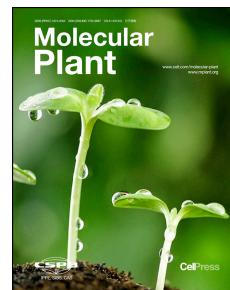


Journal Pre-proof

Large-scale Discovery of Non-conventional Peptides in Maize and Arabidopsis Through an Integrated Peptidogenomic Pipeline

Shunxi Wang, Lei Tian, Haijun Liu, Xiang Li, Jinghua Zhang, Xueyan Chen, Xingmeng Jia, Xu Zheng, Shubiao Wu, Yanhui Chen, Jianbing Yan, Liuji Wu



PII: S1674-2052(20)30147-7

DOI: <https://doi.org/10.1016/j.molp.2020.05.012>

Reference: MOLP 939

To appear in: MOLECULAR PLANT

Accepted Date: 18 May 2020

Please cite this article as: **Wang S., Tian L., Liu H., Li X., Zhang J., Chen X., Jia X., Zheng X., Wu S., Chen Y., Yan J., and Wu L.** (2020). Large-scale Discovery of Non-conventional Peptides in Maize and Arabidopsis Through an Integrated Peptidogenomic Pipeline. Mol. Plant. doi: <https://doi.org/10.1016/j.molp.2020.05.012>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

All studies published in MOLECULAR PLANT are embargoed until 3PM ET of the day they are published as corrected proofs on-line. Studies cannot be publicized as accepted manuscripts or uncorrected proofs.

© 2020 The Author

1 **Large-scale Discovery of Non-conventional Peptides in Maize and**
2 **Arabidopsis Through an Integrated Peptidogenomic Pipeline**

3

4 Shunxi Wang^{1†}, Lei Tian^{1†}, Haijun Liu^{2†}, Xiang Li², Jinghua Zhang¹, Xueyan Chen¹,
5 Xingmeng Jia¹, Xu Zheng¹, Shubiao Wu³, Yanhui Chen¹, Jianbing Yan^{2*}, Liuji Wu^{1*}

6

7 1 National Key Laboratory of Wheat and Maize Crop Science, Collaborative
8 Innovation Center of Henan Grain Crops, College of Agronomy, Henan Agricultural
9 University, Zhengzhou 450002, China.

10 2 National Key Laboratory of Crop Genetic Improvement, Huazhong Agricultural
11 University, Wuhan 430070, China.

12 3 School of Environmental and Rural Science, University of New England, Armidale
13 NSW 2351, Australia.

14 † These authors contributed equally to this work.

15

16 ***Corresponding authors**

17 Liuji Wu (wlj200120@163.com) and Jianbing Yan (yjianbing@mail.hzau.edu.cn)

18

19 **Running title:** Large-scale discovery of non-conventional peptides

20

21 **SHORT SUMMARY**

22 This study developed an integrated peptidogenomic pipeline and firstly applied it for
23 large-scale identification of non-conventional peptides (NCPs) in plant. The identified
24 NCPs, which were derived from introns, 3'UTRs, 5'UTRs, junctions and intergenic
25 regions, showed distinct characteristics compared to conventional peptides (CPs).
26 Functional analysis unveiled potential function of NCPs in plant genetic regulation of
27 complex traits and evolution.

28 **ABSTRACT**

29 Non-conventional peptides (NCPs), which include small open reading frame-encoded
30 peptides, play critical roles in fundamental biological processes. Here we developed
31 an integrated peptidogenomic pipeline using high-throughput mass spectra to probe a
32 customized six-frame translation database and applied it to large-scale identification
33 of NCPs in plants. Altogether, 1,993 and 1,860 NCPs were unambiguously identified
34 in maize and Arabidopsis, respectively. The NCPs showed distinct characteristics
35 compared to conventional peptides (CPs) and were derived from introns, 3'UTRs,
36 5'UTRs, junctions and intergenic regions. These results revealed that translation
37 events in unannotated transcripts occurred more broadly than previously thought. In
38 addition, maize NCPs were found to be enriched within regions associated with
39 phenotypic variations and domestication selection, indicating their potential function
40 in plant genetic regulations of complex traits and evolution. Summarily, this study
41 provides an unbiased and global view of plant NCPs. The identification of large-scale
42 NCPs in both monocot and dicot plants reveals that a much larger portion of the plant
43 genome can be translated to biologically functional molecules, which has important
44 implications in functional genomic studies. The present study also provides a useful
45 resource for the characterization of more hidden NCPs in other plants.

46

47 **Key words:** non-conventional peptides, small open reading frames, peptidogenomics,
48 mass spectrometry, six-frame translation, plants

49 INTRODUCTION

50 Peptides, typically composed of 2 to 100 amino acid residues, represent the small
51 biological molecules with important roles in biology (Tavormina et al., 2015). Small
52 signaling peptides (SSPs) or peptide hormones, which are a class of short peptides
53 ranging from 5 to 75 amino acid in length, also play critical roles in various biological
54 processes. For example, the discovery and application of the peptide hormone insulin
55 was one of the greatest achievements in the 20th century (Banting and Best, 2007).
56 Studies over the past few decades have mainly focused on conventional peptides (CPs)
57 derived from annotated coding sequences (CDSs) or conventional open reading
58 frames. Recently, a novel class of peptides, now defined as non-conventional peptides
59 (NCPs) in this study, has caught significant attentions as functionally important
60 endogenous peptides in various organisms (Ma et al., 2014; Couso and Patraquim,
61 2017; Plaza et al., 2017; Jackson et al., 2018; Chen et al., 2020a). These NCPs are
62 derived from previously unannotated CDSs, such as intergenic regions, untranslated
63 regions (UTRs), introns and various types of junctions, as well as different reading
64 frames from annotated CDSs.

65 A primary report of the NCP was published more than two decades ago, where a
66 10 amino acid peptide was identified to be translated from *ENOD40*, a gene
67 previously annotated as untranslated (van de Sande et al., 1996). Thereafter, the
68 *ENOD40* was further proved to play a key role in regulating the response to auxin in
69 the flowering plants (Rohrig et al., 2002). In animals and humans, NCPs are known to
70 play important roles in a diverse range of cellular processes, such as calcium transport
71 (Magny et al., 2013), embryogenesis (Kondo et al., 2010), muscle performance
72 (Nelson et al., 2016; Matsumoto et al., 2017), translation control (Hinnebusch et al.,
73 2016; Couso and Patraquim, 2017; Plaza et al., 2017), immune response (Laumont et
74 al., 2016) and stress resistance (Khitun et al., 2019). Functional NCPs, such as
75 POLARIS (Casson et al., 2002), ROTUNDIFOLIA4 (Narita et al., 2004), KOD
76 (Blanvillain et al., 2011), OSIP108 (De Coninck et al., 2013), miPEP165a
77 (Lauressergues et al., 2015), PSEP1, PSEP3, PSEP18, PSEP25 (Fesenko et al., 2019),

78 CDC26 (Lorenzo-Orts et al., 2019) and vvi-miPEP171d1 (Chen et al., 2020b), have
79 been reported in plants. These studies have demonstrated that NCPs play essential
80 roles in plant development, environmental responses and translational control.
81 However, due to the limitations of genomic annotation and peptidomic technology, a
82 plethora of NCPs are usually dismissed from further analysis or annotation in plants
83 (Andrews and Rothnagel, 2014; Yin et al., 2019).

84 The increasing importance of NCPs has led to emerging strategies for their
85 discovery. The advent of next-generation sequencing and developments in
86 bioinformatics has boosted the research of NCPs at a genome-wide scale.
87 Computational approaches based on sequence similarities have been developed to
88 identify potential translational small open reading frames (sORFs) in noncoding
89 sequences (Hurst, 2002; Kastenmayer et al., 2006; Hanada et al., 2007; Makarewicz
90 and Olson, 2017). However, conservation and homology analysis of sORFs is difficult
91 due to the short sequence and low conservation score. Another strategy is to use
92 ribosome profiling by sequencing ribosome-protected fragments that enables mapping
93 of a genome-wide set of transcripts that are being translated (Ingolia et al., 2009;
94 Ingolia et al., 2011; Ingolia, 2016; Shiber et al., 2018). In recent years, ribosome
95 profiling has been widely used to confirm the translation of non-annotated ORFs in
96 various species (Ruiz-Orera et al., 2014; Wu et al., 2019; Kurihara et al., 2020). While
97 ribosome profiling itself is an experimental approach, the evaluation of the coding
98 potential of an identified region of interest is in fact mostly computational
99 (Makarewicz and Olson, 2017). Existing ribosome profiling techniques have
100 undergone significant modifications and enhancements, which have improved reliably
101 in protein-coding transcript identification (Hsu et al., 2016; Bazin et al., 2017). As a
102 different strategy from ribosome profiling, mass spectrometry (MS)-based method is
103 able to detect peptides that are translated from a sORF and can thereby directly
104 validate the protein-coding potential of the transcript (Castellana et al., 2008;
105 Makarewicz and Olson, 2017). Recently, a new method referred as peptidogenomics,
106 which integrates peptidomics (based on high throughput MS/MS) and genomics, has
107 emerged as a promising strategy for deep analysis of the endogenous NCPs (Kersten

108 et al., 2011; Harvey et al., 2015). As an efficient strategy, peptidogenomics has been
109 successfully used in microorganisms and humans (Liu et al., 2011; Slavoff et al., 2013;
110 Mohimani and Pevzner, 2016; Mohimani et al., 2018). However, owing to
111 experimental and computational issues, such as endogenous peptide enrichment,
112 nonspecific protease digestion and lack of complete peptide reference databases, the
113 identification of NCPs using peptidogenomics in plant is still challenging.

114 Here, we developed an integrated peptidogenomic pipeline for large-scale
115 identification of NCPs in monocot and dicot plants. High-throughput mass spectra of
116 endogenous peptides were used to probe Ensembl protein database and the
117 customized peptidogenomic database derived from the six-frame translation of
118 genomic sequences. Our results revealed that NCPs could be derived from not only
119 coding sequences but also allegedly noncoding sequences. NCPs showed a distinct
120 distribution pattern from that of CPs. In addition, we found that the NCPs were
121 enriched within the genomic regions associated with phenotypic variations and
122 domestication, indicating their potential functions in regulating phenotypes and
123 shaping the evolution of the plants. These results represent a large-scale identification
124 of endogenous NCPs in plants through the integrated peptidogenomic pipeline and
125 thus provide valuable information towards the understanding of the biological
126 function of these hidden molecules.

127 **RESULTS**

128 **An Integrated Peptidogenomic Pipeline for NCPs Identification in**
129 **Plants**

130 Directly detecting NCPs is the most definitive evidence of their existence. To
131 facilitate plant NCPs discovery, we developed and applied an integrated
132 peptidogenomic pipeline for large-scale identification of plant NCPs (Figure 1A). For
133 sample preparation, an acid extraction buffer consisting of 1% trifluoroacetic acid
134 (TFA) was utilized based on a previous study (Chen et al., 2014). In addition, heat
135 stabilization by water bath combined with plant protease inhibitors was applied to
136 diminish nonspecific protease digestion. Trichloroacetic acid (TCA)-acetone

137 precipitation was also applied to establish an optimized sample preparation protocol.
138 Then, plant endogenous peptides were enriched from larger protein fragments by
139 centrifugation through 10 kDa cutoff filters before they were analyzed with liquid
140 chromatography tandem mass spectrometry (LC-MS/MS).

