

1

2 **Dynamic post-transcriptional regulation during embryonic stem cell** 3 **differentiation**

4

5 Patrick R. van den Berg¹, Bogdan Budnik², Nikolai Slavov^{3,*}, Stefan Semrau^{1,*†}

6

7

8 ¹ Leiden Institute of Physics, Leiden University, Leiden, Zuid-Holland, 2333 CC, The Netherlands

9 ² Mass Spectrometry and Proteomics Resource Laboratory, Harvard University, Cambridge,

10 Massachusetts, MA 02138, USA

11 ³ Department of Bioengineering, Northeastern University, Boston, Massachusetts, MA 02115,

12 * corresponding author, †, lead contact

13 **Summary**

14 During *in vitro* differentiation, pluripotent stem cells undergo extensive remodeling of their gene
15 expression profile. While studied extensively at the transcriptome level, much less is known about
16 protein dynamics. Here, we measured mRNA and protein levels of 7459 genes during differentiation of
17 embryonic stem cells (ESCs). This comprehensive data set revealed pervasive discordance between
18 mRNA and protein. The high temporal resolution of the data made it possible to determine protein
19 turnover rates genome-wide by fitting a kinetic model. This model further enabled us to systematically
20 identify dynamic post-transcriptional regulation. Moreover, we linked different modes of regulation to
21 the function of specific gene sets. Finally, we showed that the kinetic model can be applied to single-
22 cell transcriptomics data to predict protein levels in differentiated cell types. In conclusion, our
23 comprehensive data set, easily accessible through a web application, is a valuable resource for the
24 discovery of post-transcriptional regulation in ESC differentiation.

25 **Keywords**

26 embryonic stem cells, *in vitro* differentiation, gene regulation, transcriptomics, proteomics, kinetic
27 modeling, integration with single-cell transcriptomics, web application for data exploration

28 Introduction

29 Much of the medical potential of pluripotent stem cells is due to their ability to differentiate *in vitro* into
30 all tissue types of the adult body (Soldner and Jaenisch, 2012). While tremendous progress has been
31 made in guiding cells through successive lineage decisions, the gene regulatory mechanisms
32 underlying these decisions remain largely unknown. This gap in knowledge hampers the streamlining
33 and acceleration of differentiation protocols. A large body of work has focused on transcriptional
34 regulation, charting transcriptome changes during differentiation, most recently down to the single-cell
35 level (Klein et al., 2015; Loh et al., 2016; Semrau et al., 2016) These studies assumed implicitly that
36 mRNA levels are a good proxy for protein levels. Mounting evidence suggests that this is not a good
37 assumption for mammalian systems, where mRNA and protein levels were found to correlate only
38 moderately (Lu et al., 2009) (Kristensen et al., 2013; Peshkin et al., 2015; Schwanhäusser et al.,
39 2011). Where the discordance between protein and mRNA expression originates and what the
40 biological function might be are long-standing and controversially discussed issues (Liu et al., 2016;
41 Vogel and Marcotte, 2012). Here we study the relationship between mRNA and protein expression in
42 the context of *in vitro* differentiation, a highly dynamic process in which gene regulation at the protein
43 level likely plays an important role (Sampath et al., 2008).

44 Results

45 *Measurement of transcriptome and proteome dynamics during retinoic acid driven differentiation*

46 We used retinoic acid (RA) differentiation of mESCs as a generic model for *in vitro* differentiation.
47 Previously, we characterized this differentiation assay in detail at the transcriptional level by single-cell
48 RNA-seq (Semrau et al., 2016). In particular, we have shown that within 96 h of RA exposure, mESCs
49 bifurcate into an extraembryonic endoderm-like and an ectoderm-like cell type (XEN and ECT
50 respectively). Here we collected samples of the mixed population during an RA differentiation time
51 course as well as the two final, FACS-purified differentiated cell types at 96 h (Fig. 1a). For each time
52 point or cell type we quantified poly(A) RNA by RNA-seq and protein expression by tandem mass tag
53 (TMT) labeling followed by tandem mass spectrometry (MS/MS). In total, we obtained both RNA and
54 protein expression for 7459 genes (Supplementary Fig. 1a). Protein levels were quantified with low
55 technical error (Supplementary Fig. 1a) and high reproducibility between protein fold changes
56 measured in biological replicates (Pearson's $r = 0.92$, Supplementary Fig. 1b). Moreover, the XEN-like

57 cells measured here were similar to embryo derived XEN cells in their proteome (Mulvey et al., 2015)
58 ($r = 0.65$, Supplementary Fig. 1c).

59 *Correlation between mRNA and protein levels is moderate*

60 To explore the relationship between mRNA and protein levels we first correlated the two expression
61 levels across genes for individual time points or cell types (sample-wise correlation). In mESCs (0h
62 time point) Pearson correlation between mRNA and protein was 0.57 (Fig. 1b). Similar values have
63 been reported in other mammalian systems (de Sousa Abreu et al., 2009; Jovanovic et al., 2015;
64 Schwanhäusser et al., 2011). Sample-wise correlation was approximately the same for all samples,
65 including the purified differentiated cell types (Fig. 1c). Low mRNA-protein correlation was thus not cell
66 state dependent. Importantly, a low sample-wise correlation does not exclude the possibility that
67 relative changes in protein levels during differentiation closely follow relative changes in mRNA levels.
68 To quantify the concordance between mRNA and protein dynamics we calculated their correlation
69 across time for individual genes (gene-wise correlation, Fig. 1d-e). Some genes, like the pluripotency
70 factor *Rex1* (*Zfp42*) indeed exhibited a high correlation between mRNA and protein ($r = 0.93$ for
71 *Rex1*). Numerous genes, like the ribosomal protein *Rps6*, for example, did not exhibit any strong
72 correlation between protein or mRNA ($r = 0$ for *Rps6*). Strikingly, we also observed many genes with
73 anti-correlated profiles, like *Arpc1a* ($r = -0.91$) or *Arvcf* ($r = -0.90$). Such highly negative correlations
74 do not seem to be a result of technical noise in protein quantification, since multiple distinct peptides of
75 the same protein show similar trends (Supplementary Fig. 1d). Overall, the distribution of gene-wise
76 correlations, while peaking close to 1, had a long tail towards -1 (Fig. 1e). This result clearly shows
77 that mRNA dynamics are in general not a good predictor for protein dynamics during differentiation.

78 *Classification by dominant temporal trends visualizes widespread discordance between mRNA and* 79 *protein*

80 Having discovered that mRNA and protein dynamics are in general dissimilar we wanted to reveal the
81 main trends in expression dynamics and study how they differ between mRNA and protein. To that
82 end we used singular value decomposition (SVD) to decompose an expression profile into a weighted
83 sum of generic profiles, called eigengenes (Fig. 2a). In contrast to other classification methods, SVD
84 allows us to discriminate systematically between the main trend (the dominant eigengene) and
85 smaller, additional fluctuations (Fig. 2b). The first three eigengenes, which corresponded to monotonic,
86 transient or oscillatory trends, explained 76% and 85% of the variance in mRNA and protein
87 expression, respectively (Fig. 2c). mRNA eigengenes were more dynamic than protein eigengenes

88 (Supplementary Fig. 1e), which reflects the buffering of mRNA dynamics by protein synthesis and
89 degradation (Liu et al., 2016) (Jovanovic et al., 2015). Classification of all genes by their dominant
90 mRNA and protein eigengenes (which reflect the main temporal trends) revealed widespread
91 discordance (Fig. 2d). While there was a statistically significant enrichment of genes with similar
92 dominant mRNA and protein eigengenes (p -value $< 1E-5$), most genes (60%) had discordant mRNA
93 and protein dynamics.

94 *A simple kinetic model partially explains the mRNA-protein discordance for the majority of genes*

95 The temporal delay between mRNA and protein eigengenes (Fig. 2a) sparked the hypothesis that the
96 delay inherent to protein synthesis and degradation might cause much of the observed discordance.

97 To pursue this hypothesis we modeled protein turnover using a simple birth-death process with
98 constant protein synthesis and degradation rates (Tchourine et al., 2014) (Peshkin et al., 2015)
99 (Methods, Fig. 3a). In our model the synthesis rate k_s lumps all processes related to protein production
100 (translation initiation, elongation, etc.) while the degradation rate k_d represents all processes leading to
101 a reduction in protein levels (dilution due to cell division, active degradation, etc.). To avoid over-fitting,
102 we also considered simpler models, which correspond to cases in which a protein is only synthesized,
103 only degraded or completely constant (Fig. 3b). To select among these models, we employed the
104 Bayesian Information Criterion (BIC), a score that penalizes the fit according to the number of
105 parameters (Methods). To reveal whether there is a connection between a certain model and specific
106 molecular functions, we performed GO term enrichment analysis. This analysis revealed that the
107 “degradation only” model was enriched for genes with a role in blastocysts development and inner cell
108 mass proliferation (Supplementary Fig. 2a). These genes are likely involved in preserving the
109 pluripotent state, as exemplified by the pluripotency factor *Nanog*. Degradation of the corresponding
110 proteins is crucial for the timely exit from pluripotency. GO term enrichment analysis also showed that
111 the “synthesis only” model was enriched for genes involved in neuron development and mesenchymal
112 cell development. These genes thus likely have specific functions in differentiated cell types and hence
113 must be synthesized quickly to ensure proper function. An example of such a gene is *Lamb1*, which is
114 highly expressed in XEN cells. This analysis shows that the different regulatory modes identified by
115 our model correspond to specific functions in differentiation.

116 We next wanted to evaluate the validity of our model by comparison with relevant data sets from the
117 literature. Protein half-lives (Supplementary Fig. 2b) calculated from the degradation rates were in the
118 same range as previously reported values for other systems (Peshkin et al., 2015; Schwanhäusser et

119 al., 2011). Synthesis rates were positively correlated with translational efficiencies determined from
120 ribosome profiling in mESCs (Supplementary Fig. 2c) (Ingolia et al., 2011). The inferred kinetic rates
121 are thus biologically meaningful.

122 In order to assess how far our kinetic model can explain the observed protein-mRNA discordance we
123 calculated the correlation between measured and predicted protein levels (Fig. 3c). These correlations
124 were sharply peaked close to one, which means that our simple model is able to explain a large
125 portion of the observed mRNA-protein discordance. This discordance is likely only transient since
126 protein-to-mRNA ratios differed most from their equilibrium value ($k_{eq} = k_s/k_d$) in the beginning but
127 approached it over time (Supplementary Fig. 2d). This observation supports our conclusion that the
128 observed mRNA-protein discordance during differentiation is largely a transient, dynamic imbalance
129 caused by delayed protein synthesis and degradation.

130 *The CDS/ 3'UTR mRNA expression ratio is a modulator of the synthesis rate*

131 We next sought to further refine our kinetic model and explore whether we could find predictors of
132 protein abundance. In that respect we were intrigued by a recent report that connected the ratio of
133 mRNA expression from the coding sequence (CDS) and 3' untranslated region (UTR) to protein
134 abundance (Kocabas et al., 2015). In our data sets, the CDS/3'UTR mRNA expression ratio w also
135 had a non-trivial relationship with protein levels (Supplementary Fig. 2e). Consequently, we included w
136 in our model as a modulator of the synthesis rate (Fig. 3d, Methods). Again, using the BIC to
137 determine whether using an additional free parameter is warranted by the improvement of the fit, we
138 found that 492 genes were fit optimally by the extended kinetic model (Fig. 3e). In the cases where it
139 was optimal the extended model provided a substantial improvement over the basic model (Fig. 3f).
140 For roughly half of those genes, w has a positive effect on protein synthesis and a negative effect on
141 the other half (Supplementary Fig. 2f). While the molecular mechanism relating w to the protein
142 synthesis rate is not yet known, our analysis shows that w is an interesting predictor that should be
143 explored in future studies of protein dynamics.

144 *Failure of the kinetic model reveals dynamic post-transcriptional regulation*

145 Despite its success in explaining the mRNA-protein discordance overall, our kinetic model does not fit
146 the dynamics of all quantified proteins. We identified 1232 genes with a poor mRNA-protein correlation
147 that is not appreciably improved by any of the kinetic models (Supplementary Fig. 3a). Due to the
148 buffering of mRNA dynamics when synthesis and degradation rates are constant, the model fails in
149 particular when the protein profile is more dynamic than the mRNA profile (Supplementary Fig. 3b).

