

1 MICA: A multi-omics method to predict gene regulatory networks in early human
2 embryos

3
4 Gregorio Alanis-Lobato^{1,*}, Thomas E. Bartlett^{2,*}, Qiulin Huang^{1,3*}, Claire Simon¹, Afshan
5 McCarthy¹, Kay Elder⁴, Phil Snell⁴, Leila Christie⁴ & Kathy K. Niakan^{1,3,5,6}

6
7 1 Human Embryo and Stem Cell Laboratory, The Francis Crick Institute, London, UK.

8 2 Department of Statistical Science, University College, London WC1E 7HB, UK.

9 3 The Centre for Trophoblast Research, Department of Physiology, Development and
10 Neuroscience, University of Cambridge, Cambridge CB2 3EG, UK.

11 4 Bourn Hall Clinic, Bourn, Cambridge, CB23 2TN.

12 5 Wellcome – Medical Research Council Cambridge Stem Cell Institute, University of

13 Cambridge, Jeffrey Cheah Biomedical Centre, Puddicombe Way, Cambridge CB2 0AW, UK

14 6 Epigenetics Programme, Babraham Institute, Cambridge CB22 3AT, UK

15 *Equal contribution

16 Corresponding authors: kkn21@cam.ac.uk; g.alanis.lobato@gmail.com;

17 thomas.bartlett.10@ucl.ac.uk

18
19
20 **Abstract**

21
22 Recent advances in single-cell -omics have transformed characterisation of cell types in
23 challenging to study biological contexts. In contexts with limited single-cell samples, such as
24 the early human embryo inference of transcription factor-gene regulatory network (GRN)
25 interactions is especially difficult. Here we assessed application of different linear or non-linear
26 GRN predictions to single-cell simulated and human embryo transcriptome datasets. We also
27 compared how expression normalisation impacts on GRN predictions, finding that transcripts
28 per million reads outperformed alternative methods. GRN inferences were more reproducible
29 using a non-linear method based on mutual information (MI) applied to single-cell
30 transcriptome datasets refined with chromatin accessibility (CA) (called MICA), compared to
31 alternative network prediction methods tested. MICA captures complex non-monotonic
32 dependencies and feedback loops. Using MICA, we generated the first GRN inferences in
33 early human development. MICA predicted co-localisation of the AP-1 transcription factor
34 subunit proto-oncogene JUND and the TFAP2C transcription factor AP-2 γ in early human
35 embryos. Overall, our comparative analysis of GRN prediction methods defines a pipeline that
36 can be applied to single-cell multi-omics datasets in especially challenging contexts to infer
37 interactions between transcription factor expression and target gene regulation.

38 Introduction

39

40 Following the fusion of the oocyte and sperm, the zygote undergoes a series of cell divisions
41 until it forms a blastocyst prior to implantation into the uterus. A human blastocyst is formed of
42 a fluid-filled cavity and approximately 200 cells that comprise three distinct cell types: the
43 trophoctoderm (TE), which gives rise to fetal components of the placenta; the primitive
44 endoderm (PE), which forms the yolk sac; and the pluripotent epiblast (EPI), which gives rise
45 to the embryo proper (Blakeley et al., 2015). The specification of these three lineages
46 represents the earliest cell-fate decisions in humans. Understanding the molecular
47 mechanisms that regulate these decisions is important for applications including stem cell
48 biology, regenerative medicine, and reproductive technologies (Niakan et al., 2012).

49

50 We do not yet understand how cell fate specification is regulated in the human embryo.
51 Transcription factor (TF) and target gene regulatory interactions associated with cell fate
52 specification in this context would be informative and has not been elucidated. Defining the
53 gene regulatory networks (GRNs) associated with a given cell type at a distinct time in
54 development facilitates characterisation of cell type identity and prediction of the
55 transcriptional regulators and cis-regulatory DNA sequences that may underlie cell fate
56 specification (Davidson and Erwin, 2006; Materna and Davidson, 2007; Peter and Davidson,
57 2011). Challenges to determining GRNs in early embryos include the small number of cells
58 that contribute to the distinct cell lineages of the human embryo at these early stages of
59 development. Moreover, while single cell -omics technologies facilitate the characterisation of
60 genome-wide gene expression and chromatin accessibility changes in human blastocysts
61 (Blakeley et al., 2015; Li et al., 2018; Liu et al., 2019; Petropoulos et al., 2016; Yan et al.,
62 2013), it is unclear if resolution at this level would allow for accurate predictions of GRNs
63 contributing to cell type specific identity.

64

65 The most common computational approach to infer TF-gene regulatory interactions is to model
66 the expression of each target gene as either a linear or non-linear combination of the
67 expression of a set of potential regulators (e.g., the TFs in the dataset of interest) (Dobra et
68 al., 2004; Marbach et al., 2012). The 'target gene' method uses advanced regression methods,
69 such as sparse penalised regression or random forests, to model the expression level of a
70 particular 'target' gene (the response), conditional on the expression levels of a set of other
71 genes (predictors, such as TFs) (Haury et al., 2012; Huynh-Thu et al., 2010). By combining
72 these local network model fits genome-wide, with each gene taking a turn as the response or
73 'target gene' (Dobra et al., 2004), the GRN is constructed. The linear regression model has
74 previously been used in combination with chromatin accessibility to infer GRNs in human brain
75 organoids (Fleck et al., 2022). Bayesian networks also model the expression of each gene
76 conditional on a set of regulating genes (i.e., parents in the network) by computing how likely
77 it is that each TF in the dataset was responsible for the expression of a certain gene (Kamimoto
78 et al., 2020; Marbach et al., 2012; Pe'er, 2005; Tipping, 2001). The Bayesian/Bagging Ridge
79 model in combination with chromatin accessibility has been used effectively to predict the
80 consequence of gene perturbations and to understand cell fate transitions in the context of
81 development and cellular reprogramming (Argelaguet et al., 2022; Kamimoto et al., 2023). Of
82 note, when fitting Bayesian models it is usually necessary to assume global compatibility (of
83 the local network models), which in turn typically requires that the graph contains no loops
84 (acyclicity). This an assumption which is challenging for biological systems, where feedback
85 loops are a persistent feature (Brandman and Meyer, 2008). In another important group of
86 methods, TF-gene (or more generally gene-gene) interactions are ranked based on a pairwise
87 similarity measure of their expression across cell samples, such as a variant of correlation or
88 mutual information (Daub et al., 2004; Krishnaswamy et al., 2014). Again, alternatives are
89 available that consider conditional statistics (e.g., partial correlation or partial information
90 decomposition) (Chan et al., 2017). Finally, a range of 'black-box' machine learning methods
91 are available that model the expression of each gene in hard to visualise ways using for

92 example sophisticated non-linear combinations of TF expression values (Shu et al., 2021;
93 Wang et al., 2020).

94
95 Although GRN inference methods have been comprehensively compared and developed for
96 whole-tissue or bulk transcriptome analysis (Cahan et al., 2014; Chai et al., 2014; Marbach et
97 al., 2012; Morris et al., 2014; Thompson et al., 2015), the interest in cell type-specific
98 regulatory networks has prompted the application of these methods to single cell RNA-seq
99 (scRNA-seq) datasets (Aibar et al., 2017; Badia-I-Mompel et al., 2023; Chen and Mar, 2018;
100 Fiers et al., 2018; Iacono et al., 2019; Kang et al., 2021; Mochida et al., 2018; Nguyen et al.,
101 2021; Pratapa et al., 2020; Stone et al., 2021). These previous studies showed that the
102 approaches with the best overall performance belong to the regression and mutual information
103 categories and highlighted the challenges of performing GRN inference on single cell datasets.
104 This is attributed mainly to the heterogeneous and stochastic nature of gene expression in
105 individual cells (Nguyen et al., 2021). Of note is that even state-of-the-art methods that have
106 been evaluated on well-defined benchmark datasets predict a considerable number of false
107 positive interactions (Kang et al., 2021; Marbach et al., 2012; Nguyen et al., 2021; Pratapa et
108 al., 2020; Stone et al., 2021). Most of these false predictions originate from indirect
109 associations, for example a path $a \rightarrow b \rightarrow c$ can result in the prediction of $a \rightarrow c$ even if there
110 is no direct link between those nodes (i.e., a 'transitive edge' in the network). Several methods
111 to eliminate such effects have been proposed (Barzel and Barabási, 2013; Chan et al., 2017;
112 Feizi et al., 2013; Margolin et al., 2006; Wang et al., 2018) but they can be computationally
113 demanding and sensitive to the choice of hyper-parameters (Feizi et al., 2013).

114
115 Given the stage of development of human preimplantation embryos and their precious nature,
116 together with the restrictions on such research in some countries, the -omics datasets from
117 human blastocysts are very small compared to those from other biological contexts. This
118 makes it challenging to mine these datasets using GRN inference methods, which require a
119 sufficiently large number of cells to produce reproducible results (Kang et al., 2021; Pratapa
120 et al., 2020). This tight restriction on sample sizes places corresponding restrictions on
121 statistical power, and means that the optimal statistical network inference methodology may
122 be specific to this context. Here, we assessed whether the integration of other types of -omics
123 datasets with transcriptomic-based predictions could help to reduce indirect TF-gene
124 relationships and thereby produce more reliable gene regulatory networks in this setting,
125 especially with restricted sample sizes. We utilised a low number of cell samples in high quality
126 -omics data from early human embryos at the blastocyst stage (6-7 days post-fertilisation)
127 (Blakeley et al., 2015; Petropoulos et al., 2016; Yan et al., 2013) to evaluate the plausibility of
128 our predictions: this has the advantage of being a biological context with only three well-
129 defined cell types to evaluate the plausibility of our predictions.

130
131 The aims of this work were (i) to evaluate whether it is possible to infer reliable cell type-
132 specific regulatory networks for each one of the cell types in the early human embryo in spite
133 of the size of the available -omics data, (ii) to determine if the integration of chromatin
134 accessibility data with transcriptome analysis would better inform predictions of GRNs, and
135 (iii) to predict and validate previously unidentified gene regulations in the human blastocyst.
136 Overall, we demonstrate that the available single cell transcriptomic data was most robustly
137 analysed by a non-linear mutual information-based inference method which had been refined
138 with chromatin accessibility data (MICA). The resulting analysis predicted the first GRN in
139 human preimplantation development and showed that the interactions were consistent with
140 the transcripts and proteins that are known to be enriched in specific lineages. MICA predicted
141 a novel putative regulatory interaction between the TFAP2C transcription factor AP-2 γ and the
142 AP-1 transcription factor subunit proto-oncogene JUND in human preimplantation embryos.
143 Overall, we propose that MICA will be an informative method to make gene regulatory network
144 predictions in other challenging-to-study biological contexts with limitations in the number of
145 cells that can be analysed.