141 To capture the endogenous peptides globally present in maize, the Mascot search
142 engine was used to match the resulting mass spectrum data set against Ensembl
143 protein database and customized peptidogenomic database, respectively. The
144 customized peptidogenomic database was constructed using the six-frame translation
145 of maize genomic sequences (Figure 1B). As a result, we obtained a ~5.2-gigabase
146 (Gb) customized peptidogenomic database (containing ~136 million sequences). To
147 avoid an inflated search space for the spectral sequences, we stored the information
148 collected for every peptide (including the encoding schemes and genomic locus) in an
149 index file with the peptide's data. This reduced the digital memory required to store
150 our sequence data significantly. In addition, based on the locus-tracking approach, we
151 used an automated process to map the peptide spectrum to their genomic loci, which
152 enabled the pipeline for large-scale discovery of NCPs more effectively.

153 **Large-scale Identification of CPs and NCPs in Maize**

154 All the reliably identified peptides from Ensembl protein and customized
155 peptidogenomic databases were combined and used to identify both CPs and NCPs. In
156 total, 748 and 3,932 non-redundant peptides were identified based on Ensembl protein
157 database and customized peptidogenomic database, respectively (Figure 2A; Tables
158 S1 and S2). Of these, 3,315 peptides were specifically identified by the customized
159 peptidogenomic database (Figure 2A). Then, by mapping these peptides to genome
160 loci and applying series filtering steps (see Methods), a total of 2,837 endogenous
161 peptides were unambiguously assigned to a single genomic locus for each of the
162 peptides. Among them, 1,993 (70.3%) NCPs (Table S3) and 844 (29.7%) CPs (Figure
163 2B; Table S4) were identified. The median length of CPs was 16 amino acids, while
164 that of NCPs was 12 amino acids, with significant difference (Figure 2C), and
165 approximately 90% of the peptides were less than 23 amino acids for CPs and 16

166 amino acids for NCPs (Figure S1). Furthermore, the average molecular weight of
 167 NCPs was 1325.22 Da, with 99.25% (1,978) of peptides having a molecular weight
 168 less than 2500 Da. By contrast, the average molecular weight of CPs was 1742.16 Da,
 169 with 91.94% (776) of peptides having a molecular weight less than 2500 Da (Figure
 170 2D and 2E). These results indicated that NCPs constituted a significant portion of
 171 plant peptidome, and showed different characteristics compared with CPs.

172 Distribution Patterns of CPs and NCPs

173 Both CPs and NCPs were found unevenly distributed on the chromosomes of maize
 174 (Figure 3A). For CPs, most peptides were distributed near the telomeres, whereas
 175 NCPs were homogeneously located between centromeres and telomeres of each maize
 176 chromosome (Figure 3B). Furthermore, a total of 138 hot regions (defined by 6 Mb
 177 windows; see Methods) were discovered (Figure 3A). A total of 58 CPs hot regions
 178 containing 446 (52.84%) peptides were observed, whereas 81 NCPs hot regions
 179 containing 545 (27.35%) peptides were present (Figure 3A). Among these hot regions,
 180 one hot region located in chromosome 5 was common for both CPs and NCPs.
 181 Additionally, the number of NCPs in each chromosome was positively correlated with
 182 the chromosomal length ($r=0.07$; $p=0.0099$), but no correlation between the number
 183 of CPs and chromosomal length was detected (Figure 3C).

184 The interval between two adjacent peptides could be used to accurately define
 185 peptides coverage over the genome. We found that 74.88% (632) of CPs were less
 186 than 500 kb apart, whereas only 39.74% (792) of NCPs were within 500 kb of each
 187 other (Figure 3D). We then compared the locations of these peptides to gene models,
 188 798 (94.55%) CPs were found to be located in regions less than 2 kb from canonical
 189 translation start site (TSS), in contrast, this value was 336 (16.86%) for NCPs (Figure
 190 3E). These results reveal the widespread existence of NCPs translation along the
 191 genome and the distinct distribution patterns of CPs and NCPs.

192 To gain further insights into the mechanisms responsible for the generation of
 193 CPs and NCPs, we analyzed the nucleotide sequences of CPs and NCPs source
 194 transcripts to predict their translation start sites. We observed a preponderance of

195 non-AUG translation start sites in both CPs and NCPs (Tables S3 and S4). Although it
 196 was long thought that eukaryotic translation almost always initiates at the AUG start
 197 codon, our results reveal that non-AUG start codons are used at an astonishing
 198 frequency. This finding is consistent with the results of previous peptidomics studies
 199 that more than 90% endogenous peptides started with non-AUG codon (Chen et al.,
 200 2014; Secher et al., 2016; Corbiere et al., 2018). This result also support those of
 201 ribosome profiling and mass spectrometric studies, which demonstrate that most
 202 ORFs contain non-AUG start sites (Ingolia et al., 2011; Slavoff et al., 2013; Na et al.,
 203 2018).

204 **NCPs Derived from both Coding and Noncoding Sequences**

205 By analyzing their origins, 952 (47.77%) NCPs were assigned to the reverse strand in
 206 maize (Figure 4A). Next, by analyzing the location of the NCPs within their
 207 respective gene sources, 1,708 (85.70%) NCPs were derived from intergenic regions,
 208 139 (6.97%) from introns, 89 (4.47%) from out-of-frame exons, 25 (1.25%) from
 209 3'UTRs, 18 (0.90%) from 5'UTRs and 14 (0.70%) from junctions (5'UTR-exon or
 210 intron-exon) (Figure 4B). These results highlight the translation evidence of these
 211 allegedly noncoding sequences.

212 Length analysis showed that the average lengths of NCPs derived from
 213 intergenic regions and out-of-frame exons were longer than that derived from
 214 junctions (Figure 4C). The average lengths of NCPs derived from 3'UTRs and
 215 5'UTRs were the two shortest (Figure 4C). Molecular weight distribution analysis
 216 showed that more than 70% (1,407) of NCPs were less than 1500 Da. The average
 217 molecular weight of NCPs derived from intergenic regions was higher than that
 218 derived from introns, out-of-frame exons, 5'UTRs and 3'UTRs (Figure 4D and Figure
 219 S2A). There was no significant difference among the average isoelectric points (PI)
 220 values of NCPs derived from 3'UTRs, introns, intergenic regions, 5'UTRs,
 221 out-of-frame exons and junctions (Figure 4E and Figure S2B). Taken together, these
 222 results indicated that the identified NCPs represented a wide range of
 223 physicochemical properties and NCPs derived from different gene elements showed

224 different characteristics.

225 **Verification and Validation of NCPs**

226 To verify these identified NCPs, we assigned these peptides to their respective source
227 genomic locus. For example, NCP RMDAHLR was derived from the 5'UTR of
228 gene *Zm00001d029555* (Figure 5A), and NCP ILTVNLKP was derived from the
229 3'UTR of gene *Zm00001d050172* (Figure 5B). Besides NCPs derived from UTRs, we
230 also found a large number of NCPs from intergenic regions and introns. For example,
231 NCP QISVELPGVV was derived from the intergenic region between genes
232 *Zm00001d024336* and *Zm00001d024337* (Figure 5C). NCP EGTPKAVGHRQ was
233 derived from the intron of gene *Zm00001d008363* (Figure 5D). Next, 115 NCPs were
234 synthesized experimentally. The mass spectrometer analysis was performed under the
235 same conditions as were used for peptidogenomic analysis in this study. As shown in
236 Figure 5A-D, the spectra of synthetic peptides RMDAHLR, ILTVNLKP,
237 QISVELPGVV and EGTPKAVGHRQ agreed with the spectral data generated from
238 the peptidogenomic analysis. Verification of the other 111 NCPs was shown in
239 Supplemental Dataset 1.

240 In addition, we performed transcriptomic analyses using published RNA-seq data
241 from maize. These RNA-seq data include circular RNAs, lncRNAs, mRNAs and
242 small RNAs. Most NCPs (1,806, 90.62%) identified in the current study received
243 support from these published databases (Table S3). Among these NCPs, 1,652 were
244 from lncRNA and 859 from circular RNA (Table S3). The results indicated that these
245 identified NCPs were likely produced from allegedly noncoding sequences.

246 Lastly, to validate the identified NCPs with independent methods, the available
247 ribosome profiling datasets of maize were analyzed. Ribosome profiling, also known
248 as Ribo-seq (ribosome sequencing), is a method based on deep sequencing of
249 ribosome-protected fragments. In agreement with translation being the intermediate
250 step between transcription and the proteome, ribosome profiling has a higher
251 predictive value of final protein than mRNA-seq (van Heesch et al., 2019). The
252 ribosome profiling analysis showed that 732 (36.73%) NCPs detected by

253 peptidogenomics were also uncovered by ribosome profiling (Figure 5E; Table S5).
 254 This validation rate of 36.73% between these two methods is consistent with previous
 255 reports (Samandi et al., 2017; van Heesch et al., 2019; Chen et al., 2020a). Among
 256 these NCPs, 564 derived from intergenic regions, 82 from out-of-frame exons, 49
 257 from introns, 15 from 5'UTRs, 14 from 3'UTRs and eight from junctions. The
 258 proportions of the NCPs detected by both methods out of numbers detected by
 259 peptidogenomic analysis were: 33.02% from intergenic regions, 92.13% from
 260 out-of-frame exons, 35.25% from introns, 83.33% from 5'UTRs, 56.00% from
 261 3'UTRs and 57.14% from the junctions (Figure 5F). These NCPs, which were
 262 detected by two different methods, provide a high-confidence collection of NCPs for
 263 further studies. We speculate that those NCPs, which were detected only by
 264 peptidogenomics, were either erroneous calls or stable peptides from unstable RNAs.

265 **NCPs are Enriched in Regions Associated with Phenotypic Variations 266 and Domestication Selection**

267 In maize, coding regions only comprise a small fraction of the whole genome, and the
 268 vast majority of the genome has been considered noncoding regions. Genome-wide
 269 association study and quantitative trait locus (QTLs) analysis have identified a lot of
 270 functional elements in the noncoding regions in maize (Liu et al., 2017). The fact that
 271 1,993 (70.3%) NCPs were derived from noncoding sequences prompts us to believe
 272 that they are of significant functional relevance. Therefore, we examined the
 273 enrichment of these NCPs with identified QTLs underlying various traits, and with
 274 those regions presumed under domestication selection.