150 Importantly, the genes that are not fit well by our model are very similar to the full data set in their
151 protein reliabilities (medians: 0.970 versus 0.972) and measurement errors (median SEM: 0.121
152 versus 0.115). Hence, technical noise is in general not the reason for the lack of a good fit. Rather, the
153 model fails due to the assumption that kinetic rates are constant. Consequently, we consider genes
154 that are not fit well by the model to be dynamically regulated. We sought to find sets of such genes
155 that potentially share regulatory features. To this end we again used the classification by dominant
156 eigengenes (Supplementary Fig. 3c). As an example, we focused on a class of genes with relatively
157 simple dynamics: monotonically increasing mRNA and a transient increase in protein expression
158 (highlighted in Supplementary Fig. 3c). Notably, we discovered that genes belonging to the MAPK
159 pathway were enriched in this particular class (ConsensusPathDB, adjusted p-value = $1.8E-3$,
160 Supplementary Fig. 3d). This suggests that genes of the MAPK pathway, which is highly relevant for
161 the differentiation of mESCs (Kunath et al., 2007), are regulated dynamically at the protein level. This
162 analysis exemplifies that we can systematically identify sets of genes that are dynamically regulated at
163 the protein level, likely by common mechanisms.

164 *Sets of genes with different functions in differentiation show distinct regulatory modes*

165 We next wanted to concentrate further on the regulation of gene sets that are relevant for embryonic
166 stem cell differentiation. To that end, we defined sets of markers for the pluripotent state, XEN cells,
167 and ECT cells based on differential mRNA expression (Supplementary Fig. 4a), which were confirmed
168 by GO term enrichment (Supplementary Fig. 4b). As a fourth gene set we considered ribosomal
169 proteins since it has been shown previously that the translational state changes dramatically during
170 differentiation (Sampath et al., 2008). For these 4 gene sets we calculated the average mRNA and
171 protein profiles, correlation between mRNA and protein, classification by dominant eigengene and
172 inferred synthesis and degradation rates for the genes that are fit optimally by the full kinetic model
173 (Fig. 4a). This analysis of gene sets is also available on the companion website. Pluripotency markers
174 were in general down regulated at the mRNA level (per definition) but also at the protein level.
175 Correspondingly, we found this set to be enriched in the “degradation only” kinetic model while the
176 “synthesis only” model is underrepresented (Supplementary Fig. 4c). This observation is consistent
177 with the fact that pluripotency genes have to be down-regulated quickly to allow for a timely exit from
178 pluripotency. Nevertheless, there were some genes that showed a substantial increase in protein
179 expression and consequently had a negative correlation between measured mRNA and protein (see
180 Supplementary Fig. 4d for examples). XEN and ECT markers were in general upregulated, where ECT

181 markers came up before XEN markers, as shown by us previously (Semrau et al., 2016). In contrast to
182 the set of pluripotency markers, XEN and ECT genes showed a high level of concordance between
183 mRNA and protein, as immediately obvious from the eigengene classification. Correspondingly, both
184 gene sets were enriched for high correlation between mRNA and protein. Additionally, XEN markers
185 were enriched for the “synthesis only” model (Supplementary Fig. 4b). This might be related to the fact
186 that XEN cells have to produce high levels of extracellular matrix proteins (Mulvey et al., 2015), like
187 laminin (*Lamb1*) or collagen (*Col4a2*). Consequently, these proteins must be synthesized in a timely
188 manner to ensure the proper function of the XEN cells. All in all, it seems that cell type specific
189 markers defined at the mRNA level could be confirmed at the level of protein and that for these genes
190 protein expression closely follows mRNA expression. Compared to the gene sets discussed so far,
191 ribosomal protein (RP) genes showed a remarkable extent of discordance between mRNA and protein
192 expression. Eigengene classification revealed that many RP genes had protein profiles that were more
193 dynamic than their mRNA counterparts. Correspondingly, RP genes were enriched for low correlation
194 between mRNA and protein ($p\text{-value} = 3.3E\text{-}2$). As cells differentiated, the protein levels of RP genes
195 decreased, consistent with reduced cell division rates. The rate of decrease in abundance, however,
196 was RP specific. Thus, it will be interesting to isolate ribosomes and analyze the extent to which these
197 RP dynamics reflect ribosome remodeling and specialization (Slavov et al., 2015). In summary, we
198 have shown that the 4 analyzed gene sets follow distinct regulatory modes that can be related to
199 biological functions.

200 *The kinetic model can be applied to single-cell transcriptomics data to predict protein levels in*
201 *differentiated cell types*

202 In the experiment presented here, the existence of good antibodies for highly expressed surface
203 markers allowed us to purify differentiated cells at 96 h and profile their proteome. For earlier time
204 points or many other differentiation assays such an approach is difficult or even impossible. By
205 contrast, single-cell transcriptomics methods can be applied to any differentiation system. Hence, we
206 would like to use such data sets to predict protein levels in subpopulations. To that end, we extracted
207 cell type specific mRNA dynamics during differentiation from our earlier single-cell RNA-seq
208 measurement of the system (Semrau et al., 2016). We then applied our kinetic model to this data set
209 to predict protein levels in the differentiated cell types at 96 h (Fig. 4b, Methods). Our prediction was
210 clearly superior to a prediction that used only bulk RNA-seq measurements and protein-to-mRNA
211 ratios (Edfors et al., 2016) (Fig. 4c). We have thus demonstrated that our kinetic model with

212 parameters learned from bulk measurements can be applied to single-cell transcriptomics data to
213 predict cell type specific protein levels.

214 We finally compared the differentiated cell types directly with each other. Overall, the correlation
215 between mRNA and protein changes was poor and we identified a few outlier genes in particular that
216 showed extreme behavior (Fig. 4d). These outliers had comparable protein expression in XEN and
217 ECT cells (at most 2-fold difference) but mRNA expression was much lower in XEN cells (up to 19-
218 fold). Notably, these outliers are strongly enriched for imprinted genes (hypergeometric test, p -value =
219 $2.3E-10$). It is a well-known fact that some imprinted genes are mono-allelically expressed in extra-
220 embryonic tissues (Miri and Varmuza, 2009). Yet, the observed down-regulation goes well beyond a
221 two-fold change expected for mono-allelic expression. This observation demonstrates that our data set
222 can be used to discover significant differences in gene regulation between differentiated cell types.

223 Discussion

224 Here we systematically analyzed the dynamics of mRNA and protein expression during mESC
225 differentiation. We observed that absolute levels of protein and mRNA are only moderately correlated
226 in the steady (pluripotent) state, consistent with results in other mammalian systems (Schwanhäusser
227 et al., 2011) (Wilhelm et al., 2014) (Edfors et al., 2016). Importantly, low correlation does not
228 immediately imply a significant role of gene-specific regulation as technical noise tends to reduce the
229 observed correlation and conventional correction schemes typically ignore the effect of systematic,
230 correlated errors (Csárdi et al., 2015). Edfors et al. showed recently that the protein-to-mRNA ratio
231 (PTR) for a specific gene is constant across several tissues (Edfors et al., 2016). While the PTR might
232 allow the prediction of absolute protein levels, it is unable to capture relative changes over time or
233 relative differences between tissues (Franks et al., 2017; Silva and Vogel, 2016).

234 In this study we found widespread discordance between mRNA and protein dynamics during mESCs
235 differentiation. Such discordance has been observed recently in several systems, in particular:
236 *Xenopus* development (Peshkin et al., 2015), *C. elegans* development (Grün et al., 2014),
237 macrophage differentiation (Kristensen et al., 2013) and mESC differentiation (Lu et al., 2009). While
238 this discordance is typically interpreted as a sign of (post) translational regulation (Grün et al., 2014)
239 (Lu et al., 2009), theoretical work showed that a simple delay between mRNA and protein production
240 can lead to a reduction in gene-wise correlation (Gedeon and Bokes, 2012) (Munsky and Neuert,
241 2015). Here we showed here that a simple model with constant kinetic rates, substantially reduces the

242 discordance for 63% of discordant genes (Supplementary Fig. 3a). The same kinetic model explained
243 protein dynamics of a third of all genes during stress response in yeast (Tchourine et al., 2014) and of
244 75% of all genes in *Xenopus* development (Peshkin et al., 2015). Consistently, this simple model thus
245 explains discordance for significant proportions of the genome. We also found that the dynamics of
246 48% of all genes are best fit by a model that either includes only protein synthesis or degradation. A
247 similar observation was made analyzing the stress response in yeast (Tchourine et al., 2014). We
248 speculate that the different reduced models correspond to different regulatory mechanisms, as
249 suggested by the enrichment of different GO terms and gene sets reported here. We further showed
250 that protein-mRNA ratios were transiently out-of-steady-state on the way to a new equilibrium in the
251 differentiated cell types. The observed discordance between mRNA and protein thus most likely
252 reflects a transient, dynamic imbalance due to delayed protein synthesis and degradation. We further
253 extended the basic kinetic model by adding the CDS-3'UTR mRNA expression ratio as a useful new
254 predictor for the protein synthesis rate. We speculate that the underlying molecular mechanism is
255 related to a change in the abundances of mRNA isoforms, which are believed to have different
256 translation rates (Wong et al., 2016). Genes that were not fit well by the kinetic model, are by our
257 definition dynamically regulated at the protein level, as constant synthesis and degradation rates are
258 insufficient to describe the observed kinetics. This approach is complementary to the recently
259 developed PECA method that can be used to reveal regulatory events at the mRNA and protein level
260 (Cheng et al., 2016).

261 Our in-depth analysis of several gene sets revealed that cell type specific genes show a high
262 concordance between mRNA and protein dynamics, while for RP genes the correlation is much lower.
263 This result is reminiscent of a recent report that studied the stimulation of dendritic cells (Jovanovic et
264 al., 2015). Jovanovic et al. found that mRNA levels explain 90% of protein fold changes after
265 stimulation and proteins involved in the induced immune response were particularly enriched for this
266 regulatory mode. The dynamics of "housekeeping proteins" (including RPs), on the other hand, were
267 dominated by changes in protein synthesis and degradation rates. Similarly, Kristensen et al. reported
268 that mRNA abundance was the best predictor for proteins that were upregulated during differentiation
269 of monocytes to macrophage-like cells (Kristensen et al., 2013). Together with these previous reports
270 our study supports a model in which mRNA fold changes set the level of newly produced proteins that
271 have crucial, specific functions in the new cell state or cell type. Regulation at the level of protein
272 turnover, on the other hand, is used to adapt the existing proteome. Importantly, we also showed that
273 some pluripotency genes, defined as such by being down-regulated at the mRNA level, showed

274 increasing protein expression. This result cautions against defining markers for cell states or cell types
275 solely based on mRNA expression.

276 Finally, we applied our kinetic model, with model parameters learned in this study, to our earlier single-
277 cell transcriptomics measurement of RA differentiation. Our model successfully predicted the
278 proteomes of differentiated cell types that arise during RA differentiation. This approach thus makes it
279 possible to measure the proteomes of cell types that cannot be purified, for example due to the lack of
280 suitable antibodies.

281 In summary, this study provided the first in-depth, integrated analysis of mRNA and protein dynamics
282 during mESC differentiation. All measured data are provided in a convenient web application. We
283 hope that this application will facilitate future studies of specific gene sets or global relationships, for
284 example between sequence features and protein regulation (Vogel et al., 2010).

285 Author contributions

286 Conceptualization, S.S. and N.S.; Investigation, P. vd B., S.S., B.B. and N.S.; Resources, B.B.; Formal
287 analysis, P. vd B. and N.S.; Software, P. vd B.; Data curation, P. vd B. and N.S.; Writing – original
288 draft, S.S. and P. vd B.; Writing – review and editing, P. vd B., N.S. and S.S.; Supervision, S.S. and
289 N.S.

290 Acknowledgements

291 P. vd B. and S.S. were supported by the Netherlands Organisation for Scientific Research
292 (NWO/OCW), as part of the Frontiers of Nanoscience (NanoFront) program. Data analysis was carried
293 out on the Dutch national e-infrastructure with the support of SURF Foundation. N.S. was supported
294 by a New Innovator Award from the NIGMS of the NIH under Award number DP2GM123497. We
295 would like to thank Rudolf Jaenisch for reagents and valuable advice.

