

1 **Epigenomic co-localization and co-evolution reveal a key role for**
2 **5hmC as a communication hub in the chromatin network of ESCs**

3 *Network approaches to decipher the epigenetic communication of embryonic stem*
4 *cells*

5

6 David Juan^{1,5}, Juliane Perner^{2,5}, Enrique Carrillo de Santa Pau^{1,5}, Simone
7 Marsili^{1,5}, David Ochoa³, Ho-Ryun Chung⁴, Martin Vingron², Daniel Rico^{1,6} and
8 Alfonso Valencia^{1,6}

9

10 ¹Structural Biology and BioComputing Programme, Spanish National Cancer
11 Research Center - CNIO, Melchor Fernandez Almagro 3, 28029 Madrid, Spain.

12 ²Computational Molecular Biology, Max Planck Institute for Molecular Genetics,
13 Ihnestrasse 63-73, 14195 Berlin, Germany.

14 ³European Bioinformatics Institute (EMBL-EBI), European Molecular Biology
15 Laboratory, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD,
16 United Kingdom.

17 ⁴Otto-Warburg-Laboratories Epigenomics, Max Planck Institute for Molecular
18 Genetics, Ihnestrasse 63-73, 14195 Berlin, Germany.

19 ⁵Co-first authors.

20 ⁶Corresponding authors: valencia@cniio.es, drico@cniio.es

21

22 **Summary**

23 Epigenetic communication through histone and cytosine modifications is essential for
24 gene regulation and cell identity. Here, we propose a framework that is based on a
25 chromatin communication model to get insight on the function of epigenetic
26 modifications in ESCs. The epigenetic communication network was inferred from
27 genome-wide location data plus extensive manual annotation. Notably, we found that
28 5-hydroxymethylcytosine (5hmC) is the most influential hub of this network,
29 connecting DNA demethylation to nucleosome remodeling complexes and to key
30 transcription factors of pluripotency. Moreover, an evolutionary analysis revealed a
31 central role of 5hmC in the co-evolution of chromatin-related proteins. Further
32 analysis of regions where 5hmC colocalizes with specific interactors shows that each
33 interaction points to chromatin remodelling, stemness, differentiation or metabolism.
34 Our results highlight the importance of cytosine modifications in the epigenetic
35 communication of ESCs.

36

37

38 **Introduction**

39 Intracellular and intercellular communication between proteins and/or other elements
40 in the cell is essential for homeostasis and to respond to stimuli. Communication may
41 originate through multiple sources and it can be propagated through different
42 compartments, including the cell membrane, the cytoplasm, the nuclear envelope or
43 chromatin. Indeed, a cell's identity is defined by complex communication networks,
44 involving chemical processes that ultimately modify the DNA, histones and other
45 chromatin proteins.

46

47 It has been proposed that multiple histone modifications confer robustness and
48 adaptability to the chromatin signaling network (Schreiber & Bernstein, 2002). In
49 fact, it is now clear that the combination of different histone marks defines the
50 epigenomic scaffolds that affect the binding and function of other epigenetic elements
51 (e.g., different protein complexes). The increasing interest in characterizing the
52 epigenomic network of many biological systems has led to an impressive
53 accumulation of genome-wide experimental data from distinct cell types. This
54 accumulation of experimental data has meant that the first chromatin signaling co-
55 localization networks of histone marks and chromatin remodelers could be inferred in
56 the fly (van Bemmelen *et al*, 2013) and at promoters in human (Perner *et al*, 2014). In
57 addition, a variety of cytosine modifications have emerged as potentially important
58 pieces of this 'chromatin puzzle', such as 5-methylcytosine (5mC), 5-
59 hydroxymethylcytosine (5-hmC), 5-formylcytosine (5-fC) and 5-carboxylcytosine (5-
60 caC: Ficiz *et al*, 2011; Pastor *et al*, 2011; Williams *et al*, 2011; He *et al*, 2011; Ito *et al*,
61 *et al*, 2011; Raiber *et al*, 2012). However, Nevertheless, the role of these modifications
62 in epigenetic signaling is not yet clear (Pfeifer *et al*, 2013; Liyanage *et al*, 2014; Moen
63 *et al*, 2015). We are still far from understanding the epigenomic "syntax" and how
64 these and the other elements involved in epigenomic communication shape the
65 functional landscape of mammalian genomes.

66

67 Evolutionary information can be used to discern the basis of meaningful
68 communication in animals (Smith & Harper, 2003). Communication frequently

69 occurs among mutualistic and symbiotic species, as the evolution of communicative
70 strategies requires co-adaptation between signal production/emission and signal
71 reception/interpretation (Smith & Harper, 2003; Scott-Phillips, 2008). Similarly, the
72 continuous adaptation of living organisms to different scenarios requires a fine-tuning
73 of molecular communication. As a consequence, the conservation of communication
74 pathways is often challenged by ever-changing selection pressures potentially leading
75 to molecular co-evolution between intercommunicating proteins. Interestingly, long-
76 standing protein co-evolution can be reliably detected through directly correlated
77 evolutionary histories. In fact, co-evolutionary analysis has successfully identified
78 biologically relevant molecular interactions at different levels of detail (de Juan *et al*,
79 2013).

80

81 Here, we establish a new framework to rationalize and study epigenomic
82 communication. This framework combines network-based analyses and an
83 evolutionary characterization of the interactions of chromatin components derived
84 from high-throughput data and literature mining. In particular, we followed a systems
85 biology approach to investigate the functional interdependence between chromatin
86 components in mouse embryonic stem cells (ESCs), whereby changes to their
87 epigenome control a very broad range of cell differentiation alternatives. We
88 constructed the epigenetic signaling network of ESCs as a combination of high-
89 quality genomic co-localization networks of 77 different epigenomic features:
90 cytosine modifications, histone marks and chromatin-related proteins (CrPs) extracted
91 from a total of 139 ChIP-seq experiments. We labeled histone marks and cytosine
92 modifications as *signals* and we classified the proteins that co-localize with them as
93 their *emitters* (writers or erasers) or *receivers* (readers) based on information in the
94 literature (**Figure 1**). To our knowledge the resulting communication network is the
95 most complete global model of epigenetic signaling currently available and therefore,
96 we propose it to be a valuable tool to understand such processes in ESCs.

97

98 By analyzing this network, we found 5hmC to be a key node that mediates
99 communication between different regions of the network. In addition, our co-
100 evolutionary analysis of this network identified 5hmC as a central node that connects

101 most co-evolving CrPs. Exploration of 5hmC-centered communication revealed that
102 specific co-localization of 5hmC with the TET1, OGT, ESRRB and LSD1 produces
103 alternative partner-specific activity, such as chromatin remodeling, cell stemness and
104 differentiation, and energy metabolism. Thus, we propose that 5hmC acts as a central
105 signal in ESCs for the self-regulation of epigenetic communication.
106

107 **Results**

108 **Inference of the chromatin signaling network in mouse ESCs**

109 We built an epigenetic signaling network in mESCs through a two-step process. First,
110 we inferred the network connectivity based on co-localization in the genome-wide
111 distribution of chromatin components. In this analysis, we included 139 ChIP-Seq,
112 MEDIP and GLIB assays for 77 epigenetic features (3 cytosine modifications, 13
113 histone marks and 61 CrPs: **Table S1**). Accordingly, we employed a method
114 described recently (Perner *et al*, 2014) that reveals putative direct co-dependence
115 between factors that cannot be “explained” by any other factor included in the
116 network. Thus, we detected only relevant interactions in different functional
117 chromatin domains (see **Experimental Procedures** for details).

118

119 Second, we annotated the direction of the interactions in the network (as shown in
120 **Figure 1**). For this, we relied on previously reported experimental evidence. This
121 evidence can be roughly summarized within two possible scenarios: (1) Protein A is a
122 known *writer* or *eraser* of signal B; (2) Alterations to the genome-wide distribution of
123 protein A (e.g., through its knock-out) affect the distribution of signal B in the
124 genome. In the absence of any such evidence, proteins were defined as receivers of
125 the interacting signal.

126

127 We recovered an epigenetic communication network (**Figure 2**) with 236 connections
128 between 68 nodes, the latter represented by cytosine modifications, histone marks or
129 CrPs. The network contains 192 positive interactions (co-localizing features, 81.4%)
130 and 44 negative interactions (mutually exclusive features, 18.6%). A web interactive

131 browser of the global co-localization network enables users to explore the interactions
132 among these chromatin components in more detail (see
133 <http://dogcaesar.github.io/epistemnet>).

134

135 Our approach detected 115 direct CrP-CrP interactions that are mostly due to protein
136 complexes given that these components coincide at chromatin. These include
137 complexes such as Polycomb (RYBP/CBX7/PHF19/SUZ12/EZH2), Cohesin
138 (RAD21/SMC1/SMC3), Mediator (MED1/MED12/NIPBL), the nucleosome
139 remodeling deacetylase MI2/NuRD complex (MI2B/LSD1/HDAC1/HDAC2) and
140 CoREST/Rest (Rest/CoREST/RYPB: **Figure 2**).

141

142 In order to understand the epigenetic interaction network and its activity as a
143 communication system, we focused our analyses on directional *emitter-signal* and
144 *signal-receiver* associations. Based on the experimental information extracted from
145 the literature, we established “communication arrows” from “emitter-CrPs” to their
146 signals and from the signals to their epigenetic “receiver-CrPs”. We established
147 (52.5%) directional interactions. Of those, 56 edges involve an epigenetic emitter and
148 a signal (all experimentally supported) and 68 edges connect a receiver and a signal
149 (with 27 directions supported experimentally). In total we identified 8 emitter-CrPs,
150 17 receiver-CrPs and 18 CrP nodes that can act simultaneously as emitters and
151 receivers of different signals.

152

153 The hubs of a network are highly connected nodes that facilitate the networking of
154 multiple components. Directional edges allowed us to distinguish between two types
155 of hubs: in-hubs (nodes with a large number of incoming arrows) and out-hubs (with
156 a large number of outgoing arrows). Not surprisingly, the main in-hub was RNA
157 polymerase II with S2 phosphorylation of the C-terminal (RNAPII_S2P). Indeed, 9
158 out of 16 signals in the network pointed to this form of RNAPII, which is involved in
159 transcriptional elongation and splicing (**Figure S1A**).

160

161 By contrast, we found two main out-hubs in the network revealing a different aspect
162 of epigenetic regulation. The main hubs that accumulated connections with receivers

163 were H3K79me2 (12) and 5hmC (10: **Figures 3A-B and S1B**). H3K79me2 is
164 involved in transcription initiation and elongation, as well as promoter and enhancer
165 activity, suggesting that it is a key signal for different aspects of transcriptional
166 regulation. Interestingly, two groups of transcription factors (TFs) were connected to
167 H3K79me2: one composed of TCF3, OCT4, SOX2 and NANOG; and another that
168 contains CMYC, NMYC, STAT3, KLF4, TCF2L1 and E2F1.

169

170 5hmC is particularly interesting as it is thought to be a key element in different
171 processes even though its function in gene regulation remains controversial (Pfeifer *et al*,
172 2013; Liyanage *et al*, 2014). Whereas initially related to gene activation (Song *et al*,
173 2011), others claimed that 5hmC associates with weakly expressing poised
174 promoters (Pastor *et al*, 2011; Williams *et al*, 2011), while both roles were elsewhere
175 claimed to be possible depending on the context (Wu *et al*, 2011). In addition, 5hmC
176 was shown to play a major role in enhancer activation (Stroud *et al*, 2011; Szulwach
177 *et al*, 2011) or silencing (Choi *et al*, 2014). This apparent controversy could be
178 explained by the role of 5hmC as a central node of the communication network.
179 Indeed, 5hmC is the node traversed by the highest number of paths between nodes,
180 which implies that this node concentrates the information flow of the mESC network
181 (**Figure S2**).

182

183 We further confirmed this influential role of 5hmC in our chromatin network applying
184 an algorithm originally devised to rank the relevance of web pages in the internet by
185 using *global* link information. In brief, influential nodes are those from which
186 information easily spreads out to the rest of the network, while popular nodes gather
187 information from many regions of the network. Comparing the nodes' influence and
188 popularity, we can clearly identify and distinguish between very influential nodes and
189 very popular nodes (**Figure 3C**). These results highlight the importance of
190 directionality in the network structure. The most popular node is RNAPII_S2P,
191 suggesting that transcription is the main outcome.

192

193 Conversely, 5hmC shows the highest influence score, meaning that it is a signal
194 transmitted to many receivers that, in turn, emit signals with a strong outflow to the

195 network. The most influential CrPs, TET1 and LSD1, are also emitters of 5hmC. The
196 relevance of 5hmC and RNAPII_S2P is robust to biological and methodological
197 issues, as verified by measuring the effect of directionality miss-assignments (**Figure**
198 **S5**) and random edges removal (see **Figure 3D**). As the importance of these nodes on
199 epigenetic communication appears to be so clear in the specific case of mESCs, we
200 investigated to what extent it could have constrained the evolution of the related CrPs
201 in metazoans.

