# Molecular Recording of Mammalian Embryogenesis

## Citation

## Published Version

## Permanent link

## Terms of Use

# Share Your Story

# Molecular recording of mammalian embryogenesis

Michelle M. Chan[1,2]*, Zachary D. Smith[4,5,6]*, Stefanie Grosswendt[7], Helene Kretzmer[7], Thomas Norman[1,2], Britt Adamson[1,2], Marco Jost[1,2,3], Jeffrey J. Quinn[1,2], Dian Yang[1,2], Matthew G. Jones[1,2,8], Alex Khodaverdian[9,10], Nir Yosef[9,10,11,12], Alexander Meissner[4,5,7], Jonathan S. Weissman[1,2]


[1]Department of Cellular and Molecular Pharmacology, University of California, San Francisco, San Francisco, CA, USA

[2]Howard Hughes Medical Institute, University of California, San Francisco, San Francisco, CA, USA

[3]Department of Microbiology and Immunology, University of California, San Francisco, San Francisco, CA, USA

[4]Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

[5]Department of Stem Cell and Regenerative Biology, Harvard University, Cambridge, Massachusetts, USA

[6]Department of Molecular and Cellular Biology, Harvard University, Cambridge, Massachusetts, USA

[7]Department of Genome Regulation, Max Planck Institute for Molecular Genetics, Berlin 14195, Germany

[8]Integrative Program in Quantitative Biology, University of California, San Francisco, San Francisco, California, USA

[9]Department of Electrical Engineering and Computer Science, University of California, Berkeley, Berkeley, California, USA

24    [10]Center for Computational Biology, Berkeley, California, USA

25    [11]Chan Zuckerberg Biohub, San Francisco, California, USA

26    [12]Ragon Institute of Massachusetts General Hospital, MIT and Harvard, Cambridge,

27    Massachusetts, USA

28

29    *These authors contributed equally to this work

30    Correspondence: jonathan.weissman@ucsf.edu (J.S.W.), meissner@molgen.mpg.de (A.M.)

31

32    **Ontogeny describes the emergence of complex multicellular organisms from single**

33    **totipotent cells. In mammals, this field is particularly challenging due to the indeterminate**

34    **relationship between self-renewal and differentiation, variation of progenitor field sizes,**

35    **and internal gestation. Here, we present a flexible, high information, multi-channel**

36    **molecular recorder with a single cell (sc) readout and apply it as an evolving lineage tracer**

37    **to define a mouse cell fate map from fertilization through gastrulation. By combining**

38    **lineage information with scRNA-seq profiles, we recapitulate canonical developmental**

39    **relationships between different tissue types and reveal the nearly complete transcriptional**

40    **convergence of endodermal cells from extra-embryonic and embryonic origins. Finally, we**

41    **apply our cell fate map to estimate the number of embryonic progenitor cells and their**

42    **degree of asymmetric partitioning during specification. Our approach enables massively**

43    **parallel, high-resolution recording of lineage and other information in mammalian systems**

44    **to facilitate a quantitative framework for understanding developmental processes.**

45

46    Development of a multicellular organism from a single cell is an astonishing process.

47    Classic lineage tracing experiments using *C. elegans* revealed surprising outcomes, including

48    deviations between lineage and functional phenotype, but nonetheless benefited from the highly

49    deterministic nature of this organism's development[1].  Alternatively, more complex species

50    generate larger, more elaborate structures that progress through multiple transitions, raising

51    questions regarding the coordination between specification and commitment to ensure faithful

52    recapitulation of an exact body plan[2,3]. Single cell RNA-sequencing (scRNA-seq) has permitted

53    unprecedented explorations into cell type heterogeneity, producing profiles of developing

54    flatworms[4,5], frogs[6], zebrafish[7,8], and mice[9,10].  More recently, CRISPR-Cas9-based technologies

55    have been applied to record cell lineage[11-13], and combined with scRNA-seq to generate fate

56    maps in zebrafish[14-16].  However, these technologies include only one or two bursts of barcode

57    diversity generation, which may be limiting for other applications or organisms.

58    An ideal molecular recorder for these questions would possess the following

59    characteristics: 1) minimal impact on cellular phenotype; 2) high information content to account

60    for hundreds of thousands of cells; 3) a single cell readout for simultaneous profiling of

61    functional state[14-16]; 4) flexible recording rates that can be tuned to a broad temporal range; and

62    5) continuous generation of diversity throughout the experiment.  The last point is especially

63    relevant for mammalian development, where spatial plans are gradually and continuously

64    specified and may originate from small, transient progenitor fields. Moreover, scRNA-seq has

65    revealed populations of cells with a continuous spectrum of phenotypes, implying that

66    differentiation does not occur instantaneously, further motivating the need for an evolving

67    recorder[17].

68  Here, we generated and validated a method for simultaneously reporting cellular state and

69  lineage history in mice. Our CRISPR-Cas9-based recorder is capable of high information content

70  and multi-channel recording with readily tunable mutation rates. We employ the recorder as a

71  continuously evolving lineage tracer to observe the fate map underlying embryogenesis through

72  gastrulation, recapitulating canonical paradigms and illustrating how lineage information may

73  facilitate the identification of novel cell types.

74

75

76    **Results**

77    **A transcribed, multi-channel, and continuously evolving molecular recorder**

78    To achieve our goal of a tunable, high information content molecular recorder, we

79    utilized Cas9 to generate insertions or deletions (indels) upon repair of double-stranded breaks,

80    which are inherited in the next generation of cells[11-16].  We record within a 205 base pair,

81    synthetic DNA "target site" containing three "cut sites" and a static 8 base pair "integration

82    barcode" (intBC), which are delivered in multiple copies via piggyBac transposition (**Fig. 1a, b**).

83    We embedded this sequence into the 3'UTR of a constitutively transcribed fluorescent protein to

84    enable profiling from the transcriptome.  A second cassette encodes three independently

85    transcribed and complementary guide RNAs to permit recording of multiple, distinct signals

86    (**Fig. 1a, b**)[18].

87    Our system is capable of high information storage due to the diversity of heritable repair

88    outcomes, and the large number of targeted sites, which can be distinguished by the intBC (**Fig.**

89    **1c**).  DNA repair generates hundreds of unique indels, and the distribution for each cut site is

90    different and nonuniform: some produce highly biased outcomes while others create a diverse

91    series (**Fig. 1c**, **Extended Data Fig. 1**)[19-21]. To identify sequences that can tune the mutation rate

92    of our recorder for timescales that are not pre-defined, and may extend from days to months, we

93    screened several guide RNA series containing mismatches to their targets[22] by monitoring their

94    activity on a GFP reporter over a 20-day timecourse and selected those that demonstrated a broad

95    dynamic range  (**Fig. 1d**). Slower cutting rates may improve viability *in vivo*, as frequent Cas9-

96    mediated double-strand breaks can cause cellular toxicity[23,24]. To demonstrate information

97    recovery from single cell transcriptomes, we stably transduced K562 cells with our technology

98    and generated a primary, cell-barcoded cDNA pool via the 10x Genomics platform, allowing us

99    to assess global transcriptomes and specifically amplify mutated target sites (**Extended Data**

100    **Fig. 1c**).

101

102    **Tracing cell lineages in mouse development**

103    We next applied our technology to map cell fates during mouse early development from

104    totipotency onwards. We integrated multiple target sites into the genome, delivered constitutive

105    Cas9-GFP encoding sperm into oocytes to initiate cutting, and isolated embryos for analysis at

106    ~embryonic day (E)8.5 or E9.5 (**Fig. 2a**, **Methods**). To confirm our lineage tracing capability,

107    we amplified the target site from bulk placenta, yolk sac, and three embryonic fractions from an

108    E9.5 embryo and recapitulated their expected relationships using the similarity of their indel

109    proportions (**Fig. 2b, Extended Data Figure 2**).

110    Following this *in vivo* proof of principle, we generated single cell data from additional

111    embryos (**Extended Data Figure 3**). We collected scRNA-seq data for 7,364 – 12,990 cells

112    from 7 embryos (~15.8% – 61.4% of the total cell count) and recovered 167 – 2,461 unique

113    lineage identities (≥1 target site recovered for 15% – 75% of cells from 3 to 15 intBCs, **Fig. 2c,**

114    **Extended Data Figure 4**). Many target sites are either lowly or heterogeneously represented,

115    which we improved by changing the promoter from a truncated form of Ef1α to an intron-

116    containing version (see embryo 7, **Extended Data Figure 4**)[25].

117    We estimated the indel likelihood distribution by combining data from all seven embryos.

118    Many indels are shared with K562 cells, though their likelihoods differ, suggesting that cell type

119    or developmental status may influence repair outcomes (**Fig. 2d**, **Extended Data Figure 1, 4f**)[19].

120    Our ability to independently measure and control the rate of cutting across the target site is

121    preserved *in vivo*, with minimal interference between cut sites except when using combinations

122    of the fastest guides that may lead to end-joining between simultaneous double strand breaks

123    (**Fig. 2e**). The fastest cutters result in higher proportions of cells with identical indels, indicating

124    earlier mutations in development, which correspondingly reduce indel diversity (**Fig. 2f, g**).

125    Importantly, the lineage tracer retains additional recording capacity beyond the temporal interval

126    studied here, as most embryos still have unmodified cut sites (**Fig. 2f**).

127

128    **Assigning cellular states by simultaneous scRNA-seq**

129    Next, to ascertain cell function, we utilized annotations from a compendium of wild-type

130    mouse gastrulation (E6.5 – E8.5). We assigned cells from lineage-traced embryos by their

131    proximity to each cell state expression signature and aged each embryo by their tissue

132    proportions compared to each stage (**Fig. 3a-c**)[26]. We proceeded with six of our seven embryos,

133    as they appeared to be morphologically normal and included every expected tissue type: two

134    mapped most closely to E8.5, and the remaining four mapped to E8.0 (**Extended Data Fig. 5**).

135    Placenta was not specifically isolated, but is present in four of six embryos, serving as a valuable

136    outgroup to establish our ability to track transitions to the earliest bifurcation.

137    We also developed breeder mice that would enable facile exploration of all stages of

138    development by injecting target sites into Cas9 negative backgrounds. This approach

139    substantially increases the number of stably integrated target sites (~20). Resulting mice can be

140    crossed with Cas9 expressing strains to yield viable Cas9$^+$ F1 litters that maintain continuous,

141    stochastic indel generation into adulthood, demonstrating that cutting does not noticeably

142    interfere with normal animal development (**Extended Data Fig 6**).

143    **Single cell lineage reconstruction of mouse embryogenesis**

144        We developed phylogenetic reconstruction strategies to specifically exploit the

145    characteristics of our lineage tracer, namely categorical indels, irreversibility of mutations, and

146    presence of missing values (**Extended Data Figure 7**, **Methods**).  We determined the best

147    reconstruction by summing the log-likelihoods for all indels that appear in the tree using

148    likelihoods estimated from embryo data (**Extended Data Figures 4 and 7**). When cell type

149    identity from scRNA-seq is overlaid onto the tree, we observe functional restriction during

150    development, with fewer cell types represented as we move from root to leaves (**Fig. 4a, b,**

151    **Extended Data Figure 8**).

152        scRNA-seq-based strategies for ordering cells, such as trajectory inference, typically

153    assume that functional similarity reflects close lineage[17].  To investigate this question directly,

154    we used a modified Hamming distance to measure pairwise lineage distance and compared them

155    to RNA-seq correlation.  Generally, cells separated by a smaller lineage distance have more

156    similar transcriptional profiles, though this relationship is clearer for some embryos than others

157    (**Fig. 4c**, **Extended Data Figure 9**). This result is consistent with the notion of continuous

158    restriction of potency as cells differentiate into progressively differentiated types.

159        We also developed a shared progenitor score that estimates the degree of common

160    ancestry between different tissues by evaluating the number and specificity of shared nodes in

161    the tree (**Methods**).  Despite the stochastic timing of indel formation, this approach can

162    reproducibly recover emergent tissue relationships, such as possible shared origins between

163    anterior somites and paraxial mesoderm or neuromesodermal progenitors and the future spinal

164    cord (**Fig. 4d**).  The full map of shared progenitor scores can be clustered to create a

165    comprehensive picture of tissue relationships during development (**Extended Data Fig. 8d**).