147 Results

148

149 Assessing the conditions for optimal GRN inference on synthetic data

150

151 We generated synthetic gene-expression data with known ground-truth GRN structure for a
152 variety of sample sizes (see Methods) to determine the characteristics of the transcriptomic
153 data on which GRN inference methods perform best and, importantly, to compare these
154 conditions with our human blastocyst scRNA-seq datasets (see Fig. 1A). We then applied four
155 different GRN prediction strategies to the simulated transcriptomes (see Methods for more
156 details). 1) We applied GENIE3 (Aibar et al., 2017; Huynh-Thu et al., 2010), a random forests
157 regression approach which takes into account expression levels of the regulated or target
158 gene and is used for regulatory network inference in the SCENIC GRN prediction pipeline
159 (Aibar et al., 2017; Kang et al., 2021; Marbach et al., 2012; Nguyen et al., 2021; Pratapa et
160 al., 2020; Stone et al., 2021). 2) We compared this to the target gene approach using the LOL2
161 sparse penalised regression model (Hazimeh and Mazumder, 2018; Tibshirani, 1996) which
162 minimizes the total size of the linear model coefficients so that a minimal set of regulatory TFs
163 receive the largest linear model coefficients. 3) We also compared with the correlation
164 coefficient as a reference: we used Spearman's rank correlation coefficient (Fieller et al., 1957)
165 as we expect non-linear associations in the data. 4) Lastly, we applied a non-linear alternative
166 measure of pairwise association based on mutual information (MI) that we speculated may be
167 more appropriate for non-linear response functions in single-cell data (Faith et al., 2007;
168 Krishnaswamy et al., 2014) (Fig. S1A). For the MI method, we used the empirical distribution
169 of the MI values for each gene (Faith et al., 2007). We also varied the number of transcriptional
170 regulators of each target gene in the synthetic data, to study the impact of this parameter on
171 their performance. Further details of the analysis methods and step-by-step code are also
172 available at: https://github.com/galanisl/early_hs_embryo_GRNs

173

174 Two main conclusions emerged from these simulations. First, using synthetic datasets, the
175 performance of all methods increases from the sample size of $n = 10$ to $n = 1000$, as assessed
176 statistically by the area under the receiver operating characteristic curve (AUROC) (Fig. 1B).
177 However, we noted that the prediction accuracy of the GRNs only marginally improves as the
178 number of cell samples surpasses $n = 100$ (Fig. 1B). Therefore, while sample size is important,
179 increasing this beyond a threshold sample size does not further improve inference of the
180 predicted transcriptional regulations. We further note that low sequencing depth (such as in
181 10x Genomics single-cell studies) may increase this optimal value of n . Second, limiting the
182 number of potential transcriptional regulators of a target gene positively impacts the ability of
183 a chosen inference method to recover the ground-truth GRN (Figs. 1B,C). From our
184 simulations, 50 or fewer TFs yield the highest AUROC values for $n < 100$ (Fig. 1B). This would
185 be a reasonable prediction for a given gene in a specific cell type based on our analysis of low
186 input chromatin accessibility (CA) data from the human blastocyst (Liu et al., 2019), where the
187 median number of TF motifs per gene is 35 for the inner cell mass (ICM) and 40 for the TE
188 (Fig. S1B).

189

190 Based on these results and given the size of our human blastocyst datasets (Fig. 1C EPI: 26
191 cells, PE: 33 cells, TE: 45 cells with 1,366 TFs among 25,098 genes (Blakeley et al., 2015;
192 Petropoulos et al., 2016; Yan et al., 2013)), we reasoned that TF-gene interactions predicted
193 from experimentally collected gene expression data could be refined with complementary,
194 context-specific epigenomic datasets. Specifically, we used low-input CA analysis of human
195 blastocyst TE cells and the inner cell mass (ICM), comprised of EPI and PE cells, to refine the
196 GRNs of the respective cell types with putative cis-regulatory interactions (Liu et al., 2019).
197 Peak calling and annotation were performed with nf-core/atacseq to identify regions of open
198 chromatin (Ewels et al., 2020). TF motif enrichment analysis in these open regions of open
199 chromatin were identified using rgt-hint with TF binding models from HOCOMOCO and
200 JASPAR (Li et al., 2019). Finally, the open chromatin regions enriched for TF motifs along with
201 the predicted downstream target genes were determined based on the nearest transcriptional

202 start sites (most distances ranging from -5kbp to 10kbp). This narrowed down the number of
203 genes (to 12,780 for EPI and PE, 12,981 for TE) and TFs (514) in the datasets, thus bringing
204 our sample sizes and potential TFs per gene to the ranges suggested by our simulations.

205

206

207 **Statistical evaluation of inferred human blastocyst GRNs**

208

209 We next integrated the single-cell transcriptome data (Fig. S1C) from human blastocysts
210 generated using the SMART-seq2 library preparation protocol (Blakeley et al., 2015;
211 Petropoulos et al., 2016; Picelli et al., 2014; Yan et al., 2013) where fewer cells are sequenced
212 but better transcript coverage and sequencing depth is obtained (~15,000 genes and ~7
213 million reads per cell) compared to more conventional scRNA-seq methods, such as 10x
214 Genomics (~3,000 genes and ~50,000 reads per cell) (Wang et al., 2021). Due to the lack of
215 genome-scale experimentally derived GRNs in the human embryo context, we evaluated the
216 robustness of the predictions made by four inference approaches (L0L2 regression, GENIE3,
217 Spearman correlation and MI) using three different strategies (see Methods). Since we did not
218 find reports about the impact of gene expression normalisation choice on GRN predictions,
219 we also assessed this parameter and considered the application of GRN predictors to
220 $\log(\text{TPM}+1)$, $\log(\text{FPKM}+1)$, log-count and batch-corrected expression data.

221

222 We calculated a reproducibility score R for each putative regulatory interaction following the
223 application of the network inference methods to human blastocyst data (see Methods) and
224 investigated the distribution of the reproducibility values for the top-100,000 predicted edges
225 (Figs. 2A and S2). The reproducibility estimator R estimates the posterior probability of seeing
226 a network edge given the data; it quantifies the robustness or stability of the inference of this
227 network edge. The accuracy quantified by the R reproducibility statistic relates to the stability
228 of the model predictions to perturbation of the data. Fig. 2A reports the difference between the
229 median of this distribution and the median of the distribution produced by a predictor that
230 generates a random ranking of all possible TF-gene interactions (ΔR). We found that the most
231 reproducible regulatory interactions were inferred by MI followed by a filtering process in which
232 only TF-gene associations supported by CA data were considered in the final network (MICA)
233 (Fig. 2A). As predicted by our simulations, the size of each cell type-specific dataset had a
234 clear impact on GRN inference with methods such as GENIE3 producing more reproducible
235 interactions in the TE. Interestingly, most inference methods produced better results with gene
236 expression values following $\log(\text{TPM}+1)$ or $\log(\text{FPKM}+1)$ normalisation (Fig. 2A).

237

238 We also assessed robustness at the level of network structures, features, or subnetworks. To
239 do so, we randomly split each dataset (EPI, PE and TE) into two groups with the same number
240 of cells, then applied the inference methods to each one (with or without CA refinement), and
241 finally benchmarked them either focussing on the top regulatory interactions (early
242 recognition) or on all regulatory interactions (see Methods). We repeated this 10 times for 10
243 random splits of the data. For the first comparison, we identified the top 1% (by network score)
244 of all possible regulatory interactions from each group, and then quantified the overlap of these
245 top interactions between the groups using the area under the precision-recall curve (AUPRC)
246 (Figs. 2B and S3). By comparing the differences between the median of ten resulting AUPRCs
247 for each method to the random predictions described above, we found that MICA
248 outperformed the other network prediction approaches and gave the highest difference in
249 AUPRC when applied to $\log(\text{TPM}+1)$ - and $\log(\text{FPKM}+1)$ -normalised expression values (Fig.
250 2B). For the second comparison, we divided the data into two portions and then calculated a
251 normalised L2 loss between the network scores over the whole network, rather than just the
252 top 1% of interactions (Figs. 2C and S4). We found that the inverse value of this metric also
253 confirms the robustness of MICA (Figs. 2C and S4).

254

255 Overall, our statistical evaluation strategies showed that the most robust GRN inference
256 method to analyse the limited number of human blastocyst cells was MICA. Importantly, if

257 transcriptome-based predictions are not refined with CA data, most methods perform just
258 slightly better than random (Figs. 2 and S2-S4), underscoring the importance of integrating
259 multi-omics analysis in inference models. It was also important to assess the impact of gene
260 expression normalisation on GRN prediction because we observed apparent effects in our
261 benchmarks depending on the normalisation method used, with $\log(\text{TPM}+1)$ or $\log(\text{FPKM}+1)$
262 being the most suitable gene expression units in this context.

263

264

265 **Association of inferred GRNs to human blastocyst cell lineages**

266

267 We next evaluated the inferred GRNs to determine if they could recapitulate interactions of
268 molecular markers of the three cell types that comprise the human blastocyst. First, we
269 computed the overlap between the GRN edges predicted by each inference method for each
270 blastocyst cell type to identify interactions associated with the EPI, PE and TE, as well as the
271 interactions common to the three cell types (see Methods). We then used the out-degree of
272 NANOG, GATA4 and CDX2, TFs which respectively mark the EPI, PE and TE (Figs. 3A,B),
273 as a proxy for their activity in each one of the four networks. Our prediction is that these TFs
274 should be actively regulating genes in the cell type-specific networks and participate in only a
275 few interactions in the common GRN, based on their known expression pattern in the
276 blastocyst and function in other mammalian contexts such as the mouse (Arceci et al., 1993;
277 Chambers et al., 2003; Dietrich and Hiiragi, 2007; Mitsui et al., 2003; Niakan and Eggan, 2012;
278 Roode et al., 2012; Strumpf et al., 2005). The common network did not show marker activity
279 in the predicted GRNs (Figs. 3B and S5-S8), which is expected given that the expression of
280 these markers is known to be mutually exclusive at this stage. The expected pattern for
281 NANOG (active in the EPI and inactive in the PE and TE, Fig. 3A), was observed in the GRNs
282 with TF-gene interactions supported by CA data and inferred with GENIE3, Spearman
283 correlation and MI when applied to batch-corrected, FPKM and TPM data with some instances
284 of the non-refined LOL2, GENIE3 and Spearman GRNs also matching the expected pattern
285 (Fig. 3B). In the PE, the expected activity pattern for GATA4 (active in the PE and inactive in
286 EPI and TE, Fig. 3A) and lack of detectable networks for NANOG and CDX2, were predicted
287 by applying GENIE3+CA, Spearman correlation+CA and MICA to the FPKM and TPM
288 datasets (Fig. 3B). Finally, the TE-expected pattern for CDX2 (active in the TE and inactive in
289 the EPI and PE) and lack of detectable network for NANOG and GATA4, Fig. 3A) was only
290 observed when using Spearman correlation or MICA on log-counts and batch-corrected data,
291 MICA or MICA+CA on FPKM or TPM data (Fig. 3B). LOL2, with or without CA refinement,
292 consistently performed poorly at predicting GRNs compared to the other methods used. Log-
293 count normalisation failed to recover cell type specific GRNs for these TFs.

294

295 We next manually curated a list of gene sets representing the most relevant pathways and
296 biological processes in the EPI (e.g., regulation of pluripotency) and the TE (e.g., embryonic
297 placenta development) to perform a gene set enrichment analysis with the genes involved in
298 the 500 regulatory interactions with the best prediction scores in the EPI and TE GRNs of the
299 25 most active TFs (Tables S1, S2, Methods). We restricted this analysis to the EPI and TE,
300 because similar lists of gene-sets were not available for PE, where it is currently unclear which
301 pathways are most relevant to this cell type. Next, we computed a validation score V as the
302 number of relevant gene sets that were enriched at a significance level of 10% ($p < 0.1$) in the
303 EPI and TE GRNs (Fig. S9). We note that a low significance level was set because the
304 resulting inferences are aggregated, mitigating the effect of false positive hits for individual
305 gene groups. To facilitate the comparison between V scores across GRNs from different
306 inference and normalisation methods, we normalised this score to the maximum score
307 attained at the cell type level (Fig. 3C). We found that the CA-refined data agree better with
308 the cell type-specific gene sets for both the top-predicted EPI and TE interactions, i.e. the
309 interactions with the highest prediction scores (Fig. 3B). The GRNs inferred by Spearman+CA
310 produced the highest V scores in the EPI and TE, and these were not impacted by the
311 normalisation method employed. The second-best inference method was MICA, with

312 consistent V scores across normalisation methods. We note that both MICA and
313 Spearman+CA are non-linear methods which do not involve regression, indicating that
314 advanced regression methods may not be the most effective choice for biological discovery
315 with these restricted sample-sizes. Taken together, both biological evaluation metrics that we
316 considered confirmed the importance of refining single-cell transcriptome-based GRN
317 inferences with CA data and underscored the robustness of MICA in predicting GRNs from
318 the analysed human blastocyst datasets. Based on these results, we decided to focus on the
319 GRNs predicted by this method for the analyses presented in the following sections.