296

297 The authors declare no competing financial interests.

298 References

- 299 Cheng, Z., Teo, G., Krueger, S., Rock, T.M., Koh, H.W., Choi, H., Vogel, C., 2016. Differential
300 dynamics of the mammalian mRNA and protein expression response to misfolding stress. *Mol*
301 *Syst Biol* 12, 855–855. doi:10.15252/msb.20156423
- 302 Csárdi, G., Franks, A., Choi, D.S., Airoidi, E.M., Drummond, D.A., 2015. Accounting for Experimental
303 Noise Reveals That mRNA Levels, Amplified by Post-Transcriptional Processes, Largely
304 Determine Steady-State Protein Levels in Yeast. *PLoS Genet* 11, e1005206.
305 doi:10.1371/journal.pgen.1005206
- 306 de Sousa Abreu, R., Penalva, L.O., Marcotte, E.M., Vogel, C., 2009. Global signatures of protein and
307 mRNA expression levels. *Mol Biosyst* 5, 1512–1526. doi:10.1039/B908315D
- 308 Edfors, F., Danielsson, F., Hallström, B.M., 2016. Gene-specific correlation of RNA and protein levels
309 in human cells and tissues. *Molecular Systems* doi:10.15252/msb.20167325
- 310 Franks, A., Airoidi, E., Slavov, N., 2017. Post-transcriptional regulation across human tissues. *bioRxiv*
311 020206. doi:10.1101/020206
- 312 Gedeon, T., Bokes, P., 2012. Delayed Protein Synthesis Reduces the Correlation between mRNA
313 and Protein Fluctuations. *Biophys J* 103, 377–385.
- 314 Grün, D., Kirchner, M., Thierfelder, N., Stoeckius, M., Selbach, M., Rajewsky, N., 2014. Conservation
315 of mRNA and Protein Expression during Development of *C. elegans*. *Cell Reports* 6, 565–577.
316 doi:10.1016/j.celrep.2014.01.001
- 317 Ingolia, N.T., Lareau, L.F., Weissman, J.S., 2011. Ribosome Profiling of Mouse Embryonic Stem Cells
318 Reveals the Complexity and Dynamics of Mammalian Proteomes. *Cell* 147, 789–802.
319 doi:10.1016/j.cell.2011.10.002
- 320 Jovanovic, M., Rooney, M.S., Mertins, P., Przybylski, D., Chevrier, N., Satija, R., Rodriguez, E.H.,
321 Fields, A.P., Schwartz, S., Raychowdhury, R., Mumbach, M.R., Eisenhaure, T., Rabani, M.,
322 Gennert, D., Lu, D., Delorey, T., Weissman, J.S., Carr, S.A., Hacohen, N., Regev, A., 2015.
323 Immunogenetics. Dynamic profiling of the protein life cycle in response to pathogens. - PubMed -
324 NCBI. *Science* 347, 1259038–1259038. doi:10.1126/science.1259038
- 325 Klein, A.M., Mazutis, L., Akartuna, I., Tallapragada, N., Veres, A., Li, V., Peshkin, L., Weitz, D.A.,
326 Kirschner, M.W., 2015. Droplet Barcoding for Single-Cell Transcriptomics Applied to Embryonic
327 Stem Cells. *Cell* 161, 1187–1201. doi:10.1016/j.cell.2015.04.044
- 328 Kocabas, A., Duarte, T., Kumar, S., Hynes, M.A., 2015. Widespread Differential Expression of Coding

- 329 Region and 3' UTR Sequences in Neurons and Other Tissues. *Neuron* 88, 1149–1156.
330 doi:10.1016/j.neuron.2015.10.048
- 331 Kristensen, A.R., Gsponer, J., Foster, L.J., 2013. Protein synthesis rate is the predominant regulator
332 of protein expression during differentiation. *Mol Syst Biol* 9, 689–689. doi:10.1038/msb.2013.47
- 333 Kunath, T., Saba-El-Leil, M.K., Almousaillekh, M., Wray, J., Meloche, S., Smith, A., 2007. FGF
334 stimulation of the Erk1/2 signalling cascade triggers transition of pluripotent embryonic stem cells
335 from self-renewal to lineage commitment. *Development* 134, 2895–2902. doi:10.1242/dev.02880
- 336 Liu, Y., Beyer, A., Aebersold, R., 2016. On the Dependency of Cellular Protein Levels on mRNA
337 Abundance. *Cell* 165, 535–550. doi:10.1016/j.cell.2016.03.014
- 338 Loh, K.M., Chen, A., Koh, P.W., Deng, T.Z., Sinha, R., Tsai, J.M., Barkal, A.A., Shen, K.Y., Jain, R.,
339 Morganti, R.M., Shyh-Chang, N., Fernhoff, N.B., George, B.M., Wernig, G., Salomon, R.E.A.,
340 Chen, Z., Vogel, H., Epstein, J.A., Kundaje, A., Talbot, W.S., Beachy, P.A., Ang, L.T., Weissman,
341 I.L., 2016. Mapping the Pairwise Choices Leading from Pluripotency to Human Bone, Heart, and
342 Other Mesoderm Cell Types. *Cell* 166, 451–467. doi:10.1016/j.cell.2016.06.011
- 343 Lu, R., Markowetz, F., Unwin, R.D., Leek, J.T., Airoidi, E.M., MacArthur, B.D., Lachmann, A., Rozov,
344 R., Ma'ayan, A., Boyer, L.A., Troyanskaya, O.G., Whetton, A.D., Lemischka, I.R., 2009. Systems-
345 level dynamic analyses of fate change in murine embryonic stem cells. *Nature* 462, 358–362.
346 doi:10.1038/nature08575
- 347 Miri, K., Varmuza, S., 2009. Chapter 5 Imprinting and Extraembryonic Tissues—Mom Takes Control,
348 in: *International Review of Cell and Molecular Biology*. Elsevier, pp. 215–262.
349 doi:10.1016/S1937-6448(09)76005-8
- 350 Mulvey, C.M., Schröter, C., Gatto, L., Dikicioglu, D., Fidaner, I.B., Christoforou, A., Deery, M.J., Cho,
351 L.T.Y., Niakan, K.K., Martinez Arias, A., Lilley, K.S., 2015. Dynamic Proteomic Profiling of Extra-
352 Embryonic Endoderm Differentiation in Mouse Embryonic Stem Cells. *STEM CELLS* 33, 2712–
353 2725. doi:10.1002/stem.2067
- 354 Munsky, B., Neuert, G., 2015. From analog to digital models of gene regulation. *Phys. Biol.* 12,
355 045004. doi:10.1088/1478-3975/12/4/045004
- 356 Peshkin, L., Wühr, M., Pearl, E., Haas, W., Freeman, R.M., Gerhart, J.C., Klein, A.M., Horb, M., Gygi,
357 S.P., Kirschner, M.W., 2015. On the Relationship of Protein and mRNA Dynamics in Vertebrate
358 Embryonic Development. *Developmental cell* 35, 383–394. doi:10.1016/j.devcel.2015.10.010
- 359 Sampath, P., Pritchard, D., Pabon, L., Reinecke, H., Schwartz, S., Morris, D., Murry, C., 2008. A

360 hierarchical network controls protein translation during murine embryonic stem cell self-renewal
361 and differentiation. *Cell Stem Cell* 2, 448–460.

362 Schwanhäusser, B., Busse, D., Li, N., Dittmar, G., Schuchhardt, J., Wolf, J., Chen, W., Selbach, M.,
363 2011. Global quantification of mammalian gene expression control. *Nature* 473, 337–342.
364 doi:10.1038/nature10098

365 Semrau, S., Goldmann, J., Soumillon, M., Mikkelsen, T.S., Jaenisch, R., van Oudenaarden, A., 2016.
366 Dynamics of lineage commitment revealed by single-cell transcriptomics of differentiating
367 embryonic stem cells. *bioRxiv* 068288. doi:10.1101/068288

368 Silva, G.M., Vogel, C., 2016. Quantifying gene expression: the importance of being subtle. *Mol Syst*
369 *Biol.* doi:10.15252/msb.20167144

370 Slavov, N., Semrau, S., Airoidi, E., Budnik, B., van Oudenaarden, A., 2015. Differential Stoichiometry
371 among Core Ribosomal Proteins. *Cell Reports* 13, 865–873. doi:10.1016/j.celrep.2015.09.056

372 Soldner, F., Jaenisch, R., 2012. iPSC Disease Modeling. *Science* 338, 1155–1156.
373 doi:10.1126/science.1227682

374 Storey, J.D., 2005. Significance analysis of time course microarray experiments. *Proceedings of the*
375 *National Academy of Sciences* 102, 12837–12842. doi:10.1073/pnas.0504609102

376 Tchourine, K., Poultney, C.S., Wang, L., Silva, G.M., Manohar, S., Mueller, C.L., Bonneau, R., Vogel,
377 C., 2014. One third of dynamic protein expression profiles can be predicted by a simple rate
378 equation. *Mol Biosyst* 10, 2850–2862. doi:10.1039/C4MB00358F

379 Vogel, C., de Sousa Abreu, R., Ko, D., Le, S.Y., Shapiro, B.A., Burns, S.C., Sandhu, D., Boutz, D.R.,
380 Marcotte, E.M., Penalva, L.O., 2010. Sequence signatures and mRNA concentration can explain
381 two-thirds of protein abundance variation in a human cell line. *Mol Syst Biol* 6, 400.
382 doi:10.1038/msb.2010.59

383 Vogel, C., Marcotte, E.M., 2012. Insights into the regulation of protein abundance from proteomic and
384 transcriptomic analyses. *Nat Rev Genet.* doi:10.1038/nrg3185

385 Wilhelm, M., Schlegl, J., Hahne, H., Gholami, A.M., Lieberenz, M., Savitski, M.M., Ziegler, E.,
386 Butzmann, L., Gessulat, S., Marx, H., Mathieson, T., Lemeer, S., Schnatbaum, K., Reimer, U.,
387 Wenschuh, H., Mollenhauer, M., Slotta-Huspenina, J., Boese, J.-H., Bantscheff, M., Gerstmair,
388 A., Faerber, F., Kuster, B., 2014. Mass-spectrometry-based draft of the human proteome. *Nature*
389 509, 582–587. doi:10.1038/nature13319

390 Wong, Q.W.-L., Vaz, C., Lee, Q.Y., Zhao, T.Y., Luo, R., Archer, S.K., Preiss, T., Tanavde, V., Vardy,

391 L.A., 2016. Embryonic Stem Cells Exhibit mRNA Isoform Specific Translational Regulation. PLoS
392 ONE 11, e0143235. doi:10.1371/journal.pone.0143235
393 Ying, Q.-L., Wray, J., Nichols, J., Battle-Morera, L., Doble, B., Woodgett, J., Cohen, P., Smith, A.,
394 2008. The ground state of embryonic stem cell self-renewal. Nature 453, 519–523.
395 doi:10.1038/nature06968
396

397 Figure captions

398 *Figure 1 mRNA and protein expression correlate poorly during mESC differentiation*

399 (A) Experimental setup. (B) mRNA versus protein expression of 7459 genes in mESCs. Each data
400 point is an individual gene. Red lines indicate contour lines of equal density. (C) Sample-wise Pearson
401 correlation between mRNA and protein for all samples. The solid line indicates the average of all time
402 course samples. The grey area indicates the 5% rejection region for all samples being identical (see
403 Methods). Error bars: SEM. (D) mRNA versus protein expression at all time points for nine example
404 genes. Pearson's correlation r is indicated for each gene. The line and grey area indicate the linear
405 regression fit and 95% CI, respectively. Error bars: SEM. (E) Distribution of the gene-wise Pearson
406 correlation between mRNA and protein. Numbered arrows indicate the position of the examples shown
407 in D. See also Supplementary Figure 1.