275 Compared to randomly selected genomic sequence with same distance
 276 distribution and number (see Methods), it was revealed that significant single
 277 nucleotide polymorphisms (SNPs) associated with plant traits appeared to be
 278 significantly enriched within the regions of NCPs ($P < 0.02$, Upper-tail test; Figure
 279 6A; Table S6). Considering the presence of genetic linkage in association mapping,
 280 we further extended the positions of associated SNPs to the flanking 20 kb regions.
 281 Statistical analysis showed that these NCPs were more significantly enriched at the

282 QTL regions compared to the random regions ($P < 7.4e-06$; Figure 6B; Table S7).
 283 Among the significant enriched SNPs, several were found exactly located within
 284 NCPs, which showed associations to various phenotypes including kernel length,
 285 disease (maize rough dwarf virus, MRDV), oil and amino acid contents (Figure 6C;
 286 Table S6). For instance, an isoleucine-threonine transition at one significant SNP
 287 (chr1.s_244454699, A > G; $P < 9.42e-5$) associated with kernel length, was located
 288 within the NCP KTYSIIIYFIHVGH, which was mapped to 13 kb upstream
 289 noncoding regions of gene *Zm00001d032949* (Uncharacterized) (Figure 6D). Another
 290 significant SNP (chr3.s_136872577, C > T; $P < 2.09e-07$) related to oil content,
 291 resulting in a transition from proline to leucine was associated with the NCP
 292 LELKLIHSHPN, which was mapped to 5 kb upstream noncoding regions of gene
 293 *Zm00001d041769* (Figure 6E). These results reveal the potential functions of these
 294 NCPs in the regulation of plant phenotypes.

295 The relationship between domestication and NCPs was also investigated.
 296 Compared to randomly selected genomic sequences with the same distance
 297 distribution and number, it was found that the NCPs were enriched within the
 298 candidate regions that are associated with domestication selection ($p < 7.3e-6$,
 299 Upper-tail test; Figure 6F). A total of 55 NCPs were identified within the
 300 domestication candidate regions (Table S8). While further validations are highly
 301 needed to explore which domesticated traits are exactly affected and what's the indeed
 302 mechanism, this result, for the first time as far as we know, unveils the likely
 303 inclusion of NCPs during domestication, providing another hidden layer of functional
 304 importance of NCPs.

305 **The Applicability of the Peptidogenomics Pipeline to Arabidopsis**

306 To extend this pipeline to other plants, the dicot model plant Arabidopsis was used to
 307 test the wider applicability of peptidogenomic method. As a result, 2,353 and 3,871
 308 non-redundant peptides were identified by the Ensembl protein database and
 309 customized peptidogenomic database (Tables S9 and S10), respectively. Of these,
 310 2,270 peptides were specifically identified by the customized peptidogenomic

311 database (Figure 7A). In total, 1,860 (44.04%) NCPs (Table S11) and 2,363 (55.96%)
 312 CPs were obtained in Arabidopsis (Table S12). The median length of NCPs was 11
 313 amino acids, which was shorter than that of CPs (13 amino acids) (Figure 7B).
 314 Furthermore, the average molecular weight of NCPs (1208.34 Da) was lower than that
 315 of CPs (1420.89 Da) (Figure S3). In addition, we found that the NCPs identified in
 316 Arabidopsis have shorter peptide length and lower molecular weight than that in
 317 maize (Table 13).

318 By analyzing the origins of NCPs, 943 (50.70%) NCPs were from the reverse
 319 strand (Figure 7C). By analyzing the locations of the NCPs within their respective
 320 gene sources, 666 (35.81%) NCPs were derived from intergenic regions, 239 (12.85%)
 321 from introns, 651 (35.00%) from out-of-frame exons, 91 (4.89%) from 3'UTRs, 63
 322 (3.39%) from 5'UTRs and 150 (8.06%) from junctions (Figure 7D). The number of
 323 NCPs derived from intergenic regions in Arabidopsis was lower than that in maize,
 324 whereas the number of NCPs from other gene elements in Arabidopsis were higher
 325 than that in maize (Table S13). Length analysis showed that the average length of
 326 NCPs derived from 3'UTRs was the longest and that from introns the shortest (Figure
 327 7E). The average molecular weight of NCPs derived from out-of-frame exons was
 328 higher than that from 5'UTRs and intergenic regions (Figure 7F and Figure S4A). The
 329 average PI value of NCPs derived from out-of-frame exons and junctions were higher
 330 than that from introns (Figure 7G and Figure S4B).

331 Taken together, these results show that the developed peptidogenomic pipeline
 332 can also be used in dicot plants such as Arabidopsis. The translation of unannotated
 333 transcripts is widespread in both monocot and dicot plants, though they may have
 334 different translation patterns.

335 **DISCUSSION**

336 Endogenous peptides are formed mainly by protein degradation, gene-encoding and
 337 gene-independent enzymatic formation *in vivo* (Peng et al., 2020). The emergence of
 338 peptidomics makes it possible for large-scale identification of endogenous peptides
 339 extracted from tissues (Slavoff et al., 2013; Secher et al., 2016). However, the study of

340 the peptidomics can be particularly challenging due to nonspecific protease digestion
341 during sample preparation (Farrokhi et al., 2008; Secher et al., 2016). Despite the
342 wide use of protease inhibitors in plant peptide extraction, studies in animals and
343 humans have demonstrated that protease inhibitors are not effective enough in
344 preventing peptide degradation (Svensson et al., 2003; Parkin et al., 2005). Recently,
345 heat stabilization, such as focused microwave radiation, integrated with protease
346 inhibitors has been successfully used in animals to minimize proteolytic activity prior
347 to peptide isolation (Secher et al., 2016). However, similar attempt has not been
348 experimented in plants so far.

349 Plant cells are more complex than animal cells due to the presence of additional
350 components such as cell wall, large vacuoles and chloroplast, making the isolation of
351 complete endogenous peptides in plants more challenging. In this study, in addition to
352 the combination of heat stabilization by water bath and plant protease inhibitors to
353 minimize nonspecific protease digestion in the peptide extraction, TCA-acetone
354 precipitation was also included in the extraction protocol. TCA/acetone precipitation
355 is very useful for removing interfering compounds, such as polysaccharides,
356 polyphenols, pigments and lipids in plants (Mechin et al., 2007). Therefore, this step
357 can help limit the interference of non-protein or non-peptide compounds during
358 endogenous peptides extraction. We speculate that the protease associated nonspecific
359 degradation during peptide extraction will be a long-lasting issue as there is no
360 effective extraction protocol to completely prevent this from occurring. Therefore,
361 more efforts should be made to develop a more effective peptide extraction protocol
362 that can retain endogenous peptides in the states as they were *in vivo* for peptidomics
363 study. In addition, it should be noted that the peptides from protein degradation within
364 the cell is also another type of endogenous peptides in addition to those produced
365 from gene-encoding (Peng et al., 2020). Protein degradation ubiquitously occurs in
366 living organisms and the enzymatic degradation behavior of proteins is closely related
367 to precursor protein status and enzyme activity in living organisms (Rubinsztein,
368 2006). Therefore, peptidomic data is also a good resource for the assessment of the
369 potential protease/peptidase activity involving in the hydrolysis process, though this

370 topic is beyond the scope of this study.

371 Standard peptidomics approaches identify peptides by matching experimentally
372 observed spectra to databases of predicted spectra based on annotated genes. However,
373 such approach would not identify NCPs. The most effective strategy to do so is to
374 integrate peptidomics with the six-frame translation of genome, which is referred as
375 peptidogenomics (Kersten et al., 2011; Slavoff et al., 2013). Database derived from the
376 six-frame translation of the entire genome can be used to identify peptides encoded in
377 any genomic region (Castellana et al., 2014; Nesvizhskii, 2014; Yang et al., 2018).
378 Peptidogenomics has already proven its value in identifying peptides at the
379 genome-scale in microorganisms and humans (Kersten et al., 2011; Liu et al., 2011;
380 Nguyen et al., 2013; Slavoff et al., 2013; Mohimani and Pevzner, 2016; Mohimani et
381 al., 2018). In this study, we combined the peptidomics with a customized
382 peptidogenomic database derived from six-frame translation and Ensembl protein
383 databases to generate a peptidogenomic pipeline for both maize and Arabidopsis. To
384 the best of our knowledge, this is the first report on a peptidogenomic pipeline to
385 analyze NCPs in plants. With this strategy, 1,993 and 1,860 NCPs have been
386 identified in maize and Arabidopsis, respectively. The present study demonstrates that
387 integrative peptidogenomic strategies can provide a more holistic overview of the
388 peptidome to not only identify CPs but also NCPs. The results showed that a sizeable
389 proportion of peptides was found to be NCPs, indicating that many previously alleged
390 noncoding sequences, including 5'UTRs, 3'UTRs, intergenic regions and introns are
391 actually translatable.

392 Recently, the translation of lncRNAs has gained increasing attention (Kim et al.,
393 2014; Saghatelian and Couso, 2015; Ransohoff et al., 2018). For example, a peptide
394 encoded by a lncRNA was identified as myoregulin, which acts as an important
395 regulator of calcium uptake in skeletal muscle (Anderson et al., 2015). A peptide
396 encoded from a lncRNA epithelial cell program regulator (*EPR*) controls epithelial
397 proliferation (Rossi et al., 2019). In addition, by overexpression and mutation analysis,
398 peptides encoded by lncRNAs were shown to be involved in the regulation of growth
399 and differentiation in moss (Fesenko et al., 2019). In the present study, 1,652 NCPs

400 derived from lncRNA have been identified in maize, and future characterization of
 401 these NCPs will be an important milestone in understanding the function of plant
 402 lncRNAs.

403 Upstream ORFs (uORFs) and their encoded peptides have been intensively
 404 investigated due to their potential to regulate the translation of downstream main
 405 ORFs (mORFs) (Hellens et al., 2016; Hsu and Benfey, 2018). The translation of these
 406 uORFs can also be regulated in response to developmental or environmental cues
 407 (Starck et al., 2016; Yin et al., 2019). In this study, we identified 18 and 63 NCPs
 408 derived from 5'UTRs in maize and Arabidopsis, respectively. Among the 18 maize
 409 NCPs, 15 NCPs were also uncovered by previous ribosome profiling studies (Lei et
 410 al., 2015; Chotewutmontri and Barkan, 2016; Zoschke et al., 2017; Chotewutmontri
 411 and Barkan, 2018; Jiang et al., 2019), which further supports the results of the
 412 peptidogenomic analysis in the present study. In contrast to NCPs derived from the
 413 5'UTRs of genes, NCPs from 3'UTRs have attracted little attention because they have
 414 been considered to be noncoding for a long time (Ingolia et al., 2011). Until only
 415 recently, the presence of peptides assigned to 3'UTRs was identified, for example, in
 416 moss (Fesenko et al., 2019). In our study, we identified 25 and 91 NCPs that derived
 417 from 3'UTRs in maize and Arabidopsis, which further suggests that 3'UTRs encoded
 418 peptides deserve much more attention as these peptides may have vital biological
 419 roles in organisms.

