

1

2 SUMO maintains the chromatin environment of human induced pluripotent stem cells.

3

4

5

6 Barbara Mojsa<sup>1</sup>, Michael H. Tatham<sup>1</sup>, Lindsay Davidson<sup>2</sup>, Magda Liczmanska<sup>1</sup>, Jane E. Wright<sup>1</sup>,

7 Nicola Wiechens<sup>1</sup>, Marek Gierlinski<sup>3</sup>, Tom Owen-Hughes<sup>1</sup> and Ronald T. Hay<sup>1\*</sup>

8

9

10 <sup>1</sup>Division of Gene Regulation and Expression, <sup>2</sup>Division of Cell and Developmental Biology,

11 <sup>3</sup>Division of Computational Biology, School of Life Sciences, University of Dundee, Dundee, UK

12

13 \*Corresponding author [r.t.hay@dundee.ac.uk](mailto:r.t.hay@dundee.ac.uk)

14

## 15 **Abstract**

16 Pluripotent stem cells represent a powerful system to identify the mechanisms governing cell  
17 fate decisions during early mammalian development. Covalent attachment of the Small  
18 Ubiquitin Like Modifier (SUMO) to proteins has emerged as an important factor in stem cell  
19 maintenance. Here we show that SUMO is required to maintain stem cells in their pluripotent  
20 state and identify many chromatin-associated proteins as bona fide SUMO substrates in  
21 human induced pluripotent stem cells (hiPSCs). Loss of SUMO increases chromatin  
22 accessibility and expression of long non-coding RNAs and human endogenous retroviral  
23 elements, indicating a role for the SUMO modification of SETDB1 and a large TRIM28 centric  
24 network of zinc finger proteins in silencing of these elements. While most protein coding  
25 genes are unaffected, the Preferentially Expressed Antigen of Melanoma (PRAME) gene locus  
26 becomes more accessible and transcription is dramatically increased after inhibition of SUMO  
27 modification. When PRAME is silent, a peak of SUMO over the transcriptional start site  
28 overlaps with ChIP-seq peaks for cohesin, RNA pol II, CTCF and ZNF143, with the latter two  
29 heavily modified by SUMO. These associations suggest that silencing of the PRAME gene is  
30 maintained by the influence of SUMO on higher order chromatin structure. Our data indicate  
31 that SUMO modification plays an important role in hiPSCs by repressing genes that disrupt  
32 pluripotency networks or drive differentiation.

33

## 34 **Introduction**

35 Pluripotent cells display the property of self-renewal and have the capacity to  
36 generate all of the different cells required for the development of the adult organism. The  
37 pluripotent state is defined by the gene expression programme of the cells and is driven by  
38 expression of the core transcription factors OCT4, SOX2 and NANOG<sup>1</sup> that sustain their own  
39 expression by virtue of a positive linked autoregulatory loop while activating genes required  
40 to maintain the pluripotent state and repressing expression of the transcription factors for  
41 lineage specific differentiation<sup>2</sup>. Once terminally differentiated, somatic cell states are  
42 remarkably stable. However, forced expression of key pluripotency transcription factors that  
43 are highly expressed in embryonic stem cells (ESCs), including OCT4, SOX2 and NANOG, leads  
44 to reprogramming back to the pluripotent state<sup>3-5</sup>. As the efficiency of reprogramming is very  
45 low, it is clear that there are roadblocks to reprogramming designed to safeguard cell fates<sup>6</sup>.

46 7. Small Ubiquitin like Modifier (SUMO) has emerged as one such roadblock and reduced  
47 SUMO expression decreases the time taken and increases the efficiency of reprogramming in  
48 mouse cells<sup>8-10</sup>. Three SUMO paralogues, known as SUMO1, SUMO2 and SUMO3 are  
49 expressed in vertebrates. Based on almost indistinguishable functional and structural  
50 features SUMO2 and SUMO3 are collectively termed SUMO2/3 and share only about 50%  
51 amino acid sequence identity with SUMO1. SUMOs are conjugated to lysine residues in a large  
52 number of target proteins and as a consequence influence a wide range of biological  
53 processes. SUMOs are initially translated as inactive precursors that require a precise  
54 proteolytic cleavage carried out by SUMO specific proteases (SENPs) to expose the terminal  
55 carboxyl group of a Gly-Gly sequence that ultimately forms an isopeptide bond with the  $\epsilon$ -  
56 amino group of a lysine residue in the substrate protein. The heterodimeric E1 SUMO  
57 Activating Enzyme (SAE1/SAE2) uses ATP to adenylate the C-terminus of SUMO before  
58 forming a thioester with a cysteine residue in a second active site of the enzyme and releasing  
59 AMP. SUMO is then trans-esterified on to a cysteine residue in the single E2 SUMO  
60 conjugating enzyme Ubc9. Assisted by a small group of E3 SUMO E3 ligases, including the PIAS  
61 proteins, RanBP2 and ZNF451 the SUMO is transferred directly from Ubc9 onto target  
62 proteins<sup>11</sup>. Modification of target proteins may be short-lived, with the SUMO being removed  
63 by SENPs. Together this creates a highly dynamic SUMO cycle where the net SUMO  
64 modification status of individual proteins is determined by the rates of SUMO conjugation  
65 and deconjugation<sup>12</sup>. Preferred sites of SUMO modification conform to the consensus  $\psi$ KxE,  
66 where  $\psi$  represents a large hydrophobic residue<sup>13, 14</sup>. A conjugation consensus is present in  
67 the N-terminal sequence of SUMO2 and SUMO3 and thus permits self-modification and the  
68 formation of SUMO2/3 chains<sup>15</sup>. As a strict consensus is absent from SUMO1, it does not form  
69 chains as readily as SUMO2/3<sup>12</sup>. Once linked to target proteins SUMO allows the formation of  
70 new protein-protein interactions as the modification can be recognised by proteins  
71 containing a short stretch of hydrophobic amino acids termed a SUMO interaction motif<sup>16</sup>.

72 Stem cell lines are excellent models to study mechanisms controlling self-renewal and  
73 pluripotency. Mouse ESCs have been widely used as they can also be used for *in vivo* studies  
74 by making chimeras in mouse blastocysts, however, they do display different characteristics  
75 from human ESCs. Mouse ESCs require leukemia inhibitory factor (LIF) and bone  
76 morphogenetic protein (BMP) signalling to maintain their self-renewal and pluripotency<sup>17, 18</sup>.  
77 In contrast LIF does not support self-renewal and BMPs induce differentiation in hESCs<sup>19-21</sup>.

78 The maintenance of the pluripotent state of hESCs requires basic fibroblast growth factor  
79 (bFGF, FGF2) and activin/nodal/TGF- $\beta$  signalling along with inhibition of BMP signalling<sup>22, 23</sup>.  
80 These differences may reflect the developmental stages at which ESC lines are established *in*  
81 *vitro* from mouse and human blastocysts, or may be due to differences in early embryonic  
82 development<sup>24</sup>. As hESCs are derived from embryos their use is limited, in contrast human  
83 Induced Pluripotent Stem Cells (hiPSCs) are derived by reprogramming normal somatic cells  
84 and display most of the characteristics of hESCs<sup>3</sup>. As a result, hiPSCs are now widely used to  
85 study self-renewal and pluripotency in humans.

86 To determine the role of SUMO modification in hiPSCs we made use of ML792, a highly  
87 potent and selective inhibitor of the SUMO Activating Enzyme<sup>25</sup>. Treatment of hiPSCs with  
88 this inhibitor rapidly blocks *de novo* SUMO modification allowing endogenous SENPs to strip  
89 SUMO from targets. When used over the course of 48 hours hiPSCs treated with ML792 lose  
90 the majority of SUMO conjugation but show no large-scale changes to the cellular proteome  
91 nor loss of viability, although markers of pluripotency are reduced. By inhibiting SUMO  
92 conjugation ML792 reduces chromatin-associated SUMO and increases DNA accessibility.  
93 This results in increased transcription of a group of long non-coding RNAs (lncRNAs) and  
94 human endogenous retrovirus (HERV) elements, while protein coding genes are largely  
95 unaffected. One important exception is the Preferentially Expressed Antigen of Melanoma  
96 (PRAME) gene. SUMO modification inhibition increases the accessibility of the PRAME locus  
97 and leads to a large increase in transcription and accumulation of PRAME protein. SUMO site  
98 and paralogue specific proteomic analysis of hiPSCs reveals extensive SUMO modification of  
99 proteins involved in transcriptional repression, RNA splicing and ribosome biogenesis.  
100 Specifically, SUMO modification of the boundary and looping elements CTCF and ZNF143 and  
101 their colocalisation with cohesin components suggest an important role for SUMO in  
102 organising the higher order chromatin architecture in hiPSCs.

103

## 104 **Results**

105

106 **Inhibition of SUMO modification leads to loss of select pluripotency markers in hiPSCs.** To  
107 determine the role of SUMO modification in the maintenance of pluripotency in hiPSCs we  
108 used ML792<sup>25</sup>. This inhibitor has been reported to block proliferation of cancer cells,  
109 particularly those overexpressing Myc<sup>25</sup>, but has not been evaluated in hiPSCs. To address the

110 role of SUMO modification in ChiPS4, we established that 400nM ML792 effectively reduced  
111 SUMO modification after 4 hrs with minimal effects on cell viability in longer treatments. We  
112 restricted our analyses to ML792 treatment times that did not exceed 48 hrs, such that the  
113 immediate effects of SUMO modification inhibition could be evaluated. Microscopic  
114 examination revealed that ML792 treatment caused morphological changes with the ChiPS4  
115 cells becoming larger and flatter (Fig. 1a). The rate of proliferation was unchanged after 24  
116 hrs but was slightly reduced after 48 hrs (Fig. 1b). DNA staining of the cells and analysis by  
117 flow cytometry indicated that the cell cycle distribution after 24hrs was unaltered by ML792  
118 treatment but after 48 hrs displayed an increased proportion of cells in G2 phase and cells  
119 with increased DNA content suggesting endoreplication (Fig. 1c). Western blot analysis  
120 revealed a loss of high molecular weight SUMO1 and SUMO2 conjugates and concomitant  
121 increase in free SUMOs in ChiPS4 cells (Fig. 1d) caused by rapid removal of SUMO from  
122 modified proteins by SENPs. Analysis of the protein levels of key pluripotency markers  
123 indicated that while OCT4 and SOX2 were unchanged, inhibition of SUMOylation resulted in  
124 a decrease in NANOG protein (Fig. 1d). This appeared to be a consequence of reduced  
125 transcription as determination of mRNA levels by reverse transcriptase quantitative  
126 polymerase chain reaction (RT-qPCR) after ML792 treatment demonstrated a reduction in  
127 *NANOG* mRNA. This was also apparent for *KLF4*, but consistent with Western blotting, the  
128 levels of *OCT4* and *SOX2* mRNA were unchanged (Fig. 1e). To further investigate the nature  
129 and causes of the observed morphological changes ChiPS4 ML792 treated cells were analysed  
130 by phenotypic screening using cell painting<sup>26</sup> (Fig. 2a). Principle component analysis (PCA)  
131 indicated that there are clear differences between cells treated with ML792 for 48h and  
132 untreated/vehicle (DMSO) treated. The main differences were found in the nuclear  
133 compartment (Supplementary Fig. 1). Feature extraction identified changes in the global size  
134 (Area of nuclei) and shape (Nuclei form factor) of the nucleus and the structure of the  
135 nucleolus (Nuclei: FITC texture correlation) (Fig. 2a). These findings were validated using  
136 traditional immunofluorescence (IF) approaches. NANOG expression as well as the size and  
137 shape of the nucleus are both affected by ML792 treatment. NOP58 was used as a marker of  
138 the nucleolus, which undergoes a dramatic increase in size and shape. The classic punctate  
139 nuclear localisation pattern of SUMO1 and SUMO2 is altered by their deconjugation from  
140 substrates, becoming more diffuse and less tightly associated with the nucleus (Fig. 2b).

141 To evaluate the effect of inhibition of SUMO modification on the global proteome in  
142 hiPSCs, proteins from ChiPS4 cells either untreated or treated with ML792 for 24 or 48 hours  
143 were analysed by label-free quantitative proteomics. 4741 proteins were identified and  
144 quantified in all replicates of at least one experimental group (Supplementary Data File 1).  
145 PCA of the proteomic data showed a progressive trend of changing cellular proteome during  
146 SUMOylation inhibition (Supplementary Fig. 2a), although individual protein fold changes  
147 compared at both time-points showed little evidence for large-scale global shifts in protein  
148 abundance (Supplementary Fig. 2b, c). Known pluripotency markers were progressively  
149 reduced during ML792 exposure (Fig. 2c), and linker histones were reduced in abundance  
150 after 48 hours (Fig. 2d). Linker histones and the GOCC group 'Collagen-associated extracellular  
151 matrix' were the only two categories to show any significant co-regulation according to  
152 STRING analysis, (Supplementary Fig. 2d, e). Thus, global proteome changes do not seem to  
153 provide an explanation for an observed morphological change of this magnitude. It therefore  
154 seems likely that the observed changes in nuclear structure are due to direct consequences  
155 of removal of SUMO from critical factors that contribute to chromatin structure and function.

156

157 **Removal of SUMO in hiPSCs increases chromatin accessibility.** To determine the  
158 chromosomal landscape of SUMO1 modification Chromatin Immunoprecipitation coupled to  
159 high throughput sequencing (ChIP-seq) was conducted on ChiPS4 cells. SUMO1 bound  
160 chromatin was enriched from cross-linked cell extracts using an antibody with previously  
161 confirmed specificity for SUMO1<sup>27</sup> and utility in ChIP analysis<sup>28</sup>. Total genomic DNA was  
162 sequenced to obtain reference input profiles. SUMO peaks were usually less than 1kb and  
163 typically 300bp in length (Supplementary Fig. 3a) and were clearly enriched above the  
164 background (Supplementary Fig. 3b) and in promoter, introns and intergenic regions  
165 (Supplementary Fig. 3c). Similar to the situation reported in mouse ESCs<sup>10, 29</sup>, SUMO peaks  
166 were over-represented on endogenous retroviruses (ERVs) and other non-viral long terminal  
167 repeats (LTRs) (Supplementary Fig. 3d). Of note, SUMO1 accounts for a high proportion of  
168 protein SUMOylation in ChiPS4 cells when assessed by mass spectrometry and Western  
169 blotting. SUMO peaks overlapped with over 10% of KAP1/TRIM28 and SETDB1 peaks, but over  
170 50% of peaks for CTCF, ZNF143 and cohesion components RAD21 and SMC3 (Supplementary  
171 Fig. 3e, f).

172           Having used SUMO1 ChIP-seq to determine the precise location of SUMO1-modified  
173 proteins on chromatin in ChiPS4 cells we used ML792 to facilitate the removal of SUMO from  
174 these sites and used ATAC-seq to monitor changes in chromatin accessibility over time. For  
175 that, ChiPS4 cells were treated with DMSO vehicle or ML792 for 4, 8, 24, 48h leading to time  
176 dependent release of SUMO from high molecular weight material and a reduction in NANOG  
177 expression (Supplementary Fig. 4a). The loss of SUMO was accompanied by a general increase  
178 in chromatin accessibility across the genome as determined by an increase in number of  
179 ATAC-seq peaks over time (Fig. 3a, Supplementary Data File 5). ATAC-seq peaks that are lost  
180 do not increase in the same way and after 48h exposure to the inhibitor there are at least  
181 three times more ATAC-seq peaks gained than lost (Fig. 3a). Comparing the ATAC-seq and  
182 SUMO1 ChIP-seq data suggests that removing SUMO leads to an increase in chromatin  
183 accessibility at the sites previously occupied by SUMO as around 20% of gained ATAC-seq  
184 peaks at all time points overlap with a pre-existing SUMO1 ChIP-seq peak, while less than 5%  
185 of lost ATAC-seq peaks overlap with a pre-existing SUMO1 ChIP-seq peak (Fig. 3b).

186           To determine if specific genomic regions change their accessibility after removal of  
187 SUMO, the ATAC-seq data was analysed using HOMER<sup>30</sup>, allowing various types of genomic  
188 regions to be annotated and classified into functionally related groups. Chromatin  
189 accessibility was gained mainly in repetitive DNA sequences, such as non-LTR and LTR  
190 retrotransposons which together account for around 80% of all gained ATAC-seq peaks (Fig.  
191 3c). When compared to the types of genomic regions represented in non-changing ATAC-seq  
192 peaks, LTR-retrotransposons are strongly enriched in gained ATAC-seq peaks or gained ATAC-  
193 seq peaks overlapping with SUMO1 peaks throughout the treatment with ML792 (Fig. 3c, d).  
194 It thus suggests an important role for SUMO in maintaining these viral elements of the hiPSCs  
195 genome in a compact chromatin environment. There is little enrichment for any particular  
196 type of genomic region in the ATAC-seq peaks or ATAC-seq peaks overlapping with SUMO1  
197 that are lost after removal of SUMO (Supplementary Fig. 4b, c).

198           The dynamic behaviour of each opening locus can be assessed using density plots (Fig.  
199 3e). Regions that lost chromatin accessibility did not show any significant enrichments or  
200 patterns (Supplementary Fig. 4d). As indicated above, 20% of gained ATAC-seq peaks at all  
201 time points after treatment with ML792 overlap with a pre-existing SUMO ChIP-seq peak  
202 present in untreated cells (Fig. 3b). These are also the regions that respond quickly to ML792  
203 mediated chromatin opening (higher intensity SUMO1 ChIP peaks overlap with ATAC-seq

204 peaks that appear sooner during the time course) and once open the chromatin state is  
205 maintained throughout the time course (Fig. 3e). It is also evident that the SUMO overlapping  
206 ATAC-seq peaks that are classified as ‘nascent gained’ at a later time point (e.g. 48h) already  
207 demonstrate a trend for chromatin relaxation at earlier time points, but do not reach the  
208 necessary threshold. These SUMO overlapping ATAC-seq peaks are strongly associated with  
209 repetitive DNA sequences, particularly retrotransposons (Fig. 3d). Based on the association of  
210 SUMO with repressed chromatin and the proteins associated with repression of viral  
211 elements in the genome, a similar peak overlap procedure was performed with transcription  
212 factor ChIP-seq data obtained from the ENCODE project<sup>31</sup>. An unbiased analysis of 161 factors  
213 was undertaken, but those showing the highest overlap with SUMO ChIP-seq and gained  
214 ATAC-seq peaks were TRIM28, SETDB1, CBX3 (Fig. 3f) that are known SUMO-modified  
215 silencers of viral DNA elements. The overlap between those factors and SUMO1 ChIP-seq or  
216 gained ATAC-seq peaks is significantly higher than that calculated for non-changing or lost  
217 ATAC-seq peaks (Fig. 3f). These data suggest an important role for SUMO in maintaining a  
218 compact chromatin environment around LTR elements in hiPSCs.

219

220 **Inhibition of SUMO modification in hiPSCs selectively alters transcription.** To determine if  
221 chromatin associated SUMO regulates transcription in hiPSCs, ChiPS4 cells were treated as  
222 for ATAC-seq and analysed by RNA-seq (Supplementary Fig. 5). After 48h of treatment with  
223 ML792 996 RNAs displayed an increase in transcription while 281 RNAs were decreased (Fig.  
224 4a, Supplementary Data File 6). Although the observed effect on protein coding mRNAs was  
225 rather modest, the expression of lncRNAs was significantly affected (both increased and  
226 decreased) by inhibition of SUMO modification (Fig. 4b, Supplementary Fig. 6a). Interestingly,  
227 stage specific expression of certain lncRNAs e.g. LINC-ROR is important for the maintenance  
228 of the pluripotency network<sup>32</sup>. While protein coding genes as a group were not significantly  
229 enriched during the time-course, it is possible that the aggregation of incremental changes  
230 could exert a specific biological response. For example, consistent with the protein level  
231 analysis (Fig. 2c), mRNA levels of various pluripotency markers decrease globally during the  
232 treatment (Supplementary Fig. 6b). To determine if there were any functional patterns to the  
233 protein-coding genes regulated by ML792 treatment, a data-dependent clustering analysis  
234 was undertaken (Supplementary Data File 4). Hierarchical clustering of protein-coding RNAs  
235 based on response to ML792 allowed separation into 9 clusters (Fig. 4c). Cluster A contained



236 the least responsive RNAs and represented over 96% of the entries (Fig. 4c-f). A number of  
237 GSEA categories were depleted in cluster A with high significance, all of which were previously  
238 identified in ES cells as being regulated by histone methylation or by protein complexes  
239 themselves regulating histone methylation such as PRC2 (Fig. 4f). Indeed, these categories  
240 were among the most significantly enriched in the two largest clusters of ML792-sensitive  
241 RNAs (clusters B and C). Clusters B and C show opposing responses to ML792 treatment,  
242 therefore these data imply that although promoter methylation is a common feature of  
243 regulated genes the outcomes are not qualitatively the same for all of them (Supplementary  
244 Fig. 7). Thus, inhibition of SUMO modification increases transcription of a group of lncRNA  
245 genes, but with notable exceptions, has limited impact on transcription of protein coding  
246 genes.

247

248 **SUMO silences the PRAME gene in human stem cells.** One important exception is the PRAME  
249 gene, transcription of which significantly increases in response to ML792 treatment in a time-  
250 dependent manner. This increase is evident at the very earliest time point (4h) and continues  
251 to increase up to 48h (Fig. 5a). Inspection of the SUMO1 ChIP-seq data from untreated ChiPS4  
252 cells where *PRAME* expression is silenced reveals a prominent SUMO peak located over the  
253 transcriptional start site (TSS) (Fig. 5b). Analysis of the ATAC-seq data indicates that while this  
254 region is minimally accessible in untreated cells, removal of SUMO leads to a time dependent  
255 increase in chromatin accessibility. At the earliest time point (4h) this increase is confined to  
256 the TSS, but over time the region of accessibility spreads in towards the coding body of the  
257 gene (Fig. 5b). This is highly specific for the *PRAME* gene as the same SUMO1 peak overlaps  
258 with a TSS of the divergently transcribed *LL22NCO3-63E9.3* lncRNA gene which displays  
259 neither an increase in chromatin accessibility, nor an increase in transcription in response to  
260 ML792 treatment (Fig. 5b). Analysis of ChIP-seq data from other hESC lines reveals that CTCF,  
261 ZNF143 and the cohesion subunit RAD21 almost precisely overlap with this SUMO ChIP-seq  
262 peak (Supplementary Fig. 8a). All of these factors are known to be associated with higher  
263 order chromatin interactions. Comparison of the protein level and transcriptional changes  
264 after 48 hours ML792 exposure shows a single outlier in PRAME (Fig. 5c, Supplementary Fig.  
265 8d). PRAME is both the single most upregulated transcript and the most elevated protein  
266 upon SUMOylation inhibition. IF analysis further revealed that in untreated cells PRAME  
267 expression was not above background, but after 48 hours of ML792 treatment PRAME was

268 highly expressed and localised in the nuclei of ChiPS4 cells (Supplementary Fig. 8b). Likewise,  
269 Western blotting demonstrated that PRAME was accumulated after SUMO modification  
270 inhibition (Fig. 5d). To establish that the SUMO regulated expression of the *PRAME* gene was  
271 not unique to ChiPS4 cells a number of other hiPSC and hESC lines were exposed to ML792.  
272 PRAME expression was also robustly induced after 48h of treatment with ML792  
273 (Supplementary Fig. 8c), indicating that SUMO-dependent silencing of PRAME is a common  
274 feature of all tested human pluripotent stem cells (hPSCs). While ML792 is a highly selective  
275 and potent inhibitor of SAE<sup>25</sup> it was important to assess SUMO regulated expression of *PRAME*  
276 by an orthogonal approach. Rather than inhibiting SUMO modification, conjugated SUMO can  
277 be directly removed by expression of an exogenous SUMO specific protease. We previously  
278 used such an approach to demonstrate SUMO dependent regulation of an integrated reporter  
279 gene<sup>33</sup>. Capped and polyadenylated mRNA encoding the catalytic domain of SENP1 was  
280 electroporated into ChiPS4 cells and PRAME expression was monitored by Western blotting  
281 and RT-qPCR. Expression of the protease effectively reduced global SUMO modification and  
282 increased both PRAME protein and mRNA (Fig 5e, f). Thus, silencing of the *PRAME* gene in  
283 hiPSCs appears to be directly mediated by SUMO modification. Enrichment analysis indicated  
284 that transcription of metallothionein genes was reduced in response to SUMOylation  
285 inhibition, while transcription of five further members of the *PRAME* gene family was elevated  
286 (Supplementary Fig. 8e, f), revealing a common link between SUMO and transcription of  
287 *PRAME* genes.

288

289 **SUMO modification restricts HERV expression in hiPSCs.** ATAC-seq and ChIP-seq analyses  
290 revealed that repetitive DNA elements including LTR retrotransposons are primary targets of  
291 SUMOylation-mediated repression. Indeed, expression of lncRNAs is often controlled by LTRs,  
292 which have been hijacked by cellular machinery to function as stage specific promoters.  
293 Typically for RNA-seq experiments such repetitive DNA sequences are mostly removed during  
294 data processing. To analyse HERV expression, an independent data alignment file was created  
295 based on the available [Human Endogenous Retrovirus Database](#), which contains two major  
296 data sets: elements (contiguous sequences) and entities (loci in the human genome consisting  
297 of one or more elements)<sup>34,35</sup>. Inhibition of SUMO modification leads to a general increase in  
298 HERV expression with the number of significantly increased elements being about three times  
299 higher than those reduced at each time point tested (Fig. 6a, Supplementary Fig. 9a, b). One

300 of the best examples of a complex LTR-containing HERV element that is highly and rapidly  
301 induced by deSUMOylation is ERV\_4326325 (Fig. 6b). Inspection of the SUMO1 ChIP-seq  
302 indicates that this chromosomal location contains a pre-existing SUMO1 peak. Indeed, ATAC-  
303 seq data show that increase in chromatin accessibility is initiated from the site of the SUMO1  
304 peak and spreads in towards the HERV locus (Fig. 6c). The increase in RNA expression follows  
305 the changes in chromatin structure and is already obvious at 24h (Fig. 6b, c). To investigate  
306 the global landscape of these changes, significantly affected HERVs were used for clustering  
307 analysis. Three independent clusters were obtained, in which all HERVs with a rapid increase  
308 in expression were found in cluster 1 (Fig. 6d, e, Supplementary Fig. 9c, d). Correlation of the  
309 SUMO ChIP-seq data with the HERV loci in each of the three clusters revealed that HERVs in  
310 cluster 1 have a significantly higher overlap with SUMO1 ChIP-seq peaks when compared to  
311 the proportion observed for all HERVs (Fig. 6f). These data suggest that SUMO modification  
312 maintains HERV loci in a compact chromatin state that facilitates transcriptional repression  
313 of these viral elements in hiPSCs.