166

167 **Transcriptional state and developmental origin do not always correspond**

168      While our reconstructed tissue relationships generally recapitulate canonical knowledge,

169 extra-embryonic and embryonic endoderm display consistent and unexpectedly close ancestry

170 despite their independent origins from the hypoblast and embryo-restricted epiblast (**Fig. 5a,**

171 **Extended Data Figure 9**). Manual inspection of the trees revealed a subpopulation of cells that

172 appear transcriptionally as embryonic endoderm but that lineage analysis places within extra-

173 embryonic branches (**Fig 4c, blue**). Consistent with this finding, an earlier, targeted study using

174 marker-directed lineage tracing identified latent extra-embryonic contribution to the developing

175 hindgut during gastrulation, although it was not possible to broadly evaluate their

176 transcriptomes[27].

177      Here, scRNA-seq profiles collected in tandem with the lineage readout allow us to assess

178 the degree of convergence towards a functional endoderm signature and identify distinguishing

179 genes. Endoderm-classified cells derived from extra-embryonic origin are most similar to the

180 endoderm cell type, but do share slightly higher similarity with yolk sac that is not apparent

181 within the t-sne projection of the full embryo (**Fig. 5b, Extended Data Figure 10**). Given these

182 independent origins, we might expect a subtle, but persistent, transcriptional signature reflecting

183 their developmental history. Strikingly, when we separate endoderm cells according to their

184 lineage, we identify two X-linked genes, Trap1a and Rhox5, general markers for extra-

185 embryonic tissue[28,29] that are consistently upregulated in the extra-embryonic origin endoderm

186 across embryos (K–S test, Bonferroni corrected $P$-value <0.05, **Fig. 5d, e**). Notably, in other

187 RNA-seq studies, these relationships are not captured by whole embryo clustering, and are only

188 found by specific examination of the hindgut (**Extended Data Figure 10**) [9,30]. These

189    observations confirm that our lineage tracer can successfully pinpoint instances of convergent

190    transcriptional regulation.

191

192    **Towards a quantitative fate map**

193        Simultaneous single cell lineage tracing with phenotype provides the unique opportunity

194    to infer the cellular potency and specification biases of ancestral cells as reconstructed by our

195    fate map[31,32]. Each node within the tree represents a unique lineage identity stemming from a

196    single reconstructed progenitor cell, allowing us to estimate lower boundaries of their field size

197    (**Methods**). We investigated the founding number of progenitors during the earliest transitions in

198    cellular potential.  We defined totipotency as a node that gives rise to both embryonic and extra-

199    embryonic ectodermal/placental cell types and tiered pluripotency into "early" and "late"

200    according to the presence of extra-embryonic endoderm (**Fig. 6a**)[33].  The contributions of these

201    founders to extant lineages are asymmetric, suggesting that even though a progenitor may be

202    biased towards a specific fate, it retains the ability to generate other cell types.  Lower bound

203    estimates from our data suggest a range of 1–6 totipotent cells, 10–20 early, and 18–51 late

204    pluripotent progenitors (**Fig. 6b**).  The variable number of multipotent cells at these stages may

205    reflect an encoded robustness that ensures successful assembly of the functioning organism,

206    particularly given that a single pluripotent cell can generate all somatic lineages in an embryo[34].

207    Future studies using more replicates generated by breeding may enable statistical approaches to

208    evaluate these organism-scale developmental considerations.

209

210    **Discussion**

211       In this study, we present cell fate maps underlying mammalian gastrulation using a

212    technology for high information and continuous recording. Several key ideas have emerged,

213    including the transformative nature of CRISPR-Cas9-directed mutation with a single cell RNA-

214    seq readout[14-16], how information about a cell's history recorded by this technology can

215    complement RNA-seq profiles to characterize cell type, and an early framework for

216    quantitatively understanding stochastic transitions during mammalian development.

217       The modularity of our recorder allows for substitutions that will increase its breadth of

218    applications. Here, we use three constitutively expressed guide RNAs to record continuously

219    over time, but future modifications could employ environmentally-responsive promoters that

220    sense stress, neuronal action potentials, or cell-to-cell contacts[35], or combine these approaches

221    for multifactorial recording. Similarly, Cas9-derived base editors[36], including those that create

222    diverse mutations[37] could allow for content-recording in cells that are particularly sensitive to

223    nuclease-directed DNA double strand breaks[23,24].

224       Our cell fate map identifies phenotypic convergence of independent cell lineages,

225    showcasing the power of unbiased organism-wide lineage tracing to separate populations that

226    appear similar in scRNA-seq alone. Specifically, we substantiate the extra-embryonic origin of a

227    subset of cells that resemble embryonic endoderm. While the initial specification of these

228    lineages are known to rely on redundant regulatory programs, they are temporally separated by

229    several days, emerge from transcriptionally and epigenetically distinct progenitors, and form

230    terminal cell types with highly divergent functions. The identification of highly predictive

231    markers that segregate by origin, such as Trap1a, provides a clear outline for further exploration

232    through spatial transcriptomics[38,39,40]. More generally, our approach can be used to investigate

233    other convergent processes or to discriminate heterogeneous cell states that represent persistent

234  signatures of a cell's past, which will be critical for the assembly of a comprehensive cell atlas[41].

235  The scope of transdifferentiation within mammalian ontogenesis remains largely unexplored, but

236  can be practically inventoried using our system.

237      Ultimately, our technology is designed to quantitatively address previously opaque

238  questions in ontogenesis.  Higher order issues of organismal regulation, such as the location,

239  timing, and stringency of developmental bottlenecks, as well as the corresponding likelihoods of

240  state transitions to different cellular phenotypes, can be modeled from the assembly of historical

241  relationships. Our hope is that characterization of these attributes will lead to new insights that

242  connect large-scale developmental phenomena to the molecular regulation of cell fate decision-

243  making.

244

**References**

245 **References**

246 1.    Sulston, J. E., Schierenberg, E., White, J. G. & Thomson, J. N. The embryonic cell lineage

247       of the nematode Caenorhabditis elegans. *Developmental Biology* **100,** 64–119 (1983).

248 2.    Pijuan-Sala, B., Guibentif, C. & Göttgens, B. Single-cell transcriptional profiling: a

249       window into embryonic cell-type specification. *Nat. Rev. Mol. Cell Biol.* **19,** 399–412

250       (2018).

251 3.    Zernicka-Goetz, M. Patterning of the embryo: the first spatial decisions in the life of a

252       mouse. *Development* **129,** 815–829 (2002).

253 4.    Fincher, C. T., Wurtzel, O., de Hoog, T., Kravarik, K. M. & Reddien, P. W. Cell type

254       transcriptome atlas for the planarian Schmidtea mediterranea. *Science* **360,** eaaq1736

255       (2018).

256 5.    Plass, M. *et al.* Cell type atlas and lineage tree of a whole complex animal by single-cell

257       transcriptomics. *Science* **360,** eaaq1723 (2018).

258 6.    Briggs, J. A. *et al.* The dynamics of gene expression in vertebrate embryogenesis at

259       single-cell resolution. *Science* **360,** eaar5780 (2018).

260 7.    Farrell, J. A. *et al.* Single-cell reconstruction of developmental trajectories during

261       zebrafish embryogenesis. *Science* **360,** eaar3131 (2018).

262 8.    Wagner, D. E. *et al.* Single-cell mapping of gene expression landscapes and lineage in the

263       zebrafish embryo. *Science* **360,** 981–987 (2018).

264 9.    Ibarra-Soria, X. *et al.* Defining murine organogenesis at single-cell resolution reveals a

265       role for the leukotriene pathway in regulating blood progenitor formation. *Nat. Cell Biol.*

266       **20,** 127–134 (2018).

267 10.    Han, X. *et al.* Mapping the Mouse Cell Atlas by Microwell-Seq. *Cell* **172,** 1091–1107.e17

268       (2018).

269    11.    Perli, S. D., Cui, C. H. & Lu, T. K. Continuous genetic recording with self-targeting

270       CRISPR-Cas in human cells. *Science* **353,** aag0511–aag0511 (2016).

271    12.    Kalhor, R. *et al*. Developmental barcoding of whole mouse via homing CRISPR. *Science*

272       **361,** eaat9804 (2018).

273    13.    Frieda, K. L. *et al*. Synthetic recording and in situ readout of lineage information in single

274       cells. *Nature* **541,** 107–111 (2017).

275    14.    Raj, B. *et al*. Simultaneous single-cell profiling of lineages and cell types in the vertebrate

276       brain. *Nat*. *Biotechnol*. **36,** 442–450 (2018).

277    15.    Alemany, A., Florescu, M., Baron, C. S., Peterson-Maduro, J. & van Oudenaarden, A.

278       Whole-organism clone tracing using single-cell sequencing. *Nature* **556,** 108–112 (2018).

279    16.    Spanjaard, B. *et al*. Simultaneous lineage tracing and cell-type identification using

280       CRISPR-Cas9-induced genetic scars. *Nat*. *Biotechnol*. **36,** 469–473 (2018).

281    17.    Tanay, A. & Regev, A. Scaling single-cell genomics from phenomenology to mechanism.

282       *Nature* **541,** 331–338 (2017).

283    18.    Adamson, B. *et al*. A Multiplexed Single-Cell CRISPR Screening Platform Enables

284       Systematic Dissection of the Unfolded Protein Response. *Cell* **167,** 1867–1882.e21

285       (2016).

286    19.    van Overbeek, M. *et al*. DNA Repair Profiling Reveals Nonrandom Outcomes at Cas9-

287       Mediated Breaks. *Mol*. *Cell* **63,** 633–646 (2016).

288    20.    Schimmel, J., Kool, H., van Schendel, R. & Tijsterman, M. Mutational signatures of non-

289       homologous and polymerase theta-mediated end-joining in embryonic stem cells. *EMBO*

290       *J*. **36,** 3634–3649 (2017).

291  21.  Lemos, B. R. *et al*. CRISPR/Cas9 cleavages in budding yeast reveal templated insertions

292      and strand-specific insertion/deletion profiles. *Proc*. *Natl*. *Acad*. *Sci*. *U*.*S*.*A*. **115,** E2040–

293      E2047 (2018).

294  22.  Gilbert, L. A. *et al*. Genome-Scale CRISPR-Mediated Control of Gene Repression and

295      Activation. *Cell* **159,** 647–661 (2014).

296  23.  Ihry, R. J. *et al*. p53 inhibits CRISPR-Cas9 engineering in human pluripotent stem cells.

297      *Nat*. *Med*. **337,** 816 (2018).

298  24.  Haapaniemi, E., Botla, S., Persson, J., Schmierer, B. & Taipale, J. CRISPR-Cas9 genome

299      editing induces a p53-mediated DNA damage response. *Nat*. *Med*. **19,** 1 (2018).

300  25.  Kim, S.-Y., Lee, J.-H., Shin, H.-S., Kang, H.-J. & Kim, Y.-S. The human elongation

301      factor 1 alpha (EF-1 alpha) first intron highly enhances expression of foreign genes from

302      the murine cytomegalovirus promoter. *J*. *Biotechnol*. **93,** 183–187 (2002).

303  26.  Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell

304      transcriptomic data across different conditions, technologies, and species. *Nat*. *Biotechnol*.

305      **36,** 411–420 (2018).

306  27.  Kwon, G. S., Viotti, M. & Hadjantonakis, A.-K. The endoderm of the mouse embryo

307      arises by dynamic widespread intercalation of embryonic and extraembryonic lineages.

308      *Dev*. *Cell* **15,** 509–520 (2008).

309  28.  Eakin, G. S. & Hadjantonakis, A.-K. Sex-specific gene expression in preimplantation

310      mouse embryos. **7,** 205 (2006).

311  29.  Li, C.-S. *et al*. Trap1a is an X-linked and cell-intrinsic regulator of thymocyte

312      development. *Cell*. *Mol*. *Immunol*. **14,** 685–692 (2017).

313  30.  Pijuan-Sala, B. *et al*. A single-cell molecular map of mouse gastrulation and early

314          organogenesis. *Nature* **566,** 490–495 (2019).

315    31.    Soriano, P. & Jaenisch, R. Retroviruses as probes for mammalian development: allocation

316          of cells to the somatic and germ cell lineages. *Cell* **46,** 19–29 (1986).

317    32.    Jaenisch, R. Mammalian neural crest cells participate in normal embryonic development

318          on microinjection into post-implantation mouse embryos. *Nature* **318,** 181–183 (1985).

319    33.    Nichols, J. & Smith, A. Naive and primed pluripotent states. *Cell Stem Cell* **4,** 487–492

320          (2009).