420 Many maize QTLs have been found to be highly associated with noncoding
 421 regions (Clark et al., 2006; Silvio et al., 2007; Studer et al., 2011; Castelletti et al.,
 422 2014; Huang et al., 2018). Recently, we also examined several cases of intergenic
 423 QTLs that regulate traits by chromatin loops (Li et al., 2019; Peng et al., 2019).
 424 Apparently, it is important to study the regulatory elements in the noncoding
 425 sequences for a better understanding of the biological mechanisms underlying
 426 phenotypic traits. In this study, we found that NCPs were significantly enriched within
 427 QTLs regions. For example, NCPs were enriched within regions associated with
 428 disease resistance, kernel length, amino acid and oil contents, indicating the important
 429 functionality of NCPs in regulating these traits. Domestication is a tractable system

430 for subsequent evolutionary changes. Identification of genes involved in
 431 domestication will help us to understand the process of domestication and to
 432 accelerate the process of domesticating new crops (Wang et al., 2018). Several recent
 433 studies have used morphological, genetic, genomic and archaeological techniques to
 434 determine the progressive fixation of different domestication genes in maize (da
 435 Fonseca et al., 2015; Liu et al., 2015; Vallebueno-Estrada et al., 2016). However, to
 436 date, the molecular genetic architecture of maize domestication remains unclear. The
 437 result of statistical analysis in this study showed significant enrichment of NCPs in
 438 the domestication selection regions, which may uncover the underlying functional
 439 sites for the evolution of the maize due to selection.

440 Taken together, in contrast to previous attempts of using computational
 441 approaches or ribosome profiling strategy to discover unannotated plant coding
 442 sequences, we directly and successfully identified large-scale plant NCPs based on the
 443 integrated peptidogenomic pipeline. The identification of NCPs reveals that many
 444 5'UTRs, 3'UTRs, intergenic regions, introns, and junctions are translated and some
 445 likely express functional peptides. These findings also provide insights into the
 446 discovery of novel functional genes or proteins through the characterization of NCPs
 447 in a wider array of plants.

448 MATERIALS AND METHODS

449 Sample Preparation

450 The maize inbred line B73 was grown in a greenhouse under a 15-h light (28 °C)/9-h
 451 (25 °C) dark photoperiod to 3 leaf stage. *Arabidopsis thaliana* (Columbia-0) was
 452 grown in a greenhouse under a 16-h light (22 °C)/8-h (21 °C) dark photoperiod to 4
 453 leaf stage. Three replicates were applied for each species. The collected leaves were
 454 quickly frozen in liquid nitrogen and stored at -80 °C until analyzed.

455 Peptide Extraction

456 Maize and *Arabidopsis* leaves (2 g) as described above were quickly grounded in
 457 liquid nitrogen, respectively. The powder was firstly heated in the water at 95 °C for 5
 458 min. The samples were then precipitated in 10% (w/v) trichloroacetic acid/acetone

459 solution at -20 °C for 1 h, and the precipitate was washed with cold acetone until the
460 supernatant was colorless. The supernatant was discarded, the vacuum-dried
461 precipitate transferred to 1% TFA solution containing plant protease inhibitor cocktail
462 (Sigma, America), and incubated for an hour at 4 °C. It should be noted that TFA can't
463 be added before heat stabilization, because TFA is a strongly irritating liquid which
464 decomposes and emits toxic fluoride gas when heated. The fractions were
465 ultrasonicated on ice (40 W, 6 s ultrasonic at a time, every 8 s, and 5 times) and then
466 centrifuged at 10,000 ×g for 20 min at 4 °C. The supernatants were filtered through
467 10-kDa molecular weight cutoff centrifuge filter (Millipore, MA, USA) according to
468 the manufacturer's instructions. Peptide mixtures were desalted using C18 Cartridges
469 (Empore, SPE Cartridges C18, 7 mm inner diameter, 3 mL volume, Sigma). The
470 peptide fractions were vacuum-evaporated using a vacuum centrifugation
471 concentrator and reconstituted in 40 µl of 0.1% TFA solution for LC-MS/MS analysis.

472 **LC-MS/MS Analysis**

473 For endogenous peptide profiling, MS experiments were performed on a Q Exactive
474 mass spectrometer as described previously (Wang et al., 2019). Five µg of peptide
475 mixture was loaded onto a C18-reversed phase column (Thermo Scientific Easy
476 Column, 10 cm length, 75 µm inner diameter, 3 µm resin) in buffer A (2% acetonitrile
477 and 0.1% formic acid) and separated with a linear gradient of buffer B (80%
478 acetonitrile and 0.1% formic acid) at a flow rate of 250 nL/min controlled by
479 IntelliFlow technology over 120 min. MS data were acquired using a data-dependent
480 top10 method by dynamically choosing the most abundant precursor ions from the
481 survey scan (300-1800 m/z) for higher-energy collisional dissociation (HCD)
482 fragmentation. The determination of the target value was based on predictive
483 Automatic Gain Control. The dynamic exclusion duration was 25 s. Survey scans
484 were acquired at a resolution of 70,000 at m/z 200 and resolution for HCD spectra
485 was set to 17,500 at m/z 200. The normalized collision energy was 30 eV and the
486 underfill ratio, which specified the minimum percentage of the target value likely to
487 be reached at maximum fill time, was defined as 0.1%. The instrument was run with
488 peptide recognition mode enabled.

489 Peptide Database Construction

490 The complete genomes of maize and Arabidopsis were downloaded from Ensembl
491 Plants (ftp://ftp.ensemblgenomes.org/pub/plants/release-41/fasta/zea_mays/dna/; and
492 ftp://ftp.ensemblgenomes.org/pub/plants/release-45/fasta/arabidopsis_thaliana/dna/)
493 in FASTA format. The putative peptide database was derived from the six-frame
494 translation of genomic sequences using EMBOSS:6.6.0. Peptides were terminated
495 whenever a stop codon was encountered. Then the next peptide was started at the next
496 nucleotide following the previous stop codon. Instances of ambiguous nucleotides
497 (represented by 'N' in the genome sequence) were replaced with random nucleotides;
498 other ambiguous characters were also replaced with random nucleotides depending
499 upon their symbol. The genomic coordinates and orientation were recorded for each
500 peptide. Resulting amino acid sequences for each chromosome were recorded in a
501 FASTA formatted sequence file.

502 Peptide Identification by Mascot

503 The Mascot search engine (Matrix Science) was used to search against both the
504 Ensembl protein for maize
505 (ftp://ftp.ensemblgenomes.org/pub/plants/release-41/fasta/zea_mays), and
506 Arabidopsis
507 (ftp://ftp.ensemblgenomes.org/pub/plants/release-45/fasta/arabidopsis_thaliana/pep/),
508 and the customized peptidogenomic databases to identify peptides. Mass tolerances
509 on precursor and fragment ions were set to 5 ppm and 0.02 Da, respectively. The
510 Mascot score (≥ 25) and false discovery rate (FDR < 0.05) were applied to achieve
511 final peptides for the Ensembl protein database. The same Mascot score was then
512 applied to the peptide list identified with the customized peptidogenomic database as
513 described previously (Laumont et al., 2016). Raw data files were converted to peptide
514 maps comprising m/z values, charge states, retention time and intensity for all
515 detected ions above a threshold of 8,000 counts.

516 In order to obtain quantitative information for the peptides, the MS data were
517 analyzed using MaxQuant software (version 1.3.0.5). The MS data were searched
518 against the identified peptide sequences. An initial search was set at a precursor mass

519 window of 6 ppm, followed by an enzymatic cleavage rule of none and a mass
 520 tolerance of 20 ppm for fragment ions. The cutoff of global FDR for peptide
 521 identification was set to 0.01. Peptide intensities were used to indicate quantitative
 522 information of peptide.

523 **Identification of CPs and NCPs**

524 Peptides identified from Ensembl protein and customized peptidogenomic databases
 525 were combined and filtered with the stringent FDR cutoff (score ≥ 25 ; FDR < 0.05).
 526 The resulting peptides were assigned to their respective source genes and their
 527 MS/MS spectra were manually verified. Then, we mapped the subset of
 528 peptide-encoding regions to discard peptides coming from multiple locations in the
 529 genome (1,207 peptides for maize and 410 peptides for Arabidopsis). To determine
 530 the type of sequence (within the source gene) generating each peptide, we used the
 531 intersect function of the BEDTools suite to the bed file of the candidates as well as the
 532 Ensembl gff file. Peptides derived from annotated CDSs or conventional open reading
 533 frames were classified as CPs. Peptides derived from intergenic regions, UTRs,
 534 different reading frames from annotated CDSs, introns and various types of junctions
 535 (UTR-exon or exon-intron) were classified as NCPs.

536 **Peptide Distribution at the Genome Level**

537 Peptide density was calculated using a sliding window of 6 Mb with 3 Mb steps. Hot
 538 regions were defined as the peptide count of more than 10. We downloaded the
 539 annotated maize genome from <https://plants.ensembl.org/index.html> and extracted the
 540 physical coordinates of TSSs. We searched for the closest TSS for each peptide to
 541 draw a frequency plot of distance between each peptide and its TSS. To accurately
 542 estimate the peptide number at the chromosome level, position of both CPs and NCPs
 543 was divided by chromosome arm length.

544 **Verification of NCPs Using Synthetic Peptides**

545 The peptide sequences were chosen from different categories of NCPs identified by
 546 the peptidogenomic analysis and synthesized by GL Biochem (Shanghai) Ltd. Dried
 547 peptides were diluted with 0.1% formic acid (Yang et al., 2018), and each synthetic
 548 peptide was separately subjected to Q Exactive mass spectrometer for MS analysis

549 with the same parameters as those used for the peptidogenomic analysis.

550 **RNA-seq and Ribosome Profiling Analysis**

551 RNA-seq datasets were retrieved from the NCBI short Read Archive database
552 (<https://www.ncbi.nlm.nih.gov/sra>). These datasets including circular RNAs (Jeck et
553 al., 2013), lncRNAs (Lv et al., 2016; Zhu et al., 2017), mRNAs (Lei et al., 2015; Han
554 et al., 2019), small RNAs (He et al., 2019). In addition, the publicly available
555 ribosome profiling datasets of maize (Lei et al., 2015; Chotewutmontri and Barkan,
556 2016; Zoschke et al., 2017; Chotewutmontri and Barkan, 2018; Jiang et al., 2019)
557 were also analyzed. The maize genome sequences and annotation files were obtained
558 from the Ensembl Plants (https://plants.ensembl.org/Zea_mays/Info/Index). After
559 filtering out the low-quality reads, the remaining reads were mapped to the maize
560 genome. Then, the read count was calculated for each NCP.