314

315 **Identification of SUMO1 and SUMO2 targets in hiPSCs.** Our data suggest an important role  
316 for SUMO in maintaining the chromatin state of hiPSCs. To identify the proteins responsible  
317 and establish the sites of SUMO modification on these factors we used a SUMO site proteomic  
318 approach that allows sites modified by SUMO1 and SUMO2/3 to be identified<sup>36</sup>. To enable this  
319 analysis in hiPSCs the ChiPS4 cell line was engineered to stably express 6His-SUMO-mCherry  
320 constructs for either SUMO1 or SUMO2 (Supplementary Fig. 10a, b) that incorporated the  
321 TGG to KGG mutations to facilitate GlyGly-K peptide immunoprecipitation and  
322 identification<sup>36</sup>. As mCherry is linked to the C-terminus of SUMO the expressed fusion protein  
323 will be processed by endogenous SUMO proteases, release free mCherry and expose the C-  
324 terminal GlyGly sequence for conjugation. Western blotting of single cell clones indicated that  
325 His-tagged SUMO-KGG paralogues were conjugated to substrates in response to heat shock  
326 (Supplementary Fig. 10c). Cells expressing SUMO1-KGG and SUMO2-KGG had normal cell  
327 cycle profiles (Supplementary Fig. 11a), expressed levels of pluripotency markers comparable  
328 to wild type ChiPS4 cells (Supplementary Fig. 11b, c) and retained the ability to differentiate  
329 into endoderm, ectoderm and mesoderm (Supplementary Fig. 11d). Analysis by proteomics  
330 (Supplementary Fig. 12a, b) identified the expected exogenous mCherry, SUMO1 and SUMO2  
331 peptides (Supplementary Fig. 12c) while analysis of common peptides suggested that the

332 exogenous versions of SUMO were conjugated to substrates at roughly similar levels to their  
333 endogenous counterparts (Supplementary Fig. 12d). Whole cell proteomics (Supplementary  
334 Fig. 12e) confirmed that the engineered cell lines did not significantly change their expressed  
335 proteome (Supplementary Fig. 12f, Supplementary Data File 2). Thus, expression of SUMO  
336 mutants did not disrupt the normal pluripotent state or differentiation potential of ChiPS4  
337 cells.

338         The workflow for the identification of SUMO targets incorporates proteomic analysis  
339 at three levels (Supplementary Fig. 13a). The experiment involves analysis of whole cell  
340 extracts (Supplementary Fig. 13b) to monitor total protein levels, analysis of Nickel NTA-  
341 affinity purified proteins (Supplementary Fig. 13c) to monitor SUMO modified proteins and  
342 analysis of GG-K immunoprecipitations (Supplementary Fig. 13d) to identify sites of SUMO  
343 modification. Across the two experimental runs a total of 976 SUMO sites were identified in  
344 427 proteins. Approximately 84% of these had already been described in at least one of four  
345 large-scale SUMO2 site proteomics studies totalling 49768 unique sites of non-STEM cell  
346 origin (Fig. 7a, Supplementary Data file 3). DNA methyl transferase DNMT3B and the key  
347 embryonic stem cell transcription factor SALL4 were among a small group of proteins with at  
348 least three novel sites in this study (Fig. 7a). Based on GG-K peptide intensity SALL4 is the 6<sup>th</sup>  
349 most modified SUMO substrate in hiPSCs and contains 17 sites of modification (Fig. 7b).  
350 DNMT3B contains 12 sites and is the 7<sup>th</sup> most modified substrate while the methyl DNA  
351 binding protein MBD1 contains 8 sites and is also in the top 10 SUMO substrates (Fig. 7b).  
352 TRIM28 and TRIM24 are highly modified substrates and SUMO appears to play an important  
353 role in their ability to repress retroviral elements<sup>29</sup>. CTCF is heavily modified with SUMO and  
354 this is consistent with the overlap of SUMO and CTCF ChIP-seq peaks over the TSS of the  
355 PRAME gene that appears to be silenced by SUMO modification in ChiPS4 cells (Fig. 7b). There  
356 is also evidence for extensive SUMO chain formation as the branch points from SUMO2/3  
357 chains are amongst the most abundant GG-K peptides (Fig. 7b). Indeed, SUMOylation of a  
358 number of these heavily SUMO modified proteins could be detected directly in total cell  
359 lysates from control ChiPS4 cells, but not ML792 treated hiPSCs (Fig. 7c).

360         An advantage to using the SUMO1 and SUMO2 KGG mutants for site-level proteomics  
361 is that both paralogues leave the same Gly-Gly remnant on substrates after LysC digestion.  
362 Thus, site-specific SUMO preference can be compared. To date this has not been undertaken  
363 on a large scale and remains an important question in the SUMO field. The proteomic

364 experimental design allowed these comparisons at multiple stages of the purification process  
365 (Supplementary Fig. 13a-d): SUMO1/SUMO2 ratios from crude hiPSC extracts shows there to  
366 be few differences at the whole proteome level (0.03% significant - Supplementary Fig. 13e).  
367 There are also surprisingly few differences between NiNTA purifications from the two cell  
368 types (7.8% significant - Supplementary Fig. 13f). Exceptions include the well-documented  
369 SUMO1 substrate RanGAP1, along with TRIM24 and TRIM33 which all show similar levels of  
370 SUMO1 preference (Supplementary Fig. 13f). In contrast, over half of the GGK-containing  
371 peptides quantified showed large and significant difference between SUMO1 and SUMO2  
372 cells (Supplementary Fig. 13g). Extreme examples of SUMO1 preferential sites include  
373 RanGAP1 K524 and TRIM33 K776. Conversely, TRIM28 contains two of the most SUMO2-  
374 preferential sites at K507 and K779, and lysine 48 and 63 from ubiquitin are among the extreme  
375 SUMO2 acceptors. This was confirmed by Western blot analysis from NiNTA purifications  
376 (Supplementary Fig. 13h). Thus, when considering net modification of a protein, the bulk of  
377 SUMO modified proteins do not appear to display SUMO paralogue specificity, while this  
378 difference is clear at the site level.

379 STRING enrichment analysis of the 427 modified proteins created a network  
380 consisting of 3 clusters of proteins that could be broadly categorised as having functions in  
381 ribosome biogenesis, RNA splicing, and regulation of gene expression (Fig. 7d). Despite  
382 forming extensive protein networks (Fig. 7e and f), proteins involved in ribosome biogenesis  
383 and RNA splicing represented only approximately 5% of the total GGK peptide intensity (Fig.  
384 7b insert). The majority of the remainder have roles in transcription and chromatin structure  
385 or are closely linked to these functions (Fig. 7b insert and d). There is a prominent network of  
386 zinc-finger transcription factors, closely associated with TRIM28 (Fig. 7g) which contains many  
387 of the most heavily SUMOylated proteins identified, which play a key role in silencing  
388 retroviral elements. Histone proteins themselves, including H1, form a small cluster of SUMO  
389 substrates (Fig. 7h) in the centre of the gene regulation region of the whole network (Fig. 7d),  
390 and could potentially act as a direct link between SUMO and chromatin structure. The  
391 transcriptional regulators themselves form a bipolar network with the smaller sub-cluster  
392 consisting mainly of apparently weakly modified ribosomal proteins and the larger sub-cluster  
393 containing many heavily modified chromatin associated proteins (Fig. 7i). Strikingly, many  
394 members of chromatin remodelling complexes such as PRC2 (Fig. 7j), BAF (Fig. 7k) and NURD

395 (Fig. 7I) are among this group, potentially providing a link between SUMO and chromatin  
396 structure and remodelling.

397

## 398 **Discussion**

399 Our studies highlight an important role for SUMO in maintaining the pluripotent state  
400 of hiPSCs. Using a potent and highly specific inhibitor of the SUMO E1 ML792<sup>25</sup> we blocked *de*  
401 *nov*o SUMO modification and allowed endogenous SENPs to remove SUMO from previously  
402 modified proteins. In response to short-term SUMOylation inhibition ChiPS4 cells showed no  
403 loss of viability, but underwent clear morphological changes, losing markers of pluripotency  
404 (NANOG, KLF4, LINC-ROR), without displaying large-scale changes to the cellular proteome  
405 (Figs. 1 and 2). Upon SUMOylation inhibition ATAC-seq analysis demonstrated that sites  
406 previously occupied by SUMO (SUMO1 ChIP-seq) became more accessible (Fig. 3,  
407 Supplementary Fig. 3), indicating a role for SUMO in maintaining a compact chromatin  
408 environment. About 80% of these sites were associated with non-LTR and LTR  
409 retrotransposons (Fig. 3), while RNA-seq analysis indicated that a subset of HERV increased  
410 their transcription in response to SUMO modification inhibition (Fig. 6). ChIP-seq analysis  
411 indicated that the peak of SUMO located close to these HERVs also overlapped with the ChIP-  
412 seq derived locations of TRIM28, SETDB1 and CBX3<sup>31</sup>. These proteins along with SUMO have  
413 previously been shown to function in HERV silencing in mESCs<sup>29, 37</sup> and adult human cells<sup>38, 39</sup>  
414 and this is consistent with our proteomic analysis that indicates that all three of these proteins  
415 are heavily SUMO modified (Fig. 7). Moreover, TRIM28 co-repressor functions by interacting  
416 with DNA bound Kruppel type zinc finger proteins, which are also heavily SUMO modified in  
417 our proteomic studies. In fact, they form a large TRIM28 centric network of SUMO modified  
418 proteins that also includes the histone methyl transferase SETDB1. It is suggested that a  
419 number of developmental genes are repressed by TRIM28/KRAB-ZNFs through deposition of  
420 H3K9me3 and *de novo* DNA methylation of their promoter regions<sup>40</sup>, thus making  
421 TRIM28/ZNFs a crucial link in maintenance of pluripotency in human stem cells.

422 Several families of HERVs have been found to show stage specific expression in the  
423 preimplantation embryo and in hESCs *in vitro*<sup>41</sup>. These HERVs have been implicated in the  
424 maintenance of pluripotency in hESCs, are associated with the binding sites of pluripotency  
425 associated transcription factors (including OCT4, SOX2 and NANOG), and produce stage-  
426 specific lncRNAs that are required for the maintenance of the pluripotent state<sup>42-44</sup>.

427 Furthermore, HERV-H expression is dynamically regulated during transcription factor-  
428 mediated reprogramming and the acquisition of appropriate stage-specific expression of  
429 HERV-H is required for the re-establishment of pluripotency in hiPSCs<sup>45</sup>. Recently, HERVs have  
430 also been implicated in the regulation of Topologically Associating Domains (TADs) in hPSCs  
431 as deletion of HERV-H elements eliminates their corresponding boundaries and reduces the  
432 expression of upstream genes, while *de novo* insertion of HERV-H sequences can create new  
433 TAD boundaries<sup>46</sup>. These observations suggest that proper control of HERV expression is  
434 required for the maintenance of pluripotency in hESCs and hiPSCs, and our data suggest that  
435 SUMO modification may play a role in defining the HERVs that are expressed in these cells.

436 Consistent with this observation, the RNA-seq analysis revealed major changes in the  
437 expression of lncRNAs, but with limited changes to expression of protein coding genes in  
438 response to ML792 treatment (Fig. 4). However, a notable exception to this was the *PRAME*  
439 gene, which showed a massive increase in transcription after treatment with ML792.  
440 Increased transcription could be detected at the earliest time point analysed (4h) suggesting  
441 this was a direct effect of inhibiting SUMO modification (Fig. 5). An increase in *PRAME*  
442 expression was also observed when SUMO was removed from substrates by expression of the  
443 catalytic domain of SENP1. After removal of SUMO the chromatin around the *PRAME* locus  
444 becomes accessible and transcription of the gene increases over 10,000 fold, which would be  
445 to date, the clearest example of a protein coding gene that is negatively regulated by  
446 SUMOylation. Although SUMO has long been implicated in the repression of protein coding  
447 genes, most of the previous work has used artificial promoters<sup>33</sup> and data showing SUMO  
448 mediated regulation of endogenous genes is rather sparse.

449 Inspection of the *PRAME* gene locus indicates that the TSS of *PRAME* is adjacent to the  
450 TSS of *LL22NCO3-63E9.3* lncRNA gene that is transcribed in the opposite direction. While  
451 SUMOylation inhibition increases transcription of the *PRAME* gene it does not lead to  
452 transcription of the neighbouring gene. CHIP-seq data indicates that in absence of ML792  
453 treatment a peak of SUMO is present over the TSS of *PRAME* and overlaps with CHIP-seq  
454 peaks for cohesin, CTCF, ZNF143 and RNA pol II. Our SUMO site proteomic data indicate that  
455 CTCF is heavily modified by SUMO as is ZNF143. Cohesin, CTCF and ZNF143 are all associated  
456 with maintaining higher order chromatin structure<sup>47-49</sup> particularly loops or TADs and  
457 supports the hypothesis that silencing of the *PRAME* gene is maintained by the influence of  
458 SUMO on higher order chromatin structure. ZNF143 is thought to control transcription from

459 bidirectional promoters<sup>50</sup> and regulate the density of promoter-proximal paused RNA  
460 polymerase<sup>51</sup>, which has been associated with silenced genes that can be rapidly activated.  
461 The co-location of RNA pol II with cohesion, CTCF, ZNF143 and SUMO is suggestive of such a  
462 scenario at the PRAME gene locus. This idea is supported by the observation that the most  
463 highly SUMO modified protein in our proteomic analysis is GTF2i/TFII-I, is a component of the  
464 RNA pol II transcriptional complex. Thus, while transcriptional repression of the PRAME gene  
465 also involves SUMO, it appears to be mediated by a very different mechanism from silencing  
466 of HERV genes as TRIM28, SETDB1 and CBX that are associated with HERV silencing, are not  
467 present at the PRAME locus. While our data at HERV entities are consistent with a recent  
468 SUMO proteomic analysis in mESCs<sup>52</sup>, the tight SUMO mediated regulation of the PRAME  
469 locus appears to be specific to hPSCs. Although the *PRAME* gene is frequently over-expressed  
470 in tumours<sup>53,54</sup> it appears to play a role in the differentiation of hPSCs into mesenchymal stem  
471 cells<sup>55</sup>. PRAME appears to function as the substrate adapter of a Cul2 E3 ubiquitin ligase that  
472 is targeted to chromatin and associates with active NFY promoters<sup>56</sup>. While the targets of  
473 PRAME mediated ubiquitination have yet to be identified it has been shown to associate with  
474 the highly conserved EKC/KEOPS complex on chromatin<sup>57</sup>.

475         Analysis of the SUMO proteome of ChiPS4 cells shows that well-defined groups of  
476 protein are modified. Aside from the TRIM28/ZNF network mentioned above, proteins  
477 involved in “ribosome biogenesis” and “splicing” are SUMO modified and this likely impacts  
478 on the normal growth and self-renewal of the hiPSCs. However, the largest network of  
479 proteins falls into the category of ‘negative regulation of transcription’ with many chromatin  
480 remodellers, chromatin modification and DNA modification enzymes identified as SUMO  
481 substrates. The increases in transcription observed after SUMO modification inhibition  
482 indicate that SUMO modification plays an important role in maintaining pluripotency of  
483 hiPSCs by repressing genes that either disrupt pluripotency or drive differentiation.

484



## 485 **Acknowledgements**

486 We thank Linnan Shen for the kind gift of purified Tn5 transposase, Adel Ibrahim for GNB-  
487 SENP1 DNA and Alwyn Dady for the paPX1 vector. We would like to acknowledge the  
488 invaluable help from various facilities at the University of Dundee: Flow Cytometry and Cell  
489 Sorting, National Phenotypic Screening Centre, Dundee Imaging Facility, Human Pluripotent  
490 Stem Cell Facility. This work was supported by an Investigator Award from Wellcome  
491 (217196/Z/19/Z) and a Programme grant from Cancer Research UK (C434/A21747) to RTH.  
492 ML is part of the Ubi-Code European Training Network that received funding from the  
493 European Union Horizon 2020 research and innovation programme under the Marie Curie  
494 grant agreement No. 765445. BM was supported by the European Union's Horizon 2020  
495 research and innovation programme under the Marie Skłodowska-Curie grant agreement No.  
496 704989. JEW was supported by a Marie Skłodowska-Curie Actions Individual Fellowship from  
497 the European Commission grant agreement No. 625253. TOH and NW were supported by the  
498 MRC grant (MR/S021647/1).

499

## 500 **Author Contributions**

501 BM cloned expression vectors, generated cell lines, designed and performed most  
502 experiments using hiPSCs and interpreted data. JEW contributed to the initial design of the  
503 research and performed SUMO1 ChIP-seq. ML performed NiNTA purifications and Western  
504 blotting analyses. LD was in charge of cell culture, cell line derivation and quality control and  
505 was consulted over experimental design. RTH generated capped and polyadenylated RNAs.  
506 MHT consulted over proteomic experimental design, acquired and processed MS data and  
507 conducted bioinformatic and statistical analyses. NW performed the ATAC-seq experiments,  
508 processed and analysed the data. MG performed the bioinformatic analysis of RNA-seq and  
509 ChIP-seq data sets. BM, MG, NW, MHT, LD, TOH and RTH contributed to data analysis. BM,  
510 RTH, MHT and LD wrote the paper. RTH conceived the project.

511

## 512 **Data availability**

513 Data underlying all Figures and Supplementary Figures are available in the source data file.  
514 The mass spectrometry proteomics data have been deposited to the ProteomeXchange  
515 Consortium via the PRIDE<sup>58</sup> partner repository with following datasets identifiers:  
516 PXD023241, PXD023257. Source data for ATAC-seq, RNA-seq and SUMO1 ChIP-seq are

517 available from EBI at <https://www.ebi.ac.uk/> under accession numbers: E-MTAB-9961, E-  
518 MTAB-9962 and E-MTAB-9960 respectively. All other data are available from the  
519 corresponding authors on reasonable request.

520

#### 521 **Supplementary data files**

522 **Supplementary data file 1.** Summary of the quantitative data from the proteomics  
523 experiment to study changes to the cellular proteome during ML792 treatment of ChiPS4  
524 cells.

525 **Supplementary data file 2.** Summary of the quantitative data from the proteomics  
526 experiment to study differences in the cellular proteome among wild type ChiPS4 cells and  
527 cells expressing 6His-SUMO1-KGG-mCherry or 6His-SUMO2-KGG-mCherry.

528 **Supplementary data file 3.** Summary of the quantitative data from the proteomics  
529 experiment to identify SUMO1 and SUMO2 targets from ChiPS4 cells.

530 **Supplementary data file 4.** Merge of the RNA-seq and proteomic data for ChiPS4 cells treated  
531 with ML792 for 0h, 24h and 48h.

532 **Supplementary data file 5.** ATAC-seq data file for ChiPS4 cells treated with ML792 for 4h, 8h,  
533 24h and 48h.

534 **Supplementary data file 6.** RNA-seq edger data file for ChiPS4 cells treated with ML792 for  
535 4h, 8h, 24h and 48h.

536

#### 537 **Code availability**

538 Code used for data analysis is included in the source data file and can be found at  
539 [https://github.com/bartongroup/MG\\_SumoDiff2](https://github.com/bartongroup/MG_SumoDiff2).

540

#### 541 **Author Information**

542 Correspondence and requests for materials should be addressed to R.T.H  
543 (r.t.hay@dundee.ac.uk)

544

545 **Methods**

546 **Antibodies and inhibitors.** Rabbit antibodies against TRIM28 (4124S, 4123S), CTCF (3418S),  
547 OCT4A (2890S), SOX2 (23064S), NANOG (3580S), KLF4 (12173S) and mouse antibodies against  
548 TRA-1-60 (4746T), TRA-1-81(4745T), SSEA-4 (4755T) and SMA (D4K9N) were from Cell  
549 Signalling Technology. The anti-SALL4 (ab29112), anti-TRIM24 (ab70560), anti-NOP58  
550 (ab155556), anti-PRAME (ab219650), anti-TRIM24 (ab70560), anti-NESTIN (ab196908) were  
551 from Abcam. Mouse antibody against  $\alpha$ -Tubulin was from Bethyl Laboratories and mouse  
552 anti-LaminA/C antibody was from Sigma (SAB4200236), rabbit anti-mCherry (PA5-34974),  
553 rabbit anti-PRAME (PA5-83761), rabbit anti-TRIM33 (PA5-82152) and mouse anti-HIS (34650)  
554 were from Invitrogen and Qiagen respectively. Anti-Cytokeratin17 was a gift from R.  
555 Hickerson (University of Dundee). Sheep antibodies against SUMO1, SUMO2, and SENP1<sup>27</sup>  
556 and chicken antibodies against PML<sup>59</sup> were generated in-house. Secondary antibodies  
557 conjugated with HRP and Alexa fluorophores were from Sigma and Invitrogen, respectively.  
558 MG132 (474787) and N-ethylmaleimide (E3876) were from Sigma Aldrich. ML792 was from  
559 UbiQ. Protease Inhibitor cocktail (11836170001) was from Roche. Propidium iodide, Cy5 Cell  
560 Mask and DAPI were from Life Technologies.

561

562 **Cloning.** SUMO1-KGG-mCherry and SUMO2-KGG-mCherry PiggyBac expression vectors were  
563 generated by GATEWAY cloning. Briefly, SUMO1, SUMO2 and mCherry fragments were PCR  
564 amplified using the following resources: 6His SUMO1 T95K (300nt) from pSCAI88 and 6His  
565 SUMO2 T90K (300nt) from pSCAI89 with a common forwards primer (5'-  
566 CACCatgcatcatcatcatcatcatgct-3') and set of specific mCherry fusing primers (5'-  
567 TCACCATACCCCCTTTTGTTCCTG-3' and 5'-TCACCATACCTCCCTTCTGCTGCT-3'); mCherry from  
568 pRHAI4 CMV-OsTIR1-mCherry2-PURO (700 nt) with a set of common overlapping oligos (5'-  
569 GGTATGGTGAGCAAGGGCG-3' and 5'-TTATTACTTGACAGCTCGTCCATG-3'). Subsequently,  
570 PCR fragments were fused together using overlap extension PCR and TOPO cloned into  
571 pENTR™/D-TOPO™ (Invitrogen) and verified by DNA sequencing. The assembled SUMO1-  
572 KGG-mCherry and SUMO2-KGG-mCherry sequences were then sub-cloned from the pENTR  
573 vector into the destination PiggyBac GATEWAY expression vector paPX1 using LR clonase II  
574 (ThermoFisher Scientific).

575 **mRNA synthesis and purification.** The catalytic domain of the SUMO specific protease  
576 SENP1<sup>60</sup> was fused at its N-terminus to a nano body directed against GFP<sup>61</sup> to form GNB-  
577 SENP1. The DNA encoding GNB-SENP1 was amplified by PCR using a 5' primer containing the  
578 sequence of the T7 RNA polymerase promoter. Amplified DNA was purified on a MinElute Gel  
579 Extraction kit (Qiagen). 4 µg of the eluted DNA was used as template for the production of  
580 capped and poly adenylated mRNA by in vitro transcription using an mMessage mMachine T7  
581 Ultra kit (ThermoFisher) as described by the manufacturer. RNA was purified using a  
582 MegaClear kit (ThermoFisher) as described. Purified RNA was quantified by NanoDrop and  
583 analysed on a TapeStation (Agilent).

584 **Human Induced pluripotent stem cells (hiPSCs) culture and transfection protocols.** Human  
585 ESC lines (SA121 and SA181) were obtained from Cellartis / Takara Bio Europe. All work with  
586 hESCs was approved by the UK Stem cell bank steering committee (Approval reference:  
587 SCSC17-14). Human iPSC lines were obtained from Cellartis / Takara Bio Europe (ChiPS4) or  
588 the HipSci consortium (bubh3, oaqd3, ueah1 and wibj2). Cell lines were maintained in TESR  
589 medium<sup>62</sup> containing FGF2 (PeproTech, 30 ng/ml) and noggin (PeproTech, 10 ng/ml) on  
590 growth factor reduced geltrex basement membrane extract (Life Technologies, 10 µg/cm<sup>2</sup>)  
591 coated dishes at 37°C in a humidified atmosphere of 5% CO<sub>2</sub> in air. Cells were routinely  
592 passaged twice a week as single cells using TrypLE select (Life Technologies) and replated in  
593 TESR medium that was further supplemented with the Rho kinase inhibitor Y27632 (Tocris,  
594 10 µM). Twenty four hours after replating Y27632 was removed from the culture medium. To  
595 make SUMO1-KGG-mCherry and SUMO2-KGG-mCherry expressing stable cell lines ChiPS4  
596 cells were transfected using a Neon electroporation system (Thermo Fisher Scientific) using  
597 10 µl tips. Briefly, ChiPS4 cells were dispersed to single cells as described above then 1x10<sup>6</sup>  
598 cells were collected by centrifugation at 300xg for 2 minutes and resuspended in 11 µl of  
599 electroporation buffer R containing 1 µg of either paPX1-SUMO1-KGG-mCherry or paPX1-  
600 SUMO2-KGG-mCherry PiggyBac expression vectors along with 0.2 µg of Super PiggyBac  
601 transposase (System Biosciences). Electroporation was performed at 1150 V, 1 pulse, 30 mSec  
602 and cells plated in mTESR containing Y27632. 5 days after electroporation, mCherry positive  
603 cells were selected by fluorescence activated cell sorting (FACS) using an SH800 cell sorter  
604 (Sony). Monoclonal cell lines were prepared from the bulk sorted population by plating at low  
605 density on geltrex coated dishes and individual clones picked using 3.2 mm cloning discs

606 (Sigma Aldrich) soaked in TrypLE select. Cell lines were then expanded and analysed to check  
607 for expression of mCherry and His-SUMO1/2. Transfection of SENP1 mRNAs was performed  
608 using the same protocol.

609

#### 610 ***In vitro* differentiation assay**

611 For assessment of pluripotency *in vitro*,  $1 \times 10^4$  hiPSCs were seeded into the wells of v-  
612 bottomed 96 well plates in TEHR medium supplemented with Y27632 and centrifuged at  
613 300xg for 5 minutes. After 48 hours the resultant embryoid bodies were picked from the v-  
614 bottom plates using a pipette and seeded on gelatin coated dishes in knockout DMEM  
615 medium supplemented with 20% knockout serum replacement, 1x non-essential amino acids,  
616 1x glutamax, 100  $\mu$ M 2-mercaptoethanol. The medium was changed every 3 - 4 days and the  
617 cells were fixed with 4% formaldehyde and the expression of germ layer markers analysed by  
618 IF on day 20 of differentiation.

619

620 **Flow cytometry for cell cycle assessment and pluripotency markers.** For cell cycle analysis  
621 and staining for pluripotency markers ChiPS4 cells were harvested using standard procedures,  
622 washed and fixed with ice cold 70% ethanol or 4% formaldehyde for the analysis of cell cycle  
623 or NANOG staining respectively. Next cells were stained with propidium iodide or anti-  
624 NANOG primary antibody, followed by Alexa 488 conjugated secondary antibody and  
625 analysed by flow cytometry using a Canto analyser (Becton Dickson). Data was then analysed  
626 using FlowJo 10.

627

628 **Immunofluorescence, cell painting assay and high content microscopy.** For IF assays ChiPS4  
629 cells were seeded on  $\mu$ -Slide 8 Well (ibidi) or 96 well plates suitable for high content  
630 microscopy (Nunc). Standard IF procedure was used where appropriate. Briefly, following  
631 treatments cells were washed with PBS, fixed with 4% formaldehyde, blocked in 5% BSA in  
632 PBS-T and incubated with primary and Alexa conjugated secondary antibodies and co-stained  
633 with DAPI and/or Cy5 Cell Mask (Life Technologies). Cell painting was performed as  
634 described<sup>26</sup>. Imaging and subsequent analysis was performed using INCell Analyzer systems  
635 (GE Healthcare) and Spotfire (Tibco).