321    34.    Wang, Z. & Jaenisch, R. At most three ES cells contribute to the somatic lineages of

322          chimeric mice and of mice produced by ES-tetraploid complementation. *Developmental*

323          *Biology* **275,** 192–201 (2004).

324    35.    Baeumler, T. A., Ahmed, A. A. & Fulga, T. A. Engineering Synthetic Signaling Pathways

325          with Programmable dCas9-Based Chimeric Receptors. *Cell Rep* **20,** 2639–2653 (2017).

326    36.    Komor, A. C., Kim, Y. B., Packer, M. S., Zuris, J. A. & Liu, D. R. Programmable editing

327          of a target base in genomic DNA without double-stranded DNA cleavage. *Nature* **533,**

328          420–424 (2016).

329    37.    Hess, G. T. *et al*. Directed evolution using dCas9-targeted somatic hypermutation in

330          mammalian cells. *Nat. Methods* **13,** 1036–1042 (2016).

331    38.    Hou, J. *et al*. A systematic screen for genes expressed in definitive endoderm by Serial

332          Analysis of Gene Expression (SAGE). *BMC Dev. Biol.* **7,** 92 (2007).

333    39.    Wang, G., Moffitt, J. R. & Zhuang, X. Multiplexed imaging of high-density libraries of

334          RNAs with MERFISH and expansion microscopy. *Sci Rep* **8,** 4847 (2018).

335    40.    Shah, S., Lubeck, E., Zhou, W. & Cai, L. In Situ Transcription Profiling of Single Cells

336          Reveals Spatial Organization of Cells in the Mouse Hippocampus. *Neuron* **92,** 342–357

337       (2016).

338    41.    Regev, A. *et al*. The Human Cell Atlas. *Elife* **6,** 503 (2017).

339    42.    Tzouanacou, E., Wegener, A., Wymeersch, F. J., Wilson, V. & Nicolas, J.-F. Redefining

340        the progression of lineage segregations during mammalian embryogenesis by clonal

341        analysis. *Dev. Cell* **17,** 365–376 (2009).

342

343    **Supplementary References**

344    43.    Schiml, S., Fauser, F. & Puchta, H. The CRISPR/Cas system can be used as nuclease for

345        in plantagene targeting and as paired nickases for directed mutagenesis in Arabidopsis

346        resulting in heritable progeny. *Plant J* **80,** 1139–1150 (2014).

347    44.    Ren, X. *et al*. Performance of the Cas9 Nickase System in Drosophila melanogaster. *G3* **4,**

348        1955–1962 (2014).

349    45.    Kimura, Y., Hisano, Y., Kawahara, A. & Higashijima, S.-I. Efficient generation of knock-

350        in transgenic zebrafish carrying reporter/driver genes by CRISPR/Cas9-mediated genome

351        engineering. *Sci Rep* **4,** 206–7 (2014).

352    46.    Dong, Z., Dong, X., Jia, W., Cao, S. & Zhao, Q. Improving the efficiency for generation

353        of genome-edited zebrafish by labeling primordial germ cells. *International Journal of*

354        *Biochemistry and Cell Biology* **55,** 329–334 (2014).

355    47.    Bao, Z. *et al*. Homology-Integrated CRISPR–Cas (HI-CRISPR) System for One-Step

356        Multigene Disruption in Saccharomyces cerevisiae. *ACS Synth. Biol.* **4,** 585–594 (2014).

357    48.    Zhou, H., Liu, B., Weeks, D. P., Spalding, M. H. & Yang, B. Large chromosomal

358        deletions and heritable small genetic changes induced by CRISPR/Cas9 in rice. *Nucleic*

359        *Acids Research* **42,** 10903–10914 (2014).

360 49. Jiang, W. *et al*. Demonstration of CRISPR/Cas9/sgRNA-mediated targeted gene

361    modification in Arabidopsis, tobacco, sorghum and rice. *Nucleic Acids Research* **41,**

362    e188–e188 (2013).

363 50. Feng, Z. *et al*. Efficient genome editing in plants using a CRISPR/Cas system. *Nature*

364    *Publishing Group* **23,** 1229–1232 (2013).

365 51. DiCarlo, J. E. *et al*. Genome engineering in Saccharomyces cerevisiae using CRISPR-Cas

366    systems. *Nucleic Acids Research* **41,** 4336–4343 (2013).

367 52. Jacobs, J. Z., Ciccaglione, K. M., Tournier, V. & Zaratiegui, M. Implementation of the

368    CRISPR-Cas9 system in fission yeast. *Nature Communications* **5,** 1–5 (1AD).

369 53. Wang, H. *et al*. One-step generation of mice carrying mutations in multiple genes by

370    CRISPR/Cas-mediated genome engineering. *Cell* **153,** 910–918 (2013).

371 54. Platt, R. J. *et al*. CRISPR-Cas9 knockin mice for genome editing and cancer modeling.

372    *Cell* **159,** 440–455 (2014).

373 55. Yoshida, N. & Perry, A. C. F. Piezo-actuated mouse intracytoplasmic sperm injection

374    (ICSI). *Nat Protoc* **2,** 296–304 (2007).

375 56. Gusfield, D. Efficient algorithms for inferring evolutionary trees. *Networks* **21,** 19–28

376    (1991).

377 57. Jahn, K., Kuipers, J. & Beerenwinkel, N. Tree inference for single-cell data. *Genome Biol*.

378    **17,** 86 (2016).

379

**Author Contributions**

M.M.C., Z.D.S., A.M. and J.S.W. were responsible for the conception, design, and interpretation of the experiments and wrote the manuscript. M.M.C. and Z.D.S. conducted experiments and M.M.C. developed the analysis, with input from Z.D.S. S.G. and H.K. provided annotations for RNA-seq data and assisted in experimental and analytical optimization. B.A., T.M.N, and M.J. provided vectors, experimental protocols, and advice. J.Q. and D.Y. prepared several sequencing libraries and were engaged in discussion. M.G.J, A.K, and N.Y provided phylogenetic reconstruction strategies.

**Author Information**

The authors declare no competing interests. Correspondence and request for materials should be addressed to Jonathan Weissman (jonathan.weissman@ucsf.edu) and Alexander Meissner (meissner@molgen.mpg.de).

403     **Figure Legends**

404     **Figure 1: Optimization of a multi-purpose molecular recorder**

405         a.  Target site (top) and three guide (bottom) cassettes. The target site consists of an

406             integration barcode (intBC) and three cut sites for Cas9-based recording. Three different

407             single guide RNAs (sgRNAs) are each controlled by independent promoters (in this

408             study, mU6, hU6, and bU6).

409         b.  Molecular recording principle.  Each cell contains multiple genomic, intBC-

410             distinguishable target site integrations. sgRNAs direct Cas9 to cognate cut sites to

411             generate insertion (red) or deletion mutations.  Here, Cas9 is either ectopically delivered

412             or induced by doxycycline.

413         c.  Percentage of uniquely marked reads recovered after recording within a K562 line with

414             10 intBCs for 6 days using the following information: site 1 only with intBCs masked,

415             sites 1-3 (All) with intBCs masked, and sites 1-3 (All) with intBCs considered.

416             Information content scales with number of sites and presence of the intBC.

417         d.  sgRNA mismatches alter mutation rate. Seven protospacers were integrated into the

418             coding sequence of a GFP reporter to infer mutation rate by the fraction of positive cells

419             over a 20 day time course.  Single or dual mismatches were made in guides according to

420             proximity to the PAM: region 1 (proximal), region 2, and region 3 (distal).  Guides

421             against Gal4-4 and the GFP coding sequence act as negative and positive controls. Bold

422             sequences were incorporated into the target site.

423

424     **Figure 2: Lineage tracing in mouse from fertilization through gastrulation**

425  a.  Lineage tracing in mouse experiments. The target site (within mCherry's 3'UTR) and the

426      three guide cassettes are encoded into a single piggyBac transposon vector (ITRs,

427      inverted terminal repeats).  The vector, transposase mRNA, and Rosa26::Cas9:EGFP

428      sperm are injected into oocytes to ensure early integration and tracing in all subsequent

429      cells after zygotic genome activation. Transferred embryos are then recovered after

430      gastrulation.

431  b.  Pearson correlation coefficient heatmap of indel proportions recovered from bulk tissue

432      of an E9.5 embryo (see also **Extended Data Figure 2**).

433  c.  Indel frequency distribution estimated from 40 independent target sites from all embryos.

434      Each site produces hundreds of outcomes for high information encoding.  See **Extended**

435      **Data Figure 4** and **Methods** for frequency calculation. The indel code along the x-axis is

436      as follows: "Alignment Coordinate: Indel Size Indel type (**I**nsertion or **D**eletion)."

437  d.  Proportion of indels that span one, two, or three sites, shown per site. Each dot denotes

438      one of 40 independent intBCs and sums to one across site-spanning indels. Colors

439      indicate the guide array: P = no mismatches; 1 = mismatch in region 1; 2 = mismatch in

440      region 2.

441  e.  Percentage of cells with mutations according to guide complementarity.  Indel

442      proportions within one mouse depend on timing: mutations that happen earlier in

443      development are propagated to more cells.  Dots represent site 1 measurements from

444      independent intBCs; N = 4, 24, and 18 for P, 2, and 1 region mismatches.

445  f.  Indel diversity is inversely related to cutting efficiency for site 1 as in **e**. Early mutations

446      due to fast cutting are propagated to more cells, leading to smaller numbers of unique

447      indels.

448 **Figure 3: Assigning cellular phenotype by scRNA-seq**

449    a.  Images of a lineage-traced E8.5 embryo (embryo 2 of 7 for which single cell data was

450        collected, see **Extended Data Figure 3**), including for Cas9:EGFP and the

451        mCherry:target site.

452    b.  t-sne plot of scRNA-seq from embryo in **a**. Only large or spatially distinct clusters are

453        labeled. (Inset) Pie chart of germ layers. Lighter and darker shades represent embryonic

454        and extra-embryonic components, respectively.  Mesoderm is further separated to include

455        blood (red).  See **Extended Data Figure 5b** for additional embryos.

456    c.  Dot plot of canonical tissue-specific markers.  Grouping clusters of diverse tissue types

457        into germ layers reduces the fraction of marker positive cells, but the specificity to their

458        respective states remains high, especially when considered combinatorially. Size: fraction

459        of marker-positive cells, color intensity: normalized expression (cluster mean). XEcto,

460        extra-embryonic ectoderm/placenta; XEndo, extra-embryonic endoderm/yolk sac; PGC,

461        primordial germ cell; Endo, embryonic endoderm; Ecto, embryonic ectoderm; Meso,

462        embryonic mesoderm; XMeso, extra-embryonic mesoderm.

463

464 **Figure 4: Single cell lineage reconstruction of mouse embryogenesis**

465    a.  Reconstructed lineage tree comprised of 1,732 nodes for embryo 2 with three lineages

466        highlighted.  Each branch represents an indel generation event.

467    b.  Example paths from tree in **a** highlighted by color.  Cells for each node in the path are

468        overlaid onto the plot from **Figure 3b**, with tissue proportions as a pie chart.  Tissue

469        representation decreases with increased tree depth, indicating functional restriction.

470        Bifurcating sublineages are included for the top and bottom paths.  In the top (red) path,

471        this bifurcation occurs within the final branch after primitive blood specification. In the

472        bottom (blue) path, bifurcation happens early within bipotent cells that become either gut

473        or visceral endoderm.

474    c.  Violin plots of the pairwise relationship between lineage and expression for single cells.

475        Lineage distance uses a modified Hamming distance normalized to the number of shared

476        cut sites. Pearson correlation decreases with increasing lineage distance, showing that

477        closely related cells are more likely to share function. Red dot highlights the median,

478        edges the interquartile range, and whiskers the full range.

479    d.  Comparison of shared progenitor scores ($\log_2$-transformed) between our two most

480        information-dense embryos (Embryo 2, n = 1,400 alleles; Embryo 6, n = 2,461 alleles).

481        Cells from closely related transcriptional clusters (ex. primitive blood or visceral

482        endoderm, which have early and late states) derive from common progenitors and score

483        as highly related in both embryos. We also observe a close link between mesoderm and

484        ectoderm that may reflect shared heritage between neuromesodermal progenitors (NMPs)

485        and more posterior neural ectodermal tissues, such as the future spinal cord[42].