202

203 **Co-evolution among chromatin components**

204 Cell stemness evolved very early in metazoan evolution and it is a critical
205 phenomenon that enhances the viability of multicellular animals (Hemrich *et al*,
206 2012). Thus, it can be assumed that CrP-mediated communication in stem cells has
207 also been essential for metazoan evolution. As co-evolution consistently reflects
208 important functional interactions among conserved proteins (de Juan *et al*, 2013), we
209 studied the signatures of protein co-evolution within the context of the epigenetic
210 communication network in stem cells. We focused our analysis on the CrPs in the
211 network for which there is sufficient sequence and phylogenetic information in order
212 to perform a reliable analysis of co-evolution (de Juan *et al*, 2013). We extracted
213 evolutionary trees for 59 orthologous CrPs in our epigenetic communication network
214 and calculated their degree of co-evolution. To disentangle the direct and
215 uninformative indirect evolutionary correlations, we developed a method that
216 recovers protein evolutionary partners based on a maximum-entropy model of
217 pairwise interacting proteins (see **Experimental Procedures**).

218

219 Using this approach, we retrieved 34 significant co-evolutionary interactions among
220 54 CrPs (**Table S3**). A total of 27 co-evolved relationships were identified as
221 functional interactions by independent experimental evidence from external
222 databases, from the literature and/or from our communication network (**Table S3**).
223 These co-evolutionary associations reflected the evolutionary relevance of different
224 epigenetic communication pathways that might be at play in essential, evolutionary
225 maintained cell types like ESCs.

226

227 We identified epigenetic signals that connect CrPs related by co-evolution (i.e.: those
228 connecting co-evolving pairs) and we considered the historically influential signals as
229 those that were best connected in a co-evolutionary filtered network. This co-
230 evolutionary filtered network was obtained by maintaining the pairs of CrPs that both
231 co-evolve and that are included in a protein/signal/protein triplet (**Figure 4**). Co-
232 evolving CrP pairs are not evenly distributed in the epigenetic communication
233 network but rather, we found a statistically significant correspondence between
234 signal-mediated communication and co-evolution for H3K4me2, H3K4me3 and
235 5hmC (p-value < 0.05, see **Experimental Procedures**). Of these, 5hmC mediates
236 communication between four different co-evolving pairs that involve seven different
237 CrPs (**Figure 4**), clearly standing out as the epigenetic signal connecting more co-
238 evolving CrPs. Notably, the three positively co-occurring emitters of 5hmC (TET1,
239 OGT and LSD1) co-evolved with three different receivers (MBD2, TAF1 and
240 SIN3A). Thus, from the combination of the 5hmC interactors (**Figure 3A**), three
241 specific emitter/signal/receiver triplets with coordinated evolution were identified:
242 LSD1-5hmC-SIN3A, TET1-5hmC-MBD2 and OGT-5hmC-TAF1. In addition, we
243 detected co-evolution between the 5fC-emitter BRG1 and the 5fC-receiver NIPBL.

244

245 The case of MBD2 and TET1 is particularly interesting given the biological activities
246 of these proteins. One of the key functions of TET1 is the oxidation of 5mC, while
247 MBD2 is a methyl-binding domain protein (MBD) that shows higher binding affinity
248 to 5mC than to 5hmC (Baubec *et al*, 2013). In addition, MBDs are thought to
249 modulate 5hmC levels, inhibiting TET1 by their binding to 5mC (Hashimoto *et al*,
250 2012). The co-evolution of MBD2 and TET1 suggests certain dependence between
251 the mechanisms that maintain 5mC and 5hmC at different epigenomic locations in
252 ESCs.

253

254 The well-known TET1 interactors OGT and SIN3A each co-evolved with a different
255 CrP: TAF1 and LSD1, respectively. OGT co-occurs with 5hmC while TAF1 binding
256 is significantly enriched in 5hmC depleted regions. Similarly, LSD1 positively
257 interacts with 5hmC while its co-evolving partner SIN3A was found in a pattern that

258 is mutually exclusive to 5hmC. As in the case of TET1 and MBD2, these results
259 suggest the remarkable influence of 5hmC on the differential binding of CrPs to
260 distinct genomic regions in the ESC epigenome during metazoan evolution.

261

262 Accordingly, these results confirmed our working hypothesis that some chromatin
263 proteins interconnected via epigenetic signals have evolved in a concerted manner.
264 Interestingly, our results also suggest that 5hmC is a communication hub as it
265 connects processes that have been coordinated during metazoan evolution.

266

267 **Functional modularization of the network reveals protein complexes and** 268 **star-shaped structures**

269 We have shown that 5hmC is the most influential signal in the ESC epigenetic
270 communication network and that it mediates the communication between CrPs that
271 have co-evolved in Metazoa. Recent research has shown that the genomic localization
272 of certain combinations of core epigenetic features allows different chromatin states
273 associated with functional processes to be reliably identified (Filion *et al*, 2010; Ernst
274 & Kellis, 2010). Here, we examined how the positive interactions in the network are
275 distributed in relation to these different functional contexts. In particular, we focused
276 on the modules of co-localizing chromatin components with similar peak frequencies
277 that were associated with the diverse chromatin states in ESCs (**Figure S6**).

278

279 We found 15 groups of interactions that yielded sub-networks associated with
280 distinctive functional chromatin profiles (**Figure 5**). These **chromatin** context-
281 specific **networks** (*chromnets*) were made up of CrPs and epigenetic signals that
282 tended to co-exist in the different chromatin states at a similar frequency in ESCs. We
283 found that most chromnets could be classified into two groups: protein complexes and
284 communication chromnets. Specific examples of protein complexes chromnets were
285 Polycomb (CBX7/PHF19/SUZ12/EZH2) in chromnet-5, Cohesin
286 (RAD21/SMC1/SMC3) in chromnet-10 or Mediator (MED1/MED12/NIPBL) in
287 chromnet-11 (**Figure 5A**). These chromnets had high clustering coefficients and a
288 high proportion of CrP-CrP interactions, and their frequency in different chromatin

289 states was coherent with their known function. For example, chromnet-5 (Polycomb)
290 was strongly enriched in the two chromatin states enriched in H3K27me3.

291

292 We also noted the presence of star-like chromnets with very low clustering
293 coefficients. These star-like chromnets are mostly generated by emitter/signal and
294 signal/receiver interactions, suggesting that these are communication modules that
295 connect different protein complexes. For example, chromnet-3 contains two central
296 connectors (5fC and RYBP) connecting Polycomb, Mediator and TET1-SIN3A
297 complexes, and this chromnet is enriched in active transcription states and regulatory
298 elements.

299

300 Interestingly, chromnet-2 was a star-like module centered on 5hmC (the most central
301 hub in the network) and it contained all its positively co-localizing interactors: LSD1,
302 RYBP, ESRRB, KDM2A, TET1, OGT, G9A, and MBD2T (**Figure 5B**). In addition,
303 5hmC indirectly connects to H3K4me1 via TET1, and with 5mC via MBD2T. This
304 chromnet was clearly enriched in regulatory elements.

305

306 In summary, we have decomposed the communication network into communication
307 chromnets, functional modules of interactions with similar frequencies in the different
308 chromatin contexts. The components, structure and genomic distribution of these
309 chromnets provided information about their functional role. In particular, we detected
310 several star-like chromnets that are important to distribute epigenetic information to
311 different regions of the communication network. The wide range of functional
312 chromatin states that were enriched in these chromnets further supports their potential
313 role in mediating communication between distinct processes.

314

315 **Independent co-localization of 5hmC with ESRRB, LSD1, OGT and TET1 was** 316 **associated with different biological activities**

317

318 We have found that 5hmC is a very influential node for epigenetic communication
319 and the center of a star-like chromnet with similar enrichment associated with

320 chromatin states. We further characterized the genomic regions where 5hmC co-
321 localized independently with the stemness factor ESRRB and with the three
322 independent emitters of 5hmC, LSD1, OGT and TET1, which were also identified in
323 our co-evolutionary analysis (see above).

324

325 Remarkably, we found 6,307 genomic regions where 5hmC co-localized with its
326 receiver ESRRB in the absence of TET1, and with the rest of its interactors (**Figure**
327 **6A**). ESRRB is a transcription factor that is essential for the maintenance of ESCs
328 (Papp & Plath, 2012; Zwaka, 2012), yet to our knowledge the binding of ESRRB to
329 DNA has not been previously associated with the presence of 5hmC. However, the
330 ESRRB gene locus is known to be strongly enriched in 5hmC in ESCs (Doege *et al*,
331 2012), suggesting that 5hmC and ESRRB form a regulatory loop. **Gene ontology**
332 **analysis** carried out with the genes closest to these specific regions (McLean *et al*,
333 2010) identified stem cell maintenance, MAPK and Notch cell signaling cascades as
334 the most enriched functions (**Figure 6E**), highlighting the importance of ESRRB for
335 stemness maintenance. Surprisingly, the expression of the ESRRB gene is not ESC-
336 specific but rather it is expressed ubiquitously in most differentiated cell types
337 (Zwaka, 2012). Thus, its specific role in stemness probably requires ESC-specific
338 interactions with other components of the communication network and our results
339 suggested that 5hmC might be the key signal connecting ESRRB function with
340 stemness.

341

342 LSD1 is a H3K4- and H3K9-demethylase that can act as either a transcriptional co-
343 activator or co-repressor (Wang *et al*, 2007). To our knowledge, this was the first time
344 5hmC and LSD1 were found to coincide in the epigenome of ESCs (**Figure 6B**).
345 Interestingly, it is well known that there is a functional co-dependence between
346 histone demethylation and DNA methylation (Vaissière *et al*, 2008; Ikegami *et al*,
347 2009). Indeed, we consider LSD1 is an emitter of 5hmC because there is a global loss
348 of DNA methylation in the LSD1 knockout (Wang *et al*, 2009, 1). Remarkably, we
349 found that the 9,714 5hmC-LSD1 specific regions are significantly enriched with
350 specific terms associated with histone acetylation and DNA modification (**Figure**
351 **6E**), strengthening the dependent relationship between histone and DNA

352 modifications. Indeed, LSD1 not only functions as a histone demethylase by itself but
353 also, in association with 5hmC it can regulate the expression of proteins that modify
354 both histone acetylation and DNA methylation. These results suggest the presence of
355 a second regulatory loop involving 5hmC.

356

357 TET1 and OGT are two of the best known emitters of 5hmC (**Figure 6C-D**), with
358 TET1 a DNA demethylase that catalyzes the conversion of 5mC to 5hmC and OGT a
359 regulator of TET1 (Vella *et al*, 2013; Balasubramani & Rao, 2013). In fact, the role of
360 OGT in DNA demethylation was associated to its co-localization with TET1.
361 However, OGT is a N-acetylglucosaminyltransferase that can also bind to different
362 TFs independently of TET1 (Bond & Hanover, 2015). Notably, we observed different
363 functional enrichment of the 5hmC-TET1 and 5hmC-OGT regions (**Figure 6E**).
364 While the 27,721 5hmC-TET1 regions were enriched in stem cell maintenance and
365 morphogenesis, highlighting the role of both 5hmC and TET1 in stemness, the 1,017
366 5hmC-OGT regions were related with the metabolism of glycerophospholipids and
367 carbohydrates. Interestingly, OGT is known to bind phosphatidylinositol-3,4,5-
368 trisphosphate, regulating insulin responses and gluconeogenesis through glycosylation
369 of different proteins (Yang *et al*, 2008). Our results suggest that the alternative role of
370 OGT in gene regulation is also associated to 5hmC (but not to TET1). As the presence
371 of 5hmC requires the action of TET1, our results suggest that OGT might remain in
372 certain locations after TET1 removal, probably associated to the presence of specific
373 TFs in order to regulate the metabolism of glycerophospholipids and carbohydrates.
374 In this scenario, OGT would act as an emitter regulating 5hmC production and as a
375 receiver by acting with other proteins in the presence of 5hmC to regulate gene
376 expression.

377

378 In summary, the analysis of specific genomic regions revealed that different processes
379 and functions could be regulated and may be interconnected via 5hmC interactions
380 with other proteins. These processes include functions as relevant as epigenetic self-
381 regulation, cell signaling, maintaining stemness, morphogenesis and metabolism.

382

383 **Discussion**

384 ESCs constitute an ideal model to explore the epigenomic communication that
385 directly influences the phenotype of cells. Cytosine modifications, certain histone
386 marks and CrPs contribute to the plasticity required for the induction and maintenance
387 of pluripotency. Thus, the abundant epigenomic data from mouse ESCs has enabled
388 us to investigate how the different chromatin components communicate with each
389 other within a complex network. Using high-throughput genome-wide data and
390 information from the literature, we reconstructed the epigenetic communication
391 network of ESCs. In addition to the rigorously established co-occurrence and mutual
392 exclusion, we also annotated the directions of the CrP interactions mediated by
393 epigenetic signals (cytosine modifications and histone marks) based on information
394 extracted manually from the literature. This information allows CrPs to be classified
395 as emitters or receivers of these more basic epigenetic signals. This conceptual
396 framework constitutes the first explicit formulation to study chromatin as a biological
397 communication system.

