- 1 Mass-spectrometry-based draft of the Arabidopsis proteome
- Julia Mergner¹, Martin Frejno¹, Markus List², Michael Papacek³, Xia Chen⁴, Ajeet
 Chaudhary⁴, Patroklos Samaras¹, Sandra Richter⁵, Hiromasa Shikata^{6,7}, Maxim Messerer⁸,
 Daniel Lang⁸, Stefan Altmann⁹, Philipp Cyprys¹⁰, Daniel P. Zolg¹, Toby Mathieson¹¹,
 Marcus Bantscheff¹¹, Rashmi R. Hazarika^{12,13}, Tobias Schmidt¹, Corinna Dawid¹⁴, Andreas
 Dunkel¹⁴, Thomas Hofmann¹⁴, Stefanie Sprunck¹⁰, Pascal Falter-Braun^{9,15}, Frank
 Johannes^{12,13}, Klaus F. X. Mayer^{8,16}, Gerd Jürgens⁵, Mathias Wilhelm¹, Jan Baumbach²,
- 8 Erwin Grill³, Kay Schneitz⁴, Claus Schwechheimer⁶ and Bernhard Kuster^{1,13,17} *
- 9 ¹Chair of Proteomics and Bioanalytics, Technical University of Munich (TUM), Freising,
 10 Germany
- ²Chair of Experimental Bioinformatics, Technical University of Munich (TUM), Freising, Germany
- ³Chair of Botany, Technical University of Munich (TUM), Freising, Germany
- ⁴*Plant Developmental Biology, Technical University of Munich (TUM), Freising, Germany*
- ⁵Center for Plant Molecular Biology, University of Tübingen, Tübingen, Germany
- ⁶Chair of Plant Systems Biology, Technical University of Munich (TUM), Freising, Germany
- ⁷Devision of Plant Environmental Responses, National Institute for Basic Biology, Okazaki,
 Japan
- ⁸Plant Genome and Systems Biology, Helmholtz Center Munich, German Research Center for
- 19 Environmental Health, Munich-Neuherberg, Germany
- ⁹Institute of Network Biology (INET), Helmholtz Center Munich, German Research Center for
- 21 Environmental Health, Munich-Neuherberg, Germany
- ¹⁰Cell Biology and Plant Biochemistry, University of Regensburg, Regensburg, Germany
- 23 ¹¹Cellzome GmbH, Heidelberg, Germany
- ¹²Population epigenetics and epigenomics, Technical University of Munich (TUM), Freising,
 Germany
- ¹³Institute of Advanced Study (IAS), Technical University of Munich (TUM), Freising, Germany
- 27 ¹⁴Chair of Food Chemistry and Molecular Sensory Science, Technical University of Munich
- 28 (TUM), Freising, Germany
- ¹⁵Chair of Micobe-Host Interactions, Ludwigs-Maximilians-University (LMU), Munich, Germany
- ¹⁶*Plant Genome Biology, Technical University of Munich (TUM), Freising, Germany*
- ¹⁷Bavarian Biomolecular Mass Spectrometry Center (BayBioMS), TUM, Freising, Germany

32 Abstract

Plants are indispensable for life on earth and represent organisms of extreme biological diversity with unique molecular capabilities ¹. Here, we present a quantitative atlas of the transcriptomes, proteomes and phosphoproteomes of 30 tissues of the model plant *Arabidopsis* 36 thaliana. It provides initial answers to how many genes exist as proteins (>18,000), where they 37 are expressed, in which approximate quantities (>6 orders of magnitude dynamic range) and to 38 what extent they are phosphorylated (>43,000 sites). We present examples for how the data may be used, for instance, to discover proteins translated from short open reading frames, to 39 40 uncover sequence motifs involved in protein expression regulation, to identify tissue-specific protein complexes or phosphorylation-mediated signaling events to name a few. Interactive 41 access to this unique resource for the plant community is provided via ProteomicsDB and 42 ATHENA which include powerful bioinformatics tools to explore and characterize Arabidopsis 43 44 proteins, their modifications and interplay.

45 *Main*

46 The plant model organism Arabidopsis thaliana (AT) has revolutionized our understanding of plant biology and influenced many other areas of the life sciences¹. Knowledge derived from 47 48 Arabidopsis has also provided mechanistic understanding of important agronomic traits in crop species². The Arabidopsis genome was sequenced 20 years ago and hundreds of natural 49 variants have since been analyzed at the genome and epigenome level ^{3,4}. In contrast, the 50 51 Arabidopsis proteome as the main executer of most biological processes is far less comprehensively characterized. To address this gap, we used state-of-the-art mass 52 spectrometry and RNA sequencing (RNA-seq) to provide the first integrated proteomic, 53 phosphoproteomic and transcriptomic atlas of Arabidopsis. Illustrated by selected examples, we 54 show how this rich molecular resource can be used to explore the function of single proteins or 55 56 entire pathways across multiple omics levels.

57 Multi-omics atlas of Arabidopsis

We generated an expression atlas covering, on average, $17,603 \pm 1,317$ transcripts, 58 59 14,430 ± 911 proteins and 14,689 ± 2,509 phosphorylation sites (p-sites) per tissue, using a reproducible biochemical and analytical approach (Fig. 1a,b; Extended Data Fig. 1a-c; 60 Supplementary Data 1,2). In total, the protein expression data covers 18,210 of the 27,655 61 protein-coding genes (66%) annotated in Araport11⁵. This is a substantial increase compared 62 to the percentage of genes with protein level evidence reported in UniProt (27%)⁶ and more 63 than double the number of proteins identified in an earlier tissue proteome analysis ⁷ (Fig. 1c, 64 65 Extended Data Fig. 1d-f). In addition, we report tissue-resolved quantitative evidence for a total 66 of 43,903 p-sites making this study the most comprehensive single Arabidopsis 67 phosphoproteome published to date (Fig. 1c). 47% of the expressed proteome was found to be 68 phosphorylated in at least one instance, confirming earlier analyses of individual 69 phosphoproteomes or information assembled in the PhosPhAT database (Extended Data Fig. 1g,h) ^{8,9}. The actual figures *in plantae* are likely considerably higher as there is limited overlap 71 between these data sources. This can be readily explained by technical factors (e. g. insufficient 72 sequence coverage) and, more importantly, differences in plant genotypes or stimuli that 73 influence phosphorylation. The above figures are similar to those reported for mammalian 74 systems, underscoring the strong conservation of protein post-translational modifications as a 75 means to create functional diversity from a limited set of genes and proteins ¹⁰.

76 The atlas can be explored using the new web portal ATHENA (Arabidopsis Thaliana ExpressioN Atlas; athena.proteomics.wzw.tum.de) and ProteomicsDB (www.proteomicsdb.org) ¹¹. Analysis 77 78 and visualization options include co-expression analysis, the exploration of tissue-specific 79 interaction networks and a tool for predicting tandem mass spectra of peptide sequences, which can be used to validate peptide identifications reported in this study ¹² (Fig. 2a,b; Supplementary 80 Table 1). ProteomicsDB further provides spectral libraries for >328,000 unmodified and >43,000 81 82 phosphorylated peptides. These can facilitate the development of quantitative mass spectrometry-based protein assays which is particularly useful for research in Arabidopsis, 83 where comparatively few antibodies exist. 84

85 Proteomic annotation of the genome

Using Araport11, the proteomic data corroborated a substantial number of annotated open 86 reading frame (ORF) borders based on the detection of 2,776 N-terminal and 2,656 C-terminal 87 88 peptides (Extended Data Fig. 2a). As expected, N-terminal peptides often showed cleavage of the initiator methionine and N-terminal acetylation which was strongly dependent on the amino 89 acid adjacent to the initiator methionine (Extended Data Fig. 2b,c)¹³. The MS data covered, on 90 average, 44% of each protein sequence, enabling the detection of unique peptides for 14,115 91 92 protein isoforms. In 297 cases, the data distinguished also two or three splice variants (Extended Data Fig. 2d-f: Supplementary Data 3). A selection of these isoform-specific peptides 93 94 was validated using synthetic peptides (confirmation rate 80%; Extended Data Fig. 2g). To 95 investigate the potential translation of predicted short open reading frames (sORF), we searched the proteomic data against a published and an in-house compiled sORF collection 96 (ARA-PEP¹⁴; ATSO) leading to the identification of 51 distinct sORFs that were subsequently 97 98 confirmed using synthetic peptides (Extended Data Fig. 2g,h; Supplementary Data 3). These 99 results demonstrate the potential of the atlas to refine Arabidopsis gene models and protein 100 sequences.

101 *Quantitative expression landscapes*

102 The dynamic range of protein and transcript expression spanned six and four orders of 103 magnitude, respectively (Extended Data Fig. 2i,j). As described before, protein evidence was underrepresented for low abundant transcripts ¹⁵, indicating that proteome coverage was not yet 104 exhaustive. Low abundant transcripts were enriched for gene ontology (GO) terms such as cell 105 106 signaling or gene expression regulation (Supplementary Data 4). Still, the protein data cover ~50% of all annotated transcription factors and ~75% of all transcriptional regulators, kinases 107 and phosphatases, many of which are themselves post-translationally regulated (Extended Data 108 Fig. 2k) ^{16,17}. Interestingly, protein phosphorylation was detected across the entire dynamic 109 110 range of protein expression, highlighting the efficiency of the employed phosphopeptide enrichment (Extended Data Fig. 2j). 111

A breakdown of the molecular data for individual tissues or morphologically similar groups (e.g. leaf, seed or root) showed that only few transcripts or proteins are expressed in a truly specific manner (Extended Data Fig. 3a). Instead, tissue types are characterized by distinct quantitative abundance patterns. For instance, the floral organs express a similar set of genes but at levels that are characteristic for each flower compartment. Here, the MADS domain transcription factors showed the expected differential expression, thus recapitulating established models of floral organ identity ¹⁸ (Extended Data Fig. 3b,c).

119 The strong variation in guantitative gene expression between tissue groups was evident on both 120 protein and transcript level (Extended Data Fig. 3d,e). Still, we found surprisingly little overlap in the rank order or even identity of the most abundant proteins and mRNAs within any given 121 tissue (Extended Data Fig. 3e,f). Extreme expression ranges were detected for e.g. seeds, 122 123 where the storage protein CRA1 accounted for ~10% of the total protein and the ten most abundant gene products already represented more than 30% of the total amount of protein or 124 125 transcript, respectively (Extended Data Fig. 3e,g). Similarly, the RuBisCO complex, renowned as the most abundant plant protein on earth, accounted for 4-7% of total leaf protein content ¹⁹ 126 (Extended Data Fig. 3g). We note that despite the dominance of some proteins in certain 127 tissues, we show that proteomes of plant tissues can be probed very deeply using current 128 technology, countering the long-standing notion that plant proteomics is particularly difficult ²⁰. 129 As expected, photosynthetic activity was one of the main factors for expression pattern 130 variations between tissues, which can also be related to plastid localized proteins and their 131 132 abundance (Extended Data Fig. 3h,i). The most divergent expression pattern on both transcript 133 and protein level, however, was found for pollen (Extended Data Fig. 3h).

134 The comparison of expression patterns between tissues can be used to ascertain biological 135 functions of individual proteins or to disentangle closely related protein family members. A closer look at the AGCVIII kinase family (23 members in Arabidopsis)²¹ revealed high protein 136 expression of members of the D6PK subfamily in embryos as well as of AGC1.5 and AGC1.7 137 kinases in flowers and pollen (Fig. 2c). In line with these results, follow-up experiments showed 138 promotor activity for AGC1.5 and AGC1.7 but not for the closely related AGC1.6 in pollen, 139 strongly supporting the contribution of the former two in pollen tube growth ²² (Fig. 2d,e). We 140 further uncovered an important role for the D6PKs in embryo development, where a combined 141 142 knockout synergistically increased the occurrence of aberrant embryo phenotypes in loss-of-143 function mutants (Fig. 2f-i).

144 **Regulation of protein amount in tissues**

The amount of protein for a given gene is determined by a regulatory system comprising 145 146 anabolic and catabolic transcriptional, translational and post-translational processes. This has been extensively studied in unicellular organisms and, more recently, in human²³. While this 147 broad area of research cannot be covered in detail here, consistent observations were made in 148 the current study. We found positive correlation between transcript and protein levels in most 149 tissues (Pearson's R 0.28-0.7, Extended Data Fig. 4a), with the majority of high or low abundant 150 transcripts resulting in high or low abundant proteins, respectively ²⁴ (Fig. 3a). To identify further 151 molecular determinants of protein abundance, we implemented a model selection approach and 152 tested each of the selected features for the potential to explain protein level variations on a 153 tissue-specific or global level (Extended Data Fig. 4b-e; Supplementary Data 5). In addition to 154 well-known predictors like transcript levels and codon usage ²⁵, evolutionary conservation, 155 mRNA sequence motifs and the number of protein interactions were significant predictors of 156 157 protein levels (Fig. 3b). However, a large proportion (48%) of protein abundance variation still 158 remains unexplained, suggesting that many additional molecular factors with great individual but 159 lower global impact may be at work which will require careful case by case experiments to 160 resolve.

Protein versus transcript abundance plots clearly show that similarly abundant transcripts often lead to proteins with >100-fold differences in abundance. These differences in the protein-tomRNA ratio (PTR) of genes might be due to differences in translation efficiency of a given transcript, transcript stability and/or the stability of the respective protein (Fig. 3c) ²⁶. High PTRs were detected for genes involved in photosynthesis or energy metabolism (e.g. petA and rbcL) (Supplementary Data 4). This implies optimized translation and stability of those proteins (and

transcripts) that are needed in large quantities. Conversely, genes with low PTRs were enriched
 in processes such as hormone-mediated signaling pathways, which, in plants, often involve
 protein degradation ²⁷. Members of the auxin-labile AUX/IAA protein family (e. g. IAA8, IAA13)
 for example show low protein signal despite high transcript levels and are known to be under
 tight control of auxin-dependent proteasomal degradation ²⁸ (Fig. 3a).

We also observed inter-tissue variation of PTRs, especially for seed and pollen (Fig. 3d). Genes 172 173 with low PTRs in seeds mainly show differential abundance at the protein level and proteome-174 wide time-resolved measurements in germinating seeds in response to cycloheximide 175 (translation block) or MG132 (proteasome block) treatment suggest that low PTR proteins are rapidly translated from stored mRNA upon germination. High PTR proteins do not show this 176 177 behavior implying that they are stored in seed to be readily available for germination (Fig. 3e,f; Extended Data Fig. 4f,g). This is consistent with reports that seeds and pollen accumulate 178 179 storage molecules including mRNAs and proteins to allow rapid development at the onset of germination and likely explains the observed uncoupling of transcription and translation in these 180 durable plant tissues ²⁹. Fluctuating PTRs of genes between tissues indicate tissue-specific 181 regulation on either transcript or protein level. Conversely, genes with stable PTRs, are likely 182 under similar regulation in different tissues and the extensive transcriptomic resources available 183 for Arabidopsis may therefore be used to estimate relative protein abundances (Extended Data 184 Fig. 5a,b). Similarly, stable ratios between p-site abundance and the corresponding protein 185 levels, indicates that, under steady-state conditions, the majority of phosphorylation abundance 186 187 changes between tissues can be attributed to fluctuations of the underlying protein abundance rather than phosphorylation stoichiometry (Extended Data Fig. 5c,d). 188

189 **Co-expression of paralogs and complexes**

190 Based on the 'guilt by association' idea, gene co-expression analysis is widely used to ascribe common functionality between genes and can be performed on transcript or protein level. We 191 observed that co-expressed gene pairs also have high STRING ³⁰ scores, indicative of physical 192 or functional interactions (Extended Data Fig. 6a). To illustrate this point, we focused on co-193 expression analyses of gene paralogs and interacting proteins. Arabidopsis, like most plants, 194 has undergone genome duplications and retains a considerable number of highly conserved 195 paralogs in the genome ³¹. Compared to randomly selected gene pairs, duplicated genes show 196 197 a clear shift towards positive expression correlations (both transcript and protein). This indicator 198 for redundant gene function is supported by the finding that knock-out mutants of the more abundant paralog are more likely to display a measureable plant phenotype ³² (Extended Data 199

Fig. 6b-d; Supplementary Data 6). Such comparative expression analyses can be useful to prioritize knock-out combinations for phenotype analyses of functionally redundant genes.

Analyzing proteins that engage in physical interactions recorded in AtPIN ³³ for co-expression 202 across tissues suggests that ~26% of AtPIN pairs may be stable rather than transient complex 203 partners (online methods; Extended Data Fig. 6e,f). As expected, protein-based correlations 204 205 were consistently higher than transcript correlations. Because co-expression analyses tend to 206 generate a large number of interaction candidates, we sought to prioritize these using data from 207 size-exclusion chromatography-mass spectrometry experiments (SEC-MS) as an independent experimental approach ³⁴ (Extended Data Fig. 7a-c; Supplementary Data 7,8). Indeed, we found 208 209 that co-detection in SEC-MS, in addition to co-expression in the tissue atlas, substantially 210 reduced the list of candidates, illustrated here by the detection of the coatomer complex in the 211 intersection of the SEC-MS and tissue atlas data (Fig. 4a,b). The identification of known protein 212 complexes was consistently improved by using this approach, which combines information from SEC and the tissue atlas datasets (Extended Data Fig. 7d, Supplementary Data 9, see online 213 214 methods).

215 Interestingly, in seed tissues, the ζ -subunit of the coatomer complex was almost exclusively provided by the ζ -1 paralog, while this was almost absent in all other tissues (Extended Data 216 217 Fig. 7e). Subunit abundance comparisons also provided information about complex stoichiometry and similar results were obtained for stable complexes when using the SEC-MS 218 219 or tissue atlas dataset (Fig. 4c, Extended Data Fig. 7f,g). We, therefore, propose that our tissue 220 atlas resource can provide an initial approximation of the relative subunit frequency of known 221 complexes with unresolved stoichiometry. This could be particularly useful for membrane- or cell 222 wall-associated complexes, which are difficult to study biochemically.

223 The Arabidopsis phosphoproteome

224 The number of kinases (642 ± 55) and phosphatases (119 ± 6) as well as the ratios of individual 225 families were comparable across most tissues. Pollen and egg cell-like callus, however, stood 226 out implying extensive signaling activity in these tissues (Extended Data Figure 8a-c). In contrast to the wealth of available phosphoproteomic data to date, information on kinase-227 substrate relationships, critical for understanding signaling cascades or pathway activities, is still 228 sparse. Because these relationships are difficult to discover experimentally, scientists often 229 initially take a computational approach. We used the motif-X algorithm ³⁵ to identify 266 230 phosphorylation motifs that grouped into 'proline-directed', 'acidic', 'basic' or 'other' motif classes 231 ³⁶ (Extended Data Figure 8d-f; Supplementary Table 2, Supplementary Data 10). Together with 232

information about the co-occurrence of kinases and p-sites in tissues, as well as external
 information like co-localization or interaction, we anticipate that this dataset can help untangle
 kinase-substrate and kinase-p-site relationships in the future.

236 Overall, the number of p-sites per protein showed vast variations. For example, members of the LEA protein family ³⁷ are phosphorylated at almost every serine, threonine or tyrosine residue in 237 238 their sequences (Fig. 5a; Extended Data Fig. 8g,h). Phosphorylation of these unstructured 239 proteins involved in seed maturation and desiccation may be a mechanism to regulate their conformational state or phase transitions ³⁸. Conversely, other proteins showed p-site clusters in 240 regulatory domains, including the juxtamembrane domain of receptor-like kinases ³⁹ (Extended 241 Data Figure 8i), implying a role in recruiting interaction partners akin to human receptor tyrosine 242 243 kinases.