636

637 **Protein sample preparation and Western blotting (WB).** ChiPS4 were maintained in a stable  
638 culture as described before and treated with inhibitors for a stated time and dose, usually  
639 400nM ML792 was used for 24h or 48h. For WB cells were washed with PBS +/- and directly  
640 lysed in an appropriate volume of 2x Laemmli buffer (approximately 200  $\mu$ l of buffer was used  
641 per  $0.5 \times 10^6$  cells) (LD; [4% SDS; 20% Glycerol; 120mM 1 M Tris-Cl (pH 6.8); 0.02% w/v  
642 bromophenol blue]) and subsequently sonicated using Bioruptor Twin (Diagenode). Protein  
643 content was assessed using BCA protein assay (ThermoFisher Scientific) and for most  
644 purposes 15  $\mu$ g of total protein was loaded per lane on SDS-PAGE gel (NuPage 4-12%  
645 polyacrylamide, Bis-Tris with MOPS buffer). Proteins were transferred to PVDF membrane  
646 using iBlot™ 2 Gel Transfer Device (Invitrogen). Membranes were blocked for 1h in 5% milk in  
647 TBS-T and incubated overnight with primary antibodies and 1h with secondary HRP  
648 conjugated antibodies before being developed using enhanced chemiluminescence  
649 (ThermoFisher Scientific).

650

651 **NiNTA purification.** Cells were washed with PBS and scraped in PBS containing 1mM N-  
652 ethylmaleimide. The cells were then collected by centrifugation at 300 xg for 5 minutes and  
653 the pellets weighed. An aliquot of the cells was lysed in 1.2x NuPage sample buffer  
654 (ThermoFisher Scientific) for analysis by Western blotting. The remaining cell pellets  
655 (approximately 1 g) were lysed with 5x the pellet weight of lysis buffer (6 M guanidine-HCl,  
656 100 mM sodium phosphate buffer (pH 8.0), 10 mM Tris-HCl (pH 8.0), 10 mM imidazole and 5  
657 mM 2-mercaptoethanol). DNA was sheared by sonication using a probe sonicator (3min, 35%  
658 amplitude, 20sec pulses, 20sec intervals on ice and the samples centrifuged at 4000 rpm for  
659 15min at 4°C to remove insoluble material). The protein concentration of the lysate was  
660 determined using BCA assay and 6.5mg of total protein from each sample was then incubated  
661 overnight at 4°C with 50  $\mu$ l of packed pre-equilibrated Ni-NTA agarose beads. After the  
662 overnight incubation the supernatant was removed and the beads were washed once with  
663 10 resin volumes of lysis buffer, followed by 1 wash with 10 resin volumes of 8 M urea, 100  
664 mM sodium phosphate buffer (pH 8.0), 10 mM Tris-HCl (pH 8.0), 10 mM imidazole and 5 mM  
665 2-mercaptoethanol, and then 6 washes with 10 resin volumes of 8 M urea, 100 mM sodium  
666 phosphate buffer (pH 6.3), 10 mM Tris-HCl (pH 8.0), 10 mM imidazole and 5 mM 2-

667 mercaptoethanol. Proteins were eluted from Ni-NTA agarose beads with 125  $\mu$ L 1.2x NuPAGE  
668 sample buffer for SDS-PAGE.

669

### 670 **Mass Spectrometry based proteomics and quantitative data analysis**

671 Three proteomic experiments are described in this study;

672 (1) Changes in total proteome of ChiPS4 cells during ML792 treatment

673 ChiPS4 cells were either DMSO treated (0 hours condition), or treated with 400nM ML792 for  
674 24 hours or 48 hours. Four replicates of each condition were prepared. Crude cell extracts  
675 were made to a protein concentration of between 1 and 2 mg/ml by addition of 1.2x NuPAGE  
676 sample buffer to PBS washed cells followed by sonication. For each replicate 25 $\mu$ g protein  
677 was fractionated by SDS-PAGE (NuPage 10% polyacrylamide, Bis-Tris with MOPS buffer—  
678 Invitrogen) and stained with Coomassie blue. Each lane was excised into four roughly equally  
679 sized slices and peptides were extracted by tryptic digestion<sup>63</sup> including alkylation with  
680 chloroacetamide. Peptides were resuspended in 35  $\mu$ L 0.1% TFA 0.5% acetic acid and 10 $\mu$ L of  
681 each sample was analysed by LC-MS/MS. This was performed using a Q Exactive mass  
682 spectrometer (Thermo Scientific) coupled to an EASY-nLC 1000 liquid chromatography system  
683 (Thermo Scientific), using an EASY-Spray ion source (Thermo Scientific) running a 75  $\mu$ m x 500  
684 mm EASY-Spray column at 45 $^{\circ}$ C. A 240 minute elution gradient with a top 10 data-dependent  
685 method was applied. Full scan spectra (m/z 300–1800) were acquired with resolution R =  
686 70,000 at m/z 200 (after accumulation to a target value of 1,000,000 ions with maximum  
687 injection time of 20 ms). The 10 most intense ions were fragmented by HCD and measured  
688 with a resolution of R = 17,500 at m/z 200 (target value of 500,000 ions and maximum  
689 injection time of 60 ms) and intensity threshold of  $2.1 \times 10^4$ . Peptide match was set to  
690 'preferred', a 40 second dynamic exclusion list was applied and ions were ignored if they had  
691 unassigned charge state 1, 8 or >8. Data analysis used MaxQuant version 1.6.1.0<sup>64</sup>. Default  
692 settings were used except the match between runs option was enabled, which matched  
693 identified peaks among slices from the same position in the gel as well as one slice higher or  
694 lower. The uniprot human proteome database (downloaded 24/02/2015 - 73920 entries)  
695 digested with Trypsin/P was used as search space. LFQ intensities were required for each slice  
696 but LFQ normalization was switched off. Manual LFQ normalization was done by calculating  
697 the LFQ ratio for each protein in each slice compared to the average LFQ for the same protein  
698 across all equivalent slices in the other lanes. This was done only for proteins with LFQ

699 intensities reported in all 12 equivalent slices. The median protein Slice LFQ/Average LFQ ratio  
700 was used to normalize all protein LFQ values for each slice. The final protein LFQ intensity per  
701 lane (and therefore sample) was calculated by the sum of LFQ values for that protein intensity  
702 in all four slices. Downstream data processing used Perseus v1.6.1.1<sup>65</sup>. Proteins were only  
703 carried forward if an LFQ intensity was reported in all four replicates of at least one condition.  
704 Zero intensity values were replaced from log2 transformed data (default settings) and outliers  
705 were defined by 5% FDR from Student's t-test using an S0 value of 0.1. A summary of these  
706 data can be found in Supplementary Data File 1

707 (2) Characterisation of ChiPS4 cells stably expressing 6His-SUMO1-KGG-mCherry and  
708 6His-SUMO2-KGG-mCherry.

709 Crude cell extracts were prepared in triplicate from ChiPS4 cells, ChiPS4-6His-SUMO1-KGG-  
710 mCherry and ChiPS4-SUMO2-KGG-mCherry cells and fractionated by SDS-PAGE as described  
711 above. In an almost identical manner gels were sectioned into four slices per lane, tryptic  
712 peptides prepared, peptides analysed by LC-MS/MS, and the resultant raw data processed by  
713 MaxQuant. The only exceptions being the inclusion of a second sequence database containing  
714 the two 6His-SUMO-KGG-mCherry constructs, and the use of MaxQuant LFQ normalization.  
715 Two MaxQuant runs were performed; the first aggregating all slices per lane into a single  
716 output ("by lane"), and the second considering each slice separately ("by slice"). The former  
717 was used to determine cell-specific changes in protein abundance from the proteinGroups.txt  
718 file, and the latter used the peptides.txt file to monitor differences in abundance of SUMO-  
719 specific peptides between samples, to infer overexpression levels. For the whole cell  
720 proteome change analysis only proteins with data in all three replicates of at least one  
721 condition were carried forward. In Perseus zero intensity values were replaced from log2  
722 transformed data (default settings) and outliers were defined by 5% FDR from Student's t-  
723 test using an S0 value of 0.1. A summary of these data can be found in Supplementary Data  
724 File 2.

725 (3) Identification of SUMO1 and SUMO2 modified proteins from ChiPS4 cells.

726 Two repeats of this experiment were performed using approximately  $0.5 \times 10^8$  cells of ChiPS4-  
727 6HisSUMO1-KGG-mCherry and ChiPS4-6HisSUMO2-KGG-mCherry per replicate. Samples  
728 were taken at different steps of the protocol to assess different fractions. These were; crude  
729 cell extracts, NiNTA column elutions and GlyGly-K immunoprecipitations. The last being the  
730 source of SUMO-substrate branched peptides. The whole procedure was carried out as



731 described previously<sup>66</sup>. In brief, crude cell lysates were prepared of which approximately 100  
732  $\mu\text{g}$  was retained for whole proteome analysis as described for experiments 1 and 2 above.  
733 The remaining lysate ( $\sim 20$  mg protein) was used for NiNTA chromatographic enrichment of  
734 6His-SUMO conjugates. Elutions from the NiNTA columns were digested consecutively with  
735 LysC then GluC, of which 7% of each was retained for proteomic analysis and the remainder  
736 for GlyGly-K immunoprecipitation. The final enriched fractions of LysC and LysC/GluC GG-K  
737 peptides were resuspended in a volume of 20  $\mu\text{l}$  for proteomic analysis. Peptides from whole  
738 cell extracts were analysed once by LC-MS/MS using the same system and settings as  
739 described for experiments 1 and 2 above except a 180 minute gradient was used with a top  
740 12 data dependent method. NiNTA elution peptides were analysed identically except a top  
741 10 data dependent method was employed and maximum MS/MS fill time was increased to  
742 120ms. GG-K immunoprecipitated peptides were analysed twice. Firstly, 4  $\mu\text{l}$  was fractionated  
743 over a 90 minute gradient and analysed using a top 5 data dependent method with a  
744 maximum MS/MS fill time of 200ms. Secondly, 11  $\mu\text{l}$  of sample was fractionated over a 150  
745 minute gradient and analysed using a top 3 method with a maximum MS/MS injection time  
746 of 500ms. Data from WCE and NiNTA elutions were processed together in MaxQuant using  
747 Trypsin/P enzyme specificity (2 missed cleavages) for WCE samples and LysC (3 missed  
748 cleavages), or LysC+GluC\_D/E (considering cleavage after D or E and 8 missed cleavages) for  
749 NiNTA elutions. GlyGly (K) and phospho (STY) modifications were selected. The human  
750 database and sequences of the two exogenous 6His-SUMO-KGG-mCherry constructs  
751 described above were used as search space. In all cases every raw file was treated as a  
752 separate 'experiment' in the design template such that protein or peptide intensities in each  
753 peptide sample were reported, allowing for manual normalization. Matching between runs  
754 was allowed but only for peptide samples from the same cellular fraction (WCE, NiNTA elution  
755 or GG-K IP), the same or adjacent gel slice, the same protease and the same LC elution  
756 gradient. For example, spectra from adjacent gel slices in the WCE fraction across all lanes  
757 were matched, and spectra from all GG-K IPs that were digested by the same enzymes were  
758 matched. Normalization followed a similar method as described above where 'equivalent'  
759 peptide samples (i.e. those from the same gel slice or 'equivalent' peptide samples) from  
760 different replicates were compared with one another. Manual normalization used a similar  
761 method described above. For each protein or peptide common to all equivalent peptide

762 samples the ratio of intensity in that sample to the average across all equivalent samples was  
763 calculated. The median of that ratio of was used to normalize all protein or peptide intensities  
764 for each sample. The final protein or peptide intensity per replicate was calculated by the sum  
765 of all normalized intensities in samples derived from that replicate. Importantly, peptide  
766 samples derived from SUMO1 and SUMO2 cells were considered equivalent for normalization  
767 purposes, which assumes largely similar abundances of proteins or peptides across cell types.  
768 Zero intensity values were replaced from log2 transformed data (default settings) and outliers  
769 were defined by 5% FDR from Student's t-test using an S0 value of 0.1. A summary of these  
770 data can be found in Supplementary Data File 3.

771

772 **Bioinformatic analysis of the SUMO site proteomics.** 429 proteins identified with at least one  
773 SUMO1 or SUMO2 modification site were uploaded to STRING<sup>67</sup> for network analysis. Only  
774 proteins associated by a minimum STRING interaction score of 0.7 (high confidence) were  
775 included in the final network. Disconnected nodes were removed. Selected groups of  
776 functionally related proteins were resubmitted to STRING to create smaller sub-networks.  
777 These were visualised in Cytoscape v 3.7.2<sup>68</sup> allowing the graphical display of numbers of sites  
778 identified and total GG-K peptide intensity into the protein networks.

779

780 **SUMO1 ChIP-seq.** Cells were dispersed with TrypLE select as previously described, cross-  
781 linked with 1% formaldehyde for 10 min, then quenched for 5 min with 125 mM glycine at  
782 room temperature. Fixed cells were washed twice with ice-cold PBS, frozen in liquid nitrogen  
783 and stored at -80°C. Frozen cell pellets were thawed on ice and resuspended in lysis buffer (5  
784 mM PIPES pH 8; 85 mM KCl; 0.5% Igepal CA-630) supplemented with complete Protease  
785 Inhibitor Cocktail without EDTA (Roche) and 20mM iodoacetamide. Nuclear extraction by  
786 sonication<sup>69</sup> was performed using a Bioruptor Twin sonicator on low power with 2-3 cycles  
787 (15 secs on and 15secs off). Nuclei were collected by centrifugation at 2500 xg for 5 min at  
788 4°C, washed once in 1 ml of lysis buffer, and visualised by microscopy to confirm the release  
789 of nuclei from cells. Nuclei were resuspended to a final concentration of 2 x 10<sup>7</sup> nuclei/ml in  
790 nuclei lysis buffer (50 mM Tris-HCl, pH 8.0, 10 mM EDTA, 1% SDS) plus protease inhibitors.  
791 To shear chromatin to fragments of approximately 500bp (range 100–800 bp) in length,  
792 samples were sonicated in a volume of 300 µl for 20 cycles (10 min total sonication time)  
793 using the Bioruptor Twin on high power. Sonicated lysates were then clarified by

794 centrifugation for 15000 xg for 10 minutes at 4°C. Input DNA was purified using an IPURE kit  
795 (Diagenode) according to the manufacturer's instructions. For each IP, 15- $\mu$ L aliquots of  
796 Protein G Dynabeads (10 mg/mL, Dynal) were washed in 500  $\mu$ l PBS then 5  $\mu$ g of sheep anti-  
797 SUMO-1 antibody was bound in 1ml of PBS containing 0.1% IgG free BSA and protease  
798 inhibitor cocktail for 1 - 4hr at 4°C with agitation. 25  $\mu$ g of chromatin was premixed with 8  
799 volumes of IP dilution buffer (20 mM Tris-HCl at pH 7.6, 0.625% Triton X-100, 182.5 mM NaCl)  
800 then added to the antibody on Dynabeads and left overnight at 4°C with agitation. Samples  
801 were then centrifuged for 1min at 500 xg, placed in a magnetic separation rack, and the beads  
802 were washed once with 1ml of IP wash buffer 1 (20 mM Tris-HCl at pH 8.0, 2 mM EDTA, 0.2%  
803 SDS, 0.5% Triton X-100, and 150 mM NaCl), once with 1ml of IP wash buffer 2 (20 mM Tris-  
804 HCl at pH 8.0, 2 mM EDTA, 0.2% SDS, 0.5% Triton X-100, and 500 mM NaCl) once with 1ml of  
805 IP wash buffer 3 (0.25 M LiCl, 1 % NP-40, 1 % deoxycholate, 1 mM EDTA, 10 mM Tris pH 8.1),  
806 and once with 1ml of TE buffer (1 mM EDTA, 10 mM Tris pH 8.1). Immunoprecipitated  
807 material was eluted from the beads and DNA purified using Diagenode IPURE kit according to  
808 manufacturer's protocol. The resulting DNA was used for Illumina library preparations.

809 Libraries from SUMO1 ChIP DNA and Input DNA were prepared using the NEBNext  
810 Ultra II DNA library prep kit with sample purification beads according to the manufacturer's  
811 instruction. Barcoding of the samples was performed using NEBNext Multiplex oligos Index  
812 set 1 and 2. Chromatin size distribution was measured on Agilent TapeStation and sample  
813 concentration quantified with Qubit® dsDNA HS (High Sensitivity) Assay Kit and Qubit®  
814 Fluorometer (ThermoFisher Scientific). The sequencing was performed on a NextSeq 500  
815 (Illumina): paired end, high throughput 2x75bp run. Chromatin input seq and SUMO1 ChIP-  
816 seq data analysis were performed as described here:  
817 <http://www.compbio.dundee.ac.uk/user/mgierlinski/sumodiff/doc/analysis.html>. Briefly,  
818 ChIP-seq reads were mapped to human genome reference GRCh38 (repeat-masker filtered)  
819 using bwa version 0.7.15. Then, peak calling was done with MACS2 version 2.1.0. Data was  
820 also mapped and processed for the hg17 genome reference to be used for overlaps and  
821 alignments with ENCODE TF database.

822

823 **ATAC-seq.** ATAC-seq libraries were generated following the Omni-ATAC protocol<sup>70</sup> without  
824 enrichment for viable cells. Amplified barcoded DNA fragments were purified using NEB  
825 Monarch PCR purification kit. The DNA concentration was measured with the Qubit High

826 Sensitivity assay (ThermoFisher) and the fragment sizes were determined on the TapeStation  
827 Bioanalyser (Agilent). Samples were pooled in equimolar ratios. The pooled library was  
828 subjected to a dual size selection using Promega Pronex beads 1.2/0.4 x beads:sample ratio  
829 to enrich for fragments between 180 bp and 800 bp. Multiplexed libraries were sequenced  
830 with 2x150 bp paired-end reads by Novogene with ~20 Mio reads per sample using NovaSeq  
831 6000 S4 (Illumina). Fastq files were trimmed using trimmomatic-0.36 (CROP: 66) and aligned  
832 to the human GRCh38 genome using bowtie2 with the parameter  $-X$  1000. Peaks were called  
833 using MACS2 *callpeak* function with the following parameters:  $-t$  "\$1".bam  $-f$  BAMPE  $-n$   
834 "\$1"\_MACS  $-g$  2.7e9  $-q$  0.05  $--broad$   $-B$ . Differential peaks were obtained using *DiffBind*,  
835 doing pair-wise comparison of two time points to control DMSO treated samples. When  
836 performing *dba.count*, a *minOverlap* was set to 3, requiring a peak to be observed in at least  
837 3 datasets in order to be retained. Differential peaks were called using the *edgeR* method  
838 during *dba.analyze*. Of the differentially called peaks, a second filtering step was performed  
839 to retain only peaks that met an  $FDR < 0.00001$  and a *scores.fold*  $> 0.58$  (equal to a fold change  
840  $> 1.5$ ). A summary of these data can be found in Supplementary Data File 5. Non-changing  
841 peaks were obtained from the *DiffBind* consensus peak set, with all differential peaks  
842 removed (all timepoints, no extra thresholding). The non-differential peaks were randomly  
843 subsampled to the same sample-size as differential peaks. The *bedtools intersect* function was  
844 used to call overlap of ATAC-seq peaks with CHIP-seq data and genomic regions were  
845 annotated using the *annotatePeaks.pl* command from the *HOMER* software.

846

847 **RNA preparation and real-time quantitative PCR (RT-qPCR).** Total RNA was extracted using  
848 RNeasy Mini Kit (Qiagen) and treated with the on-column RNase-Free DNase Set (Qiagen)  
849 according to the manufacturer's instructions. RNA concentration was then measured using  
850 NanoDrop and 1 $\mu$ g of total RNA per sample was subsequently used to perform a two-step  
851 reverse transcription polymerase chain reaction (RT-PCR) using random hexamers and First  
852 Strand cDNA Synthesis Kit (ThermoFisher Scientific). Each qPCR reaction contained PerfeCTa  
853 SYBR Green FastMix ROX (Quantabio), forward and reverse primer mix (200 nM final  
854 concentration) and 6 ng of analysed cDNA and was set up in triplicates in MicroAmp™ Fast  
855 Optical 96-Well or 384-Well Reaction Plates with Barcodes (Applied Biosystems™). The  
856 sequences of primers used were as follows: NANOG (hNANOG\_FOR624  
857 ACAGGTGAAGACCTGGTTCC; hNANOG\_REV722 GAGGCCTTCTGCGTCACA), SOX2 (hSOX2

858 \_FOR907 TGGACAGTTACGCGCACAT; hSOX2\_REV1121 CGAGTAGGACATGCTGTAGGT), OCT4A  
859 (hOCT4A\_FOR825 CCCACACTGCAGCAGATCA and hOCT4A\_REV1064  
860 ACCACACTCGGACCACATCC), KLF4 (hKLF4\_FOR1630 GGGCCCAATTACCCATCCTT and  
861 hKLF4\_REV1706 GGCATGAGCTCTTGTAATGG), TBP (hTBP\_FOR896  
862 TGTGCTCACCCACCAACAAT; hTBP\_REV1013 TGCTCTGACTTTAGCACCTGTT), PRAME  
863 (hPRAME\_F1661 TACCTGGAAGCTACCCACCT and hPRAME\_R1892  
864 GTGCCTGAGCAACTGATCCA). Data were collected using QuantStudio™ 6 Flex Real-Time PCR  
865 Instrument and analysed using a corresponding software (Applied Biosystems™). Relative  
866 amounts of specifically amplified cDNA were calculated using TBP amplicons as normalizers.  
867

868 **RNA-seq.** RNA samples were collected and prepared as for standard RNA extraction  
869 procedure. Samples were quality controlled using Qubit (Thermo Fisher Scientific) and  
870 TapeStation (Agilent) and sent for further analysis to Novogen, who prepared the Illumina  
871 Library using NEB Next® Ultra™ RNA Library Prep Kit. These libraries were sequenced using  
872 the Illumina NovaSeq 6000 S4 (Illumina PE150, Q30 ≥ 80% delivering 6G raw data per sample).  
873 Following data QC, RNA-seq reads were mapped to human genome reference GRCh38 using  
874 STAR version 2.7.3a. Ensembl gene annotations release 99 were used. For HERV expression  
875 annotations of 519,060 loci from Human Endogenous Retrovirus Database  
876 (<https://herv.img.cas.cz>) were used (database accessed on 26 June 2020). Read counts per  
877 gene/HERV were found in the same STAR run. Features with at least 10 counts in at least one  
878 sample were selected. Downstream analysis was performed in RStudio using R version 4.0.2.  
879 The code is available at GitHub ([https://github.com/bartongroup/MG\\_SumoDiff2](https://github.com/bartongroup/MG_SumoDiff2)).  
880 Differential expression was performed using edgeR version 3.30.3. Gene/HERV profiles were  
881 calculated as a log<sub>2</sub> ratio of normalised counts at time point 4, 8, 24 and 48 hours versus  
882 DMSO. These profiles were used for clustering. A summary of these data can be found in  
883 Supplementary Data File 6. A notebook with details of RNA-seq, ChIP-seq and ATAC-seq data  
884 analyses is available at  
885 <http://www.compbio.dundee.ac.uk/user/mgierlinski/sumodiff2/doc/analysis.html>.

886

887 **RNA seq protein-coding genes clustering analysis.** Protein-coding gene data were cross-  
888 referenced to protein names using UniProt mapping (uniprot.org) leaving 15333 entries.  
889 These along with fold change values at each of the time-points of ML792 treatment were

890 uploaded to Perseus (v1.6.1.1) and each protein coding gene annotated with GOBP, GOMF,  
891 GOCC, KEGG, Pfam, GSEA, Keywords, Corum, PRINTS, Prosite, SMART and Reactome terms  
892 using the UniProt reference. Based on the entire time-course quantitative data hierarchical  
893 clustering was performed using an unconstrained Euclidian distance method pre-processed  
894 with k-means. 300 clusters were considered with 20 iterations and 5 restarts. Based on the  
895 hierarchical clustering, multiple rounds of cluster definition were made using a range of  
896 fbetween 5 and 10 clusters. 9 clusters gave the highest number of significantly enriched or  
897 depleted categorical terms (ontologies) according to Fisher's exact test employing a 2% FDR  
898 truncation. They were labelled A-I and carried forward for functional enrichment analysis. To  
899 reduce the quantitative data to a single metric (for STRING analysis) a slope value ( $\log_2$  fold  
900 change per hour) for each entry was calculated based on all time-point fold change values in  
901 addition to a zero value at 0h. These RNA-seq data and the ML792 proteomics data  
902 (experiment 1) were combined using gene names as cross-reference. This gave 4526 entries  
903 with data in both the proteomics and RNA-seq experiments. A summary of these data can be  
904 found in Supplementary Data File 4.  
905

906 **References**

- 907 1. Chambers, I. *et al.* Nanog safeguards pluripotency and mediates germline  
908 development. *Nature* **450**, 1230-1234 (2007).
- 909 2. Young, R.A. Control of the embryonic stem cell state. *Cell* **144**, 940-954 (2011).
- 910 3. Takahashi, K. *et al.* Induction of pluripotent stem cells from adult human fibroblasts  
911 by defined factors. *Cell* **131**, 861-872 (2007).
- 912 4. Takahashi, K. & Yamanaka, S. Induction of pluripotent stem cells from mouse  
913 embryonic and adult fibroblast cultures by defined factors. *Cell* **126**, 663-676 (2006).
- 914 5. Yu, J. *et al.* Induced pluripotent stem cell lines derived from human somatic cells.  
915 *Science* **318**, 1917-1920 (2007).
- 916 6. Vierbuchen, T. *et al.* Direct conversion of fibroblasts to functional neurons by defined  
917 factors. *Nature* **463**, 1035-1041 (2010).
- 918 7. Apostolou, E. & Hochedlinger, K. Chromatin dynamics during cellular reprogramming.  
919 *Nature* **502**, 462-471 (2013).
- 920 8. Cheloufi, S. *et al.* The histone chaperone CAF-1 safeguards somatic cell identity.  
921 *Nature* **528**, 218-224 (2015).
- 922 9. Borkent, M. *et al.* A Serial shRNA Screen for Roadblocks to Reprogramming Identifies  
923 the Protein Modifier SUMO2. *Stem Cell Reports* **6**, 704-716 (2016).
- 924 10. Cossec, J.C. *et al.* SUMO Safeguards Somatic and Pluripotent Cell Identities by  
925 Enforcing Distinct Chromatin States. *Cell Stem Cell* **23**, 742-757 e748 (2018).
- 926 11. Flotho, A. & Melchior, F. Sumoylation: a regulatory protein modification in health and  
927 disease. *Annu Rev Biochem* **82**, 357-385 (2013).
- 928 12. Hay, R.T. SUMO: a history of modification. *Mol Cell* **18**, 1-12 (2005).
- 929 13. Rodriguez, M.S., Dargemont, C. & Hay, R.T. SUMO-1 conjugation in vivo requires both  
930 a consensus modification motif and nuclear targeting. *J Biol Chem* **276**, 12654-12659  
931 (2001).
- 932 14. Sampson, D.A., Wang, M. & Matunis, M.J. The small ubiquitin-like modifier-1 (SUMO-  
933 1) consensus sequence mediates Ubc9 binding and is essential for SUMO-1  
934 modification. *J Biol Chem* **276**, 21664-21669 (2001).
- 935 15. Tatham, M.H. *et al.* Polymeric chains of SUMO-2 and SUMO-3 are conjugated to  
936 protein substrates by SAE1/SAE2 and Ubc9. *J Biol Chem* **276**, 35368-35374 (2001).
- 937 16. Song, J., Durrin, L.K., Wilkinson, T.A., Krontiris, T.G. & Chen, Y. Identification of a  
938 SUMO-binding motif that recognizes SUMO-modified proteins. *Proc Natl Acad Sci U S*  
939 *A* **101**, 14373-14378 (2004).
- 940 17. Williams, R.L. *et al.* Myeloid leukaemia inhibitory factor maintains the developmental  
941 potential of embryonic stem cells. *Nature* **336**, 684-687 (1988).
- 942 18. Ying, Q.L., Nichols, J., Chambers, I. & Smith, A. BMP induction of Id proteins suppresses  
943 differentiation and sustains embryonic stem cell self-renewal in collaboration with  
944 STAT3. *Cell* **115**, 281-292 (2003).
- 945 19. Xu, R.H. *et al.* BMP4 initiates human embryonic stem cell differentiation to  
946 trophoblast. *Nat Biotechnol* **20**, 1261-1264 (2002).
- 947 20. Xu, R.H. *et al.* Basic FGF and suppression of BMP signaling sustain undifferentiated  
948 proliferation of human ES cells. *Nat Methods* **2**, 185-190 (2005).
- 949 21. Humphrey, R.K. *et al.* Maintenance of pluripotency in human embryonic stem cells is  
950 STAT3 independent. *Stem Cells* **22**, 522-530 (2004).
- 951 22. Xu, C. *et al.* Basic fibroblast growth factor supports undifferentiated human embryonic  
952 stem cell growth without conditioned medium. *Stem Cells* **23**, 315-323 (2005).