398

399 We highlight the importance of using information taken from the literature. This
400 biological knowledge allowed us to understand the network of co-localization patterns
401 from high-throughput data, permitting us to obtain the first global picture of the
402 information flow that could take place in the ESC epigenome. Using an algorithm that
403 was originally proposed to evaluate the importance of internet web pages, we
404 identified the most influential nodes (those from which information spreads out) and
405 the most popular ones (those that collect information from many sources). Not
406 surprisingly, active RNA polymerase II was identified as the most popular node, as
407 many components of the epigenetic network regulate transcription. Our analysis
408 revealed that 5hmC is the most influential node in this network. In fact, 5hmC is a
409 signal received by eleven different CrPs, explaining its influential role for chromatin
410 communication in ESCs.

411

412 The elements that drive epigenetic communication constitute an intricate and dynamic
413 network that produces responses that range from stable programs defining cell-

414 identity to fast cellular responses. In this context, the fine-tuning of epigenetic
415 communication pathways is likely to have been a key aspect in the evolution of
416 multicellular organisms, such as metazoans. Co-evolutionary analyses point to
417 interactions that are conserved by evolutionarily coordinated changes. In fact, these
418 analyses can reveal strong functional links in the context of complex and dynamic
419 protein interactions (de Juan et al 2013). Co-evolution can occur between proteins that
420 interact directly or that participate in the same communication processes – for
421 example, via chromatin interactions mediated by histone marks or cytosine
422 modifications.

423

424 The majority of the co-evolutionary associations related to epigenetic communication
425 are triplets formed by an emitter, a signal and a receiver. Unexpectedly, four different
426 co-evolutionary associations were found between proteins interacting with 5hmC:
427 SIN3A with LSD1, TET1 with MBD2, MBD2 with MLL2, and OGT with TAF1.
428 Strikingly, all three co-occurring 5hmC emitters (TET1, OGT and LSD1) co-evolve
429 with three different 5hmC receivers, forming different emitter-5hmC-receiver triplets.
430 These associations do not reflect direct physical interactions of the protein pairs but
431 rather, complementary roles in the control of cytosine modifications and gene
432 regulation. Thus, we speculate that the balance between 5mC, 5hmC and other
433 cytosine modifications has been very important in fine-tuning epigenomic
434 communication during the evolution of metazoans.

435

436 Our results suggest that the alterations in the levels of cytosine modifications might be
437 driving important changes in the communication of chromatin components.
438 Interestingly, levels of 5hmC have been shown to be higher in stem cells and brain
439 compared to other mammalian tissues (Tahiliani *et al*, 2009; Kriaucionis & Heinz,
440 2009). Alterations in the levels and genomic location profiles of 5hmC have been
441 related to aging, neural diseases (CNSd, Song *et al*, 2011) and cancer (Pfeifer *et al*,
442 2013; Liyanage *et al*, 2014; Moen *et al*, 2015). It will be interesting to study the
443 networks of cancer and CNSd cells and evaluate the effect of 5hmC alterations in
444 their chromatin communication.

445

446 Identifying modules in networks helps to better understand their distinct components
447 (Mitra *et al*, 2013). Here, we followed a simple approach to identify functional sub-
448 networks of chromatin communication, or chromnets, clustering positive interactions
449 in function of their relative frequency in different chromatin states. This analysis
450 revealed the functional structure of the communication network and we were able to
451 automatically recover known protein-complexes, such as Polycomb and Mediator. By
452 contrast, we found that 5hmC and 5fC establish two different star-shaped chromnets,
453 suggesting that they might be involved in communication between distinct epigenetic
454 components and processes in distinct locations of the ESC epigenome.

455

456 While further experiments will be needed to reveal the functional roles of the different
457 independent interactions of 5hmC, our results generate some interesting hypotheses
458 about the possible independent functions played by 5hmC in ESCs. We propose that
459 the stem-specific role of ESRRB in ESCs could be linked to its co-occurrence with
460 5hmC, as this cytosine modification is less common in most differentiated cell types
461 (Zwaka, 2012). Our results also show that LSD1-5hmC might be specifically involved
462 in the regulation of histone modifications and DNA methylation, while the TET1-
463 5hmC interaction is associated with stem cell maintenance and morphology.
464 Furthermore, our data suggest a TET1-independent interaction between 5hmC and
465 OGT that might participate in the regulation of energy metabolism, and an interaction
466 between 5hmC and LSD1 regulates histones and DNA methylation.

467

468 The combination of genome-wide location data, prior knowledge from the literature
469 and protein co-evolution highlights conserved functional relationships between
470 5hmC-interacting CrPs that have been dynamically coordinated during evolution.
471 Based on our co-evolution analysis, we hypothesize that the different cytosine
472 modifications in different regions of the genome might have been important during
473 metazoan evolution. Our results suggest that the interaction of 5hmC with specific
474 emitters is involved in regulating different specific and critical functions.

475

476 In conclusion, network architecture conveys relevant contextual information that
477 cannot be easily obtained from analyses that focus on only a few epigenetic features.

478 The computational framework introduced here represents the basis to explore this vast
479 space and it provides the first integrated picture of the different elements involved in
480 epigenetic regulation. Accordingly, this analysis enables us to attain an integrated
481 vision of epigenetic communication in ESCs that highlighted the relevance of 5hmC
482 as a central signal.
483

484 **Experimental Procedures**

485 **ChIP-Seq, MeDIP and GLIB data processing**

486 We retrieved data for 139 Chromatin Immunoprecipitation Sequencing (ChIP-Seq),
487 Methylated DNA immunoprecipitation (MeDIP) and GLIB (glucosylation, periodate
488 oxidation and biotinylation) experiments described in **Table S1**. The sra files were
489 transformed into fastq files with the sra-toolkit (v2.1.12) and aligned to the reference
490 mm9/NCBI37 genome with bwa v0.5.9-r16 (Li & Durbin, 2009) allowing 0-1
491 mismatches. Unique reads were converted to BED format.

492

493 **Genome segmentation**

494 The input information used to segment the genome into different chromatin states was
495 that derived from the 3 cytosine modifications, the 13 histone marks and the insulator
496 protein CTCF - which has been previously shown to define a particular chromatin
497 state *per se* (Ernst & Kellis, 2010). We used the ChromHmm software (Ernst &
498 Kellis, 2012: v1.03) to define a 20 chromatin states model consistent with prior
499 knowledge regarding the function of these features (**Figure S3**). Only, intervals with a
500 probability higher than 0.95 were considered for further analysis.

501

502 **Co-location network inference**

503 We used the ChromHMM segments with a probability higher than 0.95 as samples for
504 the network inference. For a description of reads and samples filtering see Extended
505 Experimental Procedures. We applied the method described in (Perner *et al*, 2014)
506 that aims to unravel the direct interactions between factors that cannot be “explained”
507 by the other observed factors and thus, this is a more specific approach than an
508 analysis of simple pairwise correlations. Consequently, the more complete the number

509 of factors included in the analysis, the higher the certainty that inferred direct
510 interactions correspond to actual co-dependences. We inferred an interaction network
511 for each chromHMM state. Briefly, an Elastic Net was trained in a 10-fold Cross-
512 validation to predict the HM/CTCF/DNA methylation of the CrPs or to predict each
513 CrP from all other CrPs. Furthermore, the sparse partial correlation network (SPCN)
514 was obtained using all the samples available. We selected the interactions between
515 Histone marks/cytosine modifications and CrPs that obtained a high coefficient ($w \geq$
516 $2 * sd(all_w)$) in the Elastic Net prediction and that have a non-zero partial correlation
517 coefficient in the SPCN.

518

519 We counted the overlapping ChIP-Seq reads for the genomic segments using
520 Rsamtools. Using hierarchical clustering with $1 - cor(x,y)$ as a distance measure, we
521 find that most replicates or functionally related samples fall into the same branch
522 (**Figure S4A**). Given this consistency, we selected one experiment for those features
523 that are available from more than one dataset. To further test the robustness of this
524 choice, we generated 10 alternative networks by randomly selecting other replicates.
525 Our results show that the retrieved network is very robust to replicate selection
526 (**Figure S4B and S4C**).

527

528 **Influence/popularity analysis of the co-location network**

529 The popularity of a node in the chromatin network coincides with the standard
530 PageRank centrality score (Brin & Page, 1998) as computed from the (asymmetric)
531 adjacency matrix of the epigenetic communication network (see above). The
532 influence of a node has been computed as its PageRank score after inverting the
533 directions of the edges in the original network (influence-PageRank, Chepelianskii,
534 2010). For a detailed description of this analysis and the evaluations of its robustness
535 see Extended Experimental Procedures and **Figure S5**.

536

537 **Co-evolutionary network inference**

538 We retrieved protein trees of sequences at the metazoan level from eggNOG v4.0
539 (Powell *et al*, 2014). We removed tree inconsistencies using a previous pipeline (Juan
540 *et al*, 2013) and extracted only-unique-orthologous protein trees for each mouse

541 protein. The inter-orthologs evolutionary distances for the 58 mouse proteins with
542 ChIP-seq data analyzed in this study were mapped to distance bins and organized in a
543 data matrix. From this data we inferred the parameters (Besag, 1977; Aurell &
544 Ekeberg, 2012) of a pairwise model of interacting proteins in the space of species-
545 species evolutionary distances. Co-evolutionary scores were finally computed from
546 these parameters. For a detailed description of obtaining inter-orthologs evolutionary
547 distances and network inference see Extended Experimental Procedures.

548

549 **Identification of epigenetic signals with a statistically significant co-evolutionary** 550 **effect**

551 For each epigenetic signal (histone mark/cytosine modification), we identified all the
552 pairs of CrPs that satisfy the following two conditions: 1) the proteins in the pair are
553 co-evolutionary coupled (see above); and 2) each of the proteins in the pair directly
554 interacts with the epigenetic signal. We then used the number of unique CrPs in the
555 resulting set of pairs (Co-evolutionary Filtered Centrality, CFC) as a measure of the
556 influence of the signal on co-evolution between the CrPs in the epigenetic signaling
557 network. The statistical significance of each CFC was evaluated by computing a p-
558 value that corresponded to the probability of obtaining a CFC greater or equal to that
559 observed in a network model with randomly-generated edges among the CrPs in the
560 co-evolutionary analysis. This procedure identified three signals with a significant
561 CFC (p-value < 0.05): 5hmC (p-value approx. 0.04), H3K4me2 (0.01), H3K4me3
562 (0.02).

563

564 **Functional Modularization of the Co-localization Network**

565 The co-localization network was decomposed into local networks of positive
566 interactions. First, we calculated the frequency of each positive interaction for every
567 chromatin state using ChromHMM peaks, considering that an interaction is present if
568 both interactors are 'present' in the same 200 bp genomic window. The frequencies of
569 the interactions were standardized separately for every state. These vectors were
570 clustered by hierarchical clustering (Pearson correlation, average linkage) and the
571 largest statistically supported clusters (p-value < 0.05, n=10,000) according to Pvcust
572 (Suzuki & Shimodaira, 2006) were defined as chromnets (**Figure S6**).

573

574 **Gene Ontology enrichment analysis**

575 Gene Ontology enrichment analyses were carried out with GREAT v3.0.0 (McLean *et*
576 *al*, 2010). The genomic regions were associated to genes with a minimum distance of
577 5Kb upstream and 1Kb downstream, with the whole genome as the background. The
578 False Discovery Rate (FDR) considered was 0.05 (**Table S4**).

579

580 **URLs**

581 UCSC Trackhub with chromatin states, cytosine modifications, histone marks and
582 CrPs

583 <http://genome.ucsc.edu/cgi->

584 bin/hgTracks?db=mm9&hubUrl=http://epistemnet.bioinfo.cnio.es/mESC_CNIO_hub2/hub.txt

585

586 EpiStemNet web interface: visualization of co-location networks in ESCs

587 <http://dogcaesar.github.io/epistemnet>

588

589 Scripts available at <https://github.com/EpiStemNet>

590

591 Processed data available at <http://epistemnet.bioinfo.cnio.es>

592

593 **Acknowledgements**

594 We thank Prof. Stunnenberg and all members of BLUEPRINT Consortium for
595 critical comments.

596

597 **Author contributions**

598 DJ, EC, JP, DR and AV were responsible for conception and design. DJ, EC, JP and
599 SM analyzed the data and DR guided all the analyses. EC processed the ChIP-seq
600 samples. JP performed the network reconstruction. DJ, JP, SM, EC and DR
601 analyzed the network. DO developed the network visualization web tool. DJ and

602 SM performed the co-evolutionary analyses. DJ, EC, JP and DR compiled the
603 literature knowledge. MV, HRC, DR and AV supervised research. All authors were
604 responsible for interpretation of the data. DJ, EC, JP, SM and DR wrote the initial
605 draft. All authors read and approved the final manuscript.
606

607 **Conflict of interest**

608 The authors declare that they have no conflict of interest.

609 **Funding**

610 The research leading to these results has received funding from the European
611 Union's Seventh Framework Programme (FP7/2007-2013) under grant
612 agreement number 282510 (BLUEPRINT).
613