While the mere detection of a phosphorylation event does not directly imply a functional 244 consequence, it provides important starting points ⁴⁰. For example, the abscisic acid (ABA) 245 receptor RCAR10, involved in the ABA signaling cascade, was found to be phosphorylated at 246 247 four sites. We generated phosphomimic mutants (S113D, S32D) that displayed reduced ABA 248 response compared to the wild type receptor (Fig. 5b). While this might be linked to an effect on ABA binding or PP2C interaction for pS113 which is part of the interaction surface for protein 249 phosphatase 2 (PP2Cs) co-receptors, it is more likely due to altered interactions with other 250 regulatory proteins for pS32, which is localized at the receptor periphery ⁴¹ (Extended Data Fig. 251 9a,b). Interestingly, both RCAR10 phosphomimic variants altered the ABA response of 252 253 RCAR10-PP2C co-receptor combinations, suggesting changes in the interaction profile and 254 phosphorylation-based fine-tuning of ABA signaling during development and stress response 255 (Extended Data Fig. 9c,d).

As a second example for functional consequences of protein phosphorylation, we studied QKY, 256 a protein localized to plasmodesmata and involved in floral and silique morphogenesis ⁴². Since 257 one of its p-sites (S262) was detected in all assayed tissues and further p-sites were found 258 259 between the first and second C2 domains in other MCTP family members, we speculated about a possible effect of pS262 on QKY function (Extended Data Fig. 9e,f). To test this hypothesis, 260 261 we generated plant lines expressing phosphomimic and phosphomutant QKY transgenes in the strong qky-9 mutant background ⁴³. The data show that phosphomimic but not phosphomutant 262 protein constructs rescued both flower and silique mutant phenotypes, suggesting that 263 264 phosphorylation of S262 is required for normal QKY function (Fig. 5c-f, Extended Data Fig. 9g-265 j).

266 Taken together, we have generated the most comprehensive, albeit still incomplete, draft of the 267 Arabidopsis (phospho)proteome to date and highlighted some of the many uses that can be 268 envisioned for this resource. Yet, much remains to be accomplished including a more 269 systematic coverage of protein sequence variants or post-translational modifications. The tens 270 of thousands of newly discovered phosphorylation sites that await functional characterization 271 present a particular future challenge. The examples we provide demonstrate that investigating 272 different levels of omics data can lead to new insights into biological processes. We therefore anticipate that this resource along with the provided online computational tools will enable the 273 274 Arabidopsis research community to perform many further types of systems-level analyses not covered here. We also expect that the study will more broadly impact plant research as the work 275 276 demonstrates that mass spectrometry-based quantitative protein assays are a veritable way to 277 overcome the disabling lack of antibodies for plant research.

279 References

- Kramer, U. Planting molecular functions in an ecological context with Arabidopsis thaliana. *Elife* 4, doi:10.7554/eLife.06100 (2015).
- 282 2 Peng, J. *et al.* 'Green revolution' genes encode mutant gibberellin response modulators. *Nature* 283 400, 256-261, doi:10.1038/22307 (1999).
- 2843Arabidopsis Genome, I. Analysis of the genome sequence of the flowering plant Arabidopsis285thaliana. Nature 408, 796-815, doi:10.1038/35048692 (2000).
- 2864Kawakatsu, T. et al. Epigenomic Diversity in a Global Collection of Arabidopsis thaliana287Accessions. Cell 166, 492-505, doi:10.1016/j.cell.2016.06.044 (2016).
- 2885Cheng, C. Y. *et al.* Araport11: a complete reannotation of the Arabidopsis thaliana reference289genome. *Plant J* **89**, 789-804, doi:10.1111/tpj.13415 (2017).
- The UniProt, C. UniProt: the universal protein knowledgebase. *Nucleic Acids Res* 45, D158-D169, doi:10.1093/nar/gkw1099 (2017).
- 2927Baerenfaller, K. *et al.* Genome-scale proteomics reveals Arabidopsis thaliana gene models and293proteome dynamics. *Science* **320**, 938-941, doi:10.1126/science.1157956 (2008).
- van Wijk, K. J., Friso, G., Walther, D. & Schulze, W. X. Meta-Analysis of Arabidopsis thaliana
 Phospho-Proteomics Data Reveals Compartmentalization of Phosphorylation Motifs. *Plant Cell*296 26, 2367-2389, doi:10.1105/tpc.114.125815 (2014).
- 297 9 Durek, P. *et al.* PhosPhAt: the Arabidopsis thaliana phosphorylation site database. An update.
 298 *Nucleic Acids Res* 38, D828-834, doi:10.1093/nar/gkp810 (2010).
- 29910Sharma, K. *et al.* Ultradeep human phosphoproteome reveals a distinct regulatory nature of Tyr300and Ser/Thr-based signaling. *Cell Rep* **8**, 1583-1594, doi:10.1016/j.celrep.2014.07.036 (2014).
- 30111Schmidt, T. et al. ProteomicsDB. Nucleic Acids Res 46, D1271-D1281, doi:10.1093/nar/gkx1029302(2018).
- 30312Gessulat, S. *et al.* Prosit: proteome-wide prediction of peptide tandem mass spectra by deep304learning. Nat Methods 16, 509-518, doi:10.1038/s41592-019-0426-7 (2019).
- 30513Bienvenut, W. V. et al. Comparative large scale characterization of plant versus mammal306proteins reveals similar and idiosyncratic N-alpha-acetylation features. Mol Cell Proteomics 11,307M111 015131, doi:10.1074/mcp.M111.015131 (2012).
- 30814Hazarika, R. R. et al. ARA-PEPs: a repository of putative sORF-encoded peptides in Arabidopsis309thaliana. BMC Bioinformatics 18, 37, doi:10.1186/s12859-016-1458-y (2017).
- 31015Wilhelm, M. et al. Mass-spectrometry-based draft of the human proteome. Nature 509, 582-311587, doi:10.1038/nature13319 (2014).
- 31216Zheng, Y. et al. iTAK: A Program for Genome-wide Prediction and Classification of Plant313Transcription Factors, Transcriptional Regulators, and Protein Kinases. Mol Plant 9, 1667-1670,314doi:10.1016/j.molp.2016.09.014 (2016).
- 31517Yang, M. et al. A comprehensive analysis of protein phosphatases in rice and Arabidopsis. Plant316Systematics and Evolution 289, 111-126, doi:10.1007/s00606-010-0336-8 (2010).
- 31718Litt, A. & Kramer, E. M. The ABC model and the diversification of floral organ identity. Semin Cell318Dev Biol **21**, 129-137, doi:10.1016/j.semcdb.2009.11.019 (2010).
- Bar-On, Y. M. & Milo, R. The global mass and average rate of rubisco. *Proc Natl Acad Sci U S A* **116**, 4738-4743, doi:10.1073/pnas.1816654116 (2019).
- 321 20 Gupta, R. *et al.* Time to dig deep into the plant proteome: a hunt for low-abundance proteins.
 322 *Front Plant Sci* 6, 22, doi:10.3389/fpls.2015.00022 (2015).
- 32321Galvan-Ampudia, C. S. & Offringa, R. Plant evolution: AGC kinases tell the auxin tale. Trends324Plant Sci 12, 541-547, doi:10.1016/j.tplants.2007.10.004 (2007).

- 32522Zhang, Y., He, J. & McCormick, S. Two Arabidopsis AGC kinases are critical for the polarized326growth of pollen tubes. *Plant J* 58, 474-484, doi:10.1111/j.1365-313X.2009.03792.x (2009).
- Eraslan, B. *et al.* Quantification and discovery of sequence determinants of protein-per-mRNA
 amount in 29 human tissues. *Mol Syst Biol* 15, e8513, doi:10.15252/msb.20188513 (2019).
- 32924Liu, Y., Beyer, A. & Aebersold, R. On the Dependency of Cellular Protein Levels on mRNA330Abundance. Cell 165, 535-550, doi:10.1016/j.cell.2016.03.014 (2016).
- Hanson, G. & Coller, J. Codon optimality, bias and usage in translation and mRNA decay. *Nat Rev Mol Cell Biol* 19, 20-30, doi:10.1038/nrm.2017.91 (2018).
- 333 26 Schwanhausser, B. *et al.* Global quantification of mammalian gene expression control. *Nature*334 473, 337-342, doi:10.1038/nature10098 (2011).
- Santner, A. & Estelle, M. The ubiquitin-proteasome system regulates plant hormone signaling.
 Plant J 61, 1029-1040, doi:10.1111/j.1365-313X.2010.04112.x (2010).
- Luo, J., Zhou, J. J. & Zhang, J. Z. Aux/IAA Gene Family in Plants: Molecular Structure, Regulation,
 and Function. *Int J Mol Sci* **19**, doi:10.3390/ijms19010259 (2018).
- 33929Bai, B. *et al.* Seed stored mRNAs that are specifically associated to monosome are translationally340regulated during germination. *Plant Physiol*, doi:10.1104/pp.19.00644 (2019).
- 34130Szklarczyk, D. *et al.* The STRING database in 2017: quality-controlled protein-protein association342networks, made broadly accessible. *Nucleic Acids Res* **45**, D362-D368, doi:10.1093/nar/gkw937343(2017).
- 34431Wang, Y., Tan, X. & Paterson, A. H. Different patterns of gene structure divergence following345gene duplication in Arabidopsis. BMC Genomics 14, 652, doi:10.1186/1471-2164-14-652 (2013).
- 34632Lloyd, J. & Meinke, D. A comprehensive dataset of genes with a loss-of-function mutant347phenotype in Arabidopsis. *Plant Physiol* **158**, 1115-1129, doi:10.1104/pp.111.192393 (2012).
- 348 33 Brandao, M. M., Dantas, L. L. & Silva-Filho, M. C. AtPIN: Arabidopsis thaliana protein interaction
 349 network. *BMC Bioinformatics* 10, 454, doi:10.1186/1471-2105-10-454 (2009).
- 35034Kristensen, A. R., Gsponer, J. & Foster, L. J. A high-throughput approach for measuring temporal351changes in the interactome. Nat Methods 9, 907-909, doi:10.1038/nmeth.2131 (2012).
- 35235Schwartz,D & Gygi, SP. An iterative statistical approach to the identification of protein353phosphorylation motifs from large-scale data sets. Nat Biotechnology 231391-1398,354doi:10.1038/nbt1146 (2005).
- 35536Villen, J., Beausoleil, S. A., Gerber, S. A. & Gygi, S. P. Large-scale phosphorylation analysis of356mouse liver. *Proc Natl Acad Sci U S A* **104**, 1488-1493, doi:10.1073/pnas.0609836104 (2007).
- 35737Battaglia, M., Olvera-Carrillo, Y., Garciarrubio, A., Campos, F. & Covarrubias, A. A. The enigmatic358LEA proteins and other hydrophilins. *Plant Physiol* **148**, 6-24, doi:10.1104/pp.108.120725 (2008).
- 35938Bah, A. *et al.* Folding of an intrinsically disordered protein by phosphorylation as a regulatory360switch. *Nature* **519**, 106-109, doi:10.1038/nature13999 (2015).
- 36139Mitra, S. K. *et al.* An autophosphorylation site database for leucine-rich repeat receptor-like362kinases in Arabidopsis thaliana. *Plant J* 82, 1042-1060, doi:10.1111/tpj.12863 (2015).
- 40 Landry, C. R., Levy, E. D. & Michnick, S. W. Weak functional constraints on phosphoproteomes.
 364 *Trends Genet* 25, 193-197, doi:10.1016/j.tig.2009.03.003 (2009).
- Hauser, F., Li, Z., Waadt, R. & Schroeder, J. I. SnapShot: Abscisic Acid Signaling. *Cell* 171, 1708 1708 e1700, doi:10.1016/j.cell.2017.11.045 (2017).
- Vaddepalli, P. *et al.* The C2-domain protein QUIRKY and the receptor-like kinase STRUBBELIG
 localize to plasmodesmata and mediate tissue morphogenesis in Arabidopsis thaliana.
 Development 141, 4139-4148, doi:10.1242/dev.113878 (2014).
- Fulton, L. *et al.* DETORQUEO, QUIRKY, and ZERZAUST represent novel components involved in
 organ development mediated by the receptor-like kinase STRUBBELIG in Arabidopsis thaliana.
 PLoS Genet 5, e1000355, doi:10.1371/journal.pgen.1000355 (2009).

374 Figure legends

375 Figure 1 | Tissue map and multi-omics dataset. a. Schematic representation of tissue 376 samples analysed in this study and coloured according to morphology group: Flower (light grey), seed (dark brown), pollen (vellow), stem (dark green), leaf (light green), root (dark grey), fruit 377 378 (light brown), callus (magenta), cell culture (blue). b, Number of identifications at the protein, 379 transcript and phosphorylation site (p-site) level for all tissues (n = 1 measurements per tissue). 380 A dashed line indicates the number of core proteins, transcripts or p-sites detected in all tissues. Tissue-enhanced proteins or transcripts are marked in darker colour. P-sites with high 381 confidence amino acid localization (class I,>0.75 localization probability) are shown in pink, 382 383 ambiguous site localizations in purple. The number of p-sites exclusively detected in one tissue is shown as circles. c, Total number and overlap of identified gene loci in the transcriptome, 384 385 proteome and phosphoproteome datasets compared to Araport11 (left panel) and total number of identified p-sites and proportion of class I sites (right panel). 386

387 Figure 2 | Data exploration in ATHENA and ProteomicsDB. a, Schematic representation of data analysis options in the web tool ATHENA (Arabidopsis THaliana ExpressioN Atlas) and 388 389 ProteomicsDB. b, Mirror plot showing near identical observed and predicted tandem mass 390 spectra (n = 1 acquired spectra) for a peptide from the 'uncertain' protein AT4G13955, thus validating its existence as a protein. SA denotes the normalized spectral contrast angle, a 391 392 measure for spectrum similarity. c, Upper panel: summed absolute protein expression of AGC1 and AGC3 subfamilies of AGCVIII kinases. Lower panel: relative protein expression profiles of 393 D6PK family members and the closely related kinases AGC1.5, AGC1.6 and AGC1.7. Arrow 394 395 heads indicate embryo and pollen tissue with high expression of D6PKL3 or AGC1.5 and AGC1.7, respectively. d-e, Representative images of GUS-stained flowers expressing AGC1.6 396 397 (n = 20 independent transgenic lines) and AGC1.7 (n = 12 independent transgenic lines) 398 promotor-GUS fusions. Arrow-heads indicate anthers containing pollen grains. Scale bars 1 399 mm. f-i, Representative images of embryo phenotypes (at heart stage or age matched 400 equivalent) and their respective frequency for wild type (WT), d6pkl3 single, d6pk012 triple (d6pk d6pkl1 d6pkl2) and d6pk0123 (d6pk d6pkl1 d6pkl2 d6pkl3) guadruple knock-out mutants. 401 402 Scale bars 20 µm.

Figure 3 | Protein and mRNA expression. a, Scatter plot of protein versus transcript 403 404 abundance (median across 30 tissues). Arrows mark examples for genes with high protein to 405 transcript ratio (PTR, rbcL, petA) or low PTR (IAA8, IAA13). Pearson correlation (R). b, Relative 406 contribution of molecular features (averaged across tissues) able to predict protein abundance. UTR, untranslated region, Dn/Ds, ratio of non-synonymous to synonymous nucleotide 407 substitutions. c, Schematic representation of the steps where modulation of PTRs may occur. 408 For example, the higher the transcription rate the lower the PTR. d, PTR distribution of genes 409 410 for selected tissues. Genes with high and low PTRs are defined as being outside two standard deviations of the median PTR distribution. e, Analysis of the proportion of genes with high or low 411

PTRs that are at least 4-fold differently abundant on transcript and/or protein level in seed compared to all other tissues (depicted as arrows). The percentage of genes below the 4-fold criterion are placed at the origin of the plot. n denotes the number of genes with high or low PTRs in seed as defined in panel d. f, Relative protein expression changes of proteins with high or low PTRs (from panel e) in seeds incubated for 8h, 16h and 24h with 100 μM cycloheximide (translation block), MG132 (proteasome block) or DMSO (control).

418 Figure 4 | Protein complex characterization by protein co-expression and SEC-MS. a, Upper panel: 221 proteins that form a module in the tissue atlas based on their co-expression 419 420 across all tissues. Lower panel: co-elution of 57 proteins in a size-exclusion chromatography (SEC)-mass spectrometry experiment using flower tissue. b, Overlap of proteins (n=14) 421 422 between the groups of co-expressed and co-eluting proteins in a. Ten of the proteins in the 423 intersection are subunits of the coatomer complex. c, Tissue-resolved absolute abundance (top 424 panel) and relative proportions of coatomer complex subunits based on the tissue atlas protein 425 expression data (lower panel) showing heterogeneous expression but a fairly homogeneous composition of the complex across tissues. Paralogs of the α , β , β' , ϵ and ζ coatomer subunits 426 were combined for the analysis. Tissue groups are coloured as in Figure 1. 427

428 Figure 5 | Ascribing function to protein phosphorylation. a, Scatter plot of the proportion of 429 potential phosphorylation acceptor sites in a protein sequence (limited to S, T and Y) versus the 430 proportion of STY residues found to be phosphorylated in this study. While most proteins 431 including RCAR10 and QKY show only a small proportion of phosphorylated residues, the LEA proteins are phosphorylated at almost every possible residue. b, Abscisic acid (ABA) response 432 (average ± SD; n = 3) of ABA treated protoplast cells transformed with wild type or 433 phosphomimic variants of RCAR10 (S32D and S113D) suggesting that S32 and S113 434 435 phosphorylation are functionally involved in RCAR10 activity. c-f, Representative images of flowers and siliques of wild type (WT), qky-9 mutant and qky-9 expressing phospho-dead 436 437 QKY_{S262A} or phosphomimic QKY_{S262E} constructs showing that QKY_{S262E} but not QKY_{S262A} 438 rescues the *qky*-9 phenotype. Scale bars: flower 1 mm, silique 0.5 cm.

440 Methods

441 Sample preparation

Inflorescence and seed samples: Arabidopsis thaliana wild type Columbia-0 (Col-0) plants 442 were grown under continuous white light conditions at 22°C. Samples for flower parts and 443 siliques were harvested from mature plants. Stage 15 flowers ⁴⁶ were dissected into petal. 444 sepal, stamen and carpel. Fully grown green siliques were separated into septum, valves and 445 green seeds (stage 10). Surface-sterilized mature dry seeds were stored for two days at 4°C 446 and subsequently imbibed for 24 h or kept dry. For pollen collection, plants were grown on soil 447 under a long photoperiod (16 h of light, 21°C, 65% humidity). Mature pollen was bulk-harvested 448 from open flowers at developmental stage 13⁴⁶. A vacuum cleaner was modified with three 449 subsequent nylon meshes (80 micron, 35 micron, 10 micron mesh) for large scale pollen 450 isolation as described ⁴⁷. Pollen was collected in a 1.5 ml reaction tube, snap-frozen in liquid 451 452 nitrogen and stored at -80°C until further use.