408 *Figure 2 Classification of temporal mRNA and protein expression profiles by dominant trends reveals* 409 *widespread discordance*

410 (A) First six eigengenes of mRNA and protein expression profiles. (B) Reconstruction of mRNA and
411 protein expression profiles from the top three eigengenes of an example gene. (C) Cumulative
412 variance explained by the eigengenes for mRNA and protein profiles. (D) Classification of all genes by
413 their dominant mRNA eigengene (columns) and protein eigengene (rows).
414 See also Supplementary Figure 1.

415

416 *Figure 3 Simple kinetic models of protein synthesis and degradation explain mRNA-protein* 417 *discordance.*

418 (A) Kinetic model. k_s = synthesis rate constant; k_d = degradation rate constant. (B) Example fits of the
419 *full model* ($k_s > 0$, $k_d > 0$) and the three reduced models: *synthesis only* ($k_s > 0$, $k_d = 0$), *degradation only*
420 ($k_s = 0$, $k_d > 0$) and *degenerate* ($k_s = k_d = 0$). Percentages indicate the fraction of genes fit best by the
421 respective model. (C) Distribution of Pearson correlation between measured protein expression and
422 mRNA expression or predicted protein expression. (D) Extended kinetic model. $k_s(t)$ = time-dependent
423 synthesis rate. (E) mRNA expression, log ratio of expression from CDS and 3'UTR and protein
424 expression profiles of two example genes with fits of the extended model (solid line) or the basic
425 model (dashed line). (F) Distribution of Pearson correlation between measured protein expression

426 and: mRNA expression, protein expression predicted by the basic model or the extended model. Error
427 bars in (B) and (E): SEM.

428 See also Supplementary Figures 2 and 3.

429

430 *Figure 4. Classification and kinetic modelling reveal differences between gene sets involved in*
431 *differentiation and between differentiated cell types.*

432 (A) Comparison of four gene sets that are relevant for differentiation. Log₂ fold change (L2FC) of
433 mRNA and protein expression are shown for individual genes (colored) and the set average (black).
434 The p-value in the classification matrix is based on picking genes at random from all genes (chi-
435 squared test). (B) mRNA expression of XEN and ECT subpopulations (from single cell data) and the
436 mixed populations (bulk sample). Protein expression in XEN and ECT is predicted by applying the
437 kinetic model to the single cell data. Alternatively, at 96 h we also predicted protein based on the
438 protein-to-mRNA (PTR) ratio. MPI = Mean peptide intensity. (C) Sum of squared residuals (SSR) of
439 the kinetic model-based prediction compared to the PTR-based prediction for the XEN and ECT
440 marker genes. (D) mRNA and protein expression in XEN cells relative to ECT cells. Outlier genes are
441 highlighted with a dark background and imprinted genes are shown in red (obtained from
442 www.geneimprint.com, Oct-11-2016). Imprinted genes are significantly enriched in the outlier gene set
443 (hypergeometric test: p-value = 2.72e-10).

444 See also Supplementary Figure 4.

445 **Methods**

446 **Cell culture**

447 E14 mouse embryonic stem cells were cultured as previously described (Semrau et al., 2016). Briefly,
448 cells were grown in modified 2i medium (Ying et al., 2008): DMEM/F12 (Life technologies)
449 supplemented with 0.5x N2 supplement, 0.5x B27 supplement, 4mM L- glutamine (Gibco), 20 µg/ml
450 human insulin (Sigma-Aldrich), 1x 100U/ml penicillin/streptomycin (Gibco), 1x MEM Non-Essential
451 Amino Acids (Gibco), 7 µl 2-Mercaptoethanol (Sigma-Aldrich), 1 µM MEK inhibitor
452 (PD0325901, Stemgent), 3 µM GSK3 inhibitor (CHIR99021, Stemgent), 1000 U/ml mouse LIF
453 (ESGRO). Cells were passaged every other day with Accutase (Life technologies) and replated on
454 gelatin coated tissue culture plates (Cellstar, Greiner bio-one).

455 **Differentiation and sample collection**

456 Retinoic acid induced differentiation was carried out exactly as describe before (Semrau et al., 2016).
457 Prior to differentiation cells were grown in 2i medium for at least 2 passages. Cells were seeded at 2.5
458 $\times 10^5$ per 10 cm dish and grown over night (12 h). Cells were then washed twice with PBS and
459 differentiated in basal N2B27 medium (2i medium without the inhibitors, LIF and the additional insulin)
460 supplemented with 0.25 µM all-trans retinoic acid (RA, Sigma-Aldrich). Spent medium was exchanged
461 with fresh medium after 48 h.
462 To collect samples, cells were dissociated with Accutase. RNA was extracted from half of the sample
463 (RNeasy, Qiagen) and the purified RNA was stored at -80C until RNA-sequencing was performed. The
464 other half of the sample was flash frozen in liquid nitrogen and stored at -80C until mass spectrometry
465 was performed.

466 **Fluorescence-activated cell sorting**

467 FACS sorting of the differentiated cell types and quantification of the cell type frequencies was carried
468 out exactly as described previously (Semrau et al., 2016).

469 **RNA sequencing and mRNA quantification**

470 *Library preparation and RNA sequencing*

471 The libraries for RNA sequencing were prepared under standard conditions using Illumina's TruSeq
472 stranded mRNA sample preparation kit. The libraries were sequenced using Illumina HiSeq 3000 ; 40

473 basepair long, stranded single-end reads were sequenced at an average read depth of 40 million
474 reads per sample. The data is available through GEO.

475 *Read alignment*

476 An RSEM-reference was created using RSEM v1.2.28 (Li and Dewey, 2011) with the Illumina
477 iGenome GRCh38 reference using the standard settings. Next, the Illumina adapter was trimmed
478 from the reads with *cutadapt* v1.8.3 (Martin, 2011) and low quality bases with *sickle* v1.33 (Joshi et al.,
479 2011). Finally the reads were aligned with RSEM v1.2.28 (Li and Dewey, 2011) and Bowtie 2 v2.2.6
480 (Langmead and Salzberg, 2012) using standard settings except for “*--sampling-for-bam --fragment-*
481 *length-mean 40*”. The option “*--sampling-for-bam*” was applied so each read appears in the BAM file
482 once. This enabled the estimation of the CDS and 3’UTR counts by *summarizeOverlaps* from the
483 package *GenomicAlignment* v1.8.4 (Lawrence et al., 2013).

484 *Gene quantification*

485 mRNA expression was quantified by several different methods depending on the application.
486 Transcripts per million (TPM) was calculated by *RSEM* and was used when comparing between genes
487 since it is corrected for gene length. The more variance stabilized regularized log counts (rLC) were
488 determined by applying the *rlog* function from *DESeq2* v1.12.3 (Love et al., 2014) on rounded
489 expected counts obtained from *RSEM*. From this regularized counts (rC) were obtained by: $rC = 2^{rLC}$.
490 rLC and rC are corrected for overdispersion in low-read genes and are therefore used when
491 comparing one gene across multiple samples. CDS and 3’UTR counts were determined by splitting
492 the gene annotation file (GTF) with the *GenomicFeatures* package v1.26.0 (Lawrence et al., 2013) into
493 CDS and 3’UTR for every Ensembl gene ID. Next, the number of reads on the CDS and 3’UTR
494 features from the aligned BAM files were counted with *summarizeOverlaps* with default options.
495 “Union”, the default option for *mode*, discards reads, if they overlap with both CDS and 3’UTR. The
496 ratio w (CDS / 3’UTR) was only calculated for genes with at least 10 reads for CDS and 3’UTR in
497 every sample.

498 *Differentially expressed genes*

499 Differentially expressed genes (DEGs) were determined by *DESeq2* v1.12.3 (Love et al., 2014) on the
500 rounded expected counts obtained from RSEM at a false discovery rate (FDR) of 10%. The gene set
501 ‘pluripotency genes’ were DEGs that were down-regulated when comparing the samples 0h (n=2) and
502 96h (n=2). XEN- and ECT-marker gene sets were DEGs that were up-regulated when comparing the

503 samples 0h (n=2) with XEN (n=1) or ECT (n=1) respectively. Additionally, XEN- and ECT-markers
504 have at least a 2-fold difference in expression between the two cell types.

505 **Mass spectrometry and protein quantification**

506 *Sample preparation*

507 Pelleted cells were lysed in 400 μ l RIPA buffer, except for the sorted cells, which were lysed in 200 μ l
508 RIPA buffer. Volumes of cell lysate corresponding to 100 μ g protein per sample were digested with
509 trypsin using a modified FASP protocol (Wiśniewski et al., 2009). Subsequently each sample was
510 labeled with TMT 10-plex reagent (Prod# 90061, Thermo Fisher, San Jose, CA) according to the
511 manufacturer's protocol. All labeled samples were combined into a set-sample.

512 *Mass spectrometry*

513 The labeled set-sample was fractionated by electrostatic repulsion-hydrophilic interaction
514 chromatography chromatography (ERLIC) run on an HPLC 1200 Agilent system using PolyWAX LP
515 column (200x2.1 mm, 5 μ m, 30nm, PolyLC Inc, Columbia, MD) and a fraction collector (Agilent
516 Technologies, Santa Clara, CA). Set-samples were fractionated into a total of 40 ERLIC fractions.
517 Each ERLIC fraction was subsequently further separated by online nano-LC and submitted for tandem
518 mass spectrometry analysis to both LTQ OrbitrapElite or Q exactive high field (HF). One third of each
519 fraction was injected from an auto-sampler into the trapping column (75 μ m column ID, 5 cm length
520 packed with 5 μ m beads with 20 nm pores, from Michrom Bioresources, Inc.) and washed for 15 min;
521 the sample was eluted to analytic column with a gradient from 2 to 32 % of buffer B (0.1 % formic acid
522 in ACN) over 180 min gradient and fed into LTQ OrbitrapElite or Q exactive HF. The instruments were
523 set to run in TOP 20 MS/MS mode method with dynamic exclusion. After MS1 scan in Orbitrap with
524 60K resolving power, each ion was submitted to an HCD MS/MS with 60K resolving power and to CID
525 MS/MS scan subsequently. All quantification data were derived from HCD spectra.

526 *Protein quantification*

527 Relative peptide levels were estimated from reporter ion intensities measured at MS2 level. Only
528 peptides with co-isolation below 40 % were used for quantification. The intensities of all peptides
529 belonging to a Uniprot ID were averaged to form mean peptide intensity (MPI) for every protein. When
530 comparing different protein samples mean peptide intensities were normalized to the sample-mean to
531 form protein expression. Standard error of the mean (SEM) was calculated for every protein as

532 follows: 1) for every peptide the intensities were averaged across the samples, 2) the SEM was
533 calculated from these mean-centered peptide intensities for every protein and sample.

534 *Protein reliability*

535 The protein reliability was calculated for genes with at least two peptides quantified. For each gene,
536 the peptides were randomly split into two groups and the MPI was calculated for each group as
537 described above. The correlation between the MPIs of the two peptide groups across the different
538 samples is defined as the reliability of the measurement of that protein.

539 **Transcriptomics and proteomics integration**

540 While transcripts were identified by Ensembl gene IDs, Uniprot IDs were used for proteins. To
541 integrate the two, we mapped 7681 out of 8515 Uniprot IDs to Ensembl gene IDs present in the RNA-
542 seq data using the *idmapping* file from the Uniprot website (15-Sept-2016). An additional set of Uniprot
543 IDs were mapped to Ensembl IDs using *biomaRt* v2.28.0 (Durinck et al., 2009). Some proteins have
544 more than one Ensembl ID mapping to it, therefore 33 Uniprot IDs were removed. Moreover, 92
545 Uniprot IDs mapped non-uniquely to Ensembl IDs and for these the protein intensities were
546 reevaluated based on Ensembl IDs. Finally, some genes were not considered because they were not
547 detected in all samples. This resulted in a total of 7489 genes based on Ensembl gene IDs, for which
548 we have matched mRNA and protein expression data in all samples. Additionally, we observed 3770
549 genes with at least 10 mRNA reads in every sample but no detected protein.

550 *Sample-wise correlation*

551 We tested if the sample-wise correlation is constant during the differentiation time course using a
552 resampling approach. For each bootstrap a *pseudo-sample* was constructed consisting of every gene,
553 but with mRNA and protein expression randomly sampled from the different time points. The
554 correlations of 10,000 *pseudo-samples* were calculated to obtain a null distribution. Samples have
555 significantly different correlation if it falls below or above the 0.36 and 99.64 percentiles of the null
556 distribution respectively ($\alpha = 0.05$, Bonferroni correction, grey area in Figure 1c).