- 953 23. James, D., Levine, A.J., Besser, D. & Hemmati-Brivanlou, A. TGFbeta/activin/nodal  
954 signaling is necessary for the maintenance of pluripotency in human embryonic stem  
955 cells. *Development* **132**, 1273-1282 (2005).
- 956 24. Rossant, J. Mouse and human blastocyst-derived stem cells: vive les differences.  
957 *Development* **142**, 9-12 (2015).
- 958 25. He, X. *et al.* Probing the roles of SUMOylation in cancer cell biology by using a selective  
959 SAE inhibitor. *Nature Chemical Biology* **13**, 1164-1171 (2017).
- 960 26. Bray, M.A. *et al.* Cell Painting, a high-content image-based assay for morphological  
961 profiling using multiplexed fluorescent dyes. *Nat Protoc* **11**, 1757-1774 (2016).
- 962 27. Tatham, M.H. *et al.* RNF4 is a poly-SUMO-specific E3 ubiquitin ligase required for  
963 arsenic-induced PML degradation. *Nat Cell Biol* **10**, 538-546 (2008).
- 964 28. Seifert, A., Schofield, P., Barton, G.J. & Hay, R.T. Proteotoxic stress reprograms the  
965 chromatin landscape of SUMO modification. *Sci Signal* **8**, rs7 (2015).
- 966 29. Yang, B.X. *et al.* Systematic identification of factors for provirus silencing in embryonic  
967 stem cells. *Cell* **163**, 230-245 (2015).
- 968 30. Heinz, S. *et al.* Simple combinations of lineage-determining transcription factors prime  
969 cis-regulatory elements required for macrophage and B cell identities. *Mol Cell* **38**,  
970 576-589 (2010).
- 971 31. Gerstein, M.B. *et al.* Architecture of the human regulatory network derived from  
972 ENCODE data. *Nature* **489**, 91-100 (2012).
- 973 32. Wang, Y. *et al.* Endogenous miRNA sponge lincRNA-RoR regulates Oct4, Nanog, and  
974 Sox2 in human embryonic stem cell self-renewal. *Dev Cell* **25**, 69-80 (2013).
- 975 33. Girdwood, D. *et al.* P300 transcriptional repression is mediated by SUMO modification.  
976 *Mol Cell* **11**, 1043-1054 (2003).
- 977 34. Paces, J., Pavlicek, A. & Paces, V. HERVd: database of human endogenous retroviruses.  
978 *Nucleic Acids Res* **30**, 205-206 (2002).
- 979 35. Paces, J. *et al.* HERVd: the Human Endogenous RetroViruses Database: update. *Nucleic*  
980 *Acids Res* **32**, D50 (2004).
- 981 36. Tammsalu, T. *et al.* Proteome-wide identification of SUMO2 modification sites. *Sci*  
982 *Signal* **7**, rs2 (2014).
- 983 37. Wolf, D. & Goff, S.P. TRIM28 Mediates Primer Binding Site-Targeted Silencing of  
984 Murine Leukemia Virus in Embryonic Cells. *Cell* **131**, 46-57 (2007).
- 985 38. Tie, C.H. *et al.* KAP1 regulates endogenous retroviruses in adult human cells and  
986 contributes to innate immune control. *EMBO Rep* **19** (2018).
- 987 39. Schmidt, N. *et al.* An influenza virus-triggered SUMO switch orchestrates co-opted  
988 endogenous retroviruses to stimulate host antiviral immunity. *Proceedings of the*  
989 *National Academy of Sciences* **116**, 17399-17408 (2019).
- 990 40. Oleksiewicz, U. *et al.* TRIM28 and Interacting KRAB-ZNFs Control Self-Renewal of  
991 Human Pluripotent Stem Cells through Epigenetic Repression of Pro-differentiation  
992 Genes. *Stem Cell Reports* **9**, 2065-2080 (2017).
- 993 41. Goke, J. *et al.* Dynamic transcription of distinct classes of endogenous retroviral  
994 elements marks specific populations of early human embryonic cells. *Cell Stem Cell* **16**,  
995 135-141 (2015).
- 996 42. Glinsky, G.V. Transposable Elements and DNA Methylation Create in Embryonic Stem  
997 Cells Human-Specific Regulatory Sequences Associated with Distal Enhancers and  
998 Noncoding RNAs. *Genome Biol Evol* **7**, 1432-1454 (2015).

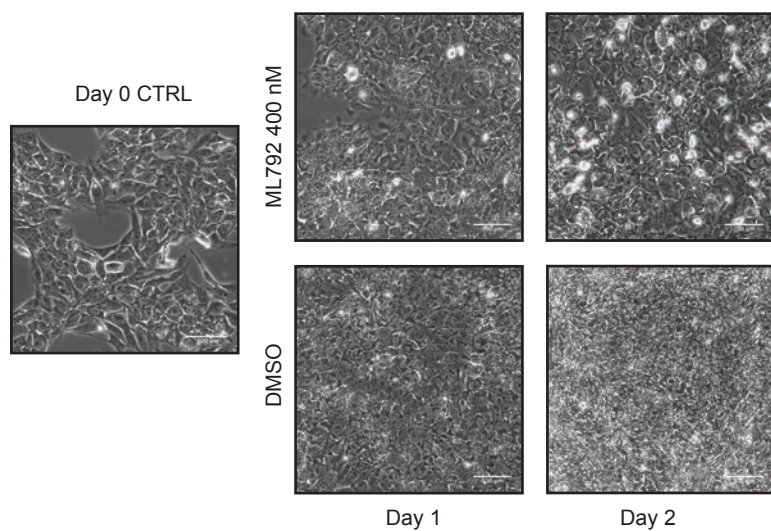


- 999 43. Santoni, F.A., Guerra, J. & Luban, J. HERV-H RNA is abundant in human embryonic stem  
1000 cells and a precise marker for pluripotency. *Retrovirology* **9**, 111 (2012).
- 1001 44. Lu, X. *et al.* The retrovirus HERVH is a long noncoding RNA required for human  
1002 embryonic stem cell identity. *Nat Struct Mol Biol* **21**, 423-425 (2014).
- 1003 45. Ohnuki, M. *et al.* Dynamic regulation of human endogenous retroviruses mediates  
1004 factor-induced reprogramming and differentiation potential. *Proc Natl Acad Sci U S A*  
1005 **111**, 12426-12431 (2014).
- 1006 46. Zhang, Y. *et al.* Transcriptionally active HERV-H retrotransposons demarcate  
1007 topologically associating domains in human pluripotent stem cells. *Nat Genet* **51**,  
1008 1380-1388 (2019).
- 1009 47. Merckenschlager, M. & Nora, E.P. CTCF and Cohesin in Genome Folding and  
1010 Transcriptional Gene Regulation. *Annu Rev Genomics Hum Genet* **17**, 17-43 (2016).
- 1011 48. Rowley, M.J. & Corces, V.G. Organizational principles of 3D genome architecture. *Nat*  
1012 *Rev Genet* **19**, 789-800 (2018).
- 1013 49. Bailey, S.D. *et al.* ZNF143 provides sequence specificity to secure chromatin  
1014 interactions at gene promoters. *Nat Commun* **2**, 6186 (2015).
- 1015 50. Anno, Y.N. *et al.* Genome-wide evidence for an essential role of the human  
1016 Staf/ZNF143 transcription factor in bidirectional transcription. *Nucleic Acids Res* **39**,  
1017 3116-3127 (2011).
- 1018 51. Sathyan, K.M. *et al.* An improved auxin-inducible degron system preserves native  
1019 protein levels and enables rapid and specific protein depletion. *Genes Dev* **33**, 1441-  
1020 1455 (2019).
- 1021 52. Theurillat, I. *et al.* Extensive SUMO Modification of Repressive Chromatin Factors  
1022 Distinguishes Pluripotent from Somatic Cells. *Cell Rep* **32**, 108146 (2020).
- 1023 53. Ikeda, H. *et al.* Characterization of an antigen that is recognized on a melanoma  
1024 showing partial HLA loss by CTL expressing an NK inhibitory receptor. *Immunity* **6**, 199-  
1025 208 (1997).
- 1026 54. Kilpinen, S. *et al.* Systematic bioinformatic analysis of expression levels of 17,330  
1027 human genes across 9,783 samples from 175 types of healthy and pathological tissues.  
1028 *Genome Biol* **9**, R139 (2008).
- 1029 55. Zhang, L. *et al.* MSX2 Initiates and Accelerates Mesenchymal Stem/Stromal Cell  
1030 Specification of hPSCs by Regulating TWIST1 and PRAME. *Stem Cell Reports* **11**, 497-  
1031 513 (2018).
- 1032 56. Costessi, A. *et al.* The tumour antigen PRAME is a subunit of a Cul2 ubiquitin ligase and  
1033 associates with active NFY promoters. *EMBO J* **30**, 3786-3798 (2011).
- 1034 57. Costessi, A. *et al.* The human EKC/KEOPS complex is recruited to Cullin2 ubiquitin  
1035 ligases by the human tumour antigen PRAME. *PLoS One* **7**, e42822 (2012).
- 1036 58. Perez-Riverol, Y. *et al.* The PRIDE database and related tools and resources in 2019:  
1037 improving support for quantification data. *Nucleic Acids Res* **47**, D442-D450 (2019).
- 1038 59. Hands, K.J., Cuchet-Lourenco, D., Everett, R.D. & Hay, R.T. PML isoforms in response  
1039 to arsenic: high-resolution analysis of PML body structure and degradation. *J Cell Sci*  
1040 **127**, 365-375 (2014).
- 1041 60. Shen, L. *et al.* SUMO protease SENP1 induces isomerization of the scissile peptide  
1042 bond. *Nature Structural & Molecular Biology* **13**, 1069-1077 (2006).
- 1043 61. Ibrahim, A.F.M. *et al.* Antibody RING-Mediated Destruction of Endogenous Proteins.  
1044 *Mol Cell* **79**, 155-166 e159 (2020).

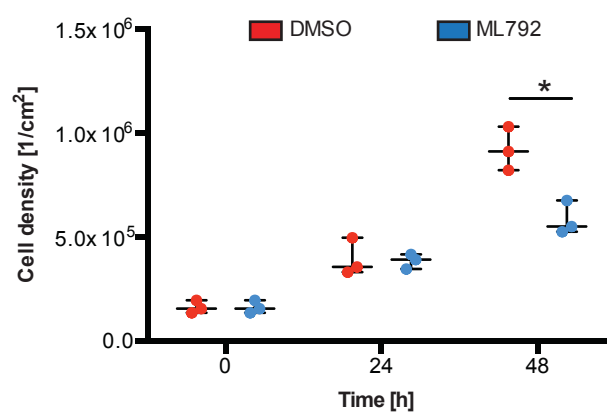
- 1045 62. Ludwig, T.E. *et al.* Derivation of human embryonic stem cells in defined conditions.  
1046 *Nat Biotechnol* **24**, 185-187 (2006).
- 1047 63. Shevchenko, A., Tomas, H., Havlis, J., Olsen, J.V. & Mann, M. In-gel digestion for mass  
1048 spectrometric characterization of proteins and proteomes. *Nat Protoc* **1**, 2856-2860  
1049 (2006).
- 1050 64. Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized  
1051 p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nat*  
1052 *Biotechnol* **26**, 1367-1372 (2008).
- 1053 65. Tyanova, S. *et al.* The Perseus computational platform for comprehensive analysis of  
1054 (prote)omics data. *Nat Methods* **13**, 731-740 (2016).
- 1055 66. Tammsalu, T. *et al.* Proteome-wide identification of SUMO modification sites by mass  
1056 spectrometry. *Nat Protoc* **10**, 1374-1388 (2015).
- 1057 67. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with increased  
1058 coverage, supporting functional discovery in genome-wide experimental datasets.  
1059 *Nucleic Acids Res* **47**, D607-D613 (2019).
- 1060 68. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of  
1061 biomolecular interaction networks. *Genome Res* **13**, 2498-2504 (2003).
- 1062 69. Arrigoni, L. *et al.* Standardizing chromatin research: a simple and universal method for  
1063 ChIP-seq. *Nucleic Acids Res* **44**, e67 (2016).
- 1064 70. Corces, M.R. *et al.* Omni-ATAC-seq: Improved ATAC-seq protocol. *Protocol Exchange*  
1065 (2017).  
1066

# Figure 1.

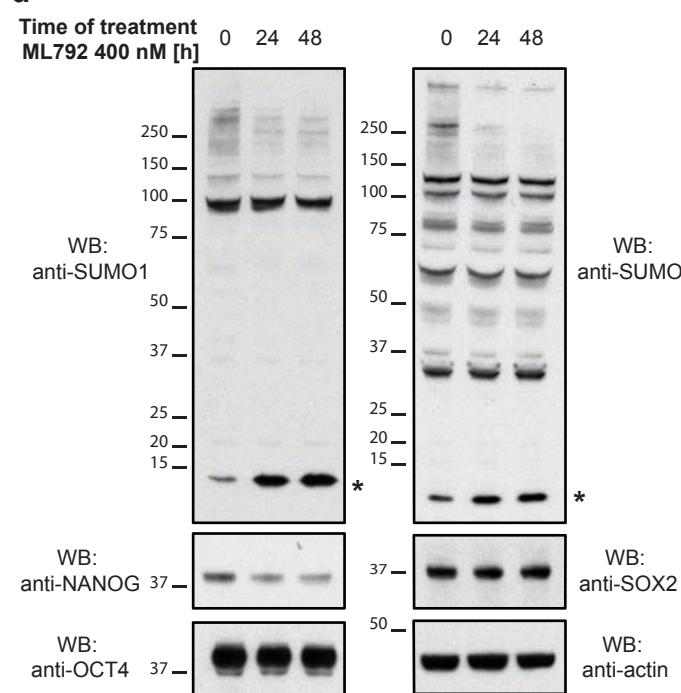
**a**



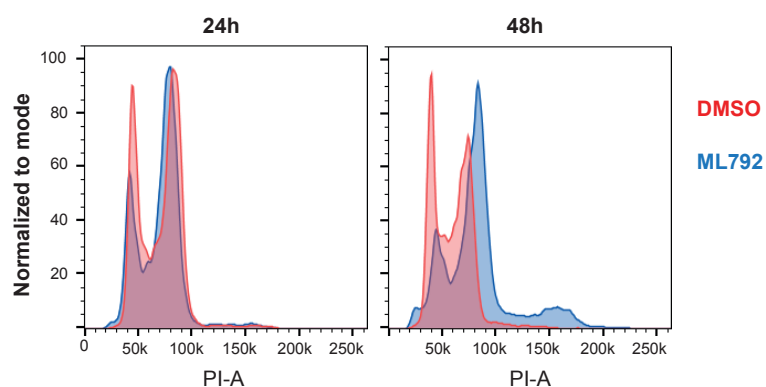
**b**



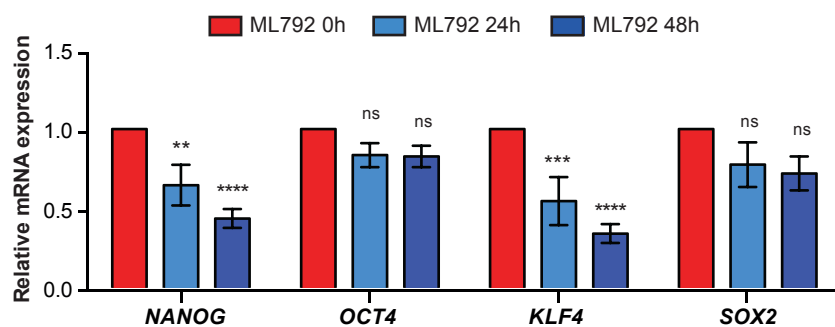
**d**



**c**



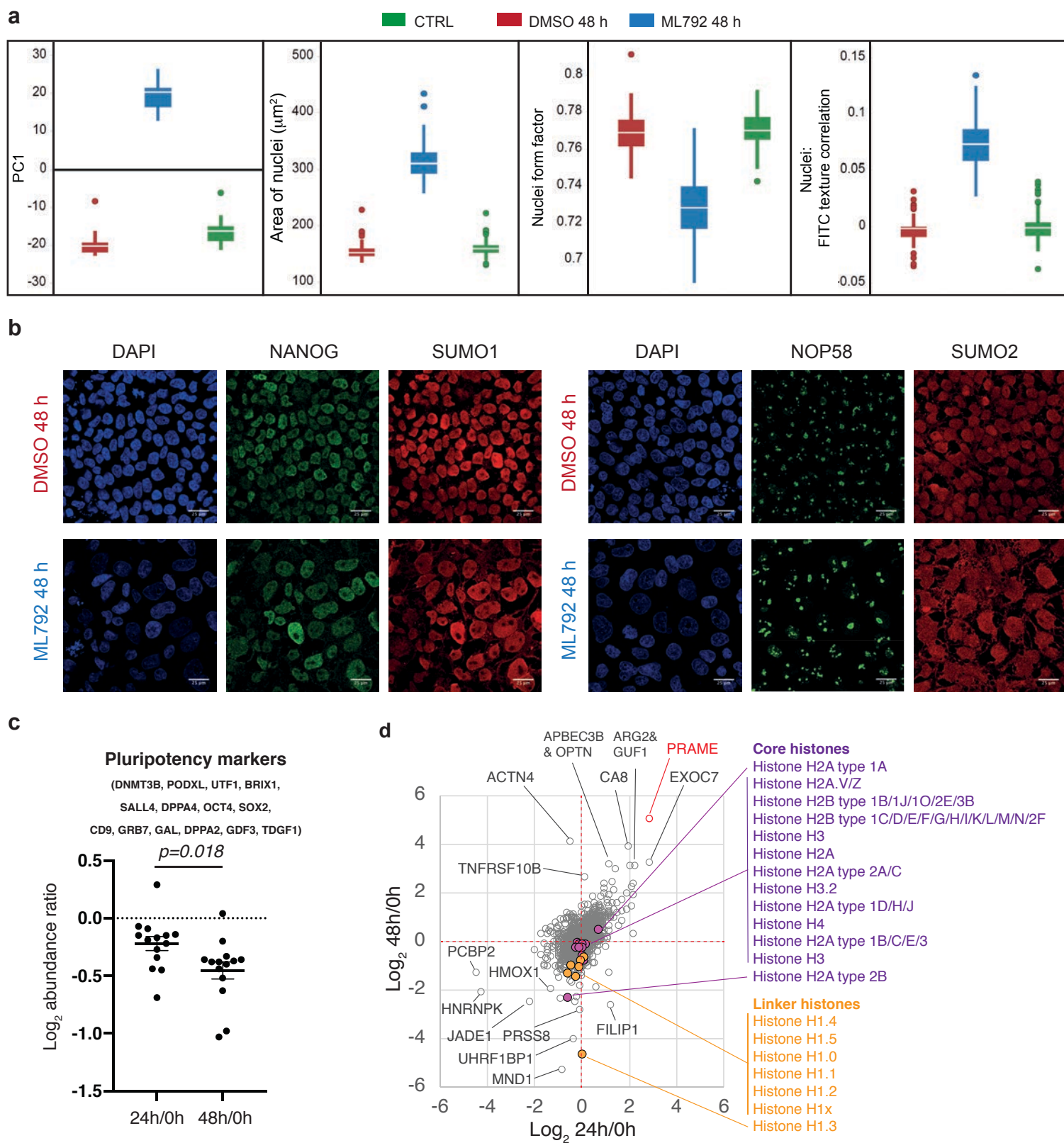
**e**



**Figure 1. Inhibition of SUMO modification leads to loss of select pluripotency markers.**

ChiPS4 cells were treated with ML792 (400 nM) or DMSO vehicle for the indicated time and analysed by various approaches. **a.** Cell morphology was assessed using phase contrast microscopy (all images contain a 100  $\mu\text{m}$  scale bar). **b.** For proliferation assessment, ChiPS4 cells were seeded at a standard density of  $3 \times 10^5$  cells/cm<sup>2</sup> in triplicate for each time point. The following day cells were treated with DMSO vehicle or ML792 and every 24h they were harvested using TrypLE select and counted. Data are plotted as mean cell density (line) with individual replicates (dots) shown N=3. Statistical significance was calculated with t-tests corrected for multiple comparisons using Holm-Sidak's method (\* P<0.05 significantly different from the corresponding DMSO control) **c.** To analyse cell cycle distribution, cells were collected as in **b**, fixed, stained with propidium iodide (PI) and analysed by flow cytometry. Plots are a representative of three independent experiments **d.** Protein samples were analysed by Western blotting to determine conjugation levels of SUMO1, SUMO2/3 and abundance of key pluripotency markers NANOG, SOX2 and OCT4 using appropriate antibodies. Anti-Actin western blot was used as a loading control. \* represents a band corresponding to free SUMO1 or SUMO2/3. **e.** Cells were treated with ML792 and after the indicated time they were lysed and total RNA was extracted. mRNA levels of *NANOG*, *OCT4*, *SOX2* and *KLF4* were determined by qPCR. Relative mRNA expression levels normalized to *TBP* were plotted as means  $\pm$  SEM of four independent experiments. \*\*P <0.01; \*\*\*P<0.001; \*\*\*\*P<0.0001 significantly different from the corresponding value for untreated control (two-way ANOVA followed by Sidak's multiple comparison test).

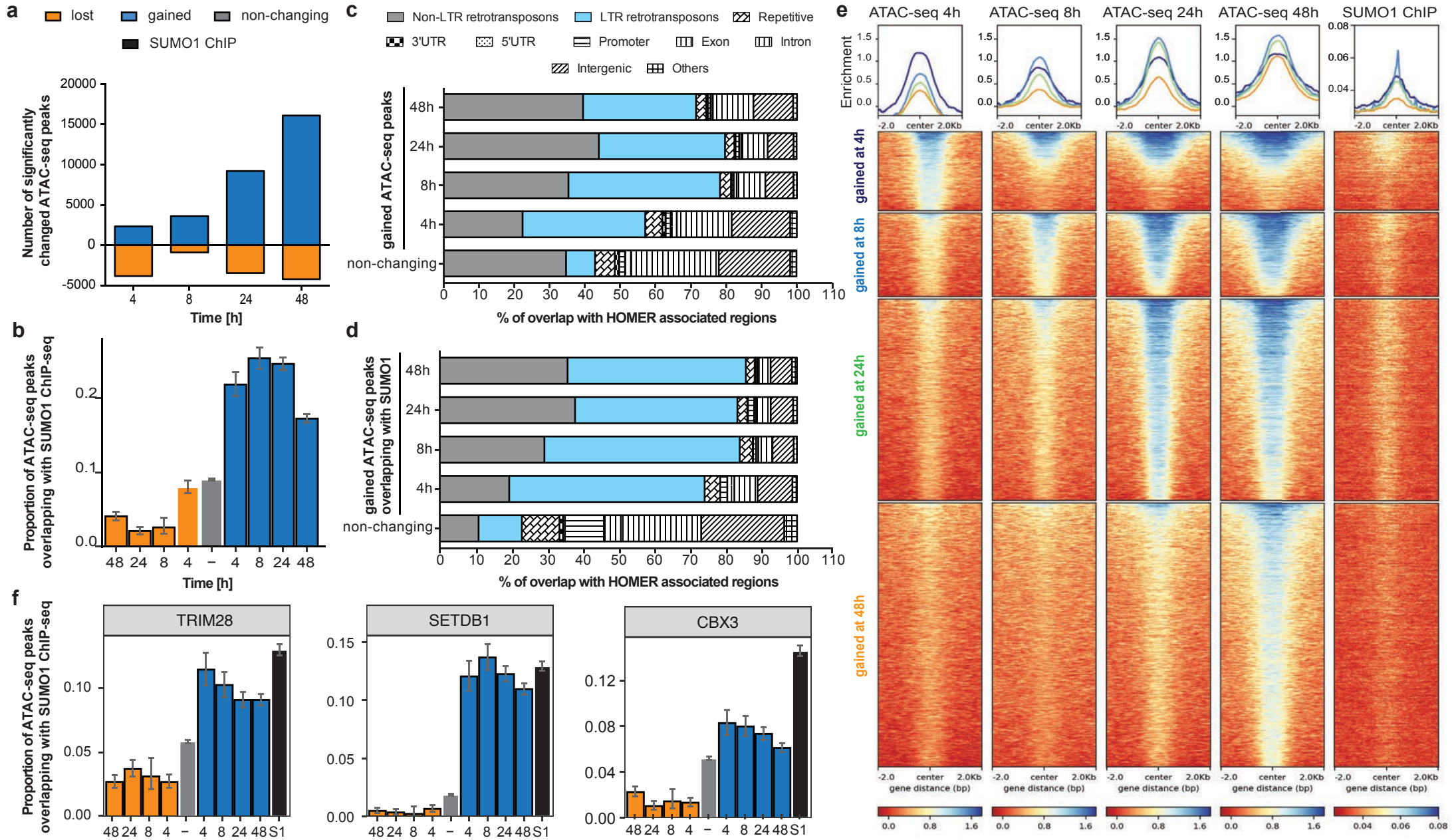
**Figure 2.**



**Figure 2. Change in morphology, but unchanged proteome in hiPSCs in presence of ML792.**

**a.** Cell painting analysis. ChiPS4 cells were treated with PBS, DMSO vehicle or 400 nM ML792 for 48 h. Cells were then stained, fixed and analysed using high content microscopy. The experiment was performed three times with 8 replicates per condition. Information extracted from cell painting analysis was focused on subcellular compartments most affected by ML792 treatment. Most variation between treatments and controls in PCA was captured by PC1. Selected graphs represent quantitation of individual measures contributing to the difference observed in PCA: area of nuclei; nuclei form factor (size and shape of nucleus); nuclei FITC texture correlation (size of nucleolar structures). **b.** ChiPS4 cells were treated for 48 h with DMSO vehicle or 400 nM ML792, fixed and stained with DAPI (blue), anti-SUMO1 or anti-SUMO2 (red) and anti-NANOG or anti-NOP58 (green) antibodies. IF images were obtained using a Leica SP8 confocal microscope and a 60x water immersion lens. All images contain 25  $\mu\text{m}$  scale bar. **c.**  $\text{Log}_2$  abundance ratio data extracted from whole cell proteomic analysis for the 14 indicated markers of pluripotency comparing 24 h and 48 h ML792 exposure to untreated cells. The plot shows individual data points and mean with standard error of the mean for the entire set. The result of a paired two-tailed student's t-test is shown. **d.** Scatter plot of  $\text{Log}_2$  24h/0h and  $\text{log}_2$  48h/0h abundance change for the entire 4741 protein whole cell proteomic dataset. Extreme outliers are indicated. All identified core and linker histones are represented by coloured markers, others are in grey. Linker histones were identified by STRING analysis as a functionally related group of proteins that are significantly reduced in abundance at 48h compared with 0h. Core histone proteins are indicated for reference.