453 Cell culture and callus samples: Cell culture samples (root cells, Col-0) were grown in 454 medium composed of 4.3 g/l Murashige and Skoog basal medium (MS), 30 g/l sucrose and 0.33 mg/ml KH₂PO₄, supplemented with 2.4-D (final concentration 1 mg/l; pH adjusted to 5.8 with 455 KOH) at 22°C under continuous light and harvested either three or ten days after sub-culturing. 456 457 To generate callus inducing medium (CIM) callus, root explants were harvested from two 458 weeks-old seedlings (Col-0) grown in sterile culture on MS plates. 5 to 10 mm long root segments were cultured on CIM medium composed of 1x Gamborg's B-5 salts, 20 g/l glucose, 459 460 0.5 g/I MES, 1x Gamborg's vitamin solution and 1% Phytoagar supplemented with 2,4-D (500 µg/l) and kinetin (50 µg/l). CIM calli appeared after 7 to 10 days and were propagated in two 461 462 weeks-subculture intervals on CIM medium. To generate the callus line with an egg cell-like 463 expression profile, the coding sequence of RKD2 (AT1G74480) was amplified from pistil cDNA using the primer pair RKD2fw and RKD2rev (Supplementary Data 11). Pistils from flowers at 464 developmental stage 12⁴⁷ were harvested to purify mRNA and generate cDNA as described ⁴⁸. 465 The PCR fragment was cloned into pENTR/D-TOPO (Invitrogen, Carlsbad, USA) and 466 subsequently transferred into the GATEWAY-compatible destination vector pH7FWG2.0⁴⁹ with 467 LR clonase. The resulting expression vector 35Sp:RKD2-GFP was used for floral dip 468 transformation of Arabidopsis ⁵⁰. Seeds of transformed plants were surface-sterilized and grown 469 on ½ MS medium with 2% (w/v) sucrose, 1% Phytoagar, and Hygromycin (30 µg/ml). RKD2-470 induced calli had formed after 20 to 30 days and were propagated in two weeks-subculture 471

472 intervals on $\frac{1}{2}$ MS, 2% (w/v) sucrose, 1% Phytoagar, Hygromycin (30 μ g/ml). Calli were 473 collected with a sterile scalpel blade and immediately frozen in liquid nitrogen.

Leaf and root samples: Arabidopsis plants (Col-0) were grown under continuous light 474 conditions at 22°C. Senescent rosette leaves were collected from 35 days-old plants. Samples 475 for stem, first node, second internode and first cauline leaf were collected from 30 days-old 476 plants. Rosette leaf sections (Leaf7: distal; proximal; petiole) were harvested from 22 days-old 477 plants before bolting. Seedlings were grown on 1/2 MS plates under continuous light for 478 479 seven days and separated into cotyledons, hypocotyl, root tip, root maturation zone or seed 480 apical meristem including cotyledons and first leaves. Whole roots were harvested from 22 days-old plants grown under continuous light on 1/2 MS plates. 481

482 Classification of growth stage and plant section was done as described ^{44,51-53}. Harvested 483 material from at least three individual plants was combined for each sample, frozen in liquid 484 nitrogen and stored at -80°C until further use.

485 *Protein lysis and digest*

Frozen plant material was homogenized with a tissue lyzer (Qiagen, Hilden, Germany) or with 486 mortar and pestle in liquid nitrogen. Proteins were precipitated over night with 10% 487 tricholoroacetic acid in acetone at -20° C and subsequently washed two times with ice-cold 488 acetone. Dry samples were incubated with urea digestion buffer (8 M urea, 50 mM Tris-HCl pH 489 7.5, 1 mM DTT, cOmplete[™] EDTA-free protease inhibitor cocktail (PIC) [Roche, Basel, 490 Switzerland], Phosphatase inhibitor [PI-III; in-house, composition resembling Phosphatase 491 492 inhibitor cocktail 1,2 and 3 from Sigma-Aldrich, St. Louis, USA]) for 1 h. Protein concentration was determined with a Bradford assay ⁵⁴. 1 mg of protein was reduced (10 mM DTT, 1h, room 493 494 temperature), alkylated (55 mM chloroacetamide, 30min, room temperature) and subsequently diluted 1:8 with digestion buffer (50 mM Tris-HCl pH 8.0, 1 mM CaCl₂). In-solution pre-digestion 495 with trypsin (Roche, Basel, Switzerland) was performed for 4 hours at 37°C (1:100 496 497 protease:protein ratio), followed by overnight digestion with trypsin (1:100 protease:protein ratio). Samples were acidified to pH 3 using trifluoroacetic acid (TFA) and centrifuged at 498 499 14,000 g for 15 min at 4°C. The supernatant was desalted on 50 mg SepPAC SPE catridges 500 (Waters, Milford, USA). Peptides were eluted with 0.1 % TFA in 50% acetonitrile (ACN) and vacuum-dried in a Thermo Savant SPD SpeedVac (Thermo Fisher Scientific, Waltham, USA). 501

502 **Peptide enrichment and off-line fractionation**

Fe³⁺-IMAC was performed as described previously with some adjustments ⁵⁵. Briefly, desalted 503 504 peptide samples were re-suspended in ice-cold IMAC loading buffer (0.1% TFA, 40% ACN). For 505 guality control, 1.5 nmol of a synthetic library of phosphopeptides and their corresponding nonphosphorylated counterpart sequence (B2 and F1) ⁵⁶ were spiked into each sample prior to 506 loading onto a Fe³⁺-IMAC column (Propac IMAC-10 4x50 mm, Thermo Fisher Scientific, 507 Waltham, USA). The enrichment was performed with Buffer A (0.07% TFA, 30% ACN) as wash 508 buffer and Buffer B (0.315% NH₄OH) as elution buffer. Collected full proteome and 509 phosphopeptide fractions were vacuum-dried and stored at -80°C until further use. 510

511 For the full proteome fraction, hydrophilic strong anion exchange chromatography (hSAX) peptide separation was performed as described previously ¹⁵. Briefly, an equivalent of 300 µg 512 513 protein digest was reconstituted in hSAX solvent A (5 mM Tris-HCl, pH 8.5) and separated using a Dionex Ultimate 3000 HPLC system (Dionex Cor., Idstein, Germany) equipped with an 514 515 IonPac AG24 guard column (2x50 mm) and an IonPac AS24 stong anion exchange column (2x 250 mm, Thermo Fisher Scientific, Waltham, USA). Fractions were collected in 96 well format 516 and subsequently pooled to 24. Individual fractions were acidified with formic acid (FA), 517 desalted on self-packed StageTips (five disks, Ø 1.5 mm C18 material, 3M Empore[™], elution 518 519 solvent 0.1% FA in 50% ACN) and dried down prior to LC-MS/MS analysis. Phosphopeptide fractions were separated into four fractions using a StageTip (five disks, Ø 1.5 mm C18 520 material, 3M Empore[™]) based high pH reversed phase protocol as described previously ⁵⁷. 521 522 Phophosphopeptides were eluted with 2.5%, 7.5%, 12.5% and 50% ACN in 25 mM NH₄FA pH 10. The flow through and 50% ACN fraction were combined and all fractions were dried down 523 524 prior to LC-MS/MS analysis.

For the cycloheximide and MG132 chase experiments, 300 µg peptides for each sample were 525 526 reconstituted in high pH reversed phase loading buffer (2.5 mM NH₄HCO₃, pH 8) and fractionated using a Dionex Ultimate 3000 HPLC system (Dionex Corp., Idstein, Germany) and 527 528 a Waters XBridge column (BEH130 C18, 3.5 um, 2.1x150mm; Waters, Milford, USA). Peptides 529 were separated by a linear gradient from 4% Buffer D to 32% D in 45 min, followed by a linear gradient from 32% D to 80% D in 6 min. The proportion of Buffer A was kept at 10% during 530 fractionation (Buffer A: 25 mM NH₄HCO₃, pH 8; Buffer C: 100% ACN; Buffer D: 100% ultrapure 531 532 water). Fractions were collected in 96 well format, subsequently pooled to 48, acidified and dried in a SpeedVac. 533

534 Seed treatment with cycloheximide and MG132

Aliguots of surface-sterilized mature dry seeds ((10 mg, Col-0) were stored for two days at 4°C 535 536 in the dark and subsequently imbibed with 2 ml liquid ½ MS at 22°C under constant light. After 4 537 h incubation, cycloheximide (CHX) or MG132 (N-(benzyloxycarbonyl)-Leu-Leu-Leu-al) was added to a concentration of 100 µM. Since both CHX and MG132 were dissolved at the same 538 539 stock concentration in 100% dimethylsufloxide (DMSO), the respective volume of 100% DMSO was added to the control samples. As a baseline sample, one seed aliguot was dried and 540 541 immediately frozen in liquid nitrogen after the 4 h incubation in 1/2 MS medium. Seeds for CHX, MG132 and DMSO control were removed after 8 h, 16 h and 24 h treatment, dried and frozen in 542 543 liquid nitrogen. Germination of treated seeds was visually checked after 4 days of incubation. Both compounds CHX and MG132 were active and led to complete or partial inhibition of seed 544 germination as compared to a DMSO treated control sample (Extended Data Fig. 4i). 545

Protein extraction and digest was performed as described above. Peptide quantification prior to
high pH reversed phase fractionation was done using the Pierce[™] BCA protein assay kit
(Thermo Fisher Scientific, Waltham, USA)⁵⁸.

549 Size exclusion chromatography (SEC)

Homogenized flower (stage 15), leaf (rosette leaf 7) and root (whole root) samples were lysed in 550 ice-cold 50 mM Tris-HCl pH 7.4, 100 mM NaCl, 10% glycerin, PIC, PI-III and 0.1% Triton-X100. 551 552 After filtering (Ø 0.2µm), 0.25 ml lysate containing 1 mg of protein was injected on a Superose 6 553 10/30 GL column (GE Healthcare, Chicago, USA) and separated at a flow rate of 250 µl/min on 554 an Äkta pure 25 (GE Healthcare, Chicago, USA). Molecular weight calibration of the column 555 was performed with the high molecular weight gel filtration calibration kit (GE Healthcare, 556 Chicago, USA). After the void volume, 80 fractions of 125 µl each were collected, vacuum dried 557 and re-solubilized in urea digestion buffer. In-solution digestion with trypsin and sample 558 desalting on self-packed StageTips was performed as described above. For quality control, PROCAL peptide standard ⁵⁹ was spiked into each SEC fraction prior to LC-MS/MS analysis. 559

560 LC-MS/MS analysis

Nanoflow LC-MS/MS was performed by coupling a Dionex 3000 (Thermo Fisher Scientific, Waltham, USA) to a QExactive Orbitrap HF (Thermo Fisher Scientific, Waltham, USA). Samples for the proteome and phosphoproteome analysis were re-suspended in loading buffer containing 0.1% formic acid (FA) or 50 mM citrate and 1% FA, respectively. Peptide loading and washing were done on a trap column (100 µm i.d. x 2 cm, packed in-house with Reprosil-Pur C18-GOLD, 5 µm resin, Dr. Maisch, Ammerbuch, Germany) at a flow rate of 5 µl/min in 100% 567 loading buffer (0.1% FA) for 10 min. Peptide separation was performed on an analytical column 568 (75 µm i.d. x 40 cm packed in-house with Reprosil-Pur C18, 3 µm resin, Dr. Maisch, 569 Ammerbuch, Germany) at a flow rate of 300 nl/min using a 110 min gradient from 4% to 32% solvent B (solvent A: 0.1% FA, 5% DMSO in HPLC grade water; solvent B: 0.1% FA, 5% DMSO 570 in acetonitrile) for the tissue map full proteome analysis and a two-step 110 min gradient from 571 2% to 27% solvent B for the phosphoproteome analysis ⁶⁰. Peptides were ionized using a spray 572 voltage of 2.2 kV and a capillary temperature of 275°C. The instrument was operated in data-573 dependent mode, automatically switching between MS and MS2 scans. For the full proteome 574 samples of the tissue atlas experiment, full scan MS spectra (m/z 360 - 1300) were acquired 575 with a maximum injection time of 10 ms at 60,000 resolution and an automatic gain control 576 (AGC) target value of 3e6 charges. For the top 12 precursor ions, high resolution MS2 spectra 577 were acquired in the orbitrap with a maximum injection time of 50 ms at 15,000 resolution 578 (isolation window 1.7 m/z), an AGC target value of 2e5 and normalized collision energy of 25%. 579 580 The underfill ratio was set to 1% with a dynamic exclusion of 30 s. Only precursors with charge states between 2 and 6 were selected for fragmentation. For the phosphoproteome analysis, the 581 582 MS2 spectra were acquired with a resolution of 30,000 (isolation window 1.7 m/z) and a 583 maximum injection time of 120 ms. Dynamic exclusion was set to 35 s.

584 SEC samples were measured using a 60 min 4% to 32% gradient. MS2 spectra were generated 585 for the top 20 precursors, with a maximum injection time of 50 ms at 30,000 resolution (isolation 586 window 1.7 m/z), a normalized collision energy of 25% and a dynamic exclusion of 20 s.

587 The samples for cycloheximide and MG132 treated seed samples were analyzed by a micro-588 flow LC system coupled online to an Orbitrap Fusion Lumos mass spectrometer (Thermo Fisher Scientific, Waltham, USA)) as described previously ⁶¹. The micro LC system was built by 589 590 combining a modified Vanquish pump with the auto-sampler of the Dionex UltiMate 3000 RSLCnano System. The sample was directly loaded onto a Thermo Fisher Scientific Acclaim 591 592 PepMap 100 C18 LC column (2 μ m particle size, 1 mm ID x 150 mm), the flow rate was 50 µl/min and column temperature was maintained at 55 °C. A 15 min linear gradient of 7%-32% 593 594 solvent B (solvent A: 0.1% FA, 3% DMSO in HPLC grade water; solvent B: 0.1% FA, 3% DMSO 595 in acetonitrile) was used to separate all the samples, followed by 1 min 95% solvent B washing 596 and 1 min 0.5% solvent B equilibrium time at a flow rate of 100 µl/min. The Orbitrap Fusion 597 Lumos mass spectrometer was operated with the Ion Max API Source installed with a HESI-II probe (50 µm ID). The detailed acquisition parameters are: Positive polarity: spray voltage 3.5 598 kV, funnel RF lens value at 40, capillary temperature of 325 °C, auxiliary gas heater 599

600 temperature of 300 °C. The flow rates for sheath gas, aux gas and sweep gas were set to 32, 5, 601 and 0, respectively. Full scan MS spectra (m/z 360 – 1300) were acquired in the orbitrap with a 602 maximum injection time of 10 ms at 120,000 and an AGC target of 4E5. MS2 spectra were acquired in the linear ion trap (rapid scan mode) after collision-induced dissociation (CID) 603 604 fragmentation with collision energy set at 35%, an AGC target of 1E4 and maximum IT of 10 ms. The isolation window was set to 0.4 m/z and scans were recorded with a maximum duty 605 606 cycle of 0.6 s and the option 'inject ions for all parallelizable time' enabled. Only precursors with charge states between 2 and 6 were selected for fragmentation and the lowest scan range of 607 608 MS2 was fixed at 100 m/z. The intensity threshold was set to 5E3 and the dynamic exclusion to 609 12 s.

610 **Peptide and protein identification and quantification**

Raw data files were searched with MaxQuant software (v. 1.5.8.3) using standard settings unless otherwise described ⁶². MS/MS spectra were searched against Araport11 ⁵ protein coding genes (Araport11_genes.201606.pep.fasta; download 06/2016), with trypsin as protease and up to two allowed missed cleavages. Carbamidomethylation of cysteines was set as fixed modification and oxidation of methionines and N-terminal acetylation as variable modifications. The 'match-between-runs' function was enabled for corresponding fractions within one parameter set, where applicable.

For the tissue map analysis, full proteome and phosphoproteome samples were processed together as two separate parameter groups, with phosphorylation of serine, threonine or tyrosine defined as variable modification only for the phosphoproteome group. Here we also added the spike-in phosphopeptide library sequences ⁵⁶ to the database search space.

For sORF identification two sORF databases (ATSO and ARA-PEP) were added to the search 622 space in addition to Araport11. ARA-PEP is a previously described repository of putative sORF-623 encoded peptides in Arabidopsis¹⁴. ATSO (Arabidopsis thaliana sORFs) was generated by 624 using sORFfinder ⁶³ and our RNA-seg analysis data to identify putative new sORFs in non-625 coding intergenic regions (ATSO.non_coding.pep.nr.cd-hit_aS1 aL0.3 c1.out.nr.fasta). Three 626 627 targets were used as input for sORFfinder. Sequences of intergenic and non-coding regions 628 relative to the Araport11 annotation and sequences resulting from a Trinity ⁶⁴ de novo transcriptome assembly which are non-overlapping with the Araport11 annotations. Afterwards, 629 the sORF nucleotide sequences were translated to amino acid sequences. CD-HIT ⁶⁵ was used 630 with the parameters -aS 1 - aL 0.3 - c 1 to reduce sequence redundancy. 631

632 SEC experiment raw files for flower, leaf and root were all searched together, with each SEC 633 fraction designated as individual experiment and the 'match-between-runs' function enabled.

The seed experiment raw files were searched using standard settings as described above.

Peak list files generated by Baerenfaller et al. (2008) ⁷ were downloaded from Pride 635 (PRD000044) and searched with the Mascot⁶⁶ search engine (version 2.4.1. Matrix science) 636 against the Araport11 database. The target-decoy option of Mascot was enabled and search 637 638 parameters included a precursor tolerance of 1.3 Da and a fragment ion tolerance 0.45 Da. Enzyme specificity was set to trypsin and up to two missed cleavages were allowed. The 639 Mascot ¹³C option, which accounts for the misassignment of the monoisotopic precursor peak, 640 was set to 1 and oxidation of methionine and carbamidomethylation of cysteines were included 641 as variable and fixed modification, respectively. The isobarQuant workflow was used to 642 generate protein identification files ⁶⁷. 643

644 Spectra validation

sORF peptides identified by database searching that could not be mapped to an existing 645 Araport11 gene model were synthethized at JPT Peptide Technologies (Berlin, Germany) using 646 Fmoc-based SPOT synthesis on membranes ⁵⁹ and measured on the same mass spectrometer 647 648 that was used for data acquisition (Tissue atlas full proteome method; see above for method description). Experimental and synthetic peptide spectra were extracted from the raw files and 649 used for similarity calculation without prior processing. Normalized spectral contrast angle (SA) 650 comparison between spectra of the tissue sample and synthetic peptides was performed using 651 in-house R scripts ⁶⁸ factoring in peaks which are either shared between spectra or exclusive to 652 the synthetic peptide spectra. Peaks exclusive to experimental spectra were ignored. 653 654 Candidates were selected with SA angle cutoff larger than 0.7 and BLASTed against the UniProt database for further validation (Supplementary Data 3). 655

Protein identifications for the 'uncertain' evidence category of UniProt were validated by 656 comparing the experimental spectra to in-silico predicted fragment spectra (Supplementary 657 12 Table 1) For this. а spectral library was obtained 658 from Prosit (https://www.proteomicsDB.org/prosit)¹² by uploading all sequence-charge combinations of 659 peptides that were identified for 'uncertain' proteins. To obtain the best matching spectra, the 660 661 collision energy of Prosit was calibrated using a standard quality control run from the same 662 mass spectrometer that was used for data acquisition. The resulting predicted spectral library

was visualized and compared to the experimentally acquired data analog to the sORF candidatepeptides. SA values larger than 0.7 were used as cutoff.