614 **References**

- 615 Aurell, E & Ekeberg, M (2012). Inverse Ising inference using all the data. *Physical*
616 *review letters*, **108**: 090201
- 617 Balasubramani A & Rao A (2013) O-GlcNAcylation and 5-Methylcytosine
618 Oxidation: An Unexpected Association between OGT and TETs. *Mol. Cell* **49**:
619 618–619
- 620 Baubec T, Ivánek R, Lienert F & Schübeler D (2013) Methylation-Dependent and
621 -Independent Genomic Targeting Principles of the MBD Protein Family. *Cell*
622 **153**: 480–492
- 623 Brin, S., & Page, L. (1998). The anatomy of a large-scale hypertextual web search
624 engine. *Comput. Networks ISDN* **30**: 107–117
- 625 Van Bommel JG, Filion GJ, Rosado A, Talhout W, de Haas M, van Welsem T, van
626 Leeuwen F & van Steensel B (2013) A network model of the molecular
627 organization of chromatin in *Drosophila*. *Mol. Cell* **49**: 759–771
- 628 Besag J (1977) Efficiency of pseudo-likelihood estimation for simple gaussian
629 fields. *Biometrika* **64**: 616–618
- 630 Bond MR & Hanover JA (2015) A little sugar goes a long way: the cell biology of
631 O-GlcNAc. *J. Cell Biol.* **208**: 869–880
- 632 Chepelianskii, A. D. (2010). Towards physical laws for software architecture.
633 arXiv preprint arXiv:1003.5455.
- 634 Choi I, Kim R, Lim H-W, Kaestner KH & Won K-J (2014) 5-hydroxymethylcytosine
635 represses the activity of enhancers in embryonic stem cells: a new epigenetic

- signature for gene regulation. *BMC Genomics* **15**: 670
- 637 Doege CA, Inoue K, Yamashita T, Rhee DB, Travis S, Fujita R, Guarnieri P, Bhagat
638 G, Vanti WB, Shih A, Levine RL, Nik S, Chen EI & Abeliovich A (2012) Early-
639 stage epigenetic modification during somatic cell reprogramming by Parp1
640 and Tet2. *Nature* **488**: 652–655
- 641 Ernst J & Kellis M (2010) Discovery and characterization of chromatin states for
642 systematic annotation of the human genome. *Nat. Biotechnol.* **28**: 817–825
- 643 Ernst J & Kellis M (2012) ChromHMM: automating chromatin-state discovery
644 and characterization. *Nat. Methods* **9**: 215–216
- 645 Ficz G, Branco MR, Seisenberger S, Santos F, Krueger F, Hore TA, Marques CJ,
646 Andrews S & Reik W (2011) Dynamic regulation of 5-hydroxymethylcytosine
647 in mouse ES cells and during differentiation. *Nature* **473**: 398–402
- 648 Filion GJ, van Bemmell JG, Braunschweig U, Talhout W, Kind J, Ward LD, Brugman
649 W, de Castro IJ, Kerkhoven RM, Bussemaker HJ & van Steensel B (2010)
650 Systematic protein location mapping reveals five principal chromatin types
651 in Drosophila cells. *Cell* **143**: 212–224
- 652 Hashimoto H, Liu Y, Upadhyay AK, Chang Y, Howerton SB, Vertino PM, Zhang X &
653 Cheng X (2012) Recognition and potential mechanisms for replication and
654 erasure of cytosine hydroxymethylation. *Nucleic Acids Res.* **40**: 4841–4849
- 655 Hemmrich G, Khalturin K, Boehm A-M, Puchert M, Anton-Erxleben F, Wittlieb J,
656 Klostermeier UC, Rosenstiel P, Oberg H-H, Domazet-Lošo T, Sugimoto T,
657 Niwa H & Bosch TCG (2012) Molecular Signatures of the Three Stem Cell
658 Lineages in Hydra and the Emergence of Stem Cell Function at the Base of
659 Multicellularity. *Mol. Biol. Evol.* **29**: 3267–3280
- 660 He Y-F, Li B-Z, Li Z, Liu P, Wang Y, Tang Q, Ding J, Jia Y, Chen Z, Li L, Sun Y, Li X,
661 Dai Q, Song C-X, Zhang K, He C & Xu G-L (2011) Tet-Mediated Formation of 5-
662 Carboxylcytosine and Its Excision by TDG in Mammalian DNA. *Science* **333**:
663 1303–1307
- 664 Ikegami K, Ohgane J, Tanaka S, Yagi S & Shiota K (2009) Interplay between DNA
665 methylation, histone modification and chromatin remodeling in stem cells
666 and during development. *Int. J. Dev. Biol.* **53**: 203–214
- 667 Ito S, Shen L, Dai Q, Wu SC, Collins LB, Swenberg JA, He C & Zhang Y (2011) Tet
668 Proteins Can Convert 5-Methylcytosine to 5-Formylcytosine and 5-
669 Carboxylcytosine. *Science* **333**: 1300–1303
- 670 de Juan D, Pazos F & Valencia A (2013) Emerging methods in protein co-
671 evolution. *Nat. Rev. Genet.* **14**: 249–261
- 672 Juan D, Rico D, Marques-Bonet T, Fernández-Capetillo O & Valencia A (2013)
673 Late-replicating CNVs as a source of new genes. *Biol. Open* **2**: 1402–1411
- 674 Kriaucionis S, Heintz N (2009) The nuclear DNA base 5-hydroxymethylcytosine
675 is present in Purkinje neurons and the brain. *Science*. **324**:929-930
- 676 Li H & Durbin R (2009) Fast and accurate short read alignment with Burrows-
677 Wheeler transform. *Bioinformatics* **25**: 1754–1760
- 678 Liyanage VRB, Jarmasz JS, Murugesan N, Del Bigio MR, Rastegar M & Davie JR
679 (2014) DNA Modifications: Function and Applications in Normal and Disease
680 States. *Biology* **3**: 670–723
- 681 McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, Wenger AM &
682 Bejerano G (2010) GREAT improves functional interpretation of cis-

- 683 regulatory regions. *Nat. Biotechnol.* **28**: 495–501
- 684 Mitra K, Carvunis A-R, Ramesh SK & Ideker T (2013) Integrative approaches for
685 finding modular structure in biological networks. *Nat. Rev. Genet.* **14**: 719–
686 732
- 687 Moen EL, Mariani CJ, Zullo H, Jeff-Eke M, Litwin E, Nikitas JN & Godley LA
688 (2015) New themes in the biological functions of 5-methylcytosine and 5-
689 hydroxymethylcytosine. *Immunol. Rev.* **263**: 36–49
- 690 Papp B & Plath K (2012) Pluripotency re-centered around Esrrb. *EMBO J.* **31**:
691 4255–4257
- 692 Pastor WA, Pape UJ, Huang Y, Henderson HR, Lister R, Ko M, McLoughlin EM,
693 Brudno Y, Mahapatra S, Kapranov P, Tahiliani M, Daley GQ, Liu XS, Ecker JR,
694 Milos PM, Agarwal S & Rao A (2011) Genome-wide mapping of 5-
695 hydroxymethylcytosine in embryonic stem cells. *Nature* **473**: 394–397
- 696 Perner J, Lasserre J, Kinkley S, Vingron M & Chung H-R (2014) Inference of
697 interactions between chromatin modifiers and histone modifications: from
698 ChIP-Seq data to chromatin-signaling. *Nucleic Acids Res.* **42**: 13689–13695
- 699 Pfeifer GP, Kadam S & Jin S-G (2013) 5-hydroxymethylcytosine and its potential
700 roles in development and cancer. *Epigenetics Chromatin* **6**: 10
- 701 Powell S, Forslund K, Szklarczyk D, Trachana K, Roth A, Huerta-Cepas J, Gabaldón
702 T, Rattei T, Creevey C, Kuhn M, Jensen LJ, von Mering C & Bork P (2014)
703 eggNOG v4.0: nested orthology inference across 3686 organisms. *Nucleic
704 Acids Res.* **42**: D231–239
- 705 Raiber E-A, Beraldi D, Ficuz G, Burgess HE, Branco MR, Murat P, Oxley D, Booth MJ,
706 Reik W & Balasubramanian S (2012) Genome-wide distribution of 5-
707 formylcytosine in embryonic stem cells is associated with transcription and
708 depends on thymine DNA glycosylase. *Genome Biol.* **13**: R69
- 709 Schreiber SL & Bernstein BE (2002) Signaling Network Model of Chromatin. *Cell*
710 **111**: 771–778
- 711 Scott-Phillips TC (2008) Defining biological communication. *J. Evol. Biol.* **21**:
712 387–395
- 713 Smith T late JM & Harper D (2003) Animal Signals Oxford Series in Ecology and
714 Evolution
- 715 Song C-X, Szulwach KE, Fu Y, Dai Q, Yi C, Li X, Li Y, Chen C-H, Zhang W, Jian X,
716 Wang J, Zhang L, Looney TJ, Zhang B, Godley LA, Hicks LM, Lahn BT, Jin P &
717 He C (2011) Selective chemical labeling reveals the genome-wide
718 distribution of 5-hydroxymethylcytosine. *Nat. Biotechnol.* **29**: 68–72
- 719 Stroud H, Feng S, Morey Kinney S, Pradhan S & Jacobsen SE (2011) 5-
720 Hydroxymethylcytosine is associated with enhancers and gene bodies in
721 human embryonic stem cells. *Genome Biol.* **12**: R54
- 722 Suzuki R & Shimodaira H (2006) Pvcust: an R package for assessing the
723 uncertainty in hierarchical clustering. *Bioinforma. Oxf. Engl.* **22**: 1540–1542
- 724 Szulwach KE, Li X, Li Y, Song C-X, Han JW, Kim S, Namburi S, Hermetz K, Kim JJ,
725 Rudd MK, Yoon Y-S, Ren B, He C & Jin P (2011) Integrating 5-
726 Hydroxymethylcytosine into the Epigenomic Landscape of Human
727 Embryonic Stem Cells. *PLoS Genet* **7**: e1002154
- 728 Tahiliani M, Koh KP, Shen Y, Pastor WA, Bandukwala H, Brudno Y, Agarwal S,
729 Iyer LM, Liu DR, Aravind L, Rao A (2009) Conversion of 5-methylcytosine to

- 730 5-hydroxymethylcytosine in mammalian DNA by MLL partner TET1. *Science*.
731 **324**:930-935
- 732 Vaissière T, Sawan C & Herceg Z (2008) Epigenetic interplay between histone
733 modifications and DNA methylation in gene silencing. *Mutat. Res.* **659**: 40–48
- 734 Vella P, Scelfo A, Jammula S, Chiacchiera F, Williams K, Cuomo A, Roberto A,
735 Christensen J, Bonaldi T, Helin K & Pasini D (2013) Tet Proteins Connect the
736 O-Linked N-acetylglucosamine Transferase Ogt to Chromatin in Embryonic
737 Stem Cells. *Mol. Cell* **49**: 645–656
- 738 Wang J, Hevi S, Kurash JK, Lei H, Gay F, Bajko J, Su H, Sun W, Chang H, Xu G,
739 Gaudet F, Li E & Chen T (2009) The lysine demethylase LSD1 (KDM1) is
740 required for maintenance of global DNA methylation. *Nat. Genet.* **41**: 125–
741 129
- 742 Wang J, Scully K, Zhu X, Cai L, Zhang J, Prefontaine GG, Kronen A, Ohgi KA, Zhu P,
743 Garcia-Bassets I, Liu F, Taylor H, Lozach J, Jayes FL, Korach KS, Glass CK, Fu
744 X-D & Rosenfeld MG (2007) Opposing LSD1 complexes function in
745 developmental gene activation and repression programmes. *Nature* **446**:
746 882–887
- 747 Williams K, Christensen J, Pedersen MT, Johansen JV, Cloos PAC, Rappsilber J &
748 Helin K (2011) TET1 and hydroxymethylcytosine in transcription and DNA
749 methylation fidelity. *Nature* **473**: 343–348
- 750 Wu H, D'Alessio AC, Ito S, Wang Z, Cui K, Zhao K, Sun YE & Zhang Y (2011)
751 Genome-wide analysis of 5-hydroxymethylcytosine distribution reveals its
752 dual function in transcriptional regulation in mouse embryonic stem cells.
753 *Genes Dev.* **25**: 679–684
- 754 Yang X, Ongusaha PP, Miles PD, Havstad JC, Zhang F, So WV, Kudlow JE, Michell
755 RH, Olefsky JM, Field SJ & Evans RM (2008) Phosphoinositide signalling links
756 O-GlcNAc transferase to insulin resistance. *Nature* **451**: 964–969
- 757 Ye T, Krebs AR, Choukrallah M-A, Keime C, Plewniak F, Davidson I & Tora L
758 (2011) seqMINER: an integrated ChIP-seq data interpretation platform.
759 *Nucleic Acids Res.* **39**: e35
- 760 Zwaka TP (2012) Pluripotency network in embryonic stem cells: maybe Leibniz
761 was right all along. *Cell Stem Cell* **11**: 441–442
762

763

764

765 **Figure legends**

766

767 **Figure 1. A framework to study communication among chromatin components**

768 Our network approach is based on a classification of the epigenomic features (see
769 **Table S1**) into three component classes, where histone and cytosine modifications
770 are always considered to be signals and the chromatin-related proteins (CrPs) can
771 be either co-occurring (or mutually exclusive) emitters (writers/erasers) or
772 receivers (readers) of those epigenetic signals.