320
321

322 **Predicted TF networks for NANOG, GATA4 and CDX2**

323

324 Using MICA, we constructed GRNs for TFs expressed in the EPI, PE and TE. All MICA GRN
325 predictions can be found on FigShare: doi.org/10.6084/m9.figshare.21968813. For
326 visualization, the predicted GRN for each of the TFs is separated into target and regulator TF
327 networks. Target networks contain a maximum of 25 top potential target TFs of the hub (or
328 central) TF, whereas regulator networks include a maximum of 25 top TFs that potentially
329 regulate the hub. The average expression of the network members across samples of the cell
330 type of interest is represented by the size of the node. MI scores are represented by the
331 thickness of the edges in the network and edge colour highlights. To further refine the MICA-
332 predictions we used the Spearman's rank correlation coefficient between the expression levels
333 of the source and target nodes across samples of the cell type of interest to define correlated
334 or anti-correlated expression. Correlated or anti-correlated node pairs correspond to positive
335 or negative Spearman's rank coefficient with p-value smaller than 0.1, whereas node pairs
336 having p-value equals or larger than 0.1 were defined as uncorrelated.

337

338 Among the top NANOG targets, TFs RREB1, NCAO3, ZNF343, ZFP42, and NME2 are
339 predicted to be positively regulated by NANOG (Fig. 4A). ZFP42 is a pluripotency marker
340 encoding the REX1 protein and has been shown to be a direct target of NANOG in mouse
341 pluripotent stem cells (Shi et al., 2006). ZNF343 is a less well-characterized TF, but multiple
342 NANOG ChIP-seq datasets in both naïve and primed human embryonic stem cells (ESCs)
343 showed high binding score (MACS2 score > 1000) in the proximal region of the transcription
344 start site (TSS) of ZNF343 (Fig. S10A; (Barakat et al., 2018; Chovanec et al., 2021; Lyu et al.,
345 2018)), which suggests direct regulation of ZNF343 by NANOG. Interestingly, NME2 was
346 previously predicted to be a regulator instead of a target of NANOG in mouse pluripotent stem
347 cells using the TENET GRN inference method (Kim et al., 2021). This inconsistency may be
348 due to the lack of chromatin accessibility data for robust directionality inference in the TENET
349 method, or a species difference. Indeed, some mouse naïve pluripotency regulators, such as
350 ESRRB, which is also a direct target of NANOG in mouse pluripotent stem cells (Festuccia et
351 al., 2012), were not expressed in the human EPI (Blakeley et al., 2013). Putative regulators of
352 NANOG predicted from our MICA network analysis contain multiple KLF factors, including
353 KLF3, KLF5, KLF9 and KLF16 (Fig. 4B). Spearman correlation analysis suggest that KLF9
354 and KLF16 potentially down-regulate NANOG expression, based on their anti-correlated
355 expression, whereas KLF3 and KLF5 are identified as potential regulators of NANOG by
356 MICA, but not significantly correlated by Spearman correlation.

357

358 In the PE, GATA4 was predicted to positively regulate the expression of SP8, TET1, and SKIL
359 while repressing the expression of ELF3, TFDP2 POGK, ZNF770, and NR2F2 (Fig. 4C). It is
360 also predicted to be a target of HNF1B and SALL4 and repressed by ETV4 and E2F6 (Fig.
361 4D). These interactions have not been experimentally validated or inferred. However, NR2F2
362 was previously identified as a maturation marker of polar TE (Meistermann et al., 2021). The
363 repression of NR2F2 by GATA4 predicted in the PE suggests a role for GATA4 in maintaining
364 the PE cell identity by inhibiting the polar TE program. Furthermore, it has been shown that
365 SALL4 is required for mouse PE-derived extra-embryonic endoderm cell derivation and
366 knockout of SALL4 in these cells cause down-regulation of GATA4 (Lim et al., 2008). In

367 addition, multiple endoderm genes such as GATA4, GATA6, and SOX17 were shown to be
368 SALL4-bound genes by ChIP-seq (Lim et al., 2008). These findings are consistent with the
369 inferred GATA4 regulatory network.

370

371 Interestingly, TBX3 is a target of CDX2 (Fig. 4E), and the role of TBX3 has been implicated in
372 trophoblast cell differentiation (Lv et al., 2019). CDX2 is predicted to be regulated by both
373 KLF5 and KLF6, which are also molecular markers of the trophectoderm lineage (Fig. 4F).
374 KLF5 is necessary for trophectoderm formation in the mouse pre-implantation embryo and is
375 required for CDX2 expression in mouse ESCs (Ema et al., 2008; Lin et al., 2010). Interestingly,
376 further analysis of the CDX2 interactions showed that TBX3 and FOXH1 are also potential
377 positive regulators of KLF5 and KLF6. Positive feedback loops are predicted between CDX2,
378 TBX3 and KLF5, as well as CDX2, FOXH1 and KLF6 (Fig. S10B). Other key developmental
379 regulators including HAND1, GATA3 and GATA6, are also predicted to regulate CDX2
380 expression (Fig. 4F). This is consistent with the timing of GATA3 protein expression preceding
381 that of CDX2 (Gerri et al., 2020). Further investigation will be needed to understand the
382 differences between the CDX2-high and CDX2-low TE cells and how the positive feedback
383 loops formed within the CDX2 network enhance and stabilize CDX2 expression in the CDX2-
384 high population.

385

386

387 **Maintenance of TFAP2C expression by JUND in all lineages of human blastocyst**

388

389 We next sought to determine if the MICA network modelling predicts novel interactions or
390 associations between TFs in early human embryo development. We focused on TFAP2C as
391 an example for the network comparison. TFAP2C is a molecular marker that is initially
392 expressed in all the cells at the morula stage in mouse embryos and later specifically restricted
393 to the TE at the blastocyst stage and it is not expressed in other lineages (Cao et al., 2015;
394 Gerri et al., 2020). However, in human embryos, TFAP2C is expressed in the TE, EPI and PE
395 at the blastocyst stage (Blakeley et al., 2015) and has been shown to maintain pluripotency in
396 naïve human embryonic stem cells by regulating OCT4 expression (Chen et al., 2018; Pastor
397 et al., 2018). By performing comparisons of TFAP2C networks in human EPI, TE, and PE
398 cells, we found a conserved putative interaction consisting of TFAP2C, JUND, SOX4, and
399 GCM1 (Fig. 5A). In addition, all four factors showed significant positive correlations in their
400 expression (Figs. 5B, S11A-B). MICA predicts that JUND and SOX4 regulate TFAP2C, while
401 TFAP2C targets GCM1 in all three lineages. In the EPI and PE, SOX4 and GCM1 formed a
402 feedback loop. Interestingly, in the TE, the correlation between GCM1 and SOX4 is absent,
403 and SOX4 is predicted to target JUND, which in turn may regulate TFAP2C. Our network
404 predicts that interactions between these four TFs maintains TFAP2C expression in all three
405 lineages.

406

407 To determine if TFAP2C and JUND protein expression is positively correlated in embryos, as
408 predicted by MICA, we performed immunofluorescence analysis of human preimplantation
409 blastocysts approximately 6.5 days after fertilisation. We observed that TFAP2C and JUND
410 were co-expressed in cells of the human blastocysts (Fig. 5C, S13A; n=4). We next performed
411 3D nuclear segmentation and calculated the Pearson correlation coefficient between TFAP2C
412 and JUND based on the DAPI-normalized protein intensity (Fig. 5D). TFAP2C intensity
413 showed a linear increase with the increase of JUND intensity in all analysed human embryos,
414 and their intensities are highly correlated (0.68-0.85 Pearson correlation; $p < 0.001$; Fig. 5D).
415 Overall, this demonstrates the informativeness of the MICA network analysis to predict
416 correlations and possible regulatory relationships between transcription factors that can be
417 experimentally tested.

418

419 In addition to the MICA predictions, we also performed Spearman correlation analysis of the
420 significant edges identified by the MICA analysis. Surprisingly, around half of the interactions
421 are in non-linear fashion, highlighting the informativeness of MICA to capture complex non-

422 monotonic dependencies. From the TFAP2C network predicted in the EPI lineage, we found
423 that TFAP2C potentially regulates FOXO3 and ZFP42 in a non-linear manner (Fig. S12). The
424 expression of FOXO3 and ZFP42 seems to fit better with the expression of TFAP2C on an
425 exponential curve rather than a linear line between TFAP2C and JUND (Fig. S13B). This
426 suggests that when analysing scRNAseq or low input multi-omics analyses similar non-linear
427 correlations may exist and this may have biological significance. It would be interesting to
428 know whether and which types of non-linear interaction predominate and the biological
429 significance of these non-linear regulations.

430
431

432 Discussion

433

434 The relationship between genome-wide transcriptomic and epigenomic changes and cell fate
435 specification in human embryogenesis is unclear. Studies of human pre-implantation
436 development rely on the donation of surplus embryos derived from assisted reproduction
437 technologies, and the use of such embryos for research is tightly regulated and subject to
438 significant limitations, such as a lack of ability to conduct such research in some jurisdictions
439 (Niakan et al., 2012). In addition, the collection of single cells from such precious embryos is
440 technically challenging and requires specialist expertise and micropipettes to disaggregate
441 microscopic embryos. Therefore, currently available -omics datasets from human blastocysts
442 comprise only a few tens of samples per cell type and are therefore very limited. This contrasts
443 with single-cell analyses in other cellular and developmental contexts that are based on tens
444 of thousands of samples (or more) per cell type (Zheng et al., 2017).

445

446 Sample size is one of the most important considerations when selecting or designing statistical
447 methodology, for example to infer networks of regulations of TFs and their target genes.
448 Hence, GRN inference in -omics data from human pre-implantation embryos presents unique
449 statistical challenges. In particular, methodology that can leverage information about gene
450 regulations from small sample-sizes is required for this context. On the other hand, the lack of
451 heterogeneity in early human embryos compared to adult tissue makes this a good context in
452 which to assess GRN inference methodologies, because there is less unmeasured variability
453 arising from environmental factors. To assess whether GRN inference method can be
454 informative in this challenging-to-study context, we have systematically compared several
455 popular methodologies. Furthermore, we have tested how incorporation of complementary
456 cis-regulatory epigenomic data from ATAC-seq improves GRN inferences. Consistent with
457 other contexts (Argelaguet et al., 2022; Kamimoto et al., 2020), we found that incorporating
458 chromatin accessibility/TF motif analysis together with transcriptional inferences improves the
459 accuracy of GRN inference, by first narrowing down the choice of TF targets from which to
460 infer the mRNA transcript co-expression network, a principle that is likely to be especially
461 important with small sample sizes. Notably, we analyse RNA-seq data generated using
462 SMART-seq2 and the sequencing method and depth of sequencing will likely impact on the
463 choice of GRN inference method.

464

465 Here we showed that incorporating complementary epigenomic data with transcriptomic data
466 improves the reproducibility of inferred GRNs. Furthermore, it has enabled us to make
467 predictions about GRNs operational in early human embryos that are consistent with an
468 understanding of the function and association of the regulators in other developmental and
469 stem cell contexts. We suggest that for network inference using advanced regression
470 methods, it may be preferable to pre-filter using epigenomic data to first narrow down the TF
471 targets in scRNA-seq datasets, because this gives the network inference algorithm an easier
472 time of selecting the regulators from transcriptome data, which is consistent with the
473 experience of others using alternative GRN prediction methods (Badia-I-Mompel et al., 2023;
474 Fleck et al., 2022; Kamimoto et al., 2020; Kartha et al., 2022). Alternative methods incorporate
475 epigenomic data after inference of the mRNA transcript co-expression network (González-
476 Blas et al., 2022) and this has been applied successfully to infer GRNs in the fly brain

477 (Janssens et al., 2022). It will be interesting to determine how the type and sequential order
478 of incorporating multi-omics datasets impacts on GRN predictions.