665 **RNA sequencing**

Total RNA was isolated with the NucleoSpin RNA Plant kit (Macherey-Nagel, Düren, Germany)
 according to the manufacturer's instructions. DNA was removed by on-column treatment with
 rDNAse (Macherey-Nagel, Düren, Germany). For recalcitrant samples (seed, silique, pollen), a
 LiCI-based protocol was adopted with minor modifications ⁶⁹. After LiCI precipitation, the RNA
 pellet was dissolved in rDNAse buffer and treated with rDNAse (Macherey-Nagel, Düren,
 Germany) at 37°C for 10 min. The final pellet was re-suspended in 35 μl DEPC-treated water.

RNA was guantified (Nanodrop[™], Thermo Fisher Scientific, Waltham, USA) and guality 672 checked with a Bioanalyzer 2100 (Agilent Technologies, Santa Clara, USA). RNA integrity 673 number (RIN) values between 6.4 and 10 were accepted for further analysis. cDNA libraries 674 were prepared using the TruSeq Stranded mRNA Sample Preparation kit (Illumina, San Diego, 675 676 USA) according to the instructions. Clusters were generated in two batches and sequenced on 677 a High throughput flow cell with the HiSeg 2500 platform (Illumina, San Diego, USA) to a depth of 36 million reads per sample. Quality assessment of raw and trimmed 75 bp paired RNA-seq 678 679 reads was performed with FastQC. Raw RNA-seq reads were trimmed to remove adapter contamination and poor quality base calls using Trimmomatic version 0.35 with parameters 680 (ILLUMINACLIP: Illumina-PE.fasta:2:30:10; LEADING:3; TRAILING:3; SLIDINGWINDOW:4:20; 681 MINLEN:36)⁷⁰. Trimmed RNA-seq reads were mapped to the Araport11⁴ transcriptome with 682 Kallisto version 0.43.1 (default parameters) ⁷¹. 683

684 Data processing

MaxQuant output tables were filtered for non-plant contaminants, reversed sequences and proteins which were only identified based on modified peptides. Protein abundance estimation was based on either intensity-based absolute quantification (iBAQ) ²⁶, or top3 quantification ⁷², depending on the analysis. MaxQuant ProteinGroups containing several gene loci were filtered out in order to retain only unambiguously identified gene loci for further analyses. In case multiple protein isoforms were identification in distinct ProteinGroups, only the isoform with the higher number of razor+unique peptides was retained.

692 mRNA quantities are displayed as transcripts-per-kilobase-million (TPM) and a cutoff of 1 TPM 693 was used as lower limit of detection across all samples. Gene ontology ⁷³ (GO)-term analysis for 694 the transcript abundance range was performed using the 1D enrichment function from Perseus

^{74,75} with Benjamini-Hochberg FDR threshold set to 0.01. In this function, a two-sided Wilcoxon Mann-Whitney test is employed to test for systematically larger or smaller transcript abundance
 within a GO category as compared to the global distribution of all values (Supplementary Data
 4).

699 For further gualitative and guantitative analyses, all transcript or protein isoform information was 700 collapsed onto gene level. Unless otherwise stated, displayed abundances for protein, transcript 701 and phosphorylation sites (p-site) were log₂ transformed. Protein, peptide and transcript 702 datasets for the tissue map and seed treatment experiments were median centered to the 703 overall median of the respective dataset. No normalization was performed for the p-site dataset, 704 since total p-site intensity variations between tissues are also due to biological sample 705 differences. Instead, the spike-in phosphopeptide library was used, to assess reproducible enrichment efficiency and MS measurement quality of phosphoproteome samples ⁵⁶. 706

Phosphorylation sites (p-sites) which were reported in the MaxQuant output table (Phospho(STY)sites.txt file) with a 'localization probability' larger than 0.75 were designated as 'high confidence localizations' or class I sites ⁷⁶. P-sites were considered exclusive if they were only detected in a single tissue. The number and identity of p-sites and phosphoproteins detected in this study were compared to datasets available through PhosPhAT4.0 and a published meta-analysis ^{8,9} (Extended Data Fig. 1h).

713 Data analysis

714 Genome annotations: Chromosome information contained in the Arabidopsis locus identifiers (AGI codes: AT [Arabidopsis thaliana]; 1, 2, 3, 4, 5, M, C [chromosome number, M for 715 mitochondrial, C for chloroplast]; G [gene]; 12300 [five-digit code for position on chromosome]), 716 717 was used to assign genes to their respective chromosomes (Extended Data Fig.1c). Araport11 718 gene identifiers (Arabidopsis gene identifiers [AGI]) were mapped to the UniProt Arabidopsis thaliana reference proteome (taxon identifier 3702; UP000006548; download 11/2018) based on 719 720 protein sequence. Swiss-Prot, TremBL as well as protein existence criteria annotations (level1: evidence on protein level; level2: evidence on transcript level; level 3: inferred from homology, 721 722 level 4: predicted, level 5: uncertain) were subsequently obtained from UniProt (Extended Data 723 Fig.1d).

N- and C-terminal peptide sequences were extracted from the MaxQuant peptides.txt file and filtered for zero missed cleavages (n = 2,776) (Extended Data Fig. 2a). N-terminal peptides were divided into groups with (n = 1,707) or without (n = 1,069) cleavage of the initiator

methionine. The frequency of the 20 genetically encoded amino acids at the position after the start codon was subsequently calculated separately for both groups and displayed as a pie chart. For N-terminal peptides with the same amino acids at the position after the start codon, the percentage of peptides with N-terminal acetylation was also calculated and displayed as bar plot. This analysis was also performed separately for the groups with or without cleavage of the initiator methionine (Extended Data Fig. 2b,c).

Gene isoforms annotated in Araport11 were considered as 'identified', if they were detected with a TPM intensity larger than our limit of detection cutoff (1 TPM; transcript level) or with an isoform-specific peptide (protein level) in at least one tissue sample (Extended Data Fig. 2f, Supplementary Data 3). A selection of isoform-specific peptides was synthesized together with the sORF peptide collection at JPT Peptide Technologies (Berlin, Germany) and the synthetic peptides were used for validation of the experimental spectra identifications as described above.

Transcription factor, transcription regulator, kinase and phosphatase family annotations and classifications are based on reports by Zheng et al. ¹⁶ and Yang et al. ¹⁷. Proportional coverage of these families within our dataset was calculated by counting how many of them could be identified on transcript, protein or phosphoprotein level, respectively (Extended Data Fig. 2k).

Tissue groups and tissue characteristics: Tissue samples were assigned to tissue groups as follows, based on their origin or common morphology: flower (sepal, petal, stamen, carpel, silique, flower); fruit (silique septum, silique valves); seed (embryo, seed, seed imbibed); pollen (pollen); leaf (cauline leaf, leaf distal, leaf proximal, leaf petiole, senescent leaf, cotyledons, shoot tip); stem (node, internode, flower pedicle, hypocotyl); root (root, root tip, root upper zone); callus (egg-cell like callus, callus); cell culture (cell culture early, cell culture late) (Fig. 1a,b).

750 To provide a measure of tissue similarity, the Pearson correlation coefficient was calculated for the gene expression of all pair-wise tissue combinations. The Pearson correlation coefficient is 751 752 a measure of the linear correlation between two variables, in this case the expression levels in 753 two tissue samples. Correlations were computed both on protein level and transcript level and 754 displayed as separate heatmaps for both omics levels (Extended Data Fig. 1a). For three 755 examples of morphologically highly similar tissue pairs, the gene expression levels were also displayed as scatter plots using protein abundance or transcript abundance measurements, 756 757 respectively (Extended Data Fig. 1b).

758 Tissue or tissue group specificity of genes was calculated based on iBAQ or TPM abundance 759 values, respectively ⁷⁷. Genes were assigned to categories in the following order: 'tissue 760 specific', 'tissue enhanced', 'group-specific', 'group-enhanced', 'shared' and 'mixed'. Genes were considered 'specific', if they were only detected in a particular tissue or tissue-group and 761 762 'enhanced' if their abundance was at least five-fold higher in a particular tissue or tissue-group as compared to the average levels in all other tissues. Genes that were detected in all 30 763 764 tissues but did not show enhanced expression in a tissue or tissue group were classified as 'shared'. All remaining genes are contained in the 'mixed' category (Extended Data Fig. 3a). 765

Genes detected in all tissues on either protein or transcript level were assigned to the core datasets (core_transcript n= 8,405; core_protein n = 7,734; core_intersection n = 5,043). This classification does not consider gene expression levels and should not be confused with the tissue-specificity classification described in the previous paragraph.

770 Hierarchical clustering analysis for tissues representing flower parts or the whole flower was 771 performed on log₂-transformed, z-scored protein intensities (iBAQ values) in Perseus using 772 Euclidean distance and average linkage (Extended Data Fig. 3b). The flower organ identity 773 model was restricted to the simplified ABC model, thus only displaying homeotic genes for the 774 A, B and C classes ⁷⁸, namely the expression of MADS-box transcription factors APETALA1 (AP1) for class A, APETALA3 (AP3) and PISTILLATA (PI) for class B and AGAMOUS (AG) for 775 class C. Class A expression specifies sepal formation, the combination of class A and B 776 specifies petal formation, the combination of class B and C specifies stamen formation and 777 778 class C expression specifies carpel formation in developing flowers.

779 For the cumulative abundance calculation of five representative tissues (flower; pollen; root tip; 780 leaf proximal; seed), proteins (iBAQ) and transcripts (TPM) were first ranked from highest to 781 lowest individual intensity (Extended Data Fig. 3e). The running total was then plotted against 782 the rank order. The names or identifiers of the five most abundant transcripts or proteins (rank 1 783 to 5) are listed in descending order for the respective tissue. Shared IDs among the top 100 784 most abundant transcripts (TPM) and proteins (iBAQ) were calculated for each individual tissue (Extended Data Fig. 3f). For the genes representing the highest abundant protein in at least one 785 tissue, the contribution to the total protein amount in an individual tissue was calculated by 786 dividing the iBAQ intensity (not log transformed) of the respective protein by the summed total 787 788 iBAQ intensity (not log transformed) in this tissue (Extended Data Fig. 3g).

Principal component analysis (PCA) of proteome and transcriptome data was performed in Perseus for the intersection of the core datasets (n = 5,043) on log₂-transformed and z-scored 791 iBAQ and TPM intensities (Extended Data Fig. 3h). Subcellular localization information was 792 downloaded from SUBA (download 11/2016; unambiguous localizations n = 3,506)⁴⁴ 793 (Supplementary Data 2) and used to calculate the intensity contribution of proteins assigned to a specific compartment for the different tissue groups. For this, the percentage of proteins with 794 the same SUBA annotation was calculated by dividing the summed protein intensity (iBAQ, not 795 log transformed) of each SUBA category by the total summed iBAQ intensity (not log 796 797 transformed) within each tissue. To compare tissue groups, the respective protein intensity proportion was then averaged for all tissues within one group. These averaged proportions were 798 799 subsequently plotted for each subcellular compartment (Extended Data Fig. 3i).

Protein/mRNA relation: The Pearson correlation value between median protein (median iBAQ 800 801 of 30 tissues) and median transcript (median TPM of 30 tissues) abundances was calculated 802 using the set of genes with abundance measurements on both protein and transcript level in 803 more than ten tissues (more than 10 pairwise complete observations; n = 14,069). The 804 scatterplot of median transcript (TPM) and protein abundance (iBAQ) was displayed together 805 with their marginal histograms (Fig. 2a). Tissue-specific Pearson correlation values between 806 protein and transcript abundance were calculated using all genes from the core proteome and core transcriptome intersection dataset (n = 5,043; Extended Data Fig. 4a). The core dataset 807 808 was used here to base the Pearson correlation value comparison between tissues on the same 809 set of genes in all tissues.

810 The data subset with more than 10 pairwise complete observations on protein and transcript level (n = 14,069) was also used for all further calculations involving protein-to-mRNA ratios 811 812 (PTR) (Fig. 3d, Extended Data Fig. 5). PTR values are calculated by building the ratio between protein abundance (iBAQ) and the corresponding transcript abundance (TPM) for each gene 813 814 and were calculated separately for each tissue. The tissue PTR values were then used to calculate the median and median absolute deviation (MAD) of PTR values across all 30 tissues 815 (Extended Data Fig. 5). The variation of PTR values across tissues indicates, whether protein 816 817 and transcript levels of a given gene are regulated in a similar (small PTR variation) or different way (high PTR variation) between individual tissues. To estimate the proportion of genes with 818 819 rather stable or variable PTR values, the PTR MAD values were distributed into five equal parts, 820 so that 20% of all genes from this analysis are contained in one part (5 MAD guantiles: Q1-Q5). 821 For each gene, the median PTR values were plotted against their MAD across tissues and the 822 MAD guantiles indicated in the plot (Extended Data Fig. 5a). GO-term analysis for the median 823 PTR distribution (median PTR across all tissues) was performed using the 1D enrichment function from Perseus ⁷⁴ with Benjamini-Hochberg FDR threshold set to 0.01 (Supplementary Data 4). In this function, a two-sided Wilcoxon-Mann-Whitney test is employed to test for systematically larger or smaller PTR values within a GO category as compared to the global distribution of all values.

828 In a similar way, the variation of phospho ratios across tissues indicates, whether p-site levels 829 are driven by protein abundance (small phospho ratio variation) or p-site stoichiometry (high 830 phospho ratio variation) changes, respectively. Phospho ratios were calculated by building the 831 ratio between a p-site abundance (intensity) and the abundance of its corresponding protein 832 (iBAQ). For this, we used the set of p-sites, with abundance measurements on both p-site and protein level in more than ten tissues (more than 10 pairwise complete observations; n =833 834 13,793). Phospho ratios were calculated separately for each tissue and subsequently combined 835 to calculate the median phospho ratio and phospho ratio MAD across all 30 tissues. To estimate 836 the proportion of p-sites which closely resemble the protein abundance profile, the phospho 837 ratio MAD values were again distributed into five equal parts (5 MAD guantiles: Q1-Q5), and the 838 median phospho ratio values plotted against their MAD across tissues each p-site. MAD quantiles were indicated in the plot (Extended Data Fig. 5c). 839

To provide further information about inter-tissue variability, MAD quantiles were also calculated and displayed for the expression levels on protein (iBAQ; n = 14,069), transcript (TPM; n = 14,069) and p-site (peptide intensity; n = 13,793) levels, respectively (Extended Data Fig. 5b,d).

843 Features for protein level prediction models: Since a substantial proportion of variation in 844 protein levels remained unexplained when using transcript level information alone (Pearson 845 correlation between protein and transcript abundance in tissues 0.28-0.79), additional molecular 846 features, that could explain protein abundance variations, were tested for their influence on 847 protein abundance in a model selection approach. Predictors selected for this analysis were: 848 mRNA levels, codon usage, non-synonymous to synonymous substitutions (Dn/Ds) ratios, 849 which are a measure of evolutionary conservation, gene/coding sequence (CDS) length, exon 850 number, gene position on the chromosome, cytosine methylation, the number of putative protein interactions and mRNA sequence motifs (kmers of size 3 to 7 nucleotides): 851

Codon usage statistics for the *A. thaliana* genome were obtained from Kazusa (www.kazusa.or.jp/codon/current/species/3702) and parsed to extract NCBI gene identifiers. These identifiers were mapped to their corresponding UniProt entries and Ensembl/TAIR10 annotations using the UniProt Retrieve/ID mapping tool. The extracted TAIR10 annotation was

856 merged with the Kazusa codon usage dataset. Codon frequencies were calculated for each 857 gene by dividing the count (x 3) of a given codon by the full length of the coding sequence.

The *Dn/Ds* substitution rates were calculated from CDS pairs between *Arabidopsis thaliana* and its closest relative *Arabidopsis lyrata*. Reciprocal best BLAST with an e-value cutoff of \leq 1E-08 was used to identify ortholog sequences. Individual CDS pairs were aligned using PRANK⁷⁹ and Gblocks⁸⁰ was applied to eliminate poorly aligned positions in an alignment with a cutoff of 8c contiguous non-conserved positions and without allowing for gap positions. The yn00 package in the program PAML⁸¹ for pairwise sequence comparison was used to estimate substitution rates, *Dn* and *Ds*, respectively.

The total number of exons and the total gene lengths were obtained from Araport11. The distance of each gene from the chromosomal centromeres were calculated to capture potential position-specific effects.

A. thaliana (Col-0) Whole Genome Bisulphite Sequencing (WGBS) data was obtained from van
 der Graaf et al. 2015 ⁸². For each gene, methylation levels were calculated for contexts CG,
 CHG and CHG (where H = adenine, cytosine, thymine) separately. Per gene methylation levels
 were defined as:

$$g_i = \frac{1}{max(j)} x \sum \frac{Nm_j}{N_j}$$

Here, max(j) is the total number of cytosines, Nm_j is the number of methylated reads and N_j the total number of reads at the jth cytosine.

Arabidopsis protein-protein interactions were downloaded from STRING ³⁰ and the number of protein interaction partners extracted for each gene.

All of these features have previously been associated with transcriptional activity and/or protein abundance levels in Arabidopsis and/or other species ⁸³⁻⁸⁸. For a detailed feature set description see Supplementary Data 5.

De novo motif identification: Motifs in the mRNA CDS, 3' or 5'UTR region were identified as previously described ²³. Briefly, protein expression levels were log₁₀ transformed and median centered. For genes with two or more transcript isoforms, the transcript isoform reported to have the largest summed iBAQ values across all tissues was defined as the major transcript isoform per gene and used to compute all sequence feature and mRNA levels. Raw RNA-seq reads were trimmed to remove adapter contaminations and poor quality base calls using Trimmomatic 885 0.39 with parameters (ILLUMINACLIP:Illumina-PE.fasta:2:30:10; LEADING:3; TRAILING:3; 886 SLIDINGWINDOW:4:20; MINLEN:36). Trimmed RNA-seq reads were mapped to the Araport11 887 transcriptome annotation with Kallisto version 0.46.0 using default parameters. To estimate the levels of mature mRNA, the number of reads mapping to exonic and intronic regions of the 888 transcripwere counted separately and then normalized by the total exonic and intronic region 889 lengths, respectively. Normalized intronic counts were subtracted from the normalized exonic 890 counts to obtain the mature mRNA counts. The resulting normalized exonic counts per sample 891 were corrected by the DESeg2⁸⁹ library size factor and log₁₀ transformed. Transcripts with 10 892 reads per 1kb were treated as transcribed. Tissue-specific protein-to-transcript ratios (PTR) 893 were computed using the normalized protein and transcript levels. GEMMA software ⁹⁰ was 894 used to identify de-novo motifs in 5'UTR, CDS and 3'UTR regions using the tissue-specific 895 PTRs as response variable. GEMMA uses a linear mixed model, in which the effect of each 896 individual kmer on the median PTR across tissues is assessed while controlling for the effect of 897 898 other kmers (random effects), region length and GC percent (fixed effects). The motif search was performed for kmers ranging from 3-7 nucleotides. Obtained p-values were adjusted for 899 900 multiple testing with Benjamini-Hochberg's false discovery rate (FDR) and jointly computed 901 across the p-values of all tissues. Gemma was run using the median PTR with FDR<0.1 and 902 covariates set to 'false'. 82 significant putative motifs were obtained based on their sequence 903 (5'UTR n = 32; 3'UTR n = 38; CDS n = 12) and sub-sequence (initial, all, end) region. 39 motifs 904 that lie in the 'initial' and 'end' sub-sequences were further selected. The presence or absence 905 of each enriched motif with respect to each gene was extracted in form of a binary matrix and 906 used for downstream multivariate feature selection analysis.

