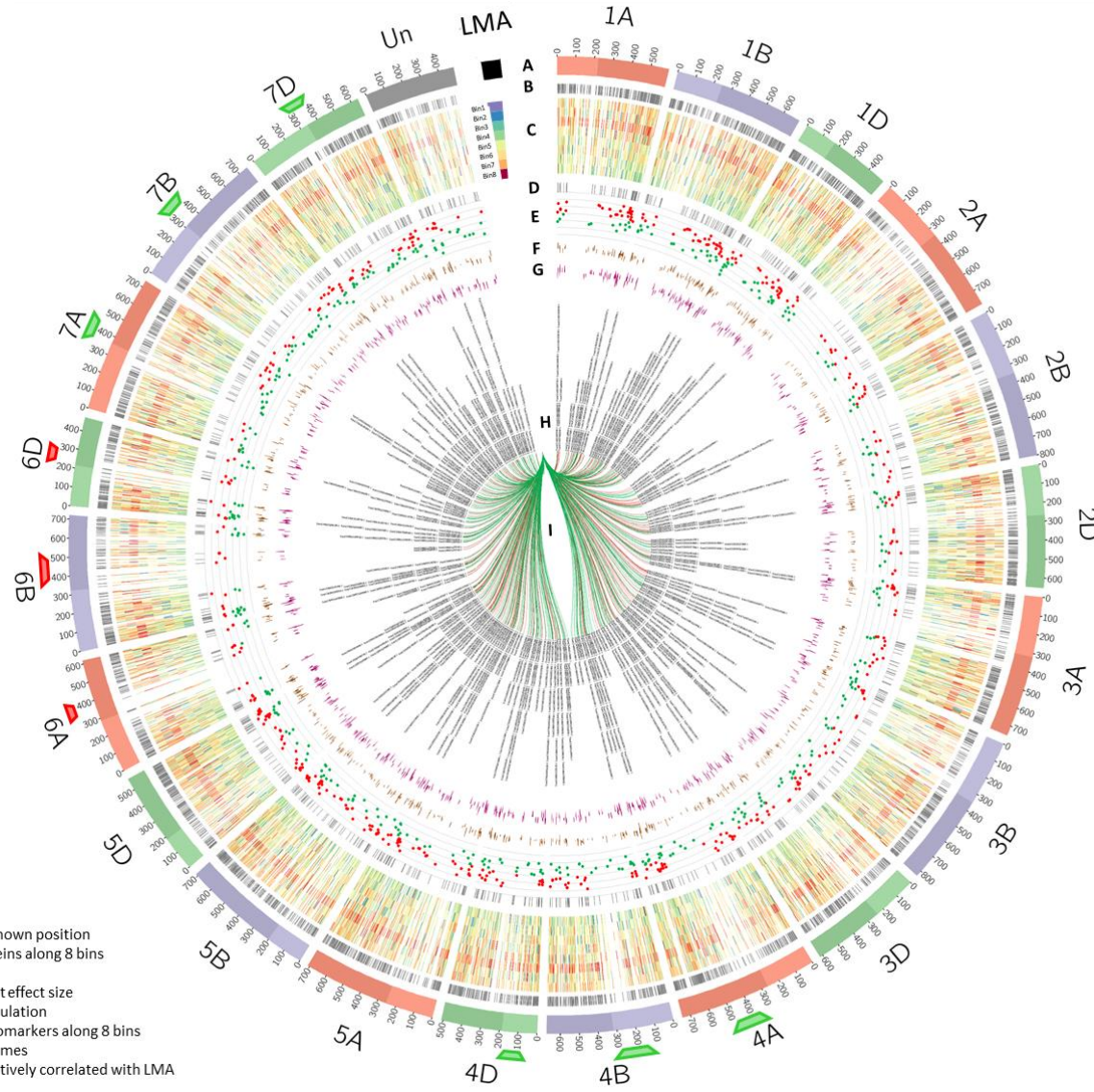


Figure 7



**Figure legend:**  
 A karyotype with centromere  
 B all identified proteins with known position  
 C expression profile of all proteins along 8 bins  
 D LMA biomarkers  
 E LMA biomarkers Volcano plot effect size  
 F LMA biomarker up/down regulation  
 G expression profile of LMA biomarkers along 8 bins  
 H LMA biomarker accession names  
 I biomarkers positively or negatively correlated with LMA



Figure 8