486   **Figure 5: Disparities between transcriptional identity and lineage history within the extra-**

487   **embryonic endoderm**

488    a.  Shared progenitor score heatmap for embryo 2 reconstructs expected relationships. The

489        number of nodes that include cells from different lineages is highlighted (Heterogeneous

490        nodes). See **Extended Data Figure 9** for additional embryos.

491    b.  For cells from embryo 2, the relative distance from the mean expression profile of either

492        the endoderm or the extra-embryonic endoderm cluster according to origin (Endo or

493        XEndo).

494      c.  Endoderm cell lineage tree from embryo 2 with expression heatmap for two extra-

495           embryonic marker genes. Middle bar indicates lineage: dark blue, extra-embryonic; light

496           blue, embryonic; grey, ambiguous.

497      d.  Expression boxplots for Trap1a and Rhox5 confirms consistent differential expression

498           across lineage-traced embryos according to their embryonic or extra-embryonic ancestry.

499           Red line highlights median, edges the interquartile range, whiskers the Tukey Fence, and

500           crosses outliers.  N's, the number of recovered XEndo origin cells of either embryonic

501           (E) or Extraembryonic (X) function per embryo.

502   **Figure 6: Lineage bias and estimated size of progenitor pools**

503      a.  Relative tissue distribution of cells descended from reconstructed or profiled pluripotent

504           progenitor cells for embryo 2.  "Profiled" is a unique lineage identity of multiple cells

505           directly observed in the data.  Pluripotent cells form all germ layers, but show

506           asymmetric propensities towards different cell fates, possibly reflecting positional biases.

507           Nodes highlighted in grey with asterisk overlasy give rise to primordial germ cells

508           (lineages 1, 4, and 5 include 9, 1, and 1 PGCs each).  Color assignments as in **Figures 3.**

509      b.  Estimated progenitor field sizes for three types of early developmental potency.

510           Totipotent cells give rise to all cells of the developing embryo, including trophectodermal

511           (TE) lineages. Pluripotent progenitors are partitioned into early and late by generation of

512           extra-embryonic endoderm (XEndo) in addition to epiblast (Epi).  Dots represent single

513           embryos; solid grey line connects estimates from the same embryo.

514

515

516

**Extended Data Figure 1: Target site indel likelihoods from *in vitro* experiments**

    a.  Histograms for the relative indel frequency for protospacer sites 1, 2, and 2b within the

        target region.  In this experiment, single guide RNA expressing vectors respective to each

        position were delivered into K562 cells.  Repair outcomes and frequencies are different

        for each site, but every site produces hundreds of discrete outcomes.  The top 20 most

        frequent indels for each site are shown. The indel code along the x-axis is as follows:

        "Alignment Coordinate: Indel Size Indel type (**I**nsertion or **D**eletion)."  Site 3 was not

        profiled in this experiment.

    b.  Histograms representing the likelihood that any specific base in the target site is deleted

        (blue) or has an insertion (red) which begins at that position, for sites 1 and 2,

        respectively.  The position of the integration barcode (intBC) and protospacer sequences

        (sites) within the target site is represented as a schematic along the bottom, with the PAM

        for each site proximal to the intBC.  Indels, specifically insertions, start at the double

        strand break point 3-bases upstream of the PAM sequence.

    c.  Simultaneous and continuous molecular recording of multiple clonal populations in K562

        cells.  We transduced K563 cells with a high complexity library of unique intBCs, sorted

        them into wells of 10 cells each and propagated them for 18 days.  At the end of the

        experiment, we detected two populations by their intBCs, implying that only two clonal

        lineages expanded from the initial population of 10, and confirmed generation of target

        site mutations. (Left) Strategy for partitioning a multi-clonal population into their clonal

        populations.  Target sites are amplified from a single cell barcoded cDNA library and the

539     intBCs in each cell is identified as present or absent.  (Middle) Heatmap of the percent

540     overlap of intBCs between all cells.  The cells segregate into two populations

541     representing the descendants of two progenitor cells from the beginning of the

542     experiment. (Right) Table summarizing results of the experiment, including the

543     generation of indels over the experiment duration. These data additionally showcase our

544     ability to combine dynamic recording with tracing based on traditional static barcodes.

545

546  **Extended Data Figure 2: Capturing early differentiation by pooled sequencing of indels**

547  **generated within an E9.5 embryo**

548  Scatterplots of indel proportions from dissected, bulk tissue of an E9.5 embryo.  Placenta is the

549  most distantly related from embryonic tissues, followed by the yolk sac, with the three

550  embryonic compartments sharing the highest similarity.

551

552  **Extended Data Figure 3: Experimental overview**

553     a.  Schematic of platform used for generation of single cell RNA-seq libraries and

554         corresponding target site amplicon libraries, adapted from Adamson et al., 2016 (Ref 18).

555         The barcoded and amplified cDNA library is split into two fractions prior to shearing:

556         one fraction is used to generate a global transcription profile and the other is used to

557         specifically amplify the target site.

558     b.  Summary table of lineage traced embryos detailing the type of guides used, the sampling

559         proportion, and sequencing results.  Embryo 4 was omitted from further analysis due to

560         the absence of cells identified as primitive heart tube.

561

562

**Extended Data Figure 4: Target site capture in mouse embryos**

a. Percentage of cells with at least one target site captured.

b. Scatterplot showing the relationship between the mean number of unique molecular identifiers (UMIs, a proxy for expression level) sequenced per target site and the percentage of cells in which the target site is detected, which we refer to as "target site capture." Generally, as the mean number of UMIs increases, the percentage of cells also increases. Using a full length, intron-containing Ef1a promoter in mouse embryos leads to a higher number of UMIs, which generally results in better target site capture.

c. Percent of cells for which a given integration barcode (intBC) is detected across all seven embryos profiled in this study.

d. Target site capture and expression level across tissues for Embryo 5, which utilizes a truncated Ef1a promoter to direct transcription of the target site. Each row corresponds to a different intBC, indicated in the top left of the histogram. (Left) The percentage of cells in each tissue for which the target site is captured. (Right) Violin plots represented the distribution of UMIs for the target site in each tissue. Dashed line refers to a 10 UMI threshold. The target site may be expressed at different levels in a tissue-specific manner, leading to higher likelihoods of capture in certain tissues. Examples such as the target sequences carrying the intBCs AGGACAAA and ATTGCTTG may also be explained by mosaic integration after the first cell cycle, as these follow a developmental logic and are preferentially expressed in extraembryonic tissues. White dot indicates the median UMI count for cells from a given germ layer, edges the interquartile range, and whiskers the full range of the data.

585     e.  Target site capture and expression level across tissues for embryo 7, which drives the

586         target site expression from a full length Ef1a promoter.  Each row corresponds to a

587         different intBC, indicated in the top left of the histogram.  (Left) The percentage of cells

588         in each tissue for which the target site is captured.  (Right) Violin plots represented the

589         distribution of UMIs for the target site in each tissue as in **d**.  Dashed line is a visual

590         threshold for 10 UMIs.  While tissue specific expression may explain some discrepancy

591         in target site capture, high expression (as estimated from number of UMIs) may still

592         correspond to low capture rates, as observed for the intBC TGGCGGGG.  One possibility

593         is that certain indels may destabilize the transcript and lead to either poor expression or

594         capture.

595     f.  Scatterplots showing the relationship between estimated relative indel frequency and the

596         median number of cells that carry the indel.  Since the indel frequency within a mouse is

597         dependent on the timing of the mutation, we cannot calculate the underlying indel

598         frequency distribution using the fraction of cells within embryos that carry a given indel.

599         Instead, we estimate this frequency by the presence or absence of an indel using all of the

600         target site integrations across mice, which reduces biases from cellular expansion but still

601         assumes that any given indel occurs only once in the history of each intBC.  Since the

602         number of integrations is small (<50), we might expect our estimates to be poor.  Here we

603         see that the number of cells marked with an indel increases with indel frequency,

604         suggesting that our frequency estimates are under-estimated for particularly frequent

605         indels.  This is likely due to the fact that we cannot distinguish between identical indels in

606         the same target site that may have resulted from multiple repair outcomes (convergent

607         indels).  The most frequent insertions are of a single base and tend to be highly biased

608    towards a single nucleotide (eg. 92:1I is uniformly an "A" in 5 out of 7 embryos and

609    never < 88%).

610

611    **Extended Data Figure 5: single cell RNA-seq tissue assignment and wild type comparison**

612    a.  Boxplots representing tissue proportions from E8.0 (top) and E8.5 (bottom) wild type

613        embryos (n = 10 each) with lineage-traced embryos mapping to each state overlaid as

614        dots.  Wild type embryos display large variance in the proportions of certain tissues and

615        our lineage-traced embryos generally fall within the range of those recovered from wild

616        type.   Large circles indicate embryos that were scored as either E8.0 or E8.5,

617        respectively, and the bold red overlay highlights embryo 2, which is used throughout the

618        text.  Note that many processes are continuous or ongoing between E8.0 to E8.5, such as

619        somitogenesis and neural development.  For example, from E8.0 to E8.5, the embryonic

620        proportions of anterior neural ectoderm and fore/midbrain are inversely correlated as one

621        cell type presumably matures into the other.  Many of our embryos scored as E8.0 exhibit

622        intermediate proportions for both tissue types, supporting the possibility that these

623        embryos are somewhat less developed than E8.5 but more developed than E8.0.  For

624        boxplots, center line indicates the median, edges the interquartile range, whiskers the

625        Tukey Fences, and crosses the outliers.

626    b.  Plots (t-sne) of single cell RNA-seq with corresponding tissue annotations for the six

627        lineage traced embryos used in this study.  (Inset) Pie chart of the relative proportions for

628        different germ layers.  Mesoderm is further separated to include blood (red).  While 36

629        different states are observed during this developmental interval, only broad classifications

630        of certain groups (eg. "neural ectoderm" or "lateral plate mesoderm") are overlaid to

631        provide a frame of reference.  In general, the relative spacing and coherence of different

632        cell states are consistent across different embryos.

633    c.  Boxplots of the Euclidean distance between single cell transcriptomes and the average

634        transcriptional profile of their assigned cluster (cluster center) in comparison to their

635        distance from the average of the next closest possible assignment.  Comparison is to the

636        same 712 informative marker genes used to assign cells to states and includes all cells

637        used in this study.  Middle bar highlights the median, edges the interquartile range,

638        whiskers the Tukey Fences, and grey dots the outliers.  N's refer to the cumulative

639        number of cells assigned to each state across all 7 embryos for which single cell data was

640        collected, including for embryo 4.

641

642    **Extended Data Figure 6. Continuous indel generation by breeding**

643    a.  Strategy for generating lineage traced mice through breeding.  The target site and guide

644        array cassette are integrated into mouse zygotes as in **Figure 2a** using C57Bl/6J sperm to

645        generate $P_0$ breeder mice, which are capable of transmitting high copy genomic

646        integrations of the technology.  Then, $P_0$ animals are crossed with homozygous,

647        constitutively expressing Cas9 transgenic animals to enable continuous cutting from

648        fertilization onwards in $F_1$ progeny.  Shown is Sibling 2 of a cross between a $P_0$ male and

649        a Cas9:EGFP female.

650    b.  Bar charts showing the degree of mutation (% cut, red) for a $P_0$ male (top row) and 4 $F_1$

651        offspring generated by breeding with a Cas9:EGFP female prior to weaning (21 days post

652        partum). Each row represents a mouse and each column represents a target site.  Each

653      sibling inherits its own subset of the 23 parental target site integrations, and demonstrates

654      different levels of mutation throughout gestation and maturation.

655    c.   Indel frequencies for the 10 most frequent indels from 3 siblings in a common target site

656      integration (column 1 in **b**). Each mouse shows a large diversity of indels and the

657      different frequencies observed in each animal demonstrates an independent mutational

658      path.

659

660 **Extended Data Figure 7: Performance of tree building algorithms used on embryonic data**

661    a.   Table summarizing contemporary Cas9-based lineage tracers that have been applied to

662      vertebrate development highlighting attributes that differ between the studies. Refer to

663      **Methods** for a more detailed overview of key characteristics of our technology. * Study

664      reports the average fraction recovered by tissue for integrations that cannot be

665      distinguished, such that percentages reported here are effectively equivalent to our "≥1

666      intBC" metric. ** Reports a plate-based DNA-sequencing approach that can be applied to

667      all methods to improve target site recovery. *** Range of cells where at least one intBC

668      is confidently detected and scored. **** Presents a tree reconstruction method, but results

669      predominantly on clonal analysis.