773

774 **Figure 2. Chromatin communication network in ESCs**

775 Full chromatin communication network in which the edges represent positive or
776 negative interactions that indicate genomic co-localization or mutual exclusion,
777 respectively. Arrows associated with the directional edges represent
778 communication flux for emitter-signal or signal-receiver pairs retrieved from the
779 literature (see **Table S2**). The colors indicate membership of known protein
780 complexes.

781 See also **Figures S1, S2, S3 and S4**.

782

783 **Figure 3. Influence and popularity of chromatin communication nodes**

784 A: Emitters and receivers of the H3K79me2 signal. B: Emitters and receivers of
785 the 5hmC signal. C: Influence vs. popularity plot for the chromatin
786 communication network. The popularity and the influence of each node
787 correspond respectively to the PageRank and the influence-PageRank values (see
788 Experimental Procedures). The size of the nodes increases linearly with the scores
789 sum, while the color reflects the scores difference. D: Influence vs. popularity plot
790 for the perturbed chromatin network after removal of 10% of the original edges.
791 The values in the plot are averaged over 2000 networks with randomly removed
792 edges. The vertical and horizontal bars show respectively the standard deviation
793 of the influence of 5hmC and of the popularity of RNAPII_S2P in the ensemble of
794 perturbed networks.

795 See also **Figure S5**.

796

797

798 **Figure 4. Co-evolution of CrPs**

799 Coupling analysis of the phylogenetic histories of CrPs revealed significant co-
800 evolution between emitters and receivers of 5hmC and 5fC. Co-evolving pairs are
801 indicated by thick colored dashed lines. The grey lines indicate co-localization or
802 mutual exclusion in the chromatin communication network (see **Figure 2** for
803 more details).

804 See also **Table S3**.

805 **Figure 5. Chromnets recover known protein complexes and star-shaped**
806 **structures**

807 The chromnets are sub-networks of interactions with similar co-occurrence across
808 different chromatin states. Each bar plot indicates the overall enrichment of the
809 chromatin states in each chromnet along (see B for details of the chromatin
810 states). B Star-like 5hmC sub-network and the overall enrichment of chromatin
811 states.

812 See also **Figure S6**

813

814 **Figure 6. 5hmC genomic regions have different functional enrichment depending**
815 **on the co-localizing partner**

816 A-D Read densities over a 10Kb windows centered on the 5hmC-ESRRB (A),
817 5hmC-LSD1 (B), 5hmC-TET1 (C) and 5hmC-OGT (D) peaks. We calculated the
818 read density of 5hmC, ESRRB, LSD1, TET1 (N- and C-terminal ChIP-seqs) and
819 OGT in 10Kb windows centered on the genomic bins (200 bp), where 5hmC co-
820 localizes exclusively with each specific partner (i.e.: the rest of the 5hmC
821 interactors are not present). The read density plots were obtained with the
822 SeqMINER platform v1.3.3e (Ye *et al*, 2011). The average density of the reads in
823 50 bp bins was plotted from the center of the 5hmC independent genomic regions
824 to +/-5000 bp. E Gene Ontology enrichment analysis of peaks in A-D using
825 GREAT (see **Table S4**).

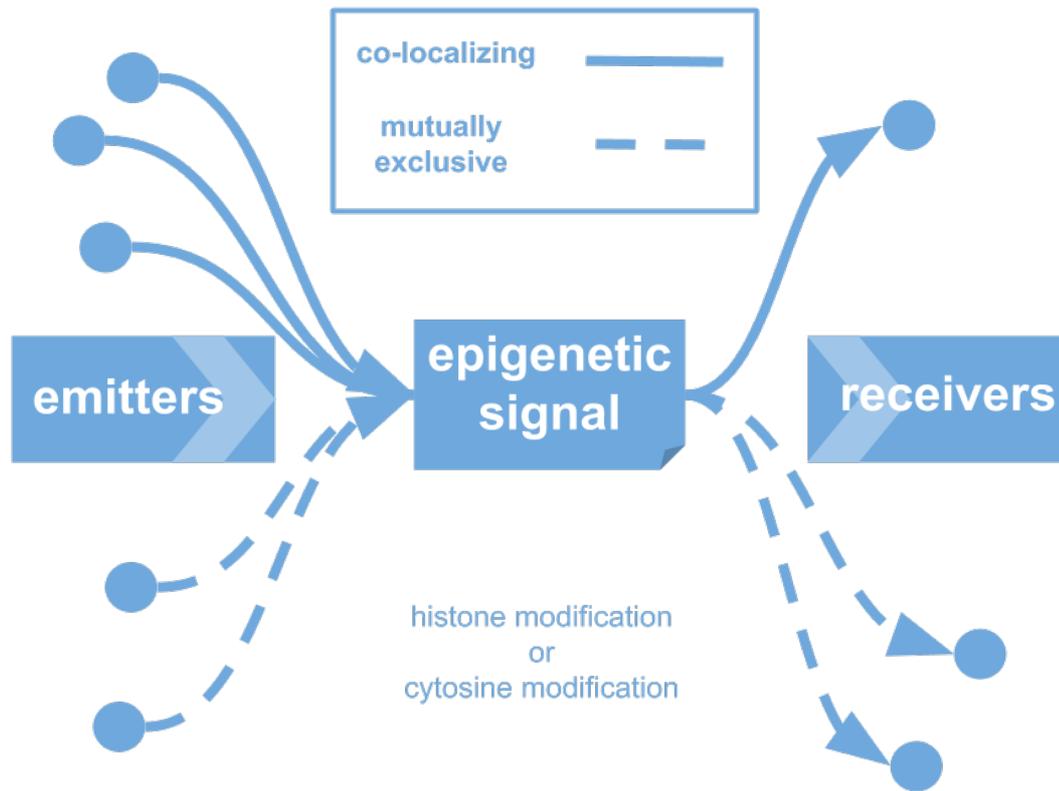


Figure 1

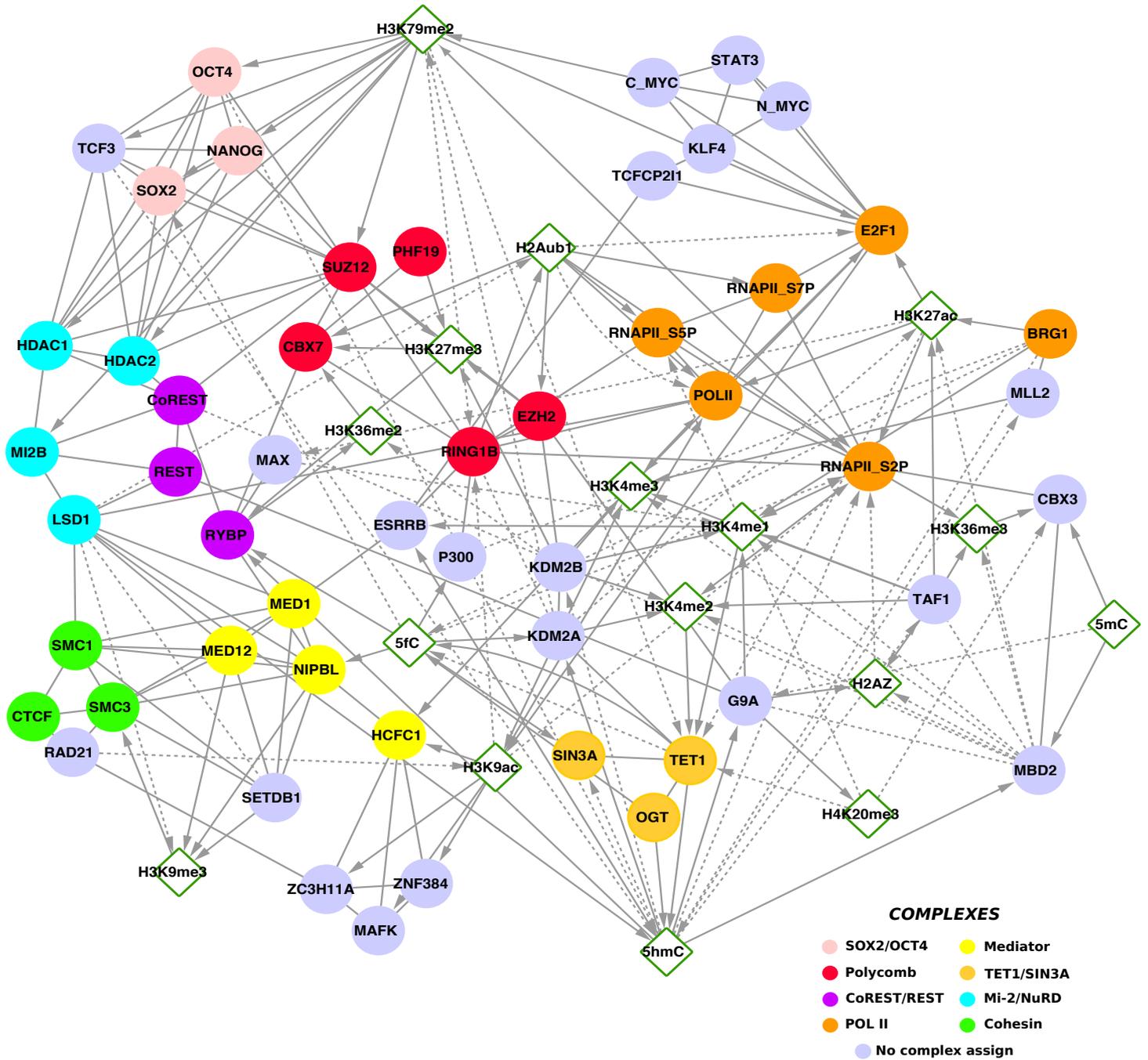
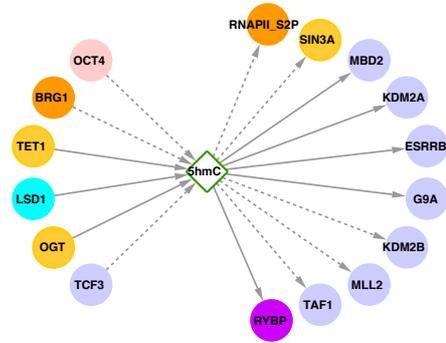
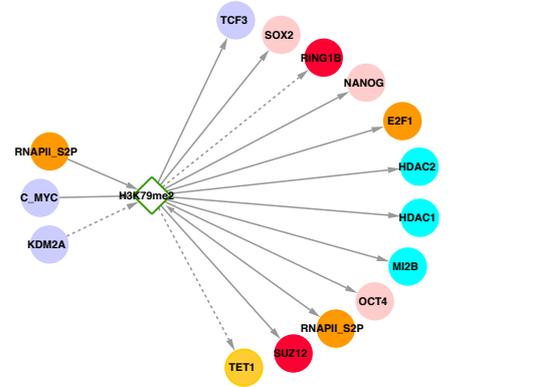