Model-based feature selection: Tissue-specific protein expression data was merged to the feature matrix. Preliminary pair-wise correlation analysis showed only weak to moderate correlation between individual features. In addition, Variance Inflation Factors (VIF) were calculated for each feature using the fmsb R package ⁹¹. The result showed low VIFs, suggesting that multicollinearity was not an issue for downstream analyses.

Two model selection approaches were implemented, stepwise regression and Lasso regression ^{92,93}. To select the most predictive features for protein abundance in each tissue, we used a forward-backward model selection approach in a multiple regression framework. The method was implemented using stepwiseAIC() function in R ⁹⁴, which compares the fit of nested models. To ensure that the comparison of model AICs were not affected by unequal sample size, missing data were removed prior to the analysis.

Because stepwise regression can occasionally lead to over-fitting ⁹⁵, a Lasso regression was 918 919 implemented as an alternative model selection procedure. Lasso regression performs L1 920 regularization, which adds a penalty equivalent to the absolute of the magnitude of regression coefficients and tries to minimize them. The strength of the penalty was controlled via the tuning 921 parameter λ^{93} . Lasso was implemented using the glmnet package ⁹⁶ in R. Model training was 922 performed on a random selection of genes (50% of the dataset) and implemented using the 923 924 cv.glmnet() function. The optimal value for λ was extracted and used to re-build the model using the glmnet() function. Finally, the fitted model was used to obtain predictions in the remaining 925 926 50% of the genes.

927 For each tissue, stepwise and Lasso regression models were compared. Stepwise regression 928 yielded only slightly higher coefficients of determination (R^2) values compared to Lasso, suggesting over-fitting was not an issue (Supplementary Data 5). As stepwise regression 929 930 yielded more parsimonious models in general, we used this approach for further analysis. Selected features here could explain on average 52% of protein abundance variation (37% in 931 pollen, 62% in cell culture early) (Supplementary Data 5). For each tissue, features from the 932 best fitting models were summarized in an incidence matrix along with the effect direction 933 (positive or negative effect on protein levels). To determine the importance of each feature to 934 the overall model fits, R^2 variance decomposition was performed using the 'genici' metric which 935 is implemented in the relaimpo R package ⁹⁷ (Extended Data Fig. 4c). Relative feature 936 contributions were averaged across all tissues (Fig. 3b). 937

Since many of the detected motifs appear in a tissue-specific manner (predictive only in tissue 938 subset), we clustered tissues according to the presence or absence status of 5'UTR motifs in 939 each tissue model (Extended Data Fig. 4d). Seed and pollen tissues form a distinct subcluster, 940 941 which might be connected to the increased/decreased PTR levels in these tissues compared to the other tissues (Fig. 3c;Extended Data Fig. 4a). In contrast to the motif analysis, Dn/Ds ratios 942 were consistently (and positively) associated with protein abundance in all 30 tissues 943 944 (Supplementary Data 5). To explore the relationship between Dn/Ds ratios and protein abundance in more detail, groups of genes with low or high (bottom 5% or top 5% of Dn/Ds 945 distribution) were compared to each other with regard to protein levels (Extended Data Fig. 4e). 946

Seed tissue PTR regulation: The PTR value distribution was plotted for the median PTR (median across all tissues, see above) and selected tissues (seed, seed imbibed, pollen, cell culture young). Seed and pollen tissues show particularly low Pearson correlation between protein and transcript abundance levels (Extended Data Fig. 4a). Cell culture early in

951 comparison shows the highest Pearson correlation among the analyzed tissues. Tissue-specific 952 outliers' with especially high or low PTR as compared to all other tissues were selected if their 953 PTR value was outside two standard deviations of the mean of the median PTR distribution. For seed, this resulted in 469 genes with 'high' and 571 genes with 'low' PTR. In order to interpret, 954 whether this change in seed PTR values was caused by differential regulation on transcript 955 and/or protein level, the fold change between seed protein and transcript abundance and the 956 957 median protein and transcript abundance of all tissues was calculated. Genes with at least 4fold-change in either protein or transcript abundance were considered regulated (less/more 958 959 transcript and/or less/more protein in seed compared to all tissues). The percentage of genes with high or low PTR levels in the seed tissue, which are regulated in the same manner 960 (less/more transcript and/or protein) was calculated and displayed as arrow plot (Fig. 3e). About 961 a guarter of each gene set (high PTR 27%; low PTR 24%) show no regulation based on our 962 4fold-change cutoff. 963

964 For the low PTR gene set, a high proportion of genes with lower protein abundance in the seed tissue was observed. A protein level chase experiment was performed, to investigate whether 965 966 this reduction was due to translational inhibition or increased protein degradation in seeds. Protein level (iBAQ, not log transformed) fold changes in CHX, MG132 and DMSO-treated 967 samples were calculated relative to the baseline protein expression at 0h (control sample). The 968 average protein level fold change for the subset of genes with either high or low PTR in seeds 969 970 (see above) was subsequently determined for each time point and plotted separately for 'high' 971 and 'low' PTR genes (Fig. 3f). To test for significant fold change differences between treatment 972 time points, a one-way analysis of variance (ANOVA) and post-hoc Tukey HSD test was performed using the stats package in R (Extended Data Fig. 4i). For this, the data was 973 974 normalized for changes in the proteome caused by seed germination, by calculating the \log_2 975 fold change in protein abundance (iBAQ) between seed samples treated with cycloheximide 976 (CHX) or MG132 and the time-point matched DMSO control. Outlier values were removed for 977 the boxplot visualization but are included in the ANOVA and post-hoc Tukey analysis.

Co-expression analysis: To compare gene co-expression information in our dataset to prior knowledge about various associations between genes pairs, the Pearson correlations between all pair-wise gene combinations on both transcript and protein level (core dataset intersection; n = 5,043) was calculated. The correlation coefficient value for a given gene pair on transcript level was then plotted against the correlation coefficient value of the same pair on protein level (Extended Data Fig. 6a). Protein-protein association data was downloaded from STRING ³⁰

984 (www.string-db.org; 3702.protein.links.detailed.v10.5.txt.gz; download 07/2018). the 985 experimental STRING subscores calculated for each gene pair and the information merged to 986 the correlation value matrix. The scatterplot of transcript level versus protein level correlations was then divided into 50 x 50 bins along the x- and y-axis and the mean experimental STRING 987 988 subscore for all gene pairs within a bin was calculated. Mean subscores were log₂ transformed and visualized as a heatmap together with their marginal histograms (Extended Data Fig. 6a). 989

990 For a subset of genes, pair-wise co-expression was further investigated. Pairs used in this analysis were genes, that are either annotated as duplicated ³¹ or as protein interactors (AtPIN 991 992 ³³, download 08/2017), respectively. The Pearson correlation between the expression levels of 993 gene-pairs across tissues was calculated for gene pairs with abundance measurements in more 994 than ten matching tissues on both protein and transcript level (more than 10 pairwise complete 995 observations). The duplicated gene set was further filtered for 3 unique peptide identification of 996 each paralog. Top3 intensities were used as protein abundance measure for the duplicate gene 997 set since the subsequent ratio comparison was also performed with top3 intensities. Paralog 998 genes often have high sequence identity on amino acid level, which can lead to distorted ratios 999 due to uneven assignment of shared (razor) peptides to only one of the paralogs. Duplicated 1000 genes (n = 3.612) were divided into three subsets: duplications caused by a whole genome 1001 duplication event (WGD, n = 2,104), local duplication (local, n = 408) or transposon-mediated duplication (transposed, n = 1,100). For each subset the distribution of Pearson correlation 1002 1003 values between duplicates on either transcript or protein levels was plotted (Extended Data Fig. 1004 6b). A set of random gene pairs (n = 27,353) was generated for comparison, filtered like the duplicated gene set (number of pairwise complete observation, 3 unique peptides) and the 1005 1006 Pearson correlation value distribution displayed in the same way. To compare relative protein 1007 expression levels between paralogs, the ratio between the protein abundance (top3) of one 1008 paralog and the protein abundance (top3) of the other paralog was calculated (Extended Data 1009 Fig. 6c). As an example for the intensity proportion of paralog genes in different tissues, the 1010 top3 intensity (not log transformed) of MAC5A and MAC5B was plotted as proportion of the summed top3 intensity (not log transformed) of both paralogs (Extended Data Fig. 6c). For 57 1011 1012 paralog gene pairs with phenotypic information about asymmetric loss-of-function mutant combinations ³², we build the average top3 intensity proportion of the first paralog 1013 1014 (dupl1/(dupl1+dupl2)) across 30 tissues and plotted them in descending order together with the 1015 phenotype information (Extended Data Fig. 6d).

A set of well-studied, stable protein complexes (26S proteasome, COP9 signalosome, 1016 1017 Chaperonin, Cellulose synthase) was selected to estimate the protein level (iBAQ) co-1018 expression expected for stable interaction partners. For this, the Pearson correlation coefficient 1019 values between all pair-wise subunit combinations of an individual complex were calculated and 1020 displayed in a density plot (Extended Data Fig. 6e). Based on this analysis, a Pearson correlation coefficient cutoff of >0.5 was subsequently used as indication for stable protein 1021 1022 interactions. The protein interactor gene set (AtPIN) was used to test the recovery of annotated protein-protein interactions based on co-expression data from this study. Again, the distribution 1023 1024 of Pearson correlation values between gene pairs (here interactors) on either transcript or 1025 protein level (iBAQ) was plotted. This was done for the whole AtPIN dataset (n = 57,152) and the following subsets: interactions identified by yeast-two hybrid experiments (Y2H, n = 7,621), 1026 by affinity-purification mass spectrometry (AP-MS, n = 17,982) or by both Y2H and AP-MS (n =1027 829). 1028

1029 **Size exclusion chromatography and complex analysis:** Reproducibility between the three 1030 SEC experiments was tested by comparing the elution behavior of the same proteins. The 1031 coefficient of determination between protein peaks in the different experiments was above 0.7 1032 (flower-root $R^2 = 0.82$; flower-leaf $R^2 = 0.78$; root-leaf $R^2 = 0.7$). Protein peak elution between 1033 the samples was also consistent across the whole SEC gradient range (Extended Data Fig. 7b).

1034 To assign proteins to potential complexes, peak correlation profiling was performed. For this, 1035 the R package CCprofiler v0.1 was used to analyze the peptides.txt output table from MaxQuant ^{62,98}. Peptide table entries were restricted to the maximum molecular weight per gene name for 1036 1037 the guantification of protein groups. This step removed peptides mapping to smaller isoforms of 1038 each gene. The restricted proteinGroups.txt served as the trace annotation input table for 1039 CCprofiler. The calibration table of the SEC column was based on the high molecular weight gel filtration calibration kit (GE Healthcare, Chicago, USA). The peptides table was split into the 1040 1041 different tissues and the resulting tables converted into the input format for CCprofiler. All 1042 subsequent steps were performed separately for each tissue table.

The peptides table was imported, converted into a traces object, where a trace object is the SEC-elution profile of a specific peptide and annotated with the corresponding fractions and additional external information (e.g. molecular weight of the respective fraction, molecular weight of the respective protein) based on the aforementioned calibration and restricted proteinGroups tables. Traces not containing at least three consecutive peptide identifications across fractions, as well as traces with a sibling peptide correlation of less than 0.2 between

peptides originating from the same protein were filtered out. Based on the remaining traces, protein features (elution peak positional information) consisting of highly correlated peptides (correlation higher than 0.9) were detected using CCprofiler's sliding window algorithm. Peptides were randomly assigned to proteins in order to control the false discovery rate of protein features at 5% (random decoy model). Subsequently, protein traces were calculated by summing up the intensities of the top two most abundant peptide traces.

- 1055 Next, complex hypotheses consisting of target and decoy protein complexes were constructed 1056 based on a previously-published mapping of Arabidopsis orthologues to mammalian protein 1057 complexes from CORUM ^{99,100}. Proteins were only mapped to the same decoy complex if they 1058 were not interacting with each other directly (minimum path length of two). These target and 1059 decoy complexes were used during the detection of complex features consisting of highly 1060 correlated proteins using the same sliding window algorithm as above in order to control the 1061 false discovery rate of complex features at 5% (target-decoy model).
- 1062 Both protein features (elution peak positional information) and protein traces (quantitative 1063 information) were then used to generate an input matrix for weighted gene correlation network analysis (WGCNA)¹⁰¹, which was used to find clusters of proteins with correlated SEC elution 1064 profiles. Protein traces were restricted to the ones, which were part of at least one protein 1065 feature. Each protein trace was duplicated for each protein feature it was part of and all 1066 1067 intensities outside of the respective protein feature were set to zero. With this, protein traces 1068 were effectively split into separate elution peaks corresponding to distinct (sub)complexes. Afterwards, elution peaks, which met either of the following criteria were filtered out: (1) The 1069 1070 absolute difference between the molecular weight at the apex of the elution peak and the 1071 monomer molecular weight of the corresponding protein was smaller than twice the monomer 1072 molecular weight. (2) The absolute difference between the molecular weight at the apex of the 1073 elution peak and the monomer molecular weight of the corresponding protein was less than 200 1074 kDa. (3) The apex of the elution peak was in fraction 77 (void fraction). (4) The apex of the 1075 elution peak was in fraction 4 (monomer range).
- The remaining elution peaks were restricted to the intersection with the core protein dataset across tissues described above (see section on tissue groups), generating the SEC WGCNA input dataset. At the same time, an additional WGCNA input dataset was generated based on the core protein dataset across tissues and restricted to the intersection with the SEC WGCNA input dataset. Proteins with several SEC elution peaks were duplicated to match the SEC WGCNA input dataset. WGCNA was carried out separately for each of these datasets. Signed

1082 co-expression similarities were calculated between all pairs of proteins with at least seven1083 pairwise-complete observations using the following formula:

$$s_{ij} = \frac{1 + cor(x_i, x_j)}{2}$$

Here, s_{ij} is the signed co-expression similarity between two proteins x_i and x_j (based on their Pearson correlation across tissues). Adjacency matrices were calculated with $A = [s_{ij}]$ for several values of $\beta \in [1,30]$.

$$a_{ij} = s_{ij}^{\beta}$$

Here, a_{ij} is the adjacency between two proteins x_i and x_j . The adjacency function parameter β was selected to be 30 for the construction of signed protein co-expression networks. The topological overlap matrix $\Omega = [\omega_{ij}]$ of the two networks was calculated with:

$$\omega_{ij} = \frac{l_{ij} + a_{ij}}{\min\{k_i, k_j\} + 1 - a_{ij}}$$

Here, $l_{ij} = \sum_{u} a_{iu} a_{uj}$ and $k_i = \sum_{u} a_{iu}$. Hierarchical clustering was performed on the topological-1090 overlap-based dissimilarity matrices $D = [1 - \omega_{ij}]$ using average linkage. Correlation clusters 1091 ('Modules') were detected using adaptive branch pruning ¹⁰² using the 'Dynamic Hybrid' method 1092 1093 set to respect the dendrogram topology during PAM operations and requiring at least 30 members per module. Similar clusters in the SEC or tissue dataset were merged, if the 1094 dissimilarity of their module eigengenes was smaller than 1-MPC or 1-MTC, respectively. MPC 1095 1096 is the median peak correlation of the FDR-filtered complex features from CCprofiler (known 1097 complexes), while MTC is the median tissue correlation of the very same complex features 1098 (pairwise correlation between proteins mapping to a specific complex feature). The resulting 1099 modules for both datasets were then mapped to each other using a combination of Fisher's Exact test and manual curation. Enrichment of functional annotations from Corum ¹⁰⁰, GO ⁷³, 1100 KEGG¹⁰³ and Reactome¹⁰⁴ were calculated in each cluster using Fisher's Exact Test. All p-1101 values were adjusted for multiple-testing by calculating FDRs ¹⁰⁵ (Supplementary Data 8). 1102

1103 In order to quantify how well complexes can be detected using the tissue atlas (TA) or SEC 1104 WGCNA output, we calculated a summary statistic termed 'complex index'. Protein interactor 1105 information was downloaded from UniProt ⁶ (https://uniprot.org; 12/2018) and manually curated 1106 for the presence of large (> 4 subunits) and small (\leq 4 subunits) complexes (Extended Data Fig. 1107 7e; Supplementary Data 9). The complex index (C) was calculated using the following formula:

$$C = \frac{Sm}{S} \times \frac{Sm}{M}$$

Here, S_m is the number of subunits detected as present in the same module. *S* is the total number of subunits identified in the restricted datasets for either SEC or TA experiments and *M* is the total number of entries in the respective module. C is a measure of how well complex subunits co-occur (either by co-elution or by co-expression) and thus are detected in the same module and also quantifies the resolution of the detection method (module size). The complex index is 1, when all subunits of a complex are identified in the same module and no other proteins are contained in the module.

1115 Complex stoichiometry analysis: Absolute SEC protein elution traces of selected complexes 1116 (chaperonin, 26S Proteasome, CSN, CESA) were plotted using top3 intensities (not log 1117 transformed). Relative subunit proportions or stoichiometry for selected protein complexes within individual tissues of the TA dataset were calculated using top3 intensities (based on the 3 1118 1119 most abundant unique peptides). Top3 intensity rather than iBAQ intensity was used, in order to 1120 avoid distorted ratios caused by shared (razor) peptide assignment. In the case of paralog genes with high protein sequence similarity, peptides were not required to be unique for one 1121 gene, but rather for one 'paralog group'. The three most abundant peptides within a paralog 1122 group were then used for the top3 calculations. 1123

1124 For the example of the coatomer complex, the relative proportion of paralog genes in different 1125 tissues was calculated by plotting the top3 intensity (not log transformed, unique for gene) of each paralog as proportion of the summed top3 intensity (not log transformed) of all paralogs 1126 1127 (Extended Data Fig. 7d). For the calculation of subunit stoichiometry ratios of selected 1128 complexes (chaperonin, 26S Proteasome, CSN, CESA), the top3 intensity (not log transformed) of a subunit was divided by the average top3 intensity (not log transformed) of all complex 1129 subunits. Subunit ratios were first calculated for each tissue and subsequently averaged across 1130 all 30 tissues (chaperonin, 26S Proteasome, CSN) or only across tissues where the respective 1131 1132 subunits are mainly expressed (CESA4/7/8: node, internode, silique septum, silique valves; 1133 CESA1/3/6: all other) (Extended Data Fig. 7g).