670    b.   Table of allele complexity, number of nodes, and log-likelihood scores for embryos.

671      Tree likelihoods are calculated using indel frequencies estimated from all embryo data

672      (see **Extended Data Figure 5** and **Methods**). Bold scores indicate the reconstruction

673      algorithm selected for each embryo (see **Figure 4**, and **Extended Data Figures 8** and **9**).

674    c.   Log likelihood of trees generated using either the greedy or biased sampling approach as

675      a function of complexity, which is measured as the number of unique alleles. There is

676        near equivalent performance of the two algorithms for low complexity embryos, but the

677        greedy algorithm produces higher likelihood trees for embryos with larger numbers of

678        unique alleles.

679

680        **Extended Data Figure 8: Single cell lineage reconstruction of early mouse development for**

681        **embryo 6**

682        a.  Reconstructed lineage tree comprised of 2,690 nodes generated from our most

683            information-dense embryo (embryo 6), which we used to compare shared progenitor

684            scores with embryo 2 in **Figure 4d**.  Each branch represents an independent indel

685            generation event, and each node contains a pie chart of the germ layer proportions for the

686            cells contained within it (colors are as in **Figure 3b**).

687        b.  Example paths from root to leaf from the selected tree (highlighted by color).  Cells for

688            each node in the path are overlaid onto the t-sne representation in **Extended Data Figure**

689            **5**, with the tissue proportion at each node in the tree included as a pie chart. In the top

690            most path (pink), the lineage bifurcates into two independently fated progenitors that

691            either generate mesoderm (secondary heart field/splanchnic plate mesoderm and

692            primitive heart tube) or neural ectoderm (anterior neural ectoderm and neural crest).

693            Note that the middle path (green) also represents an earlier bifurcation from the same tree

694            and eventually produces neural ectoderm (neural crest and future spinal cord).  These

695            paths begin with a pluripotent node that can generate visceral endoderm but subsequently

696            lose this potential.  Alternatively, the bottom path (dark blue) begins in an equivalently

697            pluripotent state but becomes restricted towards the extraembryonic visceral endoderm

698            fate.

699      c. Violin plots representing the relationship between lineage and expression for individual

700          pairs of cells as calculated for embryo 2 in **Figure 4c**. Expression Pearson correlation

701          decreases with increasing lineage distance, showing that closely related cells are more

702          likely to share function. Red dot highlights the median, edges the interquartile range, and

703          whiskers the full range.

704      d. Comprehensive clustering of shared progenitor scores for Embryo 6, which has the

705          greatest number of unique alleles and samples multiple extraembryonic tissues. Shared

706          progenitor score is calculated as the sum of shared nodes between cells from two tissues

707          normalized by the number of additional tissues that are also produced (a shared

708          progenitor score is calculated as $2^{-(n-1)}$ where n is the number of clusters present within

709          that node). In general, extraembryonic tissues that are specified before implantation, such

710          as extraembryonic endoderm or ectoderm, co-cluster away from embryonic tissues and

711          within their own groups, while the amnion and allantois of the extraembryonic mesoderm

712          cluster with other mesodermal products of the posterior primitive streak. The co-

713          clustering of anterior paraxial mesoderm and somites may reflect the continuous nature of

714          somitogenesis from presomitic mesoderm during this period, with production of only the

715          most anterior somites by E8.5. Note that the gut endoderm cluster has been further

716          portioned according to embryonic or extra embryonic lineage (see **Figure 5**).

717

718 **Extended Data Figure 9: Summary of results from additional mouse embryos**

719 Representative highest likelihood tree analyses for additional embryos, including:

720      a. Reconstructed trees as shown in **Figure 4a**.

721    b.  Shared progenitor score heatmaps as shown in **Figure 5a**, normalized to the highest score

722        for each embryo to account for differences in total node numbers.  Here, the shared

723        progenitor score is calculated as the number of nodes that are shared between tissues

724        scaled by the number of number of tissues within each node (a shared score is calculated

725        as $2^{-(n-1)}$ where n is the number of clusters present within that node).  In general,

726        clustering of shared progenitors is recapitulated across embryos, with mesoderm and

727        ectoderm sharing the highest relationship and either extra-embryonic ectoderm or extra-

728        embryonic endoderm representing the most deeply rooted and distinct outgroup, though

729        these scores are sensitive to the number of target sites and the rate of cutting.  By shared

730        progenitor, PGCs are also frequently distant from other embryonic tissues, but this often

731        reflects the rarity of these cells, which restricts them to only a few branches of the tree in

732        comparison to more represented germ layers.  The number of heterogeneous nodes from

733        which scores are derived is included for each heatmap.

734    c.  Violin plots representing the pairwise relationship between lineage distance and

735        transcriptional profile as shown for embryo 2 in **Figure 4c**.  Lineage distance is

736        calculated using a modified Hamming distance and transcriptional similarity by Pearson

737        correlation.  The exact dynamic range for lineage distance depends on the number of

738        intBCs included and the cutting rate of the three guide array.  Here, distances are binned

739        into perfect (0), close ($0 > x > 0.5$), intermediate ($0.5 \leq x < 1$), and distant ($x \geq 1$)

740        relationships for all cells containing either 3 or 6 cut sites, depending on the embryo.  As

741        lineage distance increases, transcriptional similarity decreases, consistent with functional

742        restriction over development. Red dot highlights the median, edges the interquartile

743        range, and whiskers the full range.

744

**Extended Data Figure 10: Expression characteristics of extra-embryonic and embryonic**

**endoderm**

    a.  Violin plots representing the pairwise scRNA-seq Pearson correlation coefficients for within or across group comparisons according to lineage (X, extra-embryonic; E, embryonic) and cluster assignment (light blue, gut endoderm; dark blue, visceral endoderm). Within group comparisons for cells with the same lineage and transcriptional cluster identity are shown on the left, while across group comparisons are presented on the right. Notably, extraembryonic cells with gut endoderm identities show higher pairwise correlations to embryonic cells with gut endoderm identities (column 4) than they do to visceral endoderm cells, with which they share a closer lineage relationship (column 5). Red dot highlights the median, edges the interquartile range, and whiskers the full range.

    b.  Plots (t-sne) of scRNA-seq data for embryo 2, with gut endoderm cells highlighted. Endoderm cells segregate from the rest of the embryo, and cannot be distinguished by embryonic (light blue) or extraembryonic (dark blue) origin.

    c.  Expression boxplots for the extra-embryonic markers Trap1a and Rhox5 from an independent single cell RNA-seq survey of E8.25 embryos (Ibarra-Soria et al., 2018, Ref [9]). Both genes are heterogeneously present in cells identified as mid/hindgut but uniformly present in canonical extra-embryonic tissues, consistent with a subpopulation of cells of extra-embryonic origin residing within this otherwise embryonic cluster. Red lines highlights the median, edges the interquartile range, and whiskers the Tukey Fence. Outliers were removed for clarity.

767

768 **Methods**

769

770 **Plasmid design and construction**

771 Because the principles governing Cas9 efficiency and subsequent indel generation are not

772 absolute, we screened fourteen protospacers for potential inclusion in our target site, including

773 nine protospacers known to function with moderate efficiency and five additional protospacers

774 hypothesized to function[43-52]. Each protospacer was checked against the human and mouse

775 genomes using bowtie to limit off target effects. A gene block library of the fourteen

776 protospacers (no additional bases between sequences) with an 8 base pair randomer was ordered

777 from IDT representing target site version 0.0.

778

779 The target site (tS) v0.0 vector backbone was derived from a previously described Perturb-seq

780 lentiviral vector (pBA439, Addgene, Cat#85967)[18] with the following changes: the cassette for

781 mU6-sgRNA-EF1a-PURO-BFP was removed and replaced with EF1a-tSv0.0-sfGFP using

782 Gibson assembly with the target site in the coding sequence of sfGFP for use in the fluorescent

783 reporter assay (PCT10, sequence available upon request).

784

785 A gene block library of five protospacers (ade2-whiteL-bam3-bri1-whiteB; no additional bases

786 between sequences) with an 8 base pair randomer was ordered from IDT representing target site

787 version 0.1. Protospacers in positions 1 (ade2), 3 (bam3), and 5 (whiteB), are used for cutting in

788 subsequent experiments and are referred to as sites 1, 2, and 3.

789

790 Target site (tS) v0.1 was also cloned into pBA439 with the following changes: the cassette for

791 mU6-sgRNA-EF1a-PURO-BFP was removed and replaced with EF1a-sfGFP-tSv0.1, followed

792    by BGH pA on the original backbone (PCT12).  Here the target site sits in the 3' UTR of GFP.

793    To improve the delivery of multiple targets into the same cell, we swapped the v0.1 target site

794    cassette into a commercially available piggyBac transposon vector (Systems Biosciences,

795    #PB533A-2) with the following changes: IRES-Neo was swapped for either GFP (PCT16) or

796    mCherry (PCT29).  The backbone was digested with restriction enzymes and target site v0.1

797    gene block was PCR-amplified to add Gibson arms.  Following Gibson assembly, the plasmids

798    were transformed into at least 100uL of Stbl2 competent cells (Thermo Fisher, Cat#10268019),

799    and plated onto 1-2 large plates (Fisher, #NC9372402) with LB/Carbenicillin to generate high

800    complexity target site libraries (PCT17, and PCT30, respectively).

801

802    The three-guide expression vector design and cloning protocol were adapted from [18] to utilize

803    guides against the three sites in the target site.  The guide for site 1 (ade2) is under the control of

804    the mU6 promoter, site 2 (bam3) under the control of hU6 promoter, and site 3 (whiteB) under

805    the control of bU6-2 promoter.  All guides are constitutively expressed in this system.

806    Additionally, the triple-guide cassette was moved onto the piggyBac backbone described above.

807

808    Two further modifications of the plasmids described above were used in this study.  First, in an

809    attempt to decrease the cutting percentage variation between embryos, we cloned the triple-guide

810    expression cassette without BFP into PCT29, and then cloned in the target site with intBCs to

811    generate the resulting vectors (PCT41-43, for guide combinations (P,1,P), (1,1,1), and (2,1,2),

812    respectively).  In the second modification, we changed the truncated form of Ef1a in PCT29 to a

813    promoter sequence comprised of the ubiquitous chromatic opening element (UCOE) and a full-

814    length, intron-containing Ef1$\alpha$ and cloned in a triple-guide expression cassette for the guide

815    combination (2,1,P), followed by cloning in of the target site to make PCT60. In these

816    modifications, target site plasmid libraries (PCT41-43, PCT60) were transformed and expanded

817    in 1-2L of liquid LB/Carb culture rather than on large plates.

818

819    A new target site design, v1.1, was utilized for further experiments to generate $P_0$ breeders (see

820    below). A gene block library of three protospacers (ade2-bri1-whiteB; 30-60 bases between

821    sequences) with a 14 base pair randomer was ordered from IDT representing target site version

822    v1.1. For this target site, site1 is ade2, site2 is bri1, and site3 is whiteB. We cloned v1.1 into the

823    same backbone as PCT60 with guide combinations (2,3,3) or (2,1,2) to make PCT61 and PCT62,

824    respectively.

825

826    **Cell culture, DNA transfections, and viral production**

827    The production of lentiviral particles or transfection of plasmids as is as described in[18].

828

829    **K562 GFP reporter assay**

830    To construct the target site GFP reporter cell line, a doxycycline(Dox)-inducible Cas9 K562 cell

831    line was stably transduced with PCT10 (8% infected, <0.1 MOI), and GFP positive cells were

832    sorted using fluorescence activated cell sorting on a BD FACSAria2. For each protospacer in the

833    target site, 1-4 guides was designed to achieve a series of mutation efficiencies and cloned into

834    single guide expression vectors[22]. On Day -4, the reporter cell line was plated into wells and

835    stably transduced with a different guide against target site v0.0, GFP-targeting protospacer

836    EGFP-NT2 (positive control), or Gal4-targeting protospacer (negative control) in each well. On

837    Day -2, cells were selected for guide cassette integration using 3 ug/mL puromycin. On Day 0,

838    50ng/mL Dox was added to induce Cas9 expression, and maintained through the course of the

839    experiment.  GFP fluorescence was recorded on a LSR-II flow cytometer (BD Biosciences) on

840    every $2^{nd}$ day starting at day 0, except day 13 was recorded in place of day 12.  Data was

841    analysed in Python using FlowCytometryTools (http://eyurtsev.github.io/FlowCytometryTools/).