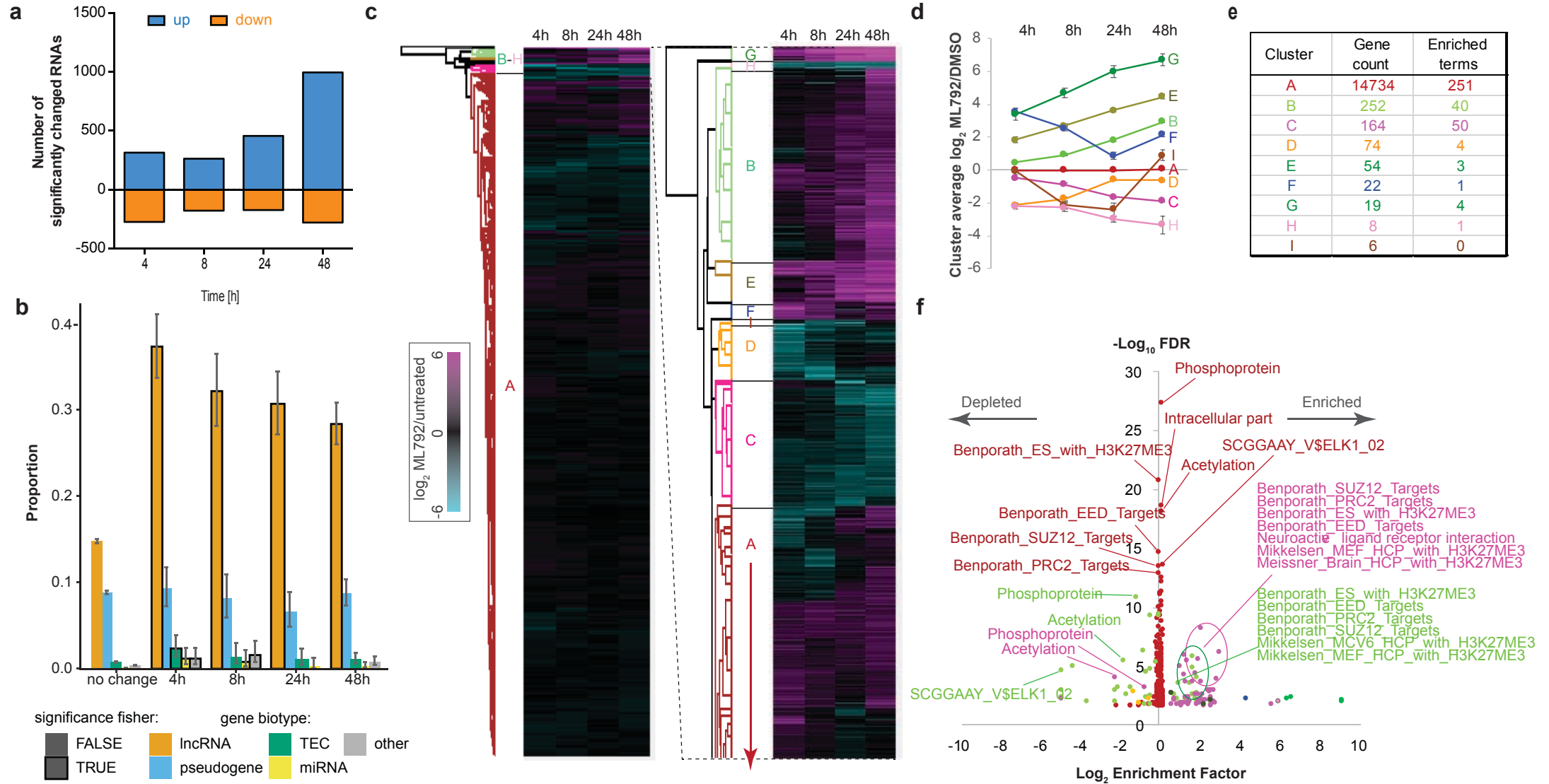
**Figure 3.**



**Figure 3. Removal of SUMO in hiPSCs increases chromatin accessibility.** **a.** Numbers of significantly changed ATAC-seq peaks (applied criteria:  $|\log_2FC| > 1.5$ , minimal overlap of 3) at each time point in ChIP4 cells treated with ML792 when compared to DMSO vehicle. Gained peaks are shown in blue, lost in orange. **b.** Proportion of ATAC-seq peaks gained or lost at different times of ML792 treatment overlapping with SUMO-1 ChIP-seq peaks (found in untreated cells). Overlap between non-changing ATAC-seq peaks and SUMO1 ChIP-seq is shown as a reference (grey). Thick black borders mark statistically significant changes (Fisher's exact test, corrected for multiple tests using Benjamini-Hochberg method), error bars are 95% confidence intervals (CI) of the proportion. **c.** Percentage of overlap between ATAC-seq peaks gained at each time point of ML792 treatment or non-changing ATAC-seq peaks found in DMSO vehicle control with HOMER-based annotations of chromatin regions (indicated). **d.** Percentage of overlap between peaks common for SUMO1 ChIP-seq/gained ATAC-seq peaks at each time point of ML792 treatment and HOMER-based annotations of chromatin regions. Different chromatin regions are represented as in the legend in **c.** **e.** Density plots for gained ATAC-seq changes at each time point following ML792 treatment. The sites are ordered by the time at which a change of  $>1.5$  fold is first detected. The same order of genomic locations has been plotted for SUMO ChIP-seq peaks. Scale used for each density plot is based on the  $\log_2$  ratio of ATAC-seq signal at each time point to the signal detected in DMSO control. Graphs at the top represent summary plots for the ATAC-seq signal at changing sites for the indicated time points. **f.** Proportion of ATAC-seq peaks gained or lost at different time points of ML792 treatment overlapping with various transcription factor ChIP-seq peaks (TRIM28, SETDB1, CBX3; data obtained from ENCODE database). Overlap between non-changing ATAC-seq peaks or SUMO1 ChIP-seq peaks are shown as references (grey and black bars respectively). Thick black borders mark statistically significant changes (Fisher's exact test, corrected for multiple tests using Benjamini-Hochberg method), error bars are 95% CI of the proportion.



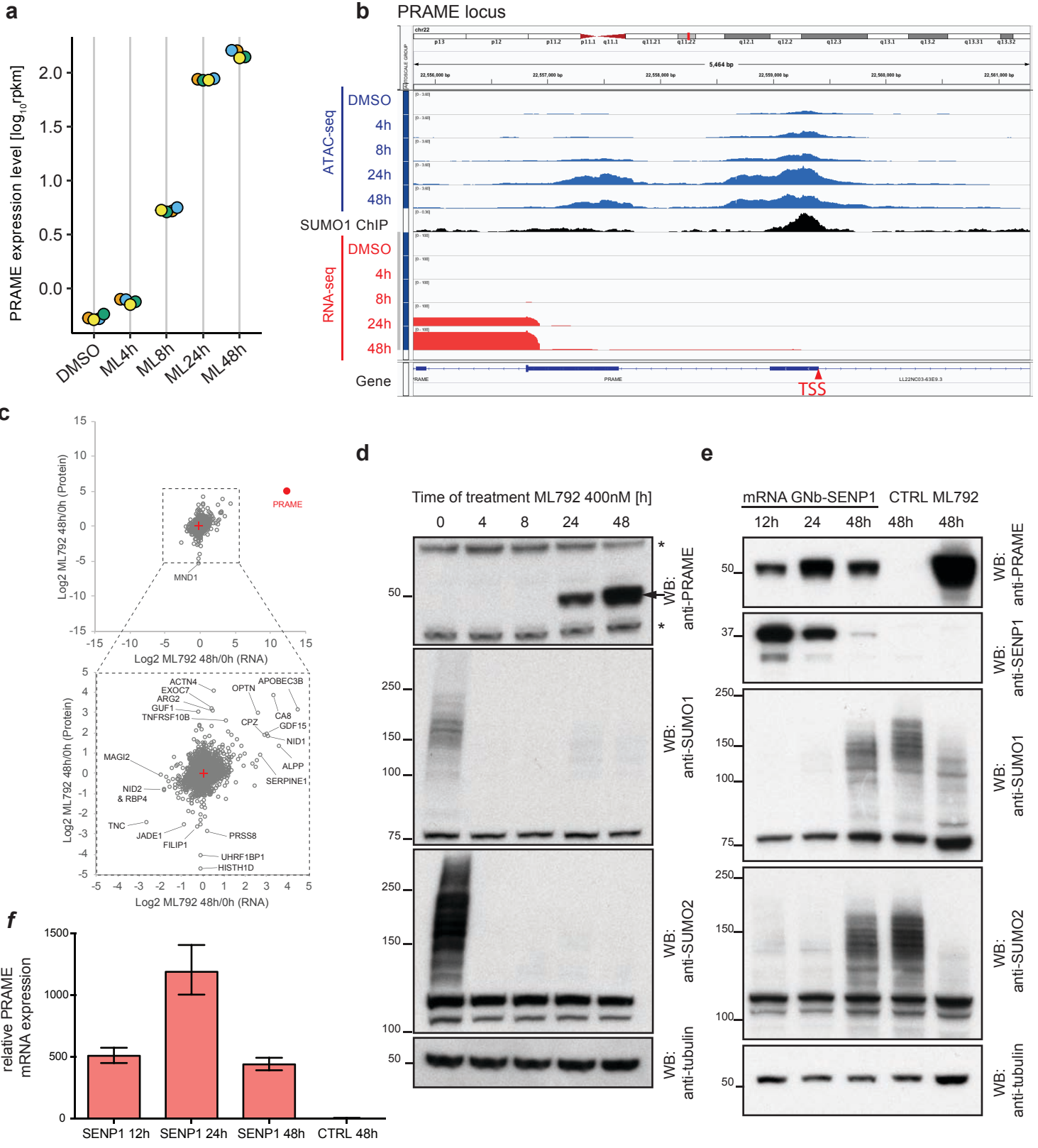
**Figure 4.**



**Figure 4. Inhibition of SUMO modification in hiPSCs selectively alters transcription. a.**

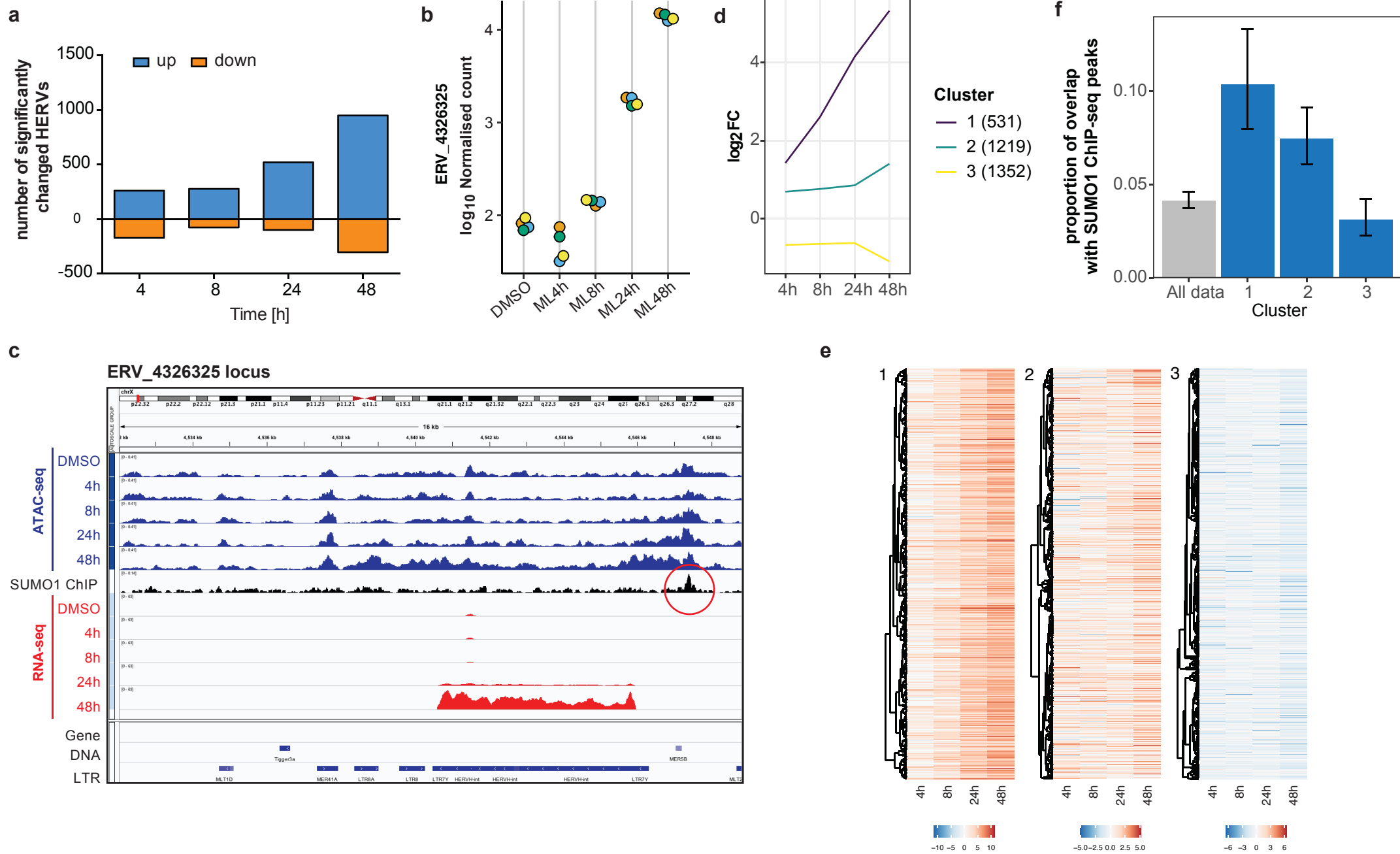
Numbers of significantly changed RNAs in ChiPS4 cells (applied criteria:  $|\log_2FC| > 1.5$ ,  $FDR < 0.05$ ,  $P < 0.01$ ) at each time point of ML792 treatment (up regulated – blue, down regulated – orange) when compared to DMSO vehicle treated cells. **b.** Distribution of biotypes among non-protein coding genes that are significantly changed between DMSO and ML792 at a given time point and genes that do not change at any time point, with respect to DMSO. Error bars are 95% CI of a proportion. Black outlines indicate proportions significantly different (Fisher's exact test,  $p < 0.05$ ) from the "no change" group. **c.** Hierarchical clustering analysis of protein-coding gene mRNAs during ML792 treatment. For each time-point data were represented as  $\log_2$  fold change compared to DMSO treatment and clustered using a Euclidean distance function with linkage based on averages and k-means pre-processing. Data were binned into 9 row clusters (labelled A-H). The entire data set of 15333 entries is shown (left) and a zoom view of clusters B-H (right). **d.** Cluster-specific data shown as average and SEM at each time-point for all members of each group. **e.** Overview of functional group enrichment for each cluster relative to the entire dataset. Gene IDs were converted to protein IDs and the enrichment of different functional annotations was calculated by Fisher Exact Test with Benjamini-Hochberg truncation at 2% FDR. **f.** Scatter plot representation of all enriched categories for all clusters using  $\log_2$  enrichment factor and  $-\log_{10}$  FDR as co-ordinates.

**Figure 5.**



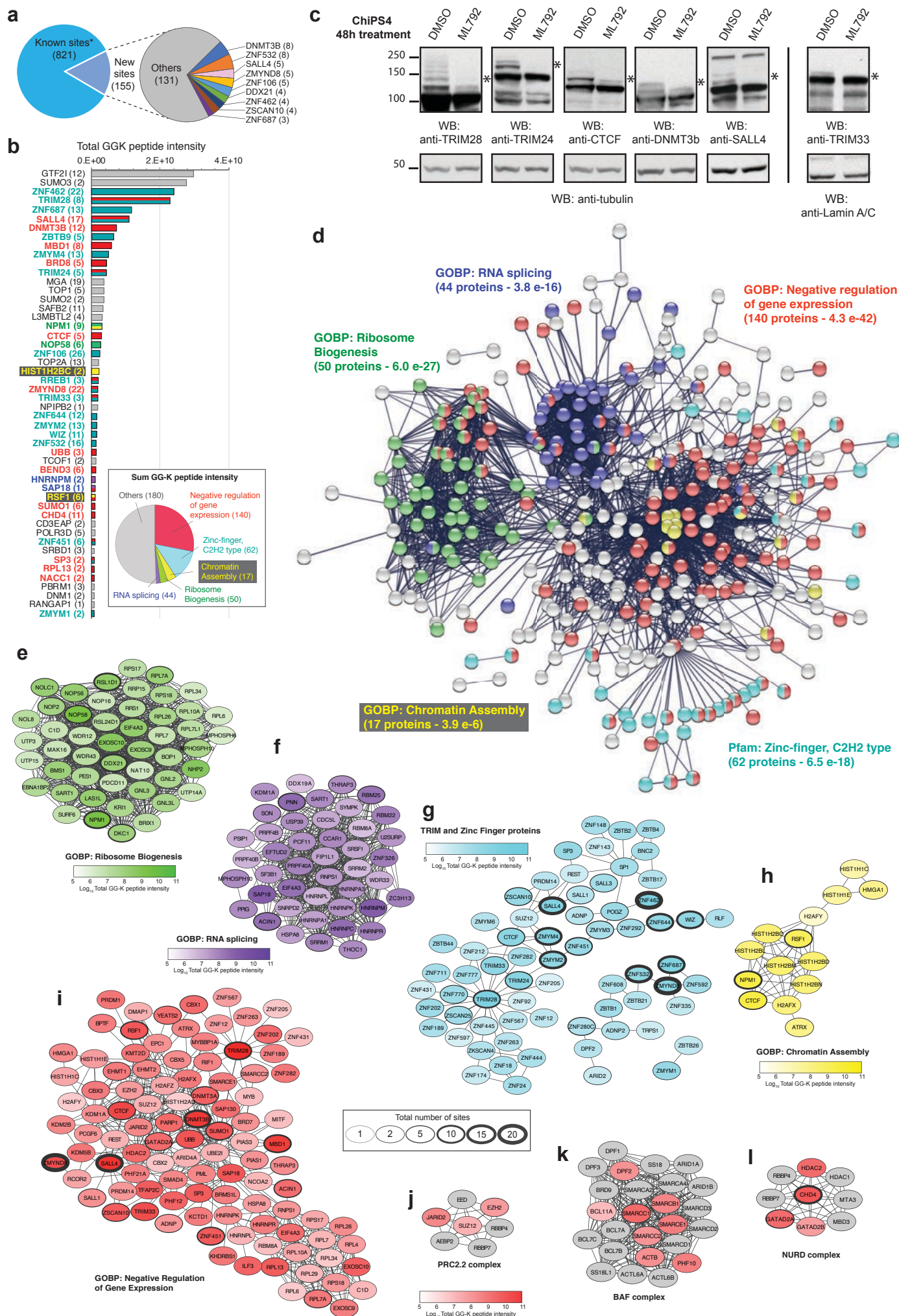
**Figure 5. Inhibition of SUMOylation leads to expression of PRAME in hiPSCs.** **a.** PRAME expression levels in RNA-seq samples at different time points were plotted as  $\log_{10}$ rpkm (reads per kilobase of transcript, per million mapped reads). Individual replicates are represented in different colours. **b.** Integrative Genomic Viewer was used to visualize changes in ATAC-seq (blue) and RNA-seq (red) occurring at the PRAME locus in response to ML792 treatment. SUMO1 ChIP-seq signal was also aligned and represented in black. All traces of the same type (ATAC-seq or RNA-seq) were normalized and scaled in the same way. **c.** Scatter-plot of RNA change versus protein-level change during 48 hours ML792 treatment (4498 common entries). Full scale is shown in the upper panel, and a smaller scale is shown below. Selected outliers are indicated. **d.** Western blot analysis of PRAME protein expression and SUMO1, SUMO2/3 conjugation levels after treatment of ChiPS4 cells with ML792 for the indicated times. Anti-tubulin was used as a loading control. \* represents nonspecific bands, while the arrow indicates PRAME. **e.** ChiPS4 cells were transfected with capped and polyadenylated RNA encoding the catalytic domain of SUMO specific protease SENP1 (GNb-SENP1). Samples were collected at the times indicated after transfection and analysed by Western blotting for PRAME expression, SENP1 expression and conjugation levels of SUMO1 and SUMO2/3. Tubulin was used as loading control. **f.** ChiPS4 cells transfected with GNb-SENP1 were analysed by RT-qPCR to assess the relative expression of *PRAME* mRNA using *TBP* as a normalizing gene.

Figure 6.



**Figure 6. SUMO modification regulates HERV expression in hiPSCs.** **a.** Numbers of significantly changed HERVs in ChiPS4 cells (applied criteria:  $|\log_2FC| > 1.5$ ,  $FDR < 0.05$ ,  $P < 0.01$ ) at each time point of ML792 treatment (up regulated – blue, down regulated – orange) when compared to DMSO vehicle treated cells. **b.** Expression levels of ERV\_4326325 in RNA-seq samples at different time points following ML792 treatment of ChiPS4 cells plotted as  $\log_{10}$  normalized counts. Individual replicates are represented in different colours. **c.** Integrative Genomic Viewer display of changes in ATAC-seq (blue) and RNA-seq (red) occurring at the ERV\_4326325 locus in response to ML792 treatment. SUMO1 ChIP-seq signal was also aligned and represented in black. All traces of the same type (ATAC-seq or RNA-seq) were normalized and scaled in the same way. **d.** Cluster centroids from clustering HERV profiles into 3 clusters, using k-means clustering. Profiles were created as  $\log_2$  ratio between a given time point and DMSO normalised counts. Only HERVs with at least one statistically significant change,  $FDR < 0.05$ , between any time point and DMSO were selected for clustering. Numbers in brackets show the number of HERVs in each cluster. **e.** Content of each of the HERV clusters indicated in **d.** represented as a heatmap of  $\log_2$  fold change between a given time point and DMSO. **f.** Overlap between HERVs with at least 10 counts detected in at least one sample and SUMO1 ChIP-seq peaks, for all HERV data and for each of the HERV clusters.

**Figure 7.**



**Figure 7. Identification of SUMO1 and SUMO2 targets in hiPSCs.** **a.** 976 SUMO sites identified from 6His-SUMO1-KGG and 6His-SUMO2-KGG ChIPs4 cells, of which 155 were novel compared with previous high-throughput SUMO site proteomics studies (Supplementary Data File 2). Proteins with three or more novel sites are highlighted. **b.** Summary of the top 50 SUMO substrates by total GGK-peptide intensity for all identified sites. Gene names are shown with numbers of sites in brackets. Bars are colour coded by category shown in panel **d**. **d.** The insert shows contribution to total GGK peptide intensity of proteins from the categories shown in **d** (note categories are not mutually exclusive). **c.** Western blot analysis of ChIPs4 cells treated with ML792 or DMSO for 48h. Total protein extracts were probed with anti-TRIM28, anti-TRIM24, anti-CTCF, anti-DNMT3b, anti-SALL4, anti-TRIM33 and anti-tubulin or anti-Lamin A/C antibodies (loading controls). SUMO-modified proteins present above the band for unmodified proteins and disappearing in samples treated with ML792 are labelled with \*. **d.** STRING interaction network of the 427 hiPSCs SUMO substrates. Only high confidence interactions were considered from 'Text mining', 'Experiments' and 'Databases' sources. Network PPI enrichment p-value  $<1.0 \times 10^{-16}$ . Nodes are coloured by functional or structural group as indicated. **e.-l.** Protein interaction networks derived from **d**. for the indicated functional groups. Node shade is proportional to  $\log_{10}$  total GGK peptide intensity and border thickness indicates numbers of sites found. **j.-l.** shows individual network clusters for selected chromatin remodelling complexes. Grey nodes were not identified in the present study. TRIM proteins were included in **g**. to allow more complete network interactions.