557 *Gene-wise correlation*

558 To define a threshold for low gene-wise correlation we applied a shuffling approach (Tchourine et al.,
559 2014). We determined the Pearson correlation for all possible permutations of the mRNA and protein
560 expression for every gene. More than 95% of all Pearson correlation values obtained in this way were
561 lower than 0.7, which we therefore set as the threshold between low and high correlation.

562 *Expression profile classification*

563 mRNA and protein expression were arranged in matrix form rows corresponding to genes and the
564 columns corresponding to time course samples. These matrices were standardized by rows. Next,
565 standard singular value decomposition (SVD) was performed separately for mRNA and protein (Wall
566 et al., 2003). From this analysis, we obtain n eigengenes \vec{V}_k where $k \in 1, \dots, n$ and n is the number of
567 time points. Using these eigengenes we can reconstruct the standardized expression of gene i , as
568 follows: $\vec{X}_i = \sum_k M_{ik} \vec{V}_k$, where M_{ik} is the contribution of eigengene k to the standardized expression of
569 gene i . We defined the eigengene with the biggest contribution to \vec{X}_i as the dominant eigengene. To
570 determine if there is an enrichment of genes with concordant mRNA and protein eigengenes, we
571 calculated an empirical p-value based on a null distribution generated by bootstrapped (number of
572 bootstraps = 100,000). This null distribution was constructed under the assumption that the marginal
573 eigengene distributions of mRNA and protein are independent. Moreover, we defined a confident set
574 of genes with a bigger than median fold-change between the contribution of the dominant eigengene
575 and the second most contributing eigengene for both mRNA and protein.

576 **Kinetic models of protein synthesis and degradation**

577 *Approximation of mRNA and CDS/3'UTR expression by natural cubic splines*

578 To describe the mRNA, CDS and 3'UTR behavior in the kinetic model of protein synthesis and
579 degradation we approximated the expression with natural cubic splines. These splines were fit on the
580 mRNA expression and on the \log_2 fold change (L2FC) of w , which we call ω . The number of degrees
581 of freedom p used for the fits of every gene was 4 for mRNA expression and 3 for ω expression.
582 These values were automatically determined as described by Storey *et al.* (Storey, 2005). Briefly, an
583 SVD was performed on the expression matrices of mRNA and ω and the first n eigengenes that
584 explain at least 60% of the variance were selected. For each of these eigengenes the optimal number
585 of degrees of freedom p_i was selected by leave one out cross validation (LOOCV) and the largest p_i
586 was used as the number of degrees of freedom p to fit the natural cubic splines for all the genes of the
587 expression matrix. The nodes of the cubic splines were equally spaced across the time course.

588 *Kinetic rate parameters estimation*

589 We model protein turnover as a birth-death process

$$590 \quad \frac{dP(t)}{dt} = k_s \cdot R(t) - k_d \cdot P(t)$$

591 where $P(t)$ and $R(t)$ are protein and mRNA expression respectively. The solution of this ordinary
592 differential equation (ODE) is given by:

604
$$P(t) = P_0 e^{-k_d t} + k_s \int_0^t R(\tau) e^{-k_d \cdot (t - \tau)} d\tau$$

593 where P_0 is the protein expression at $t = 0$ hours. The integral of this equation was estimated
594 numerically in R using the spline fits described above. We fit the model using gene specific
595 parameters P_0 , k_s and k_d with the Levenberg – Marquardt non-linear least squares algorithm, which is
596 implemented in the R package *minpack.lm* v1.2-0. Additionally, we fit models where we set $k_d = 0$,
597 $k_s = 0$ or $k_d = k_s = 0$. For each successful fit we determined the Bayesian Information Criterion:

605
$$BIC = -2 \ln(\hat{L}) + k \cdot \ln(n)$$

598 where \hat{L} is the posterior likelihood of the fit, k is number of parameters in the model and n is the
599 number of time points. \hat{L} is determined by:

606
$$\hat{L} = \prod_{j=1}^n p(P(t_j) | \hat{\theta})$$

600 where $\hat{\theta}$ is the vector of inferred model parameters. The probabilities are estimated by assuming a
601 normal distribution around the observed protein expression with a standard deviation equal to the
602 SEM of the peptide intensities. The kinetic model with the lowest BIC was selected as the optimal
603 model.

607 Additionally, for the subset of genes for which we could determine ω we constructed a model with a
608 time-dependent synthesis rate:

612
$$\frac{dP(t)}{dt} = k_s(t) \cdot R(t) - k_d \cdot P(t) = \kappa_s (1 + \beta \omega(t)) \cdot R(t) - k_d \cdot P(t)$$

609 where κ_s describes the constant synthesis rate and β parameterizes the time-dependent modulation of
610 the synthesis rate by ω . The solution of this ODE:

613
$$P(t) = P_0 e^{-k_d t} + \int_0^t ((\kappa_s + \beta \omega(\tau)) \cdot R(\tau) \cdot e^{-k_d \cdot (t - \tau)}) d\tau$$

611 was fit to the data in the same manner as above.

614 95% confidence region estimation

615 To estimate the 95% confidence intervals (CIs) for k_s and k_d we applied Wilk's theorem:

616
$$\ln(L(\theta)) \geq \ln(L(\hat{\theta})) - \frac{1}{2} \chi_{1,1-\alpha}^2$$

617 where α is 0.05 and $\chi_{1,1-\alpha}^2$ is the value at which the cumulative chi-squared distribution with 1 degree of
618 freedom reaches 0.95. We varied k_s and k_d around the obtained fit $\hat{\theta}$ to find the edges where Wilk's
619 theorem holds. These edges were determined at 24 directions in the $k_s - k_d$ solution plane to obtain a
620 crude 95% confidence region. The projection of this region on k_s and k_d defined $CI_{k_s}^{95\%}$ and $CI_{k_d}^{95\%}$, their
621 respective 95% CIs. Note that these intervals are typically much larger than the intervals obtained
622 when searching one parameter at a time. Genes with the full model (as determined by BIC), and with a
623 small $CI_{k_s}^{95\%}$ and $CI_{k_d}^{95\%}$ (each spanning less than a 10-fold range) were defined as the high-confidence
624 gene set. Additionally, for genes in this set we determined the protein half-life τ_p as

$$625 \quad \tau_p = \frac{\ln 2}{k_d}$$

626 *Protein prediction of sorted populations*

627 We applied our kinetic model to single-cell transcriptomics data of RA driven differentiation, which we
628 obtained previously (Semrau et al., 2016). We determined the mean expression of all cells, as well as
629 XEN and ECT subpopulations starting from the lineage bifurcation at 36 h. All three datasets thus
630 have identical expression up to 36 h. We then scaled the subpopulation data to the bulk data
631 measured here for every gene in the following way: 1) We standardized the single cell time course
632 data using the mean and standard deviation of the pooled single cell data, and 2) we scaled the
633 standardized single cell data to the bulk data using the mean and standard deviations of the bulk time
634 course. Next, we fit a natural cubic spline to the single cell data as before and applied the kinetic
635 model using P_0 , k_s and k_d learned from the bulk mRNA and protein measurements. We evaluated the
636 model performance by calculating the residuals between the predicted XEN and ECT protein
637 expression at 96 h and the bulk measurements of protein in the purified cell types.

638 An alternative way of predicting protein expression is by simply multiplying a gene's protein-to-mRNA
639 ratio (PTR) with the gene's mRNA expression. We defined the PTR as the mean protein expression
640 divided by the mean mRNA expression during the time course. We predicted the protein expression of
641 the XEN and ECT populations at 96 h using the bulk mRNA of the respective sorted populations. We
642 used the sorted bulk data rather than the single cell data, because it is more accurate and we
643 therefore expect this to perform better. Like with the single cell predictions, we evaluated model
644 performance using the residuals of the PTR-predictions relative to the measured protein expression of
645 the sorted bulk data.

646 **Ribosomal protein gene list**

647 The list of RPs was compiled as all Swiss-Prot proteins curated as ribosomal proteins in their
648 descriptions.

649 **Eigengene dynamics**

650 We quantified the dynamics of the eigengene profiles as the mean of the squared second derivatives
651 (roughness). The second derivatives were estimated numerically from three unequally spaced points
652 by this formula:

653
$$\frac{d^2y}{dx^2} = \frac{2y_1}{(x_2 - x_1)(x_3 - x_1)} - \frac{2y_2}{(x_3 - x_2)(x_2 - x_1)} + \frac{2y_3}{(x_3 - x_2)(x_3 - x_1)}$$

654 where x_1 , x_2 and x_3 are adjacent time points and y_1 , y_2 and y_3 are the respective eigengene
655 intensities.

656 **GO term enrichment**

657 GO term enrichment was performed with the R package *topGO* v2.24.0 (Alexa et al., 2006) with the
658 *classic* algorithm. The genes were ranked using Fisher's exact test and deemed significant with an
659 FDR of 10%.

660

661 **Accession numbers**

662 The RNA-seq data has been deposited in GEO (ID: GSE9563). The raw MS data has been deposited
663 in MassIVE (ID: MSV000080461). A web application complementing this publication, which allows
664 convenient access to all data can be found here:

665 <https://home.physics.leidenuniv.nl/~semrau/proteomics/>

666 user name: upon request

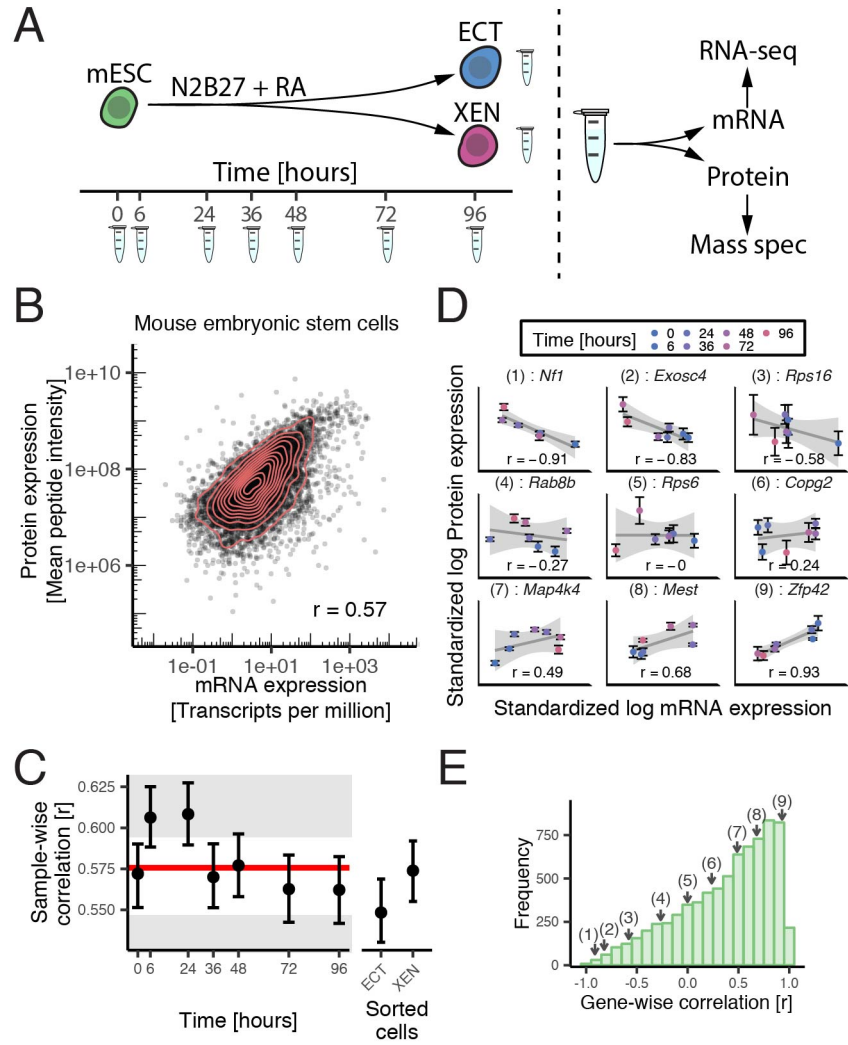
667 password: upon request

668 Additional references

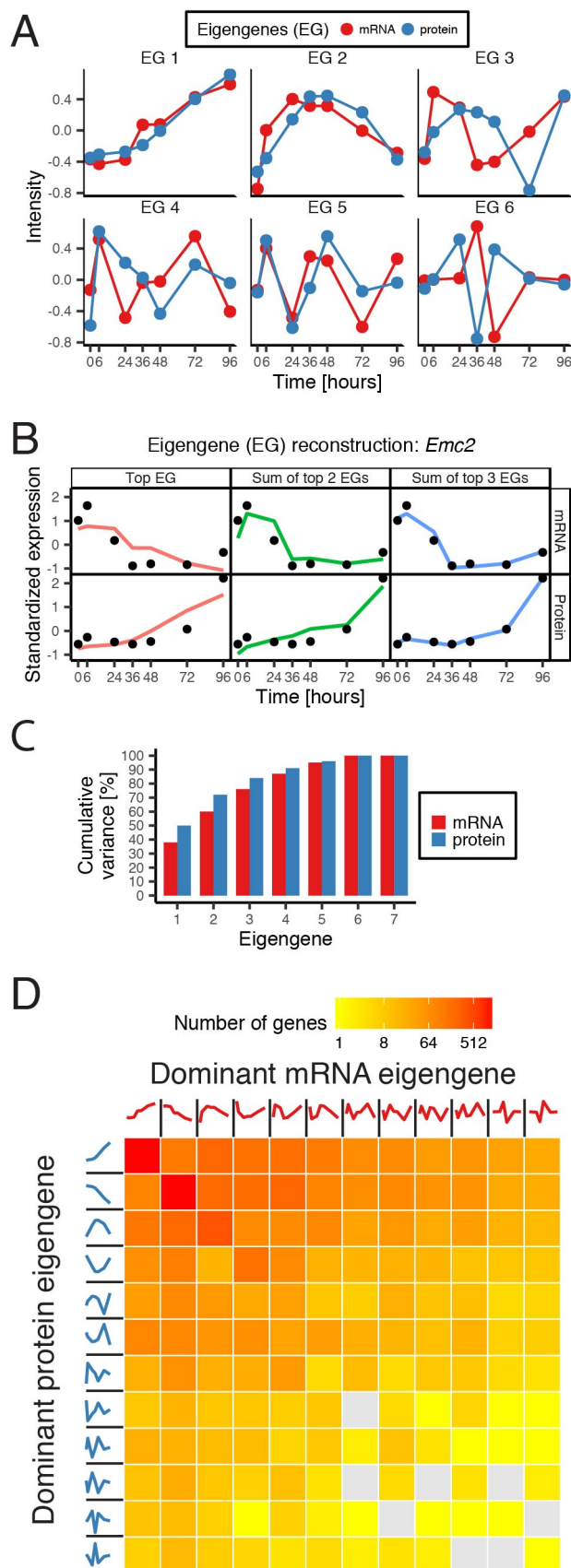
- 669 Alexa, A., Rahnenführer, J., Lengauer, T., 2006. Improved scoring of functional groups from gene
670 expression data by decorrelating GO graph structure. *Bioinformatics* 22, 1600–1607.
671 doi:10.1093/bioinformatics/btl140
- 672 Langmead, B., Salzberg, S.L., 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9, 357–
673 359. doi:10.1038/nmeth.1923
- 674 Durinck, S., Spellman, P.T., Birney, E., Huber, W., 2009. Mapping identifiers for the integration of
675 genomic datasets with the R/Bioconductor package biomaRt. *Nature Protocols* 4, 1184–1191.
676 doi:10.1038/nprot.2009.97
- 677 Lawrence, M., Huber, W., Pagès, H., Aboyoun, P., Carlson, M., Gentleman, R., Morgan, M.T., Carey,
678 V.J., 2013. Software for Computing and Annotating Genomic Ranges. *PLoS Comp Biol* 9,
679 e1003118–10. doi:10.1371/journal.pcbi.1003118
- 680 Li, B., Dewey, C.N., 2011. RSEM: accurate transcript quantification from RNA-Seq data with or
681 without a reference genome. *BMC Bioinformatics* 12, 1. doi:10.1186/1471-2105-12-323
- 682 Love, M.I., Huber, W., Anders, S., 2014. Moderated estimation of fold change and dispersion for
683 RNA-seq data with DESeq2. *Genome Biol* 15, 31. doi:10.1186/s13059-014-0550-8
- 684 Martin, M., 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.
685 *EMBnet.journal* 17, pp. 10–12. doi:10.14806/ej.17.1.200
- 686 Wall, M., Rechtsteiner, A., Rocha, L., 2003. Singular value decomposition and principal component
687 analysis. *A practical approach to microarray data analysis* 91–109.
- 688 Wiśniewski, J.R., Zougman, A., Nagaraj, N., Mann, M., 2009. Universal sample preparation method
689 for proteome analysis. *Nat Methods* 6, 359–362. doi:10.1038/nmeth.1322
- 690
- 691

692 **Figure 1**

693



694 **Figure 2**



695

696

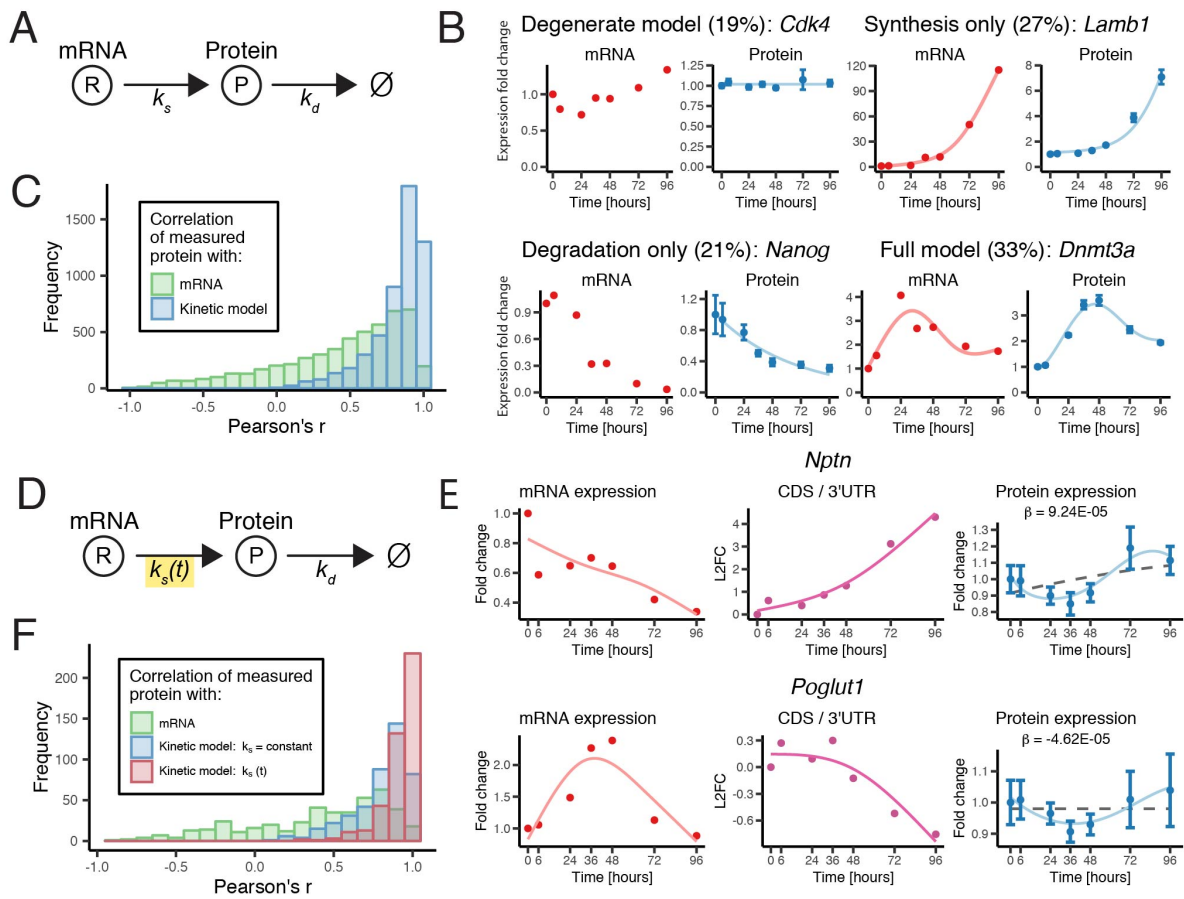
697

698

699

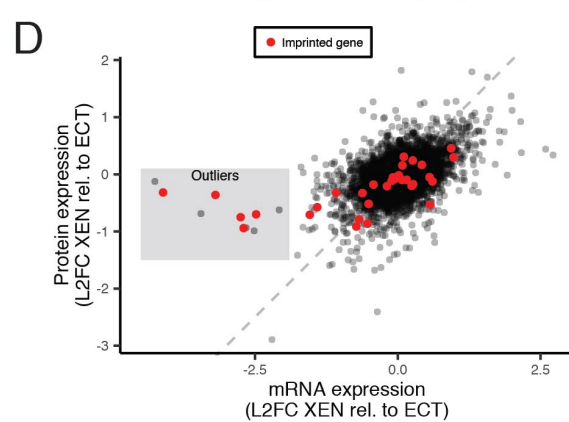
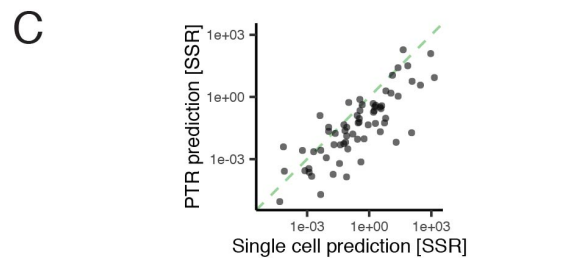
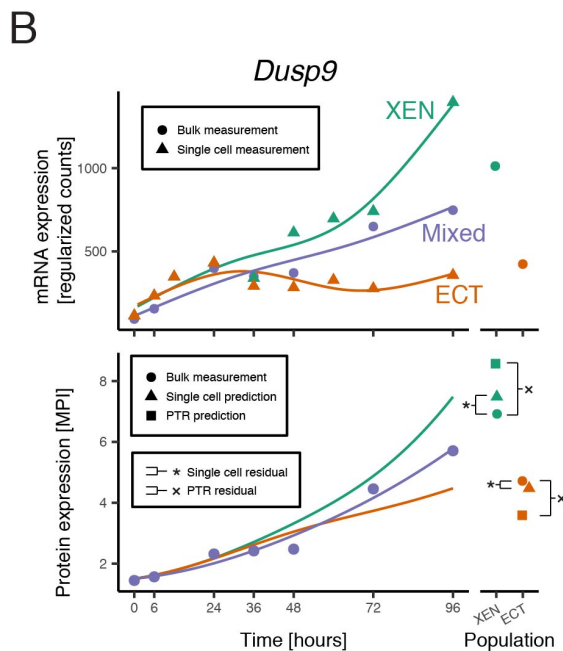
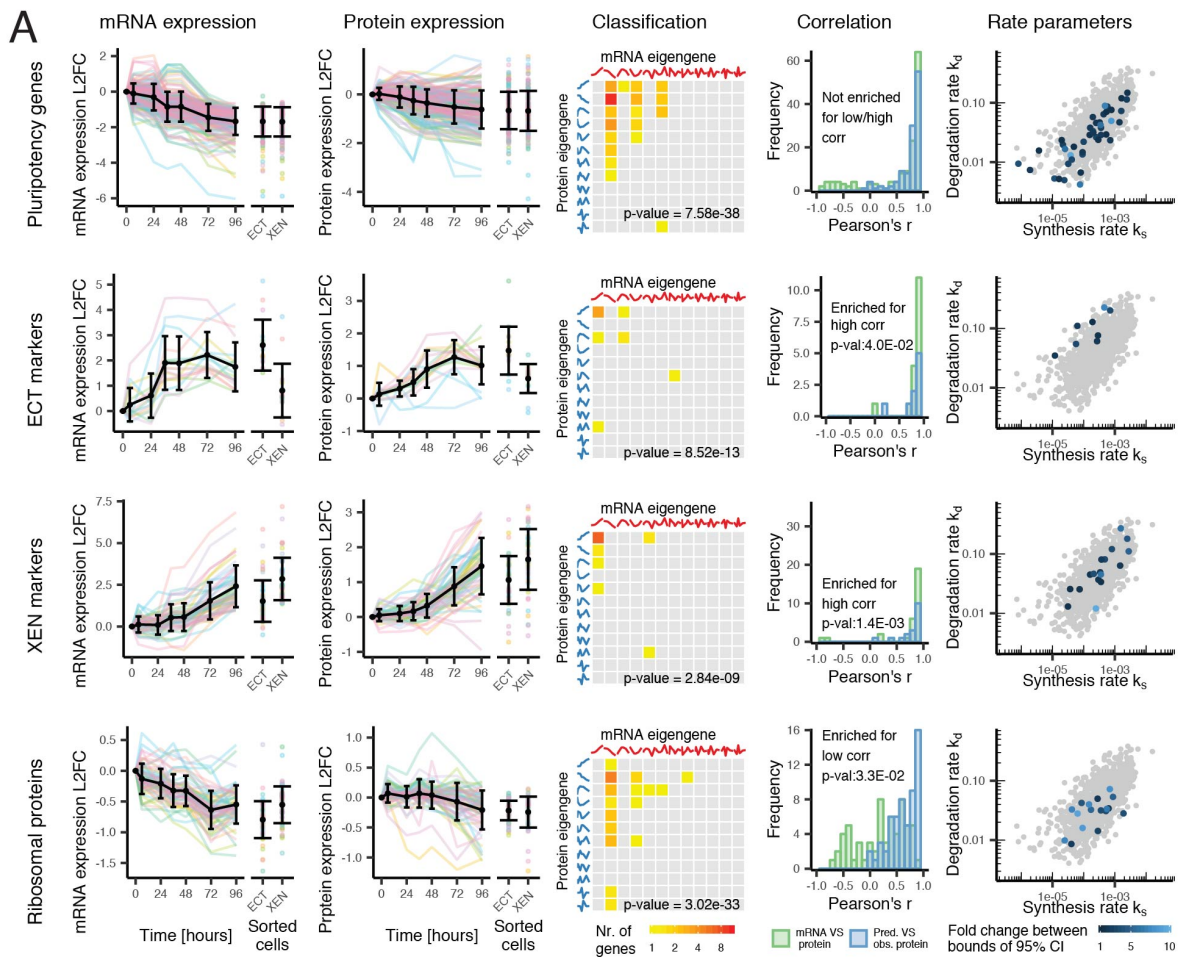
700 **Figure 3**