561 **Association Analysis of NCPs with SNP/regions Associated with a Collection of 562 Traits and the Regions Under Domestication Selection**

563 A genome-wide association study was performed using a global germplasm collection
564 of 527 elite maize inbred lines (Li et al., 2013) using the mixed-linear-model based on
565 previously reported traits, including kernel-related yield traits (Liu et al., 2017),
566 diseases (Chen et al., 2015), as well as kernel oil (Li et al., 2013) and amino acid
567 contents (Deng et al., 2017). SNPs called from the whole-genome shotgun (~20x for
568 each line) sequences generated by a recent study (Yang et al., 2019) were used in
569 association analysis. We generated 100 random genomic sets as background, each
570 assigned with the same features as NCPs, including the total number, the number
571 along different chromosomes, and the peptide length distribution (Figure S5). The 100
572 random sets were used to estimate the mean and standard deviation of the normal
573 distribution for background overlapping ratios. The *p*-values of enrichment of the
574 observed ratio compared to the normal background distribution were calculated using
575 the “pnorm” function (with lower.tail = FALSE) of R, representing the upper tail
576 *p*-value of the test statistic and indicating the probability of observed value exceeding
577 the expected distribution. Candidate regions associated with domestication were
578 identified by comparing the 527 maize inbred lines to 183 teosinte samples, and the

579 test of enrichment was estimated using the aforementioned test as QTL analysis
580 (Figure S6).

581 **Data Analysis and Visualization**

582 Unless stated otherwise, analysis and visualization were performed using R. All code
583 are available on request to the corresponding author.

584 **ACCESSION NUMBERS**

585 All raw mass spectrometry data from this study have been deposited in the
586 ProteomeXchange Consortium via the PRIDE partner repository with dataset
587 identifiers PXD017080 and PXD017081.

588 **AUTHOR CONTRIBUTIONS**

589 L. W. and J. Y. designed the project. S. Wang, J. Z., L. T., X. C. and X. J. conducted
590 experiments. S. Wang, J. Z., H. L., X. L., X. Z., Y. C., L. T. and S. Wu analyzed the
591 data. S. Wang, L. T., H. L., S. Wu, J. Y. and L. W. wrote the manuscript. L. W.
592 supervised the project. All authors read and approved the manuscript.

593 **ACKNOWLEDGMENTS**

594 We thank Dr. Steven P. Briggs for helpful discussions. We thank Dr. Anguo Sun and
595 Dr. Yanwen Xiang for technical assistance. This work is supported by the National
596 Natural Science Foundation of China (no. 31872872 and U1804113), National Key
597 Research and Development Program of China (no. 2016YFD0101003) and Henan
598 Association for Science and Technology.

599 **COMPETING INTERESTS**

600 The authors declare no competing interests.

601 **REFERENCES**

- 602 Anderson, D.M., Anderson, K.M., Chang, C.L., Makarewich, C.A., Nelson, B.R.,
 603 McAnally, J.R., Kasaragod, P., Shelton, J.M., Liou, J., Bassel-Duby, R., et al.
 604 (2015). A micropeptide encoded by a putative long noncoding RNA regulates
 605 muscle performance. *Cell* 160:595-606.
- 606 Andrews, S.J., and Rothnagel, J.A. (2014). Emerging evidence for functional peptides
 607 encoded by short open reading frames. *Nat. Rev. Genet.* 15:193-204.
- 608 Banting, F.G., and Best, C.H. (2007). The internal secretion of the pancreas
 609 (Reprinted from the Journal of Laboratory and Clinical Medicine, vol 7, pg
 610 251-266, 1922). *Indian J. Med. Res.* 125:A251-A266.
- 611 Bazin, J., Baerenfaller, K., Gosai, S.J., Gregory, B.D., Crespi, M., and Bailey-Serres, J.
 612 (2017). Global analysis of ribosome-associated noncoding RNAs unveils new
 613 modes of translational regulation. *Proc. Natl. Acad. Sci. USA*
 614 114:E10018-E10027.
- 615 Blanvillain, R., Young, B., Cai, Y.M., Hecht, V., Varoquaux, F., Delorme, V., Lancelin,
 616 J.M., Delsenay, M., and Gallois, P. (2011). The *Arabidopsis* peptide kiss of
 617 death is an inducer of programmed cell death. *Embo J.* 30:1173-1183.
- 618 Casson, S.A., Chilley, P.M., Topping, J.F., Evans, I.M., Souter, M.A., and Lindsey, K.
 619 (2002). The POLARIS gene of *Arabidopsis* encodes a predicted peptide
 620 required for correct root growth and leaf vascular patterning. *Plant Cell*
 621 14:1705-1721.
- 622 Castellana, N.E., Payne, S.H., Shen, Z., Stanke, M., Bafna, V., and Briggs, S.P. (2008).
 623 Discovery and revision of *Arabidopsis* genes by proteogenomics. *Proc. Natl.*
 624 *Acad. Sci. USA* 105:21034-21038.
- 625 Castellana, N.E., Shen, Z., He, Y., Walley, J.W., Cassidy, C.J., Briggs, S.P., and Bafna,
 626 V. (2014). An automated proteogenomic method uses mass spectrometry to
 627 reveal novel genes in *Zea mays*. *Mol. Cell. Proteomics* 13:157-167.
- 628 Castelletti, S., Tuberosa, R., Pindo, M., and Salvi, S. (2014). A MITE Transposon
 629 insertion is associated with differential methylation at the maize flowering

- 630 time QTL *Vgt1*. G3 (Bethesda) 4:805-812.
- 631 Chen, G., Wang, X., Hao, J., Yan, J., and Ding, J. (2015). Genome-wide association
632 implicates candidate genes conferring resistance to maize rough dwarf disease
633 in maize. PLoS One 10:e0142001.
- 634 Chen, J., Brunner, A.D., Cogan, J.Z., Nunez, J.K., Fields, A.P., Adamson, B., Itzhak,
635 D.N., Li, J.Y., Mann, M., Leonetti, M.D., et al. (2020a). Pervasive functional
636 translation of noncanonical human open reading frames. Science
637 367:1140-1146.
- 638 Chen, Q.J., Deng, B.H., Gao, J., Zhao, Z.Y., Chen, Z.L., Song, S.R., Wang, L., Zhao,
639 L.P., Xu, W.P., Zhang, C.X., et al. (2020b). An miRNA-encoded small peptide,
640 vvi-miPEP171d1, regulates adventitious root formation. Plant Physiol.
641 <https://doi:10.1104/pp.20.00197>
- 642 Chen, Y.L., Lee, C.Y., Cheng, K.T., Chang, W.H., Huang, R.N., Nam, H.G., and Chen,
643 Y.R. (2014). Quantitative peptidomics study reveals that a wound-induced
644 peptide from PR-1 regulates immune signaling in tomato. Plant Cell
645 26:4135-4148.
- 646 Chotewutmontri, P., and Barkan, A. (2016). Dynamics of chloroplast translation
647 during chloroplast differentiation in maize. PLoS Genet. 12:e1006106.
- 648 Chotewutmontri, P., and Barkan, A. (2018). Multilevel effects of light on ribosome
649 dynamics in chloroplasts program genome-wide and psbA-specific changes in
650 translation. PLoS Genet. 14:e1007555.
- 651 Clark, R.M., Tina Nussbaum, W., Pablo, Q., and John, D. (2006). A distant upstream
652 enhancer at the maize domestication gene tb1 has pleiotropic effects on plant
653 and inflorescent architecture. Nat. Genet. 38:594-597.
- 654 Corbiere, A., Walet-Balieu, M.L., Chan, P., Basille-Dugay, M., Hardouin, J., and
655 Vaudry, D. (2018). A peptidomic approach to characterize peptides involved in
656 cerebellar cortex development leads to the identification of the neurotrophic
657 effects of nociceptin. Mol. Cell. Proteomics 17:1737-1749.
- 658 Couso, J.P., and Patraquim, P. (2017). Classification and function of small open
659 reading frames. Nat. Rev. Mol. Cell Biol. 18:575-589.

- 660 da Fonseca, R.R., Smith, B.D., Wales, N., Cappellini, E., Skoglund, P., Fumagalli, M.,
661 Samaniego, J.A., Caroe, C., Avila-Arcos, M.C., Hufnagel, D.E., et al. (2015).
662 The origin and evolution of maize in the Southwestern United States. Nat.
663 Plants 1:14003.
- 664 De Coninck, B., Carron, D., Tavormina, P., Willem, L., Craik, D.J., Vos, C., Thevissen,
665 K., Mathys, J., and Cammue, B.P.A. (2013). Mining the genome of
666 *Arabidopsis thaliana* as a basis for the identification of novel bioactive
667 peptides involved in oxidative stress tolerance. J. Exp. Bot. 64:5297-5307.
- 668 Deng, M., Li, D., Luo, J., Xiao, Y., Liu, H., Pan, Q., Zhang, X., Jin, M., Zhao, M., and
669 Yan, J. (2017). The genetic architecture of amino acids dissection by
670 association and linkage analysis in maize. Plant Biotechnol. J. 15:1250-1263.
- 671 Farrokhi, N., Whitelegge, J.P., and Brusslan, J.A. (2008). Plant peptides and
672 peptidomics. Plant Biotechnol. J. 6:105-134.
- 673 Fesenko, I., Kirov, I., Kniazev, A., Khazigaleeva, R., Lazarev, V., Kharlampieva, D.,
674 Grafskaia, E., Zgoda, V., Butenko, I., Arapidi, G., et al. (2019). Distinct types
675 of short open reading frames are translated in plant cells. Genome Res.
676 29:1464-1477.
- 677 Han, L., Mu, Z., Luo, Z., Pan, Q., and Li, L. (2019). New lncRNA annotation reveals
678 extensive functional divergence of the transcriptome in maize. J. Integr. Plant
679 Biol. 61:394-405.
- 680 Hanada, K., Zhang, X., Borevitz, J.O., Li, W.H., and Shiu, S.H. (2007). A large
681 number of novel coding small open reading frames in the intergenic regions of
682 the *Arabidopsis thaliana* genome are transcribed and/or under purifying
683 selection. Genome Res. 17:632-640.
- 684 Harvey, A.L., Edrada-Ebel, R., and Quinn, R.J. (2015). The re-emergence of natural
685 products for drug discovery in the genomics era. Nat. Rev. Drug Discov.
686 14:111-129.
- 687 He, J., Jiang, Z., Gao, L., You, C., Ma, X., Wang, X., Xu, X., Mo, B., Chen, X., and
688 Liu, L. (2019). Genome-wide transcript and small RNA profiling reveals
689 transcriptomic responses to heat stress. Plant Physiol. 181:609-629.