1 **Supplementary information**

2

3

4

5

6

7 **SUMO maintains the chromatin environment of human induced pluripotent stem cells.**

8

9

10 Barbara Mojsa<sup>1</sup>, Michael H. Tatham<sup>1</sup>, Lindsay Davidson<sup>2</sup>, Magda Liczmanska<sup>1</sup>, Jane E. Wright<sup>1</sup>,

11 Nicola Wiechens<sup>1</sup>, Marek Gierlinski<sup>3</sup>, Tom Owen-Hughes<sup>1</sup> and Ronald T. Hay<sup>1</sup>

12

13

14

15

16

17 <sup>1</sup>Division of Gene Regulation and Expression, <sup>2</sup>Division of Cell and Developmental Biology,

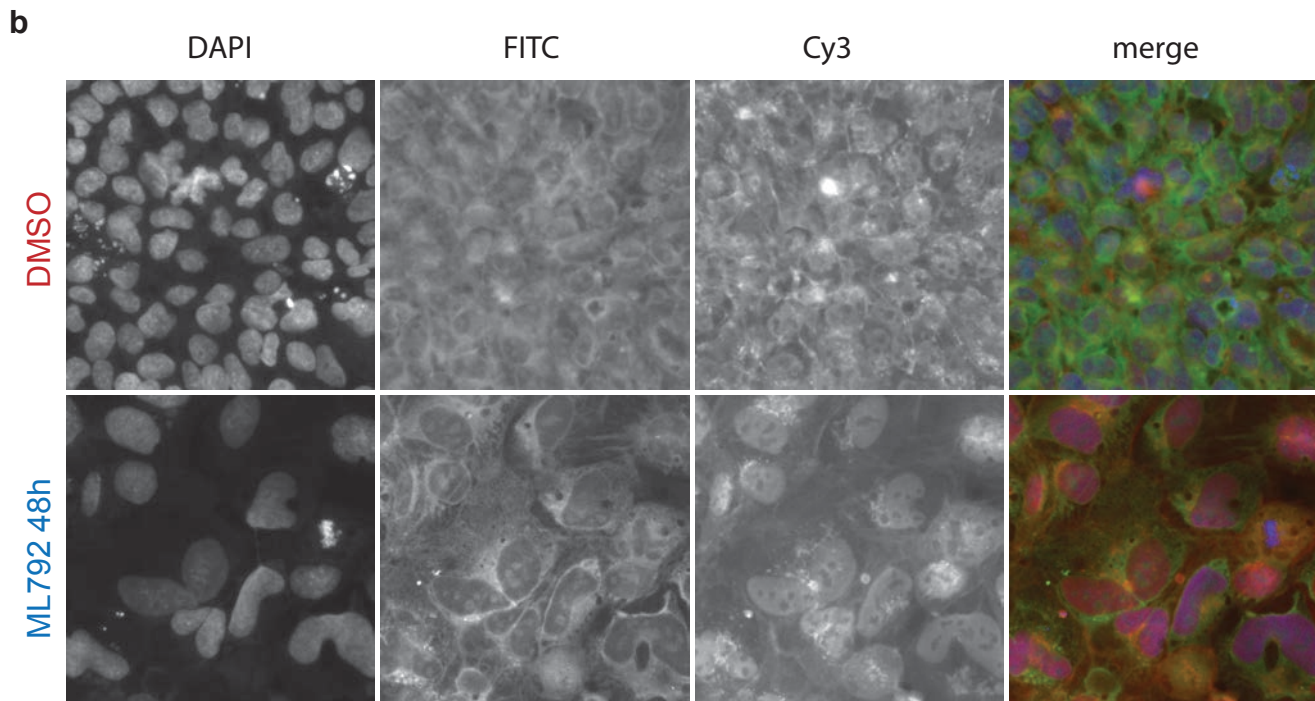
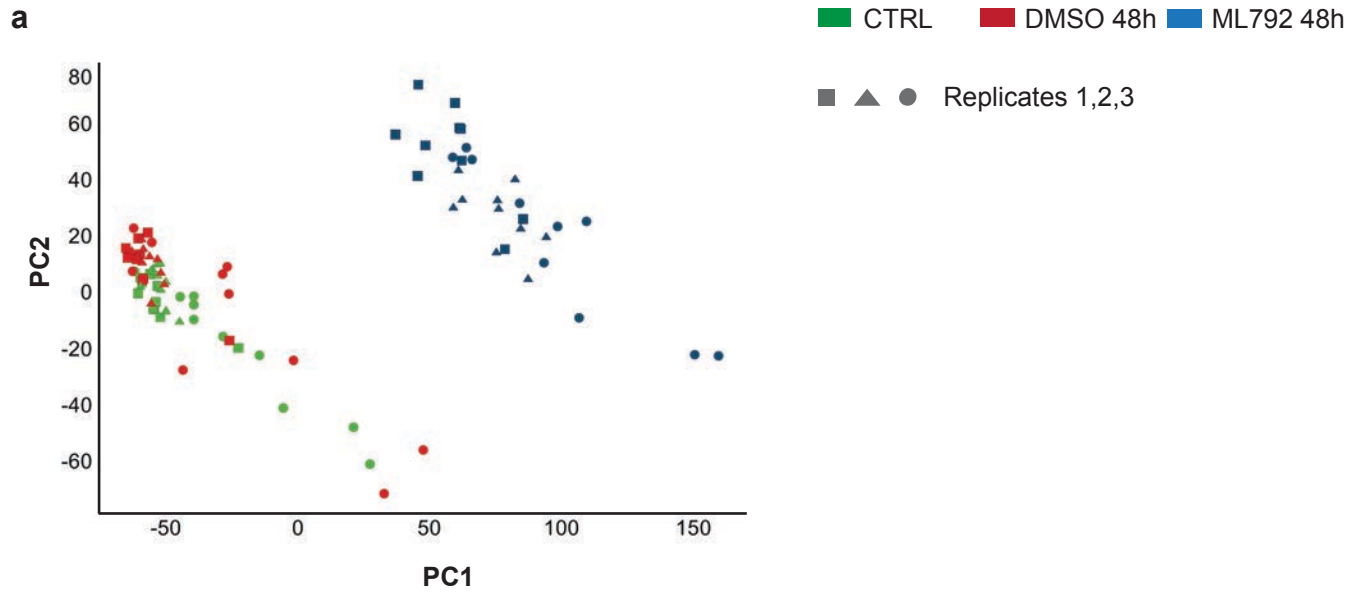
18 <sup>3</sup>Division of Computational Biology, School of Life Sciences, University of Dundee, Dundee, UK

19

20

21

## Supplementary Figure 1. (S1)



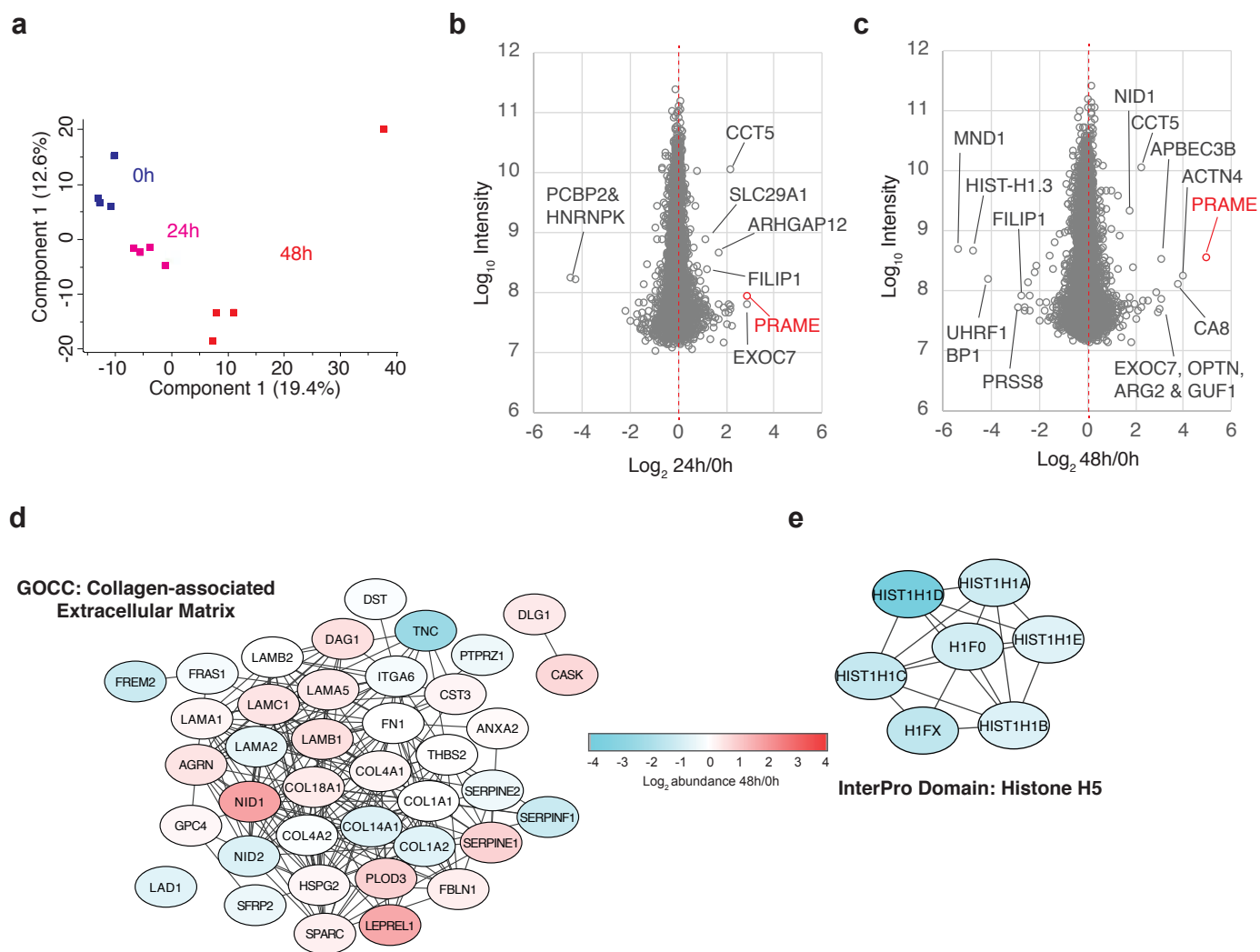
22 **Supplementary Figure 1. Nuclear and nucleolar phenotypes related to ML792 treatment.**

23 **a.** Principal component analysis of three independent cell painting experiments (replicates  
24 1,2,3) of ChiPS4 cells treated with 400nM ML792 (blue), DMSO (red) or untreated (green) for  
25 48h. **b.** Representative sample images from the cell painting experiments showing the DAPI,  
26 FITC and Cy3 channels. Cells were stained with Hoechst 33342, to reveal nuclei, Concanavalin  
27 A, to reveal endoplasmic reticulum, SYTO14 to reveal nucleoli and cytoplasmic RNAs,  
28 Phalloidin to reveal F-actin, wheat-germ agglutinin to reveal Golgi and plasma membrane and  
29 MitoTracker to reveal mitochondria.

30

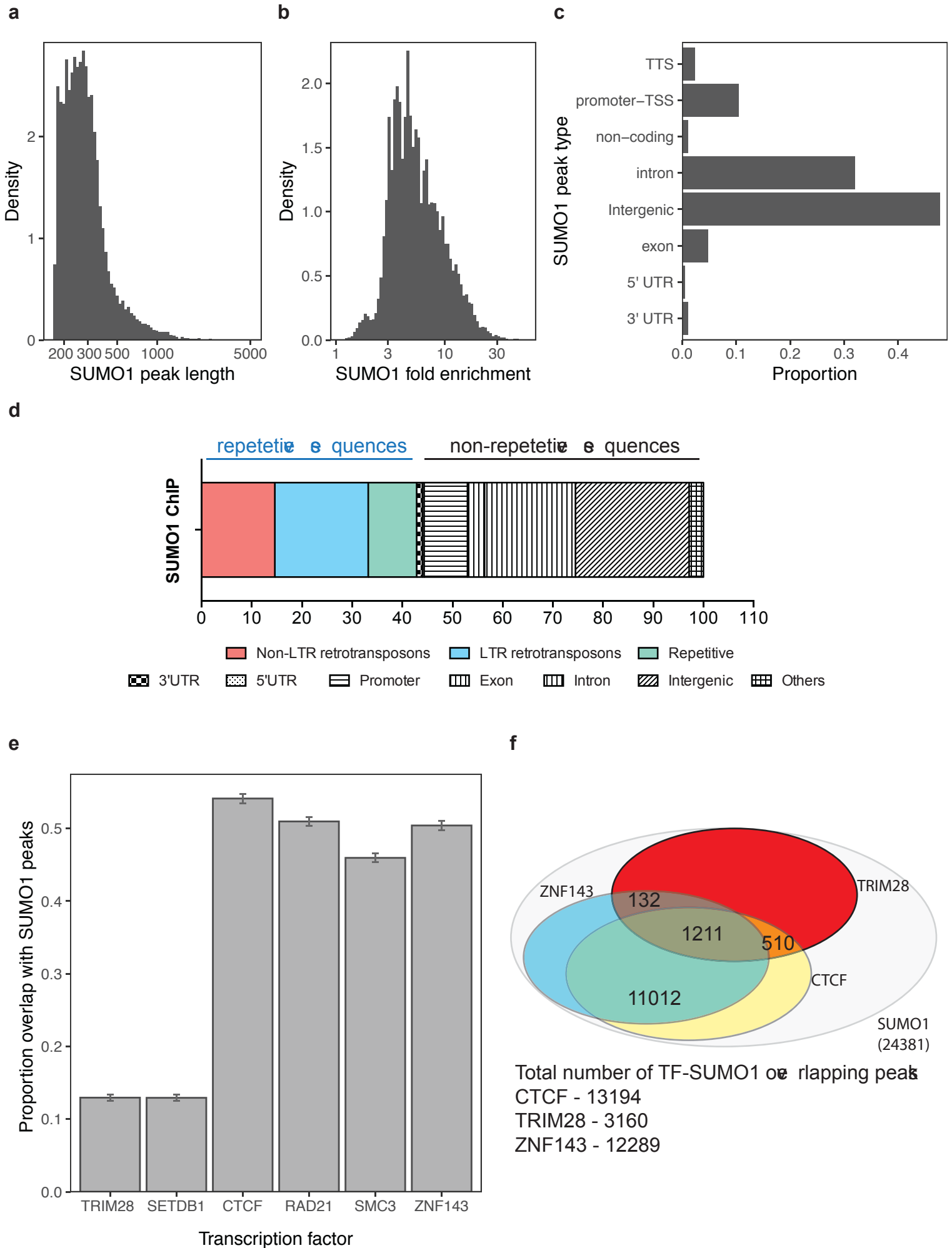
31

## Supplementary Figure 2. (S2)



32 **Supplementary Figure 2. Global deSUMOylation does not induce large changes in protein**  
33 **abundance in ChiPS4 cells over 48h. a.** Principal Component Analysis of proteomic analysis  
34 using  $\log_2$  intensity values for 4741 proteins identified from crude cell extracts from ChiPS4  
35 cells treated in quadruplicate with ML792 for the indicated times. **b.**  $\log_2$  ratio and  $\log_{10}$   
36 protein intensity data for 4741 proteins identified and quantified in crude extracts.  
37 Comparisons between untreated cells and 24h ML792. **c.** As in **b** but comparison between  
38 untreated cells and 48 h ML792. Selected outliers are indicated. **d.** Protein interaction  
39 network derived from the GOCC term 'Collagen-associated extracellular matrix', which was  
40 the only GOCC term significantly affected by 48 h ML792 treatment. All members are shown  
41 in the network which is colour-coded by  $\log_2$  48h/0h ratio. No functional group clustering was  
42 identified by STRING for 24h/0h ratio data. **e.** Protein interaction network for the InterPro  
43 group 'Histone H5'.  $\log_2$  48h/0h abundance ratio shown by colour as shown in the key.  
44  
45

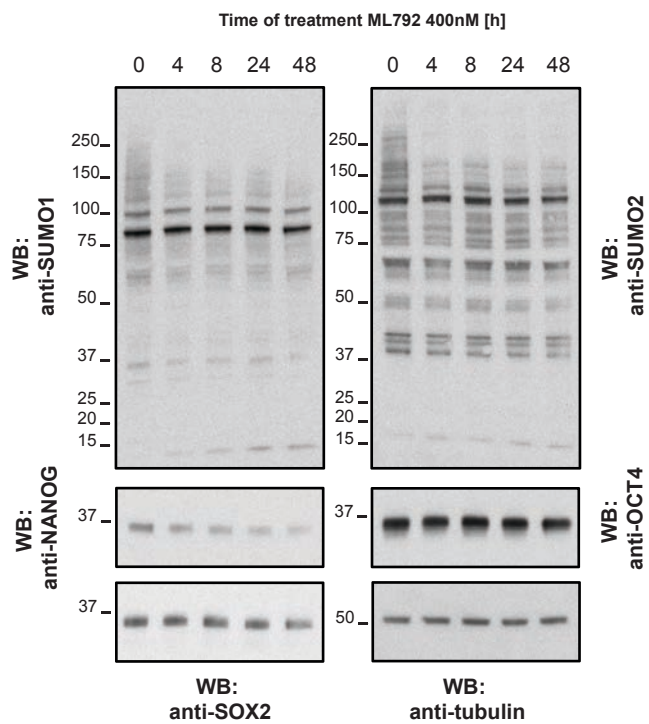
## Supplementary Figure 3. (S3)



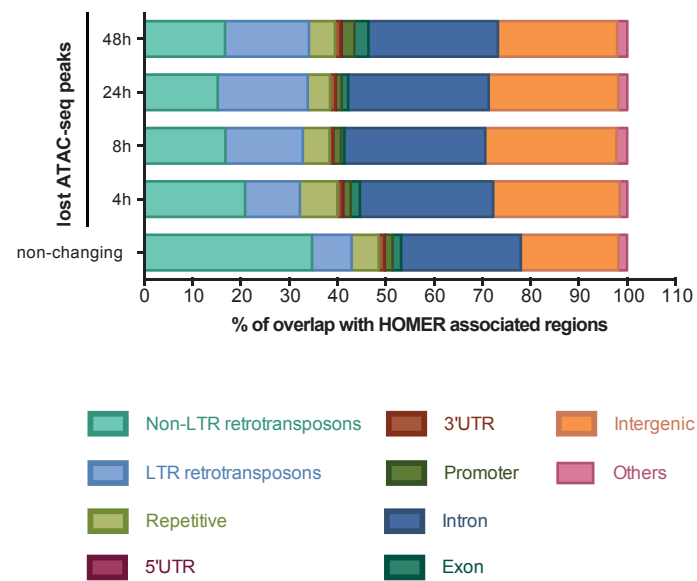
46 **Supplementary Figure 3. SUMO1 ChIP-seq peaks in untreated ChiPS4 cells. a.** Density plotted  
47 against SUMO1 peak lengths. **b.** Density plotted against fold enrichments calculated against  
48 the input samples. **c.** Proportion of SUMO1 ChIP peaks associated with various types of  
49 genomic locations were plotted based on HOMER annotations. **d.** Detailed analysis of  
50 repetitive and non-repetitive sequences using HOMER was used to calculate a percentage of  
51 overlap between those and SUMO1 ChIP-seq peaks. Different chromatin regions are  
52 represented by colours (repetitive) or grey patterns (non-repetitive) as indicated. **e.**  
53 Proportion of overlap of SUMO1 ChIP-seq peaks with various chromatin factor ChIP-seq peaks  
54 (TRIM28, SETDB1, CTCF, RAD21, SMC3, ZNF143 data for H1 hESCs were obtained from  
55 ENCODE [http://genome.ucsc.edu/cgi-](http://genome.ucsc.edu/cgi-bin/hgTrackUi?db=hg19&g=wgEncodeRegTfbsClusteredV3)  
56 [bin/hgTrackUi?db=hg19&g=wgEncodeRegTfbsClusteredV3](http://genome.ucsc.edu/cgi-bin/hgTrackUi?db=hg19&g=wgEncodeRegTfbsClusteredV3)). Error bars are 95% of confidence  
57 intervals of the proportion. **f.** Overlaps between ZNF143, CTCF and TRIM28 peaks overlapping  
58 with SUMO1 peaks were calculated using an intersect function in *bedtools* (at least 50%  
59 overlap) and plotted using Venn diagram. The exact number of peaks for each category are  
60 shown.  
61  
62

# Supplementary Figure 4. (S4)

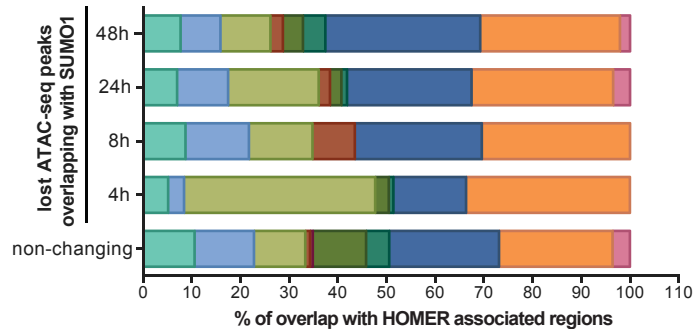
**a**



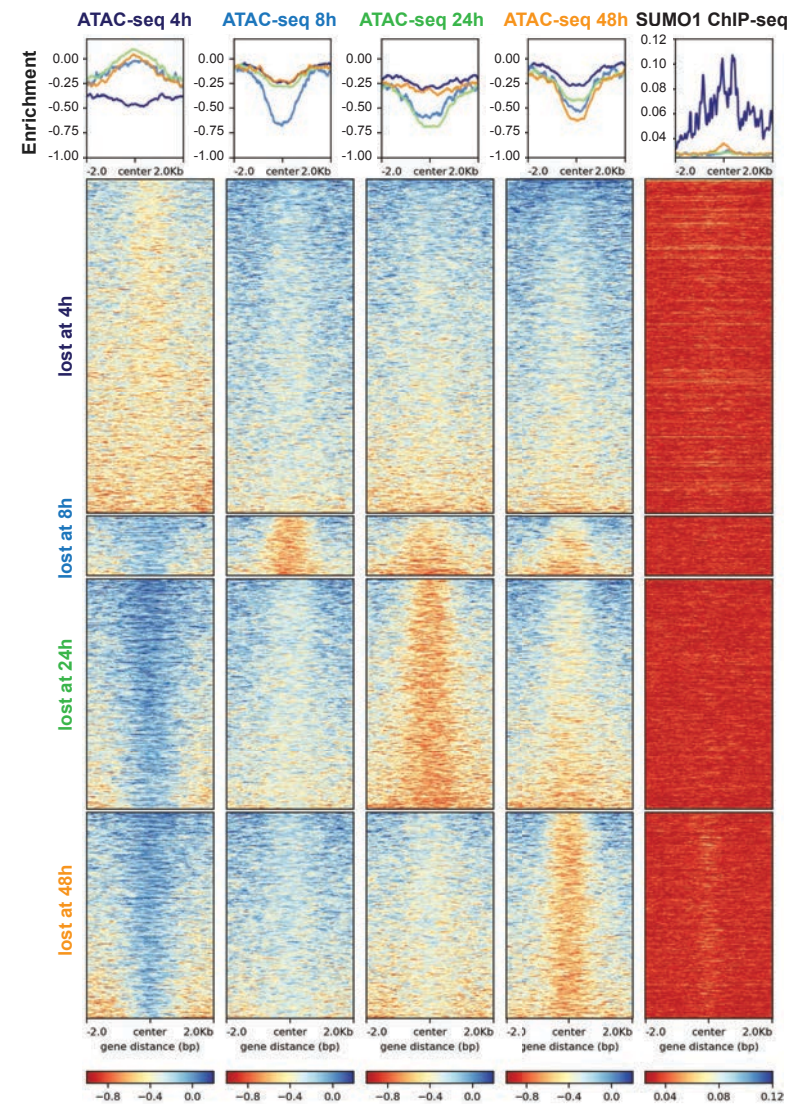
**b**



**c**



**d**





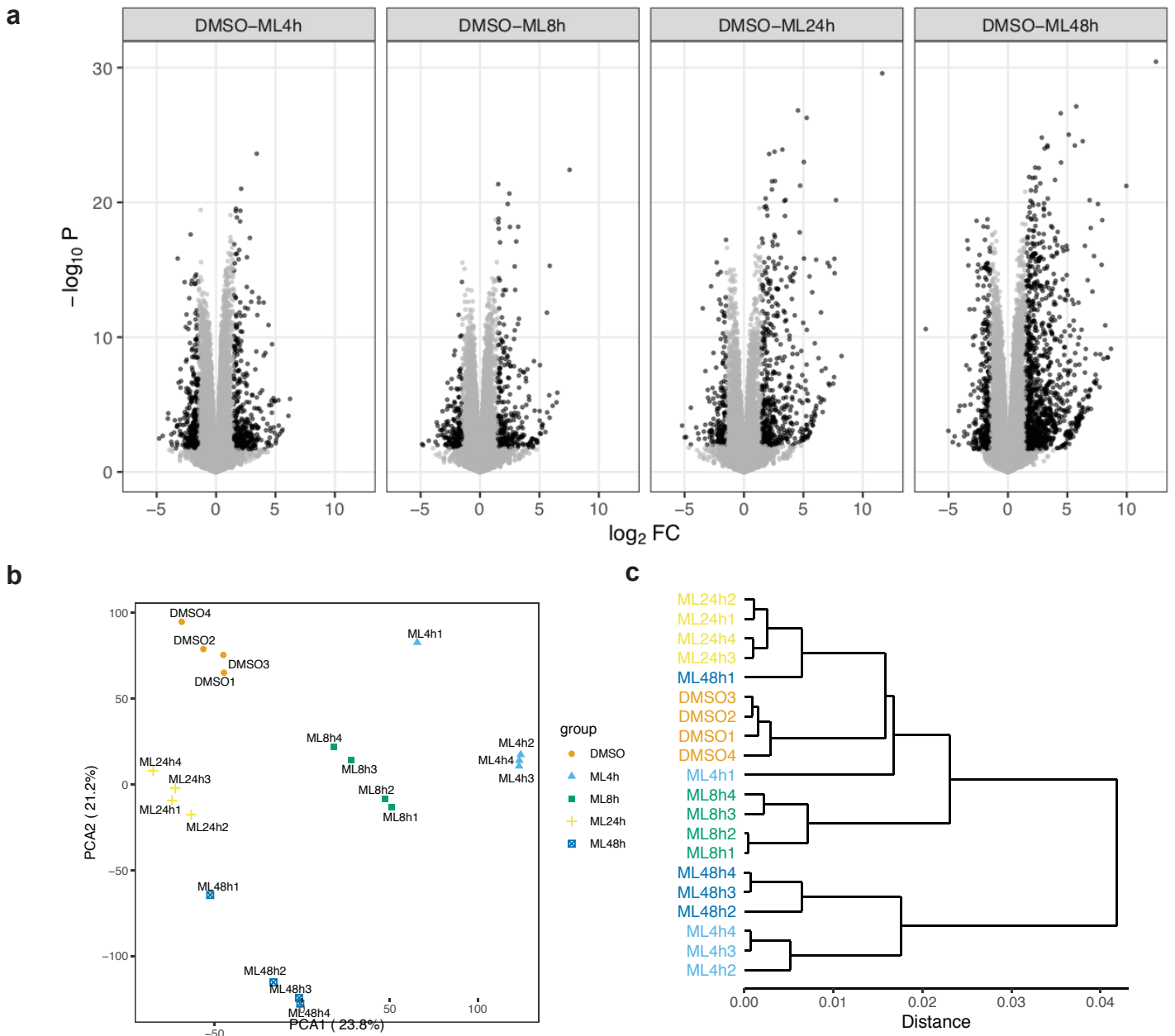
63 **Supplementary Figure 4. ATAC-seq analysis of ChIPS4 cells treated with 400nM ML792. a.**

64 Western blot analysis of total protein samples from ChIPS4 cells treated with 400nM ML792  
65 for the indicated times using antibodies against SUMO1, SUMO2, NANOG, SOX2, OCT4 and  
66 tubulin (loading control). **b.** Percentage of overlap between HOMER-based annotations of  
67 chromatin regions with lost ATAC-seq peaks at various time points and **c.** Percentage of  
68 overlap between HOMER-based annotations of chromatin regions with peaks common for  
69 SUMO1 CHIP-seq and lost ATAC-seq peaks at each time point. Different chromatin regions are  
70 represented by colours as indicated **d.** Density plots for lost ATAC-seq changes at each time  
71 point following ML792 treatment. The sites are ordered by the time at which a change of >1.5  
72 fold is first detected. The same order of genomic locations has been plotted of SUMO ChIP-  
73 seq signal. Graphs at the top represent summary plots for the ATAC-seq peaks at changing  
74 sites for the indicated time points.