701



702

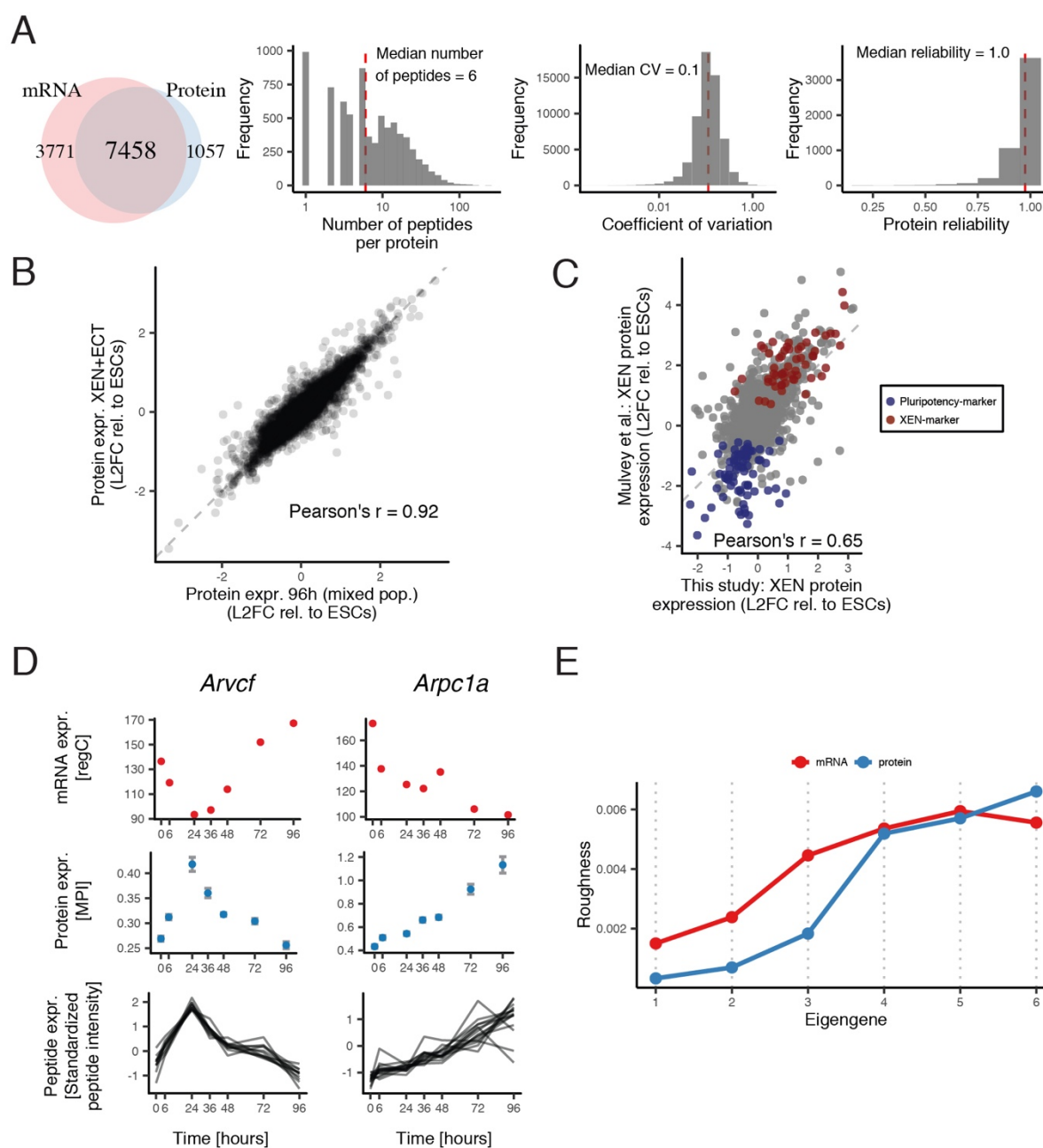
703 **Figure 4**



704

705 Supplementary Figures

706



707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

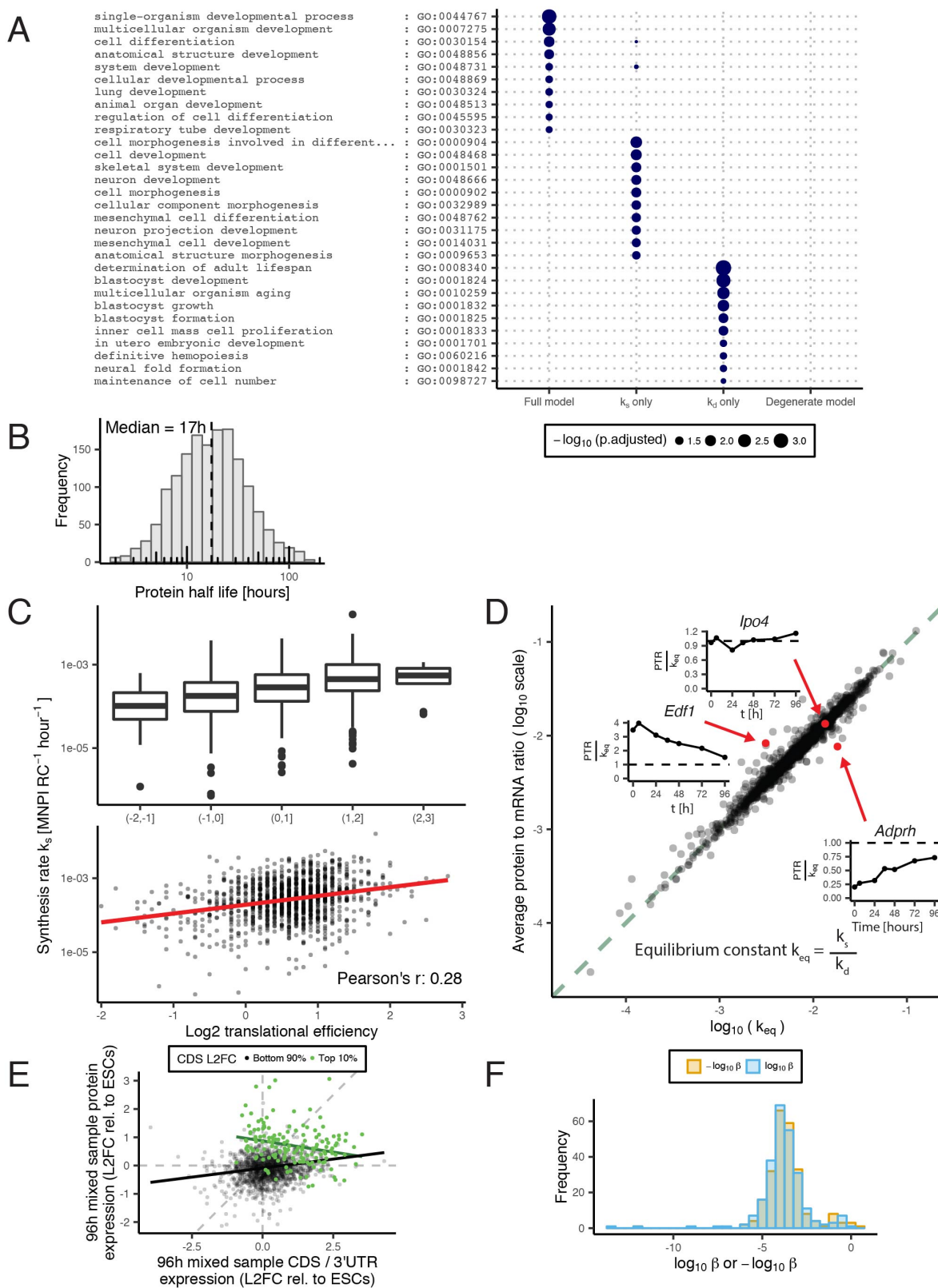
Supplementary Figure 1. Related to figures and 1 and 2. Protein quantification using TMT labeling is robust and reproduces previous results on embryo-derived XEN cells. mRNA eigengenes are more dynamic than protein eigengenes.

(A) From left to right: Venn diagram of the number of genes with quantified mRNA and protein levels (see Methods), distribution of the number of peptides used to quantify protein expression, distribution of the coefficient of variation (CV, SD/mean) of the mean-centered peptide intensities, distribution of the gene-wise protein reliability (Franks et al., 2017). The 7459 genes in the intersection are detected in all mRNA and protein samples.

(B) Protein expression of the 96 h sample (consisting of both XEN and ECT cells) compared with a sample mixed *in silico* from the independently generated purified XEN and ECT cell samples. L2FC: log₂ fold-change.