479

480 We also note some limitations on the interpretation of the GRNs predicted. For example, one
481 ATAC-seq peak could cover multiple transcription factor binding sites in a region. In such
482 cases, we include all TFs with motifs mapping to this region as potential regulators of the gene
483 with the closest TSS (transcriptional start site) to the region, for subsequent transcript co-
484 expression network modelling. However, it will be interesting to determine if this can be
485 refined, by applying the recently developed chromVAR-Multiome method to human
486 blastocysts to generate an *in silico* ChIP-seq library for this context (Argelaguet et al., 2022).
487 Moreover, the ATAC-seq data applied in this study analysed the ICM data in bulk, without
488 distinguishing between EPI and PE cells. Therefore, cell-type specific chromatin accessibility
489 could not be considered, and specific interactions may have been missed due to the
490 heterogeneity of the data or if the ICM CA data failed to reflect developmentally cis-regulations
491 of more developmentally progressed EPI and PE cells. In addition, predictions based on the
492 nearest TSS will miss long-range enhancers that are known to be important for gene regulation
493 (Schoenfelder and Fraser, 2019). In the future, integration of single-cell or nuclei matched
494 transcriptome and ATAC-seq chromatin accessibility data, like recent studies in the mouse
495 (Argelaguet et al., 2022), would be preferable to apply in the human blastocyst context. We
496 note that the GRN inference methods predicted edges that overlap between the four
497 inference methods. Identifying overlapping inferences by comparing more than one GRN
498 inference method may be a strategy to identify network edges with more confidence due to
499 the agreement between several inference methodologies. However, this strategy may also
500 miss some edges which can only be detected by one method and not another.

501

502 We also note that so far, we have separately modelled GRN structure to specific pre-
503 implantation embryonic cell types constrained to a single developmental time-point. In the
504 future, we seek to model dynamic GRN structure in transcriptionally distinct human blastocyst
505 lineages. Moreover, we seek to integrate CUT&RUN or CUT&Tag TF-DNA binding analysis
506 (Kaya-Okur et al., 2019; Meers et al., 2019; Skene and Henikoff, 2017) for key putative
507 developmental regulators, such as GATA3, to further narrow down experimentally validated
508 occupancy from the ATAC-seq predictions we used in this study, similarly to a recent
509 application in mouse blastocysts (Hainer et al., 2019; Hayashi and Inoue, 2023), though we
510 note that in this context TF occupancy studies will be restricted to a few TFs with good quality
511 antibodies. It will also be important to determine which cis regulatory regions are required for
512 target gene regulation through systematic perturbation studies. As the topic of GRN inference
513 is currently receiving much attention from the computational biology community, it will also be
514 important in subsequent work to compare our pipeline with the latest alternatives beyond
515 SCENIC, such as scMTNI (Zhang et al., 2023).

516

517 In summary, the MICA network analysis pipeline we developed is a tool that can be applied to
518 challenging-to-study developmental contexts with limited sample size, such as the human
519 blastocyst, to make predictions about TF interactions that can be experimentally tested in the
520 future. As more datasets become available, we anticipate that the networks predicted will be
521 further refined.

522

523

524 **Methods**

525

526 **Ethics statement**

527

528 All human embryo experiments followed all relevant institutional and national guidelines and
529 regulations.

530

531 This study was approved by the UK Human Fertilisation and Embryology Authority (HFEA):
532 research licence number R0162, and the Health Research Authority's Research Ethics
533 Committee (Cambridge Central reference number 19/EE/0297).

534

535 The process of licence approval entailed independent peer review along with consideration by
536 the HFEA Licence and Executive Committees. Our research is compliant with the HFEA code
537 of practice and has undergone inspections by the HFEA since the licence was granted.
538 Research donors were recruited from patients at Bourn Hall Clinic, Homerton University
539 Hospital, The Bridge Centre and IVF Hammersmith.

540

541 Informed consent was obtained from all couples that donated surplus embryos following IVF
542 treatment. Before giving consent, people donating embryos were provided with all of the
543 necessary information about the research project, an opportunity to receive counselling and
544 the conditions that apply within the licence and the HFEA code of practice. Donors were
545 informed that, in the experiments, embryo development would be stopped before 14 days
546 post-fertilization, and that subsequent biochemical and genetic studies would be performed.
547 Informed consent was also obtained from donors for all the results of these studies to be
548 published in scientific journals. No financial inducements were offered for donation. Consent
549 was not obtained to perform genetic tests on patients and no such tests were performed. The
550 patient information sheets and consent document provided to patients are publicly available
551 ([https://www.crick.ac.uk/research/a-z-researchers/researchers-k-o/kathy-niakan/hfea-](https://www.crick.ac.uk/research/a-z-researchers/researchers-k-o/kathy-niakan/hfea-licence/)
552 [licence/](https://www.crick.ac.uk/research/a-z-researchers/researchers-k-o/kathy-niakan/hfea-licence/)). Donated embryos surplus to the IVF treatment of the patient were cryopreserved
553 and were transferred to the Francis Crick Institute where they were thawed and used in the
554 research project.

555

556 **ATAC-seq data processing and analysis**

557

558 Chromatin accessibility profiles from the inner cell mass (ICM) and TE were obtained from the
559 data produced by Liu *et al.* with the LiCAT-seq protocol (Liu *et al.*, 2019). Alignment to the
560 reference genome (GRCh38), peak calling and annotation were performed with nf-
561 core/atacseq v1.1.0 (Ewels *et al.*, 2020). Then, we carried out footprinting followed by TF motif
562 enrichment analysis in the regions of open chromatin using rgt-hint v0.13.0 with TF binding
563 models from HOCOMOCO and JASPAR (Li *et al.*, 2019). Finally, we associated the TFs that
564 exhibited over-represented motifs in these regions with their closest transcription starting sites
565 (most distances ranging from -5kbp to 10kbp).

566

567 **scRNA-seq data processing and analysis**

568

569 We integrated scRNA-seq data from three different studies (Blakeley *et al.*, 2015; Petropoulos
570 *et al.*, 2016; Yan *et al.*, 2013) focusing on the three cell types present at the late blastocyst
571 stage (Fig. 1, Fig. S3A). Cell type annotations were taken from the work by Stirparo *et al.*
572 Alignment to the reference genome (GRCh38) and calculation of gene counts and TPM-
573 normalised counts were performed on each dataset separately with nf-core/rnaseq v1.4.2
574 (Ewels *et al.*, 2020). The resulting gene expression matrices were integrated and normalised
575 (log-counts and batch-corrected counts) using Bioconductor tools (Amezquita *et al.*, 2019).
576 The final set of cells was manually curated based on the UMAP representation of the batch-
577 corrected data (see Fig. S3 and Table S1). The list of 25,098 genes in the expression matrices
578 was reduced to the unique set of TFs and transcription starting sites (TSS) derived from the
579 motif enrichment analysis applied to the ATAC-seq data. We used the ICM TFs and TSS for
580 the EPI and PE matrices (12,780 genes) and the TE TFs and TSS for the TE matrices (12,981
581 genes).

582

583 **Network inference methods**

584

585 For network inference, we compared the best performing strategy in the DREAM5 challenge,
586 GENIE3 (Huynh-Thu et al., 2010), with a non-linear alternative based on MI (Faith et al., 2007),
587 the Spearman's rank correlation coefficient, and the LOL2 sparse regression method that we
588 applied using recent advances in sparse multivariate statistical modelling (Bartlett et al., 2019;
589 Hazimeh and Mazumder, 2018), as follows.

590

591 For GENIE3, transcription factors are ranked according to the degree of variability in their
592 expression level and how the expression of the putative regulator correlates with a target gene.
593 This ranking is then used to construct the co-expression network for all genes and transcription
594 factors, by thresholding the algorithm's variance reduction score at the tenth percentile of its
595 empirical distribution. The GRN is then inferred as the intersection of the edges in this co-
596 expression network with the edges in the network of all possible gene regulations derived from
597 the chromatin accessibility / DNA binding motif data. This network of all possible gene
598 regulations is defined as all network edges from regulating TF to regulated gene, where an
599 edge represents a DNA binding motif for the TF in regulatory DNA within open chromatin in
600 the regulated gene. For GENIE3, all default settings were used.

601

602 With the Spearman correlation coefficient, a weighted co-expression network is inferred as
603 the absolute value of the correlation coefficient. The gene regulatory network is then inferred
604 as the intersection of the edges in this co-expression network with the edges in the network of
605 all potential gene regulations derived from the chromatin accessibility / DNA binding motif
606 data.

607

608 For LOL2 regression (described in more detail below), the model automatically chooses a
609 ranked subset of TFs from those predicted as regulators of the target gene in the chromatin
610 accessibility / DNA binding motif data inferences. In this way, each target gene takes its turn
611 for a model to be fitted around that target gene. After the model has been fitted to every target
612 gene, the global gene regulatory network can be constructed by combining the local networks
613 fitted around each target gene. For LOL2 regression, sparsity hyperparameters were chosen
614 using the L0Learn package's internal cross-validation.

615

616 In more detail, for LOL2 regression we start with a linear model of the expression level y of the
617 regulated target gene (Dobra et al., 2004), in terms of the expression levels x_1, x_2, \dots, x_p , of
618 p transcription factors. We want to use the size of the fitted model coefficients b_1, b_2 , etc. to
619 measure the strength of regulation of the target gene by transcription factors (TFs) 1, 2, etc.
620 We use sparse regression to find b_1, b_2 , etc., as this specifically minimises the number of non-
621 zero model coefficients, by requiring that the coefficients b_q are set to zero as much as
622 possible. This leads to a more parsimonious model, in which a relatively small number $p' \ll p$
623 of transcription factors is inferred to be regulating the target gene, as a result of non-zero b_q .
624 Sparse regression minimizes:

625

$$626 \quad [y - (a + b_1x_1 + b_2x_2 + \dots + b_px_p)]^2 + \psi, \quad (1)$$

627

628 where ψ 'penalises' models with values of b_q ($q=1, \dots, p$) further from zero. The most popular
629 choices for ψ include $\psi = \lambda \sum b_j^2$, which is called 'ridge regression' (L2 regression), and $\psi = \lambda$
630 $\sum |b_j|$, which is called 'the lasso' (L1 regression). It can be shown that Eq.1 is equivalent to
631 applying the constraint $\sum b_j^2 < t^2$ for $\psi = \lambda \sum b_j^2$ (L2 regression) or $\sum |b_j| < t$ for $\psi = \lambda \sum |b_j|$ (L1
632 regression). Then, 'best-subset regression' (L0 regression) specifies the number of
633 transcription factors that can have non-zero b_j . It does this by using the constraint $\sum b_j^0 \leq k$,
634 which is equivalent to $\psi = \lambda \sum b_j^0$ in Eq.1. Combinations of these constraints are also often
635 more effective, such as L1L2 regression (also called 'the elastic net' (Zou and Hastie, 2005),
636 which was very successful in genomics), with penalty term $\psi = \gamma \sum |b_j| + \lambda \sum b_j^2$. Recently,
637 LOL2 regression has been proposed as an improvement (Hazimeh and Mazumder, 2018), with
638 penalty term $\psi = \gamma \sum b_j^0 + \lambda \sum b_j^2$. We use LOL2 regression in this study, for reasons as follows.
639 Sparse regression using the L0 penalty is an ideal model for inferring a minimal set of

640 regulating transcription factors, because it specifically selects the best set of k transcription
 641 factors for the model (i.e., the ‘best subset’, of the available transcription factors). Combining
 642 with the L2 penalty in LOL2 makes the model better specified for the data, by minimizing the
 643 total size of the linear model coefficients so that the most important TFs receive the largest
 644 linear model coefficients. We use sparse LOL2 regression to infer the best subset of regulators
 645 of each target gene, from the full list of TF-gene associations supported by the chromatin
 646 accessibility data.