Figure 2

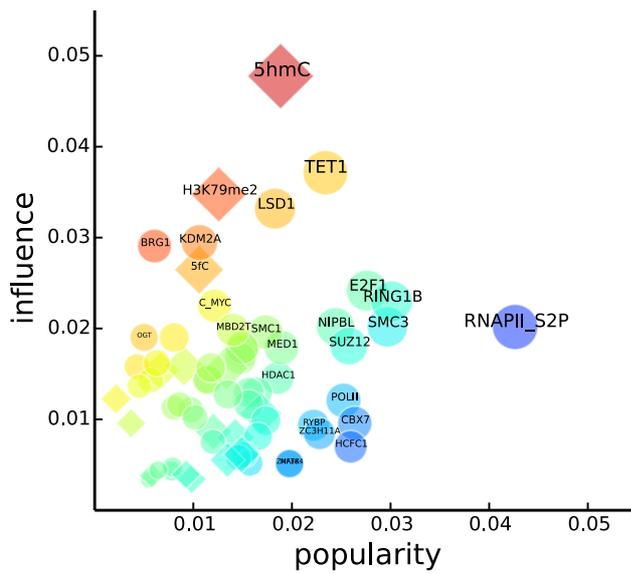
A



B



C



D

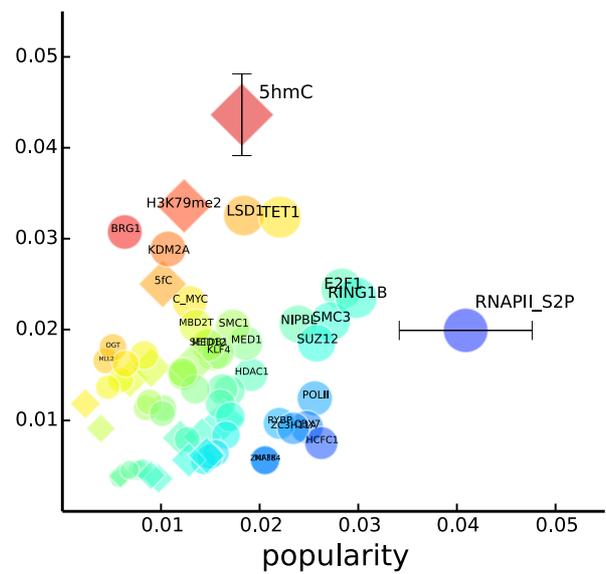


Figure 3

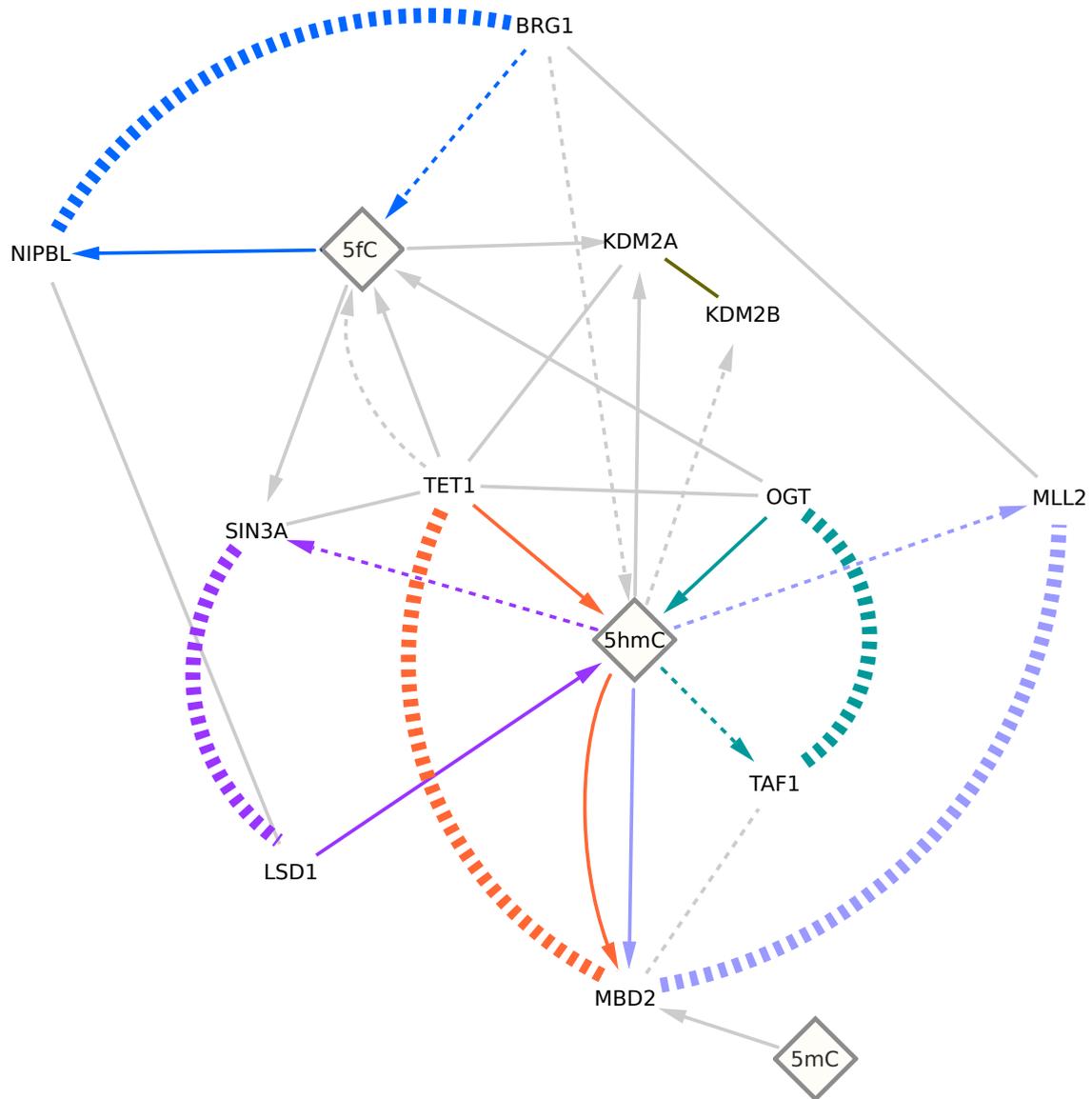


Figure 4

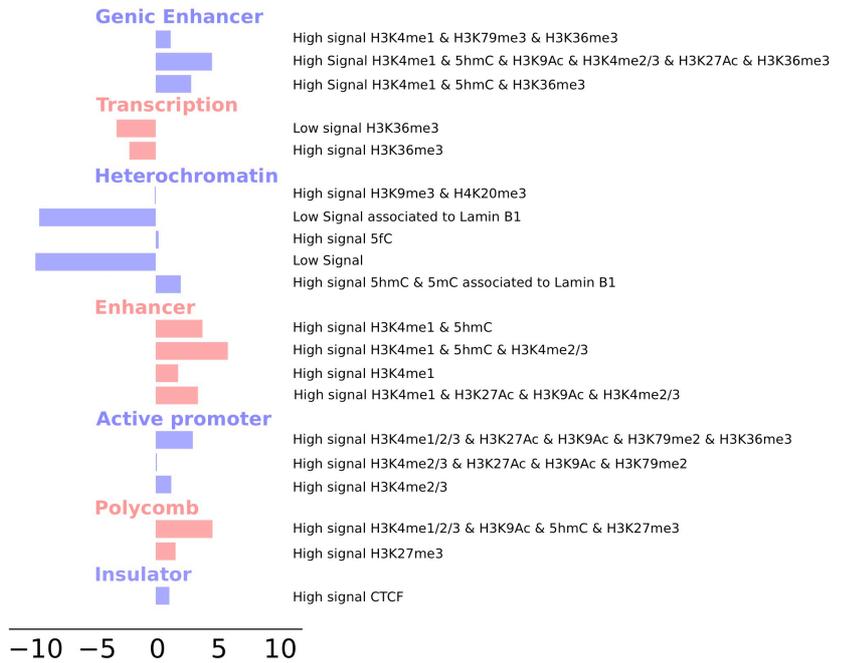
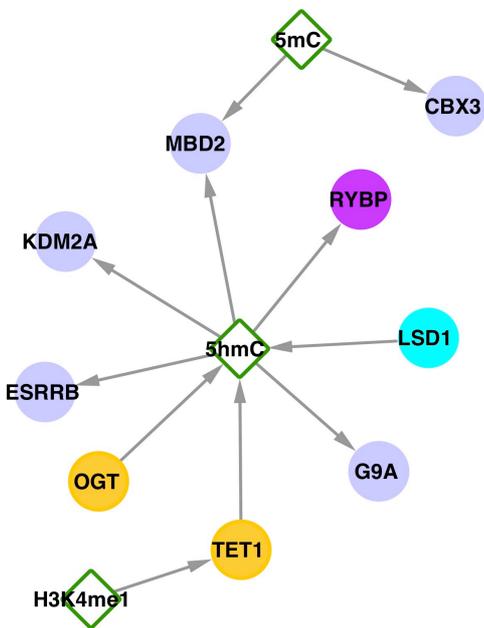
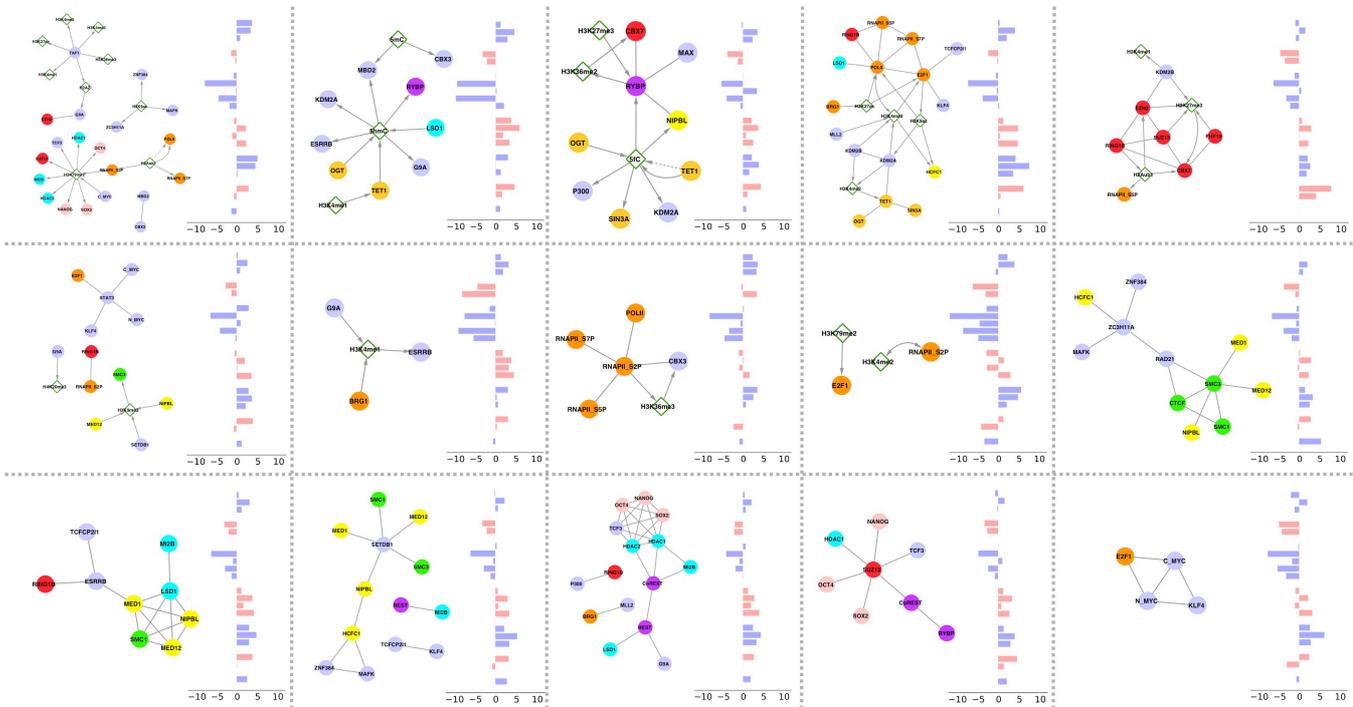
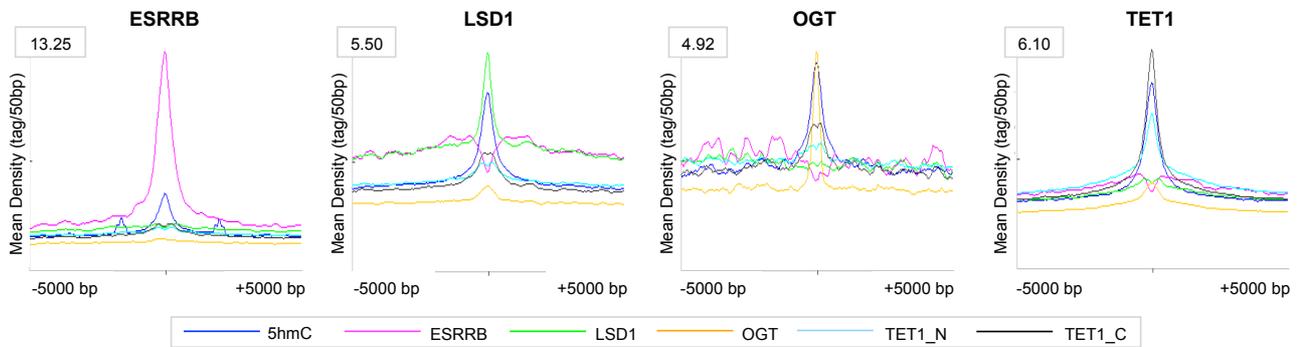


Figure 5



GO term	ESRRB	LSD1	TET1	OGT
Maturation of SSU-rRNA	■			
Regative regulation of MAPK cascade	■			
Plasma membrane organization	■			
Intrinsic apoptotic signaling pathway	■			
Cell-cell junction organization	■			
Epithelial cell morphogenesis	■			
Somatic stem cell maintenance	■		■	
Regulation of Notch signaling pathway	■		■	
Mesoderm morphogenesis	■			
Histone acetylation	■	■		
Regulation of gene expression, epigenetic	■	■		
DNA methylation or demethylation		■		
Histone H3-H4 acetylation		■		
Placenta development		■		
Blastocyst formation		■		
Proteasomal protein catabolic process		■		
Regulation of translation		■		
Regulation of stem cell maintenance			■	
Regulation of TGF-beta receptor signaling			■	
Apoptotic process involved in morphogenesis			■	
Cell-cell signaling involved in cell fate commitment			■	
Developmental induction			■	
Mesoderm development			■	
Body and head morphogenesis			■	■
Carbohydrate biosynthetic process				■
Glycerophospholipid metabolic process				■
Regulation of lipid storage				■
Maintenance of location in cell				■
Modulation by virus of host morphology or physiology				■
Regulation of organelle organization				■
Regulation of chromatin organization				■
Regulation of histone modification				■
Regulation of cytokinesis				■

Figure 6

Supplemental Information

Epigenomic co-location and co-evolution reveal a key role for 5hmC as a communication hub in the chromatin network of ESCs

Table of contents

Supplemental Figures

Figure S1. Related to Figure 2. Epigenetic communication network in- and out-degrees of nodes

Figure S2. Related to Figure 2. Epigenetic communication network betweenness

Figure S3. Related to Figure 2. Chromatin States definitions, emission probabilities and genomic annotation enrichments.

Figure S4. Related to Figure 2. Effect of sample selection.

Figure S5. Related to Figure 3. Inferred signal-receiver directions.