Subnetwork extraction: *De novo* network enrichment was performed to identify tissue- or tissue group-specific subnetworks and to link them to a molecular function. In contrast to gene set overrepresentation or gene set enrichment analysis ¹⁰⁶, this approach is suited to identify previously uncharacterized functions from large-scale molecular interaction networks. KeyPathwayMiner ¹⁰⁷ was used to extract tissue- and tissue group-specific subnetworks from

transcriptomics and proteomics datasets. To enrich subnetworks with proteins that deviate in their expression from other tissues, we employed z-scored expression values. These reflect how many standard deviations the expression in a given tissue is away from the expression range found in all other tissues:

$$z_{x,g} = \frac{x - \mu_{X \setminus x}}{\sigma_{X \setminus x}}$$

Here x is a specific tissue or tissue group of interest, g a gene μ the mean and σ the standard deviation. Based on the following rule, a one-column indicator matrix I(x,g) was constructed as input for KeyPathwayMiner for each tissue or tissue group x:

$$I(x,g) = \{ \begin{smallmatrix} 1 & if & |z_{x,g}| > 2 \\ 0 & else \end{smallmatrix}$$

Here each row corresponds to a gene g. STRING {Szklarczyk, 2017 #513} *A. thaliana* network (v. 10.5, download 10/2018) was used as the molecular interaction network, considering only high confidence interactions with a score > 900. Subnetworks extracted via KeyPathwayMiner were made available for individual tissues or tissue groups as part of the ATHENA web application.

1151 **Phosphorylation sites and motif analysis:** Serine (S), threonin (T) and tyrosine (Y) content 1152 and p-site number of individual proteins was calculated based on the longest sequence in case 1153 of ambiguous isoform identifications (Fig. 5a, Extended Data Fig. 8g). Likewise, the schematic 1154 depiction of phosphorylation (p-site) localization in LATE EMBRYOGENESI ABUNDANT (LEA) 1155 and receptor-like kinase proteins was based on the p-site localization in the longest isoform 1156 sequence (Extended Data Fig. 8h,i). Functional domain assignment for receptor-like protein 1157 kinase (RLKs) sequences done using smart (http://smart.emblwas heidelberg.de/smart/batch.pl)¹⁰⁶ and p-sites were manually assigned to specific domains based 1158 on their localization in the protein (Extended Data Fig. 8i). 1159

For motif analysis within our dataset, sequence windows of 15 amino acids centered on the identified class I p-sites (15 mers) were assigned to a motif class of either 'proline-directed', 'acidic', 'basic' or 'other' by following a binary decision tree ³⁶ (Extended Data Fig. 8d). Assignment to these categories was done sequentially as follows: proline amino acid at position -1 (Proline-directed), 5 or more aspartic acid (D) or glutamic acid (E) at position +1 to +7 (acidic), arginine (R) or lysine (K) at position -3 (basic), D or E at position +1,+2 or +3 (acidic), 2 or more R or K at position -6 to -1 (basic) and otherwise (other).

P-sites with high confident localization scores (class I) were divided into S, T and Y p-site 1167 1168 datasets. Motif extraction was performed separately for each motif class category (proline-1169 directed, acidic, basic, other) using the motif-X algorithm, implemented in the rmotifx R package 1170 ³⁵ (Supplementary Data 10). Cutoff settings were: min-seg.= 30, pval.cutoff = 1e-6 (S and T datasets); min-seq.= 10, pval.cutoff = 1e-5 (Y dataset). Motif-X calculates the fold-increase of a 1171 respective motif in the 'foreground' (p-site dataset) compared to the 'background' dataset (non-1172 1173 phosphorylated peptides). To avoid mass spectrometry-based residue biases, all peptides from our dataset that contained an STY amino acid and had only been identified in non-1174 1175 phosphorylated form were used as the background for all motif-X analyses (258,395 sequences). These peptide sequences were centered on STY and extended where necessary 1176 along the N- and C-terminal window to generate 15 mers. Position weight matrix sequence 1177 logos were drawn using the 'bits' method for amino acid sequences in the ggseglogo R package 1178 ¹⁰⁹ (Extended Data Fig. 8e, Supplementary Table 2). The motif-X fold increase for S motifs was 1179 plotted for motifs with 2, 3 and 4 fixed amino acid positions (Extended Data Fig. 8f). 1180

Published p-site motifs were retrieved from PhosPhAT4.0⁹ and the Human Protein Reference Database ¹¹⁰ (HPRD; www.hprd.org) and divided into groups depending on the number of fixed amino acid positions within the motif sequence. These motifs were subsequently matched to the p-site or background dataset in order to calculate detection frequency and fold increase (psite/background dataset) of a reported motif (Supplementary Data 10).

1186 AGCVIII protein expression and mutant phenotypic characterization

1187 For the relative expression analysis of AGC kinases across tissues, the AGC1 and AGC3 subfamilies of Arabidopsis AGCVIII kinases were selected ²¹. The summed total intensity of 1188 1189 AGC1 and AGC3 subfamily kinases (top3, not log transformed) was calculated for each tissue. 1190 Relative protein amounts of individual kinases in different tissues were calculated as the top3 intensity (not log transformed, unique for gene) of each kinase in proportion to the summed top3 1191 1192 intensity (not log transformed) of all AGC1 and AGC3 kinases (Fig. 2c). For clear visualization only the D6PK family (D6PK, D6PKL1, D6PKL2, D6PKL3), AGC1.5, AGC1.6 and AGC1.7 were 1193 1194 colored in the plot.

1195 Mutant alleles of the *D6PK* family kinase genes, *d6pk-1*, *d6pkl1-1*, *d6pkl2-2* and *d6pkl3-2*, and 1196 their combinations were previously described ¹¹¹. Embryos were prepared omitting fixation with 1197 ethanol/acetic acid as previously described ¹¹². Embryos were analyzed with a Zeiss Axio 1198 Imager.M2, Axiocam 512 camera and 20x/0,8 Plan Apochromat objective using differential 1199 interference contrast (DIC).

For promoter: GUS constructs of AGC1.5, AGC1.6 or AGC1.7, 2916 bp, 1207 bp, or 2214 bp 1200 1201 fragments upstream from the respective start codon were cloned into the Sall and EcoRI sites of pCAMBIA1391Z. Flowers from transgenic plants, obtained by the floral-dip method ⁵⁰, were first 1202 1203 incubated in 90% cold acetone for 15 min followed by a β -glucuronidase (GUS)-staining solution 1204 (50 mM sodium phosphate pH 7.0, 10 mM EDTA, 2 mM potassium ferricyanide, 2 mM potassium ferrocyanide, 0.1% Triton-X100, 0.5 mg/ml 5-bromo-4-chloro-3-indolyl-β-D-1205 1206 glucuronide) over-night in the dark. Following a wash in 70% ethanol, flowers were incubated in ethanol/acetate (6:1) for removing of chlorophyll and then rehydrated in a graded ethanol series 1207 1208 (90%, 70%, 50%, 30%). Samples were mounted in chloral hydrate/glycerol/water (8 g:1 ml: 2 ml) for imaging. The experiment was performed two times and GUS signal was observed for 1209 12/18 AGC1.5p::GUS, 0/20 AGC1.6p::GUS and 12/18 AGC1.7p::GUS independent lines. 1210

1211 ABA response of RCAR phosphomutants

Preparation and analysis of Arabidopsis protoplasts was performed as described ¹¹³. Briefly, 1212 protoplasts from three weeks-old Col-0 plants (10⁵ protoplasts: 0.1 ml) were transfected with 5 1213 1214 µg of reporter construct (pRD29B::LUC), 3 µg of p35S::GUS plasmid as internal control and 3 µg of effector plasmid. The effector plasmids drive expression of respective RCAR and PP2C 1215 cDNAs under control of the 35S promoter ¹¹⁴. RCAR10 phosphomimic reporter constructs 1216 (RCAR10_{S32D}; RCAR10_{S113D}) were obtained by site-directed mutagenesis using primer pairs 1217 S32D F/S32D R and S113D F/S113D R, respectively (Supplementary Data 11). The 1218 protoplast suspension was incubated at various levels of ABA as indicated and reporter 1219 expression was determined after 18 h of incubation at 25°C. Three biological replicates per data 1220 point were performed for each assay. All data points were normalized to the empty vector 1221 control. Structural modeling for RCAR10 was performed with SWISS-MODEL¹¹⁵ using RCAR11 1222 as template model (PDB ID 3k3k ^{45,116}) and modified with UCSF CHIMERA ¹¹⁷. 1223

1224 Phenotypic characterization of QKY phosphomutants

1225 The pCAMBIA2300-based pQKY::mCherry:QKY construct was described previously ⁴². The 1226 pQKY::mCherry:QKY_{S262A} and pQKY::mCherry:QKY_{S262F} plasmids were obtained using the Q5 1227 site-directed mutagenesis kit (NEB, #E0554S) according to the manufacturers recommendation. 1228 Primers are Q5SDM S262A F and Q5SDM S262A R for pQKY::mCherry:QKY_{S262A} and Q5SDM S262E_F and Q5SDM S262E_R for pQKY::mCherry:QKY_{S262E} (Supplementary Data 1229 1230 11). All PCR-based constructs were confirmed by sequencing. Floral tissue for quantitative realtime PCR (gPCR) was harvested from plants grown under long day conditions. With minor 1231 changes, tissue collection, RNA extraction and qPCR analysis were performed as described 1232

^{118,119}. For detection of *mCherry:QKY*, expression by gPCR, primers mCherry-gRT-For/mCherry-1233 1234 gRT-Rev were employed (Table). Average mCherry:QKY expression was calculated relative to 1235 the expression of three control genes (AT4G33380; AT2G28390; AT5G46630) (Supplementary Data 11) and normalized to the wild type control for visualization. A. thaliana (L.) Heynh. var. 1236 1237 Landsberg (erecta) (Ler) was used as wild type strain. The likely null allele qky-9 (Ler) has been described previously ⁴³. *qky*-9 mutant plants were transformed using Agrobacterium strain 1238 GV3101/pMP90¹²⁰ and the floral dip method ⁵⁰. Transgenic T1 plants were selected on 1239 Kanamycin (50 µg/ml) and transferred to soil for further inspection. Plants were grown as 1240 described earlier ⁴³. Four out of 17 independent T1 transformants showed a wild type phenotype 1241 for the S262A transgenic line, while 13/17 T1 transformants displayed no (7/17) or partial (6/17) 1242 rescue, a notable decrease in functionality compared to the wild type construct (14/17 rescue, 1243 2/17 partial rescue, 1/17 no rescue). For the S262E transgene, 11 out of 13 independent T1 1244 transformants showed a wild type phenotype (1/13 partial rescue, 1/13 no rescue). 1245

Floral organs were imaged using a Leica SAPO stereo microscope equipped with a digital MC 1246 170 HD camera (Leica Microsystems GmbH, Wetzlar, Germany). Images were adjusted for 1247 color and contrast using ImageJ/Fiji software ¹²¹. Confocal laser scanning microscopy of *qky*-9 1248 pQKY::mCHerry:QKY, qky-9 pQKY::mCHerry:QKY_{S262A} and qky-9 pQKY::mCHerry:QKYS_{262E} 1249 six days-old seedling roots was performed with an Olympus FV1000 setup using an inverted 1250 IX81 stand and FluoView software (FV10-ASW version 01.04.00.09) (Olympus Europa GmbH, 1251 Hamburg, Germany) equipped with a water-corrected 40x objective (NA 0.9) at 3x digital zoom. 1252 Confocal high sensitivity detection (HSD) was employed involving two gallium arsenide 1253 1254 phosphide photomultipliers (GaAsP PMTs) mounted equidistantly to the probe. The experiment was performed two times independently with 2 and 12 roots for qky-9 pQKY::mCHerry:QKY, 6 1255 1256 and 15 roots for *qky*-9 *pQKY::mCHerry:QKY*_{S262A} and 15 roots for qky-9 1257 pQKY::mCHerry:QKYS_{262F}.

1258 Locus identifiers

1259 Gene locus identifiers are listed for all gene names mentioned in the manuscript 1260 (Supplementary Data 11).

1261 Online methods references

126244Heazlewood, J. L., Verboom, R. E., Tonti-Filippini, J., Small, I. & Millar, A. H. SUBA: the1263Arabidopsis Subcellular Database. Nucleic Acids Res 35, D213-218, doi:10.1093/nar/gkl8631264(2007).

126545Nishimura, N. *et al.* Structural mechanism of abscisic acid binding and signaling by dimeric PYR1.1266Science **326**, 1373-1379, doi:10.1126/science.1181829 (2009).

- 1267
 46
 Smyth, D. R., Bowman, J. L. & Meyerowitz, E. M. Early flower development in Arabidopsis. *Plant*

 1268
 Cell **2**, 755-767, doi:10.1105/tpc.2.8.755 (1990).
- 126947Johnson-Brousseau, S. A. & McCormick, S. A compendium of methods useful for characterizing1270Arabidopsis pollen mutants and gametophytically-expressed genes. *Plant J* **39**, 761-775,1271doi:10.1111/j.1365-313X.2004.02147.x (2004).
- 127248Sprunck, S. *et al.* Egg cell-secreted EC1 triggers sperm cell activation during double fertilization.1273Science **338**, 1093-1097, doi:10.1126/science.1223944 (2012).
- 127449Karimi, M., Inze, D. & Depicker, A. GATEWAY vectors for Agrobacterium-mediated plant1275transformation. *Trends Plant Sci* 7, 193-195 (2002).
- 1276 50 Clough, S. J. & Bent, A. F. Floral dip: a simplified method for Agrobacterium-mediated 1277 transformation of Arabidopsis thaliana. *Plant J* **16**, 735-743 (1998).
- 127851Schmid, M. *et al.* A gene expression map of Arabidopsis thaliana development. *Nat Genet* **37**,1279501-506, doi:10.1038/ng1543 (2005).
- 128052Boyes, D. C. *et al.* Growth stage-based phenotypic analysis of Arabidopsis: a model for high1281throughput functional genomics in plants. *Plant Cell* **13**, 1499-1510 (2001).
- 1282 53 Bowman, J. L. *Arabidopsis : an atlas of morphology and development*. (Springer-Verlag, 1994).
- 128354Bradford, M. M. A rapid and sensitive method for the quantitation of microgram quantities of1284protein utilizing the principle of protein-dye binding. Anal Biochem 72, 248-254 (1976).
- 128555Ruprecht, B. et al. Optimized Enrichment of Phosphoproteomes by Fe-IMAC Column1286Chromatography. Methods Mol Biol 1550, 47-60, doi:10.1007/978-1-4939-6747-6_5 (2017).
- 128756Marx, H. et al. A large synthetic peptide and phosphopeptide reference library for mass1288spectrometry-based proteomics. Nat Biotechnol **31**, 557-564, doi:10.1038/nbt.2585 (2013).
- 128957Ruprecht, B., Zecha, J., Zolg, D. P. & Kuster, B. High pH Reversed-Phase Micro-Columns for1290Simple, Sensitive, and Efficient Fractionation of Proteome and (TMT labeled) Phosphoproteome1291Digests. Methods Mol Biol 1550, 83-98, doi:10.1007/978-1-4939-6747-6_8 (2017).
- 129258Smith, P. K. *et al.* Measurement of protein using bicinchoninic acid. Anal Biochem 150, 76-85,1293doi:10.1016/0003-2697(85)90442-7 (1985).
- 129459Zolg, D. P. et al. PROCAL: A Set of 40 Peptide Standards for Retention Time Indexing, Column1295Performance Monitoring, and Collision Energy Calibration. Proteomics 17,1296doi:10.1002/pmic.201700263 (2017).
- 129760Hahne, H. et al. DMSO enhances electrospray response, boosting sensitivity of proteomic1298experiments. Nat Methods 10, 989-991, doi:10.1038/nmeth.2610 (2013).
- 129961Bian, Y. *et al.* Robust, reproducible and quantitative analysis of thousands of proteomes by1300micro-flow LC-MS/MS. *Nat Commun* **11**, 157, doi:10.1038/s41467-019-13973-x (2020).
- 130162Tyanova, S., Temu, T. & Cox, J. The MaxQuant computational platform for mass spectrometry-1302based shotgun proteomics. Nat Protoc 11, 2301-2319, doi:10.1038/nprot.2016.136 (2016).
- 130363Hanada, K. *et al.* sORF finder: a program package to identify small open reading frames with high1304coding potential. *Bioinformatics* **26**, 399-400, doi:10.1093/bioinformatics/btp688 (2010).
- 130564Grabherr, M. G. et al. Full-length transcriptome assembly from RNA-Seq data without a1306reference genome. Nat Biotechnol 29, 644-652, doi:10.1038/nbt.1883 (2011).
- 130765Li, W., Jaroszewski, L. & Godzik, A. Clustering of highly homologous sequences to reduce the size1308of large protein databases. *Bioinformatics* **17**, 282-283, doi:10.1093/bioinformatics/17.3.2821309(2001).
- Perkins, D. N., Pappin, D. J., Creasy, D. M. & Cottrell, J. S. Probability-based protein identification
 by searching sequence databases using mass spectrometry data. *Electrophoresis* 20, 3551-3567,
 doi:10.1002/(SICI)1522-2683(19991201)20:18<3551::AID-ELPS3551>3.0.CO;2-2 (1999).

- 131367Franken, H. *et al.* Thermal proteome profiling for unbiased identification of direct and indirect1314drug targets using multiplexed quantitative mass spectrometry. *Nat Protoc* **10**, 1567-1593,1315doi:10.1038/nprot.2015.101 (2015).
- 131668Toprak, U. H. *et al.* Conserved peptide fragmentation as a benchmarking tool for mass1317spectrometers and a discriminating feature for targeted proteomics. *Mol Cell Proteomics* 13,13182056-2071, doi:10.1074/mcp.O113.036475 (2014).
- 131969Onate-Sanchez, L. & Vicente-Carbajosa, J. DNA-free RNA isolation protocols for Arabidopsis1320thaliana, including seeds and siliques. *BMC Res Notes* 1, 93, doi:10.1186/1756-0500-1-93 (2008).
- 132170Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence1322data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170 (2014).
- 132371Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq1324quantification. Nat Biotechnol 34, 525-527, doi:10.1038/nbt.3519 (2016).
- 132572Silva, J. C., Gorenstein, M. V., Li, G. Z., Vissers, J. P. & Geromanos, S. J. Absolute quantification of1326proteins by LCMSE: a virtue of parallel MS acquisition. *Mol Cell Proteomics* 5, 144-156,1327doi:10.1074/mcp.M500230-MCP200 (2006).
- 132873The Gene Ontology, C. Expansion of the Gene Ontology knowledgebase and resources. Nucleic1329Acids Res 45, D331-D338, doi:10.1093/nar/gkw1108 (2017).
- 133074Cox, J. & Mann, M. 1D and 2D annotation enrichment: a statistical method integrating1331quantitative proteomics with complementary high-throughput data. BMC Bioinformatics 131332Suppl 16, S12, doi:10.1186/1471-2105-13-S16-S12 (2012).
- 133375Tyanova, S. et al. The Perseus computational platform for comprehensive analysis of1334(prote)omics data. Nat Methods 13, 731-740, doi:10.1038/nmeth.3901 (2016).
- 133576Olsen, J. V. et al. Global, in vivo, and site-specific phosphorylation dynamics in signaling1336networks. Cell 127, 635-648, doi:10.1016/j.cell.2006.09.026 (2006).
- 1337 77 Uhlen, M. *et al.* Transcriptomics resources of human tissues and organs. *Mol Syst Biol* 12, 862, doi:10.15252/msb.20155865 (2016).
- 133978Rijpkema, A. S., Vandenbussche, M., Koes, R., Heijmans, K. & Gerats, T. Variations on a theme:1340changes in the floral ABCs in angiosperms. Semin Cell Dev Biol 21, 100-107,1341doi:10.1016/j.semcdb.2009.11.002 (2010).
- 134279Loytynoja, A. Phylogeny-aware alignment with PRANK. Methods Mol Biol 1079, 155-170,1343doi:10.1007/978-1-62703-646-7_10 (2014).
- 134480Castresana, J. Selection of conserved blocks from multiple alignments for their use in1345phylogenetic analysis. Mol Biol Evol 17, 540-552, doi:10.1093/oxfordjournals.molbev.a0263341346(2000).
- 134781Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. Mol Biol Evol 24, 1586-1591,1348doi:10.1093/molbev/msm088 (2007).
- 134982van der Graaf, A. et al. Rate, spectrum, and evolutionary dynamics of spontaneous1350epimutations. Proc Natl Acad Sci U S A 112, 6676-6681, doi:10.1073/pnas.1424254112 (2015).
- 135183Gebert, D., Jehn, J. & Rosenkranz, D. Widespread selection for extremely high and low levels of1352secondary structure in coding sequences across all domains of life. Open Biol 9, 190020,1353doi:10.1098/rsob.190020 (2019).
- 135484Camiolo, S., Melito, S. & Porceddu, A. New insights into the interplay between codon bias1355determinants in plants. DNA Res 22, 461-470, doi:10.1093/dnares/dsv027 (2015).
- Drummond, D. A., Bloom, J. D., Adami, C., Wilke, C. O. & Arnold, F. H. Why highly expressed 1356 85 1357 proteins evolve slowly. Proc Natl Acad Sci U S A 102, 14338-14343, 1358 doi:10.1073/pnas.0504070102 (2005).