647
 648 For MI-based inference, the co-expression network is estimated first, as follows. We use an
 649 empirical estimate of the distribution of the MI values for each gene (Daub et al., 2004; Faith
 650 et al., 2007). Writing the estimated MI between the expression levels of genes x and y as
 651 M_{xy} , we calculate the approximately $N(0,1)$ distributed variable z_{xy} according to the equation:

$$652 \quad z_{xy}^2 = F^{-1} [\Phi_{y|x}(M_{xy})]^2 + F^{-1} [\Phi_{x|y}(M_{xy})]^2 ,$$

653
 654 where F is the $N(0,1)$ cumulative distribution function (c.d.f.), and $\Phi_{x|y}$ and $\Phi_{y|x}$ are the empirical
 655 c.d.f.s of M_{xy} conditioned on x and y , respectively. From these $N(0,1)$ distributed variables z_{xy} ,
 656 we calculate p -values, and then threshold at a false discovery rate of 5%. Again, the gene
 657 regulatory network is inferred as the intersection of the edges in this co-expression network
 658 with the edges in the network of all possible gene regulations derived from the chromatin
 659 accessibility / DNA binding motif data.

660
 661 We also considered a “random predictor” that outputs a random ordering of all the possible
 662 TF-gene interactions.

663
 664 The putative regulatory links predicted by each of these methods using the scRNA-seq data
 665 were evaluated as *is* but also subjected to a filtering process in which only TF-gene
 666 associations supported by the chromatin accessibility data were considered in the final
 667 network (for GENIE3, MI and Spearman correlation methods). In the case of LOL2 sparse
 668 regression, only TF-gene associations supported by the chromatin accessibility data were
 669 considered as potential regulators of the target gene in the regression model. To identify these
 670 refined predictions, we added the +CA suffix to the name of the GRN inference methods. In
 671 both cases, to generate the final network, we selected the top 100,000 edges by ranking edges
 672 according to their network score. We define network score as the absolute correlation
 673 coefficient or linear regression coefficient (in the case of Spearman’s correlation and sparse
 674 regression respectively), or as $-\log_{10}p$ for mutual information, or the GENIE3 score.

675 676 677 **Simulation study**

678
 679 For our simulation study, we used the `daggitty` package in R to generate synthetic gene
 680 expression data based on pre-defined GRNs with pre-determined network structure. These
 681 GRNs were generated with network edge density of $\rho=0.07$, i.e., 7% of all possible edges (or
 682 gene regulations), are present in the network (this value was estimated from the available
 683 ATAC-seq+TF motif data). The synthetic datasets were generated with a range of sample size
 684 n and number of potential regulators of each node (gene) p , using linear models, as follows.
 685 For each combination of n and p , we generated 100 GRNs; then corresponding observed
 686 gene-expression data-sets were generated for each of the GRNs by specifying that the
 687 expression level of a downstream gene should depend on a linear combination of the
 688 expression levels of its upstream regulators:

$$689 \quad x_d = b_1 * x_{u1} + b_2 * x_{u2} + b_3 * x_{u3} + \dots + b_g * x_{ug} + e ,$$

690
 691 where x_d is the expression level of the regulated (downstream) gene, and $x_{u1} \dots x_{ug}$ are the
 692 expression levels of the g regulating (upstream) TFs, $b_1, b_2 \dots b_g$ are linear model coefficients,
 693 and e is random noise. The different network inference methods (GENIE3, LOL2 sparse
 694

695 regression, MI, Spearman correlation) were then applied to these generated data, and the
696 results were compared with the known pre-determined network structure. This comparison
697 was assessed by the AUC statistic, the 'area under the ROC (receiver-operator characteristic)
698 curve'.

699

700 **Statistical evaluation metrics**

701

702 We calculated a reproducibility score R as a bootstrap estimate of the posterior probability of
703 observing an edge E given the dataset D , $R=P(E|D)$ (Pe'er, 2005), for the top 100,000
704 predicted edges. We also carried out 2-fold cross-validation in two ways as follows. We
705 randomly split each dataset (EPI, PE and TE) in two groups with the same number of cells ten
706 times (repeated 2-fold cross-validation), applied the inference methods (with or without
707 chromatin accessibility refinement), took the top-1% of all possible regulatory interactions from
708 one group and quantified the extent to which the top interactions inferred from the other group
709 coincide with that reference using the area under the precision-recall curve (AUPRC). In
710 addition, we calculated a normalised L2 loss comparing all the network inferences from the
711 two network fits on the 2 folds of the data (at the level of the network scores). The normalised
712 L2 loss is then defined as the L2 loss comparing the network scores from each of the 2 data
713 folds, divided by the product of the square roots of the L2 norms of the network scores of each
714 of the 2 data folds.

715

716 **Biological evaluation metrics**

717

718 For each combination of inference method and normalisation approach, we computed the
719 intersection between the GRN edges predicted to identify interactions that are specific to the
720 EPI, PE and TE, as well as the interactions common to the three cell types. In set notation this
721 corresponds to EPI-specific = $(EPI \setminus PE) \setminus TE$, PE-specific = $(PE \setminus EPI) \setminus TE$, TE-specific = $(TE$
722 $\setminus EPI) \setminus PE$ and Common = $EPI \cap PE \cap TE$. Then, we used the out-degree of NANOG, GATA4
723 and CDX2 (marker TFs of the EPI, PE and TE, respectively) as a proxy for their activity in
724 each one of the four networks. We considered that a marker was active if its normalised out-
725 degree (i.e., the proportion of all genes it regulates in that GRN) was at least the median
726 normalised out-degree across all networks.

727

728 We calculated a validation score V as the number of relevant gene sets that were enriched at
729 a significance level of 10% ($p < 0.1$) in a gene-set enrichment analysis performed with the genes
730 regulated by the 25 most connected TFs and that were part of the 500 interactions with the
731 highest prediction scores produced by the GRN inference methods for the EPI and TE. The
732 list of gene sets for the EPI and TE were manually curated and represent the most relevant
733 pathways and biological processes in each one of these cell types (see Tables S1 and S2).
734 The starting list of gene sets comprised all the pathways and Gene Ontology Biological
735 Process (GO BP) terms collated by the Bader Lab at the University of Toronto on 01 October
736 2020

737 ([http://download.baderlab.org/EM_Genesets/October_01_2020/Human/symbol/Human_GO](http://download.baderlab.org/EM_Genesets/October_01_2020/Human/symbol/Human_GO_BP_AllPathways_no_GO_iea_October_01_2020_symbol.gmt)
738 [BP_AllPathways_no_GO_iea_October_01_2020_symbol.gmt](http://download.baderlab.org/EM_Genesets/October_01_2020/Human/symbol/Human_GO_BP_AllPathways_no_GO_iea_October_01_2020_symbol.gmt)). To facilitate the comparison

739 between V scores from different GRN inference methods within a cell type, we normalised V
740 to the maximum score attained by a network predictor for that cell type (i.e., EPI or TE).

741

742 **Human Embryo Culture**

743

744 Human embryos were cultured as previously described (Gerri et al., 2020). Vitrified human
745 blastocyst stage embryos were thawed using Kitazato Thawing Media (VT602, order number
746 91121) following the manufacturer's instructions.

747

748 Human embryos were cultured in drops of pre-equilibrated Global medium (LifeGlobal; LGGG-
749 20) supplemented with 10% human serum albumin (LifeGlobal; GHSA-125) and overlaid with

750 mineral oil (Origio; ART-4008-5P). Preimplantation embryos were incubated at 37 °C and 5%
751 CO₂ in an EmbryoScope+ time-lapse incubator (Vitrolife) and cultured for up to 24hrs for
752 human blastocyst.

753

754 **Immunofluorescence staining**

755

756 Embryos were fixed with freshly prepared 4% paraformaldehyde in PBS. Fixation was
757 performed for 20 min at room temperature for embryos. Embryos were then washed 3 times
758 in 1× PBS with 0.1% Tween-20 to remove residual paraformaldehyde. Permeabilization was
759 performed with 1× PBS with 0.5% Tween-20 and followed by blocking in blocking solution
760 (10% donkey serum in 1× PBS with 0.1% Tween-20) for 1 h at room temperature on a rotating
761 shaker. Then, antibody incubation was performed with primary antibodies diluted in blocking
762 solution overnight at 4 °C on rotating shaker. The following day, embryos and cell cultures
763 were washed in 1× PBS with 0.1% Tween-20 for 3 times, and then incubated with secondary
764 antibodies diluted in blocking solution for 1 h at room temperature on a rotating shaker in the
765 dark. Next, embryos and cell cultures were washed in 1× PBS with 0.1% Tween-20 for 3 times.
766 Finally, embryos were placed in 1× PBS with 0.1% Tween-20 with Vectashield and DAPI
767 mounting medium (Vector Lab; H-1200) (1:30 dilution). Embryos were placed on μ-Slide 8-
768 well dishes (Ibidi; 80826) for confocal imaging. The antibodies and concentrations used are
769 reported in Table S3.

770

771 **Confocal imaging**

772

773 Confocal immunofluorescence images were taken with Leica SP8 confocal microscope. 2-μm-
774 thick optical sections were collected for embryo Z-stack imaging.

775

776 **Nuclei segmentation and quantification**

777

778 Nuclei segmentation and quantification were performed on Z-stack confocal images taken at
779 2-μm-thick optical sections. Stardist (Weigert et al., 2020) was used for nuclei segmentation
780 followed by CellProfiler (Stirling et al., 2021) for nuclei tracking and fluorescence intensity
781 quantification based on a customized pipeline modified from Lea et al., 2021. First, LIF format
782 confocal images were exported into multi-channel Z-stack TIF format for each Z-stack image
783 using ImageJ (Schneider et al., 2012). Then, Stardist was used to identify the nuclei based on
784 the DAPI channel (358 nm). Finally, Stardist outputs were split into single-channel single-stack
785 images and loaded into the CellProfiler v4.2.5 to track nuclei across image slices and quantify
786 the fluorescence intensity. Customized pipeline and scripts can be found here:
787 https://github.com/galanisl/early_hs_embryo_GRNs

788

789 **Data and code availability**

790

791 scRNA-seq and ATAC-seq data were obtained from the European Nucleotide Archive
792 [accession numbers: PRJNA153427, PRJNA277181, PRJEB11202 and PRJNA494280].
793 Data pre-processing and analysis scripts are available at
794 https://github.com/galanisl/early_hs_embryo_GRNs.

795

796

797 **Acknowledgements**

798

799 We thank Zimeng Chen for contributing to the initial analysis of the synthetic datasets. We
800 thank members of the Niakan lab for helpful discussions and feedback on the manuscript.
801 Work in the laboratory of KKN was supported by the Wellcome (221856/Z/20/Z) and the
802 Wellcome Human Developmental Biology Initiative (215116/Z/18/Z). Work in the laboratory of
803 KKN was also supported by the Francis Crick Institute, which receives its core funding from
804 Cancer Research UK (CC2074), the UK Medical Research Council (CC2074) and the

805 Wellcome Trust (CC2074). The work carried out by TEB was partly supported by MRC grant
806 (MR/P014070/1). For the purpose of Open Access, the authors have applied a CC BY public
807 copyright licence to any Author Accepted Manuscript version arising from this submission.

808
809

810 **Competing Interests**

811

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

813

814

815 **Author contributions**

816 GAL, TEB and KKN conceived the project. GAL and TEB developed the analytical network
817 interference and statistical methods and analysed data. QH developed the network model
818 representation pipeline, analysed data and performed immunostaining analysis and
819 quantification. AM performed human embryo thaw and culture. CS provided advice on image
820 quantification and analysed data. KKN analysed data and funded the project. KE, PS and LC
821 coordinate the donation of embryos to the research project. GAL, TEB, QH, CS and KKN
822 wrote the manuscript.