75

76

## Supplementary Figure 5. (S5)



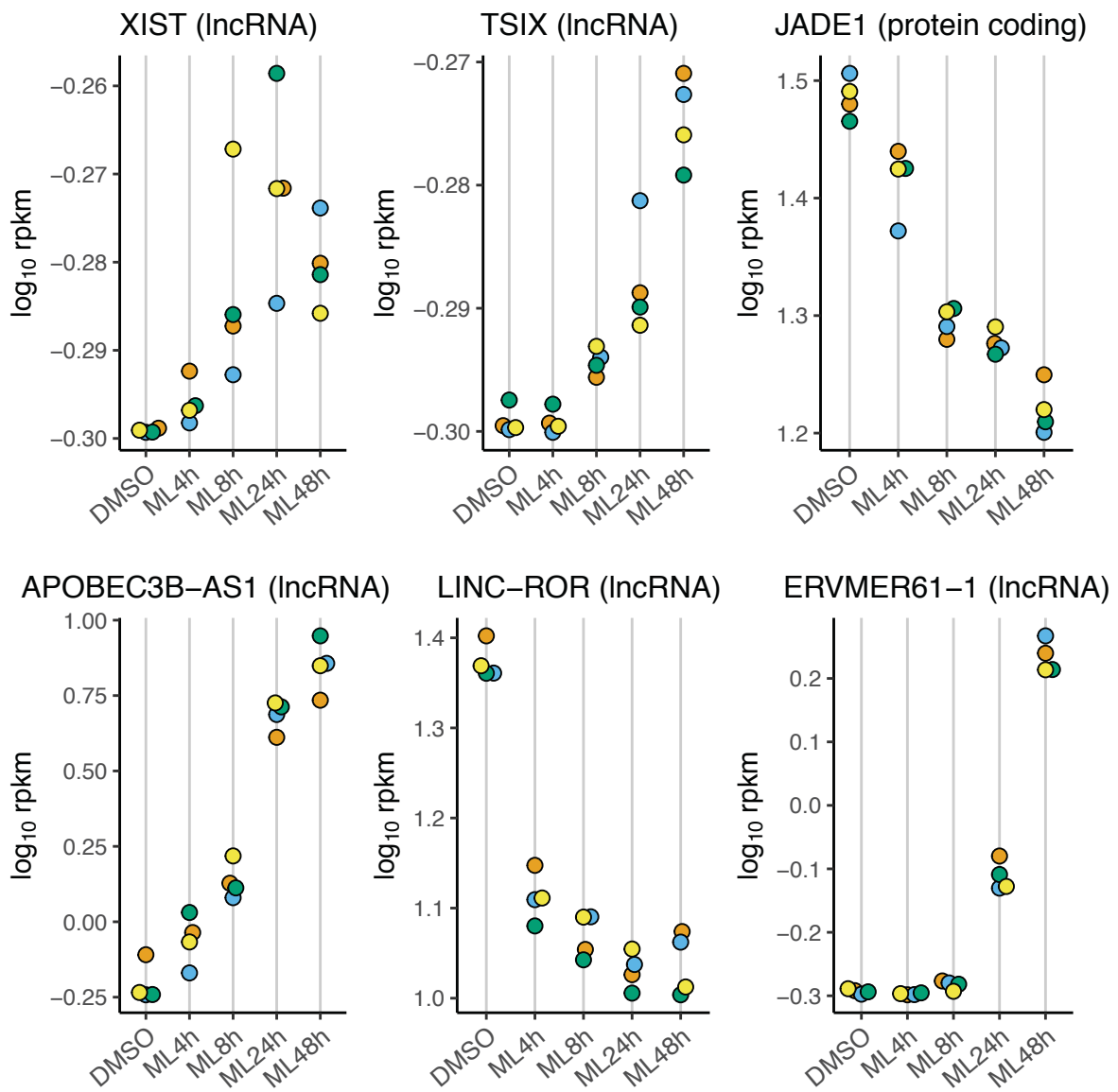
77 **Supplementary Figure 5. RNA-seq global data analysis. a.** Volcano plots showing differential  
78 expression of RNA-seq data for each time point versus DMSO. Black dots indicate  
79 “differentially expressed” genes, defined by  $FDR < 0.05$  and  $|\log_2FC| > 1.5$ . **b.** PCA  
80 decomposition of RNA-seq RPKM data. **c.** Hierarchical clustering of RNA-seq RPKM data using  
81 correlation distance.

82

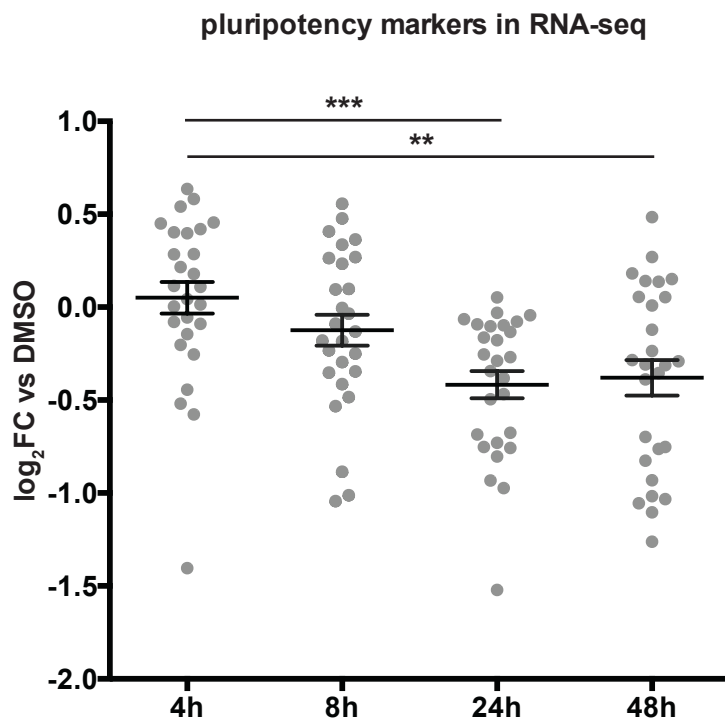
83

## Supplementary Figure 6. (S6)

**a**

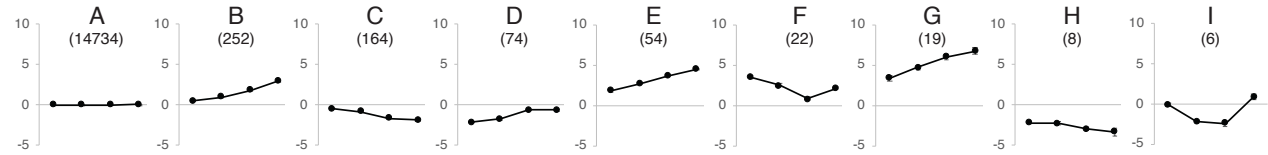
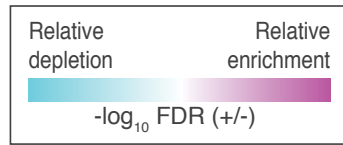


**b**



84 **Supplementary Figure 6. Expression of RNAs associated with pluripotency and ML792**  
85 **response. a.** Expression profiles of selected genes, biotype indicated in brackets. Graphs  
86 represent  $\log_{10}$ rpkm at a given time point with all replicates shown in different colours. **b.**  
87  $\log_2$  abundance ratio data for 27 markers of pluripotency comparing ML792 exposure to  
88 DMSO treated cells. The plot shows individual data points and mean with standard error of  
89 the mean for the entire set. The result of a one-way ANOVA adjusted for multiple comparisons  
90 using Holm-Sidak's method is shown. \*\*  $P < 0.01$ ; \*\*\*  $P < 0.001$ .  
91  
92

# Supplementary Figure 7. (S7)



Functional group	-log FDR_A	-log FDR_B	-log FDR_C	-log FDR_D	-log FDR_E	-log FDR_F	-log FDR_G	-log FDR_H	-log FDR_I
Keywords:Phosphoprotein	27.45	-10.93	-3.23	-2.95	-1.54	-1.55	-0.62	-0.32	-0.12
GSEA:BENPORATH_ES_WITH_H3K27ME3	-20.78	5.91	8.34	0.88	0.00	0.00	0.00	0.88	0.00
GOCC name:intracellular part	18.71	-9.36	-0.93	-2.33	-0.76	-0.13	-0.18	0.00	-0.12
Keywords:Acetylation	18.24	-5.55	-4.12	-1.45	-1.09	-0.27	0.00	0.00	0.00
GSEA:BENPORATH_EED_TARGETS	-14.80	4.99	5.67	0.49	0.04	0.00	0.00	0.00	0.00
GSEA:BENPORATH_SUZ12_TARGETS	-13.57	3.59	4.48	0.47	0.39	0.00	0.00	0.95	0.00
Keywords:Signal	-11.17	3.70	6.05	0.00	1.03	0.12	0.00	0.00	0.00
GSEA:BENPORATH_PRC2_TARGETS	-13.00	4.15	2.72	1.04	0.14	0.00	0.00	0.12	0.00
GSEA:SCGGAAGY_V\$ELK1_02	13.65	-4.74	-2.37	-0.10	-0.13	0.00	0.00	0.00	0.00
Keywords:Disulfidebond	-11.33	2.38	5.56	0.09	0.99	0.28	0.16	0.00	0.00
GOCC name:cytoplasmic part	12.22	-6.44	-0.56	-1.05	-0.47	0.00	0.00	0.00	0.00
Keywords:Glycoprotein	-9.75	2.72	5.09	0.00	0.69	0.00	0.05	0.00	0.44
GSEA:PILON_KLF1_TARGETS_DN	12.62	-2.64	-1.66	-1.42	-0.28	0.00	0.00	0.00	0.00
GSEA:DIAZ_CHRONIC_MEYLOGENOUS_LEUKEMIA_UP	11.37	-5.12	-0.81	-0.28	-0.31	0.00	0.00	0.00	0.00
GOCC name:organelle part	10.60	-4.58	-1.03	-0.74	0.00	0.00	-0.21	0.00	0.00
Keywords:Cellmembrane	-9.36	1.28	4.33	0.62	0.96	0.00	0.42	0.00	0.00
GSEA:MIKKELSEN_MEF_HCP_WITH_H3K27ME3	-10.12	2.71	3.77	0.00	0.00	0.10	0.00	0.11	0.00
GOCC name:intracellular organelle part	10.21	-4.58	-0.90	-0.65	0.00	0.00	-0.18	0.00	0.00
Keywords:Receptor	-9.38	2.34	2.04	1.57	1.00	0.00	0.00	0.00	0.00
KEGG name:Neuroactive ligand-receptor interaction	-9.08	0.06	6.29	0.00	0.05	0.00	0.00	0.00	0.00
Keywords:Referenceproteome	5.53	-9.46	0.00	-0.10	0.00	0.00	0.00	0.00	0.00
Keywords:Completeproteome	5.51	-9.40	0.00	-0.10	0.00	0.00	0.00	0.00	0.00
GSEA:MIKKELSEN_MCV6_HCP_WITH_H3K27ME3	-8.39	2.82	2.39	0.11	0.00	0.21	0.00	0.00	0.00
GSEA:GRAESSMANN_APOPTOSIS_BY_DOXORUBICIN_DN	8.88	-3.24	-0.63	-0.27	0.00	0.00	0.00	0.00	0.00
Keywords:Secreted	-6.17	1.98	2.47	-0.05	1.36	0.49	0.00	0.00	0.00
Keywords:Transducer	-7.71	2.34	0.99	0.57	0.45	0.00	0.18	0.00	0.00
GSEA:JOHNSTONE_PARVB_TARGETS_3_DN	8.22	-3.04	-0.77	-0.16	0.00	0.00	0.00	0.00	0.00
GSEA:PUJANA_BRCA1_PCC_NETWORK	8.00	-2.99	-0.37	-0.39	-0.09	0.00	0.00	0.00	0.00
Keywords:G-proteincoupledreceptor	-7.53	2.08	1.36	0.34	0.20	0.00	0.25	0.01	0.00
GSEA:REACTOME_GPCR_LIGAND_BINDING	-6.97	1.42	2.87	0.00	0.02	0.00	0.06	0.00	0.00
Keywords:Alternativesplicing	5.26	-4.83	-0.10	-0.57	0.01	-0.26	-0.19	-0.03	-0.06
GSEA:ZWANG_TRANSIENTLY_UP_BY_2ND_EGF_PULSE_ONLY	-7.35	1.34	2.40	0.07	0.00	0.00	0.09	0.00	0.00
GOCC name:organelle	6.57	-2.51	-0.06	-2.02	0.00	0.00	0.00	0.00	0.00
GOCC name:intracellular organelle	6.55	-2.58	-0.03	-1.97	0.00	0.00	0.00	0.00	0.00
GSEA:MARTENS_TRETINOIN_RESPONSE_UP	-6.66	0.51	1.46	0.28	2.20	0.00	0.00	0.00	0.00
GSEA:KEGG_NEUROACTIVE_LIGAND_RECEPTOR_INTERACTION	-6.05	0.33	3.95	0.00	0.46	0.00	0.00	0.00	0.00
GSEA:PUJANA_CHEK2_PCC_NETWORK	7.71	-2.43	-0.51	-0.08	0.00	0.00	0.00	0.00	0.00
GOBP name:cellular macromolecule metabolic process	6.95	-1.66	-1.49	-0.51	0.00	0.00	0.00	0.00	0.00
GOCC name:cytosol	6.59	-3.25	-0.47	0.00	-0.13	0.00	0.00	0.00	0.00
GOBP name:cellular metabolic process	5.77	-1.91	-0.48	-0.86	0.00	0.00	0.00	0.00	0.00
GOCC name:cell part	3.49	-5.48	0.00	0.00	0.00	0.00	0.00	0.00	0.00
GSEA:DODD_NASOPHARYNGEAL_CARCINOMA_DN	6.06	-1.68	-0.83	-0.16	-0.20	0.00	0.00	0.00	0.00
Keywords:Palmitate	-4.95	1.64	1.40	0.81	0.00	0.00	0.00	0.00	0.00
Keywords:Transmembranehelix	-4.74	0.41	2.28	0.33	0.08	0.08	0.78	0.00	0.00
GSEA:GRAESSMANN_RESPONSE_TO_MC_AND_DOXORUBICIN_DN	6.20	-1.46	-0.93	-0.06	0.00	0.00	0.00	0.00	0.00
Keywords:Transmembrane	-4.63	0.39	2.24	0.33	0.08	0.08	0.78	0.00	0.00
GOMF name:receptor activity	-6.06	1.02	1.18	0.00	0.21	0.00	0.00	0.00	0.00
GSEA:LASTOWSKA_NEUROBLASTOMA_COPY_NUMBER_DN	6.37	-1.51	-0.47	-0.07	0.00	0.00	0.00	0.00	0.00
Keywords:Cytoplasm	5.02	-2.40	-0.30	-0.32	-0.08	-0.02	-0.06	-0.02	0.00
GSEA:REACTOME_GPCR_DOWNSTREAM_SIGNALING	-5.21	0.20	2.63	0.15	0.00	0.00	0.00	0.00	0.00

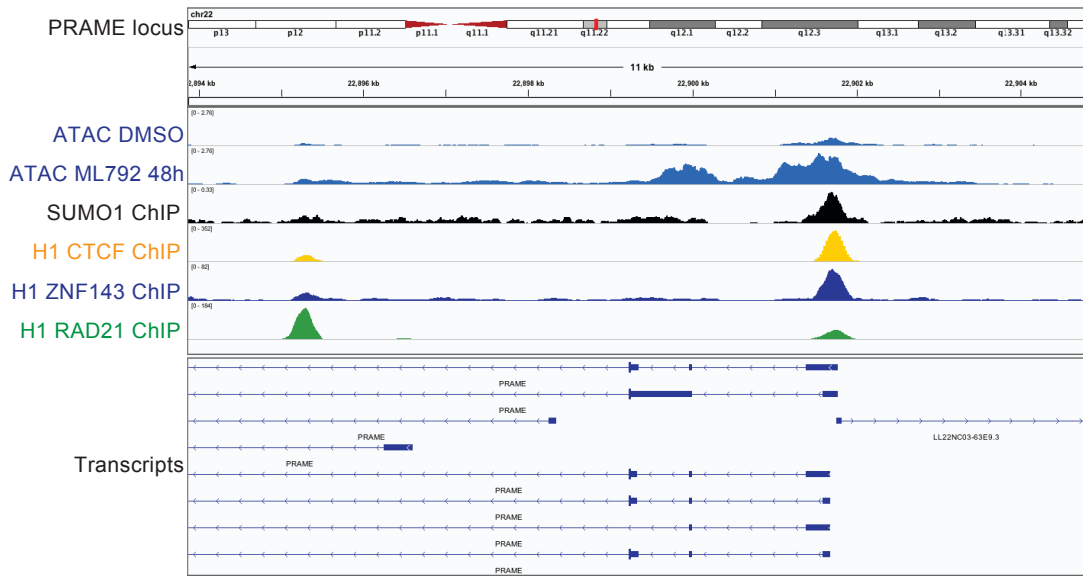
93 **Supplementary Figure 7. Summary of functional group enrichment in different RNA-seq**  
94 **clusters.** Summary of functional group enrichment in the RNAseq clusters sorted by the  
95 groups showing the most enrichment or depletion across all categories. This is the top 50  
96 most enriched/depleted groups. Negative values are depleted and positive enriched. Colour  
97 coded by degree of depletion or enrichment within each group. This analysis only uses  
98 'protein coding' genes selected from the entire RNA-seq experiment. Source data can be  
99 found in Supp. Data File 4.

100

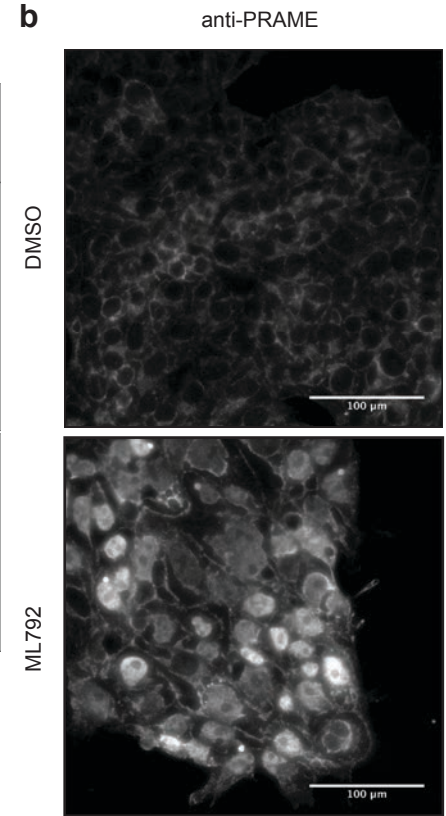
101

**Supplementary Figure 8. (58)**

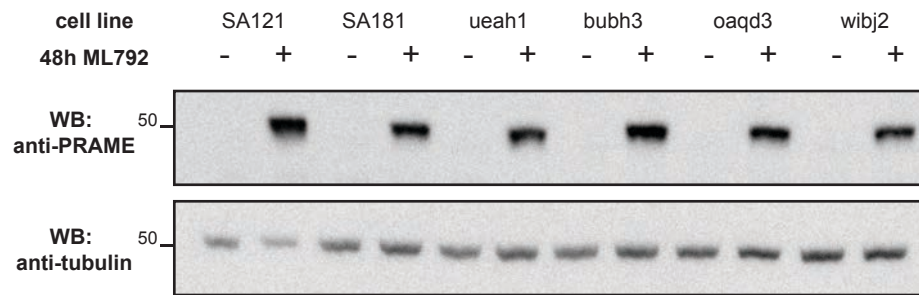
**a**



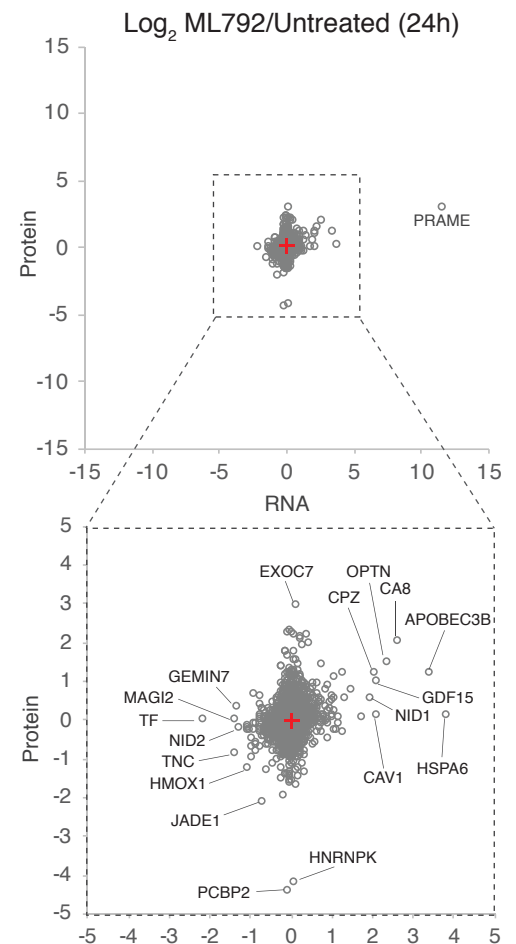
**b**



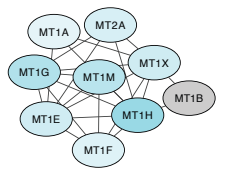
**c**



**d**

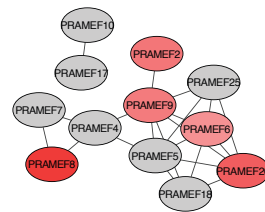


**e**

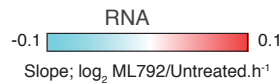


Interpro: Metallothionein domain superfamily

**f**



Interpro: PRAME family

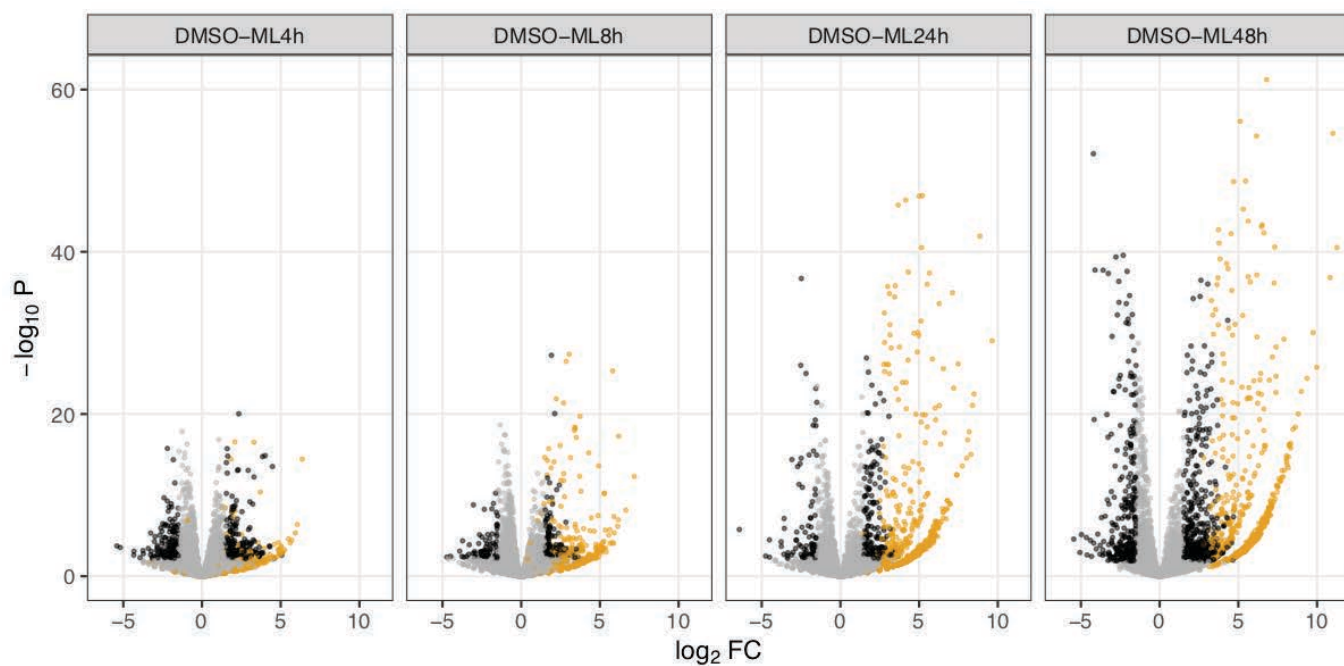




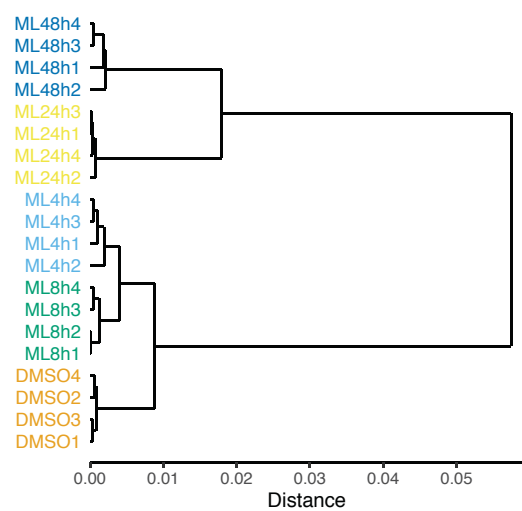
102 **Supplementary Figure 8 (S8). PRAME is strongly induced by ML792 treatment in hiPSCs. a.**  
103 Integrative Genomic Viewer was used to visualize changes in ATAC-seq (blue) at the PRAME  
104 locus in response to ML792 treatment. SUMO1 ChIP-seq signal in untreated ChiPS4 was also  
105 aligned and represented in black. ChIP-seq signals for CTCF (yellow), ZNF143 (dark blue) and  
106 RAD21 (green) from H1 hESC line were imported from ENCODE dataset and aligned with the  
107 same detailed genomic annotations. All traces of the same type were normalized and scaled  
108 in an identical way. **b.** ChiPS4 cells were treated for 48h with DMSO or 400 nM ML792, fixed  
109 and stained using DAPI and anti-PRAME antibody. IF images were obtained using a Leica DM-  
110 IRB microscope equipped with a Hamamatsu CCD camera and 20x 0.3C-Plan lens. All images  
111 contain 100  $\mu\text{m}$  scale bar. **c.** Various human stem cell lines (hESC lines: SA121, SA181; hiPSC  
112 lines: ueah1, bubh3, oaqd3, wibj2) were treated for 48h with DMSO or 400 nM ML792 and  
113 analysed by Western blot using anti-PRAME and anti-tubulin (loading control) antibodies. **d.**  
114 Scatter-plot of RNA change versus protein-level change during 24 hours ML792 treatment  
115 (4498 common entries). Full scale is shown in the upper panel, and a smaller scale is shown  
116 below. Selected outliers are indicated. **e.-f.** STRING networks for the indicated functional  
117 groups with colouring based on slopes of the RNA-seq data only. Grey nodes were not  
118 measured. Slopes were calculated using all 4 time-point ratios plus 0h = 0 fold change.  
119  
120