- 690 Hellens, R.P., Brown, C.M., Chisnall, M.A.W., Waterhouse, P.M., and Macknight,
691 R.C. (2016). The emerging world of small ORFs. Trends Plant Sci.
692 21:317-328.
- 693 Hinnebusch, A.G., Ivanov, I.P., and Sonenberg, N. (2016). Translational control by
694 5'-untranslated regions of eukaryotic mRNAs. Science 352:1413-1416.
- 695 Hsu, P.Y., and Benfey, P.N. (2018). Small but mighty: Functional peptides encoded by
696 small ORFs in plants. Proteomics 18:e1700038.
- 697 Hsu, P.Y., Calviello, L., Wu, H.L., Li, F.W., Rothfels, C.J., Ohler, U., and Benfey, P.N.
698 (2016). Super-resolution ribosome profiling reveals unannotated translation
699 events in Arabidopsis. Proc. Natl. Acad. Sci. USA 113:E7126-E7135.
- 700 Huang, C., Sun, H.Y., Xu, D.Y., Chen, Q.Y., Liang, Y.M., Wang, X.F., Xu, G.H., Tian,
701 J.G., Wang, C.L., Li, D., et al. (2018). *ZmCCT9* enhances maize adaptation to
702 higher latitudes. Proc. Natl. Acad. Sci. USA 115:E334-E341.
- 703 Hurst, L.D. (2002). The *Ka/Ks* ratio: diagnosing the form of sequence evolution.
704 Trends Genet. 18:486.
- 705 Ingolia, N.T. (2016). Ribosome footprint profiling of translation throughout the
706 genome. Cell 165:22-33.
- 707 Ingolia, N.T., Ghaemmaghami, S., Newman, J.R.S., and Weissman, J.S. (2009).
708 Genome-wide analysis *in vivo* of translation with nucleotide resolution using
709 ribosome profiling. Science 324:218-223.
- 710 Ingolia, N.T., Lareau, L.F., and Weissman, J.S. (2011). Ribosome profiling of mouse
711 embryonic stem cells reveals the complexity and dynamics of mammalian
712 proteomes. Cell 147:789-802.
- 713 Jackson, R., Kroehling, L., Khitun, A., Bailis, W., Jarret, A., York, A.G., Khan, O.M.,
714 Brewer, J.R., Skadow, M.H., Duizer, C., et al. (2018). The translation of
715 non-canonical open reading frames controls mucosal immunity. Nature
716 564:434-438.
- 717 Jeck, W.R., Sorrentino, J.A., Wang, K., Slevin, M.K., Burd, C.E., Liu, J., Marzluff,
718 W.F., and Sharpless, N.E. (2013). Circular RNAs are abundant, conserved, and
719 associated with ALU repeats. RNA 19:141-157.

- 720 Jiang, J., Chai, X., Manavski, N., Williams-Carrier, R., He, B., Brachmann, A., Ji, D.,
721 Ouyang, M., Liu, Y., Barkan, A., et al. (2019). An RNA chaperone-like protein
722 plays critical roles in chloroplast mRNA stability and translation in
723 *Arabidopsis* and maize. *Plant Cell* 31:1308-1327.
- 724 Kastenmayer, J.P., Ni, L., Chu, A., Kitchen, L.E., Au, W.C., Yang, H., Carter, C.D.,
725 Wheeler, D., Davis, R.W., Boeke, J.D., et al. (2006). Functional genomics of
726 genes with small open reading frames (sORFs) in *S. cerevisiae*. *Genome Res.*
727 16:365-373.
- 728 Kersten, R.D., Yang, Y.L., Xu, Y., Cimermancic, P., Nam, S.J., Fenical, W., Fischbach,
729 M.A., Moore, B.S., and Dorrestein, P.C. (2011). A mass spectrometry-guided
730 genome mining approach for natural product peptidogenomics. *Nat. Chem.
731 Biol.* 7:794-802.
- 732 Khitun, A., Ness, T.J., and Slavoff, S.A. (2019). Small open reading frames and
733 cellular stress responses. *Mol. Omics* 15:108-116.
- 734 Kim, M.S., Pinto, S.M., Getnet, D., Nirujogi, R.S., Manda, S.S., Chaerkady, R.,
735 Madugundu, A.K., Kelkar, D.S., Isserlin, R., Jain, S., et al. (2014). A draft map
736 of the human proteome. *Nature* 509:575-581.
- 737 Kondo, T., Plaza, S., Zanet, J., Benrabah, E., Valenti, P., Hashimoto, Y., Kobayashi, S.,
738 Payre, F., and Kageyama, Y. (2010). Small peptides switch the transcriptional
739 activity of shavenbaby during drosophila embryogenesis. *Science*
740 329:336-339.
- 741 Kurihara, Y., Makita, Y., Shimohira, H., Fujita, T., Iwasaki, S., and Matsui, M. (2020).
742 Translational landscape of protein-coding and non-protein-coding RNAs upon
743 light exposure in *Arabidopsis*. *Plant Cell Physiol.* 61:536-545.
- 744 Laumont, C.M., Daouda, T., Laverdure, J.P., Bonneil, E., Caron-Lizotte, O., Hardy,
745 M.P., Granados, D.P., Durette, C., Lemieux, S., Thibault, P., et al. (2016).
746 Global proteogenomic analysis of human MHC class I-associated peptides
747 derived from non-canonical reading frames. *Nat. Commun.* 7:10238.
- 748 Lauressergues, D., Couzigou, J.M., Clemente, H.S., Martinez, Y., Dunand, C., Becard,
749 G., and Combier, J.P. (2015). Primary transcripts of microRNAs encode

- 750 regulatory peptides. *Nature* 520:90-93.
- 751 Lei, L., Shi, J., Chen, J., Zhang, M., Sun, S., Xie, S., Li, X., Zeng, B., Peng, L., Hauck,
752 A., et al. (2015). Ribosome profiling reveals dynamic translational landscape
753 in maize seedlings under drought stress. *Plant J.* 84:1206-1218.
- 754 Li, H., Peng, Z.Y., Yang, X.H., Wang, W.D., Fu, J.J., Wang, J.H., Han, Y.J., Chai, Y.C.,
755 Guo, T.T., Yang, N., et al. (2013). Genome-wide association study dissects the
756 genetic architecture of oil biosynthesis in maize kernels. *Nat. Genet.* 45:43-50.
- 757 Li, K., Wen, W.W., Alseekh, S., Yang, X.H., Guo, H., Li, W.Q., Wan, L.X., Pan, Q.C.,
758 Zhan, W., Liu, J., et al. (2019). Large-scale metabolite quantitative trait locus
759 analysis provides new insights for high-quality maize improvement. *Plant J.*
760 99:216-230.
- 761 Liu, H.J., Luo, X., Niu, L.Y., Xiao, Y.J., Chen, L., Liu, J., Wang, X.Q., Jin, M.L., Li,
762 W.Q., Zhang, Q.H., et al. (2017). Distant eQTLs and non-coding sequences
763 play critical roles in regulating gene expression and quantitative trait variation
764 in maize. *Mol. Plant* 10:414-426.
- 765 Liu, L., Du, Y.F., Shen, X.M., Li, M.F., Sun, W., Huang, J., Liu, Z.J., Tao, Y.S., Zheng,
766 Y.L., Yan, J.B., et al. (2015). *KRN4* controls quantitative variation in maize
767 kernel row number. *Plos Genet.* 11:e1005670.
- 768 Liu, W.T., Kersten, R.D., Yang, Y.L., Moore, B.S., and Dorrestein, P.C. (2011).
769 Imaging mass spectrometry and genome mining via short sequence tagging
770 identified the anti-infective agent arylomycin in streptomyces roseosporus. *J.*
771 *Am. Chem. Soc.* 133:18010-18013.
- 772 Lorenzo-Orts, L., Witthoeft, J., Deforges, J., Martinez, J., Loubery, S., Placzek, A.,
773 Poirier, Y., Hothorn, L.A., Jaillais, Y., and Hothorn, M. (2019). Concerted
774 expression of a cell cycle regulator and a metabolic enzyme from a bicistronic
775 transcript in plants. *Nat. Plants* 5:184-193.
- 776 Lv, Y., Liang, Z., Ge, M., Qi, W., Zhang, T., Lin, F., Peng, Z., and Zhao, H. (2016).
777 Genome-wide identification and functional prediction of nitrogen-responsive
778 intergenic and intronic long non-coding RNAs in maize (*Zea mays L.*). *BMC*
779 *Genomics* 17:350.

- 780 Ma, J., Ward, C.C., Jungreis, I., Slavoff, S.A., Schwaid, A.G., Neveu, J., Budnik, B.A.,
781 Kellis, M., and Saghatelian, A. (2014). Discovery of human sORF-encoded
782 polypeptides (SEPs) in cell lines and tissue. *J. Proteome Res.* 13:1757-1765.
- 783 Magny, E.G., Pueyo, J.I., Pearl, F.M.G., Cespedes, M.A., Niven, J.E., Bishop, S.A.,
784 and Couso, J.P. (2013). Conserved regulation of cardiac calcium uptake by
785 peptides encoded in small open reading frames. *Science* 341:1116-1120.
- 786 Makarewicz, C.A., and Olson, E.N. (2017). Mining for micropeptides. *Trends Cell
787 Biol.* 27:685-696.
- 788 Matsumoto, A., Pasut, A., Matsumoto, M., Yamashita, R., Fung, J., Monteleone, E.,
789 Saghatelian, A., Nakayama, K.I., Clohessy, J.G., and Pandolfi, P.P. (2017).
790 mTORC1 and muscle regeneration are regulated by the LINC00961-encoded
791 SPAR polypeptide. *Nature* 541:228-232.
- 792 Mechlin, V., Damerval, C., and Zivy, M. (2007). Total protein extraction with
793 TCA-acetone. *Methods Mol. Biol.* 355:1-8.
- 794 Mohimani, H., Gurevich, A., Shlemov, A., Mikheenko, A., Korobeynikov, A., Cao, L.,
795 Shcherbin, E., Nothias, L.F., Dorrestein, P.C., and Pevzner, P.A. (2018).
796 Dereplication of microbial metabolites through database search of mass
797 spectra. *Nat. Commun.* 9 :4035.
- 798 Mohimani, H., and Pevzner, P.A. (2016). Dereplication, sequencing and identification
799 of peptidic natural products: from genome mining to peptidogenomics to
800 spectral networks. *Nat. Prod. Rep.* 33:73-86.
- 801 Na, C.H., Barbhuiya, M.A., Kim, M.S., Verbruggen, S., Eacker, S.M., Pletnikova, O.,
802 Troncoso, J.C., Halushka, M.K., Menschaert, G., Overall, C.M., et al. (2018).
803 Discovery of noncanonical translation initiation sites through mass
804 spectrometric analysis of protein N termini. *Genome Res.* 28:25-36.
- 805 Narita, N.N., Moore, S., Horiguchi, G., Kubo, M., Demura, T., Fukuda, H., Goodrich,
806 J., and Tsukaya, H. (2004). Overexpression of a novel small peptide
807 ROTUNDIFOLIA4 decreases cell proliferation and alters leaf shape in
808 *Arabidopsis thaliana*. *Plant J.* 38:699-713.
- 809 Nelson, B.R., Makarewicz, C.A., Anderson, D.M., Winders, B.R., Troupes, C.D., Wu,