823

824 **References**

825

826 Aibar, S., González-Blas, C. B., Moerman, T., Huynh-Thu, V. A., Imrichova, H., Hulselmans, G.,
827 Rambow, F., Marine, J.-C., Geurts, P., Aerts, J., et al. (2017). SCENIC: single-cell regulatory
828 network inference and clustering. *Nature Methods* 14, 1083–1086.

829 Amezquita, R. A., Lun, A. T. L., Becht, E., Carey, V. J., Carpp, L. N., Geistlinger, L., Marini, F., Rue-
830 Albrecht, K., Risso, D., Soneson, C., et al. (2019). Orchestrating single-cell analysis with
831 Bioconductor. *Nature Methods* 17, 137–145.

832 Arceci, R. J., King, A. A., Simon, M. C., Orkin, S. H. and Wilson, D. B. (1993). Mouse GATA-4: a retinoic
833 acid-inducible GATA-binding transcription factor expressed in endodermally derived tissues
834 and heart. *Mol Cell Biol* 13, 2235–2246.

835 Argelaguet, R., Lohoff, T., Li, J. G., Nakhuda, A., Drage, D., Krueger, F., Velten, L., Clark, S. J. and Reik,
836 W. (2022). Decoding gene regulation in the mouse embryo using single-cell multi-omics.
837 *bioRxiv* 2022.06.15.496239.

838 Badia-i-Mompel, P., Wessels, L., Muller-Dott, S., Trimbour, R., Ramirez Flores, R.O., Argelaguet, R.
839 and Saez-Rodriguez, J (2023). Gene regulatory network inference in the era of single-cell
840 multi-omics. *Nature Reviews Genetics* doi: 10.1038/s41576-023-00618-5.

841 Barakat, T. S., Halbritter, F., Zhang, M., Rendeiro, A. F., Perenthaler, E., Bock, C. and Chambers, I.
842 (2018). Functional dissection of the enhancer repertoire in human embryonic stem cells. *Cell*
843 *Stem Cell* 23, 276-288.e8.

844 Bartlett, T. E., Kosmidis, I. and Silva, R. (2019). Two-way sparsity for time-varying networks, with
845 applications in genomics. *arXiv* arXiv:1802.08114.

846 Barzel, B. and Barabási, A.-L. (2013). Network link prediction by global silencing of indirect
847 correlations. *Nature Biotechnology* 31, 720–725.

848 Blakeley, P., Fogarty, N. M. E., del Valle, I., Wamaitha, S. E., Hu, T. X., Elder, K., Snell, P., Christie, L.,
849 Robson, P. and Niakan, K. K. (2015). Defining the three cell lineages of the human blastocyst
850 by single-cell RNA-seq. *Development* 142, 3151–3164.

851 Brandman, O. and Meyer, T. (2008). Feedback Loops Shape Cellular Signals in Space and Time.
852 *Science* 322, 390–395.

853 Boroviak, T., Loos, R., Lombard, P., Okahara, J., Behr, R., Sasaki, E., Nichols, J., Smith, A., Bertone, P.
854 (2015). Lineage-specific profiling delineates the emergence and progression of naïve
855 pluripotency in mammalian embryogenesis. *Developmental Cell* 35, 366–82.

856 Butte, A. J. and Kohane, I. S. (1999). Mutual information relevance networks: functional genomic
857 clustering using pairwise entropy measurements. In *Biocomputing 2000*, pp. 418–429.
858 Honolulu, Hawaii, USA: WORLD SCIENTIFIC.

859 Cahan, P., Li, H., Morris, S. A., Lummertz da Rocha, E., Daley, G. Q. and Collins, J. J. (2014). CellNet:
860 Network biology applied to stem cell engineering. *Cell* 158, 903–915.

861 Cao, Z., Carey, T. S., Ganguly, A., Wilson, C. A., Paul, S. and Knott, J. G. (2015). Transcription factor
862 AP-2 γ induces early Cdx2 expression and represses HIPPO signaling to specify the
863 trophoctoderm lineage. *Development* 142, 1606–1615.

864 Chai, L. E., Loh, S. K., Low, S. T., Mohamad, M. S., Deris, S. and Zakaria, Z. (2014). A review on the
865 computational approaches for gene regulatory network construction. *Comput Biol Med* 48,
866 55–65.

867 Chambers, I., Colby, D., Robertson, M., Nichols, J., Lee, S., Tweedie, S. and Smith, A. (2003).
868 Functional expression cloning of Nanog, a pluripotency sustaining factor in embryonic stem
869 cells. *Cell* 113, 643–655.

870 Chan, T. E., Stumpf, M. P. H. and Babbie, A. C. (2017). Gene Regulatory Network Inference from
871 Single-Cell Data Using Multivariate Information Measures. *Cell Systems* 5, 251-267.e3.

872 Chen, S. and Mar, J. C. (2018). Evaluating methods of inferring gene regulatory networks highlights
873 their lack of performance for single cell gene expression data. *BMC Bioinformatics* 19, 232.

874 Chen, D., Liu, W., Zimmerman, J., Pastor, W. A., Kim, R., Hosohama, L., Ho, J., Aslanyan, M., Gell, J. J.,
875 Jacobsen, S. E., et al. (2018). The TFAP2C-Regulated OCT4 Naive Enhancer Is Involved in
876 Human Germline Formation. *Cell Rep* 25, 3591-3602.e5.

877 Chovanec, P., Collier, A. J., Krueger, C., Várnai, C., Semprich, C. I., Schoenfelder, S., Corcoran, A. E.
878 and Rugg-Gunn, P. J. (2021). Widespread reorganisation of pluripotent factor binding and
879 gene regulatory interactions between human pluripotent states. *Nat Commun* 12, 2098.

880 Daub, C. O., Steuer, R., Selbig, J. and Kloska, S. (2004). Estimating mutual information using B-spline
881 functions - an improved similarity measure for analysing gene expression data. *BMC*
882 *Bioinformatics* 5, 118.

883 Davidson, E. H. and Erwin, D. H. (2006). Gene regulatory networks and the evolution of animal body
884 plans. *Science* 311, 796–800.

885 Dietrich, J.-E. and Hiiragi, T. (2007). Stochastic patterning in the mouse pre-implantation embryo.
886 *Development* 134, 4219–4231.

887 Dobra, A., Hans, C., Jones, B., Nevins, J. R., Yao, G. and West, M. (2004). Sparse graphical models for
888 exploring gene expression data. *Journal of Multivariate Analysis* 90, 196–212.

889 Ema, M., Mori, D., Niwa, H., Hasegawa, Y., Yamanaka, Y., Hitoshi, S., Mimura, J., Kawabe, Y., Hosoya,
890 T., Morita, M., et al. (2008). Krüppel-like factor 5 Is Essential for Blastocyst Development and
891 the Normal Self-Renewal of Mouse ESCs. *Cell Stem Cell* 3, 555–567.

892 Ewels, P. A., Peltzer, A., Fillinger, S., Patel, H., Alneberg, J., Wilm, A., Garcia, M. U., Di Tommaso, P.
893 and Nahnsen, S. (2020). The nf-core framework for community-curated bioinformatics
894 pipelines. *Nature Biotechnology* 38, 276–278.

895 Faith, J. J., Hayete, B., Thaden, J. T., Mogno, I., Wierzbowski, J., Cottarel, G., Kasif, S., Collins, J. J. and
896 Gardner, T. S. (2007). Large-Scale Mapping and Validation of Escherichia coli Transcriptional
897 Regulation from a Compendium of Expression Profiles. *PLoS Biology* 5, e8.

898 Feizi, S., Marbach, D., Médard, M. and Kellis, M. (2013). Network deconvolution as a general method
899 to distinguish direct dependencies in networks. *Nature Biotechnology* 31, 726–733.

900 Festuccia, N., Osorno, R., Halbritter, F., Karwacki-Neisius, V., Navarro, P., Colby, D., Wong, F., Yates,
901 A., Tomlinson, S. R. and Chambers, I. (2012). Esrrb Is a Direct Nanog Target Gene that Can
902 Substitute for Nanog Function in Pluripotent Cells. *Cell Stem Cell* 11, 477–490.

903 Fieller, E. C., Hartley, H. O. and Pearson, E. S. (1957). Tests for rank correlation coefficients. I.
904 *Biometrika* 44, 470–481.

905 Fiers, M. W. E. J., Minnoye, L., Aibar, S., Bravo González-Blas, C., Kalender Atak, Z. and Aerts, S.
906 (2018). Mapping gene regulatory networks from single-cell omics data. *Brief Funct Genomics*
907 17, 246–254.

908 Fleck, J. S., Jansen, S. M. J., Wollny, D., Zenk, F., Seimiya, M., Jain, A., Okamoto, R., Santel, M., He, Z.,
909 Camp, J. G., et al. (2022). Inferring and perturbing cell fate regulomes in human brain
910 organoids. *Nature*.

911 Gerri, C., McCarthy, A., Alanis-Lobato, G., Demtschenko, A., Bruneau, A., Loubersac, S., Fogarty, N.
912 M. E., Hampshire, D., Elder, K., Snell, P., et al. (2020). Initiation of a conserved
913 trophoctoderm program in human, cow and mouse embryos. *Nature* 587, 443–447.

914 González-Blas, C. B., De Winter, S., Hulselmans, G., Hecker, N., Matetovici, I., Christiaens, V.,
915 Poovathingal, S., Wouters, J., Aibar, S. and Aerts, S. (2022). SCENIC+: single-cell multiomic
916 inference of enhancers and gene regulatory networks. *bioRxiv* 2022.08.19.504505.

917 Hainer, S. J., Bošković, A., McCannell, K. N., Rando, O. J. and Fazio, T. G. (2019). Profiling of
918 Pluripotency Factors in Single Cells and Early Embryos. *Cell* 177, 1319–1329.e11.

919 Haury, A.-C., Mordelet, F., Vera-Licona, P. and Vert, J.-P. (2012). TIGRESS: Trustful Inference of Gene
920 REgulation using Stability Selection. *BMC Systems Biology* 6, 145.

921 Hayashi, R. and Inoue, A. (2023). Low-Input CUT&RUN for Mouse Oocytes and Preimplantation
922 Embryos. *Methods Mol Biol* 2577, 83–92.

923 Hazimeh, H. and Mazumder, R. (2018). Fast Best Subset Selection: Coordinate Descent and Local
924 Combinatorial Optimization Algorithms. *arXiv* arXiv:1803.01454,.

925 Huynh-Thu, V. A., Irrthum, A., Wehenkel, L. and Geurts, P. (2010). Inferring Regulatory Networks
926 from Expression Data Using Tree-Based Methods. *PLoS ONE* 5, e12776.

927 Iacono, G., Massoni-Badosa, R. and Heyn, H. (2019). Single-cell transcriptomics unveils gene
928 regulatory network plasticity. *Genome Biology* 20, 110.

929 Janssens, J., Aibar, S., Taskiran, I. I., Ismail, J. N., Gomez, A. E., Aughey, G., Spanier, K. I., De Rop, F. V.,
930 González-Blas, C. B., Dionne, M., et al. (2022). Decoding gene regulation in the fly brain.
931 *Nature* 601, 630–636.

932 Kamimoto, K., Stringa, B., Hoffmann, C. M., Jindal K., Solnica-Krezel, L., and Morris, S. A. (2023).
933 Dissecting cell identity via network inference and in silico gene perturbation. *Nature* 614,
934 742–751.

935 Kang, Y., Thieffry, D. and Cantini, L. (2021). Evaluating the Reproducibility of Single-Cell Gene
936 Regulatory Network Inference Algorithms. *Front. Genet.* 12, 617282.