842    For guide virus produced in this experiment, labels were systematically shifted during production

843    resulting in incorrect ordering of guide effect on GFP fluorescence, which was corrected for

844    presentation in the manuscript.  We confirmed the activity order of the guide series for three

845    guides (ade2, bam3, and bri1) in sequencing experiments where new virus was prepared.

846

847    **K562 single cutting pooled assay**

848    To construct the cell line used here, a Dox-inducible Cas9 K562 cell line was stably transduced

849    with PCT12 (6% infected, <0.1 MOI), and GFP positive cells were sorted on a BD FACSAria2.

850    On Day -5, the cell line was plated and stably transduced with a different guide against target site

851    v0.1, or GFP-targeting protospacers in each well.  On Day -2, cells were selected for guide

852    cassette integration using 3 ug/mL puromycin.  On Day 0, 50ng/mL Dox was added to induce

853    Cas9 expression, and maintained through the course of the experiment.  Wells were sampled

854    every 3-6 days for 20 days with cell pellets frozen down.  Genomic DNA was isolated from

855    frozen cell pellets, and the target site was PCR-amplified to make sequencing libraries (refer to

856    **Pooled embryo library preparation** below for library prep protocol), which were sequenced on

857    the Illumina Miseq.  Timepoint samples were pooled and reads with no indels were removed to

858    calculate relative indel frequencies.

859

860    **K562 multiple target site integration cell line**

861 To construct a cell line with multiple integrations, we nucleofected 200,000 Dox-inducible Cas9

862 K562 cells with 1500ng PCT17 and 200ng piggyBAC transposase using set program T-016

863 (Lonza #V4SC-2096; Systems Biosciences, #PB210PA-1).

864

865 **K562 triple guide cutting assay, and multi-clonal lineage tracing experiment**

866 Multiple-integration cells described above were stably transduced with a triple guide expression

867 vector (Perfect-Perfect-Perfect; fastest cutting) and recovered for 2 days. GFP (target site) and

868 BFP (triple guide) double positive cells were sorted using fluorescence activated cell sorting on a

869 BD FACSAria2. For the multi-clonal lineage tracing experiment, 10 cells were sorted into wells

870 containing 200uL of pre-conditioned media on a 96 well plate (12 wells total). At day 18, wells

871 were inspected under the microscope and the 3 wells with the largest populations were selected

872 for single cell analysis on the 10x Chromium. Two of the lanes suffered wetting failures, and the

873 remaining sample was taken through library preparation described below (refer to **Target site**

874 **amplicon library preparation**). The library was sequenced on the Illumina Miseq and would

875 benefit from additional sequencing.

876

877 For the pooled experiment, ~112,000 cells were sorted into a tube, spun down, resuspended in

878 fresh media, split into two wells with 50ng/mL Dox added to one of the wells. Cells were

879 collected 6 days post-sort, genomic DNA was isolated, and the target site was PCR-amplified to

880 make sequencing libraries (refer to **Pooled embryo library preparation**), which were

881 sequenced on the Illumina Miseq. The 10 intBCs with the most reads were used for analysis.

882

883 **Embryo and $P_0$ breeder generation**

884 Protocols are adapted from those described in ref [53] To enable *in vivo* lineage tracing, B6D2F1

885 strain female mice (age 6 to 8 weeks, Jackson Labs) were superovulated by sequential

886 intraperitoneal injection of Pregnant Mare Serum Gonadotropin (5IU per mouse, Prospec Protein

887 Specialists) and Human Chorionic Gonadotropin (5IU, Millipore) 46 hours apart. Twelve hours

888 after delivery of the second hormone, MII stage oocytes were isolated and injected with *in vitro*

889 transcribed piggyBAC transposase mRNA (100 ng/ul) prepared in an injection buffer (5 mM

890 Tris buffer, 0.1 mM EDTA, pH = 7.4). Decapitated sperm isolated from an 8 week old

891 *Gt(ROSA)26Sortm1.1(CAG=cas9\*,EGFP)Fezh/J* strain mouse (Jackson labs, ref [54]) was

892 resuspended with the purified piggyBAC library in the same injection buffer at concentrations

893 ranging from 0.5 to 1.4 ug/uL.

894 Transposase-injected oocytes were then fertilized by piezo-actuated intracytoplasmic

895 sperm injection (ICSI) as previously described ref [55]. Injected embryos were cultured in 25 uL

896 EmbryoMax® KSOM drops (Millipore) covered in mineral oil (Irvine Scientific) in batches of

897 25-50 embryos. After 84 or 96 hours, successfully cavitated blastocysts were screened for

898 uniform fluorescence of the target sequence cassette and transferred into one uterine horn of 6-10

899 week old pseudopregnant CD-1 strain female mice (Charles River). Uterine transfer results in an

900 ~24 hour lag, so the day of transfer was scored as E2.5 and embryos were dissected from

901 euthanized animals 6 or 7 days later at ~E8.5 or E9.5, depending on the experiment. All

902 techniques utilized standard micromanipulation equipment, including a Hamilton Thorne XY

903 Infrared laser, Eppendorf Transferman NK2 and Patchman NP2 micromanipulators, and a Nikon

904 Ti-U inverted microscope.

905 The generation of breeders was conducted identically by coinjecting target design v1.1

906 piggyBAC plasmids with sperm from C57BL6/J strain males (Jackson labs), transferring

907    uniformly bright mCherry blastocysts into CD-1 strain mice, and allowing live pups to be

908    brought to term. Genotyping was conducted using tail tip genomic DNA purified using the

909    Quick DNA Miniprep Plus kit (Zymogen) isolated prior to weaning. Animals with large intBC

910    counts (n=23 for the male used in **Extended Data Figure 6**) were then bred into either male or

911    female *Gt(ROSA)26Sortm1.1(CAG=cas9\*,EGFP)Fezh/J* strain animals to generate live pups

912    with continuous cutting. Fluorescence of live animals was confirmed and documented using a

913    dual fluorescent protein flashlight (Nightsea).

914

915    **Pooled embryo library preparation**

916    RNeasy Mini Kit (Qiagen, #74104) was used to isolate RNA from whole embryos or dissected

917    tissue for embryonic tissue. Alternatively, genomic tail tip DNA was used for $P_0$ breeders or

918    Cas9+ $F_1$ animals. Following purification and/or first strand synthesis of cDNA from 1 ug of

919    RNA (Promega), the target site was amplified using a 2-stage PCR protocol. In the $1^{st}$ stage,

920    <100ng of diluted DNA template was amplified using 0.6 uM forward and reverse primers and

921    Kapa HiFi HotStart ReadyMix according to the following PCR protocol: (1) 98C for 3 min, (2)

922    98C for 30 s, 69C for 30 s, 72C for 15 s (16 cycles for cDNA, 24 cycles for genomic DNA), (3)

923    72C for 5 min. Following 0.7X SPRI selection, the elute served as template for $2^{nd}$ stage PCR,

924    using 0.6uM barcoded P5 and P7 secondary primers and Kapa HiFi HotStart ReadyMix

925    according to the following PCR protocol: (1) 98C for 3 min, (2) 98C for 30 s, 60C for 30 s, 72C

926    for 30 s (4-6 cycles), (3) 72C for 5 min. PCR products underwent 0.6X SPRI-selection and were

927    eluted in 20-40uL of elution buffer to produce the final library. Libraries were sequenced on the

928    Illumina HiSeq 2500 (Rapid Run) or Miseq, with the following run parameters: Read 1: 175

929    cycles, i7 index: 8 cycles, i5 index: 8 cycles, Read 2: 175 cycles.

930

931    For v1.0 target sites, the following primary primers were used:

932    MC38:

933    CGTCGGCAGCGTCAGATGTGTATAAGAGACAGTGCAGGAGCGGATTGCTTCGAACC

934    MC39:

935    TCTCGTGGGCTCGGAGATGTGTATAAGAGACAGACAACCACTACCTGAGCACCCAG

936    TC

937    For v1.1 target sites, the following primary primers were used:

938    P5_PCT48-49_F:

939    <u>TCGTCGGCAGCGTC</u>**AGATGTGTATAAGAGACA**<u>**GAATCCAGCTAGCTGTGCAGC**</u>

940    ODY120_PCT48_R_PB:

941    GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGCGATGGACGATTGCGGAAGAC

942    AG

943    Secondary amplification was conducted using the following primers:

944    P5 primer:

945    AATGATACGGCGACCACCGAGATCTACAC[ILLUMINA

946    INDEX]*TCGTCGGCAGCGTC*AGATGTGTA)

947    P7 primer

948    CAAGCAGAAGACGGCATACGAGAT[ILLUMINA

949    INDEX]*GTCTCGTGGGCTCGG*AGATGTGTATAAG

950

951    **Single cell embryo dissociation**

952    Embryos are washed through several drops of PBS after isolation to reduce debris and put into

953    ~100 uL PBS droplets on a microscope slide and screened for uniform fluorescence of the target

954    site cassette on an Olympus IX71 inverted microscope running Metamorph.  Selected embryos

955    were dissociated to single cell suspensions by adding 100 uL of TrypLE (Invitrogen, #12605010)

956    and pipetting the embryo or embryo pieces every 5 minutes for ~30 minutes until complete

957    dissociation was visually confirmed.  Trypsin was deactivated by adding 100 uL  PBS+BSA is

958    added to the droplet and moving cells into a 1.5 mL eppendorf tube, followed by several rounds

959    of additional collection with 100 to 200 uL of PBS+BSA to a final volume of 1 mL.  The

960    dissociated cells are filtered through a Flowmi filter tip (Bel-Art Products, #H13680-0040) into a

961    new tube, and spun down for 5 minutes at 1200 rpm on a tabletop centrifuge.  Following the

962    spin, 900uL of PBS+BSA is removed and the remaining volume is resuspended with an

963    additional 900uL of PBS+BSA.  The suspension is spun for 5 minutes at 1,200 rpm, 800 uL of

964    PBS+BSA is removed, the remaining volume is spun for 5 minutes at 1,200 rpm, and PBS+BSA

965    is removed until only ~30 uL of volume remains. 2 uL of the final resuspended cells were used

966    for counting using a hemocytometer.  We load ~17,000 cells into the 10x machine (Chromium

967    Single Cell 3' Library & Gel Bead Kit v2) for a targeted recovery of 10,000 cells.

968

969    **scRNA-seq library preparation and sequencing**

970    Single cell RNA-seq libraries were prepared according to the 10x user guide, except for the

971    following modification.  After cDNA amplification, the cDNA pool is split into two fractions.

972    15uL of EB buffer is added to one of the fractions of 20uL of the cDNA pool, and scRNA-seq

973    library construction proceeds as directed in the 10x user guide.  RNA-seq libraries were

974    sequenced on the Illumina HiSeq 4000 system.

975

**Target site amplicon library preparation**

The target site-specific amplification protocol was adapted from [11]. 50-100 ng of template from

the cDNA pool, 0.3 uM P5-truseq-long

(AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATC

T), 0.6 uM MC63

(TCTCGTGGGCTCGGAGATGTGTATAAGAGACAGTGCAGGAGCGGATTGCTTCGAAC

C) was split across four parallel PCR reactions, and was amplified using Kapa HiFi HotStart

ReadyMix according to the following PCR protocol: (1) 95C for 3 min, (2) 98C for 15 s, then

69C for 15 s (8-12 cycles). Reactions were re-pooled during 0.9X SPRI selection, and eluted

into 60 uL. A second PCR with the elute as the template, 0.3 uM P5

(AATGATACGGCGACCACCGA), 0.6 uM barcoded P7

(CAAGCAGAAGACGGCATACGAGAT[ILLUMINA

INDEX]GTCTCGTGGGCTCGGAGATGTGTATAAG) was split across four parallel PCR

reactions, and amplified using Kapa HiFi HotStart ReadyMix according to the following PCR

protocol: (1) 95C for 3 min, (2) 98C for 15 s, then 69C for 15 s (6 cycles). Reactions were re-

pooled during 0.9X SPRI selection and then fragments of length 200-600bp were selected using

the BluePippin. Target site libraries were sequenced on the Illumina HiSeq 2500 (Rapid Run),

with the following run parameters: Read 1: 26 cycles, i7 index: 8 cycles, i5 index: 0 cycles, Read

2: 350 cycles.

**scRNA-seq library data processing**

997    scRNA-seq data was processed and aligned using 10x Cell Ranger v2. The filtered gene-barcode

998    matrices were then processed in Seurat (https://satijalab.org/seurat/) for data normalization

999    (global scaling method "LogNormalize"), dimensionality reduction (PCA), and generation of t-

1000   sne plots, which use the first 16 principal components.