## Supplementary Figure 9. (59)

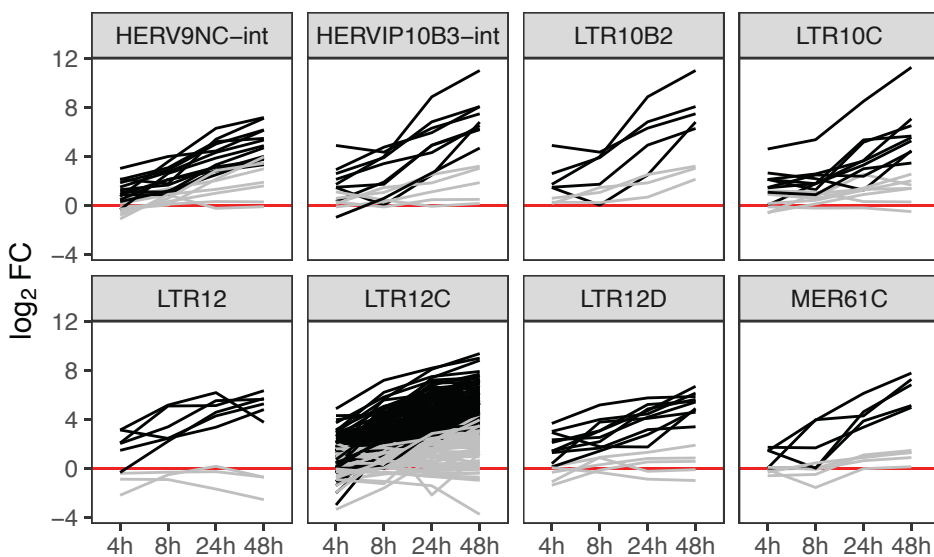
**a**



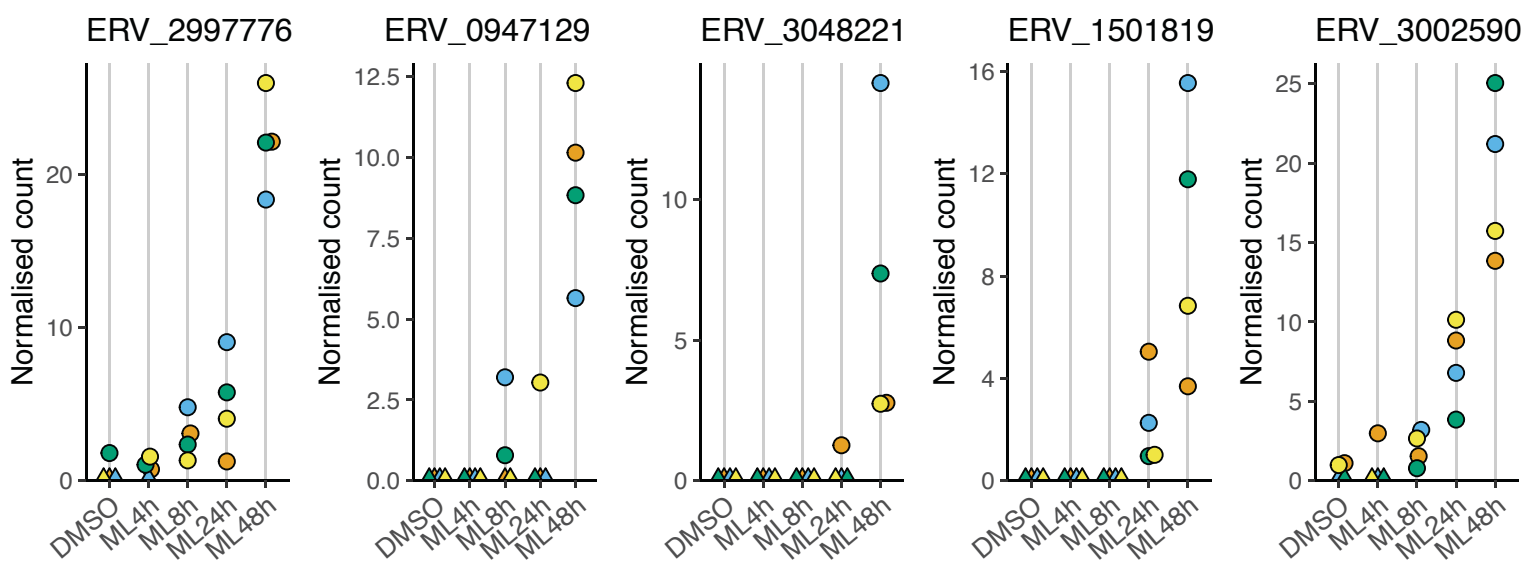
**b**



**c**



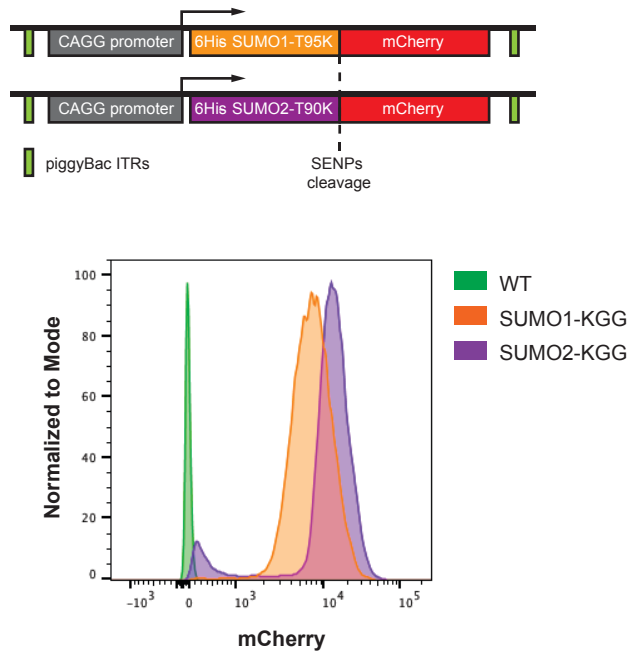
**d**



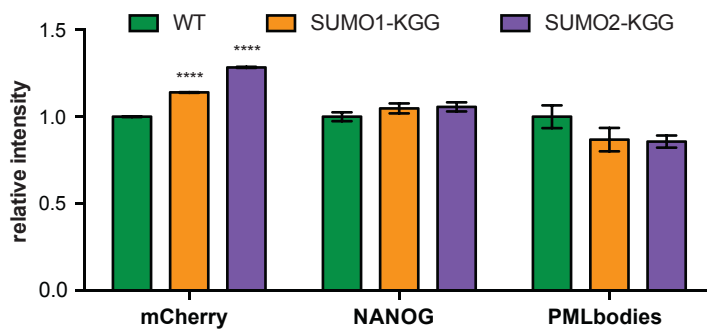
121 **Supplementary Figure 9. RNA-seq analysis of HERVs expression.** The GTF file was created  
122 with HERV loci and *STAR* was used to map RNA-seq reads to the genome without repeat  
123 masker filtering generating count reads per HERV. Of all mapped loci, 108,607 have non-zero  
124 count in at least one sample but in our analysis only HERVs (8,422) with at least 10 counts in  
125 at least one sample were considered. Differential expression of *STAR* count data was  
126 performed with *edgeR*. **a.** Volcano plots showing differential expression of HERV data for each  
127 time point versus DMSO. Black dots indicate “differentially expressed” genes, defined by FDR  
128  $< 0.05$  and  $|\log_2FC| > 1.5$ . Orange dots indicate HERVs belonging to cluster 1 (see Fig. 6d). **b.**  
129 Hierarchical clustering of HERV data using correlation distance. **c.** Time profiles ( $\log_2$  ratio  
130 between a given time point and DMSO normalised counts) of 8 selected HERV elements. Each  
131 panel shows expression from all detected loci containing this element. Black lines indicate  
132 HERVs belonging to cluster 1 (see **Fig. 6d**). **d.** Expression from five selected HERV loci. Colours  
133 indicate replicates.  
134  
135

## Supplementary Figure 10. (S10)

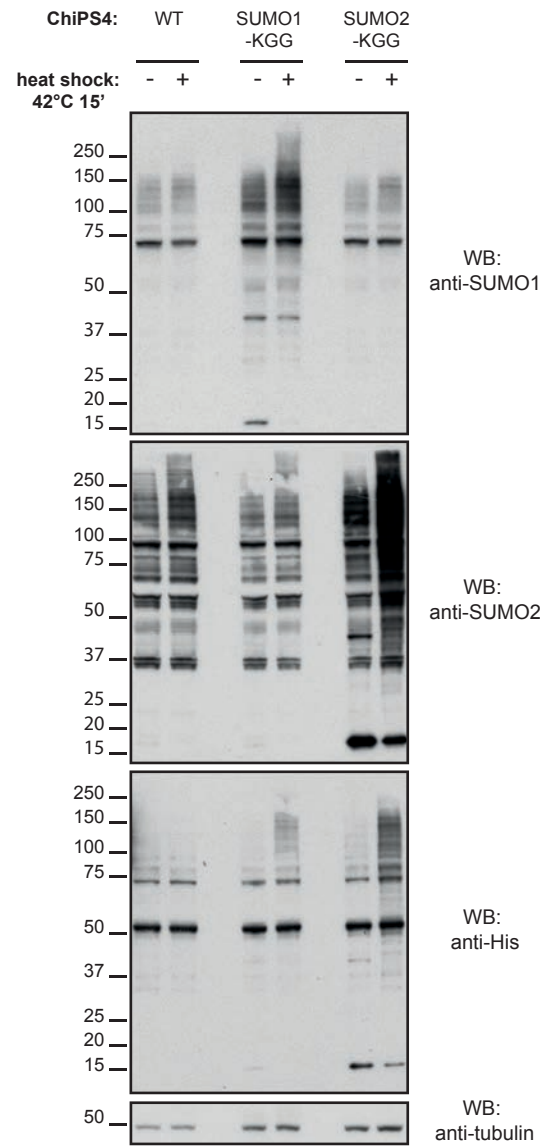
**a**



**b**



**c**



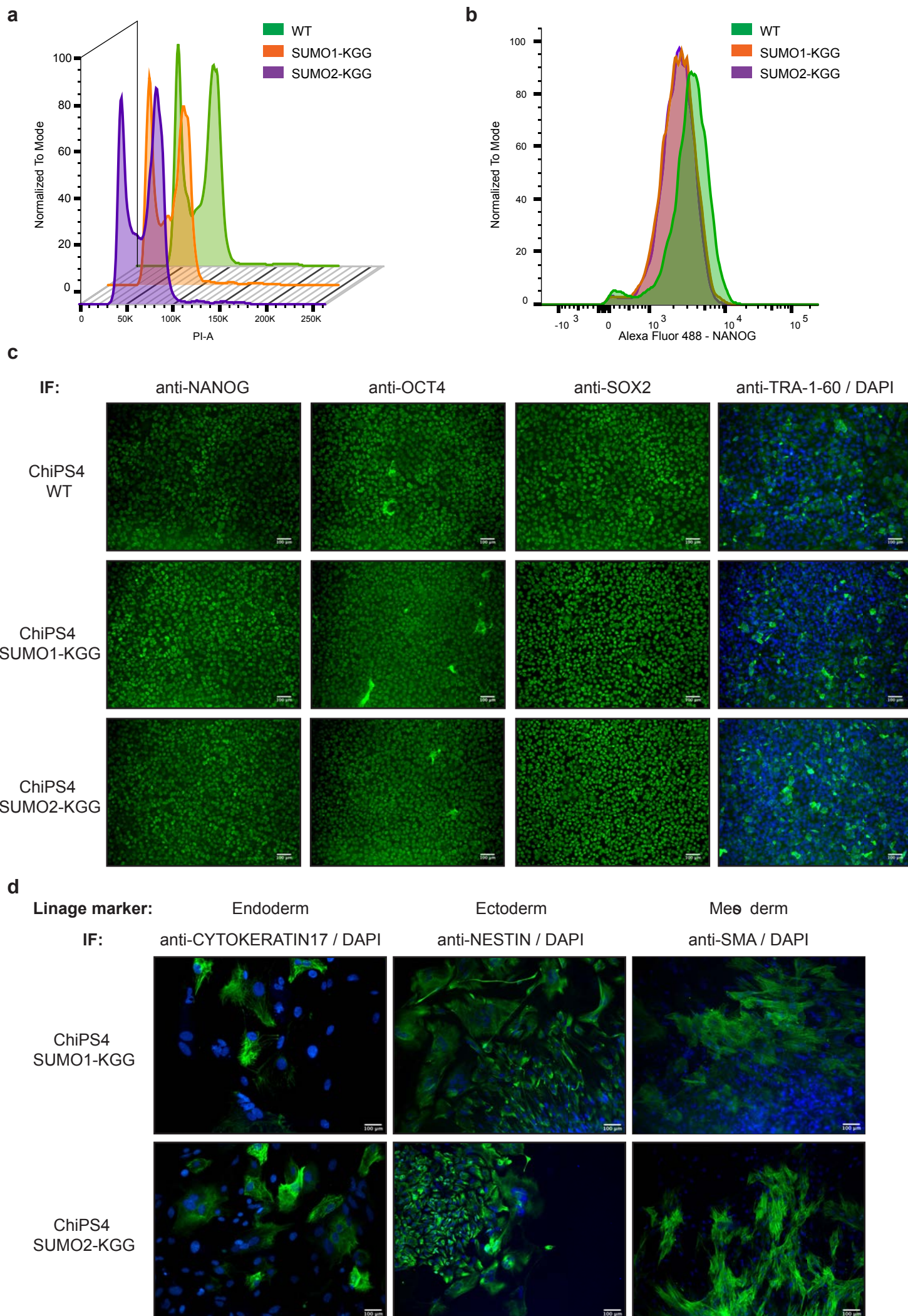
136 **Supplementary Figure 10. Generation of 6xHis-SUMO1/2-mCherry ChiPS4 cell lines. a.**

137 Design of the piggyBac constructs used for generation of ChiPS4 cell lines expressing SUMO1  
138 or SUMO2. Single cell clones were selected based on the expression levels of mCherry using  
139 flow cytometry with a view of having a similar level of His-tagged SUMO1 or SUMO2 for SUMO  
140 site proteomics experiments. **b.** Selected single cell clonal lines were further validated using  
141 High Content Screening microscopy. Cells were fixed and stained using DAPI, Cy5 Cell Mask  
142 as well as anti-NANOG and anti-PML antibodies and further assessed for mCherry expression.  
143 The result of a one-way ANOVA adjusted for multiple comparisons using Holm-Sidak's method  
144 is shown. \*\*\*\*P < 0.001 **c.** ChiPS4 WT, SUMO1-KGG and SUMO2-KGG expressing cell lines  
145 were exposed to heat shock for 15 minutes at 42°C and total protein lysates were analysed  
146 by Western blot using anti-SUMO1, anti-SUMO2/3, anti-His and anti-tubulin (loading control)  
147 antibodies.

148

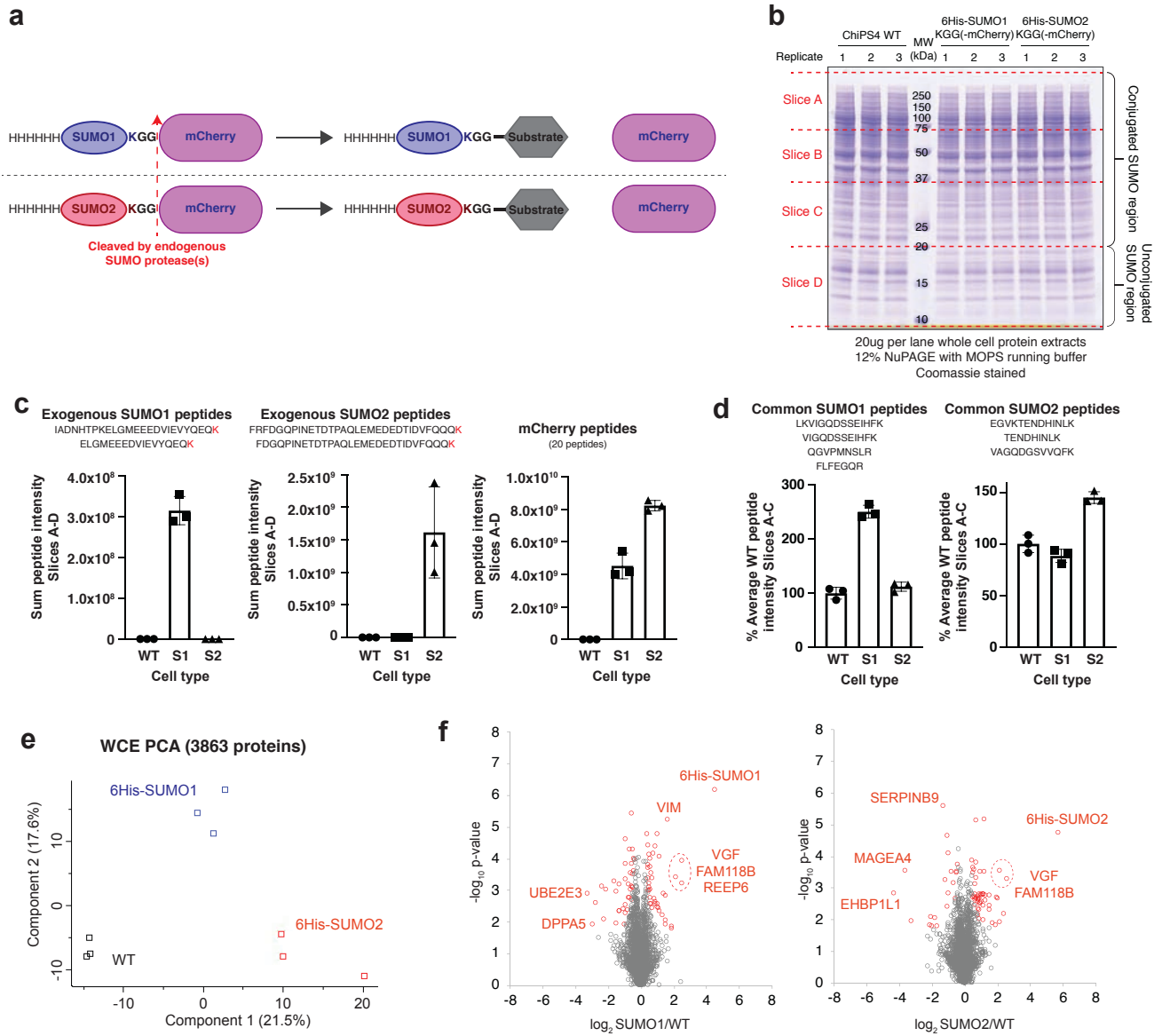
149

## Supplementary Figure 11. (S11)



150 **Supplementary Figure 11. Validation of pluripotency status of 6xHis-SUMO1/2-mCherry**  
151 **ChiPS4 cell lines. a.** Flow cytometry analysis of cell cycle and **b.** NANOG expression. **c.** IF  
152 analysis of the expression of pluripotency associated markers (NANOG, SOX2, OCT4, TRA-1-  
153 60) in ChiPS4 WT, SUMO1-KGG and SUMO2-KGG cell lines. **d.** *In vitro* differentiation potential  
154 of ChiPS4 SUMO1-KGG and SUMO2-KGG cell lines was assessed by IF staining with DAPI and  
155 specific antibodies against CYTOKERATIN 17 (Endoderm), NESTIN (Ectoderm) and SMA  
156 (Mesoderm). IF images were obtained using a Leica DM-IRB microscope equipped with a  
157 Hamamatsu CCD camera and 20x 0.3C-Plan lens. All images contain 100  $\mu$ m scale bar.  
158  
159

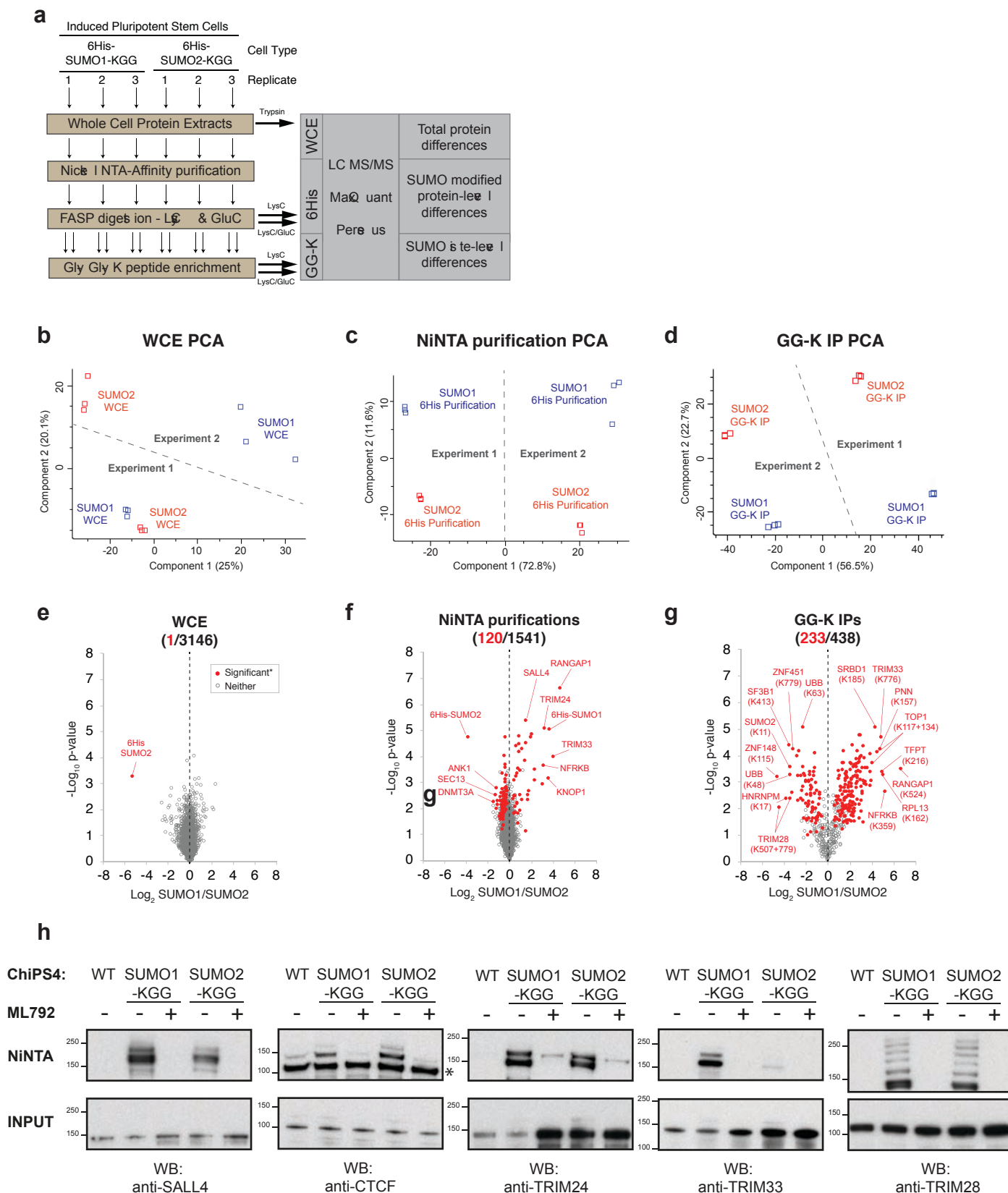
## Supplementary Figure 12. (S12)





160 **Supplementary Figure 12. Mass spectrometry-based validation of 6xHis-SUMO1/2-mCherry**  
161 **ChiPS4 cell lines.** **a.** Schematic overview of the 6His-SUMO-KGG-mCherry overexpression  
162 constructs stably expressed in ChiPS4 cells. The C-terminal mCherry protein used for cell  
163 selection is cleaved from 6His-SUMO by endogenous SUMO proteases. **b.** Coomassie stained  
164 SDS-PAGE gel fractionating whole cell protein extracts from parental ChiPS4 cells (WT) and  
165 the selected 6His-SUMO-KGG-mCherry clones. Samples were prepared in triplicate. Each lane  
166 was excised into 4 sections allowing differentiation between conjugated (slices A-C) and  
167 unconjugated (Slice D) SUMO forms. Tryptic peptides from each slice were analysed by LC-  
168 MS/MS and data processed by MaxQuant. **c.** Two peptides each from 6His-SUMO1-KGG (left)  
169 and 6His-SUMO2-KGG (centre) are specific to the exogenous construct and not the  
170 endogenous proteins. The sum of the MaxQuant Lfq peptide intensity is shown for each  
171 replicate in each cell type. Data for 20 mCherry peptides is also shown (right). **d.** Four peptides  
172 from SUMO1 (left) and three from SUMO2 (right) are common to both the endogenous and  
173 exogenous forms of the proteins. These intensities can be used to assess over-expression  
174 levels of the 6His-SUMO-KGG constructs relative to their endogenous counterparts and are  
175 presented relative to parental (WT) cell intensity. Data from slice D was omitted to allow  
176 comparisons in context of the conjugated forms of the proteins. **e.** Quantitative data from  
177 3863 proteins identified from the gel shown in **b.** were compared by principal component  
178 analysis. **f.** Numerical ratio and unpaired student's t-test results comparing WT parental cells  
179 with 6His-SUMO1-KGG-mCherry cells (left), and WT with 6His-SUMO2-KGG-mCherry cells  
180 (right). Outliers (red markers) were defined in Perseus by 5% FDR with an S0 value of 0.1 (79  
181 outliers from WT vs SUMO1 cells and 73 from WT vs SUMO2 cells - 22 common). Gene names  
182 from extreme outliers are indicated. \*MaxQuant assigned all mCherry peptides to the 6His-  
183 SUMO1-KGG-mCherry protein group, so mCherry peptides derived from 6His-SUMO2-KGG-  
184 mCherry falsely shows enrichment of the SUMO1 construct in SUMO2 cells.  
185  
186

## Supplementary Figure 13. (S13)



187 **Supplementary Figure 13. Design and quality control of SUMO site proteomics experiments.**

188 **a.** Overview of a proteomics experiment to identify IPS-specific SUMO1 and SUMO2  
189 substrates. Two experimental runs were performed with two different hiPSC lines (expressing  
190 6His-SUMO1-KGG or 6His-SUMO2-KGG), each one was performed in triplicate. Three protein  
191 fractions were analysed; whole cell extracts (WCE), NiNTA column elutions (6HIS), GlyGly-K  
192 immunoprecipitated peptide elutions (GG-K IP). All peptides were analysed by LC-MS/MS and  
193 data processed by MaxQuant. **b.-d.** Principal component analysis of MS data from the three  
194 different cell fractions. **e.-g.** Scatter plots of  $\log_2$  SUMO1/SUMO2 ratio and  $-\log_{10}$  t-test p-  
195 value for proteins or peptides detected in different cellular fractions as indicated in **a.**: **e.** WCE  
196 - Whole cell extract (measuring total protein abundance difference), **f.** NiNTA elutions  
197 (difference in proteins abundance in 6His-SUMO purifications), **g.** GGK-IP (site-level SUMO  
198 preference). \*Red markers were found to be significantly different in both experimental runs.  
199 Selected outliers are indicated. Numbers of significantly differing proteins or peptides  
200 compared to the entire set of proteins or peptides quantified are shown. **h.** NiNTA purification  
201 of His-SUMO modified proteins was performed using WT, 6His-SUMO1-KGG or 6His-SUMO2-  
202 KGG ChiPS4 cell lines that were treated with DMSO vehicle or 400 nM ML792 for 48h. Input  
203 and NiNTA elutions were analysed by Western blot using specific antibodies directed against  
204 following protein targets: SALL4, TRIM24, TRIM28, TRIM33, and CTCF. Bands for unmodified  
205 proteins that are non-specifically pulled down on NiNTA resin are labelled with an \*.