937 Kartha, V. K., Duarte, F. M., Hu, Y., Ma, S., Chew, J. G., Lareau, C. A., Earl, A., Burkett, Z. D., Kohlway,
938 A. S., Lebofsky, R., et al. (2022). Functional inference of gene regulation using single-cell
939 multi-omics. *Cell Genom* 2, 100166.

940 Kaya-Okur, H. S., Wu, S. J., Codomo, C. A., Pledger, E. S., Bryson, T. D., Henikoff, J. G., Ahmad, K. and
941 Henikoff, S. (2019). CUT&Tag for efficient epigenomic profiling of small samples and single
942 cells. *Nat Commun* 10, 1930.

943 Kim, J., T. Jakobsen, S., Natarajan, K. N. and Won, K.-J. (2021). TENET: gene network reconstruction
944 using transfer entropy reveals key regulatory factors from single cell transcriptomic data.
945 *Nucleic Acids Research* 49, e1–e1.

946 Krishnaswamy, S., Spitzer, M. H., Mingueneau, M., Bendall, S. C., Litvin, O., Stone, E., Pe'er, D. and
947 Nolan, G. P. (2014). Conditional density-based analysis of T cell signaling in single-cell data.
948 *Science* 346, 1250689.

949 Li, L., Guo, F., Gao, Y., Ren, Y., Yuan, P., Yan, L., Li, R., Lian, Y., Li, J., Hu, B., et al. (2018). Single-cell
950 multi-omics sequencing of human early embryos. *Nature Cell Biology* 20,.

951 Li, Z., Schulz, M. H., Look, T., Begemann, M., Zenke, M. and Costa, I. G. (2019). Identification of
952 transcription factor binding sites using ATAC-seq. *Genome Biology* 20,.

953 Lim, C. Y., Tam, W.-L., Zhang, J., Ang, H. S., Jia, H., Lipovich, L., Ng, H.-H., Wei, C.-L., Sung, W. K.,
954 Robson, P., et al. (2008). Sall4 Regulates Distinct Transcription Circuitries in Different
955 Blastocyst-Derived Stem Cell Lineages. *Cell Stem Cell* 3, 543–554.

956 Lin, S.-C. J., Wani, M. A., Whitsett, J. A. and Wells, J. M. (2010). Klf5 regulates lineage formation in
957 the pre-implantation mouse embryo. *Development* 137, 3953–3963.

958 Liu, L., Leng, L., Liu, C., Lu, C., Yuan, Y., Wu, L., Gong, F., Zhang, S., Wei, X., Wang, M., et al. (2019). An
959 integrated chromatin accessibility and transcriptome landscape of human pre-implantation
960 embryos. *Nature Communications* 10, 364.

961 Lv, B., An, Q., Zeng, Q., Zhang, X., Lu, P., Wang, Y., Zhu, X., Ji, Y., Fan, G. and Xue, Z. (2019). Single-cell
962 RNA sequencing reveals regulatory mechanism for trophoblast cell-fate divergence in human
963 peri-implantation conceptuses. *PLoS Biol* 17, e3000187.

964 Lyu, X., Rowley, M. J. and Corces, V. G. (2018). Architectural Proteins and Pluripotency Factors
965 Cooperate to Orchestrate the Transcriptional Response of hESCs to Temperature Stress.
966 *Molecular Cell* 71, 940-955.e7.

967 Marbach, D., Costello, J. C., Küffner, R., Vega, N. M., Prill, R. J., Camacho, D. M., Allison, K. R., The
968 DREAM5 Consortium, Kellis, M., Collins, J. J., et al. (2012). Wisdom of crowds for robust gene
969 network inference. *Nature Methods* 9, 796–804.

970 Margolin, A. A., Nemenman, I., Basso, K., Wiggins, C., Stolovitzky, G., Favera, R. D. and Califano, A.
971 (2006). ARACNE: An Algorithm for the Reconstruction of Gene Regulatory Networks in a
972 Mammalian Cellular Context. *BMC Bioinformatics* 7, S7.

973 Materna, S. C. and Davidson, E. H. (2007). Logic of gene regulatory networks. *Curr Opin Biotechnol.* 4,
974 351–354.

975 Meers, M. P., Bryson, T. D., Henikoff, J. G. and Henikoff, S. (2019). Improved CUT&RUN chromatin
976 profiling tools. *eLife* 8,.

977 Meistermann, D., Bruneau, A., Loubersac, S., Reignier, A., Firmin, J., François-Campion, V., Kilens, S.,
978 Lelièvre, Y., Lammers, J., Feyeux, M., et al. (2021). Integrated pseudotime analysis of human
979 pre-implantation embryo single-cell transcriptomes reveals the dynamics of lineage
980 specification. *Cell Stem Cell* 28, 1625-1640.e6.

981 Mitsui, K., Tokuzawa, Y., Itoh, H., Segawa, K., Murakami, M., Takahashi, K., Maruyama, M., Maeda,
982 M. and Yamanaka, S. (2003). The homeoprotein Nanog is required for maintenance of
983 pluripotency in mouse epiblast and ES cells. *Cell* 113, 631–642.

984 Mochida, K., Koda, S., Inoue, K. and Nishii, R. (2018). Statistical and Machine Learning Approaches to
985 Predict Gene Regulatory Networks From Transcriptome Datasets. *Front Plant Sci* 9, 1770.

986 Morris, S. A., Cahan, P., Li, H., Zhao, A. M., San Roman, A. K., Shivdasani, R. A., Collins, J. J. and Daley,
987 G. Q. (2014). Dissecting engineered cell types and enhancing cell fate conversion via CellNet.
988 *Cell* 158, 889–902.

989 Nguyen, H., Tran, D., Tran, B., Pehlivan, B. and Nguyen, T. (2021). A comprehensive survey of
990 regulatory network inference methods using single cell RNA sequencing data. *Briefings in*
991 *Bioinformatics* 22, bbaa190.

992 Niakan, K.K. and Eggan K. (2012). Analysis of human embryos from zygote to blastocyst reveals
993 distinct gene expression patterns relative to the mouse. *Developmental Biology* 375, 54-64.
994

995 Niakan, K. K., Han, J., Pedersen, R. A., Simon, C. and Pera, R. A. R. (2012). Human pre-implantation
996 embryo development. *Development* 139, 829–841.

997 Pastor, W. A., Liu, W., Chen, D., Ho, J., Kim, R., Hunt, T. J., Lukianchikov, A., Liu, X., Polo, J. M.,
998 Jacobsen, S. E., et al. (2018). TFAP2C regulates transcription in human naive pluripotency by
999 opening enhancers. *Nat Cell Biol* 20, 553–564.

1000 Pe’er, D. (2005). Bayesian Network Analysis of Signaling Networks: A Primer. *Science Signaling* 2005,
1001 p14–p14.

1002 Peter, I. S. and Davidson, E. H. (2011). Evolution of Gene Regulatory Networks Controlling Body Plan
1003 Development. *Cell* 144, 970–985.

1004 Petropoulos, S., Edsgård, D., Reinius, B., Deng, Q., Panula, S. P., Codeluppi, S., Plaza Reyes, A.,
1005 Linnarsson, S., Sandberg, R. and Lanner, F. (2016). Single-Cell RNA-Seq Reveals Lineage and X
1006 Chromosome Dynamics in Human Preimplantation Embryos. *Cell* 165, 1012–1026.

1007 Picelli, S., Faridani, O. R., Björklund, A. K., Winberg, G., Sagasser, S. and Sandberg, R. (2014). Full-
1008 length RNA-seq from single cells using Smart-seq2. *Nat Protoc* 9, 171–181.

1009 Pratapa, A., Jaliha, A. P., Law, J. N., Bharadwaj, A. and Murali, T. M. (2020). Benchmarking
1010 algorithms for gene regulatory network inference from single-cell transcriptomic data.
1011 *Nature Methods* 17, 147–154.

1012 Roode, M., Blair K., Snell, P., Elder, K., Marchant, S., Smith, A., Nichols, J. (2012). Human hypoblast
1013 formation is not dependent on FGF signalling. *Developmental Biology* 361, 358-363.

1014

1015 Schneider, C. A., Rasband, W. S. and Eliceiri, K. W. (2012). NIH Image to ImageJ: 25 years of image
1016 analysis. *Nat Methods* 9, 671–675.

1017 Schoenfelder, S. and Fraser, P. (2019). Long-range enhancer-promoter contacts in gene expression
1018 control. *Nat Rev Genet* 20, 437–455.

1019 Shi, W., Wang, H., Pan, G., Geng, Y., Guo, Y. and Pei, D. (2006). Regulation of the Pluripotency
1020 Marker Rex-1 by Nanog and Sox2. *Journal of Biological Chemistry* 281, 23319–23325.

1021 Shu, H., Zhou, J., Lian, Q., Li, H., Zhao, D., Zeng, J. and Ma, J. (2021). Modeling gene regulatory
1022 networks using neural network architectures. *Nature Computational Science* 1, 491–501.

1023 Skene, P. J. and Henikoff, S. (2017). An efficient targeted nuclease strategy for high-resolution
1024 mapping of DNA binding sites. *eLife* 6, e21856.

1025 Stirling, D. R., Swain-Bowden, M. J., Lucas, A. M., Carpenter, A. E., Cimini, B. A. and Goodman, A.
1026 (2021). CellProfiler 4: improvements in speed, utility and usability. *BMC Bioinformatics* 22,
1027 433.

1028 Stone, M., Li, J., McCalla, S. G., Siahpirani, A. F., Periyasamy, V., Shin, J. and Roy, S. (2021). Identifying
1029 strengths and weaknesses of methods for computational network inference from single cell
1030 RNA-seq data. *bioRxiv* 2021.06.01.446671,.

1031 Strumpf, D., Mao, C.-A., Yamanaka, Y., Ralston, A., Chawengsaksophak, K., Beck, F. and Rossant, J.
1032 (2005). Cdx2 is required for correct cell fate specification and differentiation of
1033 trophoblast in the mouse blastocyst. *Development* 132, 2093–2102.

1034 Thompson, D., Regev, A. and Roy, S. (2015). Comparative analysis of gene regulatory networks: from
1035 network reconstruction to evolution. *Annu Rev Cell Dev Biol* 31, 399–428.

1036 Tibshirani, R. (1996). Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal
1037 Statistical Society: Series B (Methodological)* 58, 267–268.

1038 Tipping, M. E. (2001). Sparse Bayesian Learning and the Relevance Vector Machine. *Journal of
1039 Machine Learning Research* 1, 211–244.

1040 Wang, B., Pourshafeie, A., Zitnik, M., Zhu, J., Bustamante, C. D., Batzoglou, S. and Leskovec, J. (2018).
1041 Network enhancement as a general method to denoise weighted biological networks.
1042 *Nature Communications* 9,.

1043 Wang, J., Ma, A., Ma, Q., Xu, D. and Joshi, T. (2020). Inductive inference of gene regulatory network
1044 using supervised and semi-supervised graph neural networks. *Computational and Structural
1045 Biotechnology Journal* 18, 3335–3343.

1046 Wang, X., He, Y., Zhang, Q., Ren, X. and Zhang, Z. (2021). Direct Comparative Analyses of 10X
1047 Genomics Chromium and Smart-seq2. *Genomics Proteomics Bioinformatics* 19, 253–266.

1048 Weigert, M., Schmidt, U., Haase, R., Sugawara, K. and Myers, G. (2020). Star-convex Polyhedra for 3D
1049 Object Detection and Segmentation in Microscopy. In *2020 IEEE Winter Conference on
1050 Applications of Computer Vision (WACV)*, pp. 3655–3662. Snowmass Village, CO, USA: IEEE.

1051 Yan, L., Yang, M., Guo, H., Yang, L., Wu, J., Li, R., Liu, P., Lian, Y., Zheng, X., Yan, J., et al. (2013).
1052 Single-cell RNA-Seq profiling of human preimplantation embryos and embryonic stem cells.
1053 *Nature Structural & Molecular Biology* 20, 1131–1139.