1001

1002   **scRNA-seq tissue assignment**

1003   An independent project conducting scRNA-seq profiling of gastrulation identified 42 distinct

1004   tissues in wild type mice. We utilized the mean expression profile for each tissue and the list of

1005   712 marker genes used for assignment of cells to tissues (see instructions for assignment in

1006   GSE122187). For each cell in lineage traced mouse embryos, we calculated the Euclidean

1007   distance between the cell's expression profile and the mean expression profile for each tissue

1008   using the 712 marker gene set, and assigned the cell to the tissue identity with the minimum

1009   distance. Expression values were transformed to log space using log(normalized UMI count + 1)

1010   before calculating the Euclidean distance. Comparisons between the best matched tissue to the

1011   next best match are presented for all data collected here in **Extended Data Figure 5** to highlight

1012   the precision of this approach.

1013

1014   **Embryo gastrulation stage assignment**

1015   The wild type mouse gastrulation compendium consists of five time points, profiling every 0.5

1016   days from E6.5 to E8.5 with at least 10 embryos collected for each time point. Tissue proportion

1017   is calculated as the number of cells assigned to the tissue divided by the total number of cells in

1018   the embryo. The median tissue proportion was calculated for each time point treating each tissue

1019   independently. For each lineage-traced embryo, the Euclidean distance between its tissue

1020  proportions and the median tissue proportion for each time point was calculated and the embryo

1021  was assigned to the time point with the minimum cumulative distance. All lineage-traced

1022  embryos were assigned to either E8.0 or E8.5 stages.

1023

1024  **Target site data processing**

1025  A custom software pipeline was built to align and call indels in the target site. The logic is as

1026  follows: (1) Assign cell barcode and UMIs to each read, (2) find the consensus sequence for each

1027  UMI, (3) align the consensus sequence to the target site reference sequence, (4) identify most

1028  likely integration barcodes (intBC) and create custom reference sequences, (5) repeat alignment

1029  against all reference sequences and select highest scoring alignment for each UMI, (6) call intBC

1030  and indels in the target site, (7) correct the intBC and allele using UMIs which appear in the

1031  same cell, (8) remove doublets. Details appear below:

1032

1033  (1) Assigning cell barcode and UMIs to each read. Specific amplification libraries of the target

1034  site amplicon were processed using 10x Cell Ranger software to assign cell barcodes and UMIs

1035  to each read. The target site is designed to be orthogonal to the human and mouse genome, and

1036  does not align in Cell Ranger processing. Unaligned reads from the Cell Ranger output bam file

1037  are parsed into fastq format with the cell barcode and UMI identifiers appended to the read

1038  name.

1039

1040  (2) Finding the consensus sequence for each UMI. To potentially increase the speed of consensus

1041  sequence finding, we attempt to trim reads to the same length for each UMI. The read is

1042  trimmed to remove sequence beyond the polyA tail using cutadapt software

48

1043 (http://cutadapt.readthedocs.io/en/stable/) with the following parameters: [–a AAAAAAAAAA –

1044 e 0.1 –o trimmedFile.fq –untrimmed-output=untrimmedFile.fq –m 20 –max-n=0.3 –trim-n].

1045 Reads that do not contain polyA sequence appear in the untrimmed file and are subjected to a

1046 second round of read trimming using a sequence which appears in the 3' end of the target site

1047 assuming the sequence has not been deleted from DNA repair, with cutadapt run using the

1048 following parameters: [-a GCTTCGTACGCGAAACTAGCGT -e 0.1 -o trimmedFile2.fq --

1049 untrimmed-output=untrimmedFile2.fq -m 20 --max-n=0.3 --trim-n --no-indels].  The adapter

1050 sequence used in the last round of trimming is then concatenated back on to the trimmed

1051 sequence to improve target site alignment in the next step.  If >=60% of trimmed sequences for a

1052 given UMI are the same sequence, then the sequence is taken as the consensus sequence.

1053 Otherwise, a multiple sequence alignment is performed using BioPython and the consensus

1054 sequence is extracted from the alignment.  Ambiguous bases are reported if there is <50%

1055 agreement for any position in the alignment.

1056

1057 (3) Aligning to the target site reference. We use the emboss implementation of the smith-

1058 waterman algorithm to align sequences to the target site reference sequence with the following

1059 parameters, which were determined empirically: [emboss water –asequence targetSiteRef.fa –

1060 sformat1 fasta –bsequence consensusUMI.fa –sformat2 fasta –gapopen 15.0 –gapextend 0.05 –

1061 outfile sam –aformat sam].  In this first alignment, the ambiguous sequence NNNNNNNN is

1062 used to represent the intBC.  A minor bug had to be corrected in the emboss implementation to

1063 successfully output sam format.  For target site v1.1, the gapopen penalty was increased to 20

1064 and the gapextend penalty to 1.

1065

1066    (4) Identifying the most likely intBC. A perl script is used to parse the intBC from the alignment.

1067    The intBCs with the highest number of UMIs are substituted into the target site reference

1068    sequence to make custom reference sequences.  This step was included because upon manual

1069    inspection, there were obvious misalignments due to the ambiguous intBC sequence, which were

1070    corrected upon substitution of a real sequence.

1071

1072    (5) Selecting the highest scoring alignment for each UMI. Repeat smith waterman alignment

1073    against all custom reference sequences and select alignment with the highest score for each UMI.

1074

1075    (6) Calling indels and intBCs. A perl script is used to parse the intBC and indels from the

1076    alignment using the CIGAR string.  The boundaries for each site is defined and indels

1077    overlapping site boundaries are called as an indel in that site.  Sequence of the indel is not

1078    considered.

1079

1080    (7) Correcting indels using multiple reads with the same UMI from the same cell.  UMIs are

1081    filtered for alignment score and only cells that are in the matched scRNA-seq data set are kept.

1082    An intBC is corrected to an intBC with a higher UMI count in the cell if the intBCs are within an

1083    edit distance of 2 and the alleles are the same. An allele is the combination of indels in sites 1, 2,

1084    and 3.  An allele is corrected to an allele with a higher UMI count in the cell if the intBC is the

1085    same and the allele is within a 1-indel difference.  Only UMIs with greater than or equal to 3

1086    UMIs are kept.

1087

1088 (8) Eliminating doublets. Cells that report two alleles for the same intBC are removed if the

1089 dominant allele is <80% of the total UMI count for the intBC. This removes 4.1-18.3% of cells

1090 in our embryos.

1091

1092 **Tree reconstruction strategies**

1093 **1. Biased search through phylogenetic space**

1094 We simulate the evolutionary process leading from a collection of uncut target sites to the final

1095 data set. The set of mutations (including "no mutation") across all target sites in a cell is referred

1096 to as an allele. In the final tree, each branch represents a mutation, and each node represents the

1097 allele of a cell, which may be a reconstructed ancestral allele, i.e. it is not present in the data set.

1098 Input: table of unique alleles

1099 - each allele may represent multiple cells

1100 - we cannot distinguish between identical indels in the same position that may result from

1101    independent mutation events (convergent indels) if they appear with an identical set of co-

1102    segregating indels

1103 Algorithm:

1104 - Create root node in tree representing an allele with 0 mutations (c_allele)

1105 - remove alleles in the table that match c_allele

1106

1107 - While alleles remain in table:

1108    - choose indel from table that can be added to current allele

1109       - can only add indels in positions that have no mutation

1110    - create new node by adding indel into c_allele (c_allele2)

1111    - draw directed edge labeled with indel between nodes from c_allele to c_allele2

1112    - remove alleles in table that match c_allele2

1113        - includes alleles that match c_allele2 with missing values for positions that have no

1114            mutations

1115    - if indels in table can be added to c_allele2, then c_allele = c_allele2; else, c_allele does

1116        not change

1117    - when indels cannot be added to c_allele, traverse up edges to ancestral nodes until an

1118        allele to which an indel can be added is found

1119

1120    We presented two methods that are used to choose indels. The first method, "Random," selects a

1121    position where an indel can be added, and then selects an indel from the data set for that position;

1122    both selections occur in an unbiased manner. The second method, "Frequency Normalized

1123    Weighted" (FNW), identifies all of the indels that can be added to the current allele and scores

1124    them according to the fraction of alleles they are found in divided by the expected independent

1125    frequency of the indel (see **Fig. 2c**). These scores are used as weights to bias selection of the

1126    indel. The reasoning behind FNW is that indels that are found in many cells (or alleles) are more

1127    likely to have occurred early, but this has to be balanced against their expected likelihood of

1128    occurring. FNW biases the search towards more likely trees. To further increase the search for

1129    good trees, we first remove all indels that are unique to a single allele since we can assume that

1130    these indels occur at the leaves of the tree. The indels are added as branches leading to leaves in

1131    the final tree before the final tree likelihood is calculated.

1132

1133    The log likelihood of the tree is calculated as the sum of the likelihoods of all the indels that

1134    appear in the tree.  The likelihood of each mutation is estimated from the embryo data set (**Fig.**

1135    **2c**).

1136

1137    It is worth noting that the number of trees that are possible grossly exceeds 30,000; however, the

1138    search is biased towards finding good trees and performs markedly better than those that are

1139    randomly generated.  Using high scoring trees to direct the search towards better ones, such as by

1140    grafting high scoring branches, could further improve our algorithm's ability to identify high

1141    scoring trees.

1142

1143    **2. Greedy search to reconstruct larger trees**

1144    Our greedy algorithm consists of building the tree top-down, recursively splitting cells into

1145    mutually exclusive groups based on the presence or absence of a specific mutation. In particular,

1146    these splits are prioritized by selecting mutations that appear frequently in the dataset, but are

1147    known to be an improbable outcome from a Cas9 mutagenesis event. This transform is equal to

1148    the product of the observed frequency of the mutation and the log prior-probability. The

1149    mutations prioritized this way, we reason, are very likely to have occurred only once and

1150    relatively early in the experiment. Under this assumption, these mutations are useful to a top-

1151    down approach as they efficiently create maximally informative tree-splits. In practice, we can

1152    calculate the prior-probabilities of mutations several ways but while describing this algorithm we

1153    assume the priors are provided (**Fig. 2c**).

1154

1155    To deal with missing values, we first split cells based on the presence or absence of a mutation.

1156    Then, for each cell that reports a missing value for this cut site, we assign the cell to the group

1157    with which it shares the greatest similarity. Here, we define similarity as the average number of

1158    mutations it shares with the cells in each group.  We follow this procedure until only one cell

1159    remains. Note that for application to the dataset described in this manuscript, we filled missing

1160    values with unique indels to force the algorithm to choose splits based on the presence of

1161    mutations rather than absence.

1162

1163    Theoretically, building the tree in this fashion is possible due to the special case of multistate

1164    compatibility afforded by our model of evolution, namely that mutations can only arise once at a

1165    particular site (i.e. Cas9 cannot re-cut a site). This context allows one to consider a traditional

1166    Gusfield algorithm[56] in which one infers phylogenies by selecting character-splits based on the

1167    most frequently occurring mutations. In a special regime of "perfect-phylogeny" (where every

1168    mutation arose exactly once), this algorithm is provably optimal and extremely efficient as

1169    compared to other algorithms (linear in the number of cells and mutations, or O(|number cells| *

1170    |number of mutations|)). In the case of multi-state characters, the notion of compatibility often

1171    breaks down as this typically implies that a character can mutate many times to different states.

1172    Yet, as described previously, in our system this cannot happen – namely, once a mutation is

1173    obtained at a site, it cannot be changed again along that evolutionary path. In this way, we can

1174    apply an approximated Gusfield algorithm to reconstruct trees, where perfect phylogeny is

1175    possible although still confounded in cases where the same mutation arises twice independently.

1176

1177    Trees are visualized using the Python ete library (http://etetoolkit.org/).

1178

1179 **Pairwise single cell lineage distance measure used for violin plots**

1180 A cut site can take 2 forms, uncut or indel. The distance is a modification of hamming distance

1181 where uncut is considered a special state.

1182 Distance = (2*(sites with different indels) + 1*(sites with indel vs uncut))/(number of sites

1183 recovered in both cells)

1184 Pairwise expression correlation was estimated using the same 712 marker genes used to assign

1185 cell states and was only included if two single cell transcriptomes shared ≥10 gene

1186 measurements.