- 135986Das, S. & Bansal, M. Variation of gene expression in plants is influenced by gene architecture1360and structural properties of promoters. PLoS One 14, e0212678,1361doi:10.1371/journal.pone.0212678 (2019).
- 136287Celaj, A. *et al.* Quantitative analysis of protein interaction network dynamics in yeast. *Mol Syst*1363*Biol* **13**, 934, doi:10.15252/msb.20177532 (2017).
- 136488Niederhuth, C. E. *et al.* Widespread natural variation of DNA methylation within angiosperms.1365Genome Biol 17, 194, doi:10.1186/s13059-016-1059-0 (2016).
- 136689Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for1367RNA-seq data with DESeq2. Genome Biol 15, 550, doi:10.1186/s13059-014-0550-8 (2014).
- 136890Zhou, X. & Stephens, M. Genome-wide efficient mixed-model analysis for association studies.1369Nat Genet 44, 821-824, doi:10.1038/ng.2310 (2012).
- 137091Nakazawa, N. fmsb: Functions for Medical Statistics Book with some Demographic Data. R1371package version 0.6.3. https://CRAN.R-project.org/package=fmsb (2018).
- 137292R Core Team. R: A language and environment for statistical computing. R Foundation for1373Statistical Computing (https://www.R-project.org/.) (2014).
- 137493Tibshirani, R. Regression Shrinkage and Selection via the Lasso. Journal of the Royal Statistical1375Society 58, 267-288 (1996).
- 137694Zhang, Z. Variable selection with stepwise and best subset approaches. Ann Transl Med 4, 136,1377doi:10.21037/atm.2016.03.35 (2016).
- 137895Knecht, W. Pilot Willingness to take Off Into Marginal Weather, Part II: Antecedent Overfitting1379With Forward Stepwise Logistic Regression. Final Report, Federal Aviation Administration (2005).
- 138096Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear Models via1381Coordinate Descent. J Stat Softw 33, 1-22 (2010).
- 138297Groemping, U. Relative importance for linear regression in R: The package relaimpo. Journal of1383Statistical Software 17, 1-27, doi: 10.18637/jss.v017.i01 (2007).
- 138498Heusel, M. *et al.* Complex-centric proteome profiling by SEC-SWATH-MS. *Mol Syst Biol* 15,1385e8438, doi:10.15252/msb.20188438 (2019).
- 138699McBride, Z., Chen, D., Reick, C., Xie, J. & Szymanski, D. B. Global Analysis of Membrane-1387associated Protein Oligomerization Using Protein Correlation Profiling. *Mol Cell Proteomics* 16,13881972-1989, doi:10.1074/mcp.RA117.000276 (2017).
- 1389 100 Ruepp, A. *et al.* CORUM: the comprehensive resource of mammalian protein complexes--2009.
 1390 *Nucleic Acids Res* 38, D497-501, doi:10.1093/nar/gkp914 (2010).
- 1391101Zhang, B. & Horvath, S. A general framework for weighted gene co-expression network analysis.1392Stat Appl Genet Mol Biol 4, Article17, doi:10.2202/1544-6115.1128 (2005).
- 1393102Langfelder, P., Zhang, B. & Horvath, S. Defining clusters from a hierarchical cluster tree: the1394DynamicTreeCutpackageforR.Bioinformatics24,719-720,1395doi:10.1093/bioinformatics/btm563 (2008).
- 1396103Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new perspectives on1397genomes, pathways, diseases and drugs. Nucleic Acids Res 45, D353-D361,1398doi:10.1093/nar/gkw1092 (2017).
- 1399104Fabregat, A. et al. The Reactome pathway Knowledgebase. Nucleic Acids Res 44, D481-487,1400doi:10.1093/nar/gkv1351 (2016).
- 1401105Hochberg, Y. B. a. Y. Controlling the false discovery rate: a practical and powerful approach to1402multiple testing. Journal of the Royal Statistical Society 57, 289-300 (1995).
- 1403106Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for1404interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-15550,1405doi:10.1073/pnas.0506580102 (2005).

- 1406107List, M. *et al.* KeyPathwayMinerWeb: online multi-omics network enrichment. Nucleic Acids Res140744, W98-W104, doi:10.1093/nar/gkw373 (2016).
- 1408108Letunic, I. & Bork, P. 20 years of the SMART protein domain annotation resource. Nucleic Acids1409Res 46, D493-D496, doi:10.1093/nar/gkx922 (2018).
- 1410109Wagih, O. ggseqlogo: a versatile R package for drawing sequence logos. *Bioinformatics* **33**, 3645-14113647, doi:10.1093/bioinformatics/btx469 (2017).
- 1412110Goel, R., Harsha, H. C., Pandey, A. & Prasad, T. S. Human Protein Reference Database and1413Human Proteinpedia as resources for phosphoproteome analysis. *Mol Biosyst* 8, 453-463,1414doi:10.1039/c1mb05340j (2012).
- 1415111Zourelidou, M. *et al.* The polarly localized D6 PROTEIN KINASE is required for efficient auxin1416transport in Arabidopsis thaliana. *Development* **136**, 627-636, doi:10.1242/dev.028365 (2009).
- 1417112Mayer, U. B., G.; Jurgens, G. Apical-basal pattern formation in the Arabidopsis embryo: studies1418on the role of the gnom gene. *Development* **177**, 149-162 (1993).
- 1419 113 Moes, D., Himmelbach, A., Korte, A., Haberer, G. & Grill, E. Nuclear localization of the mutant
 protein phosphatase abi1 is required for insensitivity towards ABA responses in Arabidopsis. *Plant J* 54, 806-819, doi:10.1111/j.1365-313X.2008.03454.x (2008).
- 1422114Tischer, S. V. *et al.* Combinatorial interaction network of abscisic acid receptors and coreceptors1423from Arabidopsis thaliana. *Proc Natl Acad Sci U S A* **114**, 10280-10285,1424doi:10.1073/pnas.1706593114 (2017).
- 1425115Waterhouse, A. *et al.* SWISS-MODEL: homology modelling of protein structures and complexes.1426Nucleic Acids Res 46, W296-W303, doi:10.1093/nar/gky427 (2018).
- 1427 116 Berman, H. M. *et al.* The Protein Data Bank. *Nucleic Acids Res* 28, 235-242 (2000).
- 1428117Pettersen, E. F. *et al.* UCSF Chimera--a visualization system for exploratory research and analysis.1429J Comput Chem 25, 1605-1612, doi:10.1002/jcc.20084 (2004).
- 1430118Box, M. S., Coustham, V., Dean, C. & Mylne, J. S. Protocol: A simple phenol-based method for 96-1431well extraction of high quality RNA from Arabidopsis. *Plant Methods* 7, 7, doi:10.1186/1746-14324811-7-7 (2011).
- 1433119Enugutti, B. *et al.* Regulation of planar growth by the Arabidopsis AGC protein kinase UNICORN.1434*Proc Natl Acad Sci U S A* **109**, 15060-15065, doi:10.1073/pnas.1205089109 (2012).
- 1435120Koncz, C. & Schell, J. The promoter of TL-DNA gene 5 controls the tissue-specific expression of1436chimaeric genes carried by a novel type of Agrobacterium binary vector. Molecular and General1437Genetics MGG 204, 383-396, doi:10.1007/bf00331014 (1986).
- 1438121Schindelin, J. *et al.* Fiji: an open-source platform for biological-image analysis. Nat Methods 9,1439676-682, doi:10.1038/nmeth.2019 (2012).
- 1440 122 Vizcaino, J. A. *et al.* 2016 update of the PRIDE database and its related tools. *Nucleic Acids Res*1441 44, D447-456, doi:10.1093/nar/gkv1145 (2016).

1443 Acknowledgements

We thank the NGS@tum core facility for RNA sequencing, Rachele Tofanelli for help with imaging the ovules, Robert J. Schmitz for providing data access for the feature analysis and Maria Reinecke, Florian Bayer and Stefanie Galinec for MS measurements. This work was in part funded by the German Science Foundation (DFG, SFB924), a research fellowship to HS by the Japan Society for the Promotion of Sciences and a research fellowship to XC by the Chinese Research Council.

1450 Author contributions

JM performed (phosho)proteomic and transcriptomic experiments under the supervision of BK.
SR and HS performed AGC kinase experiments in plants under the supervision of GJ and CS.
MP, AC and XC performed phosphomutant analysis under the supervision of EG and KS. PC
and SS generated and provided plant material. JM, MF, MM, DL, SA, DPZ, TM, CD, AD and
RRH performed data analysis under the supervision of BK, KFXM, PF, MB, TH and FJ. ML, PS
and TS generated Arabidopsis resource databases under supervision of MW and JB. JM, CS
and BK conceptualized the project and wrote the manuscript. All authors edited the manuscript.

1458 **Competing interests**

1459 MW and BK are founders and shareholders of OmicScouts GmbH and msAld GmbH. They 1460 have no operational role in the companies. MF and DPZ are founders and shareholders of 1461 msAld GmbH. TM and MB are employees and/or shareholders of Cellzome GmbH. The 1462 remaining authors declare no competing interests.

1463 Data and software availability

The data supporting the findings of this study are available within the paper, the supplementary information and the public repositories. Source data for main Figure 1-5 and Extended Data Figure 1-9 are included. Transcriptome sequencing and quantification data are available at ArrayExpress (www.ebi.ac.uk/arrayexpress) under the identifier E-MTAB-7978. The raw mass spectrometric data and MaxQuant result files have been deposited to the ProteomeXchange Consortium via PRIDE ¹²², with the dataset identifier PXD013868.

1470 **Corresponding author:** Bernhard Kuster (kuster@tum.de)

1472 Extended data figure legends

1473 Extended Data Figure 1 | Descriptive analysis of the multi-omic tissue atlas.

a, Pairwise global Pearson expression correlation analysis of all 30 tissues (n = 1 measurement
per tissue) on the transcriptome level (lower triangle) and proteome level (upper triangle) using
all identified gene loci. The data shows that proteins correlate more strongly between tissues
than transcripts. Turquoise squares mark examples for morphologically highly similar tissues.
Tissues are coloured as in Figure 1: Flower (light grey), seed (dark brown), pollen (yellow), stem
(dark green), leaf (light green), root (dark grey), fruit (light brown), callus (magenta), cell culture
(blue).

- b, Scatter plots showing highly reproducible abundance measurements for transcript (upper panels) and protein (lower panels) in morphologically similar tissues that were marked in panel a, namely node (ND) versus internode (IND), leaf distal (LFD) versus leaf proximal (LFP) and root (RT) versus root upper zone (RTUZ). *R* denotes the Pearson correlation coefficient and n denotes the number of transcripts or proteins shown in the plots.
- 1486 c, Percentage of genes encoded by a specific chromosome that were identified at the 1487 transcriptome, proteome or phosphoproteome level.
- d, Percentage of Swiss-Prot and TrEMBL protein database entries as well as protein evidence categories from UniProt that were identified at the transcriptome, proteome or phosphoproteome level. Evidence level: (1) 'protein evidence', (2) 'transcript evidence', (3) 'homology', (4) 'predicted', (5) 'uncertain'.
- e, Comparison of protein identifications between an earlier Arabidopsis proteome study by
 Baerenfaller et al. based on 12 tissues ⁷, this study (30 tissues) and the number of proteincoding genes in Araport11.
- f, iBAQ intensity distribution of proteins identified in this study. Proteins also identified byBaerenfaller et al. are projected into the same plot.
- g, Left panel: Proportion of identified p-sites on S, T or Y residues with highly confident
 localization of the phosphorylation site within the identified peptide sequence (termed class I psites if the localization score is >0.75). Right panel: distribution of proteins for which
 phosphorylated S, T or Y residues were identified.
- h, Left panel: Venn diagram comparing phosphoprotein datasets from van Wijk et al. ⁸,
 PhosPhAT4.0 and this study. Right panel: Venn diagram comparing p-site localization
 confidence between class I sites identified in this study and the low and high confidence
 datasets reported by van Wijk et al.

1505 Extended Data Figure 2 | Proteogenomics and dynamic range of transcript and protein 1506 expression.

1507 a, Number of identified N-terminal (NT) or C-terminal (CT) peptides of proteins in either 1508 unmodified or phosphorylated form.

b, Frequency of amino acids following the initiator methionine in N-terminal peptides with (-X) or
without cleavage of the initiator methionine (M-X). X denotes the amino acid following the start
codon.

- c, Frequency of protein N-terminal acetylation for amino acids in (b). Because trypsin was used
 for protein digestion, the frequencies for arginine (R) and lysine (K) could not be determined
 (n.d.).
- d, Distribution of peptide-based sequence coverage of proteins in individual tissues and for the
 combined dataset (tissue abbreviations as in Figure 1). Boxes contain 50% of the data and
 show the median as a black line. The upper and lower quartile ranges are shown as whiskers.
 The number of proteins is indicated for each tissue.
- e, Pie charts showing the percentage of proteins identified by <3, 3-10 or >10 peptides eitherallowing shared (razor) peptides or restricting to unique peptides only.
- 1521 f, Left panel: number of protein isoforms detected at the transcript and protein level compared to 1522 the number of all annotated isoforms in Araport11. Right panel: Number of multiple isoforms of 1523 the same gene distinguished at the peptide level.
- 1524 g, Validation of protein isoform and short open reading frame (sORF) identification by 1525 comparing the tandem mass spectra from the tissue atlas to those of synthetic peptide 1526 reference standards. The normalized spectral contrast angle (SA) was used as a similarity 1527 metric (see methods). Candidate isoforms and sORFs were considered valid if the SA of the 1528 spectra was >0.7. This data is reported in Supplementary Data 3.
- h, Amino acid sequence and mirror plots of tandem mass spectra for two peptides of the sORF BIP138_4. The spectra pointing upwards were collected from tissue digests, those pointing downwards were collected from synthetic peptides. The normalized spectral contrast angle (SA) and Pearson correlation coefficient (R) were used as similarity metrics (see methods) and indicate that both high scoring spectra (n=1 acquired spectra) are near identical, thus validating the identification of this sORF as a expressed protein.
- i, Dynamic range of transcript expression (grey) and proportion of transcripts that were also
 identified at the protein level projected into this plot (blue). OM: orders of magnitude. Note that
 for lower abundance transcripts, fewer proteins were detected.

j, Dynamic range of protein expression and proportion of proteins with phosphorylation
 evidence. Note that protein expression spans 6 OM whereas transcript expression only spans 4
 OM (panel i). Further note that phosphorylation was detected across the entire protein
 expression range.

k, Percentage of all annotated kinases (K), phosphatases (P), transcription factors (TF) and
transcription regulators (TR) detected at the transcript, protein or phosphoprotein level.
Numbers below the x-axis denote the number of genes for these protein classes in the AT
genome.

1546 Extended Data Figure 3 | Descriptive analysis of transcript and protein expression in 1547 tissues.

a. Distribution of expression specificity categories for protein and transcript identifications. See online methods for the definition of these categories. Briefly, there are very few transcripts and proteins that are only expressed in a single tissue. Note that the quantities of the shared transcripts or proteins can differ vastly between tissues (see below).

b, Left panel: protein identifications shared between flower (FL) and flower organs showing an
almost complete qualitative overlap of proteins. Sepal (SP), petal (PT), stamen (ST), carpel
(CP). Right panel: clustering of z-scored protein intensities showing distinct quantitative
expression differences between flower organs.

c, Expression analysis of flower organ identity marker at the protein and transcript level.
PISTILLATA (PI, green), APETALA3 (AP3, red), APETALA1 (AP1, orange), AGAMOUS (AG,
blue). The expression of these markers is in line with the model of flower organ identity (AP1
expression marking sepal, AP1, AP3, PI marking petal, AG, AP3, PI marking stamen and AG
marking carpel).

d, Total number of transcripts plotted against the total number of proteins detected in each individual tissue (n = 30 tissues) showing that the more genes are expressed as mRNAs, the more proteins can be detected in a tissue (Pearson correlation R = 0.79). Tissues are coloured according to tissue groups as in Figure 1.

e, Cumulative abundance plots of intensity-ranked identifications of transcripts and proteins for five representative tissues. The five most abundant transcripts and proteins are listed in descending order for each tissue. Note that these are generally not the same. Further note that the characteristics of the plots are not the same for all tissues. In flower, the protein line rises more quickly than the transcript line. The opposite is true for pollen and a more even characteristic is observed in seed.

1571 f, Distribution of shared and unique identifications among the 100 most abundant transcripts and 1572 proteins in each tissue. Note that relatively few proteins and transcripts are found together on 1573 the list of the 100 most abundant transcripts and proteins. This demonstrates that the 1574 quantitative differences in transcript and protein expression are more important in defining a 1575 tissue than the qualitative expression of transcripts or proteins.

1576 g, List of 11 proteins which were found as the most abundant protein (in at least one tissue) and 1577 their proportion of the total iBAQ intensity in each tissue. Note that individual proteins can 1578 represent up to 9% of the total protein in a given tissue.