1054 Chang S., Pyne S., Pietrzak S., Halbert S., McCalla, S.G., Siahpirani A.F., Sridharan, R. and Roy S.
1055 (2023). Inference of cell type-specific gene regulatory networks on cell lineages from single cell
1056 omics datasets. *Nat Commun* 14, 3064

1057 Zheng, G. X. Y., Terry, J. M., Belgrader, P., Ryvkin, P., Bent, Z. W., Wilson, R., Ziraldo, S. B., Wheeler, T.
 1058 D., McDermott, G. P., Zhu, J., et al. (2017). Massively parallel digital transcriptional profiling
 1059 of single cells. *Nat Commun* 8, 14049.
 1060 Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *J Royal*
 1061 *Statistical Soc B* 67, 301–320.
 1062
 1063
 1064
 1065
 1066

1067 **Figure legends**
 1068

1069 **Figure 1. An overview of the regulatory network inference, evaluation, and validation framework and inference on simulated scRNA-seq data.**

1070 **(A)** Footprinting analysis was applied to ATAC-seq data from human embryos at the late blastocyst stage to identify
 1071 potential transcription factor (TF) binding sites and TF-gene regulations. The TFs and genes
 1072 in this list of chromatin accessibility (CA)-predicted regulations were used to refine the size
 1073 of scRNA-seq data from each of the cell types in embryos where cells were collected at the
 1074 same developmental stage (EPI: Epiblast, 26 cells; PE: Primitive Endoderm, 33 cells; TE:
 1075 Trophectoderm, 45 cells). Four gene regulatory network (GRN) inference approaches and a
 1076 random approach were applied to the scRNA-seq data. The scRNA-seq data were
 1077 normalised using four different methods for comparison ($\log(\text{TPM}+1)$, $\log(\text{FPKM}+1)$, batch-
 1078 corrected and log-counts). The reproducibility and biological context of the predicted GRNs
 1079 were evaluated using several statistical tests. GRNs refined by keeping CA-predicted
 1080 regulations only were also evaluated. Squares represent genes that regulate other genes,
 1081 and circles represent genes that are regulated by other genes. **(B)** Comparison of different
 1082 GRN inference methods (L0L2, GENIE3, Spearman correlation and MI) to recover the
 1083 ground-truth GRN structure from simulated gene expression data as measured with the
 1084 average area under the receiver operating characteristic curve (AUROC) from 100
 1085 simulations. The average AUROC is shown as a function of different sample sizes ($n=10$ to
 1086 1000) and the number of potential regulators of each gene in the simulated datasets is also
 1087 varied from 10 to 100. The range of sizes of the three human blastocyst datasets that we
 1088 analysed is highlighted in grey as a reference. Error bars correspond to standard errors of
 1089 the mean. **(C)** Box-and-whisker plots show comparison of AUROC values for all simulations
 1090 for each method with number of samples $n=100$ and number of regulators $p=20, 50$ and 100.
 1091
 1092

1093 **Figure 2. Statistical comparison of four different methods to infer early human embryo GRNs.**

1094 Robustness of the GRNs predicted by four inference methods with or without
 1095 chromatin accessibility (CA) refinement was evaluated by three different metrics. **(A)** ΔR
 1096 measures the difference between the median edge-level reproducibility for a GRN inference
 1097 method and a random predictor. **(B)** $\Delta \text{AU} \text{PRC}$ quantifies the extent to which the interactions
 1098 inferred from half of a dataset coincide with the top-1% interactions inferred from the second
 1099 half. This 2-fold cross-validation experiment was repeated ten times and compared to a
 1100 random predictor. **(C)** The inverse of the normalised L2 loss between the network scores from
 1101 the two folds was computed. L0L2 did not converge when applied to the TE dataset. All data
 1102 used was generated using the SMART-seq2 sequencing technique.
 1103

1104 **Figure 3. Evaluation of the biological relevance of inferred GRNs in the context of early human development.**

1105 **(A)** Two-dimensional UMAP representation of the scRNA-seq data from each of the three cell types in the early human at the blastocyst stage. The expression
 1106 of marker genes for each cell type is also represented in UMAP space: NANOG is used as an
 1107 Epiblast marker, CDX2 as a Trophectoderm marker and GATA4 as a Primitive endoderm
 1108 marker. **(B)** Marker activity in the cell-type specific GRNs (EPI: Epiblast, PE: Primitive
 1109 endoderm; TE: Trophectoderm) and the network common to the three cell types in the human
 1110

1111 blastocyst. The expected pattern, i.e. the marker that was expected to be active in each GRN,
1112 is shown for reference. **(C)** Normalised V-statistic across datasets and inference methods. V
1113 corresponds to the number of gene sets that were enriched at a significance level of 10% in a
1114 gene set enrichment analysis performed with the genes involved in the top-500 regulations by
1115 the 25 most-connected TFs in the EPI or TE GRNs. All sequencing was performed using the
1116 SMART-seq2 technique.

1117
1118 **Figure 4. TF GRNs predicted by MICA to be regulated by, or to regulate, NANOG, GATA4**
1119 **or CDX2.** Target networks contain a maximum of 25 potential target TFs of the hub (or central)
1120 TF, whereas regulator networks include a maximum of 25 TFs that potentially regulate the
1121 hub. The average expression of the transcript is represented by the size of the TF node. MI
1122 scores are represented by the thickness of the line in the predicted network. To further refine
1123 the MICA-predictions we used the Spearman's rank correlation coefficient between the
1124 expression levels of the source and target nodes to define correlated or anti-correlated
1125 expression. Correlated or anti-correlated node pairs correspond to positive or negative
1126 Spearman's rank coefficient with p-value smaller than 0.1. Node pairs with a p-value equal to
1127 or larger than 0.1 were defined as uncorrelated, though they are predicted by MICA. **(A)**
1128 Network of TFs predicted to be targeted by NANOG. **(B)** Network of TFs that are predicted to
1129 regulate NANOG. **(C)** Network of TFs predicted to be targeted by GATA4. **(D)** Network of TFs
1130 that that are predicted to regulate GATA4. **(E)** Network TFs predicted to be targeted by CDX2.
1131 **(F)** Network of TFs that are predicted to regulate CDX2.

1132
1133
1134 **Figure 5. MICA predicted conserved TFAP2C interactions in human early embryo EPI,**
1135 **TE, and PE cells. (A)** Network composed of TFAP2C, JUND, SOX4, and GCM1 in EPI
1136 (left), TE (middle), and PE (right). **(B)** Correlation plots between TFAP2C, JUND, SOX4, and
1137 GCM1 in TE cells. Bottom left: log transformed RNA expression of genes in single cells.
1138 Diagonal: Distribution of log transformed RNA expression for TFAP2C, JUND, SOX4 and
1139 GCM1. Top right: Spearman correlation between TF pairs. (***: $p < 0.001$). **(C)**
1140 Immunofluorescence staining of TFAP2C and JUND in E6.5 human blastocysts. **(D)**
1141 Correlation plot of TFAP2C and JUND protein expression intensity in nuclei of human E6.5
1142 blastocysts. Numbers on the plot are Pearson correlation value. (***: $p < 0.001$).

1143
1144
1145
1146

1147 **Supplementary Figure legends**

1148

1149 **Figure S1. (A)** Regulated genes may not 'switch on' until the transcription factor expression
1150 level reaches a certain threshold. This non-linear regulatory effect may be detected more
1151 efficiently using mutual information-based methods, compared to linear models. **(B)**
1152 Distribution of enriched transcription factor (TF) motifs per gene promoter. The median
1153 values are highlighted. These data were computed using the motifs enriched in footprints
1154 found within open chromatin regions in the low-input ATAC-seq from the inner cell mass
1155 (ICM) and trophectoderm (TE). **(C)** UMAP analysis of scRNA-seq data early human
1156 blastocyst stage embryos shown in Figure 3A. Here cells were coloured according to the
1157 primary data from the respective publication indicated to demonstrate that batch effects have
1158 been minimised and the scRNA-seq data is integrated irrespective of the primary study.

1159

1160 **Figure S2.** R distribution evaluation of the GRN prediction methods with and without
1161 chromatin accessibility (CA) refinement. The reproducibility score R is the bootstrap estimate
1162 of the posterior probability of observing an edge E given the dataset D , $R=P(E|D)$, for the top
1163 100,000 predicted edges.

1164

1165 **Figure S3.** Repeated 2-fold cross-validation of the GRN prediction methods with and without
1166 chromatin accessibility (CA) refinement. The area under the precision-recall curve (AUPRC)
1167 quantifies the extent to which the interactions inferred from the first fold coincide with the top-
1168 1% interactions inferred from the second, reference fold. Cross-validation was repeated ten
1169 times. The LOL2 method didn't converge in the TE dataset.

1170
1171 **Figure S4.** Normalised L2 loss evaluation of the GRN prediction methods with and without
1172 chromatin accessibility (CA) refinement. As in Fig. S5, 2-fold cross-validation was employed
1173 ten times, comparing the scores obtained from the application of a GRN inference method to
1174 each data split with the L2 loss function. The resulting losses were divided by the product of
1175 the square roots of the L2 norms of the network scores of each of the 2 data folds.

1176
1177 **Figure S5.** The number of genes regulated by the markers of each one of the cell types in the
1178 human blastocyst (EPI: Epiblast, PE: Primitive Endoderm and TE: Trophectoderm) is used as
1179 a proxy for marker activity in the networks predicted by the different network inference
1180 methods evaluated in this work. The results in this figure were generated with the gene
1181 expression after TPM normalisation.

1182
1183 **Figure S6.** The number of genes regulated by the markers of each one of the cell types in the
1184 human blastocyst (EPI: Epiblast, PE: Primitive Endoderm and TE: Trophectoderm) is used as
1185 a proxy for marker activity in the networks predicted by the different network inference
1186 methods evaluated in this work. The results in this figure were generated with the gene
1187 expression after FPKM normalisation.

1188
1189 **Figure S7.** The number of genes regulated by the markers of each one of the cell types in the
1190 human blastocyst (EPI: Epiblast, PE: Primitive Endoderm and TE: Trophectoderm) is used as
1191 a proxy for marker activity in the networks predicted by the different network inference
1192 methods evaluated in this work. The results in this figure were generated with the gene
1193 expression after log-count normalisation.

1194
1195 **Figure S8.** The number of genes regulated by the markers of each one of the cell types in the
1196 human blastocyst (EPI: Epiblast, PE: Primitive Endoderm and TE: Trophectoderm) is used as
1197 a proxy for marker activity in the networks predicted by the different network inference
1198 methods evaluated in this work. The results in this figure were generated with the batch-
1199 corrected gene expression values.

1200
1201 **Figure S9.** The V-statistics represent the number of relevant gene-sets that were enriched at
1202 a significance level of 10% ($p < 0.1$) in a gene-set enrichment analysis performed with the genes
1203 regulated in the 25 top-predicted epiblast or trophectoderm interactions. in each of the
1204 networks predicted by the different network inference methods evaluated in this work.

1205
1206 **Figure S10. (A)** MACS2 binding score to ZNF343 (± 10 kb from TSS) from NANOG ChIP-seq
1207 studies in naive and primed hESC H9 cell lines. **(B)** Feedback network formed between CDX2
1208 targets and receptors.

1209
1210 **Figure S11.** Correlation plots between TFAP2C, JUND, SOX4 and GCM1 RNA expression
1211 in the **(A)** EPI or **(B)** PE lineage of human blastocysts.

1212
1213 **Figure S12.** Target networks and regulator networks of TFAP2C in TE, EPI and PE cells.

1214
1215 **Figure S13. (A)** Immunofluorescence staining of TFAP2C and JUND in E6.5 human
1216 blastocysts. **(B)** scRNA-seq expression plots of TFAP2C versus ZFP42 (left), FOXO3
1217 (middle), or JUND (right).