1187

1188 **Estimating ancestral tissue relationships**

1189

1190 Each node, including leaves, that includes more than one tissue type is considered a

1191 "progenitor." Progenitors found at the leaves are not reconstructed or inferred but result from the

1192 lack of new indels that distinguish between tissues (ie. the lineage tracer does not produce new

1193 indels past the progenitor stage).

1194

1195 The shared progenitor score is calculated between two tissues as the number of shared

1196 progenitors scaled by the number of tissues each progenitor contributes to, and is calculated

1197 using the following algorithm:

1198

1199 For each progenitor,

1200        tList = list of tissues progenitor contributes to

1201        pScore = 1/(2^len(tList)-1)

1202           for each pair of tissues in tList:

1203               progenitorScoreForPairOfTissues += pScore

1204    Example for a single progenitor:

1205        tList = [Endo, Meso, XMeso]

1206        pScore = 1/(2^(3-1)) = ¼

1207        ProgenitorScoreEndoMeso += ¼

1208        ProgenitorScoreEndoXMeso += ¼

1209        ProgenitorScoreMesoXMeso += ¼

1210

1211    The resulting matrix is a shared progenitor score matrix.  To transform the similarity matrix to a

1212    distance matrix, we use 1-(matrix/maxScoreInMatrix).  The distance matrix is then hierarchically

1213    clustered using either ward or average as the cluster joining criteria..

1214

1215    To account for the potential effect of cluster sizes (for example, if we assume that differentiation

1216    occurs for all tissues instantaneously, then the larger cluster sizes for mesoderm and ectoderm

1217    would increase the likelihood of detecting a progenitor between the two tissues), we

1218    downsampled each tissue before calculating the shared progenitor score: 150 cells were

1219    randomly sampled from each tissue and the tree was pruned to only include the sampled cells.

1220    For tissues with less than 150 cells, all cells were included.  For embryo 2, we downsampled to

1221    300 cells since it is a merger of two biological replicates and is therefore doubly sampled.  The

1222    shared progenitor score was calculated from the pruned tree and the process was repeated 1000

1223    times for each embryo.  The median progenitor score is presented in the heatmap.  For higher

1224   resolution clusters (**Fig. 4d**, **Extended Data Fig. 8**), we downsampled 500 times instead of 1000

1225   times.

1226

1227   Note that the number of nodes reported below the heatmaps in **Extended Data Figure 8**

1228   represents the number of progenitors that are found in the complete tree.  The number of nodes

1229   used to calculate the shared progenitor score depends on the sampled set of cells chosen.

1230        For high resolution shared progenitor scores calculated for embryos 2 and 6 (**Fig. 4d** and

1231   **Extended Data Fig. 8**), we bolstered some populations prior to calculating shared progenitor

1232   scores by merging some cluster identities if they represent the linear maturation of one tissue

1233   type to another, are primarily one cluster versus the other at the assigned developmental time

1234   point, and have very close transcriptional profiles.  Specifically, we merged node with

1235   notochord, amnion mesoderm (early) with amnion mesoderm (late), primitive blood progenitor

1236   with primitive blood (early), and anterior paraxial with pharyngeal (arch) mesoderm.  We also

1237   merged surface and preplacodal ectoderm due to the similarity of their transcriptional profiles

1238   and omitted "similar to neural crest 2" as this transcriptional cluster is ambiguously determined

1239   (the cluster is globally most similar to neural crest but not obviously comprised of specific

1240   marker genes).

1241

1242   **Endoderm lineage assignment and differentially regulated gene identification**

1243   Endoderm cells can have one of three origins based upon our tree: extra-embryonic, embryonic,

1244   or ambiguous.  Cells are considered extra-embryonic if there is a progenitor in its lineage whose

1245   descendants include ≥ 40% extra-embryonic cells.  Cells have ambiguous origin if they descend

1246 directly from the root node. Otherwise, cells are considered to be from embryonic origin. We

1247 identified endoderm cells of extra-embryonic origin in all embryos but embryo 7.

1248

1249 We use the Kolmogorov-Smirnov test (Python scipy.stats.ks_2samp) to identify differentially

1250 regulated genes between embryonic and extra-embryonic origin endoderm cells. Only highly

1251 variable genes in the embryo are considered for testing, and genes are significant if they have a

1252 Bonferroni corrected p-value under 0.05.

1253

1254 **Multipotent field size estimation and asymmetry**

1255 Progenitors are considered pluripotent if their descendants include at least one mesoderm (Meso

1256 or XMeso or Blood) cell, one ectoderm (Ecto) cell, and one endoderm (Endo) cell. A pluripotent

1257 progenitor are considered early pluripotent if it also has at least one extra-embryonic endoderm

1258 descendant, and further considered totipotent if it has at least one extra-embryonic ectoderm

1259 descendant. To estimate the lower bound for the number of multipotent cells, we start at the

1260 bottom level of the tree and count the number of multipotent cells at that level. If multipotent

1261 cells exist, then the number of multipotent cells is propagated to its ancestor in the above level,

1262 otherwise we count 0 for that level. We add one progenitor for every level that includes a

1263 multipotent cell and other cells to represent the progenitor that lead those non-multipotent cells at

1264 that level. The number of multipotent cells is then the number of cells propagated to the root of

1265 the tree. Progenitor asymmetry is simply the proportion of descendants from each of the tissues

1266 for that node.

1267 **Comparison to other technologies**

1268    Several CRISPR-Cas9 based lineage tracers have been developed, each with distinct strengths

1269    and weaknesses. In **Extended Data Figure 7**, we present a table summarizing the different

1270    technologies, and elaborate on the attributes that, in combination, distinguish our strategy here:

1271    1. Target sites are marked with a unique integration barcode (intBC). The intBC allows us to

1272        phase our target sites and perform a direct comparison for each target site across cells. This

1273        greatly improves the information content of our system as it allows us to distinguish between

1274        the same indel if it appears in different target sites (**Fig. 1c**).

1275    2. Guide RNAs are integrated into genomic DNA and constitutively expressed from

1276        totipotency, which enables our lineage tracer to be truly evolving over multiple cell

1277        generations. In technologies applied to zebrafish development, guideRNAs are expressed as

1278        a pulse, which leads to the generation of a large diversity of barcodes at one or two

1279        timepoints.

1280    3. Multiple integrations of multi-cutsite target sites are distributed throughout the genome.

1281        Technologies that integrate a single target site with many cut sites or have tandem

1282        integrations are subject to collapse of information when one indel may affect neighboring cut

1283        sites or alternatively, simultaneous cutting of several cut sites remove large portions of their

1284        lineage tracer. While our technology is also vulnerable to these effects, we are better

1285        buffered against them by distributing the target sites throughout the genome. We also

1286        highlight that indel generation is largely independent within target sites when slower cutters

1287        are used (**Fig. 2d-f**).

1288    4. Simultaneous, multi-population lineage tracing (**Extended Data Figure 1c**). Since target

1289        sites are labeled with integration barcodes, we can use the identity of these barcodes to

1290        deconvolute pools of cells upon sequencing. Alternatively, independent samples, such as

1291    embryos that have unique sets of integration barcodes, can be pooled onto a single 10x lane

1292    to decrease the cost of reagents.

1293  5.  Multi-channel recording using our triple guide vector.  In our current manuscript, we use the

1294      three channels for lineage tracing but different types of sensors can be developed to record

1295      multiple independent inputs.

1296  6.  Ability to trace over different time scales by tuning the mutation rate of the system through

1297      mismatches in the guide RNA.

1298  To fully utilize the information captured in our data set, we developed custom reconstruction

1299  strategies to identify the maximum likelihood tree (see **Tree reconstruction strategies**

1300  above).  We estimate indel likelihoods using all of our embryo data (**Fig. 2c**), which allows us to

1301  estimate tree likelihoods rather than utilize maximum parsimony criteria.  Phylogenetic

1302  algorithms developed for tumor evolution, such as SCITE[57], offer conceptual frameworks that

1303  are compelling to adapt for our technology.

1304

1305  **Code availability**

1306  Code will be shared upon request.

1307  **Data Availability**

1308  The data is available in the GEO database under accession numbers GSE117542 for lineage

1309  traced embryos and GSE122187 for the gastrulation compendium.

1310

1311
1312

**Target site cassette**

Ef1 α   FP   intBC   3 2 1   Sites   pA

Target site

**Three guide cassette**

U6   1   2   3

Uncut   + Cas9   Indels
         + 3x sgRNA

**c**   Data incorporated
n = 10 intBCs

Unique reads (%)

Site    1    All   All
intBC   −    −    +

**d**

% GFP

Controls        bri1        bam3        ade2
Gal4-4                                   site 2      site 1
sgGFP

white-O    white-B    white-L    ChrI_106412
           site 3

Days

Rate

5′ sgRNA

3
2
1

P
Constant Region

■ = mismatch

a

TPase mRNA

Rosa26::Cas9:EGFP

piggyBAC library
Target site    gRN A array

ITR                    ITR

Fertilization

Target site

E8.5–E9.5

b

E9.5
Body
Head
Tail
Yolk sac
Placenta

Indel proportion

0    Correlation    1

c

Relative indel frequency (%)

Site 1: 510 unique indels

Site 2: 545 unique indels

Site 3: 265 unique indels

d

Guide array
● 1,1,1  ● 2,1,2  ● 2,1,P  ● P,1,P

% indels

Site    1  2  3    1  2  3    1  2  3
Span    1 site      2 sites     3 sites

e

% cells cut

P    2    1
sgRN A match to cut site 1

f

Unique indels

P    2    1

Embryo 2

0.2 mm

Cas9:EGFP

mCherry

b

Ecto          PGC
                    Endo
Meso

Notochord
Future spinal cord
Fore/Midbrain          Primitive blood
                              NMPs
Neural crest     Anterior neural
                          ectoderm
Surface and        PGCs
preplacodal                    Presomitic          Visceral
ectoderm                       mesoderm           endoderm
                    Posterior lateral
Gut endoderm            plate
            Allantois                    Somites
Amnion    Splanchnic lateral
                    plate              Anterior paraxial
                                          mesoderm
Primitive heart tube     Hematopoietic and
                                  Endothelial precursors
            Angioblasts

n = 22,264 cells

t-sne 2

t-sne 1

XEcto
XEndo
XMeso
Blood
Meso
Ecto
Endo
PGC

% positive
0   50  100

Log(UMI+1)
0          3

Tfap2c Gata2 Elf5 Afp Spink1 Apoa4 Hnf4a Foxf1 Bmp4 Plagl1 Gata1 Klf1 Hoxb1 T Cdx1 Meox1 Hand1 Twist1 Twist2 Sox2 Otx2 Pax6 Ptn Hes3 Foxa2 Sox17 Dppa3 Pou5f1 Klf5

**b**

Root → Leaf

| n = 3,481 | n = 306 | n = 16 | n = 4, n = 2 |
| n = 1,899 | n = 29 | n = 2 | |
| n = 560 | n = 65, n = 435 | n = 11, n = 7 | n = 3, n = 2 |

n = 90,943 pairwise comparisons

**Correlation** vs **Lineage distance (9 cut sites)**

0.0  0.11  0.22  0.33  0.44  0.55  0.66  0.77  0.88  1.0  1.11  1.22  1.33  1.44  1.55

**d**

Embryo 2 vs Embryo 6

Visceral endoderm: early vs late
Primitive blood: early vs late
Anterior paraxial vs Somites
Future spinal cord vs NMP (early)
Angioblasts vs Hematopoietic and endothelial progenitors

**Shared progenitor score**

Meso      n = 6,239
Ecto      n = 6,025
XMeso      n = 967
Blood      n = 1,755
Endo      n = 643
XEndo      n = 292
PGC      n = 41
XEcto      n = 1

Total nodes: 1,732
Heterogeneous nodes: 314

−1   0   1
$\mathrm{Log}_{10}$

c

Lineage     Trap1a   Rhox5

Count

300
200
100
0

Lineage
Endo
XEndo

Endo     0     XEndo
Distance to RNA cluster center

**Trap1a**        **Rhox5**

Norm. RNA count

1   2   3   5   6     1   2   3   5   6

3
2
1
0

X ( n = 21)
E ( n = 75)
X ( n = 31)
E ( n = 186)
X ( n = 61)
E ( n = 36)
X ( n = 50)
E ( n = 240)
X ( n = 10)
E ( n = 145)

Pluripotent progenitors

Ecto
Meso
Blood
XMeso
Endo
PGC