Figure S6. Related to Figure 5. Co-localization pairs clustered according to their distribution in the chromatin states

Supplemental Tables

Table S1. Related to Figure 1. ChIP-Seq experiments and annotations included in the study (In Table_S1.xls)

Table S2. Related to Figure 2. Co-localization couplings with emissor or receptor directions (In Table_S2.xls)

Table S3. Related to Figure 4. Co-evolutionary couplings (In Table_S3.xls)

Table S4. Related to Figure 6. Gene Ontology enrichment analysis in 5hmC independent genomic segments (In Table_S4.xls)

Supplemental Experimental Procedures

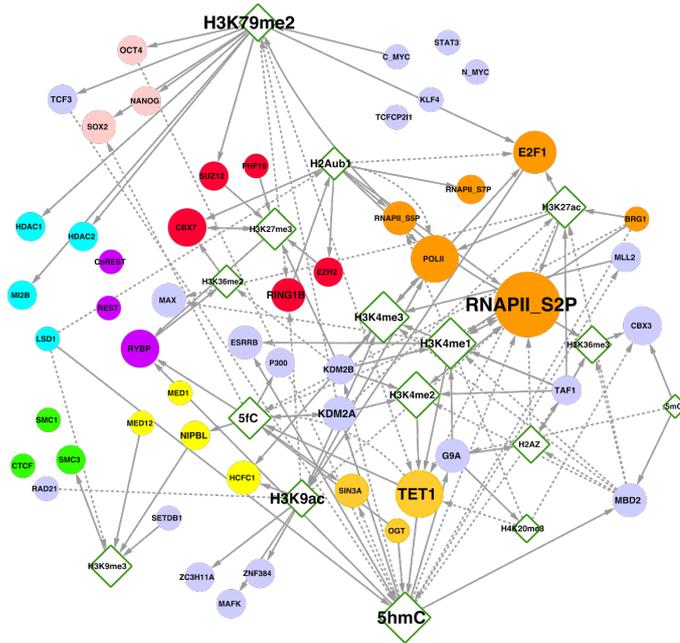
Supplemental References

Supplemental Figures

Figure S1. Related to Figure 2. Epigenetic communication network in- and out-degrees of nodes

Epigenetic communication network, as in Figure 2 of the main manuscript, showing only directional interactions mediated by epigenetic signals. Node size is proportional to its A) in-degree (number of incoming edges) and B) out-degree (number of outgoing edges).

A



B

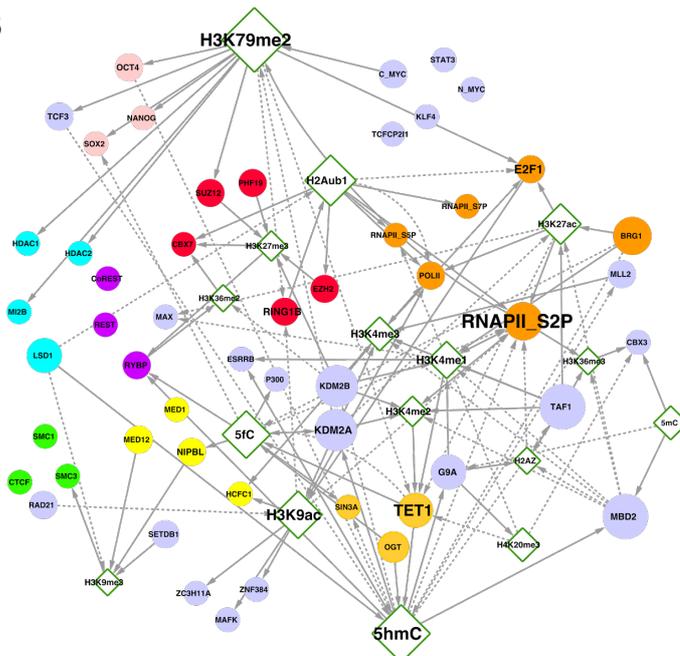


Figure S2. Related to Figure 2. Epigenetic communication network betweenness

Epigenetic communication network, as in Figure 2 of the main manuscript, showing only directional interactions mediated by epigenetic signals. Node size is proportional to its betweenness (number of shortest paths in the network that are mediated by a node).

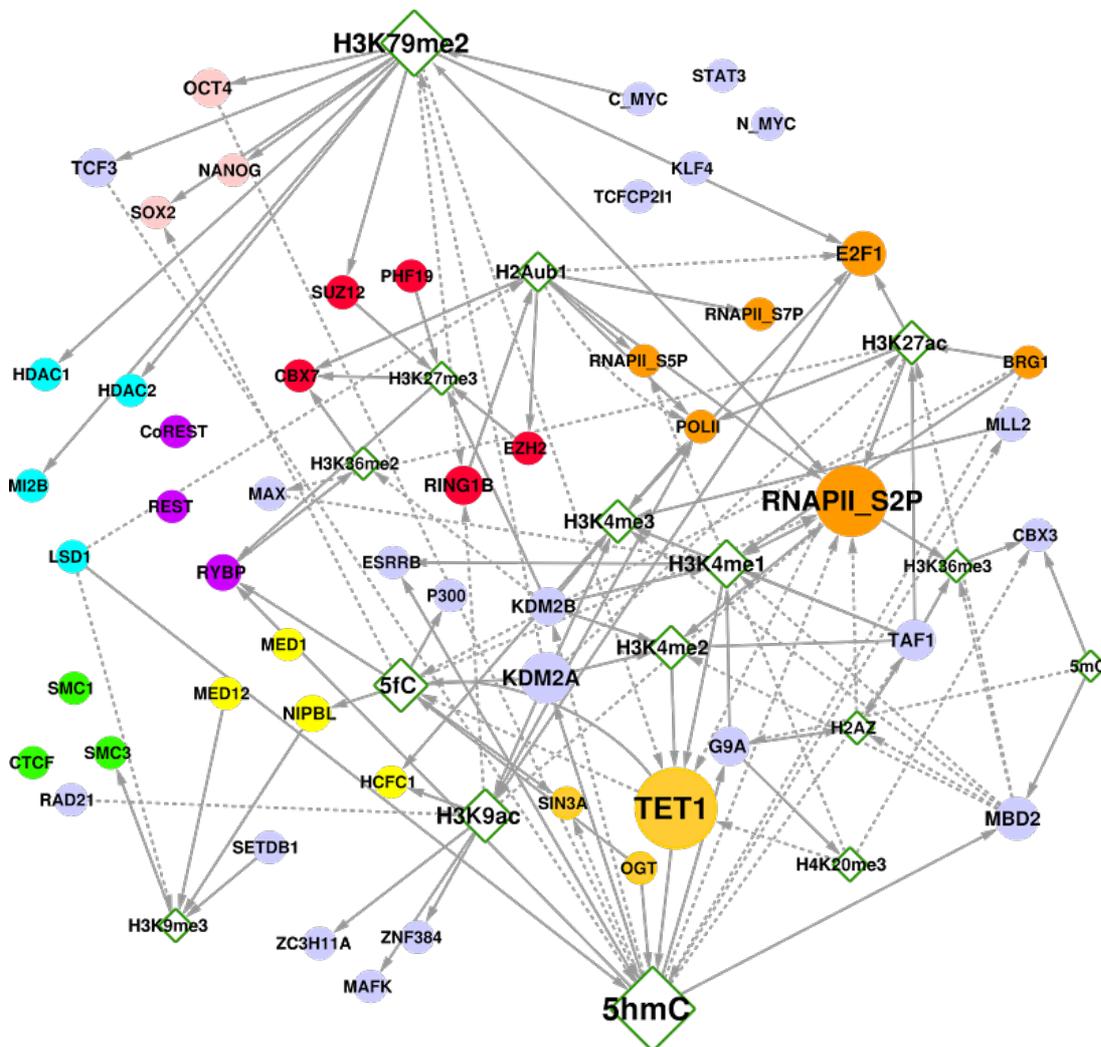


Figure S3. Related to Figure 2. Chromatin States definitions, emission probabilities and genomic annotation enrichments.

A) Emission probabilities of core epigenomic features (left) and genomic annotation enrichments (right) in the 20

chromatin states model. **B)** Chromatin states labels based on core epigenomic features combinations and genomic annotation enrichments.

CAGE_NUC, CAGE_cyto and CAGE_NAST (Fort *et al*, 2014) correspond to CAGE in nuclear compartment, cytoplasmic compartment and non-annotated stem transcripts respectively.

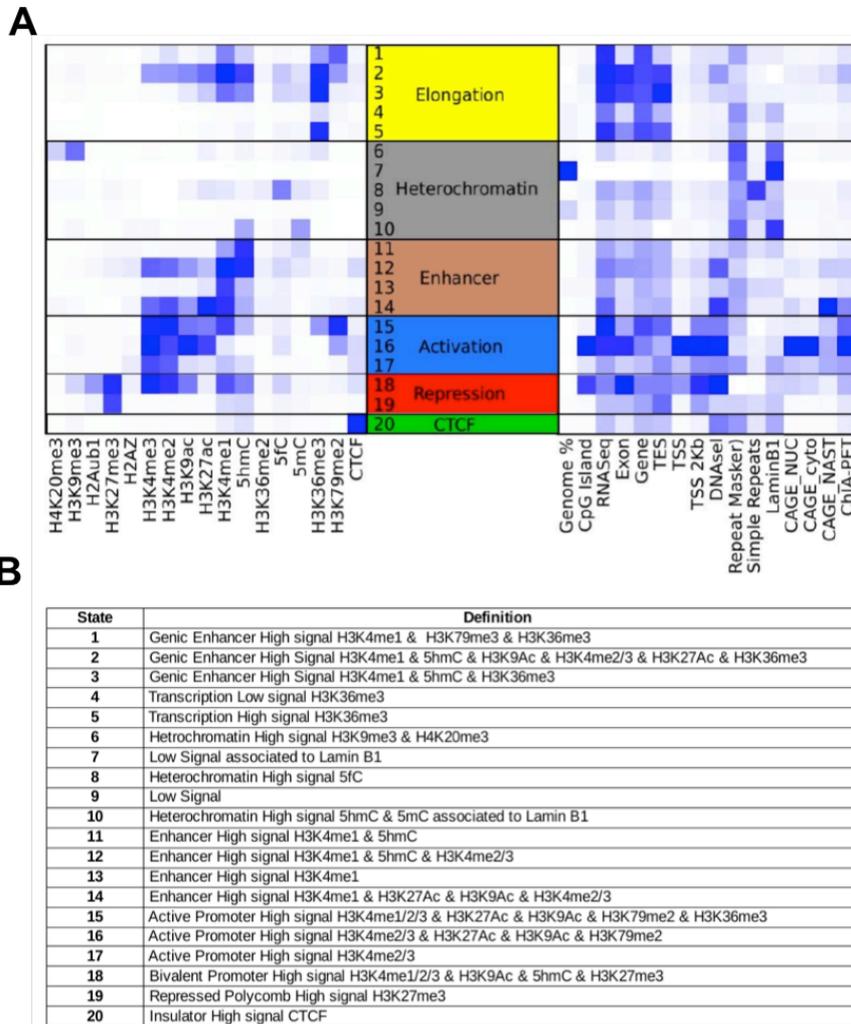


Figure S4. Related to Figure 2. Effect of sample selection.

A) Hierarchical clustering of all available samples based on the similarity measure $1-c(x,y)$. The color-coding in the 'study'-row indicates the experimental origin of the sample by accession-ID. The 'sample'-row indicates experiments detecting the same type of CrP, histone mark or methylation type. **B)** Edges overlap between reference co-location network discussed in the main text and 10 alternative networks built on random selection of a sample out of available ones for each epigenomic features. **C)** Comparison of node degrees for the reference network and mean node degrees for the 10 networks with a randomly selected sample for every epigenomic features ($r = 0.843$, slope = 1.008).

Figure S5. Related to Figure 3. Inferred signal-receiver directions.

A) Inferred signal-receiver directions (blue edges) in the network of epigenetic communication. **B)** Average PageRank-based influences and popularities in 2,000 networks where 10% randomly selected inferred signal-receiver directions were reversed. **C)** Average PageRank-based influences and popularities in 2,000 networks where 50% randomly selected inferred signal-receiver directions were reversed.