- 810 F.F., Reese, A.L., McAnally, J.R., Chen, X.W., Kavalali, E.T., et al. (2016). A
811 peptide encoded by a transcript annotated as long noncoding RNA enhances
812 SERCA activity in muscle. *Science* 351:271-275.
- 813 Nesvizhskii, A.I. (2014). Proteogenomics: concepts, applications and computational
814 strategies. *Nat. Methods* 11:1114-1125.
- 815 Nguyen, D.D., Wu, C.H., Moree, W.J., Lamsa, A., Medema, M.H., Zhao, X.L.,
816 Gavilan, R.G., Aparicio, M., Atencio, L., Jackson, C., et al. (2013). MS/MS
817 networking guided analysis of molecule and gene cluster families. *Proc. Natl.
818 Acad. Sci. USA* 110:E2611-E2620.
- 819 Parkin, M.C., Wei, H., O'Callaghan, J.P., and Kennedy, R.T. (2005).
820 Sample-dependent effects on the neuropeptidome detected in rat brain tissue
821 preparations by capillary liquid chromatography with tandem mass
822 spectrometry. *Anal. Chem.* 77:6331-6338.
- 823 Peng, J., Zhang, H., Niu, H., and Wu, R.a. (2020). Peptidomic analyses: The progress
824 in enrichment and identification of endogenous peptides. *Trends in Analyt
825 Chem* 125:115835.
- 826 Peng, Y., Xiong, D., Zhao, L., Ouyang, W., Wang, S., Sun, J., Zhang, Q., Guan, P., Xie,
827 L., Li, W., et al. (2019). Chromatin interaction maps reveal genetic regulation
828 for quantitative traits in maize. *Nat. Commun.* 10:2632.
- 829 Plaza, S., Menschaert, G., and Payre, F. (2017). In search of lost small peptides. *Annu.
830 Rev. Cell Dev. Biol.* 33:391-416.
- 831 Ransohoff, J.D., Wei, Y.N., and Khavari, P.A. (2018). The functions and unique
832 features of long intergenic non-coding RNA. *Nat. Rev. Mol. Cell Bio.*
833 19:143-157.
- 834 Rohrig, H., Schmidt, J., Miklashevichs, E., Schell, J., and John, M. (2002). Soybean
835 ENOD40 encodes two peptides that bind to sucrose synthase. *Proc. Natl. Acad.
836 Sci. USA* 99:1915-1920.
- 837 Rossi, M., Bucci, G., Rizzotto, D., Bordo, D., Marzi, M.J., Puppo, M., Flinois, A.,
838 Spadaro, D., Citi, S., Emionite, L., et al. (2019). LncRNA EPR controls
839 epithelial proliferation by coordinating *Cdkn1a* transcription and mRNA decay

- 840 response to TGF-beta. *Nat. Commun.* 10:1969.
- 841 Rubinsztein, D.C. (2006). The roles of intracellular protein-degradation pathways in
842 neurodegeneration. *Nature* 443:780-786.
- 843 Ruiz-Orera, J., Messeguer, X., Subirana, J.A., and Alba, M.M. (2014). Long
844 non-coding RNAs as a source of new peptides. *Elife* 3:e03523.
- 845 Saghatelian, A., and Couso, J.P. (2015). Discovery and characterization of
846 smORF-encoded bioactive polypeptides. *Nat. Chem. Biol.* 11:909-916.
- 847 Samandi, S., Roy, A.V., Delcourt, V., Lucier, J.F., Gagnon, J., Beaudoin, M.C.,
848 Vanderperre, B., Breton, M.A., Motard, J., Jacques, J.F., et al. (2017). Deep
849 transcriptome annotation enables the discovery and functional characterization
850 of cryptic small proteins. *Elife* 6:e27860.
- 851 Secher, A., Kelstrup, C.D., Conde-Friboes, K.W., Pyke, C., Raun, K., Wulff, B.S.,
852 and Olsen, J.V. (2016). Analytic framework for peptidomics applied to
853 large-scale neuropeptide identification. *Nat. Commun.* 7:11436.
- 854 Shiber, A., Doring, K., Friedrich, U., Klann, K., Merker, D., Zedan, M., Tippmann, F.,
855 Kramer, G., and Bukau, B. (2018). Cotranslational assembly of protein
856 complexes in eukaryotes revealed by ribosome profiling. *Nature* 561:268-272.
- 857 Silvio, S., Giorgio, S., Michele, M., Dwight, T., Xiaomu, N., Fengler, K.A., Robert,
858 M., Ananiev, E.V., Sergei, S., and Edward, B. (2007). Conserved noncoding
859 genomic sequences associated with a flowering-time quantitative trait locus in
860 maize. *Proc. Natl. Acad. Sci. USA* 104:11376-11381.
- 861 Slavoff, S.A., Mitchell, A.J., Schwaid, A.G., Cabili, M.N., Ma, J., Levin, J.Z., Karger,
862 A.D., Budnik, B.A., Rinn, J.L., and Saghatelian, A. (2013). Peptidomic
863 discovery of short open reading frame-encoded peptides in human cells. *Nat.*
864 *Chem. Biol.* 9:59-64.
- 865 Starck, S.R., Tsai, J.C., Chen, K.L., Shodiya, M., Wang, L., Yahiro, K., Martins-Green,
866 M., Shastri, N., and Walter, P. (2016). Translation from the 5' untranslated
867 region shapes the integrated stress response. *Science* 351:3867.
- 868 Studer, A., Zhao, Q., Ross-Ibarra, J., and Doebley, J. (2011). Identification of a
869 functional transposon insertion in the maize domestication gene *tb1*. *Nat.*

- 870 Genet. 43:1160-1164.
- 871 Svensson, M., Skold, K., Svenningsson, P., and Andren, P.E. (2003).
872 Peptidomics-based discovery of novel neuropeptides. J. Proteome Res.
873 2:213-219.
- 874 Tavormina, P., De Coninck, B., Nikonorova, N., De Smet, I., and Cammue, B.P.
875 (2015). The plant peptidome: An expanding repertoire of structural features
876 and biological functions. Plant Cell 27:2095-2118.
- 877 Vallebueno-Estrada, M., Rodriguez-Arevalo, I., Rougon-Cardoso, A., Gonzalez, J.M.,
878 Cook, A.G., Montiel, R., and Vielle-Calzada, J.P. (2016). The earliest maize
879 from San Marcos Tehuacan is a partial domesticate with genomic evidence of
880 inbreeding. Proc. Natl. Acad. Sci. USA 113:14151-14156.
- 881 van de Sande, K., Pawlowski, K., Czaja, I., Wieneke, U., Schell, J., Schmidt, J.,
882 Walden, R., Matvienko, M., Wellink, J., van Kammen, A., et al. (1996).
883 Modification of phytohormone response by a peptide encoded by ENOD40 of
884 legumes and a nonlegume. Science 273:370-373.
- 885 van Heesch, S., Witte, F., Schneider-Lunitz, V., Schulz, J.F., Adami, E., Faber, A.B.,
886 Kirchner, M., Maatz, H., Blachut, S., Sandmann, C.L., et al. (2019). The
887 Translational Landscape of the Human Heart. Cell 178:242-260.
- 888 Wang, M., Li, W., Fang, C., Xu, F., Liu, Y., Wang, Z., Yang, R., Zhang, M., Liu, S., Lu,
889 S., et al. (2018). Parallel selection on a dormancy gene during domestication
890 of crops from multiple families. Nat. Genet. 50:1435-1441.
- 891 Wang, S., Chen, Z., Tian, L., Ding, Y., Zhang, J., Zhou, J., Liu, P., Chen, Y., and Wu,
892 L. (2019). Comparative proteomics combined with analyses of transgenic
893 plants reveal ZmREM1.3 mediates maize resistance to southern corn rust.
894 Plant Biotechnol. J. 17:2153-2168.
- 895 Wu, H.L., Song, G., Walley, J.W., and Hsu, P.Y. (2019). The tomato translational
896 landscape revealed by transcriptome assembly and ribosome profiling. Plant
897 Physiol. 181:367-380.
- 898 Yang, M.K., Lin, X.H., Liu, X., Zhang, J., and Ge, F. (2018). Genome annotation of a
899 model diatom *Phaeodactylum tricornutum* using an integrated proteogenomic

- 900 pipeline. Mol. Plant 11:1292-1307.
- 901 Yang, N., Liu, J., Gao, Q., Gui, S.T., Chen, L., Yang, L.F., Huang, J., Deng, T.Q., Luo,
902 J.Y., He, L.J., et al. (2019). Genome assembly of a tropical maize inbred line
903 provides insights into structural variation and crop improvement. Nat. Genet.
904 51:1052-1059.
- 905 Yin, X., Jing, Y., and Xu, H. (2019). Mining for missed sORF-encoded peptides.
906 Expert Rev. Proteomics 16:257-266.
- 907 Zhu, M., Zhang, M., Xing, L., Li, W., Jiang, H., Wang, L., and Xu, M. (2017).
908 Transcriptomic Analysis of Long non-coding RNAs and coding genes
909 uncovers a complex regulatory network that is involved in maize seed
910 development. Genes (Basel) 8:274.
- 911 Zoschke, R., Chotewutmontri, P., and Barkan, A. (2017). Translation and
912 co-translational membrane engagement of plastid-encoded
913 chlorophyll-binding proteins are not influenced by chlorophyll availability in
914 maize. Front. Plant Sci. 8:385.
- 915

916 **FIGURE CAPTIONS**917 **Figure 1. Peptidogenomic Workflow for Plant NCPs Identification.**

918 (A) Peptidogenomic workflow for plant NCPs identification. Endogenous peptides
 919 from plant leaves were extracted using an optimized protocol. Heat stabilization by
 920 water bath at 95 °C combined with an acid extraction buffer containing 1% TFA and
 921 plant protease inhibitors was applied to minimize the peptide degradation during
 922 peptide extraction. Plant endogenous peptides were enriched from larger protein
 923 fragments by centrifugation through 10 kDa cutoff filters. The peptides were analyzed
 924 on a high-resolution and high-accuracy mass spectrometer. MS/MS spectra data were
 925 searched against the customized peptidogenomic database and Ensembl protein
 926 database using Mascot searching engine. The resulting peptides were used to filter out
 927 the CPs and thus obtain the NCPs. (B) Customized peptidogenomic database
 928 construction. The complete maize genomic sequence was downloaded from Ensembl
 929 Plants in FASTA format, and then translated into six-frame using EMBOSS:6.6.0
 930 package. The translation of the genomic DNA started from the first, second, and third
 931 nucleotides on each strand of each chromosome and ended when a stop codon was
 932 encountered. Triplets were translated according to the standard genetic code to assign
 933 a one letter symbol for each amino acid and a '*' symbol for a stop codon. A peptide
 934 index file containing genomic coordinates and orientations
 935 (e.g. >7:150140249-150140647+|p2) was assigned to each peptide sequence.

936 **Figure 2. Overview of the Peptidogenomic Results.**

937 (A) Venn diagram showing the number of peptides identified by Ensembl protein and
 938 customized peptidogenomic databases. The areas shown in the diagram are not
 939 proportional to the number of peptides in each group. (B) The number of CPs and
 940 NCPs identified through peptidogenomic analysis. (C) Length of CPs and NCPs.
 941 Boxes represent the second and third quartiles, whiskers represent 1.5 × interquartile
 942 range. Fisher's exact test was used for hypothesis testing, * p < 0.05. (D) The
 943 molecular weight distribution of CPs (n=844). (E) Molecular weight distribution of
 944 NCPs (n=1,993). The rug plot above the x-axis represents the frequency at each

945 exposure level.