- 1579 h, Principal component analysis (PCA) of the core tissue proteomes and transcriptomes (i.e. the proteins and transcripts that were identified in every tissue) using z-scored abundances. Note 1580 that only about 30% of all protein and 20% of all mRNAs were detected in every of the 30 1581 tissues despite the fact that all tissues were deeply profiled at both protein and transcript level. 1582 1583 This shows that strong qualitative and quantitative expression differences exist between tissues. 1584 The PCA separates tissues into photosynthetically active versus inactive tissues (component 1) 1585 and separates pollen from all other tissues (component 2) implying that the molecular composition of pollen is particularly different from all other tissues. 1586
- i, Proportion of the total summed protein intensity for genes with specific subcellular 1587 compartment annotation (from SUBA ⁴⁴, see methods) in the different tissue groups. The 1588 1589 comparison of photosynthetically active and inactive tissues shows that e.g. a majority of the 1590 protein content in photosynthetically active tissues are contained in the plastids, whereas most 1591 protein is found in the cytosol for photosynthetically inactive tissues. Proteins with only one 1592 single subcellular compartment annotation were selected for the plot and the proportion of their 1593 iBAQ intensities were averaged for each tissue group. Nucleus (n = 1,393), endoplasmatic reticulum (n = 58), golgi (n = 68), peroxisome (n = 67), plastid (n = 525), mitochondrion (n = 1594 317), vacuole (n = 71), cytosol (n = 385), cytoskeleton (n = 1), plasma membrane (n = 268), 1595 1596 extracellular (n = 351).
- 1597 Extended Data Figure 4 | Relationships between transcript and protein levels.
- a, Pearson correlation (*R*) of transcriptome and proteome expression (core datasets; n = 5,043)
 for each tissue.
- b, Pearson correlation (*R*) between measured and predicted protein abundance levels in all
 tissues. Predicted protein abundance levels were obtained from the best fitting feature selection
 model for each tissue (see methods). The number of genes used for the correlation analysis is
 indicated for each tissue.
- c, Violin plots showing the spread in relative contribution of selected features to the prediction of
 gene-level protein abundance across tissues (n = 30 tissues) using our model. Violin shapes
 show the kernel density estimation of the data distribution and the median as white dot. Thick
 black bars denote the interquartile range.

d, Specific nucleotide sequence motifs in 5'UTRs of mRNAs contribute to the prediction of protein levels in a subset of tissues. Clustering tissues based on the presence or absence of detected 5'UTR motifs shows that several features are repeatedly selected for inclusion in the model while others appear to be more tissue-specific.

1612 e, Based on the observation that the ratio of non-synonymous to synonymous nucleotide 1613 substitutions (Dn/Ds) between orthologous of Arabidopsis thaliana and Arabidopsis lyrata 1614 contributed to the prediction of protein levels (see above), we analysed this feature in more detail. Left panel: Distribution of the ratio of non-synonymous to synonymous nucleotide 1615 1616 substitutions (Dn/Ds) for orthologous genes in Arabidopsis thaliana and Arabidopsis lyrata. The distribution is plotted for the example of 'leaf distal' (n = 6,447 genes). To compare evolutionarily 1617 1618 conserved genes (defined by low non-synonymous to synonymous substitutions [Dn/Ds] ratios) and genes that evolve neutrally or are under positive selection (high Dn/Ds), we selected the 1619 bottom 5% and top 5% of the Dn/Ds ratio distribution, respectively. Right panel: evolutionarily 1620 conserved genes (low Dn/Ds ratio) show 10-20x higher protein abundance than genes under 1621 1622 evolutionary pressure. Boxes contain 50% of the data and show the median as a black line. 1623 Whiskers denote 1.5 times the interguartile range. Outliers were omitted from the plot for clarity.

1624 f, Time course analysis of median protein abundance changes upon cycloheximide (translation block, CHX) or MG132 (proteasome block) treatment vs time-matched DMSO control samples 1625 1626 (see methods). Boxes contain 50% of the data and show medians as black lines. Whiskers 1627 denote 1.5 times the interguartile range. Outliers were omitted from the plot for clarity but were included in the statistical tests below. Shown are either all proteins in the experiment (n = 8,920, 1628 grey), proteins that have a high (n = 425, red) or low (n = 254, blue) PTR in seed (defined as in 1629 main Figure 3d). Differences between time points were tested for significance within each 1630 1631 subset (all; high PTR; low PTR) using one-way ANOVA (analysis of variance) and the post-hoc 1632 Tukey HSD (Honestly significant difference) test. *** HSD p-Value < 0.001 (all_CHX8-CHX16: p-Value < 1e-7; all_CHX8-CHX24: p-Value < 1e-7; all_CHX16-CHX24: p-Value = 0.0002; 1633 1634 highPTR_CHX8-CHX24: p-Value = 0.0003; lowPTR_CHX8-CHX16: p-Value = 0.0000004; lowPTR_CHX8-CHX24: p-Value < 1e-7; lowPTR_CHX16-CHX24: p-Value < 1e-7). 1635

1636 g, Representative images of seeds after 4 days of incubation with CHX, MG132 or DMSO 1637 control medium (n = 1). Germination was completely inhibited by CHX and partially inhibited by 1638 MG132 showing that the drug treatments were effective.

1639 Extended Data Figure 5 | Correlations between transcriptomes, proteomes and 1640 phosphoproteomes.

a, Median protein to mRNA ratios (PTRs) across tissues plotted against the inter-tissue
 variation of these PTRs (expressed as median absolute deviation, MAD; proteins and
 transcripts had to be detected in at least 10 matching tissues to be included in the analysis).

Examples for genes with low transcript/high protein (rbcL and petA), and high transcript/low protein (IAA8 and IAA13) are marked by arrows. The bar blot shows the MAD range segmented into 5 quantiles each containing the same number of genes (coloured bars and dashed lines). It is evident that most genes have reasonably stable PTRs across tissues.

b, Same as panel a (dataset n = 14,069) but for transcript (left plot) and protein measurements (right plot). As can be seen, there is somewhat more variation in protein levels across tissues than there is mRNA variation (80% of all transcripts show a MAD of <1; 80% of all proteins show a MAD of 1.2). There is also more variation in the protein levels across tissues for low abundant proteins. This may in part due to technical limitations as low abundance proteins can generally be less accurately quantified.

1654 c, Same as panel a but for the ratio of phosphorylation site versus protein abundance. P-sites
1655 and proteins had to be detected in at least 10 matching tissues to be included in the analysis
1656 (n = 13,793).

d, Same as panel b (dataset n = 13,793) but for phosphorylation site abundance. Note that psite abundance shows somewhat greater variation across tissues than protein abundance (60%
of all p-sites show MAD<1 compared to 80% of all proteins, see panel b). Again, this may in part
due to technical limitations as p-site quantification is performed on a peptide level and does not
benefit from aggregating multiple peptide quantifications into one value for protein quantification.

1662 Extended Data Figure 6 | Inferring redundant gene function and physical interactions 1663 from co-expression analysis.

a, Scatter plot of Pearson correlation coefficients (*R*) as a measure for co-expression across tissues for all pairs of proteins (x-axis) and all pairs of transcripts (y-axis) (core dataset only, n =5,043) along with their marginal histograms. Colours denote the log₂-normalized STRING scores of individual gene pairs as a measure of known or predicted direct (physical) or indirect (functional) associations. The data show that strong co-expression of transcripts or proteins or both are more strongly related (physically or functionally) than transcripts and proteins that are not.

b, Co-expression analysis of duplicated genes (pairs had to be detected in at least 10 matching tissues to be included in the analysis). The density plots show the distribution of Pearson correlation coefficients (*R*) of co-expressed transcripts (grey) or proteins (blue) for genes that arose by whole genome duplications (WGD), local duplications (local) or transposon-mediated duplications (transposed). Randomly selected gene pairs are shown as control (random). Medians (\bar{x}) are given and displayed as dotted lines. The data shows that there is substantial co-expression of duplicated genes implying that these genes likely have redundant functions. 1678 c, Left panel: Protein level Pearson correlation coefficient (*R*) values (from panel b) for all 1679 duplicate gene pairs (WGD, local, transposed) plotted against the protein abundance ratio of 1680 each pair (average across 30 tissues, see methods). Blue arrows point out one example each 1681 for a high or low ratio of protein expression for the duplicated genes. Right panel: Example for 1682 tissue-resolved protein intensity proportions (top3, see methods) for the duplicate pair: MAC5A 1683 and MAC5B. The data shows that irrespective of the tissue, MAC5A is always much higher 1684 expressed than MAC5B. Tissues are coloured as in Figure 1.

- d, Upper panel: ranked protein abundance ratio for selected duplicate pairs (average ± SD; n = 30) and annotated for phenotypic effects (lower panel) in the loss-of-function mutant for either duplicate 1 or duplicate 2 (+). Absence of a phenotypic effect is marked by (-). The data shows that asymmetric protein expression within duplicate pairs can be associated with the occurrence of a phenotype in the loss-of-function mutant of the higher expressed duplicate protein implying a dominant functional role of the more highly expressed protein. Blue arrows point out MAC5A/MAC5B and PHB3/PHB4 as examples.
- e, Inference of physical protein-protein interactions from co-expression data. Distribution of pairwise Pearson correlation coefficients (*R*) of co-expressed proteins across (at least 10) tissues that are subunits of selected protein complexes. An *R* value of >0.5 (shaded in grey) was chosen as a cut-off for the selection of proteins for subsequent analysis in order to make sure that proteins present in well characterized protein complexes are retained. CONSTITUTIVE PHOTOMORPHOGENESIS9 SIGNALOSOME (CSN), CELLULOSE SYNTHASE (CESA).
- 1698 f, Recovery of annotated protein-protein interactions by co-expression analysis. Distribution of Pearson correlation coefficients (R) of pairs of transcripts (grey) or protein (blue) that are 1699 annotated to interact physically in the AtPIN database ³³ (pairs had to be detected in at least 10 1700 matching tissues to be included in the analysis). Subsets of the AtPIN database, namely 1701 1702 interactions detected by the yeast two-hybrid (Y2H) method, by affinity purification - mass 1703 spectrometry (AP-MS) or by both (Y2H+AP-MS) are displayed separately. R values >0.5 are 1704 shaded in blue (protein). Median (\bar{x} , dotted lines). The data shows that co-expression only 1705 recovers a minority of annotated physical interactions and that interactions supported by more than one line of experimental evidence also tend to show stronger co-expression. 1706

Extended Data Figure 7 | Inferring protein complexes and subunit stoichiometry from proteome correlation profiling using size-exclusion chromatography – mass spectrometry.

a, Molecular weights of monomeric proteins (MW, determined from sequence) plotted against
the MW determined from the apex of the elution profile for proteins identified by mass
spectrometry in size-exclusion chromatography (SEC) fractions of flower tissue (sFL). The inset
shows the MW calibration of the SEC column using a protein calibration standard (MW between

44 and 690 kDa). The MW distribution of proteins annotated in Araport11 is shown on top of the scatterplot. It is apparent, that a large number of proteins show a much higher apparent MW than what would be expected from their sequences (data points above the x=y line). This implies that these proteins engage in physical protein interactions that are sufficiently stable during SEC separation.

b, SEC traces of proteins from five well characterized protein complexes for flower, leaf and root tissue. Even though the resolution of SEC separations is not very high, it is apparent that the complex subunits show very strong co-elution behaviour and that the SEC separations of the 5 complexes are reproducible between tissues. Acetyl-CoA carboxylase n = 4 proteins; Cell division control protein 48 (CDC48) n = 3 proteins; Ribulose bisphosphate carboxylase oxygenase (RubisCO) n = 4 proteins; Prefoldin n = 6 proteins; Succinate-CoA synthetase (SCS) n = 3 proteins.

c, Intensity normalized SEC elution profile of proteins for flower tissue. Proteins are ordered
based on the SEC fraction in which their intensity peaks and the data is displayed as heat map
(n = 2,485 protein traces). Co-eluting proteins were grouped into so-called trace modules (see
methods for details). Proteins in trace modules may represent members of protein complexes
and thus serve as candidates for further experimental validation.

1731 d. In order to quantify how well protein complexes can be detected using co-expression analysis 1732 from data in the tissue atlas (TA) or by SEC-MS a summary statistic termed 'complex index' was 1733 calculated (see methods). The complex index is 1, when all subunits of a complex are identified 1734 in the same module and no other proteins are contained in the module. The bar plots show examples for complex indices obtained from the different data sets and are divided into large 1735 1736 (>4 subunits) and small (<4 subunits) protein complexes (according to UniProt). The data shows that co-expression alone generates many candidates of interactors but that combining co-1737 1738 expression and SEC-MS analysis is an efficient way to prioritize candidates for follow-up 1739 experiments.

e, Subunit heterogeneity within the coatomer complex. The coatomer complex consists of 7 subunits, five of which (α , β , β ', ϵ and ζ) can be provided by 12 paralogs of these 5 genes. The plots show the protein proportions of these paralogs in all 30 tissues (data from tissue atlas). The data implies that the coatomer complex has a similar composition in most tissues. A notable exception are seed tissues in which subunit ζ -1 protein expression dominates over the two other paralogous proteins suggesting that the coatomer complex in seed tissue also preferentially contains the ζ -1 subunit. Tissues are coloured as in Figure 1.

f, Absolute SEC intensity traces of individual complex subunits for determining subunit
stoichiometry. Examples from left to right: the chaperonin complex (flower, 8 proteins, ratio of all
subunits: 1:1), the 26S proteasome core and lid (flower, 14+17 proteins, ratio of all subunits:
1:1), the COP9 Signalosome (flower, CSN; 8 proteins, ratio of all subunits: 1:1) and the

1751 CESA1/3/6 complex (root, 3 proteins, ratio of all subunits: 1:1). Note that CSN3 and CSN5 were 1752 detected both as part of the CSN complex and in monomeric form.

1753 g, Upper panels: total intensity of protein complex subunits across all tissues for the complexes shown in panel f (subunit intensities from the tissue atlas). Middle panels: Relative proportion of 1754 the complex subunits within each tissue. Lower panels: Estimation of subunit stoichiometry 1755 1756 using the average ± SD (n = 30 tissues) of the proportions of subunits across tissues (see methods). For the CESA complex, the ratios were calculated for the subunit combinations 1757 CESA1/3/6 and CESA4/7/8. The data shows that the stoichiometries determined from the tissue 1758 1759 expression data are generally well aligned with the expected 1:1 ratio of subunits in these complexes. As noted above, a substantial amount of CSN5 was detected as a monomer in the 1760 1761 SEC analysis and the tissue expression atlas also shows higher relative expression of this protein compared to all other complex partners. This suggests that this protein is produced in 1762 excess over what is required for the COP9 complex (as has been observed by others before) 1763 and may therefore imply an additional function within the cell. 1764

1765 Extended Data Figure 8 | Kinases, phosphatases and phosphorylation motifs.

a, Percentage of annotated kinases and phosphatases family members detected at the protein
or phosphoprotein level. Brackets denote the number of genes in each family in the Arabidopsis
genome

b, Tissue-resolved combined intensity (i.e. protein expression) of families of kinases (left) and
phosphatases (right). Tissues are coloured as in Figure 1. Note that several tissues (notably
pollen) stand out in terms of the expression of kinases and phosphatases implying that these
tissues are particularly active in phosphorylation-mediated dynamic signalling.

c, Upper panel: Pie chart of protein expression specificity categories (see methods for definition)
for kinases and phosphatases. Lower panels: Distribution of tissue-enhanced kinases and
phosphatases across the 30 tissues showing that several tissues (notably pollen) stand out in
terms of the expression of certain kinases and phosphatases implying tissue-specific signalling.

d, Pie charts showing the proportion of 'proline-directed', 'acidic', 'basic' and 'other' motif
categories for phosphorylated serine (pS), threonine (pT) and tyrosine (pY) residues. Only class
I p-sites (localization score > 0.75, methods) were considered in this analysis.

e, Example motif logo plots for 'proline-directed', 'acidic', 'basic' and 'other' motifs. P-site motifs were identified using the motif-X algorithm (see Supplementary Table 2 for all 266 motifs). n denotes the number of phosphorylation sites that contain the respective motif; fc denotes the fold change (i. e. enrichment) of the motif in phosphorylated vs unmodified peptides (see methods).

- 1785 f, Enrichment of 'proline-directed' (yellow), 'acidic' (red), 'basic' (blue) and 'other' (grey) 1786 sequence motifs (circles) in the serine p-site dataset versus the same motifs detected in the 1787 background dataset of unmodified peptides (see methods). Motifs are shown for 2, 3 and 4 fixed 1788 amino acid positions. The p-site in each motif example is underlined. X denotes any amino acid.
- g, Number of identified p-sites for a given protein plotted against the sequence lengths of the
 same protein. LATE EMBRYOGENESIS ABUNDANT (LEA) proteins are marked by a circle and
 arrow.
- h, Schematic representation of the LEA protein sequences (black bars). Pink marks denote
 phosphorylated and blue marks unphosphorylated STY residues. It is evident, that almost all
 STY residues in LEA proteins can be phosphorylated.
- i, Schematic representation the sequences and domain topology of the receptor-like kinases
 STRUBBELIG RECEPTOR FAMILY4 (SRF4), FERONIA (FER) and CHITIN RECEPTOR
 KINASE1 (CERK1). It is apparent that p-sites often preferentially occur in specific domains,
 notably the juxtamembrane domain. Protein sequence regions covered by identified peptides
 are marked in blue and p-sites are marked in pink.
- 1800 Extended Data Figure 9 | Functional analysis of phosphorylation mutants of RCAR10 and1801 QKY.
- a, Phosphorylation site (p-site) localization within the structure of REGULATORY
 COMPONENTS OF ABA RECEPTOR10 (RCAR10). The RCAR10 structure (blue) was
 modelled using the RCAR11 protein crystal structure (cornflowerblue) as a template ⁴⁵. ABA binding loops are shown in turquoise, p-sites in pink and ABA ligand in yellow.
- b, RCAR10 expression across tissues at the protein (blue, iBAQ), transcript (grey, TPM) and p-site (pink, intensity) level.
- c, Tissue-resolved total protein intensity and relative proportions of the members of the protein
 phosphatase 2C (PP2C) co-receptor family. Note that seed tissues stand out in terms of overall
 expression as well as the dominance of AHG1 in these tissues.
- d, Measurement of ABA response when expressing RCAR or RCAR10 phosphomimetic mutant variants in combination with different PP2C co-receptors in protoplasts (see methods). Columns display the average ABA response (\pm SD, n = 3) and grey dots indicate individual measurements. The data shows that co-expression of the phosphatases HAI1-3 lead to similar responses of the phophomimic mutants, whereas other co-expressed phosphatases show diverse responses.
- 1817 e, QUIRKY (QKY) expression across tissues at the protein (blue, iBAQ), transcript (grey, TPM)1818 and p-site (pink, intensity) level.

- 1819 f, Members of the MULTIPLE C2 AND TRANSMEMBRANE DOMAIN-CONTAINING PROTEIN 1820 (MCTP) family clustered by sequence similarity (left) and schematic representation of their 1821 domain structures along with detected p-sites (right). MCTP11a, 12 and 13 were not detected in 1822 this study (n.d.). QKY (MCTP15) is marked in bold.
- 1823 g, Number of independent transgenic plant lines (*qky*-9 mutant background) transformed with 1824 WT (QKY), phosphomutant (SA) or phosphomimic (SE) constructs that show complete, partial 1825 or no rescue of the mutant phenotype. qPCR results (average \pm SD; individual data points as 1826 grey dots) show the relative transgene expression in wild type (WT), *qky*-9 mutant and selected 1827 transgenic lines. QKY wild type (grey), QKY_{S262A} (SA; blue), QKY_{S262E} (SE, purple).
- 1828 h-j, Representative confocal images of six days-old *qky-9 pQKY::mCherry:QKY* (WT QKY; n =
- 1829 14 roots), qky-9 pQKY::mCherry:QKYS262A (phosphomutant; n = 21 roots), qky-9
- 1830 *pQKY::mCherry:QKYS262E* (phosphomimic; n = 15 roots) root epidermal cells of the
- 1831 meristematic zone. The punctate signal along the cell circumference shows the expected
- localization of QKY protein. Arrows indicate punctate structures. Scale bars 5 µm.

















Phosphomutant (QKY_{S262A})



Phosphomimic (QKY_{S262E})