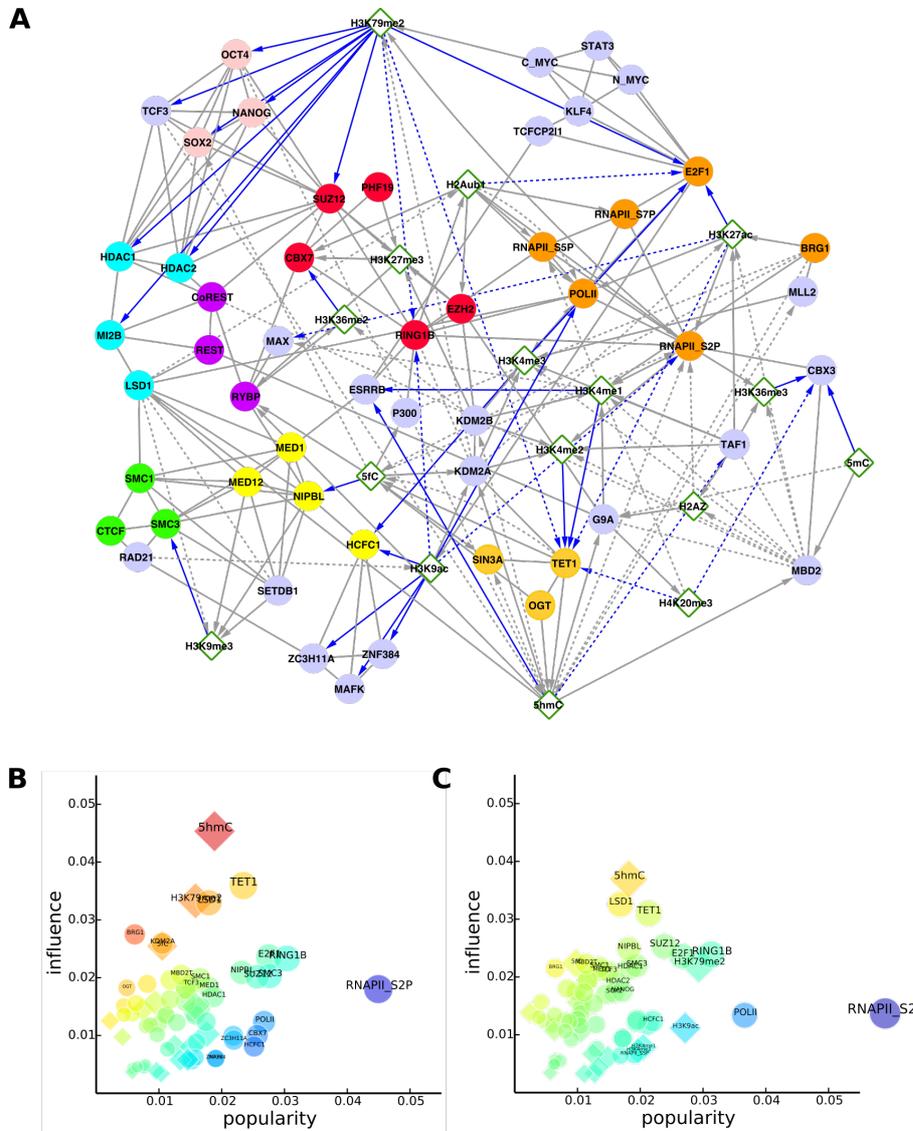
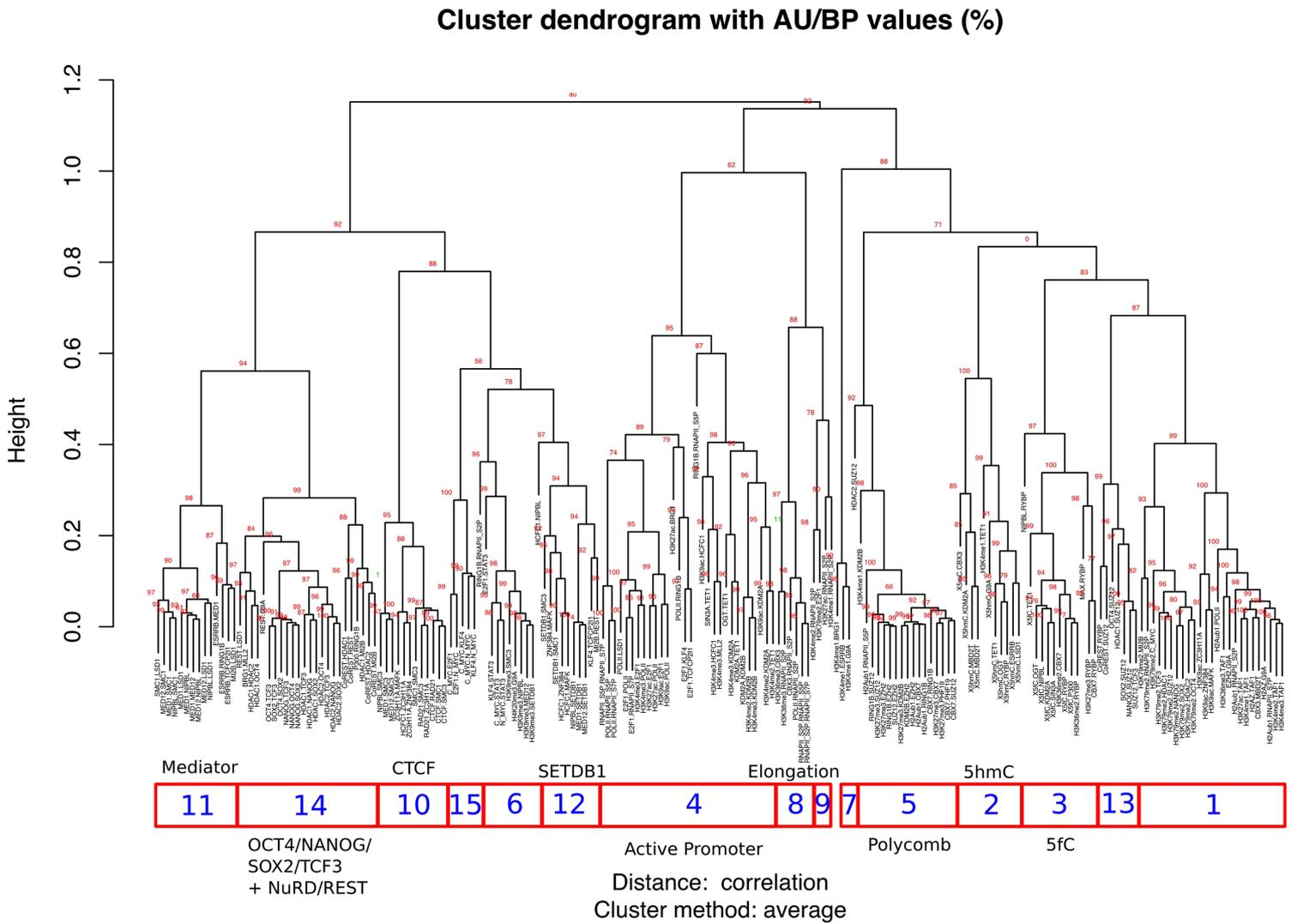


Figure S6. Related to Figure 5. Co-localization pairs clustered according to their distribution in the chromatin states

Hierarchical clustering of the positive interactions in the network by their frequencies along the 20 different chromatin states. The corresponding clusters of interactions (1-15) correspond to the subnetworks (chromnets).



Supplemental Tables

Table S1. Related to Figure 1. ChIP-Seq experiments and annotations included in the study (In Table_S1.xls)

Table S2. Related to Figure 2. Co-localization couplings with emitter or receptor directions (In Table_S2.xls)

Table S3. Related to Figure 4. Co-evolutionary couplings (In Table_S3.xls)

Table S4. Related to Figure 6. Gene Ontology enrichment analysis in 5hmC independent genomic segments (In Table_S4.xls)

Supplemental Experimental Procedures

Read counts and pre-processing for the co-location network inference

We used the ChromHMM segments with a probability higher than 0.95 as samples for the network inference. We filtered all bins for each state that were unexpectedly large (the upper 1% for each state) because they might produce outliers in the data and it is hard to justify where the signal occurs within the region. We counted the overlapping ChIP-Seq reads for the resulting segments using Rsamtools, although some of the ChIP-experiments had to be excluded from the network inference due to the low number of reads per bin, or the low number of bins with signal or study dependent artifacts, including: CTCF_GSE11431, NANOG_GSE11431, LAMIN1B and H3K27me3_GSE36114, SMAD1_GSE11431, MBD1A_GSE39610, MBD1B_GSE39610, MBD2A_GSE39610, MBD3A_GSE39610, MBD4_GSE39610, and MECP2_GSE39610 (as MBD2A was not used, the MBD2 co-localization data corresponds to MBD2T).

Robustness of co-localization network to sample selection

We performed a hierarchical clustering on the number of overlapping ChIP-Seq reads in the selected genomic segments with $1-\text{cor}(x,y)$ as a distance measure. We find that most replicates or functionally related samples fall into the same branch (**Figure S4A**). To further evaluate the effect of sample selection on co-localization network we compared our co-localization network to 10 alternative networks by randomly selecting other replicates. Our results show that the retrieved network is very robust to replicate selection in terms of retrieved interactions (**Figure S4B**). In order to ensure that minor discrepancies between these networks don't affect to network structure we compare average node degrees for random networks to node degrees of our reference network. **Figure S4C** shows good agreement on node degrees ($r = 0.843$, slope = 1.008).

Influence/popularity analysis of the co-location network

The popularity of a node in the chromatin network coincides with the standard PageRank centrality score (Brin and Page, 1998) as computed from the (asymmetric) adjacency matrix of the epigenetic communication network. The influence of a node has been computed as its PageRank score after inverting edges directions in the original network (influence-PageRank, Chepelianskii, 2010). In both cases, the damping factor was set to 0.85. All network analysis were conducted using the NetworkX library (<http://networkx.lanl.gov/>). Robustness of the results of this protocol were evaluated on 2,000 networks where 10% of the edges in the reference network were removed randomly (**Figure S5B**). Similarly, we tested robustness of these results to changes in the more uncertain directions (signal-receiver inferred directions, **Figure S5A**). For this, PageRank analyses were performed on 2,000 networks where 50% of these inferred directed edges were reversed (**Figure S5C**).

Co-evolutionary network inference

We retrieved 46,041 protein trees of sequences at the metazoan level from egglog v4.0 (Powell *et al*, 2014), including over a million protein sequences. We extracted only-unique-orthologous protein trees for each mouse

protein using a species-tree reconciliation approach (Nenadic & Greenacre, 2007) and a previously developed pipeline to deal with tree inconsistencies (Juan *et al*, 2013). Unique orthologs were defined as those with the shortest bidirectional evolutionary distance in the tree.

The inter-orthologs evolutionary distances for the mouse proteins with ChIP-seq data analyzed in this study (NP=58) were encoded in a data matrix, such that each row of this matrix represents the distances in the set of proteins for a given pair of species. For each row, distances were ranked and binned into five equally populated intervals. An additional state, 'NA', was used for any missing values in the distance matrix. We inferred the coupling parameters of a six-state Potts model in which each protein corresponds to a variable restricted to an alphabet of six letters (the five distance bins plus the 'NA' state). To this aim, we maximized an l_2 -regularized version of the (log) pseudo-likelihood (Besag, 1977) of the model parameters, $\{\theta_k^*\} = \text{argmax}_{\theta} [l_{\text{pseudo}}(\{\theta_k\}) - \lambda \sum_k \theta_k^2]$ where θ_k denotes a generic parameter, using a fast asymmetric approach (Aurell *et al*, 2012) and $\lambda = 0.01$. The final coupling parameters $J_{p,q}(d_p, d_q)$ are of special interest since they regulate the interactions between proteins in the model. For example, a strongly positive parameter $J_{p,q}(\text{short}, \text{short})$ can be interpreted as the direct interaction between the two proteins p and q , favoring the co-occurrence of short distances in the respective trees. Co-evolutionary interactions between proteins were ranked using a score function introduced for contact prediction in protein structural analysis (Ekeberg *et al*, 2013). For each pair p, q , we double-centered the sub-matrix $J_{p,q}$ and computed the Frobenius norm $F_{p,q} = [\sum_{a,b=1,5} J_{p,q}(a,b)^2]^{1/2}$. Couplings for the interaction with the 'NA' state were not included in the sum.

Finally, we applied an average product correction (Dunn *et al*, 2008) and obtained the co-evolutionary coupling between proteins p and q as $C_{p,q} = F_{p,q} - F_p F_q / F$ where F is the mean value of $F_{p,q}$ across all the pairs and F_p and F_q respectively the mean values for the proteins p and q .

In order to assess the statistical significance of co-evolutionary couplings, we repeated the analysis on 10,000 sets of randomly selected mouse proteins. The randomized sets contain the same number of proteins (58) as our set of chromatin modifiers. P-values were assigned from the resulting null distribution and associations supported by p values < 0.05 were regarded as statistically significant (Table S3).

Supplemental References

- Dunn, SD, Wahl, LM & Gloor, GB (2008) Mutual information without the influence of phylogeny or entropy dramatically improves residue contact prediction. *Bioinformatics* **24**,333–340
- Ekeberg M, Lovkvist C, Lan Y, Weigt M & Aurell E (2013) Improved contact prediction in proteins: Using pseudolikelihoods to infer Potts models. *Phys. Rev. E* **87**: 012707
- Fort A, Hashimoto K, Yamada D, Salimullah M, Keya CA, Saxena A, Bonetti A, Voineagu I, Bertin N, Kratz A, Noro Y, Wong C-H, de Hoon M, Andersson R, Sandelin A, Suzuki H, Wei C-L, Koseki H, Hasegawa Y, Forrest ARR, et al (2014) Deep transcriptome profiling of mammalian stem cells supports a regulatory role for retrotransposons in pluripotency maintenance. *Nat. Genet.* **46**: 558–566
- Nenadic O & Greenacre M (2007) Correspondence Analysis in R, with Two- and Three-dimensional Graphics: The ca Package. *J. Stat. Softw.* **20**: 1–13