(C) Protein expression in XEN cells relative to ESCs as measured in this study compared with *in vivo* derived XEN cells measured by Mulvey et al. (2015). Pluripotency- and XEN-marker

723 gene sets were defined using a support vector machine learning algorithm. The pluripotency
724 set is significantly enriched in genes that are downregulated in our data (p-value = 4.0E-4) and
725 the XEN-marker gene set is enriched in genes that are upregulated (p-value = 1.4E-4, gene
726 set enrichment analysis).
727 (D) mRNA, protein and peptide expression for two genes with negative time-wise correlation:
728 *Arcvf* ($r = -0.90$) and *Arpc1a* ($r = -0.91$). 15 and 11 peptides, respectively, were quantified for
729 each gene. regC = regularized counts; MPI = mean peptide intensity.
730 (E) Roughness of mRNA and protein expression eigengenes. The roughness of a profile is
731 defined as the average squared second derivative.
732



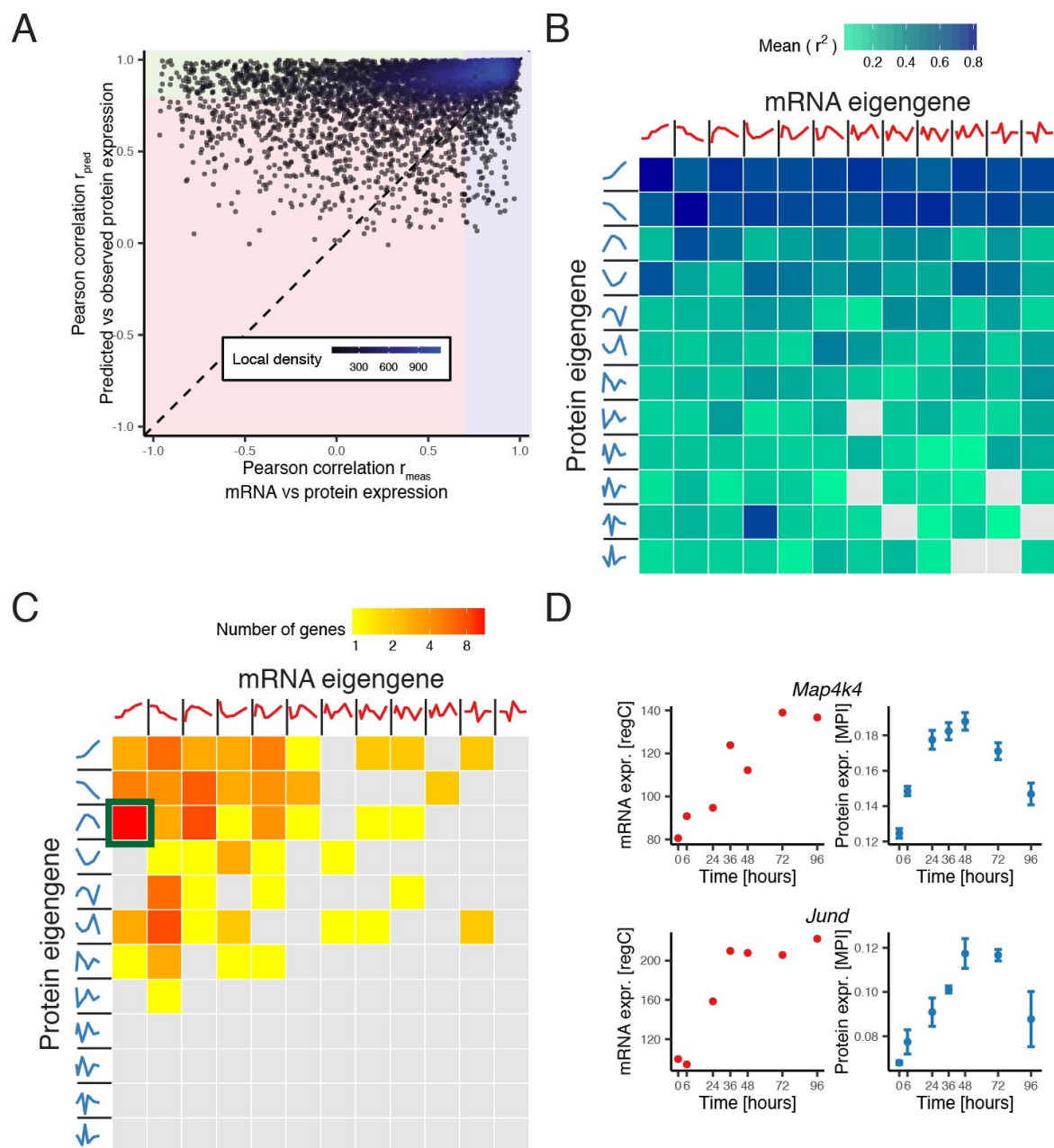
733
734
735
736
737
738
739
740
741
742

Supplementary Figure 2. Related to figure 3. The kinetic models can be related to biological functions and the inferred kinetic rates are biologically meaningful.

(A) Union of the top 10 significantly enriched *cellular differentiation* GO terms for genes fit best by each of the four kinetic models. False discovery rate = 10%.

(B) Protein half life distribution for 1554 genes that were fit best by the full model (according to the BIC) and have precise estimates of the rates (upper and lower bound of the 95% confidence intervals (CIs) fall within a 10-fold range)

743 (C) Translational efficiency (TE) in mESCs from Ingolia et al. (2011) versus our synthesis
744 rates. We show the rates for 1284 genes (intersection between data from Ingolia et al. (2011)
745 and the 1554 genes shown in B). Boxplots represent the binned TE with whiskers indicating
746 $1.5 \times \text{IQR}$.
747 (D) Log_{10} protein to mRNA ratio (PTR) versus equilibrium constant ($k_{eq} = k_s / k_d$) for the 1554
748 genes described in B. Each data point is an individual gene. Genes that are at equilibrium
749 ($\text{PTR} = k_{eq}$) are on the 1:1 line (green). Inserts: PTR relative to k_{eq} across time are shown for
750 three example genes that are above, approximately on and below the 1:1 line.
751 (E) Ratio of CDS and 3'UTR expression versus protein expression in the 96h sample relative
752 to ESCs. The genes with the highest CDS expression fold change are indicated in green. Solid
753 lines indicate linear regression fits. CDS = coding DNA sequence, 3'UTR = 3' untranslated
754 region.
755 (F) Distribution of the parameter β of the extended model, which sets the strength of the
756 influence of the CDS-3'UTR ratio on the synthesis rate. Shown are the values of β for the 492
757 genes that are improved by the extended kinetic model (according to the BIC).
758

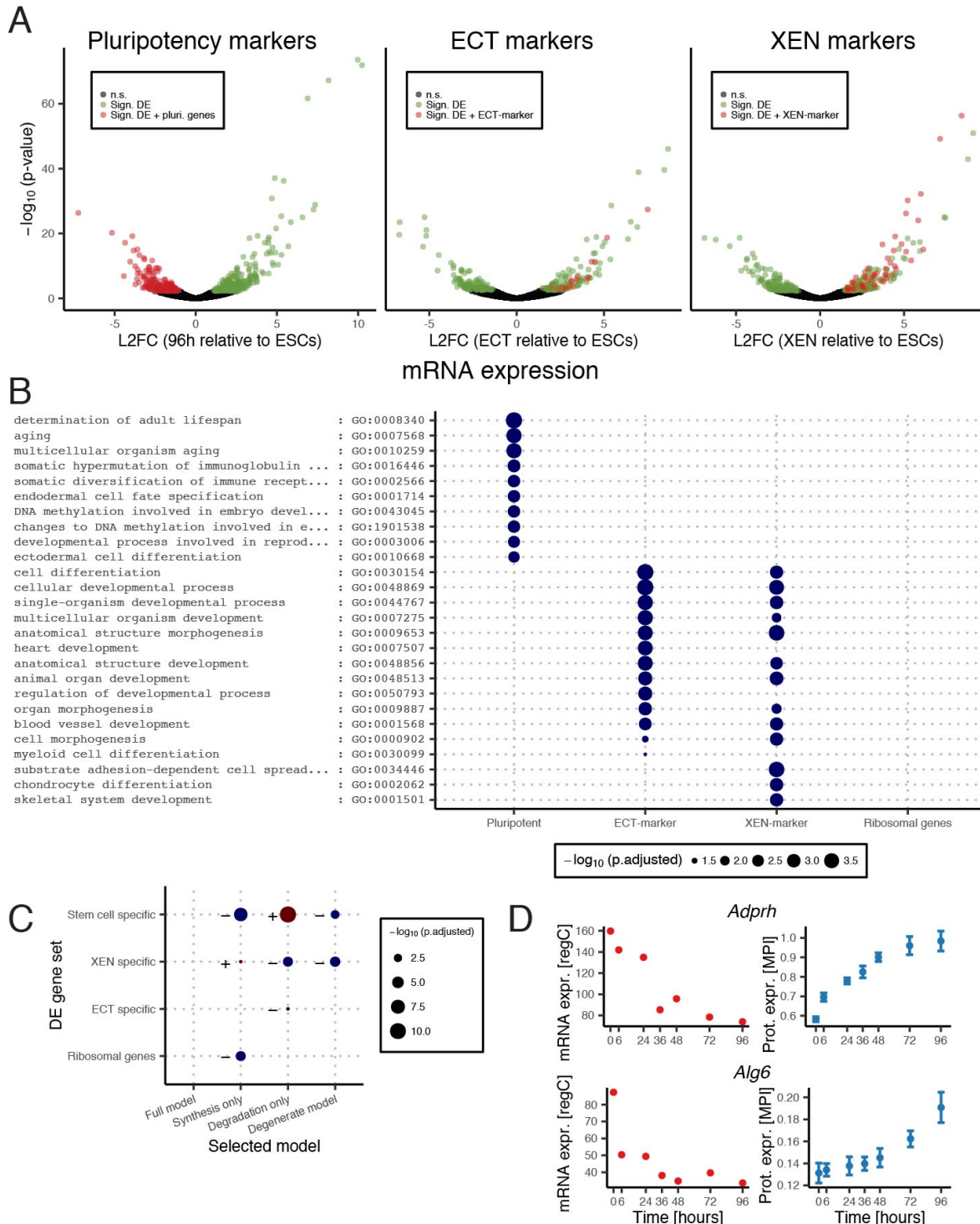


759
760 **Supplementary Figure 3. Related to figure 3. Genes in the MAPK signaling pathway are**
761 **regulated dynamically at the protein level during differentiation.**
762

763 (A) Pearson correlation between measured protein and mRNA (r_{meas}) versus Pearson
764 correlation between measured and predicted protein (r_{pred}). Background coloring indicates:
765 concordant genes with (high r_{meas} , blue), discordant genes that are not well-fit (low r_{meas} , low
766 r_{pred} , red) and discordant genes that are well-fit (low r_{meas} , high r_{pred} , green). Here we
767 consider genes with $r_{\text{meas}} < 0.7$ to be discordant (see Methods). To assure the the model
768 prediction correlates substantially better with the measured protein than the measured mRNA
769 we require $r_{\text{pred}} \geq 0.8$ for a gene to be considered well-fit.
770 (B) Dominant eigengene classification of all 7459 genes. The color of a tile indicates the mean
771 fraction of variance explained (mean r^2) by the best-fitting kinetic model for genes with a
772 particular combination of dominant mRNA and protein eigengene.
773 (C) Dominant eigengene classification of the 368 genes that are not well-fit by the basic kinetic
774 model (red area of A) and exhibit a bigger than median fold-change between the contribution
775 of the dominant eigengene and the second most contributing eigengene. The color of a tile
776 indicates the number of genes with a particular combination of dominant mRNA and protein
777 eigengene. Enrichment analysis revealed an enrichment of MAPK signaling pathway genes in
778 the tile highlighted in green (q-value = $1.8e-3$).

779 (D) mRNA and protein expression profiles of two genes from the tile highlighted in C. Error
780 bars: SEM. regC = regularized counts; MPI = mean peptide intensity.
781
782

783



784

785

786

787

788

789

790

791

792

793

794

795

796

Supplementary Figure 4. Related to figure 4. The different subtypes of the kinetic model are enriched in gene sets defined by the differentiation process

(A) Volcano plots (mRNA relative expression versus p-value for differential expression) for the 96 h sample, the ECT sample and the XEN sample. mRNA expression is always relative to the 0 h sample (ESCs). Genes colored in both red or green are significantly differentially expressed with a false discovery rate (FDR) of 10%. Only genes colored red are considered marker genes: pluripotency markers are down regulated in the 96 h sample, ECT and XEN markers are upregulated and have a minimum fold change of 2 compared with the other purified sample (see Methods).

(B) Union of the top 10 significantly enriched *cellular differentiation* GO terms for genes in each of the three DE gene sets and the ribosomal genes. FDR = 10 %.

797 (C) Overrepresentation (+ / blue) and underrepresentation (– / red) of the various subtypes of
798 the basic kinetic model in the gene sets from B. (D) Genes in pluripotency gene set with
799 upregulated protein expression. regC = regularized counts; MPI = mean peptide intensity.
800
801