946 **Figure 3. CPs and NCPs Distribution in Maize.**

947 (A) The genome-wide distribution of CPs (green) and NCPs (red). For each
 948 chromosome, the peptide distribution pattern includes three columns. Left: CPs (green)
 949 and NCPs (red) mapped onto chromosomes. Black circles are the centromeres.
 950 Middle: CPs (green) and NCPs (red) distribution patterns by using a window size of 6
 951 Mb and 3 Mb steps based on B73 reference genome. Right: hot region distribution of
 952 CPs (green) and NCPs (red). Hot regions were defined as more than 10 peptides in a
 953 window size of 6 Mb. (B) The normalized distribution of CPs and NCPs was shown
 954 along the chromosomal arms. The x-axis represents the normalized length of each arm
 955 with the centromere set to “0” and the telomere to “1”. The y-axis reports the number
 956 of both CPs (green) and NCPs (red). (C) Correlations between CP or NCP counts and
 957 chromosomal length (Pearson correlation: CPs, $r=0.09$, $p=0.7948$; NCPs, $r=0.77$,
 958 $p=0.0099^{**}$). (D) The histogram of the distances between two of adjacent CPs or
 959 NCPs. (E) The histograms showing the distance from each CP or NCP to the closest
 960 TSS.

961 **Figure 4. Characteristics of NCPs.**

962 (A) Number of NCPs derived from both forward and reverse strands. (B) Number of
 963 NCPs derived from different gene elements. (C) Length of NCPs derived from
 964 different gene elements. Boxes represent the second and third quartiles, whiskers
 965 represent $1.5 \times$ the interquartile ranges. Fisher’s exact test was used for hypothesis
 966 testing, * $p < 0.05$. Violin plots that combine box plot and kernel density trace to
 967 describe the distribution patterns of molecular weight (D) and isoelectric point (E).
 968 Tomato: NCPs derived from 3’UTRs (n=25); beige: NCPs derived from introns
 969 (n=139); lilac: NCPs derived from intergenic regions (n=1,708); yellow: NCPs
 970 derived from 5’UTRs (n=18); green: NCPs derived from out-of-frame exons (n=139);
 971 light blue: NCPs derived from junctions (n=14). The black bars and thin lines within
 972 the violin plots represent the interquartile ranges and the entire data ranges,
 973 respectively. White dots in the center indicate the average values. The width of the
 974 violin plot represents the density of the distribution. Fisher’s exact test was used for

975 hypothesis testing, * $p < 0.05$.

976 **Figure 5. Verification and Validation of NCPs.**

977 (A) NCP RMDAHLR mapped to the 5'UTR of a gene in chromosome 1 (left).
 978 Verification of this NCP by comparing the spectra of the peptide identified by the
 979 integrative peptidogenomic pipeline (middle) to that of synthetic peptide (right). (B)
 980 NCP ILTVNLKP mapped to the 3'UTR of a gene in chromosome 4 (left). Verification
 981 of this NCP by comparing the spectra of the peptide identified by the integrative
 982 peptidogenomic pipeline (middle) to that of synthetic peptide (right). (C) NCP
 983 QISVELPGVV mapped to the intergenic region between two genes in chromosome
 984 10 (left). Verification of this NCP by comparing the spectra of the peptide identified
 985 by the integrative peptidogenomic pipeline (middle) to that of synthetic peptide (right).
 986 (D) NCP EGTPKAVGHRQ mapped to the intron of a gene in chromosome 8 (Left).
 987 Verification of this NCP by comparing the spectra of the peptide identified by the
 988 integrative peptidogenomic pipeline (middle) to that of synthetic peptide (right). (E)
 989 Percentages of NCPs detected by peptidogenomics and ribosome profiling. (F)
 990 Percentages of NCPs derived from different gene elements detected by
 991 peptidogenomics and ribosome profiling.

992 **Figure 6. Quantitative Trait Loci (QTLs) Associated Significantly with
 993 Phenotypic Traits Linked to NCPs.**

994 (A) The enrichment of NCPs within QTLs. (B) The enrichment of NCPs located
 995 within 20 kb flanking regions of significant SNPs. (C) Diagram showing the
 996 distribution of significant SNPs associated with plant traits within NCPs, one SNP
 997 associated with kernel length, one with disease (maize rough dwarf virus, MRDV),
 998 two with oil content, and four with amino acid content. (D) An isoleucine–threonine
 999 transition caused by a SNP (chr1.s_244454699, A > G; $P < 9.42\text{e-}5$) associated with
 1000 kernel length. Significant SNPs are indicated by red dotted lines. The black arrow
 1001 indicates the NCP derived from the reverse strand. (E) The oil content associated
 1002 significant SNP (chr3.s_136872577, C > T; $P < 2.09\text{e-}07$) that leads to a proline to
 1003 leucine substitution in the NCP. The black arrow shows that the NCP was derived
 1004 from the forward strand. (F) The enrichment of NCPs within regions under positive

1005 selection during maize domestication. The x-axis shows the ratio of overlapping
 1006 between the associated SNPs and the NCPs (Obs), and that between the associated
 1007 SNPs and randomly generated regions (Random). *P*-values for upper tail test were
 1008 calculated using the “pnorm” function implemented in R (lower.tail = FALSE).

1009 **Figure 7. Identification of NCPs in Arabidopsis.**

1010 (A) Venn diagram showing the number of peptides identified by Ensembl protein and
 1011 customized peptidogenomic databases. (B) Length of CPs and NCPs in Arabidopsis.
 1012 Boxes represent the second and third quartiles, whiskers represent $1.5 \times$ the
 1013 interquartile ranges. Fisher’s exact test was used for hypothesis testing, * $p < 0.05$. (C)
 1014 Number of NCPs derived from the forward and reverse strands. (D) Number of NCPs
 1015 derived from different gene elements. (E) Length of NCPs derived from different gene
 1016 elements. Boxes represent the second and third quartiles, whiskers represent $1.5 \times$ the
 1017 interquartile ranges. Fisher’s exact test was used for hypothesis testing, * $p < 0.05$.
 1018 Violin plot combines box plot and kernel density trace to describe the distribution
 1019 patterns of molecular weight (F) and isoelectric point (G). Tomato: NCPs derived
 1020 from 3’UTRs (n=91); beige: NCPs derived from introns (n=239); lilac: NCPs derived
 1021 from intergenic regions (n=666); yellow: NCPs derived from 5’UTRs (n=63); green:
 1022 NCPs derived from out-of-frame exons (n=651); light blue: NCPs derived from
 1023 junctions (n=150). The black bars and thin lines within the violin plots represent the
 1024 interquartile ranges and the entire data ranges, respectively. White dots in the center
 1025 indicate the average values. The width of the violin plot represents the density of the
 1026 distribution. Fisher’s exact test was used for hypothesis testing, * $p < 0.05$.

1027

1028 **SUPPLEMENTAL INFORMATION**

1029 **Figure S1. Length Distribution of CPs (A) and NCPs (B) in Maize.**

1030 **Figure S2. Molecular Weight and Isoelectric Point Distribution of NCPs in**
 1031 **Maize.**

1032 (A) Molecular weight of NCPs derived from different gene elements in maize. (B)
 1033 Isoelectric point distribution of NCPs derived from different gene elements in maize.

1034 Tomato: NCPs derived from 3'UTR (n=25); beige: NCPs derived from introns
 1035 (n=139); lilac: NCPs derived from intergenic regions (n=1,708); yellow: NCPs
 1036 derived from 5'UTRs (n=18); green: NCPs derived from out-of-frame exons (n=139);
 1037 light blue: NCPs derived from junctions (n=14). The rug plot above the x-axis
 1038 represents the frequency at each exposure level.

1039 **Figure S3. Molecular Weight Distribution of CPs and NCPs in Arabidopsis.**

1040 (A) Molecular weight distribution of CPs (n=2,363). (B) Molecular weight
 1041 distribution of NCPs (n=1,860). The rug plot above the x-axis represents the
 1042 frequency at each exposure level.

1043 **Figure S4. Molecular Weight and Isoelectric Point Distribution of NCPs in**
 1044 **Arabidopsis.**

1045 (A) Molecular weight of NCPs derived from different gene elements in Arabidopsis.
 1046 (B) Isoelectric point distribution of NCPs derived from different gene elements in
 1047 Arabidopsis. Tomato: NCPs derived from 3'UTR (n=91); beige: NCPs derived from
 1048 introns (n=239); lilac: NCPs derived from intergenic regions (n=666); yellow: NCPs
 1049 derived from 5'UTRs (n=63); green: NCPs derived from out-of-frame exons (n=651);
 1050 light blue: NCPs derived from junctions (n=150). The rug plot above the x-axis
 1051 represents the frequency at each exposure level.

1052 **Figure S5. Enrichment Analysis of SNPs within NCPs in Maize.**

1053 Analysis of the ratios of the numbers of NCPs containing SNPs associated with plant
 1054 traits. As a control, we also collected similar ratios in the “random regions” by
 1055 randomly shifting the genomic sequence 100 times in the same chromosome with the
 1056 same number and same distance distribution. Statistics analysis was conducted for the
 1057 within ratio between NCPs and random sequences.

1058 **Figure S6. Enrichment Analysis of NCPs within the Domestication Selection**
 1059 **Regions in Maize.**

1060 To further explore the relationship between NCPs and domestication, we selected the
 1061 NCPs with at least 1 bp overlapped with the domestication candidate region. As the
 1062 control, we randomly shifted the genomic sequence 100 times in the same
 1063 chromosome to generate random sequences with the same number and same distance

1064 distribution.

1065

1066 **Table S1. Non-redundant Peptides Identified by Protein Database in Maize.**

1067 **Table S2. Non-redundant Peptides Identified by the Customized Peptidogenomic
1068 Database in Maize.**

1069 **Table S3. NCPs Identified in Maize.**

1070 **Table S4. CPs Identified in Maize.**

1071 **Table S5. NCPs Detected by Both Peptidogenomics and Ribosome Profiling in
1072 Maize.**

1073 **Table S6. SNPs Significantly Associated with NCPs in Maize.**

1074 **Table S7. SNPs Located within the 20 kb Flanking Regions of NCPs in Maize.**

1075 **Table S8. Colocalization of Domestication and NCPs Regions in Maize.**

1076 **Table S9. Non-redundant Peptides Identified by Protein Database in
1077 Arabidopsis.**

1078 **Table S10. Non-redundant Peptides Identified by the Customized
1079 Peptidogenomic Database in Arabidopsis.**

1080 **Table S11. NCPs Identified in Arabidopsis.**

1081 **Table S12. CPs Identified in Arabidopsis.**

1082 **Table S13. Comparisons between the NCPs Identified in Maize and Arabidopsis.**

1083

1084 **Supplemental Dataset 1. Verification of the other 111 NCPs.**

1085 The other 111 NCPs were verified by comparing the spectra of the endogenous NCPs
1086 identified by the integrative peptidogenomic pipeline to that of synthetic peptides.

